# Lecture 9: Generative Models

# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data**: (x, y)

x is data, y is label

**Goal**: Learn a *function* to map x -> y

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

Classification



Cat

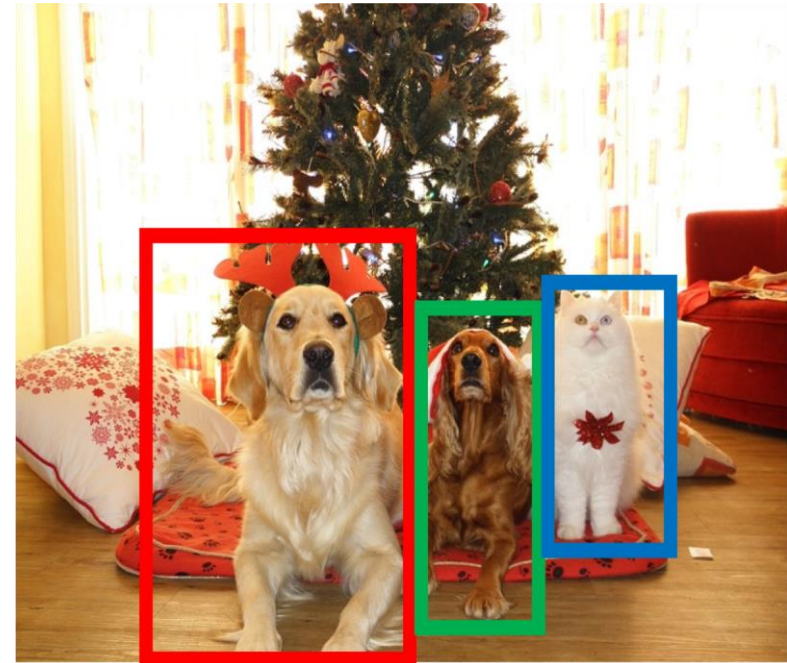# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data**: (x, y)

x is data, y is label

**Goal**: Learn a *function* to map x -> y

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

Object Detection



**DOG**, **DOG**, **CAT**

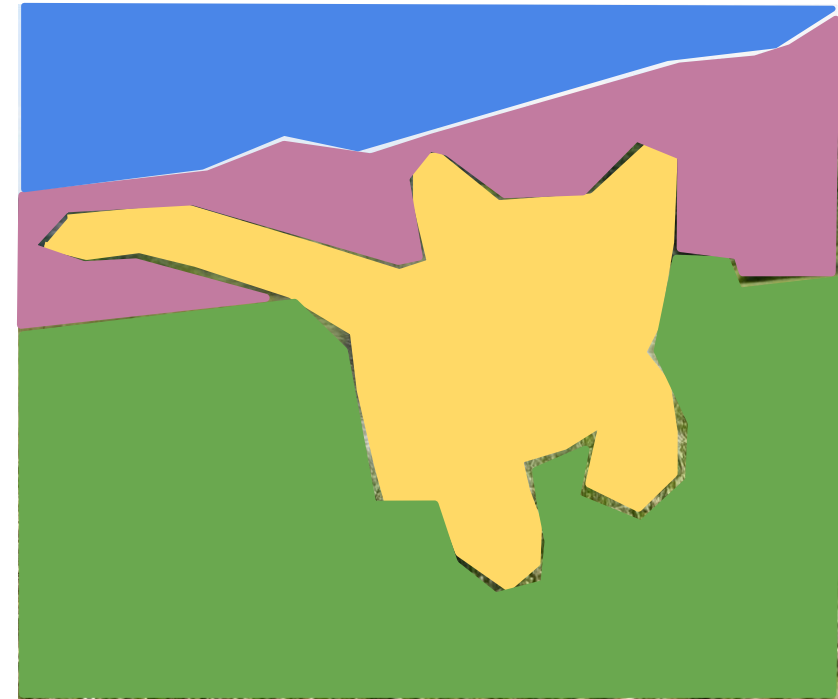# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data**: (x, y)

x is data, y is label

**Goal**: Learn a *function* to map x -> y

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

Semantic Segmentation



**GRASS**, **CAT**, **TREE**, **SKY**

# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data**: (x, y)

x is data, y is label

**Goal**: Learn a *function* to map x -> y

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

Image captioning



*A cat sitting on a suitcase on the floor*

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data**: (x, y)

x is data, y is label

**Goal**: Learn a *function* to map x -> y

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

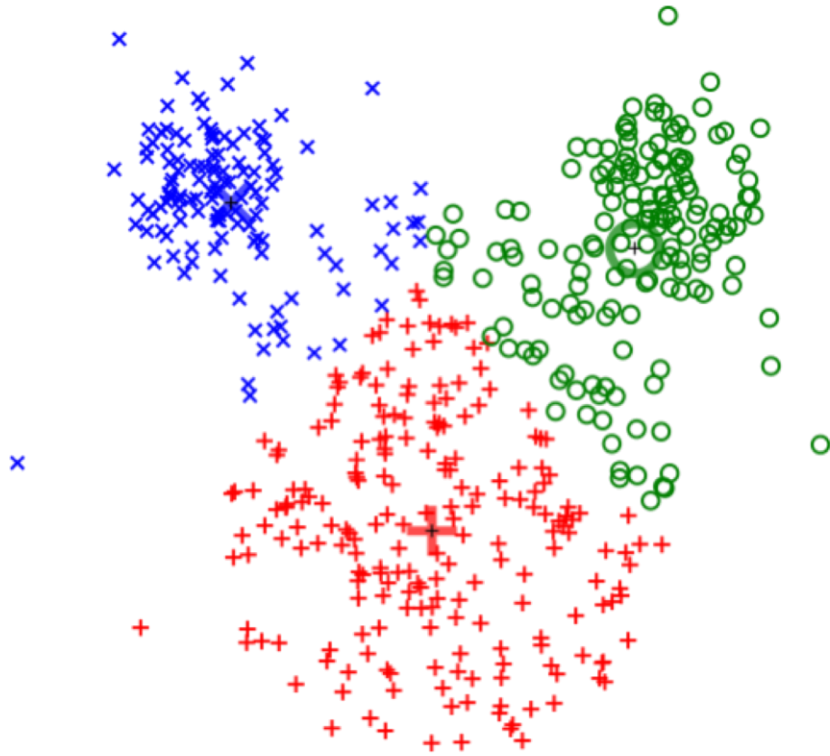## Unsupervised Learning

**Data**: x

Just data, no labels!

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

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

# Supervised vs Unsupervised Learning

## Clustering
## (e.g. K-Means)

**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

## Dimensionality Reduction
## (e.g. Principal Components Analysis)



original data space

PCA →

component space

3D  →  2D

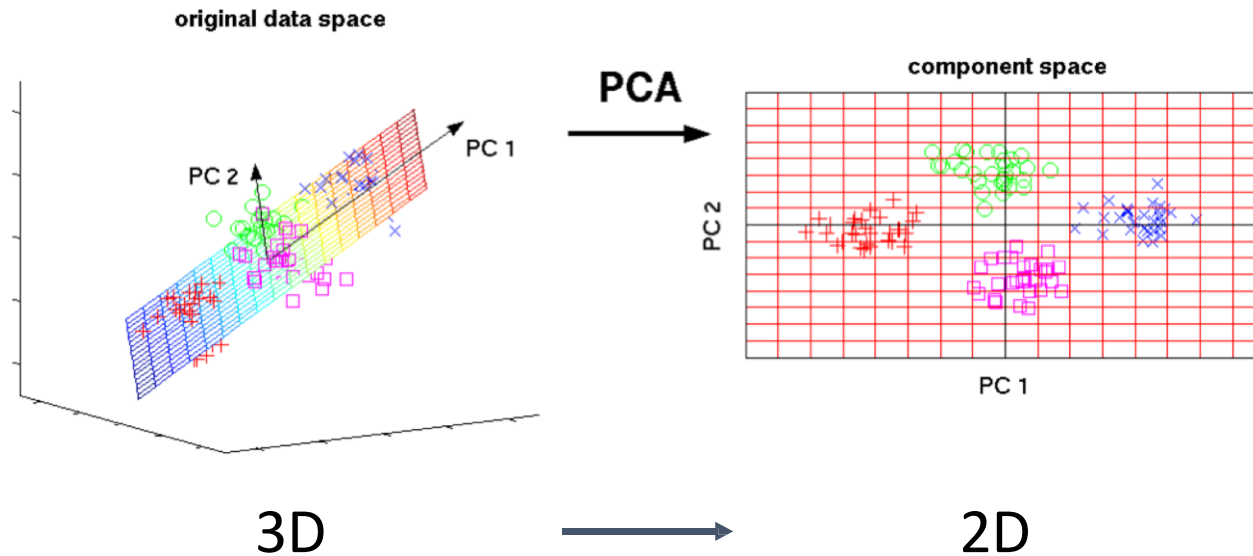**Unsupervised Learning**

**Data**: x

Just data, no labels!

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

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

0

# Supervised vs Unsupervised Learning

## Feature Learning
## (e.g. autoencoders)

L2 Loss function:
$$\|x - \hat{x}\|^2$$

Reconstructed input data — $\hat{x}$

Decoder

**Features** — $z$

Encoder

Input data — $x$

Reconstructed data

**Encoder**: 4-layer conv
**Decoder**: 4-layer upconv
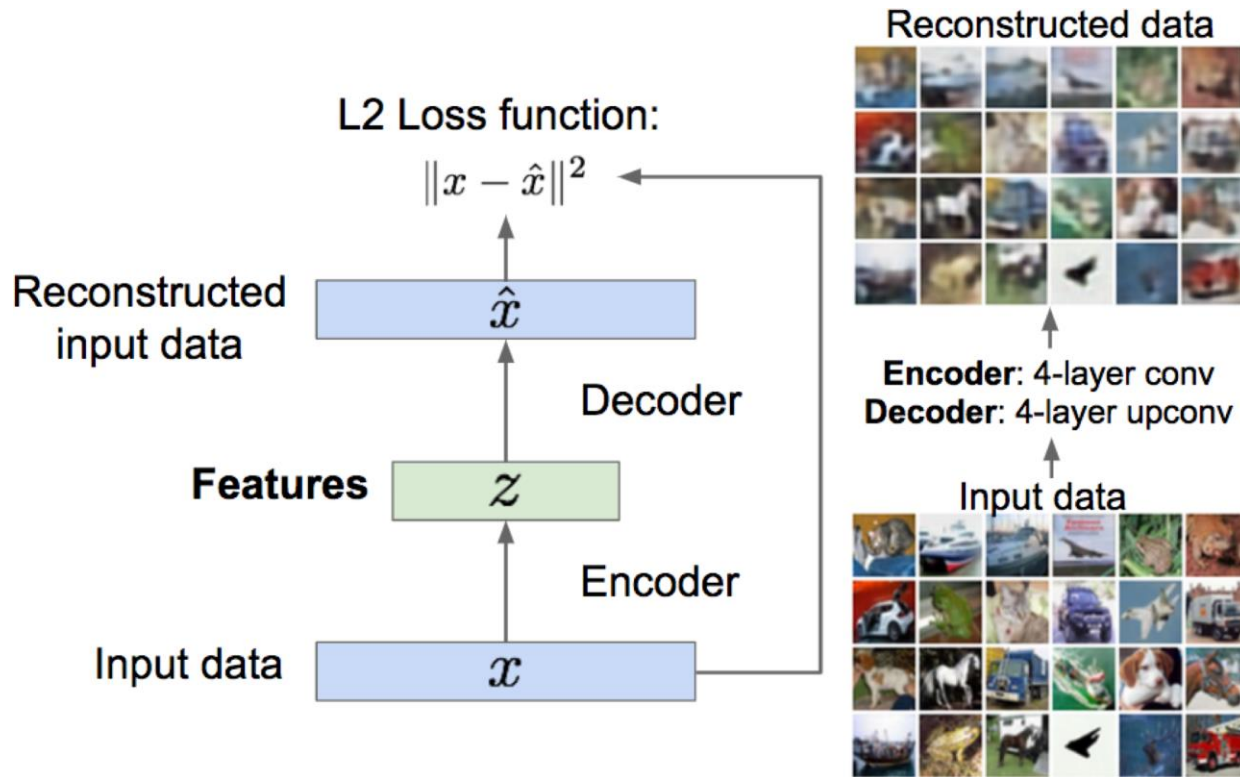
Input data

**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

## Density Estimation


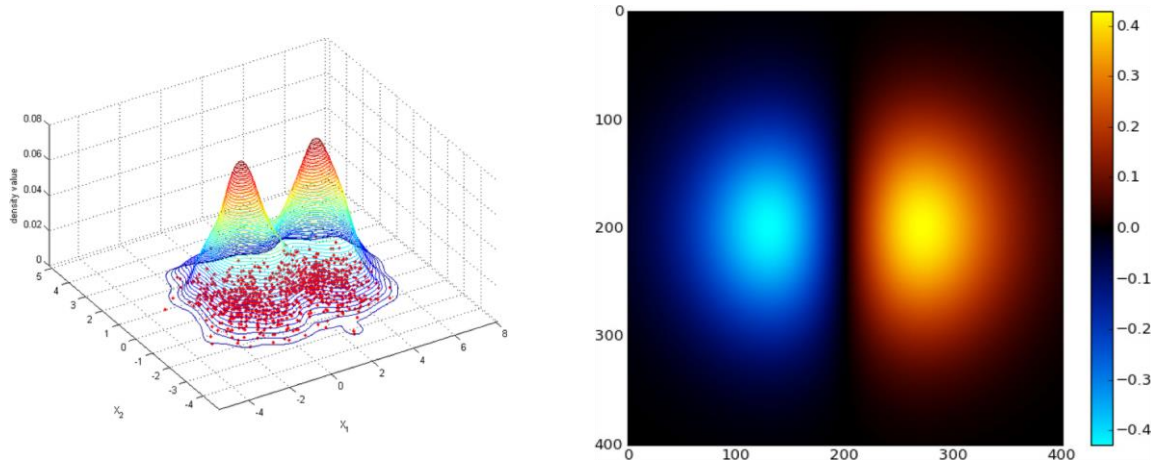
**Unsupervised Learning**

**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

## Supervised Learning

**Data**: (x, y)

x is data, y is label

**Goal**: Learn a *function* to map x -> 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.

# 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)

**Data**: x



**Label**: y

Cat

# Discriminative vs Generative Models

## Probability Recap:

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

**Data**: x



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

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

**Label**: y
## Cat

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

Density functions are **normalized**:

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

Different values of x **compete** for density

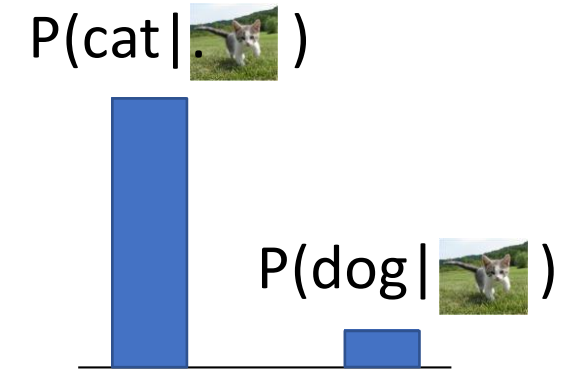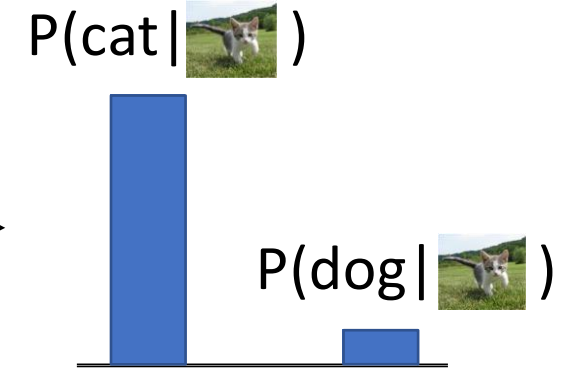# Discriminative vs Generative Models

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

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

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

**Data**: x



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

P(cat |  )

P(dog |  )

Density functions are **normalized**:

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

Different values of x **compete** for density

# Discriminative vs Generative Models

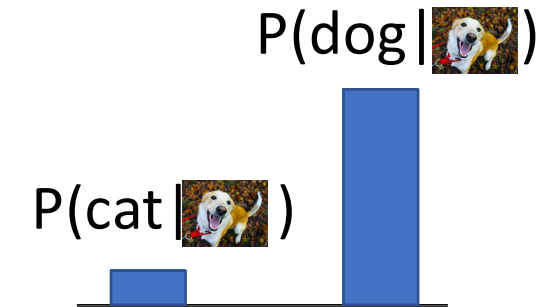**Discriminative Model:**
Learn a probability
distribution p(y|x)

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

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



P(cat| )

P(dog| )

P(dog| )

P(cat| )

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)

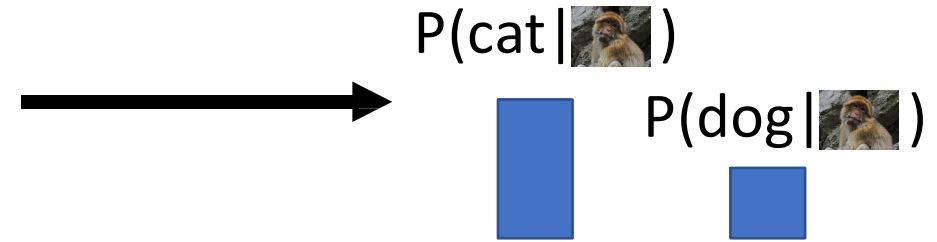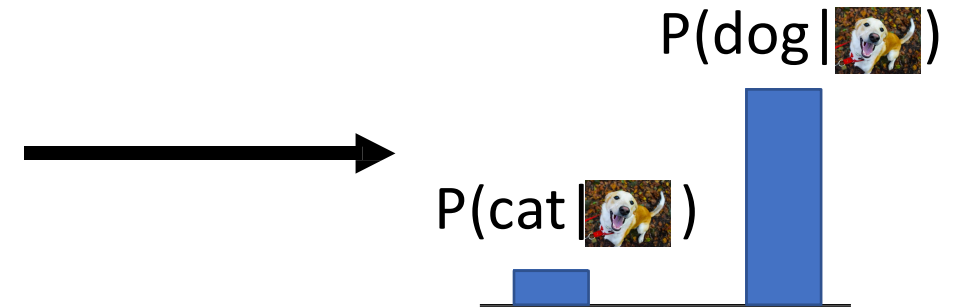**Generative Model**:
Learn a probability
distribution p(x)

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



P(cat| ) 

P(dog| )

P(dog| )

P(cat| )

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



P(cat| )
P(dog| )

P(dog| )
P(cat| )

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

# Discriminative vs Generative Models

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

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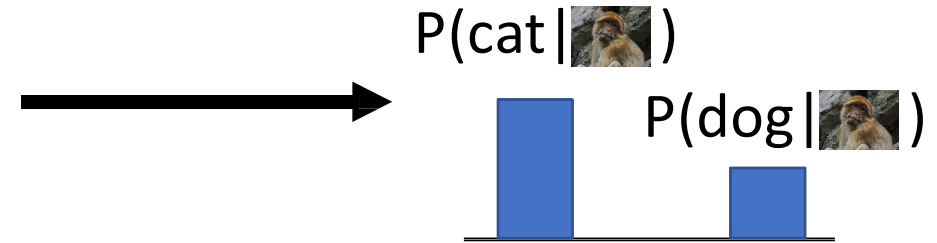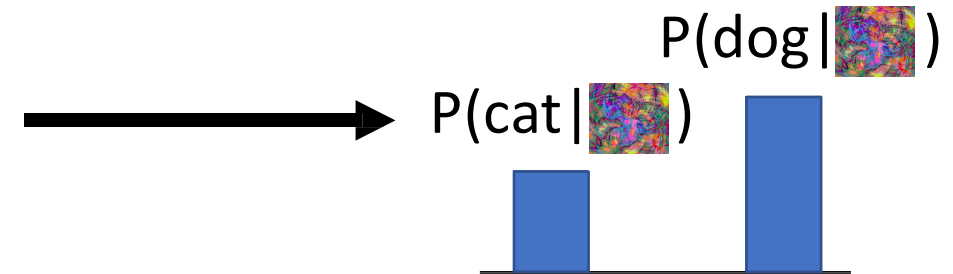Learn a probability
distribution p(x)

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



Generative model: All possible images compete
with each other for probability mass
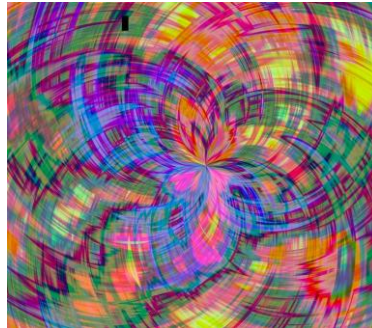
# Discriminative vs Generative Models

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

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

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



$P(\text{cat})$  $P(\text{dog})$  $P(\text{monkey})$  $P(\text{abstract})$

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

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

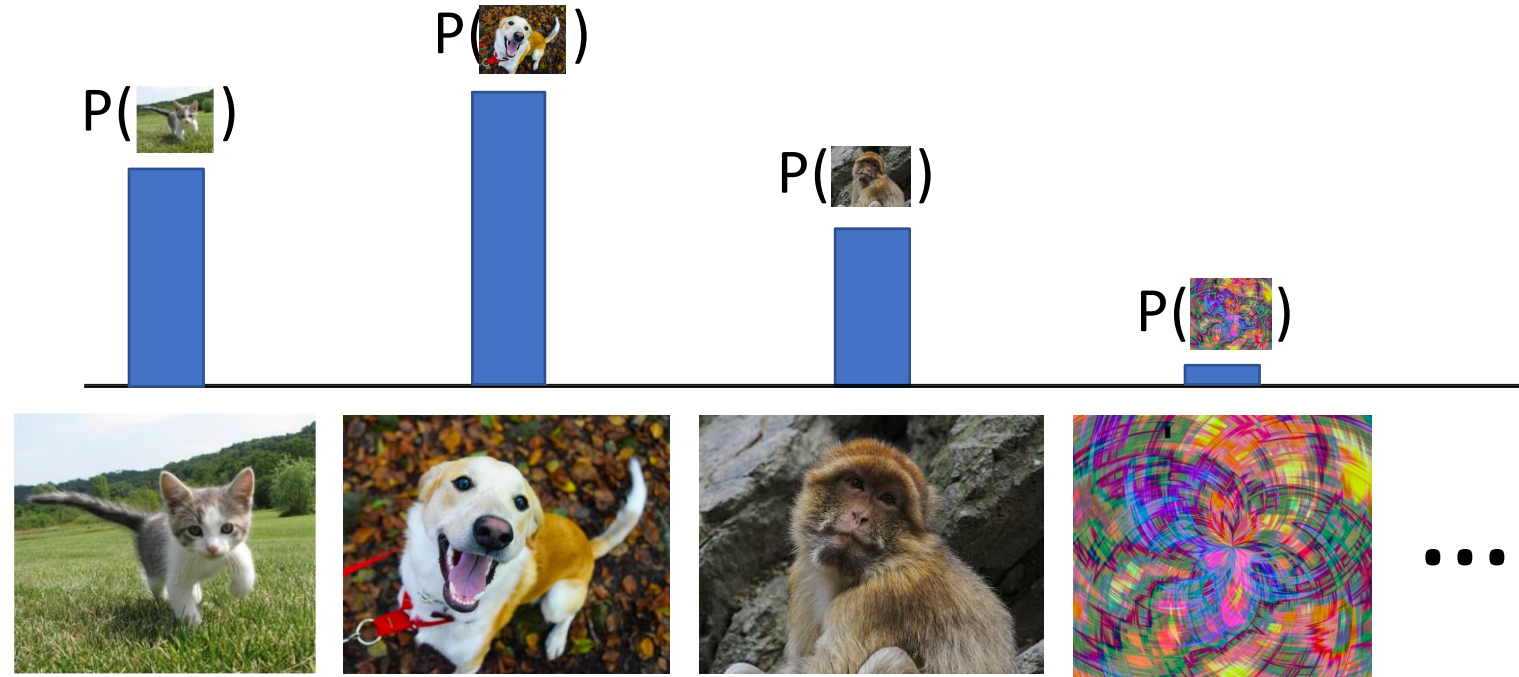# Discriminative vs Generative Models

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

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

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



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

Model can "reject" unreasonable inputs by assigning them small values

22
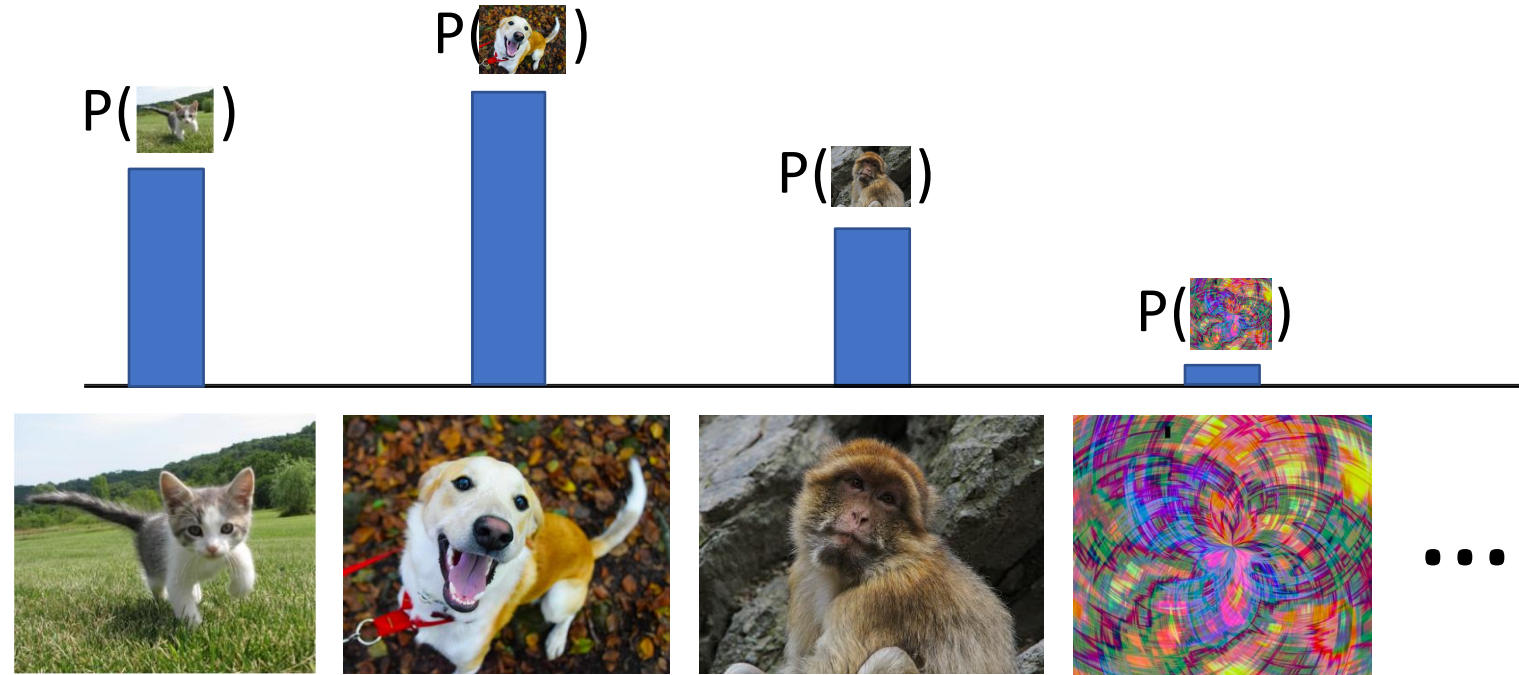
# Discriminative vs Generative Models

