

# M2177.0043 Introduction to Deep Learning

## Lecture 13: Deep Generative Models<sup>1</sup>

Hyun Oh Song<sup>1</sup>

<sup>1</sup>Dept. of Computer Science and Engineering, Seoul National University

April 28, 2020

---

<sup>1</sup>Many slides and figures adapted Justin Johnson

## Last time

- ▶ Architectures
- ▶ Recurrent networks

# Outline

Unsupervised learning

Autoencoders

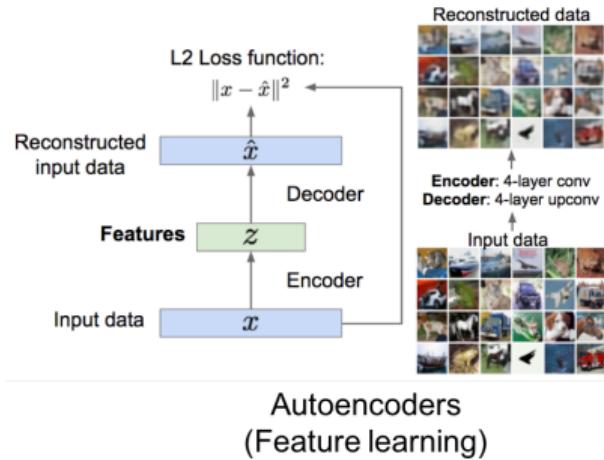
Variational autoencoders

VAE ELBO derivations

VAE training & inference

# Unsupervised learning

- ▶ **Data:**  $x$ , Just data, no labels.
- ▶ **Goal:** Learn some underlying hidden structure of the data
- ▶ **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc

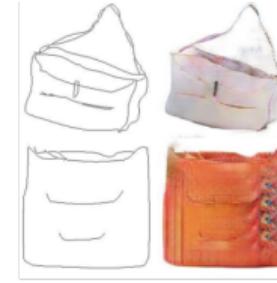


## Supervised vs unsupervised

- ▶ **Data:**  $(x,y)$ ,  $x$  is data,  $y$  is label
- ▶ **Goal:** Learn a function to map  $x$  to  $y$ .
- ▶ **Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.
- ▶ **Data:**  $x$ , Just data, no labels.
- ▶ **Goal:** Learn some underlying hidden structure of the data
- ▶ **Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc

# Why Generative models?

- ▶ Realistic samples for artwork, super-resolution, colorization, etc.
- ▶ Data augmentation
- ▶ Feature learning



# Outline

Unsupervised learning

Autoencoders

Variational autoencoders

VAE ELBO derivations

VAE training & inference

Autoencoders

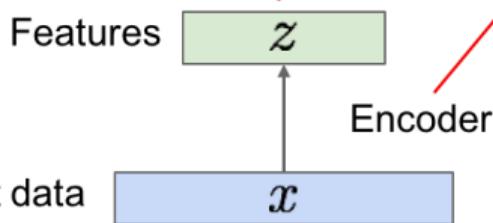
# Autoencoders

**$z$**  usually smaller than  **$x$**   
(dimensionality reduction)

Q: Why  
dimensionality  
reduction?

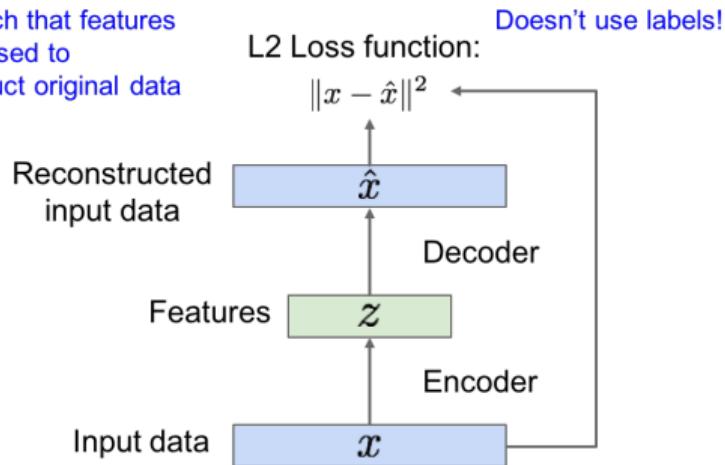
A: Want features  
to capture  
meaningful factors of  
variation in data

Originally: Linear +  
nonlinearity (sigmoid)  
Later: Deep, fully-connected  
Later: ReLU CNN



# Autoencoders

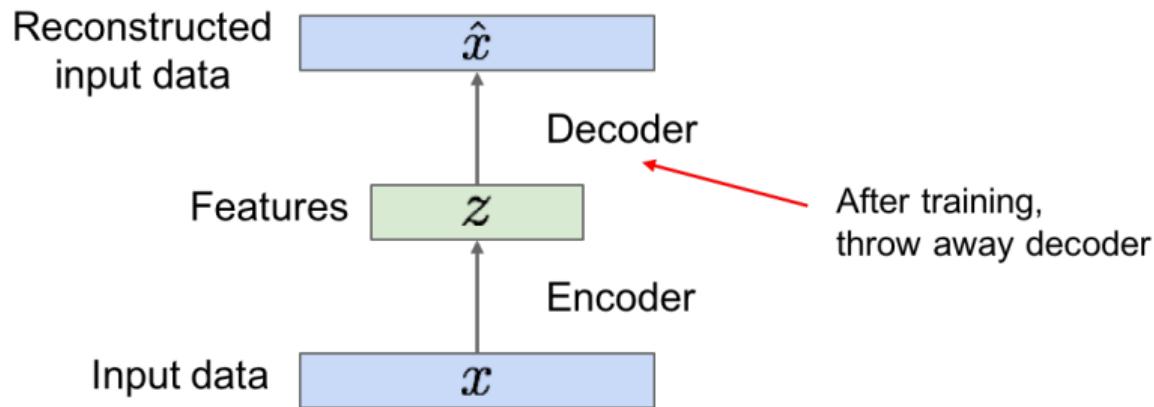
Train such that features can be used to reconstruct original data



