

Variational Autoencoders

1

1

References

These notes are based on the following:

- Kingma, D. and Welling, M. (2014). Auto-
Encoding Variational Bayes. In ICLR 2014.
<https://arxiv.org/pdf/1312.6114.pdf>
- Carl Doersch's tutorial on Variational Auto-
encoders. <https://arxiv.org/pdf/1606.05908.pdf>
- Stefano Ermon's CS 228 notes:
[https://ermongroup.github.io/cs228-
notes/extras/vae/](https://ermongroup.github.io/cs228-notes/extras/vae/)

2

2

Introduction

- Suppose you had a dataset with N instances $\mathbf{X} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$ where N is very big
- Each instance has d features i.e. $\mathbf{x}^i = \{x_1^i, \dots, x_d^i\}$
- E.g. each data instance \mathbf{x}_i is a $20 \times 28 = 560$ dimensional image of a face as shown below.



Source: Frey Face dataset.
https://cs.nyu.edu/~roweis/data/frey_rawface.mat

3

3

Latent Variable Models

- We would like to learn the joint distribution $P(\mathbf{X})$ that generates this dataset (i.e. learn a **generative model**)
- We also assume that there are K unobserved variables (called **latent factors**)
 $\mathbf{Z} = \{z_1, z_2, \dots, z_K\}$
- The latent variables control variation in the features e.g. smiling/frowning, facing left/center/right, etc.

4

4

Latent Variable Models

We assume the following generative process:

1. Sample $\mathbf{z}^i \sim P_{\theta^*}(\mathbf{z})$ where \mathbf{z} is an unobserved continuous random variable and θ^* are the true parameters. Note: \mathbf{z}^i is k dimensional i.e. $\mathbf{z}^i = \{z_1^i, \dots, z_k^i\}$
2. Sample $x_i \sim P_{\theta^*}(x|\mathbf{z})$. Note: x^i is d dimensional i.e. $x^i = \{x_1^i, \dots, x_d^i\}$.

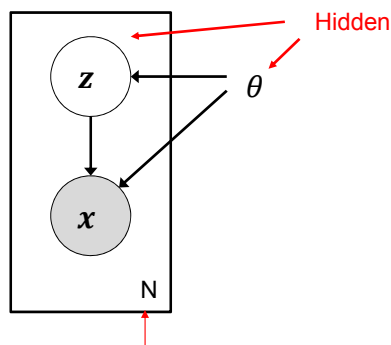
Note: θ^* and \mathbf{z}_i are unknown but we observe x_i

5

5

Latent Variable Models

This leads to the following probabilistic graphical model:



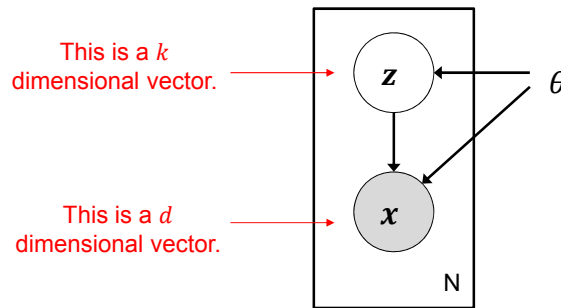
This box is a plate, meaning you replicate the contents by the number in the lower right corner. Each replicate corresponds to a data instance x_i and its corresponding latent factor \mathbf{z}_i .

6

6

Latent Variable Models

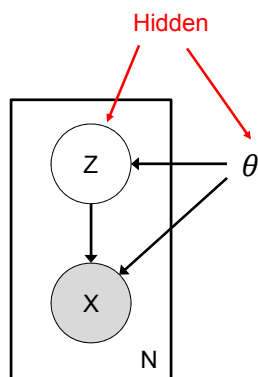
Note: This is multivariate data. The diagram is deliberately vague to abstract away the complex relationships between the dimensions of z and x



7

7

Latent Variable Models



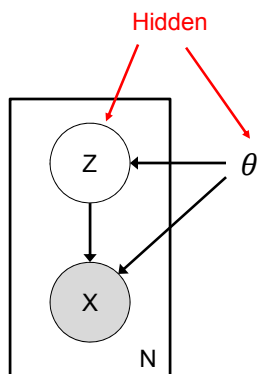
There are 3 tasks of interest:

