

DSC190: Machine Learning with Few Labels

Unsupervised Learning

Zhiting Hu

Lecture 17, November 6, 2024

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

Unsupervised learning: Variational Inference

Presentations

- **So Hirota:** Mastering Chess with a Transformer Model
- **Jack Kai Lim:** GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models
- **Bingyan Liu:** REaLTabFormer: Generating Realistic Relational and Tabular Data using Transformers
- **Stephanie Wang:** Toolformer: Language Models Can Teach Themselves to Use Tools
- **Bobby Zhu:** Visualizing Data using t-SNE

Recap: EM and Variational Inference

- The EM algorithm:

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$
 $= p(\mathbf{z}|\mathbf{x}, \theta^t) = \frac{p(\mathbf{z}, \mathbf{x}|\theta^t)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta^t)}$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \\ &= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta)) \end{aligned}$$

Recap: EM and Variational Inference

- The EM algorithm:

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$

Intractable when model $p(\mathbf{z}, \mathbf{x}|\theta)$ is complex

$$= p(\mathbf{z}|\mathbf{x}, \theta^t) = \frac{p(\mathbf{z}, \mathbf{x}|\theta^t)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta^t)}$$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$

Approximate $p(\mathbf{z}|\mathbf{x}, \theta^t)$:

- find a tractable $q(\mathbf{z}|\mathbf{x}, \mathbf{v}^*)$ that is closest to $p(\mathbf{z}|\mathbf{x}, \theta^t)$

$$q(\mathbf{z}|\mathbf{x}, \mathbf{v}^*) = \min_{\mathbf{v}} \text{KL}(q(\mathbf{z}|\mathbf{x}, \mathbf{v}) || p(\mathbf{z}|\mathbf{x}, \theta^t))$$

$$= \min_{\mathbf{v}} F(q(\mathbf{z}|\mathbf{x}, \mathbf{v}), \theta^t) + \text{const.}$$

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

$$= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

Recap: EM and Variational Inference

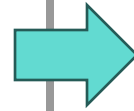
- The EM algorithm:

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t)$

Intractable when model $p(\mathbf{z}, \mathbf{x}|\theta)$ is complex

$$= p(\mathbf{z}|\mathbf{x}, \theta^t) = \frac{p(\mathbf{z}, \mathbf{x}|\theta^t)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta^t)}$$

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$



Approximate $p(\mathbf{z}|\mathbf{x}, \theta^t)$:

- find a tractable $q(\mathbf{z}|\mathbf{x}, \mathbf{v}^*)$ that is closest to $p(\mathbf{z}|\mathbf{x}, \theta^t)$

$$q(\mathbf{z}|\mathbf{x}, \mathbf{v}^*) = \min_{\mathbf{v}} \text{KL}(q(\mathbf{z}|\mathbf{x}, \mathbf{v}) || p(\mathbf{z}|\mathbf{x}, \theta^t))$$

$$= \min_{\mathbf{v}} F(q(\mathbf{z}|\mathbf{x}, \mathbf{v}), \theta^t) + \text{const.}$$

Question: what is the difference?

A: VI assume a simple form of q to ensure tractability

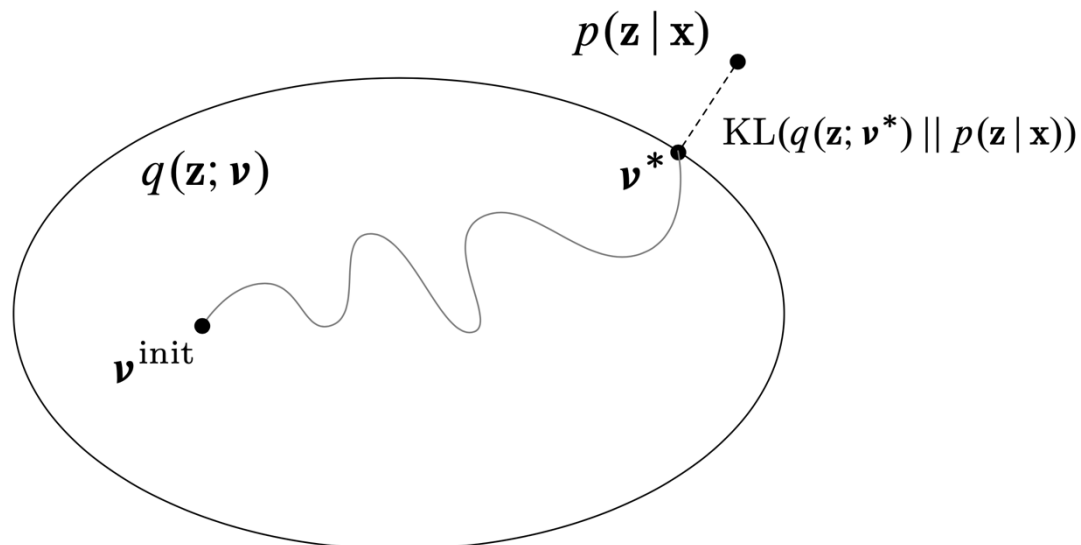
Recap: EM and Variational Inference

- The EM algorithm:

- E-step: $q^{t+1} = \underset{q}{\operatorname{arg\,min}} F(q, \theta^t)$

$$= p(\mathbf{z}|\mathbf{x}, \theta^t) = \frac{p(\mathbf{z}, \mathbf{x}|\theta^t)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta^t)}$$

Intractable when model $p(\mathbf{z}, \mathbf{x}|\theta)$ is complex



Approximate $p(\mathbf{z}|\mathbf{x}, \theta^t)$:

- find a tractable $q(\mathbf{z}|\mathbf{x}, \nu^*)$ that is closest to $p(\mathbf{z}|\mathbf{x}, \theta^t)$

$$q(\mathbf{z}|\mathbf{x}, \nu^*) = \min_{\nu} \operatorname{KL}(q(\mathbf{z}|\mathbf{x}, \nu) || p(\mathbf{z}|\mathbf{x}, \theta^t))$$

$$= \min_{\nu} F(q(\mathbf{z}|\mathbf{x}, \nu), \theta^t) + \text{const.}$$

Question: what is the difference?

A: VI assume a simple form of q to ensure tractability

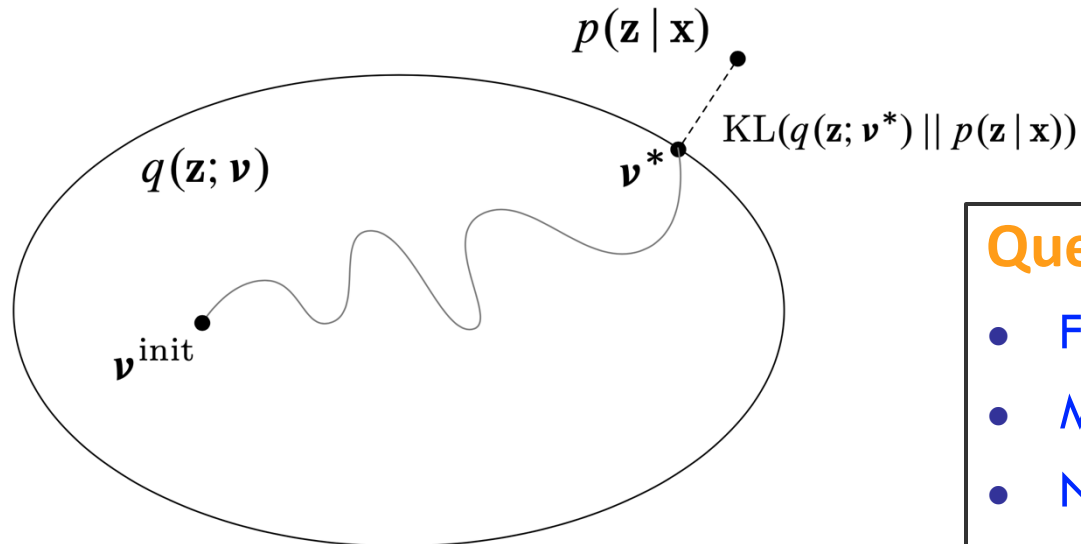
Recap: EM and Variational Inference

- The EM algorithm:

- E-step: $q^{t+1} = \underset{q}{\arg \min} F(q, \theta^t)$

$$= p(\mathbf{z}|\mathbf{x}, \theta^t) = \frac{p(\mathbf{z}, \mathbf{x}|\theta^t)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta^t)}$$

Intractable when model $p(\mathbf{z}, \mathbf{x}|\theta)$ is complex



Approximate $p(\mathbf{z}|\mathbf{x}, \theta^t)$:

- find a tractable $q(\mathbf{z}|\mathbf{x}, \mathbf{v}^*)$ that is closest to $p(\mathbf{z}|\mathbf{x}, \theta^t)$

$$q(\mathbf{z}|\mathbf{x}, \mathbf{v}^*) = \min_{\mathbf{v}} \text{KL}(q(\mathbf{z}|\mathbf{x}, \mathbf{v}) || p(\mathbf{z}|\mathbf{x}, \theta^t))$$

$$= \min_{\mathbf{v}} F(q(\mathbf{z}|\mathbf{x}, \mathbf{v}), \theta^t) + \text{const.}$$

Question: What forms of $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$ shall we choose?

