CSC413/2516 Lecture 10: Generative Models & Reinforcement Learning

Jimmy Ba

Jimmy Ba

CSC413/2516 Lecture 10: Generative Model

Four modern approaches to generative modeling:

- Autoregressive models (Lectures 3, 7, and 8)
- Generative adversarial networks (last lecture)
- Reversible architectures (this lecture)
- Variational autoencoders (this lecture)

All four approaches have different pros and cons.

Autoencoders

- An autoencoder is a feed-forward neural net whose job it is to take an input **x** and predict **x**.
- To make this non-trivial, we need to add a bottleneck layer whose dimension is much smaller than the input.



3

Autoencoders

Why autoencoders?

- Map high-dimensional data to two dimensions for visualization
- Compression (i.e. reducing the file size)
 - Note: this requires a VAE, not just an ordinary autoencoder.
- Learn abstract features in an unsupervised way so you can apply them to a supervised task
 - Unlabled data can be much more plentiful than labeled data
- Learn a semantically meaningful representation where you can, e.g., interpolate between different images.

- 4 目 ト - 4 日 ト - 4 日 ト

Principal Component Analysis (optional)

 The simplest kind of autoencoder has one hidden layer, linear activations, and squared error loss.

$$\mathcal{L}(\mathbf{x}, \tilde{\mathbf{x}}) = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2$$

- This network computes $\tilde{\mathbf{x}} = \mathbf{UVx}$, which is a linear function.
- If K ≥ D, we can choose U and V such that UV is the identity. This isn't very interesting.
 - But suppose K < D:
 - **V** maps **x** to a *K*-dimensional space, so it's doing dimensionality reduction.
 - The output must lie in a *K*-dimensional subspace, namely the column space of **U**.



(人間) システン イラン

Principal Component Analysis (optional)

- Review from CSC311: linear autoencoders with squared error loss are equivalent to Principal Component Analysis (PCA).
- Two equivalent formulations:
 - Find the subspace that minimizes the reconstruction error.
 - Find the subspace that maximizes the projected variance.
- The optimal subspace is spanned by the dominant eigenvectors of the empirical covariance matrix.



"Eigenfaces"

Deep Autoencoders

- Deep nonlinear autoencoders learn to project the data, not onto a subspace, but onto a nonlinear manifold
- This manifold is the image of the decoder.
- This is a kind of nonlinear dimensionality reduction.



Deep Autoencoders

 Nonlinear autoencoders can learn more powerful codes for a given dimensionality, compared with linear autoencoders (PCA)



くほと くほと くほと

Deep Autoencoders

- Some limitations of autoencoders
 - They're not generative models, so they don't define a distribution
 - How to choose the latent dimension?

イロト 不得下 イヨト イヨト

Observation Model

• Consider training a generator network with maximum likelihood.

$$p(\mathbf{x}) = \int p(\mathbf{z}) p(\mathbf{x} \,|\, \mathbf{z}) \, \mathrm{d}\mathbf{z}$$

- One problem: if z is low-dimensional and the decoder is deterministic, then p(x) = 0 almost everywhere!
 - The model only generates samples over a low-dimensional sub-manifold of $\ensuremath{\mathcal{X}}.$
- Solution: define a noisy observation model, e.g.

$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x}; G_{\boldsymbol{\theta}}(\mathbf{z}), \eta \mathbf{I}),$$

where G_{θ} is the function computed by the decoder with parameters θ .



Observation Model

- At least $p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) d\mathbf{z}$ is well-defined, but how can we compute it?
- Integration, according to XKCD:



11 / 28

Jimmy Ba

Observation Model

- At least $p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) d\mathbf{z}$ is well-defined, but how can we compute it?
 - The decoder function $G_{\theta}(\mathbf{z})$ is very complicated, so there's no hope of finding a closed form.
- Instead, we will try to maximize a lower bound on $\log p(\mathbf{x})$.
 - The math is essentially the same as in the EM algorithm from CSC411.

