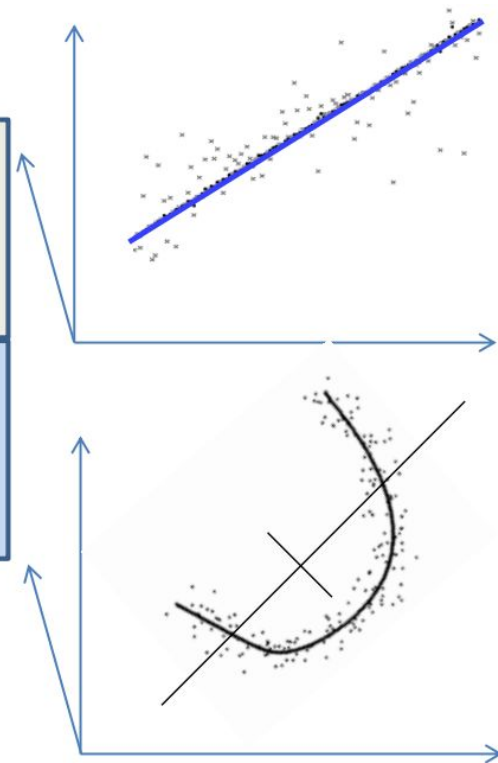
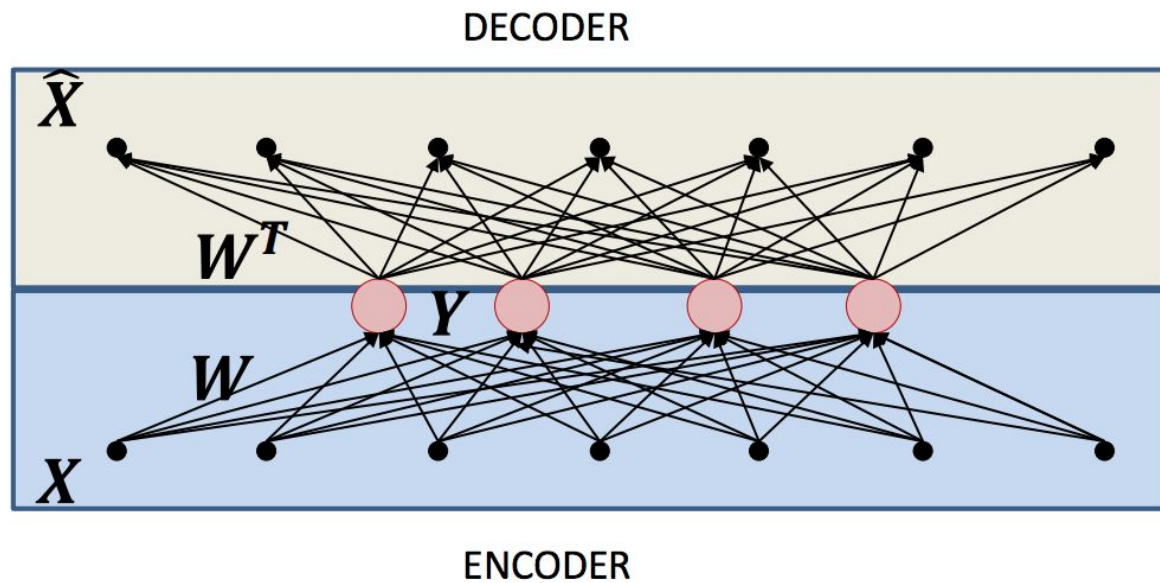


