

On discretisation drift and smoothness regularisation in neural network training

Mihaela Claudia Rosca

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
of
University College London.

Department of Computer Science
University College London

June 14, 2023

I, Mihaela Claudia Rosca, confirm that the work presented in this thesis is my own. Where information has been derived from other sources, I confirm that this has been indicated in the work.

Abstract

The deep learning recipe of casting real-world problems as mathematical optimisation and tackling the optimisation by training deep neural networks using gradient-based optimisation has undoubtedly proven to be a fruitful one. The understanding behind *why* deep learning works, however, has lagged behind its practical significance. We aim to make steps towards an improved understanding of deep learning with a focus on optimisation and model regularisation. We start by investigating gradient descent (GD), a discrete-time algorithm at the basis of most popular deep learning optimisation algorithms. Understanding the dynamics of GD has been hindered by the presence of discretisation drift, the numerical integration error between GD and its often studied continuous-time counterpart, the negative gradient flow (NGF). To add to the toolkit available to study GD, we derive novel continuous-time flows that account for discretisation drift. Unlike the NGF, these new flows can be used to describe learning rate specific behaviours of GD, such as training instabilities observed in supervised learning and two-player games. We then translate insights from continuous time into mitigation strategies for unstable GD dynamics, by constructing novel learning rate schedules and regularisers that do not require additional hyperparameters. Like optimisation, smoothness regularisation is another pillar of deep learning's success with wide use in supervised learning and generative modelling. Despite their individual significance, the interactions between smoothness regularisation and optimisation have yet to be explored. We find that smoothness regularisation affects optimisation across multiple deep learning domains, and that incorporating smoothness regularisation in reinforcement learning leads to a performance boost that can be recovered using adaptations to optimisation methods. We end by showing that optimisation can also affect smoothness, as discretisation drift can act as an implicit smoothness regulariser in neural network training.

Impact Statement

Our work focuses on understanding and improving components of deep learning systems. As such, it provides two avenues for impact: first, by providing insights that can be useful to other researchers in the field and second, by helping downstream applications that use such systems.

Inside deep learning academic research, this work provides a set of novel tools and insights concerning two main building blocks of deep learning systems, namely optimisation and models. Our insights include novel continuous-time flows to analyse optimisation dynamics, as well as a rethinking of the effect of smoothness constraints imposed on deep learning models. We have accompanied theoretical results with experimental frameworks for analysis and validation, including approaches that can stabilise training and reduce hyperparameter sensitivity and thus computational costs. Future theoretical research directions based on our line of work include specialising our theoretical results to specific classes of neural networks, further expanding the family of available continuous-time flows capturing optimisation dynamics, and finding new beneficial or detrimental implicit regularisation forces in deep learning optimisation. Promising avenues for practical impact include new methods for training deep learning models, and finding new approaches for model selection. Throughout this work, we provide insights across multiple deep learning domains, including supervised learning, two-player games, and reinforcement learning. The work included here has been presented at multiple conferences and published in journals.

While we do not tackle specific applications directly, this work can have impact outside the research domain by improving the stability of deep learning models used in applications across multiple domains; examples include image classification and generation, as well as game play using reinforcement learning.

Acknowledgements

One of the biggest opportunities I have been provided with throughout my PhD studies was to have not one, but two incredible supervisors. My primary supervisor, Marc Deisenroth, has been an wonderful source of guidance and strength. Marc, you have raised my bar for rigour in thinking and writing, and have shown me that technical excellence does not come in spite of kindness, but due to kindness; for that, I will be eternally grateful. My second supervisor, Shakir Mohamed, has been my champion and supporter for many years. Shakir, working with you has been the privilege of a life-time; your breadth of technical knowledge is outstanding, and your careful and deliberate application of it through a meticulous choice of impactful research directions should serve as a guide to us all.

I have also been lucky to have an amazing set of collaborators throughout my PhD: Andriy Mnih, Arthur Gretton, Benoit Dherin, Chongli Qin, Claudia Clopath, David G.T. Barrett, Florin Gogianu, Lucian Busoniu, Michael Figurnov, Michael Munn, Razvan Pascanu, Theophane Weber, Tudor Berariu, and Yan Wu. Special thanks to Benoit Dherin, for many evening chats on everything to do with deep learning optimisation and regularisation; to Yan Wu, for a long term fruitful collaboration; to Theophane Weber for many chats on reinforcement learning; to Tudor Berariu and Florin Gogianu for being so much fun while doing research; and to Arthur Gretton for guidance early in my PhD.

This work was done as part of the DeepMind-UCL PhD program, and I am grateful to those organising it and championing this program. Thanks also Frederic Besse, and to Sarah Hodgkinson and Claudia Pope for their support.

During my PhD I was lucky to write a book chapter on Implicit Generative models with Balaji Lakshminarayanan and Shakir Mohamed, for Kevin P. Murphy's second edition of 'Machine Learning: a Probabilistic Perspective'. Beyond allowing

me to share my enthusiasm for this area of machine learning, a great clarity of thinking followed the writing process, which helped fledge some of the new ideas presented in this work. I thus must thank Kevin, Balaji, and Shakir for this opportunity.

Thanks also go to my examination committee, Ferenc Huszár and Patrick Rebeschini, for their time and an incredible and thoughtful discussion on my work.

I also have many reasons to be grateful to those who have supported me personally during this time and beyond. I am here because my parents have prioritised my education and encouraged me throughout the way. I also had the most loving of grandmothers, and I am sure she would have loved to see me pursue this degree.

I have been fortunate to have a great set of friends, especially in the tough times of coronavirus lockdowns. Niklas, you have been there for me since the times of the Imperial labs, and have always had faith in me; there are no words that can express how grateful I am. Fabio, you are a consistent, reliable, and wise friend; you are also raising my favourite humans, who are a great source of happiness for me. Many thanks also go to Ada, Dhruva, David, George, James, Ira, Sascha, Siqi, Slemi, Victor, and many other friends, for the time spent together and the great memories.

Contents

1	Introduction	2
1.1	Optimisation in deep learning	4
1.2	Smoothness w.r.t. inputs	7
1.3	Interactions between optimisation and smoothness regularisation	9
2	Optimisation in deep learning	14
2.1	The direction of steepest descent	14
2.1.1	Computing the gradient	15
2.2	Algorithms: from gradient descent to Adam	16
2.3	Continuous-time methods	24
2.3.1	A brief overview of stability analysis	24
2.3.2	Backward error analysis	26
2.3.3	BEA proofs	27
2.4	Multiple objective optimisation	29
2.5	Conclusion	32
3	A new continuous-time model of gradient descent	33
3.1	Introduction	33
3.2	Continuous time models of gradient descent	35
3.2.1	Limitations of existing continuous-time flows	36
3.3	The principal flow	40
3.3.1	The PF and the Hessian eigen-decomposition	44
3.3.2	The stability analysis of the PF	47
3.4	Predicting neural network gradient descent dynamics with the PF	50

3.4.1	Predicting $(\nabla_{\theta} E)^T \mathbf{u}_0$ using the PF	51
3.4.2	Around critical points: escaping sharp local minima and saddles	53
3.5	The PF, stability coefficients and edge of stability results	55
3.6	Non-principal terms can stabilise training	63
3.7	Stabilising training by accounting for discretisation drift	64
3.7.1	$\nabla_{\theta}^2 E \nabla_{\theta} E$ determines discretisation drift	65
3.7.2	Drift adjusted learning (DAL)	67
3.7.3	The trade-off between stability and performance	69
3.8	Future work	73
3.9	Related work	76
3.10	Conclusion	79
4	Continuous time models of optimisation in two-player games	80
4.1	Background	81
4.1.1	Generative adversarial networks	82
4.2	Discretisation drift in two-player games	84
4.2.1	DD for simultaneous Euler updates	85
4.2.2	DD for alternating Euler updates	86
4.2.3	Sketch of the proofs	87
4.2.4	Visualising trajectories	88
4.3	The stability of DD	88
4.4	Common-payoff games	92
4.5	Analysis of zero-sum games	93
4.5.1	Dirac-GAN: an illustrative example	94
4.6	Experimental analysis of GANs	97
4.6.1	Does DD affect training?	98
4.6.2	The importance of learning rates in DD	99
4.6.3	Improving performance by explicit regularisation	100
4.6.4	Extension to non-zero-sum GANs	103
4.7	A comment on different learning rates	104
4.8	Extending the PF for two-player games	109

4.8.1	Stability analysis	112
4.8.2	An empirical example and different learning rates	114
4.9	Related work	115
4.10	Conclusion	117
5	Finding new implicit regularisers by revisiting backward error analysis	118
5.1	Revisiting BEA	119
5.2	Implicit regularisation in multiple stochastic gradient descent steps .	121
5.3	Implicit regularisation in generally differentiable two-player games . .	126
5.3.1	The effect of the interaction terms: a GAN example	129
5.4	Conclusion	131
6	The importance of model smoothness in deep learning	132
6.1	Why model smoothness?	132
6.2	Measuring smoothness	133
6.3	Smoothness regularisation in deep learning	134
6.4	The benefits of learning smooth models	136
6.5	Challenges with smoothness regularisation	142
6.5.1	Too much smoothness hurts performance	142
6.5.2	Interactions between smoothness regularisation and optimisation	145
6.5.3	Sensitivity to data scaling	147
6.6	Conclusion	147
7	Smoothness and optimisation effects of Spectral Normalisation	149
7.1	Introduction	149
7.2	Reinforcement learning	150
7.2.1	Deep reinforcement learning	152
7.2.2	Challenges with optimisation in DRL	153
7.2.3	Reinforcement learning environments	154
7.3	Spectral Normalisation in RL	155
7.3.1	Applying Spectral Normalisation to a subset of layers	155

7.3.2	Experimental results on the Atari benchmark	156
7.3.3	An optimisation effect	158
7.4	Revisiting Spectral Normalisation in GANs	166
7.4.1	Spectral Normalisation and the discriminator output	167
7.4.2	Interactions with loss functions	170
7.4.3	Hypothesis: Adam in the low-curvature regime	174
7.5	Related work	177
7.6	Conclusion	179
8	Geometric Complexity: a smoothness complexity measure implicitly regularised in deep learning	180
8.1	Introduction	182
8.2	Geometric Complexity	184
8.2.1	GC for deep linear models	186
8.2.2	GC for rectifier feed-forward networks	186
8.2.3	GC for convolutional and residual layers	187
8.3	GC increases as training progresses	188
8.4	GC and regularisation	190
8.4.1	Explicit regularisation	191
8.4.2	Implicit regularisation via discretisation drift	193
8.5	Double descent and GC	196
8.6	Related work	197
8.7	Limitations and future work	199
8.8	Conclusion	200
9	Conclusion	201
A	On a new continuous-time model of gradient descent	243
A.1	Proofs	243
A.1.1	BEA proof structure	243
A.1.2	Third-order flow	244
A.1.3	Higher order terms of the form $(\frac{df}{d\theta})^n f$ or $(\nabla_{\theta}^2 E)^n \nabla_{\theta} E$	246

A.1.4	Linearisation results around critical points	255
A.1.5	The solution of the PF for quadratic losses	256
A.1.6	The Jacobian of the PF at critical points	257
A.1.7	Jacobians of the IGR flow and NGF at critical points	259
A.1.8	Multiple gradient descent steps	259
A.1.9	Approximations to per-iteration drift for gradient descent and momentum	261
A.2	Comparison with a discrete-time approach	262
A.2.1	Changes in loss function	262
A.2.2	The dynamics of $\nabla_{\theta} E^T \mathbf{u}_i$	263
A.2.3	The connection between DAL and Taylor expansion optimal learning rate	263
A.3	Experimental details	264
A.4	Additional experimental results	267
A.5	Figure reproduction details	281
B	Continuous time models of gradient descent in two-player games	286
B.1	Proof of the main theorems	286
B.1.1	Simultaneous updates (Theorem 4.2.1)	287
B.1.2	Alternating updates (Theorem 4.2.2)	290
B.2	Proof of the main corollaries	296
B.2.1	Common-payoff alternating two-player games	296
B.2.2	Zero-sum simultaneous two-player games	298
B.2.3	Zero-sum alternating two-player games	299
B.2.4	Self and interaction terms in zero-sum games	299
B.3	Discretisation drift in Runge–Kutta 4	300
B.3.1	Runge–Kutta 4 for one player	300
B.3.2	Runge–Kutta 4 for two players	300
B.4	Stability analysis	304
B.4.1	Simultaneous Euler updates	304
B.4.2	Alternating Euler updates	306

B.5	SGA in zero-sum games	308
B.6	DiracGAN - an illustrative example	310
B.6.1	Reconciling discrete and continuous updates in Dirac-GAN	310
B.6.2	DD changes the convergence behaviour of Dirac-GAN	312
B.6.3	Explicit regularisation stabilises Dirac-GAN	313
B.7	The PF for games: linearisation results around critical points	314
B.8	GAN Experimental details	316
B.8.1	Implementing explicit regularisation	316
B.9	Additional experimental results	317
B.9.1	Additional results using zero-sum GANs	317
B.9.2	GANs using the non-saturating loss	317
B.9.3	Explicit regularisation in zero-sum games trained using simultaneous gradient descent	318
B.9.4	Explicit regularisation in zero-sum games trained using alternating gradient descent	321
B.10	Individual figure reproduction details	329
C	Finding new implicit regularisers by revisiting backward error analysis	331
C.1	Two consecutive steps of SGD	332
C.2	Multiple steps of SGD	334
C.3	Multiple steps of full-batch gradient descent	338
C.3.1	Expectation over all shufflings	339
C.4	Two-player games	341
C.4.1	Effects on the non-saturating GAN	341
D	The importance of model smoothness in deep learning	343
D.1	Individual figure reproduction details	343
E	Spectral Normalisation in Reinforcement learning	345
E.1	Additional experimental results	345
E.1.1	The weak correlation between smoothness and performance	345

E.1.2	Other regularisation methods	345
E.1.3	DIVOUT, DIVGRAD and MULEPS	347
E.2	Experimental details	350
E.2.1	Application of Spectral Normalisation	350
E.2.2	Measuring smoothness: the norm of the Jacobians	350
E.2.3	Evaluation	351
E.2.4	Minatar	351
E.2.5	Atari experiments	352
E.3	Additional hyperparameter sweeps	353
E.4	Individual figure reproduction details	356
F	Geometric Complexity: a smoothness complexity measure implicitly regularised in deep learning	358
F.1	Figure and experiments reproduction details	358
F.2	Additional experimental results	359

List of Figures

- 1.1 Overview of optimisation in deep learning. 6
- 1.2 The importance of smoothness on the model fit. 8
- 2.1 The importance of the learning rate in gradient descent. 18
- 2.2 The importance of the learning rate and decay rate in gradient descent
with momentum. 21
- 2.3 The effect of hyperparameters in Adam. 24
- 3.1 Defining discretisation drift. 37
- 3.2 Motivation for the Principal Flow (PF): existing flows fail to capture
the oscillatory or unstable behaviour of gradient descent. 38
- 3.3 Complex flows can capture oscillations and divergence of gradient
descent around local minima. 38
- 3.4 Per-iteration error between the NGF, IGR and gradient descent. . . . 39
- 3.5 Comparing the coefficients α_{NGF} and α_{PF} based on the value of αh .
While α_{NGF} is always -1 , the value of α_{PF} depends on αh and can
be complex, both with positive and negative real parts. 45
- 3.6 The behaviour of the PF on $E(z) = \frac{1}{2}z^2$ 46
- 3.7 The PF captures the behaviour of gradient descent exactly for
quadratic losses. 48
- 3.8 The behaviour of the PF when the loss is the banana function. 48
- 3.9 Small neural networks: comparing continuous flows for multiple itera-
tions. 50
- 3.10 Predictions of $\nabla_{\theta} E^T \mathbf{u}_0$ according to continuous-time flows for a VGG
model trained on CIFAR-10. 53

3.11 Sharp minima are not attractive to gradient descent, a neural network experiment.	54
3.12 The edge of stability behaviour in neural networks.	55
3.13 Comparison of how continuous-time models capture the behaviour of gradient descent at edge of stability.	56
3.14 Using the PF to understand stability and instability in a 4 layer MLP trained on MNIST.	58
3.15 Using the PF and stability coefficient to understand instabilities in deep learning.	59
3.16 Understanding the behaviour of the loss through the PF and the behaviour of λ_0	59
3.17 Assessing whether 1 dimension is enough to cause instability in a continuous-time flow.	60
3.18 Decreasing the learning rate in the edge of stability area leads to increased stability.	61
3.19 The unstable dynamics of $\nabla_{\theta} E^T \mathbf{u}$ in the edge of stability area.	62
3.20 The value of non-principal terms; a Resnet-18 trained on CIFAR-10 sweep.	64
3.21 $\ \nabla_{\theta}^2 E \nabla_{\theta} E\ $ and the per-iteration drift measured during training.	66
3.22 Correlation between $\ \nabla_{\theta}^2 E \nabla_{\theta} E\ $ and the per-iteration drift.	67
3.23 Models trained using a learning rate sweep or DAL on CIFAR-10 and Imagenet. Models are VGG and Resnet-50.	68
3.24 Learning rate and update norms in DAL compared to fixed learning rate training; a Resnet-18 model trained on CIFAR-10.	69
3.25 DAL-0.5 results on a CIFAR-10 and Imagenet using VGG, and Resnet-50 model respectively.	70
3.26 DAL- p sweep on full-batch training on CIFAR-10 with VGG and Resnet-18 models and a Resnet-50 model trained on Imagenet.	70
3.27 VGG model trained on CIFAR-10 using full batch gradient descent, either with a fixed learning rate in a sweep or a DAL- p sweep.	71

3.28	DAL and gradient descent learned landscapes; batch size 64 on CIFAR-10.	72
3.29	DAL- p and momentum Resnet-50; model trained on Imagenet.	74
3.30	Per-parameter DAL- p ; batch size 8192 on Imagenet.	75
4.1	Visualisation of our approach using backward error analysis (BEA) to find new continuous-time flows describing the behaviour of gradient descent in two-player games.	85
4.2	Modified flows given by BEA better capture discrete dynamics on a two dimensional example.	89
4.3	Discretisation drift can change the stability of a game.	91
4.4	In common-payoff games alternating updates lead to higher gradient norms and unstable training when equal learning rates are used.	93
4.5	Discretisation drift can cause instability in DiracGAN; cancelling the sources of instabilities leads to convergence.	96
4.6	The effect of discretisation drift on zero-sum games.	98
4.7	Alternating gradient descent performs better for learning rate ratios which reduce the adversarial nature of discretisation drift.	99
4.8	Explicit regularisation cancelling the interaction terms of discretisation drift in simultaneous gradient descent improves performance in zero-games.	100
4.9	Comparison with Symplectic Gradient Adjustment and Consensus Optimisation.	101
4.10	Explicit regularisation informed by discretisation drift for alternating updates in zero-sum games.	103
4.11	The effect of discretisation drift depends on the game: its effect is less strong for the non-saturating loss.	104
4.12	An alternative interpretation to different learning rates leads to a new set of flows for simultaneous gradient descent updates.	109
4.13	An alternative interpretation to different learning rates leads to a new set of flows for alternating gradient descent updates.	109

4.14	A visualisation of a simple example of the extension of the PF to two-games.	114
5.1	Visualising the standard approach to backward error analysis alongside an approach which constructs a different flow per gradient descent iteration.	121
6.1	Double descent in deep learning compared to the traditional U-shaped complexity curve.	137
6.2	The importance of critic smoothness when estimating divergences and distances.	141
6.3	Decision surfaces on <i>two moons</i> under different regularisation methods.	143
6.4	Lipschitz constant of each layer of an MLP trained on the two moons dataset.	144
6.5	The effect of regularisation on <i>local</i> smoothness.	144
6.6	Smoothness constraints interact with learning rates.	146
6.7	Smoothness constraints interact with momentum decay rates.	147
7.1	Adding Spectral Normalisation to C51 improves the agent’s performance on Atari.	157
7.2	Selectively applying Spectral Normalisation to one layer leads to significant improvements over a DQN agent trained with Adam.	158
7.3	Spectral Normalisation shows gains for all model sizes when applied to a DQN agent on MinAtar.	159
7.4	Spectral Normalisation can significantly decrease hyperparameter sensitivity of DQN when the Adam optimiser is used.	161
7.5	The effect of DIVOUT, DIVGRAD and MULEPS on RL performance on MinAtar.	163
7.6	Optimisation changes based on Spectral Normalisation do not recover the performance of Spectral Normalisation in GANs.	167
7.7	The effect of Spectral Normalisation on discriminator predictions.	168

7.8	Spectral Normalisation increases the entropy of the discriminator’s predictions on data and samples throughout training.	169
7.9	Spectral Normalisation increases the entropy of the discriminator’s predictions on data and samples in the early stages of training.	169
7.10	Spectral Normalisation increases the entropy of the discriminator’s predictions on data and samples in the late stages of training.	170
7.11	Understanding the landscape of GAN losses and gradients: comparing the saturating and non-saturating generator loss.	171
7.12	Only applying the Jacobian change to the discriminator only leads to extremely poor performance when the saturating loss is used.	172
7.13	DIVOUT obtains similar performance to Spectral Normalisation for GANs on CIFAR-10 when the non-saturating loss is used.	173
7.14	The effect of the learning rate and ϵ on the magnitude of the update in the low-curvature regime in Adam.	175
7.15	The effect of ϵ on Adam performance in GAN training.	176
8.1	Neural networks initialise close to the 0 constant function.	189
8.2	Geometric Complexity (GC) at initialisation decreases as the number of layers in the network increases.	190
8.3	Geometric Complexity (GC) increases as training progresses.	191
8.4	Common explicit regularisation methods lead to decreased GC when using vanilla stochastic gradient descent.	192
8.5	Explicit minimising GC leads to increased generalisation.	193
8.6	The implicit regularisation effect of large learning rates on GC.	194
8.7	The implicit regularisation effect of small batch sizes on GC.	195
8.8	Double descent and GC.	197
A.1	Backward error analysis proof structure.	243
A.2	The effects of the approximation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ on DAL; results on full-batch training on CIFAR-10.	266

A.3	The effects of the approximation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ on DAL; results with batch size 512 on CIFAR-10.	266
A.4	The effects of the approximation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ on DAL; Imagenet results.	267
A.5	The effects of the approximation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ on DAL-0.5; Imagenet results.	267
A.6	Continuous time flows approximating gradient descent for a 1-D non quadratic function.	268
A.7	Global parameter error between gradient descent and continuous-time flows on an MLP trained on the UCI breast cancer dataset.	269
A.8	The eigenspectrum obtained along the gradient descent path of a 5 layer small MLP trained on the UCI Iris dataset.	269
A.9	One eigendirection is sufficient to cause instability in continuous-time; MNIST results.	270
A.10	Peak areas of loss increase and the stability coefficient for the largest eigendirection of the PF.	271
A.11	Understanding changes to λ_0 using the PF.	271
A.12	The connection between the PF, λ_0 and loss instabilities with Resnet-18 trained on CIFAR-10.	272
A.13	Assessing whether the value of the DAL learning rate for a fixed learning rate sweep has magnitudes which can be used as a learning rate in training.	272
A.14	DAL- p sweep with a VGG model trained on CIFAR-10 for small and large batch sizes.	273
A.15	Models trained using a learning rate sweep or DAL on CIFAR-10 using a Resnet-18 model.	274
A.16	DAL on Imagenet, a learning rate sweep across batch sizes.	274
A.17	DAL-0.5: on Imagenet, a learning rate sweep across batch sizes.	275
A.18	DAL- p sweep on VGG, Resnet-18 and Imagenet.	276
A.19	DAL results with a least square loss.	276

A.20	DAL and gradient descent learned landscapes; full batch on CIFAR-10.	277
A.21	λ_0 in DAL and gradient descent; full batch on CIFAR-10.	278
A.22	DAL and gradient descent learned landscapes; batch size 8192 on ImageNet.	278
A.23	VGG trained on CIFAR-10 with batch size 1024: connection between drift, test error and λ_0	279
A.24	DAL- p with momentum 0.9 on Imagenet. The model is Resnet-50 trained with batch size 1024.	280
A.25	The global error in parameter space between the NGF and gradient descent; MNIST results.	280
B.1	The effect of discretisation drift on zero-sum games; using the saturat- ing generator loss.	318
B.2	The effect of discretisation drift on zero-sum games; least square losses.	319
B.3	The effect of discretisation drift depends on the game.	320
B.4	Using explicit regularisation to cancel the interaction terms of discreti- sation drift leads to a substantial performance improvement.	321
B.5	Using explicit regularisation to cancel the interaction terms of discreti- sation drift obtains the same peak performance as Adam.	322
B.6	The form of discretisation drift allows us to construct efficient regu- larisers without the need of a hyperparameter sweep.	323
B.7	Cancelling interaction terms increases performance across multiple percentages used to select top models.	323
B.8	Cancelling interaction terms increases performance across batch sizes.	324
B.9	Comparing cancelling interactions terms with Symplectic Gradient Adjustment.	324
B.10	Comparing hyperparameter sensitivity of cancelling interactions terms and Symplectic Gradient Adjustment.	325
B.11	Comparing cancelling interactions terms with Consensus Optimisation.	325
B.12	Comparing hyperparameter sensitivity of cancelling interactions terms and Consensus Optimisation.	326

B.13	Variability across seeds when cancelling interaction terms.	326
B.14	Gradient clipping can reduce variability across seeds when cancelling interaction terms.	327
B.15	Cancelling interaction terms of discretisation drift in alternating gra- dient descent updates.	327
B.16	Variability when cancelling interaction terms of discretisation drift in alternating gradient descent updates.	328
B.17	Cancelling interaction terms in alternating updates can increase per- formance across multiple percentages used to select top models. . . .	328
E.1	Applying Spectral Normalisation on a subset of layers does not always result in smoother networks.	346
E.2	Other normalisation and regularisation methods do not recover Spec- tral Normalisation performance when applied to the DQN critic. . . .	347
E.3	Assessing the effect of the Lipschitz constant K when applying Spectral Normalisation to a DQN agent.	347
E.4	Spectral radii in a long training run on MinAtar using a 4-layer critic.	348
E.5	MinAtar Normalised Scores for the four architectures in Table E.6. . .	349
F.1	The implicit regularisation effect of increased learning rates and de- creased batch sizes leads to decreased Geometric Complexity (GC) when momentum is used.	360
F.2	Gradient norm regularisation leads to decreased GC when using stochastic gradient descent.	361
F.3	Spectral regularisation leads to decreased Geometric Complexity when using stochastic gradient descent.	362
F.4	L_2 regularisation leads to decreased GC when using stochastic gradient descent.	363

List of Tables

- 3.1 Writing the continuous-time flows discussed in terms of the eigen-decomposition of the Hessian. 44

- 4.1 Contrasting the implicit regularisation effects of the continuous-time flows proposed in zero-sum games. 108
- 4.2 Comparing cancelling discretisation drift terms with existing explicit regularisation methods in zero-sum games. 116

- 7.1 The percentage of training settings for which applying Spectral Normalisation to more layers leads to smoother networks, as measured by the model’s Jacobian at peak performance. 160
- 7.2 Spearman Correlation coefficient between the negative norm of the Jacobian and peak performance for each game. 160

- E.1 MinAtar default hyperparameter settings, used for all MinAtar experiments unless otherwise specified. 352
- E.2 MinAtar maximum and random scores used for computing the MinAtar Normalised Score. 353
- E.3 Network used for MinAtar experiments. The last fully connected layer has an output size given by the number of actions. 353
- E.4 DQN-Adam hyperparameters. 354
- E.5 Network used for Atari experiments. 354
- E.6 Details of neural architecture sweeps. 355

Chapter 1

Introduction

Why has deep learning been so successful in recent decades at solving a myriad of problems, such as scientific discovery and applications [1–3], language and image generation [4–7], and super level human game play in intricate games [8–10]? A common answer to the question of why deep learning has been so successful is that it scales well with increases in amount of data and computational resources, both of which have drastically increased in recent years. But this answer raises further questions, since many of the reasons why deep learning scales well have yet to be fully uncovered. Particularly, we still do not know why the recipe behind deep learning—using simple gradient-based algorithms to optimise deep neural networks—is so effective on large datasets. One valid response to this observation is to state that it might not matter *why* deep learning works, but *that* it does and to continue exploring its benefits on a wide range of domains, as well as continue the trend of scaling up data, architectures, and compute. This approach works, can lead to progress, and has been widely used in other sciences. This is not the approach we take here, however. Throughout this thesis, we take the point of view that understanding why a system works is not only desirable, but necessary in order to ensure targeted research progress, by exploring the areas most likely to lead to an acceleration of new methods, increased resource efficiency, and safe deployment. In taking this point of view, we follow the footsteps of a growing body of work aiming to uncover, examine, and understand deep learning systems [11–21].

There are many avenues worthy of study inside the deep learning domain, ranging from investigating why certain model architectures work well [22–25], to

studying why certain types of generative models or learning principles outperform others [26–28]. Here, we choose to focus on the study of optimisation and model regularisation in the form of smoothness regularisation. We are motivated by their generality and applicability across deep learning domains from supervised learning, probabilistic modelling, and reinforcement learning; their significant impact on training performance and test set and out of distribution generalisation [29–31]; and the opportunity for analysis due to the vast gap between their aforementioned practical performance and the understanding of their mechanisms. Unanswered questions include: What causes instability in deep learning? How can instability be mitigated? Why is there not more instability? What are the implicit regularisation effects of popular optimisers? What are the effects of smoothness regularisation? What are the interactions between model regularisation and optimisation in deep learning? Which optimisation and regularisation effects are domain specific and which transfer across deep learning domains? These are some of the questions we will tackle in this thesis.

A motivating example. Generative Adversarial Networks (GANs) [32] are a type of generative model that has led to many image generation breakthroughs [7, 33, 34]. Many GAN variants have been proposed, some based on minimising different distributional divergences or distances [27, 35–38] while others use convergence analysis around a Nash equilibrium [15, 17]; each of these variants are theoretically appealing and have different properties worth studying. Yet, the biggest successes in GANs have consistently come from changes to optimisation and regularisation: the choice of optimiser [31], large batch sizes [33], and incorporating smoothness regularisation [30] are the key ingredients of all successful GANs, while the loss functions used are often those closest to the original formulation [30, 33, 39, 40]. Similar optimisation and regularisation techniques have also allowed the expansion of GANs from continuous input domains, such as images, to discrete input domains, such as text [41]. The GAN example showcases what opportunities are available by understanding and improving optimisation and smoothness regularisation, opportunities we aim to explore here.

1.1 Optimisation in deep learning

Machine learning is the science of casting real world problems such as prediction, generation, and action into optimisation problems. Thus, many machine learning problems can be formulated as

$$\min_{\boldsymbol{\theta}} E(\boldsymbol{\theta}), \quad (1.1)$$

where E is a loss function suitably chosen for the underlying task and $\boldsymbol{\theta} \in \mathbb{R}^D$. For many machine learning tasks, E is an expected value of an unknown data distribution $p^*(\mathbf{x})$; E gets estimated given a training dataset \mathcal{D} of unbiased samples from $p^*(\mathbf{x})$: $E(\boldsymbol{\theta}) = \mathbb{E}_{p^*(\mathbf{x})} E(\boldsymbol{\theta}; \mathbf{x}) \approx \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x}_i \in \mathcal{D}} E(\boldsymbol{\theta}; \mathbf{x}_i)$. A distinction is made between the model and the objective function components of $E(\boldsymbol{\theta}; \mathbf{x})$, with the model $f(\boldsymbol{\theta}; \mathbf{x})$ depending on the parameters $\boldsymbol{\theta}$, and the objective function depending only on the model's output. In deep learning, f is a neural network, $\boldsymbol{\theta}$ are its parameters, and E is non-convex. Due to the high dimensionality of $\boldsymbol{\theta}$ or large dataset size, many optimisation algorithms are too costly or inefficient to be used in deep learning, as they scale unfavourably with parameter or dataset size. Luckily, however, we can use the compositional structure of neural networks to compute the gradient $\nabla_{\boldsymbol{\theta}} E \in \mathbb{R}^D$ using the backpropagation algorithm [42–44]. From there, we can use first-order gradient-based iterative algorithms, such as gradient descent and its variants [45–48].

Studying the behaviour of optimisation algorithms can be challenging, even for algorithms as seemingly simple as gradient descent. Typical questions of study include: Will the algorithm converge to a local minima? Does the algorithm get stuck in saddle points? When does the algorithm fail to converge and when it does converge, how quickly does it converge? The complexity of these issues gets compounded when studying the optimisation of deep neural networks with millions to billions of parameters that do not satisfy commonly used assumptions such as convexity; this often creates a discrepancy between what theory analyses and what is used in practice. Fundamentally, the deep learning optimisation community is interested in answering the questions: *Why do first-order optimisers work so well in training neural networks* and *How can we improve optimisation in deep learning?*

Progress in this area has been made recently both through theoretical and empirical means of studying neural network optimisation dynamics [18, 49–53]. Novel insights challenge commonly held assumptions, such as the belief that local minima and saddle points are a major challenge for neural network optimisation due to high parameter dimensionality [49–51, 54]. There has also been a growing body of work showing the importance of the Hessian $\nabla_{\theta}^2 E$ in supervised learning optimisation, and specifically on the connection between the learning rate and the largest Hessian eigenvalue and observed instabilities in training [18, 52, 53]; we will examine the importance of the Hessian in a new light in Chapter 3.

Many analyses of discrete-time methods, such as gradient descent, take a continuous-time approach [11, 17, 54–61] as continuous-time proofs tend to be easier to construct [54, 62]. The downside of taking a continuous-time and not discrete-time approach to optimisation analysis is that one does not directly analyse the update of interest; the discrepancy between the discrete-time trajectory and its continuous counterpart, which we will call discretisation drift, can lead to conclusions that do not transfer from continuous-time analysis to discrete-time algorithms [63, 64]. Discretisation drift is still not greatly understood and has only recently come into attention in deep learning [12, 13, 65]. The incorporation of techniques from the numerical integration community that construct continuous-time flows accounting for discretisation drift has shed light on the implicit regularisation effects of gradient descent and highlighted the importance of the learning rate in generalisation [12, 13]. We will use the same numerical integration techniques in Chapter 3, where we find a new continuous-time flow that, unlike existing flows, can capture instabilities induced by gradient descent. We then use insights from our continuous-time analysis to devise an automatic learning rate schedule that can trade-off training stability and generalisation performance.

While many studies focus on optimisation dynamics in the single-objective case, such as supervised learning, strides have also been made towards understanding optimisation in two-player games [15, 17, 61, 66, 67]. Unlike in the single-objective setting, in two-player games not all parameters minimise the same objective, as each

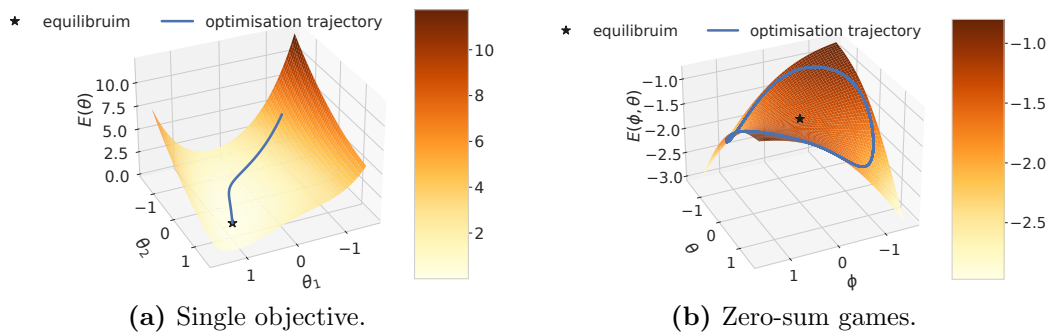


Figure 1.1: Overview of optimisation in deep learning. Many optimisation problems in deep learning can be seen as either single-objective problems, where all parameters are updated to minimise the same loss, or multi objective problems where different sets of parameters are updated to minimise different objectives; when the number of objectives is two, these are called two-player games. A special case of two-player games are adversarial games such as zero-sum games, where the objective of one player is to minimise a function E , while the objective of the other player is to maximise the same function. This adversarial structure can lead to challenges in optimisation, with approaches that lead to converge to local minima in single objective optimisation (a) resulting in cyclic trajectories in zero sum games (b).

player has its own set of parameters, here denoted by ϕ and θ , and its own objective:

$$\min_{\phi} E_{\phi}(\phi, \theta) \quad (1.2)$$

$$\min_{\theta} E_{\theta}(\phi, \theta). \quad (1.3)$$

Two-player games can be adversarial, where the two players have opposing objectives, and even zero sum, where one player's gain is another's loss: $E_{\phi}(\theta, \phi) = -E_{\theta}(\theta, \phi)$. We visualise a simple two-player zero-sum game in Figure 1.1b, showing that adversarial dynamics can make it challenging to reach a local equilibrium compared to a single-objective counterpart, shown in Figure 1.1a.

Interest in the analysis of two-player game optimisation has been fuelled by interest in GANs [32], a generative model trained via an adversarial two player game. While GANs have been a successful generative model, they have a reputation of being notoriously hard to train, and this has served as a motivation for finding approaches to stabilise and improve the dynamics of adversarial two-player games. Analysing the behaviour and convergence of two-player games has led to progress

in understanding the sources of instability and divergence in games and has led to many a regulariser that can stabilise training [17, 61, 67–70]. We add to this body of work by further investigating sources of instability in two-player games, coming from discretisation drift rather than the original continuous-time system, in Chapter 4. By modelling discretisation drift, we devise explicit regularisers that stabilise training without requiring any additional hyperparameter sweep.

1.2 Smoothness w.r.t. inputs

Optimisation provides an approach to approximate an unknown decision surface required for prediction, generation, or action from an available dataset. That is, given an optimal predictor f^* , optimisation provides a recipe to find a parametrised function $f(\cdot; \boldsymbol{\theta})$ close to f^* on the training dataset. But there are often many such optimisation solutions, each with different properties outside the training data. How useful this approximation is in the underlying application depends on how well it performs on *new data*. This is an essential desideratum of machine learning algorithms, and it is relevant for generalisation, continual learning [71], and transfer learning [72]. Since we know from the no-free-lunch theorem [73] that no machine learning algorithm will perform better than all other algorithms for any data distribution and prediction problem, to provide guarantees beyond the training data requires making assumptions about the data distribution, the structure of f^* , or both. One approach to encode beliefs about f^* into $f(\cdot; \boldsymbol{\theta})$ is through what is often referred to in machine learning as an *inductive bias* [74]¹.

In accordance with Occam’s razor [75], we often want to encode a simplicity inductive bias: we are searching for the least complex functions that can explain the data. One approach to formalising simplicity is by measuring the effect changes in function input have on the function’s output. If small changes in function input do not result in large changes in function output, we call that function *smooth*. A smoothness inductive bias encodes a preference for learning smooth functions:

¹When we talk about smoothness with respect to inputs, we often write the predictor as a function of data, given parameters. When discussing optimisation and the focus is on parameters, we write the reverse for clarity.

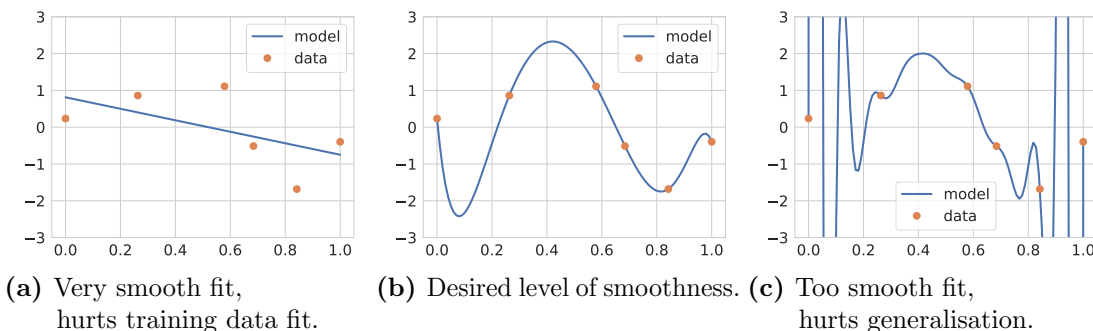


Figure 1.2: The importance of smoothness with respect to inputs on the model fit: too much smoothness can hurt capacity (a), while not enough smoothness can hurt generalisation (c).

between two functions that fit the data equally well, the smoothest one should be preferred. We show an intuitive example highlighting the importance of smoothness in Figure 1.2. Figure 1.2a exemplifies how too much smoothness can hurt training performance by reducing model capacity, while Figure 1.2c shows how the lack of smoothness can lead to overfitting.

The formal definition of smoothness is often taken to be Lipschitz smoothness or related measures. A function $f(\cdot; \theta) : \mathcal{X} \rightarrow \mathcal{Y}$ is K -Lipschitz if

$$\|f(\mathbf{x}_1; \theta) - f(\mathbf{x}_2; \theta)\|_{\mathcal{Y}} \leq K \|\mathbf{x}_1 - \mathbf{x}_2\|_{\mathcal{X}} \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \quad (1.4)$$

where K is known as the Lipschitz constant of function $f(\cdot; \theta)$. Throughout this manuscript, when we refer to the smoothness of a neural network, we refer to the smoothness w.r.t. data \mathbf{x} , and not w.r.t. parameters θ .

For many a machine learning task smoothness with respect to inputs constitutes a reasonable prior; this can be easily illustrated in the image domain, where the prediction or action that ought to be taken does not change when a few pixels in the inputs change: an image of a dog still encodes a dog if a few pixels are modified. Thus, it does not come as a surprise that despite having vastly different motivations and implementations, many existing deep learning regularisation methods target a smoothness inductive bias, including L_2 regularisation, dropout [76], early stopping. We will discuss in detail how these regularisation approaches are connected to

smoothness in Chapter 6, formalise these connections in Chapter 8, as well as link smoothness to important phenomena in deep learning, such as double descent [77–79].

As the importance of smoothness is beginning to be recognised by the deep learning community, methods that directly target smoothness regularisation have been developed [29, 30, 80]. These regularisation methods have been incorporated in supervised learning to help generalisation [29, 81] and adversarial robustness [82–84], as well as into generative models leading to improvements in performance and stability [30, 39, 85]. We investigate the importance of smoothness regularisation and its effects in Chapter 6 as well as incorporate it into a new problem domain, reinforcement learning, in Chapter 7.

1.3 Interactions between optimisation and smoothness regularisation

We have thus far discussed optimisation and smoothness regularisation as separate pillars of a deep learning system. Optimisation—concerned with changes in *parameters*—affects training performance, speed, and stability. Smoothness—concerned with changes in *inputs*—can be a powerful inductive bias avoiding learning overly complex models. What is often overlooked is the *implicit bias optimisation has on the model class being learned*. While neural architectures and model regularisers define the class of functions we can represent with our models, optimisation methods determine *which* functions in the function class we can learn [12]. Likewise, different model architectures are easier to learn than others, with ResNets [86] providing a prime example of an architecture that aids optimisation [22], followed by model normalisation techniques such as batch normalisation [87], which have been shown to primarily benefit optimisation [23].

While not previously studied, we show that smoothness regularisation and optimisation heavily interact in deep learning. We start by showing examples of this interaction in Chapter 6, where we observe that smoothness regularisation methods can interact with hyperparameter choices such as learning rates. We then highlight how in some problem domains smoothness regularisation leads to increased

performance by improving optimisation dynamics in Chapter 7. In Chapter 8, we show how the implicit regularisation induced by the discretisation drift of an optimisation method can lead to a smoothness inductive bias, and that the strength of this inductive bias is dependent on optimisation hyperparameters such as learning rate and batch size.

Throughout this thesis, we will see examples of how optimisation can affect the learned model class and how the model smoothness can affect optimisation. While it is tempting to think of each component of a deep learning system individually, it can hinder progress and misguide our intuitions. We will use these observations to argue for a cohesive view between optimisation and regularisation in deep learning.

Relationship to published papers

The presented thesis is based on the following published works:

- *On a continuous-time model of gradient descent dynamics and instability in deep learning.* **Mihaela Rosca**, Yan Wu, Chongli Qin, Benoit Dherin. Transactions of Machine Learning Research, 2023.
- *Discretization Drift in Two-Player Games.* **Mihaela Rosca**, Yan Wu, Benoit Dherin and David G. T. Barrett. International Conference for Machine Learning, 2021.
- *A case for new neural network smoothness constraints.* **Mihaela Rosca**, Theophane Weber, Arthur Gretton, Shakir Mohamed. 1st I Can't Believe It's Not Better Workshop (ICBINB@NeurIPS 2020), Vancouver, Canada.
- *Spectral Normalisation for Deep Reinforcement Learning: an Optimisation Perspective.* Florin Gogianu, Tudor Berariu, **Mihaela Rosca**, Claudia Clopath, Lucian Busoniu, Razvan Pascanu. International Conference for Machine Learning, 2021.
- *Why neural networks find simple solutions: The many regularizers of geometric complexity.* Benoit Dherin, Michael Munn, **Mihaela Rosca**, David G.T. Barrett. Neural Information Processing Systems, 2022.

Other work performed as part of the PhD program that is not discussed here:

- Chapter on implicit generative models in the second edition of ‘Machine Learning: a Probabilistic Perspective’ by Kevin P. Murphy. **Mihaela Rosca**, Balaji Lakshminarayanan and Shakir Mohamed.
- *Monte Carlo gradient estimation in machine learning.* Shakir Mohamed, **Mihaela Rosca**, Michael Figurnov, Andriy Mnih. Equal contributions. Journal of Machine Learning Research, 21(132), pp.1–62. 2020.
- Proposed and advised a project for a master student at UCL. The topic of study was the optimisation landscape in offline reinforcement learning.
- Co-organised the ICML 2021 workshop on ‘Continuous Time Perspectives in Machine Learning’ <https://icml.cc/virtual/2022/workshop/13452>.

Overview of thesis

- Chapter 2 is an overview of the optimisation concepts required for the rest of this thesis.
- Chapter 3 consists of the work published in *On a continuous-time model of gradient descent dynamics and instability in deep learning* [88]. My contributions to this paper were: I started the project and investigated the idea of finding a new continuous-time flow for gradient descent; did all the proofs (except the justification of Thm 3.3.2); I wrote all the code, ran all the experiments, and did all the visualisations; I wrote the paper.
- Chapter 4 consists of the work published in *Discretization Drift in Two-Player Games* [89], together with novel insights presented in this thesis regarding the use of different learning rates for the two-players (Section 4.7) as well as expanding the flow derived in Chapter 3 to games (Section 4.8). My contributions to the paper were: I initiated this project by suggesting to expand the work of Barrett and Dherin [12] to two-player games; derived the main results; did the relevant literature survey; I wrote all the code, ran all the experiments and did all the visualisations; and I wrote most of the paper.
- Chapter 5 consists of unpublished work aiming to develop a new framework for constructing modified flows and implicit regularisers in deep learning.
- Chapter 6 consists of the work published in the workshop paper *A case for new neural network smoothness constraints* [90]. This work is the result of the early reading I have done as part of the PhD program. I wrote the code and performed all the experiments for this work, and wrote the paper.
- Chapter 7 is based on the paper *Spectral Normalisation for Deep Reinforcement Learning: an Optimisation Perspective* [91] together with the application of the insights from the paper applied to GANs (Section 7.4), which is unpublished novel work I did while writing this thesis. The paper contents were rewritten and the text is mine. My contributions to the paper were: I added to the view that Spectral Normalisation helps RL through optimisation, helped derive the

optimisers that allow us to recover the performance of Spectral Normalisation, and suggested experiments and analysis such as plotting the change in Jacobian of the critic against performance as well as provide intuition and advice on the application on Spectral Normalisation and other smoothness regularisation methods; I did not contribute with any coding or experiments to the paper.

- Chapter 8 is based on *Why neural networks find simple solutions: The many regularizers of Geometric Complexity* [92]. The paper contents were rewritten and the text is mine. My contributions to the paper were: I provided a broader context and related work on smoothness regularisation in deep learning; I suggested, analysed, and visualised the double descent experiments that led to showing that Geometric Complexity (GC) follows the traditional U-shape curve associated with complexity measures; provided coding and debugging help throughout the project, which enabled the use of GC explicit regularisation, improved evaluation and training. I also contributed to writing and the rebuttal. I did not run experiments for this project, but provided guidance for what experiments to run. I did not derive the main results in the paper, but was involved in related discussions throughout the project.

Chapter 2

Optimisation in deep learning

In this chapter we describe the optimisation tools and methods we will require in the rest of this thesis. When describing the commonly used algorithms in deep learning we often take a continuous-time approach. Compared to a discrete-time approach, a continuous-time approach tends to be more amenable to proof construction [54], having many tools available, including stability analysis, and can be used to construct conserved quantities [65]. Indeed, the ease of analysis of continuous-time approaches has led to a general uptake into continuous-time methods in deep learning recently, with use in generative models, architectures, and optimisation [67, 93–97].

2.1 The direction of steepest descent

Once we have formulated our machine learning problem into an optimisation problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^D} E(\boldsymbol{\theta}), \quad (2.1)$$

we have to choose an optimisation algorithm with which to solve the above problem. Solving the problem in Eq (2.1) in closed form is intractable for the problems we are interested in. Thus, the optimisation algorithms used are iterative, and aim to answer the question: given a current set of parameters $\boldsymbol{\theta}$, what is a direction of descent of function E ? The simplest choice of the direction is given by the negative gradient $-\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$, known as *the direction of steepest descent*. To see why, consider the continuous-time flow, often referred to as the *gradient flow* or *negative gradient*

flow (NGF)

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E. \quad (2.2)$$

If parameters $\boldsymbol{\theta}$ follow the NGF in continuous-time, $E(\boldsymbol{\theta})$ will decrease until a stationary point $\nabla_{\boldsymbol{\theta}} E = \mathbf{0}$ is reached, since

$$\frac{dE}{dt} = \frac{d\boldsymbol{\theta}^T}{dt} \nabla_{\boldsymbol{\theta}} E = -(\nabla_{\boldsymbol{\theta}} E)^T \nabla_{\boldsymbol{\theta}} E = -\|\nabla_{\boldsymbol{\theta}} E\|^2 \leq 0. \quad (2.3)$$

If E is convex, following the NGF will find the global minimum. If E is not convex, following the NGF in continuous-time will reach a local minimum (a proof will be provided in Corollary 2.1). A point $\boldsymbol{\theta}^* \in \mathbb{R}^D$ is a local minimum if there is a neighbourhood \mathcal{V} around $\boldsymbol{\theta}^*$ such that $\forall \boldsymbol{\theta}' \in \mathcal{V} \quad E(\boldsymbol{\theta}') \geq E(\boldsymbol{\theta}^*)$; a more amenable characterisation is that local minima are stationary points $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*) = \mathbf{0}$ with a positive semi-definite Hessian, i.e. $\nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*)$ has only non-negative eigenvalues. It is often convenient to also define strict local minima, i.e. local minima where eigenvalues are strictly positive.

We have already found two insights that would follow us throughout the rest of this thesis: the ease of analysis in continuous-time, and the importance of the negative gradient. Before exploring the optimisation algorithms using these insights, we briefly turn to the question of *computing* the gradient.

2.1.1 Computing the gradient

To compute the gradient $\nabla_{\boldsymbol{\theta}} E$ we have to consider the specification of the loss function E . In machine learning, E often takes the form of an expectation or sum of expectations; if the expectation is under the data distribution $p^*(\mathbf{x})$ and we can write $E(\boldsymbol{\theta}) = \mathbb{E}_{p^*(\mathbf{x})} l(\boldsymbol{\theta}; \mathbf{x})$, we have

$$\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{p^*(\mathbf{x})} l(\boldsymbol{\theta}; \mathbf{x}) = \mathbb{E}_{p^*(\mathbf{x})} \nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}; \mathbf{x}). \quad (2.4)$$

We can now compute an unbiased estimate of the gradient via Monte Carlo estimation:

$$\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) = \mathbb{E}_{p^*(\mathbf{x})} \nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}; \mathbf{x}) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}; \mathbf{x}_i), \quad \mathbf{x}_i \sim p^*(\mathbf{x}). \quad (2.5)$$

In deep learning, the gradient $\nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}; \mathbf{x}_i)$ is often computed via backpropagation, i.e. using the compositional structure of neural networks to compute gradients in a recursive fashion via the chain rule. If the entire dataset is used to approximate the gradient in Eq (2.5), this corresponds to *full-batch training*; if an unbiased sample from the dataset is used, the training procedure is referred to as *mini-batch training*.

Additional challenges arise for objectives of the form $E(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} l(\mathbf{x})$, as is the case for some reinforcement learning and generative modelling objectives. In such situations, the change of expectation and gradient in Eq (2.4) is no longer possible. For certain family of distributions $p(\mathbf{x}; \boldsymbol{\theta})$, however, gradient estimators are available. In later chapters we will use the pathwise estimator, where the random variable X with distribution $p(\mathbf{x}; \boldsymbol{\theta})$ can be written as $X = g(Z; \boldsymbol{\theta})$, in which case we can use the change of variable formula for probability distributions:

$$\mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} l(\mathbf{x}) = \mathbb{E}_{p(\mathbf{z})} l(g(Z; \boldsymbol{\theta})). \quad (2.6)$$

From here we can write

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} l(\mathbf{x}) = \nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{z})} l(g(Z; \boldsymbol{\theta})) = \mathbb{E}_{p(\mathbf{z})} \nabla_{\boldsymbol{\theta}} l(g(Z; \boldsymbol{\theta})), \quad (2.7)$$

which is again amenable to Monte Carlo estimation. We refer to our general overview for more details on the pathwise estimator as well as other estimators [98].

2.2 Algorithms: from gradient descent to Adam

We have seen how following the negative gradient in continuous-time leads to convergence to a local minimum. Implementing infinitely small updates in continuous time on our discrete computers is not feasible, however, and thus we need a discrete-time algorithm for optimisation. Discretising continuous-time flows is a heavily studied topic with its own area of research in applied mathematics, numerical integration [99]. Numerical integrators provide approaches to approximate the solution of the flow at time h and initial conditions $\boldsymbol{\theta}(0)$, which we denote as $\boldsymbol{\theta}(h; \boldsymbol{\theta}(0))$.

A simple discretisation method is Euler discretisation, which when applied to the NGF leads to

$$\boldsymbol{\theta}(h; \boldsymbol{\theta}(0)) = \boldsymbol{\theta}(0) + \int_0^h \dot{\boldsymbol{\theta}}(t) dt = \boldsymbol{\theta}(0) - \int_0^h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}(t)) dt \quad (2.8)$$

$$\approx \boldsymbol{\theta}(0) - \int_0^h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}(0)) dt = \boldsymbol{\theta}(0) - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}(0)). \quad (2.9)$$

Setting $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_{t-1}$ in the above equation leads to the familiar gradient descent update:

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}). \quad (2.10)$$

By assuming that the vector field of the flow—the gradient $\nabla_{\boldsymbol{\theta}} E$ —does not change over the time interval h we obtained a fast discrete-time algorithm. However, due to the approximation in Eq (2.9), following gradient descent is not the same as following the NGF: the numerical integration error leads to a difference between the gradient descent and NGF trajectories. Throughout this thesis we will refer to this difference in trajectories as *discretisation drift* (often also known as *discretisation error*). Due to discretisation drift, there are no guarantees that gradient descent will reach a local minimum. Unlike the NGF, gradient descent can diverge or converge to saddle points (stationary points that have at least one negative eigenvalue of the Hessian).

We visualise the effects of discretisation drift in Figure 2.1. Discretisation drift can speed up training (Figure 2.1a), destabilise it (Figure 2.1b) or lead to divergence (Figure 2.1c). Importantly, when using gradient descent, the discrete counterpart of the NGF, for large h the function E can *increase*, leading to training instabilities. Understanding *when* gradient descent leads to instabilities can be challenging, due to intricate dependencies between the learning rate and the shape of E ; this gets compounded in the case of deep learning due to the very high dimensional nature of the parameter space. We will tackle this question later in this thesis.

Recent work has made great strides in understanding the behaviour of gradient descent in deep learning. Questions of study include whether gradient descent converges to local or global minima [49, 100–102], the prevalence and effect of saddle points [51, 103], the interaction between neural architectures and optimisa-

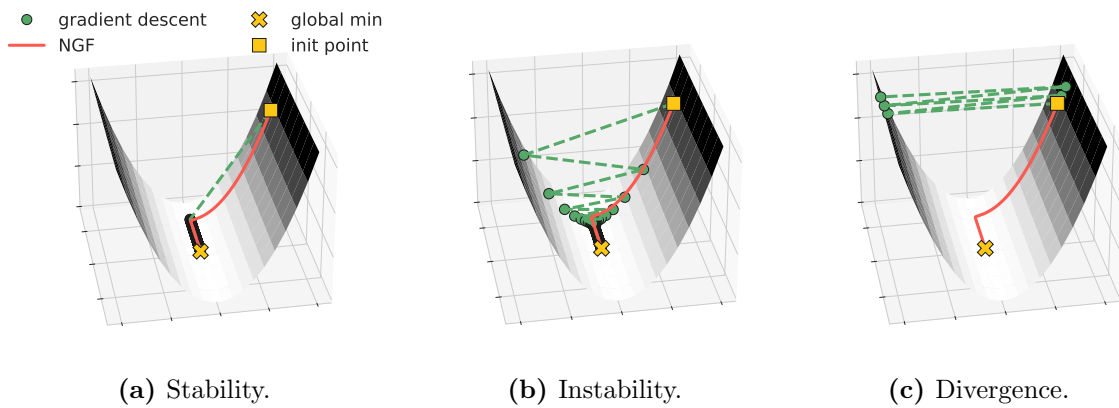


Figure 2.1: The importance of the learning rate in gradient descent: small learning rates lead to convergence to local minima but slow to converge (a), while larger learning rates can lead to instability (b) and divergence (c).

tion [22, 23, 104], the connection between learning rates and Hessian eigenvalues and instabilities [52], the effect of the learning rate on generalisation [12, 105].

Rprop and RMSProp. We have derived gradient descent as a discretisation of the NGF; we used the NGF since it decreases the value of the objective function E . We note, however, that only the sign of the gradient for each parameter is relevant in order to achieve a descent direction locally. To see why, consider constants $c_i > 0$ and the system:

$$\dot{\theta}_i = -c_i \nabla_{\theta_i} E. \quad (2.11)$$

This flow, like the NGF, also decreases E , since as $c_i > 0$

$$\frac{dE}{dt} = \sum_i \frac{dE}{d\theta_i} \frac{d\theta_i}{dt} = \sum_i \frac{dE}{d\theta_i} (-c_i \nabla_{\theta_i} E) = - \sum_i c_i (\nabla_{\theta_i} E)^2 \leq 0. \quad (2.12)$$

This observation is the motivation for Rprop [106], which instead of using the gradient as the parameter update, uses only its sign: $\text{sign}(\nabla_{\theta_i} E) = \frac{\nabla_{\theta_i} E}{\sqrt{(\nabla_{\theta_i} E)^2}}$. We can think of Rprop as the discretisation of the following flow at each iteration:

$$\dot{\theta}_i = - \frac{1}{|\nabla_{\theta_i} E(\theta_0)|} \nabla_{\theta_i} E. \quad (2.13)$$

While this can solve issues with exploding gradients and an imbalance between the magnitude of the updates of different parameters, it can be crude: Rprop cannot

distinguish between areas of space where gradient are positive, but very small or areas of space with very large gradients; the update in both of these cases would be the learning rate h . The solution proposed to this issue was RMSprop [48], which instead of normalising the gradient of each parameter by $(\nabla_{\theta_i} E)^2$, it normalises the gradient by a moving average of $(\nabla_{\theta_i} E)^2$ obtained from previous iterations. By using element-wise operations, we can write the RMSprop update as

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \nabla_{\theta} E(\boldsymbol{\theta}_{t-1})^2; \quad \boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \frac{\nabla_{\theta} E(\boldsymbol{\theta}_{t-1})}{\sqrt{\mathbf{v}_t}}, \quad (2.14)$$

where β is a hyperparameter often set to values in $[0.9, 0.999]$. Since the denominator is positive, the RMSprop update has the same sign as the gradient, but a different magnitude. The use of moving averages can dampen the effect of a large gradient, which would lead to instabilities in gradient descent, but still accounts for the the magnitude of the gradient, not only its sign, unlike Rprop. RMSprop and its variants are often used in reinforcement learning [107–109].

Momentum. Moving averages are an important staple of optimisation algorithms, beyond adjusting the magnitude of the local gradient, as we have seen with RMSprop. The idea behind momentum algorithms is to speed up or stabilise training by using previous iteration gradients.

To provide intuition regarding momentum, we note that along a trajectory towards a local minimum the sign of the gradient will be the same throughout that trajectory. Let $\boldsymbol{\theta}^*$ be a local minimum of $E(\boldsymbol{\theta})$ and $\boldsymbol{\theta}_{t+1}$ and $\boldsymbol{\theta}_t$ iterates along a trajectory towards $\boldsymbol{\theta}^*$ and consider parameter index i such that $\text{sign}(\boldsymbol{\theta}_{t+1,i} - \boldsymbol{\theta}_i^*) = \text{sign}(\boldsymbol{\theta}_{t,i} - \boldsymbol{\theta}_i^*)$; this encodes the assumption that the trajectory goes towards a $\boldsymbol{\theta}^*$ without escaping it in dimension i . But $\text{sign}(\boldsymbol{\theta}_{t+1,i} - \boldsymbol{\theta}_i^*) = \text{sign}(\nabla_{\theta} E(\boldsymbol{\theta}_t)_i)$ as they are both the direction which minimises E . Thus, $\text{sign}(\nabla_{\theta} E(\boldsymbol{\theta}_t)_i) = \text{sign}(\nabla_{\theta} E(\boldsymbol{\theta}_{t+1})_i)$. This tells us that as long as a trajectory is not jumping over a local minimum in a direction, the sign of the gradient in that direction does not change. Hence, one can speed up training by taking bigger steps in that direction by accounting for the previous updates; alternatively this is often framed as an approach to speed up training in areas of low curvature (where by definition the gradient does not change

substantially) [110]. This intuition suggests the following algorithm [111]:

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \sum_{i=1}^t \beta^{i-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-i}), \quad (2.15)$$

where the update is given by a geometrical average of the gradients at previous iterations, and β is the decay rate, a hyperparameter with $\beta < 1$. We can rewrite the momentum algorithm in its commonly known form:

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}); \quad \boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + \mathbf{v}_t. \quad (2.16)$$

Traditionally momentum has been seen as a discretisation of the second order flow

$$\ddot{\boldsymbol{\theta}} + c_v \dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}), \quad (2.17)$$

from where the vector \mathbf{v} can be introduced:

$$\dot{\mathbf{v}} = -c_v \mathbf{v} - \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}); \quad \dot{\boldsymbol{\theta}} = \mathbf{v} \quad (2.18)$$

which, when setting $c_v = \frac{1-\beta}{h}$ and Euler discretised with learning rates h and 1, leads to the discrete updates above. Momentum can also be seen as a discretisation of the first-order continuous-time flow, by combining implicit and explicit Euler discretisation [65]:

$$\dot{\boldsymbol{\theta}} = -\frac{1}{1-\beta} \nabla_{\boldsymbol{\theta}} E. \quad (2.19)$$

Since $\beta < 1 \implies \frac{1}{1-\beta} > 1$, this is in line with our previous intuition of speeding up training. As with gradient descent, however, this intuition does not account for the effect of large learning rates leading to discretisation errors; while following the flow in Eq (2.19) always decreases E , following the trajectory of the momentum optimiser instead loses that guarantee. We show the effect of β and h on the optimisation trajectory in Figure 2.2. While momentum can speed up training and reach a neighbourhood of a local minimum quicker than vanilla gradient descent, for large decay rates β it can lead to large updates that move further away from

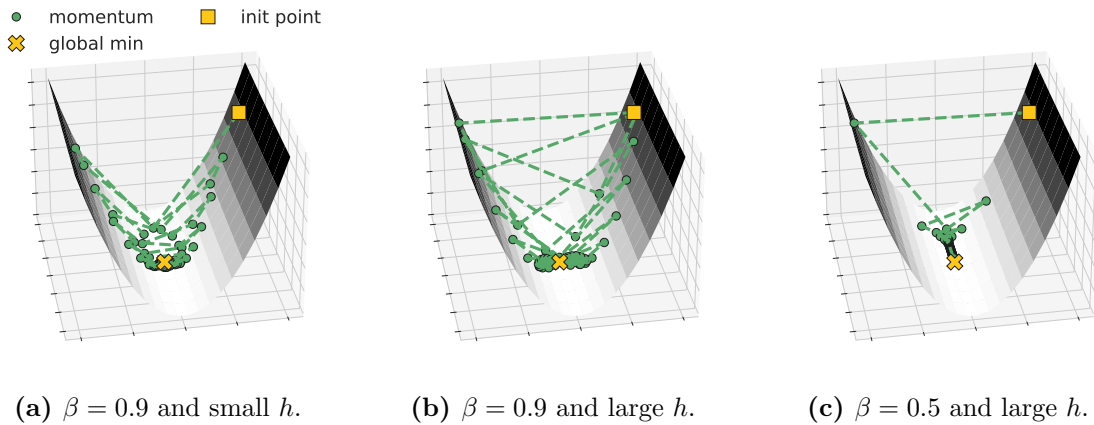


Figure 2.2: The importance of the learning rate h and decay rate β in gradient descent with momentum. (a): momentum can speed up training but can also introduce instabilities when the gradient changes direction; the learning rate used is the same as in Figure 2.1a, where gradient descent converged to the minimum without instabilities. (b) and (c): for learning rates where gradient descent diverges, as we observed in Figure 2.1c, momentum can stabilise training by dampening the effect of local large gradient magnitudes and large learning rates; from this simple example one can already observe the intricacies of the interactions of decay rates and learning rates when momentum is used.

the equilibrium. We show this effect in Figure 2.2a, where we observe that despite using a small learning rate, with a high decay rate β momentum can oscillate around the minimum significantly before reaching it, while gradient descent with the same learning rate does not oscillate around the minimum (Figure 2.1a). This effect can be mitigated by using a smaller decay rate, as exemplified in Figure 2.2c. Dampening the strength of the current gradient update through momentum can also stabilise training, as we show in Figure 2.2b. While for the same learning rate gradient descent diverges (Figure 2.1c), momentum does not, since the large local gradients moving away from the equilibrium are dampened by the moving averages obtained from past gradients. For quadratic functions, this difference between gradient descent and momentum can be seen in the smallest learning rate which leads to divergence: while gradient descent diverges if $h > 2/\lambda_0$ where λ_0 is the largest eigenvalue of the Hessian $\nabla_{\theta}^2 E$, for momentum divergence occurs when $h > 2/\lambda_0 + 2\beta/\lambda_0 \geq 2/\lambda_0$ since $\beta \geq 0$ [52]. We will come back to the importance of the learning rate relative to the Hessian eigenvalues in later chapters.

Gradient descent with momentum is still very much used for deep learning opti-

misation and has been credited with increased generalisation performance compared to other optimisers [112, 113], though recent work has shown that other well tuned methods can outperform momentum if well tuned [47].

The Adam optimiser [47], perhaps the most popular optimisation algorithm in deep learning, combines some of the previously discussed insights: it uses momentum and an RMSprop-like approach of dividing by a moving average of the square gradient (all operations below are element-wise):

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) \quad (2.20)$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})^2 \quad (2.21)$$

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \frac{\frac{1}{1-\beta_1^t} \mathbf{m}_t}{\sqrt{\frac{1}{1-\beta_2^t} \mathbf{v}_t + \epsilon}}, \quad (2.22)$$

where β_1, β_2 and ϵ are hyperparameters usually set to $\beta_2 = 0.999$ and $\beta_1 \in [0, 0.9]$ and $\epsilon \in [10^{-1}, 10^{-8}]$ with the default at 10^{-8} . The division in Eq (2.22) is obtained using a bias correction argument for the moving averages. To see this, we write:

$$\mathbf{m}_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{i-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-i}). \quad (2.23)$$

The authors make the argument that if the gradients at iteration i $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_i)$ are drawn from the distribution $p_{g,i}$

$$\mathbb{E}_{p_m}[\mathbf{m}_t] = (1 - \beta_1) \sum_{i=1}^t \beta_1^{i-1} \mathbb{E}_{p_{g,i}}[\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-i})] \quad (2.24)$$

$$= (1 - \beta_1^t) \mathbb{E}_{p_g}[\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})] + \boldsymbol{\xi}, \quad (2.25)$$

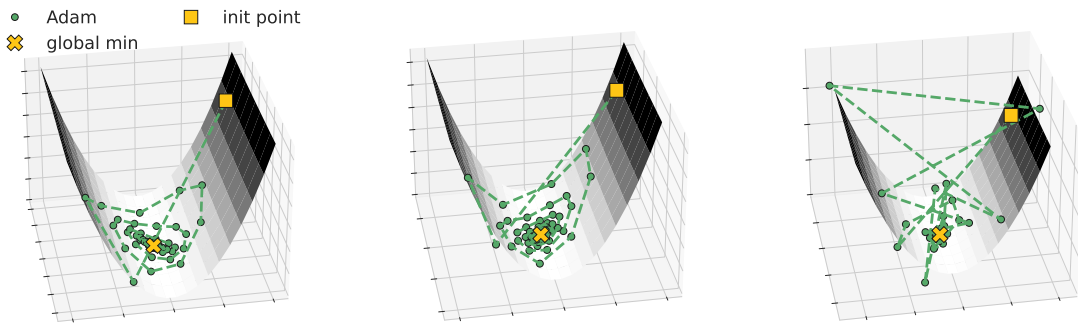
where $\boldsymbol{\xi}$ is a bias term accounting for the changes in the expected value of the gradients at different iterations; it will be $\mathbf{0}$ if the gradient mean does not change between iterations otherwise it is expected to be small due to the decaying contribution of older gradients. Thus, to ensure the expected value of the moving average stays the same as the gradient and does not decrease steadily as training progresses, the

division in Eq (2.22) is used; a similar argument can be made for the second moment. This ensures that in expectation the update (if a small ϵ is used), is approximately of magnitude h for each parameter. We note that while the assumption that the bias ξ is small does not hold for certain areas of training, such as when the gradient changes direction due to instability around an equilibrium, the above argument holds in areas of low curvature (i.e. areas of space where the gradient does not change substantially); this is something we will use later.

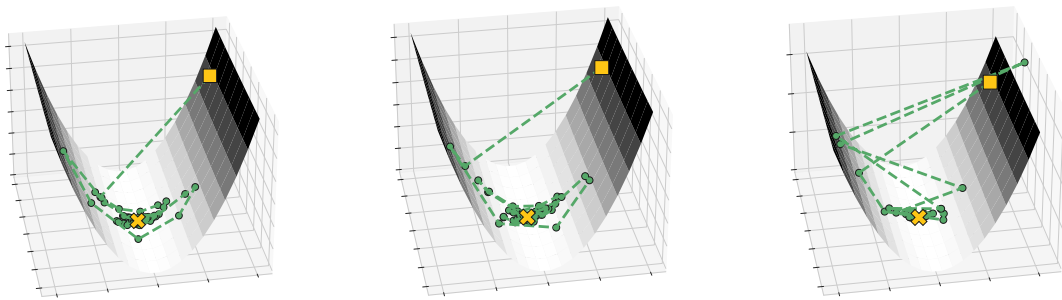
By investigating the commonly used Adam hyperparameters, $\beta_1 \in [0, 0.9]$ and $\beta_2 = 0.999$, we get intuition about its behaviour. Since β_1 is smaller than β_2 , Adam is more sensitive to the sign of the gradient (controlled by β_1) than its magnitude (controlled by β_2). Figures 2.3a and 2.3b show that the exponential moving averages in the denominator of the Adam update can dampen large updates compared to momentum (the corresponding momentum figures are shown in Figure 2.2a and 2.2b), but oscillations around the minimum still occur. When a low decay rate for the nominator is used, such as $\beta_1 = 0.5$, in conjunction with a large h , instabilities can occur if small gradients are followed by a large gradient, since \mathbf{m}_t will be large but \mathbf{v}_t will remain small due to $\beta_2 = 0.999$. We see this in Figure 2.3c, which we can contrast with momentum for $\beta = \beta_1$ and the same learning rate h in Figure 2.2c: momentum exhibits stable training, and does not exit the area around the equilibrium. Relatedly, the use of moving averages in the denominator \mathbf{v}_t has been linked to the lack of convergence in Adam for simple online convex optimisation problems [114].

The hyperparameter ϵ does not change the sign of the update, but it can change its magnitude. Specifically, a high ϵ decreases the magnitude of the parameter update. We show the effect of ϵ in Figure 2.3, where we show that using $\epsilon = 10^{-2}$ (bottom row) can lead to more stable trajectories compared to $\epsilon = 10^{-8}$ (top row).

Adam is a very popular the optimiser in deep learning, and has been attributed with progress in various fields including GANs [31], and training modern models such as transformers [115, 116]. The ϵ hyperparameter in Adam has been noted to be important for certain machine learning applications [116], including reinforcement learning, where the ϵ used is large compared to supervised learning [109, 117].



(a) $\beta_1 = 0.9$, small h , $\epsilon = 10^{-8}$. (b) $\beta_1 = 0.9$, large h , $\epsilon = 10^{-8}$. (c) $\beta_1 = 0.5$, large h , $\epsilon = 10^{-8}$.



(d) $\beta_1 = 0.9$, small h , $\epsilon = 10^{-2}$. (e) $\beta_1 = 0.9$, small h , $\epsilon = 10^{-2}$. (f) $\beta_1 = 0.5$, large h , $\epsilon = 10^{-2}$.

Figure 2.3: The importance of the Adam hyperparameters: learning rate h and decay rate β_1 and ϵ . We use the default $\beta_2 = 0.999$. Top row: the update normalisation in Adam can stabilise training but can also lead to instability compared to momentum; this row uses $\epsilon = 10^{-8}$, the default value and the other hyperparameters are the same as in Figure 2.2 which shows the corresponding momentum trajectories. Bottom row: while keeping other hyperparameters fixed, changing ϵ to 10^{-2} leads to vastly different behaviours and can stabilise training by avoiding dividing by small values.

2.3 Continuous-time methods

We now describe the continuous-time approaches we will use in the rest of this thesis.

2.3.1 A brief overview of stability analysis

Stability analysis is a tool for continuous-time systems that helps decide whether $\theta^* \in \mathbb{R}^D$ is a stable fixed point for a continuous-time flow

$$\dot{\theta} = f(\theta), \quad (2.26)$$

where $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is the vector field of the flow. The most used definitions of stability are Lyapunov stability and asymptotic stability. Lyapunov stability requires

$$\forall \epsilon > 0, \exists \delta > 0 \quad \|\boldsymbol{\theta}(0) - \boldsymbol{\theta}^*\| < \delta \implies \forall t \geq 0 \quad \|\boldsymbol{\theta}(t) - \boldsymbol{\theta}^*\| < \epsilon. \quad (2.27)$$

Asymptotic stability requires Lyapunov stability and additionally

$$\exists \delta > 0 \quad \|\boldsymbol{\theta}(0) - \boldsymbol{\theta}^*\| < \delta \implies \lim_{t \rightarrow \infty} \|\boldsymbol{\theta}(t) - \boldsymbol{\theta}^*\| = 0. \quad (2.28)$$

The rate of convergence of asymptotic stability is exponential if

$$\exists \delta, C, \alpha > 0 \quad \|\boldsymbol{\theta}(0) - \boldsymbol{\theta}^*\| < \delta \implies \forall t \geq 0 \quad \|\boldsymbol{\theta}(t) - \boldsymbol{\theta}^*\| \leq C \|\boldsymbol{\theta}(0) - \boldsymbol{\theta}^*\| e^{-\alpha t}. \quad (2.29)$$

Exponential asymptotic stability can be established using the following criterion:

Remark 2.3.1 (From [118]) *A fixed point $\boldsymbol{\theta}^*$ with $f(\boldsymbol{\theta}^*) = \mathbf{0}$ where the Jacobian of the vector field $\frac{df}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*)$ has eigenvalues with strictly negative real part is a stable fixed point of the flow $\dot{\boldsymbol{\theta}} = f(\boldsymbol{\theta})$, with asymptotic exponential convergence. In contrast, if a fixed point has a Jacobian with an eigenvalue with strictly positive real part, that is an unstable equilibrium and the flow will not converge to it.*

Remark 2.3.1 provides us with a useful test to establish convergence based on a fixed point's Jacobian eigenvalues. The test is inconclusive for equilibria with non-strictly negative real part Jacobian eigenvalues, i.e. Jacobians where a subset of eigenvalues are have strictly negative real part while some have 0 real part. Throughout this thesis, when performing stability analysis of a continuous-time system we will be testing exponential asymptotic stability using the above criterion.

Corollary 2.1 (The NGF and local minima.) *The NGF $\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E$ is only locally attracted to local minima of E , and all strict local minima are exponentially attractive under the NGF.*

This follows from the application of Remark 2.3.1. If a stationary point $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*) = \mathbf{0}$ is attractive under the NGF, the eigenvalues of the Jacobian $\frac{d(-\nabla_{\boldsymbol{\theta}} E)}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*)$ have

strictly negative or 0 real part. Since $\frac{d(-\nabla_{\boldsymbol{\theta}} E)}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) = -\nabla_{\boldsymbol{\theta}}^2 E$ is a symmetric matrix, this condition is equivalent to the eigenvalues of $\nabla_{\boldsymbol{\theta}}^2 E$ being non-negative, the local minimum condition. Conversely, since strict local minima have strictly positive eigenvalues, the Jacobian has strictly negative eigenvalues at strict local minima, leading to exponential local stability under the NGF.

2.3.2 Backward error analysis

Backward error analysis (BEA) is a tool in numerical analysis developed to understand the discretisation error of numerical integrators. We now present an overview of how to use BEA in the context of the gradient descent update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$ with $\boldsymbol{\theta} \in \mathbb{R}^D$; for a general overview see Hairer et al. [99]. BEA provides a modified vector field

$$\tilde{f}_n(\boldsymbol{\theta}) = -\nabla_{\boldsymbol{\theta}} E + hf_1(\boldsymbol{\theta}) + \dots + h^n f_n(\boldsymbol{\theta}) \quad (2.30)$$

by finding functions $f_1 : \mathbb{R}^D \rightarrow \mathbb{R}^D, \dots, f_n : \mathbb{R}^D \rightarrow \mathbb{R}^D$, such that the solution of the modified flow at order n , that is,

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E + hf_1(\boldsymbol{\theta}) + \dots + h^n f_n(\boldsymbol{\theta}) \quad (2.31)$$

follows the discrete dynamics of the gradient descent update with an error $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}(h)\|$ of order $\mathcal{O}(h^{n+2})$, where $\boldsymbol{\theta}(h)$ is the solution of the modified equation truncated at order n at time h and $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_{t-1}$. The full modified vector field with all orders ($n \rightarrow \infty$)

$$\tilde{f}(\boldsymbol{\theta}) = -\nabla_{\boldsymbol{\theta}} E + hf_1(\boldsymbol{\theta}) + \dots + h^n f_n(\boldsymbol{\theta}) + \dots \quad (2.32)$$

is usually divergent and only forms an asymptotic expansion. What BEA provides is the Taylor expansion in h of an unknown h -dependent vector field $f_h(\boldsymbol{\theta})$ developed at $h = 0$:

$$\tilde{f}(\boldsymbol{\theta}) = \text{Taylor}_{h=0} f_h(\boldsymbol{\theta}). \quad (2.33)$$

Thus, a strategy for finding f_h is to find a series of the form in Eq (2.32) via BEA and then find the function f_h such that its Taylor expansion in h at 0 results in the found series. Using this approach we can find the flow $\dot{\boldsymbol{\theta}} = f_h(\boldsymbol{\theta})$ which describes

the gradient descent step $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ exactly.

While flows obtained using BEA are constructed to approximate one gradient descent step, the same flows can be used over multiple gradient descent steps as shown in Section A.1.8 in the Appendix.

Barrett and Dherin [12] used this technique to find the $\mathcal{O}(h^2)$ correction term of gradient descent in the single-objective setting. They showed that for a model with parameters $\boldsymbol{\theta}$ and loss $E(\boldsymbol{\theta})$, optimised with gradient descent $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta})$, the first-order modified equation is

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}}E - \frac{h}{2}\nabla_{\boldsymbol{\theta}}^2E\nabla_{\boldsymbol{\theta}}E, \quad (2.34)$$

which can be written as

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}}\tilde{E}(\boldsymbol{\theta}) = -\nabla_{\boldsymbol{\theta}}\left(E(\boldsymbol{\theta}) + \frac{h}{4}\|\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta})\|^2\right). \quad (2.35)$$

Thus, gradient descent can be seen as implicitly minimising the modified loss

$$\tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}) + \frac{h}{4}\|\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta})\|^2. \quad (2.36)$$

This shows that when training models with gradient descent, there is an implicit regularisation effect, dependent on learning rate h , which biases learning towards paths with low gradient norms. The authors refer to this phenomenon as ‘implicit gradient regularisation’; we will thus refer to the flow in Eq (2.34) as the IGR flow.

2.3.3 BEA proofs

The general structure of BEA proofs is as follows: start with a Taylor expansion in h of the modified flow in Eq (2.31); write each term in the Taylor expansion as a function of $\nabla_{\boldsymbol{\theta}}E$ and the desired f_i using the chain rule; group together terms of the same order in h in the expansion; and identify f_i such that all terms of $\mathcal{O}(h^p)$ are 0 for $p \geq 2$, as is the case in the gradient descent update. A formal overview of BEA proofs can be found in Section A.1.1 in the Appendix.

We now exemplify how to use BEA to find the IGR flow we showed in

Eq (2.34) [12]. Since we are only looking for the first correction term, we only need to find f_1 . We perform a Taylor expansion to find the value of $\boldsymbol{\theta}(h)$ up to order $\mathcal{O}(h^3)$ and then identify f_1 from that expression such that the error $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}(h)\|$ is of order $\mathcal{O}(h^3)$. We have

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}_{t-1} + h\boldsymbol{\theta}^{(1)}(\boldsymbol{\theta}_{t-1}) + \frac{h^2}{2}\boldsymbol{\theta}^{(2)}(\boldsymbol{\theta}_{t-1}) + \mathcal{O}(h^3). \quad (2.37)$$

We know by the definition of the modified vector field in Eq (2.31) that

$$\boldsymbol{\theta}^{(1)} = -\nabla_{\boldsymbol{\theta}}E + hf_1(\boldsymbol{\theta}). \quad (2.38)$$

We can then use the chain rule to obtain

$$\boldsymbol{\theta}^{(2)} = \frac{d(-\nabla_{\boldsymbol{\theta}}E + hf_1(\boldsymbol{\theta}))}{dt} = -\frac{d\nabla_{\boldsymbol{\theta}}E}{dt} + \mathcal{O}(h) = \nabla_{\boldsymbol{\theta}}^2E\nabla_{\boldsymbol{\theta}}E + \mathcal{O}(h). \quad (2.39)$$

Thus,

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) + h^2f_1(\boldsymbol{\theta}_{t-1}) + \frac{h^2}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta}_{t-1})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) + \mathcal{O}(h^3). \quad (2.40)$$

We can then write

$$\boldsymbol{\theta}_t - \boldsymbol{\theta}(h) = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}(h) \quad (2.41)$$

$$= h^2f_1(\boldsymbol{\theta}_{t-1}) + \frac{h^2}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta}_{t-1})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) + \mathcal{O}(h^3). \quad (2.42)$$

For the error to be of order $\mathcal{O}(h^3)$ the terms of order $\mathcal{O}(h^2)$ have to be $\mathbf{0}$. This entails

$$f_1(\boldsymbol{\theta}_{t-1}) = -\frac{1}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta}_{t-1})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}), \quad (2.43)$$

from which we conclude to $f_1 = -\frac{1}{2}\nabla_{\boldsymbol{\theta}}^2E\nabla_{\boldsymbol{\theta}}E$ leading to Eq (2.34).

2.4 Multiple objective optimisation

We have thus far discussed optimisers for single-objective problems, which encompasses many machine learning use cases, including supervised learning and certain formulations of reinforcement learning and generative modelling. There are, however, use cases in machine learning beyond single-objective problems, including GANs and types of reinforcement learning, such as actor-critic frameworks. These settings are often formulated via nested minimisation problems:

$$\min_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\arg \min_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}(\boldsymbol{\phi}, \boldsymbol{\theta}), \boldsymbol{\theta}). \quad (2.44)$$

Solving nested optimisation problems is computationally expensive, as it would require solving the inner optimisation problem for each parameter update of the outer optimisation problem. Thus, in deep learning practice nested optimisation problems are often translated into problems of the form

$$\min_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}(\boldsymbol{\phi}, \boldsymbol{\theta}) \quad (2.45)$$

$$\min_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\boldsymbol{\phi}, \boldsymbol{\theta}). \quad (2.46)$$

Multi-objective optimisation problems of this kind are often cast as a game. Each player has a set of actions to take, in the form of an update for their parameters, $\boldsymbol{\phi}$ for the first player, and $\boldsymbol{\theta}$ for the second. The loss functions are the negative of the player's respective payoff. As it is common in the literature, we will often adapt the game theory nomenclature in this setting; this connection allows for a bridge between the machine learning and the game theory communities.

When translating problems of the form in Eq (2.44) into optimisation procedures that preserve the nested structure of a game it is common to use one of the single-objective algorithms discussed Section 2.2 to update the first player's parameters, followed by the use of the same algorithm to update the second player's parameters. This approach is often referred to as *alternating updates*; we summarise it in Algorithm 1. Optimisation approaches that further account for the nested game

Algorithm 1 Alternating updates

$k \geq 0$; the number of first player updates for each second player update
 Choice of optimiser_routine(E , current_value); e.g: $gd(E, \boldsymbol{\theta}_t) = \boldsymbol{\theta}_t - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_t)$
 Initialise $\boldsymbol{\phi}_0, \boldsymbol{\theta}_0$
 $t \leftarrow 0$
while training **do**
 $\boldsymbol{\phi} \leftarrow \boldsymbol{\phi}_t$
 for i in $\{1, \dots, k\}$ **do**
 $\boldsymbol{\phi} \leftarrow \text{optimiser_routine}(E_{\boldsymbol{\phi}}(\cdot, \boldsymbol{\theta}_t), \boldsymbol{\phi})$
 end for
 $\boldsymbol{\phi}_{t+1} \leftarrow \boldsymbol{\phi}$
 $\boldsymbol{\theta}_{t+1} \leftarrow \text{optimiser_routine}(E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_{t+1}, \cdot), \boldsymbol{\theta}_t)$
 $t \leftarrow t + 1$
end while

Algorithm 2 Simultaneous updates

Initialise $\boldsymbol{\phi}_0, \boldsymbol{\theta}_0$
 Choice of optimiser_routine(E , current_value); e.g: $gd(E, \boldsymbol{\theta}_t) = \boldsymbol{\theta}_t - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_t)$
while training **do**
 $\boldsymbol{\phi}_{t+1} \leftarrow \text{optimiser_routine}(E_{\boldsymbol{\phi}}(\cdot, \boldsymbol{\theta}_t), \boldsymbol{\phi}_t)$
 $\boldsymbol{\theta}_{t+1} \leftarrow \text{optimiser_routine}(E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_t, \cdot), \boldsymbol{\theta}_t)$
 $t \leftarrow t + 1$
end while

structure do exist [119], but they tend to be more computationally prohibitive and thus less used in practice. In some settings, there is no nested structure or the nested structure is further discarded and both players are updated simultaneously, as we show in Algorithm 2.

Simultaneous gradient descent updates derived for solving the problem in Eqs (2.45) and (2.46)—using gradient descent as the optimiser in Algorithm 2—can be seen as the Euler discretisation of the continuous-time flow

$$\dot{\boldsymbol{\phi}} = -\nabla_{\boldsymbol{\phi}}E_{\boldsymbol{\phi}} \quad (2.47)$$

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}}E_{\boldsymbol{\theta}}. \quad (2.48)$$

Analysis of these game dynamics is more challenging than in the single objective case (shown in Eq (2.3)), as we can no longer ascertain that following such flows

minimises the corresponding loss functions since

$$\frac{dE_\phi}{dt} = \frac{d\phi^T}{dt} \nabla_\phi E_\phi + \frac{d\theta^T}{dt} \nabla_\theta E_\phi = -\|\nabla_\phi E_\phi\|^2 - \nabla_\theta E_\theta^T \nabla_\theta E_\phi, \quad (2.49)$$

which need not be negative. Furthermore, it is unclear how to transfer the alternating structure that preserves player order from Algorithm 1 into a continuous-time flow; we will address this in a later chapter. Nonetheless, the flows in Eqs (2.47) and (2.48) are often studied to understand the behaviour of gradient descent in nested optimisation and games [17, 61, 67, 120].

The challenges pertaining to analysing optimisation in games get further compounded when considering the notions of equilibria. The natural translation of local minimum from single-objective optimisation to games is a local Nash equilibrium. A local Nash equilibrium is a point in parameter space (ϕ^*, θ^*) such that no player has an incentive to deviate from (ϕ^*, θ^*) to another point in its local neighbourhood; that is ϕ^* is a local minimum for $E_\phi(\cdot, \theta^*)$ and θ^* is a local minimum for $E_\theta(\phi^*, \cdot)$. Equivalently, $\nabla_\phi E_\phi(\phi^*, \theta^*) = \mathbf{0}$ and $\nabla_\theta E_\theta(\phi^*, \theta^*) = \mathbf{0}$ and $\nabla_\phi^2 E_\phi(\phi^*, \theta^*)$ and $\nabla_\theta^2 E_\theta(\phi^*, \theta^*)$ are positive semi-definite, i.e. do not have negative eigenvalues.

While in the single-objective case, local minima are the only exponentially attractive equilibria of the NGF (see Corollary 2.1), this *correspondence between desired equilibrium and attractive dynamics does not generally occur in games*. By applying stability analysis (Remark 2.3.1), we know that the system in Eqs (2.47) and (2.48) is attracted to equilibria for which the Jacobian

$$\mathbf{J} = \begin{bmatrix} -\frac{dE_\phi}{d\phi d\phi} & -\frac{dE_\phi}{d\theta d\phi} \\ -\frac{dE_\theta}{d\phi d\theta} & -\frac{dE_\theta}{d\theta d\theta} \end{bmatrix} \quad (2.50)$$

evaluated at (ϕ^*, θ^*) has only eigenvalues with strictly negative real part. In contrast, a Nash equilibrium requires that the block diagonals of \mathbf{J} , namely $-\nabla_\phi^2 E_\phi(\phi^*, \theta^*)$ and $-\nabla_\theta^2 E_\theta(\phi^*, \theta^*)$, have negative eigenvalues. Since for an unconstrained matrix there is no relationship between its eigenvalues and the eigenvalues of its block diagonals, there is no inclusion relationship between Nash equilibria and locally

attractive equilibria of the underlying continuous game dynamics that applies to all two-player games.

For zero-sum games, i.e. games where $E_\phi = -E_\theta = E$, a strict Nash equilibrium will be locally attractive when following the flow in Eqs (2.47) and (2.48). To see why, we have to examine the eigenvalues of the flow's Jacobian

$$\mathbf{J} = \begin{bmatrix} -\frac{dE}{d\phi d\phi} & -\frac{dE}{d\theta d\phi} \\ \frac{dE}{d\phi d\theta} & \frac{dE}{d\theta d\theta} \end{bmatrix} \quad (2.51)$$

evaluated at a strict Nash equilibrium of such game. Due to the zero-sum formulation, the off-diagonal blocks of \mathbf{J} satisfy $-\frac{dE}{d\theta d\phi} = -(\frac{dE}{d\phi d\theta})^T$, and we have

$$\mathbf{J} + \mathbf{J}^T = \begin{bmatrix} -\nabla_\phi^2 E & \mathbf{0} \\ \mathbf{0} & \nabla_\theta^2 E \end{bmatrix}. \quad (2.52)$$

Thus, $\mathbf{J} + \mathbf{J}^T$ at the Nash equilibrium (ϕ^*, θ^*) is a symmetric block diagonal matrix with negative definite blocks, $-\nabla_\phi^2 E(\phi^*, \theta^*)$ and $\nabla_\theta^2 E(\phi^*, \theta^*)$, and thus has negative eigenvalues. From here it follows that \mathbf{J} has eigenvalues with strictly negative real part, and thus satisfies the condition of Remark 2.3.1. We note that the reverse does not hold, as being locally attractive in this case is a weaker condition than that of a Nash equilibrium. This observation together with the realisation that in deep learning certain games might not have Nash equilibria [121], has lead the search for other equilibrium measures including local stability and Stackelberg equilibria [69, 122, 123].

2.5 Conclusion

In this section we have introduced the main optimisation algorithms used in deep learning, together with their main challenges and the analytical tools we will employ in the rest of this thesis.

Chapter 3

A new continuous-time model of gradient descent

The recipe behind the success of deep learning has been training neural networks with gradient-based optimisation on large datasets. Understanding the behaviour of gradient descent, however, and particularly its instability, has lagged behind its empirical success. To add to the theoretical tools available to study gradient descent we propose *the principal flow* (PF), a continuous-time flow that approximates gradient descent dynamics. To the best of our knowledge, the PF is the only continuous flow that captures the divergent and oscillatory behaviours of gradient descent, including escaping local minima and saddle points. Through its dependence on the eigen-decomposition of the Hessian the PF sheds light on the recently observed edge of stability phenomena in deep learning [52]. Using our new understanding of instability we propose a learning rate adaptation method that enables us to control the trade-off between training stability and test set evaluation performance.

3.1 Introduction

Our goal is to use continuous-time models to understand the behaviour of gradient descent. Using continuous-time dynamics to understand discrete-time systems opens up tools from dynamical systems, such as stability analysis, and has a long history in optimisation and machine learning [11–13, 17, 54–60]. Most theoretical analysis of gradient descent using continuous-time systems uses the negative gradient flow, but this has well-known limitations, such as not being able to explain any

behaviour contingent on the learning rate. To mitigate these limitations we find a new continuous-time flow that reveals important new roles of the Hessian in gradient descent training. To do so, we use backward error analysis (BEA), a method with a long history in the numerical integration community [99] that has only recently been used in the deep learning context [12, 13].

We find that the proposed flow sheds new light on gradient descent stability, including, but not limited to, divergent and oscillatory behaviour around a fixed point. Instability—areas of training where the loss consistently increases—and edge of stability behaviours [52]—areas of training where the loss does not behave monotonically but decreases over long time periods—are pervasive in deep learning and occur for all learning rates and architectures [52, 53, 124, 125]. We use our novel insights to understand and mitigate these instabilities.

The structure of this chapter is as follows:

- We discuss the advantages of a continuous-time approach in Section 3.2, where we also highlight the limitations of existing continuous-time flows.
- We introduce **the principal flow** (PF), a flow in complex space defined by the eigen-decomposition of the Hessian (Section 3.3). To the best of our knowledge the PF is the first continuous-time flow that captures that gradient descent can diverge around local minima and saddle points. We show that using a complex flow is crucial in understanding instabilities in gradient descent.
- We show the PF is better than existing flows at modelling neural network training dynamics in Section 3.4. In Section 3.5, we use the PF to shed new light on edge of stability behaviours in deep learning. We do so by connecting changes in the loss and Hessian eigenvalues with core quantities exposed by the PF and neural network landscapes explored through the behaviour of gradient flows.
- Through a continuous-time perspective we demonstrate empirically how to control the trade-off between stability and performance in deep learning in Section 3.7. We do so using DAL (Drift Adjusted Learning), an approach to

setting the learning rate dynamically based on insights on instability derived from the PF.

- We end by showcasing the potential of integrating our continuous-time approach with other optimisation schemes and highlighting how the PF can be used as a tool for existing continuous-time analyses in Section 3.8.

Notation: We denote as E the loss function, $\boldsymbol{\theta}$ the parameter vector of dimension D , $\nabla_{\boldsymbol{\theta}}^2 E \in \mathbb{R}^{D \times D}$ the loss Hessian and λ_i the Hessian's i th largest eigenvalue with \mathbf{u}_i the corresponding eigenvector. Since if \mathbf{u}_i is an eigenvector of $\nabla_{\boldsymbol{\theta}}^2 E$ so is $-\mathbf{u}_i$, we always use \mathbf{u}_i such that $\Re[\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i] \geq 0$; this has no effect on our results and is only used for convenience and we used the notation $\Re(z)$ to be the real part of z . For a continuous-time flow $\boldsymbol{\theta}(h)$ refers to the solution of the flow at time h .

3.2 Continuous time models of gradient descent

The aim of this work is to understand the dynamics of gradient descent updates with learning rate h

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) \quad (3.1)$$

from the perspective of continuous dynamics. When using continuous-time dynamics to understand gradient descent it is most common to use *the negative gradient flow* (NGF)

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E. \quad (3.2)$$

Gradient descent can be obtained from the NGF through Euler numerical integration, with an error of $\mathcal{O}(h^2)$ after one gradient descent step, see Section 2.2. Studying gradient descent and its behaviour around equilibria and beyond has thus taken two main approaches: directly studying the discrete updates of Eq (3.1) [15, 49, 64, 102, 126–130], or the continuous-time NGF of Eq (3.2) [11, 17, 54–61]. The appeal of continuous-time systems lies in their connection with dynamical systems and the plethora of tools that thus become available, such as stability analysis; the

simplicity by which conserved quantities can be obtained [60, 131]; and analogies that can be constructed through similarities with physical systems [60]. Because of the availability of tools for the analysis of continuous-time systems, it has been previously noted that discrete-time approaches are often more challenging and discrete-time proofs are often inspired from continuous-time ones [54, 62]. We use an example to showcase the ease of continuous-time analyses: when following the NGF the loss E decreases since $\frac{dE}{dt} = -\|\nabla_{\theta} E\|^2$, as we derived in Eq 2.3. Showing that and *when* following the discrete-time gradient descent update in Eq (3.1) is more challenging and requires adapting the analysis on the form of the loss function E . Classical convergence guarantees associated with other optimisation approaches, such as natural gradient, are also derived in continuous-time [132–134]. By analysing the properties of continuous-time systems one can also determine whether optimisers should more closely follow the underlying continuous-time flow [67, 135], what regularisers should be constructed to ensure convergence or stability [17, 61, 136], construct convergence guarantees in functional space for infinitely wide networks [100, 137].

3.2.1 Limitations of existing continuous-time flows

The well-known discrepancy between Euler integration and the NGF, often called *discretisation error* or *discretisation drift* (Figure 3.1) leads to certain limitations when using the NGF to describe gradient descent: the NGF cannot explain divergence around a local minimum for high learning rates or convergence to flat minima as often seen in the training of neural networks. Critically, since the NGF does not depend on the learning rate, it cannot explain any learning rate dependent behaviour.

The appeal of continuous-time methods together with the limitations of the NGF have inspired the machine learning community to look for other continuous-time systems that may better approximate the gradient descent trajectory. One approach to constructing continuous-time flows approximating gradient descent that takes into account the learning rate is backward error analysis (BEA); for an overview of BEA, see Section 2.3.2. Using this approach, Barrett and Dherin [12] introduce the

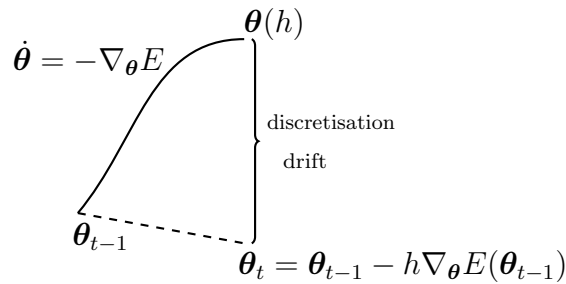


Figure 3.1: Discretisation drift. Using continuous-time flows to understand gradient descent is limited by the gap between the discrete and continuous dynamics. In the case of the negative gradient flow $\dot{\theta} = -\nabla_{\theta} E$, we call this gap *discretisation drift*. Since the negative gradient flow always minimises E , instabilities in gradient descent—areas where the loss E increases—are due to discretisation drift.

Implicit Gradient Regularisation flow (IGR flow):

$$\dot{\theta} = -\nabla_{\theta} E - \frac{h}{2} \nabla_{\theta}^2 E \nabla_{\theta} E, \quad (3.3)$$

which tracks the dynamics of the gradient descent step $\theta_t = \theta_{t-1} - h\nabla_{\theta} E(\theta_{t-1})$ with an error of $\mathcal{O}(h^3)$, thus reducing the order of the error compared to the NGF. Unlike the NGF flow, the IGR flow depends on the learning rate h . This dependence explains certain properties of gradient descent, such as avoiding trajectories with high gradient norm; the authors connect this behaviour to convergence to flat local minima.

Like the NGF flow, however, the IGR flow does not explain the instabilities of gradient descent, as we illustrate in Figure 3.2. Indeed, Barrett and Dherin [12] (their Remark 3.4) show that performing stability analysis around local minima using the IGR flow does not lead to qualitatively different conclusions from those using the NGF: both NGF and the IGR flow predict gradient descent to be always locally attractive around a local minimum (proofs in Section A.1.7), contradicting the empirically observed behaviour of gradient descent.

To understand why both the NFG and the IGR flow cannot capture oscillations and divergence around a local minimum, we note that stationary points $\nabla_{\theta} E = \mathbf{0}$ are fixed points for both flows. We visualise an example in Figure 3.3a: since to

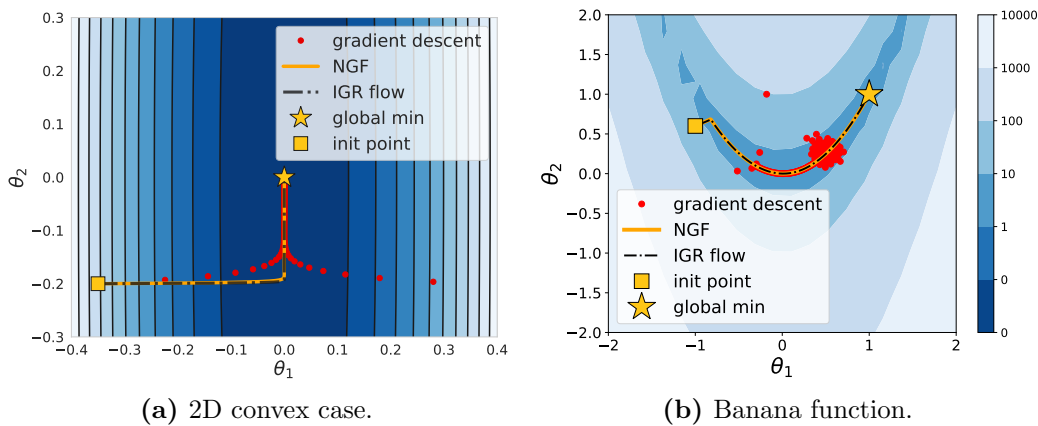


Figure 3.2: Motivation. Existing continuous-time flows fail to capture the oscillatory or unstable behaviour of gradient descent.

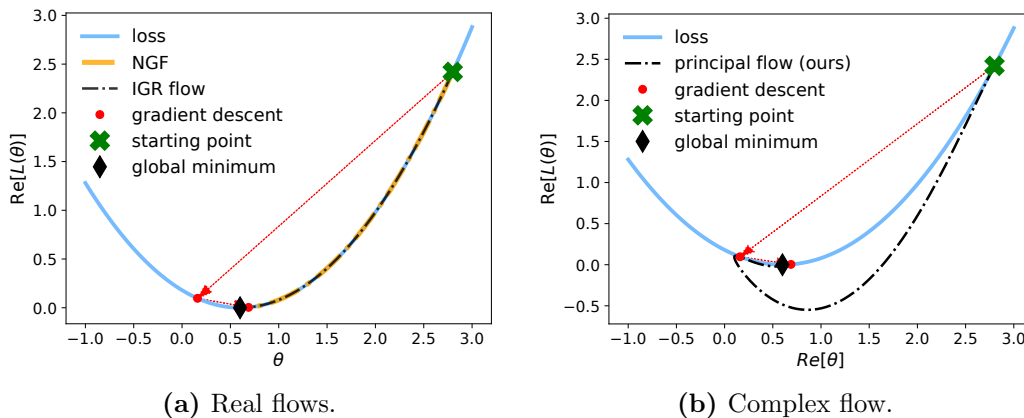


Figure 3.3: Complex flows capture oscillations and divergence around local minima. In the real space, the trajectory going from the starting point to the second gradient descent iterate goes through the global minima, and real flows stop there. In complex space however, that need not be the case.

go from the initial point to the gradient descent iterates requires passing through the local minimum, both flows would stop at the local minimum and never reach the following gradient descent iterates. In the case of neural networks we show in Figure 3.4 that while the IGR flow is better than the NGF at describing gradient descent, a substantial gap remains.

The lack of ability of existing continuous-time flows to model instabilities empirically observed in gradient descent such as those shown in Figure 3.2 has been used as a motivation to use discrete-time methods instead [63, 64]. The goal of our work is to overcome this issue by introducing a novel continuous-time flow that

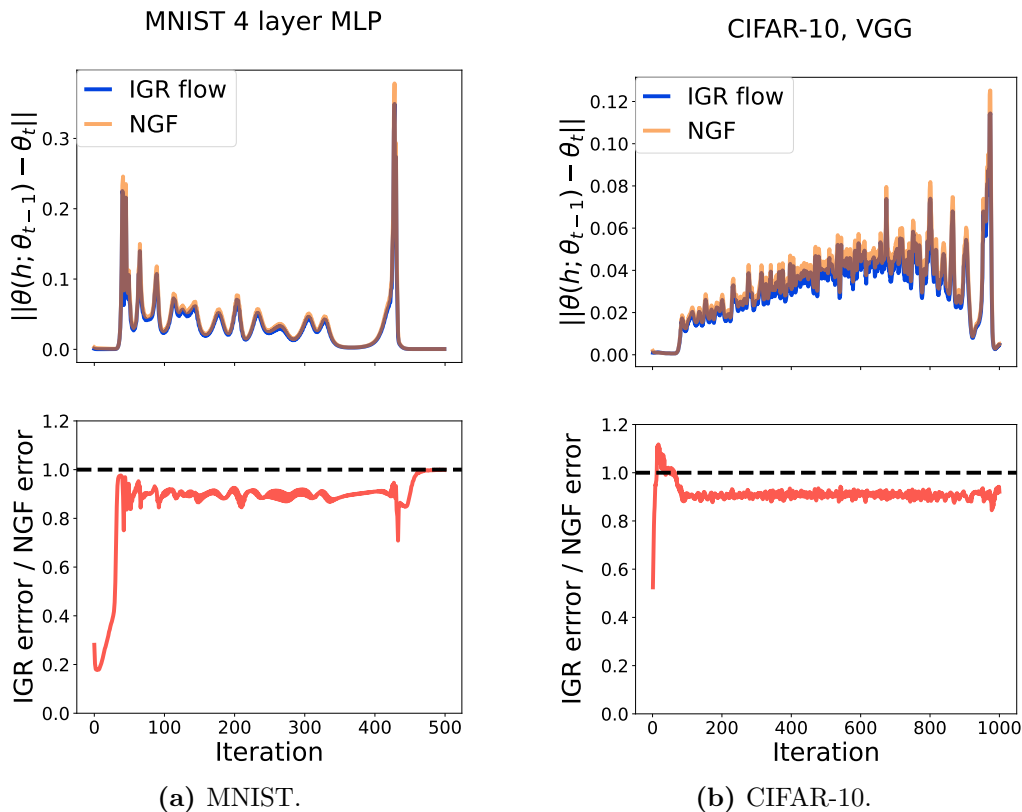


Figure 3.4: Motivation. While the IGR flow captures aspects of the per-iteration discretisation drift for neural networks trained with gradient descent, a significant gap remains. To measure how well each flow captures gradient descent dynamics at one iteration, we train a model using gradient descent $\theta_t = \theta_{t-1} - h\nabla_{\theta} E(\theta_{t-1})$ and at each iteration t we approximate the respective flows for time h . We then compute the difference in norms between the resulting flow parameters and the gradient descent parameters at the next iteration: $\|\theta_t - \theta(h)\|$ with $\theta(0) = \theta_{t-1}$, which we report in the top row. The bottom row reports the ratio between these error norms obtained with the IGR and NGF flow, in order to provide a relative error scale. We observe that in the two experiments we perform here, the IGR flow captures approximately 10% of discretisation drift at the majority of training iterations.

captures instabilities observed in gradient descent. To do so, we follow the footsteps of Barrett and Dherin [12] and use BEA. By using a continuous-time flow we can leverage the tools and advantages of continuous-time methods discussed earlier in this section; by incorporating discretisation drift into our model of gradient descent we can increase their applicability to explain unstable training behaviour. Indeed, we show in Figure 3.3b that the flow we propose captures the training instabilities; a key reason why is that, unlike existing flows, it operates in complex space. In

Section 3.3, we show the importance of operating in complex space in order to understand oscillatory and instability behaviours of gradient descent.

3.3 The principal flow

In the previous section, we have seen how BEA can be used to define continuous-time flows that capture the dynamics of gradient descent up to a certain order in learning rate. We have also explored the limitations of these flows, including the lack of ability to explain oscillations observed empirically when using gradient descent. To further expand our understanding of gradient descent via continuous-time methods, we would like to get an intuition for the structure of higher order modified vector fields provided by BEA. We start with the following modified vector field, which we will call *the third-order flow* (proof in Section A.1.2):

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E - \frac{h}{2} \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E - h^2 \left(\frac{1}{3} (\nabla_{\boldsymbol{\theta}}^2 E)^2 \nabla_{\boldsymbol{\theta}} E + \frac{1}{12} \nabla_{\boldsymbol{\theta}} E^T (\nabla_{\boldsymbol{\theta}}^3 E) \nabla_{\boldsymbol{\theta}} E \right), \quad (3.4)$$

with $\nabla_{\boldsymbol{\theta}}^3 E \in \mathbb{R}^{D \times D \times D}$ and $(\nabla_{\boldsymbol{\theta}} E^T (\nabla_{\boldsymbol{\theta}}^3 E) \nabla_{\boldsymbol{\theta}} E)_k = \sum_{i,j} (\nabla_{\boldsymbol{\theta}} E)_i (\nabla_{\boldsymbol{\theta}}^3 E)_{i,k,j} (\nabla_{\boldsymbol{\theta}} E)_j$. The third-order flow tracks the dynamics of the gradient descent step $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$ with an error of $\mathcal{O}(h^4)$, thus further reducing the order of the error compared to the IGR flow. Like the IGR flow and the NGF, the third-order flow has the property that $\dot{\boldsymbol{\theta}} = \mathbf{0}$ if $\nabla_{\boldsymbol{\theta}} E = \mathbf{0}$ and thus will exhibit the same limitations observed in Figure 3.3. The third-order flow allows us to spot a pattern: the correction term of order $\mathcal{O}(h^n)$ in the BEA modified flow describing gradient descent contains the term $(\nabla_{\boldsymbol{\theta}}^2 E)^n \nabla_{\boldsymbol{\theta}} E$ and terms that contain higher order derivatives with respect to parameters, terms which we will denote as $\mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E)$.