- Factorized distribution -> mean field VI
- Mixture of Gaussian distribution -> black-box VI
- Neural-based distribution -> Variational Autoencoders

Black-box Variational Inference (BBVI)

- Probabilistic model: \mathbf{x} -- observed variables, \mathbf{z} -- latent variables
- Variational distribution $q_{\lambda}(\mathbf{z}|\mathbf{x})$ with parameters λ , e.g.,
 - Gaussian mixture distribution:
 - “A Gaussian mixture model is a **universal approximator** of densities, in the sense that any smooth density can be approximated with any specific nonzero amount of error by a Gaussian mixture model with enough components.” (Deep Learning book, pp.65)
 - Deep neural networks
- ELBO:

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

- Want to compute the gradient w.r.t variational parameters λ

The General Problem: Computing Gradients of Expectations

- When the objective function \mathcal{L} is defined as an expectation of a (differentiable) test function $f_\lambda(\mathbf{z})$ w.r.t. a probability distribution $q_\lambda(\mathbf{z})$

$$\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$$

- Computing exact gradients w.r.t. the parameters λ is often infeasible
- Need stochastic gradient estimates
 - The score function estimator (a.k.a log-derivative trick, REINFORCE)
 - The reparameterization trick (a.k.a the pathwise gradient estimator)

Computing Gradients of Expectations w/ score function

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$
- Log-derivative trick: $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$
- **Question:** show that the gradient of \mathcal{L} w.r.t. λ is:

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

Computing Gradients of Expectations w/ score function

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$
- Log-derivative trick: $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$
- Gradient of \mathcal{L} w.r.t. λ :

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- **score function**: the gradient of the log of a probability distribution
- **Monte Carlo estimation** of the expectation:
 - Compute noisy unbiased gradients with Monte Carlo samples from q_λ

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S f_\lambda(\mathbf{z}_s) \nabla_\lambda \log q_\lambda(\mathbf{z}_s) + \nabla_\lambda f_\lambda(\mathbf{z}_s) \quad \text{where } \mathbf{z}_s \sim q_\lambda(\mathbf{z})$$

Computing Gradients of Expectations w/ score function

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$
- Log-derivative trick: $\nabla_\lambda q_\lambda = q_\lambda \nabla_\lambda \log q_\lambda$
- Gradient of \mathcal{L} w.r.t. λ :

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- **score function**: the gradient of the log of a probability distribution
- **Monte Carlo estimation** of the expectation:
 - Compute noisy unbiased gradients with Monte Carlo samples from q_λ

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S f_\lambda(\mathbf{z}_s) \nabla_\lambda \log q_\lambda(\mathbf{z}_s) + \nabla_\lambda f_\lambda(\mathbf{z}_s) \quad \text{where } \mathbf{z}_s \sim q_\lambda(\mathbf{z})$$

- Pros: generally applicable to any distribution $q(\mathbf{z}|\lambda)$
- Cons: empirically has high variance \rightarrow slow convergence

Computing Gradients of Expectations w/ reparametrization trick

- Loss: $\mathcal{L} = \mathbb{E}_{q_{\lambda}(\mathbf{z})}[f_{\lambda}(\mathbf{z})]$
- Assume that we can express the distribution $q_{\lambda}(\mathbf{z})$ with a transformation

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ \mathbf{z} &= t(\epsilon, \lambda) \end{aligned} \iff \mathbf{z} \sim q(\mathbf{z}|\lambda)$$

- E.g.,

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ \mathbf{z} &= \epsilon\sigma + \mu \end{aligned} \iff \mathbf{z} \sim \text{Normal}(\mu, \sigma^2)$$

- Reparameterization gradient:

$$\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[f_{\lambda}(\mathbf{z}(\epsilon, \lambda))]$$

- **Question:** what's the gradient of \mathcal{L} w.r.t. λ ?

Computing Gradients of Expectations w/ reparametrization trick

- Loss: $\mathcal{L} = \mathbb{E}_{q_{\lambda}(\mathbf{z})}[f_{\lambda}(\mathbf{z})]$
- Assume that we can express the distribution $q_{\lambda}(\mathbf{z})$ with a transformation

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ \mathbf{z} &= t(\epsilon, \lambda) \end{aligned} \iff \mathbf{z} \sim q(\mathbf{z}|\lambda)$$

- E.g.,

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ \mathbf{z} &= \epsilon\sigma + \mu \end{aligned} \iff \mathbf{z} \sim \text{Normal}(\mu, \sigma^2)$$

- Reparameterization gradient:

$$\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[f_{\lambda}(\mathbf{z}(\epsilon, \lambda))]$$

- **Question:** what's the gradient of \mathcal{L} w.r.t. λ ?

$$\nabla_{\lambda} \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_{\lambda}(\mathbf{z}) \nabla_{\lambda} t(\epsilon, \lambda)]$$

Computing Gradients of Expectations w/ reparametrization trick

- Loss: $\mathcal{L} = \mathbb{E}_{q_{\lambda}(\mathbf{z})}[f_{\lambda}(\mathbf{z})]$
- Assume that we can express the distribution $q_{\lambda}(\mathbf{z})$ with a transformation

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ \mathbf{z} &= t(\epsilon, \lambda) \end{aligned} \iff \mathbf{z} \sim q(\mathbf{z}|\lambda)$$

- E.g.,

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ \mathbf{z} &= \epsilon\sigma + \mu \end{aligned} \iff \mathbf{z} \sim \text{Normal}(\mu, \sigma^2)$$

- Reparameterization gradient

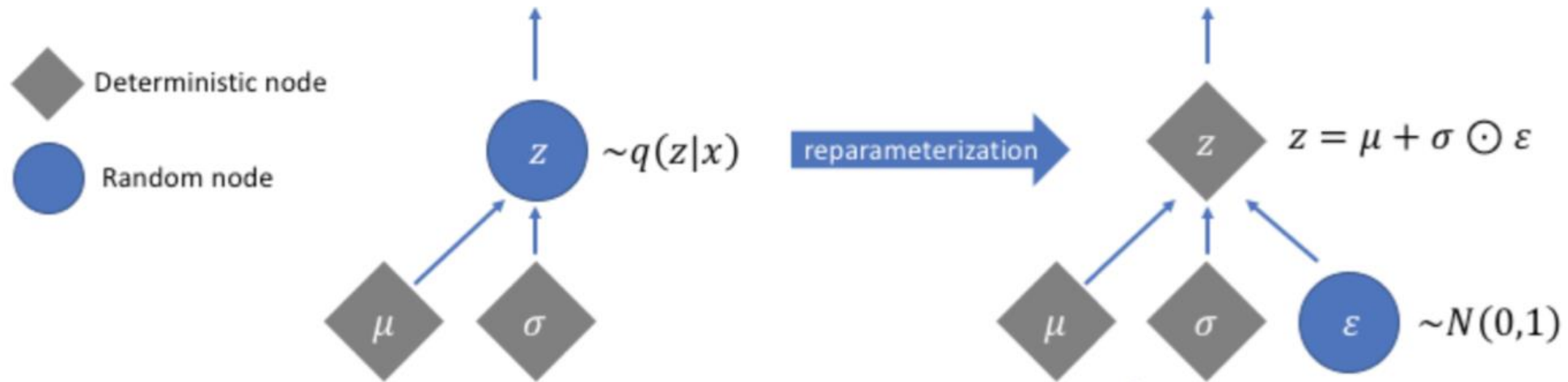
$$\begin{aligned} \mathcal{L} &= \mathbb{E}_{\epsilon \sim s(\epsilon)}[f_{\lambda}(\mathbf{z}(\epsilon, \lambda))] \\ \nabla_{\lambda} \mathcal{L} &= \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_{\lambda}(\mathbf{z}) \nabla_{\lambda} t(\epsilon, \lambda)] \end{aligned}$$

