Machine Learning Lecture 10: Variational Autoencoder

Nevin L. Zhang lzhang@cse.ust.hk

Department of Computer Science and Engineering The Hong Kong University of Science and Technology

This set of notes is based on internet resources and Auto-encoding variational bayes DP Kingma, M Welling (2013). Auto-encoding variational bayes. https://arxiv.org/abs/1312.6114 C Doersch (2016). Tutorial on Variational Autoencoders. https://arxiv.org/abs/1606.05908

Outline

1 Introduction to Unsupervised Learning

- 2 The Task
- 3 The Objective function
- 4 Optimization
- 5 Generating Examples
- 6 Discussions

Introduction

- So far, supervised learning
 - Discriminative methods:

$$\{\mathbf{x}_i, y_i\}_{i=1}^N \to p(y|\mathbf{x})$$

Generative methods:

$$\{\mathbf{x}_i, y_i\}_{i=1}^N \to P(y), p(\mathbf{x}|y)$$

- Next, unsupervised learning:
 - Finite mixture models for clustering [Skipped]

 $\{\mathbf{x}_i\}_{i=1}^N \to P(z), p(\mathbf{x}|z)$

- Varitional autoencoder for data generation and representation learning $\{\mathbf{x}_i\}_{i=1}^N, p(\mathbf{z}) \rightarrow p(\mathbf{x}|\mathbf{z}) \quad q(\mathbf{z}|\mathbf{x})$ used in inference
- Generative adversarial networks for data generation

 $\{\mathbf{x}_i\}_{i=1}^N, p(\mathbf{z}) \rightarrow \mathbf{x} = g(\mathbf{z})$

Outline

- 1 Introduction to Unsupervised Learning
- 2 The Task
- 3 The Objective function
- 4 Optimization
- 5 Generating Examples
- 6 Discussions

The Task



- Suppose we have an **unlabeled dataset X** = {**x**⁽ⁱ⁾}^N_{i=1}, where each training example **x**⁽ⁱ⁾ is a vector that represents an image and each component of **x**⁽ⁱ⁾ represents a pixel in the image.
- We would like to learn a distribution $p(\mathbf{x})$ from the dataset so that we can generate more images that are similar (but different) to those in the dataset.
- If we can solve this task, then we have the ability to learn very complex probabilistic model for high dimensional data.
- The ability to generate realistic looking images would be useful for video game designers.

The Generative Model

- We assume that each image has a label **z** that is not observed. **z** is a vector of much lower dimension that **x**.
- We further assume that the images are generated as follows:
 - $z ~\sim~ p(z) = \mathcal{N}(0,I)$ where I is the identity matrix
 - ${\bf x} ~\sim~ \rho_{\theta}({\bf x}|{\bf z})$ where θ denotes model parameters
- Then we have

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



The Generative Model

In addition, we assume that the conditional distribution is a Gaussian

$$p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mu_{x}(\mathbf{z},\boldsymbol{\theta}),\sigma_{x}^{2}(\mathbf{z},\boldsymbol{\theta})\mathbf{I})$$

- With mean vector is $\mu_x(\mathbf{z}, \boldsymbol{\theta})$ and diagonal covariance matrix is $\sigma_x(\mathbf{z}, \boldsymbol{\theta})\mathbf{I}$.
- The mean vector μ_x(z, θ) and the vector σ_x(z, θ) of sd's are deterministically determined by z via a deep neural network with parameters θ.
- So, we make use of the ability of neural network in representing complex functions to learn complicated probabilistic models.



Outline

- 1 Introduction to Unsupervised Learning
- 2 The Task
- 3 The Objective function
- 4 Optimization
- 5 Generating Examples
- 6 Discussions

The Likelihood Function

• To learn the model parameters, we need maximize the following likelihood function:

$$\log p_{\theta}(\mathbf{X}) = \sum_{i=1}^{N} \log p_{\theta}(\mathbf{x}^{(i)})$$

where

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

- We want to use gradient ascent to maximize the likelihood function, which requires the gradient $\nabla_{\theta} \log p_{\theta}(\mathbf{x}^{(i)})$.
- The gradient is intractable because of the integration.

Naive Monte Carlo Gradient Estimator

- Here is a naive method to estimate $p_{\theta}(\mathbf{x}^{(i)})$ and hence the gradients.
- Sample L points $\mathbf{z}^{(1)}, \ldots, \mathbf{z}^{(L)}$ from $p(\mathbf{z})$, and estimate $p_{\theta}(\mathbf{x}^{(i)})$ using

$$p_{\theta}(\mathbf{x}^{(i)}) \approx \frac{1}{L} \sum_{l=1}^{L} p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(l)})$$

- Then we can compute $\nabla_{\theta} \log p_{\theta}(\mathbf{x}^{(i)})$.
- Unfortunately, this would not work.
 - The reason is that \mathbf{x} is high-dimensional (thousands to millions of dimensions). Given \mathbf{z} , $p_{\theta}(\mathbf{x}|\mathbf{z})$ is highly skewed, taking non-negligible values only in a very small region.
 - To state it another way, for a given data point $x^{(i)}$, $p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})$ takes non-negligible values only for \mathbf{z} from a very small region. As such, L needs to be extremely large for the estimate to be accurate.