Our approach. We will use the terms of the form $(\nabla_{\boldsymbol{\theta}}^2 E)^n \nabla_{\boldsymbol{\theta}} E$ to construct a new continuous-time flow. We will take a three-step approach. First, for an arbitrary order $\mathcal{O}(h^n)$ we will find the terms containing only first and second order derivatives in the modified vector field given by BEA and show they are of the form $(\nabla_{\boldsymbol{\theta}}^2 E)^n \nabla_{\boldsymbol{\theta}} E$ (Theorem 3.3.1). Second, we will use all orders to create a series (Corollary 3.1). Third, we will use the series to find the modified flow given by BEA (Theorem 3.3.2). All proofs are provided in Section A.1 of the Appendix.

Theorem 3.3.1 *The modified vector field with an error of order $\mathcal{O}(h^{n+2})$ to the gradient descent update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ has the form:*

$$\dot{\boldsymbol{\theta}} = \sum_{p=0}^n \frac{-1}{p+1} h^p (\nabla_{\boldsymbol{\theta}}^2 E)^p \nabla_{\boldsymbol{\theta}} E + \mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E), \quad (3.5)$$

where $\mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E)$ denotes the family of functions that can be written as a sum of terms, each term containing a derivative of higher order than 3 with respect to parameters.

The result is proven by induction, with the full proof provided in Section A.1.3 of the Appendix. The base cases for $n = 1, 2$, and 3 follow from the NGF, IGR, and third-order flows. For higher order terms, the proof uses induction to find the term in f_i depending on $\nabla_{\boldsymbol{\theta}}^2 E$ and $\nabla_{\boldsymbol{\theta}} E$ only and follows the BEA proof structure highlighted in Section 2.3.2, but Step 3 is modified to not account for terms in $\mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E)$. From the above, we can obtain the following corollary by using all orders n and the eigen-decomposition of $\nabla_{\boldsymbol{\theta}}^2 E$:

Corollary 3.1 *The full order modified flow obtained by performing BEA on gradient descent updates is of the form:*

$$\dot{\boldsymbol{\theta}} = \sum_{p=0}^{\infty} \frac{-1}{p+1} h^p (\nabla_{\boldsymbol{\theta}}^2 E)^p \nabla_{\boldsymbol{\theta}} E + \mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E) \quad (3.6)$$

$$= \sum_{p=0}^{\infty} \frac{-1}{p+1} h^p \left(\sum_{i=0}^{D-1} \lambda_i^p \mathbf{u}_i \mathbf{u}_i^T \right) \nabla_{\boldsymbol{\theta}} E + \mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E) \quad (3.7)$$

$$= \sum_{i=0}^{D-1} \left(\sum_{p=0}^{\infty} \frac{-1}{p+1} h^p \lambda_i^p \right) (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i + \mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E), \quad (3.8)$$

where λ_i and \mathbf{u}_i are the respective eigenvalues and eigenvectors of the Hessian $\nabla_{\boldsymbol{\theta}}^2 E$.

If $\lambda_0 > 1/h$ the BEA series above diverges. Generally BEA series are not convergent and approximate the discrete scheme only by truncation [99]. When the series in Eq (3.8) diverges, truncating it up to any order n , however, will result in a flow that will not be able to capture instabilities, even in the quadratic case. Such flows (including the IGR flow) will always predict the loss function will decrease for a quadratic loss where a minimum exists, since:

$\frac{dE}{dt} = \nabla_{\theta} E^T \left(\sum_{p=0}^n \frac{-1}{p+1} h^p (\nabla_{\theta}^2 E)^p \nabla_{\theta} E \right) = - \sum_{p=0}^n \frac{1}{p+1} h^p \sum_{i=0}^{D-1} (\lambda_i^p) (\nabla_{\theta} E^T \mathbf{u}_i)^2$ which is never positive for any quadratic loss where a minimum exists (i.e. when $\lambda_i \geq 0, \forall i$). The above also entails that the flows always predict convergence around a local minimum, which is not the case for gradient descent which can diverge for large learning rates. To further track instabilities we can use the BEA series to find the following flow:

Definition 3.3.1 We define the *principal flow (PF)* as

$$\dot{\theta} = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i. \quad (3.9)$$

We note that $\lim_{\lambda \rightarrow 0} \frac{\log(1-h\lambda)}{h\lambda} = -1$ and thus the PF is well defined when the Hessian $\nabla_{\theta}^2 E$ is not invertible. Unlike the NGF and the IGR flow, the modified vector field of the PF cannot be always written as the gradient of a loss function in \mathbb{R} , and can be complex valued.

Theorem 3.3.2 The Taylor expansion in h at $h = 0$ of the PF vector field coincides with the series coming from the BEA of gradient descent (Eq (3.8)).

Proof: Using the Taylor expansion at $z = 0$, $\text{Taylor}_{z=0} \frac{\log(1-z)}{z} = \sum_{p=0}^{\infty} \frac{-1}{p+1} z^p$ we obtain:

$$\text{Taylor}_{h=0} \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i = \sum_{i=0}^{D-1} \left(\sum_{p=0}^{\infty} \frac{-1}{p+1} h^p \lambda_i^p \right) (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i. \quad (3.10)$$

□

We have used BEA to find the flow that when Taylor expanded at $h = 0$ leads to the series in Eq (3.8). When the BEA series in Eq (3.8) converges, namely $\lambda_0 < 1/h$, the PF and the flow given by the BEA series are the same. When $\lambda_0 > 1/h$, however, the PF is complex and the BEA series diverges. While in this case any BEA truncated flow will not be able to track gradient descent closely, we show that for quadratic losses the PF will track gradient descent exactly, and that it is a good model of

gradient descent around fixed points. We show examples of the PF tracking gradient descent exactly in the quadratic case in Figures 3.3b and 3.7.

Remark 3.3.1 For quadratic losses of the form $E = \frac{1}{2}\boldsymbol{\theta}^T \mathbf{A}\boldsymbol{\theta} + \mathbf{b}^T \boldsymbol{\theta}$, the PF captures gradient descent exactly. This case has been proven in Hairer et al. [99]. The solution of the PF can also be computed exactly in terms of the eigenvalues of $\nabla_{\boldsymbol{\theta}}^2 E$: $\boldsymbol{\theta}(t) = \sum_{i=0}^{D-1} e^{\frac{\log(1-h\lambda_i)}{h}t} \boldsymbol{\theta}_0^T \mathbf{u}_i \mathbf{u}_i + t \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i$.

Remark 3.3.2 In a small enough neighbourhood around a critical point (where higher-order derivatives can be ignored) the PF can be used to describe gradient descent dynamics closely. We show this also using a linearisation argument in Section A.1.4 in the Appendix.

Definition 3.3.2 The terms $\mathcal{C}(\nabla_{\boldsymbol{\theta}}^3 E)$ are called **non-principal terms**. The term $\frac{1}{12} \nabla_{\boldsymbol{\theta}} E^T (\nabla_{\boldsymbol{\theta}}^3 E) \nabla_{\boldsymbol{\theta}} E$ in Eq (3.4) is a non-principal term (we will call this term non-principal third-order term).

Definition 3.3.3 We define the **principal flow with third-order non-principal term** as

$$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i)}{h\lambda_i} (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i - \underbrace{\frac{h^2}{12} \nabla_{\boldsymbol{\theta}} E^T (\nabla_{\boldsymbol{\theta}}^3 E) \nabla_{\boldsymbol{\theta}} E}_{\text{third-order non-principal term}}. \quad (3.11)$$

General theoretical bounds on the error between continuous time flows and gradient descent are challenging to construct in the case of a general parametrised $E(\boldsymbol{\theta})$ as the error will be determined by the shape of E . We know the conditions which determine when certain flows follow gradient descent exactly. The NGF and gradient descent will follow the same trajectory in areas where $\nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E = 0$ (see Theorem 3.7.1) and thus E has a constant gradient in time, since $\frac{d\nabla_{\boldsymbol{\theta}} E}{dt} = \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E$. The PF generalises the NGF, in that it follows the same trajectory as gradient descent not only for trajectories where $\nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E = 0$, but also when E is quadratic. Informally, we can state that the closer we are to these exact conditions, the more likely the flows are to capture the dynamics of gradient descent. Formally, bounds on the error between GD and NGF can be provided by the Fundamental Theorem

NGF	$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \underbrace{-}_{\alpha_{NGF}(h\lambda_i)=-1} (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$
IGR Flow	$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \underbrace{-(1 + \frac{h}{2}\lambda_i)}_{\alpha_{IGR}(h\lambda_i)=-(1+\frac{h}{2}\lambda_i)} (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$
PF	$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \underbrace{\frac{\log(1 - h\lambda_i)}{h\lambda_i}}_{\alpha_{PF}(h\lambda_i)=\frac{\log(1-h\lambda_i)}{h\lambda_i}} (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$

Table 3.1: Understanding the differences between the flows discussed in terms of the eigen-decomposition of the Hessian. All flows have the form $\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \alpha(h\lambda_i)(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$ with different α summarised here.

(Theorem 10.6 in Wanner and Hairer [138]) which has recently been adapted to a neural network parametrisation by Elkabetz and Cohen [54]; this bound depends on the magnitude of the smallest Hessian eigenvalue along the NGF trajectory. We hope that future work can expand the Fundamental Theorem such that error bounds between the PF and gradient descent can be constructed for deep neural networks. Here we take an empirical approach and show that although not exact outside the quadratic case the PF captures key features of the gradient descent dynamics in stable or unstable regions of training, around and outside critical points, for small examples or large neural networks.

3.3.1 The PF and the Hessian eigen-decomposition

All flows considered here have the form $\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \alpha(h\lambda_i)(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$, where α is a function computing the corresponding coefficient; we will denote the one associated with each flow as α_{NGF} , α_{IGR} and α_{PF} respectively. For a side-by-side comparison between the NGF, IGR flow, and the PF as functions of the Hessian eigen-decomposition see Table 3.1. Since $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i \geq 0$, the α function determines the sign of a modified vector field in the direction \mathbf{u}_i . For brevity it will be useful to define the coefficient of \mathbf{u}_i in the vector field of the PF:

Definition 3.3.4 We call $sc_i = \frac{\log(1-h\lambda_i)}{h\lambda_i} (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) = \alpha_{PF}(h\lambda_i) \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i$ the *stability coefficient* for eigendirection i . $sign(sc_i) = sign(\alpha_{PF}(h\lambda_i))$.

In order to understand the PF and how it is different from the NGF we explore

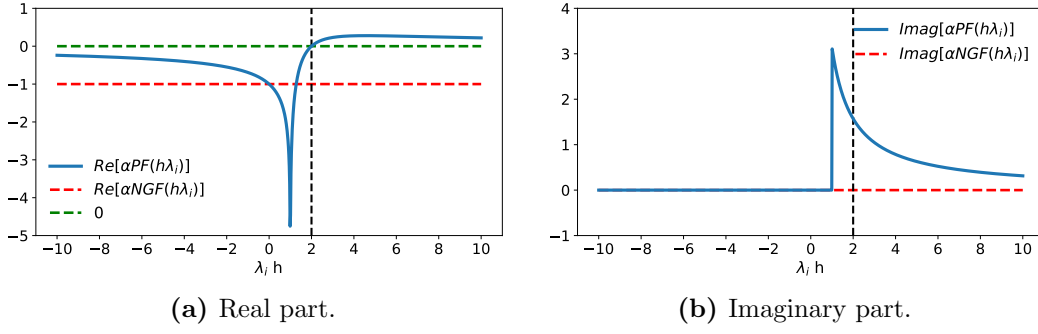


Figure 3.5: Comparing the coefficients α_{NGF} and α_{PF} based on the value of αh . While α_{NGF} is always -1 , the value of α_{PF} depends on αh and can be complex, both with positive and negative real parts.

the change in each eigendirection \mathbf{u}_i and we perform case analysis on the relative value of the eigenvalues λ_i and the learning rate h . To do so, we will compare $\alpha_{NGF}(h\lambda_i)$ and $\alpha_{PF}(h\lambda_i)$ since the sign of $\alpha_{NGF}(h\lambda_i)$ determines the direction that minimises E given by \mathbf{u}_i . Since our goal is to understand the behaviour of gradient descent, we perform the case by case analysis of what happens at the start of a gradient descent iteration and thus use real values for λ_i and \mathbf{u}_i even when the PF is complex valued. We visualise α_{NGF} and α_{PF} in Figure 3.5 and we use Figure 3.6 to show examples of each case using a simple function.

Real stable case: $\lambda_i < 1/h$.

- $\text{sign}(\alpha_{NGF}(h\lambda_i)) = \text{sign}(\alpha_{PF}(h\lambda_i)) = -1$.
- $\alpha_{NGF}(h\lambda_i) = -1$ and $\alpha_{PF}(h\lambda_i) = \frac{\log(1-h\lambda_i)}{h\lambda_i} < 0$.
- The coefficients of the NGF and PF in eigendirection \mathbf{u}_i are negative and real.
- This case is visualised in Figure 3.6a.

Complex stable case: $1/h < \lambda_i < 2/h$.

- $\text{sign}(\alpha_{NGF}(h\lambda_i)) = \text{sign}(\Re[\alpha_{PF}(h\lambda_i)]) = -1$. $\alpha_{PF}(h\lambda_i) \in \mathbb{C}$.
- $\alpha_{NGF}(h\lambda_i) = -1$ and $\alpha_{PF}(h\lambda_i) = \frac{\log(1-h\lambda_i)}{h\lambda_i} = \frac{\log(-1+h\lambda_i)+i\pi}{h\lambda_i} \in \mathbb{C}$ and $\Re[\alpha_{PF}(h\lambda_i)] = \frac{\log(-1+h\lambda_i)}{h\lambda_i} < 0$.
- The real part of the coefficient of the NGF and PF in eigendirection \mathbf{u}_i are both negative. The imaginary part of α_{PF} can introduce instability and oscillations.
- This case is visualised in Figure 3.6b.

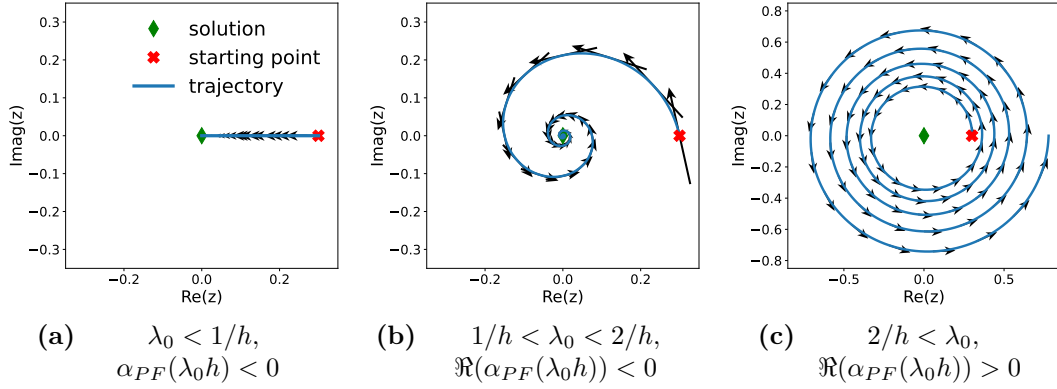


Figure 3.6: The behaviour of PF on the quadratic function $E(z) = \frac{1}{2}z^2$ with solution $z(t) = e^{\log(1-h)/h}z(0)$. When $\lambda_0 < 1/h$, $z(t) = (1-h)^{t/h}z(0)$ which is in real space and converges to the equilibrium. When $\lambda_0 > 1/h$, $z(t) = (h-1)^{t/h}(\cos(\pi t/h) + i \sin(\pi t/h))z(0)$. This exhibits oscillatory behaviour, and when $\lambda_0 > 2/h$, diverges.

Unstable complex case: $2/h < \lambda_i$.

- $\text{sign}(\alpha_{NGF}(h\lambda_i)) \neq \text{sign}(\Re[\alpha_{PF}(h\lambda_i)])$. $\alpha_{PF}(h\lambda_i) \in \mathbb{C}$.
- $\alpha_{NGF}(h\lambda_i) = -1$ and $\alpha_{PF}(h\lambda_i) = \frac{\log(1-h\lambda_i)}{h\lambda_i} = \frac{\log(-1+h\lambda_i)+i\pi}{h\lambda_i} \in \mathbb{C}$ and $\Re[\alpha_{PF}(h\lambda_i)] = \frac{\log(-1+h\lambda_i)}{h\lambda_i} > 0$.
- The real part of the coefficient of the NGF in eigendirection \mathbf{u}_i is negative, while the real part of the coefficient of the PF is positive. The PF goes in the opposite direction of the NGF which minimises E ; this change in sign can cause instabilities. The imaginary component can still introduce oscillations, however, the larger $\lambda_i h$, the smaller the imaginary part of α_{PF} .
- This case is visualised in Figure 3.6c.

The importance of the largest eigenvalue λ_0 . The largest eigenvalue λ_0 plays an important part in the PF. Since $h\lambda_0 \geq h\lambda_i \quad \forall i$, λ_0 determines where in the above cases the PF is situated and thus whether there are oscillations and unstable behaviour in training. For all flows of the form we consider we can write:

$$\frac{dE(\boldsymbol{\theta})}{dt} = \nabla_{\boldsymbol{\theta}} E^T \frac{d\boldsymbol{\theta}}{dt} = \nabla_{\boldsymbol{\theta}} E^T \sum_{i=0}^{D-1} \alpha(h\lambda_i) \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i \mathbf{u}_i = \sum_{i=0}^{D-1} \alpha(h\lambda_i) (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)^2 \quad (3.12)$$

and thus if $\alpha(h\lambda_i) \in \mathbb{R}$ and $\alpha(h\lambda_i) < 0 \quad \forall i$ then $\frac{dE(\boldsymbol{\theta})}{dt} \leq 0$ and following the correspond-

ing flow minimises E . In the case of the PF this gets determined by λ_0 . If $\lambda_0 < \frac{1}{h}$ then $\alpha_{PF}(h\lambda_i) < 0 \forall i$ (real stable case above) and the PF minimises E . If $1/h < \lambda_0 < \frac{2}{h}$ then $\Re[\alpha_{PF}(h\lambda_i)] < 0 \forall i$ (complex stable case above) close to a gradient descent iteration $\lambda_i, \mathbf{u}_i \in \mathbb{R}$ we can write that $\frac{d\Re[E(\boldsymbol{\theta})]}{dt} = \sum_{i=0}^{D-1} \Re[\alpha_{PF}(h\lambda_i)](\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)^2$ and thus the real part of the loss function decreases. If $\lambda_0 > \frac{2}{h}$ then $\Re[\alpha_{PF}(h\lambda_0)] > 0$ (unstable complex case above) and if $(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0)^2$ is sufficiently large we can no longer ascertain the behaviour of E . We present a discrete-time argument for this observation in Section A.2.1.

Building intuition. For a quadratic objective $E(\boldsymbol{\theta}) = \frac{1}{2}\boldsymbol{\theta}^T \mathbf{A}\boldsymbol{\theta}$ the PF describes gradient descent exactly. We show examples Figures 3.3 and 3.7. Unlike the NGF or the IGR flow, the PF captures the oscillatory and divergent behaviour of gradient descent. Importantly, to capture the unstable behaviour that occurs when $\lambda_0 > 1/h$ the imaginary part of the PF is needed. To expand intuition outside the quadratic case, we show the PF for the banana function [139] in Figure 3.8 and an additional example in 1D with a non-quadratic function in Figure A.6 in the Appendix. While in this case, the PF no longer follows the gradient descent trajectory exactly, we still observe the importance of the PF in capturing instabilities of gradient descent. In Figure 3.8c, we also observe that adding non-principal terms—described in Definition 3.3.3—can restabilise the trajectory when $\lambda_0 \gg 2/h$.

Remark 3.3.3 *For the banana function, the principal terms have a destabilising effect when $h > 2/\lambda_0$ while the non-principal terms can have a stabilising effect.*

3.3.2 The stability analysis of the PF

We now perform stability analysis on the PF, to understand how it can be used to predict certain behaviours of gradient descent around critical points of the loss function E ; for an overview of stability analysis, see Section 2.3.1. Consider $\boldsymbol{\theta}^*$ such a critical point, i.e. $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*) = \mathbf{0}$. For a critical point $\boldsymbol{\theta}^*$ to be exponentially asymptotically attractive, all eigenvalues of the Jacobian evaluated at $\boldsymbol{\theta}^*$ need to have strictly negative real part.

The PF has the following Jacobian at critical points (proof in Section A.1.6 in

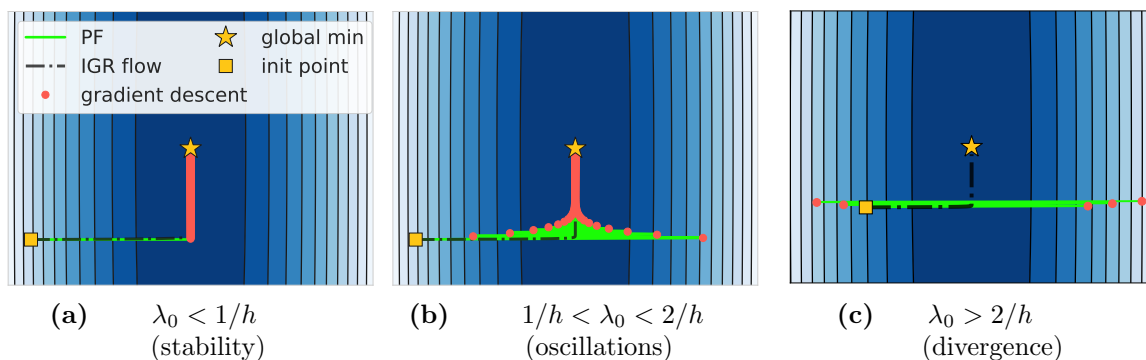


Figure 3.7: Quadratic losses in 2 dimensions. The PF captures the behaviour of gradient descent exactly for quadratic losses, including oscillatory behaviour (b) and divergence (c).

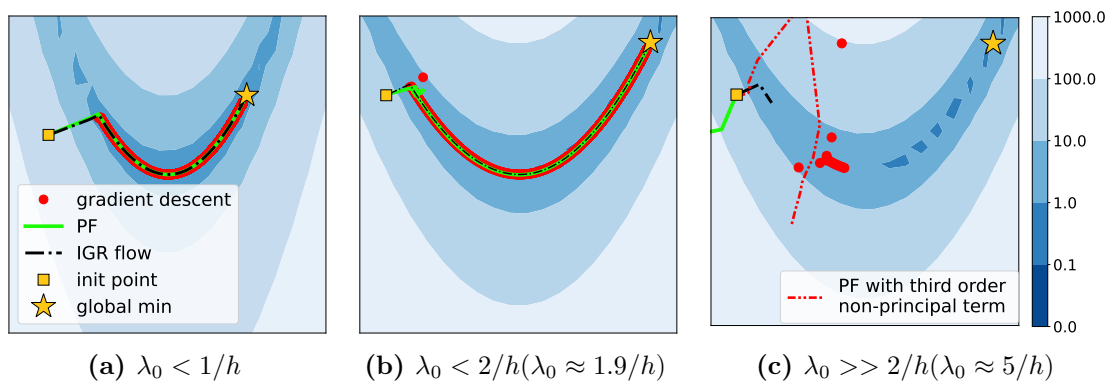


Figure 3.8: Banana function. The PF can capture instability and the gradient descent trajectory over many iterations when λ_0 is close to $2/h$. When $\lambda_0 \gg 2/h$ (c) the PF does not track the GD trajectory over many gradient descent steps, but when including a non-principal term the flow is able to capture the general trajectory of gradient descent and unstable behaviour of gradient descent.

the Appendix):

$$\mathbf{J}_{\text{PF}}(\boldsymbol{\theta}^*) = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i^*)}{h} \mathbf{u}_i^* \mathbf{u}_i^{*T}, \quad (3.13)$$

where λ_i^* , \mathbf{u}_i^* are the eigenvalues and eigenvectors of the Hessian $\nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*)$. We thus have that the eigenvalues of the Jacobian $\mathbf{J}_{\text{PF}}(\boldsymbol{\theta}^*)$ at the critical point $\boldsymbol{\theta}^*$ are $\frac{1}{h} \log(1 - h\lambda_i^*)$ for $i = 1, \dots, D$.

Local minima. Suppose that $\boldsymbol{\theta}^*$ is a local minimum. Then all Hessian eigenvalues are non-negative $\lambda_i^* \geq 0$. We perform the stability analysis in cases given

by the value of λ_i^* , corresponding to the cases in Section 3.3.1:

- $h < 1/\lambda_i^*$. The corresponding eigenvalue of the Jacobian $\frac{1}{h} \log(1 - h\lambda_i^*)$ is negative, since $0 < 1 - h\lambda_i^* < 1$. The PF is attractive in the corresponding eigenvector direction.
- $h \in [1/\lambda_i^*, 2/\lambda_i^*)$. The corresponding eigenvalue of the Jacobian $\frac{1}{h} \log(1 - h\lambda_i^*) = \frac{1}{h} \log(h\lambda_i^* - 1) + i\frac{\pi}{h}$ is complex, with negative real part since $h\lambda_i^* - 1 < 1$. The PF is attractive in the corresponding eigenvector direction.
- $h \geq 2/\lambda_i^*$. The corresponding eigenvalue of the Jacobian $\frac{1}{h} \log(1 - h\lambda_i^*) = \frac{1}{h} \log(h\lambda_i^* - 1) + i\frac{\pi}{h}$ is complex, with non-negative real part, since $h\lambda_i^* - 1 \geq 1$. If $h > 2/\lambda_i^*$, the PF is repelled in the corresponding eigenvector direction.

The last case tells us that the PF is not always attracted to local minima, as it is repelled in eigen-directions where $h > 2/\lambda_i^*$. Thus **like gradient descent, the PF can be repelled around local minima for large learning rates**. This is in contrast to the NGF and the IGR flow, which always predict convergence around a strict local minimum: the eigenvalues of the NGF Jacobian are $-\lambda_i^*$, and for the IGR flow the eigenvalues are $-\lambda_i^* - \frac{h^2}{2}\lambda_i^{*2}$, both are strictly negative when λ_i^* is strictly positive. For derivations see Section A.1.7 in the Appendix.

Remark 3.3.4 *For quadratic losses, where the PF is exact, the results above recover the classical gradient descent result for quadratic losses namely that gradient descent converges if $\lambda_0 < 2/h$, otherwise diverges.*

Saddle points. Suppose that θ^* is a strict saddle point. In this case there exists λ_s^* , such that $\lambda_s^* < 0$. We want to analyse the behaviour of the PF in the direction of the corresponding eigenvector \mathbf{u}_s^* . In that case, $\log(1 - h\lambda_s^*) > 0$ which entails that the PF is repelled in the eigendirections of strict saddle points. Note that this is also the case for the NGF since the corresponding eigenvalues of the Jacobian of the NGF would be $-\lambda_s^*$, also positive. Unlike the NGF, however, the subspace of eigendirections that the PF is repelled by can be larger since it includes also eigendirections where $\lambda_i^* > 2/h > 0$.

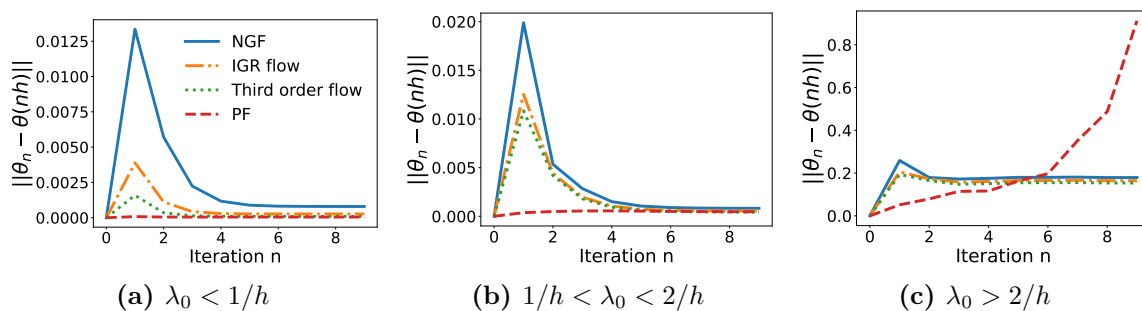


Figure 3.9: Error between gradient descent parameters and parameters obtained following continuous-time flows for multiple iterations: $\|\theta_n - \theta(nh)\|$ with $\theta(0) = \theta_0$. For small n , the PF is better at capturing the behaviour of gradient descent across all cases.

3.4 Predicting neural network gradient descent dynamics with the PF

Computing the PF on large neural networks during training is computationally prohibitive, as it requires finding all eigenvalues of the Hessian matrix once for each step of the flow simulation, corresponding to many eigen-decompositions per gradient descent step. To build intuition about the PF for neural networks, we start with a small MLP for a two-dimensional input regression problem, with random inputs and labels. Here, we can compute the PF’s vector field exactly and compare it with the behaviour of gradient descent. We show results in Figure 3.9, where we visualise the norm of the difference between gradient descent parameters at each iteration and the parameters produced by the continuous-time flows we compare with. We observe that *short term the PF is better than all other flows at tracking the behaviour of gradient descent*. As the number of iterations increases, however, the PF accumulates error in the case of $\lambda_0 > 2/h$; this is likely due to the fact that while gradient descent parameters are real, this is not the case for the PF, as discussed in Remark 3.4.1. Since we are primarily concerned with using the PF to understand gradient descent for a small number of iterations this will be less of a concern in our experimental settings. Additional results that confirm the PF is better than the other flows at tracking gradient descent on a bigger network trained the UCI breast cancer dataset [140] are shown in Figure A.7 in the Appendix.

Remark 3.4.1 Multiple iteration behaviour of the PF. While gradient descent parameters are real for any iteration $\boldsymbol{\theta}_t, \boldsymbol{\theta}_{t+1}, \dots, \boldsymbol{\theta}_{t+n}$ when we approximate the behaviour of gradient descent by initialising $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_t$ and running the PF for time nh , there is nothing enforcing that $\boldsymbol{\theta}(h), \dots, \boldsymbol{\theta}(nh)$ will be real when the PF is complex valued ($\lambda_0 > 1/h$). Outside the quadratic case the Hessian is not Hermitian, the eigenvalues and eigenvectors of the Hessian will not be real, and the eigenvectors are no longer guaranteed to form a basis¹. For long-term trajectories (larger n), this can have an effect on the long-term error between gradient descent and PF trajectories, through an accumulating effect of the imaginary part in the PF. This can be mitigated by using the PF to understand the short-term behaviour of gradient descent (small n).

3.4.1 Predicting $(\nabla_{\boldsymbol{\theta}} E)^T \mathbf{u}_0$ using the PF

For large neural networks, instead of simulating the PF describing how the entire parameter vector changes in time we can use the PF to approximate changes in a scalar quantity only. This will allow us to compare the predictions of the PF against the predictions of the NGF and IGR flow on realistic settings, where the models have millions of parameters. To do so, we first have to compute how the gradient $\nabla_{\boldsymbol{\theta}} E$ changes in time:

Corollary 3.2 *If $\boldsymbol{\theta}$ follows the PF, then:*

$$(\nabla_{\boldsymbol{\theta}} \dot{E}) = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h} (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i. \quad (3.14)$$

This follows from applying the chain rule and using the definition of the PF. We now contrast this result with how the gradient evolves if the parameters follow the NGF or the IGR flow.

¹To avoid the concern around the eigenvectors of the Hessian no longer forming a basis, one can use the Jordan normal form instead, as we will do later, in Section 4.8, where we expand the PF to games. We don't take this approach here as most of our following analysis is not affected, and is concerned with the behaviour of the PF around one gradient descent iteration. Furthermore, support of the Jordan normal form in code libraries is limited (especially for complex matrices), and we did not find this to be a significant issue in the experiments where we simulate the PF outside the quadratic case for a few iterations. We note, however, that mathematical analysis of long-term PF trajectories for general functions should use the Jordan normal form.

Corollary 3.3 *If θ follows the NGF, then:*

$$(\nabla_{\theta} \dot{E}) = \sum_{i=0}^{D-1} -\lambda_i (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i. \quad (3.15)$$

Corollary 3.4 *If θ follows the IGR flow, then:*

$$(\nabla_{\theta} \dot{E}) = \sum_{i=0}^{D-1} - \left(\lambda_i + \frac{h}{2} \lambda_i^2 \right) (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i. \quad (3.16)$$

We would like to use the above to assess how $\nabla_{\theta} E^T \mathbf{u}_i$ changes in time under the above flows and check their predictions empirically against results obtained when training neural networks with gradient descent. Since \mathbf{u}_i is an eigenvector of the Hessian it also changes in time according to the changes given by the corresponding flow, making $(\nabla_{\theta} \dot{E}^T \mathbf{u}_i) = \frac{d(\nabla_{\theta} E^T \mathbf{u}_i)}{dt}$ difficult to calculate. Even when if we wrote an exact flow for $\frac{d(\nabla_{\theta} E^T \mathbf{u}_i)}{dt}$, it would be computationally challenging to simulate it since finding the new values of \mathbf{u}_i would depend on the full Hessian and would lead to the same computational issues we are trying to avoid in the case of large neural networks. In order to mitigate these concerns, we will make the additional approximation that λ_i and \mathbf{u}_i do not change inside an iteration, which will allow us to approximate changes to $\nabla_{\theta} E^T \mathbf{u}_i$ and compare them against empirical observations. We note that we will not use this approximation for any other results.

Remark 3.4.2 *If we assume that λ_i , \mathbf{u}_i do not change between iterations, if θ follows the PF then $\frac{d(\nabla_{\theta} E^T \mathbf{u}_i)}{dt} = \frac{\log(1-h\lambda_i)}{h} \nabla_{\theta} E^T \mathbf{u}_i$.*

Remark 3.4.3 *If we assume that λ_i , \mathbf{u}_i do not change between iterations, if θ follows the NGF we can write $\frac{d(\nabla_{\theta} E^T \mathbf{u}_i)}{dt} = -\lambda_i \nabla_{\theta} E^T \mathbf{u}_i$.*

Remark 3.4.4 *If we assume that λ_i , \mathbf{u}_i do not change between iterations, if θ follows the IGR flow we can write $\frac{d(\nabla_{\theta} E^T \mathbf{u}_i)}{dt} = - \left(\lambda_i + \frac{h}{2} \lambda_i^2 \right) \nabla_{\theta} E^T \mathbf{u}_i$.*

The above flows have the form $\dot{x} = cx$, with solution $x(t) = x(0)e^{ct}$. We can thus test these solutions empirically by training neural networks with gradient descent with learning rate h and at each step compute $\nabla_{\theta} E(\theta_t)^T (\mathbf{u}_i)_{t-1}$ and compare it

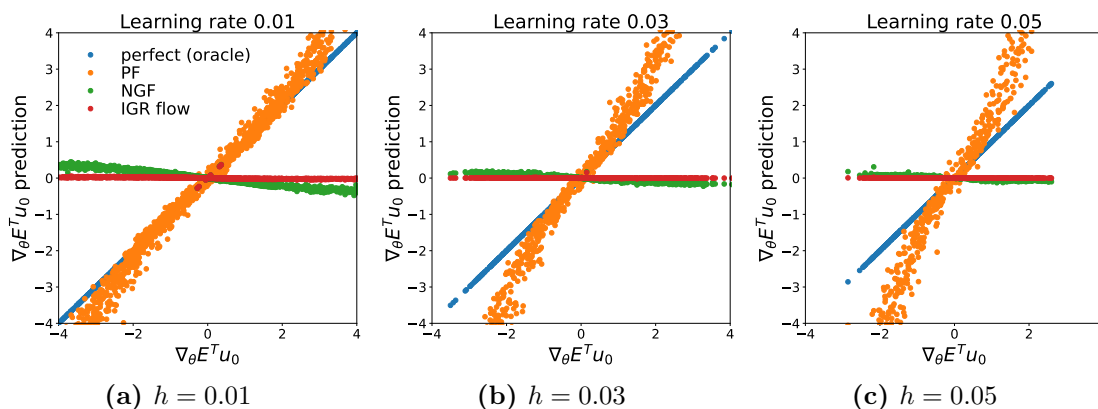


Figure 3.10: Predictions of $\nabla_{\theta} E^T \mathbf{u}_0$ according to the NGF, IGR flow and the PF on full-batch training of a VGG model on CIFAR-10. On the x axis we plot the value of $\nabla_{\theta} E^T \mathbf{u}_0$ as measured empirically in training, and on the y axis we plot the corresponding prediction according to the flows from the value of the dot product at the previous iteration. The ‘exact match’ line indicates a perfect prediction, the upper bound of performance. The PF performs best from all the compared flows, however, for higher learning rates its performance degrades when $\nabla_{\theta} E^T \mathbf{u}_0$ is large; this is due to the fact that the higher the learning rate and the higher the gradient norm, the more likely it is that the additional assumption we used that λ_i, \mathbf{u}_i do not change does not hold.

with the prediction $x(h)$ obtained from the solution from each flow initialised at the previous iteration, i.e. $x(0) = \nabla_{\theta} E(\boldsymbol{\theta}_{t-1})^T(\mathbf{u}_i)_{t-1}$. We show results with a VGG model trained on CIFAR-10 in Figure 3.10. The results show that the PF is substantially better than the NGF and IGR flow at predicting the behaviour of $\nabla_{\theta} E^T \mathbf{u}_0$. Since the NGF and the IGR flow solutions scale the initial value by the inverse of an exponential of magnitude given by λ_0 , for large λ_0 this leads to a small prediction, which is not aligned with what is observed empirically. We also note that the higher the value of $\nabla_{\theta} E^T \mathbf{u}_0$, the worse the prediction of the PF; these are the areas where the approximations made in the above remarks are likely not to hold due to large gradient norms.

3.4.2 Around critical points: escaping sharp local minima and saddles

The stability analysis we performed in Section 3.3.2 showed the PF is repelled by local minima where $\lambda_0^* > 2/h$: that is, even if the model is close to a sharp local minima

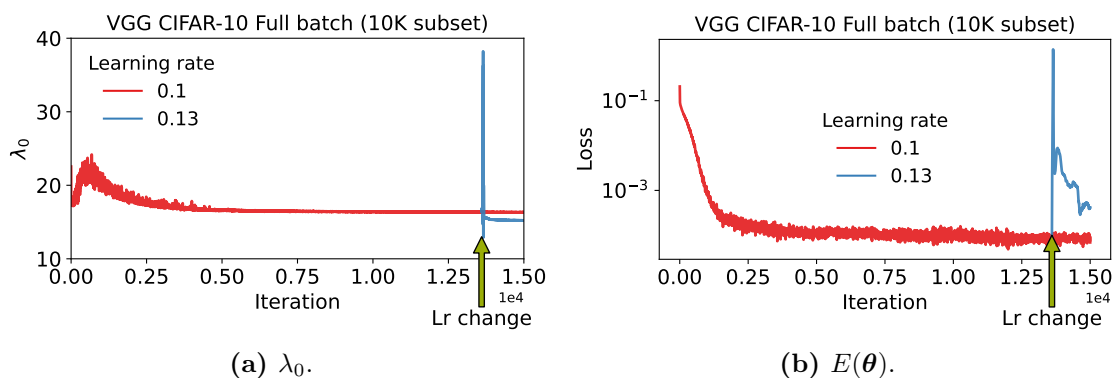


Figure 3.11: Local minima are not attractive to gradient descent if $\lambda_0^* > 2/h$. Since we cannot easily find local minima in neural network landscapes, we perform an experiment where we train a neural network with learning rate $h' < h$ close to convergence, and then change the learning rate to h , close to but above $2/\lambda_0^*$. Increasing the learning rate slightly leads to instabilities and exiting the stable training area, despite being in a late stage of training. We note that while this is *not* the same as assessing local attraction in a stability sense, since we do not know if the point of the learning rate change is in the neighbourhood provided by the existence conditions of local stability, it is the closest practical approximation that we can provide given the complexity of neural network landscapes.

(with $\lambda_0^* > 2/h$), that local minima will not be attractive and training will continue until a shallow minima is reached. We provide experimental evidence to support that hypothesis in neural network training in Figure 3.11. We train a neural network close to convergence, and then increase the learning rate slightly to a learning rate that exceeds the sharpness threshold $h > 2/\lambda_0^*$, and observe that the model promptly exits the area, despite being in a late stage of training. These results are consistent with observations in the deep learning literature [52, 141]. Furthermore, while saddle points have long been considered a challenge with high dimensional optimisation [103] in practice gradient descent has not been observed to converge to saddles [142]. Our analysis suggests that saddles will be repelled not only in the direction of strictly negative eigenvalues, but also in the eigendirections with large positive eigenvalues when large learning rates are used; this can explain why neural networks do not converge to non-strict saddles which exist in deep neural landscapes [143] but need not be repelling for the NGF and IGR flow (existing analyses of escaping saddle points by gradient descent apply only to strict saddles [51, 142]).

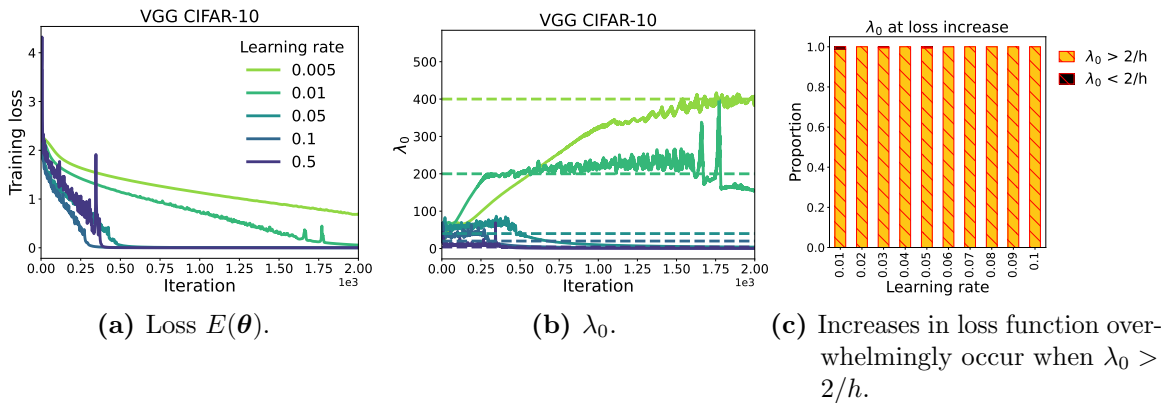


Figure 3.12: Edge of stability in neural networks. Cohen et al. [52] observed that when training neural networks with full-batch gradient descent, training is stable and λ_0 increases until $\lambda_0 > 2/h$, when local loss instabilities occur and λ_0 starts oscillating around $2/h$.

3.5 The PF, stability coefficients and edge of stability results

Edge of stability results. Cohen et al. [52] did a thorough empirical study to show that when training deep neural networks with full-batch gradient descent the largest eigenvalue of the Hessian, λ_0 , keeps growing until reaching approximately $2/h$ (a phase of training they call *progressive sharpening*), after which it remains in that area; for mean-squared losses this continues indefinitely while for cross entropy losses they show it further decreases later in training. They also show that instabilities in training occur when $\lambda_0 > 2/h$. Their empirical study spans neural architectures, data modalities, and loss functions. We visualise the edge of stability behaviour they observe in Figure 3.12; since we use a cross entropy loss λ_0 decreases later in training. We also visualise that iterations where the loss increases compared to the previous iteration overwhelmingly occur when $\lambda_0 > 2/h$. Cohen et al. [52] also empirically observe that $\theta^T \mathbf{u}_0$ has oscillatory behaviour in the edge of stability area but is 0 or small outside it.

Continuous-time models of gradient descent at edge of stability. To investigate if existing continuous-time flows and the PF capture gradient descent behaviour at the edge of stability we train a 5 layer MLP on the toy UCI Iris

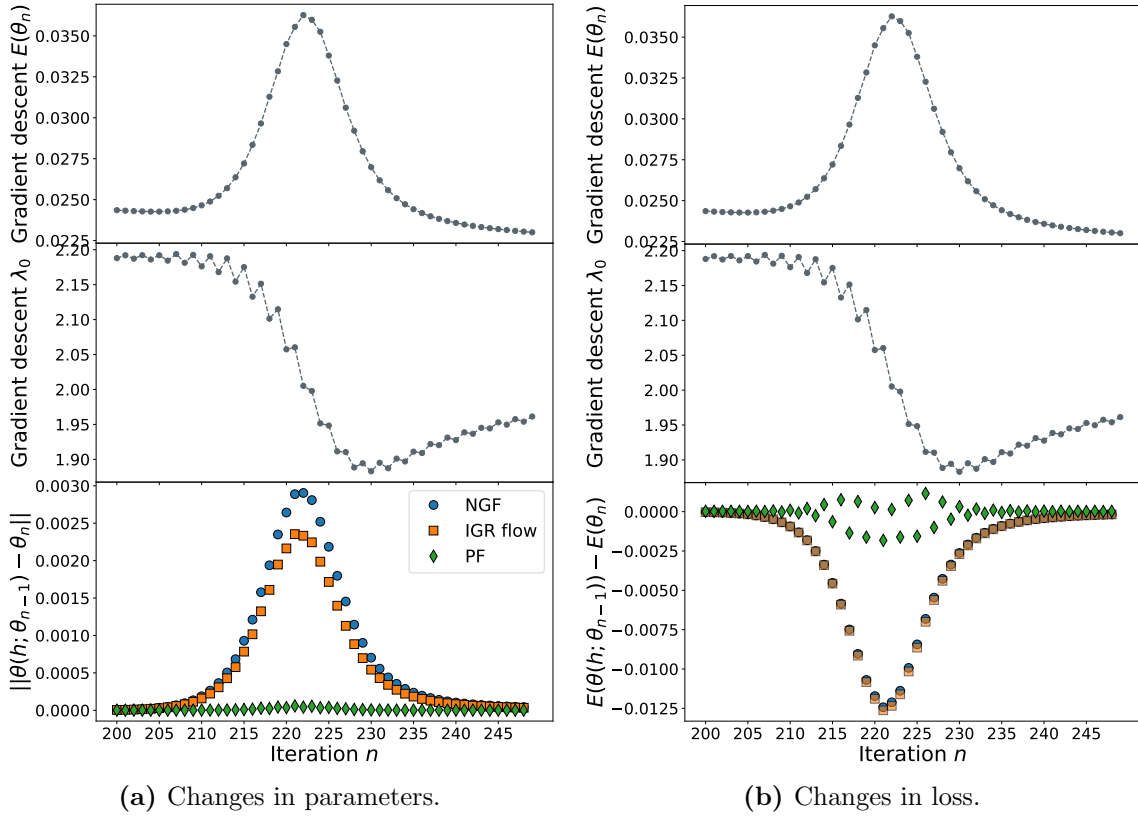


Figure 3.13: Comparing different continuous-time models of gradient descent at the edge of stability area on a small 5 layer MLP, with 10 units per layer. We show the local parameter prediction error $\|\theta_n - \theta(h; \theta_{n-1})\|$ for the NGF, IGR and PF flows (a), as well as the loss prediction error $E(\theta(h; \theta_{n-1})) - E(\theta_n)$ (b); when plotting the loss, for the PF we take the real part of the loss value.

dataset [140]; this simple setting allows for the computation of the full eigenspectrum of the Hessian. We show results in Figure 3.13: the NGF and IGR flow have a larger error compared to the PF when predicting the parameters at the next gradient descent iteration in the edge of stability regime; the NGF and IGR flow predict the loss will decrease, while the PF captures the loss increase observed when following gradient descent. As we remarked in Section 3.2, the NGF and the IGR flow do not capture instabilities when the eigenvalues of the Hessian are positive, which has been remarked to be largely the case for neural network training through empirical studies [144–146] and we observe here (Figure A.8 in the Appendix). We spend the rest of the section using the PF to understand and model edge of stability phenomena using a continuous-time approach.

Connection with the principal flow: stability coefficients. The PF captures the key quantities observed in the edge of stability phenomenon: the eigenvalues of the Hessian λ_i and the threshold $2/h$. These quantities appear in the PF via the stability coefficient $sc_i = \frac{\log(1-h\lambda_i)}{h\lambda_i} \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i = \alpha_{PF}(\lambda_i h) \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i$ of eigendirection \mathbf{u}_i . Through the PF, by connecting the case analysis in Section 3.3.1 with existing and new empirical observations, we can shed light on the edge of stability behaviour in deep learning.

First phase of training (progressive sharpening): $\lambda_0 < 2/h$. This entails $\Re[sc_i] = \Re[\alpha_{PF}(h\lambda_i)] \leq 0, \forall i$ (Real stable and complex stable cases of the analysis in Section 3.3.1). $\text{sign}(\alpha_{NGF}) = \text{sign}(\alpha_{PF}) = -1$ and following the PF minimises E or its real part (Eq (3.12)). To understand the behaviour of λ_0 , we now have to make use of empirical observations about the behaviour of the NGF early in the training of neural networks. It has been empirically observed that in early areas of training, λ_0 increases here when following the NGF [52]; we further show this in Figure A.25 in the Appendix. Since in this part of training gradient descent follows closely the NGF, it exhibits similar behaviour and λ_0 increases. We show this case in Figure 3.14a.

Second phase of training (edge of stability) $\lambda_0 \geq 2/h$. This entails $\Re[sc_0(\boldsymbol{\theta})] = \Re[\alpha_{PF}(h\lambda_i)] \geq 0$. (Unstable complex case of the analysis in Section 3.3.1). We can no longer say that following the PF minimises E . $\text{sign}(\alpha_{NGF}(h\lambda_0)) \neq \text{sign}(\Re[(\alpha_{PF}(h\lambda_0))])$, since $\alpha_{NGF}(h\lambda_0) = -1$ and $\text{sign}(\Re[(\alpha_{PF}(h\lambda_0))]) > 0$ meaning that in that direction gradient descent resembles the positive gradient flow $\dot{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}} E$ rather than the NGF. The positive gradient flow component can cause instabilities, and the strength of the instabilities depends on the stability coefficient $sc_0 = \alpha_{PF}(h\lambda_0) \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0$. We show in Figures 3.14b and 3.16 how the behaviour of the loss and λ_0 are affected by the behaviour of the positive gradient flow when $\lambda_0 > 2/h$.

More than λ_0 : the importance of stability coefficients. While the sign of the real part of the stability coefficient sc_0 is determined by λ_0 , its magnitude is modulated by the dot product $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0$, since $sc_0 = \alpha_{PF}(h\lambda_0) \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0$. The

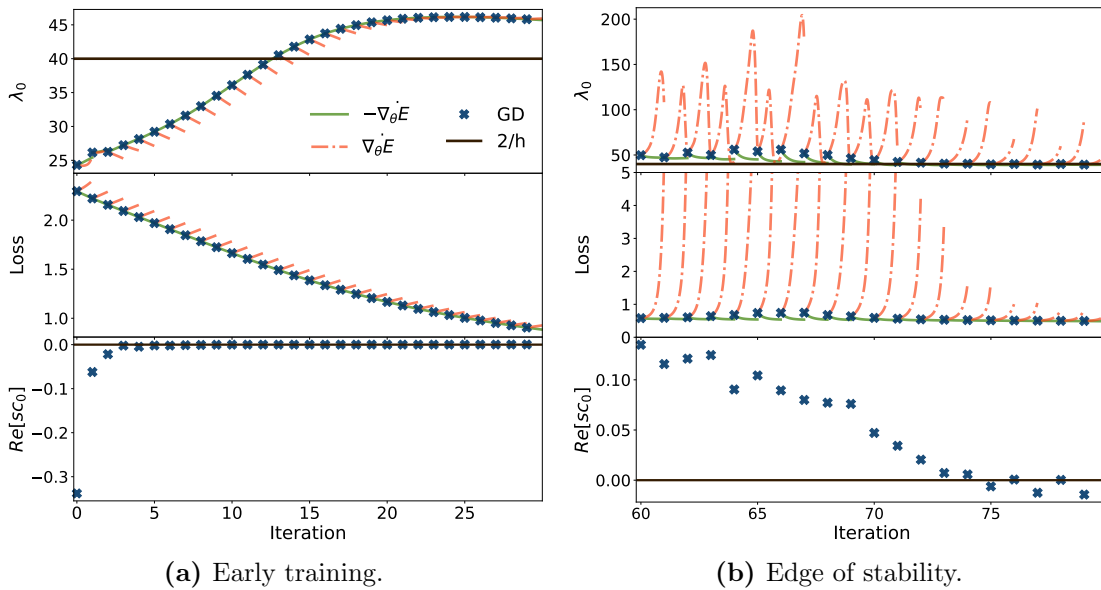


Figure 3.14: Understanding the edge of stability results using the PF on a 4 layer MLP: we plot the behaviour of the NGF $\dot{\theta} = -\nabla_{\theta}E$ and the positive gradient flow $\dot{\theta} = \nabla_{\theta}E$ initialised at each gradient descent iteration parameters, and see that the behaviour of gradient descent is connected to the behaviour of the respective flow through the stability coefficient $sc_0(\theta) = \alpha_{PF}(h\lambda_0)\nabla_{\theta}E^T\mathbf{u}_0$. (a) shows that even when $\lambda_0 > 2/h$, if the real part of the stability coefficient sc_0 is negative or close to 0, there are no instabilities in the loss and the eigenvalue λ_0 keeps increasing, as it does when following the NGF in that region.

magnitude of $\nabla_{\theta}E^T\mathbf{u}_0$ plays an important role, since if λ_0 is the only eigenvalue greater than $2/h$ training is stable if $\nabla_{\theta}E^T\mathbf{u}_0 = 0$, as we observe in Figure 3.14. *To understand instabilities, we have to look at stability coefficients, not only eigenvalues.* We show in Figure 3.15 how the instabilities in training can be related with the stability coefficient sc_0 : the increases in loss occur when the corresponding $\Re[sc_0]$ is positive and large. In Figure 3.16, we show results with the behaviour of λ_0 : λ_0 increases or decreases based on the behaviour of the corresponding flow and the strength of the stability coefficient and that gets reflected in instabilities in the loss function; specifically when $\lambda_0 > 2/h$, we use the positive gradient flow and see how the strength of its fluctuations affect the changes both in the loss value and λ_0 of gradient descent. We show additional results in Figures A.11 and A.12 in the Appendix.

Is one eigendirection enough to cause instability? One question that

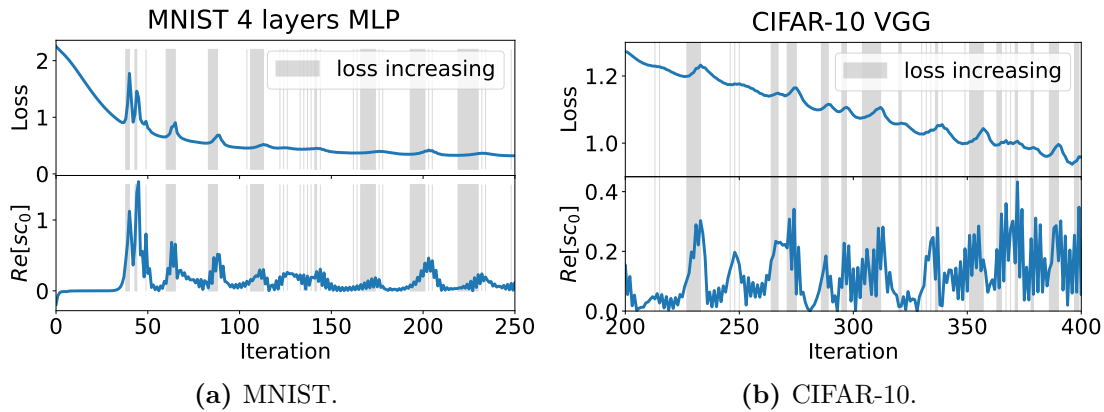


Figure 3.15: The value of the loss function and stability coefficients: areas where the loss increases corresponds to areas where sc_0 is large. The highlighted areas correspond to regions where the loss increases.

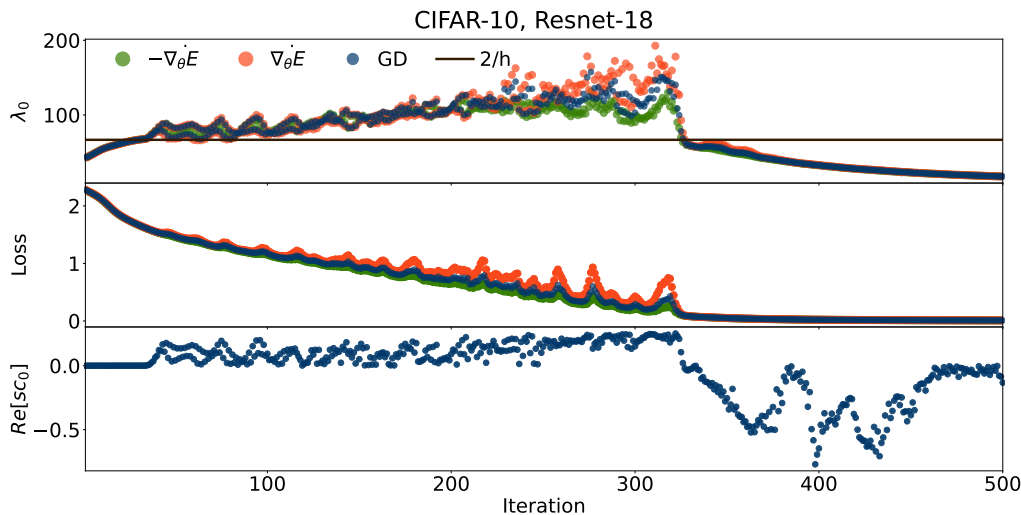


Figure 3.16: Loss instabilities, λ_0 and stability coefficients for CIFAR-10. Together with the behaviour of gradient descent, we plot the behaviour of the NGF and positive gradient flow initialised at θ_t and simulated for time h for each iteration t . The analysis we performed based on the PF suggests that when $\Re[sc_0] > 0$ and large we should expect gradient descent to exhibit behaviours close to those of the positive gradient flow. What we observe empirically is that increases in loss value of gradient descent are proportional to the increase of the positive gradient flow in that area (can be seen best between iterations 200 and 350); the same behaviour can be seen in relation to the eigenvalue λ_0 .

arises from the PF is whether the leading eigendirection \mathbf{u}_0 can be sufficient to cause instabilities, especially in the context of deep networks with millions of parameters. To assess this we train a model with gradient descent until it reaches

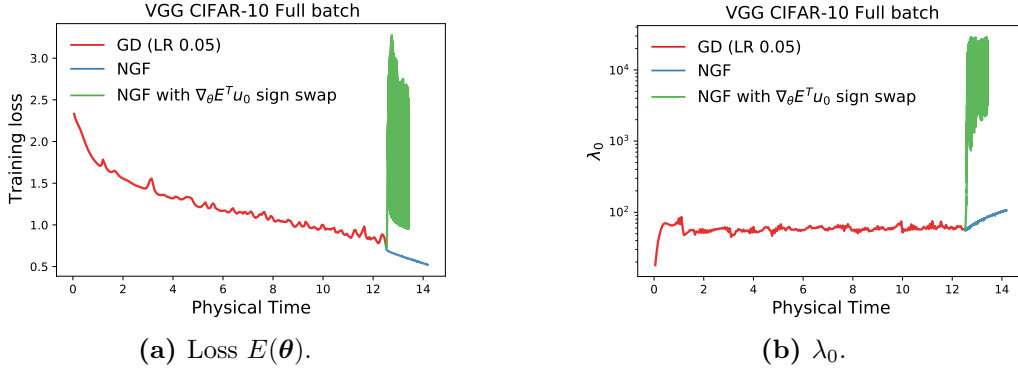


Figure 3.17: One eigendirection is sufficient to lead to instabilities. To create a situation similar to that of the PF in neural network training, we construct a flow given by the NGF in all eigendirections but \mathbf{u}_0 ; in the direction of \mathbf{u}_0 , we change the sign of the flow’s vector field. This leads to the flow $\dot{\boldsymbol{\theta}} = (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0) \mathbf{u}_0 + \sum_{i=1}^{D-1} -(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$. We show this flow can be very unstable when initialised in an edge of stability area, with increases in loss (a) and λ_0 (b).

the edge of stability ($\lambda_0 \approx 2/h$), after which we simulate the continuous flow $\dot{\boldsymbol{\theta}} = (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0) \mathbf{u}_0 + \sum_{i=1}^{D-1} -(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$. The coefficients of the modified vector field of this flow are negative for all eigendirections except from \mathbf{u}_0 , which is positive; this is also the case for the PF when λ_0 is the only eigenvalue greater than $2/h$. In Figure 3.17 we empirically show that a positive coefficient for \mathbf{u}_0 can be responsible for an increase in loss value and a significant change in λ_0 in neural network training.

Decreasing the learning rate. Cohen et al. [52] show that if the edge of stability behaviour is reached and the learning rate is decreased, the training stabilises and λ_0 keeps increasing; we visualise this phenomenon in Figure 3.18. The PF tells us that decreasing the learning rate entails going from $\Re[sc_0] \geq 0$ to $\Re[sc_0] \leq 0$ since $\lambda_0 < 2/h$ after the learning rate change. Since all stability coefficients are now negative, this reduces instability. The increase in λ_0 is likely due to the behaviour of the NGF in that area (as can be seen in Figure 3.17 when changing from gradient descent training to the NGF in an edge of stability area leads to an increase of λ_0).

The behaviour of $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0$. The PF also allows us to explain the unstable behaviour of $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0$ around edge of stability areas. As done in Section 3.4.1, we assume that λ_i, \mathbf{u}_i do not change substantially between iterations and write $\frac{d(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)}{dt} = \frac{\log(1-h\lambda_i)}{h} \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i$ under the PF, with solution

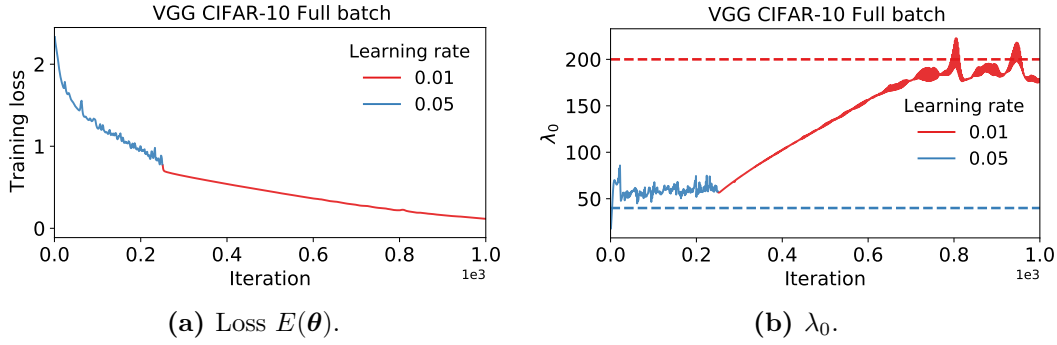


Figure 3.18: Decreasing the learning rate in the edge of stability area leads to increased stability, as discovered by Cohen et al. [52]. The lack of loss instability with a smaller learning rate (a) is explained by the PF now being in a situation where all stability coefficients are negative, and the increase in eigenvalue λ_0 (b) is consistent with what occurs when the PF more closely follows the NGF, since the NGF tends to increase eigenvalues in that part of training, as observed in Figure 3.17.

$(\nabla_{\theta} E^T \mathbf{u}_i)(t) = (\nabla_{\theta} E^T \mathbf{u}_i)(0) e^{\frac{\log(1-h\lambda_i)}{h} t}$. This solution has different behaviour depending on the value of λ_0 relative to $2/h$: decreasing below $2/h$ and increasing above $2/h$. We show this theoretically predicted behaviour in Figure 3.19, alongside empirical behaviour showcasing the fluctuation of $\nabla_{\theta} E^T \mathbf{u}_0$ in the edge of stability area, which confirms the theoretical prediction. We also compute the prediction error of the proposed flow and show it can capture the dynamics of $\nabla_{\theta} E^T \mathbf{u}_0$ closely in this setting. We present a discrete-time argument for this observation in Section A.2.2. We note that the stable behaviour early in training together with the oscillatory behaviour of $\nabla_{\theta} E^T \mathbf{u}_0$ in the edge of stability area, which we predict and observe can explain the results of Cohen et al. [52] on the behaviour of $\theta^T \mathbf{u}_0$, since θ accumulates changes given by gradient updates.

Why not more instability? To determine why there isn't more instability in the edge of stability area we have to consider that neural networks are not quadratic, which has two effects. Firstly, when following the PF the landscape changes slightly locally; this leads to changes in stability coefficients and thus the behaviour of gradient descent as we have consistently seen in the experiments in this section. Secondly, non-principal terms can have an effect; while we do not know all non-principal terms in Section 3.6 we provide a justification for why the non-principal term we do know

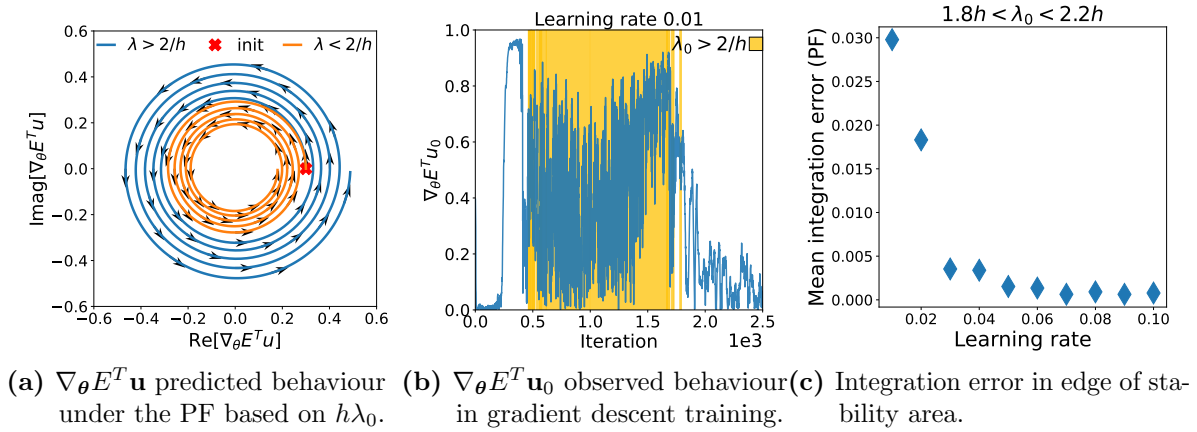


Figure 3.19: Predicting the unstable dynamics of $\nabla_{\theta} E^T \mathbf{u}$ in the edge of stability area ($\lambda \approx 2/h$) using the PF. (a): the predicted behaviour of $\nabla_{\theta} E^T \mathbf{u}$ under $\frac{d(\nabla_{\theta} E^T \mathbf{u}_i)}{dt} = \frac{\log(1-h\lambda_i)}{h} \nabla_{\theta} E^T \mathbf{u}_i$, with an inflection point at $\lambda = 2/h$. (b): empirical behaviour of $\nabla_{\theta} E^T \mathbf{u}$ for a model shows instabilities in the edge of stability area (highlighted). (c): the approximation made to derive the flow is suitable around $\lambda \approx 2/h$. (b) and (c) plots are obtained from training VGG models on CIFAR-10 with full-batch training.

(Eq (3.11)) can have a stabilising effect by inducing a regularisation pressure to minimise $\lambda_i(\nabla_{\theta} E^T \mathbf{u}_i)^2$ in certain parts of the training landscape.

In this section we have shown the PF closely predicts the behaviour of gradient descent in neural network training. This has led to additional insights, including the importance of stability coefficients in determining instabilities in gradient descent (Figures 3.14, 3.15, 3.16), causally showing one eigendirection is sufficient to cause instability (Figure 3.17) and being able to closely predict the behaviour of the dot product between the gradient and the largest eigenvector (Figure 3.19). This evidence suggests that the PF captures significant aspects of the behaviour of gradient descent in deep learning; this is likely due to the specific structure of neural network models. While we take a continuous-time approach, a discrete-time approach can be used to motivate some of our observations (Section A.2); this is complementary to our approach but nonetheless related, since it also does not account for higher order derivatives of the loss and further suggests the strength of a quadratic approximation of the loss in the case of neural networks, as observed by Cohen et al. [52].

3.6 Non-principal terms can stabilise training

This chapter focuses on understanding the effects of the PF on the behaviour of gradient descent. The principal terms, however, are not the only terms in the discretisation drift: we have found one non-principal term of the form $\nabla_{\theta} E^T (\nabla_{\theta}^3 E) \nabla_{\theta} E$ (Eq (3.4)) and have seen that it can have a stabilising effect (Figure 3.8). To provide some intuition, while $\nabla_{\theta} E^T (\nabla_{\theta}^3 E) \nabla_{\theta} E$ cannot be written as a gradient operator we can write:

$$\nabla_{\theta} E^T \nabla_{\theta}^3 E \nabla_{\theta} E = \sum_{i=0}^{D-1} (\mathbf{u}_i^T \nabla_{\theta}^3 E \mathbf{u}_i) (\nabla_{\theta} E^T \mathbf{u}_i)^2 \approx \sum_{i=0}^{D-1} \nabla_{\theta} \lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2, \quad (3.17)$$

where we used that \mathbf{u}_i does not change substantially around a GD iteration to write the derivative of λ_i , namely $\nabla_{\theta} \lambda_i = \mathbf{u}_i^T \nabla_{\theta}^3 E \mathbf{u}_i$, since $\nabla_{\theta}^2 E$ is real and symmetric around a gradient descent iteration [147]. We can now write, again assuming \mathbf{u}_i is locally constant:

$$\sum_{i=0}^{D-1} \nabla_{\theta} \lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2 = \sum_{i=0}^{D-1} \nabla_{\theta} (\lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2) - \sum_{i=0}^{D-1} 2\lambda_i^2 (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i. \quad (3.18)$$

From here we can see that if λ_i or $\nabla_{\theta} E^T \mathbf{u}_i$ are close to 0, the non-principal term leads to a pressure to minimise $\lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2$, since the non-gradient term $\lambda_i^2 (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i$ in Eq (3.18) diminishes. This creates an incentive to keep the value of $\lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2$ around 0. We note that the work of Damian et al. [148] is what inspired us to write the third-order non-principal term in the form of Eq (3.18), after we had previously noted its stabilising properties.

The incentive to keep $\lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2$ close to 0 can have a stabilising effect, since $\lambda_i = 0$ or $\nabla_{\theta} E^T \mathbf{u}_i = 0$ results in $\alpha_{PF}(h\lambda_i) = \alpha_{NGF}(h\lambda_i) = -1$. Thus, minimising $\lambda_i (\nabla_{\theta} E^T \mathbf{u}_i)^2$ can reduce instability from the PF, since in those directions the PF has the same behaviour as the NGF. We think this can shed light on observations about neural network training behaviour: the NGF tends to increase the eigenvalues, but as they reach close to 0 there is a pressure to minimise the above, which leads to λ_i staying around 0; this is consistent with observations in the literature [144–146]. This

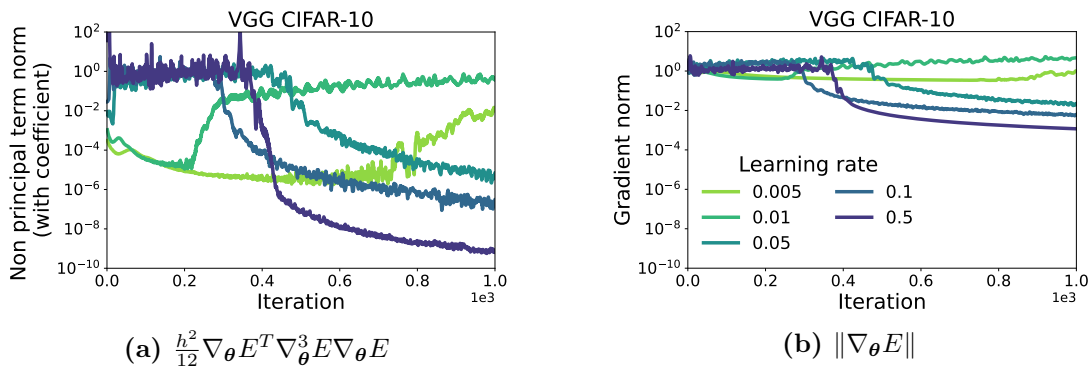


Figure 3.20: The value of the non-principal third-order term in training: outside the edge of stability areas (best seen for learning rates 0.005 and 0.01), the non-principal third-order term is very small compared to the gradient. Inside the edge of stability areas, the non-principal term has a higher magnitude. Corresponding loss and λ_0 plots are provided in Figure 3.12.

observation also explains the pressure for the dot product $\nabla_{\theta} E^T \mathbf{u}_0$ to stay small if it is initialised around 0, as we see is the case in the early in neural network training (see Figures 3.14a and 3.15). Furthermore, as we show in Figure 3.19 and previously argued, the dot product $\nabla_{\theta} E^T \mathbf{u}_0$ fluctuates around 0 in the edge of stability areas, making it likely for the minimisation effect induced by Eq (3.18) to kick in.

We plot the value of this non-principal term in neural network training in Figure 3.20. These results show that the non-principal third-order term is very small outside the edge of stability area, but has a larger magnitude around the edge of stability, where the largest eigenvalues stop increasing but the magnitude of $(\nabla_{\theta} E(\theta_0)^T \mathbf{u}_i)^2$ fluctuates and can be large (Figure 3.19). This observation is consistent with the above interpretation, but more theoretical and empirical work is needed to understand the effects of non-principal terms on training stability and generalisation.

3.7 Stabilising training by accounting for discretisation drift

The PF allows us to understand not only how gradient descent differs from the trajectory given by the NGF, but also when they follow each other closely. Understanding when gradient descent behaves like the NGF flow reveals when the existing analyses

of gradient descent using the NGF discussed in Section 3.2 are valid. It also has practical implications, since, in areas where gradient descent follows the NGF closely, training can be sped up by increasing the learning rate. Prior works have empirically observed that gradient descent follows the NGF early in neural network training [52], and this observation can be used to explain why decaying learning rates [149] or learning rate warm up [150] are successful when training neural networks: having a high learning rate in areas where the drift is small will not cause instabilities and can speed up training while decaying the learning rate avoids instabilities later in training when the drift is larger.

3.7.1 $\nabla_{\theta}^2 E \nabla_{\theta} E$ determines discretisation drift

In previous sections we have seen that the Hessian plays an important role in defining the PF and in training instabilities. We now want to quantify the difference between the NGF and the PF in order to understand when the NGF can be used as a model of gradient descent. We find that:

Remark 3.7.1 *In a region of the space where $\nabla_{\theta}^2 E \nabla_{\theta} E = \mathbf{0}$ the PF is the same as the NGF.*

To see why, we can expand

$$\nabla_{\theta}^2 E \nabla_{\theta} E = \sum_{i=0}^{D-1} \lambda_i \nabla_{\theta} E^T \mathbf{u}_i \mathbf{u}_i. \quad (3.19)$$

If $\nabla_{\theta}^2 E \nabla_{\theta} E = \mathbf{0}$ we have that $\lambda_j \nabla_{\theta} E^T \mathbf{u}_j = 0, \forall j \in \{1, \dots, D\}$, thus either $\lambda_j = 0$ leading to $\alpha_{NGF}(h\lambda_j) = \alpha_{PF}(h\lambda_j) = -1$ or $\nabla_{\theta} E^T \mathbf{u}_j = 0$. Then $\dot{\theta} = \sum_{i=0}^{D-1} \alpha_{PF}(h\lambda_i) (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i = \sum_{i=0}^{D-1} \alpha_{NGF}(h\lambda_i) (\nabla_{\theta} E^T \mathbf{u}_i) \mathbf{u}_i$.

Thus comparing the PF with the NGF reveals an important quantity: $\nabla_{\theta}^2 E \nabla_{\theta} E$. Further investigating this quantity reveals it has a connection with the total drift, since:

Theorem 3.7.1 *The discretisation drift (error between gradient descent and the NGF) after one iteration $\theta_t = \theta_t - h \nabla_{\theta} E(\theta_{t-1})$ is $\frac{h^2}{2} \nabla_{\theta}^2 E(\theta') \nabla_{\theta} E(\theta')$ for a set of parameters θ' in the neighborhood of θ_{t-1} .*

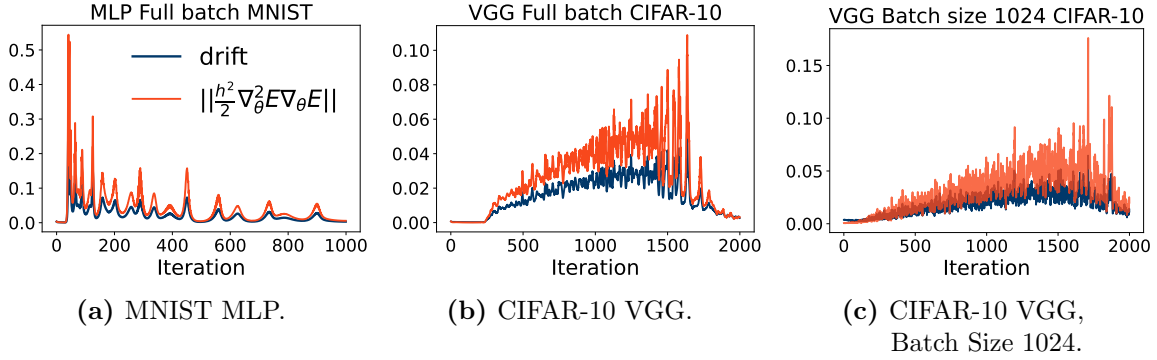


Figure 3.21: $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ and the per-iteration drift have similar behaviour during training. As suggested by Theorem 3.7.1, the norm of the per-iteration drift $\|\theta(h; \theta_{n-1}) - \theta_n\|$ can be approximated by $\|\nabla_{\theta}^2 E(\theta_{n-1}) \nabla_{\theta} E(\theta_{n-1})\|$. The error in the approximation is due to the theorem providing an existence argument for a set of parameters θ' in the neighbourhood of θ_{n-1} such that $\|\nabla_{\theta}^2 E(\theta') \nabla_{\theta} E(\theta')\|$ is equal to the per-iteration drift; since we do not know those parameters we use the iteration parameters θ_{n-1} instead.

This follows from the Taylor reminder theorem in mean value form (proof in Section A.1.9). This leads to:

Corollary 3.5 *In a region of space where $\nabla_{\theta}^2 E \nabla_{\theta} E = \mathbf{0}$ gradient descent follows the NGF.*

Thus the PF revealed $\nabla_{\theta}^2 E \nabla_{\theta} E$ as a core quantity in the discretisation drift of gradient descent. To further see the connection between with the PF consider that $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|^2 = \left\| \sum_{i=0}^{D-1} \lambda_i \nabla_{\theta} E^T \mathbf{u}_i \mathbf{u}_i \right\|^2 = \sum_{i=0}^{D-1} \|\lambda_i \nabla_{\theta} E^T \mathbf{u}_i\|^2$; the higher each term in the sum, the higher the difference between the NGF and the PF. To measure the connection between per-iteration drift and $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ in neural network training we approximate it via $\left\| \theta_t - \widetilde{NGF}(\theta_{t-1}, h) \right\|$ where \widetilde{NGF} is the numerical approximation to the NGF initialised at θ_{t-1} . Results in Figures 3.21 and 3.22 show the strong correlation between per-iteration drift and $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ throughout training and across learning rates. Since Theorem 3.7.1 tells us the form of the drift but not the exact value of θ' , we have used θ_{t-1} instead to evaluate $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ and thus some error exists.

Understanding this connection is advantageous since computing discretisation drift is computationally expensive as it requires simulating the continuous-time NGF

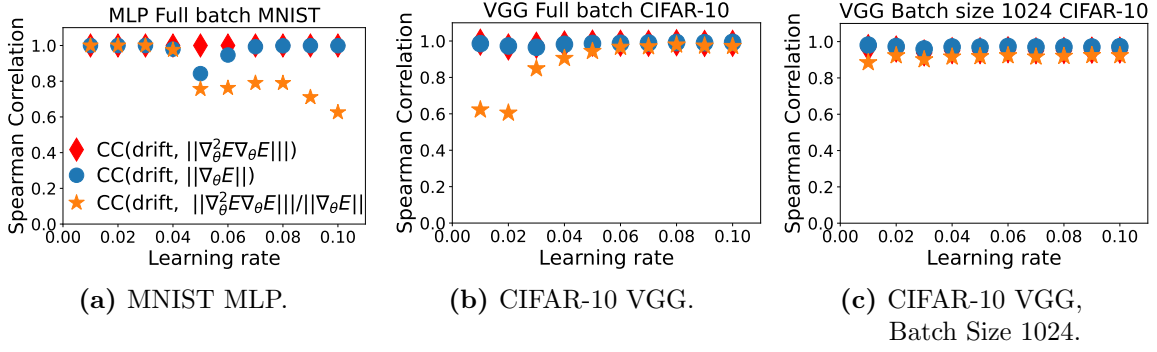


Figure 3.22: Correlation between $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ and the per-iteration drift. Since $\|\nabla_{\theta}^2 E \nabla_{\theta} E\| = (\|\nabla_{\theta}^2 E \nabla_{\theta} E\|) / \|\nabla_{\theta} E\| \|\nabla_{\theta} E\|$, we measure the correlation between per-iteration drift and $(\|\nabla_{\theta}^2 E \nabla_{\theta} E\|) / \|\nabla_{\theta} E\|$ as well as $\|\nabla_{\theta} E\|$. As Theorem 3.7.1 provides an existence argument for an unknown set of parameters θ' in the neighbourhood of θ_{n-1} such that $\|\nabla_{\theta}^2 E(\theta') \nabla_{\theta} E(\theta')\|$ is equal to the per-iteration drift $\|\theta(h; \theta_{n-1}) - \theta_n\|$, we evaluate $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ at parameters θ_{n-1} .

but computing $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ via Hessian-vector products is cheaper and approximations are available, such as $\nabla_{\theta}^2 E \nabla_{\theta} E \approx \frac{\nabla_{\theta} E(\theta + \epsilon \nabla_{\theta} E) - \nabla_{\theta} E(\theta)}{\epsilon}$ which only requires an additional backward pass [151].

3.7.2 Drift adjusted learning (DAL)

A natural question to ask is how to use the correlation between $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ and the per-iteration drift to improve training stability; $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ captures all the quantities we have shown to be relevant to instability highlighted by the PF: λ_i and $\nabla_{\theta} E^T \mathbf{u}_i$ (Eq. (3.19)). One way to use this information is to adapt the learning rate of the gradient descent update, such as using $\frac{2}{\|\nabla_{\theta}^2 E \nabla_{\theta} E\|}$ as the learning rate. This learning rate slows down training when the drift is large—areas where instabilities are likely to occur—and it speeds up training in regions of low drift—areas where instabilities are unlikely to occur. Computing the norm of the update provided by this learning rate shows a challenge, however, since $\frac{2\|\nabla_{\theta} E\|}{\|\nabla_{\theta}^2 E \nabla_{\theta} E\|} \geq \frac{2\|\nabla_{\theta} E\|}{\lambda_0 \|\nabla_{\theta} E\|} = \frac{2}{\lambda_0}$; this implies that when using this learning rate the norm of the gradient descent update will never be 0 and thus training will not result in convergence. Furthermore, the magnitude of the parameter update will be independent of the gradient norm. To

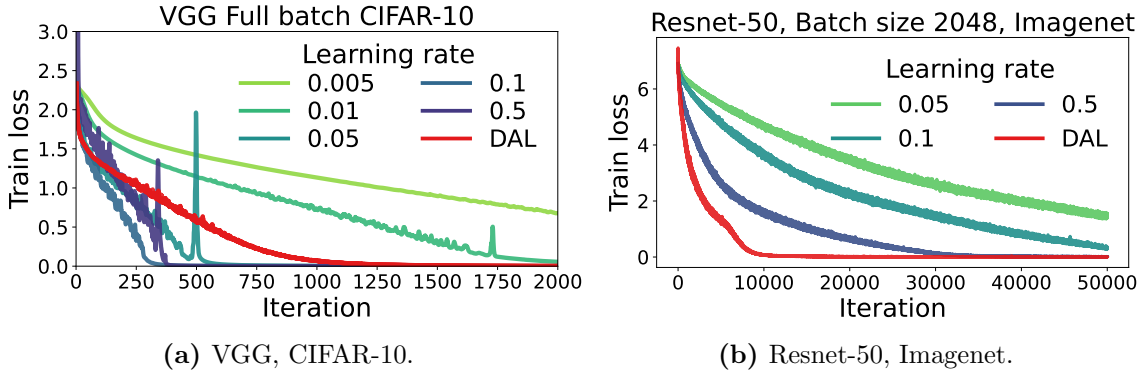


Figure 3.23: DAL: using the learning rate $\frac{2}{\|\nabla_{\theta}^2 E \hat{\mathbf{g}}(\theta)\|}$ results in improved stability without requiring a hyperparameter sweep; Resnet-18 results on CIFAR-10 in Figure A.15 in the Appendix.

reinstate the gradient norm, we propose using the learning rate

$$h(\theta) = \frac{2}{\|\nabla_{\theta}^2 E \nabla_{\theta} E\| / \|\nabla_{\theta} E\|} = \frac{2}{\|\nabla_{\theta}^2 E \hat{\mathbf{g}}(\theta)\|}, \quad (3.20)$$

where $\hat{\mathbf{g}}(\theta)$ is the unit normalised gradient $\nabla_{\theta} E / \|\nabla_{\theta} E\|$. We will call this learning rate **DAL** (Drift Adjusted Learning). As shown in Figure 3.21, $\|\nabla_{\theta}^2 E \hat{\mathbf{g}}(\theta)\|$ has a strong correlation with the per-iteration drift. Another interpretation of DAL can be provided through a signal-to-noise perspective: the size of the learning signal obtained by minimising E is that of the update $h \|\nabla_{\theta} E\|$, while the norm of the noise coming from the drift can be approximated as $\frac{h^2}{2} \|\nabla_{\theta}^2 E \nabla_{\theta} E\|$. Thus, the ‘signal-to-noise ratio’ can be approximated as $h \|\nabla_{\theta} E\| / (\frac{h^2}{2} \|\nabla_{\theta}^2 E \nabla_{\theta} E\|) = 2 / (h \|\nabla_{\theta}^2 E \hat{\mathbf{g}}(\theta)\|)$, which when using DAL (Eq (3.20)) is 1; DAL can be seen as balancing the gradient signal and the regularising drift noise in gradient descent training.

We use DAL to set the learning rate and show results across architectures, models, and datasets in Figure 3.23 (with additional results in Figure A.16 in the Appendix). *Despite not requiring a learning rate sweep, DAL is stable compared to using fixed learning rates.* To provide intuition about DAL, we show the learning rate and the update norm in Figure 3.24: for DAL the learning rate decreases in training after which it slowly increases when reaching areas with low drift. Compared to larger learning static learning rates where the update norm can increase in the edge

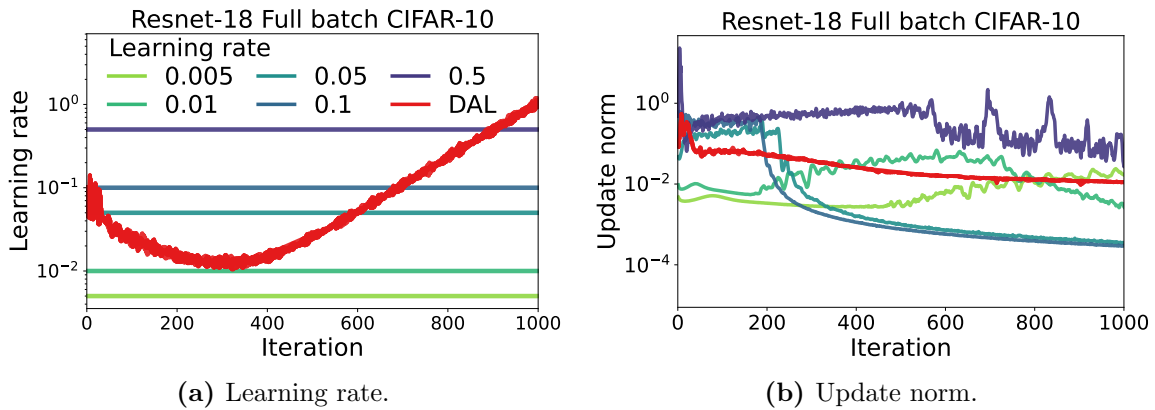


Figure 3.24: Key quantities in DAL versus fixed learning rate training: (a) shows the learning rate $h(\boldsymbol{\theta})$ which is either constant or $\frac{2}{\|\nabla_{\boldsymbol{\theta}}^2 E \hat{\mathbf{g}}(\boldsymbol{\theta})\|}$ for DAL, and (b) shows update norms $h(\boldsymbol{\theta}) \|\nabla_{\boldsymbol{\theta}} E\|$. Unlike using a fixed learning rate, DAL adjusts to the training landscape and even though later in training the learning rate increases, this does not lead to an increase in the update norm.

of stability area with DAL the update norm steadily decreases in training.

3.7.3 The trade-off between stability and performance

Since we are interested in understanding the optimisation dynamics of gradient descent, we have so far focused on training performance. We now try to move our attention to test performance and generalisation. Previous works [12, 105, 152] have shown that higher learning rates lead to better generalisation performance. We now try to further connect this information with the per-iteration drift and DAL. To do so, we use learning rates with various degrees of sensitivity to per-iteration drift using DAL- p :

$$h_p(\boldsymbol{\theta}) = \frac{2}{(\|\nabla_{\boldsymbol{\theta}}^2 E \hat{\mathbf{g}}(\boldsymbol{\theta})\|)^p}. \quad (3.21)$$

The higher p , the slower the training and less drift there is; the lower p , there is more drift. We start with extensive experiments with $p = 0.5$, which we show in Figure 3.25, and show more results in Figure A.17. Compared to $p = 1$ (DAL), there is faster training but at times also more instability. Performance on the test set shows that DAL-0.5 performs as well or better than when using fixed learning rates.

Remark 3.7.2 We find that across datasets and batch sizes, DAL-0.5 performs best

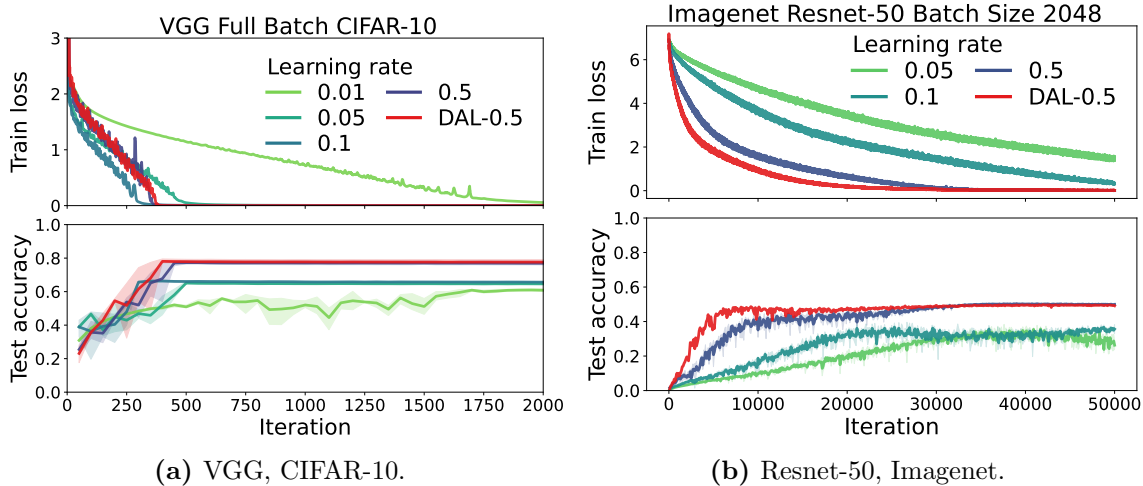


Figure 3.25: DAL-0.5 leads to increased training speed and generalisation compared to a sweep of fixed learning rates. Resnet-18 results on CIFAR-10 in Figure A.15 in the Appendix.

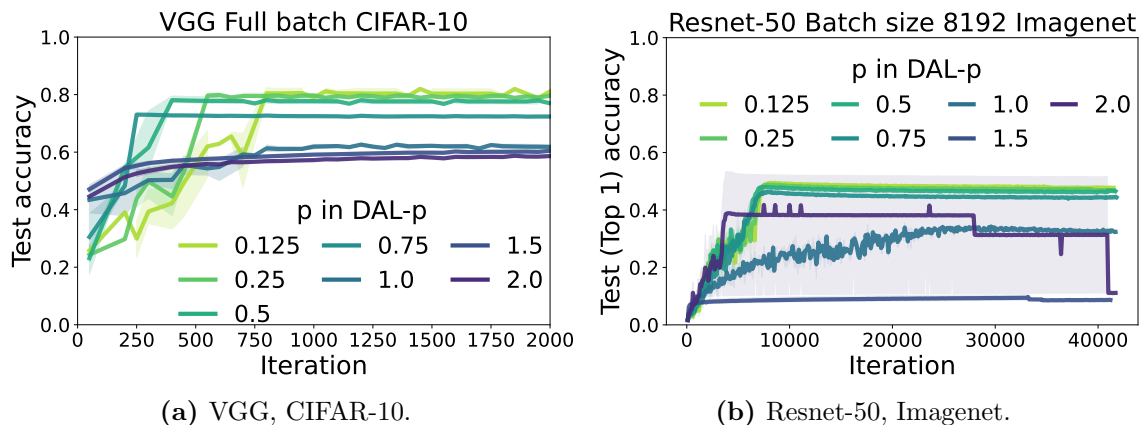
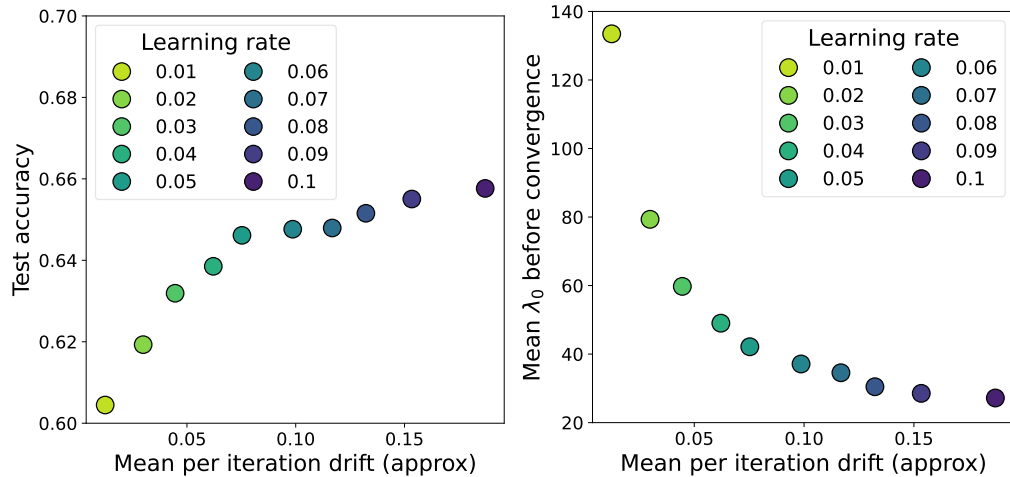


Figure 3.26: DAL- p sweep: discretisation drift helps test performance at the cost of stability. Corresponding training curves and loss functions are present in the Figure A.18 in the Appendix; results showing the same trends across various batch sizes are shown in Figure A.14.

in terms of the stability generalisation trade-off and in these settings it can be used as a drop-in replacement for a learning rate sweep.

To further investigate the connection between drift and test set performance, we perform a set of sweeps over the power p and show results in Figure 3.26. These results show that the higher the drift (the smaller p), the more generalisation; additional results across batch sizes showing the same trend are shown in Figure A.14 in the Appendix. We also show in Figure 3.27 the correlation between mean per-



(a) Fixed learning rate sweep.

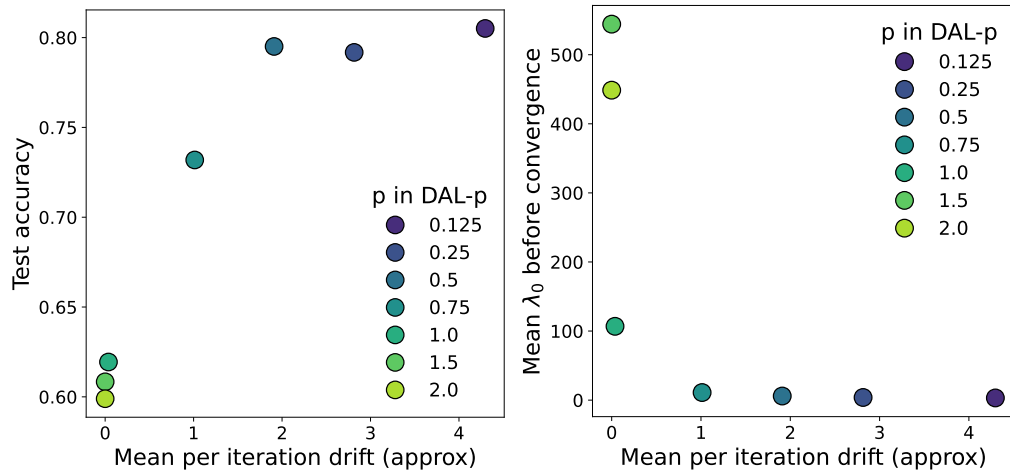
(b) DAL- p sweep.

Figure 3.27: The correlation between per-iteration drift, test set performance and λ_0 in full batch training on CIFAR-10. With fixed learning rates (a), the drift increases as the learning rate increases which leads to improved test performance and lower λ_0 on average in training. When using DAL (b), the drift is controlled through the dynamic learning rate which adjusts to the magnitude of the drift; here too we observe that with high drift leads to increased performance and lower λ_0 on average in training. The same pattern can be seen in results with mini-batch training in Figure A.23.

iteration drift and test accuracy both for learning rate and DAL- p sweeps. The results consistently show that the higher the mean per-iteration drift, the higher the test accuracy. We also show that the mean per-iteration drift has a connection to the largest eigenvalue λ_0 : the higher the drift, the smaller λ_0 . These results add further evidence to the idea that discretisation drift is beneficial for generalisation

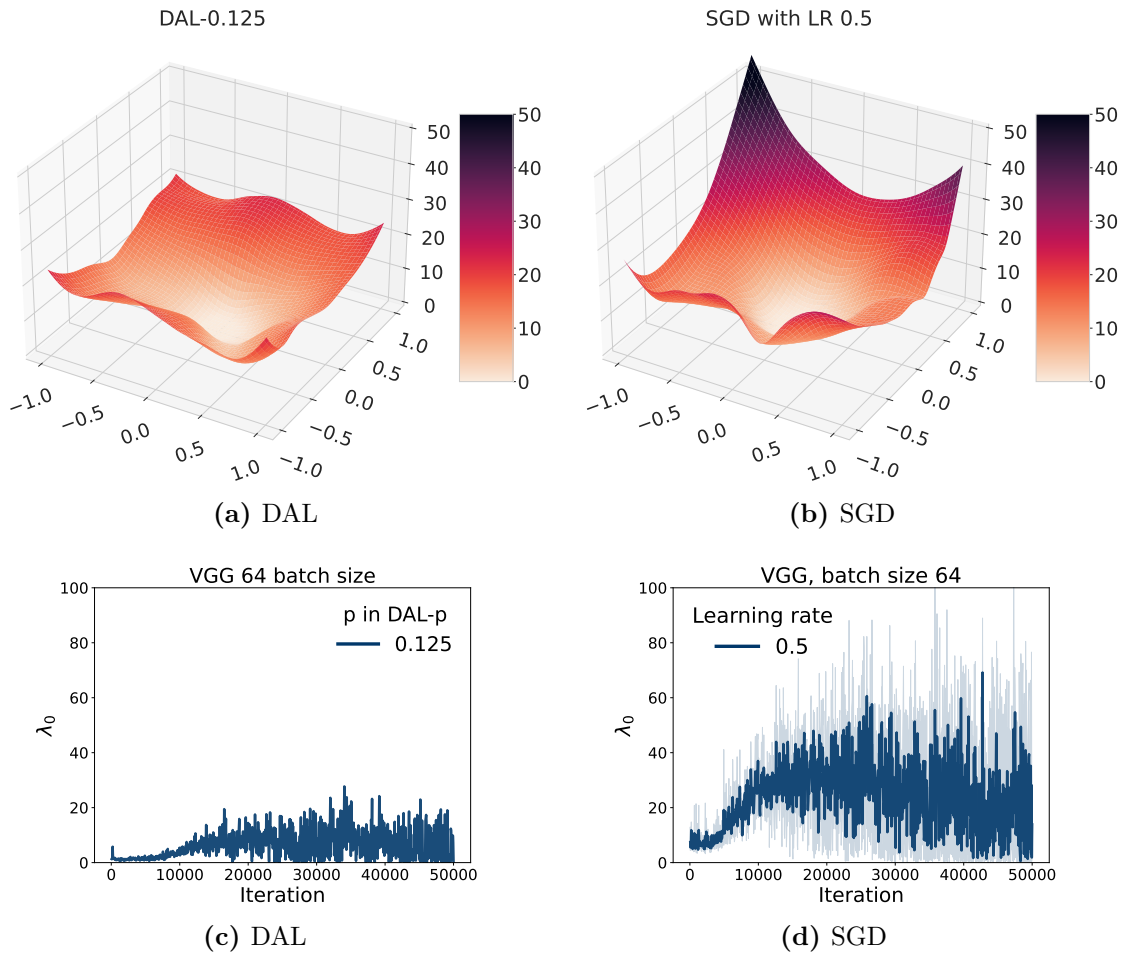


Figure 3.28: CIFAR-10, batch size 64. The 2D projection of the DAL- p (a) and SGD (b) learned landscapes on CIFAR-10 using a VGG model. The visualisation is made using the method of Li et al. [22]. *Though both models achieve a 86% accuracy on the test dataset, DAL- p has learned a flatter model.* Bottom row: We also show the trajectory of λ_0 for both models, which shows that throughout the learning trajectory, DAL (c) travels through flatter landscapes compared to SGD (d), as measured by λ_0 .

performance in deep learning. We also notice that DAL- p with smaller values of p leads to a small λ_0 compared to vanilla gradient descent even when large learning rates are used for the latter; this could explain its generalisation capabilities as lower sharpness has been connected to generalisation in previous works [18, 141, 153]. To consolidate these results, we use the method of Li et al. [22] to visualise the loss landscape learned by DAL- p compared to that learned using gradient descent in Figure 3.28, and observe that *even when reporting similar accuracies, DAL- p converges to a flatter landscape* and traverses a flatter area of space throughout

training. These results are consistent across small and large batch sizes; we show additional results in Figures A.20 and A.22 in the Appendix.

Inspired by understanding when the PF is close to the NGF, in this section we investigated the total discretisation drift of gradient descent. This led us to DAL- p , a method to automatically set the learning rate based on approximation to the per-iteration drift of gradient descent; we have seen that DAL produces stable training and further connected discretisation drift, generalisation, and flat landscapes as measured by leading Hessian eigenvalues.

3.8 Future work

Beyond gradient descent. In this chapter, we focused on understanding vanilla gradient descent. Understanding discretisation drift via the PF can be beneficial for improving other gradient based optimisation algorithms as well, as we briefly illustrate for momentum updates with decay β and learning rate h (for a discussion on momentum, see Section 2.2):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}); \quad \boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + \mathbf{v}_t. \quad (3.22)$$

We can scale $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$ in the above not by a fixed learning rate h , but by adjusting the learning rate according to the approximation to the drift. This has two advantages: it removes the need for a learning rate sweep and it uses local landscape information in adapting the moving average, such that in areas of large drift the contribution is decreased, while it is increased in areas where the drift is small (a more formal justification is provided in Section A.1.9). This leads to the following updates:

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} - \frac{1}{2 \|\nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}_{t-1}) \hat{\mathbf{g}}(\boldsymbol{\theta})(\boldsymbol{\theta}_{t-1})\|} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}); \quad \boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + \mathbf{v}_t. \quad (3.23)$$

As with DAL- p , we can use powers to control the stability performance trade-off: the lower p , the more the current update contribution is reduced in high drift (instability) areas. We tested this approach on Imagenet and show results in Figure 3.29. The

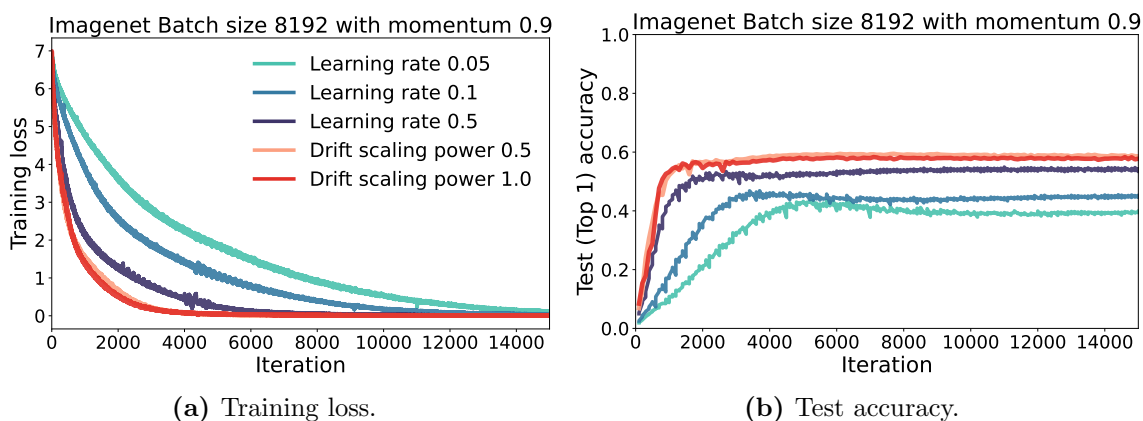


Figure 3.29: DAL with momentum: integrating drift information results in faster and more stable training compared to a fixed learning rate sweep. Compared to vanilla gradient descent there is also a significant performance and convergence speed boost.

results show that integrating drift information improves the speed of convergence compared to standard gradient descent (Figure 3.26), and leads to more stable training compared to using a fixed learning rate. We present additional experimental results in the Appendix.

Just as momentum is a common staple of optimisation algorithms, so are adaptive schemes such as Adam [47] and Adagrad [154], which adjust the step taken for each parameter independently. We can also use the knowledge from the PF to set a per-parameter learning rate: instead of using $\|\nabla_{\theta}^2 E(\theta_{t-1})\nabla_{\theta} E(\theta_{t-1})\|$ to set a global learning rate, we can use the per-parameter information provided by $\nabla_{\theta}^2 E(\theta_{t-1})\nabla_{\theta} E(\theta_{t-1})$ to adapt the learning rate of each parameter. We present preliminary results in the Figure 3.30. The above two approaches (momentum and per-parameter learning rate adaptation) can be combined, bringing us closer to the most commonly used deep learning optimisation algorithms, such as Adam (see Section 2.2). While we do not explore this avenue here, we are hopeful that this understanding of discretisation drift can be leveraged further to stabilise and improve deep learning optimisation.

Non-principal terms. This work focuses on understanding the effects of the PF on the behaviour of gradient descent. Beyond the PF, we found one non-principal term (Eq (3.4)) and have seen that it can have a stabilising effect (Figure 3.8). We

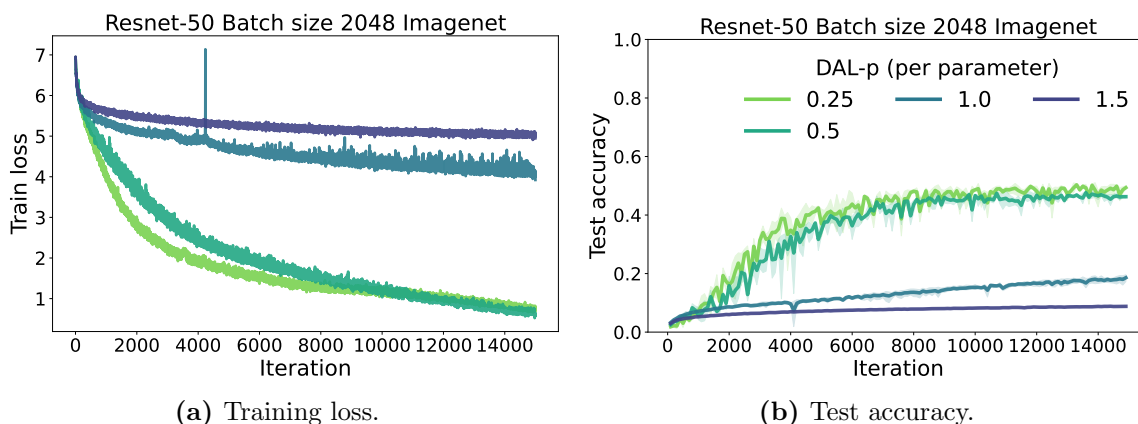


Figure 3.30: Imagenet results when using a DAL like method of scaling per-parameter. Instead of using the global learning rate $2/||\nabla_{\theta} E^2 \hat{g}(\theta)||$, using the per-parameter θ_i learning rate $h(\theta_i) = \frac{2}{(\nabla_{\theta} E^2 \hat{g}(\theta) / \sqrt{D})_i^p}$. These preliminary results show that using the per parameter information from discretisation drift can be used as an information to guide training.

provided a preliminary explanation for the stabilising effect of this non-principal term in Section 3.6. One promising avenue of non-principal terms is theoretically modelling the change of the eigenvalues λ_i in time; another promising direction is that of implicit regularisation: while existing work that uses BEA in deep learning has found important implicit regularisation effects [12, 13, 136], we have shown here that effects of $\mathcal{O}(h^3)$ are not sufficient to capture the intricacies of gradient descent, suggesting that other implicit regularisation effects could be uncovered using the non-principal terms.

Neural network theory. Many theoretical works studying at gradient descent in the neural network context use the NGF [54, 65, 100, 131]. We posit that replacing NGF in these theoretical contexts with PF may yield interesting results. In contrast to the NGF, the PF allows for the incorporation of the learning rate into the analysis, and unlike existing continuous-time models of gradient descent, it can model unstable behaviours observed in the discrete case. An example can be seen using the Neural Tangent Kernel: Jacot et al. [100] model gradient descent using the NGF to show that in the infinite-width limit gradient descent for neural networks follows kernel gradient descent. The PF can be incorporated in this analysis either by replacing the NGF with the PF as a model of gradient descent or by studying the difference in

the PF for infinitely wide and finite width networks, since discretisation drift could be responsible for the observed gap between finite and infinite networks in the large learning rate case [155].

3.9 Related work

Modified flows for deep learning optimisation. Barrett and Dherin [12] found the first order correction modified flow for gradient descent using BEA and uncovered its regularisation effects; they were the first to show the power of BEA in the deep learning context. Smith et al. [13] find the first order error correction term in expectation during one epoch of stochastic gradient descent. Modified flows have also been used for other optimisers than vanilla gradient descent: França et al. [60], Shi et al. [156] compare momentum and Nesterov accelerated momentum; Kunin et al. [65] study the symmetries of deep neural networks and use modified flows to show commonly used discrete updates break conservation laws present when using the NGF (for gradient descent they use the IGR flow while for momentum and weight decay they introduce different flows); Kovachki and Stuart [157] use modified flows to understand the behaviour of momentum by approximating Hamiltonian systems; França et al. [158] construct optimisers controlling their stability and convergence rates while Li et al. [159] construct optimisers with adaptive learning rates in the context of stochastic differential equations.