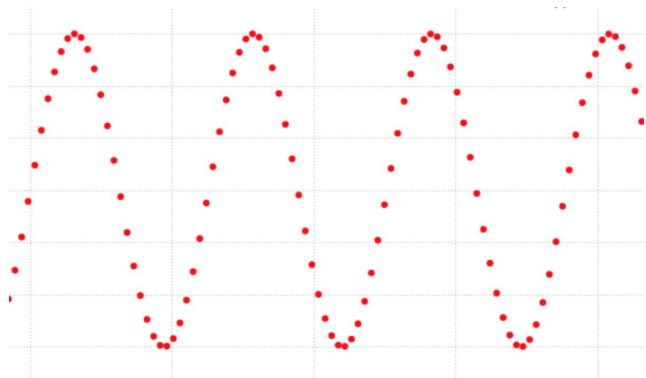
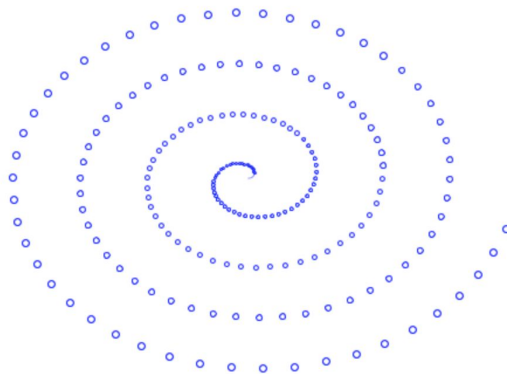
Variational Autoencoders

Recitation 9

Problem Setup



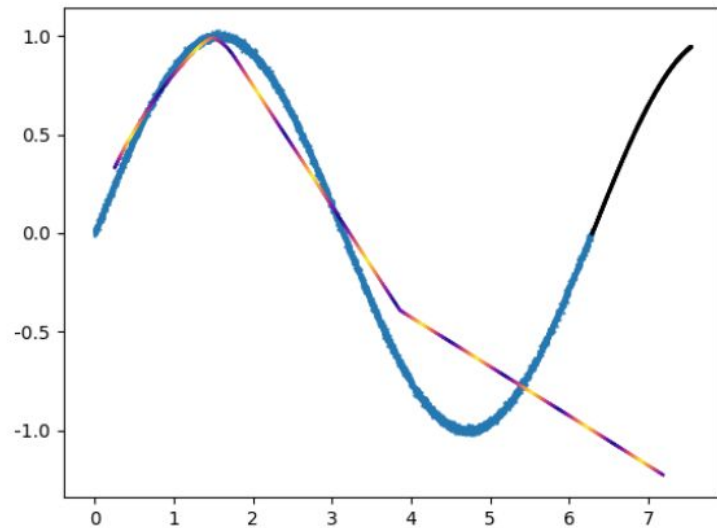
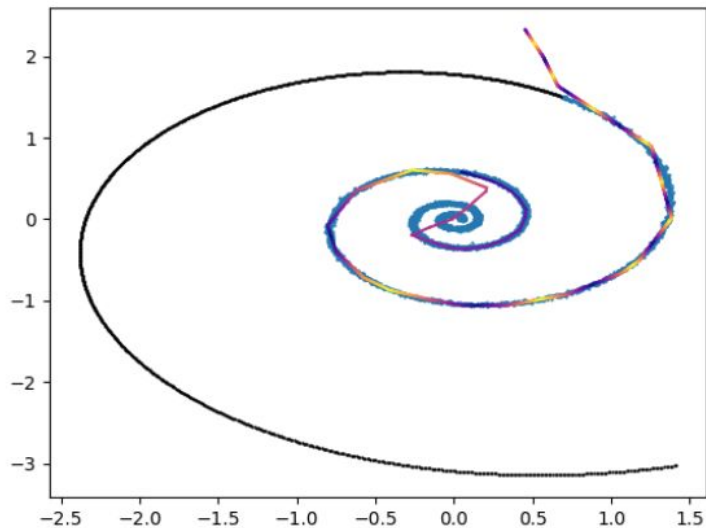
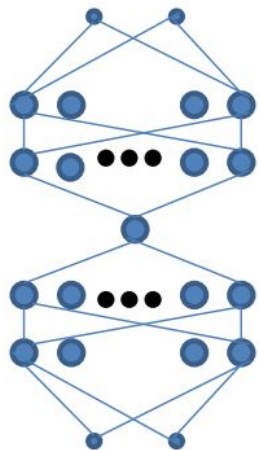
Problem Setup



