

CS 4803 / 7643: Deep Learning

Topics:


- Variational Auto-Encoders (VAEs)
- Variational Inference, ELBO

Dhruv Batra
Georgia Tech

Administrativa

- Project submission instructions released
 - Due: 12/03, 11:55pm
 - Last deliverable in the class
 - Can't use late days
 - https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/

Recap from last time



Variational Autoencoders (VAE)



So far...

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

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

$$p(\vec{x}) \quad D = \{ \vec{x} \}$$

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

“latent”
“hidden”

conditional prior

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 :

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

z continuous

$$\sum_z p_{\theta}(z) p_{\theta}(x|z)$$

z discrete

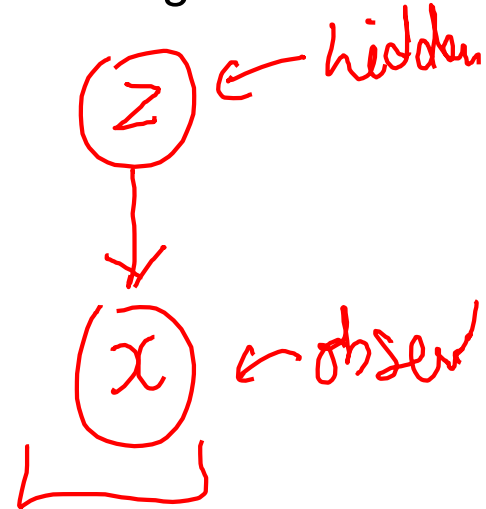
So far...

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 :

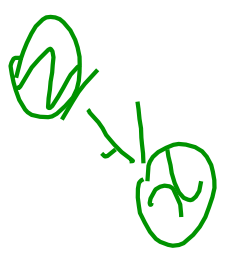
$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$



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

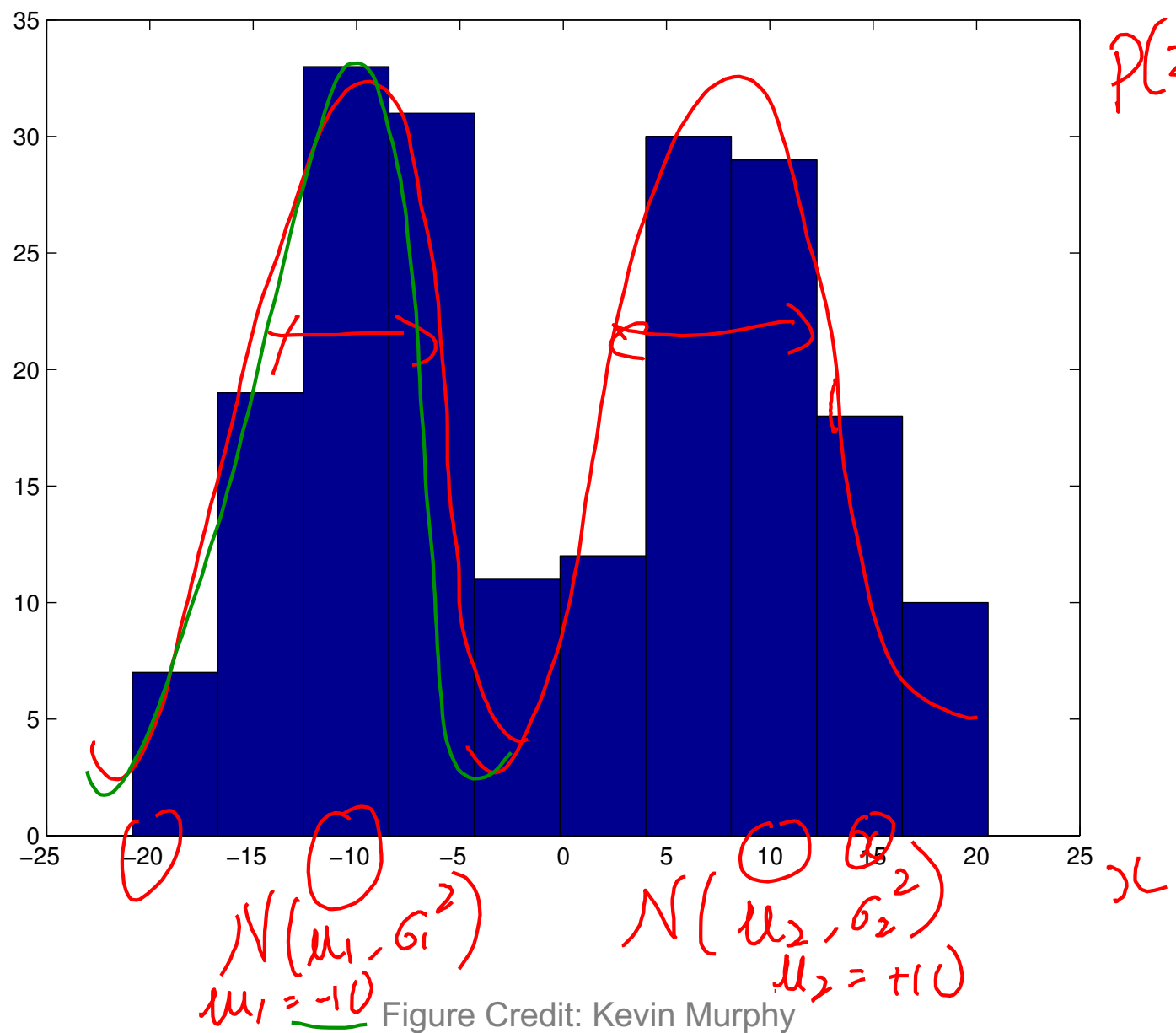
GMM Gaussian Mixture Model

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



$P(z)$

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



Gaussian Mixture Model



$$\underline{Z} \sim \text{Cat}(\vec{\pi})$$

$$\begin{bmatrix} \pi_1 \\ \vdots \\ \pi_k \end{bmatrix}$$

$z \in \{1, \dots, k\}$

$$\pi_c = P(Z=c)$$

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

$$\frac{p(\vec{x})}{p(x, z)} = \underbrace{\sum_{z} p(z)}_{\vec{\pi}_z} \underbrace{p(\vec{x} | z)}_{N}$$

Gaussian Mixture Model

$$P(z) = \pi_z$$

$$P(x|z) = N(\mu_z, \sigma_z^2)$$

available
from
model

$$P(\vec{x}) = \sum_z P(z) P(x|z) \equiv \text{Marginalization}$$

$$P(z|\vec{x}) = \frac{P(z, \vec{x})}{P(\vec{x})} = \frac{P(\vec{x}|z) P(z)}{\sum_z P(\vec{x}|z) P(z)} \equiv \text{"Inference"}$$

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

2. Variational Approximation

- Variational Lower Bound / ELBO

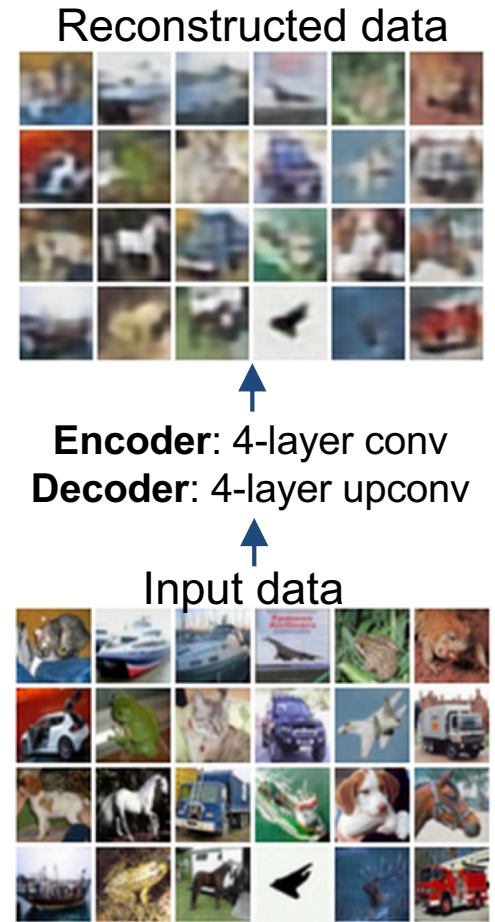
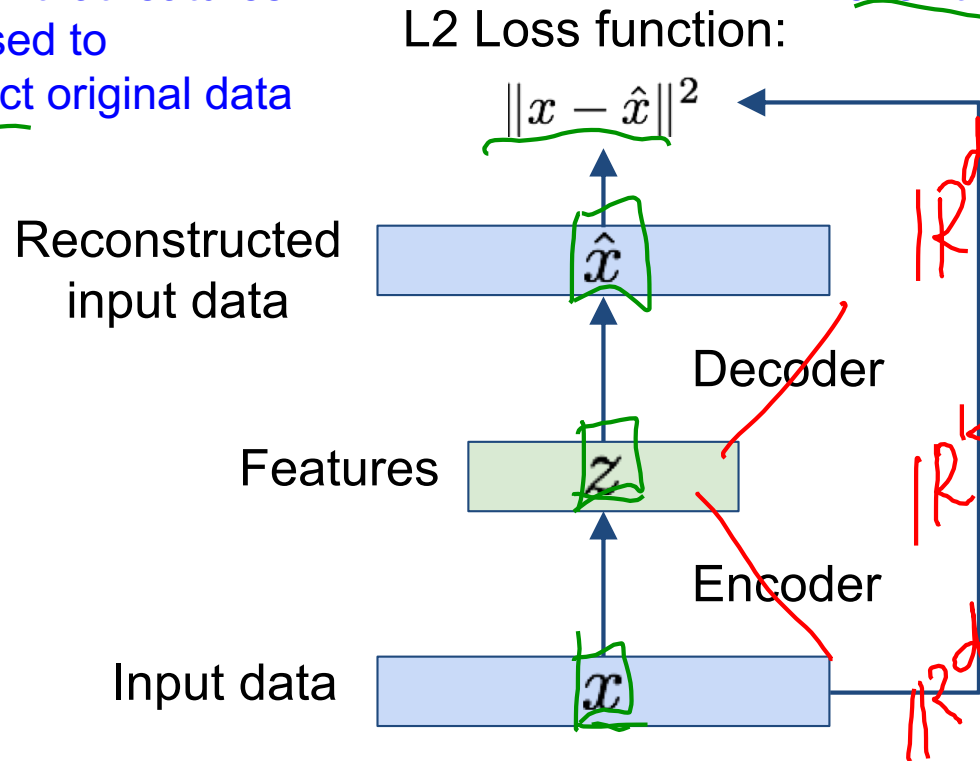
3. Amortized Inference Neural Networks

~~4.~~ “Reparameterization” Trick

Autoencoders

Train such that features can be used to reconstruct original data

Doesn't use labels!



