

DSC291: Machine Learning with Few Labels

Variational Inference
Variational Autoencoders

Zhiting Hu

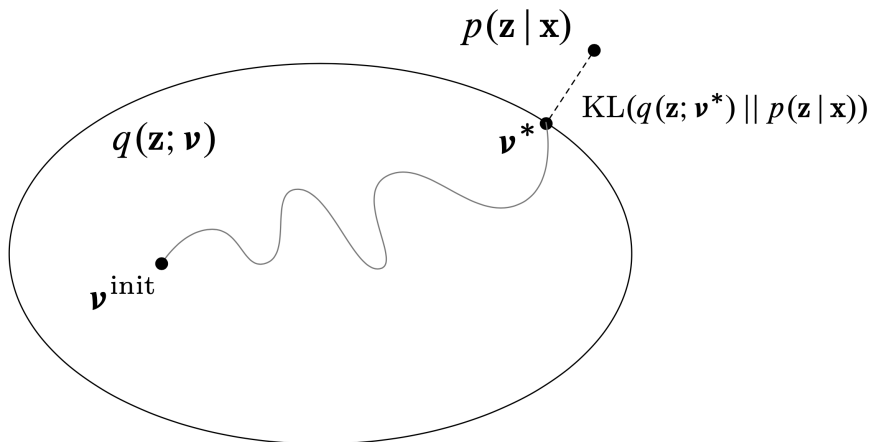
Lecture 8, January 27, 2023

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Recap: VI

- We often cannot compute posteriors, and so we need to approximate them, using variational methods.
- In variational Bayes, we'd like to find an approximation within some family that minimizes the KL divergence to the posterior, but we can't directly minimize this
- Therefore, we defined the ELBO, which we can maximize, and this is equivalent to minimizing the KL divergence.



Evidence Lower Bound (ELBO)

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

Recap: Mean-Field VI

- We defined a family of approximations called “mean field” approximations, in which there are no dependencies between latent variables

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

- We optimize the ELBO with coordinate ascent updates to iteratively optimize each local variational approximation under mean field assumptions

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(\mathbf{x}, \mathbf{z})] \right\}$$

Recap: VI with coordinate ascent

Example: Bayesian mixture of Gaussians

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

- Initialize the global variational distributions $q(\mu_k)$ and $q(\pi)$
- **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 $q(\mu_k)$ and $q(\pi)$
- **Until** ELBO converges

- What if we have millions of data examples? This could be very slow.

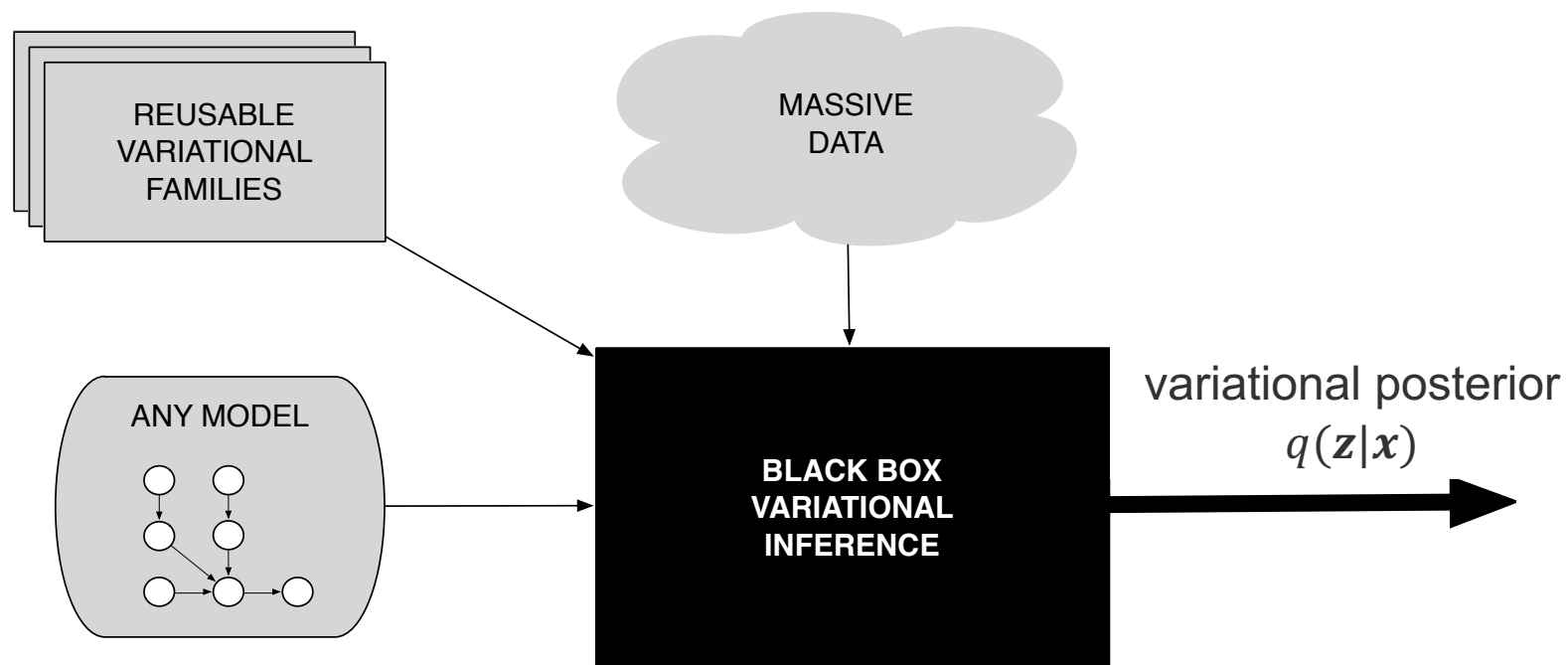
Recap: Stochastic VI

Example: Bayesian mixture of Gaussians

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

- Initialize the global variational distributions $q(\mu_k)$ and $q(\pi)$
- **Repeat:**
 - Sample a data example $i \in \{1, 2, \dots, D\}$
 - Update the local variational distribution $q(z_i)$
 - Update the global variational distributions $q(\mu_k)$ and $q(\pi)$ with **natural gradient ascent**
- **Until** ELBO converges
- (Setting natural gradient = 0 gives the traditional mean-field update)

Recap: Black-box Variational Inference (BBVI)



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

Recap: 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

Recap: 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 unfeasible
- 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)

Recap: 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 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
- 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
 - To reduce variance: Rao-Blackwellization, control variates, importance sampling, ...