In concurrent work Miyagawa [160] use BEA to find a modified flow coined ‘Equations of Motion’ (EOM) to describe gradient descent and find higher order terms, including non-principal terms; their focus is, however, on EOM(1), which is the IGR flow, which they use to understand scale and translation invariant layers. Their approach does not expand to complex space and does not capture the instabilities studied here (see also the discussion on the difference between the full modified flow provided by BEA and the PF in Section 3.3).

Edge of stability and the importance of the Hessian. There have been a number of empirical studies on the Hessian in gradient descent. Cohen et al. [52] observed the edge of stability behaviour and performed an extensive study which led

to many empirical observations used in this work. Jastrzebski et al. [141] performed a similar study in the context of stochastic gradient descent. Sagun et al. [144], Ghorbani et al. [145], Papyan [146] approximate the entire spectrum of the Hessian, and show that there are only a few negative eigenvalues, plenty of eigenvalues centered around 0, and a few positive eigenvalues with large magnitude. Similarly, Gur-Ari et al. [124] discuss how gradient descent operates in a small subspace. Lewkowycz et al. [53] discuss the large learning rate catapult in deep learning when the largest eigenvalue exceeds $2/h$. Gilmer et al. [125] assess the effects of the largest Hessian eigenvalue in a large number of empirical settings.

There have been a series of concurrent works aimed at theoretically explaining the empirical results above. Ahn et al. [161] connect the edge of stability behaviour with what they coin as the ‘relative progress ratio’: $\frac{E(\boldsymbol{\theta}-h\nabla_{\boldsymbol{\theta}}E)-E(\boldsymbol{\theta})}{h\|\nabla_{\boldsymbol{\theta}}E\|^2}$, which they empirically show is 0 in stable areas of training and 1 in the edge of stability areas. To see the connection between the relative progress ratio and the quantities discussed in this chapter, one can perform a Taylor expansion on $\frac{E(\boldsymbol{\theta}-h\nabla_{\boldsymbol{\theta}}E)-E(\boldsymbol{\theta})}{h\|\nabla_{\boldsymbol{\theta}}E\|^2} \approx \frac{-h\nabla_{\boldsymbol{\theta}}E^T\nabla_{\boldsymbol{\theta}}E+h^2/2\nabla_{\boldsymbol{\theta}}E^T\nabla_{\boldsymbol{\theta}}^2E\nabla_{\boldsymbol{\theta}}E}{h\|\nabla_{\boldsymbol{\theta}}E\|^2} = -1 + \frac{h}{2} \frac{\nabla_{\boldsymbol{\theta}}E^T\nabla_{\boldsymbol{\theta}}^2E\nabla_{\boldsymbol{\theta}}E}{\|\nabla_{\boldsymbol{\theta}}E\|^2}$. While this ratio is related to the quantities we discuss, we also note significant differences: it is a scalar, and not a parameter-length vector and thus does not capture per eigendirection behaviour as we see with the stability coefficients (Section 3.5). Arora et al. [162] prove the edge of stability result occurs under certain conditions either on the learning rate or on the loss function. Ma et al. [163] empirically observe the multi-scale structure of the loss landscape in neural networks and use it to theoretically explain the edge of stability behaviour of gradient descent. Chen and Bruna [164] use low dimensional theoretical insights around a local minima to understand the edge of stability behaviour. Damian et al. [148] use a cubic Taylor expansion to show that gradient descent follows the trajectory of a projected method which ensures that $\lambda_0 < 2/h$ and $\nabla_{\boldsymbol{\theta}}E^T\mathbf{u} = 0$; their work is what inspired us to write the third-order non-principal term in the form of Eq (3.18), after we had previously noted its stabilising properties. These important works are complementary to our own work; they do not use continuous-time approaches and tackle primarily the

edge of stability problem or its subcases, while we focus on understanding gradient descent and applying that understanding broadly, including but not limited to the edge of stability phenomenon.

Discrete models of gradient descent. The desire to understand learning rate specific behaviour in gradient descent has been a motivation in the construction of discrete-time analyses. These analyses have provided great insights, from studying noise in the stochastic gradient descent setting [64, 165], the study of overparametrised neural models and their convergence [49, 127, 128], providing examples of when gradient descent can converge to local maxima [130], the importance of width for proving convergence in deep linear networks [129]. We differ from these studies both in motivation and execution: we are looking for a continuous-time flow that will increase the applicability of continuous-time analysis of gradient descent. We do so by incorporating discretisation drift using BEA and showing that the resulting flow is a useful model of gradient descent, which captures instabilities and escape of local minima and saddle points.

Understanding the difference between the negative gradient flow and gradient descent. Elkabetz and Cohen [54] recently examined the differences between gradient descent and the NGF in the deep learning context; their work examines the importance of the Hessian in determining when gradient descent follows the NGF. Their theoretical results show that neural networks are roughly convex and thus for reasonably sized learning rates one can expect that gradient descent follows the NGF flow closely. Their results complement ours and their approach might be extended to help us understand why the PF is sufficient to shed light on many instability observations in the neural network training.

Second-order optimisation. By using second order information (or approximations thereof) to set the learning rate, DAL is related to second-order approaches used in deep learning. Many second-order methods can be seen as approximates of Newton’s method $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \nabla_{\boldsymbol{\theta}}^2 E^{-1}(\boldsymbol{\theta}_{t-1}) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$. Since computing the Hessian inverse, or storing a matrix of the size of the Hessian as required by Quasi-Newton methods [166], can be prohibitively expensive for large models, many practical

methods approximate it with tractable alternatives [167]. Foret et al. [153] propose an optimisation scheme directly aimed at minimising sharpness and show this can improve generalisation.

Connection between drift and generalisation. We have made the connection between discretisation drift and generalisation. This connection was first made by Barrett and Dherin [12] through the IGR flow. Generalisation has also been connected to the largest eigenvalue λ_0 [18, 53, 141, 168]; recently Kaur et al. [169] however showed a more complex picture, primarily in the context of stochastic gradient descent. The largest eigenvalue could be a confounder to the drift as we have observed in Section 3.7.3; we hope that future work can deepen these connections.

3.10 Conclusion

We have expanded on previous works that used backward error analysis in deep learning to find a new continuous-time flow, called the **Principal Flow** (PF), to analyse the behaviour of gradient descent. Unlike existing flows, the PF operates in complex space, which enables it to better capture the behaviour of gradient descent compared to existing flows, including but not limited to instability and oscillatory behaviour. We use the form of the PF to find new quantities relevant to the stability of gradient descent, and shed light on newly observed empirical phenomena, such as the edge of stability results. After understanding the core quantities connected to instabilities in deep learning we devised an automatic learning rate schedule, DAL, which exhibits stable training. We concluded by cementing the connection between large discretisation drift and increased generalisation performance. We ended by highlighting future work avenues including incorporating the PF in existing theoretical analyses of gradient descent which use the negative gradient flow, incorporating our understanding of the drift of gradient descent in other optimisation approaches and specialising the PF for neural network function approximators.

Chapter 4

Continuous time models of optimisation in two-player games

The fusion of deep learning with two-player games has produced a wealth of breakthroughs in recent years, from Generative Adversarial Networks (GANs) [32] through to model-based reinforcement learning [170, 171]. Gradient descent methods are widely used across these settings, partly because these algorithms scale well to high-dimensional models and large datasets. However, much of the recent progress in our theoretical understanding of two-player differentiable games builds upon the analysis of continuous differential equations that model the dynamics of training [17, 172, 173], leading to a gap between theory and practice. Our aim is to take a step forward in our understanding of two player games by finding continuous systems which better match the gradient descent updates used in practice.

Our work builds upon Barrett and Dherin [12], who use backward error analysis (BEA) to quantify the discretisation drift induced by using gradient descent in supervised learning. We extend their work and use BEA to understand the impact of discretisation in the training dynamics of two-player games. More specifically, we quantify *Discretisation Drift* (DD), the difference between the solutions of the original flows defining the game and the discrete steps of the numerical scheme used to approximate them. To do so, we construct modified continuous systems that closely follow the discrete updates. While in supervised learning DD has a beneficial regularisation effect [12] (see also Section 2.3.2), we find that the interaction between players in DD can have a destabilising effect in adversarial games.

Contributions: Our primary contribution, Theorems 4.2.1 and 4.2.2, provides the continuous modified systems that quantify the discretisation drift in simultaneous and alternating gradient descent for general two-player differentiable games. Both theorems are novel in their scope and generality, as well as their application toward understanding the effect of the discretisation drift in two-player games parametrised with neural networks. Theorems 4.2.1 and 4.2.2 allow us to use dynamical system analysis to describe GD without ignoring discretisation drift, which we then use to:

- Provide new stability analysis tools (Section 4.3).
- Motivate explicit regularisation methods that drastically improve the performance of simultaneous gradient descent in GAN training (Section 4.6.3).
- Pinpoint optimal regularisation coefficients and shed new light on existing explicit regularisers (Table 4.2).
- Pinpoint the best performing learning rate ratios for alternating updates (Sections 4.4 and 4.6.2).
- Explain previously observed but unexplained phenomena such as the difference in performance and stability between alternating and simultaneous updates in GAN training (Section 4.5).

4.1 Background

A well-developed strategy for understanding two-player games in gradient-based learning is to analyse the continuous dynamics of the game [172, 173]. Tools from dynamical systems theory have been used to explain convergence behaviours and to improve training using modifications of learning objectives or learning rules [17, 61, 70]. While many insights from continuous analysis carry over to two-player games trained with discrete updates, discrepancies are often observed between the continuous analysis conclusions and discrete game training [68]. In this chapter, we extend BEA from the single-objective case to the more general two-player game setting, thereby bridging the gap between the continuous systems that are often analysed in theory

and the discrete numerical methods used in practice. For an overview of optimisation in two-player games, we refer the reader to Section 2.4. We provided an overview of BEA in Section 2.3.2, and used it in the previous chapter in the single-objective case.

4.1.1 Generative adversarial networks

Generative adversarial networks (GANs) [32] learn an implicit generative model through a two-player game. An implicit generative model does not have an explicit likelihood associated with it, as it is often the case with other machine learning models. Instead, an implicit model can be thought of as a simulator: it provides a sampling path that can be used to generate samples from its distribution. In the case of GANs the sampling path is defined via a latent variable model:

$$G(\mathbf{z}; \boldsymbol{\theta}), \quad \mathbf{z} \sim p(\mathbf{z}) \quad (4.1)$$

where $p(\mathbf{z})$ is a prior distribution and $G(\mathbf{z}; \boldsymbol{\theta})$ is the model with parameters $\boldsymbol{\theta}$, often a deep neural network. This sampling process induces the implicit distribution

$$p(\mathbf{x}; \boldsymbol{\theta}) = \int_{G(\mathbf{z}; \boldsymbol{\theta})=\mathbf{x}} p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})p(\mathbf{z})d\mathbf{z} = \int \mathbb{I}(G(\mathbf{z}; \boldsymbol{\theta}), \mathbf{x})p(\mathbf{z})d\mathbf{z}, \quad (4.2)$$

where $\mathbb{I}(G(\mathbf{z}; \boldsymbol{\theta}), \mathbf{x})$ is 1 if its two arguments are equal and 0 otherwise. The above integral is intractable to compute or expensive to approximate via approaches such as annealed importance sampling [174] for many cases of interest, such as when $G(\cdot, \boldsymbol{\theta})$ is given by a deep neural network. This makes standard methods, such as maximum likelihood, unavailable to learn implicit models. To learn the parameters $\boldsymbol{\theta}$, GANs introduce another model, the discriminator $D(\cdot, \boldsymbol{\phi})$, which aims to distinguish between real data—samples from the dataset—and generated data—samples from the model $G(\cdot, \boldsymbol{\theta})$. The generative model $G(\cdot, \boldsymbol{\theta})$, called *the generator*, aims to generate samples that cannot be distinguished from real data according to the discriminator.

In the original GAN paper [32] the discriminator $D(\cdot, \boldsymbol{\phi})$ is a binary classifier trained using the binary cross entropy loss; the generator’s aim is to ensure that the discriminator does not classify generated data as fake. This can be formalised as the

zero-sum game:

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} E = \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} \log D(\mathbf{x}; \boldsymbol{\phi}) + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} \log(1 - D(\mathbf{x}; \boldsymbol{\phi})) \quad (4.3)$$

$$= \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} \log D(\mathbf{x}; \boldsymbol{\phi}) + \frac{1}{2} \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})). \quad (4.4)$$

where Eq (4.4) uses the change-of-variable in the pathwise estimator discussed in Section 2.1.1. To connect this zero-sum game to generative modelling objectives used to learn probability distributions, Goodfellow et al. [32] set the functional derivative in Eq (4.4) to 0 to find the optimal discriminator D^* :

$$\frac{1}{2} p^*(\mathbf{x}) \frac{1}{D^*(\mathbf{x})} - \frac{1}{2} p(\mathbf{x}; \boldsymbol{\theta}) \frac{1}{1 - D^*(\mathbf{x})} = 0 \implies D^*(\mathbf{x}) = \frac{2p^*(\mathbf{x})}{p^*(\mathbf{x}) + p(\mathbf{x}; \boldsymbol{\theta})}. \quad (4.5)$$

Replacing the above in Eq (4.4) we have:

$$E = \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} \log \frac{2p^*(\mathbf{x})}{p^*(\mathbf{x}) + p(\mathbf{x}; \boldsymbol{\theta})} + \frac{1}{2} \mathbb{E}_{p(\mathbf{x}; \boldsymbol{\theta})} \frac{2p(\mathbf{x}; \boldsymbol{\theta})}{p^*(\mathbf{x}) + p(\mathbf{x}; \boldsymbol{\theta})} \quad (4.6)$$

$$= \mathbb{KL}(p^* \parallel \frac{1}{p^* + p(\cdot; \boldsymbol{\theta})}) + \mathbb{KL}(p(\cdot; \boldsymbol{\theta}) \parallel \frac{1}{p^* + p(\cdot; \boldsymbol{\theta})}) \quad (4.7)$$

$$= \mathbb{JSD}(p^* \parallel p(\cdot; \boldsymbol{\theta})), \quad (4.8)$$

where \mathbb{KL} and \mathbb{JSD} denote the KL and Jensen–Shannon divergences respectively. This entails that if the discriminator is optimal, the GAN generator is minimising the Jensen–Shannon divergence between the model and data distribution.

The connection between distributional divergences and two-player games has fuelled a line of research motivated by finding the best objective to train GANs. This led to new GANs motivated by Integral Probability Metrics, such as the Wasserstein distance or the Maximum Mean Discrepancy [36, 37]; f -divergences [35]; and others [175]. Practitioners have observed, however, that what is most relevant in increasing GAN performance often comes from changes architectures, optimisation, and regularisation [31, 33, 34], rather than underlying divergence used. An explanation for this observation is that the GAN discriminator is not optimal, both because it is modelled using a parametric model and because the discriminator is not

optimised to convergence for each generator update, as required to obtain the best estimate of the underlying divergence [28]. Instead, as discussed in Section 2.4, due to computational considerations the nested optimisation problem in Eq (4.4) gets translated into the alternating updates described in Algorithm 1.

Since GANs are implicit generative models, evaluating them via likelihoods is computationally intractable; approaches to approximate the likelihood do exist, but they have strong limitations [176] and are rarely used in practice. Instead, the main avenue for GAN evaluation is to measure aspects of sample quality and diversity via a pretrained classifier, as done by the Inception Score [177] and Fréchet Inception Distance [173]. The Inception Score uses label information from the dataset, together with a pre-trained classifier, to assess a generative model. It does so by using the pre-trained classifier to obtain class posterior probabilities from samples $p(y|\mathbf{x})$, and from there the marginal distribution induced over class labels by the model distribution $p(y; \boldsymbol{\theta}) = \int p(y|\mathbf{x})p(\mathbf{x}; \boldsymbol{\theta})d\mathbf{x}$, which is estimated using Monte-Carlo estimation. The Inception Score then measures average KL divergence between $p(y; \boldsymbol{\theta})$ and the distribution $p(y|\mathbf{x})$ obtained from model samples: $e^{\mathbf{E}_{p(\mathbf{x}; \boldsymbol{\theta})} \mathbb{KL}(p(y|\mathbf{x})||p_{\boldsymbol{\theta}}(y))}$. When a model produces high quality samples from all classes in the dataset $p(y; \boldsymbol{\theta})$ will be a broad distribution capturing the marginal distribution over classes in the dataset, while $p(y|\mathbf{x})$ will be a sharp distribution pre-training to the class associated with \mathbf{x} ; this leads to a high $\mathbb{KL}(p(y|\mathbf{x})||p(y; \boldsymbol{\theta}))$ and thus a high Inception Score score.

4.2 Discretisation drift in two-player games

We denote by $\boldsymbol{\phi} \in \mathbb{R}^m$ and $\boldsymbol{\theta} \in \mathbb{R}^n$ the parameters of the first and second player, respectively. The players update functions will be denoted correspondingly by $f(\boldsymbol{\phi}, \boldsymbol{\theta}) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^m$ and by $g(\boldsymbol{\phi}, \boldsymbol{\theta}) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^n$. The Jacobian $\frac{df}{d\boldsymbol{\theta}}(\boldsymbol{\phi}, \boldsymbol{\theta})$ is the $m \times n$ matrix $\frac{df}{d\boldsymbol{\theta}}(\boldsymbol{\phi}, \boldsymbol{\theta})_{i,j} = (\partial_{\theta_j} f_i)$ with $i = 1, \dots, m$ and $j = 1, \dots, n$.

We aim to understand the impact of discretising the flows

$$\dot{\boldsymbol{\phi}} = f(\boldsymbol{\phi}, \boldsymbol{\theta}) \tag{4.9}$$

$$\dot{\boldsymbol{\theta}} = g(\boldsymbol{\phi}, \boldsymbol{\theta}), \tag{4.10}$$

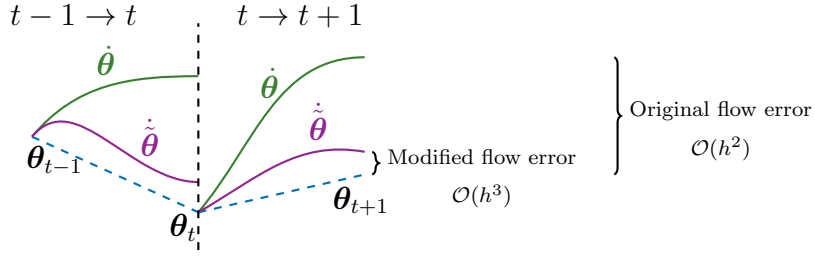


Figure 4.1: We use BEA to find the modified flow $\dot{\tilde{\theta}}$ that captures the change in parameters introduced by the discrete updates with an error of $\mathcal{O}(h^3)$. The modified flow follows the discrete update more closely than the flow $\dot{\theta}$, which has an error of $\mathcal{O}(h^2)$. Here, we show only θ for simplicity.

using either simultaneous or alternating Euler updates. To do so, we derive a modified continuous system of the form

$$\dot{\phi} = f(\phi, \theta) + hf_1(\phi, \theta) \quad (4.11)$$

$$\dot{\theta} = g(\phi, \theta) + hg_1(\phi, \theta), \quad (4.12)$$

which closely follows the discrete Euler update steps; the local error between a discrete update and the modified system is of order $\mathcal{O}(h^3)$ instead of $\mathcal{O}(h^2)$ as is the case for the continuous system given by Eqs (4.9) and (4.10). We visualise our approach in Figure 4.1.

We can specialise these modified equations using $f = -\nabla_{\phi} E_{\phi}$ and $g = -\nabla_{\theta} E_{\theta}$, where E_{ϕ} and E_{θ} are the loss functions for the two players. We will use this setting later to investigate the modified dynamics of *simultaneous* or *alternating gradient descent*. We can further specialize the form of these updates for common-payoff games ($E_{\phi} = E_{\theta} = E$) and zero-sum games ($E_{\phi} = -E$, $E_{\theta} = E$).

4.2.1 DD for simultaneous Euler updates

The *simultaneous Euler updates* with learning rates $v_{\phi}h$ and $v_{\theta}h$ respectively are given by

$$\phi_t = \phi_{t-1} + v_{\phi}hf(\theta_{t-1}, \phi_{t-1}) \quad (4.13)$$

$$\theta_t = \theta_{t-1} + v_{\theta}hg(\theta_{t-1}, \phi_{t-1}). \quad (4.14)$$

By applying BEA to these discrete updates, we obtain:

Theorem 4.2.1 *The discrete simultaneous Euler updates in (4.13) and (4.14) follow the continuous system*

$$\dot{\phi} = f - \frac{v_\phi h}{2} \left(\frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \quad (4.15)$$

$$\dot{\theta} = g - \frac{v_\theta h}{2} \left(\frac{dg}{d\phi} f + \frac{dg}{d\theta} g \right) \quad (4.16)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Remark 4.2.1 *A special case of Theorem 4.2.1 for zero-sum games with equal learning rates can be found in Lu [178].*

4.2.2 DD for alternating Euler updates

For *alternating Euler updates*, the players take turns to update their parameters, and can perform multiple updates each. We denote the number of alternating updates by m and k for the first player and second player, respectively. We scale the learning rates by the number of updates, leading to the following updates $\phi_t := \phi_{m,t}$ and $\theta_t := \theta_{k,t}$ where

$$\phi_{i,t} = \phi_{i-1,t} + \frac{v_\phi h}{m} f(\phi_{i-1,t}, \theta_{t-1}), \quad i = 1, \dots, m, \quad (4.17)$$

$$\theta_{j,t} = \theta_{j-1,t} + \frac{v_\theta h}{k} g(\phi_{m,t}, \theta_{j-1,t}), \quad j = 1, \dots, k. \quad (4.18)$$

By applying BEA to these discrete updates, we obtain:

Theorem 4.2.2 *The discrete alternating Euler updates in (4.17) and (4.18) follow the continuous system*

$$\dot{\phi} = f - \frac{v_\phi h}{2} \left(\frac{1}{m} \frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \quad (4.19)$$

$$\dot{\theta} = g - \frac{v_\theta h}{2} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} f + \frac{1}{k} \frac{dg}{d\theta} g \right) \quad (4.20)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Remark 4.2.2 *Equilibria of the original continuous systems (i.e., points where $f = \mathbf{0}$ and $g = \mathbf{0}$) remain equilibria of the modified continuous systems.*

Definition 4.2.1 The discretisation drift components we identify for each player have two terms: one term containing a player’s own update function only—terms we will call *self terms*—and a term that also contains the other player’s update function—which we will call *interaction terms*.

4.2.3 Sketch of the proofs

Following BEA, the idea is to modify the original continuous system by adding corrections in powers of the learning rate: $\tilde{f} = f + hf_1 + h^2f_2 + \dots$ and $\tilde{g} = g + hg_1 + h^2g_2 + \dots$, where for simplicity in this proof sketch, we use the same learning rate for both players (detailed derivations can be found in the Appendix Section B.1; the general structure of BEA proofs can be found in Section A.1.1). We want to find corrections f_i, g_i such that the modified system $\dot{\phi} = \tilde{f}$ and $\dot{\theta} = \tilde{g}$ follows the discrete update steps exactly. To do that we proceed in three steps:

Step 1: We expand the numerical scheme to find a relationship between ϕ_t and ϕ_{t-1} and θ_t and θ_{t-1} up to order $\mathcal{O}(h^2)$. In the case of the simultaneous Euler updates this does not require any change to Eqs (4.13) and (4.14), while for alternating updates we have to expand the intermediate steps of the integrator using Taylor series.

Step 2: We compute the Taylor series of the modified equations solution:

$$\phi(h) = \phi_{t-1} + hf + h^2 \left(f_1 + \frac{1}{2} \left(\frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \right) + \mathcal{O}(h^3) \quad (4.21)$$

$$\theta(h) = \theta_{t-1} + hg + h^2 \left(g_1 + \frac{1}{2} \left(\frac{dg}{d\phi} f + \frac{dg}{d\theta} g \right) \right) + \mathcal{O}(h^3), \quad (4.22)$$

where all the f s, g s, and their derivatives are evaluated at $(\phi_{t-1}, \theta_{t-1})$.

Step 3: We match the terms of equal power in h so that the solution of modified equations coincides with the discrete update after one step. This amounts to finding the corrections, f_i s and g_i s, so that all the terms of order higher than $\mathcal{O}(h)$ in the modified equation solution above will vanish; this yields the first order drift terms f_1

and g_1 in terms of f , g , and their derivatives. For simultaneous updates we obtain $f_1 = -\frac{1}{2}(\frac{df}{d\phi}f + \frac{df}{d\theta}g)$ and $g_1 = -\frac{1}{2}(\frac{dg}{d\phi}f + \frac{dg}{d\theta}g)$, and by construction, we have obtained the modified truncated system which follows the discrete updates exactly up to order $\mathcal{O}(h^3)$, leading to Theorem 4.2.1.

Remark 4.2.3 *The modified equations in Theorems 4.2.1 and 4.2.2 closely follow the discrete updates only for learning rates where errors of size $\mathcal{O}(h^3)$ can be neglected. Beyond this, higher order corrections are likely to contribute to the DD.*

4.2.4 Visualising trajectories

To illustrate the effect of DD in two-player games, we use a simple example adapted from Balduzzi et al. [61]:

$$\dot{\phi} = f(\phi, \theta) = -\epsilon_1\phi + \theta; \quad \dot{\theta} = g(\phi, \theta) = \epsilon_2\theta - \phi. \quad (4.23)$$

In Figure 4.2 we validate our theory empirically by visualising the trajectories of the discrete Euler steps for simultaneous and alternating updates, and show that they closely match the trajectories of the corresponding modified continuous systems that we have derived. To visualise the trajectories of the original unmodified continuous system, we use Runge–Kutta 4 (a fourth-order numerical integrator that has no DD up to $\mathcal{O}(h^5)$ in the case where the same learning rates are used for the two players—see Hairer and Lubich [179] and Appendix Section B.3). We will use Runge–Kutta 4 as a low drift method for the rest of this chapter, and compare it with Euler updates in order to understand the effect of DD.

4.3 The stability of DD

The long-term behaviour of gradient-based training can be characterised by the stability of its equilibria. Using stability analysis of a continuous system to understand discrete dynamics in two-player games has a fruitful history; however, prior work ignored discretisation drift, since they analyse the stability of the unmodified game flows in Eq (4.9) and 4.10. A general overview of stability analysis is provided in Section 2.3.1. Using the modified flows given by BEA provides two benefits here:

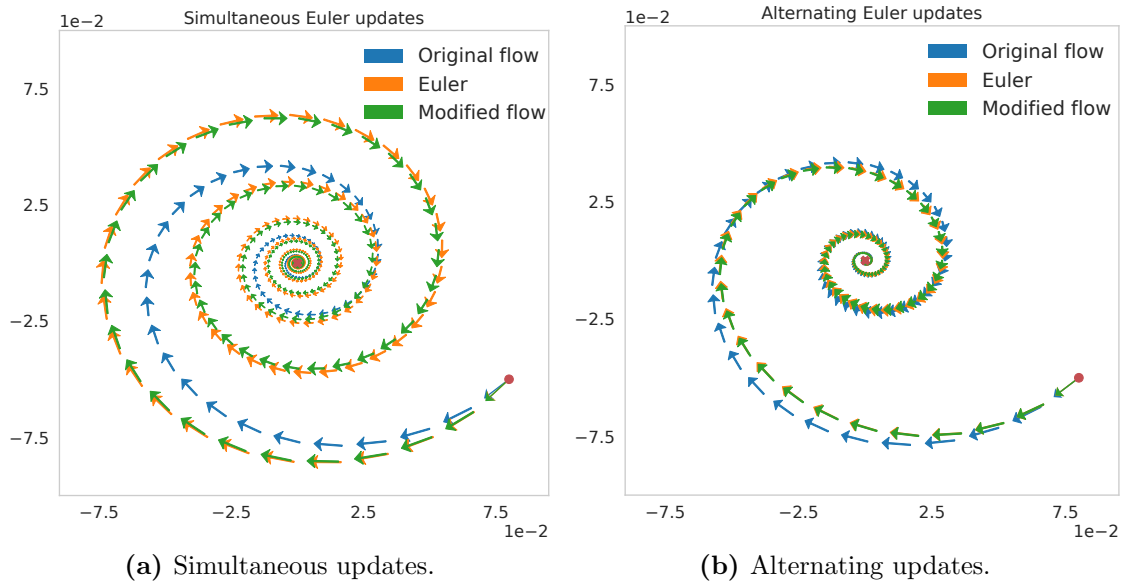


Figure 4.2: By capturing the first order DD, the modified continuous flows we derive better capture the dynamics of their corresponding discrete update than the original game continuous dynamics. By doing so, the flows also capture the difference between simultaneous (a) and alternating Euler updates (b).

first, they account for the discretisation drift, and second, they provide *different* flows for simultaneous and for alternating updates, capturing the specificities of the two optimisers.

The modified flow approach gives us a method to analyse the stability of the discrete updates: 1) Choose the system and the update type—simultaneous or alternating—to be analysed; 2) write the modified flows for the chosen system as given by Theorems 4.2.1 and 4.2.2; 3) write the corresponding modified Jacobian, and evaluate it at an equilibrium; 4) determine the stability of the equilibrium by computing the eigenvalues of the modified Jacobian, using Remark 2.3.1.

Steps 1 and 2 are easy, and step 4 is required in the stability analysis of any continuous system. For step 3, we provide a general form of the modified Jacobian at the equilibrium point of the original system, where $f = \mathbf{0}$ and $g = \mathbf{0}$

$$\tilde{\mathbf{J}} = \begin{bmatrix} \frac{d\tilde{f}}{d\phi} & \frac{d\tilde{f}}{d\theta} \\ \frac{d\tilde{g}}{d\phi} & \frac{d\tilde{g}}{d\theta} \end{bmatrix} = \mathbf{J} - \frac{h}{2}\mathbf{K}, \quad (4.24)$$

where \mathbf{J} is the Jacobian of the un-modified system in Eqs (4.9)-(4.10) and \mathbf{K} depends

on the update type; proofs are provided in Section B.4. For simultaneous updates

$$\mathbf{K}_{\text{sim}} = \begin{bmatrix} v_\phi \left(\frac{df}{d\phi}\right)^2 + v_\phi \frac{df}{d\theta} \frac{dg}{d\phi} & v_\phi \frac{df}{d\phi} \frac{df}{d\theta} + v_\phi \frac{df}{d\theta} \frac{dg}{d\theta} \\ v_\theta \frac{dg}{d\theta} \frac{dg}{d\phi} + v_\theta \frac{dg}{d\phi} \frac{df}{d\phi} & v_\theta \left(\frac{dg}{d\theta}\right)^2 + v_\theta \frac{dg}{d\phi} \frac{df}{d\theta} \end{bmatrix}. \quad (4.25)$$

For alternating updates, we have:

$$\mathbf{K}_{\text{alt}} = \begin{bmatrix} \frac{v_\phi}{m} \left(\frac{df}{d\phi}\right)^2 + v_\phi \frac{df}{d\theta} \frac{dg}{d\phi} & \frac{v_\phi}{m} \frac{df}{d\phi} \frac{df}{d\theta} + v_\phi \frac{df}{d\theta} \frac{dg}{d\theta} \\ \frac{v_\theta}{k} \frac{dg}{d\theta} \frac{dg}{d\phi} + v_\theta \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\phi} & \frac{v_\theta}{k} \left(\frac{dg}{d\theta}\right)^2 + v_\theta \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\theta} \end{bmatrix}. \quad (4.26)$$

Using the method above we show that, in two-player games, the drift term of order $\mathcal{O}(h^2)$ can change a stable equilibrium into an unstable one. This contrasts games and supervised learning, as the stability analysis of the IGR flow shows it is attracted to strict local minima (see Barrett and Dherin [12] and Section 3.3.2). For simultaneous updates with equal learning rates, this recovers a result of Daskalakis and Panageas [180] derived for zero-sum games. We show this by example: consider the game given by the system of flows in Eq (4.23). By replacing the values of f and g into the results for the Jacobian of the modified flows obtain above, we obtain the corresponding Jacobians. For simultaneous updates

$$\tilde{\mathbf{J}}_{\text{sim}} = \begin{bmatrix} -\epsilon_1 - h/2\epsilon_1^2 + h/2 & 1 + h/2\epsilon_1 - h/2\epsilon_2 \\ -1 - h/2\epsilon_1 + h/2\epsilon_2 & \epsilon_2 + h/2 - h/2\epsilon_2^2 \end{bmatrix}, \quad (4.27)$$

while for alternating updates

$$\tilde{\mathbf{J}}_{\text{alt}} = \begin{bmatrix} -\epsilon_1 - h/2\epsilon_1^2 + h/2 & 1 + h/2\epsilon_1 - h/2\epsilon_2 \\ -1 + h/2\epsilon_1 + h/2\epsilon_2 & \epsilon_2 - h/2 - h/2\epsilon_2^2 \end{bmatrix}. \quad (4.28)$$

If we consider the system with $\epsilon_1 = \epsilon_2 = 0.09$, the stability analysis of its modified flows for the simultaneous Euler updates shows they diverge when $v_\phi h = v_\theta h = 0.2$, since $\text{Tr}(\tilde{\mathbf{J}}_{\text{sim}}) = 0.19 > 0$, and thus at least one eigenvalue has positive real part. For alternating updates, we obtain $\text{Tr}(\tilde{\mathbf{J}}_{\text{alt}}) = -0.0016 < 0$ and $|\tilde{\mathbf{J}}_{\text{alt}}| = 0.981 > 0$, leading to a stable system. In both cases, the results obtained using the stability

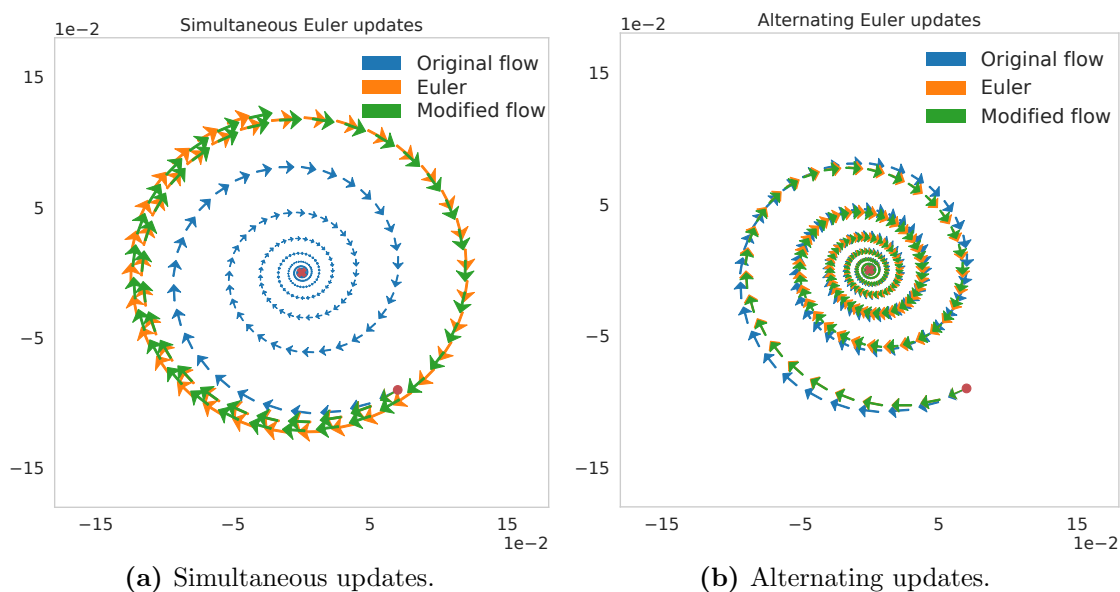


Figure 4.3: Discretisation drift can change the stability of a game. We show here an example where the flow given by the original game dynamics converges, but simultaneous Euler updates do not (a); for the same system, alternating updates converge (b). In both cases, the modified flows capture the convergent or divergent behaviour of Euler updates. $\epsilon_1 = \epsilon_2 = 0.09$ and a learning rate $v_\phi h = v_\theta h = 0.2$.

analysis of the modified flows is consistent with empirical outcomes obtained by following the corresponding discrete updates, as shown in Figure 4.3; this would not have been the case had we used the original system to do stability analysis, which would have always predicted convergence to an equilibrium and would not have been able to distinguish between simultaneous and alternating updates.

Benefits and caveats of using the modified flows for stability analysis. The modified flows help us bridge the gap between theory and practice: they allow us to extend the reach of stability analysis to a wider range of techniques used for training, such as alternating gradient descent. We hope the method we provide will be used in the context of GANs, to expand prior work, such as that of Nagarajan and Kolter [17], to alternating updates. However, the modified flows we propose here are not without limitations: they ignore discretisation errors smaller than $\mathcal{O}(h^3)$, and thus they are not equivalent to the discrete updates. We have seen in Chapter 3 how in supervised learning higher order terms can be crucial in determining instabilities of gradient descent; we will briefly investigate the effect of these higher order terms

in games in Section 4.8. To fully account for DD, methods that directly assess the convergence of discrete updates (e.g. Mescheder et al. [15]) remain an indispensable tool for understanding discrete systems.

4.4 Common-payoff games

When the players share a common loss, as in *common-payoff games*, we recover supervised learning with a single loss E , but with the extra-freedom of training the weights corresponding to different parts of the model with possibly different learning rates and update strategies (see for instance You et al. [181] where a per-layer learning rate is used to obtain extreme training speed-ups at equal levels of test accuracy). A special case occurs when two players with equal learning rates ($v_\phi = v_\theta$) perform simultaneous gradient descent. In this case, the modified losses recover the Implicit Gradient Regularisation flow in Eq (2.36). Barrett and Dherin [12] argue that minimising the loss-gradient norm, in this case, has a beneficial effect.

In this section, we instead focus on alternating gradient descent. We partition a neural network into two sets of parameters, ϕ for the parameters closer to the input and θ for the parameters closer to the output. This procedure freezes one part of the network while training the other part and alternating between the two parts. This scenario may arise in a distributed training setting, as a form of block coordinate descent. In the case of common-payoff games, we have the following as a special case of Theorem 4.2.2 by substituting $f = -\nabla_\phi E$ and $g = -\nabla_\theta E$:

Corollary 4.4.1 *In a two-player common-payoff game with common loss E , alternating gradient descent—as described in Eqs (4.17) and (4.18)—with one update per player follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = E + \frac{v_\phi h}{4} (\|\nabla_\phi E\|^2 + \|\nabla_\theta E\|^2) \quad (4.29)$$

$$\tilde{E}_\theta = E + \frac{v_\theta h}{4} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \|\nabla_\phi E\|^2 + \|\nabla_\theta E\|^2 \right) \quad (4.30)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

The term $(1 - \frac{2v_\phi}{v_\theta}) \|\nabla_\phi E\|^2$ in Eq. (4.30) is negative when the learning rates are

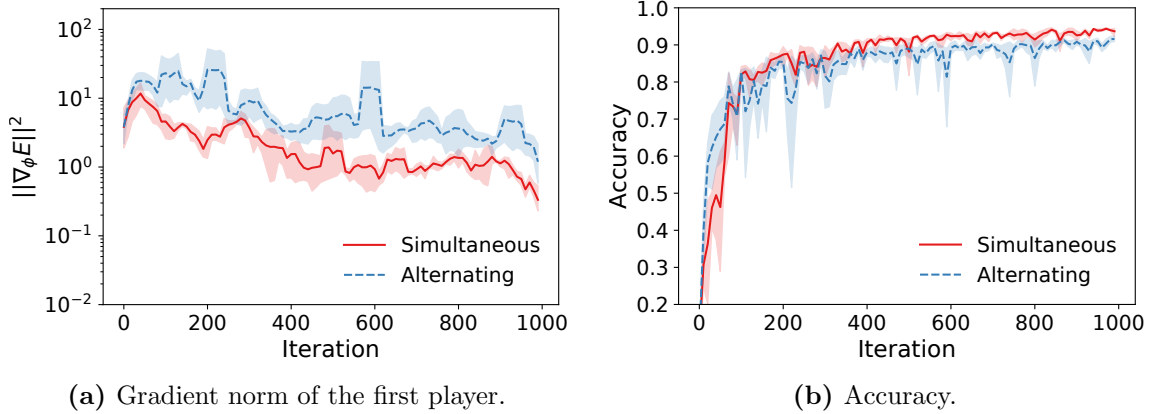


Figure 4.4: In common-payoff games alternating updates lead to higher gradient norms and unstable training when equal learning rates $v_\phi h = v_\theta h$ are used. As predicted by Corollary 4.4.1, the maximisation of the first player’s gradient norm $\|\nabla_\phi E\|$ by the second player’s modified loss in alternating updates leads to an increased gradient norm for the first player compared to simultaneous updates (a), which also results to loss instability and decrease in performance (b). The common-payoff objective here is classification of MNIST digits.

equal, seeking to maximise the gradient norm of player ϕ . According to Barrett and Dherin [12], we expect less stable training and worse performance in this case. This prediction is confirmed in Figure 4.4, where we compare simultaneous and alternating gradient descent for a MLP trained on MNIST with a common learning rate.

4.5 Analysis of zero-sum games

We now study zero-sum games, where the one player’s gain is another player’s loss: $-E_\phi = E_\theta = E$. We substitute the updates $f = \nabla_\phi E$ and $g = -\nabla_\theta E$ in the Theorems in Section 4.2 and obtain:

Corollary 4.5.1 *In a zero-sum two-player differentiable game, simultaneous gradient descent updates—as described in Eqs (4.13) and (4.14)—follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = -E + \frac{v_\phi h}{4} \|\nabla_\phi E\|^2 - \frac{v_\phi h}{4} \|\nabla_\theta E\|^2 \quad (4.31)$$

$$\tilde{E}_\theta = E - \frac{v_\theta h}{4} \|\nabla_\phi E\|^2 + \frac{v_\theta h}{4} \|\nabla_\theta E\|^2, \quad (4.32)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Corollary 4.5.1 shows that the modified losses for simultaneous gradient descent preserve the adversarial structure of the game, with the *interaction terms* maximising the gradient norm of the opposite player. We note, however, that the modified losses of zero-sum games trained with simultaneous gradient descent with different learning rates are no longer zero-sum, due to the different implicit regularisation coefficients introduced by DD.

Corollary 4.5.2 *In a zero-sum two-player differentiable game, alternating gradient descent—as described in Eqs (4.17) and (4.18)—follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = -E + \frac{v_\phi h}{4m} \|\nabla_\phi E\|^2 - \frac{v_\phi h}{4} \|\nabla_\theta E\|^2 \quad (4.33)$$

$$\tilde{E}_\theta = E - \frac{v_\theta h}{4} \left(1 - \frac{2v_\phi}{v_\theta}\right) \|\nabla_\phi E\|^2 + \frac{v_\theta h}{4k} \|\nabla_\theta E\|^2 \quad (4.34)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Corollary 4.5.2 shows that the modified losses of zero-sum games trained with alternating gradient descent are not zero-sum.

Remark 4.5.1 *Since $(1 - \frac{2v_\phi}{v_\theta}) < 1$ in the alternating case there is always less weight on the term encouraging maximizing $\|\nabla_\phi E\|^2$ compared to the simultaneous case under the same learning rates. For $\frac{v_\phi}{v_\theta} > \frac{1}{2}$ both players minimise $\|\nabla_\phi E\|^2$.*

Corollaries 4.32 and 4.34 show that in zero-sum games interaction terms can induce a pressure to maximise the gradient norm of the other player. We postulate that this maximisation effect can be a source of instability, something we will study throughout this chapter.

4.5.1 Dirac-GAN: an illustrative example

Mescheder et al. [68] introduce the Dirac-GAN as an example to illustrate the often complex training dynamics in zero-sum games. We follow this example to provide an intuitive understanding of DD. Dirac-GAN aims to learn a Dirac delta distribution with mass at zero; the generator is modelling a Dirac with parameter θ : $G(z; \theta) = \theta$

and the discriminator is a linear model on the input with parameter ϕ : $D(x; \phi) = \phi x$. This results in the zero-sum game given by

$$E = l(\theta\phi) + l(0), \quad (4.35)$$

where l depends on the GAN formulation used ($l(z) = -\log(1 + e^{-z})$ for instance). As in Mescheder et al. [68], we assume l is continuously differentiable with $l'(x) \neq 0$ for all $x \in \mathbb{R}$. The partial derivatives of the loss function

$$\frac{\partial E}{\partial \phi} = l'(\theta\phi)\theta, \quad \frac{\partial E}{\partial \theta} = l'(\theta\phi)\phi, \quad (4.36)$$

lead to the underlying continuous dynamics

$$\dot{\phi} = f(\theta, \phi) = l'(\theta\phi)\theta, \quad \dot{\theta} = g(\theta, \phi) = -l'(\theta\phi)\phi. \quad (4.37)$$

The unique equilibrium point is $\theta = \phi = 0$.

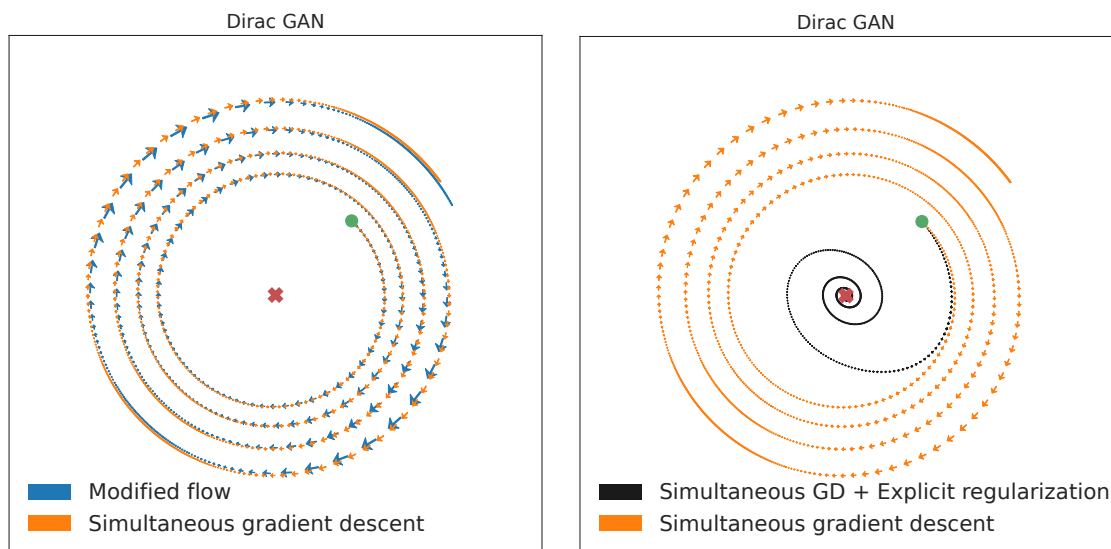
Reconciling discrete and continuous updates in Dirac-GAN. The original continuous dynamics in Eq (4.37) preserve $\theta^2 + \phi^2$, since

$$\frac{d(\theta^2 + \phi^2)}{dt} = 2\theta \frac{d\theta}{dt} + 2\phi \frac{d\phi}{dt} = -2\theta l'(\theta\phi)\phi + 2\phi l'(\theta\phi)\theta = 0. \quad (4.38)$$

Mescheder et al. [68] show, however, that for simultaneous gradient descent $\theta^2 + \phi^2$ increases with each update: different conclusions are reached when analysing the dynamics of the original continuous system versus the discrete updates. We show that, by using the modified flows given by Eqs (4.31) and (4.32) instead of the original game dynamics, we can resolve this discrepancy. When following the modified flow induced by simultaneous gradient descent in DiracGAN, we have that (proof in Section B.6.1 in the Appendix), unless we start at the equilibrium $\theta = \phi = 0$,

$$\frac{d(\theta^2 + \phi^2)}{dt} = h\theta^2 l'(\phi, \theta)^2 + h\phi^2 l'(\phi, \theta^2) > 0. \quad (4.39)$$

We thus obtain *consistent conclusions from the modified continuous flows and the*



(a) Reconciling discrete and continuous dynamics. (b) Explicit regularisation leads to convergence.

Figure 4.5: DiracGAN. (a): The dynamics of simultaneous gradient descent updates match the continuous dynamics given by BEA for DiracGAN. (b): Explicit regularisation cancelling the interaction terms of DD stabilises the DiracGAN trained with simultaneous gradient descent.

discrete dynamics used in practice; we visualise this in Figure 4.5a.

Explicit regularisation stabilises Dirac-GAN. We show that the instability of simultaneous updates in DiracGAN can be seen as a result of the effect of interaction terms in zero-sum games, as we postulated as a result of Corollary 4.32. We counteract the norm maximisation effect of the interaction terms in simultaneous updates by cancelling the interaction terms using explicit regularisation: $E_\phi = -E + \gamma \|\nabla_\theta E\|^2$ and $E_\theta = E + \zeta \|\nabla_\phi E\|^2$ where γ, ζ are of $\mathcal{O}(h)$. We find (proofs are provided in Section B.6.3 of the Appendix) that the modified Jacobian of simultaneous updates induced by these new losses is negative definite if $\gamma > hv_\phi/4$ and $\zeta > hv_\theta/4$ —which are exactly the coefficients required to cancel the interaction terms shown in Corollary 4.32—so the system is asymptotically stable and converges to the optimum. Unlike the modified flow of the original zero-sum dynamics, which diverge as we have seen in Eq 4.39, the flow of the system induced by cancelling the interaction terms converges. We visualise this behaviour in Figure 4.5b. Notably, by quantifying DD, we are able to find the regularisation coefficients that guarantee convergence and show that they depend on learning rates.

4.6 Experimental analysis of GANs

To understand the effect of DD on more complex adversarial games, we analyse GANs trained for image generation on the CIFAR-10 dataset. We follow the model architectures from Spectral-Normalised GAN [30]. Both players have millions of parameters. Unlike Spectral-Normalised GAN, we employ the original GAN formulation, as described in Eq (4.4).

When it comes to gradient descent, GAN practitioners often use alternating, not simultaneous updates: the discriminator is updated first, followed by the generator; we described alternating updates and simultaneous updates in Algorithms 1 and 2, respectively. Recent work, however, shows higher-order numerical integrators can work well with simultaneous updates [67]. We will show that DD can be seen as the culprit behind some of the challenges in simultaneous gradient descent in zero-sum GANs, indicating ways to improve training performance in this setting. For a clear presentation of the effects of DD, we employ a minimalist training setup. Instead of using popular adaptive optimisers, such as Adam [47], we train all models with vanilla stochastic gradient descent, without momentum or variance reduction methods.

We use the Inception Score (IS) [177] for evaluation, and we report Fréchet Inception Distance (FID) results [173] in the Appendix; we described the computation of the Inception Score in Section 4.1. Our training curves contain a horizontal line at the Inception Score of 7.5, obtained with the same architectures we use, but with the Adam optimiser (the score reported by Miyato et al. [30] is 7.42). We report box plots showing the performance quantiles across all hyperparameter and seeds, together with *learning curves corresponding to the best 10% models from a sweep over learning rates and 5 seeds*. This is due to the large variability observed across hyperparameters and seeds, but we observe the same relative results with the top 20% and 30% of models, as we show in the Appendix Section B.9, which also has additional discussions on robustness. For SGD we use learning rates $\{5 \times 10^{-2}, 1 \times 10^{-2}, 5 \times 10^{-3}, 1 \times 10^{-3}\}$ for each player; for Adam, we use learning rates $\{1 \times 10^{-4}, 2 \times 10^{-4}, 3 \times 10^{-4}, 4 \times 10^{-4}\}$, which have been widely used in the literature [30]. When comparing to Runge–Kutta (RK4) to assess the effect of DD, we always use the same learning rates

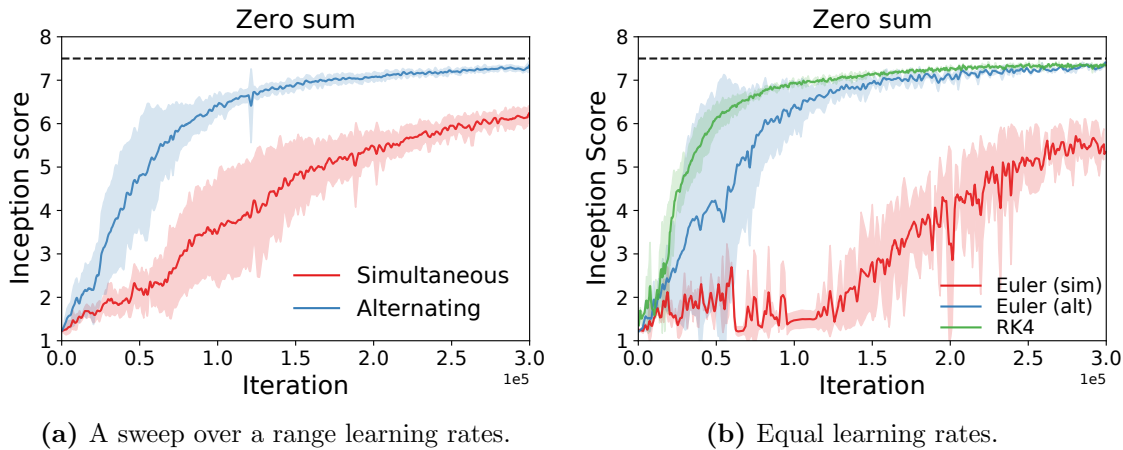


Figure 4.6: The effect of DD on zero-sum games. Over a sweep of learning rates, alternating updates perform better than simultaneous updates (a), and show lower variance. When using equal learning rates RK4 has $\mathcal{O}(h^5)$ drift, and thus we use it as a baseline to understand the effect of the drift in (b); we observe that reducing drift leads to improved performance. These results show that DD has a strong effect on these games, with a substantial difference in performance between simultaneous and alternating updates, and an increased speed early in training when using RK4. While the modified flows we study only capture DD up to $\mathcal{O}(h^3)$, these results are consistent with what we expect from Corollaries 4.5.1 and 4.5.2, which show that the drift of the first player leads to an incentive to maximise the second player’s gradient norm; this can lead to instability compared to methods with reduced drift such as RK4 as we have seen in (b). Similarity when comparing alternating and simultaneous updates, we expect less instability for alternating updates as the second player does not always have the incentive to maximise the first player’s gradient norm, which is always the case for simultaneous updates; this is consistent with what we observe in (a).

for both players. We present results using additional losses, via LS-GAN [38], and experimental details in the Appendix. The code associated with this work can be found at https://github.com/deepmind/discretisation_drift.

4.6.1 Does DD affect training?

We start our experimental analysis by showing the effect of DD on zero-sum games. We compare simultaneous gradient descent, alternating gradient descent, and Runge–Kutta 4 updates, since they follow different continuous dynamics given by the modified equations we have derived up to $\mathcal{O}(h^3)$ error. Figure 4.6 shows simultaneous gradient descent performs substantially worse than alternating gradient descent. When compared to Runge–Kutta 4, which has a DD error of $\mathcal{O}(h^5)$ when the two

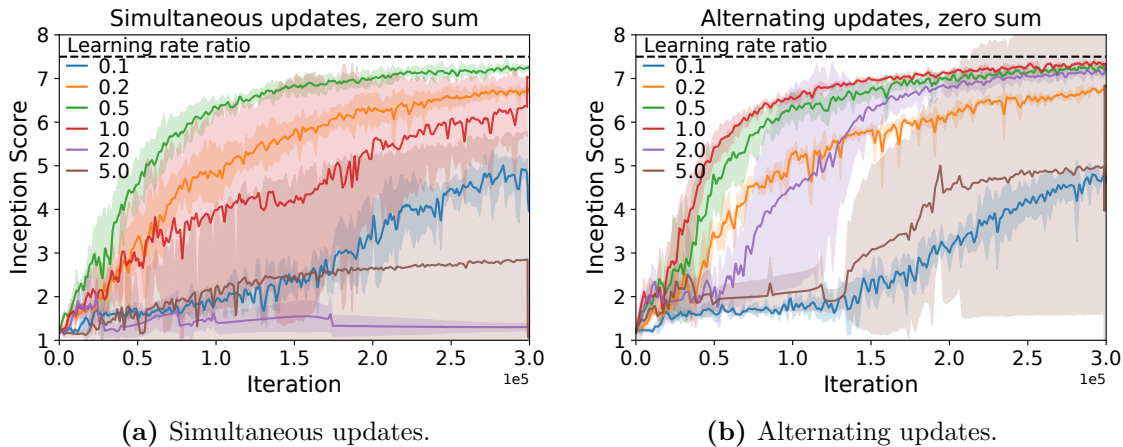


Figure 4.7: Alternating gradient descent (b) performs better for learning rate ratios which reduce the adversarial nature of DD, namely learning rates where $v_\phi/v_\theta \geq 0.5$ and thus the second player’s (the generator) modified loss minimises, instead of maximises the discriminator’s gradient norm. The same learning rate ratios show no advantage in the simultaneous case (a).

players have equal learning rates, we see that Runge–Kutta performs better, and that removing the drift improves training. Multiple updates also affect training, either positively or negatively, depending on learning rates—see Appendix Section B.3.

4.6.2 The importance of learning rates in DD

While in simultaneous updates (Eqs (4.31) and (4.32)) the interaction terms of both players maximise the gradient norm of the other player, alternating gradient descent (Eqs (4.33) and (4.34)) exhibits less pressure on the second player (generator) to maximise the norm of the first player (discriminator). In alternating updates, when the ratio between the discriminator and generator learning rates exceeds 0.5, both players are encouraged to minimise the discriminator’s gradient norm. To understand the effect of learning rate ratios in training, we performed a sweep where the discriminator learning rates v_ϕ are sampled uniformly between $[0.001, 0.01]$, and the learning rate ratios v_ϕ/v_θ are in $\{0.1, 0.2, 0.5, 1, 2, 5\}$, with results shown in Figure 4.7. Learning rate ratios greater than 0.5 perform best for alternating updates and enjoy a substantial increase in performance compared to simultaneous updates, which is consistent with our expectations based on the sign of the interaction terms.

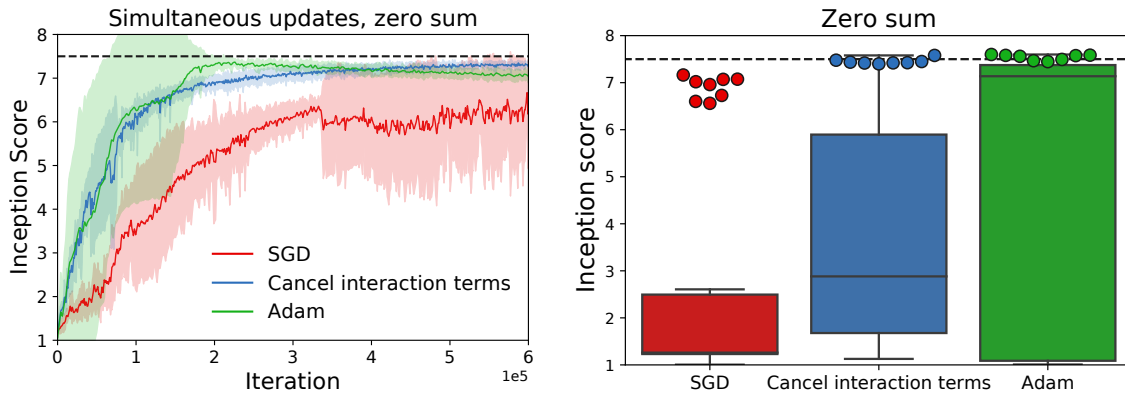


Figure 4.8: Simultaneous updates in zero-sum games: Explicit regularisation cancelling the interaction terms of DD improves performance and stability, both for the best performing models (a) and across a learning rate sweep (b).

4.6.3 Improving performance by explicit regularisation

We have postulated that the interaction terms can hinder performance and stability in zero-sum games, as they can lead to an incentive to maximise the gradient norm of one, or both players. Now, we investigate whether *cancelling the interaction terms* between the two players can improve training stability and performance in zero-sum games trained using simultaneous gradient descent. We train models using the losses

$$E_\phi = -E + c_1 \|\nabla_\theta E\|^2 \quad (4.40)$$

$$E_\theta = E + c_2 \|\nabla_\phi E\|^2. \quad (4.41)$$

If c_1, c_2 are $\mathcal{O}(h)$ we can ignore the DD from these regularisation terms, since their effect on DD will be of order $\mathcal{O}(h^3)$. We can set coefficients to be the negative of the coefficients present in DD, namely $c_1 = v_\phi h/4$ and $c_2 = v_\theta h/4$ to cancel the interaction terms, which maximise the gradient norm of the other player, while keeping the self terms, which minimise the player’s own gradient norm. We show results in Figure 4.8: cancelling the interaction terms leads to substantial improvement compared to SGD, obtains the same peak performance as Adam (though requires more training iterations) and recovers the performance of Runge–Kutta 4 (Figure 4.6). Unlike Adam, we do not observe a decrease in performance later in training but report higher variance in performance across seeds—see Figure 4.8b and additional

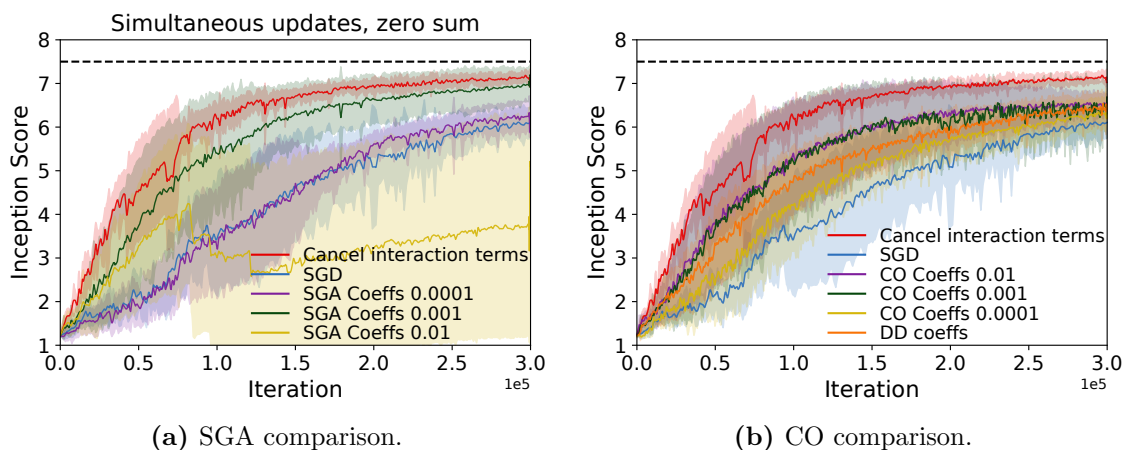


Figure 4.9: Comparison with Symplectic Gradient Adjustment (SGA) and Consensus Optimisation (CO) in simultaneous gradient descent for zero-sum games: DD motivates the form of explicit regularisation and provides does not require a larger hyperparameter sweep for the regularisation coefficient, as the explicit regularisation coefficients are determined by the player’s learning rates. Cancelling the interaction terms leads to improved performance compared to both SGA and CO.

experiments in Appendix Section B.9.

Connection with Symplectic Gradient Adjustment (SGA): Balduzzi et al. [61] proposed SGA to improve the dynamics of gradient-based method for games, by counter-acting the rotational force of the vector field (see Section B.5 for a discussion). Adjusting the gradient field can be viewed as modifying the losses as in Eq (4.41); the modification from SGA cancels the interaction terms we identified. However, it is unclear whether the fixed coefficients of $c_1 = c_2 = \frac{1}{2}$ used by SGA are always optimal, while our analysis shows the strength of DD changes with the exact discretisation scheme, such as the step size. Indeed, as our experimental results in Figure 4.9a show, adjusting the coefficient in SGA strongly affects training, and that cancelling the interaction terms outperforms SGA.

Connection with Consensus Optimisation (CO): Mescheder et al. [15] analyse the discrete dynamics of gradient descent in zero-sum games and prove that, under certain assumptions, adding explicit regularisation that encourages the players to minimise the gradient norm of both players guarantees convergence to a local Nash equilibrium. Their approach includes cancelling the interaction terms, but also

requires *strengthening the self terms*, using losses:

$$E_\phi = -E + c_1 \|\nabla_{\theta} E\|^2 + s_1 \|\nabla_{\phi} E\|^2 \quad (4.42)$$

$$E_{\theta} = E + s_2 \|\nabla_{\theta} E\|^2 + c_2 \|\nabla_{\phi} E\|^2, \quad (4.43)$$

where they use $s_1 = s_2 = c_1 = c_2 = \gamma$ where γ is a hyperparameter. In order to understand the effect of the self and interaction terms, we compare to CO, as well as a similar approach where we use coefficients proportional to the drift, namely $s_1 = v_{\phi}h/4$ and $s_2 = v_{\theta}h/4$; this effectively doubles the strength of the self terms in DD. We show results in Figure 4.9b. We first notice that CO can improve results over vanilla SGD. However, similarly to what we observed with SGA, the regularisation coefficient is important and thus requires a hyperparameter sweep, unlike our approach, which uses the coefficients provided by the DD. We further notice that strengthening the norms using the DD coefficients can improve training, but performs worse compared to only cancelling the interaction terms. This shows the importance of finding the right training regime, and that strengthening the self terms does not always improve performance.

Alternating updates: We perform the same exercise for alternating updates, where $c_1 = v_{\phi}h/4$ and $c_2 = v_{\theta}h/4(1 - \frac{2v_{\phi}}{v_{\theta}})$. We also study the performance obtained by only cancelling the discriminator interaction term, since when $v_{\phi}/v_{\theta} > 0.5$ the generator interaction term minimises, rather than maximises, the discriminator gradient norm and thus the generator interaction term might not have a strong destabilising force. We observe that adding explicit regularisation only brings the benefit of reduced variance when cancelling the discriminator interaction term (Figure 4.10). As for simultaneous updates, we find that knowing the form of DD guides us to a choice of explicit regularisation: for alternating updates cancelling both interaction terms can hurt training, but the form of the modified losses suggests that we should only cancel the discriminator interaction term, with which we can obtain some gains.

Does cancelling the interaction terms help for every choice of learning rates? The substantial performance and stability improvement we observe applies to the performance obtained across a learning rate sweep. For individual learning rate

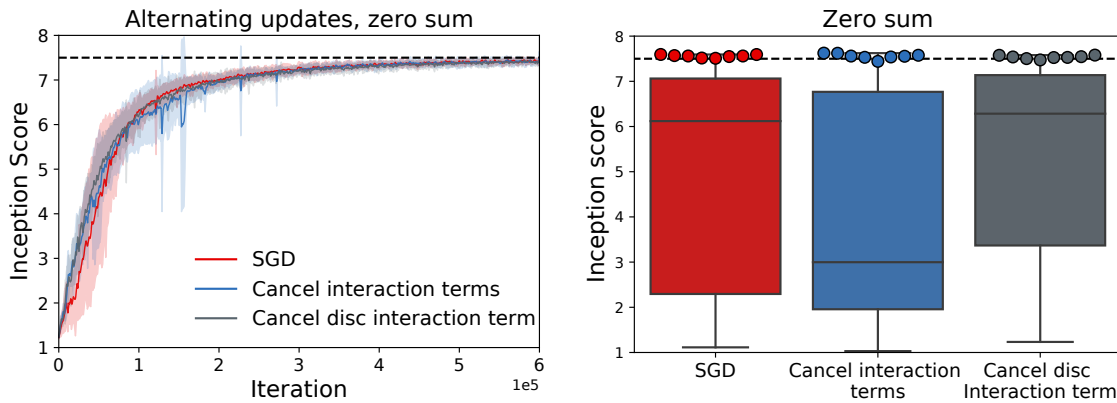


Figure 4.10: Alternating updates in zero-sum games: the form of DD guides us in finding explicit regularisation terms. Since for alternating updates the generator’s interaction term can be beneficial by minimising the discriminator’s gradient norm, cancelling only the discriminator interaction term improves sensitivity across hyperparameters (b).

choices, cancelling the interaction terms is not guaranteed to improve learning.

4.6.4 Extension to non-zero-sum GANs

Finally, we extend our analysis to GANs with the non-saturating loss for the generator

$$E_{\theta} = -\log D(G(\mathbf{z}; \theta); \phi) \quad (4.44)$$

introduced by Goodfellow et al. [32], while keeping the discriminator loss unchanged as that from the zero-sum formulation in Eq (4.4). The non-saturating loss has been observed to perform better than its zero-sum counterpart, and has often been used in practice. In contrast with the dynamics from zero-sum GANs we analysed earlier, changing from simultaneous to alternating updates results in little change in performance, as can be seen in Figure 4.11. Despite having the same adversarial structure and the same discriminator loss, changing the generator loss changes the relative performance of the different discrete update schemes.

Since the effect of DD strongly depends on the game, we recommend analysing the performance of discrete numerical schemes on a case by case basis. Indeed, while for general two-player games we cannot always write modified losses as for the zero-sum case—see the Section 5.3 for a discussion—we can use Theorems 4.2.1

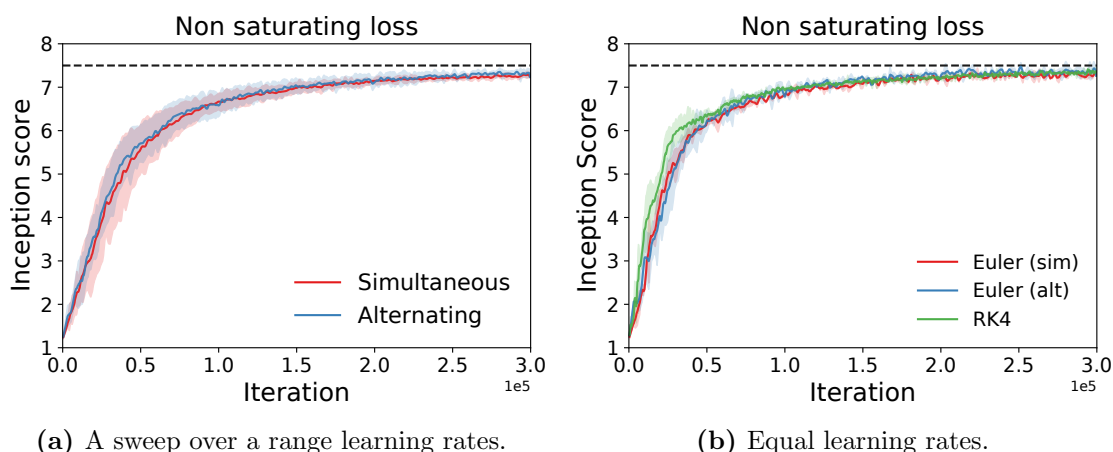


Figure 4.11: The effect of DD depends on the game: its effect is less strong for the non-saturating loss. Since when using the non-saturating loss GAN losses are no longer zero-sum, estimating the impact of the drift becomes more challenging, since we cannot easily write modified losses (though we provide novel insights into this problem in the next chapter, Chapter 5). When using equal learning rates RK4 has $\mathcal{O}(h^5)$ drift, and thus we use it as a baseline to understand the effect of the drift in (b).

and 4.2.2 to understand the effect of the drift for specific choices of loss functions. We will tackle the question of modified losses for generally differentiable two-player games in the Chapter 5, and further contrast the non-saturating and the saturating generator losses in GANs with Spectral Normalisation in Chapter 7.

4.7 A comment on different learning rates

When considering different learning rates for the two players, we have thus far kept the consistency between physical time and learning rates. That is, we assumed the gradient descent updates are obtained by Euler discretisation with learning rates $v_\phi h$ and $v_\theta h$ of the flow

$$\dot{\phi} = f(\phi, \theta) \quad (4.45)$$

$$\dot{\theta} = g(\phi, \theta). \quad (4.46)$$

We then used BEA to fine modified flow such that after time $v_\phi h$ and $v_\theta h$ the error between the modified flow and the discrete updates is $\mathcal{O}(h^3)$.

Alternatively, we can consider the gradient descent updates as the result of

Euler discretisation with learning rate h of the flow

$$\dot{\boldsymbol{\phi}} = v_{\phi} f(\boldsymbol{\phi}, \boldsymbol{\theta}) \quad (4.47)$$

$$\dot{\boldsymbol{\theta}} = v_{\theta} g(\boldsymbol{\phi}, \boldsymbol{\theta}). \quad (4.48)$$

While in this approach the learning rate loses the connection with physical time for each player, it keeps the consistency of time between two players: after one iteration the same amount of time, h , has passed for both players. With this in mind, we will call this ‘the same-physical time between players’ approach.

Considering gradient descent as a discretisation of Eqs (4.47) and (4.48) instead of Eqs (4.45) and (4.46) has implications for the analysis of the behaviour of gradient descent using continuous-time tools such as stability analysis, as the two systems can lead to different results based on the two learning rates. Using Eqs (4.47) and (4.48) can also lead to a different perspective from the loss minimisation view. For the game with $f = -\nabla_{\phi} E$ and $g = \nabla_{\theta} E$, the second system leads to the interpretation that the first player is minimising function $v_{\phi} E(\boldsymbol{\phi}, \boldsymbol{\theta})$ while the second player is maximising function $v_{\theta} E(\boldsymbol{\phi}, \boldsymbol{\theta})$. This is different than the zero-sum formulation of the first player minimising $E(\boldsymbol{\phi}, \boldsymbol{\theta})$ while the second player is maximising it, resulting from our previous interpretation. The same-physical time between players approach also presents the following challenge:

Remark 4.7.1 (The identifiability problem.) *Given learning rates $v_{\phi}h$ and $v_{\theta}h$ for the two-players, h cannot be identified exactly and has to be chosen by the user. The choice of h leads to different flows describing the dynamics of the discrete updates. This is in contrast with the approach taken previously, which only depends on v_{ϕ} and v_{θ} through the learning rates $v_{\phi}h$ and $v_{\theta}h$.*

Following the same-physical time between players approach also has implications for the use of BEA to find flows with error $\mathcal{O}(h^3)$ after one Euler update. Luckily, to find the modified flows using BEA under the same-physical time approach, we simply have to apply our results from Section 4.2 for the system in Eqs (4.47) and (4.48) by using their corresponding vector fields and learning rate h . We do so, and obtain:

Theorem 4.7.1 *The discrete simultaneous Euler updates in Eqs (4.13) and (4.14) follow the continuous system*

$$\dot{\phi} = v_\phi f - \frac{h}{2} \left(v_\phi^2 \frac{df}{d\phi} f + v_\phi v_\theta \frac{df}{d\theta} g \right) \quad (4.49)$$

$$\dot{\theta} = v_\theta g - \frac{h}{2} \left(v_\phi v_\theta \frac{dg}{d\phi} f + v_\theta^2 \frac{dg}{d\theta} g \right) \quad (4.50)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Theorem 4.7.2 *The discrete alternating Euler updates in (4.17) and (4.18) follow the continuous system*

$$\dot{\phi} = v_\phi f - \frac{h}{2} \left(\frac{v_\phi^2}{m} \frac{df}{d\phi} f + v_\phi v_\theta \frac{df}{d\theta} g \right) \quad (4.51)$$

$$\dot{\theta} = v_\theta g - \frac{h}{2} \left(-v_\theta v_\phi \frac{dg}{d\phi} f + \frac{v_\theta^2}{k} \frac{dg}{d\theta} g \right) \quad (4.52)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

We observe that both modified flows each player's vector field depends on the learning rate of the other player, unlike in our previous results where for simultaneous updates, the vector field for each player depended only on its learning rate. When interpreting these results to obtain implicit regularisation effects as we have done previously in zero-sum and common-payoff games, the resulting coefficients of the implicit regularisers are different than those obtained in Section 4.5. For zero-sum games, we obtain the following corollaries:

Corollary 4.7.1 *In a zero-sum two-player differentiable game, simultaneous gradient descent updates—as described in Eqs (4.13) and (4.14)—follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = -v_\phi E + \frac{v_\phi^2 h}{4} \|\nabla_\phi E\|^2 - \frac{v_\phi v_\theta h}{4} \|\nabla_\theta E\|^2 \quad (4.53)$$

$$\tilde{E}_\theta = v_\theta E - \frac{v_\phi v_\theta h}{4} \|\nabla_\phi E\|^2 + \frac{v_\theta^2 h}{4} \|\nabla_\theta E\|^2 \quad (4.54)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Corollary 4.7.2 *In a zero-sum two-player differentiable game, alternating gradient descent—as described in Eqs (4.17) and (4.18)—follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = -v_\phi E + \frac{v_\phi^2 h}{4m} \|\nabla_\phi E\|^2 - \frac{v_\phi v_\theta h}{4} \|\nabla_\theta E\|^2 \quad (4.55)$$

$$\tilde{E}_\theta = v_\theta E + \frac{v_\phi v_\theta h}{4} \|\nabla_\phi E\|^2 + \frac{v_\theta^2 h}{4k} \|\nabla_\theta E\|^2 \quad (4.56)$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

From Corollary 4.7.2 one concludes for alternating updates that the second player’s interaction *always minimises the gradient norm of the first player*, which again helps explain the empirically observed stability of alternating updates over simultaneous updates; this is in contrast with our previous interpretation which suggested that this occurs only for certain learning rate ratios (Corollary 4.5.2).

If we would like to implement explicit regularisation methods to cancel the implicit regularisers obtained above using the same-physical time between players approach, we have to discretise the modified losses with learning rate h ; this is opposed to using learning rates $v_\phi h$ and $v_\theta h$ to discretise the modified flows obtained in Section 4.5. For the *self terms*, no difference is obtained compared to the previous interpretation; under both approaches of tackling different learning rates, the strength of the self terms in the parameters update is proportional to $v_\phi^2 h^2$ and $v_\theta^2 h^2$, respectively. For interaction terms, however, the coefficients required in the parameter update are changed. We highlighted these changes in Table 4.1; these coefficients also show that while in order to find the modified flow in the same-physical time between player approach we require to choose v_ϕ , v_θ and h , for explicit regularisation *we only need the learning rates of the two players*.

We thus can investigate the effect of explicit regularisation to cancel the interaction terms under the same-physical time approach in zero-sum games. As before, we do so by using the generator saturating loss in GAN training. We show results in Figures 4.12 and 4.13 and compare with our previous results from Section 4.6.3. While for the simultaneous updates using the same physical time approach leads

Simultaneous updates	First player (ϕ)	Second player (θ)
Corollary 4.5.1	$\frac{v_\phi^2 h^2}{4} \nabla_\phi \ \nabla_\theta E\ ^2$	$\frac{v_\theta^2 h^2}{4} \nabla_\theta \ \nabla_\phi E\ ^2$
Corollary 4.7.1	$\frac{v_\phi v_\theta h^2}{4} \nabla_\phi \ \nabla_\theta E\ ^2$	$\frac{v_\phi v_\theta h^2}{4} \nabla_\theta \ \nabla_\phi E\ ^2$
Alternating updates	First player (ϕ)	Second player (θ)
Corollary 4.5.2	$\frac{v_\phi^2 h^2}{4} \nabla_\phi \ \nabla_\theta E\ ^2$	$\left(\frac{v_\theta^2 h}{4} - \frac{2v_\phi v_\theta h^2}{4}\right) \nabla_\theta \ \nabla_\phi E\ ^2$
Corollary 4.7.2	$\frac{v_\phi v_\theta h^2}{4} \nabla_\phi \ \nabla_\theta E\ ^2$	$-\frac{v_\phi v_\theta h^2}{4} \nabla_\theta \ \nabla_\phi E\ ^2$

Table 4.1: The strength of the interaction terms in DD under the different modified continuous-time flows we find for simultaneous and alternating gradient descent in zero-sum games. We obtained the different flows via different approaches of handling different learning rates for the two players, namely $v_\phi h \neq v_\theta h$. We do not write the self terms here as they are the same for both approaches. We write the implicit regularisation coefficients as observed in the continuous-time flow displacement for one gradient descent update, rather than the modified loss, in order to account for the different learning rates under the two interpretations (which are $v_\phi h$ and h for the first player, respectively). For simultaneous updates, the insight that the interaction terms for each player leads to a maximisation of the other player’s gradient norm remains true regardless of interpretation, though the strength of the regularisation changes in the two interpretations. For alternating updates, the different interpretations suggest not only different regularisation strengths but also signs, with the result from Corollary 4.7.2 showing that for alternating updates, the interaction term of the second player always minimises the gradient norm of the first player, as opposed to the result in Corollary 4.5.2, for which this effect depends on the learning rate ratio v_ϕ/v_θ . Both results are consistent in explaining why alternating updates are more stable than simultaneous updates, since in both cases for alternating updates the second player’s interaction term has a smaller coefficient when maximising the first player’s gradient norm compared to that of simultaneous updates, or is minimising the first player’s gradient norm.

to more hyperparameters with increased performance, no significant difference is obtained for alternating updates.

We conclude by noting that when studying two-player games trained with different learning rates, one can choose between two approaches of constructing corresponding continuous dynamics and modified losses. Both are valid, in that both lead to modified flows with the same error in learning rate, and can be used to understand implicit regularisation effects and the discrepancies between simultaneous and alternating Euler updates. Additional examination is needed to understand how these approaches explain continuous-time behaviour over a larger number of steps; we leave that as future work.

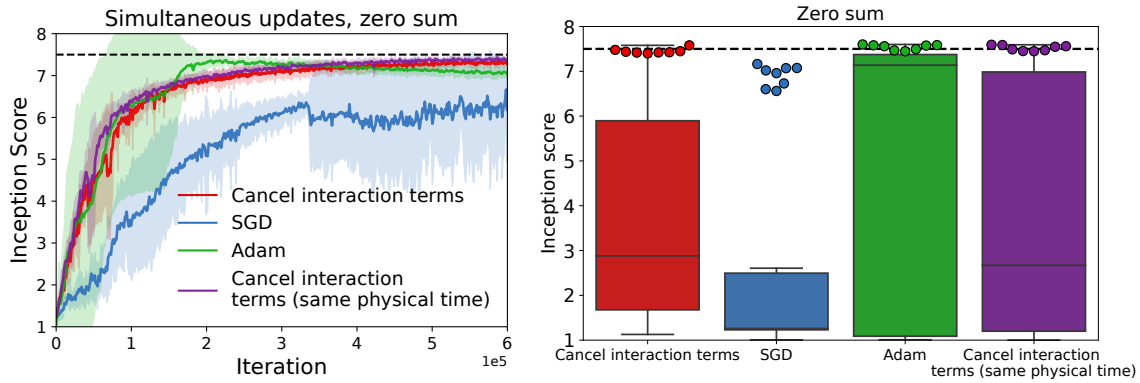


Figure 4.12: Simultaneous updates in zero-sum games: cancelling the interaction terms under the different interpretations of tackling different learning rates (Corollaries 4.5.1 and 4.7.1) does not result in a significant empirical difference. We denote the approach of Corollary 4.7.1 as ‘same physical time’, as per its motivation of using the same physical time between the two players.

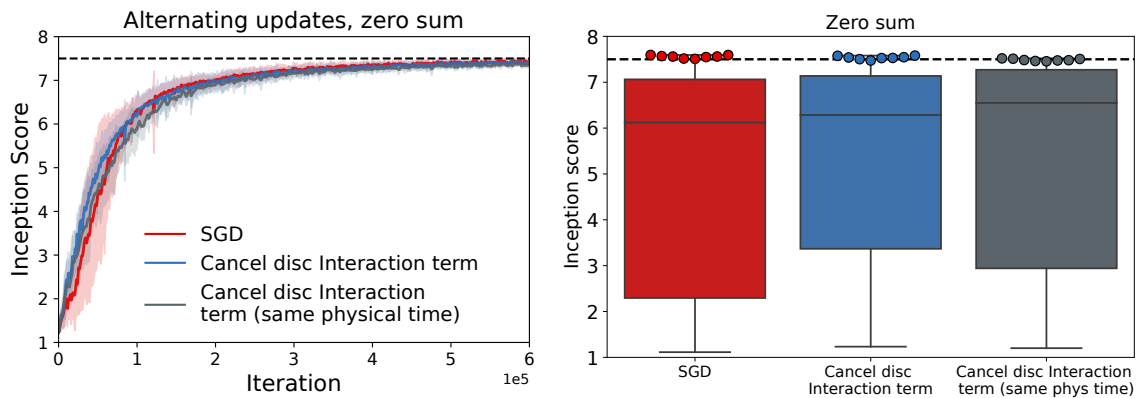


Figure 4.13: Alternating updates in zero-sum games: cancelling the interaction terms under the different interpretations of tackling different learning rates (Corollaries 4.5.2 and 4.7.2) does not result in a significant empirical difference. We denote the approach of Corollary 4.7.2 as ‘same physical time’, as per its motivation of using the same physical time between the two players.

4.8 Extending the PF for two-player games

In Chapter 3, we introduced the principal flow (PF) as a model of single-objective gradient descent. We now briefly show how the same approach can be generalised to simultaneous Euler updates in two-player games, which will allow us to capture discretisation drift effects beyond those of order $\mathcal{O}(h^3)$ we explored in the rest of this chapter. For simplicity, we use assume equal learning rates for the two-players, though one can use the ‘same-physical time approach’ from the previous section to

adapt the results to different learning rates.

The approach we took in the single-objective setting had two steps: first, we developed the principal series with BEA (Theorem 3.3.1) and second, we used the principal series given by BEA to find the PF based on the eigen-decomposition of the Hessian (Corollary 3.1 and Theorem 3.3.2). The first step can be readily translated to games as it applies to any vector field, using Theorem A.1.2 in the Appendix, which Theorem 3.3.1 is a corollary of. For completeness, we reproduce Theorem A.1.2 here:

Theorem 4.8.1 *The modified flow given by BEA with an error of order $\mathcal{O}(h^{p+1})$ to the Euler update $\boldsymbol{\psi}_t = \boldsymbol{\psi}_{t-1} + hu(\boldsymbol{\psi}_{t-1})$ has the form:*

$$\dot{\boldsymbol{\psi}} = \sum_{n=0}^p \frac{(-1)^n}{n+1} \left(\frac{du}{d\boldsymbol{\psi}} \right)^n u + \mathcal{C}\left(\frac{du}{d^2\boldsymbol{\psi}} \right) \quad (4.57)$$

where $\mathcal{C}\left(\frac{du}{d^2\boldsymbol{\psi}} \right)$ denotes a class of functions defined as a sum of individual terms each containing higher than first order derivatives applied to u .

To use the above result in games, consider the two-player game with the original the dynamics in Eqs (4.9)-(4.10):

$$\begin{bmatrix} \dot{\boldsymbol{\phi}} \\ \dot{\boldsymbol{\theta}} \end{bmatrix} = \begin{bmatrix} f \\ g \end{bmatrix}. \quad (4.58)$$

If we define

$$\boldsymbol{\psi} = \begin{bmatrix} \boldsymbol{\phi} \\ \boldsymbol{\theta} \end{bmatrix} \quad u(\boldsymbol{\psi}) = \begin{bmatrix} f(\boldsymbol{\phi}, \boldsymbol{\theta}) \\ g(\boldsymbol{\phi}, \boldsymbol{\theta}) \end{bmatrix} \quad (4.59)$$

and the negative of the Jacobian as

$$\mathbf{H}(\boldsymbol{\phi}, \boldsymbol{\theta}) = - \begin{bmatrix} \frac{df}{d\boldsymbol{\phi}} & \frac{df}{d\boldsymbol{\theta}} \\ \frac{dg}{d\boldsymbol{\phi}} & \frac{dg}{d\boldsymbol{\theta}} \end{bmatrix}, \quad (4.60)$$

we can replace this choice vector field and \mathbf{H} in Theorem 4.8.1, and obtain:

Corollary 4.1 *The full series of BEA constructed from simultaneous Euler updates (Eqs 4.13 and 4.14) with learning rate h exactly is of the form*

$$\begin{bmatrix} \dot{\phi} \\ \dot{\theta} \end{bmatrix} = \sum_{p=0}^{\infty} \frac{1}{p+1} h^p \mathbf{H}(\phi, \theta)^p \begin{bmatrix} f \\ g \end{bmatrix} + \mathcal{C}\left(\frac{df}{d^2\theta}\right). \quad (4.61)$$

We can no longer use the eigen-decomposition of \mathbf{H} , as unlike in the single-objective case ($f = -\nabla_{\phi}E$ and $g = -\nabla_{\theta}E$) we examined in Chapter 3, $\mathbf{H}(\phi, \theta)$ need not be symmetric and thus we can no longer conclude Corollary 3.1. We instead consider the Jordan normal form of \mathbf{H}

$$\mathbf{H}(\phi, \theta) = \mathbf{P}^{-1} \mathbf{J} \mathbf{P}, \quad (4.62)$$

where \mathbf{J} is a block diagonal matrix. We then have

$$\mathbf{H}(\phi, \theta)^p = \mathbf{P}^{-1} \mathbf{J}^p \mathbf{P}, \quad (4.63)$$

leading to

$$\begin{bmatrix} \dot{\phi} \\ \dot{\theta} \end{bmatrix} = \sum_{p=0}^{\infty} \frac{-1}{p+1} \mathbf{P}^{-1} (h\mathbf{J})^p \mathbf{P} \begin{bmatrix} -f \\ -g \end{bmatrix} + \mathcal{C}\left(\frac{df}{d^2\theta}\right), \quad (4.64)$$

where we use the negative of the vector field in order to obtain consistency with the PF results in Chapter 3, and recover single-objective results if $f = -\nabla_{\phi}E$ and $g = -\nabla_{\theta}E$. We are interested in finding a form for the series

$$\sum_{p=0}^{\infty} \frac{-1}{p+1} (h\mathbf{J})^p, \quad (4.65)$$

which if converges can be written as another block diagonal matrix,

$$\frac{1}{h} \log(\mathbf{I} - h\mathbf{J}) \mathbf{J}^{-1}. \quad (4.66)$$

This result leads to the game equivalent of the PF:

Definition 4.8.1 When expanding the PF to simultaneous Euler updates (Eqs 4.13 and 4.14) with learning rate h in two-player games, we obtain:

$$\begin{bmatrix} \dot{\phi} \\ \dot{\theta} \end{bmatrix} = \frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) \mathbf{J}^{-1} \mathbf{P} \begin{bmatrix} -f \\ -g \end{bmatrix} + \mathcal{C}\left(\frac{df}{d^2\theta}\right), \quad (4.67)$$

where $\mathbf{H}(\phi, \theta) = \mathbf{P}^{-1} \mathbf{J} \mathbf{P}$ is the Jordan normal form of \mathbf{H} .

While the Jordan normal form of \mathbf{H} is not unique due to permutations of the order of the blocks in \mathbf{J} , $\mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) \mathbf{J}^{-1} \mathbf{P}$ will be the same regardless of the block order.

4.8.1 Stability analysis

We now perform stability analysis on the modified flow obtained in Eq (4.67), and hope it sheds light on the behaviour of Euler updates—and thus gradient descent—in games. We start by estimating the eigenvalues of the flow's Jacobian:

$$\frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) \mathbf{P}. \quad (4.68)$$

As \mathbf{P} is invertible, the Jacobian's eigenvalues are the same as the eigenvalues of

$$\frac{1}{h} \log(\mathbf{I} - h\mathbf{J}). \quad (4.69)$$

Since $\log(\mathbf{I} - h\mathbf{J})$ is block diagonal, its eigenvalues are the eigenvalues of its blocks $\frac{1}{h} \log(\mathbf{I}_{n_k} - h\mathbf{J}_k)$, where \mathbf{I}_{n_k} is the identity matrix of dimension n_k and $\mathbf{J}_k \in \mathbb{C}^{n_k, n_k}$. Each block \mathbf{J}_k corresponds to a unique eigenvalue λ_k of \mathbf{H} with multiplicity n_k . Given

$$\mathbf{Z}_{n_k} = \begin{bmatrix} 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \quad (4.70)$$

a matrix with all 0s apart from the one off diagonal, \mathbf{J}_k can be written as

$$\mathbf{J}_k = \lambda_k \mathbf{I}_{n_k} + \mathbf{Z}_{n_k}. \quad (4.71)$$

From here

$$\log(\mathbf{I} - h\mathbf{J}_k) = \log(\mathbf{I} - h\lambda_k\mathbf{I}_{n_k} - h\mathbf{Z}_{n_k}), \quad (4.72)$$

and since $\mathbf{I} - h\mathbf{J}_k$ is upper triangular, its eigenvalues are its diagonal entries, which are $1 - h\lambda_k$ and the eigenvalues of $\log(\mathbf{I} - h\mathbf{J}_k)$ are $\log(1 - h\lambda_k)$. Thus, to perform stability analysis we have to assess whether $\Re[\frac{1}{h}\log(1 - h\lambda_k)] < 0$ for all λ_k . Since $h > 0$, this is equivalent to $\Re[\log(1 - h\lambda_k)] < 0$ for all λ_k . Unlike the single-objective case, λ_k might be complex. If we write $\lambda_k = x_k + iy_k$ we have:

$$\Re[\log(1 - h(x_k + iy_k))] = \log \sqrt{(1 - hx_k)^2 + h^2y_k^2}, \quad (4.73)$$

which is negative if $(1 - hx_k)^2 + (hy_k)^2 < 1$.

We now contrast this condition to the exponential stability convergence condition of the original continuous-time game dynamics. Under the original continuous-time dynamics, for a fixed point to be attractive, the real part of the \mathbf{H} 's eigenvalues would need to be strictly positive since \mathbf{H} is the negative of the flow's Jacobian, i.e. $x_k > 0$. When accounting for DD, however, the importance of the imaginary part y_k becomes apparent. Moreover, if we rewrite the condition $(1 - hx_k)^2 + (hy_k)^2 < 1$ as a condition on the maximum learning rate to be convergent, we obtain

$$(1 - hx_k)^2 + (hy_k)^2 < 1 \implies h < \frac{1}{x_k} \frac{2}{1 + \left(\frac{y_k}{x_k}\right)^2}, \quad (4.74)$$

which is the condition Mescheder et al. [15] (their Lemma 4) obtained from a discrete-time perspective. Thus in the game setting too, we recover existing fundamental results using a continuous-time perspective. We note that using this intuition Mescheder et al. [15] derive the algorithm Consensus Optimisation (CO), which we empirically assessed in Section 4.6.

If $f = -\nabla_\phi E$ and $g = -\nabla_\theta E$ we have a single objective, with $y_k = 0 \forall k$. In this case we recover the convergence conditions we are familiar with from the PF, which we derived in Section 3.3.2.

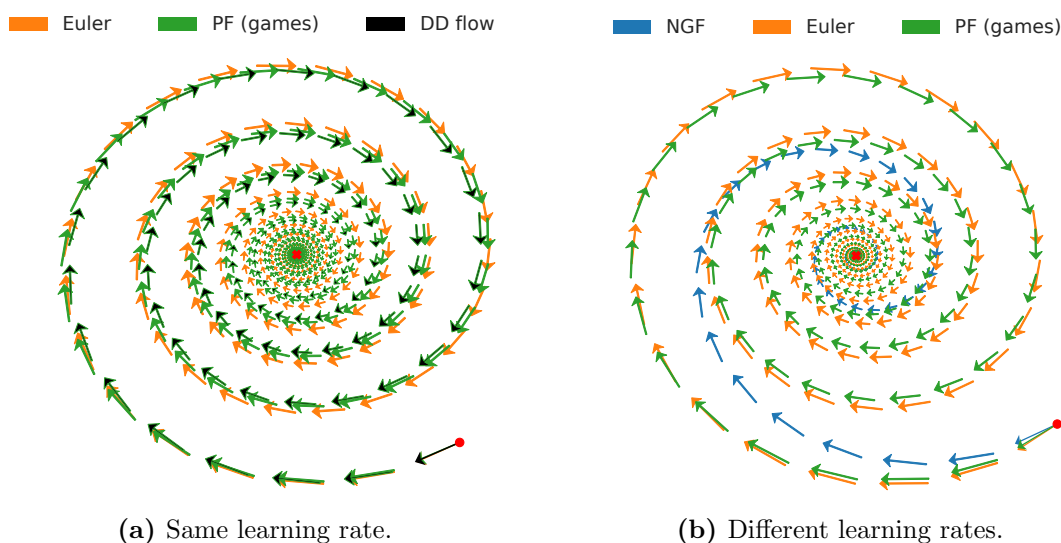


Figure 4.14: A visualisation of a simple example of the extension of the PF to two-games. Here we observe that the PF better tracks the gradient descent update compared to the flow obtained from using the $\mathcal{O}(h^3)$ error, and the NGF.

In the zero-sum case, where $f = -\nabla_{\phi}E$ and $g = \nabla_{\theta}E$, we have seen in Section 2.4 that strict local Nash equilibria are local stable fixed points under the original game dynamics. That is, at a strict local Nash equilibrium we have that $x_k > 0, \forall k$. Under the continuous-time dynamics we derived in this section, which account for DD, this is no longer true, since x_k being positive does not lead to $(1 - hx_k)^2 + (hy_k)^2 < 1$. We thus obtain an analogous result in zero-sum games to that of the single-objective case in Chapter 3, where we saw that unlike the NGF, whether or not the PF is attracted to strict local minima depends on the learning rate h .

4.8.2 An empirical example and different learning rates

The game PF describes the Euler discretisation of the system in Eq (4.23) exactly (see Section B.7 for a full proof). We visualise this example in Figure 4.14a, where we observe that the PF tracks the behaviour of the discrete updates better than the modified third-order correction previously derived. We note that the Jordan normal form (here implemented using ‘sympy’ [182]) is not always numerically stable and the observed (small) difference between the discrete updates and the PF is due to numerical issues.

If different learning rates are used for the two-players, we can use the approach

outlined in Section 4.7 to construct vector fields that account for the each player’s learning rate, and then apply the same analysis as above to the corresponding Jacobian to construct a flow and stability analysis conditions. We use this approach to visualise the same example from Eq (4.23) in Figure 4.14b, and observe that here too the PF tracks the Euler updates closely.

4.9 Related work

Backward error analysis: There have been a number of recent works applying BEA to machine learning in the one-player case, which can potentially be extended to the two-players settings. In particular, Smith et al. [13] extended the analysis of Barrett and Dherin [12] to stochastic gradient descent. Kunin et al. [65] used modified gradient flow equations to show that discrete updates break certain conservation laws present in the gradient flows of deep learning models. França et al. [60] compare momentum and Nesterov acceleration using their modified equations. França et al. [158] used BEA to help devise optimisers with a control on their stability and convergence rates. Li et al. [159] used modified equation techniques in the context of stochastic differential equations to devise optimisers with adaptive learning rates. Other recent works such as Unlu and Aitchison [183] and Jastrzebski et al. [184] noticed the strong impact of implicit gradient regularisation in the training of over-parametrised models with SGD using different approaches. These works provide evidence that DD has a significant impact in deep learning, and that BEA is a powerful tool to quantify it. At last, note that a special case of Theorem 4.2.1 for zero-sum games with equal learning rates can be found in Lu [178].

Two-player games: As one of the best-known examples at the interface of game theory and deep learning, GANs have been powered by gradient-based optimisers as other deep neural networks. The idea of an implicit regularisation effect induced by simultaneous gradient descent in GAN training was first discussed in Schäfer et al. [185]; the authors show empirically that the implicit regularisation induced by gradient descent can have a positive effect on GAN training, and take a game theoretic approach to strengthening it using competitive gradient descent [66]. By

	First player (ϕ)	Second player (θ)
Cancel DD interaction terms (S)	$\frac{v_\phi h}{4} \ \nabla_{\theta} E\ ^2$	$\frac{v_\theta h}{4} \ \nabla_{\phi} E\ ^2$
Cancel DD interaction terms (A)	$\frac{v_\phi h}{4} \ \nabla_{\theta} E\ ^2$	$-\frac{(2v_\phi - v_\theta)h}{4} \ \nabla_{\phi} E\ ^2$
SGA (S)	$\frac{1}{2} \ \nabla_{\theta} E\ ^2$	$\frac{1}{2} \ \nabla_{\phi} E\ ^2$
CO (S)	$\zeta \ \nabla_{\phi} E\ ^2 + \zeta \ \nabla_{\theta} E\ ^2$	$\zeta \ \nabla_{\phi} E\ ^2 + \zeta \ \nabla_{\theta} E\ ^2$
Locally stable GAN (S)	X	$\zeta \ \nabla_{\phi} E\ ^2$
ODE-GAN (S)	$\zeta \ \nabla_{\theta} E\ ^2$	X

Table 4.2: Comparing cancelling discretisation drift terms with explicit regularisation methods in zero-sum games. SGA (without alignment, see the Appendix Section B.5 and Balduzzi et al. 61), Consensus Optimisation (CO) [15], Locally Stable GAN [17], ODE-GAN [67]. We assume a $\min_{\theta} \min_{\phi}$ game with learning rates $v_\phi h$ and $v_\theta h$ and number of updates m and k , respectively. ζ denotes the regularisation coefficient hyperparameter in methods which require one. S and A denote simultaneous and alternating updates. Unlike other approaches, the drift provides us with the coefficients to use to cancel the interaction terms and improve performance, and tackles both alternating and simultaneous updates.

quantifying the DD induced by gradient descent using BEA, we shed additional light on the regularisation effect that discrete updates have on GAN training, and show it has both beneficial and detrimental components (as also shown in Daskalakis and Panageas [180] with different methods in the case of simultaneous GD). Moreover, we show that the explicit regularisation inspired by DD results in drastically improved performance. The form of the modified losses we have derived are related to explicit regularisation for GANs, which has been one of the most efficient methods for stabilising GANs as well as other games. Some of the regularisers constrain the complexity of the players [30, 33, 36], while others modify the dynamics for better convergence properties [17, 61, 67–70]. Our approach is orthogonal in that, we start from understanding the most basic gradient steps, underpinning any further modifications of the losses or gradient fields. Importantly, we discovered a relationship between learning rates and the underlying regularisation. Since merely cancelling the effect of DD is insufficient in practice (as also observed by Qin et al. [67]), our approach complements regularisers that are explicitly designed to improve convergence.

We leave to future research to further study other regularisers in combination with our analysis. We summarise these regularisers in Table 4.2, and contrast them to

our results which cancel the interaction terms of DD, as suggested by Corollaries 4.5.1 and 4.5.2 respectively.

To our knowledge, this is the first work towards finding continuous systems that better match the gradient descent updates used in two-player games. Studying discrete algorithms using continuous systems has a rich history in optimisation [186, 187]. Recently, models that directly parametrise differential equations demonstrated additional potential from bridging the discrete and continuous perspectives [94, 96].

4.10 Conclusion

We have shown that using modified continuous systems to quantify the discretisation drift (DD) induced by gradient descent updates can help bridge the gap between the discrete optimisation dynamics and those of continuous methods. This allowed us to cast a new light on the stability and performance of games trained using gradient descent, and guided us towards explicit regularisation strategies inspired by these modified systems. We note however that DD merely modifies a game’s original dynamics, and that the DD terms alone can not fully characterise a discrete scheme. In this sense, our approach complements works analysing the underlying game [188].

We have focused our empirical analysis on a few classes of two-player games, but the effect of DD will be relevant for all games trained using gradient descent. Our method can be expanded beyond gradient descent to other optimisation methods such as Adam [47], as well as to the stochastic setting as shown in Smith et al. [13]. We hope that our analysis provides a useful building block to further the understanding and performance of two-player games.

Chapter 5

Finding new implicit regularisers by revisiting backward error analysis

In previous chapters, we used Backward Error Analysis (BEA) to construct flows that approximate a discrete optimiser update and shed light on optimisation dynamics, both for single-objectives and two-player games trained with gradient descent. The current usage of BEA is not without limitations, however, since not all the vector fields of continuous-time flows obtained using BEA can be written as a gradient, hindering the construction of modified losses revealing implicit regularisers. Implicit regularisers uncover quantities minimised by following discrete updates and thus can be used to improve performance and stability across problem domains, from supervised learning to two-player games such as Generative Adversarial Networks [12, 13, 151, 185]; we have seen an example in common-payoff and zero-sum games in the previous chapter. In this chapter, we provide a novel approach to use BEA, and show how our approach can be used to construct continuous-time flows with vector fields that can be written as gradients. We then use this to find previously unknown implicit regularisation effects, such as those induced by multiple stochastic gradient descent steps while accounting for the exact data batches used in the updates, and in generally differentiable two-player games.

5.1 Revisiting BEA

In Section 2.3.2, we described BEA and showed how Barrett and Dherin [12] use BEA to find the IGR flow

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E - \frac{h}{2} \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E = -\nabla_{\boldsymbol{\theta}} \left(E(\boldsymbol{\theta}) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})\|^2 \right), \quad (5.1)$$

which has an error of $\mathcal{O}(h^3)$ after one gradient descent update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$. Gradient descent can thus be seen as implicitly minimising the modified loss $E(\boldsymbol{\theta}) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})\|^2$. This showcases an implicit regularisation effect induced by the discretisation drift of gradient descent, dependent on learning rate h , which biases learning towards paths with low gradient norms.

We would like to go beyond full-batch gradient descent and model the implicit regularisation induced by the discretisation of one stochastic gradient descent (SGD). Given the SGD update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}_t)$ —where we denoted $E(\boldsymbol{\theta}_{t-1}; \mathbf{X}) = \frac{1}{B} \sum_{i=1}^B E(\boldsymbol{\theta}_{t-1}; \mathbf{x}_i)$ as the average loss over batch \mathbf{X} and B is the batch size—we can use the IGR flow induced at this time step:

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}_t) - \frac{h}{4} \nabla_{\boldsymbol{\theta}} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}_t)\|^2. \quad (5.2)$$

Since the vector field in Eq (5.2) is the negative gradient of the modified loss $E(\boldsymbol{\theta}; \mathbf{X}_t) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}_t)\|^2$, this reveals a local implicit regularisation which minimises the gradient norm $\|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}_t)\|^2$. It is not immediately clear, however, how to combine the IGR flows obtained for each SGD update in order to model the *combined* effects of multiple SGD updates, each using a different batch. What are, if any, the implicit regularisation effects induced by two SGD steps

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}_t) \quad (5.3)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - h\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t; \mathbf{X}_{t+1})? \quad (5.4)$$

Smith et al. [13] find a modified flow in expectation over the shuffling of batches in an epoch, and use it to find implicit regularisation effects specific to SGD and

study the effect of batch sizes in SGD. Since their approach works in expectation over an epoch, however, it does not account for the implicit regularisation effects of a smaller number of SGD steps, or account for the exact data batches used in the updates; we return to their results in the next section. We take a different approach, and introduce a novel way to find implicit regularisers in SGD by revisiting the BEA proof structure and the assumptions made thus far when using BEA. Our approach can be summarised as follows:

Remark 5.1.1 *Given the discrete update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$, BEA constructs $\dot{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}E + hf_1(\boldsymbol{\theta}) + \dots + h^p f_p(\boldsymbol{\theta})$, such that $\|\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}_t\| \in \mathcal{O}(h^{p+2})$ for a choice of $p \in \mathbb{N}, p \geq 1$. This translates into a constraint on the value of $\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})$. Thus, BEA asserts only what the value of the correction terms f_i in the vector field of the modified flow is at $\boldsymbol{\theta}_{t-1}$. Given the constraints on $f_i(\boldsymbol{\theta}_{t-1})$, we can choose $f_i : \mathbb{R}^D \rightarrow \mathbb{R}^D$ in the vector field of the modified flow $\dot{\boldsymbol{\theta}}$ to depend on the initial condition $\boldsymbol{\theta}_{t-1}$.*

For example, in the proof by construction of the IGR flow (Section 2.3.3), we obtained that if $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ and we want to find $\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}) + hf_1(\boldsymbol{\theta})$ such that $\|\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}_t\|$ is of order $\mathcal{O}(h^3)$, then $f_1(\boldsymbol{\theta}_{t-1}) = -\frac{1}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta}_{t-1})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ (Eq (2.43)). From there, following Barrett and Dherin [12], we concluded that $f_1(\boldsymbol{\theta}) = -\frac{1}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta})$. But notice how BEA only sets a constraint on the value of the vector field at the initial point $\boldsymbol{\theta}_{t-1}$. If we allow the modified vector field to depend on the initial condition, an equally valid choice for $f_1(\boldsymbol{\theta}) = -\frac{1}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ or $f_1(\boldsymbol{\theta}) = -\frac{1}{2}\nabla_{\boldsymbol{\theta}}^2E(\boldsymbol{\theta}_{t-1})\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta})$. By construction, the above flows will also have an error of $\mathcal{O}(h^3)$ after one gradient descent step of learning rate h with initial parameters $\boldsymbol{\theta}_{t-1}$. The latter vector fields only describe the gradient descent update with initial parameters $\boldsymbol{\theta}_{t-1}$ and thus they only apply to this specific gradient descent step, though as previously noted that is also the case with the IGR flow due to the dependence on the data batch—see Eq (5.2). The advantage of constructing modified vector fields that depend on initial parameters lies in the ability to write modified losses when a modified vector field depending only on $\boldsymbol{\theta}$ cannot be written as a gradient operator, as we shall see in the next sections.

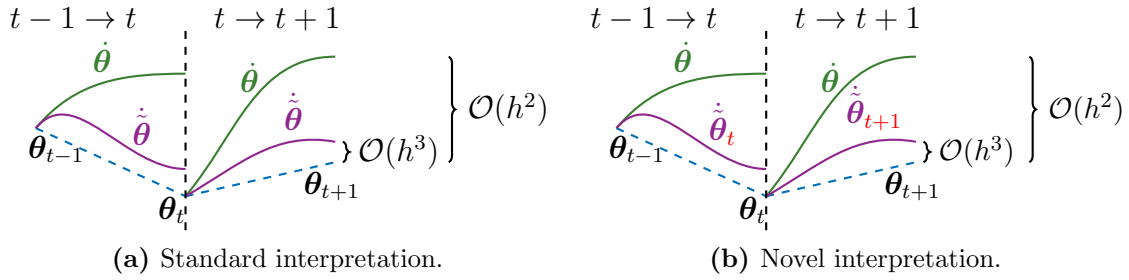


Figure 5.1: Visualising the standard approach to BEA alongside an approach which constructs a different flow per gradient descent iteration. In our previous use of BEA (a), we constructed modified flows $\tilde{\theta}$ to capture the discretisation error of gradient descent; these flows did not depend on the initial iteration parameters. In this chapter, we take the second approach (b), allowing us to construct additional flows $\tilde{\theta}_t$, which depend on initial parameters, and showcase additional implicit regularisation effects.

