

CS 4803 / 7643: Deep Learning

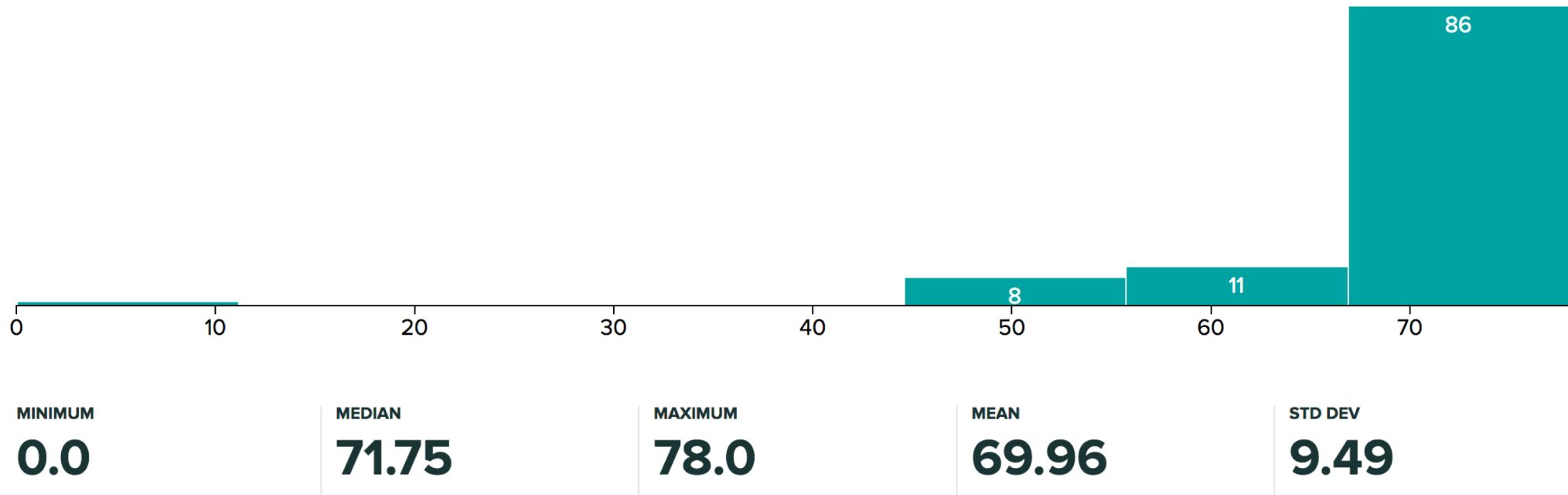
Topics:

- Variational Auto-Encoders (VAEs) |
- Key Ideas
 - AEs, Variational Inference, ELBO, Reparameterization

Dhruv Batra
Georgia Tech

Administrativa


- HW2 Grades Released
 - Max regular points: 54 (4803), 78 (7643)
 - Regrade requests close: 11/20, 11:55pm



Recap from last time

Overview

- Unsupervised Learning
- Generative Models
 - PixelRNN and PixelCNN
 - Variational Autoencoders (VAE)
 - Generative Adversarial Networks (GAN)



Variational Autoencoders (VAE)

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^n p_{\theta}(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent z :

$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

$\sum p(x|z)p(z)$

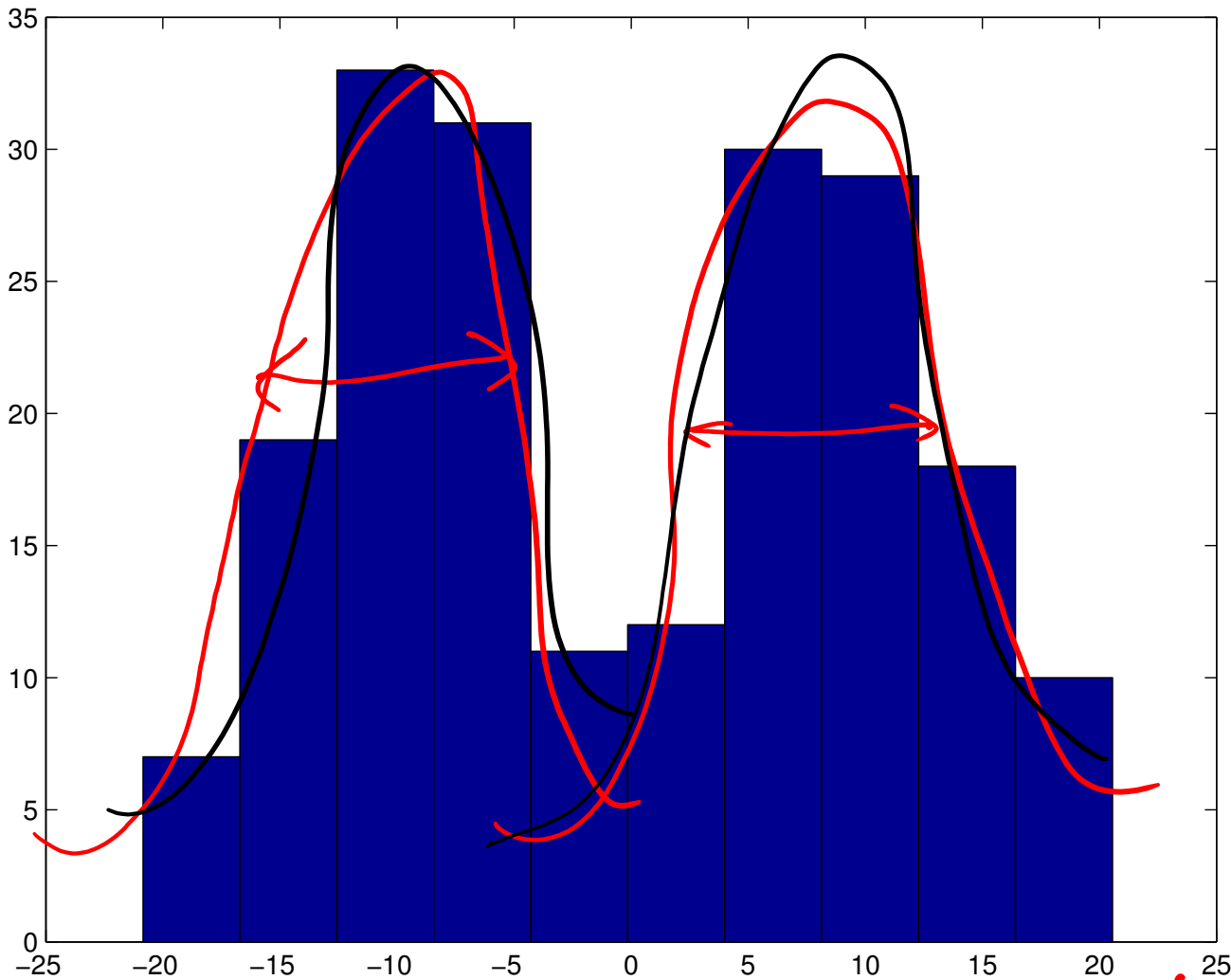
$P(\vec{x})$

$P(\vec{x}, z)$

$= P(\vec{x}|z)P(z)$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

GMM



$p(x)$

$p(x)$

$p(x|z)$

$z \in \{1, 2\}$

$p(z) = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

$p(x|z=1) \sim \mathcal{N}(\underline{\mu}_1, \underline{\sigma}_1^2)$
 $p(x|z=2) \sim \mathcal{N}(\underline{\mu}_2, \underline{\sigma}_2^2)$

Gaussian Mixture Model

$$\rightarrow \underline{z} \sim \text{Cat}(\underline{\pi}) \quad \left[\begin{array}{c} \pi_1 \\ \vdots \\ \pi_k \end{array} \right] \quad \left\{ \begin{array}{l} \pi_c \geq 0 \\ \sum_{c=1}^k \pi_c = 1 \end{array} \right\} \quad \pi_c = P(z=c)$$

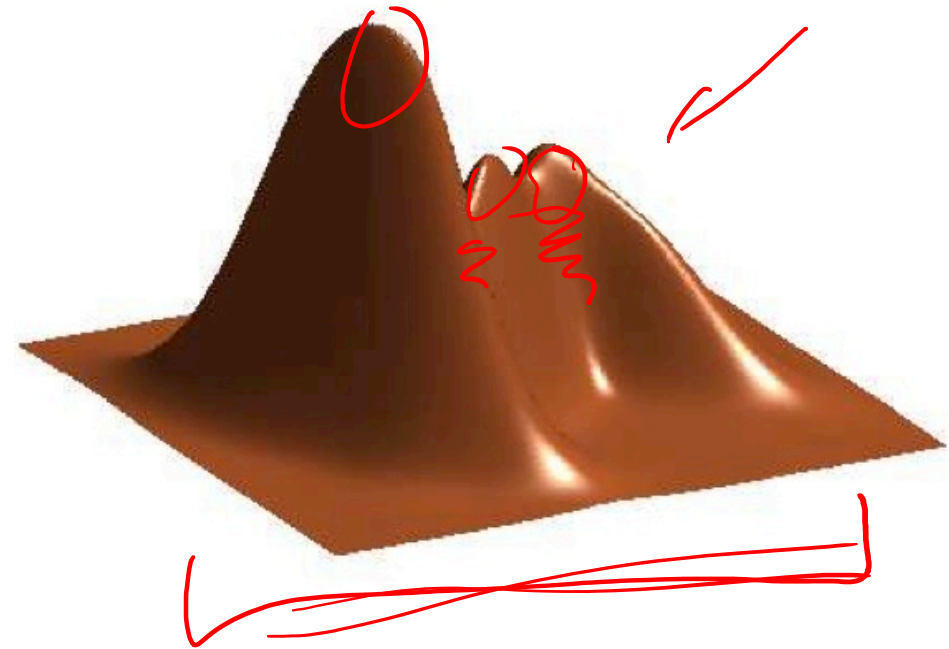
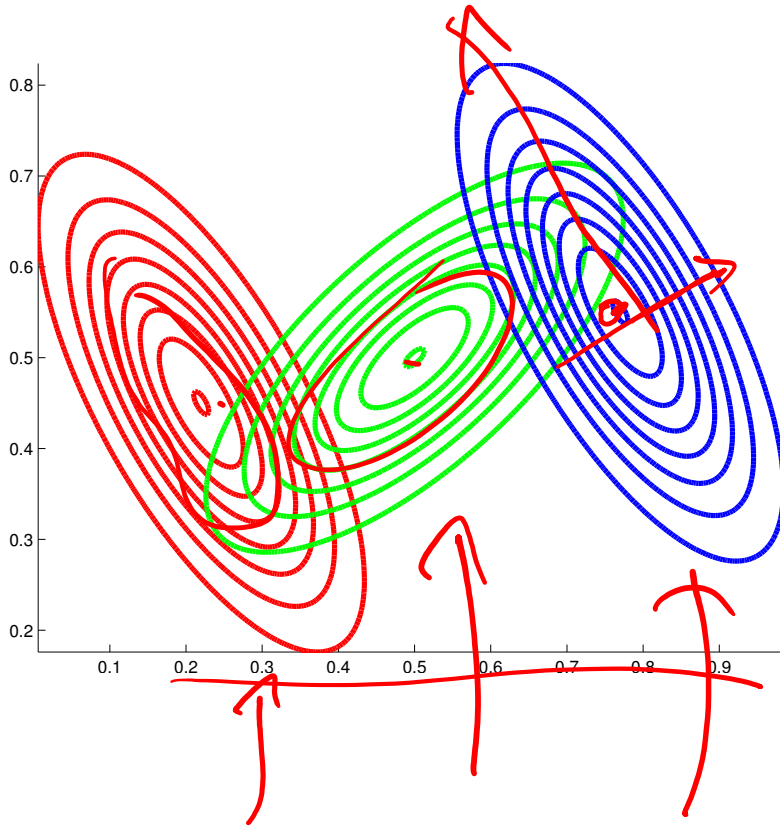
$$\underline{x} | z=c \sim N(\mu_c, \sigma_c^2) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(x-\mu_c)^2}{2\sigma_c^2}}$$

$$\left[\begin{array}{l} p(z) \\ p(x|z) \end{array} \right] = p(x, z)$$
$$\pi_c \cdot N(\mu_c, \sigma_c^2)$$

GMM

$$\vec{\mu} \in \mathbb{R}^d$$
$$\Sigma \in \mathbb{R}^{d \times d}$$

$$N(\vec{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1}(\vec{x}-\vec{\mu})}$$



