# CS 4803 / 7643: Deep Learning

Topics:
- Variational Auto-Encoders (VAEs)
  - AEs, Variational Inference
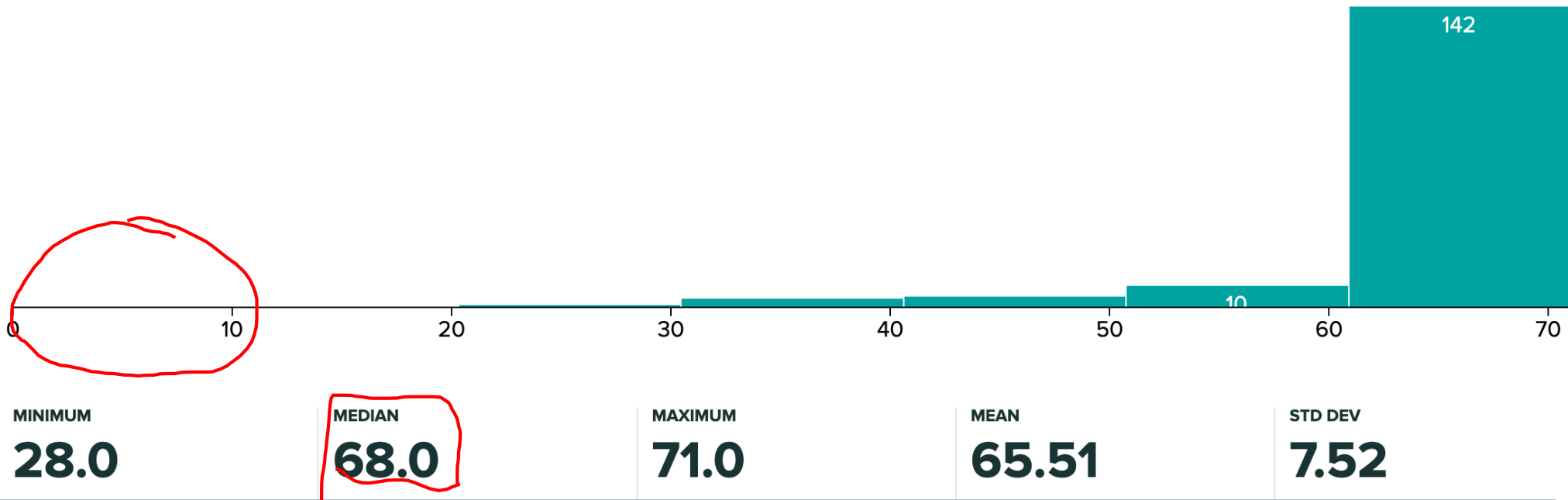
Dhruv Batra

Georgia Tech

# Administrativia

- HW4 Reminder
  - Due: 11/07, 11:55pm
  - Reinforcement Learning
  - Last HW. Focus on project after that.

- Final project
  - No poster session
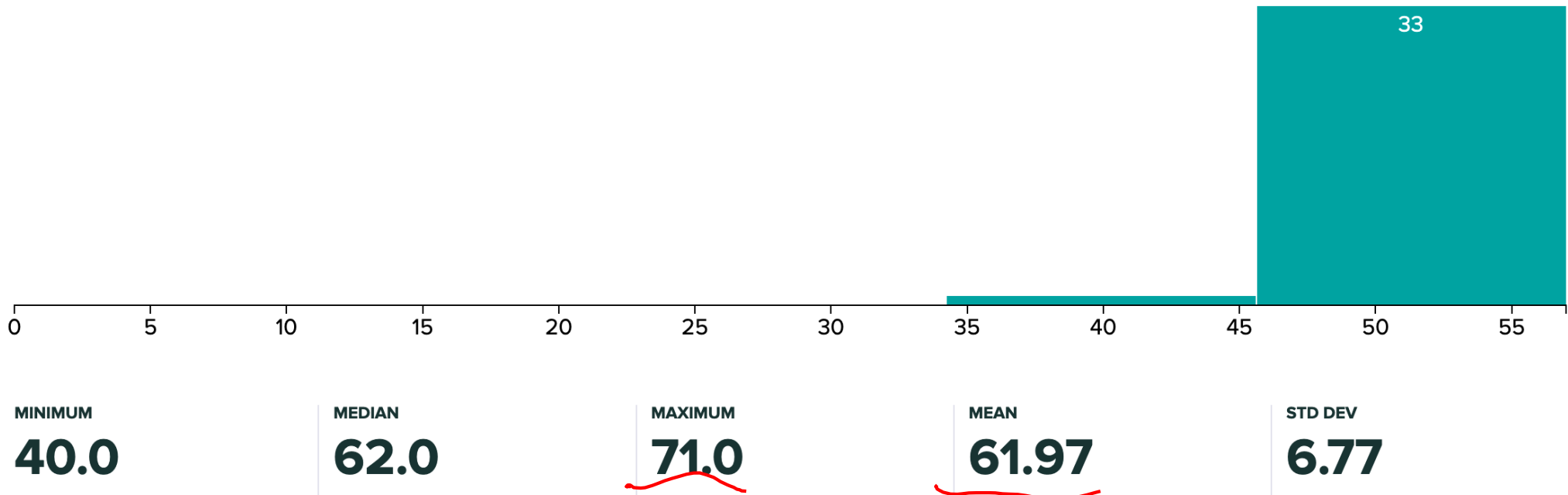  - PDF Report submission
    - Details out soon

# Administrativia

- HW3 Grades Released
  - Regrade requests close: 11/15, 11:55pm
  - Please check solutions first!

- Grade histogram: 7643
  - Max possible: 71 (regular credit) + 0 (extra credit)



| MINIMUM | MEDIAN | MAXIMUM | MEAN | STD DEV |
|---------|--------|---------|------|---------|
| 28.0 | 68.0 | 71.0 | 65.51 | 7.52 |

# Administrativia

- HW3 Grades Released
  - Regrade requests close: 11/15, 11:55pm
  - Please check solutions first!

- Grade histogram: 4803
  - Max possible: 55 (regular) + 14 (extra credit)



| MINIMUM | MEDIAN | MAXIMUM | MEAN | STD DEV |
|---------|--------|---------|------|---------|
| 40.0 | 62.0 | 71.0 | 61.97 | 6.77 |

# Recap from ~~last time~~ 2 lectures ago

# Supervised vs Reinforcement vs Unsupervised Learning

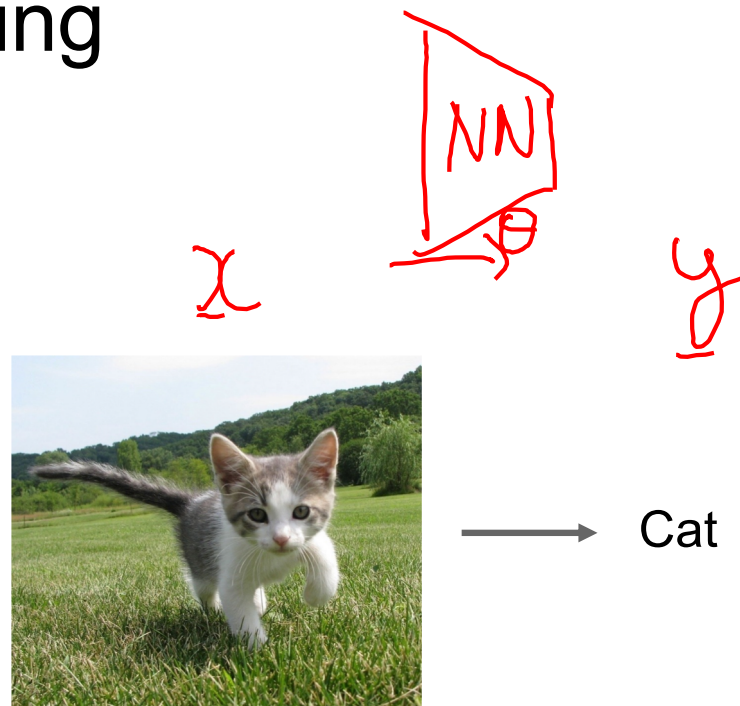# Supervised vs Reinforcement vs Unsupervised Learning

**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.

$NN$

$x$

$y$

$\longrightarrow$ Cat

Classification

# Supervised vs Reinforcement vs Unsupervised Learning
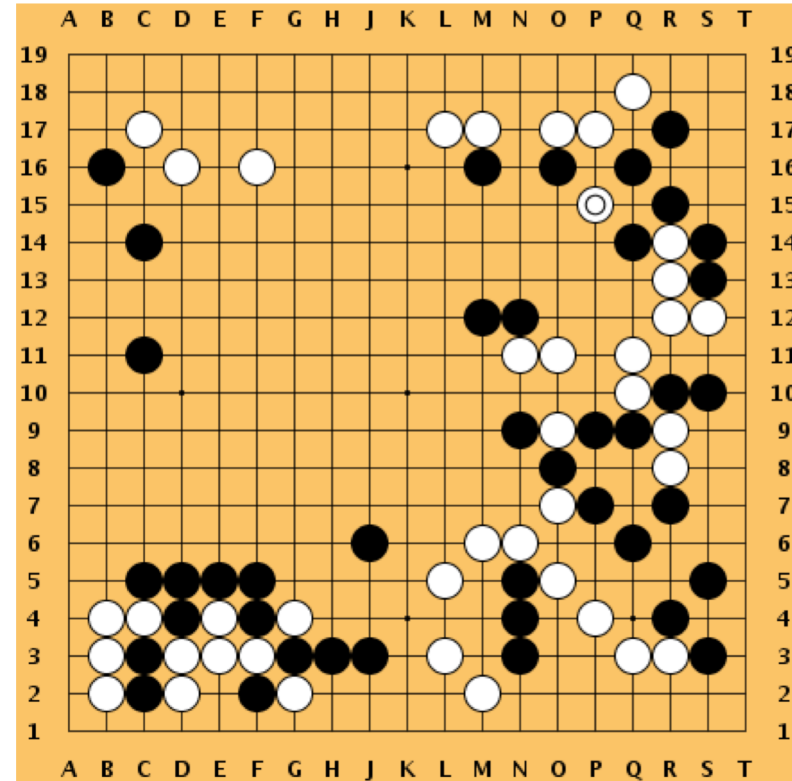
$S \rightarrow a^{st}$

**Reinforcement Learning**

**Given**: (e, r)
Environment e, Reward function r
  (evaluative feedback)

**Goal**: Maximize expected reward

**Examples**: Robotic control, video games, board games, etc.

# Supervised vs Reinforcement 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 Reinforcement 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 Learning
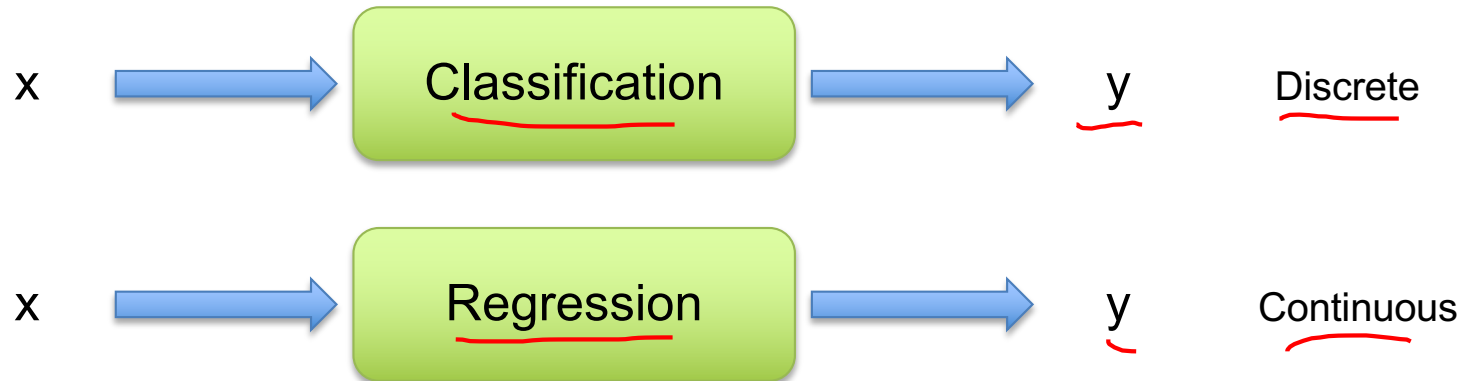
**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x → y