(本語)と (本語)と (本語)と

Variational Inference

• We obtain the lower bound using Jensen's Inequality: for a convex function *h* of a random variable *X*,

 $\mathbb{E}[h(X)] \ge h(\mathbb{E}[X])$

Therefore, if h is concave (i.e. -h is convex),

 $\mathbb{E}[h(X)] \leq h(\mathbb{E}[X])$

• The function log *z* is concave. Therefore,

 $\mathbb{E}[\log X] \le \log \mathbb{E}[X]$





< 3 > < 3 >

Variational Inference

- Suppose we have some distribution q(z). (We'll see later where this comes from.)
- We use Jensen's Inequality to obtain the lower bound.

$$\begin{split} \log p(\mathbf{x}) &= \log \int p(\mathbf{z}) \, p(\mathbf{x}|\mathbf{z}) \, \mathrm{d}\mathbf{z} \\ &= \log \int q(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} p(\mathbf{x}|\mathbf{z}) \, \mathrm{d}\mathbf{z} \\ &\geq \int q(\mathbf{z}) \log \left[\frac{p(\mathbf{z})}{q(\mathbf{z})} \, p(\mathbf{x}|\mathbf{z}) \right] \, \mathrm{d}\mathbf{z} \qquad \text{(Jensen's Inequality)} \\ &= \mathbb{E}_q \left[\log \frac{p(\mathbf{z})}{q(\mathbf{z})} \right] + \mathbb{E}_q \left[\log p(\mathbf{x}|\mathbf{z}) \right] \end{split}$$

• We'll look at these two terms in turn.

- The first term we'll look at is $\mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})]$
- Since we assumed a Gaussian observation model,

$$\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; G_{\theta}(\mathbf{z}), \eta \mathbf{I})$$

=
$$\log \left[\frac{1}{(2\pi\eta)^{D/2}} \exp \left(-\frac{1}{2\eta} \|\mathbf{x} - G_{\theta}(\mathbf{z})\|^2 \right) \right]$$

=
$$-\frac{1}{2\eta} \|\mathbf{x} - G_{\theta}(\mathbf{z})\|^2 + \text{const}$$

• So this term is the expected squared error in reconstructing **x** from **z**. We call it the reconstruction term.

Variational Inference

- The second term is $\mathbb{E}_q\left[\log \frac{p(z)}{q(z)}\right]$.
- This is just $-D_{KL}(q(z)||p(z))$, where D_{KL} is the Kullback-Leibler (KL) divergence

$$\mathbb{D}_{\mathrm{KL}}(q(\mathsf{z}) \| p(\mathsf{z})) riangleq \mathbb{E}_q \left[\log rac{q(\mathsf{z})}{p(\mathsf{z})}
ight]$$

- KL divergence is a widely used measure of distance between probability distributions, though it doesn't satisfy the axioms to be a distance metric.
- More details in tutorial.
- Typically, $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. Hence, the KL term encourages q to be close to $\mathcal{N}(\mathbf{0}, \mathbf{I})$.
- We'll give the KL term a much more interesting interpretation when we discuss Bayesian neural nets.

Variational Inference

• Hence, we're trying to maximize the variational lower bound, or variational free energy:

$$\log p(\mathbf{x}) \geq \mathcal{F}(\boldsymbol{\theta}, q) = \mathbb{E}_q \left[\log p(\mathbf{x}|\mathbf{z})\right] - \mathrm{D}_{\mathrm{KL}}(q\|p).$$

- The term "variational" is a historical accident: "variational inference" used to be done using variational calculus, but this isn't how we train VAEs.
- We'd like to choose q to make the bound as tight as possible.
- It's possible to show that the gap is given by:

$$\log p(\mathbf{x}) - \mathcal{F}(\boldsymbol{\theta}, q) = \mathrm{D}_{\mathrm{KL}}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x})).$$

Therefore, we'd like q to be as close as possible to the posterior distribution $p(\mathbf{z}|\mathbf{x})$.

