

DSC291: Machine Learning with Few Labels

Self-supervised Learning

Zhitong Hu

Lecture 7, April 22, 2025

Outline

- Variational Auto-Encoders (VAEs)
- Self-Supervised Learning (SSL)

PRML

Recap: EM and Variational Inference

- The EM algorithm:

- E-step:

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

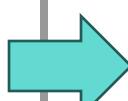
Intractable when model $p(z, x | \theta)$ is complex

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

- M-step:

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

~~$p(x, z | \theta)$~~



Need to approximate $p(z|x, \theta^t)$ with Variational Inference (VI)

① MCMC
Markov Chain Monte Carlo
Sampling
optimization

HMC
Langevin dynamics



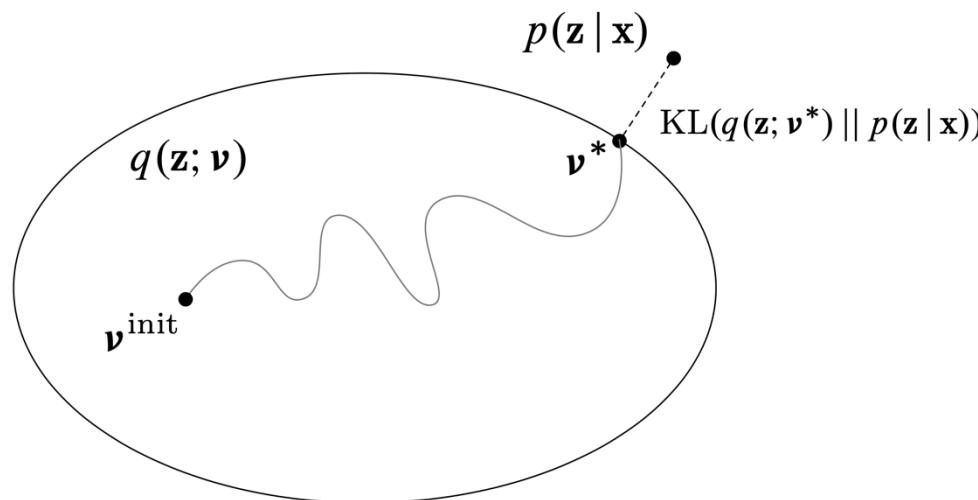
$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \underline{\text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))}$$

Recap: Variational Inference

Maximizing KL(q(z|x, v) || p(z|x))

- 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 $p(\mathbf{z}|\mathbf{x})$

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

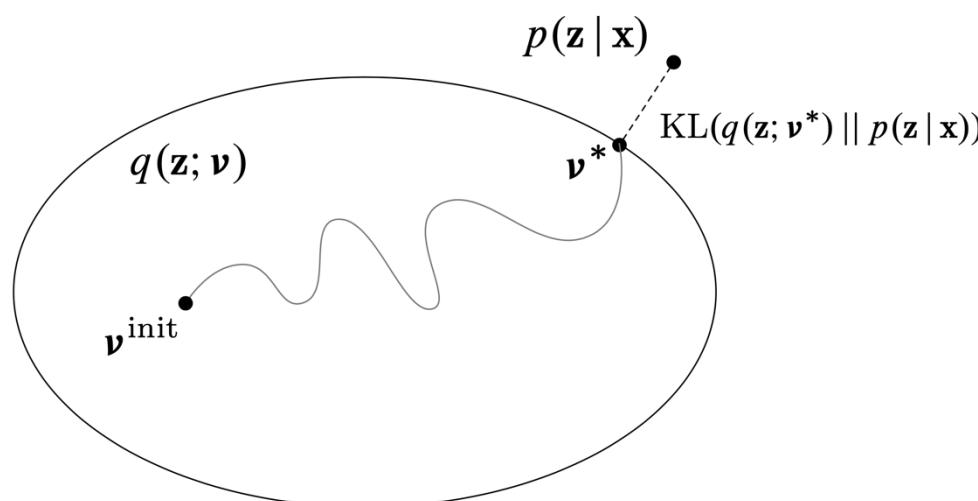


$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

Recap: Variational Inference

- Choose a family of distributions over the latent variables \mathbf{z} with its own set of variational parameters ν , i.e. $q(\mathbf{z}|\mathbf{x}, \nu) \xrightarrow{\phi}$
- 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)

Variational Auto-Encoders (VAEs)

- Model $p_{\theta}(x, z) = p_{\theta}(x|z)p(z)$ θ
○ $p_{\theta}(x|z)$: a.k.a., generative model, generator, (probabilistic) decoder, ...
○ $p(z)$: prior, e.g., Gaussian
- Assume variational distribution $q_{\phi}(z|x)$ ϕ
○ E.g., a Gaussian distribution parameterized as **deep neural networks**
○ a.k.a, recognition model, inference network, (probabilistic) encoder, ...
- ELBO:
$$\mathcal{L}(\theta, \phi; x) = \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z)]}_{\text{Reconstruction}} + \underbrace{\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x) || p(z))}_{\text{Divergence from prior}}$$

Diagram illustrating the ELBO components:

 - Reconstruction:** $\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)]$ (highlighted in yellow)
 - Divergence from prior:** $-\text{KL}(q_{\phi}(z|x) || p(z))$ (highlighted in orange)
 - Other terms:** $\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z)]$ and $\text{KL}(q_{\phi}(z|x) || p(z))$ are also shown.