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

# Supervised vs Reinforcement vs Unsupervised Learning

$$p(\vec{x}) \qquad p(x_i \mid y_j)$$

## Unsupervised Learning

Training data is cheap

**Data**: x
Just data, no labels!

Holy grail: Solve unsupervised learning => understand structure of visual world

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

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

## Supervised Learning

**Data**: (x, y)
x is data, y is label

**Goal**: Learn a *function* to map x → y

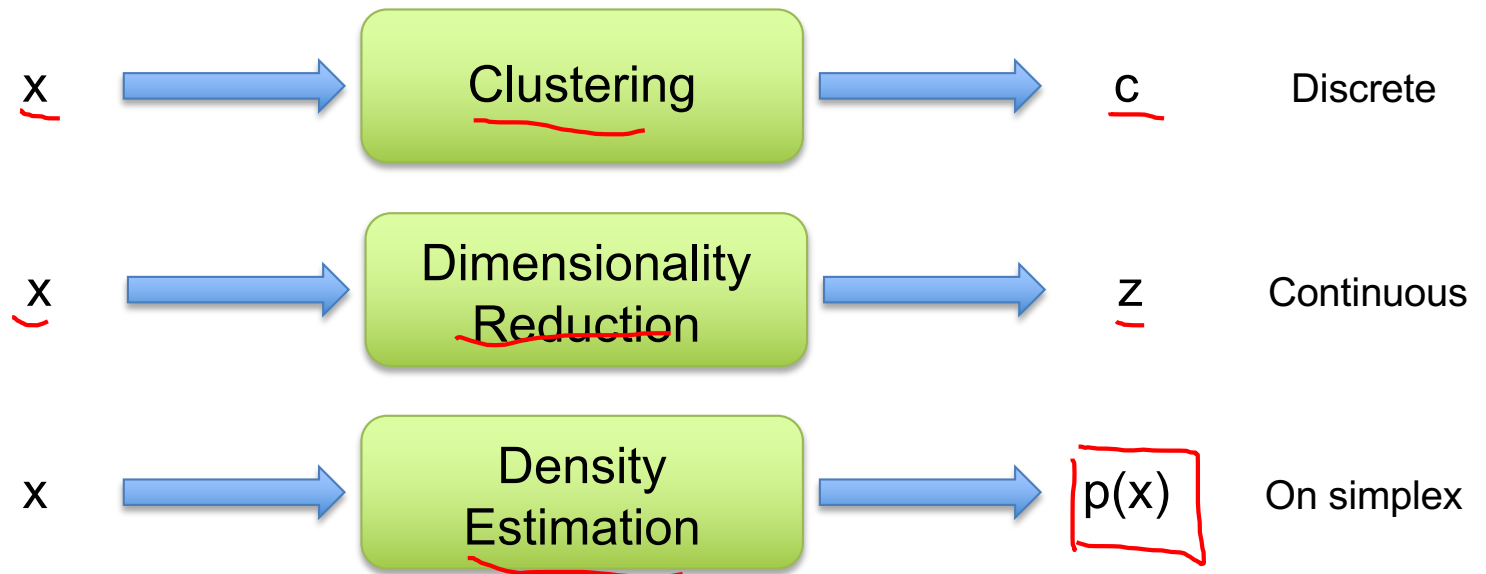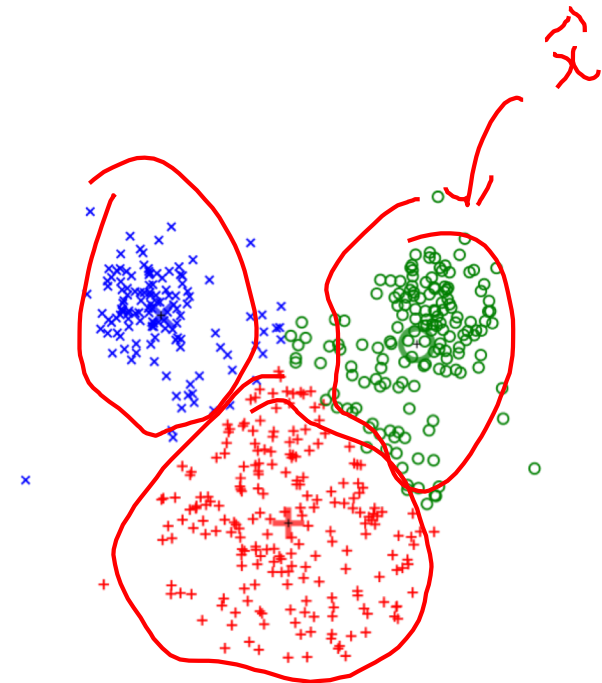**Examples**: Classification, regression, object detection, semantic segmentation, image captioning, etc.

# Tasks

## Supervised Learning

$x$ → **Classification** → $y$    Discrete

$x$ → **Regression** → $y$    Continuous

## Unsupervised Learning

$x$ → **Clustering** → $c$    Discrete

$x$ → **Dimensionality Reduction** → $z$    Continuous

$x$ → **Density Estimation** → $p(x)$    On simplex

# Supervised vs Reinforcement vs Unsupervised Learning

**Unsupervised Learning**

**Data**: x
Just data, no labels!

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

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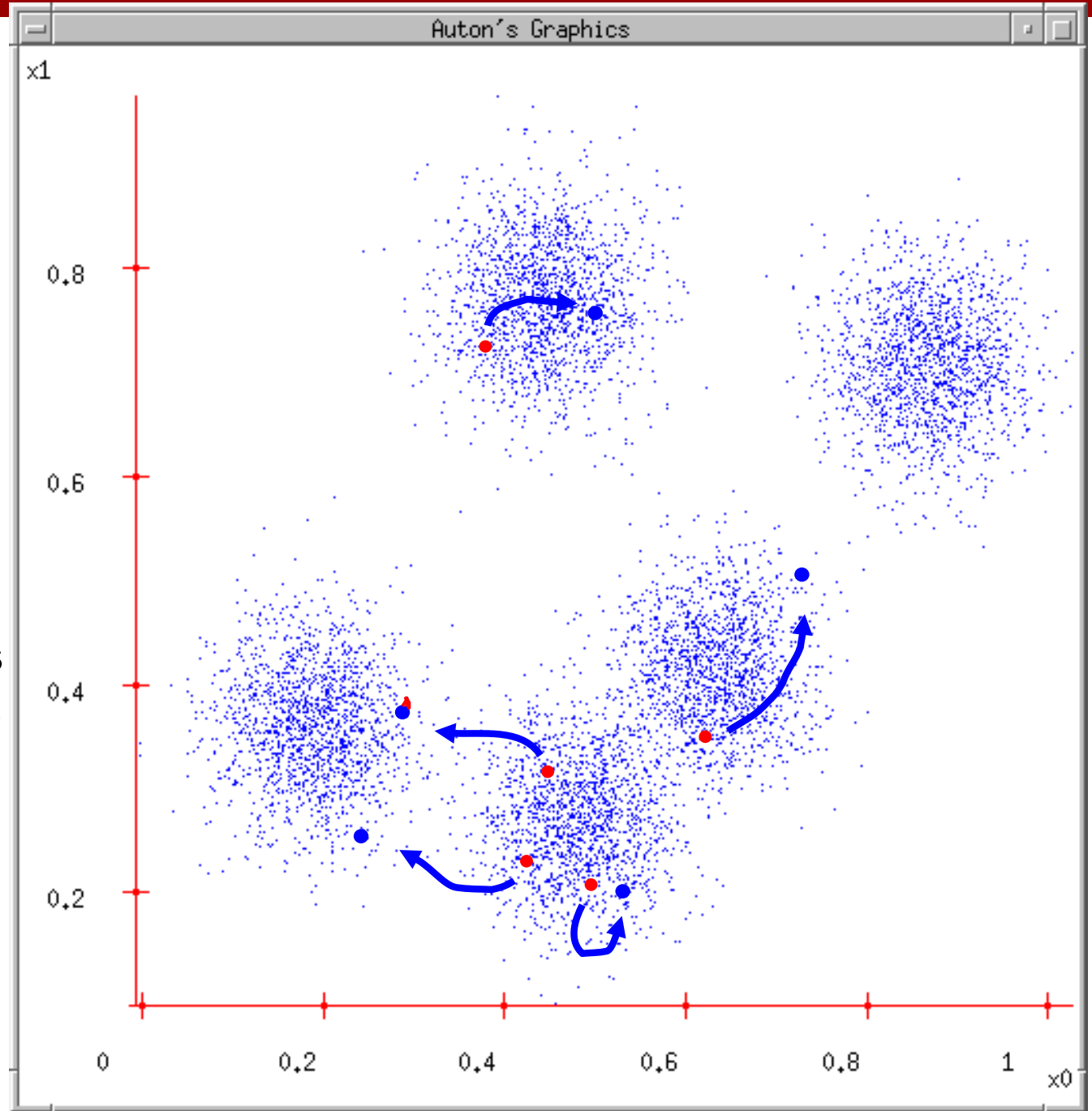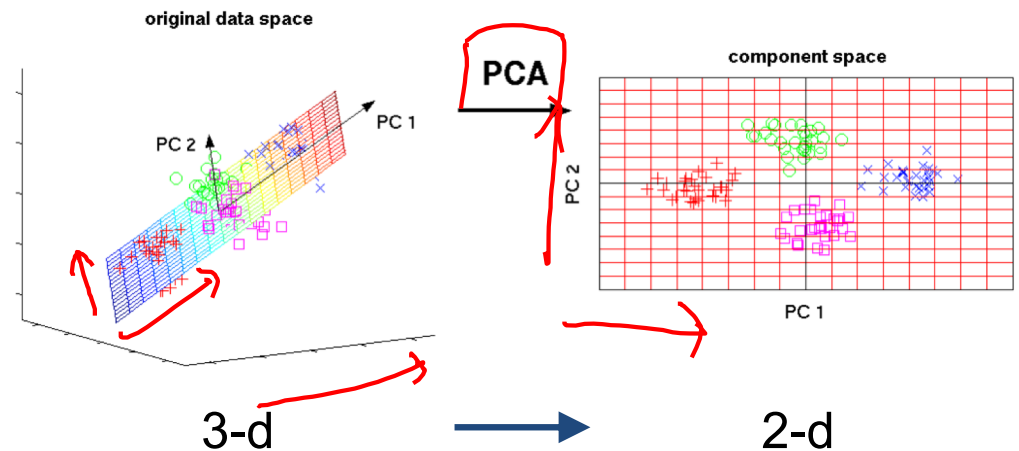    feature learning, density
    estimation, etc.

