

UNSUPERVISED LEARNING

L REPRESENTATION LEARNING.

Goal: Given a set of "data points", "understand"
Something "useful" about them.

Data points \rightarrow vectors in \mathbb{R}^d $\begin{bmatrix} \text{height} \\ \text{weight} \\ \text{age} \end{bmatrix} \in \mathbb{R}^3$

Running theme: "COMPREHENSION IS COMPRESSION"
 \hookrightarrow understanding
 \hookrightarrow learning
(George Church)

Problem: Input: $\{x_1, x_2, \dots, x_n\}$ $x_i \in \mathbb{R}^d$ \leftarrow # of features

Output: Some "Compressed" representation of the dataset

Example: $\{x_1, x_2, x_3, x_4\}$

$$\left\{ \begin{bmatrix} -7 \\ -14 \end{bmatrix}, \begin{bmatrix} 2.5 \\ 5 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}$$

Question: How many real numbers are needed to store this dataset 8

Representative

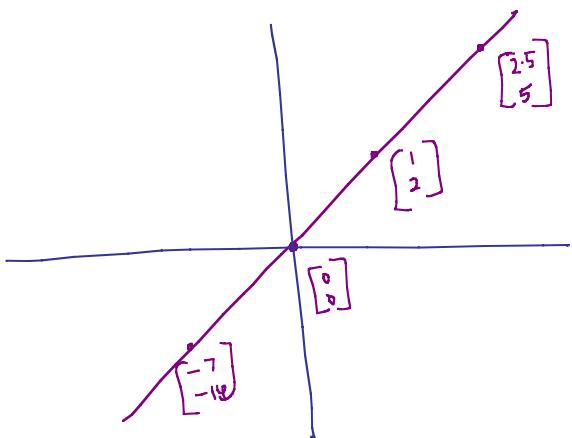
$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \in \mathbb{R}^2$$

co-efficients

$$\{-7, 2.5, 0.5, 0\}$$

6

NOTE:
using representative
ie co-efficients
can "RECONSTRUCT"
the dataset exactly.



Rep:

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

to eff
 $\left\{ -7, 2.5, 0.5, 0 \right\}$

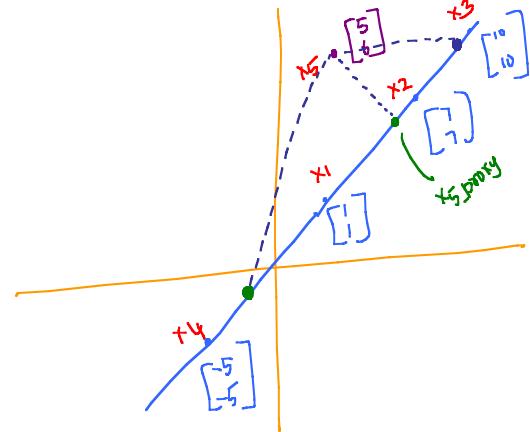
$$\begin{bmatrix} \sqrt{5} \\ 2\sqrt{5} \end{bmatrix}$$

$$\left\{ -7\sqrt{5}, 2.5\sqrt{5}, 0.5\sqrt{5}, 0 \right\}$$

NOTE:
 Any vector along
 the purple line can
 be chosen as
 a representative
 [except $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$]

Original: # real numbers = $d n$

real numbers in
 compressed representation = $d + n$



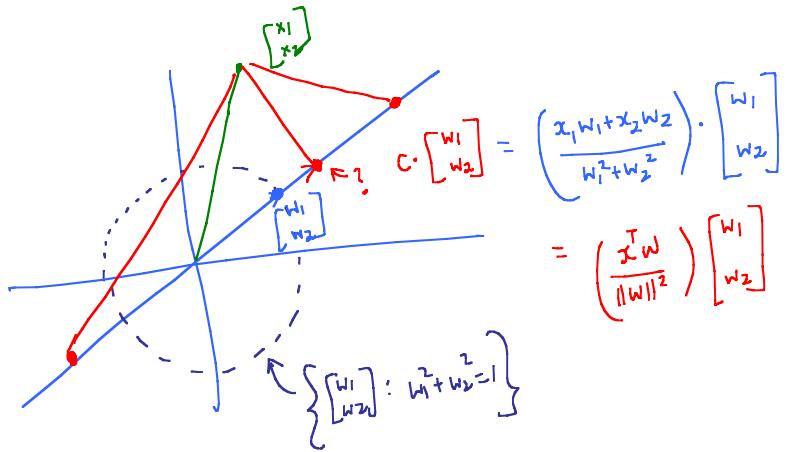
Q: Who can "pretend" to be a proxy for x_5 along the blue line?

Ans: Projection of x_5 onto the blue line.

$$\text{Rep} \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$

$$\left\{ (-5, -5), (1, 1), (1, 7), (5, 1), \underline{(10, 10)} \right\}$$

$$x_3 = 10 \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 10 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 10 \\ 10 \end{bmatrix}$$



Inner product/
dot product of
 \mathbf{x} and \mathbf{w} .

$$\min_c \text{length}^2(\text{error vector}) \rightarrow \begin{bmatrix} x_1 - c w_1 \\ x_2 - c w_2 \end{bmatrix}$$

$$\min_c (x_1 - c w_1)^2 + (x_2 - c w_2)^2$$

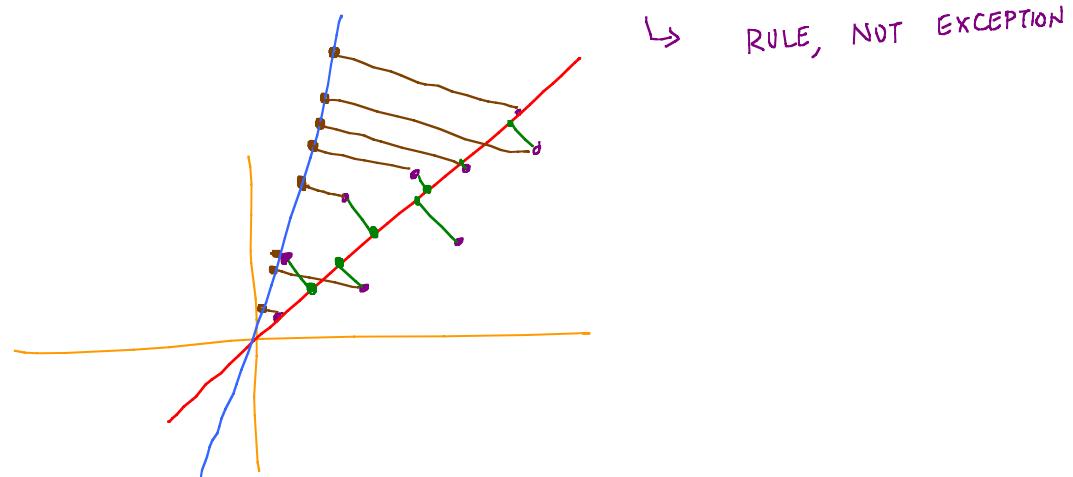
$$c^* = \left(\frac{x_1 w_1 + x_2 w_2}{w_1^2 + w_2^2} \right) \quad (\text{scalar})$$

$$\text{length}^2 \left(\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \right) \leftarrow$$

NOTE: Can always pick $\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ s.t.
 $\text{length} \left(\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \right) = 1$

$$\Rightarrow c^* = (\mathbf{x}^T \mathbf{w}) \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

Goal: Develop a way to find a "compressed" representation
of data when data points not-necessarily fall on line



Goal: Find the line that has the least "reconstruction" error.

$$\|z\|^2 = z^T z$$

Dataset: $\{x_1, x_2, \dots, x_n\} \quad x_i \in \mathbb{R}^d$

$$\begin{aligned} \text{ERROR}(\text{line}, \text{dataset}) &= \sum_{i=1}^n \text{error}(\text{line}, x_i) \\ &= \sum_{i=1}^n \text{length}^2(x_i - (x_i^T w)w) \\ &= \sum_{i=1}^n \|x_i - (x_i^T w) \cdot w\|^2 \end{aligned}$$

Represented
 using w
 s.t. $\|w\|=1$

$$\begin{aligned} f(w) &= \frac{1}{n} \sum_{i=1}^n \left\| \underbrace{x_i - (x_i^T w) \cdot w}_{\text{error}} \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - (x_i^T w) \cdot w)^T (x_i - (x_i^T w) \cdot w) \end{aligned}$$

$$= \frac{1}{n} \sum_{i=1}^n \left[x_i^\top x_i - (x_i^\top \omega)^2 - (x_i^\top \omega)^2 + (x_i^\top \omega)^2 \cdot 1 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i^\top x_i - (x_i^\top \omega)^2)$$

$$\min_{\substack{\omega \\ \|\omega\|^2=1}} g(\omega) = \frac{1}{n} \sum_{i=1}^n -(x_i^\top \omega)^2$$

$$\max_{\substack{\omega \\ \|\omega\|^2=1}} \frac{1}{n} \sum_{i=1}^n (x_i^\top \omega)^2 = \frac{1}{n} \sum_{i=1}^n (\underbrace{x_i^\top \omega}_{d \times d \text{ matrix}})(\underbrace{x_i^\top \omega}_{d \times 1})$$

$$= \frac{1}{n} \sum_{i=1}^n \underline{\omega^\top} (x_i x_i^\top) \underline{\omega}$$

$$= \underline{\omega^\top} \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right) \underline{\omega}$$

C $d \times d$ matrix

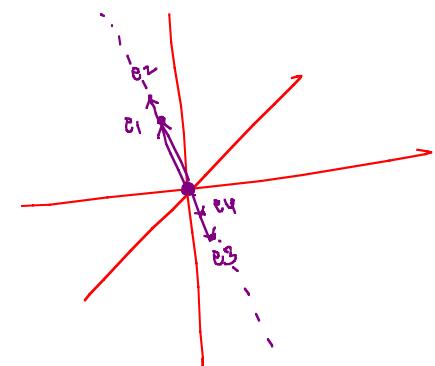
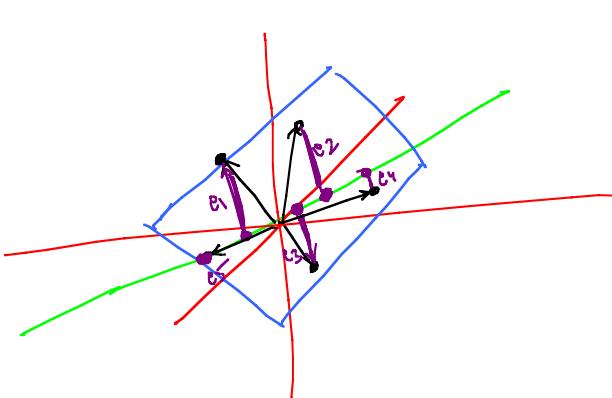
equivalently

$$\max_{\mathbf{w}} \mathbf{w}^T \mathbf{C} \mathbf{w}$$
$$\text{subject to } \|\mathbf{w}\|^2 = 1$$

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

↳ Co-variance matrix

Soln!: \mathbf{w} is the eigenvector corresponding to the maximum eigenvalue of \mathbf{C}



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

↓ Find w

$(x^T w) \cdot w$

↓ Residue/error

$x - (x^T w) \cdot w$

Might not be
error but
has "information"

Possible Algorithm

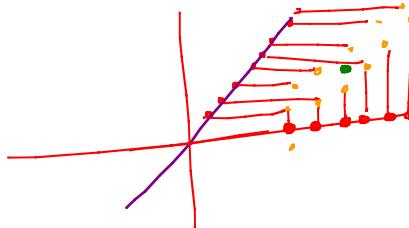
Input: $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$

→ $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ $x_i = x_i - \bar{x}$ + ϵ
→ Find "best" line $w_1 \in \mathbb{R}^d$

→ Replace $x_i \leftarrow x_i - (x_i^T w)$

→ Repeat to obtain w_2

Issue: Data may not
be centered.



Questions

→ How to solve $\max_w w^T C w$?
 $\|w\|^2 = 1$

→ How many times to repeat the procedure?

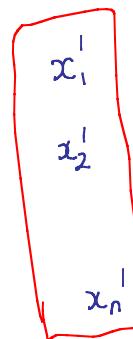
→ Where exactly is "compression" is happening?