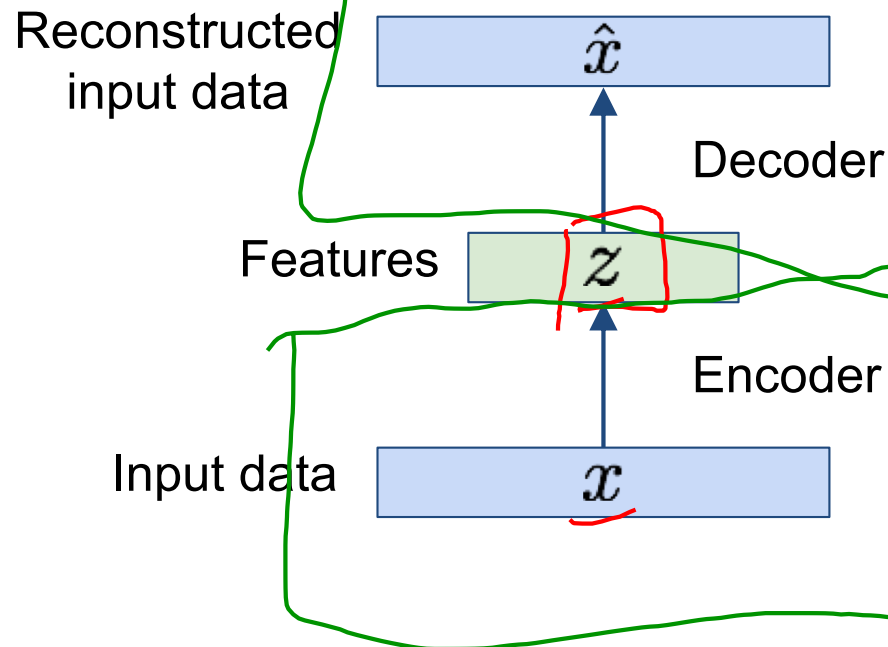
Autoencoders

$$p(y|x)$$

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

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

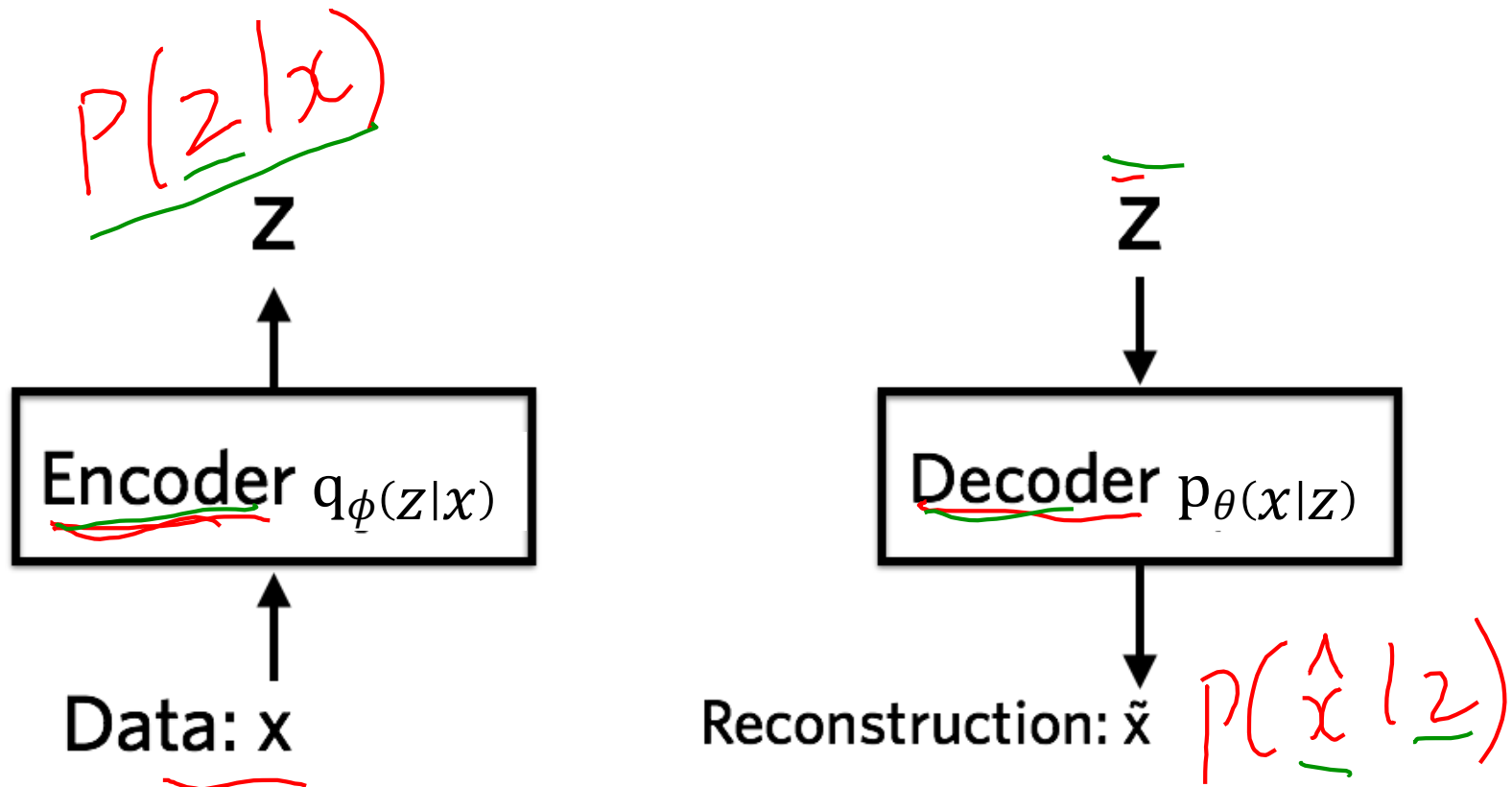
Autoencoders can reconstruct data, and can learn features to initialize a supervised model



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

• $P(z|x) =$
↓
 $q_i(z)$

$$\frac{P(z, x)}{P(x)} = \frac{P(x|z) p(z)}{\int_z P(x|z) p(z)}$$

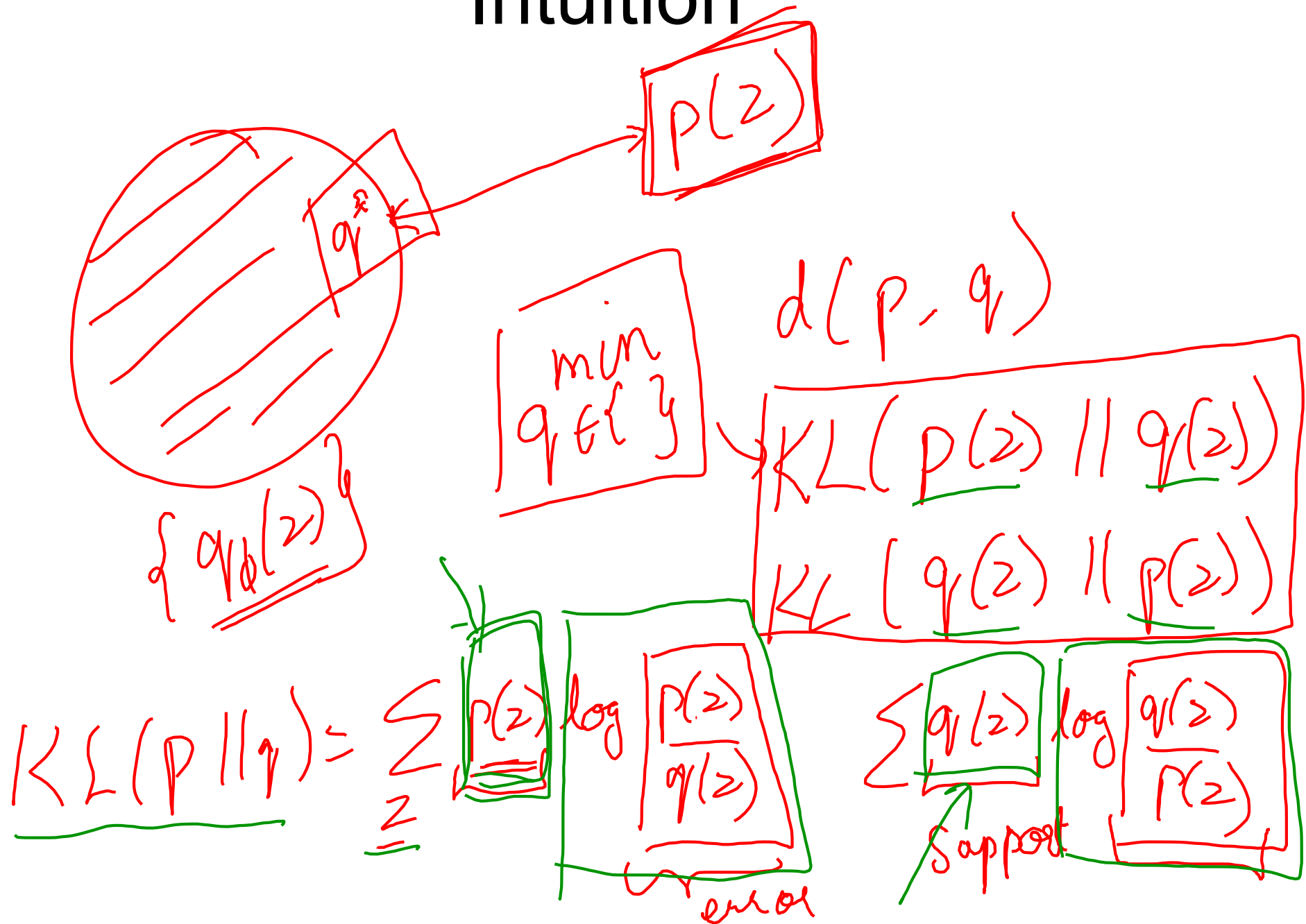
Hard



What is Variational Inference?

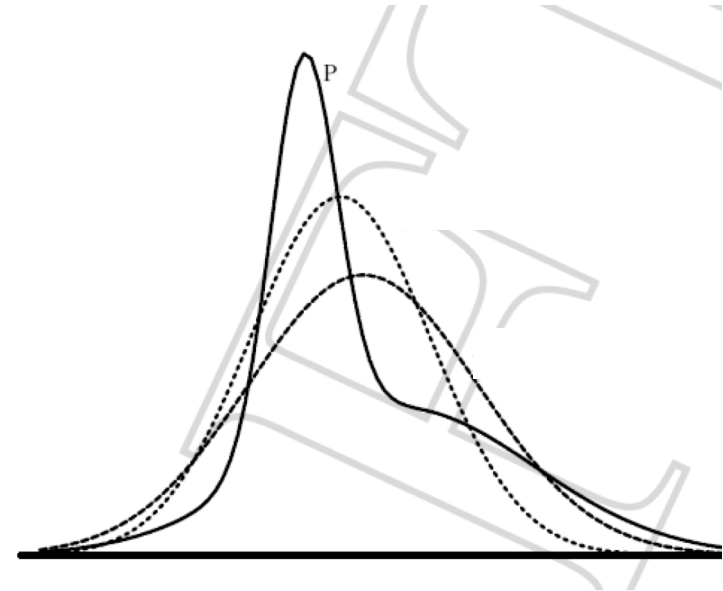
- Key idea
 - Reality is complex
 - Can we approximate it with something “simple”?
 - Just make sure simple thing is “close” to the complex thing.

Intuition



Find simple approximate distribution

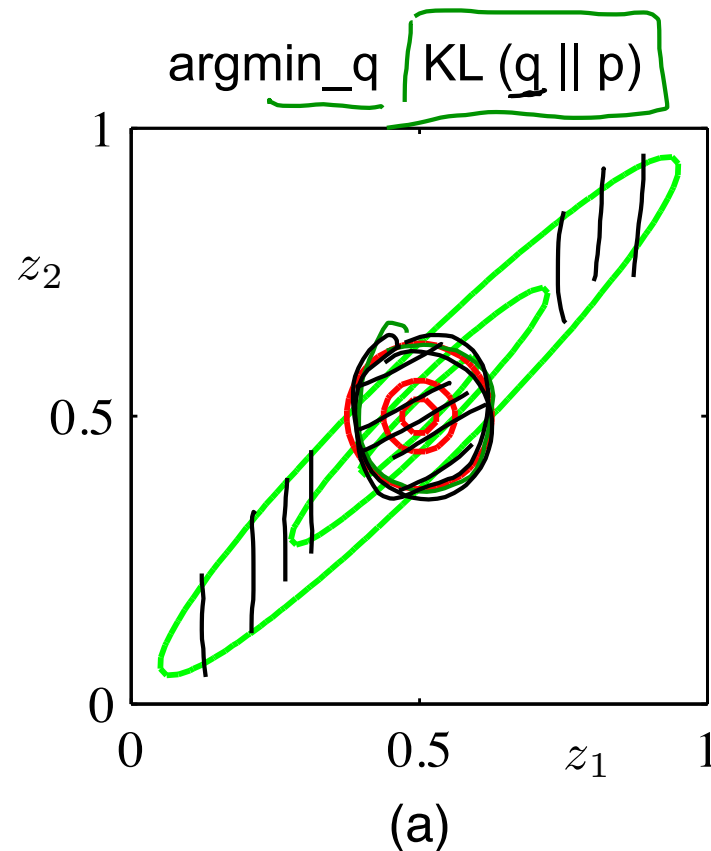
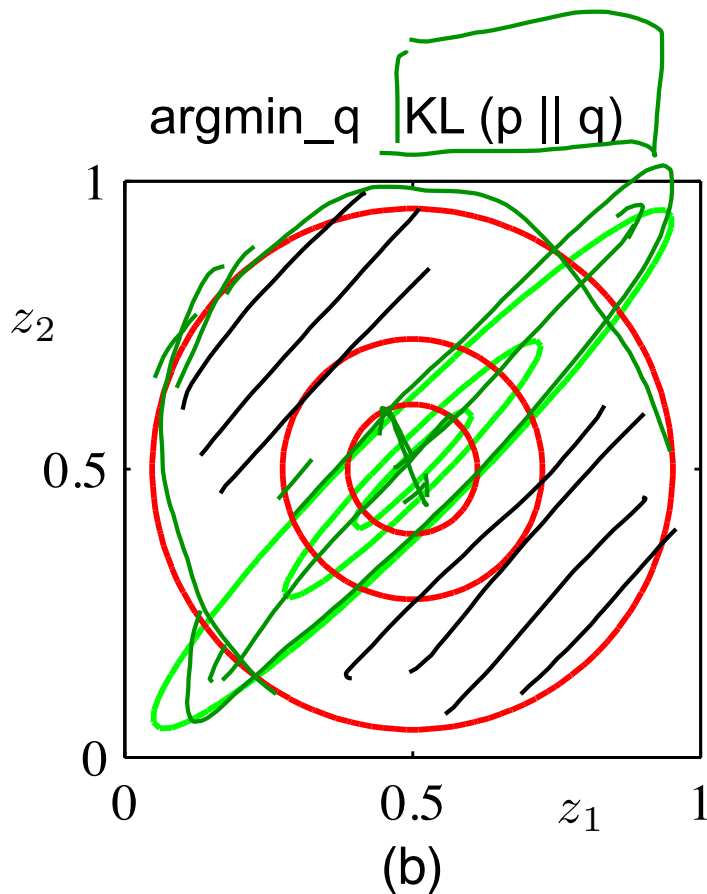
- Suppose p is intractable posterior
- Want to find simple q that approximates p
- KL divergence not symmetric
- $D(p||q)$
 - true distribution p defines support of diff.
 - the “correct” direction
 - will be intractable to compute
- $D(q||p)$
 - approximate distribution defines support
 - tends to give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



