

# Generative Models and GANs

1

Yuan YAO

HKUST

Based on Feifei Li, Aleksander Mqdry, Tong Zhang, et al.



# Overview

- ▶ Unsupervised Learning vs. Supervised Learning
- ▶ Generative Models
  - ▶ Variational Autoencoders (VAE)
  - ▶ Generative Adversarial Networks (GAN)

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



→ Cat

Classification

This image is CC0 publicdomain

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



DOG, DOG, CAT

Object Detection

This image is CC0 publicdomain

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



GRASS, CAT,  
TREE, SKY

Semantic Segmentation

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



A cat sitting on a suitcase on the floor

Image captioning

Caption generated using [neuraltalk2](#)  
[Image](#) is CC0 Public domain.

# Supervised vs Unsupervised Learning

## 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 Learning

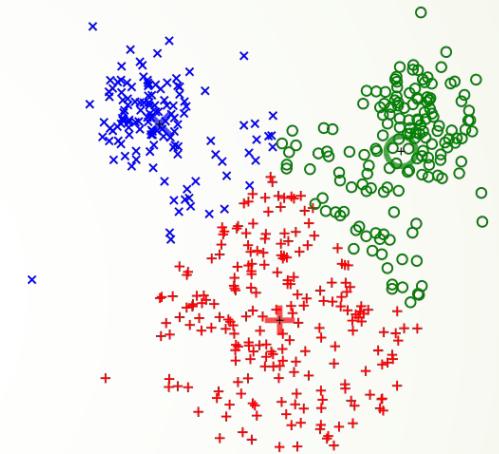
## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



K-means clustering

This image is [CC0 publicdomain](#)

# Supervised vs Unsupervised Learning

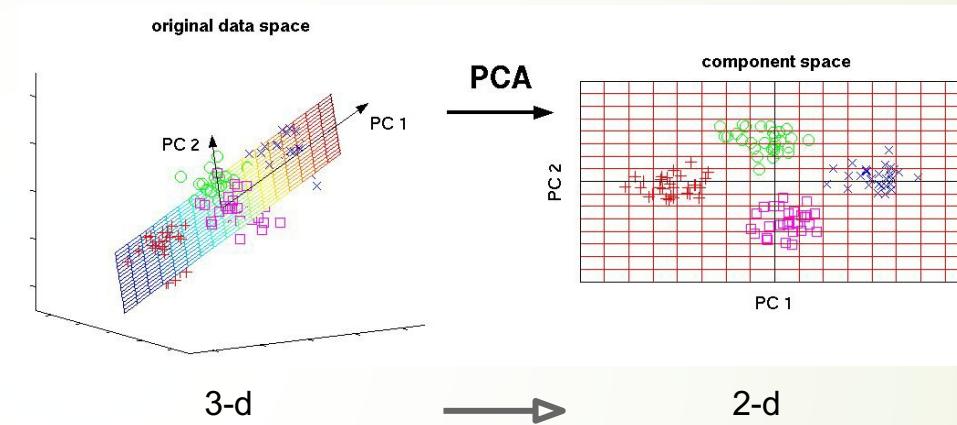
## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



Principal Component Analysis,  
Multidimensional Scaling,  
and Manifold Learning  
(Dimensionality reduction)

This image from Matthias Scholz  
is CC0 publicdomain

# Supervised vs Unsupervised Learning

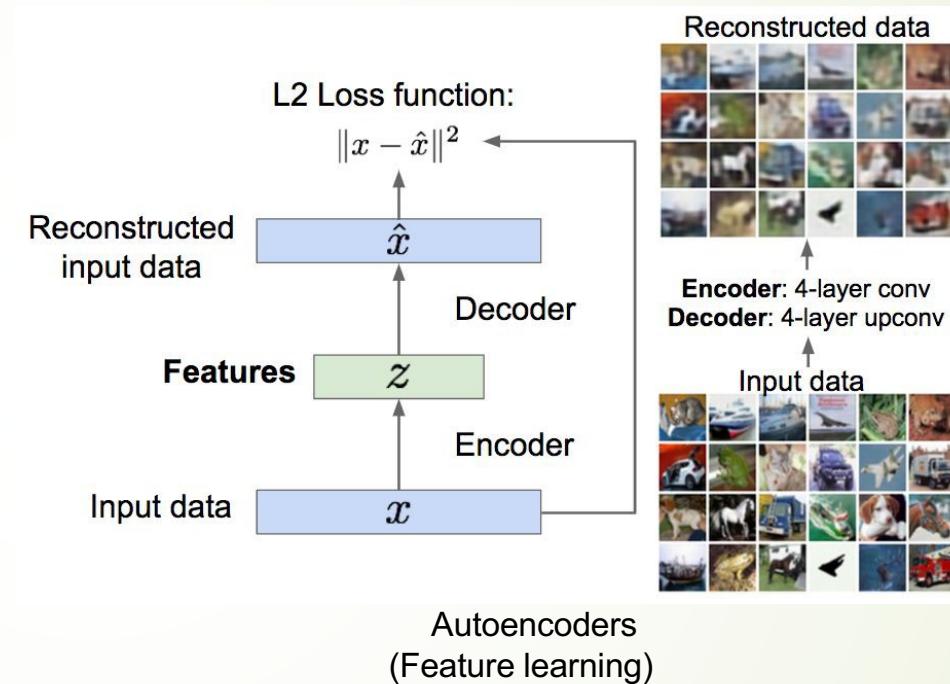
## Unsupervised Learning

**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 Learning

## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

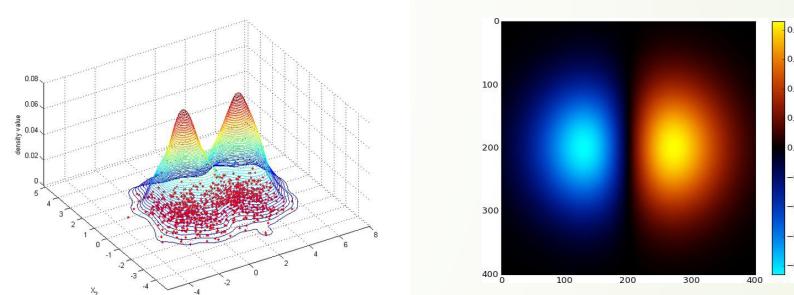
**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



2-d density estimation

2-d density images [left](#) and [right](#) are [CC0 publicdomain](#)

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.

## 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.

Lecture 13 -

# Generative Models

Given training data, generate new samples from same distribution



Training data  $\sim p_{\text{data}}(x)$

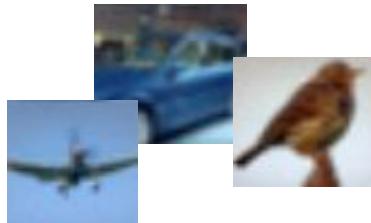


Generated samples  $\sim p_{\text{model}}(x)$

Want to learn  $p_{\text{model}}(x)$  similar to  $p_{\text{data}}(x)$

# Generative Models

Given training data, generate new samples from same distribution



Training data  $\sim p_{\text{data}}(x)$



Generated samples  $\sim p_{\text{model}}(x)$

Want to learn  $p_{\text{model}}(x)$  similar to  $p_{\text{data}}(x)$

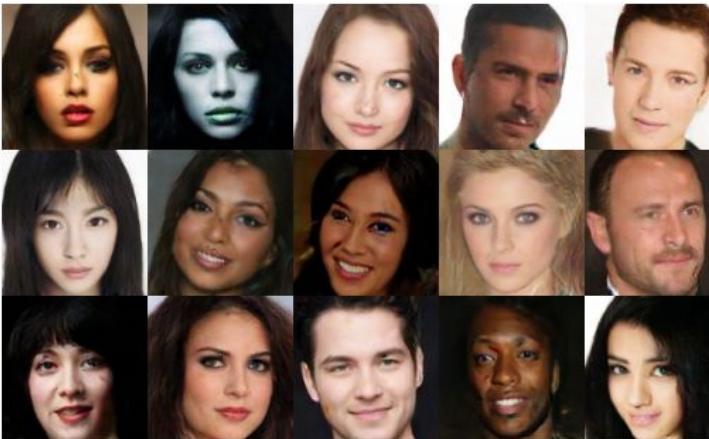
Addresses density estimation, a core problem in unsupervised learning

## Several flavors:

- Explicit density estimation: explicitly define and solve for  $p_{\text{model}}(x)$
- Implicit density estimation: learn model that can sample from  $p_{\text{model}}(x)$  w/o explicitly defining it