Recognition Model



■ To overcome the aforementioned difficulty, we introduce a recognition model q_φ(z|x)

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_{z}(\mathbf{x},\phi),\sigma_{z}^{2}(\mathbf{x},\phi)\mathbf{I})$$

- The mean vector μ_z(x, φ) and the vector σ_z(x, φ) of sd's are deterministically determined by z via a deep neural network with parameters φ.
- We hope to get from $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$ samples of \mathbf{z} for which $p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})$ has non-negligible values.
- The question now is: How to make use of $q_{\phi}(\mathbf{z}|\mathbf{x})$ when maximize the likelihood log $p_{\theta}(\mathbf{X}) = \sum_{i=1}^{N} \log p_{\theta}(\mathbf{x}^{(i)})$.

The Variational Lower Bound

$$\log p_{\theta}(\mathbf{x}^{(i)}) = E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\theta}(\mathbf{x}^{(i)}) \right]$$

$$= E_{\mathbf{z} \sim q_{\phi}} \left[\log \frac{p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})} \right]$$

$$= E_{\mathbf{z} \sim q_{\phi}} \left[\log \frac{p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})} \frac{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})}{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \right]$$

$$= E_{\mathbf{z} \sim q_{\phi}} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - E_{\mathbf{z} \sim q_{\phi}} \left[\frac{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})}{p_{\theta}(\mathbf{z})} \right] + E_{\mathbf{z} \sim q_{\phi}} \left[\frac{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})} \right]$$

$$= E_{\mathbf{z} \sim q_{\phi}} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - \mathcal{D}_{KL} [q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})] + E_{\mathbf{z} \sim q_{\phi}} \left[\frac{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})} \right]$$

$$= \mathcal{L}(\mathbf{x}^{(i)}, \theta, \phi) + \mathcal{D}_{KL} [q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})]$$

So, we have the following variational lower bound on loglikelihood, which is tight if q has high capacity.

$$\log p_{oldsymbol{ heta}}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\mathbf{x}^{(i)},oldsymbol{ heta},\phi)$$

Nevin L. Zhang (HKUST)

The Variational Lower Bound: Alternative Perspective

 \blacksquare Our new objective is to maximize the variational bound w.r.t both θ and ϕ

$$\mathcal{L}(\mathbf{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi}) = E_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}^{(i)})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \mathcal{D}_{KL}[q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}^{(i)}) || p_{\boldsymbol{\theta}}(\mathbf{z})]$$

- Interpretation
 - The recognition model $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$ can be viewed as a **encoder** that takes a data point $\mathbf{x}^{(i)}$ and probabilistically encodes it into a latent vector \mathbf{z} .
 - The decoder p_θ(x|z) then takes the latent representation and probabilistically decodes it into a vector x in the data space.
 - The first term in *L* measure how well (the distribution of) the decoded output match the input **x**^(*i*). It is the **reconstruction error**.
 - The second term is a regularization terms that encourages the posterior distribution $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$ of the encoding \mathbf{z} to be close to the prior $p_{\theta}(\mathbf{z})$.
 - So, the method is called variational autoencoder (VAE).

Illustration of Variational Autoencoder

$$\mathcal{L}(\mathbf{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi}) = E_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}^{(i)})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right] - \mathcal{D}_{KL}[q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}^{(i)}) || p_{\boldsymbol{\theta}}(\mathbf{z})]$$



Illustration of Variational Autoencoder



- The encoder maps the data distribution, which is complex, to approximately an Gaussian distribution.
- The decoder maps a Gaussian distribution to the data distribution.

Illustration of Variational Autoencoder



• Fake images generated by picking points in the latent space and map them back to the data space using the decoder.

Outline

- 1 Introduction to Unsupervised Learning
- 2 The Task
- 3 The Objective function

4 Optimization

- 5 Generating Examples
- 6 Discussions

Optimization

The Need For Reparameterization



• The computation of the first term \mathcal{L}_1 of \mathcal{L} requires **sampling**:

$$\mathcal{L}_{1} = E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] \approx \frac{1}{L} \sum_{l=1}^{L} \log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$$
where $\mathbf{z}^{(i,l)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$.

- But sampling looses gradient ∇_{ϕ}
 - While the LHS depends on ϕ , the RHS does not.
 - So, the stochastic connections from μ_z and σ_z to z makes backpropagation impossible.

