Unsupervised learning

EEML 2020 Mihaela Rosca DeepMind & University College London

Mihaela Rosca, EEML 2020

Disclaimer

- This is one view of the literature in the field, there are many others.
- Each citation is meant to be used as a representative example. Further related work can be found using <u>connectedpapers.com</u>.
- Email me with pointers or suggestions and I will update the slides.

Unsupervised learning

Aim: learn structure from data.

Types of learning

Supervised learning

Learn a mapping from input **x** to output **y**.

Challenge: generalization, having a flexible enough parametrization to learn the mapping.



Reinforcement learning

Learn behaviours to maximize rewards.

Challenge: finding rewarding behaviour (exploration), generalization, transfer.



Unsupervised learning

Learn structure from data.

Challenge: No labels, no rewards. Generalization.



Unsupervised learning is hard



Unsupervised learning is hard



Mihaela Rosca, EEML 2020

Unsupervised learning is hard





Types of unsupervised learning

- Clustering
- Generative modeling
- Representation learning

Often the lines can be blurry.

Clustering



Mihaela Rosca, EEML 2020

Generative modeling

Learn a model of the true underlying data distribution $p^*(x)$ from samples

Generative modeling

Learn a model of the true underlying data distribution $p^*(x)$ from samples



Representation learning



The success of deep learning tells us about the importance of learning representations.

Easier for downstream tasks to work with learned representations rather than high dimensional data.

Representation learning



How to do unsupervised learning: a recipe

- Find an objective
- Find a model
- Find a way to learn the model using your objective
- Find an evaluation metric

The problem

Case study: generative modeling

Learn a model of the true underlying data distribution $p^*(x)$ from samples



What can we do with a distribution?

Query it



Sample from it



Types of data

Continuous data:

- Images
- Audio/ Speech
- Health data:
 - Age
 - Weight

Discrete data:

- Text
- Images
- Health data:
 - # of times someone was admitted to hospital
 - is smoker?

The objective

Measuring distances between distributions



How can we measure the distance between these two distributions?

Measuring distances between distributions





Caveat: we only have samples from the true distribution.

Mihaela Rosca, EEML 2020

Monte Carlo estimation

How can we incorporate the data distribution in the objective if we only have samples from it?



Divergence and distance minimization

- The objective of generative models is often to minimize a divergence or distance.
- Most common: Maximum likelihood (KL divergence).

Why divergence/distance minimization?

$$D(p^*||p_\theta) = 0 \implies p_\theta = p^*$$

KL divergence - maximum likelihood

$$\begin{aligned} \operatorname{KL}(p^*(\mathbf{x})||p_{\theta}(\mathbf{x})) &= \int p^*(\mathbf{x}) \log \frac{p^*(\mathbf{x})}{p_{\theta}(\mathbf{x})} d\mathbf{x} \\ &= C - \int p^*(\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

KL divergence - maximum likelihood

min
$$\operatorname{KL}(p^*(\mathbf{x})||p_{\theta}(\mathbf{x})) = \int p^*(\mathbf{x}) \log \frac{p^*(\mathbf{x})}{p_{\theta}(\mathbf{x})} d\mathbf{x}$$

= $C - \int p^*(\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{x}$

max
$$\mathbb{E}_{p^*(\mathbf{x})} \log p_{\theta}(\mathbf{x})$$

Mihaela Rosca, EEML 2020

Effects of the choice of divergence



Want to learn more?

Arxiv (2016)

Goodfellow, et al. NIPS 2016 Tutorial: Generative Adversarial Networks

Jensen Shannon divergence





Given a model family, different divergences can lead to a vastly different distribution.

Optimising Reverse KL & JSD - not so easy

$$\begin{split} \mathrm{KL}(p^*(\mathbf{x})||p_{\theta}(\mathbf{x})) &= \int p^*(\mathbf{x}) \log \frac{p^*(\mathbf{x})}{p_{\theta}(\mathbf{x})} d\mathbf{x} \\ &= C - \int \underbrace{p^*(\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{x}}_{\mathbf{Monte Carlo estimation}} \\ \mathrm{KL}(p_{\theta}(\mathbf{x})||p^*(\mathbf{x})) &= \int p_{\theta}(\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x})}{p^*(\mathbf{x})} d\mathbf{x} \\ &= \int \underbrace{p_{\theta}(\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{x}}_{\mathbf{Fntropy hard to estimate}} - \int \underbrace{p_{\theta}(\mathbf{x}) \log p^*(\mathbf{x}) d\mathbf{x}}_{\mathbf{Need access to the true data distribution}} \\ \end{split}$$

Beyond divergence minimization - two player games

Discriminator

Learns to distinguish between real and generated data.

VS

Generator

Learns to generate data to "fool" the discriminator.





Created by Minus icons from Noun Project



Created by Minus icons from Noun Project

Generative adversarial networks



Want to learn more?