Autoencoders are able to find lower dimensional representations of data

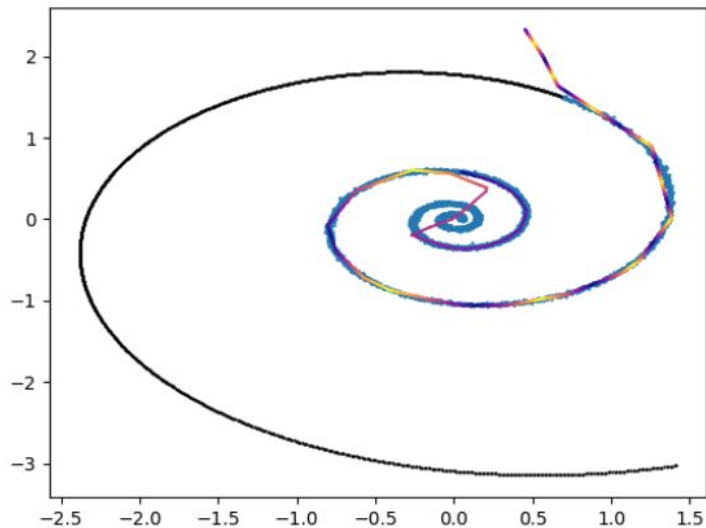
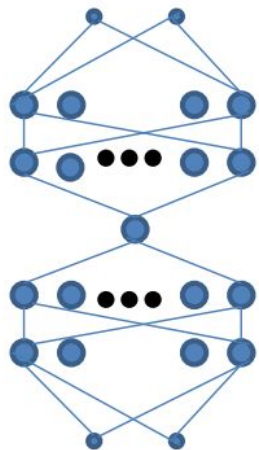
We want the representations to not only reconstruct the inputs, but also be able to interpolate and/or extend the data

Problem Setup



Data that are evidently generated from a single parameter cannot be 'learnt' by a vanilla AE

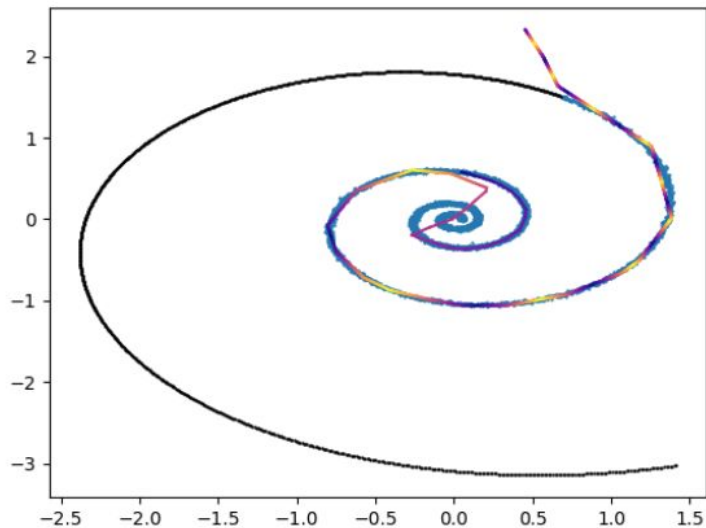
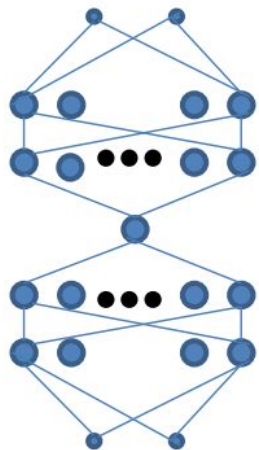
Problem Setup



Spiral fails to smoothly interpolate between training data; unable to follow the curve beyond the last training example.

Data that are evidently generated from a single parameter cannot be learnt by a vanilla AE