$p(z)$ Gaussian Mixture Model

$$p(x, z)$$

$$p(x) = \sum_c p(z) p(x|z) \leftarrow \text{Marginalization}$$

$$p(z|x) = \frac{p(z, x)}{p(x)} = \frac{p(x|z) p(z)}{\sum_z p(x|z) p(z)} \leftarrow \begin{array}{l} \text{'Inference'} \\ \text{Posterior} \end{array}$$

$$\sim p(x|z)$$

~~$p(x)$~~

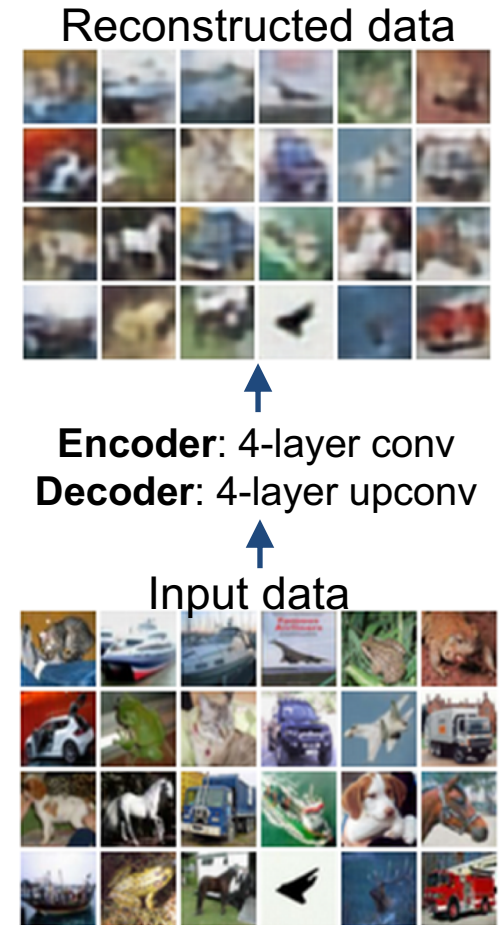
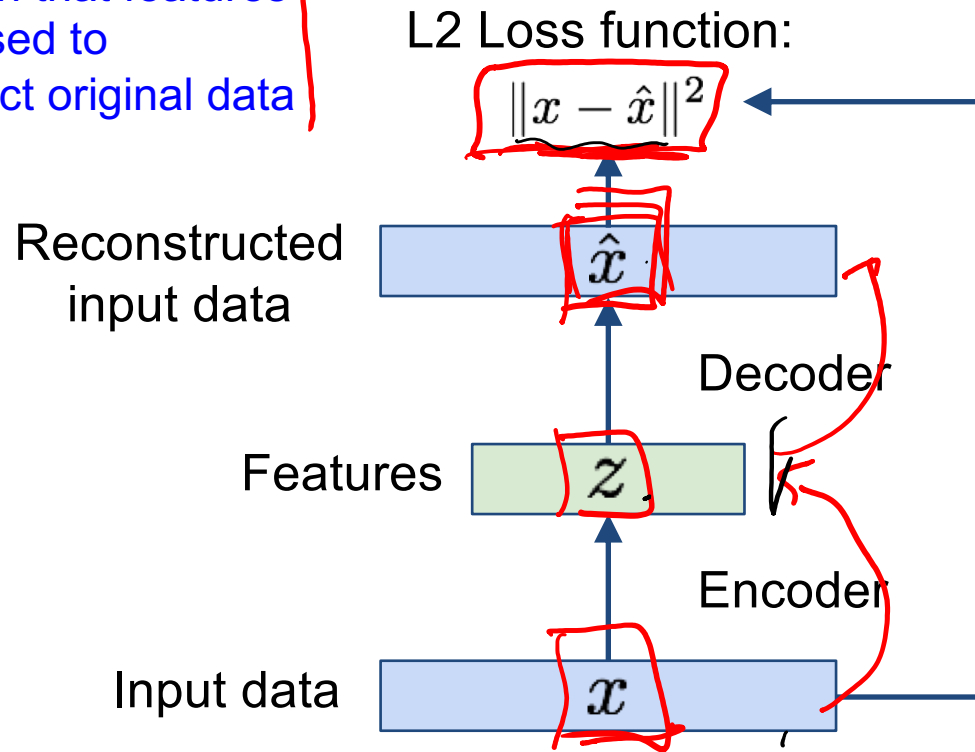
Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

Autoencoders

Train such that features can be used to reconstruct original data

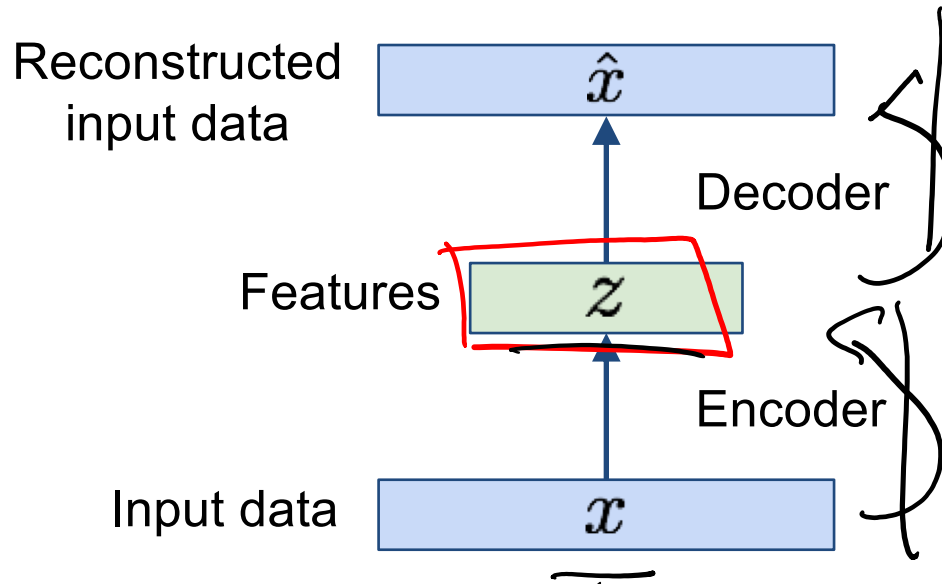


Autoencoders

$$z = f_{\phi}(x)$$
$$\hat{x} = g_{\theta}(z)$$

$$p(z)$$
$$p(\hat{x} | z)$$

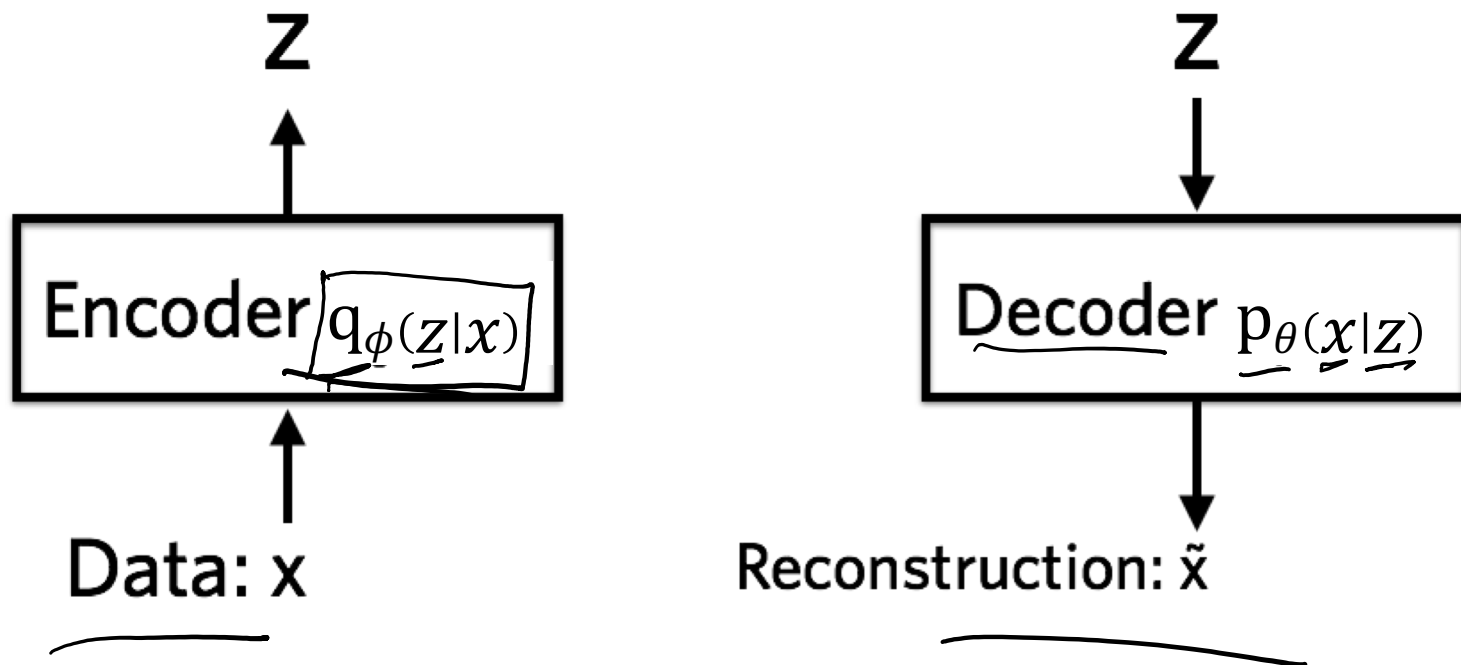
Autoencoders can reconstruct data, and can learn features to initialize a supervised model



Features capture factors of variation in training data. Can we generate new images from an autoencoder?

Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!



Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

Key problem

$$\bullet \boxed{P(z|x)} = \frac{P(z,x)}{\boxed{P(z)}} = \frac{P(x|z)P(z)}{\boxed{\sum_z P(x|z)P(z)}} \quad |$$

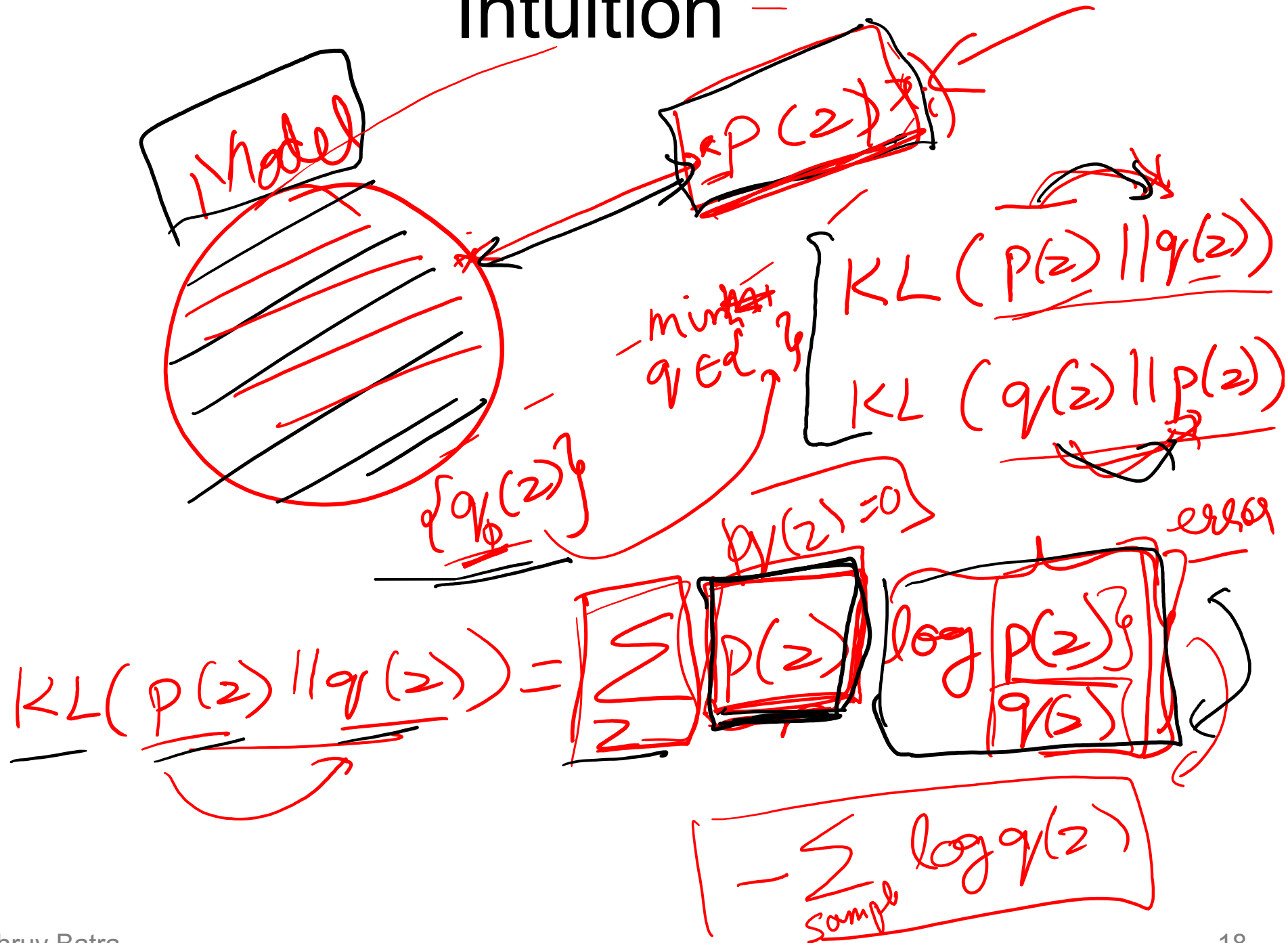
\downarrow

$$\boxed{q(z)}$$

What is Variational Inference?

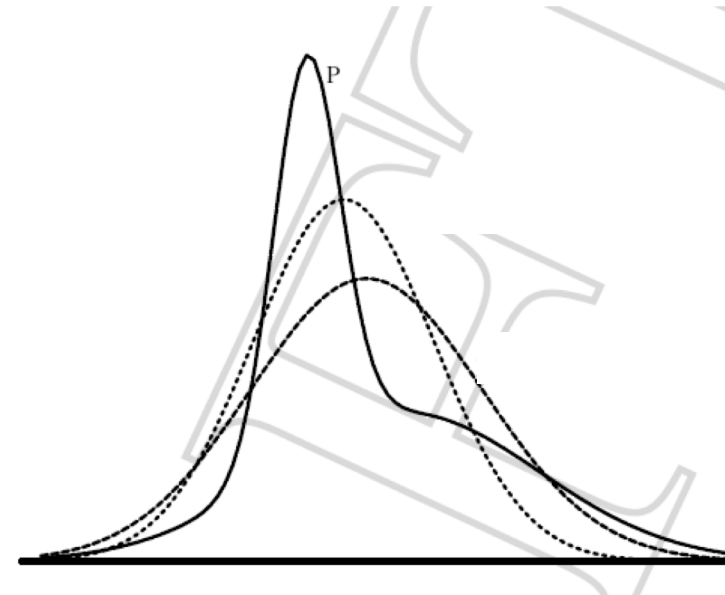
- Key idea
 - Reality is complex
 - Can we approximate it with something “simple”?
 - Just need to make sure the simple thing is “close” to the complex thing.