This observation about BEA leads us to the following remarks:

Remark 5.1.2 *There are multiple flows that lead to the same order in learning rate error after one discrete update. Many of these flows depend on the initial conditions of the system, i.e. the initial parameters of the discrete update. We visualise this approach in Figure 5.1 and contrast it with the existing BEA interpretation.*

Remark 5.1.3 *Implicit in the choice of the IGR flow [12], as well as the flows we have introduced thus far in this work using BEA, namely the PF (Chapter 3) and the modified flows of two-player games (Chapter 4), make an implicit assumption: that we are looking for the modified flows that hold for every training iteration and do not depend on the initial parameters. This is however challenged already in the case of stochastic gradient descent, where each implementation of the flows requires dependence on data batches, as shown in Eq (5.2).*

5.2 Implicit regularisation in multiple stochastic gradient descent steps

We now use the above observations to build modified flows that capture multiple gradient descent updates with an error of $\mathcal{O}(h^3)$ after n SGD steps. We are interested in modified flows that can be used to construct modified losses by writing the vector

field of the flow as the negative gradient of a function; this enables us to capture the implicit regularisation effects of taking multiple SGD steps. We analyse n SGD steps, starting at iteration t

$$\boldsymbol{\theta}_{t+\mu} = \boldsymbol{\theta}_{t+\mu-1} - \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t+\mu-1}; \mathbf{X}^{t+\mu}), \quad \mu \in \{0, \dots, n-1\} \quad (5.5)$$

where for simplicity we denoted $E(\boldsymbol{\theta}_{t-1}; \mathbf{X}) = \frac{1}{B} \sum_{i=1}^B E(\boldsymbol{\theta}_{t-1}; \mathbf{x}_i)$, with elements \mathbf{x}_i forming the batch \mathbf{X} . We start with the following remark:

Lemma 5.2.1 *Denote $E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) = \frac{1}{n} \sum_{\mu=0}^{n-1} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})$, i.e. the average loss obtained from the n data batches. Then the trajectory obtained by taking n steps of stochastic gradient descent follows the trajectory of minimising the loss in continuous-time with $\mathcal{O}(h^2)$ error:*

$$\tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}). \quad (5.6)$$

Perhaps surprisingly, the above shows that effects from mini-batches appear only at higher-order terms in learning rate h . Thus, to capture implicit regularisation effects that account for the mini-batches used in the SGD updates, we use BEA. We find a modified flow that describes the SGD update with an error of $\mathcal{O}(h^3)$, with a vector field that can be written as a negative gradient. We provide proofs and the flows that construct the regularisers in the Appendix C. We use the same steps used for our other BEA proofs: 1) we expand the n discrete updates in Eq (5.5) to write a relation between the parameters at time $t-1$ and $t+n-1$ up to second order in learning rate; 2) expand the changes in continuous-time up to $\mathcal{O}(h^2)$; and 3) match the $\mathcal{O}(h^2)$ terms to find the terms which quantify the drift. The significant difference with our previous approaches is that we allow the vector field to depend on the initial parameters, and choose a vector field that can be written as a gradient in order to construct a modified loss.

Theorem 5.2.2 *Denote $E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) = \frac{1}{n} \sum_{\mu=0}^{n-1} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})$, i.e. the average loss obtained from the n data batches. Then the trajectory obtained by taking n steps of stochastic gradient descent follows the trajectory of minimising the loss in*

continuous-time with $\mathcal{O}(h^3)$ error

$$\tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{nh}{4} \underbrace{\|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2}_{\text{full batch norm regularisation}} \quad (5.7)$$

$$- \frac{h}{n} \sum_{\mu=1}^{n-1} \underbrace{\left[\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})^T \left(\sum_{\tau=0}^{\mu-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}) \right) \right]}_{\text{mini-batch gradient alignment}}. \quad (5.8)$$

We thus find the implicit regularisation effects induced by n steps of SGD, capturing the *importance of exact batches used and their order* in Eq (5.8). We note that without making use the observations regarding BEA in the previous section, and thus without the parameters $\boldsymbol{\theta}_{t-1}$, a modified loss could not have been constructed outside of the full-batch case where one can recover the IGR flow. The following remark immediately follows by setting $n = 2$ in Eq (5.8):

Remark 5.2.1 *When taking a second stochastic gradient descent step, there is an implicit regularisation term maximising the dot product between the gradient at the current step and the gradient at the previous step: $\nabla_{\boldsymbol{\theta}} E^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$. This can be achieved by aligning the direction of the gradients between the two iterations or increasing the gradient norm, but we note that increasing gradient norm is counter other implicit regularisation effect in Eq (5.7).*

We now compare this novel modified loss with the modified loss we obtained by ignoring stochasticity and assuming all n updates have been done with a full-batch; this entails using the IGR loss (proof for multiple steps in Section C.3)

$$\tilde{E} = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2. \quad (5.9)$$

The above modified losses show that both full-batch gradient descent and SGD have a pressure to minimise the gradient norm $\|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|$. SGD leads to an additional regularisation effect capturing the importance of the order in which mini-batches are presented in training: maximising the dot product between gradients computed at the current parameters given a batch and the gradients computed at the

initial parameters *for all batches presented before the given batch*. While this can be achieved both by increasing the norm of the gradients or by aligning gradients with those at the initial iteration, we note that increasing the gradient norm is counter to the other regulariser induced by SGD, the gradient norm minimisation effect shown in Eq (5.7).

The important role of learning rates in BEA appears here too: for our approximations to hold nh has to be sufficiently small. If we adjust the learning rate by the number of updates, i.e. if we set the learning rate for SGD equal to h/n , where h is the learning rate used by full-batch gradient descent, we obtain the same implicit regularisation coefficient for the gradient norm $\|\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|$ minimisation as the IGR flow in Eq 5.9, and the main difference between the two modified losses is given by the mini-batch gradient alignment terms present in Eq (5.8).

The number of updates, n , also plays an important role. While the mini-batch gradient regularisation term in Eq (5.8) has a coefficient of $\frac{h}{n}$, there are $\frac{n(n-1)}{2}$ terms composing the term. Thus the magnitude of the mini-batch gradient regularisation term can grow with n , but its effects strongly depend on the distribution of the gradients computed at different batches. For example, if gradients $\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+i})$ are normally distributed with the mean at the full-batch gradient, i.e. $\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+i}) \sim \mathcal{N}(\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}), \sigma^2)$, as the number of updates grows the regularisation effect in Eq 5.8 will result in a pressure to align mini-batch gradients with the full-batch gradient at the initial parameters, $\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$. Since our results hold for multiple values of n , empirical assessments need to be made to understand the interplay between the number of updates and the strength of mini-batch gradient regularisation on training.

The approach provided here is complementary to that of that of Smith et al. [13], who obtain a relationship similar to Eq (5.8) but in expectation—they describe expected value of the modified loss $\mathbb{E}_{\sigma} [E_{sgd}(\boldsymbol{\theta}; \{\mathbf{X}^{\sigma(t)}, \dots, \mathbf{X}^{\sigma(t+n-1)}\})]$ where the expectation is taken over all possible data batch shufflings σ , *but not the elements in the batch*. As before we denote $E(\boldsymbol{\theta}; \mathbf{X}^k)$ as the loss given by mini-batch k (the equivalent to \hat{C}_k in their notation, see their Eqs (1) and (2)), and write the modified

loss they obtain in expectation as

$$\mathbb{E}_\sigma [E_{sgd}(\boldsymbol{\theta})] = E(\boldsymbol{\theta}; \{\mathbf{X}^0, \dots, \mathbf{X}^{n-1}\}) + \frac{h}{4n} \sum_{k=0}^{n-1} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^k)\|^2 \quad (5.10)$$

$$= E(\boldsymbol{\theta}; \{\mathbf{X}^0, \dots, \mathbf{X}^{n-1}\}) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^0, \dots, \mathbf{X}^{n-1}\})\|^2 \quad (5.11)$$

$$+ \frac{h}{4n} \sum_{k=0}^{n-1} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^k) - \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^0, \dots, \mathbf{X}^{n-1}\})\|^2. \quad (5.12)$$

Our proposed approach does not require working in expectations and accounts for the exact batches sampled from the dataset; thus it directly describes SGD as used in practice. If we make the same assumptions as Smith et al. [13] and take expectation over all possible batch shufflings in an epoch in Eq (5.8) we obtain (proof in Section C.3.1):

$$\mathbb{E}_\sigma [E_{sgd}(\boldsymbol{\theta})] = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (5.13)$$

$$+ \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (5.14)$$

$$- \frac{h}{2n} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (5.15)$$

$$- \frac{h}{2n} \left[\sum_{k=0}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+k})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+k}) \right]. \quad (5.16)$$

As expected, we obtain a different result than that of Smith et al. [13]: while the gradient norm minimisation is still present in our results, the per-batch regularisation in their formulation gets translated into a dot product term, where both the full-batch and per-batch gradients are regularised to be aligned with those at the beginning of the epoch. Both approaches share the limitation that nh needs to be suitably small for approximations to be relevant; we note that since we do not require n to be the number of updates in an epoch—and we obtain interestingly regularisation effects for $n = 2$, see Remark 5.2.1—this is less of an issue for our approach. Their approach has the advantage of finding an implicit regularisation effect that does not depend on the initial parameters. By depending on initial parameters, however, our approach does not require working in expectations and accounts for the exact batches used in

the SGD updates. We hope this can be used to stabilise SGD over multiple steps, as we have done with a single GD step in Section 3.7, that it can be used in continual and transfer learning [71, 72, 189], as well as understanding the effects of the order of examples on model optimisation in online learning.

5.3 Implicit regularisation in generally differentiable two-player games

In Chapter 4, we presented a framework for the quantification of discretisation drift in two-player games. We found distinct modified flows that describe *simultaneous* Euler updates and *alternating* Euler updates. In both cases, we found corrections f_1 and g_1 to the original system

$$\dot{\phi} = f(\phi, \theta) \tag{5.17}$$

$$\dot{\theta} = g(\phi, \theta), \tag{5.18}$$

such that the modified continuous system

$$\dot{\phi} = f(\phi, \theta) + hf_1(\phi, \theta) \tag{5.19}$$

$$\dot{\theta} = g(\phi, \theta) + hg_1(\phi, \theta), \tag{5.20}$$

follows the discrete steps of the method with a local error of order $\mathcal{O}(h^3)$. More precisely, if (ϕ_t, θ_t) denotes the discrete step of the method at time t and $(\phi(h), \theta(h))$ corresponds to the continuous solution of the modified system above starting at $(\phi_{t-1}, \theta_{t-1})$, $\|\phi_t - \phi(h)\|$ and $\|\theta_t - \theta(h)\|$ are of order $\mathcal{O}(h^3)$. In this section, we assume for simplicity that both players use the same learning rate h and simultaneous updates, but the same arguments can be made when they use different learning rates or alternating updates.

Using this framework, we constructed *modified loss functions* in the case of zero-sum (Section 4.5) and common-payoff games (Section 4.4). However, using the aforementioned modified flows, *we cannot always write the vector fields of the modified*

flows as a gradient for differentiable two-player games and thus we cannot construct modified losses. To see why, consider the case where we have two loss functions for the two players respectively $E_\phi(\phi, \theta) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$ and $E_\theta(\phi, \theta) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}$. This leads to the update functions $f = -\nabla_\phi E_\phi$ and $g = -\nabla_\theta E_\theta$. By using $f = -\nabla_\phi E_\phi$ and $g = -\nabla_\theta E_\theta$ in Theorem 4.2.1, we have that for simultaneous gradient descent

$$f_1 = -\frac{1}{2} \frac{df}{d\phi} f - \frac{1}{2} \frac{df}{d\theta} g \quad (5.21)$$

$$= \underbrace{-\frac{1}{4} \nabla_\phi \|\nabla_\phi E_\phi\|^2}_{\text{self term}} - \underbrace{\frac{1}{2} \frac{d\nabla_\phi E_\phi}{d\theta} \nabla_\theta E_\theta}_{\text{interaction term}} \quad (5.22)$$

and similarly

$$g_1 = \underbrace{-\frac{1}{4} \nabla_\theta \|\nabla_\theta E_\theta\|^2}_{\text{self term}} - \underbrace{\frac{1}{2} \frac{d\nabla_\theta E_\theta}{d\phi} \nabla_\phi E_\phi}_{\text{interaction term}}. \quad (5.23)$$

Thus, it is not always possible to write f_1 and g_1 as gradient functions, since the interaction terms—see Definition 4.2.1 for the definition of self and interaction terms—cannot always be written as a gradient. Thus, the modified flows cannot be used to construct modified losses leading to implicit regularisers in generally differentiable two-player games, as we have done for zero-sum games in Chapter 4.

We now use the interpretation of BEA we provided in this chapter, to *choose* two other functions f_1 and g_1 and we will use them to construct another set of modified flows. These flows will still satisfy the value constraints required by the BEA proofs and by construction after one gradient descent update the error between the modified flows and the discrete updates is $\mathcal{O}(h^3)$. Indeed, what the BEA proofs provide for simultaneous gradient descent is (see Eq (B.34) in the Appendix)

$$f_1(\phi_{t-1}, \theta_{t-1}) = -\frac{1}{2} \frac{df}{d\phi}(\phi_{t-1}, \theta_{t-1}) f(\phi_{t-1}, \theta_{t-1}) - \frac{1}{2} \frac{df}{d\theta}(\phi_{t-1}, \theta_{t-1}) g(\phi_{t-1}, \theta_{t-1}) \quad (5.24)$$

$$g_1(\phi_{t-1}, \theta_{t-1}) = -\frac{1}{2} \frac{dg}{d\phi}(\phi_{t-1}, \theta_{t-1}) f(\phi_{t-1}, \theta_{t-1}) - \frac{1}{2} \frac{dg}{d\theta}(\phi_{t-1}, \theta_{t-1}) g(\phi_{t-1}, \theta_{t-1}). \quad (5.25)$$

From here, we can choose f_1 and g_1 , now depending on the iteration number t :

$$f_{1,t}(\boldsymbol{\phi}, \boldsymbol{\theta}) = -\frac{1}{2} \frac{df}{d\boldsymbol{\phi}}(\boldsymbol{\phi}, \boldsymbol{\theta}) f(\boldsymbol{\phi}, \boldsymbol{\theta}) - \frac{1}{2} \frac{df}{d\boldsymbol{\theta}}(\boldsymbol{\phi}, \boldsymbol{\theta}) g(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) \quad (5.26)$$

$$g_{1,t}(\boldsymbol{\phi}, \boldsymbol{\theta}) = -\frac{1}{2} \frac{dg}{d\boldsymbol{\phi}}(\boldsymbol{\phi}, \boldsymbol{\theta}) f(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) - \frac{1}{2} \frac{dg}{d\boldsymbol{\theta}}(\boldsymbol{\phi}, \boldsymbol{\theta}) g(\boldsymbol{\phi}, \boldsymbol{\theta}). \quad (5.27)$$

Here, we treat $g(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})$ and $f(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})$ as constants. We replace $f = -\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}$ and $g = -\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}$ in $f_{1,t}$ and $g_{1,t}$ and write the drift terms as gradient functions:

$$f_{1,t}(\boldsymbol{\phi}, \boldsymbol{\theta}) = -\frac{1}{2} \frac{d\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}}{d\boldsymbol{\phi}} \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}} - \frac{1}{2} \frac{d\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}}{d\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) \quad (5.28)$$

$$= -\frac{1}{2} \nabla_{\boldsymbol{\phi}} \|\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}\|^2 - \frac{1}{2} \nabla_{\boldsymbol{\phi}} (\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\phi}}^T \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})) \quad (5.29)$$

$$= -\nabla_{\boldsymbol{\phi}} \left(\frac{1}{2} \|\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}\|^2 + \frac{1}{2} \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\phi}}^T \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) \right) \quad (5.30)$$

$$g_{1,t}(\boldsymbol{\phi}, \boldsymbol{\theta}) = -\nabla_{\boldsymbol{\theta}} \left(\frac{1}{2} \|\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}\|^2 + \frac{1}{2} \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\theta}}^T \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) \right). \quad (5.31)$$

Replacing the above in the modified flows in Eqs (5.19) and (5.20), we obtain

$$\dot{\boldsymbol{\phi}} = f(\boldsymbol{\phi}, \boldsymbol{\theta}) + hf_1(\boldsymbol{\phi}, \boldsymbol{\theta}) \quad (5.32)$$

$$= -\nabla_{\boldsymbol{\phi}} \left(E_{\boldsymbol{\phi}} + h \left(\frac{1}{2} \|\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}\|^2 + \frac{1}{2} \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\phi}}^T \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) \right) \right) \quad (5.33)$$

$$\dot{\boldsymbol{\theta}} = g(\boldsymbol{\phi}, \boldsymbol{\theta}) + hg_1(\boldsymbol{\phi}, \boldsymbol{\theta}) \quad (5.34)$$

$$= -\nabla_{\boldsymbol{\theta}} \left(E_{\boldsymbol{\theta}} + h \left(\frac{1}{2} \|\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}\|^2 + \frac{1}{2} \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\theta}}^T \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) \right) \right). \quad (5.35)$$

We can now write modified losses for each gradient descent iteration of a general two-player differentiable game (which depend on the iteration t), which describe the local trajectory of simultaneous gradient descent up to $\mathcal{O}(h^3)$:

$$\tilde{E}_{\boldsymbol{\phi},t} = E_{\boldsymbol{\phi}} + h \left(\underbrace{\frac{1}{2} \|\nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}\|^2}_{\text{self term}} + \underbrace{\frac{1}{2} \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\phi}}^T \nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})}_{\text{interaction term}} \right) \quad (5.36)$$

$$\tilde{E}_{\boldsymbol{\theta},t} = E_{\boldsymbol{\theta}} + h \left(\underbrace{\frac{1}{2} \|\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}}\|^2}_{\text{self term}} + \underbrace{\frac{1}{2} \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\theta}}^T \nabla_{\boldsymbol{\phi}} E_{\boldsymbol{\phi}}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})}_{\text{interaction term}} \right). \quad (5.37)$$

The *self terms* result in the implicit gradient regularisation found in supervised learning [12] and zero-sum games (Chapter 4): each player has an incentive to minimise its own gradient norm. The *interaction term* for each player encourages minimising the dot product between the gradients of its loss with respect to the other player’s parameters and the previous gradient update of the other player. Consider the first player’s interaction term, equal to $(-\nabla_{\boldsymbol{\theta}} E_{\phi}(\boldsymbol{\phi}_t, \cdot))^T (-\nabla_{\boldsymbol{\theta}} E_{\theta}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}))$, where $-\nabla_{\boldsymbol{\theta}} E_{\theta}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})$ is the previous update direction of $\boldsymbol{\theta}$ aimed at minimising E_{θ} . Its implicit regularisation effect depends on the functional form of E_{ϕ} and E_{θ} : if $\nabla_{\boldsymbol{\theta}} E_{\theta}$ and $\nabla_{\boldsymbol{\theta}} E_{\phi}$ have aligned directions by construction, then the implicit regularisation effect nudges the first player’s update towards a point in space where the second player’s update changes direction; the opposite is true if $\nabla_{\boldsymbol{\theta}} E_{\theta}$ and $\nabla_{\boldsymbol{\theta}} E_{\phi}$ are misaligned by construction.

5.3.1 The effect of the interaction terms: a GAN example

We work through the regularisation effect of interaction terms for a pair of commonly used GAN losses that do not form a zero-sum game (using the generator non-saturating loss [190]) and contrast it with the zero-sum case (the saturating loss); we have previously investigated these losses in Chapter 4. As before, we denote the first player, the discriminator, as D , parametrised by $\boldsymbol{\phi}$, and the generator as G , parametrised by $\boldsymbol{\theta}$. We denote the data distribution as $p^*(\mathbf{x})$ and the latent distribution $p(\mathbf{z})$. Given the non-saturating GAN loss Goodfellow et al. [32]:

$$E_{\phi}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \mathbb{E}_{p^*(\mathbf{x})} \log D(\mathbf{x}; \boldsymbol{\phi}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})) \quad (5.38)$$

$$E_{\theta}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{z})} -\log D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi}), \quad (5.39)$$

we can then write the interaction term $\nabla_{\boldsymbol{\theta}} E_{\phi}^T \nabla_{\boldsymbol{\theta}} E_{\theta}(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})$ in Eq (5.36) at iteration t as (derivation in Appendix C.4.1)

$$\frac{1}{B^2} \sum_{i,j=1}^B c_{i,j}^{non-sat} \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})^T \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1}), \quad \text{with} \quad (5.40)$$

$$c_{i,j}^{non-sat} = \frac{1}{1 - D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})} \frac{1}{D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})}, \quad (5.41)$$

where \mathbf{z}_t^i is the latent variable with index i in the batch at time t with batch size B . Thus, the strength of the regularisation— $c_{i,j}^{non-sat}$ —depends on how *confident* the discriminator is. In particular, this implicit regularisation encourages the discriminator update into a new set of parameters where the gradient $\nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})$ *points away* from the direction of $\nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})$ when $c_{i,j}$ is large. This occurs for \mathbf{z}_t^i where the discriminator is fooled by the generator—i.e. $1 - D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})$ is close to 0 and $\frac{1}{1 - D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})}$ is large—and samples \mathbf{z}_{t-1}^j where the discriminator was correct at the previous iteration— $D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})$ low and thus $\frac{1}{D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})}$ is large. This can be seen as beneficial regularisation for the generator, by ensuring the update direction $-\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{z})} \frac{1}{D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})} \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})$ is adjusted accordingly to the discriminator’s output. We note, however, that regularisation might not have a strong effect, as gradients $\nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})$ that have a high weight in the generator’s update are those where the discriminator is correct and classifies generated data as fake—i.e. $\frac{1}{D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})}$ is large—but there the factor $\frac{1}{1 - D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi})}$ in the interaction term coefficient in Eq (5.41) will be low. This is inline with our empirical results in Section 4.6.4, where we saw that for the non-saturating loss discretisation drift does not have a strong effect on performance.

We can contrast this with what regularisation we obtain when using the saturating loss [32], where $E_{\boldsymbol{\theta}}(\boldsymbol{\phi}, \boldsymbol{\theta}) = -\mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi}))$. We used the saturating loss extensively in experiments in Section 4.6, and we have shown it performs poorly in comparison to the non-saturating loss, in-line with the literature [32]. If we perform the same analysis as above for the saturating loss we obtain the following implicit regulariser

$$\frac{1}{B^2} \sum_{i,j=1}^B c_{i,j}^{sat} \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})^T \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})), \quad \text{with} \quad (5.42)$$

$$c_{i,j}^{sat} = \frac{1}{1 - D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})} \frac{1}{1 - D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})}. \quad (5.43)$$

Here, $c_{i,j}^{sat}$ is high for \mathbf{z}_t^i and \mathbf{z}_{t-1}^j where the generator was fooling the discriminator—i.e. low $1 - D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})$ and $1 - D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})$. Thus, instead of moving $\nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})$ away from directions where the generator was doing

poorly previously as is the case for the non-saturating loss, *it is moving it away from directions where the generator was performing well*, which could lead to instabilities or loss of performance. Moreover, unlike for the non-saturating loss, the implicit regularisation will have a strong effect on the generator’s update $-\nabla_{\boldsymbol{\theta}} E_{\boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{z})} \frac{1}{1-D(G(\mathbf{z};\boldsymbol{\theta});\boldsymbol{\phi})} \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z};\boldsymbol{\theta});\boldsymbol{\phi})$, since gradients $\nabla_{\boldsymbol{\theta}} D(G(\mathbf{z};\boldsymbol{\theta});\boldsymbol{\phi})$ that have a high weight in the generator’s update are those where the interaction term coefficient in Eq (5.43) is high. This is inline with our results in Section 4.6, showing that in zero-sum games (such as the one induced by the saturating loss), the regularisation induced by the discretisation error of gradient descent is strong, and can hurt performance and stability.

5.4 Conclusion

We provided a novel approach of interpreting backward error analysis and used it to find implicit regularisers induced by gradient descent in the single objective setting and in two-player games. In the single objective case, we found implicit regularisation terms revealing importance of the alignment of gradients at the exact data batches used in multiple steps of stochastic gradient descent, while in two-player games we highlighted the need to examine the game structure in order to determine the effects of implicit regularisation. We hope future work can empirically verify the effects of these implicit regularisers in deep learning. We believe that our observations in this chapter might unlock future research, including verifying whether the dot product alignment regularisation we found explains the benefits of using stochastic gradient descent in deep learning, and what the empirical effect of the interaction terms is on generally differentiable two-player games.

Chapter 6

The importance of model smoothness in deep learning

In previous chapters, we used continuous-time approaches to model gradient descent optimisation dynamics and empirically assessed the effect of these dynamics on neural network training. We have seen how improvements to optimisation lead to increased stability and better performance across problem domains, from supervised learning to two-player games such as GANs. Optimisation also interacts with the choice of model through neural architectures and model regularisation: some neural architectures are easier to optimise than others [115, 191], and model regularisation techniques can improve training stability [23, 30, 85]. For the rest of this thesis, we explore interactions between optimisation and models, and in particular model smoothness and smoothness regularisers. While smoothness quantifies sensitivity to changes in model *inputs* and optimisation is concerned with changes to model *parameters*, we show that in neural network training smoothness regularisation can improve optimisation and that implicit regularisation effects induced by optimisation—such as those discussed in previous chapters—can increase model smoothness.

6.1 Why model smoothness?

How certain should a classifier be when it is presented with out of distribution data? How much mass should a generative model assign around a datapoint? How much should an agent’s behaviour change when its environment changes slightly? Answering these questions shows the need to quantify the manner in which the

output of a model varies with changes in its input, a quantity we will intuitively call the smoothness of the model. In deep learning, neural network smoothness has been shown to boost generalisation and robustness of classifiers and increase the stability and performance of generative models, as well as provide better priors for representation learning and reinforcement learning agents [30, 81–85, 192–194]. After providing an overview of the benefits of smoothness in deep learning and motivated by its impact across problem domains, later in the chapter we investigate some of the previously unexplored effects of smoothness regularisation: strong interactions with optimisation, reduced model capacity, and interactions with data scaling.

6.2 Measuring smoothness

Neural network “smoothness” is a broad, vague, catch all term. We use it to convey formal definitions, such as differentiable, bounded, Lipschitz, as well as intuitive concepts such as invariant to data dimensions or projections, robust to input perturbations, and others. One definition states that a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is n smooth if is n times differentiable with the n -th derivative being continuous. The differentiability of a function is not a very useful inductive bias for a model, as it is both very local and constructed according to the metric of the space where limits are taken. What we are looking for is the ability to choose both the distance metric and how local or global our smoothness inductive biases are. With this in mind, Lipschitz continuity is appealing as it defines a *global* property and provides the choice of distances in the domain and co-domain of f . It is defined as:

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_{\mathcal{Y}} \leq K \|\mathbf{x}_1 - \mathbf{x}_2\|_{\mathcal{X}} \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \quad (6.1)$$

where K is denoted as the Lipschitz constant of function f . To avoid learning trivially smooth functions and maintain useful variability, it is often beneficial to constrain the function variation both from above and below. This leads to bi-Lipschitz continuity:

$$K_1 \|\mathbf{x}_1 - \mathbf{x}_2\|_{\mathcal{X}} \leq \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_{\mathcal{Y}} \leq K_2 \|\mathbf{x}_1 - \mathbf{x}_2\|_{\mathcal{X}}. \quad (6.2)$$

Enforcing Eq (6.1) can be difficult, but according to Rademacher’s theorem [195] if $\mathcal{X} \subset \mathbb{R}^I$ is an open set and $\mathcal{Y} = \mathbb{R}^O$ and f is K -Lipschitz then $\left\| \frac{df(\mathbf{x})}{d\mathbf{x}} \right\|_{op} \leq K$ wherever the total derivative $D_{\mathbf{x}}f(\mathbf{t}) = \frac{df(\mathbf{x})}{d\mathbf{x}}\mathbf{t}$ exists, which is almost everywhere. Conversely, a function f that is differentiable everywhere with $\left\| \frac{df(\mathbf{x})}{d\mathbf{x}} \right\|_{op} \leq K$ is K -Lipschitz. Since the operator norm is not always easily computable, at times the Frobenius norm of the Jacobian $\frac{df(\mathbf{x})}{d\mathbf{x}}$ is constrained instead, since the Frobenius norm bounds the operator norm from above: $\left\| \frac{df(\mathbf{x})}{d\mathbf{x}} \right\|_{op} \leq \left\| \frac{df(\mathbf{x})}{d\mathbf{x}} \right\|_2$; for $O = 1$, $\|\nabla_{\mathbf{x}}f(\mathbf{x})\|_{op} = \|\nabla_{\mathbf{x}}f(\mathbf{x})\|_2$. Thus, a convenient strategy to make a differentiable function K -Lipschitz is to ensure $\left\| \frac{df(\mathbf{x})}{d\mathbf{x}} \right\|_2 \leq K, \forall \mathbf{x} \in \mathcal{X}$.

If f and g are Lipschitz with constants K_f and K_g , $f \circ g$ is Lipschitz with constant $K_f K_g$. Since commonly used activation functions are 1-Lipschitz, a neural network can be made Lipschitz by constraining each learnable layer to be Lipschitz. Many neural networks layers are linear operators (linear and convolutional layers, BatchNormalisation [87]), and to compute their Lipschitz constant we can use that the Lipschitz constant of a linear operator A under common norms such as l_1, l_2, l_∞ is $\sup_{\mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}$ [147]. For the l_2 norm, most commonly used in practice, this leads to finding the spectral norm of A .

6.3 Smoothness regularisation in deep learning

Smoothness regularisers have long been part of the toolkit of the deep learning practitioner: early stopping encourages smoothness by stopping optimisation before the model overfits the training data; dropout [76] makes the network more robust to small changes in the input by randomly masking hidden activations; max pooling encourages smoothness with respect to local changes; L_2 weight regularisation and weight decay [196] discourage large changes in output by not allowing individual weight norms to grow; data augmentation allows us to specify what changes in the input should not result in large changes in the model prediction and thus is also closely related to smoothness and invariance to input transformations. These smoothness regularisation techniques are often introduced as methods that directly target generalisation and other beneficial effects of smoothness discussed in Section 6.4,

instead of being seen through the lens of smoothness regularisation.

Methods that explicitly target smoothness on the entire input space focus on restricting the learned model family. A common approach is to ensure Lipschitz smoothness with respect to the l_2 metric by individually restricting the Lipschitz constant of each layer; for layers that are linear operators, this entails restricting their spectral norm. Spectral regularisation [80] uses the sum of the spectral norms—the largest singular value—of each layer as a regularisation loss to encourage Lipschitz smoothness. Spectral Normalisation [30] ensures the learned models are 1-Lipschitz by adding a node in the computational graph of the model layers by replacing the weights with their normalised version:

$$\mathbf{W} \rightarrow \mathbf{W} / \|\mathbf{W}\|_{op}, \quad (6.3)$$

where $\|\mathbf{W}\|_{op}$ is the spectral norm of \mathbf{W} . Both Spectral Normalisation and Spectral Regularisation use power iteration [197] to compute the spectral norm of weight matrices $\|\mathbf{W}\|_{op}$ for layers which are linear operators, such as convolutional or linear layers. Gouk et al. [198] use a projection method by dividing the weights by the spectral norm after a gradient update. This is unlike Spectral Normalisation, which backpropagates through the normalisation operation. The majority of this line of work has focused on constraints for linear and convolutional layers, and only recently attempts to expand to other layers, such as self attention, have been made [199]. Efficiency is always a concern and heuristics are often used even for popular layers such as convolutional layers [30] despite more accurate algorithms being available [198, 200]. Parseval networks [84] ensure weight matrices of linear and convolutional layers are 1-Lipschitz by enforcing a stronger constraint, an extension of orthogonality to non-square matrices. Bartlett et al. [126] show that any bi-Lipschitz function can be written as a composition of residual layers [191].

Instead of restricting the learned function on the entire space, another approach of targeting smoothness constraints is to regularise the norm of the gradients of a predictor $f : \mathbb{R}^I \rightarrow \mathbb{R}$ with respect to inputs of the network $\|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})\|$, at different *regions of the space* [28, 36, 82, 201, 202]. This is often enforced by adding a gradient

penalty to the loss function $\mathcal{L}(\boldsymbol{\theta})$:

$$\mathcal{L}(\boldsymbol{\theta}) + \zeta \mathbb{E}_{p_{reg}(\mathbf{x})} (\|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})\|^2 - K^2)^2, \quad (6.4)$$

where ζ is a regularisation coefficient, $p_{reg}(\mathbf{x})$ is the distribution at which the regularisation is applied, which can either be the data distribution [82, 201] or around it [28, 202], or, in the case of generative models, at linear interpolations between data and model samples [36]. Gradient penalties encourage the function to be smooth around the support of $p_{reg}(\mathbf{x})$ either by encouraging Lipschitz continuity ($K \neq 0$) or by discouraging drastic changes of the function as the input changes ($K = 0$).

Smoothness for classification tasks is defined by Lassance et al. [192] as preserving features similarities within the same class as we advance through the layers of the network. The penalty used is $\sum_{l=1}^L \sum_{c=1}^C |\sigma^l(s_c) - \sigma^{l+1}(s_c)|$, where $\sigma^l(s_c)$ is the signal of features belonging to class c computed using the Laplacian of layer l . The Laplacian of a layer is defined by constructing a weighted symmetric adjacency matrix of the graph induced by the pairwise most similar layer features in the dataset.

6.4 The benefits of learning smooth models

Generalisation. Learning models that generalise beyond training data is the goal of machine learning. The study of generalisation has long been connected with the study of model complexity, as models with small complexity generalise better [203]. Despite this, we have seen that deep, over-parametrised neural networks tend to generalise better than their shallow counterparts [83] and that for Bayesian methods, Occam’s razor does not apply to the number of parameters used, but to the complexity of the function [204]. A way to reconcile these claims is to incorporate smoothness into definitions of model complexity and to show that smooth, over-parametrised neural networks generalise better than their less smooth counterparts. Methods that encourage smoothness—such as weight decay, dropout, and early stopping—have been long shown to aid generalisation [76, 79, 205–207]. Data augmentation has been shown to increase robustness to random noise or to modality specific transformations, such as image cropping and rotations [208–210]. Sokolić et al. [82]

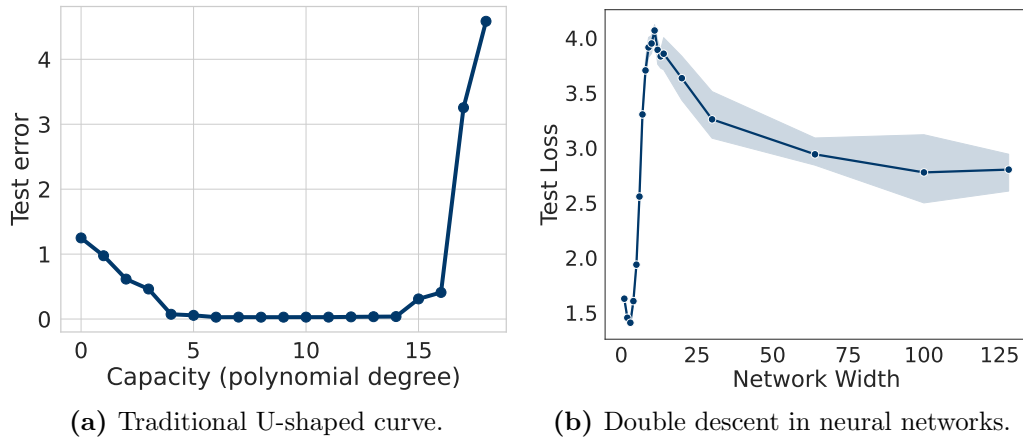


Figure 6.1: Double descent in deep learning compared to the traditional U-shaped complexity curve. (a): traditional U-shaped curve associated with increased capacity; as we increase capacity, the test error first decreases after which it increases due to overfitting—the example we show here is that from Figure 1.2. Figure 6.1b: when the complexity of neural networks is measured using the number of model parameters a double descent is observed, where a first descent is followed by a second descent.

show that the generalisation error of a network with linear, softmax, and pooling layers is bounded by the classification margin in input space. Since classifiers are trained to increase classification margins in output space, smoothing by bounding the spectral norm of the model’s Jacobian increases generalisation performance; this leads to empirical gains on standard image classification tasks. Bartlett et al. [81] provide a generalisation bound depending on classification margins and the product of spectral norms of the model’s weights and show how empirically the product of spectral norms correlates with excess risk.

Generalisation in deep learning has been recently re-examined under the light of double descent [79, 211], a phenomenon named after the shape of the generalisation error plotted against the size of a deep neural network: as the size of the network increases the generalisation error decreases (first descent), then increases, after which it decreases again (second descent). This is unlike the traditional U-shaped curve associated with non-neural models, where the second descent does not occur. We plot the traditional U-shaped curve in Figure 6.1a and contrast it with the double descent phenomenon shown in Figure 6.1b. We postulate there is a strong connection between double descent and smoothness: in the first descent, the generalisation

error is decreasing as the model is given extra capacity to capture the decision surface; the increase happens when the model has enough capacity to fit the training data, but it cannot do so and retain smoothness; the second descent occurs as the capacity increases and smoothness can be retained. This view of double descent is supported by empirical evidence that shows that its effect is most pronounced on clean label datasets and when early stopping and other regularisation techniques are not used [79]. Later in this chapter, we show that smoothness constraints heavily interact with optimisation. This further suggests that empirical investigations into the impact of smoothness on the double descent phenomenon are needed, since implicit regularisation effects induced by optimisation might act as smoothness regularisers affecting model complexity. We will come back to this question in Chapter 8.

Reliable uncertainty estimates. Neural networks trained to minimise classification losses provide notoriously unreliable uncertainty estimates; an issue which gets compounded when the networks are faced with out of distribution data. However, one can still leverage the power of neural networks to obtain reliable uncertainty estimates by combining smooth neural feature learners with non-softmax decision surfaces [212, 213]. The choice of smoothness regularisation or classifier can vary, from using gradient penalties on the neural features with a Radial Basis Function classifier [212], to using Spectral Normalisation on neural features and a Gaussian Process classifier [213]. These methods are competitive with standard techniques used for out-of-distribution detection [214] on both vision and language understanding tasks. The importance of smoothness regularising neural features indicates that having a smooth decision surface, such as a Gaussian Process, is not sufficient to compensate for sharp feature functions when learning models for uncertainty estimation [214].

Robustness to adversarial attacks. Adversarial robustness has become an active area of research in recent years [190, 215–217]. Early works have observed that the existence of adversarial examples is related to the magnitude of the gradient of the hidden network activation with respect to its input, and suggested that constraining the Lipschitz constant of individual layers can make networks more robust to attacks [215]. Initial approaches to combating adversarial attacks, however, focused

on data augmentation methods [190, 218–220], and only more recently smoothness has come into focus [82–84, 192]. We can see the connection between smoothness and robustness by looking at the desired robustness properties of classifiers, which aim to ensure that inputs in the same ϵ -ball result in the same function output:

$$\|\mathbf{x} - \mathbf{x}'\| \leq \epsilon \implies \arg \max f(\mathbf{x}) = \arg \max f(\mathbf{x}'). \quad (6.5)$$

The aim of adversarial defences and robustness techniques is to have ϵ be as large as possible without affecting classification accuracy. Robustness against adversarial examples has been shown to correlate with generalisation [221], and with the sensitivity of the network output with respect to the input as measured by the Frobenius norm of the Jacobian of the learned function [82, 83]. Lassance et al. [192] show that robustness to adversarial examples is enhanced when the function approximator is smooth as defined by the Laplacian smoothness signal discussed in Section 6.3. Tsuzuku et al. [200] show that Eq (6.5) holds when the l_2 norm is used if ϵ is smaller than the ratio of the classification margin and the Lipschitz constant of the network times a constant, and thus they increase robustness by ensuring the margin is larger than the Lipschitz constant.

Improved generative modelling performance. Smoothness constraints through gradient penalties or Spectral Normalisation have become a recipe for obtaining state of the art generative models. In Generative Adversarial Networks (GANs) [32], smoothness constraints on the discriminator and the generator have played a big part in scaling up training on large, diverse image datasets at high resolution [33] and a combination of smoothness constraints has been shown to be a requirement to get GANs to work on discrete data such as text [41]. The latest variational autoencoders [222, 223] incorporate Spectral Regularisation to boost performance and stability [85]. Explicit likelihood tractable models like normalizing flows [224] benefit from smoothness constraints through powerful invertible layers built using residual connections $g(\mathbf{x}) = \mathbf{x} + f(\mathbf{x})$ where f is Lipschitz [225].

More informative critics. Critics, learned approximators to intractable decision functions, have become a fruitful endeavour in generative modelling, repre-

sentation learning, and reinforcement learning.

Critics are used in generative modelling to approximate divergences and distances between the learned model and the true unknown data distribution, and have been mainly popularised by GANs. A critic in a function class \mathcal{F} can be used to approximate the KL divergence by minimizing the bound [226, 227]:

$$\text{KL}(p||q) = \mathbb{E}_{p(\mathbf{x})} \log \frac{p(\mathbf{x})}{q(\mathbf{x})} = \sup_f \mathbb{E}_{p(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q(\mathbf{x})} e^{f(\mathbf{x})-1} \quad (6.6)$$

$$\geq \sup_{f \in \mathcal{F}} \mathbb{E}_{p(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q(\mathbf{x})} e^{f(\mathbf{x})-1}. \quad (6.7)$$

Due to the density ratio $p(\mathbf{x})/q(\mathbf{x})$ in its definition, the KL divergence provides no learning signal when the model and data distributions do not have overlapping support. Choosing \mathcal{F} to be a family of smooth functions in Eq (6.7), however, results in a bound on the KL that provides useful gradients and can be used to train a model, even when the two distributions do not have overlapping support [28, 228]. We show an illustrative example in Figure 6.2a: the true decision surface jumps from zero to infinity, while the approximation provided by the MLP is smooth. Similarly, training the critic more and making it better at estimating the true decision surface but less smooth can hurt training [185]. It’s not surprising that imposing smoothness constraints on critics has become part of many flavours of GANs [27, 28, 33, 36, 80, 201, 229].

The same conclusions have been reached in unsupervised representation learning, where parametric critics are trained to approximate another intractable quantity, the mutual information, using the Donsker–Varadhan or similar bounds [230, 231]. An extensive study on representation learning techniques based on mutual information showed that tighter bounds do not lead to better representations [194]. Instead, the success of these methods is attributed to the inductive biases of the critics employed to approximate the mutual information. In reinforcement learning, neural function approximators or neural “critics” approximate state-value functions or action-state value functions and are then used to train a policy to maximise the expected reward. Directly learning a neural network parametric estimator of the action value *gradients*

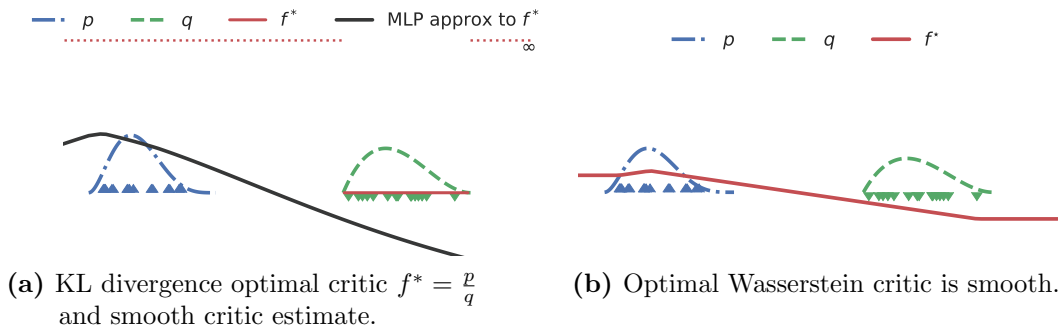


Figure 6.2: The importance of critic smoothness when estimating divergences and distances. (a): When the two distributions do not have overlapping support, the KL divergence provides no learning signal, while a smooth approximation via a learned critic does. (b): The optimal Wasserstein critic has a smoothness Lipschitz constraint in its definition.

—the gradients of the action value with respect to the action—results in more accurate gradients (Figure 3 in [193]), but also makes gradients smoother. This provides an essential exploration prior in continuous control, where similar actions likely result in the same reward and observing the same action twice is unlikely due to size of the action space; encouraging the policy network to extrapolate from the closest seen action improves performance over both model free and model based continuous control approaches [193].

Distributional distances. Including smoothness constraints in the definition of distributional distances by using optimal transport has seen a uptake in machine learning applications in recent years, from generative modelling [27, 36, 232, 233] to reinforcement learning [117, 234], neural flows [235] and fairness [236, 237]. Optimal transport is connected to Lipschitz smoothness as the Wasserstein distance can be computed via the Kantorovich-Rubinstein duality [238]:

$$W_1(p(\mathbf{x}), q(\mathbf{x})) = \sup_{f: \|f\|_{\text{Lip}} \leq 1} \mathbb{E}_{p(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{q(\mathbf{x})} f(\mathbf{x}). \quad (6.8)$$

The Wasserstein distance finds the critic that can separate the two distributions in expectation, but constraints that critic to be Lipschitz in order to avoid pathological solutions. The importance of the Lipschitz constraint on the critic can be seen in Figure 6.2b: unlike the KL divergence, the optimal Wasserstein critic is well defined

when the two distributions do not have overlapping support, and does not require an approximation to provide useful learning signal for a generative model.

6.5 Challenges with smoothness regularisation

Despite the many benefits of learning smooth models we outlined in the previous section, many of the effects of smoothness regularisation have not been extensively studied. In this section, we lay the ground work for the next chapters by looking at some of the unexplored effects of smoothness regularisation, from reduced model capacity to interactions with optimisation. Figure reproduction details for this chapter are available in Appendix Chapter D.1.

6.5.1 Too much smoothness hurts performance

Needlessly limiting the capacity of our models by enforcing smoothness constraints is a significant danger: a constant function is very smooth, but not very useful. Beyond trivial examples, Jacobsen et al. [239] show that one of the reasons neural networks are vulnerable to adversarial perturbations is invariance to task relevant changes—too much smoothness with respect to the wrong metric. A neural network can be “too Lipschitz”: methods aimed at increasing robustness to adversarial examples do indeed decrease the Lipschitz constant of a classifier, but once the Lipschitz constant becomes too low, accuracy drops significantly [240].

There are two main avenues for being too restrictive in the specification of smoothness constraints, depending on *where* and *how* smoothness is encouraged. Smoothness constraints can be imposed on the entire input space or only in certain pockets, often around the data distribution. Methods that impose constraints on the entire space throw away useful information about the input distribution and restrict the learned function needlessly by forcing it to be smooth in areas of the space where there is no data. This is especially problematic when the input lies on a small manifold in a large dimensional space, such as in the case of natural images, which are a tiny fraction of the space of all possible images. Model capacity can also be needlessly restrained by imposing strong constraints on the individual components of the model, often the network layers, instead of allowing the network to allocate

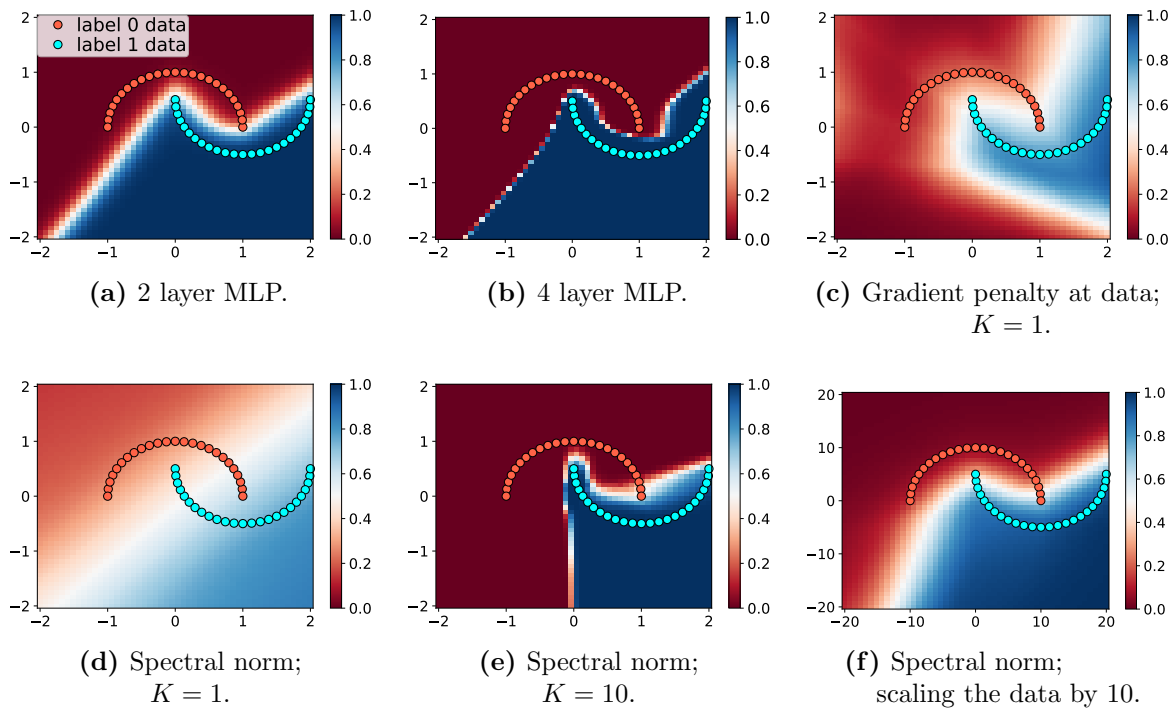


Figure 6.3: Decision surfaces on *two moons* under different regularisation methods.

Unless otherwise specified the model architecture is a 4 layer MLP. We observe that network depth can decrease smoothness (b); that data dependent smoothness penalties can increase smoothness (c); that hard smoothness constraints can lead to smooth surfaces but hurt performance (d), but that this effect can be mitigated by increasing the Lipschitz constant of the constraint (e) or scaling the range of the input data (f).

capacity as needed.

We can exemplify the importance of where and how constraints are imposed with an example, by contrasting gradient penalties—end to end regularisation applied around the training data—and Spectral Normalisation—layer-wise regularisation applied to the entire space. Figure 6.3d shows that using Spectral Normalisation to restrict the Lipschitz constant of an MLP to be 1 decreases the capacity of the network and severely affects accuracy compared to the baseline MLP—Figure 6.3b—or the MLP regularised using gradient penalties—Figure 6.3c. Further insight comes from Figure 6.4, which shows that the gradient penalty only enforces a weak constraint on the model and does not heavily restrict the spectral norms of individual layers; this is in stark contrast with Spectral Normalisation, which, by construction, ensures each network layer has spectral norm equal to 1. To show the effect of data dependent

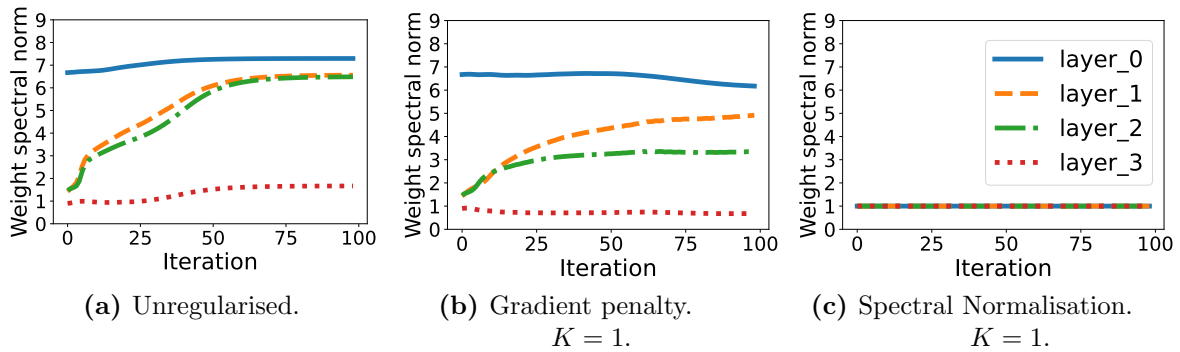


Figure 6.4: Lipschitz constant of each layer of an MLP trained on the two moons dataset. The decision surfaces for the same models can be seen in Figure 6.3. Smaller means smoother. Soft regularisation via the gradient penalty leads to a decrease in spectral norm of the model’s weight matrices, with each layer adjusting differently to the constraint. In contrast, for Spectral Normalisation the spectral norm of the weight matrix for every layer is equal to 1.

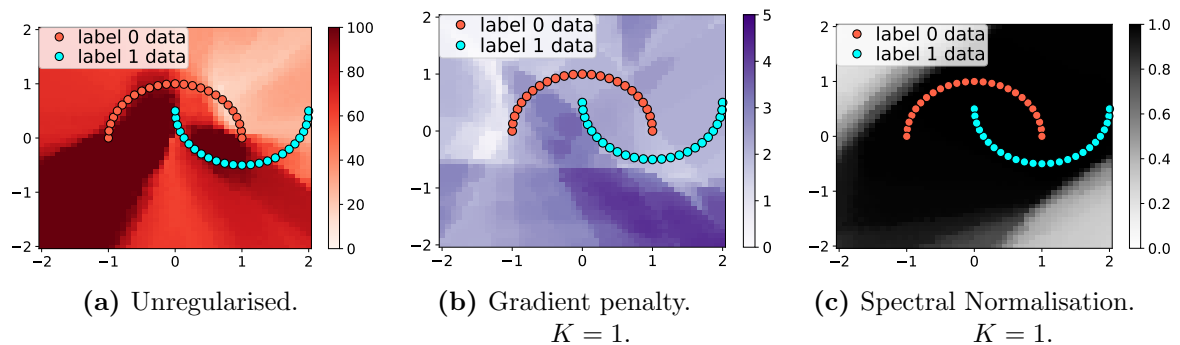


Figure 6.5: The effect of regularisation on *local* smoothness. We plot the local Lipschitz constants obtained using an exhaustive grid search in *local neighbourhoods*, instead of loose bounds. We use different colours to emphasise the different scale of the constants for the different methods. While gradient penalties increase smoothness, they do so primarily around the data, where the constraint is active, while Spectral Normalisation affects a large part of the input space.

regularisation on *local smoothness* we plot the Lipschitz constants of the model at neighbourhoods spanning the entire space in Figure 6.5. Each Lipschitz constant is computed using an exhaustive grid search inside each local neighbourhood rather than a bound—details are provided in Appendix D.1. As expected, gradient penalties impose stronger constraints around the training data, while Spectral Normalisation has a strong effect on the smoothness around points in the entire space. This simple example suggests that the search for better smoothness priors needs to investigate

where we want functions to be smooth and re-examine *how* smoothness constraints should account for the compositional aspect of neural networks, otherwise we run the risk of learning trivially smooth functions. In the next chapter, we will first hand encounter the need to avoid enforcing too strong smoothness constraints when applying Spectral Normalisation to reinforcement learning.

6.5.2 Interactions between smoothness regularisation and optimisation

We show that viewing smoothness only through the lens of the model is misleading, as smoothness constraints have a strong effect on optimisation. The interaction between smoothness and optimisation has been mainly observed when training generative models; encouraging the smoothness of the encoder through Spectral Regularisation increased the stability of hierarchical VAEs and led variational inference models to the state of the art of explicit likelihood non autoregressive models [85], while smoothness regularisation of the critic (or discriminator) has been established as an indispensable stabiliser of GAN training, independently of the training criteria used [28, 33, 80, 201, 241].

Some smoothness regularisation techniques affect optimisation by changing the loss function (gradient penalties, Spectral Regularisation) or the optimisation regime directly (early stopping, projection methods). Even if they don't explicitly change the loss function or optimisation regime, smoothness constraints affect the path the model takes to reach convergence. We use a simple example to show why smoothness regularisation interacts with optimisation in Figure 6.6. We use different learning rates to train two unregularised MLP classifiers on MNIST [242] and observe that the learning rate used affects its smoothness throughout training, without changing testing accuracy. Here, we measured the model's smoothness as the upper bound on its Lipschitz constant given by product of spectral norms of learnable layers. These results show that imposing similar smoothness constraints on two models that share the same architecture but are trained with different learning rates would lead to very different strengths of regularisation and drastically change the trajectory of optimisation. This is particularly true for methods such as Spectral Normalisation,

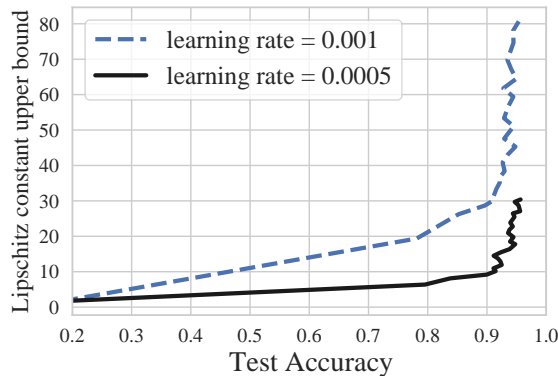


Figure 6.6: Smoothness constraints interact with learning rates. MNIST classifiers trained with different learning rates exhibit different smoothness behaviour throughout training, even at the same test accuracy. We measured the model’s smoothness as the upper bound on its Lipschitz constant given by product of spectral norms of learnable layers. This example shows that imposing the same smoothness constraint on the same model architecture trained with different learning rates will have different effects on optimisation dynamics. This can be particularly the case for methods that impact the spectral norms of individual layers we measured here.

which directly impact the spectral norms of individual layers we measured here.

Beyond learning rates, smoothness constraints also interact decay rate hyperparameters in optimisers that use momentum. When the Adam optimiser is used, this is the β_1 hyperparameter; for an overview of Adam and momentum we refer the reader to Section 2.2. When training GANs, Gulrajani et al. [36] observed that weight clipping in the Wasserstein GAN critic requires low to no momentum. Weight clipping has since been abandoned in the favour of other methods, but as we show in Figure 6.7, current methods like Spectral Normalisation applied to GAN critics trained with low momentum decrease sensitivity to learning rates but perform poorly in conjunction with high momentum, leading to slower convergence and higher hyperparameter sensitivity.

We have shown that smoothness constraints interact with optimisation parameters, such as learning rates and momentum, and argue that we need to reassess our understanding of smoothness constraints, not only as constraints on the final model, but as methods that influence the *optimisation path*. We will study this in detail in the next chapter.

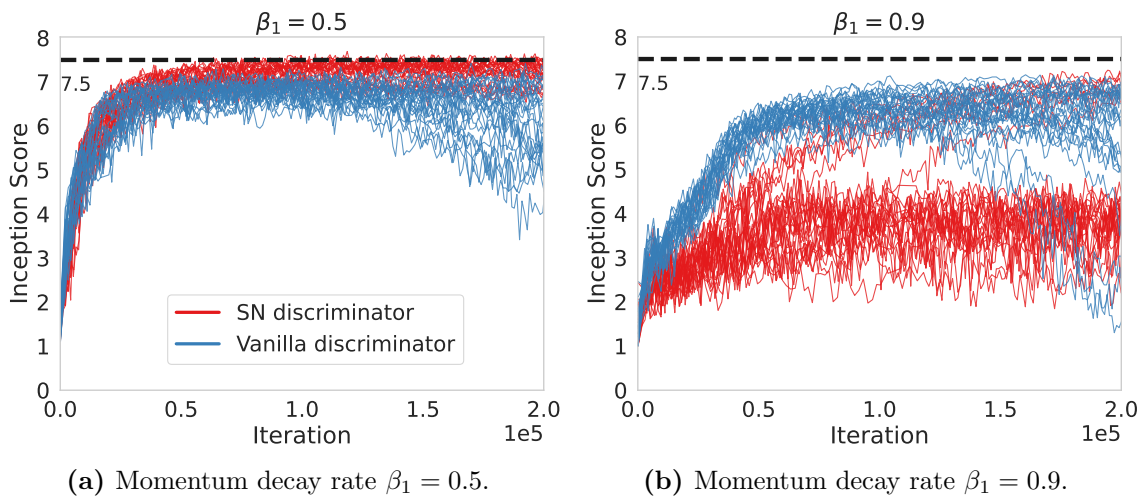


Figure 6.7: Smoothness constraints interact with momentum decay rates: Spectral Normalisation requires low momentum in GAN training. The effects of momentum decay rate β_1 in the Adam optimiser on Spectral Normalisation applied to the GAN discriminator on CIFAR-10, over a sweep of learning rates. Higher is better. Individual learning curves are obtained from a sweep over seeds and learning rates.

6.5.3 Sensitivity to data scaling

Sensitivity to data scaling of smoothness constraints can make training neural network models sensitive to additional hyperparameters. Let f^* be the optimal decision function for a task obtained from using i.i.d. samples from random variable X , and f_c^* obtained similarly from i.i.d. samples obtained from cX . Since f^* and f_c^* can be highly non linear, the relationship between the smoothness of the two functions is unclear. This gets further complicated when we consider their closest approximators under a neural family. The effect of data scaling on the smooth constraints required to fit a model can be exemplified using the two moons dataset: with a Lipschitz constraint of 1 on the model the data is poorly fit—Figure 6.3d—but a much better fit can be obtained by changing the Lipschitz constant to 10—Figure 6.3e—or scaling the data—Figure 6.3f.

6.6 Conclusion

In this chapter, we outlined the many benefits of neural network smoothness across multiple problem domains. Motivated by these benefits and the wide use of smooth-

ness regularisation in deep learning, we highlighted areas where the effects of smoothness have not been fully understood, from connections to the double descent phenomenon to interactions between smoothness regularisation and optimisation. We continue this line of inquiry in the following chapters.

Chapter 7

Smoothness and optimisation effects of Spectral Normalisation

Despite the many benefits of smoothness in supervised learning and generative modelling outlined in the previous chapter, smoothness regularisation has only been briefly studied in deep reinforcement learning (DRL). In this chapter, we study the effects of Spectral Normalisation, a common smoothness regularisation method described in Section 6.3, on DRL. We use Spectral Normalisation as it has shown tremendous success in GANs, and DRL and GANs share common characteristics, such as operating in non-stationary settings and challenges with optimisation [15, 30, 243, 244]. We find that Spectral Normalisation can lead to improved performance in DRL, as long as it is only applied selectively to a few layers in order to avoid the reduced capacity effects discussed in Section 6.5.1. Building upon results in Section 6.5.2, we find interactions between Spectral Normalisation and optimisation dynamics, and show that in DRL the improvements brought by Spectral Normalisation can be obtained by optimisation only changes. Equipped with this knowledge, we reinvestigate the effect of Spectral Normalisation on GANs, and observe that there Spectral Normalisation plays a compounded role, with model and optimisation changes interacting with the choice of loss function in intricate ways.

7.1 Introduction

Reinforcement Learning (RL) has made astounding progress in recent years, showcasing the ability to play games where alternatives, such as handcrafted search or

brute force approaches, are intractable [8, 10, 245]. While this progress could not be achieved without the feature learning aspect of DRL due to a prohibitively large state space or number of actions, research focus has been on RL-specific algorithmic improvements, such as objective functions, exploration strategies, and replay prioritisation [117, 246–251]. We take a different approach, and focus on the *deep* aspect of DRL, by investigating the interplay between training neural networks and RL objectives. We do so by exploring the use Spectral Normalisation in DRL; as we have done in previous chapters, we hone in on interactions between smoothness regularisation and optimisation. To this effect, we make the following contributions:

- We show that Spectral Normalisation, *when applied to a subset of layers*, can improve performance of DRL agents and match the performance of Rainbow [109], a competitive agent combining multiple RL algorithmic advances.
- We show that the effect of selectively applying of Spectral Normalisation to a subset of layers is not solely a smoothness one, but an optimisation one, and can decrease sensitivity to optimisation hyperparameters.
- We introduce a set of adaptations to optimisation methods based on Spectral Normalisation, which *without changing the model*, can recover and sometimes improve the results of Spectral Normalisation in DRL.
- We use our new insights to re-examine the effects of Spectral Normalisation in GANs, where we observe that normalising the discriminator avoids loss areas of low or high gradient magnitudes, and that Spectral Normalisation interacts with optimisation hyperparameters in important ways.

7.2 Reinforcement learning

RL translates the problem of acting sequentially in an environment in order to maximise reward into a formal optimisation problem, usually defined via Markov Decision Processes (MDP). An MDP is a tuple $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ with \mathcal{S} a set of available states, \mathcal{A} a set of available actions (which throughout this thesis we will assume to be discrete), P a conditional transition probability or density $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^+$,

such that $P(s'|s, a)$ defines the density to transition to state s' from state s if action a is taken by the agent, as well as the unconditional $P(s_0)$ defining the density over initial states. For our use cases, we assume that P is not available in closed form but can be sampled from by obtaining $s' \sim P(\cdot|s, a)$ from an environment simulator. The reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ defines the reward $r(s, a)$ obtained by the agent if it takes action a from state s and the discount factor $\gamma \in \mathbb{R}$, with $\gamma \in (0, 1)$. The goal of RL is to find a policy π defined as the conditional probability distribution $\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^+$, where $\pi(a|s)$ encodes the probability of taking action a if in state s under policy π , such that the time discounted reward over an episode is maximised:

$$\max_{\pi} \mathbb{E}_{P(s_0)} \mathbb{E}_{\pi(a_0|s_0)} \left[r(s_0, a_0) + \sum_{t=1}^{\infty} \mathbb{E}_{P(s_t|s_{t-1}, a_{t-1})} \mathbb{E}_{\pi(a_t|s_t)} \gamma^t r(s_t, a_t) \right]. \quad (7.1)$$

Directly optimising Eq (7.1) via Monte Carlo estimation can be challenging as it can lead to high variance estimates of the objectives, as well as slow learning in the case of long episodes [170]. Instead, it is common to define *state-action value functions* as

$$Q_{\pi}(s, a) = r(s, a) + \sum_{\substack{t=1, \\ a_0=a, s_0=s}}^{\infty} \mathbb{E}_{P(s_t|s_{t-1}, a_{t-1})} \mathbb{E}_{\pi(a_t|s_t)} \gamma^t r(s_t, a_t). \quad (7.2)$$

State-action value functions satisfy the recurrence relation

$$Q_{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{P(s'|s, a)} \mathbb{E}_{\pi(a'|s')} Q_{\pi}(s', a'), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}, \quad (7.3)$$

known as the Bellman equation. The objective in Eq (7.1) can now be written as

$$\max_{\pi} \mathbb{E}_{P(s_0)} \mathbb{E}_{\pi(a_0|s_0)} Q_{\pi}(s_0, a_0). \quad (7.4)$$

If we denote the optimal policy as π^* , another Bellman recurrence equation can also be used to find the Q_{π^*} , namely the Bellman optimality equation

$$Q_{\pi^*}(s, a) = r(s, a) + \gamma \mathbb{E}_{P(s'|s, a)} \max_{a'} Q_{\pi^*}(s', a'), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}. \quad (7.5)$$

Iterative updates ensuring Bellman optimality consistency of the current estimate Q of Q_{π^*} form the basis for Q -learning algorithms [252]. Since $P(\cdot|s, a)$ is unknown, the expectation in Eq (7.5) is replaced with a Monte Carlo estimate obtained by performing action a in state s and obtaining s' from the environment; actions are often obtained using the ϵ -greedy policy

$$\pi(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a = \arg \max_{a'} Q(s, a') \\ \epsilon / (|\mathcal{A}| - 1), & \text{otherwise} \end{cases} \quad (7.6)$$

and thus the Bellman optimality updates are evaluated at tuples $(s, a \sim \pi(\cdot|s), s' \sim P(\cdot|s, a), r = r(s, a))$.

We note that while here we focus on value based, model-free, offline algorithms, other popular formulations of RL exist, and investigating the effect of smoothness for other algorithm families, including actor-critic and model-based RL, is a worthy area of study; for an overview of RL algorithms, see Sutton and Barto [170].

7.2.1 Deep reinforcement learning

While in traditional RL, often called tabular RL, when estimating a state-action value function individual values are approximated for each $s \in \mathcal{S}$ and $a \in \mathcal{A}$, in DRL $Q(s, a)$ is modelled using a neural network. The advantage of this approach is two fold: first, it makes large state and action spaces—which otherwise would be prohibitive due to memory requirements of storing a value for each state and each action—tractable; and second, it benefits from the generalisation provided by neural network features by learning a similarity measure between states.

One approach to DRL is DQN [107], which adapts the Q -learning algorithm to the neural setting, with the objective

$$\min_{\theta} E(\theta) = \mathbb{E}_{(s,a,s',r) \sim \mathcal{R}} \left(Q(s, a; \theta) - \left(r + \gamma \max_{a'} Q(s', a'; \theta) \right) \right)^2, \quad (7.7)$$

where \mathcal{R} is a replay buffer that stores experiences obtained by acting according to an ϵ greedy policy induced by Q . Since computing the maximum operator requires computing $Q(s, a) \quad \forall a \in \mathcal{A}$, the neural critic implementation takes as input

a state s and predicts using a forward pass all state-action value functions, thus implementing $Q(\cdot; \boldsymbol{\theta}) : \mathcal{S} \rightarrow \mathbb{R}^{|\mathcal{A}|}$. Due to instabilities in training, the target value $r + \gamma \max_{a'} Q(s', a'; \boldsymbol{\theta})$ is often replaced with a target network $r + \gamma \max_{a'} Q(s', a'; \boldsymbol{\theta}_{old})$, where $\boldsymbol{\theta}_{old}$ are a set of parameters obtained previously in training. The target network parameters $\boldsymbol{\theta}_{old}$ get updated at regular intervals in training.

Many improvements in DRL have focused on changing RL specific aspects of training, such as finding better approaches to sample from the replay buffer \mathcal{R} , by weighing experiences either according to recency or estimated importance [246] [233, 247]. Significant progress has also been achieved by changes in loss functions and parametrisations, particularly through a distributional approach to value based learning [117, 253], which updates the Bellman operators to act on a distributional space rather than in expectation as shown in Eq (7.3). The resulting DQN style agent obtained from the distributional perspective is called CategoricalDQN. Rainbow [251] is an agent collating many of the above improvements, by combining categorical losses, improved prioritised replay strategies [246], exploration bonuses [247], and other RL specific advances such as double Q -learning [254] and disambiguation of state-action value estimates [255]. We will show how similar performance improvements can be obtained by focusing on the neural training aspects via changes to the critic modelling the state-action value, instead of taking an RL centric view.

7.2.2 Challenges with optimisation in DRL

DRL presents its own unique set of optimisation challenges within the deep learning family. One challenge stems from non-stationarity and the lack of a unique objective throughout training: due to the expectation under the policy in loss functions, the loss function itself changes as training progresses; in this aspect, RL is similar to GANs, as the GAN loss also changes as the distribution induced by the generator changes. In DRL the i.i.d. assumption required to obtain an unbiased estimate of the gradient (see Section 2.1.1) is also often violated, as data instances might come from the same episode; furthermore it can be shown that in RL the parameter updates might not form a gradient vector field [243, 244]. Both of these aforementioned issues are mitigated, but not fully resolved, through the use of replay buffers, which

accumulate transition pairs required to perform the Q -learning update (Eq 7.7) across a long horizon of agent experience. The loss formulation in RL can present an additional challenge: consistency losses where the target value of the function approximator depends on *itself*, as can be seen in Eq (7.7), are harder to optimise than standard regression losses where the regression target is provided by the dataset, as is the case in supervised learning.

Despite these specificities of DRL, the optimisation algorithms used are those in standard deep learning and described in Section 2.2, which assume i.i.d. data and the existence of a unique function to optimise. Many a practical implication result as a consequence, including sensitivity to hyperparameterameters [256], the lack of performance of regularisation methods that have been shown to improve optimisation in supervised learning [23], the heavy reliance on adaptive optimisers such as Adam and RMSprop compared to supervised learning where similar results can be obtained with gradient descent with momentum [109, 257].

Inspired by the certain similarities between DRL and GANs, such as non-stationarity, together with the observed common challenges such as sensitivity to hyperparameters, we investigate the use of Spectral Normalisation to alleviate some of the optimisation challenges in DRL. In GANs, Spectral Normalisation has been shown to improve sensitivity to hyperparameters such as learning rate ranges [30], and has been integrated in many state of the art algorithms; for a wider overview of the effects of Spectral Normalisation outside RL, we refer the reader to Chapter 6.

7.2.3 Reinforcement learning environments

We evaluate agents on the Atari ALE benchmark [258], containing 51 Atari games. The Atari benchmark is often used in DRL, and thus provides useful baselines to serve as a comparison against our methods. Here we use as baselines C51—the CategoricalDQN agent introduced by Bellemare et al. [117] trained on Atari ALE—and the Rainbow agent [109]—which improves C51 with additional enhancements, providing a strong baseline on this environment. Unless otherwise specified, experiments on Atari use the standard DQN critic architecture and the experimental setting provided by the C51 baseline, shown in the Appendix in Table E.5 and Section E.2.5,

respectively.

When performing ablation studies using a large setting such as ALE is computationally prohibitive, and we thus use the MinAtar environment [259] instead. MinAtar reproduces five games from ALE without the visual complexity and has been found to be a good testing ground for reproducing ablations made using ALE [260]. Unless otherwise specified, MinAtar experiments use the critic architecture described in Table E.3 in the Appendix.

7.3 Spectral Normalisation in RL

Our aim is to investigate whether the application of Spectral Normalisation in DRL can improve performance, as has been observed in other deep learning domains. To do so, we use Spectral Normalisation applied to neural critics in Q learning algorithms; for a description of Spectral Normalisation, see Section 6.3. The Spectral Normalisation formulation we use is largely based on the one proposed in the original paper Miyato et al. [30], with a few alterations: we do not backpropagate through the power iteration operation; we only apply the normalisation if the spectral norm of a layer is larger than 1; for convolutional layers we do not reshape the 4 dimensional filter tensors to a two-dimensional matrix and compute the spectral norm of the newly obtained linear operator but estimate the of the original convolutional operator instead [198]¹. Code available to reproduce our results is available at <https://github.com/floringogianu/snrl>.

7.3.1 Applying Spectral Normalisation to a subset of layers

When we started experimenting with applying Spectral Normalisation to the DQN critic, we applied Spectral Normalisation to all layers, as is standard in GANs [30]. We quickly noticed, however, that imposing the hard smoothness constraint on the state-action value function to be 1-Lipschitz can substantially decrease performance. As we have observed in Section 6.5.1, enforcing a 1-Lipschitz constraint can reduce a model’s

¹This only applies to the RL experiments in this section, and not to any other Spectral Normalisation experiments in the thesis. We also note that the majority of the results in this section apply normalisation to linear layers, yielding them comparable with the other presented results.

capacity to capture the decision surface and decrease performance (Figure 6.3d); relaxing the constraint by forcing the model to be 10-Lipschitz can mitigate this loss of performance (Figure 6.3e). Similar observations have been made by Gouk et al. [261], who propose introducing a hyperparameter that determines the strength of the Lipschitz constraint. We take a different approach here, where instead of changing the Lipschitz constant as a hyperparameter, we choose to only normalise certain layers of the network, most often just one. In particular we observe that normalising the last layer of the network is detrimental. This suggests that imposing a smoothness constraint on the Q function is prohibiting the model from approximating the required decision surface. In the case of Atari games we study here this is not entirely unexpected, since similar input states can lead to a drastically different game outcome. The expected reward of a particular action a can be determined by a few pixels, such as the location of the ball and paddle in *Pong*, and that of the ghost in *MS-PacMan*. Thus, even for states s, s' such that $\|s - s'\|$ is small, $\|Q(s, a; \theta) - Q(s', a; \theta)\|$ can be large, making the optimal Q function non-smooth.

Instead of normalising the entire network, we instead normalise only a few layers. Experimentally we found that Spectral Normalisation applied to a DRL critic performs best on the layer with the largest number of parameters; for the standard DQN architecture this is the penultimate layer [107, 109]. We thus denote the index of the Spectral Normalised layers from the bottom of the network, i.e. $\text{SN}[-1]$ will refer to a critic where Spectral Normalisation is applied to the output layer and $\text{SN}[-2]$ to the penultimate layer, while $\text{SN}[-2, -3]$ refers to a network where Spectral Normalisation is applied both the penultimate layer and the one preceding it.

7.3.2 Experimental results on the Atari benchmark

We assess the effect of applying Spectral Normalisation to the critic both for the standard DQN as well as its distributional counterpart (C51) when trained on the ALE benchmark containing 54 Atari games and show results in Figure 7.1. We consistently observe that $\text{SN}[-2]$ improves over the baseline significantly. Furthermore, $\text{SN}[-2]$ provides a significant improvement over Rainbow [251], a competitive RL agent based on C51 combining multiple competitive RL specific techniques. We further show

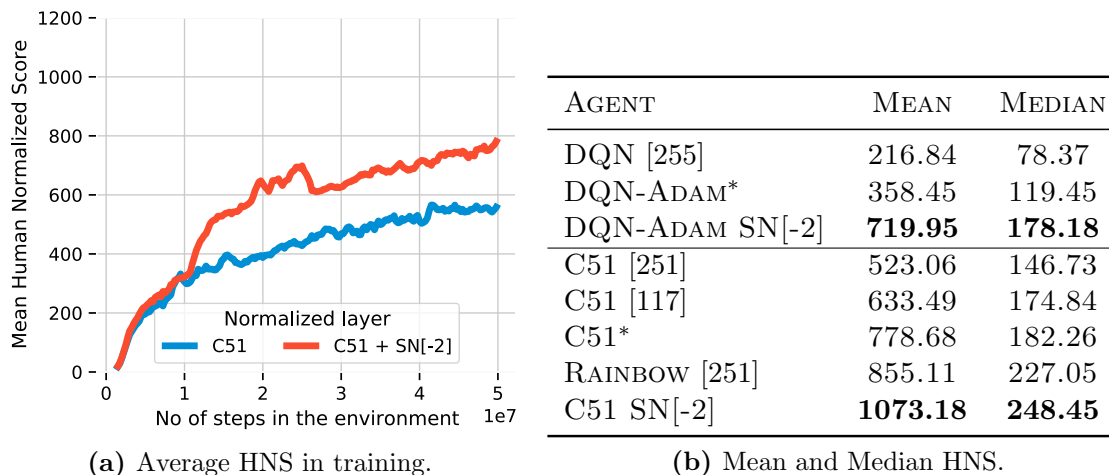


Figure 7.1: Selectively adding Spectral Normalisation to RL agents improves their performance. (a): *Average* Human Normalised Score (HNS) per time-step for C51 with Spectral Normalisation. (b): Mean and median Human Normalised Score on 54 Atari games with random starts evaluation. We observe that applying Spectral Normalisation to the penultimate layer of the critic leads to a significant performance boost both for DQN and C51, and rivals algorithm improvements obtained by RAINBOW. References indicate the sources for the scores for each algorithm. We mark our own implementations of the baseline with *. Agents are evaluated with the protocol in [251]; all results use an average over 54 ALE games.

in Figure 7.2a that SN[-2] leads to better results when trained longer, as opposed to the un-normalised DQN which plateaus earlier in training, and towards the end of training has a decrease in performance. Figures 7.2a and 7.2b show that, while normalising the layer with the largest number of parameters—SN[-2]—performs best, normalising other layers also significantly improves performance over the baseline DQN agent. For all experiments in Figure 7.1, when applying Spectral Normalisation we use the *same hyperparameters as the baseline*, without retuning hyperparameters. These results suggest that Spectral Normalisation can be added to existing performant DRL training configurations without requiring a new hyperparameter search.

In order to assess the performance of Spectral Normalisation beyond the standard DQN network, we perform a sweep over architectures and consistently observe that Spectral Normalisation improves performance when applied to a range of architectures, as shown in Figure 7.3; here we use 12 different model sizes of three different model depths and four different layer widths.

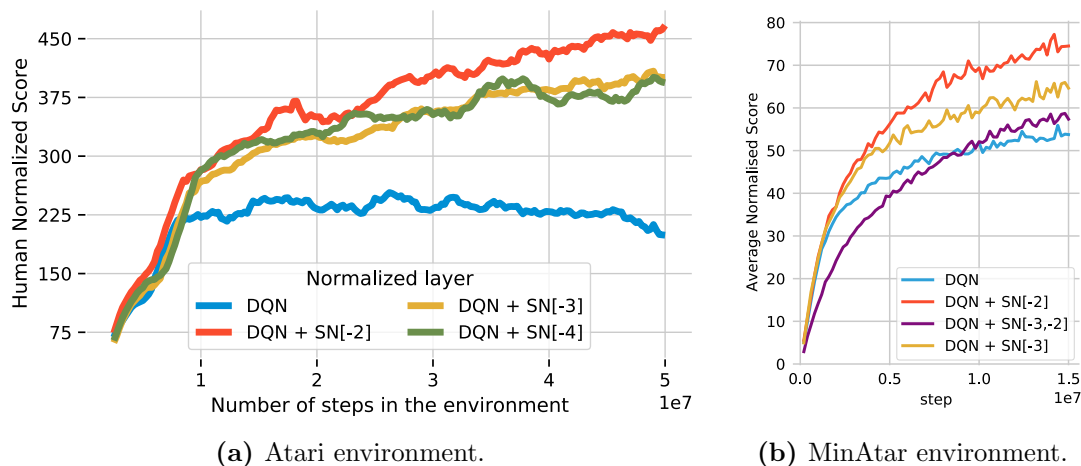


Figure 7.2: Selectively applying Spectral Normalisation to one layer leads to significant improvements over a DQN trained with Adam. While the baseline plateaus early in training, when applying Spectral Normalisation the agents continue to improve throughout training. (a): Human Normalised Score for a DQN-Adam baseline with Spectral Normalisation applied on three different layers. Average over 54 Atari games. (b): On a 15M steps training on MinAtar; each line is an average of 10 seeds over four MinAtar games and four different model sizes.

7.3.3 An optimisation effect

While selectively applying Spectral Normalisation to a subset of layers of a Q -learning critic is a simple intervention, we have seen in the previous section that it has a strong effect on performance. We now aim to understand *why* Spectral Normalisation leads to these improvements. We first start with the smoothness hypothesis, considering whether the performance improvement is due to an increased smoothness of the critic. We note, however, that unlike other applications of Spectral Normalisation which normalise the entire network, selectively applying Spectral Normalisation to a subset of layers as we do here does not lead to any guarantee regarding the smoothness of the entire critic. Indeed, since there are no constraints on the other layers, the smoothness of those layers could increase throughout training to compensate for the constraint on the normalised layers. When aggregating over the games in MinAtar we observe 75% of applications of $\text{SN}[-2]$ decrease the smoothness of the model—as approximated by the maximum norm of the critic’s Jacobian with respect to inputs over a large set of state-action pairs, for a connection between input Jacobians

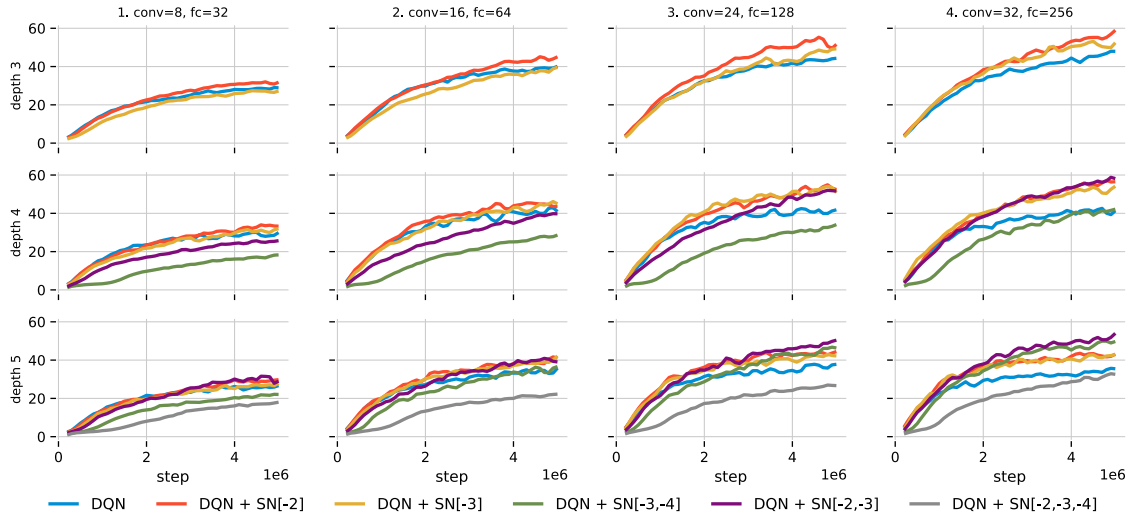


Figure 7.3: Spectral Normalisation shows gains for all model sizes. Looking at the baseline ($-$ DQN), we observe two performance regimes on MinAtar: for shallow, depth 3 models, performance increases with the width of the model; for deeper models performance generally stagnates with increasing depth and width. In both regimes applying SN on individual ($-$, $-$) or multiple ($-$) layers improves upon the baseline suggesting a regularisation effect we could not reproduce with other regularisation methods. Notice that the strong regularisation resulted from applying Spectral Normalisation to input layers ($-$) or too many layers ($-$) can however degrade performance. Each line is an average over normalised scores of 4 games \times 10 seeds.

and Lipschitz smoothness, see Section 6.2—when controlling over hyperparameters; for certain games it is as high as 90%, while for others as low as 58%. We also observe that applying Spectral Normalisation to more layers do not always increase smoothness in Table 7.1. Furthermore, we notice that while applying Spectral Normalisation increased performance for all games in MinAtar, there is not a strong effect between smoothness and performance as measured by the Spearman-rank correlation between the smoothness of the learned model and the best performance obtained by the agent in all games, as we show in Table 7.2. To further assess the effect of the critic smoothness on performance, we used gradient penalties as a soft constraint to encourage smoothness (see Section 6.3 for more details on gradient penalties), which did not result in any significant benefit compared to the baseline, as we show in Figure E.2 in the Appendix.

The above evidence suggests that the increased performance from the application

		Less or no normalisation		
		No normalisation	SN[-2]	SN[-3]
More SN	SN[-2]	74.6%	X	X
	SN[-3]	69.2%	X	X
	SN[-2, -3]	75%	57%	60%

Table 7.1: The percentage of training settings (architectures and seeds) for which applying Spectral Normalisation to more layers leads to smoother networks, as measured by the model’s Jacobian at peak performance. Selectively applying Spectral Normalisation to only a subset of layers does not consistently lead to an increase in smoothness, and neither does applying Spectral Normalisation to two layers compared to only one layer. Results are obtained from a sweep over architectures detailed in Table E.6, with 10 seeds each.

GAME	SPEARMAN RANK
Asterix	0.129
Breakout	0.199
Seaquest	0.151
Space Invaders	0.453

Table 7.2: Spearman Correlation coefficient between the negative norm of the Jacobian (our measure of decreased smoothness) and peak performance for each game. Higher means more correlation. We don’t consistently observe a strong correlation between performance and increased smoothness across games.

of Spectral Normalisation to certain critic layers in DRL might not be fully explainable due to an increased smoothness effect. To better understand why Spectral Normalisation helps DRL, we turn our attention to the possibility of an increase in performance due to an improvement in optimisation; we have already seen in Section 6.5.2 that Spectral Normalisation can interact with learning rate and momentum decay rates in supervised learning and in GANs. In RL, we show in Figure 7.2 that using Spectral Normalisation allows for longer training without early plateaus, and that applying Spectral Normalisation to existing choice of hyperparameters leads to an increase of performance without requiring an additional hyperparameter sweep, as shown in the results in Figure 7.1. To investigate whether Spectral Normalisation has an effect on hyperparameter sensitivity of DQN agents, we perform a large sweep over Adam hyperparameters learning rate h and ϵ (for an overview of Adam see

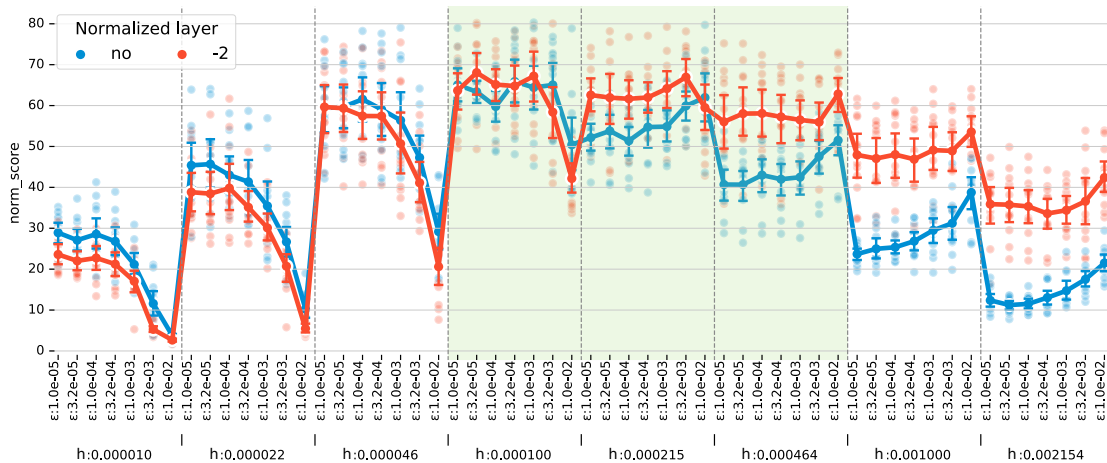


Figure 7.4: SN[-2] applied to DQN, Spectral Normalisation applied to the largest layer in the DQN critic, can significantly decrease sensitivity to the learning rate and ϵ hyperparameter of Adam. The cross product of learning rates in $h \in \{0.00001, \dots, 0.00215\}$ and in $\epsilon \in \{0.00001, \dots, 0.01\}$ is explored.

Section 2.2) and show that indeed, applying Spectral Normalisation increases the range of performant hyperparameters significantly; results are shown in Figure 7.4.

In order to further investigate the effect of Spectral Normalisation in DRL training, we develop a set of methods which independently assess its effects. To this end, consider the neural network to which Spectral Normalisation is applied, which we assume has N layers, with weights $\mathbf{W}_i \in \mathbb{R}^{d_{i+1} \times d_i}$, $\mathbf{b}_i \in \mathbb{R}^{d_{i+1}}$:

$$\mathbf{o} = \mathbf{W}_N \mathbf{h}_{N-1} + \mathbf{b}_N \quad (7.8)$$

$$\mathbf{h}_{N-1} = a(\mathbf{z}_{N-1}) = a(\mathbf{W}_{N-1} \mathbf{h}_{N-2} + \mathbf{b}_{N-1}) \quad (7.9)$$

$$\dots \quad (7.10)$$

$$\mathbf{h}_1 = a(\mathbf{z}_1) = a(\mathbf{W}_1 \mathbf{h}_0 + \mathbf{b}_1), \quad (7.11)$$

with $\mathbf{h}_0 = \mathbf{x}$ and a is an activation function applied element-wise to the input vector. The activation a has to be a rectifier, such as ReLU [262] or Leaky ReLU [263], since we rely on $a(\sigma_l \mathbf{x}) = \sigma_l a(\mathbf{x})$, with $\sigma_l > 0$. We can write our optimisation problem as:

$$\min_{\theta} E(\theta) = E(\mathbf{o}; \theta = \{\mathbf{W}_{1 \leq i \leq N}, \mathbf{b}_{1 \leq i \leq N}\}). \quad (7.12)$$

We assess the effect of applying Spectral Normalisation to layer l , but note the same

derivation applies when multiple layers are normalised. We note that for simplicity we assumed above that the model is an MLP, though the same computation carries for convolutional layers. The analysis that follows relies on the following observation:

Remark 7.3.1 *In a neural network with rectifier activation functions, if the bias term of a layer where Spectral Normalisation is applied gets adjusted to account for the normalisation, Spectral Normalisation can be seen as scaling the output of the layer by the inverse of the spectral norm σ_l .*

To see why, consider $\mathbf{h}_{l+1} = a(\mathbf{W}_l \mathbf{h}_{l+1} + \mathbf{b}_l)$ and its normalised counterpart: $\mathbf{h}'_{l+1} = a(\mathbf{W}_l/\sigma_l \mathbf{h}_{l+1} + \mathbf{b}_l)$. Since we chose the activation function such that $a(\mathbf{h}/\sigma_l) = a(\mathbf{h})/\sigma_l$, we observe that if we also scale the bias \mathbf{b}_l by $1/\sigma_l$, we can obtain $\mathbf{h}'_{l+1} = \mathbf{h}_{l+1}/\sigma_l$. Furthermore, if we apply the same bias scaling to subsequent layers $\mathbf{b}_i \rightarrow \mathbf{b}_i/\sigma_l, \forall i \geq l$, we obtain that $\mathbf{h}'_i = \mathbf{h}_i/\sigma_l$, including $\mathbf{o}' = \mathbf{o}/\sigma_l$. Thus if we assume that the bias vectors can be scaled by a constant in training, the effect of Spectral Normalisation is to scale each activation following the normalisation by the inverse of the spectral norm of the normalised layer; if multiple layers are normalised, activations are scaled by inverse of the product of the spectral norm of the normalised layers. We note that the bias scaling assumption does not change how expressive the model can be, as the bias term could learn such as scaling through training. Indeed, we test this assumption empirically and use a procedure we term DIVOUT, in which instead of applying Spectral Normalisation, we divide the model output by the product of the spectral norms of the normalised layers ρ :

$$\text{DIVOUT}[\rho] : \quad \min_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) = E\left(\frac{\mathbf{o}}{\prod_{s \in \rho} \sigma_s}; \boldsymbol{\theta}\right). \quad (7.13)$$

We investigate where DIVOUT can perform similarly to its Spectral Normalisation counterparts and show in Figure 7.5 that in offline reinforcement learning DIVOUT $[\rho]$ and SN $[\rho]$ have similar behaviours.

We now turn our attention towards the effects that Spectral Normalisation has on the model gradients. To do so, we would like to compute $\frac{dE}{d\mathbf{W}_i}$ and $\frac{dE}{d\mathbf{b}_i}$ for normalised and un-normalised networks. As before, for simplicity we consider the

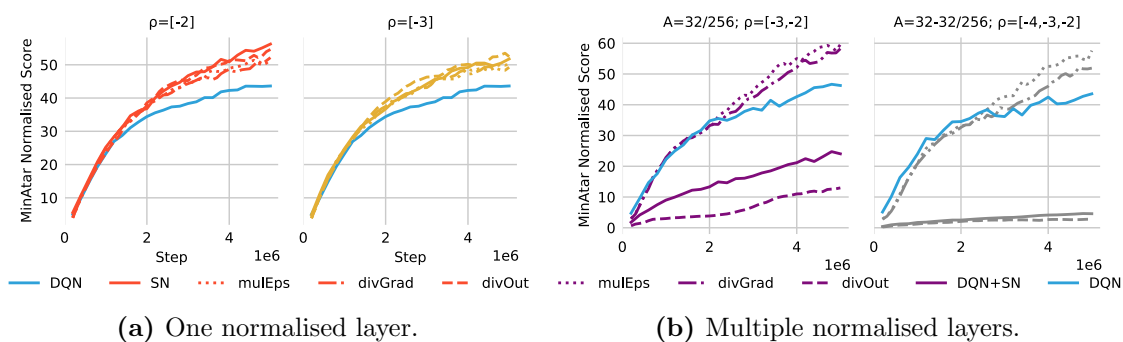


Figure 7.5: The effect of DIVOUT, DIVGRAD and MULEPS on RL performance on MinAtar. *Despite not changing the model, but only changing the optimisation procedure based on the effect Spectral Normalisation has on parameter updates, DIVGRAD and MULEPS perform comparable to, or better than, Spectral Normalisation.* (a) shows results obtained from simulating optimisation effects of Spectral Normalisation when applied to one layer, where DIVGRAD and MULEPS and DIVOUT obtain similar performance performance to Spectral Normalisation, and significantly improve over the baseline. (b) shows results obtained when looking at the optimisation effects of applying Spectral Normalisation to multiple layers; here Spectral Normalisation hurts performance, and so does DIVOUT, but DIVGRAD and MULEPS still significantly improve performance compared to the baseline. This suggests that large changes model output obtained by Spectral Normalisation and DIVOUT can be detrimental, while optimisation only changes are beneficial in that setting.

situation of one normalised layer l and we assume bias scaling. We then have

$$\mathbf{o}' = \mathbf{o}/\sigma_l; \quad \mathbf{h}'_i = \mathbf{h}_i/\sigma_l, \quad \forall i \geq l; \quad \mathbf{h}'_i = \mathbf{h}_i, \quad \forall i < l. \quad (7.14)$$

While the computational graph only changes at layer l where the $\mathbf{h}_l = a(\mathbf{W}_l \mathbf{h}_{l-1} + \mathbf{b}_l)$ is now replaced with

$$\mathbf{h}'_l = a(\mathbf{z}'_l) = a(\mathbf{W}_l/\sigma_l \mathbf{h}_{l-1} + \mathbf{b}_l/\sigma_l) = a(\mathbf{W}_l \mathbf{h}_{l-1} + \mathbf{b}_l)/\sigma_l, \quad (7.15)$$

the evaluation of the higher level derivatives will also change due to the scaling effect in Eq (7.14). To make the changes in the computational graph due to the normalisation explicit, we will call the loss with the normalisation E' . We are interested in computing $\frac{dE'}{d\mathbf{W}_i}$ and $\frac{dE'}{d\mathbf{b}_i}$, and comparing it with the backpropagation in

an un-normalised model

$$\frac{dE}{d\mathbf{W}_N} = \frac{dE}{d\mathbf{o}} \mathbf{h}_{N-1}^T \quad (7.16)$$

$$\frac{dE}{d\mathbf{h}_{N-1}} = \frac{d\mathbf{o}}{d\mathbf{h}_{N-1}}^T \frac{dE}{d\mathbf{o}} = \mathbf{W}_N^T \frac{dE}{d\mathbf{o}} \quad (7.17)$$

$$\frac{dE}{d\mathbf{W}_i} = \frac{dE}{d\mathbf{z}_i} \mathbf{h}_{i-1}^T = \left(\frac{dE}{d\mathbf{h}_i} \odot a'(\mathbf{z}_i) \right) \mathbf{h}_{i-1}^T \quad i < N \quad (7.18)$$

$$\frac{dE}{d\mathbf{h}_{i-1}} = \frac{d\mathbf{z}_i}{d\mathbf{h}_{i-1}}^T \frac{dE}{d\mathbf{z}_i} = \mathbf{W}_i^T \left(\frac{dE}{d\mathbf{h}_i} \odot a'(\mathbf{z}_i) \right) \quad i < N. \quad (7.19)$$

When normalisation is applied, the computational graph does not change for the layers following the normalisation compared to the un-normalised case, but the evaluation of the derivatives changes due to the scaling in Eq (7.14)

$$\frac{dE'}{d\mathbf{o}'}(\mathbf{o}') = \frac{dE}{d\mathbf{o}}(\mathbf{o}/\sigma_l) \quad (7.20)$$

$$\frac{dE'}{d\mathbf{W}_N} = \frac{dE'}{d\mathbf{o}'} \mathbf{h}'_{N-1}{}^T = \frac{dE'}{d\mathbf{o}'} \left(\underbrace{\mathbf{h}_{N-1}/\sigma_l}_{(7.14)} \right)^T = \frac{1}{\sigma_l} \frac{dE'}{d\mathbf{o}'} \mathbf{h}_{N-1}^T \quad (7.21)$$

$$\frac{dE'}{d\mathbf{h}'_{N-1}} = \mathbf{W}_N^T \frac{dE}{d\mathbf{o}}(\mathbf{o}/\sigma_l) \quad (7.22)$$

$$\frac{dE'}{d\mathbf{W}_i} = \left(\frac{dE'}{d\mathbf{h}'_i} \odot a'(\mathbf{z}'_i) \right) \mathbf{h}'_{i-1}{}^T = \frac{1}{\sigma_l} \left(\frac{dE'}{d\mathbf{h}'_i} \odot a'(\mathbf{z}'_i) \right) \mathbf{h}_{i-1}^T, \quad i > l. \quad (7.23)$$

For the normalised layer l , we have that due to Eq (7.15)

$$\frac{dE'}{d\mathbf{W}_l} = \frac{dE'}{d\mathbf{z}'_l} \underbrace{\mathbf{h}'_{l-1}{}^T/\sigma_l}_{\text{normalisation,(7.15)}} = \frac{1}{\sigma_l} \left(\frac{dE'}{d\mathbf{h}'_l} \odot a'(\mathbf{z}'_l) \right) \mathbf{h}'_{l-1}{}^T \quad (7.24)$$

$$\frac{dE'}{d\mathbf{h}'_{l-1}} = \underbrace{\frac{1}{\sigma_l} \mathbf{W}_l^T}_{\text{normalisation,(7.15)}} \left(\frac{dE'}{d\mathbf{h}'_l} \odot a'(\mathbf{z}'_l) \right). \quad (7.25)$$

For the layers before the normalisation, we have

$$\frac{dE'}{d\mathbf{W}_i} = \left(\frac{dE'}{d\mathbf{h}'_i} \odot a'(\mathbf{z}'_i) \right) \mathbf{h}'_{i-1}{}^T \quad i < l \quad (7.26)$$

$$\frac{dE'}{d\mathbf{h}'_{i-1}} = \mathbf{W}_i^T \left(\frac{dE'}{d\mathbf{h}'_i} \odot a'(\mathbf{z}'_i) \right) \quad i < l. \quad (7.27)$$

Thus, when comparing the gradients obtained with the normalised network with those in an un-normalised network, we observe two changes: first, a change in the output Jacobian (Eq (7.20)), and second, a scaling effect which applies to *all layers in the network*, either through the scaling of the activations for layers following the normalisation (Eq (7.23)), through the changes in the computational graph for the normalised layer (Eq (7.24)), or through the effect of backpropagation of the scaling of $\frac{dE'}{d\mathbf{h}'_{i-1}}$, for the layers preceding the normalisation (Eq (7.26)). A similar argument can be made for biases. While the change in Jacobian is difficult to obtain without affecting the model as we have done with DIVOUT, to capture the gradient scaling effect we propose DIVGRAD, a normalisation scheme based on the effects of applying Spectral Normalisation which *does not change the model, but only changes the optimisation procedure*, by scaling the gradients of the parameters:

$$\text{DIVGRAD}[\rho] : \frac{dE}{d\boldsymbol{\theta}_i} \rightarrow \frac{1}{\prod_{i \in \rho} \sigma_i} \frac{dE}{d\boldsymbol{\theta}_i}. \quad (7.28)$$

Results in Figure 7.5 show that DIVGRAD performs well and that it can outperform Spectral Normalisation or DIVOUT on MinAtar; further results supporting this hypothesis are shown in Figure E.5 in the Appendix.

The above observation regarding gradient scaling also allows us to intuit the effect of Spectral Normalisation in conjunction with optimisers like Adam (for a detailed description of the Adam optimiser, see Section 2.2). If we consider that the spectral norm of a weight matrices changes slowly in training (which we verify empirically in Figure E.4), *under the bias scaling assumption*, the effect of Spectral Normalisation on the Adam update can be seen as

$$\mathbf{m}'_t = \beta_1 \mathbf{m}'_{t-1} + (1 - \beta_1) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) / \sigma \approx \mathbf{m}'_t / \sigma \quad (7.29)$$

$$\mathbf{v}'_t = \beta_2 \mathbf{v}'_{t-1} + (1 - \beta_2) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})^2 \approx \mathbf{v}'_t / \sigma^2 \quad (7.30)$$

$$\Delta \boldsymbol{\theta}'_t = h \frac{\frac{1}{1-\beta_1^t} \mathbf{m}'_t}{\sqrt{\frac{1}{1-\beta_2^t} \mathbf{v}'_t + \epsilon}} = h \frac{\frac{1}{1-\beta_1^t} \mathbf{m}_t}{\sqrt{\frac{1}{1-\beta_2^t} \mathbf{v}_t + \sigma \epsilon}}, \quad (7.31)$$

where $\sigma = \prod_{i \in \rho} \sigma_i$, the product of the spectral norms of the normalised layers

and we used element-wise operations. Thus, the effect on the Adam update is reflected in a scaling of the hyperparameter ϵ by σ . To investigate this effect on the Adam optimiser, we introduce MULEPS, which changes the Adam update by multiplying ϵ by the product of the spectral norms of the normalised layers. We show in Figure 7.5 that MULEPS performs similarly to selectively applying Spectral Normalisation in DRL, and can improve upon it. A strong effect is particularly observed when MULEPS and Spectral Normalisation are applied to multiple layers (Figure 7.5b), where Spectral Normalisation *decreases* performance compared to the baseline, while MULEPS increases it. Thus, MULEPS, provides another example of how we can use Spectral Normalisation to derive optimisation only changes that lead to improvements in performance in DRL.

7.4 Revisiting Spectral Normalisation in GANs

By analysing the effect of Spectral Normalisation in the previous section, we have uncovered multiple effects: model effects such as output scaling (Eq (7.20)), and optimisation effects such as gradient scaling (Eq (7.28)) and interactions with optimisation hyperparameters such as ϵ in Adam (Eq (7.31)). We saw that in DRL the effects of Spectral Normalisation can be recovered through the use of optimisation schemes, such as DIVGRAD and MULEPS, derived based on the optimisation only effects of Spectral Normalisation. Since Spectral Normalisation was first introduced for GANs, and indeed is a key component of many successful GANs [33, 39], it is only natural to ask whether the same optimisation methods perform equally well when training GANs. We studied GANs previously and described them in Section 4.1, and the interactions between smoothness normalisation and optimisation in GANs have been discussed in Chapter 6. The application of Spectral Normalisation is different in the GAN setting compared to RL, in that in GANs it is applied to *every layer* of the model². This difference between the DRL and GANs is significant, since applying Spectral Normalisation on every layer ensures the model is Lipschitz

²When Spectral Normalisation was introduced, it was only applied to the discriminator, but since it has also been applied to the generator; here, for consistency with the Spectral Normalisation baseline in Miyato et al. [30], we apply it only to the discriminator.

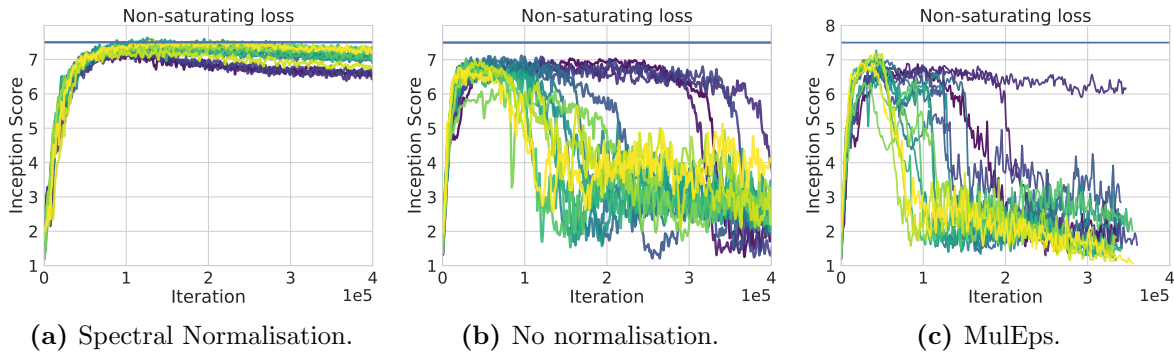


Figure 7.6: Optimisation changes based on Spectral Normalisation do not recover the performance of Spectral Normalisation in GANs. Different units indicate different learning rates for the discriminator and generator.

(at least when Spectral Normalisation is applied on the linear operator given by convolutional layers and not on the reshaped filters), while applying it to only one layer leads to no such assurances, as we have seen in the previous section. This difference in application of Spectral Normalisation suggests perhaps that the change in the model is important in GANs. To test this intuition, we experimented with optimisation only methods DIVGRAD and MULEPS in GAN training. We note that as the generator’s backpropagation requires a backward pass through the discriminator, any normalisation applied to the discriminator in MULEPS and DIVGRAD will also affect the generator, as was the case with un-normalised early layers in the RL critic (see Eqs (7.31) and (7.28)). Results in Figure 7.6 confirm that indeed, methods that only make changes to optimisation dynamics and do not account for model changes induced by Spectral Normalisation do not recover its performance in GAN training.

7.4.1 Spectral Normalisation and the discriminator output

Since optimisation only changes do not recover the performance of Spectral Normalisation in GAN training, we now investigate whether Spectral Normalisation alters the loss prediction landscape in a substantial way. We start by considering the effect of Spectral Normalisation on the range of discriminator outputs. Since Spectral Normalisation applies a *global constraint* on the entire space, this applies between pairs of data instances, pairs of samples, and pairs of data and samples. If we write $D(\cdot; \phi) = \varsigma(D_L(\cdot; \phi))$, with ς being the sigmoid function, Spectral Nor-

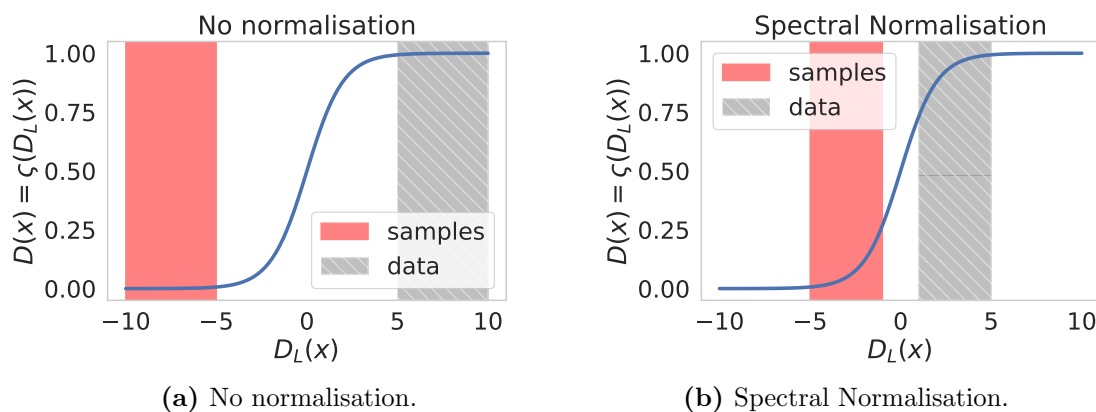


Figure 7.7: Visualisation of the effect of Spectral Normalisation on discriminator predictions: Spectral Normalisation is making the discriminator output less likely to reach the saturating areas of the sigmoid function, denoted here by ς .

malisation ensures that $D_L(\cdot; \phi)$ is 1-Lipschitz by normalising each layer in the discriminator’s network. Given a generator sample $G(\mathbf{z}; \theta)$ and a data instance \mathbf{x} , $\|D_L(G(\mathbf{z}; \theta); \phi) - D_L(\mathbf{x}; \phi)\|_2$ is bound by $\|G(\mathbf{z}; \theta) - \mathbf{x}\|_2$. Thus, while with no normalisation the discriminator has an easier task of separating the data and samples by using the saturating areas of the sigmoid function, this is harder to do when Spectral Normalisation is used; we visualise this effect in Figure 7.7.

