

Autoencoder and Variational Autoencoder

Group Meeting

Reporter: Liang Cao

2020.06.17

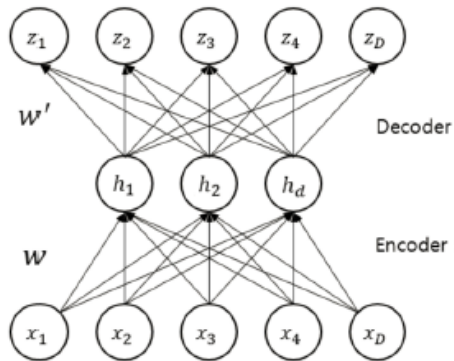
Introduction

- ✓ An **autoencoder** neural network is an unsupervised learning algorithm that applies back-propagation, setting the **outputs** to be equal to the **inputs**.
- ✓ Typically for the purpose of dimensionality reduction; was firstly introduced as a way of conducting pretraining in ANNs.
- ✓ The learned hidden layers are called the **representation(encoding)** of the input data.

Introduction

- ✓ Unsupervised learning: automatically extract meaningful features for you data; enhance the availability of unlabeled data.

- Autoencoder is a feedforward neural network trained to reproduce its input at the output layer: $\mathbf{z} \approx \mathbf{x}$



- Encoder:

$$\mathbf{h} = f_{\theta}(\mathbf{x}) = a(\mathbf{W}\mathbf{x} + \mathbf{b})$$

- Decoder:

$$\mathbf{z} = g_{\theta'}(\mathbf{h}) = a(\mathbf{W}'\mathbf{h} + \mathbf{b}')$$

$$\theta, \theta' : \{\mathbf{W}, \mathbf{b}\}, \{\mathbf{W}', \mathbf{b}'\}$$

- Parameter estimation:

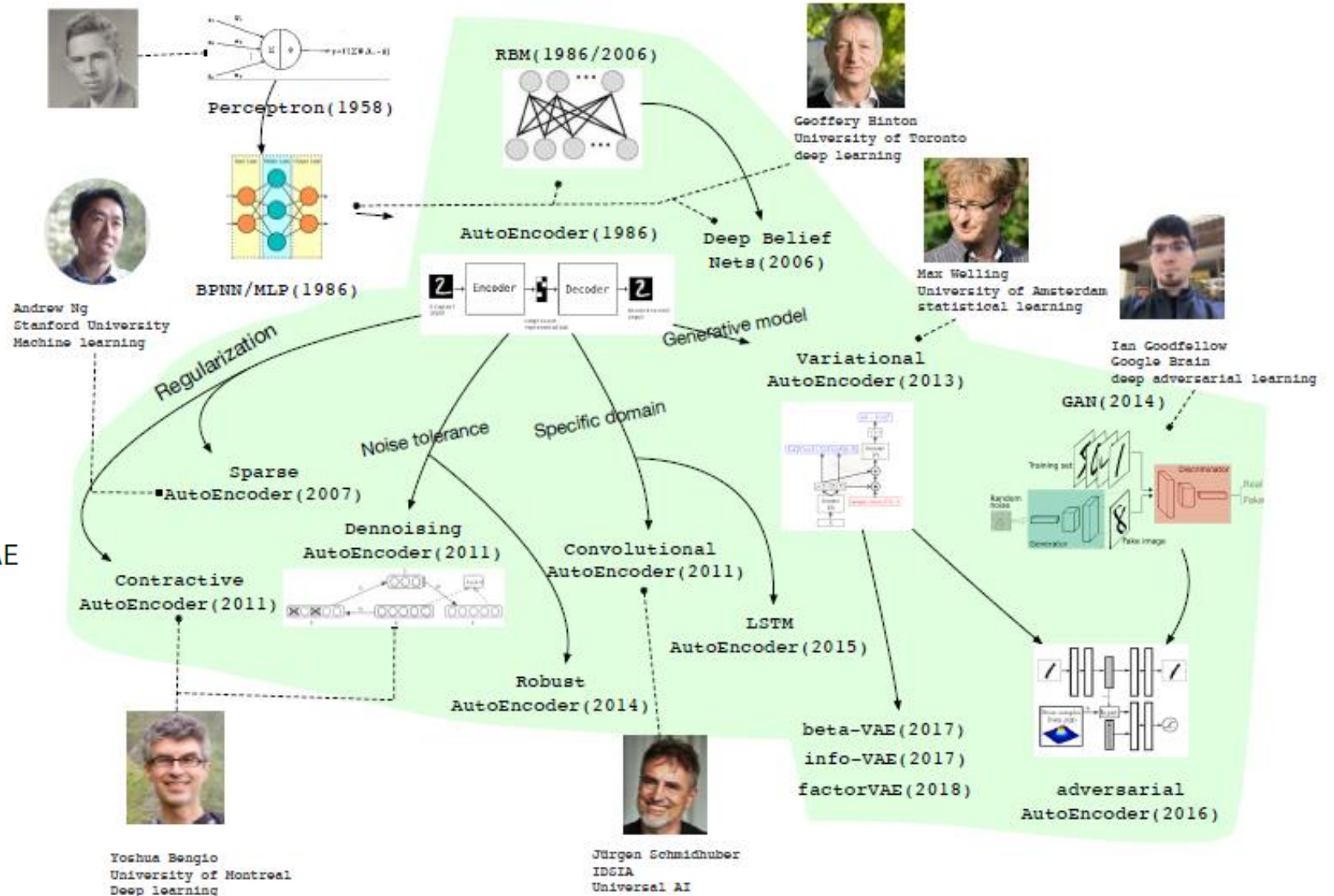
$$\begin{aligned} \theta^*, \theta'^* &= \operatorname{argmin}_{\theta, \theta'} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) \\ &= \operatorname{argmin}_{\theta, \theta'} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, g_{\theta'}(f_{\theta}(\mathbf{x}^{(i)}))) \end{aligned}$$

where L is a loss function such as *squared error*
 $L(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|^2$.

Introduction

- What is a **good** representation ?