→ What "representations" are we learning?

$$D = \{x_1, x_2, \dots, x_n\} \quad x_i \in \mathbb{R}^d.$$

$$w_1 = \underset{\substack{w: \\ \|w\|^2=1}}{\operatorname{argmax}} \quad w^T C w$$

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

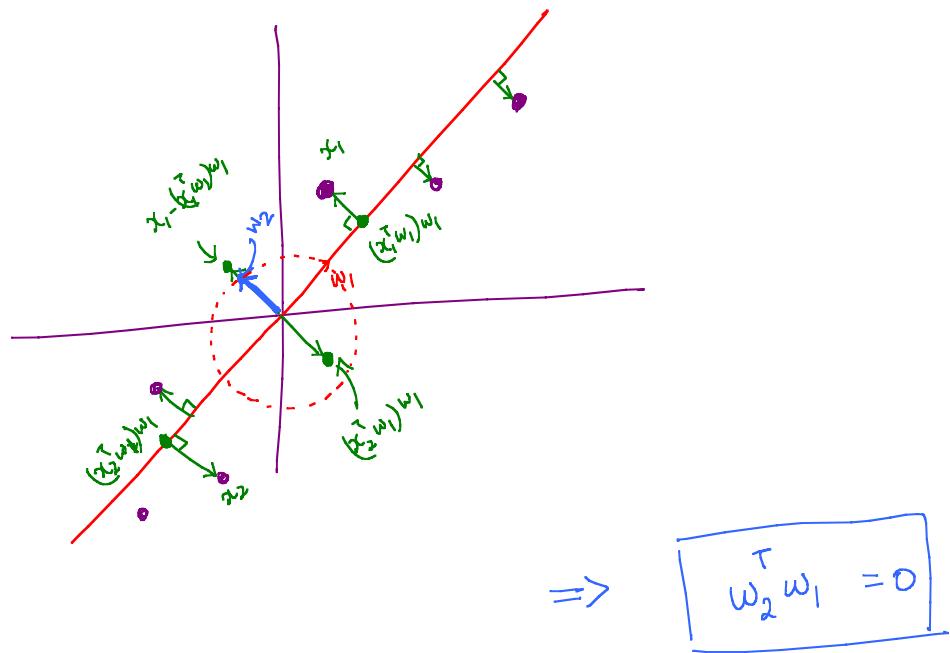


$$\begin{aligned} x_1 &\rightarrow x_1 - (x_1^T w_1) w_1 \\ x_2 &\rightarrow x_2 - (x_2^T w_1) w_1 \\ &\vdots \\ x_n &\rightarrow x_n - (x_n^T w_1) w_1 \end{aligned}$$

$$w_2 = \underset{\substack{w: \\ \|w\|_2^2=1}}{\operatorname{argmax}} \quad w^T C' w$$

$$C' = \frac{1}{n} \sum_{i=1}^n x_i^T x_i$$

Question: What can we say about $w_1 \neq w_2$?



Observation:

- All residues are orthogonal to w_1
- Any line which minimizes sum of errors w.r.t residues must also be orthogonal to w_1 [ARGUE WHY]

By continuing this procedure, we get $\{w_1, w_2, \dots, w_d\}$ s.t. $\|w_k\|_2^2 = 1 \forall k$ and $w_i^\top w_j = 0 \forall i \neq j$

\downarrow

ORTHONORMAL
VECTORS

Residue after round 1

$$\left\{ z_1 - (z_1^\top w_1)w_1, \dots, z_n - (z_n^\top w_1)w_1 \right\}$$

all vectors in
 $\in \mathbb{R}^d$

- $w_2 \rightarrow$ Best line that fits residues.
- $w_1^\top w_2 = 0$

Residues after round 2

$$\left\{ \underbrace{x_i - (x_i^T w_1) w_1}_{\text{Residue}} - \left(\underbrace{(x_i - (x_i^T w_1) w_1)^T w_2}_{\text{Residue}} \right) w_2, \dots \right\}$$

$$= \left\{ x_i - (x_i^T w_1) w_1 - \left(x_i^T w_2 - \underbrace{(x_i^T w_1) \cdot w_1^T w_2}_{=0} \right) w_2, \dots \right\}$$

$$= \left\{ x_i - (x_i^T w_1) w_1 - \underbrace{(x_i^T w_2) w_2}, \dots \right\}$$

Residues after d-rounds

$$\forall i \quad x_i - \left((x_i^T w_1) w_1 + (x_i^T w_2) w_2 + \dots + (x_i^T w_d) w_d \right) = \vec{0} \in \mathbb{R}^d$$

$$\text{for } x_i = (\underline{x}_i^T \underline{w}_1) \underline{w}_1 + (\underline{x}_i^T \underline{w}_2) \underline{w}_2 + \dots + (\underline{x}_i^T \underline{w}_d) \underline{w}_d.$$

What have we gained?

- If data lives in a "low" dimensional linear sub-space, then residues become 0 much earlier than d rounds.

Example: Say dataset is such that after 3-rounds, residues become 0.

$$\text{Dataset} = \{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^{100}$$

$$\forall i \quad x_i = \underline{(x_i^T w_1) w_1} + \underline{(x_i^T w_2) w_2} + \underline{(x_i^T w_3) w_3} \\ \{w_1, w_2, w_3\} \in \mathbb{R}^{100}$$

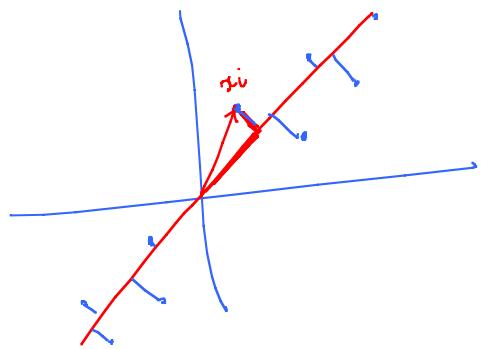
REP
 $\{w_1, w_2, w_3\}$
 $x_i \rightarrow \begin{bmatrix} x_i^T w_1 & x_i^T w_2 & x_i^T w_3 \end{bmatrix} \in \mathbb{R}^3$
 ↳ Data point specific
 ↳ Common Rep dataset

Original: $100 \times n$
 $100 \times 100 = \underline{\underline{10000}}$
 $\boxed{d \times n}$

Now!: $3 \times 100 + 3n$
 $3 \times 100 + 3 \times 100 = \underline{\underline{600}}$

$$\boxed{d \times k + k \times n}$$

Question: What if data "approximately" lies in a low-dimensional space?

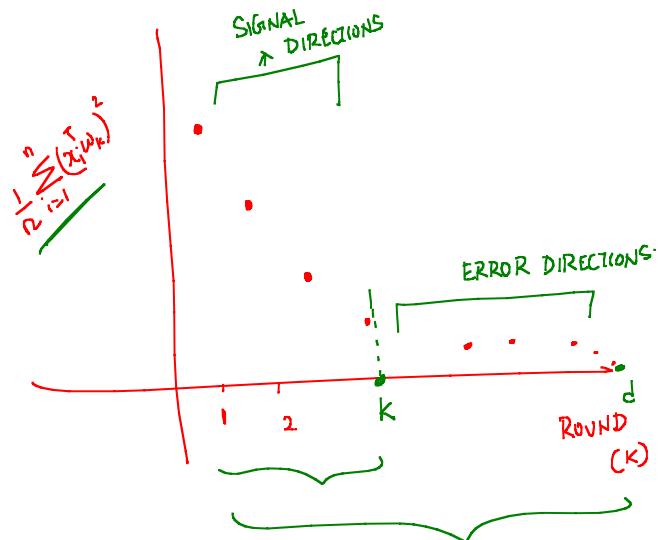


Pythagorean

$$\text{For any } w \in \mathbb{R}^d, \text{ s.t } \|w\|_2 = 1 \\ + i \|x_i\|^2 = \|x_i - (x_i^T w) w\|^2 + \| (x_i^T w) w \|^2$$

$$\frac{1}{n} \sum_{i=1}^n \|x_i\|^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - (x_i^T w) w\|^2 + \underbrace{\frac{1}{n} \sum_{i=1}^n (x_i^T w) w \|^2}_{\text{As large as possible.}}$$

"Larger the value of $\underbrace{\frac{1}{n} \sum_{i=1}^n (x_i^T w)^2}_{}$, the better the fit"



ENTER LINEAR ALGEBRA

$$\max_w \quad w^T C w$$

w: $\|w\|_2 = 1$

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

C → Covariance matrix.

Soln: w_1 is eigenvector corresponding to the "largest" eigenvalue of C [HILBERT min-max theorem]

In fact $\{w_1, \dots, w_d\}$, the eigenvectors of C form an orthonormal basis.

$w_k \rightarrow$ Best line one can obtain in round k

What do eigenvalues of C mean?

We know

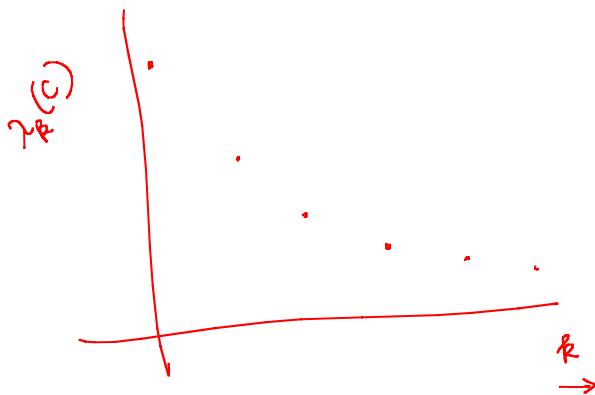
$$C w_1 = \lambda_1 w_1$$

$$w_1^T C w_1 = w_1^T (\lambda_1 w_1) = \lambda_1$$

$$\lambda_1 = w_1^T C w_1 = w_1^T \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) w_1$$

$$\boxed{\lambda_1 = \frac{1}{n} \sum_{i=1}^n (x_i^T w_1)^2}$$

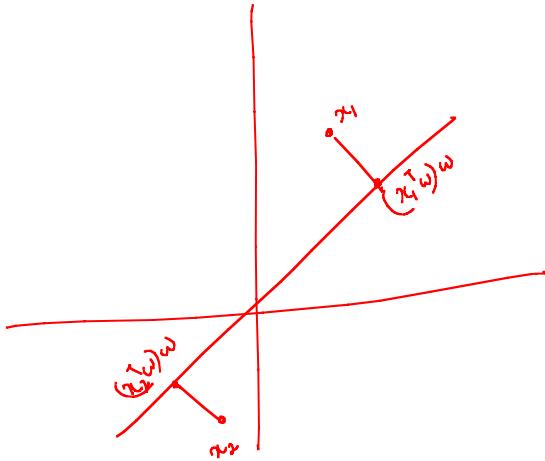
term we used earlier



Rule of thumb for # dimensions

$$\left(\frac{\sum_{i=1}^d \lambda_i(C)}{\sum_{i=1}^d \lambda_i(C)} \right) \geq 0.95$$

usually in practice



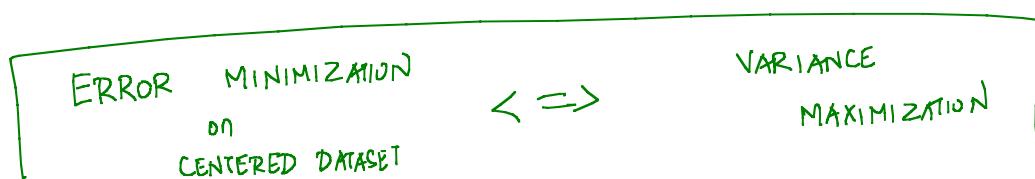
$$\{(x_1^T w), \dots, (x_n^T w)\}$$

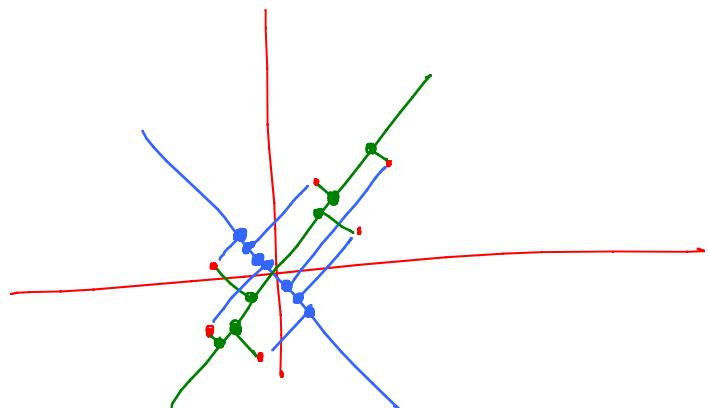
Average?

$$\frac{1}{n} \sum_{i=1}^n (x_i^T w) = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i \right)^T w}_{0^T w} \quad [\text{for a centered dataset}]$$

Variance

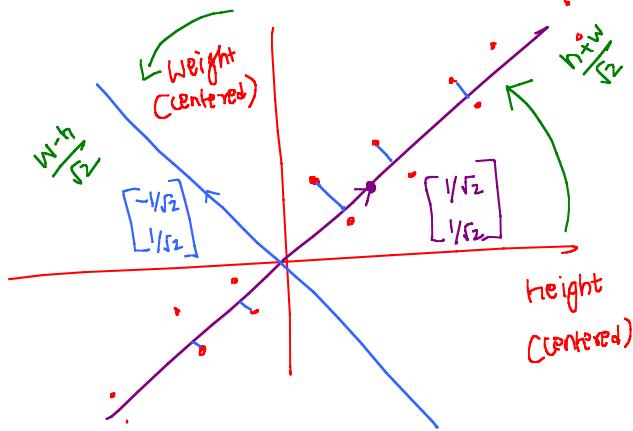
$$\frac{1}{n} \sum_{i=1}^n (x_i^T w - \underbrace{\text{mean}}_{=0})^2 = \boxed{\frac{1}{n} \sum_{i=1}^n (x_i^T w)^2}$$



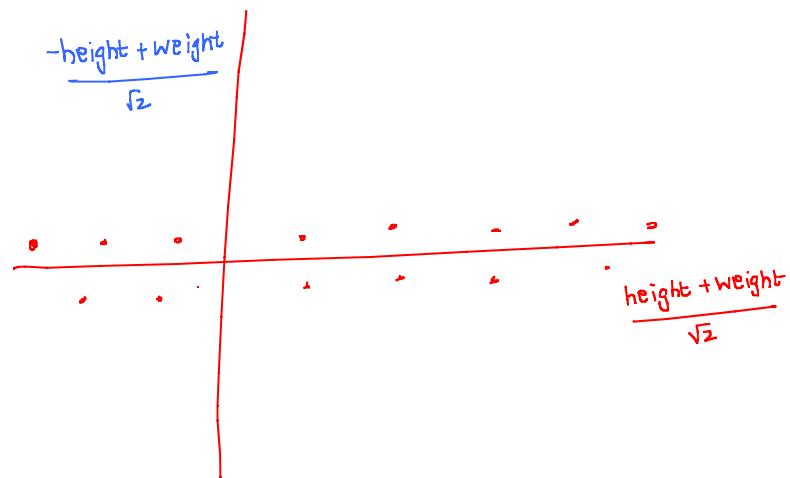


Want directions where
Projections don't "crowd up"
i.e., variance is not small.