1. Learn the parameters θ from data
2. Learn the latent factors Z given X
3. Fill in missing values in a new data point x_{new} e.g. for image inpainting

8

8

Latent Variable Models



These tasks involve computation of the following:

$$p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})}{p_{\theta}(\mathbf{x})}$$

But this is intractable:

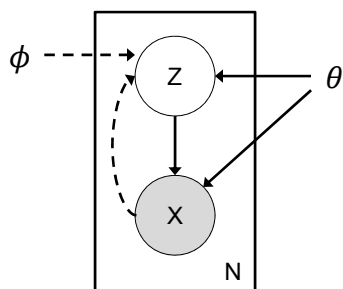
$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})d\mathbf{z}$$

We also assume there is so much data we fit it into memory and have to work with minibatches

9

9

The Variational Bound



Let $q_{\phi}(\mathbf{z}|\mathbf{x})$ be a variational approximation to the intractable posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$

Note that $q_{\phi}(\mathbf{z}|\mathbf{x})$ has parameters ϕ that we need to optimize

10

10

The Variational Bound

Recall from the variational inference section that we can write:

$$\log p_{\theta}(\mathbf{x}) = KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) + ELBO(\phi, \theta, \mathbf{x})$$

where

$$ELBO(\phi, \theta, \mathbf{x}) = E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]$$

We can also rewrite:

$$ELBO(\phi, \theta, \mathbf{x}) = E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$$

11

11

The Variational Bound

$$ELBO(\phi, \theta, \mathbf{x}) = E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$$

Reconstruction error:
Needs to reconstruct \mathbf{x}
from \mathbf{z} such that \mathbf{x} has
high probability under
generative distribution P_{θ}

Regularization term:
Makes \mathbf{z} look like the
prior (Gaussian) to
prevent learning the
identity function.

12

12

The Variational Bound

$$ELBO(\phi, \theta, \mathbf{x}) = E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$$

Think of $p_{\theta}(\mathbf{x}|\mathbf{z})$ as a **probabilistic decoder** that takes a **code \mathbf{z}** and decodes it to an instance \mathbf{x}

Think of $q_{\phi}(\mathbf{z}|\mathbf{x})$ as a **probabilistic encoder** that takes an instance \mathbf{x} and encodes it to a **code \mathbf{z}** . It approximates the true posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$.

This is the connection to autoencoders in deep learning!

13

13

The Variational Bound

$$ELBO(\phi, \theta, \mathbf{x}) = E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$$

We want to differentiate the ELBO above w.r.t:

1. The generative parameters θ (straightforward)
2. The variational parameters ϕ (not straightforward)

14

14

The Reparameterization Trick

Need to compute gradient of an expectation:

$$\begin{aligned} \nabla_{\phi} E_{z \sim q_{\phi}(z)}[f(z)] &= \nabla_{\phi} \int q_{\phi}(z) f(z) dz \\ &= \int \nabla_{\phi} q_{\phi}(z) f(z) dz \\ &= \int \nabla_{\phi} \log(q_{\phi}(z)) q_{\phi}(z) f(z) dz \\ &= E_{z \sim q_{\phi}(z)}[f(z) \nabla_{\phi} \log(q_{\phi}(z))] \\ &\approx \frac{1}{L} \sum_{l=1}^L f(z^l) \nabla_{\phi} \log(q_{\phi}(z^l)) \text{ with } z^l \sim q_{\phi}(z|x^i) \end{aligned}$$

Note:

$$\begin{aligned} \nabla_{\phi} \log(q_{\phi}(z)) q_{\phi}(z) &= \frac{1}{q_{\phi}(z)} \nabla_{\phi} q_{\phi}(z) q_{\phi}(z) \\ &= \nabla_{\phi} q_{\phi}(z) \end{aligned}$$

Estimating integral with L samples

15

15

The Reparameterization Trick

This is called the score function estimator:

$$\nabla_{\phi} E_{z \sim q_{\phi}(z)}[f(z)] = E_{z \sim q_{\phi}(z)}[f(z) \nabla_{\phi} \log(q_{\phi}(z))]$$

Reference: M. C. Fu. Gradient estimation. Handbooks in operations research and management science, 13:575–616, 2006

Suffers from high variance. To see why:

- Doesn't use information about $f(z)$ to guide sampling of $z^l \sim q_{\phi}(z)$
- Could sample z that is in low probability regions of $f(z)$
- Needs lots of samples to get a good estimate. With small # of samples, you get high variance

