



Computer
Science

CSC580: Probabilistic Graphical Models

Final Exam Review

Jason Pacheco

Final Exam

- Similar format to Midterm but take-home
- Some students need to take it early so I will release in the next day or two (Due: 12/13)
- 6+1 Questions
 - 1 of these is only for CSC580 students
 - No coding (might use minimal code for one problem)

Random Variables

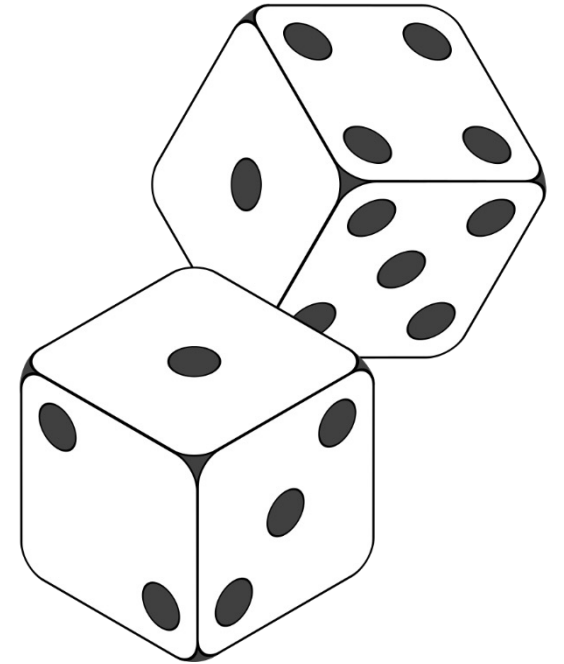
(Informally) A random variable is an unknown quantity that maps events to numeric values.

Example X is the *sum of two dice* with values,

$$X \in \{2, 3, 4, \dots, 12\}$$

Example Flip a coin and let random variable Y represent the outcome,

$$Y \in \{\text{Heads}, \text{Tails}\}$$



Random Variables and Probability

Capitol letters represent
random variables

Lowercase letters are
realized *values*

$$X = x$$

$X = x$ is the **event** that X takes the value x

Example Let X be the random variable (RV) representing the sum of two dice with values,

$$X \in \{2, 3, 4, \dots, 12\}$$

$X=5$ is the *event* that the dice sum to 5.

Probability Mass Function

A function $p(X)$ is a **probability mass function (PMF)** of a discrete random variable if the following conditions hold:

(a) It is nonnegative for all values in the support,

$$p(X = x) \geq 0$$

(b) The sum over all values in the support is 1,

$$\sum_x p(X = x) = 1$$

Intuition Probability mass is conserved, just as in physical mass. Reducing probability mass of one event must increase probability mass of other events so that the definition holds...

Joint Probability

Definition Two (discrete) RVs X and Y have a *joint PMF* denoted by $p(X, Y)$ and the probability of the event $X=x$ and $Y=y$ denoted by $p(X = x, Y = y)$ where,

(a) It is nonnegative for all values in the support,

$$p(X = x, Y = y) \geq 0$$

(b) The sum over all values in the support is 1,

$$\sum_x \sum_y p(X = x, Y = y) = 1$$

Joint Probability

Let X and Y be *binary RVs*. We can represent the joint PMF $p(X, Y)$ as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

All values are nonnegative

Joint Probability

Let X and Y be *binary RVs*. We can represent the joint PMF $p(X,Y)$ as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

**The sum over all values is 1:
 $0.04 + 0.36 + 0.30 + 0.30 = 1$**

Joint Probability

Let X and Y be *binary RVs*. We can represent the joint PMF $p(X, Y)$ as a 2x2 array (table):

		Y	
		0	1
X	0	0.04	0.36
	1	0.30	0.30

$$P(X=1, Y=0) = 0.30$$

Tabular Method

Let X, Y be binary RVs with the joint probability table

For Binomial use K-by-K probability table.

		Y	
		y_1	y_2
X	x_1	0.04	0.36
	x_2	0.30	0.30

0.4 $P(x_1)$

0.6 $P(x_2)$

$P(x)$

$P(y_1) = P(x_1, y_1) + P(x_2, y_1)$
 $P(y_2) = P(x_1, y_2) + P(x_2, y_2)$
[i.e., sum down columns]

0.34 $P(y_1)$

0.66 $P(y_2)$

$P(y)$

$P(x_1) = P(x_1, y_1) + P(x_1, y_2)$
 $P(x_2) = P(x_2, y_1) + P(x_2, y_2)$
[i.e., sum across rows]

Tabular Method

We don't care about event $Y=y_2$

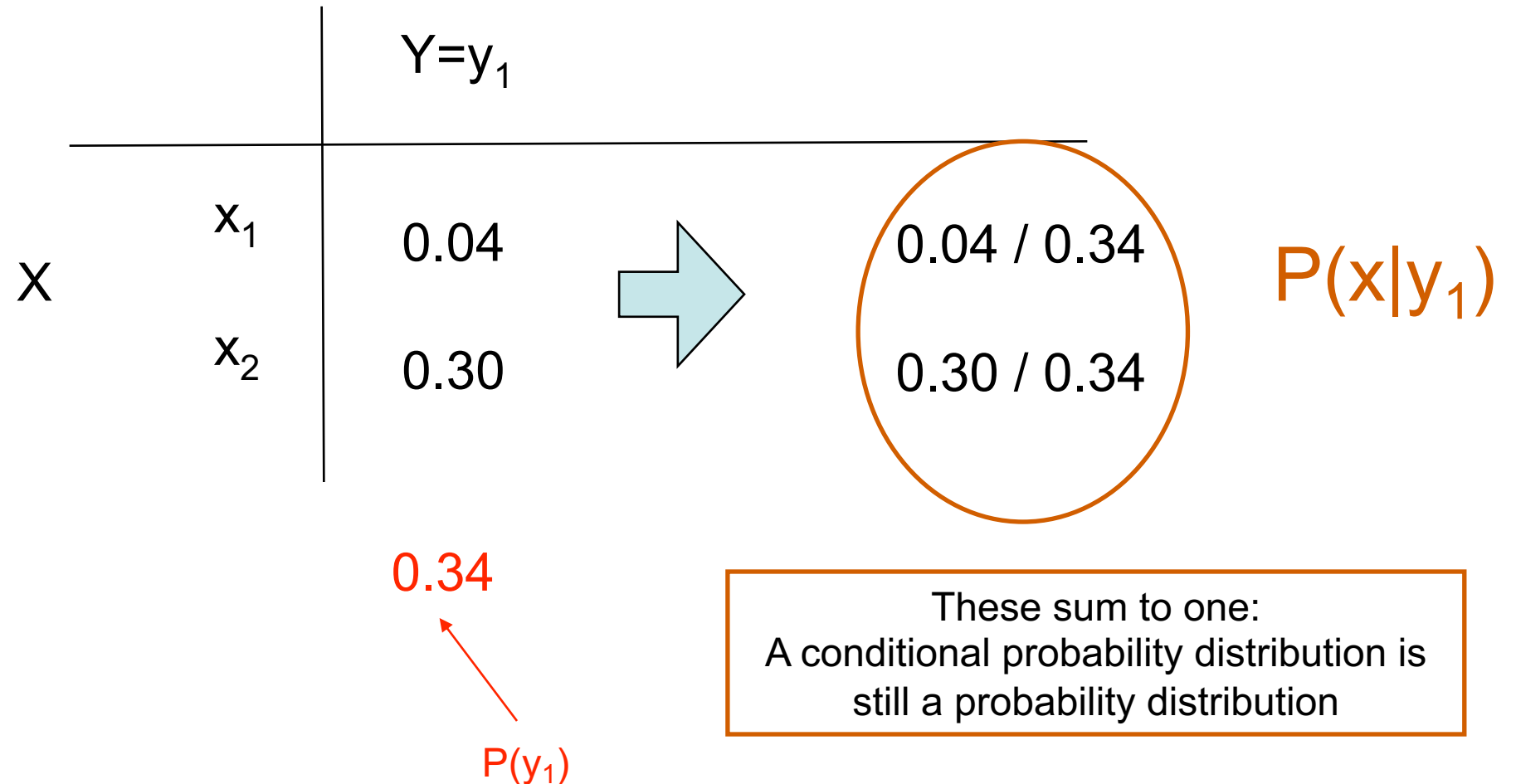
		Y	
		y_1	y_2
X	x_1	0.04	Censored!
	x_2	0.30	

$P(x|y_1)=?$

0.34

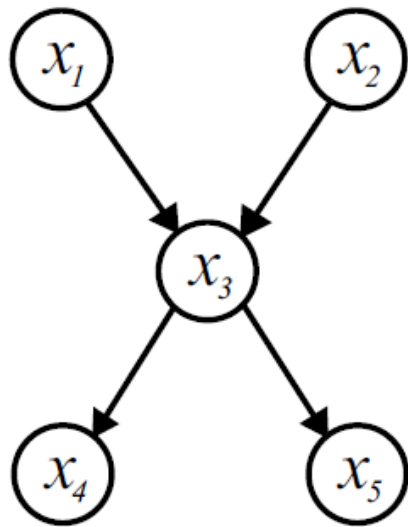
$P(y_1)$

Tabular Method



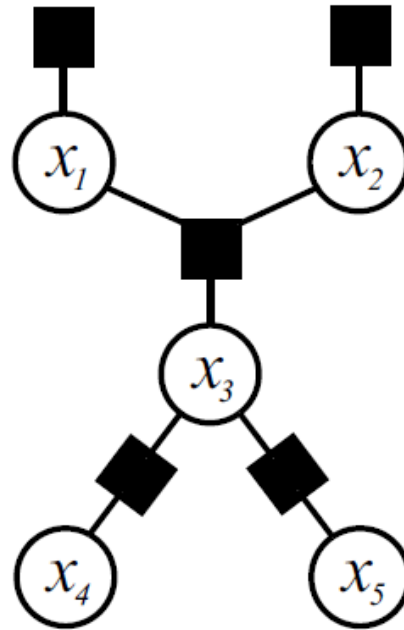
Graphical Models

A variety of graphical models can represent the same probability distribution

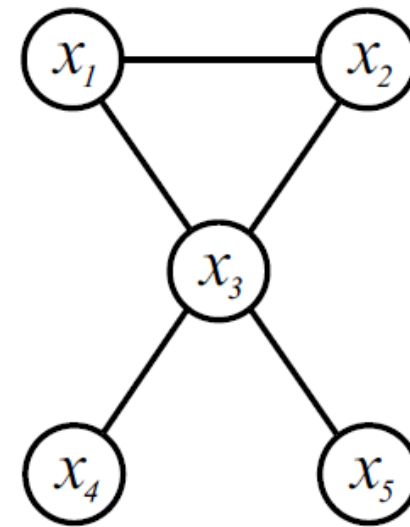


Bayes Network

Directed Models



Factor Graph

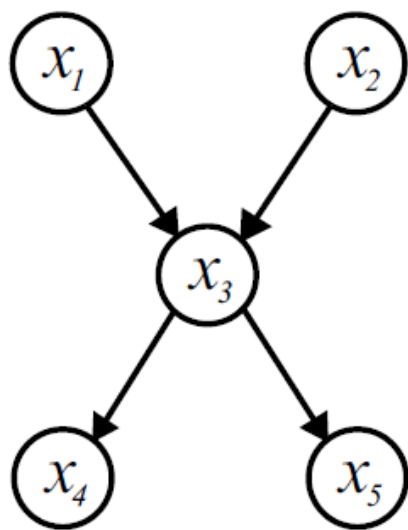


Markov Random Field

Undirected Models

Graphical Models

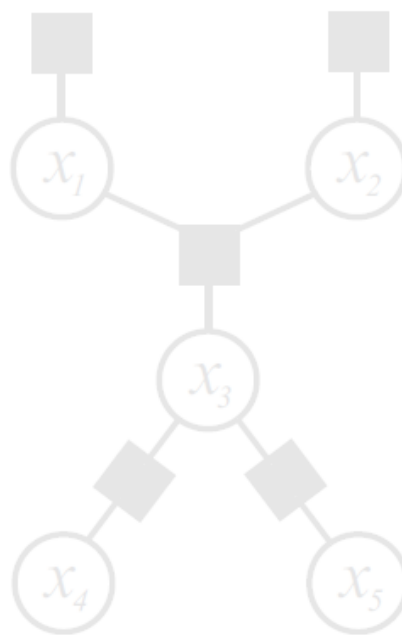
A variety of graphical models can represent the same probability distribution



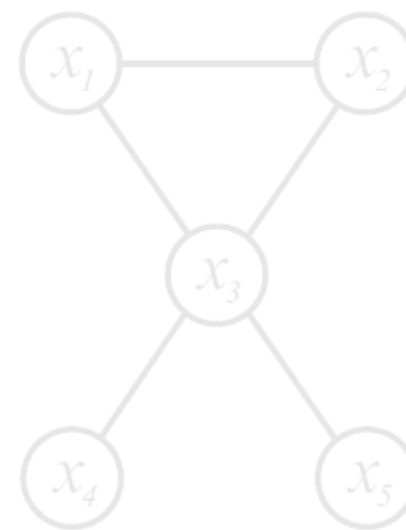
Bayes Network



Directed Models



Factor Graph



Markov Random Field

Undirected Models

From Probabilities to Pictures

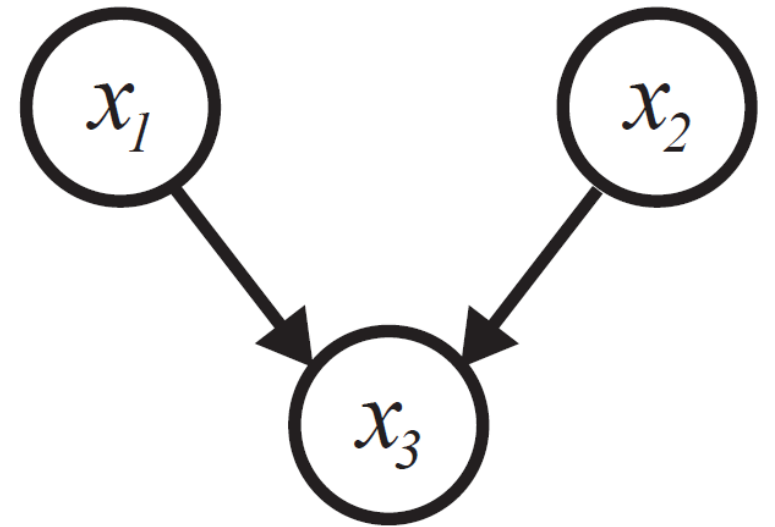
A probabilistic graphical model allows us to pictorially represent a probability distribution

Probability Model:

$$p(x_1, x_2, x_3) = p(x_1)p(x_2)p(x_3 \mid x_1, x_2)$$



Graphical Model:

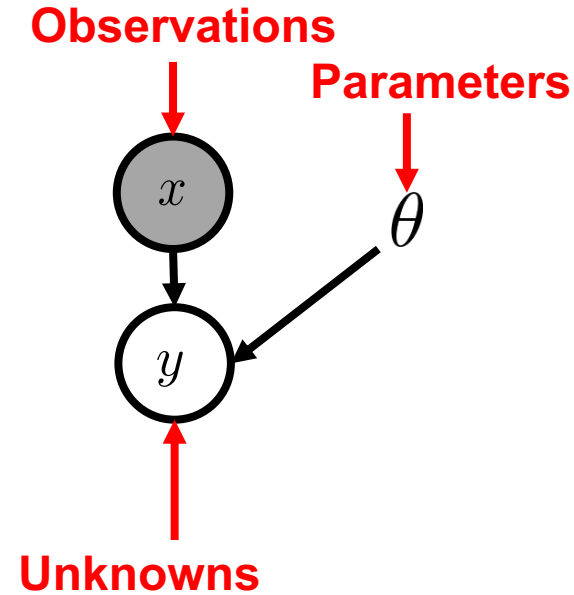
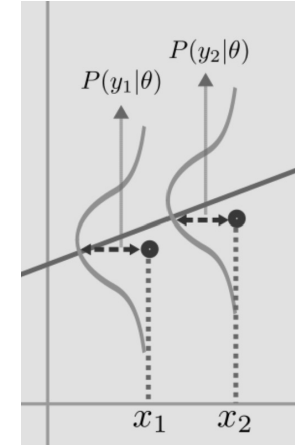


Conditional distribution on each RV is dependent on its parent nodes in the graph

Discriminative vs Generative modeling

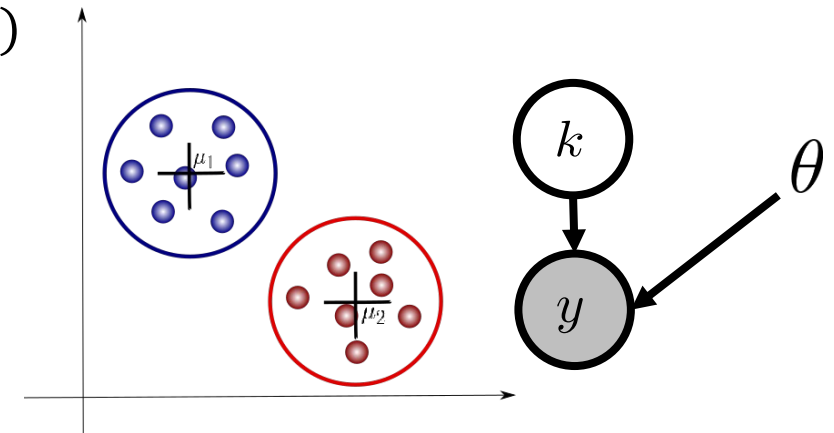
Discriminative model:

- Only models $P(y | x, \theta)$ -- i.e. *doesn't model data x*
- Recall linear regression: $y | x; \theta \sim N(x^\top \theta, \sigma^2)$
- Logistic regression: $y | x; \theta \sim \text{Bernoulli}(\sigma(x^\top \theta))$



Generative model:

- Models everything including data: $P(k, y) = P(k)P(y | k, \theta)$
- e.g., Gaussian mixture model (GMM)
 - $\theta = (\pi_k, \mu_k, \Sigma_k)_{k=1}^K$
 - $k \sim \text{Categorical}(\pi)$ (*hidden*), i.e. $P(k = l) = \pi_l$
 - $y | k \sim N(\mu_k, \Sigma_k)$



Barbershop Example

Suppose you go to a barbershop at every last Friday of the month. You want to be able to predict the waiting time. You have collected 12 data points (i.e., how long it took to be served) from the last year: $S = \{x_1, \dots, x_{12}\}$

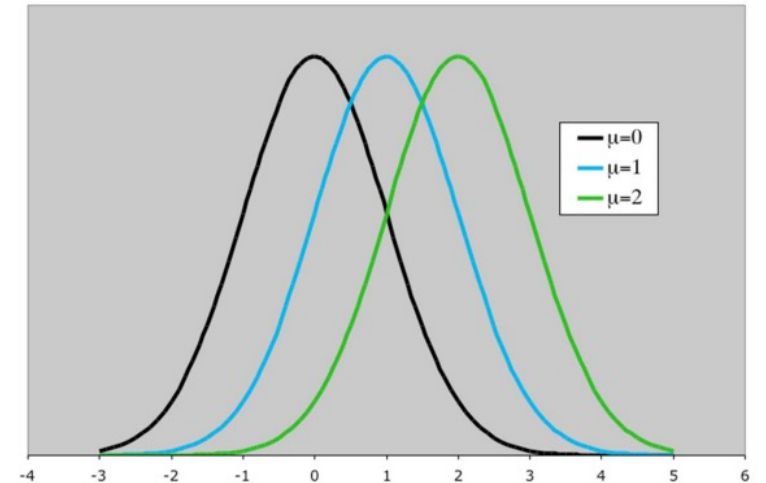
- 1. Modeling assumption: $x_i \sim$ Gaussian distribution $N(\mu, 1)$

- $p(x; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2}\right)$
- Observation: this distribution has mean μ

- 2. Find the MLE $\hat{\mu}$ from data S

- (2.1) write down the neg. log likelihood of the sample

$$L_n(\mu) = -\ln P(x_1, \dots, x_n; \mu) = 12 \ln \sqrt{2\pi} + \frac{1}{2} \sum_{i=1}^{12} (x_i - \mu)^2$$



Is this a generative or discriminative model?

Generative model: basic example I (cont'd)

2. Find the MLE $\hat{\mu}$ from data S

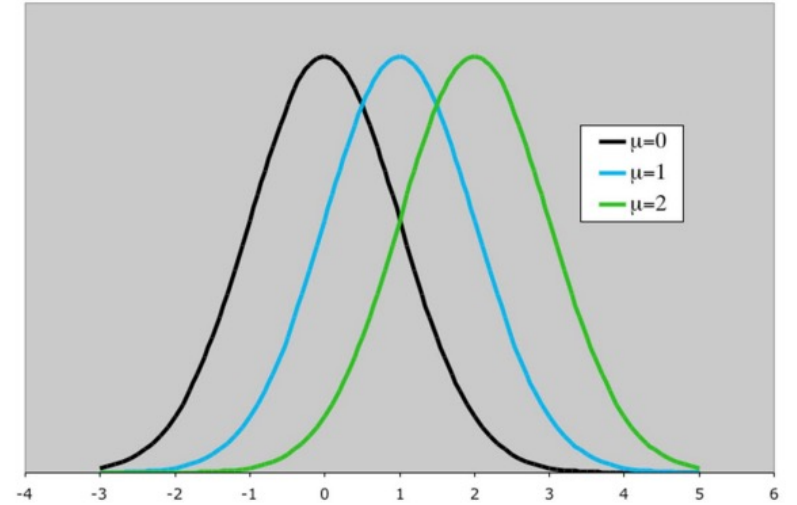
- (2.2) compute the first derivative, set it to 0, solve for λ (be sure to check convexity)

$$L'_n(\mu) = \sum_{i=1}^{12} (x_i - \mu) = 0 \Rightarrow \mu = \frac{x_1 + \dots + x_{12}}{12}$$

Sample Mean

3. The learned model $N(\hat{\mu}, 1)$ is yours!

- Simple prediction: e.g., predict the next wait time by $\mathbb{E}_{X \sim N(\hat{\mu}, 1)}[X]$
- which is $\hat{\mu} = \frac{x_1 + \dots + x_{12}}{12}$



4. (Optional: Model Checking) Generate some data... Does it look realistic?

Conditional Independence

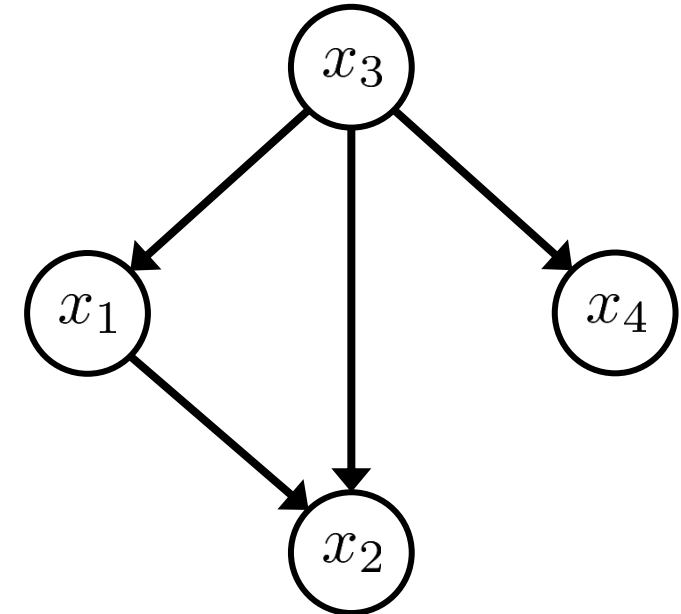
Recall two RVs X and Y are **conditionally independent** given Z (or $X \perp Y \mid Z$) iff:

$$p(X \mid Y, Z) = p(X \mid Z)$$

Idea Apply *chain rule* with ordering that exploits conditional independencies to simplify the terms

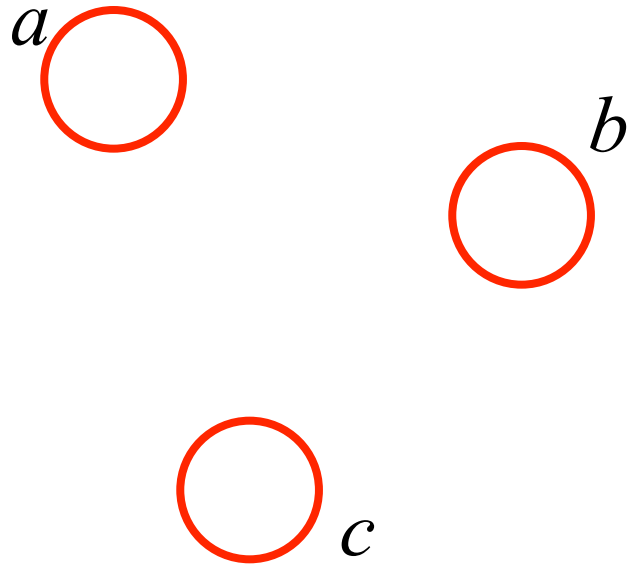
Ex. Suppose $x_4 \perp x_1 \mid x_3$ and $x_2 \perp x_4 \mid x_1$ then:

$$\begin{aligned} p(x) &= p(x_3)p(x_1 \mid x_3)p(x_4 \mid x_1, x_3)p(x_2 \mid x_1, x_3, x_4) \\ &= p(x_3)p(x_1 \mid x_3)p(x_4 \mid x_3)p(x_2 \mid x_1, x_3) \end{aligned}$$

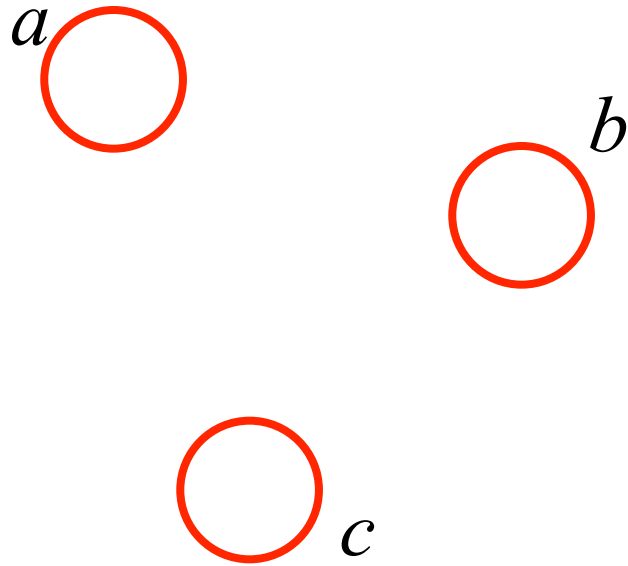


Can visualize conditional dependencies using **directed acyclic graph (DAG)**

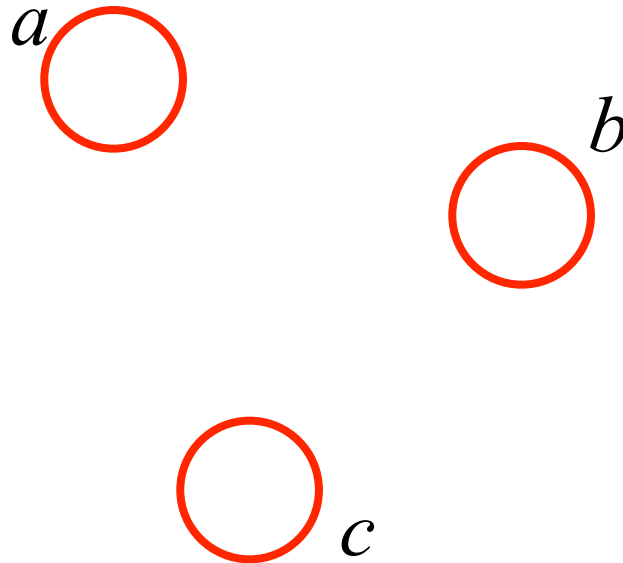
What is the joint factorization?



$$p(\mathbf{a}, \mathbf{b}, \mathbf{c}) = p(\mathbf{a})p(\mathbf{b})p(\mathbf{c})$$

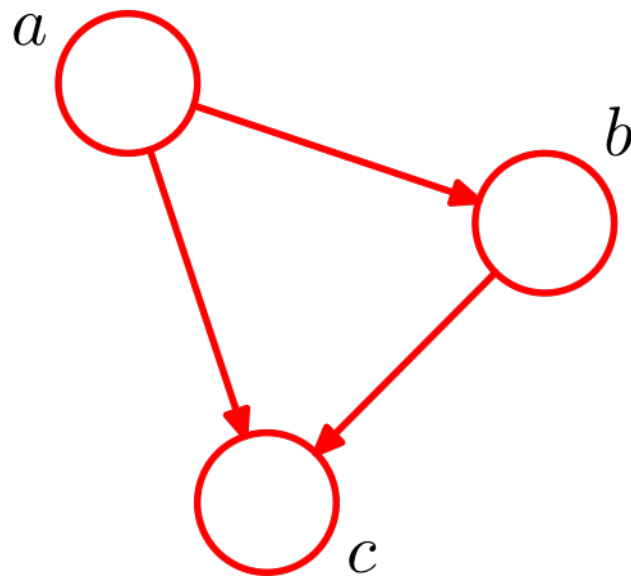


Are a and b independent ($a \perp b$)?



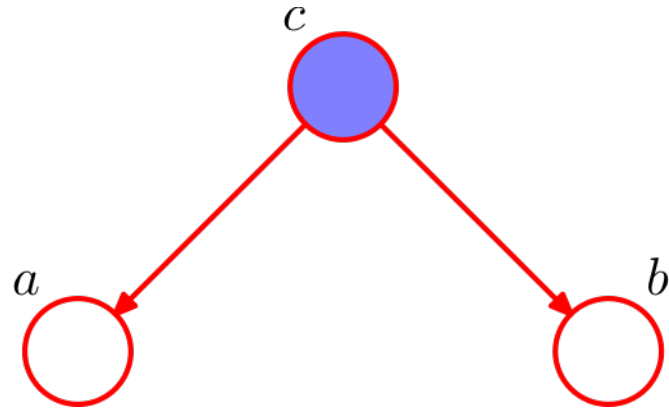
$$\mathbf{p(a,b,c) = p(a)p(b)p(c)}$$

$$p(a,b,c) = p(a)p(b|a)p(c|a,b)$$



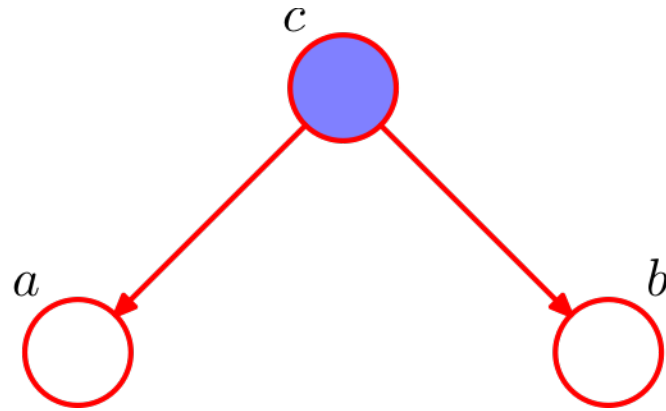
Note there are **no conditional independencies**

Case one where c is observed



Is $a \perp b \mid c$?