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

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

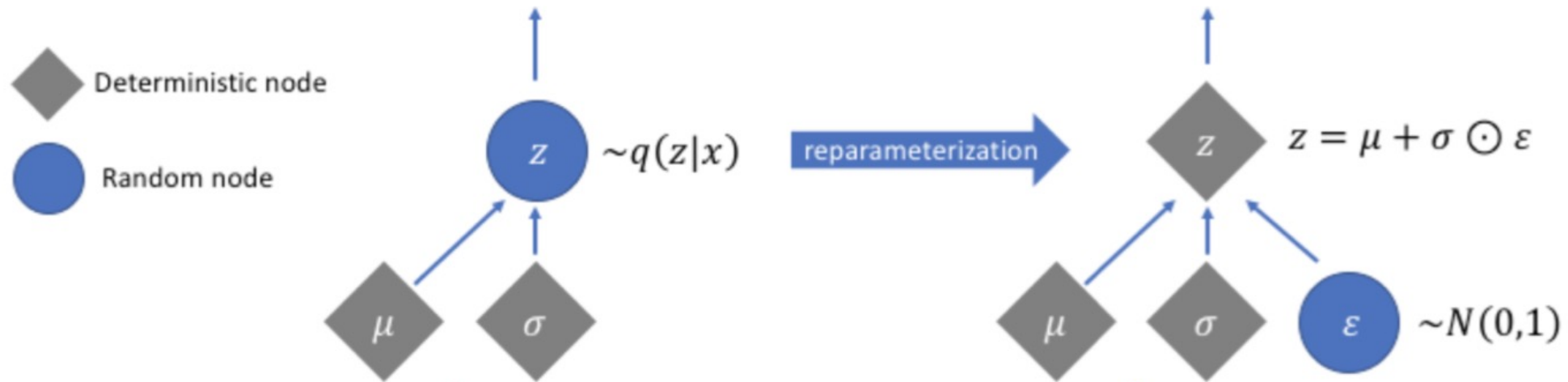
$$\nabla_\lambda \mathcal{L} = \mathbb{E}_{\epsilon \sim s(\epsilon)}[\nabla_{\mathbf{z}} f_\lambda(\mathbf{z}) \nabla_\lambda t(\epsilon, \lambda)]$$

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

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} : $z = F^{-1}(\epsilon)$
 - 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

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

- ELBO:

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

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

BBVI with the score gradient

- ELBO:

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

- Gradient w.r.t. λ (using the log-derivative trick)

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

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

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

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

BBVI with the reparameterization gradient

- ELBO:

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

- Gradient w.r.t. λ

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

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

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:

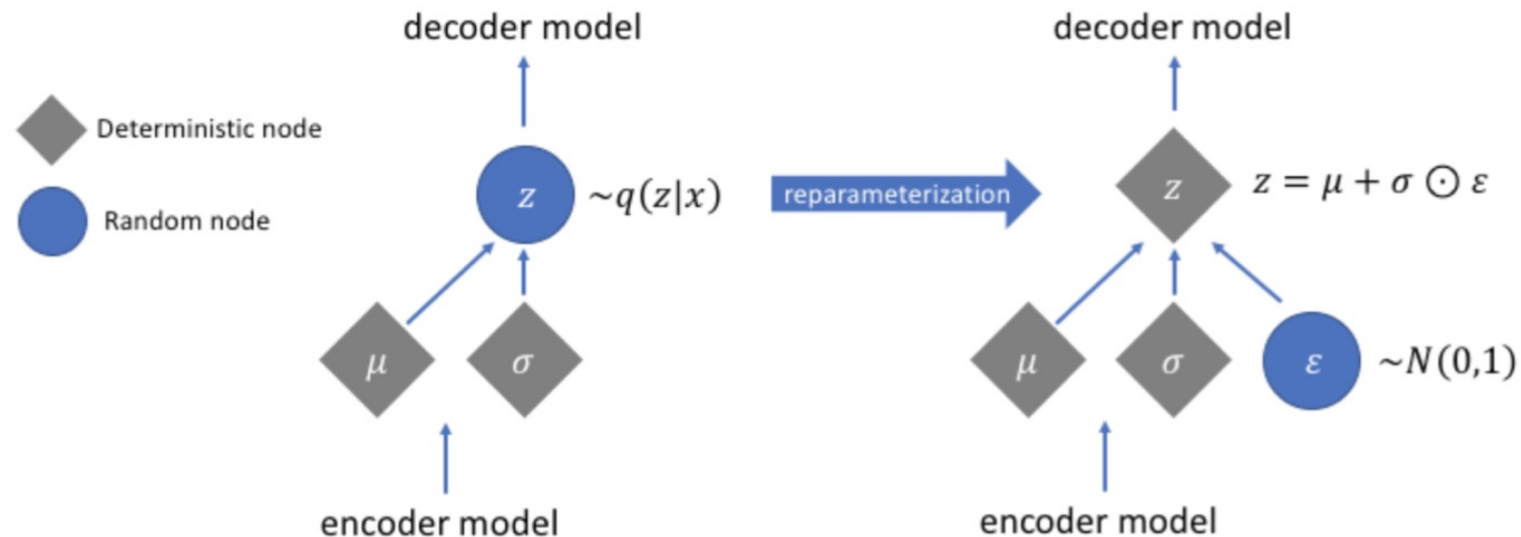
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] -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

Reconstruction

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

Variational Auto-Encoders (VAEs)

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



Variational Auto-Encoders (VAEs)

- ELBO:

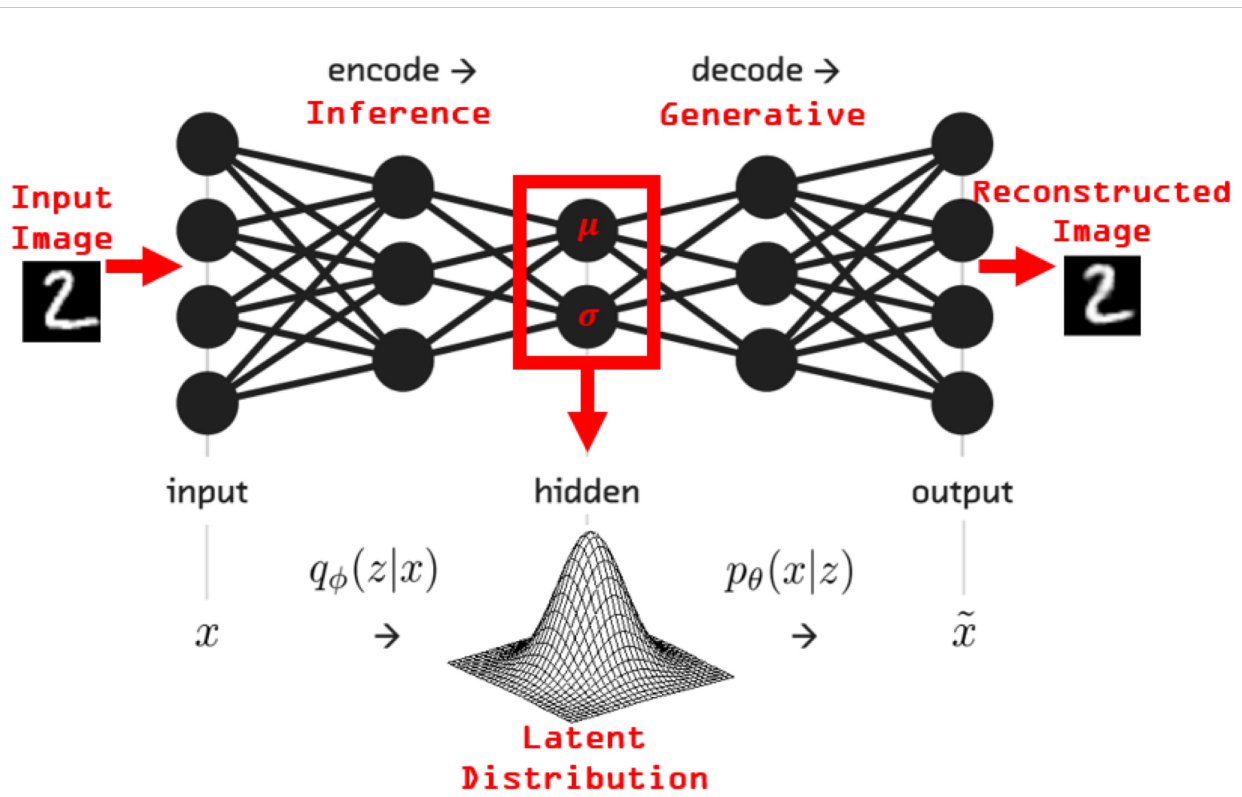
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] + H(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log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

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

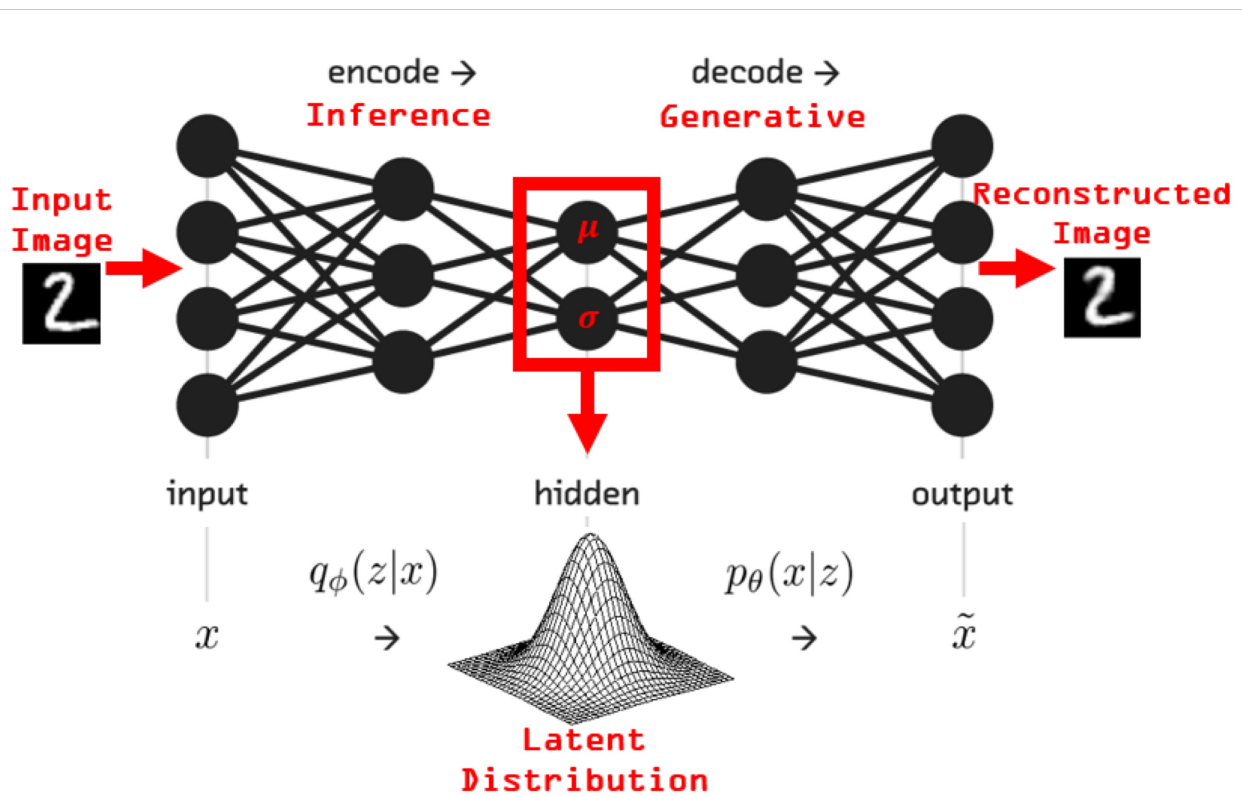
$$\nabla_{\boldsymbol{\phi}} \mathcal{L} = \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})} [\nabla_{\mathbf{z}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})] \nabla_{\boldsymbol{\phi}} \mathbf{z}(\boldsymbol{\epsilon}, \boldsymbol{\phi})]$$