16

16

The Reparameterization Trick

To get a low variance estimator, we rewrite $z \sim q_\phi(z|x)$ as a two step process:

1. Sample noise term $\varepsilon \sim p(\varepsilon)$ where $p(\varepsilon)$ is a simple distribution
2. Create a deterministic function that combines ε with x i.e. $z = g_\phi(\varepsilon, x)$

17

17

The Reparameterization Trick

Example: suppose $q_\phi(z|x)$ is a normal distribution. Previously we wrote $z \sim N(z; \mu, \sigma)$. Now we write:

- 1) Sample $\varepsilon \sim N(0,1)$
- 2) $z = \mu + \sigma\varepsilon$

Note: This produces the same distribution

18

18

The Reparameterization Trick

More generally, it allows us to do the following:

$$\begin{aligned} & \nabla_{\phi} E_{z \sim q_{\phi}(z|x)} [f(x, z)] \\ &= \nabla_{\phi} E_{\varepsilon \sim p(\varepsilon)} [f(x, g_{\phi}(\varepsilon, x))] \\ &= E_{\varepsilon \sim p(\varepsilon)} [\nabla_{\phi} f(x, g_{\phi}(\varepsilon, x))] \end{aligned}$$

The gradient moves inside the expectation. This estimator has much lower variance than the score function estimator.

See Appendix D of Rezende, D. J., Mohamed, S. and Wierstra, D. (2014). Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In Proceedings of the 31st International Conference on Machine Learning.

19

19

The Reparameterization Trick

$$\widetilde{ELBO}(\phi, \theta, \mathbf{x}) = -KL(q_{\phi}(z|x) || p_{\theta}(z)) + E_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]$$

$$E_{q_{\phi}(z|x)} [\log q_{\phi}(z|x)] \simeq \frac{1}{L} \sum_{l=1}^L (\log p_{\theta}(x|z^l))$$

Where $z^l = g_{\phi}(\varepsilon^l, x)$
and $\varepsilon^l \sim p(\varepsilon)$
and $L = \#$ of samples

20

20

The Reparameterization Trick

With minibatches:

- Sample M datapoints $\mathbf{X}^M = \{\mathbf{x}^1, \dots, \mathbf{x}^M\}$ from the full dataset of N datapoints
- Compute

$$\widehat{ELBO}^M(\phi, \theta, \mathbf{X}^M) = \frac{N}{M} \sum_{i=1}^M \widehat{ELBO}(\phi, \theta, \mathbf{x}^i)$$

21

21

The Variational Autoencoder

How do we choose p_θ and q_ϕ ?

- Could use standard probabilistic graphical models
- Or you could use a neural network to parameterize the distributions p_θ and q_ϕ

22

22

The Variational Autoencoder

Example:

- $p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}; \mu(\mathbf{z}), \text{diag}(\sigma(\mathbf{z})^2))$

Outputs of neural networks

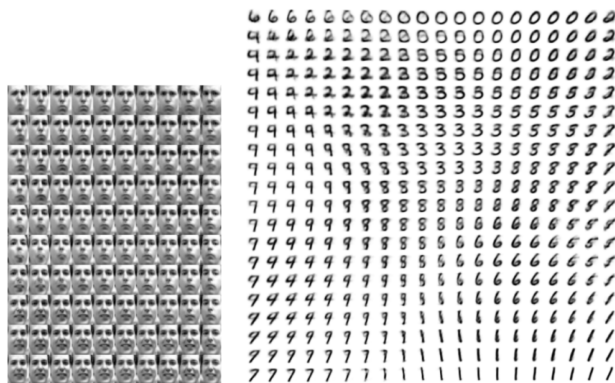
- $q(\mathbf{z}|\mathbf{x}) = N(\mathbf{z}; \mu(\mathbf{x}), \text{diag}(\sigma(\mathbf{x})^2))$

- $p(\mathbf{z}) = N(\mathbf{z}; \mathbf{0}, \mathbf{I})$

23

23

The Variational Autoencoder



(a) Learned Frey Face manifold

(b) Learned MNIST manifold

Results on 2D latent space. From: Kingma, D. and Welling, M. (2014). Auto-Encoding Variational Bayes. In ICLR 2014. <https://arxiv.org/pdf/1312.6114.pdf>

24

24