



LG Advanced Data Scientists Program

Deep Learning

[6: Generative Models]

Prof. Sungroh Yoon

Electrical & Computer Engineering | Seoul National University

© 2020 Sungroh Yoon. this material is for educational uses only. some contents are based on the material provided by other paper/book authors and may be copyrighted by them.

(last compiled at 22:47:00 on 2020/02/25)

Outline

Generative Models

Approximate Inference

Summary

References

- *Deep Learning* by Goodfellow, Bengio and Courville [▶ Link](#)
 - ▶ Chapter 19: Approximate Inference
 - ▶ Chapter 20: Deep Generative Models
- *Pattern Recognition and Machine Learning* by Bishop
 - ▶ Chapter 10: Approximate Inference
- online resources:
 - ▶ *Stanford CS231n: CNN for Visual Recognition* [▶ Link](#)
 - ▶ *CVPR 2018 GAN Tutorial* [▶ Link](#)
 - ▶ *NIPS 2016 Variational Inference Tutorial* [▶ Link](#)
 - ▶ *NIPS 2016 GAN Tutorial* [▶ Link](#)

Outline

Generative Models

Introduction

Autoregressive Models

Approximate Inference

Summary

Supervised vs unsupervised learning

	supervised	unsupervised
data	(x, y) x : data, y : label	x just data, no _____
goal	learn a <i>function</i> to map $x \mapsto y$	learn inherent <i>structure</i> of the data
examples	classification regression object detection semantic segmentation	clustering dimensionality reduction representation learning density estimation

(source: cs231n)

Discriminative vs generative models

- assume supervised learning
 - goal: learn a *function* f to map $x \mapsto y$

	discriminative	generative
goal	directly estimate $p(y x)$	estimate $p(x y)$; then deduce ¹ $p(y x)$
should evaluate	$f(x) = \operatorname{argmax}_y p(y x)$	$f(x) = \operatorname{argmax}_y p(x y)p(y)$
what's learned	decision boundary	probability _____ of data
		
examples	SVM, neural nets	Gaussian mixture, Bayes nets

(source: [Ng and Jordan, 2002]², stackoverflow)

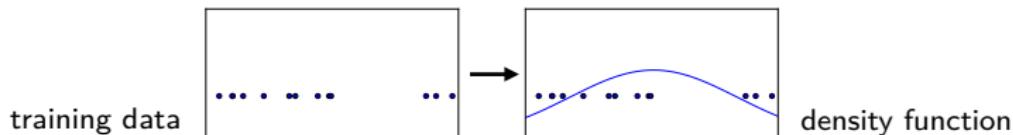
¹e.g. using Bayes rule: $p(y | x) = \frac{p(x | y)p(y)}{p(x)}$

²Ng, A. Y. and Jordan, M. I. (2002). On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In *Advances in neural information processing systems*, pages 841–848