Example 1

- $p = 2D$ Gaussian with arbitrary co-variance (Σ)
- $q = 2D$ Gaussian with isotropic co-variance ($\sigma^2 I$)

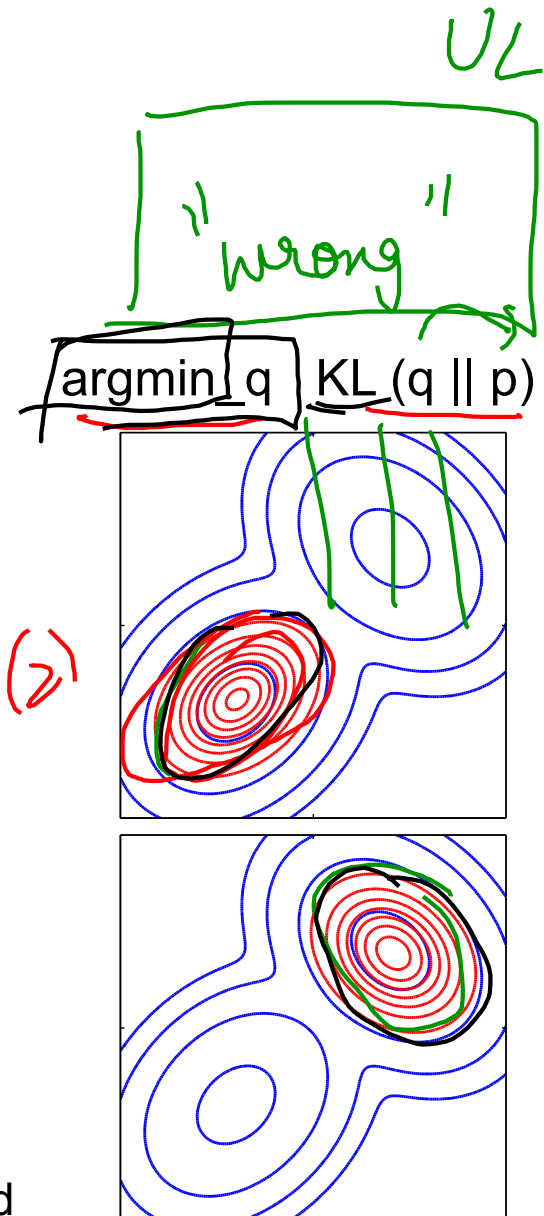
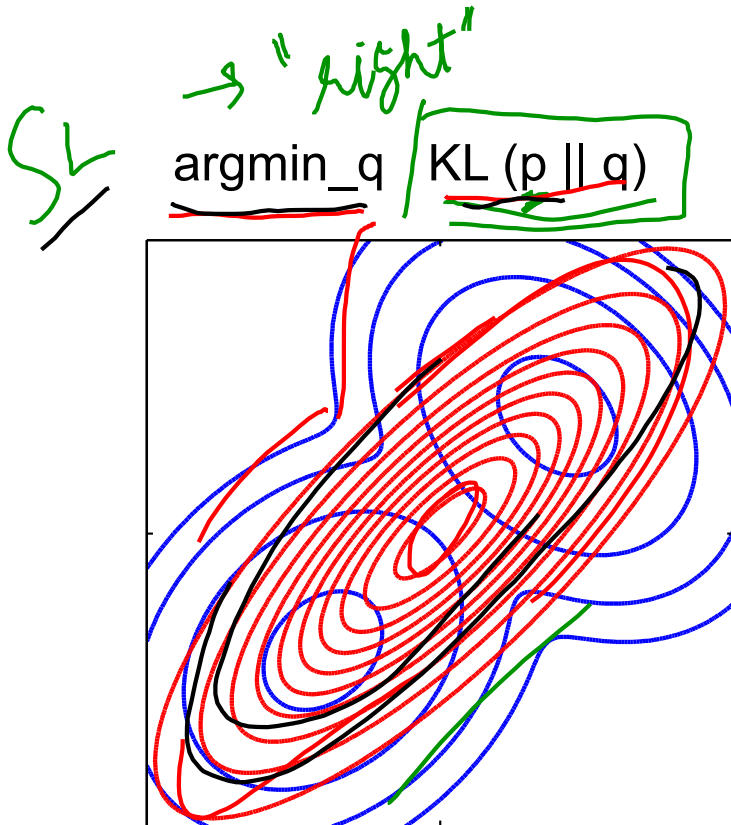
$$\begin{bmatrix} \cdot & \cdot \\ 0 & \sigma^2 \end{bmatrix}_{2 \times 2}$$



p = Green; q = Red

Example 2

- p = Mixture of Two Gaussians
- q = Single Gaussian μ, σ



Plan for Today

- VAEs
 - Variational Inference → Evidence Based Lower Bound
 - Putting it all together
- Next time:
 - Reparameterization trick for optimizing VAEs |

The general learning problem with missing data

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

$$\begin{aligned}
 \underline{ll}(\theta : \mathcal{D}) &= \underline{\log} \left[\prod_{i=1}^N P(\mathbf{x}_i | \theta) \right] \\
 &= \sum_{i=1}^N \log P(\mathbf{x}_i | \theta) \\
 &= \sum_{i=1}^N \log \left[\sum_{\mathbf{z}} P(\mathbf{x}_i, \mathbf{z} | \theta) \right]
 \end{aligned}$$

$\mathcal{D} = \{ \vec{x}_1, \dots, \vec{x}_N \}$

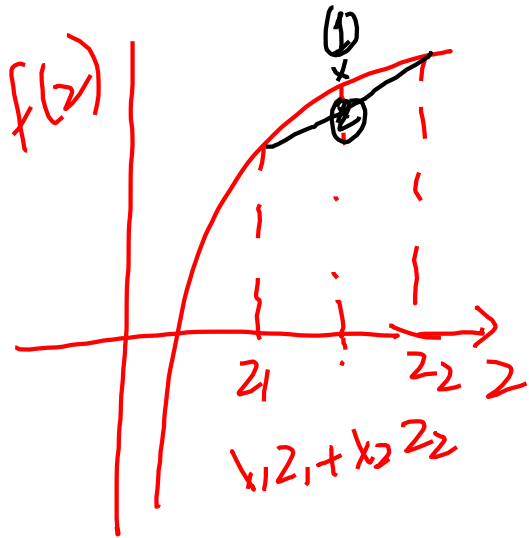
$P(\vec{x}, \mathbf{z})$

$P(x_i | \theta) P(\mathbf{z} | x_i, \theta)$

$\log \sum_{\mathbf{z}} P(\mathbf{z} | x_i, \theta) P(x_i | \theta)$
 $\approx \log \left[\mathbb{E}_{\mathbf{z}} [P(x_i | \theta)] \right]$

Applying Jensen's inequality

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



$$\textcircled{1} \geq \textcircled{2}$$

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

$$f\left(\sum_{i=1}^{k(2)} \lambda_i z_i\right) \geq \sum_{i=1}^{k(2)} \lambda_i f(z_i) \quad \leftarrow z \rightarrow k$$

$$\begin{cases} \lambda_1, \lambda_2 \geq 0 \\ \lambda_1 + \lambda_2 = 1 \end{cases}$$

$$\equiv f(E[z]) \geq E_{P(z)}[f(z)] \quad \leftarrow z \rightarrow g(z)$$

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

Applying Jensen's inequality

(z)
↓
(x)

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

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

$$\ell(\theta) \geq \underbrace{\sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\vec{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}}_{\text{"Free Energy"} F(\theta, Q_i)}$$

Variational Lower Bound
Evidence Lower Bound (ELBO)

Evidence Lower Bound

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

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

(VAEs)

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

(GMMs)

$\rightarrow P(z | x_i, \theta) P(x_i | \theta)$

ELBO: Factorization #1 (GMMs)

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

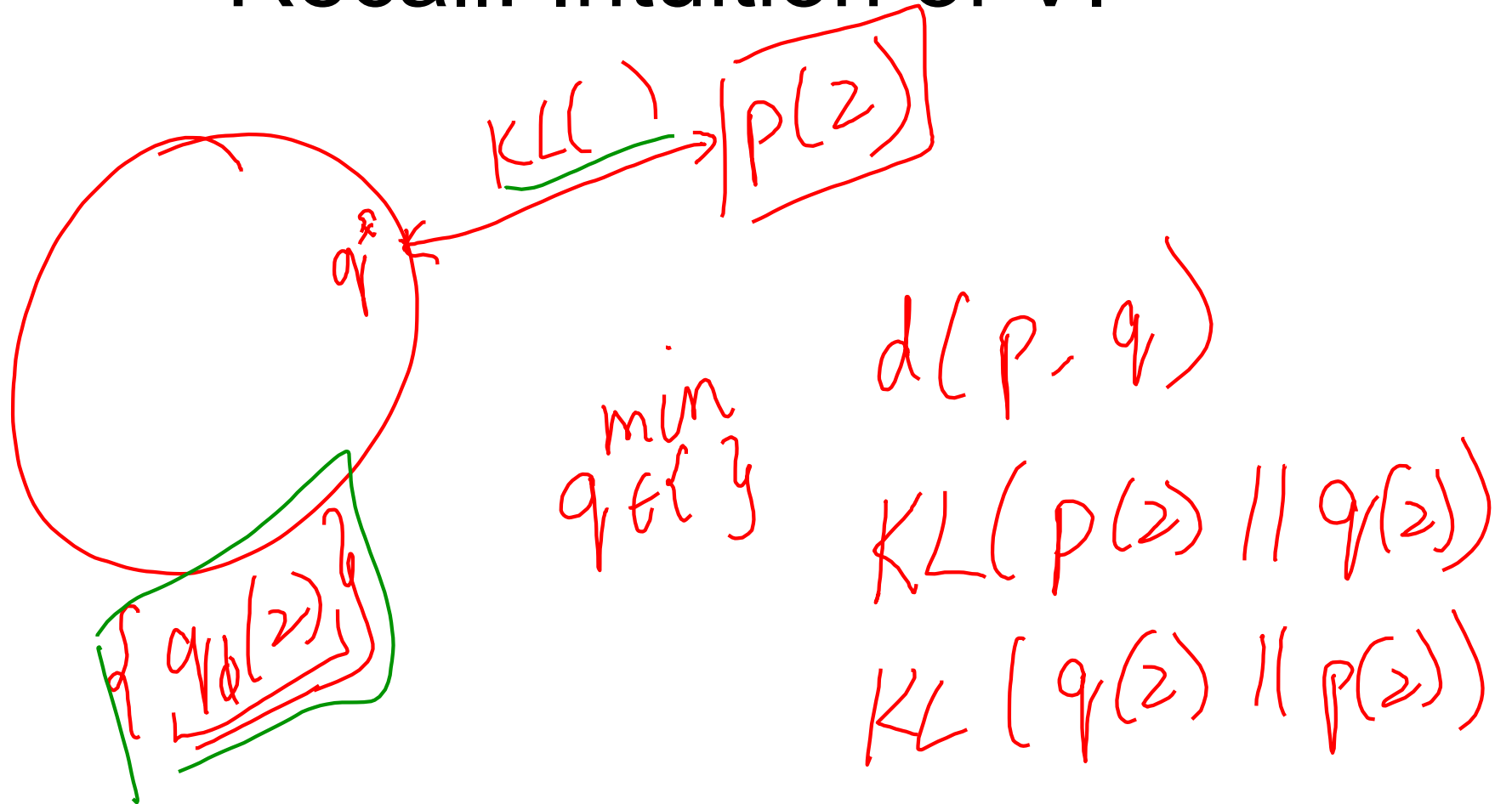
$$P(\tilde{\mathbf{x}}_i | \theta) P(\mathbf{z} | \tilde{\mathbf{x}}_i, \theta)$$

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

$$F(\theta, Q_i) = \underbrace{\log P(\tilde{\mathbf{x}}_i | \theta)}_{l(\theta)} - \text{KL} \left(\underbrace{Q_i(\mathbf{z})}_{\text{approx}} \parallel \underbrace{P(\mathbf{z} | \tilde{\mathbf{x}}_i, \theta)}_{\text{target}} \right)$$

$$l(\theta) \geq \underline{F(\theta, Q_i)} \quad \text{KL} \left(\dots \right) \downarrow$$

