

CS 4803 / 7643: Deep Learning

Topics:

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

Dhruv Batra
Georgia Tech

Administrativa

- Project submission instructions
 - Due: 11/24, 11:59pm
 - Last deliverable in the class
 - Can't use late days
 - https://www.cc.gatech.edu/classes/AY2021/cs7643_fall/
- Aware of the discussions on Piazza

Recap from ~~last time~~ 2 lectures ago

Types of Learning

$$D = \{ (\bar{x}_i, \bar{y}_i) \}$$

$$x \xrightarrow{\Delta} y$$

- Supervised learning
 - Learning from a “teacher”
 - Training data includes desired outputs

$$s \rightarrow a \quad Q(s, a)$$

- Reinforcement learning
 - Learning to act under delayed evaluative feedback (rewards)

- Unsupervised learning
 - Discover structure in data
 - Training data does not include desired outputs

$$D = \{ \bar{x}_i \}$$

Supervised vs Reinforcement vs Unsupervised Learning

Unsupervised Learning

Training data is cheap

Data: x

Just data, no labels!

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

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

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

Supervised Learning

Data: (x, y)

x is data, y is label

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

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

Supervised vs Reinforcement vs Unsupervised Learning

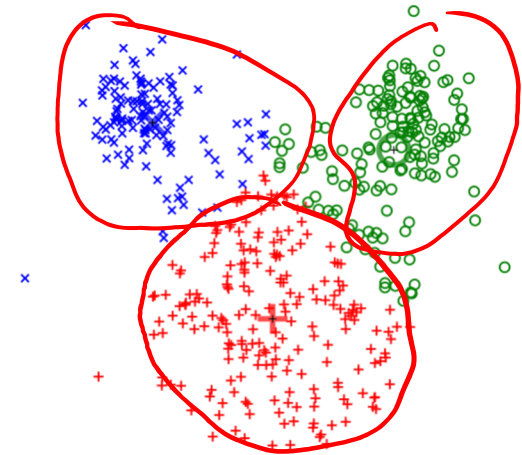
Unsupervised Learning

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

[This image is CC0 public domain](#)

Supervised vs Reinforcement vs Unsupervised Learning

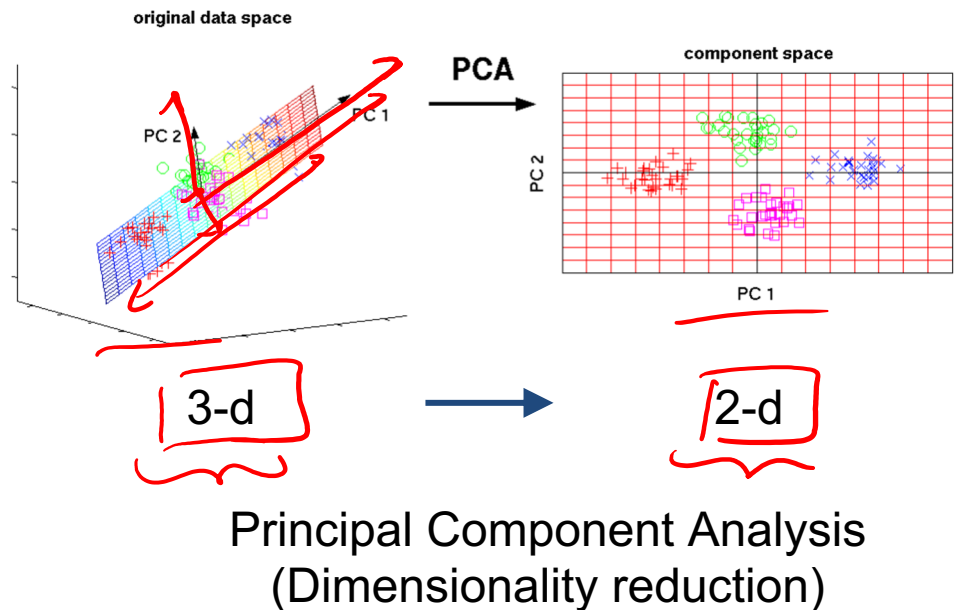
Unsupervised Learning

Data: x

Just data, no labels!

Goal: Learn some underlying hidden structure of the data

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



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

Supervised vs Reinforcement vs Unsupervised Learning

$p(x)$

$x \rightarrow p(x)$

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.

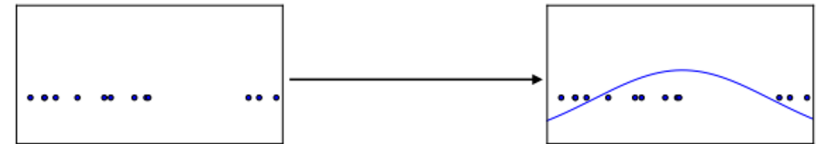
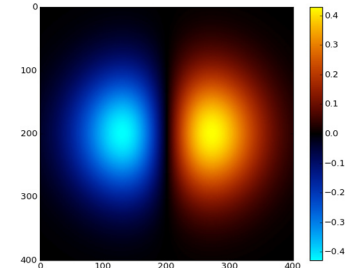
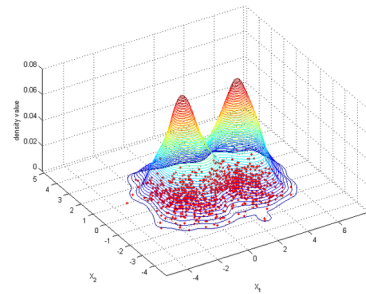


Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



2-d density estimation

$$\vec{x} = [x_1 \dots x_d]$$

$$\underline{P(x_i | x_j, \dots, x_{j+n})}$$

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

Generative Models

Given training data, generate new samples from same distribution



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



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

Goal:

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

Taxonomy of Generative Models

block box

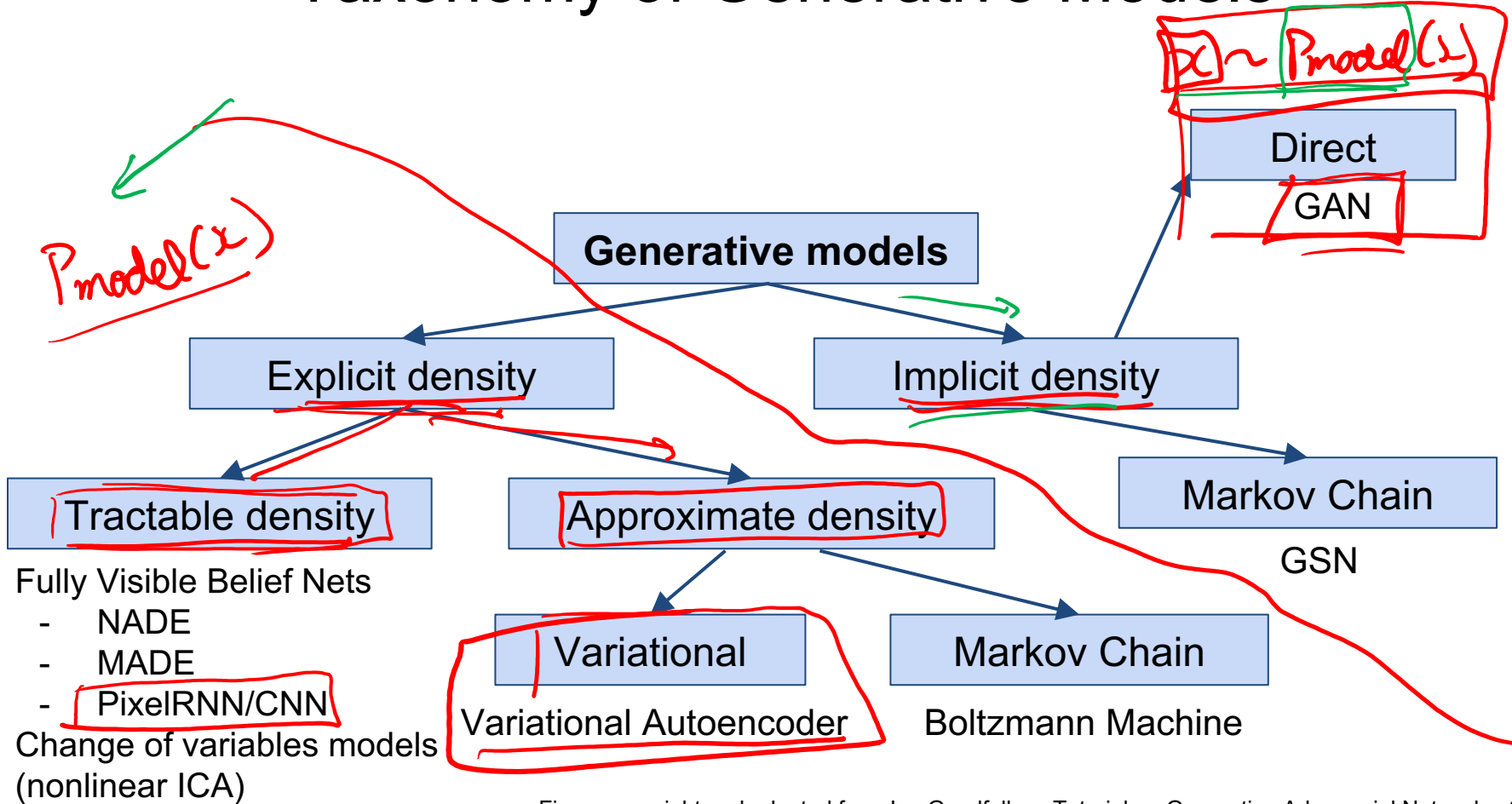


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Taxonomy of Generative Models