ONE MORE EXAMPLE



height gives some
information about
weight.



$(h+w)/\sqrt{2}$ does not give any information
about $\frac{w-h}{\sqrt{2}}$.

“de-correlated”

PRINCIPAL COMPONENT ANALYSIS - $\{w_1, \dots, w_k\}$
PRINCIPAL COMPONENT

"DIMENSIONALITY REDUCTION"

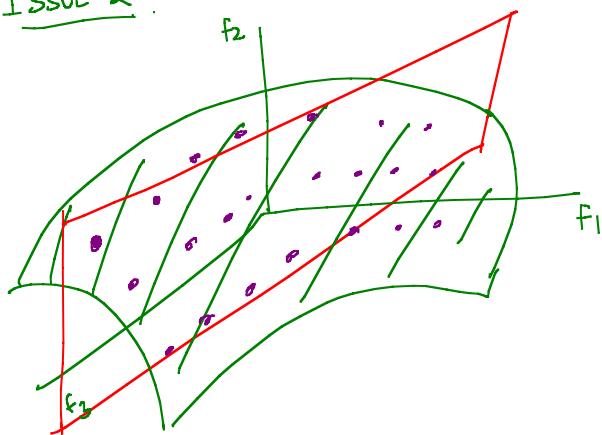
PCA finds combination of features that are
de-correlated. [loosely speaking independent of each other].

"EIGENFACES"

ISSUES / CONCERNS with PCA

- TIME COMPLEXITY - Finding the Eigen vectors and Eigen values.
 $C \in \mathbb{R}^{d \times d}$ Typically $O(d^3)$
- Issue when d is large
Example: Face recognition (Eigenfaces)

- ISSUE 2 :



Data may not necessarily live in a low-dimensional LINEAR subspace.

SURPRISING RESULT :

Same solution to both issues!

Issue 1:large d $[d \gg n]$

features

data points

$$X = \begin{bmatrix} | & | & | & & | \\ x_1 & x_2 & x_3 & \dots & x_n \\ | & | & | & & | \end{bmatrix}$$

$$X \in \mathbb{R}^{d \times n}$$

$$C = \frac{1}{n} \left(\sum_{i=1}^n x_i x_i^T \right)$$

$$XX^T = \begin{bmatrix} | & & | \\ x_1 & \dots & x_n \\ | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n x_i x_i^T$$

exercise:
Show this

$\Rightarrow C = \frac{1}{n} XX^T$

Let w_k be the eigenvector corresponding to the k^{th} largest eigenvalue of $C \xrightarrow{\text{e-value}} (\lambda_k)$

$$C w_k = \lambda_k w_k \quad [\text{by definition}]$$

$$\left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) w_k = \lambda_k w_k$$

$$w_k = \frac{\sum_{i=1}^n (\lambda_i^T w_k) \cdot x_i}{\lambda_k}$$

w_k is a LINEAR COMBINATION of data points!

$$\begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \\ | & & | \\ \alpha_{k1} & & \alpha_{kn} \end{bmatrix}$$

$$= \sum_{i=1}^n \alpha_{ki} x_i$$

$$w_k = x \alpha_k \quad \text{for some } \alpha_k \in \mathbb{R}^n$$

How to get α_k ?

Some Algebra:

$$w_k = x \alpha_k$$

* k

$$C w_k = \lambda_k w_k$$

$$\left(\frac{1}{n} X^T \right) (x \alpha_k) = \lambda_k x \alpha_k$$

$$(X^T) \times \alpha_R = n \lambda_R \times \alpha_R$$

Pre multiply by X^T

$$\underline{X^T} (\underline{(X^T)} \times \alpha_R) = \underline{X^T} (n \lambda_R \times \alpha_R)$$

$$(\underline{X^T} X) (\underline{X^T} \alpha_R) = n \lambda_R (\underline{X^T} \alpha_R)$$

$$\begin{array}{|c|}\hline X \in \mathbb{R}^{d \times n} \\ X^T X \in \mathbb{R}^{n \times n} \\ \hline\end{array}$$

Call $\underline{X^T} X := K$

$$K^2 \alpha_R = n \lambda_R K \alpha_R$$

if we can find α_R that satisfies

$$K \alpha_R = (\underline{\lambda} \underline{\alpha}_R) \alpha_R \leftarrow$$

\hookrightarrow Eigen equation

we know $w_R = X \alpha_R$

$$w_R^T w_R = (\underline{X \alpha_R})^T (\underline{X \alpha_R}) = \underline{\alpha_R^T (X^T X) \alpha_R}$$

$$1 = \alpha_R^T K \alpha_R$$

LINEAR ALGEBRA FACT : The non-zero eigenvalues of $\frac{\underline{xx}^T}{\mathbb{R}^{d \times d}}$ and $\frac{\underline{x}^T x}{\mathbb{R}^{n \times n}}$ are exactly the same!

$$\underline{C} = \frac{1}{n} \underline{xx}^T$$

$$\text{Eigenvectors} = \{ \underline{w}_1, \dots, \underline{w}_k \}$$

$$\text{Eigenvalues} = \{ \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \}$$

$$\underline{\underline{xx}^T} = n \underline{C}$$

$$\text{Eigenvectors} = \{ \underline{w}_1, \dots, \underline{w}_k \}$$

$$\text{Eigenvalues} = \{ n\lambda_1 \geq n\lambda_2 \geq \dots \geq n\lambda_k \}$$

$$\frac{X^T X}{k}$$

Eigenvalues = $\{ \beta_1, \dots, \beta_k \}$
 $\|\beta_k\|^2 = 1$ & β_k

$$\text{Eigenvalues} = \{ n\lambda_1 \geq n\lambda_2 \geq \dots \geq n\lambda_k \}$$

$$k\beta_k = (n\lambda_k)\beta_k$$

IS $\beta_k = \alpha_k$?

$$\underbrace{\beta_R^T K \beta_R}_{= n\lambda_R} = \underbrace{\beta_R^T (n\lambda_R \beta_R)}_{= n\lambda_R} = n\lambda_R \underbrace{\beta_R^T \beta_R}_{= 1}$$

Set $\alpha_R := \frac{\beta_R}{\sqrt{n\lambda_R}}$ Now ① $K\alpha_R = (n\lambda_R)\alpha_R$

② $\alpha_R^T K \alpha_R = 1$

$$= \frac{\beta_R^T K \beta_R}{n\lambda_R} = \frac{n\lambda_R}{n\lambda_R} = 1$$

Input: $D = \{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$ $\boxed{d \gg n}$

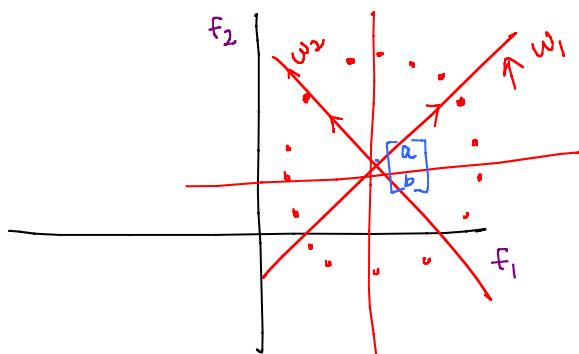
$$k_{ij} = x_i^T x_j$$

- Step 1: Compute $K = \boxed{x^T x}$ $K \in \mathbb{R}^{n \times n}$
- Step 2: Compute eigen decomposition of K corresponding to e-values
 e-vectors $\{P_1, \dots, P_k\}$
 e-values $\{n\lambda_1, \dots, n\lambda_k\}$ $\boxed{O(n^3)}$

- Step 3 set $d_k = \frac{P_k}{\sqrt{n\lambda_k}}$ $\forall k = 1, \dots, k$

- Step 4 $w_k = \boxed{x d_k}$ $\forall k$

Issue 2 → Non-linear relationships



What would PCA give? $\{w_1, w_2\}$ and both important.

Relation between features

$$(f_1 - a)^2 + (f_2 - b)^2 = r^2$$

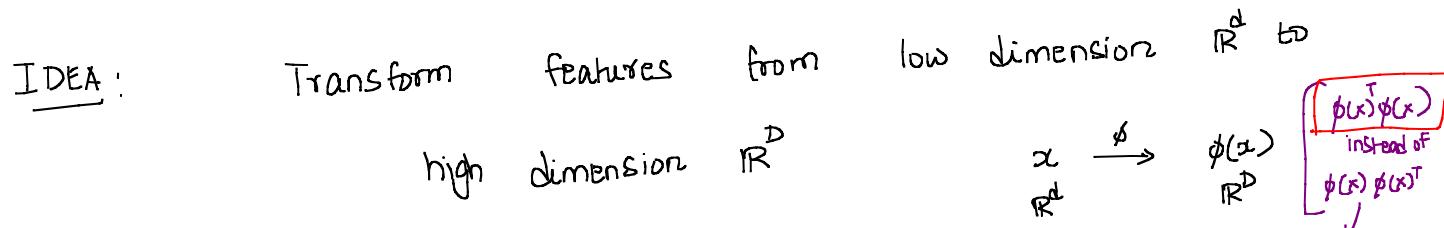
$$\Rightarrow \boxed{f_1^2 + a^2 - 2f_1 \cdot a + f_2^2 + b^2 - 2f_2 \cdot b - r^2 = 0} \quad \leftarrow \textcircled{*}$$

$$\begin{bmatrix} f_1 & f_2 \end{bmatrix}_{\mathbb{R}^2} \xrightarrow{\Phi} \begin{bmatrix} \phi(x) \\ 1 & f_1^2 & f_2^2 & f_1 f_2 & f_1 & f_2 \end{bmatrix}_{\mathbb{R}^6}$$

Let $u \in \mathbb{R}^6$

$$\begin{bmatrix} a^2 + b^2 - c^2 & 1 & 1 & 0 & -2a & -2b \end{bmatrix}$$

Each datapoint satisfies $\boxed{\phi(x)^T u = 0} \equiv \boxed{\Phi u = 0}$ \Rightarrow the datapoints lie in a LINEAR subspace of \mathbb{R}^6 !



- We already know how to handle the case when $D \gg n$

ISSUE - $\phi(x) \in \mathbb{R}^D$ may be too hard to compute.

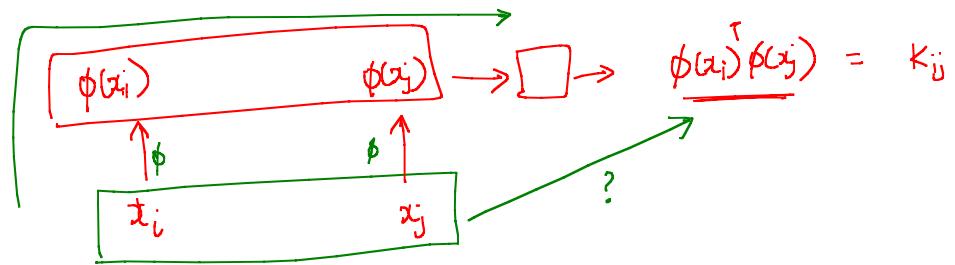
$$[f_1 \ f_2 \ f_3 \ f_4]$$

$$\phi \xrightarrow{\text{Cubic relations}} \begin{bmatrix} 1 & f_1 & f_2 & f_3 & f_4 \\ & \underbrace{f_1 f_2}_{4c_1} & \underbrace{f_1 f_3}_{4c_2} & \dots & \underbrace{f_3 f_4}_{4c_2} \\ & & & & \underbrace{f_1 f_2 f_3}_{4c_3} & \underbrace{f_1 f_2 f_4}_{4c_3} & \dots \end{bmatrix}$$

$$1 + 4 + 6 + 4$$

In general, d features, $\leq p^m$ power

$$\sum_{i=0}^p d c_i \approx O(d^p)$$



Example.

$$x = [f_1 \ f_2]$$

$$x' = [g_1 \ g_2]$$

Consider the function

$$\begin{aligned} (x^T x + 1)^2 &= \left([f_1 \ f_2] \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} + 1 \right)^2 \\ &= (f_1 g_1 + f_2 g_2 + 1)^2 \\ &= f_1^2 g_1^2 + f_2^2 g_2^2 + 1 + 2 f_1 g_1 f_2 g_2 + 2 f_1 g_1 + 2 f_2 g_2 \end{aligned}$$

$$\begin{bmatrix} x^T x + 1 \end{bmatrix}^2$$

$\xrightarrow{\quad}$

$$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \quad \begin{bmatrix} g_1 \\ g_2 \end{bmatrix}$$

$$= \begin{bmatrix} f_1^2 & f_2^2 & 1 & \sqrt{2}f_1f_2 & \sqrt{2}f_1 & \sqrt{2}f_2 \end{bmatrix} \begin{bmatrix} g_1^2 \\ g_2^2 \\ 1 \\ \sqrt{2}g_1g_2 \\ \sqrt{2}g_1 \\ \sqrt{2}g_2 \end{bmatrix}$$

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

where

$$\phi(x) = \phi\left(\begin{bmatrix} a \\ b \end{bmatrix}\right) = \begin{bmatrix} a^2 \\ b^2 \\ 1 \\ f_2ab \\ f_2a \\ \sqrt{2}b \end{bmatrix}$$

INSIGHT:

$(x^T x + 1)^2$ computes the dot-product in a "transformed space".

$$\begin{bmatrix} f_1 & f_2 \end{bmatrix} \rightarrow \begin{bmatrix} f_1^2 & f_2^2 & 1 & \sqrt{2}f_1f_2 & \sqrt{2}f_1 & \sqrt{2}f_2 \end{bmatrix}$$

$$\begin{bmatrix} g_1 & g_2 \end{bmatrix} \rightarrow \begin{bmatrix} g_1^2 & g_2^2 & 1 & \sqrt{2}g_1g_2 & \sqrt{2}g_1 & \sqrt{2}g_2 \end{bmatrix}$$

$$\begin{array}{c} x \in \mathbb{R}^2 \\ \downarrow \\ \phi(x) \in \mathbb{R}^6 \end{array}$$

WE MANAGED TO COMPUTE

$$\phi(x)^T \phi(x')$$

WITHOUT EXPLICITLY

COMPUTING $\phi(x)$

MORE EXAMPLES.