Recall: Intuition of VI



ELBO: Factorization #1 (GMMs)

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

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

- E-step: Fix $\theta^{(t)}$, maximize F over Q_i
- M-step: Fix $Q_i^{(t)}$, maximize F over θ

EM for Learning GMMs

- Simple Update Rules
- **E-step:** Fix $\theta^{(t)}$, maximize F over Q_i

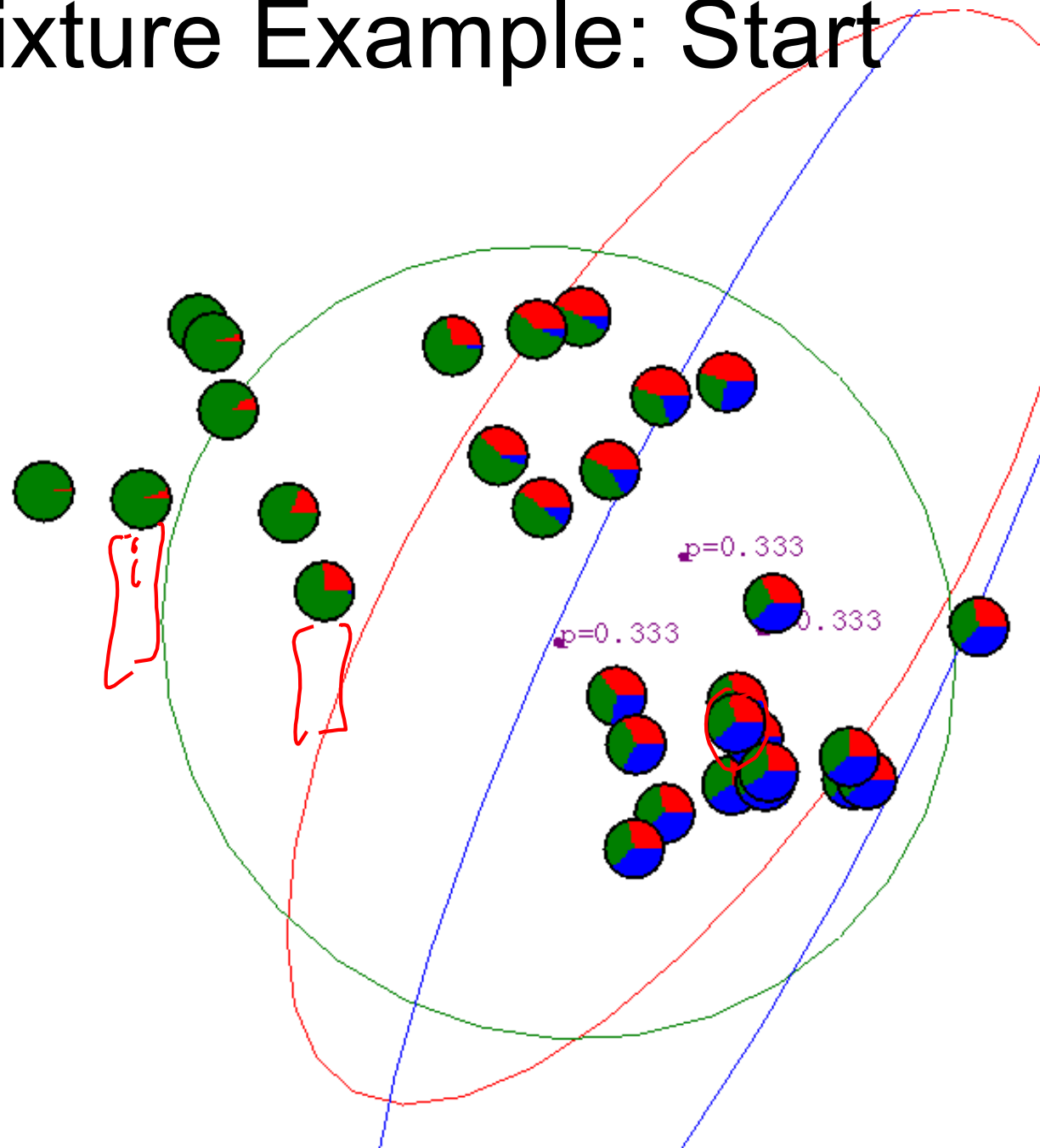
$$Q_i^{(t)}(\mathbf{z}) = P(\mathbf{z} \mid \mathbf{x}_i, \theta^{(t)})$$

- **M-step:** Fix $Q_i^{(t)}$, maximize F over θ
 - maximize expected likelihood under $Q_i(\mathbf{z})$
 - Corresponds to weighted dataset:
 - $\langle \mathbf{x}_1, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 \mid \mathbf{x}_1)$
 - $\langle \mathbf{x}_1, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2 \mid \mathbf{x}_1)$
 - $\langle \mathbf{x}_1, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3 \mid \mathbf{x}_1)$
 - $\langle \mathbf{x}_2, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 \mid \mathbf{x}_2)$
 - $\langle \mathbf{x}_2, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2 \mid \mathbf{x}_2)$
 - $\langle \mathbf{x}_2, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3 \mid \mathbf{x}_2)$

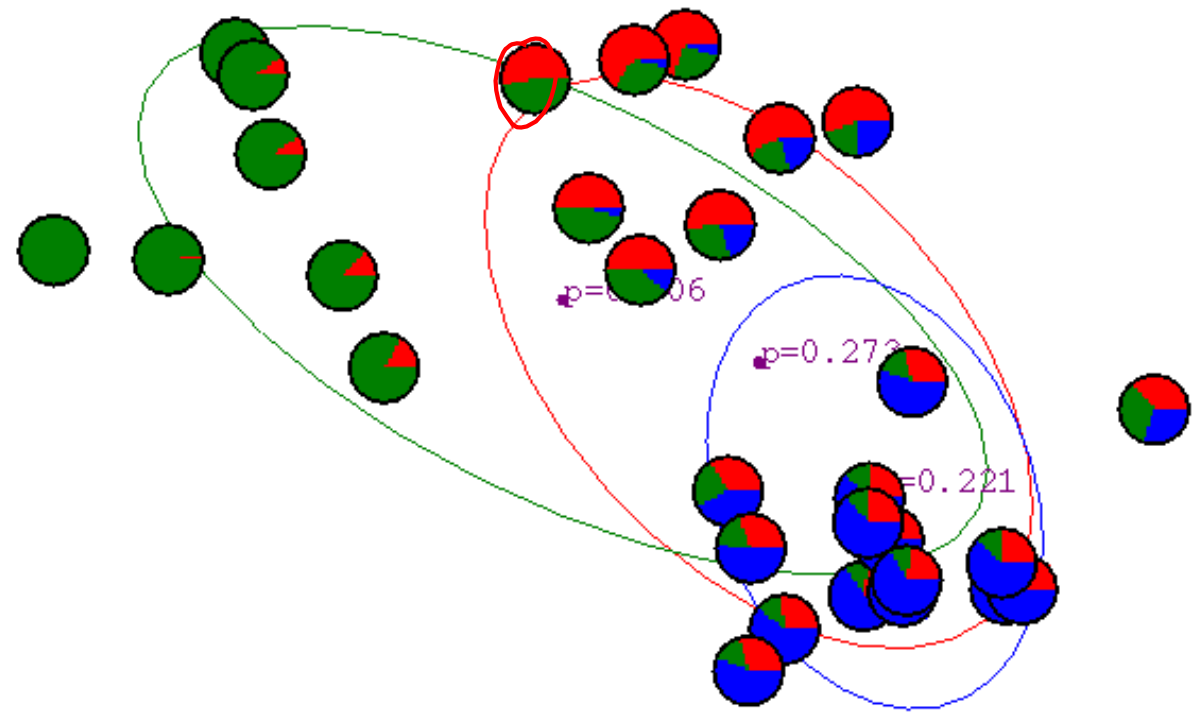
Gaussian Mixture Example: Start

$$\theta = \{\mu_i, \Sigma_i\}_{i=1}^k \quad k=3$$

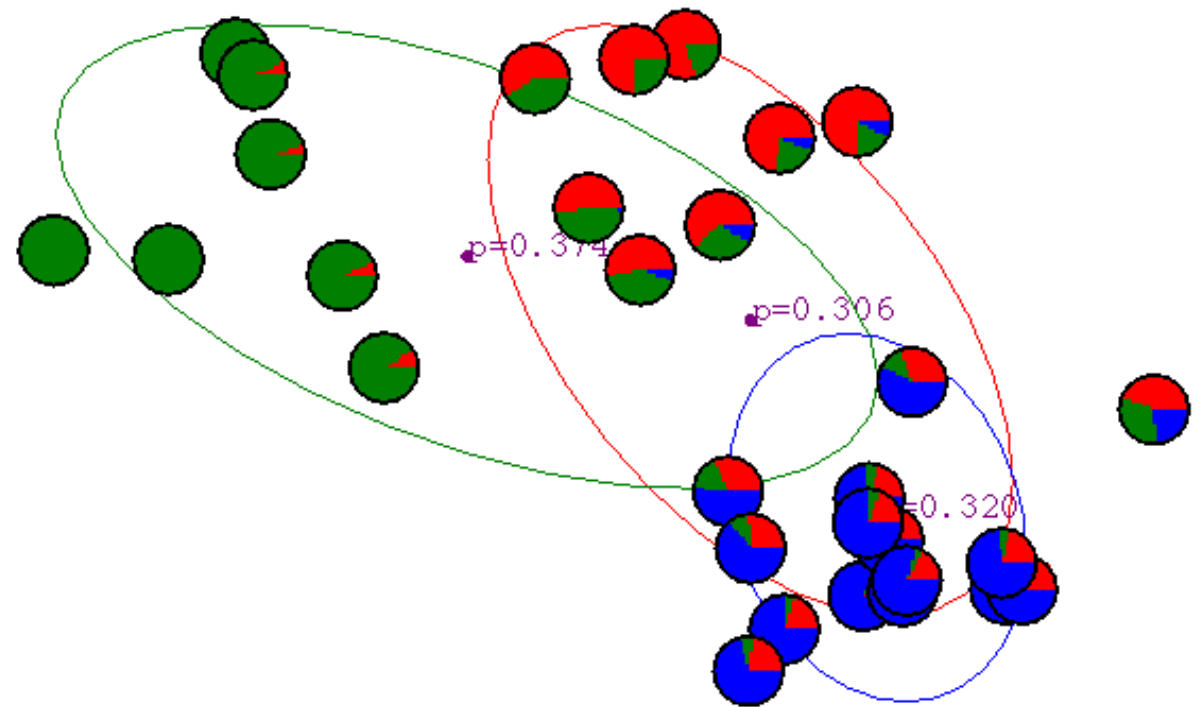
$$q(z|x_i) = \begin{cases} 0.3 \\ 0.3 \\ 0.3 \end{cases}$$