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

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

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



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

Recall **Bayes' Rule:**

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

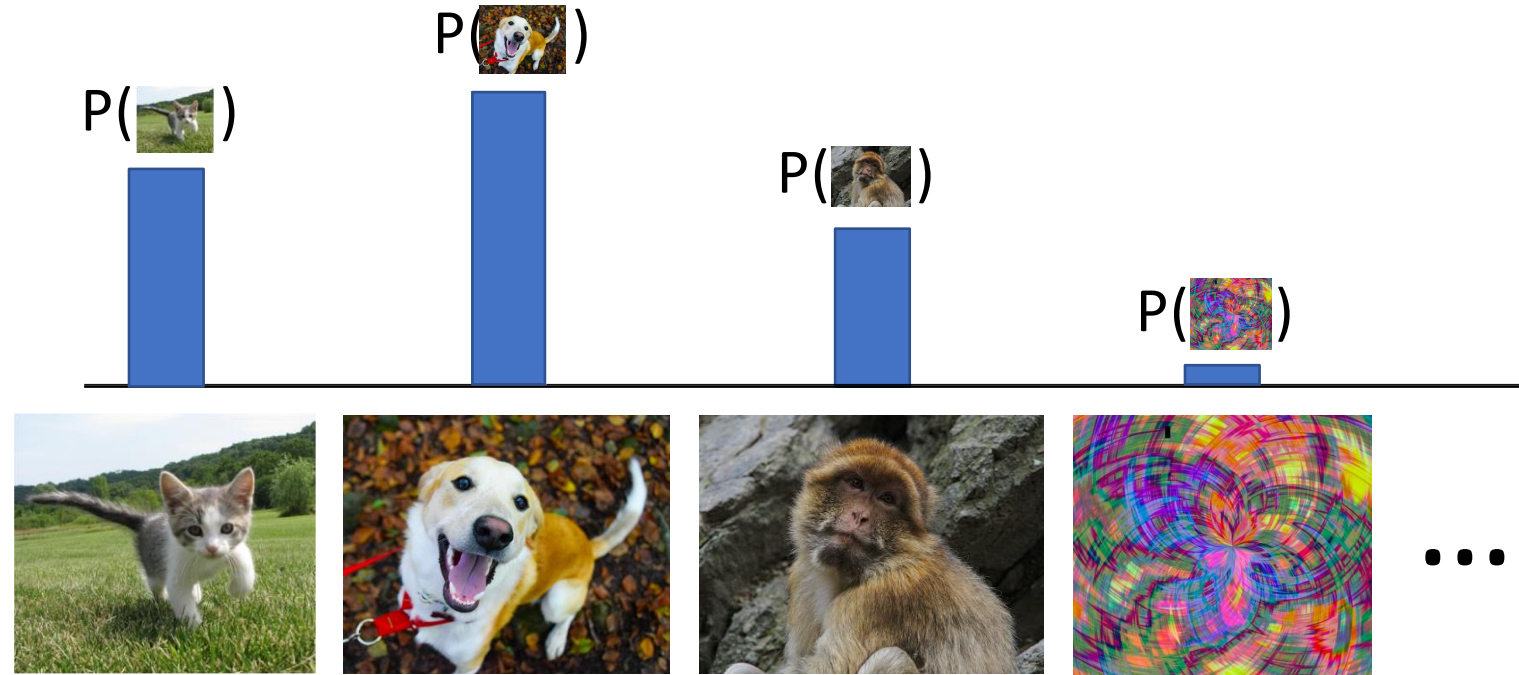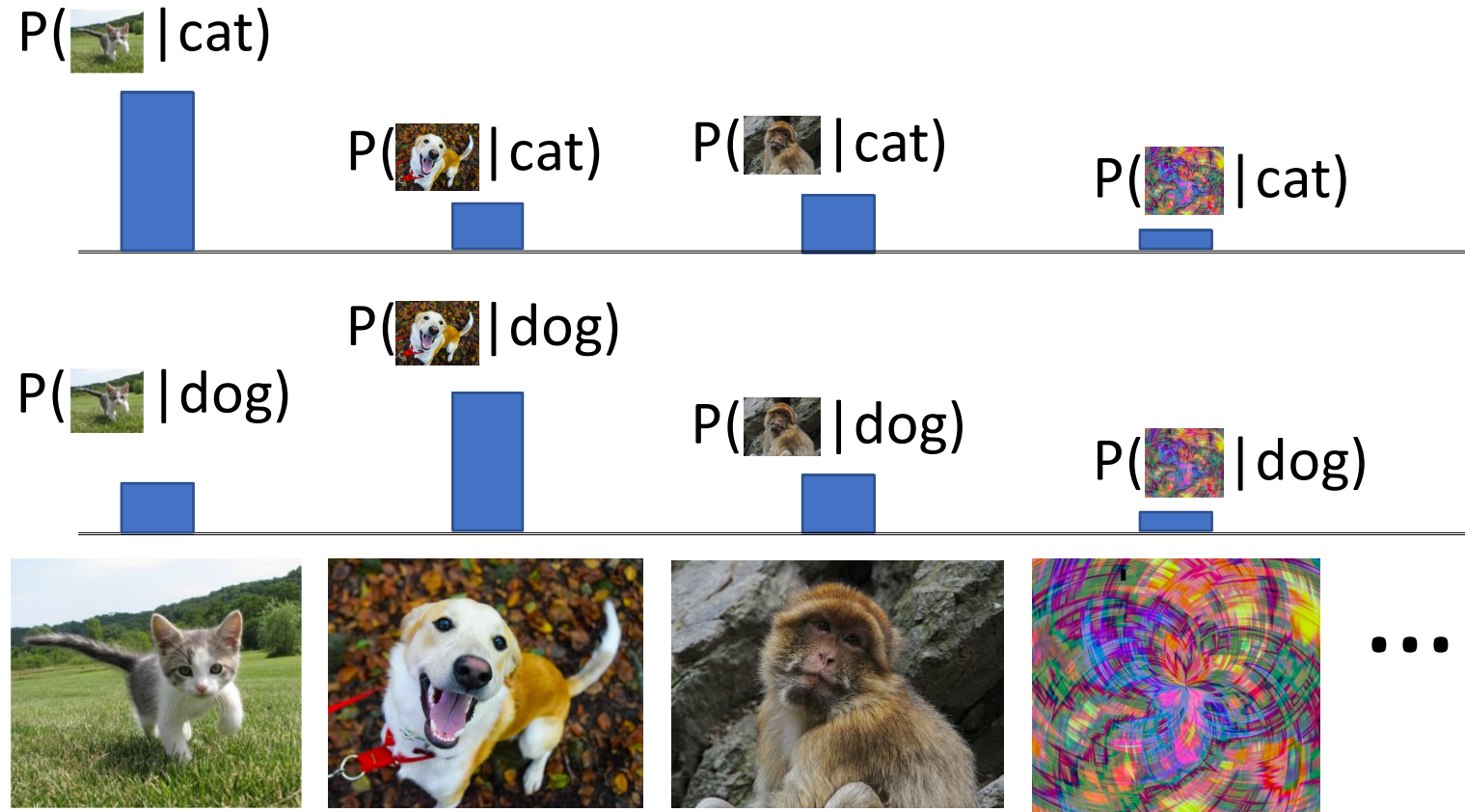# Discriminative vs Generative Models

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

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

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

Recall **Bayes' Rule:**

Discriminative Model

(Unconditional) Generative Model

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

Conditional Generative Model

Prior over labels

We can build a conditional generative model from other components!

# What can we do with a discriminative model?

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

⟶ Assign labels to data
Feature learning (with labels)

- **Generative    Model**:
Learn    a    probability
distribution p(x)

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

# What can we do with a generative model?

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

→ Assign labels to data
Feature learning (with labels)

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

→ Detect outliers
Feature learning (without labels)
**Sample** to generate new data

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

# What can we do with a generative model?

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

→ Assign labels to data
Feature learning (supervised)

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

→ Detect outliers
Feature learning (unsupervised)
Sample to **generate** new data

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

→ Assign labels, while rejecting outliers!
Generate new data conditioned on input labels

# Taxonomy of Generative Models



Generative models

# Taxonomy of Generative Models

Model can
compute p(x)

**Generative models**

Model does not explicitly
compute p(x), but can
sample from p(x)

Explicit density

Implicit density

# Taxonomy of Generative Models



Model can
compute p(x)

**Generative models**

Model does not explicitly
compute p(x), but can
sample from p(x)

Explicit density

Implicit density

Can compute
approximation to p(x)

Tractable density

Approximate density

Can compute p(x)
- Autoregressive
- NADE / MADE
- NICE / RealNVP
- Glow
- Ffjord

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

# Taxonomy of Generative Models

Model can compute p(x)

**Generative models**

Model does not explicitly compute p(x), but can sample from p(x)

Explicit density

Implicit density

Can compute approximation to p(x)

Tractable density

Approximate density

Can compute p(x)
- Autoregressive
- NADE / MADE
- NICE / RealNVP
- Glow
- Ffjord

Variational

Markov Chain

Variational Autoencoder

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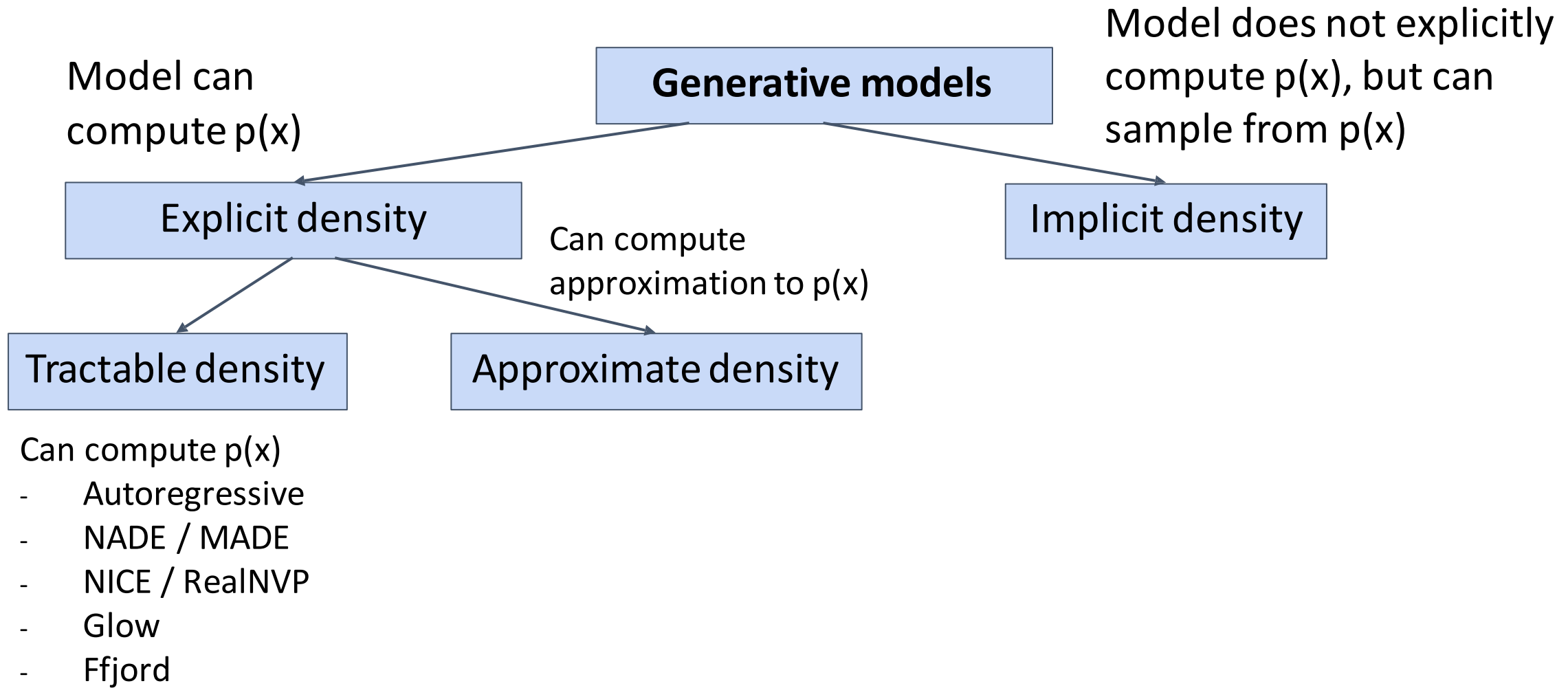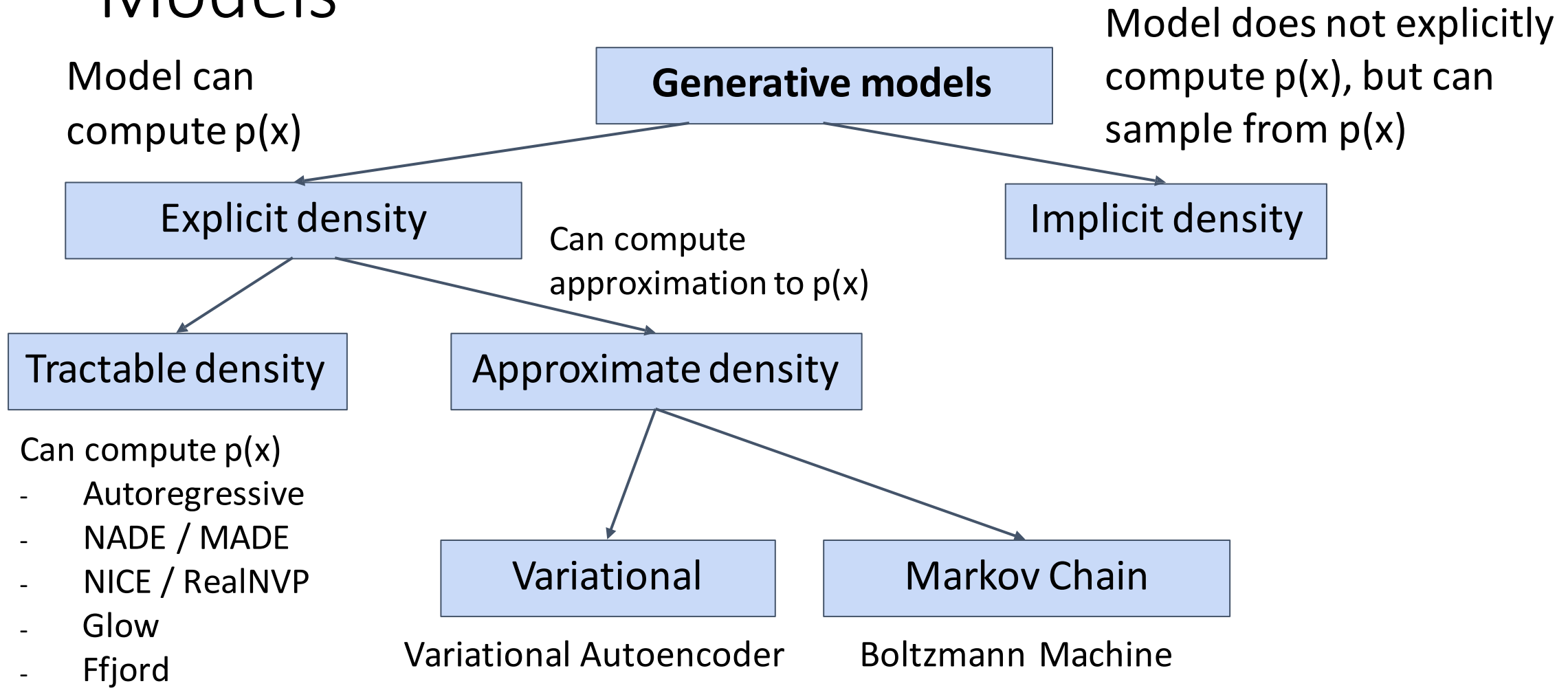
Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