Intuition



Find simple approximate distribution

- Suppose p is intractable posterior
- Want to find simple q that approximates p
- KL divergence not symmetric
- $D(p||q)$
 - true distribution p defines support of diff.
 - the “correct” direction
 - will be intractable to compute
- $D(q||p)$
 - approximate distribution defines support
 - tends to give overconfident results
 - will be tractable

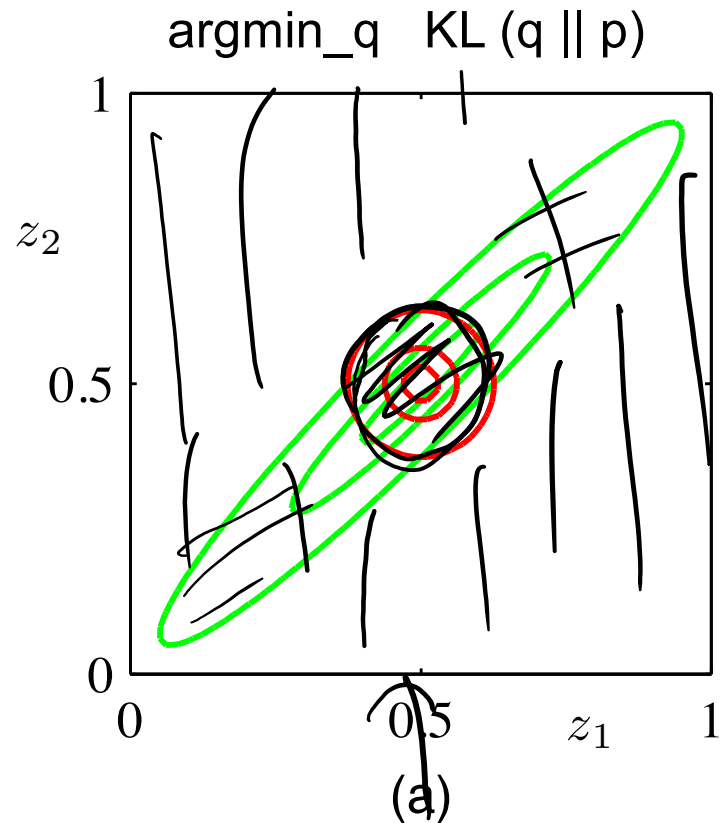
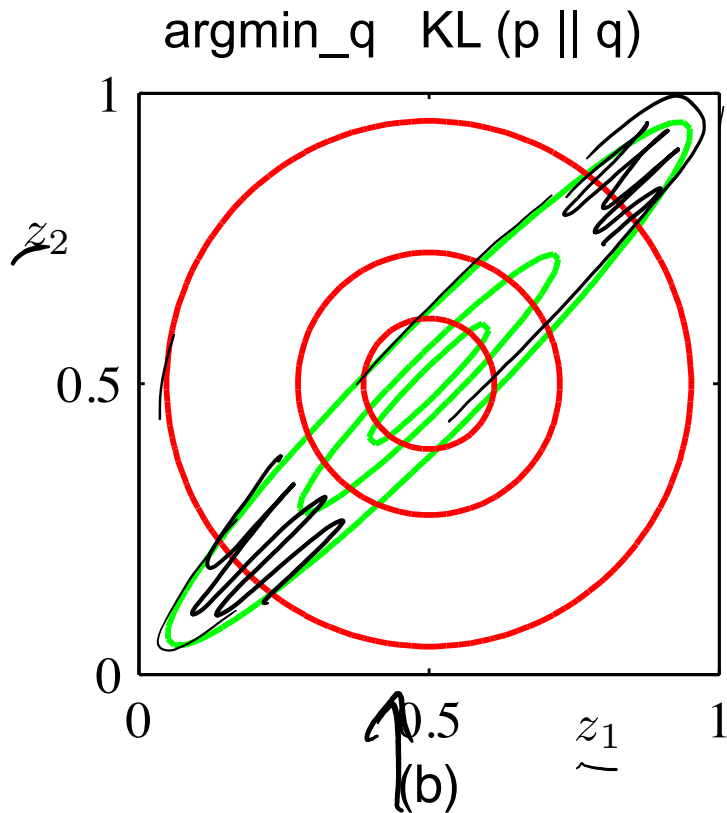


Example 1

- p = 2D Gaussian with arbitrary co-variance
- q = 2D Gaussian with diagonal co-variance

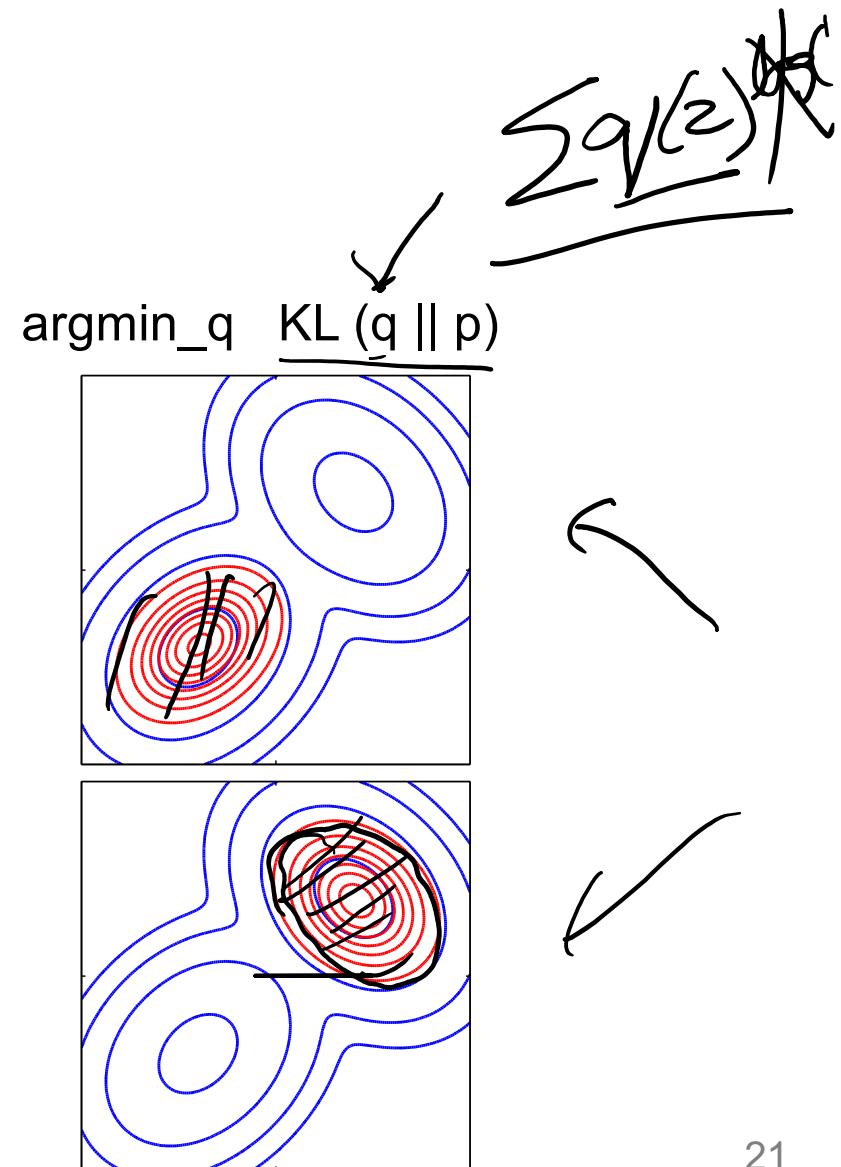
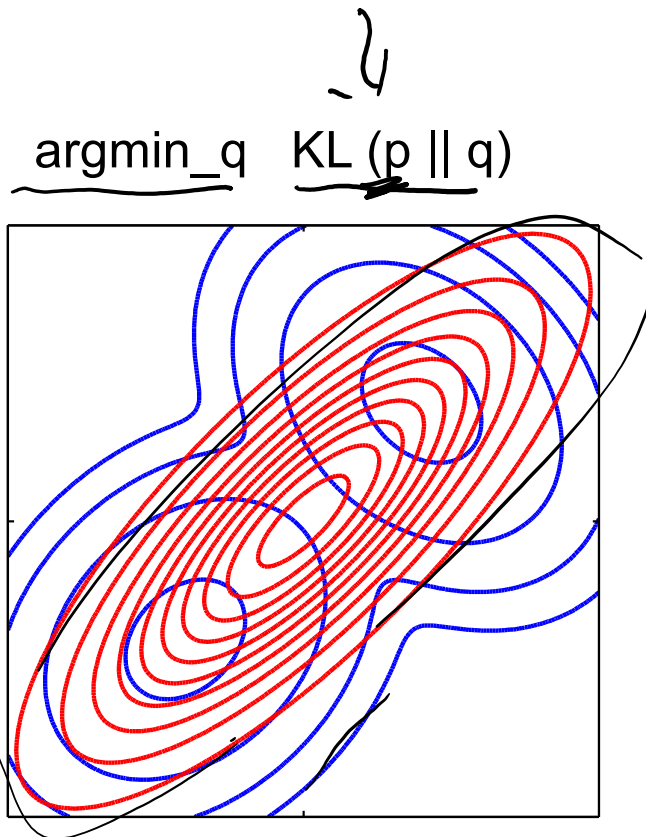
$$\Sigma_p = \begin{bmatrix} & \\ & \end{bmatrix}_{2 \times 2}$$

$$\Sigma_q = \begin{bmatrix} & 0 \\ 0 & \end{bmatrix}$$



Example 2

- p = Mixture of Two Gaussians
- q = Single Gaussian



Plan for Today

- VAEs
 - Variational Inference → Evidence Based Lower Bound
 - Reparameterization trick
 - Putting it all together
- Generative Adversarial Networks

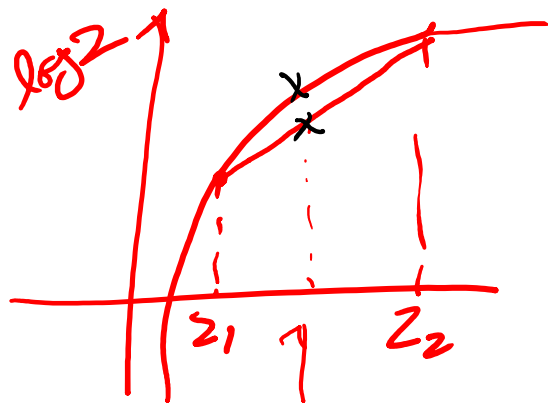
The general learning problem with missing data

- Marginal likelihood – \mathbf{x} is observed, \mathbf{z} is missing:

$$\begin{aligned} \boxed{ll(\theta : \mathcal{D})} &= \log \prod_{i=1}^N P(\mathbf{x}_i | \theta) \\ &= \sum_{i=1}^N \log P(\mathbf{x}_i | \theta) \quad \leftarrow P(\vec{x}_i, \mathbf{z}) \\ &= \cancel{\sum_{i=1}^N} \log \sum_{\mathbf{z}} P(\mathbf{x}_i, \mathbf{z} | \theta) \\ &= \log \left[\sum_{\mathbf{z}} P(\vec{x}_i | \theta) P(\mathbf{z} | \vec{x}_i, \theta) \right] \\ &= \log \left[\underset{P(\mathbf{z} | \vec{x}_i, \theta)}{\mathbb{E}} \left[P(\vec{x}_i | \theta) \right] \right] \end{aligned}$$

Applying Jensen's inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) g(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log g(\mathbf{z})$



$$\lambda z_1 + \lambda_2 z_2$$

$$\lambda_1, \lambda_2 \geq 0$$

$$\lambda_1 + \lambda_2 = 1$$

$$f(\lambda_1 z_1 + \lambda_2 z_2) \geq \lambda_1 f(z_1) + \lambda_2 f(z_2)$$

$$f\left(\sum_{i=1}^K \lambda_i z_i\right) \geq \sum_{i=1}^K \lambda_i f(z_i)$$

$$f(E[z]) \geq E[f(z)]$$

$$f(E[g(z)]) \geq E[f(g(z))]$$

$z \rightarrow g(z)$

Applying Jensen's inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) g(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log g(\mathbf{z})$

$$\rightarrow \log \sum_{\mathbf{z}} P(\vec{x}_i, \mathbf{z} | \theta) \cdot \frac{Q_i(\mathbf{z})}{Q_i(\mathbf{z})}$$

$$\geq \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\vec{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})} \quad \left. \vphantom{\sum_{\mathbf{z}}} \right] \frac{F(\theta, Q_i)}{\text{"Free Energy" Variational}}$$

VLB / ELBO

Applying Jensen's inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) g(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log g(\mathbf{z})$

$$\underline{l(\theta : \mathcal{D})} = \sum_{i=1}^N \log \sum_{\mathbf{z}} Q_i(\mathbf{z}) \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

Handwritten annotations:

- A red box around the summation $\sum_{i=1}^N$ is crossed out with a red 'X'.
- A red arrow points from the fraction $\frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$ to the expression $\underline{P(\mathbf{z} | \theta) P(\mathbf{x}_i | \mathbf{z}, \theta)}$.
- A red arrow points from the denominator $Q_i(\mathbf{z})$ to the expression $\underline{P(\mathbf{x}_i | \theta) P(\mathbf{z} | \mathbf{x}_i, \theta)}$.