We now turn to empirically investigating this effect in GAN training. We measure the output of the discriminator throughout training both on samples and data batches and create a binned-histogram of the outputs obtained from a batch. That is, for each \mathbf{x}_i in a data batch with $D(\mathbf{x}_i; \phi) \in [0, 1]$ we assign it the interval $[0.1p, 0.1(p + 1)]$ with $p \in \{0, \dots, 9\}$ if $D(\mathbf{x}_i; \phi) \in [0.1p, 0.1(p + 1)]$; we do the same for samples $D(G(\mathbf{z}_i; \theta); \phi)$ where \mathbf{z}_i are the latent vectors used to obtain the sample batch. This allows us to measure the variability of the discriminator prediction inside data and sample batches; we measure the entropy of the discriminator output across a batch with and without Spectral Normalisation by measuring the entropy of the distribution induced by the histogram. We show aggregate results throughout training over a hyperparameter sweep in Figure 7.8, and individual results with a fixed set of hyperparameters in Figures 7.9 and 7.10. As we intuited, Spectral Normalisation increases the entropy of the discriminator output throughout training, both on data (Figure 7.8b) and samples (Figure 7.8c). Early in training (Figure 7.9),

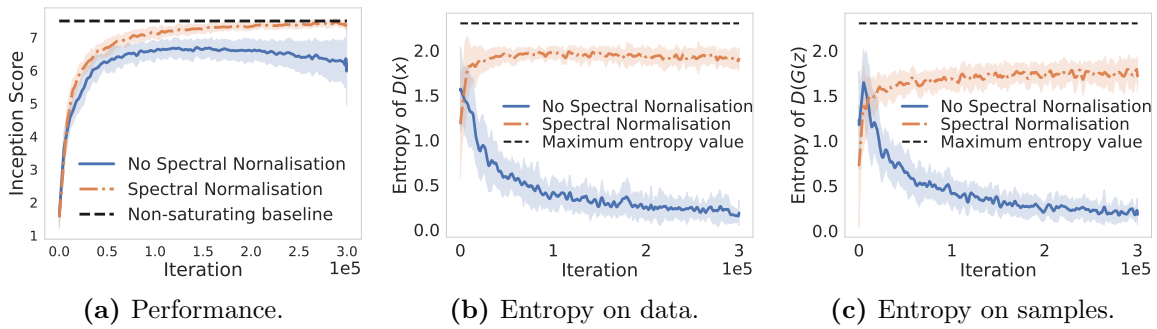


Figure 7.8: Spectral Normalisation increases the entropy of the discriminator’s predictions on data (b) and samples (c). The results are obtained over a sweep of learning rates.

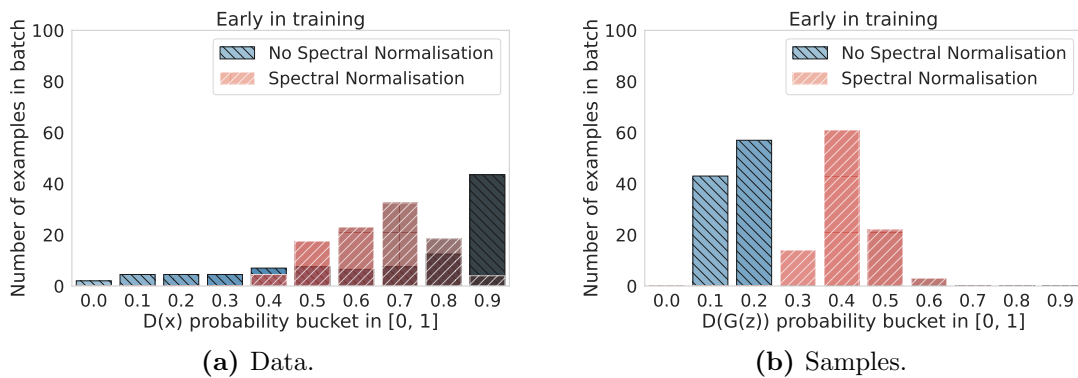


Figure 7.9: Early in training. Spectral Normalisation increases the entropy of the discriminator’s predictions on data (a) and samples (b) in the early stages of training. Without Spectral Normalisation the discriminator is confident that the generator samples are fake, with most samples being assigned less than 0.3 probability of being real.

discriminators trained with Spectral Normalisation are less likely to assign high probability to real data being real, or to samples being fake. Late in training (Figure 7.10), GANs trained with Spectral Normalisation also exhibit higher entropy of discriminator outputs; this end of training behaviour is more consistent to that of a Nash equilibrium (see Section 2.4 for a definition and discussion on Nash equilibria), where the discriminator should not be able to distinguish between samples and data.

Given the effect of Spectral Normalisation on the discriminator’s predictions, it’s natural to ask how this impacts training. We turn to this question next.

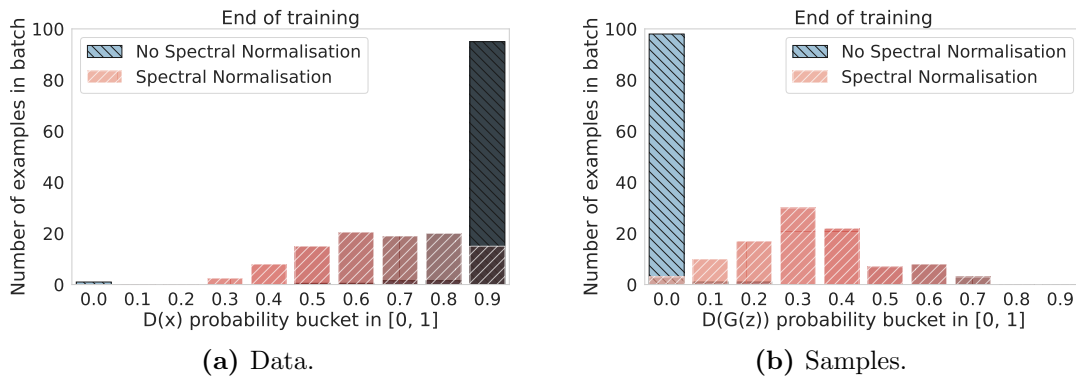


Figure 7.10: Late in training. Spectral Normalisation increases the entropy of the discriminator’s predictions on data (a) and samples (b) in the late stages of training. Without Spectral Normalisation the discriminator is confident that real data is real and generator samples are fake. When Spectral Normalisation is used, the discriminator assigns both data and samples a wide range of probability of being real.

7.4.2 Interactions with loss functions

To examine the effect of the increased entropy on the discriminator’s output, we use two formulations of GAN losses: the non-saturating loss (Eq (4.44)), more commonly used in practice and in the above experiments; and the original zero-sum formulation of GANs (Eq (4.4)), which has been known to be harder to train than other losses. The low performance of this original loss has been attributed to the lack of learning signal the generator early in training, since when the generator performs poorly and the discriminator easily classifies its data as fake, the discriminator does not provide strong gradients for the generator to improve, as we visualise in Figure 7.11; for this reason, the generator loss in this formulation is sometimes called the ‘saturating loss’. Moreover, strong gradients are provided when the generator is doing well, which can lead to unstable behaviour and exiting areas of good performance. Figure 7.11b also highlights why the increased entropy effect in the discriminator’s output induced by Spectral Normalisation—and specifically the decreased probability of assigning outputs close to 0 or 1 observed in Figures 7.9 and 7.10—can be beneficial, as *for both losses, it avoids areas with large or small gradient magnitudes*.

By analysing the interaction between the loss function and the behaviour of the discriminator, we can now make predictions about the empirical behaviours

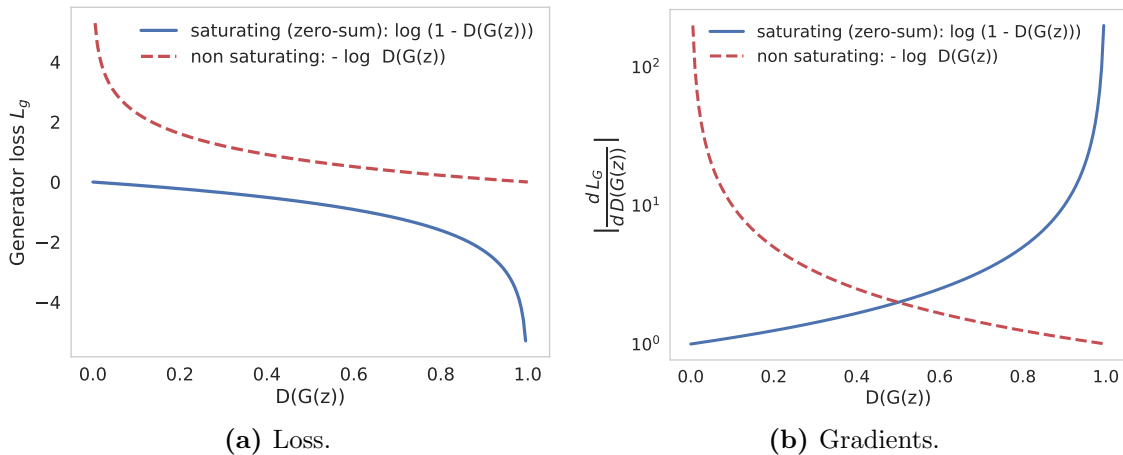


Figure 7.11: Understanding the landscape of GAN losses and gradients. (a): comparing the saturating and non-saturating generator loss. (b): unlike the saturating loss, the non-saturating loss provides stronger gradients when the discriminator can easily classify generated data as fake.

of the two losses. For the saturating loss, without Spectral Normalisation the challenge is that early in training there is no learning signal for the generator; we thus would expect that without Spectral Normalisation the performance is low, and that Spectral Normalisation improves is significantly by avoiding the saturated areas early in training. We observe this in Figure 7.12, which shows that without Spectral Normalisation the saturating loss exhibits extremely poor performance, but that adding Spectral Normalisation significantly decreases the gap between the saturating and non-saturating losses. For the non-saturating loss, we expect that without Spectral Normalisation, the discriminator can be too confident in its prediction that generated data is fake late in the game, as we observed in Figure 7.10b. This leads to large gradients—see Figure 7.11b—which can destabilise training and exit areas of good performance. This expectation is consistent with what we observe empirically in Figure 7.13: adding Spectral Normalisation improves stability later in training.

We are now ready to return to the optimisation methods we derived in Section 7.3.3 based on Spectral Normalisation. Since DIVOUT captures the changes in the model output induced by Spectral Normalisation, we would expect DIVOUT to perform well for GANs trained on CIFAR-10. In particular, by dividing the discriminator output by the product of the spectral norms of the discriminator’s

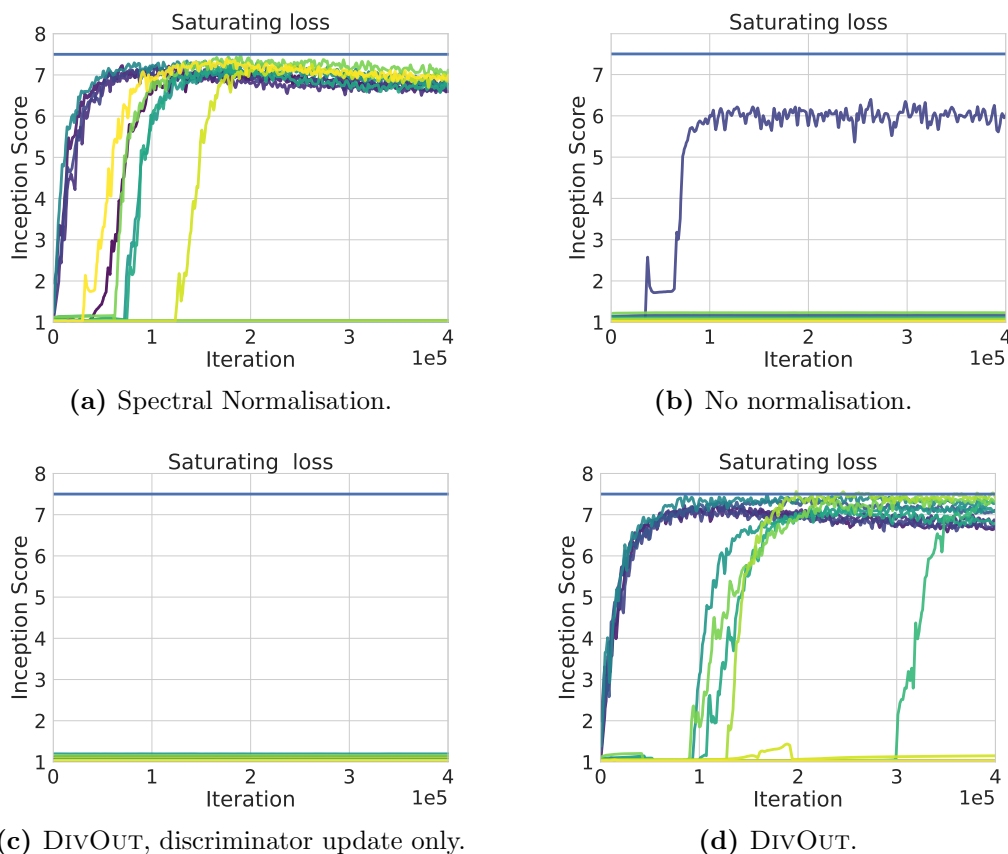


Figure 7.12: Saturating loss. Spectral Normalisation improves performance of GANs trained with the saturating loss. DIVOUT obtains similar performance as Spectral Normalisation, but only when applied to both the discriminator and the generator. Only applying the output scaling from DIVOUT to the discriminator only leads to extremely poor performance, as it does not affect the generator update, and thus does not avoid the areas with poorly conditioned generator gradients. Different units indicate different learning rates for the discriminator and generator. The training optimiser is Adam, used with simultaneous updates.

weights, DIVOUT captures the intuition present in Figure 7.7: it ensures that data and generator samples are less likely to be separated and reach the saturating areas of the sigmoid function. Indeed, as we show in Figures 7.12 and 7.13, DIVOUT obtains a similar performance to Spectral Normalisation both for the saturating and non-saturating loss. Importantly, however, when applying DIVOUT only to the discriminator update, DIVOUT obtains extremely poor performance in the case of the saturating loss as we show in Figure 7.12c, but still performs on par with Spectral Normalisation when the non-saturating loss is used, as seen in Figure 7.13c. This

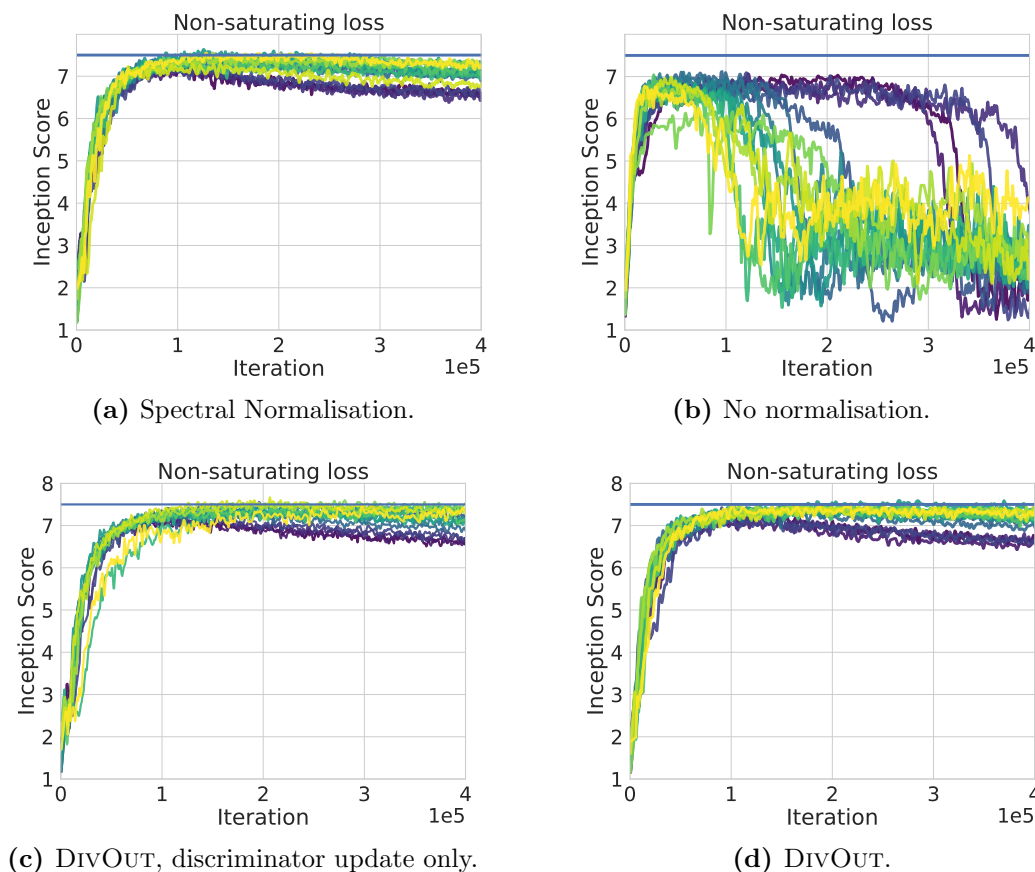


Figure 7.13: Non-saturating loss. DIVOUT obtains similar performance to Spectral Normalisation for GANs on CIFAR-10 when the non-saturating loss is used. Different units indicate different learning rates for the discriminator and generator. Only applying the output scaling from DIVOUT to the discriminator leads to similar performance as when applying it to both the discriminator and the generator, since the non-saturating loss makes the generator less susceptible to the scale of the discriminator output, unlike for the saturating loss as we saw in Figure 7.12. The training optimiser is Adam, used with simultaneous updates.

further confirms that in the case of the saturating loss, one of the effects of Spectral Normalisation is to change the conditioning of the generator by avoiding the areas with low magnitude gradients early in training. We are also ready to revisit the MULEPS results with the non-saturating loss presented in Figure 7.6; since MULEPS does not capture the changes in discriminator output, it will not be able to avoid the decrease in performance observed when no Spectral Normalisation is used.

We thus see that in the GAN case, there is a strong effect in model training brought by Spectral Normalisation, which does not come from an only optimisation

change. We now turn our attention to the following question: even if in GANs Spectral Normalisation does not lead to a primarily optimisation effect, are GANs potentially benefiting from the optimisation changes brought by Spectral Normalisation?

7.4.3 Hypothesis: Adam in the low-curvature regime

We have seen in previous sections that Spectral Normalisation can be seen as acting to multiply the hyperparameter ϵ of Adam (Eq (7.31)). Since the spectral norms of weight matrices grow in training, this can be seen as an anneal of ϵ from small to large. In the previous section we have seen that in the GAN case the scheduler derived from this observation, MULEPS, does not recover the performance of Spectral Normalisation in GANs. Now we ask whether the ϵ annealing behaviour of Spectral Normalisation can still have a positive effect on GAN training, despite not being the only effect.

We start with a few observations. First, when a small ϵ is used, for dimensions of low-curvature Adam leads to a constant update given by the learning rate. Second, GANs trained with Adam exit the area of high performance, even when Spectral Normalisation and the non-saturating losses are used, as we show in Figure 7.6. Third, this behaviour is not observed when gradient descent or Runge–Kutta4 are used (see Figure 4.8 and Qin et al. [67]), so this behaviour seems specific to the use of the Adam optimiser. These observations lead to us to postulate the following hypothesis: *Is Adam exiting the areas of high performance in GAN training because many parameters are in a low-curvature regime, and taking a step of size given by the learning rate is then leading to an exit from the high performance area?*

To first understand the effect of ϵ in Adam in the low-curvature regime, we use the arguments from Section 2.2. If we are in a low-curvature regime for parameter $\theta_{[l]} \in \mathbb{R}$, then by definition its gradient $\nabla_{\theta_{[l]}} E$ does not change between iterations. We use this to expand \mathbf{m}_t in the Adam update, defined in Eq (2.23), at iteration t :

$$(\mathbf{m}_t)_{[l]} = (1 - \beta_1) \sum_{i=1}^t \beta_1^{i-1} \nabla_{\theta_{[l]}} E(\theta_{t-i}) \approx (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \nabla_{\theta_{[l]}} E \quad (7.32)$$

$$= (1 - \beta_1^t) \nabla_{\theta_{[l]}} E. \quad (7.33)$$

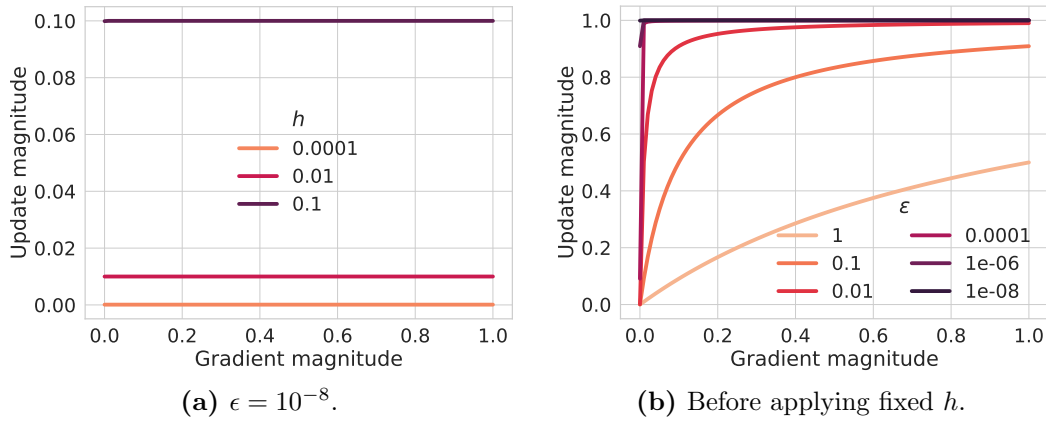


Figure 7.14: The effect of the learning rate (a) and ϵ (b) on the magnitude of the update in the low-curvature regime in Adam: while the learning rate scales the gradient by a constant, the effect of ϵ strongly depends on the magnitude of the gradient, particularly for large ϵ .

With the same arguments we have

$$(\mathbf{v}_t)_{[l]} \approx (1 - \beta_2^t)(\nabla_{\theta_{[l]}} E)^2. \quad (7.34)$$

We thus have that the Adam update for parameter $\theta_{[l]}$ is

$$\Delta\theta_{[l]} = h \frac{\frac{1}{1-\beta_1^t}(\mathbf{m}_t)_{[l]}}{\sqrt{\frac{1}{1-\beta_2^t}(\mathbf{v}_t)_{[l]} + \epsilon}} \approx h \frac{\nabla_{\theta_{[l]}} E}{\sqrt{(\nabla_{\theta_{[l]}} E)^2 + \epsilon}}. \quad (7.35)$$

This result shows the importance of ϵ in Adam in the low-curvature regime: if ϵ is small then the update is given by the learning rate h ; this makes the update independent of the gradient magnitude. We plot this in Figure 7.14b, and show how increasing ϵ can reinstate the dependence on gradient magnitudes. Importantly, even if the model is in the neighbourhood of a local equilibrium and the gradient is close to $\mathbf{0}$, the size of the Adam update with default $\epsilon = 10^{-8}$ is given by the learning rate. Thus, if there are many parameters in the low-curvature regime around an area of high performance, Adam can exit that area of high performance by taking large steps when gradients are small; we have seen examples of this behaviour in a simple setting in Figure 2.3. Since SGD and Runge–Kutta4 only use local information about the gradient, they would not exit the high performance area if gradients are small there.

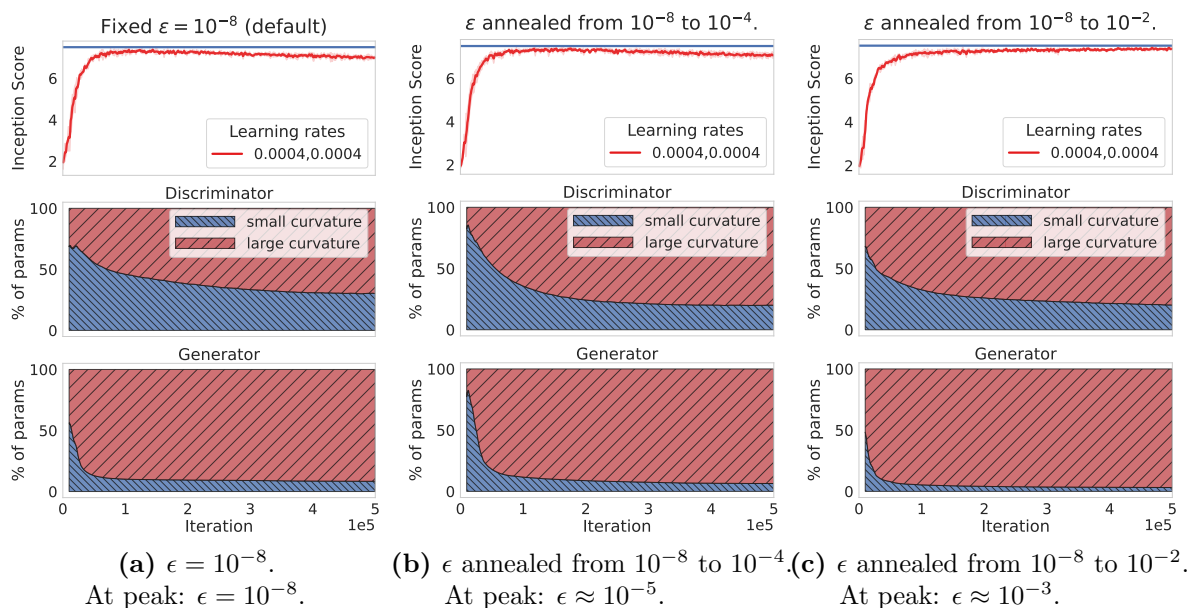


Figure 7.15: The effect of ϵ on Adam performance. Increasing ϵ in training ensures that Adam does not exit areas of good performance, even with many parameters in the low-curvature regime, particularly for the discriminator. We note that using a large ϵ from early in training leads to extremely low performance (likely due to small gradients training does not pick up) and we do not show it here.

In order to understand whether this effect occurs when training GANs in practice, we need to measure whether there are many parameters in a low-curvature regime, as well as investigate the effect of ϵ on GAN training. To answer the first question, we measure

$$\left| \nabla_{\theta_{[l]}} E(\theta_t) - \frac{1}{1 - \beta_1^t} (\mathbf{m}_t)_{[l]} \right| \quad \text{and} \quad \left| \left(\nabla_{\theta_{[l]}} E(\theta_t) \right)^2 - \frac{1}{1 - \beta_2^t} (\mathbf{v}_t)_{[l]} \right| \quad (7.36)$$

for all discriminator and generator parameters. These quantities are readily available when using Adam, we know they are 0 for parameters in areas of low curvature (Eq 7.33 and Eq 7.34), and we know how they affect the Adam update (Eq 7.35). We define a parameter as being in an area of low curvature at iteration t if $\left| \nabla_{\theta_{[l]}} E(\theta_t) - \frac{1}{1 - \beta_1^t} (\mathbf{m}_t)_{[l]} \right| < 10^{-5}$ and $\left| \left(\nabla_{\theta_{[l]}} E(\theta_t) \right)^2 - \frac{1}{1 - \beta_2^t} (\mathbf{v}_t)_{[l]} \right| < 10^{-5}$. Results in Figure 7.15 show that when training GANs with Adam on CIFAR-10 many parameters are in the low-curvature regime when the high performance area is reached, particularly for the discriminator. When using $\epsilon = 10^{-8}$, Figure 7.15a shows that

after the area of high performance is reached, the model exits it and performance degrades; this is consistent with our hypothesis that, due to large updates in areas of low curvature, Adam using a small ϵ can exit areas of high performance. By annealing ϵ throughout training we can ensure that a larger ϵ is used areas of high performance, and thus the Adam update is more sensitive to the gradient magnitude, especially for parameters with small gradients, as we show in Figure 7.14b. We show results in GAN training by annealing the Adam ϵ in Figures 7.15b and 7.15c: the rate at which the high performance area is left decreases, and for ϵ annealed to as high as 10^{-2} , the high performance area is no longer exited.

These results suggest that indeed, even though Spectral Normalisation does not have a purely optimisation effect, its annealing of ϵ when using the Adam optimiser has a beneficial effect, by avoiding large steps in areas of high performance. Furthermore, these results are consistent with what we observed in the DRL case, where Spectral Normalisation ensures that performance increases for longer in training (Figure 7.2). While here we focused on the effect of ϵ in the low-curvature regime, we want to emphasise that this is not the only effect of ϵ , but simply one easy to analyse. As ϵ dampens the size of parameter updates regardless of curvature, increasing ϵ has an effect on all parameters. From that perspective, increasing ϵ in training is akin to decreasing to learning rate. We do note significant differences, however: as we show in Figure 7.14a, a small ϵ in the low-curvature regime always leads to a lack of sensitivity to the gradient magnitude, regardless of learning rate, while changing ϵ can lead to increased sensitivity to gradient magnitudes. We believe more studies in the role of ϵ in Adam are needed, specifically in the context of DRL, where it has empirically been shown (here and elsewhere) to play an important role.

7.5 Related work

For an overview of the use of smoothness regularisation and normalisation in deep learning, together with benefits and downsides, see Chapter 6. The relaxation of the 1-Lipschitz constraint in the application of Spectral Normalisation has been previously done by Gouk et al. [261], who relax the Lipschitz constant but in contrast

to our approach in DRL, still regularise the entire network.

Other normalisation methods, such as Weight Normalisation and Batch Normalisation have been studied from an optimisation perspective [23, 264]. Unlike Spectral Normalisation however, these normalisation techniques are not hard normalisations, as their parametrisation allows recovering the un-normalised approaches. Weight Normalisation [264] decouples the learning of the Frobenius norm of the matrix from its direction; in DRL, the authors show that applying Weight Normalisation improves the performance of DQN on a subset of Atari games. When introducing Spectral Normalisation, Miyato et al. [30] compare Spectral Normalisation and Weight Normalisation empirically and in the GAN case obtain significant gains using Spectral Normalisation compared to Weight Normalisation (their Table 2). The connection between optimisation hyperparameters in Adam and normalisation—as we have done with MULEPS—has been previously made for Batch Normalisation and Weight Normalisation [265–267]; in contrast to our approach, these works require L_2 regularisation and obtain the ϵ scaling effect jointly with a learning rate scaling effect, while we obtained no change on the learning rate (Eq (7.31)). The connection between ϵ and low-curvature is connected to existing works which track curvature during training [167].

While optimisation specific issues in DRL have long been overlooked, recent works have studied the negative interplay between temporal difference learning methods such as Q -learning and adaptive optimisation methods including RMSprop and Adam [244], while others propose optimisers for DRL, but they are either restricted to linear estimators [268, 269], or they show no empirical advantages [270].

Inspired by the insights obtained when investigating the effects of Spectral Normalisation in DRL, we also analysed its effects on GANs. A study of the effects of Spectral Normalisation was similarly performed recently, with Lin et al. [271] exploring the effects of Spectral Normalisation theoretically and uncovering its connections with exploding and vanishing gradients. Their work is related to our approach, which similarly uncovered the effects of Spectral Normalisation on training dynamics, as well as explored the effect of optimisation methods based on Spectral

Normalisation, such as DIVOUT, and the interaction with loss function choices and optimisers such as Adam.

7.6 Conclusion

We investigated the use of Spectral Normalisation in DRL and showed that normalising one layer of the learned neural critic in Q -learning agents can substantially improve performance over methods which collate multiple RL specific improvements. We then investigated *why* Spectral Normalisation helps DRL. By using the compositional structure of neural networks we uncovered strong optimisation effects, including interactions with optimiser hyperparameters. Inspired by these optimisation effects, we developed novel optimisation updates that recover, or even outperform, the effects of Spectral Normalisation in DRL. Equipped with this knowledge we re-examined the effects of Spectral Normalisation in GANs, another domain where Spectral Normalisation has been applied with success, but where the normalisation procedure is applied to the entire network. We observed that for GANs optimisation only changes do not recover Spectral Normalisation performance, and that a strong effect of Spectral Normalisation is avoiding poor-conditioned areas of loss functions. Despite not recovering its performance however, optimisation changes due to Spectral Normalisation can help training when the Adam optimiser is used, particularly for parameters which exhibit low-curvature during training.

Chapter 8

Geometric Complexity: a smoothness complexity measure implicitly regularised in deep learning

Machine learning methods have long faced the challenge of finding the right balance between fitting the training data and generalising beyond said training data [272]. The need to avoid overly complex models that overfit by memorising training datasets has led to the need to measure model complexity, closely followed by the desire to construct regularisation methods that control said model complexity. The aim of complexity measures is thus to capture generalisation (often measured via test set error, where lower is better) via a U-shaped curve: when complexity is low, the model is underfitting the training data and underperforming on the test data leading to a high test error (high bias), but as complexity increases the test error decreases, followed by an increase in test error once the model is overly complex and overfitting (high variance); we have observed this behaviour in Figure 6.1a of Chapter 6. While this traditional U-shaped curve has been long associated with complexity measures in machine learning, deep learning has led to a few puzzling observations. For deep or over-parametrised networks, overfitting is less of an issue than we would expect under the traditional U-shape curve, and the models both fit and the data and

generalise [79, 273, 274]. Deep neural networks present an additional challenge for complexity measures, as many existing measures are either too simplistic, such as number of learnable parameters or weight norms, or intractable to compute for the neural network class of models, such as classical measures like Rademacher [275, 276] complexity and VC dimensions [203, 277].

Recent work has also uncovered connections between optimisation and generalisation in deep learning, while previously optimisation was viewed purely from a training objective minimisation perspective; these works broadly fall under the umbrella of implicit regularisation [12, 13, 65, 278, 279]. Since many existing complexity measures measure the capacity of a family of functions [203, 275, 277], they cannot capture implicit regularisation effects induced by optimisation inside a model class, such as the implicit regularisation induced by discretisation drift we studied in previous chapters.

Motivated by the need to reconcile complexity measures and deep learning, and emboldened by our previous study of implicit regularisation and interactions between smoothness and optimisation, we make the following contributions:

- We introduce a new measure of model complexity for deep neural networks, which we call *Geometric Complexity* (GC). GC captures the complexity of a model rather than function class, can be computed exactly during training, and has connections to Lipschitz smoothness.
- We show that GC is connected with or implicitly regularised by many aspects of deep learning, including initialisation, explicit regularisation, and implicit regularisers induced by the discretisation drift of optimisation methods.
- We show that through implicit regularisation induced by discretisation drift, the implicit minimisation of GC in neural network training can be strengthened through hyperparameter choices such as batch size or learning rate.
- We show that GC captures the double descent phenomenon, and restores the classical complexity versus performance U-shaped curve.

8.1 Introduction

Existing complexity measures. Many existing complexity measures target the complexity of a class of functions. These types of complexity measures range from the very simple, such as number of learned parameters, to more sophisticated, such as Rademacher complexity and VC dimensions [203, 275–277, 280]. While convenient and inexpensive, using the number of learned parameters as a complexity measure does not account for architectures and expressivity; when using the number of parameters to measure the capacity of a neural network, the now famous double descent phenomenon [77–79] occurs in deep learning: as the number of parameters exceeds a critical threshold, the test error of the model decreases again, as opposed to increasing as expected under the U-shaped complexity curve. Rademacher complexity [275, 276] and VC dimension do account for model class [203, 277], but are geared towards measuring capacity of fitting *random labels*. They thus do not account for the specific task and cannot leverage the inherent structure in the data present in many deep learning problems, and are also challenging to compute for specific deep learning architectures, though bounds do exist [281, 282]. Additional challenges get presented by the fact that neural networks have been shown to consistently be able to learn mappings to random labels [274], making it hard to distinguish between classes of neural networks.

Metrics for measuring the complexity of a single model rather than model family are available. They range from simple approaches, such as the norm of the learned parameters [278], to the more complex such as classification margins [283] or flatness [18, 20], Kolmogorov complexity [284] and relatedly minimum description length [285, 286], the number of linear pieces for ReLU networks [14, 287], but they are not without limitations. Metrics based on parameter norms, such as weight norms, are independent of model architectures and for large models it is not inherently clear how changing the parameter norm changes the model predictions. Norm-based metrics have also been shown to correlate poorly with generalisation, especially when stochastic optimisation is used [20]; classification margins can be challenging to compute for data with a high number of dimensions or classes;

while flatness (or equivalently low sharpness) has been shown to be connected with generalisation [12, 18, 20, 141, 168], limitations have been noted [19, 169]. We previously discussed flatness—defined there as the leading eigenvalue of the Hessian—as a measure of generalisation in Chapter 3. Complexity measures for distributions over models, such as PAC-Bayes [288], have also been adapted to neural networks [289], but they are outside of the scope of this thesis.

Neyshabur et al. [281] analysed many existing measures of complexity and showed that they increase alongside the number of the hidden units, and thus cannot explain behaviours observed in over-parametrised neural networks, where increasing the number of hidden units does not lead to increased overfitting (their Figure 5). Jiang et al. [20] perform an extensive study of existing complexity measures and their correlation with generalisation as well as aim to find causal relationships and find that flatness has the strongest positive connection with generalisation.

Explicit and implicit regularisation. Explicit regularisation in deep learning takes many forms, either through adding terms to loss functions or regularising the model architecture [15, 17, 61, 76, 80, 151, 290]. Explicit regularisation can affect model capacity in intricate ways and its effect can strongly depend on the magnitude of regularisation coefficients, with no straightforward method of translating regularisation effects into restrictions on a model class. Thus, despite their practical significance in obtaining increased generalisation performance, the effect of common explicit regularisers has not been unified by existing complexity measures.

Recent years have seen an expansion in the discovery of implicit regularisation effects in deep learning, from the implicit regularisation effects of learning rates [12, 13, 21, 105, 152, 278, 291], mini-batch optimisation noise [165, 292, 293] or initialisation [16, 50, 294–298]. We previously explored implicit regularisation induced by the discretisation drift of optimisers in Chapters 4 and 5. Implicit regularisation effects have been largely unexplained by existing complexity measures; by construction complexity measures of function classes rather than individual models cannot capture the effects of optimisation inside a function class. Novel connections between optimisation hyperparameters such as learning rates and per-model

complexity measures such as sharpness have been recently found (Cohen et al. [52] and Chapter 3), but more investigation into its connection with stochasticity is needed [169]. The role of network depth and overparametrisation as an implicit regulariser has also been studied [101, 281, 297], and has been shown to play an important role in the generalisation properties of neural networks as well as the double descent phenomenon; in Chapter 6 we postulated a connection between depth, overparametrisation, and smoothness, which we will make concrete here.

Smoothness: We highlighted the benefits of smoothness for neural networks extensively in Section 6.4 and intuited its deep connection with generalisation and double descent. Existing works have shown that generalisation bounds can be obtained from smoothness constraints [82]. The connection between smoothness and many implicit and explicit regularisers has been highlighted in Chapter 6; here we will define a smoothness based complexity measure and make this association concrete through theoretical arguments and experimentation.

8.2 Geometric Complexity

We are now ready to define a novel metric of complexity for neural models, with the aim of finding a measure that can capture some of the implicit and explicit regularisation effects previously discussed.

Definition 8.2.1 (Geometric Complexity) *For a neural network $g(\cdot; \boldsymbol{\theta}) : \mathbb{R}^I \rightarrow \mathbb{R}^O$ with parameters $\boldsymbol{\theta} \in \mathbb{R}^D$, we write $g(\cdot; \boldsymbol{\theta}) = \tau(f(\cdot; \boldsymbol{\theta}))$ where τ is the last layer activation function (such as the softmax for multi-label classification). Given a dataset \mathcal{D} with elements \mathbf{x} i.i.d. realisations of random variable \mathcal{X} , we define Geometric Complexity (GC) as*

$$GC(f; \boldsymbol{\theta}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_2^2, \quad (8.1)$$

where $\|\mathbf{v}\|_2$ denotes the Frobenius norm.

GC measures the complexity of a single model $f(\cdot; \boldsymbol{\theta})$, not a model class, can be computed efficiently, and depends on the task at hand through the training dataset \mathcal{D} .

While GC is defined based on the training data, it is an estimate of an expectation:

Remark 8.2.1 *GC is the unbiased estimate of the underlying expectation under $p(\mathbf{x})$ the density of \mathcal{X} :*

$$\mathbb{E}_{p(\mathbf{x})} \left[\left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_2^2 \right] = \int \left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_2^2 p(\mathbf{x}) d\mathbf{x}. \quad (8.2)$$

GC can be connected to Lipschitz smoothness—which we discussed at length in Chapter 6—using the following bound:

Remark 8.2.2 (Connection to Lipschitz smoothness) *If K is the Lipschitz constant of $f(\cdot; \boldsymbol{\theta}) : \mathbb{R}^I \rightarrow \mathbb{R}^O$ w.r.t. the Frobenius norm then we have that:*

$$GC(f; \boldsymbol{\theta}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_2^2 \quad (8.3)$$

$$\leq \min(I, O) \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_{op}^2 \leq \min(I, O) K^2, \quad (8.4)$$

where $\|A\|_2$ and $\|A\|_{op}$ denote the Frobenius and operator norms of A , respectively.

This result follows from $\|A\|_2^2 \leq \text{rank}(A) \|A\|_{op}^2 \leq \min(I, O) \|A\|_{op}^2$, $\forall A \in \mathbb{R}^I \times \mathbb{R}^O$ and $\left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_{op}^2 \leq K^2$. For many problems of concern we have that $\min(I, O) = O$ and thus Remark 8.2.2 exposes that the tightness of the bound between GC and Lipschitz smoothness depends on the number of outputs of the network.

While the Lipschitz smoothness of a function provides an upper bound on GC, there are a few fundamental differences between the two quantities. Firstly, GC is data dependent: a model can have low GC while having high Lipschitz constant due to the model not being smooth in parts of the space where there is no training data. Measures on the entire space, such as Lipschitz constants, are unlikely to explain implicit and explicit regularisation effects when many of the regularisers are applied around data instances. Secondly, GC can be computed exactly given a model and dataset, while for when estimating the Lipschitz constant of neural networks loose upper and lower bounds that rely on function composition are often used [299, 300]. Fazlyab et al. [240] provide an algorithm with tighter bounds by leveraging that

activation functions are derivatives of convex functions and cast finding the Lipschitz constant as the result of a convex optimisation problem; however, their most accurate approach scales quadratically with the number of neurons and only applies to feed forward networks. Sokolić et al. [82] provide an upper bound for the Lipschitz constant of a neural network with linear, softmax and pooling layers restricted to input space \mathcal{X} via $\|f(\mathbf{x}) - f(\mathbf{y})\|_2 \leq \sup_{\mathbf{z} \in \text{convex.hull}(\mathcal{X})} \left\| \frac{df}{dz} \right\|_2 \|\mathbf{x} - \mathbf{y}\|_2$, $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$, but empirically resort to layer-wise bounds. In contrast, GC can be computed exactly regardless of network architecture and can be easily implemented using modern machine learning software frameworks.

8.2.1 GC for deep linear models

Consider a deep linear network $lin(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{W}_n \dots \mathbf{W}_1 \mathbf{x} + \mathbf{b}$, with $\boldsymbol{\theta} = \{\mathbf{b}\} \cup \{\mathbf{W}_i, i \in \{1, \dots, n\}\}$ We then have that

$$GC(lin; \boldsymbol{\theta}, \mathcal{D}) = \|\mathbf{W}_n \dots \mathbf{W}_1\|_2^2. \quad (8.5)$$

This shows that the GC of an affine transformation does not depend on the data and recovers the standard Frobenius norm, which is often used as a regulariser, particularly in ridge regression.

8.2.2 GC for rectifier feed-forward networks

Networks with rectifier activations, such as ReLu [301] and Leaky Relu [302], are piecewise affine models, that is for each \mathbf{x}_i we have that $f(\mathbf{x}_i; \boldsymbol{\theta}) = \mathbf{A}_i \mathbf{x}_i + \mathbf{b}_i$, where \mathbf{A}_i and \mathbf{b}_i depend on where in the partition of the input space input \mathbf{x}_i falls in. We can then write

$$GC(f; \boldsymbol{\theta}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left\| \frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} \right\|_2^2 = \sum_{p=1}^P \frac{n_p}{|\mathcal{D}|} \|\mathbf{A}_p\|_2^2, \quad (8.6)$$

where n_p is the number of times an input is in the partitioning corresponding to the linear piece $\mathbf{A}_p \mathbf{x} + \mathbf{b}_p$ in the piecewise model. Eq (8.6) shows that for models with rectifier activations, GC can be approximated by using batches, instead of the entire dataset, as long as the batch size is large enough to include representation

from many piecewise branches of the network.

To further understand the form of GC for rectifier models, let's consider a simple two-layer network, where we have $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_2 a(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$ with $\boldsymbol{\theta} = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2\}$. Denote $\mathbf{h} = a(\mathbf{z}) = a(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$, where a is the rectifier activation function. We then have

$$\frac{df}{d\mathbf{x}} = \mathbf{W}_2 \text{diag}(a'(\mathbf{z})) \mathbf{W}_1, \quad (8.7)$$

where $\text{diag}(\mathbf{v})$ is the diagonal matrix with diagonal entries given by \mathbf{v} . Since $(\text{diag}(a'(\mathbf{z})) \mathbf{W}_1)_{i,j} = a'(\mathbf{z}_i) \mathbf{W}_{i,j}$ and since in the case of rectifiers $a'(\mathbf{z}_i)$ is a piecewise constant function (such as $a'(x) = 0$, if $x < 0$ or 1 otherwise for ReLU activations), we see the connection between L_2 regularisation and weight decay in rectifier networks, and how decreasing the norm of the weight matrices will decrease GC. We observe that unlike weight norms, GC is sensitive to the impact the weight has on the output of the network. This can be most clearly seen for the ReLU activation, where the activation derivative a' can be 0 and thus for units that do not contribute to the output of a data instance, no norm penalty is imposed. Consider the case of unit $\mathbf{h}_i = 0$, where $\mathbf{z}_i = (\mathbf{W}_1)_i^T \mathbf{x} + (\mathbf{b}_1)_i < 0$. Then the entries of row $(\mathbf{W}_1)_i$, which don't contribute to the model output, will not be penalised by $\left\| \frac{df}{d\mathbf{x}} \right\|$ in GC, since they are modulated by $a'(\mathbf{z}_i) = 0$ in Eq 8.7. The same can be said about the column i of \mathbf{W}_2 , which does not contribute to the model output as $\mathbf{o} = \mathbf{W}_2 \mathbf{z} + \mathbf{b}_2$. Thus, GC is robust to weights that do not contribute to model output, such as when a weight matrix has an entry with negative but large absolute value, which leads to a negative pre-activation and a zero hidden unit for all inputs in the dataset. In contrast, the weight norm will be heavily skewed by such a weight due to its large absolute value; this insight might help explain why parameter norm measures are not indicative of generalisation performance [20].

8.2.3 GC for convolutional and residual layers

Since convolutional layers are linear operators, the form of GC when convolutional layers are used is that of Eq (8.6), where the linear operator \mathbf{A}_p will be determined

by the linear operators of the convolutional layers. A significant difference between convolutional and linear layers will be that each parameter of a convolutional filter will appear repeatedly in the linear operator present in GC, which further highlights the difference between GC and weight norms discussed above.

Residual layers [86] have the form of $a(\mathbf{M}\mathbf{x} + \mathbf{x})$, where \mathbf{M} is a composition of convolutional layers followed by the activation function a , usually a rectifier [262, 302]. Alternatively, they can use an additional projection \mathbf{R} leading to $a(\mathbf{M}\mathbf{x} + \mathbf{R}\mathbf{x})$. Thus, networks rectifier activations with residual layers will also be piecewise linear models and Eq (8.6) holds.

8.3 GC increases as training progresses

We now focus on the value of GC at network initialisation, and we show that for deep networks with rectifier activations and commonly used initialisers, the GC of the model is close to 0 at initialisation. To see why, we again leverage the piecewise nature of networks with rectifier activations, together with the fact that network biases are initialised to be 0: at initialisation rectifier MLPs can be written as $f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_n \mathbf{P}_{n-1} \mathbf{W}_{n-1} \dots \mathbf{P}_1 \mathbf{W}_1 \mathbf{x}$ where $\mathbf{P}_i = \text{diag}(a'(\mathbf{z}_i)) = \text{diag}(a'(\mathbf{W}_i \mathbf{z}_{i-1}))$. For ReLu networks, since the diagonal entries of \mathbf{P}_i will be 0 if \mathbf{z}_{i-1} is negative, as the depth of the network increases, the chance that the output is not 0 diminishes. Consider the 1-dimensional case where $f(x; \boldsymbol{\theta}) = w_n a(w_{n-1} \dots a(w_1 x))$ and a is the ReLu activation function. If we assume standard Gaussian initialisations of the form $w_i \sim \mathcal{N}(0, \sigma_i^2)$ we have that $P(w_1 x > 0) = \frac{1}{2}$, thus $P(a(w_1 x) \neq 0) = \frac{1}{2}$. With similar arguments, $P(a(w_2 z_1) \neq 0 | z_1 \neq 0) = \frac{1}{2}$ leading to $P(a(w_2 z_1) \neq 0) = \frac{1}{2}^2$. By induction we can then show that for n layers $P(f(x; \boldsymbol{\theta}) \neq 0) = \frac{1}{2^n}$; this shows that as we add one layer, the probability of the output of the network not being 0 is exponentially decreasing. When we consider higher dimensional weight matrices, the chance of hitting exactly 0 decreases, but the obtained function values will be small due to the decreased chance of consistently obtaining positive numbers under the random initialisation. Thus, we expect GC, as the average norm of $\frac{df(\mathbf{x}; \boldsymbol{\theta})}{d\mathbf{x}} = \mathbf{W}_n \mathbf{P}_{n-1} \mathbf{W}_{n-1} \dots \mathbf{P}_1 \mathbf{W}_1$, to be low at initialisation and decrease as the

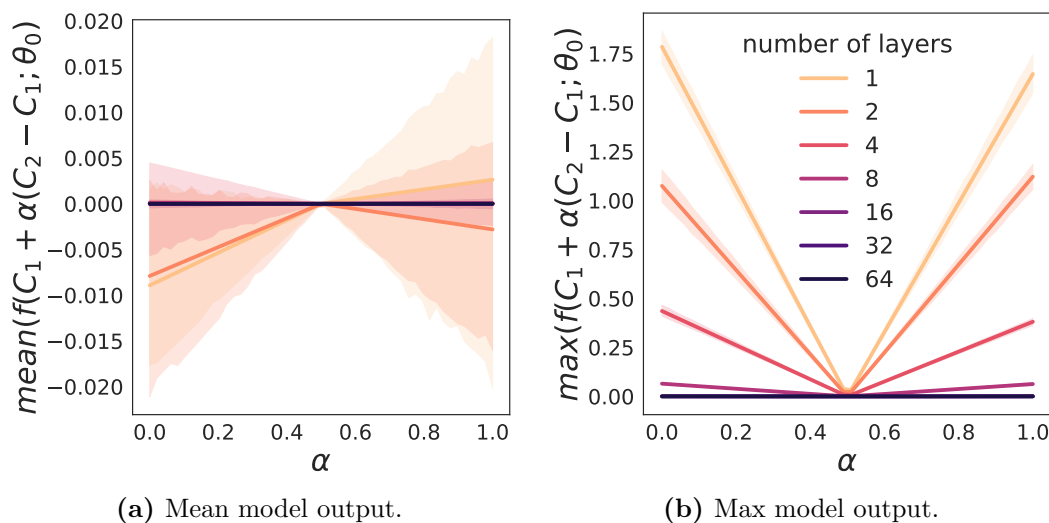


Figure 8.1: Initialisation of MLPs. Outputs given by $f(\mathbf{C}_1 + \alpha(\mathbf{C}_2 - \mathbf{C}_1); \theta)$ where $\mathbf{C}_1 \in \mathbb{R}^I$ with every entry equal to 1 and $\mathbf{C}_2 = -\mathbf{C}_1$ where f is an MLP with 500 neurons per layer; the input dimension $I = 150528 = 224 \times 224 \times 3$. The mean model outputs (a) and max model outputs (b) are closer to 0 at initialisation as depth increases. Uncertainty is obtained by using multiple seeds to initialise the network.

number of layers in the network increases. In Figure 8.1 we confirm this by looking at the effect of depth on the output of an MLP with input dimensions given by those of input dimension $I = 150528 = 224 \times 224 \times 3$ and $O = 1000$ (as from the Imagenet dataset), and show how both the mean and the maximum output obtained from 100 samples are decreasing as the number of layers increases. To empirically assess the effect of the number of layers on GC at initialisation, we use MLPs initialised either with the default initialisation in many deep learning packages, which we will call ‘standard initialisation’ (with weights sampled from a truncated normal distribution bounded between -2 and 2 and scaled by the square root of number of incoming units), and Glorot initialisation [303]. We show in Figure 8.2 that the value of GC at initialisation decreases as the number of layers increases. These results show deep MLPs are smoother at initialisation, in that they are less sensitive to the value of the input and are close to the constant function 0; this is a useful initialisation to start with as otherwise the training procedure might have to undo the effects of a bad initialisation (where for example inputs from the same class might be initialised to have very different output logits), and furthermore increases reliability as training

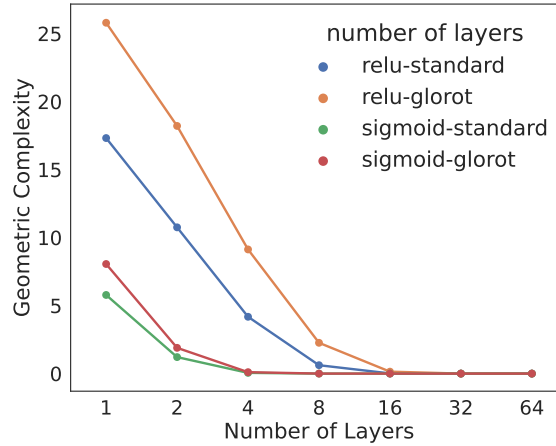


Figure 8.2: Geometric Complexity (GC) at initialisation decreases as the number of layers in an MLP increases; this reflects the initialisation phenomenon showcased in Figure 8.1. This effect shows that the deeper the MLP, the smoother it is at initialisation.

schemes become less dependent on initialisations. We also observe that GC captures this effect at initialisation, and helps explain why deeper models can better model the training data.

We have seen how deep MLPs initialised with commonly used initialisation schemes start training at low (close to 0) GC. We also observe that as training progresses and the model captures the training data, its GC increases; we see this in Figure 8.3 and we will consistently observe this when training larger models as well (see Figures 8.6 and 8.7). During training however, deep networks do not learn overly complex functions, even as GC increases; as shown in Figure 8.3, even when using a large MLP to train 10 data points, the model fits the data without overfitting. We now turn to explore the potential mechanisms for this observation.

8.4 GC and regularisation

We now highlight how many implicit and explicit regularisation methods regularise GC, despite not targetting it directly. To disentangle the effects studied here from optimisation choices, we use stochastic gradient descent without momentum to train all models unless otherwise specified. We also study the impact of a given training heuristic on GC in isolation of other techniques to avoid masking effects. For this reason, we do not use data augmentation or learning rate schedules.

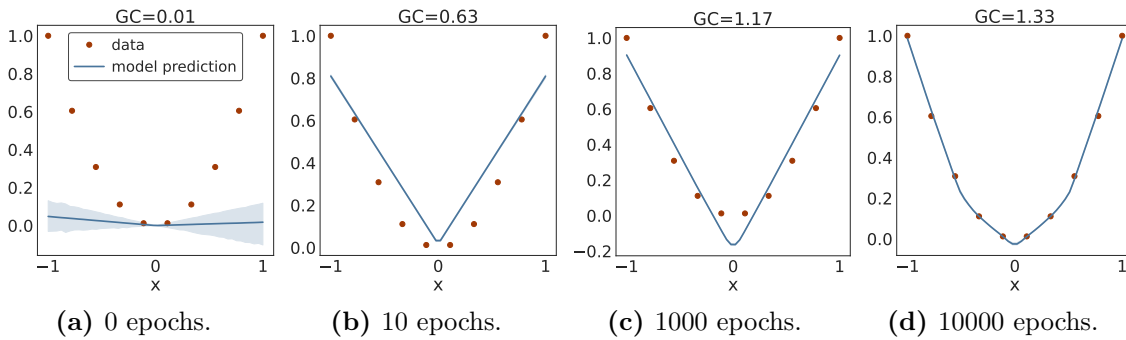


Figure 8.3: For a large MLP fitting 10 points in 1 dimension, the GC of the function being learned gradually grows in training. Despite having capacity to overfit, the model learns a simple function while fitting the data.

8.4.1 Explicit regularisation

L_2 norm regularisation. Section 8.2.2 has shown the connection between GC and weight norms in the case of rectifier neural networks. L_2 norm regularisation adds the regulariser $\zeta \sum_{l=1}^N \|\mathbf{W}_l\|_2^2$ with $\zeta > 0$ to the loss function, minimising the magnitude of each weight parameter. This results in a decreased GC in the case of rectifier neural networks where the GC can be written using the form in Eq (8.6), where the value of each piecewise linear function will be determined by the magnitude of the network’s weights. While this is not a formal argument, we confirm in Figure 8.4a that as the L_2 regularisation strength increases, GC decreases.

Gradient norm regularisation. Another popular form of regularisation is to add an explicit regulariser with the norm of the gradient $\|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})\|_2^2$ to the loss function [13, 15, 17, 61, 67, 151]. While GC is concerned with gradients with respect to inputs, and the gradient norm regulariser penalises the gradient norm with respect to parameters, the two are connected via backpropagation; we previously used neural network structure and backpropagation to connect the effect of an input smoothness regularisation technique—Spectral Normalisation—to scaling of gradients with respect to parameters in Chapter 7. As we see in backpropagation Eqs (7.16)—(7.19), decreasing the gradient norm $\|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})\|_2^2$ can be achieved by decreasing the magnitude of the weight matrix entries or that of the activation derivative $a'(\mathbf{z})$, which for ReLU networks can be achieved with increased sparsity; both of these will lead to a decrease in GC. We show in Figure 8.4b that indeed,

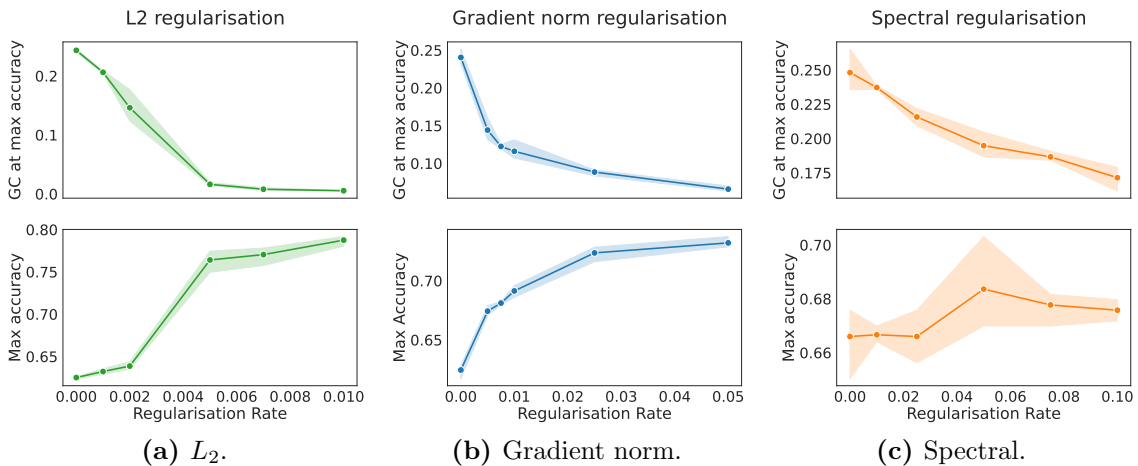
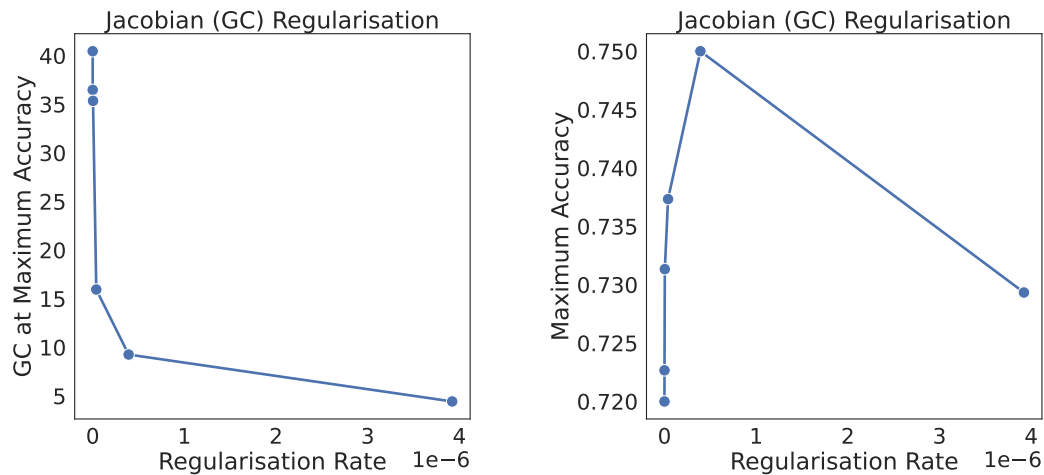


Figure 8.4: Common explicit regularisation methods lead to decreased GC when using vanilla stochastic gradient descent to train a Resnet-18 model on CIFAR-10. GC decreases with increased regularisation strength. Similar trends are observed when momentum is used, as shown in Section F.2 of the Appendix, where we also show loss curves.

increased gradient norm regularisation consistently leads to a decrease in GC.

Spectral Normalisation/regularisation. GC is deeply connected to Lipschitz smoothness; as we have shown in Section 8.2, the Lipschitz constant of the model is an upper bound on its GC. Hard Lipschitz normalisation methods, such as Spectral Normalisation, will thus impose a fixed upper bound on the GC of the model, while regularisation methods such as Spectral Regularisation induce a regularisation pressure to decrease that bound; for an overview of these methods as well as a contrast between soft and hard constraints, see Chapter 6. We show in Figure 8.4c that minimising the spectral norm of each layer, as done by Spectral Regularisation, leads to a decrease in GC. Furthermore, as the strength of Spectral Regularisation increases—by increasing the regularisation coefficient—GC decreases.

Jacobian regularisation. Jacobian regularisation is an *existing* form of explicit regularisation which *directly minimises GC*, i.e. it adds $\zeta GC(f; \theta, \mathcal{D}) = \frac{\zeta}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \left\| \frac{df(\mathbf{x}; \theta)}{d\mathbf{x}} \right\|_2^2$ to the loss function. Jacobian regularisation has been correlated with increased generalisation but also robustness to input distributional shift [29, 80, 82, 304]. Sokolić et al. [82] show that adding Jacobian regularisation to the loss function can lead to an increase in test set accuracy (their Tables III, IV, and



(a) GC decreases with increased regularisation strength.

(b) Accuracy increases with GC regularisation strength, unless too much regularisation is used.

Figure 8.5: Jacobian (GC) regularisation on CIFAR-10. Explicit regularisation minimising GC (already existing in literature as Jacobian regularisation) leads to increased generalisation, unless the regularisation coefficient is too large. This effect reflects what we have observed in other settings, where too much smoothness can hurt performance, from an illustrative example in Figure 1.2, to the Spectral Normalisation example in Figure 6.3d and the reinforcement learning examples in Figure 7.5. Results obtained using 1 seed.

V). We confirm their findings, namely an increase in test accuracy and a decrease in GC with Jacobian/GC regularisation in Figure 8.5. As we have noted previously, too much smoothness can hurt model capacity and decrease performance, which is something that we observe here too for large regularisation coefficients.

8.4.2 Implicit regularisation via discretisation drift

We now show how the implicit regularisation induced by discretisation drift in gradient descent optimisation, previously discussed in Chapters 2, 4 and 5, can lead to a decrease in GC due the connection between gradient norm minimisation and GC in neural networks. Since the strength of implicit regularisation is affected by optimisation hyperparameters such as learning rate and batch size, we empirically observe how changing these hyperparameters leads to direct effect on GC.

Learning rate. Large learning rates have been consistently shown to have an implicit regularisation effect and lead to increased generalisation [12, 13, 105, 305]. Attempts at formalising this include the Implicit Gradient Regularisation (IGR) flow

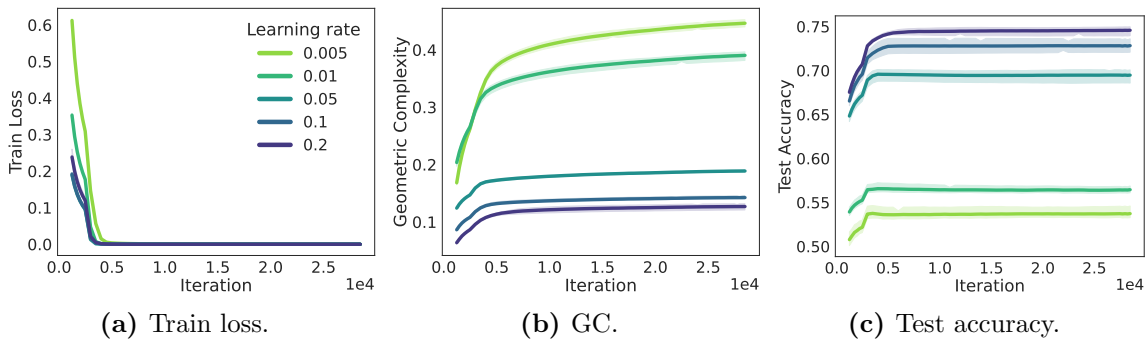


Figure 8.6: The implicit regularisation effect of large learning rates on GC: increasing the learning rate results in decreased GC and increased performance. The results here use vanilla gradient descent; a similar trend is observed when momentum is used, as shown in Section F.2 of the Appendix.

discussed in Chapter 2, which suggests that gradient descent implicitly minimises the loss function $\tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E\|^2$, where h is the learning rate. In the previous section, we described how backpropagation reveals shared components between gradients w.r.t. parameters— $\nabla_{\boldsymbol{\theta}} E$ present in IGR—and gradients w.r.t. inputs— $\nabla_{\mathbf{x}} f$ present in GC. We then showed that explicitly regularising the gradient norm $\|\nabla_{\boldsymbol{\theta}} E\|^2$ leads to a decrease in GC. Based on these results, we expect that increasing the learning rate h leads to a decrease in GC due to increased implicit gradient regularisation pressure $\frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E\|^2$, which we empirically observe in Figure 8.6.

Batch size. The increased generalisation performance of lower batch sizes has been long studied in deep learning, often under the name ‘the generalisation gap’ [18, 292, 305–307]. While traditionally this increased generalisation effect has been attributed to gradient noise due to the sampling process, recent work has revisited this assumption: the generalisation gap was recently bridged with Geiping et al. [151] showing that full-batch training can achieve the same test set performance as stochastic gradient descent. Geiping et al. [151] use an explicit regulariser inspired by IGR, with the intuition that the strength of the implicit gradient norm regulariser induced by discretisation drift $\|\nabla_{\boldsymbol{\theta}} E\|^2$ is stronger for stochastic gradient descent compared to full-batch gradient descent, as described by Smith et al. [13]; we previously discussed the implicit regularisation induced by discretisation drift in stochastic gradient descent and the results of Smith et al. [13] in Chapter 5. Thus,

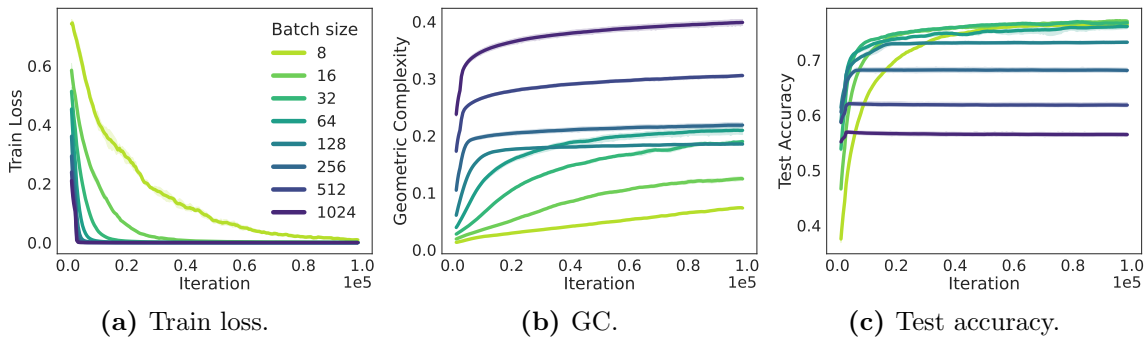


Figure 8.7: The implicit regularisation effect of small batch sizes on GC: decreasing the batch size results in decreased GC and increased performance. The results here use vanilla gradient descent; a similar trend is observed when momentum is used, as shown in Section F.2 of the Appendix.

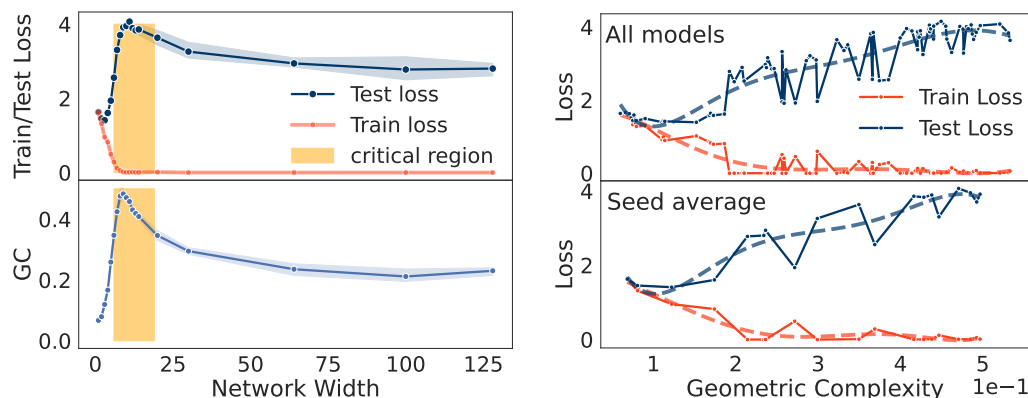
Geiping et al. [151] use $\frac{h}{4|\mathcal{B}|} \sum_{B \in \mathcal{B}} \left\| \frac{1}{B} \sum_{\mathbf{x} \in B} \nabla_{\boldsymbol{\theta}} E(\mathbf{X}; \boldsymbol{\theta}) \right\|^2$ instead, where \mathcal{B} is the set of batches used for training, in contrast to $\frac{h}{4} \left\| \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} \nabla_{\boldsymbol{\theta}} E(\mathbf{X}; \boldsymbol{\theta}) \right\|^2$, the IGR regulariser in the full-batch case. Since the gradient norm is computed over a batch rather than the entire dataset, outliers, such as examples with a large gradient, are less likely to be averaged out and this leads to a stronger regularisation pressure. Their result further strengthens the connection between small batch sizes and increased gradient norm regularisation pressure, which we have consistently seen to lead to smaller GC. Coupling these results, suggesting that for smaller batch sizes there is a stronger effect due to IGR, with what we have previously observed, namely that minimising $\|\nabla_{\boldsymbol{\theta}} E\|^2$ leads to a decrease in GC, we should expect that decreasing batch sizes leads to a decrease in GC. We empirically assess this question and show results in Figure 8.7, observing a consistent decrease in GC with smaller batch sizes. These results show the implicit regularisation effect of batch sizes on model smoothness.

Results and justifications in this section use vanilla gradient descent, but we note that the gradient norm minimisation effect induced by the discretisation drift of momentum has recently been mathematically formalised by Ghosh et al. [308], and thus we expect our arguments to hold for gradient descent with momentum too. We confirm this with empirical results in Appendix F.

8.5 Double descent and GC

When neural network complexity is measured using a simple parameter count, a double descent phenomenon has been consistently observed: as the number of parameters increases, the test loss decreases at first, before increasing again, followed by a second descent toward a low test error value [77–79]. An excellent overview of the double descent phenomena in deep learning can be found in Belkin [77], together with connections to smoothness as an inductive bias of deep networks and the role of optimisation. We previously discussed the double descent phenomenon in Section 6.4, where we contrasted it with the U-shaped curve associated with standard machine learning methods and complexity measures in Figure 6.1. We discussed the double descent phenomenon from the smoothness perspective and intuited that smoothness likely plays an important role in the second descent: we postulated that in the first descent, the generalisation error is decreasing as the model is given extra capacity to capture the decision surface; the increase happens when the model has enough capacity to fit the training data, but it cannot do so and retain smoothness, and thus overfits; the second descent occurs as the capacity increases and smoothness can be retained.

To explore whether GC as a measure of model complexity based on input smoothness can capture the double descent phenomena we follow the set up introduced in Nakkiran et al. [79]: we train multiple ResNet-18 networks on CIFAR-10 with increasing layer width, and show results in Figure 8.8. We make two observations: first, like the test loss, GC follows a double descent curve as width increases (Figure 8.8a); second, when plotting GC against the test loss, we observe a U-shape curve, recovering the traditional expected behaviour of complexity measures (Figure 8.8b). Importantly, we observe the connection between the generalisation properties of over-parametrised models and GC: after the critical region, increase in model size leads to a decrease in GC. These results provide further evidence that GC is able to capture model capacity in a meaningful way, suggestive of a reconciliation of traditional complexity theory and deep learning, as well as provide an empirical basis for our previously highlighted connection between smoothness and double descent.



(a) GC and train/test losses as the network width increases. (b) Train and test losses as GC increases.

Figure 8.8: Double descent and GC. (a): GC captures the double descent phenomenon. (b): GC captures the traditional U-shape curve, albeit with some noise, showing (top) GC compared to train and test losses for all models and (bottom) GC compared to train and test losses averaged across different seeds. We fit a 6 degree polynomial to the curves to showcase the trend.

8.6 Related work

Complexity measures for deep neural networks. There has been a recent effort to create complexity measures that capture both the generalisation phenomenon of deep learning and are computationally tractable. Lee et al. [309] introduce ‘Neural Complexity Measures’ in which they use a neural network to approximate the excess risk—the difference between empirical and true risk, which the author call ‘the generalisation gap’—from validation data and use it as a regulariser for the model; Nagarajan and Kolter [16] argue for an initialisation dependent measure of complexity; Jiang et al. [283] approximate classification decision margins for neural networks and show their predictive generalisation gap correlates to the true generalisation gap. Maddox et al. [310] use the effective dimensionality [75] of the Hessian matrix $\nabla_{\theta}^2 E$ and show that it can capture the double descent phenomenon; since effective dimensionality depends on the eigenspectrum of the Hessian which is computationally expensive to compute for large networks, the authors select the leading eigenvalues to estimate the metric.

Most recently, Grant and Wu [311] modernised the Generalised Degrees of Freedom (GDof) [312, 313] complexity measure by providing a tractable approach

to estimate it online during training. GDof is a similar measure to GC in that it computes the model sensitivity with respect to the dataset labels (GC instead computes the sensitivity to inputs). Grant and Wu [311] showed that GDof can resolve the double descent phenomenon; we hope that future work will assess the connection between GC and GDof both empirically and theoretically.

Huh et al. [314] argue that deep neural networks exhibit a low-rank bias, measured as the rank of the Gram matrix of feature embeddings. By measuring a relationship between data points, a low-rank measure is less local than GC, which assess how sensitive the model is to individual input changes. Huh et al. [314] connect low-rank bias with depth and initialisation, but contend that depth plays a more significant role than initialisation and optimisation. Instead, we have associated the implicit bias induced by GC with optimisation and optimisation hyperparameters.

Jacobian regularisation. Using GC as an explicit regulariser has existed previously as Jacobian regularisation, used to induce model robustness with respect to inputs and increase generalisation [29, 80, 82, 304]. Our goal here was not to introduce a new regulariser, but to show that many existing implicit or explicit regularisers minimise GC, and that GC can be used as a neural network complexity measure.

Gradient penalties. Gradient penalties are used to regularise $f : \mathbb{R}^I \rightarrow \mathbb{R}$ using a penalty of the form $\sum_{\mathbf{x}} (\|\nabla_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta})\| - K)^2$ with K a hyperparameter; instead of minimising the Jacobian norm, these methods ensure it is close K . For additional discussion on gradient penalties, see Section 6.3. Gradient penalties have been primarily used for GAN training [28, 36, 202]; we leave the investigation of the importance of GC outside supervised learning for future work.

Smoothness and generalisation. Novak et al. [83] is perhaps the closest work to ours: they show the correlation across multiple datasets between generalisation and

$\frac{1}{\mathcal{D}_{test}} \sum_{\mathbf{x}_i \in \mathcal{D}_{test}} \left\| \left. \frac{d(\tau \circ f(\cdot; \boldsymbol{\theta}))}{d\mathbf{x}} \right|_{\mathbf{x}_i} \right\|_2$ where τ is the last layer activation function (such as the softmax). They also briefly show that their Jacobian norm metric captures the regularisation effect of stochastic gradient descent compared to full-batch training using L-BFGS [315]. Since our goal is to find a complexity measure we do not use the test set but the training set instead; our main focus is to show the connections

between GC implicit and explicit regularisation, as well as initialisation in deep learning.

Bartlett et al. [81] provide a generalisation bound depending on classification margins and the product of spectral norms of the model’s weights and show empirically that the product of spectral norms correlates with excess risk. Bubeck and Sellke [316] connect Lipschitz smoothness with over-parametrisation, by providing a probabilistic upper bound on the model’s Lipschitz constant based on its training performance and the inverse of its number of parameters.

Smoothness and double descent. In concurrent work, Gamba et al. [317] discuss the connection between the Jacobian and the Hessian of the *loss* with respect to inputs, namely the empirical average of $\|\nabla_{\mathbf{x}}E\|_2$ and $\|\nabla_{\mathbf{x}}^2E\|_2$ over the training set and show that they follow a double descent curve. Like our work with GC, Gamba et al. [317] mention how their work suggests our postulate in Chapter 3: ‘*Finally, Rosca et al. (2020) investigate how to encourage smoothness, either on the entire input space or around training data points. Interestingly, they postulate a connection between model-wise double descent and smoothness: during the first ascent the model size is large enough to fit the training data at the expense of smoothness, while the second descent happens as the model size becomes large enough for smoothness to increase. [...] Our work provides empirical evidence supporting the postulate of Rosca et al. (2020).*’

8.7 Limitations and future work

Our experiments use the Resnet-18 and MLP architectures, and the ReLu activation function. We hope that future work will further validate these results in a variety of settings, including for other architectures, activation functions, and input domains.

Our results suggest that GC is able to capture a model’s complexity and its generalisation abilities on a test set obtained using i.i.d. sampling from the same underlying distribution as the training set. We hope future work will investigate the properties of GC when the test time distribution differs from the training distribution, as is the case in transfer learning [72, 189], or when the training set gradually changes

in time, as in continual learning [71]. We intuit that the results will depend on the distance between the training and test distribution in the case of transfer learning, or the rate of change of the input in the case of continual learning.

We believe incorporating GC for model selection has the potential to construct criteria analogous to the Bayesian information criterion [318], adapted to the deep learning domain. For architecture search GC could prove to be useful by excluding architectures that have high GC at initialisation on the training data. While we did not investigate this here, we hope this is a fruitful avenue for future work.

8.8 Conclusion

We introduced Geometric Complexity (GC), a measure of model complexity defined as the average Jacobian norm of a network’s logits with respect to inputs from the training dataset. We showed that GC is closely connected to Lipschitz smoothness, but unlike the Lipschitz constant of a deep network, GC is tractable and can be measured during training. We have empirically shown how common explicit regularisers that target other quantities, such as weight norm regularisation and gradient norm regularisation, implicitly regularise GC. We then connected GC with the implicit regularisation induced by the discretisation drift of gradient descent. We saw how mechanisms that induce higher implicit regularisation, such as decreased batch size and increased learning rate, lead to a decrease in GC. Last but not least, we showed how GC can resolve the double descent challenge associated with traditional complexity measures and recovers the expected U-shaped curve.

Chapter 9

Conclusion

Throughout the thesis we have seen that despite being fundamental ingredients in the melting pot behind deep learning’s recipe for success, the effects of gradient-based optimisation and smoothness regularisation are far from being understood. Our aim was to take a few steps towards a practical theory of deep learning, inspired by great recent works making strides in that direction. To this end, we started with investigating gradient descent, a perhaps deceptively simple algorithm, which forms the basis of all modern deep learning optimisation. We took a continuous-time perspective to understanding gradient descent, either by quantifying discretisation drift between existing continuous-time flows and gradient descent or by finding new flows that better describe the gradient descent trajectory. We then constructed novel continuous-time analyses to determine sources of instability in supervised learning and two-player games such as Generative Adversarial Networks, which led to mitigation strategies either in the form of adaptive learning rates or explicit regularisers. We continued our investigation by studying interactions between optimisation and models, in particular model smoothness with respect to inputs and related regularisers. We expanded the use of smoothness regularisation to a new domain, reinforcement learning, where we observed that Spectral Normalisation leads to a consistent gain in agent performance that can be recovered through changes to optimisation. When revisiting the use of Spectral Normalisation in Generative Adversarial Networks, however, we encountered a mixture of model regularisation and optimisation effects. We concluded by introducing a measure of model complexity based on smoothness and explored its interactions with other ingredients of the deep learning melting pot,

such as implicit regularisation effects induced by discretisation drift, initialisation, and explicit regularisation.

Motivated by our desire to dig deeper into why deep learning works so well, we employed theoretical and empirical tools to perform an investigation into individual deep learning components and their many interactions. Using our novel intuitions and understanding of deep learning, we constructed methods that led to increased training stability or improved evaluation performance across multiple problem domains. We hope both the approach and the results presented here can form the fruitful basis of further lines of inquiry, for the author as well as the wider community.

Bibliography

- [1] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold. In *Nature*, 2021.
- [2] Suman Ravuri, Karel Lenc, Matthew Willson, Dmitry Kangin, Remi Lam, Piotr Mirowski, Megan Fitzsimons, Maria Athanassiadou, Sheleem Kashem, Sam Madge, Rachel Prudden, Amol Mandhane, Aidan Clark, Andrew Brock, Karen Simonyan, Raia Hadsell, Niall Robinson, Ellen Clancy, Alberto Arribas, and Shakir Mohamed. Skilful precipitation nowcasting using deep generative models of radar. In *Nature*, 2021.
- [3] Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Francisco J. R. Ruiz, Julian Schrittwieser, Grzegorz Swirszcz, David Silver, Demis Hassabis, and Pushmeet Kohli. Discovering faster matrix multiplication algorithms with reinforcement learning. In *Nature*, 2022.
- [4] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.

- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.
- [6] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. In *Advances in Neural Information Processing Systems*, 2022.
- [7] Axel Sauer, Tero Karras, Samuli Laine, Andreas Geiger, and Timo Aila. Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. *arXiv preprint arXiv:2301.09515*, 2023.
- [8] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. In *Nature*, 2017.
- [9] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017.

- [10] OpenAI: Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d.O. Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- [11] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *International Conference on Machine Learning*, 2018.
- [12] David GT Barrett and Benoit Dherin. Implicit gradient regularization. In *International Conference on Learning Representations*, 2021.
- [13] Samuel L Smith, Benoit Dherin, David GT Barrett, and Soham De. On the origin of implicit regularization in stochastic gradient descent. In *International Conference on Learning Representations*, 2021.
- [14] Raman Arora, Amitabh Basu, Poorya Mianjy, and Anirbit Mukherjee. Understanding deep neural networks with rectified linear units. In *International Conference on Learning Representations*, 2018.
- [15] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. The numerics of gans. In *Advances in Neural Information Processing Systems*, 2017.
- [16] Vaishnavh Nagarajan and J Zico Kolter. Generalization in deep networks: The role of distance from initialization. 2019.
- [17] Vaishnavh Nagarajan and J Zico Kolter. Gradient descent gan optimization is locally stable. In *Advances in Neural Information Processing Systems*, 2017.
- [18] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning:

- Generalization gap and sharp minima. *International Conference for Learning Representations*, 2017.
- [19] Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning*, 2017.
- [20] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2019.
- [21] Samuel L. Smith, Pieter-Jan Kindermans, and Quoc V. Le. Don't decay the learning rate, increase the batch size. In *International Conference on Learning Representations*, 2018.
- [22] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems*, 2018.
- [23] Nils Bjorck, Carla P Gomes, Bart Selman, and Kilian Q Weinberger. Understanding batch normalization. In *Advances in Neural Information Processing Systems*, 2018.
- [24] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. In *International Conference on Learning Representations*, 2021.
- [25] Zhiyan Ding, Shi Chen, Qin Li, and Stephen J Wright. Overparameterization of deep resnet: zero loss and mean-field analysis. In *Journal of machine learning research*, 2022.
- [26] Bharath K Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert RG Lanckriet. On integral probability metrics, ϕ -divergences and binary classification. *arXiv preprint arXiv:0901.2698*, 2009.

- [27] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, 2017.
- [28] William Fedus, Mihaela Rosca, Balaji Lakshminarayanan, Andrew M Dai, Shakir Mohamed, and Ian Goodfellow. Many paths to equilibrium: Gans do not need to decrease a divergence at every step. In *International Conference on Learning Representations*, 2018.
- [29] Judy Hoffman, Daniel A Roberts, and Sho Yaida. Robust learning with jacobian regularization. *arXiv preprint arXiv:1908.02729*, 2019.
- [30] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018.
- [31] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *International Conference on Learning Representations*, 2015.
- [32] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- [33] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2018.
- [34] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [35] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems*, 2016.

- [36] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, 2017.
- [37] Mikołaj Bińkowski, Dougal J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations*, 2018.
- [38] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *International Conference on Computer Vision*, 2017.
- [39] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International Conference on Machine Learning*, 2019.
- [40] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Advances in Neural Information Processing Systems*, 2021.
- [41] Cyprien de Masson d’Autume, Shakir Mohamed, Mihaela Rosca, and Jack Rae. Training language gans from scratch. In *Advances in Neural Information Processing Systems*, 2019.
- [42] Paul J Werbos. Applications of advances in nonlinear sensitivity analysis. In *System modeling and optimization*, 1982.
- [43] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. In *Nature*, 1986.
- [44] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. In *Neural Computation*, 1989.
- [45] Augustin Cauchy et al. Méthode générale pour la résolution des systemes d’équations simultanées. In *Comp. Rend. Sci. Paris*, number 1847.

- [46] Herbert Robbins and Sutton Monro. A stochastic approximation method. In *The annals of mathematical statistics*, 1951.
- [47] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [48] Tijmen Tieleman and Geoffrey Hinton. Rmsprop: Divide the gradient by a running average of its recent magnitude. coursera: Neural networks for machine learning. *Coursera Neural Networks Machine Learning course*, 2012.
- [49] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, 2019.
- [50] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep relu networks. In *Machine Learning*. Springer, 2020.
- [51] Simon S Du, Chi Jin, Jason D Lee, Michael I Jordan, Aarti Singh, and Barnabas Poczos. Gradient descent can take exponential time to escape saddle points. In *Advances in Neural Information Processing Systems*, 2017.
- [52] Jeremy M Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021.
- [53] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- [54] Omer Elkabetz and Nadav Cohen. Continuous vs. discrete optimization of deep neural networks. In *Advances in Neural Information Processing Systems*, 2021.
- [55] Paul Glendinning. *Stability, instability and chaos: an introduction to the theory of nonlinear differential equations*. Cambridge university press, 1994.

- [56] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations*, 2014.
- [57] Andrew K Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer learning in deep linear networks. In *International Conference on Learning Representations*, 2019.
- [58] Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. In *Neural Networks*, 2020.
- [59] Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In *Conference on Learning Theory*, 2021.
- [60] Guilherme França, Jeremias Sulam, Daniel Robinson, and René Vidal. Conformal symplectic and relativistic optimization. In *Advances in Neural Information Processing Systems*, 2020.
- [61] David Balduzzi, Sebastien Racaniere, James Martens, Jakob Foerster, Karl Tuyls, and Thore Graepel. The mechanics of n-player differentiable games. In *International Conference on Machine Learning*, 2018.
- [62] Robert M May. Simple mathematical models with very complicated dynamics. In *Nature*, 1976.
- [63] Sho Yaida. Fluctuation-dissipation relations for stochastic gradient descent. In *International Conference on Learning Representations*, 2018.
- [64] Kangqiao Liu, Liu Ziyin, and Masahito Ueda. Noise and fluctuation of finite learning rate stochastic gradient descent. In *International Conference on Machine Learning*, 2021.
- [65] Daniel Kunin, Javier Sagastuy-Brena, and Hidenori Tanaka Ganguli, Surya Daniel L.K. Yamins. Symmetry, conservation laws, and learning dynamics in neural networks. In *International Conference on Learning Representations*, 2021.

- [66] Florian Schäfer and Anima Anandkumar. Competitive gradient descent. In *Advances in Neural Information Processing Systems*, 2019.
- [67] Chongli Qin, Yan Wu, Jost Tobias Springenberg, Andy Brock, Jeff Donahue, Timothy Lillicrap, and Pushmeet Kohli. Training generative adversarial networks by solving ordinary differential equations. In *Advances in Neural Information Processing Systems*, 2020.
- [68] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International Conference on Machine Learning*, 2018.
- [69] Yuanhao Wang, Guodong Zhang, and Jimmy Ba. On solving minimax optimization locally: A follow-the-ridge approach. In *International Conference on Learning Representations*, 2019.
- [70] Eric V Mazumdar, Michael I Jordan, and S Shankar Sastry. On finding local nash equilibria (and only local nash equilibria) in zero-sum games. *arXiv preprint arXiv:1901.00838*, 2019.
- [71] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. In *Neural Networks*, 2019.
- [72] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. In *Proceedings of the Institute of Electrical and Electronics Engineers*, 2020.
- [73] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [74] Tom M Mitchell. *The need for biases in learning generalizations*. Citeseer, 1980.
- [75] David MacKay. Bayesian model comparison and backprop nets. In *Advances in Neural Information Processing Systems*, 1991.

- [76] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. In *The Journal of Machine Learning Research*, 2014.
- [77] Mikhail Belkin. Fit without fear: remarkable mathematical phenomena of deep learning through the prism of interpolation. In *Acta Numerica*. Cambridge University Press, 2021.
- [78] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. In *National Academy of Sciences*, 2019.
- [79] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations*, 2019.
- [80] Yuichi Yoshida and Takeru Miyato. Spectral norm regularization for improving the generalizability of deep learning. *arXiv preprint arXiv:1705.10941*, 2017.
- [81] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, 2017.
- [82] Jure Sokolić, Raja Giryes, Guillermo Sapiro, and Miguel RD Rodrigues. Robust large margin deep neural networks. In *Transactions on Signal Processing*, 2017.
- [83] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. In *International Conference on Learning Representations*, 2018.
- [84] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier. Parseval networks: improving robustness to adversarial examples. In *International Conference on Machine Learning*, 2017.

- [85] Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. In *Advances in Neural Information Processing Systems*, 2020.
- [86] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- [87] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, 2015.
- [88] Mihaela Rosca, Yan Wu, Chongli Qin, and Benoit Dherin. On a continuous-time model of gradient descent dynamics and instability in deep learning. In *Transactions on Machine Learning Research*, 2022.
- [89] Mihaela Rosca, Yan Wu, Benoit Dherin, and David GT Barrett. Discretization drift in two-player games. In *International Conference on Machine Learning*, 2021.
- [90] Mihaela Rosca, Theophane Weber, Arthur Gretton, and Shakir Mohamed. A case for new neural network smoothness constraints. 2020.
- [91] Florin Gogianu, Tudor Berariu, Mihaela Rosca, Claudia Clopath, Lucian Busoniu, and Razvan Pascanu. Spectral normalisation for deep reinforcement learning: an optimisation perspective. In *International Conference on Machine Learning*, 2021.
- [92] Benoit Dherin, Michael Munn, Mihaela Rosca, and David GT Barrett. Why neural networks find simple solutions: The many regularizers of geometric complexity. In *Advances in Neural Information Processing Systems*, 2022.
- [93] Changyou Chen, Chunyuan Li, Liqun Chen, Wenlin Wang, Yunchen Pu, and Lawrence Carin Duke. Continuous-time flows for efficient inference and density estimation. In *International Conference on Machine Learning*, 2018.

- [94] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems*, 2018.
- [95] Patrick Kidger. On neural differential equations. *arXiv preprint arXiv:2202.02435*, 2022.
- [96] Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord: Free-form continuous dynamics for scalable reversible generative models. In *International Conference on Learning Representations*, 2019.
- [97] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. In *Advances in Neural Information Processing Systems*, 2021.
- [98] Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient estimation in machine learning. In *Journal of Machine Learning Research*, 2020.
- [99] Ernst Hairer, Marlis Hochbruck, Arieh Iserles, and Christian Lubich. Geometric numerical integration. *Oberwolfach Reports*, 2006.
- [100] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- [101] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2018.
- [102] Peter Bartlett, Dave Helmbold, and Philip Long. Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks. In *International Conference on Machine Learning*, 2018.

- [103] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in Neural Information Processing Systems*, 2014.
- [104] Zhiyuan Li, Yi Zhang, and Sanjeev Arora. Why are convolutional nets more sample-efficient than fully-connected nets? In *International Conference on Learning Representations*, 2021.
- [105] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large learning rate in training neural networks. In *Advances in Neural Information Processing Systems*, 2019.
- [106] Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning: The rprop algorithm. In *International Conference on Neural Networks*, 1993.
- [107] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. In *Nature*, 2015.
- [108] Alex Kearney, Vivek Veeriah, Jaden B Travnik, Richard S Sutton, and Patrick M Pilarski. Tidbd: Adapting temporal-difference step-sizes through stochastic meta-descent. *arXiv preprint arXiv:1804.03334*, 2018.
- [109] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *Association for the Advancement of Artificial Intelligence*, 2018.
- [110] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the

- importance of initialization and momentum in deep learning. In *International Conference on Machine Learning*, 2013.
- [111] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr computational mathematics and mathematical physics*, 1964.
- [112] Aman Gupta, Rohan Ramanath, Jun Shi, and S Sathiya Keerthi. Adam vs. sgd: Closing the generalization gap on image classification. In *Workshop on Optimization for Machine Learning, Advances in Neural Information Processing Systems*, 2021.
- [113] Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Hoi, and Weinan E. Towards theoretically understanding why sgd generalizes better than adam in deep learning. In *Advances in Neural Information Processing Systems*, 2020.
- [114] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2019.
- [115] Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. Understanding the difficulty of training transformers. *arXiv preprint arXiv:2004.08249*, 2020.
- [116] Dami Choi, Christopher J Shallue, Zachary Nado, Jaehoon Lee, Chris J Maddison, and George E Dahl. On empirical comparisons of optimizers for deep learning. *arXiv preprint arXiv:1910.05446*, 2019.
- [117] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, 2017.
- [118] Philip Hartman. A lemma in the theory of structural stability of differential equations. In *Proceedings of the American Mathematical Society*, 1960.

- [119] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks. In *International Conference on Learning Representations*, 2017.
- [120] Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In *International Conference on Machine Learning*, 2020.
- [121] Farzan Farnia and Asuman Ozdaglar. Do gans always have nash equilibria? In *International Conference on Machine Learning*, 2020.
- [122] Tanner Fiez, Benjamin Chasnov, and Lillian J Ratliff. Convergence of learning dynamics in stackelberg games. In *International Conference on Machine Learning*, 2020.
- [123] Hugo Berard, Gauthier Gidel, Amjad Almahairi, Pascal Vincent, and Simon Lacoste-Julien. A closer look at the optimization landscapes of generative adversarial networks. In *International Conference on Learning Representations*, 2020.
- [124] Guy Gur-Ari, Daniel A Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace. *arXiv preprint arXiv:1812.04754*, 2018.
- [125] Justin Gilmer, Behrooz Ghorbani, Ankush Garg, Sneha Kudugunta, Behnam Neyshabur, David Cardoze, George Dahl, Zachary Nado, and Orhan Firat. A loss curvature perspective on training instability in deep learning. In *International Conference on Learning Representations*, 2022.
- [126] Peter L Bartlett, Steven N Evans, and Philip M Long. Representing smooth functions as compositions of near-identity functions with implications for deep network optimization. *arXiv preprint arXiv:1804.05012*, 2018.
- [127] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*, 2018.

- [128] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, 2019.
- [129] Simon Du and Wei Hu. Width provably matters in optimization for deep linear neural networks. In *International Conference on Machine Learning*, 2019.
- [130] Liu Ziyin, Botao Li, James B Simon, and Masahito Ueda. Sgd can converge to local maxima. In *International Conference on Learning Representations*, 2021.
- [131] Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Advances in Neural Information Processing Systems*, 2018.
- [132] Shun-Ichi Amari. Natural gradient works efficiently in learning. In *Neural computation*, 1998.
- [133] Yann Ollivier. Riemannian metrics for neural networks: feedforward networks. In *Information and Inference: A Journal of the Institute of Mathematics*, 2015.
- [134] Yann Ollivier. Riemannian metrics for neural networks ii: recurrent networks and learning symbolic data sequences. *Information and Inference: A Journal of the IMA*, 2015.
- [135] Yang Song, Jiaming Song, and Stefano Ermon. Accelerating natural gradient with higher-order invariance. In *International Conference on Machine Learning*, 2018.
- [136] Mihaela C Rosca, Yan Wu, Benoit Dherin, and David GT Barrett. Discretization drift in two-player games. In *International Conference on Machine Learning*, 2021.
- [137] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in Neural Information Processing Systems*, 2019.

- [138] Gerhard Wanner and Ernst Hairer. *Solving ordinary differential equations II*. Springer Berlin Heidelberg New York, 1996.
- [139] HoHo Rosenbrock. An automatic method for finding the greatest or least value of a function. *The computer journal*, 1960.
- [140] Arthur Asuncion and David Newman. Uci machine learning repository, 2007.
- [141] Stanisław Jastrzebski, Zachary Kenton, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. On the relation between the sharpest directions of dnn loss and the sgd step length. In *International Conference on Learning Representations*, 2019.
- [142] Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht. Gradient descent only converges to minimizers. In *Conference on learning theory*, 2016.
- [143] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, 2016.
- [144] Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the hessian of over-parametrized neural networks. In *Workshop track - International Conference on Learning Representations 2018*, 2017.
- [145] Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization via hessian eigenvalue density. In *International Conference on Machine Learning*, 2019.
- [146] Vardan Papyan. The full spectrum of deepnet Hessians at scale: Dynamics with sgd training and sample size. *arXiv preprint arXiv:1811.07062*, 2018.
- [147] Kaare Brandt Petersen and Michael Syskind Pedersen. *The matrix cookbook*. 2008.
- [148] Alex Damian, Eshaan Nichani, and Jason D Lee. Self-stabilization: The implicit bias of gradient descent at the edge of stability. In *Workshop on Optimization*

- for Machine Learning, Advances in Neural Information Processing Systems*, 2022.
- [149] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017.
- [150] Tong He, Zhi Zhang, Hang Zhang, Zhongyue Zhang, Junyuan Xie, and Mu Li. Bag of tricks for image classification with convolutional neural networks. In *Conference on Computer Vision and Pattern Recognition*, 2019.
- [151] Jonas Geiping, Micah Goldblum, Phillip E Pope, Michael Moeller, and Tom Goldstein. Stochastic training is not necessary for generalization. In *International Conference on Learning Representations*, 2022.
- [152] Stanisław Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho, and Krzysztof Geras. The break-even point on optimization trajectories of deep neural networks. In *International Conference on Learning Representations*, 2019.
- [153] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2020.
- [154] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. In *Journal of machine learning research*, 2011.
- [155] Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. In *Advances in Neural Information Processing Systems*, 2020.
- [156] Bin Shi, Simon S Du, Michael I Jordan, and Weijie J Su. Understanding the acceleration phenomenon via high-resolution differential equations. In *Mathematical Programming*, 2021.

- [157] Nikola B Kovachki and Andrew M Stuart. Continuous time analysis of momentum methods. In *Journal of Machine Learning Research*, 2021.
- [158] Guilherme França, Michael I Jordan, and René Vidal. On dissipative symplectic integration with applications to gradient-based optimization. *Journal of Statistical Mechanics: Theory and Experiment*, 2021.
- [159] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and adaptive stochastic gradient algorithms. In *International Conference on Machine Learning*, 2017.
- [160] Taiki Miyagawa. Toward equation of motion for deep neural networks: Continuous-time gradient descent and discretization error analysis. In *Advances in Neural Information Processing Systems*.
- [161] Kwangjun Ahn, Jingzhao Zhang, and Suvrit Sra. Understanding the unstable convergence of gradient descent. In *International Conference on Machine Learning*, 2022.
- [162] Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on edge of stability in deep learning. In *International Conference on Machine Learning*, 2022.
- [163] Chao Ma, Lei Wu, and Lexing Ying. The multiscale structure of neural network loss functions: The effect on optimization and origin. *arXiv preprint arXiv:2204.11326*, 2022.
- [164] Lei Chen and Joan Bruna. On gradient descent convergence beyond the edge of stability. *arXiv preprint arXiv:2206.04172*, 2022.
- [165] Liu Ziyin, Kangqiao Liu, Takashi Mori, and Masahito Ueda. Strength of mini-batch noise in sgd. In *International Conference on Learning Representations*, 2022.
- [166] Roger Fletcher. *Practical methods of optimization*. John Wiley & Sons, 2013.

- [167] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International Conference on Machine Learning*, 2015.
- [168] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. In *Neural computation*, 1997.
- [169] Simran Kaur, Jeremy Cohen, and Zachary C Lipton. On the maximum hessian eigenvalue and generalization. 2022.
- [170] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [171] Aravind Rajeswaran, Igor Mordatch, and Vikash Kumar. A game theoretic framework for model based reinforcement learning. In *International Conference on Machine Learning*, 2020.
- [172] Satinder P Singh, Michael J Kearns, and Yishay Mansour. Nash convergence of gradient dynamics in general-sum games. In *Uncertainty in Artificial Intelligence*, 2000.
- [173] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, 2017.
- [174] Radford M Neal. Annealed importance sampling. *Statistics and computing*, 2001.
- [175] Alexia Jolicoeur-Martineau. The relativistic discriminator: a key element missing from standard gan. In *International Conference on Learning Representations*, 2019.
- [176] Aditya Grover, Manik Dhar, and Stefano Ermon. Flow-gan: Combining maximum likelihood and adversarial learning in generative models. In *Proceedings of the AAAI conference on artificial intelligence*, 2018.

- [177] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, 2016.
- [178] Haihao Lu. An $o(sr)$ -resolution ode framework for understanding discrete-time algorithms and applications to the linear convergence of minimax problems. In *Mathematical Programming*, 2021.
- [179] Ernst Hairer and Christian Lubich. The life-span of backward error analysis for numerical integrators. In *Numerische Mathematik*, 1997.
- [180] Constantinos Daskalakis and Ioannis Panageas. The limit points of (optimistic) gradient descent in min-max optimization. In *Advances in Neural Information Processing Systems*, 2018.
- [181] Yang You, Zhao Zhang, Cho-Jui Hsieh, James Demmel, and Kurt Keutzer. Imagenet training in minutes. In *International Conference on Parallel Processing*, 2018.
- [182] Aaron Meurer, Christopher P. Smith, Mateusz Paprocki, Ondřej Čertík, Sergey B. Kirpichev, Matthew Rocklin, AMiT Kumar, Sergiu Ivanov, Jason K. Moore, Sartaj Singh, Thilina Rathnayake, Sean Vig, Brian E. Granger, Richard P. Muller, Francesco Bonazzi, Harsh Gupta, Shivam Vats, Fredrik Johansson, Fabian Pedregosa, Matthew J. Curry, Andy R. Terrel, Štěpán Roučka, Ashutosh Saboo, Isuru Fernando, Sumith Kulal, Robert Cimrman, and Anthony Scopatz. SymPy: symbolic computing in python. *PeerJ Computer Science*, 2017.
- [183] Ali Unlu and Laurence Aitchison. Gradient regularisation as approximate variational inference. In *Symposium on Advances in Approximate Bayesian Inference*, 2020.
- [184] Stanisław Kamil Jastrzebski, Devansh Arpit, Oliver Åstrand, Giancarlo Kerg, Huan Wang, Caiming Xiong, Richard Socher, Kyunghyun Cho, and Krzysztof J.

- Geras. Catastrophic fisher explosion: Early phase fisher matrix impacts generalization. In *International Conference on Machine Learning*, 2021.
- [185] Florian Schäfer, Hongkai Zheng, and Animashree Anandkumar. Implicit competitive regularization in gans. In *International Conference on Machine Learning*, 2020.
- [186] Weijie Su, Stephen Boyd, and Emmanuel J Candes. A differential equation for modeling nesterov’s accelerated gradient method: theory and insights. *The Journal of Machine Learning Research*, 2016.
- [187] Andre Wibisono, Ashia C Wilson, and Michael I Jordan. A variational perspective on accelerated methods in optimization. *proceedings of the National Academy of Sciences*, 2016.
- [188] Yoav Shoham and Kevin Leyton-Brown. *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge University Press, 2008.
- [189] Mohammadreza Iman, Khaled Rasheed, and Hamid R Arabnia. A review of deep transfer learning and recent advancements. *arXiv preprint arXiv:2201.09679*, 2022.
- [190] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [191] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016.
- [192] Carlos Eduardo Rosar Kos Lassance, Vincent Gripon, and Antonio Ortega. Laplacian networks: Bounding indicator function smoothness for neural network robustness. In *Transactions on Signal and Information Processing*, 2018.

- [193] Pierluca D’Oro and Wojciech Jaśkowski. How to learn a useful critic? model-based action-gradient-estimator policy optimization. In *Advances in Neural Information Processing Systems*, 2020.
- [194] Michael Tobias Tschannen, Josip Djolonga, Paul Kishan Rubenstein, Sylvain Gelly, and Mario Lučić. On mutual information maximization for representation learning. In *International Conference on Learning Representations*, 2020.
- [195] Herbert Federer. *Geometric measure theory*. Springer, 1996.
- [196] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
- [197] RV Mises and Hilda Pollaczek-Geiringer. Praktische verfahren der gleichungsaufösung. *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik*, 9(1):58–77, 1929.
- [198] Henry Gouk, Eibe Frank, Bernhard Pfahringer, and Michael J Cree. Regularisation of neural networks by enforcing lipschitz continuity. In *Machine Learning*, 2020.
- [199] Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention. In *International Conference on Machine Learning*, 2021.
- [200] Yusuke Tsuzuku, Issei Sato, and Masashi Sugiyama. Lipschitz-margin training: Scalable certification of perturbation invariance for deep neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- [201] Michael Arbel, Dougal Sutherland, Mikołaj Bińkowski, and Arthur Gretton. On gradient regularizers for mmd gans. In *Advances in Neural Information Processing Systems*, 2018.
- [202] Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira. On convergence and stability of gans. *arXiv preprint arXiv:1705.07215*, 2017.

- [203] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013.
- [204] Carl Edward Rasmussen and Zoubin Ghahramani. Occam’s razor. In *Advances in Neural Information Processing Systems*, 2001.
- [205] Peter L Bartlett. For valid generalization the size of the weights is more important than the size of the network. In *Advances in Neural Information Processing Systems*, 1997.
- [206] Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, 2018.
- [207] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, 2016.
- [208] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.
- [209] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [210] Heewoo Jun, Rewon Child, Mark Chen, John Schulman, Aditya Ramesh, Alec Radford, and Ilya Sutskever. Distribution augmentation for generative modeling. In *International Conference on Machine Learning*, 2020.
- [211] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. In *Proceedings of the National Academy of Sciences*. National Acad Sciences, 2019.

- [212] Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Simple and scalable epistemic uncertainty estimation using a single deep deterministic neural network. In *International Conference on Machine Learning*, 2020.
- [213] Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. In *Advances in Neural Information Processing Systems*, 2020.
- [214] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, 2017.
- [215] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [216] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *European symposium on Security and Privacy*, 2016.
- [217] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defences: A survey. In *Transactions on Intelligence Technology*, 2018.
- [218] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *International Conference on Learning Representations*, 2016.
- [219] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *Conference on Computer Vision and Pattern Recognition*, 2017.
- [220] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras,

- and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- [221] Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S Schoenholz, Maithra Raghu, Martin Wattenberg, and Ian Goodfellow. Adversarial spheres. In *Workshop Track, International Conference on Learning Representations*, 2018.
- [222] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [223] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 2014.
- [224] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, 2015.
- [225] Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. In *International Conference on Machine Learning*, 2019.
- [226] XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. In *Transactions on Information Theory*, 2010.
- [227] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems*, 2016.
- [228] Michael Arbel, Liang Zhou, and Arthur Gretton. Generalized energy based models. In *International Conference on Learning Representations*, 2021.
- [229] Zhiming Zhou, Jiadong Liang, Yuxuan Song, Lantao Yu, Hongwei Wang, Weinan Zhang, Yong Yu, and Zhihua Zhang. Lipschitz generative adversarial nets. In *International Conference on Machine Learning*, 2019.

- [230] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R Devon Hjelm. Mine: mutual information neural estimation. In *International Conference on Machine Learning*, 2018.
- [231] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [232] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018.
- [233] Georg Ostrovski, Will Dabney, and Remi Munos. Autoregressive quantile networks for generative modeling. In *International Conference on Machine Learning*, 2018.
- [234] Will Dabney, Georg Ostrovski, David Silver, and Remi Munos. Implicit quantile networks for distributional reinforcement learning. In *International Conference on Machine Learning*, 2018.
- [235] Chris Finlay, Jörn-Henrik Jacobsen, Levon Nurbekyan, and Adam M Oberman. How to train your neural ode. In *International Conference on Machine Learning*, 2020.
- [236] Chiappa Silvia, Jiang Ray, Stepleton Tom, Pacchiano Aldo, Jiang Heinrich, and Aslanides John. A general approach to fairness with optimal transport. In *Association for the Advancement of Artificial Intelligence*, 2020.
- [237] Ray Jiang, Aldo Pacchiano, Tom Stepleton, Heinrich Jiang, and Silvia Chiappa. Wasserstein fair classification. In *Uncertainty in Artificial Intelligence*, 2020.
- [238] C. Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008.

- [239] Joern-Henrik Jacobsen, Jens Behrmann, Richard Zemel, and Matthias Bethge. Excessive invariance causes adversarial vulnerability. In *International Conference on Learning Representations*, 2018.
- [240] Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. In *Advances in Neural Information Processing Systems*, 2019.
- [241] Karol Kurach, Mario Lučić, Xiaohua Zhai, Marcin Michalski, and Sylvain Gelly. A large-scale study on regularization and normalization in gans. In *International Conference on Machine Learning*, 2019.
- [242] Yann LeCun, Patrick Haffner, Léon Bottou, and Yoshua Bengio. Object recognition with gradient-based learning. In *Shape, contour and grouping in computer vision*. Springer, 1999.
- [243] Wojciech M Czarnecki, Razvan Pascanu, Simon Osindero, Siddhant Jayakumar, Grzegorz Swirszcz, and Max Jaderberg. Distilling policy distillation. In *International Conference on Artificial Intelligence and Statistics*, 2019.
- [244] Emmanuel Bengio, Joelle Pineau, and Doina Precup. Interference and generalization in temporal difference learning. In *International Conference on Machine Learning*, 2020.
- [245] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Ko-

- ray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. In *Nature*, 2019.
- [246] T. Schaul, John Quan, Ioannis Antonoglou, and D. Silver. Prioritized experience replay. In *International Conference for Learning Representations*, 2016.
- [247] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. In *International Conference on Learning Representations*.
- [248] Ian Osband, Charles Blundell, A. Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped dqn. In *Advances in Neural Information Processing Systems*, 2016.
- [249] Georg Ostrovski, Marc G Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration with neural density models. In *International Conference on Machine Learning*, 2017.
- [250] Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression. In *Association for the Advancement of Artificial Intelligence*, 2018.
- [251] Matteo Hessel, Joseph Modayil, H. V. Hasselt, T. Schaul, Georg Ostrovski, W. Dabney, Dan Horgan, B. Piot, Mohammad Gheshlaghi Azar, and D. Silver. Rainbow: Combining improvements in deep reinforcement learning. In *Association for the Advancement of Artificial Intelligence*, 2018.
- [252] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 1992.
- [253] Will Dabney, André Barreto, Mark Rowland, Robert Dadashi, John Quan, Marc G Bellemare, and David Silver. The value-improvement path: Towards better representations for reinforcement learning. In *Association for the Advancement of Artificial Intelligence*, 2021.

