

DeepMind

# Smoothness constraints in Deep Learning

Mihaela Rosca

Research Engineer at DeepMind, PhD student at UCL

02/03/2022



---

# A case for new neural network smoothness constraints

---

Mihaela Rosca<sup>1,2</sup> Theophane Weber<sup>1</sup> Arthur Gretton<sup>2</sup> Shakir Mohamed<sup>1</sup>  
<sup>1</sup>DeepMind <sup>2</sup>University College London  
{mihaelacr,theophane,shakir}@google.com, arthur.gretton@gmail.com

## Abstract

How sensitive should machine learning models be to *input* changes? We tackle the question of model smoothness and show that it is a useful inductive bias which aids generalization, adversarial robustness, generative modeling and reinforcement learning. We explore current methods of imposing smoothness constraints and observe they lack the flexibility to adapt to new tasks, they don't account for data modalities, they interact with losses, architectures and optimization in ways not yet fully understood. We conclude that new advances in the field are hinging on finding ways to incorporate *data*, *tasks* and *learning* into our definitions of smoothness.

# References

References are provided at the end of the talk, additional references can be found in the paper.

If I missed out any reference, please let me know.



1

# What is smoothness?



# Smoothness with respect to inputs

This talk is about smoothness with respect to  $\mu$  ~~inputs~~ not parameters.

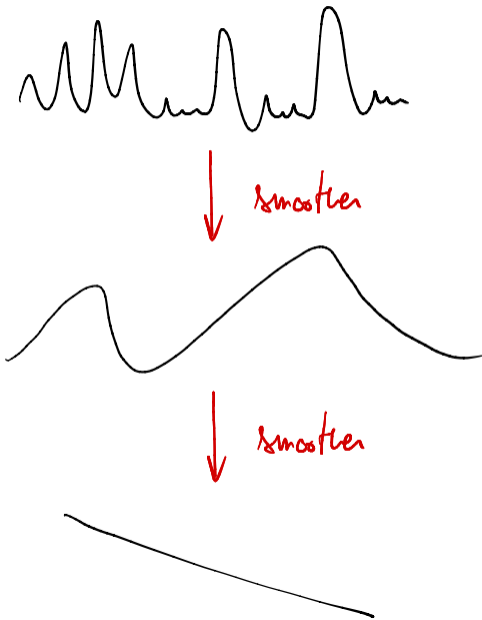
$\mu ; )$



We are talking about how the model's  $\hat{C}_a$ , not the loss changes with changes in input.



## Intuition



# How do we measure smoothness?

A few intuitive ways:

Lipschitz smoothness. A function is  $D$ -Lipschitz if

$$\| \mu(x_1) - \mu(x_2) \|_Y \leq D \| x_1 - x_2 \|_X \quad \forall x_1, x_2 \in X \quad (1)$$

Rademacher's theorem: if  $X \subseteq \mathbb{R}^d$  is an open set and  $Y = \mathbb{R}^k$  and  $\mu$  is  $D$ -Lipschitz then  $\| S\mu(x) \|_Y \leq D$  wherever the total derivative  $S\mu(x)$  exists.

The norm of the model Jacobian  $\| S\mu(x) \|_Y = \frac{\| \dot{\mu}(x) \|_Y}{\| \dot{x} \|_X}$ . Jacobian metrics account for how  $\| \dot{\mu}(x) \|_Y$  of the function output is allowed to vary as individual input dimensions change.

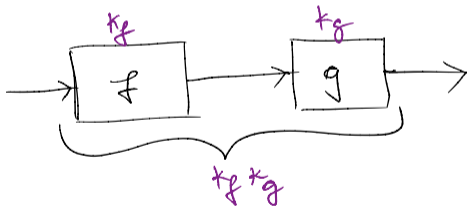




## Composition of Lipschitz functions

If  $\mu$  and  $\eta$  are Lipschitz with constants  $D_\mu$  and  $D_\eta$ ,  $\mu \circ \eta$  is Lipschitz with constant  $D_\mu D_\eta$ .

Since commonly used activation functions are 1-Lipschitz, the task of ensuring a neural network is Lipschitz reduces to constraining the learnable layers to be Lipschitz.



# The Lipschitz constant of linear operators

The Lipschitz constant of a linear operator under common norms  $(\ell_1, \ell_2, \ell_\infty)$  is  $\sup_{\|x\|=1} \|Ax\|$ .