We will discuss 3 most popular types of generative models

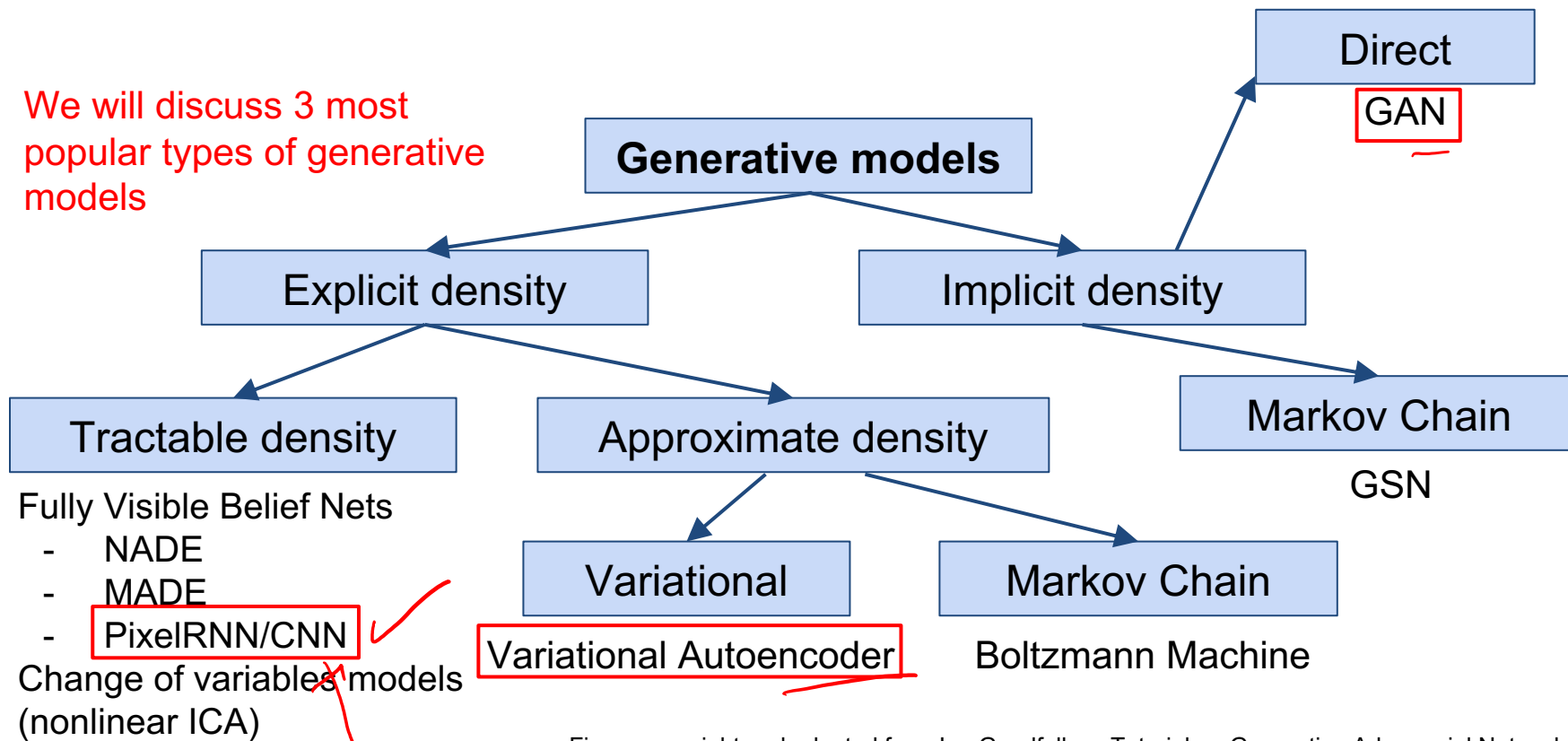


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Plan for Today

- Goal: Variational Autoencoders
- Latent variable probabilistic models
 - Example GMMs
- Autoencoders
- Variational Inference

Variational Autoencoders (VAE)

So far...

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

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

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

$N = \# \text{samples}$

"complex" $\rightarrow P(\vec{x})$

R.V



$$\rightarrow P(\vec{x}, \vec{z})$$

"Latent Variables"

Unobserved RV

$$= \underbrace{P(\vec{x} | \vec{z})}_{\text{conditional}} \underbrace{P(\vec{z})}_{\text{Prior}}$$



"simpler"

So far...

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

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

$$D = \{ \vec{x}_i \} \leftarrow$$

$$P(\vec{x}, \vec{z})$$

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

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

if z is continuous

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

if z is discrete

So far...

