



北京航空航天大學  
BEIHANG UNIVERSITY

# 生成式AI与大模型第7讲

## 生成模型

Beihang University

人工智能研究院  
黄雷

01

## 生成模型基础

# Discriminative model (判别模型)

- Learning a map from X to y



→ Cat

Classification



**DOG, DOG, CAT**

Object Detection

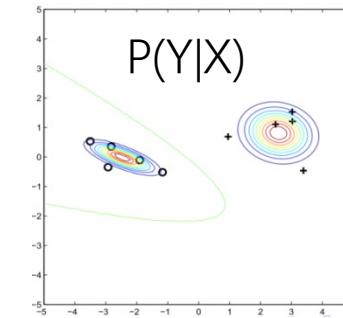
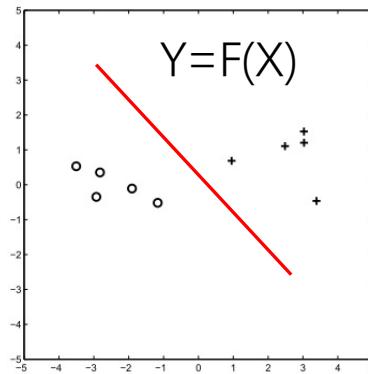


**GRASS, CAT,  
TREE, SKY**

Semantic Segmentation

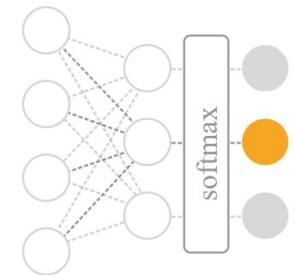
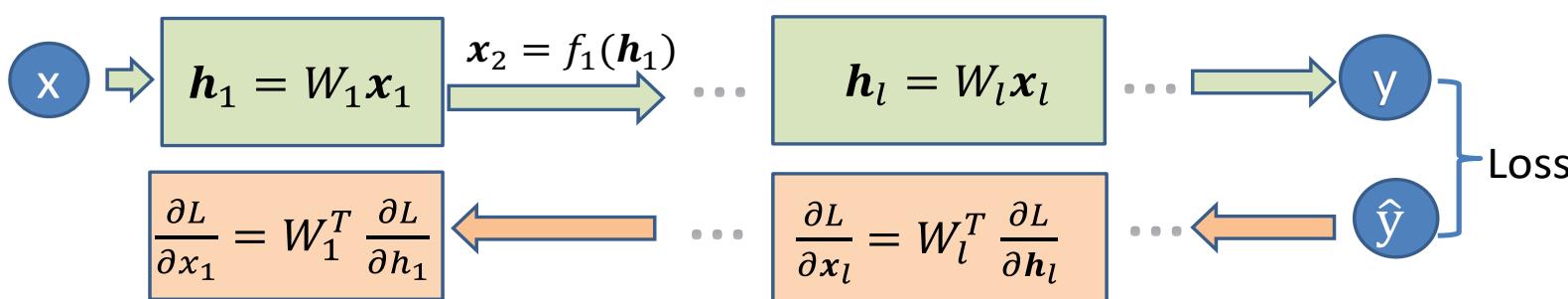
# Why probability?