Many neural networks layers are linear operators:

- linear layers

- convolutional layers

- BatchNorm



DeepMind

2

# Smoothness constraints for neural networks

02/03/2022



# Indirect smoothness constraints

A lot of common regularization methods indirectly target smoothness:

- early stopping

- dropout

- weight decay

- data augmentation

While interesting in their own right, in this talk we will primarily focus on methods which explicitly target smoothness regularisation.



## Soft constraints: gradient penalties

Soft constraints add a regularisation term to the loss function to encourage Lipschitz smoothness, by adding a gradient penalty to the loss function  $L(\theta)$ :

$$L(\theta) + \mathbb{E}_{\mathcal{B}_{\text{act}}(\theta)} \left( \kappa \left( \|\nabla_{\theta} L(\theta)\|_2 - D \right)^2 \right) \quad (2)$$

where

$\kappa$  is a regularization coefficient

$\mathcal{B}_{\text{act}}(\theta)$  is the distribution at which the regularization is applied, which can either be the data distribution or around it.



## Soft constraints - Spectral Regularisation

Spectral regularization uses the sum of the spectral norms – the largest singular value – of each layer as a regularization loss to encourage Lipschitz smoothness:

$$L(\theta) + \lambda \sum_l \|z_l\|_2 \quad (3)$$

where

$\lambda$  is a regularization coefficient

$\|z_l\|_2$  is the spectral norm of  $z_l$ , computed using power iteration.

For convolutional layers, the weights get reshaped to a 2D matrix. This technically does not compute the Lipschitz constant of the operator, but seems good enough in practice.



## Hard constraints - Spectral Normalisation

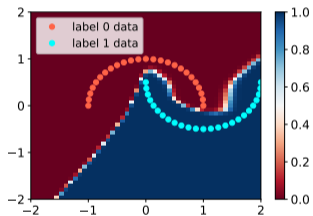
Spectral Normalization ensures the learned models are 1-Lipschitz by adding a node in the computational graph of the model layers by replacing the weights with their normalized version:

$$L(z) / L(\hat{z}) \tag{4}$$

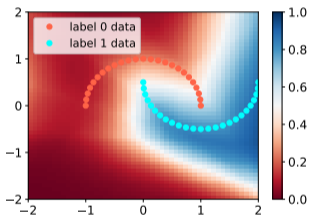
where  $\hat{z} = z / \|z\|_2$  and  $\|z\|_2$  is the spectral norm of  $z$ .



# Smoothness constraints on two moons



(a) MLP.  
No regularisation.



(b) Gradient penalty at data;  
 $D=1$ .

(c) Spectral norm;  
 $D=1$ .





# 3

## The benefits of smoothness constraints

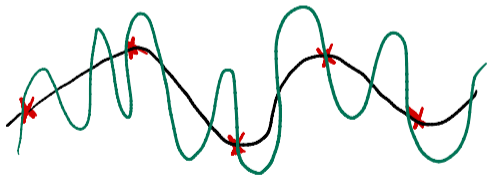
02/03/2022



# Generalisation

Methods that encourage smoothness such as weight decay, dropout, data augmentation and early stopping have been long shown to aid generalization.

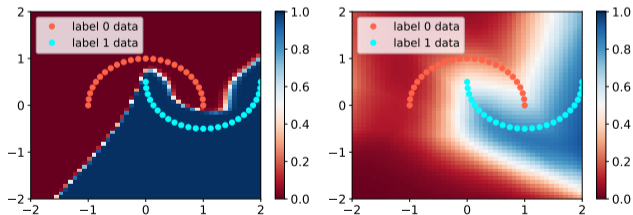
Recent works directly connect smoothness to classification margins, and use that to obtain empirical gains on standard image classification tasks.



# Reliable uncertainty estimates

Neural networks provide notoriously unreliable uncertainty estimates.

To leverage the power of neural networks to obtain reliable uncertainty estimates, by combining smooth neural feature learners with non-softmax decision surfaces.



(a) 4 layer MLP.

(b) Gradient penalty at data;  $D=1$ .



## Robustness to adversarial attacks

Robustness for classifiers can be defined by ensuring that inputs in the same  $\epsilon$ -ball result in the same function output:

$$k - \epsilon \in \mathcal{G} \Rightarrow \arg \max_{\mathcal{G}} \mu(x) = \arg \max_{\mathcal{G}} \mu(x + \delta) \quad (5)$$

This definition is directly connected with Lipschitz smoothness.

