

CS 4803 / 7643: Deep Learning

Topics:

- Variational Auto-Encoders (VAEs)
- Key Ideas
 - AEs, Variational Inference

Dhruv Batra
Georgia Tech

Administrativa

- HW3 out
 - Due: 11/06, 11:55pm
- Final project
 - No poster session
 - Webpage submission
 - Details out soon

Recap from last time

Overview

- Unsupervised Learning
- Generative Models
 - PixelRNN and PixelCNN
 - Variational Autoencoders (VAE)
 - Generative Adversarial Networks (GAN)

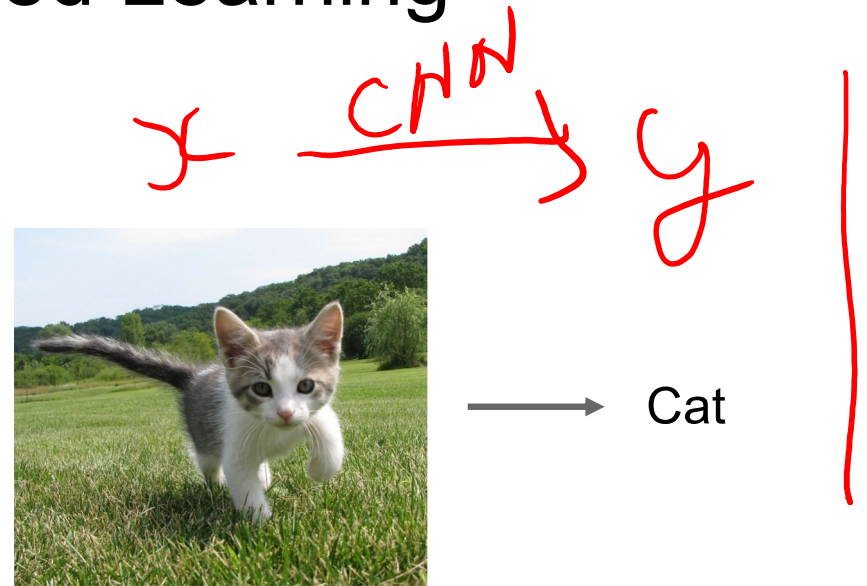
So far... Supervised Learning

Data: (x, y)

x is data, y is label

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

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



Classification

[This image](#) is [CC0 public domain](#)

Supervised vs Unsupervised Learning

Unsupervised Learning

Data: x

Just data, no labels!

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

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

Supervised vs Unsupervised Learning

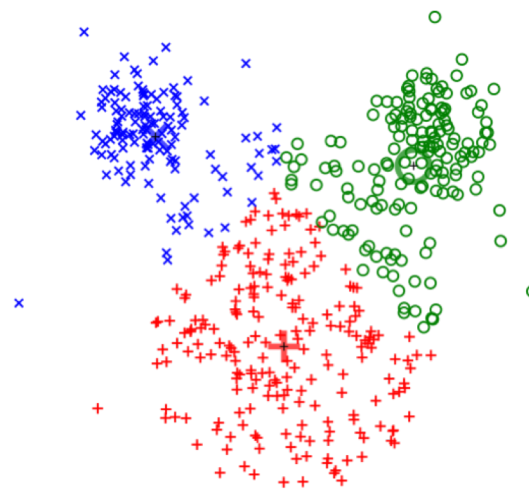
Unsupervised Learning

Data: x

Just data, no labels!

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

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



K-means clustering

[This image is CC0 public domain](#)

Supervised vs Unsupervised Learning

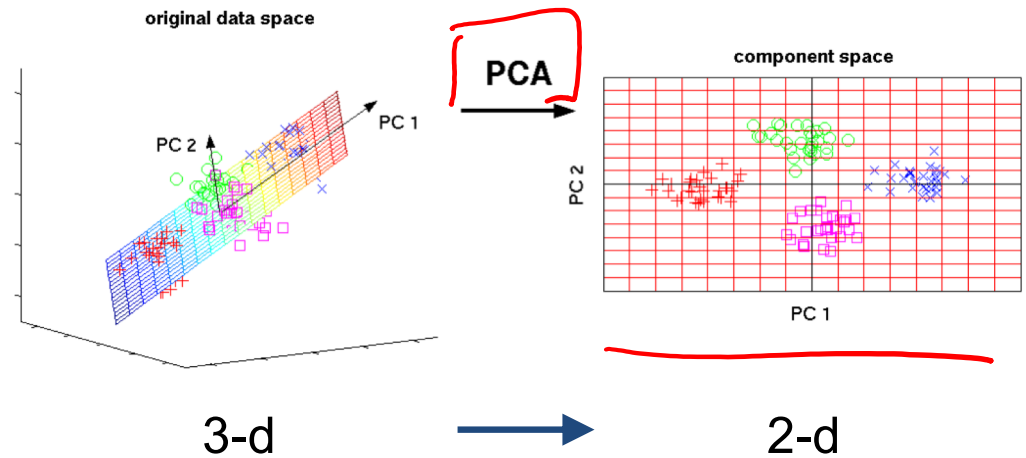
Unsupervised Learning

Data: x

Just data, no labels!

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

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



Principal Component Analysis
(Dimensionality reduction)

[This image](#) from Matthias Scholz is [CC0 public domain](#)

~~$P(x)$~~ $P(x, \dots, y)$

Supervised vs Unsupervised Learning

Unsupervised Learning

Data: x

Just data, no labels!

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

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

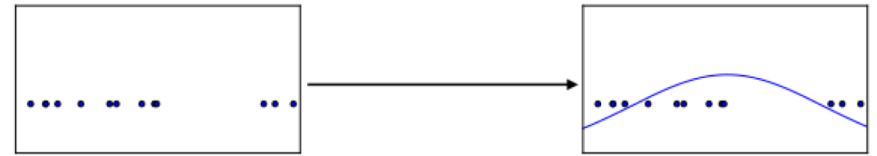
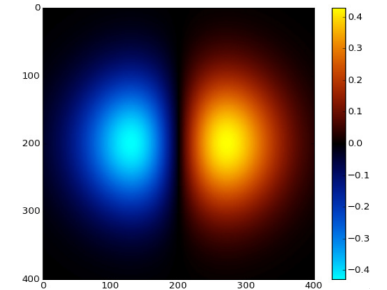
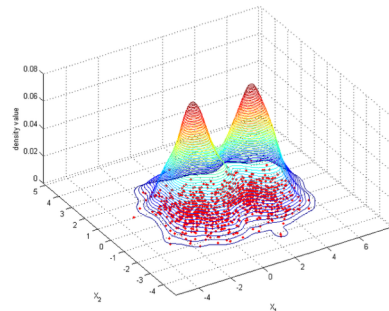


Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation

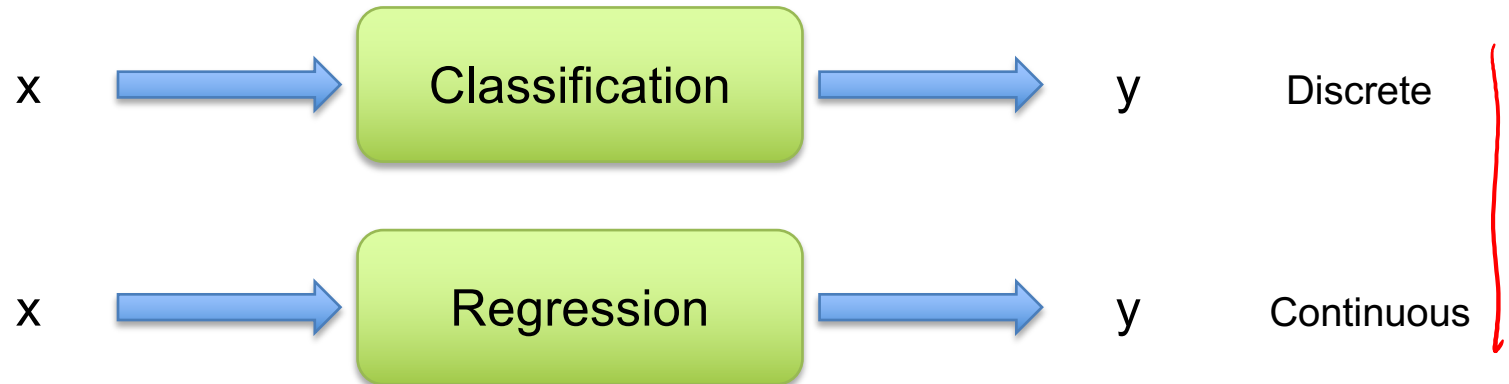


2-d density estimation

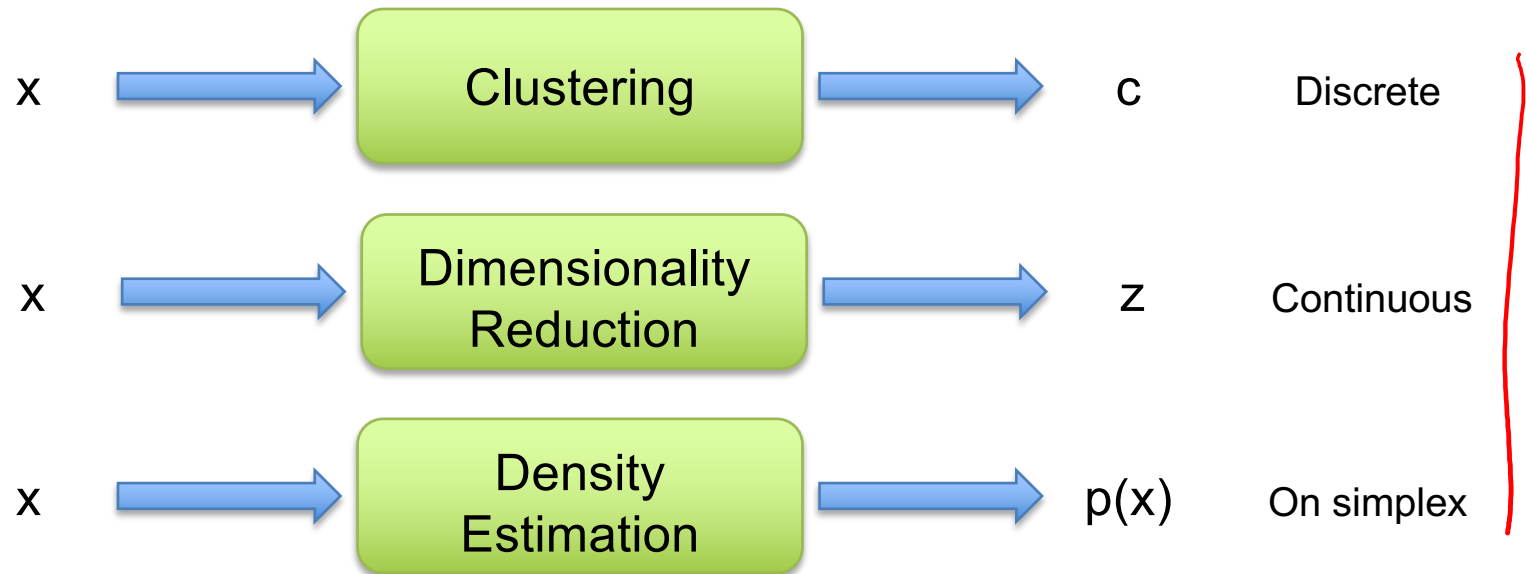
2-d density images [left](#) and [right](#) are [CC0 public domain](#).