- $Y=f(X)$



# Why distribution?

## Deep neural networks visualization



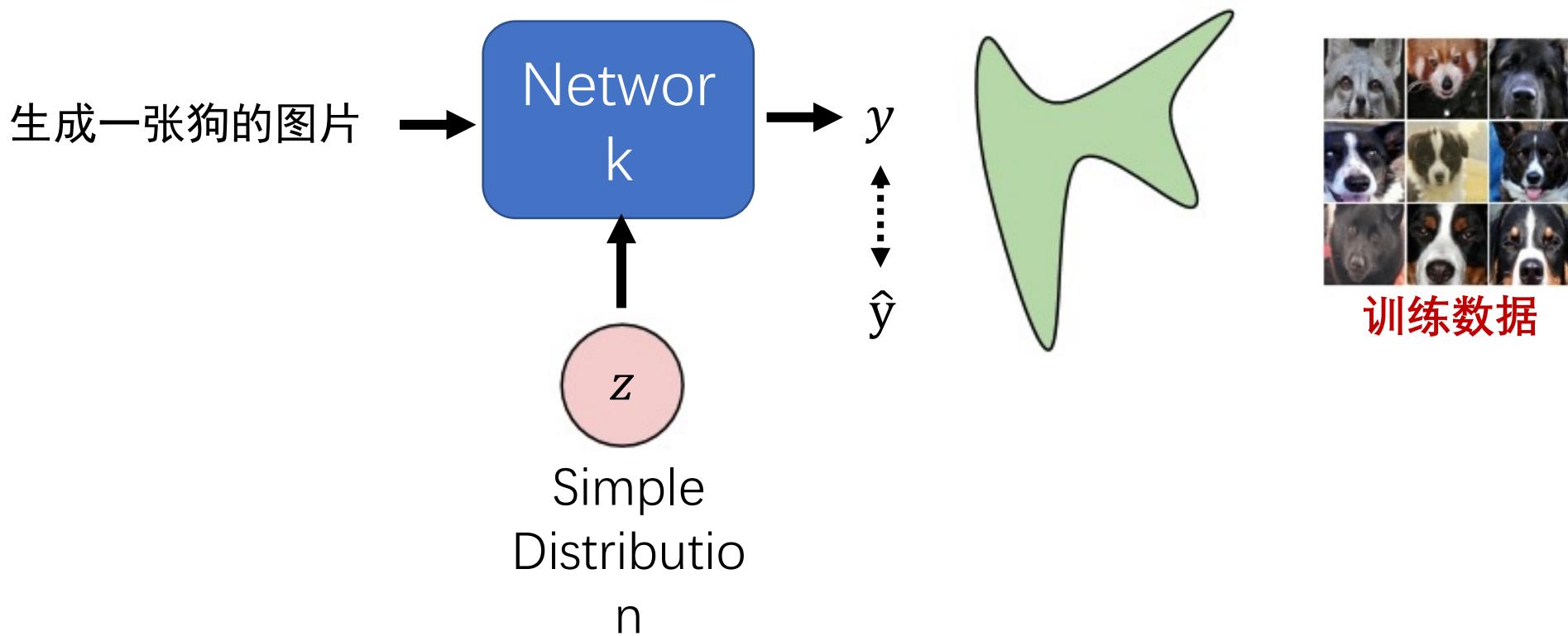
生成一张狗的图片

Networ  
k

$y$

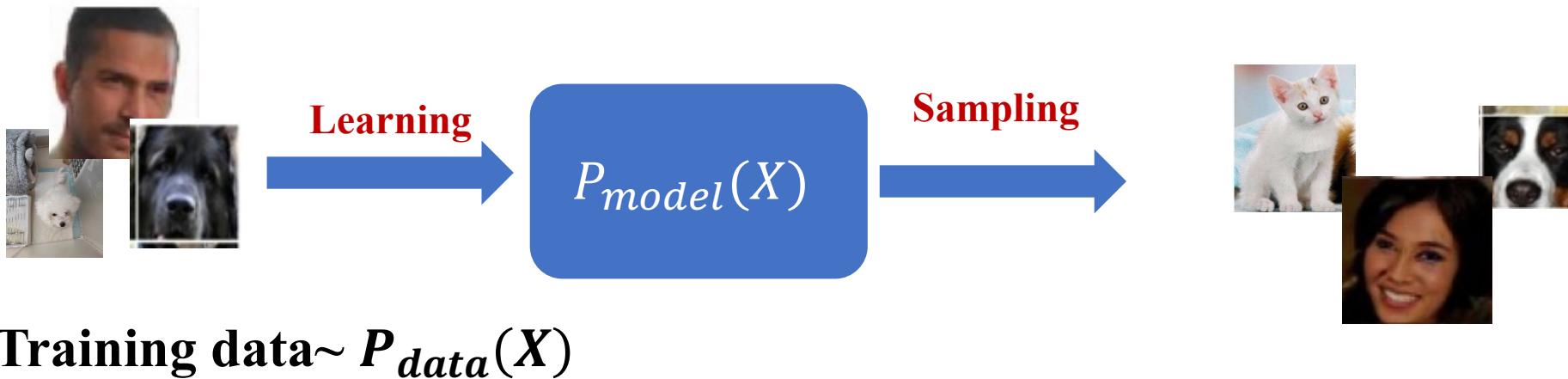


训练数据



# Generative modeling

Given training data, generate new samples from same distribution

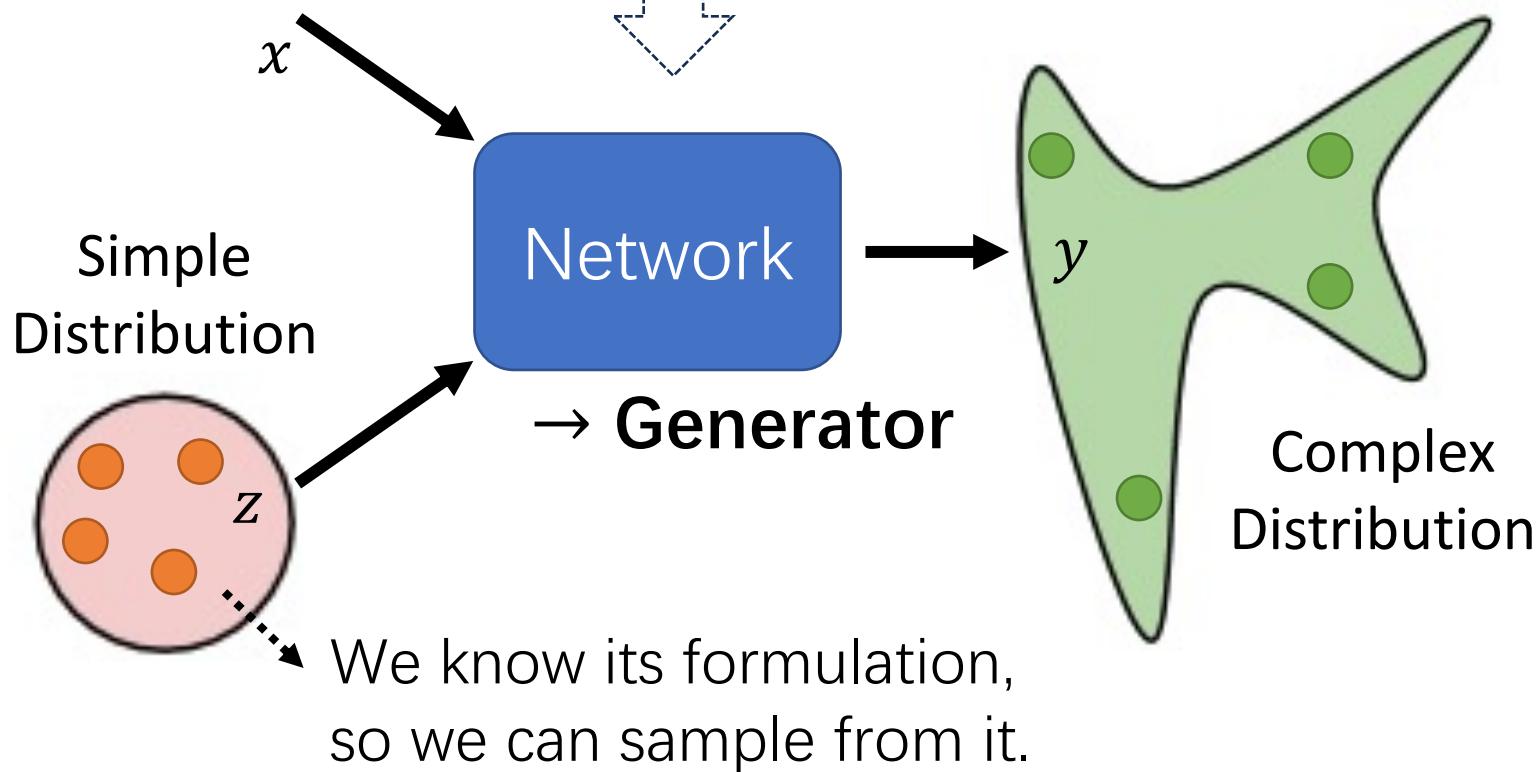
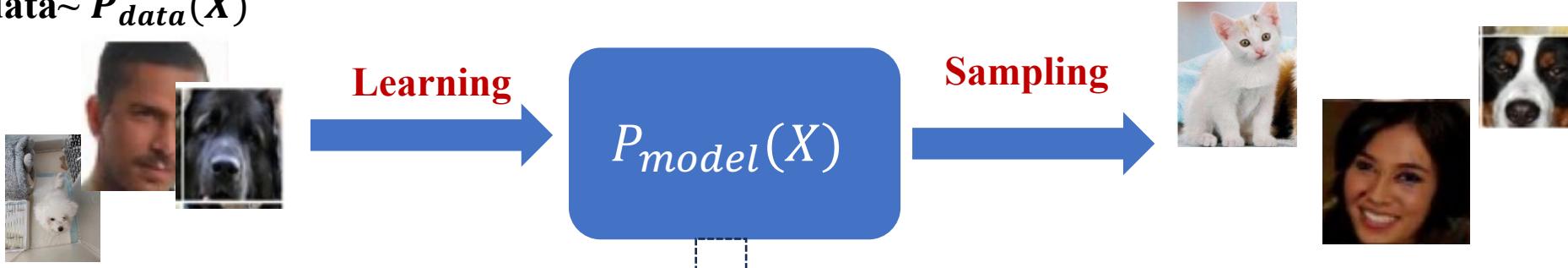


Objective:

- Train: Learn  $P_{model}(X)$  that approximates  $P_{data}(X)$
- Inference: Sampling new  $x$  from  $P_{model}(X)$

# Network as Generator

Training data~  $P_{data}(X)$



02

## 自回归模型生成

# Language Model

- Auto-Regressive (自回归)
  - GPTs

(all characters)

## One-hot Encoding

$$\text{深} = [1 \ 0 \ 0 \ 0 \ 0]$$

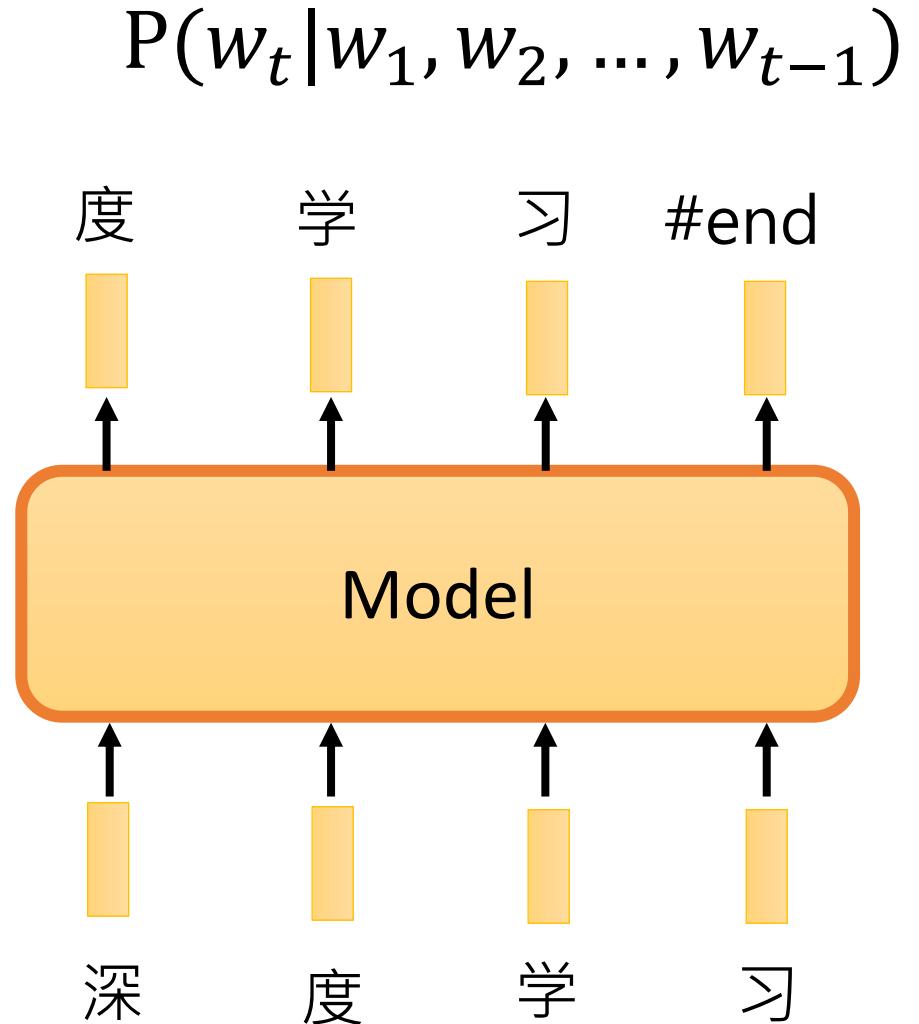
$$\text{度} = [0 \ 1 \ 0 \ 0 \ 0]$$

$$\text{学} = [0 \ 0 \ 1 \ 0 \ 0]$$

$$\text{习} = [0 \ 0 \ 0 \ 1 \ 0]$$

$$\#\text{end} = [0 \ 0 \ 0 \ 0 \ 1]$$

深	0.1
度	0.7
学	0.1
习	0.05
#end	0.05



# Vision Model

## Explicit density model

$$p(x) = p(x_1, x_2, \dots, x_n)$$



Likelihood of  
image  $x$



Joint likelihood of each  
pixel in the image

Fully visible belief network(FVBN)

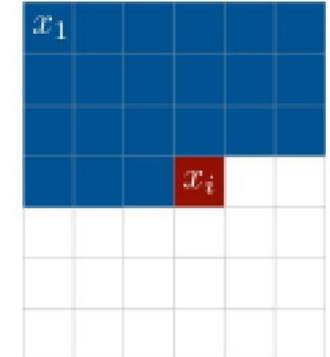
# Vision Model

Explicit density model

Use chain rule to decompose likelihood of an image  $x$  into product of 1-d distributions:

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

↑                              ↑  
Likelihood of                  Probability of i'th pixel value  
image  $x$                           given all previous pixels