Initial approaches to combating adversarial attacks focused on data augmentation methods and only more recently smoothness constraints have come into focus.



# Improved generative modelling performance

Smoothness constraints have become part of many state of the art generative models:

GANs: Spectral Normalisation or gradient penalties are present in many GANs.

Variational Autoencoders: Spectral regularization boosts performance and stability.

Normalising Flows: benefit from smoothness constraints through powerful invertible layers built using residual connections  $\mathcal{F}(x) = x + \mu(x)$  where  $\mu$  is Lipschitz.

Figure: GAN performance is improved when the discriminator uses Spectral Normalisation.



## More informative critics

Critics (e.g. GAN, Wasserstein, Jensen-Shannon, etc.) have become more and more important in machine learning:

Generative models: The GAN critic is used to approximate distributional divergences and distances.

Representation learning: parametric critics are trained to approximate another intractable quantity, the mutual information, using the Donsker–Varadhan or similar bounds.

Reinforcement learning: parametric critics are used to approximate value and state-value functions.



## More informative critics

Smooth critics provide more informative models the are training:

Generative models: Smooth approximations to decision surfaces of  $\mu$ -divergences provide useful gradients when the underlying divergence does not.

Representation learning: tighter bounds do not lead to better representations; the success of these methods is attributed to the inductive biases of the critics.

Reinforcement learning: see next talk.



# 4

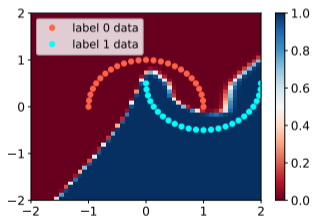
## The downsides of smoothness constraints



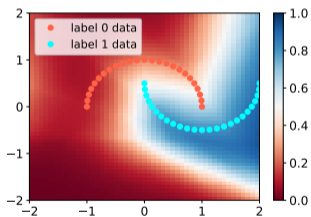


## Weak models

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.



(a) 4 layer MLP.



(b) Gradient penalty at data;  $D=1$ .



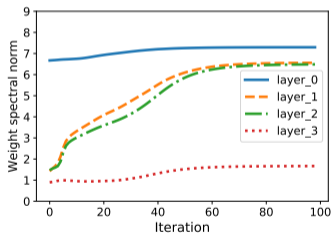
(c) Spectral normalisation  $D=1$ .

Figure: Smoothness constraints can limit model capacity and decrease performance.

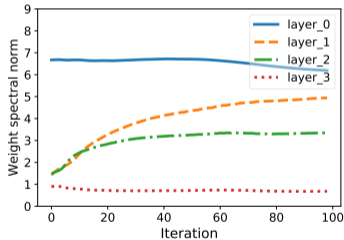


# Weak models

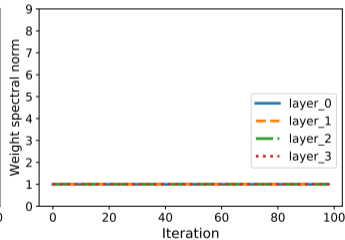
Soft, local methods like gradient penalties which only apply regularisation in one part of the space can be less restrictive.



(a) Unregularized.



(b) Gradient penalty.  $D=1$ .



(c) Spectral Normalization.  $D=1$ .

Figure: Lipschitz constant of each layer of an MLP trained on the two moons dataset. Smaller means smoother.



# Overlooked interactions with optimization

Smoothness has been traditionally seen as changing the  $\hat{I} \hat{O} \hat{E}$ . We show here that smoothness has strong interactions with optimisation (more in the next talk!).

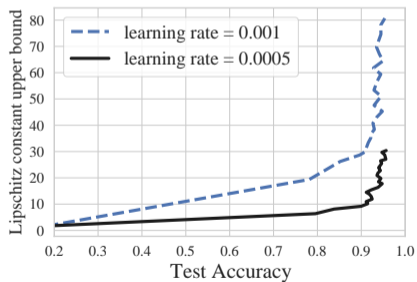
Some smoothness regularization techniques affect optimization by changing the loss function (gradient penalties, spectral regularization) or the optimization regime directly (early stopping).

Even if they don't explicitly change the loss function or optimization regime,  $\hat{I} \hat{O} \hat{E}$



# Overlooked interactions with optimization

Training with different learning rates leads to different smoothness properties of models; imposing the  $L_{\infty}$  constraint on  $\hat{I}(\hat{O})$  trained with different learning rates will have vastly different outcomes.