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

$$R(x, x') = \underline{(x^T x' + 1)^p}$$

for some $p \geq 1$

→ can be shown to be a "valid" function

i.e., $\exists \phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ such that

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

EXERCISE:

compute the explicit
 ϕ for $p=3$ and
 $p=4$

ONE MORE EXAMPLE

$$R(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \quad \text{for some } \sigma > 0$$

RADIAL BASIS
FUNCTION.

→ Can be shown to be a "Valid" map.

→ Interestingly, ϕ in this case maps x to an "infinite" dimensional space.

[Technicalities aside, can think of this as mapping a point to a "function" and dot-products between functions become integrals.]

KERNEL FUNCTION

Any function $R: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ which is a "valid" map is called a kernel function

$$R(x, x') = (x^T x' + 1)^p \Rightarrow \text{POLYNOMIAL KERNEL}$$

$$R(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \Rightarrow \text{RADIAL BASIS / GAUSSIAN KERNEL}$$

Given a function $R: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, how can we say its a valid kernel?

METHOD 1: Exhibit a map ϕ explicitly.
 [might be hard sometimes]

METHOD 2: MERCER'S THEOREM (Informal)

$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid kernel

$$\phi(x)^T \phi(x)$$

$$\underline{\phi(x) \phi(x)^T}$$

All e-values of K
 are non-negative.

if and only if

(a) K is symmetric i.e., $K(x, x') = K(x', x)$.

(b) For any dataset $\{x_1, \dots, x_n\}$, the matrix
 $K \in \mathbb{R}^{n \times n}$ where $K_{ij} = K(x_i, x_j)$ is POSITIVE SEMI DEFINITE

KERNEL PCA

- Input - $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$; Kernel $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$.

- Step 1: Compute $K \in \mathbb{R}^{n \times n}$ where
$$k_{ij} = k(x_i, x_j) \quad \forall i, j$$
- $\xrightarrow{\hspace{1cm}}$
 - Center the kernel"
- Step 2: Compute B_1, \dots, B_l eigen vectors of K .
 $n\lambda_1 \geq \dots \geq n\lambda_l$ Eigenvalues

and normalize to get

$$\alpha_u = \frac{B_u}{\sqrt{n\lambda_u}}$$

X. Step 3:

$$w_k = \underline{\phi(x)} \alpha_k \rightarrow \text{Defeats the purpose because it needs } \underline{\phi(x)}$$

We cannot "reconstruct" the eigenvectors of the covariance matrix.

But we can still compute the "compressed" representation.

$$\phi(x_i)^T w_k = \phi(x_i)^T \left(\sum_{j=1}^n \phi(x_j) \alpha_{kj} \right)$$

$$= \sum_{j=1}^n \alpha_{kj} \phi(x_i)^T \phi(x_j) = \sum_{j=1}^n \alpha_{kj} K_{ij}$$

For downstream tasks, this is usually good enough.

• Modified Step 3 : Compute $\sum_{j=1}^n \alpha_{kj} K_{ij} + b$

$x_i \in \mathbb{R}^d \rightarrow \left[\sum_{j=1}^n \alpha_{1j} K_{ij}, \sum_{j=1}^n \alpha_{2j} K_{ij}, \dots, \sum_{j=1}^n \alpha_{dj} K_{ij} \right]$

DETAILS: (Centering the kernel)

Given $K \in \mathbb{R}^{n \times n}$ $K_{ij} = R(x_i, x_j) + i, j$

$$\begin{array}{c} \{x_1, \dots, x_n\} \\ \downarrow \\ \{\phi(x_1), \dots, \phi(x_n)\} \\ \{ \phi(x_1) - \bar{x}, \dots, \phi(x_n) - \bar{x} \} \\ \bar{x} = \frac{1}{n} \sum_i \phi(x_i) \end{array} \quad \left| \begin{array}{c} (\phi(x_i) - \bar{x})^T (\phi(x_j) - \bar{x}) \\ (\phi(x_i) - \bar{x})^T \phi(x_j) \\ -\bar{x}^T \phi(x_j) \\ -\bar{x}^T \phi(x_i) \\ + \bar{x}^2 \end{array} \right\} \quad \text{Create a new kernel } K^c \leftarrow \text{centered.}$$

$$K_{ij}^c = K_{ij} - \underline{\theta_i \underline{1}_j^T} - \underline{1_i^T \theta_j} + \underline{P \underline{1}_i \underline{1}_j^T}$$

where

$$\theta_i = \frac{1}{n} \sum_{k=1}^n K_{ik} + i$$

$$P = \frac{1}{n^2} \sum_{i,j} K_{ij}$$

Unsupervised Learning

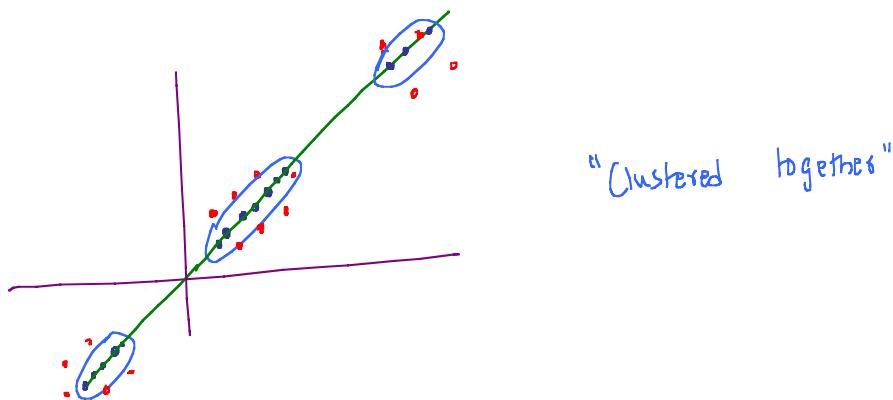
L Representation learning

L PCA

L kernel PCA

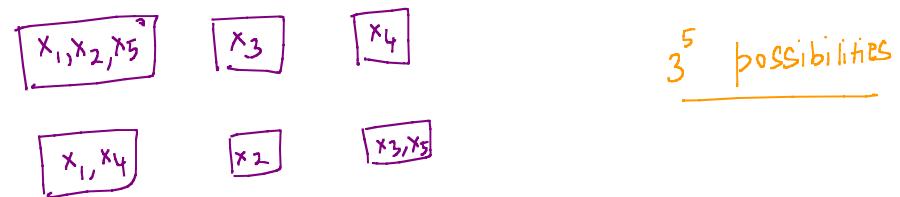
L

clustering (today)



Goal: $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$
Partition the data k different clusters

Example: $\{x_1, x_2, x_3, x_4, x_5\}$ $k = 3$



$x_1, x_2, \dots, x_n \leftarrow$ DATA POINTS

$z_1, z_2, \dots, z_n \leftarrow$ CLUSTER INDICATOR

$z_i \in \{1, \dots, k\}$

Question: Given a cluster assignment, how good is it?

$$F(z_1, \dots, z_n) = \sum_{i=1}^n \|x_i - \underline{\mu}_{z_i}\|_2^2$$

↳ Mean/average of z_i^{th} cluster

$$\underline{\mu}_k = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i = k)}{\sum_{i=1}^n \mathbb{1}(z_i = k)}$$

$$\mathbb{1}(u) = \begin{cases} 1 & \text{if } u \text{ is true} \\ 0 & \text{o/w} \end{cases}$$

Example:

$$\begin{array}{ccccc} \textcircled{x}_1 & \textcircled{x}_2 & \textcircled{x}_3 & \textcircled{x}_4 & \textcircled{x}_5 \\ z_1=1 & z_2=2 & z_3=1 & z_4=1 & z_5=2 \end{array} \quad k=2$$

$$M_1 = \frac{x_1 + x_3 + x_4}{3} \quad ; \quad M_2 = \frac{x_2 + x_5}{2}$$

Goal

$$\min_{\{z_1, \dots, z_n\}} \sum_{i=1}^n \|x_i - M_{z_i}\|^2$$

Too many possibilities! (k^n)



NP-HARD

LLOYD'S ALGORITHM / K-MEANS ALGORITHM

INITIALIZATION

$$z_1^0, z_2^0, \dots, z_n^0 \in \{1, \dots, k\}$$

→ UNTIL CONVERGENCE

- COMPUTE MEANS

$$\forall R \quad \mu_R^t = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i^t = R)}{\sum_{i=1}^n \mathbb{1}(z_i^t = R)}$$

- RE-ASSIGNMENT STEP

$$\forall i \quad z_i^{t+1} = \arg \min_R \|x_i - \mu_R^t\|_2^2$$

mean of
[if the current
assignment is
smallest, then
don't reassign]

FACT: LLOYD'S ALGORITHM CONVERGES. ← Good news.

- Converged solution may not be "optimal"
- But produces "reasonable" clusters in practice.

QUESTIONS

- CONVERGENCE?
- NATURE OF CLUSTERS?
- INITIALIZATION?
- CHOICE OF K?

LLYOD'S ALGORITHM

QUESTIONS:

1. CONVERGENCE
 2. NATURE OF CLUSTERS
 3. INITIALIZATION
 4. CHOICE OF K.
- 

CONVERGENCE

- Does Lloyd's Algorithm Converge? YES.

PROOF:

FACT 1: Let $x_1, x_2, \dots, x_\ell \in \mathbb{R}^d$

$$v^* = \arg \min_{v \in \mathbb{R}^d} \sum_{i=1}^{\ell} \|x_i - v\|^2$$

Answer: $v^* = \frac{1}{\ell} \sum_{i=1}^{\ell} x_i$

[view this objective
as a function of v ,
take derivative, set to 0
and solve]

- Say we are at iteration t of Lloyd's algorithm.

- CURRENT ASSIGNMENT.

$$z_1^t, z_2^t, \dots, z_n^t \in \{1, \dots, k\}$$

$m_k^t \leftarrow$ mean of cluster k in iteration t

- Say we update our assignments to

$$z_1^{t+1}, z_2^{t+1}, \dots, z_n^{t+1} \in \{1, \dots, k\}$$

$$\sum_{i=1}^n \|x_i - \underline{M}_{z_i^t}\|^2 < \sum_{i=1}^n \|x_i - \underline{M}_{\underline{z}_i^t}\|^2$$

[By
Algorithmic
choice]

Mean of
 cluster where
 x_i wants to
 go to

$F(z_1^t, \dots, z_n^t)$
 Mean of
 current clusters
 where x_i is
 assigned to.

$$\sum_{i=1}^n \|x_i - \underline{M}_{z_i^{t+1}}\|^2 \leq \sum_{i=1}^n \|x_i - \underline{M}_{\underline{z}_i^{t+1}}\|^2$$

$F(z_1^{t+1}, \dots, z_n^{t+1})$

$\underline{M}_{\underline{z}_i^{t+1}}$

$$= \sum_{i=1}^n \sum_{k=1}^K \|x_i - m_k^{t+1}\|^2 \cdot \mathbb{1}(z_i^{t+1} = k)$$

For every k

$$\sum_{i \in C_k} \|x_i - m_k^{t+1}\|^2 \leq \sum_{i \in C_k} \|x_i - v\|^2$$

[FACT 1]

$$\leq \sum_{i \in C_k} \|x_i - m_k^t\|^2$$

\Rightarrow The objective function strictly reduces after each re-assignment.

$$F(z_1^{t+1}, \dots, z_n^{t+1}) < F(z_1^t, \dots, z_n^t)$$

\Rightarrow There are only "FINITE" number of assignments

\Rightarrow Algorithm must converge!

NATURE OF CLUSTERS

K=2.

