

Variational Autoencoders

Spring 2020

11-785 Introduction to Deep Learning

By Lady Yuying Zhu

And Lord Christopher George

Outline

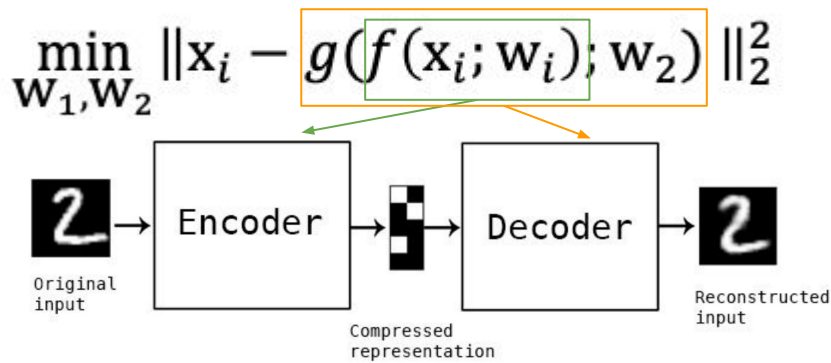
- **Autoencoders**
- Bake a cake 🤔 TBD, vanilla please
- *Probability Refresher*
- **Generative Models (basic overview)**
- Take over the universe
- Discuss our eventual doomed existence and the impact of covid-19 on society
- **Variational Autoencoders**

Auto-encoder Recap

A neural network that the output is the input itself.

Intuition:

- A good representation should keep the information well (reconstruction error)
- Deep + nonlinearity might help enhance the representation power



Why auto-encoder? Why not?

Why?

- Map high-dimensional data to low-dimensions for visualization
- Learn the salient features of the data.

Why not?

Encoded data can be decoded without loss if the autoencoder has enough degrees of freedom (leads to severe overfitting)

Without explicit regularisation, some points of the latent space are “meaningless” once decoded

Irregular latent space prevent us from using autoencoder for new content generation

Probability Refresher

- Continuous vs. Discrete
- Bayes Rule

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- Prior/ Posterior

$$P(A|B) \propto P(B|A)P(A)$$

posterior likelihood prior

- KL Divergence: measure the closeness of the two distributions

$$KL(q||p) = \mathbb{E}_{q(\mathbf{z};\lambda)}[\log q(\mathbf{z}; \lambda) - \log p(\mathbf{z}|\mathbf{x})]$$

- ELBO (evidence lower bound)

$$\text{ELBO}(\lambda) = \mathbb{E}_{q(\mathbf{z}; \lambda)}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \lambda)]$$

Generative Models, An Overview

Task: generate new samples follows the same probabilistic distribution of a given a training dataset



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

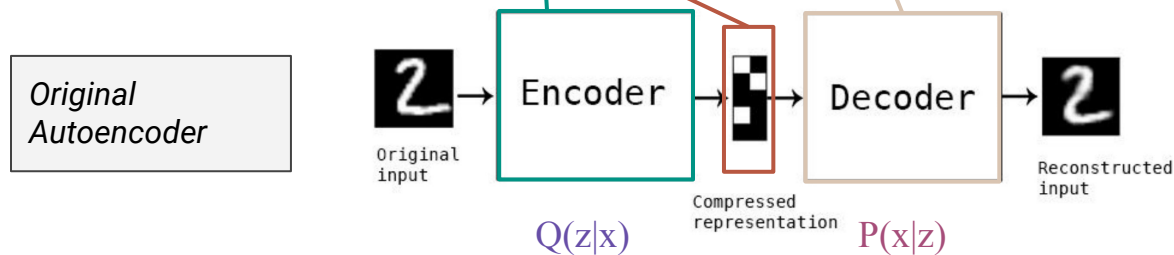


Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

Variational Autoencoders from Autoencoders

“In VAEs, the encoder becomes a variational inference network, mapping observed inputs to (approximate) posterior distributions over latent space, and the decoder becomes a generative network, capable of mapping arbitrary latent coordinates back to distributions over the original data space.”



Variational Autoencoders, (VAEs)

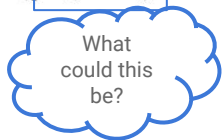
Big Idea: Maximize likelihood of seeing our data. “We are aiming maximize the probability of each X in the training set under the entire generative process.”

$$P(X) = \int P(X|z; \theta)P(z)dz$$

Where X is a datapoint, z is the latent variable, and theta are the parameters used for calculating z.