Lower means smoother.



## Overlooked interactions with optimization

Other hyperparameters, like momentum, are also affected by smoothness constraints.

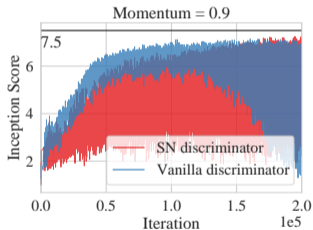
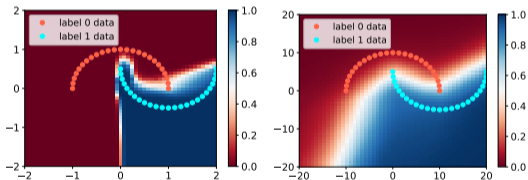


Figure: Spectral Normalization requires low momentum in GAN training. Higher is better.



# Sensitivity to data scaling

Sensitivity to data scaling of smoothness constraints can make training neural network models sensitive to additional hyperparameters.



(a) Spectral norm;  
 $D=1$ .

(b) Spectral norm;  
 $D=10$ .

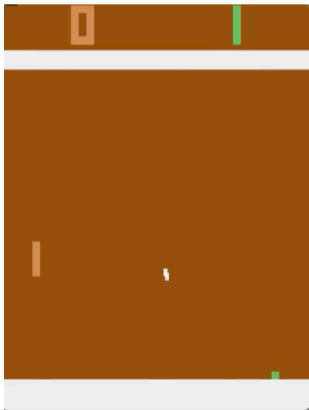
(c) Spectral norm; scaling the  
data by 10.

Figure: Changing the scale of the data or the Lipschitz constant can lead to vastly different results.



## Wrong model priors

Depending on the task, smoothness might not be the right model prior. In reinforcement learning, one pixel change might require a big change in the value function.



DeepMind

5

# The future of smoothness con- straints

02/03/2022





## New ways of defining smoothness

Improving model generalization and robustness requires specifying the right level of invariance by using task information to define smoothness constraints.

We have to ask what are the desired properties of  $\phi$  such that

$$\|\mu^{\phi}(x) - \mu^{\phi}(y)\| \leq \|\phi(x) - \phi(y)\| \quad (6)$$

To ensure the mapping  $\phi$  does not discard task relevant information in the data, maintains useful diversity and accounts for input modalities, it has to be ~~YOLO~~ and ~~ICAE~~ dependent.

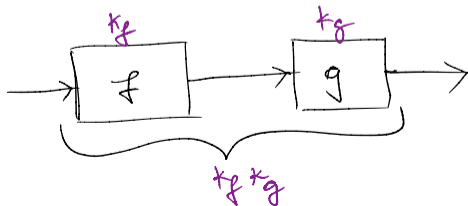


## New ways of measuring smoothness

Measuring smoothness of a function parametrized by a neural network is challenging even for the most common measure of smoothness used in machine learning, Lipschitzness.

Currently, we only have loose upper bounds available, or more accurate methods which are very costly.

To further improve the effect of smoothness regularisation methods, we have to understand them better and measure smoothness more accurately.

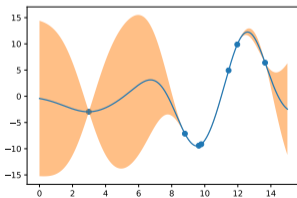


# New learning paradigms

Combining non parametric methods with feature learning is a promising approach to learning smooth decision surfaces. Requires:

- learning the right features (which themselves might have to be smooth)

- scaling non parametric methods such as Gaussian Processes, Support Vector Machines and Nearest Neighbours methods to large datasets.



DeepMind

6

# Conclusion

02/03/2022



# 7

# References

Please see the paper for a comprehensive list of references.  
If I missed anything, please let me know.

02/03/2022



## References - Methods

Gradient penalties and their use in GANs [SGSR17, GAA<sup>+</sup>17, FRL<sup>+</sup>18, ASBG18, KAHK17]

Spectral Regularisation [YM17]

Spectral Normalisation [MKKY18]



## References - Benefits

Generalisation [SGSR17, Bar97, GRS18, HRS16, NKB<sup>+</sup> 19, SHK<sup>+</sup> 14]

Reliable uncertainty estimates [vASTG20, LLP<sup>+</sup> 20]

Robustness to adversarial attacks [CBG<sup>+</sup> 17, NBA<sup>+</sup> 18, SGSR17, LGO18]