Tasks

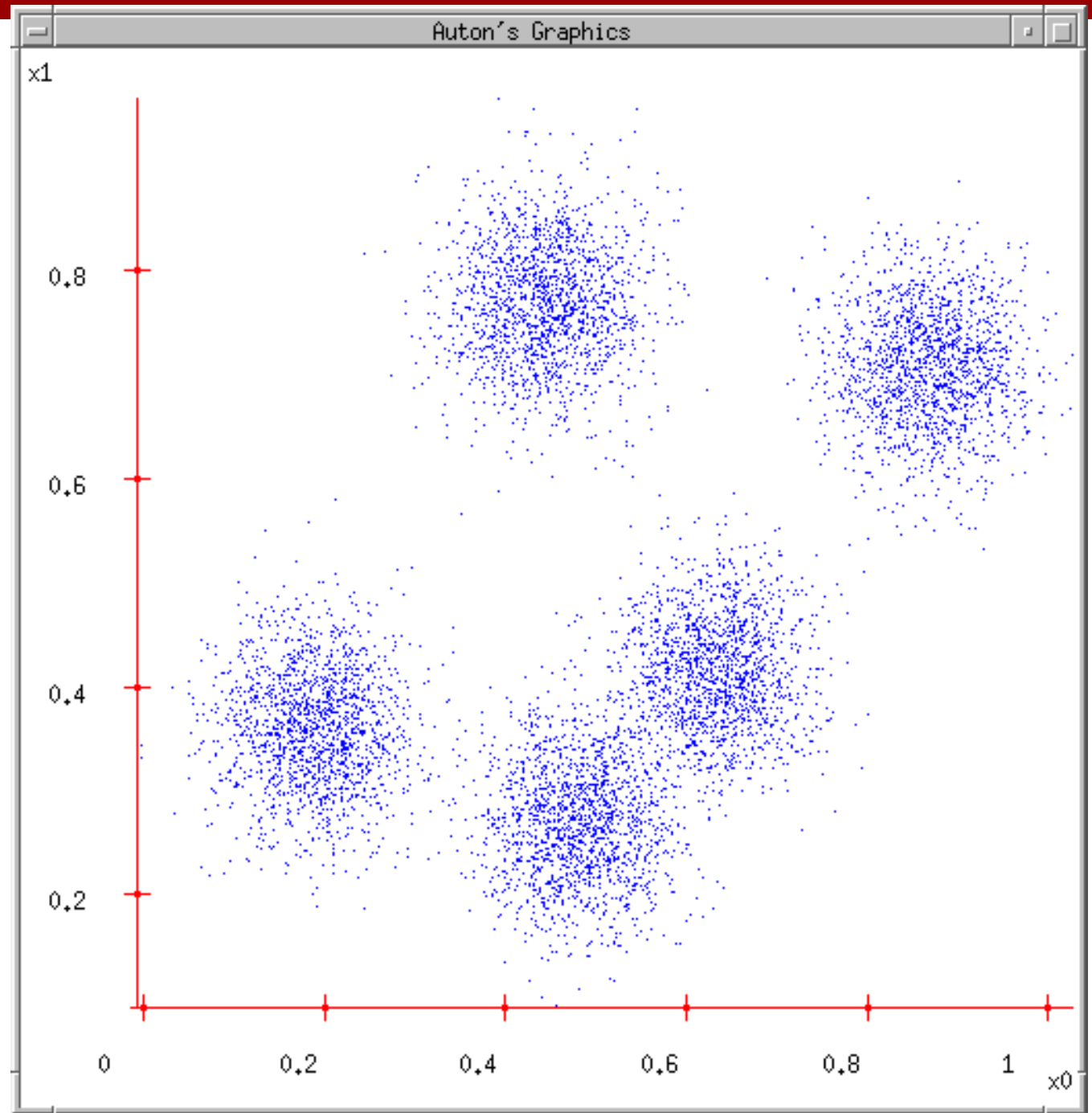
Supervised Learning



Unsupervised Learning

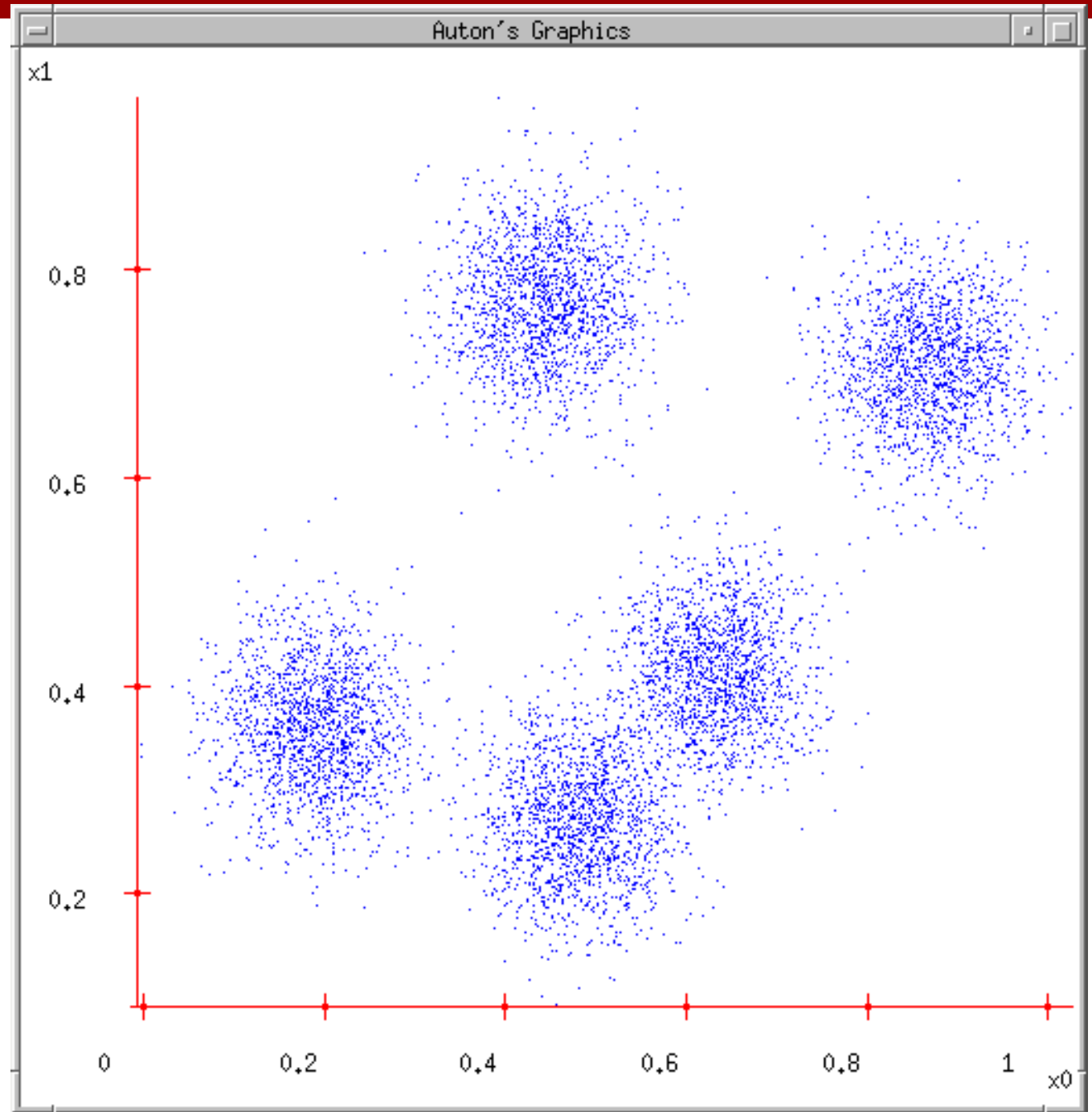


Some Data



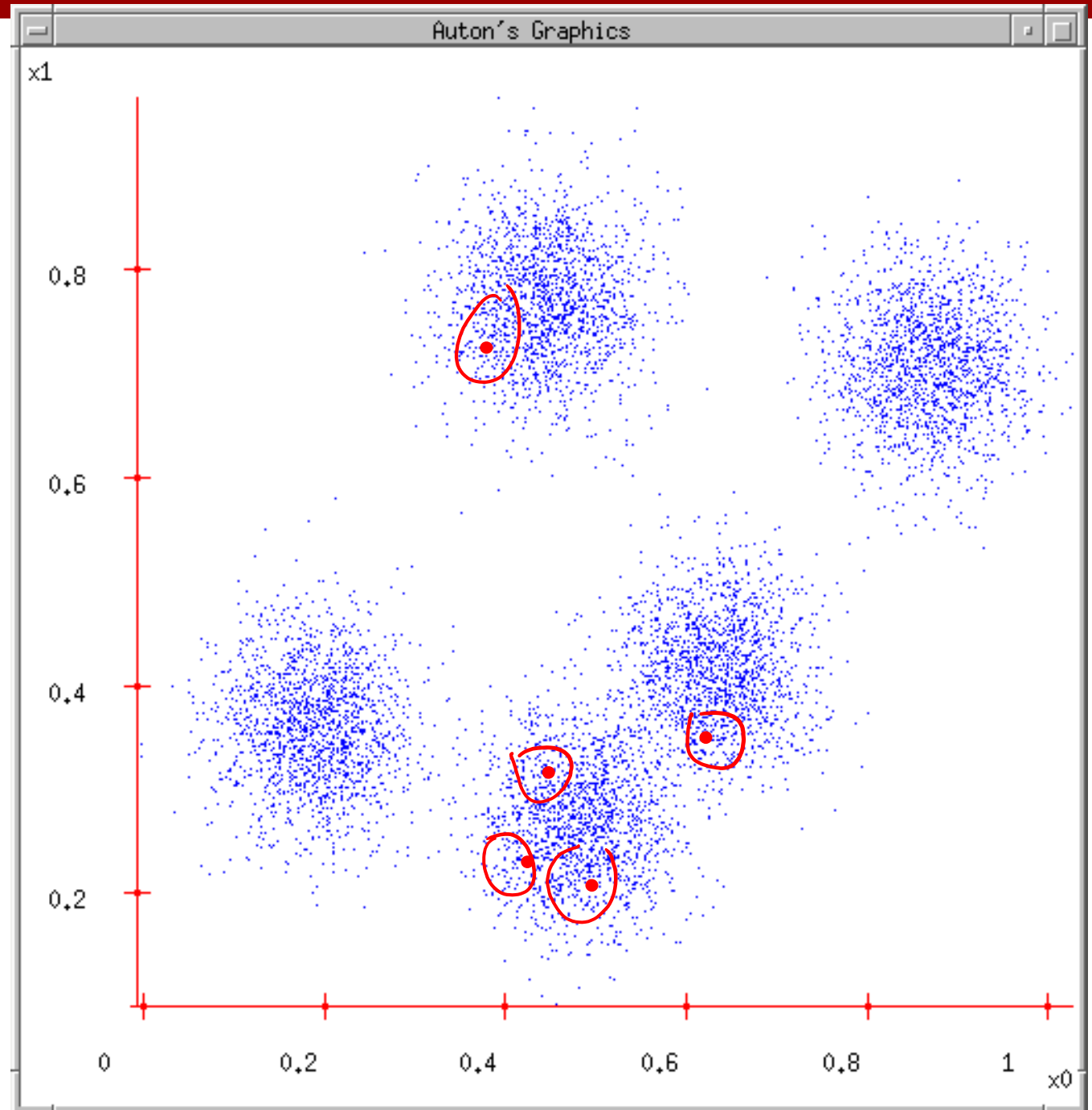
K-means

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



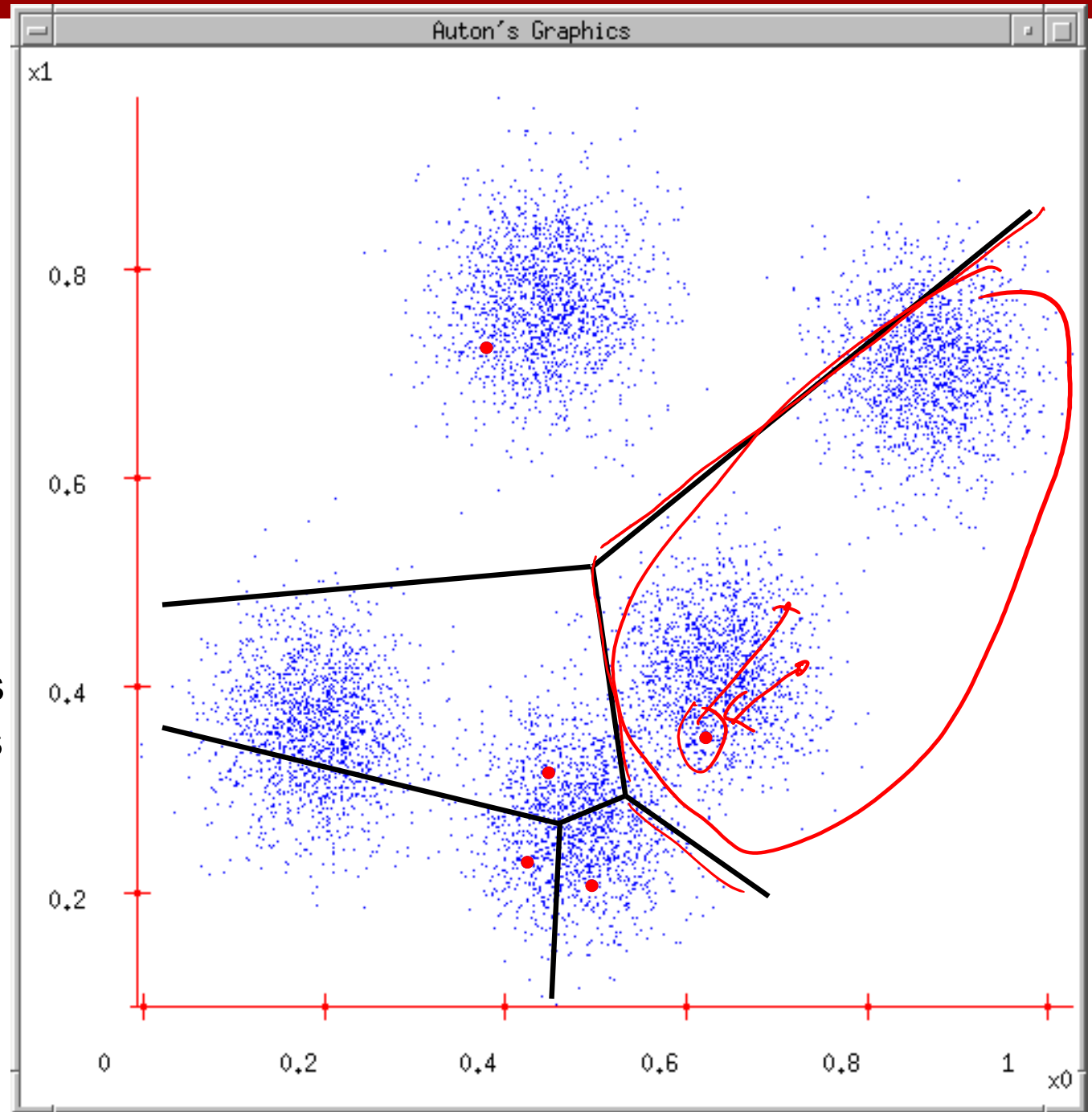
K-means

1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations



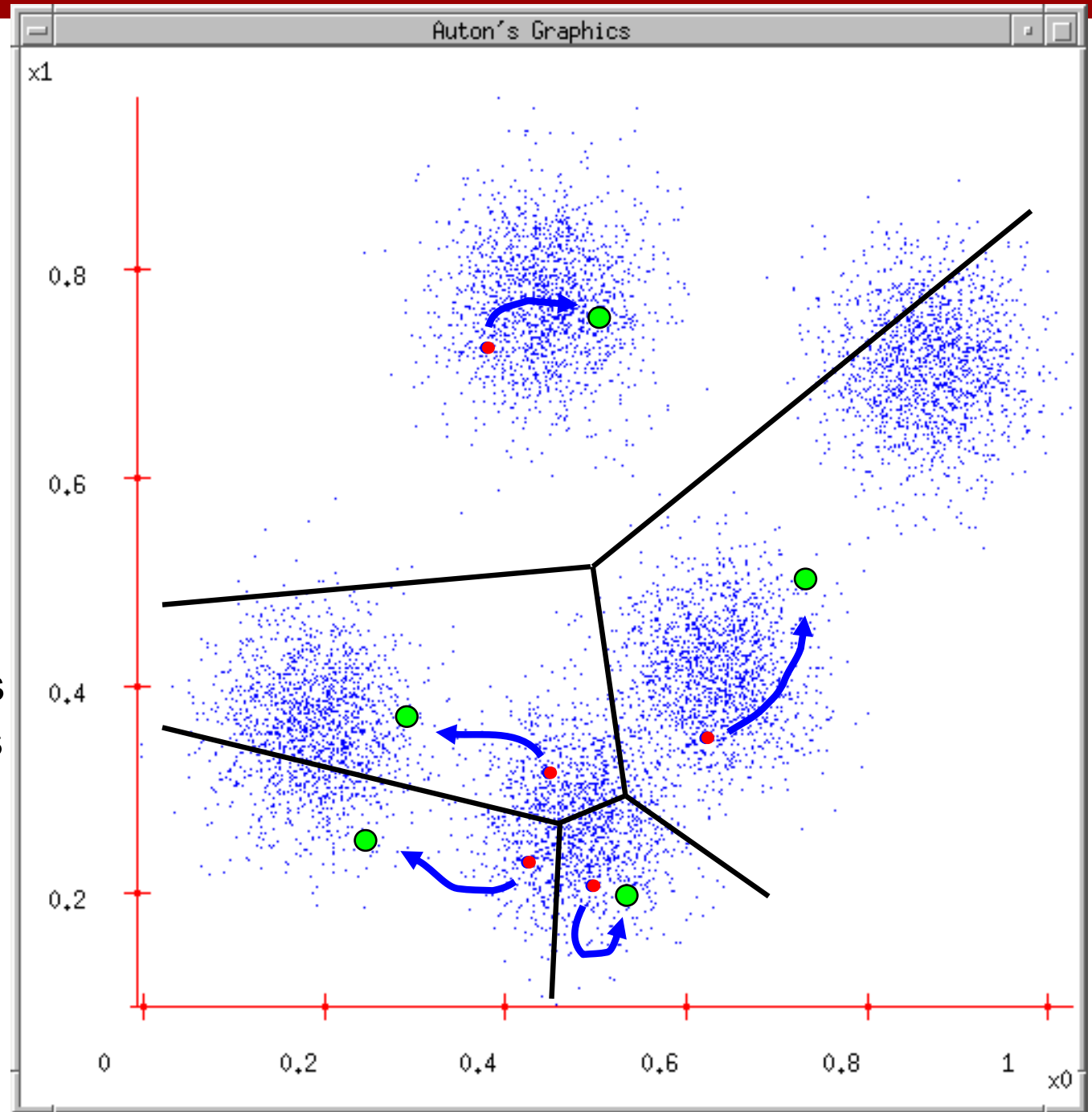
K-means

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. (Thus each Center "owns" a set of datapoints)



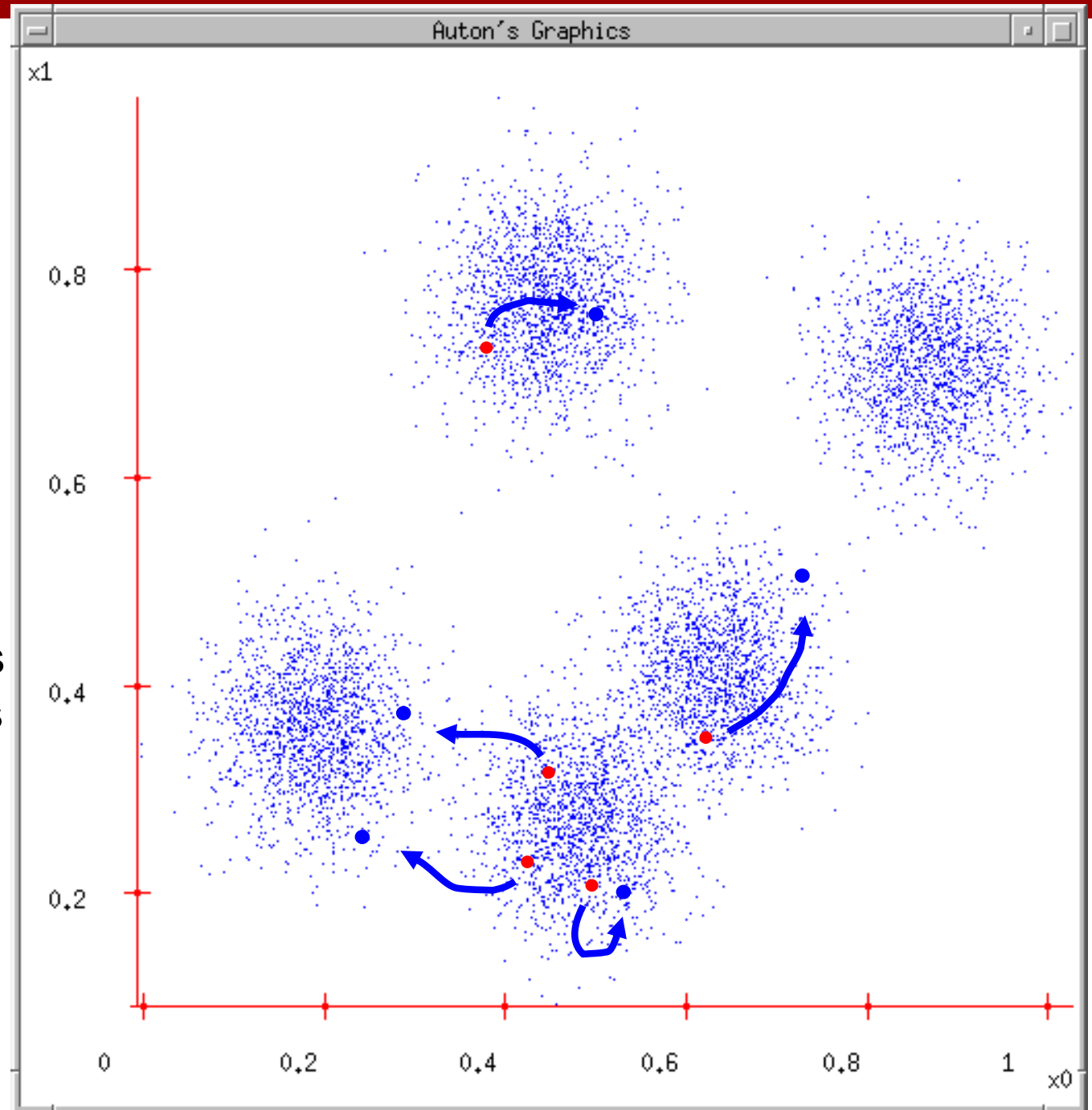
K-means

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



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!



K-means

- Randomly initialize k centers
 - $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$
- **Assign:**
 - Assign each point $i \in \{1, \dots, n\}$ to nearest center:
 - $C(i) \leftarrow \underset{j}{\operatorname{argmin}} \|\mathbf{x}_i - \mu_j\|^2$
- **Recenter:**
 - μ_j becomes centroid of its points