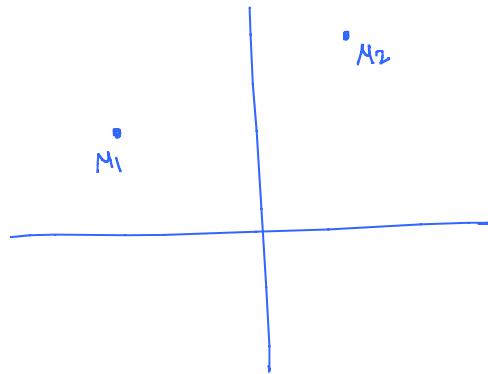
-
Lloyd's algorithm produces 2 clusters
with means μ_1 and μ_2

-
What can we say about points assigned to
Cluster 1 ?

By Algorithm's Construction

For points in Cluster 1,

$$\|x - \mu_1\|^2 \leq \|x - \mu_2\|^2$$

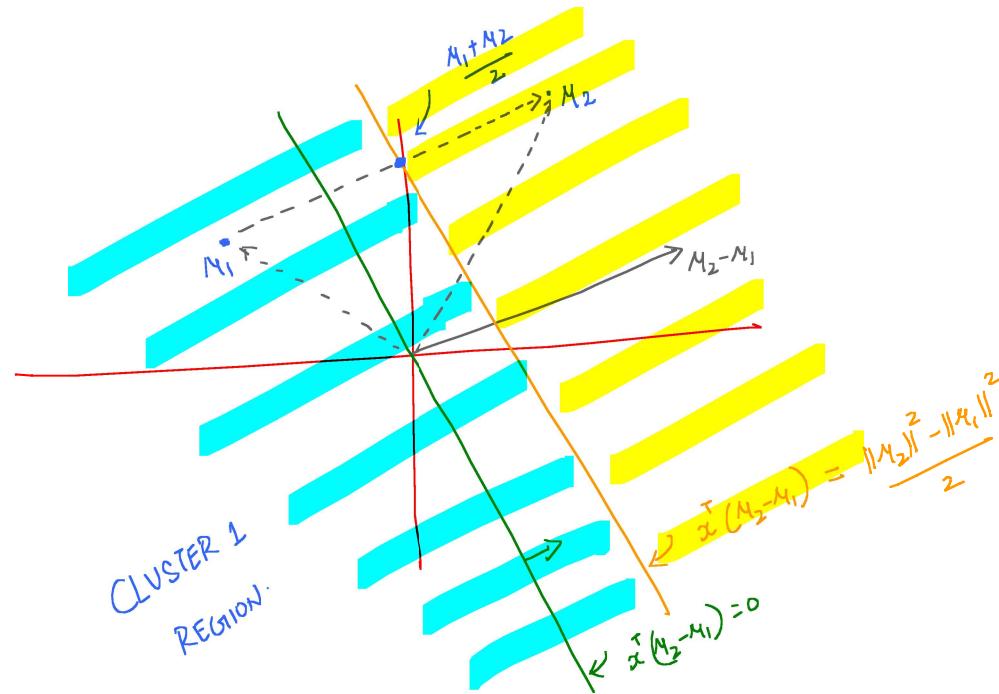


* x in
Cluster

$$\|x\|^2 + \|\mu_1\|^2 - 2x^\top \mu_1 \leq \|x\|^2 + \|\mu_2\|^2 - 2x^\top \mu_2$$

$$\boxed{x^\top (\mu_2 - \mu_1) \leq \frac{\|\mu_2\|^2 - \|\mu_1\|^2}{2}}$$

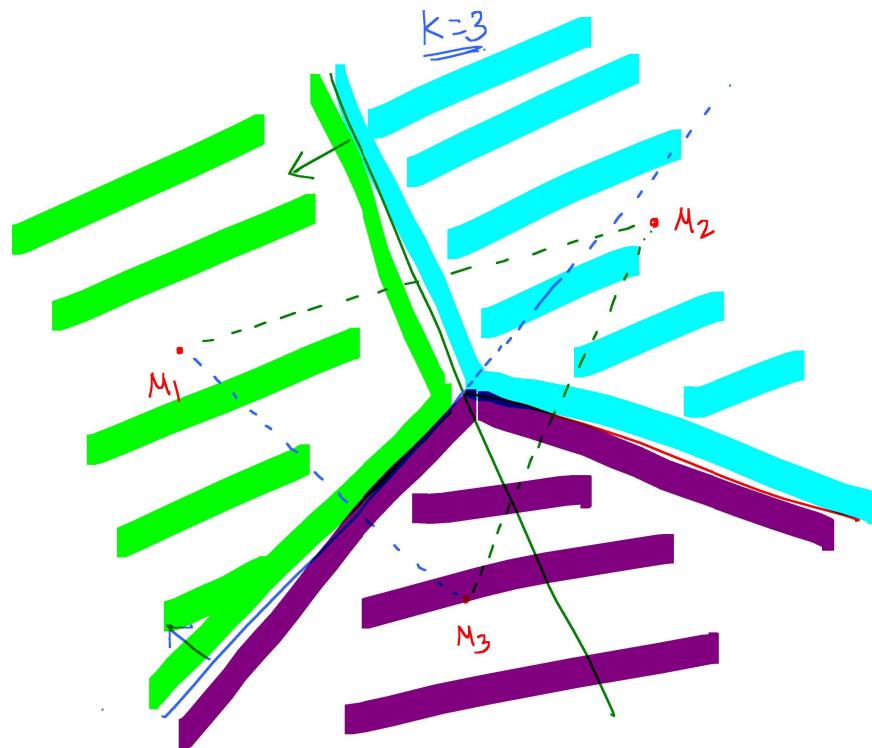
$$x^\top w \leq b$$



$$\underline{x}^T (\underline{M}_2 - \underline{M}_1) \leq \frac{\|\underline{M}_2\|^2 - \|\underline{M}_1\|^2}{2}$$

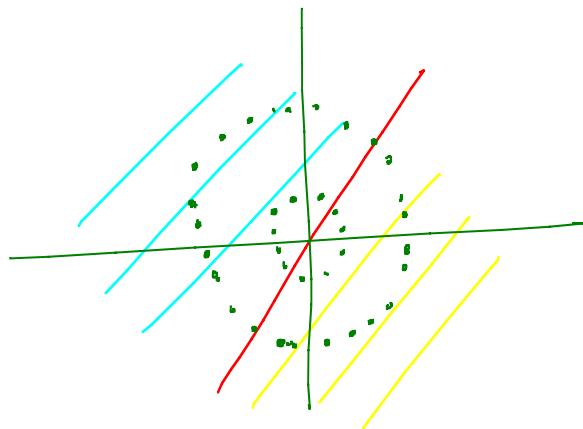
$$\underline{x} = \frac{\underline{M}_1 + \underline{M}_2}{2}$$

$$\begin{aligned} & (\underline{M}_1 + \underline{M}_2)^T (\underline{M}_2 - \underline{M}_1) \\ &= \frac{\|\underline{M}_2\|^2 - \|\underline{M}_1\|^2}{2} \end{aligned}$$



cluster regions
are
intersection of
Half Spaces

VORONOI REGIONS



How to fix?

- KERNELIZE
K-MEANS
- SPECTRAL
CLUSTERING



INITIALIZATION

POSSIBILITIES

- Pick k-means uniformly at random from the dataset

K-MEANS ++

→ Choose first mean M_1^0 uniformly at random
from $\{x_1, \dots, x_n\}$

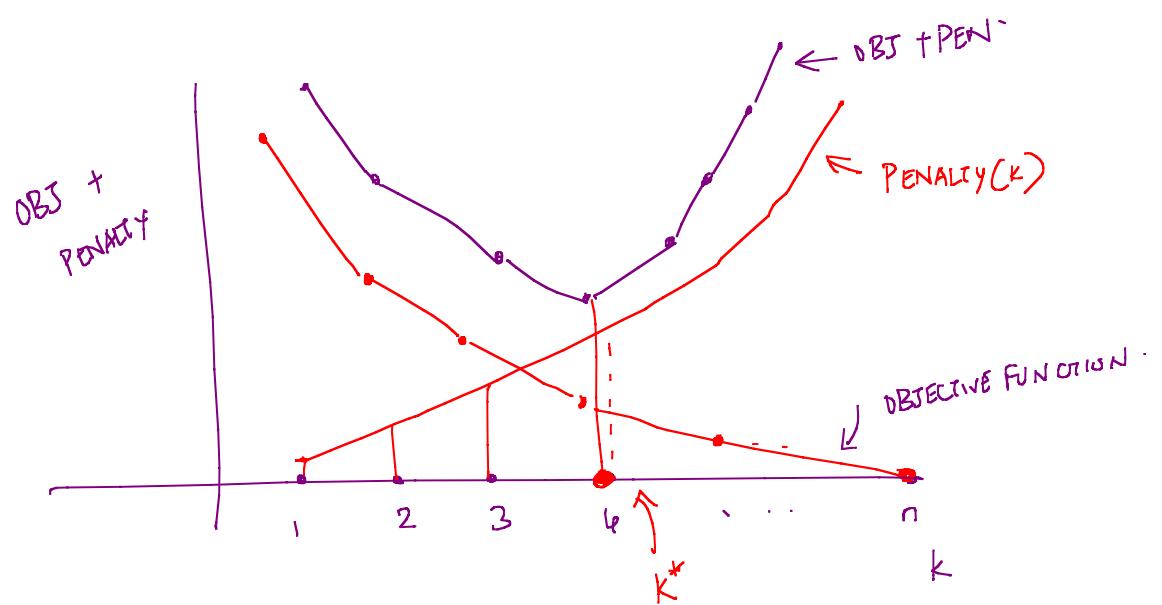
CHOICE OF K

$$\rightarrow F(z_1, \dots, z_n) = \sum_{i=1}^n \|x_i - M_{z_i}\|^2 \quad k=n$$

→ want k to be as small as possible.

→ Penalize large values of k .

Find k that has the
Smallest objective function value + Penalty(k)



Some Common Criterion

A.I.C - Akaike Information Criterion

$$[2K - 2 \underbrace{\log(L(\theta^*))}_{\text{blue bracket}}]$$

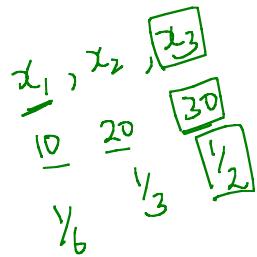
B.I.C - Bayesian Information Criterion

$$[K \underbrace{\log(n)}_{\text{blue bracket}} - 2 \underbrace{\log(L(\theta^*))}_{\text{blue bracket}}]$$

- CONVERGENCE - YES
- NATURE OF CLUSTERS - VORNOI REGIONS
- INITIALIZATION - K-MEANS++
- CHOICE OF K - OBJ + PENALTY(K).

\rightarrow for $l = 2, \dots, k$

Choose μ_l^0 probabilistically proportional to slope



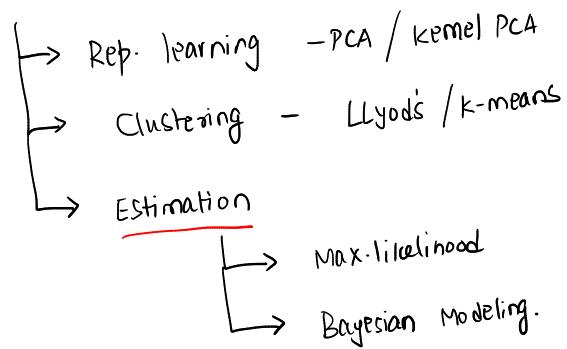
$$\text{for } x \quad S(x) = \min_{j=1, \dots, l-1} \|x - \mu_j^0\|^2$$

GUARANTEE

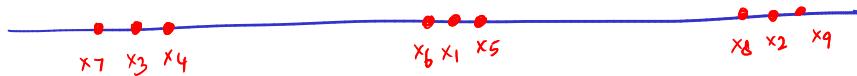
$$\mathbb{E} \left[\sum_{i=1}^n \|x_i - \mu_{z_i}\|^2 \right] \leq O(\log k) \left[\min_{z_1, \dots, z_n} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2 \right]$$

over randomness of algorithm

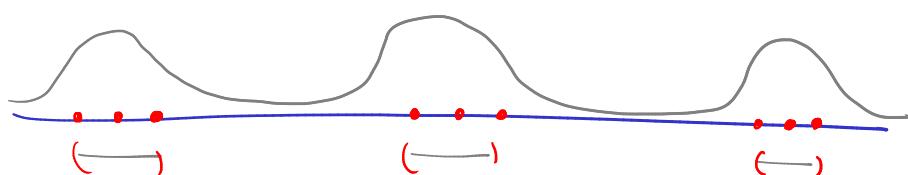
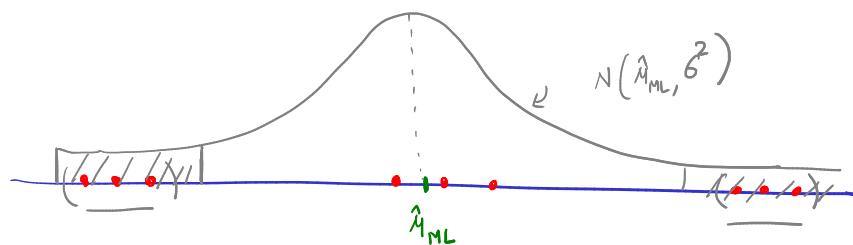
Unsupervised learning



Today: Slightly complicated data



- What could be a good "generative" story?