Evidence Lower Bound

- Define potential function $F(\theta, Q)$:

$$\underline{\mathcal{L}(\theta : \mathcal{D})} \geq \underline{F(\theta, Q_i)} = \sum_{i=1}^N \sum_{\mathbf{z}} \underline{Q_i(\mathbf{z})} \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

ELBO: Factorization #1

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q_i) = \sum_{i=1}^N \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

$\rightarrow P(\bar{x}_i | \theta) P(z | \bar{x}_i, \theta)$

$$\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \left(\frac{P(\bar{x}_i | \theta) P(z | \bar{x}_i, \theta)}{Q_i(\mathbf{z})} \right)$$

$$= \underbrace{\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log P(\bar{x}_i | \theta)}_1 + \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(z | \bar{x}_i, \theta)}{Q_i(\mathbf{z})} \right]$$

$$\underline{F(\theta, Q_i)} = \underline{\log P(\bar{x}_i | \theta)} - \left[\text{KL} \left(\frac{Q_i(\mathbf{z})}{P(z | \bar{x}_i, \theta)} \right) \right]$$

$$\max_{\theta, Q_i} F = \text{ll} - \text{KL}$$

ELBO: Factorization #1

\max_{θ, Q_i}
 $F(\theta, Q_i)$

$$l(\theta : \mathcal{D}) \geq F(\theta, Q_i) = \sum_{i=1}^N \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

$\rightarrow P(\mathbf{z} | \theta) P(\mathbf{x}_i | \mathbf{z}, \theta)$

$$= \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{z} | \theta) P(\mathbf{x}_i | \mathbf{z}, \theta)}{Q_i(\mathbf{z})}$$

$$= \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log P(\mathbf{x}_i | \mathbf{z}, \theta) \right] + \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{z} | \theta)}{Q_i(\mathbf{z})} \right]$$

decoder

$$= \left[\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log P(\mathbf{x}_i | \mathbf{z}, \theta) \right] - \text{KL} [Q_i(\mathbf{z}) || P(\mathbf{z} | \theta)]$$

encode \uparrow
Explain the data

Regularizer

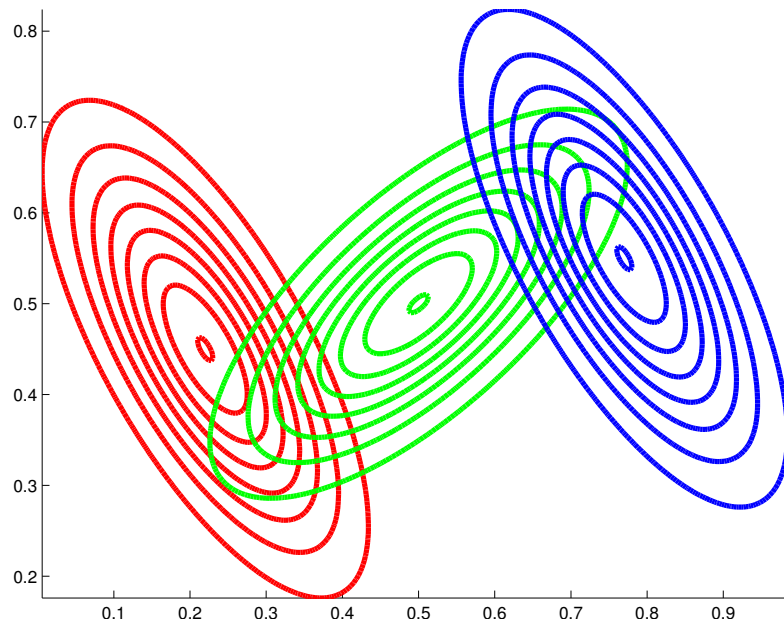
Evidence Lower Bound

- Define potential function $F(\theta, Q)$:

$$l(\theta : \mathcal{D}) \geq \underline{F(\theta, Q_i)} = \sum_{i=1}^N \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

- **EM** corresponds to coordinate ascent on F
 - Thus, maximizes lower bound on marginal log likelihood

GMM



EM for Learning GMMs

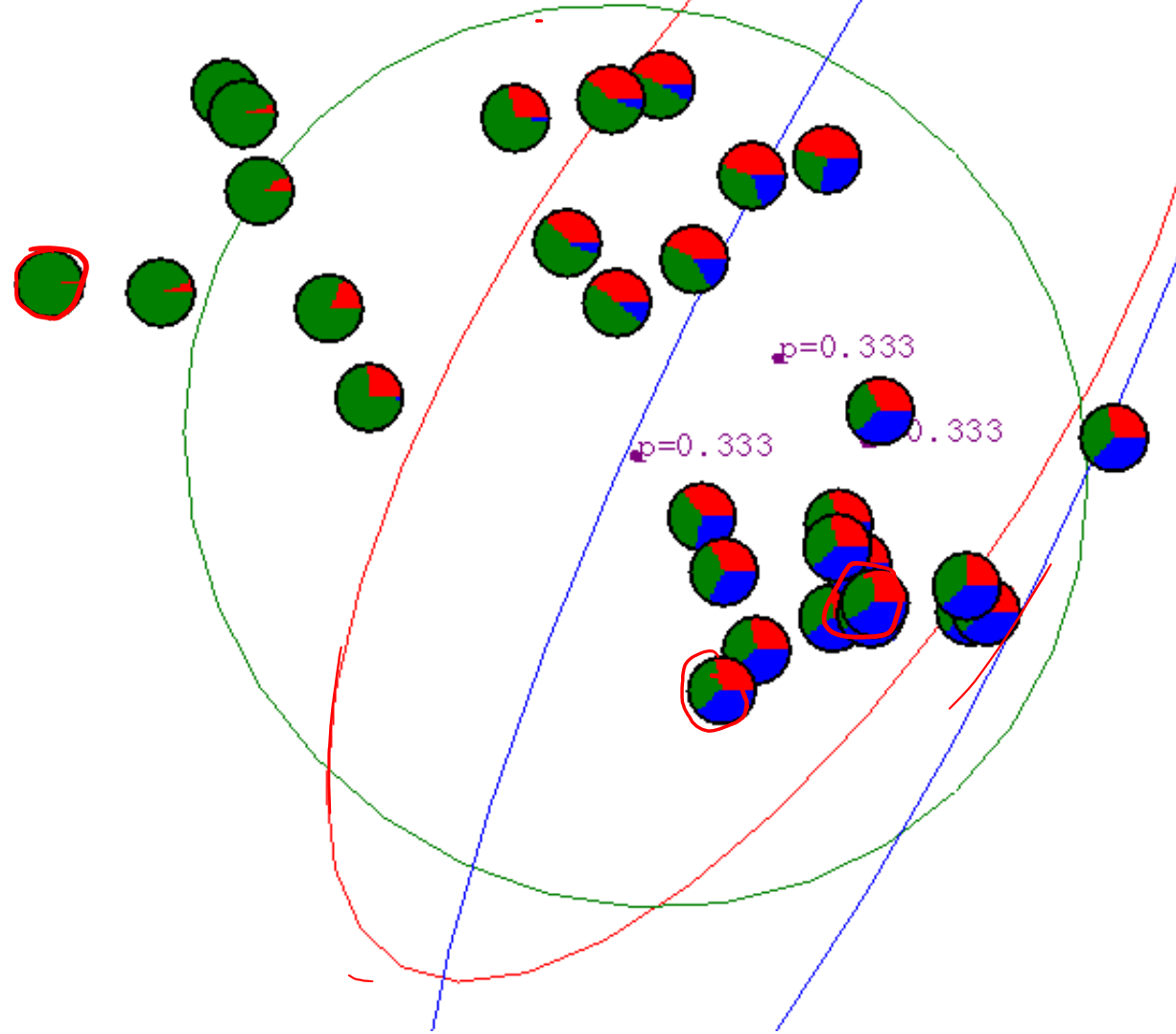
- Simple Update Rules
 - E-Step: estimate $Q_i(z) = \Pr(z = j \mid x_i)$
 - M-Step: maximize expected likelihood under $Q_i(z)$

