

# DSC291: Machine Learning with Few Labels

## Unsupervised Learning

Zhitong Hu

Lecture 5, April 15, 2025

# EM Algorithm: Quick Recap

- Observed variables  $x$ , latent variables  $z$
- To learn a model  $p(x, z|\theta)$ , we want to maximize the marginal log-likelihood

- But it's too difficult  $\ell(\theta; x) = \log p(x|\theta) = \log \sum_z p(x, z|\theta)$

- EM algorithm:
  - maximize a lower bound of  $\ell(\theta; x)$
- Key equation:

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

Variational free energy

# EM Algorithm: Quick Recap

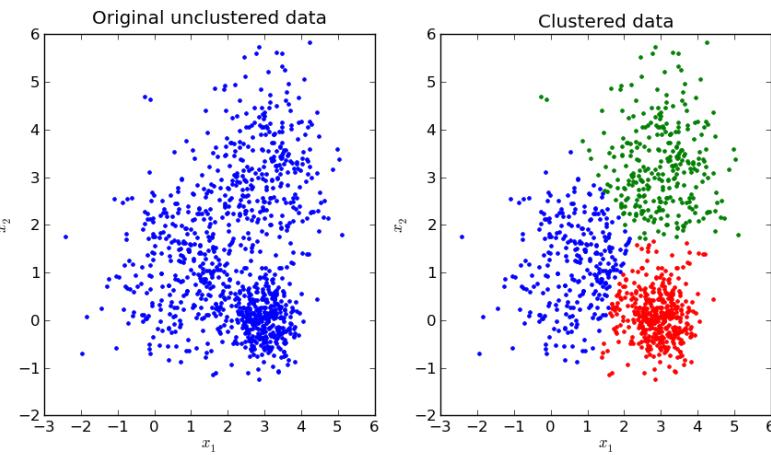
- The EM algorithm is coordinate-decent on  $F(q, \theta)$ 
  - E-step:  $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$ 
    - the posterior distribution over the latent variables given the data and the current parameters
  - M-step:  $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta^t) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

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

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

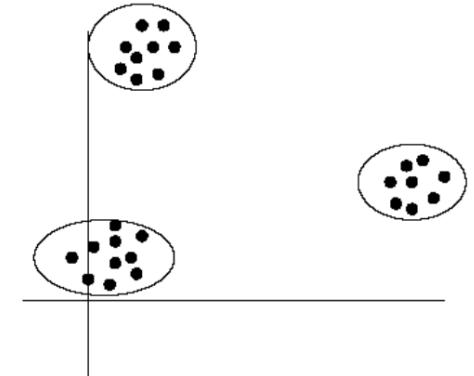
# Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$ 
  - Consider a mixture of  $K$  Gaussian components



# Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$ 
  - Consider a mixture of  $K$  Gaussian components



# Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$ 
  - Consider a mixture of  $K$  Gaussian components
  - $p(x, z|\theta) =$
  - $p(x|\theta) =$
  - The log likelihood of a sample  $x_n \in \{x_i\}_{i=1}^N$ :  
 $\log p(x_n|\theta) =$

$$q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$$

## Example: Gaussian Mixture Models

- Observed variables  $\mathbf{x}$ :
- Latent variables  $\mathbf{z}$ :
- Want to learn a model  $p(\mathbf{x}, \mathbf{z}|\theta)$
- Now we apply EM algorithm for learning the model:
  - E-step: computing the posterior of  $\mathbf{z}_n$  given the current estimate of the parameters (i.e.,  $\pi, \mu, \Sigma$ )

# Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$
- Now we apply EM algorithm for learning the model:
  - E-step: computing the posterior of  $z_n$  given the current estimate of the parameters (i.e.,  $\pi, \mu, \Sigma$ )

$$p(z_n^k = 1 | x, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} N(x_n, | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} N(x_n, | \mu_i^{(t)}, \Sigma_i^{(t)})}$$

$p(z_n^k = 1, x, \mu^{(t)}, \Sigma^{(t)})$

$p(x, \mu^{(t)}, \Sigma^{(t)})$

$$\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta^t) = \operatorname{argmax}_{\theta} \sum_z q^{t+1}(z|x) \log p(x, z|\theta)$$

## Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$
- Now we apply EM algorithm for learning the model:
  - M-step: computing the parameters given the current estimate of  $z_n$

$$\mathbb{E}_{q(z|x)} [\log p(x, z|\theta)] =$$

$$\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta^t) = \operatorname{argmax}_{\theta} \sum_z q^{t+1}(z|x) \log p(x, z|\theta)$$

## Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$
- Now we apply EM algorithm for learning the model:
  - M-step: computing the parameters given the current estimate of  $z_n$

$$\begin{aligned}\mathbb{E}_{q(z|x)} [\log p(x, z|\theta)] &= \sum_n \mathbb{E}_q [\log p(z_n | \pi)] + \sum_n \mathbb{E}_q [\log p(x_n | z_n, \mu, \Sigma)] \\ &= \sum_n \sum_k \mathbb{E}_q [z_n^k] \log \pi_k - \frac{1}{2} \sum_n \sum_k \mathbb{E}_q [z_n^k] \left( (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log |\Sigma_k| + C \right)\end{aligned}$$

$$\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta^t) = \operatorname{argmax}_{\theta} \sum_z q^{t+1}(z|x) \log p(x, z|\theta)$$

## Example: Gaussian Mixture Models

- Observed variables  $x$ :
- Latent variables  $z$ :
- Want to learn a model  $p(x, z|\theta)$
- Now we apply EM algorithm for learning the model:
  - M-step: computing the parameters given the current estimate of  $z_n$

$$\begin{aligned}\pi_k^* &= \arg \max \langle l_c(\theta) \rangle, & \Rightarrow \frac{\partial}{\partial \pi_k} \langle l_c(\theta) \rangle &= 0, \forall k, \quad \text{s.t. } \sum_k \pi_k = 1 \\ && \Rightarrow \pi_k^* &= \left. \sum_n \langle z_n^k \rangle_{q^{(t)}} \right/ N = \left. \sum_n \tau_n^{k(t)} \right/ N = \left. \langle n_k \rangle \right/ N\end{aligned}$$

$$\mu_k^* = \arg \max \langle l(\theta) \rangle, \quad \Rightarrow \mu_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} x_n}{\sum_n \tau_n^{k(t)}}$$

$$\Sigma_k^* = \arg \max \langle l(\theta) \rangle, \quad \Rightarrow \Sigma_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} (x_n - \mu_k^{(t+1)}) (x_n - \mu_k^{(t+1)})^T}{\sum_n \tau_n^{k(t)}}$$