• • = • • = •

- Let's think about the role of each of the two terms.
- The reconstruction term

$$\mathbb{E}_{q}[\log p(\mathbf{x}|\mathbf{z})] = -\frac{1}{2\sigma^{2}}\mathbb{E}_{q}[\|\mathbf{x} - G_{\theta}(\mathbf{z})\|^{2}] + \text{const}$$

is minimized when q is a point mass on

$$\mathbf{z}_* = \arg\min_{\mathbf{z}} \|\mathbf{x} - G_{\boldsymbol{\theta}}(\mathbf{z})\|^2.$$

• But a point mass would have infinite KL divergence. (Exercise: check this.) So the KL term forces q to be more spread out.

Reparameterization Trick

- To fit q, let's assign it a parametric form, in particular a Gaussian distribution: $q(z) = \mathcal{N}(z; \mu, \Sigma)$, where $\mu = (\mu_1, \dots, \mu_K)$ and $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_K^2)$.
- In general, it's hard to differentiate through an expectation. But for Gaussian *q*, we can apply the reparameterization trick:

$$\mathbf{z}_i = \mu_i + \sigma_i \epsilon_i,$$

where $\epsilon_i \sim \mathcal{N}(0, 1)$.

Hence,

$$\overline{\mu_i} = \overline{z_i} \qquad \overline{\sigma_i} = \overline{z_i} \epsilon_i.$$

 This is exactly analogous to how we derived the backprop rules for droopout.

イロト 不得下 イヨト イヨト 三日

Amortization

- This suggests one strategy for learning the decoder. For each training example,
 - Fit q to approximate the posterior for the current x by doing many steps of gradient ascent on F.
 - **2** Update the decoder parameters θ with gradient ascent on \mathcal{F} .
- **Problem:** this requires an expensive iterative procedure for every training example, so it will take a long time to process the whole training set.

- 4 週 ト - 4 三 ト - 4 三 ト

Amortization

- Idea: amortize the cost of inference by learning an inference network which predicts (μ, Σ) as a function of x.
- The outputs of the inference net are μ and log σ. (The log representation ensures σ > 0.)
- If $\sigma \approx 0$, then this network essentially computes z deterministically, by way of μ .
 - But the KL term encourages σ > 0, so in general z will be noisy.
- The notation q(z|x) emphasizes that q depends on x, even though it's not actually a conditional distribution.



< 回 > < 三 > < 三 >

Amortization

- Combining this with the decoder network, we see the structure closely resembles an ordinary autoencoder. The inference net is like an encoder.
- Hence, this architecture is known as a variational autoencoder (VAE).
- The parameters of both the encoder and decoder networks are updated using a single pass of ordinary backprop.
 - The reconstruction term corresponds to squared error $\|\mathbf{x} \tilde{\mathbf{x}}\|^2$, like in an ordinary VAE.
 - The KL term regularizes the representation by encouraging **z** to be more stochastic.



• • = • • = •

VAEs vs. Other Generative Models

- In short, a VAE is like an autoencoder, except that it's also a generative model (defines a distribution $p(\mathbf{x})$).
- Unlike autoregressive models, generation only requires one forward pass.
- Unlike reversible models, we can fit a low-dimensional latent representation. We'll see we can do interesting things with this...



Class-Conditional VAE

- So far, we haven't used the labels *y*. A class-conditional VAE provides the labels to both the encoder and the decoder.
- Since the latent code z no longer has to model the image category, it can focus on modeling the stylistic features.
- If we're lucky, this lets us disentangle style and content. (Note: disentanglement is still a dark art.)
- See Kingma et al., "Semi-supervised learning with deep generative models."



A B A A B A

Class-Conditional VAE

By varying two latent dimensions (i.e. dimensions of z) while holding y fixed, we can visualize the latent space.

З 2 7 2 7 2 2 2 3 3 2 2 7

イロト 不得下 イヨト イヨト

Class-Conditional VAE

• By varying the label y while holding z fixed, we can solve image analogies.



ヘロト 不得下 不足下 不足下

Latent Space Interpolations

 You can often get interesting results by interpolating between two vectors in the latent space:



Ha and Eck, "A neural representation of sketch drawings"

< ロ > < 同 > < 回 > < 回 > < 回 > <

Latent Space Interpolations

• Latent space interpolation of music: https://magenta.tensorflow.org/music-vae

- 4 週 ト - 4 三 ト - 4 三 ト