Case one where c is observed



$$a \perp b \mid c$$

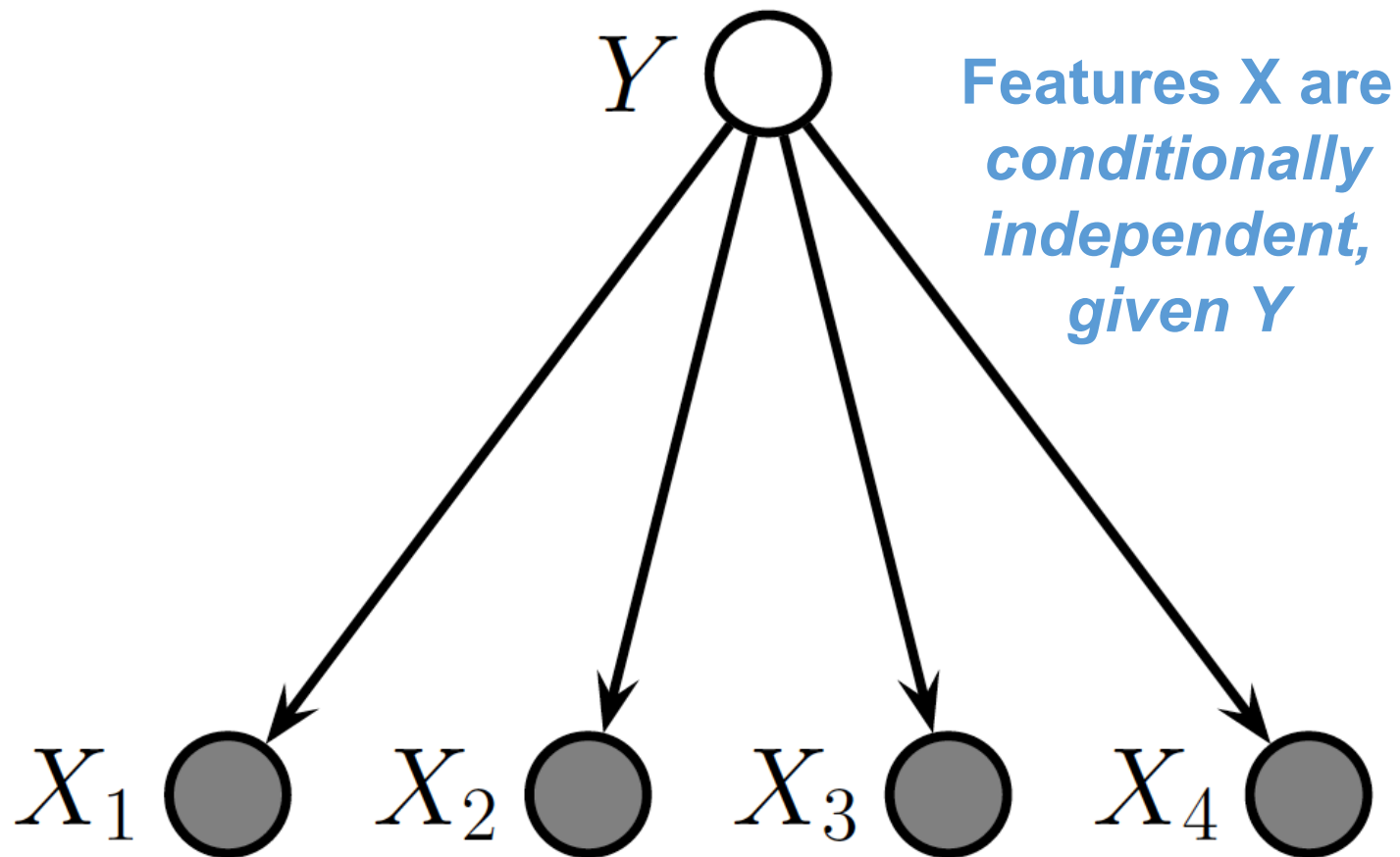
$$p(a, b, c) = p(c)p(a|c)p(b|c) \quad (\text{what the graph represents in general})$$

$$p(a, b|c) = p(a|c)p(b|c) \quad (\text{with } c \text{ observed})$$

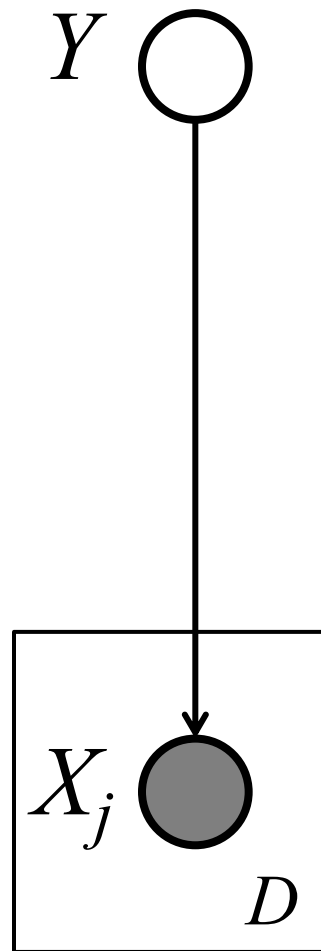
This is the definition of $a \perp b \mid c$

Shading & Plate Notation

Convention: Shaded nodes are observed, open nodes are latent/hidden/unobserved



Features X are
*conditionally
independent,
given Y*



*Plates denote
replication of
random variables*

$$p(y, \mathbf{x}) = p(y) \prod_{j=1}^D p(x_j | y)$$

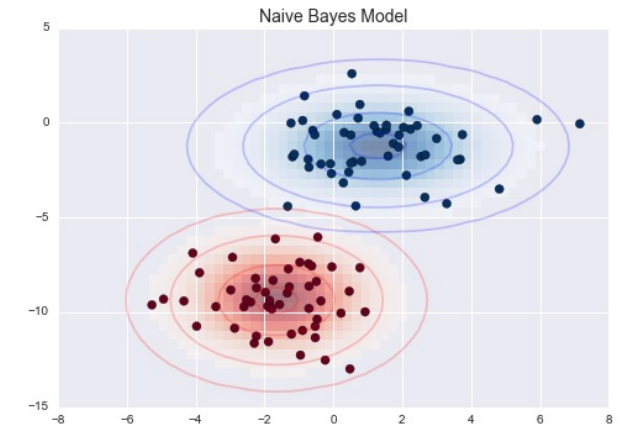
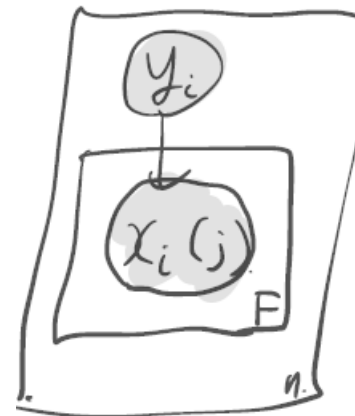
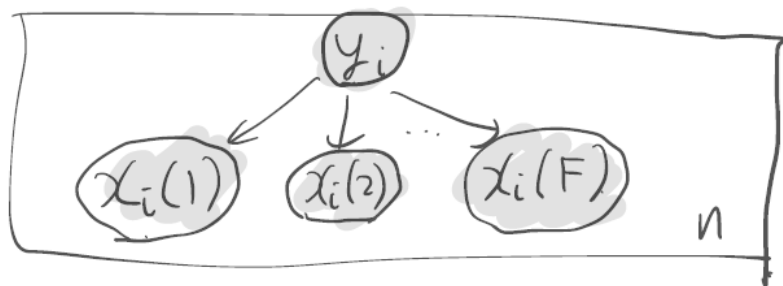
Naïve Bayes for supervised learning

- Motivation: supervised learning for classification
- high-dimensional $x = (x(1), \dots, x(F))$, modeling $P(x | y)$ can be tricky
- In general, $P(x | y) = P(x(1) | y) \cdot P(x(2) | x(1), y) \cdot \dots \cdot P(x(F) | x(1), \dots, x(F - 1), y)$
- A modeling assumption: $x(1), \dots, x(F)$ are conditionally independent given y
i.e. for all i

$$x(i) \perp\!\!\!\perp (x(1), \dots, x(i - 1), x(i + 1), \dots, x(F)) | y$$

(Conditional independence notation: $A \perp\!\!\!\perp B | C$)

- Equivalently $P(x | y) = P(x(1) | y) \cdot \dots \cdot P(x(F) | y)$



Recall : Class Preference Prediction

Define the labeled training dataset $S = \{(x_i, y_i)\}_{i=1}^m$

To make this a binary classification we set
“Liked” = $\{+2, +1, 0\}$
“Nah” = $\{-1, -2\}$

	Rating	Easy?	AI?	Sys?	Thy?	Morning?
Features	+2	y	y	n	y	n
Feature Values	+2	y	y	n	y	n
	+2	n	y	n	n	n
	+2	n	n	n	y	n
	+2	n	y	y	n	y
	+1	y	y	n	n	n
	+1	y	y	n	y	n
	+1	n	y	n	y	n
Labels	0	n	n	n	n	y
	0	y	n	n	y	y
	0	n	y	n	y	n
	0	y	y	y	y	y
	-1	y	y	y	n	y
	-1	n	n	y	y	n
	-1	n	n	y	n	y
	-1	y	n	y	n	y
	-2	n	n	y	y	n
	-2	n	y	y	n	y
Data Point	-2	y	n	y	n	n
	-2	y	n	y	n	y

Naïve Bayes: binary-valued features

Training Data $S = \{(x_i, y_i)\}_{i=1}^n$,

$$x_i \in \{0,1\}^F$$

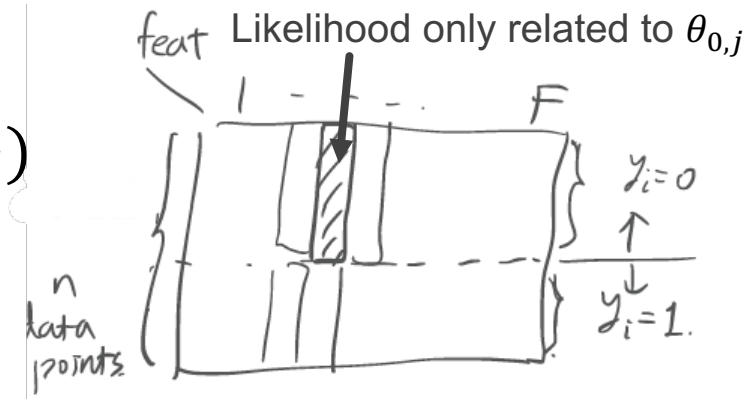
$$y_i \in \{0,1\}$$

Generative Story

$y \sim \text{Bernoulli}(\pi)$; for all $j \in [F]$, $x(j) \mid y = c \sim \text{Bernoulli}(\theta_{c,j})$

#parameters = $1 + 2F$

Training (denote by $\theta = \{\theta_{c,j}\}$)

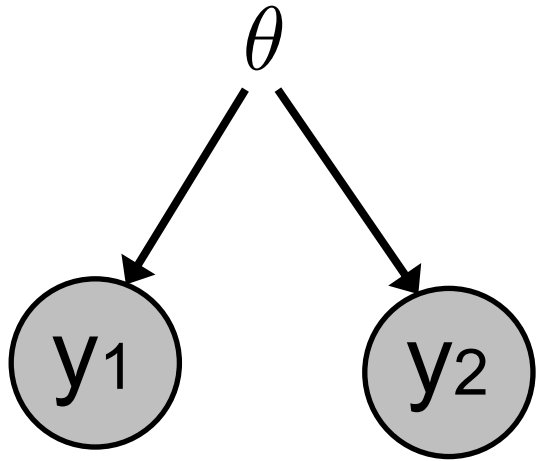


$$\begin{aligned} \max_{\pi, \theta} \sum_{i=1}^n \ln P(x_i, y_i; \pi, \theta) &= \sum_{i=1}^n \ln P(y_i; \pi) + \sum_{i=1}^n \ln P(x_i \mid y_i; \theta) \\ &= \max_{\pi} \sum_{i=1}^n \ln P(y_i; \pi) + \max_{\{\theta_{0,j}\}} \sum_{i:y_i=0} \ln P(x_i \mid y_i; \theta) + \max_{\{\theta_{1,j}\}} \sum_{i:y_i=1} \ln P(x_i \mid y_i; \theta) \end{aligned}$$

Key observation: optimal π , optimal $\{\theta_{0,j}\}$, optimal $\{\theta_{1,j}\}$ can be found separately

$$\text{Optimal } \pi: \max_{\pi} \sum_{i=1}^n \ln P(y_i; \pi) = \max_{\pi} n_0 \ln(1 - \pi) + n_1 \ln(\pi) \Rightarrow \hat{\pi} = \frac{n_1}{n}$$

Learning / Training



Model random data with hyperparameters θ :

$$y \sim p(y | \theta)$$

Sometimes we use:

$$p(y; \theta)$$

Given training data:

$$\{y_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(y | \theta)$$

Learn parameters, e.g. via *maximum likelihood estimation*:

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \log p(y_1, \dots, y_n | \theta)$$

We will talk more
about MLE in
coming weeks

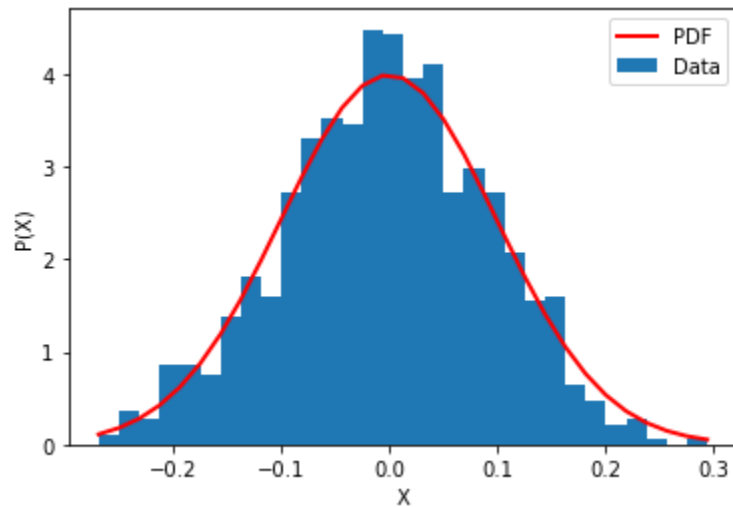
Other estimators are possible:

- *Maximum a posteriori (MAP)*
- *Minimum mean squared error (MMSE)*
- *Etc.*

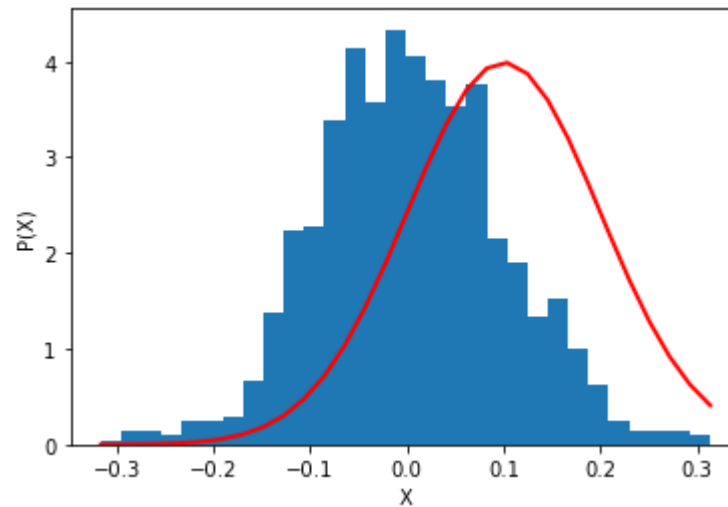
Likelihood (Intuitively)

Suppose we observe N data points from a Gaussian model and wish to estimate model parameters...

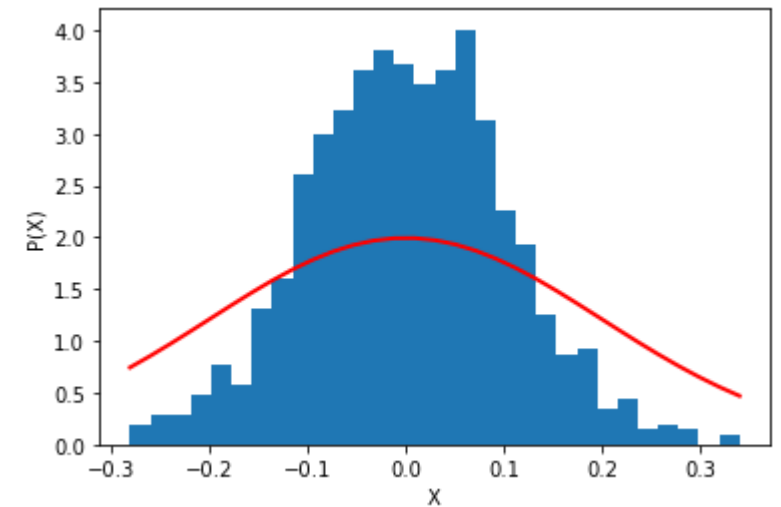
High Likelihood



Low Likelihood (mean)



Low Likelihood (variance)



Likelihood Principle *Given a statistical model, the likelihood function describes all evidence of a parameter that is contained in the data.*

Likelihood Function

Suppose $x_i \sim p(x; \theta)$, then what is the **joint probability** over N *independent identically distributed* (iid) observations x_1, \dots, x_N ?

$$p(x_1, \dots, x_N; \theta) = \prod_{i=1}^N p(x_i; \theta)$$

- We call this the **likelihood function**, often denoted $\mathcal{L}_N(\theta)$
- It is a function of the parameter θ , the data are fixed
- Measures how well parameter θ describes data (*goodness of fit*)

How could we use this to estimate a parameter θ ?

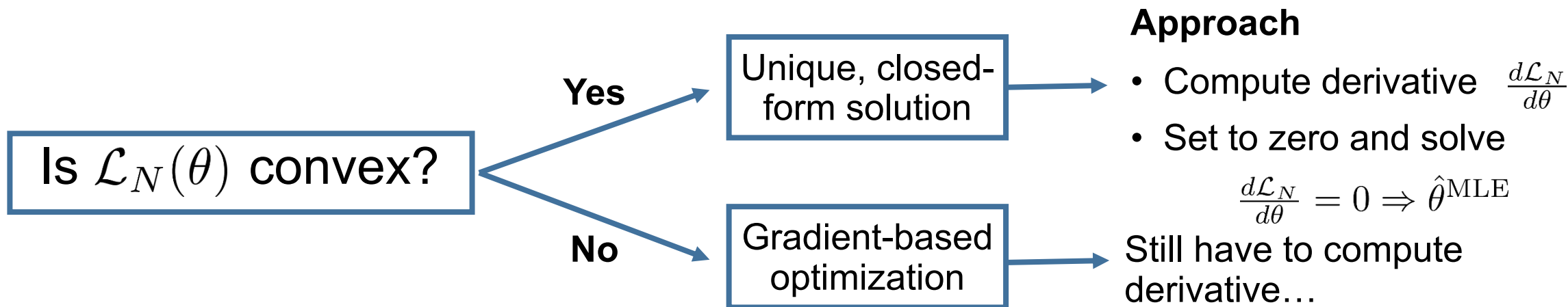
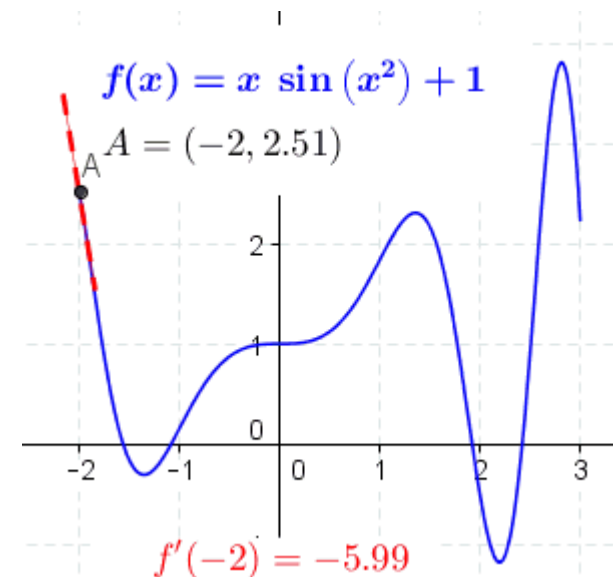
Maximum Likelihood

Maximum Likelihood Estimator (MLE) as the name suggests, maximizes the likelihood function.

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \mathcal{L}_N(\theta) = \prod_{i=1}^N p(x_i; \theta)$$

Question How do we find the MLE?

Answer Remember calculus...



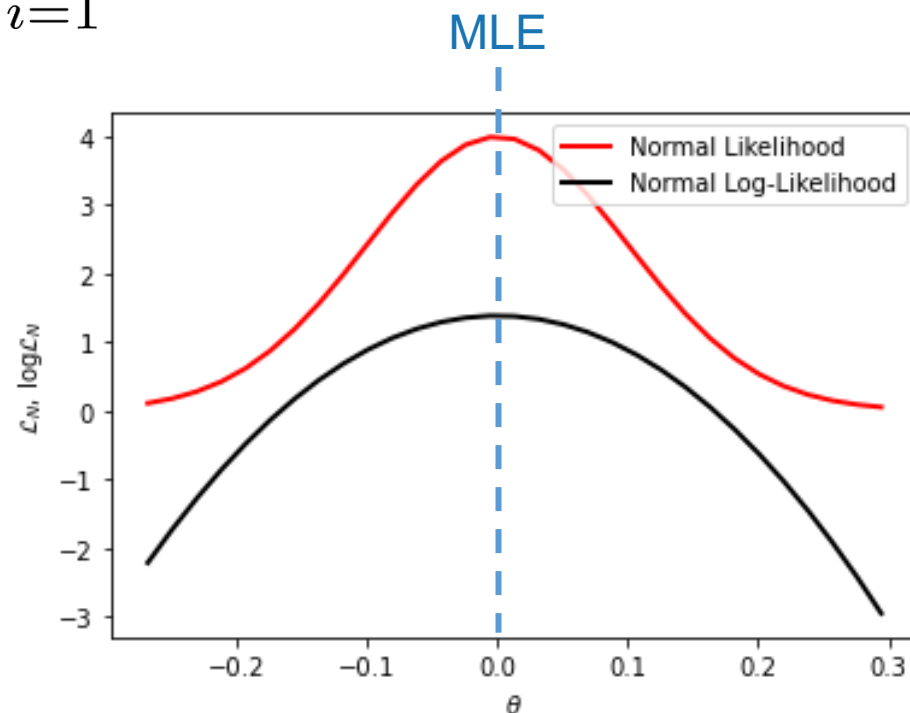
Maximum Likelihood

Maximizing log-likelihood makes the math easier (as we will see) and doesn't change the answer (logarithm is an increasing function)

$$\hat{\theta}^{\text{MLE}} = \arg \max_{\theta} \log \mathcal{L}_N(\theta) = \sum_{i=1}^N \log p(x_i; \theta)$$

Derivative is a linear operator so,

$$\frac{d}{d\theta} \log \mathcal{L}_N(\theta) = \underbrace{\sum_{i=1}^N \frac{d}{d\theta} \log p(x_i; \theta)}_{\substack{\text{One term per data point} \\ \text{Can be computed in parallel} \\ \text{(big data)}}$$



Maximum Likelihood

Example Suppose we have N coin tosses with $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ but we don't know the coin bias p . The likelihood function is,

$$\mathcal{L}_n(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^S (1-p)^{n-S}$$

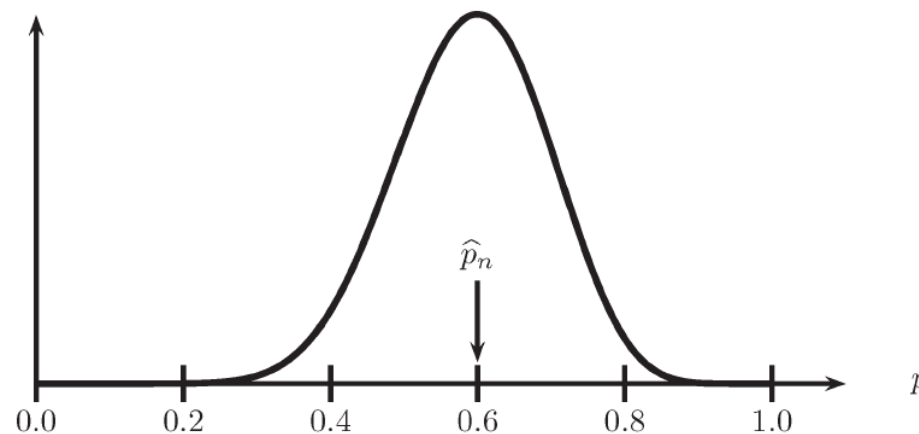
where $S = \sum_i x_i$. The log-likelihood is,

$$\log \mathcal{L}_n(p) = S \log p + (n - S) \log(1 - p)$$

Set the derivative of $\log \mathcal{L}_n(p)$ to zero and solve,

$$\hat{p}^{\text{MLE}} = S/n = \frac{1}{n} \sum_{i=1}^n x_i$$

[Source: Wasserman, L. 2004]



Likelihood function for Bernoulli with $n=20$ and $\sum_i x_i = 12$ heads

Maximum likelihood is equivalent to sample mean in Bernoulli

Marginal Likelihood

More often, we have a joint distribution with observations y , unknown variables z , and parameters θ

Marginal likelihood is normalizer of posterior:

$$p(z | y) = \frac{p(z)p(y | z)}{p(y)}$$

Bayes' Rule

$$p(z, y | \theta) = p(z | \theta)p(y | z, \theta)$$



Prior



Likelihood

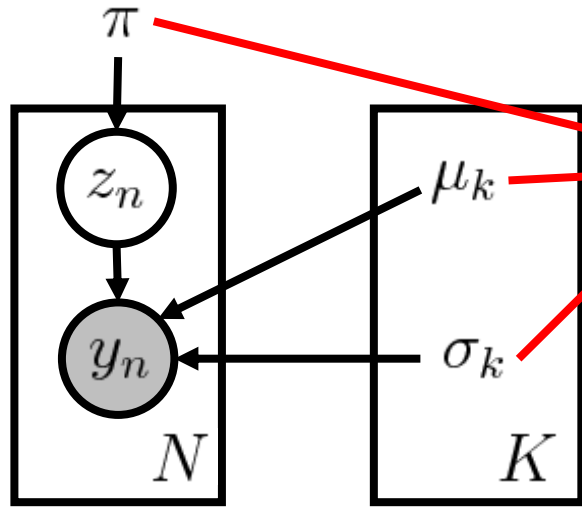
Need to *marginalize* out unknown variables, hence the name *marginal likelihood*:

$$p(y | \theta) = \int p(z | \theta)p(y | z, \theta) dz = \mathcal{L}(\theta)$$

Typically, this integral lacks a closed-form solution...so we need to compute *approximate* MLE solutions

Example: Gaussian Mixture Model

Model is often specified in terms of *unknown parameters*



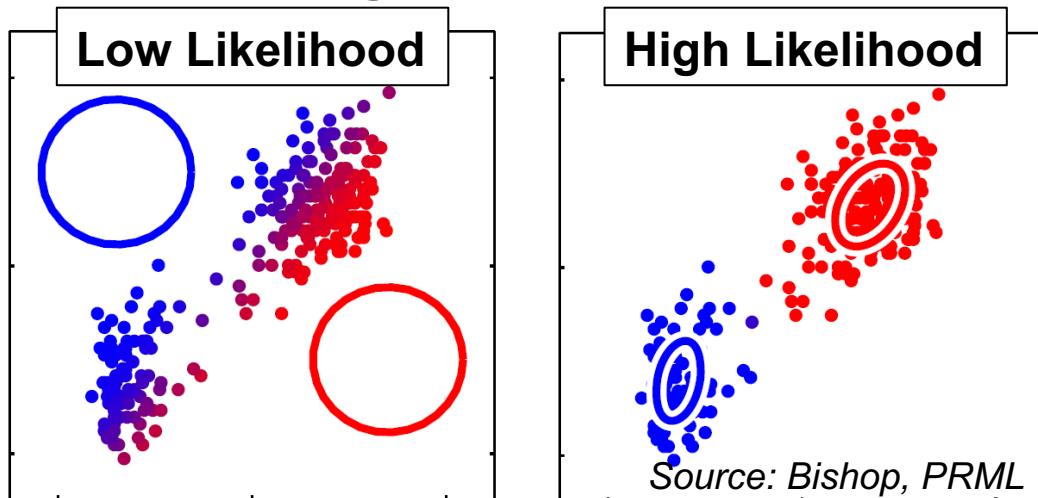
How *likely* are parameters for observed data?

$$\theta = \{\pi, \mu_1, \sigma_1, \dots, \mu_K, \sigma_K\} \quad \mathcal{Y} = \{y_1, \dots, y_N\}$$

Marginal Likelihood (likelihood function):

$$p(\mathcal{Y} | \theta) = \sum_{z_1} \dots \sum_{z_N} p(z_1, \dots, z_N, \mathcal{Y} | \theta)$$

GMM

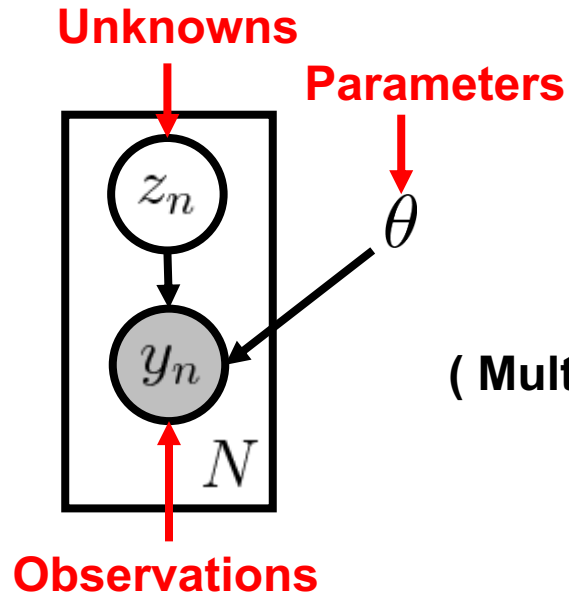


Sum over all possible K^N assignments, which we cannot compute

Intuition Learn / estimate parameters that assign highest probability (under the model) to data we've observed.