Generative modeling in unsupervised learning

- **density estimation**

- ▶ goal: learn $p(x)$



- generation

- ▶ given training data, generate new samples from the same distribution



training data $\sim p_{\text{data}}(x)$

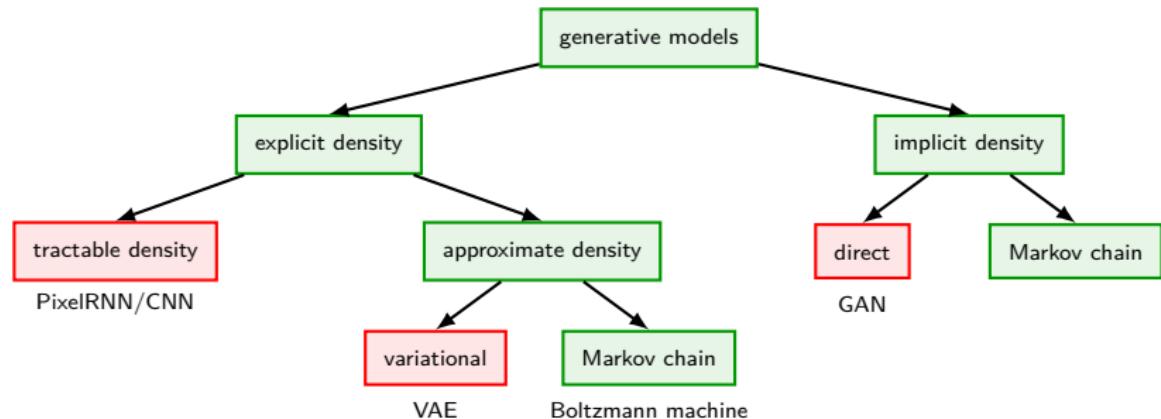


generated samples $\sim p_{\text{model}}(x)$

- ▶ want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

(source: Goodfellow (2018), Karras (2017))

Taxonomy of generative models



(source: Goodfellow, 2017)

1. explicit density estimation

- ▶ explicitly define and solve for $p_{\text{model}}(\mathbf{x})$

2. implicit density estimation

- ▶ learn a model that can sample from $p_{\text{model}}(\mathbf{x})$ w/o explicitly defining it

Trends in generative models

- conventional GM approaches have one of three drawbacks³:
 1. require strong assumptions about structure in data
 2. make approximations \Rightarrow suboptimal results
 3. rely on computationally expensive inference procedures (*e.g.* MCMC)
- recent advances
 - ▶ train **neural nets** as powerful function approximators through backprop
 \Rightarrow gives framework for **backprop-based function approximators** to build GMs

e.g. variational autoencoder: one of the most popular such frameworks

- ▶ assumptions of this model: weak
- ▶ training: fast via _____
- ▶ make an approximation but error is small given high-capacity models

³Doersch, C. (2016). Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*

- generative models based on neural nets:

1. autoregressive models (or fully visible belief nets)

- ▶ model $p(\mathbf{x})$ as $\prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$

e.g. PixelRNN, PixelCNN, WaveNet

2. Helmholtz machines

- ▶ model $p(\mathbf{x})$ as $\int p(\mathbf{z})p(\mathbf{x} | \mathbf{z})d\mathbf{z}$

- ▶ two components: _____ net (encoder) + generative net (decoder)

- ▶ use variational inference to maximize ELBO

e.g. variational autoencoder (VAE)

3. generative adversarial network (GAN)

- ▶ no explicit density modeling

- ▶ two components: generator + discriminator

- ▶ train models by solving minimax problem

Why study generative models?

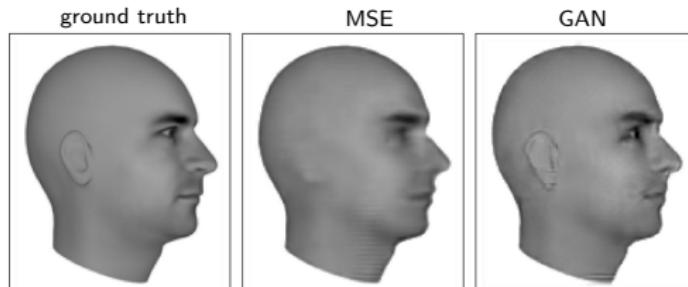
- excellent test of our ability
 - ▶ to represent/manipulate **high-dim/complicated probability distributions**
- can be incorporated into **reinforcement learning** (RL)
 - ▶ simulate possible futures for planning
- can be trained with missing data
 - ▶ can provide predictions on inputs that are missing data
 - ⇒ facilitate -supervised learning



(source: Chen+, 2018)

- enable machine learning to work with multi-modal outputs

e.g. a single input → many different correction answers



← video frame prediction

- ▶ MSE: averaging
- ⇒ blurry
- ▶ GAN: sharp

(source: Lotter+, 2015)

- enable inference of _____ representations

⇒ can be useful as general features

e.g. identity-preserving latent representations →



(source: Antipov+, 2017)

age conditional GAN

- many tasks require **realistic** generation of samples
- e.g. single image super-resolution, art creation, image-to-image translation



first AI art sold at Christie (\$432,000)

