

DSC190: Machine Learning with Few Labels

Unsupervised Learning

Zhitong Hu

Lecture 16, November 4, 2024

Outline

Unsupervised learning: Variational Inference

Presentations

- **Peiyuan Sun:** Reasoning with Language Model is Planning with World Model
- **Mingyang Yao:** Pop Music Transformer: Beat-based Modeling and Generation of Expressive Pop Piano Compositions
- **Zhaoxiang Feng:** Learning Equilibria in Matching Markets from Bandit Feedback
- **Bella Wang:** Language Models Are Realistic Tabular Data Generators

Recap: Variational Inference

Recall that in EM, we assume $q(z|x)$ can be any distribution. E-step shows the optimal $q(z|x)$ is the posterior distribution.

The main idea behind variational inference:

- Choose a family of distributions over the latent variables $z_{1:m}$ with its own set of variational parameters ν , i.e.

$$q(z_{1:m}|\nu)$$

- Then, we find the setting of the parameters that makes our approximation q closest to the posterior distribution.
 - This is where optimization algorithms come in.
- Then we can use q with the fitted parameters in place of the posterior.
 - E.g. to form predictions about future data, or to investigate the posterior distribution over the hidden variables, find modes, etc.

Recap: Variational Inference

- We want to minimize the KL divergence between our approximation $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})$ and our posterior $p(\mathbf{z}|\mathbf{x})$

$$\text{KL}(q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu}) \parallel p(\mathbf{z}|\mathbf{x}))$$

- But we can't actually minimize this quantity w.r.t q because $p(\mathbf{z}|\mathbf{x})$ is unknown
- **Question:** how can we minimize the KL divergence?

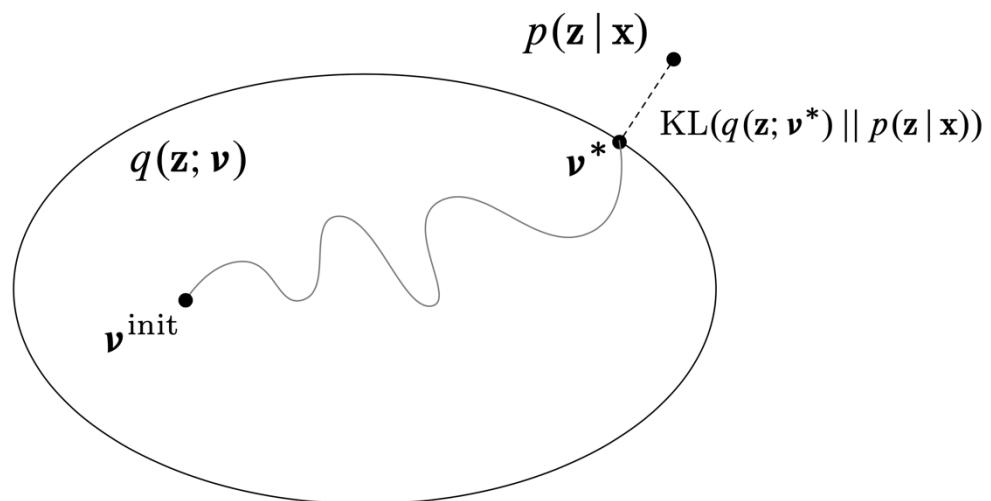
$$\ell(\theta; \mathbf{x}) = \underbrace{\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]}_{\text{Evidence Lower Bound (ELBO)}} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

- The ELBO is equal to the negative KL divergence up to a constant $\ell(\theta; \mathbf{x})$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters ν , i.e. $q(\mathbf{z}|\mathbf{x}, \nu)$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

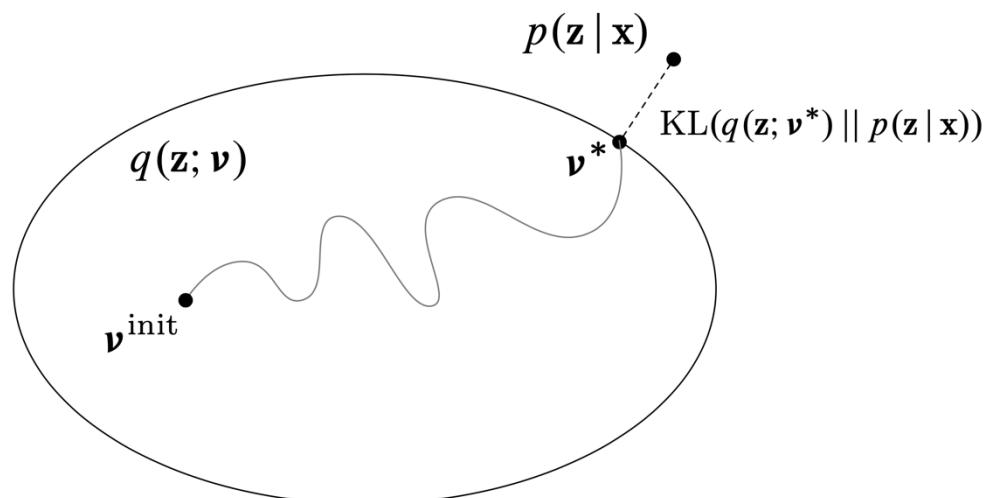
$$\begin{aligned} & \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \nu)} \right] \\ &= \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log q(\mathbf{z}|\mathbf{x}, \nu)] \end{aligned}$$



Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters ν , i.e. $q(\mathbf{z}|\mathbf{x}, \nu)$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \nu)} \right] \\ &= \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log q(\mathbf{z}|\mathbf{x}, \nu)] \end{aligned}$$

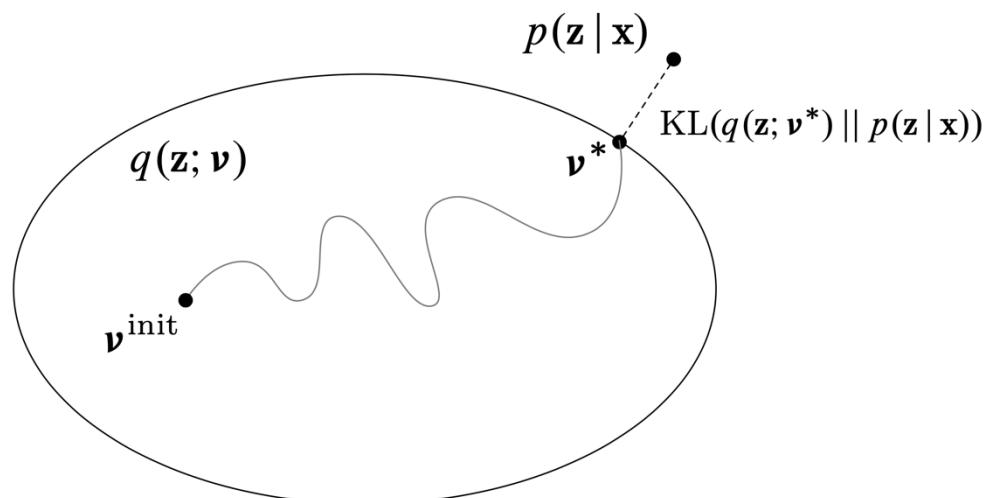


Question: How do we choose the variational family $q(\mathbf{z}|\mathbf{x}, \nu)$?

Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters ν , i.e. $q(\mathbf{z}|\mathbf{x}, \nu)$
- We maximize the ELBO over q to find an “optimal approximation” to $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \nu)} \right] \\ &= \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log q(\mathbf{z}|\mathbf{x}, \nu)] \end{aligned}$$



Question: How do we choose the variational family $q(\mathbf{z}|\mathbf{x}, \nu)$?