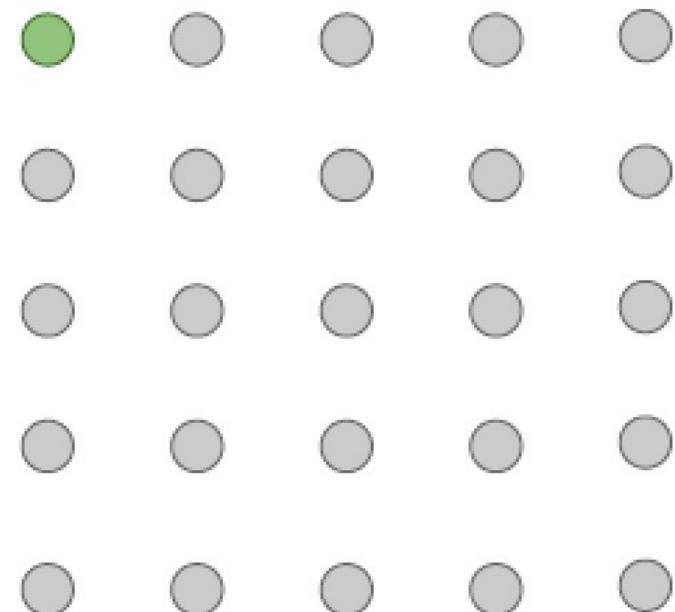
Then maximize likelihood of training data

Fully visible belief network(FVBN)

# PixelRNN

## PixelRNN *[van der Oord et al. 2016]*

Generate image pixels starting from corner



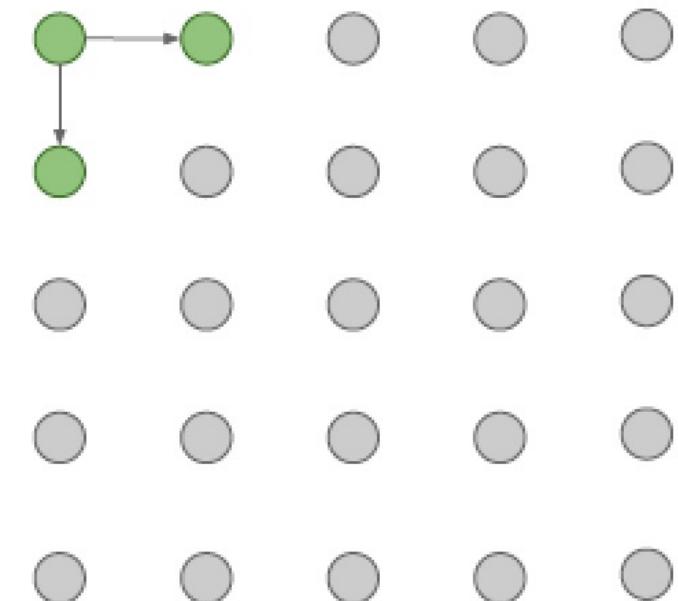
Dependency on previous pixels modeled  
using an RNN (LSTM)

# PixelRNN

## PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

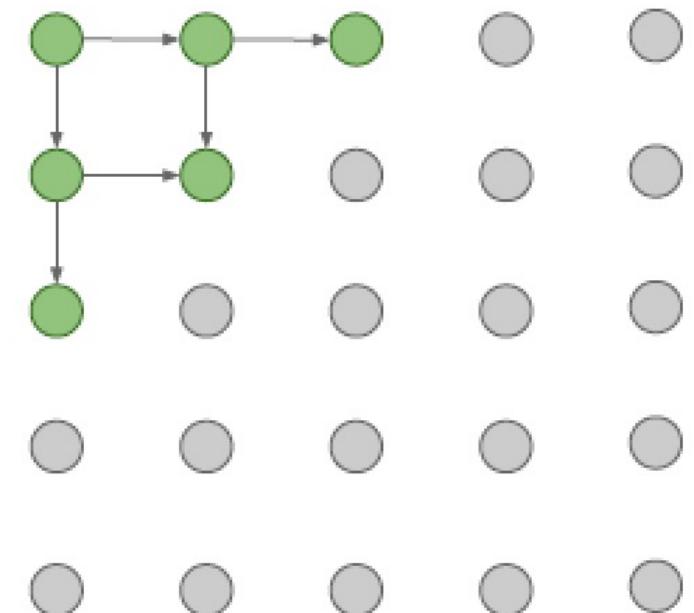


# PixelRNN

## PixelRNN *[van der Oord et al. 2016]*

Generate image pixels starting from corner

Dependency on previous pixels modeled  
using an RNN (LSTM)



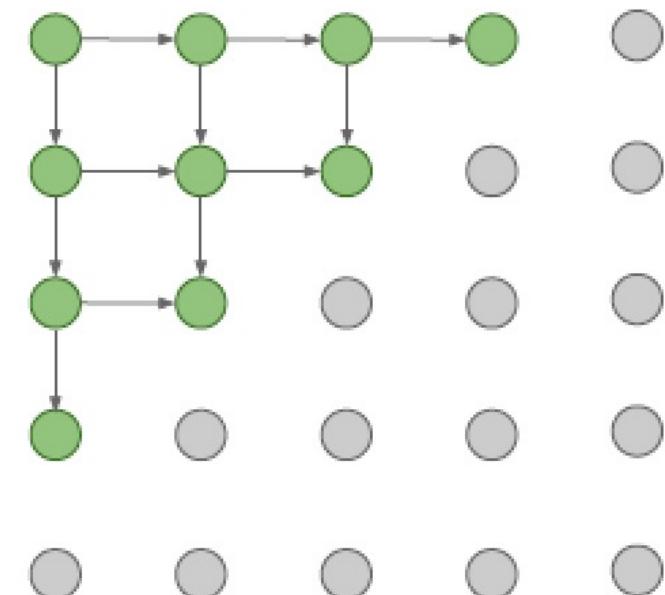
# PixelRNN

## PixelRNN *[van der Oord et al. 2016]*

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Drawback: sequential generation is slow in both training and inference!

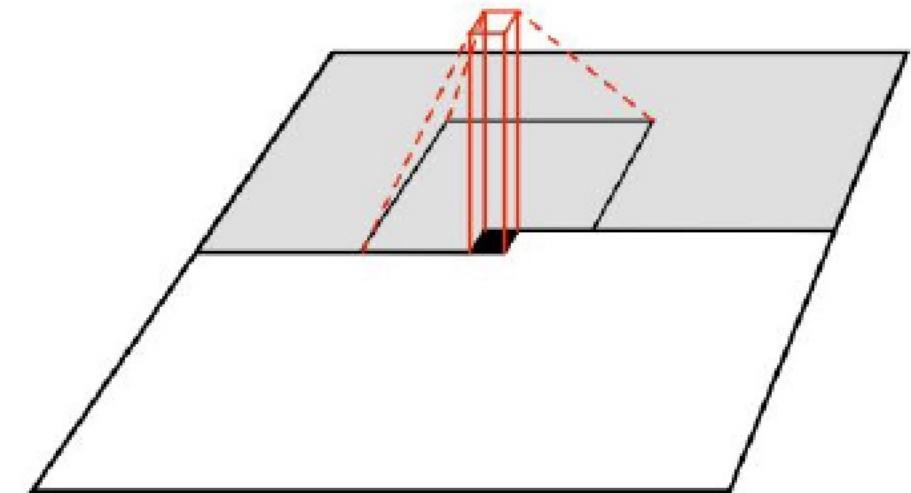


# PixelRNN

## PixelCNN [van der Oord et al. 2016]

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region  
**(masked convolution)**



## PixelCNN [van der Oord et al. 2016]

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region (masked convolution)

Training is faster than PixelRNN  
(can parallelize convolutions since context region values known from training images)

Generation is still slow:

For a 32x32 image, we need to do forward passes of the network 1024 times for a single image

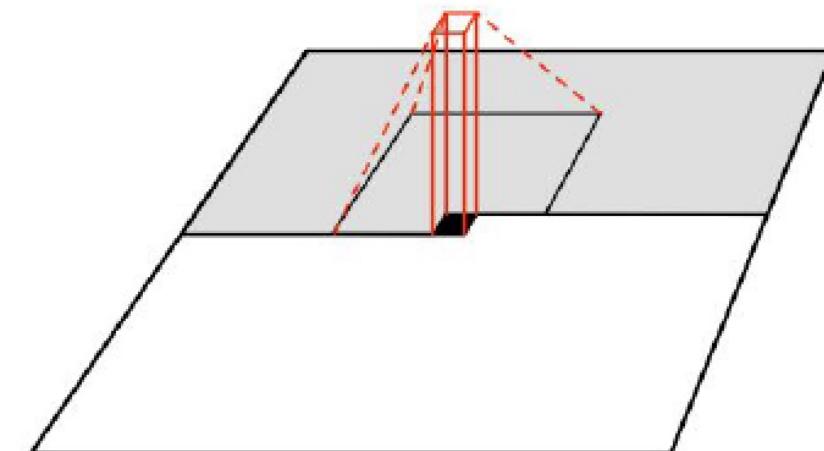
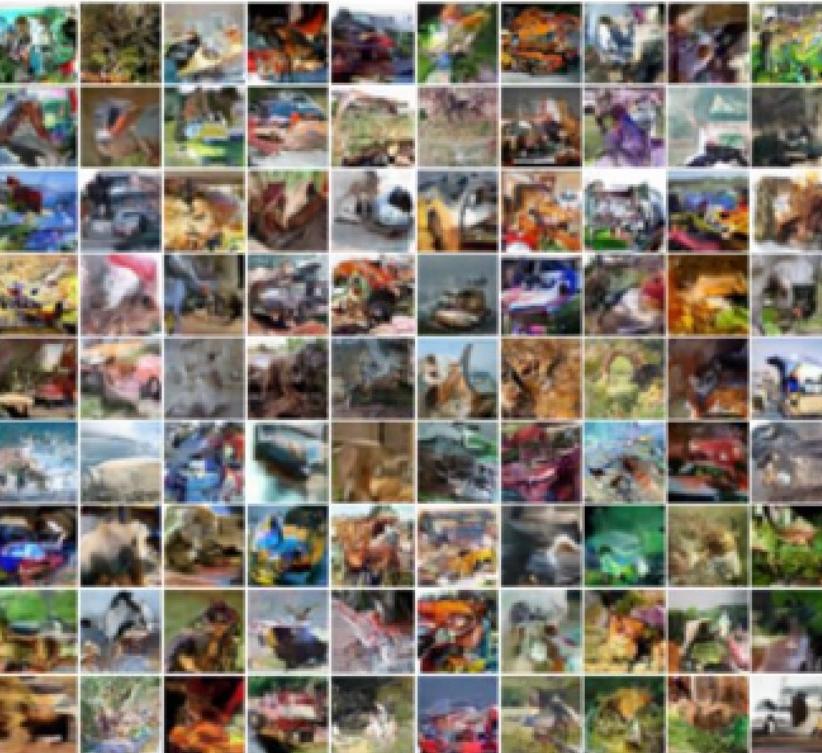


