

CS 4803 / 7643: Deep Learning

Topics:

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

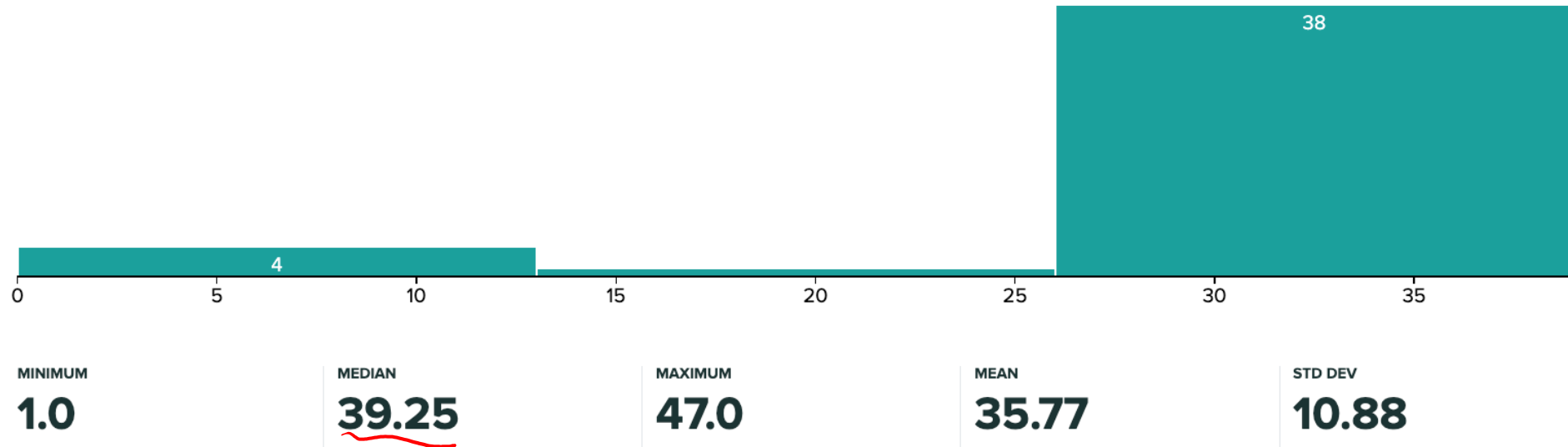
Dhruv Batra
Georgia Tech

Administrativa

- Project submission instructions released
 - Due: 11/24, 11:59pm
 - Last deliverable in the class
 - [– ~~Can't use late days~~ 8 free late days
 - https://www.cc.gatech.edu/classes/AY2021/cs7643_fall/

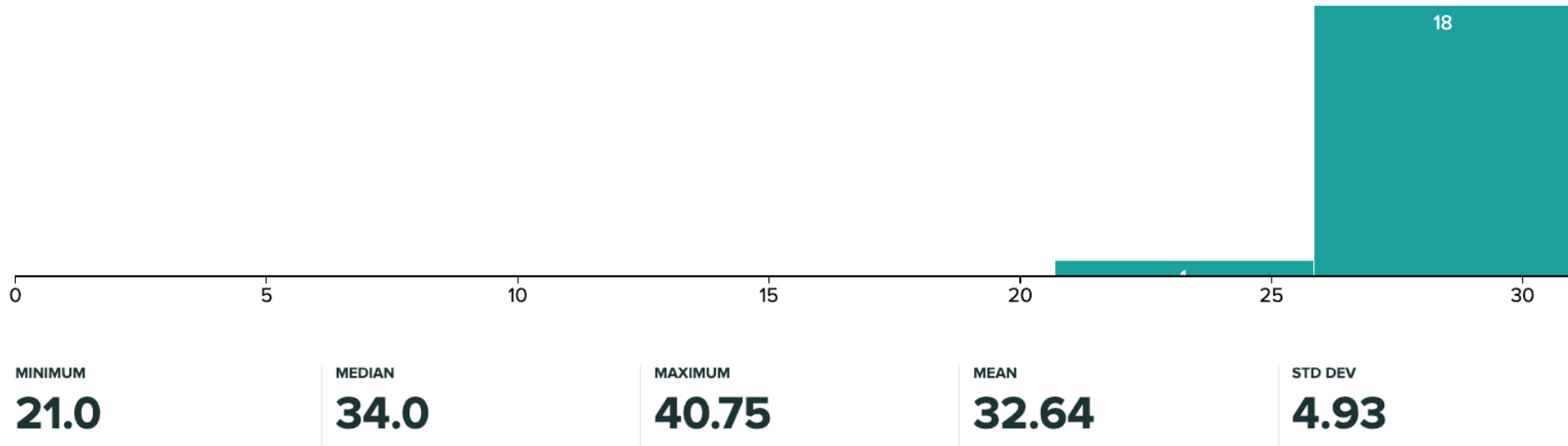
Administrativa

- HW5 Grades Released
 - Regrade requests close: 11/24, 11:59pm
- Grade histogram: 7643
 - Max possible: 39 (regular credit) + 11 (extra credit)



Administrativa

- HW5 Grades Released
 - Regrade requests close: 11/24, 11:59pm
- Grade histogram: 4803
 - Max possible: 31 (regular) + 19 (extra credit)



Recap from last time

Variational Autoencoders (VAE)

So far...

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

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

$$D = \{ \vec{x}_i \}_{i=1}^N$$

$N = \# \text{samples}$

"complex" $\rightarrow P(\vec{x})$

R.V

"Latent Variables"

Unobserved RV

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

$$= \underbrace{P(\vec{x} | \vec{z})}_{\text{conditional}} \underbrace{P(\vec{z})}_{\text{Prior}}$$

"simpler"

So far...

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

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

$$D = \{ \vec{x}_i \} \leftarrow$$

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

VAEs define intractable density function with latent \mathbf{z} :

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

if z is continuous

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

if z is discrete

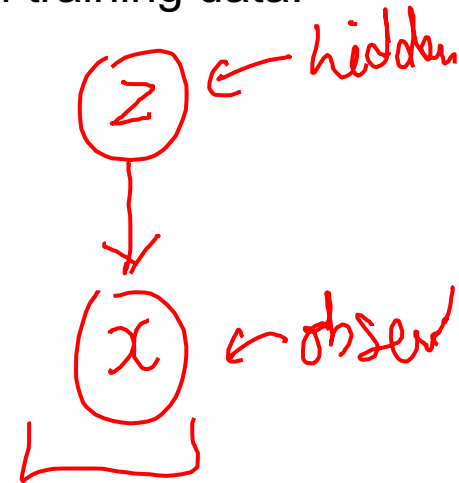
So far...

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 \mathbf{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

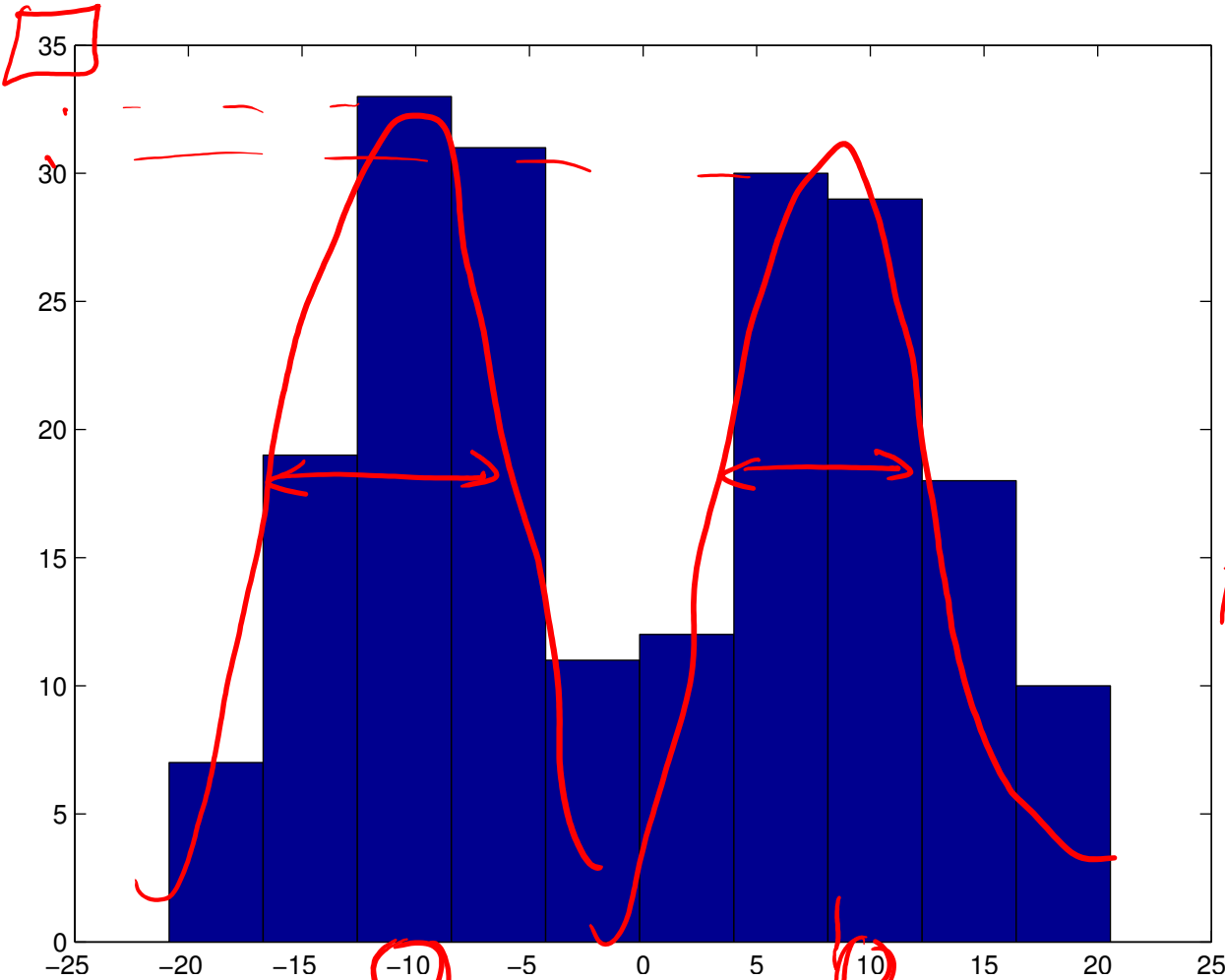
$x \in \mathbb{R}^1$

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

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

$$P(z) = \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$$

$\sim P(x)$



$p(x|z) \sim \mathcal{N}(\mu_1, \sigma_1^2)$ $\mu_1 = -10$

$\mathcal{N}(\mu_2, \sigma_2^2)$ $\mu_2 = +10$

Gaussian Mixture Model $z \in \{1, \dots, k\}$

$P(x, z)$

② Latent

↓

① Observed

$z \sim \text{Cat}(\pi)$

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

$\pi_c = P(z=c)$

$$\begin{aligned} \underline{P(x|z=c)} &= \mathcal{N}(\mu_c, \sigma_c^2) \\ &= \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(x-\mu_c)^2}{2\sigma_c^2}} \end{aligned}$$

$$\underline{P(x, z)} = \underline{P(x|z)} P(z)$$

Gaussian Mixture Model

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

$$P(x|z) = N(\quad)$$

Available
from model

$$\underline{P(\bar{x})} = \sum_z P(x, z)$$

$$= \sum_z \overbrace{P(x|z)} \overbrace{P(z)}$$

≡ Marginalization

$$\underline{P(\underline{z}|x)}$$

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

$$= \frac{P(x|z) p(z)}{\sum_z (\downarrow) (\downarrow)}$$

≡ (Inference)

Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

2. Variational Approximation

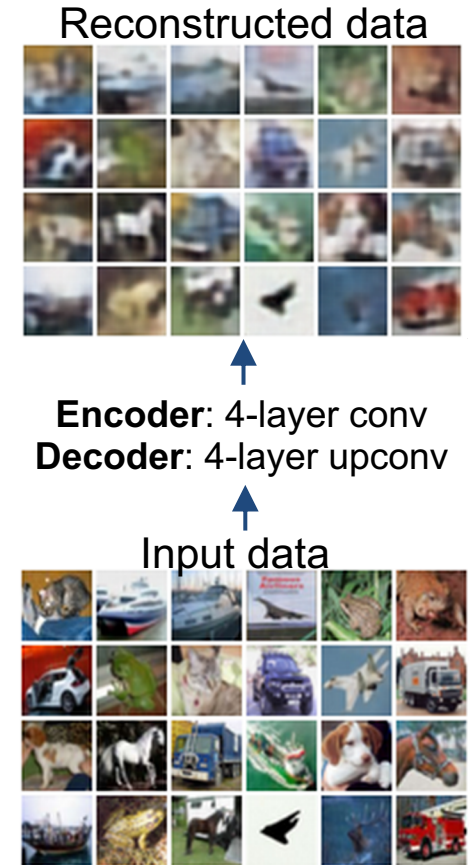
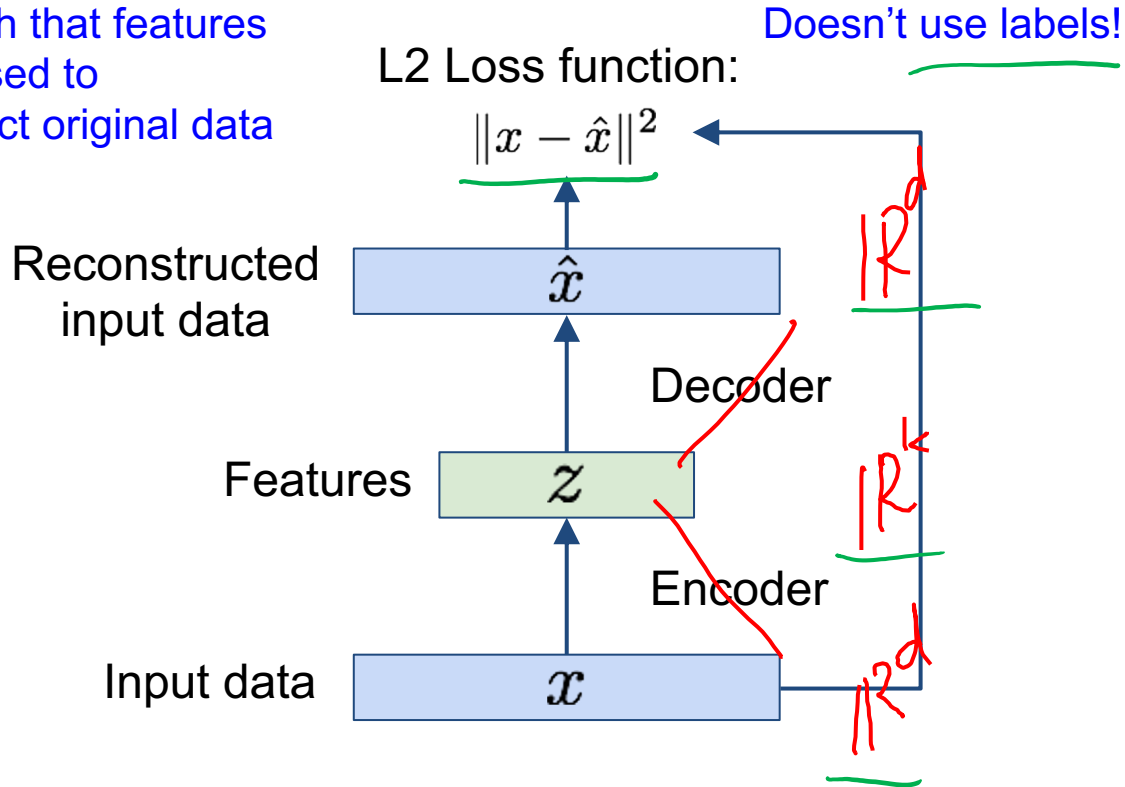
- Variational Lower Bound / ELBO

3. Amortized Inference Neural Networks

4. “Reparameterization” Trick

Autoencoders

Train such that features can be used to reconstruct original data



Autoencoders

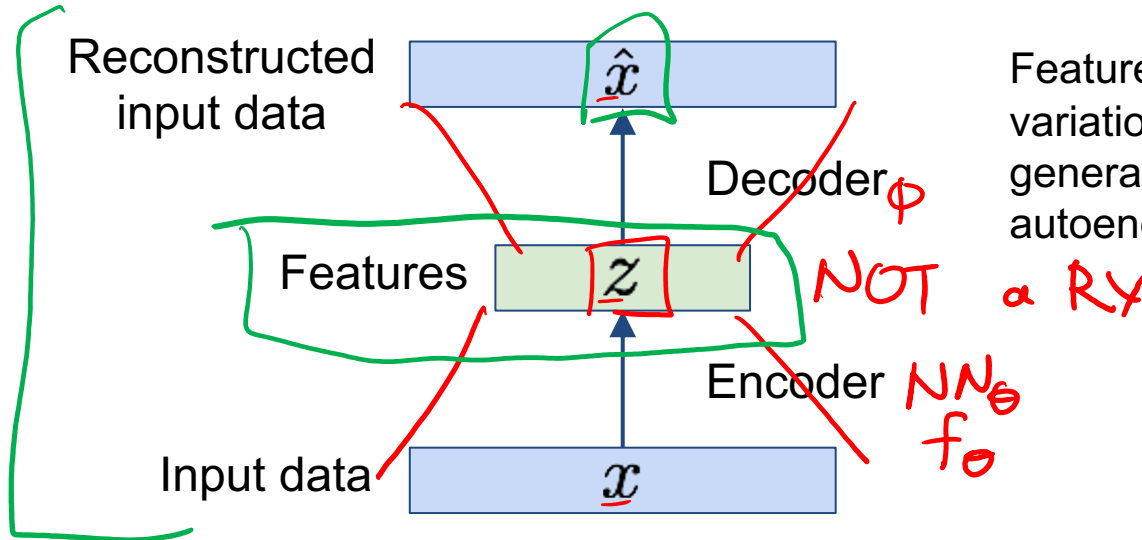
$$z = f_{\theta}(x)$$
$$\hat{x} = g_{\phi}(z)$$

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

VAE

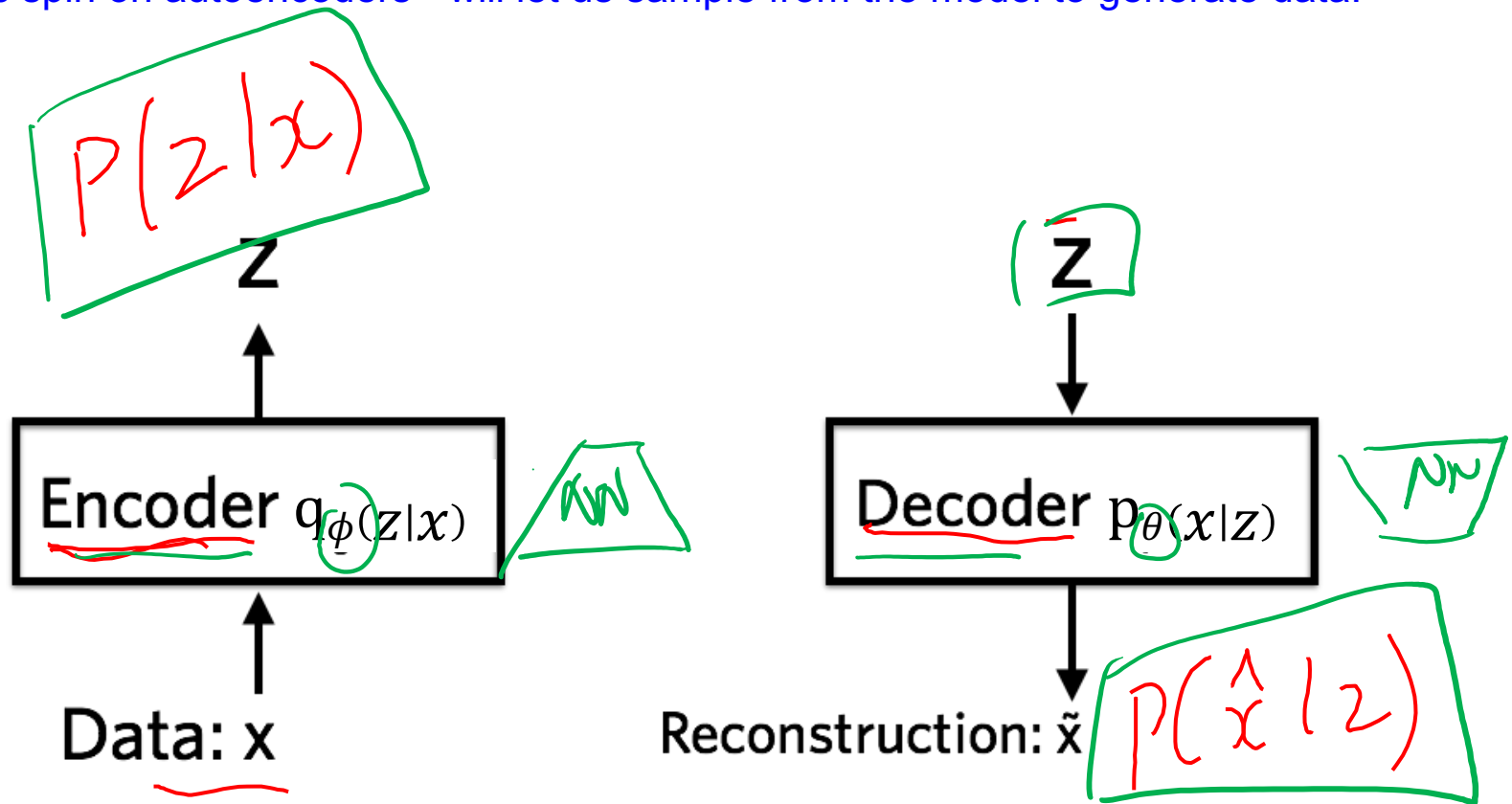
Autoencoders can reconstruct data, and can learn features to initialize a supervised model

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!



Plan for Today

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

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.

Handwritten diagram illustrating the concept of Variational Inference. It shows a complex distribution $p(z)$ being approximated by a simple distribution $q(z)$. The diagram includes the expression $E_{p(z)}[f(z)]$ and an arrow pointing to $E_{q(z)}[f(z)]$.