- Sparsity? \Rightarrow Sparse AE
- Denoise? \Rightarrow Denoising AE
- Robust to disturbance? \Rightarrow contractive AE
- High level? \Rightarrow Stacked AE
- Specific data (image) structure? Convolutional AE
- Generate new data? \Rightarrow Variational AE

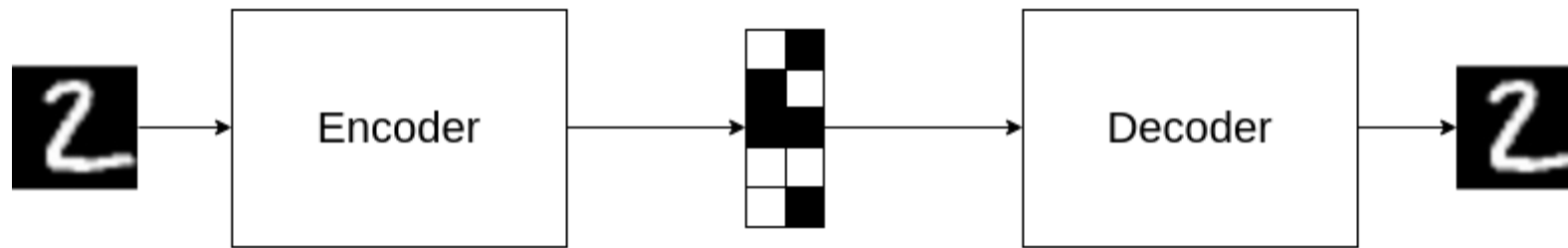


Diederik P Kingma, Max Welling. Universiteit van Amsterdam. **Auto-Encoding Variational Bayes**. December, 2013

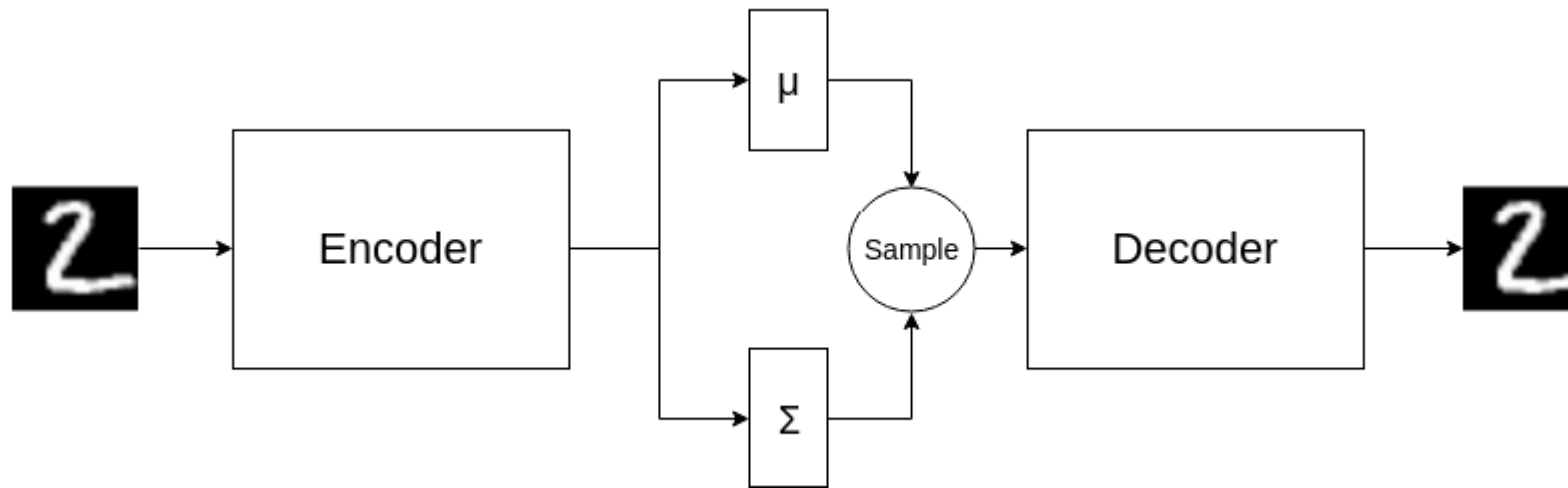


Number of citations: 9201

- ✓ VAE is a **generative model**, we can use generative model to only learn the distribution of our data $p(x)$.
- ✓ After training we can **generate new data** similar to x .
- ✓ Generated data instances should come from points with high probability in datasets distribution space.



Classical autoencoders



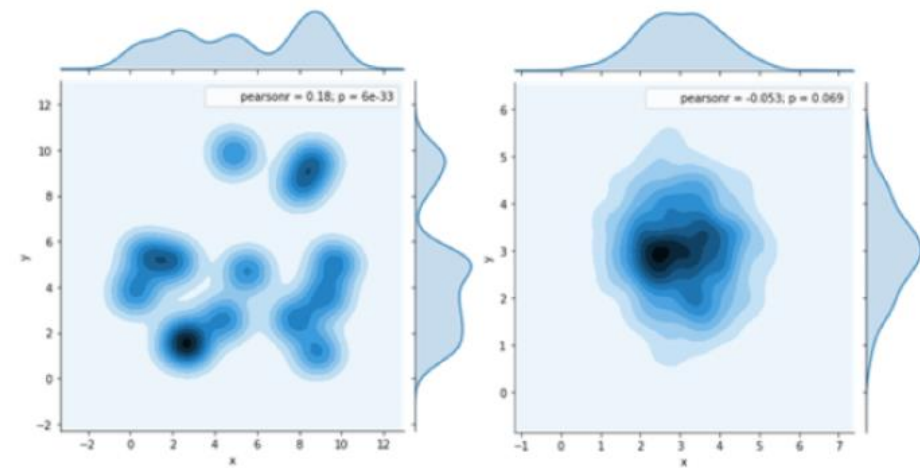
Variational autoencoders

1. Classical autoencoders minimize a reconstruction loss $\|x - \hat{x}\|^2$, they are unregularized in latent space.

2. VAEs are one approach to regularizing (impose structure) the latent distribution.

Classical autoencoders minimize a reconstruction loss $\|x - \hat{x}\|^2$, they are unregularized in latent space.

- **Cons**
 - ✓ This yields an **unstructured latent space**.
 - ✓ Examples from the data distribution are mapped to codes **scattered** in the space.
 - ✓ No constraint that similar inputs are mapped to nearby points in the latent space
 - ✓ We cannot sample codes to generate novel examples.