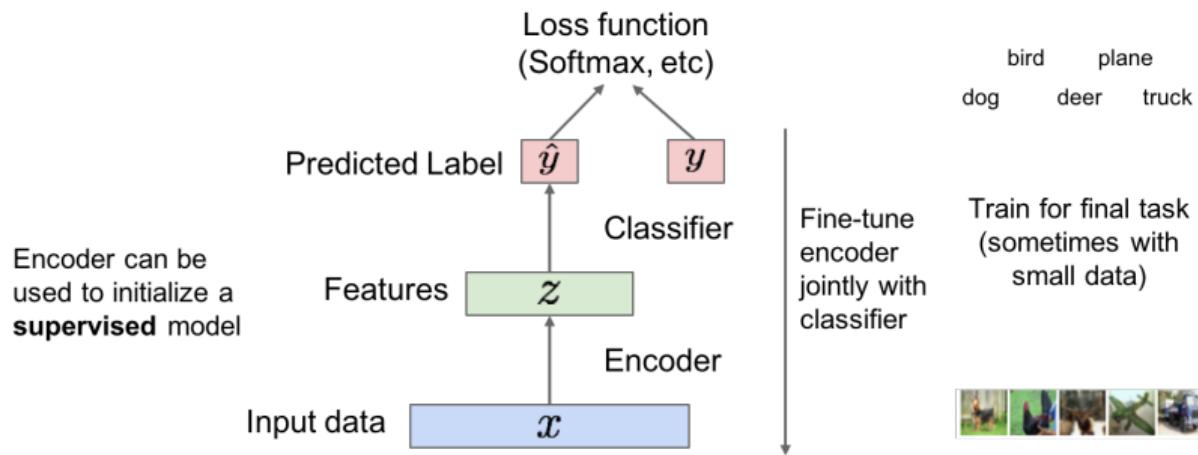
Autoencoders



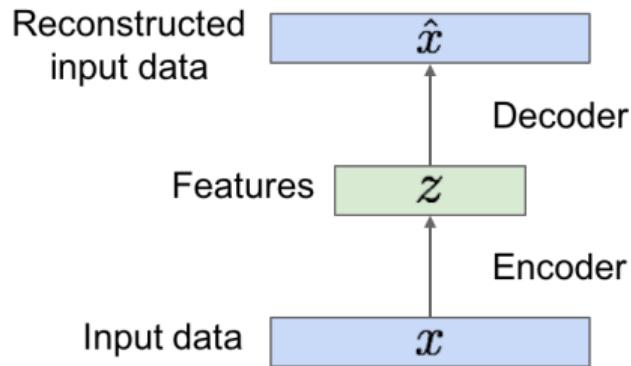
# Autoencoders



# Autoencoders



# Autoencoders



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?

# Outline

Unsupervised learning

Autoencoders

Variational autoencoders

VAE ELBO derivations

VAE training & inference

Variational autoencoders

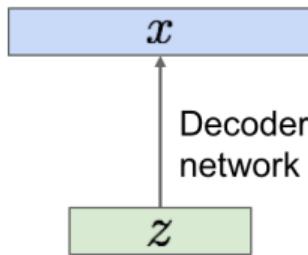
## Variational autoencoders

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

# Variational autoencoders

Sample from  
true conditional  
 $p_{\theta^*}(x | z^{(i)})$

Sample from  
true prior  
 $p_{\theta^*}(z)$



We want to estimate the true parameters  $\theta^*$  of this generative model.

How to train the model?

Learn model parameters to maximize likelihood  
of training data

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

Q: What is the problem with this?

Intractable!

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

## Variational autoencoders

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

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

# Variational autoencoders

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

↑  
Simple Gaussian prior

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

63

# Variational autoencoders

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

Decoder neural network

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

64

## Variational autoencoders

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

↑  
Intractible to compute  
 $p(x|z)$  for every  $z$ !

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

65

## Variational autoencoders

Data likelihood:  $p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$

Posterior density also intractable:  $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$



Intractable data likelihood

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

## Variational autoencoders

Data likelihood:  $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Posterior density also intractable:  $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

Solution: In addition to decoder network modeling  $p_\theta(x|z)$ , define additional encoder network  $q_\phi(z|x)$  that approximates  $p_\theta(z|x)$