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

Example: Mean Field Variational Inference

- A popular family of variational approximations
- In this type of variational inference, we assume the variational distribution over the latent variables **factorizes** as

$$q(\mathbf{z}) = q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

- (where we omit variational parameters for ease of notation)
- We refer to $q(z_j)$, the variational approximation for a single latent variable, as a “local variational approximation”
- In the above expression, the variational approximation $q(z_j)$ over each latent variable z_j is independent

Mean Field Variational Inference with Coordinate Ascent

Recap: Bayesian mixture of Gaussians

Assume mean-field $q(\mu_{1:K}, z_{1:n}) = \prod_k q(\mu_k) \prod_i q(z_i)$

- Initialize the global variational distributions $q(\mu_k)$ and parameters $\{\tau^2, \sigma^2, \pi\}$
 - **Repeat:**
 - **For** each data example $i \in \{1, 2, \dots, D\}$
 - Update the local variational distribution $q(z_i)$
 - **End for**
 - Update the global variational distributions $q(\mu_k)$
 - Update the parameters $\{\tau^2, \sigma^2, \pi\}$
 - **Until** ELBO converges
-
- The diagram illustrates the iterative process of Mean Field Variational Inference. It shows a sequence of steps: 1. Initialize global variational distributions and parameters. 2. Repeat the following steps: a. For each data example i , update the local variational distribution $q(z_i)$. b. End for. c. Update the global variational distributions $q(\mu_k)$. d. Update the parameters $\{\tau^2, \sigma^2, \pi\}$. A red bracket groups the first two steps (the E-step), and a red arrow points to the last two steps (the M-step).
- What if we have millions of data examples? This could be very slow.

Stochastic VI

Recap: Bayesian mixture of Gaussians

Assume mean-field $q(\mu_{1:K}, z_{1:n}) = \prod_k q(\mu_k) \prod_i q(z_i)$

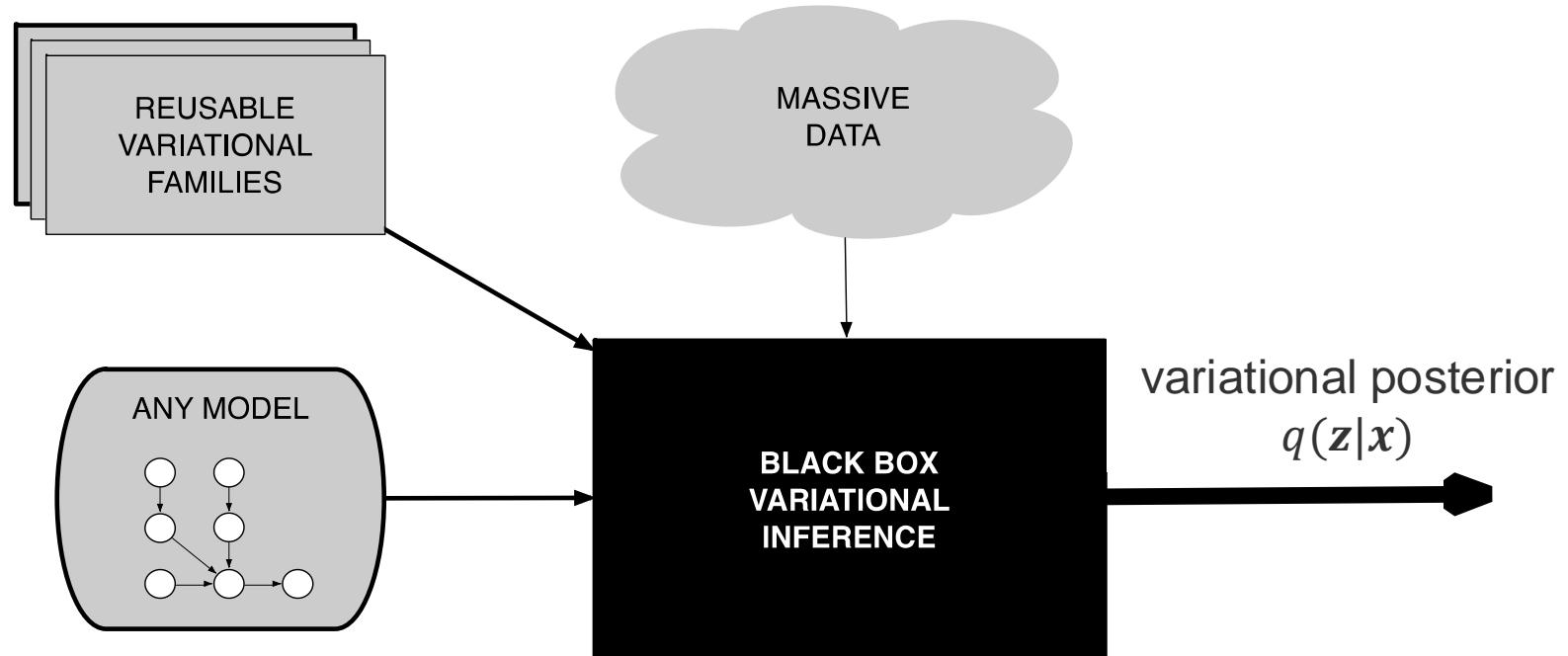
- Initialize the global variational distributions $q(\mu_k)$ and parameters $\{\tau^2, \sigma^2, \pi\}$
- **Repeat:**
 - **Sample** a data example $i \in \{1, 2, \dots, D\}$
 - Update the local variational distribution $q(z_i)$
 - Update the global variational distributions $q(\mu_k)$ with **natural gradient ascent**
 - Update the parameters $\{\tau^2, \sigma^2, \pi\}$
- **Until** ELBO converges

Black-box Variational Inference

Black-box Variational Inference (BBVI)

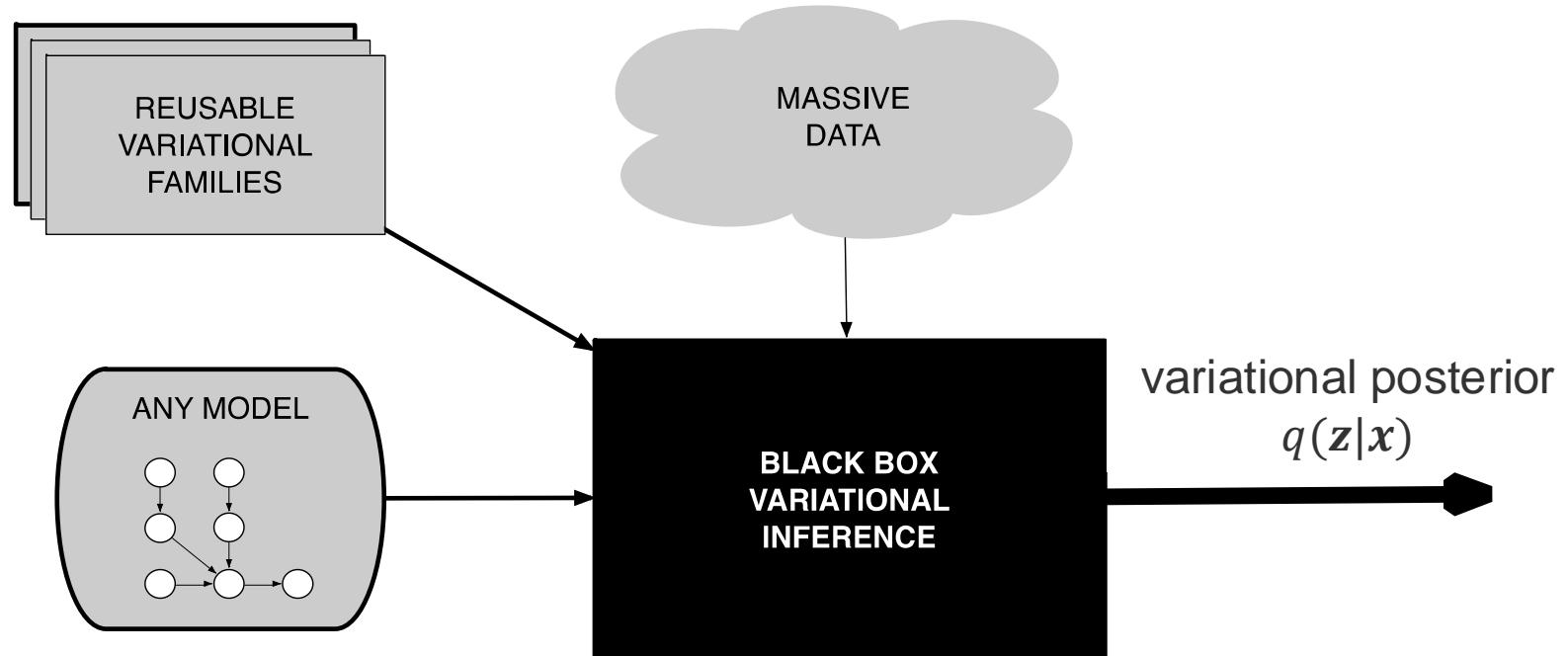
- We have derived variational inference specific for Bayesian Gaussian (mixture) models
- There are innumerable models
- Can we have a solution that does not entail model-specific work?