- Pros: empirically, lower variance of the gradient estimate
- Cons: Not all distributions can be reparameterized

Reparameterization trick

- Reparameterizing Gaussian distribution

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ z &= \epsilon\sigma + \mu \end{aligned} \iff z \sim \text{Normal}(\mu, \sigma^2)$$



Reparameterization trick

- Reparametrizing Gaussian distribution

$$\begin{aligned} \epsilon &\sim \text{Normal}(0, 1) \\ z &= \epsilon\sigma + \mu \end{aligned} \iff z \sim \text{Normal}(\mu, \sigma^2)$$

- Other reparameterizable distributions: $\epsilon \sim \text{Uniform}(\epsilon) \iff z \sim q(z)$
 - Tractable inverse CDF F^{-1} :
 - Exponential, Cauchy, Logistic, Rayleigh, Pareto, Weibull, Reciprocal, Gompertz, Gumbel, Erlang
 - Location-scale:
 - Laplace, Elliptical, Student's t, Logistic, Uniform, Triangular, Gaussian
 - Composition:
 - Log-Normal (exponentiated normal) Gamma (sum of exponentials) Dirichlet (sum of Gammas) Beta, Chi-Squared, F

Computing Gradients of Expectations: Summary

- Loss: $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$

- **Score gradient**

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- Pros: generally applicable to any distribution $q(\mathbf{z}|\lambda)$
- Cons: empirically has high variance \rightarrow slow convergence

- **Reparameterization gradient**

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_\lambda(\mathbf{z}) \nabla_\lambda t(\epsilon, \lambda)]$$

- Pros: empirically, lower variance of the gradient estimate
- Cons: Not all distributions can be reparameterized

Recall: Black-box Variational Inference (BBVI)

- Probabilistic model: \mathbf{x} -- observed variables, \mathbf{z} -- latent variables
- Variational distribution $q_\lambda(\mathbf{z}|\mathbf{x})$ with parameters λ , e.g.,
 - Gaussian mixture distribution:
 - “A Gaussian mixture model is a universal approximator of densities, in the sense that any smooth density can be approximated with any specific nonzero amount of error by a Gaussian mixture model with enough components.” (Deep Learning book, pp.65)
 - Deep neural networks

- ELBO to be maximized:
$$\mathcal{L}(\lambda) \triangleq \mathbb{E}_{q_\lambda(\mathbf{z})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})]$$

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

- Want to compute the gradient w.r.t variational parameters λ

BBVI with the score gradient

$$\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$$

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- ELBO:

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

- **Question:** what's the score gradient w.r.t. λ ?

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_q[\nabla_\lambda \log q(\mathbf{z}|\lambda)(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\lambda))]$$

BBVI with the score gradient

$$\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$$

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z}) \nabla_\lambda \log q_\lambda(\mathbf{z}) + \nabla_\lambda f_\lambda(\mathbf{z})]$$

- ELBO:

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

- **Question:** what's the score gradient w.r.t. λ ?

$$\nabla_\lambda \mathcal{L} = \mathbb{E}_q[\nabla_\lambda \log q(\mathbf{z}|\lambda)(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\lambda))]$$

- Compute noisy unbiased gradients of the ELBO with Monte Carlo samples from the variational distribution

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S \nabla_\lambda \log q(z_s|\lambda)(\log p(\mathbf{x}, z_s) - \log q(z_s|\lambda)),$$

where $z_s \sim q(\mathbf{z}|\lambda)$.

BBVI with the reparameterization gradient

$$\mathcal{L} = \mathbb{E}_{q_{\lambda}(\mathbf{z})}[f_{\lambda}(\mathbf{z})]$$
$$\nabla_{\lambda}\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}}f_{\lambda}(\mathbf{z}) \nabla_{\lambda}t(\epsilon, \lambda)]$$

- ELBO:

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

- **Question:** what's the reparameterization gradient w.r.t. λ ?

$$\begin{array}{l} \epsilon \sim s(\epsilon) \\ \mathbf{z} = t(\epsilon, \lambda) \end{array} \iff \mathbf{z} \sim q(\mathbf{z}|\lambda)$$

$$\nabla_{\lambda}\mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \nabla_{\lambda}t(\epsilon, \lambda)]$$

Variational Autoencoders (VAEs)

Variational Auto-Encoders (VAEs)

VAEs are a combination of the following ideas:

- Variational Inference
 - ELBO
- Variational distribution parametrized as neural networks
- Reparameterization trick

Variational Auto-Encoders (VAEs)

- Model $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
 - $p_{\theta}(\mathbf{x}|\mathbf{z})$: a.k.a., generative model, generator, (probabilistic) decoder, ...
 - $p(\mathbf{z})$: prior, e.g., Gaussian
- Assume variational distribution $q_{\phi}(\mathbf{z}|\mathbf{x})$
 - E.g., a Gaussian distribution parameterized as **deep neural networks**
 - a.k.a, recognition model, inference network, (probabilistic) encoder, ...
- ELBO:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$

Reconstruction

Divergence from prior
(KL divergence between two Gaussians has
an analytic form)

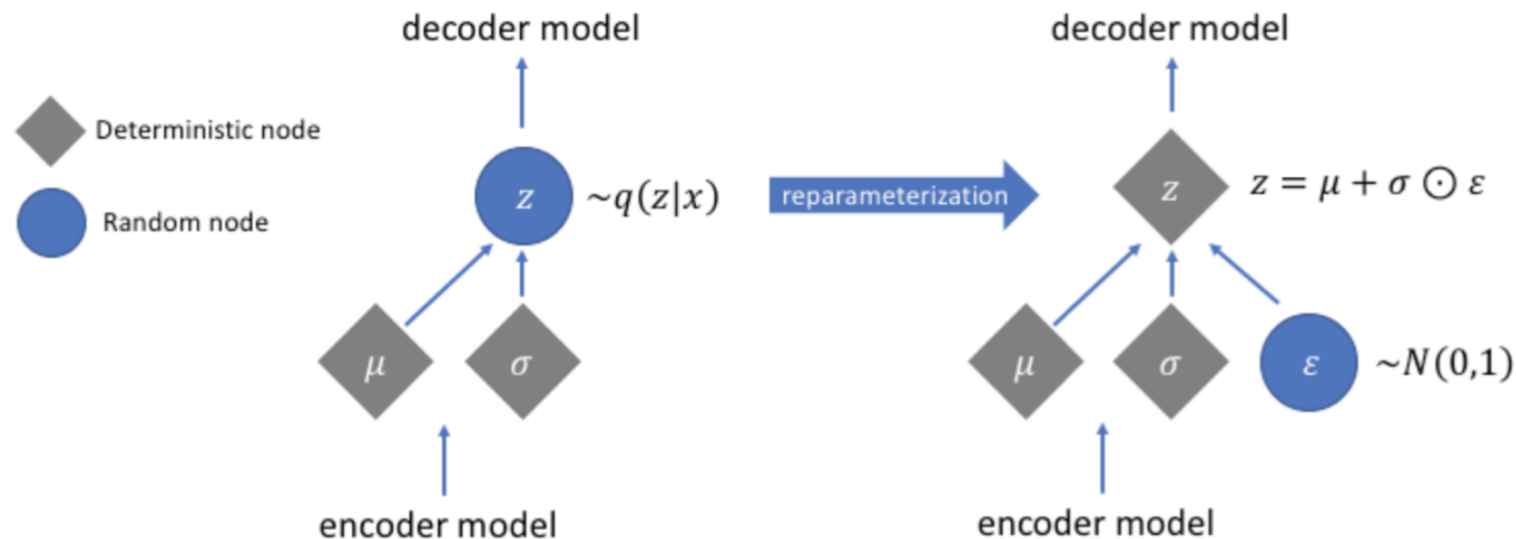
Variational Auto-Encoders (VAEs)

- ELBO:

$$\begin{aligned}\mathcal{L}(\theta, \phi; \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$

- Reparameterization:

- $[\mu; \sigma] = f_{\phi}(\mathbf{x})$ (a neural network)
- $\mathbf{z} = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(\mathbf{0}, \mathbf{1})$