Optimization

The Reparameterization Trick

■ Here is the **recognition model**

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_{z}(\mathbf{x},\phi),\sigma_{z}^{2}(\mathbf{x},\phi)\mathbf{I})$$

■ Using the reparameterization trick, we change it into the following equivalent form

$$\mathbf{z} = \mu_z(\mathbf{x}, \phi) + \sigma_z(\mathbf{x}, \phi) \odot \boldsymbol{\epsilon}, \ \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

where \odot is element-wise product.

• Note that, now, z depends on μ_z , σ_z and ϵ deterministically. ϵ is stochastic, but it is an input the network.



Optimization

The Reparameterization Trick

The reconstruction error term can be written as

$$\mathcal{L}_{1} = E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right]$$

$$= E_{\epsilon \sim p(\epsilon)} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}(\mathbf{x}^{(i)},\phi,\epsilon)) \right] \approx \frac{1}{L} \sum_{l=1}^{L} \log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$$
where $\mathbf{z}^{(i,l)} = \mathbf{z}(\mathbf{x}^{(i)},\phi,\epsilon^{(l)}) = \mu_{z}(\mathbf{x}^{(i)},\phi) + \sigma_{z}(\mathbf{x}^{(i)},\phi) \odot \epsilon^{(l)}$, and
$$\epsilon^{(l)} \sim \mathcal{N}(\mathbf{0},\mathbf{I}).$$

Gradient $\nabla_{\theta,\phi} \mathcal{L}_1$ can now be computed because, for each given ϵ , the network is deterministic.



١

The Regularization Term

- The second term of \mathcal{L} is $\mathcal{L}_2 = -\mathcal{D}_{KL}[q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})].$
- The two distributions involved are both Gaussian. Hence the term has a closed-form:

$$\mathcal{L}_{2} = \frac{1}{2} \sum_{j=1}^{J} \left(1 + \log((\sigma_{j}^{(i)})^{2}) - (\mu_{j}^{(i)})^{2} - (\sigma_{j}^{(i)})^{2} \right)$$

where J is the dimension of z, $\sigma_j^{(i)}$ is the *j*-th component of $\sigma_z(\mathbf{x}^{(i)}, \phi)$, and $\mu_j^{(i)}$ is the *j*-th component of $\mu_z(\mathbf{x}^{(i)}, \phi)$

• The gradient $\nabla_{\phi} \mathcal{L}_2$ is straightforward compute

The Final Objective Function

 Putting together, this is the objective function that we maximize using gradient ascent

$$\mathcal{L} \approx \frac{1}{L} \sum_{l=1}^{L} \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)}) + \frac{1}{2} \sum_{j=1}^{J} \left(1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right)$$

where $\mathbf{z}^{(i,l)} = \mu_z(\mathbf{x}^{(i)}, \phi) + \sigma_z^2(\mathbf{x}^{(i)}, \phi) \odot \boldsymbol{\epsilon}^{(l)}$, and $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

- We have discussed how to compute the gradient $\nabla_{\theta,\phi}\mathcal{L}$. Using the gradient, we can estimate θ and ϕ simultaneously using gradient ascent.
- The hyperparameter L is usually set to 1.

Comparison with Naive Monte Carlo

 Earlier, we mentioned a naive method for optimizing the parameters θ of the generative model that involves the following objective:

$$\log p_{\theta}(\mathbf{x}^{(i)}) \approx \log \frac{1}{L} \sum_{l=1}^{L} p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(l)})$$
(1)

where the values of **z** are sampled from the prior $p(\mathbf{z})$. Those values do not depend on $\mathbf{x}^{(i)}$ and do not give high probabilities to $\mathbf{x}^{(i)}$. The RHS is a poor approximation of log $p_{\theta}(\mathbf{x}^{(i)})$.

• Here is the our final objective function:

$$\mathcal{L} \approx \frac{1}{L} \sum_{l=1}^{L} \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)}) + \frac{1}{2} \sum_{j=1}^{J} \left(1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right)$$
(2)

where the values of \mathbf{z} are sampled in such a way that they depend upon $\mathbf{x}^{(i)}$. Those values usually give high probabilities to $\mathbf{x}^{(i)}$. The LHS is a better approximation (lower bound) of log $p_{\theta}(\mathbf{x}^{(i)})$.

Outline

- 1 Introduction to Unsupervised Learning
- 2 The Task
- 3 The Objective function
- 4 Optimization
- 5 Generating Examples
- 6 Discussions

Example Generation

During learning, both encoder and decoder are trained simultaneously.



 \blacksquare To generate examples, we need on the decoder.

z ~
$$p(z)$$
, **x** ~ $p_{\theta}(\mathbf{x}|\mathbf{z})$



Example Generation

- One way to generate examples is to sample z from N(0, I) and then sample x from p_θ(x|z).
- Here are images sampled from VAE's learned from the MNIST dataset. Note that several values are used for the dimensionality of z.