K-means clustering

# K-means

1. Ask user how many clusters they'd like. *(e.g. k=5)*

2. Randomly guess k cluster Center locations

3. Each datapoint finds out which Center it's closest to.

4. Each Center finds the centroid of the points it owns…

5. …and jumps there

6. …Repeat until terminated!

Slide Credit: Carlos Guestrin

15

# K-means as Co-ordinate Descent

- Optimize objective function:

$$\min_{\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_k} \min_{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_N} F(\boldsymbol{\mu}, \boldsymbol{a}) = \min_{\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_k} \min_{\boldsymbol{a}_1,\ldots,\boldsymbol{a}_N} \sum_{i=1}^{N} \sum_{j=1}^{k} a_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$$

*(handwritten annotations: $a_{ij} = 1$ if $x_i \equiv$ clusb $j$; distoclion)*

- Fix **μ**, optimize a
- Fix a, optimize **μ**

# Supervised vs Reinforcement 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.



3-d ➡ 2-d

Principal Component Analysis
(Dimensionality reduction)

# Supervised vs Reinforcement vs Unsupervised Learning

$p(\vec{x})$

## Unsupervised Learning

**Data**: x
Just data, no labels!

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

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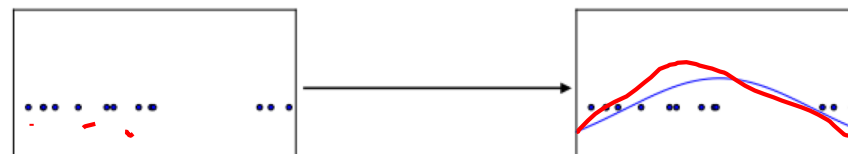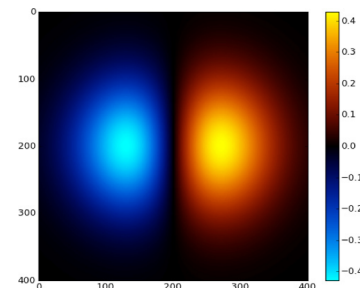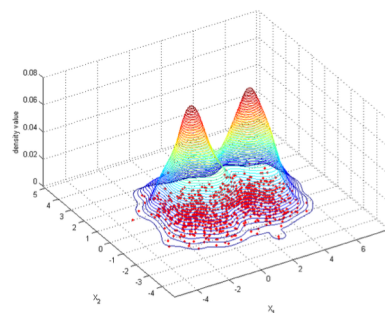feature learning, density
estimation, etc.

Figure copyright Ian Goodfellow, 2016. Reproduced with permission.
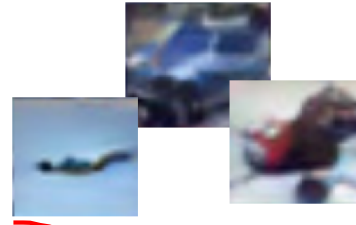
1-d density estimation

2-d density estimation

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Generative Models

Given training data, generate new samples from same distribution



Training data ~ $p_{data}(x)$

Generated samples ~ $p_{model}(x)$

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

# Generative Models

Given training data, generate new samples from same distribution

Training data ~ $p_{data}(x)$                    Generated samples ~ $p_{model}(x)$

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

Addresses density estimation, a core problem in unsupervised learning

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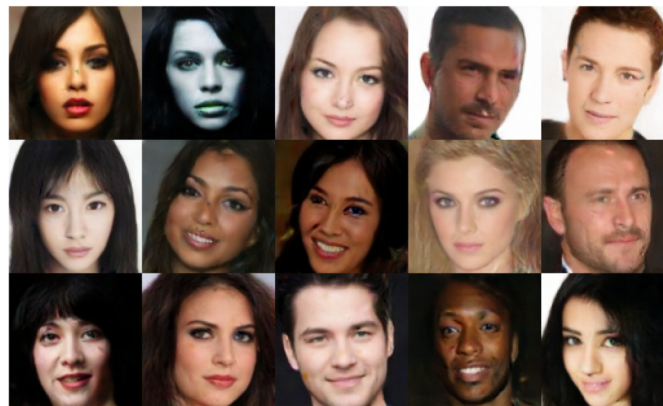- Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

# Why Generative Models?

*p(image | control)*

- 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
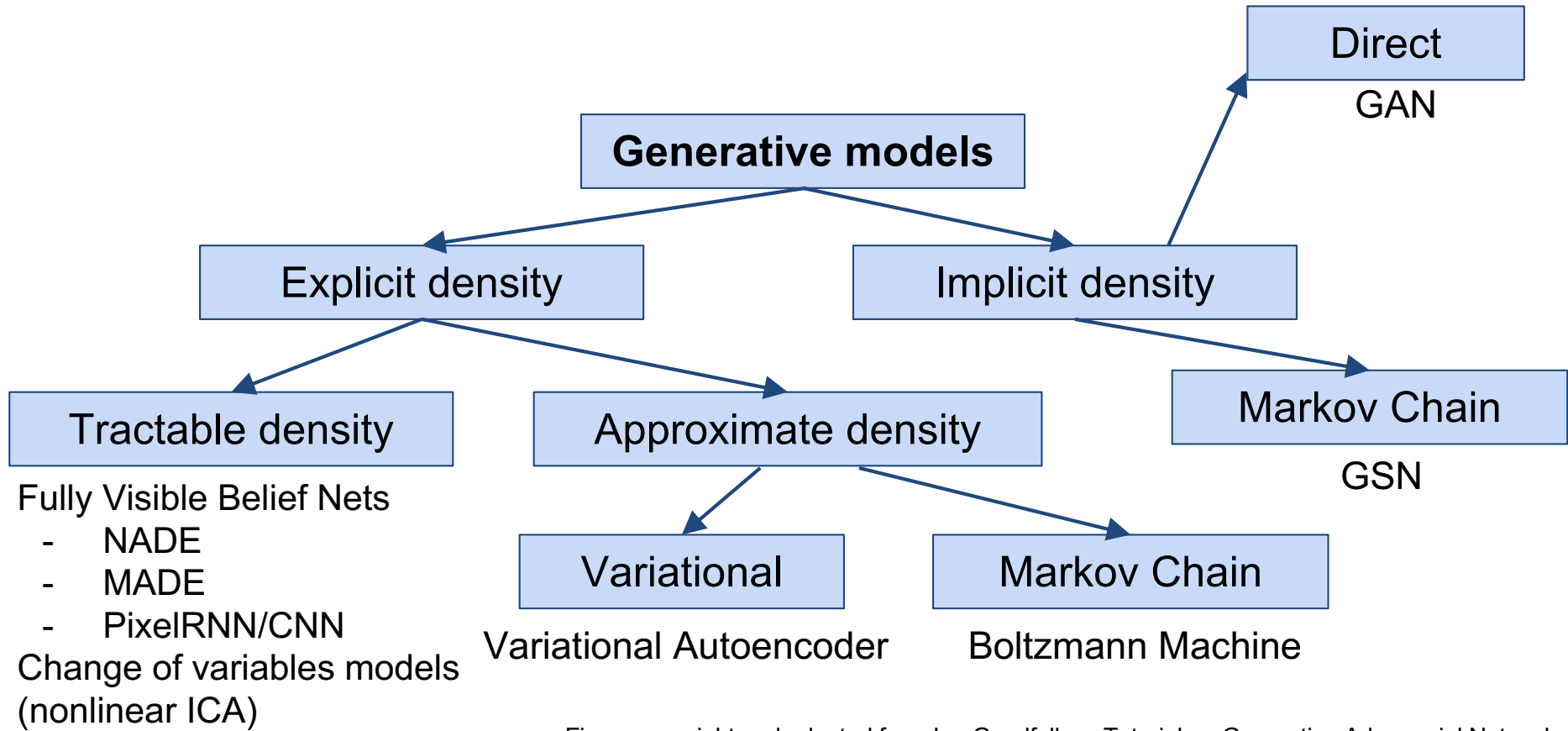
# Taxonomy of Generative Models