Will see that this allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize

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

# Variational autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic

Mean and (diagonal) covariance  
of  $z|x$

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

$$\mu_{z|x}$$

$$\Sigma_{z|x}$$

Mean and (diagonal) covariance  
of  $x|z$

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

$$\mu_{x|z}$$

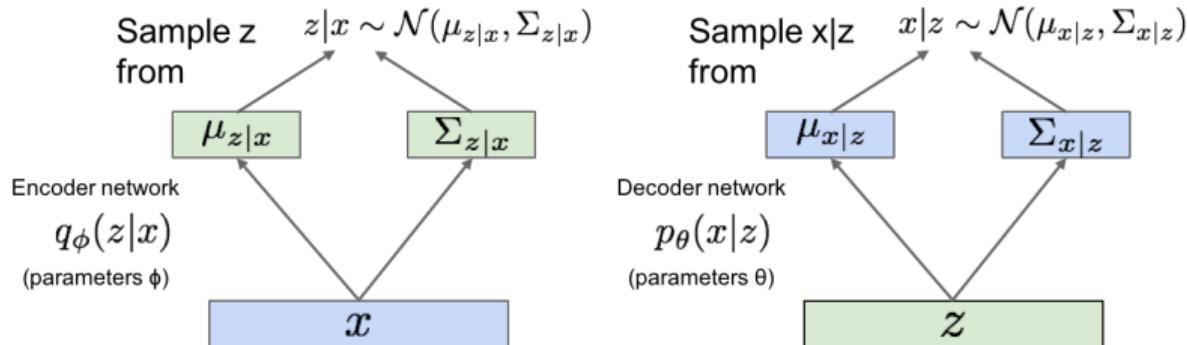
$$\Sigma_{x|z}$$

$$x$$

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

# Variational autoencoders

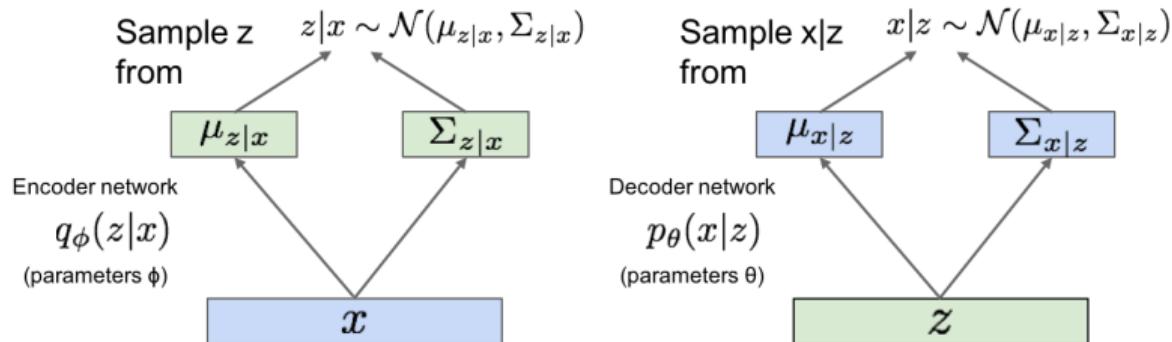
Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



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

# Variational autoencoders

Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



Encoder and decoder networks also called “recognition”/“inference” and “generation” networks

Kingma and Welling, “Auto-Encoding Variational Bayes”, ICLR 2014

## KL divergence

### Definition

For distributions  $P$  and  $Q$  of continuous random variable, the Kullback-Leibler divergence is defined to be the integral:

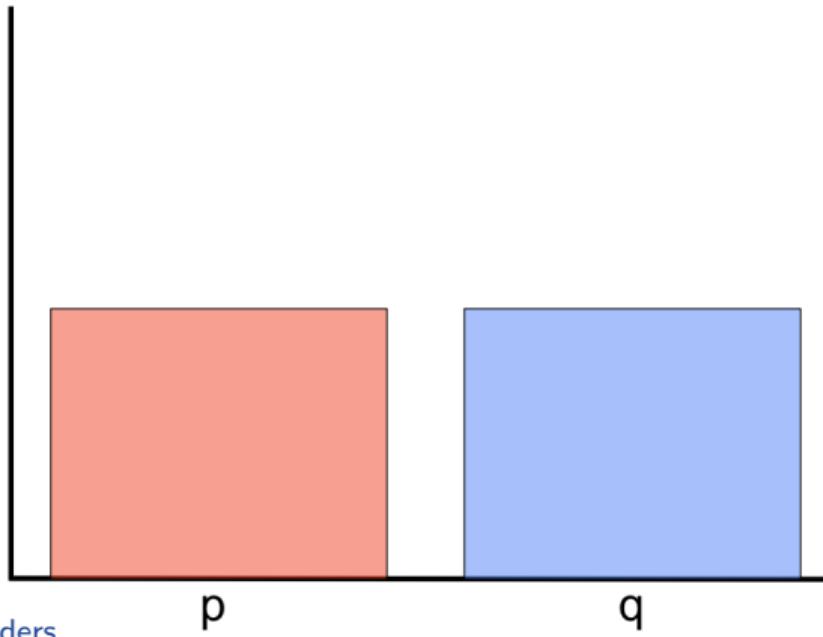
$$D_{KL}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx = \mathbb{E}_p \left[ \log \frac{p(x)}{q(x)} \right],$$

where  $p$  and  $q$  denote the densities of  $P$  and  $Q$ .

Note that despite the intuition of divergence as a distance measure,  $D(p \parallel q) \neq D(q \parallel p)$ . It is the case however, that if  $p = q$ , the  $D(p \parallel q) = 0$ .