# Taxonomy of Generative Models

Model can compute p(x)

**Generative models**

Model does not explicitly compute p(x), but can sample from p(x)

Explicit density

Can compute approximation to p(x)

Implicit density

Tractable density

Approximate density

Markov Chain

Direct

Can compute p(x)
- Autoregressive
- NADE / MADE
- NICE / RealNVP
- Glow
- Ffjord

GSN

Generative Adversarial Networks (GANs)

Variational

Markov Chain

Variational Autoencoder

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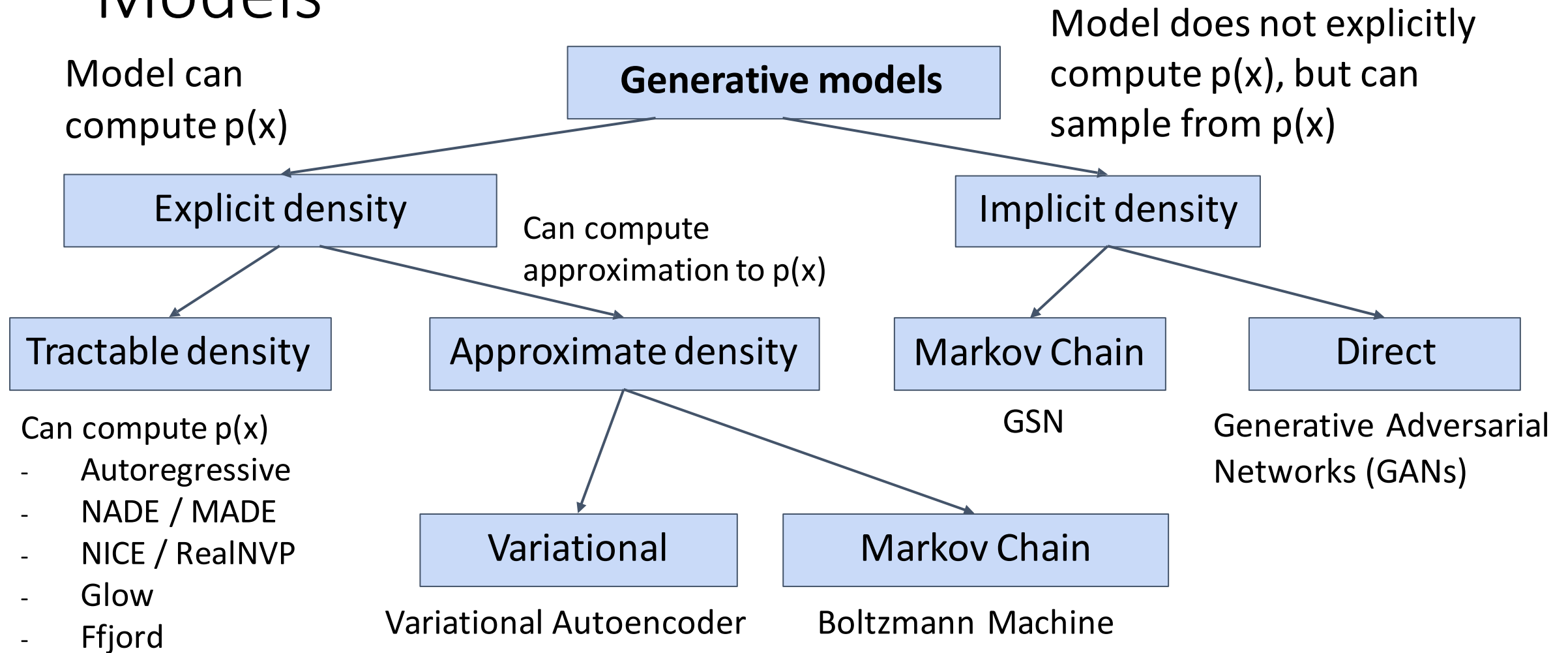
Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

# Taxonomy of Generative Models

Model can
compute p(x)

**Generative models**

Model does not explicitly
compute p(x), but can
sample from p(x)

Explicit density

Can compute
approximation to p(x)

Implicit density

Tractable density

Approximate density

Markov Chain

Direct

Can compute p(x)
- Autoregressive
- NADE / MADE
- NICE / RealNVP
- Glow
- Ffjord

GSN

Generative Adversarial
Networks (GANs)

Variational

Markov Chain

Variational Autoencoder

Boltzmann Machine
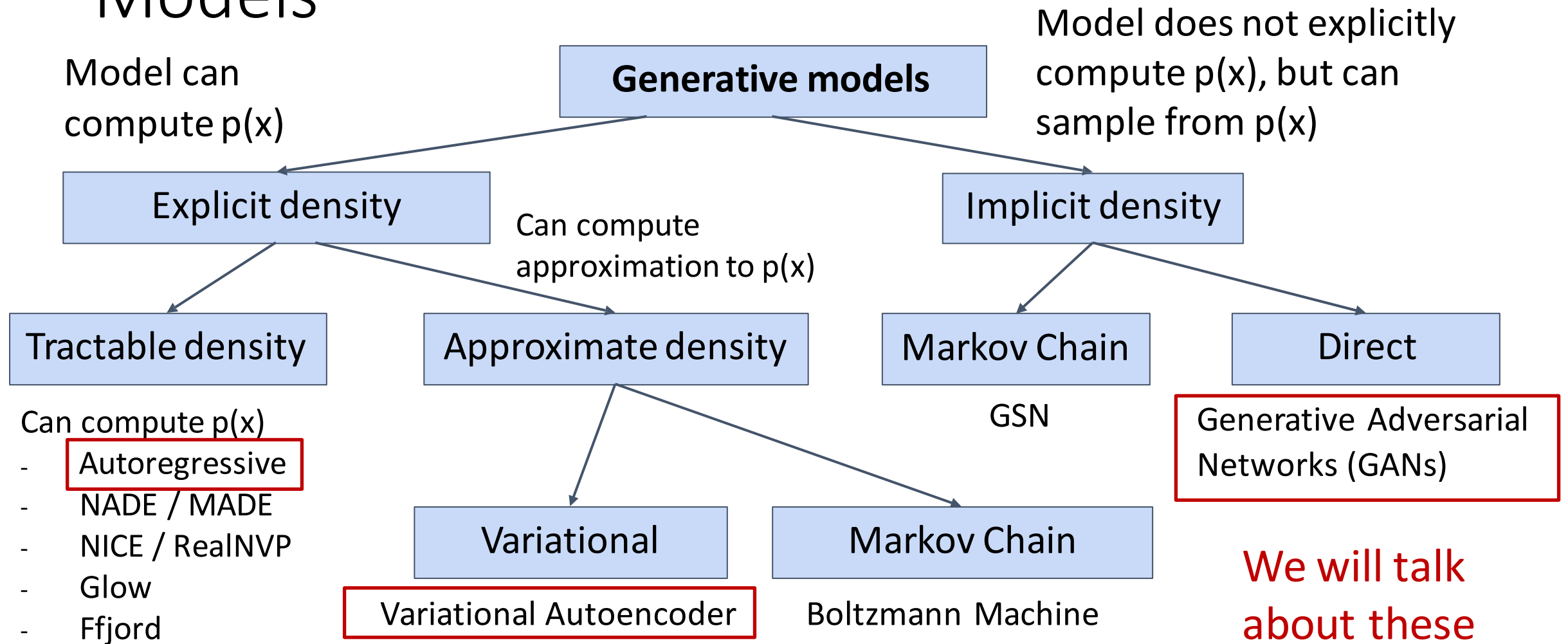
We will talk
about these

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

# Variational Autoencoders

# Variational Autoencoders

Typical density models explicitly parameterize density function with a neural network, so we can train to maximize likelihood of training data:

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

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

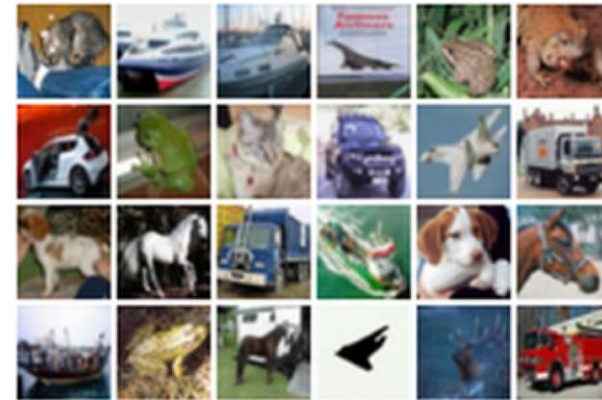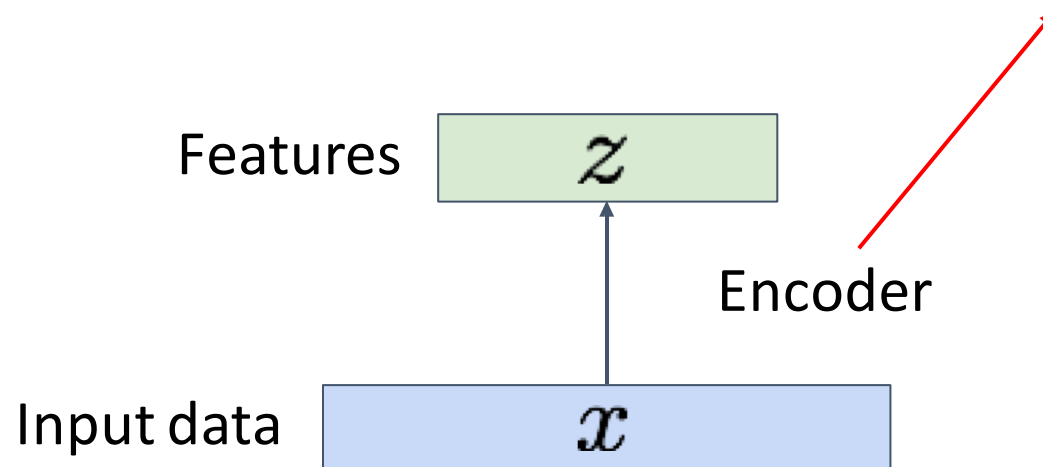# Variational Autoencoders