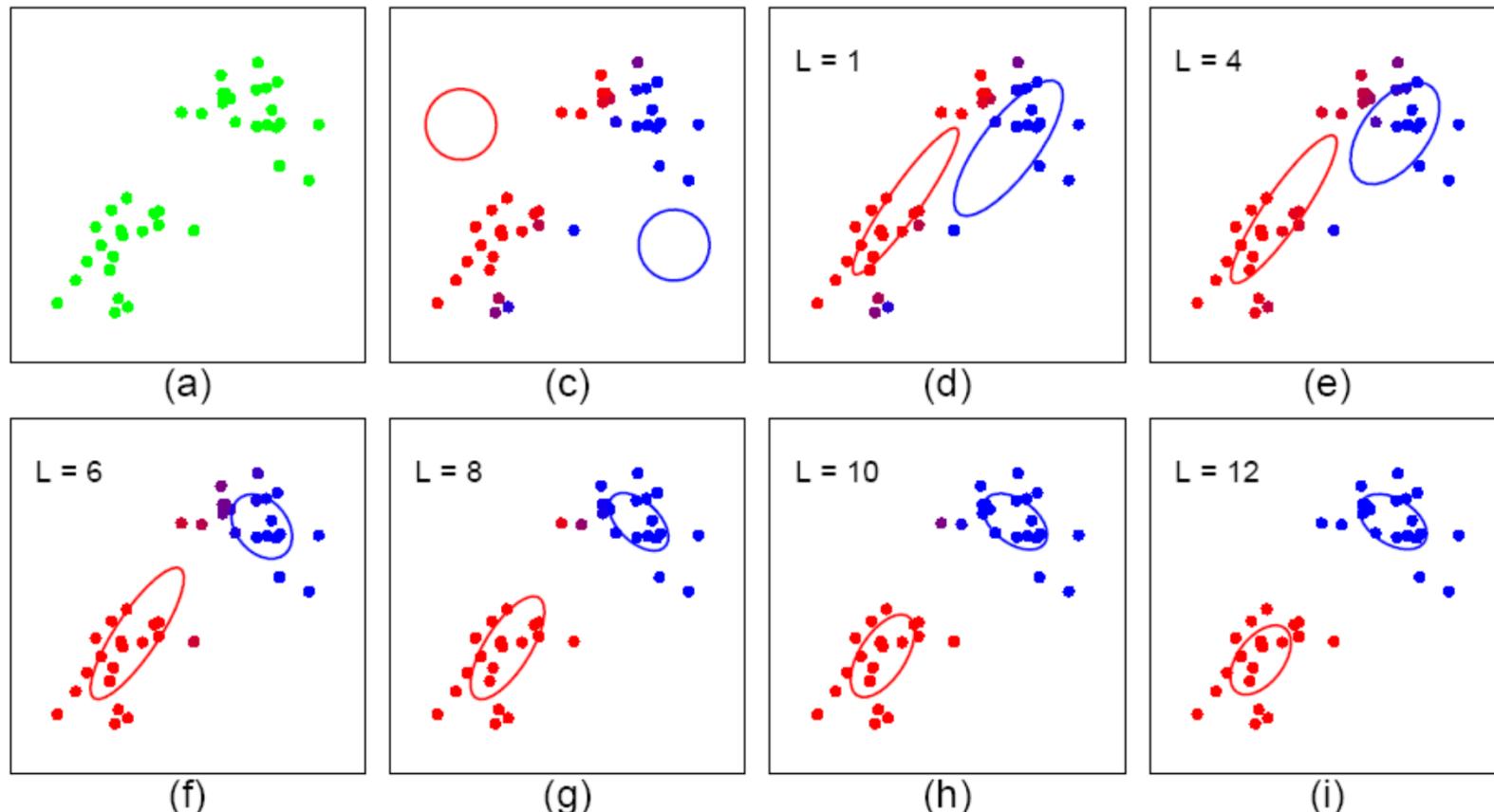
Fact:

$$\frac{\partial \log |A^{-1}|}{\partial A^{-1}} = A^T$$

$$\frac{\partial \mathbf{x}^T A \mathbf{x}}{\partial A} = \mathbf{x} \mathbf{x}^T$$

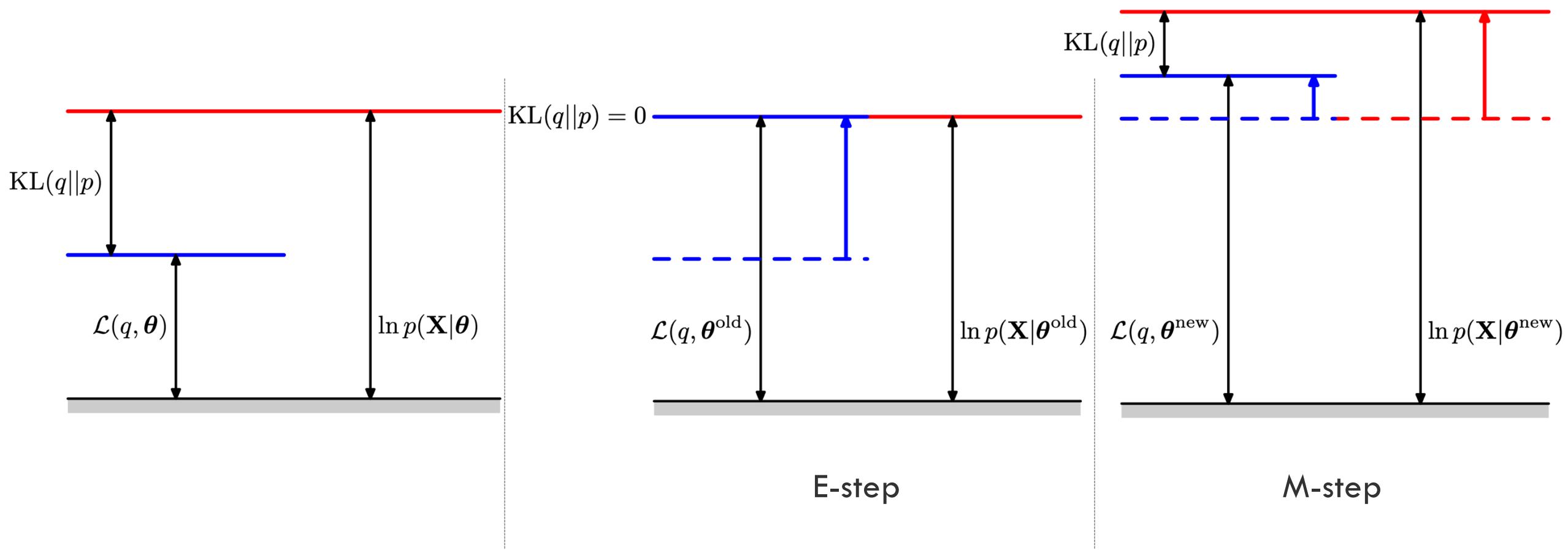
# Example: Gaussian Mixture Models (GMMs)

- Start: “guess” the centroid  $\mu_k$  and covariance  $\Sigma_k$  of each of the K clusters
- Loop:



# Each EM iteration guarantees to improve the likelihood

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



# Summary: EM Algorithm

- A way of maximizing likelihood function for latent variable models. Finds MLE of parameters when the original (hard) problem can be broken up into two (easy) pieces
  - Estimate some “missing” or “unobserved” data from observed data and current parameters.
  - Using this “complete” data, find the maximum likelihood parameter estimates.
- Alternate between filling in the latent variables using the best guess (posterior) and updating the parameters based on this guess:
  - E-step:
  - M-step:  $q^{t+1} = \arg \min_q F(q, \theta^t)$   
 $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$

# **Variational Inference / Variational Auto-Encoders (VAEs)**

# Inference

- Given a model, the goals of inference can include:
  - Computing the likelihood of observed data  $p(\mathbf{x}^*)$
  - Computing the marginal distribution over a given subset of variables in the model  $p(\mathbf{x}_A)$
  - Computing the conditional distribution over a subsets of nodes given a disjoint subset of nodes  $p(\mathbf{x}_A | \mathbf{x}_B)$
  - Computing a mode of the density (for the above distributions)  $\text{argmax}_{\mathbf{x}} p(\mathbf{x})$
  - ....

# Variational Inference

- Observed variables  $x$ , latent variables  $z$
- Variational (Bayesian) inference, a.k.a. **variational Bayes**, is most often used to **approximately** infer the **posterior distribution** over the latent variables

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

- We cannot directly compute the posterior distribution for many interesting models
  - I.e. the posterior density is in an intractable form (often involving integrals) which cannot be easily analytically solved.