Problem Setup

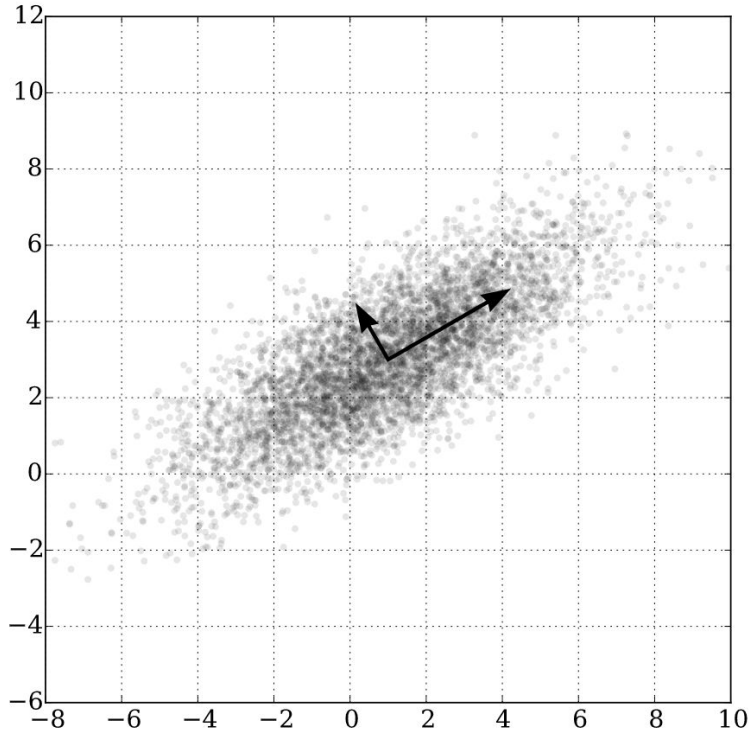


Spiral fails to smoothly interpolate between training data; unable to follow the curve beyond the last training example.

How do you find an appropriate representation?

Data that are evidently generated from a single parameter cannot be learnt by a vanilla AE

A long time ago in a semester far, far away

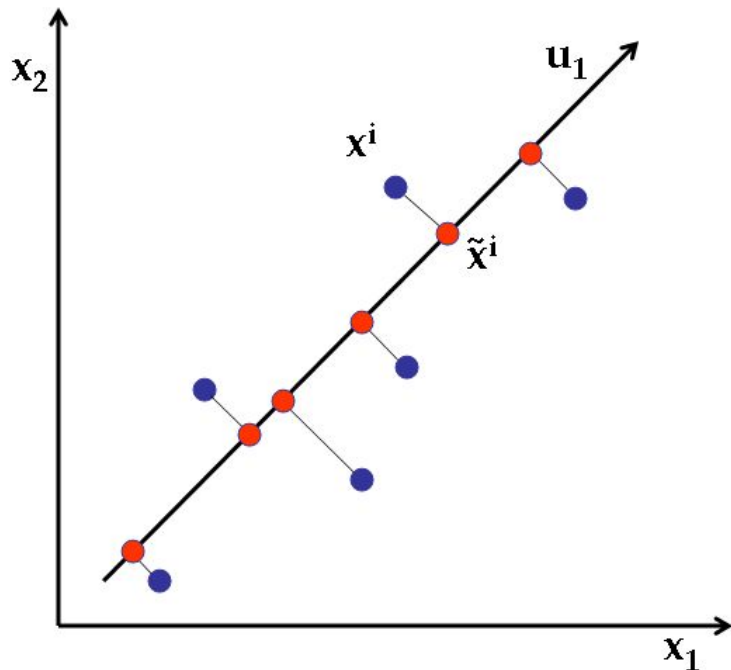


Principal Component Analysis

Fit an n-dimensional ellipsoid to the data

Minimize distortion

Principal component analysis

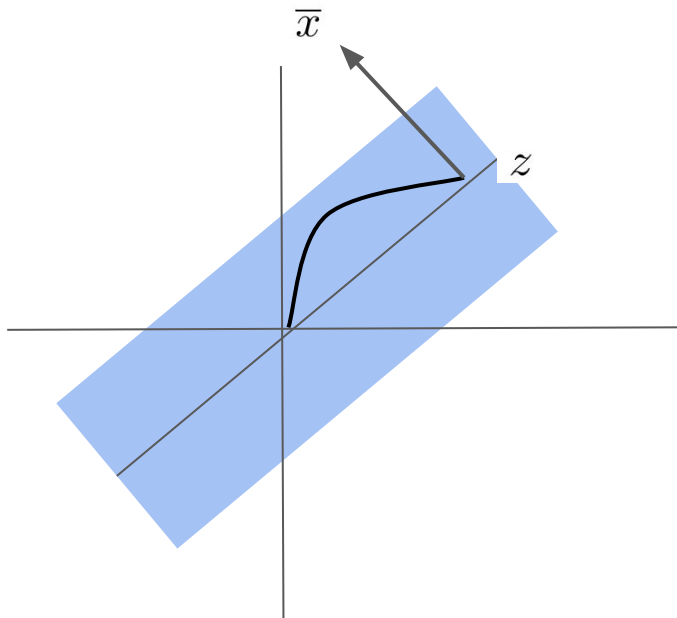


$$X = Vz + \epsilon$$

$$\operatorname{argmin}_{V,z} \|X - Vz\|^2 + \Lambda(VV^T - \mathbb{I})$$

Fits training set by minimizing reconstruction error - no discussion of generality