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

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

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

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

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

GMM

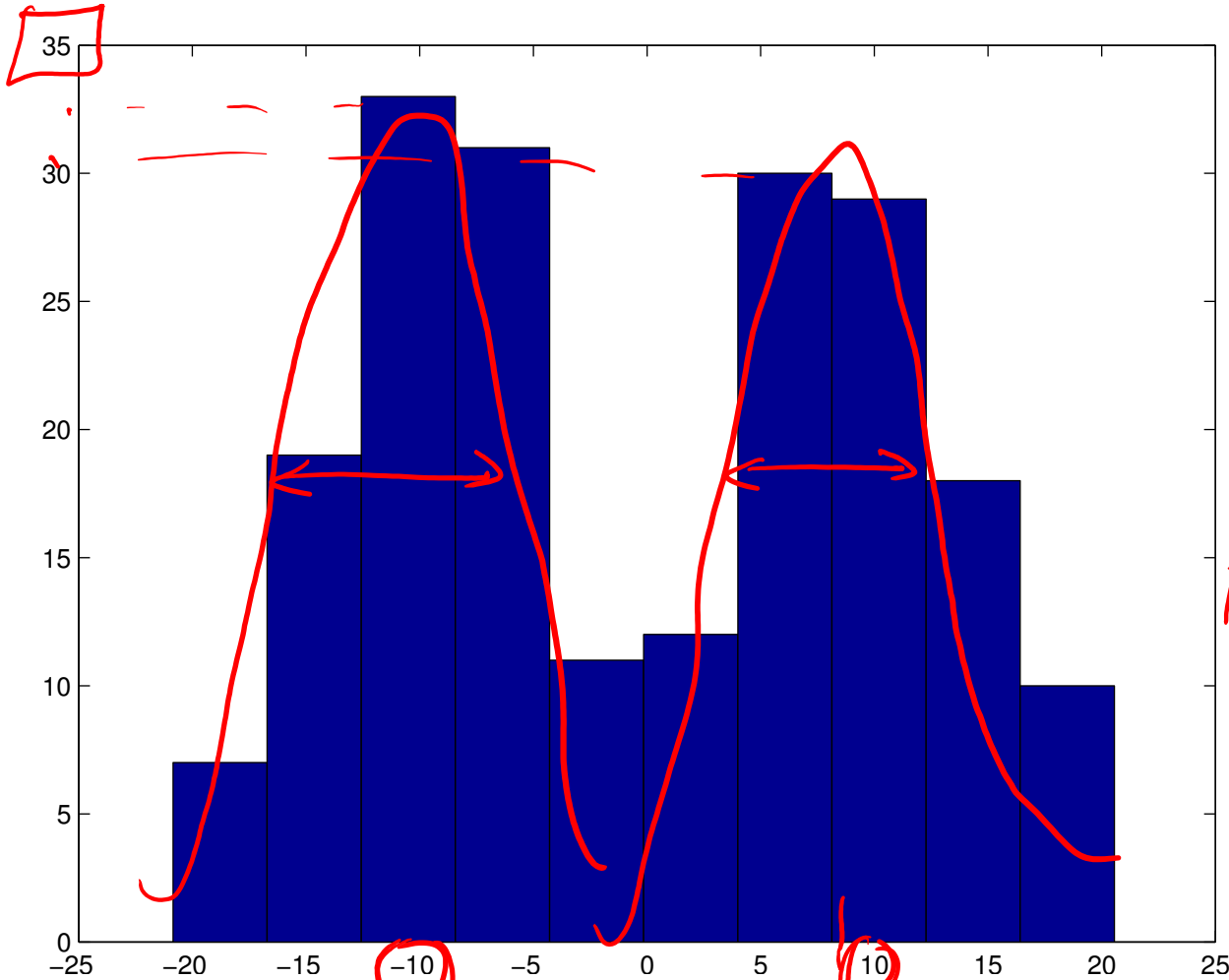
$x \in \mathbb{R}^1$

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

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

$$P(z) = \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix}$$

$\sim P(x)$



$p(x|z) \sim \mathcal{N}(\mu_1, \sigma_1^2)$ $\mu_1 = -10$

$\mathcal{N}(\mu_2, \sigma_2^2)$ $\mu_2 = +10$

Gaussian Mixture Model $z \in \{1, \dots, k\}$

$$P(x, z)$$

② Latent

$$z \sim \text{Cat}(\underline{\pi})$$

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

$$\pi_c = P(z=c)$$

↓

① Observed

$$\begin{aligned} P(x|z=c) &= \mathcal{N}(\mu_c, \sigma_c^2) \\ &= \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(x-\mu_c)^2}{2\sigma_c^2}} \end{aligned}$$

$$P(\underline{x}, z) = P(x|z) P(z)$$

Gaussian Mixture Model

$$P(Z=c) = \pi_c$$

$$P(x|z) = N(\quad)$$

Available
from model

$$P(\vec{x}) = \sum_z P(x, z)$$

$$= \sum_z \overbrace{P(x|z)} \overbrace{P(z)} \quad \equiv \text{Marginalization}$$

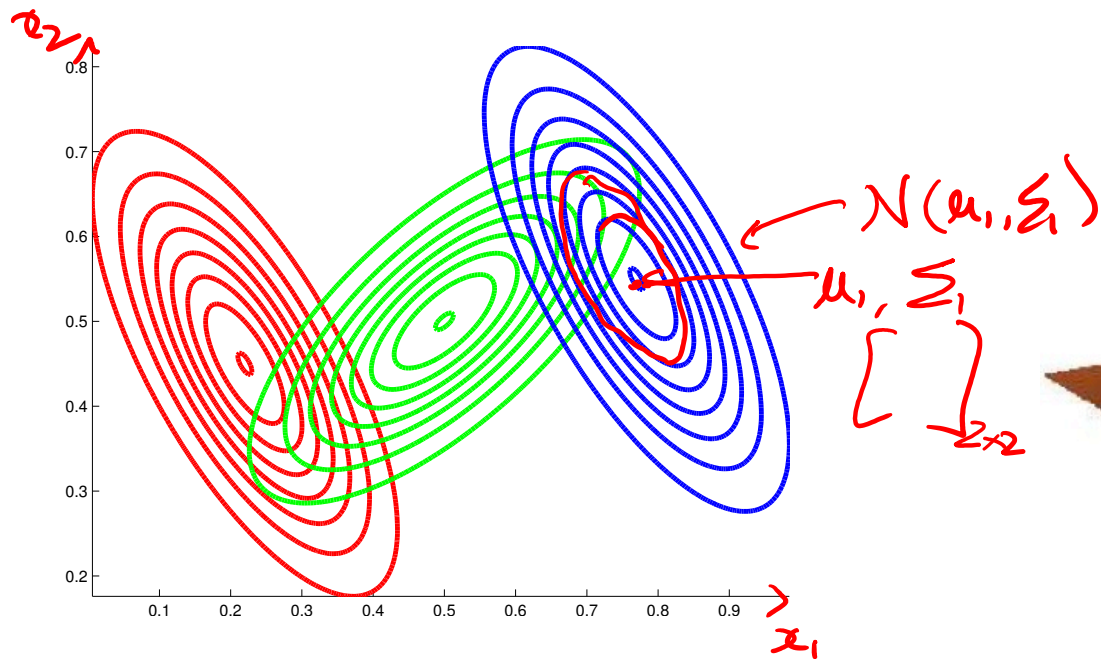
$$\boxed{P(z|x)} = \frac{P(z, x)}{P(x)} = \frac{P(x|z) p(z)}{\sum_z (\downarrow) (\downarrow)}$$

\equiv (Inference)

$$\begin{aligned} \mathbf{x} &\in \mathbb{R}^d \\ \bar{\boldsymbol{\mu}} &\in \mathbb{R}^d \\ \boldsymbol{\Sigma} &\in \mathbb{R}^{d \times d} \end{aligned}$$

GMM

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



K-means vs GMM

- K-Means
 - <http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>
- GMM
 - <https://lukapopijac.github.io/gaussian-mixture-model/>

$P(x, z)$

Hidden Data Causes Problems #1

- Fully Observed (Log) Likelihood factorizes
- Marginal (Log) Likelihood doesn't factorize
- All parameters coupled!

Parameters: $\{\pi_1, \dots, \pi_k, \vec{\mu}_1, \dots, \vec{\mu}_k, \Sigma_1, \dots, \Sigma_k\} \equiv \theta$

$$D = \{\vec{x}_i\}_{i=1}^N \quad \vec{x}_i \in \mathbb{R}^d$$

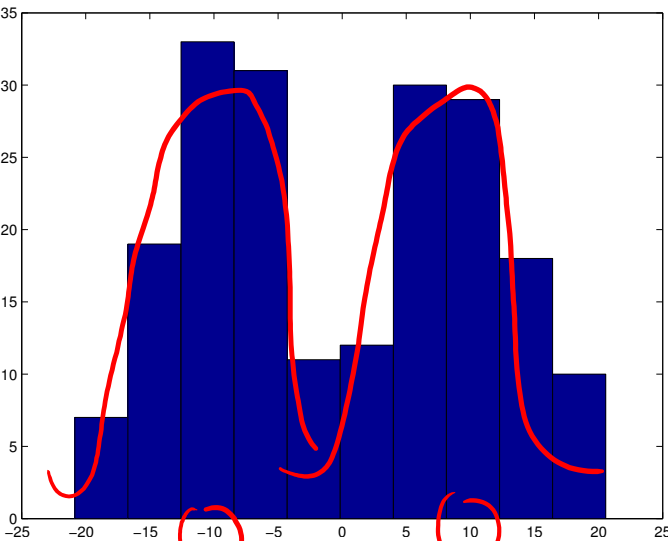
$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \log P(D | \theta)$$
$$= \sum_{i=1}^N \log P(\vec{x}_i | \theta)$$

$$\sum_{i=1}^N \log \left[\sum_{z_i} P(x_i, z_i | \theta) \right]$$

\uparrow
[S]