Key problem

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

"complex" $P(z)$



"simple" $q(z)$



"simple" distributions $\{q_\phi(z)\}$

$* P(z)$

$\{q_{\phi^*}(z)\}$

$$\min_{q \in \mathcal{Q}} \underbrace{d(p, q)}_{\text{KL}}$$

$\text{KL}(p||q)$
"right"
"hard"

$\text{KL}(q||p)$
"wrong"
"easy"

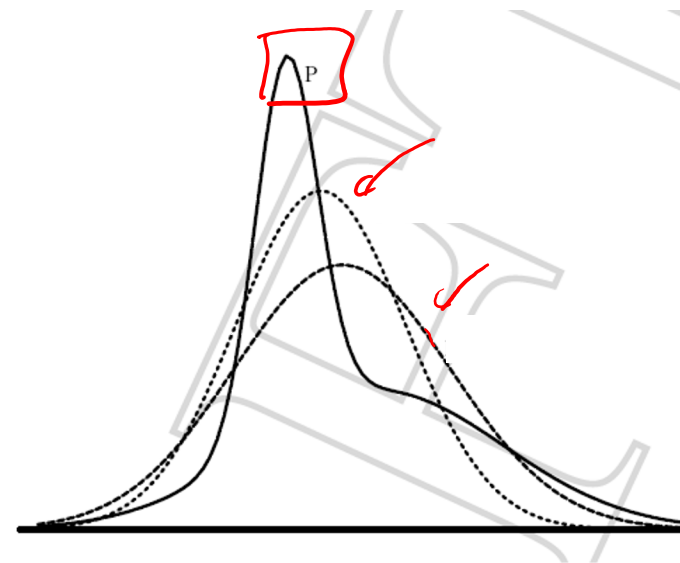
Intuition

$$\underline{KL(P \parallel Q)} = \sum_z P(z) \log \frac{P(z)}{Q(z)}$$

$$KL(Q \parallel P) = \sum_z \underbrace{Q(z)}_{\text{"Support"}} \log \underbrace{\frac{Q(z)}{P(z)}}_{\text{"error"}}$$

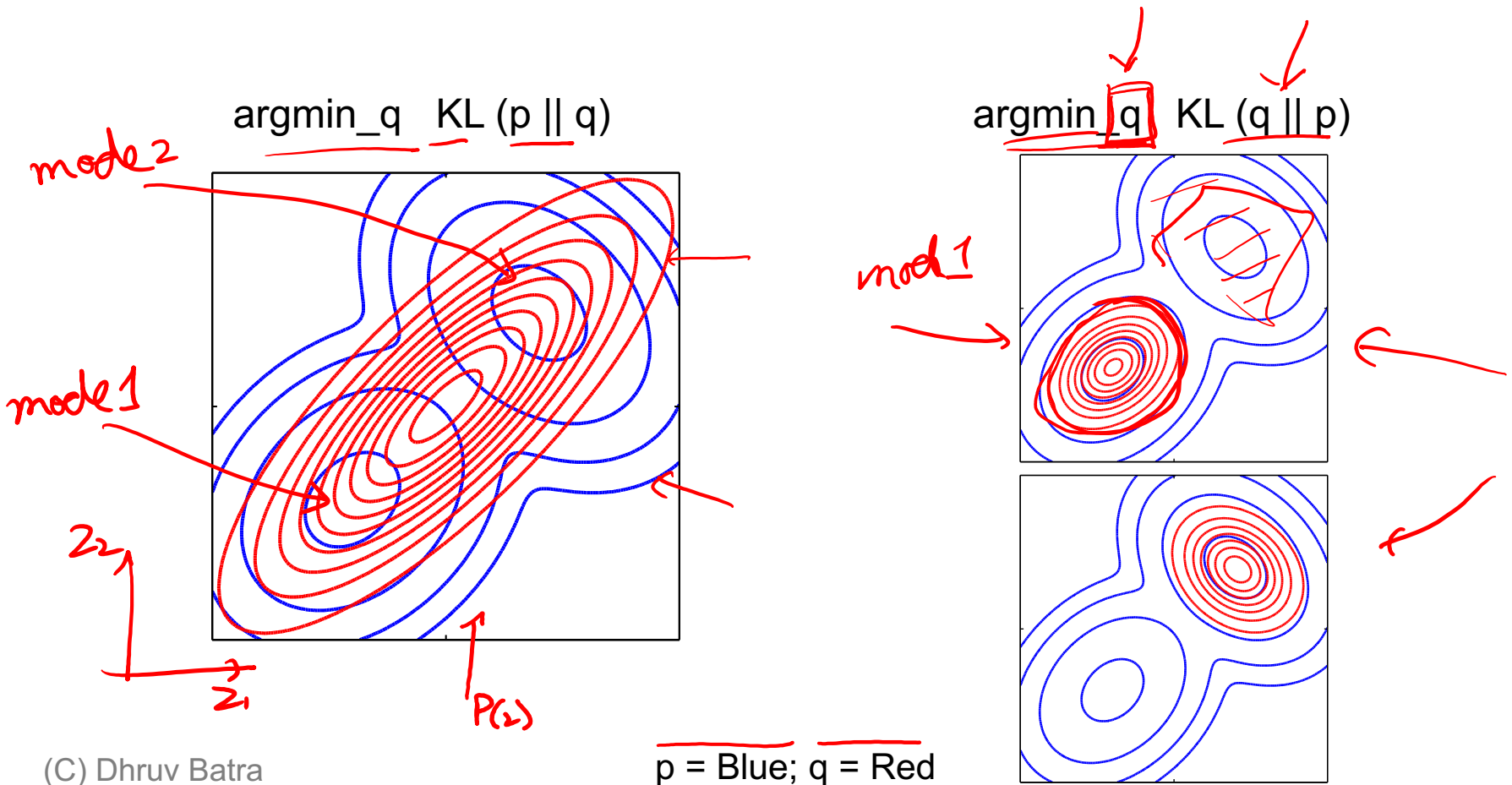
Find simple approximate distribution

- 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 2

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



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:

log. likelihood

$$\underline{ll(\theta : \mathcal{D})} = \underline{\log} \prod_{i=1}^N \underline{P(\mathbf{x}_i | \theta)}$$

$$\mathcal{D} = \{ \underline{\tilde{\mathbf{x}}_i} \}_{i=1}^N$$

$$\underline{P(\tilde{\mathbf{x}}, \mathbf{z})}$$

$$= \sum_{i=1}^N \underline{\log} \underline{P(\mathbf{x}_i | \theta)}$$

$$= \sum_{i=1}^N \log \sum_{\mathbf{z}} \underline{P(\mathbf{x}_i, \mathbf{z} | \theta)}$$

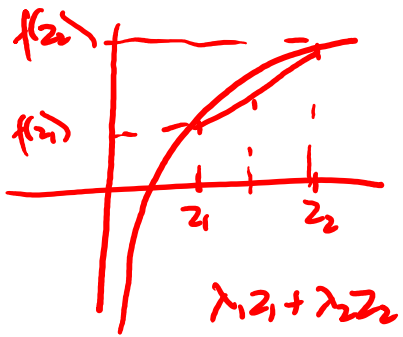
$$\log \sum_{\mathbf{z}} \underline{P(\tilde{\mathbf{x}}_i | \theta) P(\mathbf{z} | \tilde{\mathbf{x}}_i, \theta)}$$

$$\log \frac{E_{P(\mathbf{z} | \tilde{\mathbf{x}}_i, \theta)} [\underline{P(\tilde{\mathbf{x}}_i | \theta)}]}{f(\cdot)}$$

“complex” \rightarrow simple

Applying Jensen's inequality

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



$$f(\lambda_1 z_1 + \lambda_2 z_2) \geq \lambda_1 f(z_1) + \lambda_2 f(z_2) \quad \left\{ \begin{array}{l} \text{Convexity} \\ \text{of } f \end{array} \right.$$

$$f\left(\sum_{i=1}^k \lambda_i z_i\right) \geq \sum_{i=1}^k \lambda_i f(z_i) \quad \boxed{\lambda_1, \dots, \lambda_k}$$

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

$$\lambda_1 + \lambda_2 = 1$$

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

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

$z \rightarrow g(z)$

Applying Jensen's inequality

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

$$\mathcal{L}(\theta) \equiv \log P(\bar{x}_i | \theta) = \log \sum_{\mathbf{z}} \underbrace{P(x_i, \mathbf{z} | \theta)}_{Q_i(\mathbf{z})}$$

$$\geq \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(x_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

"Free Energy" $F(\theta, Q_i)$

$$\begin{array}{l} \max_{\theta} \mathcal{L}(\theta) \geq F(\theta, Q_i) \\ \hline \downarrow \\ \max_{\theta, Q_i} F(\theta, Q_i) \end{array}$$

Variational Lower Bound
Evidence-based LB (ELBO)

Evidence Lower Bound

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

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

(VAE)
objective

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

(GMM)
objective

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

ELBO: Factorization #1 (GMMs)

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

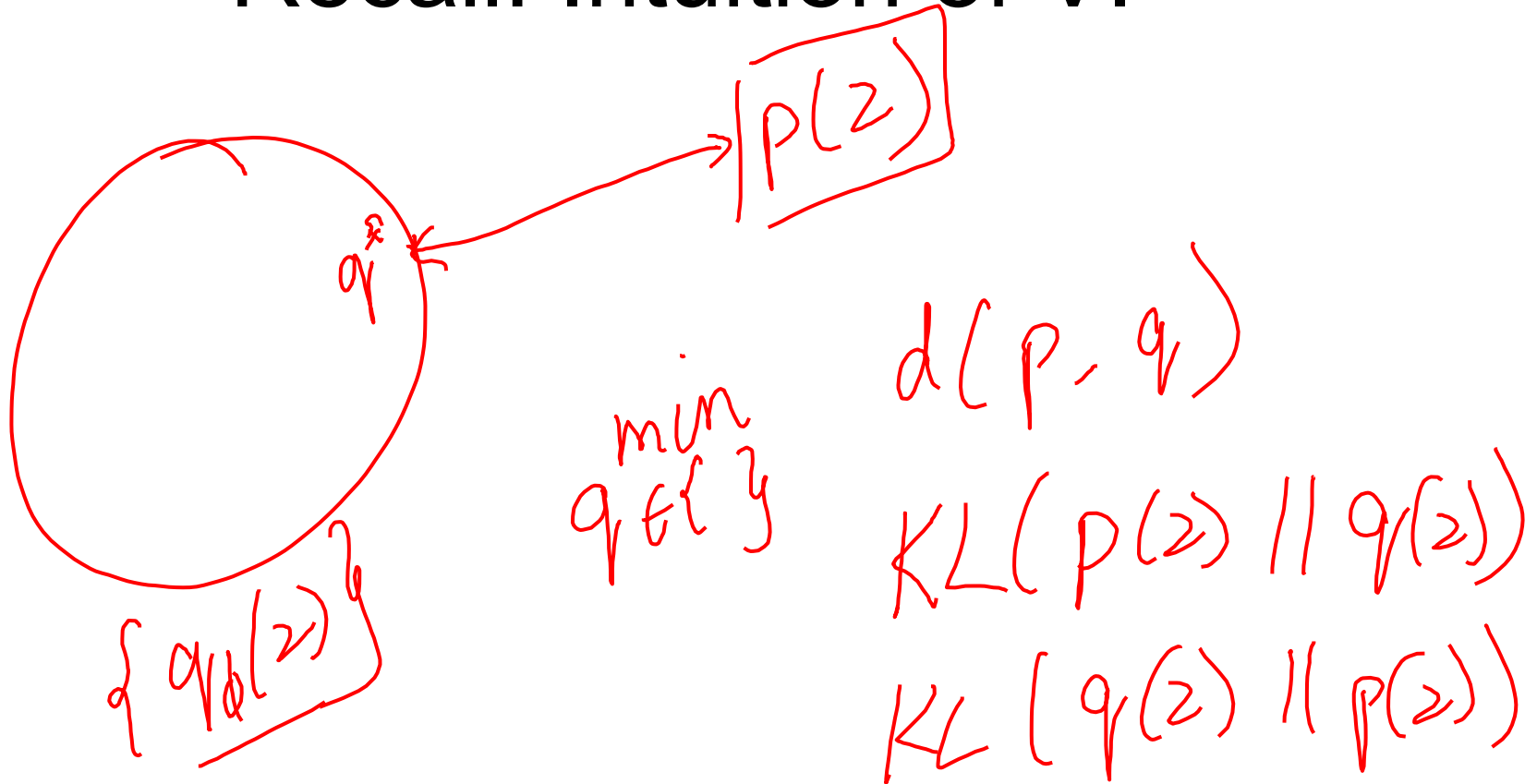
$$\underline{\underline{\ell(\theta : \mathcal{D})}} \geq \underline{\underline{F(\theta, Q_i)}} = \sum_{i=1}^N \sum_{\mathbf{z}} Q_i(\mathbf{z}) \log \frac{P(\mathbf{x}_i, \mathbf{z} | \theta)}{Q_i(\mathbf{z})}$$

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

$$F(\theta, Q_i) = \underbrace{\log P(\tilde{x}_i | \theta)}_{\ell(\theta)} - \text{KL}(Q_i(z) || P(z | \tilde{x}_i, \theta))$$

$$\max_{\theta, Q_i} \underline{\underline{F(\theta, Q_i)}} = \underline{\underline{\ell(\theta)}} - \text{KL}(\quad) \downarrow$$

Recall: Intuition of VI



ELBO: Factorization #1 (GMMs)