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 Guassians has
an analytic form)

Variational Auto-Encoders (VAEs)

- ELBO:

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

$$\nabla_\phi \mathcal{L} =$$

$$\nabla_\theta \mathcal{L} =$$

Variational Auto-Encoders (VAEs)

- ELBO:

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

$$\nabla_\phi \mathcal{L} =$$

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

Variational Auto-Encoders (VAEs)

- ELBO:

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

$$\nabla_\phi \mathcal{L} =$$

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

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

Variational Auto-Encoders (VAEs)

- ELBO:

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

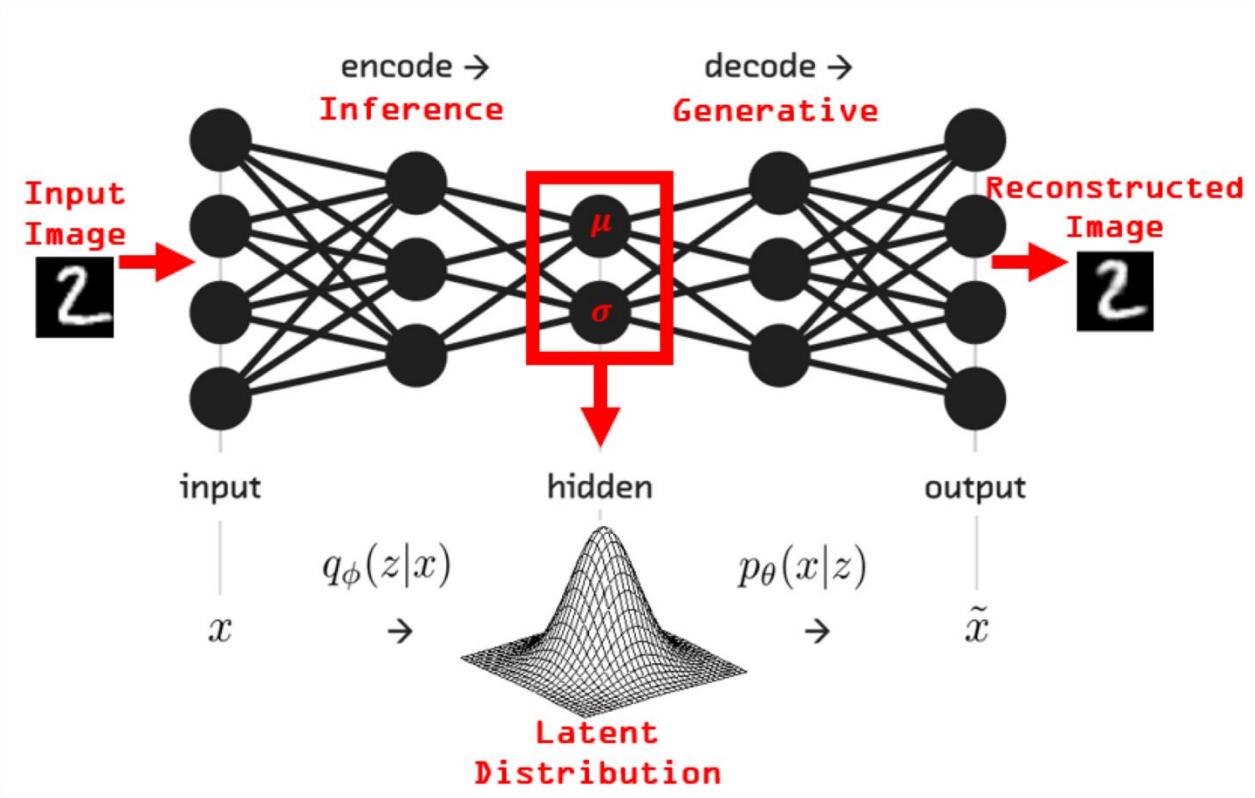
$$\nabla_\phi \mathcal{L} = \mathbb{E}_{\epsilon \sim N(0,1)} [\nabla_z [\log p_\theta(x, z) - \log q_\phi(z|x)] \nabla_\phi z(\epsilon, \phi)]$$

- Reparameterization:

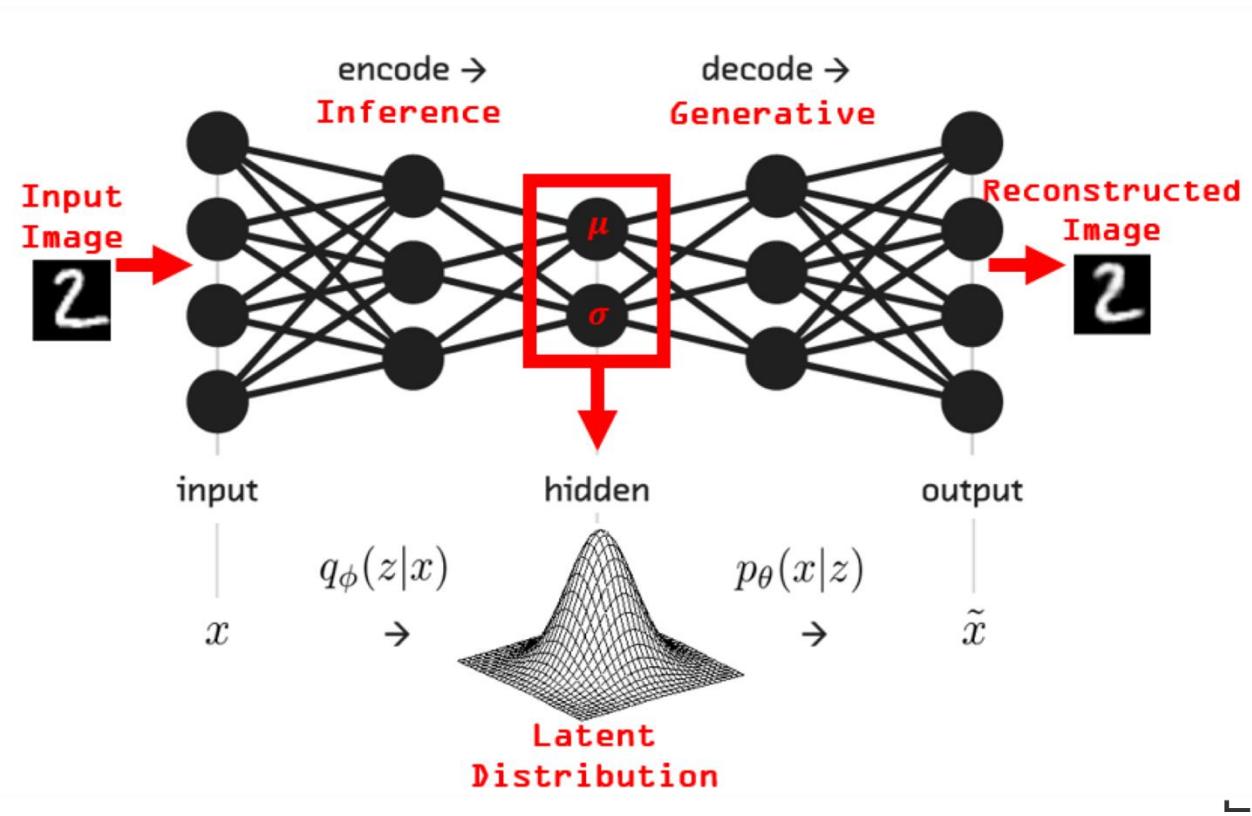
- $[\mu; \sigma] = f_\phi(x)$ (a neural network)
- $z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(0, 1)$