What is K-means optimizing?

- Objective $F(\mu, C)$: function of centers μ and point allocations C :

$$- F(\mu, C) = \sum_{i=1}^N \|\mathbf{x}_i - \mu_{C(i)}\|^2$$

- 1-of-k encoding

$$F(\mu, \mathbf{a}) = \sum_{i=1}^N \sum_{j=1}^k a_{ij} \|\mathbf{x}_i - \mu_j\|^2$$

- Optimal K-means:

$$- \min_{\mu} \min_{\mathbf{a}} F(\mu, \mathbf{a})$$

K-means as Co-ordinate Descent

- Optimize objective function:

$$\min_{\mu_1, \dots, \mu_k} \min_{a_1, \dots, a_N} F(\mu, a) = \min_{\mu_1, \dots, \mu_k} \min_{a_1, \dots, a_N} \sum_{i=1}^N \sum_{j=1}^k a_{ij} \|\mathbf{x}_i - \mu_j\|^2$$

- Fix μ , optimize a (or C)

K-means as Co-ordinate Descent

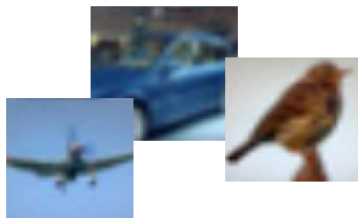
- Optimize objective function:

$$\min_{\mu_1, \dots, \mu_k} \min_{\mathbf{a}_1, \dots, \mathbf{a}_N} F(\boldsymbol{\mu}, \mathbf{a}) = \min_{\mu_1, \dots, \mu_k} \min_{\mathbf{a}_1, \dots, \mathbf{a}_N} \sum_{i=1}^N \sum_{j=1}^k a_{ij} \|\mathbf{x}_i - \underline{\mu}_j\|^2$$