Lower Bounding Marginal Likelihood

Conditionally-independent model with partial observations...



$$\log p(\mathcal{Y} \mid \theta) = \log \sum_{z_1} \dots \sum_{z_N} p(z_1, \dots, z_N, \mathcal{Y} \mid \theta)$$

(Multiply by $q(z)/q(z)=1$)

$$= \log \sum_z p(z, \mathcal{Y} \mid \theta) \left(\frac{q(z)}{q(z)} \right)$$

Shorthand

$z = z_1, \dots, z_N$

(Definition of Expected Value)

$$= \log \mathbf{E}_q \left[\frac{p(z, \mathcal{Y} \mid \theta)}{q(z)} \right]$$

$q(z)$ is any distribution with support over Z

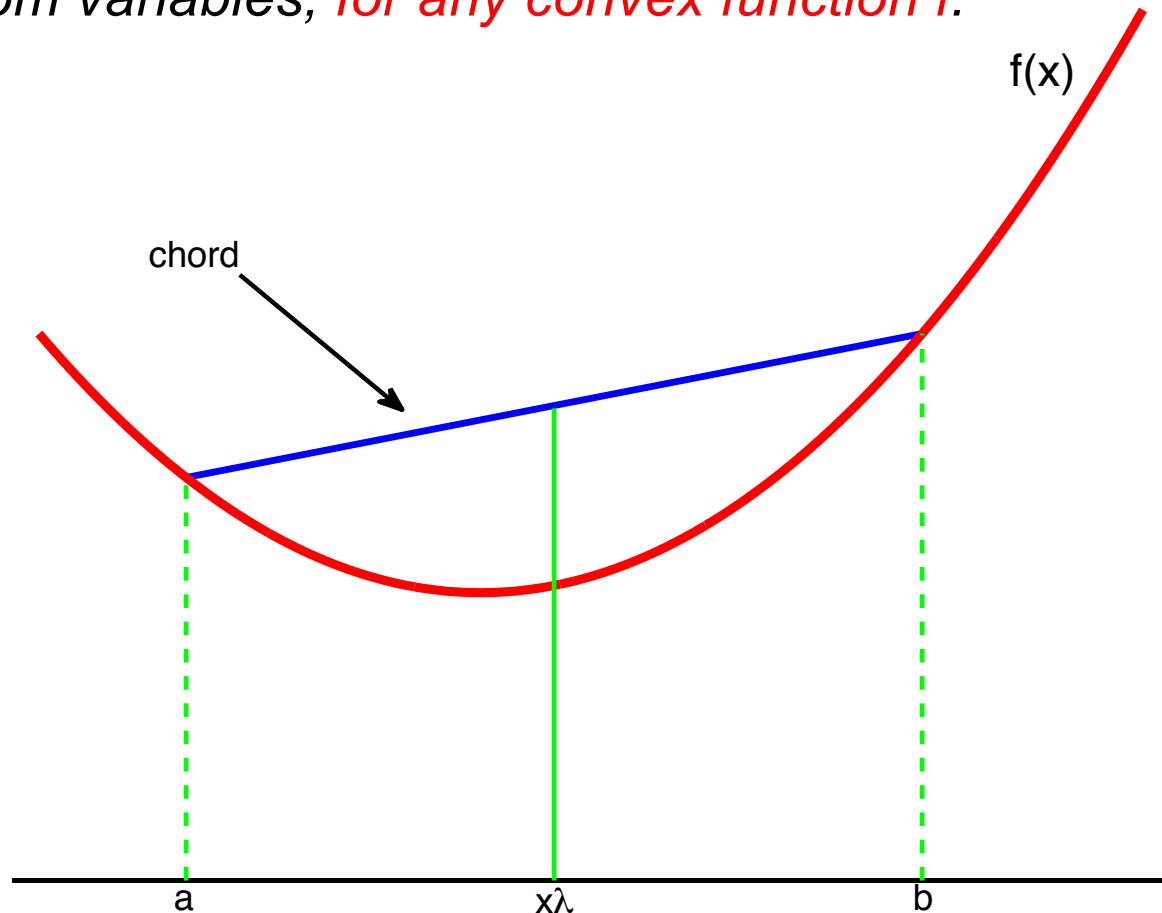
(Jensen's Inequality)

$$\geq \mathbf{E}_q \left[\log \frac{p(z, \mathcal{Y} \mid \theta)}{q(z)} \right]$$

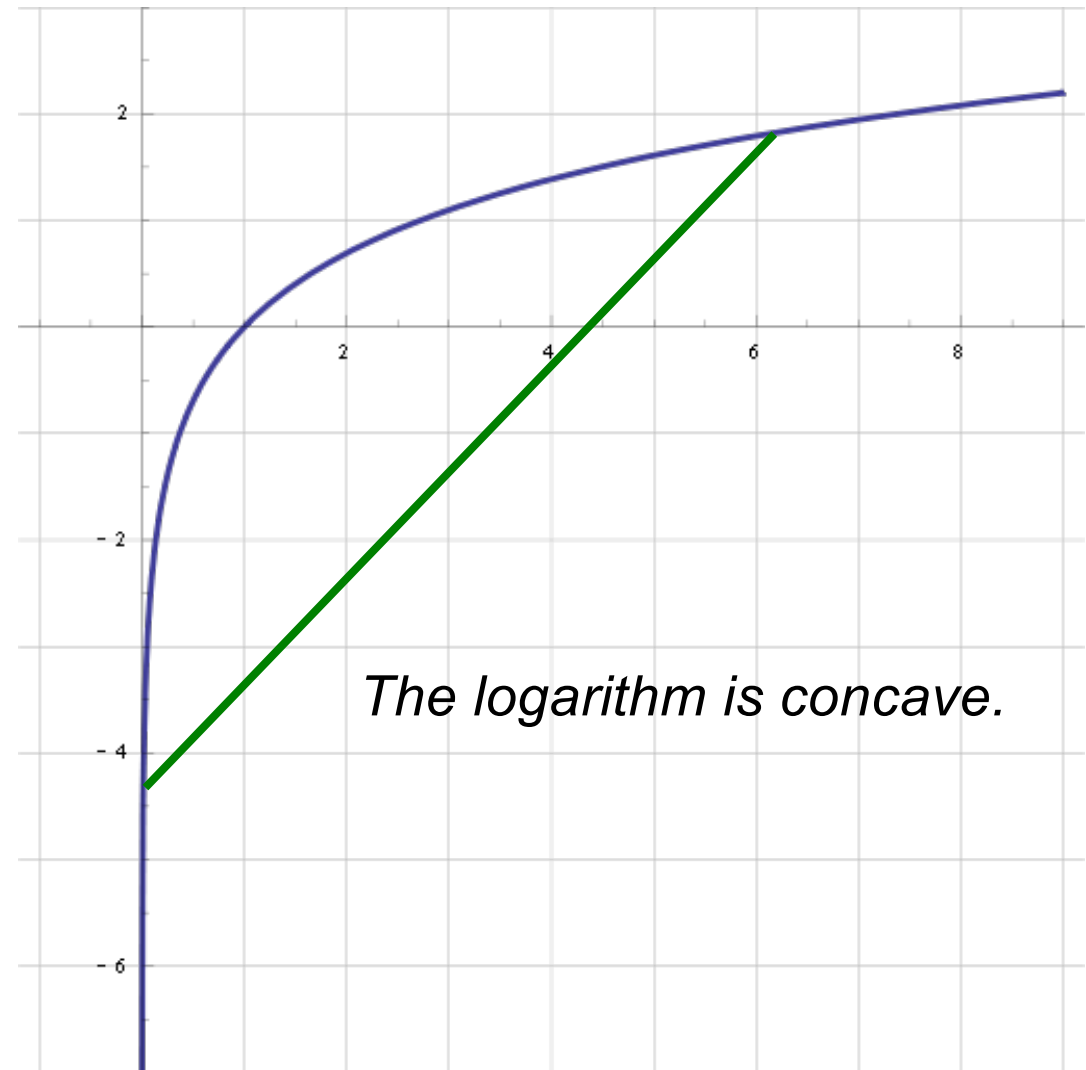
Jensen's Inequality

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$$

Valid for both discrete (expectations are sums)
and continuous (expectations are integrals)
random variables, *for any convex function f .*



$$\ln(\mathbb{E}[X]) \geq \mathbb{E}[\ln(X)]$$



Expectation Maximization

Find tightest lower bound of marginal likelihood,

$$\max_{\theta} \log p(\mathcal{Y} | \theta) \geq \max_{q, \theta} \mathbf{E}_q \left[\log \frac{p(z, \mathcal{Y} | \theta)}{q(z)} \right] \equiv \mathcal{L}(q, \theta)$$

Solve by coordinate ascent...

Initialize Parameters: $\theta^{(0)}$

At iteration t do:

E-Step: $q^{(t)} = \arg \max_q \mathcal{L}(q, \theta^{(t-1)})$

M-Step: $\theta^{(t)} = \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta)$

Until convergence

Fix θ



Fix q



E-Step

$$q^{(t)}(z) = \arg \max_q \mathcal{L}(q, \theta^{(t-1)}) \equiv \mathbf{E}_q \left[\log \frac{p(z, y | \theta^{(t-1)})}{q(z)} \right]$$

Concave (in $q(z)$) and optimum occurs at,

$$q^{(t)}(z) = p(z | y, \theta^{(t-1)})$$

Set $q(z)$ to posterior with current parameters

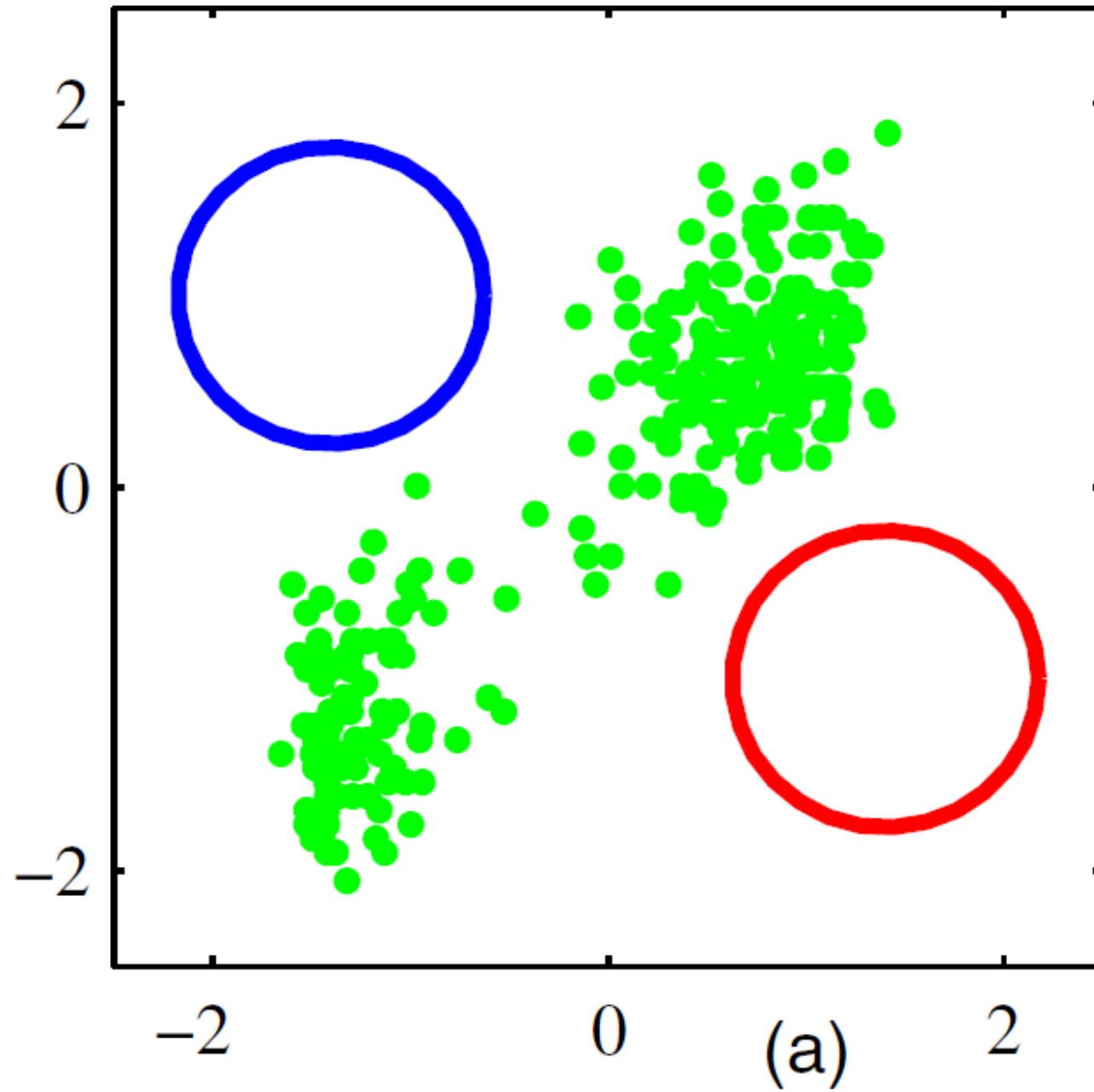
Initialize Parameters: $\theta^{(0)}$

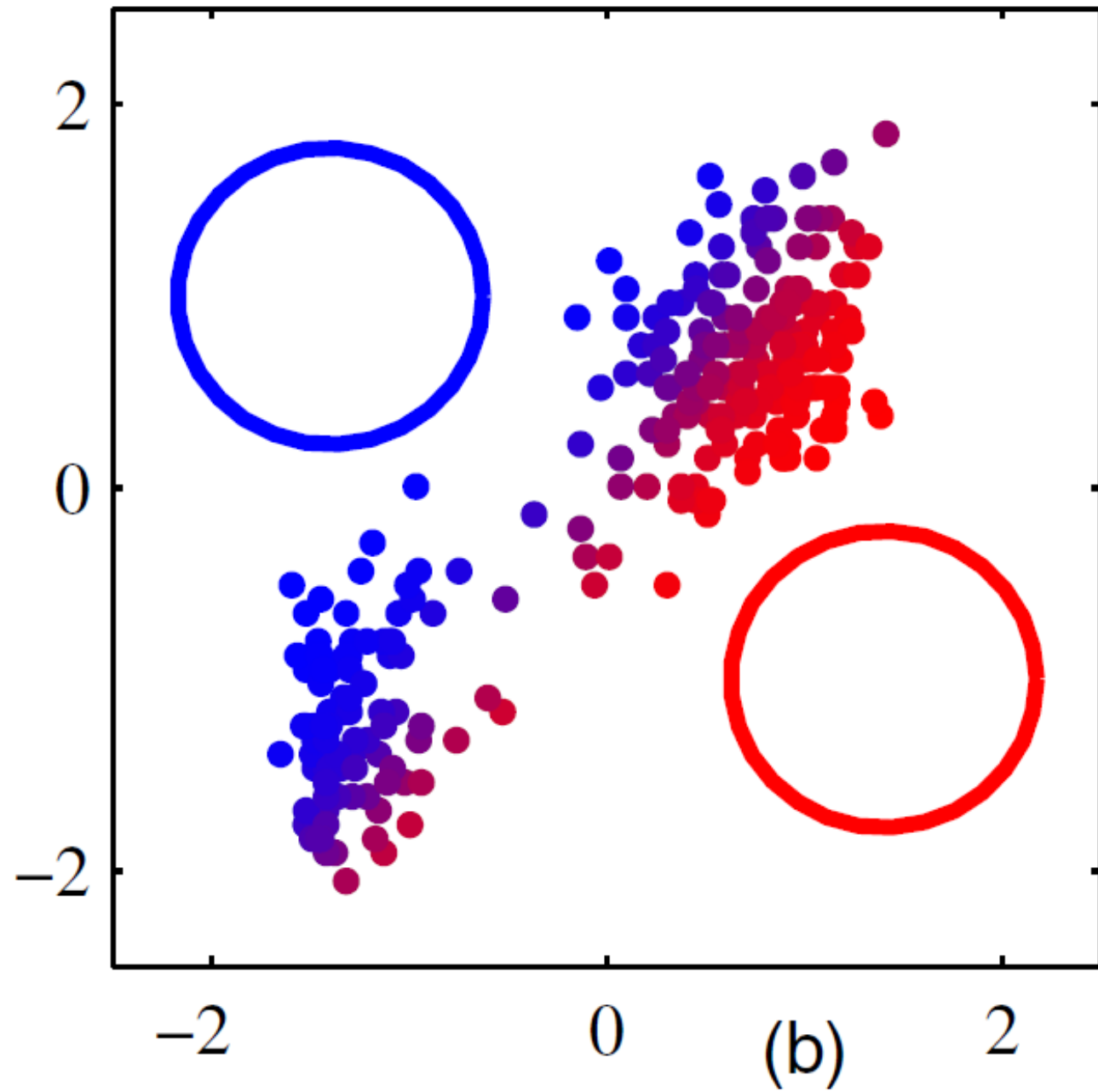
At iteration t do:

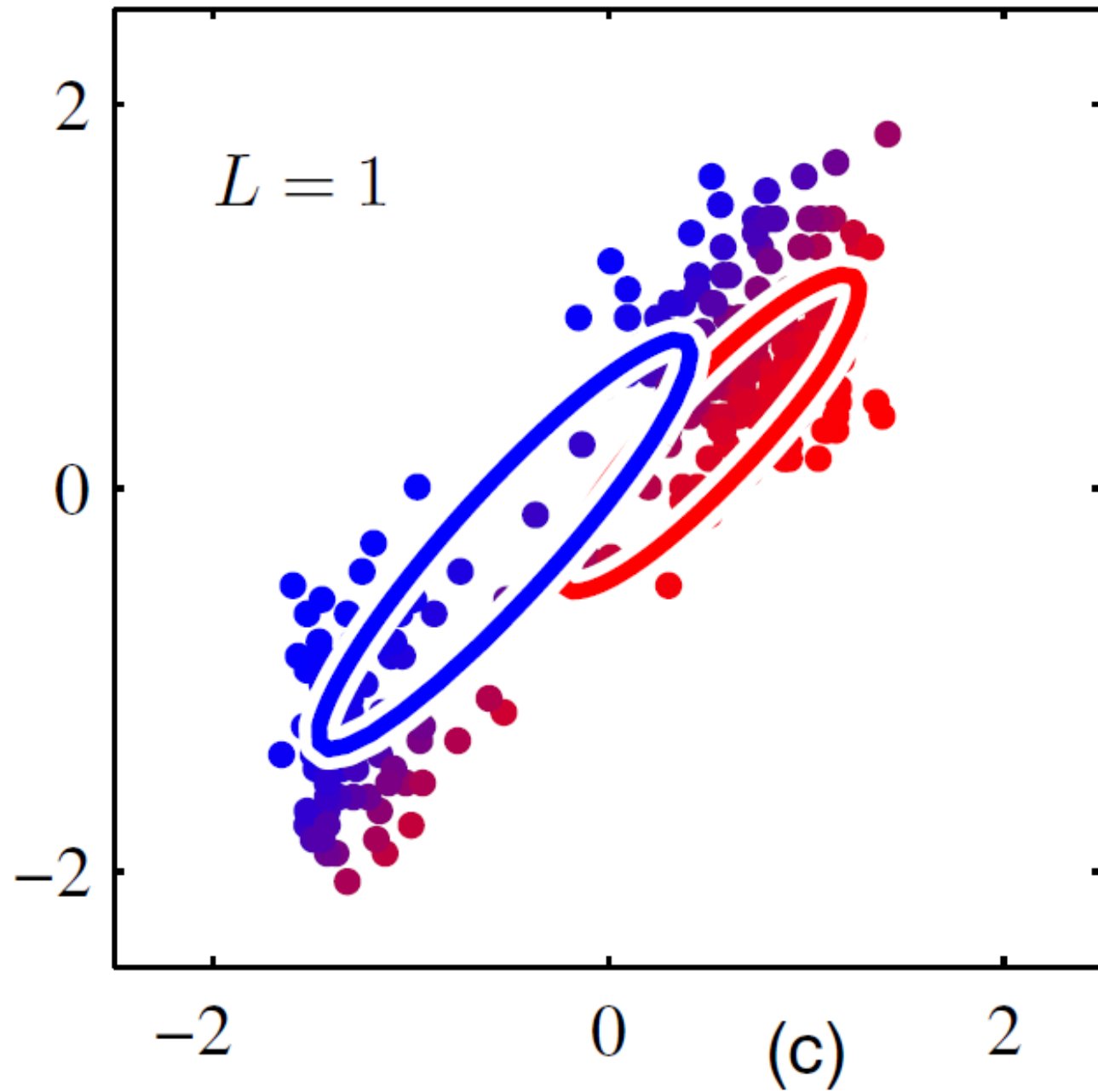
E-Step: $q^{(t)}(z) = p(z | y, \theta^{(t-1)})$

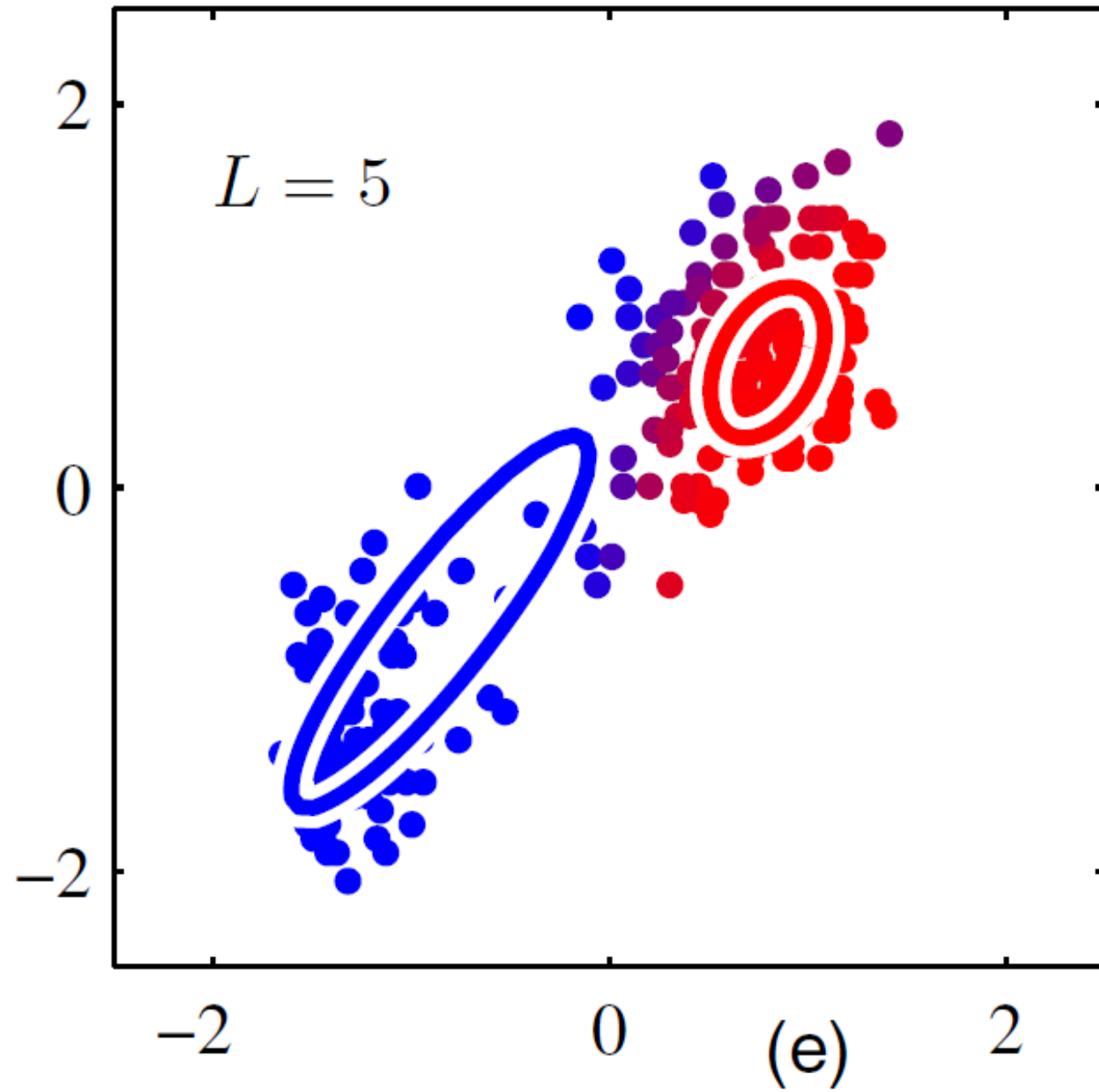
M-Step: $\theta^{(t)} = \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta)$

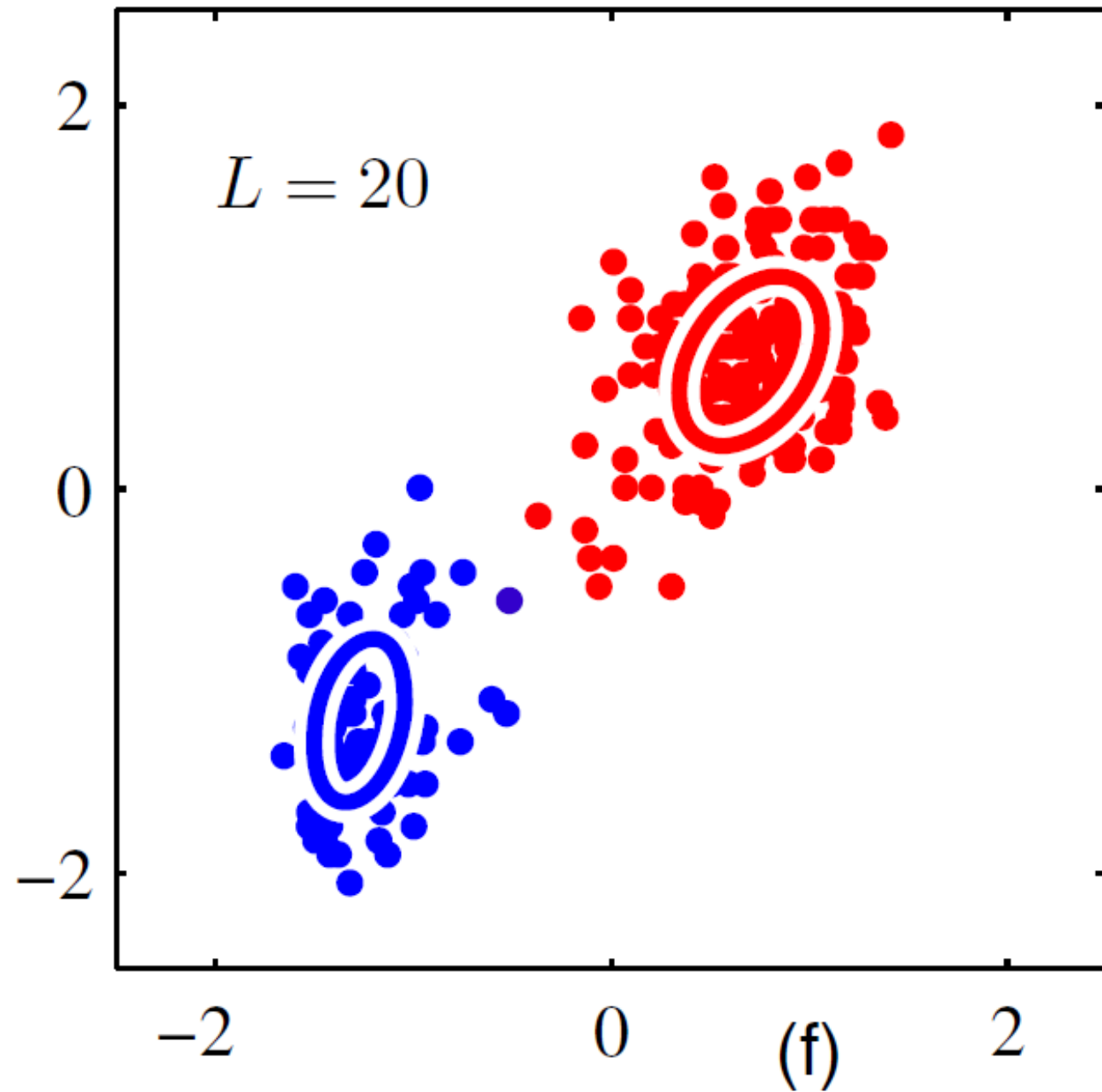
Until convergence



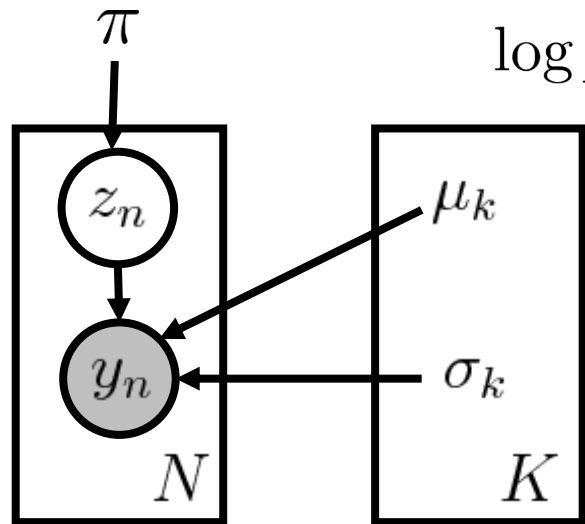








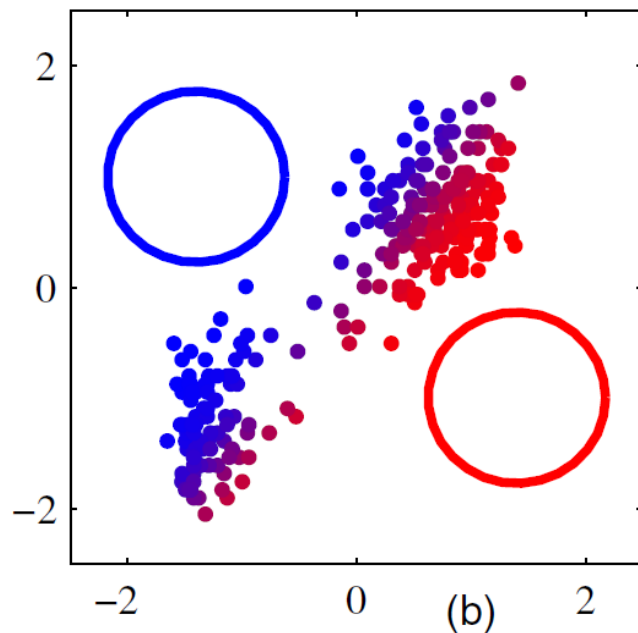
Example: Gaussian Mixture Model



$$\log p(\mathcal{Y} \mid \pi, \mu, \Sigma) \geq \sum_{n=1}^N \sum_{k=1}^K q(z_n = k) \log \{ \pi_k \mathcal{N}(y_n \mid \mu_k, \Sigma_k) \} = \mathcal{L}(q, \theta)$$