Statistical interpretation of PCA



$$\bar{x} = Vz + \epsilon$$

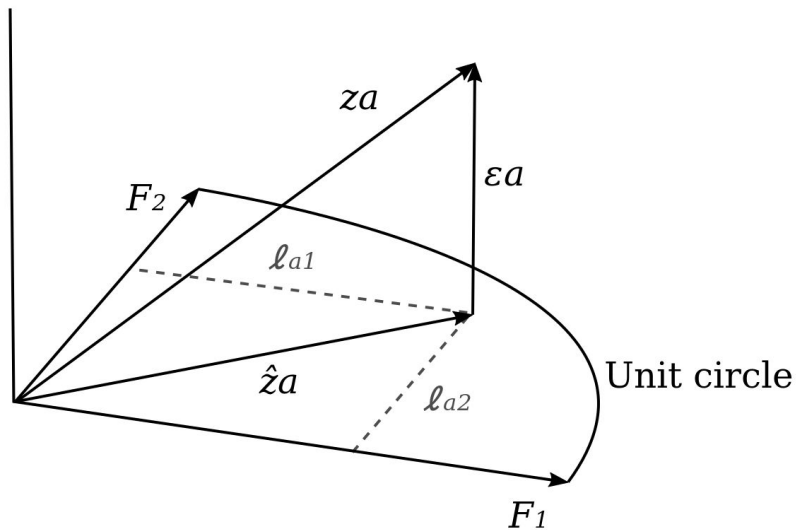
$$z \sim \mathcal{N}(0, B)$$

$$\epsilon \sim \mathcal{N}(0, E)$$

Assumption - error term is orthogonal to plane defined by V i.e. signal is purely along principal subspace, noise is purely orthogonal to principal subspace

Not true in general!

Factor Analysis



$$\bar{x} = Vz + \epsilon$$

$$z \sim \mathcal{N}(0, B)$$

$$\epsilon \sim \mathcal{N}(0, E)$$

Here, E is full rank - no longer constrained to be perpendicular to the plane

FA assumes there are underlying latent factors that affect our observed data

Autoencoders: our set of assumptions

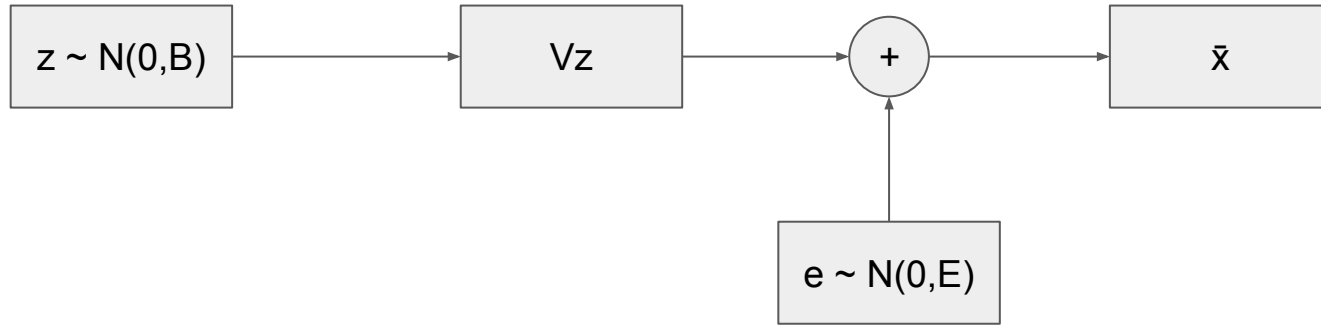
1. All data is generated as a continuous function of some latent variables
2. The exact choice of how latent variables map to data is up to us to decide
3. All our data is unlabeled, so nothing can be used as something related to the latent variables
4. (VAE's only) The function from latent variables to data is probabilistic in nature

What's the our reason behind holding each assumption?