# EM and Variational Inference

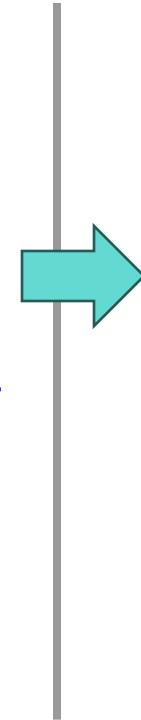
- The EM algorithm:

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



Need to approximate  $p(\mathbf{z}|\mathbf{x}, \theta^t)$   
with VI

# 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.

# 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  $\nu$ , 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.

# 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?
  - **Hint:** recall the equation that holds for any  $q$ :

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

# 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?

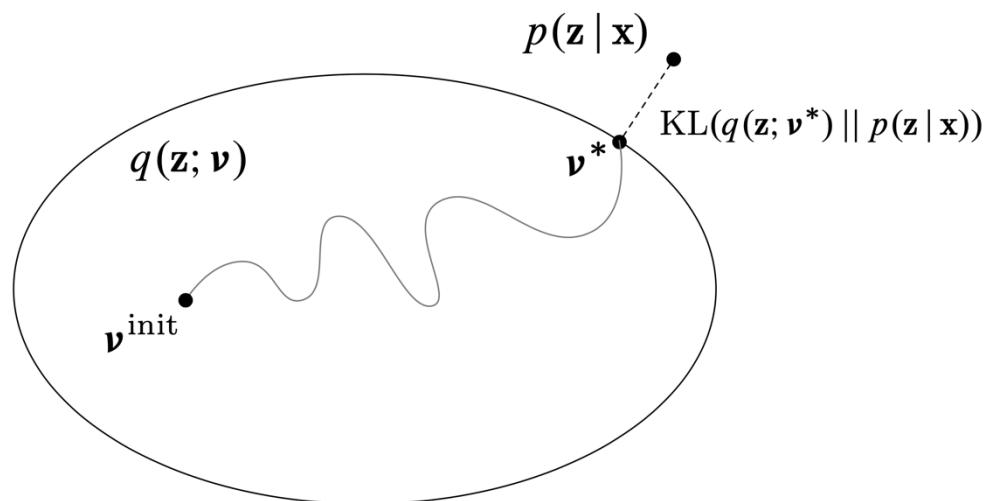
$$\ell(\theta; \mathbf{x}) = \underbrace{\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right]}_{\text{Evidence Lower Bound (ELBO)}} + \text{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}, \theta))$$

- 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  $\nu$ , i.e.  $q(\mathbf{z}|\mathbf{x}, \nu)$
- We maximize the ELBO over  $q$  to find an “optimal approximation” to  $p(\mathbf{z}|\mathbf{x})$

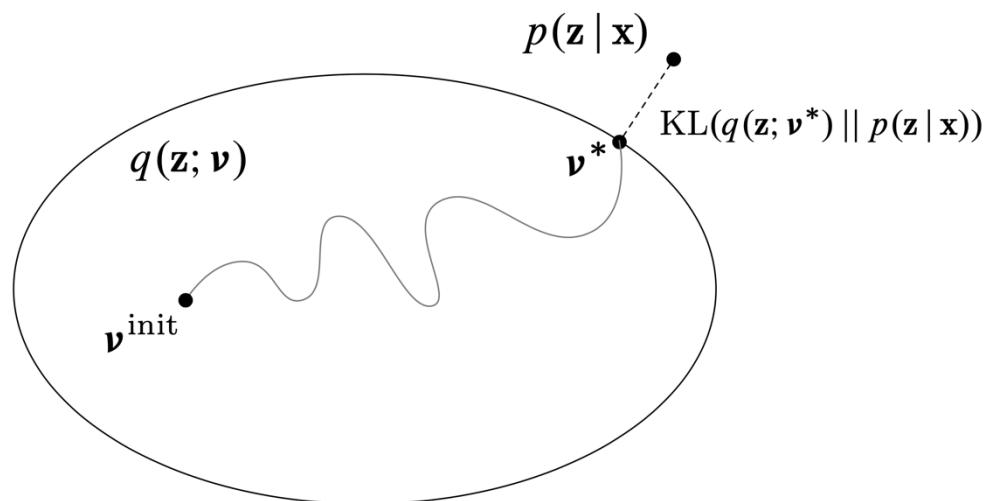
$$\begin{aligned} & \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \nu)} \right] \\ &= \operatorname{argmax}_{\nu} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \nu)} [\log p(\mathbf{x}, \mathbf{z}|\theta)] -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



# Variational Inference

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

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

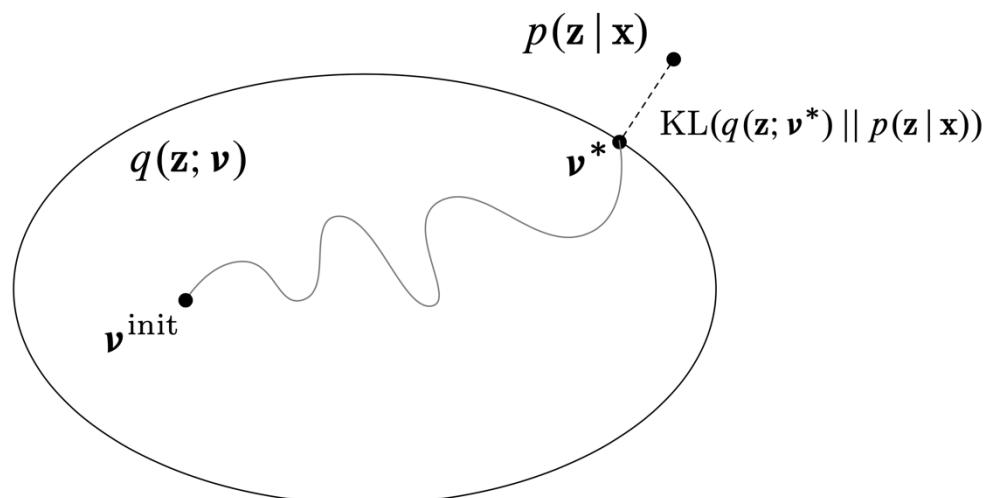


**Question:** How do we choose the variational family  $q(\mathbf{z}|\mathbf{x}, \nu)$ ?

# Variational Inference

- Choose a family of distributions over the latent variables  $\mathbf{z}$  with its own set of variational parameters  $\nu$ , 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)

# 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:

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

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

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

# 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:

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

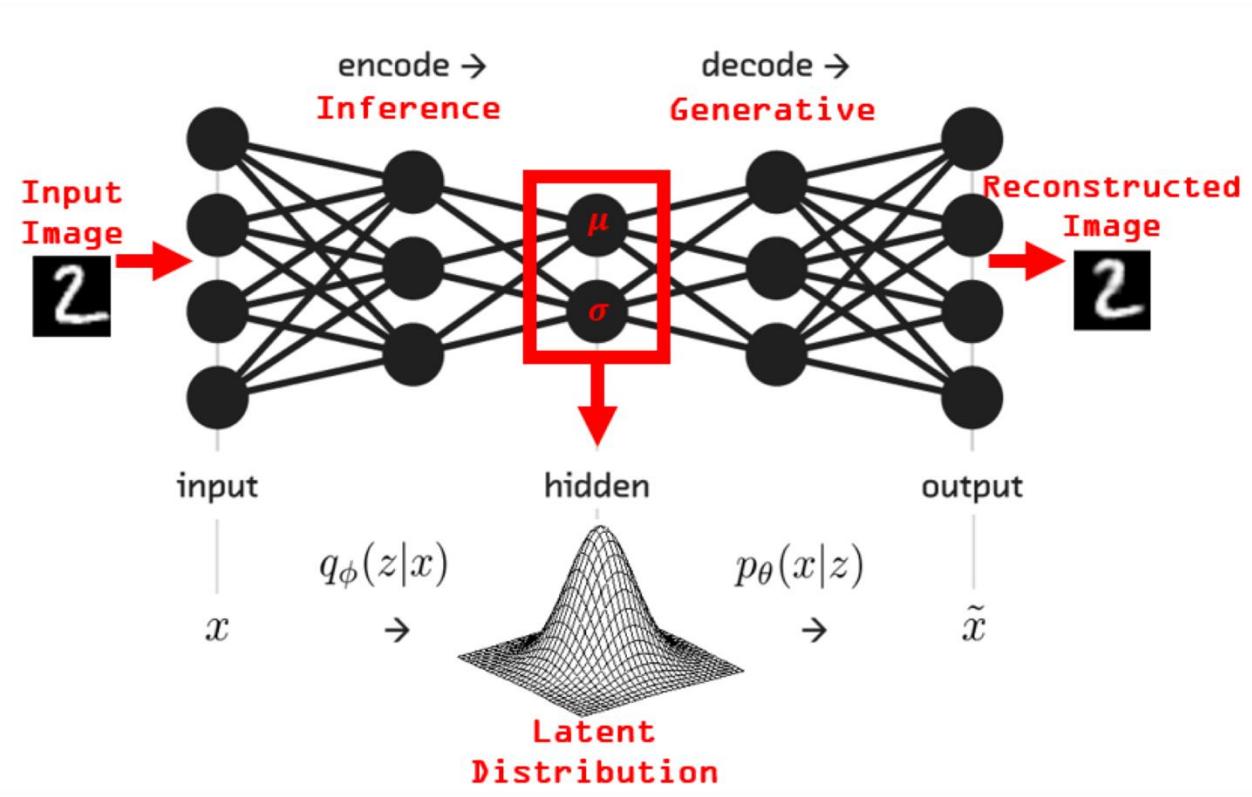
$$\nabla_\phi \mathcal{L} = \mathbb{E}_{\epsilon \sim N(0,1)} [\nabla_z [\log p_\theta(x, z) - \log q_\phi(z|x)] \nabla_\phi z(\epsilon, \phi)]$$