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

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

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

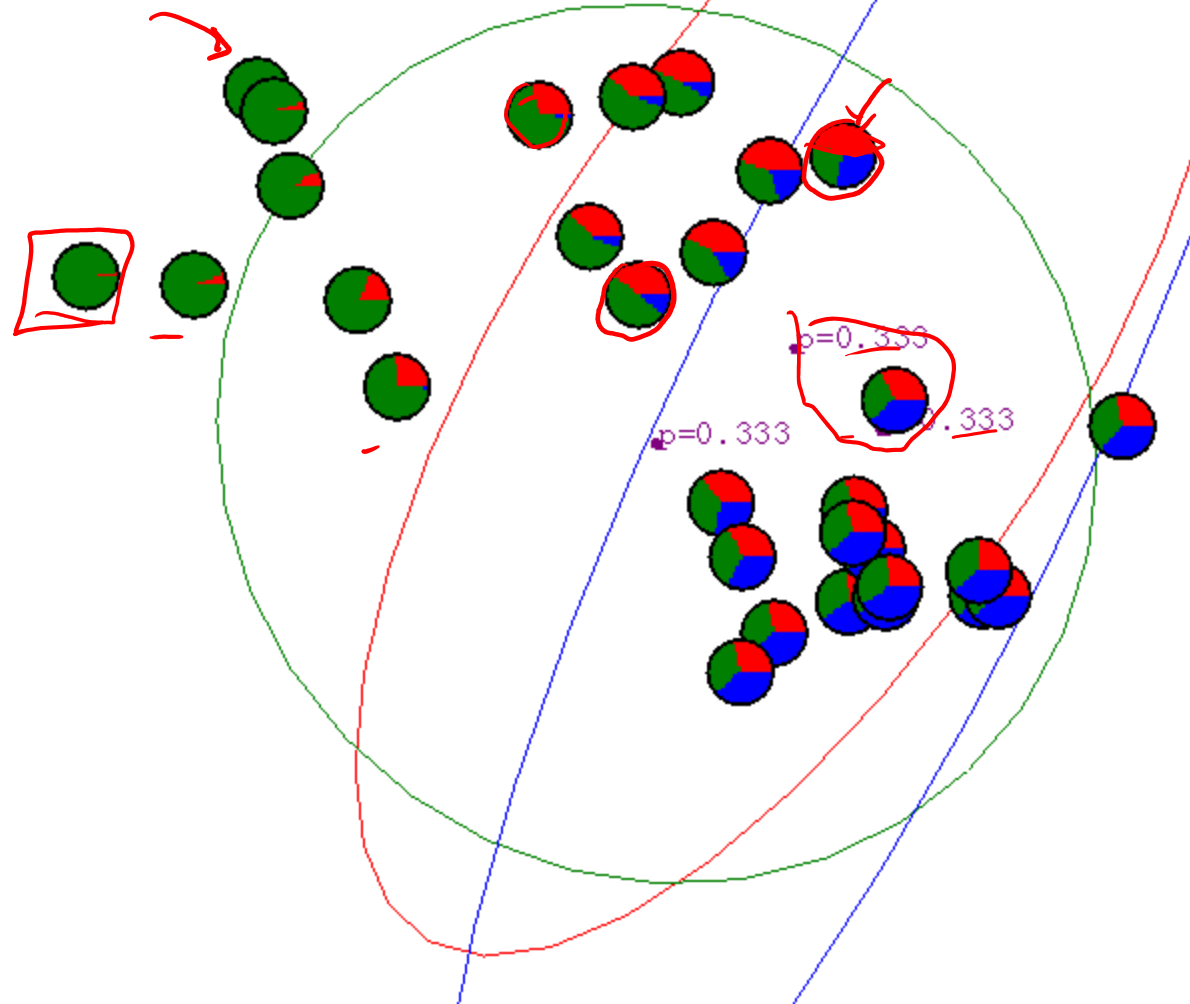
- **M-step:** Fix $Q_i^{(t)}$, maximize F over θ
 - maximize expected likelihood under $Q_i(\mathbf{z})$
 - Corresponds to weighted dataset:
 - $\langle \mathbf{x}_1, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1 \mid \mathbf{x}_1)$
 - $\langle \mathbf{x}_1, \mathbf{z}=2 \rangle$ with weight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

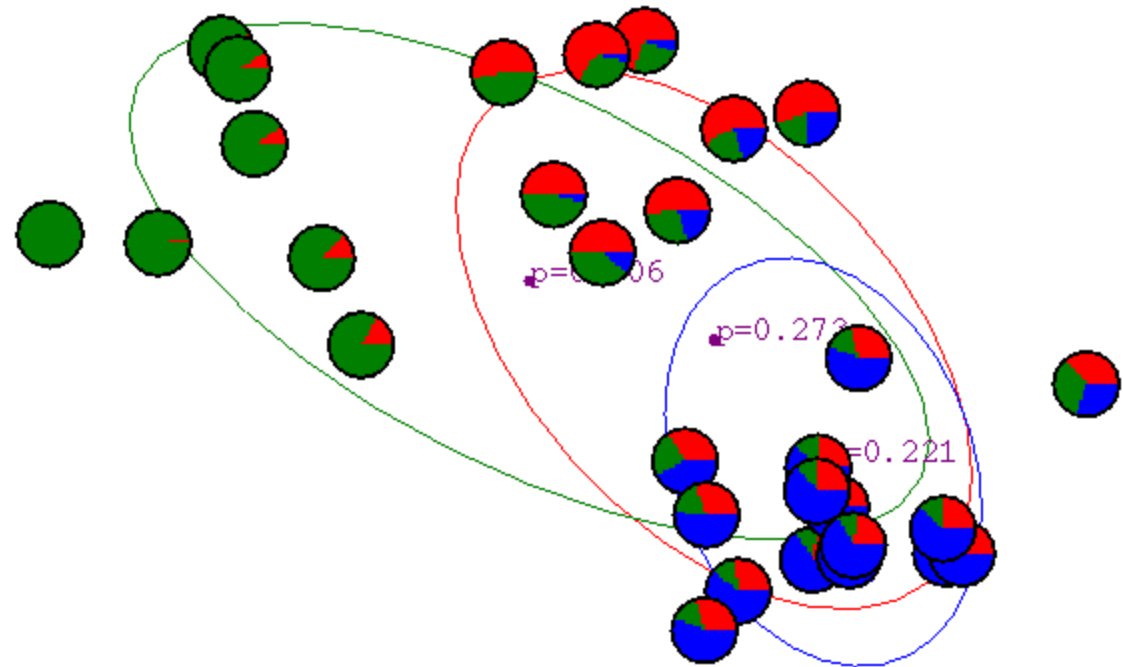
$$\Theta = \{\pi_1, \pi_2, \pi_3, \underline{\mu}_1, \underline{\mu}_2, \underline{\mu}_3, \underline{\Sigma}_1, \underline{\Sigma}_2, \underline{\Sigma}_3\}$$