Can you think of different problem setups where each assumption doesn't hold?

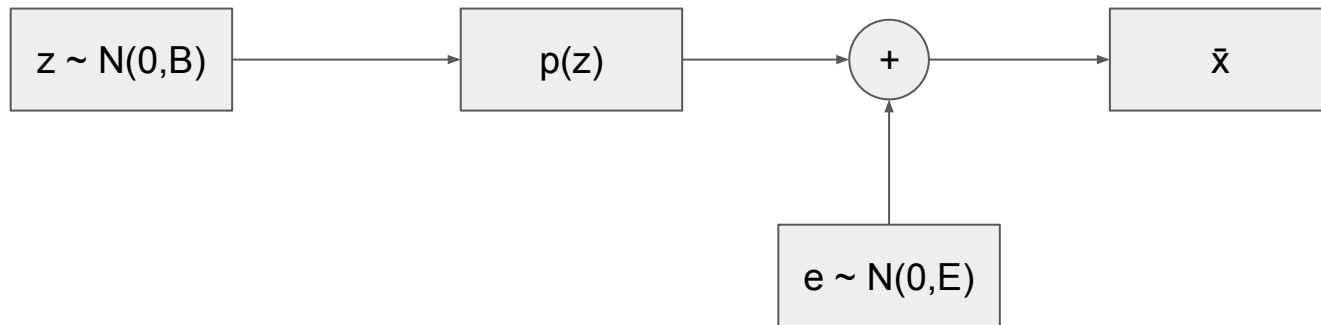
Why does having these assumptions force us into something similar to an encoder-decoder type of arrangement? How does EM address issues we face while holding these assumptions

FA as a generative model



More powerful - error need not be orthogonal to V
Still linear in z !

Non-linear FA

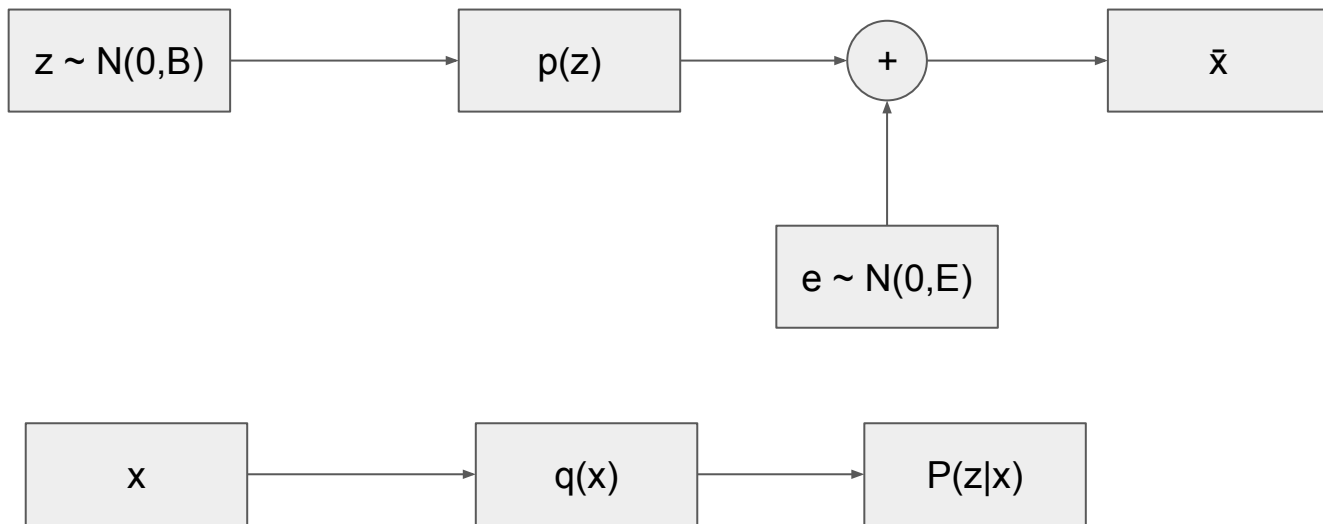


Even more powerful - Gaussian step can now be along a curved surface!