- Reparameterization:

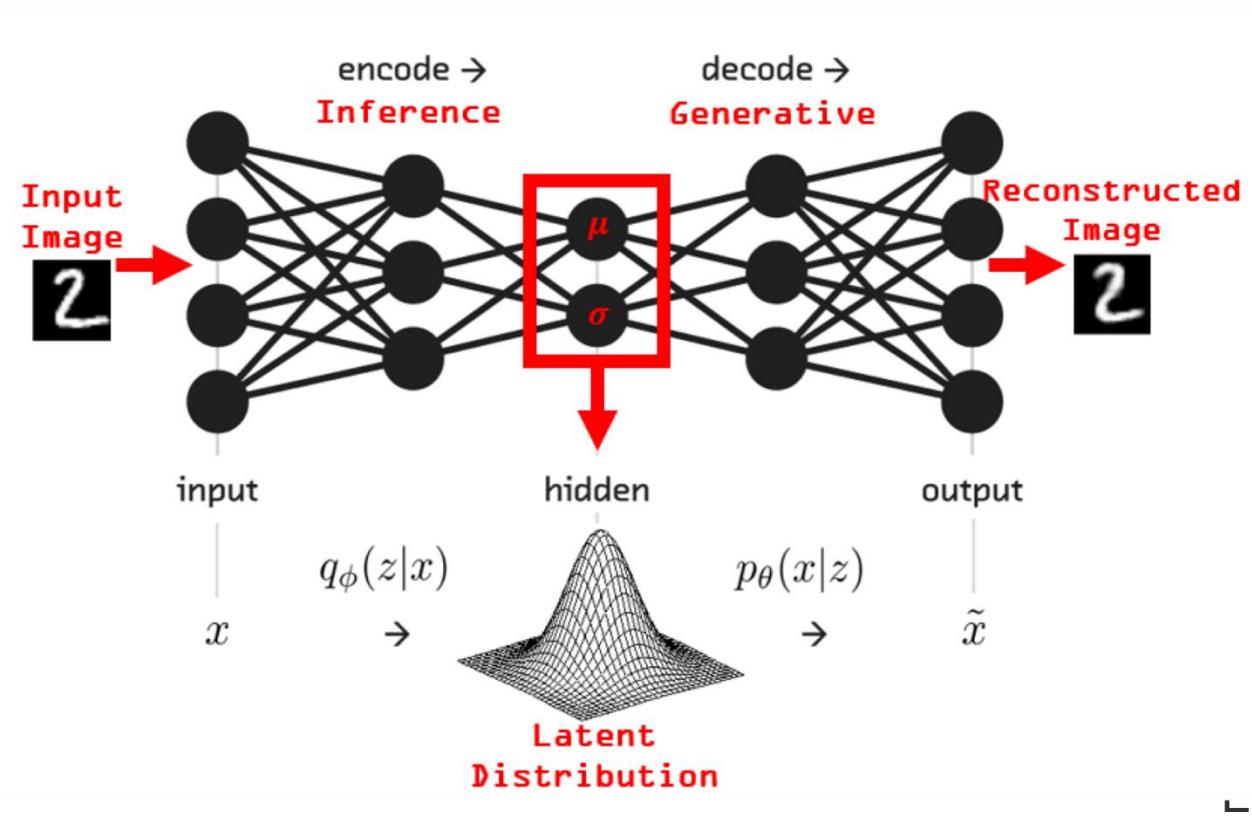
- $[\mu; \sigma] = f_\phi(x)$  (a neural network)
- $z = \mu + \sigma \odot \epsilon, \quad \epsilon \sim N(0, 1)$