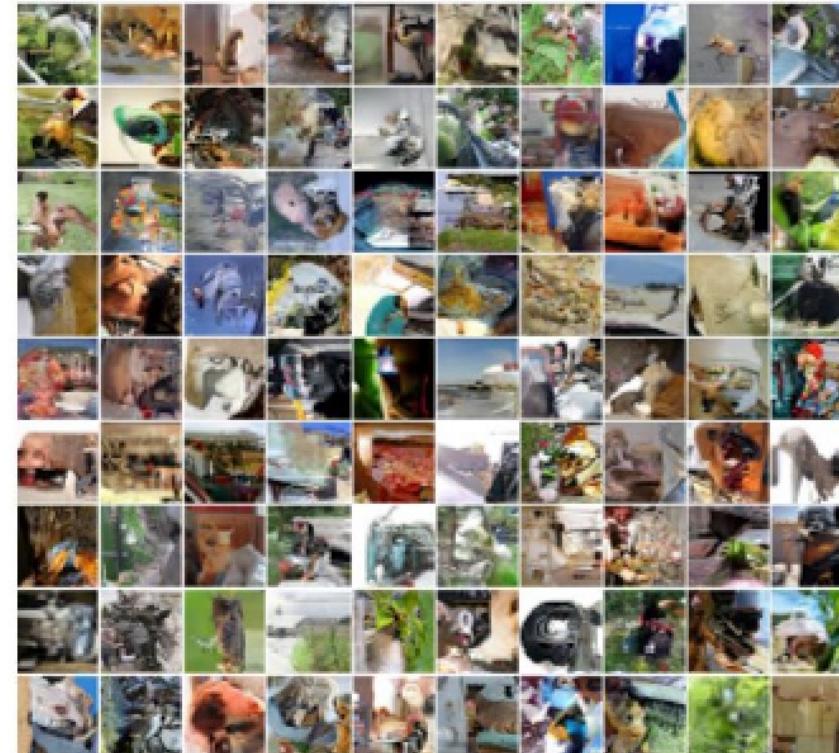
Figure copyright van der Oord et al., 2016. Reproduced with permission.

# PixelRNN

## Generation Samples



32x32 CIFAR-10



32x32 ImageNet

03

## 变分自编码器

# Variational Autoencoders (VAE)

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

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

Variational Autoencoders (VAEs) define intractable density function with latent  $\mathbf{z}$ :

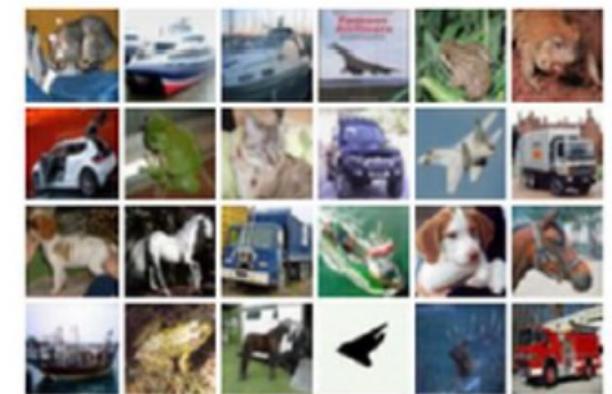
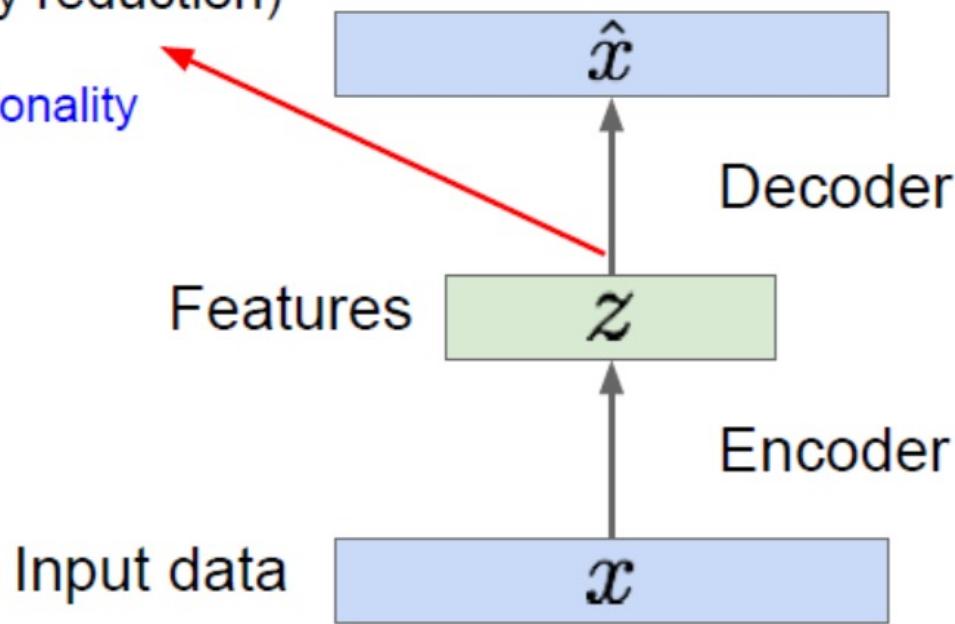
$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

# Background: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

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

Q: Why dimensionality reduction?



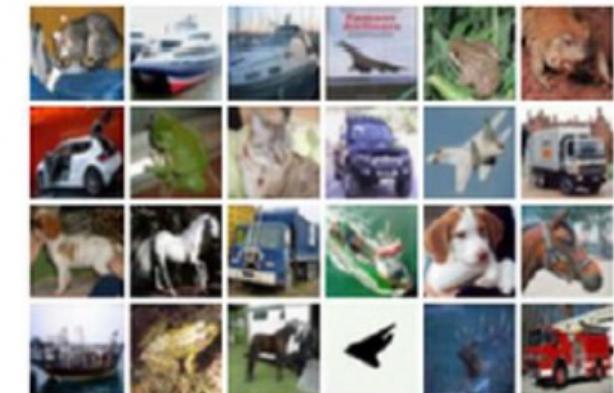
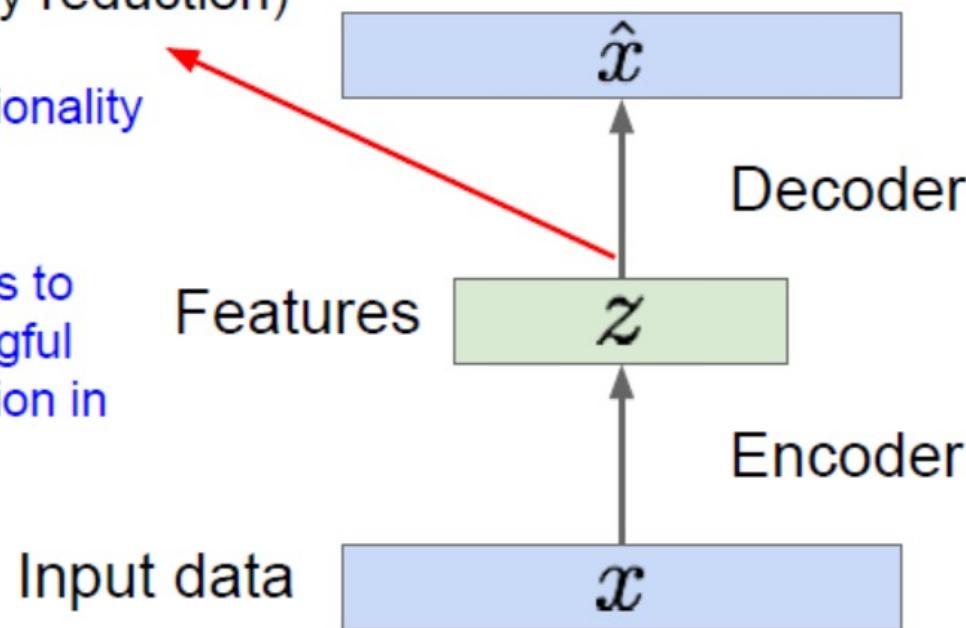
# Background: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

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

Q: Why dimensionality reduction?