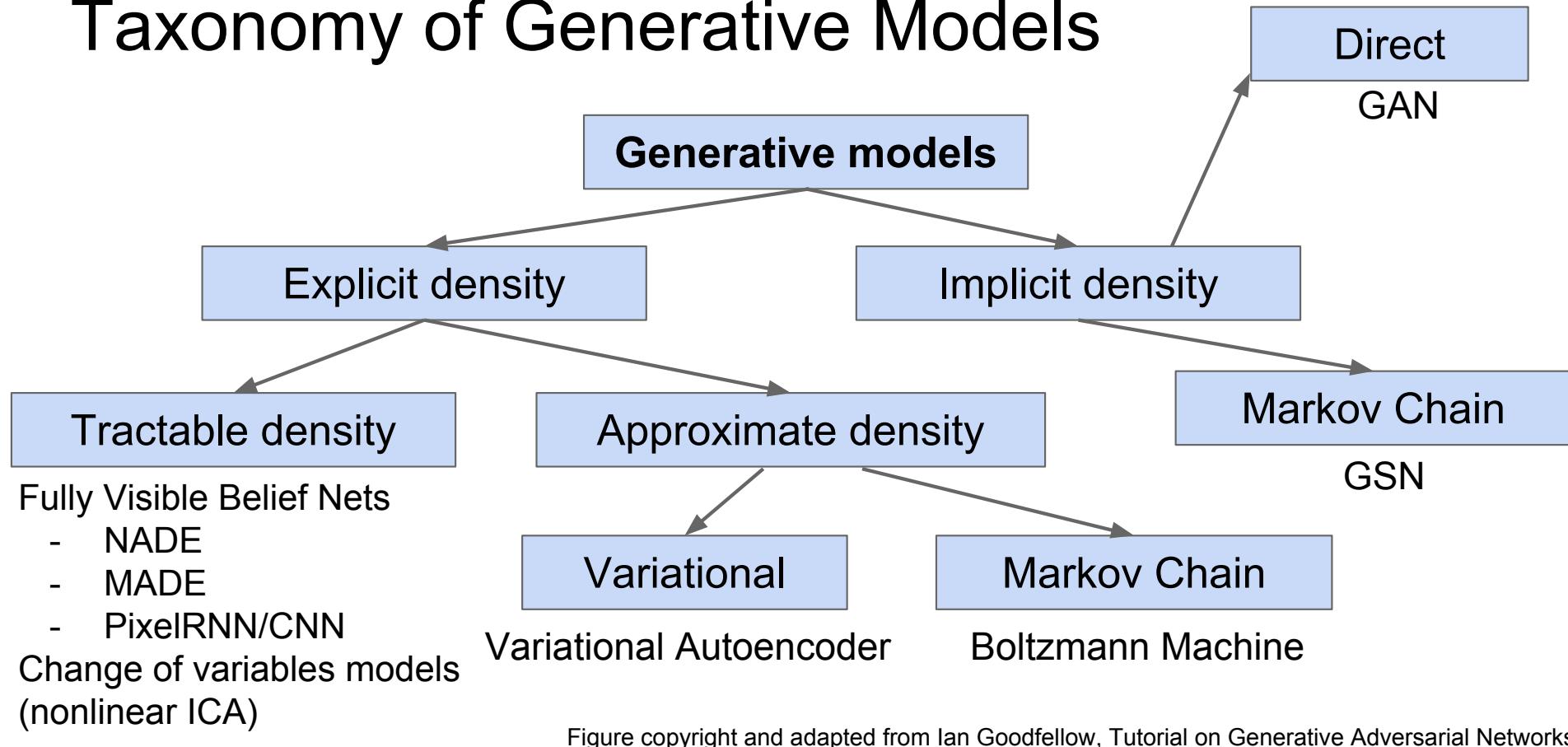
# Why Generative Models?

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for simulation and planning (reinforcement learning applications!)
- Training generative models can also enable inference of latent representations that can be useful as general features

# Taxonomy of Generative Models



# Taxonomy of Generative Models

Today: discuss 3 most popular types of generative models today

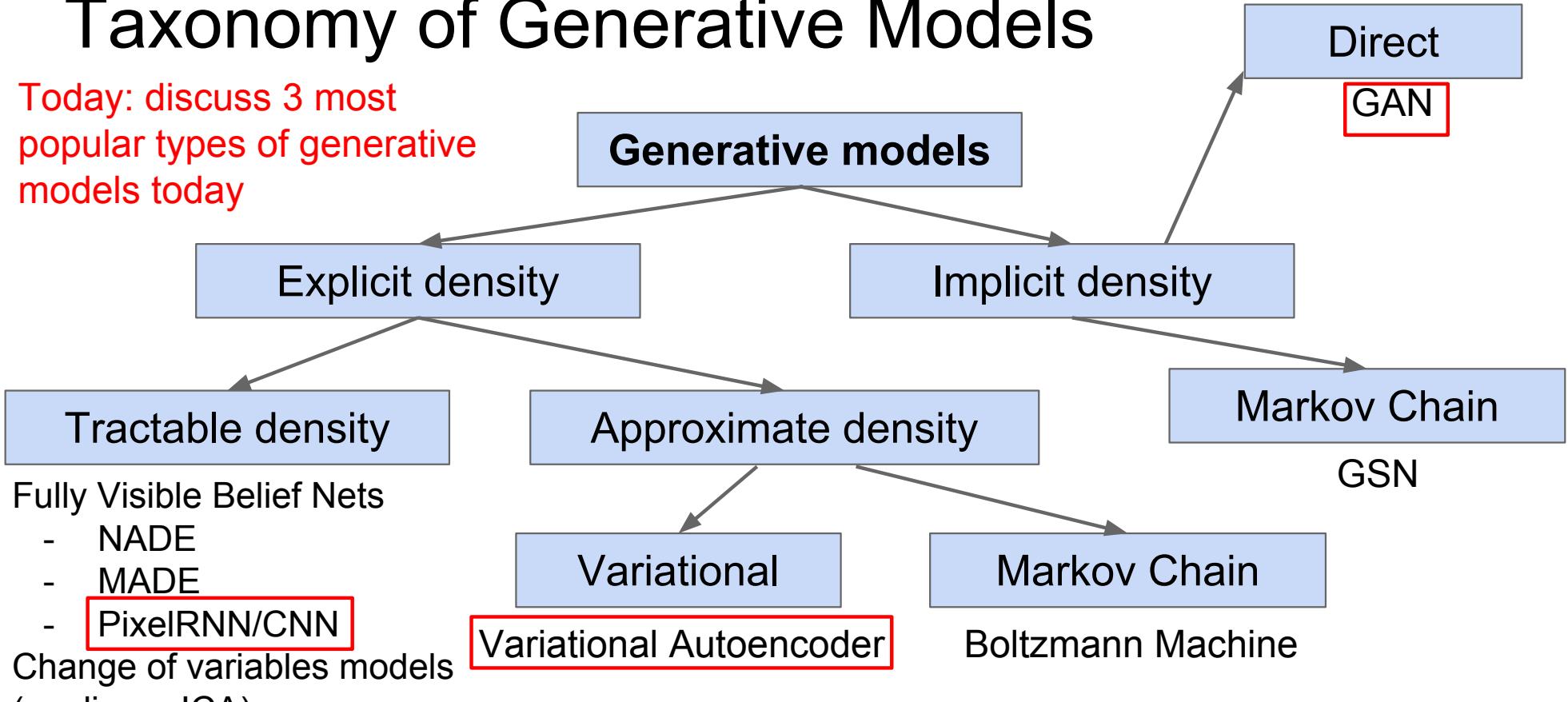


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.



# PixelRNN and Pixel CNN

# Fully visible belief network

# Explicit density model

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

Then maximize likelihood of training data

# Fully visible belief network

# Explicit density model

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

Then maximize likelihood of training data

Will need to define ordering of “previous pixels”

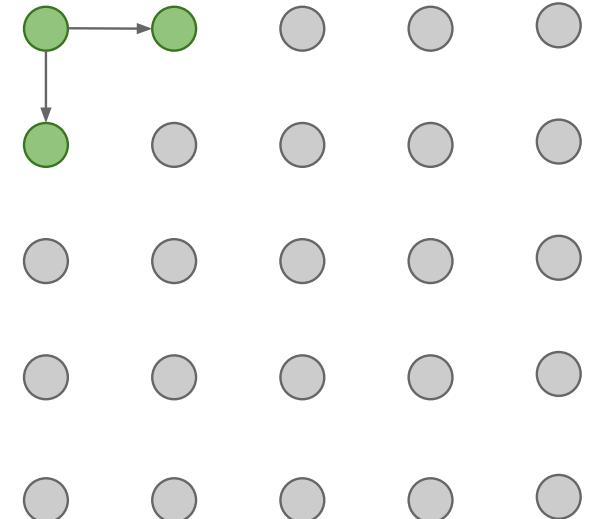
Complex distribution over pixel values => Express using a neural network!

# PixelRNN

*[van der Oord et al. 2016]*

Generate image pixels starting from corner

Dependency on previous pixels modeled  
using an RNN (LSTM)

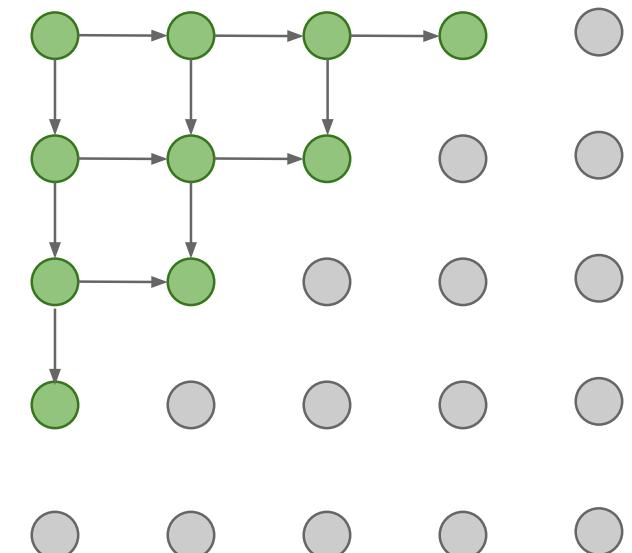


# PixelRNN *[van der Oord et al. 2016]*

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Drawback: sequential generation is slow!



# 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

Training: maximize likelihood of training images

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

Softmax loss at each pixel

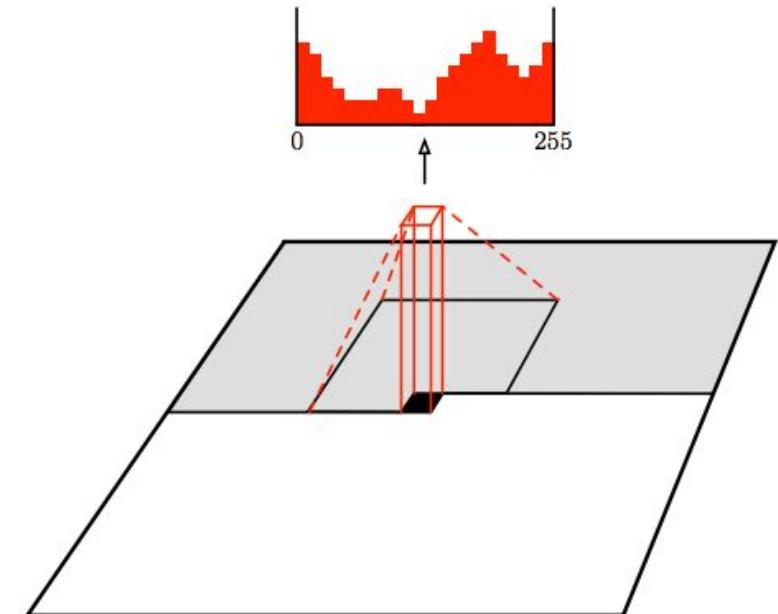


Figure copyright van der Oord et al., 2016. Reproduced with permission.