This is the motivation behind VAE. Learn $p(z)$ using nnet

$p(\cdot)$ can learn to transform random variable Z into a completely different distribution

VAE setup



Training objective for VAEs

We want to maximize log likelihood of data given model parameters.

$$\log P(\mathcal{X}) = \sum_{\mathbf{x} \in \mathcal{X}} \log P(\mathbf{x}; \theta)$$

Likelihood of each example is $P(\mathbf{x}; \theta) = \int P(\mathbf{x} | \mathbf{z}; \theta) P(\mathbf{z}) d\mathbf{z} \approx \sum_{\mathbf{z} \in S(\mathbf{x})} P(\mathbf{x} | \mathbf{z}; \theta)$

where $S(\mathbf{x}) := \{\mathbf{z}'; \mathbf{z}' \sim P(\mathbf{z})\}$

But sampling from high dimensional latent spaces is inefficient!

Sampling

$$S(\mathbf{x}) := \{\mathbf{z}'; \mathbf{z}' \sim P(\mathbf{z})\}$$

Introduce some $Q(\mathbf{z}|\mathbf{x})$ as distribution over \mathbf{z} that likely produced this \mathbf{x} to approximate the true $P(\mathbf{z}|\mathbf{x})$. From KL divergence we arrive at

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

We want to maximize LHS, and we are going to show how to optimize RHS.

Training

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

We want to optimize the above objective for X drawn from data.

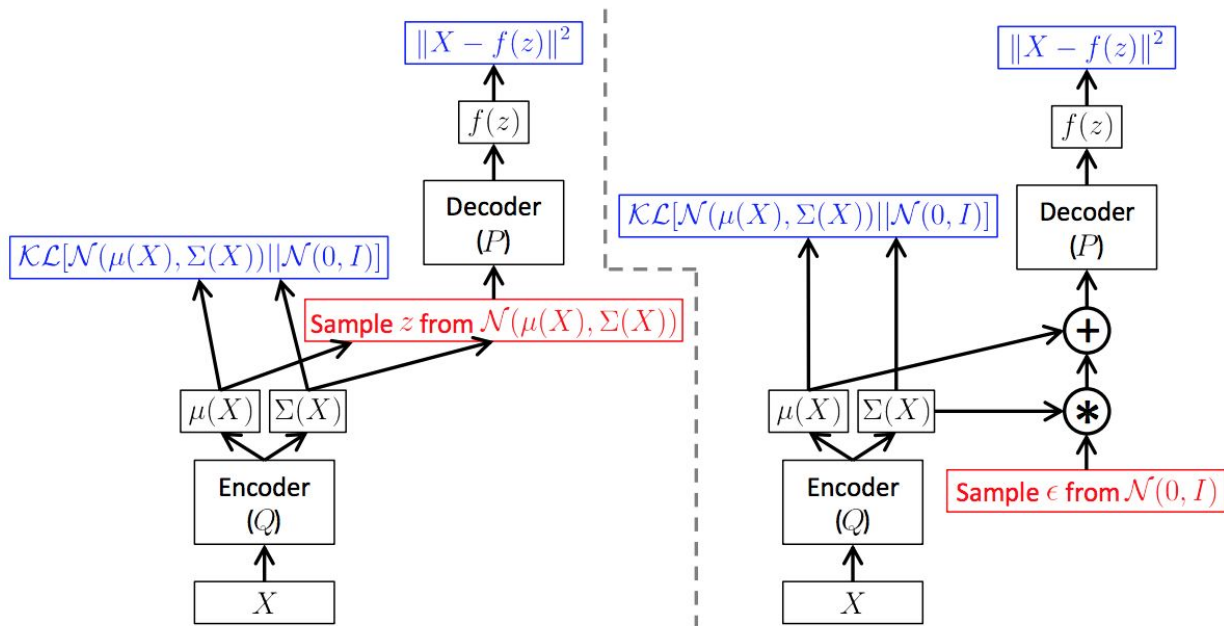
$$Q(z|X) \sim \mathcal{N}(\mu, \Sigma)$$

Sample one z , use it to approximate first term, then compute gradient of

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

‘Reparametrize’ in order to backpropagate through the sampling term by treating $\mathcal{N}(0, \mathbb{I})$ as a component of the input layer

Backpropagation



Notebook

MNIST example