E-Step: $q^{\text{new}} = \arg \max_q \mathcal{L}(q, \theta^{\text{old}})$

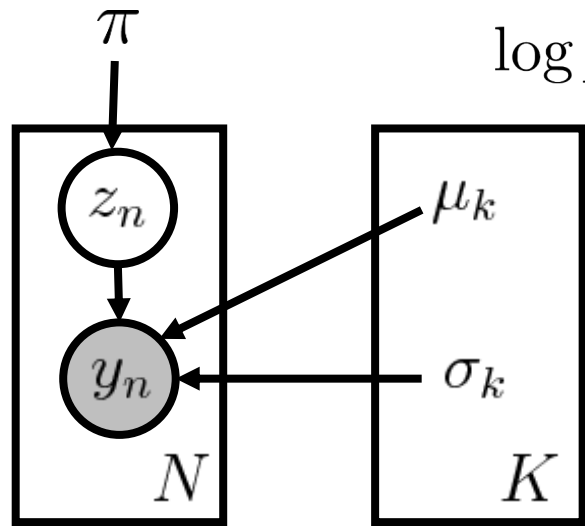
$$q^{\text{new}}(z_n = k) = p(z_n = k \mid \mathcal{Y}, \mu^{\text{old}}, \Sigma^{\text{old}}, \pi^{\text{old}})$$



$$\begin{aligned} &= \frac{p(z_n = k, \mathcal{Y} \mid \mu^{\text{old}}, \Sigma^{\text{old}}, \pi^{\text{old}})}{\sum_{j=1}^K p(z_n = j, \mathcal{Y} \mid \mu^{\text{old}}, \Sigma^{\text{old}}, \pi^{\text{old}})} \\ &= \frac{\pi_k^{\text{old}} \mathcal{N}(y_n \mid \mu_k^{\text{old}}, \Sigma_k^{\text{old}})}{\sum_{j=1}^K \pi_j^{\text{old}} \mathcal{N}(y_n \mid \mu_j^{\text{old}}, \Sigma_j^{\text{old}})} \end{aligned}$$

Commonly refer to $q(z_n)$ as *responsibility*

Example: Gaussian Mixture Model



$$\log p(\mathcal{Y} \mid \pi, \mu, \Sigma) \geq \sum_{n=1}^N \sum_{k=1}^K q(z_n = k) \log \{ \pi_k \mathcal{N}(y_n \mid \mu_k, \Sigma_k) \} = \mathcal{L}(q, \theta)$$

M-Step: $\theta^{\text{new}} = \arg \max_{\theta} \mathcal{L}(q^{\text{new}}, \theta)$

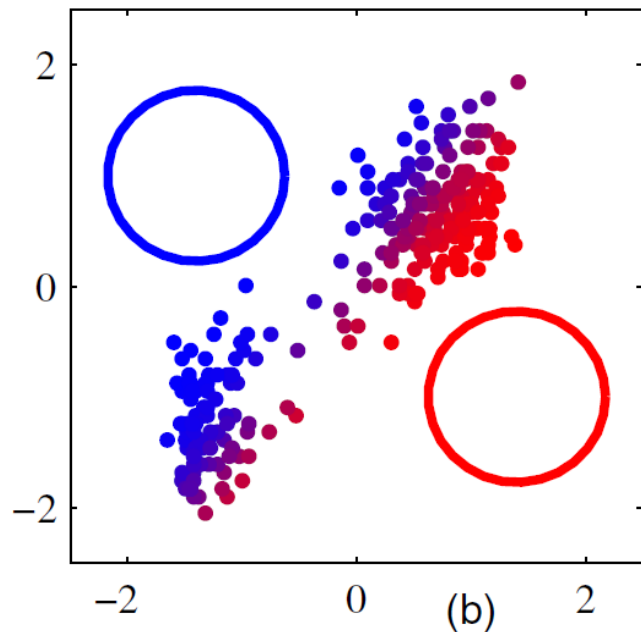
Start with mean parameter μ_k ,

$$0 = \nabla_{\mu_k} \mathcal{L}(q^{\text{new}}, \theta)$$

$$0 = \sum_{n=1}^N \nabla_{\mu_k} \mathbf{E}_{z_n \sim q^{\text{new}}} [\log \mathcal{N}(y_n \mid \mu_{z_n}, \Sigma_{z_n})]$$

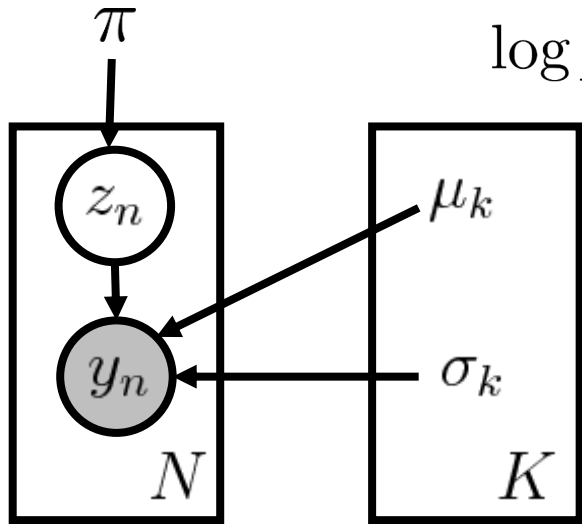
$$0 = - \sum_{n=1}^N q^{\text{new}}(z_n = k) \Sigma_k (y_n - \mu_k)$$

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N q^{\text{new}}(z_n = k) y_n \quad \text{where} \quad N_k = \sum_{n=1}^N q(z_n = k)$$



Example: Gaussian Mixture Model

$$\log p(\mathcal{Y} \mid \pi, \mu, \Sigma) \geq \sum_{n=1}^N \sum_{k=1}^K q(z_n = k) \log \{ \pi_k \mathcal{N}(y_n \mid \mu_k, \Sigma_k) \} = \mathcal{L}(q, \theta)$$

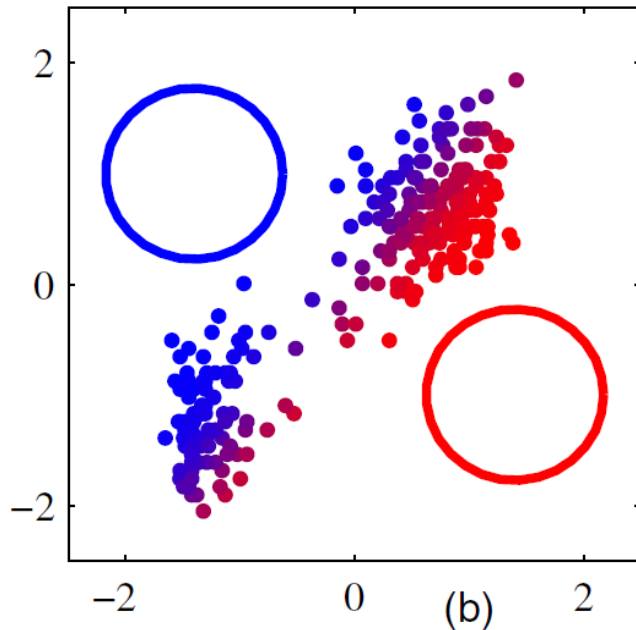


M-Step: $\theta^{\text{new}} = \arg \max_{\theta} \mathcal{L}(q^{\text{new}}, \theta)$

Repeat for remaining parameters,

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N q(z_n = k) (y_n - \mu_k^{\text{new}})(y_n - \mu_k^{\text{new}})^T$$

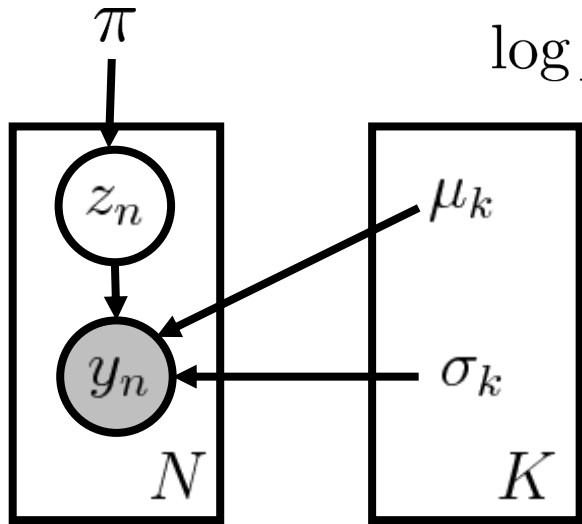
$$\pi_k^{\text{new}} = \frac{N_k}{N}$$



- Solving for mixture weights requires a bit more work
- Need constraint $\sum_k \pi_k = 1$
- Use Lagrange multiplier approach

Example: Gaussian Mixture Model

$$\log p(\mathcal{Y} \mid \pi, \mu, \Sigma) \geq \sum_{n=1}^N \sum_{k=1}^K q(z_n = k) \log \{ \pi_k \mathcal{N}(y_n \mid \mu_k, \Sigma_k) \} = \mathcal{L}(q, \theta)$$

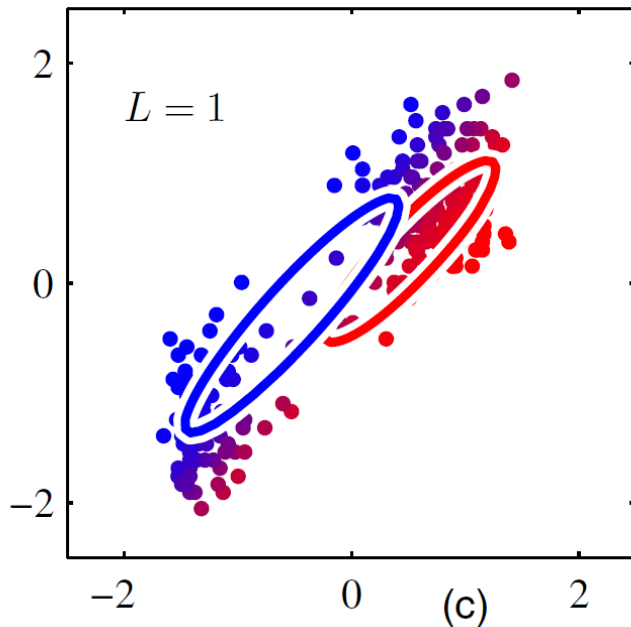


M-Step: $\theta^{\text{new}} = \arg \max_{\theta} \mathcal{L}(q^{\text{new}}, \theta)$

Repeat for remaining parameters,

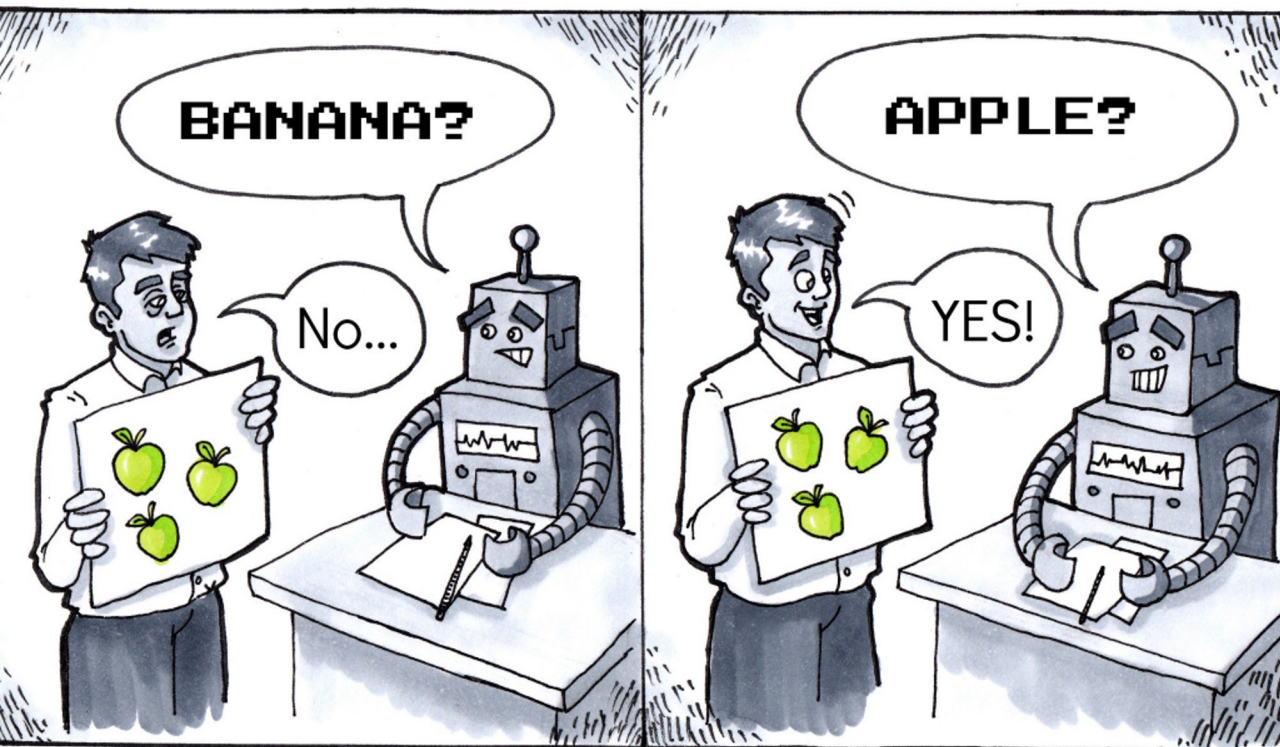
$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N q(z_n = k) (y_n - \mu_k^{\text{new}})(y_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

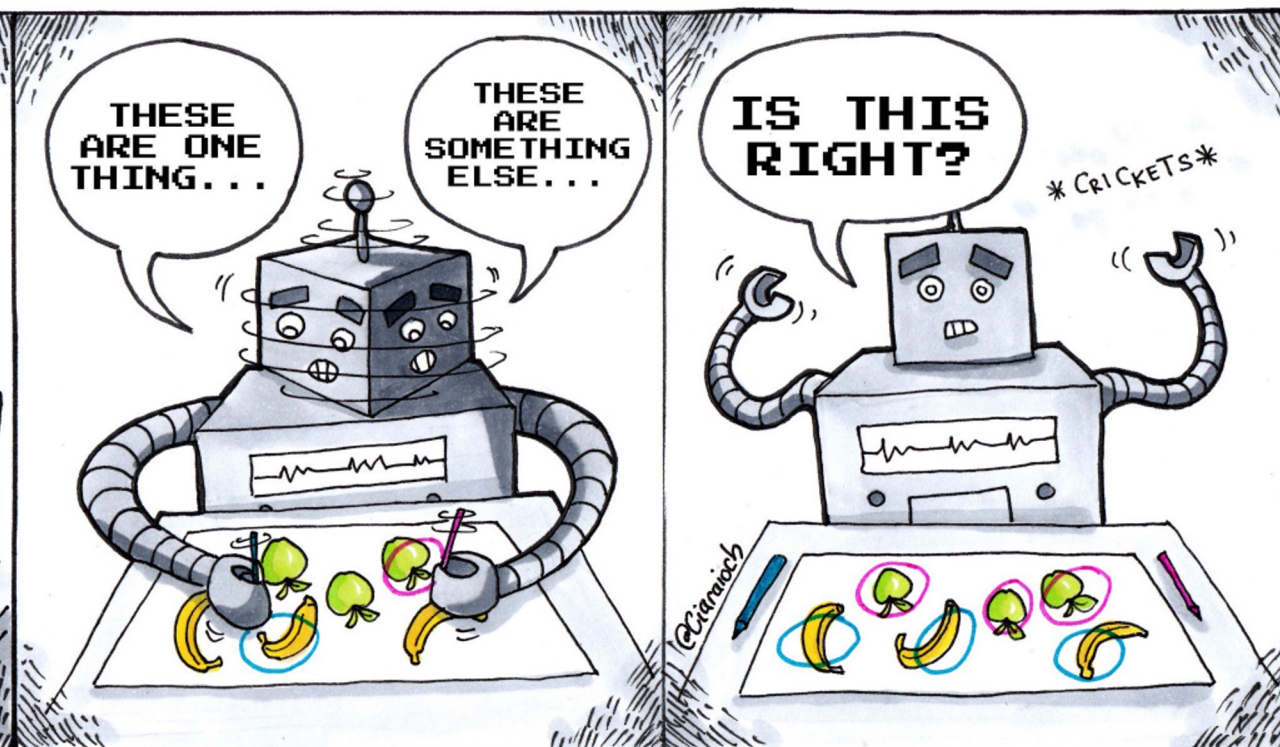


- Solving for mixture weights requires a bit more work
- Need constraint $\sum_k \pi_k = 1$
- Use Lagrange multiplier approach

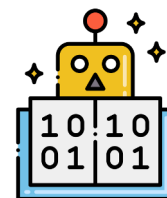
K-Means Clustering



Supervised Learning



Unsupervised Learning

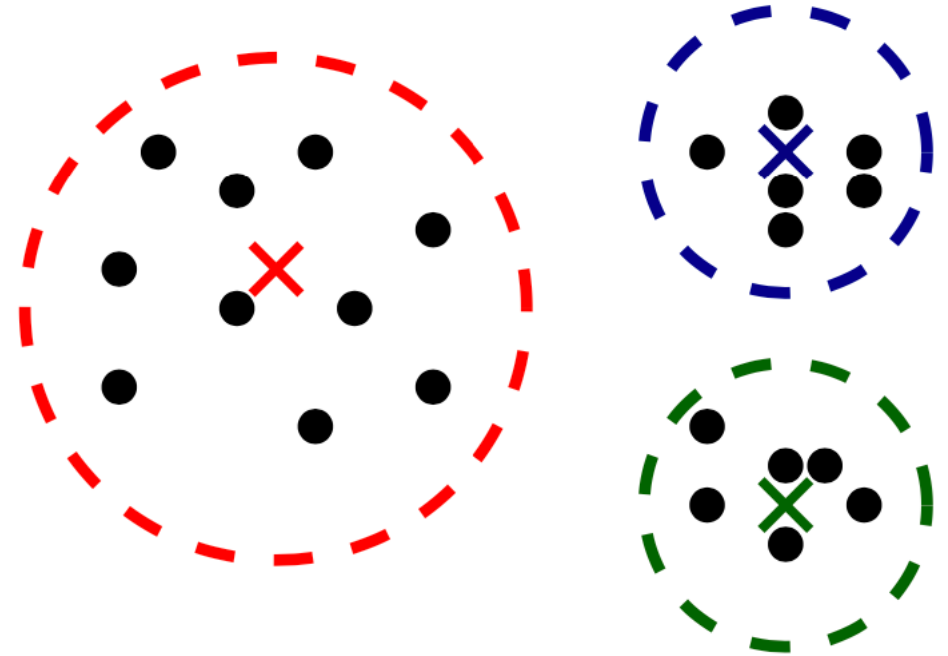


Clustering

- Input: k : the number of clusters (hyperparameter)

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

- Output
 - partition $\{G_i\}_{i=1}^k$ s.t. $S = \cup_i G_i$ (disjoint union).
 - often, we also obtain 'centroids'
- Q: what would be a reasonable definition of centroids?



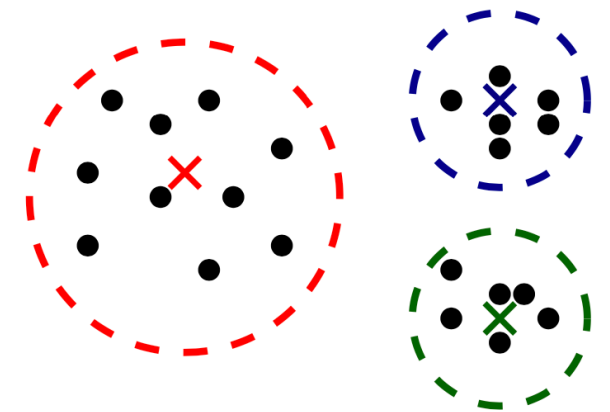
k -means clustering

- Idea: to partition the data, it would be great if someone gives us k reasonable centroids c_1, \dots, c_k , since then we can partition the data with them.

$$A(x) = \arg \min_{j \in [k]} \|x - c_j\|_2$$

- But we don't have those centroids => Let's find them with an optimization formulation.

$$\underset{c_1, \dots, c_k}{\text{minimize}} f(c_1, \dots, c_k), \text{ where } f(c_1, \dots, c_k) = \sum_{i=1}^n \min_{j \in [k]} \|x - c_j\|_2^2$$



Special case: $k=1$

- $\min_{c_1, \dots, c_k} \sum_{i=1}^n \min_{j \in [k]} \|x_i - c_j\|_2^2 \Rightarrow \min_c \sum_{i=1}^n \|x_i - c\|_2^2$
- Let $F(c) = \sum_{i=1}^n \|x_i - c\|_2^2$ convex; minimizer c^* satisfies that $\nabla F(c^*) = 0$
 $\Rightarrow \sum_{i=1}^n (x_i - c^*) = 0 \Rightarrow c^* = \frac{1}{n} \sum_{i=1}^n x_i$

For $k \geq 2$

- minimize $f(c_1, \dots, c_k)$, where $f(c_1, \dots, c_k) = \sum_{i=1}^n \min_{j \in [k]} \|x - c_j\|_2^2 \Rightarrow$ NP-hard even when $d = 2$

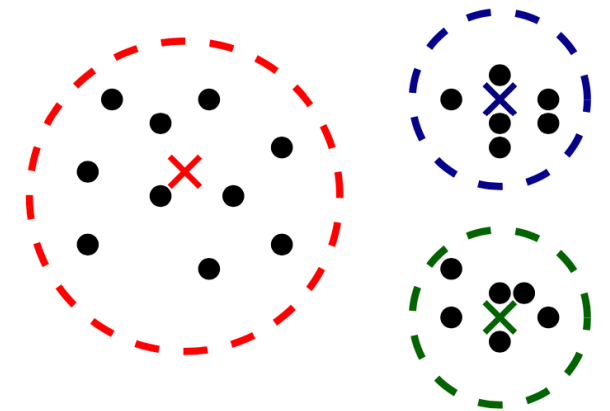
- **K-means algorithm**: solve it approximately (heuristic)

(Also called Lloyd's algorithm)

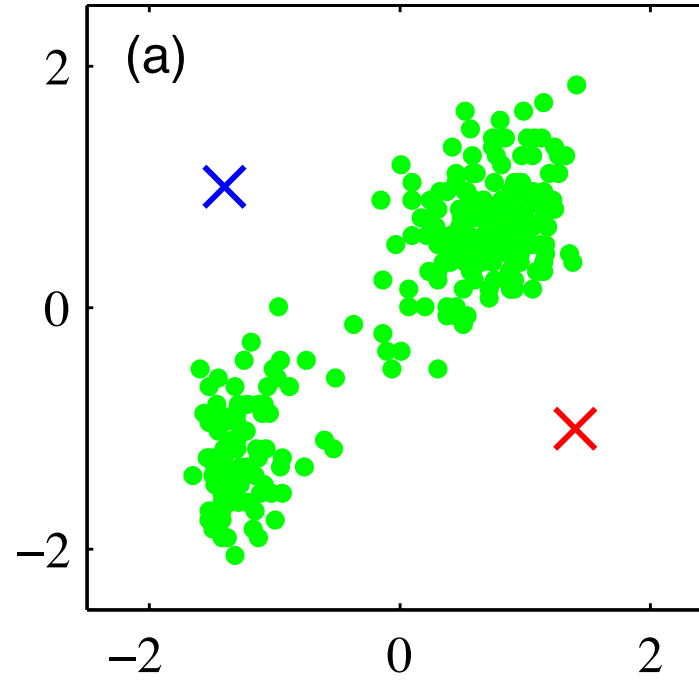
- Observation: The chicken-and-egg problem.

- Cluster center location depends on the cluster assignment
- Cluster assignment depends on cluster location

- Very common heuristic (that may or may not be the best thing to do)

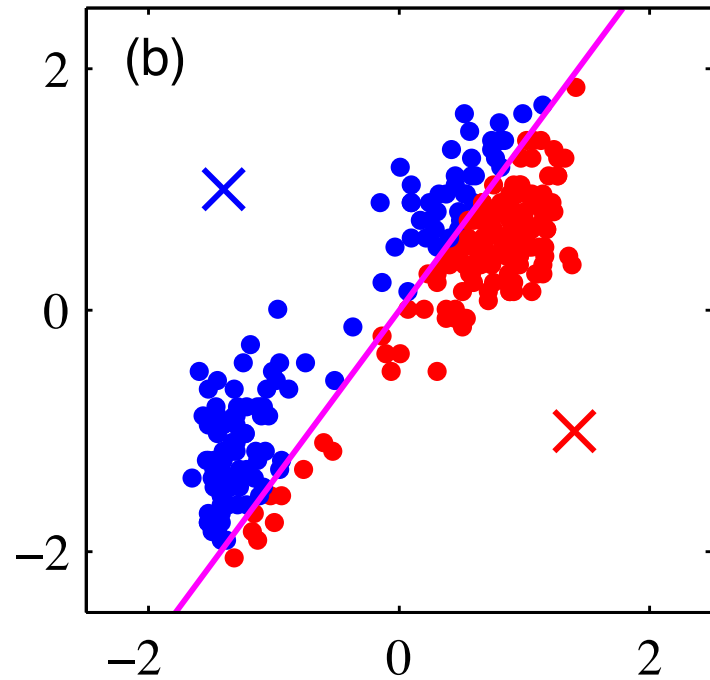


Initialization

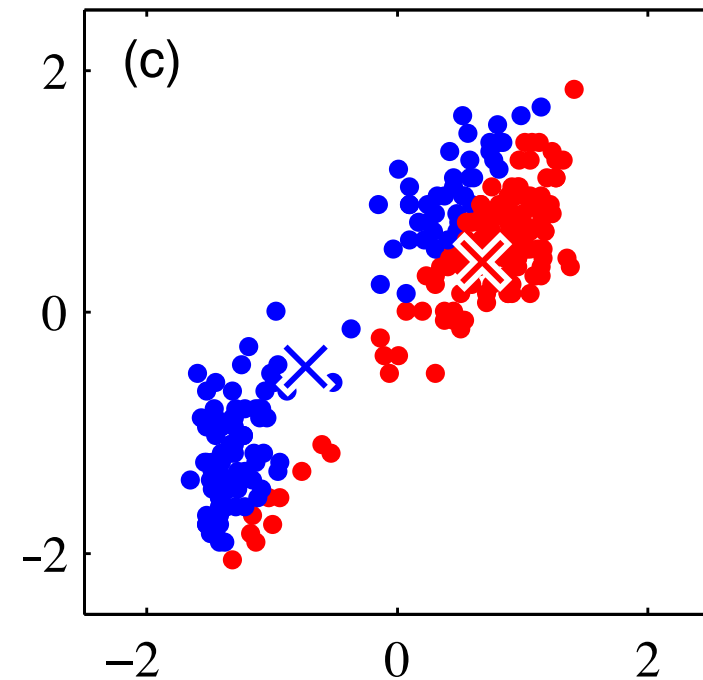


Arbitrary/random initialization of c_1 and c_2

Iteration 1

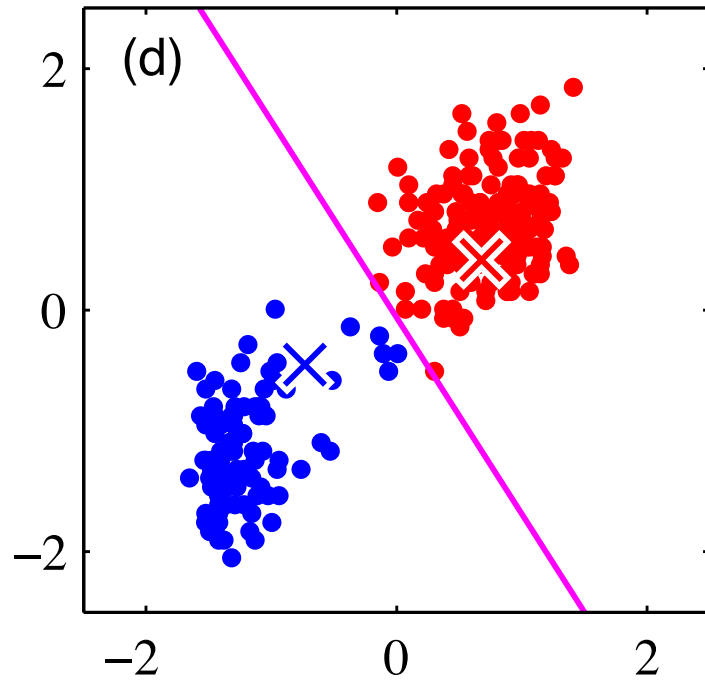


(A) update the cluster assignments.

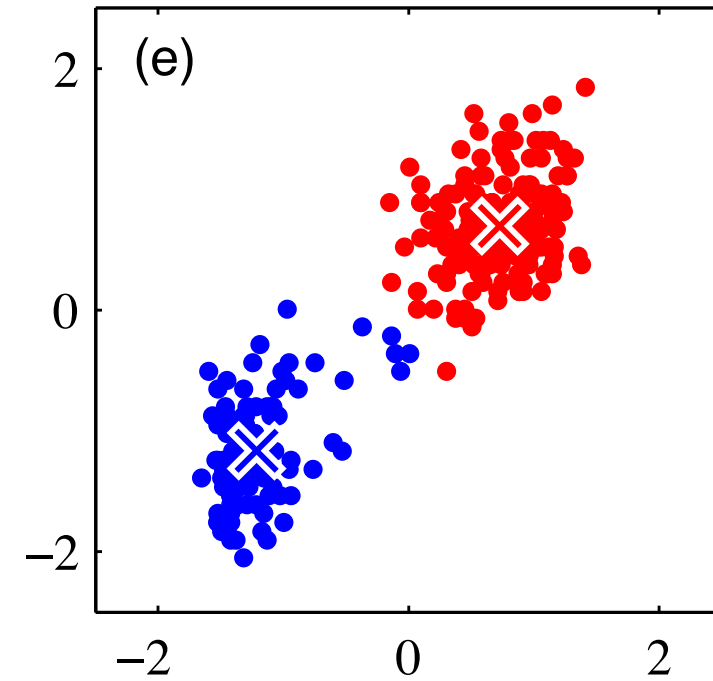


(B) Update the centroids $\{c_j\}$

Iteration 2

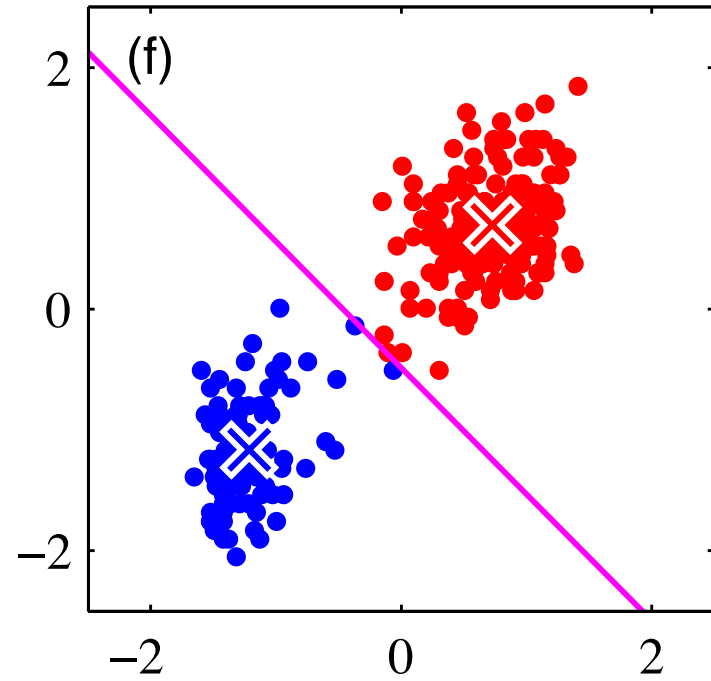


(A) update the cluster assignments.