Hidden Data Causes Problems #2

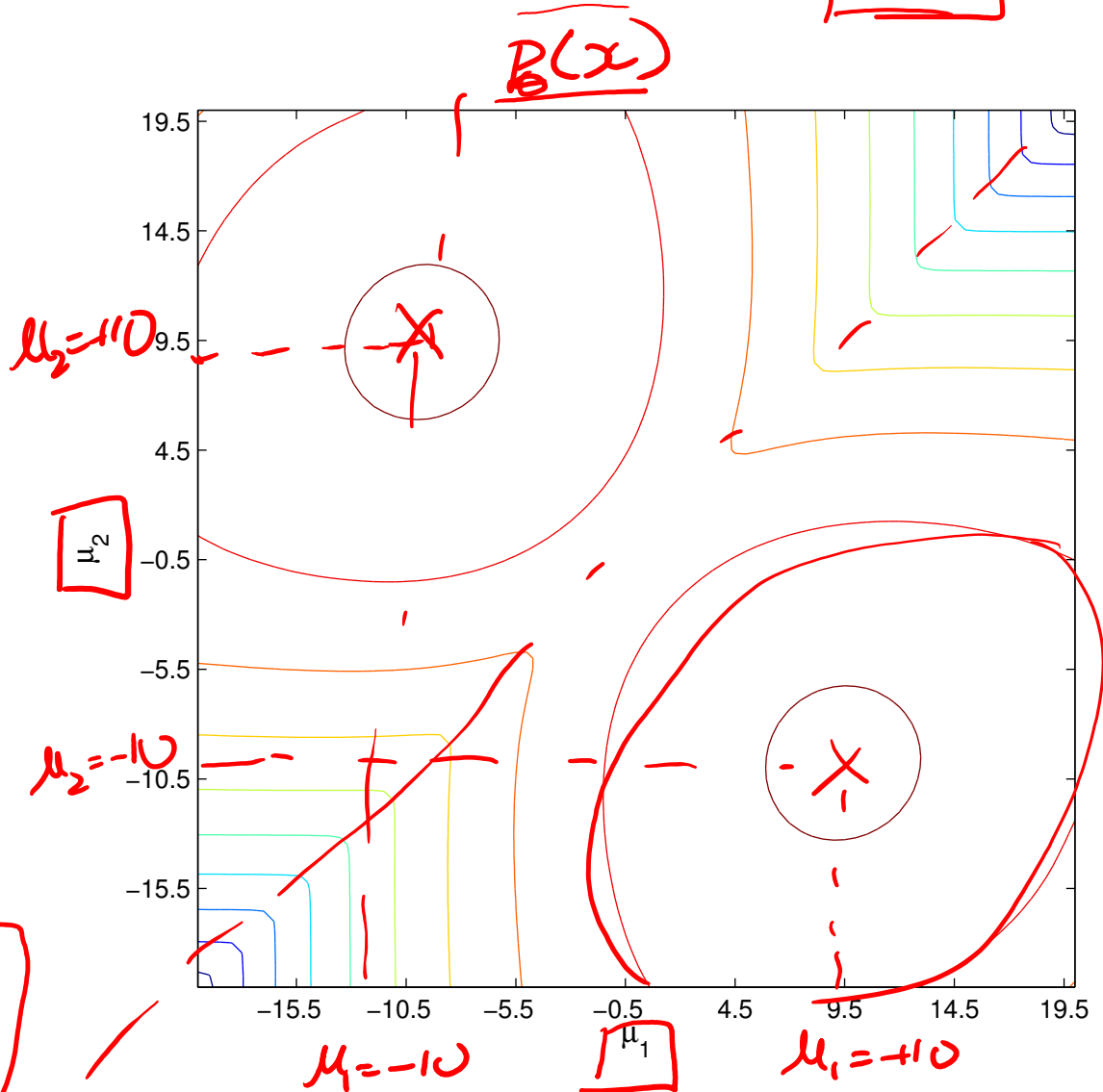
- Identifiability



μ_1

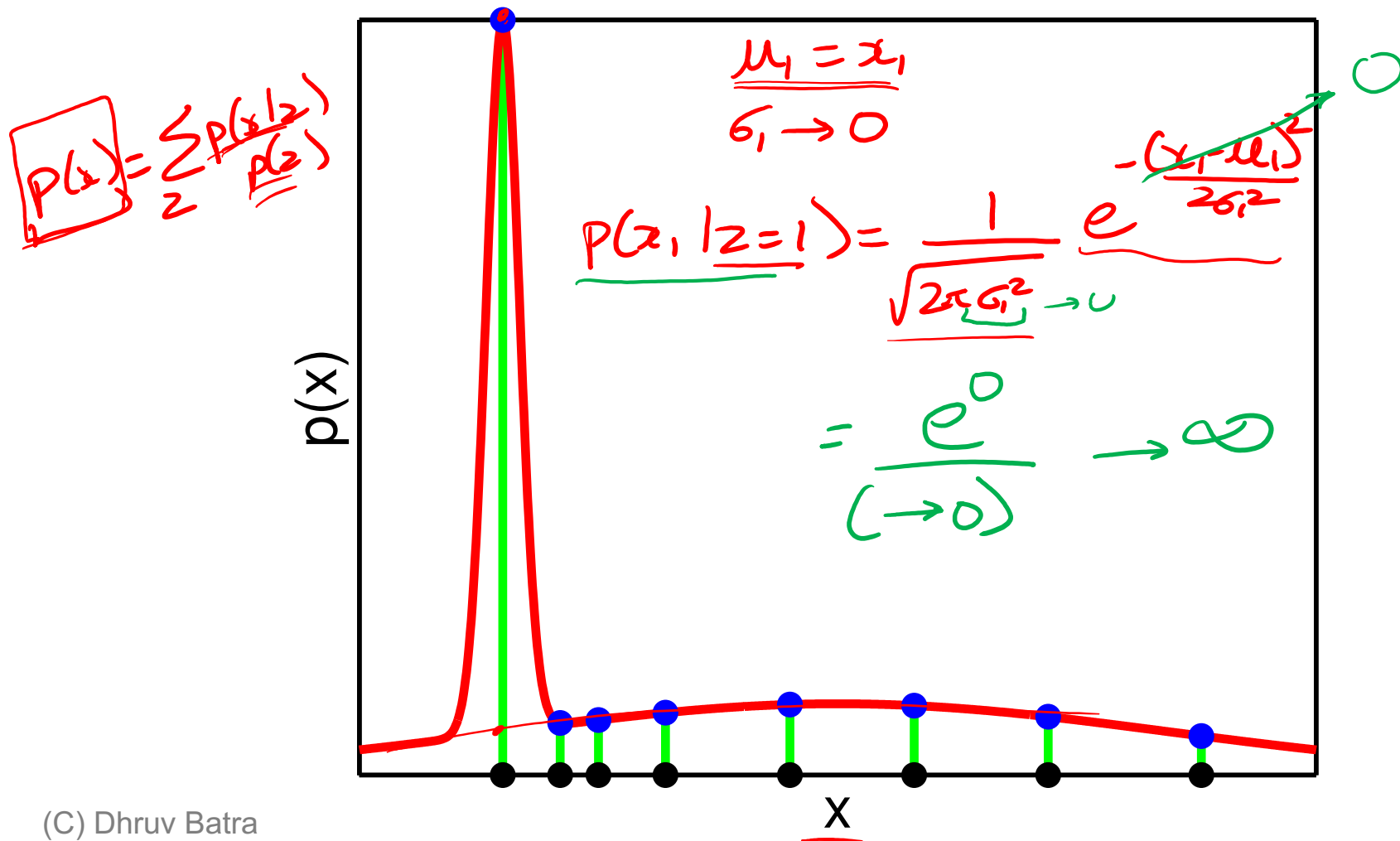
μ_2

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



Hidden Data Causes Problems #3

- Likelihood has singularities if one Gaussian “collapses”



Variational Auto Encoders

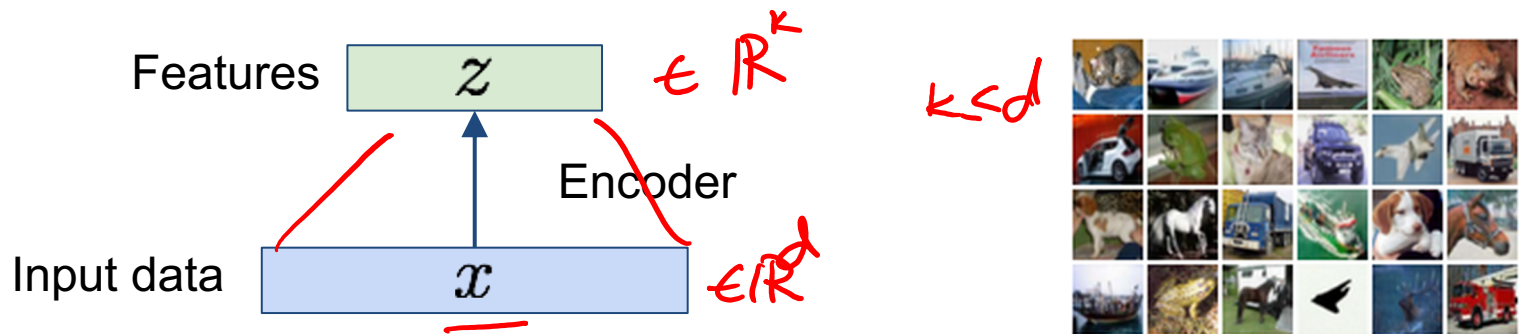
$$\begin{aligned} & \mathcal{L}_V \\ & P(x, z) \\ & \frac{P(z|x)}{P(x|z)} \end{aligned}$$