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

# Example: VAEs for images



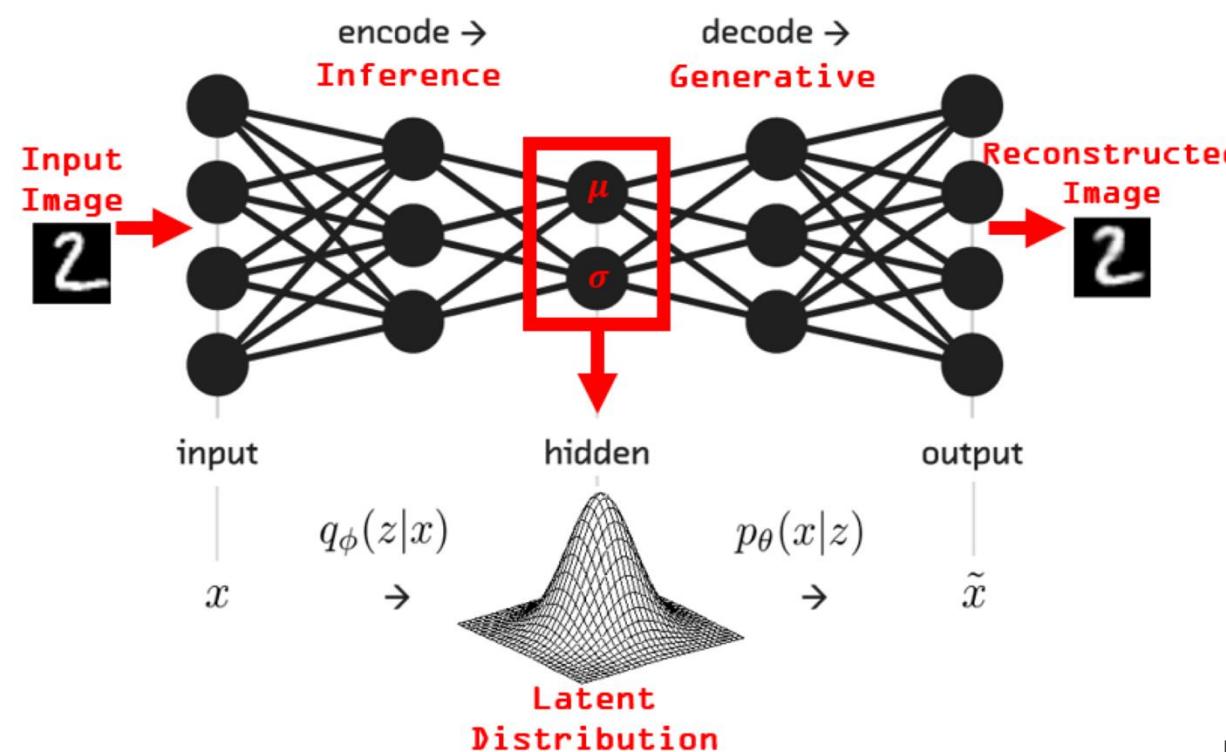
# Example: VAEs for images



Input Data

$x$

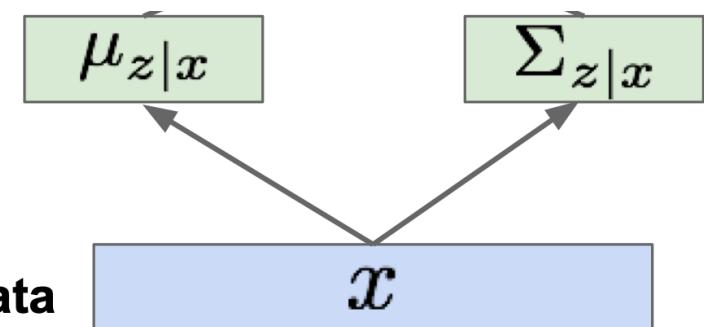
# Example: VAEs for images



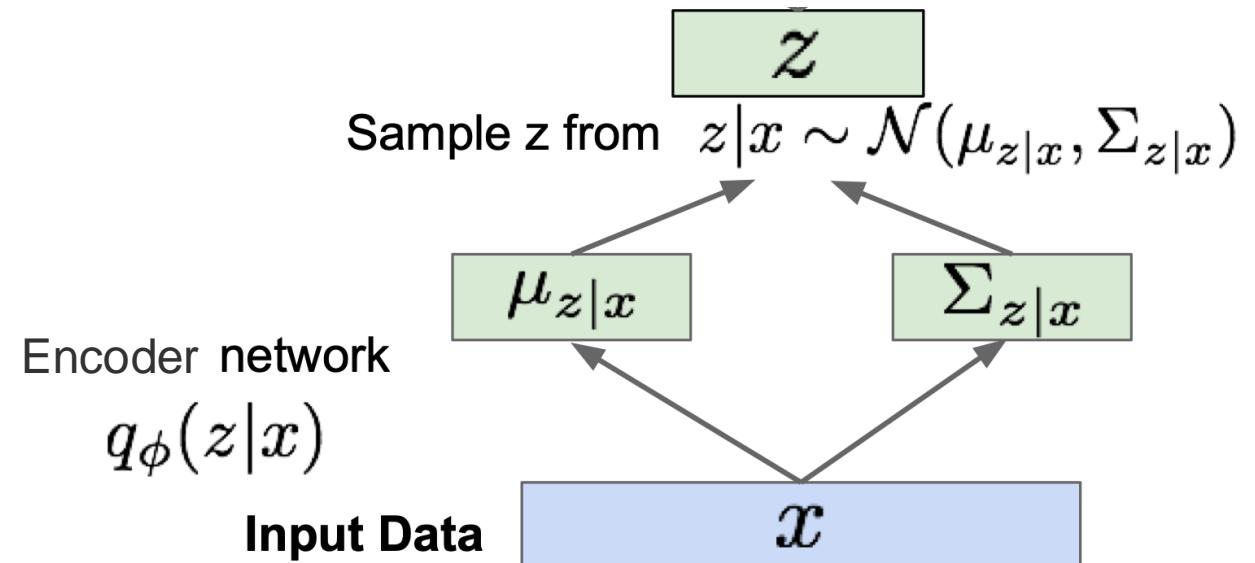
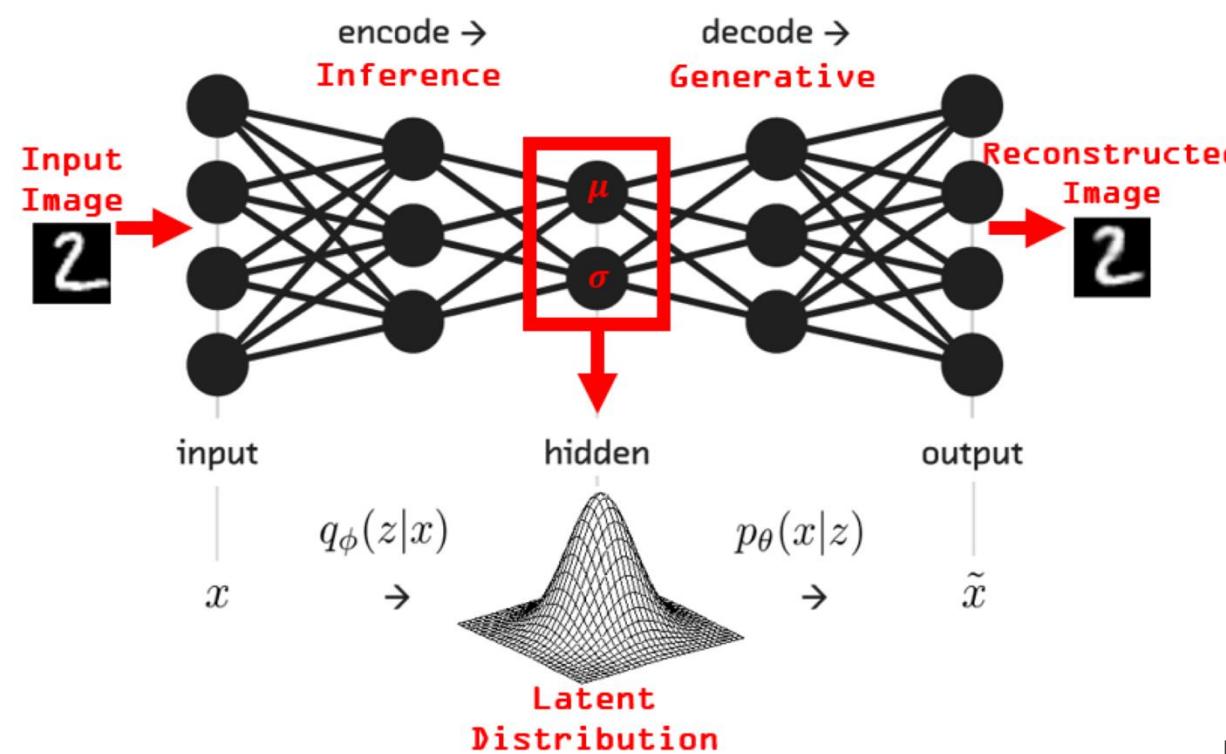
Encoder network