- Fix a (or C) optimize μ

Generative Models

Given training data, generate new samples from same distribution



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



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

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

$p_{\text{model}}(x)$

Taxonomy of Generative Models

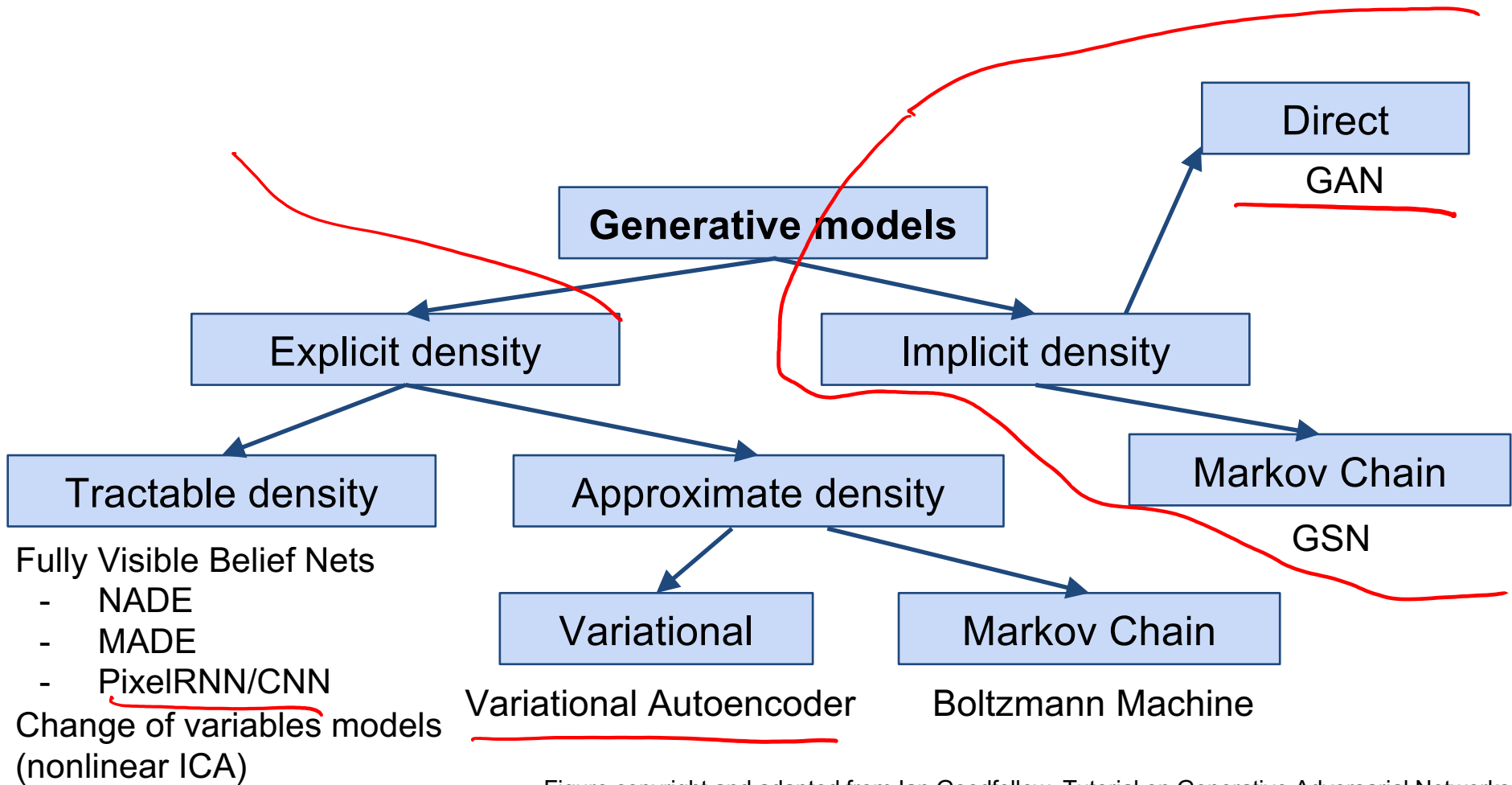


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

We will discuss 3 most popular types of generative models

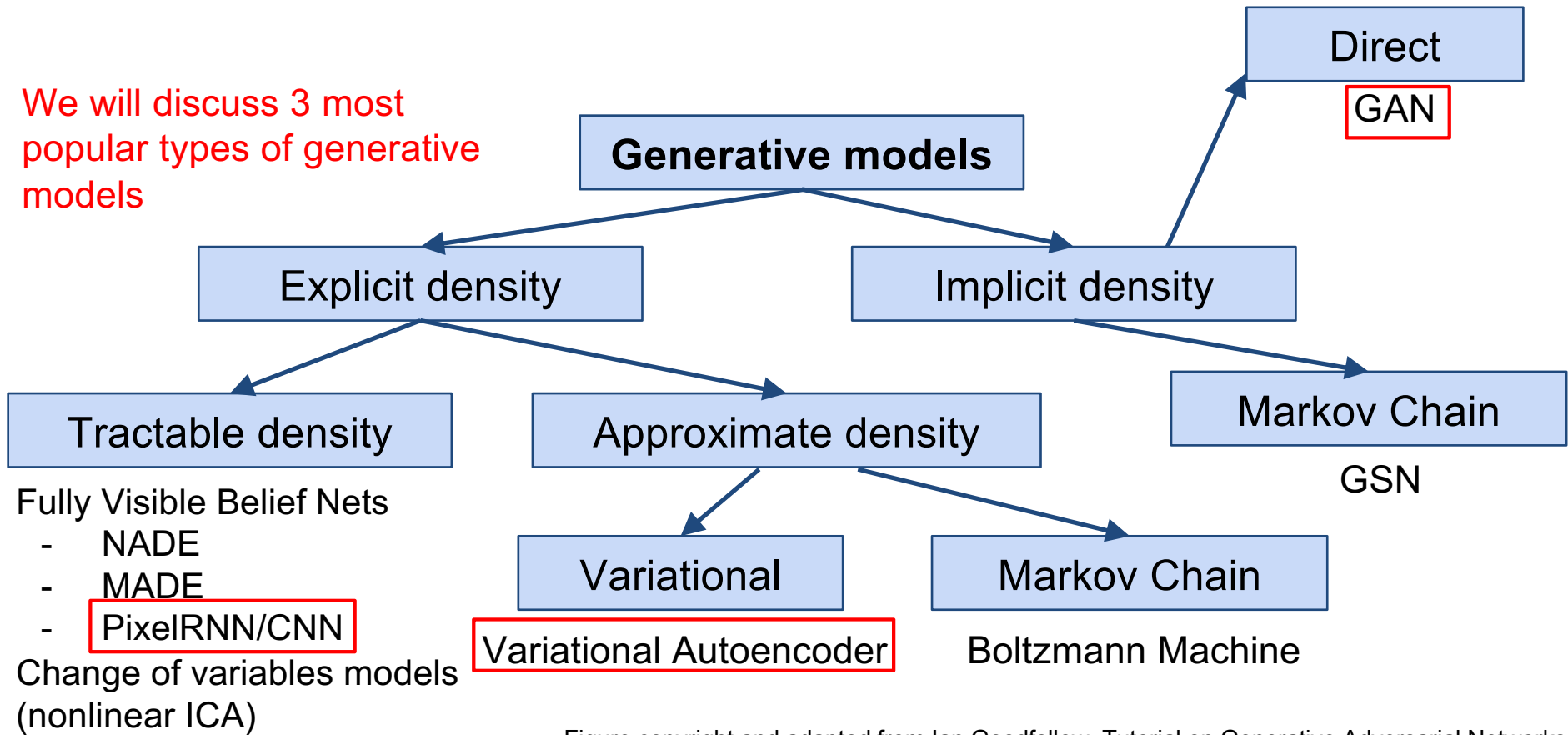
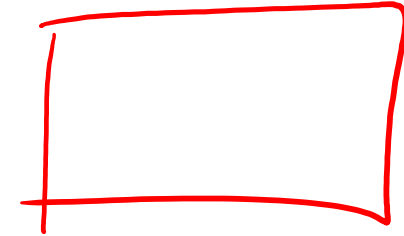


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

PixelRNN and PixelCNN

Fully Observable Model

Explicit density model



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

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

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

Then maximize likelihood of training data

$$D = \{x_1, \dots, x_n\}$$

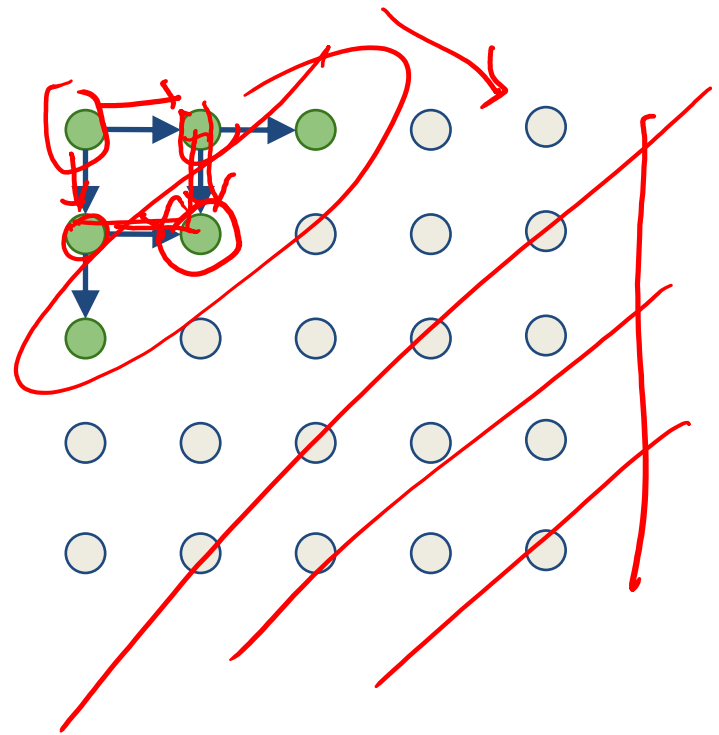
$$x_1, \dots, x_{i-1} \rightarrow \boxed{\text{NN}}_{\theta} \rightarrow x_i \quad \max_{\theta} \log P_{\theta}(x)$$

PixelRNN [van der Oord et al. 2016]

$$P(x_1 \dots x_{H \times W})$$

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)



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

Still generate image pixels starting from corner

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

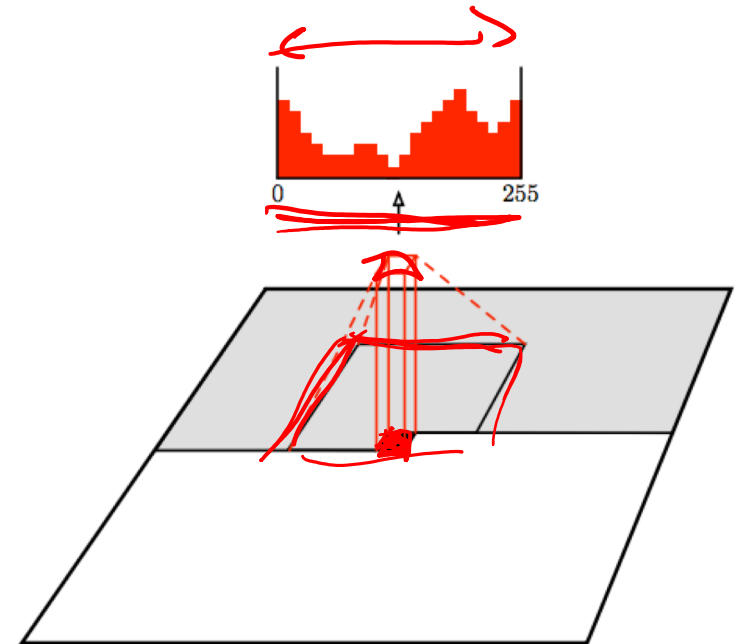


Figure copyright van der Oord et al., 2016. Reproduced with permission.



Variational Autoencoders (VAE)



So far...

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

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

So far...

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

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

$$p(\vec{x})$$

VAEs define intractable density function with latent \mathbf{z} :

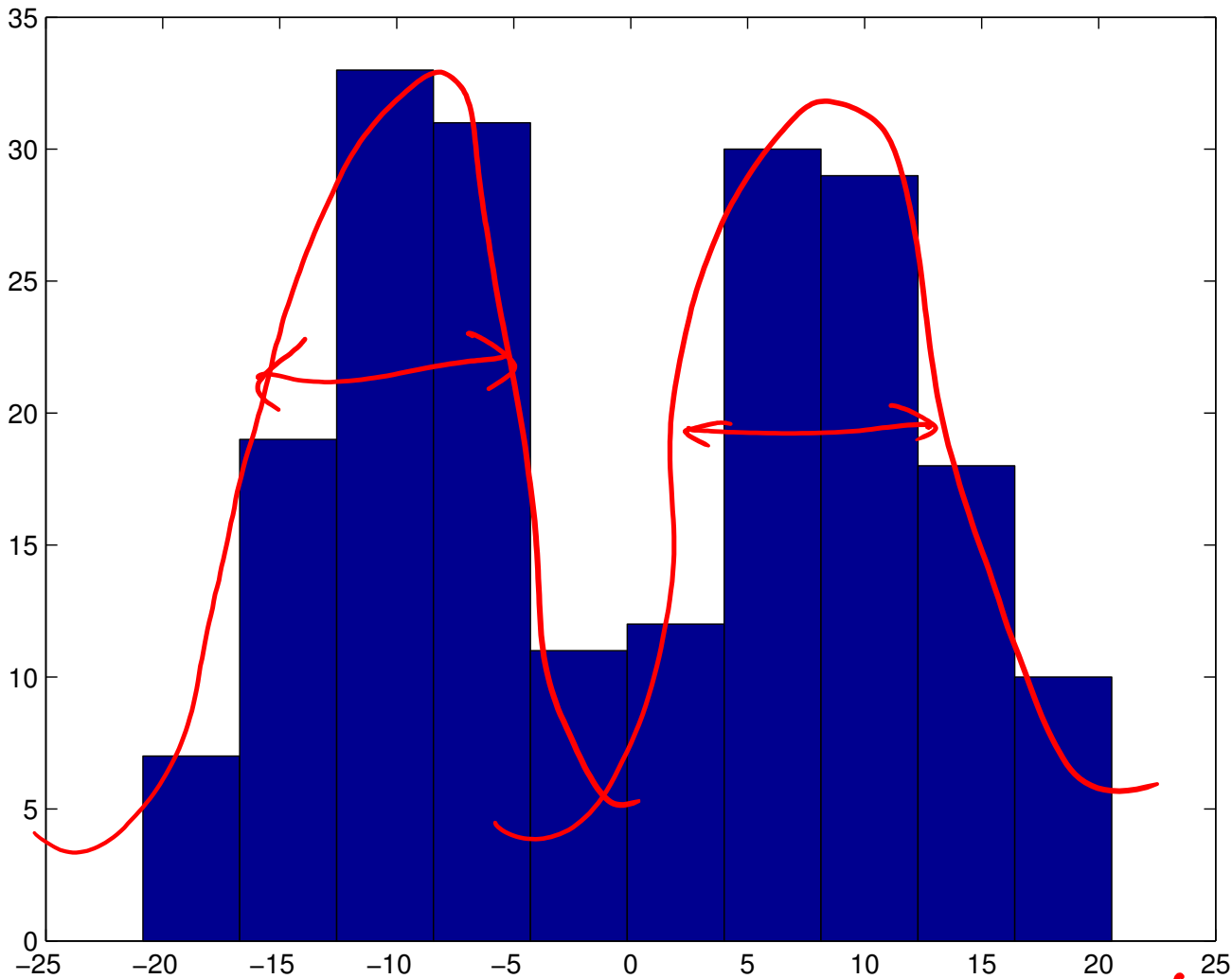
$$p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$$

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

$$p(\vec{x}, \mathbf{z}) = p(\vec{x} | \mathbf{z}) p(\mathbf{z})$$

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

GMM



$p(x)$

$p(x)$

$p(x|z)$

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

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

$p(x|z=1) \sim \mathcal{N}(\mu_1, \sigma_1^2)$
 $p(x|z=2) \sim \mathcal{N}(\mu_2, \sigma_2^2)$

Gaussian Mixture Model

$$Z \sim \text{Cat}(\vec{\pi}) \quad \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_k \end{bmatrix} \quad \left\{ \begin{array}{l} \pi_c \geq 0 \\ \sum_{c=1}^k \pi_c = 1 \end{array} \right. \quad \pi_c = P(Z=c)$$

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

$$\frac{P(Z) \cdot p(x|Z)}{\pi_c \cdot N(\mu_c, \sigma_c^2)} = p(x, z)$$

$p(z)$ Gaussian Mixture Model

$$p(x, z)$$

$$p(x) = \sum_c p(z) p(x|z) \leftarrow \text{Marginalization}$$

$$p(\underline{z} | \underline{x}) = \frac{p(z, x)}{p(x)} = \frac{p(x|z) p(z)}{\sum_z p(x|z) p(z)} \left| \begin{array}{l} \text{'Inference'} \\ \text{Posterior} \end{array} \right.$$