# 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

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

Generation must still proceed sequentially  
=> still slow

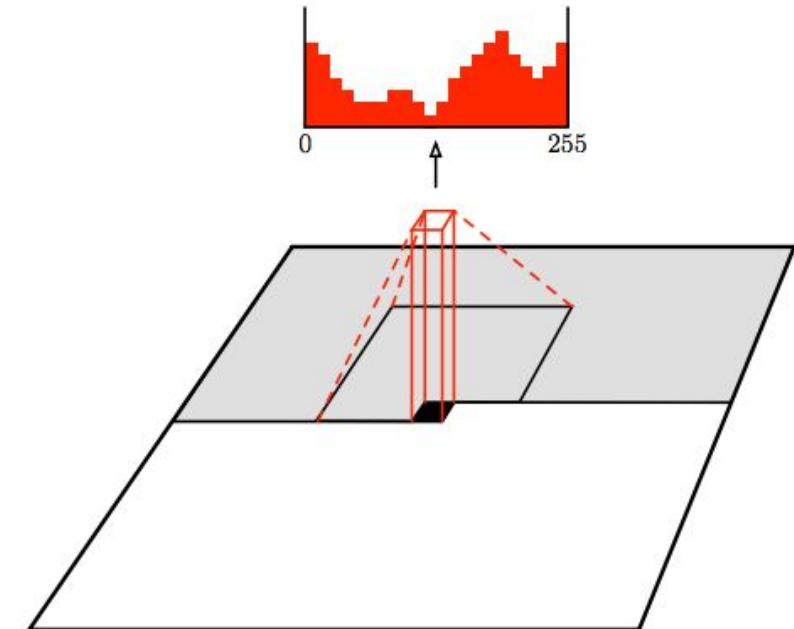
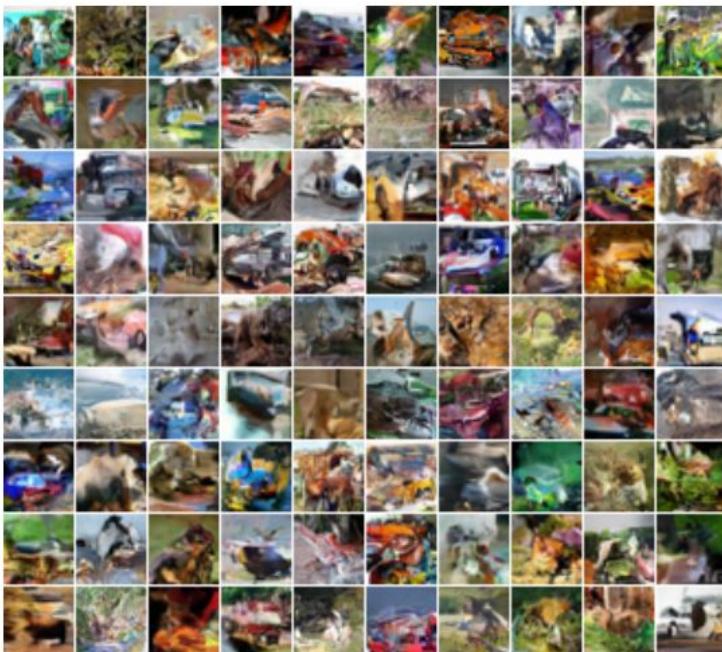
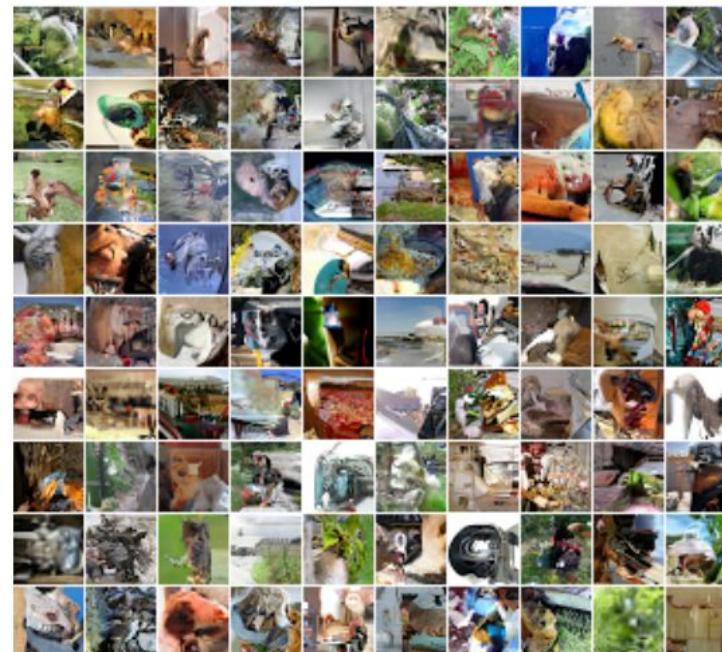


Figure copyright van der Oord et al., 2016. Reproduced with permission.

# Generation Samples



32x32 CIFAR-10



32x32 ImageNet

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

# PixelRNN and PixelCNN

- ▶ Pros:
  - ▶ Can explicitly compute likelihood  $p(x)$
  - ▶ Explicit likelihood of training: data gives good evaluation metric
  - ▶ Good samples
- ▶ Con:
  - ▶ Sequential generation => slow
- ▶ Improving PixelCNN performance
  - ▶ Gated convolutional layers
  - ▶ Short-cut connections
  - ▶ Discretized logistic loss
  - ▶ Multi-scale
  - ▶ Training tricks
  - ▶ Etc...
- ▶ See
  - ▶ Van der Oord et al. NIPS 2016
  - ▶ Salimans et al. 2017 (PixelCNN++)



Variational Autoencoders (VAE)

# From Likelihood to Mixture Likelihood

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

# Pearson's Crabs

- ▶ Karl Pearson, *Contributions to the Mathematical Theory of Evolution*, 1894
- ▶ Forehead width/Body Length measurements in population of crabs from Malta [collected by zoologist [Weldon and his wife](#)]
- ▶ Peculiar deviation from normal distribution ?!?
- ▶ Pearson writes:
  - ▶ "In the case of certain biological, sociological, and economic measurements there is, however, a well-marked deviation from this normal shape, and it becomes important to determine the direction and amount of such deviation. The asymmetry may arise from the fact that the units grouped together in the measured material are not really homogeneous. It may happen that we have a mixture of 2,3,...,n homogeneous groups, each of which deviates about its own mean symmetrically and in a manner represented with sufficient accuracy by the normal curve."
- ▶ Stipulates mixture of two Gaussians!
  - ▶ Identifies components by solving degree 9 polynomial in the first 5 moments of the distribution (uses 6<sup>th</sup> moment to filter roots)



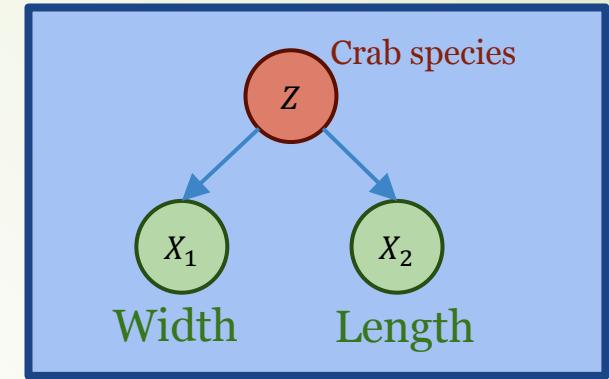
# Pearson's Crabs

- ▶ Analytical difficulties encountered by Pearson signature of partially observable models
- ▶ Here a **hidden/latent** variable  $Z$  influences  $\vec{X}$ , the forehead width and body length of the crabs
- ▶ Maximum likelihood estimation in *latent variable models*:
  - ▶ suppose  $(\vec{x}_1, \vec{z}_1), \dots, (\vec{x}_N, \vec{z}_N) \sim_{i.i.d.} p_{\vec{\theta}}(\vec{x}, \vec{z})$
  - ▶ **goal:** recover  $\vec{\theta}$ , given only  $\vec{x}_1, \dots, \vec{x}_N$

$$\text{MLE} = \arg \max_{\vec{\theta}} \sum_i \log \left( \sum_{\vec{z}} p_{\vec{\theta}}(\vec{x}_i, \vec{z}) \right)$$

samples
 $p_{\vec{\theta}}(\vec{x}_i)$

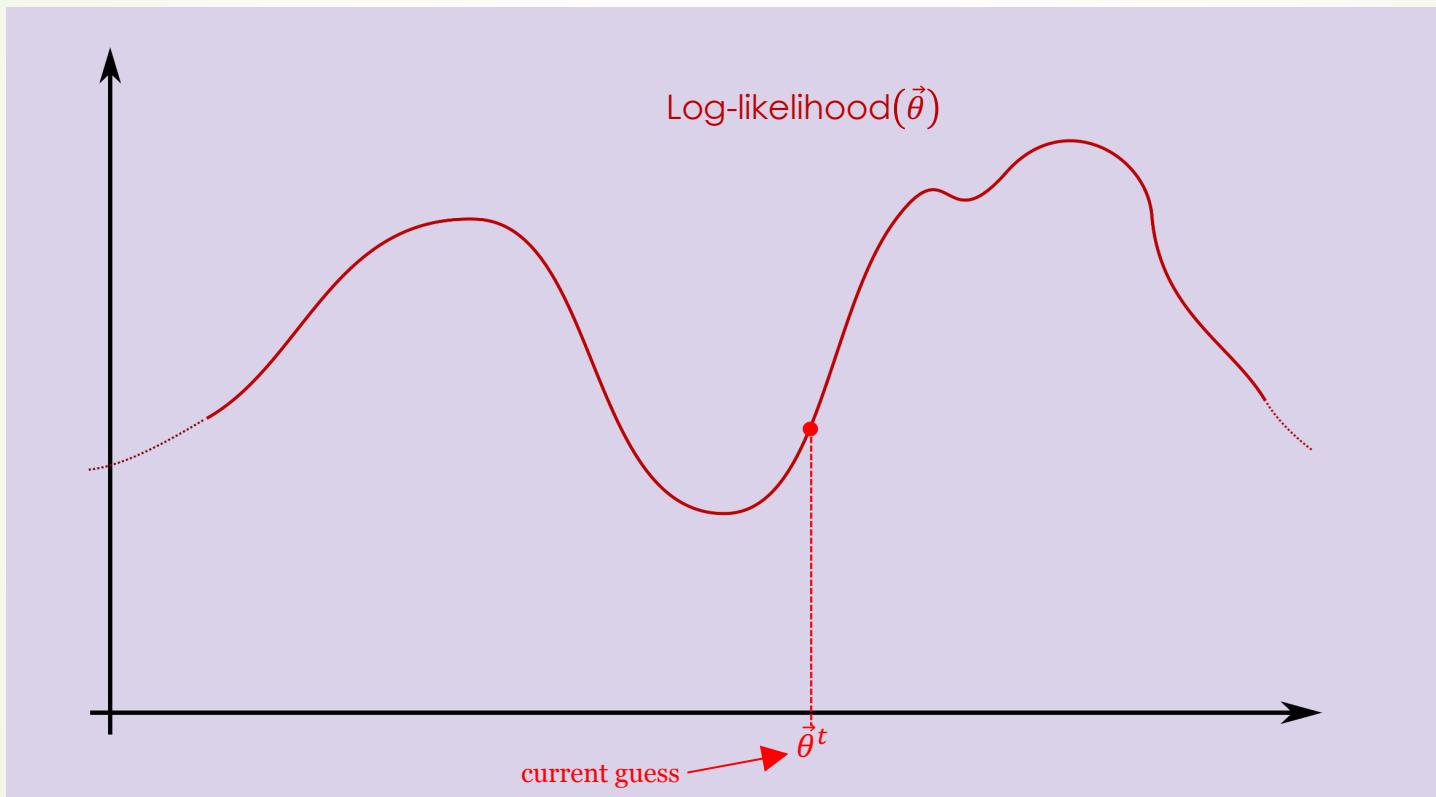
- ▶ Non-convex, **intractable** so cannot optimize directly
  - ▶ Expectation-Maximization algorithm, variational inference optimize lower bound of MLE