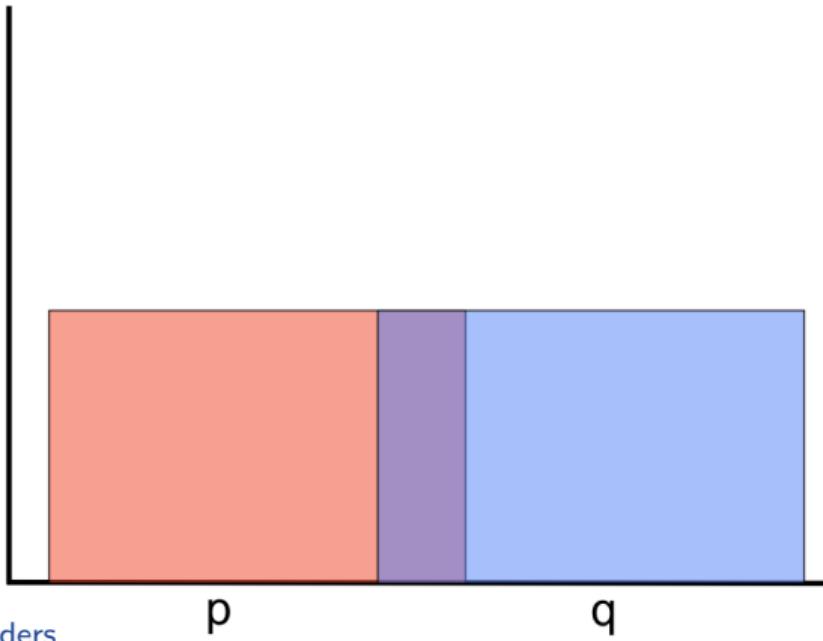
## KL divergence - support<sup>2</sup>

- ▶ KL divergence is defined iff  $Q(x) = 0 \implies P(x) = 0$ .
- ▶ Disjoint support, infinite KL



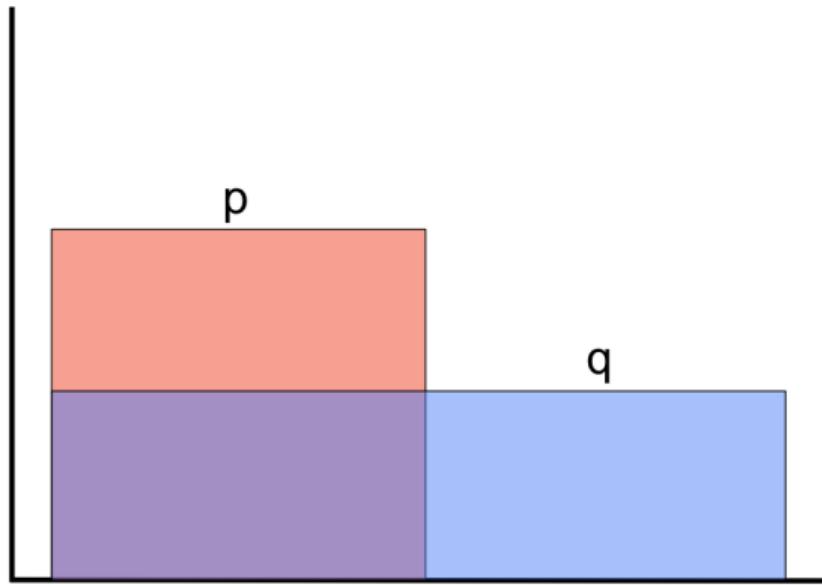
## KL divergence - support<sup>3</sup>

- ▶ KL divergence is defined iff  $Q(x) = 0 \implies P(x) = 0$ .
- ▶ Overlapping support but neither is a subset of the other. Infinite KL



## KL divergence - support<sup>4</sup>

- ▶ KL divergence is defined iff  $Q(x) = 0 \implies P(x) = 0$ .
- ▶  $\text{Supp}(p) \subseteq \text{Supp}(q)$ . Finite KL.



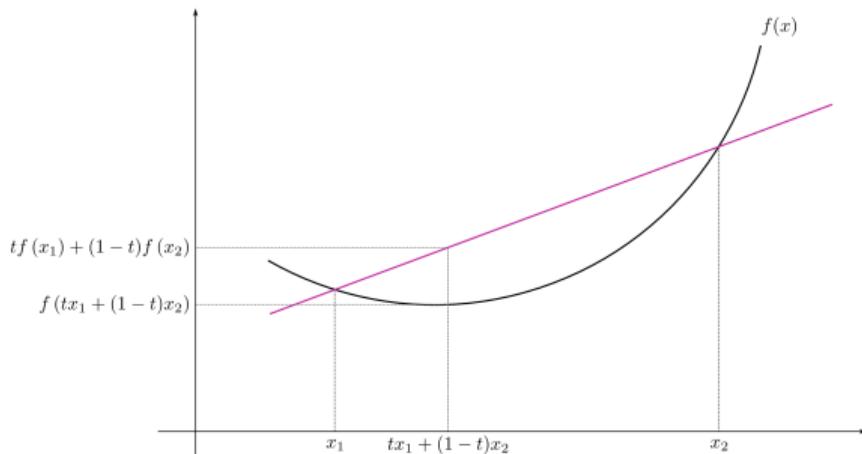
## KL divergence - summary

- ▶  $D_{KL}(p \parallel q)$  is finite only if the support of  $p$  is contained within the support of  $q$
- ▶ Definition of KL divergence uses the following conventions (Cover and Thomas, Elements of Information Theory)

$$0 \log \frac{0}{0} = 0, \quad 0 \log \frac{0}{q(x)} = 0, \quad p(x) \log \frac{p(x)}{0} = \infty$$

- ▶ Note, KL divergence can be infinite even if  $p(x)$  and  $q(x)$  are nonzero for all  $x$  (i.e. KL between Cauchy and Normal density).