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

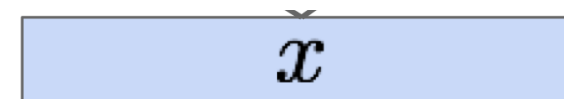
Example: VAEs for images



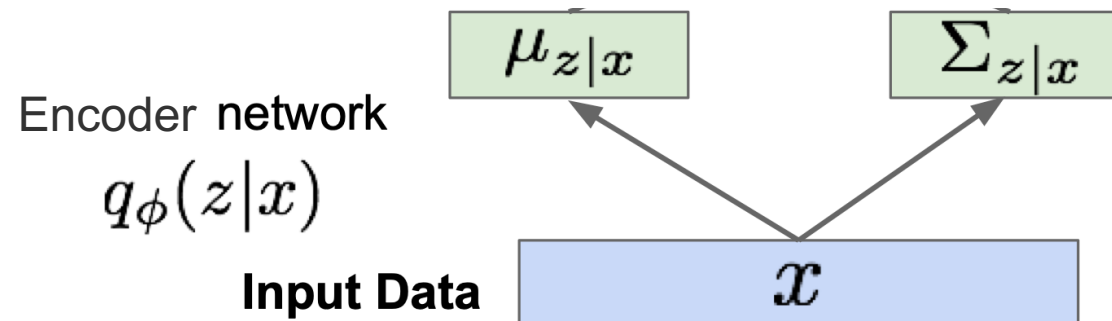
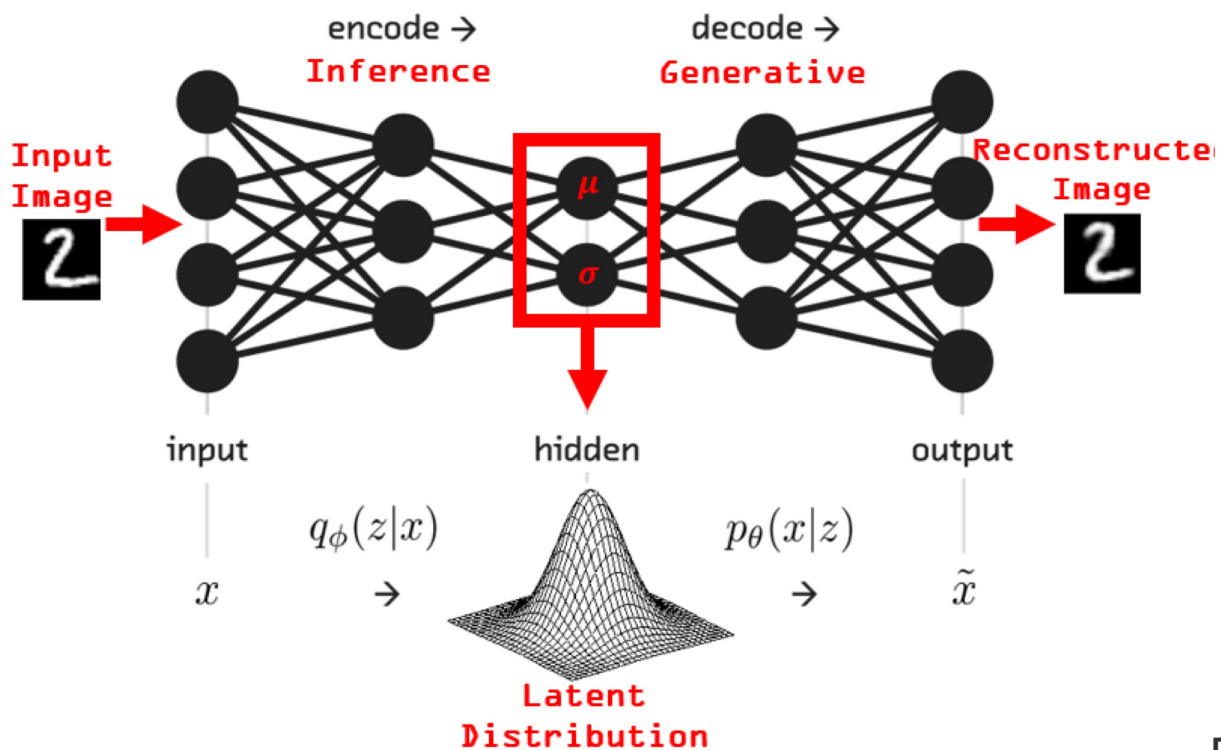
Example: VAEs for images



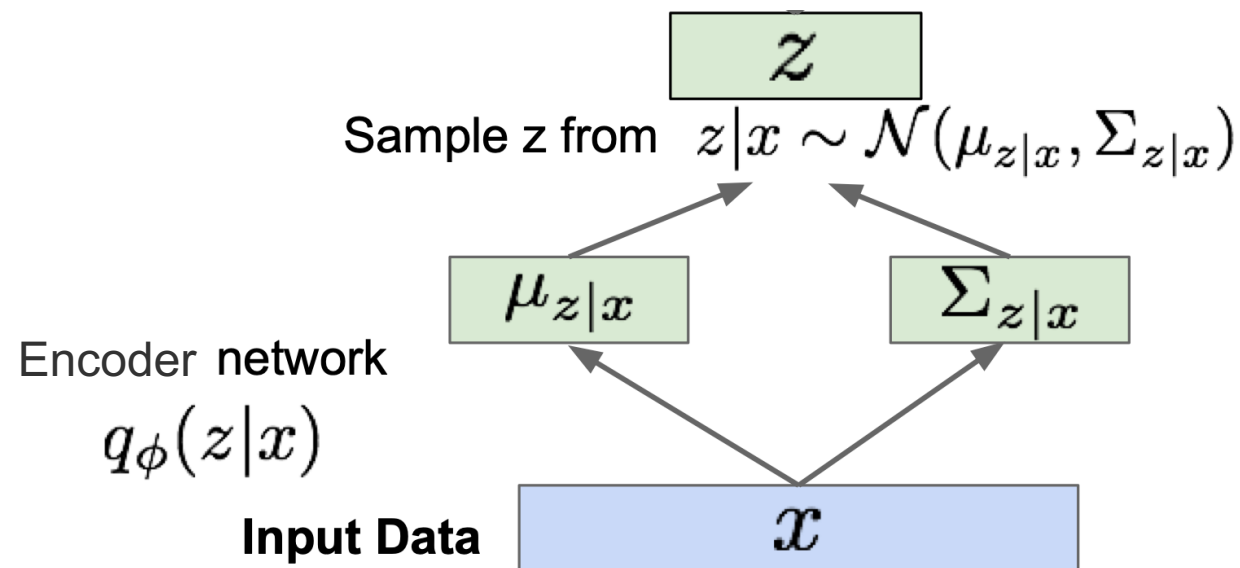
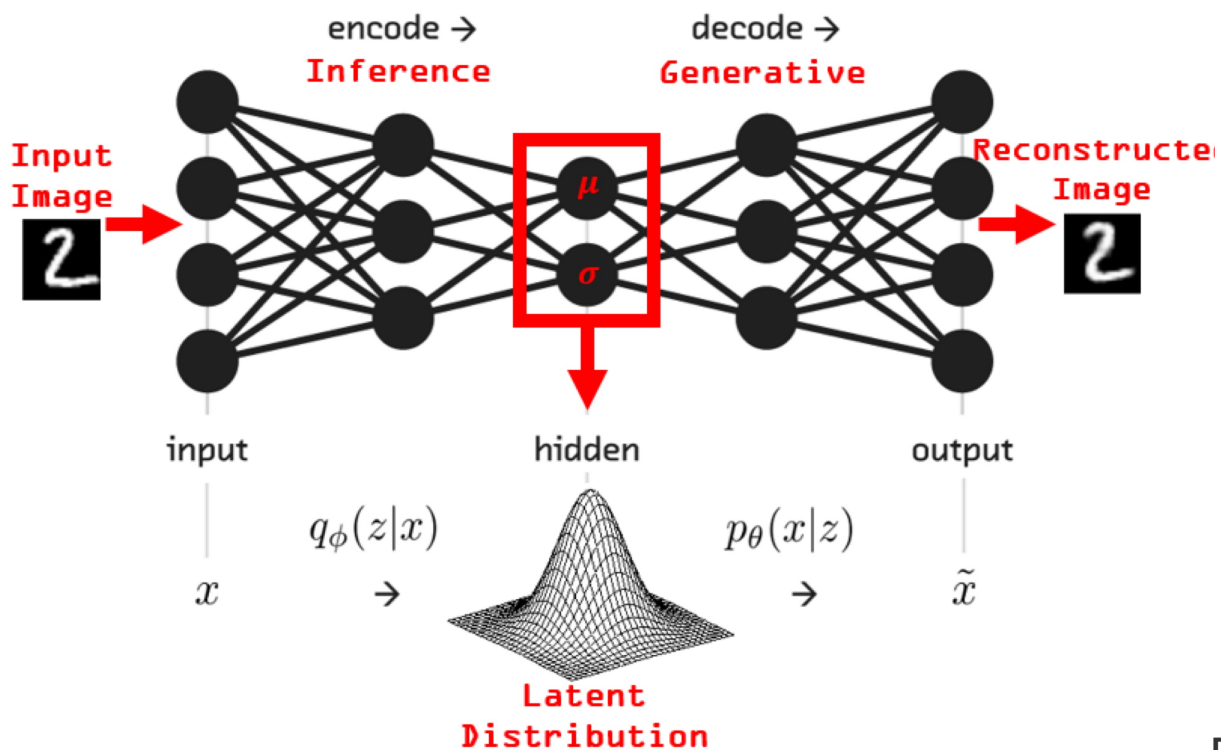
Input Data