On the left the Conventional AE latent space and on the right VAEs latent space

Problem definition

- ✓ Observable data: $X = \{x_1, x_2, x_3, \dots, x_n\}$
- ✓ Hidden variable: $Z = \{z_1, z_2, z_3, \dots, z_m\}$

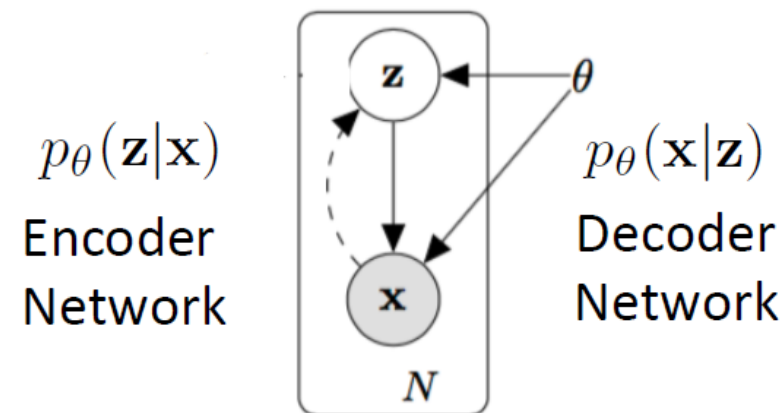
In encoder network, we need to do inference (calculate the posterior):

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Need to calculate evidence: $p(x) = \int p(x|z)p(z) dz$

For higher dimensional latent variable: $p(x) = \int \int \int \dots \int p(x|z_i)p(z_i) dz_1 dz_2 \dots dz_m$

$p(x)$ is **intractable**, which leads to $p(z|x)$ **intractable**.



$p(z|x)$ intractable.

Solutions for intractable posterior: **approximate inference**

- **Deterministic approximation (variational inference)**

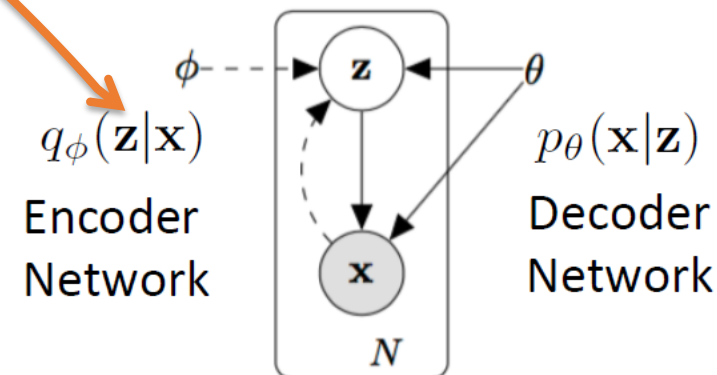
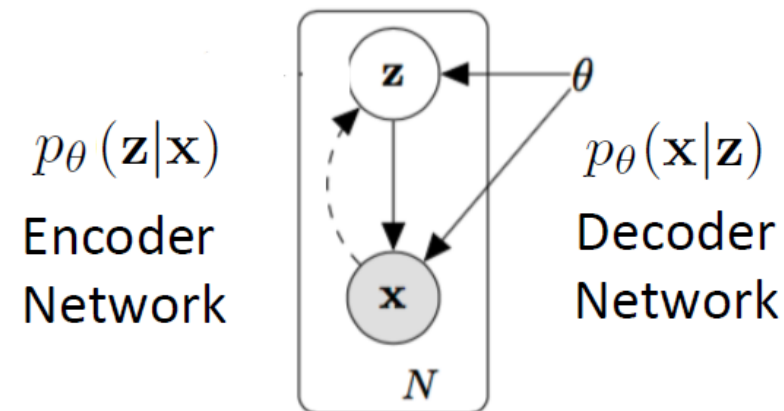
- Approximate p with "closest" distribution q from a tractable family,

$$p(z|x) \approx q(z|x)$$

- Turns inference into optimization (need to find best q).
 - Called **variational Bayes**.

- **Stochastic approximation (Markov Chain Monte Carlo)**

- Approximate p with empirical distribution over samples,
- Turns inference into sampling.



$p(z|x)$ **intractable**.

Solutions for intractable posterior: **approximate inference**

- Deterministic approximation (variational inference)**

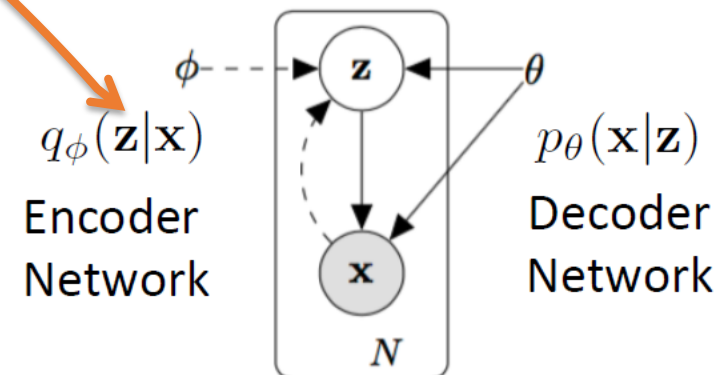
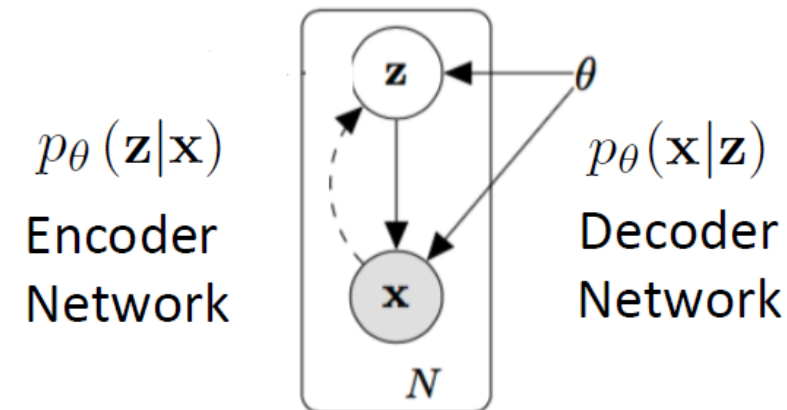
- Approximate p with "closest" **distribution q** from a tractable family,