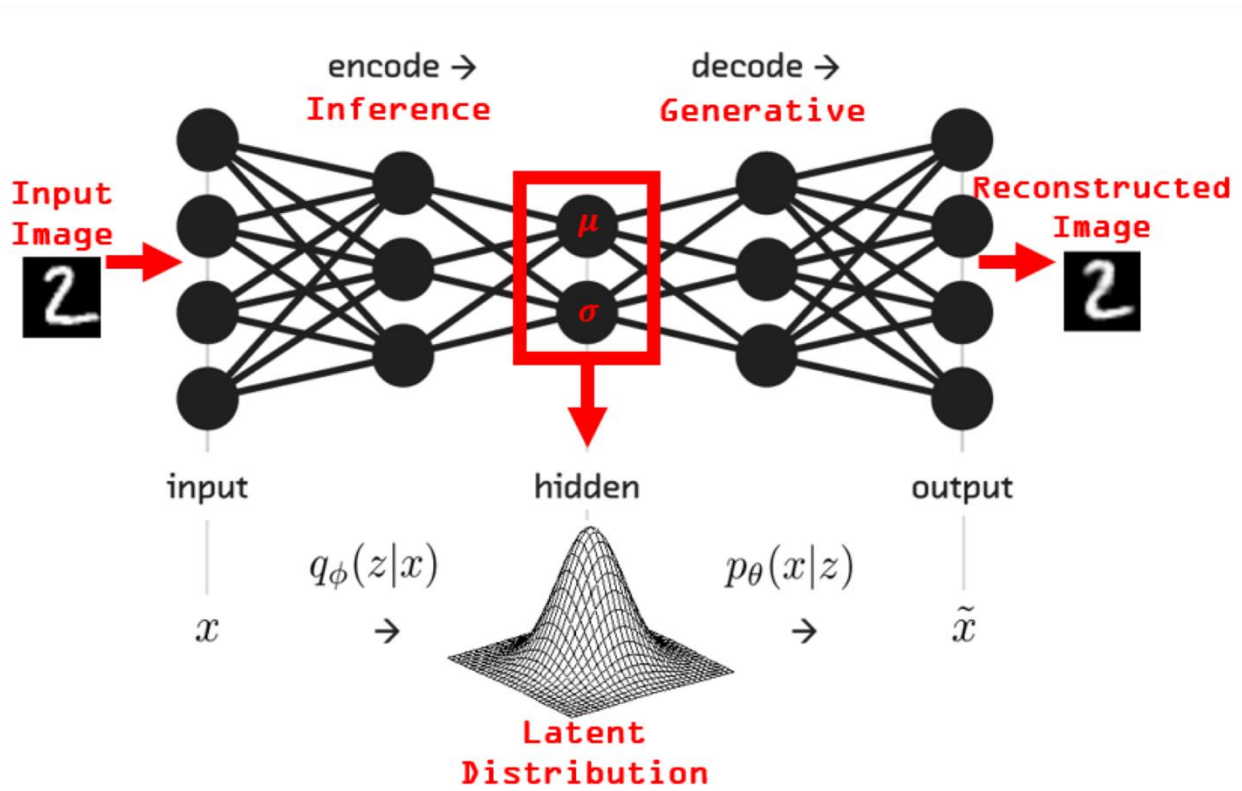
Variational Auto-Encoders (VAEs)

- ELBO:
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] + H(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$
- Reparameterization:
 - $[\boldsymbol{\mu}; \boldsymbol{\sigma}] = f_{\boldsymbol{\phi}}(\mathbf{x})$ (a neural network)
 - $\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})$

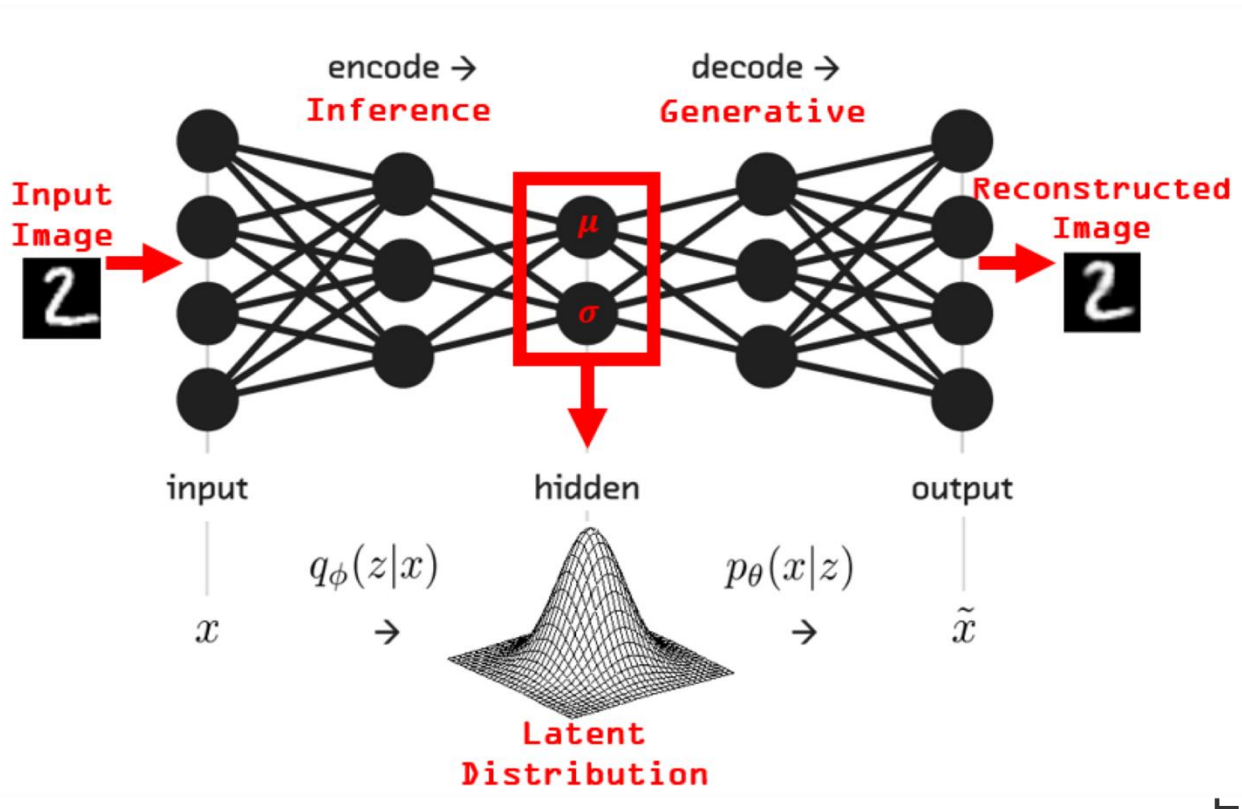
$$\nabla_{\boldsymbol{\phi}} \mathcal{L} = \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})} [\nabla_{\mathbf{z}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})] \nabla_{\boldsymbol{\phi}} \mathbf{z}(\boldsymbol{\epsilon}, \boldsymbol{\phi})]$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