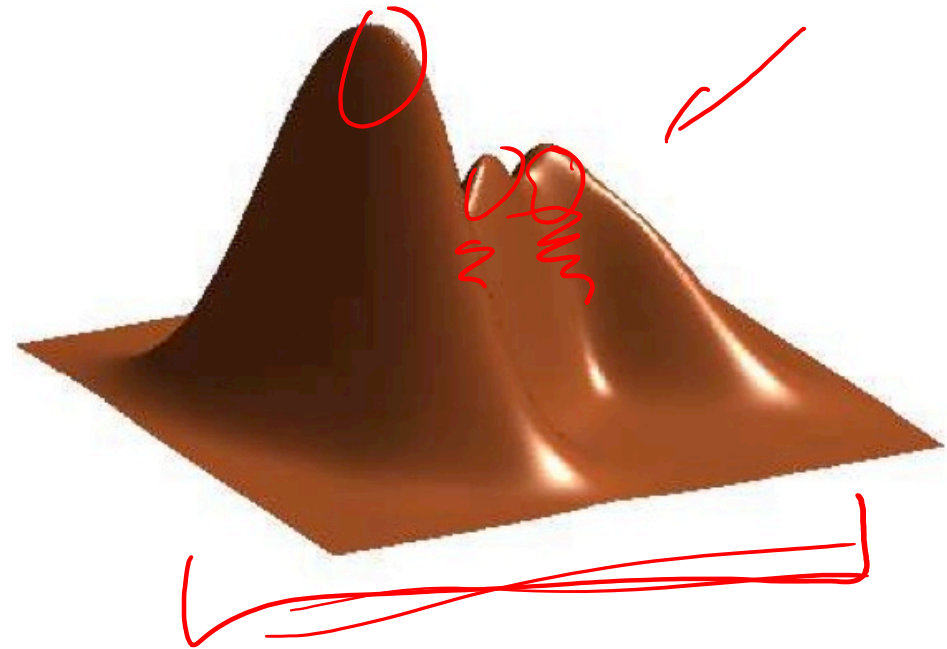
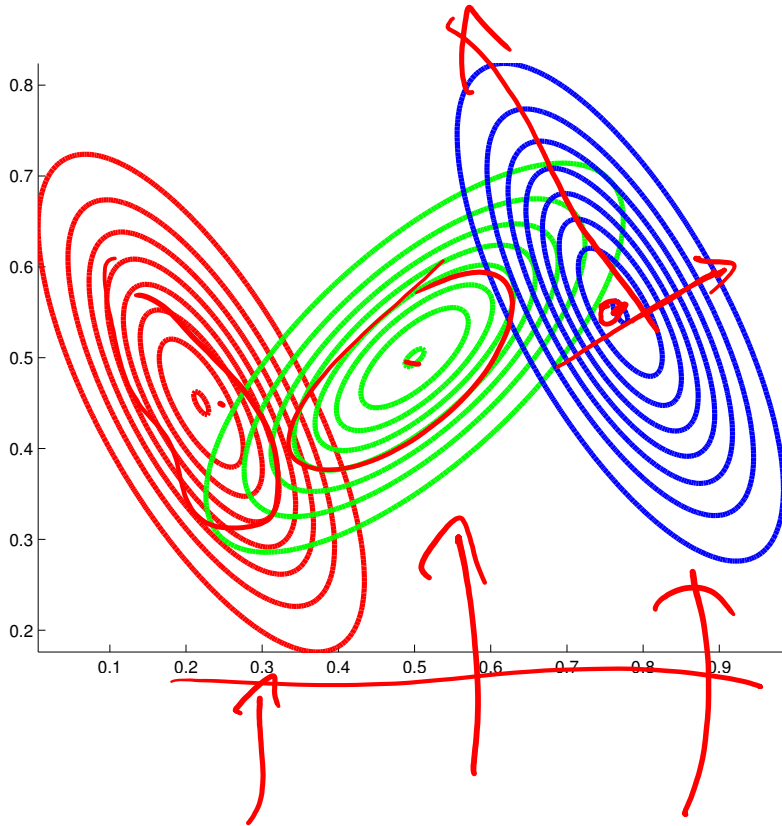
$$\sim p(\underline{x} | \underline{z})$$

~~$p(x)$~~

GMM

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

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



K-means vs GMM

- K-Means
 - <http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>
- GMM
 - <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!

$$\{\pi_1, \dots, \pi_k, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k\} \equiv \Theta$$

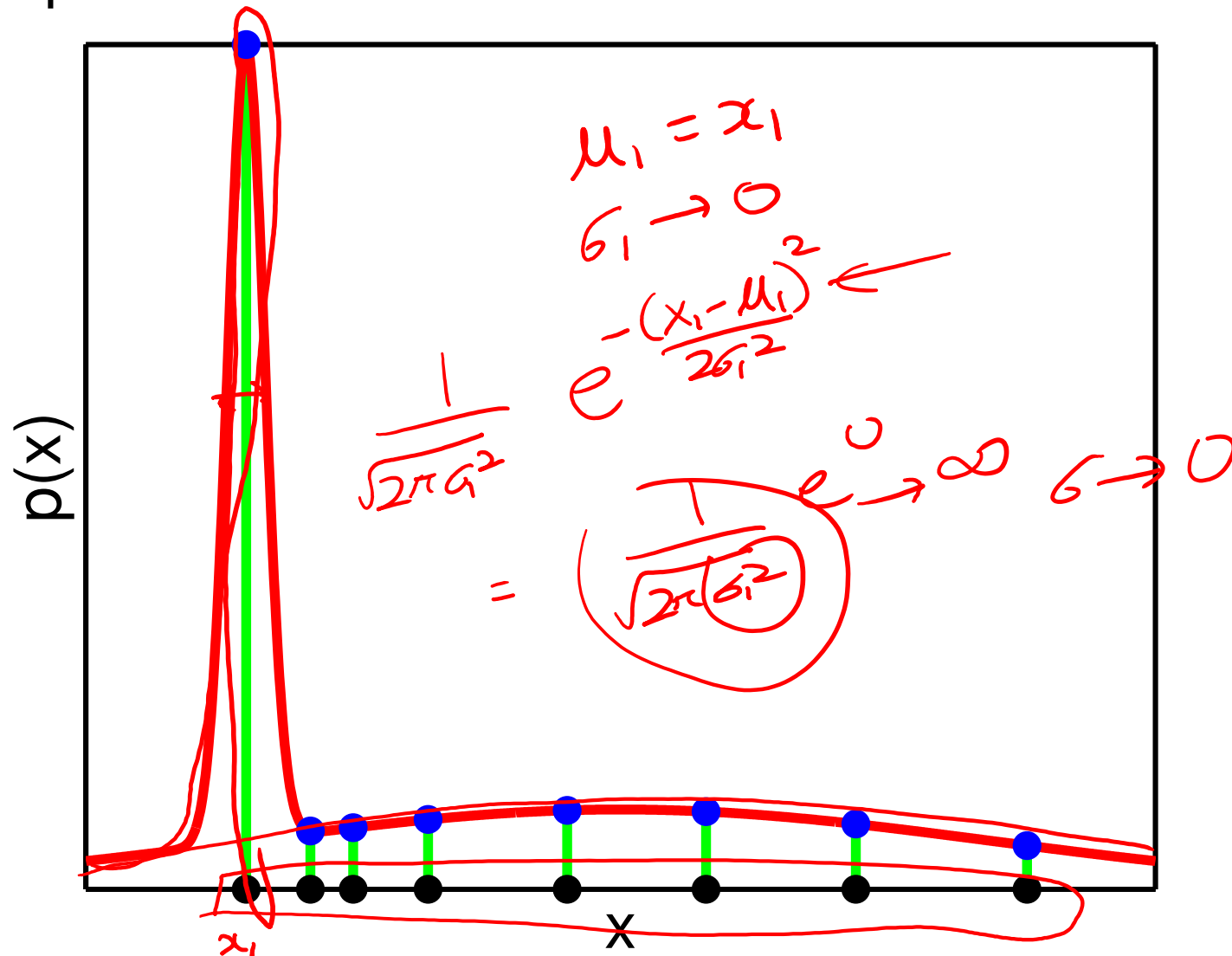
$$\{\vec{x}_1, \dots, \vec{x}_N\}$$

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} P(\underline{D} | \Theta) = \underset{\Theta}{\operatorname{argmax}} \sum_i \log \underline{P(\vec{x}_i | \Theta)}$$

$$\text{data } i \quad \sum \log \underbrace{\sum_z \underbrace{P(\vec{x}_i, z | \Theta)}_{P(x|z)P(z)}}_{\text{data } i}$$

Hidden Data Causes Problems #3

- Likelihood has singularities if one Gaussian “collapses”



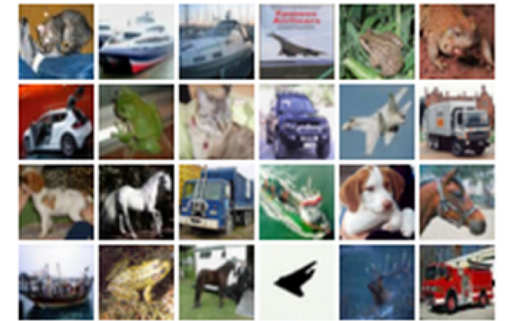
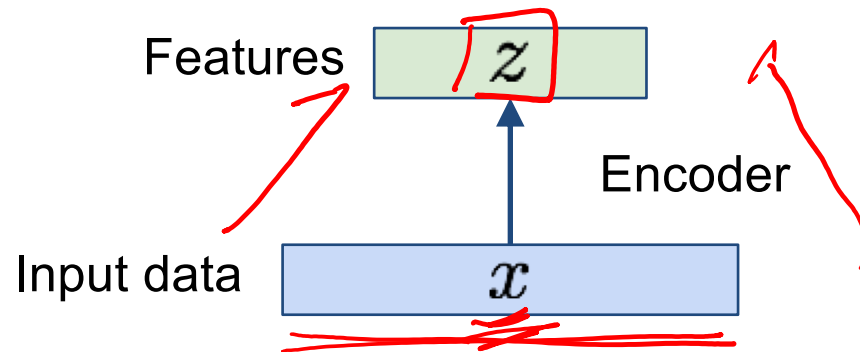
Variational Auto Encoders

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

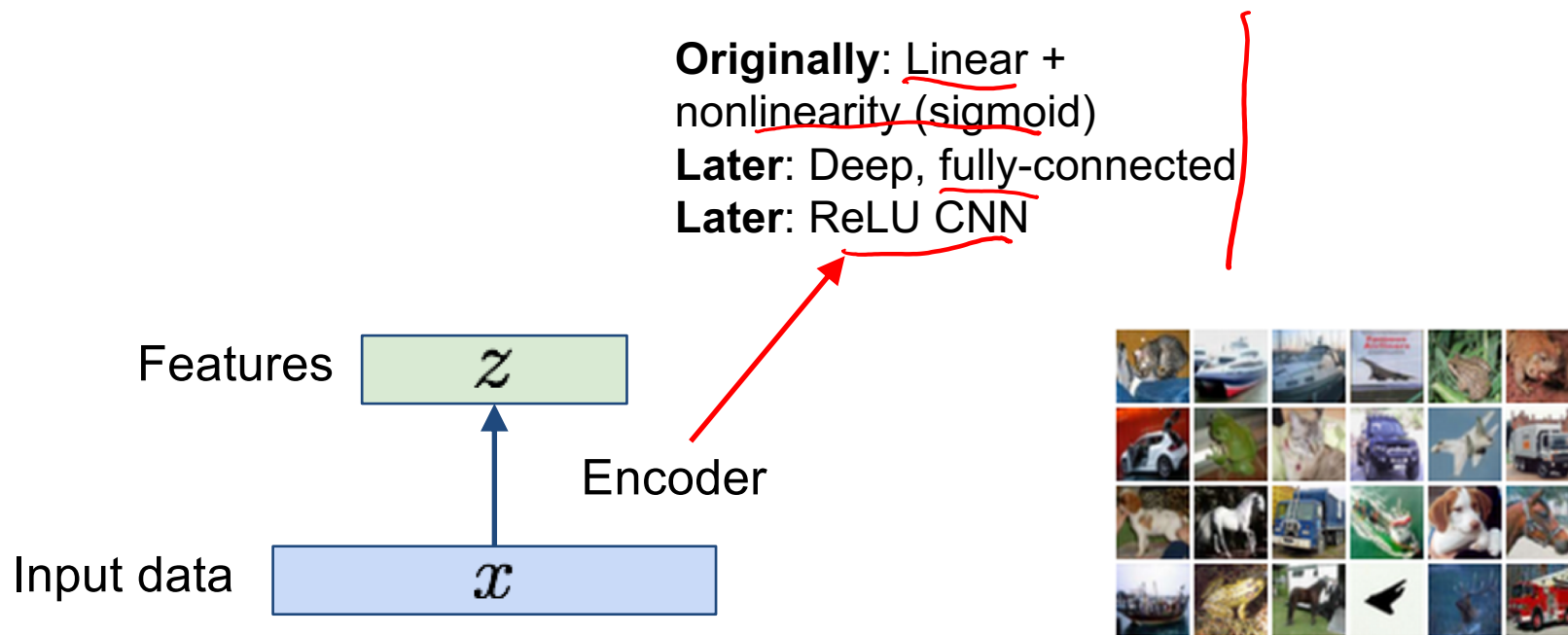
Autoencoders

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



Autoencoders

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



Autoencoders

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

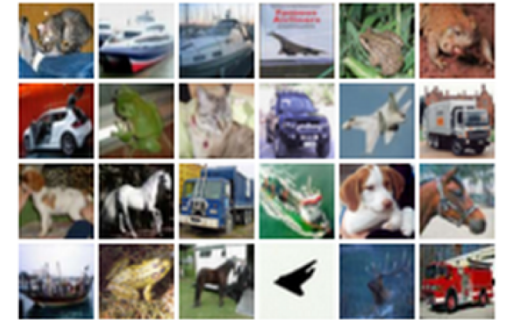
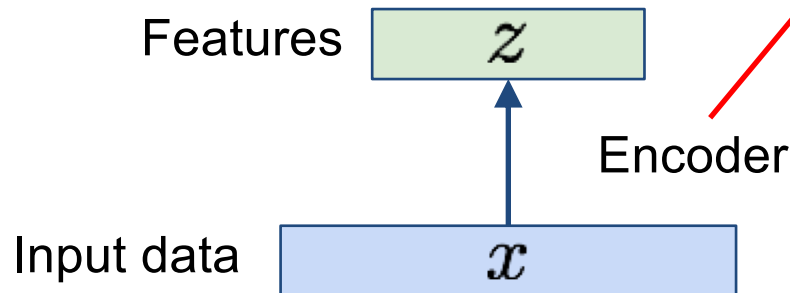
\mathbf{z} usually smaller than \mathbf{x}
(dimensionality reduction)

Q: Why dimensionality reduction?