$$p(z|x) \approx q(z|x)$$

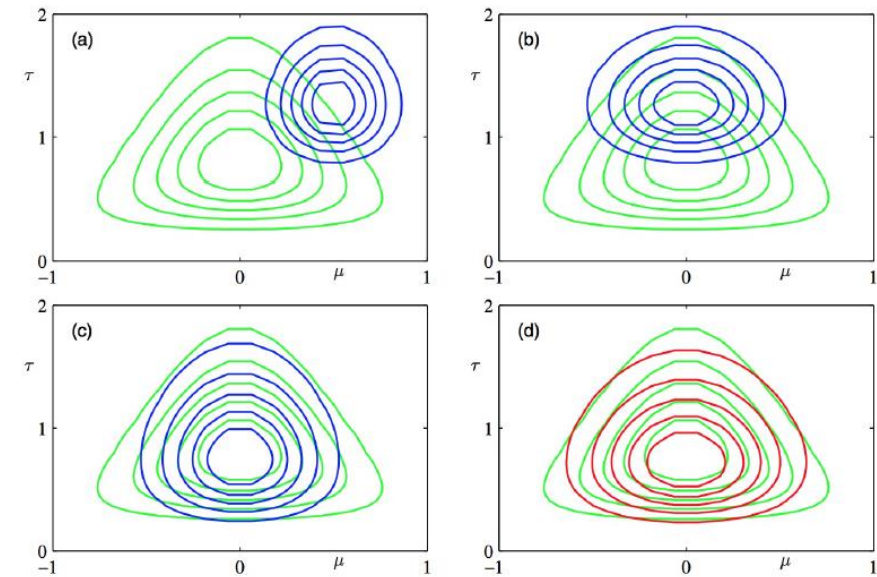
- Turns **inference into optimization** (need to find best q).
 - Called **variational Bayes**.

Why not use MCMC?

- MCMC works asymptotically, but may take forever.
- Variational methods **not consistent**, but very fast.
(trade off accuracy vs. computation)



- **Deterministic approximation (variational inference)** $p(z|x) \approx q(z|x)$
 - ✓ The main idea behind variational inference is to, first pick a tractable family of distributions over the latent variables with initial variational parameters.
 - ✓ Then to find parameters that make it as close as possible to the true posterior.
 - ✓ KL divergence measures information lost when using $q_\phi(z|x)$ to approximate $p_\theta(z|x)$
 - ✓ We want to choose ϕ to minimize $D_{KL}(q_\phi(z|x) || p_\theta(z|x))$



- KL divergence

Used to measure similarity between two probability distributions(w.r.t. one of them)

Discrete and continuous form:

- $D_{KL}(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$
- $D_{KL}(p||q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx$

Properties:

- $D_{KL}(p||q) \geq 0, \forall p, q$
- $D_{KL}(p||q) = 0 \iff p = q$
- $D_{KL}(p||q) \neq D_{KL}(q||p)$ in general

- Variational Lower Bound

$$\begin{aligned}
 D_{KL}(q(z|x)||p(z|x)) &= \int_z q(z|x) \log \frac{q(z|x)}{p(z|x)} \\
 &= - \int_z q(z|x) \log \frac{p(z|x)}{q(z|x)} & p(z|x) &= \frac{p(x,z)}{p(x)} \\
 &= - \left[\int_z q(z|x) \log \frac{p(x,z)}{q(z|x)} - \int_z q(z|x) \log p(x) \right] \\
 &= - \int_z q(z|x) \log \frac{p(x,z)}{q(z|x)} + \log p(x) \int_z q(z|x) \\
 &= -\mathcal{L} + \log p(x)
 \end{aligned}$$

- ✓ Still contains $p(x)$ term! So cannot compute directly.
- ✓ But $p(x)$ does not depend on parameter ϕ in $q_\phi(z|x)$, $p(x)$ is a constant for different ϕ , so still hope.
- ✓ Define L as variational lower bound.

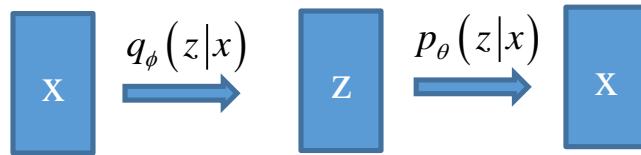
$$\log p(x) = \mathcal{L} + D_{KL}(q(z|x)||p(z|x))$$

- ✓ Minimizing the KL divergence is equal to maximizing variational lower bound L .

- Variational Lower Bound

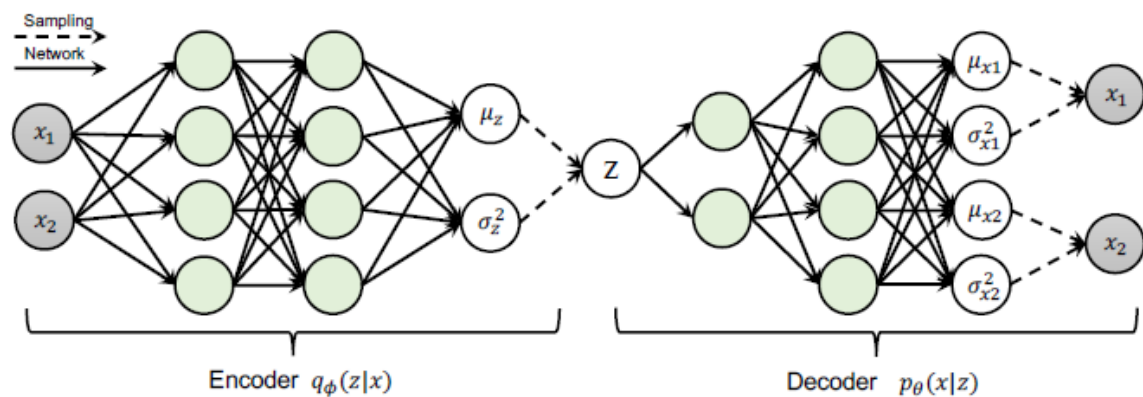
$$\begin{aligned}
 \mathcal{L} &= \int_z q(z | x) \log \frac{p(x, z)}{q(z | x)} \\
 &= \int_z q(z | x) \log \frac{p(x | z)p(z)}{q(z | x)} \\
 &= \int_z q(z | x) \log p(x | z) + \int_z q(z | x) \log \frac{p(z)}{q(z | x)}
 \end{aligned}$$

$$\mathcal{L} = \mathbb{E}_{q(z|x)} \log p(x | z) - D_{KL}(q(z || x) || p(z))$$



The first term is conceptually the negative reconstruction error and the second is regularize, makes approximate inference term close to the prior $p(z)$, the $p(z)$ is usually chosen as standard normal distribution.