Improved generative modeling performance [MKKY18, BDS18, BGC<sup>+</sup> 19, VK20]

More informative critics [FRL<sup>+</sup> 18, AZG20, SZA19, ASBG18, ACB17, GAA<sup>+</sup> 17, FRL<sup>+</sup> 18, BDS18, ASBG18, ZLS<sup>+</sup> 19, YM17, TDR<sup>+</sup> 20, DJ20]



## References - Future

Weak models [JBZB18, FRH<sup>+</sup> 19]

Interactions with optimisation [GAA<sup>+</sup> 17]

New Ways Of Measuring smoothness [VS18, CP19, FRH<sup>+</sup> 19, SGSR17]





# References I

Martin Arjovsky, Soumith Chintala, and Léon Bottou, *z*, *Proceedings of the 34th International Conference on Machine Learning–Volume 70*, 2017, pp. 214–223.

Michael Arbel, Dougal Sutherland, Mikołaj Bińkowski, and Arthur Gretton, *RĪ*, *Advances in neural information processing systems*, 2018, pp. 6700–6710.

Michael Arbel, Liang Zhou, and Arthur Gretton, *DDP*, *arXiv preprint arXiv:2003.05033* (2020).

Peter L Bartlett, *30*, *Advances in neural information processing systems*, 1997, pp. 134–140.

Andrew Brock, Jeff Donahue, and Karen Simonyan, *FCM*, *International Conference on Learning Representations*, 2018.



## References II

Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen, *Learning to Optimize: A Gradient Descent Approach*, International Conference on Machine Learning, 2019, pp. 573–582.

Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier, *Fast Neural Architecture Search*, Proceedings of the 34th International Conference on Machine Learning–Volume 70, 2017, pp. 854–863.

Patrick L Combettes and Jean-Christophe Pesquet, *Primal-Dual Splitting for Image Restoration*, arXiv preprint arXiv:1903.01014 (2019).

Pierluca D’Oro and Wojciech Jaskowski, *Fast Gradient Descent for Image Denoising*, arXiv preprint arXiv:2004.14309 (2020).

Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas, *Learning to Optimize: A Gradient Descent Approach*, Advances in Neural Information Processing Systems, 2019, pp. 11427–11438.



## References III

William Fedus, Mihaela Rosca, Balaji Lakshminarayanan, Andrew M Dai, Shakir Mohamed, and Ian Goodfellow, *Learning to generate with noisy neighbors*, International Conference on Learning Representations, 2018.

Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville, *Improved techniques for training GANs*, Advances in neural information processing systems, 2017, pp. 5767–5777.

Noah Golowich, Alexander Rakhlin, and Ohad Shamir, *Learning from noisy neighbors*, Conference On Learning Theory, PMLR, 2018, pp. 297–299.

Moritz Hardt, Ben Recht, and Yoram Singer, *Learning from noisy neighbors*, International Conference on Machine Learning, PMLR, 2016, pp. 1225–1234.

Joern-Henrik Jacobsen, Jens Behrmann, Richard Zemel, and Matthias Bethge, *Learning to generate with noisy neighbors*, International Conference on Learning Representations, 2018.



## References IV

Naveen Kodali, Jacob Abernethy, James Hays, and Zsolt Kira, *RI-YOI*  $\hat{u}$   $\hat{c}$   $\hat{a}$   $\hat{f}$   $\hat{c}$   $\hat{i}$   $\hat{y}$   $\hat{c}$   $\hat{e}$   $\hat{y}$   $\hat{a}$   $\hat{t}$   $\hat{c}$   $\hat{e}$   $\hat{b}$   $\hat{e}$   $\hat{n}$   $\hat{a}$   $\hat{t}$   $\hat{c}$   $\hat{e}$   $\hat{a}$   $\hat{r}$   $\hat{x}$   $\hat{i}$   $\hat{v}$  preprint arXiv:1705.07215 (2017).