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 ^c

Autoencoders

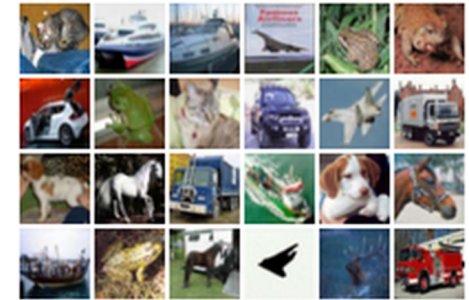
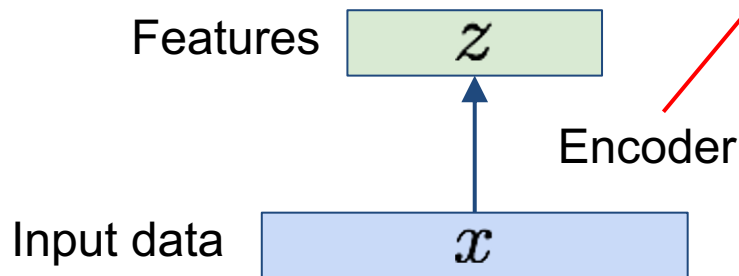
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data



Autoencoders

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

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

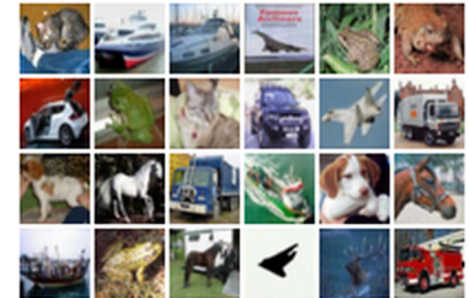
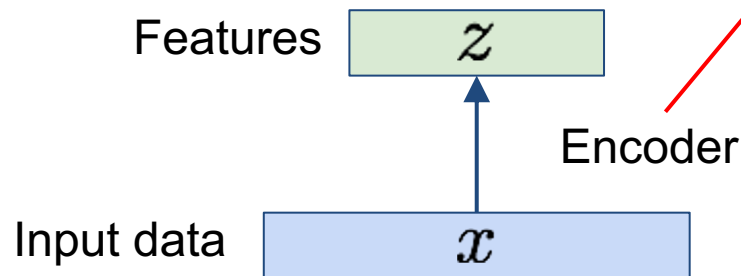
z usually smaller than x
(dimensionality reduction)

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

z usually smaller than x
(dimensionality reduction)

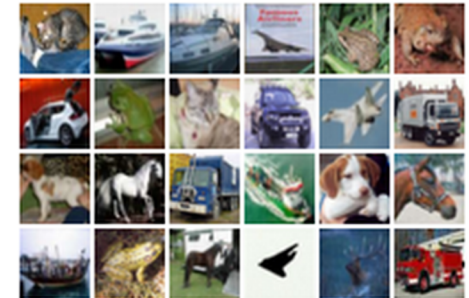
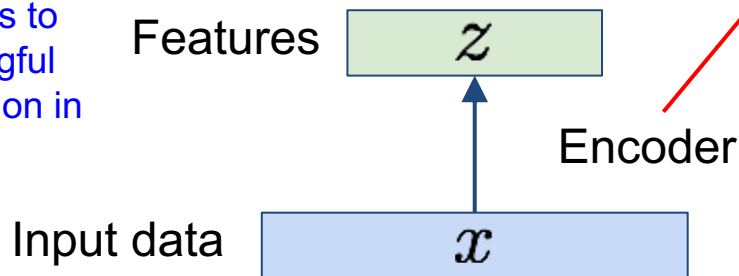
Q: Why dimensionality reduction?

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

Originally: Linear + nonlinearity (sigmoid)

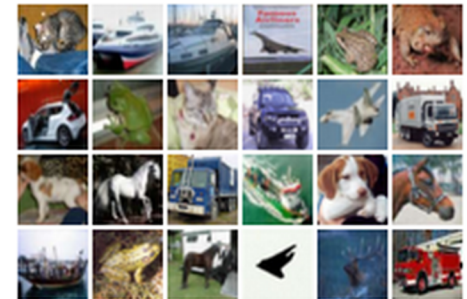
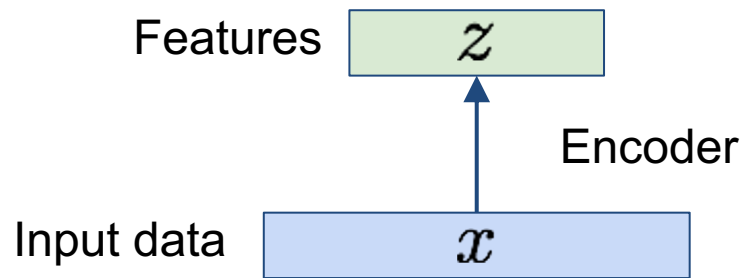
Later: Deep, fully-connected

Later: ReLU CNN



Autoencoders

How to learn this feature representation?

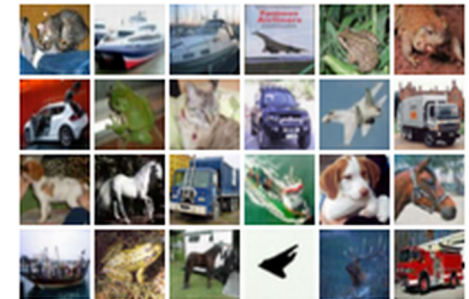
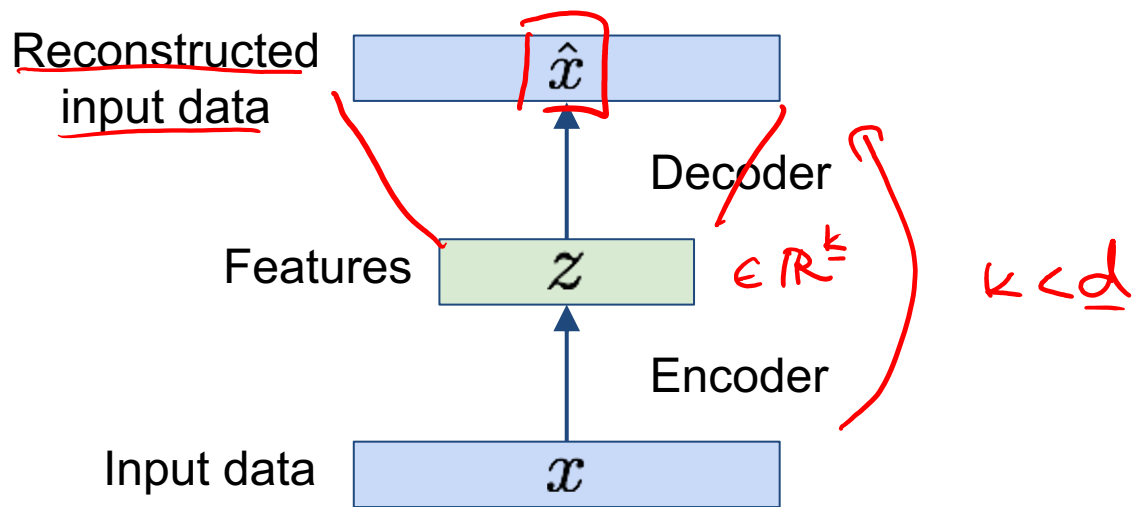


Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data

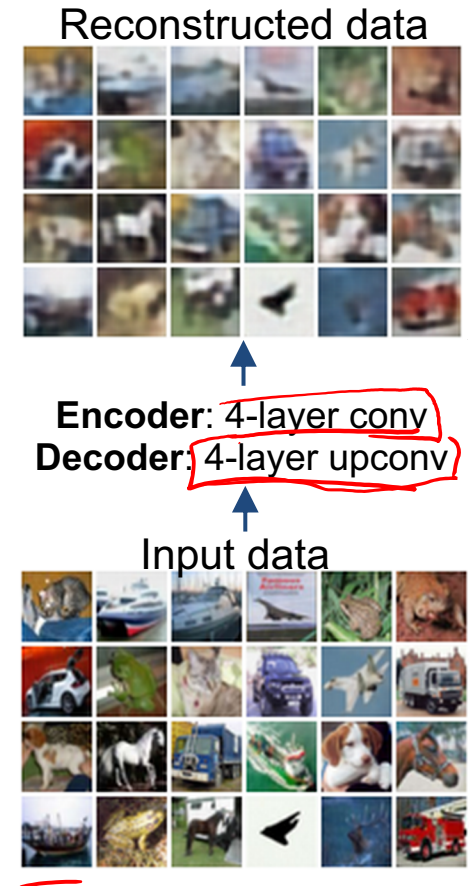
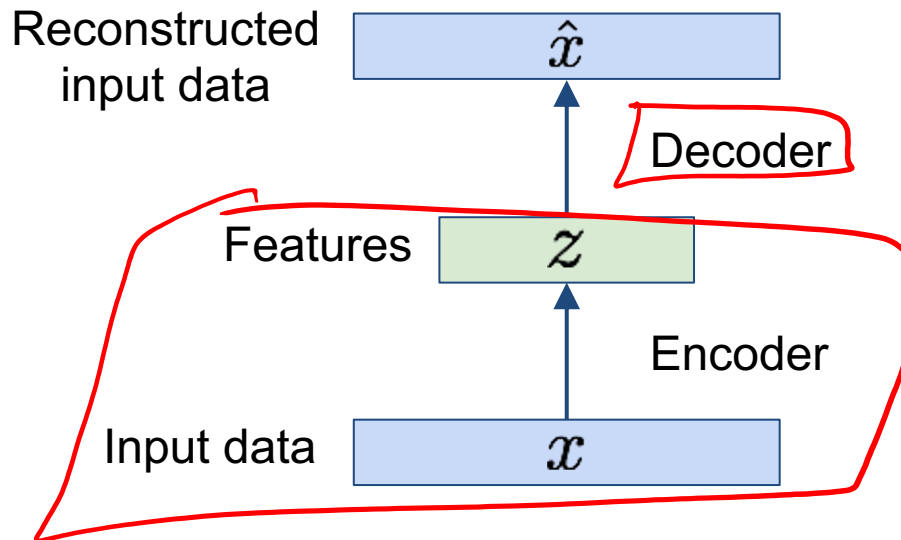
“Autoencoding” - encoding itself