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

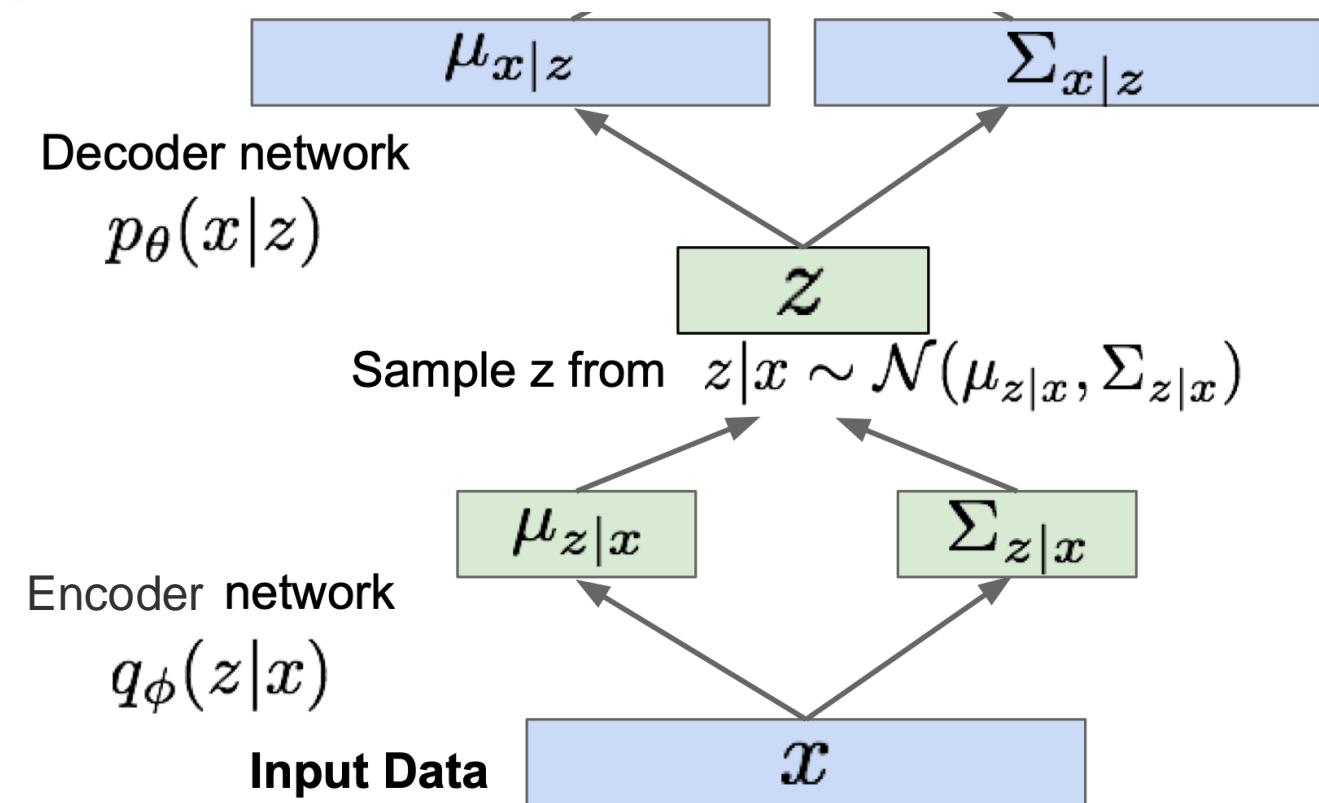
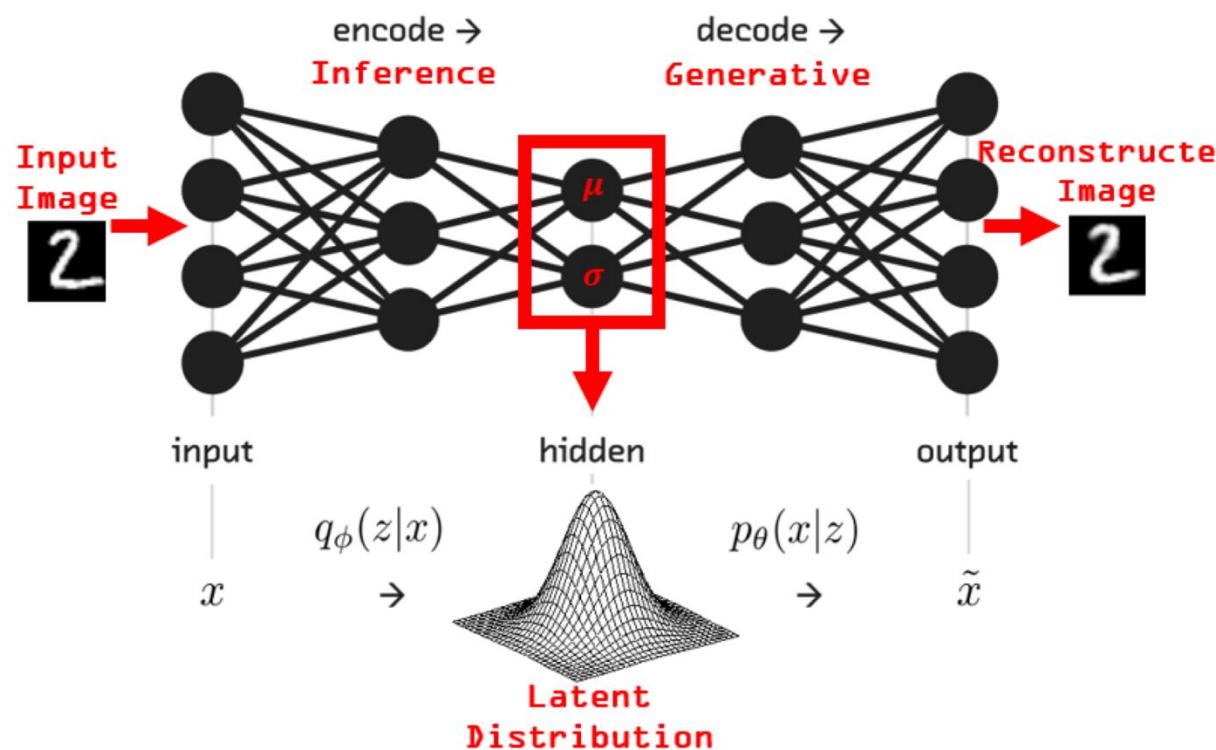
**Input Data**



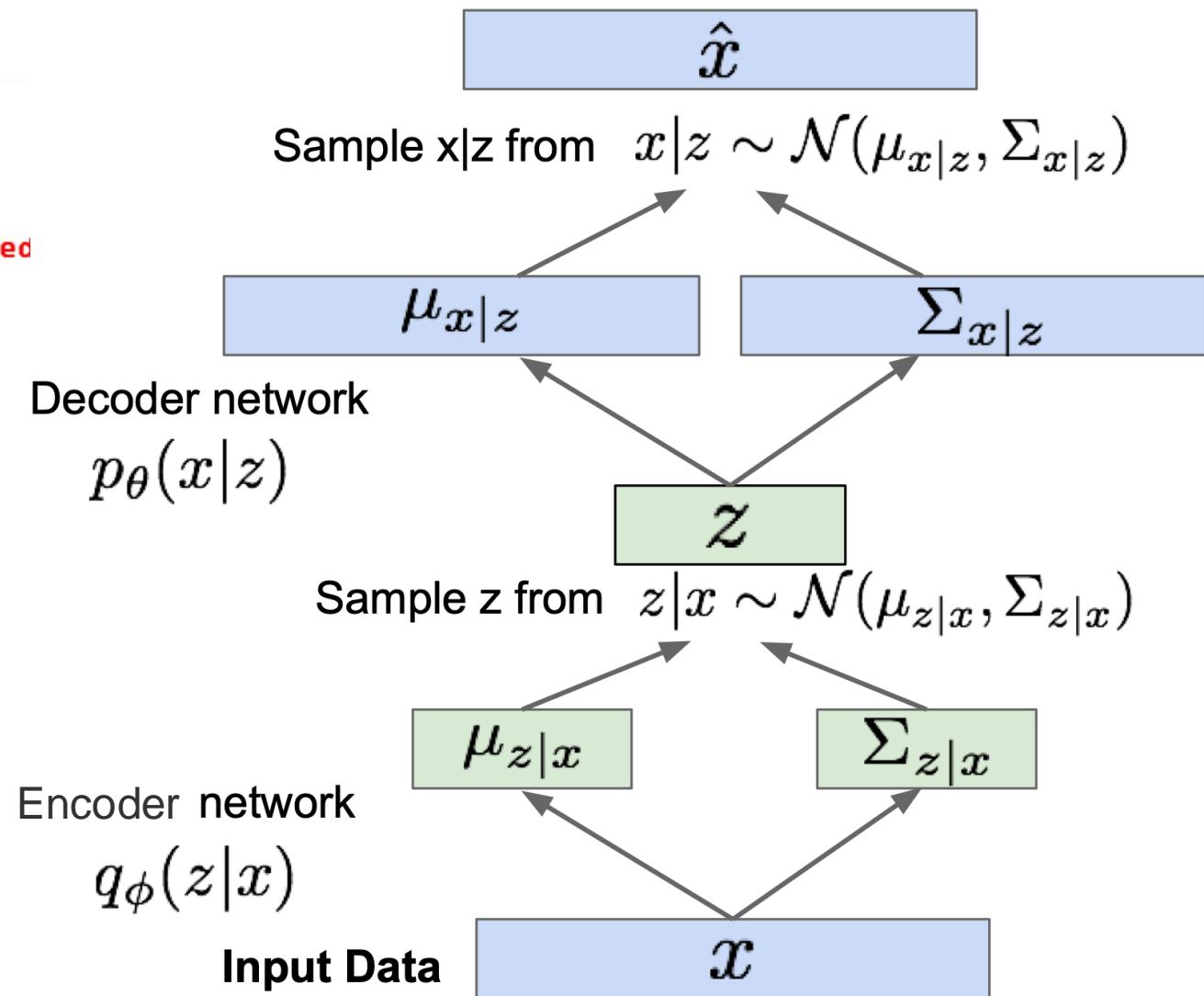
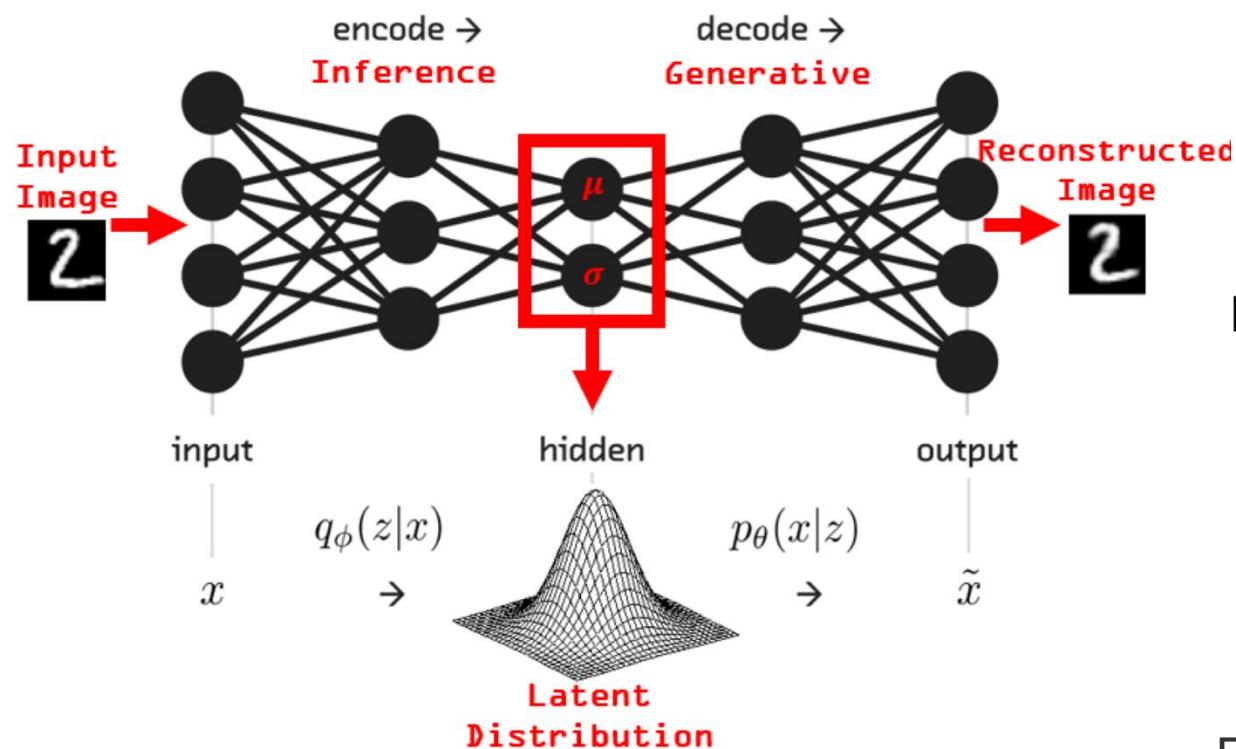
# Example: VAEs for images