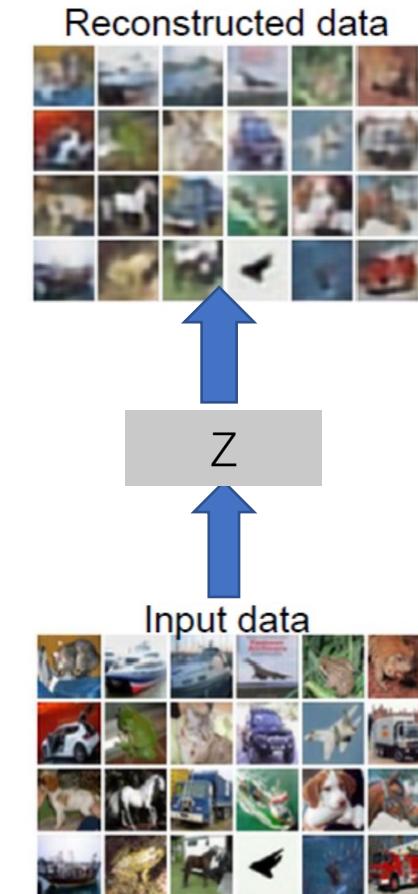
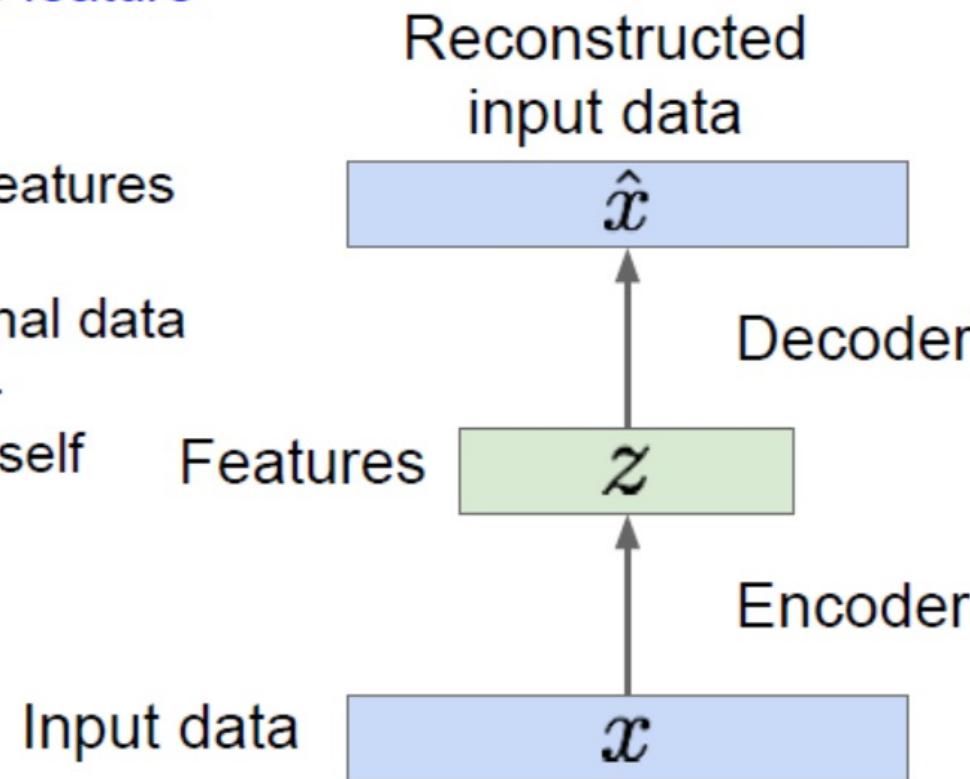
A: Want features to capture meaningful factors of variation in data



# Background: Autoencoders

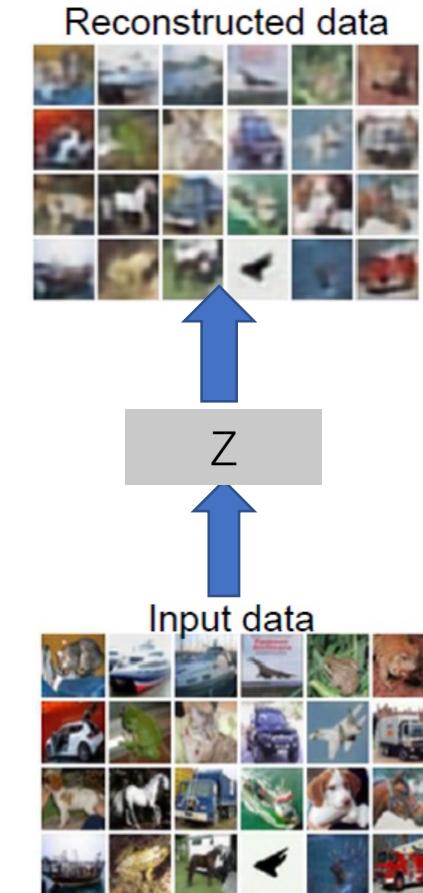
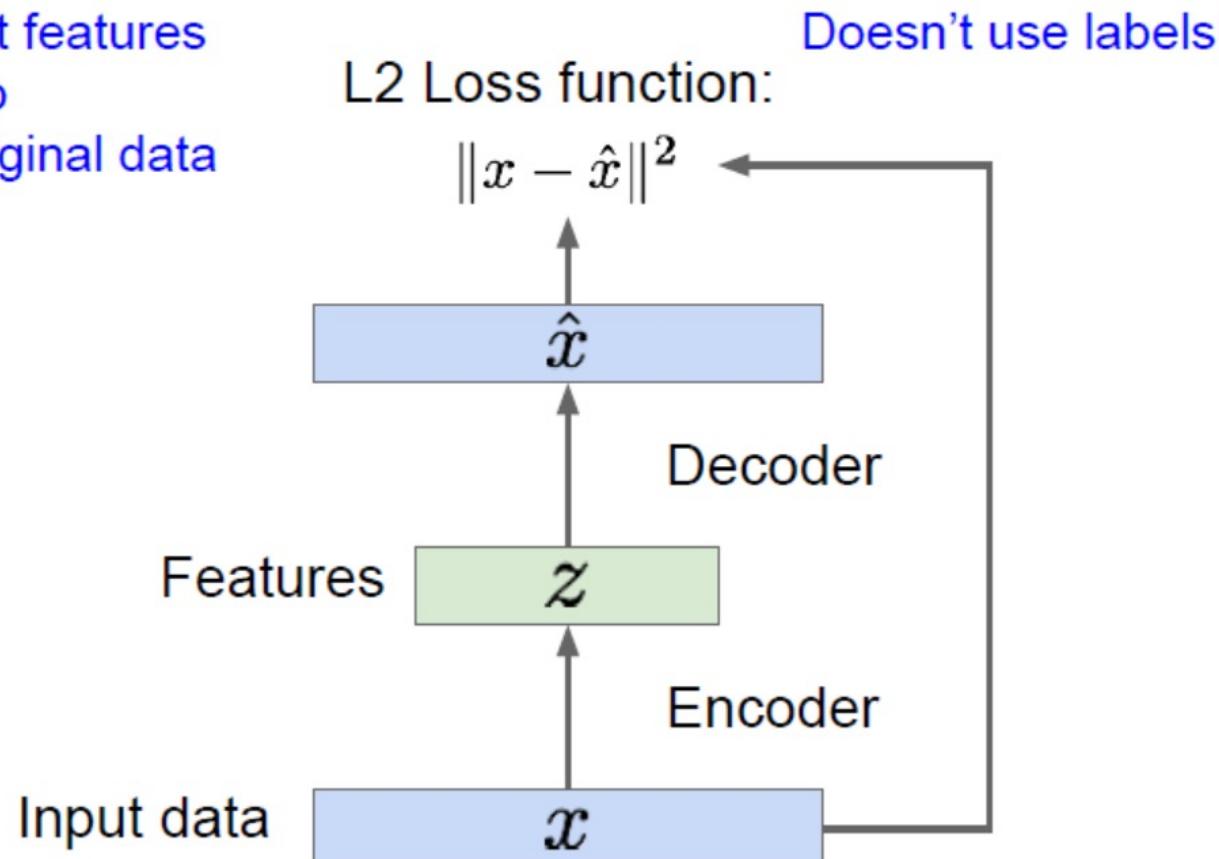
How to learn this feature representation?

Train such that features can be used to reconstruct original data  
“Autoencoding” - encoding input itself

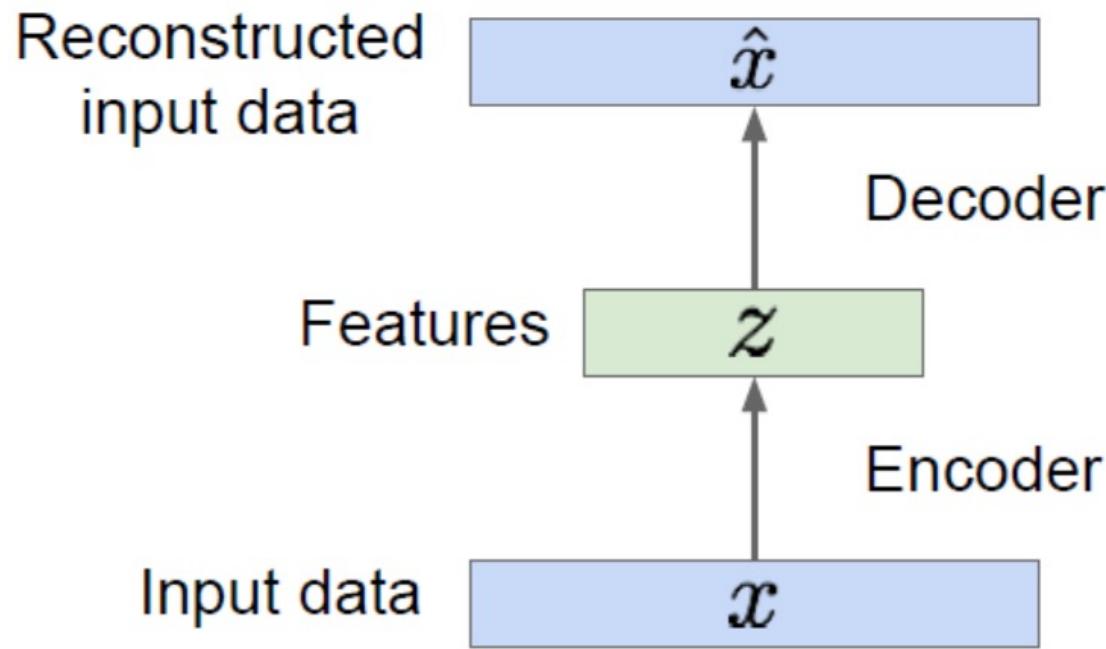


# Background: Autoencoders

Train such that features can be used to reconstruct original data



# Background: Autoencoders



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

Features capture factors of variation in training data.