Autoencoders

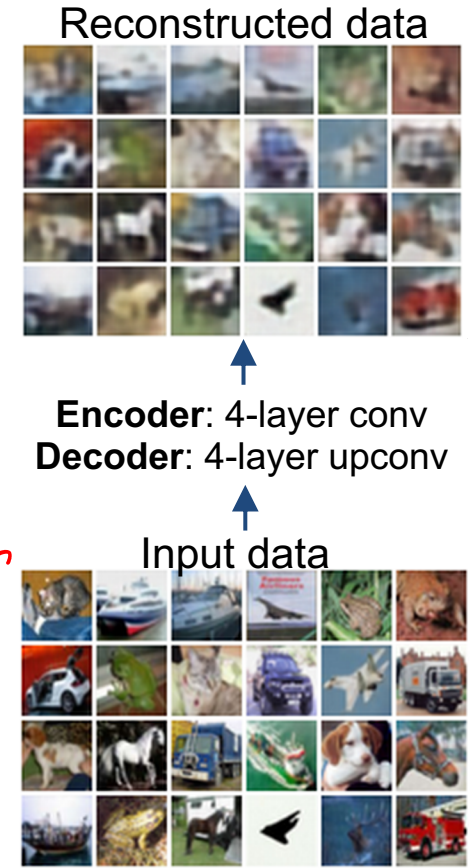
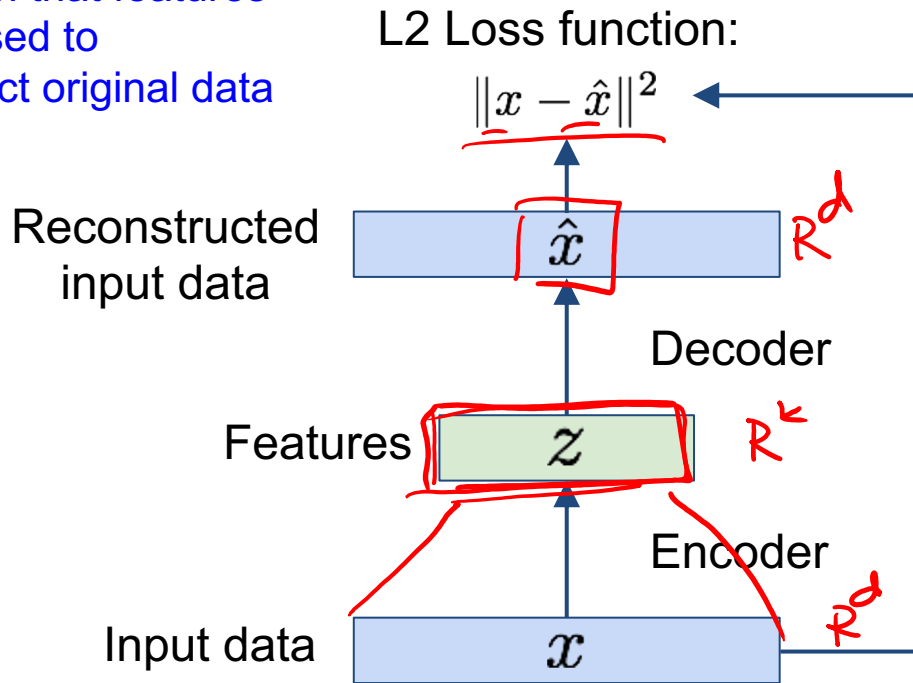
How to learn this feature representation?

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



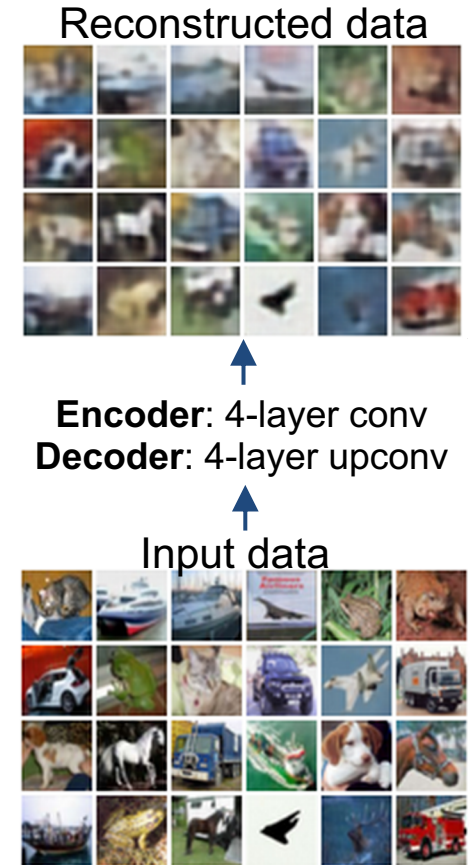
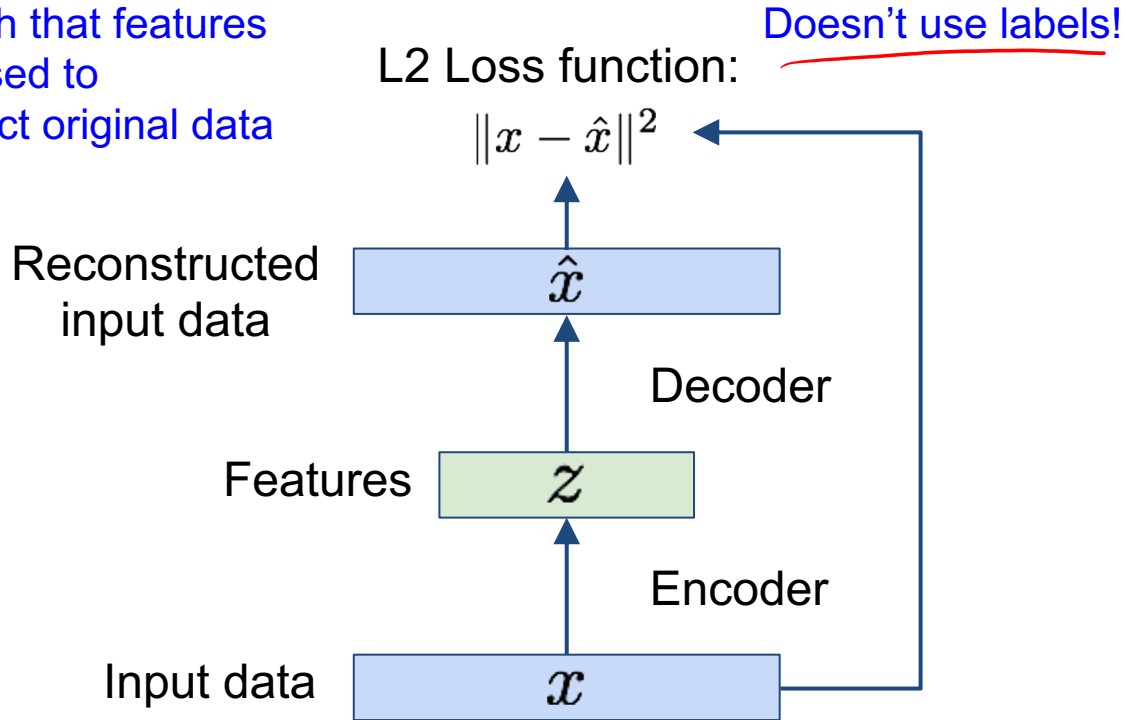
Autoencoders