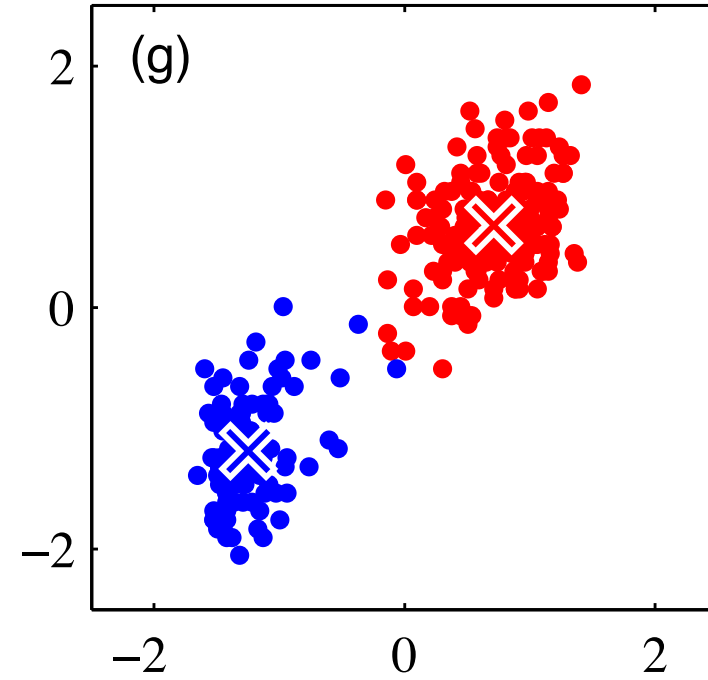


(B) Update the centroids $\{c_j\}$

Iteration 3

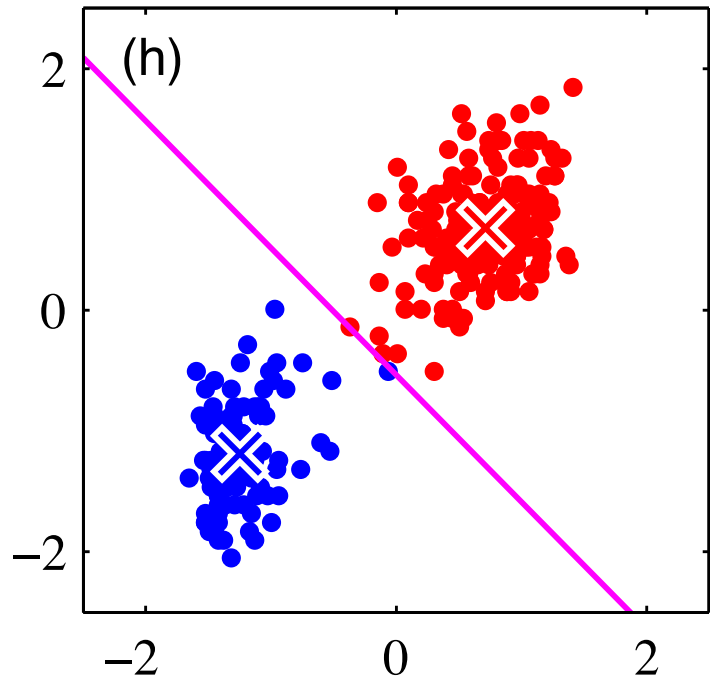


(A) update the cluster assignments.

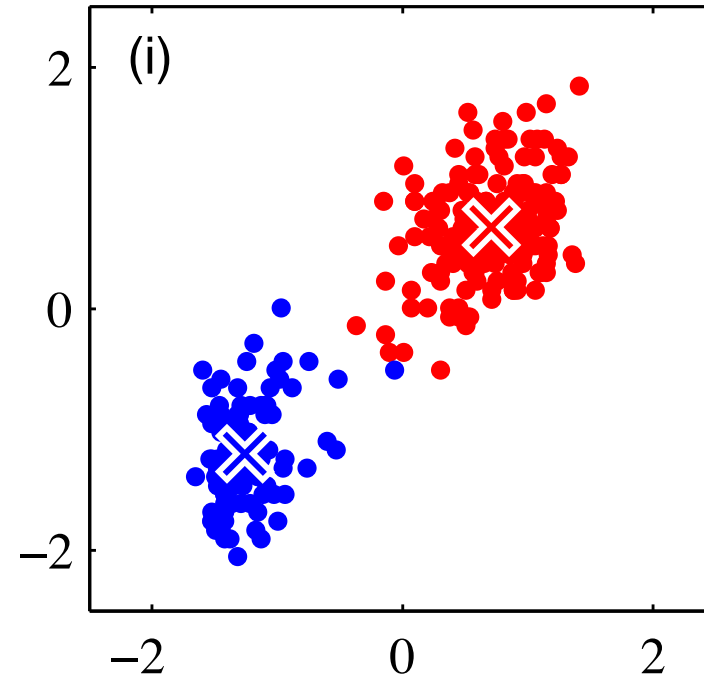


(B) Update the centroids $\{c_j\}$

Iteration 4



(A) update the cluster assignments.



(B) Update the centroids $\{c_j\}$

K-means clustering

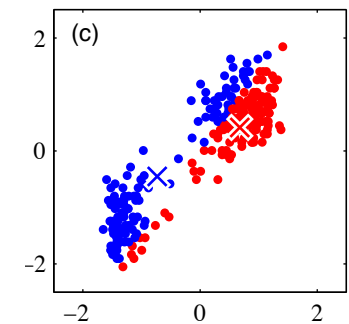
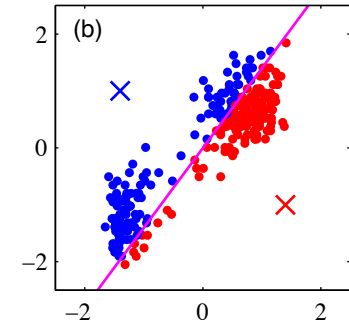
Input: k : num. of clusters, $S = \{x_1, \dots, x_n\}$

[Initialize] Pick c_1, \dots, c_k as randomly selected points from S (see next slides for alternatives)

For $t=1,2,\dots,\text{max_iter}$

- **[Assignments]** $\forall x \in S, a_t(x) = \arg \min_{j \in [k]} \|x - c_j\|_2^2$
- If $t \neq 1$ AND $a_t(x) = a_{t-1}(x), \forall x \in S$
 - break
- **[Centroids]** $\forall j \in [k], c_j \leftarrow \text{average}(\{x \in S: a_t(x) = j\})$

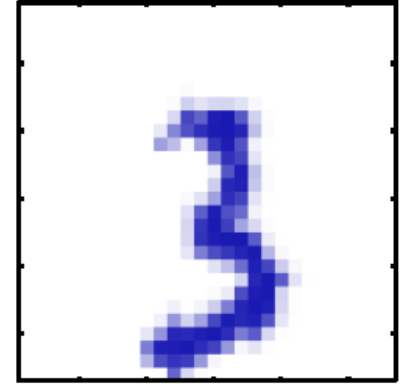
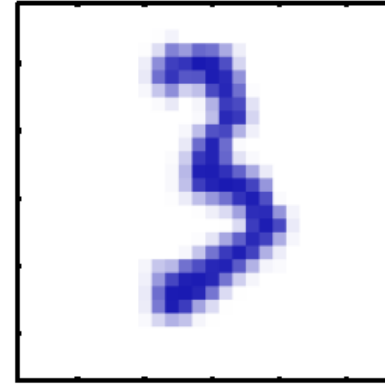
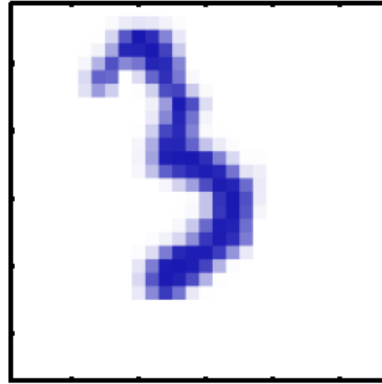
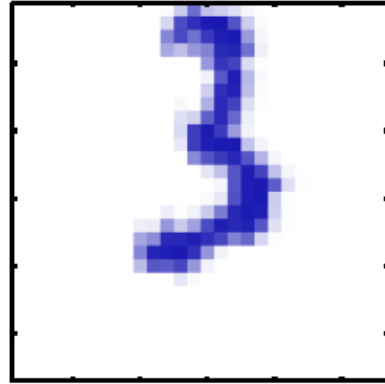
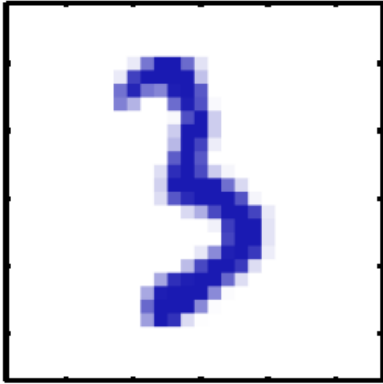
Output: c_1, \dots, c_k and $\{a_t(x_i)\}_{i \in [n]}$



Dimensionality Reduction and Principal Component Analysis (PCA)

Motivation

Data often have a lot of redundant information...



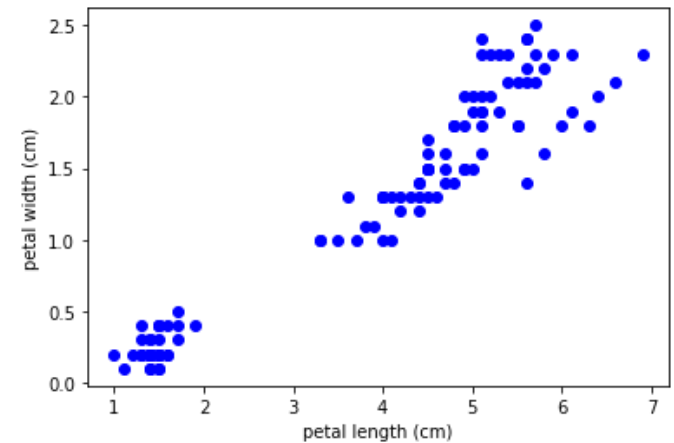
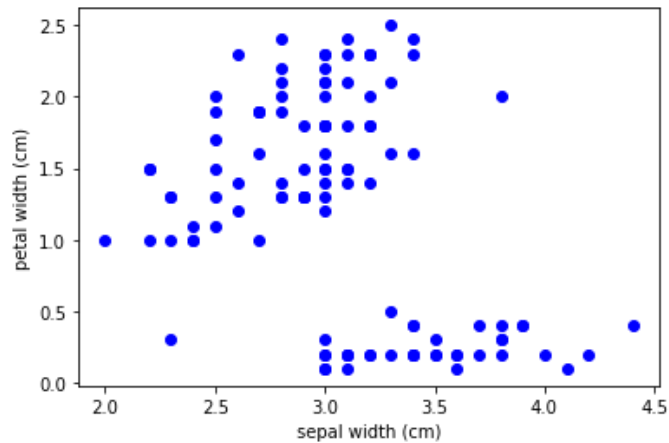
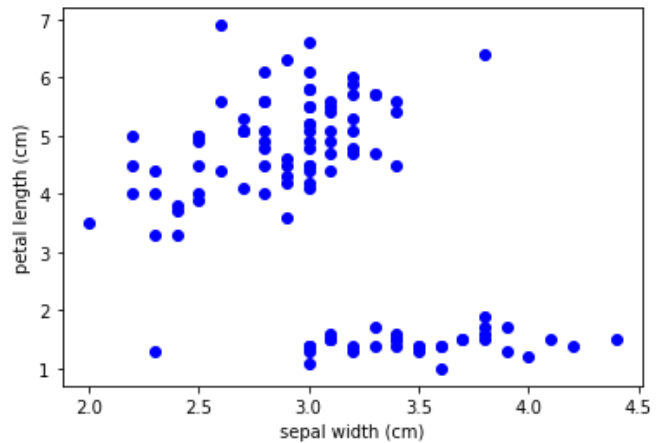
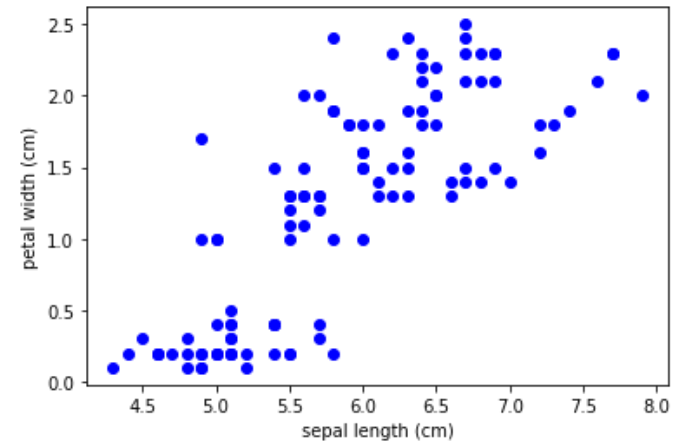
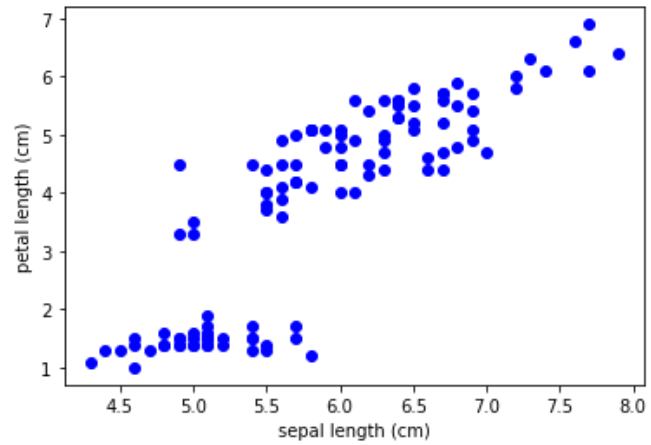
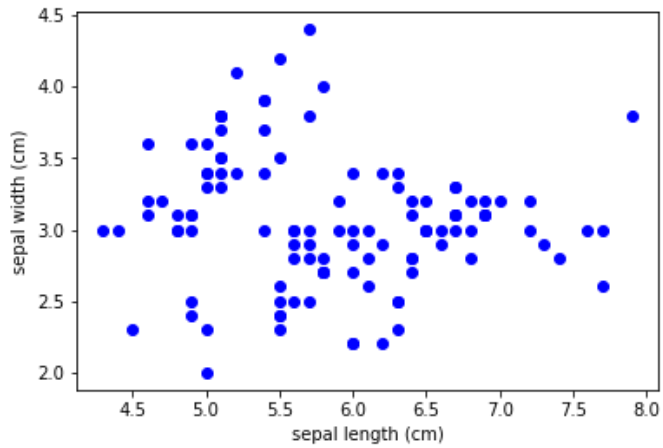
Example A dataset consisting of a hand-drawn 3 at random locations and rotations in a 100x100 pixel image.

Data Dimension $100 \times 100 = 10,000$

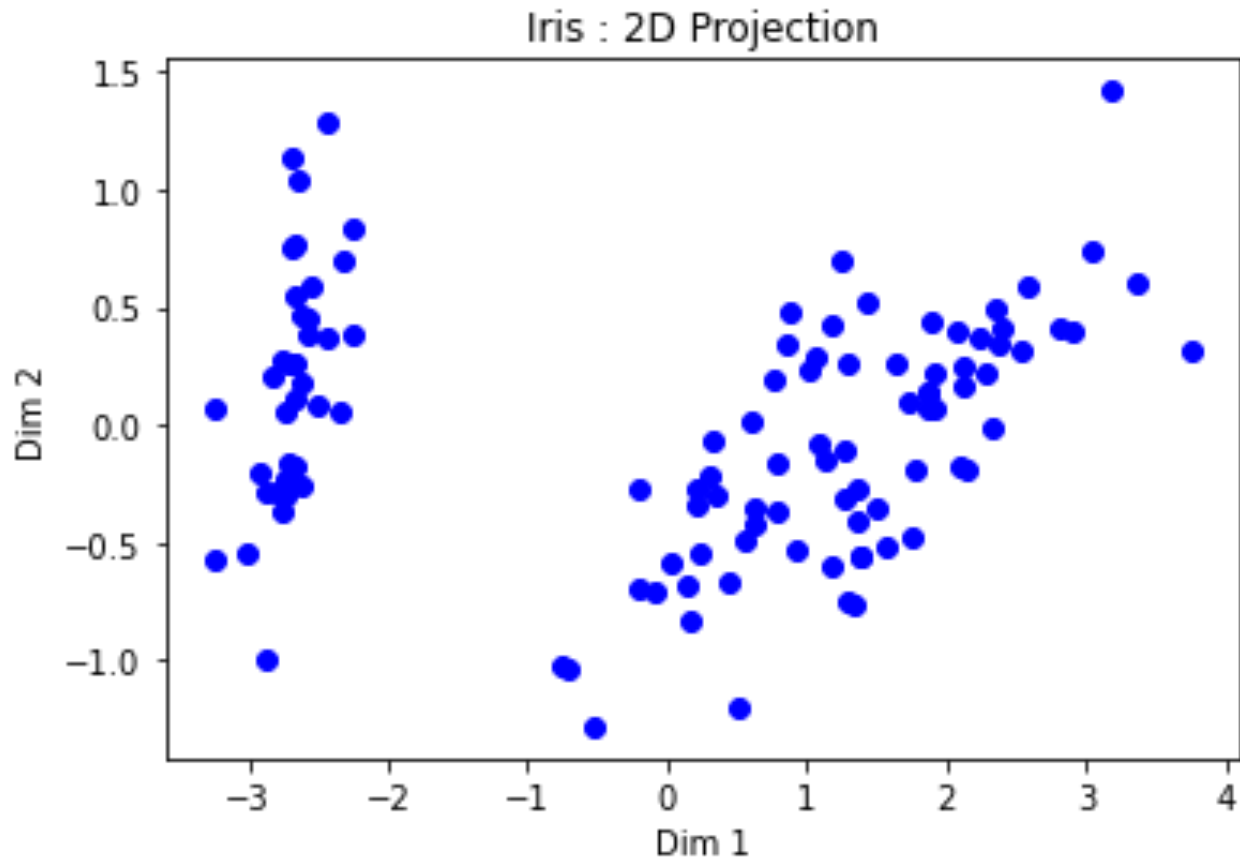
Intrinsic Dimension 3 (X-position, Y-position, Rotation)

Example : Iris Dataset

*Recall that the Iris dataset has 4 features:
sepal length / width, petal length / width...*



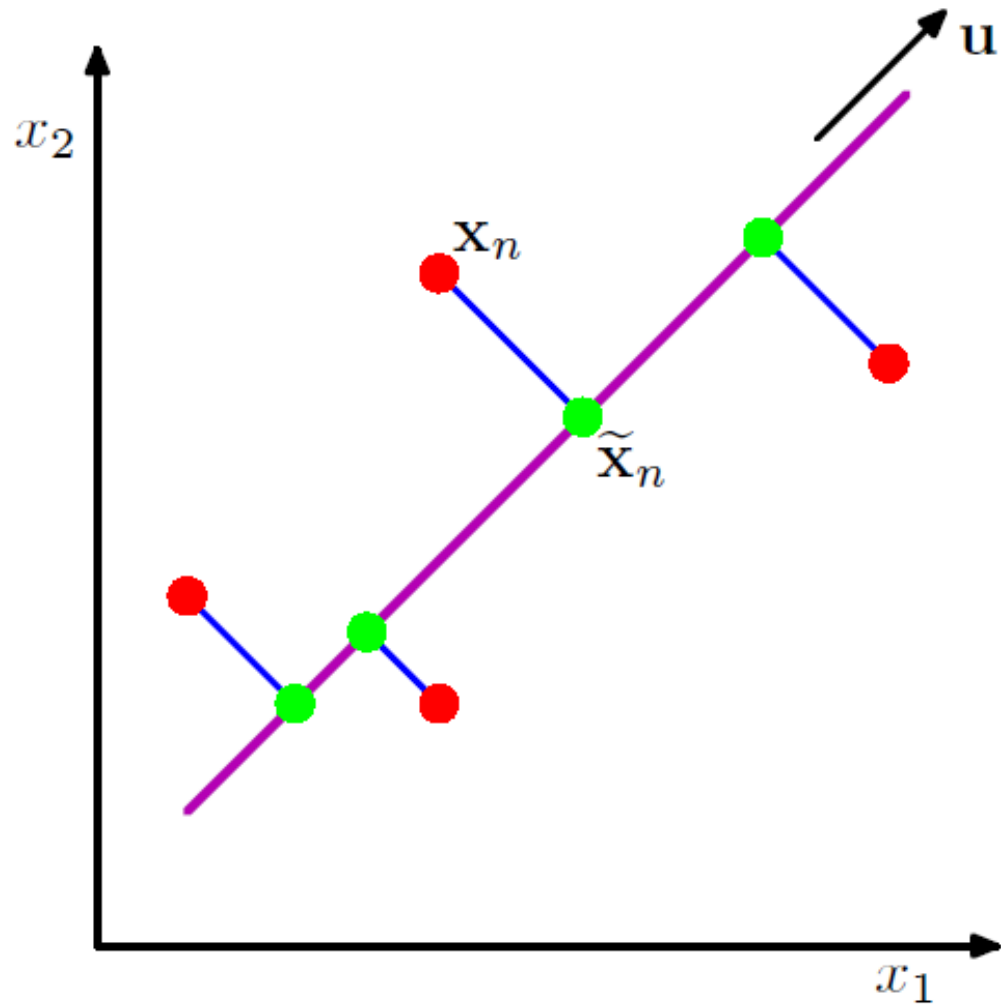
Example : Iris Dataset



Data still cluster in a two-dimensional subspace

We can fit model in 2D to reduce complexity, visualize results, etc.

Linear Dimensionality Reduction

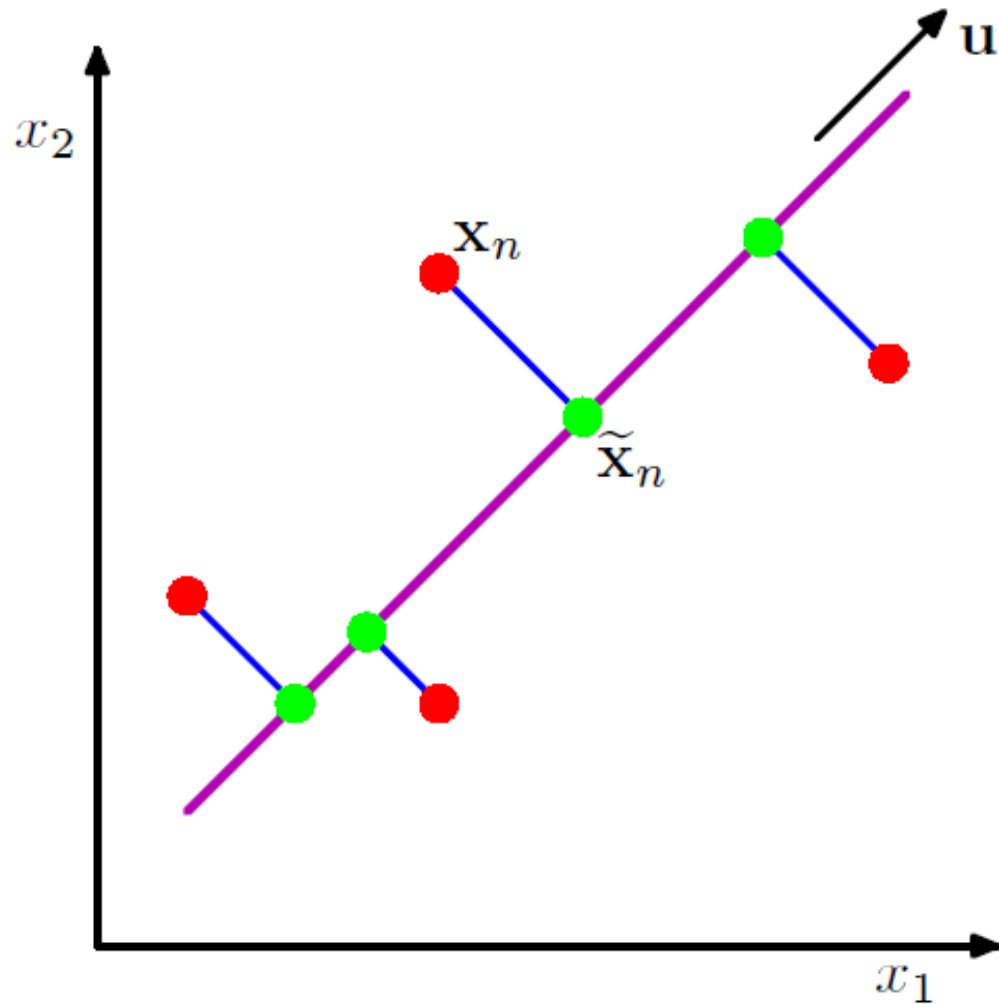


Project data onto a line or plane...

...one of the simplest dimensionality reduction approaches

First, let's review some linear algebra...

Linear Dimensionality Reduction



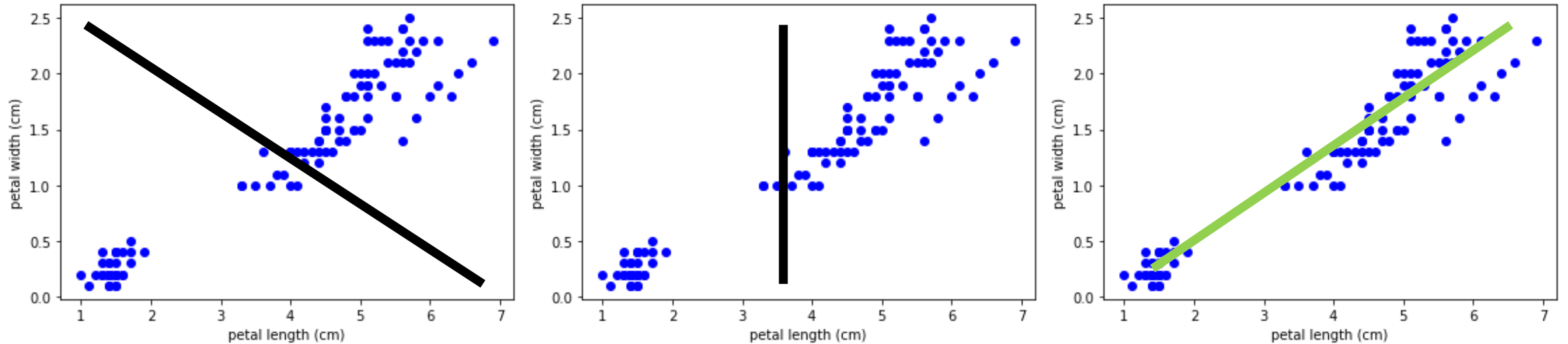
Projecting data onto a vector is a simple inner product,

$$\tilde{x}_n = u^T x_n$$

We call u the linear subspace

Linear Dimensionality Reduction

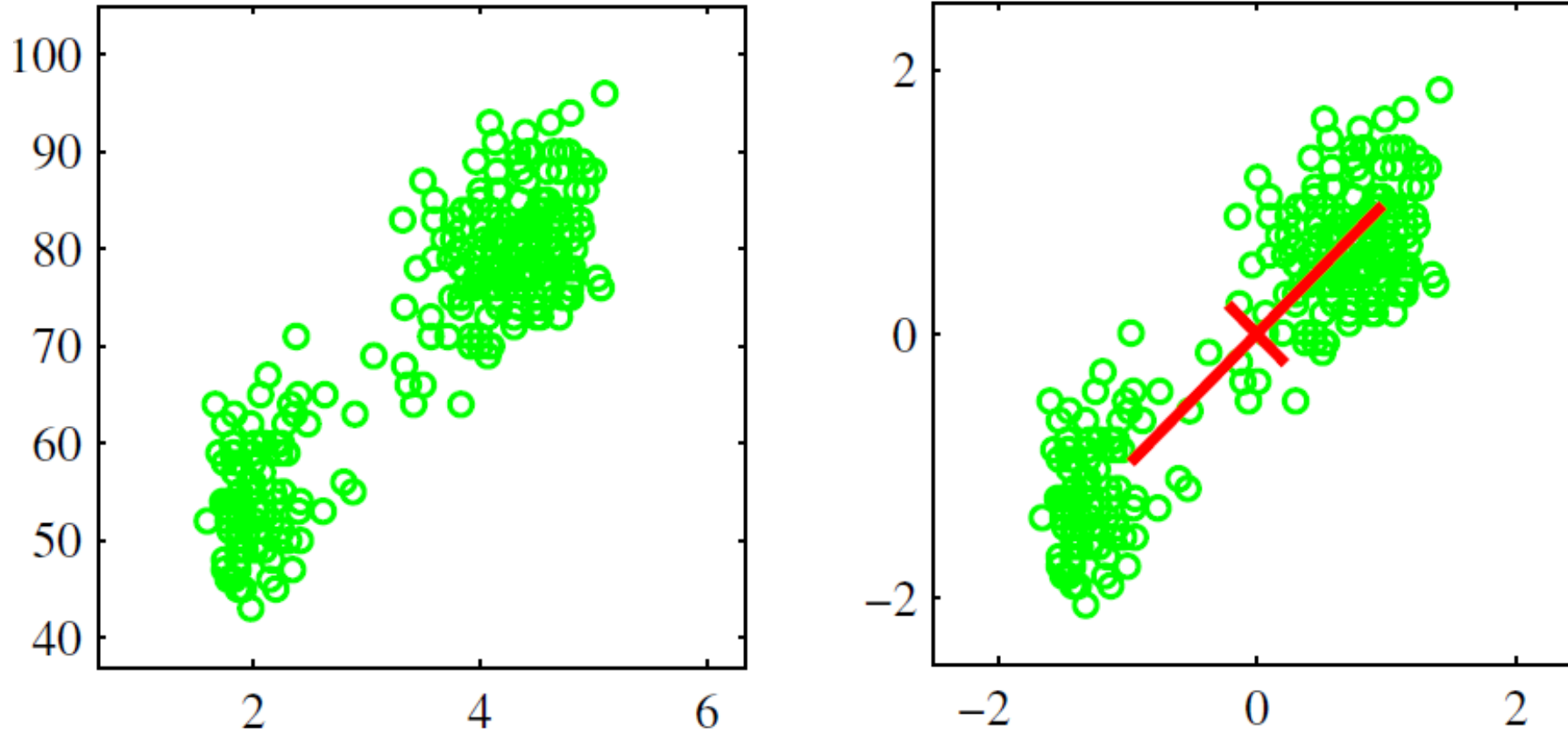
Which choice of subspace is best? And why?