Originally: Linear + nonlinearity (sigmoid)

Later: Deep, fully-connected

Later: ReLU CNN



Autoencoders

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

\mathbf{z} usually smaller than \mathbf{x}
(dimensionality reduction)

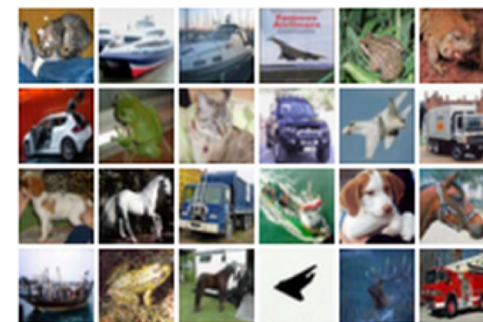
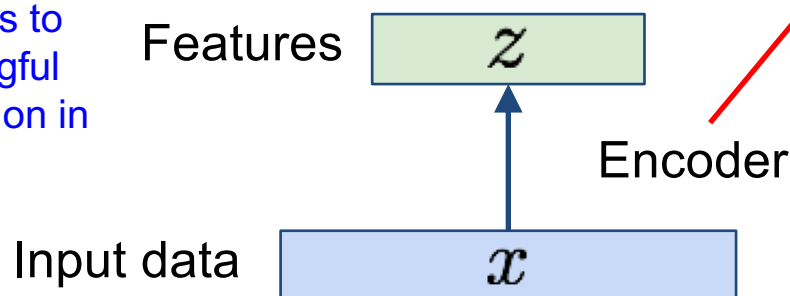
Q: Why dimensionality reduction?

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

Originally: Linear + nonlinearity (sigmoid)

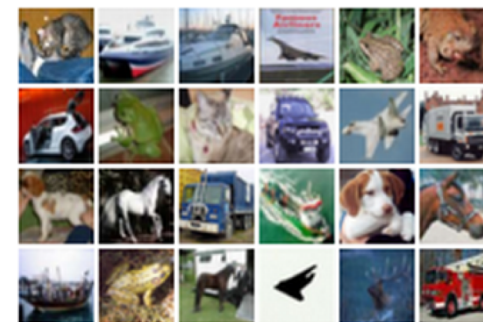
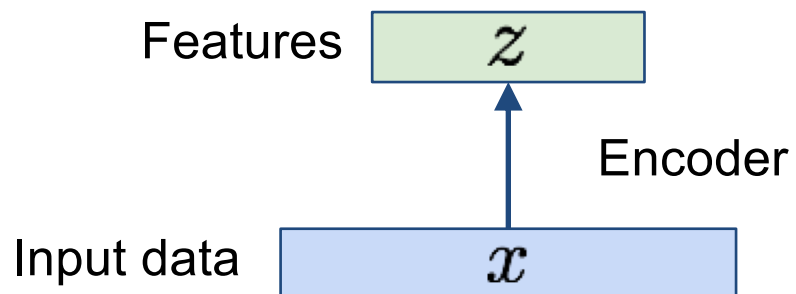
Later: Deep, fully-connected

Later: ReLU CNN



Autoencoders

How to learn this feature representation?

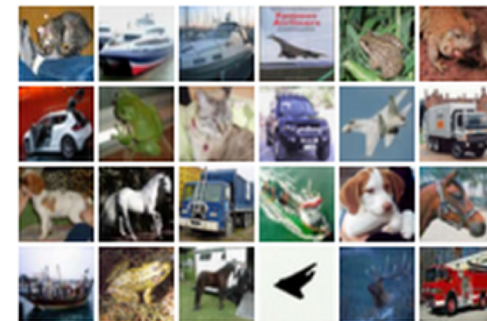
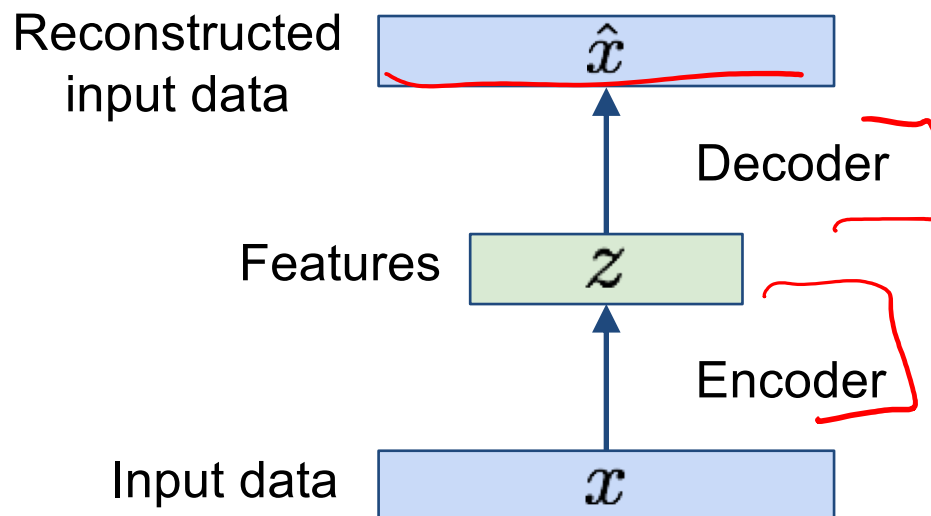


Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself

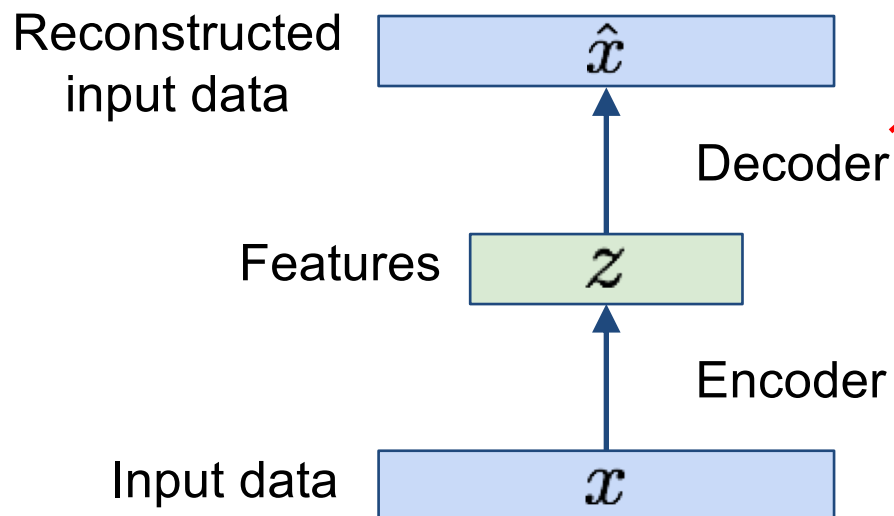


Autoencoders

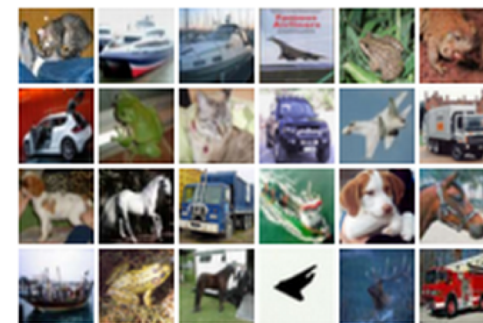
How to learn this feature representation?

Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself



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

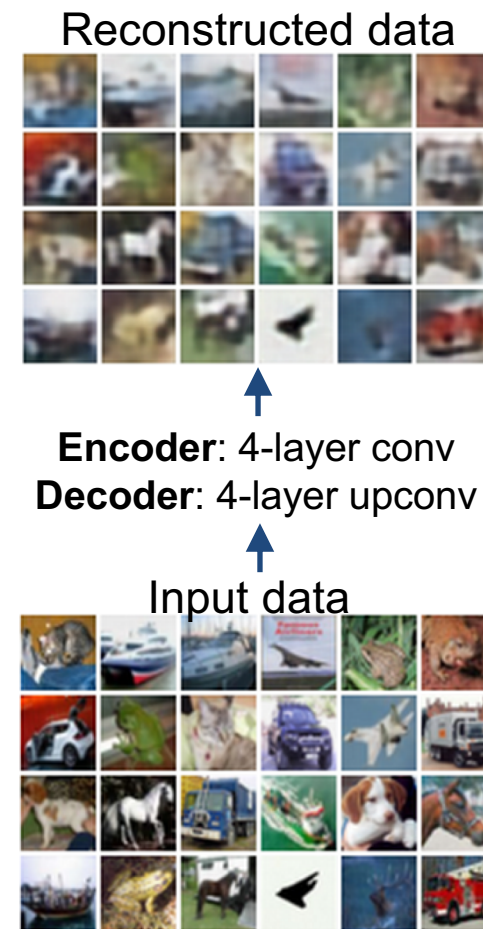
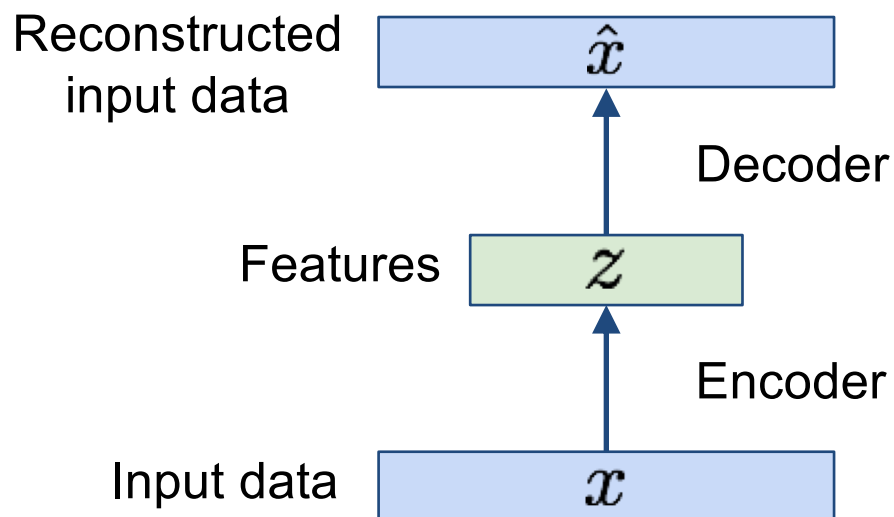


Autoencoders

How to learn this feature representation?