- [254] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Association for the Advancement of Artificial Intelligence*, 2016.
- [255] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International Conference on Machine Learning*, 2016.
- [256] Peter Henderson, Joshua Romoff, and Joelle Pineau. Where did my optimum go?: An empirical analysis of gradient descent optimization in policy gradient methods. 2018.
- [257] Matteo Hessel, Hubert Soyer, Lasse Espeholt, Wojciech Czarnecki, Simon Schmitt, and Hado van Hasselt. Multi-task deep reinforcement learning with popart. In *Association for the Advancement of Artificial Intelligence*, 2019.
- [258] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. In *Journal of Artificial Intelligence Research*, 2013.
- [259] Kenny Young and Tian Tian. Minatar: An atari-inspired testbed for thorough and reproducible reinforcement learning experiments. *arXiv preprint arXiv:1903.03176*, 2019.
- [260] Johan Samir Obando Ceron and Pablo Samuel Castro. Revisiting rainbow: Promoting more insightful and inclusive deep reinforcement learning research. In *International Conference on Machine Learning*, 2021.
- [261] Henry Gouk, Eibe Frank, Bernhard Pfahringer, and Michael J Cree. Regularisation of neural networks by enforcing lipschitz continuity. In *Machine Learning*, 2020.
- [262] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning*, 2010.

- [263] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.
- [264] Tim Salimans and Diederik P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*, 2016.
- [265] Twan Van Laarhoven. L2 regularization versus batch and weight normalization. *arXiv preprint arXiv:1706.05350*, 2017.
- [266] Elad Hoffer, Ron Banner, Itay Golan, and Daniel Soudry. Norm matters: efficient and accurate normalization schemes in deep networks. In *Advances in Neural Information Processing Systems*, 2018.
- [267] Sanjeev Arora, Zhiyuan Li, and Kaifeng Lyu. Theoretical analysis of auto rate-tuning by batch normalization. In *International Conference on Learning Representations*, 2019.
- [268] A. Givchi and M. Palhang. Quasi newton temporal difference learning. In *Asian Conference on Machine Learning*, 2014.
- [269] Tao Sun, Han Shen, Tianyi Chen, and Dongsheng Li. Adaptive temporal difference learning with linear function approximation. In *Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [270] Joshua Romoff, Peter Henderson, David Kanaa, Emmanuel Bengio, Ahmed Touati, Pierre-Luc Bacon, and Joelle Pineau. Tdprop: Does jacobi preconditioning help temporal difference learning? *arXiv preprint arXiv:2007.02786*, 2020.
- [271] Zinan Lin, Vyas Sekar, and Giulia Fanti. Why spectral normalization stabilizes gans: Analysis and improvements. In *Advances in Neural Information Processing Systems*, 2021.

- [272] Scott Fortmann-Roe. Understanding the bias-variance tradeoff. URL: <http://scott.fortmann-roe.com/docs/BiasVariance.html> (hämtad 2019-03-27), 2012.
- [273] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. In *Workshop Track, International Conference on Learning Representations*, 2014.
- [274] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. In *Communications of the Association for Computing Machinery*, 2021.
- [275] V. Koltchinskii and D. Panchenko. Rademacher processes and bounding the risk of function learning. In *High Dimensional Probability II*. Birkhäuser, 1999.
- [276] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. In *Journal of Machine Learning Research*, 2002.
- [277] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. In *Theory of Probability and its Applications*, 1971.
- [278] Behnam Neyshabur. Implicit regularization in deep learning. *arXiv preprint arXiv:1709.01953*, 2017.
- [279] Chao Ma and Lexing Ying. The sobolev regularization effect of stochastic gradient descent. *arXiv preprint arXiv:2105.13462*, 2021.
- [280] Peter L Bartlett, Stéphane Boucheron, and Gábor Lugosi. Model selection and error estimation. In *Machine Learning*, 2002.
- [281] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. The role of over-parametrization in generalization of neural networks. In *International Conference on Learning Representations*, 2018.

- [282] Martin Anthony and Peter L Bartlett. *Neural network learning: Theoretical foundations*. Cambridge university press, 1999.
- [283] Yiding Jiang, Dilip Krishnan, Hossein Mobahi, and Samy Bengio. Predicting the generalization gap in deep networks with margin distributions. In *International Conference on Learning Representations*, 2018.
- [284] J. Schmidhuber. Discovering neural nets with low kolmogorov complexity and high generalization capability. In *Neural networks : the official journal of the International Neural Network Society*, 1997.
- [285] Geoffrey E Hinton and Richard Zemel. Autoencoders, minimum description length and helmholtz free energy. In J. Cowan, G. Tesauro, and J. Alspector, editors, *Advances in Neural Information Processing Systems*. Morgan-Kaufmann, 1994.
- [286] Léonard Blier and Yann Ollivier. The description length of deep learning models. *Advances in Neural Information Processing Systems*, 2018.
- [287] Thiago Serra, Christian Tjandraatmadja, and Srikumar Ramalingam. Bounding and counting linear regions of deep neural networks. In *International Conference on Machine Learning*, 2018.
- [288] Pierre Alquier. User-friendly introduction to pac-bayes bounds. *arXiv preprint arXiv:2110.11216*, 2021.
- [289] Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *Uncertainty in Artificial Intelligence*, 2017.
- [290] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [291] Sihyeon Seong, Yegang Lee, Youngwook Kee, Dongyoon Han, and Junmo Kim. Towards flatter loss surface via nonmonotonic learning rate scheduling. In *Uncertainty in Artificial Intelligence*, 2018.

- [292] Samuel Smith, Erich Elsen, and Soham De. On the generalization benefit of noise in stochastic gradient descent. In *International Conference on Machine Learning*, 2020.
- [293] Daniel A. Roberts. Sgd implicitly regularizes generalization error. In *Advances in Neural Information Processing Systems 2018 Workshop*. 2018.
- [294] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International Conference on Artificial Intelligence and Statistics*, 2010.
- [295] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In *International Conference on Machine Learning*, 2018.
- [296] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *International Conference on Computer Vision*, 2015.
- [297] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, 2018.
- [298] Yaoyu Zhang, Zhi-Qin John Xu, Tao Luo, and Zheng Ma. A type of generalization error induced by initialization in deep neural networks. In *Mathematical and Scientific Machine Learning*, 2020.
- [299] Aladin Virmaux and Kevin Scaman. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *Advances in Neural Information Processing Systems*, 2018.
- [300] Patrick L Combettes and Jean-Christophe Pesquet. Lipschitz certificates for neural network structures driven by averaged activation operators. *arXiv preprint arXiv:1903.01014*, 2019.

- [301] Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018.
- [302] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *International Conference on Machine Learning*, 2013.
- [303] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International Conference on Artificial Intelligence and Statistics*, 2010.
- [304] Dániel Varga, Adrián Csiszárík, and Zsolt Zombori. Gradient regularization improves accuracy of discriminative models. *arXiv preprint arXiv:1712.09936*, 2017.
- [305] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems*, 2017.
- [306] Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*. Springer, 2012.
- [307] Dominic Masters and Carlo Luschi. Revisiting small batch training for deep neural networks. *arXiv preprint arXiv:1804.07612*, 2018.
- [308] Avrajit Ghosh, He Lyu, Xitong Zhang, and Rongrong Wang. Implicit regularization in heavy-ball momentum accelerated stochastic gradient descent. In *International Conference on Learning Representations*, 2023.
- [309] Yoonho Lee, Juho Lee, Sung Ju Hwang, Eunho Yang, and Seungjin Choi. Neural complexity measures. In *Advances in Neural Information Processing Systems*, 2020.
- [310] Wesley J Maddox, Gregory Benton, and Andrew Gordon Wilson. Rethinking parameter counting in deep models: Effective dimensionality revisited. *arXiv preprint arXiv:2003.02139*, 2020.

- [311] Erin Grant and Yan Wu. Predicting generalization with degrees of freedom in neural networks. In *International Conference on Machine Learning 2022 2nd AI for Science Workshop*, 2022.
- [312] Bradley Efron. Estimating the error rate of a prediction rule: improvement on cross-validation. In *Journal of the American Statistical Association*. Taylor & Francis, 1983.
- [313] Jianming Ye. On measuring and correcting the effects of data mining and model selection. *Journal of the American Statistical Association*, 1998.
- [314] Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The low-rank simplicity bias in deep networks. *Transactions on Machine Learning Research*, 2023.
- [315] Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. In *Mathematical programming*, 1989.
- [316] Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. In *Advances in Neural Information Processing Systems*, 2021.
- [317] Matteo Gamba, Erik Engleson, Mårten Björkman, and Hossein Azizpour. Deep double descent via smooth interpolation. *arXiv preprint arXiv:2209.10080*, 2022.
- [318] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Springer, 2006.
- [319] David M Grobman. Homeomorphism of systems of differential equations. *Doklady Akademii Nauk SSSR*, 1959.
- [320] Yan LeCun, Corinna Cortes, and Christopher J.C. Burges. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [321] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.

- [322] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009.
- [323] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). In *International Conference on Learning Representations*, 2015.
- [324] Ernst Hairer, Christian Lubich, and Gerhard Wanner. *Geometric numerical integration*. Springer-Verlag, Berlin, second edition, 2006.
- [325] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs. 2018. URL <http://github.com/google/jax>.
- [326] Tom Hennigan, Trevor Cai, Tamara Norman, and Igor Babuschkin. Haiku: Sonnet for JAX. 2020. URL <http://github.com/deepmind/dm-haiku>.
- [327] Matteo Hessel, David Budden, Fabio Viola, Mihaela Rosca, Eren Sezener, and Tom Hennigan. Optax: composable gradient transformation and optimisation, in jax! 2020. URL <http://github.com/deepmind/optax>.
- [328] Marco Cuturi and David Avis. Ground metric learning. In *Journal of Machine Learning Research*, 2014.
- [329] Marlos C Machado, Marc G Bellemare, Erik Talvitie, Joel Veness, Matthew Hausknecht, and Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. In *Journal of Artificial Intelligence Research*, 2018.
- [330] Marin Toromanoff, Emilie Wirbel, and Fabien Moutarde. Is deep reinforcement learning really superhuman on atari? leveling the playing field. *arXiv preprint arXiv:1908.04683*, 2019.

Appendices

	241
A On a new continuous-time model of gradient descent	243
A.1 Proofs	243
A.2 Comparison with a discrete-time approach	262
A.3 Experimental details	264
A.4 Additional experimental results	267
A.5 Figure reproduction details	281
B Continuous time models of gradient descent in two-player games	286
B.1 Proof of the main theorems	286
B.2 Proof of the main corollaries	296
B.3 Discretisation drift in Runge–Kutta 4	300
B.4 Stability analysis	304
B.5 SGA in zero-sum games	308
B.6 DiracGAN - an illustrative example	310
B.7 The PF for games: linearisation results around critical points	314
B.8 GAN Experimental details	316
B.9 Additional experimental results	317
B.10 Individual figure reproduction details	329
C Finding new implicit regularisers by revisiting backward error analysis	331
C.1 Two consecutive steps of SGD	332
C.2 Multiple steps of SGD	334
C.3 Multiple steps of full-batch gradient descent	338
C.4 Two-player games	341
D The importance of model smoothness in deep learning	343
D.1 Individual figure reproduction details	343
E Spectral Normalisation in Reinforcement learning	345
E.1 Additional experimental results	345
E.2 Experimental details	350

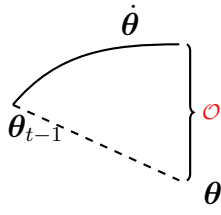
E.3	Additional hyperparameter sweeps	353
E.4	Individual figure reproduction details	356
F	Geometric Complexity: a smoothness complexity measure implicitly regularised in deep learning	358
F.1	Figure and experiments reproduction details	358
F.2	Additional experimental results	359

Appendix A

On a new continuous-time model of gradient descent

A.1 Proofs

A.1.1 BEA proof structure



$$\begin{aligned}
 \dot{\boldsymbol{\theta}} \quad \theta(h) &= \sum_{p=0}^{\infty} \frac{h^p}{p!} \boldsymbol{\theta}^{(p)} \\
 &= \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) + \underbrace{\sum_{i=2}^{n+1} h^i l_i(\boldsymbol{\theta}_{t-1})}_{\mathbf{0}} + \mathcal{O}(h^{n+2}) \\
 &= \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})
 \end{aligned}$$

Figure A.1: BEA finds continuous modified flows that describe the gradient descent update with learning rate h with an error of $\mathcal{O}(h^{n+2})$. We identify f_1, \dots, f_n so that terms of order $\mathcal{O}(h^p)$, $2 \leq p \leq n+1$ in $\boldsymbol{\theta}(h)$ are $\mathbf{0}$.

General structure. The goal of BEA is to find the functions f_1, \dots, f_n such that the flow

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E + h f_1(\boldsymbol{\theta}) + \dots + h^n f_n(\boldsymbol{\theta}) \quad (\text{A.1})$$

has an error $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}(h)\|$ of order $\mathcal{O}(h^{n+2})$ after 1 gradient descent step of learning rate h . To do so requires multiple steps (visualized in Figure A.1):

1. Expand $\boldsymbol{\theta}(h)$ via a Taylor expansion in h : $\boldsymbol{\theta}(h) = \sum_{p=0}^{\infty} \frac{h^p}{p!} \boldsymbol{\theta}^{(p)}$;
2. Expand each $\boldsymbol{\theta}^{(p)}$ up to order $\mathcal{O}(h^{n+2-p})$ as a function of f_1, \dots, f_n via the chain

rule;

3. Group together terms of the same order in h in the expansion, up to order $n + 2$.

$$\boldsymbol{\theta}(h) = \sum_{i=0}^{n+1} h^i l_i(\boldsymbol{\theta}_{t-1}) + \mathcal{O}(h^{n+2}) \quad (\text{A.2})$$

$$= \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) + \sum_{i=2}^{n+1} h^i l_i(\boldsymbol{\theta}_{t-1}) + \mathcal{O}(h^{n+2}) \quad (\text{A.3})$$

4. Compare the above update with the gradient descent update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$ and conclude that $l_i = \mathbf{0}$, $\forall i \in \{2, n + 1\}$. Use this to identify f_1, \dots, f_n .

Notation and context: all proofs below apply to general Euler updates not only gradient descent. We thus assume an update function f with the Euler step $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1})$. We can then use BEA to find the higher order correction terms describing the Euler discrete update up to a certain order, and replace $f = -\nabla_{\boldsymbol{\theta}} E$ to obtain the corresponding results for gradient descent. When we perform a Taylor expansion in h we often drop in notation the evaluation at $h = 0$ and we make that implicit.

A.1.2 Third-order flow

Theorem A.1.1 *The modified flow*

$$\dot{\boldsymbol{\theta}} = f - h \frac{1}{2} \frac{df}{d\boldsymbol{\theta}} f + h^2 \left(\frac{1}{3} \left(\frac{df}{d\boldsymbol{\theta}} \right)^2 f + \frac{1}{12} f^T \frac{df}{d^2\boldsymbol{\theta}} f \right) \quad (\text{A.4})$$

with $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_{t-1}$ follows an Euler update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1})$ with an error of $\mathcal{O}(h^4)$ after 1 gradient descent step.

Proof: We are looking for functions f_1 and f_2 such that the modified flow

$$\dot{\boldsymbol{\theta}} = f + hf_1 + h^2 f_2 \quad (\text{A.5})$$

follows the steps of GD with an error up to $\mathcal{O}(h^4)$. We perform a Taylor expansion of step size h of the above modified flow to be able to see the displacement in that time up to order $\mathcal{O}(h^4)$. We obtain (all function evaluations of f and f_i are at $\boldsymbol{\theta}_{t-1}$ which we omit for simplicity, and annotate proof steps, CR denotes Chain Rule):

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}(0) + h\dot{\boldsymbol{\theta}}(0) + \frac{1}{2}h^2\ddot{\boldsymbol{\theta}}(0) + \frac{1}{6}h^3\dddot{\boldsymbol{\theta}}(0) + \mathcal{O}(h^4) \quad (\text{A.6})$$

$$= \boldsymbol{\theta}_{t-1} + h(f + hf_1 + h^2f_2) + \frac{1}{2}h^2\frac{d}{dt}(f + hf_1 + h^2f_2) \quad (\text{A.7})$$

$$+ \frac{1}{6}h^3\ddot{\boldsymbol{\theta}} + \mathcal{O}(h^4) \quad (\text{using Eq (A.5), A.8})$$

$$= \boldsymbol{\theta}_{t-1} + hf + h^2f_1 + h^3f_2 \quad (\text{A.9})$$

$$+ \frac{1}{2}h^2\frac{d}{d\boldsymbol{\theta}}(f + hf_1 + h^2f_2)(f + hf_1 + h^2f_2) \quad (\text{A.10})$$

$$+ \frac{1}{6}h^3\ddot{\boldsymbol{\theta}} + \mathcal{O}(h^4) \quad (\text{CR, A.11})$$

$$= \boldsymbol{\theta}_{t-1} + hf + h^2f_1 + h^3f_2 + \frac{1}{2}h^2\left(\frac{df}{d\boldsymbol{\theta}}f + h\frac{df}{d\boldsymbol{\theta}}f_1 + h\frac{df_1}{d\boldsymbol{\theta}}f\right) \quad (\text{A.12})$$

$$+ \frac{1}{6}h^3\ddot{\boldsymbol{\theta}} + \mathcal{O}(h^4) \quad (\text{A.13})$$

$$= \boldsymbol{\theta}_{t-1} + hf + h^2f_1 + h^3f_2 + \frac{1}{2}h^2\left(\frac{df}{d\boldsymbol{\theta}}f + h\frac{df}{d\boldsymbol{\theta}}f_1 + h\frac{df_1}{d\boldsymbol{\theta}}f\right) \quad (\text{A.14})$$

$$+ \frac{h^3}{6}\frac{d}{dt}\left[\frac{df}{d\boldsymbol{\theta}}f + h\frac{df}{d\boldsymbol{\theta}}f_1 + h\frac{df_1}{d\boldsymbol{\theta}}f\right] + \mathcal{O}(h^4) \quad (\text{A.15})$$

$$= \boldsymbol{\theta}_{t-1} + hf + h^2f_1 + h^3f_2 + \frac{1}{2}h^2\left(\frac{df}{d\boldsymbol{\theta}}f + h\frac{df}{d\boldsymbol{\theta}}f_1 + h\frac{df_1}{d\boldsymbol{\theta}}f\right) \quad (\text{A.16})$$

$$+ \frac{1}{6}h^3\frac{d}{dt}\left[\frac{df}{d\boldsymbol{\theta}}f\right] + \mathcal{O}(h^4) \quad (\text{A.17})$$

$$= \boldsymbol{\theta}_{t-1} + hf + h^2f_1 + h^3f_2 + \frac{1}{2}h^2\left(\frac{df}{d\boldsymbol{\theta}}f + h\frac{df}{d\boldsymbol{\theta}}f_1 + h\frac{df_1}{d\boldsymbol{\theta}}f\right) \quad (\text{A.18})$$

$$+ \frac{1}{6}h^3\frac{d}{d\boldsymbol{\theta}}\left(\frac{df}{d\boldsymbol{\theta}}f\right)(f + hf_1 + h^2f_2) + \mathcal{O}(h^4) \quad (\text{CR, A.19})$$

$$= \boldsymbol{\theta}_{t-1} + hf + h^2f_1 + h^3f_2 + \frac{1}{2}h^2\left(\frac{df}{d\boldsymbol{\theta}}f + h\frac{df}{d\boldsymbol{\theta}}f_1 + h\frac{df_1}{d\boldsymbol{\theta}}f\right) \quad (\text{A.20})$$

$$+ \frac{1}{6}h^3\frac{d}{d\boldsymbol{\theta}}\left(\frac{df}{d\boldsymbol{\theta}}f\right)f + \mathcal{O}(h^4) \quad (\text{A.21})$$

If we match the terms with the Euler update, we obtain that:

$$f_1 = -\frac{1}{2}\frac{df}{d\boldsymbol{\theta}}f \quad (\text{A.22})$$

and

$$f_2 = -\frac{1}{2} \left[\frac{df}{d\boldsymbol{\theta}} f_1 + \frac{df_1}{d\boldsymbol{\theta}} f \right] - \frac{1}{6} \frac{d}{d\boldsymbol{\theta}} \left(\frac{df}{d\boldsymbol{\theta}} f \right) f \quad (\text{A.23})$$

$$f_2 = -\frac{1}{2} \left[\frac{df}{d\boldsymbol{\theta}} \left(-\frac{1}{2} \frac{df}{d\boldsymbol{\theta}} f \right) + \frac{d}{d\boldsymbol{\theta}} \left(-\frac{1}{2} \frac{df}{d\boldsymbol{\theta}} f \right) f \right] - \frac{1}{6} \frac{d}{d\boldsymbol{\theta}} \left(\frac{df}{d\boldsymbol{\theta}} f \right) f \quad (\text{A.24})$$

$$f_2 = \frac{1}{4} \left(\frac{df}{d\boldsymbol{\theta}} \right)^2 f + \frac{1}{4} \frac{d}{d\boldsymbol{\theta}} \left(\frac{df}{d\boldsymbol{\theta}} f \right) f - \frac{1}{6} \frac{d}{d\boldsymbol{\theta}} \left(\frac{df}{d\boldsymbol{\theta}} f \right) f \quad (\text{A.25})$$

$$f_2 = \frac{1}{4} \left(\frac{df}{d\boldsymbol{\theta}} \right)^2 f + \frac{1}{12} \frac{d}{d\boldsymbol{\theta}} \left(\frac{df}{d\boldsymbol{\theta}} f \right) f \quad (\text{A.26})$$

$$f_2 = \frac{1}{4} \left(\frac{df}{d\boldsymbol{\theta}} \right)^2 f + \frac{1}{12} \left(\left(\frac{df}{d\boldsymbol{\theta}} \right)^2 + f^T \frac{df}{d^2\boldsymbol{\theta}} \right) f \quad (\text{A.27})$$

$$f_2 = \frac{1}{3} \left(\frac{df}{d\boldsymbol{\theta}} \right)^2 f + \frac{1}{12} f^T \frac{df}{d^2\boldsymbol{\theta}} f \quad (\text{A.28})$$

where $f^T \frac{df}{d^2\boldsymbol{\theta}}$ is a matrix in $\mathbb{R}^{D \times D}$, with $(f^T \frac{df}{d^2\boldsymbol{\theta}})_{i,j} = \left(\frac{d \frac{df_i}{d\boldsymbol{\theta}_j}}{d\boldsymbol{\theta}} \right)^T f$ and $f^T \frac{df}{d^2\boldsymbol{\theta}} f$ is a vector in \mathbb{R}^D with $(f^T \frac{df}{d^2\boldsymbol{\theta}} f)_k = \sum_{i,j} \frac{df_k}{d\boldsymbol{\theta}_i d\boldsymbol{\theta}_j} f_i f_j$.

Replacing f_1 and f_2 in Eq (A.5) leads to the desired Eq (A.4).

□

In the case of gradient descent, $f = -\nabla_{\boldsymbol{\theta}} E$, and thus $\frac{df}{d\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}^2 E$ and we write $\frac{df}{d^2\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}^3 E \in \mathbb{R}^{D \times D \times D}$ and $(\nabla_{\boldsymbol{\theta}} E^T (\nabla_{\boldsymbol{\theta}}^3 E) \nabla_{\boldsymbol{\theta}} E)_k = \sum_{i,j} (\nabla_{\boldsymbol{\theta}} E)_i (\nabla_{\boldsymbol{\theta}}^3 E)_{i,k,j} (\nabla_{\boldsymbol{\theta}} E)_j$. Substituting in f_1 and f_2 , we obtain (Eq (3.4)):

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E - \frac{h}{2} \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E - h^2 \left(\frac{1}{3} (\nabla_{\boldsymbol{\theta}}^2 E)^2 \nabla_{\boldsymbol{\theta}} E + \frac{1}{12} \nabla_{\boldsymbol{\theta}} E^T (\nabla_{\boldsymbol{\theta}}^3 E) \nabla_{\boldsymbol{\theta}} E \right). \quad (\text{A.29})$$

A.1.3 Higher order terms of the form $\left(\frac{df}{d\boldsymbol{\theta}} \right)^n f$ or $(\nabla_{\boldsymbol{\theta}}^2 E)^n \nabla_{\boldsymbol{\theta}} E$

Theorem A.1.2 *The modified flow with an error of order $\mathcal{O}(h^{p+1})$ to the Euler update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1})$ has the form:*

$$\dot{\boldsymbol{\theta}} = \sum_{n=0}^p \frac{(-1)^n}{n+1} \left(\frac{df}{d\boldsymbol{\theta}} \right)^n f + \mathcal{C} \left(\frac{df}{d^2\boldsymbol{\theta}} \right) \quad (\text{A.30})$$

where $\mathcal{C} \left(\frac{df}{d^2\boldsymbol{\theta}} \right)$ denotes a class of functions defined as a sum of individual terms each containing higher than first order derivatives applied to f .

Proof: Consider the modified flow given by BAE:

$$\dot{\boldsymbol{\theta}} = \tilde{f} = f + hf_1 + h^2 f_2 + \cdots + h^n f_n + \dots \quad (\text{A.31})$$

For the proof we need to show that the only terms not involving higher order derivatives are of the form $(\frac{df}{d\boldsymbol{\theta}})^n f$ and we need to find the coefficients of $(\frac{df}{d\boldsymbol{\theta}})^n f$ in f_n , which we will call c_n . We already know that $c_0 = 1$, $c_1 = -\frac{1}{2}$ (Eq. (A.22)) and $c_3 = \frac{1}{3}$ (Eq. (A.28)). We use the general structure provided in Section A.1.1.

We start by expanding $\boldsymbol{\theta}(h)$ via Taylor expansion (Step 1 in the proof structure in Section A.1.1):

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}(0) + h\dot{\boldsymbol{\theta}} + \sum_{k=2}^{\infty} \frac{1}{k!} \boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}) + \sum_{n=1}^{\infty} h^{n+1} f_n + \sum_{k=2}^{\infty} \frac{h^k}{k!} \boldsymbol{\theta}^{(k)} \quad (\text{A.32})$$

We now express in $\boldsymbol{\theta}^{(k)}$ in terms of the wanted quantities f_n via Lemma A.1.3:

$$\boldsymbol{\theta}^{(k)} = \sum_{m \geq 0} h^m \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i \quad k \geq 2 \quad (\text{A.33})$$

By replacing Eq (A.33) in Eq (A.32) we obtain (Step 2 in the proof structure in Section A.1.1):

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}) + \sum_{n=1}^{\infty} h^{n+1} f_n(\boldsymbol{\theta}_{t-1}) \quad (\text{A.34})$$

$$+ \sum_{k=2}^{\infty} \frac{h^k}{k!} \sum_{m \geq 0} h^m \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \left. \frac{df_j}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_{t-1}} f_i(\boldsymbol{\theta}_{t-1}) \quad (\text{A.35})$$

$$= \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}) + \sum_{n=1}^{\infty} h^{n+1} f_n(\boldsymbol{\theta}_{t-1}) \quad (\text{A.36})$$

$$+ \sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{h^{k+m}}{k!} \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \left. \frac{df_j}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_{t-1}} f_i(\boldsymbol{\theta}_{t-1}) \quad (\text{A.37})$$

We note that due to the derivative w.r.t. t and a dependence of h in the modified vector field, there will be a dependence of h in $\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i$, since the chain

rule will expose $\frac{d\boldsymbol{\theta}}{dt}$ which depends on h (Eq (A.31)). To highlight this dependence on h we introduce the notation

$$\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i = \sum_p h^p A_{k,m}^p \quad (\text{A.38})$$

where $A_{k,m}^p : \mathbb{R}^D \rightarrow \mathbb{R}^D$ is defined to be the term of order h^p in $\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i$; thus by definition $A_{k,m}^p$ does not depend on h . Replacing the above in Eq (A.37), we conclude Step 3 outlined in the proof structure in Section A.1.1:

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}) + \sum_{n=1}^{\infty} h^{n+1} f_n(\boldsymbol{\theta}_{t-1}) + \sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{h^{k+m}}{k!} \sum_p h^p A_{k,m}^p(\boldsymbol{\theta}_{t-1}) \quad (\text{A.39})$$

Since the Euler update $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1})$ does not involve terms of order higher than h the f_n terms are obtained by cancelling the terms of order $\mathcal{O}(h^{n+1})$ in $\boldsymbol{\theta}(h)$ (Step 4 outlined in the proof structure in Section A.1.1), we have that for $k + m + p = n + 1$:

$$f_n = - \sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{1}{k!} A_{k,m}^{n+1-(k+m)} \quad (\text{A.40})$$

Since $A_{k,m}^{n+1-(k+m)}$ depends on the functions f_n , we have found a recursive relation for f_n . We now use induction to show that

$$f_n = c_n \left(\frac{df}{d\boldsymbol{\theta}} \right)^n f + \mathcal{C} \left(\frac{df}{d^2\boldsymbol{\theta}} \right) \quad (\text{A.41})$$

where $\mathcal{C} \left(\frac{df}{d^2\boldsymbol{\theta}} \right)$ denotes a class of functions defined as a sum of individual terms each containing higher than first order derivatives applied to f (Lemma A.1.4). The same proof shows the recurrence relation for c_n :

$$c_n = - \sum_{k=2}^{n+1} \frac{1}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} c_{l_2} \dots c_{l_k} \quad (\text{A.42})$$

with solution $c_n = \frac{(-1)^n}{n+1}$ (Lemma A.1.5).

We have found that each f_n can be written as $\frac{(-1)^n}{n+1} \left(\frac{df}{d\boldsymbol{\theta}}\right)^n f + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right)$. Replacing this in Eq.(A.31) concludes the proof. \square

In the case of gradient descent, where $f = -\nabla_{\boldsymbol{\theta}} E$ and $\frac{df}{d\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}^2 E$, we have that the modified flow has the form

$$\dot{\boldsymbol{\theta}} = \sum_{n=0}^p \frac{(-1)^n}{n+1} \left(-\frac{df}{d\boldsymbol{\theta}}\right)^n (-\nabla_{\boldsymbol{\theta}} E) + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right) \quad (\text{A.43})$$

$$= \sum_{n=0}^p \frac{(-1)^n}{n+1} (-\nabla_{\boldsymbol{\theta}}^2 E)^n (-\nabla_{\boldsymbol{\theta}} E) + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right) \quad (\text{A.44})$$

$$= \sum_{n=0}^p \frac{-1}{n+1} (\nabla_{\boldsymbol{\theta}}^2 E)^n (\nabla_{\boldsymbol{\theta}} E) + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right), \quad (\text{A.45})$$

which is the statement of Theorem 3.3.1 in the main thesis.

Lemma A.1.3

$$\boldsymbol{\theta}^{(k)} = \sum_{m \geq 0} h^m \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i \quad k \geq 2 \quad (\text{A.46})$$

Proof:

$$\dot{\boldsymbol{\theta}} = \tilde{f} \quad (\text{A.47})$$

$$\dot{\boldsymbol{\theta}} = \frac{d\tilde{f}}{d\boldsymbol{\theta}} \tilde{f} \quad (\text{A.48})$$

$$\boldsymbol{\theta}^{(k)} = \frac{d^{k-2}}{dt^{k-2}} \dot{\boldsymbol{\theta}} = \frac{d^{k-2}}{dt^{k-2}} \frac{d\tilde{f}}{d\boldsymbol{\theta}} \tilde{f} \quad (\text{A.49})$$

We have that

$$\frac{d\tilde{f}}{d\boldsymbol{\theta}} \tilde{f} = \left(\sum_i h^i \frac{df_i}{d\boldsymbol{\theta}} \right) \left(\sum_j h^j f_j \right) = \sum_{m \geq 0} h^m \sum_{i+j=m} \frac{df_i}{d\boldsymbol{\theta}} f_j \quad (\text{A.50})$$

which concludes the proof. We note that this is an adaptation of Lemma A.2 from [12].

We need this adaptation to avoid assuming that f is a symmetric vector field—i.e.

that $\frac{d\tilde{f}}{d\boldsymbol{\theta}}\tilde{f} = \frac{1}{2}\nabla_{\boldsymbol{\theta}}\left\|\tilde{f}\right\|^2$. □

Lemma A.1.4 *By induction we have that:*

$$f_n = c_n\left(\frac{df}{d\boldsymbol{\theta}}\right)^n f + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right) \quad (\text{A.51})$$

with

$$c_n = -\sum_{k=2}^{n+1} \frac{1}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1}c_{l_2}\dots c_{l_k} \quad (\text{A.52})$$

Proof: Base step: This is true for $n = 1, 2$ and 3 , for which we computed all terms in the BEA expansion.

Induction step: We know that (Eq (A.38)):

$$f_n = -\sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{1}{k!} A_{k,m}^{n+1-(k+m)} \quad (\text{A.53})$$

where $A_{k,m}^p$ is defined as the coefficient of h^p in:

$$\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i = \sum_p h^p A_{k,m}^p \quad (\text{A.54})$$

Under the induction hypothesis we can use that

$$f_i = c_i\left(\frac{df}{d\boldsymbol{\theta}}\right)^i f + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right) \quad \forall i < n \quad (\text{A.55})$$

And thus that the flow is of the form:

$$\dot{\boldsymbol{\theta}} = f + \sum_{i=1}^{n-1} c_i\left(\frac{df}{d\boldsymbol{\theta}}\right)^i f + \mathcal{C}\left(\frac{df}{d^2\boldsymbol{\theta}}\right) + \mathcal{O}(h^n) \quad (\text{A.56})$$

To find f_n (Eq (A.53)) we need to find $A_{k,m}^{n+1-(k+m)}$ with $k \geq 2$ and $m \geq 0$. We thus need to find $A_{k,m}^p$ with $p \leq n-1$ (since $p \leq (n+1) - (k+m)$ and $k \geq 2, m \geq 0 \implies p \leq n-1$). We do so expanding $\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\boldsymbol{\theta}} f_i$ under the induction step.

Since $i + j = m$ and $k + m \leq n + 1$ and $k \geq 2$, we have that $m \leq n - 1$ and thus $i, j \leq n - 1$ and thus we can apply the induction hypothesis for f_i and f_j .

We annotate the steps with IH (Induction Hypothesis, often by using Eq. (A.56)) and CR (Chain Rule) and S (Simplifying or grouping terms of $\mathcal{C}(\frac{df}{d^2\theta})$) by using that products and sum of terms in this function class also belong in this function class).

$$\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\theta} f_i \quad (\text{A.57})$$

$$= \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{d}{d\theta} \left(c_j \left(\frac{df}{d\theta} \right)^j f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(c_i \left(\frac{df}{d\theta} \right)^i f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \quad (\text{IH, A.58})$$

$$= \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \left(c_j \left(\frac{df}{d\theta} \right)^{j+1} + c_j f^T \frac{d}{d\theta} \left(\left(\frac{df}{d\theta} \right)^j \right) + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(c_i \left(\frac{df}{d\theta} \right)^i f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \quad (\text{CR, A.59})$$

$$= \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \left(c_j \left(\frac{df}{d\theta} \right)^{j+1} + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(c_i \left(\frac{df}{d\theta} \right)^i f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \quad (\text{S, A.60})$$

$$= \frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} c_i c_j \left(\frac{df}{d\theta} \right)^{m+1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \quad (\text{S, A.61})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \frac{d}{dt} c_i c_j \left(\frac{df}{d\theta} \right)^{m+1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \quad (\text{A.62})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \frac{d}{d\theta} \left(c_i c_j \left(\frac{df}{d\theta} \right)^{m+1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \frac{d\theta}{dt} \quad (\text{CR, A.63})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \frac{d}{d\theta} \left(c_i c_j \left(\frac{df}{d\theta} \right)^{m+1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(\sum_{l_1=0}^{n-1} h^{l_1} f_{l_1} + \mathcal{O}(h^n) \right) \quad (\text{Eq. (A.56), A.64})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \frac{d}{d\theta} \left(c_i c_j \left(\frac{df}{d\theta} \right)^{m+1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(\sum_{l_1=0}^{n-1} h^{l_1} f_{l_1} \right) + \mathcal{O}(h^n) \quad (\text{A.65})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \left(c_i c_j \left(\frac{df}{d\theta} \right)^{m+2} + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(\sum_{l_1=0}^{n-1} h^{l_1} f_{l_1} \right) + \mathcal{O}(h^n) \quad (\text{CR, S, A.66})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \left(c_i c_j \left(\frac{df}{d\theta} \right)^{m+2} + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \left(\sum_{l_1=0}^{n-1} h^{l_1} c_{l_1} \left(\frac{df}{d\theta} \right)^{l_1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) + \mathcal{O}(h^n) \quad (\text{IH, A.67})$$

Continuing by applying the induction hypothesis and the chain rule we obtain:

$$\frac{d^{k-2}}{dt^{k-2}} \sum_{i+j=m} \frac{df_j}{d\theta} f_i \quad (\text{A.68})$$

$$= \frac{d^{k-3}}{dt^{k-3}} \sum_{i+j=m} \sum_{l_1=0}^{n-1} h^{l_1} c_i c_j c_{l_1} \left(\frac{df}{d\theta} \right)^{l_1+m+2} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) + \mathcal{O}(h^n) \quad (\text{S, A.69})$$

$$= \frac{d^{k-4}}{dt^{k-4}} \sum_{i+j=m} \left(\sum_{l_1=0}^{n-1} h^{l_1} c_i c_j c_{l_1} \frac{d}{d\theta} \left(\left(\frac{df}{d\theta} \right)^{l_1+m+2} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \right) \left(\sum_{l_2=0}^{n-1} h^{l_2} f_{l_2} \right) + \mathcal{O}(h^n) \quad (\text{CR, Eq. (A.56), A.70})$$

$$= \frac{d^{k-4}}{dt^{k-4}} \sum_{i+j=m} \left(\sum_{l_1=0}^{n-1} h^{l_1} c_i c_j c_{l_1} \left(\left(\frac{df}{d\theta} \right)^{l_1+m+3} + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) \right) \quad (\text{A.71})$$

$$\left(\sum_{l_2=0}^{n-1} h^{l_2} c_{l_2} \left(\frac{df}{d\theta} \right)^{l_2} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \right) + \mathcal{O}(h^n) \quad (\text{CR, IH, A.72})$$

$$= \frac{d^{k-4}}{dt^{k-4}} \sum_{i+j=m} \sum_{l_1=0}^{n-1} \sum_{l_2=0}^{n-1} h^{l_1+l_2} c_i c_j c_{l_1} c_{l_2} \left(\frac{df}{d\theta} \right)^{l_1+l_2+m+3} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) + \mathcal{O}(h^n) \quad (\text{S, A.73})$$

$$= \sum_{i+j=m} \sum_{l_1=0}^{n-1} \sum_{l_2=0}^{n-1} \dots \sum_{l_{k-2}=0}^{n-1} h^{l_1+l_2+\dots+l_{k-2}} c_i c_j c_{l_1} c_{l_2} \dots c_{l_{k-2}} \left(\frac{df}{d\theta} \right)^{l_1+l_2+\dots+l_{k-2}+m+k-1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) + \mathcal{O}(h^n) \quad (\text{A.74})$$

$$= \sum_{i+j=m} \sum_{\substack{0 \leq l_1, \dots, l_{k-2} \leq n-1, \\ l_1+l_2+\dots+l_{k-2}=p}} h^p c_i c_j c_{l_1} c_{l_2} \dots c_{l_{k-2}} \left(\frac{df}{d\theta} \right)^{p+m+k-1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) + \mathcal{O}(h^n) \quad (\text{A.75})$$

By extracting the terms of $\mathcal{O}(h^p)$ with $p \leq n-1$ from the above we can now conclude that (note that since we are now only concerned with $p \leq n-1$ we drop the $l_i \leq n-1$ since it is implied from $l_1 + l_2 + \dots + l_{k-2} = p$):

$$A_{k,m}^p = \sum_{i+j=m} \sum_{\substack{l_1, \dots, l_{k-2} \geq 0, \\ l_1+l_2+\dots+l_{k-2}=p}} c_i c_j c_{l_1} c_{l_2} \dots c_{l_{k-2}} \left(\frac{df}{d\theta} \right)^{p+m+k-1} f + \mathcal{C} \left(\frac{df}{d^2\theta} \right) \quad (\text{A.76})$$

$$\forall p \leq n-1 \quad (\text{A.77})$$

Replacing this in f_n :

$$f_n = - \sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{1}{k!} A_{k,m}^{n+1-(k+m)} \tag{A.78}$$

$$= - \sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{1}{k!} \sum_{i+j=m} \sum_{\substack{l_1, \dots, l_{k-2} \geq 0, \\ l_1+l_2+\dots+l_{k-2}=n+1-(k+m)}} c_i c_j c_{l_1} c_{l_2} \dots c_{l_{k-2}} \left(\frac{df}{d\theta}\right)^{n+1-(k+m)+k+m-1} f + \mathcal{C}\left(\frac{df}{d^2\theta}\right) \tag{A.79}$$

$$= - \sum_{m \geq 0} \sum_{k=2}^{\infty} \frac{1}{k!} \sum_{i+j=m} \sum_{\substack{l_1, \dots, l_{k-2} \geq 0, \\ l_1+l_2+\dots+l_{k-2}=n+1-(k+m)}} c_i c_j c_{l_1} c_{l_2} \dots c_{l_{k-2}} \left(\frac{df}{d\theta}\right)^n f + \mathcal{C}\left(\frac{df}{d^2\theta}\right) \tag{A.80}$$

$$= - \sum_{k=2}^{\infty} \frac{1}{k!} \sum_{m \geq 0} \sum_{\substack{l_1, \dots, l_k \geq 0, \\ l_1+l_2+\dots+l_k=n+1-k}} c_{l_1} c_{l_2} \dots c_{l_k} \left(\frac{df}{d\theta}\right)^n f + \mathcal{C}\left(\frac{df}{d^2\theta}\right) \tag{A.81}$$

We can now say that $f_n = c_n \left(\frac{df}{d\theta}\right)^n f + \mathcal{C}\left(\frac{df}{d^2\theta}\right)$. Not only that, we have now found the recurrence relation we were seeking:

$$c_n = - \sum_{k=2}^{n+1} \frac{1}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} c_{l_2} \dots c_{l_k} \tag{A.82}$$

□

Lemma A.1.5 *The solution to the recurrence relation*

$$c_n = - \sum_{k=2}^{n+1} \frac{1}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} c_{l_2} \dots c_{l_k} \tag{A.83}$$

with $c_0 = 1$, is $c_n = \frac{(-1)^n}{n+1}$.

Proof: We will use generating functions to solve the recurrence. Let $c(x)$ be the generating function of c_n :

$$c(x) = c_0 + \sum_{n=1}^{\infty} c_n x^n \tag{A.84}$$

We have that:

$$c(x) = 1 - \sum_{n=1}^{\infty} x^n \left(\sum_{k=2}^{n+1} \frac{1}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} c_{l_2} \dots c_{l_k} \right) \quad (\text{A.85})$$

$$= 1 - \frac{1}{x} \sum_{n=1}^{\infty} x^{n+1} \left(\sum_{k=2}^{n+1} \frac{1}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} c_{l_2} \dots c_{l_k} \right) \quad (\text{A.86})$$

$$= 1 - \frac{1}{x} \sum_{n=1}^{\infty} \left(\sum_{k=2}^{n+1} \frac{x^k}{k!} x^{n+1-k} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} c_{l_2} \dots c_{l_k} \right) \quad (\text{A.87})$$

$$= 1 - \frac{1}{x} \sum_{n=1}^{\infty} \left(\sum_{k=2}^{n+1} \frac{x^k}{k!} \sum_{l_1+l_2+\dots+l_k=n-k+1} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \right) \quad (\text{A.88})$$

$$= 1 - \frac{1}{x} \sum_{k=2}^{\infty} \frac{x^k}{k!} \left(\sum_{l_1=0}^{\infty} \sum_{l_2=0}^{\infty} \sum_{l_k=0}^{\infty} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \right) \quad (\text{A.89})$$

$$= 1 - \frac{1}{x} \sum_{k=2}^{\infty} \frac{x^k}{k!} \left(\sum_{l_1=0}^{\infty} c_{l_1} x^{l_1} \right) \left(\sum_{l_2=0}^{\infty} c_{l_2} x^{l_2} \right) \dots \left(\sum_{l_k=0}^{\infty} c_{l_k} x^{l_k} \right) \quad (\text{A.90})$$

$$= 1 - \frac{1}{x} \sum_{k=2}^{\infty} \frac{x^k}{k!} c(x)^k \quad (\text{A.91})$$

$$= 1 - \frac{1}{x} (e^{xc(x)} - 1 - xc(x)) \quad (\text{A.92})$$

Where Eq (A.89) can be obtained via the following reasoning:

$$\sum_{k=2}^{\infty} \frac{x^k}{k!} \left(\sum_{l_1=0}^{\infty} \sum_{l_2=0}^{\infty} \sum_{l_k=0}^{\infty} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \right) \quad (\text{A.93})$$

$$= \sum_{k=2}^{\infty} \frac{x^k}{k!} \left(\sum_{m=0}^{\infty} \sum_{l_1+l_2+\dots+l_k=m} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \right) \quad (\text{A.94})$$

$$= \sum_{k=2}^{\infty} \frac{x^k}{k!} \left(\sum_{n=k-1}^{\infty} \sum_{l_1+l_2+\dots+l_k=n+1-k} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \right) \quad (\text{A.95})$$

$$= \sum_{k=2}^{\infty} \sum_{n=k-1}^{\infty} \frac{x^k}{k!} \sum_{l_1+l_2+\dots+l_k=n+1-k} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \quad (\text{A.96})$$

$$= \sum_{n=1}^{\infty} \sum_{k=2}^{n+1} \frac{x^k}{k!} \sum_{l_1+l_2+\dots+l_k=n+1-k} c_{l_1} x^{l_1} c_{l_2} x^{l_2} \dots c_{l_k} x^{l_k} \quad (\text{A.97})$$

Solving for $c(x)$ in Eq (A.92) we obtain:

$$c(x) = 1 - \frac{1}{x}(e^{xc(x)} - 1 - xc(x)) \quad (\text{A.98})$$

$$c(x) = 1 - \frac{1}{x}e^{xc(x)} + \frac{1}{x} + c(x) \quad (\text{A.99})$$

$$0 = 1 - \frac{1}{x}e^{xc(x)} + \frac{1}{x} \quad (\text{A.100})$$

$$e^{xc(x)} = x + 1 \quad (\text{A.101})$$

$$xc(x) = \log(x + 1) \quad (\text{A.102})$$

$$c(x) = \frac{\log(x + 1)}{x} \quad (\text{A.103})$$

We now perform the series expansion for $c(x)$:

$$c(x) = \frac{\log(x + 1)}{x} = \frac{\sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} x^n}{x} = \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} x^{n-1} = \sum_{n=0}^{\infty} \frac{(-1)^n}{n+1} x^n \quad (\text{A.104})$$

We thus have that $c_n = \frac{(-1)^n}{n+1}$ which finishes the proof.

□

A.1.4 Linearisation results around critical points

Assume we are around a critical point $\boldsymbol{\theta}^*$, i.e. $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*) = \mathbf{0}$. Via the Hartman–Grobman theorem [118, 319] we can write:

$$\nabla_{\boldsymbol{\theta}} E \approx \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*) + \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) = \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \quad (\text{A.105})$$

Replacing this in the gradient descent updates:

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) \implies \boldsymbol{\theta}_t - \boldsymbol{\theta}^* \approx (I - h \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*)) (\boldsymbol{\theta}_{t-1} - \boldsymbol{\theta}^*) \quad (\text{A.106})$$

$$\implies \boldsymbol{\theta}_t - \boldsymbol{\theta}^* \approx (I - h \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*))^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \quad (\text{A.107})$$

Since $(I - h \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*))^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) = \sum_{i=0}^{D-1} (1 - h \lambda_i^*)^t \mathbf{u}_i^* \mathbf{u}_i^{*T} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)$, gradient descent converges in the limit of $t \rightarrow \infty$ when $|1 - h \lambda_i^*| < 1$, i.e. $\lambda_i < 2/h$. Thus, we obtain the same conclusion as obtained with the PF around a local minima.

We now also show how we can use these results to derived a specialised version of the PF, one that is only valid around critical points and requires the additional assumption that the Hessian of the critical point is invertible. We start with

$$\boldsymbol{\theta}_t - \boldsymbol{\theta}^* \approx (I - h\nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*))^t (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) = \sum_{i=0}^{D-1} (1 - h\lambda_i^*)^t \mathbf{u}_i^* \mathbf{u}_i^{*T} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*). \quad (\text{A.108})$$

Under the above approximation we have:

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}^* + \sum_{i=0}^{D-1} (1 - h\lambda_i^*)^t \mathbf{u}_i^* \mathbf{u}_i^{*T} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \quad (\text{A.109})$$

which we can write in continuous form as (we use that iteration n is at time $t = nh$):

$$\boldsymbol{\theta}(t) = \boldsymbol{\theta}^* + \sum_{i=0}^{D-1} (1 - h\lambda_i^*)^{t/h} \mathbf{u}_i^* \mathbf{u}_i^{*T} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*) \quad (\text{A.110})$$

$$= \boldsymbol{\theta}^* + \sum_{i=0}^{D-1} e^{\log(1-h\lambda_i^*)t/h} \mathbf{u}_i^* \mathbf{u}_i^{*T} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*). \quad (\text{A.111})$$

Taking the derivative wrt to t :

$$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i^*)}{h} \underbrace{e^{\log(1-h\lambda_i^*)t/h} \mathbf{u}_i^* \mathbf{u}_i^{*T} (\boldsymbol{\theta}_0 - \boldsymbol{\theta}^*)}_{\text{Eq. (A.109)}} = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i^*)}{h} (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \quad (\text{A.112})$$

From here, if we assume the Hessian at the critical point is invertible we can write from the approximation above (Eq. (A.105)) that $\boldsymbol{\theta} - \boldsymbol{\theta}^* = (\nabla_{\boldsymbol{\theta}}^2 E)^{-1} \nabla_{\boldsymbol{\theta}} E = \sum_{i=0}^{D-1} \frac{1}{\lambda_i^*} \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i^* \mathbf{u}_i^{*T}$ and replacing this in the above we the PF:

$$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i^*)}{h\lambda_i^*} \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i^* \mathbf{u}_i^{*T}. \quad (\text{A.113})$$

A.1.5 The solution of the PF for quadratic losses

We derive the solution of the PF in the case of quadratic losses - Remark 3.3.1 in the main thesis. We have $E(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^t \mathbf{A} \boldsymbol{\theta} + \mathbf{b}^T \boldsymbol{\theta}$ with \mathbf{A} symmetric. $\nabla_{\boldsymbol{\theta}} E = \mathbf{A} \boldsymbol{\theta} + \mathbf{b}$

and $\nabla_{\boldsymbol{\theta}}^2 E = \mathbf{A}$. Thus λ_i and \mathbf{u}_i are the eigenvalues and eigenvectors of \mathbf{A} and do not change based on $\boldsymbol{\theta}$. Replacing this in the PF leads to:

$$\dot{\boldsymbol{\theta}} = \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} (\mathbf{A}\boldsymbol{\theta} + b)^T \mathbf{u}_i \mathbf{u}_i \quad (\text{A.114})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} (\mathbf{A}\boldsymbol{\theta})^T \mathbf{u}_i \mathbf{u}_i + \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i \mathbf{u}_i \quad (\text{A.115})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} \boldsymbol{\theta}^T \mathbf{A} \mathbf{u}_i \mathbf{u}_i + \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i \mathbf{u}_i \quad (\text{A.116})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} \lambda_i \boldsymbol{\theta}^T \mathbf{u}_i \mathbf{u}_i + \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i \mathbf{u}_i \quad (\text{A.117})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h} \boldsymbol{\theta}^T \mathbf{u}_i \mathbf{u}_i + \sum_{i=0}^{D-1} \frac{\log(1 - h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i \mathbf{u}_i \quad (\text{A.118})$$

From here

$$(\boldsymbol{\theta}^T \mathbf{u}_i) = \frac{\log(1 - h\lambda_i)}{h} \boldsymbol{\theta}^T \mathbf{u}_i + \frac{\log(1 - h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i. \quad (\text{A.119})$$

The solution of this flow is $(\boldsymbol{\theta}^T \mathbf{u}_i)(t) = e^{\frac{\log(1-h\lambda_i)}{h}t} \boldsymbol{\theta}_0^T \mathbf{u}_i + t \frac{\log(1-h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i$.

Since \mathbf{u}_i form a basis of \mathbb{R}^D , we can now write

$$\boldsymbol{\theta}(t) = \sum_{i=0}^{D-1} \boldsymbol{\theta}(t)^T \mathbf{u}_i \mathbf{u}_i \quad (\text{A.120})$$

$$= \sum_{i=0}^{D-1} \left(e^{\frac{\log(1-h\lambda_i)}{h}t} \boldsymbol{\theta}_0^T \mathbf{u}_i + t \frac{\log(1 - h\lambda_i)}{h\lambda_i} b^T \mathbf{u}_i \right) \mathbf{u}_i. \quad (\text{A.121})$$

We thus have obtained the solution described in Remark 3.3.1 in the main text of the thesis.

A.1.6 The Jacobian of the PF at critical points

We compute Jacobian of the PF at a critical point $\boldsymbol{\theta}^*$ with the eigenvalues and eigenvectors of $\nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*)$ denoted as λ_i^* and \mathbf{u}_i^* . We will use that $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*) = \mathbf{0}$.

$$\mathbf{J}_{\text{PF}}(\boldsymbol{\theta}^*) = \sum_{i=0}^{D-1} \frac{d(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)^{\frac{\log(1-h\lambda_i)}{h\lambda_i}} \mathbf{u}_i}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \quad (\text{A.122})$$

$$= \sum_{i=0}^{D-1} \frac{d(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)^{\frac{\log(1-h\lambda_i)}{h\lambda_i}}}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \mathbf{u}_i^{*T} \quad (\text{A.123})$$

$$+ \sum_{i=0}^{D-1} \underbrace{\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*)^T}_{\mathbf{0}} \mathbf{u}_i^* \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \frac{d\mathbf{u}_i}{d\boldsymbol{\theta}} \quad (\text{A.124})$$

$$= \sum_{i=0}^{D-1} \frac{d(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)^{\frac{\log(1-h\lambda_i)}{h\lambda_i}}}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \mathbf{u}_i^{*T} \quad (\text{A.125})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \frac{d(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \mathbf{u}_i^{*T} \quad (\text{A.126})$$

$$+ \sum_{i=0}^{D-1} \underbrace{(\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*)^T \mathbf{u}_i^*)}_{\mathbf{0}} \frac{d^{\frac{\log(1-h\lambda_i)}{h\lambda_i}}}{d\boldsymbol{\theta}} \mathbf{u}_i^{*T} \quad (\text{A.127})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \frac{d(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i)}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \mathbf{u}_i^{*T} \quad (\text{A.128})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \left(\frac{d\nabla_{\boldsymbol{\theta}} E}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \right)^T \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.129})$$

$$+ \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \left(\frac{d\mathbf{u}_i}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \right)^T \underbrace{\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*)}_{\mathbf{0}} \mathbf{u}_i^{*T} \quad (\text{A.130})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*) \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.131})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h\lambda_i^*} \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.132})$$

$$= \sum_{i=0}^{D-1} \frac{\log(1-h\lambda_i^*)}{h} \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.133})$$

We thus arrived at the result used in Equation 3.13 of the thesis, where we used to perform stability analysis on the PF. This results shows that the eigenvalues of the Jacobian for the PF are $\frac{\log(1-h\lambda_i^*)}{h}$, thus local minima where $\lambda_i^* \geq 0$ are only attractive if $h\lambda_i^* < 2$.

A.1.7 Jacobians of the IGR flow and NGF at critical points

For the NGF, we have

$$\mathbf{J}_{\text{NGF}}(\boldsymbol{\theta}^*) = -\nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*) = -\sum_{i=0}^{D-1} \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.134})$$

and thus the eigenvalues of the Jacobian for the NGF are $-\lambda_i^*$. Thus local minima where $\lambda_i^* > 0 \implies -\lambda_i^* < 0$ are attractive for the NGF.

For the IGR flow, we have

$$\mathbf{J}_{\text{IGR}}(\boldsymbol{\theta}^*) = -\frac{d\left(\nabla_{\boldsymbol{\theta}} E + \frac{h^2}{2} \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E\right)}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \quad (\text{A.135})$$

$$= -\frac{d\nabla_{\boldsymbol{\theta}} E}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) - \frac{h^2}{2} \frac{d\nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E}{d\boldsymbol{\theta}}(\boldsymbol{\theta}^*) \quad (\text{A.136})$$

$$= -\sum_{i=0}^{D-1} \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*T} - \frac{h^2}{2} \left(\nabla_{\boldsymbol{\theta}} E^T \nabla_{\boldsymbol{\theta}}^3 E + \nabla_{\boldsymbol{\theta}}^2 E \frac{d\nabla_{\boldsymbol{\theta}} E}{d\boldsymbol{\theta}} \right)(\boldsymbol{\theta}^*) \quad (\text{A.137})$$

$$= -\sum_{i=0}^{D-1} \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*T} - \frac{h^2}{2} \underbrace{\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}^*)^T}_{\mathbf{0}} \nabla_{\boldsymbol{\theta}}^3 E(\boldsymbol{\theta}^*) - \frac{h^2}{2} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*) \quad (\text{A.138})$$

$$= -\sum_{i=0}^{D-1} \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*T} - \frac{h^2}{2} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}^*) \quad (\text{A.139})$$

$$= -\sum_{i=0}^{D-1} \lambda_i^* \mathbf{u}_i^* \mathbf{u}_i^{*T} - \sum_{i=0}^{D-1} \frac{h^2}{2} \lambda_i^{*2} \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.140})$$

$$= -\sum_{i=0}^{D-1} \left(\lambda_i^* + \frac{h^2}{2} \lambda_i^{*2} \right) \mathbf{u}_i^* \mathbf{u}_i^{*T} \quad (\text{A.141})$$

and thus the eigenvalues of the Jacobian for the IGR flow are $-\left(\lambda_i^* + \frac{h^2}{2} \lambda_i^{*2}\right)$. Thus local minima where $\lambda_i^* > 0 \implies -\left(\lambda_i^* + \frac{h^2}{2} \lambda_i^{*2}\right) < 0$ are attractive for the IGR flow.

A.1.8 Multiple gradient descent steps

We now show that if we use BEA to construct a flow that has an error of order $\mathcal{O}(h^p)$ after 1 gradient descent step, then the error will be of the same order after 2 steps. The same argument can be applied to multiple steps, though the error scales proportionally to the number of steps and the bound becomes vacuous as the number

of steps increases. We have $\dot{\boldsymbol{\theta}}$ with $\boldsymbol{\theta}(0) = \boldsymbol{\theta}_{t-1}$ tracks $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ with an error $\mathcal{O}(h^p)$ (i.e. $\|\boldsymbol{\theta}(h) - \boldsymbol{\theta}_t\|$ is $\mathcal{O}(h^p)$), then $\dot{\boldsymbol{\theta}}$ tracks two steps of gradient descent $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1})$ and $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_t)$ with an error of the same order $\mathcal{O}(h^p)$ (i.e. $\|\boldsymbol{\theta}(2h) - \boldsymbol{\theta}_{t+1}\|$ is $\mathcal{O}(h^p)$). We prove this below. For the purpose of this proof we use the following notation $\boldsymbol{\theta}(h; \boldsymbol{\theta}')$ is the value of the flow at time h with initial condition $\boldsymbol{\theta}'$. Making the initial condition explicit is necessary for the proof.

$$\|\boldsymbol{\theta}(2h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}_{t+1}\| = \|\boldsymbol{\theta}(2h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}(h; \boldsymbol{\theta}_t) + \boldsymbol{\theta}(h; \boldsymbol{\theta}_t) - \boldsymbol{\theta}_{t+1}\| \quad (\text{A.142})$$

$$\leq \|\boldsymbol{\theta}(2h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}(h; \boldsymbol{\theta}_t)\| + \|\boldsymbol{\theta}(h; \boldsymbol{\theta}_t) - \boldsymbol{\theta}_{t+1}\| \quad (\text{A.143})$$

$$\leq \|\boldsymbol{\theta}(2h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}(h; \boldsymbol{\theta}_t)\| + \mathcal{O}(h^p) \quad (\text{A.144})$$

$$= \|\boldsymbol{\theta}(h; \boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})) - \boldsymbol{\theta}(h; \boldsymbol{\theta}_t)\| + \mathcal{O}(h^p) \quad (\text{A.145})$$

We thus have to bound how the flow changes after time h when starting with two different initial conditions $\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})$ and $\boldsymbol{\theta}_t$. We also know that $\|\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}_t\|$ is $\mathcal{O}(h^p)$. Thus by expanding the Taylor series and the mean value theorem we obtain:

$$\|\boldsymbol{\theta}(h; \boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})) - \boldsymbol{\theta}(h; \boldsymbol{\theta}_t)\| \quad (\text{A.146})$$

$$= \left\| \sum_{i=0}^{\infty} \frac{h^i}{i!} \boldsymbol{\theta}^{(i)}(\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})) - \sum_{i=0}^{\infty} \frac{h^i}{i!} \boldsymbol{\theta}^{(i)}(\boldsymbol{\theta}_t) \right\| \quad (\text{A.147})$$

$$= \left\| \sum_{i=0}^{\infty} \frac{h^i}{i!} (\boldsymbol{\theta}^{(i)}(\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})) - \boldsymbol{\theta}^{(i)}(\boldsymbol{\theta}_t)) \right\| \quad (\text{A.148})$$

$$= \left\| \sum_{i=0}^{\infty} \frac{h^i}{i!} \frac{d}{d\boldsymbol{\theta}} \boldsymbol{\theta}^{(i)}(\boldsymbol{\theta}') (\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}_t) \right\| \quad (\text{A.149})$$

$$\leq \sum_{i=0}^{\infty} \frac{h^i}{i!} \left| \frac{d}{d\boldsymbol{\theta}} \boldsymbol{\theta}^{(i)}(\boldsymbol{\theta}') \right| \underbrace{\|(\boldsymbol{\theta}(h; \boldsymbol{\theta}_{t-1})) - \boldsymbol{\theta}_t\|}_{\mathcal{O}(h^p)} = \mathcal{O}(h^p) \quad (\text{A.150})$$

This tells us that we can construct a bound:

$$\|\boldsymbol{\theta}(2h; \boldsymbol{\theta}_{t-1}) - \boldsymbol{\theta}_{t+1}\| \leq \mathcal{O}(h^p) + \mathcal{O}(h^p) = \mathcal{O}(h^p) \quad (\text{A.151})$$

Dependence on the number of steps. While the order in learning rate is the same, we note that as the number of steps increases the errors are likely to

accumulate with the number of steps (for n discrete steps we will sum n terms of order $\mathcal{O}(h^p)$ in the above bound). For example when taking the number of steps $n \rightarrow \infty$ the above no longer provides a bound.

A.1.9 Approximations to per-iteration drift for gradient descent and momentum

We now prove Thm 3.7.1, on the per-iteration drift of gradient descent. We apply the Taylor remainder theorem on the NGF initialised at gradient descent iteration parameters $\boldsymbol{\theta}_{t-1}$. We have that there exists h' such that

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}(0) + h\dot{\boldsymbol{\theta}}(0) + \frac{h^2}{2}\ddot{\boldsymbol{\theta}}(h') \quad (\text{A.152})$$

$$= \boldsymbol{\theta}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) + \frac{h^2}{2}\nabla_{\boldsymbol{\theta}'}^2E(\boldsymbol{\theta}')\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}') \quad (\text{A.153})$$

with $h' \in (0, h)$. The above proof measures the drift of gradient descent, which we then used to construct DAL.

In Section 3.8, we also used DAL for momentum, and used an intuitive justification. We now show a more theoretical justification, by measuring the total per-iteration drift of momentum. To do so, we use the following flow to describe momentum, provided by Kunin et al. [65]

$$\dot{\boldsymbol{\theta}} = -\frac{1}{1-\beta}\nabla_{\boldsymbol{\theta}}E, \quad (\text{A.154})$$

and we measure the discretisation drift of a momentum update compared to this flow. If we apply the same strategy as above, we have that

$$\boldsymbol{\theta}(h) = \boldsymbol{\theta}(0) + h\dot{\boldsymbol{\theta}}(0) + \frac{h^2}{2}\ddot{\boldsymbol{\theta}}(h') \quad (\text{A.155})$$

$$= \boldsymbol{\theta}_{t-1} - \frac{h}{1-\beta}\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) + \frac{h^2}{2(1-\beta)^2}\nabla_{\boldsymbol{\theta}'}^2E(\boldsymbol{\theta}')\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}') \quad (\text{A.156})$$

When we compare this with the discrete update

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + \beta\mathbf{v}_{t-1} - h\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}) \quad (\text{A.157})$$

and obtain that the per-iteration drift is

$$\left(\boldsymbol{\theta}_{t-1} - \frac{h}{1-\beta} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) + \frac{h^2}{2(1-\beta)^2} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}') \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}') \right) \quad (\text{A.158})$$

$$- (\boldsymbol{\theta}_{t-1} + \beta \mathbf{v}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})) \quad (\text{A.159})$$

$$= \frac{h^2}{2(1-\beta)^2} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}') \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}') + \beta \mathbf{v}_{t-1} + h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) - \frac{h}{1-\beta} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) \quad (\text{A.160})$$

$$= \frac{h^2}{2(1-\beta)^2} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}') \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}') + \beta \mathbf{v}_{t-1} - \frac{h\beta}{1-\beta} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}) \quad (\text{A.161})$$

Thus there is one part of the drift that comes from the same source as gradient descent and another part of the drift which comes from the alignment of the current gradient and the moving average obtained using momentum. Thus when using DAL- p with momentum, we only focus on one of the sources of drift (the first term).

A.2 Comparison with a discrete-time approach

A.2.1 Changes in loss function

We aim to obtain similar intuition to what we have obtained from the PF by discretising the NGF using Euler steps. We have:

$$E(\boldsymbol{\theta}_{t+1}) \approx E(\boldsymbol{\theta}_t) - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) + \frac{h^2}{2} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}_t) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) \quad (\text{A.162})$$

$$\approx E(\boldsymbol{\theta}_t) - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) + \frac{h^2}{2} \sum_i \lambda_i (\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \mathbf{u}_i)^2 \quad (\text{A.163})$$

$$\approx E(\boldsymbol{\theta}_t) - h \sum_i (\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \mathbf{u}_i)^2 + \frac{h^2}{2} \sum_i \lambda_i (\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \mathbf{u}_i)^2 \quad (\text{A.164})$$

$$\approx E(\boldsymbol{\theta}_t) + \sum_i (1 - h/(2\lambda_i)) (\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T \mathbf{u}_i)^2 \quad (\text{A.165})$$

thus under this approximation loss to decrease between iterations one requires $1 - h/(2\lambda_i) \leq 0$. This is consistent with the observations of the loss obtained using the PF in Eq (3.12).

A.2.2 The dynamics of $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i$

We now show to obtain the approximated dynamics of $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i$ obtained from the PF in Section 3.5 using a discrete-time approach.

We have that:

$$\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t+1}) \approx \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)) \quad (\text{A.166})$$

$$\approx \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) - h \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}_t) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) \quad (\text{A.167})$$

If we assume that the Hessian eigenvectors do not change between iterations, we have:

$$\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t+1})^T (\mathbf{u}_i)_{t+1} \approx \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T (\mathbf{u}_i)_{t+1} - h \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}_t) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) (\mathbf{u}_i)_{t+1} \quad (\text{A.168})$$

$$= (1 - h \lambda_i) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t)^T (\mathbf{u}_i)_{t+1} \quad (\text{A.169})$$

and obtain the same behaviour predicted in Section 3.5.

A.2.3 The connection between DAL and Taylor expansion optimal learning rate

In Section 3.7.2 introduced DAL as a way to control the drift of gradient descent; to do so, we set the learning rate of gradient descent to be inverse proportional to $\|\nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E\| / \|\nabla_{\boldsymbol{\theta}} E\| = \|\nabla_{\boldsymbol{\theta}}^2 E \hat{\mathbf{g}}\|$.

The learning rate set by DAL is similar, but distinct to that obtained by using a Taylor expansion of the loss E , as follows:

$$E(\boldsymbol{\theta} - h \nabla_{\boldsymbol{\theta}} E) = E(\boldsymbol{\theta}) - h \nabla_{\boldsymbol{\theta}} E^T \nabla_{\boldsymbol{\theta}} E + \frac{h^2}{2} \nabla_{\boldsymbol{\theta}} E^T \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E + \mathcal{O}(h^3) \quad (\text{A.170})$$

solving for the optimal h in the above leads to:

$$h = \frac{\nabla_{\boldsymbol{\theta}} E^T \nabla_{\boldsymbol{\theta}} E}{\nabla_{\boldsymbol{\theta}} E^T \nabla_{\boldsymbol{\theta}}^2 E \nabla_{\boldsymbol{\theta}} E} = \frac{1}{\hat{\mathbf{g}}^T \nabla_{\boldsymbol{\theta}}^2 E \hat{\mathbf{g}}} \quad (\text{A.171})$$

While this learning rate contains similarities to DAL, since $\hat{\mathbf{g}}^T \nabla_{\boldsymbol{\theta}}^2 E \hat{\mathbf{g}} = \sum_i \lambda_i (\hat{\mathbf{g}}^T \mathbf{u}_i)^2$ and $\|\nabla_{\boldsymbol{\theta}}^2 E \hat{\mathbf{g}}\| = |\sum_i \lambda_i \hat{\mathbf{g}}^T \mathbf{u}_i|$, we also note a few significant differences:

- Since DAL is derived from discretisation drift, we can use it in a variety of settings. We have shown in Section 3.8 how DAL can be adapted to be used with momentum, where the main idea is to increase the contribution of local gradients to the momentum moving average when the drift is large, and decrease it otherwise.
- While we focused on using one learning rate for all parameters in the main thesis, and thus use the norm operator in DAL, this is not necessary. Instead we can use the *per-parameter drift* to construct per-parameter adaptive learning rates. We show preliminary results in Figure 3.30.
- The aim of DAL is to control the discretisation drift of gradient descent and show that this can affect both the stability and performance of gradient descent, through DAL-p. Beyond providing a training tool, DAL-p is another tool to use to understand gradient descent. Importantly, the ability to control the drift also enhances the generalisation capabilities of the method.

A.3 Experimental details

Estimating continuous-time trajectories of flows. To estimate the continuous-time trajectories we use Euler integration with a small learning rate 5×10^{-5} . We reached this learning rate through experimentation: further decreasing it did not change the obtained results. We note that this approach can be computationally expensive: to estimate the trajectory of the NGF of time corresponding to one gradient descent step with learning rate 10^{-2} , we need to do 5000 gradient steps. It is common in the literature to use Runge–Kutta4 to approximate continuous-time flows, but we noticed that approximating a flow for time h using Runge–Kutta4 with learning rate h still introduced significant drift: if the learning rate was further reduced and multiple steps were taken the results were significantly different. Thus Runge–Kutta4 also needs multiple steps to estimate continuous trajectories. Given that one Runge–Kutta4 update requires computing 4 gradients, we found that using Euler integration with small learning rates is both sufficient and more efficient in practice.

Datasets. When training neural networks with primarily used three standard datasets: MNIST [320], CIFAR-10 [321] and Imagenet [322]. On the small NN example in Section 3.3, we used a dataset of 5 examples, where the input is 2D sampled from a Gaussian distribution and the regression targets are also sample randomly. We used the UCI breast cancer dataset [140] in Figure A.7.

Architectures. We use standard architectures: MLPs for MNIST, VGG [209] or Resnet-18 (Version 1) [86] for CIFAR-10, Resnet-50 for Imagenet (Version 1). We do not use any form of regularisation or early stopping. We use the Elu activation function [323] to ensure that the theoretical setting we discussed applies directly (we thus avoid discontinuities caused by ReLus [301]). We note that in the CIFAR-10 experiments, we did not adapt the Resnet-18 architecture to the dataset; doing so will likely increase performance.

Losses. Unless otherwise specified we used a cross entropy loss.

Computing eigenvalues and eigenvectors. We use the Lanczos algorithm to compute λ_i and \mathbf{u}_i .

Seeds. All test accuracies are shown averaged from 3 seeds. For training curves, we compare models across individual training runs to be able to observe the behaviour of gradient descent. We did not observe variability across seeds with any of the behaviours reported in this thesis.

DAL. When using DAL we set a maximum learning rate of 5 to avoid any potential instabilities. We did not experiment with other values.

Computing $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$. We use Hessian vector products to compute $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$. We experimented with using the approximation $\nabla_{\theta}^2 E \nabla_{\theta} E = \frac{1}{2} \nabla_{\theta} \|\nabla_{\theta} E\|^2 \approx \frac{\nabla_{\theta} E(\theta + \epsilon \nabla_{\theta} E) - \nabla_{\theta} E(\theta)}{\epsilon}$ with $\epsilon = 0.01 / \|\nabla_{\theta} E\|$ as suggested by Geiping et al. [151] and saw no decrease in performance in the full-batch setting when using it in DAL, but a slight decrease in performance when using it in stochastic gradient descent. We show experimental results in Figures A.2, A.3, A.4 and A.5. More experimentation is needed to see how to best leverage this approximation, either by trying other values of ϵ or using larger batches to approximate $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$. We note however that all the conclusions we observed in the thesis regarding the relative performance of p

values of DAL still hold, with and without the approximation and that indeed we see less of a difference when lower values of p are used.

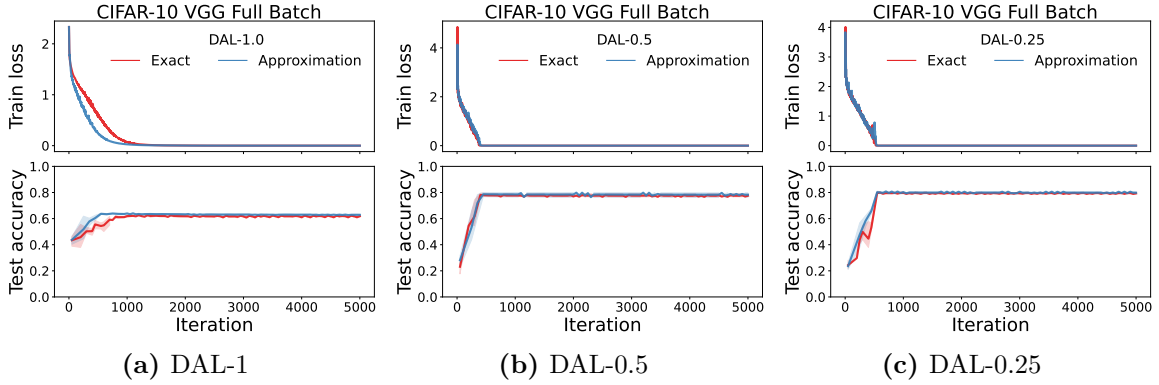


Figure A.2: CIFAR-10 DAL results with a Hessian vector product computation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ compared to an approximation. In this case, we observe no difference in performance between the two approaches.

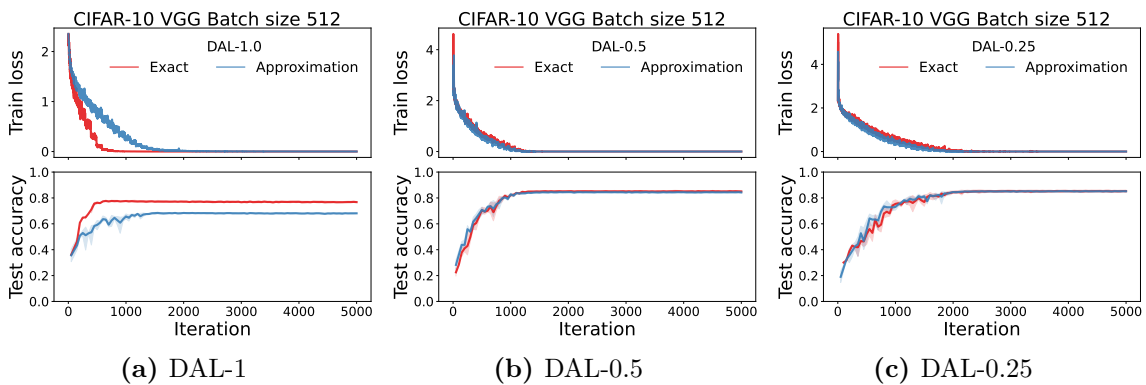


Figure A.3: CIFAR-10 DAL results with a Hessian vector product computation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ compared to an approximation. For DAL-1, using the approximation leads to a decrease in test accuracy, but the same is not observed for DAL-0.5 and DAL-0.25.

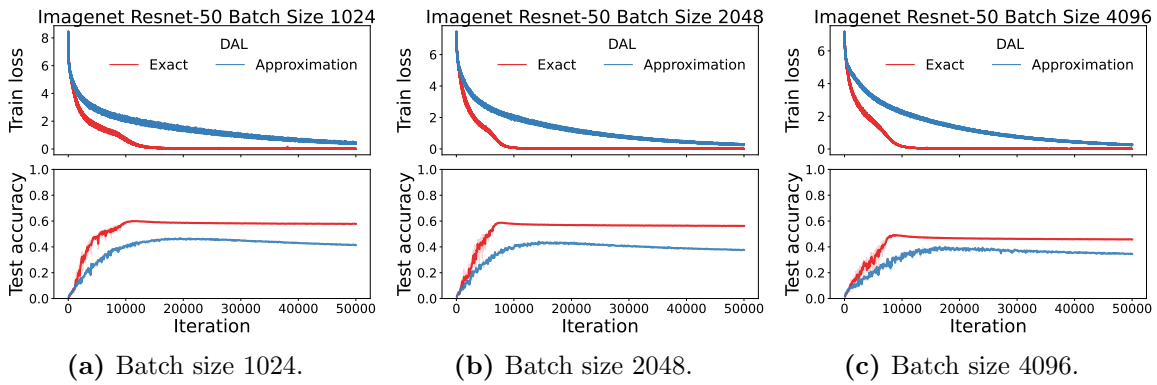


Figure A.4: Imagenet DAL results with a Hessian vector product computation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ compared to an approximation. Using the approximation results into a significant decrease in performance; we suspect that this can be resolved by using another value for ϵ .

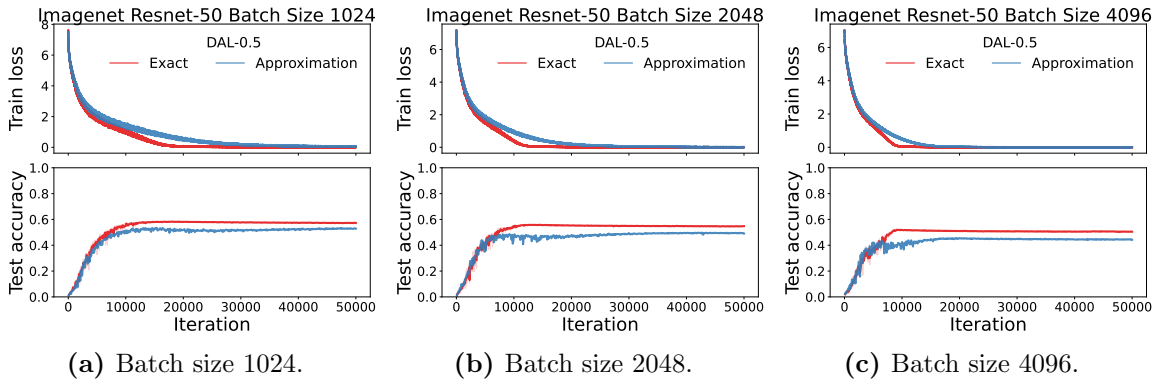


Figure A.5: Imagenet DAL-0.5 results with a Hessian vector product computation of $\nabla_{\theta}^2 E \nabla_{\theta} E$ compared to an approximation. Using the approximation results into a significant decrease in performance; we suspect that this can be resolved by using another value for ϵ .

A.4 Additional experimental results

Oscillations for non linear functions. We show that in the case of non-linear functions in Figure A.6, using \cos . Here too, the PF is better at describing the behaviour of GD than the NGF and IGR flow, which stop at local minima.

Global error for the UCI breast cancer dataset. We show in Figure A.7 that the PF is better at tracking the behaviour of GD in the case of NNs both in the stability and edge of stability areas.

Additional edge of stability experiments. We show that one dimension is sufficient to cause instability in training on MNIST in Figure A.9 (as we have done in

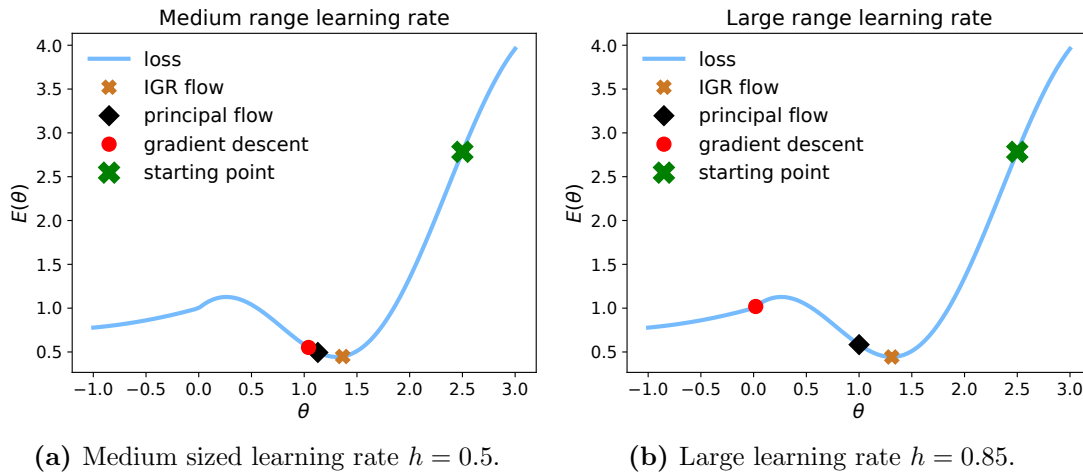


Figure A.6: Results with a non quadratic function, including cos and sin. The function used is: $E(\theta) = \cos(\theta) + \theta$, if $\theta < 0$; $(\theta/3)^2 2 + 1 + \theta/3$ otherwise. The principal flow and the IGR flow stay on the cosine branch, but the PF is able to exit the local minimum. The main reason for using this function was to show the escape of gradient descent into a flatter valley and to assess what the corresponding flows do.

Figure 3.17 for CIFAR-10). We show the connection between the stability coefficient sc_0 and areas where the loss increases in Figure A.10: this shows that when the loss increases, it is in an area where sc_0 is also increasing. We note that to get a full picture of the behaviour of the PF we would need to use continuous-time computation, but since that is computationally prohibitive we use the incomplete but easily available discrete data. Figure A.11 zooms in the behaviour of λ_0 as dependent of the behaviour of the corresponding flows on MNIST and CIFAR-10 and Figure A.12 shows additional Resnet-18 results.

DAL. We use gradient descent training with a fixed learning rate to measure the quantities we would like to use as a learning rate in DAL, to see if they have a reasonable range and show results in Figure A.13. We show sweeps across batch sizes Figures A.14, A.16, A.17. We show DAL- p sweeps with effective learning rates, training losses and test accuracies in Figure A.18. We show results with DAL on a mean square loss in Figure A.19.

DAL learned landscapes. To investigate the landscapes learned by DAL- p , we use the method of Li et al. [22] and compare against the landscapes learned using

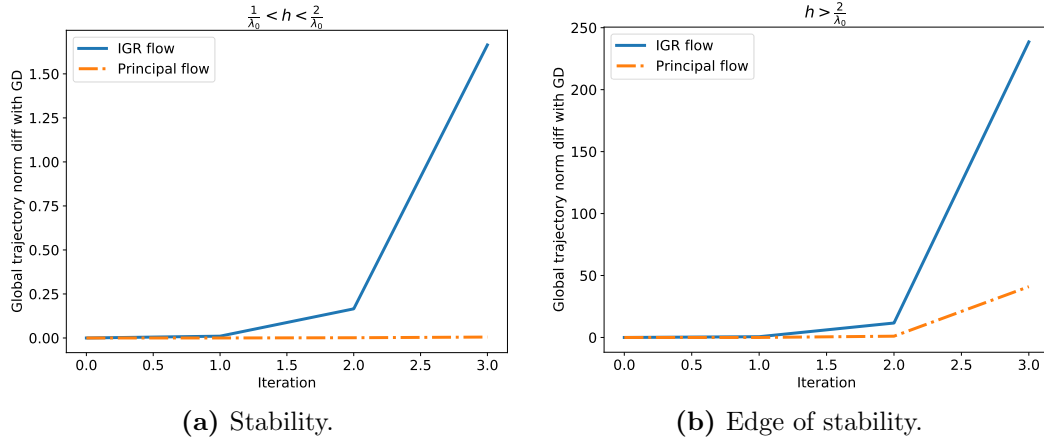


Figure A.7: Global parameter error $\|\theta(nh) - \theta_n\|$ for IGR and PF flows on an MLP trained on the UCI breast cancer dataset. We initialise both gradient descent and the IGR and PF flows at a set of initial parameters, and run them for time step $2h$ and compare the error after each time step h . The PF is better at tracking the trajectory of GD than the IGR flow by a significant margin.

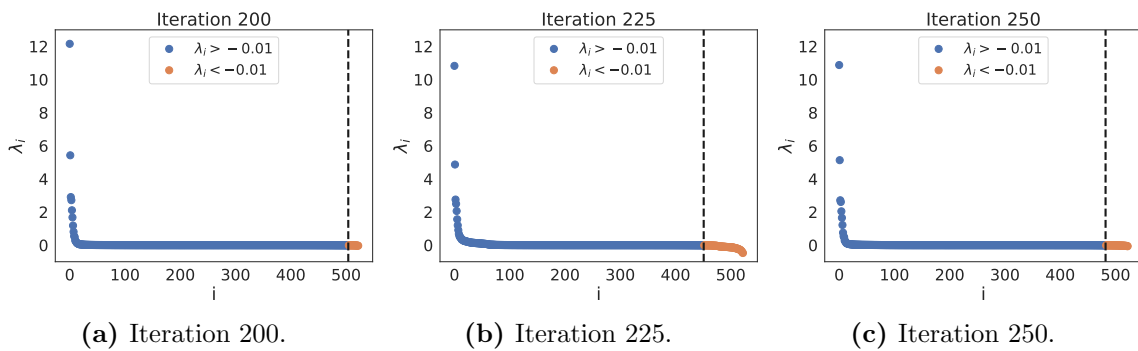


Figure A.8: The eigenspectrum obtained along the gradient descent path of a 5 layer small MLP trained on the UCI Iris dataset. The Hessian eigenvalues are largely positive. The result accompanies Figure 3.13 in the main thesis.

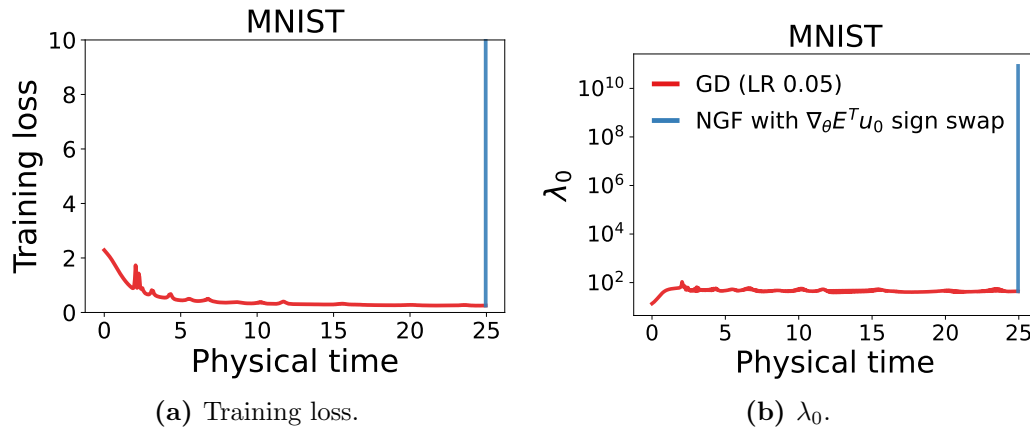
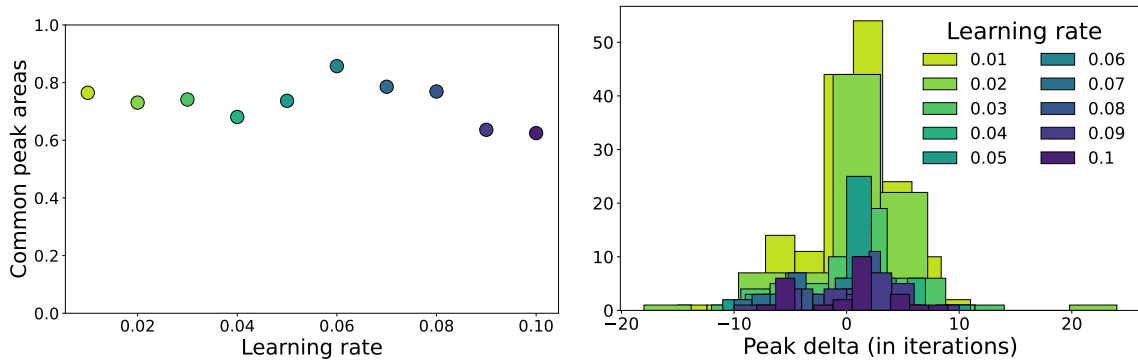


Figure A.9: One eigendirection is sufficient to lead to instabilities. To create a situation similar to that of the PF in neural network training, we construct a flow given by the NGF in all eigendirections but \mathbf{u}_0 ; in the direction of \mathbf{u}_0 , we change the sign of the flow’s vector field. This leads to the flow $\dot{\boldsymbol{\theta}} = (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0) \mathbf{u}_0 + \sum_{i=1}^{D-1} -(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$. We show this flow can be very unstable when initialised in an edge of stability area, with increases in loss (a) and λ_0 (b). MNIST results when changing one direction in continuous-time: one eigendirection is sufficient to cause instability. The model was trained with learning rate 0.05 until the edge of stability area is reached after which the flow $\dot{\boldsymbol{\theta}} = (\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0) \mathbf{u}_0 + \sum_{i=1}^{D-1} -(\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i) \mathbf{u}_i$ is used.

SGD. We have shown results with batch size 64 on CIFAR-10 in Figure 3.28 in the main text. Here present additional results in Figures A.20, A.22. Across datasets and batch sizes, we consistently observe that DAL- p learns flatter landscapes. We also consistently observe that during training, λ_0 is smaller with DAL- p than with SGD, as shown in Figure A.21.

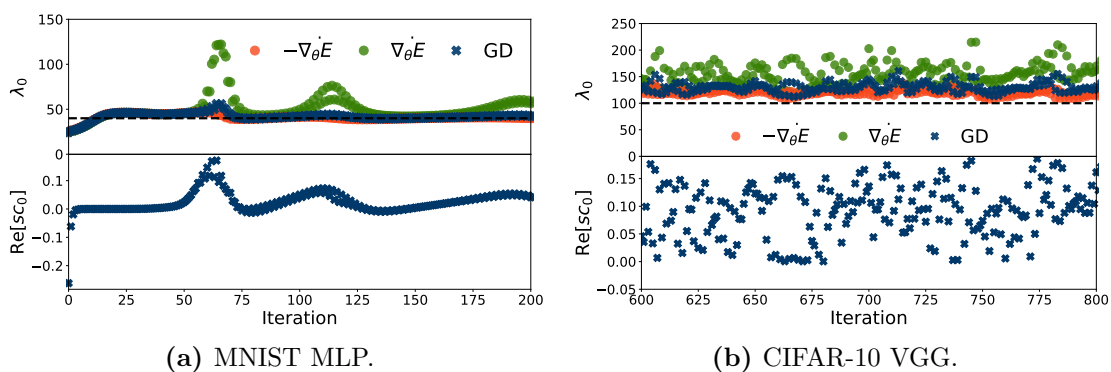
Momentum. We show additional results with learning rate adaptation and momentum in Figure A.24.

Global discretisation error. We show the global error in trajectory between the NGF and gradient descent in Figure A.25. As previously observed [52], gradient descent follows the NGF early in training, and the eigenvalues of the trajectory following the NGF keeps growing.



(a) Proportion of loss function local peaks that are also stability coefficient local peaks. (b) Delta in peak areas of the loss and the stability coefficient counted in iterations.

Figure A.10: Peak areas of loss increase and the stability coefficient for the largest eigendirection of the PF. We find peaks in our training signals via `scipy.signal.find_peaks`. Here, we normalise the coefficient by the gradient norm to remove the gradient norm as a potential confounder. Results on CIFAR-10. (a) shows that most loss function peaks are 5 iterations away from a local stability coefficient peak. (b) shows the delta in iterations between a loss function peak and the closes stability coefficient peak.



(a) MNIST MLP.

(b) CIFAR-10 VGG.

Figure A.11: Understanding changes to λ_0 using the PF: increases in λ_0 above $2/h$ correspond to increases in the flow $\dot{\theta} = \nabla_{\theta} E$. The stability coefficient reflects sc_0 the changes in λ_0 .

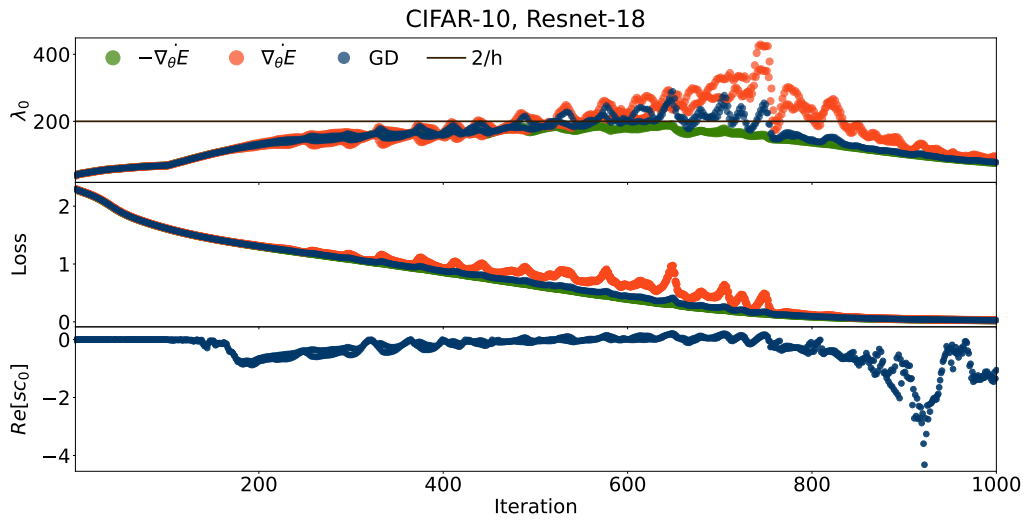


Figure A.12: Learning rate: 0.01. Edge of stability results: the connection between the PF, stability coefficients and loss instabilities with Resnet-18. Together with the behaviour of gradient descent, we plot the behaviour of the NGF and positive gradient flow initialized at θ_t and simulated for time h for each iteration t . The analysis we performed based on the PF suggests that when $\Re[s_{c_0}] > 0$ and large we should expect gradient descent to exhibit behaviours close to those of the positive gradient flow. What we observe empirically is that increases in loss value of gradient descent are proportional to the increase of the positive gradient flow in that area (can be seen best between iterations 600 and 800); the same behaviour can be seen in relation to the eigenvalue λ_0 . Results for a larger learning rate are shown in Figure 3.16.

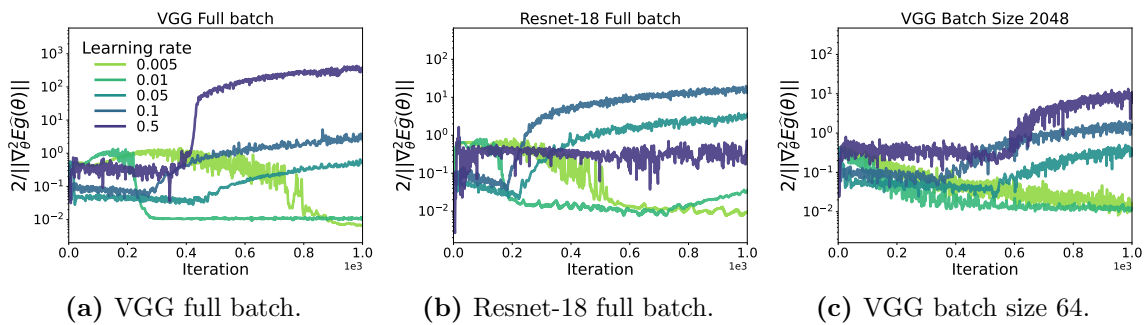


Figure A.13: Evaluating the value of the DAL learning rate for a fixed learning rate sweep to assess whether it would have reasonable magnitudes. We observe that the quantity is mainly in the ranges 10^{-2} to 10 , which can be used for model training. We note that this is a preliminary check only, as using DAL in training changes the optimisation trajectory and thus can affect these ranges. As we have seen however, DAL can be suitably used for model training.

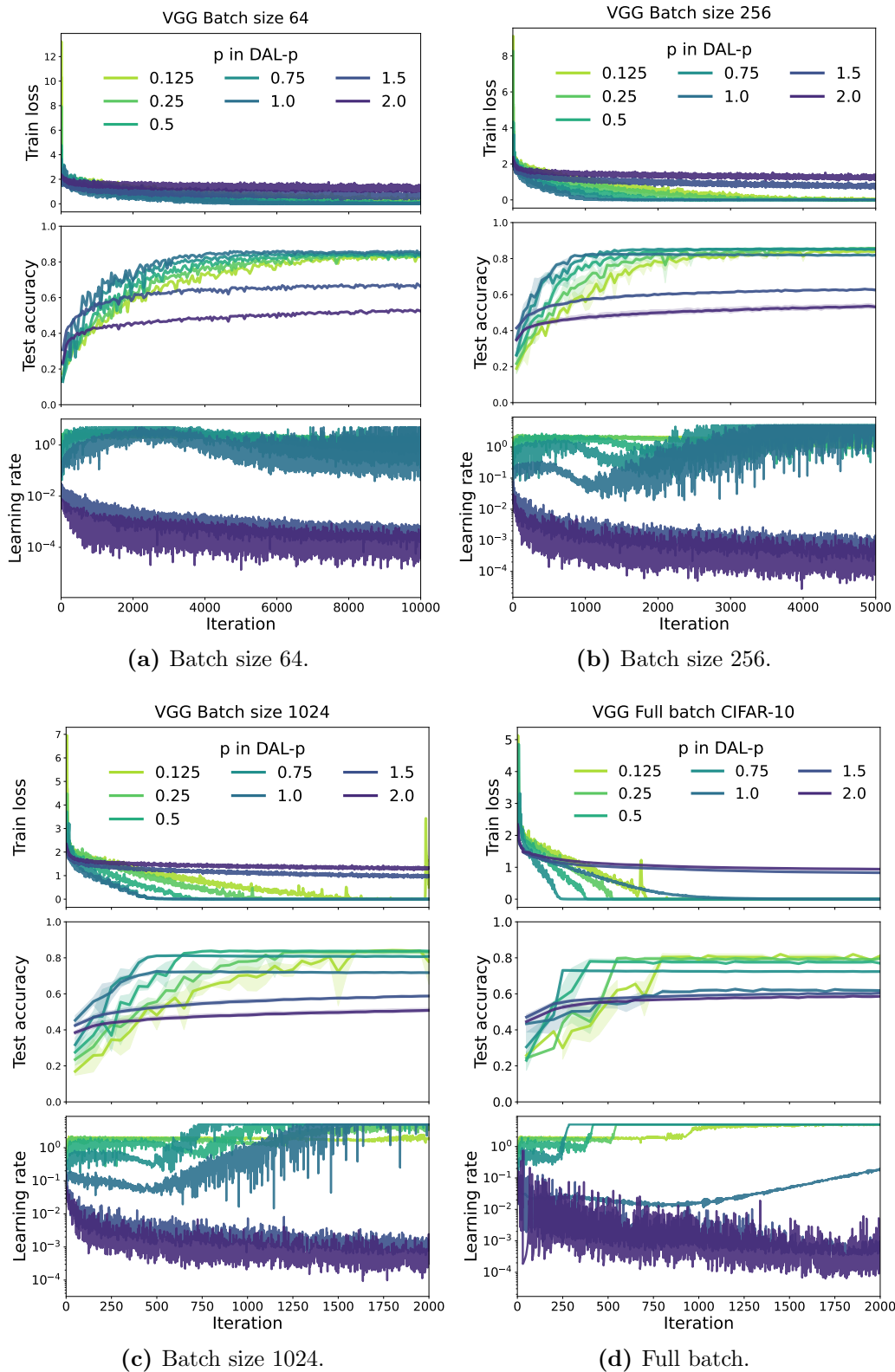


Figure A.14: DAL- p training a VGG model on CIFAR-10. Sweep across batch sizes: discretisation drift helps test performance, but at the cost of stability. We also show the effective learning rate and train losses and test accuracies.

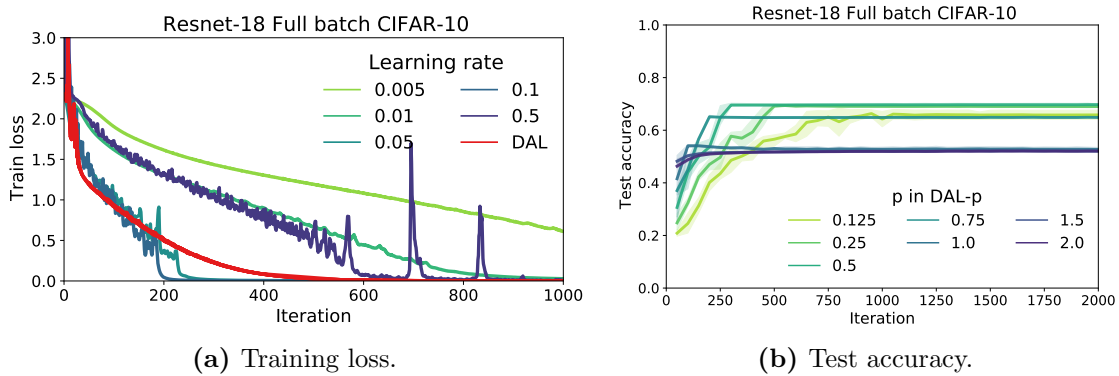


Figure A.15: DAL and DAL- p results with Resnet-18 on CIFAR-10: DAL results in improved stability without requiring a hyperparameter sweep, while DAL- p achieves high generalisation.

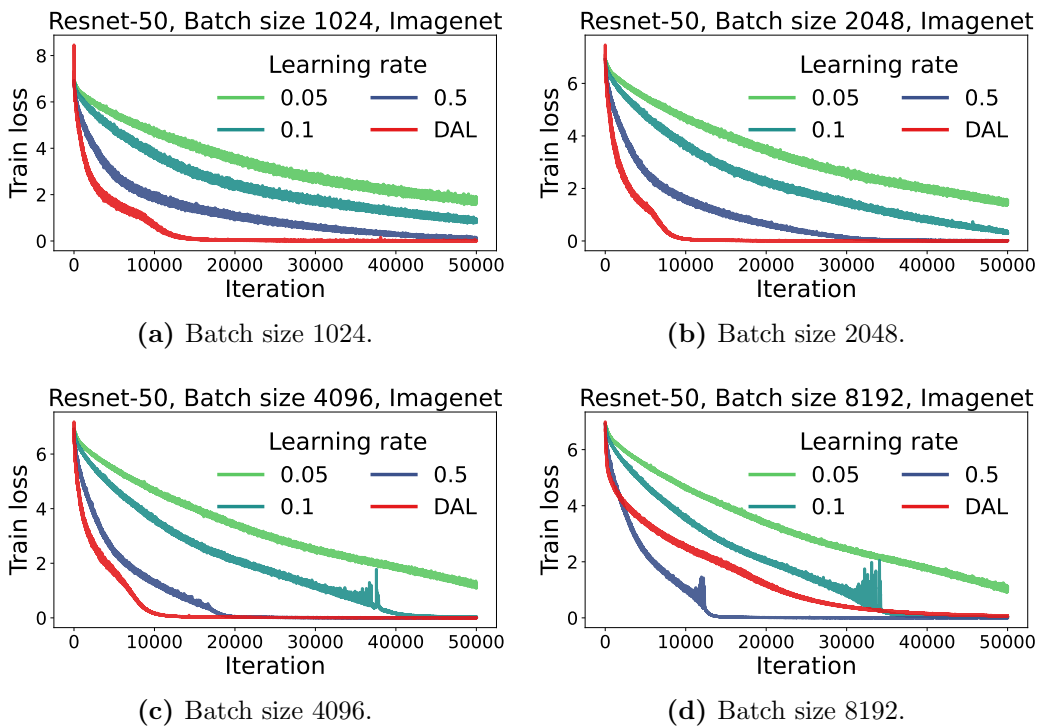


Figure A.16: DAL: Imagenet results across batch sizes. We observe quick convergence and stable training.

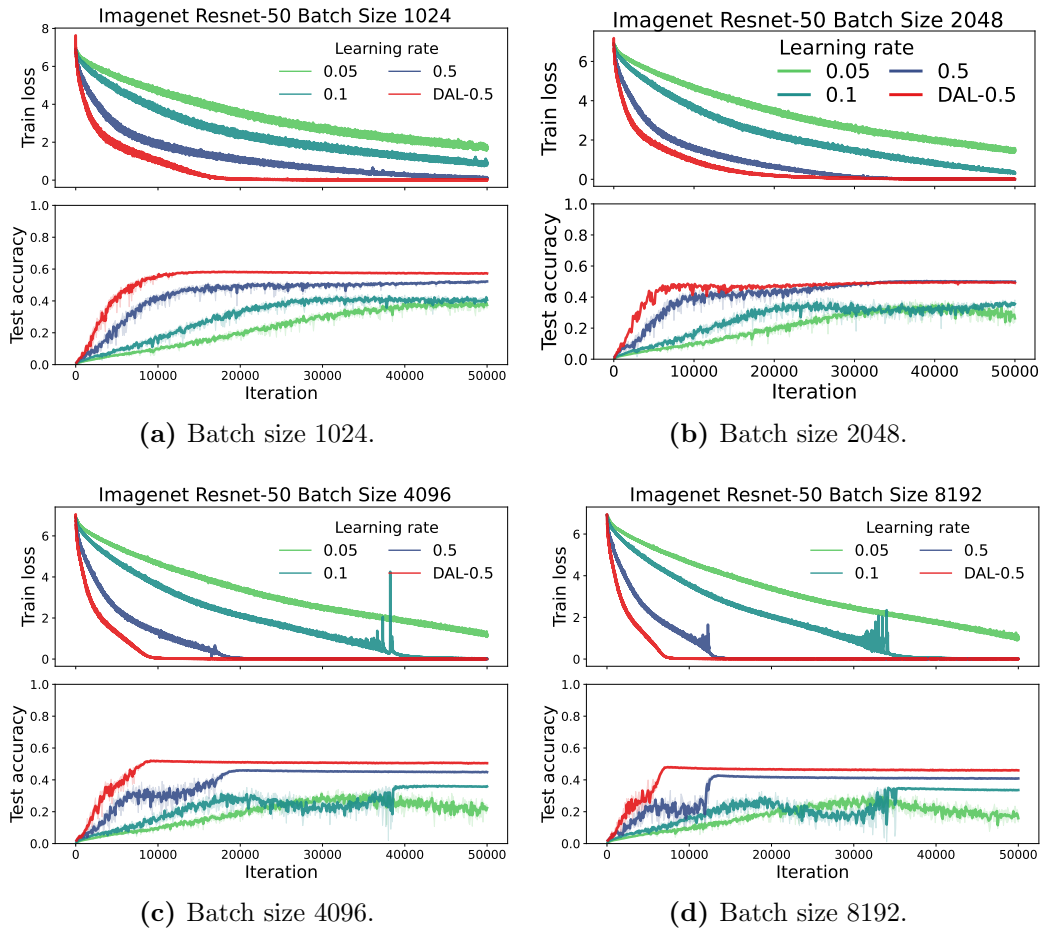


Figure A.17: DAL-0.5: Imagenet results across batch sizes. We observe quick training and good generalisation performance.

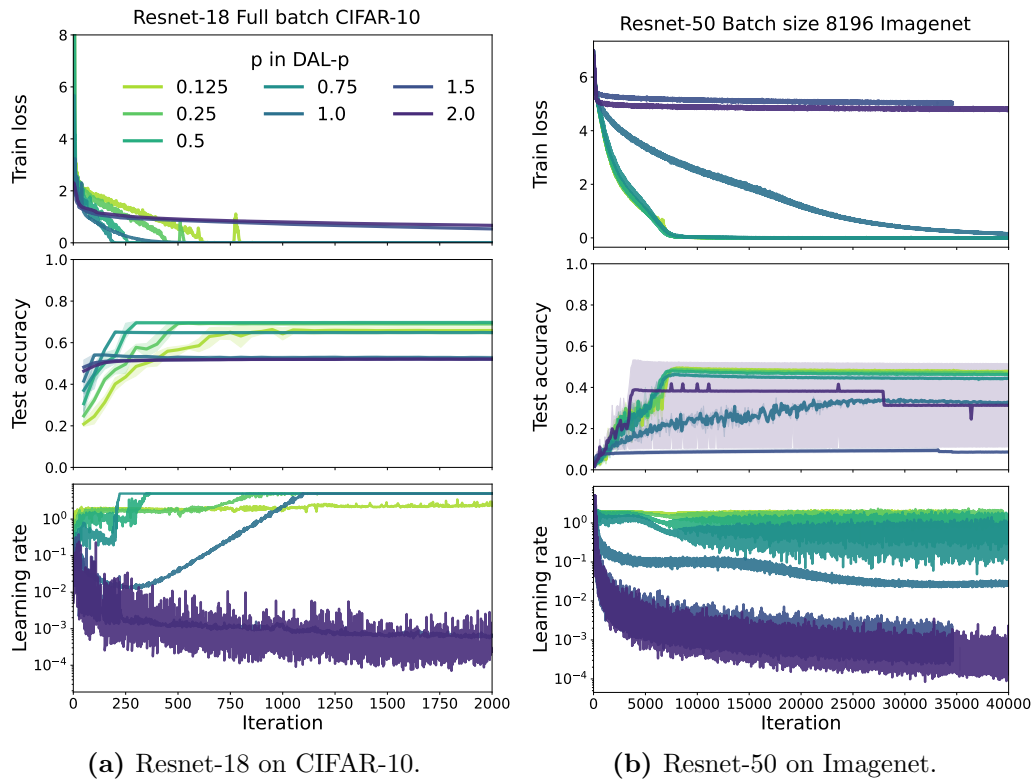


Figure A.18: DAL- p sweep: discretisation drift helps test performance, but at the cost of stability. We also show the effective learning rate and train losses and test accuracies.

