

Lecture 19: Generative Models, Part 1

Last Time: Videos

Many video models:

Single-frame CNN (Try this first!)

Late fusion

Early fusion

3D CNN / C3D

Two-stream networks

CNN + RNN

Convolutional RNN

Spatio-temporal self-attention

SlowFast networks (current SoTA)

Today: Generative Models, Part 1

Supervised vs Unsupervised Learning

Supervised Learning

Classification



Data: (x, y)

x is data, y is label

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

Cat

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

[This image is CC0 public domain](#)

Supervised vs Unsupervised Learning

Supervised Learning

Object Detection



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.

[This image is CC0 public domain](#)

DOG, DOG, CAT

Supervised vs Unsupervised Learning

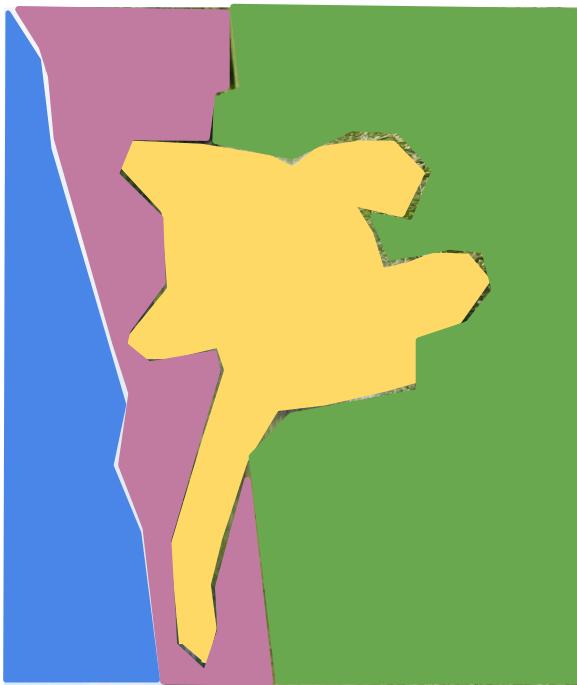
Supervised Learning

Data: (x, y)

x is data, y is label

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

Semantic Segmentation



GRASS, CAT, TREE, SKY

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 $x \rightarrow y$

Image captioning



*A cat sitting on a
suitcase on the floor*

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

Caption generated using neuraltalk2
Image is CC0 Public domain

Supervised vs Unsupervised Learning

Supervised Learning

Unsupervised Learning

Data: (x, y)

x is data, y is label

Data: x

Just data, no labels!

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

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

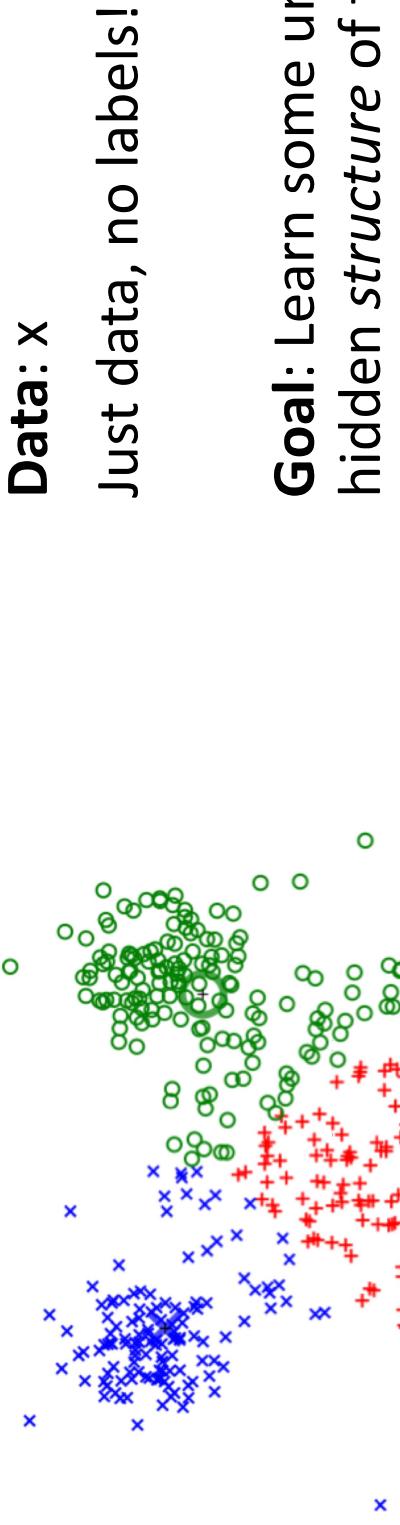
Examples: Classification, regression,
object detection, semantic
segmentation, image captioning, etc.

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

Supervised vs Unsupervised Learning

Unsupervised Learning

Clustering (e.g. K-Means)

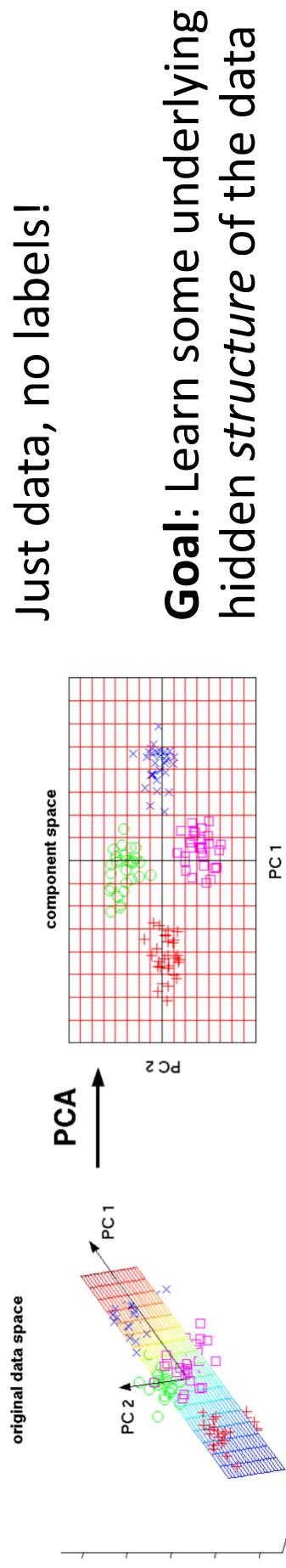


Goal: Learn some underlying hidden structure of the data

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

Supervised vs Unsupervised Learning

Dimensionality Reduction (e.g. Principal Components Analysis)



Goal: Learn some underlying hidden structure of the data

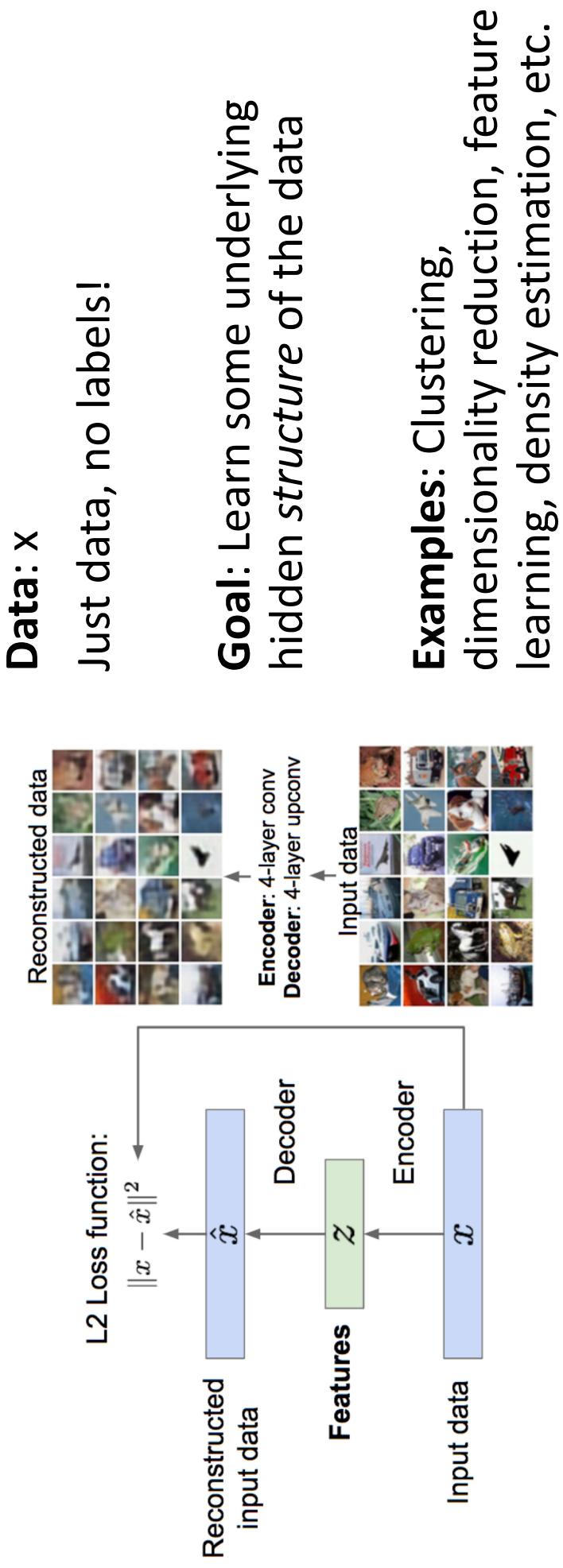
Data: x

Just data, no labels!

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

Supervised vs Unsupervised Learning

Feature Learning (e.g. autoencoders)



Supervised vs Unsupervised Learning

Unsupervised Learning

Density Estimation

Data: x

Just data, no labels!

Goal: Learn some unknown hidden structure of x .

0.4
0.3
0.2
0.1
0.0
-0.1
-0.2
-0.3
-0.4

400 300 200 100 0

400 300 200 100 0

0.08
0.06
0.04
0.02
0.00

5 4 3 2 1 0 -1 -2 -3 -4

8 6 4 2 0 -2 -4

x_1 x_2 x_3

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

Justin Johnson

Lecture 19 - 12

November 20, 2019

Images left and right are CCO public domain

Supervised vs Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Data: x

Just data, no labels!

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

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

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

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

Discriminative vs Generative Models

Discriminative Model:

Learn a probability distribution $p(y|x)$

Data: x



Generative Model:

Learn a probability distribution $p(x)$

Label: y

Cat

Conditional Generative Model:

Learn $p(x|y)$

Discriminative vs Generative Models

Probability Recap:

Discriminative Model:

Learn a probability distribution $p(y|x)$

Data: x



Density Function

$p(x)$ assigns a positive number to each possible x ; higher numbers mean x is more likely

Generative Model:

Learn a probability distribution $p(x)$

Density functions are

normalized:

$$\int_X p(x) dx = 1$$

Label: y
Cat

Conditional Generative Model: Learn $p(x|y)$

Different values of x
compete for density

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$

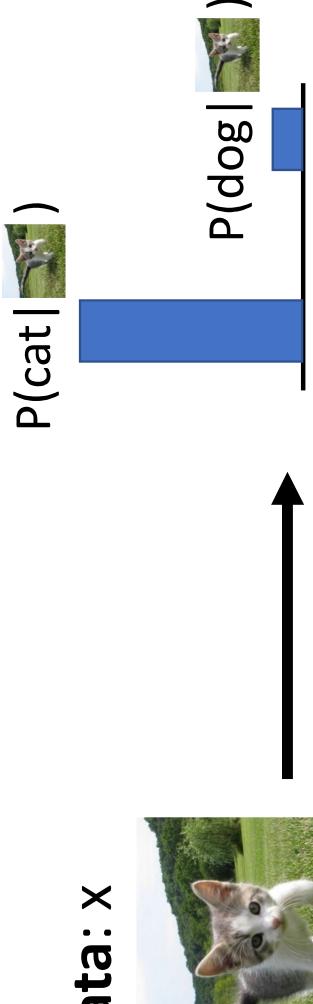


Data: x

$P(\text{cat} | \text{kitten})$

$P(\text{dog} | \text{dog})$

→



Generative Model:
Learn a probability distribution $p(x)$

Density functions are **normalized**:

$$\int_X p(x) dx = 1$$

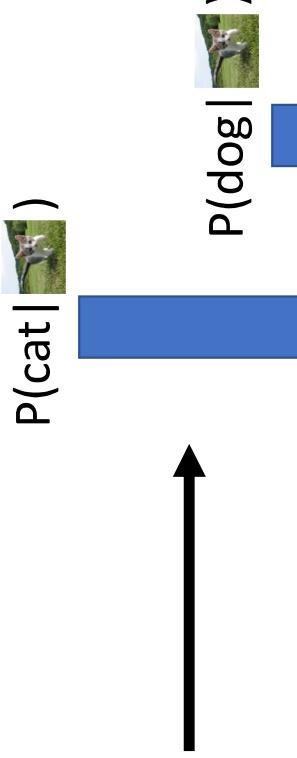
Different values of x **compete** for density

Discriminative vs Generative Models

Discriminative Model:

Learn a probability distribution $p(y|x)$

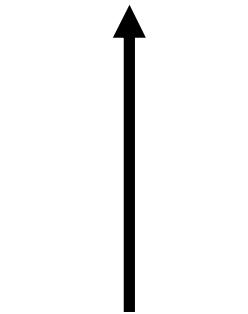
$$P(\text{cat} | \text{img})$$



Generative Model:

Learn a probability distribution $p(x)$

$$P(\text{dog} | \text{img})$$



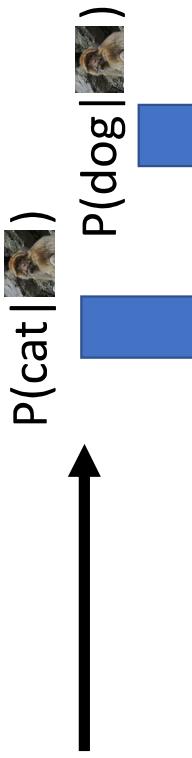
Conditional Generative Model: Learn $p(x|y)$

Discriminative model: the possible labels for each input "compete" for probability mass.
But no competition between **images**