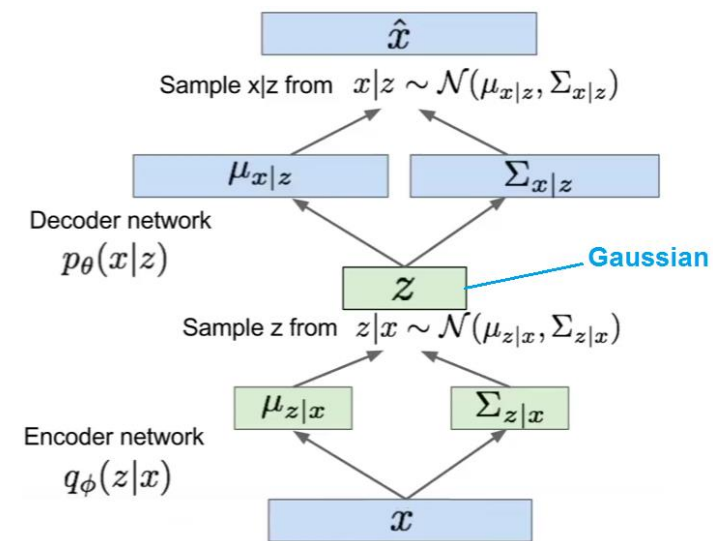
Deep Learning Perspective



- Encoder and decoder networks also called “recognition” / “inference” and “generation” networks
- We aim to learn the parameters ϕ and θ via backpropagation.

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

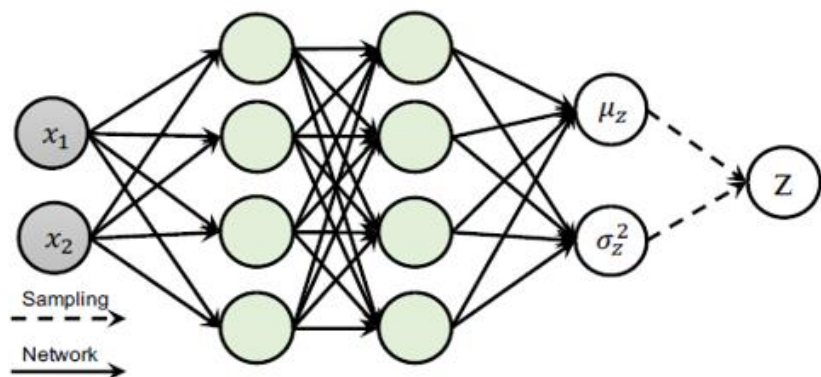
Probabilistic Model Perspective



Loss function: $L(\theta, \phi, \mathbf{x}) = -D_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})] + E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$

- Use $\mathcal{N}(0, 1)$ as prior for \mathbf{z} ; $q_\phi(\mathbf{z}|\mathbf{x})$ is Gaussian distribution $\mathcal{N}(\mu_{\mathbf{z}}(\mathbf{x}; \phi), \sigma_{\mathbf{z}}^2(\mathbf{x}; \phi))$ determined by NN.
 - The KL-divergence:

$$-D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) = \frac{1}{2}(1 + \log \sigma_{\mathbf{z}}^2 - \mu_{\mathbf{z}}^2 - \sigma_{\mathbf{z}}^2)$$



Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

The marginal likelihood is composed of a sum over the marginal likelihoods of individual datapoints $\log p_\theta(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)})$, which can each be rewritten as:

$$\log p_\theta(\mathbf{x}^{(i)}) = D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z}|\mathbf{x}^{(i)})) + \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) \quad (1)$$

The first RHS term is the KL divergence of the approximate from the true posterior. Since this KL-divergence is non-negative, the second RHS term $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$ is called the (variational) *lower bound* on the marginal likelihood of datapoint i , and can be written as:

$$\log p_\theta(\mathbf{x}^{(i)}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log q_\phi(\mathbf{z}|\mathbf{x}) + \log p_\theta(\mathbf{x}, \mathbf{z})] \quad (2)$$

which can also be written as:

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] \quad (3)$$

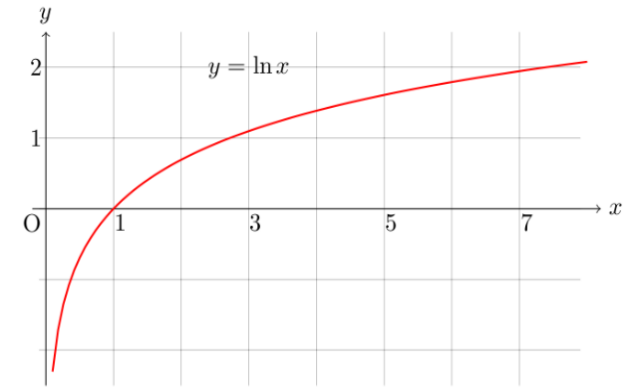
We want to differentiate and optimize the lower bound $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$ w.r.t. both the **variational parameters ϕ** and **generative parameters θ** . However, the gradient of the lower bound w.r.t. ϕ is a bit problematic. The usual (naïve) Monte Carlo gradient estimator for this type of problem is: $\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z})} \log q_\phi(\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)}) \nabla_{q_\phi(\mathbf{z}^{(l)})} \log q_\phi(\mathbf{z}^{(l)})$ where $\mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$. **This gradient estimator exhibits very high variance (see e.g. [BJP12]) and is impractical for our purposes.**

Loss function: $L(\theta, \phi, \mathbf{x}) = -D_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})] + E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$

If we use Monte Carlo gradient estimator:

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z})} \log q_\phi(\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)}) \nabla_{q_\phi(\mathbf{z}^{(l)})} \log q_\phi(\mathbf{z}^{(l)})$$

The gradient has **log term**, so it has **high variance** and needs lots of samples.