# (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

Features $\quad \boxed{z}$

Encoder

Input data $\quad \boxed{x}$

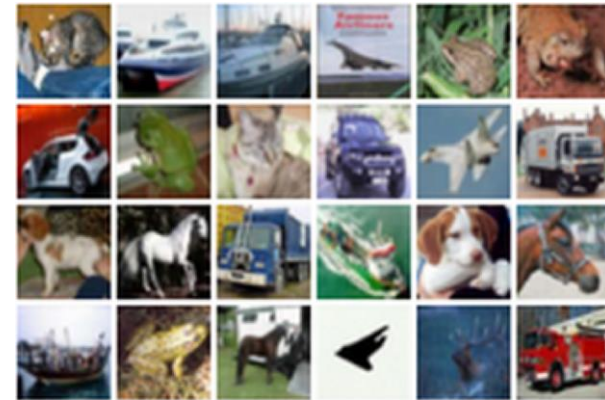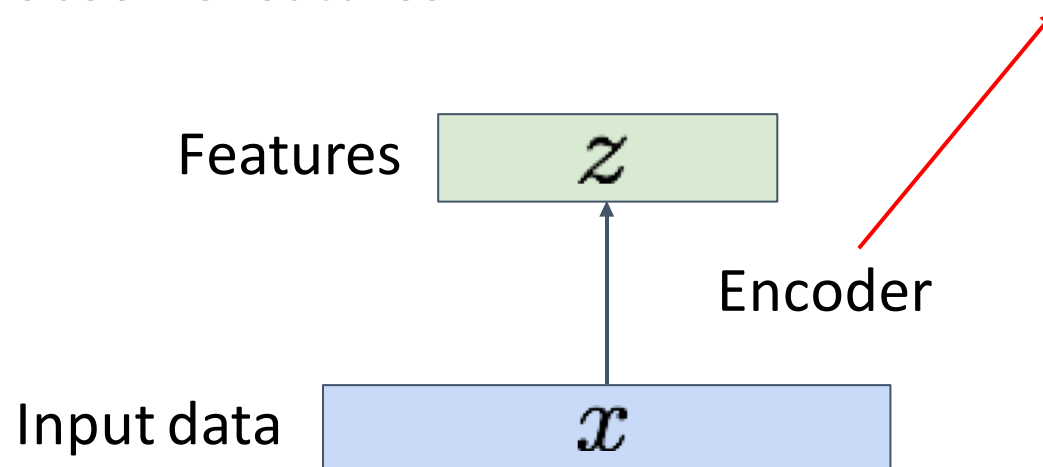# (Regular, non-variational) Autoencoders

**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
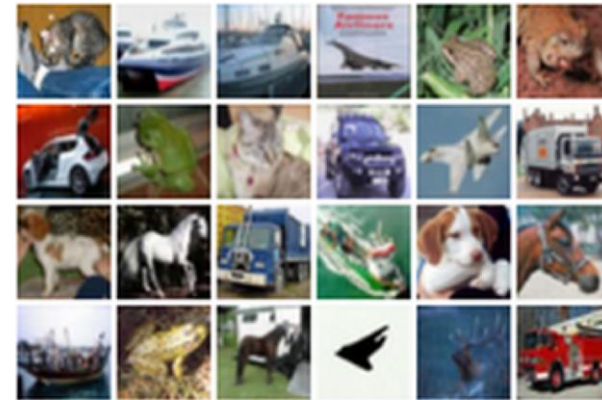
Features $z$

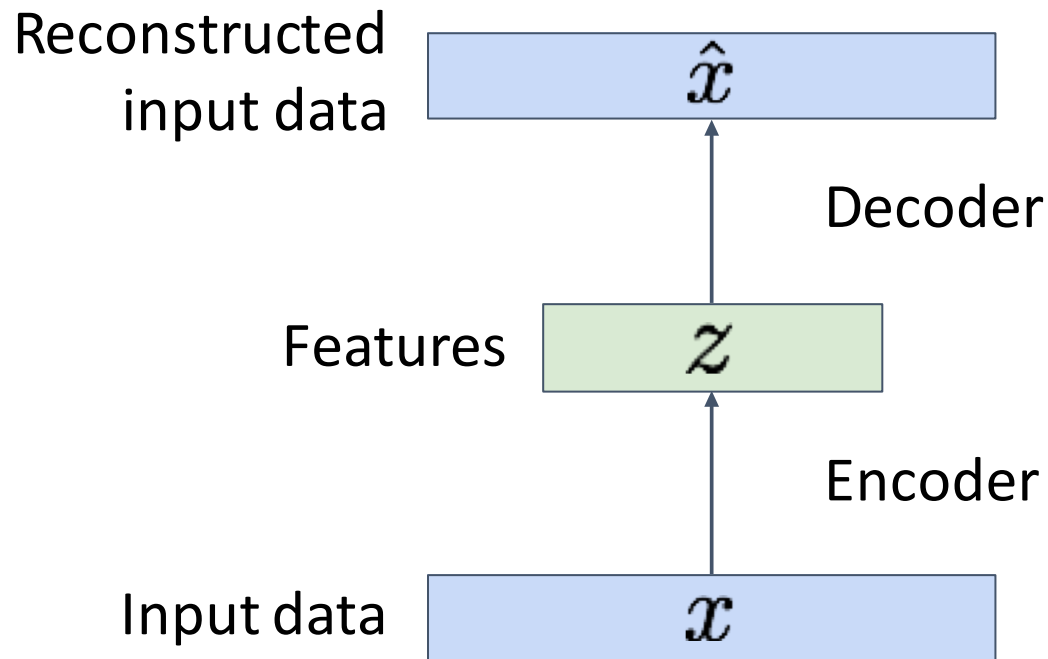Input data $x$

Encoder



Input Data

# (Regular, non-variational) Autoencoders

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

**Idea**: Use the features to <u>reconstruct</u> the input data with a **decoder**
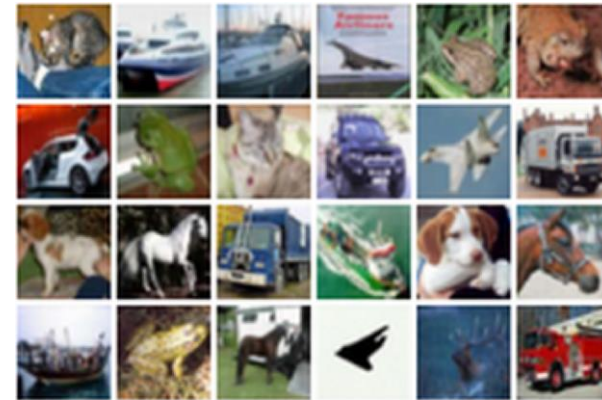
"Autoencoding" = encoding itself