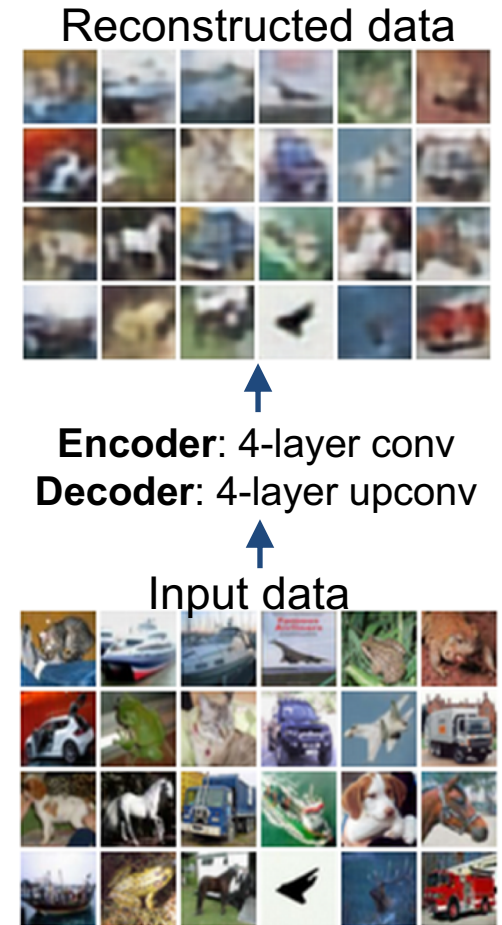
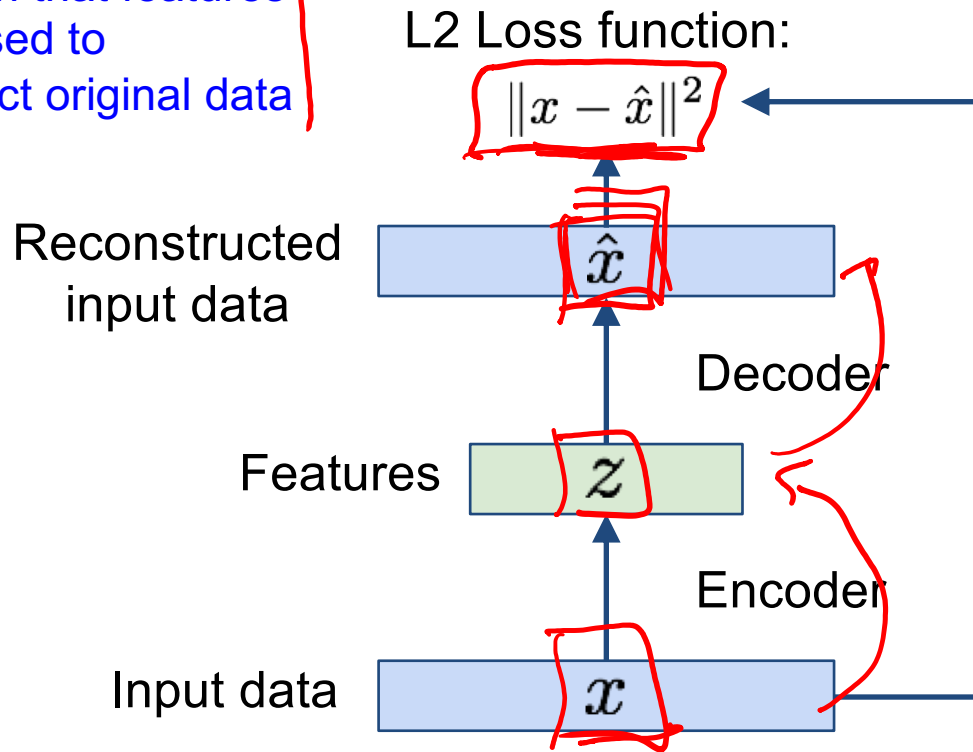
Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself



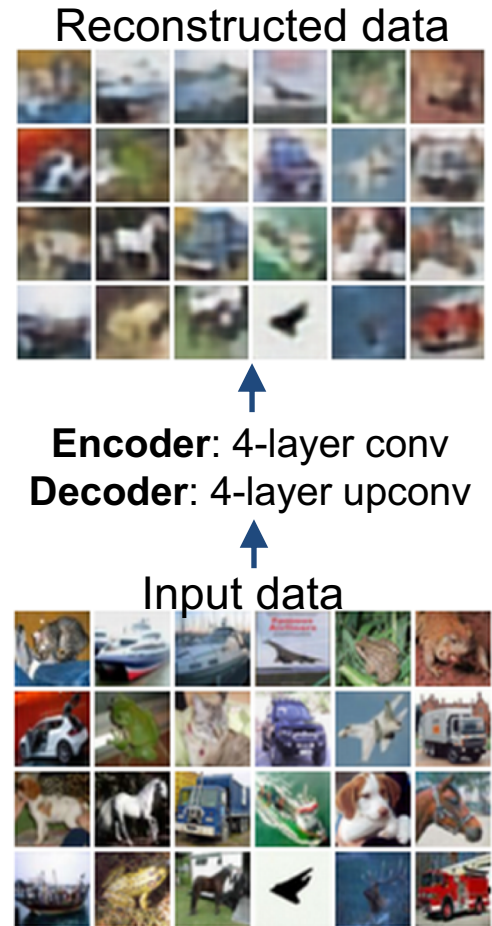
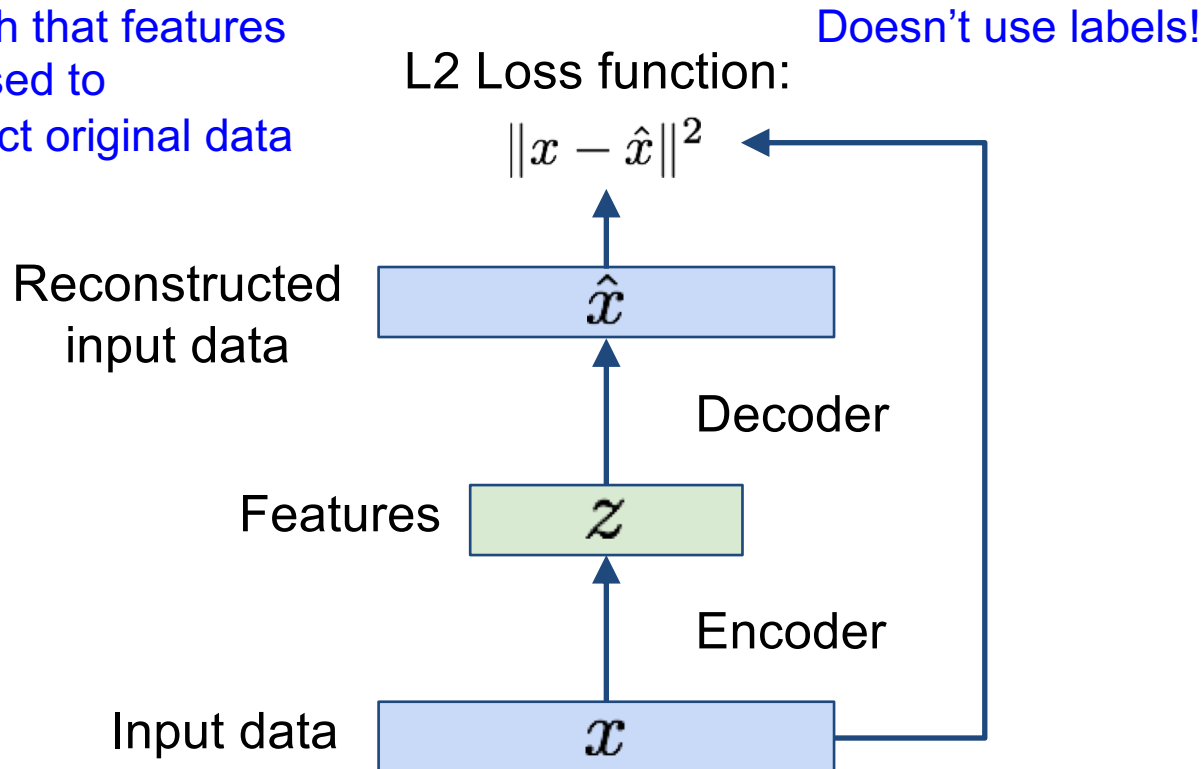
Autoencoders

Train such that features can be used to reconstruct original data



Autoencoders

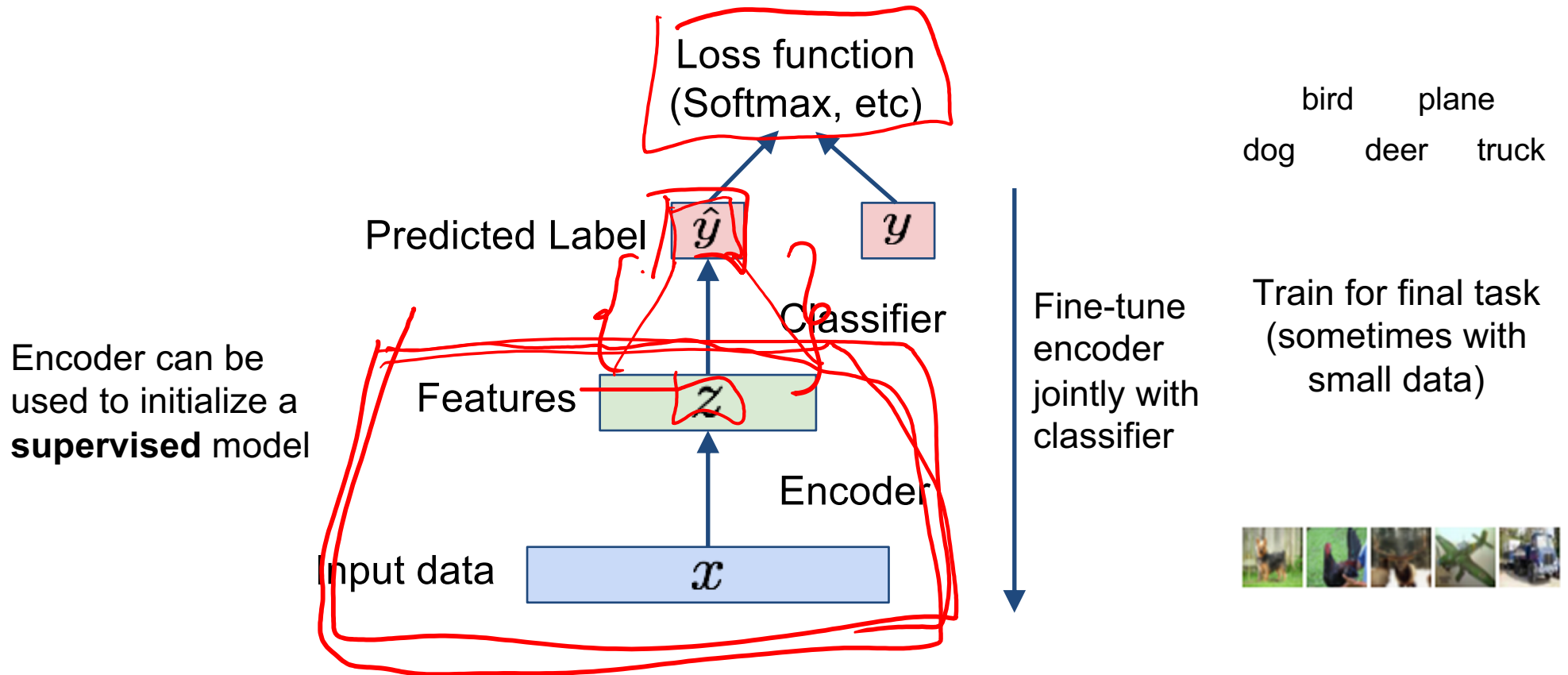
Train such that features can be used to reconstruct original data



Autoencoders

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

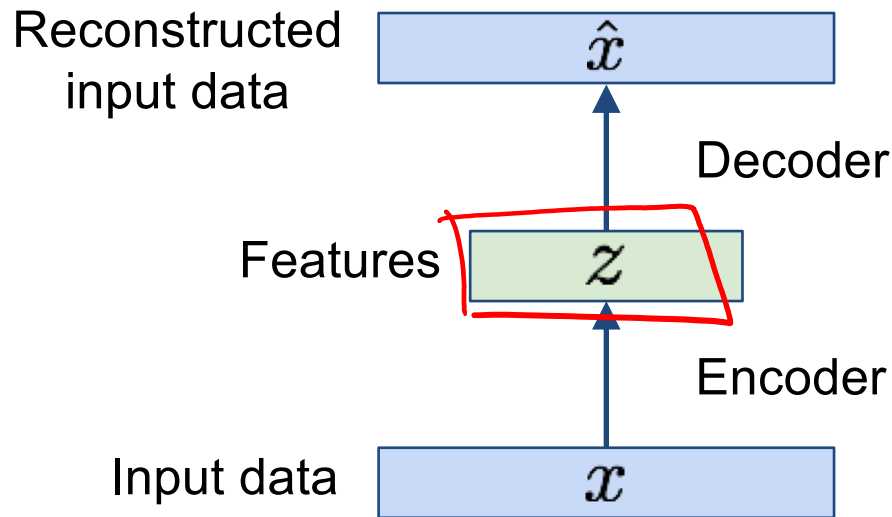
Autoencoders



Autoencoders

$$z = f_{\phi}(x)$$
$$\hat{x} = g_{\theta}(z)$$
$$p(z)$$
$$p(\hat{x}|z)$$

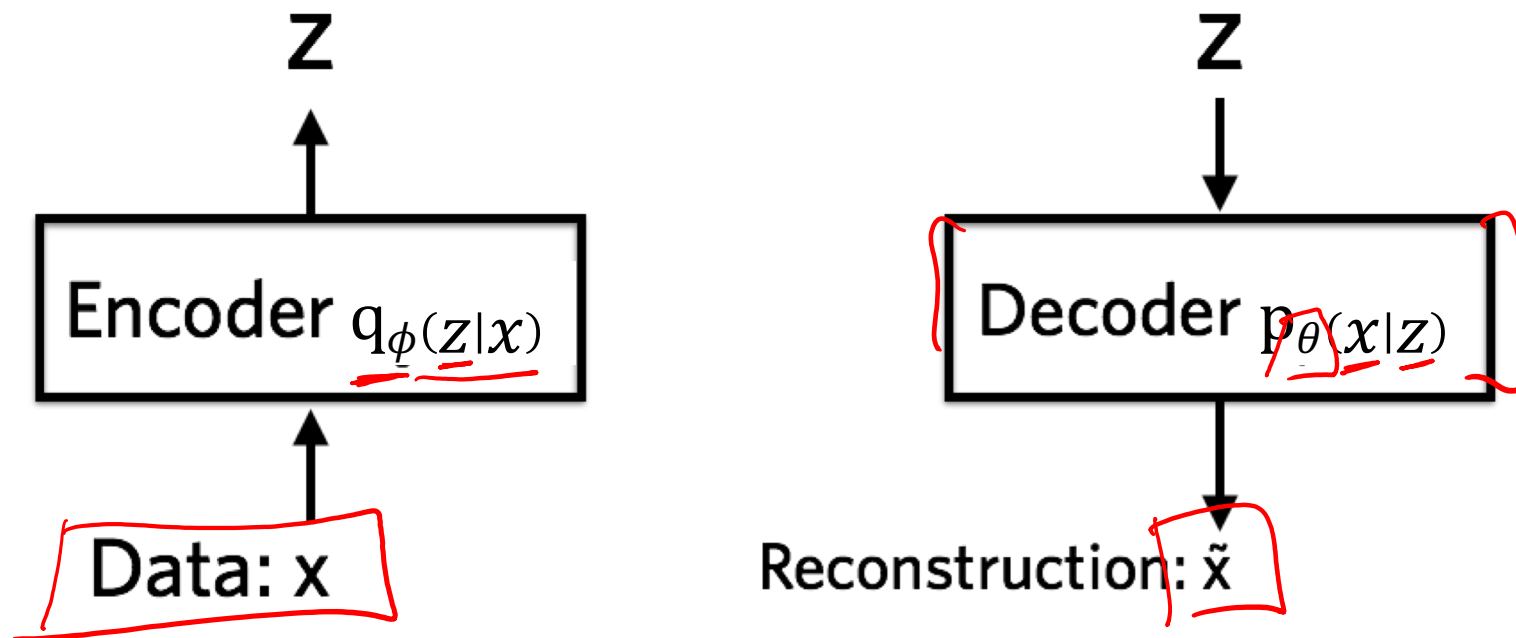
Autoencoders can reconstruct data, and can learn features to initialize a supervised model



Features capture factors of variation in training data. Can we generate new images from an autoencoder?