Gaussian Mixture Example: Start

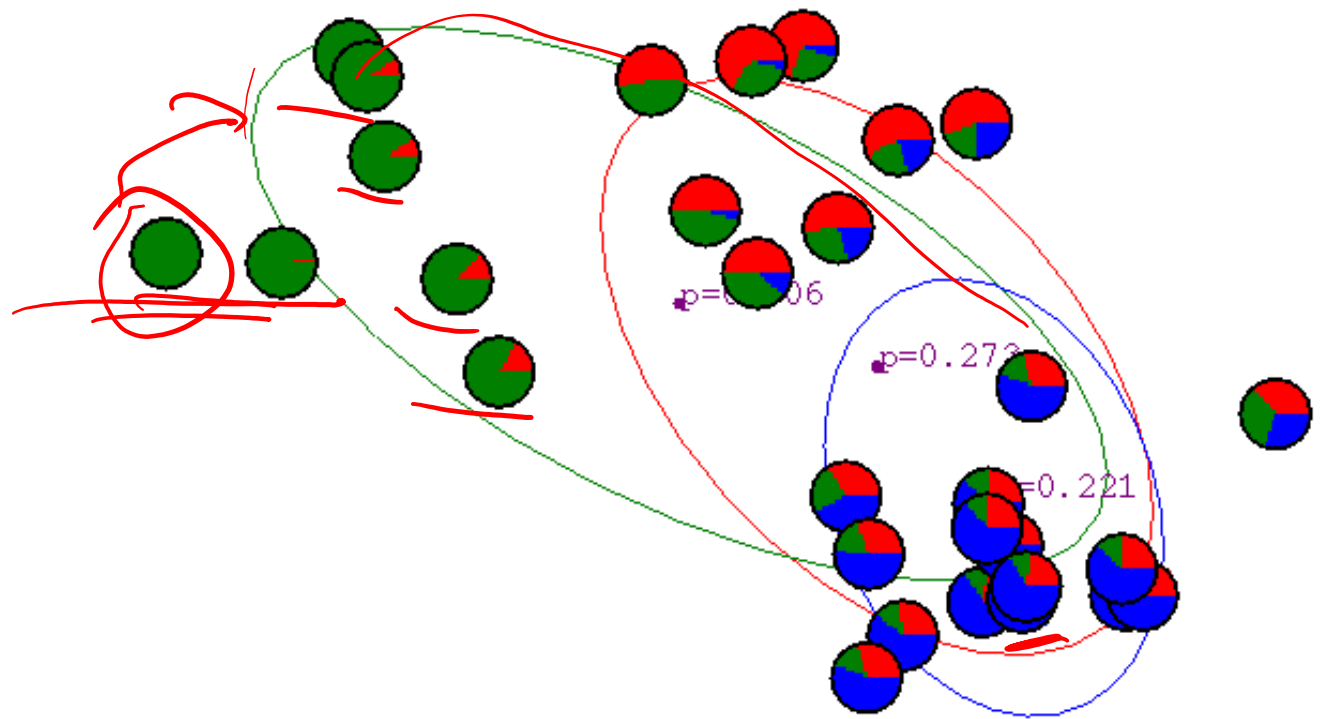
$$Q_i(z) = P(z | x_i)$$

$$\begin{pmatrix} 0.3 \\ 0.3 \\ 0.3 \end{pmatrix}$$

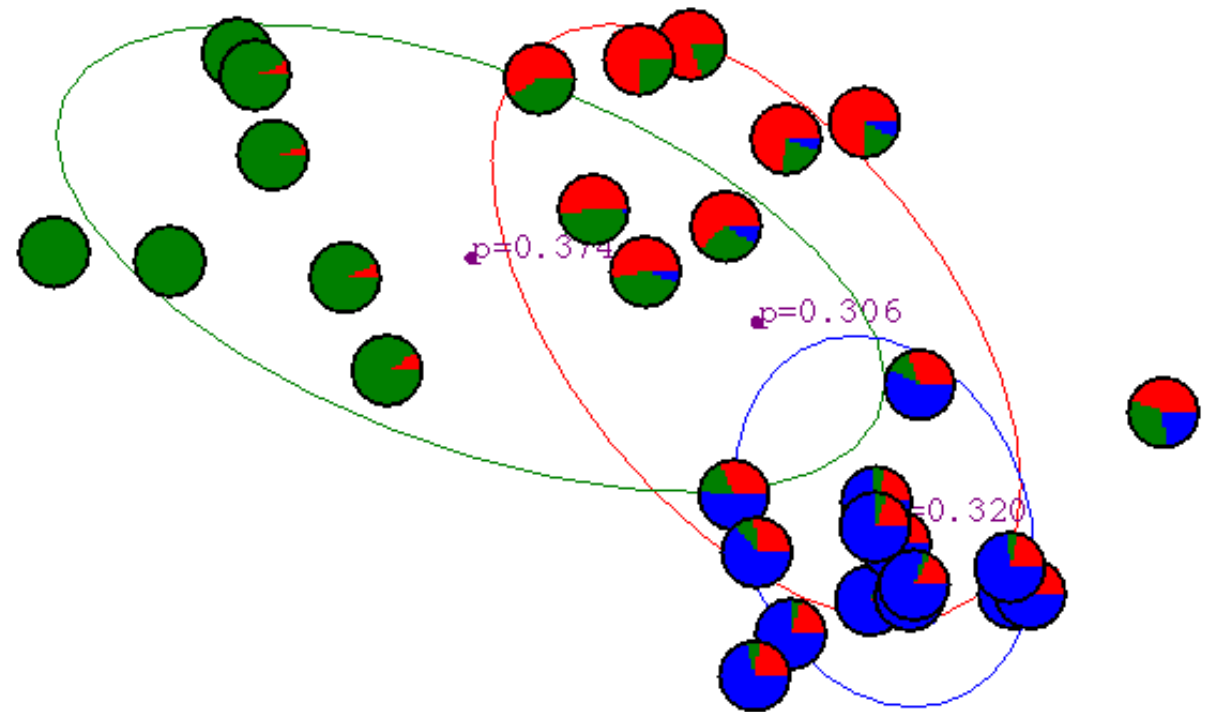
$$\Theta = \{\mu_c, \Sigma_c, \pi_c\}$$



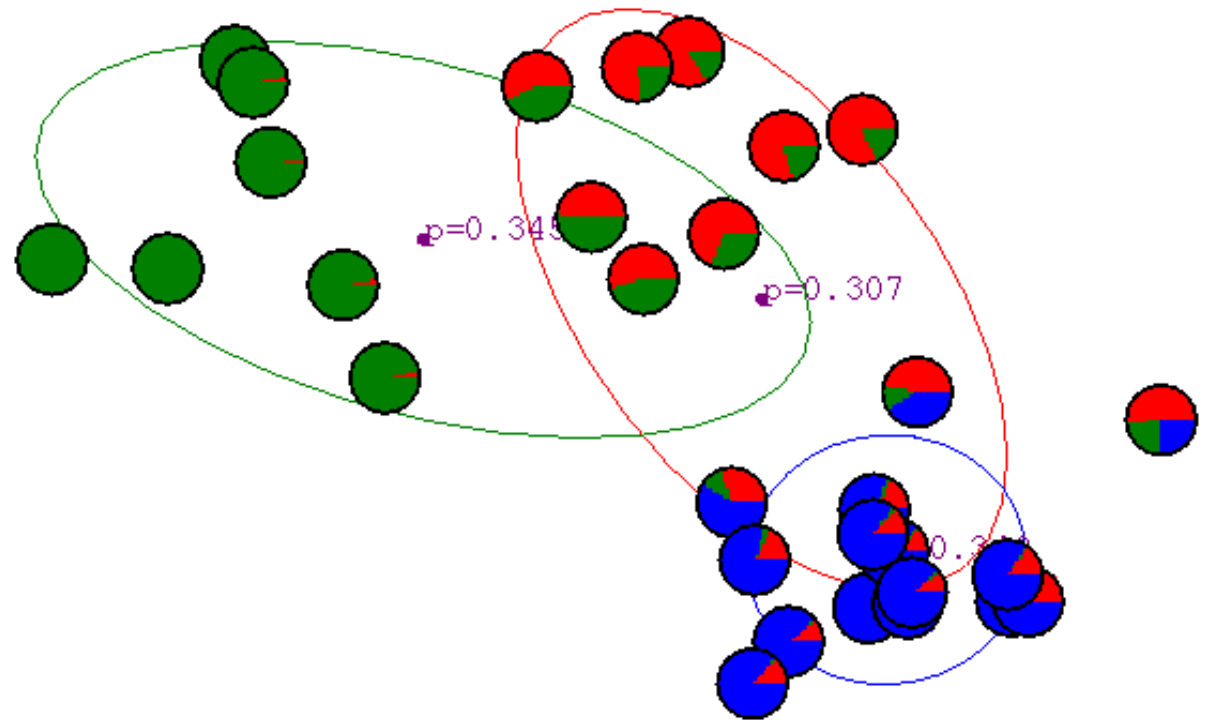
After 1st iteration



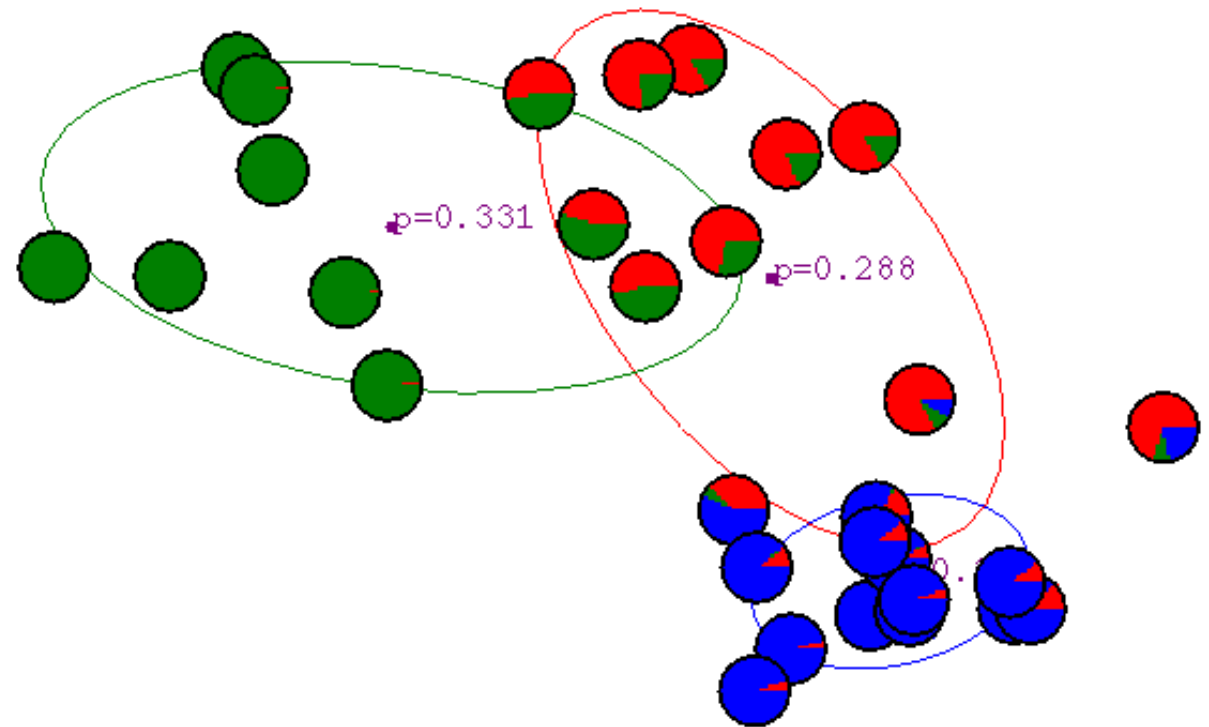
After 2nd iteration



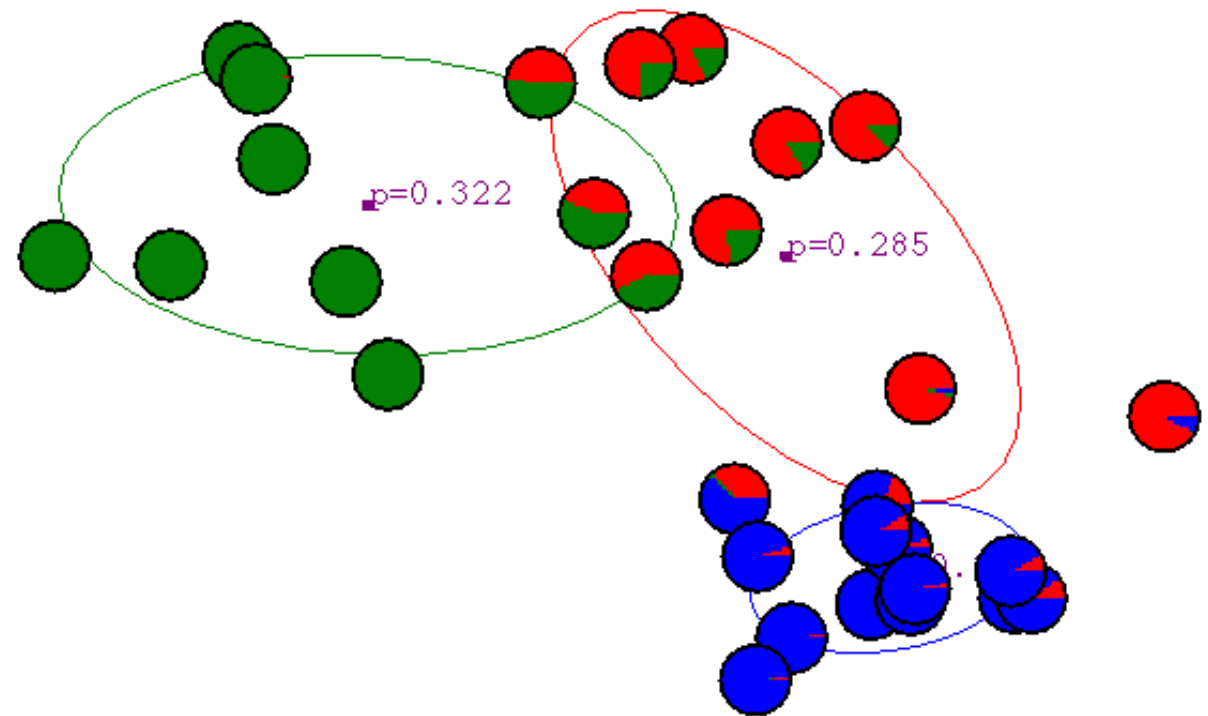
After 3rd iteration



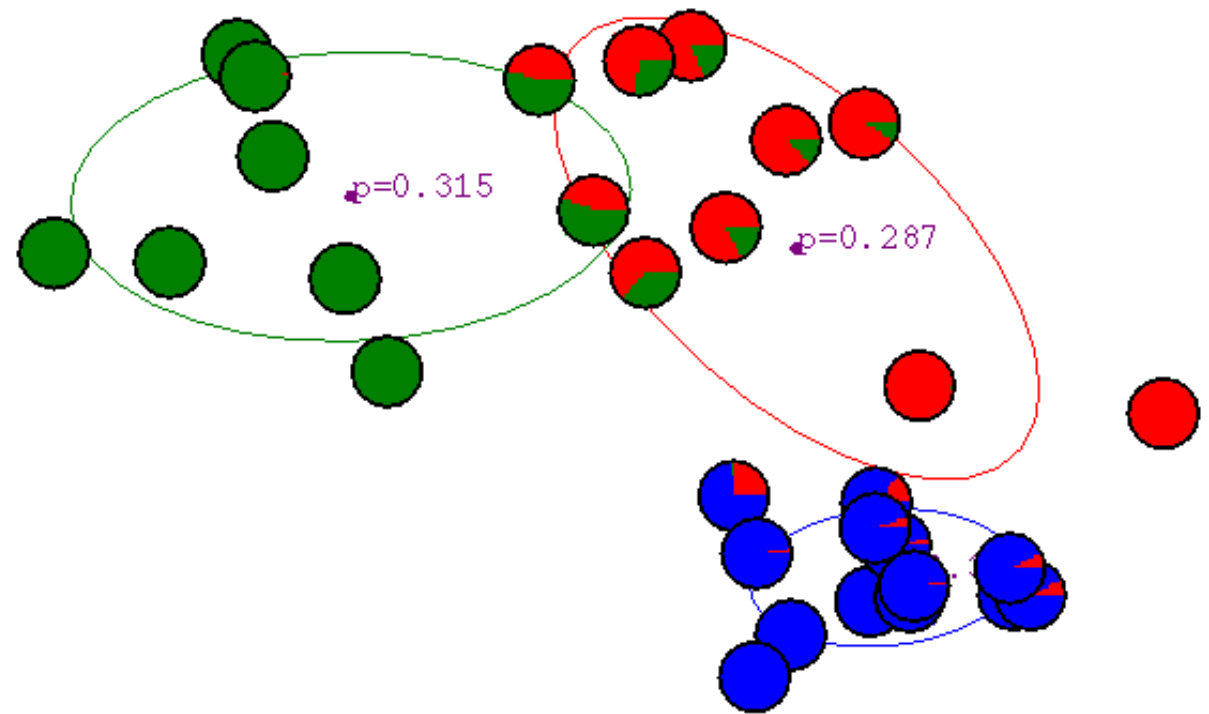
After 4th iteration



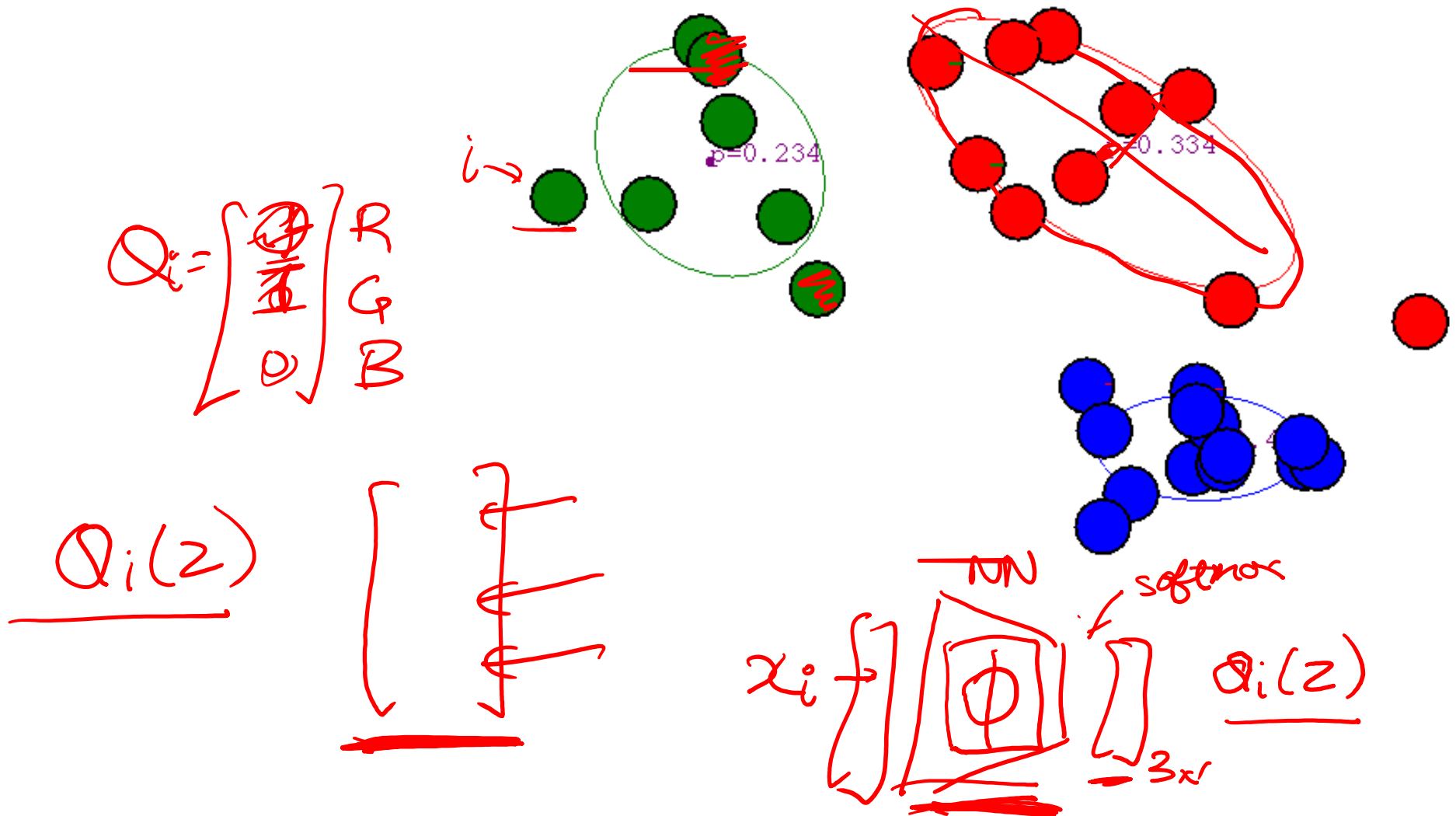
After 5th iteration



After 6th iteration

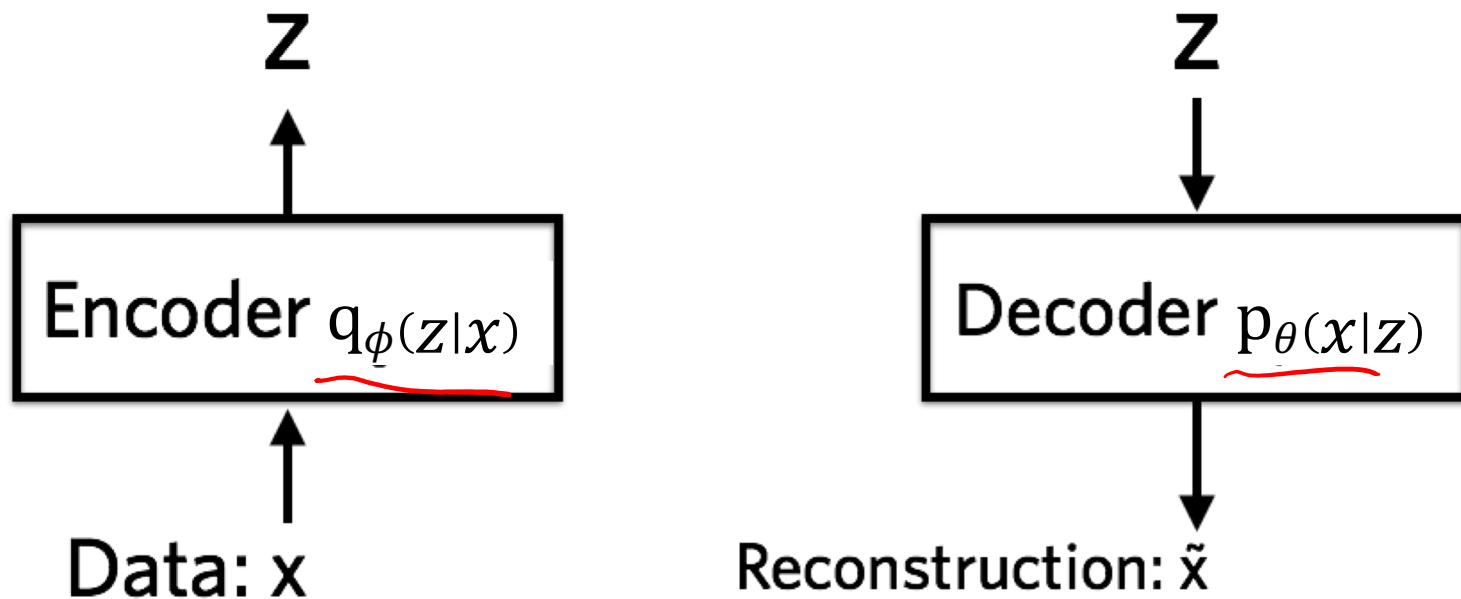


After 20th iteration



Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!