Black-box Variational Inference (BBVI)



- Easily use variational inference with **any model**
- Perform inference with **massive data**
- **No mathematical work** beyond specifying the model

Black-box Variational Inference (BBVI)



- Sample from $q(\cdot)$
- Form noisy gradients (without model-specific computation)
- Use stochastic optimization

Black-box Variational Inference (BBVI)

- Probabilistic model: x -- observed variables, z -- latent variables
- Variational distribution $q_\lambda(z|x)$ 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:

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

- 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 :

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

- **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 :

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

- **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

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

- E.g.,

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

- 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

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

- E.g.,
$$\begin{aligned}\epsilon &\sim \text{Normal}(0, 1) \\ z &= \epsilon\sigma + \mu\end{aligned}\quad \Leftrightarrow \quad z \sim \text{Normal}(\mu, \sigma^2)$$

- 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 ?

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

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

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

- E.g.,

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

- Reparameterization gradient

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

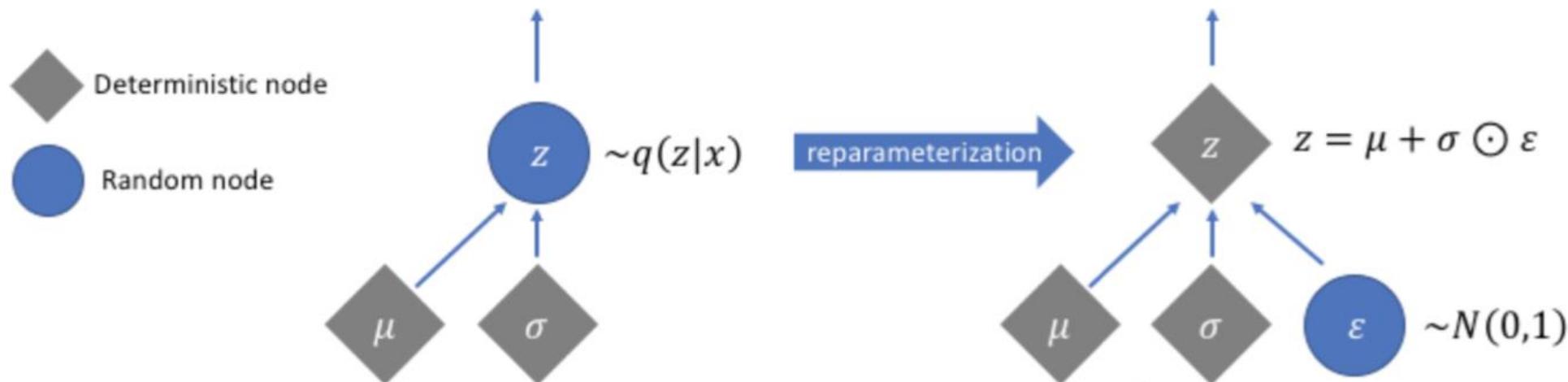
$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_z f_\lambda(\mathbf{z}) \nabla_\lambda t(\epsilon, \lambda)]$$

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

Reparameterization trick

- Reparametrizing Gaussian distribution

$$\epsilon \sim \text{Normal}(0, 1) \Leftrightarrow z \sim \text{Normal}(\mu, \sigma^2)$$
$$z = \epsilon\sigma + \mu$$



Reparameterization trick

- Reparametrizing Gaussian distribution

$$\epsilon \sim \text{Normal}(0, 1) \Leftrightarrow z \sim \text{Normal}(\mu, \sigma^2)$$
$$z = \epsilon\sigma + \mu$$

- Other reparameterizable distributions:

- 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

$$\epsilon \sim \text{Uniform}(\epsilon) \Leftrightarrow z \sim q(z)$$
$$z = F^{-1}(\epsilon)$$

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 $q(z|\lambda)$
 - 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: x -- observed variables, z -- latent variables
- Variational distribution $q_\lambda(z|x)$ 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:
$$\mathcal{L}(\lambda) \triangleq \mathbb{E}_{q_\lambda(z)}[\log p(x, z) - \log q(z)]$$

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

- 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:

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

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

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

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:

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

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

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

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

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

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

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:

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

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

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

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

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

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

Reconstruction



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

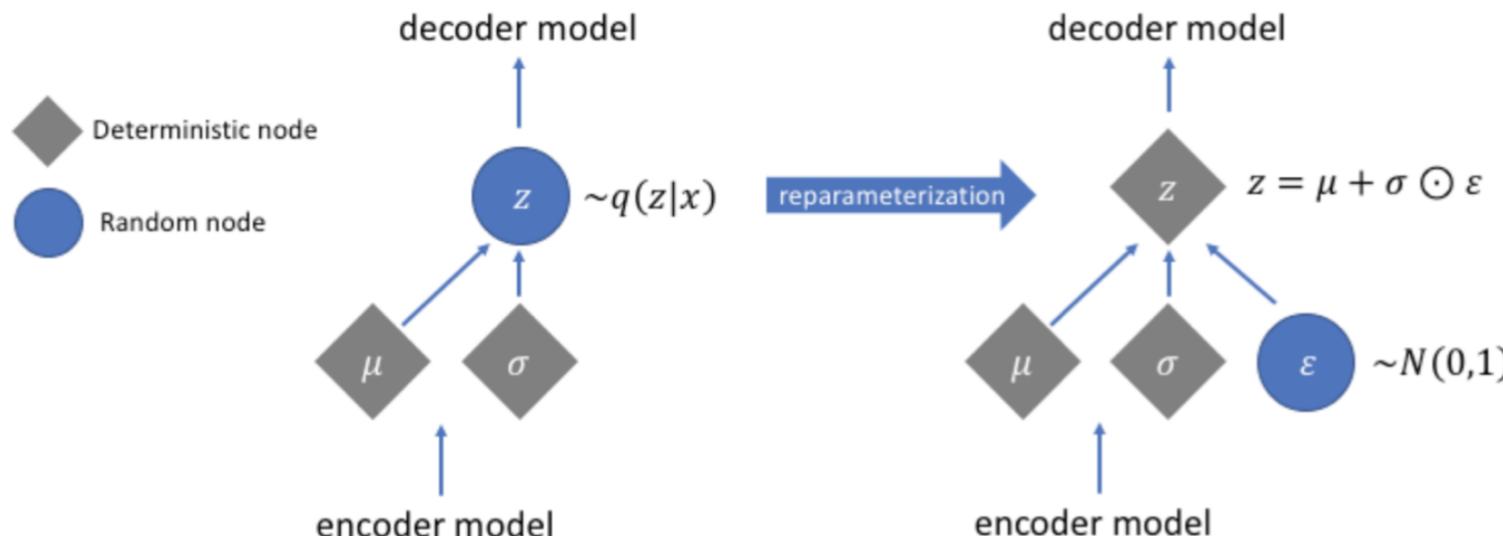
Variational Auto-Encoders (VAEs)

- ELBO:

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

- Reparameterization:

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



Variational Auto-Encoders (VAEs)

- ELBO:

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

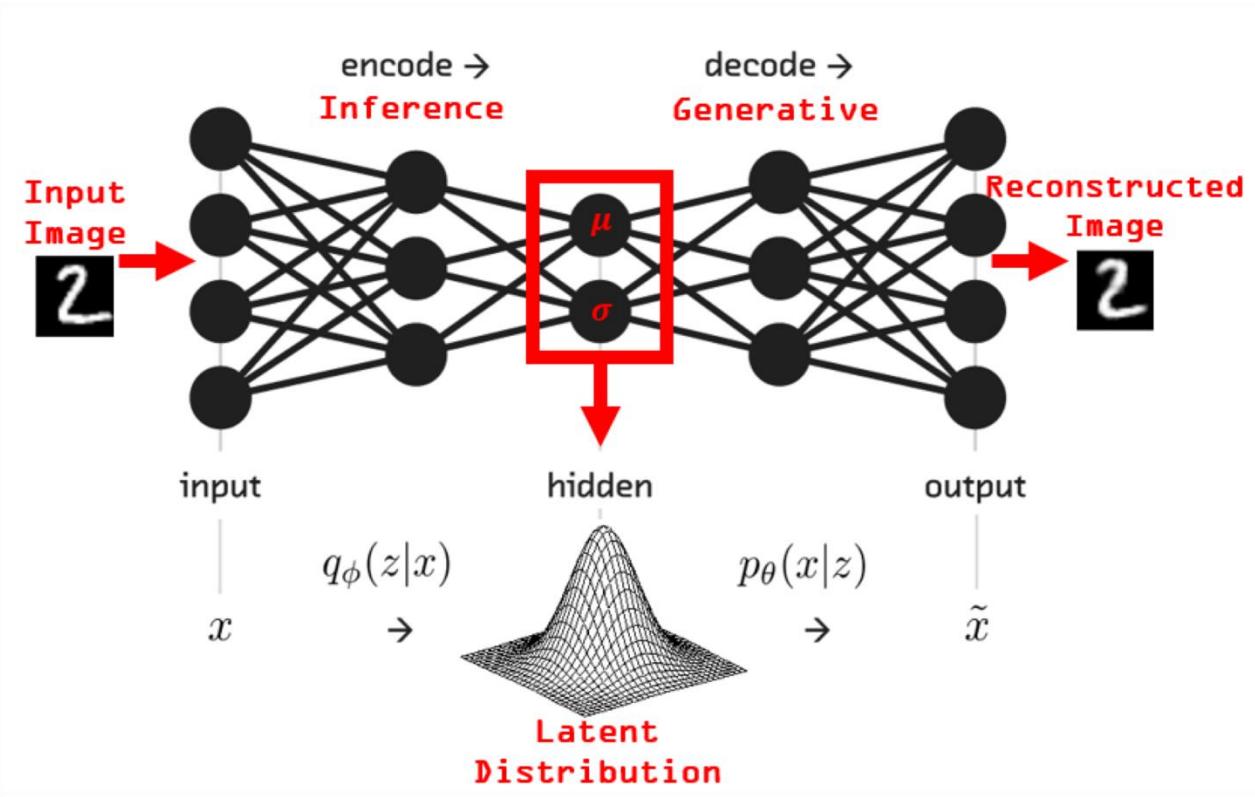
- Reparameterization:

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

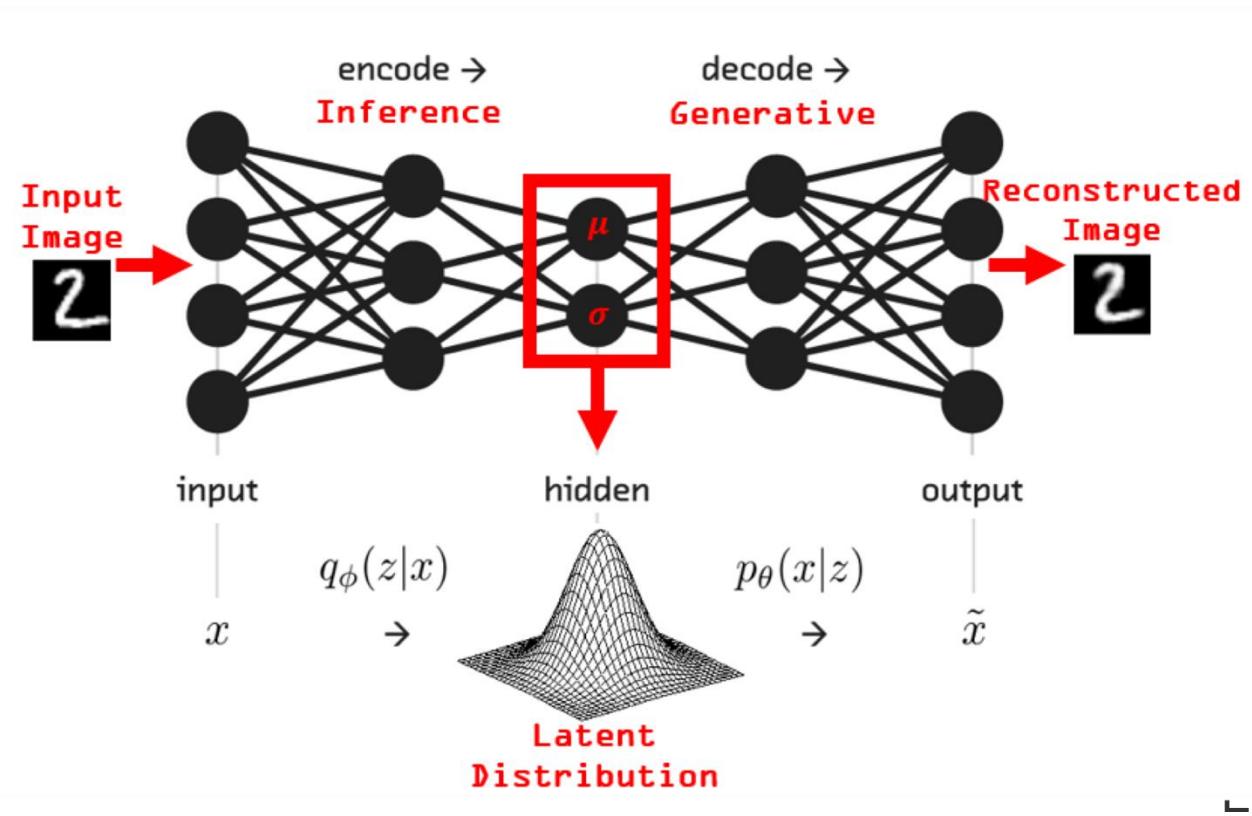
$$\nabla_\phi \mathcal{L} = \mathbb{E}_{\epsilon \sim N(0, 1)} [\nabla_z [\log p_\theta(x, z) - \log q_\phi(z|x)] \nabla_\phi z(\epsilon, \phi)]$$