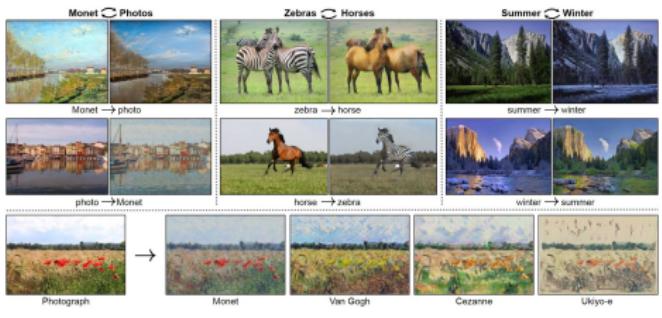


image-to-image translation



(source: Christie's, Zhu+ (2017), Brundage+ (2018))

Outline

Generative Models

Introduction

Autoregressive Models

Approximate Inference

Summary

Fully visible belief network (FVBN)

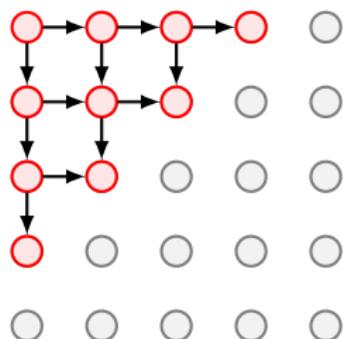
- explicit density model
- how it works:
 1. decompose likelihood of each input x into a product of 1D distributions

$$\underbrace{p(\mathbf{x})}_{\text{likelihood of input } \mathbf{x}} = \prod_{i=1}^n \underbrace{p(x_i | x_1, \dots, x_{i-1})}_{\text{probability of } i\text{-th feature given all previous features}}$$

2. maximize likelihood of training data
- complex distribution over feature values
 - ▶ we express it using a **neural net**
 - main issue:
 - ▶ need to define _____ of “previous features”

PixelRNN

- context: image (ICML 2016 best paper⁴)
- idea: the same as _____ modeling applied to image
 - ▶ generate image pixels starting from corner
 - ▶ model dependency on previous pixels using 2D LSTM
- limitation
 - ▶ sequential generation \Rightarrow slow



(source: [van den Oord et al., 2016c]⁵)

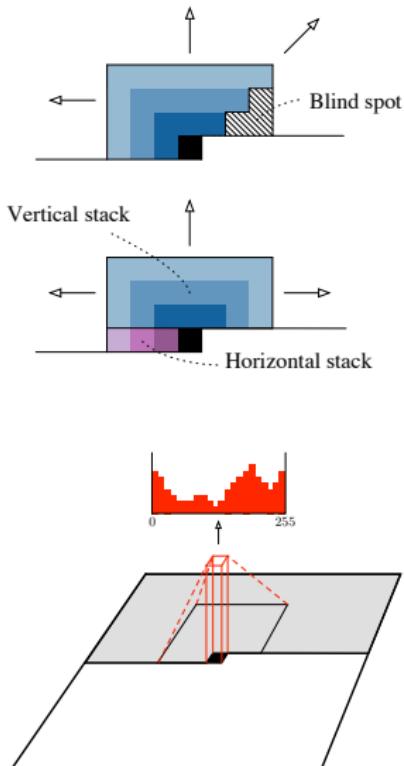
⁴ this paper also proposed a simple PixelCNN, which was revised later in their NIPS paper (next page)

⁵ van den Oord, A., Kalchbrenner, N., and Kavukcuoglu, K. (2016c). [Pixel recurrent neural networks](#). *arXiv preprint arXiv:1601.06759*

PixelCNN

- still generate image pixels starting from corner
- what's new
 - ▶ use CNN (not LSTM) to model dependency on previous pixels
 - ▶ two **conv stacks** to remove _____ spots
- training
 - ▶ maximize likelihood of training images
 - ▶ loss: softmax loss at each pixel
 - ▶ faster⁶ than PixelRNN
- image generation
 - ▶ must still proceed sequentially
 - ⇒ still slow