Reconstructed input data

$$\hat{x}$$

Decoder

Features $\quad z$

Encoder

Input data $\quad x$



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

$\hat{x}$

Decoder

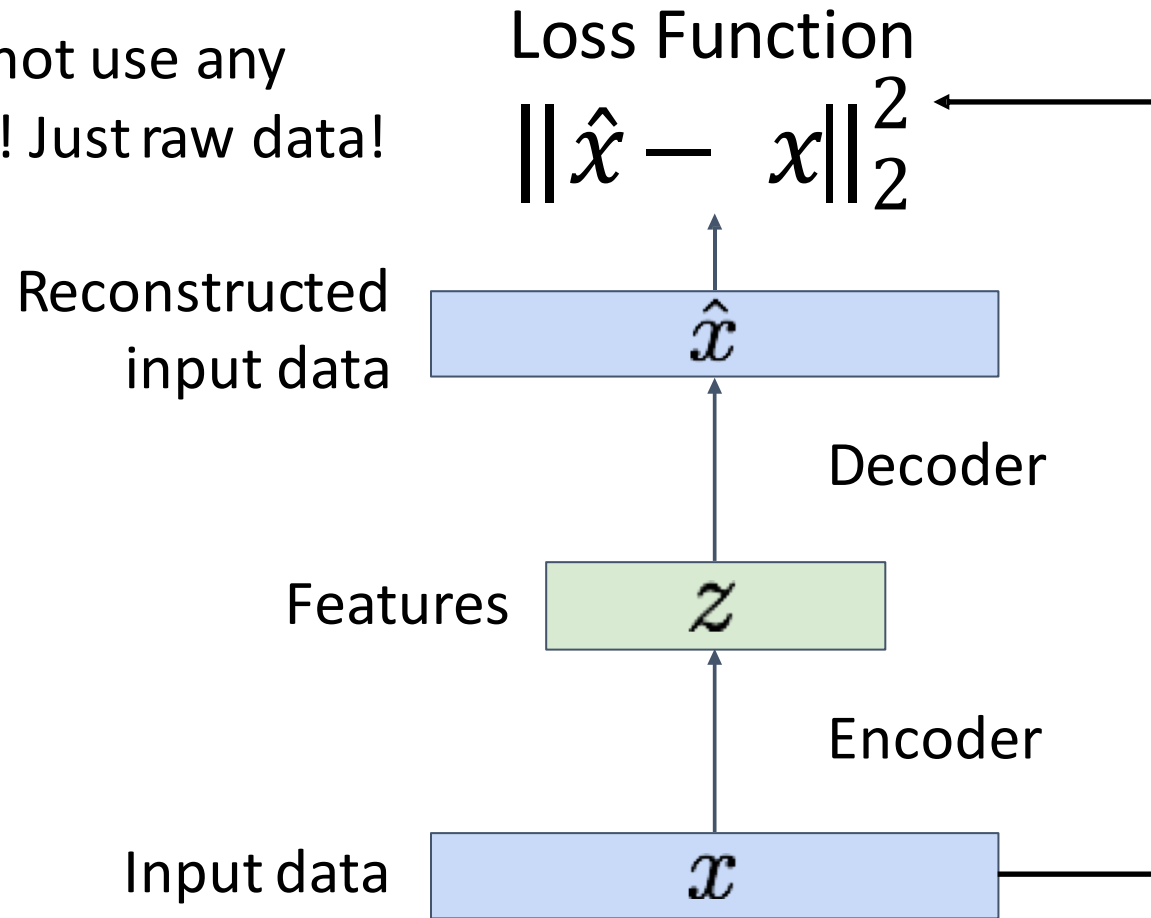Features

$z$

Encoder

Input data

$x$



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: $\hat{x}$

Decoder

Features: $z$

Encoder

Input data: $x$

Reconstructed data
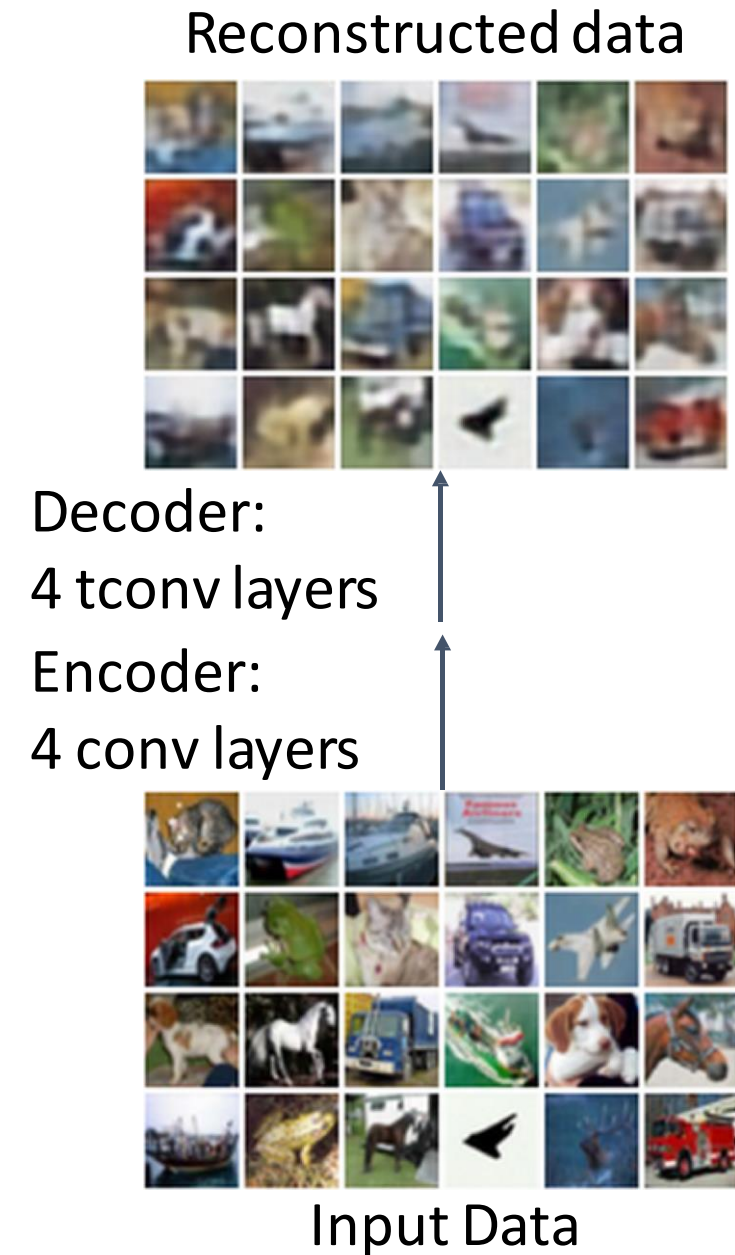


Decoder:
4 tconv layers
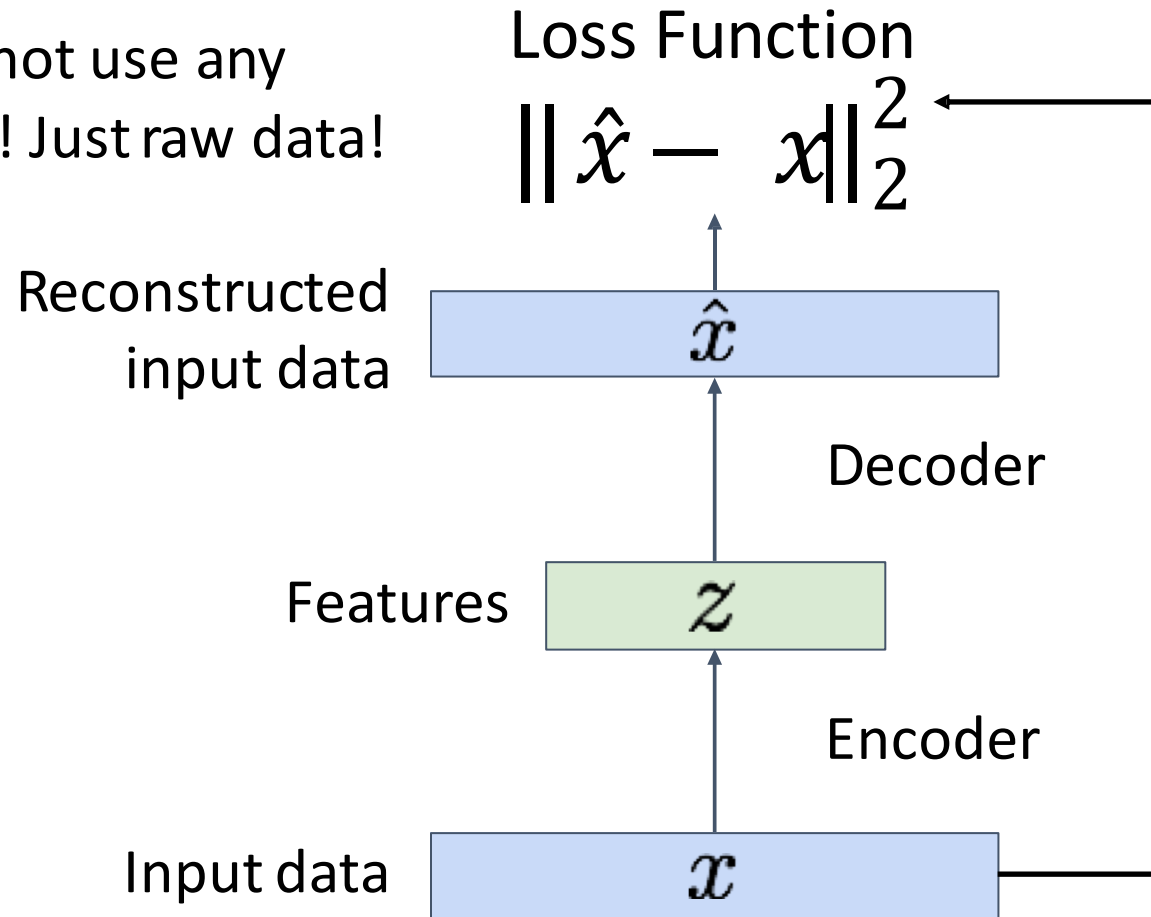
Encoder:
4 conv layers



Input Data

# (Regular, non-variational) Autoencoders

**Loss**: L2 distance between input and reconstructed data.
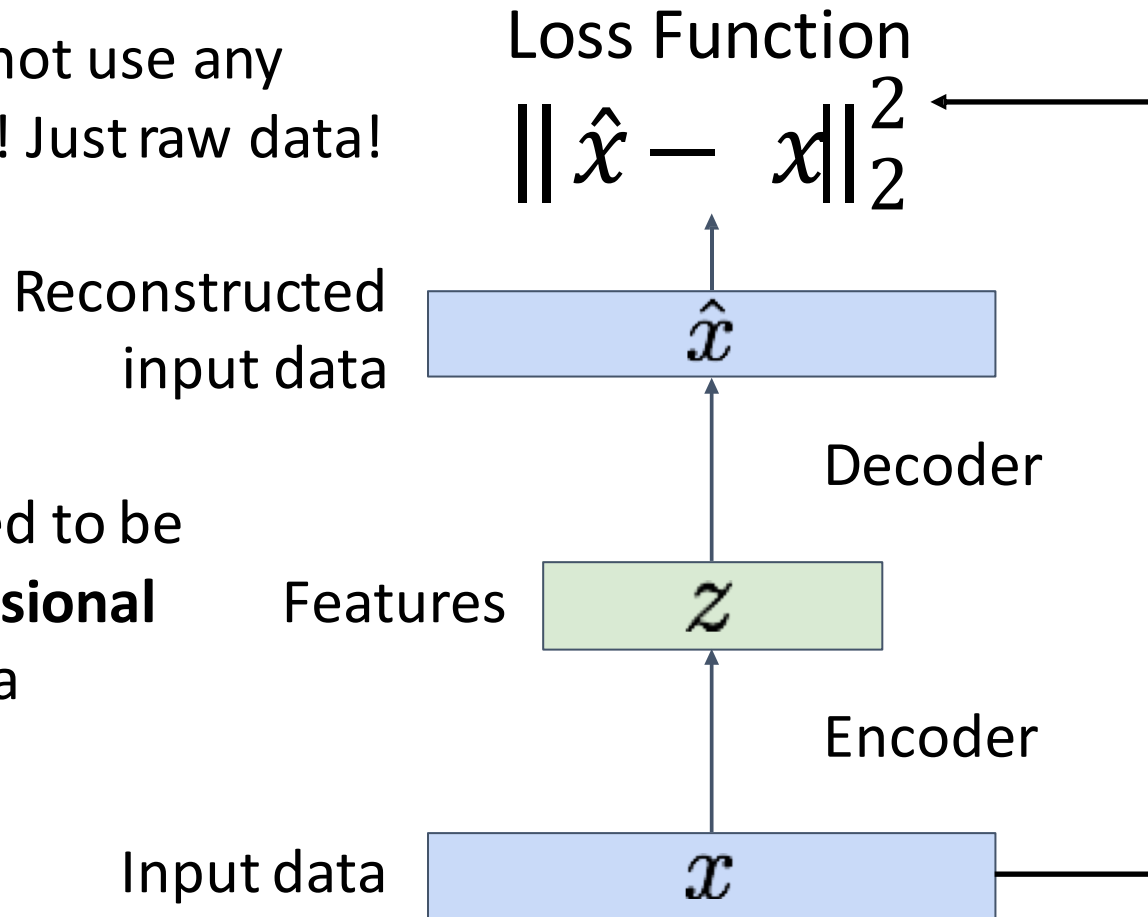
Does not use any labels! Just raw data!

Loss Function

$$\|\hat{x} - x\|_2^2$$
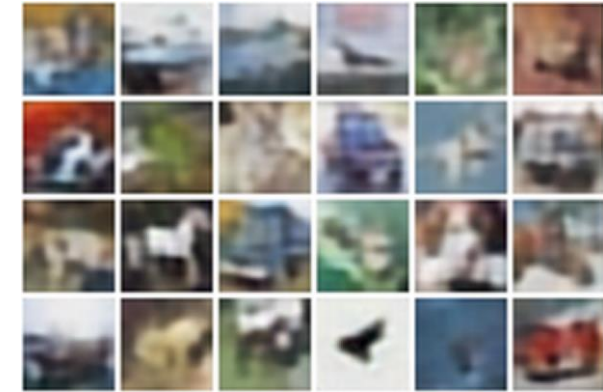
Reconstructed input data

$\hat{x}$

Decoder

Features need to be **lower dimensional** than the data
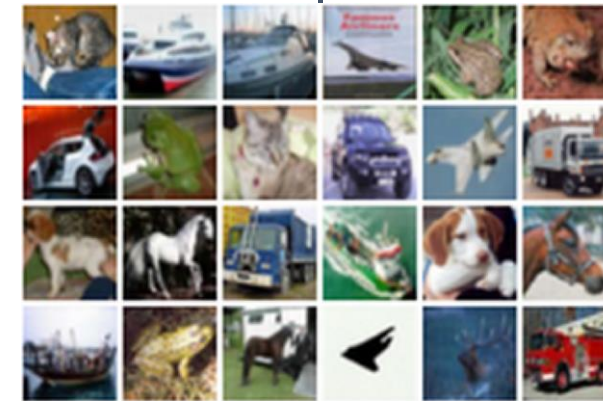
Features

$z$

Encoder

Input data

$x$

Reconstructed data



Decoder:
4 tconv layers

Encoder:
4 conv layers



Input Data

# (Regular, non-variational) Autoencoders

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

Reconstructed input data
$$\hat{x}$$

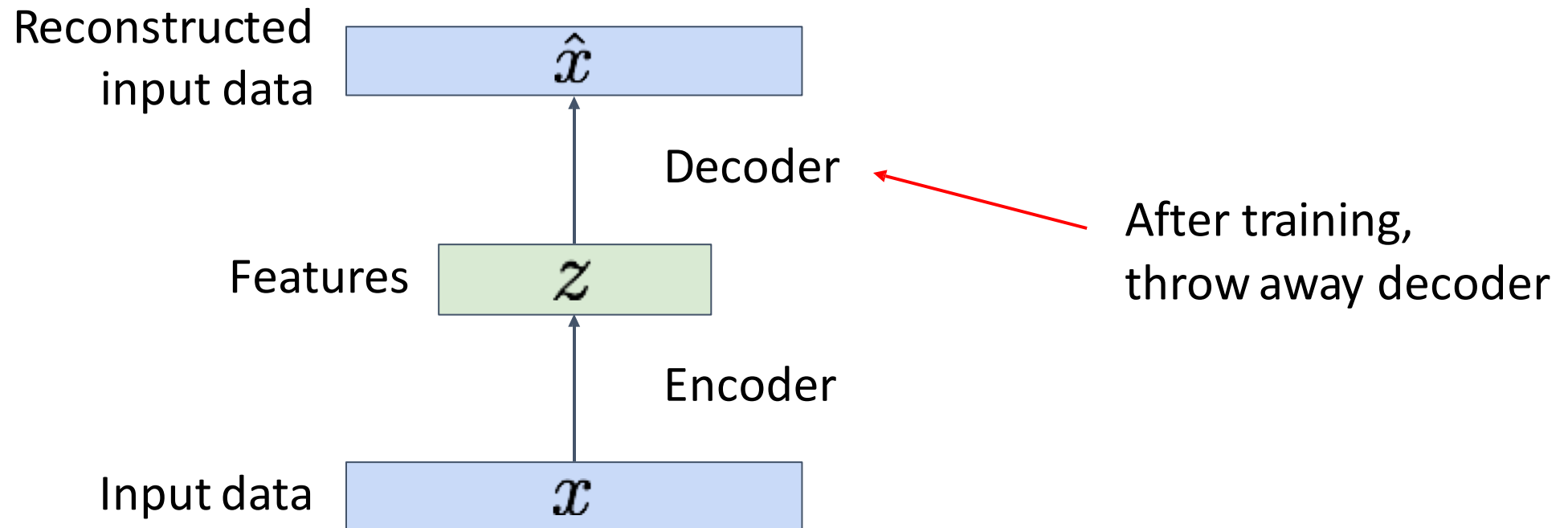Decoder ⟵ After training, throw away decoder

Features
$$z$$

Encoder

Input data
$$x$$

# (Regular, non-variational) Autoencoders

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

Loss function
(Softmax, etc)

Predicted Label $\hat{y}$    $y$

Classifier

Features $z$

Encoder

Input data $x$

Fine-tune encoder jointly with classifier

Encoder can be used to initialize a **supervised** model

bird    plane
dog    deer    truck

Train for final task (sometimes with small data)

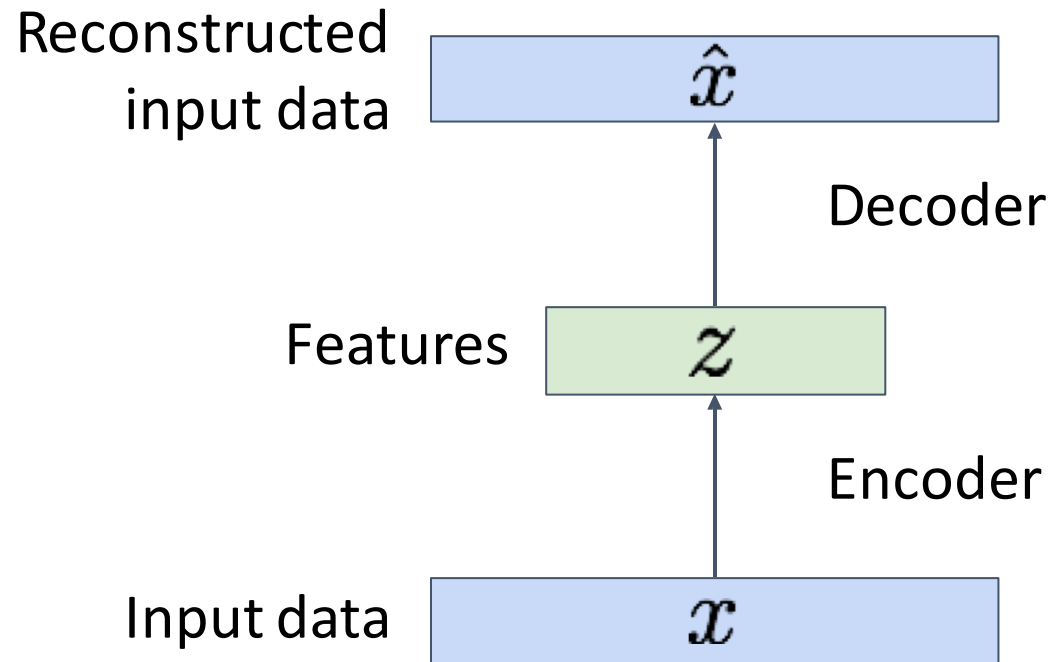# (Regular, non-variational) Autoencoders

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

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

# <u>Variational</u> Autoencoders

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

# 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 z from raw data
2. Sample from the model to generate new data

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.