Idea Choose the subspace that captures the most variation in the original data

Principal Component Analysis (PCA)

Identify directions of *maximum variation* as subspaces...



...we call each direction a *principal component*

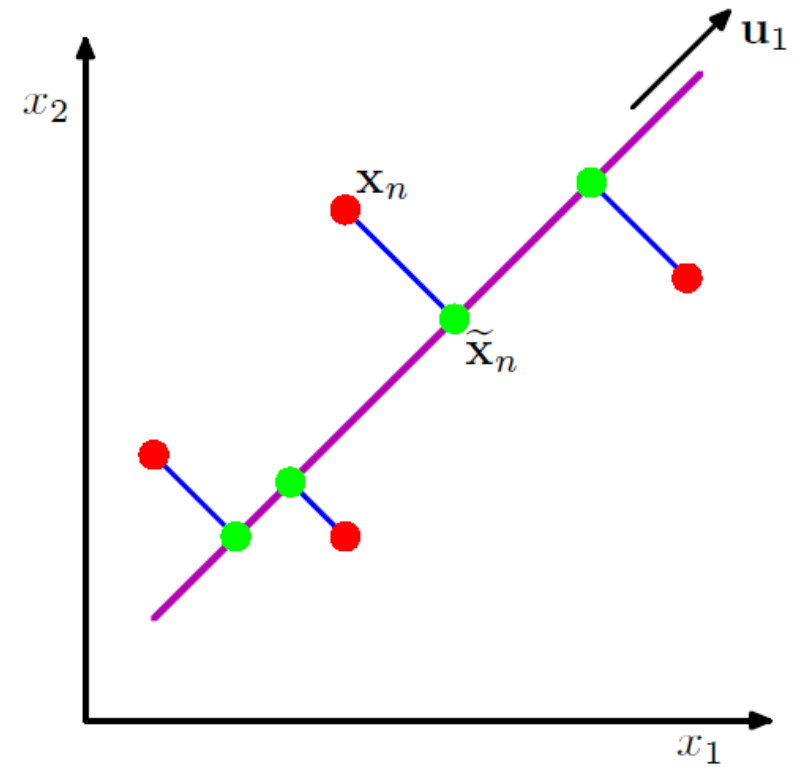
Principal Component Analysis (PCA)

First, center the data by subtracting the sample mean,

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

Variance of projected subspace,

$$\frac{1}{N} \sum_{n=1}^N \left(\underbrace{u^T x_n}_{\text{Projection of } n^{\text{th}} \text{ data point}} - \underbrace{u^T \bar{x}}_{\text{Projection of mean}} \right)^2$$




Maximum Variance Formulation

A little algebra...

$$\frac{1}{N} \sum_{n=1}^N (u^T x_n - u^T \bar{x})^2 = \frac{1}{N} \sum_{n=1}^N \{u^T (x_n - \bar{x})\}^2 \quad \text{Pull out } u$$

$$\text{Quadratic form} = \frac{1}{N} \sum_{n=1}^N u^T (x_n - \bar{x})(x_n - \bar{x})^T u$$

Define: $S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$

Then: $\frac{1}{N} \sum_{n=1}^N (u^T x_n - u^T \bar{x})^2 = u^T S u$  This is what we will optimize over u

Maximum Variance Formulation

Find u so that projected variance is maximal...


$$\max_u u^T S u$$

Don't want to *cheat* with large magnitude u , so we add penalty,

$$\max_u u^T S u - \lambda u^T u$$

Set the derivative (gradient) to zero and solve...

$$S u - \lambda u = 0$$

$$S u = \lambda u$$


What equation is this?

**u is an *eigenvector* with
eigenvalue λ**

Recap of Concepts

- Learning a reduced *intrinsic dimension* is useful for a bunch of reasons
- The easiest approach is to find a *linear subspace*
- PCA defines the linear subspace as that which maximizes variance of the projected data
- The set of subspaces are defined by the *eigenvectors*,

$$\max_u u^T S u - \lambda u^T u$$

$$S u = \lambda u$$

Eigenstuff

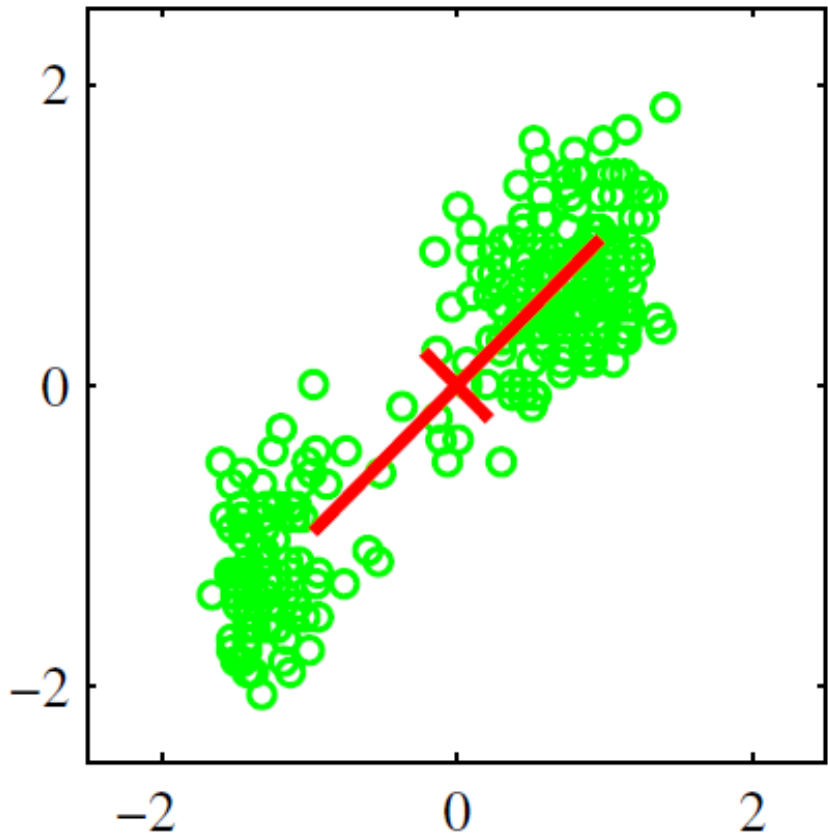
Eigenvectors / values of a matrix solve the equation

$$Su = \lambda u$$

- Matrix S may have *multiple* eigenvectors / values that solve the above equation
- For D -dimensional u can find all vectors in $O(D^3)$ time
- PCA finds $M < D$ vectors with largest eigenvalues
- Can find $M < D$ sorted eigenvectors in $O(MD^2)$ time
- Note that D can be large!

Principal Component Analysis (PCA)

How much variance is captured by just the first principal component (i.e. eigenvector with largest eigenvalue)?



Let u_1 be the first principal component, then variance of first PC is,

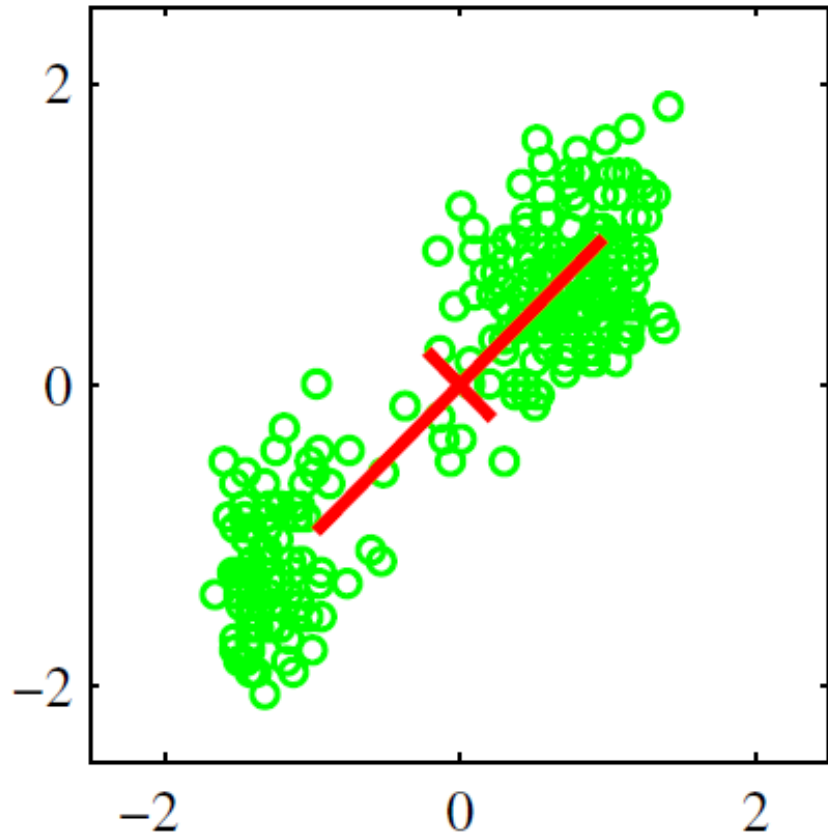
$$\frac{1}{N} \sum_n \{u_1^T (x_n - \bar{x})\}^2$$

How much in the second PC?

$$\frac{1}{N} \sum_n \{u_2^T (x_n - \bar{x})\}^2$$

Explained Variance

How much variance is captured in $M < D$ principal components?



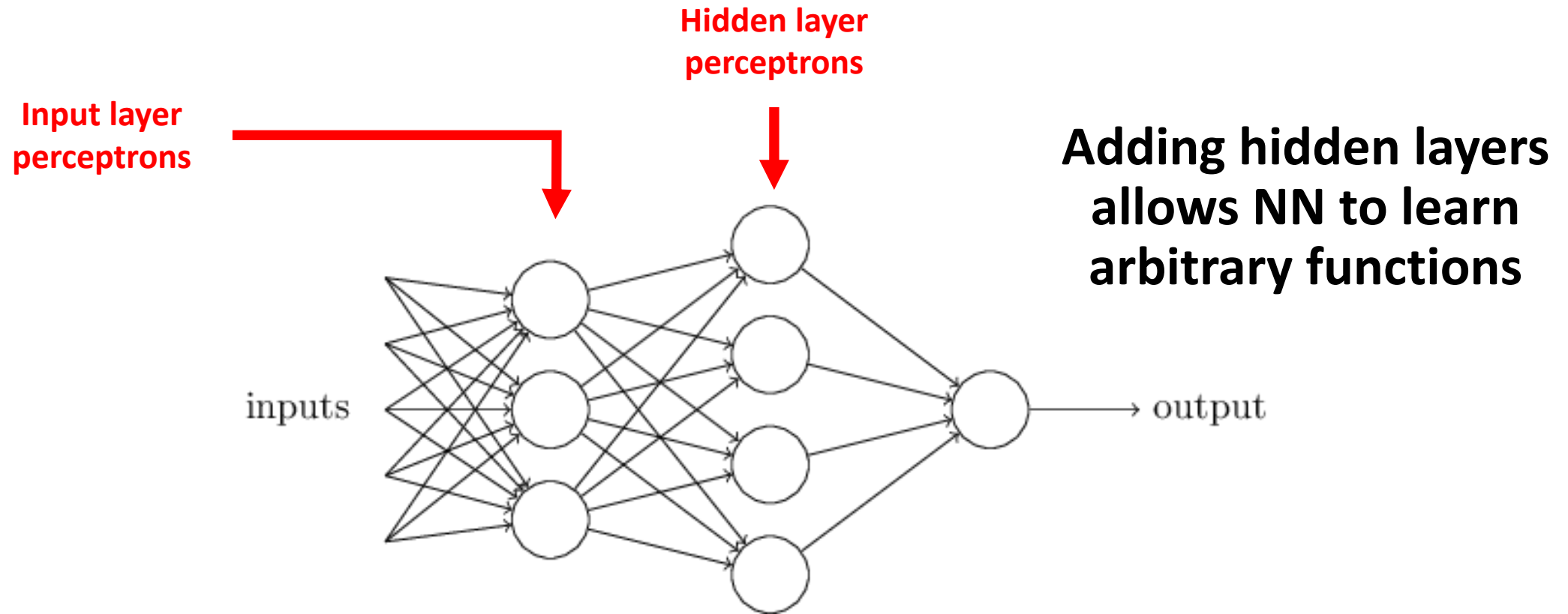
$$\frac{1}{N} \sum_{m=1}^M \sum_n \left\{ u_m^T (x_n - \bar{x}) \right\}^2$$

We call this the *explained variance* of the first M principal components

Divide by total variance to find percentage of the total variance explained by the subspace

Fully Connected Neural Networks

Multilayer Perceptron



This is the quintessential *Neural Network*...

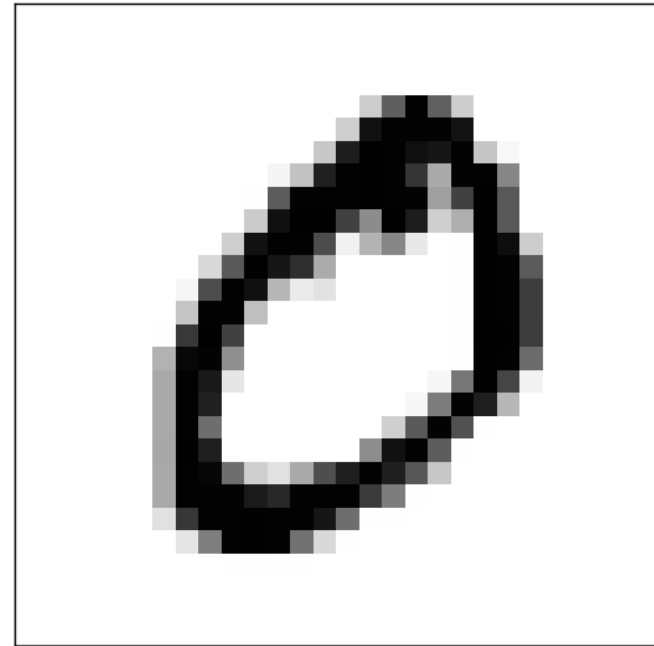
...also called *Feed Forward Neural Net* or *Artificial Neural Net*

Handwritten Digit Classification

Classifying handwritten digits is the “Hello World” of NNs



Each character is centered in a
28x28=784 pixel grayscale
image

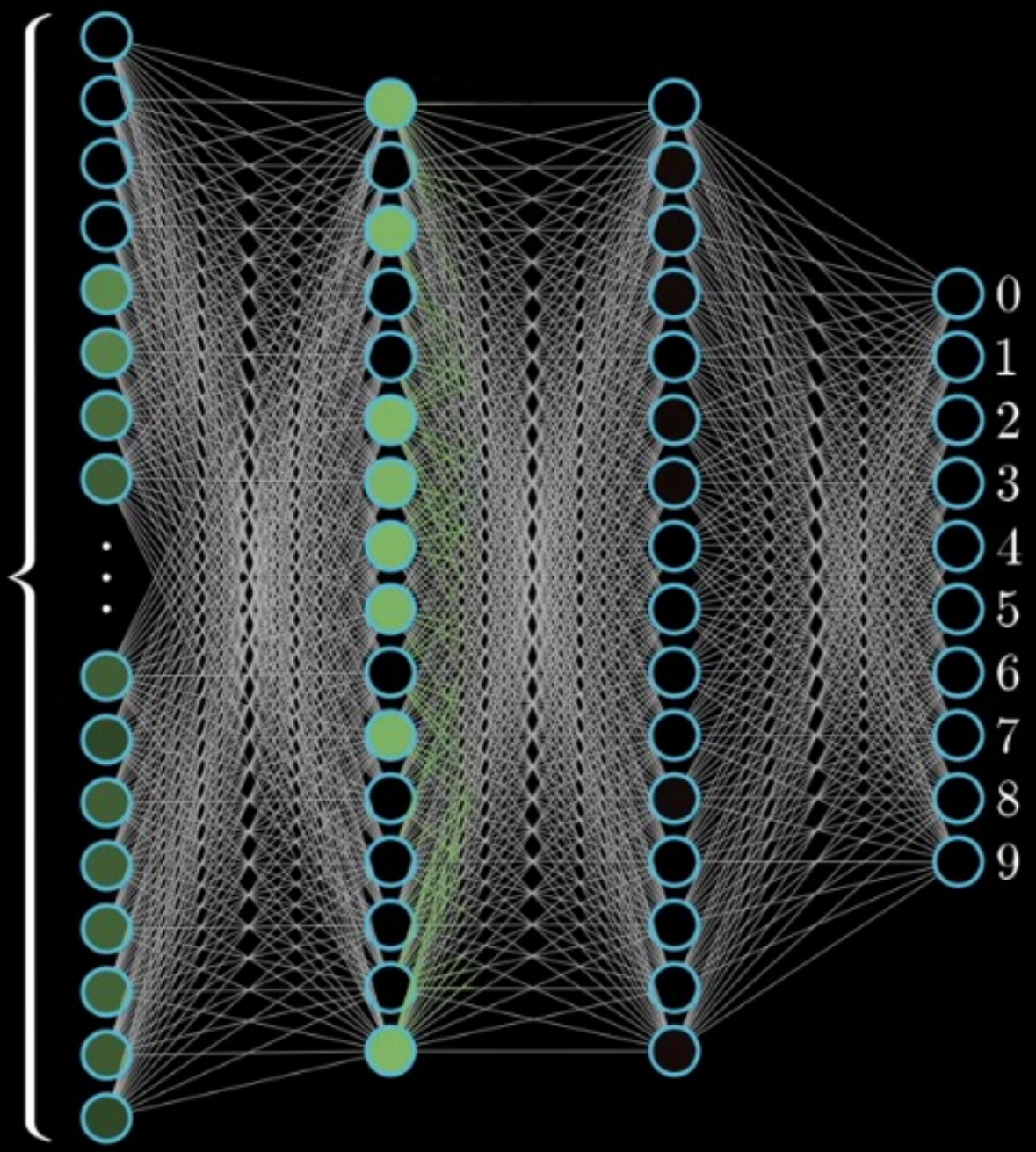


Modified National Institute of
Standards and Technology (MNIST)
database contains 60k training and
10k test images

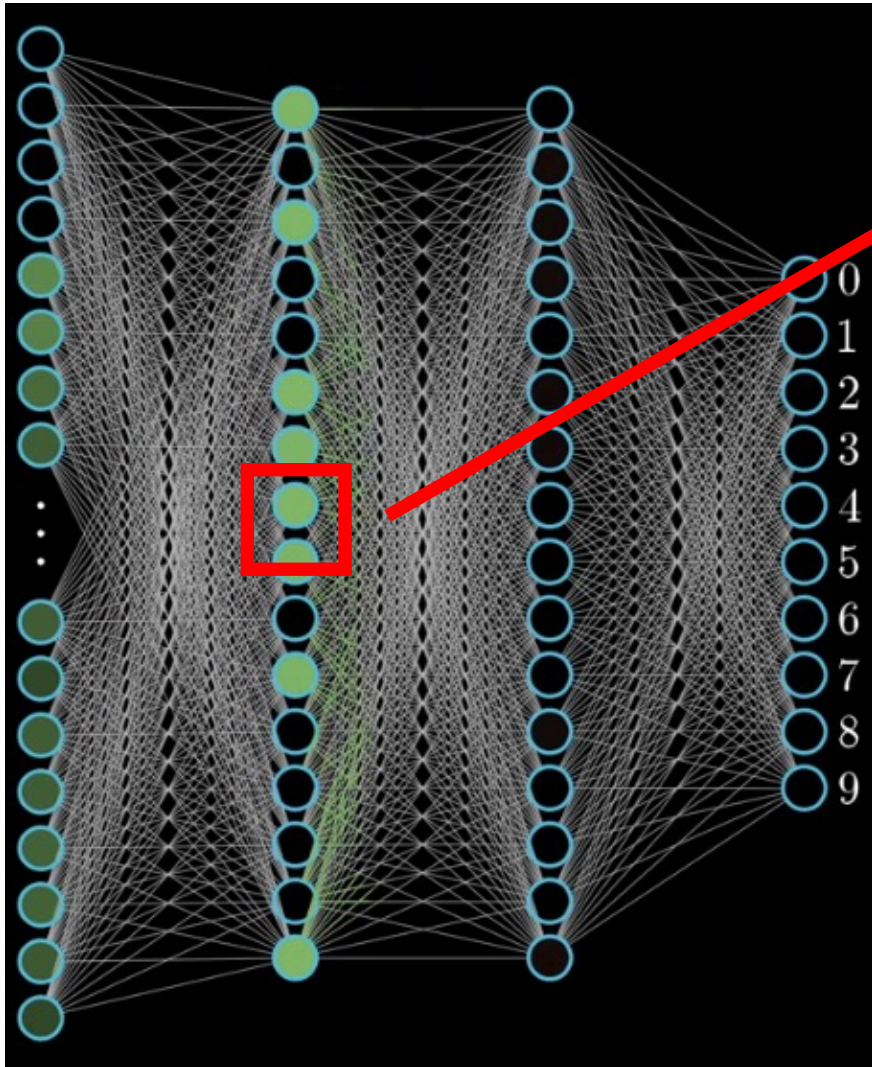


784

Each image pixel is a number in [0,1] indicated by highlighted color



Feedforward Procedure



Each node computes a *weighted combination* of nodes at the previous layer...

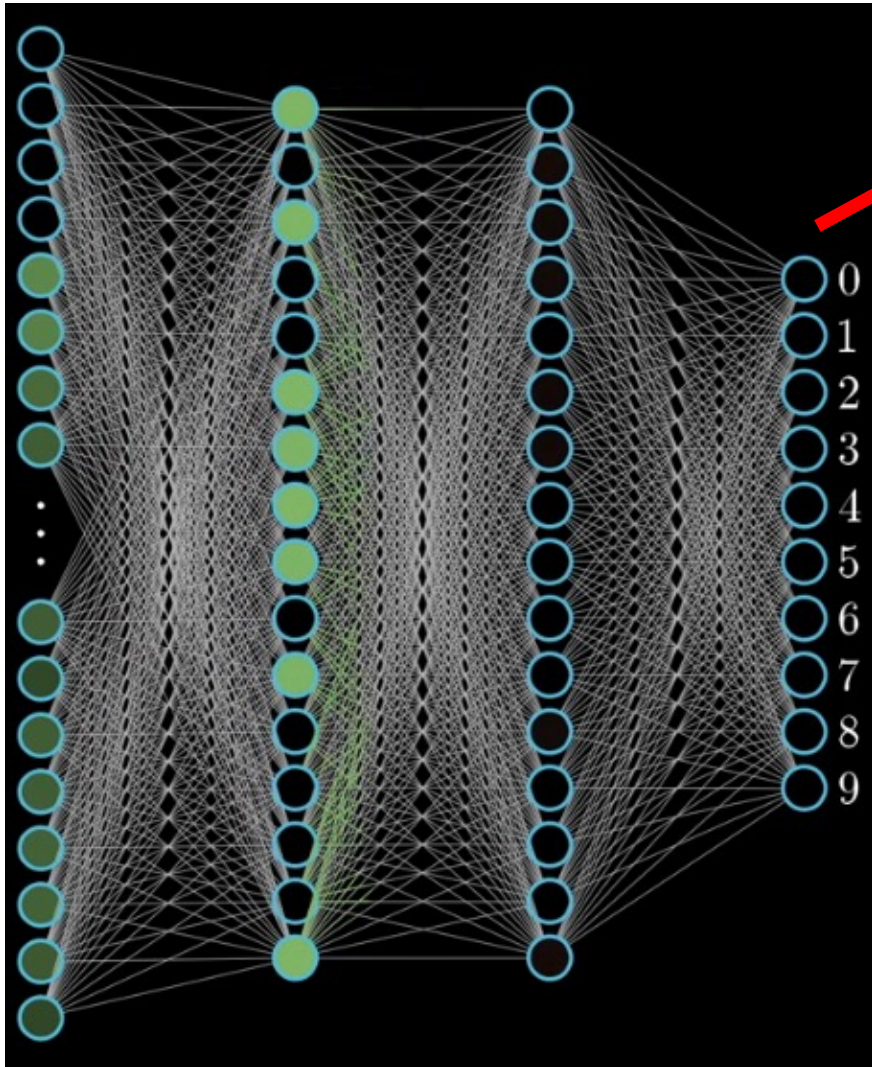
$$w_1x_1 + w_2x_2 + \dots + w_nx_n$$

Then applies a *nonlinear function* to the result

$$\sigma(w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$

Often, we also introduce a constant *bias* parameter

Multilayer Perceptron



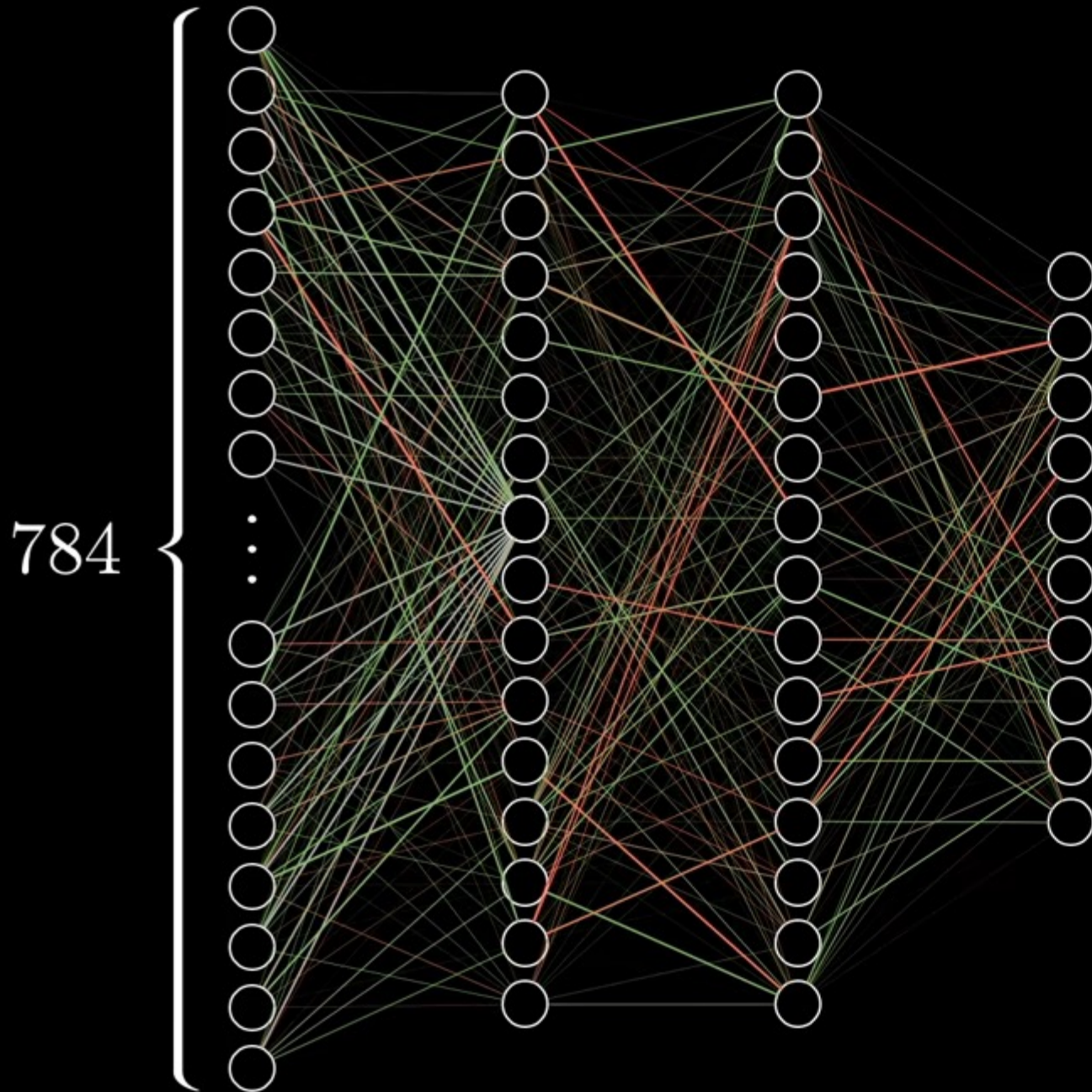
Final layer is typically a linear model...for classification this is a Logistic Regression

$$\sigma(w^T x + b) = \frac{1}{1 + e^{-(w^T x + b)}}$$

Vector of activations from previous layer

Recall that for multiclass logistic regression with K classes,

$$p(\text{Class} = k \mid x) \propto \sigma(w_k^T x + b_k)$$



$$784 \times 16 + 16 \times 16 + 16 \times 10$$

weights

$$16 + 16 + 10$$

biases

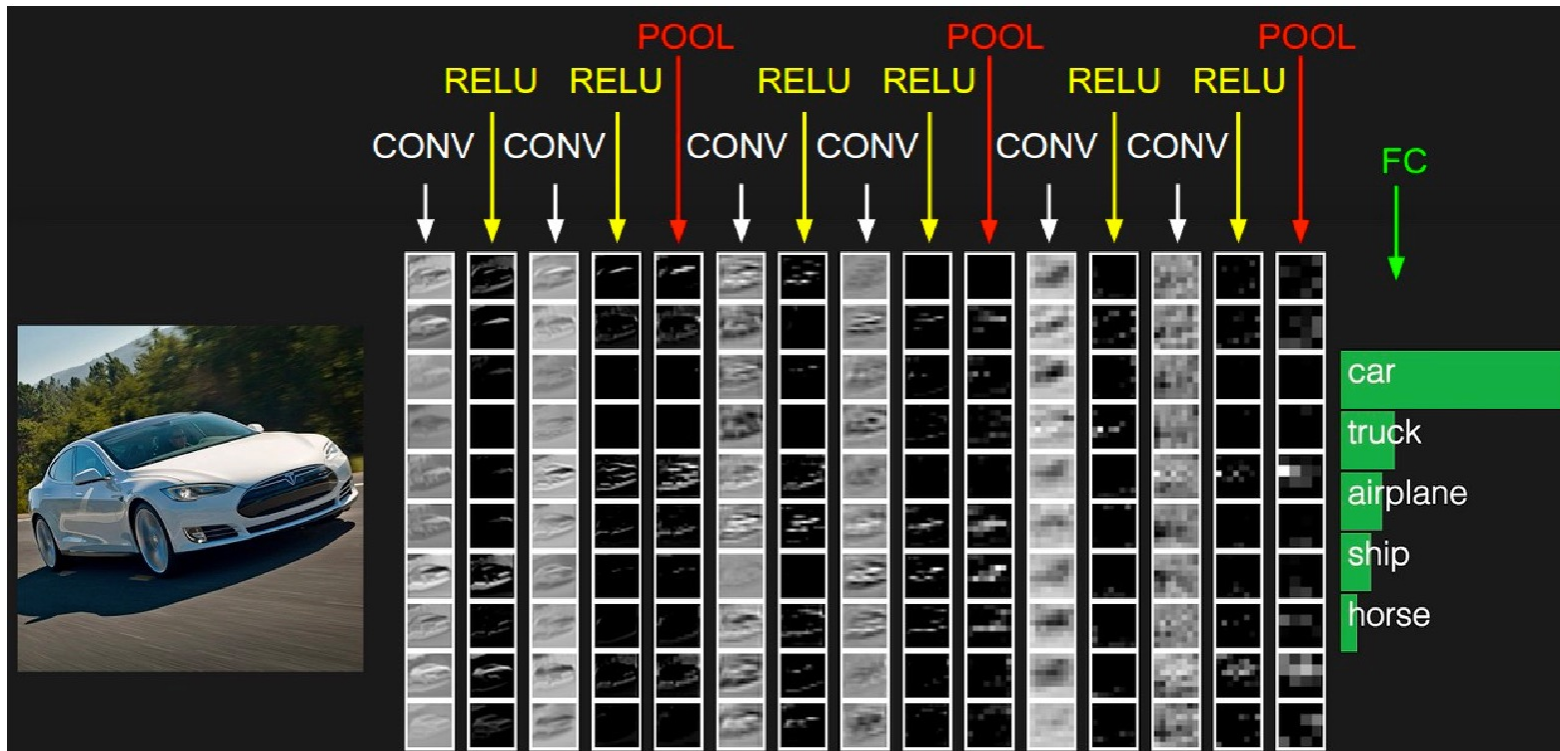
13,002

Each parameter has some impact on the output...need to tweak (learn) all parameters simultaneously to improve prediction accuracy