## Jensen's inequality<sup>5</sup>



### Definition

Let  $f$  be a convex function and  $X$  be a random variable. Then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}X).$$

---

<sup>5</sup>From wikipedia

## KL divergence properties

### Theorem 1

*Non-negativity*  $D_{KL}(p \parallel q) \geq 0$

### Proof.

$D_{KL}(p \parallel q) = \mathbb{E}_p \left[ \log \frac{p(x)}{q(x)} \right]$ . By convexity of  $f(x) = -\log(x)$  and Jensen's inequality, we have that

$$D_{KL}(p \parallel q) = \mathbb{E}_p \left[ -\log \frac{q(x)}{p(x)} \right] \geq -\log \mathbb{E}_p \frac{q(x)}{p(x)} = 0,$$

where in the last step we used that  $\int q(x)dx = 1$ . Thus we conclude that KL-divergence is always non-negative. ■

# Outline

Unsupervised learning

Autoencoders

Variational autoencoders

VAE ELBO derivations

VAE training & inference

VAE ELBO derivations

## ELBO Derivation method I

- ▶ Examine the KL divergence between the intractable posterior  $p_\theta(z | x^{(i)})$  and the variational distribution  $q_\phi(z | x^{(i)})$ .

$$\begin{aligned} D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z | x^{(i)})) &= \mathbb{E}_{z \sim q_\phi(\cdot | x^{(i)})} \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \\ &= \mathbb{E}_{z \sim q_\phi(\cdot | x^{(i)})} \left[ \log q_\phi(z | x^{(i)}) - \log p_\theta(x^{(i)} | z) - \log p_\theta(z) + \log p_\theta(x^{(i)}) \right] \end{aligned}$$

Above is from applying Bayes rule. Log prior doesn't depend on  $z$ . Omit distribution on  $z$  from below for brevity.

$$\begin{aligned} \log p_\theta(x^{(i)}) &= \mathbb{E}_z \log p_\theta(x^{(i)} | z) + \mathbb{E}_z \left[ \log p_\theta(z) - \log q_\phi(z | x^{(i)}) \right] \\ &\quad + D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z | x^{(i)})) \\ &= \mathbb{E}_z \log p_\theta(x^{(i)} | z) - D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z)) \\ &\quad + D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z | x^{(i)})) \end{aligned}$$

## ELBO Derivation method I

- ▶ Mathematical meaning of individual terms

$$\log p_\theta(x^{(i)}) = \underbrace{\mathbb{E}_z \log p_\theta(x^{(i)} | z)}_{\text{Decoder gives } p_\theta(x|z)} - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z))}_{\begin{array}{l} \text{KL between two Gaussians has closed form sol} \\ := \mathcal{L}(x^{(i)}, \theta, \phi), \text{ Differentiable lower bound} \end{array}} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z | x^{(i)}))}_{p_\theta(z|x^{(i)}) \text{ intractable but KLD is always greater than zero}}$$

- ▶  $\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$  called variational lower bound (ELBO).
- ▶ The optimization problem is to maximize the lower bound (instead of maximizing the data log likelihood)  
 $\theta^*, \phi^* = \operatorname{argmax}_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$

## ELBO Derivation method II

- ▶ Unlike the previous derivation method, start from the log likelihood

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbb{E}_{z \sim q_\phi(\cdot | x^{(i)})} \log p_\theta(x^{(i)}) \quad (p_\theta(x^{(i)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \\ &\qquad\qquad\qquad (\text{Bayes' Rule, note } p_\theta(z | x^{(i)}) \text{ intractable}) \\ &= \mathbb{E}_z \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \underbrace{\frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})}}_{\text{(Mult. by 1)}} \\ &= \mathbb{E}_z \log p_\theta(x^{(i)} | z) - \mathbb{E}_z \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} + \mathbb{E}_z \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \\ &= \underbrace{\mathbb{E}_z \log p_\theta(x^{(i)} | z) - D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z))}_{:= \mathcal{L}(x^{(i)}, \theta, \phi), \text{ Differentiable lower bound}} \\ &\quad + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z | x^{(i)}))}_{\geq 0, \text{ Intractable}}\end{aligned}$$

## ELBO Derivation method III

- ▶ Start from the joint distribution and apply Jensen's inequality

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \log \int p_{\theta}(x^{(i)}, z) dz \\&= \log \int \frac{p_{\theta}(x^{(i)}, z)}{q_{\phi}(z \mid x^{(i)})} q_{\phi}(z \mid x^{(i)}) dz \quad (\text{Mult. by 1}) \\&\geq \int q_{\phi}(z \mid x^{(i)}) \log \frac{p_{\theta}(x^{(i)}, z)}{q_{\phi}(z \mid x^{(i)})} dz \quad (\text{Jensen's inequality}) \\&= \mathbb{E}_{z \sim q_{\phi}(\cdot \mid x^{(i)})} \log \frac{p_{\theta}(x^{(i)} \mid z) p_{\theta}(z)}{q_{\phi}(z \mid x^{(i)})} \\&= \mathbb{E}_z \log p_{\theta}(x^{(i)} \mid z) - \mathbb{E}_z \log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)} \\&= \underbrace{\mathbb{E}_z \log p_{\theta}(x^{(i)} \mid z) - D_{KL}(q_{\phi}(z \mid x^{(i)}) \parallel p_{\theta}(z))}_{:= \mathcal{L}(x^{(i)}, \theta, \phi), \text{ Differentiable lower bound}}\end{aligned}$$

# Outline

Unsupervised learning

Autoencoders

Variational autoencoders

VAE ELBO derivations

VAE training & inference

## ELBO intuition



$$\mathcal{L}(x^{(i)}, \theta, \phi) = \underbrace{\mathbb{E}_{z \sim q_\phi(\cdot | x^{(i)})} \log p_\theta(x^{(i)} | z)}_{\text{Decoder gives } p_\theta(x|z)} - \underbrace{D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z))}_{\text{KL between two Gaussians has closed form sol}}$$

- ▶  $\mathbb{E}_z \log p_\theta(x^{(i)} | z)$  : Expectation of all samples of  $z$  sampled from passing  $x^{(i)}$  through the encoder network ( $q_\phi(\cdot | x^{(i)})$ ), sampling  $z$ , computing the log likelihood of decoder reconstruction from the sample  $z$ , ( $p_\theta(\cdot | z)$ ). Encourage good reconstruction.
- ▶  $D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z))$  : Make the variational approximate posterior distribution close to prior.
- ▶ Typically, we use standard normal or uniform distribution for the prior.