Generally in VAEs $P(X|z; \theta) = \mathcal{N}(X|f(z; \theta), \sigma^2 * I)$

And $P(z) = \mathcal{N}(z|0, I)$



Note: Why are we doing this? Because Variational Autoencoders can be interpreted as using variational bayesian inference where, in this bayesian view, our data is seen as being pulled from some underlying distribution. $x \sim p(x)$

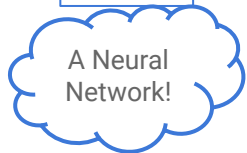
VAEs, Question 1

How do we define the latent variables z (i.e., decide what information they represent)?

* Samples of z can be drawn from a simple distribution, i.e., $N(0, I)$, where I is the identity matrix.

* Provided powerful function approximators, we can simply learn a function which maps our independent, normally-distributed z values to whatever latent variables might be needed for the model, and then map those latent variables to X .

$$P(X|z; \theta) = \mathcal{N}(X | f(z; \theta), \sigma^2 * I)$$



VAEs, Question 2

How do we deal with the integral over z ?

- * VAEs alter the sampling procedure to make it faster!
- * The key idea behind the variational autoencoder is to attempt to sample values of z that are likely to have produced X , and compute $P(X)$ just from those. This means that we need a new function $Q(z|X)$ which can take a value of X and give us a distribution over z values that are likely to produce X .
- * Given this $Q(Z|x)$, we can compute $P(X|z)$ relatively easily, so...

VAEs and KL Divergence

So now we need to relate $E_{z \sim Q} P(X|z)$ and $P(X)$, which we can do by minimizing the KL Divergence:

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(z|X)] \quad \text{By definition}$$

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(X|z) - \log P(z)] + \log P(X) \quad \text{Apply Bayes Rule}$$

$$\log P(X) - \mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z) \| P(z)] \quad \text{Rearrange}$$

$$\log P(X) - \mathcal{D} [Q(z|X) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z|X) \| P(z)] \quad \text{Reconstruct Q}$$

$$\log P(X|z) - \mathcal{D} [Q(z|X) \| P(z)] \quad \text{Take the gradient of RHS}$$

Minimizing KL Divergence = Maximizing ELBO

VAEs and KL Divergence

$$\log P(X|z) - \mathcal{D} [Q(z|X) \| P(z)]$$

Minimizing KL Divergence = Maximizing ELBO

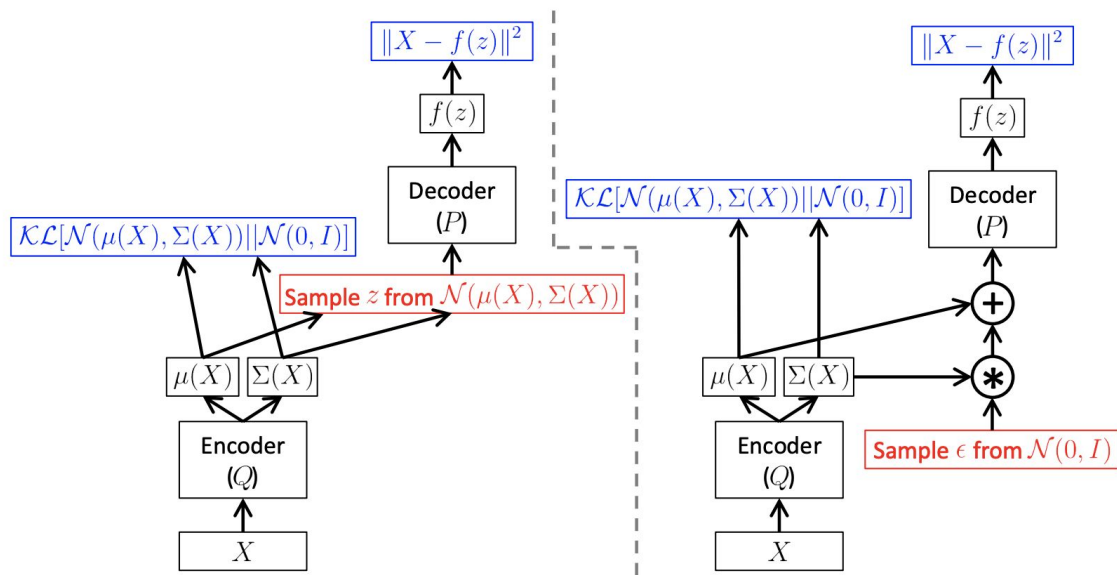
Negative Reconstruction Loss: Can take sample of z and treat $P(X|z)$ as an approximation, then backpropagate through the neural network. → Reparameterization trick.