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

# Taxonomy of Generative Models

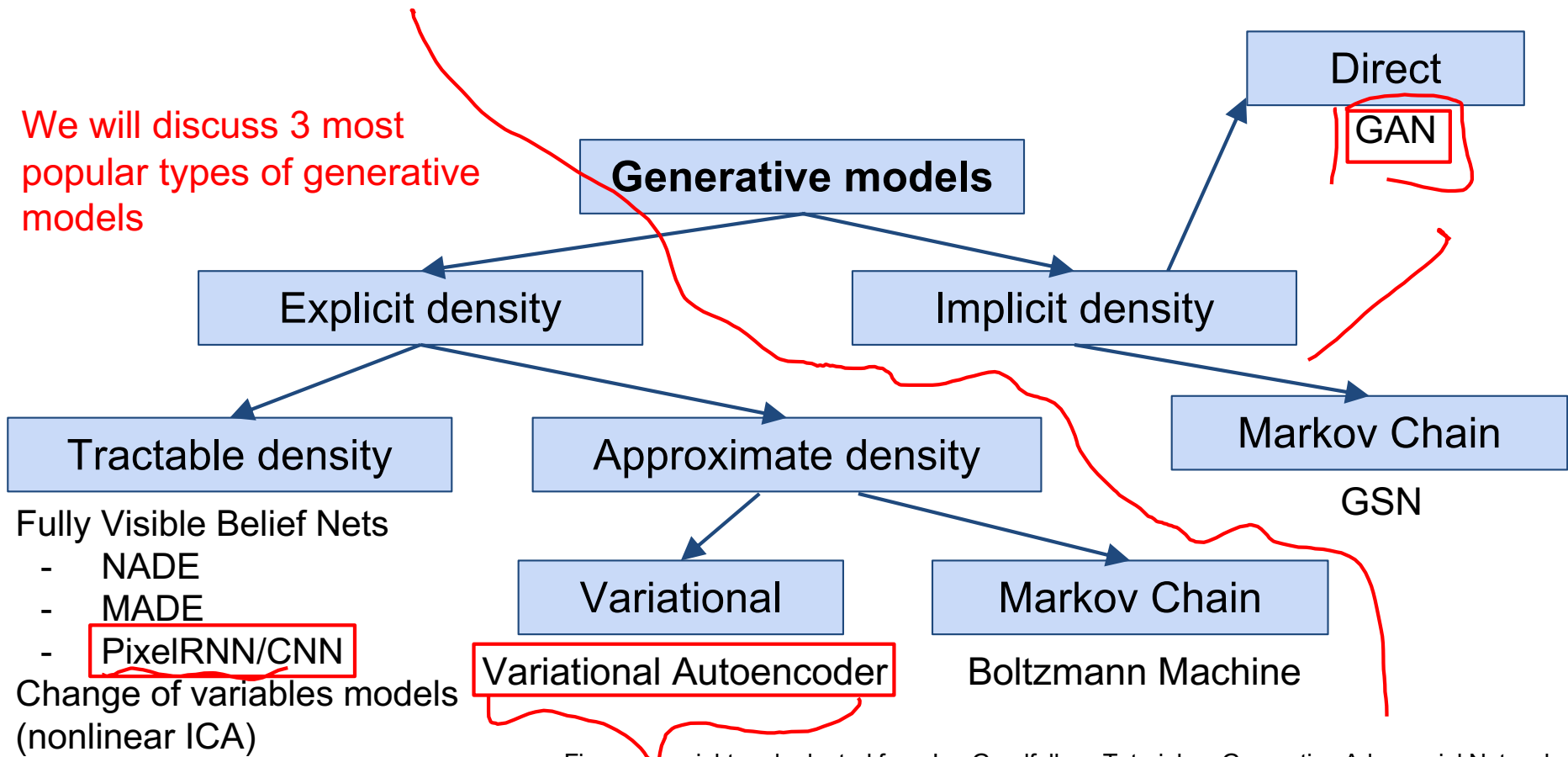We will discuss 3 most popular types of generative models

**Generative models**

Direct

GAN

Explicit density

Implicit density

Tractable density

Approximate density

Markov Chain

GSN

Fully Visible Belief Nets
- NADE
- MADE
- PixelRNN/CNN

Change of variables models (nonlinear ICA)

Variational

Markov Chain

Variational Autoencoder

Boltzmann Machine

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

# Fully Observable Model

## Explicit density model

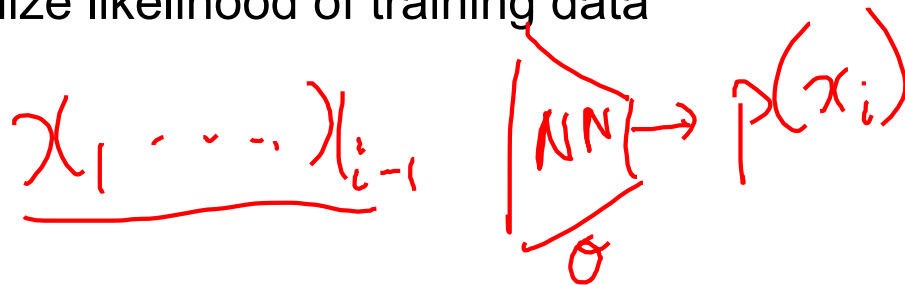Use chain rule to decompose likelihood of an image x into product of 1-d distributions:

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

$$D = \{ \vec{x}^{(i)} \}_{i=1}^{N}$$

Likelihood of image x

Probability of i'th pixel value given all previous pixels

Then maximize likelihood of training data

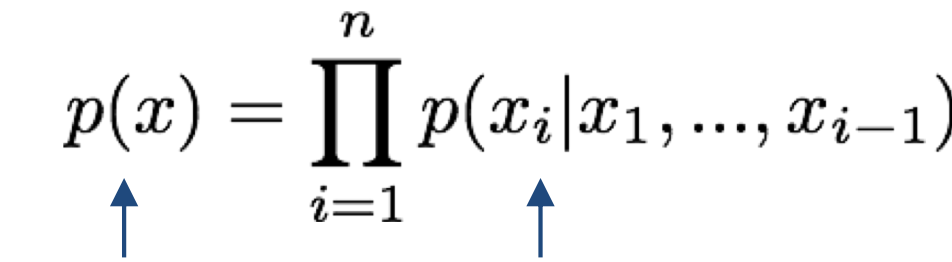$$x_1 \cdots x_{i-1} \quad | NN | \rightarrow p(x_i)$$

# Fully Observable Model

## Explicit density model

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

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

Likelihood of image x

Probability of i'th pixel value given all previous pixels

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

Then maximize likelihood of training data

# Plan for Today

- Goal: Variational Autoencoders

- Latent variable probabilistic models
  - Example GMMs

- Autoencodeders
- Variational Inference

# Variational Autoencoders (VAE)

# So far...

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

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

$$p(\vec{x}) \qquad D = \{\vec{x}\}$$

$$p(\vec{x}, \vec{z}) \quad \text{"latent"}$$
$$\text{"hidden"}$$

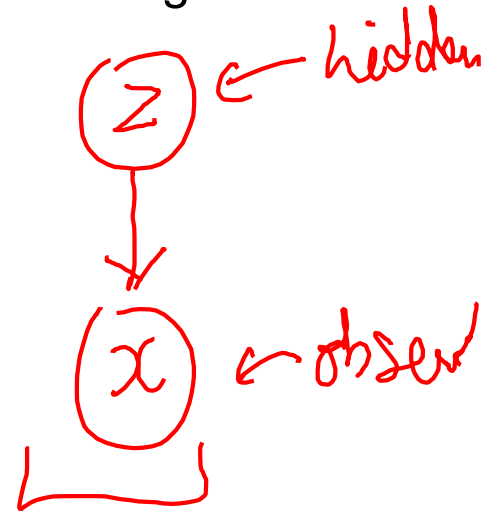$$= \underbrace{p(\vec{x}|z)}_{\text{conditional}} \underbrace{p(z)}_{\text{prior}}$$

# So far...

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

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

VAEs define intractable density function with latent **z**:

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

z continuous

$$\sum_z p_\theta(z) \, p_\theta(x|z)$$

z discrete

# So far...

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

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

VAEs define intractable density function with latent **z**:

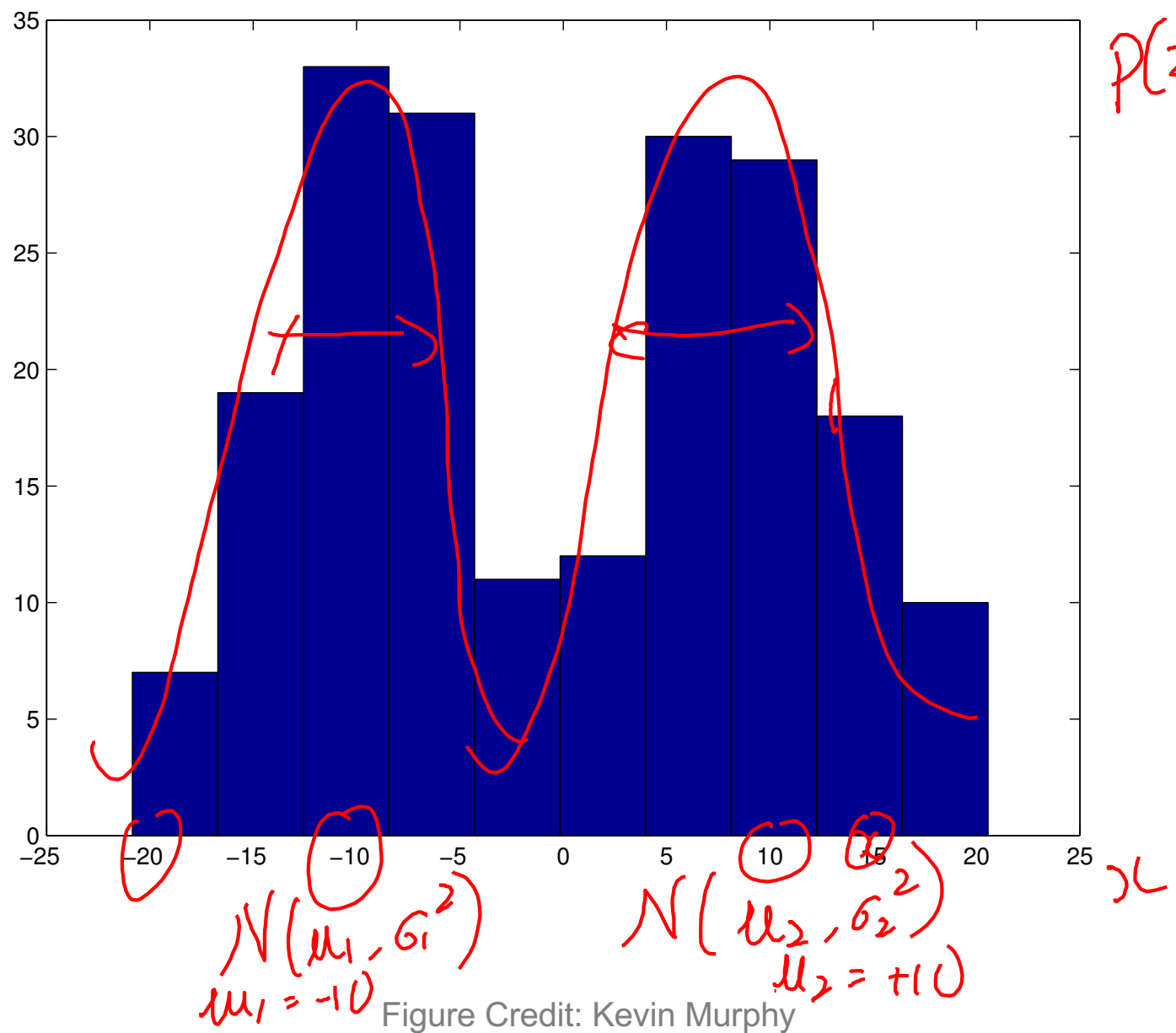$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$