Want a density like above to explain this data.

A NEW GENERATIVE MODEL

MIXTURE OF GAUSSIANS

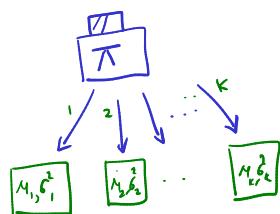
STEP 1: Pick which mixture a data point comes from.

STEP 2: Generate data point from that mixture.

STEP 1: Generate a mixture component among $\{1, \dots, k\}$ $z_i \in \{1, \dots, k\}$

$$P(z_i = l) = \pi_l \quad \left[\begin{array}{l} \sum_{i=1}^k \pi_i = 1 \\ 0 \leq \pi_i \leq 1 \end{array} \right]$$

STEP 2: Generate $x_i \sim N(\mu_{z_i}, \sigma_{z_i}^2)$



$\{x_1, \dots, x_n\} \rightarrow \text{OBSERVED}$

$\{z_1, \dots, z_n\} \rightarrow \text{UNOBSERVED/LATENT.}$

Latent variable Model

Parameters: $\pi = [\pi_1, \pi_2, \dots, \pi_k]$

Total: $2k + k - 1$

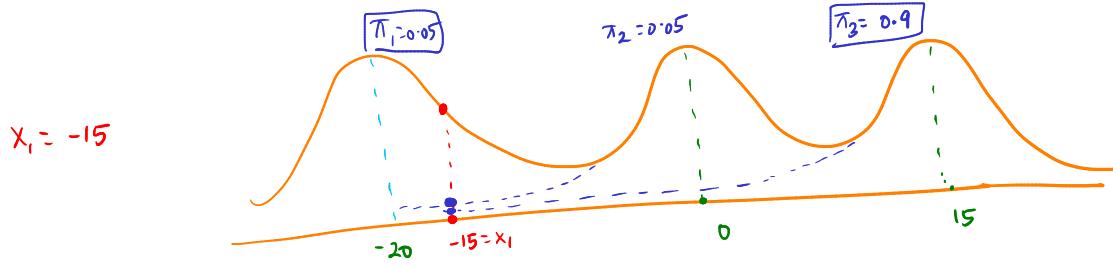
$\mu_R, (\mu_R, \sigma_R^2)$

$3k - 1$

Max. Likelihood for GMM

$$L \left(\underbrace{\mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2, \pi_1, \dots, \pi_k}_{\text{parameters}} ; x_1, \dots, x_n \right) = \prod_{i=1}^n f_{\text{mix}}(x_i ; \underbrace{\mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2}_{\text{parameters}}, \pi_1, \dots, \pi_k)$$

$$= \prod_{i=1}^n \left[\underbrace{\sum_{R=1}^k \pi_R}_{\text{Mixing proportions}} \cdot \underbrace{\frac{f(x_i; \mu_R, \sigma_R^2)}{\text{NORMAL/Gaussian Density}}} \right]$$



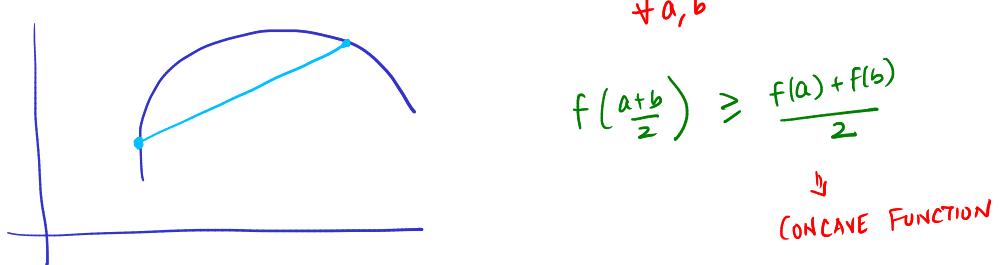
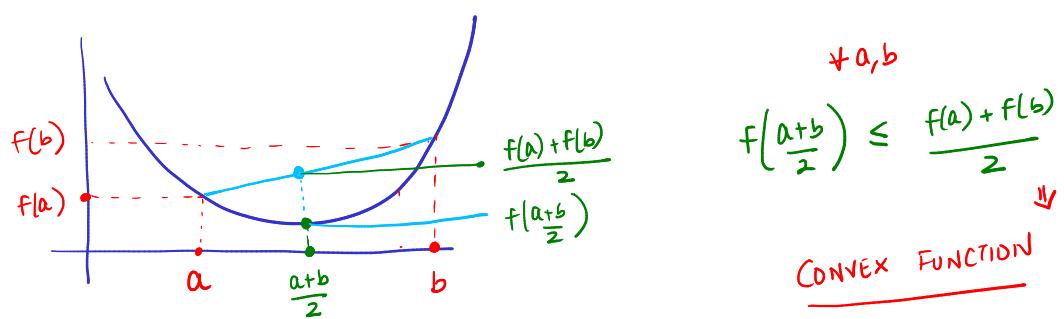
$$L(\theta) = \prod_{i=1}^n \left[\sum_{k=1}^K \pi_k \frac{e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k} \right]$$

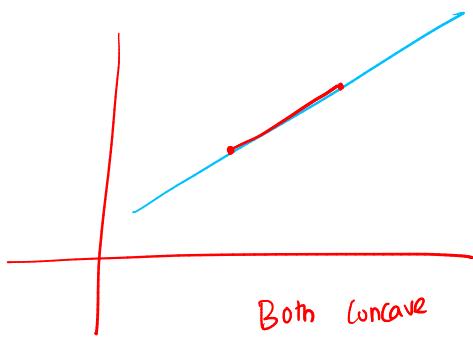
↑
all parameters

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k \frac{e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k} \right) - \star$$

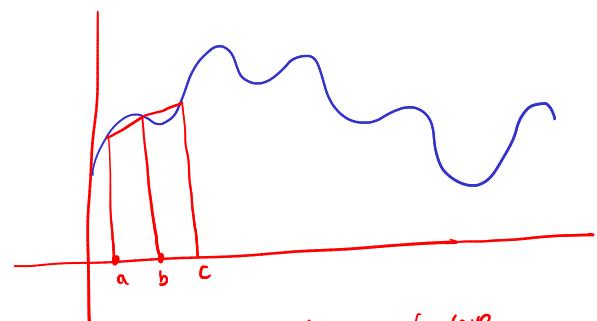
- Not possible to solve this analytically.
- Need an alternate way to solve this efficiently!

Quick detour - Convex functions



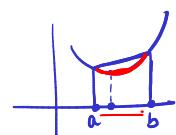


Both Concave
and
Convex



Neither Concave
nor Convex

$$f\left(\frac{1}{2} \cdot a + \frac{1}{2} b\right) \leq \frac{1}{2} f(a) + \frac{1}{2} f(b).$$



$$\Rightarrow f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda) f(b) \quad \lambda \in [0, 1]$$

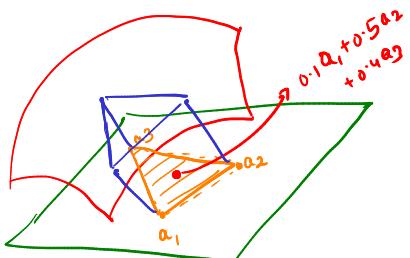
For Concave

$$f(\lambda_1 a_1 + \lambda_2 a_2 + \dots + \lambda_k a_k) \geq \lambda_1 f(a_1) + \dots + \lambda_k f(a_k)$$

$$\begin{cases} \sum_{i=1}^k \lambda_i = 1 \\ 0 \leq \lambda_i \leq 1 \end{cases}$$

JENSEN'S INEQUALITY.

$$f\left(\sum_{k=1}^K \lambda_k a_k\right) \geq \sum_{k=1}^K \lambda_k f(a_k)$$



- Log is a concave function!

[Why?
Exercise]

- How can we exploit Jensen's for performing maximum likelihood.

Recall

$$\textcircled{*} \log L(\theta) = \sum_{i=1}^n \log \left(\underbrace{\sum_{k=1}^K}_{=} \left(\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma_k} \right) \right)$$

- INTRODUCE for every data point i , two parameters
 $\{ \lambda_1^i, \dots, \lambda_K^i \}$ st $\forall i \sum_{k=1}^K \lambda_k^i = 1, 0 \leq \lambda_k^i \leq 1 + \epsilon$

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \lambda_k^i \underbrace{\left(\frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k^i} \right)}_{\text{constant}} \right)$$

By Jensen's

$$\begin{aligned} \log L(\theta) &\geq \text{modified_log } L(\theta, \lambda) \\ &= \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left(\frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k^i} \right) \end{aligned}$$

- Note that the above modified log likelihood gives a lower bound for the true log likelihood at θ

for any choice of λ

$$\left\{ \begin{array}{l} \lambda_1^1, \dots, \lambda_K^1 \\ \lambda_1^2, \dots, \lambda_K^2 \\ \vdots \\ \lambda_1^n, \dots, \lambda_K^n \end{array} \right\} \quad \left\{ \begin{array}{l} \mu_1, \dots, \mu_K \\ \sigma_1^2, \dots, \sigma_K^2 \\ \pi_1, \dots, \pi_K \end{array} \right\}$$

- But what are we gaining?

Key insight:

- if we fix λ , it is easy to maximize w.r.t θ

- if we fix θ , it is easy to maximize w.r.t λ .

Fix λ and maximize over θ

$$\max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \left[\log \left(\frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k^i} \right) \right] / \lambda_k^i$$

$$= \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \left[\lambda_k^i \log \pi_k - \lambda_k^i \frac{(x_i - \mu_k)^2}{2\sigma_k^2} - \lambda_k^i \log \sqrt{2\pi} \frac{\sigma_k}{\lambda_k^i} \right]$$

Take derivative w.r.t μ, σ to get

$$\hat{\mu}_k^{MML} = \frac{\sum_{i=1}^n \lambda_k^i x_i}{\sum_{i=1}^n \lambda_k^i} \quad | \quad \hat{\sigma}_k^{MML} = \frac{\sum_{i=1}^n \lambda_k^i (x_i - \hat{\mu}_k^{MML})^2}{\sum_{i=1}^n \lambda_k^i}$$

$$\max_{\pi_1, \dots, \pi_K} \sum_{i=1}^n \left(\sum_{k=1}^K \lambda_k^i \log \pi_k \right)$$

$$\text{s.t. } \sum_k \pi_k = 1; \pi_k \geq 0$$

Can solve using method of Lagrange multipliers

$$\hat{\pi}_k^{MML} = \frac{\sum_{i=1}^n \lambda_k^i}{n}$$

Fixing λ , we get

$$\hat{\mu}_k^{MML} = \frac{\sum_{i=1}^n \lambda_k^i x_i}{\sum_{i=1}^n \lambda_k^i}$$

$$\hat{\sigma}_k^{MML} = \frac{\sum_{i=1}^n \lambda_k^i (x_i - \hat{\mu}_k^{MML})^2}{\sum_{i=1}^n \lambda_k^i}$$

$$\hat{\pi}_k^{MML} = \frac{\sum_{i=1}^n \lambda_k^i}{n}$$

- Fix θ and maximize λ

$$\sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left(\frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k^i} \frac{1}{\sqrt{2\pi} \sigma_k} \right)$$

$$= \sum_{i=1}^n \left[\sum_{k=1}^K \lambda_k^i \log(\alpha_{ik}) - \lambda_k^i \log \lambda_k^i \right]$$

where $\alpha_{ik} = \frac{1}{\pi_k} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}$

Fix any i ,

$$\max_{\lambda_1, \dots, \lambda_K} \sum_{k=1}^K [\lambda_k^i \log(\alpha_{ik}^i) - \lambda_k^i \log \lambda_k^i]$$

st $\sum_{k=1}^K \lambda_k^i = 1 \quad 0 \leq \lambda_k^i \leq 1$

can be solved analytically

$$\hat{\lambda}_k^{i \text{ MLE}} = \frac{\left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \right) \cdot \frac{P(z_i = k | x_i)}{\sum_{l=1}^K \left(\frac{1}{\sqrt{2\pi}\sigma_l} e^{-\frac{(x_i - \mu_l)^2}{2\sigma_l^2}} \cdot \pi_l \right)}}{P(z_i)}$$

$P(z_i = k | x_i) = P(x_i | z_i = k) \cdot P(z_i = k)$

ALGORITHM - E-M ALGORITHM (1970's Dempster et al.)

iteration

→ Initialize $\theta^0 = \left\{ \begin{array}{l} \mu_1^0, \dots, \mu_K^0, \\ \sigma_1^2, \dots, \sigma_L^2, \\ \pi_1^0, \dots, \pi_K^0 \end{array} \right\}$

usually comes from Lloyd's

→ until convergence $(\|\theta^{t+1} - \theta^t\| \leq \epsilon)$

Tolerance parameter

- EM produces "soft clustering"

- EM takes variances into account.