Systems (2014)

Goodfellow, et al. Generative adversarial

networks.. Neural Information Processing

Generative adversarial networks



 $\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} \mathbb{E}_{p^*(\mathbf{x})} \left[\log \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x}) \right] + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} \left[\log (1 - \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x})) \right]$ log-probability that D correctly log-probability that D correctly predicts predicts real data x are real generated data are generated

Are GANs doing divergence minimization?

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} \mathbb{E}_{p^*(\mathbf{x})} \left[\log \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x}) \right] + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} \left[\log (1 - \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x})) \right]$$

If the discriminator (D) is optimal: the generator is minimizing the Jensen Shannon divergence between the true and generated distributions.

Connection to optimality:

$$JSD(p^*||p_\theta) = 0 \implies p_\theta = p^*$$

Other divergences and distances



Wasserstein Distance

$$W(p^*, p_{\theta}) = \sup_{\|f\|_{L} \le 1} \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{p_{\theta}(\mathbf{x})} f(\mathbf{x})$$

$$|f(x) - f(y)| \le |x - y|$$

Other divergences and distances



Wasserstein Distance

$$W(p^*, p_{\theta}) = \sup_{||f||_L \le 1} \mathbb{E}_{p^*(\mathbf{x})} f(\mathbf{x}) - \mathbb{E}_{p_{\theta}(\mathbf{x})} f(\mathbf{x})$$




Other divergences and distances



Learning

$$\min_{\substack{\theta \\ ||D_{\phi}||_{L} \leq 1}} \max_{\substack{\phi, \\ ||D_{\phi}||_{L} \leq 1}} \mathbb{E}_{p^{*}(\mathbf{x})} D_{\phi}(\mathbf{x}) - \mathbb{E}_{p_{\theta}(\mathbf{x})} D_{\phi}(\mathbf{x})$$

GANs: More than divergence minimization

In practice D is not optimal: imited computational resources we do not have access to the true data distribution (just samples)

Discriminators as learned "distances"



Arora, et al Generalization and Equilibrium in Generative Adversarial Nets. ICLR (2017)

 $\min_{G} \max_{D} V(D,G)$

We can think of D (the discriminator) as learning a "distance" between the data and model distribution that can provide useful gradients to the model.

GANs (learned distance) or divergence minimization?

GANs

good samples
learned loss function

hard to analyze dynamics (game theory)
(in practice) no optimal convergence guarantees

Divergence minimization

optimal convergence guarantees

- easy to analyze loss properties
- hard to get good samples
- loss functions don't correlate with human evaluation

The model

The importance of the model (in maximum likelihood training)



overfitting

The importance of the model (in maximum likelihood training)



where we want to be

Mihaela Rosca, EEML 2020

Explicit likelihood models

Model the density p(x).



Implicit models

Do not model the density, but the sampling path.



Observed variable models



$$p_{\theta}(\mathbf{x}) = \prod_{i=0}^{N} p_{\theta}(x_i | x_{< i})$$

Latent variable models



Generation





- Challenge: learning the factor unsupervised
- Sampling is often cheap
- Representation learning
 - Inverting the generation process = inference

Ζ

Generation





Inference



Explicit models - canonical distributions

$$p_{\theta}(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mu, \Sigma)$$

- Learn parameters of canonical distribution
- Example: Gaussian, Poisson
- Pro: Easy to learn
- Con: Not very expressive, especially in high dimensions

Explicit models - mixture models

$$p_{\theta,\pi}(\mathbf{x}) = \sum_{k=1}^{K} \pi_k p_{\theta}(\mathbf{x} | \mathbf{z} = k)$$

٨

- Pro: models multi modality.
- Con: number of modes are fixed.
- Mixture components can be simple or complex distributions.

Explicit models - autoregressive models

$$p_{\theta}(\mathbf{x}) = \prod_{i=0}^{N} p_{\theta}(x_i | x_{< i})$$

- Pro: Very expressive
- Challenge: Slow at sampling (though can be parallelize)
- Modality: great for sequential data, text, audio but have also been used for images

Explicit models - normalizing flows



1 - a 1 / x





$$p_{\mathbf{y}}(\mathbf{y}) = p_{\mathbf{x}}(f_{\theta}^{-1}(\mathbf{y})) \det \left| \frac{df_{\theta}^{-1}(\mathbf{z})}{d\mathbf{z}} \right|_{\mathbf{z}=\mathbf{y}}$$

X and Y have the same dimension!

Challenge: modeling invertible functions using neural networks

Mihaela Rosca, EEML 2020

Explicit models - normalizing flows



Composing normalizing flows leads to another flow.

Simple transformations can be used to build complex distributions.

Explicit latent variable models

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$



Mihaela Rosca, EEML 2020



Explicit latent variable models

Lower bound on maximum likelihood objective (ELBO):

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \mathrm{KL}(q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

Approximate posterior

 $q_{\eta}(\mathbf{z}|\mathbf{x})$

Mihaela Rosca, EEML 2020

Implicit models - latent variable models

Directly the sampling path, without require likelihoods explicitly (no need for the sum rule).