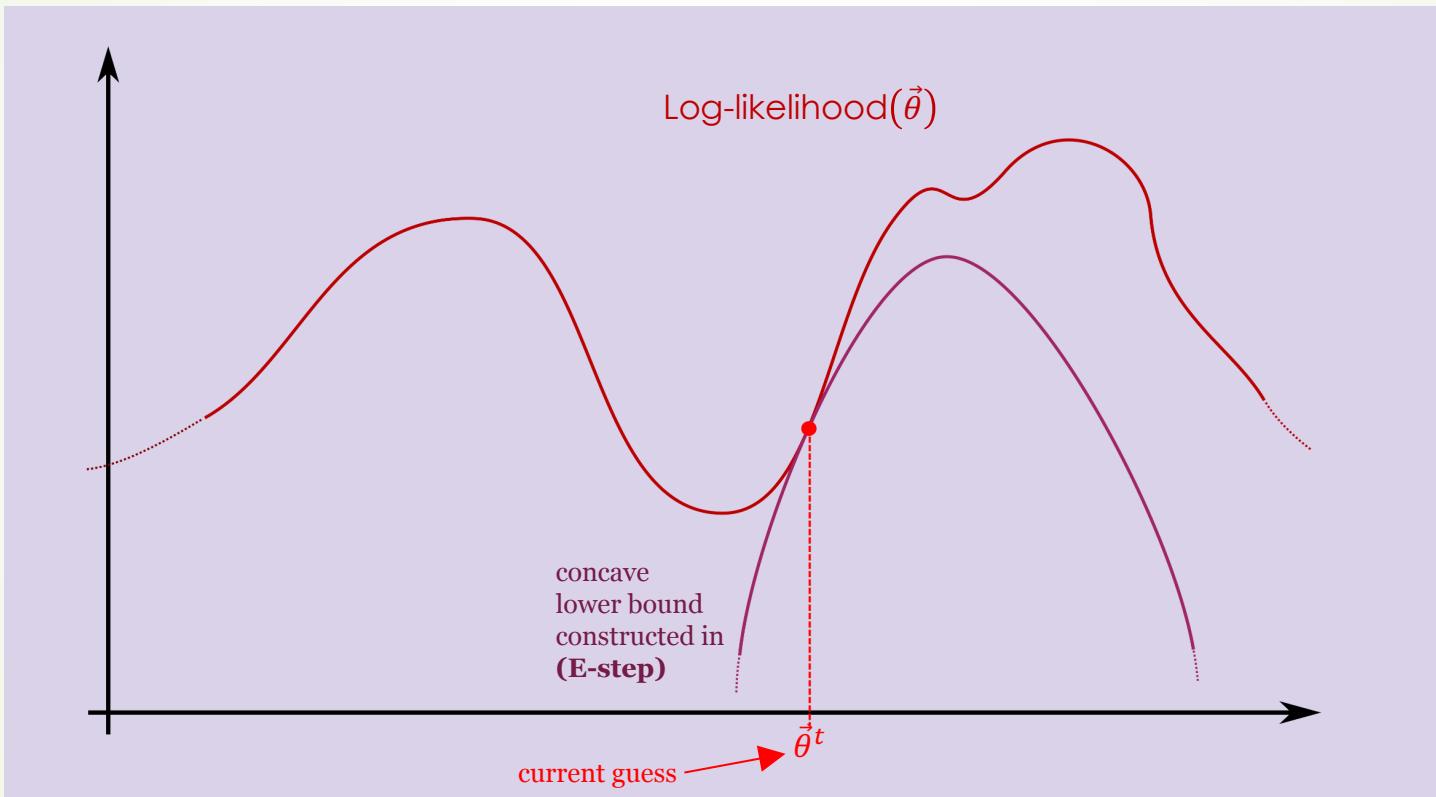
# E-M in pictures



# E-M in pictures

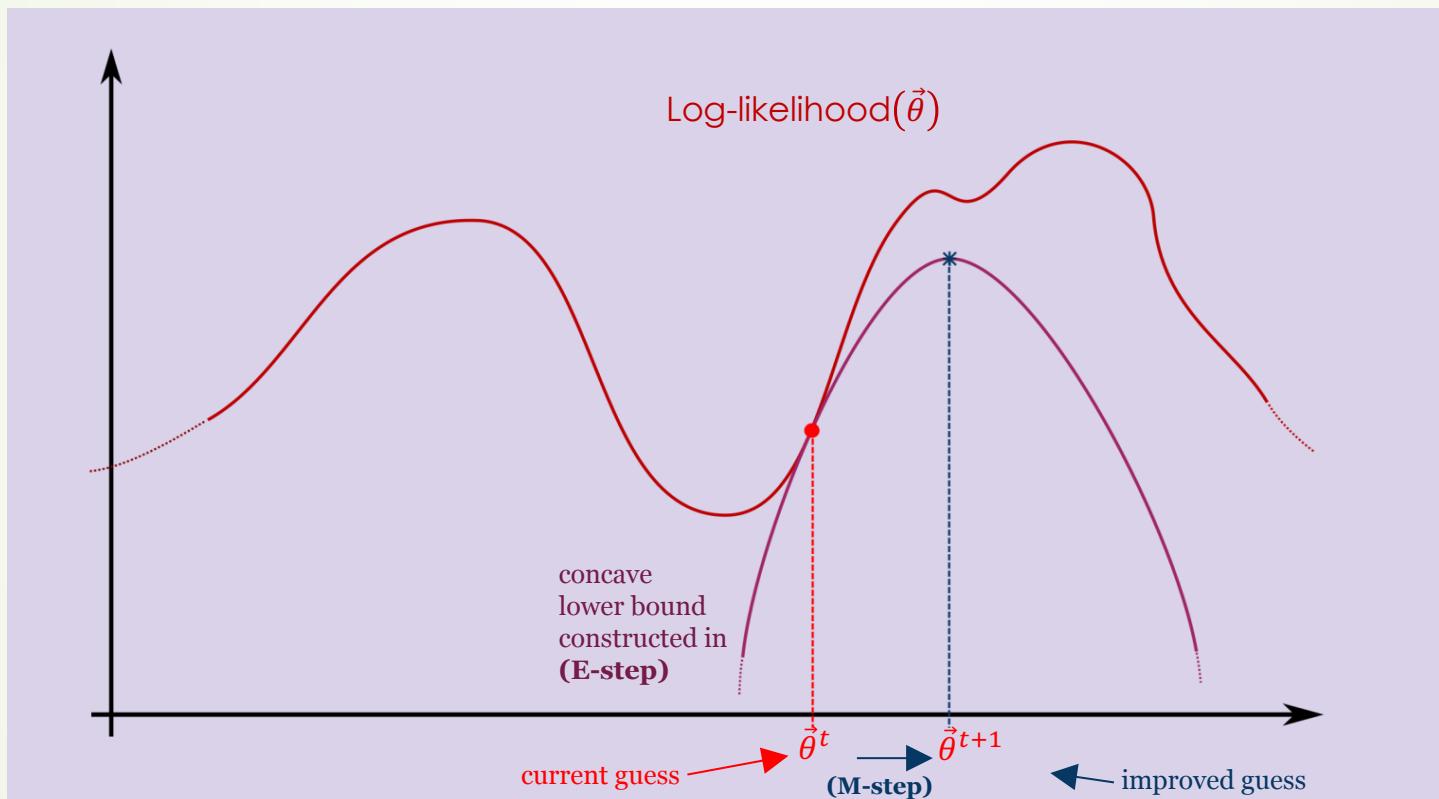


# E-M in pictures



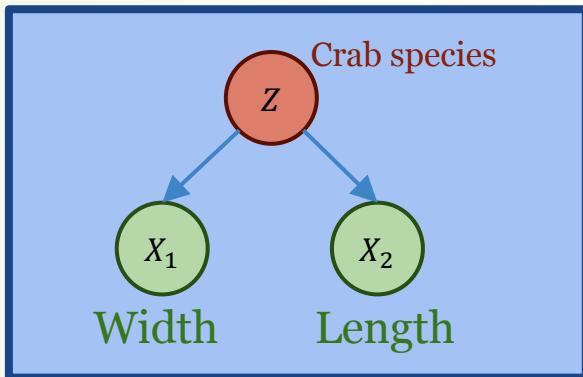
**E-step:** a beautiful example of a variational lower bound that is (i) tight at  $\vec{\theta}^t$ ; and (ii) concave, hence easy to optimize

# E-M in pictures



# So far...

Pearson's crabs: hidden variable was low-dimensional (in fact binary)