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

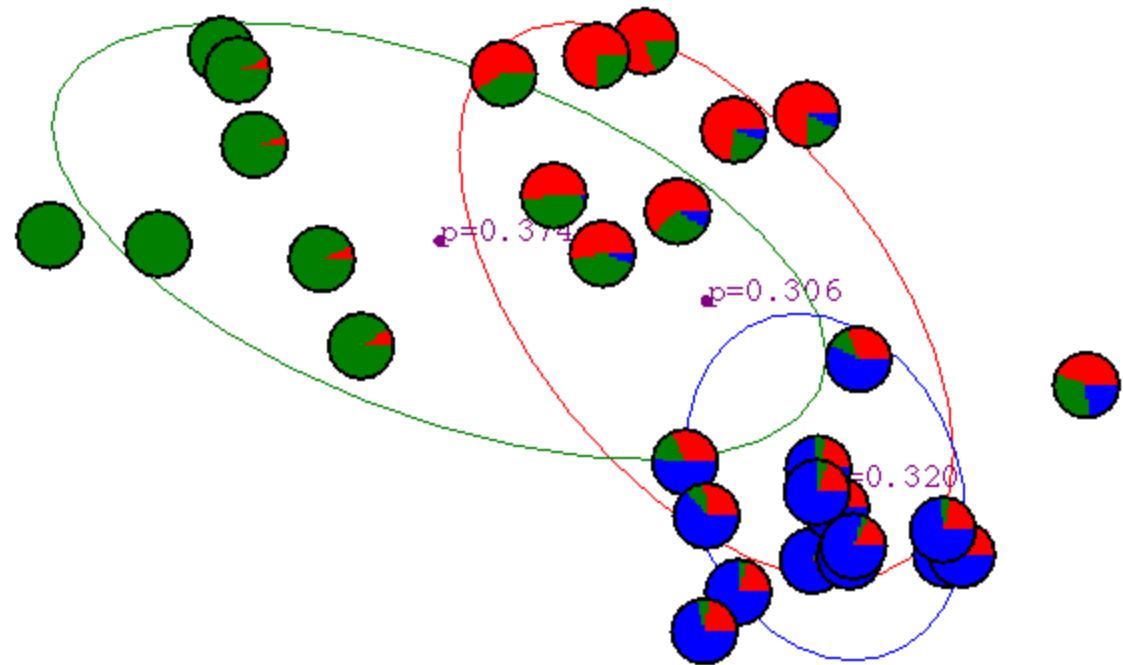
$i = 1:N$



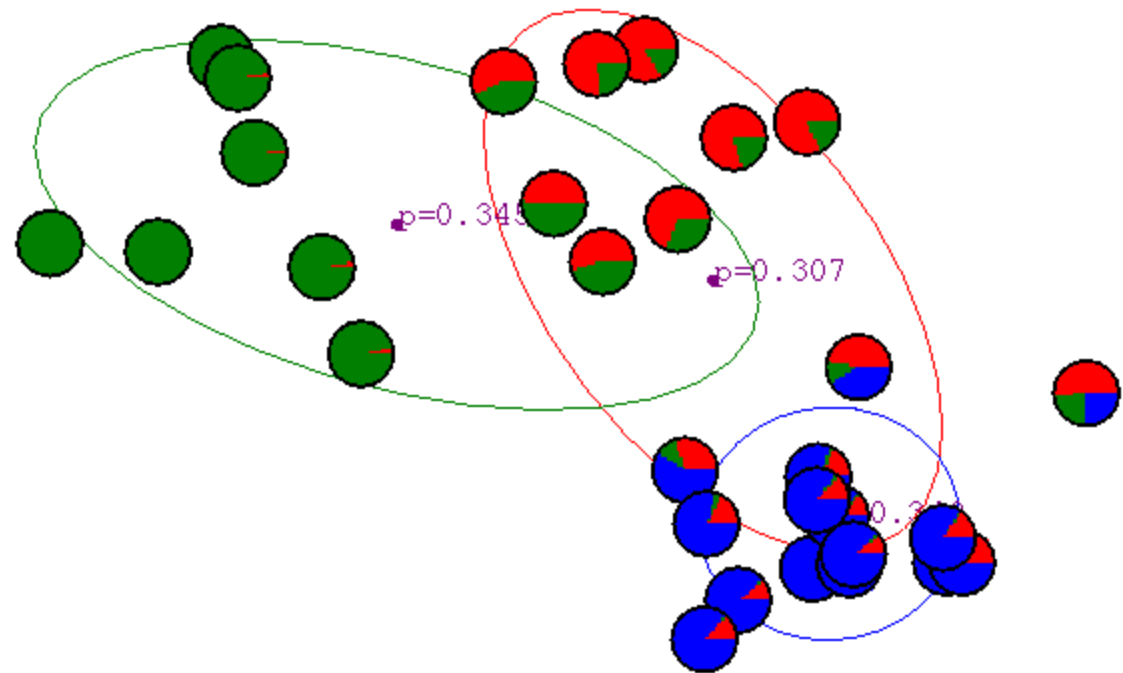
After 1st iteration



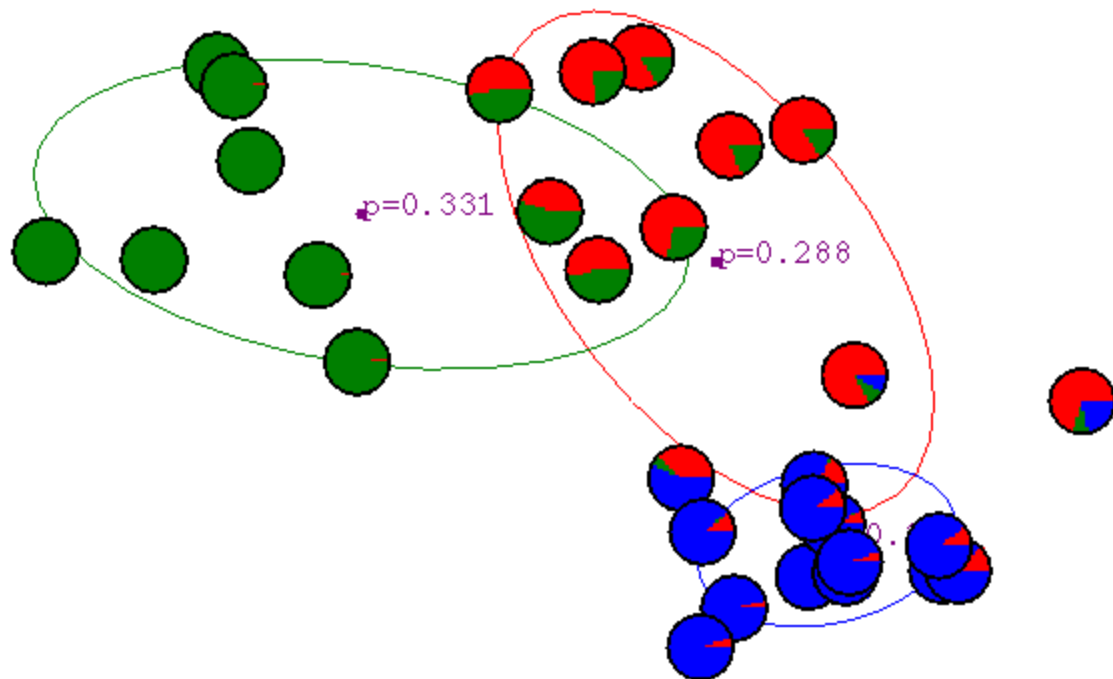
After 2nd iteration



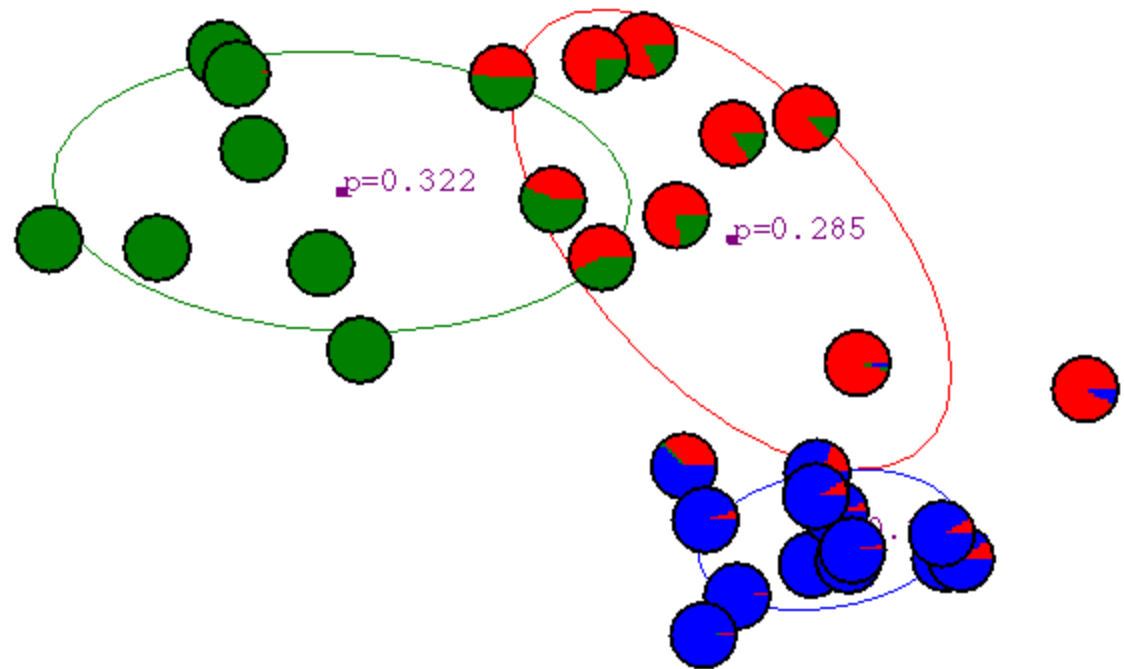
After 3rd iteration



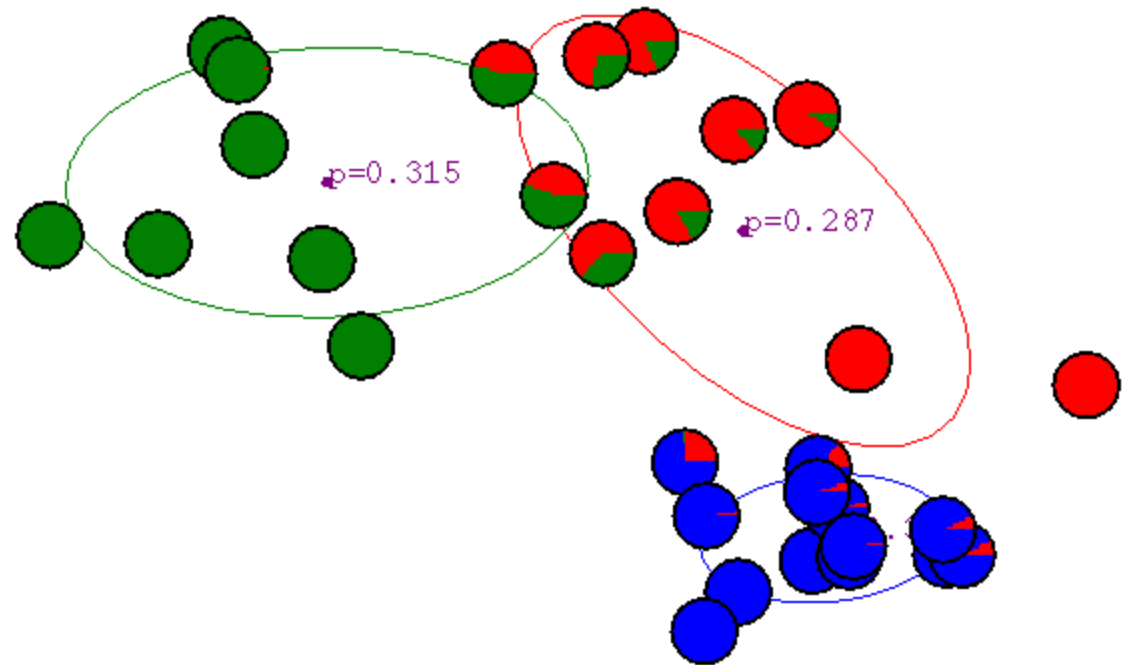
After 4th iteration



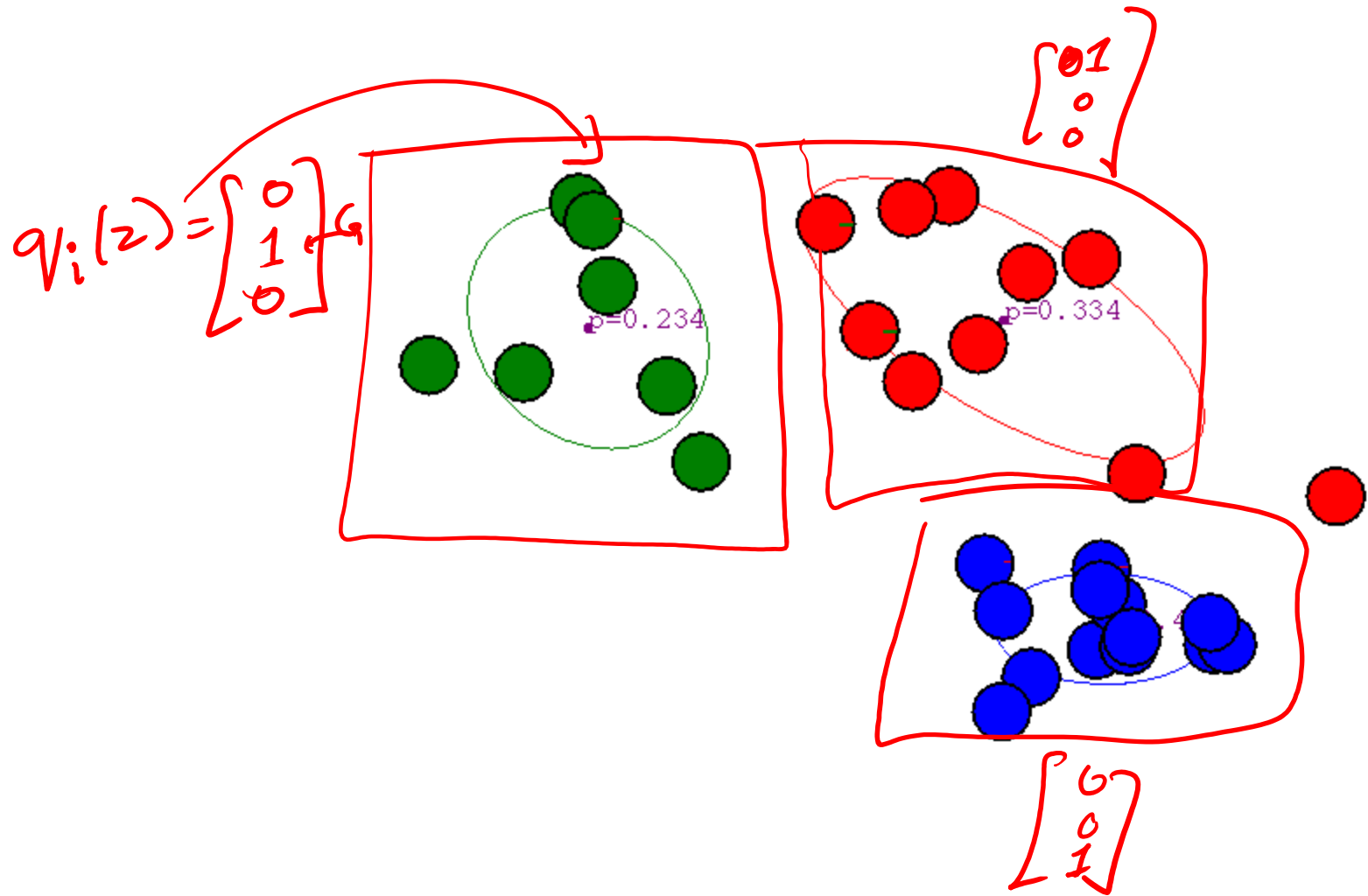
After 5th iteration



After 6th iteration



After 20th iteration



ELBO: Factorization #2 (VAEs)

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

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

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

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

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

(VAE Objective) maximize

"Explain the data"

"Regularize"
"Be simple"