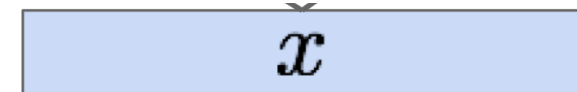
Example: VAEs for images



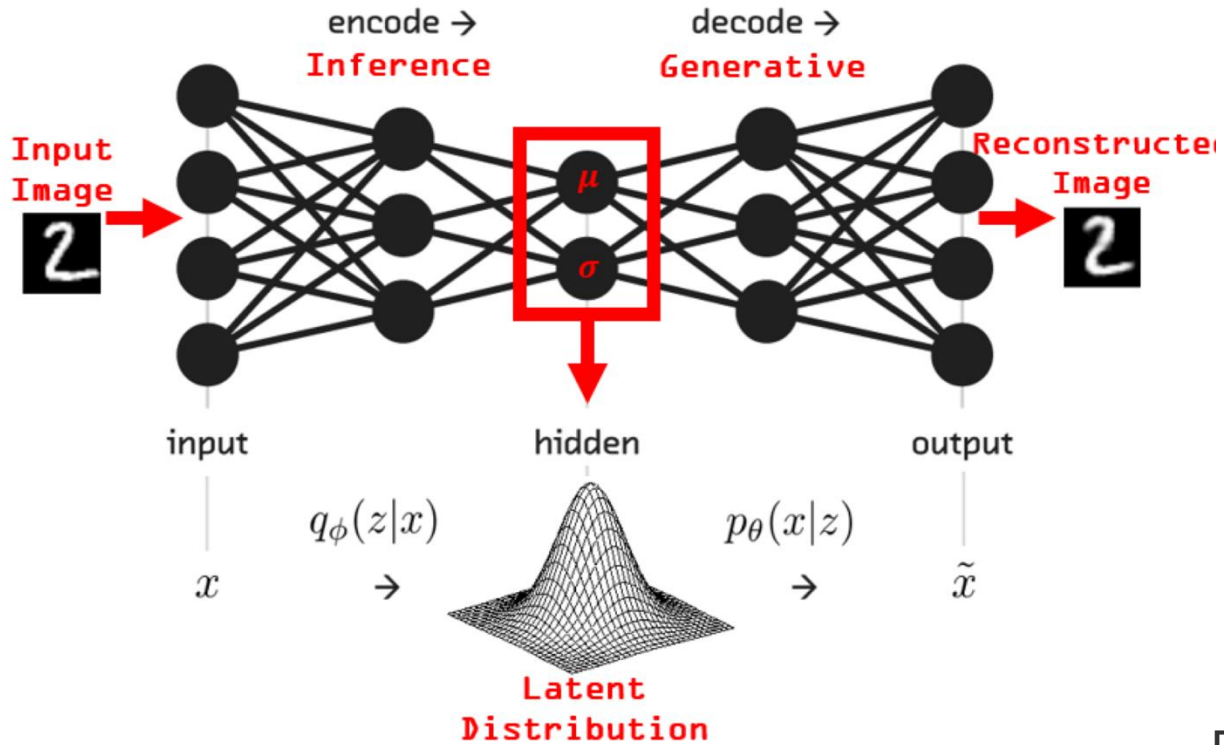
Example: VAEs for images



Input Data



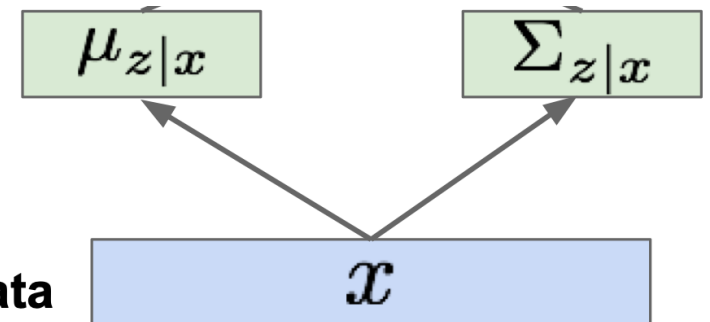
Example: VAEs for images



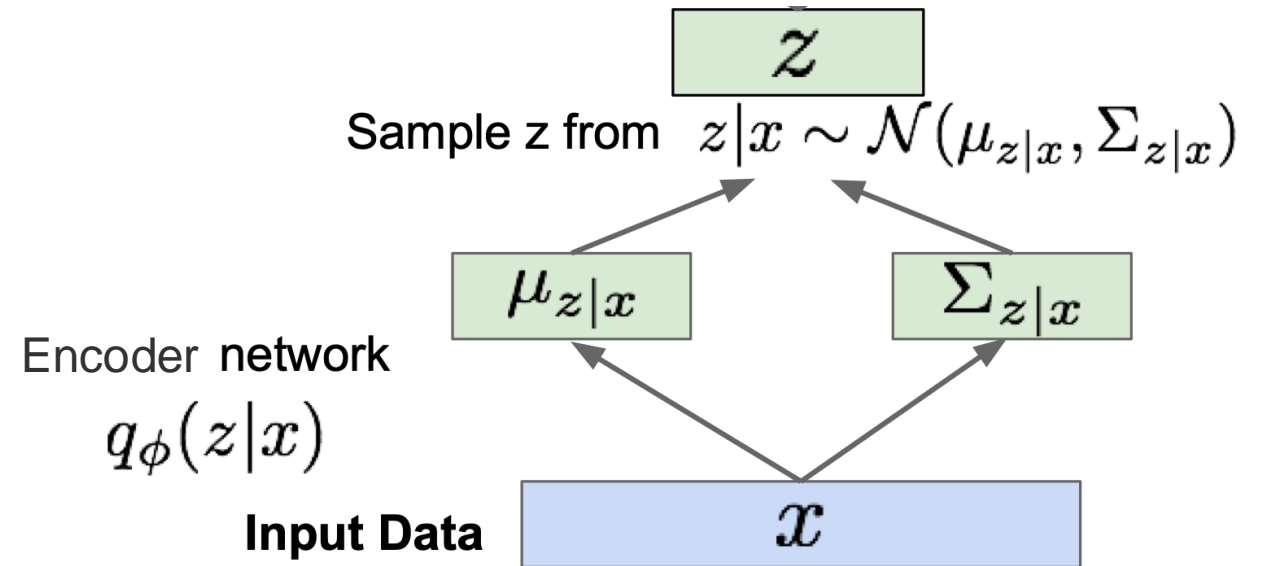
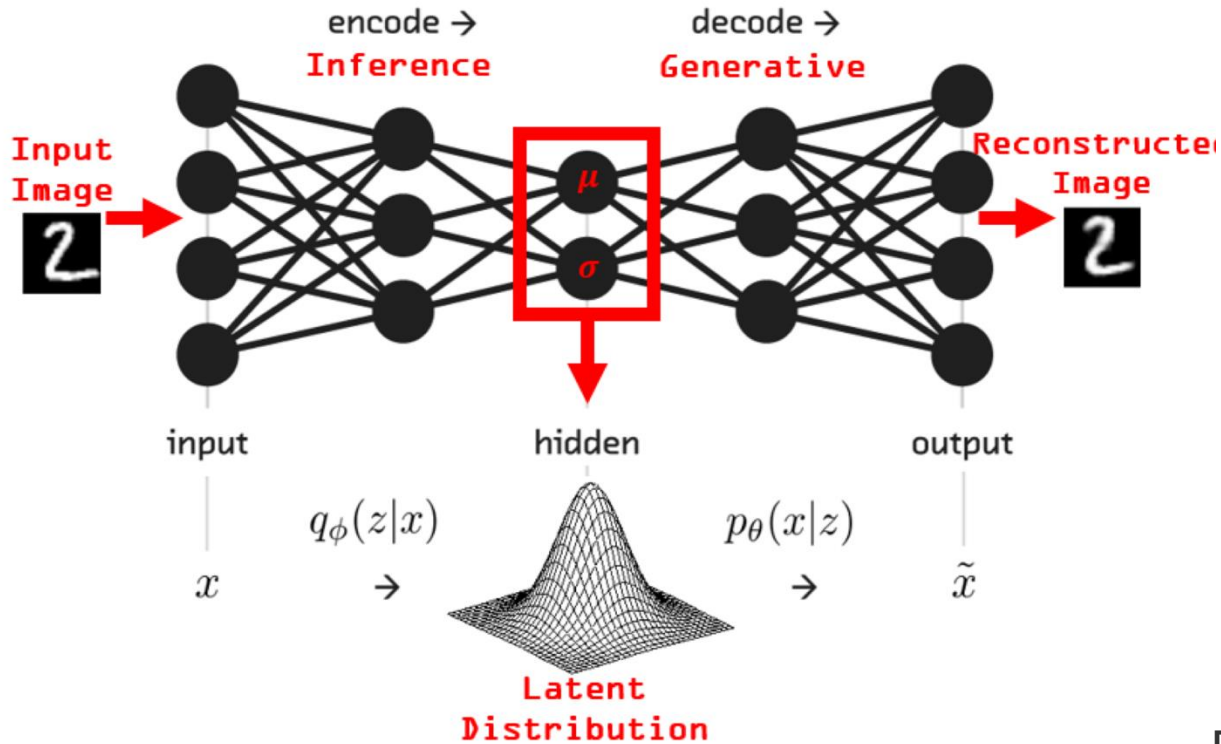
Encoder network