Convolutional Neural Networks

Convolutional neural networks (CNN)

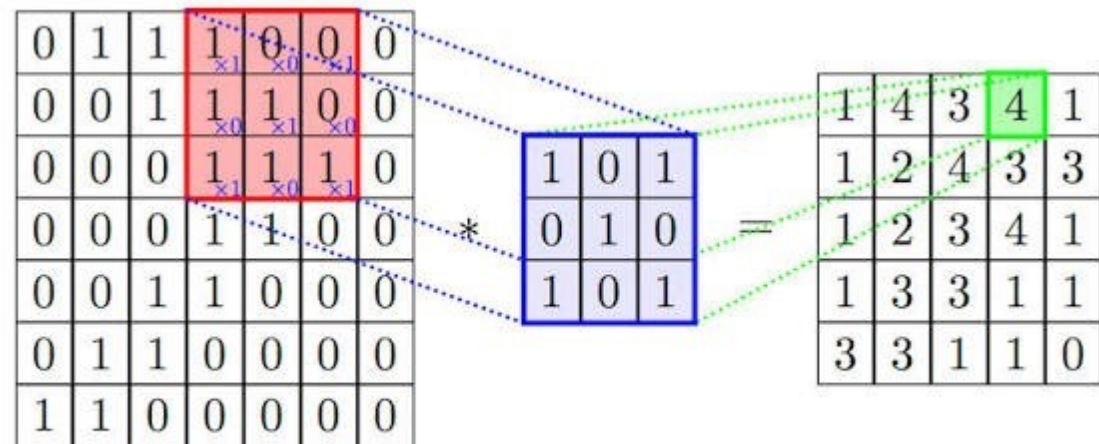
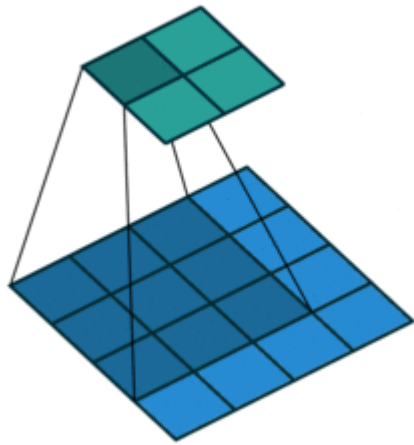
- A.K.A. ConvNet architecture
- A set of neural network architecture that consists of
 - convolutional layers
 - pooling layers
 - fully-connected (FC) layers



Convolution for single-channel images

Consider one filter with weights $\{w_{i,j}\}$ with size $F \times F$

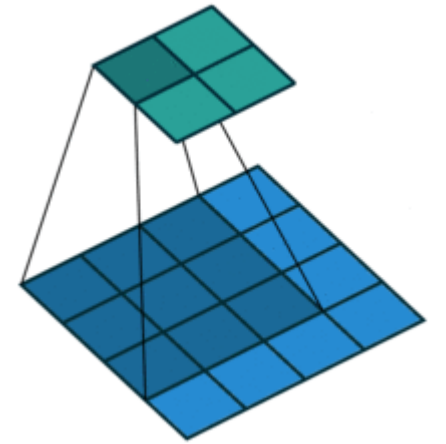
- For every $F \times F$ region of the image, perform inner product (= element wise product, then sum them all)
- Q: given a $w \times h$ image, after convolution with a $F \times F$ filter, what is the size of the resulting image?
- Terminologies: filter size, receptive field size, kernel.



Convolution: Some Intuition

Define the convolution of filter f on image I as:

$$(I * f)(x) = \sum_m \sum_n f(x - m, y - n)I(m, n)$$



Many ML libraries *actually* implement cross-correlation:

$$(f * I)(x) = \sum_m \sum_n f(x, y)I(x + m, y + n)$$

Learning finds good values for the convolution filter...

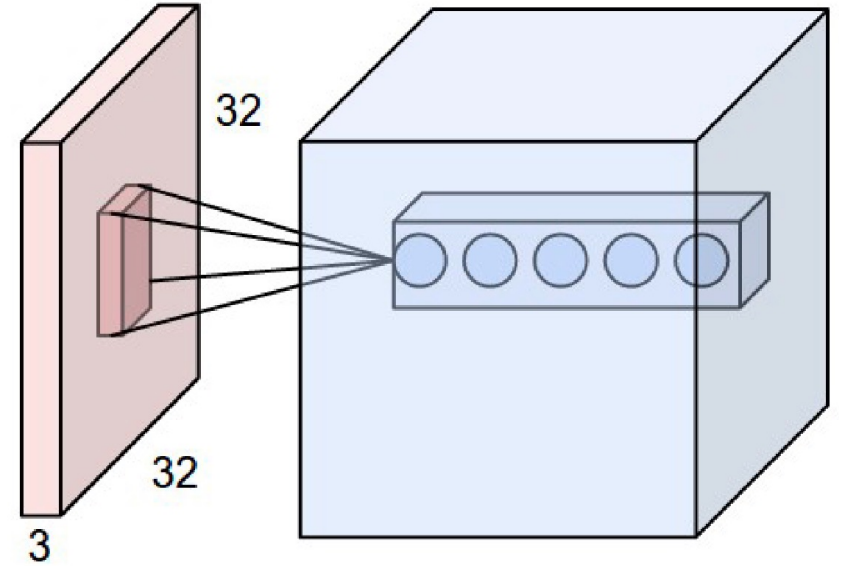
Convolutional layer for multi-channel images

Input: w (width) \times h (height) \times c (#channels)

- E.g. $32 \times 32 \times 3$
- 3 channels: R, G, and B

A convolutional filter on such image is of shape $F \times F \times c$

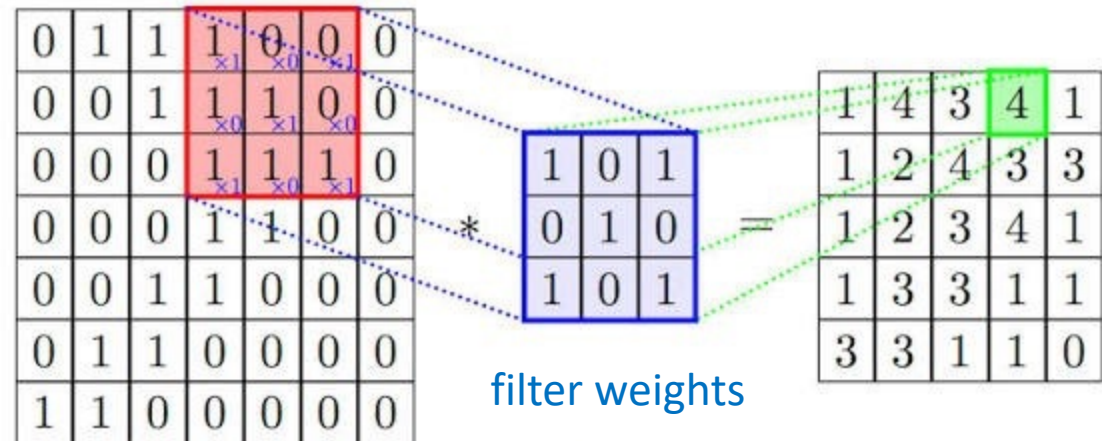
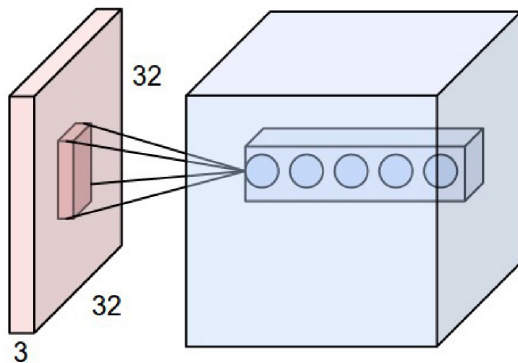
- Only spatial structure in the first two dimensions
- Denoted by $\{w_{i,j,k}\}$



Convolutional layer: visual explanation

- Consider one filter with weights $\{w_{i,j,k}\}$ with $5 \times 5 \times 3$
 - Imagine a sliding 3D window.
 - **Convolution:** For every 5×5 region of the image, perform inner product (= element wise product, then sum them all)
 - Then apply the activation function (e.g., ReLU)
- Results in $28 \times 28 \times 1$ – called **activation map**.
- Now, we can do K of these filters but with different weights $\{w_{i,j,k}^{(\ell)}\}$ for $\ell \in [K] \Rightarrow$ output is

$28 \times 28 \times K$



(depth=1 here)

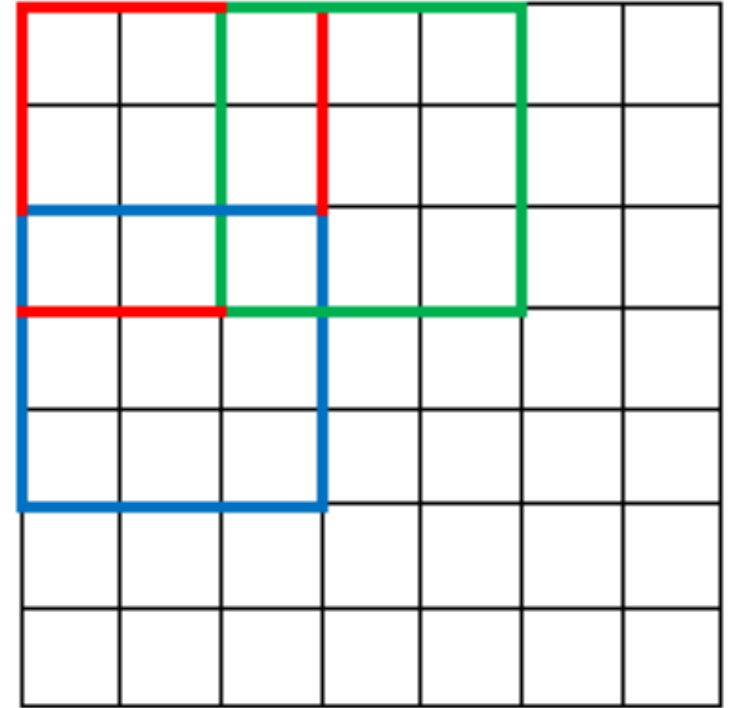
Convolutional Layer: More Details

Stride length S

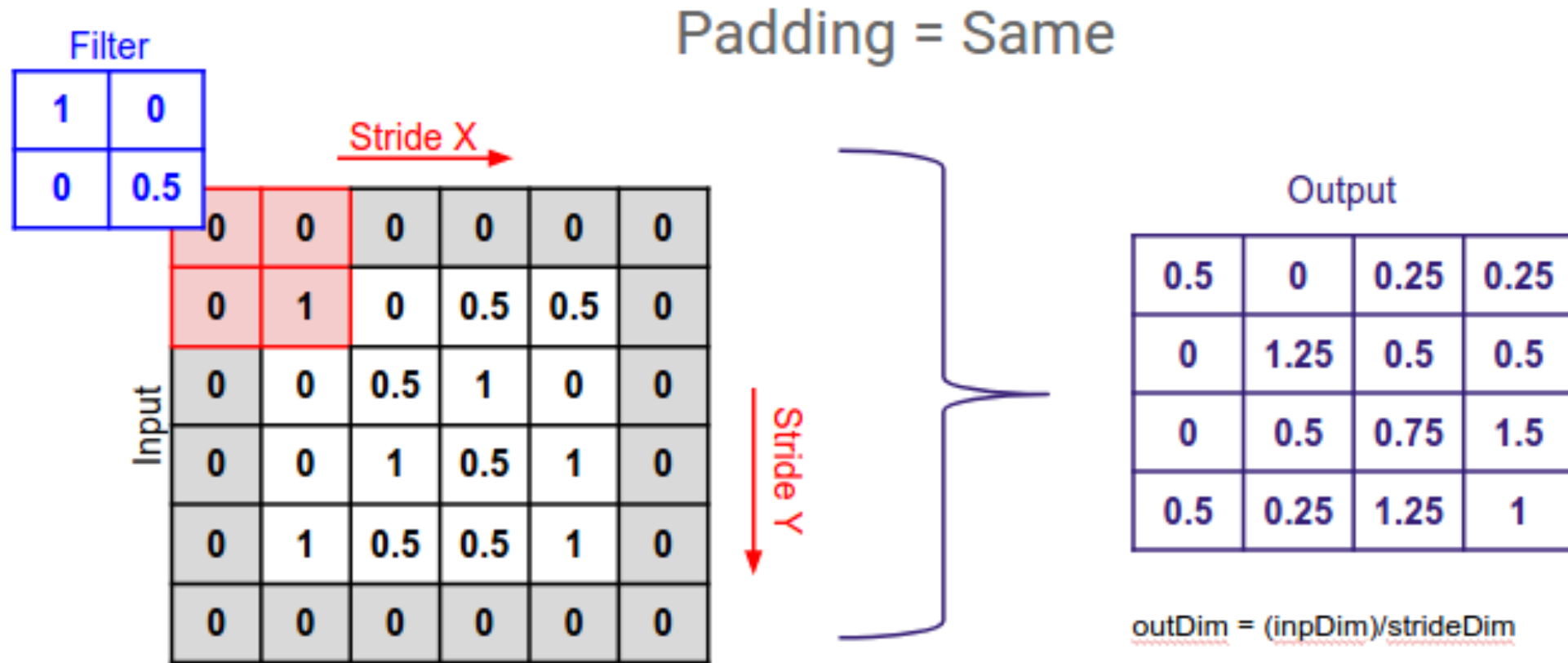
- Skip input regions; Move the sliding window of a filter not by 1 but by S .
- E.g., $S=2$ means skipping every other 5 by 5 region.

Zero-padding P : add P number of artificial pixels with value 0 around the input image on both sides

- To ensure the spatial dimension is maintained (otherwise, patterns at the corners are not detected well)
- If we use $P=1$, then the activation map will be 30×30 , not 28×28 in our example!



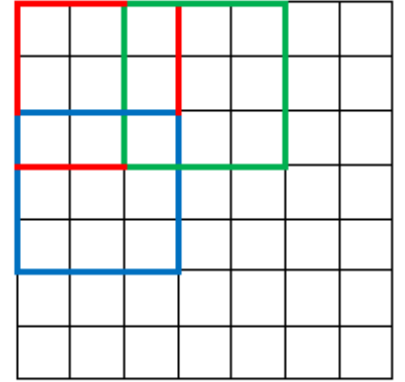
Example



Convolutional Layer: More Details

Stride length S

- Skip input regions; Move the sliding window of a filter not by 1 but by S .
- E.g., $S=2$ means skipping every other 5 by 5 region.



Zero-padding P : add P number of artificial pixels with value 0 around image.

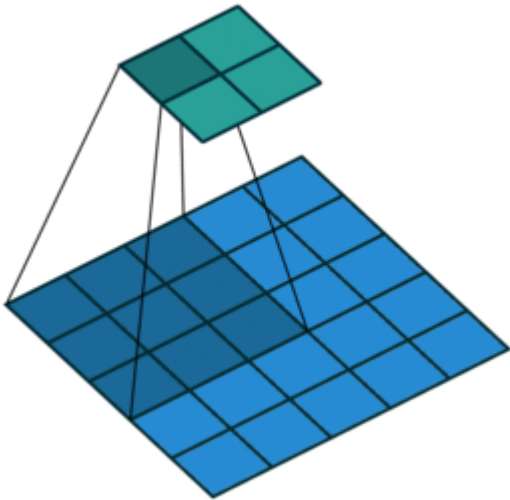
- To ensure the spatial dimension is maintained (otherwise, patterns at the corners are not detected well)
- If we use $P=2$, then the activation map will be 32 by 32 not 28 by 28 in our example!

Rules (same goes for **height**)

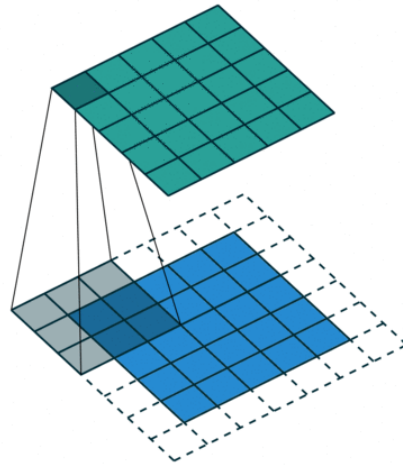
- W : input volume **width**, F : filter **width** (usually, the filter has the same width and height)
- The output **width** $K = \text{floor}((W - F + 2P)/S) + 1$
- E.g., $W=32, F=5, P=0, S=1 \Rightarrow K = 28$
- E.g., $W=32, F=5, P=2, S=1 \Rightarrow K = 32$

Strides and padding: animations

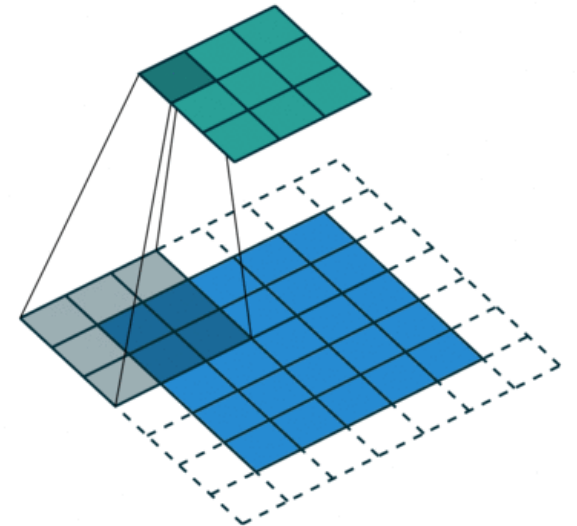
Strides only



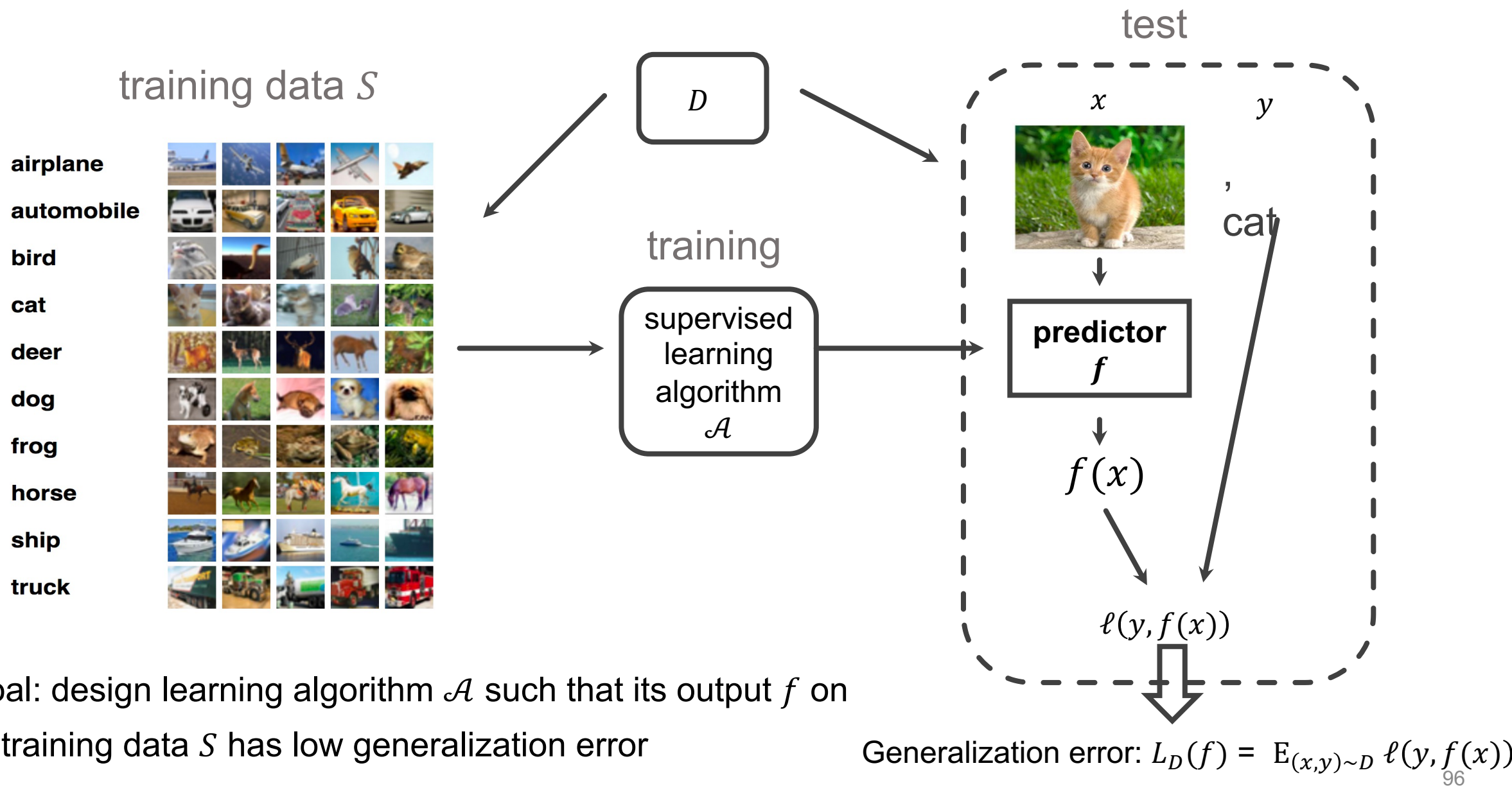
Padding only



Strides + Padding



Supervised learning setup: putting it together



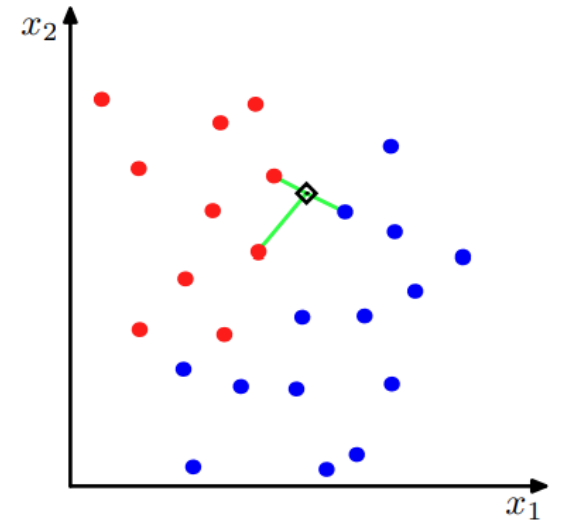
k -nearest neighbors (k -NN): main concept

Training set: $S = \{ (x_1, y_1), \dots, (x_m, y_m) \}$

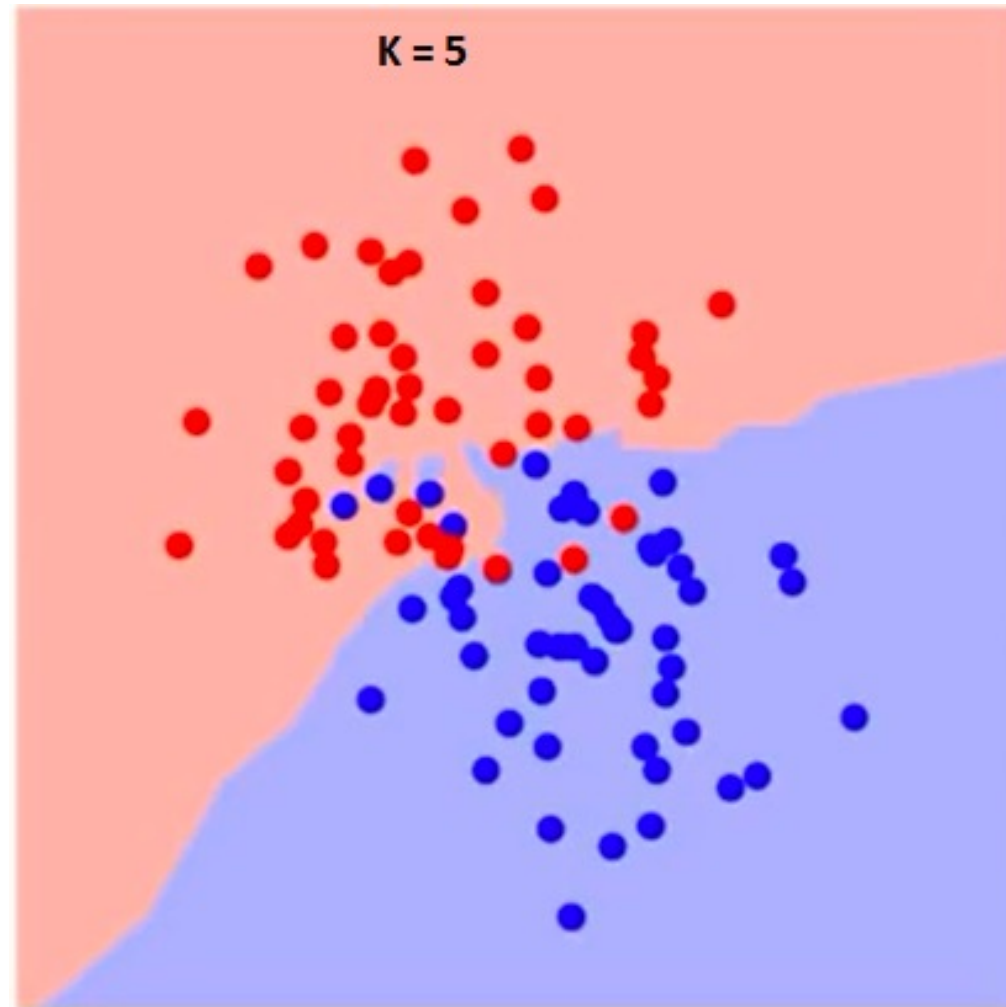
Inductive bias: given test example x , its label should resemble the labels of **nearby points**

Function

- input: x
- find the k nearest points to x from S ; call their indices $N(x)$
- output: the majority vote of $\{y_i: i \in N(x)\}$
 - For regression, the average.



k-NN classification example



decision boundary

k -NN classification: pseudocode

- Training is trivial: store the training set
- Test:

Algorithm 3 KNN-PREDICT(\mathbf{D}, K, \hat{x})

list	→	1: $S \leftarrow []$	
		2: for $n = 1$ to N do	
append to list	→	3: $S \leftarrow S \oplus \langle d(x_n, \hat{x}), n \rangle$	// store distance to training example n
		4: end for	
sort in first coordinate	→	5: $S \leftarrow \text{SORT}(S)$	// put lowest-distance objects first
		6: $\hat{y} \leftarrow 0$	
		7: for $k = 1$ to K do	
		8: $\langle \text{dist}, n \rangle \leftarrow S_k$	// n this is the k th closest data point
		9: $\hat{y} \leftarrow \hat{y} + y_n$	// vote according to the label for the n th training point
		10: end for	
Majority vote of $\{y_i: i \in N(x)\}$	→	11: return SIGN(\hat{y})	// return +1 if $\hat{y} > 0$ and -1 if $\hat{y} < 0$

- Time complexity (assuming distance calculation takes $O(d)$ time)
 - $O(m d + m \log m + k) = O(m(d + \log m))$
- Faster nearest neighbor search: k-d trees, locality sensitive hashing

Variations

- Classification

- Recall the majority vote rule: $\hat{y} = \arg \max_{y \in \{1, \dots, C\}} \sum_{i \in N(x)} 1\{y_i = y\}$

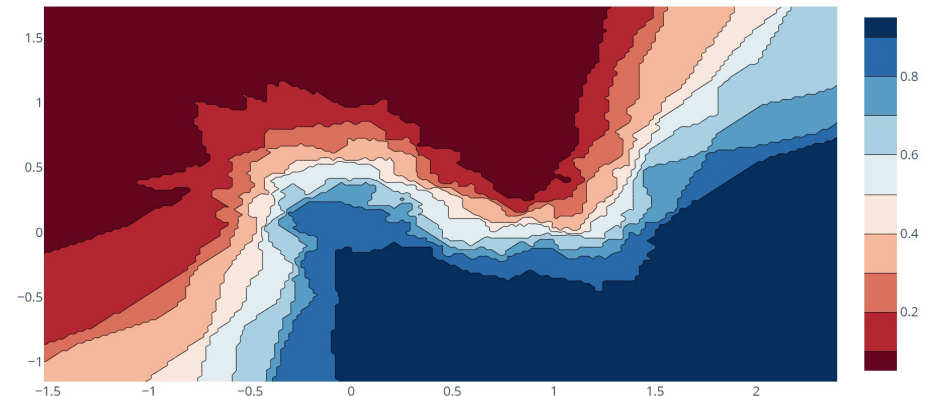
- Soft weighting nearest neighbors: $\hat{y} = \arg \max_{y \in \{1, \dots, C\}} \sum_{i=1}^m w_i 1\{y_i = y\}$,

where $w_i \propto \exp(-\beta d(x, x_i))$, or $\propto \frac{1}{1+d(x, x_i)^\beta}$

- Class probability estimates

- $\hat{P}(Y = y | x) = \frac{1}{k} \sum_{i \in N(x)} 1\{y_i = y\}$

- Useful for “classification with rejection”



Inductive Bias

Training



class A



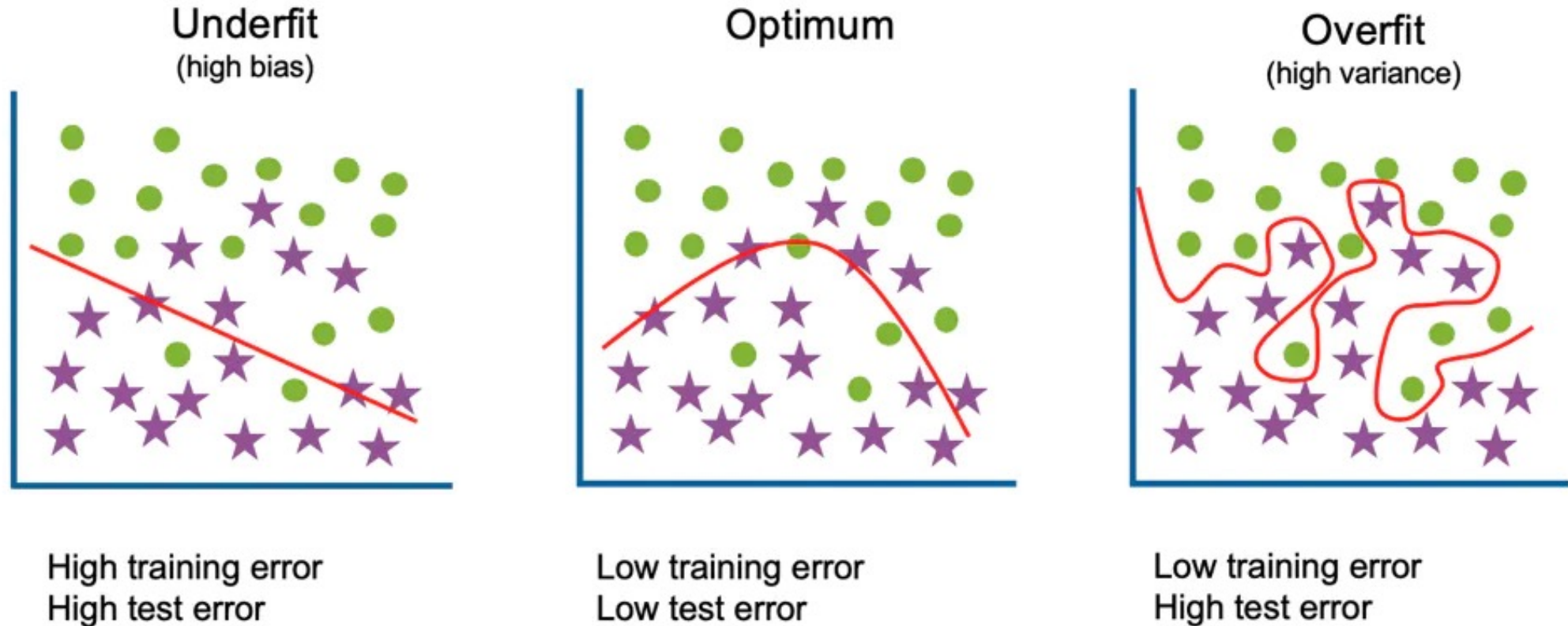
class B

Test



How would you label the test examples?