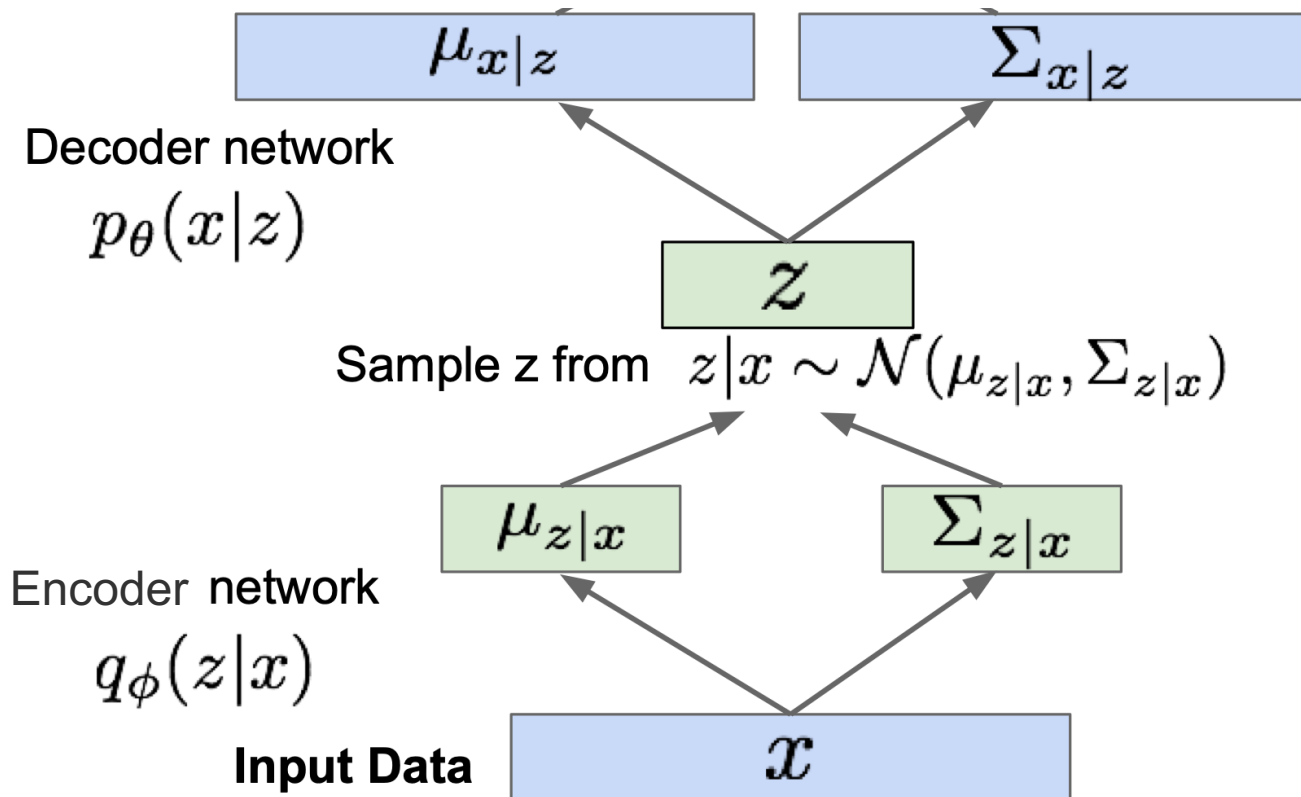
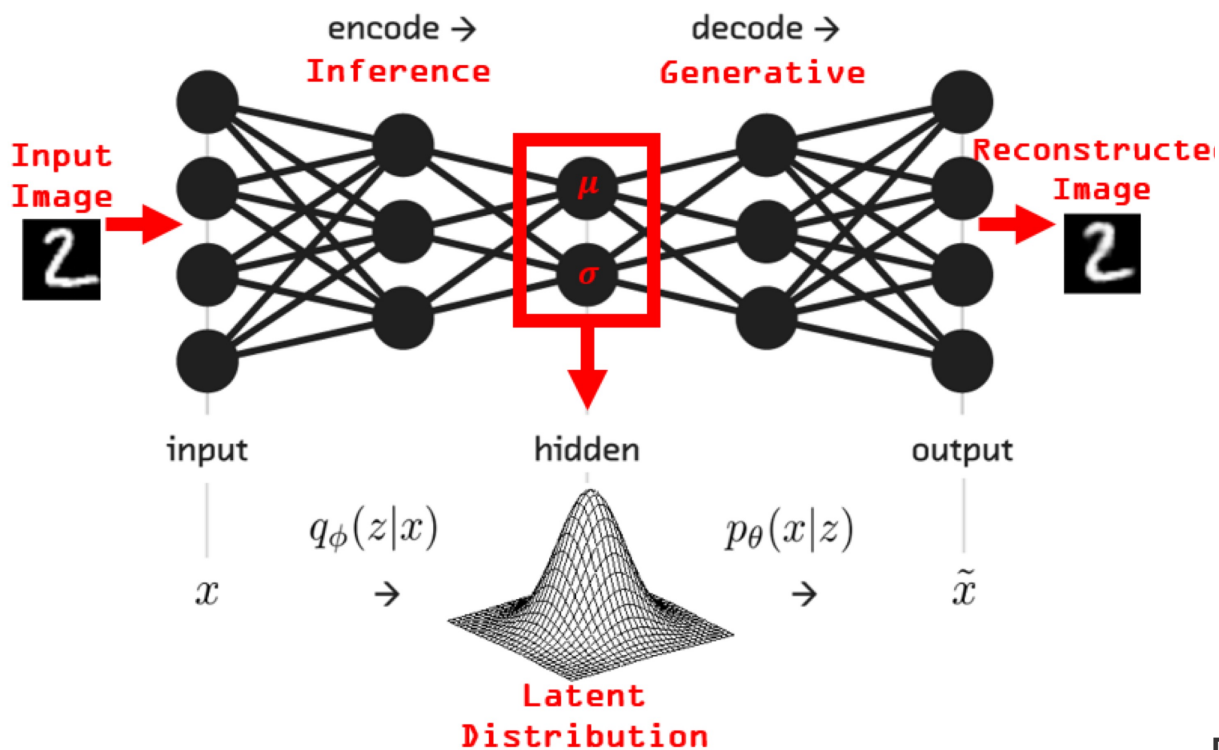
Example: VAEs for images