Discriminative vs Generative Models

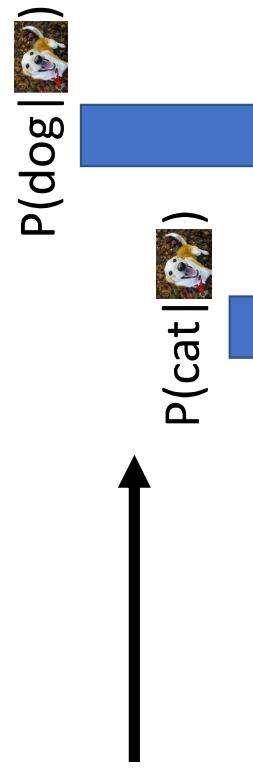
Discriminative Model:

Learn a probability distribution $p(y|x)$



Generative Model:

Learn a probability distribution $p(x)$



Conditional Generative Model:

Learn $p(x|y)$

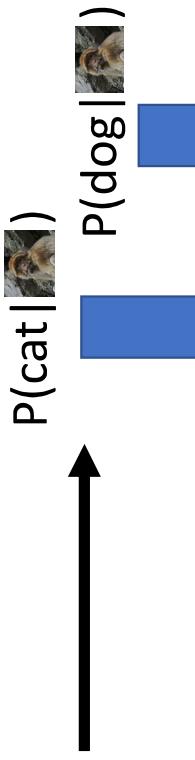
Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

Monkey image is CC0 Public Domain

Discriminative vs Generative Models

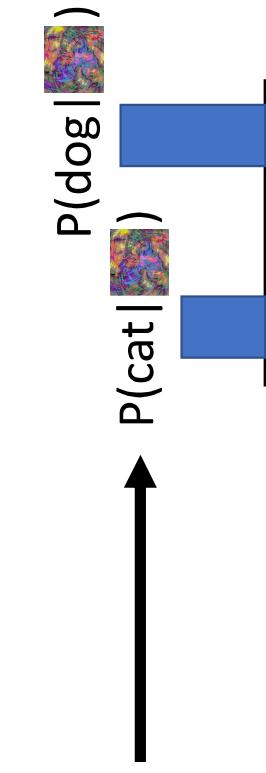
Discriminative Model:

Learn a probability distribution $p(y|x)$



Generative Model:

Learn a probability distribution $p(x)$



Conditional Generative Model:

Learn $p(x|y)$

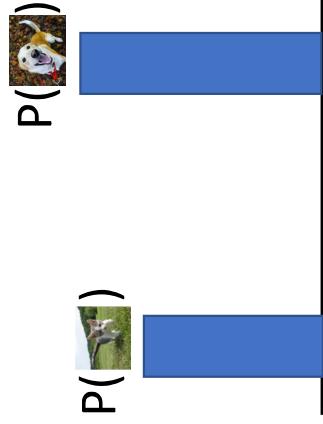
Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

Monkey image is CC0 Public Domain
Abstract image is free to use under the Pixabay license

Discriminative vs Generative Models

This image is CC0 Public Domain
This image is CC0 Public Domain
This image is CC0 Public Domain
This image is CC0 Public Domain

Discriminative Model:
Learn a probability distribution $p(y|x)$



Generative Model:
Learn a probability distribution $p(x)$

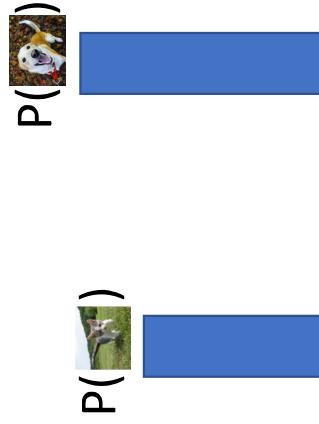


Generative model: All possible images compete with each other for probability mass

Conditional Generative Model: Learn $p(x|y)$

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$



Generative Model:
Learn a probability distribution $p(x)$



Generative model: All possible images compete with each other for probability mass

Conditional Generative Model: Learn $p(x|y)$

Requires deep image understanding! Is a dog more likely to sit or stand? How about 3-legged dog vs 3-armed monkey?

Discriminative vs Generative Models

Discriminative Model:
Learn a probability distribution $p(y|x)$



Generative Model:
Learn a probability distribution $p(x)$



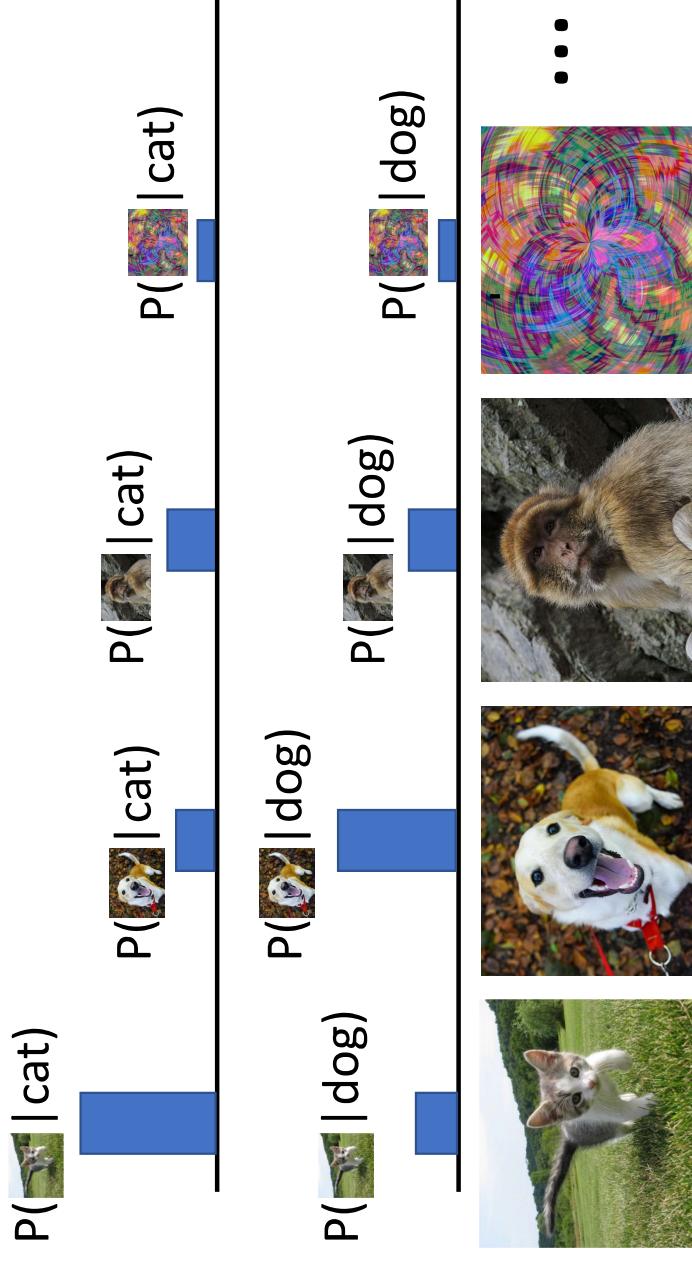
...

Generative model: All possible images compete with each other for probability mass

Conditional Generative Model: Learn $p(x|y)$

Model can “reject” unreasonable inputs by assigning them small values

Discriminative vs Generative Models



Conditional Generative Model: Learn $p(x|y)$

Conditional Generative Model: Each possible label induces a competition among all images

Discriminative vs Generative Models

Software is ©2019 public-domain
and may be freely distributed.
However, it is CC0 Public Domain
and therefore cannot be used under the
Creative Commons license.

Discriminative Model:

Learn a probability distribution $p(y|x)$

Recall Bayes' Rule:

$$P(x | y) = \frac{P(y | x)}{P(y)} P(x)$$

Generative Model:

Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$

Discriminative vs Generative Models

Discriminative Model:

Learn a probability distribution $p(y|x)$

Recall Bayes' Rule:

$$P(x|y) = \frac{P(y|x)}{P(y)}$$

Discriminative Model Conditional Generative Model

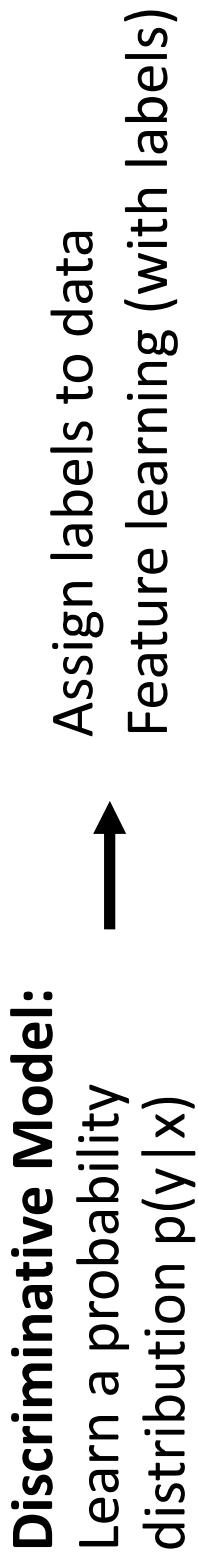
Unconditional Generative Model Prior over labels

Generative Model:
Learn a probability distribution $p(x)$

Conditional Generative Model: Learn $p(x|y)$

We can build a conditional generative model from other components!

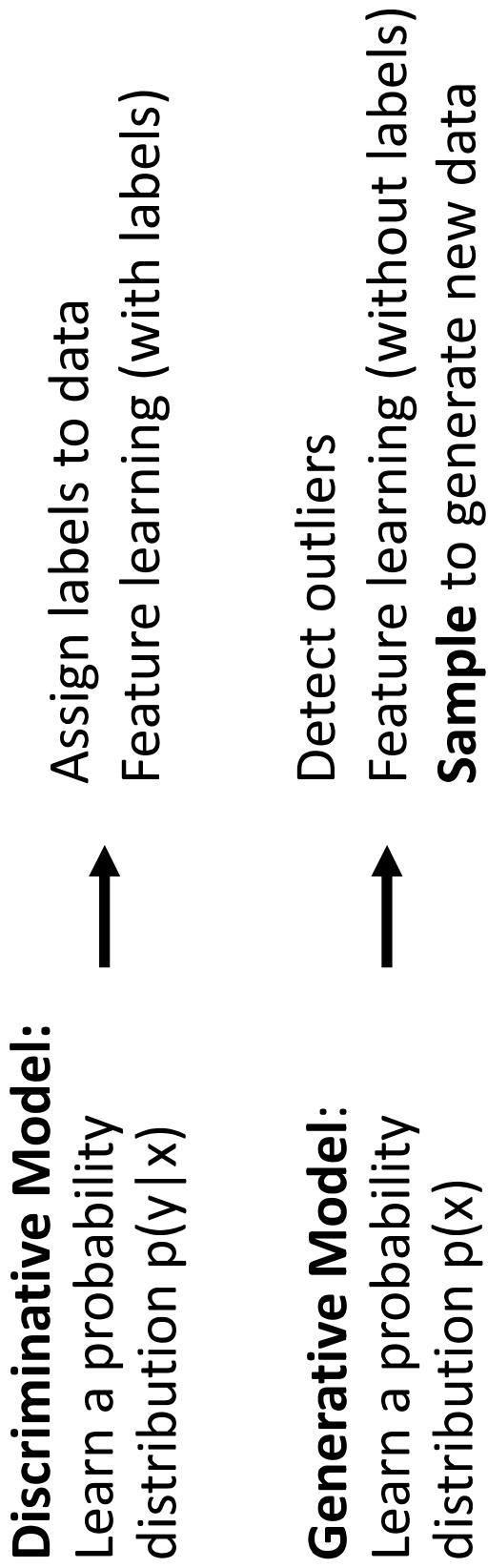
What can we do with a discriminative model?



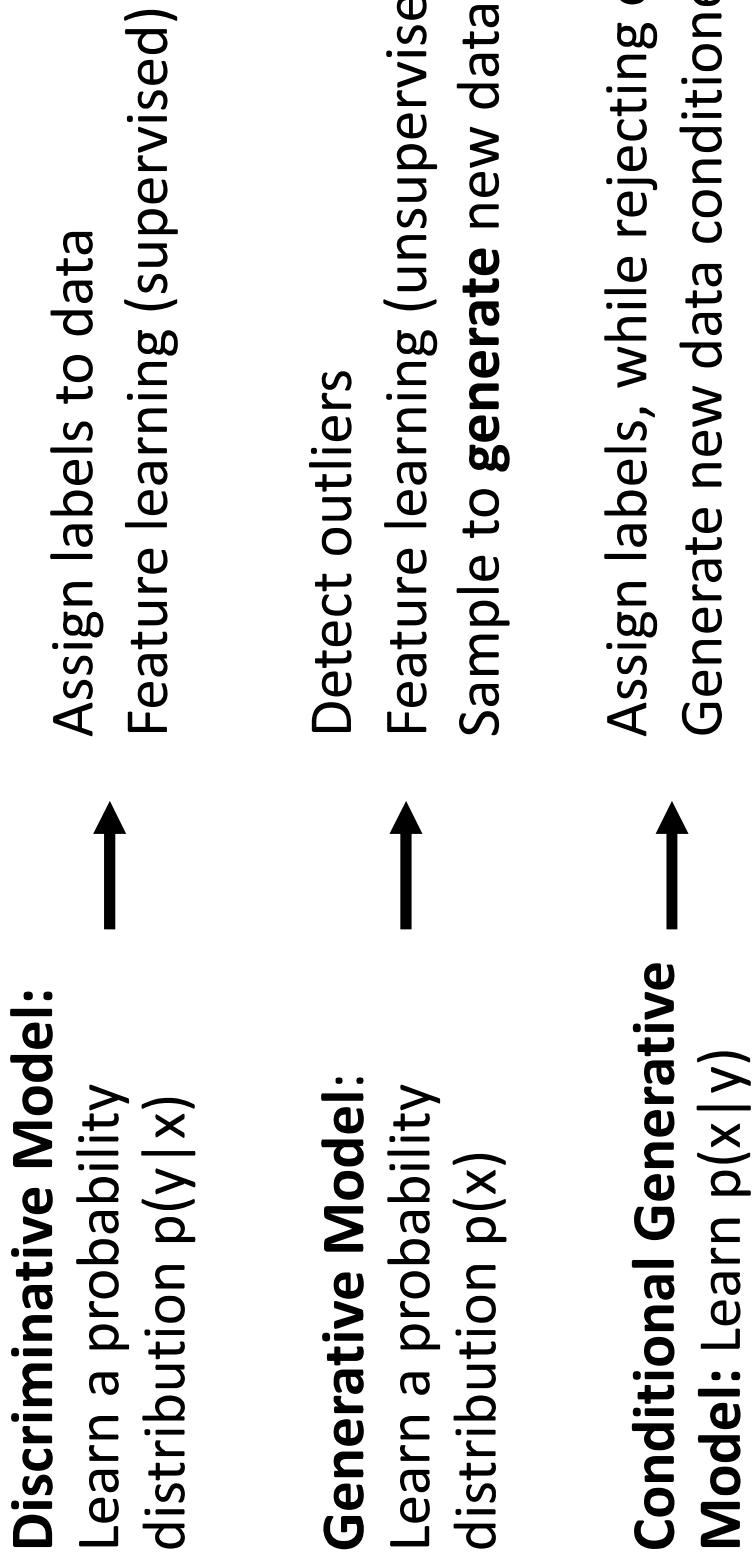
Generative Model:
Learn a probability
distribution $p(x)$

**Conditional Generative
Model:** Learn $p(x|y)$

What can we do with a generative model?