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

Example: VAEs for images



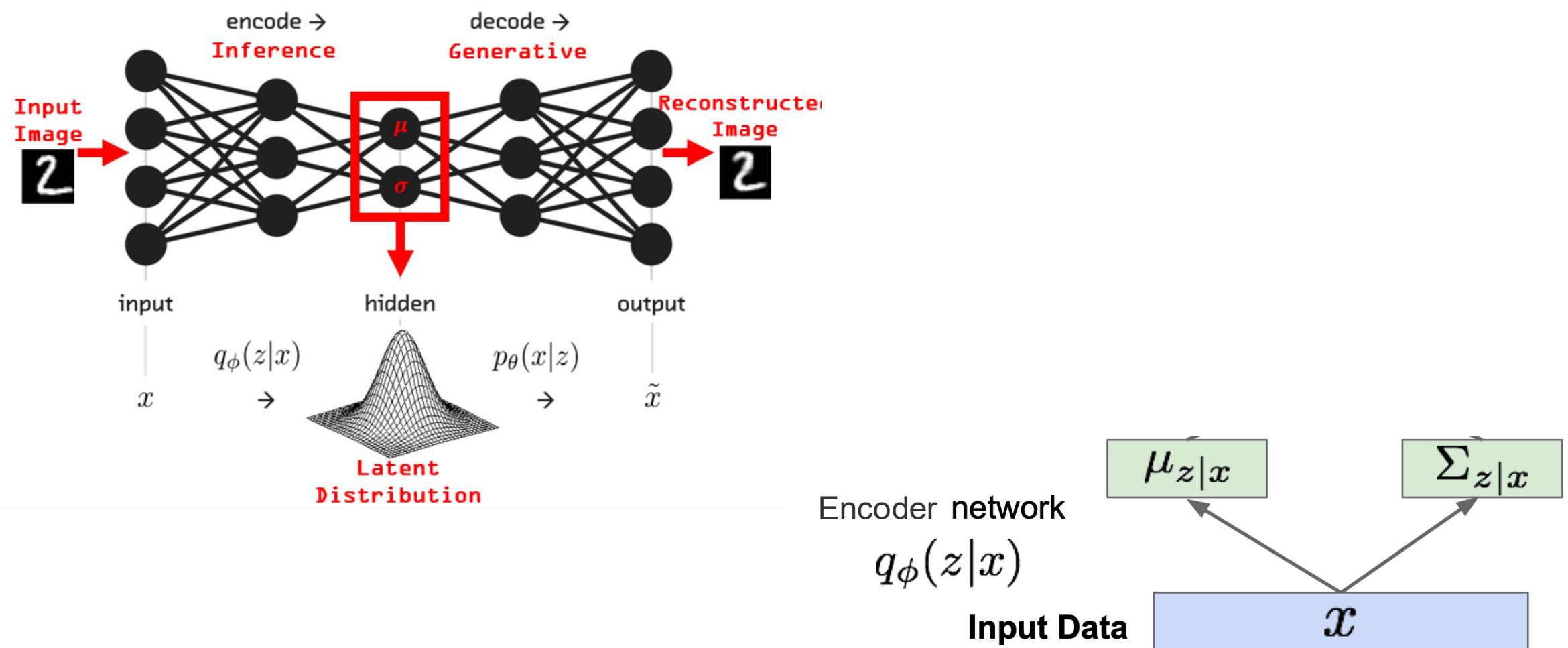