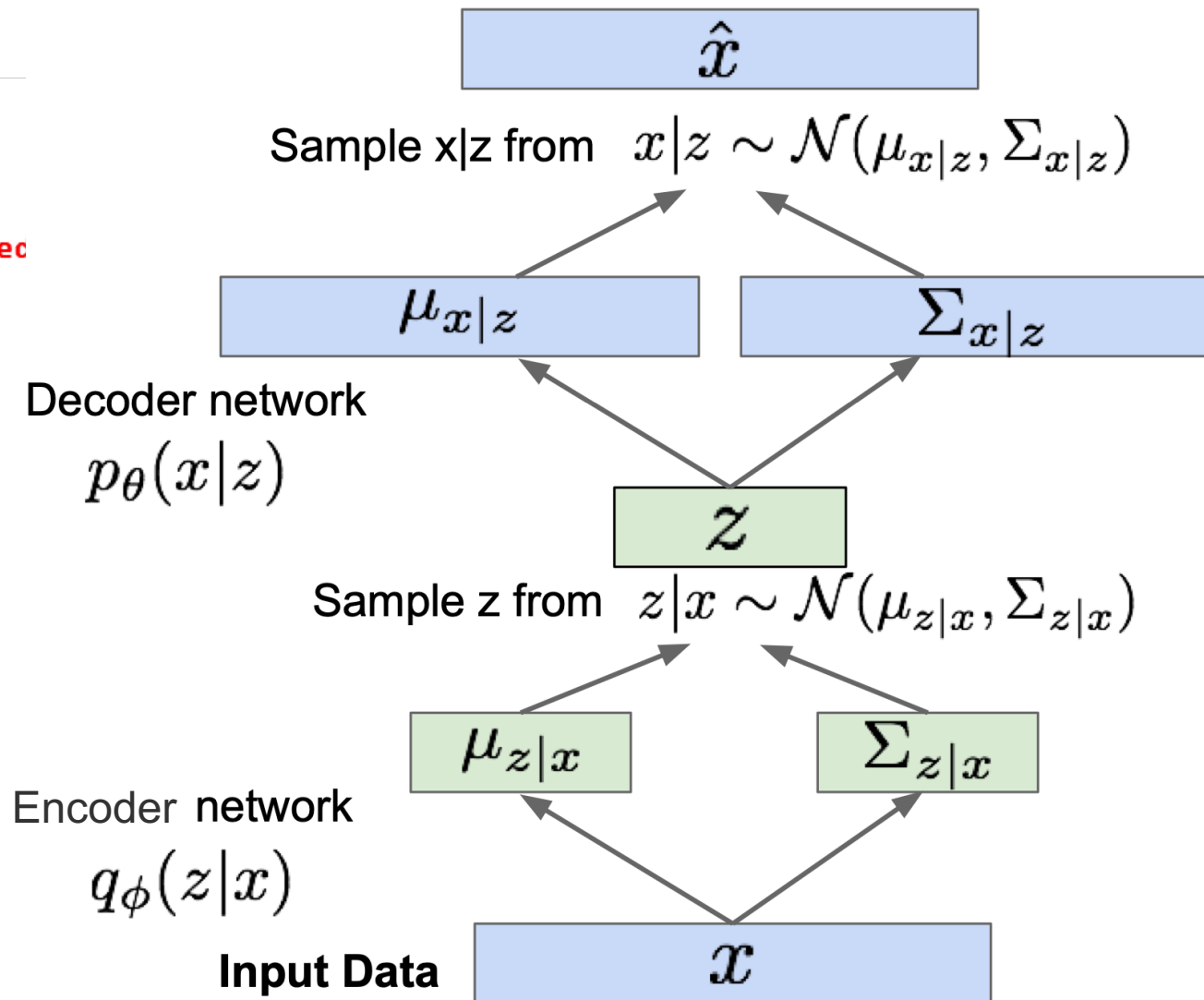
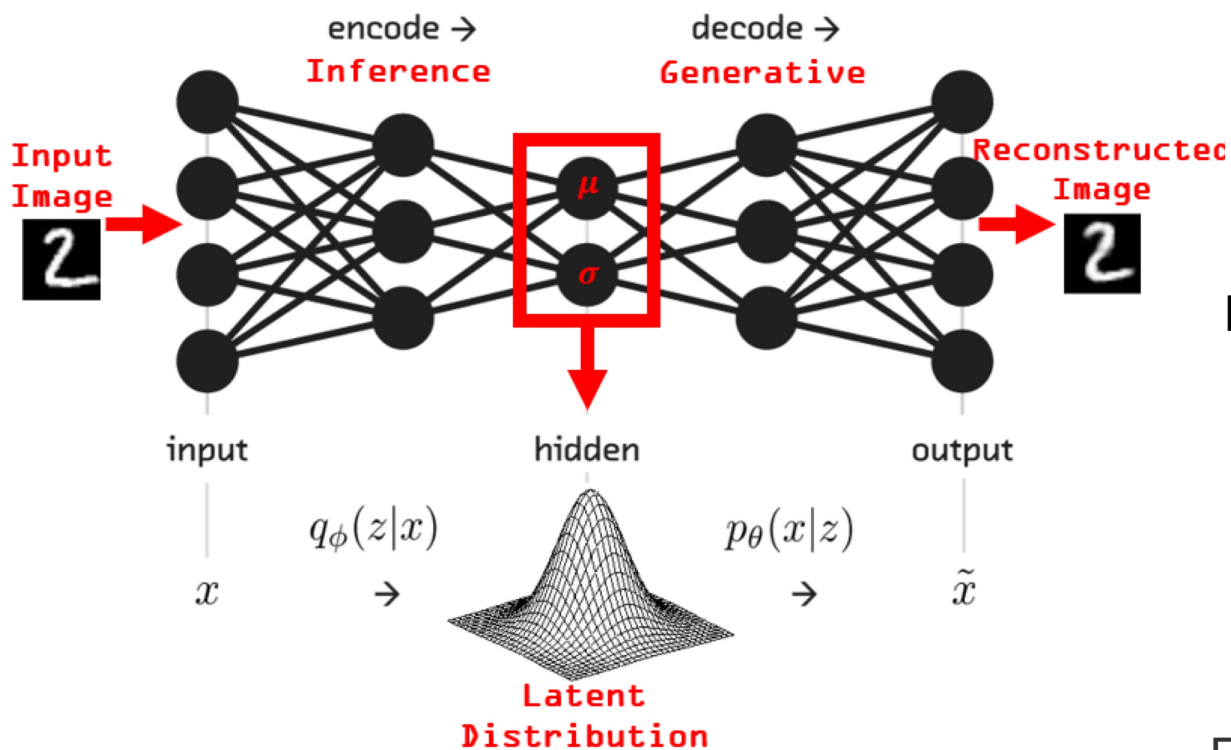
Example: VAEs for images