But we can't generate new images from an autoencoder because we don't know the space of  $z$ .

How do we make autoencoder a **generative model**?

# Variational Autoencoders (VAE)

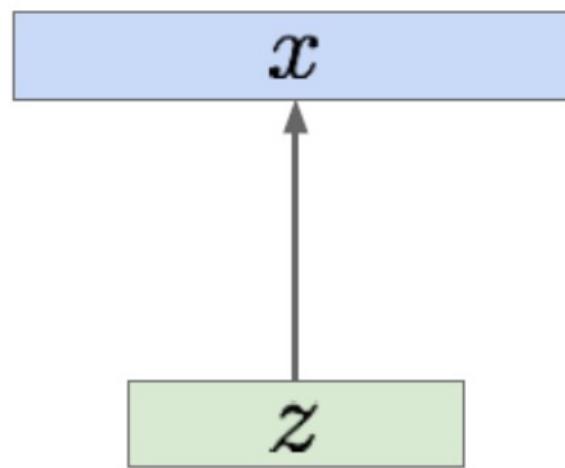
Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Assume training data  $\{x^{(i)}\}_{i=1}^N$  is generated from the distribution of unobserved (latent) representation  $z$

Sample from  
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

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



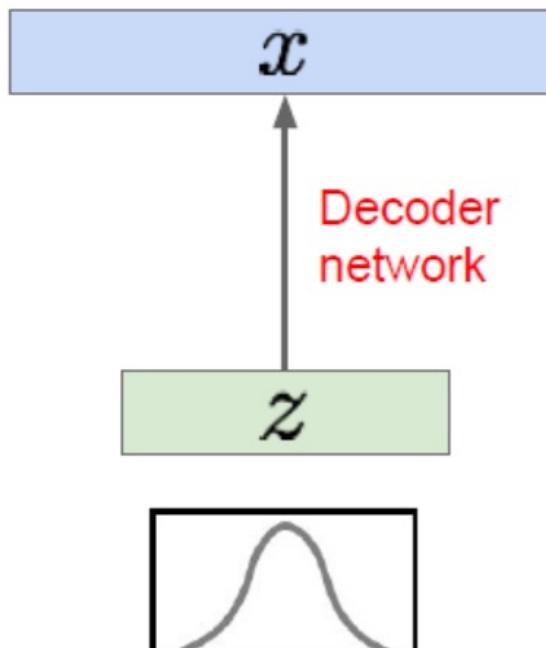
# Variational Autoencoders (VAE)

➤ How to represent



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

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



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

How should we represent this model?

Choose prior  $p(z)$  to be simple, e.g. Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

Conditional  $p(x|z)$  is complex (generates image) => represent with neural network

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

# Variational Autoencoders (VAE)

➤ How to train

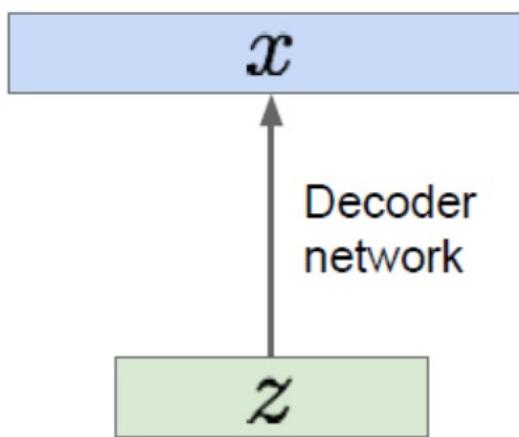
We want to estimate the true parameters  $\theta^*$  of this generative model given training data  $x$ .

Sample from  
true conditional

$$p_{\theta^*}(x \mid z^{(i)})$$

Sample from  
true prior

$$z^{(i)} \sim p_{\theta^*}(z)$$



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 (VAE)

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



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

$$\log p(x) \approx \log \frac{1}{k} \sum_{i=1}^k p(x|z^{(i)}), \text{ where } z^{(i)} \sim p(z)$$

Monte Carlo estimation is too high variance

# Variational Autoencoders (VAE)

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 modeling  $p_{\theta}(x|z)$ , learn  $q_{\phi}(z|x)$  that approximates the true posterior  $p_{\theta}(z|x)$ .

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

**Variational inference** is to approximate the unknown posterior distribution from only the observed data  $x$

# Variational Autoencoders (VAE)

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\&= \mathbf{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\&= \mathbf{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\&= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbf{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\&= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))\end{aligned}$$

The expectation wrt.  $z$  (using encoder network) let us write nice KL terms

# Variational Autoencoders (VAE)

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbf{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))\end{aligned}$$



Decoder network gives  $p_{\theta}(x|z)$ , can compute estimate of this term through sampling (need some trick to differentiate through sampling).



This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!



$p_{\theta}(z|x)$  intractable (saw earlier), can't compute this KL term :( But we know KL divergence always  $\geq 0$ .

# Variational Autoencoders (VAE)

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule})$$

Decoder:  
reconstruct  
the input data

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}$$

Encoder:  
make approximate  
posterior distribution  
close to prior

Tractable lower bound which we can take  
gradient of and optimize! ( $p_\theta(x|z)$  differentiable,  
KL term differentiable)

# Variational Autoencoders (VAE)

➤ Training objective

Putting it all together: maximizing the likelihood lower bound

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

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

$$D_{KL}(\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I))$$

Have analytical solution

Make approximate posterior distribution close to prior

Encoder network

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

Input Data

$$\mu_{z|x}$$

$$\Sigma_{z|x}$$

$$x$$

# Variational Autoencoders (VAE)

➤ Training objective

Putting it all together: maximizing the likelihood lower bound

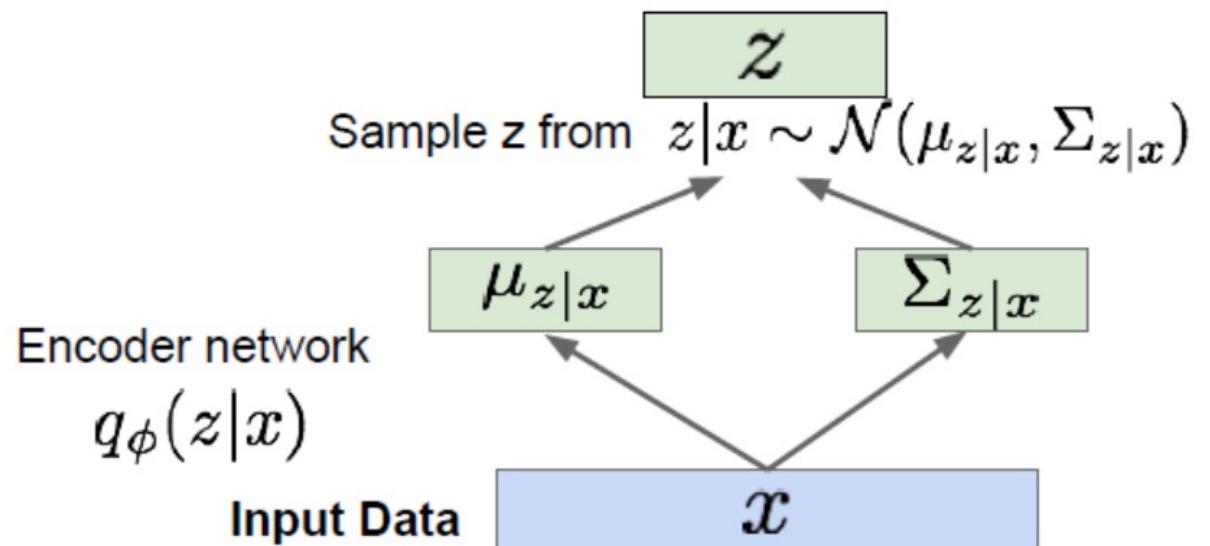
$$\underbrace{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Sample  $\epsilon \sim \mathcal{N}(0, I)$