Example: VAEs for images



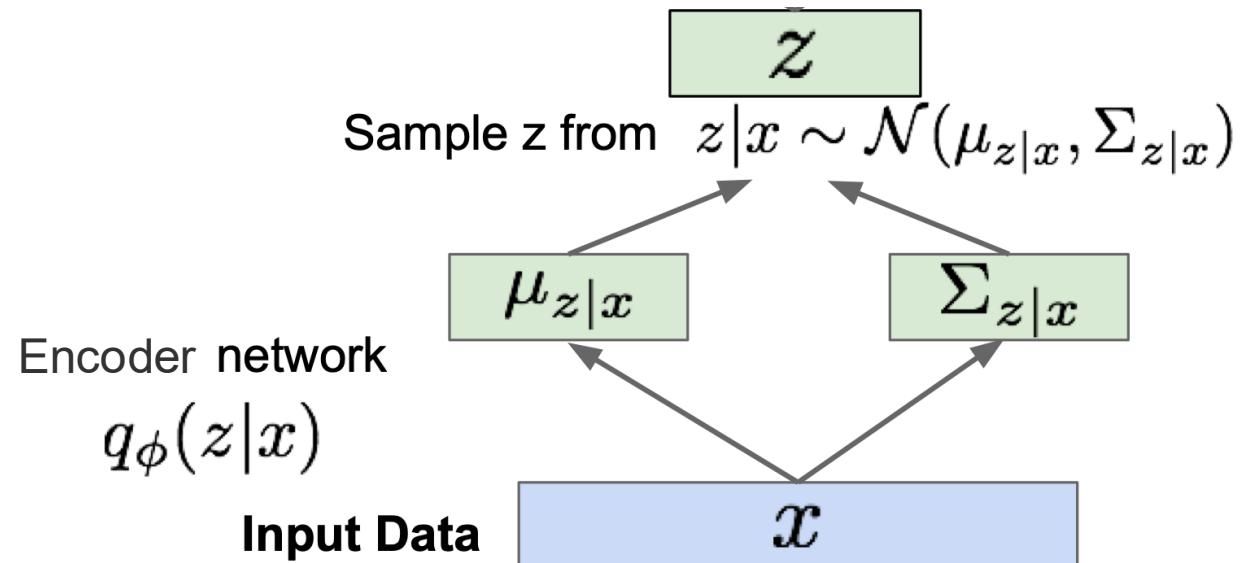
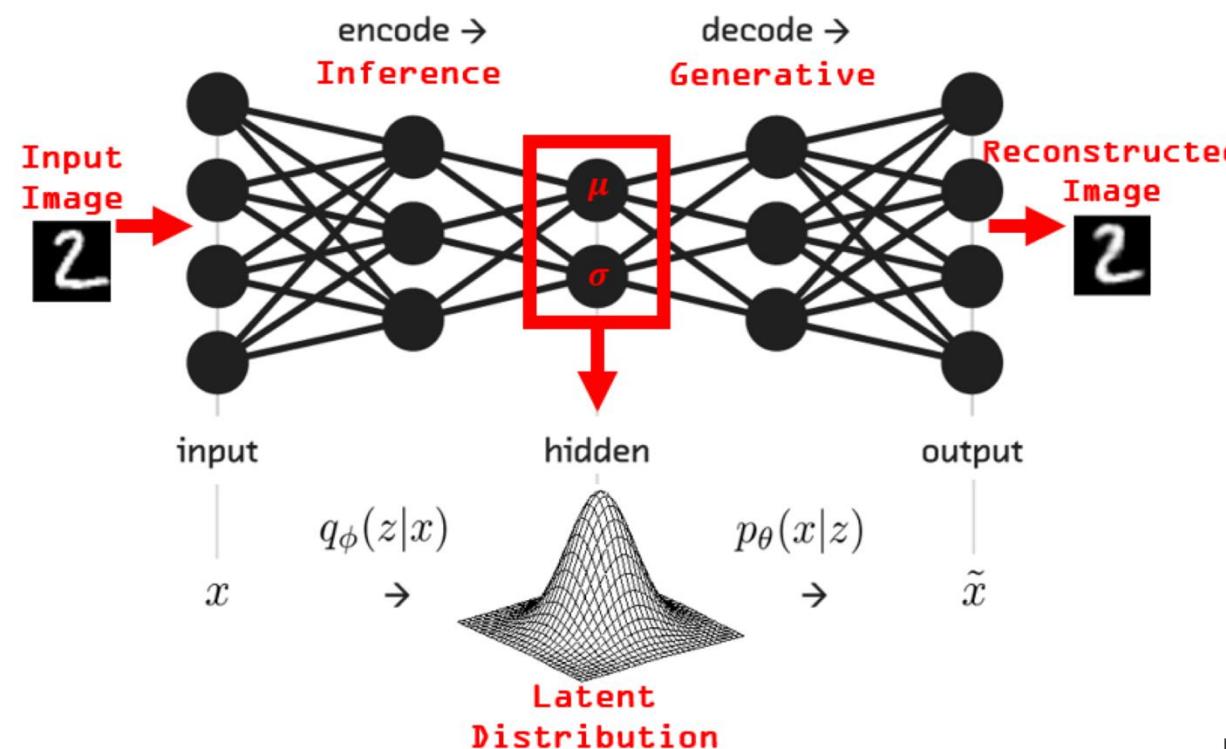
Input Data