Overfitting vs Underfitting



Bayes optimal classifier

Theorem f_{BO} achieves the smallest 0-1 error among all classifiers.

$$f_{BO}(x) = \arg \max_{y \in \mathcal{Y}} P_D(X = x, Y = y) = \arg \max_{y \in \mathcal{Y}} P_D(Y = y | X = x), \forall x \in \mathcal{X}$$

Example Iris dataset classification:



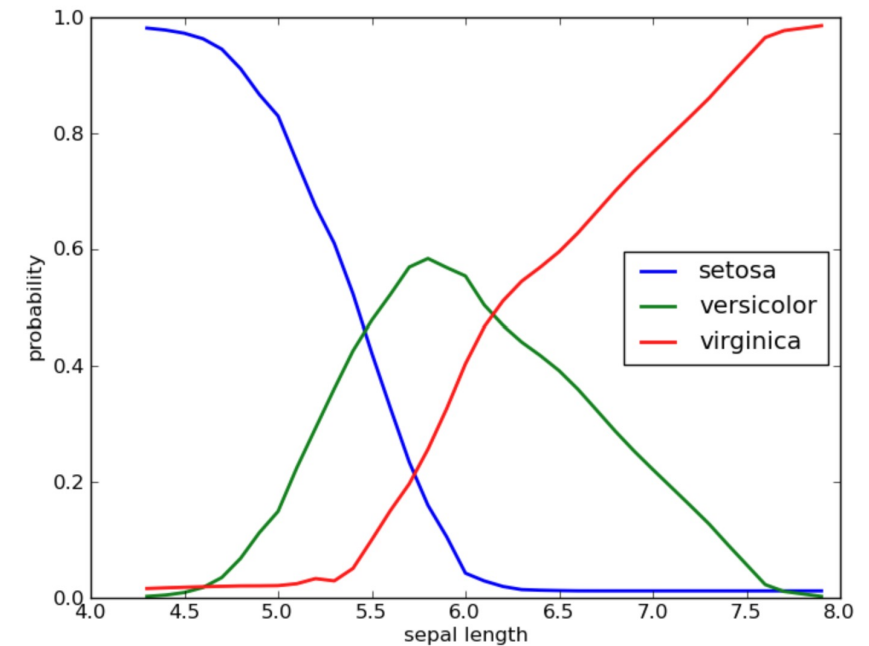
Iris Setosa



Iris Versicolor



Iris Virginica

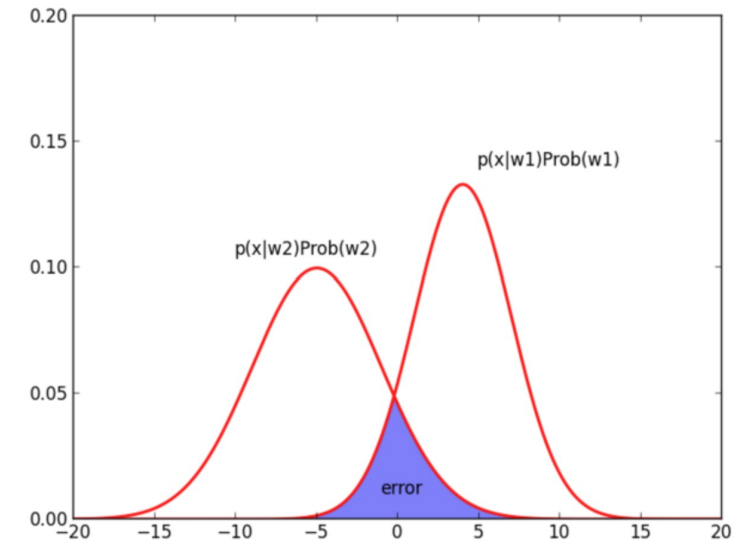
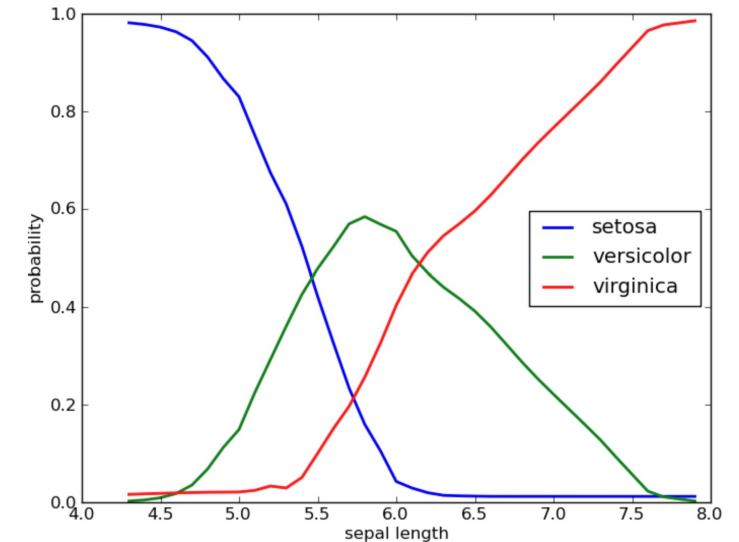


Bayes error rate: alternative form

$$\begin{aligned}L_D(f_{BO}) &= P_D(Y \neq f_{BO}(X)) \\&= \sum_x P_D(Y \neq f_{BO}(x) \mid X = x) P_D(X = x) \\&= \sum_x (1 - P_D(Y = f_{BO}(x) \mid X = x)) P_D(X = x) \\&= \sum_x \left(1 - \max_y P_D(Y = y \mid X = x)\right) P_D(X = x) \\&= E \left[1 - \max_y P_D(Y = y \mid X)\right]\end{aligned}$$

- **Special case: binary classification**

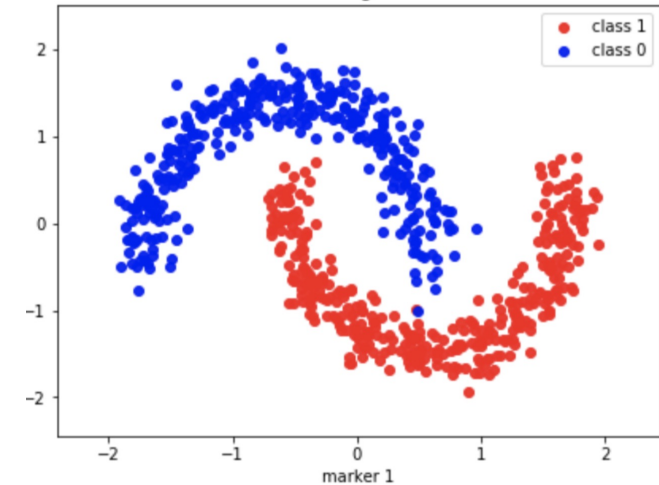
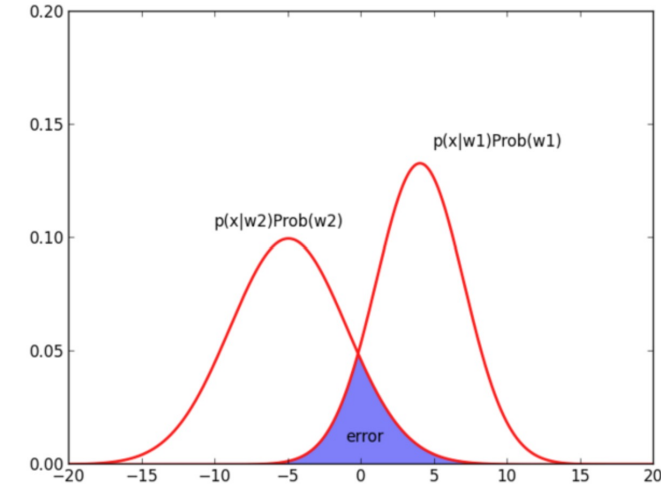
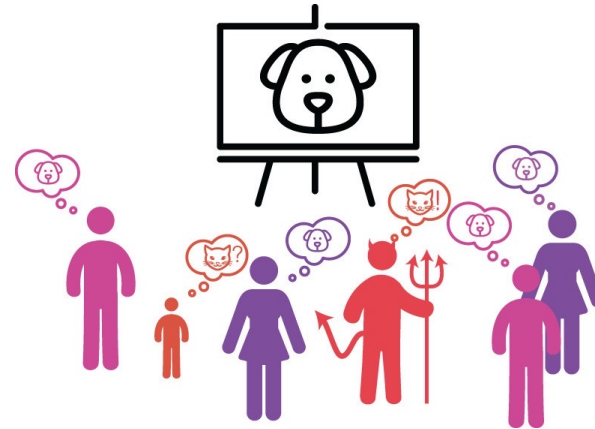
- $L_D(f_{BO}) = \sum_x P_D(Y \neq f_{BO}(x), X = x)$
 $= \sum_x \min(P_D(Y = +1, X = x), P_D(Y = -1, X = x))$



When is the Bayes error rate nonzero?

$$L_D(f_{BO}) = \sum_x \min(P_D(Y = +1, X = x), P_D(Y = -1, X = x))$$

- Limited feature representation
- Noise in the training data
 - Feature noise
 - Label noise
 - Sensor failure
 - Typo in reviews for sentiment classification
- May not be a single “correct” answer
- Inductive bias of the model / learning algorithm



New measures of classification performance

- True positive rate (TPR)

$$= \frac{TP}{P} = \frac{P(\hat{y}=+1, y=+1)}{P(y=+1)}$$

(aka recall, sensitivity)

- True negative rate (TNR) = $\frac{TN}{N}$

(specificity)

- False positive rate (FPR) = $\frac{FP}{N}$

- False negative rate (FNR) = $\frac{FN}{P}$

- Precision = $\frac{TP}{P\text{-called}} = \frac{P(\hat{y}=+1, y=+1)}{P(\hat{y}=+1)}$, P - called = TP + FP

The diagram illustrates a confusion matrix. The horizontal axis is labeled 'actual class' and is divided into 'positive' and 'negative'. The vertical axis is labeled 'predicted class' and is divided into 'positive' and 'negative'. The matrix cells are: top-left (positive predicted, positive actual) is 'true positives (TP)'; top-right (positive predicted, negative actual) is 'false positives (FP)' with a blue label 'Type I error' below it; bottom-left (negative predicted, positive actual) is 'false negatives (FN)' with a blue label 'Type II error' below it; bottom-right (negative predicted, negative actual) is 'true negatives (TN)'. Red brackets group the columns under 'actual class' and the rows under 'predicted class'.

		actual class	
		positive	negative
predicted class	positive	true positives (TP)	false positives (FP) Type I error
	negative	false negatives (FN) Type II error	true negatives (TN)

$$P = TP + FN \quad N = FP + TN$$

Linear Regression

Regression Learn a function that predicts outputs from inputs,

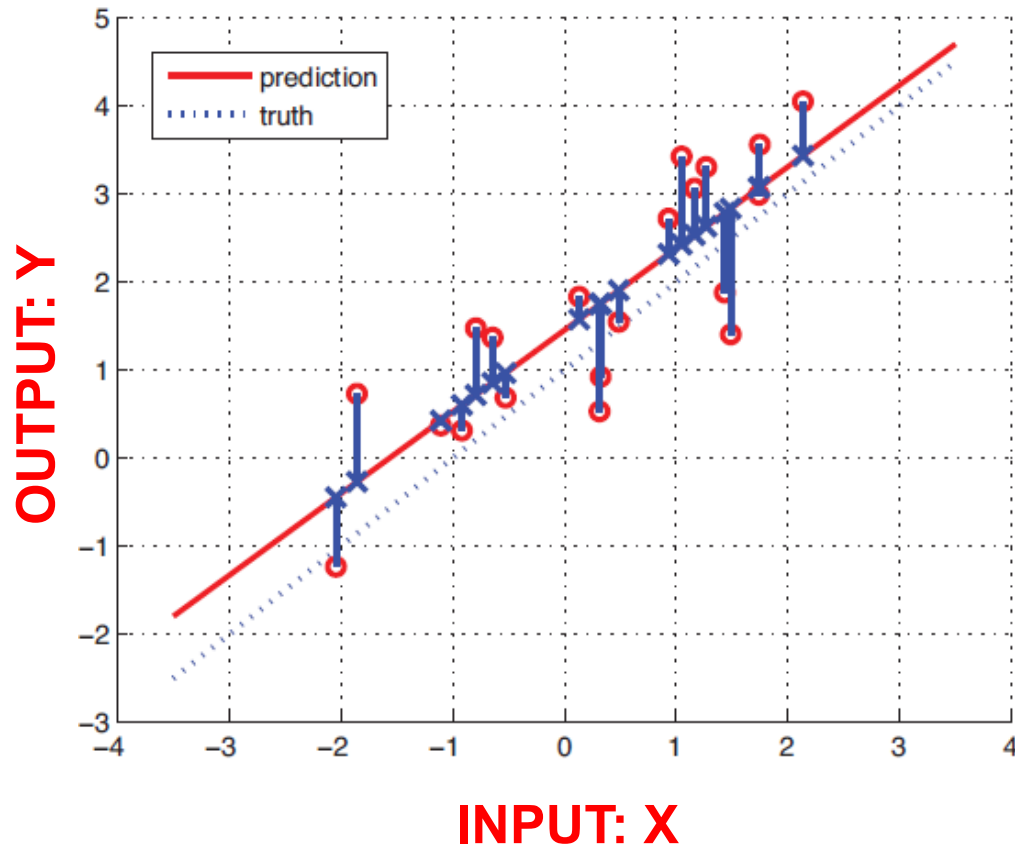
$$y = f(x)$$

Outputs y are real-valued

Linear Regression As the name suggests, uses a *linear function*:

$$y = w^T x + b$$

We will add noise later...



Linear Regression

Input-output mapping is not exact, so we will add zero-mean Gaussian noise,

$$y = w^T x + \epsilon \quad \text{where} \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

**Multivariate Normal
(uncorrelated)**

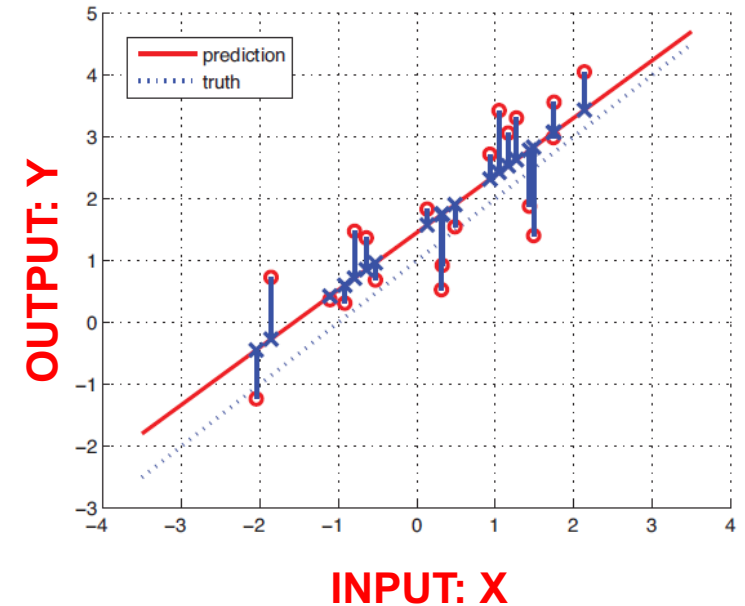
This is equivalent to the likelihood function,

$$p(y | w, x) = \mathcal{N}(y | w^T x, \sigma^2)$$

Because Adding a constant to a Normal RV is still a Normal RV,

$$z \sim \mathcal{N}(m, P) \quad z + c \sim \mathcal{N}(m + c, P)$$

In the case of linear regression $z \rightarrow \epsilon$ and $c \rightarrow w^T x$



Great, we're done right?

We need to fit it to data by learning the regression weights

Data – We have this

$$y = w^T x + \epsilon$$

Random; Can't do anything about it

**Don't know these;
need to learn them**

How to do this?
What makes *good* weights?

Learning Linear Regression Models

There are several ways to think about fitting regression:

- **Intuitive** Find a plane/line that is close to data
- **Functional** Find a line that minimizes the *least squares* loss
- **Estimation** Find maximum likelihood estimate of parameters

They are all the same thing...

Learning Linear Regression Models

There are several ways to think about fitting regression:

- **Intuitive** Find a plane/line that is close to data
- **Functional** Find a line that minimizes the *least squares* loss
- **Estimation** Find maximum likelihood estimate of parameters

They are all the same thing...

MLE for Linear Regression

Given training data $\{(x_i, y_i)\}_{i=1}^N$ likelihood function is given by,

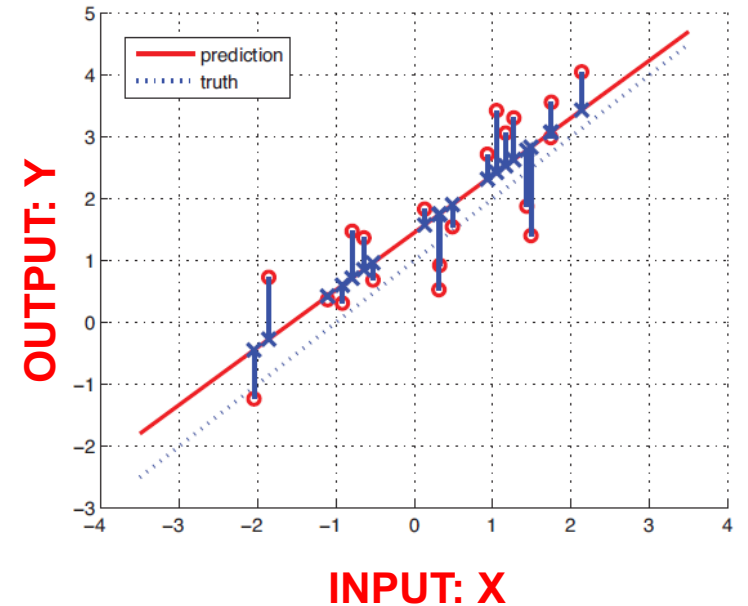
$$\log \prod_{i=1}^N p(y_i | x_i, w) = \sum_{i=1}^N \log p(y_i | x_i, w)$$

Recall that the likelihood is Gaussian:

$$p(y | w, x) = \mathcal{N}(y | w^T x, \sigma^2)$$

So MLE maximizes the log-likelihood over the whole data as,

$$w^{\text{MLE}} = \arg \max_w \sum_{i=1}^N \log \mathcal{N}(y_i | w^T x_i, \sigma^2)$$



MLE of Gaussian Mean

Assume data are i.i.d. univariate Gaussian,

$$p(\mathcal{Y} | \mu) = \prod_{i=1}^N \mathcal{N}(y_i | \mu, \sigma^2)$$

↖ Variance is known

Log-likelihood function:

$$\mathcal{L}(\mu) = \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2} (y_i - \mu)^2 \sigma^{-2} \right) \right)$$

Constant doesn't
depend on mean

↖

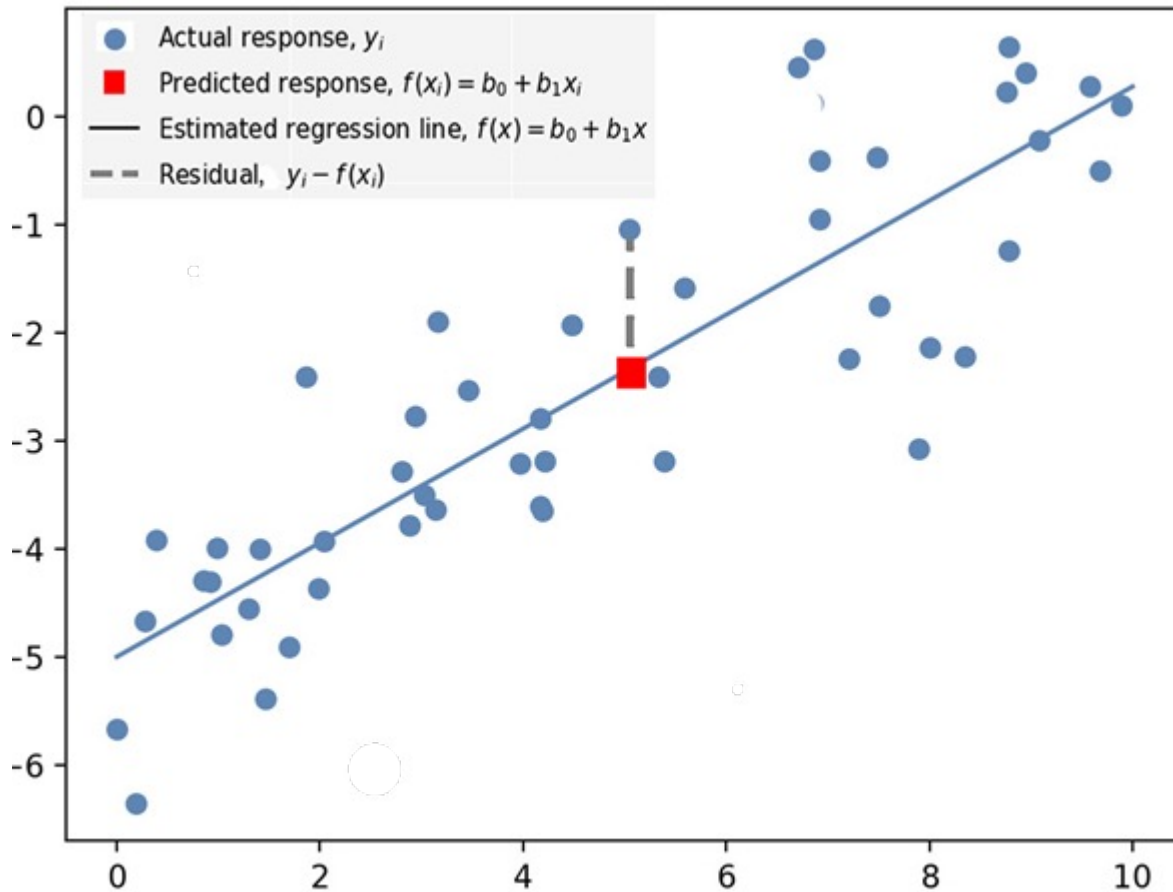
$$= \text{const.} - \frac{1}{2} \sum_{i=1}^N ((y_i - \mu)^2 \sigma^{-2})$$

MLE doesn't change when we:
1) Drop constant terms (in μ)
2) Minimize negative log-likelihood

MLE estimate is *least squares estimator*:

$$\mu^{\text{MLE}} = -\frac{1}{2\sigma^2} \arg \max_{\mu} \sum_{i=1}^N (y_i - \mu)^2 = \arg \min_{\mu} \sum_{i=1}^N (y_i - \mu)^2$$

MLE of Linear Regression



Substitute linear regression prediction into MLE solution and we have,

$$\min_w \sum_{i=1}^N (y_i - wx_i)^2$$

So for Linear Regression,
MLE = Least Squares
Estimation

MLE of Linear Regression

Using previous results, MLE is equivalent to minimizing squared residuals,

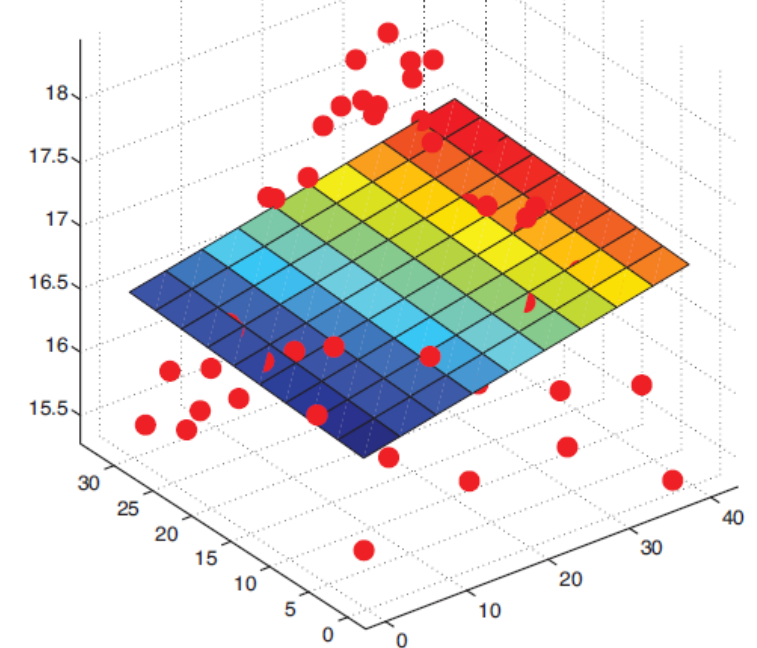
$$\min_w \sum_{i=1}^N (y_i - w^T x_i)^2 = \|\mathbf{y} - w^T \mathbf{X}\|^2$$

Some slightly more advanced linear algebra gives us a solution,

$$w = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Ordinary Least Squares (OLS) solution

[Image: Murphy, K. (2012)]



- Derivation a bit involved for lecture but...
- We know it has a closed-form and why
 - We can evaluate it
 - Generally know where it comes from

Basis Functions

- A **basis function** can be any function of the input features X
- Define a set of m basis functions $\phi_1(x), \dots, \phi_m(x)$
- Fit a linear regression model in terms of basis functions,

$$y = \sum_{i=1}^m w_i \phi_i(x) = w^T \phi(x)$$

- Regression model is *linear in the basis transformations*
- Model is *nonlinear in the data X*

Kernel Functions

*A **kernel function** is an inner-product of some basis function computed on two inputs*

$$k(x, x') = \phi(x)^T \phi(x') = \sum_{i=1}^M \phi_i(x) \phi_i(x')$$

A consequence is that kernel functions are non-negative real-valued functions over a pair of inputs,

$$\kappa(x, x') \in \mathbb{R} \qquad \kappa(x, x') \geq 0$$

Kernel functions can be interpreted as a measure of distance between two inputs

Kernel Functions

Example The *linear basis* $\phi(x) = x$ produces the kernel,

$$\kappa(x, x') = \phi(x)^T \phi(x') = x^T x'$$

It is often easier to directly specify the kernel rather than the basis function...

Example Gaussian kernel models similarity according to an unnormalized Gaussian distribution,

$$\kappa(x, x') = \exp\left(-\frac{1}{2\sigma^2}(x - x')^2\right)$$

Note Despite the name, this is **not** a Gaussian probability density.

Also called a *radial basis function* (RBF)

Kernel Functions

Given *any* set of data $\{x_i\}_{i=1}^n$ a necessary and sufficient condition of a valid kernel function is that the $n \times n$ **gram matrix**,

$$\mathbf{K} = \begin{pmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \dots & \kappa(x_1, x_n) \\ \kappa(x_2, x_1) & \kappa(x_2, x_2) & \dots & \kappa(x_2, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ \kappa(x_n, x_1) & \kappa(x_n, x_2) & \dots & \kappa(x_n, x_n) \end{pmatrix}$$

Is a *symmetric positive semidefinite matrix*.

Kernel Ridge Regression

Kernel representation requires inversion of NxN matrix

Primal

$$\Phi = \begin{pmatrix} 1 & \phi_1(x_1) & \dots & \phi_M(x_1) \\ 1 & \phi_1(x_2) & \dots & \phi_M(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \phi_1(x_N) & \dots & \phi_M(x_N) \end{pmatrix}$$

$$w = (\underbrace{\Phi^T \Phi + \lambda I}_{\text{MxM Matrix Inversion}})^{-1} \Phi^T y$$

MxM Matrix Inversion
 $O(M^3)$

Dual

$$\mathbf{K} = \begin{pmatrix} \kappa(x_1, x_1) & \kappa(x_1, x_2) & \dots & \kappa(x_1, x_n) \\ \kappa(x_2, x_1) & \kappa(x_2, x_2) & \dots & \kappa(x_2, x_n) \\ \vdots & \vdots & \vdots & \vdots \\ \kappa(x_n, x_1) & \kappa(x_n, x_2) & \dots & \kappa(x_n, x_n) \end{pmatrix}$$

$$y(x) = \mathbf{k}(x)^T (\underbrace{\mathbf{K} + \lambda I}_{\text{NxN Matrix Inversion}})^{-1} \mathbf{y}$$

NxN Matrix Inversion
 $O(N^3)$

Number of training data N greater than basis functions M