# Example: VAEs for images



# Example: VAEs for images

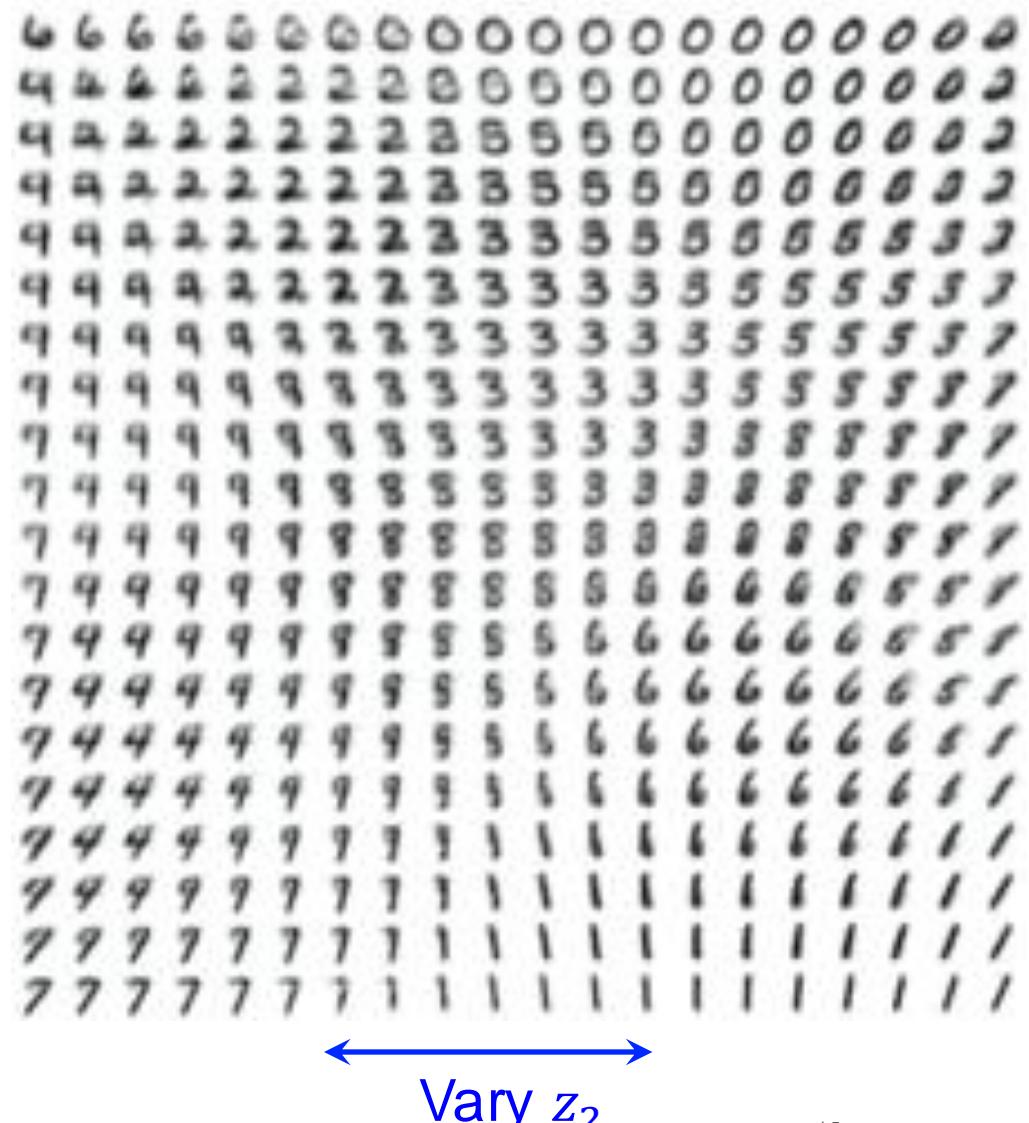
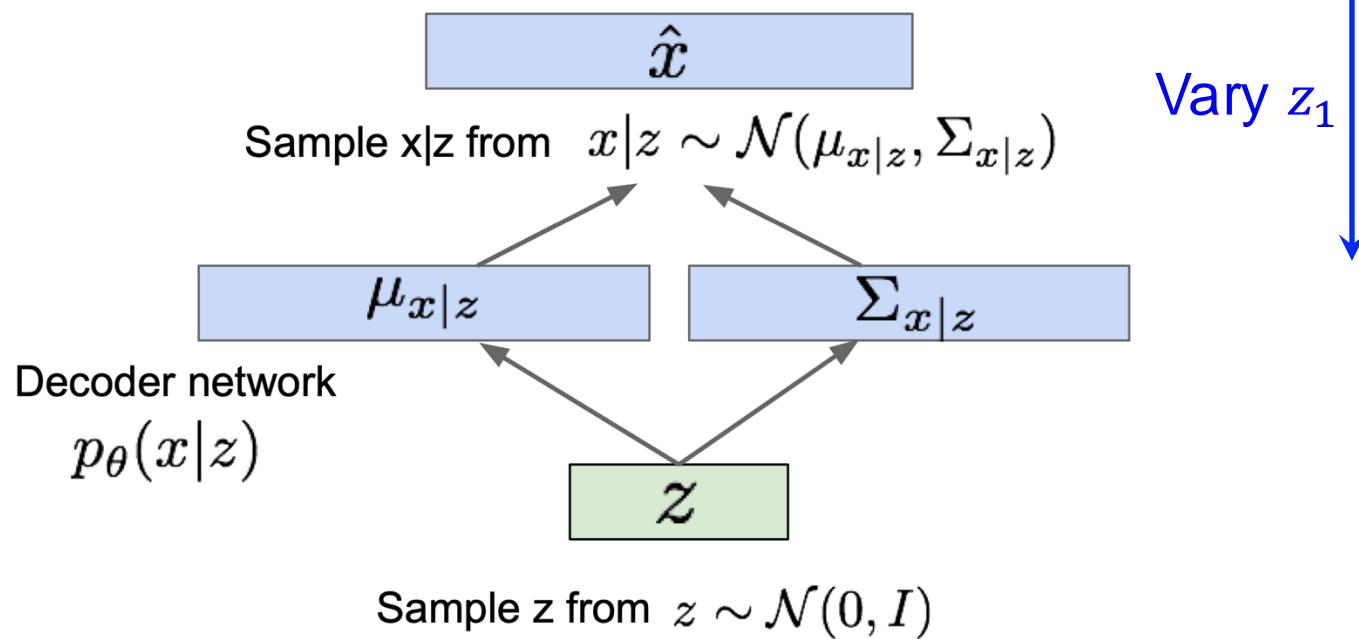