*z* ← hidden

*x* ← obser

Cannot optimize directly, derive and optimize lower bound on likelihood instead

# GMM

Gaussian Mixture Model

$z \in \{1, 2\}$

$$P(z) = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$



$p(x)$

$N(\mu_1, \sigma_1^2)$

$\mu_1 = -10$

$N(\mu_2, \sigma_2^2)$

$\mu_2 = +10$

$x$

# Gaussian Mixture Model

$Z \sim Cat(\vec{\pi})$

$\begin{bmatrix} \pi_1 \\ \vdots \\ \pi_k \end{bmatrix}$

$Z \in \{1, \ldots, k\}$

$\pi_c = P(z = c)$

$$X | \boxed{Z = c} \sim N(\mu_c, \sigma_c^2) = \frac{1}{\sqrt{2\pi \sigma_c^2}} e^{-\frac{(x - \mu_c)^2}{2\sigma_c^2}}$$

$$P(\vec{x}) = \underbrace{P(z)}_{\pi_z} \underbrace{P(\vec{x} | z)}_{N}$$

# Gaussian Mixture Model

$$P(z) = \pi_z$$
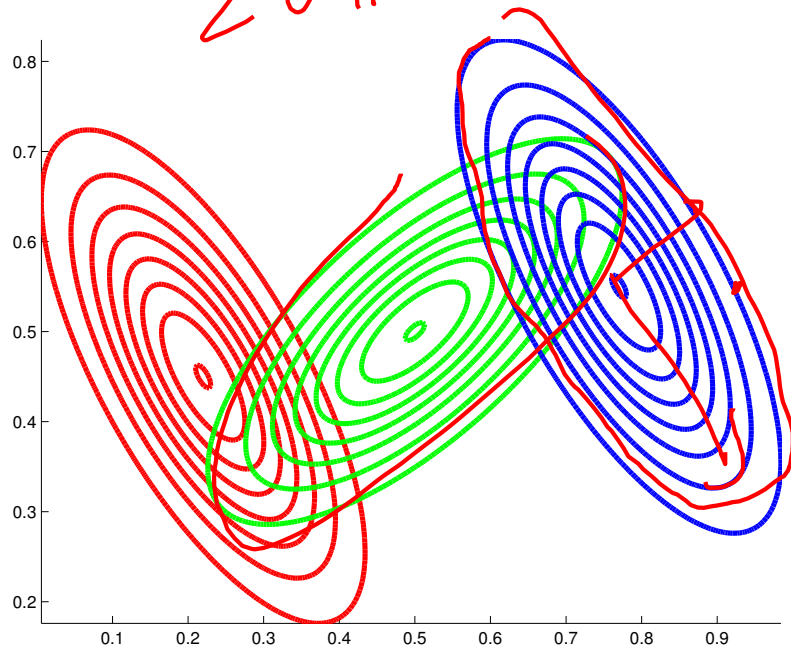
$$P(x|z) = N(\mu_z, \sigma_z^2)$$

Available from model

$$P(\vec{x}) = \boxed{\sum_z \boxed{P(z)} P(x|z)} = \text{Marginalization}$$

$$P(z|\vec{x}) = \frac{P(z, \vec{x})}{P(\vec{x})} = \frac{P(\vec{x}|z) P(z)}{\boxed{\sum_z ( \quad ) ( \quad )}} = \text{"Inference"}$$
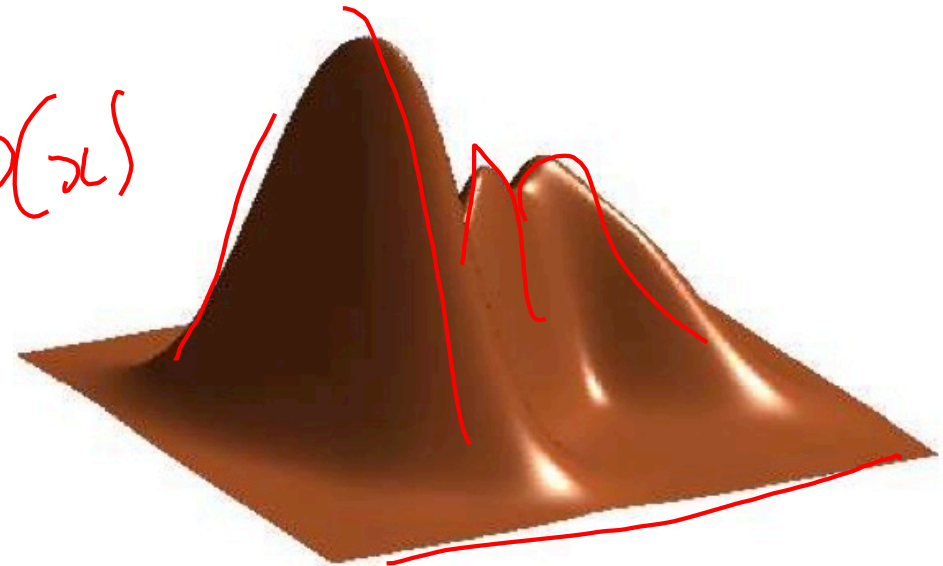
# GMM

$$\vec{x} \in \mathbb{R}^d$$
$$\vec{\mu} \in \mathbb{R}^d$$
$$\Sigma \in \mathbb{R}^{d \times d}$$

$$N(\vec{\mu}, \Sigma) = \frac{1}{\sqrt{2\pi^d |\Sigma|}} e^{-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1}(\vec{x}-\vec{\mu})}$$



$p(x)$

# K-means vs GMM

- K-Means
  - [http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html](http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html)

- GMM
  - [https://lukapopijac.github.io/gaussian-mixture-model/](https://lukapopijac.github.io/gaussian-mixture-model/)

# Hidden Data Causes Problems #1

- Fully Observed (Log) Likelihood factorizes

- Marginal (Log) Likelihood doesn't factorize

- All parameters coupled!

$$\text{Parameters} = \{ r_1, \cdots r_k, \mu_1 \cdots \mu_k, \Sigma_1 \cdots \Sigma_k \} \equiv \Theta$$

$$D = \{ \vec{x_i} \}$$

$$\hat{\Theta}_{MLE} = \underset{\Theta}{\arg\max} \; P(D \mid \Theta) \equiv \log P(D \mid \Theta)$$

$$= \sum_i \log P(\vec{x_i} \mid \Theta)$$

$$= \sum_i \log \sum_{z_i} P(x_i, z_i \mid \Theta)$$

$$\log \int \qquad P(x_i \mid z_i) P(z_i)$$

$$P(x_i, z_i \mid \Theta)$$
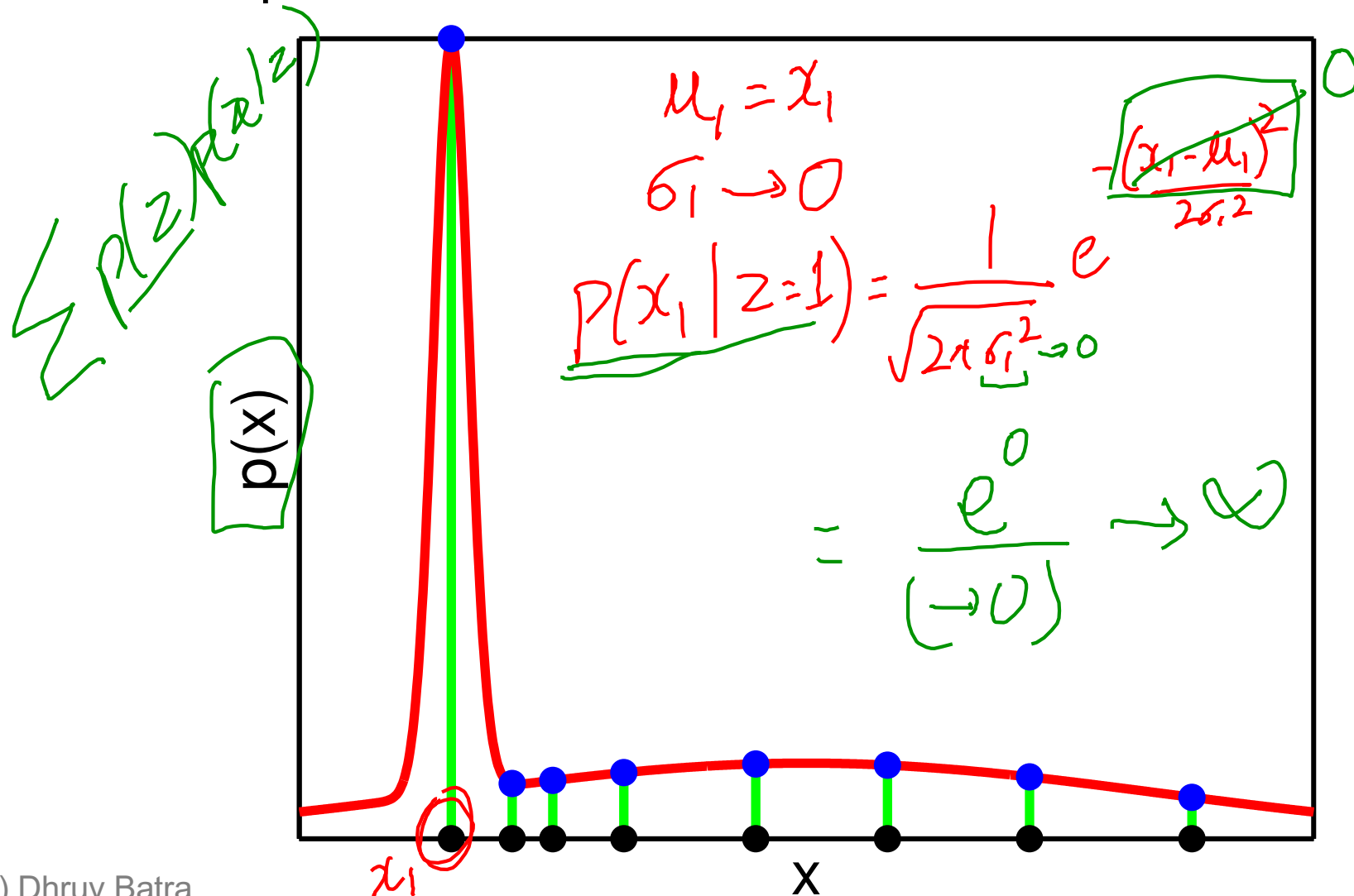$$z_i = 1$$

# Hidden Data Causes Problems #2

- Identifiability



$\mu_1 = -10$
$\mu_2 = +10$  |  $\mu_1 = +10$
$\mu_2 = -10$

$P(z)$

$\overline{P(x)}$

# Hidden Data Causes Problems #3

- Likelihood has singularities if one Gaussian "collapses"



$$\sum_z p(z) p(x|z)$$

$$\mu_1 = x_1$$

$$\sigma_1 \to 0$$

$$P(x_1 | z=1) = \frac{1}{\sqrt{2\pi \sigma_1^2} \to 0} e^{\frac{-(x_1-\mu_1)^2}{2\sigma_1^2} \to 0}$$