# Variational Autoencoders

Probabilistic spin on autoencoders:
1. Learn latent features z from raw data
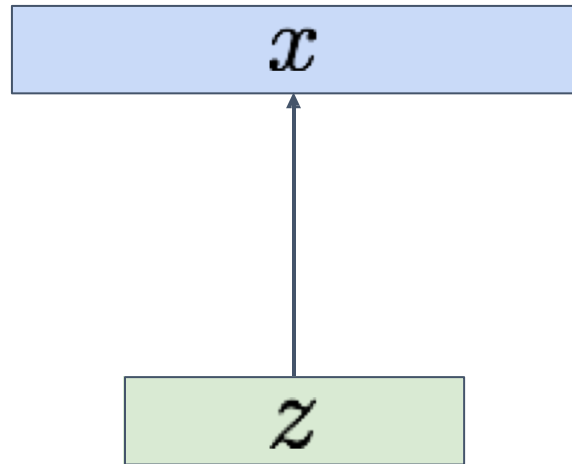2. Sample from the model to generate new data

After training, sample new data like this:

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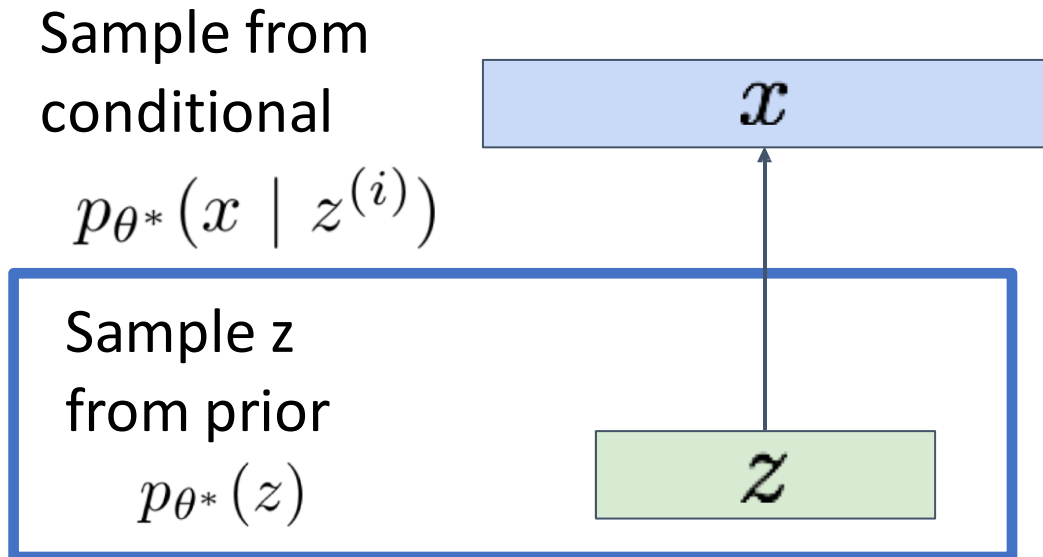
**Intuition: x** is an image, **z** is latent factors used to generate **x:** attributes, orientation, etc.

# Variational Autoencoders

Probabilistic spin on autoencoders:
1. Learn latent features z from raw data
2. Sample from the model to generate new data

After training, sample new data like this:

Sample from conditional

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

Sample z from prior

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

$$x$$

$$z$$

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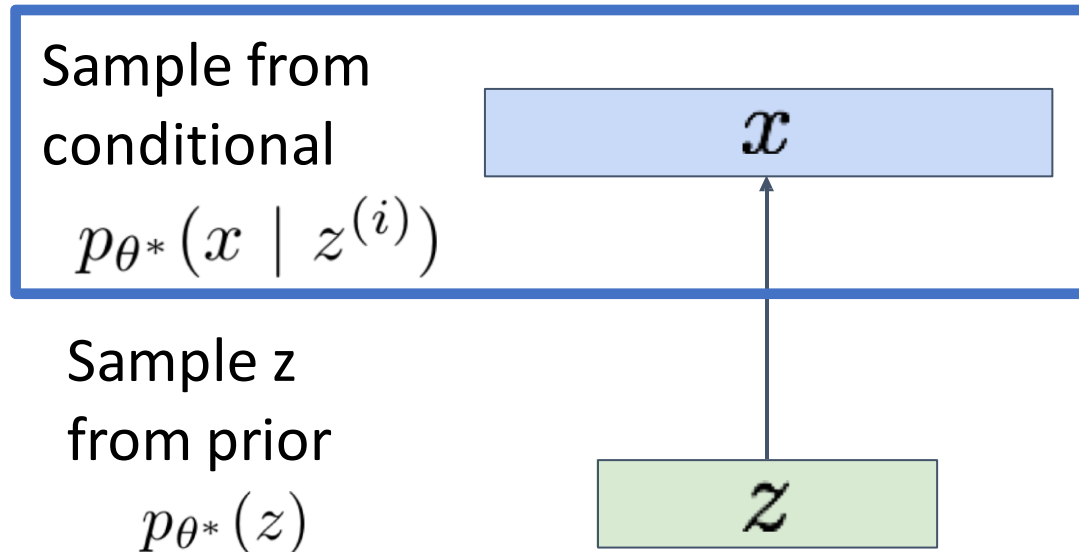
Assume simple prior p(z), e.g. Gaussian

# Variational Autoencoders

Probabilistic spin on autoencoders:
1. Learn latent features z from raw data
2. Sample from the model to generate new data

After training, sample new data like this:

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

$$x$$

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

$$z$$

- 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 autencoder)

# 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

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

$$\mu_{x|z} \qquad \Sigma_{x|z}$$

$$z$$

- 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 autencoder)

# Variational Autoencoders

Decoder must be **probabilistic**:
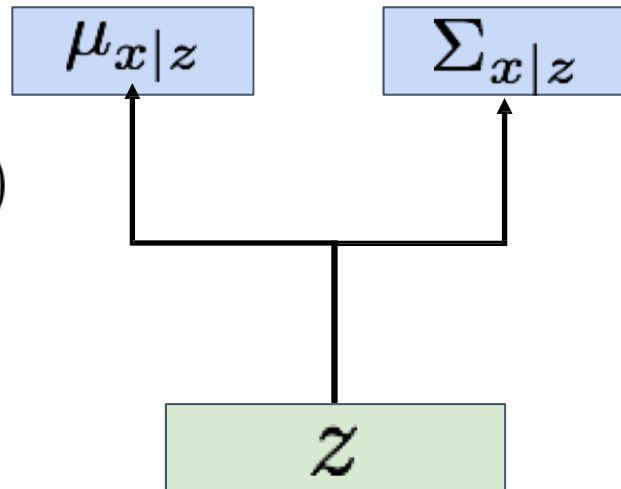Decoder inputs z, outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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

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

$$\mu_{x|z}$$

$$\Sigma_{x|z}$$

$$z$$

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 $\Sigma_{x|z}$

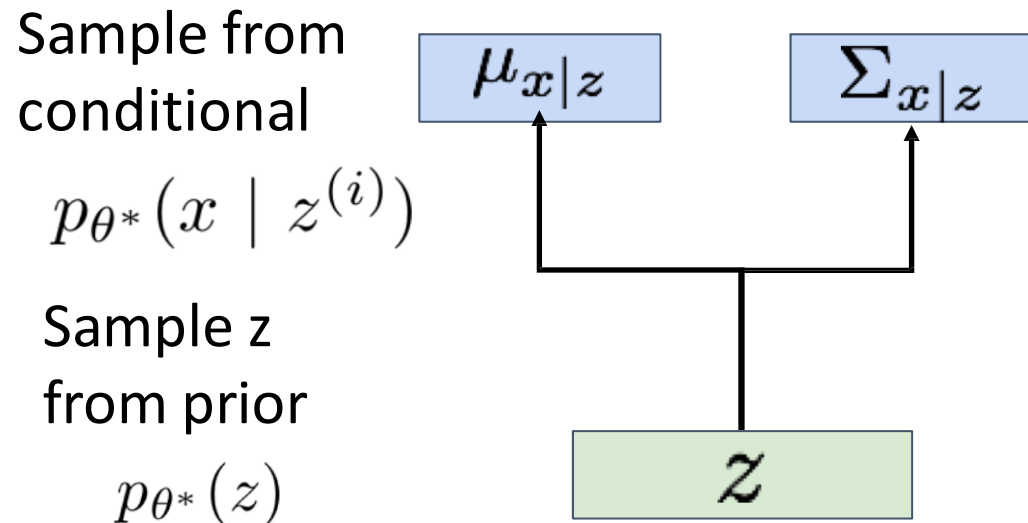Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

Sample from conditional

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

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**

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

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

$\mu_{x|z}$

$\Sigma_{x|z}$

$z$

# 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 $\Sigma_{x|z}$
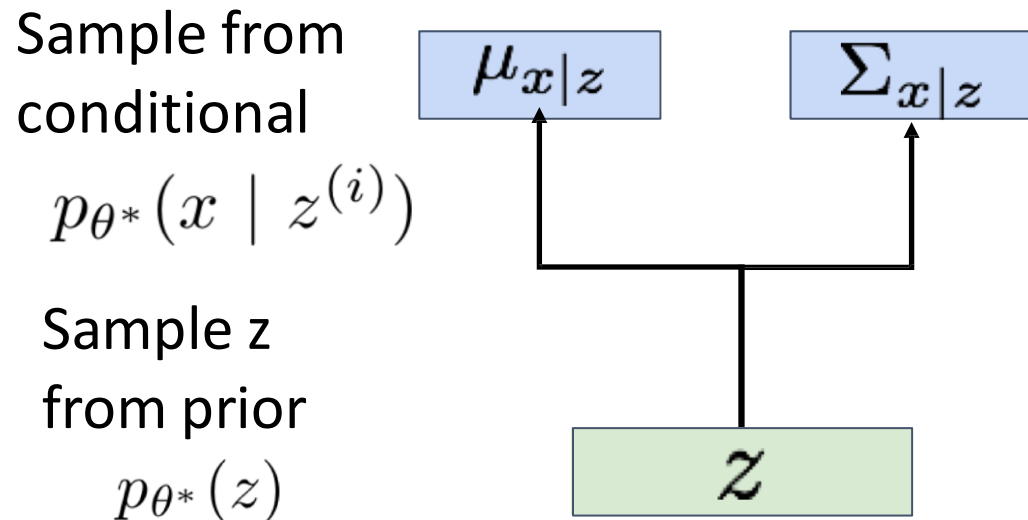
Sample from conditional

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

Sample z from prior

$p_{\theta*}(z)$



| $\mu_{x\mid z}$ | $\Sigma_{x\mid z}$ |
|---|---|

$z$

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**

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, can compute this with decoder network

# Variational Autoencoders

Decoder must be **probabilistic**:
Decoder inputs z, outputs mean $\mu_{x|z}$
and (diagonal) covariance $\Sigma_{x|z}$

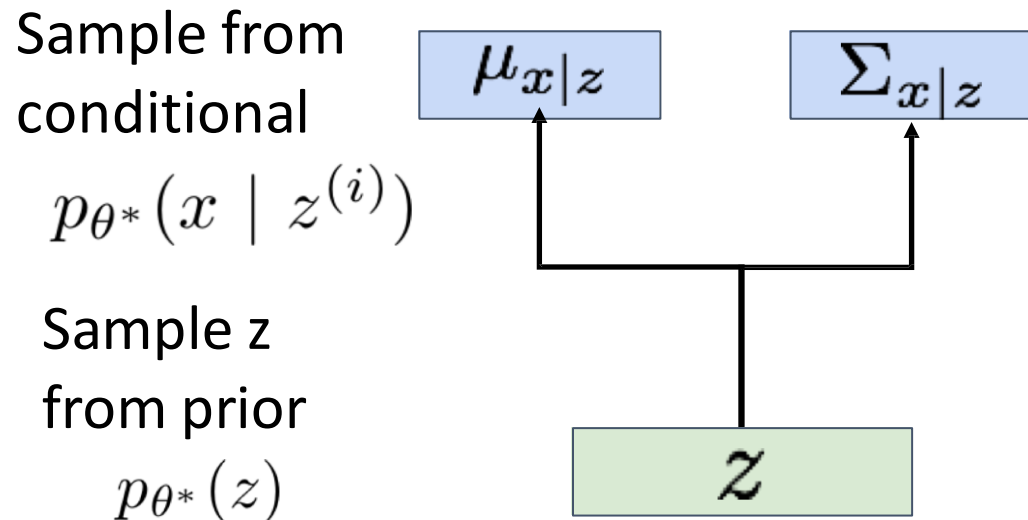Assume training data $\{x^{(i)}\}_{i=1}^{N}$ is generated from unobserved (latent) representation **z**
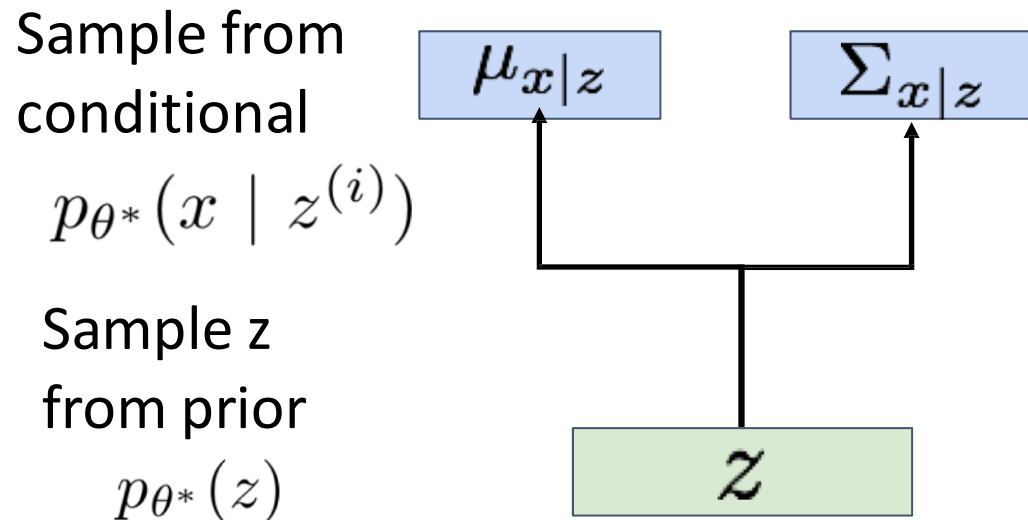
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

How to train this model?