Evidence Lower Bound

- Define potential function $F(\theta, Q)$:

$$l(\theta : \mathcal{D}) \geq F(\theta, Q_i) = \sum_{i=1}^N \sum_{\mathbf{z}} \underbrace{Q_i(\mathbf{z})}_{\text{red underline}} \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

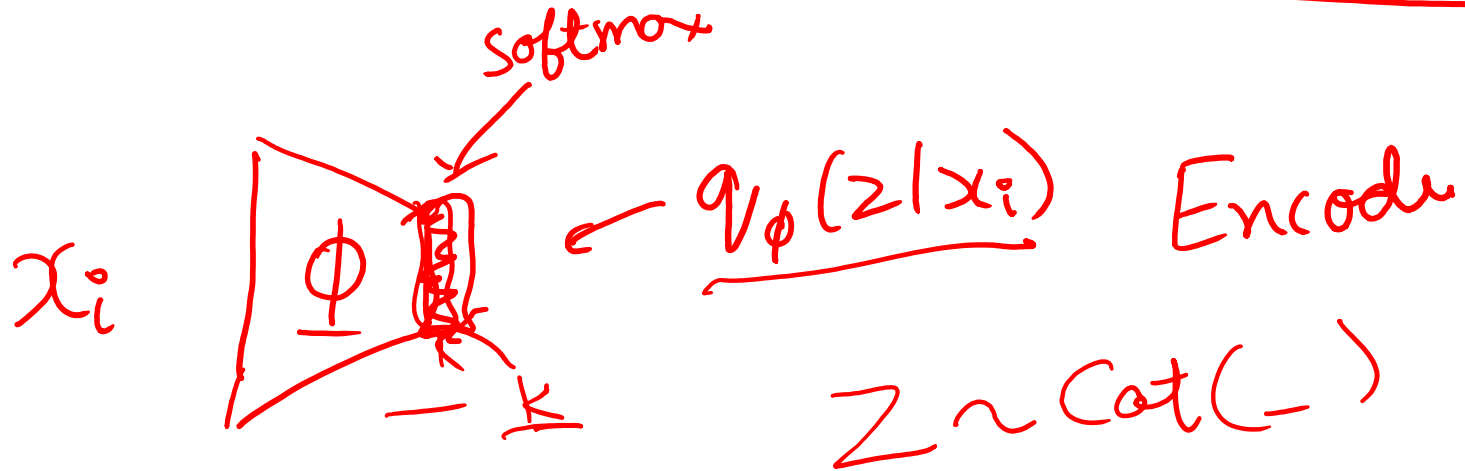
2. Variational Approximation

- Variational Lower Bound / ELBO

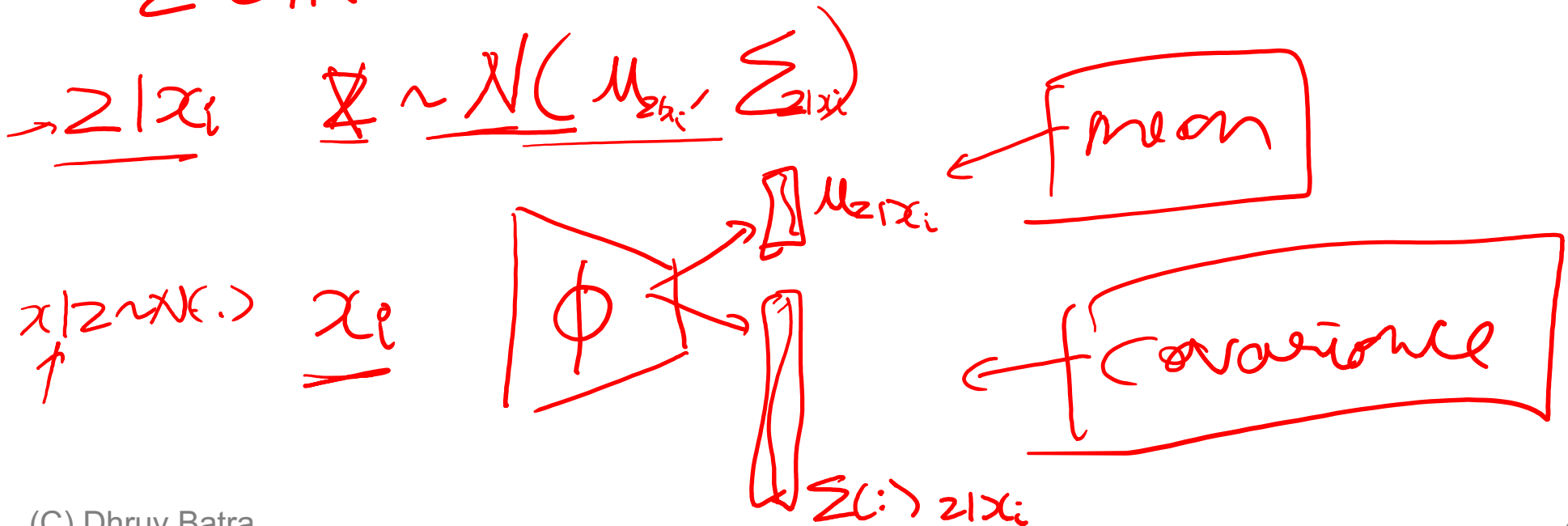
3. Amortized Inference Neural Networks ←

4. “Reparameterization” Trick

Amortized Inference Neural Networks



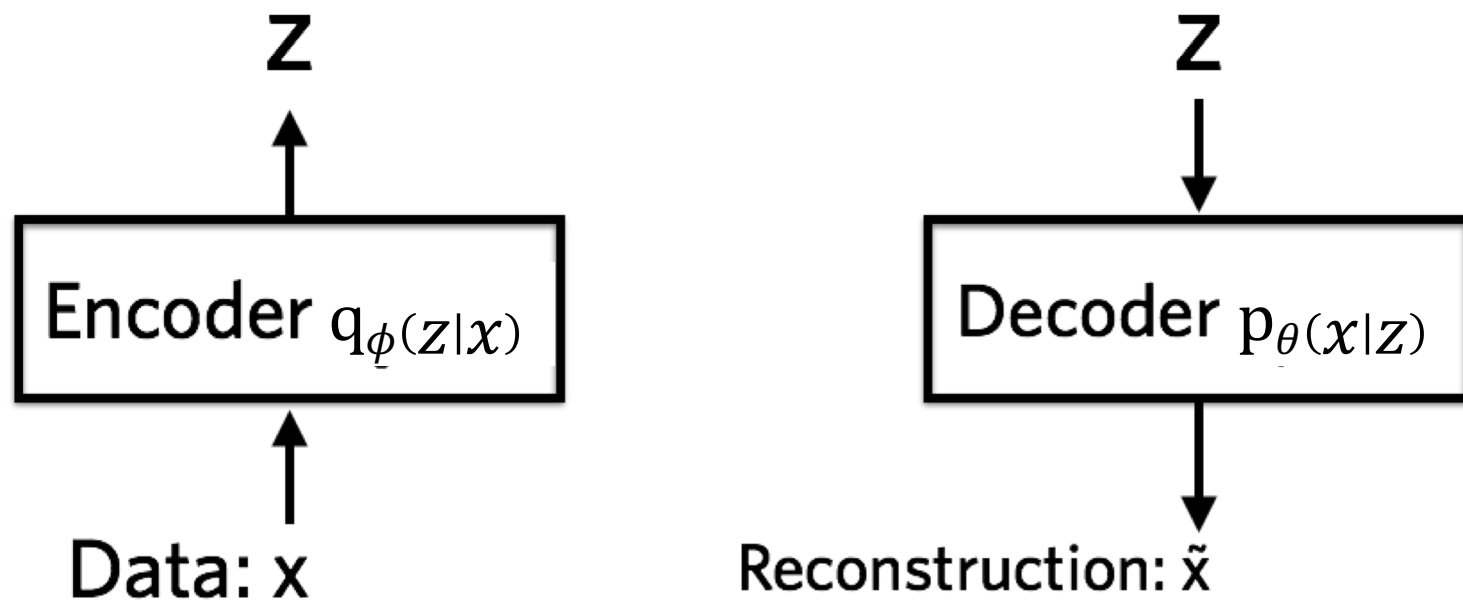
$$z \in \mathbb{R}^k$$



Amortized Inference Neural Networks

Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!

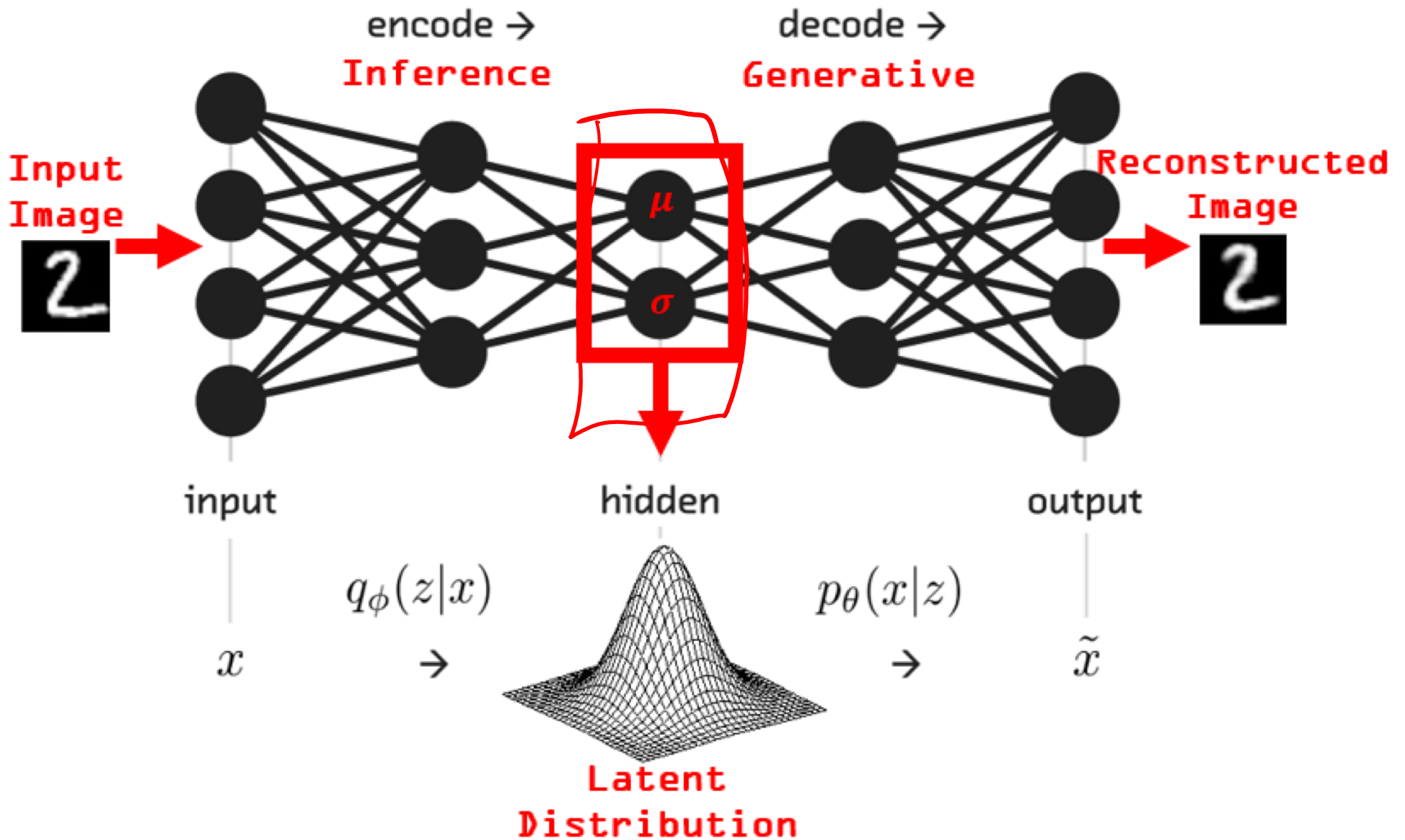


Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick

VAEs



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(\underline{x}^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Let's look at computing the bound (forward pass) for a given minibatch of input data