What can we do with a generative model?



Taxonomy of Generative Models

Generative models

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

Justin Johnson

Lecture 19 - 29

November 20, 2019

Taxonomy of Generative Models

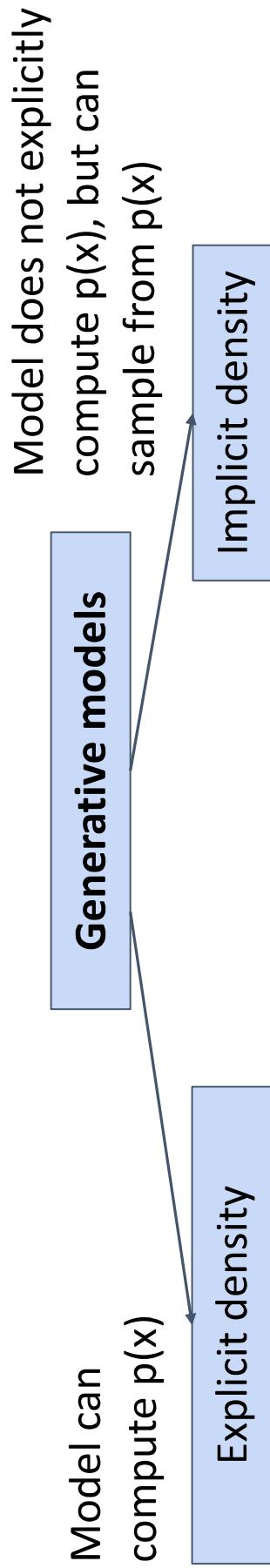


Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

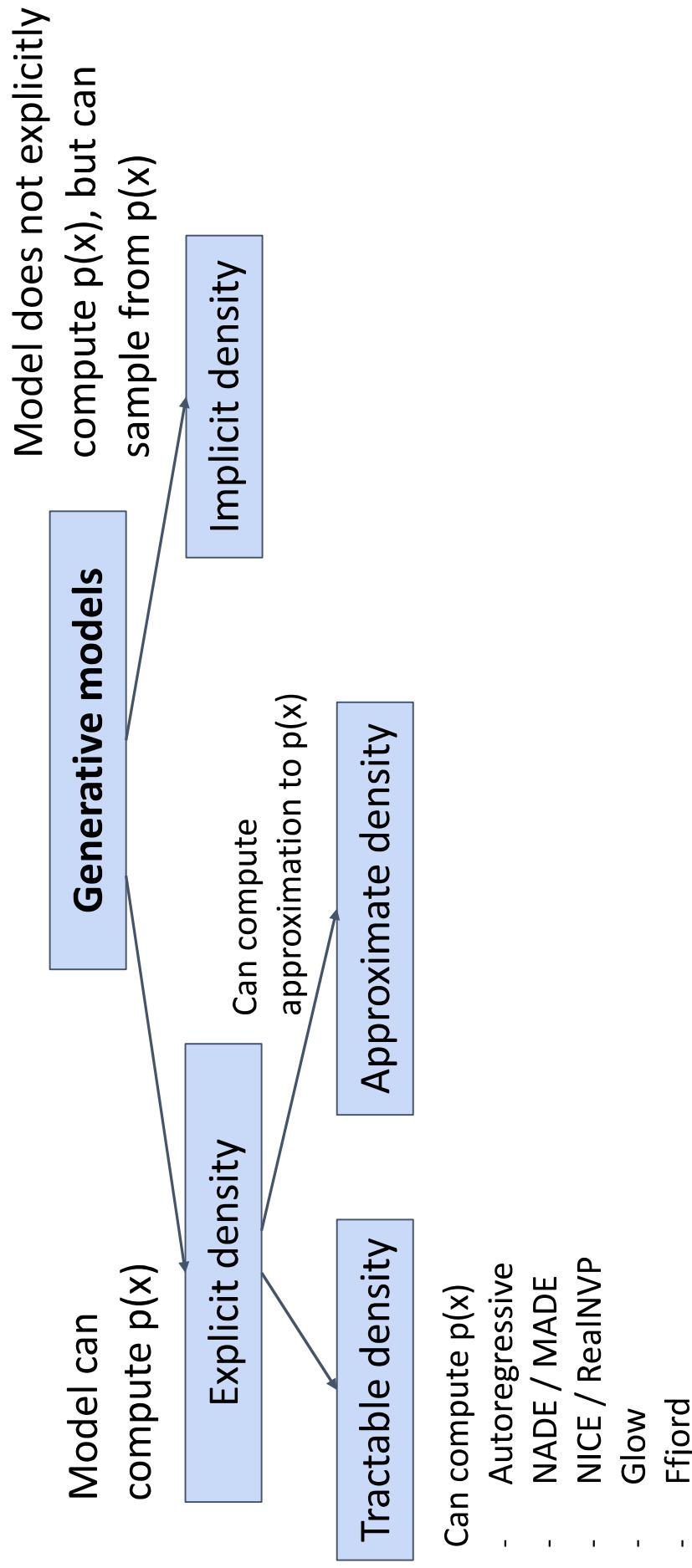


Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

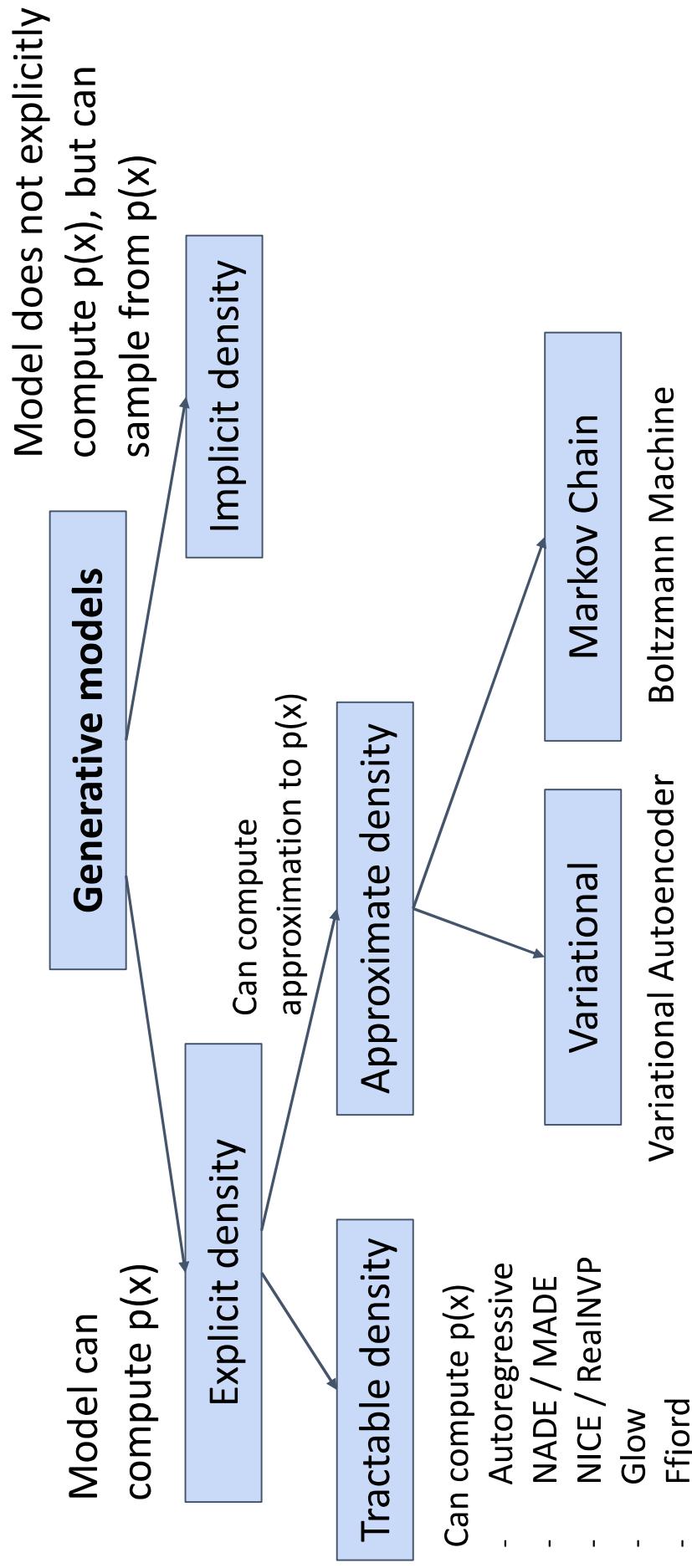


Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

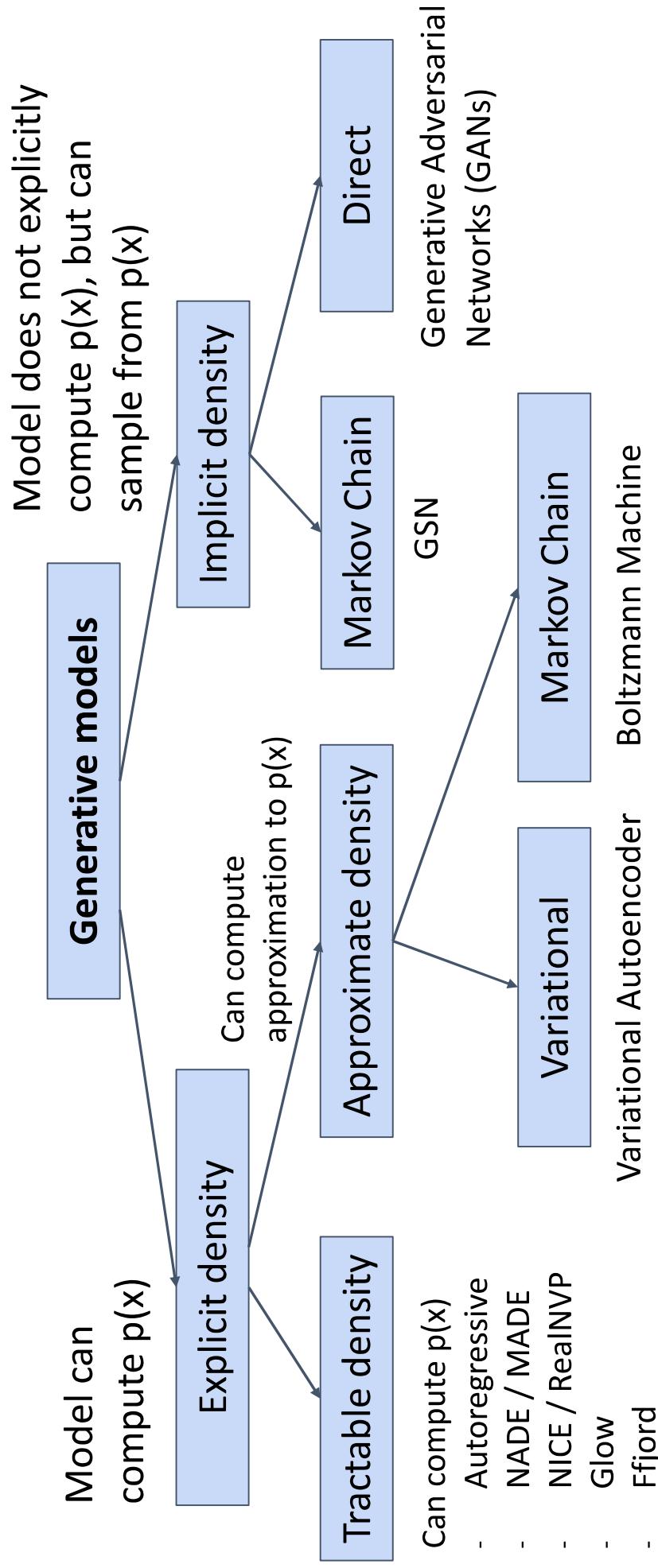


Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

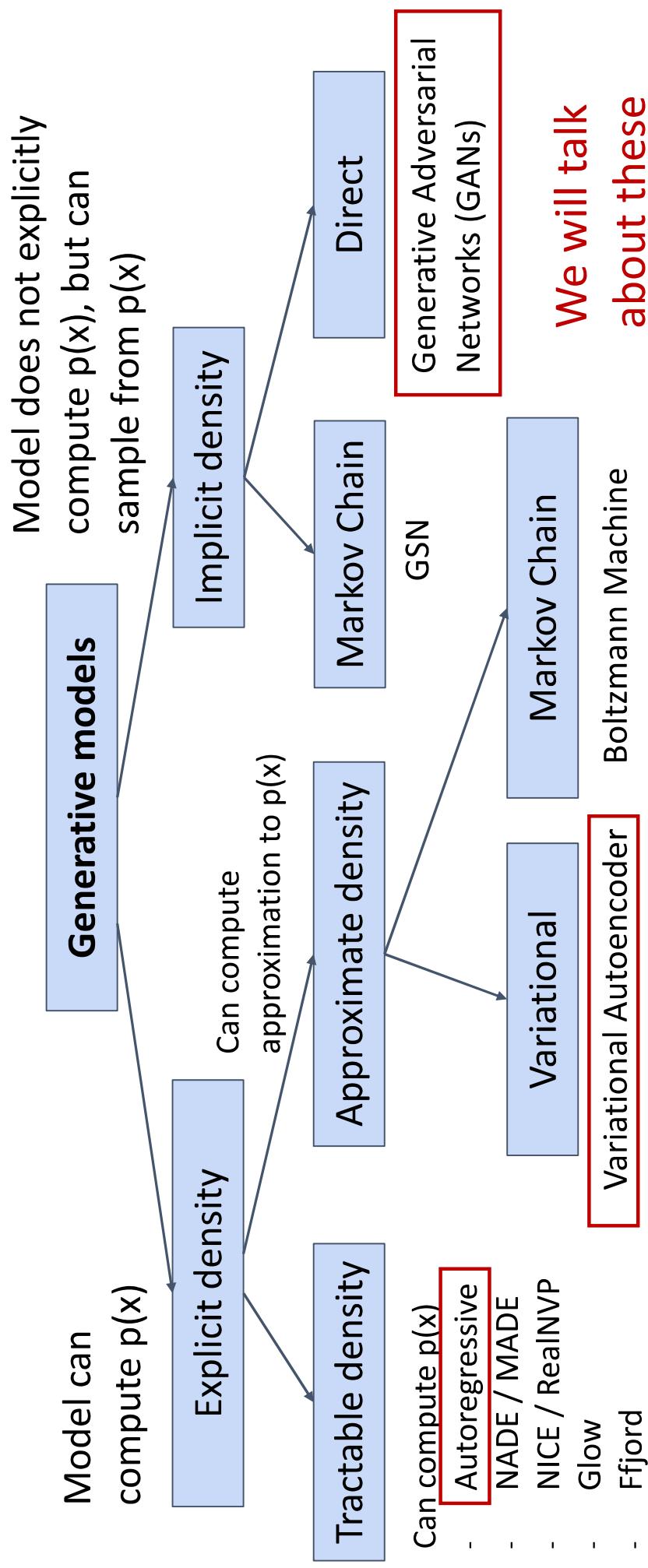


Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Autoregressive models

Explicit Density Estimation

Goal: Write down an explicit function for $p(x) = f(x, W)$

Explicit Density Estimation

Goal: Write down an explicit function for $p(x) = f(x, W)$

Given dataset $x^{(1)}, x^{(2)}, \dots, x^{(N)}$, train the model by solving:

$$W^* = \arg \max_W \prod_i p(x^{(i)})$$

Maximize probability of training data
(Maximum likelihood estimation)

Explicit Density Estimation

Goal: Write down an explicit function for $p(x) = f(x, W)$

Given dataset $x^{(1)}, x^{(2)}, \dots, x^{(N)}$, train the model by solving:

$$\begin{aligned} W^* &= \arg \max_W \prod_i p(x^{(i)}) && \text{Maximize probability of training data} \\ &= \arg \max_W \sum_i \log p(x^{(i)}) && \text{(Maximum likelihood estimation)} \end{aligned}$$

Log trick to exchange product for sum

Explicit Density Estimation

Goal: Write down an explicit function for $p(x) = f(x, W)$

Given dataset $x^{(1)}, x^{(2)}, \dots, x^{(N)}$, train the model by solving:

$$\begin{aligned} W^* &= \arg \max_W \prod_i p(x^{(i)}) && \text{Maximize probability of training data} \\ &= \arg \max_W \sum_i \log p(x^{(i)}) && \text{Log trick to exchange product for sum} \\ &= \arg \max_W \sum_i \log f(x^{(i)}, W) && \begin{array}{l} \text{This will be our loss function!} \\ \text{Train with gradient descent} \end{array} \end{aligned}$$

Explicit Density: Autoregressive Models

Goal: Write down an explicit function for $p(x) = f(x, W)$

Assume x consists of
multiple subparts:
$$x = (x_1, x_2, x_3, \dots, x_T)$$

Explicit Density: Autoregressive Models

Goal: Write down an explicit function for $p(x) = f(x, W)$

Assume x consists of
multiple subparts:

$$x = (x_1, x_2, x_3, \dots, x_T)$$

Break down probability
using the chain rule:

$$\begin{aligned} p(x) &= p(x_1, x_2, x_3, \dots, x_T) \\ &= p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2) \dots \end{aligned}$$