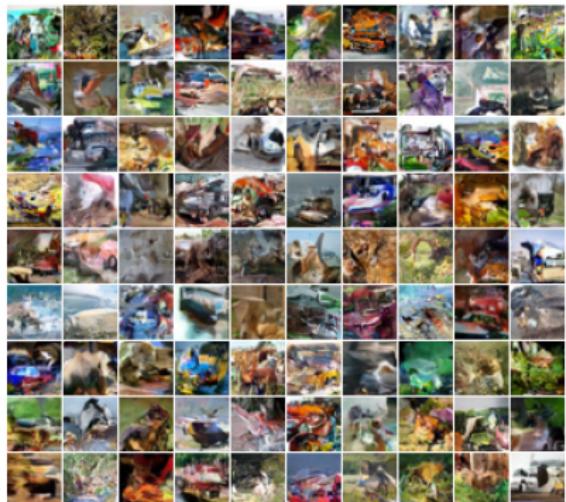
(source: [van den Oord et al., 2016b]⁷)



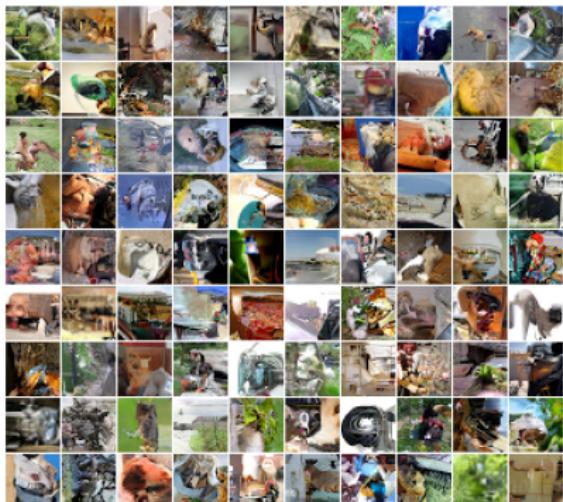
⁶ we can parallelize convolutions since context region values are known from training images

⁷ van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al. (2016b). [Conditional image generation with pixelcnn decoders](#). In *Advances in Neural Information Processing Systems*, pages 4790–4798

Examples of generated samples



32×32 CIFAR-10



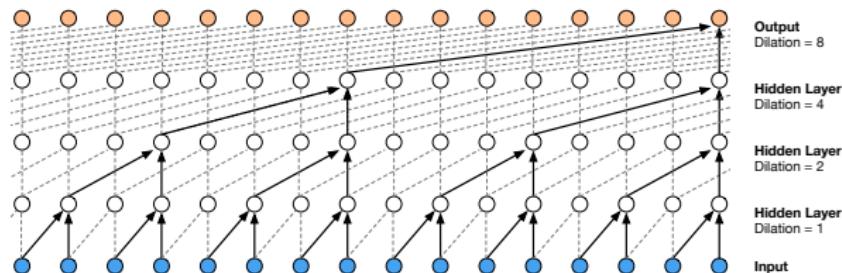
32×32 ImageNet

(source: [van den Oord et al., 2016c]⁸)

⁸ van den Oord, A., Kalchbrenner, N., and Kavukcuoglu, K. (2016c). [Pixel recurrent neural networks](#). *arXiv preprint arXiv:1601.06759*

WaveNet

- can generate raw audios (trained on raw audio waveforms)
 - ▶ similar to PixelCNN in structure, but much more successful
 - ▶ amazing quality but **slow generation** (2 min to synthesize 1 sec audio)
- architecture
 - ▶ use **dilated convolution** to vastly _____ receptive field



(source: [van den Oord et al., 2016a]⁹)

- ▶ stack many layers like the above together
- ▶ use classification (not regression) to generate next sample point
- ▶ convert continuous audio → 256-level quantized classes

⁹ van den Oord, A., Dieleman, S., et al. (2016a). *Wavenet: A generative model for raw audio*. In *SSW*, page 125