Variational Auto Encoders

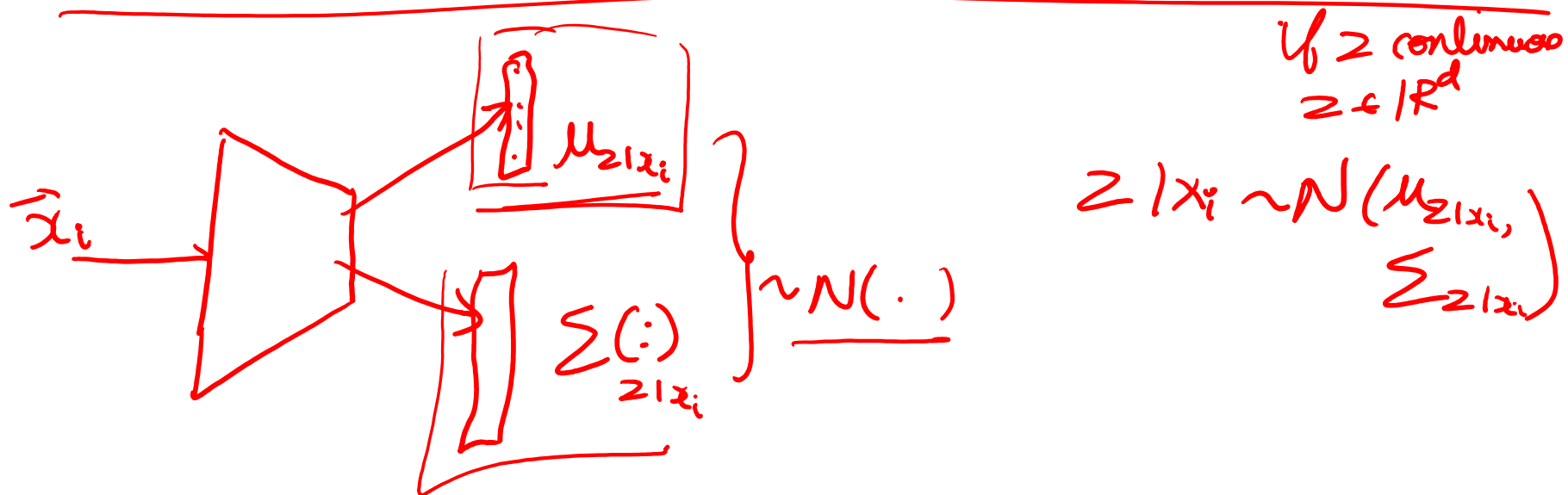
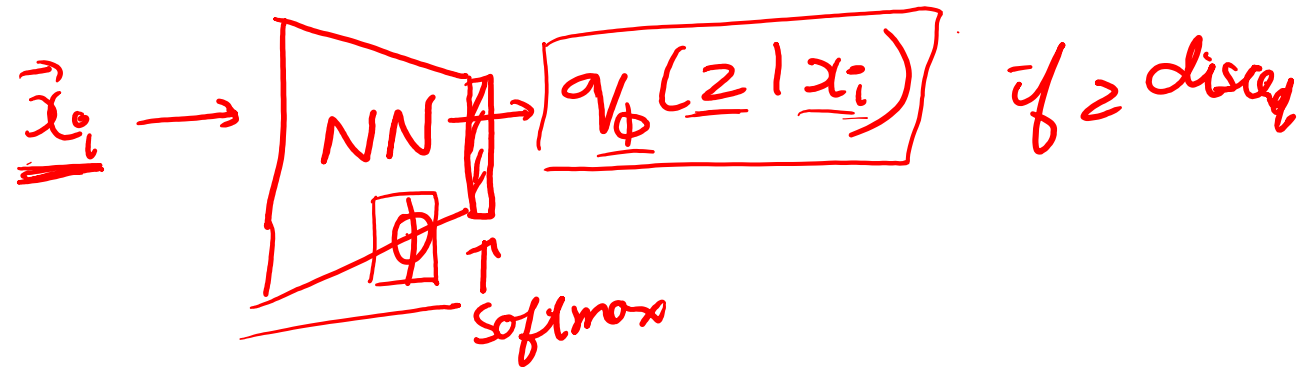
VAEs are a combination of the following ideas:

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

Amortized Inference Neural Networks

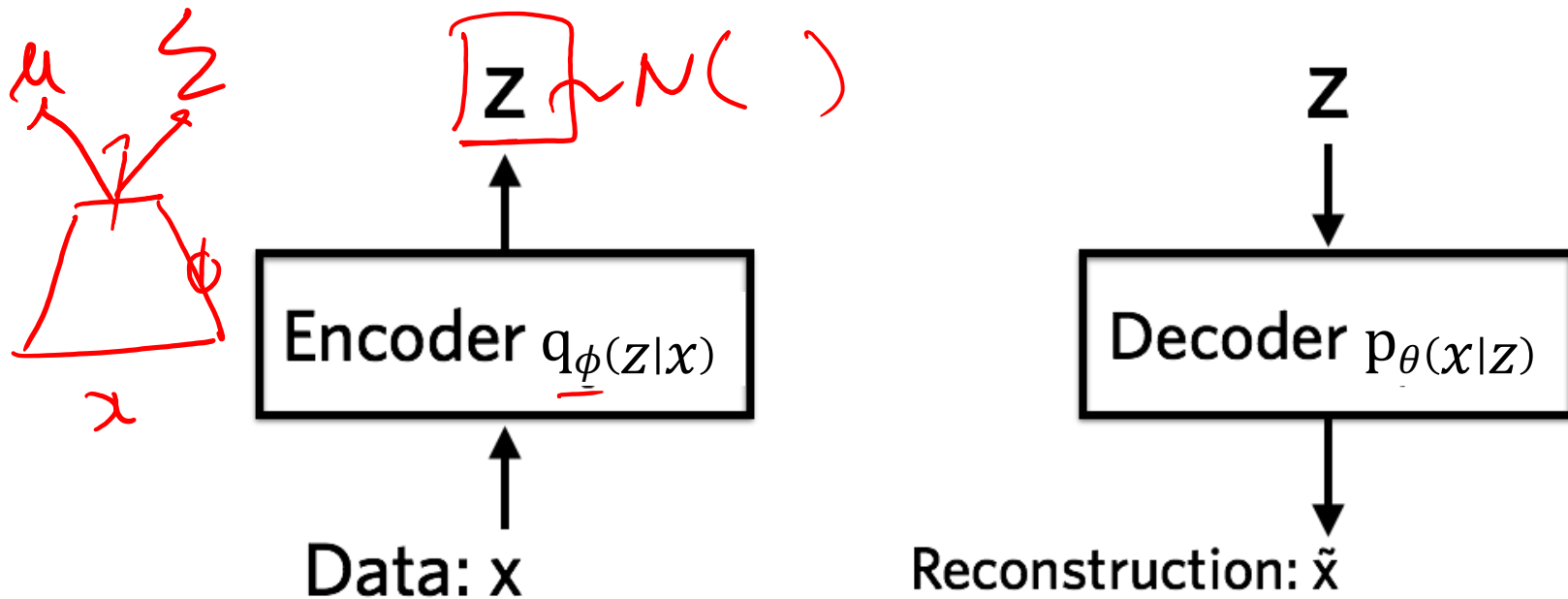
"Let's shove a NN into a Bayesian inf prob"

$$Q_{\phi}(z) = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.7 \end{bmatrix}$$