Train such that features can be used to reconstruct original data



Autoencoders

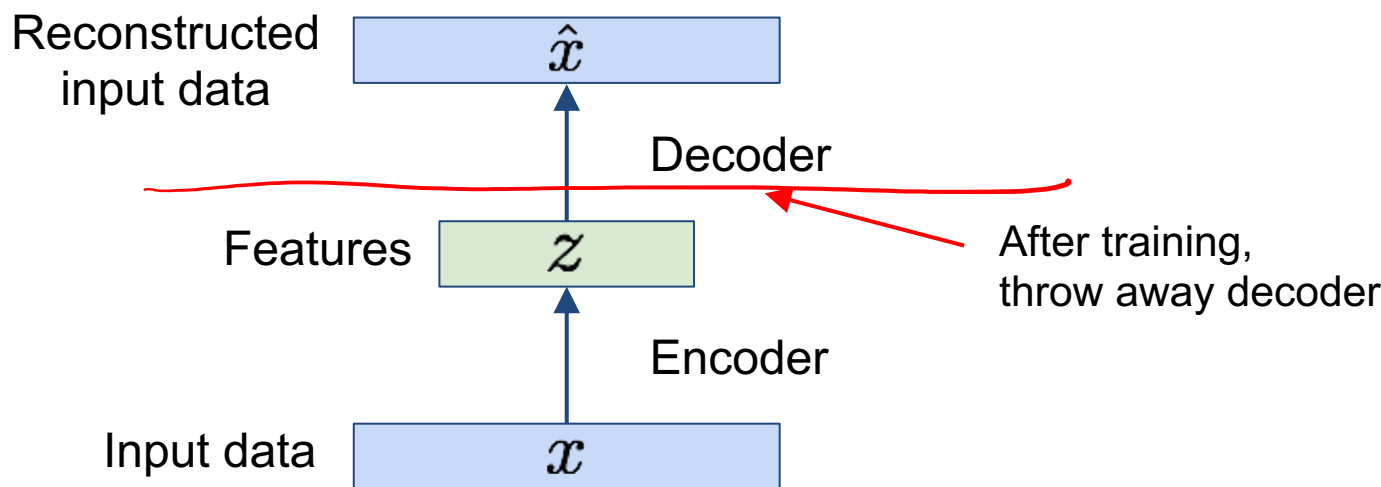
Train such that features can be used to reconstruct original data



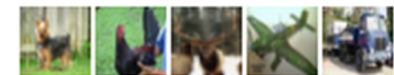
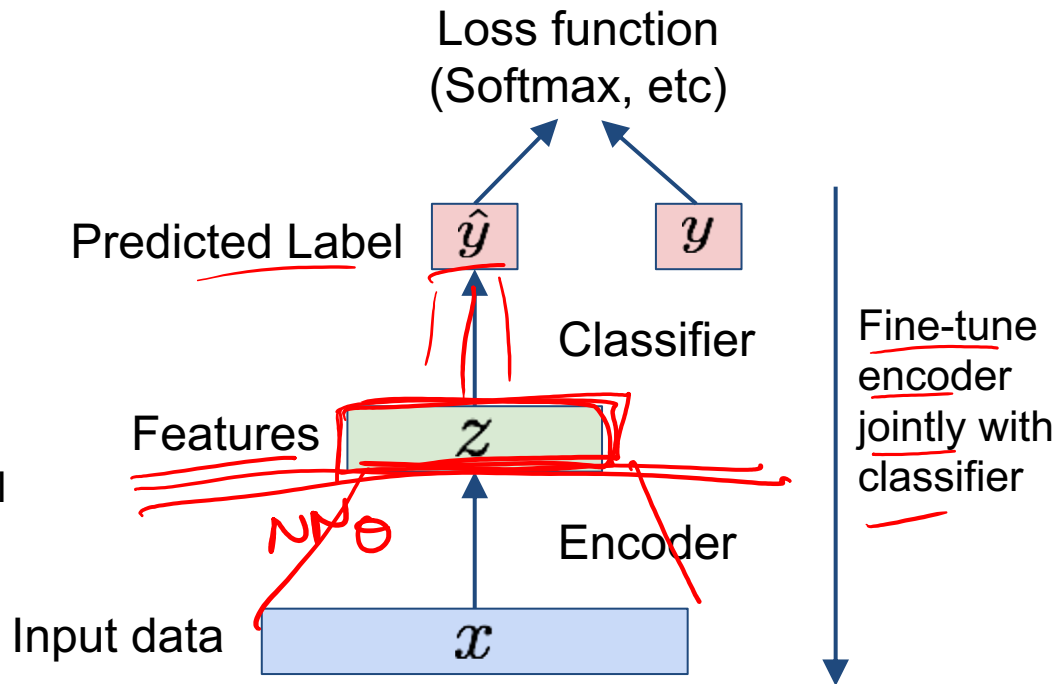
Autoencoders