Remarks

- autoregressive models: pros and cons
 - ▶ can explicitly compute $p(x)$
 - ▶ good samples
 - ▶ sequential generation: _____
- to improve performance
 - ▶ gated convolution layers
 - ▶ short-cut connections
 - ▶ discretized logistic loss
 - ▶ multi-scale
 - ▶ parallelization
 - ▶ and many more!
- more information: [▶ Link](#)

(source: cs231n, [Salimans et al., 2017]¹⁰)

PixelCNN++ samples from CIFAR-10



real CIFAR-10 images



¹⁰ Salimans, T., Karpathy, A., Chen, X., and Kingma, D. P. (2017). [Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications](#). *arXiv preprint arXiv:1701.05517*

Outline

Generative Models

Approximate Inference

Summary

Probabilistic machine learning

- a probabilistic model
 - ▶ a joint distribution of hidden variables \mathbf{z} and observed variables \mathbf{x}

$$p(\mathbf{z}, \mathbf{x})$$

- ▶ describes how (a portion of) the world works
- inference about the unknowns:
 - ▶ through **posterior** _____
i.e. conditional distribution of the hidden variables given the observations

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})} \leftarrow \begin{array}{l} \text{model} \\ \text{data} \end{array}$$

- ▶ the posterior links data and model \Rightarrow used in all downstream analyses
- posterior inference:
 - ▶ therefore a central task in probabilistic models

Posterior inference

- refers to
 - (1) computing posterior distribution $p(z | x)$ or
 - (2) taking expectations computed wrt this distribution

e.g. expectation-maximization (EM) algorithm

- ▶ evaluates expectation of **complete-data log likelihood**
wrt **posterior distribution of latent variables**:

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) \triangleq \sum_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}^{\text{old}}) \log p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$$

- for many models of practical interest
 - ▶ it is intractable to do (1)/(2) \Rightarrow called “challenge of inference”
- challenge of inference
 - ▶ makes it difficult to ____ probabilistic models

Challenge of inference

- general reasons
 - ▶ dimensionality of latent space: too high to work with directly
 - ▶ posterior: highly complex \Rightarrow expectations are not analytically tractable
- reasons in deep learning
 - ▶ interactions between _____ variables (*e.g.* connections between layers)
 - ▶ most neural nets with multiple layers of hidden variables
 - ▷ have intractable posterior distributions
- a solution: approximate posterior inference

Bayesian view

- we setup the general problem
 - ▶ consider a joint density of latent variables \mathbf{z} and observations \mathbf{x}
$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$$
- latent variables help govern distribution of data in Bayesian models
- a Bayesian model
 - ▶ draws _____ variables from a prior $p(\mathbf{z})$, and then
 - ▶ relates them to observations through likelihood $p(\mathbf{x} | \mathbf{z})$
- inference in a Bayesian model
 - ▶ conditions on data and computes posterior $p(\mathbf{z} | \mathbf{x})$

↑
this often requires approximate inference

Approximate posterior inference

1. stochastic (e.g. Markov chain Monte Carlo: MCMC)

- ▶ given infinite computational resource, can generate exact results
- ▶ approximation arises from using a finite amount of processor time
- ⚠ sampling methods can be computationally demanding
- ⇒ their use is often limited to small-scale problems
- ⚠ difficult to know whether _____ samples are being generated

2. deterministic (e.g. variational inference: VI)

- ▶ some of which scale well to large applications
 - ▶ based on analytical approximations to posterior $p(z | x)$
- e.g. assume it factorizes or has a _____ form (like Gaussian)
- ⇒ never generate exact results

Outline

Generative Models

Approximate Inference

Summary

Summary

- generative models: density estimation and sample generation
 - ▶ explicit density: PixelRNN/CNN, variational autoencoder (VAE)
 - ▶ implicit density: generative adversarial network (GAN)
- generative models are versatile and useful for many tasks
 - ▶ representation/manipulation of high-dim distributions
 - ▶ reinforcement learning, semi-supervised learning
 - ▶ multi-modal outputs, inference of latent representations
- approximate inference schemes fall into two classes: stochastic or deterministic
 - ▶ their strengths and weaknesses are complementary to each other