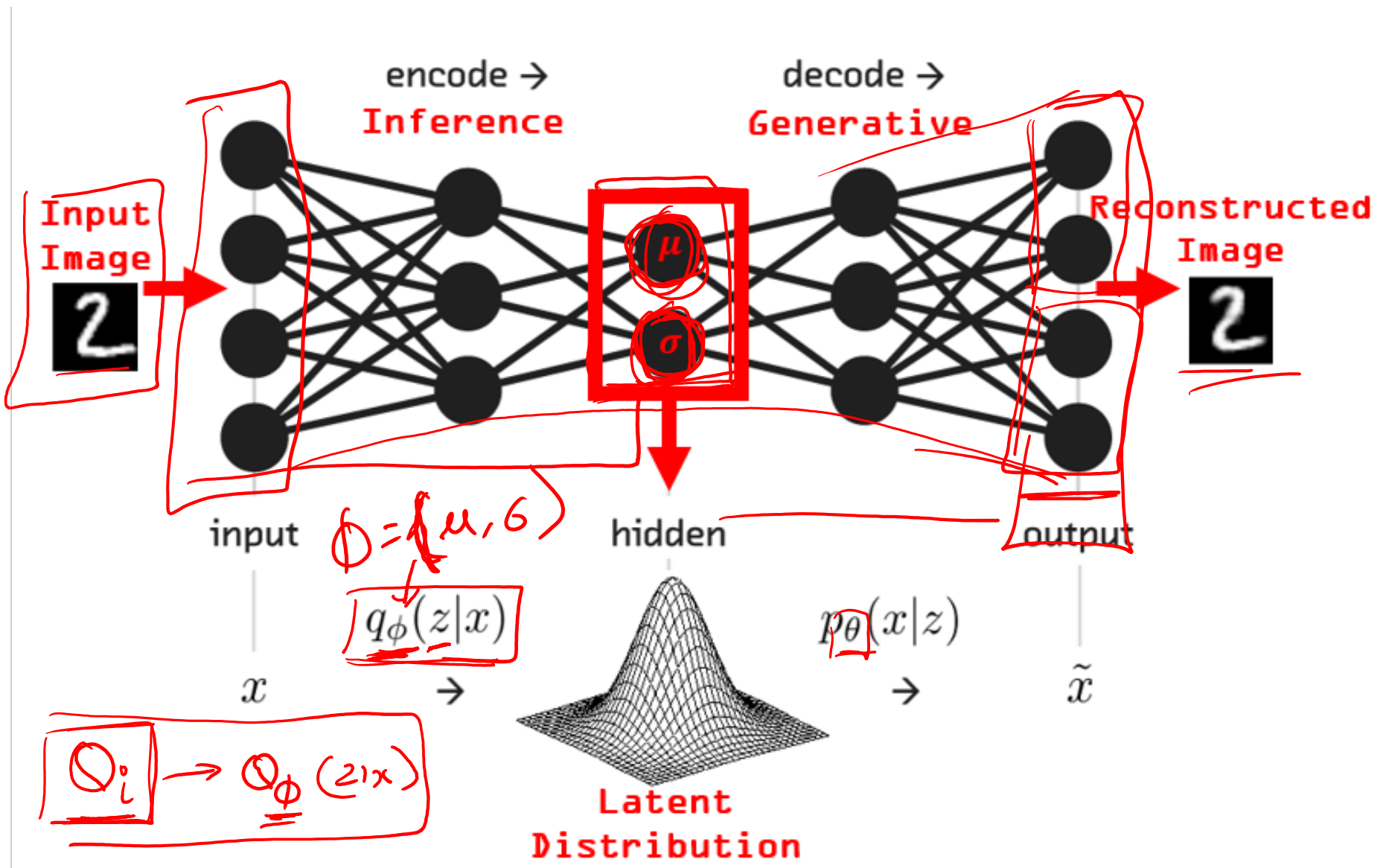


Variational Autoencoders

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

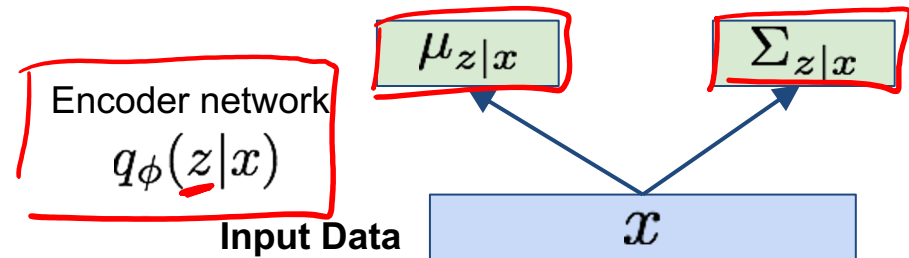
Input Data

\mathcal{X}

Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

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

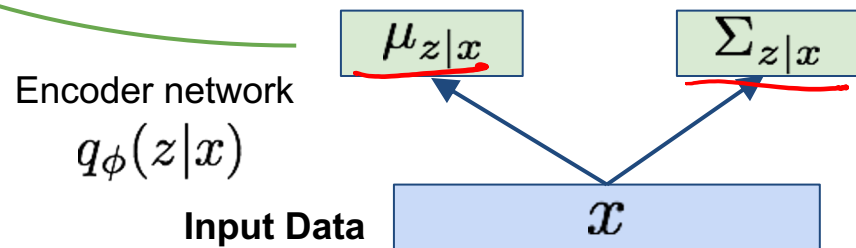


Variational Auto Encoders

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

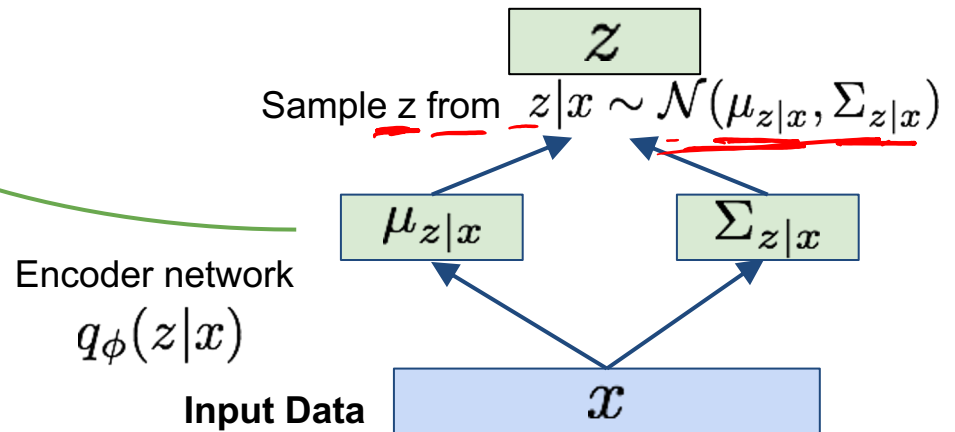


Variational Auto Encoders

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

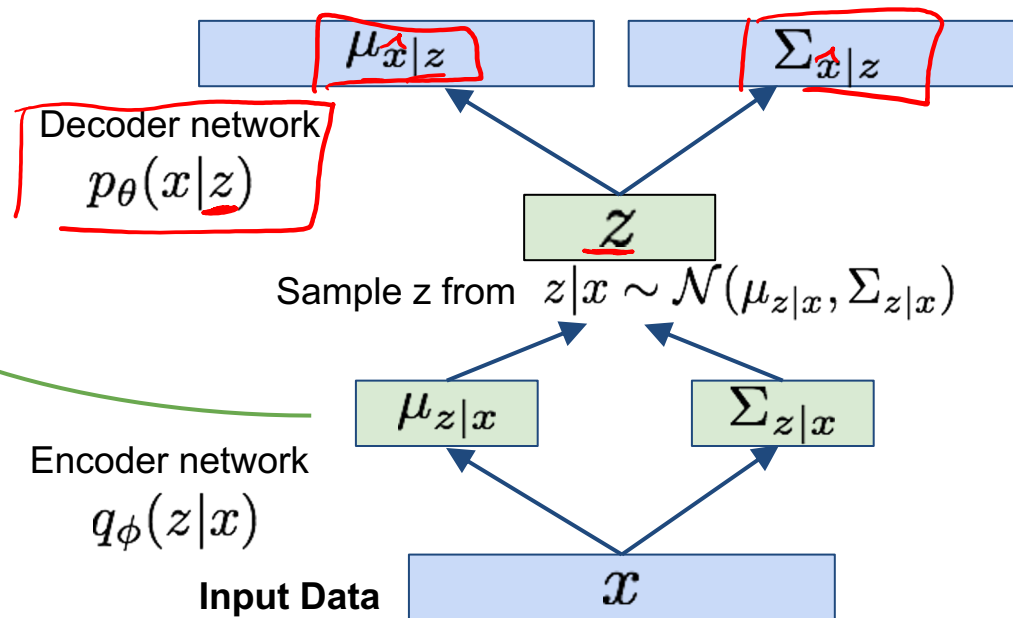


Variational Auto Encoders

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

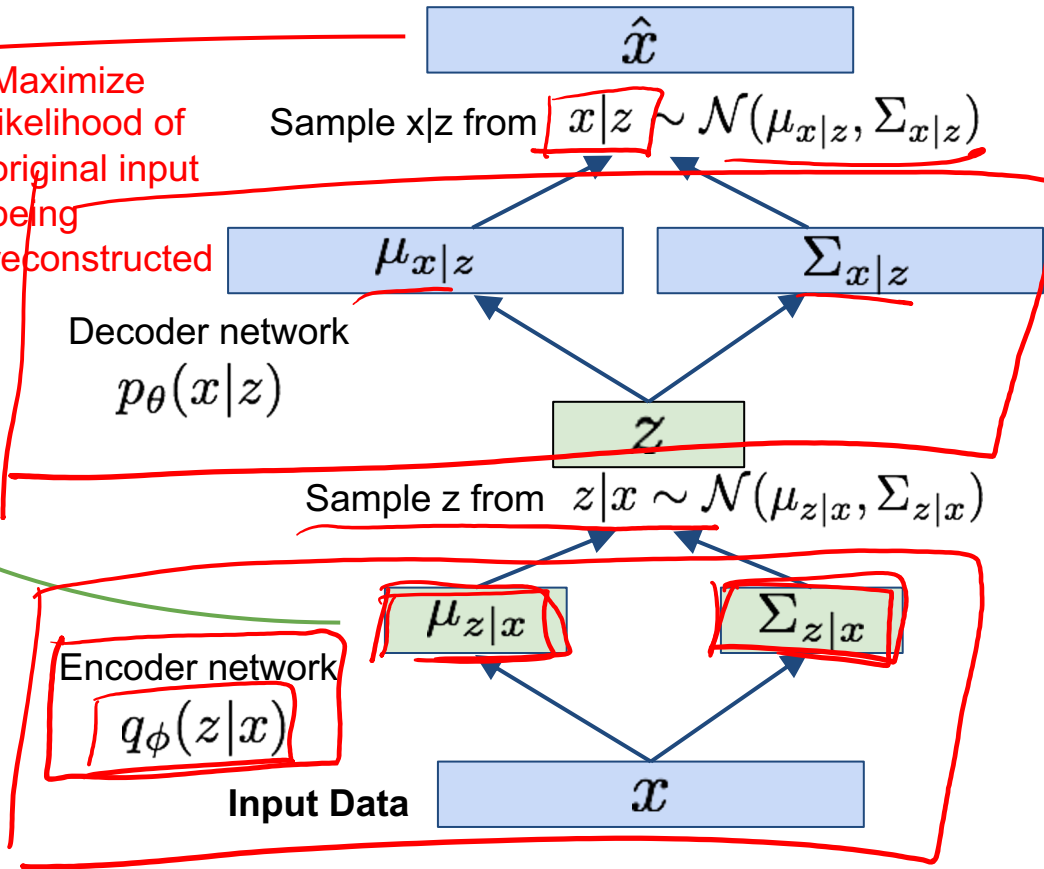
Putting it all together: maximizing the likelihood lower bound

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

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

Make approximate posterior distribution close to prior

Maximize likelihood of original input being reconstructed



Variational Auto Encoders

Putting it all together: maximizing the likelihood lower bound

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

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

$\frac{\partial \mathcal{L}}{\partial \theta}$
 $\frac{\partial \mathcal{L}}{\partial \phi}$

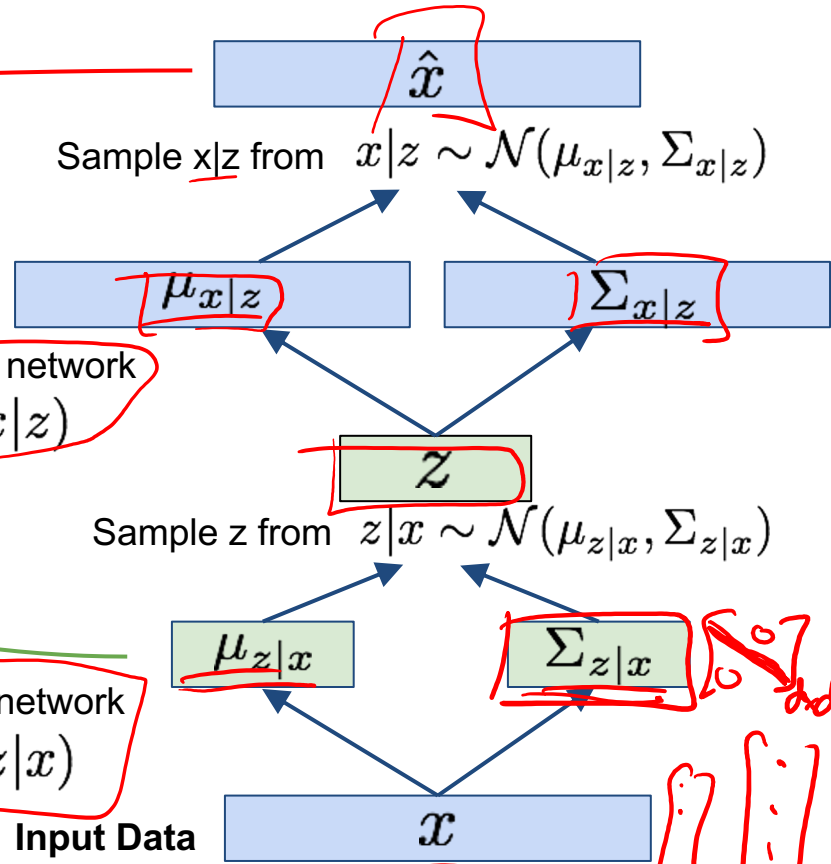
Make approximate posterior distribution close to prior