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

Autoencoders



Autoencoders



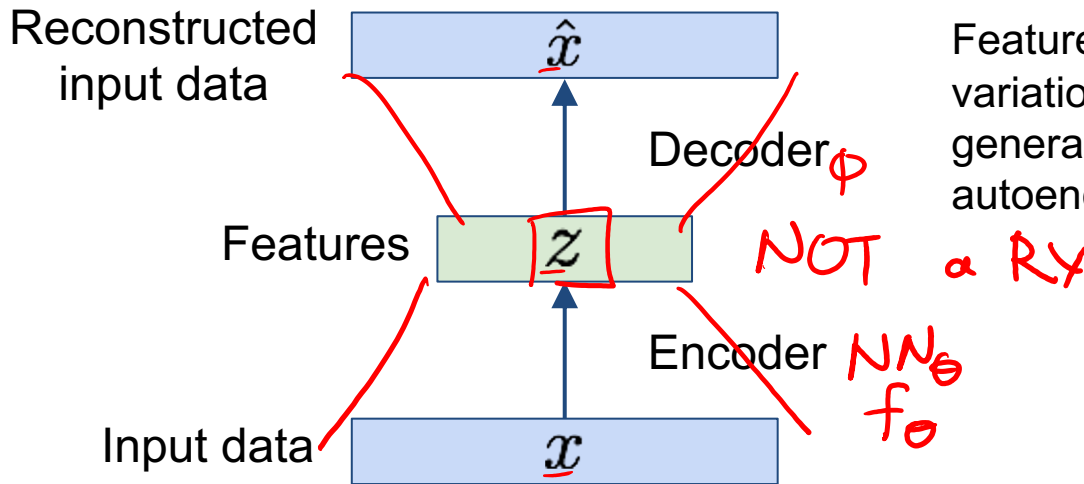
Autoencoders

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

$$\left[\begin{array}{c} p(z|x) \\ \sim p(\hat{x}|z) \end{array} \right]$$

VAE

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

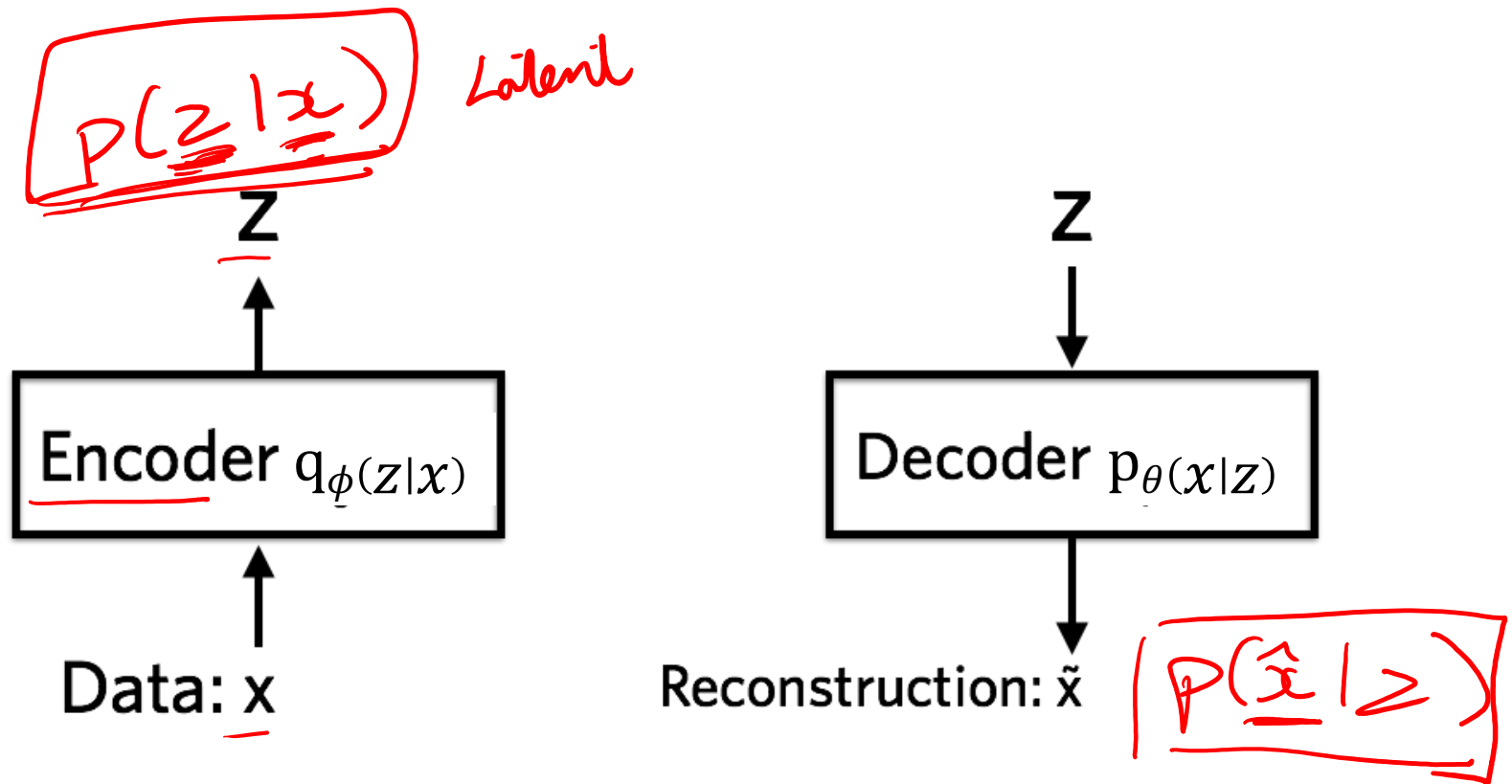


Features capture factors of variation in training data. Can we generate new images from an autoencoder?

Variational Autoencoders

z

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

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

"complex" $P(z)$



"simple" $q_\phi(z)$



$*P(z)$

$$\min_{q \in \mathcal{Q}} \underbrace{d(p, q)}_{\text{KL}(\cdot)}$$

"simple" distributions $\{q_\phi(z)\}$

$\text{KL}(p||q)$
"right"
"hard"

$\text{KL}(q||p)$
"wrong"
"easy"

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.

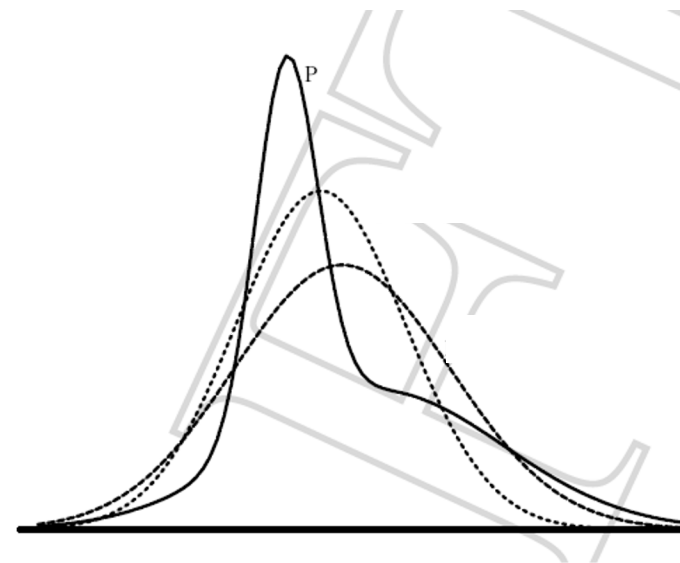
Intuition

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

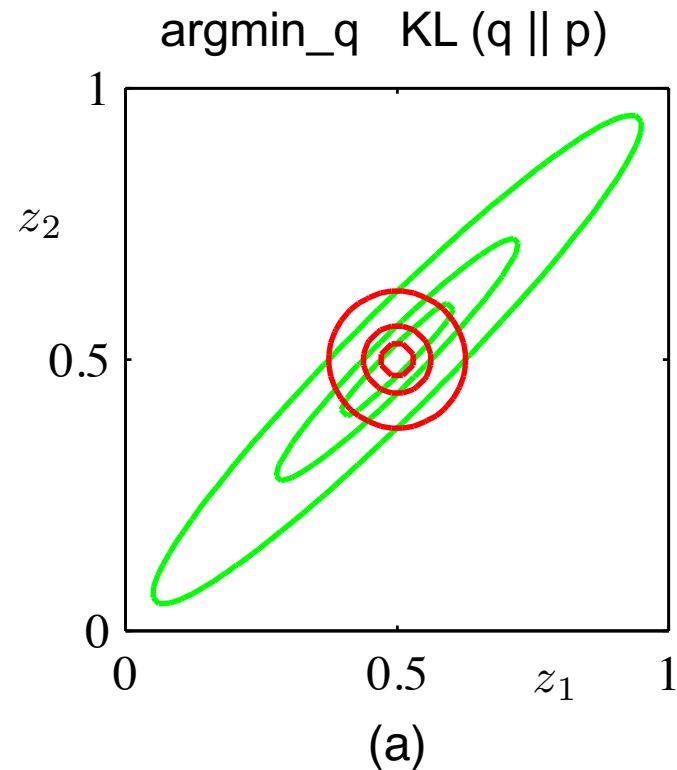
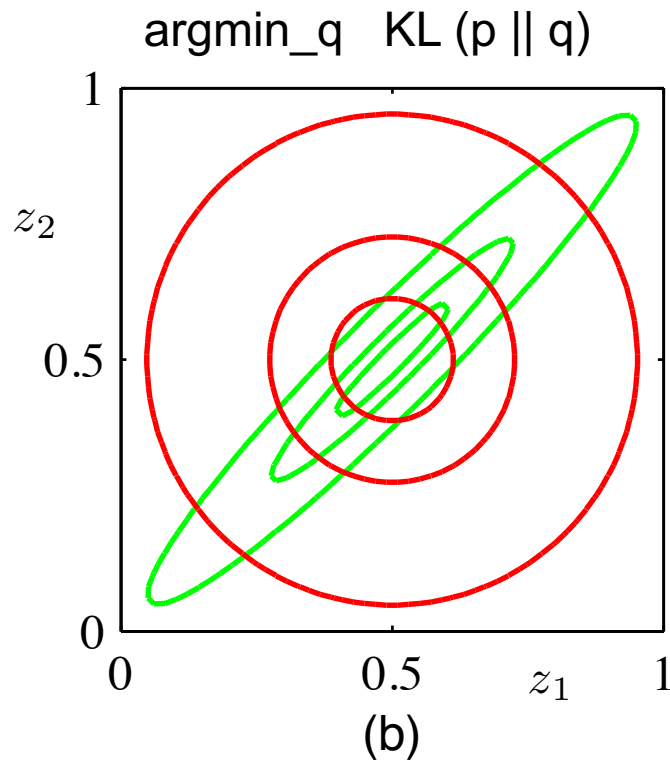
Find simple approximate distribution

- 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

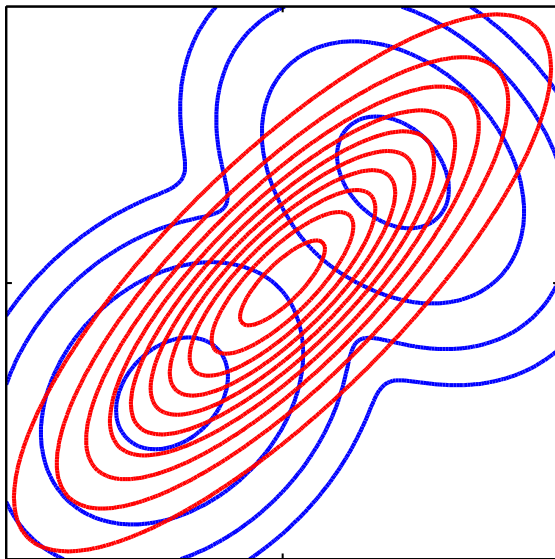
- p = 2D Gaussian with arbitrary co-variance
- q = 2D Gaussian with diagonal co-variance



Example 2

- p = Mixture of Two Gaussians
- q = Single Gaussian

argmin_q KL ($p \parallel q$)



argmin_q KL ($q \parallel p$)