Example: VAEs for images



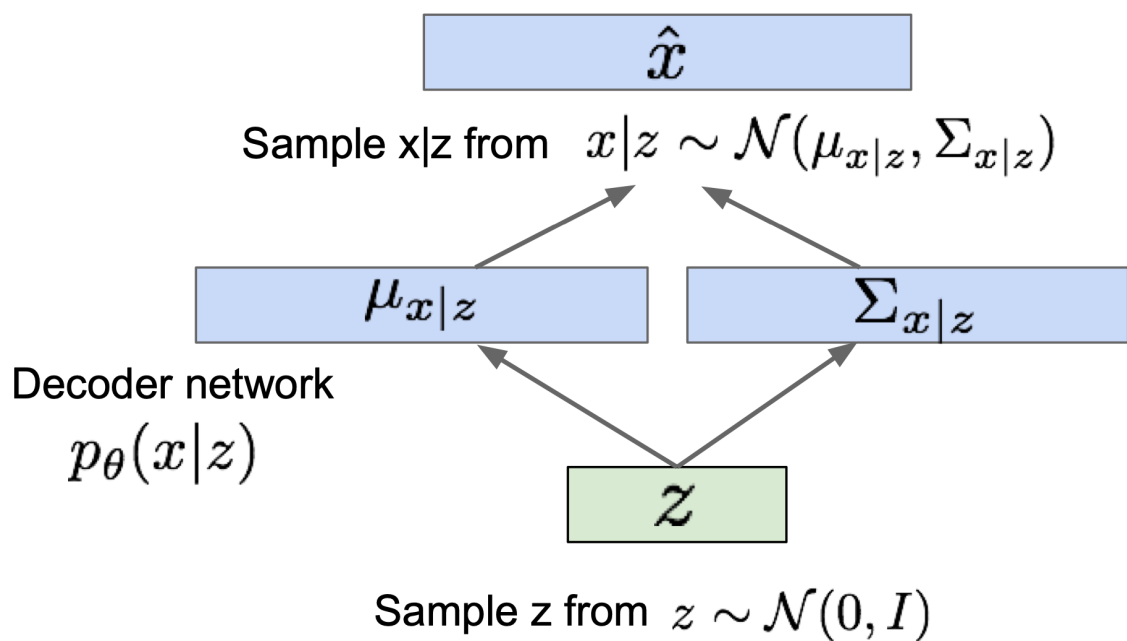
Example: VAEs for images



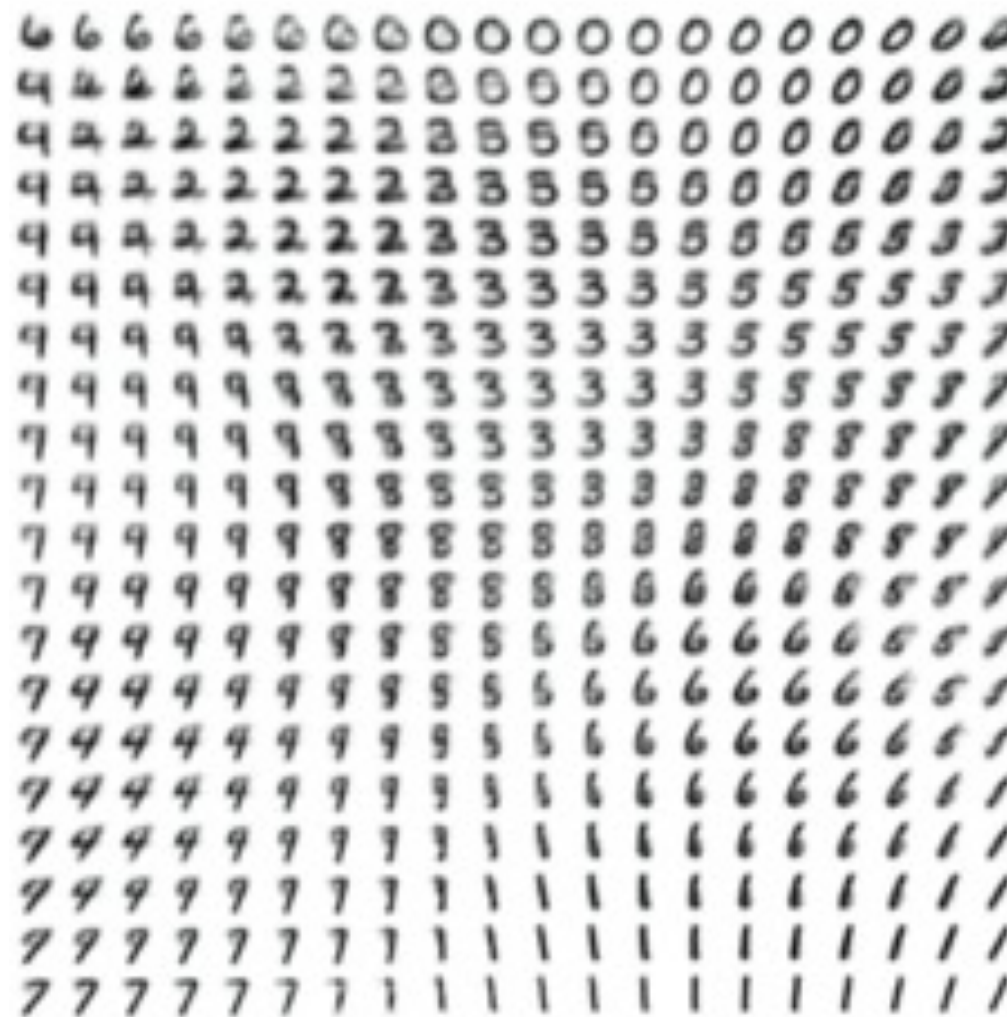
Example: VAEs for images

Generating samples:

- Use decoder network. Now sample z from prior!



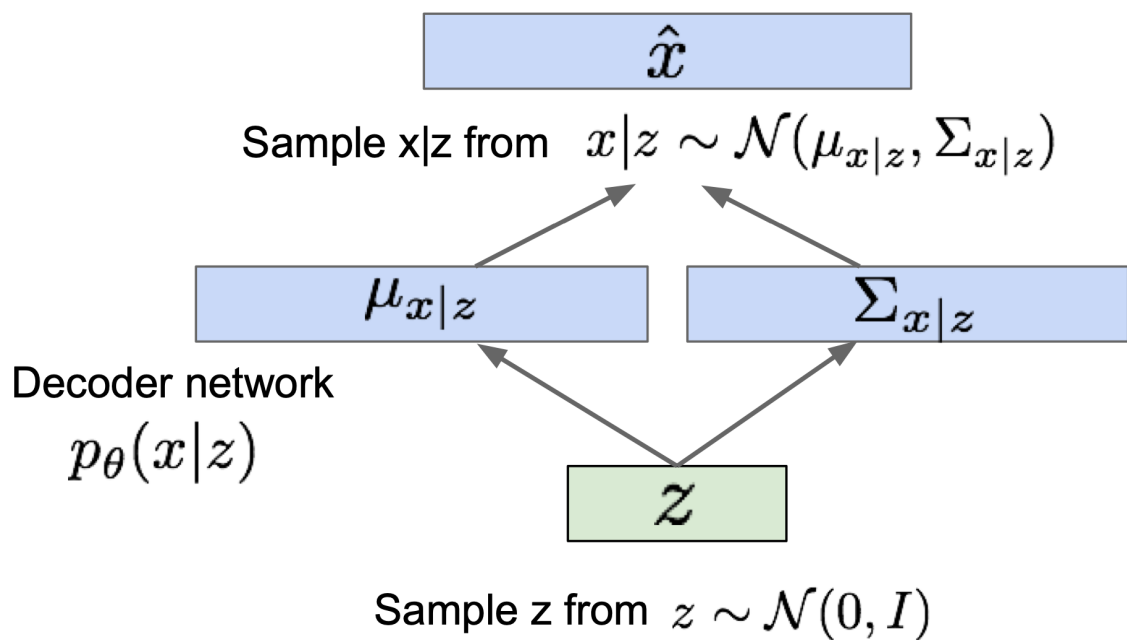
Data manifold for 2-d z



Example: VAEs for images

Generating samples:

- Use decoder network. Now sample z from prior!



Data manifold for 2-d z

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 .

Variational Auto-Encoders (VAEs)

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

[Kingma & Welling, 2014]

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

Key Takeaways

- Stochastic VI