x

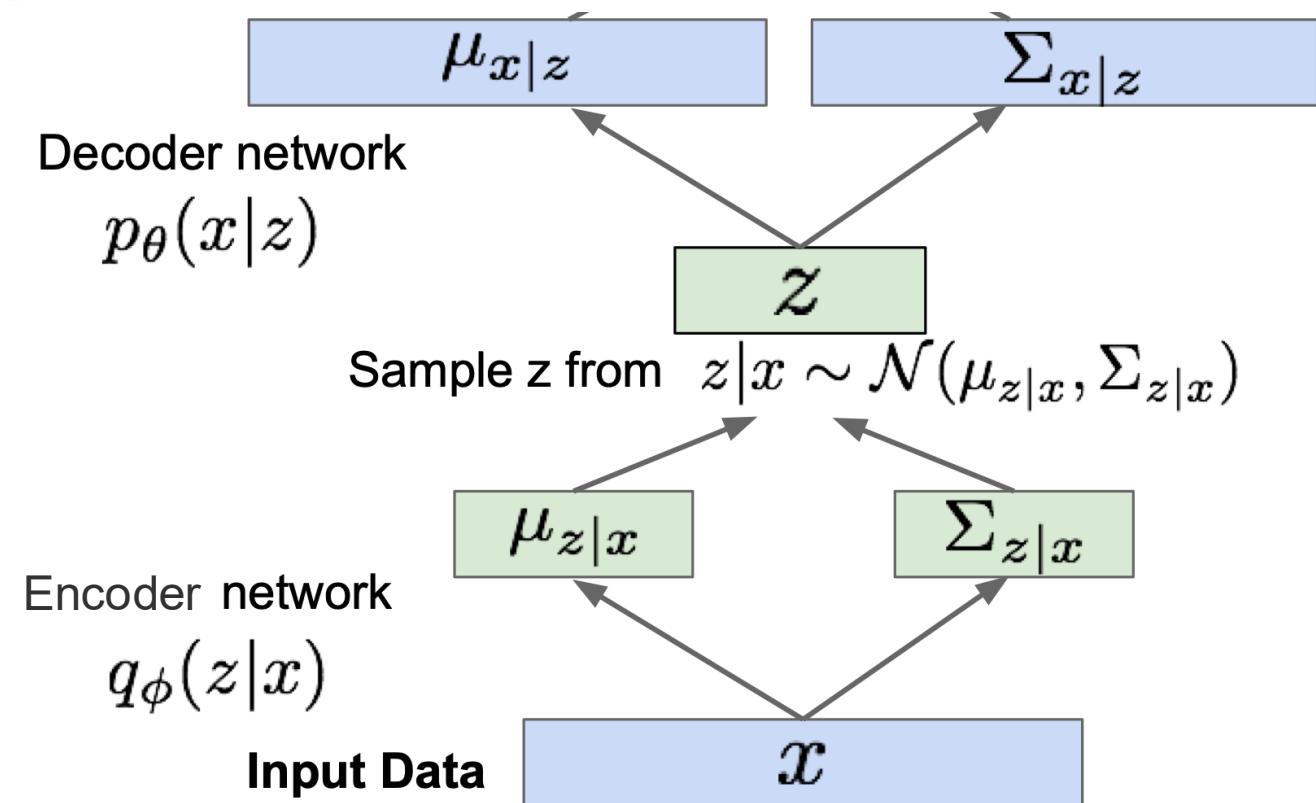
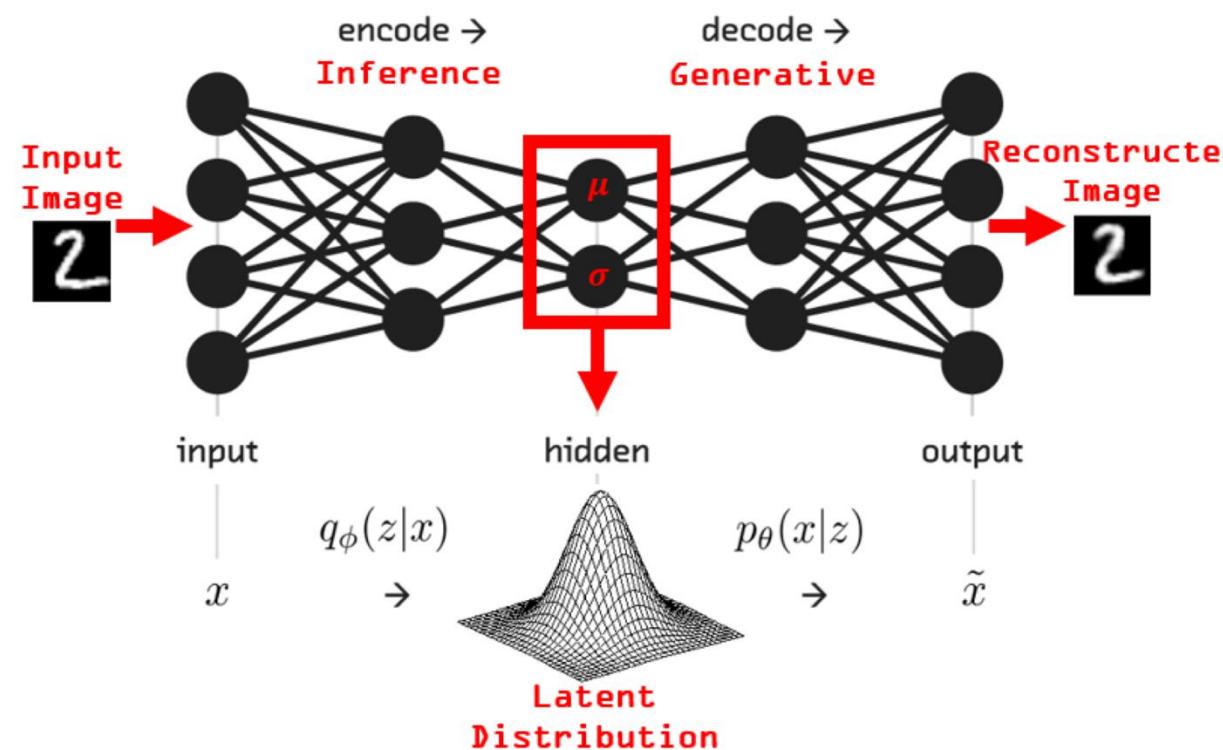
Example: VAEs for images