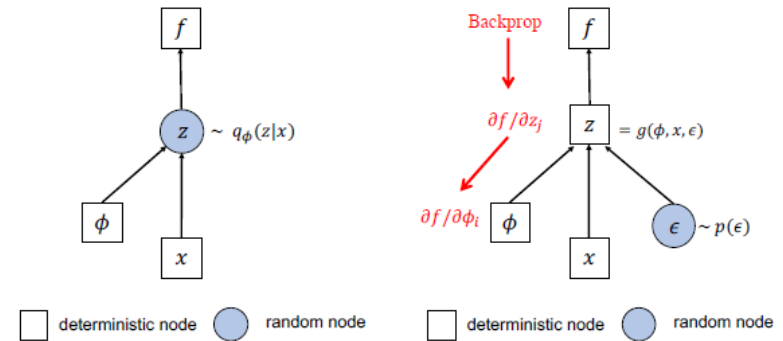


To solve this problem, we introduce **reparameterization trick**:

$$\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) \quad \Rightarrow \quad \begin{array}{l} \text{auxiliary variable} \\ \epsilon \sim p(\epsilon) \\ \text{deterministic variable} \\ \mathbf{z} = g_\phi(\epsilon, \mathbf{x}) \end{array}$$

Example:

$$z \sim p(z|x) = \mathcal{N}(\mu, \sigma^2) \quad \Rightarrow \quad \begin{array}{l} \epsilon \sim \mathcal{N}(0, 1) \\ z = \mu + \sigma\epsilon \end{array}$$



$$z^{(k)} \sim \mathcal{N}(\mu_z(x; \phi), \sigma_z^2(x; \phi))$$

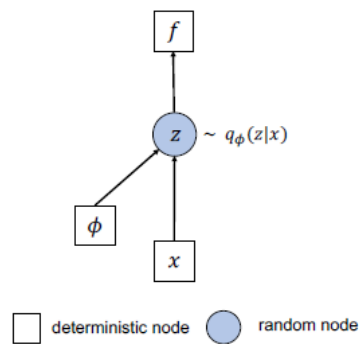
$$\begin{array}{l} \epsilon^{(k)} \sim \mathcal{N}(0, 1) \\ z^{(k)} = \mu_z(x, \phi) + \sigma_z(x, \phi) \cdot \epsilon^{(k)} \end{array}$$

Loss function: $L(\theta, \phi, \mathbf{x}) = -D_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})] + E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$

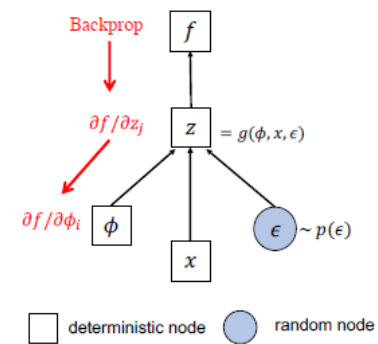
$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} [f(g_\phi(\epsilon, \mathbf{x}^{(i)}))] \simeq \frac{1}{L} \sum_{l=1}^L f(g_\phi(\epsilon^{(l)}, \mathbf{x}^{(i)})) \quad \text{where } \epsilon^{(l)} \sim p(\epsilon)$$

$$\epsilon^{(l)} \sim p(\epsilon), \quad l = 1, 2, \dots, L$$

$$\nabla_\phi \mathcal{L} \approx \frac{1}{L} \sum_{i=1}^L \left[\left(\frac{d}{dz} (\log p(x^{(i)}, z) - \log q(z|\phi)) \right) \left(\frac{d}{d\phi} g_\phi(\epsilon, \mathbf{x}^{(i)}) \right) \right] \quad \text{where } \mathbf{Z} = g_\phi(\epsilon, \mathbf{x})$$



$$z^{(k)} \sim \mathcal{N}(\mu_z(x; \phi), \sigma_z^2(x; \phi))$$



$$\epsilon^{(k)} \sim \mathcal{N}(0, 1)$$

$$z^{(k)} = \mu_z(x, \phi) + \sigma_z(x, \phi) \cdot \epsilon^{(k)}$$

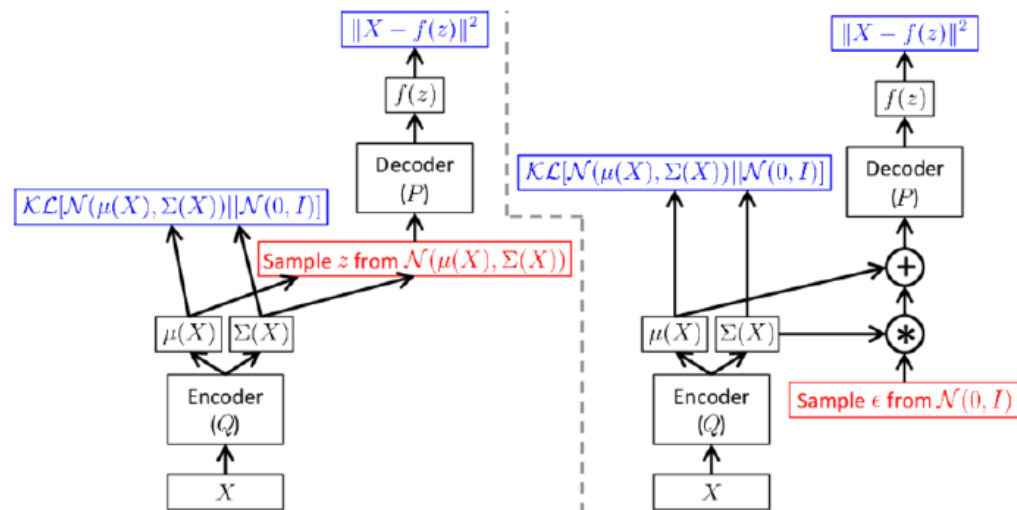
Loss function: $L(\theta, \phi, \mathbf{x}) = -D_{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})] + E_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} [f(g_\phi(\epsilon, \mathbf{x}^{(i)}))] \simeq \frac{1}{L} \sum_{l=1}^L f(g_\phi(\epsilon^{(l)}, \mathbf{x}^{(i)})) \quad \text{where } \epsilon^{(l)} \sim p(\epsilon)$$

$$\epsilon^{(l)} \sim p(\epsilon), \quad l = 1, 2, \dots, L$$

$$\nabla_\phi \mathcal{L} \approx \frac{1}{L} \sum_{l=1}^L \left[\left(\frac{d}{dz} (\log p(x^{(i)}, z) - \log q(z|\phi)) \right) \left(\frac{d}{d\phi} g_\phi(\epsilon, \mathbf{x}^{(i)}) \right) \right]$$

where $\mathbf{Z} = g_\phi(\epsilon, \mathbf{x})$



• Loss function:

$$L(\theta, \phi, \mathbf{x}) \approx \frac{1}{2} \sum_j (1 + \log \sigma_{zj}^2 - \mu_{zj}^2 - \sigma_{zj}^2) + \frac{1}{L} \sum_k \log(p_\theta(x|z^{(k)}))$$

where $z^{(k)} = \mu_z(x, \phi) + \sigma_z(x, \phi) \cdot \epsilon^{(k)}$ and $\epsilon^{(k)} \sim \mathcal{N}(0, 1)$.