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

Example: VAEs for images



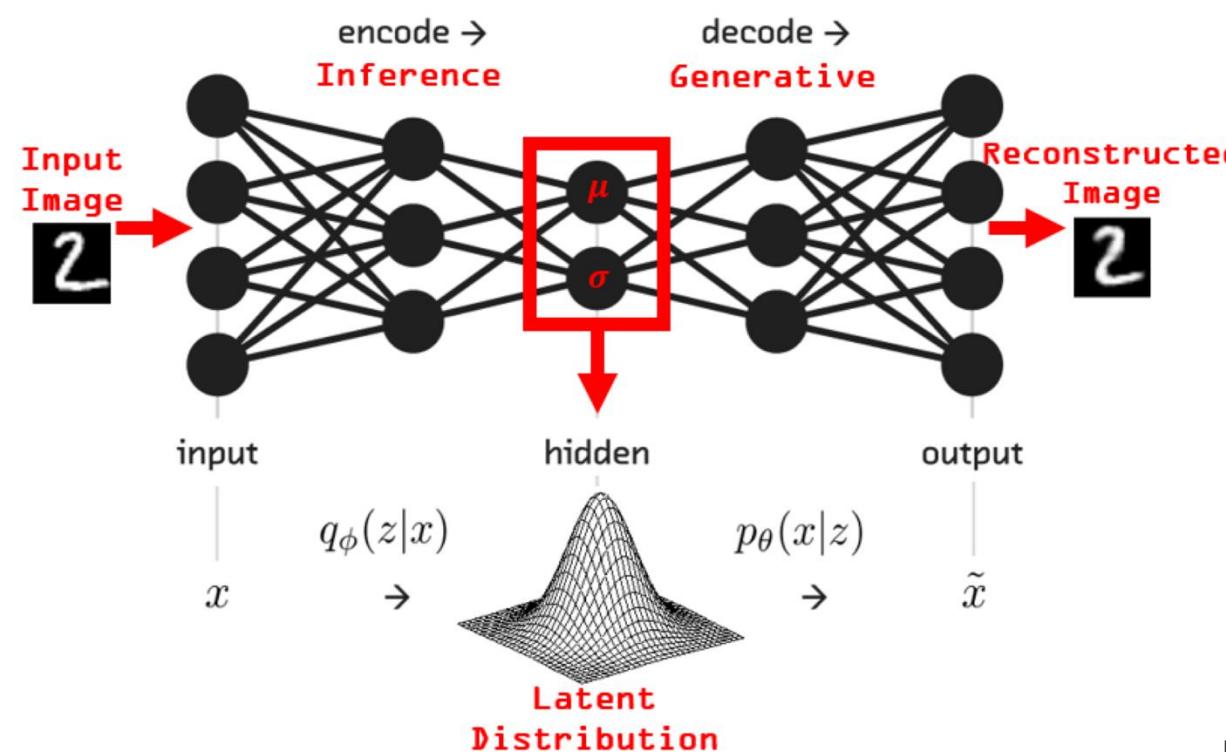
Example: VAEs for images



Input Data

x

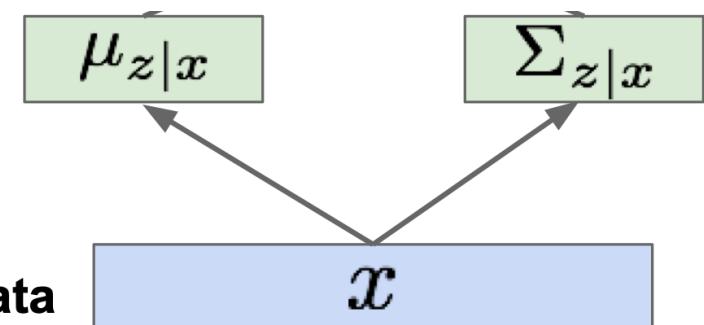
Example: VAEs for images



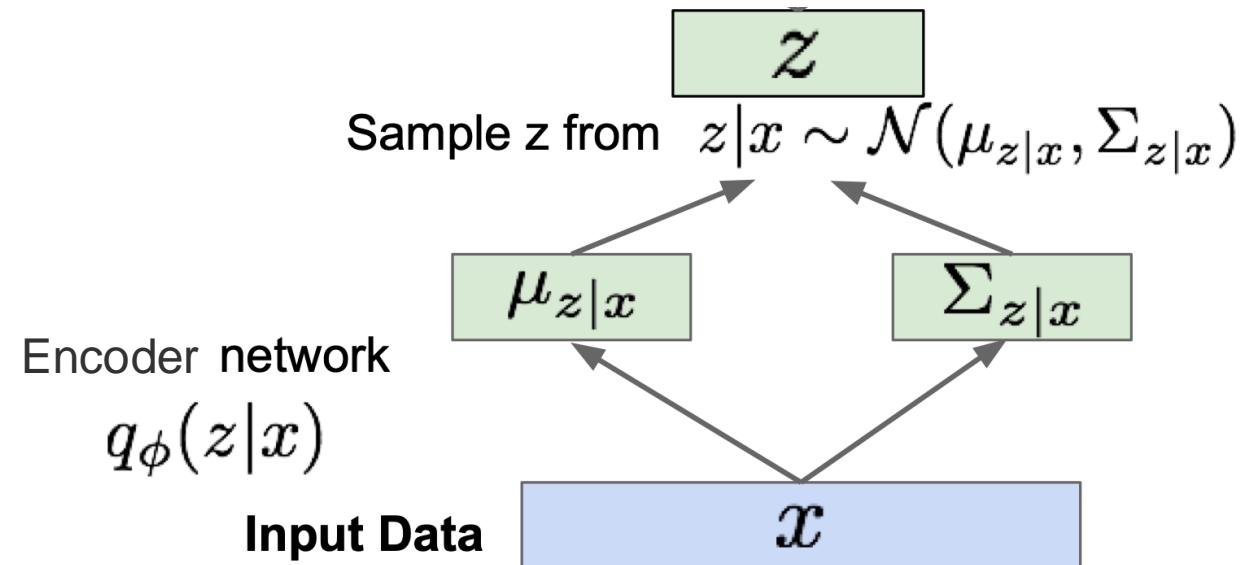
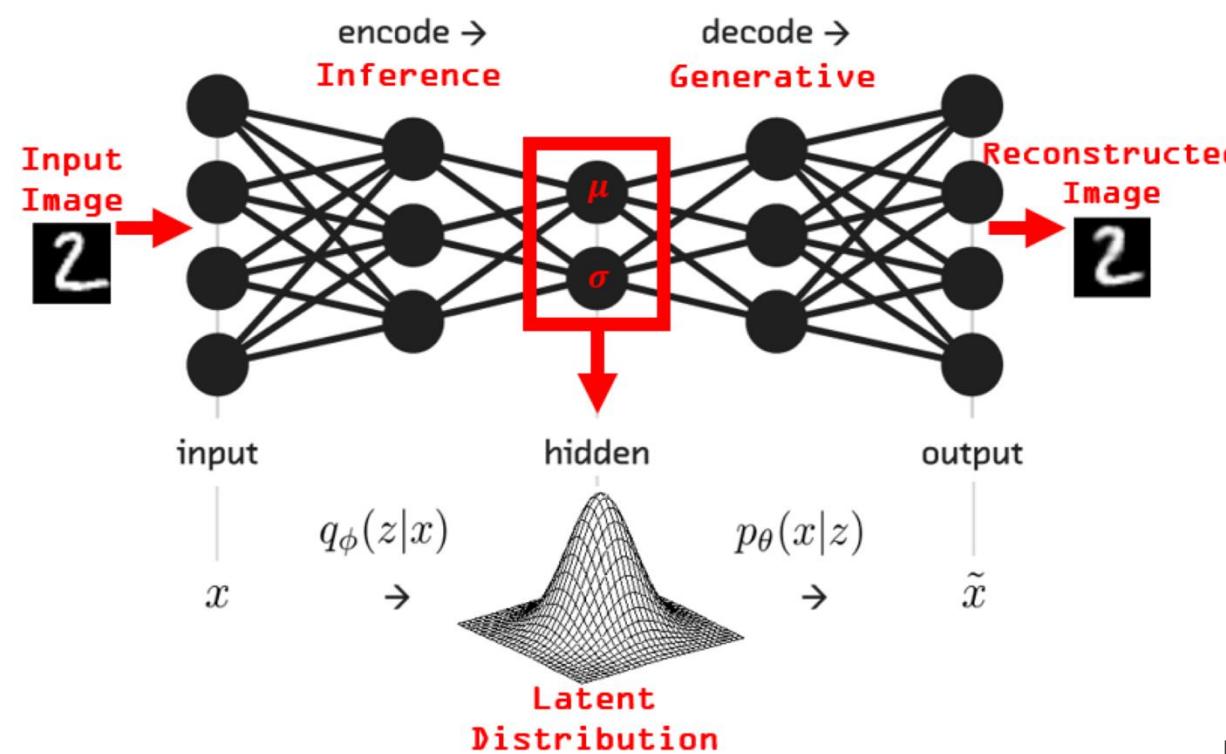
Encoder network