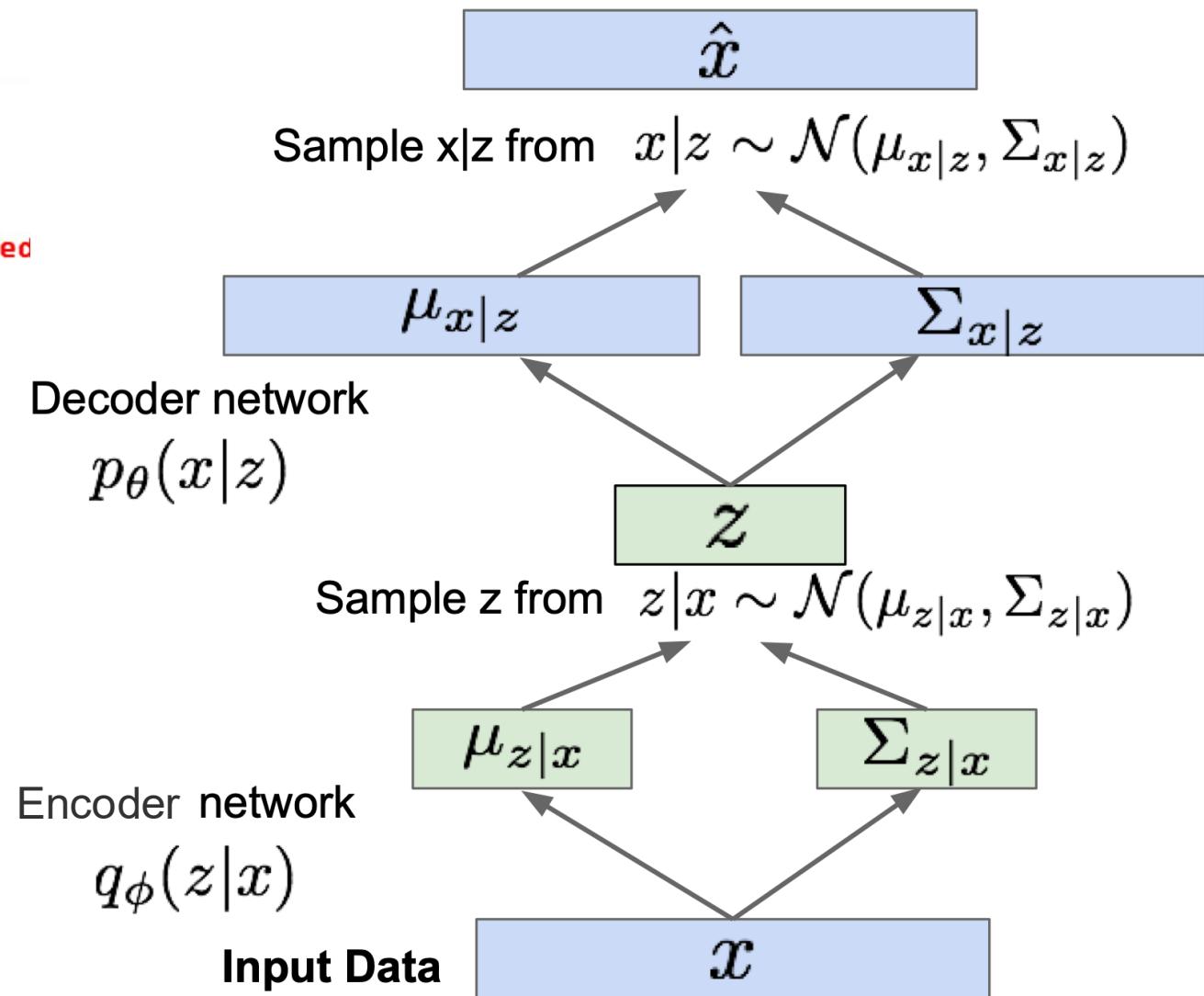
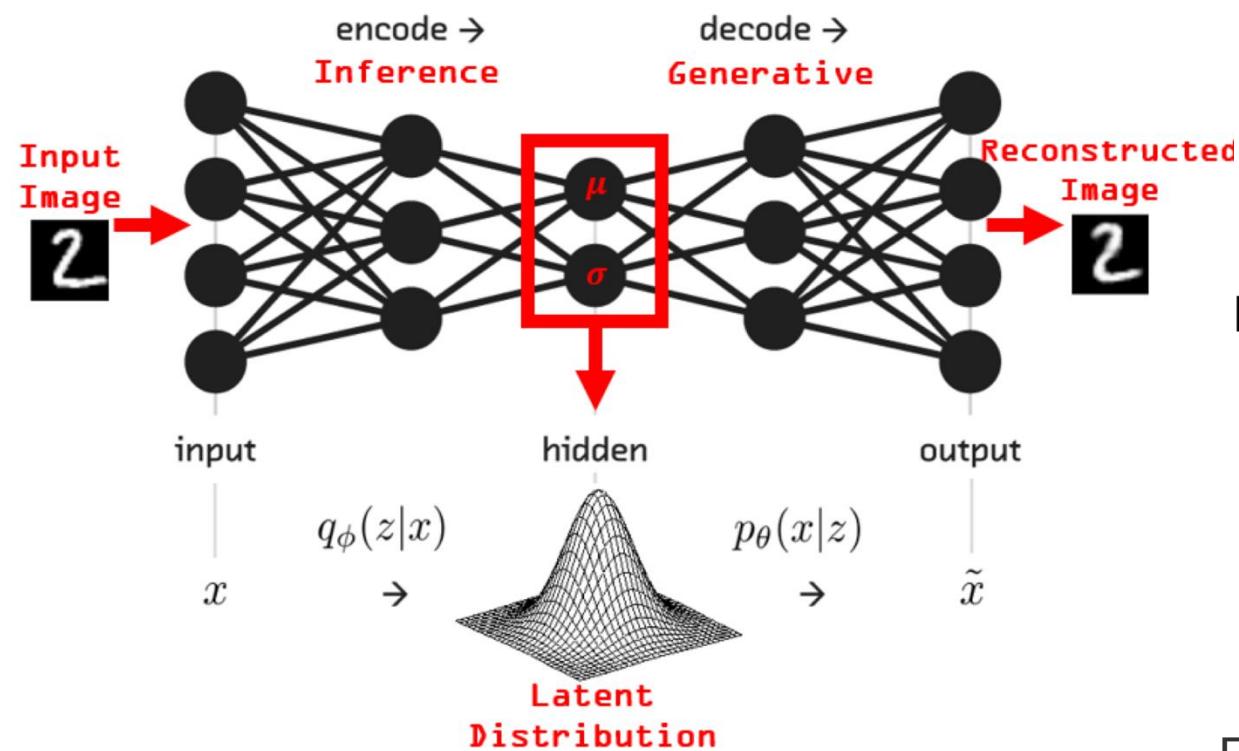
Example: VAEs for images