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

Maximize likelihood of original input being reconstructed

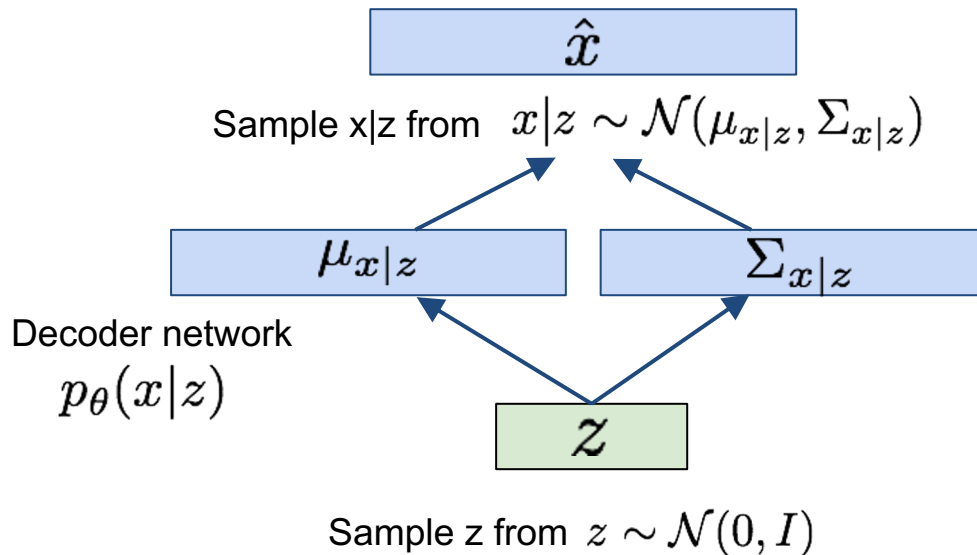
Decoder network
 $p_\theta(x|z)$

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



Variational Auto Encoders: Generating Data

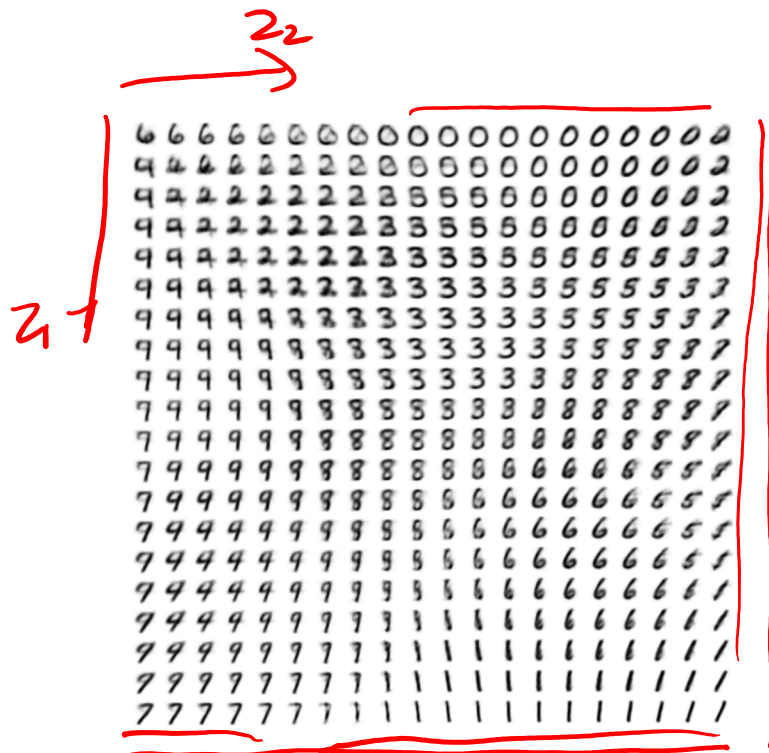
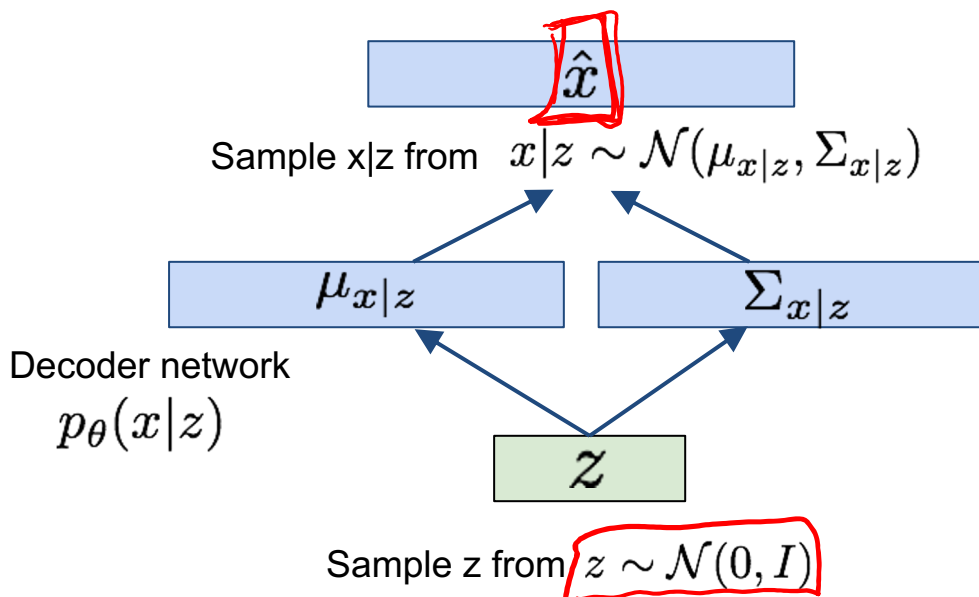
Use decoder network. Now sample z from prior!



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

Variational Auto Encoders: Generating Data

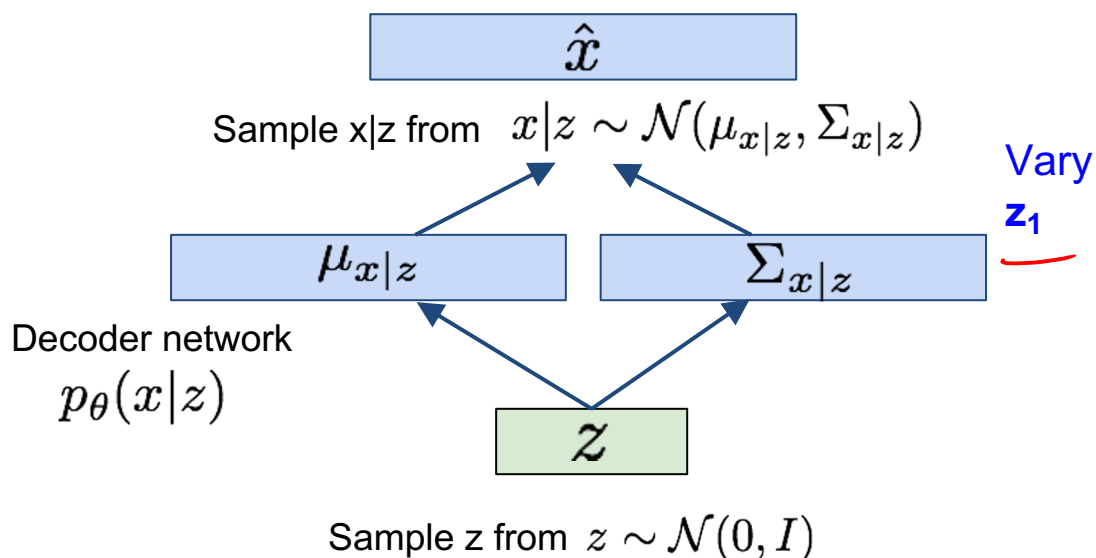
Use decoder network. Now sample z from prior!



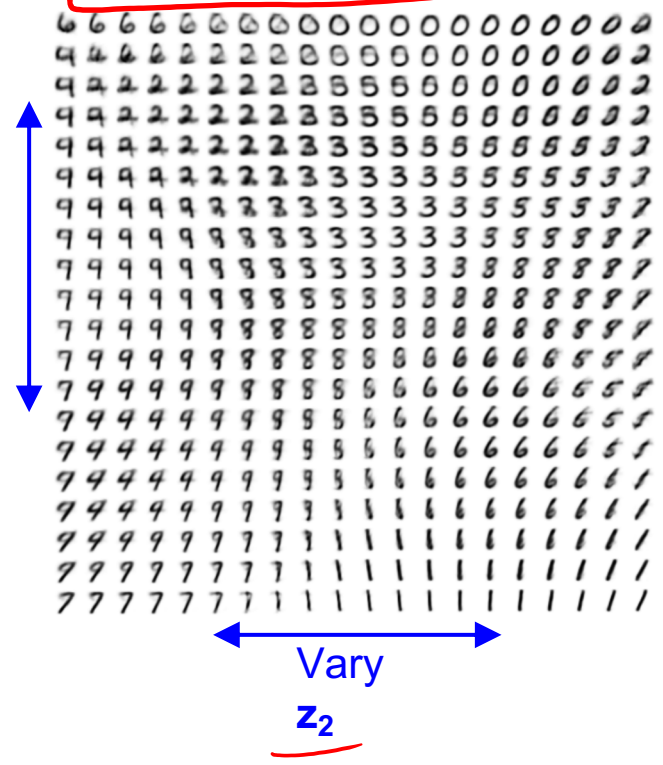
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Auto Encoders: Generating Data

Use decoder network. Now sample z from prior!



Data manifold for 2-d z

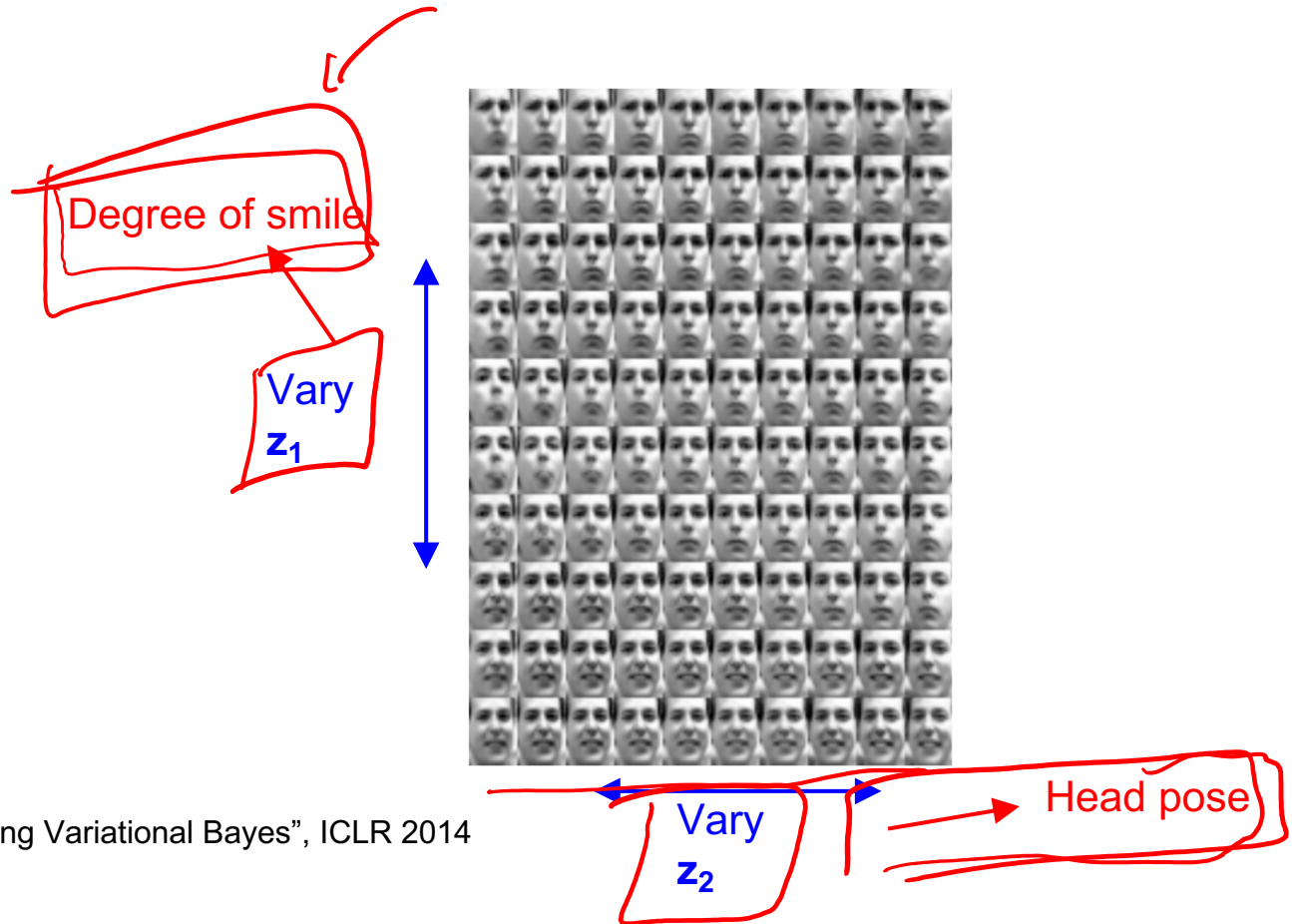


Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Auto Encoders: Generating Data

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

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



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 $q_\phi(\mathbf{z}|x)$!

Degree of smile



Vary
 z_1



Vary
 z_2

Head pose

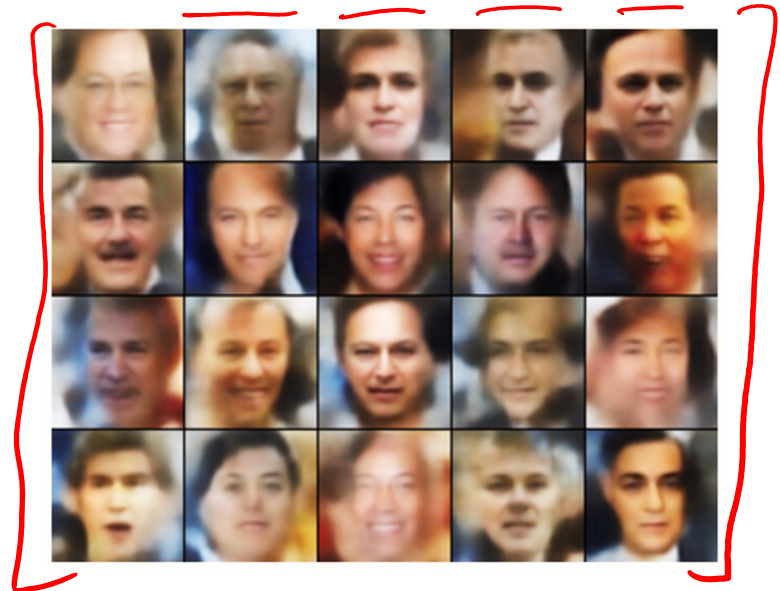


Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

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

$P(x, z)$



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