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

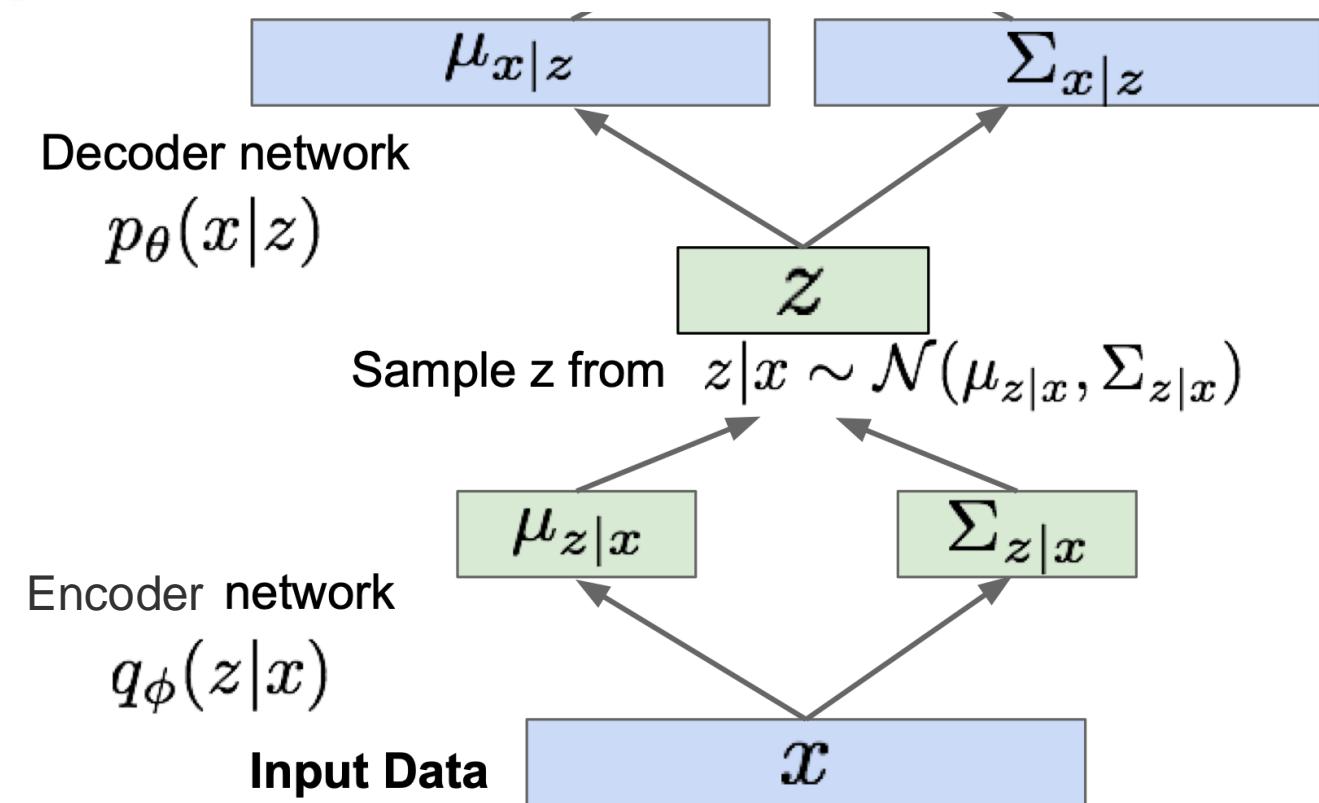
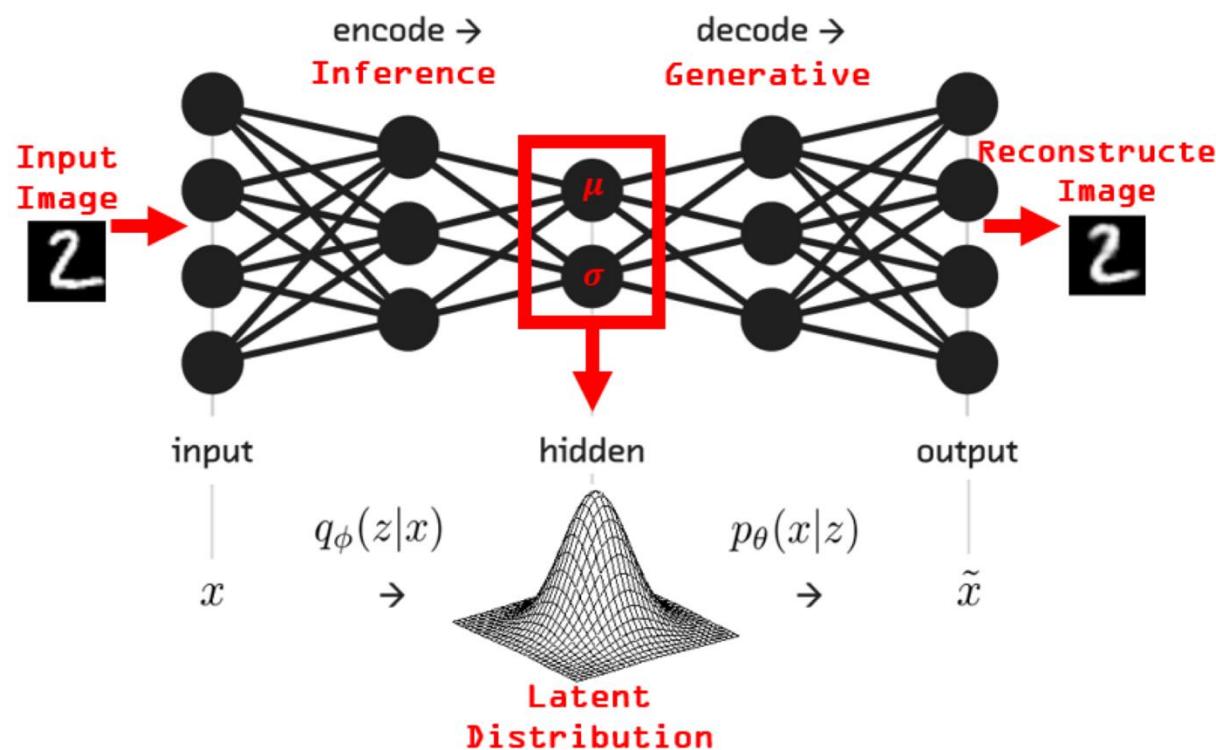
Input Data



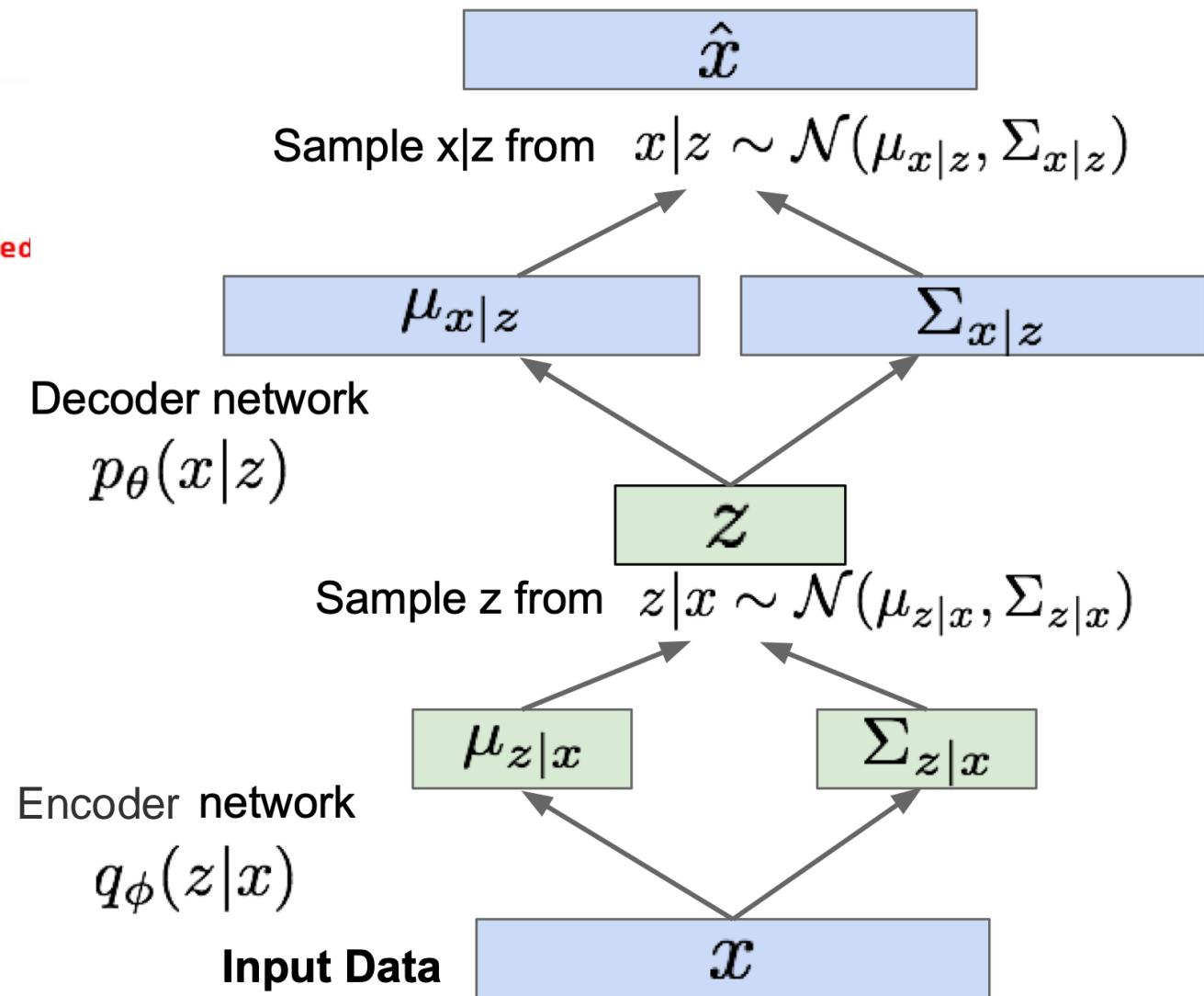
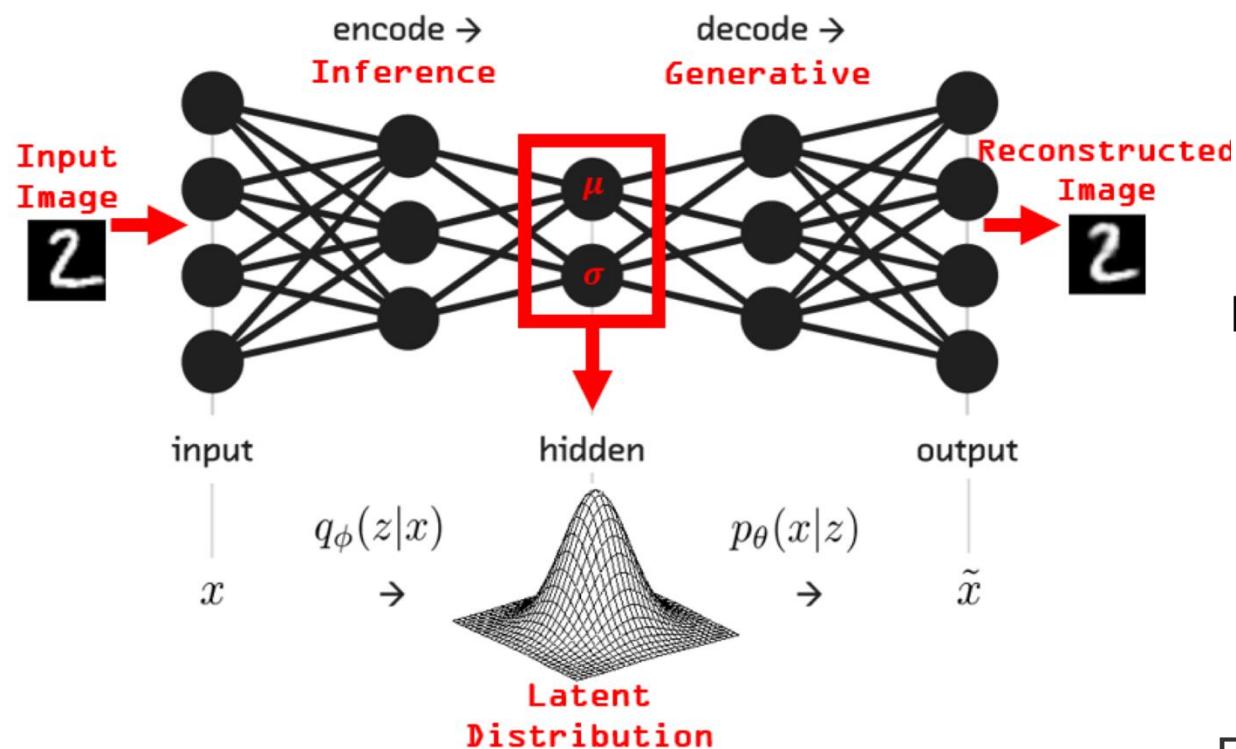
Example: VAEs for images