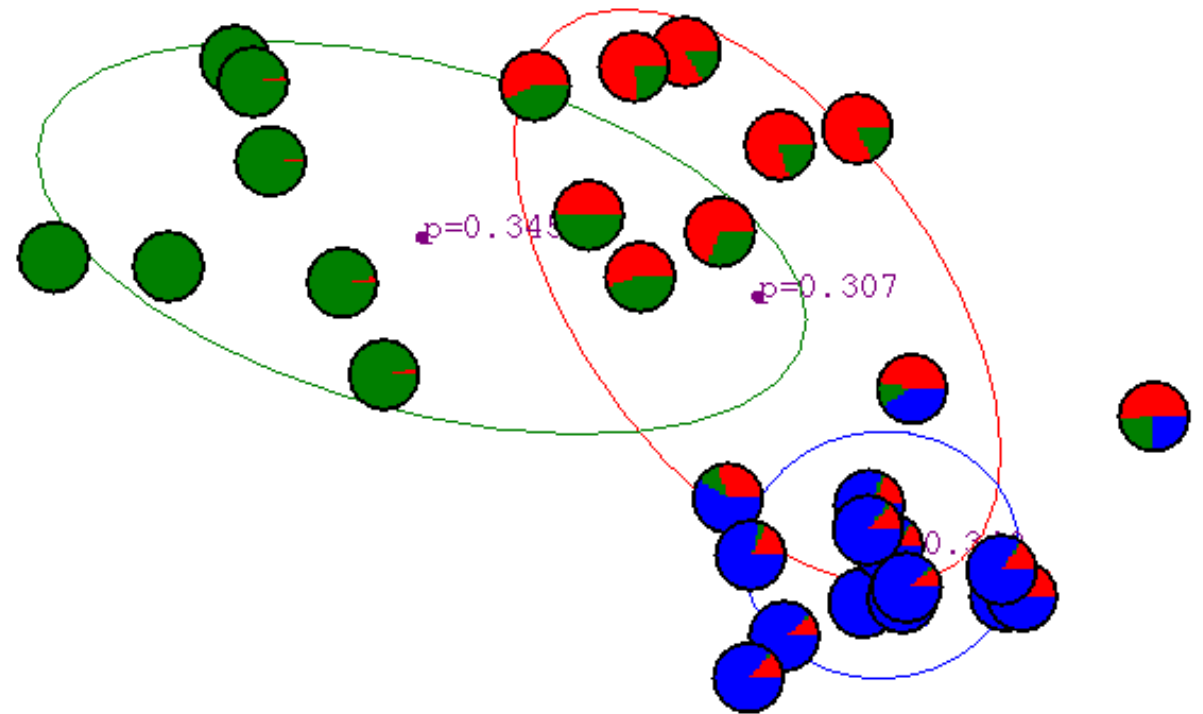
After 1st iteration



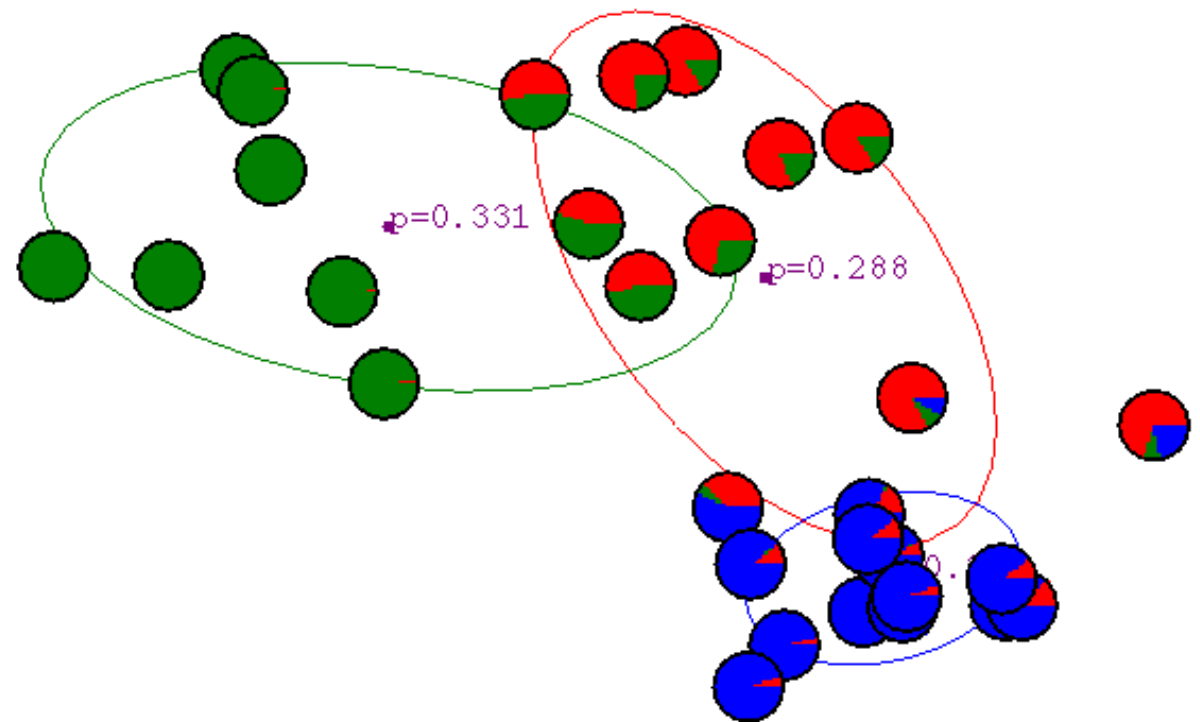
After 2nd iteration



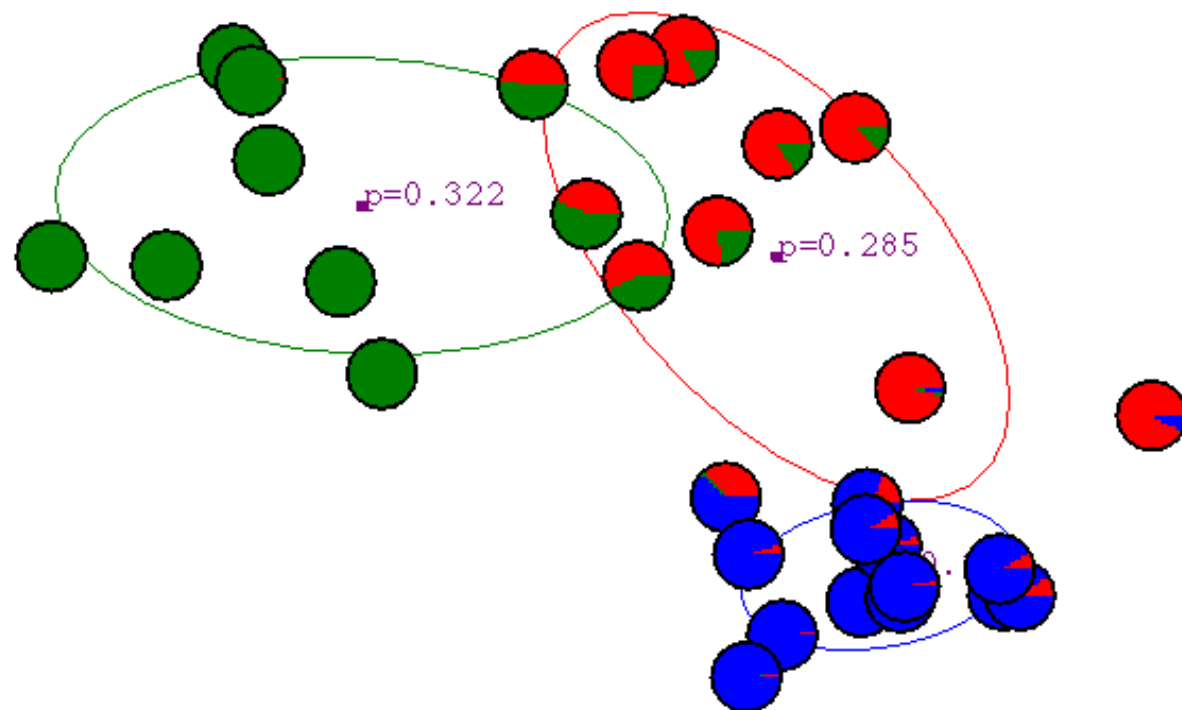
After 3rd iteration



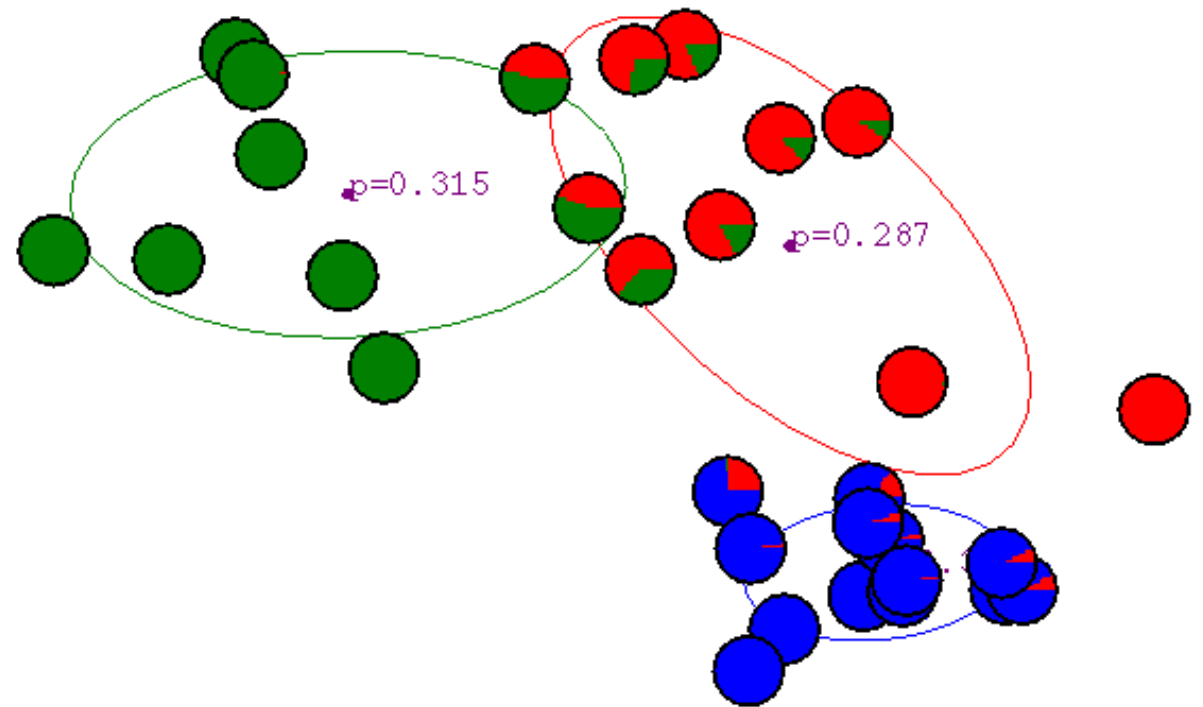
After 4th iteration



After 5th iteration

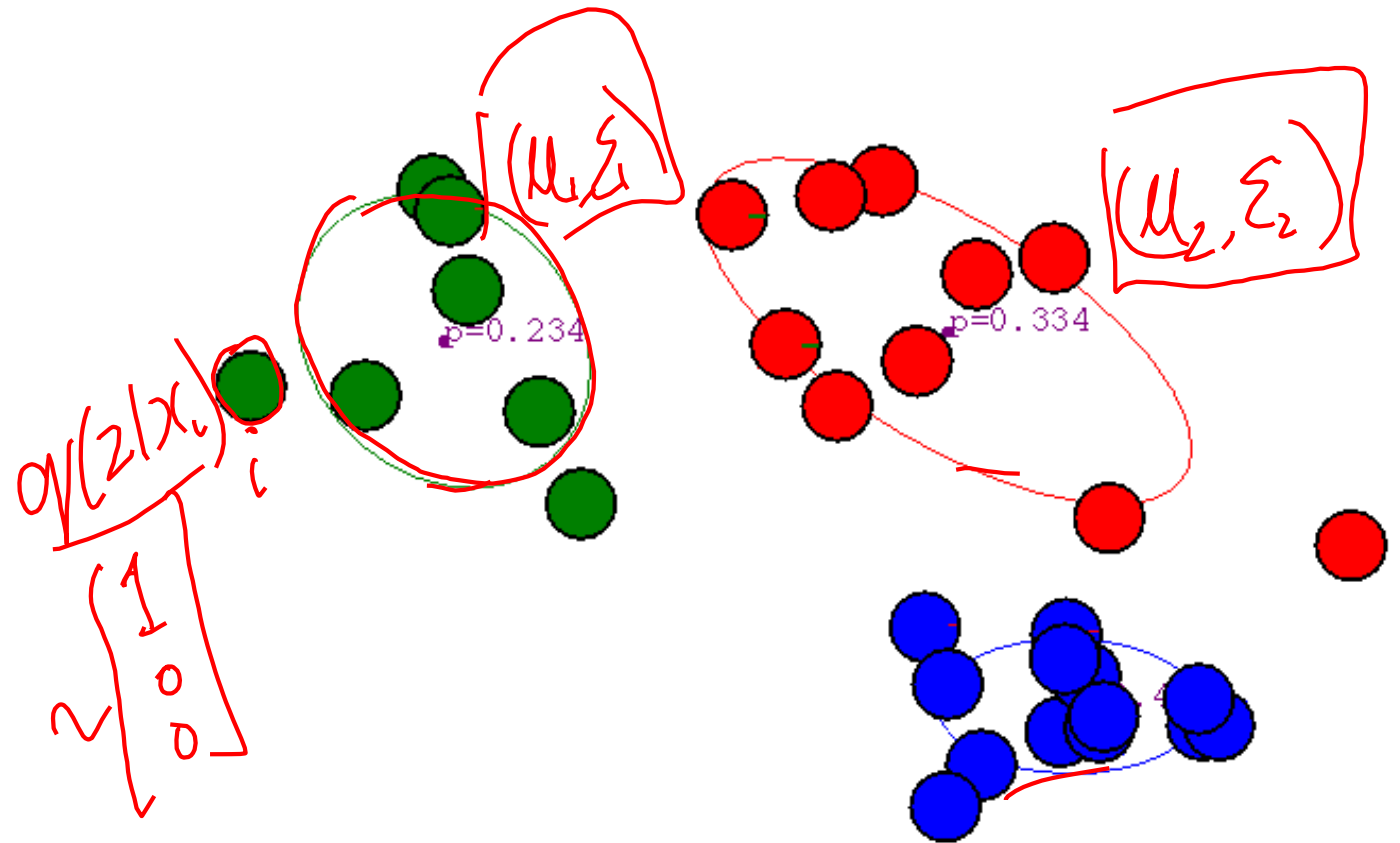


After 6th iteration



After 20th iteration

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



ELBO: Factorization #2 (VAEs)

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

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

(VAEs)

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

“Explain the data”

“Regulariser”