$$= \frac{e^0}{(\to 0)} \to \infty$$

$$\textcircled{z} \sim \text{Cat}$$

$$\textcircled{x}$$

$$x | z \sim \mathcal{N}$$

GMM

# Variational Auto Encoders

VAEs are a combination of the following ideas: $p(\bar{x})$

$\sum p(x|z)$

$p(z)$

1. Auto Encoders

2. Variational Approximation
   - Variational Lower Bound / ELBO
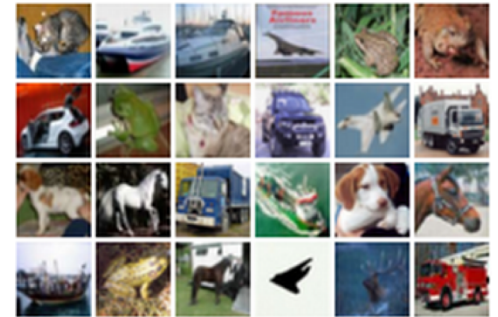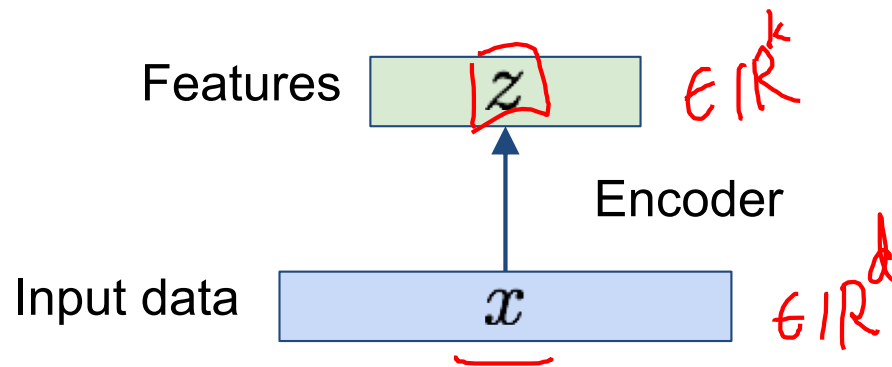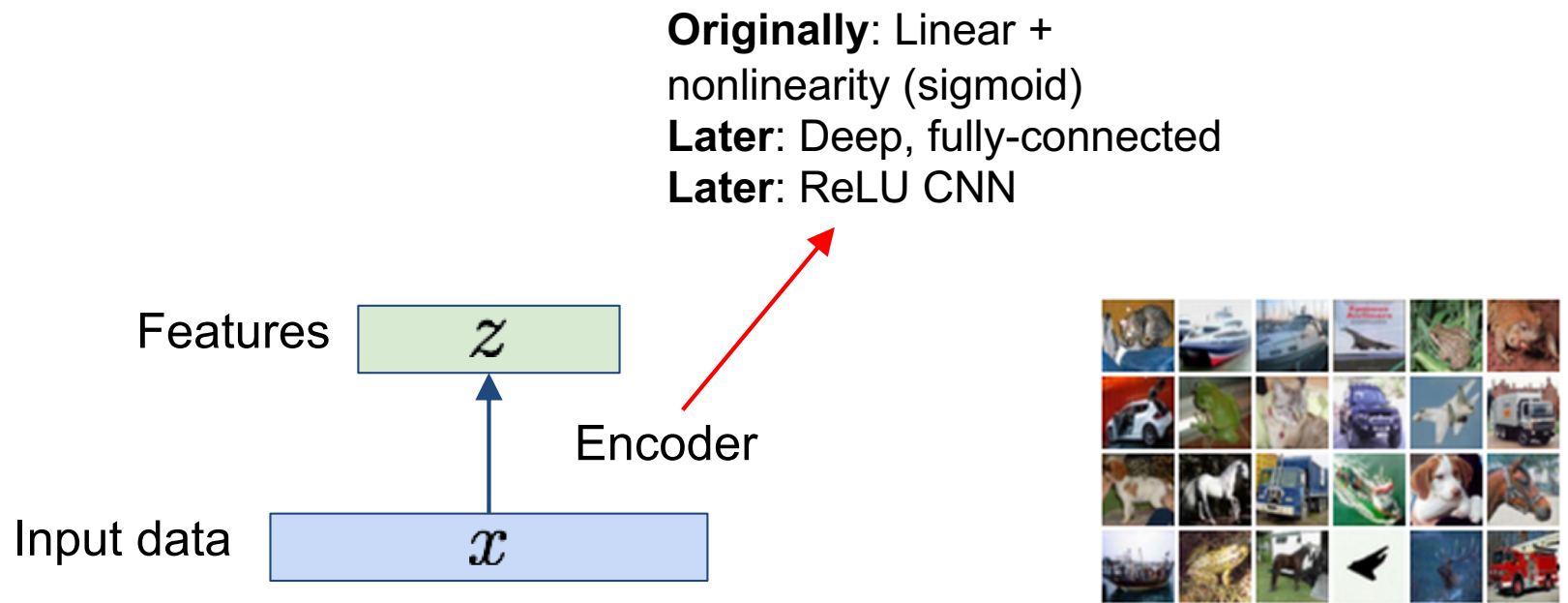
3. Amortized Inference Neural Networks

4. "Reparameterization" Trick

# Autoencoders

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

Features $z$ $\in \mathbb{R}^k$

↑ Encoder

Input data $x$ $\in \mathbb{R}^d$

# Autoencoders

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

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

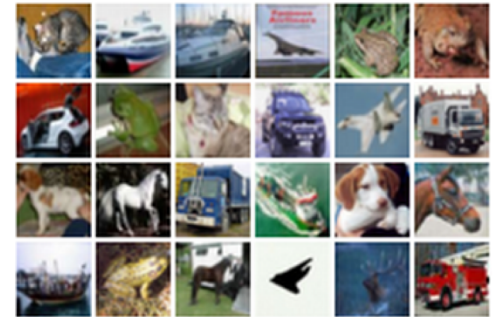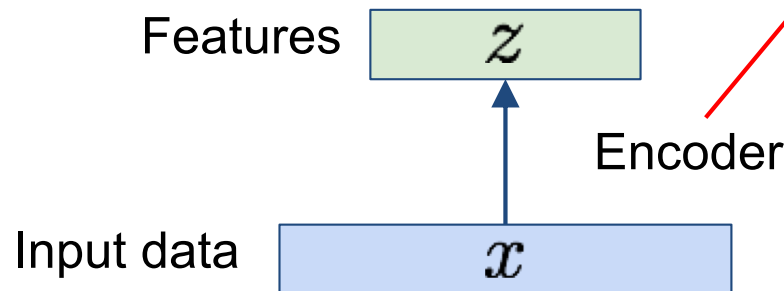Features $z$

Encoder

Input data $x$

# Autoencoders

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

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

Q: Why dimensionality reduction?

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

Features $z$

Input data $x$

Encoder

# Autoencoders

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

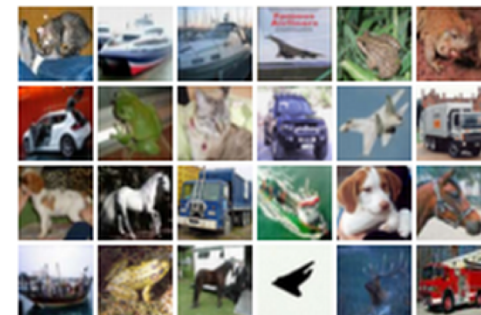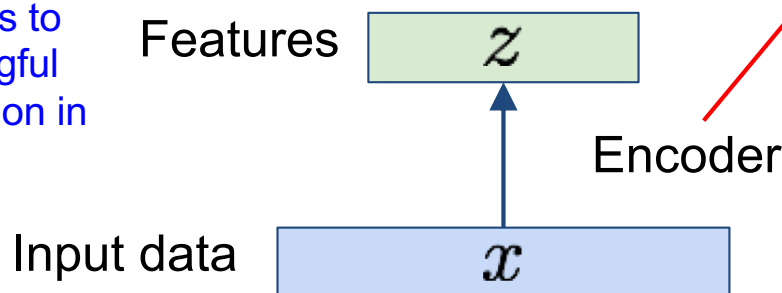**z** usually smaller than **x** (dimensionality reduction)

Q: Why dimensionality reduction?

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

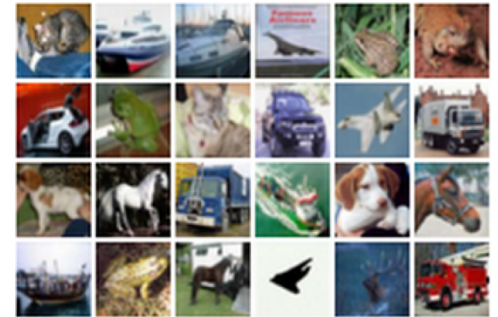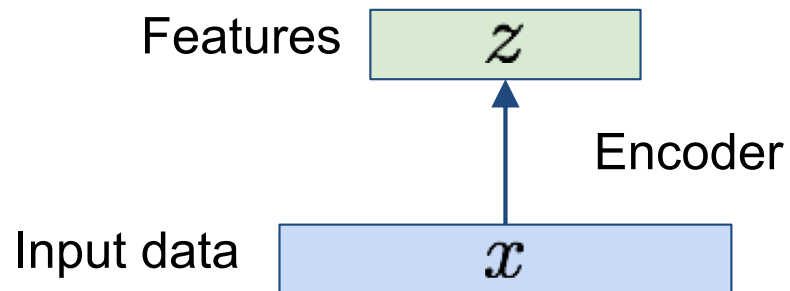**Originally**: Linear + nonlinearity (sigmoid)
**Later**: Deep, fully-connected
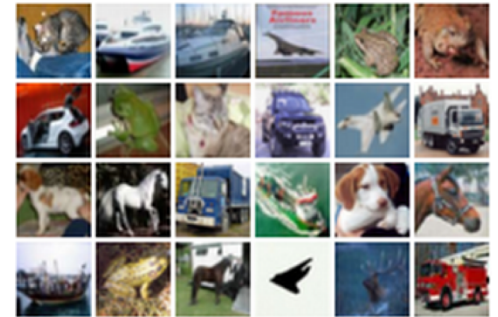**Later**: ReLU CNN

Features $z$

Input data $x$

Encoder

# Autoencoders

How to learn this feature representation?