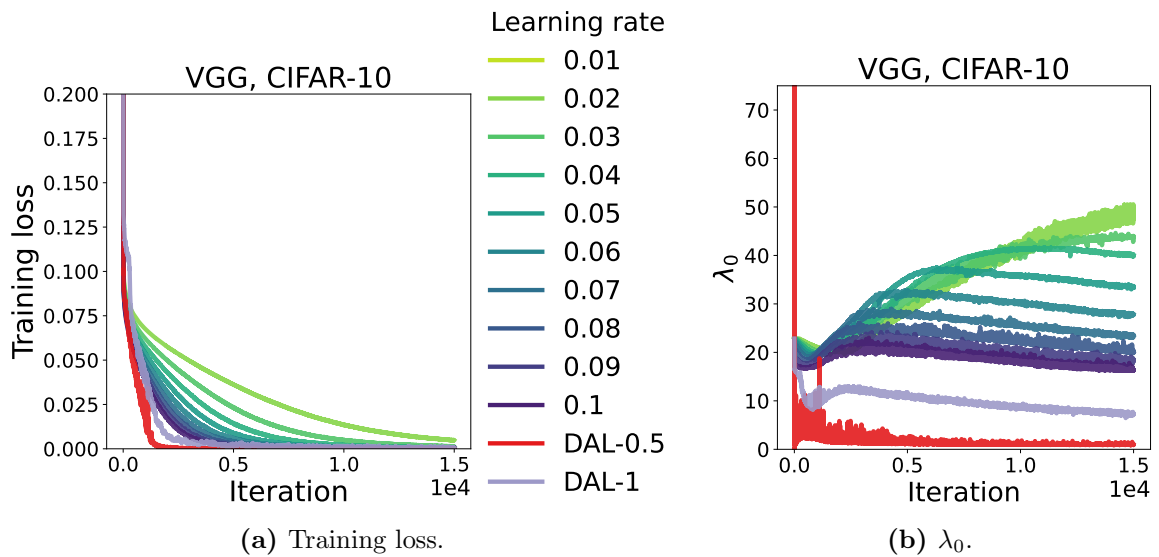


Figure A.19: Results with a least square loss. DAL and DAL-0.5 lead to quicker training and lower λ_0 in this setting as well.

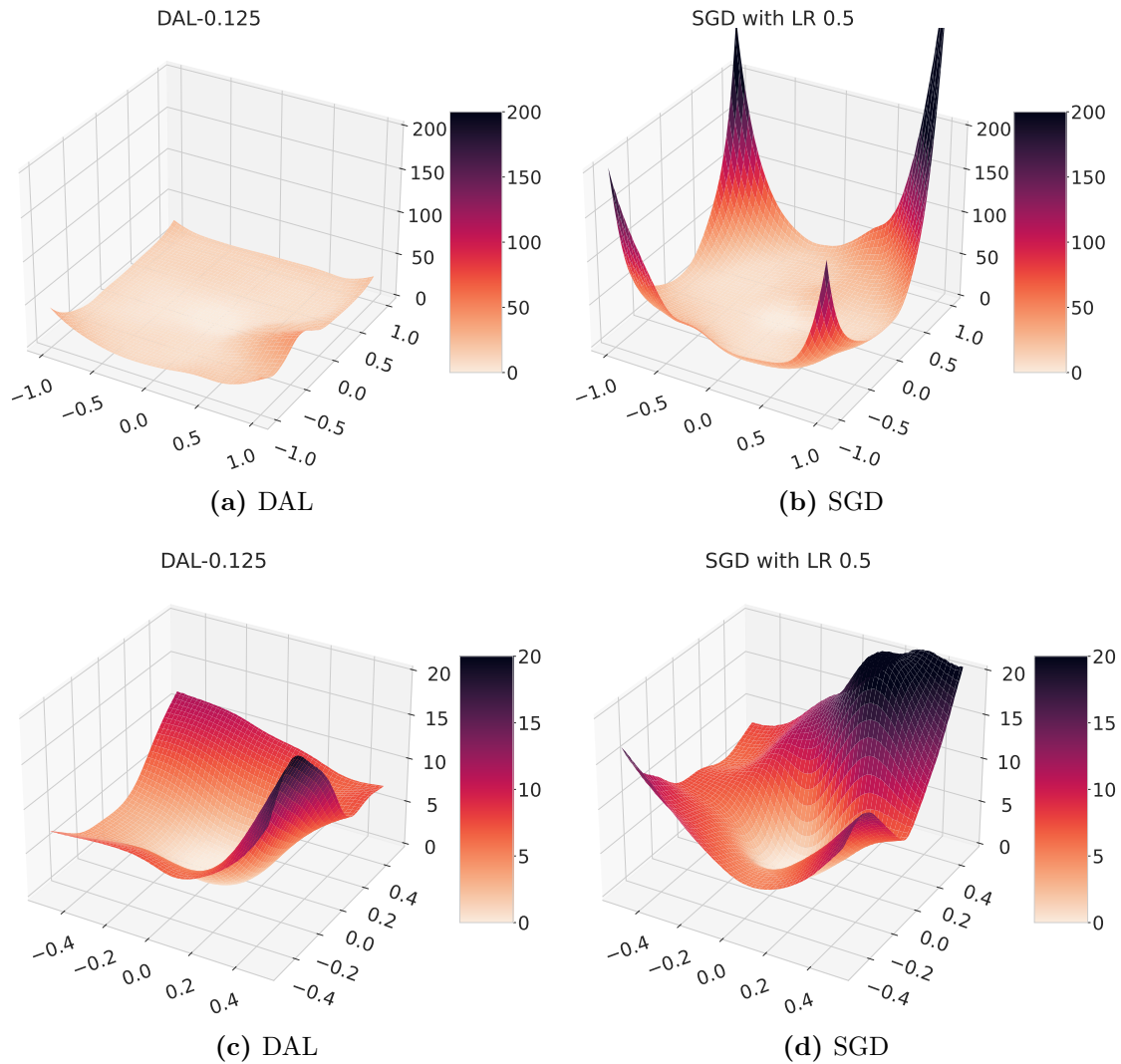


Figure A.20: CIFAR-10, full batch, at various levels of zoom into around the convergence points. The 2D projection of the DAL- p and SGD learned landscapes on CIFAR-10, using a VGG model. The visualisation is made using the method of Li et al. [22]. The DAL- p model achieves an accuracy of 82% which the SGD model achieves 77% accuracy. We also show the trajectory of λ_0 for both models in Figure A.21.

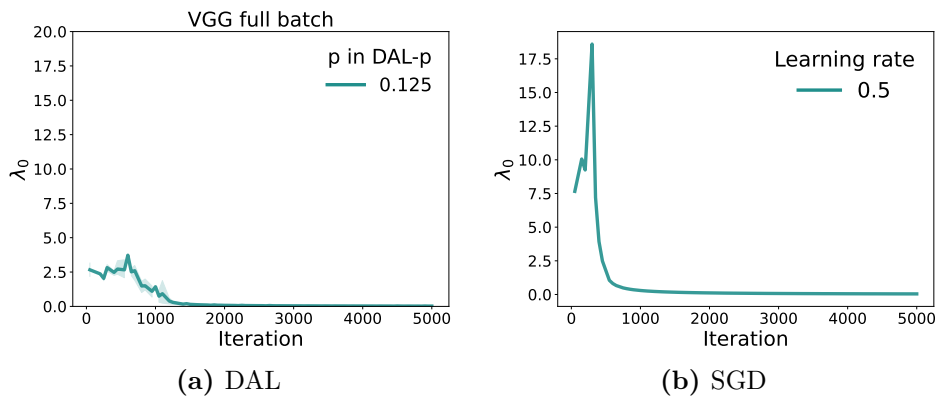


Figure A.21: λ_0 for learned models using full batch gradient descent and DAL- p on CIFAR-10: DAL traverses areas of space with smaller λ_0 .

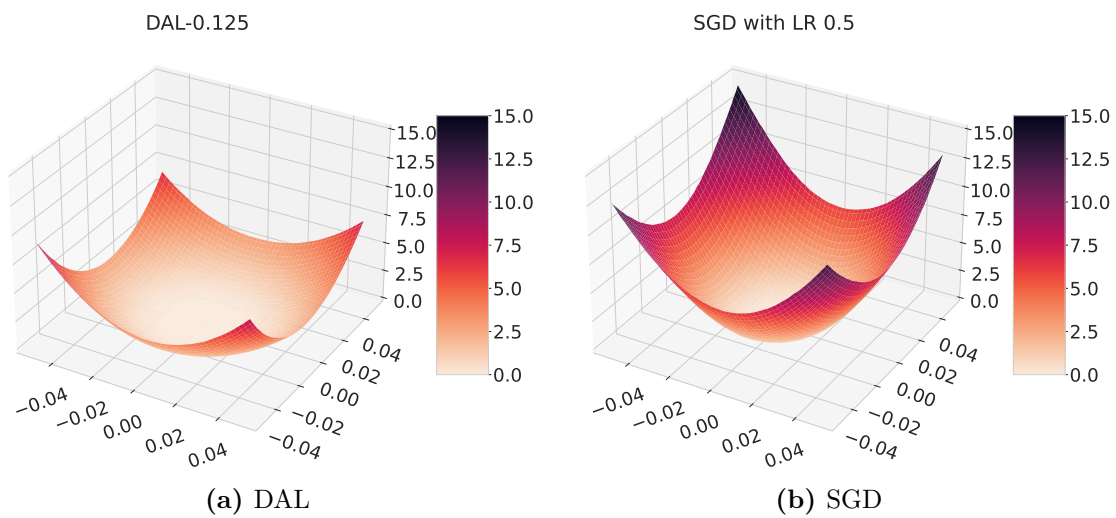
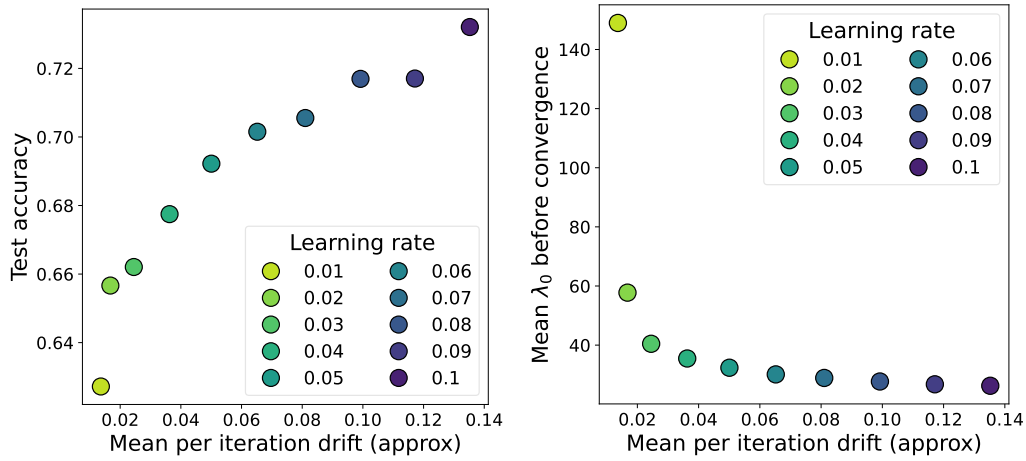
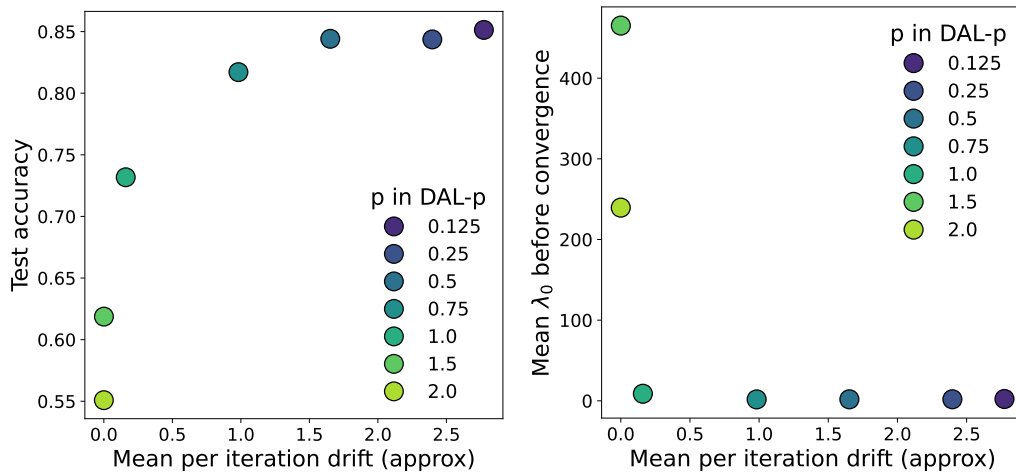


Figure A.22: Imagenet, batch size 8192. The 2D projection of the DAL- p and SGD learned landscapes, using the method of Li et al. [22]. The DAL- p model achieves an accuracy of 50% which the SGD model achieves 42% accuracy.



(a) Fixed learning rate sweep.



(b) DAL- p sweep.

Figure A.23: VGG CIFAR-10 with batch size 1024: connection between drift, test error and λ_0 . We observe the same patterns for other batch sizes, namely that the higher the drift the more generalisation and the lower λ_0 , both with fixed learning rates and DAL.

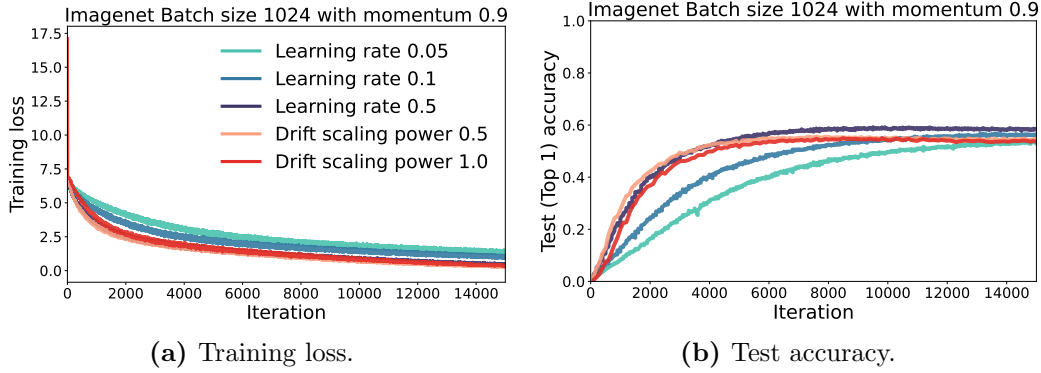


Figure A.24: Integrating information on drift with momentum: with smaller batch sizes, we observe a smaller empirical gain. Since these are preliminary results, more investigation is needed into why that is. These results are consistent with the observation obtained with vanilla gradient descent, where we say larger gains with larger batch sizes on Imagenet.

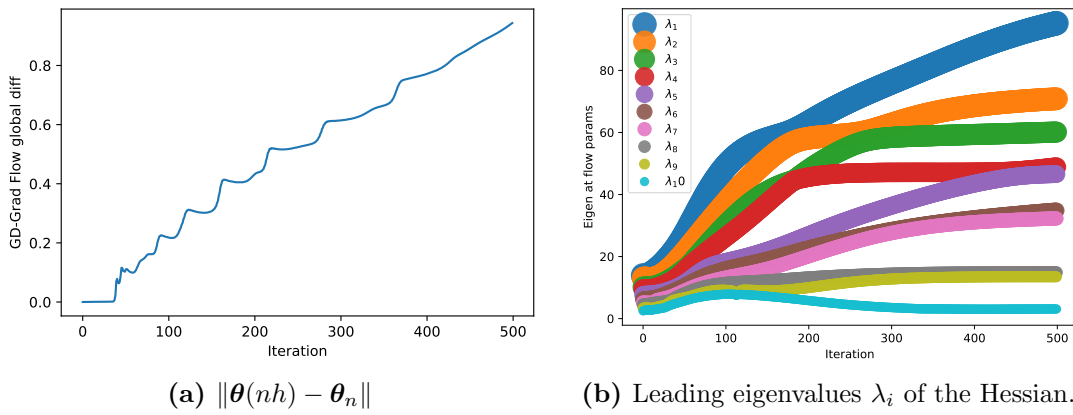


Figure A.25: The global error in parameter space $\|\theta(nh) - \theta_n\|$ between the NGF and gradient descent (a), as well as the behaviour of λ_0 of the NGF (b). We observe that early in training gradient descent follows the NGF, and then when following the NGF the leading eigenvalues grow.

A.5 Figure reproduction details

- Figure 3.2: 2D convex plot. The loss function is $E = \boldsymbol{\theta}^T A \boldsymbol{\theta}$ with $A = ((1, 0.0), (0.0, 0.01))$. The learning rate is 0.9.
- Figure 3.3: 1D example. $E = \frac{1}{2}(\boldsymbol{\theta} - 0.6)^2$ and $h = 1.2$
- Figure 3.6: $E(z) = \frac{1}{2}z^2$. Learning rates are 0.8, 1.5 and 2.1.
- Figure 3.7: Quadratic losses in 2 dimensions. $E = \boldsymbol{\theta}^T A \boldsymbol{\theta}$ with $A = ((1, 0.0), (0.0, 0.01))$, with learning rates 0.5, 0.9 and 1.05.
- Figure 3.8: Banana function. Learning rates 0.0006, 0.0017 0.005.
- Figure 3.9: Error between gradient descent parameters and parameters obtained following continuous-time flows for multiple iterations on a small MLP: $\|\boldsymbol{\theta}_{t+n} - \boldsymbol{\theta}(nh)\|$. Input size 2, MLP of with output units 10, 1. Learning rates 0.1, 0.2 and 0.25. We use a dataset with 5 examples where the input is generated from a Gaussian distribution with 2 dimensions and the targets are samples from a Gaussian distribution with 1 dimension. A mean square error loss is used.
- Figure 3.10: Predictions of $\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0$ using the PF on neural networks. Full-batch training with a VGG network on CIFAR-10. Learning rates 0.01, 0.03 and 0.05.
- Figure 3.12: Edge of stability on VGG networks on CIFAR-10. Full-batch training.
- Figure 3.13: Edge of stability on a small 5 layer MLP with 10 units per layer on the UCI Iris dataset. Full-batch training. The learning rate used is 0.18.
- Figure 3.14: Edge of stability results. MNIST MLP with 4 layers. Learning rate 0.05. For each model we approximate the NGF and the positive gradient flow initialised at the current gradient descent parameters and measure the value of the loss function and λ_0 .

- Figure 3.15: Using the PF and stability coefficient to understand instabilities in deep learning. MNIST results use a 4 layer MLP and a 0.05 learning rate. CIFAR-10 MLP use a VGG network and a 0.02 learning rate.
- Figure 3.16: Understanding the behaviour of the loss through the PF and the behaviour of λ_0 . Using a Resnet-18 trained on CIFAR-10 using a learning rate of 0.03. For each model we approximate the NGF and the positive gradient flow initialised at the current gradient descent parameters and measure the value of the loss function and λ_0 .
- Figure 3.17: Assessing whether 1 dimension is enough to cause instability in a continuous-time flow. We train a model with gradient descent until it reaches the edge of stability ($\lambda_0 \approx 2/h$), after which we use a approximate the continuous flow $\dot{\theta} = \nabla_{\theta} E^T \mathbf{u}_0 \mathbf{u}_0 + \sum_{i=1}^{D-1} -\nabla_{\theta} E^T \mathbf{u}_i \mathbf{u}_i$. The model is a VGG model trained on CIFAR-10 with learning rate 0.05 and after the edge of stability is reached the flow is approximated using Euler steps of size 5×10^{-5} .
- Figure 3.19: The unstable dynamics of $\nabla_{\theta} E^T u$ in the edge of stability area. The model results are obtained from a VGG model trained on CIFAR-10 with a learning rate of 0.01.
- Figure 3.20: Learning rate sweep showing the value of the non-principal third-order term in training, on a Resnet-18 model trained on CIFAR-10. The coefficients used for the non-principal term are those in Eq (3.4).
- Figure 3.21: $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ and the per-iteration drift as measured during training. The plots are obtained on an MNIST MLP with 4 layers and 100 units per layer, and a VGG model trained on CIFAR-10. The learning rates used are 0.05, 0.01 and 0.01 respectively.
- Figure 3.22: Correlation between $\|\nabla_{\theta}^2 E \nabla_{\theta} E\|$ and the iteration drift. The plots are obtained on an MNIST MLP with 4 layers and 100 units per layer, and a VGG model trained on CIFAR-10. Results are obtained over a learning rate sweep of 10 learning rates 0.01, 0.02, \dots , 0.1.

- Figure 3.23: Models trained using a learning rate sweep or DAL on CIFAR-10 and Imagenet. Models are VGG, Resnet-18 and Resnet-50 respectively.
- Figure 3.24: Key quantities in DAL versus fixed learning rate training: learning rate, and update norms. Results obtained on a Resnet-18 model trained on CIFAR-10.
- Figure 3.25: DAL-0.5 results on a CIFAR-10 and Imagenet using VGG, Resnet-18 and Resnet-50 model respectively.
- Figure 3.26: DAL- p sweep on full-batch training on CIFAR-10 with VGG and Resnet-18 models and a Resnet-50 model trained on Imagenet.
- Figure 3.27: VGG model trained on CIFAR-10 using full-batch gradient descent, either with a fixed learning rate in a sweep or a DAL- p sweep. The drift per-iteration is approximated using $h^2/2\nabla_{\theta}E^2\nabla_{\theta}E$ or $h(\theta)^2/2\nabla_{\theta}E^2\nabla_{\theta}E$ (the adaptive learning rate) in case of DAL- p .
- Figure 3.29: Resnet-50 results with 0.9 momentum and a learning rate sweep or DAL. Models trained on Imagenet.
- Figure A.2: CIFAR-10 DAL results with a Hessian vector product computation of $\nabla_{\theta}^2E\nabla_{\theta}E$ compared to an approximation. We used $\epsilon = 0.01$ in the approximation. Full-batch training.
- Figure A.3: CIFAR-10 DAL results with a Hessian vector product computation of $\nabla_{\theta}^2E\nabla_{\theta}E$ compared to an approximation. We used $\epsilon = 0.01$ in the approximation. Batch size 512.
- Figure A.4: Imagenet DAL results with a Hessian vector product computation of $\nabla_{\theta}^2E\nabla_{\theta}E$ compared to an approximation. We used $\epsilon = 0.01$ in the approximation.
- Figure A.5: Imagenet DAL-0.5 results with a Hessian vector product computation of $\nabla_{\theta}^2E\nabla_{\theta}E$ compared to an approximation. We used $\epsilon = 0.01$ in the approximation.

- Figure A.6: Results with a non quadratic function, including cos and sin. The function used is: $E(\theta) = \cos(\theta) + \theta$, if $\theta < 0$; $(\theta/3)^2 2 + 1 + \theta/3$ otherwise.
- Figure 3.4: Per-iteration error between the NGF, IGR and PF and gradient descent. MNIST learning rate 0.05. CIFAR learning rate 0.02. We obtain the results by training a model with gradient descent $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - h \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1})$ and at each iteration t of gradient descent approximating the respective flows for time h and computing the difference in norms between the resulting flow parameters and the gradient descent parameters at the next iteration: $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}(h)\|$ with $\boldsymbol{\theta}(h) = \boldsymbol{\theta}_{t-1}$. Results obtained on full-batch training.
- Figure A.7: Global parameter error for IGR ad PF flows compared on an MLP trained on the UCI breast cancer dataset. MLP with 10 hidden units and 1 output unit. Learning rates 1.5 and 5. respectively.
- Figure 3.11: We train a VGG model on CIFAR-10 with learning rate 0.1, and as it reaches convergence increase the learning rate slightly so that $\lambda_0 > 2/h$. We use a square loss here to avoid training the decrease in eigenvalues later in training associated with cross entropy losses.
- Figure A.8: The eigenspectrum obtained along the gradient descent path of a 5 layer small MLP trained on the UCI Iris dataset. The Hessian eigenvalues are largely positive. The result accompanies Figure 3.13 in the main thesis.
- Figure 3.18: Training a VGG model on CIFAR-10 and training with learning rate 0.05 after which switching to 0.01.
- Figure A.9: MNIST results when changing one direction in continuous-time. The model is a 4 layer MLP with 100 hidden units per layer. The model was trained with learning rate 0.05 until the edge of stability area is reached after which the flow $\dot{\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_0 \mathbf{u}_0 + \sum_{i=1}^{D-1} -\nabla_{\boldsymbol{\theta}} E^T \mathbf{u}_i \mathbf{u}_i$ is used.
- Figure A.10: Peak areas of loss increase and the stability coefficient for the largest eigendirection of the PF. Results using a VGG model trained with full-batch gradient descent on CIFAR-10 across a learning rate sweep.

- Figure A.11: Understanding changes to λ_0 using the PF. MLP with 4 layers with 100 units each trained on MNIST, and VGG model trained on CIFAR-10. Learning rate 0.05 and 0.01 respectively.
- Figure A.12: The connection between the PF, λ_0 and loss instabilities with Resnet-18 trained on CIFAR-10 with a learning rate 0.01.
- Figure A.13: Evaluating the value of the DAL learning rate for a vanilla (fixed) learning rate sweep to assess whether it would have reasonable magnitudes. Models trained on CIFAR-10 using a learning rate sweep.
- Figure A.14: DAL- p sweep with a VGG model trained on CIFAR-10 for small and large batch sizes.
- Figure A.16: DAL on Imagenet, a learning rate sweep across batch sizes.
- Figure A.17: DAL-0.5: on Imagenet, a learning rate sweep across batch sizes.
- Figure A.18: DAL- p sweep on VGG, Resnet-18 and Imagenet.
- Figure A.19: DAL results with a least square loss. Full-batch training of a VGG model on CIFAR-10.
- Figure 3.30: Resnet-50 Imagenet training with a learning rate per parameter. Batch size 8192.
- Figure A.23: VGG trained on CIFAR-10 with batch size 1024: connection between drift, test error and λ_0 .
- Figure A.24: DAL- p with momentum 0.9 on Imagenet. The model is Resnet-50 trained with batch size 1024.
- Figure A.25: MNIST results with a 4 layer MLP with 100 hidden units. The learning rate used is 0.05.

Appendix B

Continuous time models of gradient descent in two-player games

B.1 Proof of the main theorems

Notation: We use $\phi \in \mathbb{R}^m$ to denote the parameters of the first player and $\theta \in \mathbb{R}^n$ for the second player. We assume that the players parameters are updated simultaneously with learning rates $v_\phi h$ and $v_\theta h$ respectively. We consider the vector fields $f(\phi, \theta) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g(\phi, \theta) : \mathbb{R}^m \times \mathbb{R}^n \rightarrow \mathbb{R}^n$. The Jacobian $\frac{df}{d\theta}(\phi, \theta) \in \mathbb{R}^{n,m}$ denotes the matrix with entries $[\frac{df}{d\theta}(\phi, \theta)]_{i,j} = \frac{df_j}{d\theta_i}$ with $i \in \{1, \dots, n\}, j \in \{1, \dots, m\}$.

We now prove Theorem 4.2.1 and Theorem 4.2.2 in the main thesis, corresponding to the *simultaneous* Euler updates and to the *alternating* Euler updates, respectively. In both cases, our goal is to find corrections f_1 and g_1 to the original system

$$\dot{\phi} = f(\phi, \theta), \tag{B.1}$$

$$\dot{\theta} = g(\phi, \theta), \tag{B.2}$$

such that the modified continuous system

$$\dot{\phi} = f(\phi, \theta) + hf_1(\phi, \theta), \tag{B.3}$$

$$\dot{\theta} = g(\phi, \theta) + hg_1(\phi, \theta), \tag{B.4}$$

follows the discrete steps of the method with a local error of order $\mathcal{O}(h^3)$. More

precisely, if (ϕ_t, θ_t) denotes the discrete step of the method at time t , and $(\phi(h), \theta(h))$ corresponds to the continuous solution of the modified system above starting at $(\phi_{t-1}, \theta_{t-1})$, we want that the local errors for both players, i.e.,

$$\|\phi_t - \phi(v_\phi h)\| \quad \text{and} \quad \|\theta_t - \theta(v_\theta h)\|$$

to be of order $\mathcal{O}(h^3)$. In the expression above, $v_\phi h$ and $v_\theta h$ are the effective learning rates (or step-sizes) for the first and the second player respectively for both the simultaneous and alternating Euler updates.

Our proofs from BEA follow the same steps (see Section 5.1.1 for an overview):

1. Expand the discrete updates to find a relationship between ϕ_t and ϕ_{t-1} and θ_t and θ_{t-1} up to order $\mathcal{O}(h^2)$.
2. Expand the changes in continuous-time of the modified flow given by BEA.
3. Find the first order Discretisation Drift (DD) by matching the discrete and continuous updates up to second order in learning rates.

Notation: To avoid cluttered notations, we use $f_{(t)}$ to denote the $f(\phi_t, \theta_t)$ and $g_{(t)}$ to denote $g(\phi_t, \theta_t)$ for all t . If no index is specified, we use f to denote $f(\phi, \theta)$, where ϕ and θ are variables in a continuous system.

B.1.1 Simultaneous updates (Theorem 4.2.1)

Here we prove Theorem 4.2.1 in the main thesis, which we reproduce here:

The simultaneous Euler updates with learning rates $v_\phi h$ and $v_\theta h$ respectively are given by:

$$\phi_t = \phi_{t-1} + v_\phi h f(\phi_{t-1}, \theta_{t-1}) \tag{B.5}$$

$$\theta_t = \theta_{t-1} + v_\theta h g(\phi_{t-1}, \theta_{t-1}) \tag{B.6}$$

Theorem (DD for simultaneous Euler updates) (*Theorem 4.2.1*) *The discrete simultaneous Euler updates in (B.5) and (B.6) follow the continuous system*

$$\dot{\phi} = f - \frac{v_\phi h}{2} \left(\frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \quad (\text{B.7})$$

$$\dot{\theta} = g - \frac{v_\theta h}{2} \left(\frac{dg}{d\phi} f + \frac{dg}{d\theta} g \right) \quad (\text{B.8})$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Step 1: Expand the updates per player via a Taylor expansion.

We expand the numerical scheme to find a relationship between ϕ_t and ϕ_{t-1} and θ_t and θ_{t-1} up to order $\mathcal{O}(h^2)$. In the case of the simultaneous Euler updates, this does not require any change to Eqs (B.5) and (B.6).

Step 2: Expand the continuous-time changes for the modified flow given by BEA.

Lemma B.1.1 *If:*

$$\dot{\phi} = \tilde{f}(\phi, \theta) \quad (\text{B.9})$$

$$\dot{\theta} = \tilde{g}(\phi, \theta) \quad (\text{B.10})$$

where

$$\tilde{f}(\phi, \theta) = f(\phi, \theta) + \sum_{i=1} \tau_\phi^i f_i(\phi, \theta) \quad (\text{B.11})$$

$$\tilde{g}(\phi, \theta) = g(\phi, \theta) + \sum_{i=1} \tau_\theta^i g_i(\phi, \theta) \quad (\text{B.12})$$

and τ_θ and τ_ϕ are scalars. Then—for ease of notation, we drop the argument τ on the evaluations of ϕ and θ on the RHS:

$$\phi(\tau + \tau_\phi) = \phi(\tau) + \tau_\phi f + \tau_\phi^2 f_1 + \tau_\phi^2 \frac{1}{2} \frac{df}{d\phi} f + \tau_\phi^2 \frac{1}{2} \frac{df}{d\theta} g + \mathcal{O}(\tau_\phi^3) \quad (\text{B.13})$$

Proof: We perform a Taylor expansion:

$$\phi(\tau + \tau_\phi) \tag{B.14}$$

$$= \phi + \tau_\phi \dot{\phi} + \tau_\phi^2 \frac{1}{2} \ddot{\phi} + \mathcal{O}(\tau_\phi^3) \tag{B.15}$$

$$= \phi + \tau_\phi \tilde{f} + \tau_\phi^2 \frac{1}{2} \dot{\tilde{f}} + \mathcal{O}(\tau_\phi^3) \tag{B.16}$$

$$= \phi + \tau_\phi \tilde{f} + \tau_\phi^2 \frac{1}{2} \left(\frac{d\tilde{f}}{d\phi} \tilde{f} + \frac{d\tilde{f}}{d\theta} \tilde{g} \right) + \mathcal{O}(\tau_\phi^3) \tag{B.17}$$

$$= \phi + \tau_\phi (f + \tau_\phi f_1 + \mathcal{O}(\tau_\phi^2)) + \tau_\phi^2 \frac{1}{2} \left(\frac{d\tilde{f}}{d\phi} \tilde{f} + \frac{d\tilde{f}}{d\theta} \tilde{g} \right) + \mathcal{O}(\tau_\phi^3) \tag{B.18}$$

$$= \phi + \tau_\phi f + \tau_\phi^2 f_1 + \tau_\phi^2 \frac{1}{2} \left(\frac{d\tilde{f}}{d\phi} \tilde{f} + \frac{d\tilde{f}}{d\theta} \tilde{g} \right) + \mathcal{O}(\tau_\phi^3) \tag{B.19}$$

$$= \phi + \tau_\phi f + \tau_\phi^2 f_1 \tag{B.20}$$

$$+ \tau_\phi^2 \frac{1}{2} \left(\frac{d\tilde{f}}{d\phi} (f + \tau_\phi f_1 + \mathcal{O}(\tau_\phi^2)) + \frac{d\tilde{f}}{d\theta} (g + \tau_\theta g_1 + \mathcal{O}(\tau_\theta^2)) \right) + \mathcal{O}(\tau_\phi^3) \tag{B.21}$$

$$= \phi + \tau_\phi f + \tau_\phi^2 f_1 + \tau_\phi^2 \frac{1}{2} \frac{d\tilde{f}}{d\phi} f + \tau_\phi^2 \frac{1}{2} \frac{d\tilde{f}}{d\theta} g + \mathcal{O}(\tau_\phi^3) \tag{B.22}$$

$$= \phi + \tau_\phi f + \tau_\phi^2 f_1 + \tau_\phi^2 \frac{1}{2} \frac{df}{d\phi} f + \tau_\phi^2 \frac{1}{2} \frac{df}{d\theta} g + \mathcal{O}(\tau_\phi^3) \tag{B.23}$$

where we assumed that τ_ϕ and τ_θ are in the same order of magnitude. \square

Step 3: Matching the discrete and modified continuous updates.

We substitute the values given by the discrete updates, namely ϕ_{t-1} and θ_{t-1} , in Lemma B.1.1 in order to calculate the displacement according to the continuous updates:

$$\phi(\tau + \tau_\phi) = \phi_{t-1} + \tau_\phi f_{(t-1)} + \tau_\phi^2 f_1(\phi_{t-1}, \theta_{t-1}) \tag{B.24}$$

$$+ \frac{1}{2} \tau_\phi^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} \tau_\phi^2 \frac{df}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \mathcal{O}(\tau_\phi^3) \tag{B.25}$$

$$\theta(\tau + \tau_\theta) = \theta_{t-1} + \tau_\theta g_{(t-1)} + \tau_\theta^2 g_1(\phi_{t-1}, \theta_{t-1}) \tag{B.26}$$

$$+ \frac{1}{2} \tau_\theta^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} \tau_\theta^2 \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \mathcal{O}(\tau_\theta^3) \tag{B.27}$$

In order to find f_1 and g_1 such that the continuous dynamics of the modified updates

$f + v_\phi h f_1$ and $g + v_\theta h g_1$ match the discrete updates in Eqs (B.5) and (B.6), we look for the corresponding continuous increments of the discrete updates in the modified continuous system, such that $\|\phi(\tau + \tau_\phi) - \phi_t\|$ and $\|\theta(\tau + \tau_\theta) - \theta_t\|$ are $\mathcal{O}(h^3)$. The first order terms in Eqs (B.5) and (B.6) and Eqs (B.26) and (B.27) suggest that:

$$v_\phi h = \tau_\phi \tag{B.28}$$

$$v_\theta h = \tau_\theta \tag{B.29}$$

We can now proceed to find f_1 and g_1 from the second order terms

$$0 = v_\phi^2 h^2 f_1(\phi_{t-1}, \theta_{t-1}) + \frac{1}{2} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f^{(t-1)} + \frac{1}{2} v_\phi^2 h^2 \frac{df}{d\theta} \Big|_{(t-1)} g^{(t-1)} \tag{B.30}$$

$$f_1(\phi_{t-1}, \theta_{t-1}) = -\frac{1}{2} \frac{df}{d\phi} \Big|_{(t-1)} f^{(t-1)} - \frac{1}{2} \frac{df}{d\theta} \Big|_{(t-1)} g^{(t-1)}. \tag{B.31}$$

Similarly, for g_1 we obtain

$$g_1(\phi_{t-1}, \theta_{t-1}) = -\frac{1}{2} \frac{dg}{d\phi} \Big|_{(t-1)} f^{(t-1)} - \frac{1}{2} \frac{dg}{d\theta} \Big|_{(t-1)} g^{(t-1)}. \tag{B.32}$$

Leading to the first order corrections

$$f_1(\phi_{t-1}, \theta_{t-1}) = -\frac{1}{2} \frac{df}{d\phi} \Big|_{(t-1)} f^{(t-1)} - \frac{1}{2} \frac{df}{d\theta} \Big|_{(t-1)} g^{(t-1)} \tag{B.33}$$

$$g_1(\phi_{t-1}, \theta_{t-1}) = -\frac{1}{2} \frac{dg}{d\phi} \Big|_{(t-1)} f^{(t-1)} - \frac{1}{2} \frac{dg}{d\theta} \Big|_{(t-1)} g^{(t-1)}. \tag{B.34}$$

We have found the functions f_1 and g_1 such that after one discrete optimisation step the flows $\dot{\phi} = f + v_\phi h f_1$ and $\dot{\theta} = g + v_\theta h g_1$ follow the discrete updates up to order $\mathcal{O}(h^3)$, finishing the proof.

B.1.2 Alternating updates (Theorem 4.2.2)

Here we prove Theorem 4.2.2 in the main thesis, which we reproduce here. For the *alternating Euler updates*, the players take turns to update their parameters, and can perform multiple updates each. We denote the number of alternating updates of the first player (resp. second player) by m (resp. k). We scale the learning rates by the

number of updates, leading to the following updates $\phi_t := \phi_{m,t}$ and $\theta_t := \theta_{k,t}$ where

$$\phi_{i,t} = \phi_{i-1,t} + \frac{v_\phi h}{m} f(\phi_{i-1,t}, \theta_{t-1}), \quad i = 1 \dots m, \quad (\text{B.35})$$

$$\theta_{j,t} = \theta_{j-1,t} + \frac{v_\theta h}{k} g(\phi_{m,t}, \theta_{j-1,t}), \quad j = 1 \dots k. \quad (\text{B.36})$$

Theorem DD for alternating Euler updates (*Theorem 4.2.2*) *The discrete alternating Euler updates in (B.35) and (B.36) follow the continuous system*

$$\dot{\phi} = f - \frac{v_\phi h}{2} \left(\frac{1}{m} \frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \quad (\text{B.37})$$

$$\dot{\theta} = g - \frac{v_\theta h}{2} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} f + \frac{1}{k} \frac{dg}{d\theta} g \right) \quad (\text{B.38})$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Step 1: Discrete updates

In the case of alternating Euler discrete updates, we have

$$\phi_{1,t} = \phi_{t-1} + \frac{v_\phi}{m} h f(\phi_{t-1}, \theta_{t-1}) \quad (\text{B.39})$$

$$\dots \quad (\text{B.40})$$

$$\phi_{m,t} = \phi_{m-1,t} + \frac{v_\phi}{m} h f(\phi_{m-1,t}, \theta_{t-1}) \quad (\text{B.41})$$

$$\theta_{1,t} = \theta_{t-1} + \frac{v_\theta}{k} h g(\phi_{m,t}, \theta_{t-1}) \quad (\text{B.42})$$

$$\dots \quad (\text{B.43})$$

$$\theta_{k,t} = \theta_{k-1,t} + \frac{v_\theta}{k} h g(\phi_{m,t}, \theta_{k-1,t}) \quad (\text{B.44})$$

$$\phi_t = \phi_{m,t} \quad (\text{B.45})$$

$$\theta_t = \theta_{k,t}. \quad (\text{B.46})$$

Lemma B.1.2 *Given updates $\phi_{m,t} = \phi_{m-1,t} + h f(\phi_{m-1,t}, \theta_{t-1})$, we can write:*

$$\phi_{m,t} = \phi_{t-1} + m h f_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.47})$$

Proof: Proof by induction. Base step.

$$\phi_{2,t} = \phi_{1,t} + hf(\phi_{1,t}, \theta_{t-1}) \quad (\text{B.48})$$

$$= \phi_{t-1} + hf(\phi_{t-1}, \theta_{t-1}) + hf(\phi_{1,t}, \theta_{t-1}) \quad (\text{B.49})$$

$$= \phi_{t-1} + hf_{(t-1)} + hf(\phi_{t-1} + hf_{(t-1)}, \theta_{t-1}) \quad (\text{B.50})$$

$$= \phi_{t-1} + hf_{(t-1)} + h \left(f_{(t-1)} + h \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \right) \quad (\text{B.51})$$

$$= \phi_{t-1} + 2hf_{(t-1)} + h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.52})$$

Inductive step:

$$\phi_{m+1,t} = \phi_{m,t} + hf(\phi_{m,t}, \theta_{t-1}) \quad (\text{B.53})$$

$$= \phi_{t-1} + mhf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.54})$$

$$+ hf \left(\phi_{t-1} + mhf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3), \theta_{t-1} \right) + \mathcal{O}(h^3) \quad (\text{B.55})$$

$$= \phi_{t-1} + mhf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.56})$$

$$+ h \left(f_{(t-1)} + \frac{df}{d\phi} \Big|_{(t-1)} \left(mhf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \right) \right) + \mathcal{O}(h^3) \quad (\text{B.57})$$

$$= \phi_{t-1} + (m+1)hf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.58})$$

$$+ h \frac{df}{d\phi} \Big|_{(t-1)} \left(mhf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \right) + \mathcal{O}(h^3) \quad (\text{B.59})$$

$$= \phi_{t-1} + (m+1)hf_{(t-1)} + \frac{m(m-1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.60})$$

$$+ h^2 m \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.61})$$

$$= \phi_{t-1} + (m+1)hf_{(t-1)} + \frac{m(m+1)}{2} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.62})$$

□

From Lemma B.1.2 with $h = v_\phi h/m$ we have that

$$\phi_{m,t} = \phi_{t-1} + v_\phi h f_{(t-1)} + \frac{m-1}{2m} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3). \quad (\text{B.63})$$

We now turn our attention to the second player. We define $g_{t'} = g(\phi_{m,t}, \theta_{t-1})$. This captures the difference between simultaneous and alternating updates. From Lemma B.1.2 with $h = v_\theta h/k$ we have that

$$\theta_{k,t} = \theta_{t-1} + v_\theta h g_{t'} + \frac{(k-1)}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{t'} g_{t'} + \mathcal{O}(h^3). \quad (\text{B.64})$$

We now expand $g_{t'}$ by Taylor expansion:

$$g_{t'} = g(\phi_{m,t}, \theta_{t-1}) \quad (\text{B.65})$$

$$= g\left(\phi_{t-1} + v_\phi h f_{(t-1)} + v_\phi^2 \frac{m-1}{2m} h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3), \theta_{t-1}\right) \quad (\text{B.66})$$

$$= g_{(t-1)} + \frac{dg}{d\phi} \Big|_{(t-1)} \left(v_\phi h f_{(t-1)} + \frac{(m-1)}{2m} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3) \right) \quad (\text{B.67})$$

If we replace the above in Eq (B.64):

$$\theta_{k,t} = \theta_{t-1} + v_\theta h g_{t'} + \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{t'} g_{t'} + \mathcal{O}(h^3) \quad (\text{B.68})$$

$$= \theta_{t-1} + v_\theta h \left(g_{(t-1)} + \frac{dg}{d\phi} \Big|_{(t-1)} \left(v_\phi h f_{(t-1)} + \frac{m-1}{2m} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \right) \right) \quad (\text{B.69})$$

$$+ \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{t'} g_{t'} + \mathcal{O}(h^3) \quad (\text{B.70})$$

$$= \theta_{t-1} + v_\theta h g_{(t-1)} + v_\theta h \frac{dg}{d\phi} \Big|_{(t-1)} \left(v_\phi h f_{(t-1)} + \frac{m-1}{2m} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \right) \quad (\text{B.71})$$

$$+ \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{t'} g_{t'} + \mathcal{O}(h^3) \quad (\text{B.72})$$

$$= \theta_{t-1} + v_\theta h g_{(t-1)} + v_\theta v_\phi h^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{t'} g_{t'} + \mathcal{O}(h^3) \quad (\text{B.73})$$

$$\boldsymbol{\theta}_{k,t} = \boldsymbol{\theta}_{t-1} + v_\theta h g_{(t-1)} + v_\theta v_\phi h^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.74})$$

$$+ \frac{k-1}{2k} h^2 \frac{dg}{d\theta} \Big|_{t'} \left(g_{t-1} + \frac{dg}{d\phi} \Big|_{(t-1)} \left(v_\phi h f_{(t-1)} + \frac{(m-1)}{2m} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} \right) \right) \quad (\text{B.75})$$

$$+ \mathcal{O}(h^3) \quad (\text{B.76})$$

$$= \boldsymbol{\theta}_{t-1} + v_\theta h g_{(t-1)} + v_\theta v_\phi h^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{t'} g_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.77})$$

$$= \boldsymbol{\theta}_{t-1} + v_\theta h g_{(t-1)} + v_\theta v_\phi h^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.78})$$

$$+ \frac{k-1}{2k} v_\theta^2 h^2 \frac{d \left(g + \frac{dg}{d\phi} \left(v_\phi h f + \frac{(m-1)}{2m} v_\phi^2 h^2 \frac{df}{d\phi} f \right) \right)}{d\theta} \Big|_{t-1} g_{(t-1)} + \mathcal{O}(h^3) \quad (\text{using Eq (B.67), B.79})$$

$$= \boldsymbol{\theta}_{t-1} + v_\theta h g_{(t-1)} + v_\theta v_\phi h^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.80})$$

We then have:

$$\phi_{m,t} = \phi_{t-1} + v_\phi h f_{(t-1)} + \frac{m-1}{2m} v_\phi^2 h^2 \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.81})$$

$$\boldsymbol{\theta}_{k,t} = \boldsymbol{\theta}_{t-1} + v_\theta h g_{(t-1)} + v_\theta v_\phi h^2 \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \mathcal{O}(h^3) \quad (\text{B.82})$$

Step 2: Expand the continuous-time changes for the modified flow given by BEA (Identical to the simultaneous update case.)

Step 3: Matching the discrete and modified continuous updates.

The linear terms are identical to those in the simultaneous updates:

$$\tau_\phi = v_\phi h \quad (\text{B.83})$$

$$\tau_\theta = v_\theta h \quad (\text{B.84})$$

We then obtain f_1 from matching the quadratic terms in Eqs (B.26) and Eqs (B.81)—

below we denote $f_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})$ by f_1 and $g_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1})$ by g_1 , for brevity:

$$\begin{aligned} v_\phi^2 h^2 f_1 + \frac{1}{2} v_\phi^2 h^2 \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} v_\phi^2 h^2 \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \\ = \frac{m-1}{2m} v_\phi^2 h^2 \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} \end{aligned} \quad (\text{B.85})$$

$$f_1 + \frac{1}{2} \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} = \frac{m-1}{2m} \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} \quad (\text{B.86})$$

$$f_1 = \left(\frac{m-1}{2m} - \frac{1}{2} \right) \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} - \frac{1}{2} \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \quad (\text{B.87})$$

$$f_1 = -\frac{1}{2m} \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} - \frac{1}{2} \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \quad (\text{B.88})$$

For g_1 , from Eqs (B.27) and (B.82):

$$\begin{aligned} v_\theta^2 h^2 \left(g_1 + \frac{1}{2} \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \right) \\ = v_\theta v_\phi h^2 \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} v_\theta^2 h^2 \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \end{aligned} \quad (\text{B.89})$$

$$\begin{aligned} g_1 + \frac{1}{2} \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \\ = \frac{v_\phi}{v_\theta} \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \end{aligned} \quad (\text{B.90})$$

$$\begin{aligned} g_1 + \frac{1}{2} \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{1}{2} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \\ = \frac{v_\phi}{v_\theta} \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \frac{k-1}{2k} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \end{aligned} \quad (\text{B.91})$$

$$g_1 = \left(\frac{v_\phi}{v_\theta} - \frac{1}{2} \right) \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} - \frac{1}{2k} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \quad (\text{B.92})$$

We thus have that

$$f_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) = -\frac{1}{2m} \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} - \frac{1}{2} \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} \quad (\text{B.93})$$

$$g_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) = \left(\frac{v_\phi}{v_\theta} - \frac{1}{2} \right) \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} - \frac{1}{2k} \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)}. \quad (\text{B.94})$$

We found f_1 and g_1 such that after one optimisation step the flows $\dot{\boldsymbol{\phi}} = f + v_\phi h f_1$

and $\dot{\boldsymbol{\theta}} = g + v_\theta h g_1$ follow the discrete updates up to order $\mathcal{O}(h^3)$, finishing the proof.

B.2 Proof of the main corollaries

In this section, we will write the modified equations in the case of using gradient descent common-payoff games and zero-sum games. To do so, we will replace f and g with the values given by gradient descent in the first order corrections we have derived in the previous sections. In the common-payoff case $f = -\nabla_\phi E$ and $g = -\nabla_\theta E$, while in the zero-sum case $f = \nabla_\phi E$ and $g = -\nabla_\theta E$, where $E(\boldsymbol{\phi}, \boldsymbol{\theta})$ is a function of the player parameters. We will use the following identities, which are obtained from the chain rule:

$$\frac{d\nabla_\phi E}{d\boldsymbol{\phi}} \nabla_\phi E = \nabla_\phi \left(\frac{\|\nabla_\phi E\|^2}{2} \right) \quad (\text{B.95})$$

$$\frac{d\nabla_\theta E}{d\boldsymbol{\theta}} \nabla_\theta E = \nabla_\theta \left(\frac{\|\nabla_\theta E\|^2}{2} \right) \quad (\text{B.96})$$

$$\frac{d\nabla_\phi E}{d\boldsymbol{\theta}} \nabla_\theta E = \left(\frac{d\nabla_\theta E}{d\boldsymbol{\phi}} \right)^T \nabla_\theta E = \nabla_\phi \left(\frac{\|\nabla_\theta E\|^2}{2} \right) \quad (\text{B.97})$$

$$\frac{d\nabla_\theta E}{d\boldsymbol{\phi}} \nabla_\phi E = \left(\frac{d\nabla_\phi E}{d\boldsymbol{\theta}} \right)^T \nabla_\phi E = \nabla_\theta \left(\frac{\|\nabla_\phi E\|^2}{2} \right). \quad (\text{B.98})$$

B.2.1 Common-payoff alternating two-player games

Corollary common-payoff alternating two-player games (*Corollary 4.4.1*)

In a two-player common-payoff game with common loss E , alternating gradient descent – as described in Eqs (B.35) and (B.36)—with one update per player follows a gradient flow given by the modified losses

$$\tilde{E}_\phi = E + \frac{v_\phi h}{4} (\|\nabla_\phi E\|^2 + \|\nabla_\theta E\|^2) \quad (\text{B.99})$$

$$\tilde{E}_\theta = E + \frac{v_\theta h}{4} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \|\nabla_\phi E\|^2 + \|\nabla_\theta E\|^2 \right) \quad (\text{B.100})$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

In the common-payoff case, both players minimise the loss function E . Substituting $f = -\nabla_\phi E$ and $g = -\nabla_\theta E$ into the corrections f_1 and g_1 for the alternating

Euler updates in Theorem 4.2.2 and using the gradient identities above yields

$$f_1 = -\frac{1}{2} \left(\frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \quad (\text{B.101})$$

$$= -\frac{1}{2} \left(\frac{1}{m} \frac{d\nabla_\phi E}{d\phi} \nabla_\phi E + \frac{d\nabla_\phi E}{d\theta} \nabla_\theta E \right) \quad (\text{def } f, g, \text{ B.102})$$

$$= -\frac{1}{2} \left(\frac{1}{m} \nabla_\phi \frac{\|\nabla_\phi E\|^2}{2} + \nabla_\phi \frac{\|\nabla_\theta E\|^2}{2} \right) \quad (\text{via (B.95),(B.97), B.103})$$

$$= -\nabla_\phi \left(\frac{1}{4m} \|\nabla_\phi E\|^2 + \frac{1}{4} \|\nabla_\theta E\|^2 \right) \quad (\text{B.104})$$

$$g_1 = -\frac{1}{2} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} f + \frac{1}{k} \frac{dg}{d\theta} g \right) \quad (\text{B.105})$$

$$= -\frac{1}{2} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{d\nabla_\theta E}{d\phi} \nabla_\phi E + \frac{1}{k} \frac{d\nabla_\theta E}{d\theta} \nabla_\theta E \right) \quad (\text{def } f, g, \text{ B.106})$$

$$= -\frac{1}{2} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \nabla_\theta \frac{\|\nabla_\phi E\|^2}{2} + \frac{1}{k} \nabla_\theta \frac{\|\nabla_\theta E\|^2}{2} \right) \quad (\text{via (B.96),(B.98), B.107})$$

$$= -\nabla_\theta \left(\frac{1}{4} \left(1 - \frac{2v_\phi}{v_\theta}\right) \|\nabla_\phi E\|^2 + \frac{1}{k} \|\nabla_\theta E\|^2 \right) \quad (\text{B.108})$$

Now, replacing the gradient expressions for f_1 and g_1 calculated above into the modified equations $\dot{\phi} = -\nabla_\phi E + v_\phi h f_1$ and $\dot{\theta} = -\nabla_\theta E + v_\theta h g_1$ and factoring out the gradients, we obtain the modified equations in the form of the flows:

$$\dot{\phi} = -\nabla_\phi \tilde{E}_\phi \quad (\text{B.109})$$

$$\dot{\theta} = -\nabla_\theta \tilde{E}_\theta, \quad (\text{B.110})$$

with the following modified losses for each players:

$$\tilde{E}_\phi = E + \frac{v_\phi h}{4} \left(\frac{1}{m} \|\nabla_\phi E\|^2 + \|\nabla_\theta E\|^2 \right) \quad (\text{B.111})$$

$$\tilde{E}_\theta = E + \frac{v_\theta h}{4} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \|\nabla_\phi E\|^2 + \frac{1}{k} \|\nabla_\theta E\|^2 \right). \quad (\text{B.112})$$

We obtain Corollary 5.1 by setting the number of player updates to one: $m = k = 1$.

B.2.2 Zero-sum simultaneous two-player games

Corollary Zero-sum simultaneous two-player games (Corollary 4.5.1) *In a zero-sum two-player differentiable game, simultaneous gradient descent updates - as described in Eqs (B.5) and (B.6) - follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = -E + \frac{v_\phi h}{4} \|\nabla_\phi E\|^2 - \frac{v_\phi h}{4} \|\nabla_\theta E\|^2 \quad (\text{B.113})$$

$$\tilde{E}_\theta = E - \frac{v_\theta h}{4} \|\nabla_\phi E\|^2 + \frac{v_\theta h}{4} \|\nabla_\theta E\|^2 \quad (\text{B.114})$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

Substituting $f = \nabla_\phi E$ and $g = -\nabla_\theta E$ into the corrections f_1 and g_1 for the simultaneous Euler updates and using the gradient identities above yields

$$f_1 = -\frac{1}{2} \left(\frac{df}{d\phi} f + \frac{df}{d\theta} g \right) \quad (\text{B.115})$$

$$= -\frac{1}{2} \left(\frac{d\nabla_\phi E}{d\phi} \nabla_\phi E - \frac{d\nabla_\phi E}{d\theta} \nabla_\theta E \right) \quad (\text{def } f, g, \text{ B.116})$$

$$= -\frac{1}{2} \left(\nabla_\phi \frac{\|\nabla_\phi E\|^2}{2} - \nabla_\phi \frac{\|\nabla_\theta E\|^2}{2} \right) \quad (\text{via (B.95),(B.97), B.117})$$

$$= -\nabla_\phi \left(\frac{1}{4} \|\nabla_\phi E\|^2 - \frac{1}{4} \|\nabla_\theta E\|^2 \right) \quad (\text{B.118})$$

$$g_1 = -\frac{1}{2} \left(\frac{dg}{d\phi} f + \frac{dg}{d\theta} g \right) \quad (\text{B.119})$$

$$= -\frac{1}{2} \left(-\frac{d\nabla_\theta E}{d\phi} \nabla_\phi E + \frac{d\nabla_\theta E}{d\theta} \nabla_\theta E \right) \quad (\text{def } f, g, \text{ B.120})$$

$$= -\frac{1}{2} \left(-\nabla_\theta \frac{\|\nabla_\phi E\|^2}{2} + \nabla_\theta \frac{\|\nabla_\theta E\|^2}{2} \right) \quad (\text{via (B.96),(B.98), B.121})$$

$$= -\nabla_\theta \left(-\frac{1}{4} \|\nabla_\phi E\|^2 + \frac{1}{4} \|\nabla_\theta E\|^2 \right) \quad (\text{B.122})$$

Now, replacing the gradient expressions for f_1 and g_1 calculated above into the modified equations $\dot{\phi} = -\nabla_\phi(-E) + v_\phi h f_1$ and $\dot{\theta} = -\nabla_\theta E + v_\theta h g_1$ and factoring out the gradients, we obtain the modified equations in the form of the flows:

$$\dot{\phi} = -\nabla_\phi \tilde{E}_\phi \quad (\text{B.123})$$

$$\dot{\theta} = -\nabla_\theta \tilde{E}_\theta, \quad (\text{B.124})$$

with the following modified losses for each players:

$$\tilde{E}_\phi = -E + \frac{v_\phi h}{4} \|\nabla_\phi E\|^2 - \frac{v_\phi h}{4} \|\nabla_\theta E\|^2 \quad (\text{B.125})$$

$$\tilde{E}_\theta = E - \frac{v_\theta h}{4} \|\nabla_\phi E\|^2 + \frac{v_\theta h}{4} \|\nabla_\theta E\|^2. \quad (\text{B.126})$$

B.2.3 Zero-sum alternating two-player games

Corollary Zero-sum alternating two-player games (*Corollary 4.5.2*) *In a zero-sum two-player differentiable game, alternating gradient descent - as described in Eqs (B.35) and (B.36) - follows a gradient flow given by the modified losses*

$$\tilde{E}_\phi = -E + \frac{v_\phi h}{4m} \|\nabla_\phi E\|^2 - \frac{v_\phi h}{4} \|\nabla_\theta E\|^2 \quad (\text{B.127})$$

$$\tilde{E}_\theta = E - \frac{v_\theta h}{4} \left(1 - \frac{2v_\phi}{v_\theta}\right) \|\nabla_\phi E\|^2 + \frac{v_\theta h}{4k} \|\nabla_\theta E\|^2 \quad (\text{B.128})$$

with an error of size $\mathcal{O}(h^3)$ after one update step.

In this last case, substituting $f = \nabla_\phi E$ and $g = -\nabla_\theta E$ into the corrections f_1 and g_1 for the alternating Euler updates and using the gradient identities above yields the modified system as well as the modified losses exactly in the same way as for the two previous cases above. (This amounts to a single sign change in the proof of Corollary 5.1)

B.2.4 Self and interaction terms in zero-sum games

Remark B.2.1 *Throughout the Appendix, we will refer to self terms and interaction terms, as originally defined in our thesis (Definition 4.2.1), and we will also use this terminology to refer to terms in our derivations that originate from the self terms and interaction terms. While a slight abuse of language, we find it useful to emphasise the provenance of these terms in our discussion.*

For the case of zero-sum games trained with simultaneous gradient descent, the self terms encourage the minimisation of the player's own gradient norm, while the

interaction terms encourage the maximisation of the other player’s gradient norm:

$$\tilde{E}_\phi = -E + \underbrace{\frac{v_\phi h}{4} \|\nabla_\phi E\|^2}_{self} - \underbrace{\frac{v_\phi h}{4} \|\nabla_\theta E\|^2}_{interaction}, \quad (\text{B.129})$$

$$\tilde{E}_\theta = E - \underbrace{\frac{v_\theta h}{4} \|\nabla_\phi E\|^2}_{interaction} + \underbrace{\frac{v_\theta h}{4} \|\nabla_\theta E\|^2}_{self}. \quad (\text{B.130})$$

Similar terms are obtained for alternating gradient descent, with the only difference that the sign of the second player’s interaction term can be positive, depending on the learning rate ratio of the two players.

B.3 Discretisation drift in Runge–Kutta 4

B.3.1 Runge–Kutta 4 for one player

Runge–Kutta 4 (RK4) is a Runge–Kutta a method of order 4. That is, the discrete steps of RK4 coincide with the dynamics of the flow they discretise up to $\mathcal{O}(h^5)$ (i.e., the local error after one step is of order $\mathcal{O}(h^5)$). The modified equation for a method of order n starts with corrections at order h^{n+1} (i.e., all the lower corrections vanish; see [179] and [324] for further details). This means that RK4 has no DD up to order $\mathcal{O}(h^5)$, and why for small learning rates RK4 can be used as a proxy for the exact flow.

B.3.2 Runge–Kutta 4 for two players

When we use equal learning rates and simultaneous updates, the two-player game is always equivalent to the one player case, so Runge–Kutta 4 will have a local error of $\mathcal{O}(h^5)$. However, in the case of two-players games, we have the additional freedom of having different learning rates for each of the players. We now show that when the learning rates of the two players are different, *RK4 also has a drift effect of order 2 and the DD term comes exclusively from the interaction terms*. To do so, we use BEA and apply the same steps we performed for simultaneous and alternating Euler updates.

Step 1: Expand the updates per player via a Taylor expansion.

The simultaneous Runge–Kutta 4 updates for two players require defining:

$$k_{1,\phi} = f(\phi_{t-1}, \theta_{t-1}) \quad (\text{B.131})$$

$$k_{1,\theta} = g(\phi_{t-1}, \theta_{t-1}) \quad (\text{B.132})$$

$$k_{2,\phi} = f\left(\phi_{t-1} + \frac{v_\phi h}{2} k_{1,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{1,\theta}\right) \quad (\text{B.133})$$

$$k_{2,\theta} = g\left(\phi_{t-1} + \frac{v_\phi h}{2} k_{1,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{1,\theta}\right) \quad (\text{B.134})$$

$$k_{3,\phi} = f\left(\phi_{t-1} + \frac{v_\phi h}{2} k_{2,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{2,\theta}\right) \quad (\text{B.135})$$

$$k_{3,\theta} = g\left(\phi_{t-1} + \frac{v_\phi h}{2} k_{2,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{2,\theta}\right) \quad (\text{B.136})$$

$$k_{4,\phi} = f\left(\phi_{t-1} + \frac{v_\phi h}{2} k_{3,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{3,\theta}\right) \quad (\text{B.137})$$

$$k_{4,\theta} = g\left(\phi_{t-1} + \frac{v_\phi h}{2} k_{3,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{3,\theta}\right) \quad (\text{B.138})$$

From here we can define the update vectors:

$$k_\phi = \frac{1}{6} (k_{1,\phi} + 2k_{2,\phi} + 2k_{3,\phi} + k_{4,\phi}) \quad (\text{B.139})$$

$$k_\theta = \frac{1}{6} (k_{1,\theta} + 2k_{2,\theta} + 2k_{3,\theta} + k_{4,\theta}) \quad (\text{B.140})$$

And the parameter updates

$$\phi_t = \phi_{t-1} + v_\phi h k_\phi \quad (\text{B.141})$$

$$\theta_t = \theta_{t-1} + v_\theta h k_\theta \quad (\text{B.142})$$

To find the discretisation drift we expand each intermediate step:

$$k_{1,\phi} = f_{(t-1)} \quad (\text{B.143})$$

$$k_{1,\theta} = g_{(t-1)} \quad (\text{B.144})$$

$$k_{2,\phi} = f \left(\phi_{t-1} + \frac{v_\phi h}{2} k_{1,\phi}, \theta_{t-1} + \frac{v_\theta h}{2} k_{1,\theta} \right) \quad (\text{B.145})$$

$$= f_{(t-1)} + \frac{v_\theta h}{2} \frac{df}{d\theta} \Big|_{(t-1)} k_{1,\theta} + \frac{v_\phi h}{2} \frac{df}{d\phi} \Big|_{(t-1)} k_{1,\phi} + \mathcal{O}(h^2) \quad (\text{B.146})$$

$$= f_{(t-1)} + \frac{v_\theta h}{2} \frac{df}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.147})$$

$$k_{2,\theta} = g \left(\phi_{t-1} + \frac{v_\phi h}{2} k_{1,\phi}, \theta + \frac{v_\theta h}{2} k_{1,\theta} \right) \quad (\text{B.148})$$

$$= g_{(t-1)} + \frac{v_\theta h}{2} \frac{dg}{d\theta} \Big|_{(t-1)} k_{1,\theta} + \frac{v_\phi h}{2} \frac{dg}{d\phi} \Big|_{(t-1)} k_{1,\phi} + \mathcal{O}(h^2) \quad (\text{B.149})$$

$$= g_{(t-1)} + \frac{v_\theta h}{2} \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.150})$$

$$k_{3,\phi} = f_{(t-1)} + \frac{v_\theta h}{2} \frac{df}{d\theta} \Big|_{(t-1)} k_{2,\theta} + \frac{v_\phi h}{2} \frac{df}{d\phi} \Big|_{(t-1)} k_{2,\phi} + \mathcal{O}(h^2) \quad (\text{B.151})$$

$$= f_{(t-1)} + \frac{v_\theta h}{2} \frac{df}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.152})$$

$$k_{3,\theta} = g_{(t-1)} + \frac{v_\theta h}{2} \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.153})$$

$$k_{4,\phi} = f_{(t-1)} + \frac{v_\theta h}{2} \frac{df}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.154})$$

$$k_{4,\theta} = g_{(t-1)} + \frac{v_\theta h}{2} \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.155})$$

with the update direction:

$$k_\phi = \frac{1}{6} (k_{1,\phi} + 2k_{2,\phi} + 2k_{3,\phi} + k_{4,\phi}) \quad (\text{B.156})$$

$$= f_{(t-1)} + \frac{v_\theta h}{2} \frac{df}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{df}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.157})$$

$$k_\theta = \frac{1}{6} (k_{1,\theta} + 2k_{2,\theta} + 2k_{3,\theta} + k_{4,\theta}) \quad (\text{B.158})$$

$$= g_{(t-1)} + \frac{v_\theta h}{2} \frac{dg}{d\theta} \Big|_{(t-1)} g_{(t-1)} + \frac{v_\phi h}{2} \frac{dg}{d\phi} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^2) \quad (\text{B.159})$$

and thus the discrete update of the Runge–Kutta 4 for two players are:

$$\boldsymbol{\phi}_t = \boldsymbol{\phi}_{t-1} + v_\phi h f_{(t-1)} + \frac{1}{2} v_\phi v_\theta h^2 \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} + \frac{1}{2} v_\phi^2 h^2 \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3)$$
(B.160)

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + v_\theta h g_{(t-1)} + \frac{1}{2} v_\theta^2 h^2 \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} + \frac{1}{2} v_\phi v_\theta h^2 \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} + \mathcal{O}(h^3)$$
(B.161)

Step 2: Expand the continuous-time changes for the modified flow

(Identical to the simultaneous Euler updates.)

Step 3: Matching the discrete and modified continuous updates.

As in the always in Step 3, we substitute $\boldsymbol{\phi}_{t-1}$ and $\boldsymbol{\theta}_{t-1}$ in Lemma B.1.1:

$$\begin{aligned} \boldsymbol{\phi}(\tau + \tau_\phi) &= \boldsymbol{\phi}_{t-1} + \tau_\phi f_{(t-1)} + \tau_\phi^2 f_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) + \frac{1}{2} \tau_\phi^2 \frac{df}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} \\ &\quad + \frac{1}{2} \tau_\phi^2 \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} + \mathcal{O}(\tau_\phi^3) \end{aligned}$$
(B.162)

$$\begin{aligned} \boldsymbol{\theta}(\tau + \tau_\theta) &= \boldsymbol{\theta}_{t-1} + \tau_\theta g_{(t-1)} + \tau_\theta^2 g_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) + \frac{1}{2} \tau_\theta^2 \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)} \\ &\quad + \frac{1}{2} \tau_\theta^2 \frac{dg}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)} + \mathcal{O}(\tau_\theta^3) \end{aligned}$$
(B.163)

For the first order terms we obtain $\tau_\phi = v_\phi h$ and $\tau_\theta = v_\theta h$. We match the $\mathcal{O}(h^2)$ terms in the equations above with the discrete Runge–Kutta 4 updates shown in Eq (B.160) and (B.161) and notice that:

$$f_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) = \frac{1}{2} \left(\frac{v_\theta}{v_\phi} - 1 \right) \frac{df}{d\boldsymbol{\theta}} \Big|_{(t-1)} g_{(t-1)}$$
(B.164)

$$g_1(\boldsymbol{\phi}_{t-1}, \boldsymbol{\theta}_{t-1}) = \frac{1}{2} \left(\frac{v_\phi}{v_\theta} - 1 \right) \frac{dg}{d\boldsymbol{\phi}} \Big|_{(t-1)} f_{(t-1)}$$
(B.165)

Thus, if $v_\phi \neq v_\theta$ RK4 has second order drift. This is why, in all our experiments comparing with RK4, we use the same learning rates for the two players $v_\phi = v_\theta$, to ensure that we use a method which has no DD up to order $\mathcal{O}(h^5)$.

B.4 Stability analysis

In this section, we give all the details of the stability analysis results, to showcase how the modified flows we have derived can be used as a tool for stability analysis. We provide the full computation for the Jacobian of the modified vector fields for simultaneous and alternating Euler updates, as well as the calculation of their trace, and show how this can be used to determine the stability of the modified vector fields. While analysing the modified vector fields is not equivalent to analysing the discrete dynamics due to the higher order errors of the drift which we ignore, it provides a better approximation than what is often used in practice, namely the original flows, which completely ignore the drift.

B.4.1 Simultaneous Euler updates

Consider a two-player game with dynamics given by $\phi_t = \phi_{t-1} + v_\phi h f_{(t-1)}$ and $\theta_t = \theta_{t-1} + v_\theta h g_{(t-1)}$ (Eqs (B.5) and (B.6)). The modified dynamics for this game are given by $\dot{\phi} = \tilde{f}$, $\dot{\theta} = \tilde{g}$, where $\tilde{f} = f - \frac{v_\phi h}{2} \left(\frac{df}{d\phi} f + \frac{df}{d\theta} g \right)$ and $\tilde{g} = g - \frac{v_\theta h}{2} \left(\frac{dg}{d\phi} f + \frac{dg}{d\theta} g \right)$ (Theorem (DD for simultaneous Euler updates)).

The stability of this system can be characterised by the modified Jacobian matrix evaluated at the equilibria of the two-player game. The equilibria that we are interested in for our stability analysis are the steady-state solutions of Eqs (B.5) and (B.6), given by $f = \mathbf{0}$, $g = \mathbf{0}$. These are also equilibrium solutions for the steady-state modified equations¹ given by $\tilde{f} = \mathbf{0}$, $\tilde{g} = \mathbf{0}$.

The modified Jacobian can be written, using block matrix notation as:

$$\tilde{\mathbf{J}} = \begin{bmatrix} \frac{d\tilde{f}}{d\phi} & \frac{d\tilde{f}}{d\theta} \\ \frac{d\tilde{g}}{d\phi} & \frac{d\tilde{g}}{d\theta} \end{bmatrix} \quad (\text{B.166})$$

Next, we calculate each term in this block matrix. (In the following analysis, each

¹There are additional steady-state solutions for the modified equations. However, we can ignore these since they are spurious solutions that do not correspond to steady states of the two-player game, arising instead as an artifact of the $\mathcal{O}(h^3)$ error in BEA.

term is evaluated at an equilibrium solution given by $f = \mathbf{0}, g = \mathbf{0}$). We find:

$$\frac{d\tilde{f}}{d\phi} = \frac{df}{d\phi} - \frac{v_\phi h}{2} \left(\left(\frac{df}{d\phi} \right)^2 + \frac{df}{d\phi d\phi} f + \frac{df}{d\theta} \frac{dg}{d\phi} + \frac{df}{d\phi d\theta} g \right) \quad (\text{B.167})$$

$$= \frac{df}{d\phi} - \frac{v_\phi h}{2} \left(\left(\frac{df}{d\phi} \right)^2 + \frac{df}{d\theta} \frac{dg}{d\phi} \right) \quad (\text{B.168})$$

$$\frac{d\tilde{f}}{d\theta} = \frac{df}{d\theta} - \frac{v_\phi h}{2} \left(\frac{df}{d\phi} \frac{df}{d\theta} + \frac{df}{d\theta d\phi} f + \frac{df}{d\theta} \frac{dg}{d\theta} + \frac{df}{d\theta d\theta} g \right) \quad (\text{B.169})$$

$$= \frac{df}{d\theta} - \frac{v_\phi h}{2} \left(\frac{df}{d\phi} \frac{df}{d\theta} + \frac{df}{d\theta} \frac{dg}{d\theta} \right) \quad (\text{B.170})$$

$$\frac{d\tilde{g}}{d\phi} = \frac{dg}{d\phi} - \frac{v_\theta h}{2} \left(\frac{dg}{d\theta} \frac{dg}{d\phi} + \frac{dg}{d\phi d\theta} g + \frac{dg}{d\phi} \frac{df}{d\phi} + \frac{dg}{d\phi d\phi} f \right) \quad (\text{B.171})$$

$$= \frac{dg}{d\phi} - \frac{v_\theta h}{2} \left(\frac{dg}{d\theta} \frac{dg}{d\phi} + \frac{dg}{d\phi} \frac{df}{d\phi} \right) \quad (\text{B.172})$$

$$\frac{d\tilde{g}}{d\theta} = \frac{dg}{d\theta} - \frac{v_\theta h}{2} \left(\left(\frac{dg}{d\theta} \right)^2 + \frac{dg}{d\theta d\theta} g + \frac{dg}{d\phi} \frac{df}{d\theta} + \frac{dg}{d\theta d\phi} f \right) \quad (\text{B.173})$$

$$= \frac{dg}{d\theta} - \frac{v_\theta h}{2} \left(\left(\frac{dg}{d\theta} \right)^2 + \frac{dg}{d\phi} \frac{df}{d\theta} \right) \quad (\text{B.174})$$

where $\frac{df}{d\phi d\theta} g$ is the matrix in $\mathbb{R}^{m,m}$ with $\left(\frac{df}{d\phi d\theta} g \right)_{i,j} = \left(\frac{d^2 f_i}{d\phi_j d\theta} \right)^T g$.

Given these calculations, we can now write

$$\tilde{\mathbf{J}} = \begin{bmatrix} \frac{d\tilde{f}}{d\phi} & \frac{d\tilde{f}}{d\theta} \\ \frac{d\tilde{g}}{d\phi} & \frac{d\tilde{g}}{d\theta} \end{bmatrix} = \mathbf{J} - \frac{h}{2} \mathbf{K}_{\text{sim}} \quad (\text{B.175})$$

where \mathbf{J} is the Jacobian of the unmodified flow

$$\mathbf{J} = \begin{bmatrix} \frac{df}{d\phi} & \frac{df}{d\theta} \\ \frac{dg}{d\phi} & \frac{dg}{d\theta} \end{bmatrix} \quad (\text{B.176})$$

and

$$\mathbf{K}_{\text{sim}} = \begin{bmatrix} v_\phi \left(\frac{df}{d\phi} \right)^2 + v_\phi \frac{df}{d\theta} \frac{dg}{d\phi} & v_\phi \frac{df}{d\phi} \frac{df}{d\theta} + v_\phi \frac{df}{d\theta} \frac{dg}{d\theta} \\ v_\theta \frac{dg}{d\theta} \frac{dg}{d\phi} + v_\theta \frac{dg}{d\phi} \frac{df}{d\theta} & v_\theta \left(\frac{dg}{d\theta} \right)^2 + v_\theta \frac{dg}{d\phi} \frac{df}{d\theta} \end{bmatrix}. \quad (\text{B.177})$$

Using the modified Jacobian, we can use Remark 2.3.1 to determine the stability of fixed points. A necessary condition for stability is that the trace of the modified

Jacobian is less than or equal to zero (i.e. $\text{Tr}(\tilde{\mathbf{J}}) \leq 0$), since the trace is the sum of eigenvalues. Using the property of trace additivity and the trace cyclic property we see that:

$$\text{Tr}(\tilde{\mathbf{J}}) = \text{Tr}(\mathbf{J}) - \frac{h}{2} \left(v_\phi \text{Tr}\left(\left(\frac{df}{d\phi}\right)^2\right) + v_\theta \text{Tr}\left(\left(\frac{dg}{d\theta}\right)^2\right) \right) - \frac{h}{2}(v_\phi + v_\theta) \text{Tr}\left(\frac{df}{d\theta} \frac{dg}{d\phi}\right). \quad (\text{B.178})$$

Instability caused by discretisation drift. We now use the above analysis to show that the equilibria of a two-player game following the modified flow obtained simultaneous Euler updates as defined by Eqs (B.5) and (B.6) can become asymptotically unstable for some games.

There are choices of f and g that have stable equilibrium without DD, but are unstable under DD. For example, consider the zero-sum two-player game with $f = \nabla_\phi E$ and $g = -\nabla_\theta E$. In this case $\frac{df}{d\phi} = \nabla_\phi^2 E$, $\frac{dg}{d\theta} = -\nabla_\theta^2 E$ and $\frac{df}{d\theta} = \frac{dg}{d\phi} = \nabla_{\phi,\theta} E$. We can thus use $\text{Tr}(A^T A) = \|A\|_2$ where $\|\cdot\|_2$ denotes the Frobenius norm, we have

$$\text{Tr}(\tilde{\mathbf{J}}) = \text{Tr}(\mathbf{J}) - \frac{h}{2} \left(v_\phi \|\nabla_\phi^2 E\|_2^2 + v_\theta \|\nabla_\theta^2 E\|_2^2 \right) + \frac{h}{2}(v_\phi + v_\theta) \|\nabla_{\phi,\theta} E\|_2^2. \quad (\text{B.179})$$

The Dirac-GAN is an example of a zero-sum two-player game that is stable without DD, but becomes unstable under DD with $\text{Tr}(\tilde{\mathbf{J}}) = h(v_\phi + v_\theta) \|\nabla_{\phi,\theta} E\|_2^2 / 2 > 0$ (see Section B.6.2).

B.4.2 Alternating Euler updates

Consider a two-player game with dynamics given by Eqs (B.35) and (B.36). The modified dynamics for this game are given by $\dot{\phi} = f - \frac{v_\phi h}{2} \left(\frac{1}{m} \frac{df}{d\phi} f + \frac{df}{d\theta} g \right)$, $\dot{\theta} = g - \frac{v_\theta h}{2} \left(\left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} f + \frac{1}{k} \frac{dg}{d\theta} g \right)$ (Theorem DD for alternating Euler updates).

The stability of this system can be characterised by the modified Jacobian matrix evaluated at the equilibria of the two-player game. The equilibria that we are interested in for our stability analysis are the steady-state solutions of Eqs (B.5) and (B.6), given by $f = \mathbf{0}, g = \mathbf{0}$. These are also equilibrium solutions for the steady-state modified equations² given by $\tilde{f} = \mathbf{0}, \tilde{g} = \mathbf{0}$.

²There are additional steady-state solutions for the modified equations. However, we can ignore these since they are spurious solutions that do not correspond to steady states of the two-player

The modified Jacobian can be written, using block matrix notation as:

$$\tilde{\mathbf{J}} = \begin{bmatrix} \frac{d\tilde{f}}{d\phi} & \frac{d\tilde{f}}{d\theta} \\ \frac{d\tilde{g}}{d\phi} & \frac{d\tilde{g}}{d\theta} \end{bmatrix} \quad (\text{B.180})$$

Next, we calculate each term in this block matrix. (In the following analysis, each term is evaluated at an equilibrium solution given by $f = \mathbf{0}, g = \mathbf{0}$). We find:

$$\frac{d\tilde{f}}{d\phi} = \frac{df}{d\phi} - \frac{v_\phi h}{2} \left(\frac{1}{m} \left(\frac{df}{d\phi} \right)^2 + \frac{1}{m} \frac{df}{d\phi d\theta} f + \frac{df}{d\theta} \frac{dg}{d\phi} + \frac{df}{d\phi d\theta} g \right) \quad (\text{B.181})$$

$$= \frac{df}{d\phi} - \frac{v_\phi h}{2} \left(\frac{1}{m} \left(\frac{df}{d\phi} \right)^2 + \frac{df}{d\theta} \frac{dg}{d\phi} \right) \quad (\text{B.182})$$

$$\frac{d\tilde{f}}{d\theta} = \frac{df}{d\theta} - \frac{v_\phi h}{2} \left(\frac{1}{m} \frac{df}{d\theta} \frac{df}{d\phi} + \frac{1}{m} \frac{df}{d\theta d\phi} f + \frac{df}{d\theta} \frac{dg}{d\theta} + \frac{df}{d\theta d\theta} g \right) \quad (\text{B.183})$$

$$= \frac{df}{d\theta} - \frac{v_\phi h}{2} \left(\frac{1}{m} \frac{df}{d\phi} \frac{df}{d\theta} + \frac{df}{d\theta} \frac{dg}{d\theta} \right) \quad (\text{B.184})$$

$$\frac{d\tilde{g}}{d\phi} = \frac{dg}{d\phi} - \frac{v_\theta h}{2} \left(\frac{1}{k} \frac{dg}{d\theta} \frac{dg}{d\phi} + \frac{1}{k} \frac{dg}{d\phi d\theta} g + \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\phi} + \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi d\phi} f \right) \quad (\text{B.185})$$

$$= \frac{dg}{d\phi} - \frac{v_\theta h}{2} \left(\frac{1}{k} \frac{dg}{d\theta} \frac{dg}{d\phi} + \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\phi} \right) \quad (\text{B.186})$$

$$\frac{d\tilde{g}}{d\theta} = \frac{dg}{d\theta} - \frac{v_\theta h}{2} \left(\frac{1}{k} \left(\frac{dg}{d\theta} \right)^2 + \frac{1}{k} \frac{dg}{d\theta d\theta} g + \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\theta} + \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\theta d\phi} f \right) \quad (\text{B.187})$$

$$= \frac{dg}{d\theta} - \frac{v_\theta h}{2} \left(\frac{1}{k} \left(\frac{dg}{d\theta} \right)^2 + \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\theta} \right) \quad (\text{B.188})$$

where $\frac{df}{d\phi d\theta} g$ is the matrix in $\mathbb{R}^{m,m}$ with $\left(\frac{df}{d\phi d\theta} g \right)_{i,j} = \left(\frac{d^2 f_i}{d\phi_j d\theta} \right)^T g$.

Given these calculations, we can now write:

$$\tilde{\mathbf{J}} = \begin{bmatrix} \frac{d\tilde{f}}{d\phi} & \frac{d\tilde{f}}{d\theta} \\ \frac{d\tilde{g}}{d\phi} & \frac{d\tilde{g}}{d\theta} \end{bmatrix} = \mathbf{J} - \frac{h}{2} \mathbf{K}_{\text{alt}} \quad (\text{B.189})$$

where \mathbf{J} is the Jacobian of the unmodified flow:

$$\mathbf{J} = \begin{bmatrix} \frac{df}{d\phi} & \frac{df}{d\theta} \\ \frac{dg}{d\phi} & \frac{dg}{d\theta} \end{bmatrix} \quad (\text{B.190})$$

and

$$\mathbf{K}_{\text{alt}} = \begin{bmatrix} \frac{v_\phi}{m} \left(\frac{df}{d\phi}\right)^2 + v_\phi \frac{df}{d\theta} \frac{dg}{d\phi} & \frac{v_\phi}{m} \frac{df}{d\phi} \frac{df}{d\theta} + v_\phi \frac{df}{d\theta} \frac{dg}{d\theta} \\ \frac{v_\theta}{k} \frac{dg}{d\theta} \frac{dg}{d\phi} + v_\theta \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\phi} & \frac{v_\theta}{k} \left(\frac{dg}{d\theta}\right)^2 + v_\theta \left(1 - \frac{2v_\phi}{v_\theta}\right) \frac{dg}{d\phi} \frac{df}{d\theta} \end{bmatrix} \quad (\text{B.191})$$

Using the modified Jacobian, we can use Remark 2.3.1 to determine the stability of fixed points. A necessary condition for stability is that the trace of the modified Jacobian is less than or equal to zero (i.e. $\text{Tr}(\tilde{\mathbf{J}}) \leq 0$), since the trace is the sum of eigenvalues. Using the property of trace additivity and the trace cyclic property we see that:

$$\text{Tr}(\tilde{\mathbf{J}}) = \text{Tr}(\mathbf{J}) - \frac{h}{2} \left(\frac{v_\phi}{m} \text{Tr}\left(\left(\frac{df}{d\phi}\right)^2\right) + \frac{v_\theta}{k} \text{Tr}\left(\left(\frac{dg}{d\theta}\right)^2\right) \right) - \frac{h}{2} (v_\theta - v_\phi) \text{Tr}\left(\frac{df}{d\theta} \frac{dg}{d\phi}\right). \quad (\text{B.192})$$

We note that unlike for simultaneous updates, even if $\text{Tr}\left(\frac{df}{d\theta} \frac{dg}{d\phi}\right)$ is negative, if $v_\theta < v_\phi$, the trace of the modified system will stay negative, so the necessary condition for the system to remain stable is still satisfied. However, since this is not a sufficient condition, the modified system could still be unstable.

B.5 SGA in zero-sum games

For clarity, this section reproduces Symplectic Gradient Adjustment (SGA) from Balduzzi et al. [61] for zero-sum games, to showcase the appearance of interaction terms. We have two players ϕ and θ minimizing loss functions E and $-E$, respectively, then SGA defines the vector

$$\boldsymbol{\xi} = \begin{bmatrix} \nabla_\phi E \\ -\nabla_\theta E \end{bmatrix}. \quad (\text{B.193})$$

They then define \mathbf{H}_ξ the game Hessian (their Section 2.2)

$$\mathbf{H}_\xi = \begin{bmatrix} \frac{dE}{d\phi d\phi} & \frac{dE}{d\phi d\theta} \\ -\frac{dE}{d\theta d\phi} & -\frac{dE}{d\theta d\theta} \end{bmatrix}. \quad (\text{B.194})$$

Thus \mathbf{H}_ξ has the *anti-symmetric* component

$$\mathbf{A} = \frac{1}{2}(\mathbf{H}_\xi - \mathbf{H}_\xi^T) = \frac{1}{2} \begin{bmatrix} 0 & \frac{dE}{d\phi d\theta} + \frac{dE}{d\theta d\phi}^T \\ -\frac{dE}{d\theta d\phi} - \frac{dE}{d\phi d\theta}^T & 0 \end{bmatrix} = \begin{bmatrix} 0 & \frac{dE}{d\phi d\theta} \\ -\frac{dE}{d\theta d\phi} & 0 \end{bmatrix}. \quad (\text{B.195})$$

The vector field ξ is modified according to SGA as

$$\hat{\xi} = \xi + \zeta \mathbf{A}^T \xi = \begin{bmatrix} \nabla_\phi E \\ -\nabla_\theta E \end{bmatrix} + \zeta \begin{bmatrix} 0 & \frac{dE}{d\phi d\theta} \\ -\frac{dE}{d\theta d\phi} & 0 \end{bmatrix} \begin{bmatrix} \nabla_\phi E \\ -\nabla_\theta E \end{bmatrix} \quad (\text{B.196})$$

$$= \begin{bmatrix} \nabla_\phi E \\ -\nabla_\theta E \end{bmatrix} + \zeta \begin{bmatrix} -\frac{dE}{d\phi d\theta} \nabla_\theta E \\ -\frac{dE}{d\theta d\phi} \nabla_\phi E \end{bmatrix}, \quad (\text{B.197})$$

where ζ is either a hyperparameter or adjusted dynamically in training [61]. Thus, via Eqs (B.98) and (B.97) the modified vector field can be simplified to

$$\hat{\xi} = \begin{bmatrix} \nabla_\phi (E - \zeta \frac{1}{2} \|\nabla_\theta E\|^2) \\ \nabla_\theta (-E - \zeta \frac{1}{2} \|\nabla_\phi E\|^2) \end{bmatrix}. \quad (\text{B.198})$$

Therefore, since $\hat{\xi}$ defines the negative of the vector field followed by the system, the modified losses for the two players can be written respectively as

$$\tilde{L}_1 = E - \frac{\zeta}{2} \|\nabla_\theta E\|^2 \quad (\text{B.199})$$

$$\tilde{L}_2 = -E - \frac{\zeta}{2} \|\nabla_\phi E\|^2. \quad (\text{B.200})$$

Thus if $\zeta < 0$, which is what we use in our experiments, the functional form of the modified losses given by SGA is the same used to cancel the interaction terms of DD in the case of simultaneous gradient descent updates in zero sum games. We do

however highlight a few differences in our approach compared to SGA: our approach extends to alternating updates and provides the optimal regularisation coefficients; cancelling the interaction terms of the drift is different compared to SGA for general games.

B.6 DiracGAN - an illustrative example

In their work assessing the convergence of GAN training Mescheder et al. [68] introduce the example of the DiracGAN, where the GAN is trying to learn a delta distribution with mass at zero. More specifically, the generator $G(z; \theta) = \theta$ with parameter θ parametrises constant functions whose images $\{\theta\}$ correspond to the support of the delta distribution δ_θ . The discriminator is a linear model $D(x; \phi) = \phi \cdot x$ with parameter ϕ . The loss function is given by

$$E(\theta, \phi) = l(\theta\phi) + l(0), \quad (\text{B.201})$$

where l depends on the GAN used; for the standard GAN it is $l = -\log(1 + e^{-t})$. As in Mescheder et al. [68], we assume l is continuously differentiable with $l'(x) \neq 0$ for all $x \in \mathbb{R}$. The partial derivatives of the loss function

$$\frac{\partial E}{\partial \phi} = l'(\theta\phi)\theta, \quad \frac{\partial E}{\partial \theta} = l'(\theta\phi)\phi, \quad (\text{B.202})$$

lead to the underlying continuous dynamics

$$\dot{\phi} = f(\theta, \phi) = l'(\theta\phi)\theta, \quad \dot{\theta} = g(\theta, \phi) = -l'(\theta\phi)\phi. \quad (\text{B.203})$$

Thus the only equilibrium of the game is $\theta = 0$ and $\phi = 0$.

B.6.1 Reconciling discrete and continuous updates in Dirac-GAN

Mescheder et al. [68] observed a discrepancy between the continuous and discrete dynamics of DiracGAN. They show that, for the problem in Eq (B.201), the continuous dynamics preserve $\theta^2 + \phi^2$, and thus cannot converge (Lemma 2.3 in Mescheder

et al. [68]), since

$$\frac{d(\theta^2 + \phi^2)}{dt} = 2\theta \frac{d\theta}{dt} + 2\phi \frac{d\phi}{dt} = -2\theta l'(\theta\phi)\phi + 2\phi l'(\theta\phi)\theta = 0. \quad (\text{B.204})$$

They also observe that when following the discrete dynamics of simultaneous gradient descent that $\theta^2 + \phi^2$ increases in time (Lemma 2.4 in Mescheder et al. [68]). We resolve this discrepancy between conclusions reached with continuous-time and discrete-time analysis in DiracGAN, by showing that the modified continuous dynamics given by DD result in behaviour consistent with that of the discrete updates.

Proposition B.6.1 *Following the continuous-time modified flow we obtained using BEA describing simultaneous gradient descent increases $\theta^2 + \phi^2$ in DiracGAN.*

Proof: We assume both the generator and the discriminator use learning rates h , as in Mescheder et al. [68]. We first compute terms used by the zero-sum Corollaries 4.5.1 and 4.5.2 in the main thesis.

$$\|\nabla_{\theta} E\|^2 = l'(\theta\phi)^2 \phi^2 \quad (\text{B.205})$$

$$\|\nabla_{\phi} E\|^2 = l'(\theta\phi)^2 \theta^2 \quad (\text{B.206})$$

and

$$\nabla_{\theta} \|\nabla_{\theta} E\|^2 = 2\phi^3 l'(\theta\phi) l''(\theta\phi) \quad (\text{B.207})$$

$$\nabla_{\theta} \|\nabla_{\phi} E\|^2 = 2\theta l'(\theta\phi)^2 + 2\theta^2 \phi l'(\theta\phi) l''(\theta\phi) \quad (\text{B.208})$$

$$\nabla_{\phi} \|\nabla_{\theta} E\|^2 = 2\phi l'(\theta\phi)^2 + 2\phi^2 \theta l'(\theta\phi) l''(\theta\phi) \quad (\text{B.209})$$

$$\nabla_{\phi} \|\nabla_{\phi} E\|^2 = 2\theta^3 l'(\theta\phi) l''(\theta\phi). \quad (\text{B.210})$$

Thus, the modified flows are given by:

$$\dot{\phi} = l'(\theta\phi)\theta + \frac{h}{2} [-\theta^3 l'(\theta\phi) l''(\theta\phi) + \phi l'(\theta\phi)^2 + \phi^2 \theta l'(\theta\phi) l''(\theta\phi)] \quad (\text{B.211})$$

$$\dot{\theta} = -l'(\theta\phi)\phi - \frac{h}{2} [-\theta l'(\theta\phi)^2 - \theta^2 \phi l'(\theta\phi) l''(\theta\phi) + \phi^3 l'(\theta\phi) l''(\theta\phi)]. \quad (\text{B.212})$$

By denoting $l'(\theta\phi)$ by l' and $l''(\theta\phi)$ by l'' , then we have:

$$\frac{d(\theta^2 + \phi^2)}{dt} = 2\theta \frac{d\theta}{dt} + 2\phi \frac{d\phi}{dt} \quad (\text{B.213})$$

$$\begin{aligned} &= 2\theta \left(-l'\phi - \frac{h}{2} [-\theta l'^2 - \theta^2 \phi l'l'' + \phi^3 l'l''] \right) \\ &\quad + 2\phi \left(l'\theta + \frac{h}{2} [-\theta^3 l'l'' + \phi l'^2 + \phi^2 \theta l'l''] \right) \end{aligned} \quad (\text{B.214})$$

$$= 2\theta \left(-\frac{h}{2} [-\theta l'^2 - \theta^2 \phi l'l'' + \phi^3 l'l''] \right) + 2\phi \left(\frac{h}{2} [-\theta^3 l'l'' + \phi l'^2 + \phi^2 \theta l'l''] \right) \quad (\text{B.215})$$

$$= h\theta^2 l'^2 + h\theta^3 \phi l'l'' - h\phi^3 \theta l'l'' - h\phi \theta^3 l'l'' + h\phi^2 l'^2 + h\phi^3 \theta l'l'' \quad (\text{B.216})$$

$$= h\theta^2 l'^2 + h\phi^2 l'^2 > 0 \quad (\text{B.217})$$

for all $\phi, \theta \neq 0$, which shows that $\theta^2 + \phi^2$ is not preserved and it will strictly increase for all values away from the equilibrium (we have used the assumption that $l'(x) \neq 0, \forall x \in \mathbb{R}$). \square

We have thus identified a continuous system which exhibits the same behaviour as described by Lemma 2.4 in Mescheder et al. [68], where a discrete system is analysed.

B.6.2 DD changes the convergence behaviour of Dirac-GAN

The Jacobian of the unmodified Dirac-GAN evaluated at the equilibrium solution given by $\phi = 0, \theta = 0$. Since each parameter is one dimensional, we can write

$$\mathbf{J} = \begin{bmatrix} \nabla_{\phi, \phi} E & \nabla_{\theta, \phi} E \\ -\nabla_{\phi, \theta} E & -\nabla_{\theta, \theta} E \end{bmatrix} = \begin{bmatrix} 0 & l'(0) \\ -l'(0) & 0 \end{bmatrix} \quad (\text{B.218})$$

We see that $\text{Tr}(\mathbf{J}) = 0$ and the determinant $|\mathbf{J}| = l'(0)^2$. Therefore, the eigenvalues of this Jacobian are $v_{\theta\pm} = \text{Tr}(\mathbf{J})/2 \pm \sqrt{\text{Tr}(\mathbf{J})^2 - 4|\mathbf{J}|}/2 = \pm il'(0)$ (Reproduced from Mescheder et al. [68]). This is an example of a stable *center equilibrium*.

Next we calculate the Jacobian of the modified flows given by DiracGAN,

evaluated at an equilibrium solution and find $\tilde{\mathbf{J}} = \mathbf{J} - h\Delta/2$, where

$$\Delta = \begin{bmatrix} v_\phi(\nabla_{\phi,\phi}E)^2 - v_\phi\nabla_{\phi,\theta}E\nabla_{\theta,\phi}E & v_\phi\nabla_{\theta,\phi}E\nabla_{\phi,\phi}E - v_\phi\nabla_{\theta,\theta}E\nabla_{\theta,\phi}E \\ v_\theta\nabla_{\phi,\theta}E\nabla_{\theta,\theta}E - v_\theta\nabla_{\phi,\phi}E\nabla_{\phi,\theta}E & v_\theta(\nabla_{\theta,\theta}E)^2 - v_\theta\nabla_{\theta,\phi}E\nabla_{\phi,\theta}E \end{bmatrix} \quad (\text{B.219})$$

$$= \begin{bmatrix} -v_\phi\nabla_{\phi,\theta}E\nabla_{\theta,\phi}E & 0 \\ 0 & -v_\theta\nabla_{\theta,\phi}E\nabla_{\phi,\theta}E \end{bmatrix} \quad (\text{B.220})$$

$$= \begin{bmatrix} -v_\phi l'(0)^2 & 0 \\ 0 & -v_\theta l'(0)^2 \end{bmatrix} \quad (\text{B.221})$$

so

$$\tilde{\mathbf{J}} = \begin{bmatrix} hv_\phi l'(0)^2/2 & l'(0) \\ -l'(0) & hv_\theta l'(0)^2/2 \end{bmatrix}. \quad (\text{B.222})$$

Now, we see that the trace of the modified Jacobian for the Dirac-GAN is $\text{Tr}(\tilde{\mathbf{J}}) = (h/2)(v_\phi + v_\theta)l'(0)^2 > 0$, so the modified flows induced by gradient descent in DiracGAN are unstable.

B.6.3 Explicit regularisation stabilises Dirac-GAN

Here, we show that we can use our stability analysis to identify forms of explicit regularisation that can counteract the destabilizing impact of DD. Consider the Dirac-GAN with explicit regularisation of the following form: $E_\phi = -E + \gamma\|\nabla_\theta E\|^2$ and $E_\theta = E + \zeta\|\nabla_\phi E\|^2$ where $\phi_t = \phi_{t-1} - v_\phi h\nabla_\phi E_\phi$ and $\theta_t = \theta_{t-1} - v_\theta h\nabla_\theta E_\theta$ and with $\gamma, \zeta \sim \mathcal{O}(h)$. The modified Jacobian for this system is given by

$$\tilde{\mathbf{J}} = \begin{bmatrix} hv_\phi/2\nabla_{\phi,\theta}E\nabla_{\theta,\phi}E - \gamma\nabla_{\phi,\phi}\|\nabla_\theta E\|^2 & \nabla_{\theta,\phi}E \\ -\nabla_{\phi,\theta}E & hv_\theta/2\nabla_{\theta,\phi}E\nabla_{\phi,\theta}E - \zeta\nabla_{\theta,\theta}\|\nabla_\phi E\|^2 \end{bmatrix} \quad (\text{B.223})$$

$$= \begin{bmatrix} (hv_\phi/2 - 2\gamma)l'(0)^2 & l'(0) \\ -l'(0) & (hv_\theta/2 - 2\zeta)l'(0)^2 \end{bmatrix} \quad (\text{B.224})$$

The determinant of the modified Jacobian is $|\tilde{\mathbf{J}}| = (hv_\phi/2 - 2\gamma)(hv_\theta/2 - 2\zeta)l'(0)^4 + l'(0)^2$ and the trace is $\text{Tr}(\tilde{\mathbf{J}}) = (hv_\phi/2 - 2\gamma)l'(0)^2 + (hv_\theta/2 - 2\zeta)l'(0)^2$. A necessary and sufficient condition for asymptotic stability is $|\tilde{\mathbf{J}}| > 0$ and $\text{Tr}(\tilde{\mathbf{J}}) < 0$ (since this guarantees that the eigenvalues of the modified Jacobian have negative real part). Therefore, if $\gamma > hv_\phi/4$ and $\zeta > hv_\theta/4$, the system is asymptotically stable. We note however that in practice, when using discrete updates, the exact threshold for stability will have a $\mathcal{O}(h^3)$ correction, arising from the $\mathcal{O}(h^3)$ error in our BEA. Also, we see that when $\gamma = hv_\phi/4$ and $\zeta = hv_\theta/4$, the contribution of the cross-terms is cancelled out (up to an $\mathcal{O}(h^3)$ correction). In Figure 4.5b, we see an example where this explicit regularisation stabilises the DiracGAN, so that it converges toward the equilibrium solution.

B.7 The PF for games: linearisation results around critical points

Assume we are in the general setting described in Section 4.8. To simplify our notation we denote $[\phi, \theta] = \boldsymbol{\psi}$ and the vector $[\nabla_\phi E_\phi, \nabla_\theta E_\theta]$ as $u(\boldsymbol{\psi})$. Thus our discrete simultaneous gradient descent updates become

$$\boldsymbol{\psi}_t = \boldsymbol{\psi}_{t-1} - hu(\boldsymbol{\psi}_{t-1}). \quad (\text{B.225})$$

We also note that from here on we no longer use the structure of the vector field, so the below proof applies to any Euler update.

Assume we are around a critical point $\boldsymbol{\psi}^*$, i.e. $u(\boldsymbol{\psi}^*) = 0$. We can write

$$u(\boldsymbol{\psi}) \approx u(\boldsymbol{\psi}^*) + \frac{du}{d\boldsymbol{\psi}}(\boldsymbol{\psi})(\boldsymbol{\psi} - \boldsymbol{\psi}^*) = \frac{du}{d\boldsymbol{\psi}}(\boldsymbol{\psi}^*)(\boldsymbol{\psi} - \boldsymbol{\psi}^*). \quad (\text{B.226})$$

Replacing this in the discrete updates in Eq (B.225)

$$\boldsymbol{\psi}_t - \boldsymbol{\psi}^* \approx (\mathbf{I} - h \frac{du}{d\boldsymbol{\psi}}(\boldsymbol{\psi}^*))(\boldsymbol{\psi}_{t-1} - \boldsymbol{\psi}^*) \approx (\mathbf{I} - h \frac{du}{d\boldsymbol{\psi}}(\boldsymbol{\psi}^*))^t(\boldsymbol{\psi}_0 - \boldsymbol{\psi}^*). \quad (\text{B.227})$$

If we now take the Jordan normal form of $\frac{du}{d\psi}(\psi^*) = \mathbf{P}^{-1}\mathbf{J}\mathbf{P}$, we can write

$$\psi_t - \psi^* \approx (\mathbf{I} - h \frac{du}{d\psi}(\psi^*))^t (\psi_0 - \psi^*) \quad (\text{B.228})$$

$$= (\mathbf{I} - h\mathbf{P}^{-1}\mathbf{J}\mathbf{P})^t (\psi_0 - \psi^*) \quad (\text{B.229})$$

$$= \mathbf{P}^{-1}(\mathbf{I} - h\mathbf{J})^t \mathbf{P}(\psi_0 - \psi^*). \quad (\text{B.230})$$

Under the above approximation

$$\psi_t = \psi^* + \mathbf{P}^{-1}(\mathbf{I} - h\mathbf{J})^t \mathbf{P}(\psi_0 - \psi^*), \quad (\text{B.231})$$

which we can write in continuous form as (we use that iteration n is at time $t = nh$):

$$\psi(t) = \psi^* + \mathbf{P}^{-1}(\mathbf{I} - h\mathbf{J})^{t/h} \mathbf{P}(\psi_0 - \psi^*) \quad (\text{B.232})$$

$$= \psi^* + \mathbf{P}^{-1}e^{\log(\mathbf{I} - h\mathbf{J})t/h} \mathbf{P}(\psi_0 - \psi^*). \quad (\text{B.233})$$

Taking the derivative w.r.t. t

$$\dot{\psi} = \frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) e^{\log(\mathbf{I} - h\mathbf{J})t/h} \mathbf{P}(\psi_0 - \psi^*) \quad (\text{B.234})$$

$$= \frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) e^{\log(\mathbf{I} - h\mathbf{J})t/h} \mathbf{P} \mathbf{P}^{-1} (e^{\log(\mathbf{I} - h\mathbf{J})t/h})^{-1} \mathbf{P}(\psi - \psi^*) \quad (\text{B.235})$$

$$= \frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) \mathbf{P}(\psi - \psi^*). \quad (\text{B.236})$$

From here, if we assume $u(\psi^*)$ is invertible we can write from the approximation above (Eq. (B.226)) that $\psi - \psi^* = \frac{du}{d\psi}(\psi^*)^{-1} u(\psi) = \mathbf{P}^{-1}\mathbf{J}^{-1}\mathbf{P}u(\psi)$ and replacing this in the above we obtain

$$\dot{\psi} = \frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) \mathbf{P} \mathbf{P}^{-1} \mathbf{J}^{-1} \mathbf{P} u(\psi) \quad (\text{B.237})$$

$$= \frac{1}{h} \mathbf{P}^{-1} \log(\mathbf{I} - h\mathbf{J}) \mathbf{J}^{-1} \mathbf{P} u(\psi), \quad (\text{B.238})$$

which is the generalisation of the PF presented in Eq (4.67).

B.8 GAN Experimental details

Unless otherwise specified, all GAN experiments use the CIFAR-10 dataset and the CIFAR-10 convolutional SN-GAN architectures—see Table 3 in Section B.4 in Miyato et al. [30].

Libraries: We use JAX [325] to implement our models, with Haiku [326] as the neural network library, and Optax [327] for optimisation.

Computer Architectures: All models are trained on NVIDIA V100 GPUs. Each model is trained on 4 devices.

B.8.1 Implementing explicit regularisation

The loss functions we are interested are of the form

$$L = \mathbb{E}_{p(\mathbf{x})} f_{\boldsymbol{\theta}}(\mathbf{x}) \quad (\text{B.239})$$

We then have

$$\nabla_{\boldsymbol{\theta}_i} L = \mathbb{E}_{p(\mathbf{x})} \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\mathbf{x}) \quad (\text{B.240})$$

$$\|\nabla_{\boldsymbol{\theta}} L\|^2 = \sum_{i=1}^{i=|\boldsymbol{\theta}|} (\nabla_{\boldsymbol{\theta}_i} L)^2 = \sum_{i=1}^{i=|\boldsymbol{\theta}|} (\mathbb{E}_{p(\mathbf{x})} \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\mathbf{x}))^2 \quad (\text{B.241})$$

to obtain an unbiased estimate of the above using samples, we have that:

$$\|\nabla_{\boldsymbol{\theta}} L\|^2 = \sum_{i=1}^{i=|\boldsymbol{\theta}|} (\mathbb{E}_{p(\mathbf{x})} \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\mathbf{x}))^2 \quad (\text{B.242})$$

$$= \sum_{i=1}^{i=|\boldsymbol{\theta}|} (\mathbb{E}_{p(\mathbf{x})} \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\mathbf{x})) (\mathbb{E}_{p(\mathbf{x})} \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\mathbf{x})) \quad (\text{B.243})$$

$$\approx \sum_{i=1}^{i=|\boldsymbol{\theta}|} \left(\frac{1}{N} \sum_{k=1}^N \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\widehat{\mathbf{x}}_{1,k}) \right) \left(\frac{1}{N} \sum_{j=1}^N \nabla_{\boldsymbol{\theta}_i} f_{\boldsymbol{\theta}}(\widehat{\mathbf{x}}_{2,j}) \right) \quad (\text{B.244})$$

so we have to use two sets of samples $\mathbf{x}_{1,k} \sim p(\mathbf{x})$ and $\mathbf{x}_{2,j} \sim p(\mathbf{x})$ from the true distribution (by splitting the batch into two or using a separate batch) to obtain the correct norm. To compute an estimator for $\nabla_{\phi} \|\nabla_{\boldsymbol{\theta}} L\|^2$, we can compute the gradient of the above unbiased estimator of $\|\nabla_{\boldsymbol{\theta}} L\|^2$. However, to avoid computing

gradients for two sets of samples, we derive another unbiased gradient estimator, which we use in all our experiments:

$$\frac{2}{N} \sum_{i=1}^{i=|\theta|} \sum_{k=1}^N \nabla_{\phi} \nabla_{\theta_i} f_{\theta}(\widehat{\mathbf{x}}_{1,j}) \nabla_{\theta_i} f_{\theta}(\widehat{\mathbf{x}}_{2,j}) \quad (\text{B.245})$$

B.9 Additional experimental results

B.9.1 Additional results using zero-sum GANs

We show additional results showcasing the effect of DD on zero-sum games in Figure B.1. We see that not only do simultaneous updates perform worse than alternating updates when using the best hyperparameters, but that simultaneous updating is much more sensitive to hyperparameter choice. We also see that multiple updates can improve the stability of a GAN trained using zero-sum losses, but this strongly depends on the choice of learning rate.

Least squares GANs. In order to assess the robustness of our results independent of the GAN loss used, we perform additional experimental results using GANs trained with a least square loss (LS-GAN [38]). We show results in Figure B.2, where we see that for the least square loss too, the learning rate ratios for which the generator drift does not maximise the discriminator norm (learning rate ratios above or equal to 0.5) perform best and exhibit less variance.

B.9.2 GANs using the non-saturating loss

Next, we explore how the strength of the DD depends on the game dynamics. We do this by comparing the *relative effect* that numerical integration schemes have across different games. To this end, we consider the non-saturating loss introduced in the original GAN paper ($-\log D(G(\mathbf{z}; \theta); \phi)$). This loss has been extensively used since it helps to avoid problematic gradients early in training. When using this loss, we see that there is little difference between simultaneous and alternating updates (Figure B.3), unlike the saturating loss case. These results demonstrate that since DD depends on the underlying dynamics, it is difficult to make general game-independent predictions about which numerical integrator will perform best.