Input Data

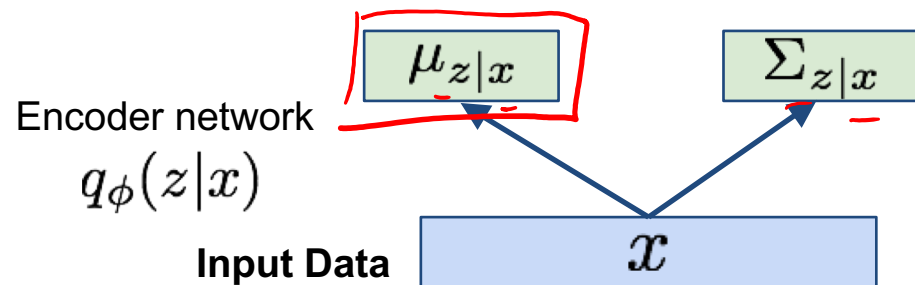


\mathcal{X}

Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

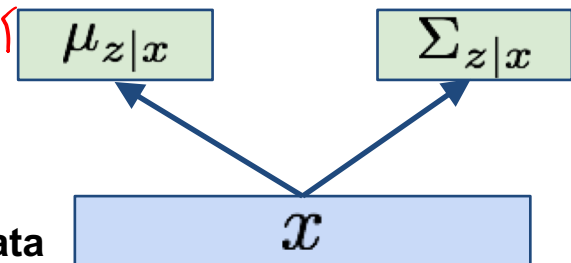
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\text{Make approximate posterior distribution close to prior}}$$

Make approximate posterior distribution close to prior

Encoder network

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

Input Data

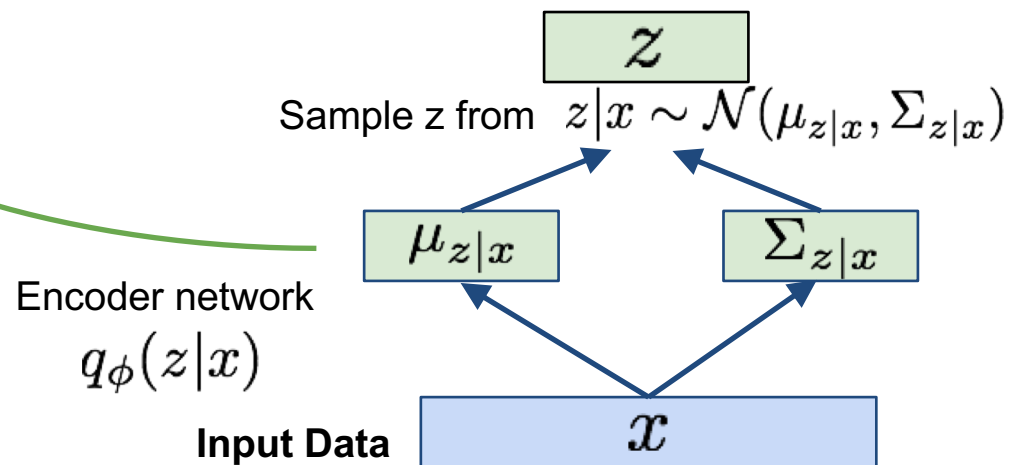


Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

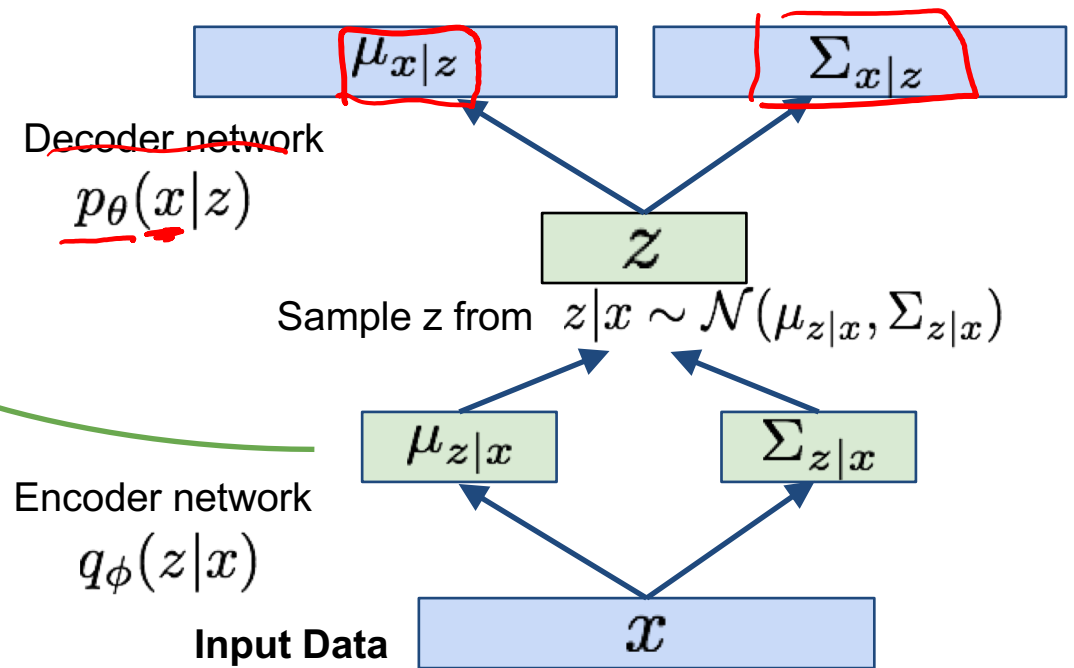


Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior



Variational Auto Encoders

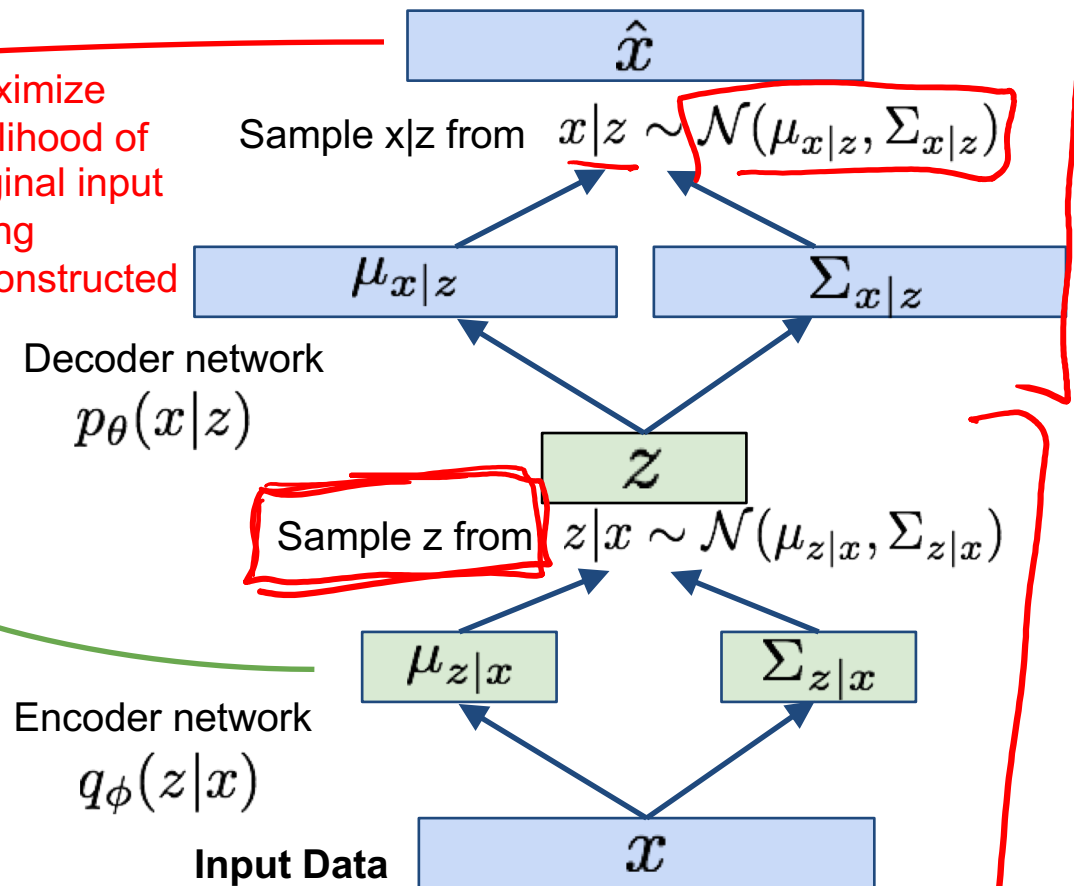
Putting it all together: maximizing the likelihood lower bound