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

Often discrete and low-dimensional hidden variable  $Z$  insufficient

Suppose both  $Z$  and  $X$  are continuous and high-dimensional, then

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

Maximum likelihood **very intractable**

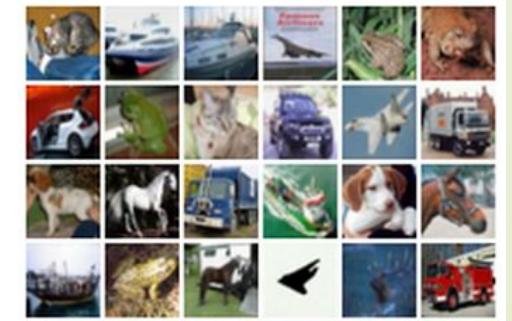
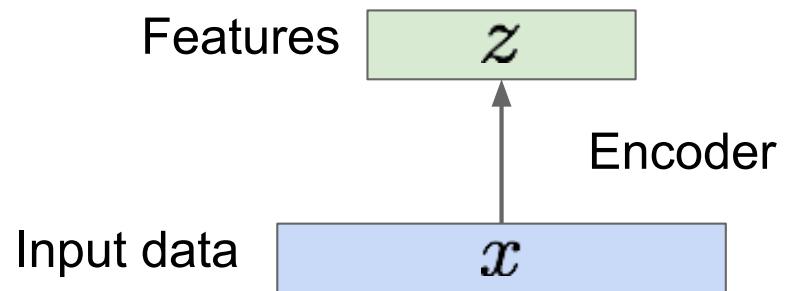
cannot optimize directly, not even analytically choose a good lower bound to it

**VAEs:** canonical way to define density function with high-dimensional latent  $Z$ , optimize lower bound of likelihood

# Some background first: Autoencoders

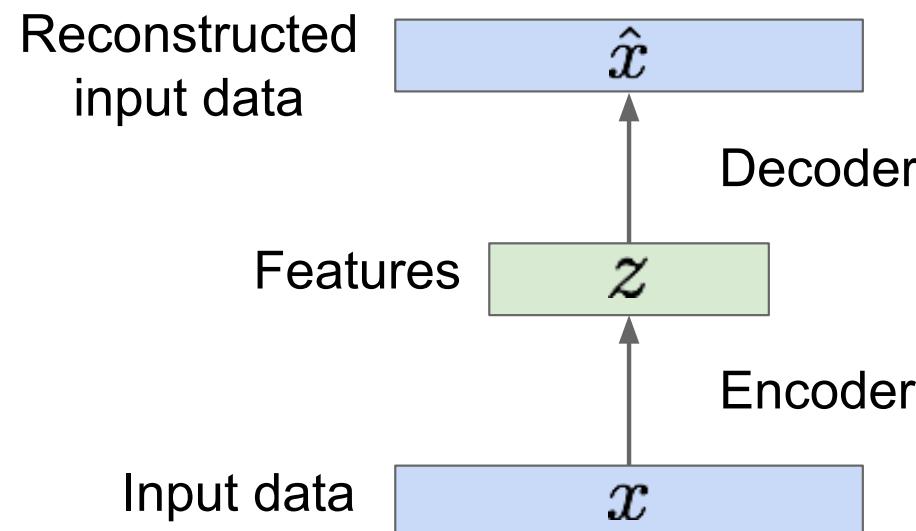
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

e.g. PCA, Manifold Learning, Dictionary Learning

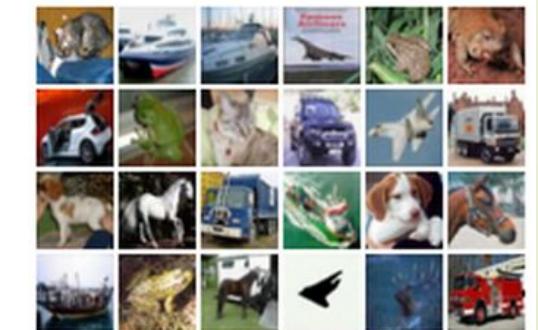


## How to learn this feature representation?

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



e.g. PCA, Manifold Learning,  
Dictionary Learning, Matrix  
Factorization:  $D = E'$



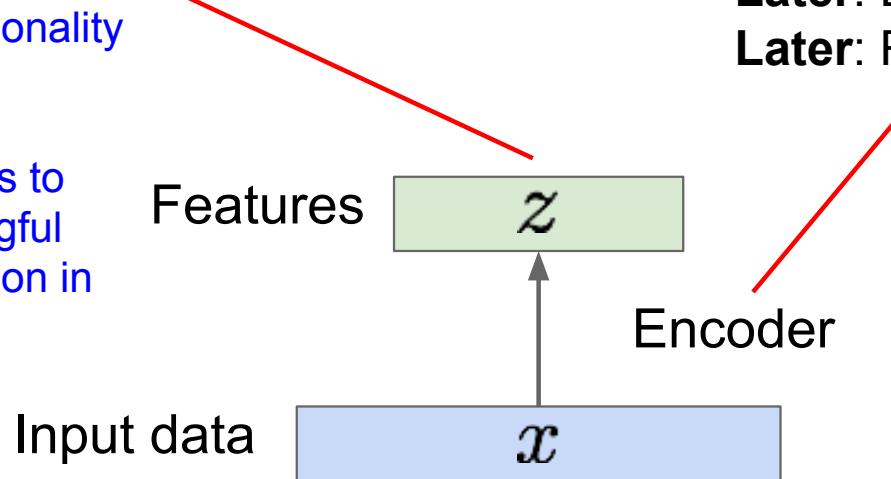
# Deep Learning for encoders

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

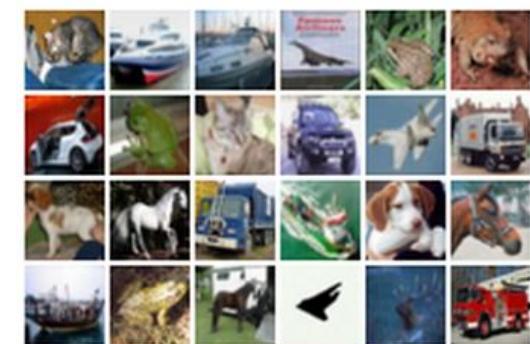
**z** usually smaller than **x**  
(dimensionality reduction)

## Q: Why dimensionality reduction?

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



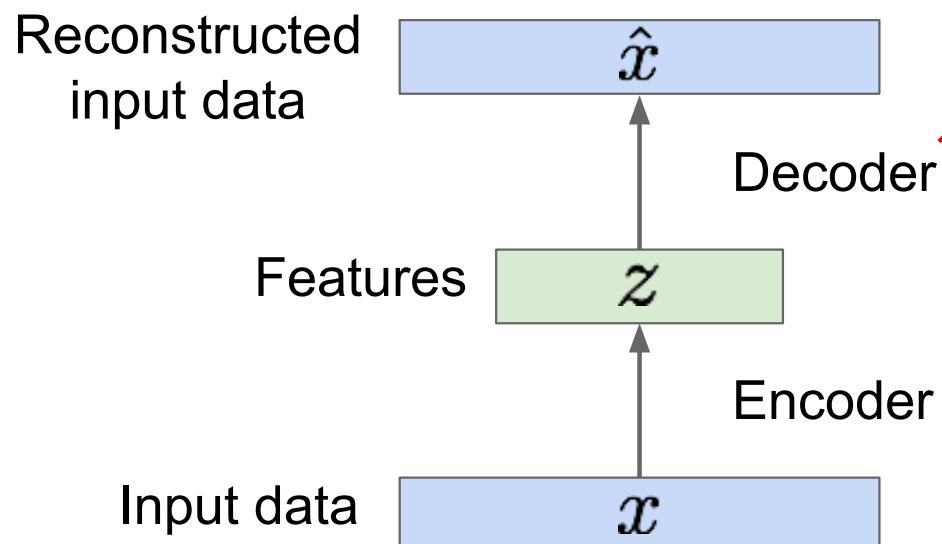
**Originally:** Linear +  
nonlinearity (sigmoid)  
**Later:** Deep, fully-connected  
**Later:** ReLU CNN



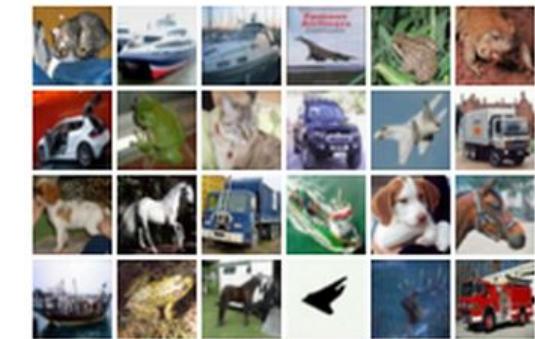
# Deep Learning for decoders

How to learn this feature representation?