## KL between Gaussians

►  $D_{KL}(\mathcal{N}(\mu_0, \Sigma_0) \parallel \mathcal{N}(\mu_1, \Sigma_1)) =$

$$\frac{1}{2} \left( \text{Tr} (\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) - k + \log \left( \frac{\det \Sigma_1}{\det \Sigma_0} \right) \right)$$

► In our case, this simplifies to:

$$D_{KL}(\mathcal{N}(\mu(x^{(i)}), \Sigma(x^{(i)})) \parallel \mathcal{N}(0, I)) =$$

$$\frac{1}{2} \left( \text{Tr} \Sigma(x^{(i)}) + \mu(x^{(i)})^\top \mu(x^{(i)}) - k - \log \det \Sigma(x^{(i)}) \right)$$

► In case of diagonal covariance matrix:

$$D_{KL}(\mathcal{N}(\mu(x^{(i)}), \boldsymbol{\sigma}(x^{(i)})) \parallel \mathcal{N}(0, I)) =$$

$$\frac{1}{2} \sum_j^d 1 + \log \sigma_j^2(x^{(i)}) - \mu_j^2(x^{(i)}) - \sigma_j^2(x^{(i)})$$

► Derive these for fun at home!

► Can you eyeball the simple case from the general case?

## VAE - training

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

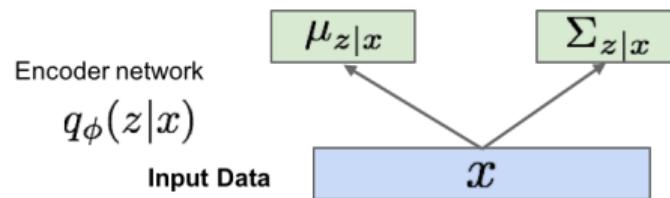
Input Data

$x$

## VAE - training

Putting it all together: maximizing the likelihood lower bound

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



## VAE - training

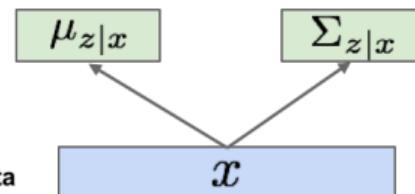
Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior

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

Input Data

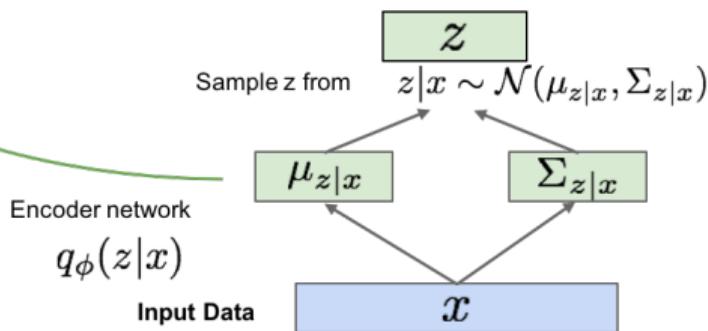


## VAE - training

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

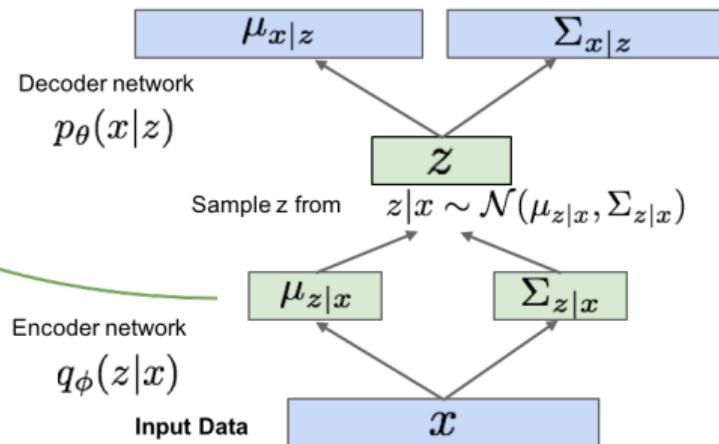


## VAE - training

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close to prior



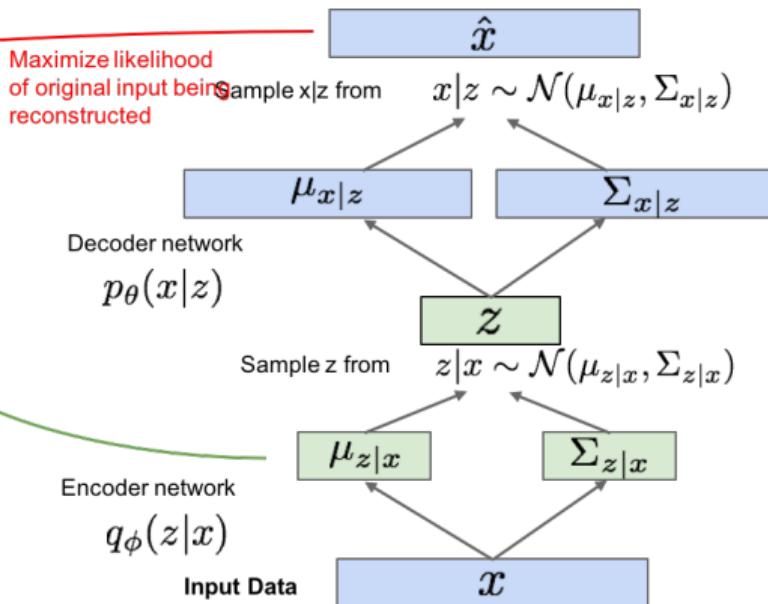
# VAE - training

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



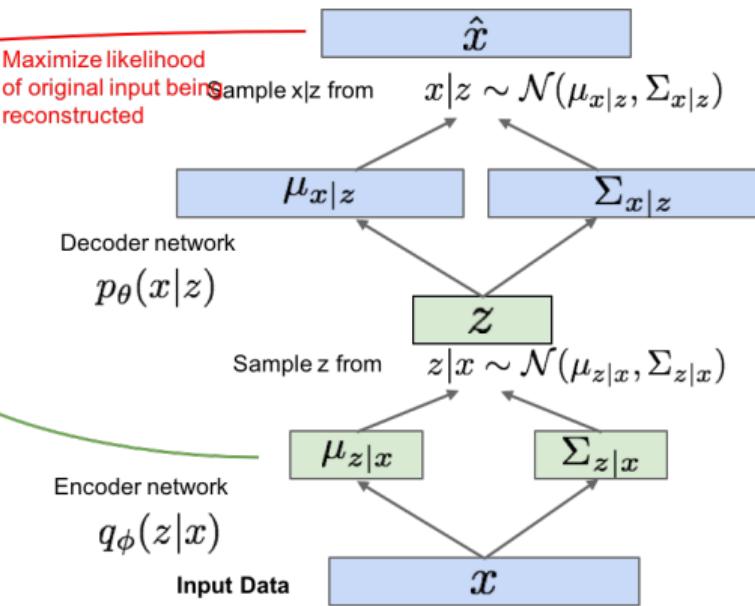
# VAE - training

Putting it all together: maximizing the likelihood lower bound

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

Make approximate posterior distribution close

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

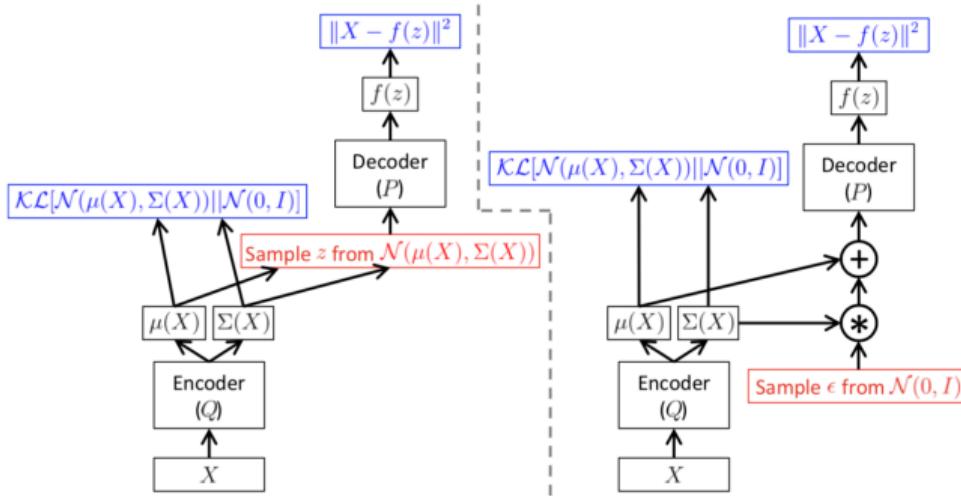


## Reparameterization trick

- ▶ Reconstruction term  $\mathbb{E}_z \log p_\theta(x^{(i)} | z)$  has sampling procedure which is not differentiable.
- ▶ Move sampling to an input layer. Given  $\mu(x^{(i)}), \Sigma(x^{(i)})$  from the encoder  $q_\phi(z | x^{(i)})$ , we can sample from  $\mathcal{N}(\mu(x^{(i)}), \Sigma(x^{(i)}))$  by first sampling  $\epsilon \sim \mathcal{N}(0, I)$ , then computing  $z = \mu(x^{(i)}) + A(x^{(i)})\epsilon$ . Remember the homework problem?
- ▶  $A$  could be a Cholesky factor or  $QD^{1/2}$  factor of  $\Sigma(x^{(i)})$
- ▶ Reconstruction is now reparameterized as  
$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \log p_\theta(x^{(i)} | z) = \mu(x^{(i)}) + A(x^{(i)})\epsilon$$

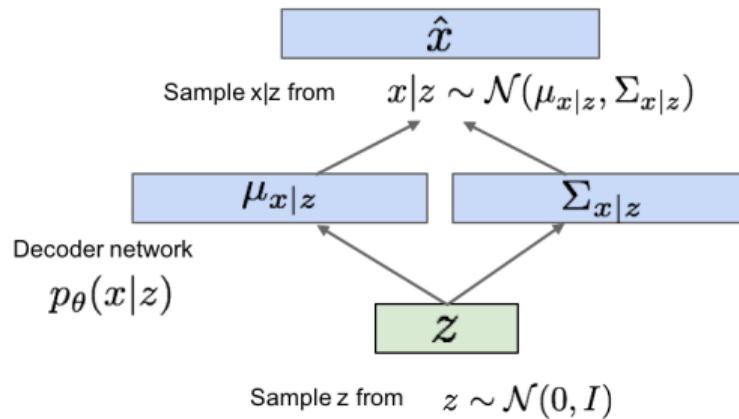
## Reparameterization trick<sup>6</sup>

- ▶ Left is without the reparameterization trick and right is with it. Red shows sampling operations that are non-differentiable. Blue shows loss layers (with Gaussian log likelihood). The feedforward behavior of these networks is identical, but backpropagation can be applied only to the right network.