Carlos Eduardo Rosar Kos Lassance, Vincent Gripon, and Antonio Ortega, *FCBLYC*  $\hat{e}$   $\hat{i}$   $\hat{c}$   $\hat{u}$   $\hat{a}$   $\hat{t}$   $\hat{e}$   $\hat{i}$   $\hat{o}$   $\hat{d}$   $\hat{i}$   $\hat{y}$   $\hat{i}$   $\hat{f}$   $\hat{i}$   $\hat{y}$   $\hat{y}$   $\hat{c}$   $\hat{o}$   $\hat{a}$   $\hat{\mu}$   $\hat{o}$   $\hat{i}$   $\hat{y}$   $\hat{o}$   $\hat{i}$   $\hat{a}$   $\hat{t}$   $\hat{o}$   $\hat{o}$   $\hat{i}$   $\hat{c}$   $\hat{a}$   $\hat{e}$   $\hat{y}$   $\hat{c}$   $\hat{a}$   $\hat{i}$   $\hat{c}$   $\hat{o}$   $\hat{a}$   $\hat{c}$   $\hat{e}$   $\hat{i}$   $\hat{c}$   $\hat{u}$   $\hat{a}$   $\hat{t}$   $\hat{e}$   $\hat{i}$   $\hat{o}$   $\hat{d}$   $\hat{a}$   $\hat{i}$   $\hat{c}$   $\hat{a}$   $\hat{r}$   $\hat{x}$   $\hat{i}$   $\hat{v}$  preprint arXiv:1805.10133 (2018).

Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshminarayanan, *dI BIC EY Bai YBEY di Yca dI i caI EO ü i ° YQ ca i i a» Y YCCB ECI i f uCY aCEY cE ECI cae* arXiv preprint arXiv:2006.10108 (2020).

Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida, *dB©YiâE*  $\hat{i}$   $\hat{c}$   $\hat{a}$   $\hat{t}$   $\hat{e}$   $\hat{i}$   $\hat{o}$   $\hat{d}$   $\hat{i}$   $\hat{y}$   $\hat{i}$   $\hat{f}$   $\hat{i}$   $\hat{y}$   $\hat{y}$   $\hat{c}$   $\hat{o}$   $\hat{a}$   $\hat{\mu}$   $\hat{o}$   $\hat{i}$   $\hat{y}$   $\hat{o}$   $\hat{i}$   $\hat{a}$   $\hat{t}$   $\hat{o}$   $\hat{o}$   $\hat{i}$   $\hat{c}$   $\hat{a}$   $\hat{e}$   $\hat{y}$   $\hat{c}$   $\hat{a}$   $\hat{i}$   $\hat{c}$   $\hat{o}$   $\hat{a}$   $\hat{c}$   $\hat{e}$   $\hat{i}$   $\hat{c}$   $\hat{u}$   $\hat{a}$   $\hat{t}$   $\hat{e}$  International Conference on Learning Representations, 2018.

Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein, *dCí a» üi CEY fCí caE EO i I cãCEI cüâAe E CÍ BãYCaòY ,* International Conference on Learning Representations, 2018.



## References V

Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever, *Scaling Vision Transformers: Efficient Data and Architectural Design*, International Conference on Learning Representations, 2019.

Jure Sokolić, Raja Giryes, Guillermo Sapiro, and Miguel RD Rodrigues, *Fast and Accurate Face Recognition in the Wild*, IEEE Transactions on Signal Processing (2017), no. 16, 4265–4280.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, *Dropout: A Simple Way to Prevent Neural Networks from Overfitting*, The journal of machine learning research (2014), no. 1, 1929–1958.

Florian Schäfer, Hongkai Zheng, and Anima Anandkumar, *Learning to Detect Faces in the Wild: A Deep Face Representation with a Hierarchical Architecture*, arXiv preprint arXiv:1910.05852 (2019).

Michael Tobias Tschannen, Josip Djolonga, Paul Kishan Rubenstein, Sylvain Gelly, and Mario Lučić, *Improving Self-Supervised Representation with Contrastive Learning*, International Conference on Learning Representations, 2020.



## References VI

Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal, *Learning to Rank with Gradient Descent*, Proceedings of the 37th International Conference on Machine Learning (2020).

Arash Vahdat and Jan Kautz, *Learning to Rank with Gradient Descent*, Advances in Neural Information Processing Systems (2020).

Aladin Virmaux and Kevin Scaman, *Learning to Rank with Gradient Descent*, Advances in Neural Information Processing Systems, 2018, pp. 3835–3844.

Yuichi Yoshida and Takeru Miyato, *Learning to Rank with Gradient Descent*, arXiv preprint arXiv:1705.10941 (2017).

Zhiming Zhou, Jiadong Liang, Yuxuan Song, Lantao Yu, Hongwei Wang, Weinan Zhang, Yong Yu, and Zhihua Zhang, *Learning to Rank with Gradient Descent*, International Conference on Machine Learning, 2019, pp. 7584–7593.