- EM clusters need not be Voronoi regions!

$$\lambda^{t+1} = \arg \max_{\lambda} \text{modified-} \log L(\theta^t, \lambda)$$

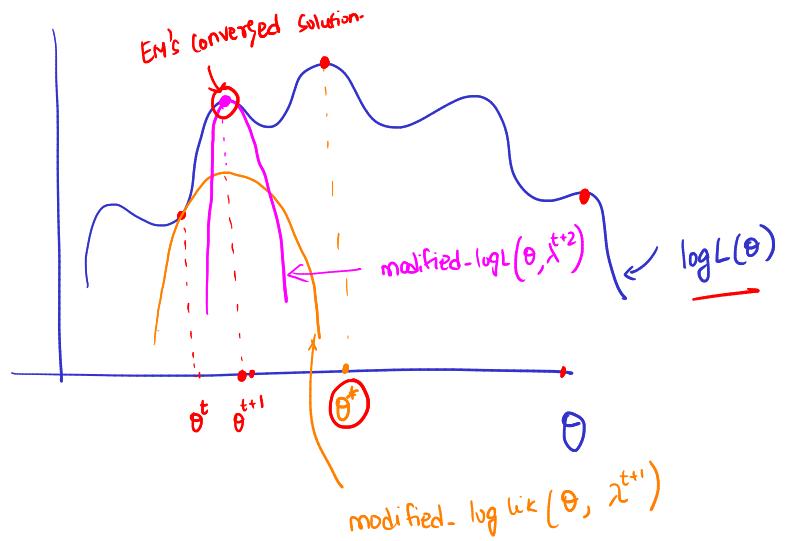
$$\theta^{t+1} = \arg \max_{\theta} \text{modified-} \log L(\theta, \lambda^{t+1})$$

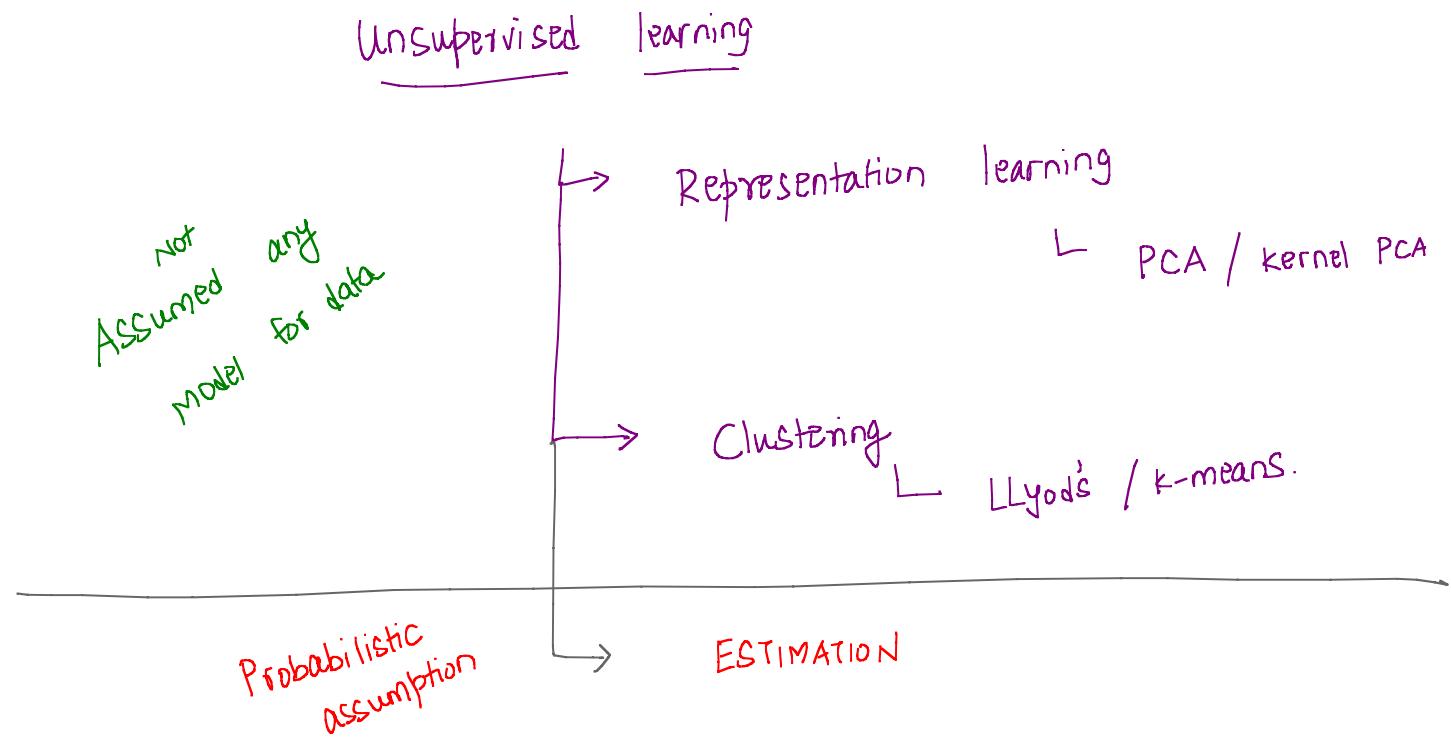
↓

Maximization Step Expectation Step

→ end.

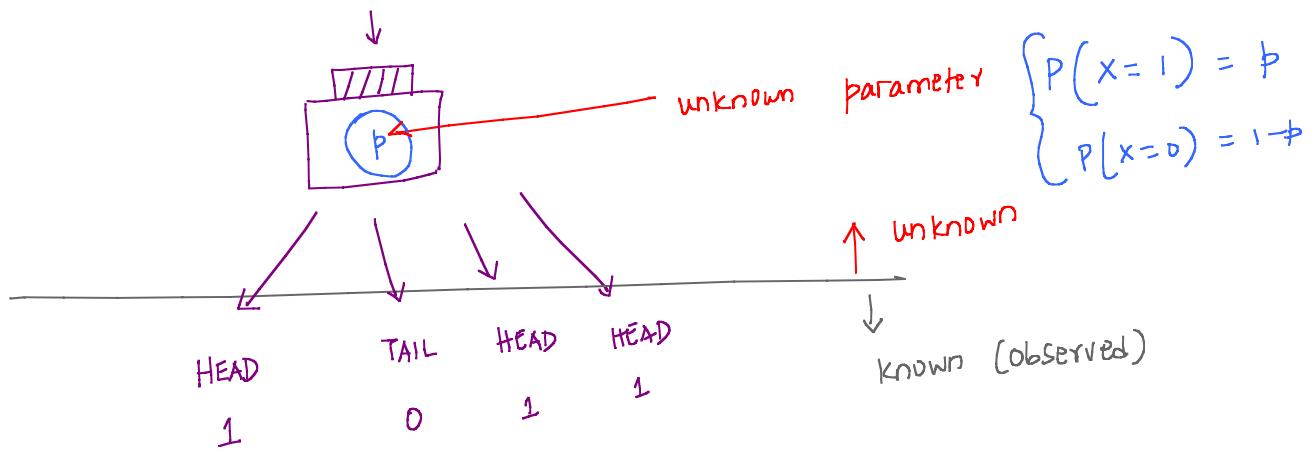
EM Converges to
a local maximum
of log likelihood.





"There is some probabilistic mechanism that generates data about which we don't know "something". Given data,
find/estimate what we don't know"

- OBSERVE data
- "ASSUME" a probabilistic model that generates data
- ESTIMATE unknown parameters using data.

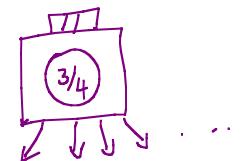


OBSERVE

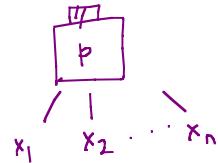
$$\{1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1\}$$

ESTIMATE:

$$a_{12} = \boxed{2/4} \leftarrow$$



ASSUMPTIONS



OBSERVATIONS ARE

(i) INDEPENDENT

(ii) IDENTICALLY DISTRIBUTED

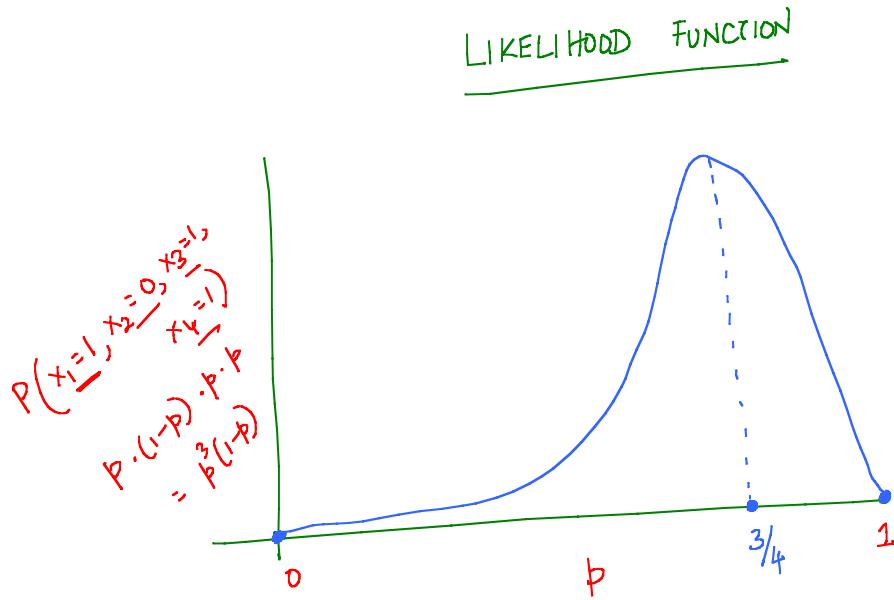
$$\begin{aligned} \text{GUESS} &= 2/3 ? \quad Y \\ &= 0.0001 ? \quad Y \\ &= 0 ? \quad N \\ &= 1 ? \quad N \end{aligned}$$

INDEPENDENCE

$$P(x_i | x_j) = P(x_i) \quad * i \neq j$$

IDENTICAL DISTRIBUTION

$$P(x_i=1) = P(x_j=1) = P \quad * i, j$$



(Eg) : $\{1, 0, 1, 1\}$

$3/4$

FISHER'S PRINCIPLE OF
MAXIMUM LIKELIHOOD

$$\begin{aligned}
 L(p; \{x_1, x_2, \dots, x_n\}) &= P(x_1, x_2, \dots, x_n; p) \\
 &= P(x_1; p) \cdot P(x_2; p) \cdots P(x_n; p)
 \end{aligned}$$

[Independence]

$$= \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

{ if $x_i = 1 \Rightarrow p(1-p) = p$
 if $x_i = 0 \Rightarrow p(1-p) = 1-p$

$$\begin{aligned}
 \hat{p}_{ML} &= \arg \max_p \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \\
 &= \arg \max_p \log \left(\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \right)
 \end{aligned}$$

[log is monotonic increasing]

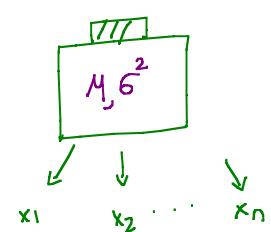
$$= \underset{p}{\operatorname{arg\ max}} \sum_{i=1}^n [x_i \log p + (1-x_i) \log (1-p)]$$

Take derivative of $\log L(p)$, set it to 0 to get

$$\hat{p}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

Fraction of 1's

$$\text{Data} = \{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}$$



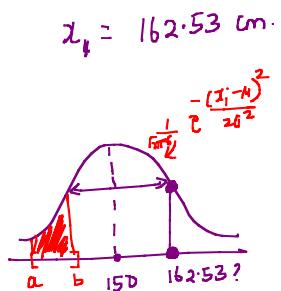
$$x_i \sim \text{Gaussian}(N, \sigma^2)$$

$\mu \rightarrow \text{unknown}; \sigma^2 \rightarrow \text{known}$.

$$\underline{L}(\underline{\mu}, \underline{\sigma^2}, \{x_1, \dots, x_n\}) = P(x_1, \dots, x_n; \underline{\mu}, \underline{\sigma^2})$$

$$= \prod_{i=1}^n P(x_i; \underline{\mu}, \underline{\sigma^2}) \quad \times$$

O



$$\underline{L}(\underline{\mu}, \underline{\sigma^2}, \{x_1, \dots, x_n\}) = f_{x_1, \dots, x_n}(\underline{x}_i, \dots, \underline{x}_n; \underline{\mu}, \underline{\sigma^2})$$

$$= \prod_{i=1}^n f_{x_i}(\underline{x}_i; \underline{\mu}, \underline{\sigma^2})$$

$$= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\underline{x}_i - \mu)^2}{2\sigma^2}} \right]$$

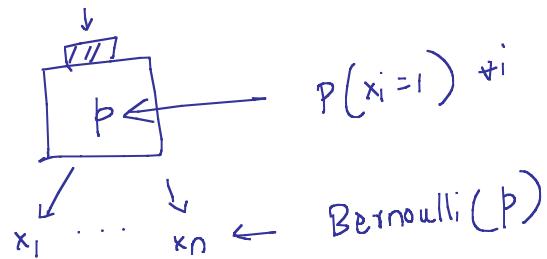
$$\log L(\mu, \sigma^2, \{x_1, \dots, x_n\}) = \sum_{i=1}^n \left[\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$\hat{\mu}_{ML} = \underset{\mu}{\operatorname{argmax}} \sum_{i=1}^n - (x_i - \mu)^2$$

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

Consider the coin examples.