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



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

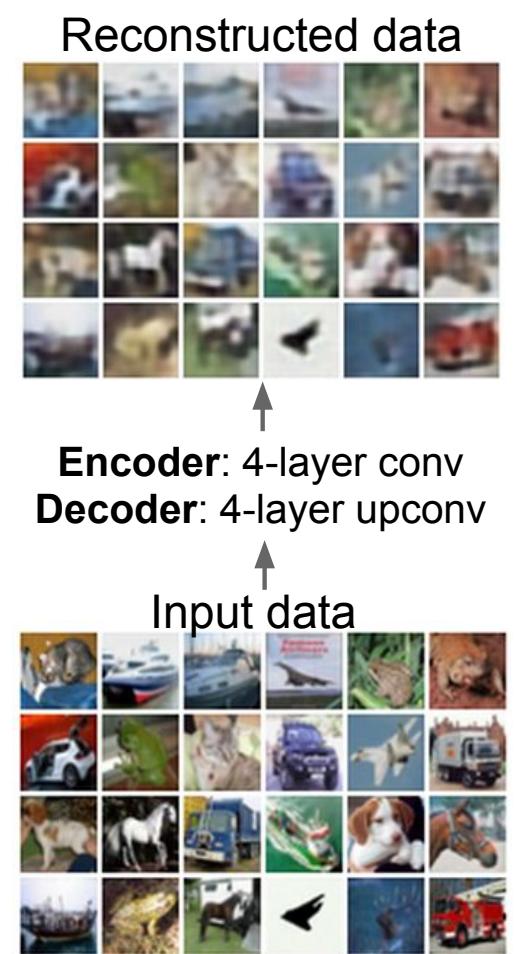
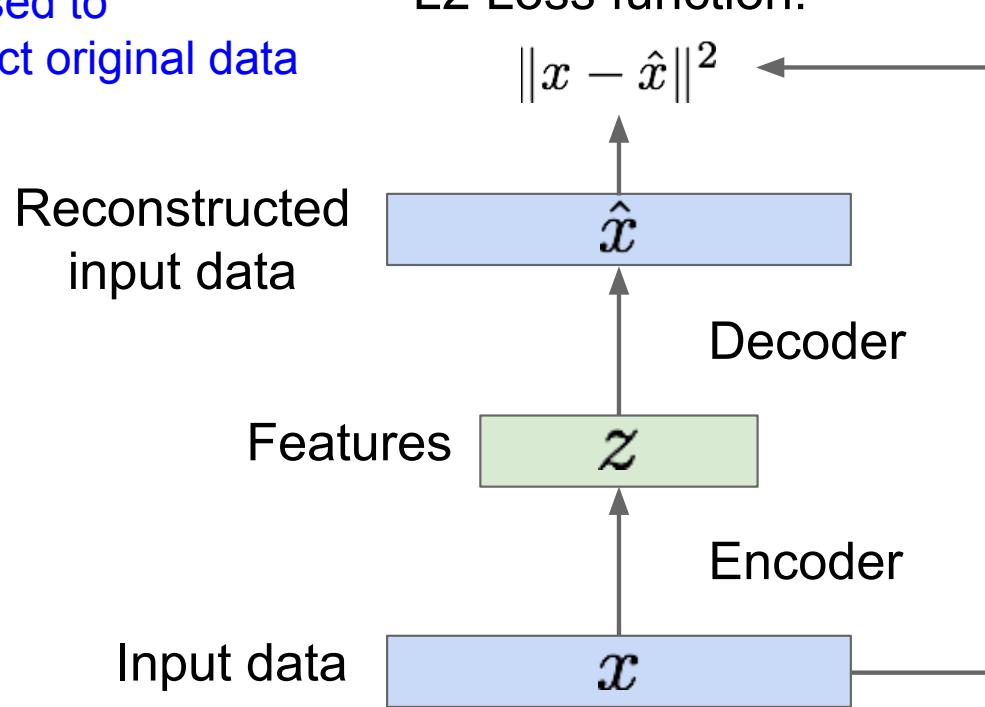


# L2 Loss functions

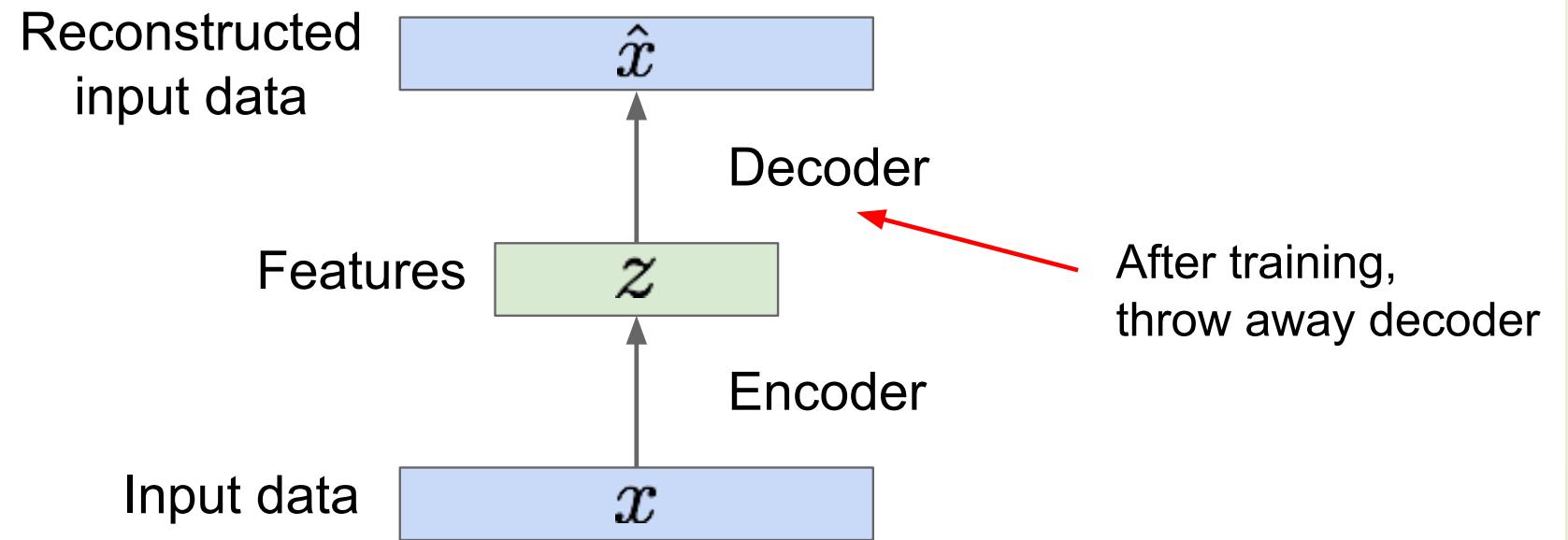
## Some background first: Autoencoders

Train such that features can be used to reconstruct original data

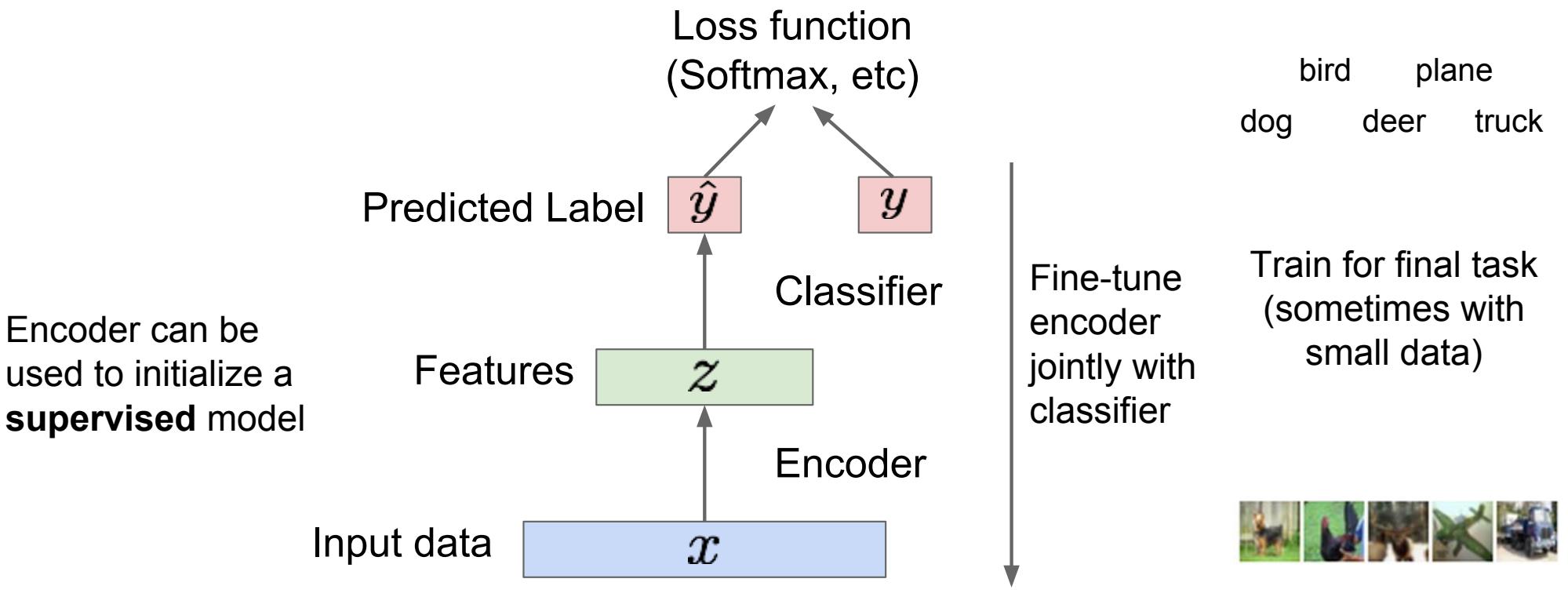
Doesn't use labels!



# Some background first: Autoencoders

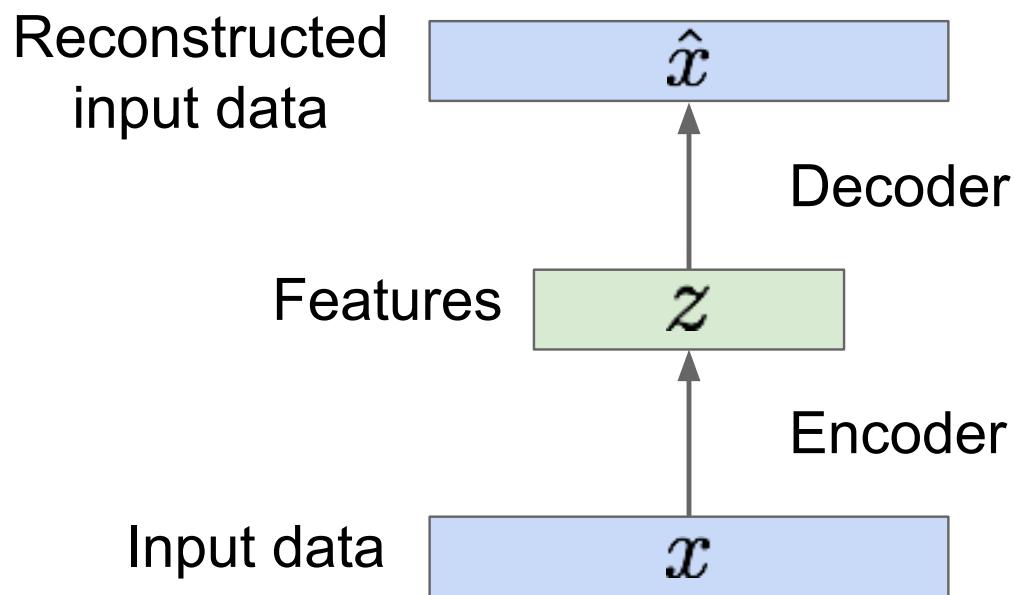


# Autoencoders for Transfer Learning





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?