KLD Loss: Has a closed form solution

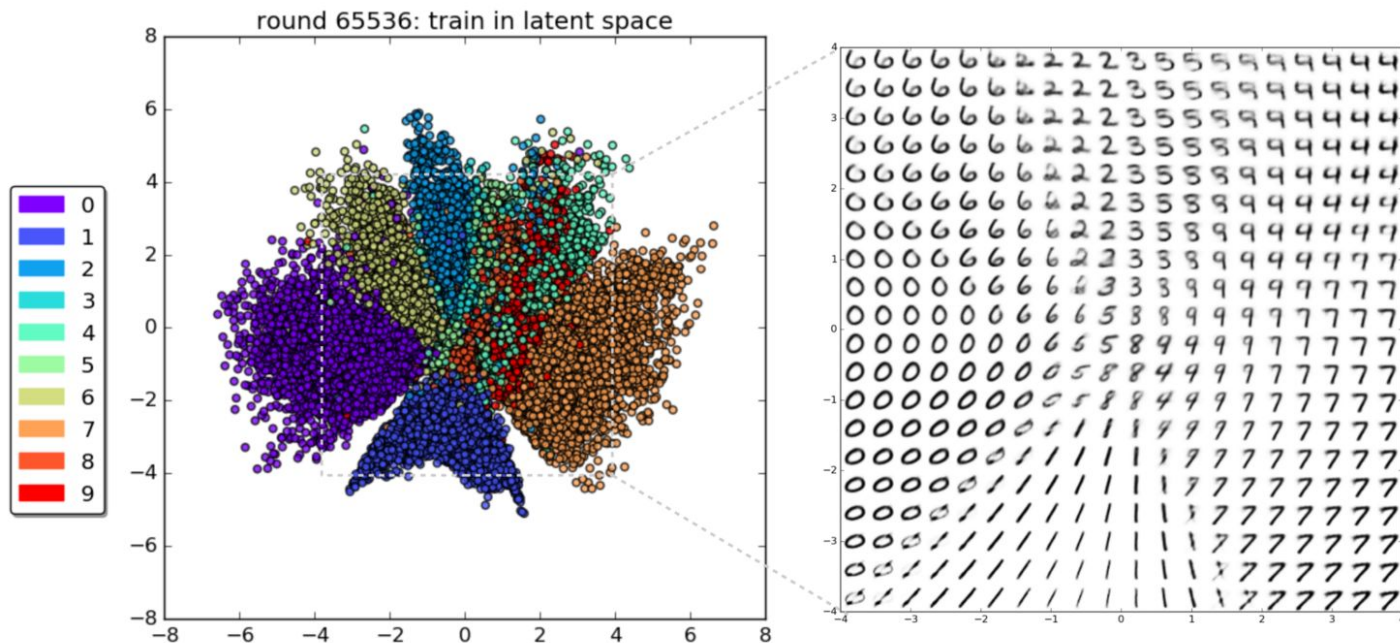
$$\mathcal{D}[\mathcal{N}(\mu(X), \Sigma(X)) \| \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma(X)) + (\mu(X))^T (\mu(X)) - k - \log \det(\Sigma(X)) \right)$$

Note: ELBO is the Evidence Lower Bound, and can be reformulated as the equation above.

VAEs, Reparameterization



VAEs, Example



References

- (1) Shenlong, Wang [Deep Generative Models](#)
- (2) [Chapter 20, Deep Generative Models](#)
- (3) [Tutorial on Variational Autoencoders](#)
- (4) [Fast Forward Labs](#), Under the Hood of the Variational Autoencoder (in Prose and Code)
- (5) [Fast Forward Labs](#), Introducing Variational Autoencoders (in Prose and Code)
- (6) [examples/main.py at master · pytorch/examples](#)
- (7) [CS231n](#), stanford