$$q_\phi(z|x)$$

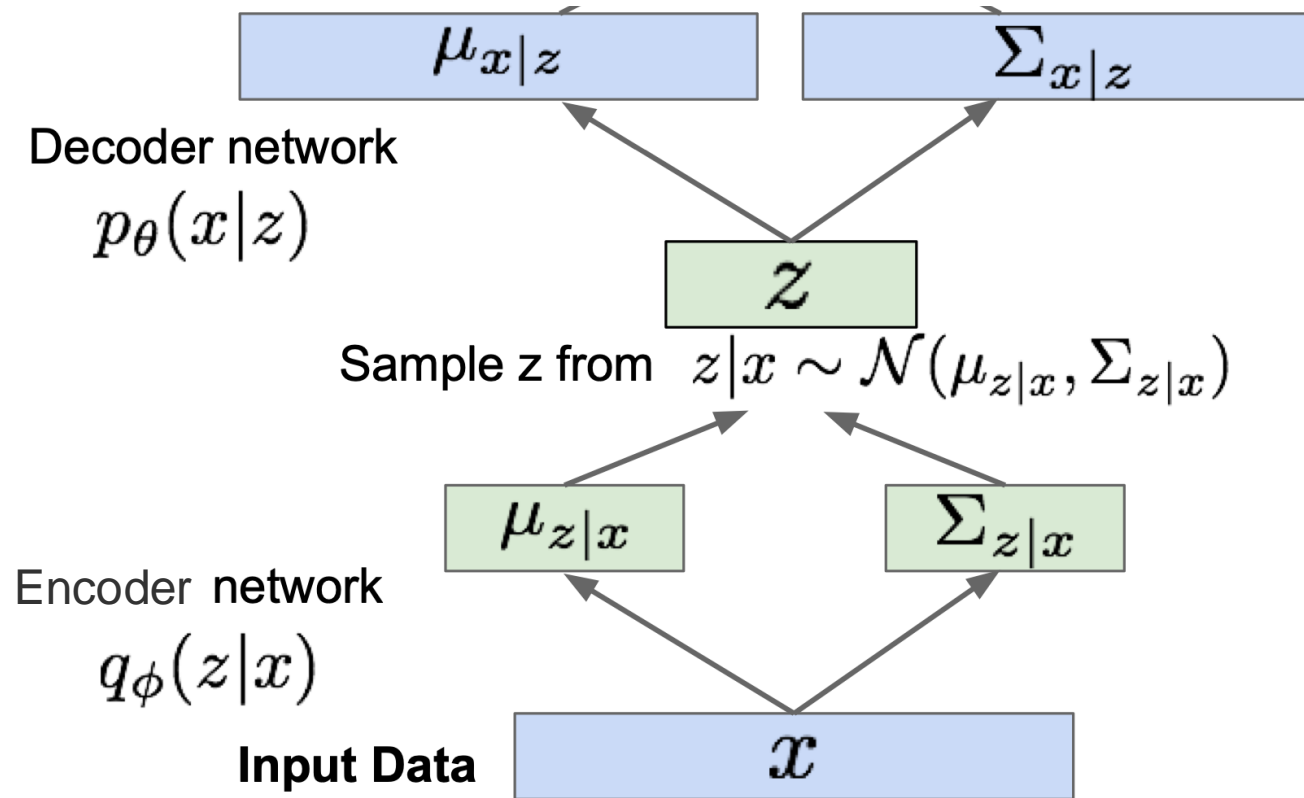
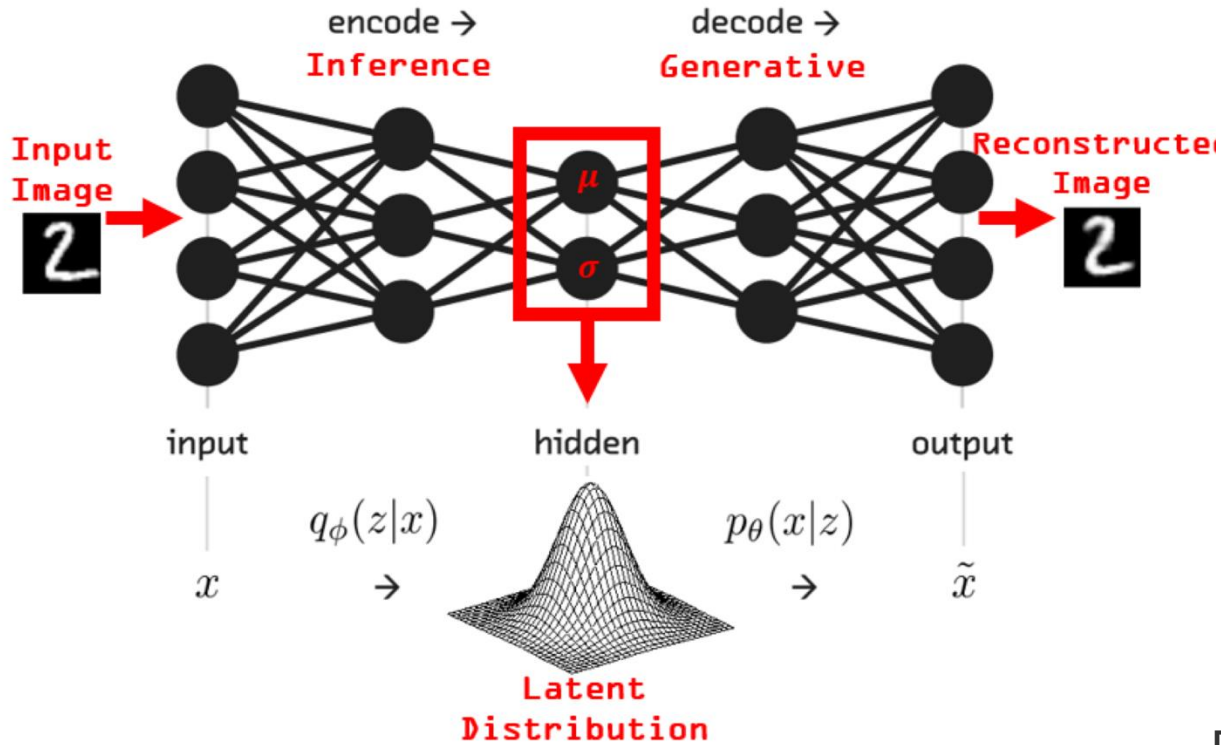
Input Data