$$\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

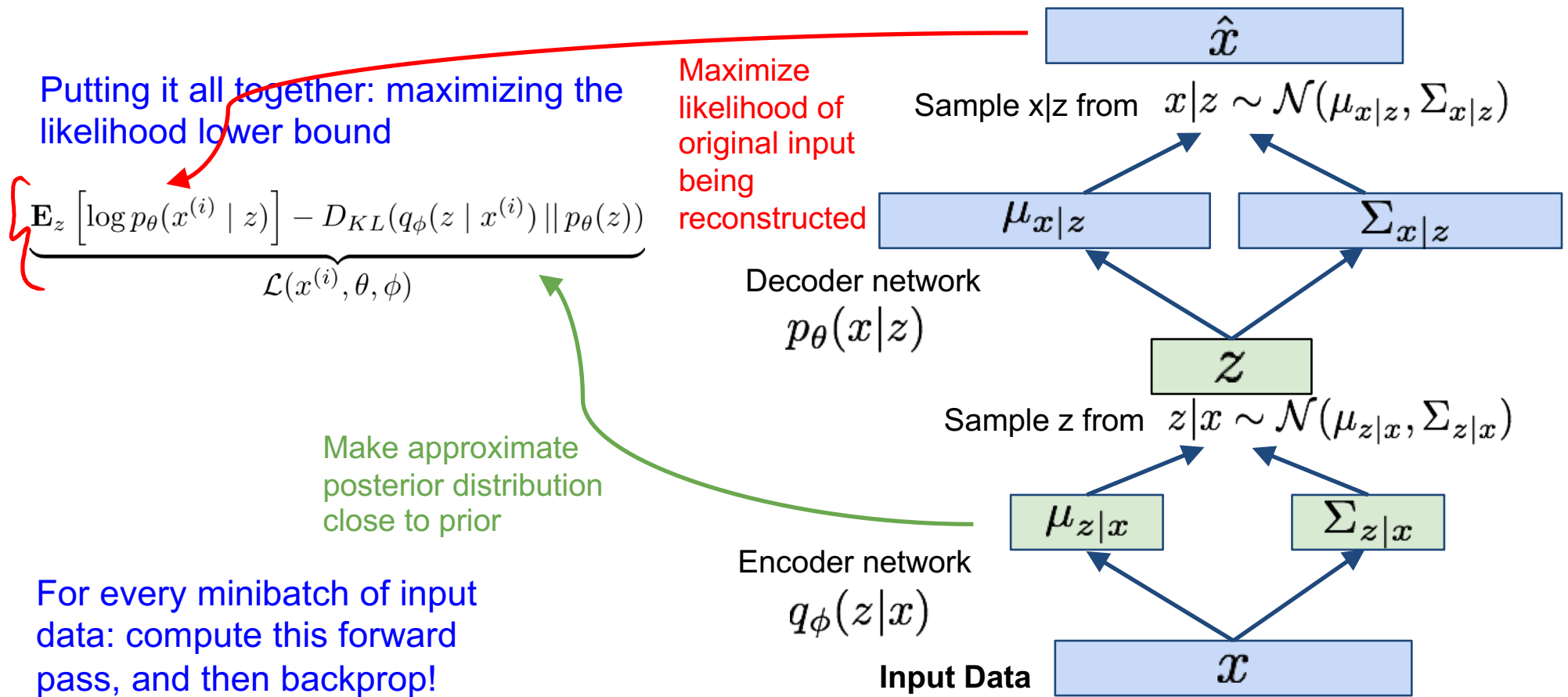
$\mathcal{L}(x^{(i)}, \theta, \phi)$

Maximize likelihood of original input being reconstructed

Make approximate posterior distribution close to prior

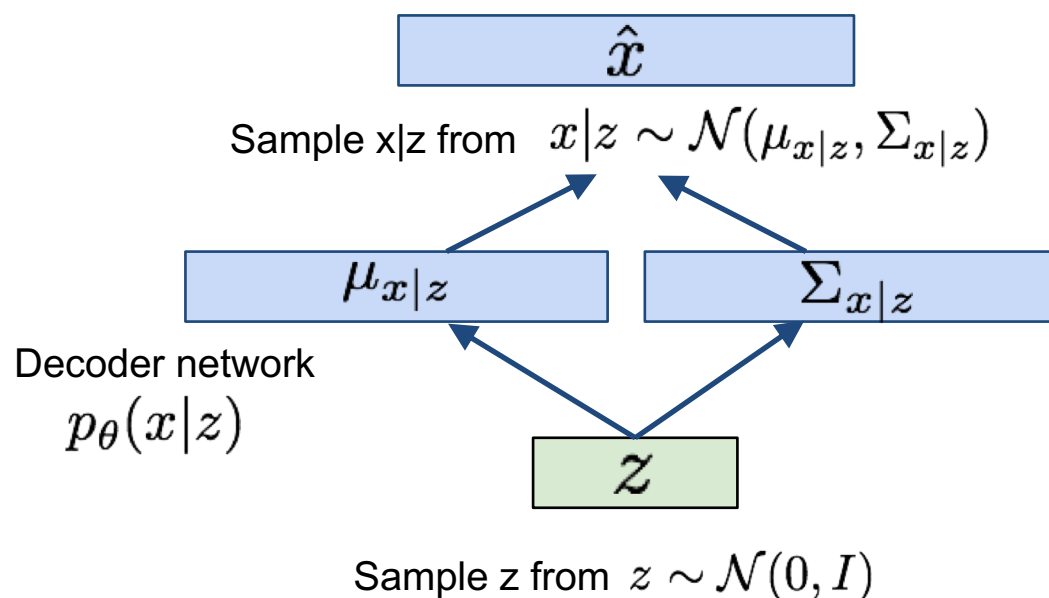


Variational Auto Encoders



Variational Auto Encoders: Generating Data

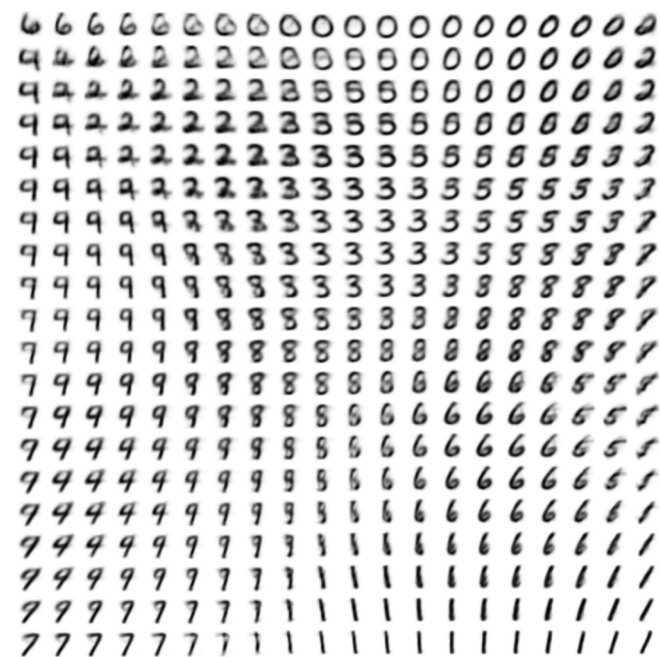
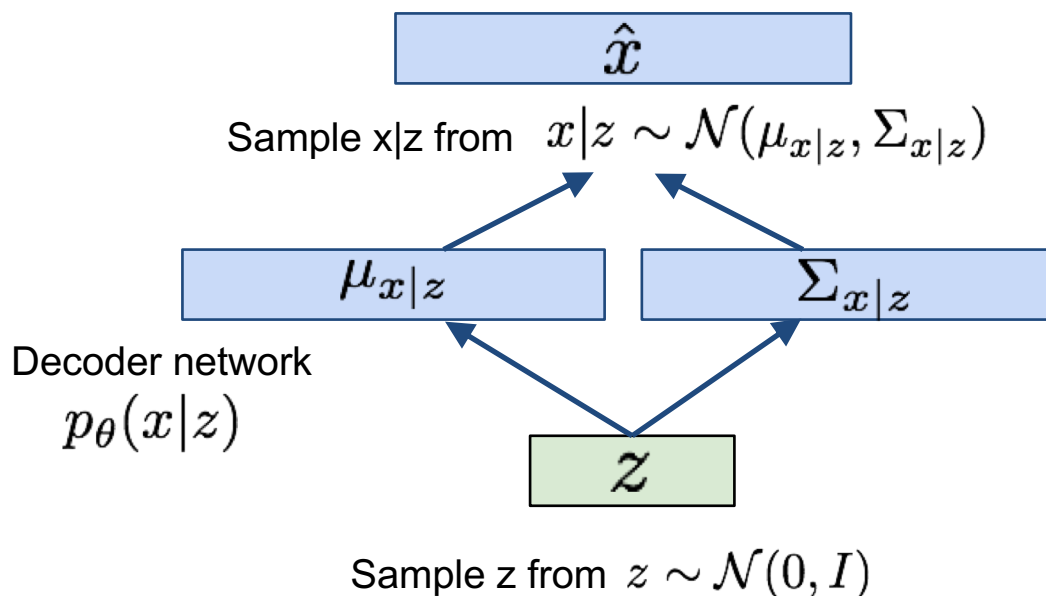
Use decoder network. Now sample z from prior!



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Auto Encoders: Generating Data

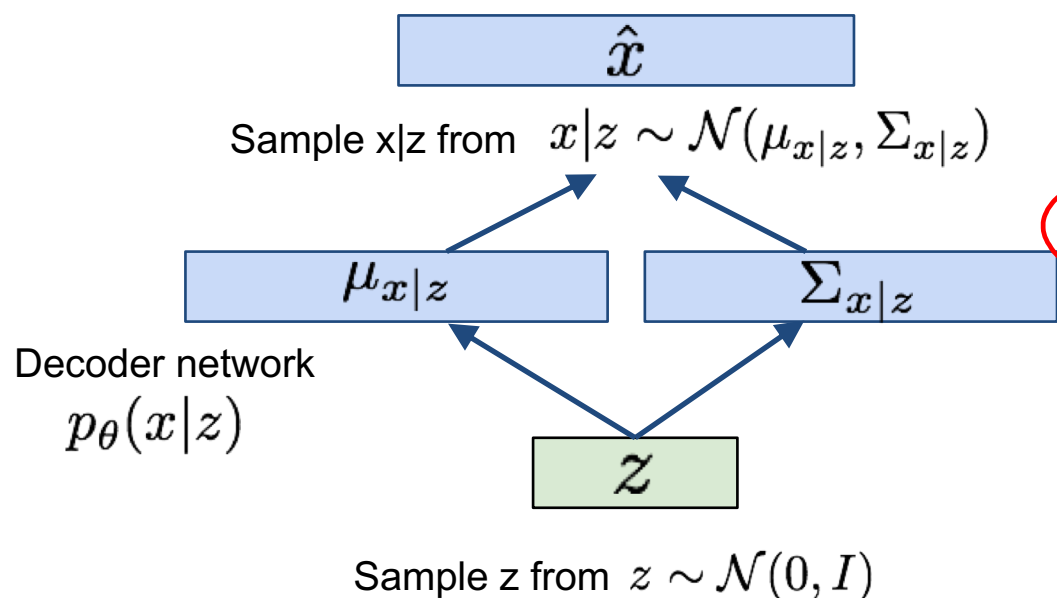
Use decoder network. Now sample z from prior!



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

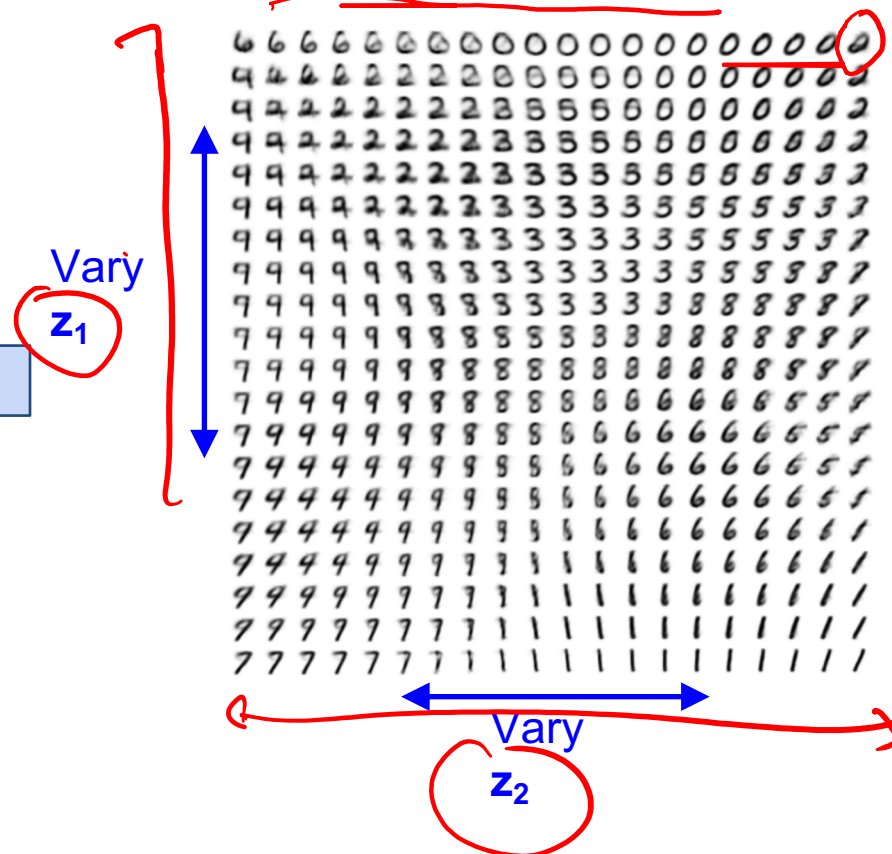
Variational Auto Encoders: Generating Data

Use decoder network. Now sample z from prior!



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

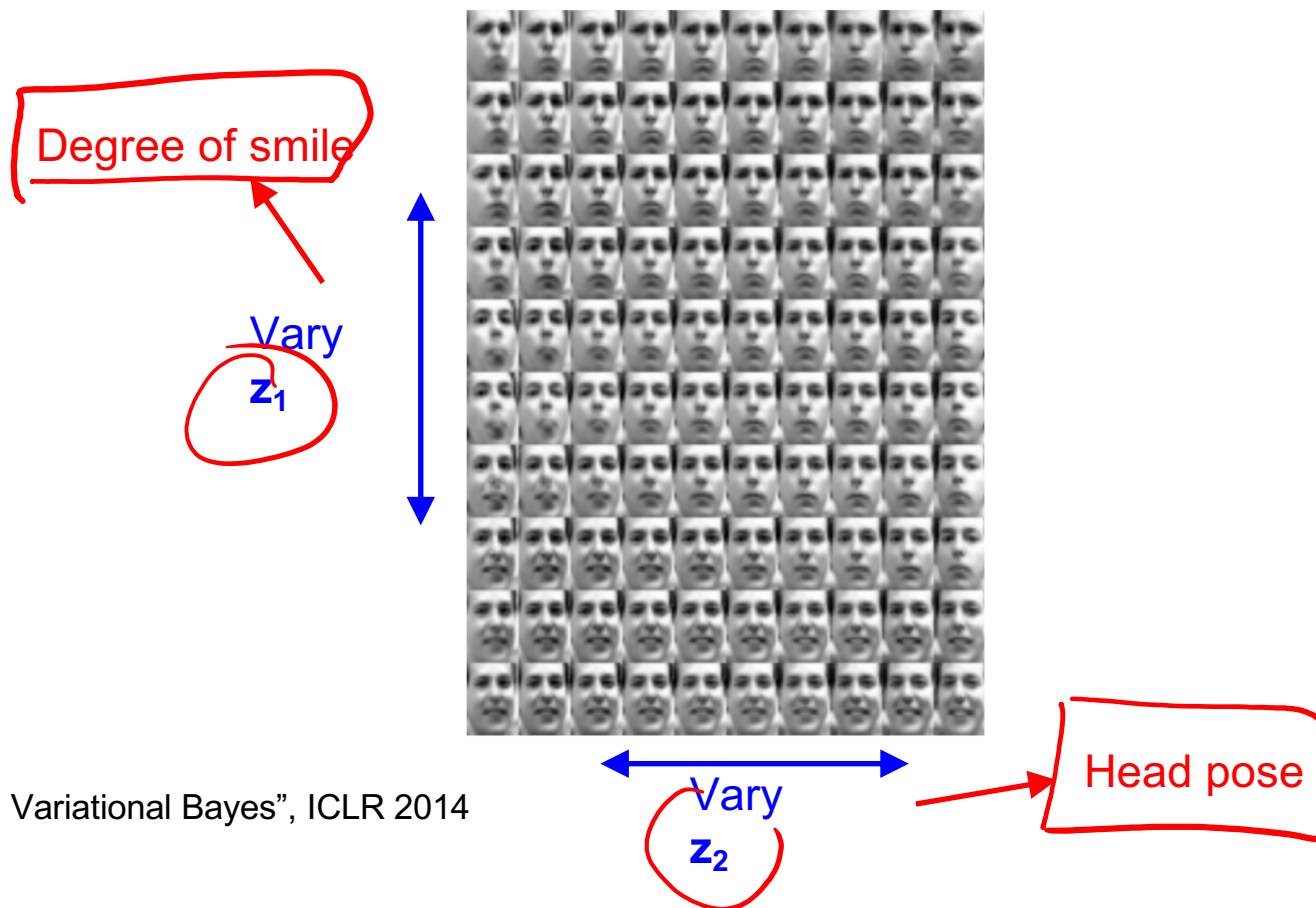
Data manifold for 2-d z



Variational Auto Encoders: Generating Data

Diagonal prior on \mathbf{z}
=> independent
latent variables

Different
dimensions of \mathbf{z}
encode
interpretable factors
of variation



Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Auto Encoders: Generating Data

Diagonal prior on \mathbf{z}
=> independent
latent variables

Different
dimensions of \mathbf{z}
encode
interpretable factors
of variation

Also good feature representation that
can be computed using $q_\phi(\mathbf{z}|\mathbf{x})!$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Degree of smile

Vary
 \mathbf{z}_1



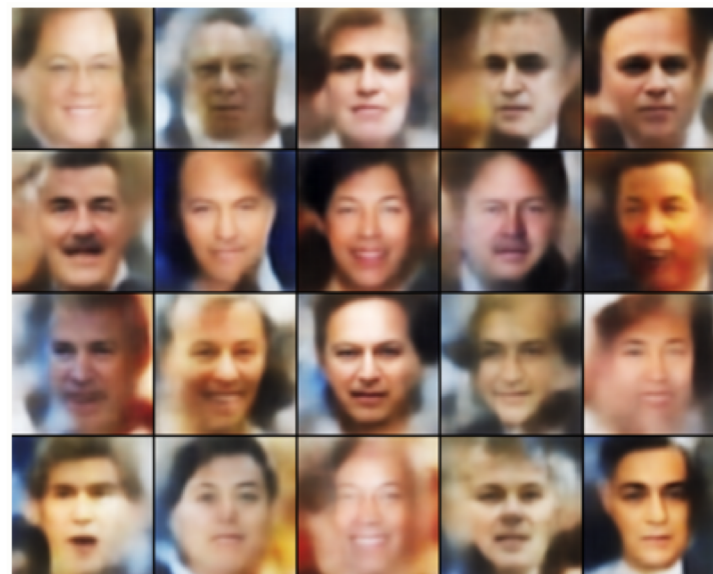
Vary
 \mathbf{z}_2

Head pose

Variational Auto Encoders: Generating Data



32x32 CIFAR-10



Labeled Faces in the Wild



Figures copyright (L) Dirk Kingma et al. 2016; (R) Anders Larsen et al. 2017. Reproduced with permission.

Variational Autoencoders

Probabilistic spin to traditional autoencoders => allows generating data
Defines an intractable density => derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Allows inference of $q(z|x)$, can be useful feature representation for other tasks

Cons:

- Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
- Samples blurrier and lower quality compared to state-of-the-art (GANs)

Active areas of research:

- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian
- Incorporating structure in latent variables

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders
2. Variational Approximation
 - Variational Lower Bound / ELBO
3. Amortized Inference Neural Networks
4. “Reparameterization” Trick