Explicit Density: Autoregressive Models

Goal: Write down an explicit function for $p(x) = f(x, W)$

Assume x consists of
multiple subparts:

$$x = (x_1, x_2, x_3, \dots, x_T)$$

Break down probability
using the chain rule:

$$\begin{aligned} p(x) &= p(x_1, x_2, x_3, \dots, x_T) \\ &= p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2) \dots \\ &= \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}) \end{aligned}$$

Probability of the next subpart
given all the previous subparts

Explicit Density: Autoregressive Models

Goal: Write down an explicit function for $p(x) = f(x, W)$

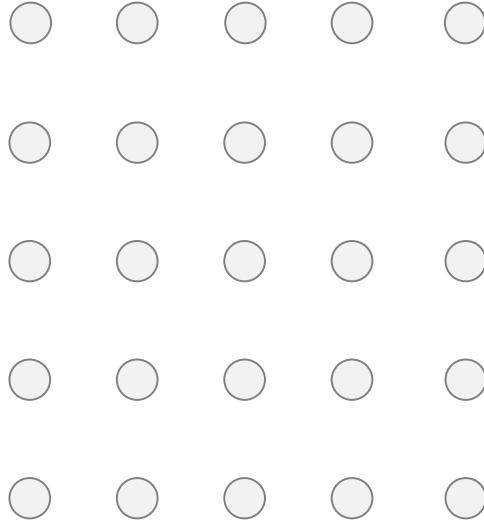
Assume x consists of
multiple subparts:

Break down probability
using the chain rule:

$$\begin{aligned} p(x) &= p(x_1, x_2, x_3, \dots, x_T) \\ &= p(x_1)p(x_2 \mid x_1)p(x_3 \mid x_1, x_2) \dots \\ &= \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}) \\ &\quad \text{Probability of the next subpart} \\ &\quad \text{given all the previous subparts} \\ p(x_1) &\quad p(x_2) \quad p(x_3) \quad p(x_4) \quad \text{We've already} \\ \uparrow &\quad \uparrow &\quad \uparrow &\quad \uparrow \quad \text{seen this!} \\ h_1 \rightarrow h_2 \rightarrow h_3 \rightarrow h_4 &\quad \rightarrow h_4 \quad \text{Language} \\ \uparrow &\quad \uparrow &\quad \uparrow &\quad \uparrow \quad \text{modeling with} \\ x_0 &\quad x_1 \quad x_2 \quad x_3 \quad \text{an RNN!} \end{aligned}$$

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner

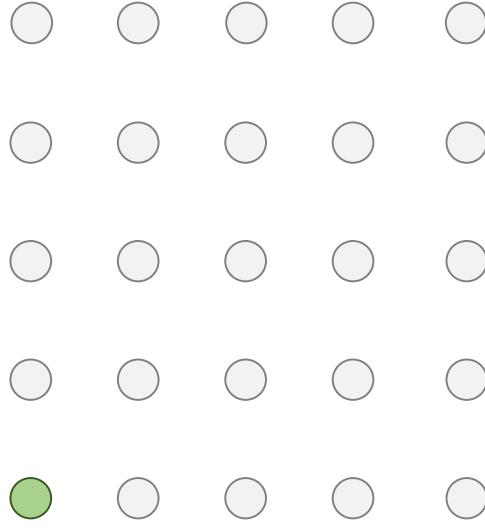


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixelRNN

Generate image pixels one at a time, starting at the upper left corner

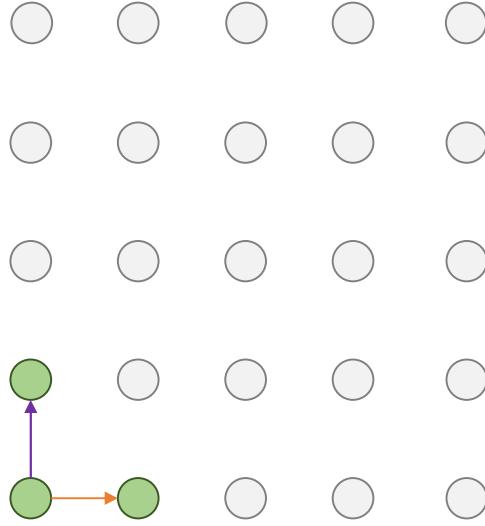


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixelRNN

Generate image pixels one at a time, starting at the upper left corner

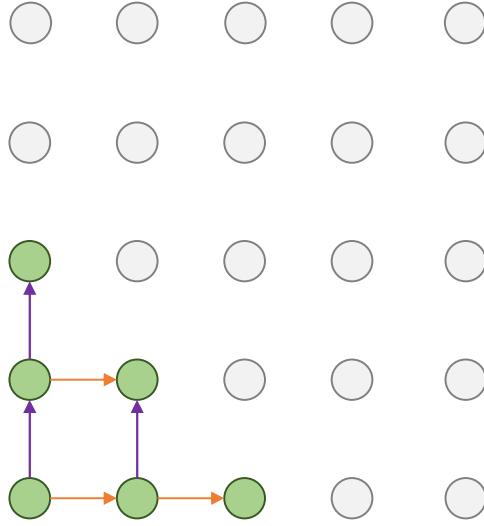


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner

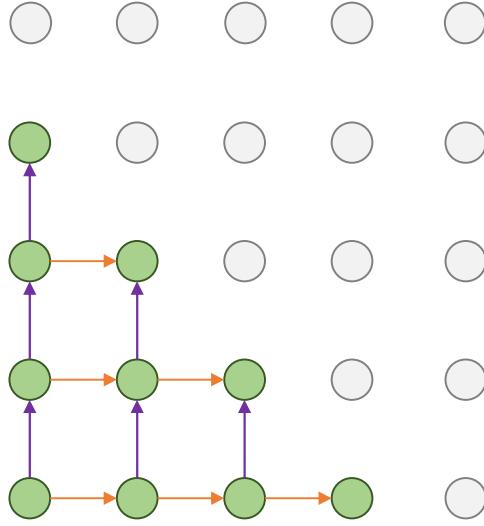


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner

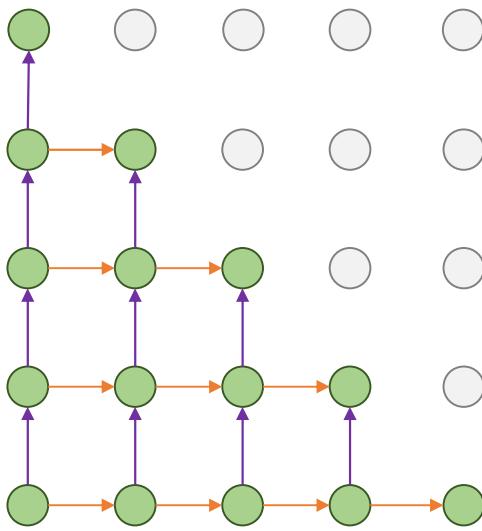


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner

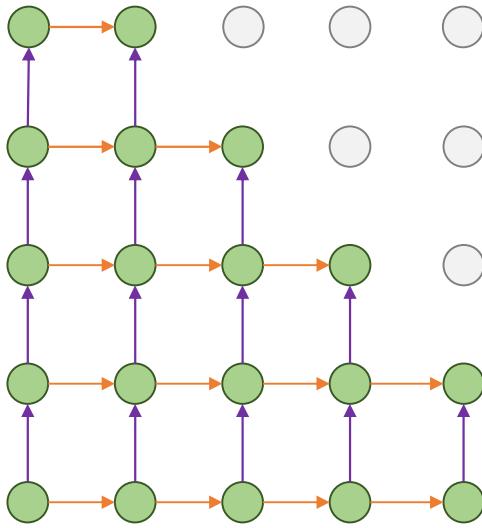


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner

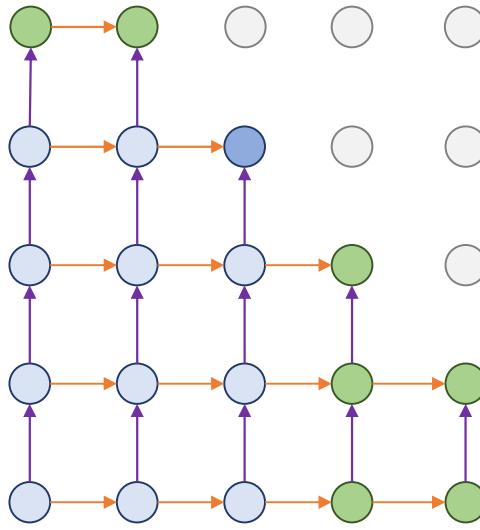


Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

PixelRNN

Generate image pixels one at a time, starting at the upper left corner



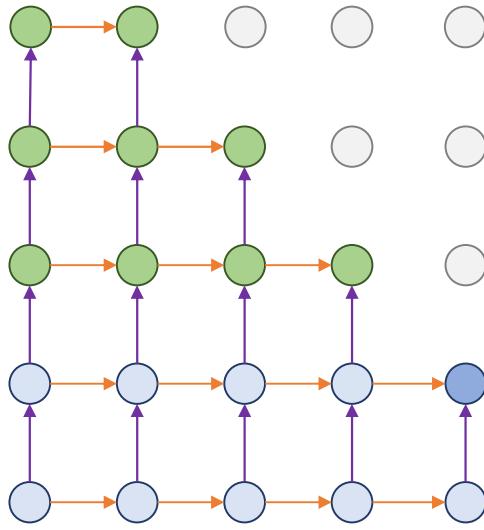
Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

Each pixel depends **implicity** on all pixels above and to the left:

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner



Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$

At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

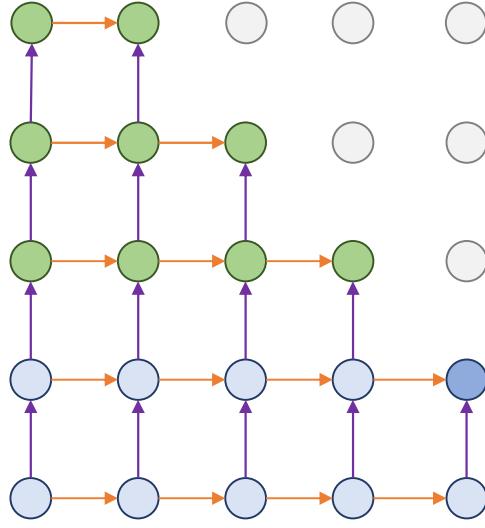
Each pixel depends **implicity** on all pixels above and to the left:

Van den Oord et al, "Pixel Recurrent Neural Networks", ICML 2016

PixeIRNN

Generate image pixels one at a time, starting at the upper left corner

Compute a hidden state for each pixel that depends on hidden states and RGB values from the left and from above (LSTM recurrence)
$$h_{x,y} = f(h_{x-1,y}, h_{x,y-1}, W)$$



At each pixel, predict red, then blue, then green:
softmax over $[0, 1, \dots, 255]$

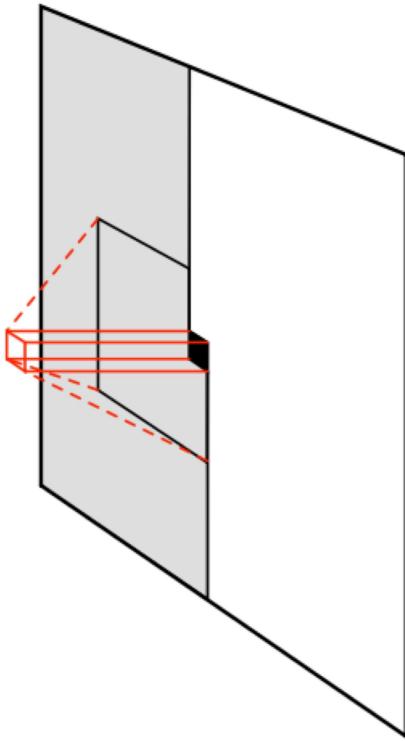
Each pixel depends **implicity** on all pixels above and to the left:

Van den Oord et al, "Pixel Recurrent Neural Networks", ICMl 2016

PixelCNN

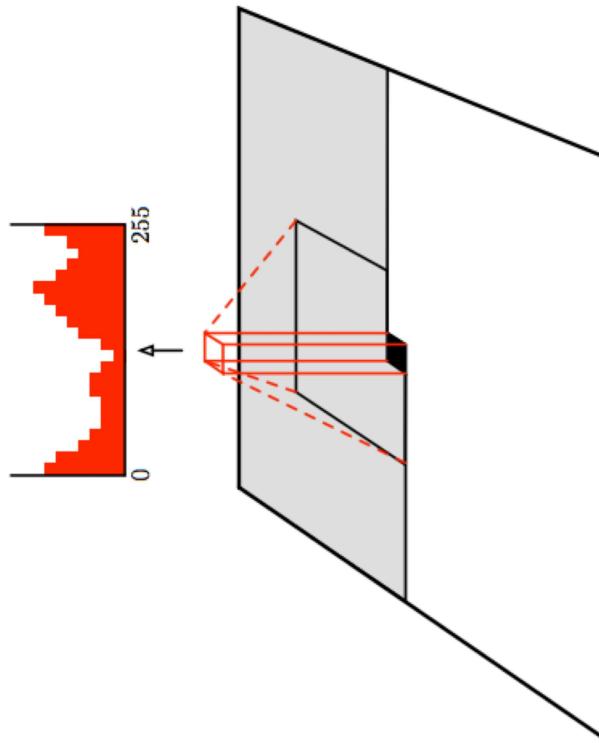
Still generate image pixels starting from corner

Dependency on previous pixels now modeled
using a CNN over context region



PixelCNN

Softmax loss
at each pixel



Still generate image pixels starting from corner

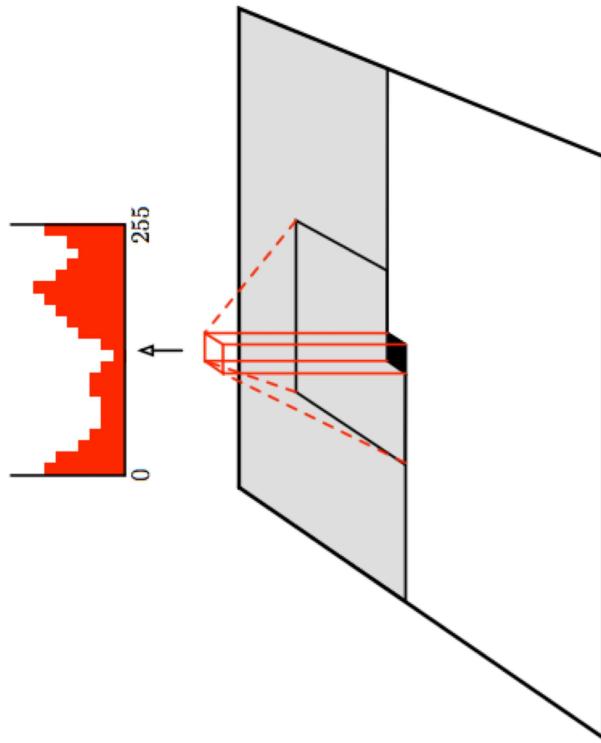
Dependency on previous pixels now modeled
using a CNN over context region

Training: maximize likelihood of training images

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

PixelCNN

Softmax loss
at each pixel



Still generate image pixels starting from corner

Dependency on previous pixels now modeled
using a CNN over context region

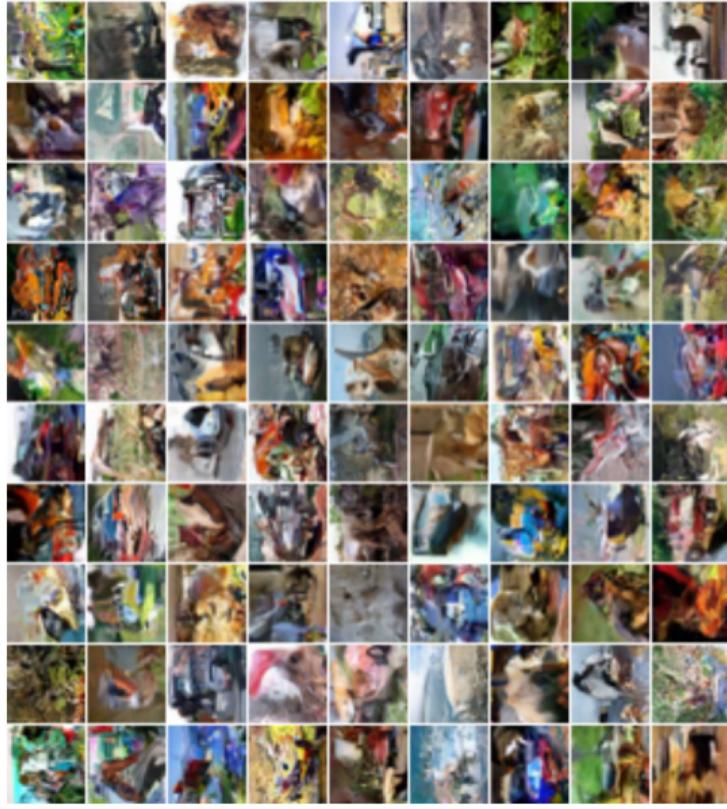
Training: maximize likelihood of training images

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

Generation must still proceed sequentially
=> still slow

Van den Oord et al, "Conditional Image Generation with PixelCNN Decoders", NeurIPS 2016

PixelRNN: Generated Samples



32x32 CIFAR-10

32x32 ImageNet

Autoregressive Models: PixelRNN and PixelCNN

Improving PixelCNN performance

Pros:

- Can explicitly compute likelihood $p(x)$
- Explicit likelihood of training data gives good evaluation metric
- Good samples

Con:

- Sequential generation => slow

See

- Van der Oord et al. NIPS 2016
- Salimans et al. 2017 (PixelCNN++)

Variational Autoencoders

Variational Autoencoders

PixelRNN / PixelCNN explicitly parameterizes density function with a neural network, so we can train to maximize likelihood of training data:

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

Variational Autoencoders (VAE) define an **intractable density** that we cannot explicitly compute or optimize

But we will be able to directly optimize a **lower bound** on the density

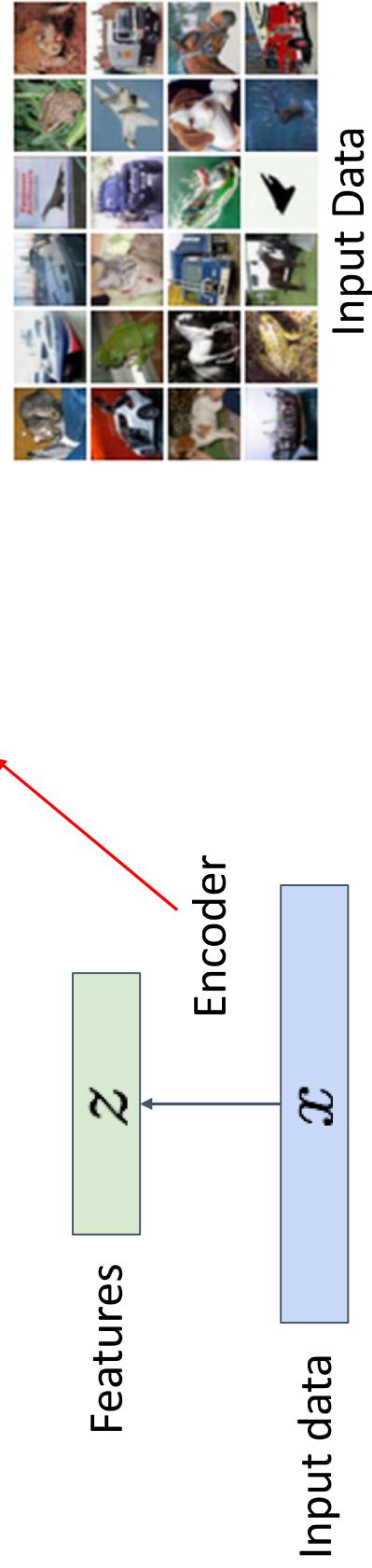
Variational Autoencoders

(Regular, non-variational) Autoencoders

Unsupervised method for learning feature vectors from raw data x , without any labels

Features should extract useful information (maybe object identities, properties, scene type, etc) that we can use for downstream tasks

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

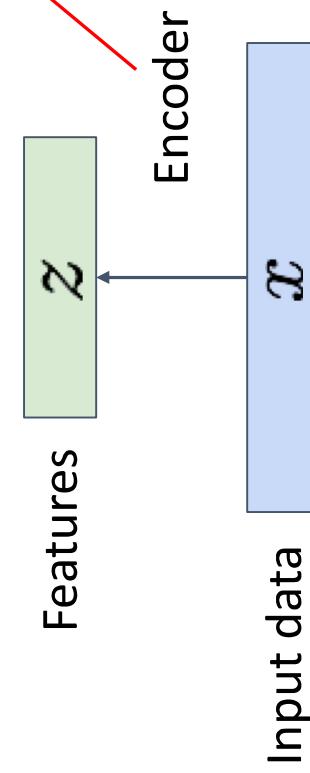
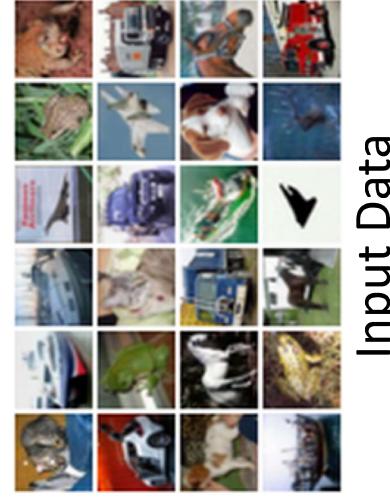


(Regular, non-variational) Autoencoders

Problem: How can we learn this feature transform from raw data?

Features should extract useful information (maybe object identities, properties, scene type, etc) that we can use for downstream tasks
But we can't observe features!

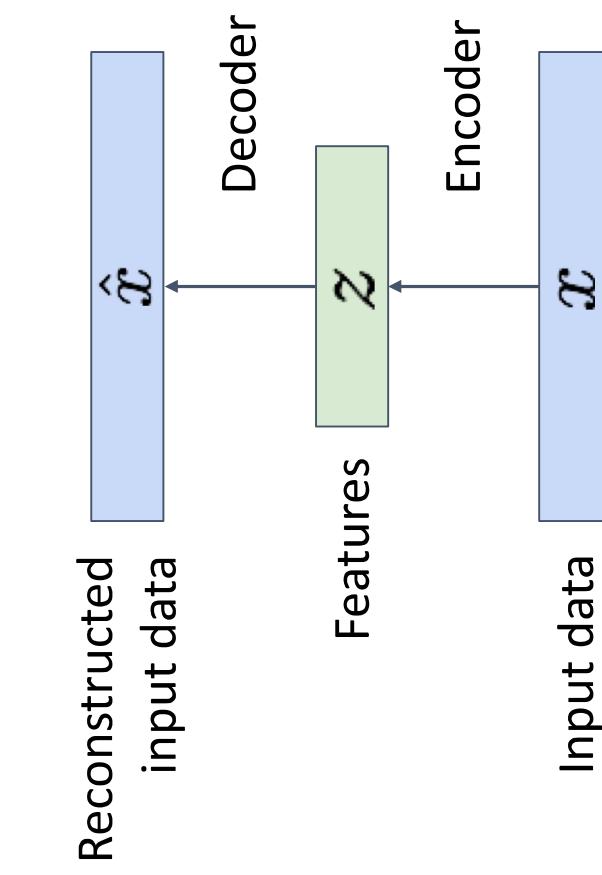
Originally: Linear + nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN



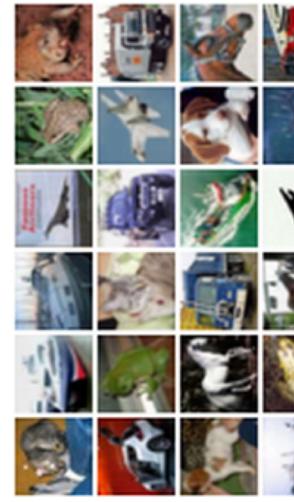
(Regular, non-variational) Autoencoders

Problem: How can we learn this feature transform from raw data?