Example: VAEs for images



Example: VAEs for images

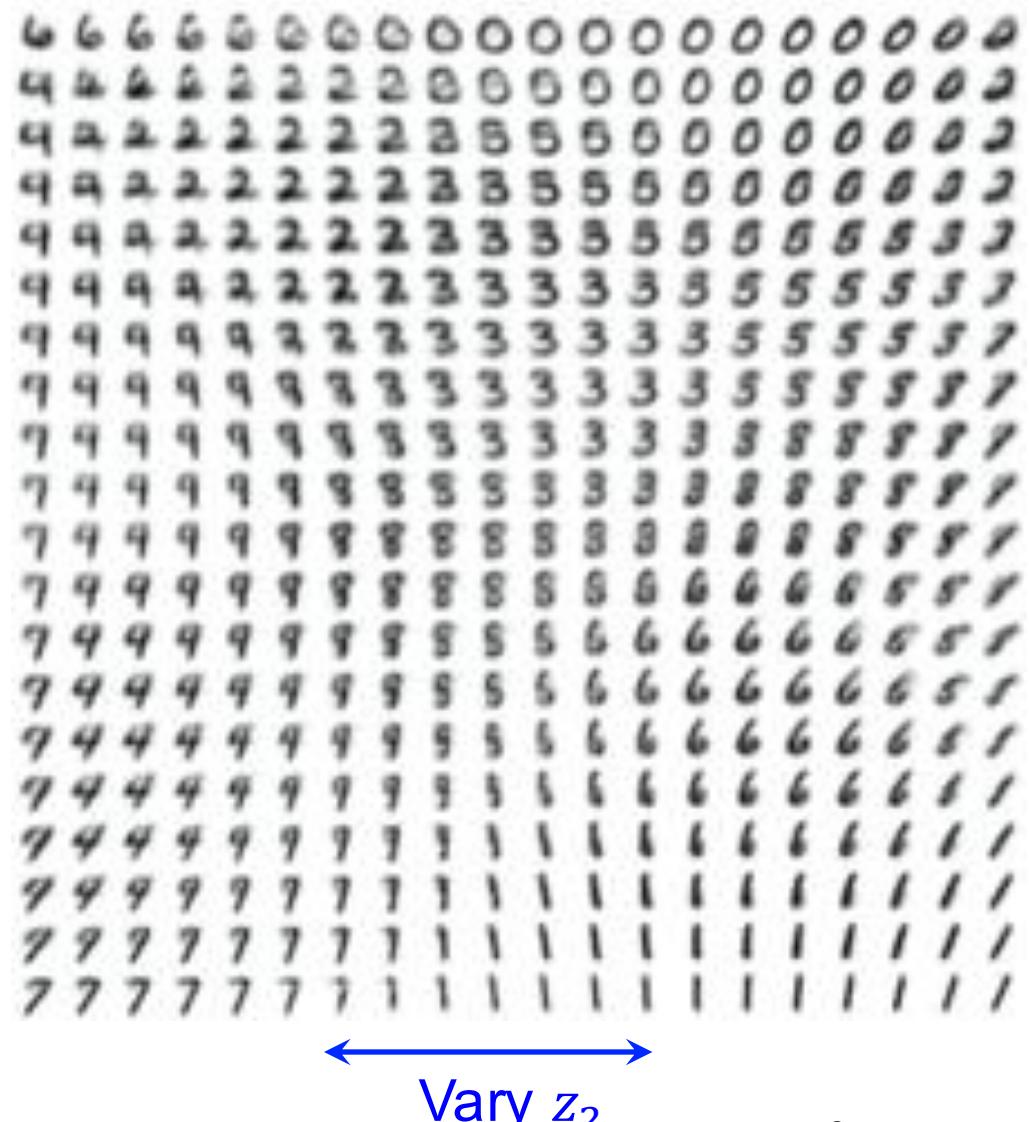
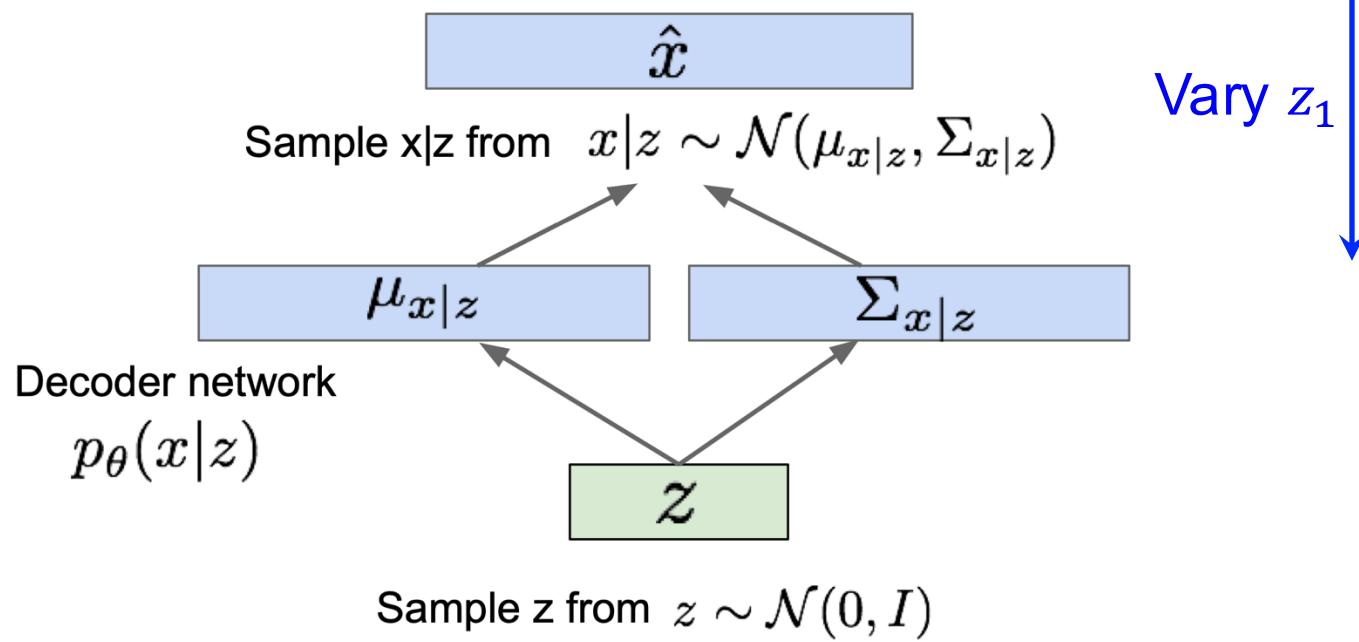


Example: VAEs for images

Data manifold for 2-d z

Generating samples:

- Use decoder network. Now sample z from prior!