Features $\quad$ $z$

Encoder

Input data $\quad$ $x$

# Autoencoders

How to learn this feature representation?
Train such that features can be used to reconstruct original data
"Autoencoding" - encoding itself

Reconstructed
input data          $\hat{x}$          $\mathbb{R}^d$

                              Decoder

Features            $z$          $\mathbb{R}^k$

                              Encoder

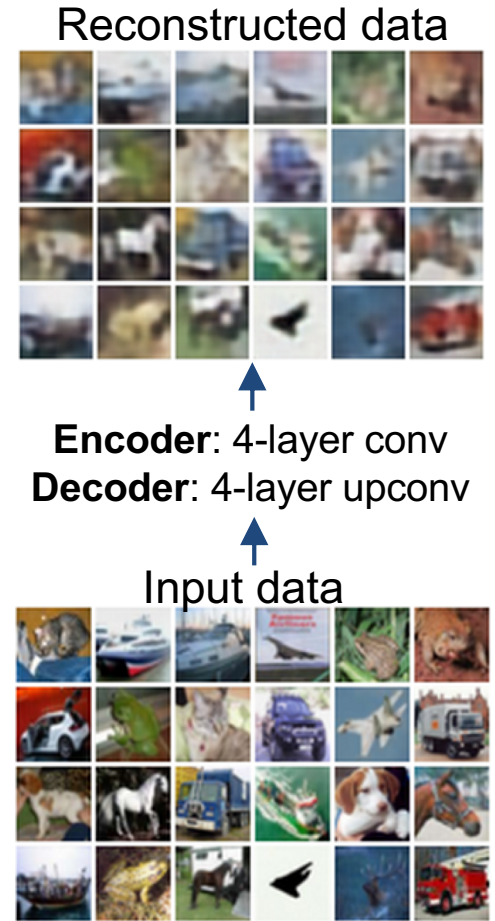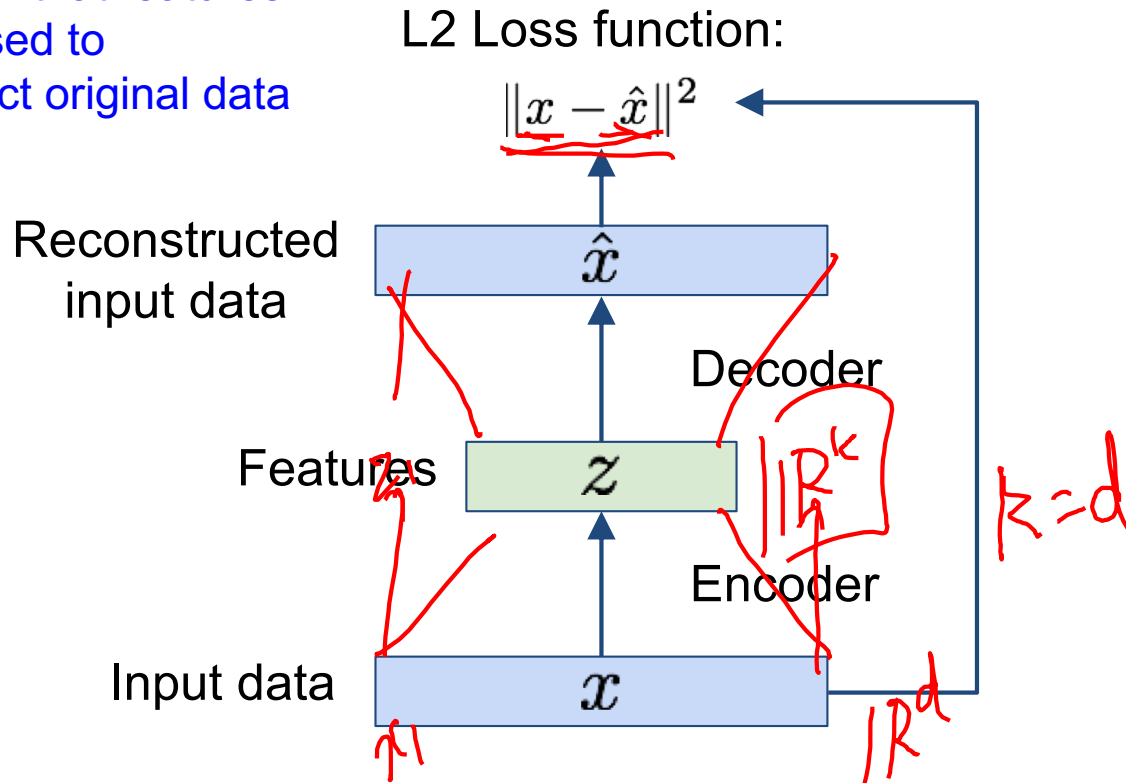Input data          $x$          $\mathbb{R}^d$

# Autoencoders

How to learn this feature representation?
Train such that features can be used to reconstruct original data
"Autoencoding" - encoding itself

Reconstructed data

Reconstructed input data — $\hat{x}$

Decoder

Features — $z$

Encoder

Input data — $x$

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

Input data

# Autoencoders

Reconstructed data

Train such that features can be used to reconstruct original data

L2 Loss function:

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

Reconstructed input data

$\hat{x}$

Decoder

Features

$z$

$\|\mathbb{R}^k$

$k = d$

Encoder

Input data

$x$

$\mathbb{R}^d$

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

Input data

# Autoencoders

Reconstructed data

Train such that features can be used to reconstruct original data

Doesn't use labels!

L2 Loss function:

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

Reconstructed input data

$\hat{x}$

Decoder

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

Features    $z$

Encoder

Input data

Input data    $x$

# Autoencoders

- Demo
  - https://cs.stanford.edu/people/karpathy/convnetjs/demo/autoencoder.html

# Autoencoders

Reconstructed
input data

$\hat{x}$

Decoder

Features

$z$

After training,
throw away decoder

Encoder

Input data

$x$

# Autoencoders

Loss function
(Softmax, etc)

bird        plane

dog        deer        truck

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

Classifier

Fine-tune
encoder
jointly with
classifier

Train for final task
(sometimes with
small data)

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

Features    $z$

Encoder

Input data    $x$

# Autoencoders
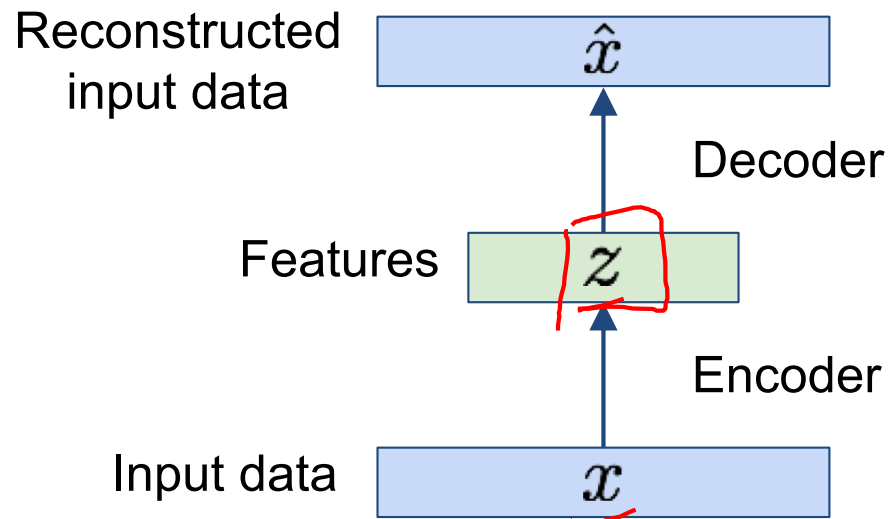
$z = f_\theta(x)$

$\hat{x} = g_\phi(z)$
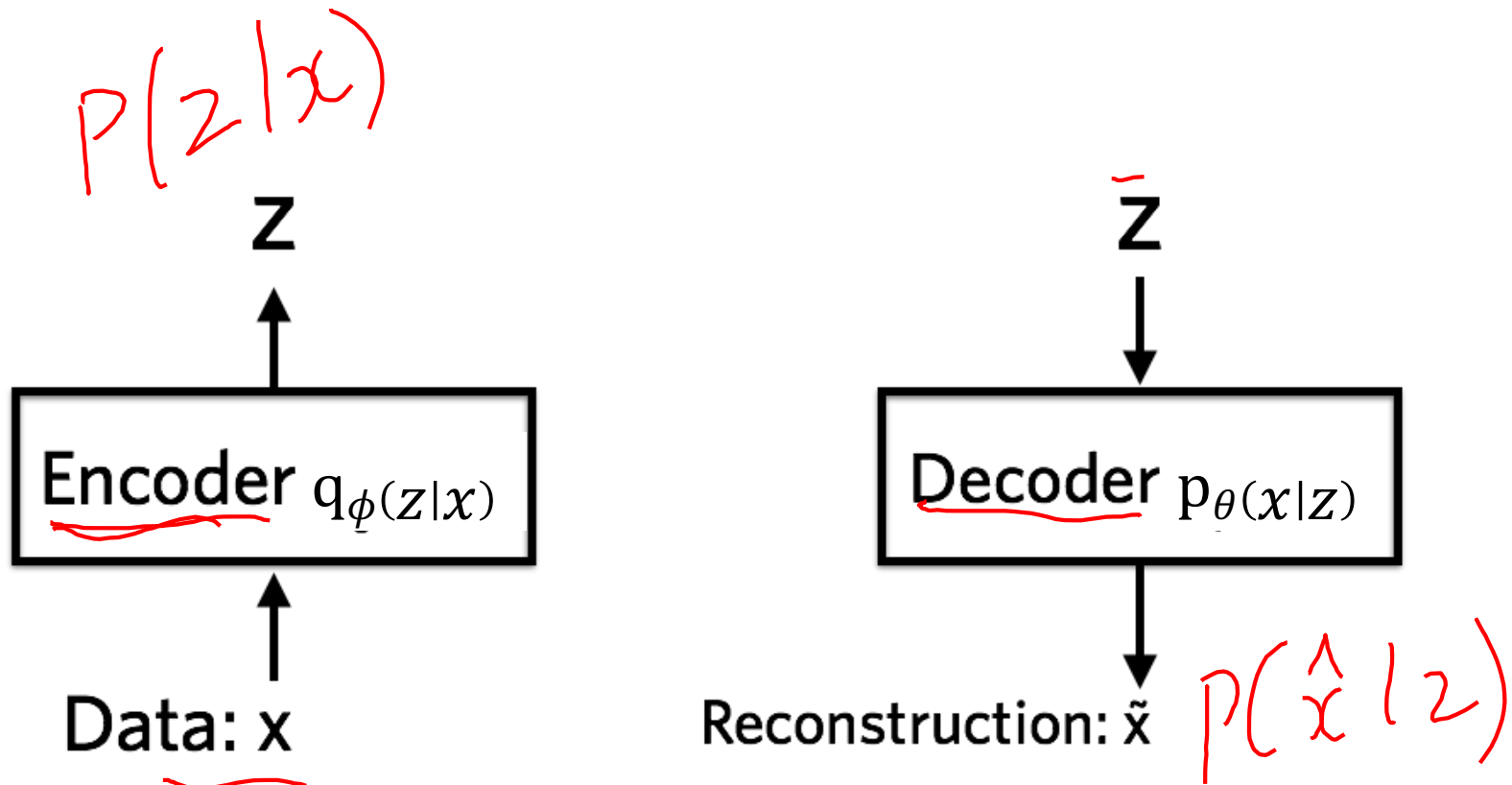
$P(z|x)$

$P(\hat{x}|z)$

$P(y|x)$

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?