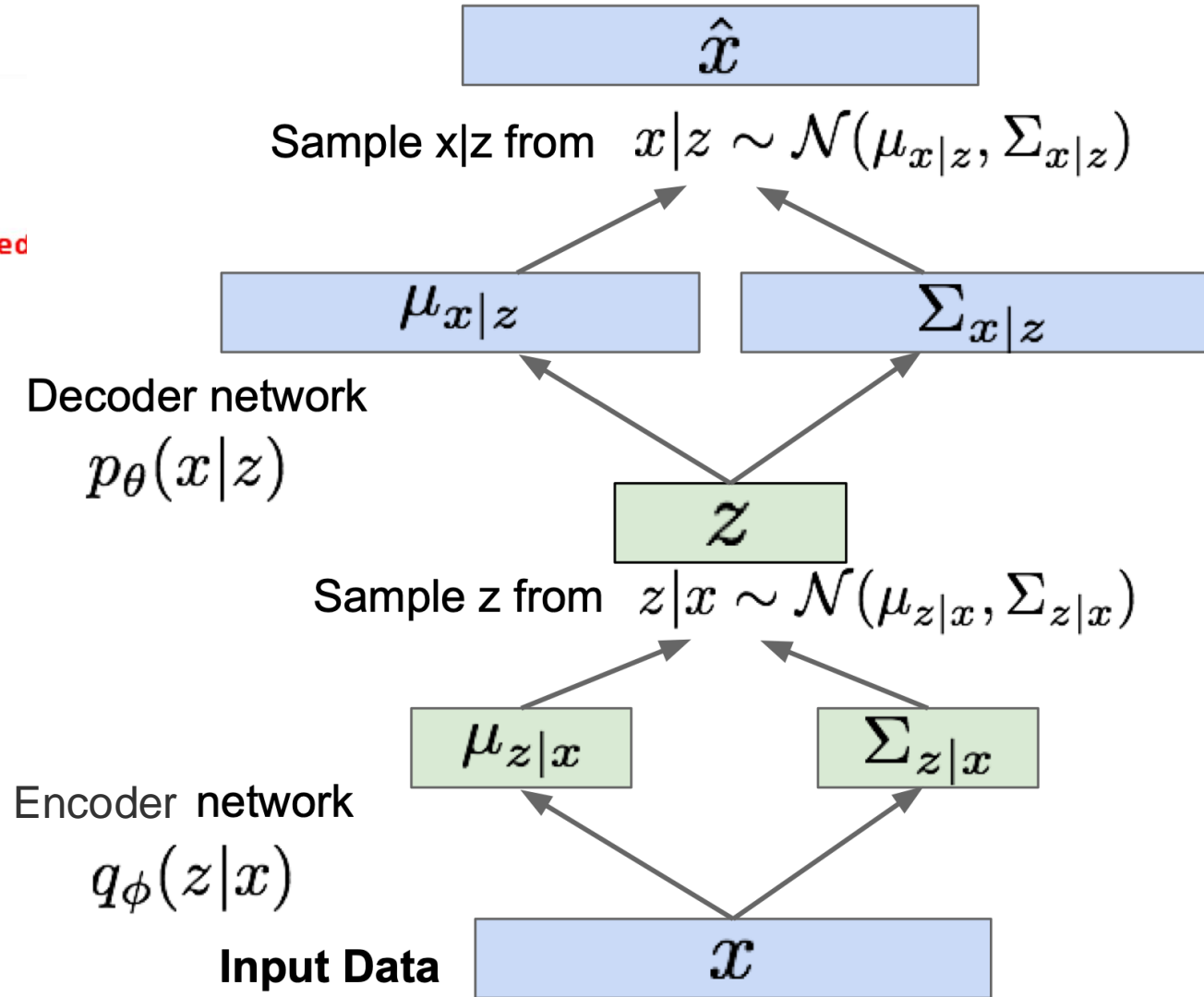
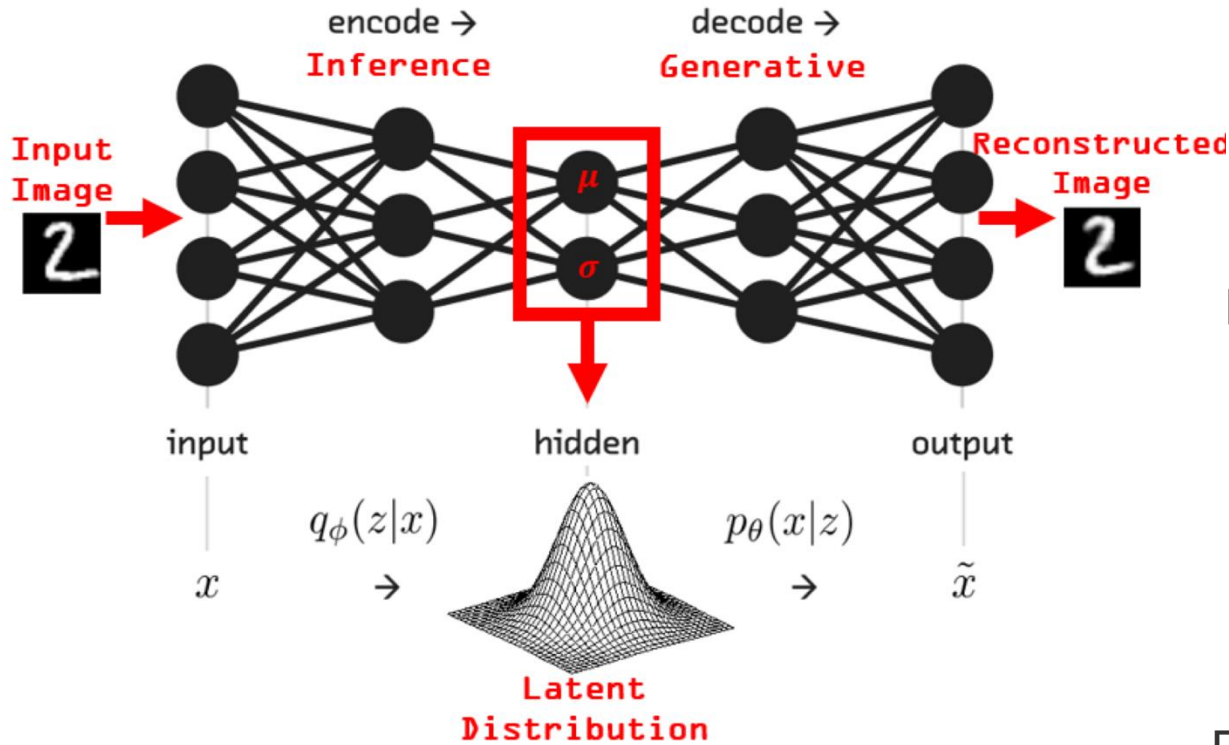
Example: VAEs for images