- Computing Gradients of Expectations $\mathcal{L} = \mathbb{E}_{q_\lambda(\mathbf{z})}[f_\lambda(\mathbf{z})]$

- Score gradient

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

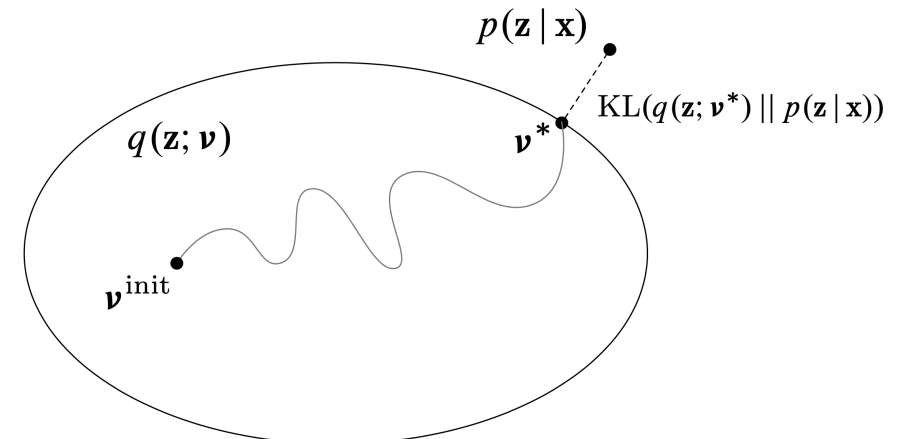
- Reparameterization gradient

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

- Black-box VI
- Variational autoencoders (VAEs)

Summary so far: Supervised Learning, Unsupervised Learning

- Supervised Learning
 - Maximum likelihood estimation (MLE)
 - Duality between MLE and Maximum Entropy Principle
- Unsupervised learning
 - Maximum likelihood estimation (MLE) with latent variables
 - Marginal log-likelihood
 - EM algorithm for MLE
 - ELBO
 - Variational Inference
 - ELBO
 - Variational distributions



Questions?