Be simple

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

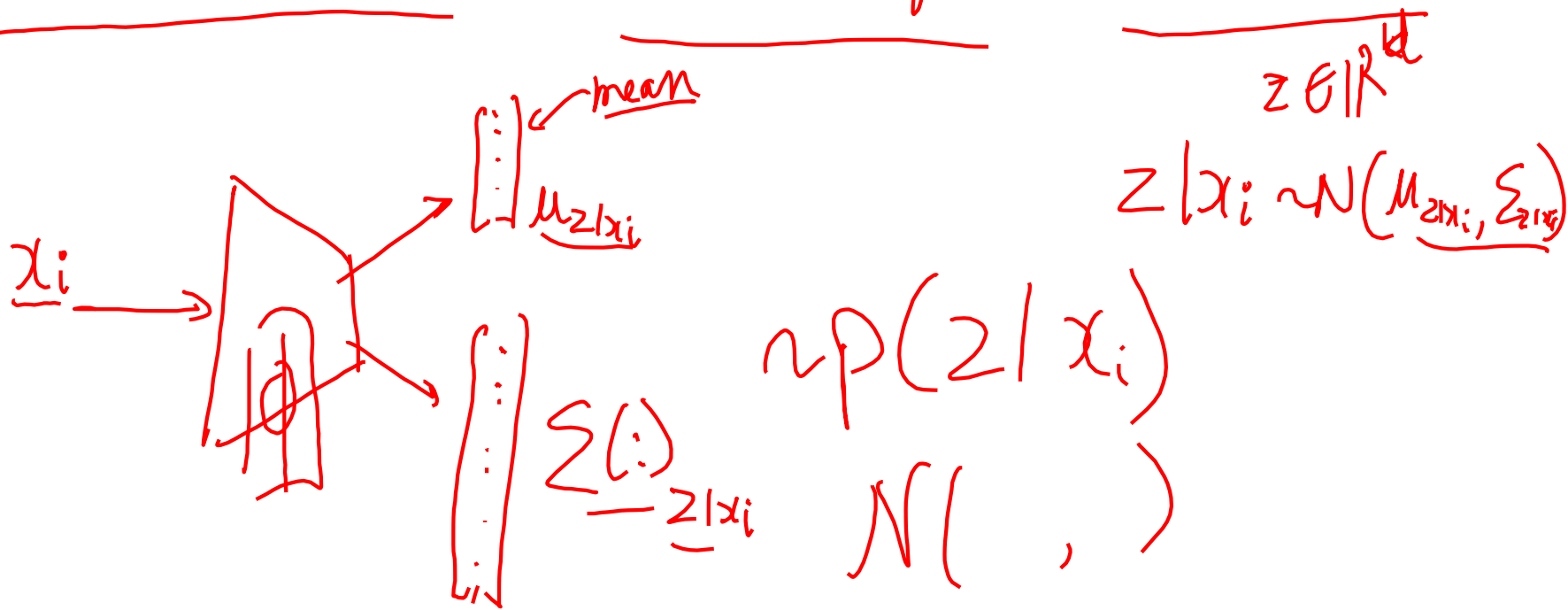
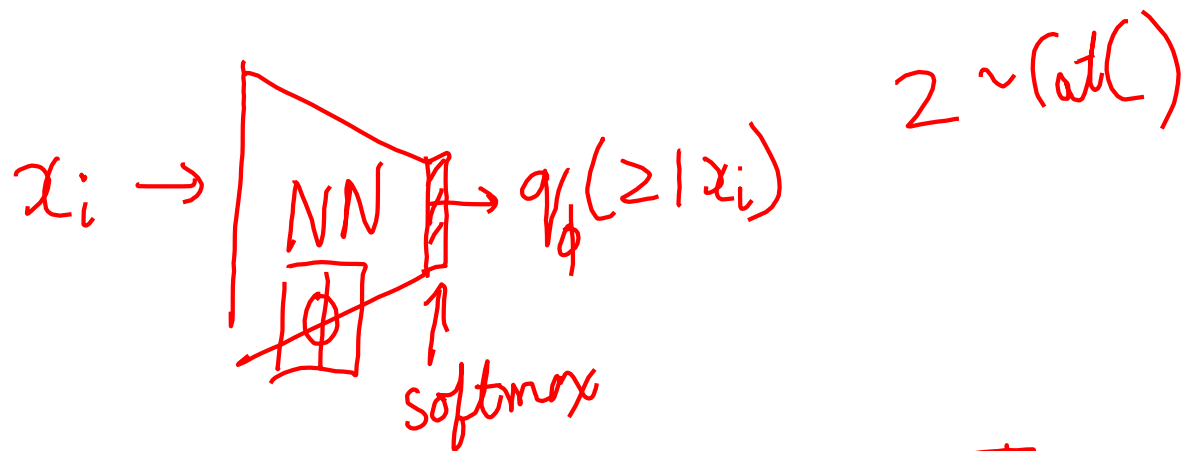
2. Variational Approximation
• Variational Lower Bound / ELBO

3. Amortized Inference Neural Networks

4. “Reparameterization” Trick

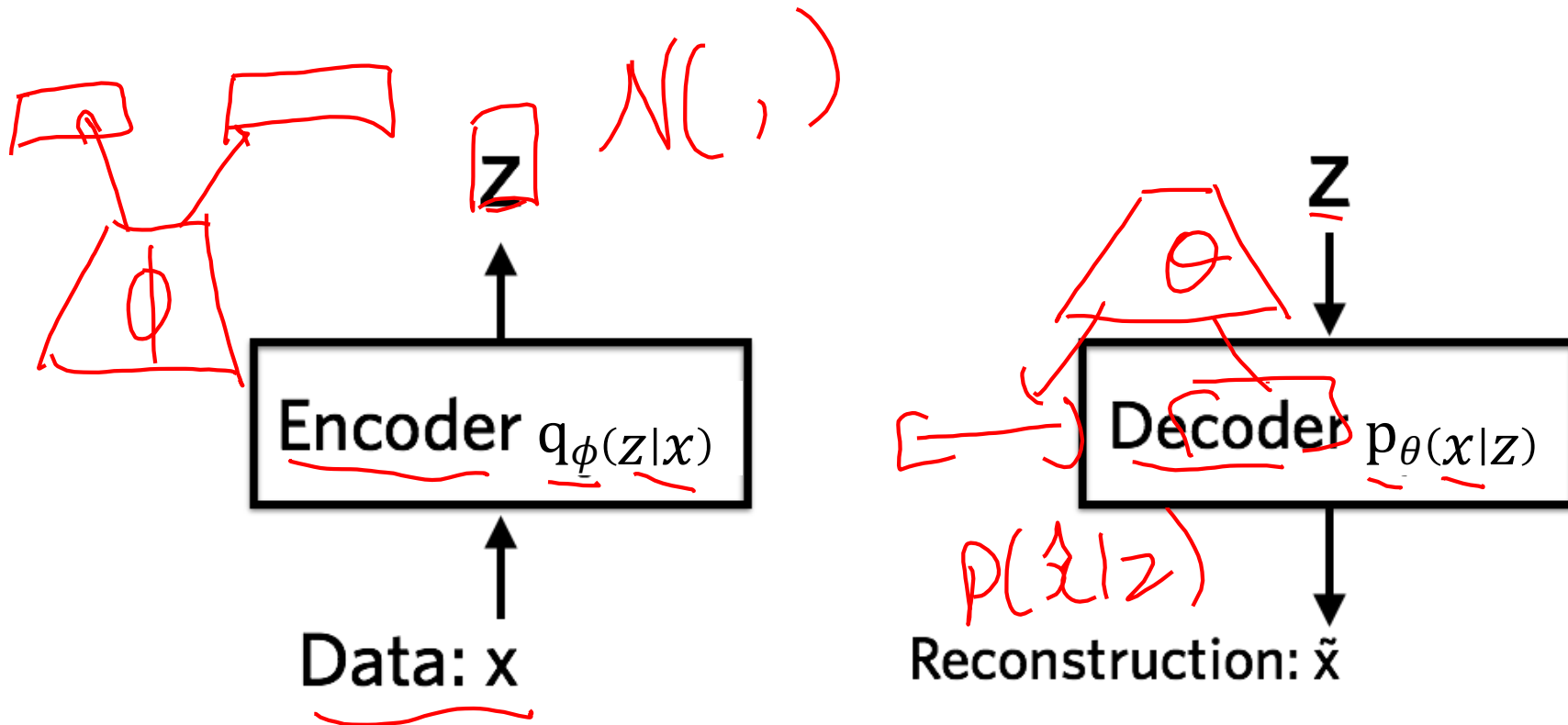
Amortized Inference Neural Networks

$$Q_i(z) = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.3 \end{bmatrix}$$

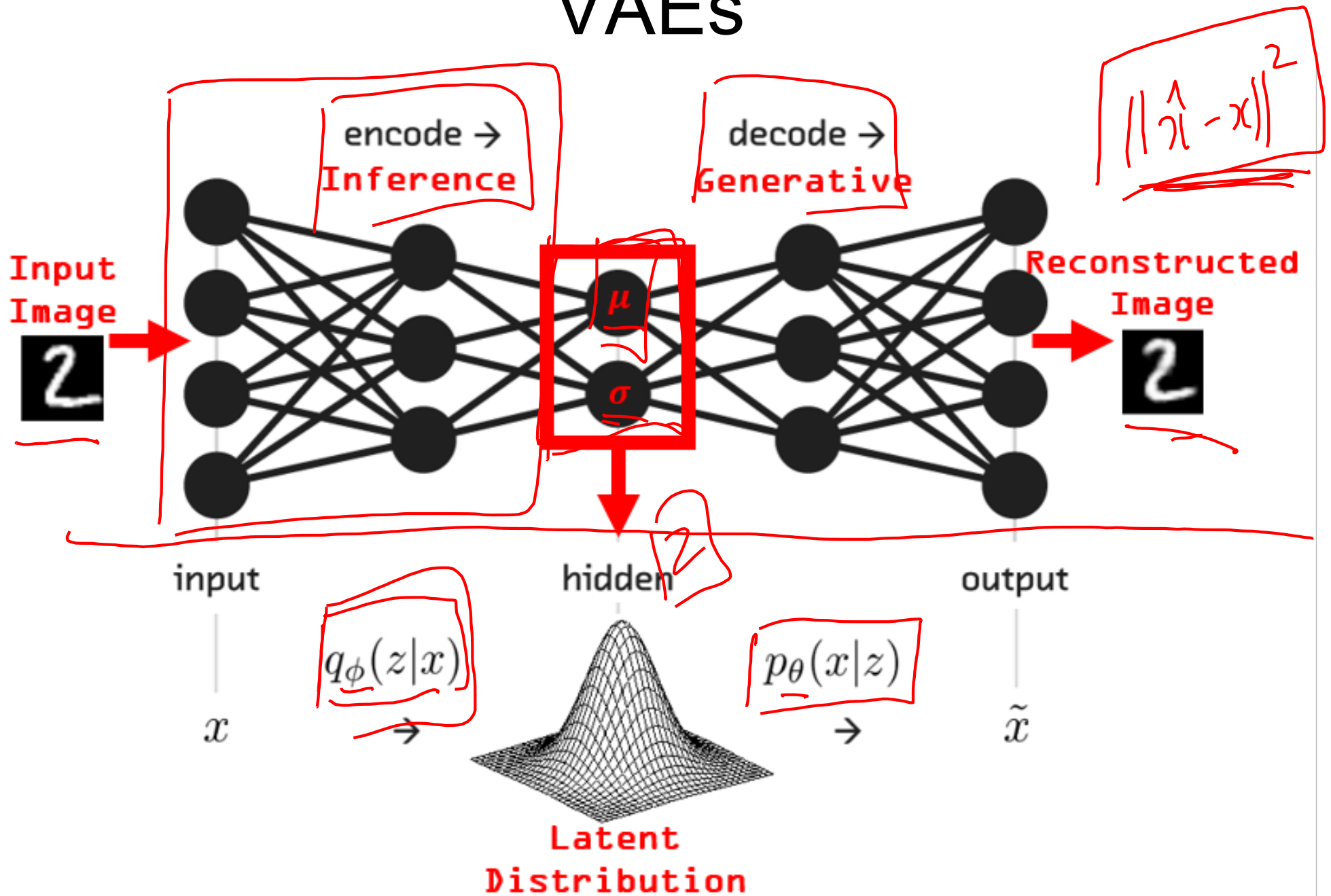


Variational Autoencoders

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



VAEs



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

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

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

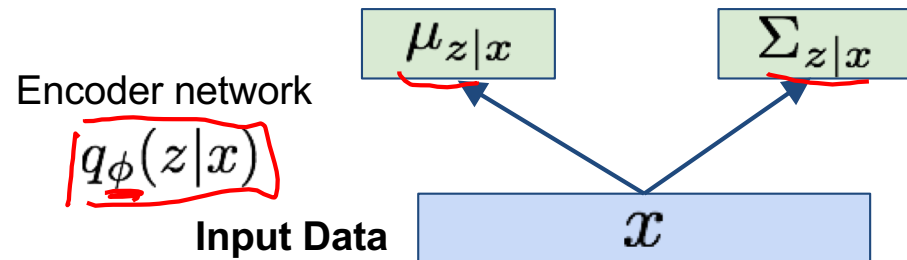
Input Data

\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

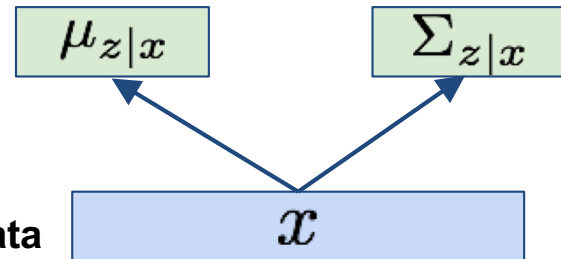
$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Encoder network