Training

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

$\theta, \phi \leftarrow$ Initialize parameters

repeat

$\mathbf{X}^M \leftarrow$ Random minibatch of M datapoints (drawn from full dataset)

$\epsilon \leftarrow$ Random samples from noise distribution $p(\epsilon)$

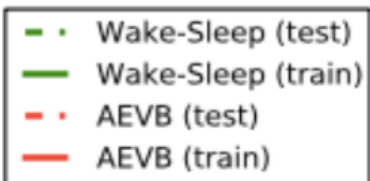
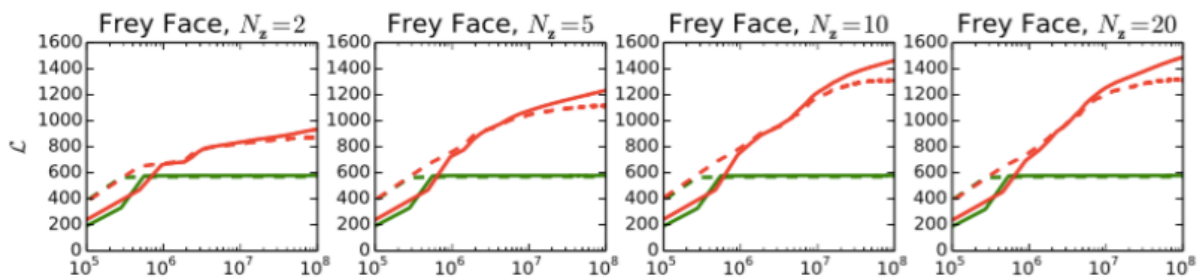
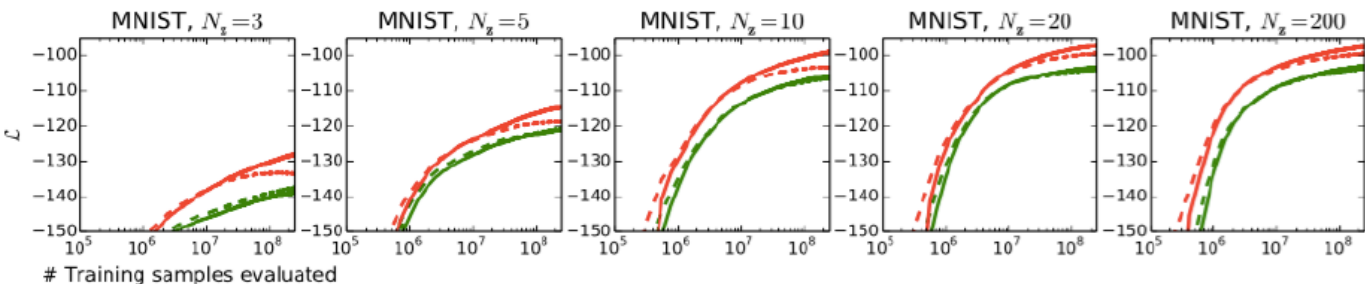
$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$ (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10])

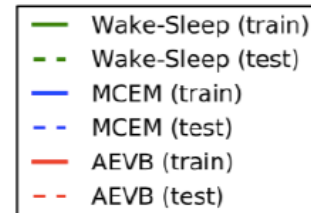
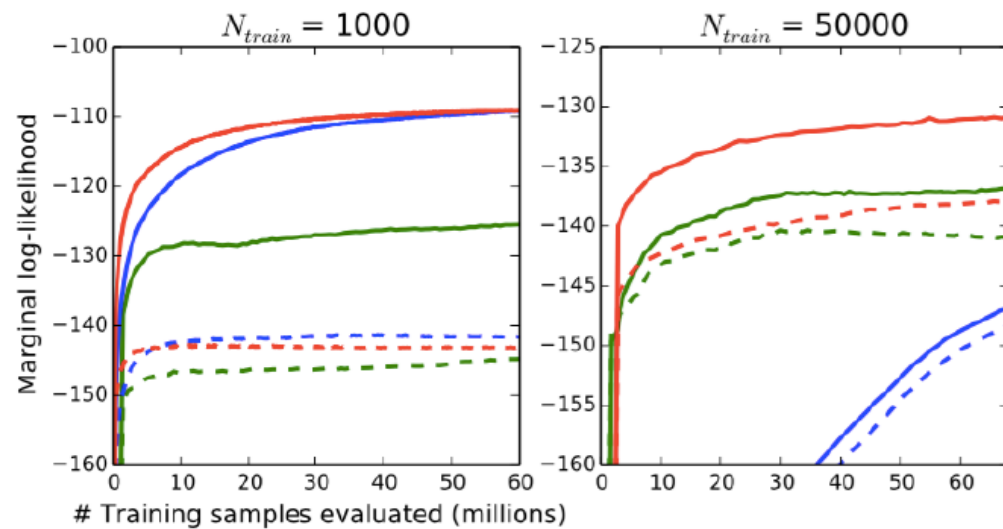
until convergence of parameters (θ, ϕ)