Often not trained with maximum likelihood, but suitable for GAN training.





Learn using divergence minimization

Maximum likelihood:

$$\mathbb{E}_{p^*(\mathbf{x})}[\log p_{\theta}(\mathbf{x})]$$

To learn parameters by gradient descent:

$$\nabla_{\theta} \mathbb{E}_{p^*(x)} [\log p_{\theta}(x)] = \mathbb{E}_{p^*(x)} \nabla_{\theta} [\log p_{\theta}(x)]$$

Monte Carlo estimation

Stochastic gradient estimation

 $\nabla_{\theta} \mathbb{E}_{p_{\theta}(\mathbf{x})} f(\mathbf{x})$

Cannot put the gradient inside the expectation. But there are other ways to leverage Monte Carlo estimation to compute gradients.

A few training criteria affected

GANs

$$\begin{split} \min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} \mathbb{E}_{p^*(\mathbf{x})} \left[\log \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x}) \right] + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} \left[\log(1 - \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x})) \right] \\ \text{Bound on ML (ELBO)} \\ \\ \max_{\boldsymbol{\theta}, \eta} \mathbb{E}_{q_{\eta}(\mathbf{z}|\mathbf{x})} \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) - \mathrm{KL}(q_{\eta}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \end{split}$$

Options

Score function

- few assumptions on cost
- no assumptions on p
- often high variance
- discrete and continuous data

Pathwise

- cost needs to be differentiable
- assumptions on p
- often low variance
- continuous data



Measure valued

- few assumptions on cost
- computationally expensive
- low variance



We are interested in having Monte Carlo estimators not only for the loss, but also to estimate gradients for learning.



Often papers present **algorithms**, which are a choice of:

- objective
- model
- learning choice (parameter update rules)

Models, training and learning criteria

Explicit models are often trained by maximum likelihood:

 $\mathbb{E}_{p^*(\mathbf{x})} \log p_{\theta}(\mathbf{x})$

Autoregressive models trained by maximum likelihood

- PixelCNN/PixelRNN (image data)
- Wavenet (audio)
- GPT (text)



Figure for van den Oord, 2016.

Want to learn more?

Radford, et al. Language **Models are Unsupervised** Multitask Learners, OpenA Blog (2019)

Want to learn more?

van den Oord , et al. Conditional Image **Generation with PixelCNN** Decoders, Neurips (2016)

Want to learn more?

van den Oord , et al. Pixel **Recurrent Neural Networks** ICML (2016)

Want to learn more?

van den Oord, et al (2016)

WaveNet: A Generative Model for Raw Audio, arxiv

Implicit latent variable models & GAN training

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} \mathbb{E}_{p^*(\mathbf{x})} \left[\log \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x}) \right] + \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x})} \left[\log (1 - \mathcal{D}_{\boldsymbol{\phi}}(\mathbf{x})) \right]$$

You will often see the GAN criteria written as:

$$\min_{\theta} \max_{\phi} \mathbb{E}_{p^{*}(\mathbf{x})} \log D_{\phi}(\mathbf{x}) + \mathbb{E}_{p(\mathbf{z})} \log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))$$

This assumes:

- The GAN *model* is an implicit latent variable model (need not be).
- The model is *learned* using the pathwise estimator (need not be).

Evaluation



Theis, et al A note on the evaluation of generative models International Conference for Learning Representations (2016)

Why evaluation is hard

No evaluation metric is able to capture all desired properties.



representation learning

Evaluate based on end goal

 \Rightarrow \Rightarrow \Rightarrow

semi supervised learning: classification accuracy reinforcement learning: agent reward

data generation: human (user) evaluation

Applications

Image generation

Implicit Latent variable + GAN

Photo realistic sample quality.

Modality matters: GANs on discrete data such as text are harder to train.

Want to learn more?



Brock et all Large scale gan training for high fidelity natural image synthesis Neurips (2018)



Generative Adversarial Imitation Learning

Learn agents to imitate the behaviour of an expert (human), using a discriminator.



Want to learn more?

Ho et all, Generative Adversarial Imitation Learning Neurips (2016)

Representation learning with explicit latent variable models (GQN)



 Eslami et all, Neural scene representation and rendering, Science (2018)



Exploration in RL

Density estimation can be used to test for out of distribution data.

In RL, this can be used to provide an exploration bonus for unseen states:

have I been here before?



Want to learn more?

Ostrovski, et al Count-Based Exploration with Neural Density Models International Conference on Machine Learning (2016)
Multi task language learning



Autoregressive text models trained by maximum likelihood can be used for multiple downstream tasks.

Key: Neural architecture, billions of parameters and large amounts of data



Figure from Radford et al. (2019)

Image to image translation - CycleGAN



Zhu et all, Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, ICCV (2017)



Conclusion



You have choices! Many choices!

Options

Objective

- Divergence minimization
- Adversarial approaches

Model

- Explicit models
- Implicit models
- Observed models
- Latent variable models

Learning

• Monte Carlo estimators