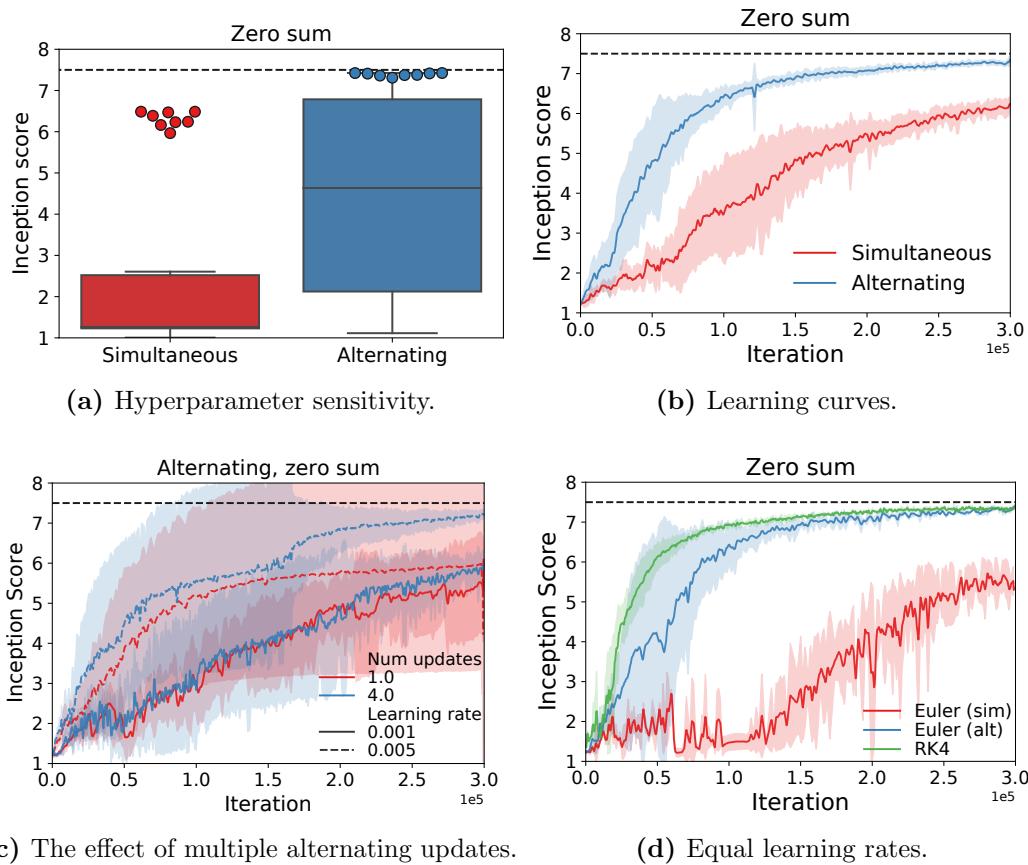


Figure B.1: The effect of discretisation drift on zero-sum games; using the saturating generator loss.

B.9.3 Explicit regularisation in zero-sum games trained using simultaneous gradient descent

We now show additional experimental results and visualisations obtained using explicit regularisation obtained using the original zero sum GAN objective, as presented in Goodfellow et al. [32].

We show the improvement that can be obtained compared to gradient descent with simultaneous updates by cancelling the interaction terms in Figure B.4. We additionally show results obtained from strengthening the self terms in Figure B.6.

In Figure B.5, we show that by cancelling interaction terms, SGD becomes a competitive optimisation algorithm when training the original GAN. We also note that in the case of Adam, while convergence is substantially faster than with SGD, we notice a degrade in performance later in training. This is something that has

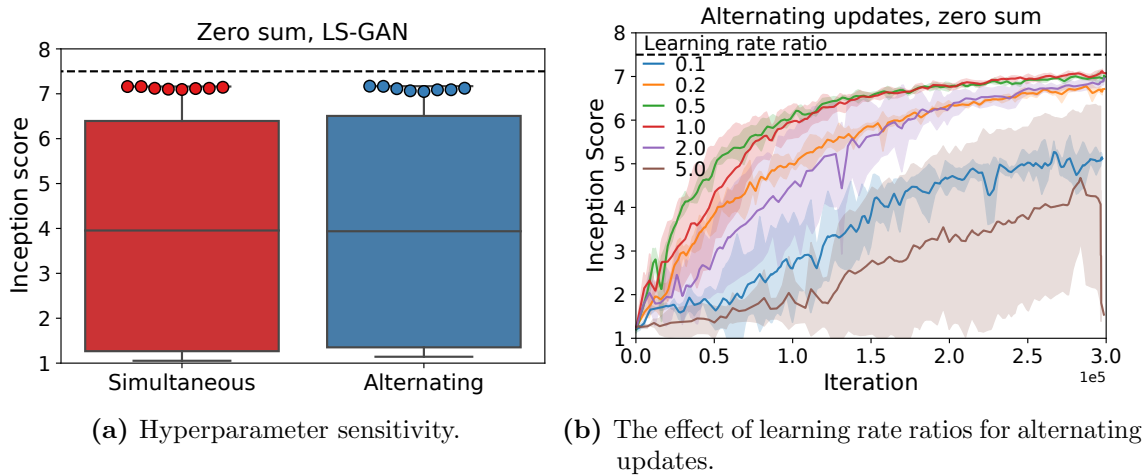


Figure B.2: The effect of discretisation drift on zero-sum games; least square losses.

been observed in other works as well (e.g. [67]).

More percentiles. Throughout the main thesis, we displayed the best 10% performing models for each optimisation algorithm used. We now expand that to show performance results across the best 20% and 30% of models in Figure B.7. We observe a consistent increase in performance obtained by cancelling the interaction terms.

Batch size comparison. We now show that the results which show the efficacy of cancelling the interaction terms are resilient to changes in batch size in Figure B.8.

Comparison with Symplectic Gradient Adjustment. We show results comparing with Symplectic Gradient Adjustment (SGA) [61] in Figure B.9 (best performing models) and Figure B.10 (quantiles showing performance across all hyperparameters and seeds). We observe that despite having the same functional form as SGA, cancelling the interaction terms of DD performs better; this is due to the choice of regularisation coefficients, which in the case of cancelling the interaction terms is provided by Corollary 4.5.1.

Comparison with Consensus Optimisation. We show results comparing with Consensus Optimisation (CO) [15] in Figure B.11 (best performing models) and Figure B.12 (quantiles showing performance across all hyperparameters and seeds). We observe that cancelling the interaction terms performs best, and that additionally

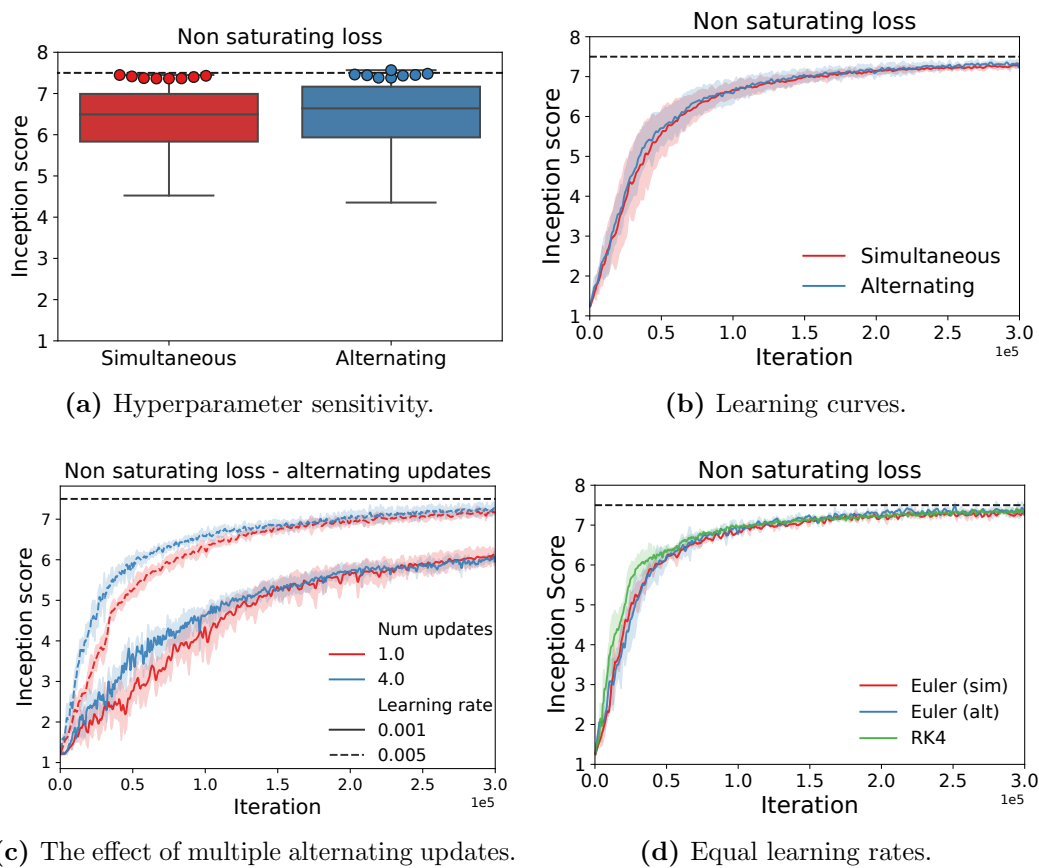


Figure B.3: The effect of discretisation drift depends on the game: with the non saturating loss, the relative performance of different numerical estimators is different compared to the saturating loss, and the effect of explicit regularisation is also vastly different.

strengthening the self terms does not provide a performance benefit.

Variance across seeds. We have mentioned in the main text the challenge with variance across seeds observed when training GANs with SGD, especially in the case of simultaneous updates in zero sum games. We first notice that performance of simultaneous updates depends strongly on learning rates, with most models not learning. We also notice variance across seeds, both in vanilla SGD and when using explicit regularisation to cancel interaction terms. In order to investigate this effect, we ran a sweep of 50 seeds for the best learning rates we obtain when cancelling interaction terms in simultaneous updates, namely a discriminator learning rate of 0.01 and a generator learning rate of 0.005, we obtain Inception Score results with mean 5.61, but a very large standard deviation of 2.34. Indeed, as shown

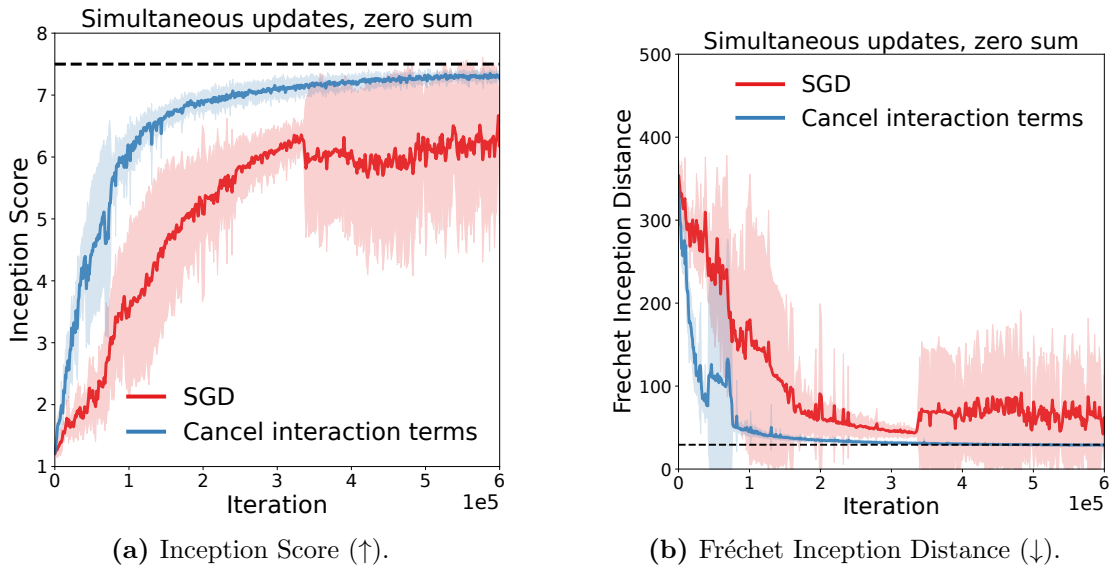


Figure B.4: Using explicit regularisation to cancel the effect of the interaction components of drift leads to a substantial improvement compared to SGD without explicit regularisation.

in Figure B.13, more than 50% of the seeds converge to an IS greater than 7. To investigate the reason for the variability, we repeat the same experiment, but clip the gradient value for each parameter to be in $[-0.1, 0.1]$, and show results in Figure B.14. We notice that another 10% of jobs converge to an IS score greater than 7, and 10% drop in the number of jobs that do not manage to learn. This makes us postulate that the reason for the variability is due to large gradients, perhaps early in training. We contrast this variability across seeds with the consistent performance we obtain by looking at the best performing *models* across learning rates, where as we have shown in the main thesis and throughout the Appendix, we obtain consistent performance which consistently leads to a substantial improvement compared to SGD without explicit regularisation, and obtains performance comparable with Adam.

B.9.4 Explicit regularisation in zero-sum games trained using alternating gradient descent

We perform the same experiments as done for simultaneous updates also with alternating updates. To do so, we cancel the effect of DD using the same explicit regularisation functional form, but updating the coefficients to be those of alternating

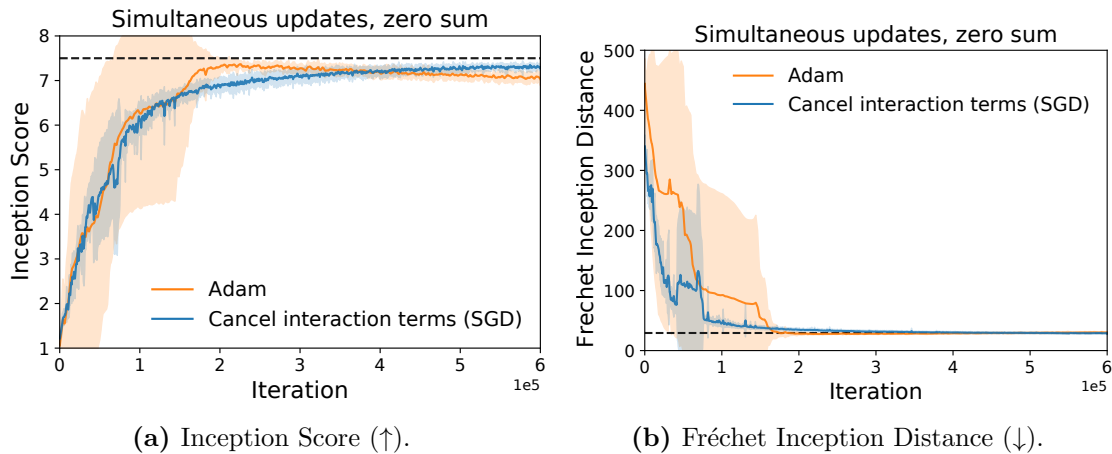


Figure B.5: Using explicit regularisation to cancel the effect of the interaction components of drift allows us to obtain the same peak performance as Adam using SGD without momentum.

updates. We show results in Figure B.15, where we see very little difference in the results compared to vanilla SGD, perhaps apart from less instability early on in training. We postulate that this could be due the effect of DD in alternating updates can be beneficial, especially for learning rate ratios for which the second player also minimises the gradient norm of the first player. We additionally show results obtained from strengthening the self terms in Figure B.16.

More percentiles. Throughout the main thesis, we displayed the best 10% performing models for each optimisation algorithm used. We now expand that to show performance results across the 20% and 30% jobs in Figure B.17. We observe a consistent increase in performance obtained by cancelling the interaction terms.

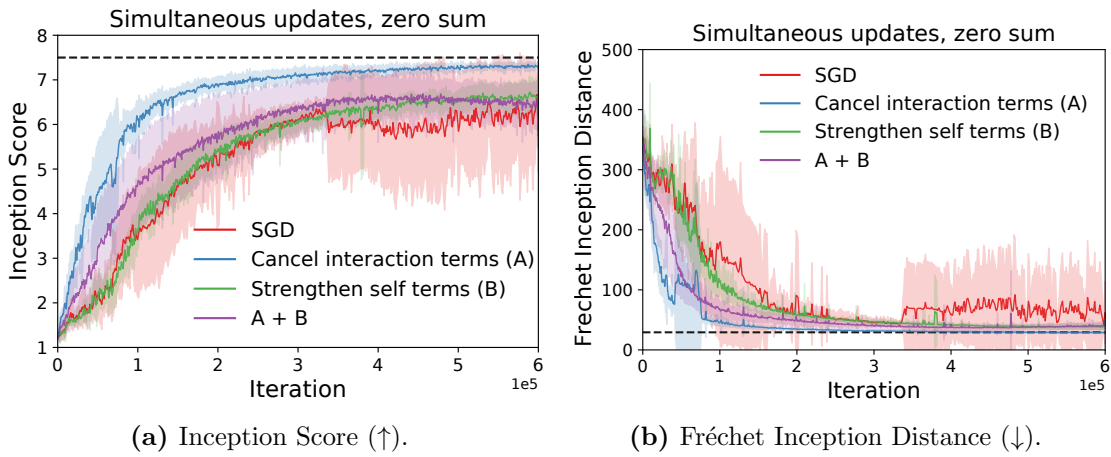


Figure B.6: Using explicit regularisation to cancel the effect of the drift leads to a substantial improvement compared to SGD without explicit regularisation. Strengthening the self terms—the terms which minimise the player’s own norm—does not lead to a substantial improvement; this is somewhat expected since while the modified flows give us the exact coefficients required to *cancel* the drift, they do not tell us how to strengthen it, and our choice of exact coefficients from the drift might not be optimal.

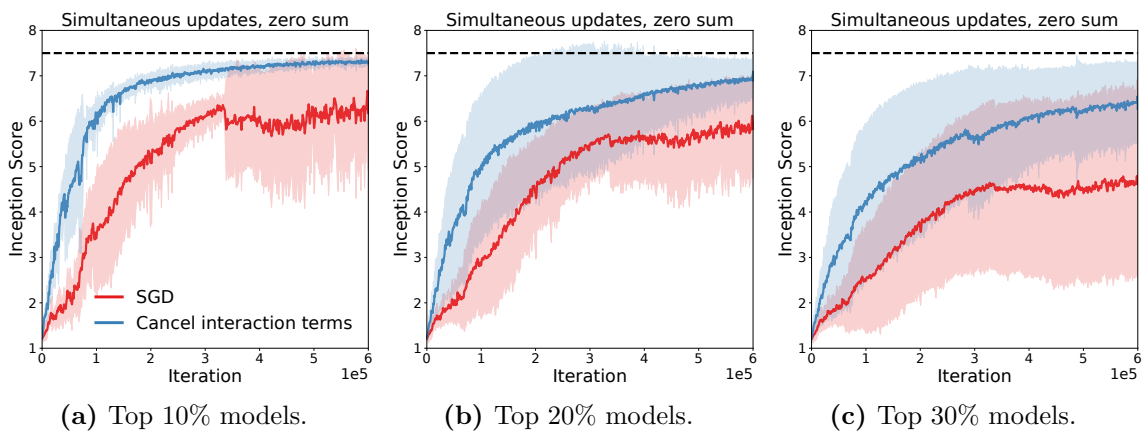


Figure B.7: Performance across the top performing models for vanilla SGD and with cancelling the interaction terms. Across all percentages, cancelling the interaction terms improves performance.

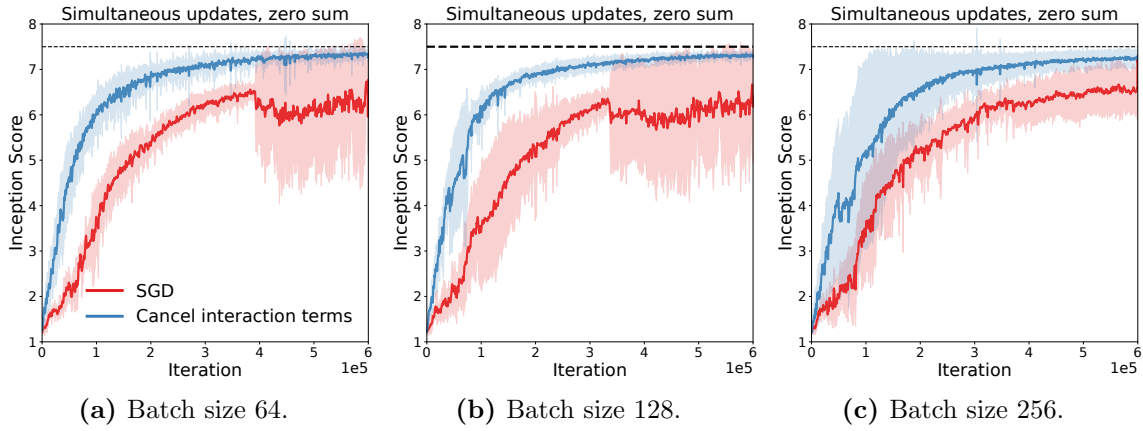


Figure B.8: Performance when changing the batch size. We consistently see that cancelling the interaction terms improves performance.

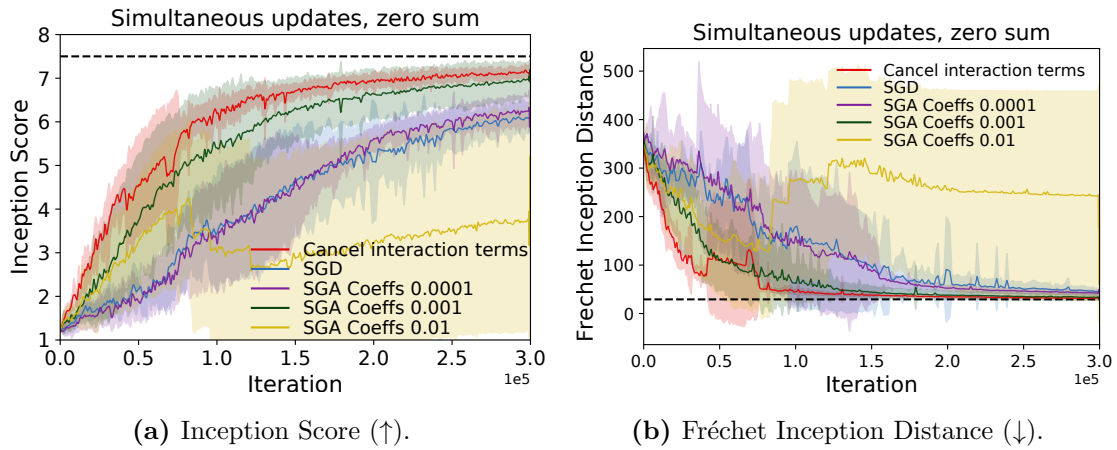


Figure B.9: Comparison with Symplectic Gradient Adjustment (SGA). Cancelling the interaction terms results in similar performance, but with less variance. The performance of SGA heavily depends on the strength of regularisation, adding another hyperparameter to the hyperparameter sweep, while cancelling the interaction terms of the drift requires no other hyperparameters, since the explicit regularisation coefficients strictly depend on learning rates.

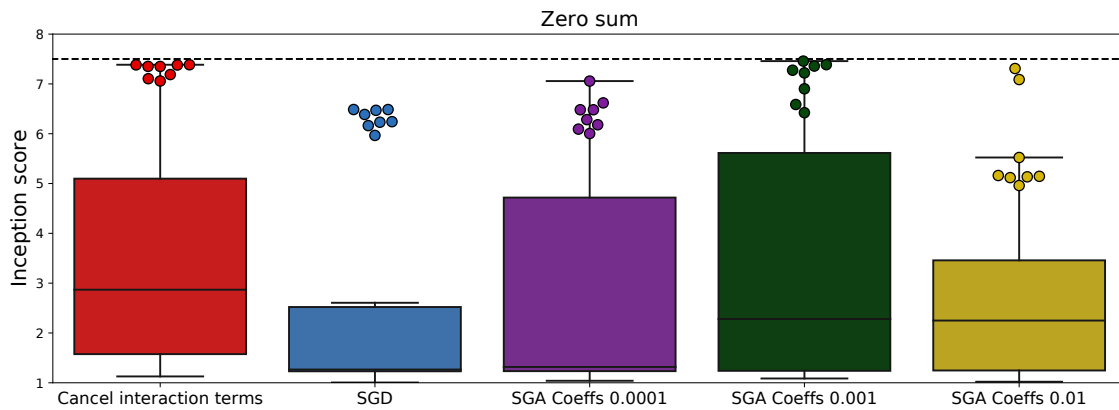
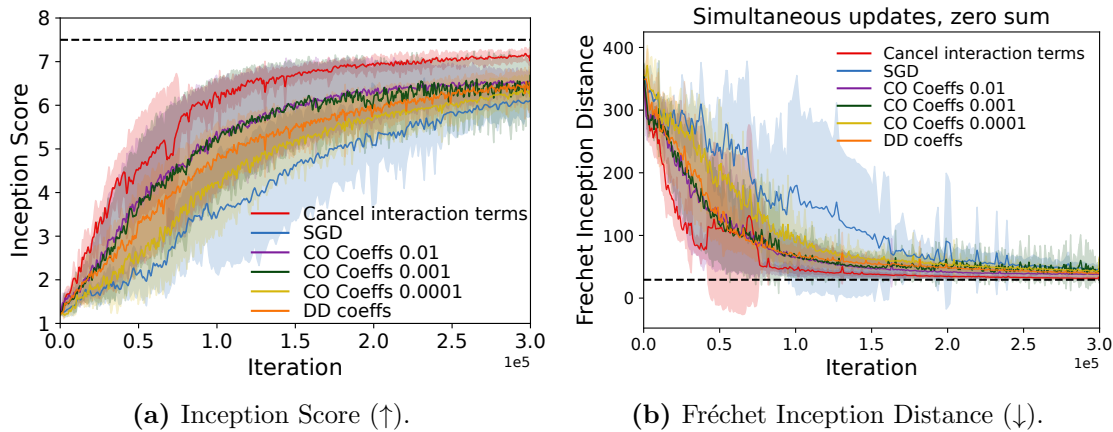


Figure B.10: Comparison with Symplectic Gradient Adjustment (SGA), obtained from *all models in the sweep*. Without requiring an additional sweep over the regularisation coefficient, cancelling the interaction terms results in better performance across the learning rate sweep and less sensitivity to hyperparameters.



(a) Inception Score (\uparrow).

(b) Fréchet Inception Distance (\downarrow).

Figure B.11: Comparison with Consensus Optimisation (CO). Despite not requiring additional hyperparameters compared to the standard SGD learning rate sweep, cancelling the interaction terms of the drift performs better than consensus optimisation. Using consensus optimisation with a fixed coefficient can perform better than using the drift coefficients when we use them to the strengthen the norms – this is somewhat expected since while the modified flows give us the exact coefficients required to *cancel* the drift, they do not tell us how to strengthen it, and our choice of exact coefficients from the drift might not be optimal.

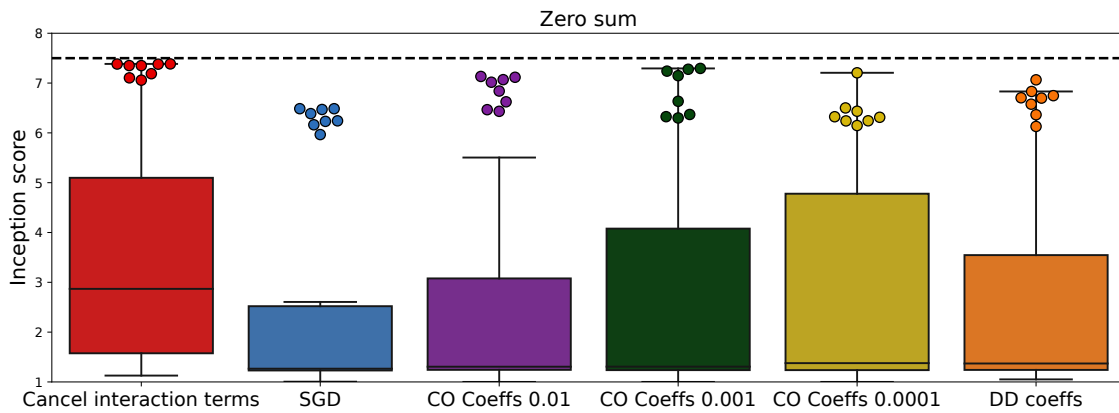
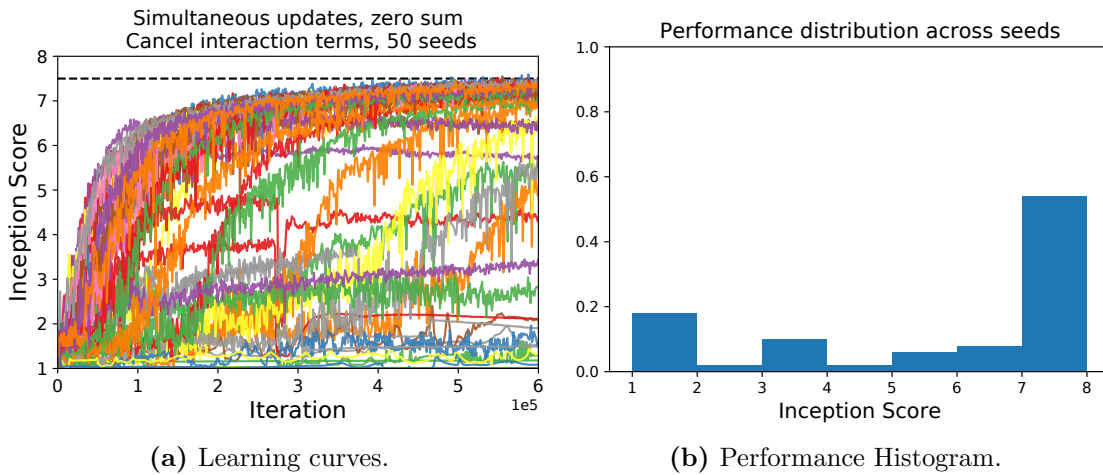


Figure B.12: Comparison with Consensus Optimisation (CO), obtained from *all models in the sweep*. Without requiring an additional sweep over the regularisation coefficient, cancelling the interaction terms results in better performance across the learning rate sweep and less sensitivity to hyperparameters.



(a) Learning curves.

(b) Performance Histogram.

Figure B.13: Variability across seeds for the best performing hyperparameters, when cancelling interaction terms in simultaneous updates for the original GAN, with a zero sum loss.

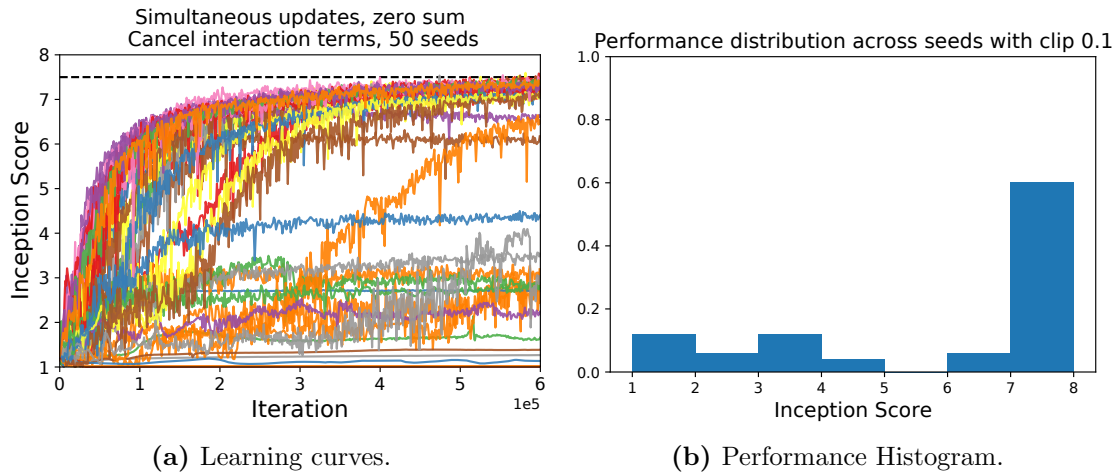


Figure B.14: *With gradient clipping.* Gradient clipping can reduce variability. This suggests that the instabilities observed in gradient descent are caused by large gradient updates.

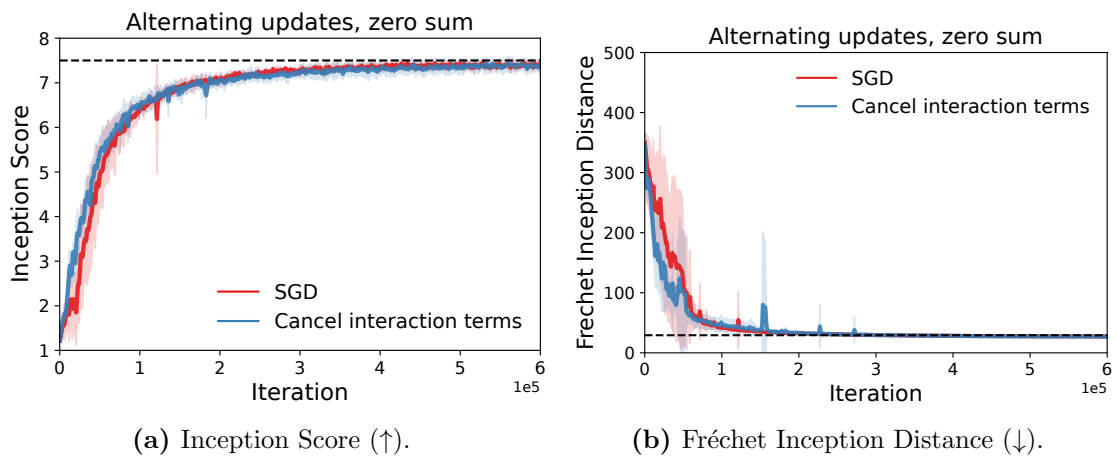


Figure B.15: In alternating updates, using explicit regularisation to cancel the effect of the interaction components of drift does not substantially improve performance compared to SGD, but can reduce variance. This is expected, given that the interaction terms for the second player in the case of alternating updates can have a beneficial regularisation effect.

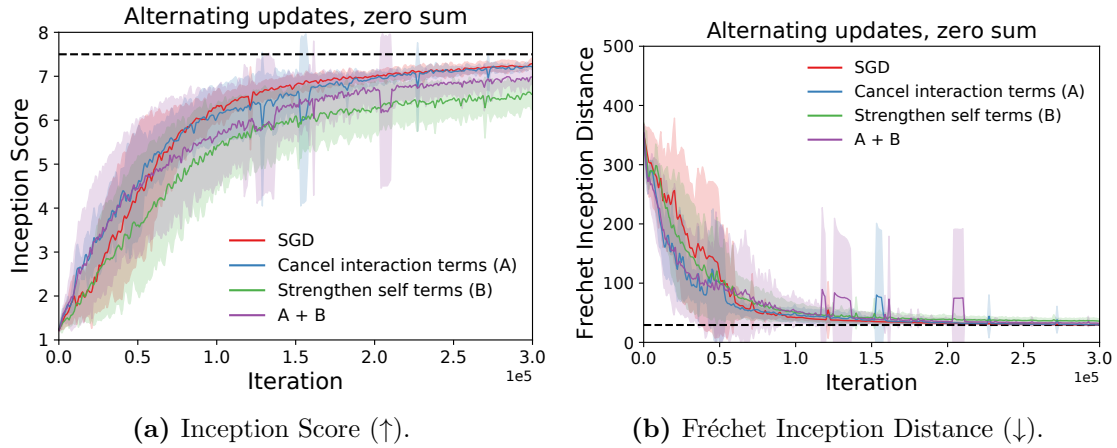


Figure B.16: In alternating updates, using explicit regularisation to cancel the effect of the interaction components of drift does not substantially improve performance compared to SGD, but can reduce variance. This is expected, given that the interaction terms for the second player in the case of alternating updates can have a beneficial regularisation effect. Strengthening the self terms—the terms which minimise the player’s own norm—does not lead to a substantial improvement; this is somewhat expected since while the modified flows give us the exact coefficients required to *cancel* the drift, they do not tell us how to strengthen it, and our choice of exact coefficients from the drift might not be optimal.

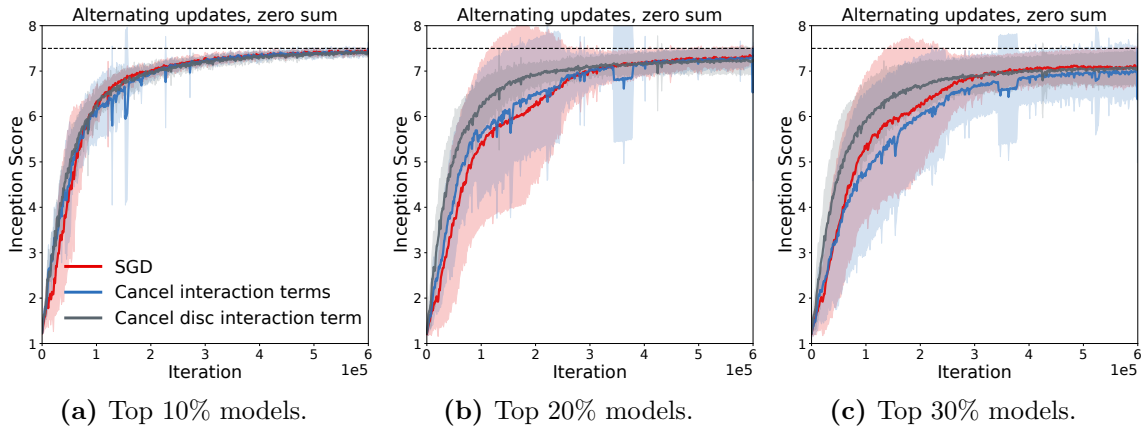


Figure B.17: Performance across the top performing models for vanilla SGD and with cancelling the interaction terms and only cancelling the discriminator interaction terms. We notice that cancelling only the discriminator interaction terms can result in higher performance across more models, and that the interaction term of the generator can play a positive role, likely due to the smaller strength of the generator interaction term compared to simultaneous updates.

B.10 Individual figure reproduction details

- Figure 4.4 performs MNIST classification results with alternating updates. We use an MLP with layers of size $[100, 100, 100, 10]$ and a learning rate of 0.08. The batch size used is 50 and models are trained for 1000 iterations. Error bars are obtained from 5 different seeds.
- Figure 4.5 uses a learning rate of 0.01 for both the discriminator and the generator. The continuous-time simulation is obtained using Runge–Kutta4.
- Figure 4.6 is obtained using sweep over $\{0.01, 0.005, 0.001, 0.0005\}$ for the discriminator and for the generator learning rates. We restrict the ratio between the two learning rates to be in the interval $[0.1, 10]$ to ensure the validity of our approximations. The architecture is Spectral Normalised GAN. Batch size of 128. When comparing with Runge–Kutta 4, we no longer perform a cross product over discriminator and generator learning rates, but instead restrict ourselves to equal learning rates.
- Figure 4.7 controls for the number of experiments which have the same learning rate ratio. To do so, we obtain 5 learning rates uniformly sampled from the interval $[0.001, 0.01]$ which we use for the discriminator, we fix the learning rate ratios to be in $\{0.1, 0.2, 0.5, 1., 2., 5.\}$ and we obtain the generator learning rate from the discriminator learning rate and the learning rate ratio. Batch size of 128.
- Figure 4.8, B.4, B.5, B.6, B.7, B.13, and B.14 use the unbiased estimator described in Section B.8.1 for explicit regularisation. SGD results are obtained from a sweep of $\{0.01, 0.005, 0.001, 0.0005\}$ for the discriminator and for the generator learning rates. We restrict the ratio between the two learning rates to be in the interval $[0.1, 10]$ to ensure the validity of our approximations. For Adam, we use the learning rate sweep $\{10^{-4}, 2 \times 10^{-4}, 3 \times 10^{-4}, 4 \times 10^{-4}\}$ for the discriminator and the same for the generator, with $\beta_1 = 0.5$ and $\beta_2 = 0.99$. All models use simultaneous updates and batch size of 128.

- Figure 4.9, B.11, B.12, B.9, B.10 use the unbiased estimator described in Section B.8.1 for explicit regularisation. Results are obtained from a sweep of $\{0.01, 0.005, 0.001, 0.0005\}$ for the discriminator and for the generator learning rates. We restrict the ratio between the two learning rates to be in the interval $[0.1, 10]$ to ensure the validity of our approximations. The SGA and CO explicit regularisation coefficients are taken from a sweep over $\{0.01, 0.001, 0.0001\}$. All models use simultaneous updates and batch size of 128.
- Figures 4.10, B.15 B.16, B.17 use the same experimental setup as Figure 4.8, with the exception of the update type: here we use alternating updates.
- Figure 4.11 uses the same experimental setting as Figure 4.6, but the non-saturating GAN loss is used.
- Figure B.1 uses the same experimental setting as Figure 4.6.
- Figure B.2 uses a least square loss. Figure B.2a uses the same experimental setting as Figure 4.6, while Figure B.2b uses the same experimental setting as Figure 4.7.
- Figure B.3 uses the same experimental setting as Figure B.1, but uses the non-saturating generator loss.
- Figure B.7 shows results from the same experiment Figure 4.8.
- Figure B.8 uses the same experimental setting as Figure 4.8, but instead of using a fixed batch size of 128, performs a sweep over batch sizes.

Appendix C

Finding new implicit regularisers by revisiting backward error analysis

We now provide the proofs required to write modified losses for stochastic gradient descent in supervised learning discussed in Chapter 5. We denote by $\boldsymbol{\theta}$ the parameters of the model and by f the update function. Since we are now interested in the stochastic setting, we write f as an argument both of data and parameters: $f(\boldsymbol{\theta}; \mathbf{x})$ denotes the application of f at input \mathbf{x} using parameters $\boldsymbol{\theta}$. We often consider $f = -\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{x})$, where E is the loss function. We will denote as

$$E(\boldsymbol{\theta}; \mathbf{X}^t) = \frac{1}{B} \sum_{i=1}^B E(\boldsymbol{\theta}; \mathbf{x}_i^t) \quad (\text{C.1})$$

and

$$f(\boldsymbol{\theta}; \mathbf{X}^t) = -\frac{1}{B} \sum_{i=1}^B \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{x}_i^t) \quad (\text{C.2})$$

$$f(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) = \frac{1}{n} \sum_{i=0}^{n-1} f(\boldsymbol{\theta}; \{\mathbf{X}^{t+i}\}) \quad (\text{C.3})$$

as averages over batches for convenience and clarity.

Our goal is to find $\hat{\boldsymbol{\theta}}$ such that the distance between n steps of stochastic gradient descent (SGD) with learning rate h and $\hat{\boldsymbol{\theta}}$ is of order $\mathcal{O}(h^3)$. We take the same

approach as in the other BEA proofs:

- We expand the n discrete SGD updates up to $\mathcal{O}(h^3)$.
- We expand change of the modified continuous updates of the form $\dot{\boldsymbol{\theta}} = f + hf_1$ in time nh , where f is the gradient function obtained using the concatenation of all batches used in the first step.
- We match the terms between the discrete and continuous updates of $\mathcal{O}(h^2)$ to find f_1 .

C.1 Two consecutive steps of SGD

Step 1: Expand the discrete updates.

From the definition of SGD:

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.4})$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + hf(\boldsymbol{\theta}_t; \mathbf{X}^{t+1}) \quad (\text{C.5})$$

We expand the gradient descent steps and obtain:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + hf(\boldsymbol{\theta}_t; \mathbf{X}^{t+1}) \quad (\text{C.6})$$

$$= \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + hf(\boldsymbol{\theta}_t; \mathbf{X}^{t+1}) \quad (\text{From (C.4), C.7})$$

$$= \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + hf(\boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t); \mathbf{X}^{t+1}) \quad (\text{From (C.5), C.8})$$

$$= \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+1}) \quad (\text{C.9})$$

$$+ h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + \mathcal{O}(h^3) \quad (\text{Taylor expansion, C.10})$$

$$= \boldsymbol{\theta}_{t-1} + 2h \left(\frac{1}{2} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + \frac{1}{2} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+1}) \right) \quad (\text{Grouping } f \text{ terms, C.11})$$

$$+ h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + \mathcal{O}(h^3) \quad (\text{C.12})$$

$$= \boldsymbol{\theta}_{t-1} + 2h \underbrace{f(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}_{\text{update with a batch obtained by concatenating the two batches}} \quad (\text{Eq (C.3), C.13})$$

$$+ h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + \mathcal{O}(h^3) \quad (\text{C.14})$$

Step 2: Taylor expansion of modified flow

We now expand what happens in a time of $2h$ in continuous-time, by using the form of $\dot{\boldsymbol{\theta}}$ given by BEA

$$\dot{\boldsymbol{\theta}} = f(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + 2hf_1(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \quad (\text{C.15})$$

thus for any τ

$$\boldsymbol{\theta}(\tau + 2h) = \boldsymbol{\theta}(\tau) + 2h\dot{\boldsymbol{\theta}}(\tau) + 4h^2\ddot{\boldsymbol{\theta}}(\tau) + \mathcal{O}(h^3) \quad (\text{C.16})$$

$$= \boldsymbol{\theta} + 2hf(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \quad (\text{C.17})$$

$$+ 4h^2 \left[f_1 + \frac{1}{2} \frac{df(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \right] + \mathcal{O}(h^3). \quad (\text{C.18})$$

Step 3: Matching terms of the second order

From the above expansion of the discrete updates (Step 1) and of the continuous updates (Step 2) we can find the value of f_1 at the last initial parameters, $\boldsymbol{\theta}_{t-1}$:

$$4h^2 \left[f_1 + \frac{1}{2} \frac{df(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \right] = h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.19})$$

This leads to:

$$f_1(\boldsymbol{\theta}_{t-1}) = \underbrace{-\frac{1}{2} \frac{df(\cdot; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}_{\text{the IGR full-batch drift term}} \quad (\text{From flow: (C.18), C.20})$$

$$+ \frac{1}{4} \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{From SGD: (C.14), C.21})$$

From here we choose a function f_1 which satisfies the above constraint

$$f_1(\boldsymbol{\theta}) = \underbrace{-\frac{1}{2} \frac{df(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}_{\text{the IGR full-batch drift term}} \quad (\text{C.22})$$

$$+ \frac{1}{4} \frac{df(\boldsymbol{\theta}; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t). \quad (\text{C.23})$$

We have now found the modified flow $\dot{\boldsymbol{\theta}}$, which by construction follows two Euler updates with an error of $\mathcal{O}(h^3)$. Note the use of the initial parameters (highlighted in red) in the flow's vector field; this choice is required to ensure we can write a modified loss function in the single objective optimisation setting by using $f = -\nabla_{\boldsymbol{\theta}}E$:

$$f_1(\boldsymbol{\theta}) = -\frac{1}{2} \frac{d\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})}{d\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \quad (\text{C.24})$$

$$+ \frac{1}{4} \frac{d\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.25})$$

$$= -\nabla_{\boldsymbol{\theta}} \left(\frac{1}{4} \|\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})\|^2 - \frac{1}{4} \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \mathbf{X}^{t+1})^T \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \right),$$

(using (B.96), C.26)

where $E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) = \frac{1}{2} (E(\boldsymbol{\theta}; \mathbf{X}^t) + E(\boldsymbol{\theta}; \mathbf{X}^{t+1}))$. We can now write

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} \left(E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \quad (\text{C.27}) \right.$$

$$\left. + 2h \left(\frac{1}{4} \|\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})\|^2 - \frac{1}{4} \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \mathbf{X}^{t+1})^T \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \right) \right),$$

(C.28)

and since the vector field is a negative gradient, this leads to the modified loss

$$\tilde{E} = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \quad (\text{C.29})$$

$$+ 2h \left(\frac{1}{4} \|\nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\})\|^2 - \frac{1}{4} \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}; \mathbf{X}^{t+1})^T \nabla_{\boldsymbol{\theta}}E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \right). \quad (\text{C.30})$$

C.2 Multiple steps of SGD

We will now derive a similar result, for n SGD steps.

Step 1: Expand discrete updates

$$\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} + hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.31})$$

...

$$\boldsymbol{\theta}_{t+n-1} = \boldsymbol{\theta}_{t+n-2} + hf(\boldsymbol{\theta}_{t+n-2}; \mathbf{X}^{t+n-1}) \quad (\text{C.32})$$

From Eq (C.14), we know that if we expand two steps we obtain

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_{t-1} + 2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + \mathcal{O}(h^3). \quad (\text{C.33})$$

Then, by expanding the third step

$$\boldsymbol{\theta}_{t+2} = \boldsymbol{\theta}_{t+1} + hf(\boldsymbol{\theta}_{t+1}; \mathbf{X}^{t+2}) \quad (\text{C.34})$$

$$= \boldsymbol{\theta}_{t-1} + 2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{Eq (C.33), C.35})$$

$$+ hf(\boldsymbol{\theta}_{t+1}; \mathbf{X}^{t+2}) + \mathcal{O}(h^3) \quad (\text{C.36})$$

$$= \boldsymbol{\theta}_{t-1} + 2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.37})$$

$$+ hf \left(\boldsymbol{\theta}_{t-1} + 2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t); \mathbf{X}^{t+2} \right) \quad (\text{Eq (C.33), C.38})$$

$$+ \mathcal{O}(h^3) \quad (\text{C.39})$$

$$= \boldsymbol{\theta}_{t-1} + 2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.40})$$

$$+ hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+2}) \quad (\text{C.41})$$

$$+ h \frac{df(\cdot; \mathbf{X}^{t+2})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} \left(2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \right) \quad (\text{Taylor expansion, C.42})$$

$$+ \mathcal{O}(h^3) \quad (\text{C.43})$$

$$= \boldsymbol{\theta}_{t-1} + 2hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.44})$$

$$+ hf(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+2}) \quad (\text{C.45})$$

$$+ 2h^2 \frac{df(\cdot; \mathbf{X}^{t+2})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) \quad (\text{Simplifying } h^3 \text{ term, C.46})$$

$$+ \mathcal{O}(h^3) \quad (\text{C.47})$$

$$\boldsymbol{\theta}_{t+2} = \boldsymbol{\theta}_{t-1} + 3hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}, \mathbf{X}^{t+2}\}) \quad (\text{Grouping } f \text{ terms, C.48})$$

$$+ h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.49})$$

$$+ 2h^2 \frac{df(\cdot; \mathbf{X}^{t+2})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}\}) + \mathcal{O}(h^3) \quad (\text{C.50})$$

$$= \boldsymbol{\theta}_{t-1} + 3hf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \mathbf{X}^{t+1}, \mathbf{X}^{t+2}\}) + h^2 \frac{df(\cdot; \mathbf{X}^{t+1})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) \quad (\text{C.51})$$

$$+ h^2 \frac{df(\cdot; \mathbf{X}^{t+2})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^t) + h^2 \frac{df(\cdot; \mathbf{X}^{t+2})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+1}) \quad (\text{Eq (C.3), C.52})$$

$$+ \mathcal{O}(h^3). \quad (\text{C.53})$$

Thus by induction we get

$$\boldsymbol{\theta}_{t+n-1} = \boldsymbol{\theta}_{t-1} + nhf(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.54})$$

$$+ h^2 \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{df(\cdot; \mathbf{X}^{t+\mu})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}) + \mathcal{O}(h^3). \quad (\text{C.55})$$

Step 2: Taylor expansion of modified flow

$$\boldsymbol{\theta}(\tau + nh) = \boldsymbol{\theta} + nhf(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.56})$$

$$+ n^2 h^2 \left[f_1 + \frac{1}{2} \frac{df(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \right] \quad (\text{C.57})$$

$$+ \mathcal{O}(h^3) \quad (\text{C.58})$$

Step 3: Matching terms

From the above expansion of the discrete updates (Step 1) and of the continuous updates (Step 2) we can find the value of f_1 at $\boldsymbol{\theta}_{t-1}$ by matching the terms of order

$\mathcal{O}(h^2)$:

$$\sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{df(\cdot; \mathbf{X}^{t+\mu})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}) \quad (\text{SGD: Eq (C.55), C.59})$$

$$= n^2 \left[f_1 + \frac{1}{2} \frac{df(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \right] \quad (\text{Flow: Eq (C.58), C.60})$$

Leading to

$$f_1(\boldsymbol{\theta}_{t-1}) = \underbrace{-\frac{1}{2} \frac{df(\cdot; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}_{\text{the usual drift term}} \quad (\text{C.61})$$

$$+ \frac{1}{n^2} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{df(\cdot; \mathbf{X}^{t+\mu})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}). \quad (\text{C.62})$$

From here we choose a function f_1 which satisfies the above constraint.

$$f_1(\boldsymbol{\theta}) = \underbrace{-\frac{1}{2} \frac{df(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}_{\text{the usual drift term}} \quad (\text{C.63})$$

$$+ \frac{1}{n^2} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{df(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}) \quad (\text{C.64})$$

We now replace $f = -\nabla_{\boldsymbol{\theta}} E$, using that in this case $\frac{df}{d\boldsymbol{\theta}} = \nabla_{\boldsymbol{\theta}}^2 E$,

$$f_1(\boldsymbol{\theta}_{t-1}) = -\frac{1}{2} \frac{d\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}{d\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.65})$$

$$+ \frac{1}{n^2} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{d\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})}{d\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}) \quad (\text{C.66})$$

$$= -\nabla_{\boldsymbol{\theta}} \left(\frac{1}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \right) \quad (\text{using (B.96), C.67})$$

$$- \frac{1}{n^2} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}). \quad (\text{C.68})$$

This leads to the modified loss

$$\tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.69})$$

$$+ \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.70})$$

$$- \frac{h}{n} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}). \quad (\text{C.71})$$

Through algebraic manipulation we can write the above in the form

$$\tilde{E}(\boldsymbol{\theta}) = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.72})$$

$$+ \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.73})$$

$$- \frac{h}{n} \sum_{\mu=1}^{n-1} \left[\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+\mu})^T \left(\sum_{\tau=0}^{\mu-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}) \right) \right]. \quad (\text{C.74})$$

We can now compare this with the modified loss we obtained by ignoring stochasticity and assuming both updates have been done with a full batch; this entails using the IGR loss:

$$\tilde{E} = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{h}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.75})$$

Thus, when using multiple batches, there is the additional pressure to maximise the dot product of the gradients obtained using the last k batches that came before the current batch, evaluated at the parameters at which we started the n iterations.

C.3 Multiple steps of full-batch gradient descent

To contrast our results with multiple steps of gradient descent *with the same batch* (full-batch case), we show that the IGR flow follows gradient descent with error of order $\mathcal{O}(h^3)$ after n gradient descent steps. The proof steps follow the same approach as above, with two main differences: first, we assume all batches are the same, and second, we no longer fix the starting parameters when choosing f_1 . With the same

steps as before we obtain Eq (C.62):

$$f_1(\boldsymbol{\theta}_{t-1}) = -\frac{1}{2} \frac{df(\cdot; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.76})$$

$$+ \frac{1}{n^2} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{df(\cdot; \mathbf{X}^{t+\mu})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_{t-1}} f(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+\tau}). \quad (\text{C.77})$$

Assuming identical batches and using that f is an empirical average, we obtain $f(\cdot; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) = f(\cdot; \mathbf{X}^t) = f(\cdot; \mathbf{X}^{t+i})$, for any choice of i . We then have

$$f_1(\boldsymbol{\theta}) = -\frac{1}{2} \frac{df(\boldsymbol{\theta}, \mathbf{X}^t)}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{X}^t) + \frac{1}{n^2} \sum_{\tau=0}^{n-1} \sum_{\mu=\tau+1}^{n-1} \frac{df(\boldsymbol{\theta}, \mathbf{X}^t)}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{X}^t) = \quad (\text{C.78})$$

$$= \left(-\frac{1}{2} + \frac{n(n-1)}{2n^2}\right) \frac{df(\boldsymbol{\theta}, \mathbf{X}^t)}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{X}^t) \quad (\text{C.79})$$

$$= -\frac{1}{2n} \frac{df(\boldsymbol{\theta}, \mathbf{X}^t)}{d\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{X}^t) \quad (\text{C.80})$$

Replacing $f = -\nabla_{\boldsymbol{\theta}} E(\cdot; \mathbf{X}^t)$ into the form of the modified flow in Eq C.58, we obtain:

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^t) - \frac{nh}{2n} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}, \mathbf{X}^t) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^t) \quad (\text{C.81})$$

$$= -\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^t) - \frac{h}{2} \nabla_{\boldsymbol{\theta}}^2 E(\boldsymbol{\theta}, \mathbf{X}^t) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^t), \quad (\text{C.82})$$

which recovers the IGR flow. Importantly, the flow in the full-batch case *does not depend on the number of iterations* n .

C.3.1 Expectation over all shufflings

We contrast our results with those of Smith et al. [13], who construct a modified loss over an epoch of stochastic gradient descent in expectation over all possible shufflings of the batches in the epoch. That is, Smith et al. [13] describe expected value of the modified loss $\mathbb{E}_{\sigma} \left[\tilde{E}(\boldsymbol{\theta}; \{\mathbf{X}^{\sigma(t)}, \dots, \mathbf{X}^{\sigma(t+n-1)}\}) \right]$, where n is the number of batches in an epoch, and t is the iteration at which the new epoch start, while σ denotes the set of all possible permutations of batches $\{1, \dots, n\}$. While our results do not require an expectation and account for the exact batches used, to contrast our results with Smith et al. [13], we also take an expectation over all possible batch

shufflings (but not the elements in the batch) in an epoch in our results shown in Eq (5.8), and obtain

$$\mathbb{E}_\sigma [E_{sgd}(\boldsymbol{\theta})] = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.83})$$

$$- \frac{h}{n} \mathbb{E}_\sigma \left[\sum_{k=0}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+k})^T \left(\sum_{i=0}^{k-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+i}) \right) \right] \quad (\text{C.84})$$

$$= E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.85})$$

$$- \frac{h}{n} \frac{1}{2} \mathbb{E}_\sigma \left[\sum_{k=0}^{n-1} \sum_{i=0, i \neq k}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+k})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+i}) \right], \quad (\text{C.86})$$

We used the symmetry of the permutation structure since for each permutation where $\sigma(i) < \sigma(j)$ there is also a permutation where $\sigma(i) > \sigma(j)$ by swapping the values of $\sigma(i)$ and $\sigma(j)$. We expand the last term

$$\mathbb{E}_\sigma [E_{sgd}] = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.87})$$

$$- \frac{h}{n} \frac{1}{2} \mathbb{E}_\sigma \left[\sum_{k=0}^{n-1} \sum_{i=0}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+k})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+i}) \right] \quad (\text{C.88})$$

$$- \frac{h}{n} \frac{1}{2} \mathbb{E}_\sigma \left[\sum_{k=0}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+k})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+k}) \right] \quad (\text{C.89})$$

From here

$$\mathbb{E}_\sigma [E_{sgd}] = E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) + \frac{nh}{4} \|\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})\|^2 \quad (\text{C.90})$$

$$- \frac{h}{2n} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \{\mathbf{X}^t, \dots, \mathbf{X}^{t+n-1}\}) \quad (\text{C.91})$$

$$- \frac{h}{2n} \left[\sum_{k=0}^{n-1} \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}; \mathbf{X}^{t+k})^T \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_{t-1}; \mathbf{X}^{t+k}) \right]. \quad (\text{C.92})$$

We obtain that the pressure to minimise individual batch gradient norms is translated into a pressure to maximise the dot product between the gradients at the end of the epoch with those at the beginning of the epoch, both for each batch and for the entire dataset.

C.4 Two-player games

C.4.1 Effects on the non-saturating GAN

We now expand of implicit regularisation terms found in Section 5.3 on GANs, specifically, the effect of the interaction terms, which takes the form of a dot product. We will denote the first player, the discriminator, as D , parametrised by ϕ , and the generator as G , parametrised by θ . We denote the data distribution as $p^*(\mathbf{x})$ and the latent distribution $p(\mathbf{z})$. Consider the non-saturating GAN loss described by Goodfellow et al. [32]:

$$E_\phi(\phi, \theta) = \mathbb{E}_{p^*(\mathbf{x})} \log D(\mathbf{x}; \phi) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}; \theta); \phi)) \quad (\text{C.93})$$

$$E_\theta(\phi, \theta) = \mathbb{E}_{p(\mathbf{z})} - \log D(G(\mathbf{z}; \theta); \phi) \quad (\text{C.94})$$

Consider the interaction term for the discriminator (see Eq (5.33)), and replace the definitions of the above loss functions:

$$\nabla_\theta E_\phi^T \nabla_\theta E_\theta(\phi_{t-1}, \theta_{t-1}) \quad (\text{C.95})$$

$$= \nabla_\theta \left(\mathbb{E}_{p^*(\mathbf{x})} \log D(\mathbf{x}; \phi) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}; \theta); \phi)) \right)^T \nabla_\theta E_\theta(\phi_{t-1}, \theta_{t-1}) \quad (\text{C.96})$$

$$= \nabla_\theta \left(\mathbb{E}_{p(\mathbf{z})} \log(1 - D(G(\mathbf{z}; \theta); \phi)) \right)^T \nabla_\theta E_\theta(\phi_{t-1}, \theta_{t-1}) \quad (\text{C.97})$$

$$= \left(\mathbb{E}_{p(\mathbf{z})} \nabla_\theta \log(1 - D(G(\mathbf{z}; \theta); \phi)) \right)^T \nabla_\theta E_\theta(\phi_{t-1}, \theta_{t-1}) \quad (\text{C.98})$$

$$= \left(-\mathbb{E}_{p(\mathbf{z})} \frac{1}{1 - D(G(\mathbf{z}; \theta); \phi)} \nabla_\theta D(G(\mathbf{z}; \theta); \phi) \right)^T \nabla_\theta E_\theta(\phi_{t-1}, \theta_{t-1}) \quad (\text{C.99})$$

$$= \left(-\mathbb{E}_{p(\mathbf{z})} \frac{1}{1 - D(G(\mathbf{z}; \theta); \phi)} \nabla_\theta D(G(\mathbf{z}; \theta); \phi) \right)^T \quad (\text{C.100})$$

$$\left(-\mathbb{E}_{p(\mathbf{z})} \frac{1}{D(G(\mathbf{z}; \theta_{t-1}); \phi_{t-1})} \nabla_\theta D(G(\mathbf{z}; \theta_{t-1}); \phi_{t-1}) \right) \quad (\text{C.101})$$

$$= \left(\mathbb{E}_{p(\mathbf{z})} \frac{1}{1 - D(G(\mathbf{z}; \theta); \phi)} \nabla_\theta D(G(\mathbf{z}; \theta); \phi) \right)^T \quad (\text{C.102})$$

$$\left(\mathbb{E}_{p(\mathbf{z})} \frac{1}{D(G(\mathbf{z}; \theta_{t-1}); \phi_{t-1})} \nabla_\theta D(G(\mathbf{z}; \theta_{t-1}); \phi_{t-1}) \right) \quad (\text{C.103})$$

where the expectations can be evaluated at the respective mini-batches used in the

updates for iterations t and $t - 1$, respectively. Consider \mathbf{z}_t^i the latent variable with index i in the batch at time t . Then the above is approximated as

$$\frac{1}{B^2} \sum_{i,j=1}^B c_{i,j}^{non-sat} \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})^T \nabla_{\boldsymbol{\theta}} D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1}), \quad \text{with (C.104)}$$

$$c_{i,j}^{non-sat} = \frac{1}{1 - D(G(\mathbf{z}_t^i; \boldsymbol{\theta}); \boldsymbol{\phi})} \frac{1}{D(G(\mathbf{z}_{t-1}^j; \boldsymbol{\theta}_{t-1}); \boldsymbol{\phi}_{t-1})}. \quad (\text{C.105})$$

Appendix D

The importance of model smoothness in deep learning

D.1 Individual figure reproduction details

- Figure 6.1: the U-shape curve is obtained by fitting polynomials of increased degree on a simple one dimensional regression problem, with the true underlying function $f(x) = \sin(5x) - \sin(10x)$. The double descent curve is obtained by training Resnet-18 models on CIFAR-10.
- Figure 6.2 uses samples from two Beta distributions, one with parameters 3 and 5, and one with parameters 2 and 3. The MLP approximation to the density ratio is obtained with an MLP with units [10, 10, 40, 1]. The Wasserstein optimal critic is approximated using a linear programming approach [328].
- Figures 6.3, 6.4, 6.5 use the an MLP with 4 layers and 100, 100, 100 and 1 output unit, respectively. The shallow MLP has 2 layers of 100 and 1 unit. All methods were trained for 100 iterations on 50 data points. We use the Adam optimiser, with learning rate 0.01, $\beta_1 = 0.9$, $\beta_2 = 0.999$. To compute the local Lipschitz function of the decision surface learned on two moons in Figure 6.5, we split the space into small neighbourhoods (2500 equally sized grids); for each grid, we sample 2500 random pairs of points in the grid and report $\max \|f(\mathbf{x}) - f(\mathbf{y})\| / \|\mathbf{x} - \mathbf{y}\|$.
- Figure 6.6 is obtained using an MLP with 4 layers of 1000, 1000, 1000 and

10 units each and are trained for 500 iterations at batch size 100, reaching an accuracy of 95% on the entire test set. The models are trained with the Adam optimiser, with $\beta_1 = 0.9$ $\beta_2 = 0.999$, and learning rates 0.001 and 0.0005.

- Figure 6.7 For the GAN CIFAR-10 experiments, we use the architectures specified in the Spectral normalization paper [30]. We use the Adam optimiser [47] with default β_2 , β_1 as specified in the respective subfigure and a learning rate sweep $\{10^{-4}, 2 \times 10^{-4}, 3 \times 10^{-4}, 4 \times 10^{-4}\}$ for the discriminator and the same for the generator.

Appendix E

Spectral Normalisation in Reinforcement learning

E.1 Additional experimental results

E.1.1 The weak correlation between smoothness and performance

In Section 7.3.3 we have made the argument that smoothness and performance are not strongly correlated. We visualise the individual results used to obtain the correlation values shown in the main thesis in Figure E.1.

E.1.2 Other regularisation methods

We investigated whether other regularisation methods imposing smoothness constraints can have similar effects on the agent’s performance. To this end we ran experiments with both gradient penalties and Batch Normalisation on several architectures (Table E.6). We consistently observe that these approaches do not recover the performance of Spectral Normalisation.

Batch Normalisation. For each of the architectures we employed Batch Normalisation after the ReLU activation of every convolutional or linear layer except the output.

Gradient Penalty. We did extensive experimentation with gradient penalty regularisation. We tried multiple forms of regularisation, including penalising the norm of the sum of the gradients of all actions $\left\| \sum_{i=1}^{|\mathcal{A}|} \frac{dQ_i}{ds} \right\|_2$, regularising the expected norm

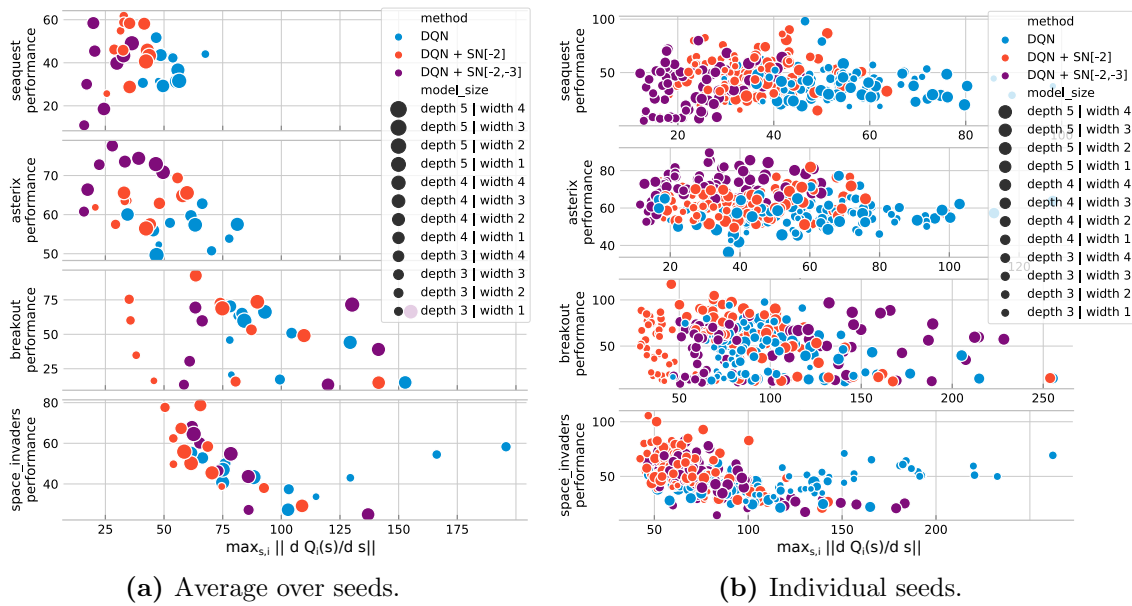


Figure E.1: Applying normalisation does not always result in smoother networks. Even if normalising a subset of the network’s layers makes the network less smooth than the baseline while performance still improves. Results averaged over 10 seeds in (a); individual seeds shown in (b). Results are obtained from a sweep over architectures detailed in Table E.6.

of each Q-value with respect to the state $\sum_{i=1}^{|\mathcal{A}|} \left\| \frac{dQ_i}{ds} \right\|_2$, and also regularising the norm of the gradient of the Q-value associated with the optimal action $\left\| \frac{d \max Q_i}{ds} \right\|_2$. In all cases we swept through a wide range of penalty coefficients ζ . In Figure E.2 we report the results of the best setting we could identify.

Relaxations to 1-Lipschitz normalisation. We now investigate whether relaxing the 1-Lipschitz constraint imposed by Spectral Normalisation to a K -Lipschitz constraint allows us to normalise more layers without a decrease in performance. We have seen in Chapter 6 that this can occur on simple classification tasks; the need for relaxing the Lipschitz constraint has also been observed by Gouk et al. [261]. Figure E.3 shows that $K > 1$ values Spectral Normalisation can be applied to multiple layers, and it further leads to an increase in performance. As expected, we also observe that for rather large K (such as $K = 10$), Spectral Normalisation no longer provides a benefit. The increased computation required when approximating the spectral norm for all the layers and the addition of one hyperparameter per layer determined us to not pursue this setup further.

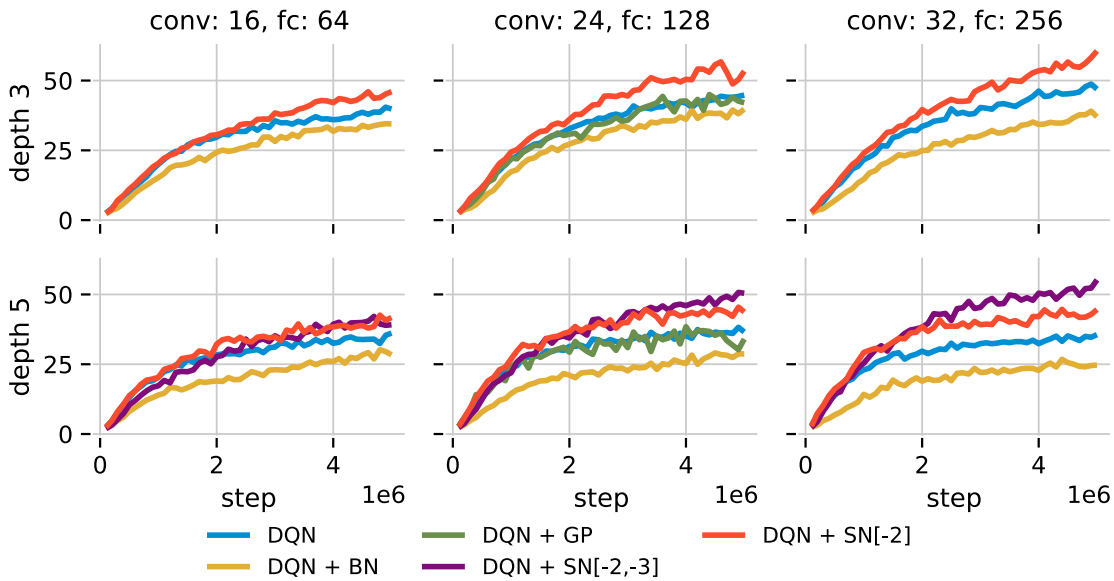


Figure E.2: Regularisation does not recover Spectral Normalisation performance. Performance on MinAtar games of Spectral Normalisation, gradient penalties and Batch Normalisation. Each line is an average over normalized scores of each game. Ten seeds for each configuration.

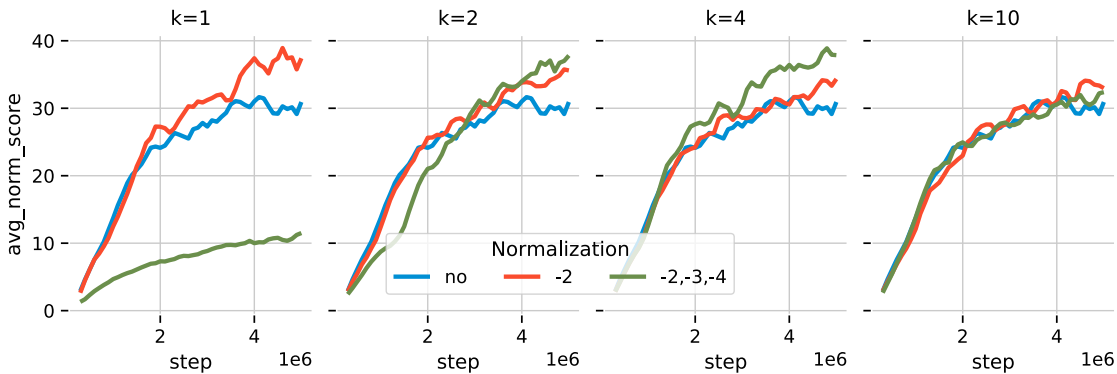


Figure E.3: A sweep over Lipschitz constants K when normalising all but the last layer of the critic. We see that while for $K = 1$, normalising 3 layers drastically reduces performance, this is not the case for $K > 1$. For large K , however, we also do not observe any performance improvement. Each line is an average over normalised scores of 4 games \times 10 seeds.

E.1.3 divOut, divGrad and mulEps

We present results comparing DIVOUT, MULEPS and DIVGRAD across multiple normalised layers and architectures in Figure E.5. DIVOUT has a close behaviour to that of Spectral Normalisation. In contrast, the MULEPS and DIVGRAD optimisers outperform Spectral Normalisation when all hidden layers are normalised.

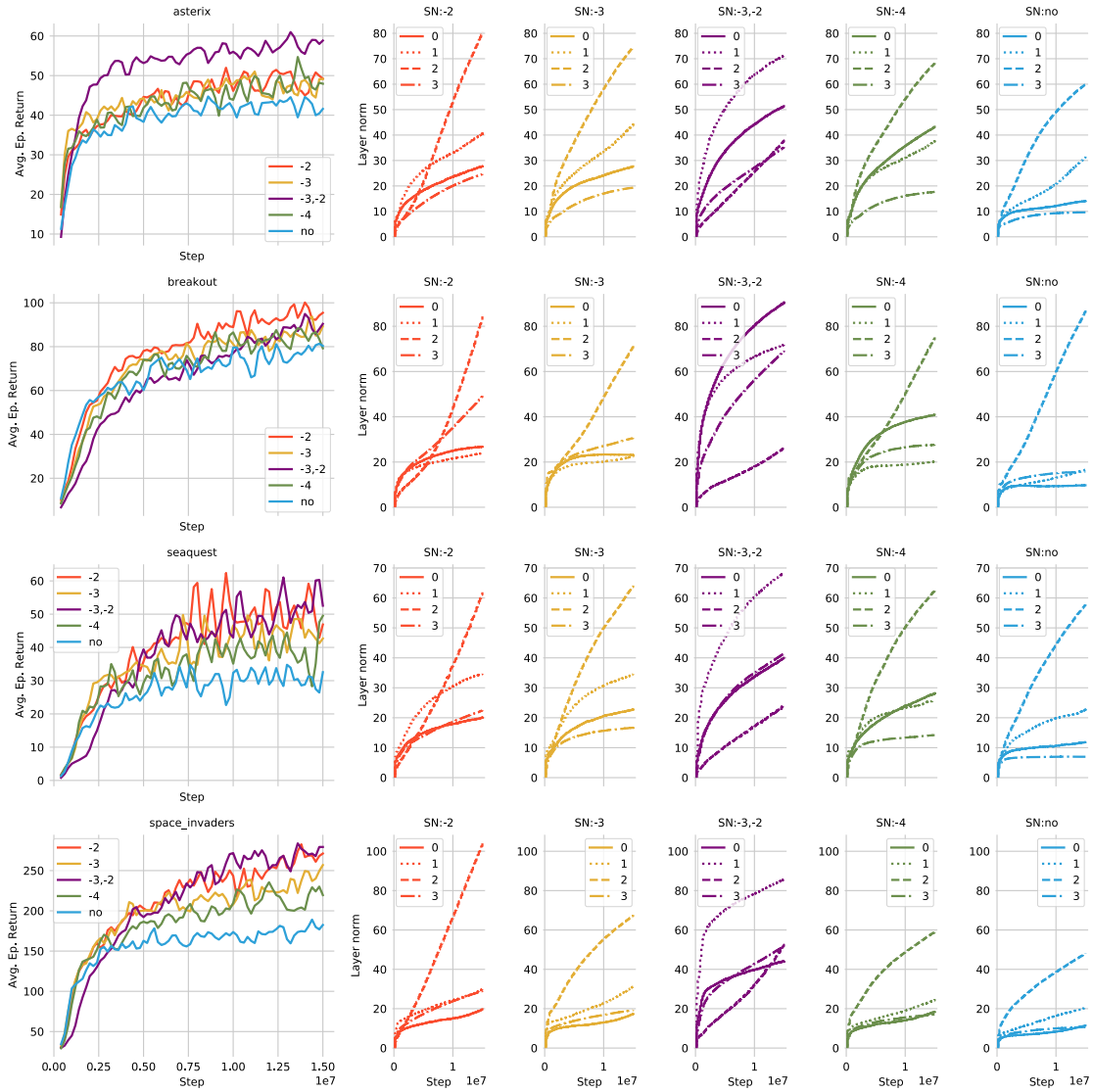


Figure E.4: All spectral radii for the 15M experiment on MinAtar using a 4-layer architecture (conv=24-24,fc=128). Colors code the subsets of layers that are normalised (consistent with the rest of the document), while line styles code the four layers. Note how the penultimate layer has the largest spectral norm across all normalisation variants. 10 seeds.

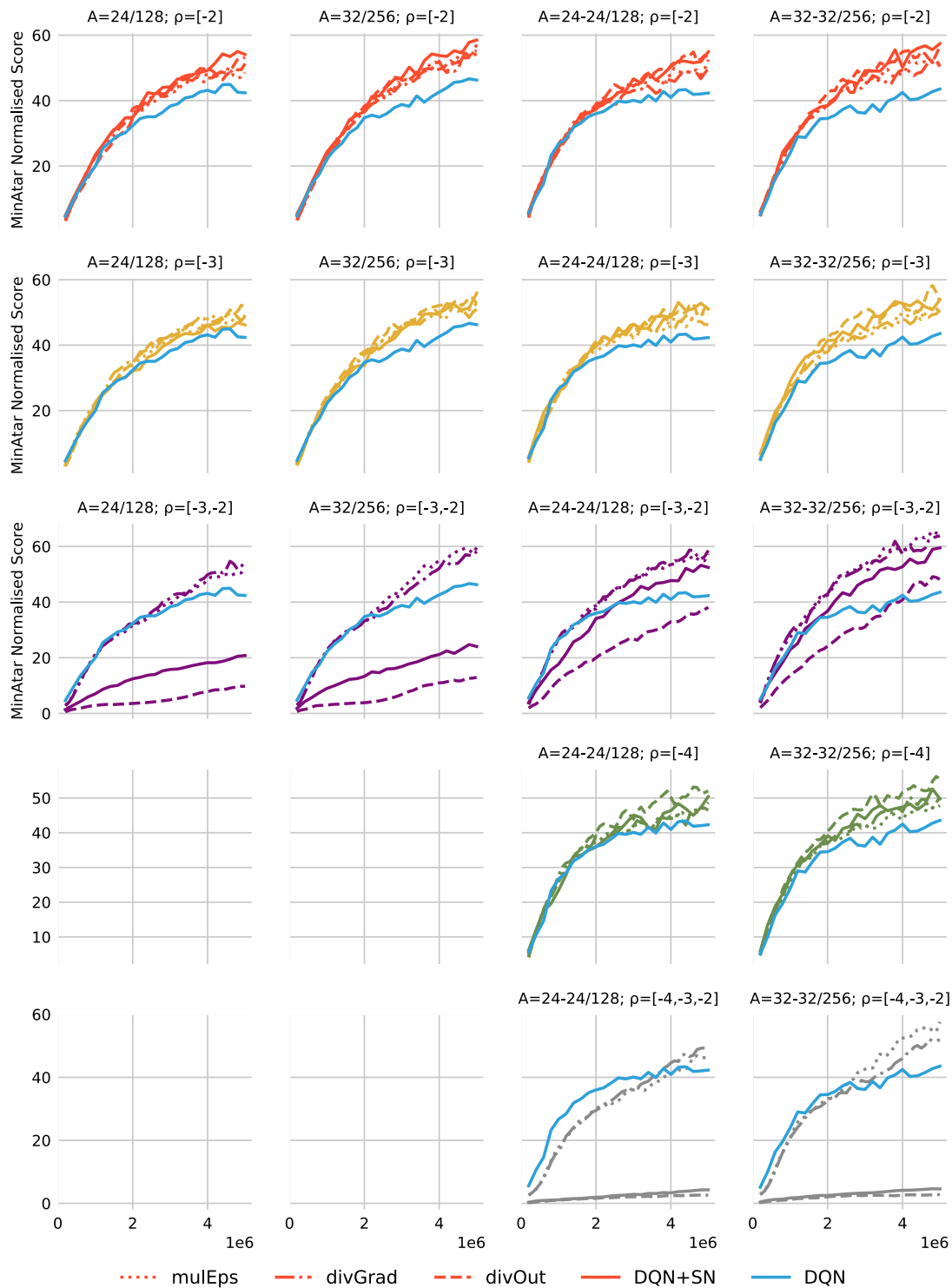


Figure E.5: MinAtar Normalised Scores for the four architectures in Table E.6 and various subsets of layers whose spectral radii are used for Spectral Normalisation or optimisation methods we derived based on Spectral Normalisation. Notice that DIVOUT behaves similarly to Spectral Normalisation (even when they fail to train), while MULEPS and DIVGRAD converge even when all hidden layers are normalised.

E.2 Experimental details

E.2.1 Application of Spectral Normalisation

We use power iteration to approximate the spectral norm of weight matrices, as suggested by Miyato et al. [30]. For convolutional layers we adapt the procedure in Gouk et al. [261] and the two matrix-vector multiplications are replaced by convolutional and transposed convolutional operations.

Backpropagating through the norm. Since the parameters tuned during optimisation are the unnormalised weights \mathbf{W}_i , we investigated whether there are any advantages in backpropagating through the power iteration step. In a set of experiments performed on a subset of games of Atari we noticed no loss in performance when treating the spectral norm as a constant from the backpropagation perspective, and thus for computational efficiency we adopt this approach.

E.2.2 Measuring smoothness: the norm of the Jacobians

Experiments in Section E.1.1 used the maximum norm of the Jacobians w.r.t. the inputs as an indirect metric for network function’s smoothness. This is to avoid using potential loose bounds discussed in Chapter 6, since the exact computation of the Lipschitz constant of a network is NP-hard [299]. For each network we collected thousands of states (*on-policy*), computing the maximum Euclidean norm of the Jacobian w.r.t. the inputs: $\max_{i,\mathbf{x}} \left\| \frac{dQ_i}{d\mathbf{s}} \right\|_2$, where $\|\mathbf{s}\|_2$ denotes the Frobenius norm and Q_i denotes the Q value corresponding to action i .

We note that what we compute is a lower bound of the Lipschitz constant. To see why consider $Q(\cdot; \boldsymbol{\theta}) : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{A}|}$ with Lipschitz constant w.r.t. the Frobenius norm K . We then we have that:

$$\max_i \left\| \frac{dQ_i}{d\mathbf{s}} \right\|_2^2 \leq \left\| \frac{dQ}{d\mathbf{s}} \right\|_2^2 \leq \min(|\mathcal{S}|, |\mathcal{A}|) \left\| \frac{dQ}{d\mathbf{s}} \right\|_{op}^2 \leq \min(|\mathcal{S}|, |\mathcal{A}|) K^2. \quad (\text{E.1})$$

Since we chose the value of \mathbf{s} with maximum norm, this ensures that the bound is as tight as possible with the available data (though this can be made tighter by computing the $\max \left\| \frac{dQ}{d\mathbf{s}} \right\|_2^2$).

E.2.3 Evaluation

The mean and median Human Normalised Score have been critiqued in the past [329, 330] because they can be dominated by some large normalised scores. However we notice that with very few exceptions (3/54 for DQN and 11/54 games for C51) normalising at most one of the network’s layers will not degrade the performance compared to the baseline but instead it will improve upon it, often substantially.

E.2.4 Minatar

Game selection. MinAtar [259] benchmark is a collection of five games that reproduce the dynamics of ALE counterparts, albeit in a smaller observational space. Out of *Asterix*, *Breakout*, *Seaquest*, *Space Invaders* and *Freeway* we excluded the latter from all our experiments since all agents performed similarly on this game.

Network architecture. All experiments on MinAtar are using a convolutional network with L_C convolutional layers with the same number of channels, a hidden linear layer, and the output layer. The number of input channels is game-dependent in MinAtar. All convolutional layers have a kernel size of 3 and a stride of 1. All hidden layers have rectified linear units. Whenever we vary depth we change the number of convolutional layers L_C , keeping the two linear layers. When the width is varied we change the width of both convolutional layers and the penultimate linear layer. All convolutional layers are always identically scaled. We list all the architectures used in various experiments described in this section in Table E.6.

General hyperparameter settings. In all our MinAtar experiments we used the same set of hyperparameters returned by a small grid search around the initial values published by Young and Tian [259]. We list the values we settled on in Table E.1. For the rest of this section we only mention how we deviate from this set of hyperparameters and settings for each of the experiments that follow.

MinAtar Normalised Score. In our work we present many MinAtar experiments as averages over the four games we tested on. Since in MinAtar the range of the expected returns is game dependent we normalise the score. Inspired by the Human Normalised Score in Mnih et al. [107] we take the largest score ever recorded by a baseline agent in our experiments and use it to compute $MNS = 100 \times (\text{score}_{\text{agent}} -$

HYPERPARAMETER	VALUE
discount γ	0.99
update frequency	4
target update frequency	4,000
starting ϵ	1.0
final ϵ	0.01
ϵ steps	250,000
ϵ schedule	linear
warmup steps	5,000
replay size	100,000
history length	1
cost function	MSE LOSS
optimiser	ADAM
learning rate h	0.00025
damping term ϵ	0.0003125
β_1, β_2	(0.9, 0.999)
validation steps	125000
validation ϵ	0.001

Table E.1: MinAtar default hyperparameter settings, used for all MinAtar experiments unless otherwise specified.

$\text{score}_{\text{random}})/(\text{score}_{\text{max}} - \text{score}_{\text{random}})$. We use the resulting MinAtar Normalised Score whenever we report performance aggregates over the games.

Network architecture. The default network architecture used for MinAtar experiments is shown in Table E.5.

E.2.5 Atari experiments

Evaluation protocols on Atari. Since we mostly compare our ALE results with the Rainbow agent, we adopt the evaluation protocol from Hessel et al. [251]. Every 250K training steps in the environment we suspend the learning and evaluate the agent on 125K steps (or 500K frames). All the agents we train on ALE follow this validation protocol, the only difference being the validation epsilon value: $\epsilon = 0.001$ for C51 and DQN-Adam that we directly compare to Rainbow, which uses the same value we use.

DQN-Adam. We list the full details in Table E.4 especially since these hyperparam-

GAME	MAX	RANDOM
Asterix	78.90	0.49
Breakout	122.88	0.52
Seaquest	93.91	0.09
Space Invaders	360.92	2.86

Table E.2: MinAtar maximum and random scores used for computing the MinAtar Normalised Score.

Network hyperparameter	Value
Number of convolutional channels	16
Filter size for convolutional channels	3×3
Stride for convolutional channels	1
Hidden units for fully connected layers	128
Number of fully connected layers	2
Number of convolutional layers	2
Number of outputs	Number of actions

Table E.3: Network used for MinAtar experiments. The last fully connected layer has an output size given by the number of actions.

eters differ considerably from the the original DQN agent; since we also wanted to be able to compare our results with those of Rainbow we use similar hyperparameters to those in Hessel et al. [251].

Network architecture. The network architecture used for Atari experiments is shown in Table E.5.

E.3 Additional hyperparameter sweeps

Network architecture results from other sweeps are shown in Table E.6.

HYPER-PARAMETER	VALUE
discount γ	0.99
update frequency	4
target update frequency	8000
starting ϵ	1.0
final ϵ	0.01
ϵ steps	250000
ϵ schedule	linear
warmup steps	20000
replay size	1M
batch size	32
history length	4
cost function	SMOOTHL1LOSS
optimiser	ADAM
learning rate h	0.00025
damping term ϵ	0.0003125
β_1, β_2	(0.9, 0.999)
validation steps	125000
validation ϵ	0.001

Table E.4: DQN-Adam hyperparameters.

Network hyperparameter	Value
Number of convolutional channels	32, 64, 64
Filter size for convolutional channels	8×8 , 4×4 , 3×3
Stride for convolutional channels	4, 2, 1
Hidden units for fully connected layers	512
Number of fully connected layers	2
Number of outputs	Number of actions

Table E.5: Network used for Atari experiments, taken from Hessel et al. [109]. The last fully connected layer has an output size given by the number of actions.

Experiment	Nr. of conv layers	Conv width	FC width
Sweep over optimisation hyperparameters (Figure 7.4.)	1	24	128
	2	24	128
	3	24	128
	4	24	128
	2	32	256
	3	32	256
Regularisation (Figure E.2)	1	16	64
	3	16	64
	1	24	128
	3	24	128
	1	32	256
	3	32	256
Optimisation methods (Figure E.5)	1	24	128
	2	24	128
	1	32	256
	2	32	256
Smoothness analysis (Figures 7.3 E.1 and E.1b, Tables 7.2 and 7.1)	1	8	32
	2	8	32
	1	32	256
	3	8	32
	2	32	256
	3	32	256
	3	24	128
	2	16	64
	2	24	128
	1	24	128
	3	16	64
	1	16	64

Table E.6: Details of neural architecture sweeps. The total number of critic layers is the number of convolutional layers specified here, followed by two linear layers.

E.4 Individual figure reproduction details

- Figure 7.1 uses the default Atari network shown in Table E.5. The DQN-Adam hyperparameters as shown in Table E.4. The evaluation protocol is described in Section E.2.5. The C51 hyperparameters are the same as those in Hessel et al. [251].
- Figure 7.3 performs a sweep over depth and width of networks on MinAtar. The optimisation and reinforcement learning hyperparameters as shown in Table E.1. The model sweep is described in Table E.6.
- Figure 7.2 presents both Atari and MinAtar results. The Atari results in Figure 7.2a use the default Atari network shown in Table E.5. The DQN-Adam hyperparameters as shown in Table E.4. The Figure 7.2b uses the hyperparameters as shown in Table E.1 and the network in Table E.3.
- Tables 7.1 and 7.2 use data from the same experiment as Figure 7.3.
- Figure 7.4 uses a MinAtar sweep over architectures detailed in Table E.6, with a sweep over learning rates $h \in \{0.00001, \dots, 0.00215\}$ and $\epsilon \in \{0.00001, \dots, 0.01\}$. The remaining hyperparameters are the ones in Table E.1.
- Figure 7.5 used the default MinAtar optimisation and reinforcement learning hyperparameters as shown in Table E.1, and the network in Table E.3.
- Figures 7.6, 7.12, 7.13 use the architectures in Spectral Normalised GAN [30], with the Adam default hyperparameters with the exception of $\beta_1 = 0.5$, and learning rates sweeps as a cross product of $\{0.0001, 0.0002, 0.0003, 0.0004\}$ for the discriminator and generator. Simultaneous updates are used.
- Figure 7.8 uses the architectures in Spectral Normalised GAN [30], with the Adam default hyperparameters with the exception of $\beta_1 = 0.5$, and learning rates sweeps as a cross product of $\{0.0001, 0.0002, 0.0003, 0.0004\}$ for the discriminator and generator. The generator uses the non-saturating loss. Alternating updates are used.

- Figures 7.9 and 7.10 use the architectures in Spectral Normalised GAN [30], with the Adam default hyperparameters with the exception of $\beta_1 = 0.5$, and learning rates 0.0001 for both the discriminator and generator. The generator is trained using the non-saturating loss. Alternating updates are used.
- Figure 7.15 uses the architectures in Spectral Normalised GAN [30], with the Adam default hyperparameters with the exception of $\beta_1 = 0.5$, and learning rates 0.0004 for both the discriminator and generator. The ϵ values are provided in each subfigure. The generator is trained using the non-saturating loss. Simultaneous updates are used.
- Figures E.1 and E.1b use a sweep over architectures detailed in Table E.6. The optimisation and reinforcement learning hyperparameters as shown in Table E.1.
- Figure E.2 used a sweep over architectures detailed in Table E.6. The optimisation and reinforcement learning hyperparameters as shown in Table E.1.
- Figure E.3 uses the hyperparameters as shown in Table E.1 and the network in Table E.3.
- Figure E.5 is a different visualisation of the same experiment shown in Figure 7.5. The architectures are those in Table E.6.

Appendix F

Geometric Complexity: a smoothness complexity measure implicitly regularised in deep learning

F.1 Figure and experiments reproduction details

- Figure 8.2: the MLP shown have 500 units per layer. GC is computed using 100 data points.
- Figure 8.3: the MLP shown have 300 units per layer and 3 layers. The learning rate used is 0.02.
- Figure 8.4: The network used is a Resnet-18, with a batch size of 512 and a learning rate of 0.02. Each experiment uses 3 seeds. The activation function is ReLu.
- Figure 8.5: The network used is a Resnet-18, with a batch size of 128 and a learning rate of 0.02. Each experiment uses 3 seeds. The activation function is ReLu.
- Figure 8.6: The network used is a Resnet-18, trained with batch size 512. Each experiment uses 3 seeds. The activation function is ReLu.

- Figure 8.7: The network used is a Resnet-18, trained with learning rate 0.02. Each experiment uses 3 seeds. The activation function is ReLu.
- Figure 8.8: The network used is a Resnet-18, trained with learning rate 0.02 and batch size 512. The activation function is ReLu.

F.2 Additional experimental results

We present additional the implicit and explicit regularisation experiments using gradient descent *with momentum*, which is widely used in practice. We observe the same phenomena as with gradient descent without momentum, namely that more implicit regularisation through decreased batch size and increased learning rate (Figure F.1) and explicit regularisation using gradient norm, spectral and L2 regularisation (shown in Figures F.2, F.3 and F.4 respectively) produces solutions with higher test accuracy and lower Geometric Complexity (GC).

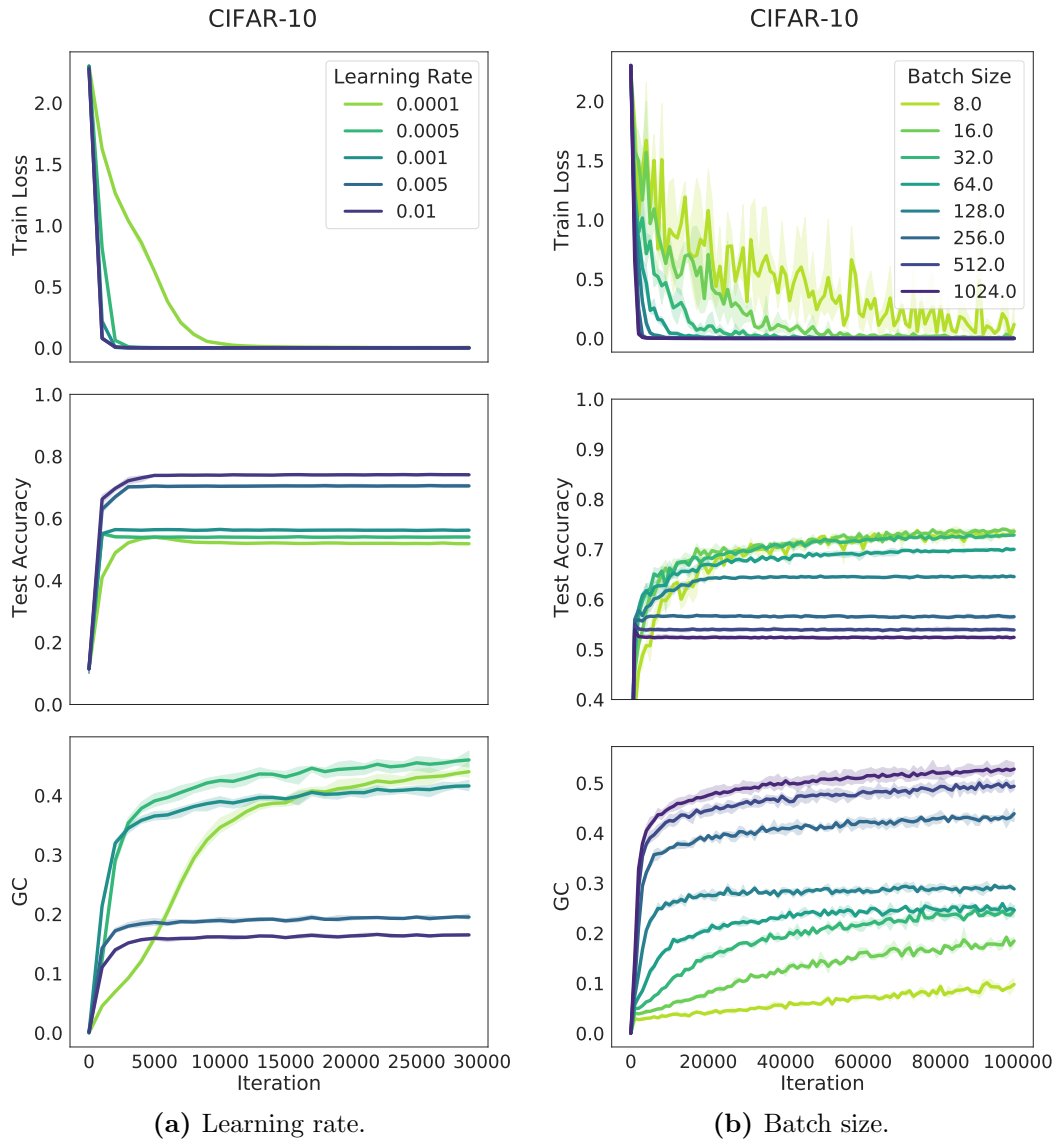


Figure F.1: The implicit regularisation effect of increased learning rates (a) and decreased batch sizes (b) leads to decreased GC when momentum is used. The network used is a Resnet-18, the momentum rate is 0.9. For the learning rate experiments, the batch size used is 512. For the batch size experiments, the learning rate used is 0.005. Each experiment uses 3 seeds.

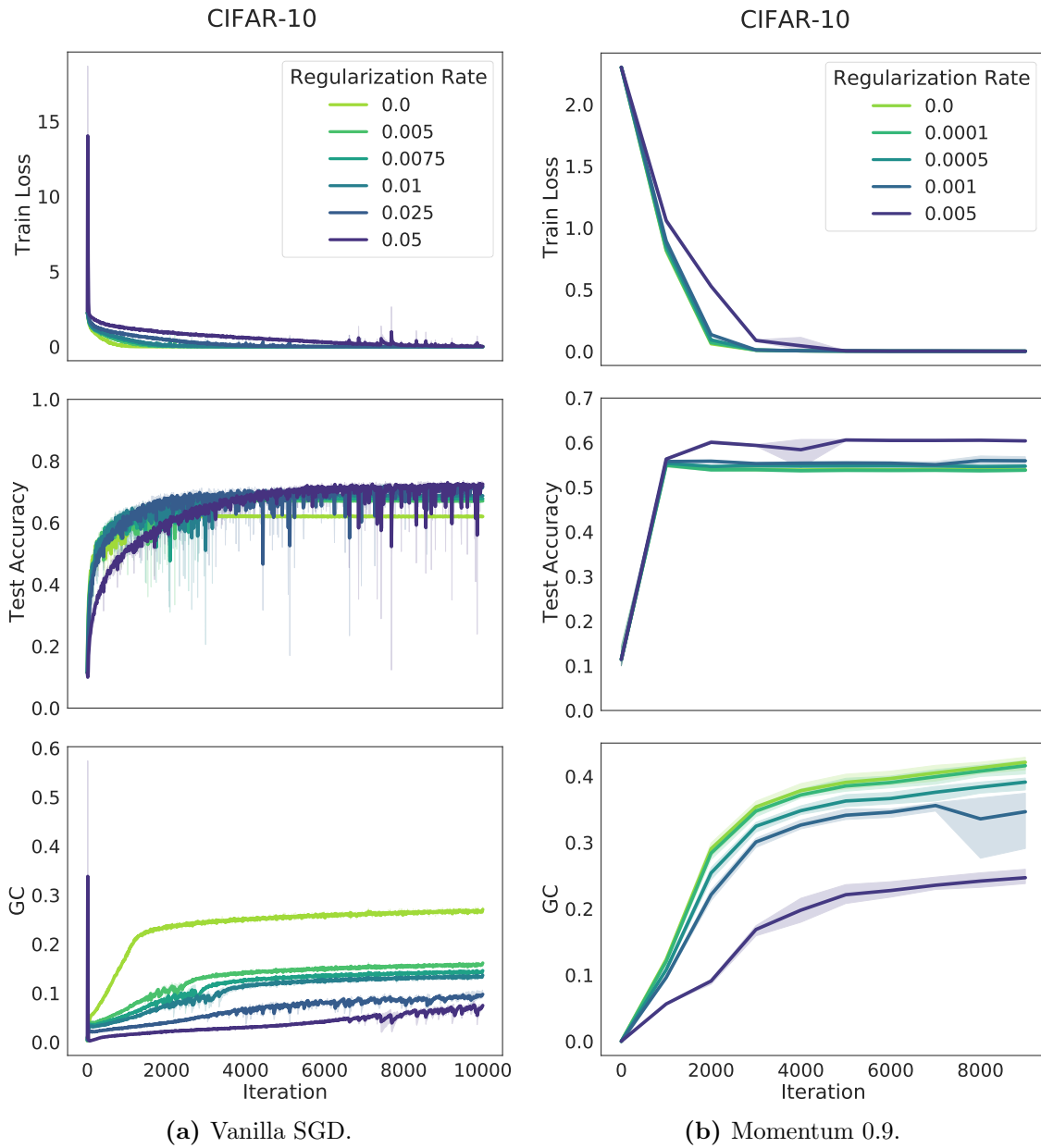


Figure F.2: Gradient norm regularisation leads to decreased Geometric Complexity when using stochastic gradient descent, with and without momentum.. The network used is a Resnet-18, with a batch size of 512 and a learning rate of 0.005 for momentum, and 0.02 for vanilla gradient descent. Each experiment uses 3 seeds.

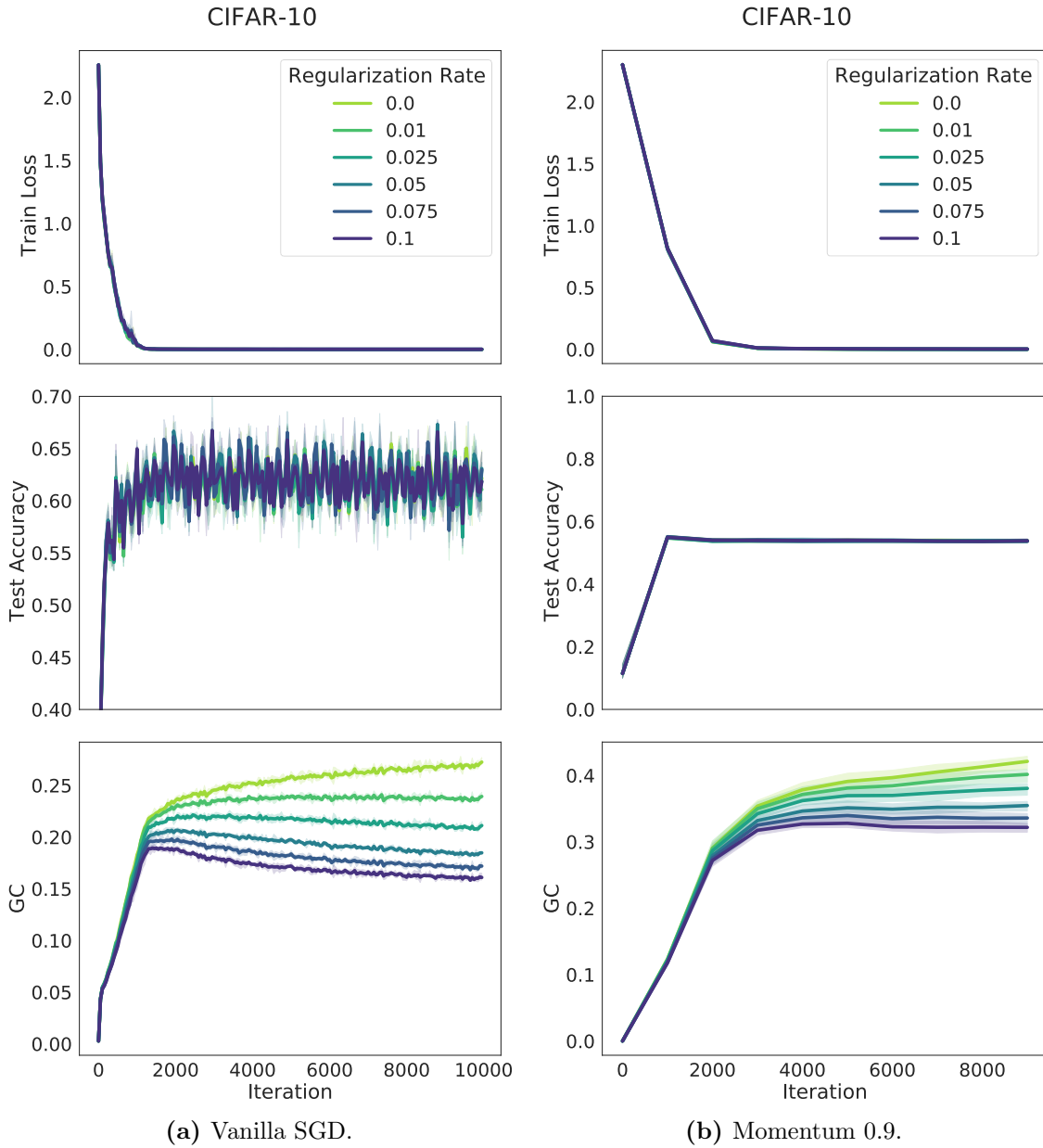


Figure F.3: Spectral regularisation leads to decreased GC when using stochastic gradient descent, with (b) and without momentum (a). The network used is a Resnet-18, with a batch size of 512 and a learning rate of 0.005 for momentum, and 0.02 for vanilla gradient descent. Each experiment uses 3 seeds.

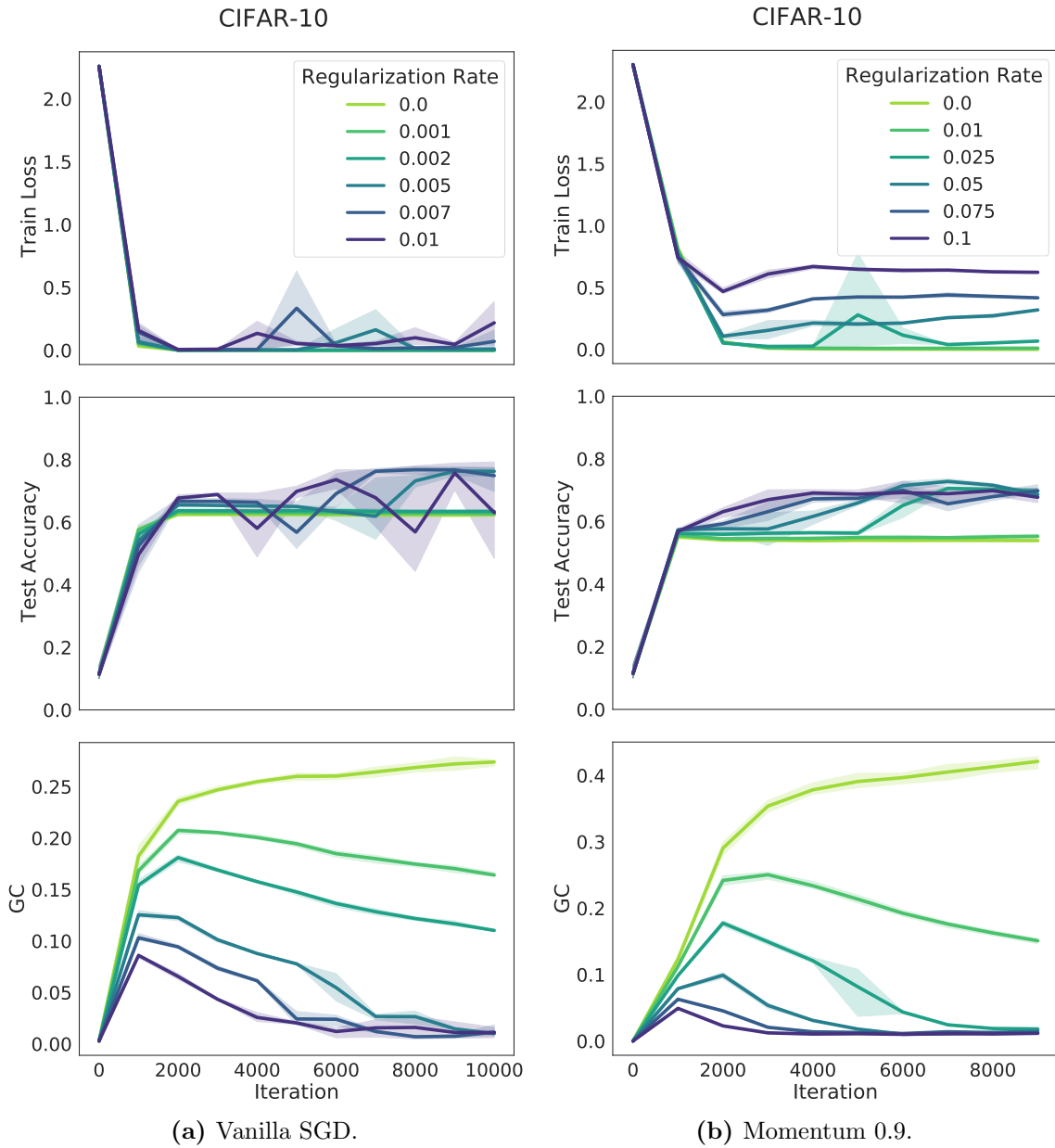


Figure F.4: L_2 regularisation leads to decreased Geometric Complexity when using stochastic gradient descent, with (b) and without momentum (a). The network used is a Resnet-18, with a batch size of 512 and a learning rate of 0.005 for momentum, and 0.02 for vanilla gradient descent. Each experiment uses 3 seeds.