return θ, ϕ

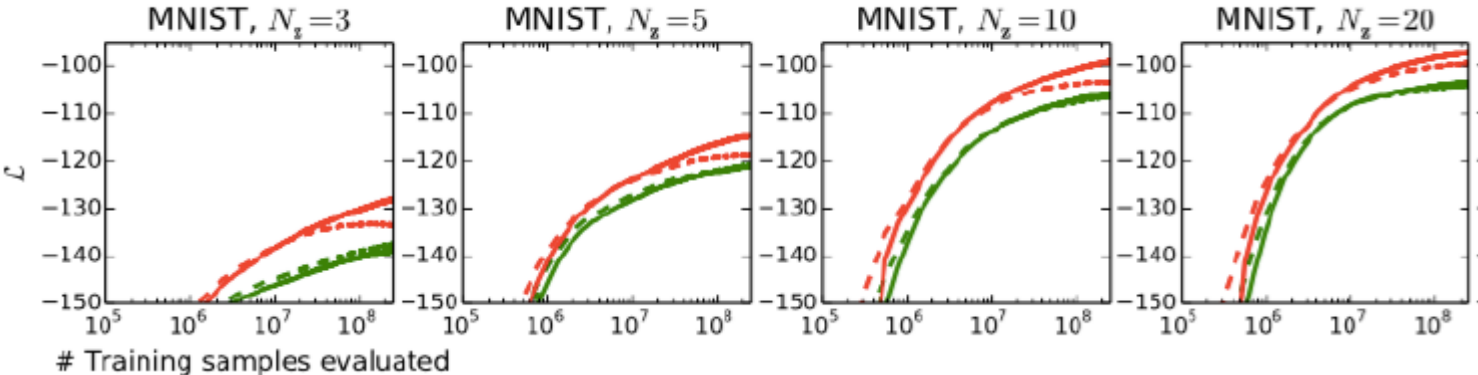
Results



Horizontal axis: size of training data
 Vertical axis: evidence Lower Bound
 N_z : dimensions of hidden variables



Results

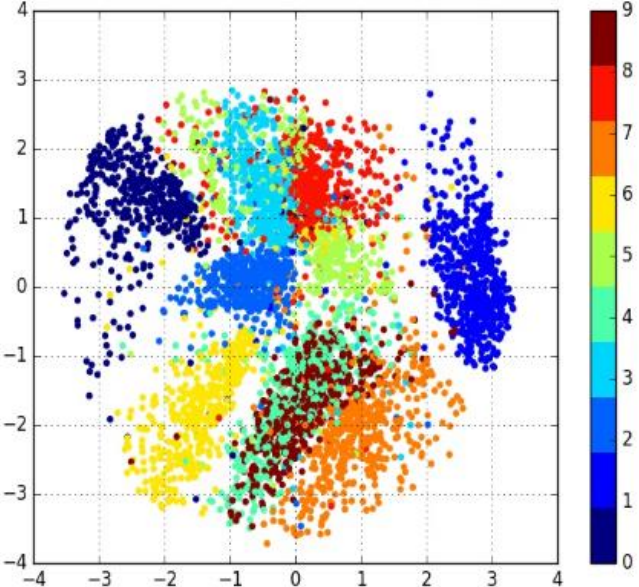
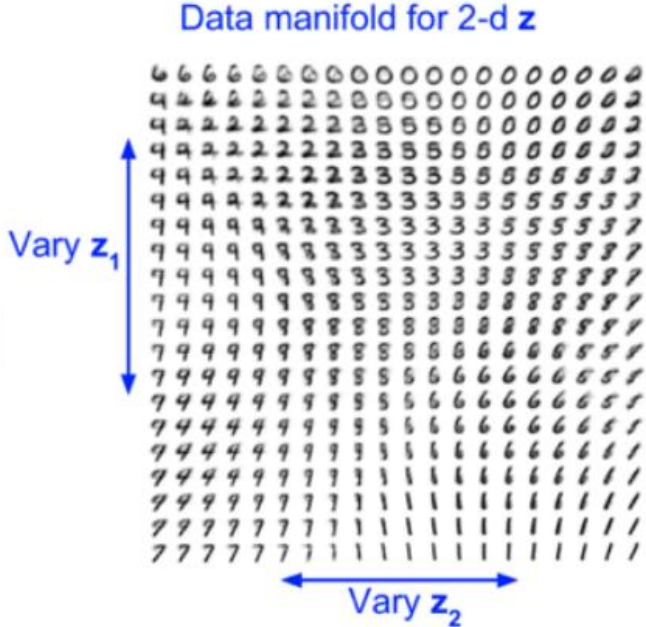
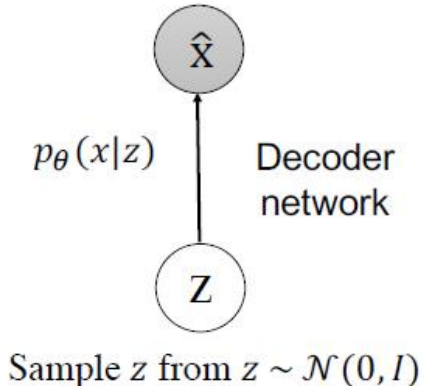


Input image	2-D latent space	5-D latent space	10-D latent space	20-D latent space
7210414959 0690159784 9665407401 3134727121 1742351244 6355604195 7893746430 7029173297 9627847361 3693141769	7210919989 0690139789 9665907901 3130727121 1792351299 6355604198 7892796430 7037193297 9627847361 3693141769	7210414999 0690159734 9665407401 3136727121 1742351294 6355604195 7893746430 7027173297 9627847361 3693141769	7210414959 0690159734 9665407401 3134727121 1742351244 6355604195 7892746430 7027173297 9627847361 3693141769	7210414959 0690159784 9665407401 3134727121 1742351244 6355604195 7893746430 7027173297 9627847361 3693141769

Well trained VAE must be able to reproduce input image, this figure shows reproduce performance of learned generative models for different dimensionalities.

Results: Generate new data

- Use decoder network, sample z from prior!



Visualizations of learned MNIST manifold for generative models with 2-dim and its latent space distribution.

<https://github.com/hwalsuklee/tensorflow-mnist-VAE>

Conclusion

1. Classical autoencoders minimize a reconstruction loss $\|x - \hat{x}\|^2$, they are unregularized in latent space.
2. VAEs are one approach to regularizing (impose structure) the latent distribution.
3. Probabilistic model: KL divergence, approximate inference, reparameterization trick.

Thanks!

Conditional VAE

By using the variational autoencoder, we do not have control over the data generation process.

E.g. we cannot generate only one specific digit from a model trained on MNIST dataset.

We want, for example, to input the character 9 to our model and get a generated image of a handwritten digit 9.

Conditional VAE

We will condition encoder and decoder on other inputs as well as the image, lets call those inputs c .

Encoder becomes: $q(z | x, c)$

Decoder becomes: $p(x | z, c)$

Now our variational lower bound objective becomes:

$$\mathcal{L} = \mathbb{E}[\log p(x | z, c)] - D_{KL}(q(z | x, c) || p(z | c))$$

Conditional VAE

Learned MNIST manifold with a condition of label 2	Learned MNIST manifold with a condition of label 3	Learned MNIST manifold with a condition of label 4