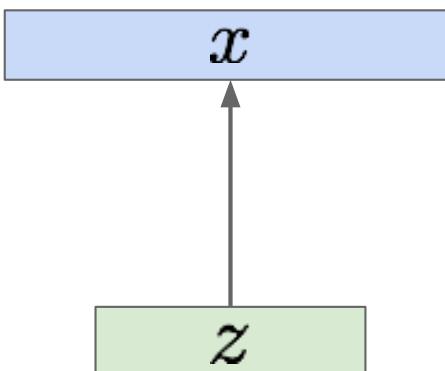
# Variational Autoencoders

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

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

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

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



**Intuition** (remember from autoencoders!):  
 $x$  is an image,  $z$  is latent factors used to  
generate  $x$ : attributes, orientation, etc.

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

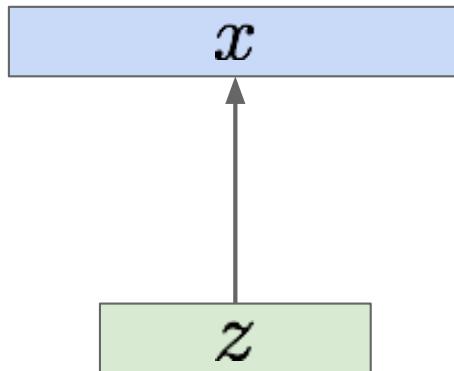
# Variational Autoencoders

Sample from  
true conditional

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

Sample from  
true prior

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



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

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.

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

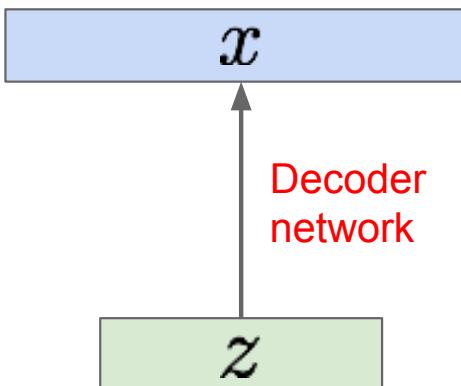
# Variational Autoencoders

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 should we represent this model?

Choose prior  $p(z)$  to be simple, e.g.  
Gaussian.

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

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

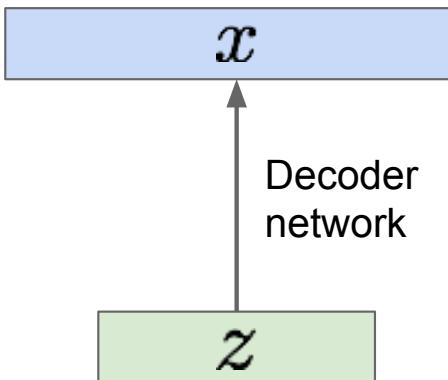
# Variational Autoencoders

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?

Remember strategy for training generative models from FVBMs. Learn model parameters to maximize likelihood of training data

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

Now with latent  $z$

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

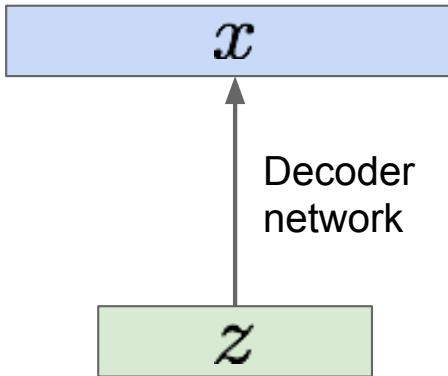
# Variational Autoencoders

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?**

Remember strategy for training generative models from FVBMs. 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: Intractability

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

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

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

Intractable data likelihood

# Variational Lower Bounds

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

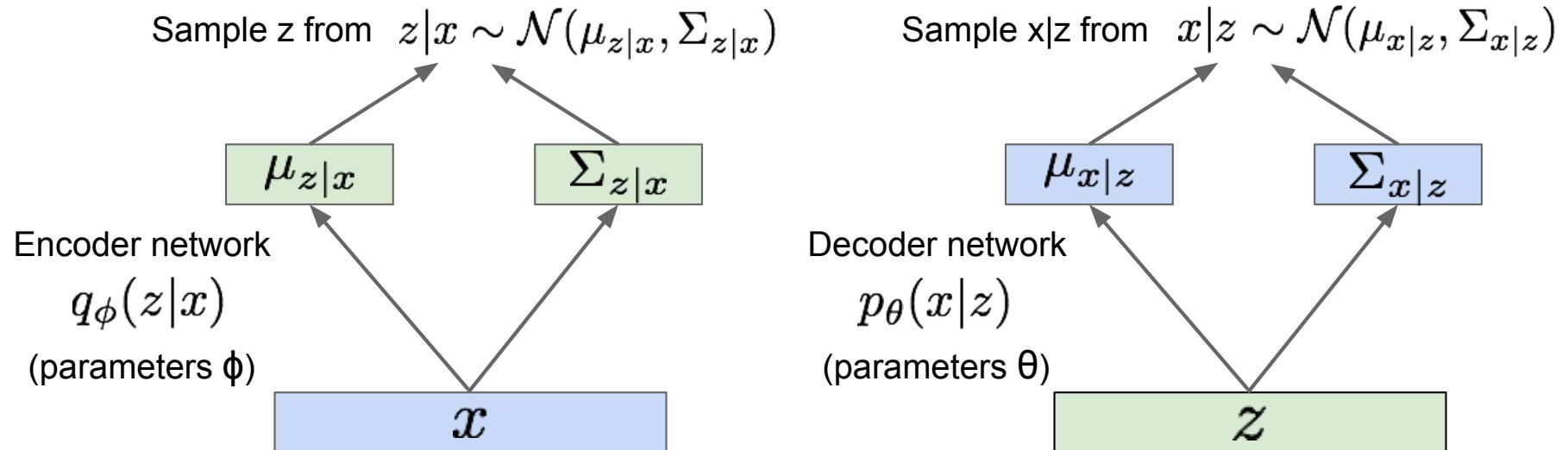
Posterior density also intractable:  $p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$

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

# 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

Assume that  $\Sigma_{x|z}$  and  $\Sigma_{z|x}$  are both diagonal, *i.e.* conditional independence.

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

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

↑  
Decoder network gives  $p_{\theta}(x|z)$ , can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

↑  
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  $\geq 0$ .

Lecture 13

The Kullback–Leibler divergence from  $\mathcal{N}_0(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  to  $\mathcal{N}_1(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ , for non-singular matrices  $\boldsymbol{\Sigma}_0$  and  $\boldsymbol{\Sigma}_1$ , is:

$$D_{KL}(\mathcal{N}_0 || \mathcal{N}_1) = \frac{1}{2} \left\{ \text{tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - k + \ln \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} \right\},$$

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}\end{aligned}$$

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

# Variational Autoencoders

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

Reconstruct the input data

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

Make approximate posterior distribution close to prior

$$\begin{aligned} &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{> 0} \end{aligned}$$

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

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

# Stage I in Forward Pass

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

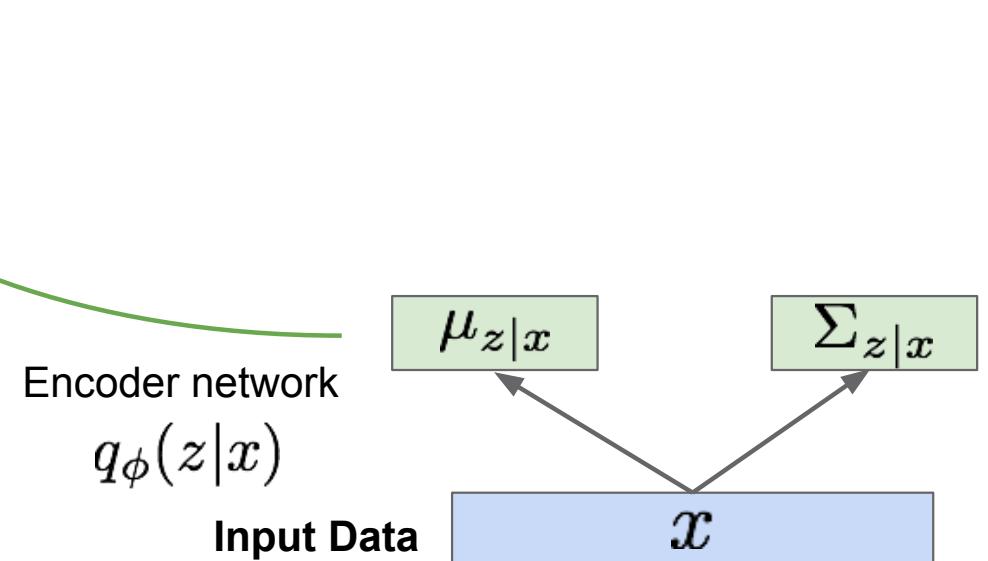
Encoder network

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

Input Data

$$\mu_{z|x}$$

$$\Sigma_{z|x}$$



# Stage II in forward pass

Putting it all together: maximizing the likelihood lower bound

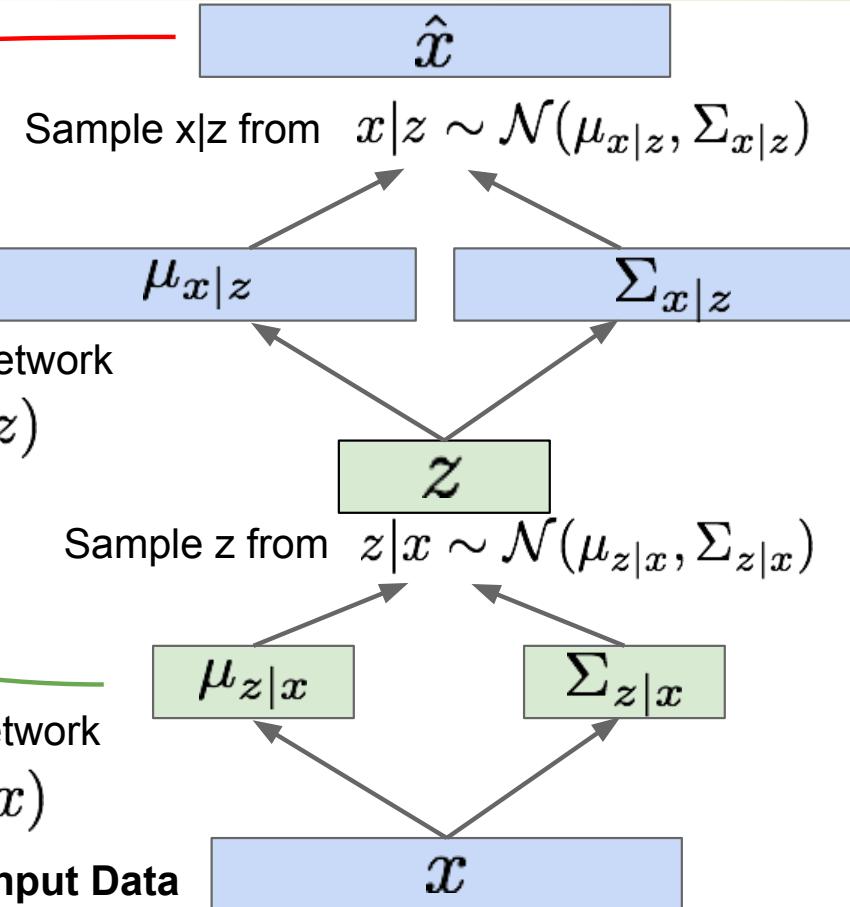
$$\underbrace{\mathbb{E}_z \left[ \log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Maximize likelihood of original input being reconstructed