Example: VAEs for images



Example: VAEs for images

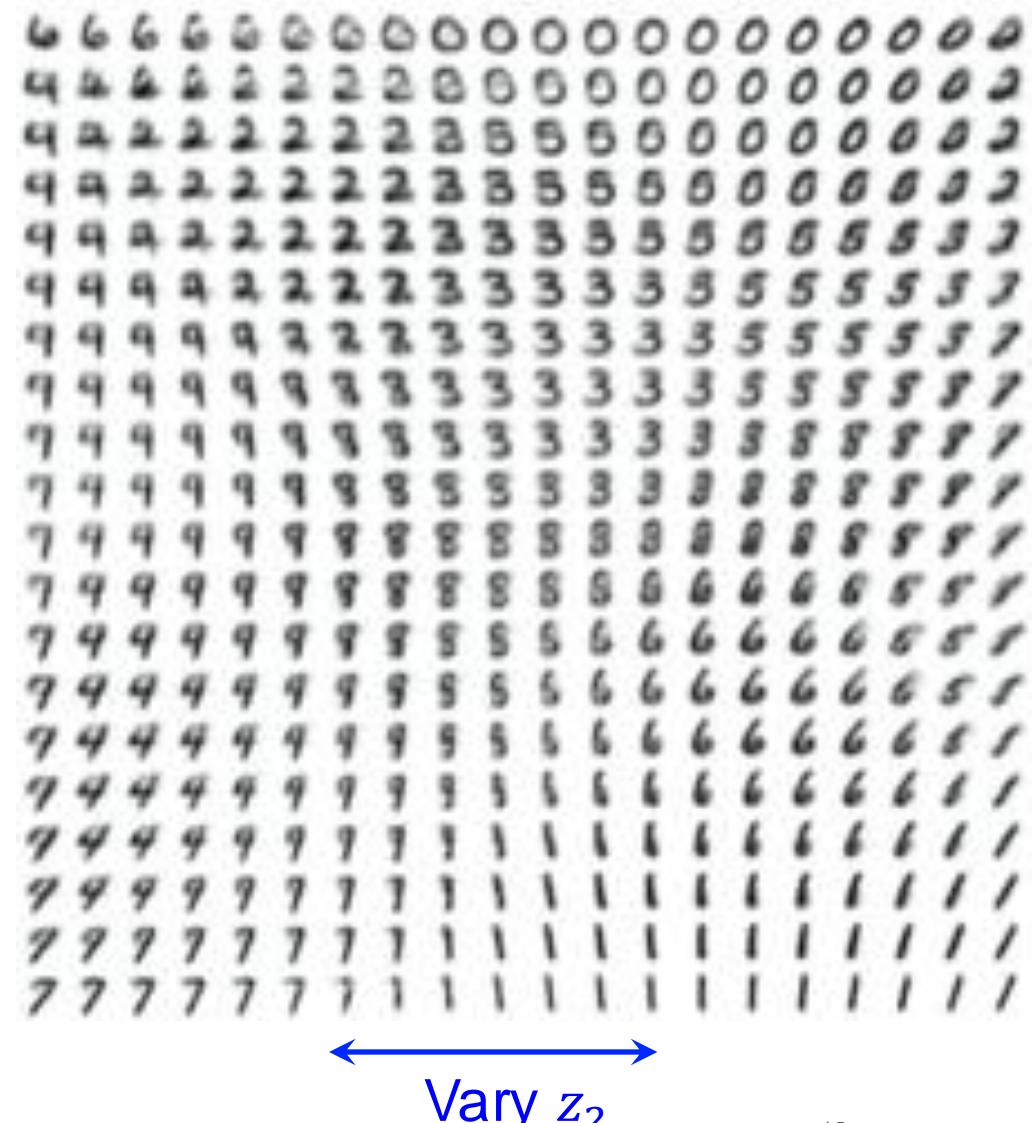
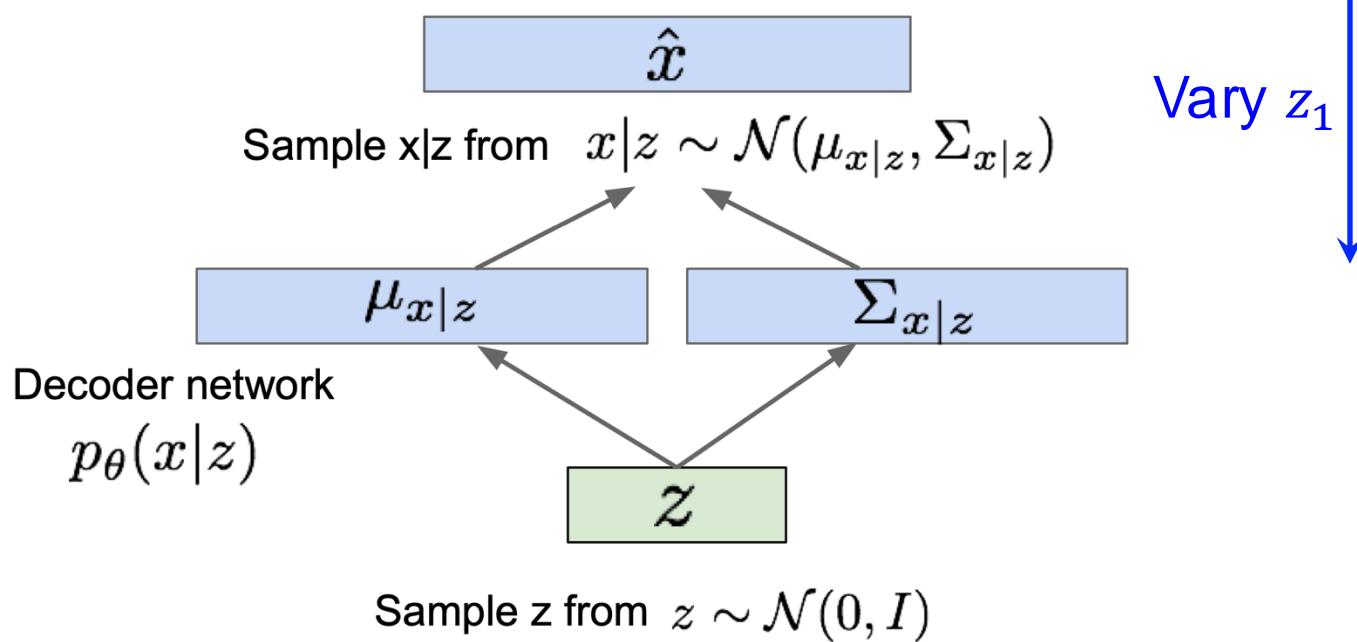