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

Input Data

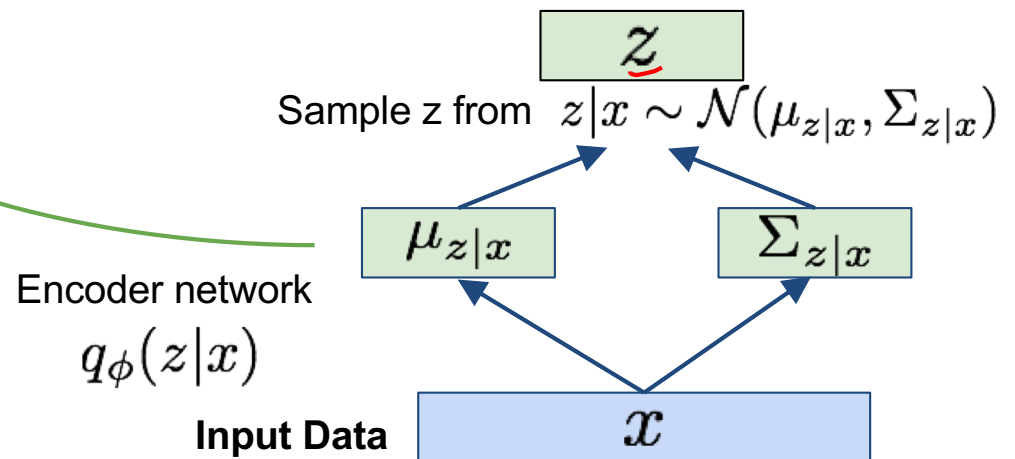


Variational Auto Encoders

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

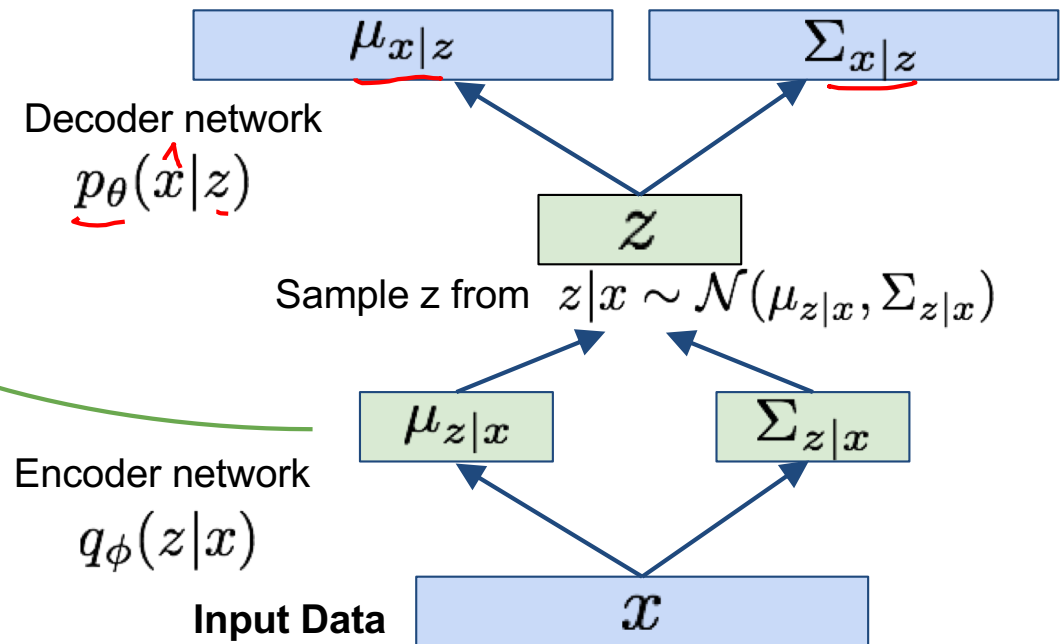


Variational Auto Encoders

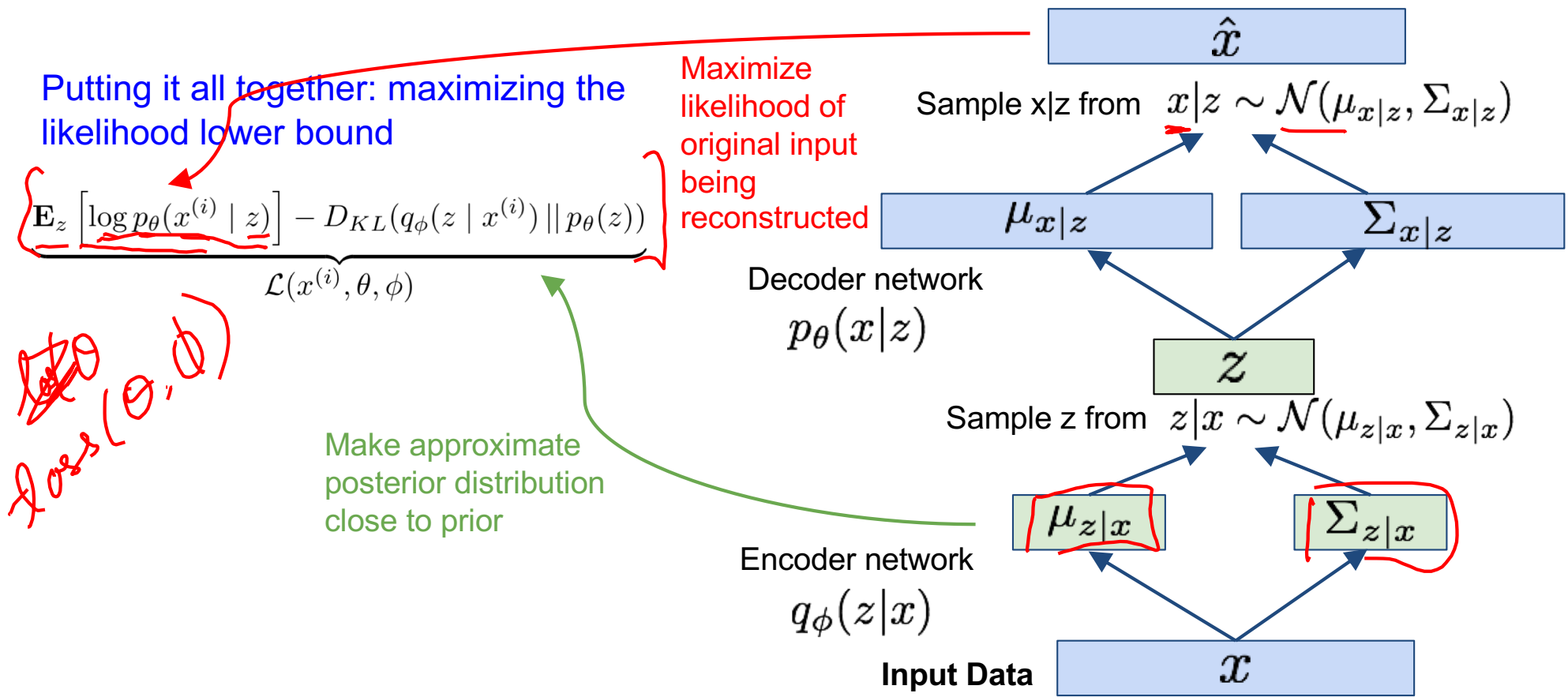
Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior



Variational Auto Encoders



Variational Auto Encoders

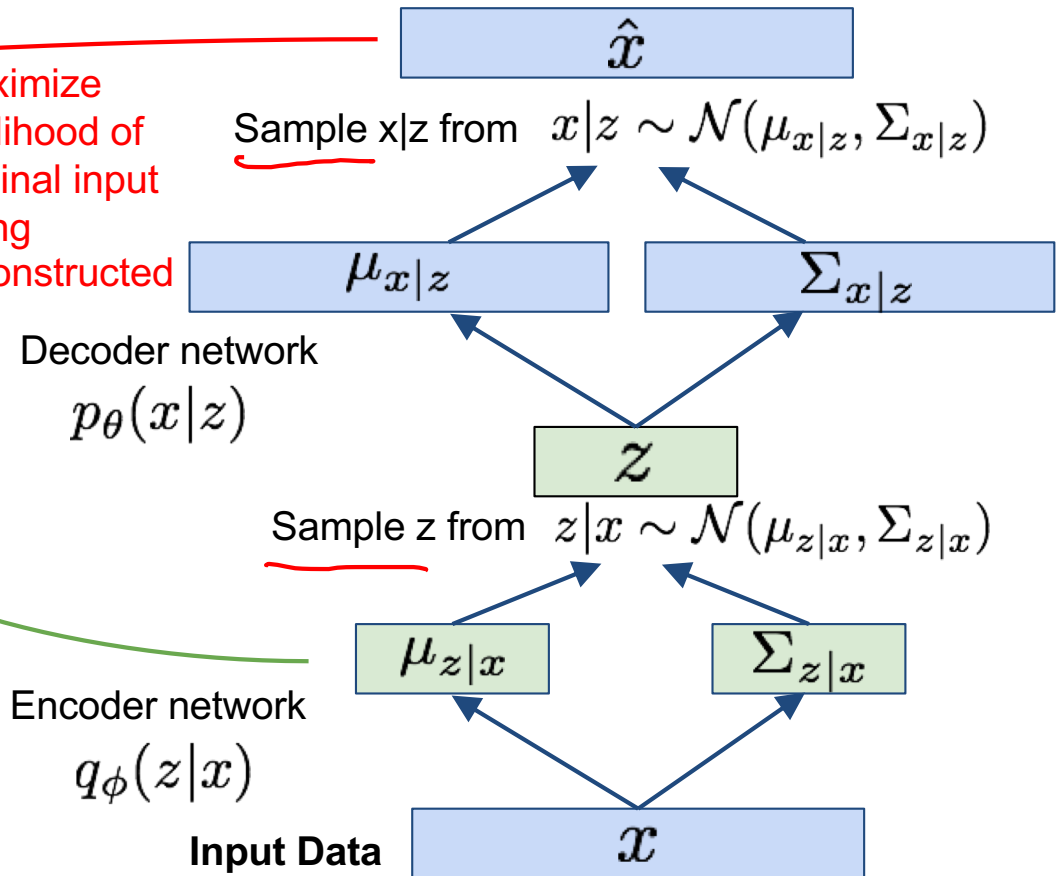
Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

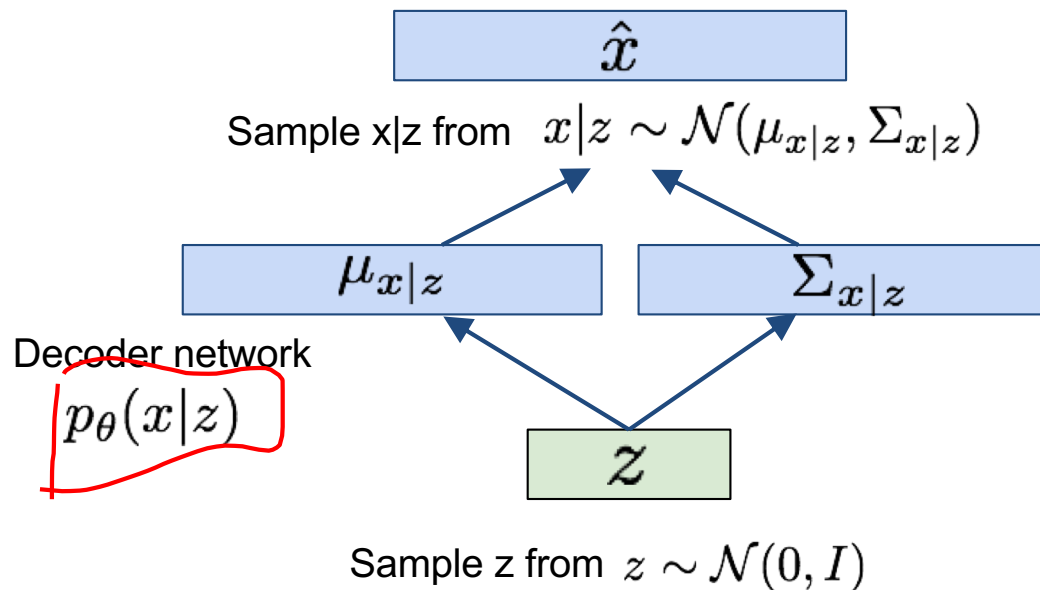
For every minibatch of input data: compute this forward pass, and then backprop!