Idea: Use the features to reconstruct the input data with a **decoder**
“Autoencoding” = encoding itself

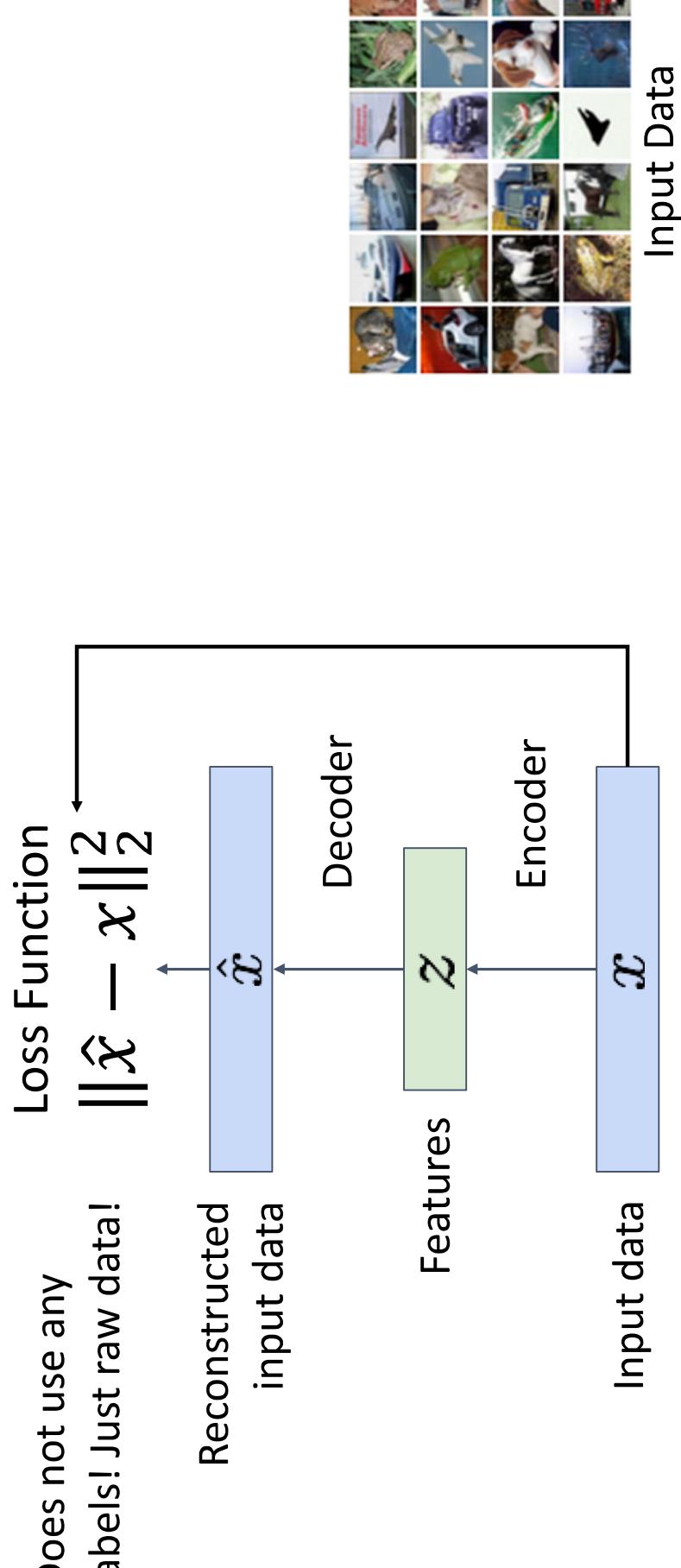


Originally: Linear +
nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN (upconv)



(Regular, non-variational) Autoencoders

Loss: L2 distance between input and reconstructed data.



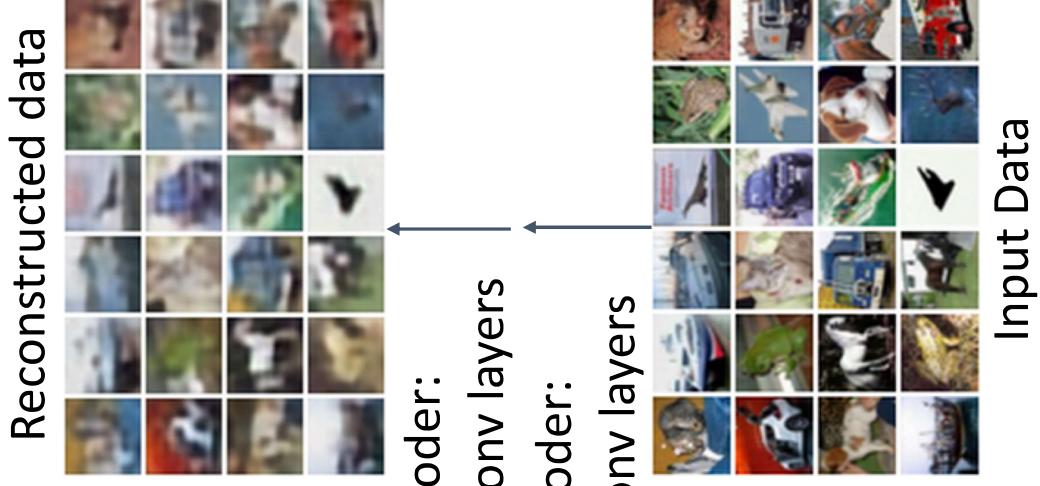
(Regular, non-variational) Autoencoders

Loss: L2 distance between input and reconstructed data.

Does not use any
labels! Just raw data!

$$\text{Loss Function} \quad \|\hat{x} - x\|_2^2$$

Reconstructed
input data



(Regular, non-variational) Autoencoders

Loss: L2 distance between input and reconstructed data.

Does not use any
labels! Just raw data!

Loss Function

$$\|\hat{x} - x\|_2^2$$

Reconstructed
input data

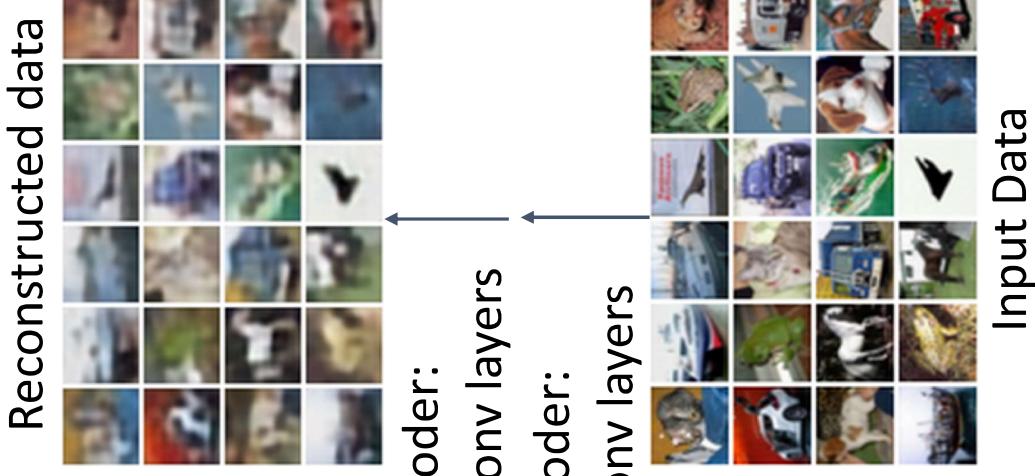
Decoder

Features need to be
lower dimensional
than the data

Decoder

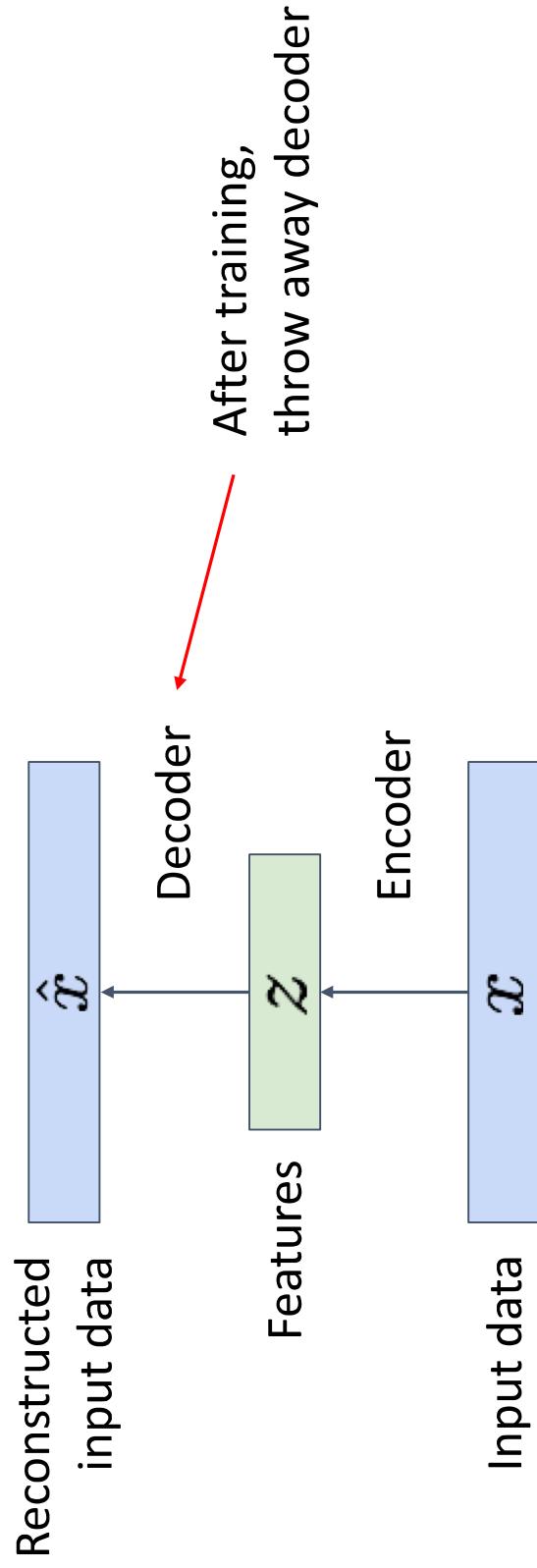
Encoder

Input data



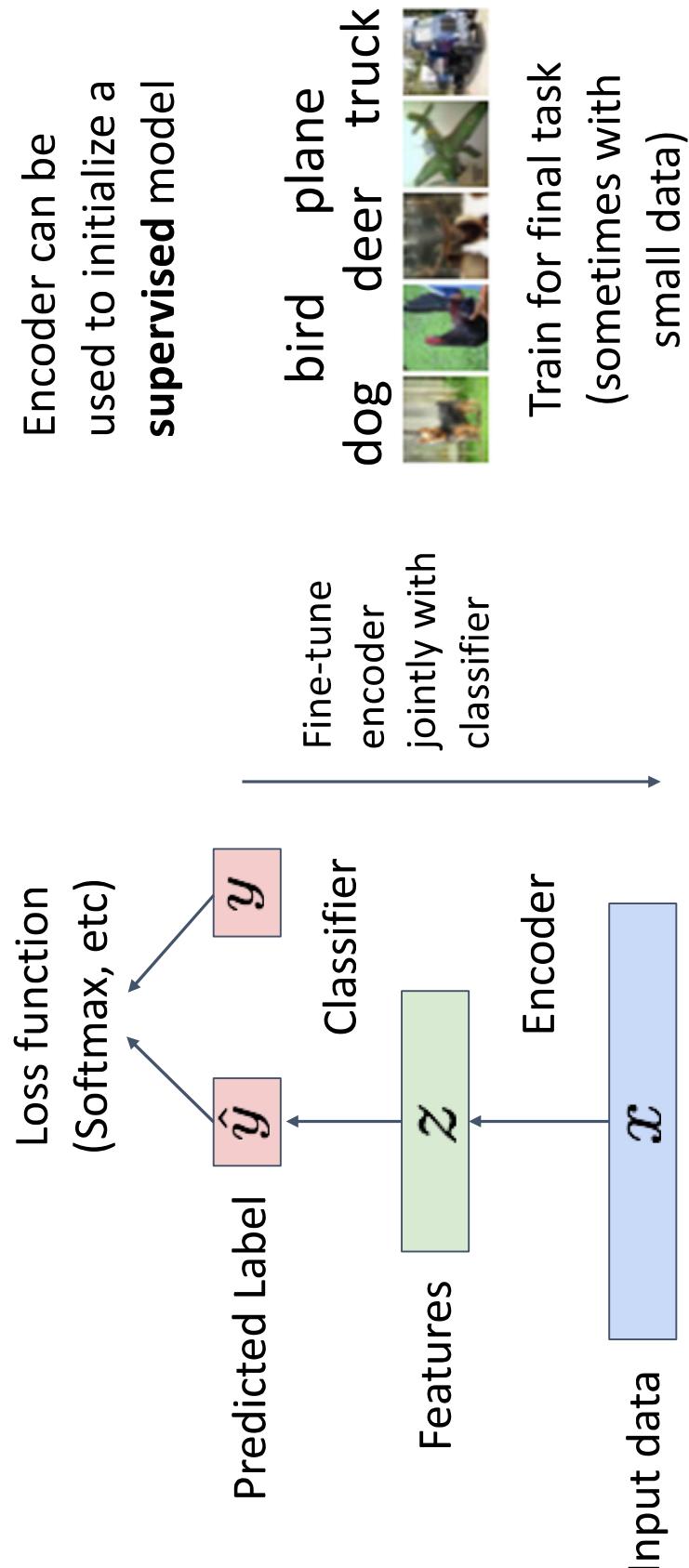
(Regular, non-variational) Autoencoders

After training, **throw away decoder** and use encoder for a downstream task



(Regular, non-variational) Autoencoders

After training, **throw away decoder** and use encoder for a downstream task

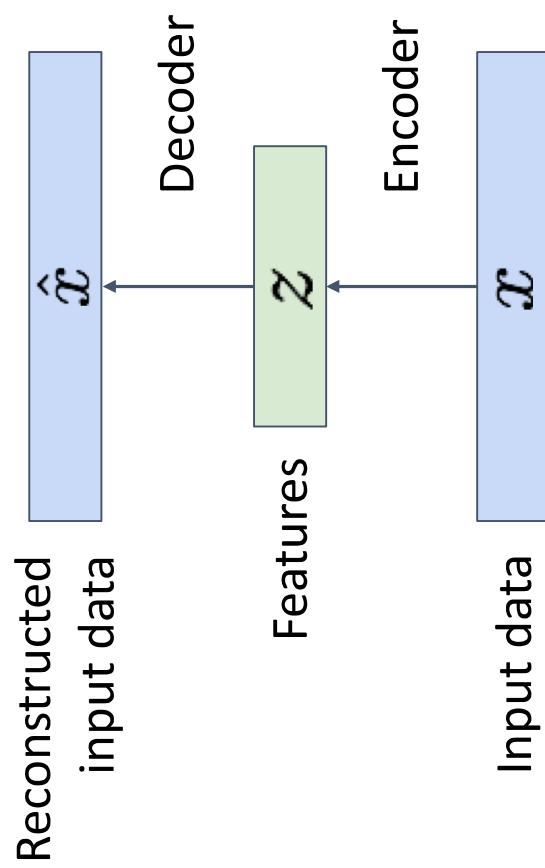


(Regular, non-variational) Autoencoders

Autoencoders learn **latent features** for data without any labels!