8 9 6 6 3 1 9 (a) 2-D latent space (b) 5-D latent space (c) 10-D latent space (d) 20-D latent space

Example Generation

- Alternatively, we can manually pick z and sample x from p_θ(x|z). This allows us to interpret each dimension of z.
- Here are images sampled from VAE's learned from the Frey Face dataset. We can see that
 - X-axis represents head pose, and
 - Y-axis represent degree of smile.



Outline

- 1 Introduction to Unsupervised Learning
- 2 The Task
- 3 The Objective function
- 4 Optimization
- 5 Generating Examples

6 Discussions

- The key reason for the excitement in VAE is that it shows that, by combining deep learning and the probabilistic approach, we can now learn complicated and high quality probabilistic models.
- In terms of specific functionality,
 - The decoder *p*(**x**|**z**) of VAE can be used to generate samples (images) that are similar (but different) to those in the training set.
 - The decoder gives a distribution $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$, which can be approximated using the variational lower bound.
 - The encoder *q*(**z**|**x**) can be used to obtain low dimension latent representation of data.

Autoencoders



While variational autoencoder is a probabilistic model, autoencoder is deterministic. The model parameters are trained by minimizing the reconstruction errors

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

- It is designed to learn a latent representation of data.
- However, it does not define a probability distribution p(x) over data space and does not to generate new samples.
- Recent work on AE: He, Kaiming, et al. "Masked autoencoders are scalable vision learners." arXiv preprint arXiv:2111.06377 (2021).

Denoising Autoencoders



Denoising autoencoder is a probabilistic model. The input x is randomly corrupted using C(x̃|x) to get x̃ and the weights are optimized to minimize the following objective function:

$$-E_{\mathbf{x} \sim p_{data}} E_{\tilde{\mathbf{x}} \sim C(\tilde{\mathbf{x}}|x)} \log p(\mathbf{x}' = \mathbf{x} | \mathbf{z}(\tilde{\mathbf{x}}))$$

- It is more robust than autoencoder as a tool for learning latent representation of data.
- However, it does not define a probability distribution p(x) over data space and does not to generate new samples.

Data Distribution in Latent Space (MNIST)



- VAE: Forces data into a normal distribution in the latent space.
- DAE: Preserves class separation better.

Recent Developments: Flow-Based Generative Models

- Dinh L, Sohl-Dickstein J, Bengio S. Density estimation using Real NVP[J]. arXiv preprint arXiv:1605.08803, 2016.
- Kingma D P, Dhariwal P. Glow: Generative flow with invertible 1x1 convolutions[C]//Advances in Neural Information Processing Systems. 2018: 10236-10245.
- Generative Model:
 - **z** $\sim p_z(z)$ (Gaussian), $\mathbf{x} = g_\theta(z)$,
 - $g_{\theta}(\mathbf{z})$ is invertible, i.e., exists $f_{\theta}(\mathbf{x})$ such that $g_{\theta}(f_{\theta}(\mathbf{x})) = \mathbf{x}$.
 - Consequently, z has the same dimensionality as x.
 - The model defines a distribution over inputs:

$$p_{ heta}(\mathbf{x}) = p_z(f_{ heta}(\mathbf{x})) |det(rac{\partial f_{ heta}(\mathbf{x})}{\partial \mathbf{x}^ op})|$$

• f_{θ} is implemented as a sequence of invertible functions, called **flows**, each represented as an CNN. See references.

Recent Developments: Flow-Based Generative Models

Objective function for learning: maximizing the loglikelihood

$$\log p_{\theta}(\mathbf{X}) = \sum_{i=1}^{N} \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^{N} \log p_{z}(f_{\theta}(\mathbf{x}^{i})) + \sum_{i=1}^{N} \log |det(\frac{\partial f_{\theta}(\mathbf{x}^{(i)})}{\partial (\mathbf{x}^{(i)})^{\top}})|$$

 Intuitively, pick θ so that the CNN map images from their original space, where they are not Gaussian distributed, to a latent space where they are Gaussian distributed.



Recent Developments: Flow-Based Generative Models

■ Image Synthesis: $\mathbf{z} \sim p_z(\mathbf{z})$ (Gaussian), $\mathbf{x} = g_\theta(\mathbf{z})$



Recent Developments: Flow-Based Generative Models

Image Interpolation: Take a pair of real images, encode them with the encoder, and linearly interpolate between the latents to obtain samples.



Recent Developments: Flow-Based Generative Models

- Semantic Manipulation:
 - **z**_{pos}: average latent vectors with an attribute (e.g., smiling)
 - **z**_{neg}: average latent vectors without the attribute.
 - Use the difference $\mathbf{z}_{pos} \mathbf{z}_{neg}$ as a direction of manipulation.



(a) Smiling



(c) Blond Hair