# Example: VAEs for images

Data manifold for 2-d  $z$

Generating samples:

- Use decoder network. Now sample  $z$  from prior!

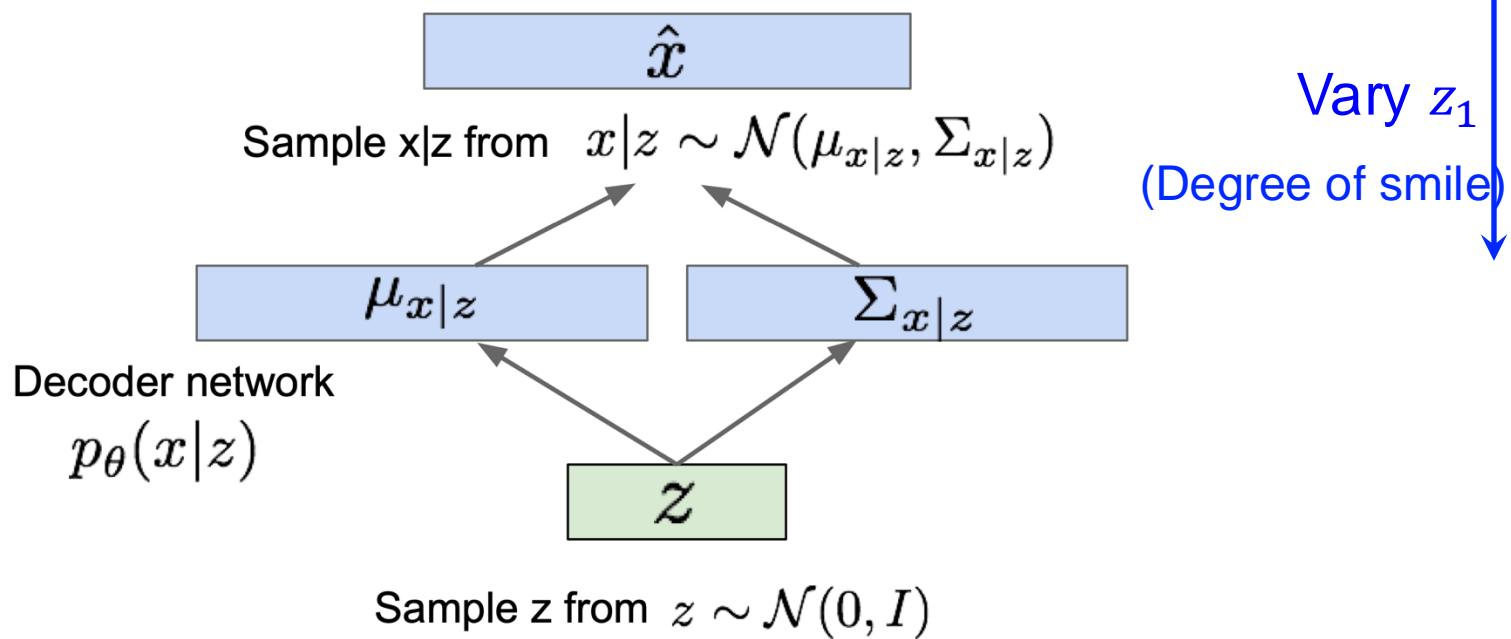


# Example: VAEs for images

Data manifold for 2-d  $z$

Generating samples:

- Use decoder network. Now sample  $z$  from prior!



Vary  $z_2$  (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: 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

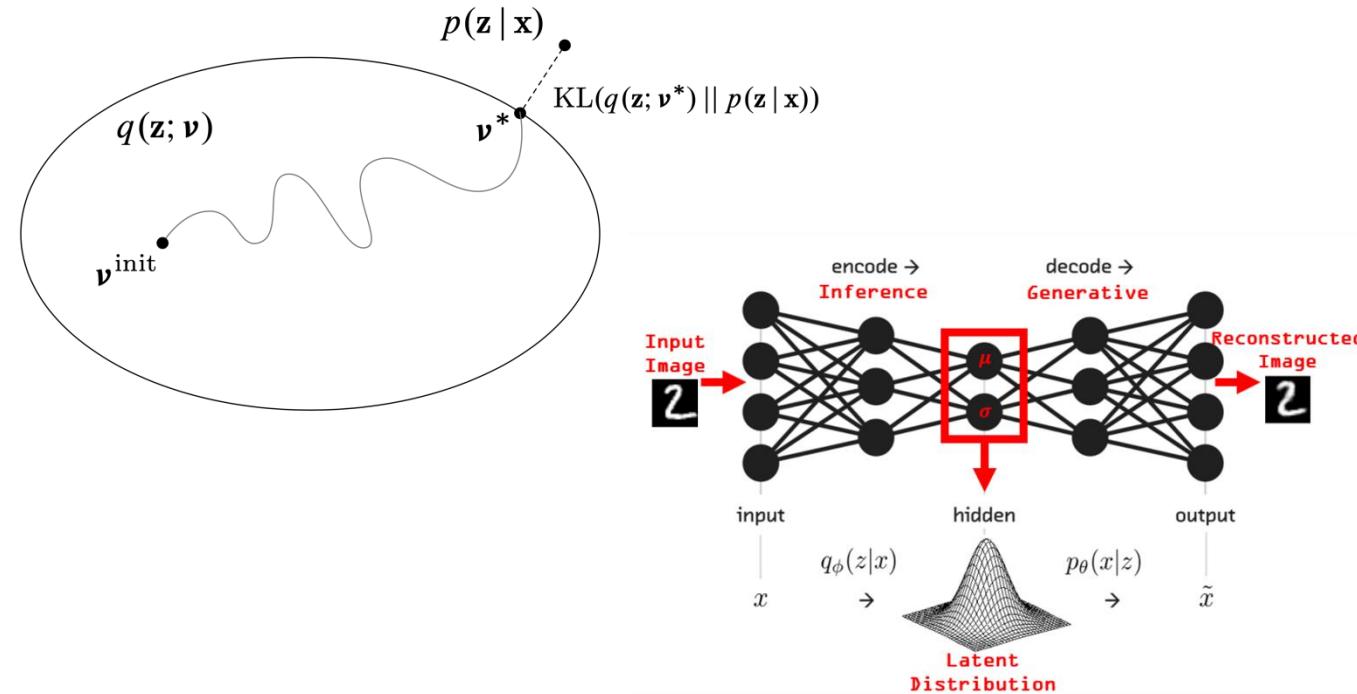


(Razavi et al., 2019)

- 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?