$z = \mu_{z|x} + \epsilon \sigma_{z|x}$

Part of computation graph

Reparametrization trick



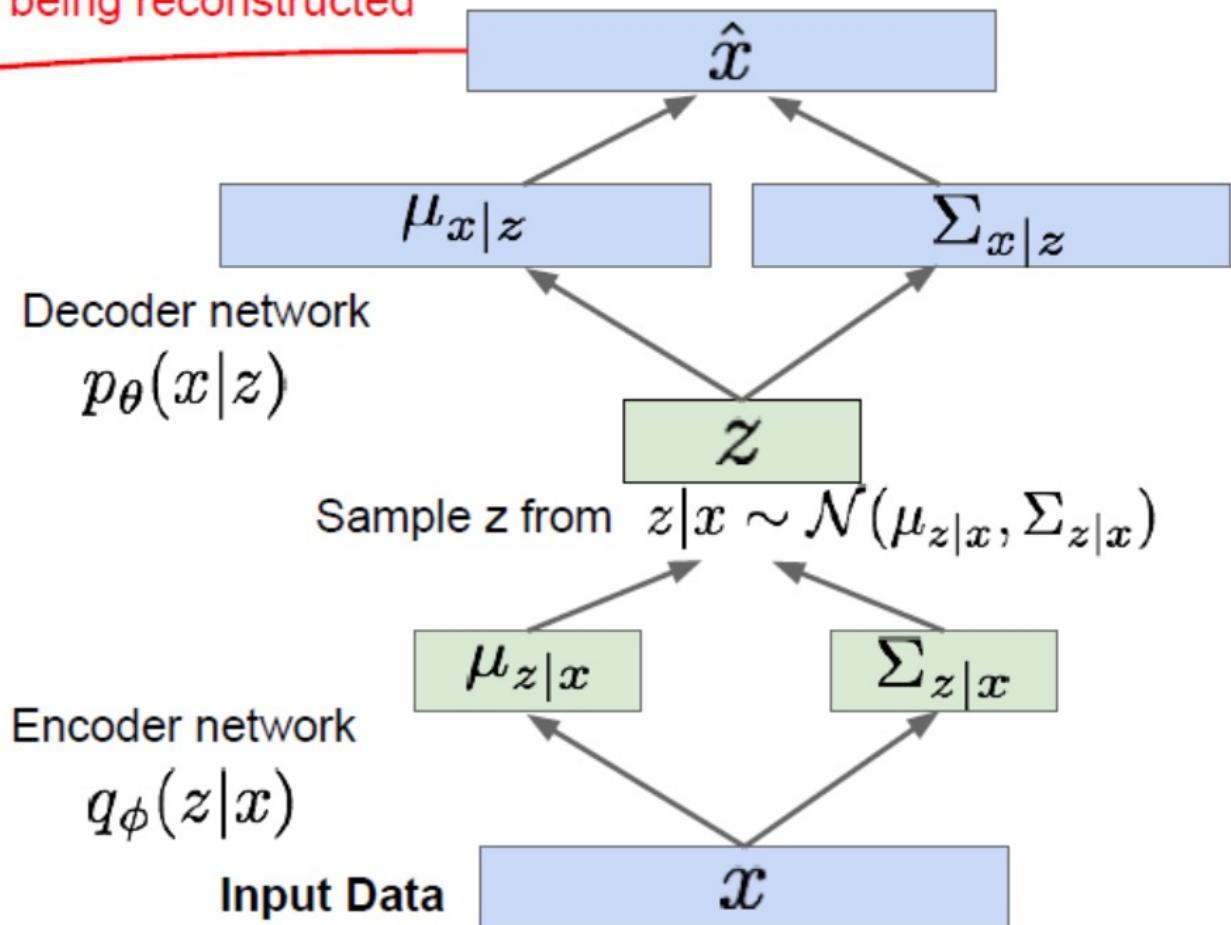
# Variational Autoencoders (VAE)