Can use features to initialize a **supervised** model

Not probabilistic: No way to sample new data from learned model



Variational Autoencoders

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

November 20, 2019

Lecture 19 - 71

Justin Johnson

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features z from raw data
2. Sample from the model to generate new data

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features \mathbf{z} from raw data
2. Sample from the model to generate new data representation \mathbf{z}

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

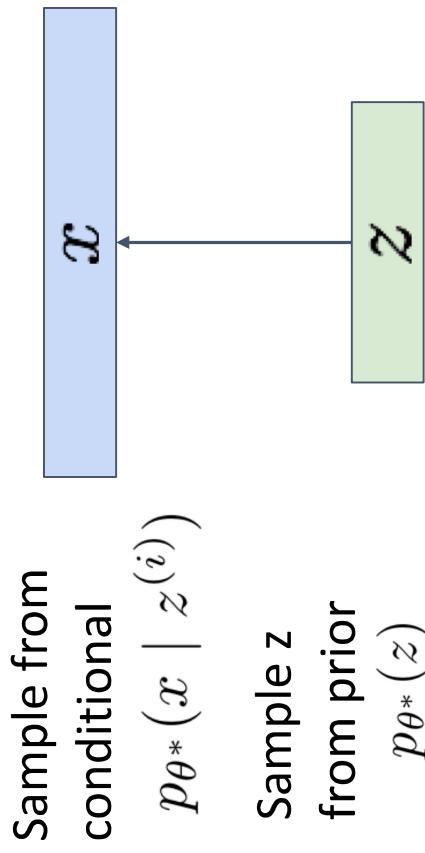
Intuition: \mathbf{x} is an image, \mathbf{z} is latent factors used to generate \mathbf{x} : attributes, orientation, etc.

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features z from raw data
2. Sample from the model to generate new data representation x

After training, sample new data like this:



Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

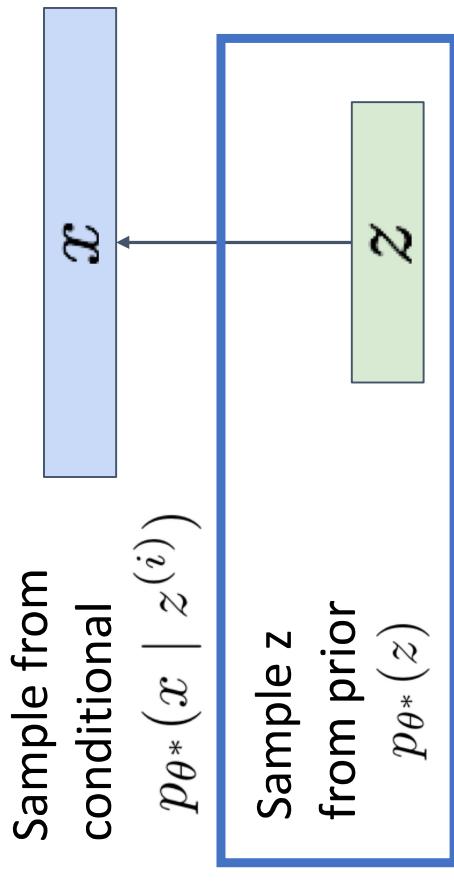
Assume training data $\{x^{(i)}\}_{i=1}^N$ is generated from unobserved (latent) representation z

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features z from raw data
2. Sample from the model to generate new data representation x

After training, sample new data like this:



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

Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

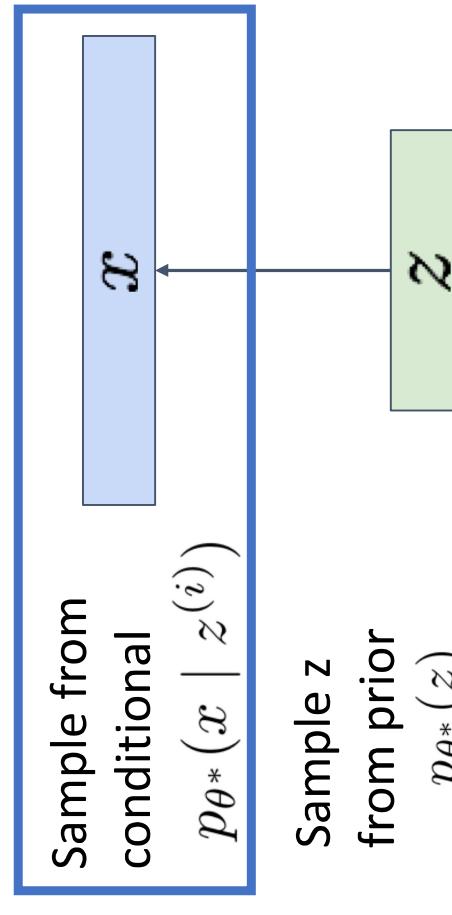
Assume simple prior $p(z)$, e.g. Gaussian

Variational Autoencoders

Probabilistic spin on autoencoders:

1. Learn latent features z from raw data
2. Sample from the model to generate new data representation z

After training, sample new data like this:



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

representation z

Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

Assume simple prior $p(z)$, e.g. Gaussian

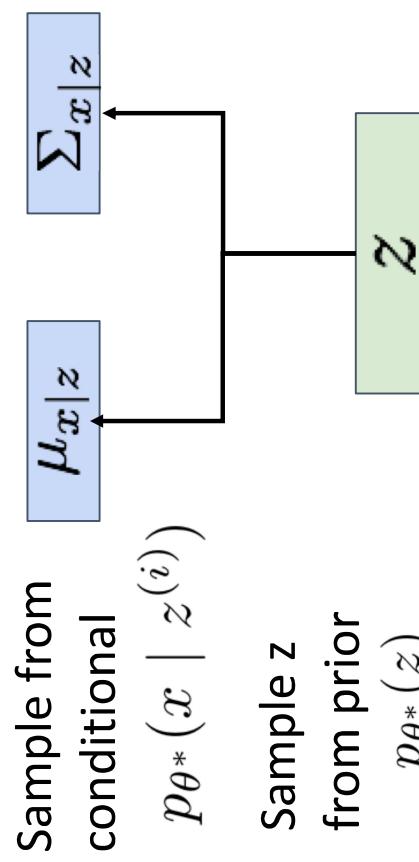
Represent $p(x|z)$ with a neural network
(Similar to **decoder** from autoencoder)

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$



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

Intuition: x is an image, z is latent factors used to generate x : attributes, orientation, etc.

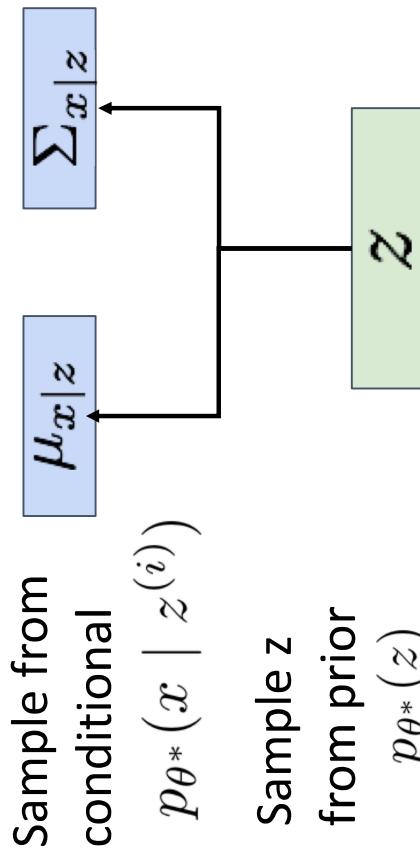
Assume simple prior $p(z)$, e.g. Gaussian
Represent $p(x|z)$ with a neural network
(Similar to **decoder** from autoencoder)

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$



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

How to train this model?

Basic idea: **maximize likelihood of data**

If we could observe the z for each x , then could train a *conditional generative model* $p(x|z)$

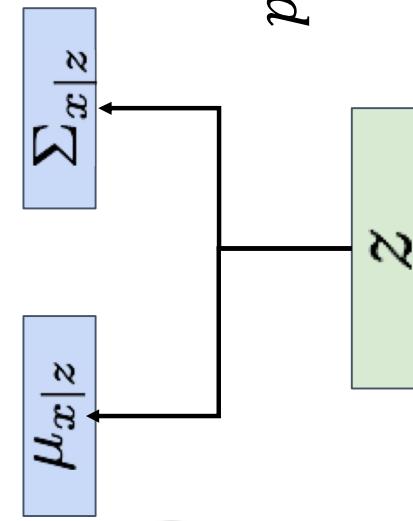
Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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



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

How to train this model?

Basic idea: **maximize likelihood of data**

We don't observe z , so need to marginalize:

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

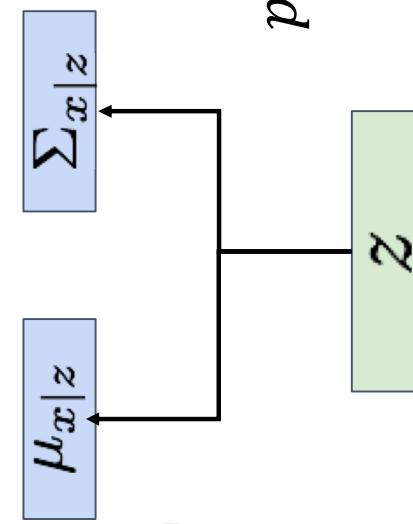
Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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



Basic idea: **maximize likelihood of data**
 $p_{\theta}(x) = \int p_{\theta}(x, z) dz = \int p_{\theta}(x|z)p_{\theta}(z)dz$

We don't observe z , so need to marginalize:

Ok, can compute this with decoder network

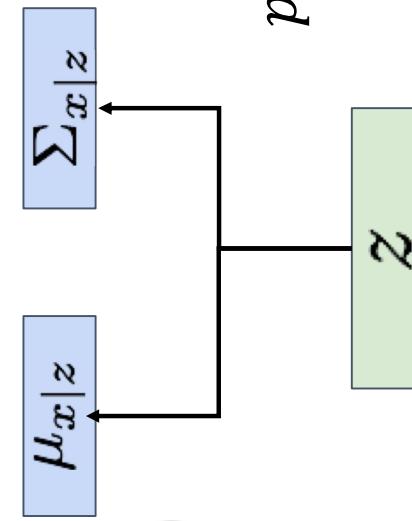
Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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



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

How to train this model?

Basic idea: **maximize likelihood of data**

We don't observe z , so need to marginalize:

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

Ok, we assumed Gaussian prior for z

Variational Autoencoders

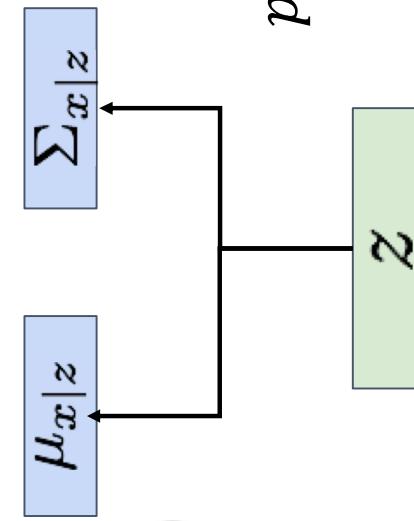
Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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

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



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

How to train this model?

Basic idea: **maximize likelihood of data**

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

We don't observe z , so need to marginalize:

Problem: Impossible to integrate over all z !

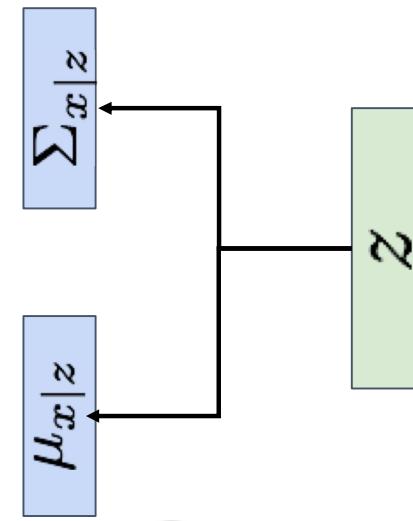
Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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



Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

$$p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)}$$

Recall $p(x, z) = p(x \mid z)p(z) = p(z \mid x)p(x)$

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

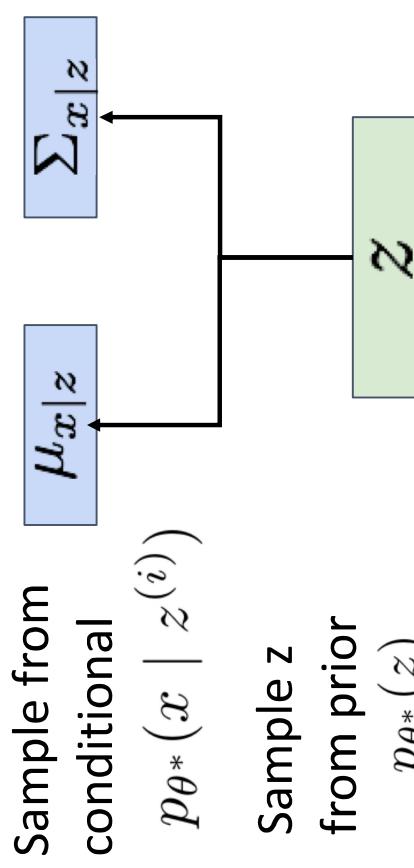
How to train this model?