Variational Autoencoders

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



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

$$\bullet \frac{P(z|x)}{P(z)} = \frac{P(z,x)}{P(z)} = \frac{P(x|z)P(z)}{\sum_z P(x|z)P(z)}$$

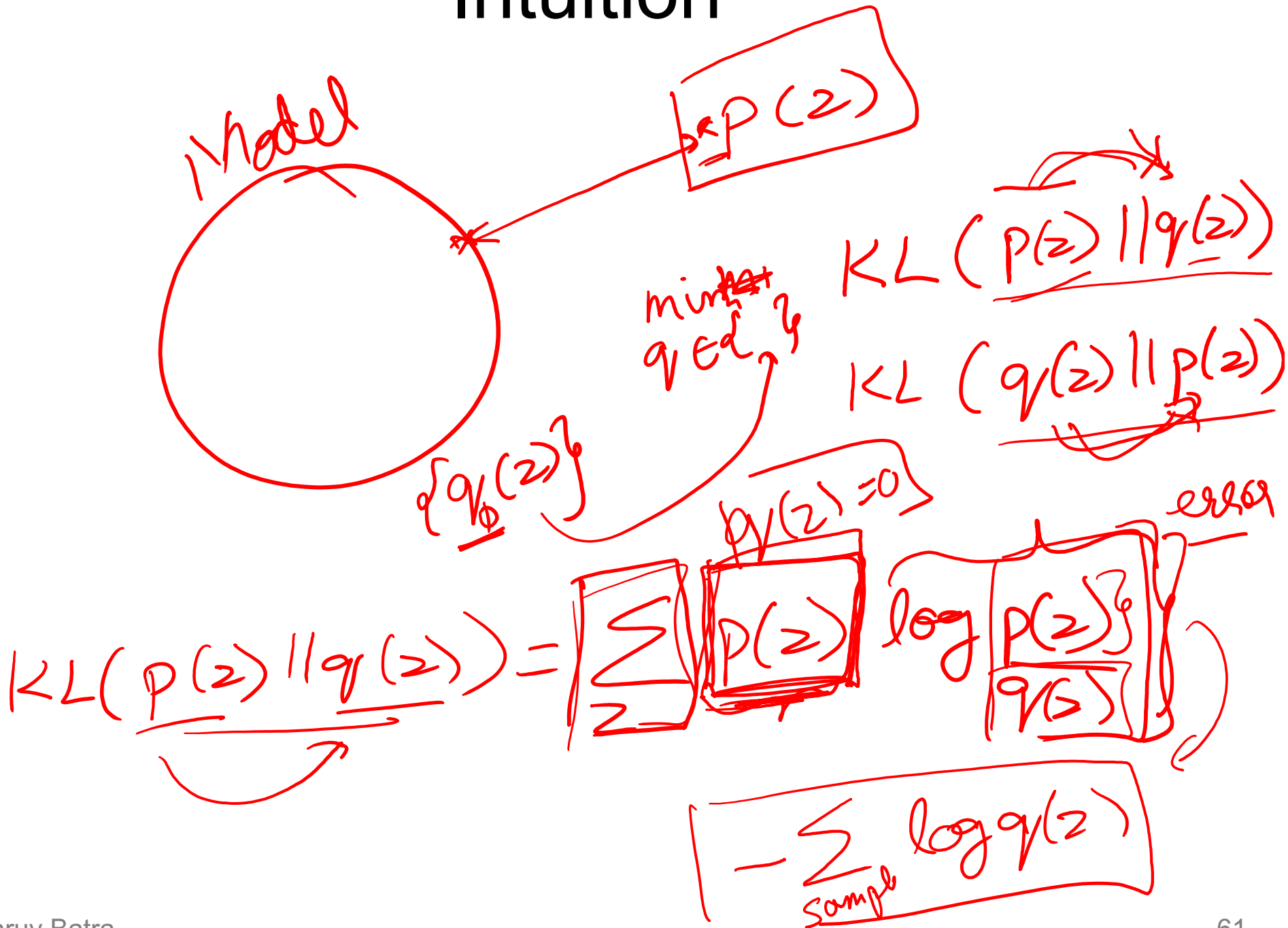
\downarrow

$q(z)$

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



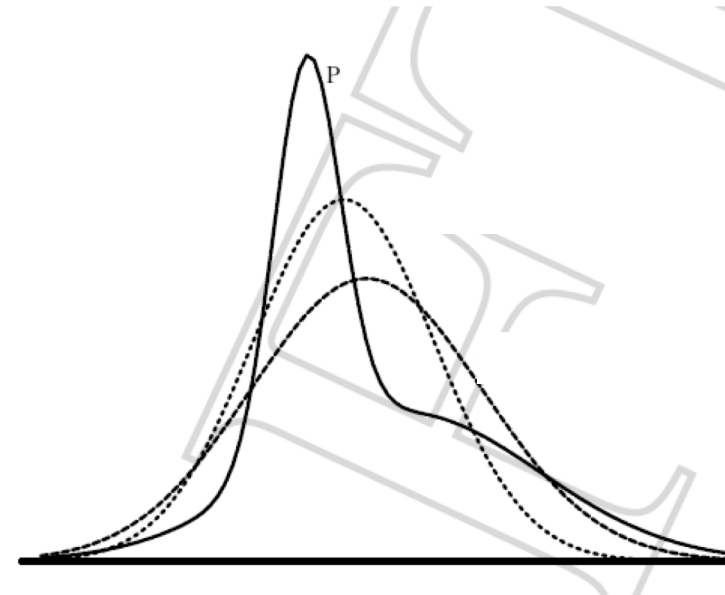
KL divergence:

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

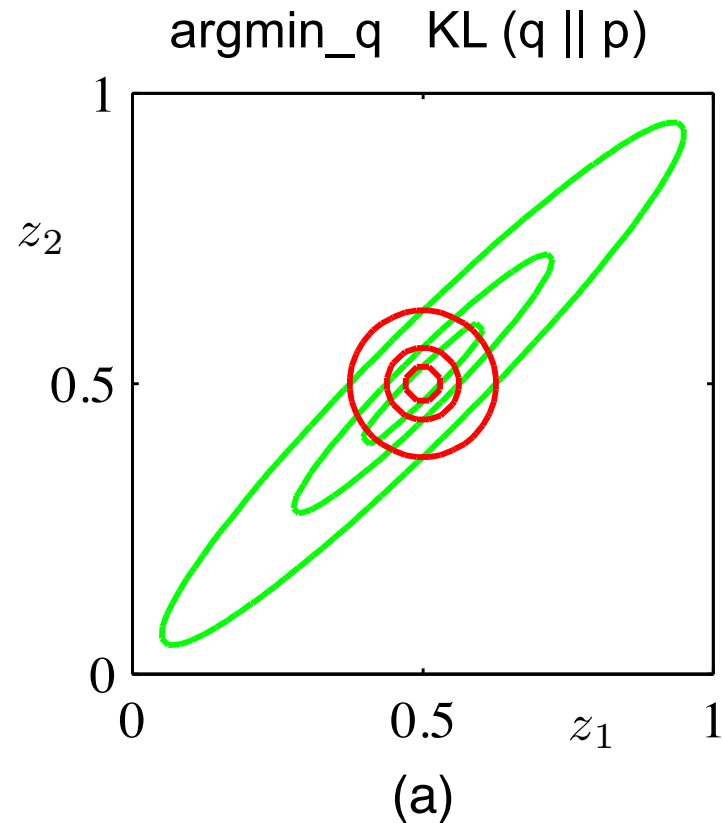
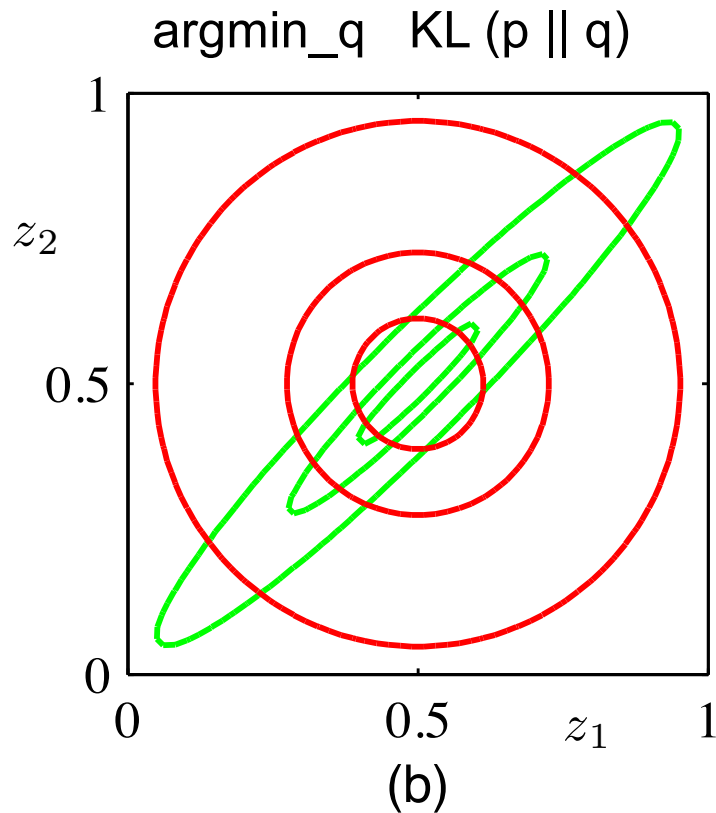
Find simple approximate distribution

- 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

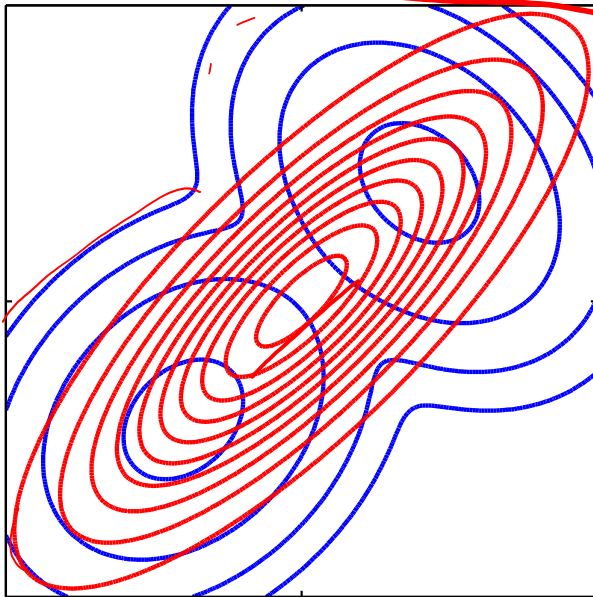
- p = 2D Gaussian with arbitrary co-variance
- q = 2D Gaussian with diagonal co-variance



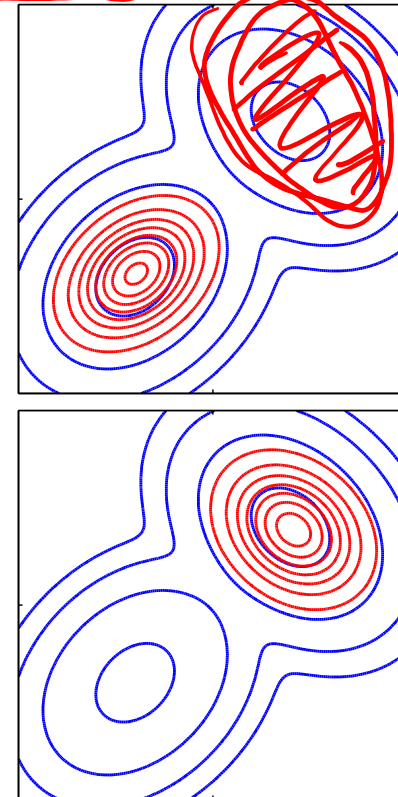
Example 2

- $p = \text{Mixture of Two Gaussians}$
- $q = \text{Single Gaussian}$

$\text{argmin}_q \text{KL}(p \parallel q)$



$\text{argmin}_q \text{KL}(q \parallel p)$



$$\sum \frac{p(z)}{q(z)} \left[\frac{q(z)}{p(z)} \right]$$