Example: VAEs for images

Data manifold for 2-d z

Generating samples:

- Use decoder network. Now sample z from prior!

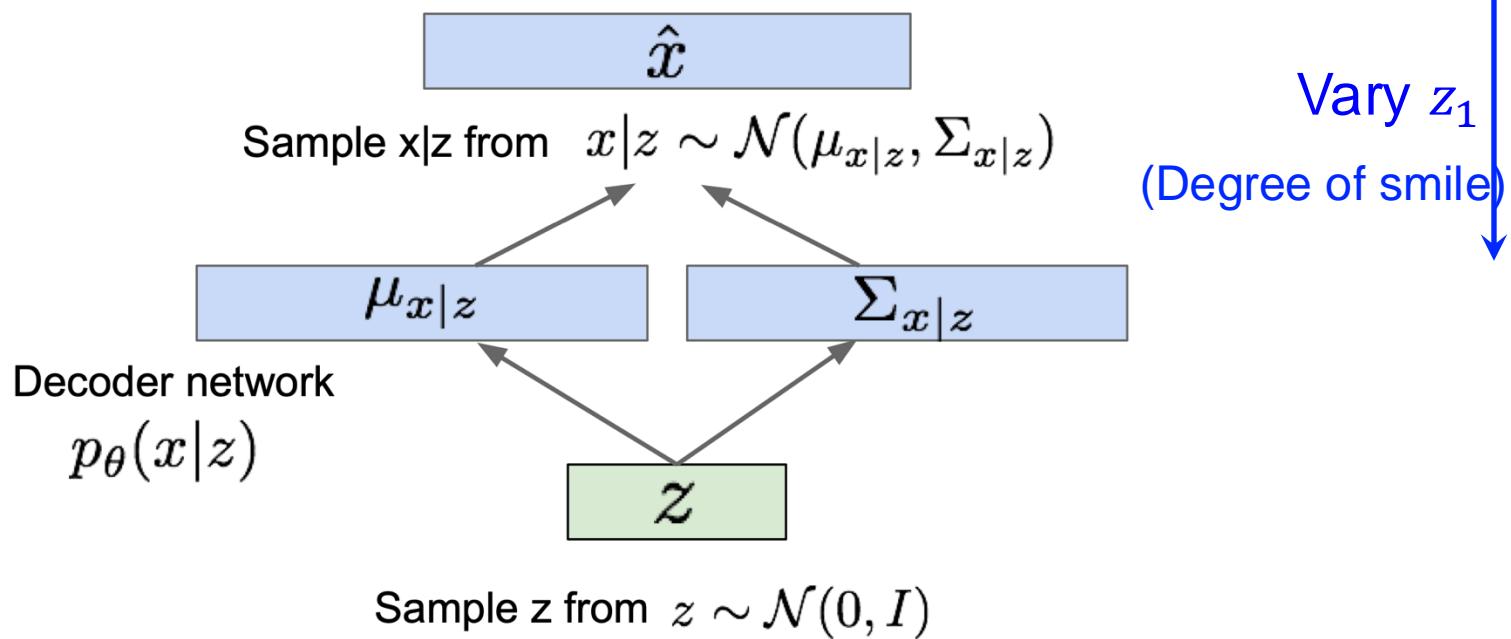


Example: VAEs for images

Data manifold for 2-d z

Generating samples:

- Use decoder network. Now sample z from prior!



Vary z_1 (Degree of smile)
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 $\boldsymbol{\phi}$ (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

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

Reconstruction



Divergence from prior

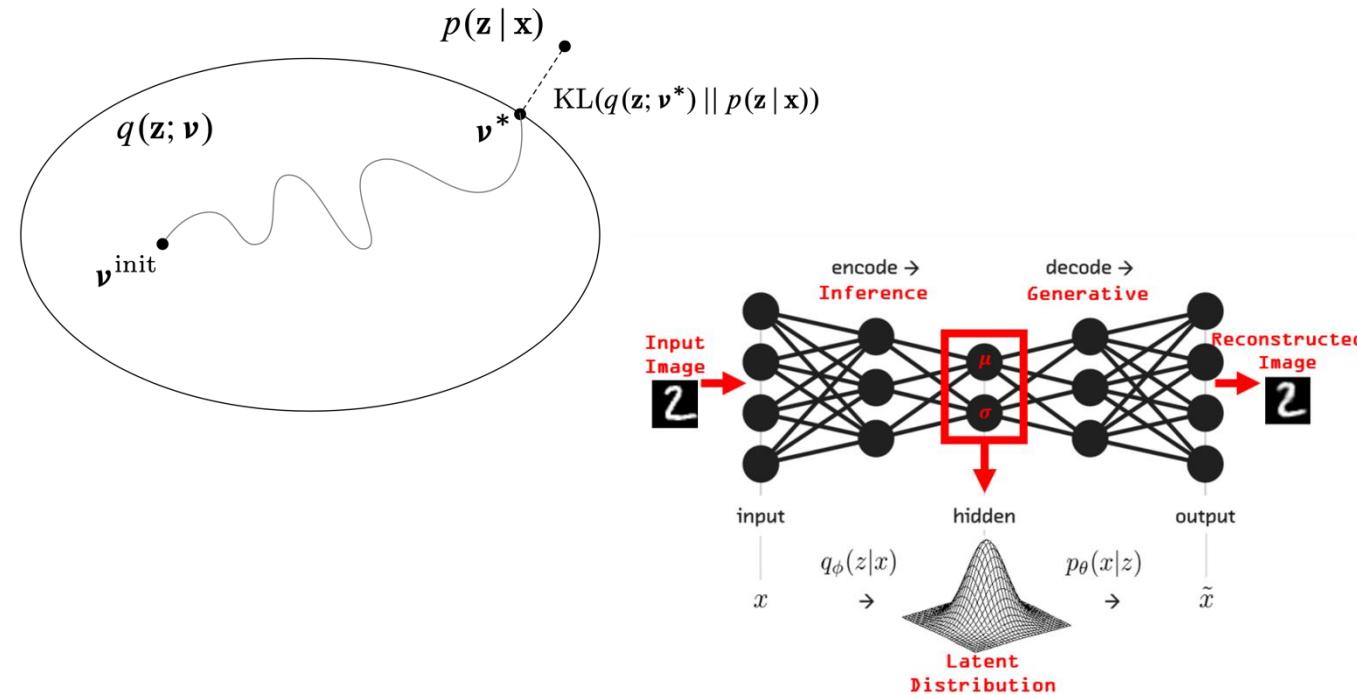


(Razavi et al., 2019)

- Pros:
 - Principled approach to generative models
 - Allows inference of $q(z|x)$, 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?