➤ Training objective

Putting it all together: maximizing the likelihood lower bound

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

Maximize likelihood of original input being reconstructed



# Variational Autoencoders (VAE)

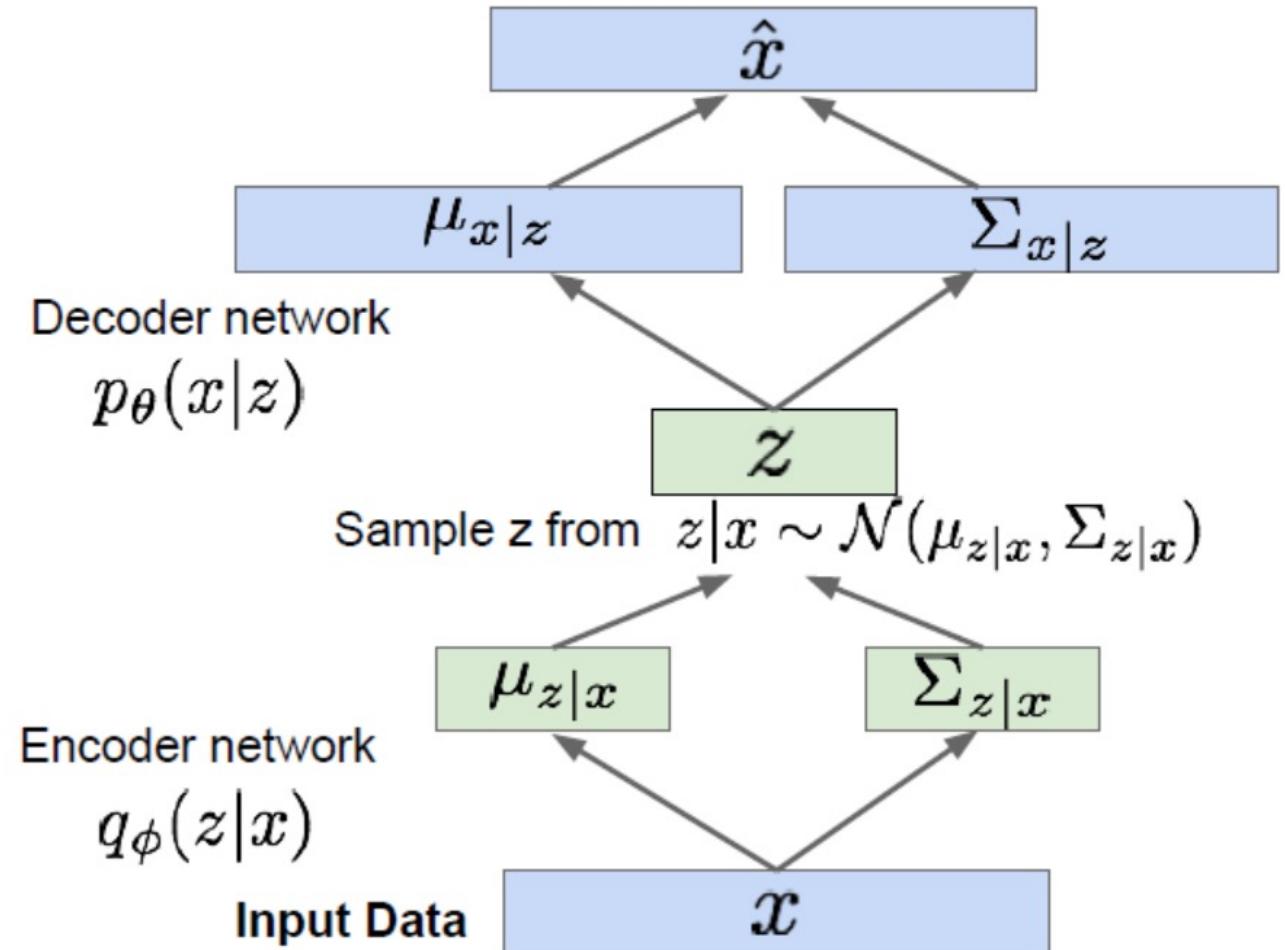
➤ Training objective

Putting it all together: maximizing the likelihood lower bound

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

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

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

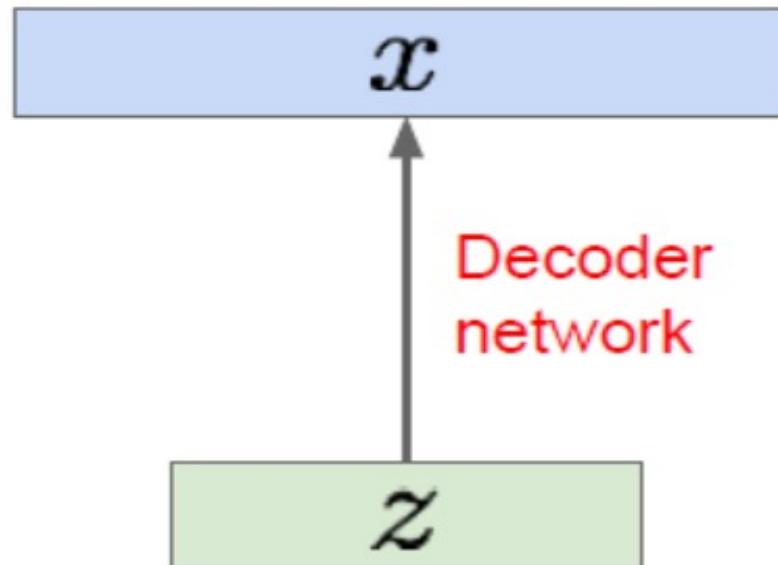


# Variational Autoencoders (VAE)

Our assumption about data generation process

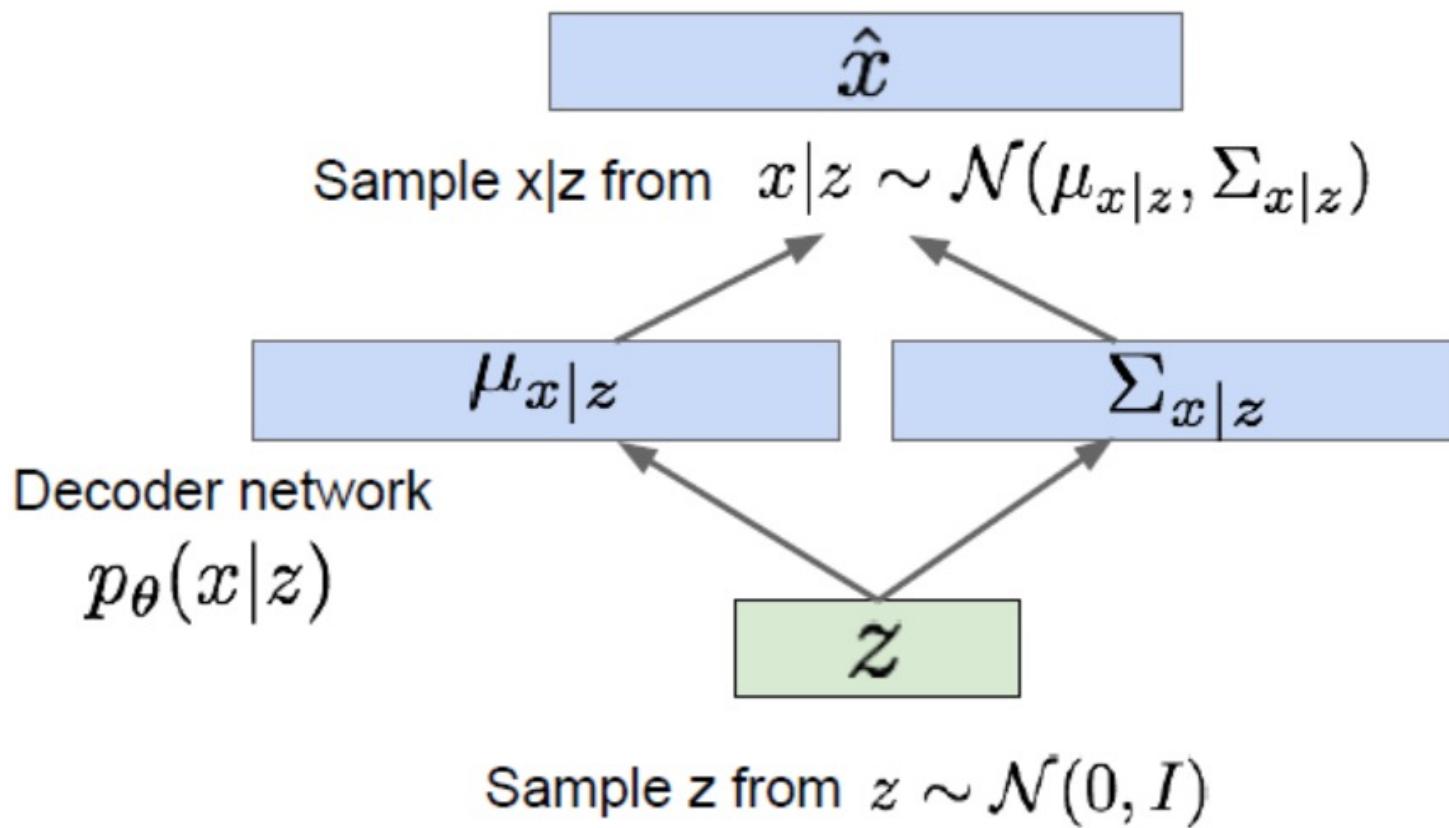
Sample from true conditional  
 $p_{\theta^*}(x \mid z^{(i)})$

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



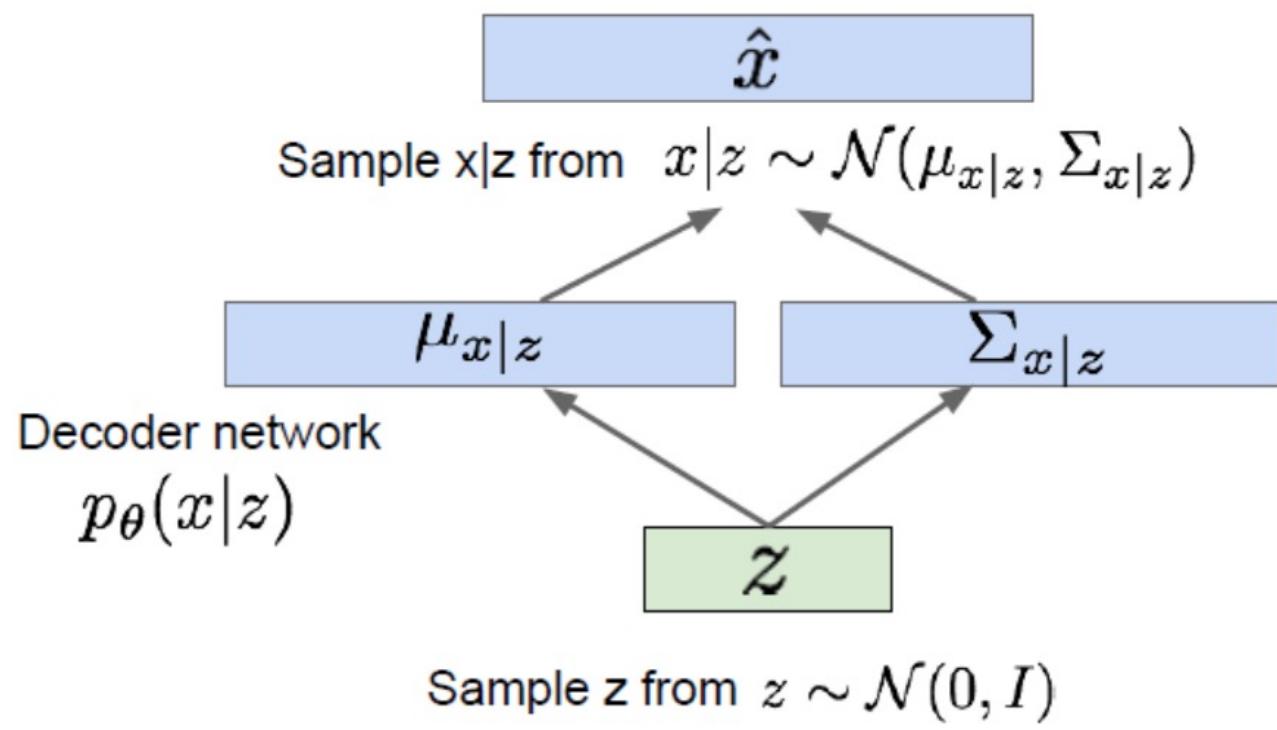
# Variational Autoencoders (VAE)

Now given a trained VAE:  
use decoder network & sample z from prior!

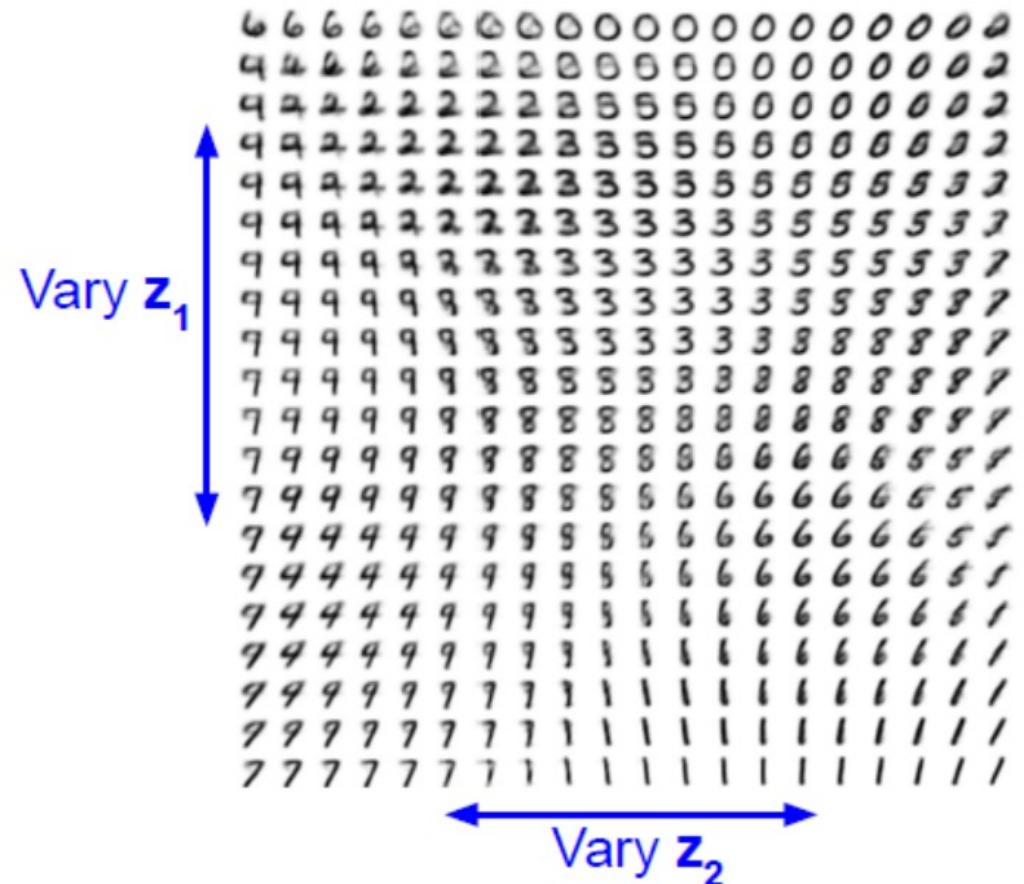


# Variational Autoencoders (VAE)

Use decoder network. Now sample z from prior!



Data manifold for 2-d  $z$



# Variational Autoencoders (VAE)



32x32 CIFAR-10



Labeled Faces in the Wild

谢谢！