Sample from conditional

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



Basic idea: **maximize likelihood of data**

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

Sample z from prior

$p_{\theta*}(z)$

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

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**

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

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

**Problem: Impossible to integrate over all z!**

$\mu_{x|z}$  $\Sigma_{x|z}$

$z$

# Variational Autoencoders

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

Decoder must be **probabilistic**:
Decoder inputs z, outputs 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**
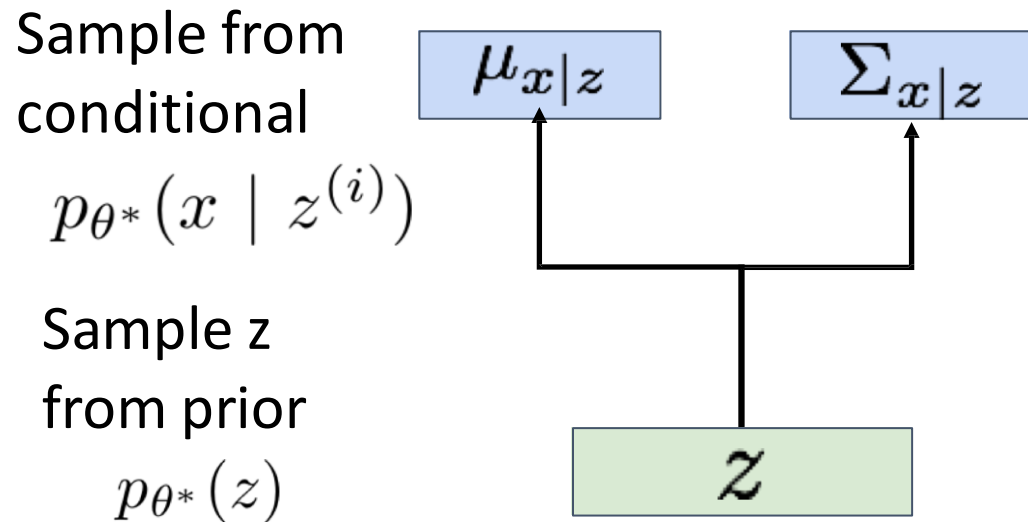
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

How to train this model?

Sample from conditional

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

Sample z from prior

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


$\mu_{x|z}$    $\Sigma_{x|z}$    $z$

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

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

# Variational Autoencoders

Decoder must be **probabilistic**:
Decoder inputs z, outputs mean $\mu_{x|z}$
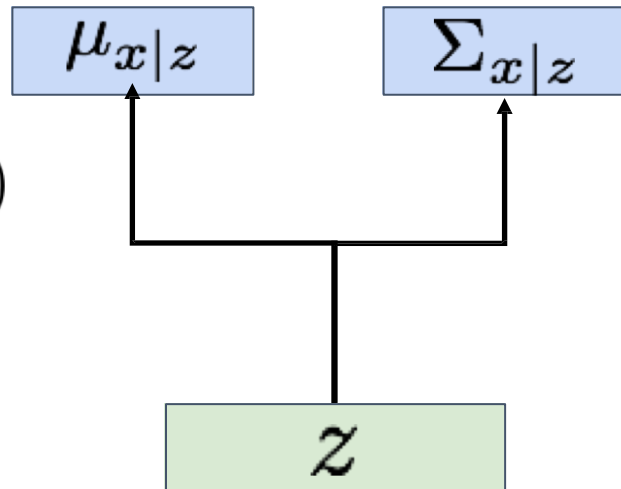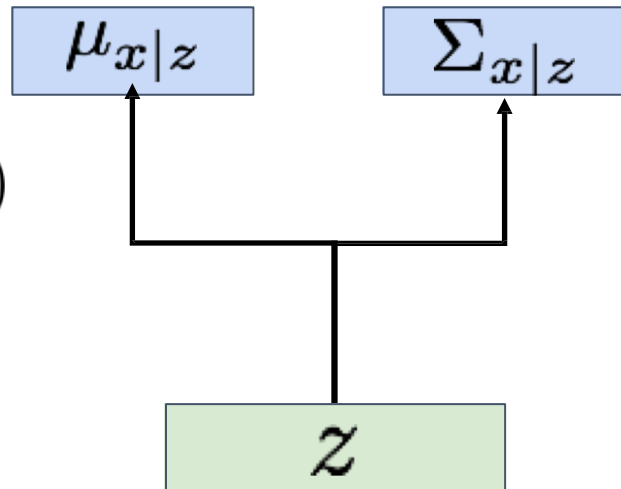and (diagonal) covariance $\Sigma_{x|z}$

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**

Another idea: Try Bayes' Rule:

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

Ok, compute with
decoder network

# Variational Autoencoders

Decoder must be **probabilistic**:
Decoder inputs z, outputs mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$
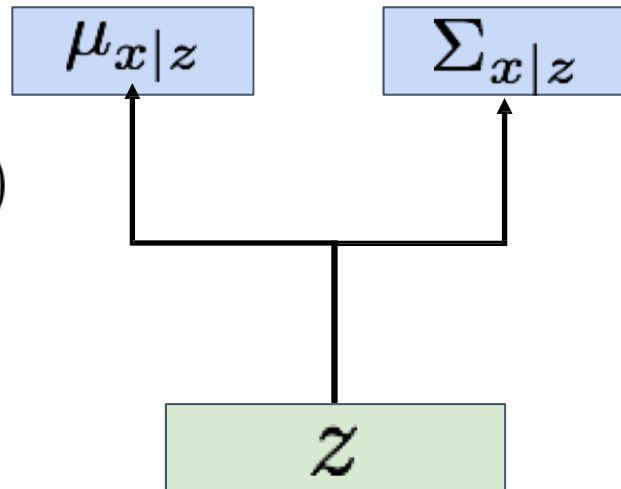
Sample x from Gaussian with mean $\mu_{x|z}$ and (diagonal) covariance $\Sigma_{x|z}$

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**

Another idea: Try Bayes' Rule:

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

Ok, we assumed Gaussian prior

$\mu_{x|z}$  $\Sigma_{x|z}$

$z$

# Variational Autoencoders

Decoder must be **probabilistic**:
Decoder inputs z, outputs mean $\mu_{x|z}$
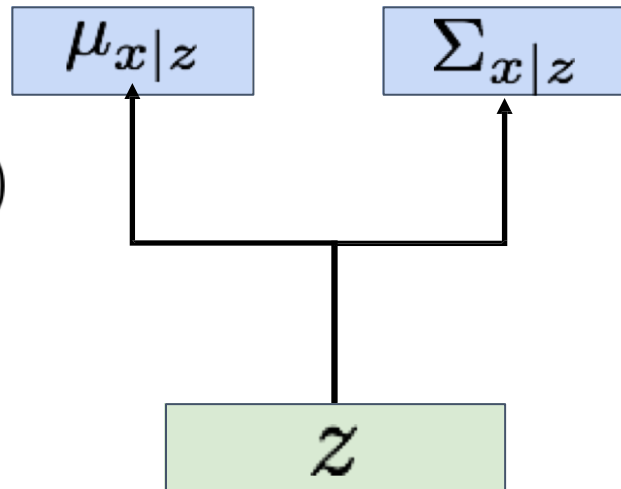and (diagonal) covariance $\Sigma_{x|z}$

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

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

How to train this model?

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



$\mu_{x|z}$    $\Sigma_{x|z}$

$z$

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

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

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

**Problem**: No way to compute this!

# Variational Autoencoders

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

Decoder must be **probabilistic**:
Decoder inputs z, outputs mean $\mu_{x|z}$
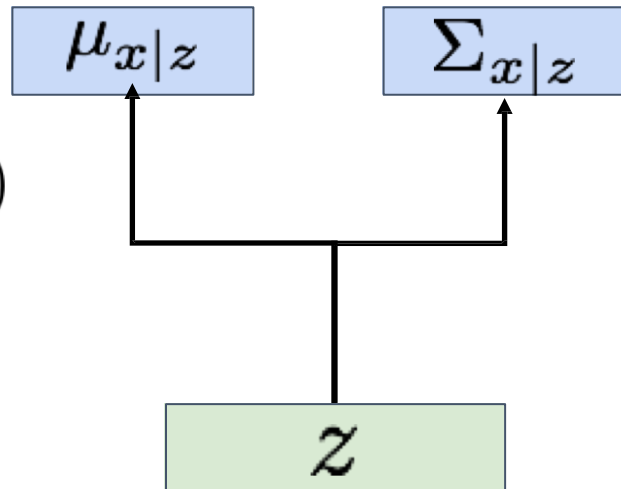and (diagonal) covariance $\Sigma_{x|z}$

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

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

How to train this model?

Sample from conditional

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



Sample z from prior

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

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

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

**Solution:** Train another network **(encoder)** that learns

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

# Variational Autoencoders

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

Decoder must be **probabilistic**:
Decoder inputs z, outputs mean $\mu_{x|z}$
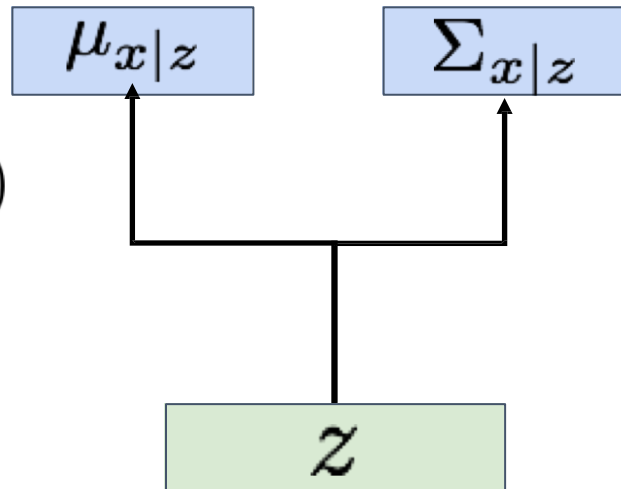and (diagonal) covariance $\Sigma_{x|z}$

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

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

How to train this model?

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

Basic idea: **maximize likelihood of data**

Another idea: Try Bayes' Rule:

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



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

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

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



**Encoder network** inputs data x, gives distribution over latent codes z

$$q_N(z \mid x) = N(\mu_{z|x}, \Sigma_{z|x})$$



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

then we can approximate

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

**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

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

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

# Variational Autoencoders

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

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

Split up using rules for logarithms

# Variational Autoencoders

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

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

Split up using rules for logarithms

# Variational Autoencoders

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

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

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

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

# Variational Autoencoders

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

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

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

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

# Variational Autoencoders

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

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

$$= E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{\text{KL}}\left(q_\phi(z|x), p(z)\right) + D_{\text{KL}}(q_\phi(z|x), p_\theta(z|x))$$

Data reconstruction

# Variational Autoencoders

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

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

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

KL divergence between prior, and samples from the encoder network

# Variational Autoencoders

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

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

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

<span style="color:red">KL divergence between encoder and posterior of decoder</span>

# Variational Autoencoders

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

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

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

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

## Variational Autoencoders

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

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

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

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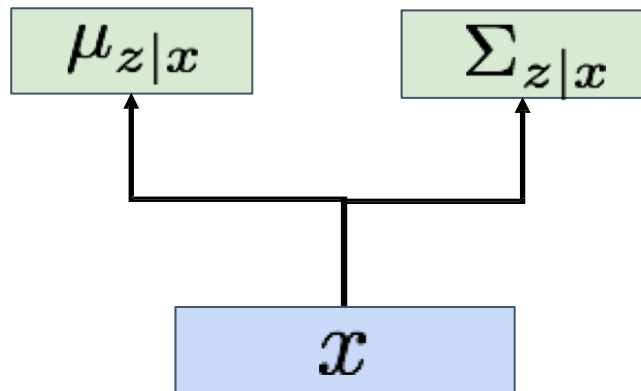
# Variational Autoencoders

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

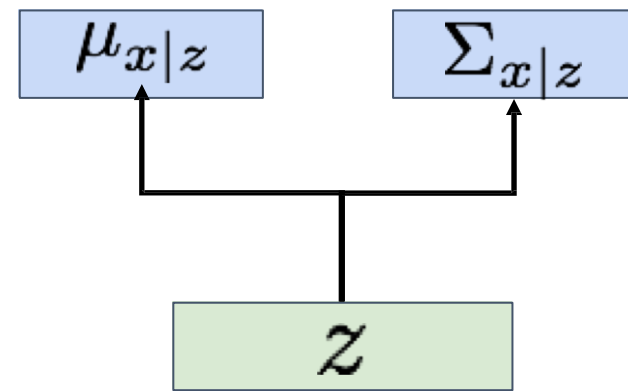$$\log p_\theta(x) \geq E_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{\text{KL}}\Big(q_\phi(z|x), p(z)\Big)$$

**Encoder Network**

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

$\mu_{z|x}$   $\Sigma_{z|x}$

$x$

**Decoder Network**

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

$\mu_{x|z}$   $\Sigma_{x|z}$

$z$

# Generative Adversarial Nets

- Review of the GAN:

$$\min_{G} \max_{D} E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim q(z)}[\log 1 - D(z)]$$

- A two-player game: $G$ wants to fool $D$ by creating an approximant distribution $q$ close to the true distribution $p_{data}$, and $D$ wants to be smart so that it can distinguish between $p_{data}$ and $q$ induced by $G$.



Code z

Generator G

Generated Image

Training Image

Discriminator D

Fake Or Real