Reconstructed input data
$\hat{x}$

Decoder

Features
$z$

Encoder

Input data
$x$

# Variational Autoencoders

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

$$p(z|x)$$

z

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

Data: x

$$\bar{z}$$

z

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

Reconstruction: x̃

$$p(\hat{x}|z)$$

# Variational Auto Encoders

VAEs are a combination of the following ideas:

1. Auto Encoders

2. Variational Approximation
   - Variational Lower Bound / ELBO

3. Amortized Inference Neural Networks

4. "Reparameterization" Trick

# Key problem

⊙2

- $P(z|x) = \dfrac{P(z, x)}{P(x)} = \dfrac{P(x|z)\, p(z)}{\boxed{\sum_z P(x|z)\, p(z)}}$

$\int$ Hard

$q(z)$

$\vec{x}$
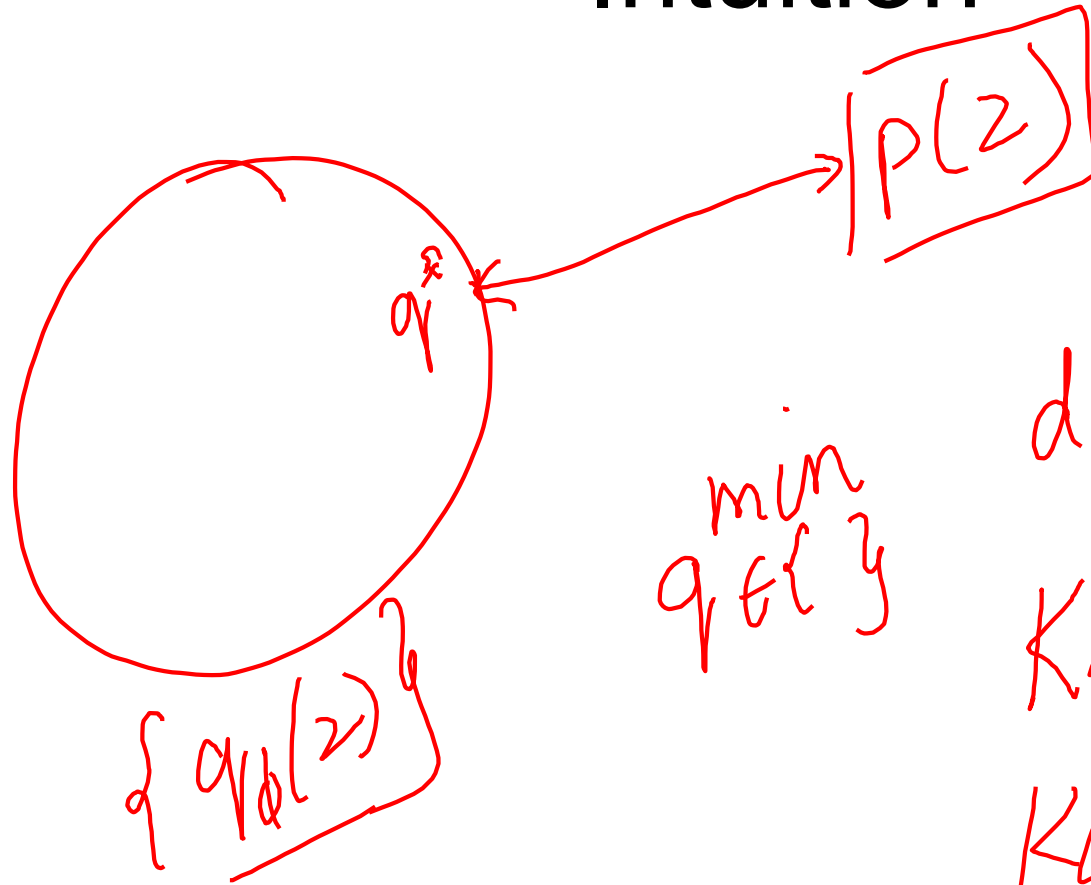
# What is Variational Inference?

- A class of methods for
  - approximate inference, parameter learning
  - and approximating integrals basically..

- Key idea
  - Reality is complex
  - Instead of performing approximate computation in something complex,
  - Can we perform exact computation in something "simple"?
  - Just need to make sure the simple thing is "close" to the complex thing.

# Intuition



$$p(z)$$

$$q^*_\phi$$

$$\{ q_\phi(z) \}$$

$$\min_{q \in \{\}}$$

$$d(p, q)$$

$$KL(p(z) \| q(z))$$

$$KL(q(z) \| p(z))$$

$$KL(p \| q) = \sum_z \underline{p(z)} \log \underbrace{\frac{p(z)}{q(z)}}_{error}$$

$$\sum \underline{q(z)} \log \frac{q(z)}{p(z)}$$
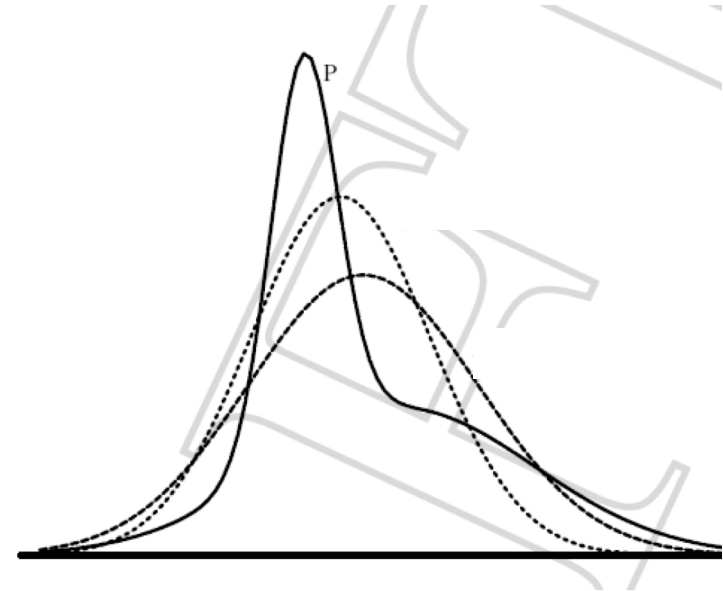
Support

# KL divergence:
# Distance between distributions

- Given two distributions $p$ and $q$ KL divergence:


- $D(p||q) = 0$ iff p=q


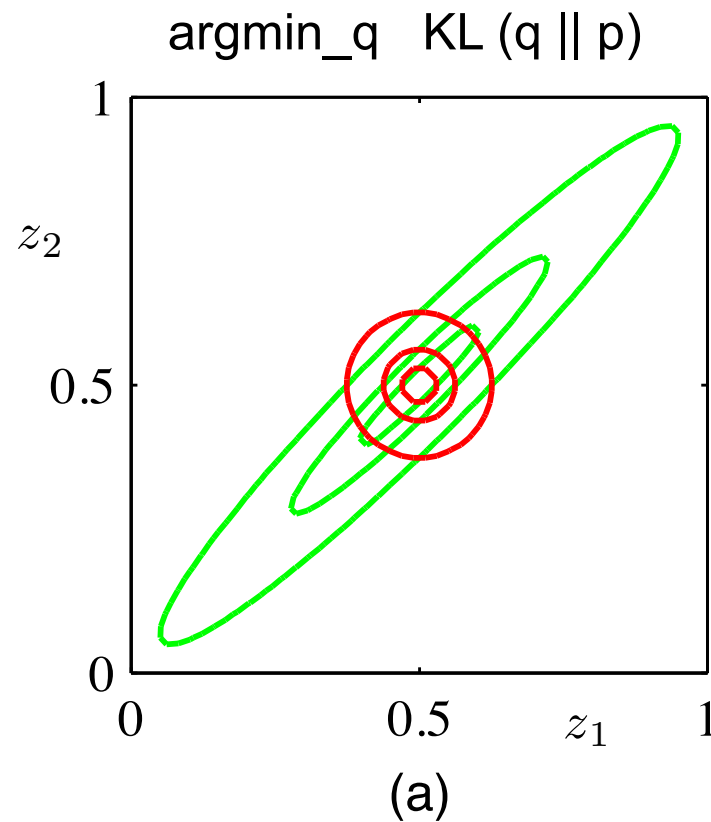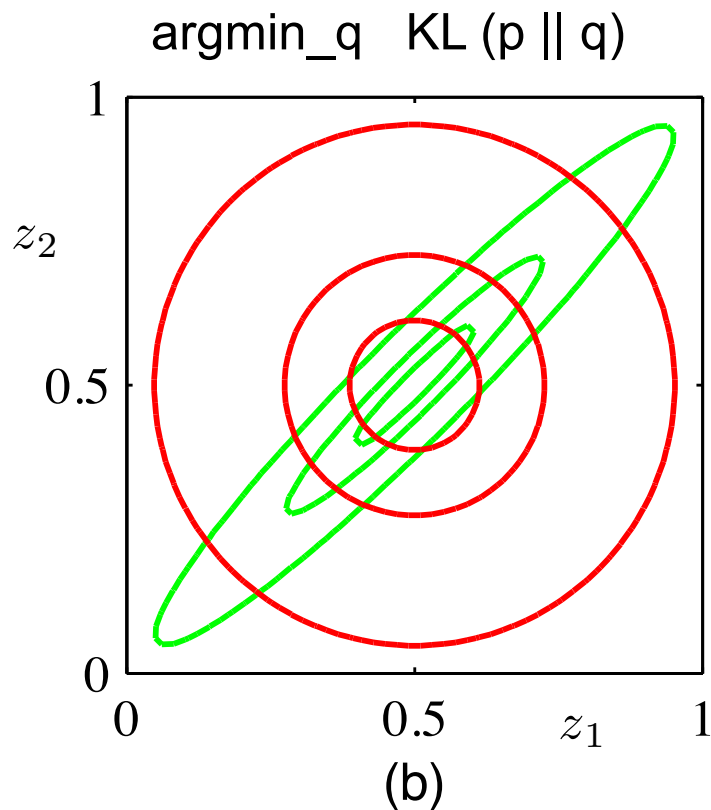- Not symmetric – p determines where difference is important

# Find simple approximate distribution

- Suppose *p* is intractable posterior
- Want to find simple *q* that approximates *p*
- KL divergence not symmetric

- D(p||q)
  – true distribution p defines support of diff.
  – the "correct" direction
  – will be intractable to compute

- D(q||p)
  – approximate distribution defines support
  – tends to give overconfident results
  – will be tractable

# Example 1

- p = 2D Gaussian with arbitrary co-variance
- q = 2D Gaussian with diagonal co-variance



argmin_q   KL (p || q)

(b)

argmin_q   KL (q || p)

(a)

# Example 2

- p = Mixture of Two Gaussians
- q = Single Gaussian  $\mu, \sigma$

UL

"wrong"

SL

→ "right"

argmin_q  KL (p || q)

argmin_q  KL (q || p)

q(z)