Example: VAEs for images



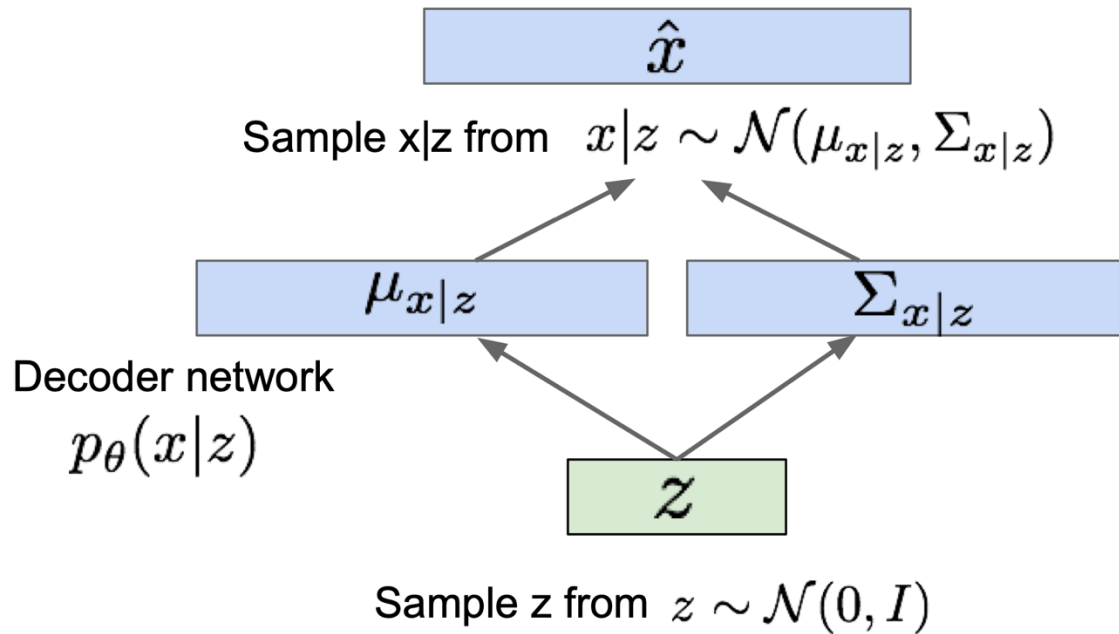
Example: VAEs for images



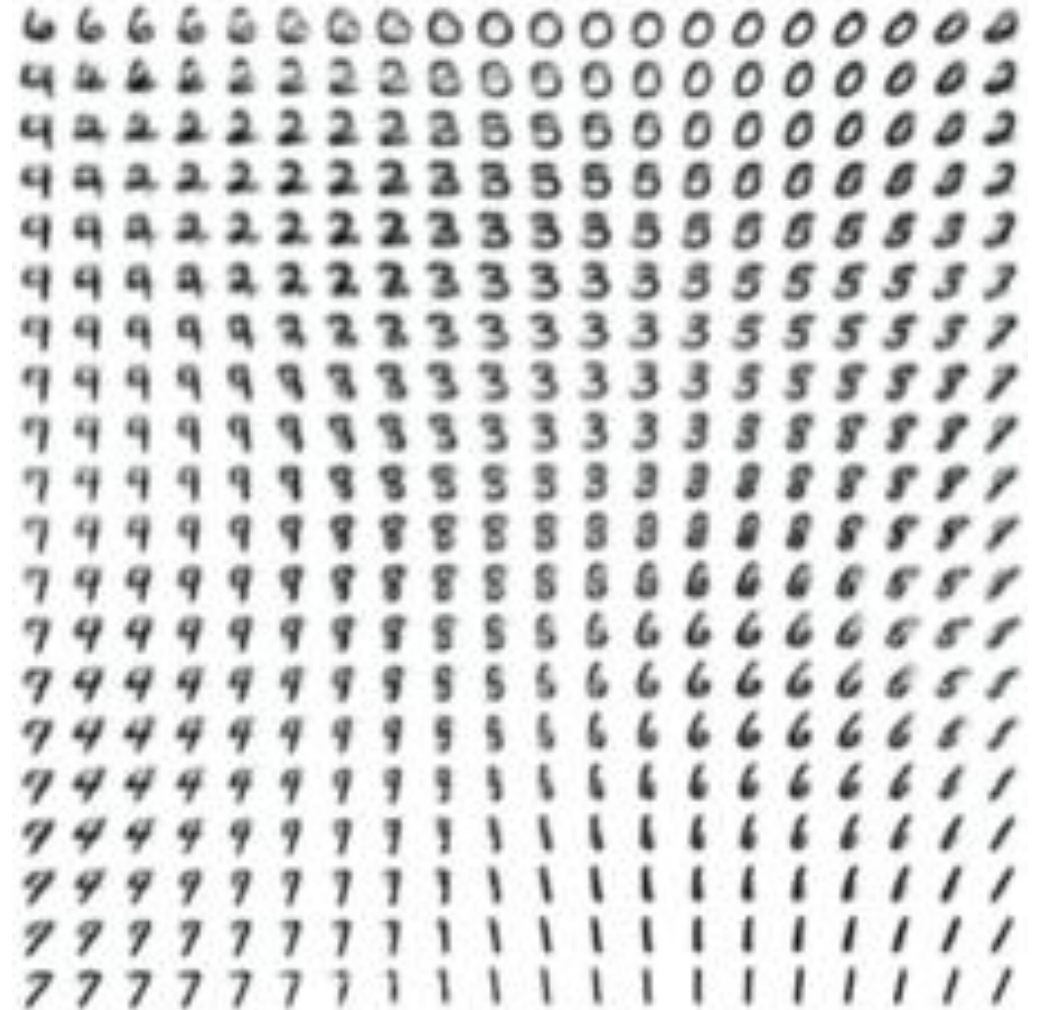
Example: VAEs for images

Generating samples:

- Use decoder network. Now sample z from prior!



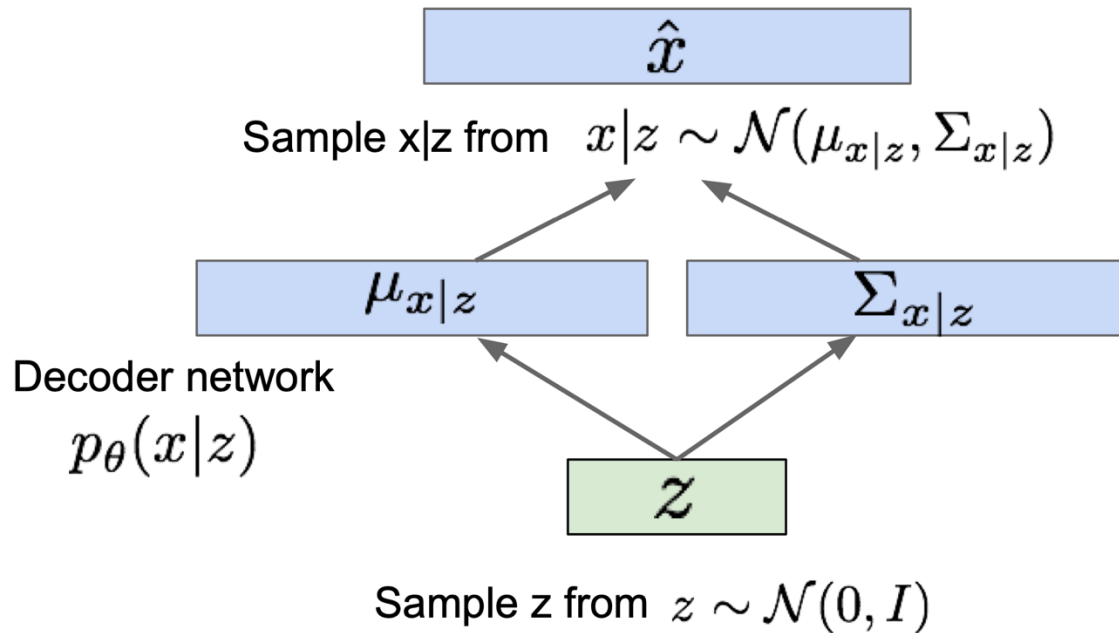
Data manifold for 2-d z



Example: VAEs for images

Generating samples:

- Use decoder network. Now sample z from prior!



Vary z_1
(Degree of smile)

Data manifold for 2-d z



Vary z_2 (head pose)

Example: VAEs for text

- Latent code interpolation and sentences generation from VAEs [Bowman et al., 2015].

“ i want to talk to you . ”

“i want to be with you . ”

“i do n’t want to be with you . ”

i do n’t want to be with you .

she did n’t want to be with him .

Note: Amortized Variational Inference

- Variational distribution as an **inference model** $q_{\phi}(\mathbf{z}|\mathbf{x})$ with parameters ϕ (which was traditionally factored over samples)
- Amortize the cost of inference by learning a **single** data-dependent inference model
- The trained inference model can be used for quick inference on new data

Variational Auto-encoders: Summary

- A combination of the following ideas:
 - Variational Inference: ELBO
 - Variational distribution parametrized as neural networks
 - Reparameterization trick

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

←
Reconstruction

↓
Divergence from prior

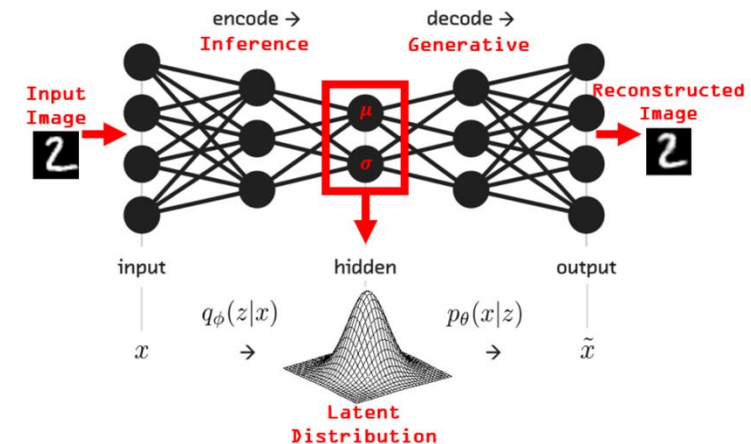
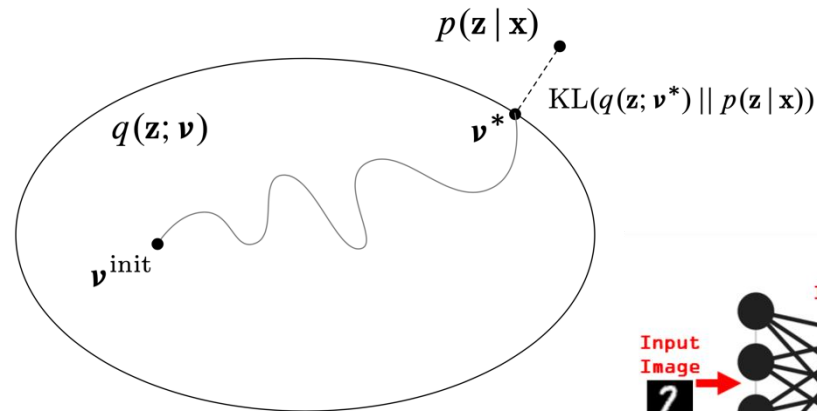


(Razavi et al., 2019)

- Pros:
 - Principled approach to generative models
 - Allows inference of $q(\mathbf{z}|\mathbf{x})$, can be useful feature representation for other tasks
- Cons:
 - Samples blurrier and lower quality compared to GANs
 - Tend to collapse on text data

Summary: Supervised / Unsupervised Learning

- Supervised Learning
 - Maximum likelihood estimation (MLE)
- Unsupervised learning
 - Maximum likelihood estimation (MLE) with latent variables
 - Marginal log-likelihood
 - EM algorithm for MLE
 - ELBO / Variational free energy
 - Variational Inference
 - ELBO / Variational free energy
 - Variational distributions
 - Factorized (mean-field VI)
 - Mixture of Gaussians (Black-box VI)
 - Neural-based (VAEs)



Presentations

Questions?