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

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



# Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

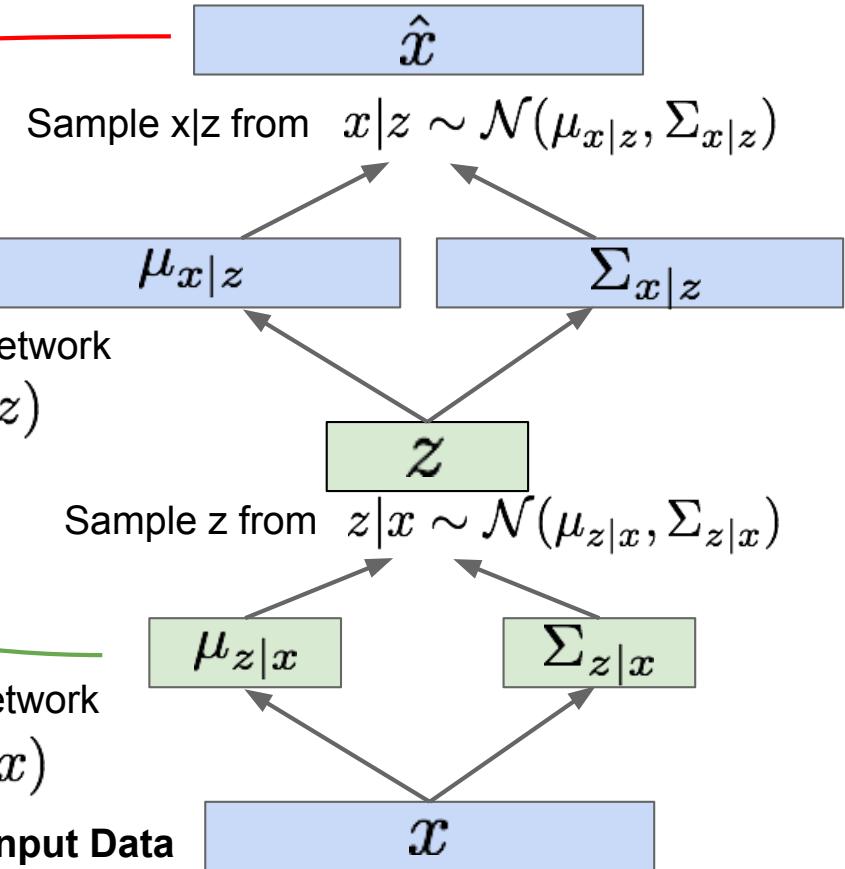
$$\underbrace{\mathbb{E}_z \left[ \log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

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

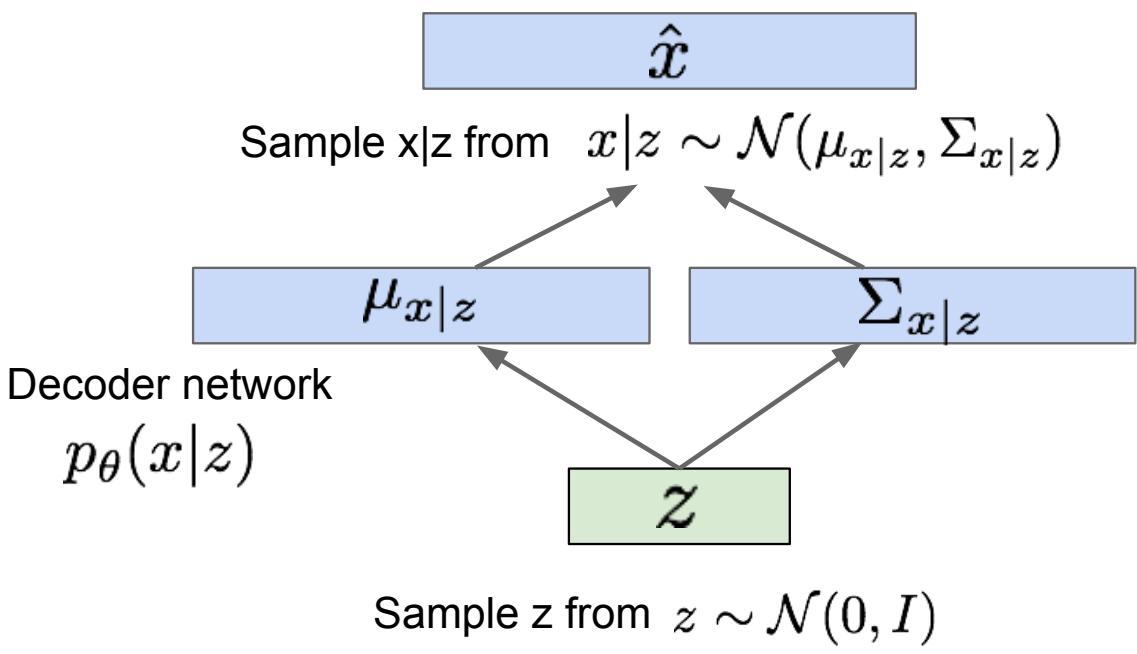
Maximize likelihood of original input being reconstructed

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



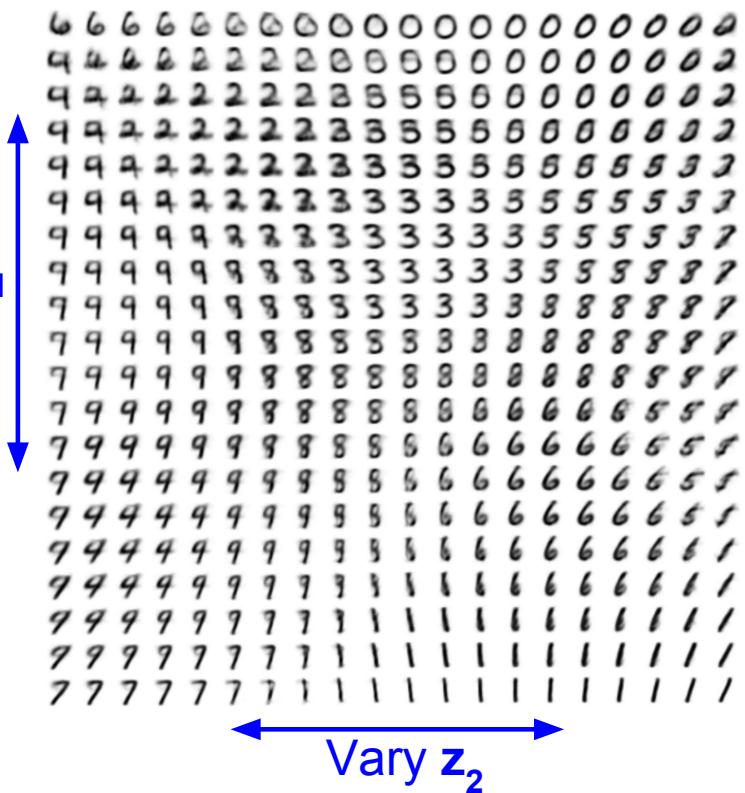
# VAE: generating data

Use decoder network. Now sample z from prior!



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

## Data manifold for 2-d $z$



# VAE: generating data

Diagonal prior on  $z$   
=> independent  
latent variables

Different  
dimensions of  $z$   
encode  
interpretable factors  
of variation

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

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

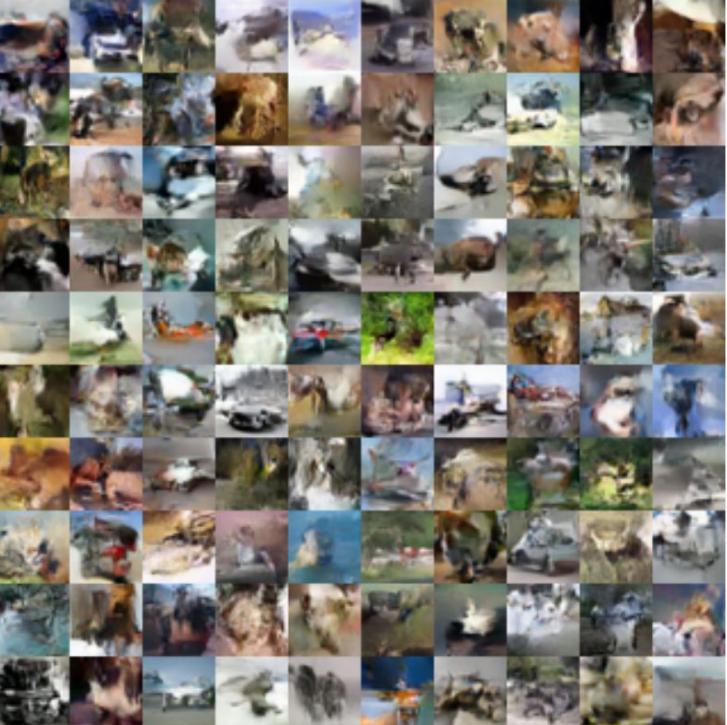
Degree of smile

Vary  $z_1$



Vary  $z_2$  Head pose

# VAE: 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

Thank you!