Maximize likelihood of original input being reconstructed



Variational Auto Encoders: Generating Data

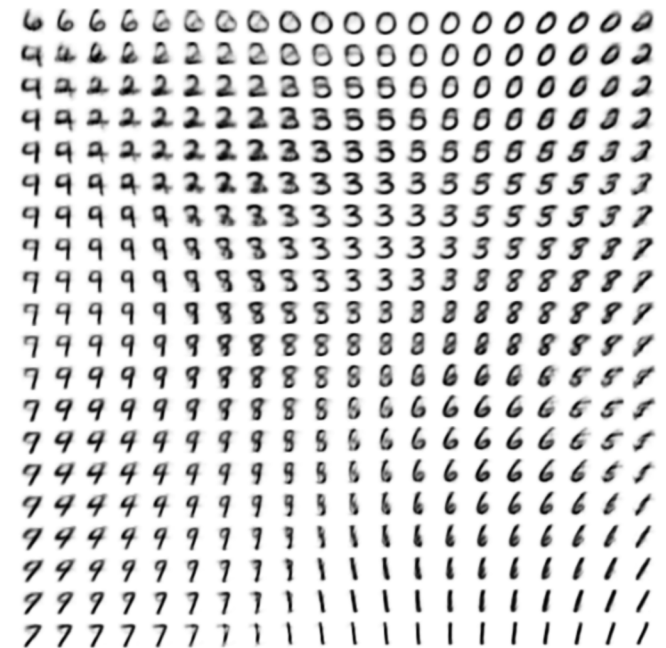
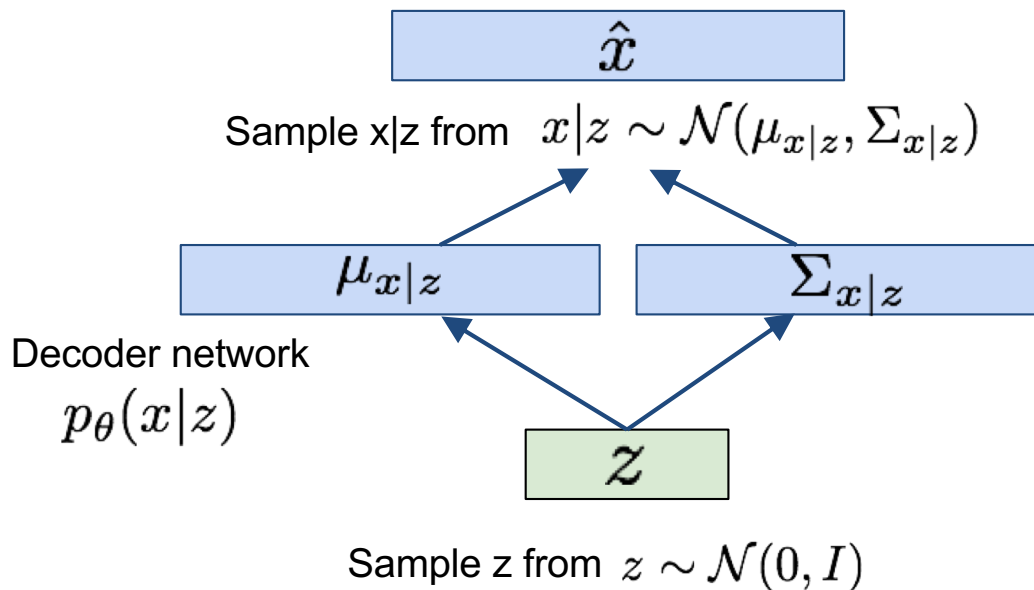
Use decoder network. Now sample z from prior!



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

Variational Auto Encoders: Generating Data

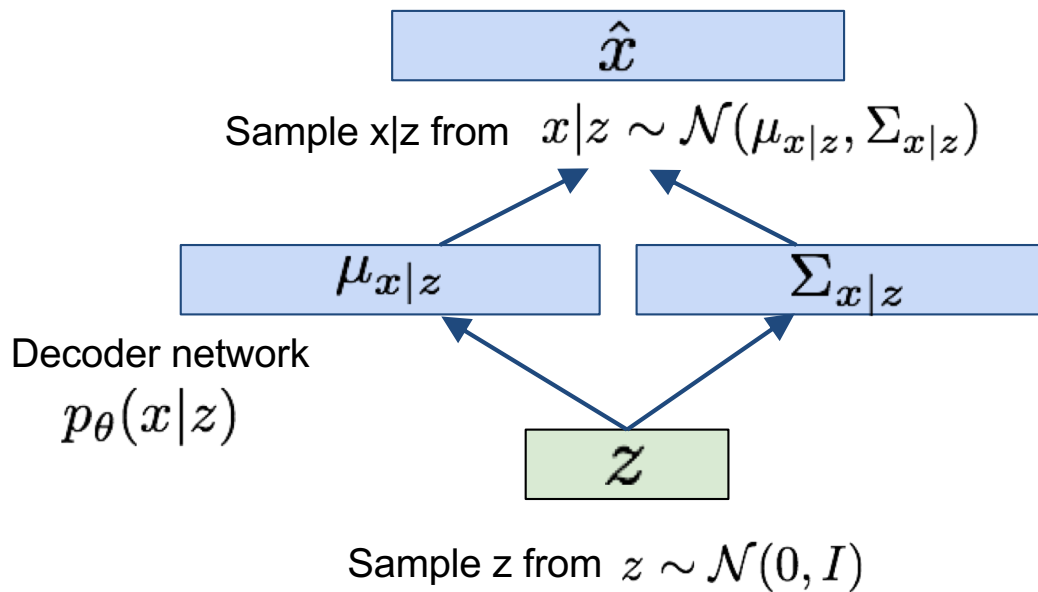
Use decoder network. Now sample z from prior!



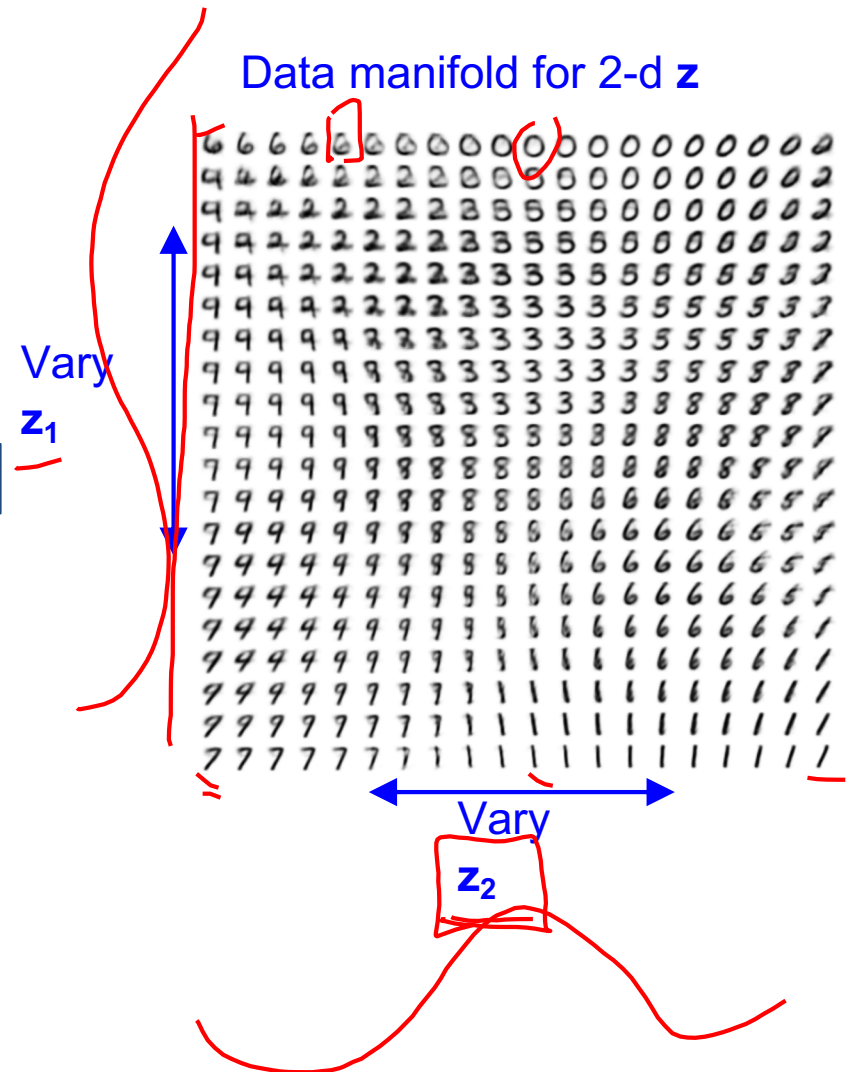
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Auto Encoders: Generating Data

Use decoder network. Now sample z from prior!



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

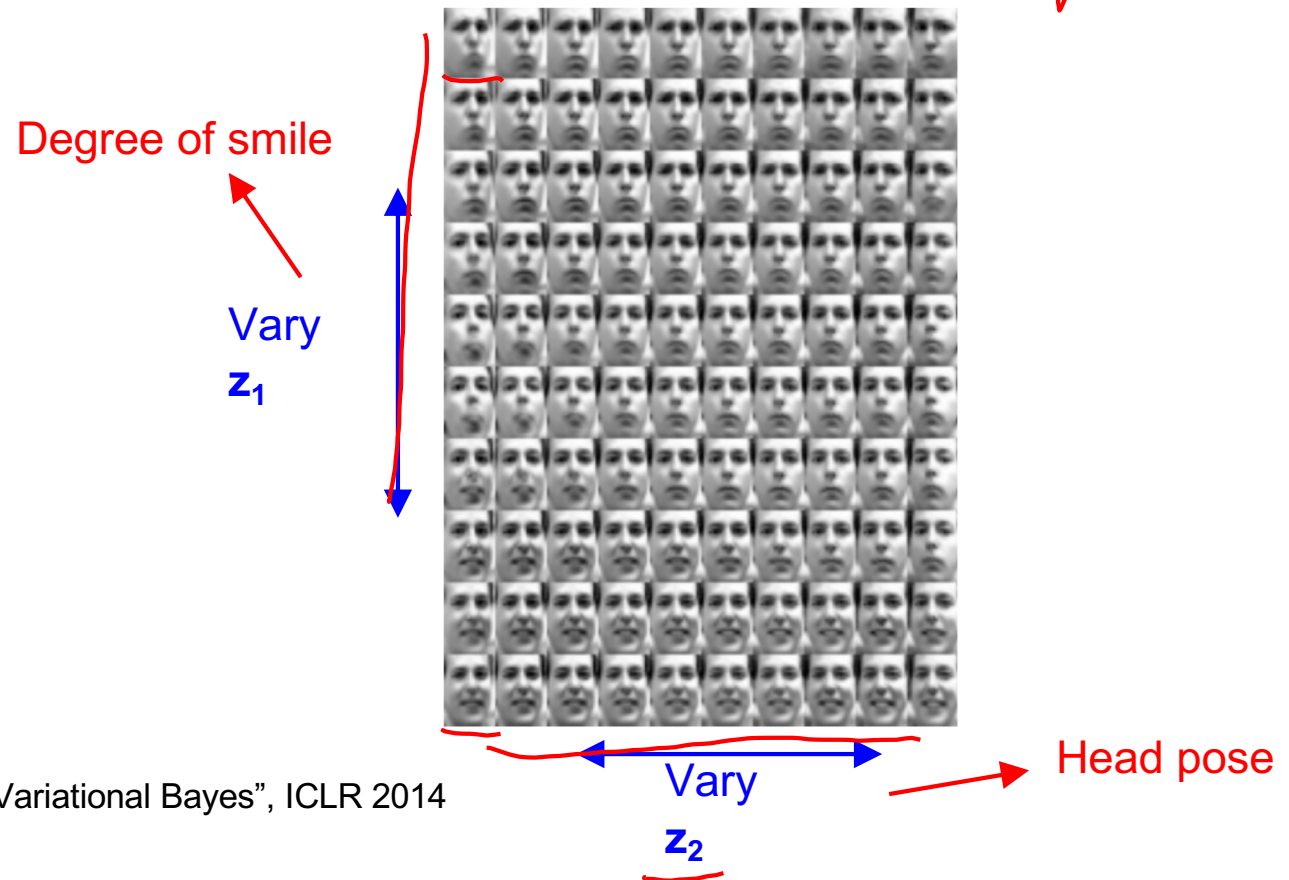


Variational Auto Encoders: Generating Data

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

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

$$p(\mathbf{z}|\theta)$$



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