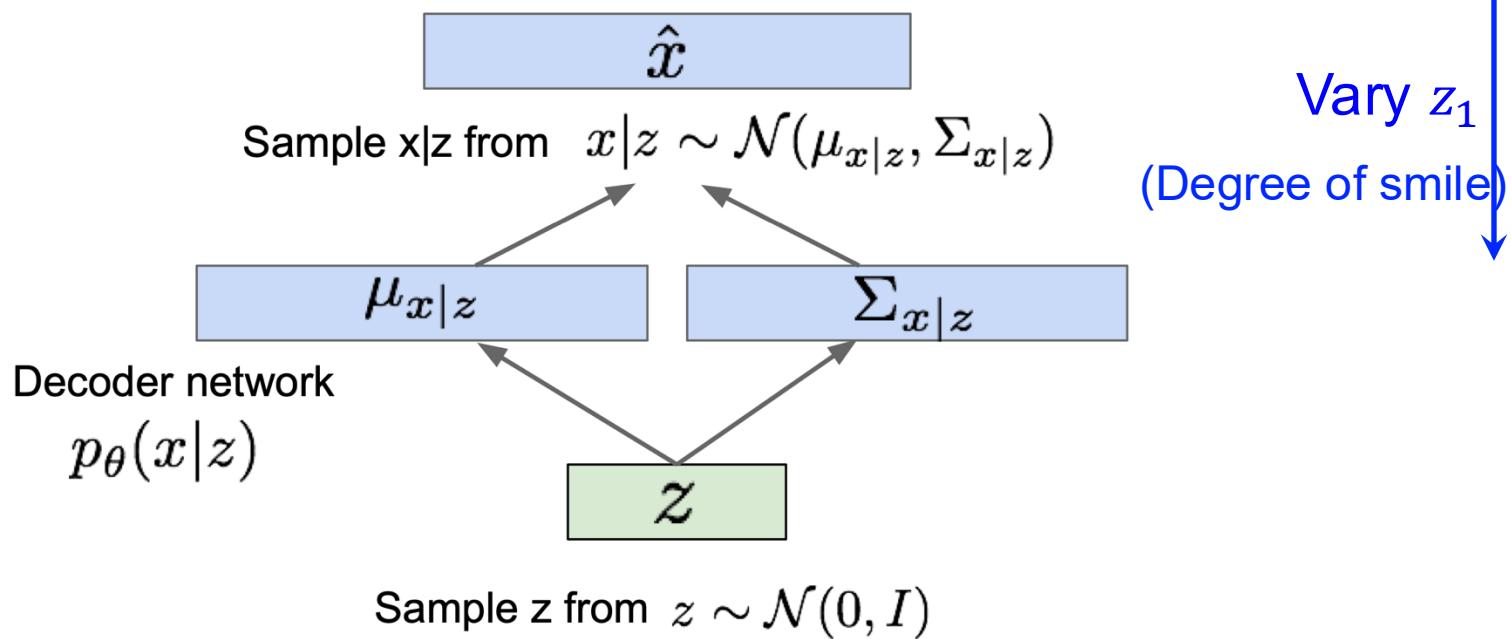


Example: VAEs for images

Data manifold for 2-d z

Generating samples:

- Use decoder network. Now sample z from prior!



Vary z_2 (head pose)
19

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: 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

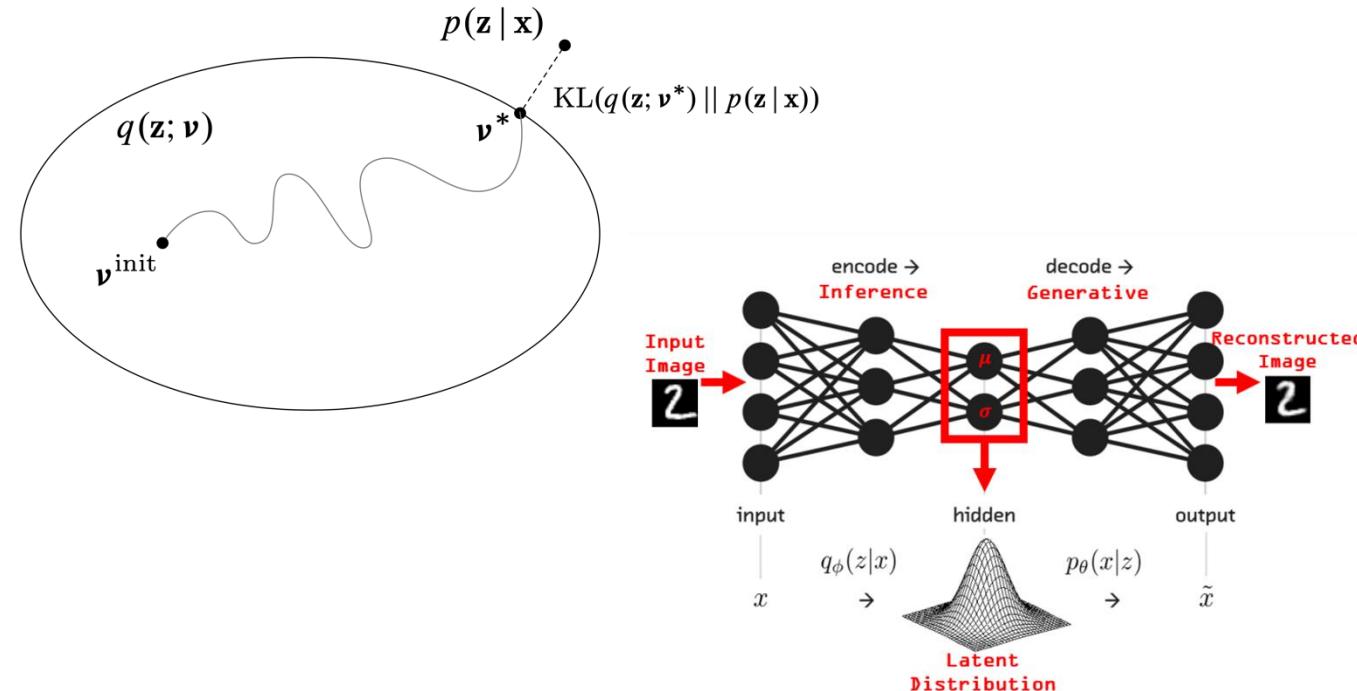


(Razavi et al., 2019)

- 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)



Questions?