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$



Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

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

How to train this model?

Basic idea: **maximize likelihood of data**

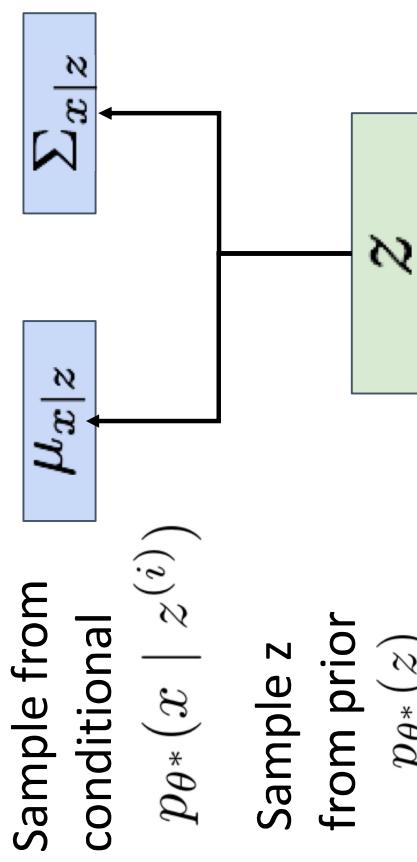
Another idea: Try Bayes' Rule:
$$p_\theta(x) = \frac{p_\theta(x | z)p_\theta(z)}{p_\theta(z | x)}$$
 Ok, compute with decoder network

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$



Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

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

How to train this model?

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:
$$p_\theta(x) = \frac{p_\theta(x | z)p_\theta(z)}{p_\theta(z | x)}$$

Ok, we assumed Gaussian prior

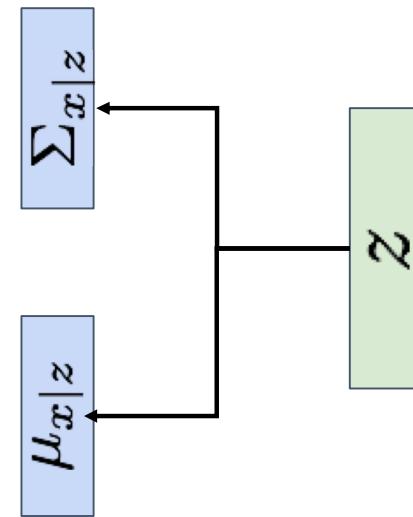
Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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



Recall $p(x, z) = p(x \mid z)p(z) = p(z \mid x)p(x)$

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

How to train this model?

Basic idea: **maximize likelihood of data**

$$p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)}$$

Another idea: Try Bayes' Rule:

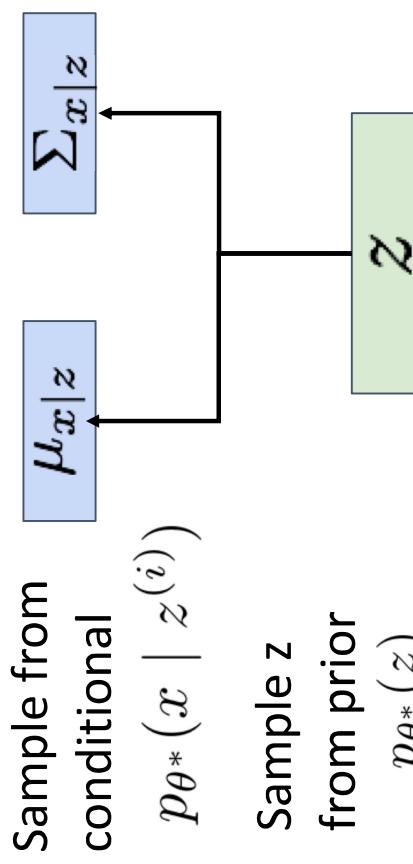
Problem: No way to compute this!

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$



Recall $p(x, z) = p(x | z)p(z) = p(z | x)p(x)$

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

How to train this model?

Basic idea: **maximize likelihood of data**

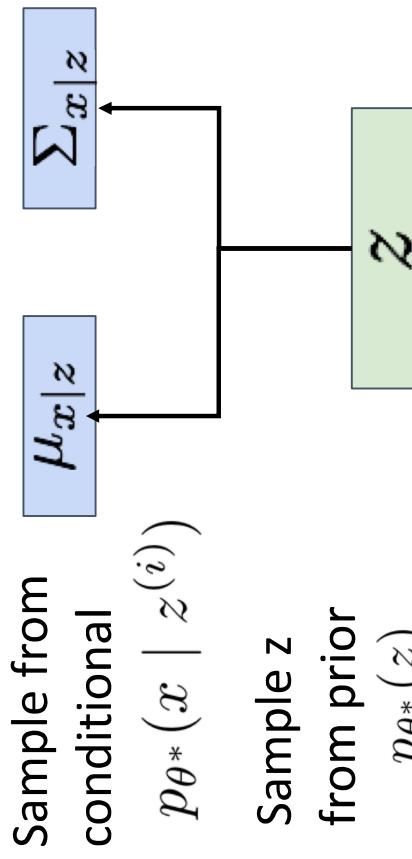
Solution: Train another network (**encoder**) that learns $q_\phi(z | x) \approx p_\theta(z | x)$

Variational Autoencoders

Decoder must be **probabilistic**:

Decoder inputs z , outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$



Recall $p(x, z) = p(x \mid z)p(z) = p(z \mid x)p(x)$

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

How to train this model?

Basic idea: **maximize likelihood of data**

$$p_{\theta}(x) = \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x)} \approx \frac{p_{\theta}(x \mid z)p_{\theta}(z)}{q_{\phi}(z \mid x)}$$

Use **encoder** to compute $q_{\phi}(z \mid x) \approx p_{\theta}(z \mid x)$

Variational Autoencoders

Decoder network inputs

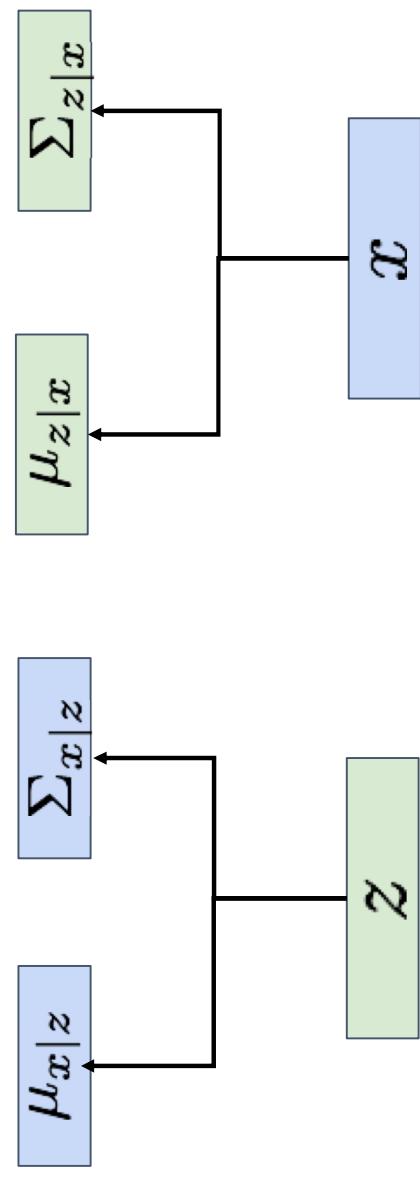
latent code z , gives distribution over data x

Encoder network inputs

data x , gives distribution over latent codes z

$$p_{\theta}(x | z) = N(\mu_x|z, \Sigma_{x|z})$$

$$q_{\phi}(z | x) = N(\mu_z|x, \Sigma_{z|x})$$



If we can ensure that
 $q_{\phi}(z | x) \approx p_{\theta}(z | x)$,

$$p_{\theta}(x) \approx \frac{p_{\theta}(x | z)p(z)}{q_{\phi}(z | x)}$$

Idea: Jointly train both encoder and decoder

Variational Autoencoders

$$\log p_{\theta}(x) = \log \frac{p_{\theta}(x \mid z)p(z)}{p_{\theta}(z \mid x)}$$

Bayes' Rule

Variational Autoencoders

$$\log p_{\theta}(x) = \log \frac{p_{\theta}(x | z)p(z)}{p_{\theta}(z | x)} = \log \frac{p_{\theta}(x | z)p(z)q_{\phi}(z | x)}{p_{\theta}(z | x)q_{\phi}(z | x)}$$

Multiply top and bottom by $q_{\Phi}(z | x)$

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= \log p_\theta(x|z) - \log \frac{q_\phi(z|x)}{p(z)} + \log \frac{q_\phi(z|x)}{p_\theta(z|x)}\end{aligned}$$

Split up using rules for logarithms

Variational Autoencoders

$$\begin{aligned} \log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)q_\phi(z | x)} \\ &= \log p_\theta(x | z) - \log \frac{q_\phi(z | x)}{p(z)} + \log \frac{q_\phi(z | x)}{p_\theta(z | x)} \end{aligned}$$

Split up using rules for logarithms

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= \log p_\theta(x|z) - \log \frac{q_\phi(z|x)}{p(z)} + \log \frac{q_\phi(z|x)}{p_\theta(z|x)}\end{aligned}$$

We can wrap in an expectation since it doesn't depend on z

$$\log p_\theta(x) = E_{z \sim q_\phi(z|x)} [\log p_\theta(x)]$$

Variational Autoencoders

$$\log p_\theta(x) = \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)q_\phi(z | x)}{p_\theta(z | x)q_\phi(z | x)}$$

$$= E_z [\log p_\theta(x | z)] - E_z \left[\log \frac{q_\phi(z | x)}{p(z)} \right] + E_z \left[\log \frac{q_\phi(z | x)}{p_\theta(z | x)} \right]$$

We can wrap in an expectation since it doesn't depend on z

$$\log p_\theta(x) = E_{z \sim q_\phi(z | x)} [\log p_\theta(x)]$$

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x|z)p(z)}{p_\theta(z|x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= E_z[\log \frac{q_\phi(z|x)}{p(z)}] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right] \\ &= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z)) + D_{KL}(q_\phi(z|x), p_\theta(z|x))\end{aligned}$$

Data reconstruction

Variational Autoencoders

$$\begin{aligned} \log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= E_z[\log \frac{q_\phi(z|x)}{p(z)}] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] + E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right] \\ &= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}\left(q_\phi(z|x), p(z)\right) + D_{KL}(q_\phi(z|x), p_\theta(z|x)) \end{aligned}$$

KL divergence between prior, and
samples from the encoder network

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= E_z [\log \frac{q_\phi(z|x)}{p(z)}] - E_z \left[\log \frac{q_\phi(z|x)}{p(z)} \right] + E_z \left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)} \right] \\ &= E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z)) + D_{KL} (q_\phi(z|x), p_\theta(z|x))\end{aligned}$$

KL divergence between encoder
and posterior of decoder

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x | z)p(z)}{p_\theta(z | x)} = \log \frac{p_\theta(x | z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= E_z[\log \frac{q_\phi(z|x)}{p(z)}] - E_z\left[\log \frac{q_\phi(z|x)}{p(z)}\right] \\ &= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z)) + D_{KL}(q_\phi(z|x), p_\theta(z|x))\end{aligned}$$

KL is ≥ 0 , so dropping this term gives a
lower bound on the data likelihood:

Variational Autoencoders

$$\begin{aligned}\log p_\theta(x) &= \log \frac{p_\theta(x|z)p(z)}{p_\theta(z|x)} = \log \frac{p_\theta(x|z)p(z)q_\phi(z|x)}{p_\theta(z|x)q_\phi(z|x)} \\ &= E_z[\log \frac{q_\phi(z|x)}{p(z)}] - E_z\left[\log \frac{q_\phi(z|x)}{p_\theta(z|x)}\right] \\ &= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z|x)) + D_{KL}(q_\phi(z|x), p_\theta(z|x)) \\ \log p_\theta(x) &\geq E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}(q_\phi(z|x), p(z))\end{aligned}$$

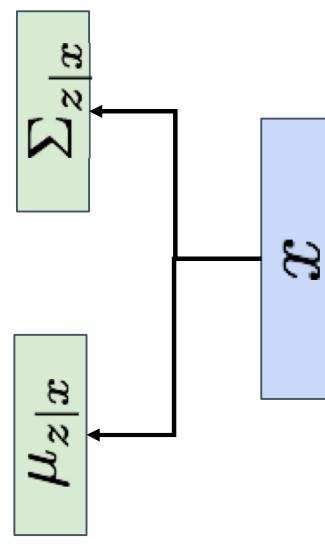
Variational Autoencoders

Jointly train **encoder** q and **decoder** p to maximize
the **variational lower bound** on the data likelihood

$$\log p_\theta(x) \geq E_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL} (q_\phi(z|x), p(z))$$

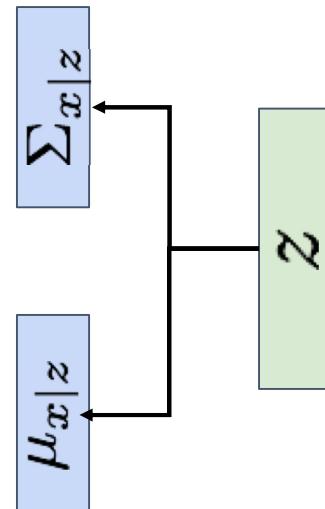
Encoder Network

$$q_\phi(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$$



Decoder Network

$$p_\theta(x | z) = N(\mu_{x|z}, \Sigma_{x|z})$$



Next Time:
Generative Models, part 2

More Variational Autoencoders,
Generative Adversarial Networks