## VAE - inference

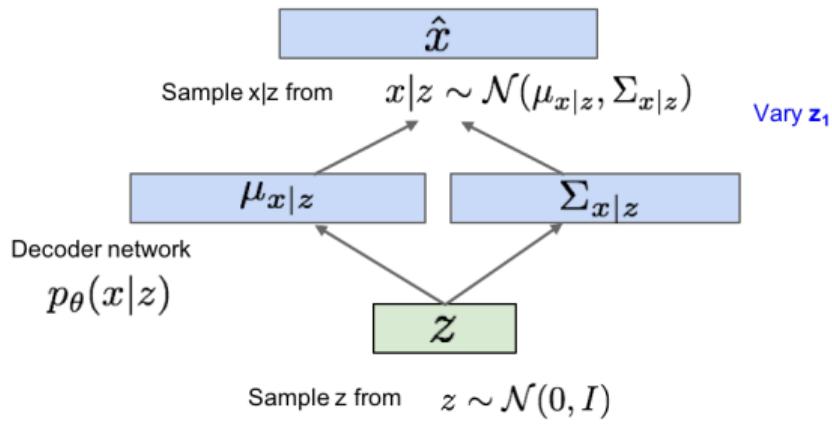
Use decoder network. Now sample z from prior!



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

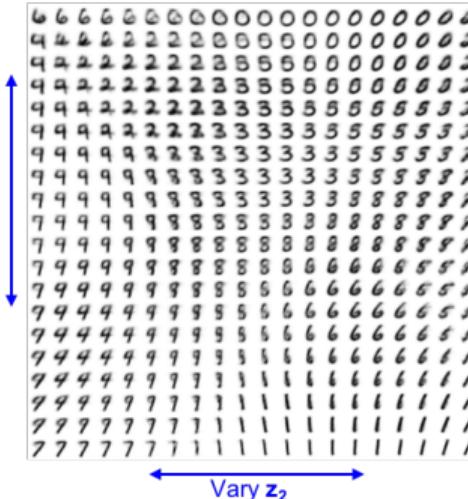
## VAE - inference

Use decoder network. Now sample z from prior!



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

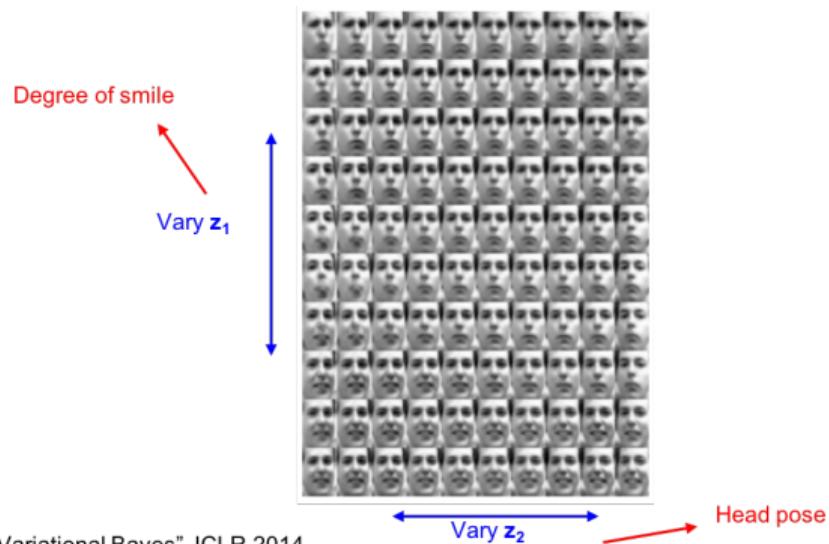
Data manifold for 2-d  $z$



## VAE - inference

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

## VAE - inference

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

Degree of smile  
Vary  $\mathbf{z}_1$

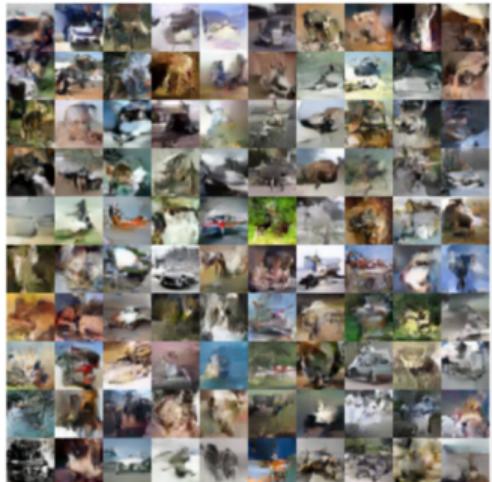


Vary  $\mathbf{z}_2$

Head pose

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

## VAE - inference



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.

## Summary

- ▶ Principled approach to generative models
- ▶ Allows inference of  $q(z|x)$  and can be useful feature representation for other tasks
- ▶ Active areas of research: More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian. Incorporating structure in latent variables (as opposed to i.i.d z).
- ▶ How about discrete latent variables? (See Jeong & Song ICML19)