BAYESIAN
MODELING



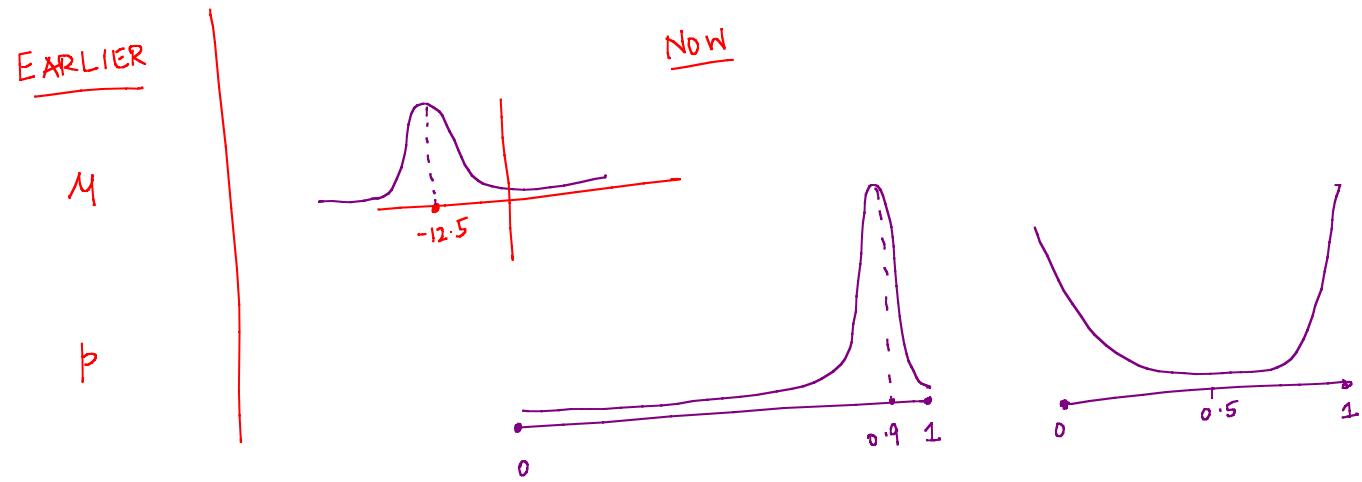
Let's say
someone says:

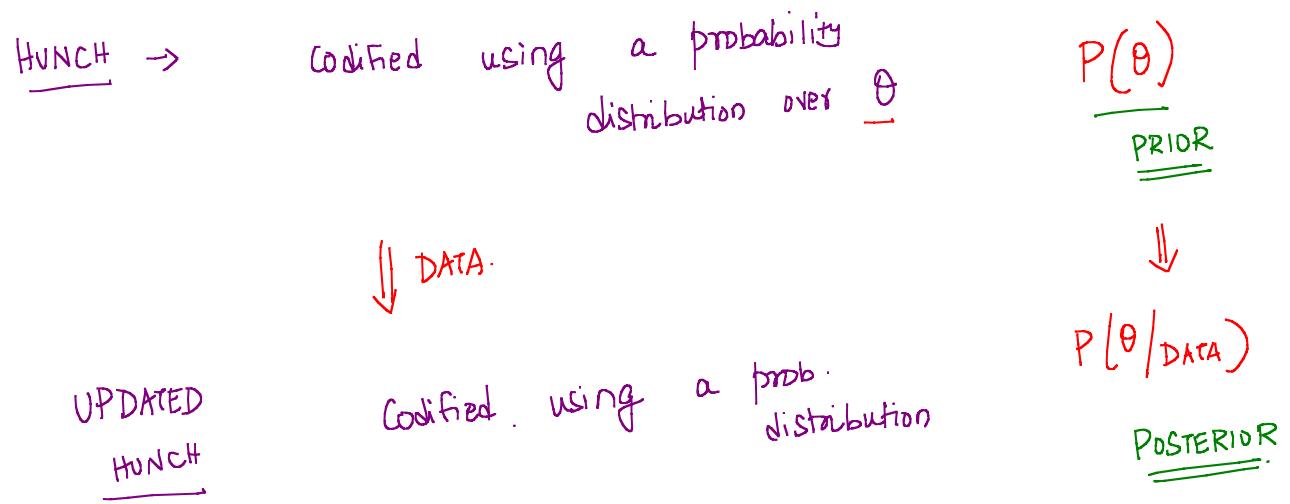
"I believe the bias p is somewhere close to 1"

We may have "HUNCH" about parameters

Goal: Incorporate "hunch/belief" about parameters of interest
into the estimation procedure.

APPROACH: Think of the parameter to estimate as a "random" variable.





Bayes law
 $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$

$A \rightarrow$ Parameters θ

$B \rightarrow$ DATA $\{x_1, \dots, x_n\}$

LIKELIHOOD

$$P(\theta | \{x_1, \dots, x_n\})$$

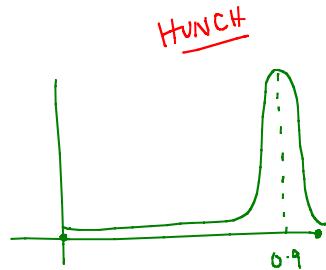
POSTERIOR

$$= \left(\frac{P(\{x_1, \dots, x_n\} | \theta)}{P(\{x_1, \dots, x_n\})} \right) \cdot \frac{P(\theta)}{\text{PRIOR}}$$

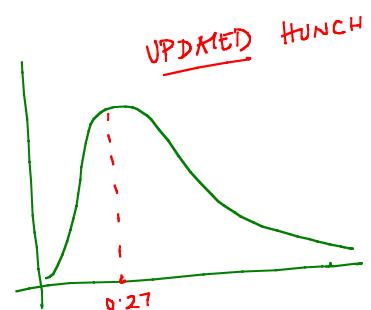
EVIDENCE

DOES not
depend on θ .

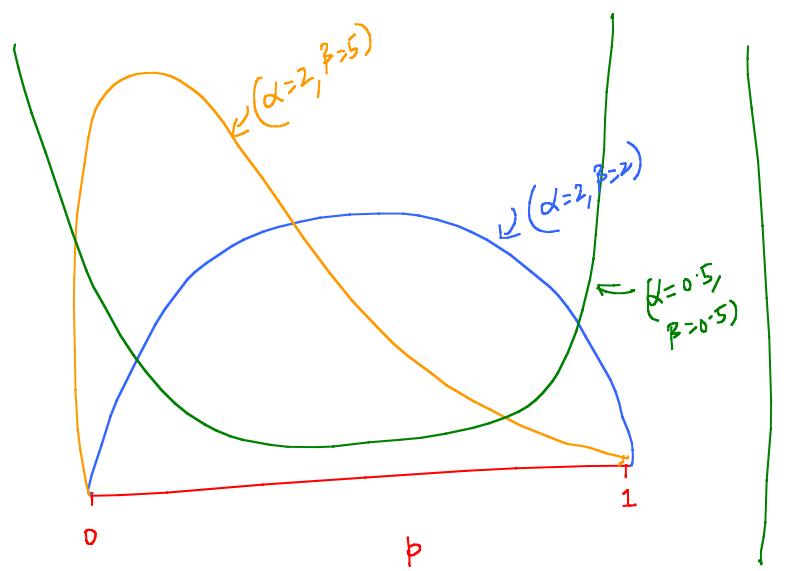
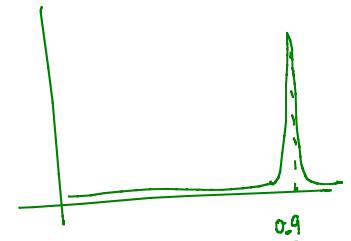
EXAMPLE



DATA \Rightarrow $\{0, 0, 0, 0, 1, 1, 0, 0, 0\}$



→ {1, 1, 1, 1, 0, 1, 1, 1}



DATA - Bernoulli (p)
PRIOR ? $P(\theta)$

BETA PRIOR
 $f(p; \alpha, \beta)$

$$\left\{ \frac{p^{\alpha-1} (1-p)^{\beta-1}}{Z} \quad p \in [0, 1] \right.$$

$$P(\theta | \text{DATA}) \propto P(\text{DATA} | \theta) \cdot P(\theta)$$

$$f_{\theta | \text{DATA}}(p) \propto \left[\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \right] \cdot \begin{bmatrix} \alpha-1 & \beta-1 \\ p & 1-p \end{bmatrix}$$

Likelihood

$$f_{\theta | \text{DATA}}(p) \propto p^{\sum x_i + \alpha - 1} (1-p)^{\sum (1-x_i) + \beta - 1}$$

↳ same functional form
as the PRIOR!!

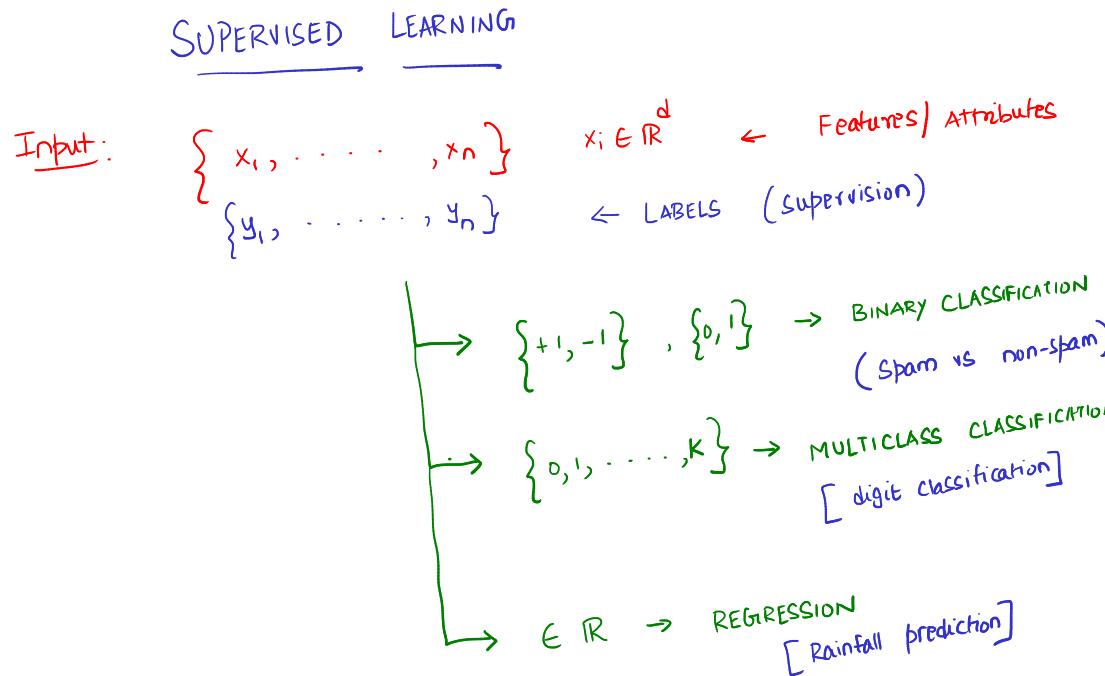
$$\text{BETA PRIOR } (\alpha, \beta) \xrightarrow[\text{Bernoulli}]{\text{DATA}} \text{BETA POSTERIOR } (\alpha + n_H, \beta + n_T)$$

one possible guess

$$\frac{\alpha + n_h}{\alpha + n_h + \beta + n_t} = \frac{\alpha + n_h}{(\alpha + \beta) + n}$$

$$\underline{\underline{E[Posterior]}} = \underline{\underline{E[\text{Beta}(\alpha+n_h, \beta+n_t)]}} =$$

MAP Estimator - Maximum A posteriori Estimator
 \hat{p}_{MAP}



REGRESSION

INPUT / TRAINING DATA $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$
 $\{y_1, \dots, y_n\}$ $y_i \in \mathbb{R}$

Goal! Learn $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$

- How do we measure "goodness" of a function $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$

- $\text{error}(f_i) = \sum_{i=1}^n (f_i(x_i) - y_i)^2$

- How small can this error be? 0

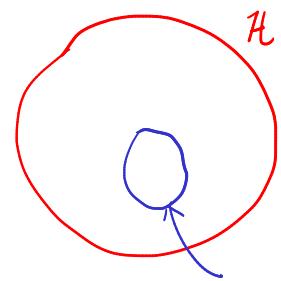
- Which f_i achieves zero error?

$$f_i(x_i) = y_i \quad \forall i$$

What is the problem?

- By "memorizing", we can get zero error on training data
- What we care is about test performance.
- Impose "STRUCTURE" to reduce search space.

SIMPLEST STRUCTURE - LINEAR STRUCTURE



Set of all functions
 $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$.

$$\mathcal{H}_{\text{linear}} = \left\{ \begin{array}{l} f_{\mathbf{w}}: \mathbb{R}^d \rightarrow \mathbb{R} \\ f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \end{array} \right. \text{ s.t. } \text{ for } \mathbf{w} \in \mathbb{R}^d \}$$

GOAL:

$$\min_{f_w \in \mathcal{H}_{\text{linear}}} \sum_{i=1}^n (f_w(x_i) - y_i)^2$$

(or) equivalently -

$$\boxed{\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2}$$

LINEAR REGRESSION

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 = \|\mathbf{x}^T \mathbf{w} - \mathbf{y}\|_2^2$$

$$\mathbf{x}^T = \begin{bmatrix} \mathbf{x}_1 & \cdots \\ \mathbf{x}_2 & \cdots \\ \vdots & \ddots \\ \mathbf{x}_n & \cdots \end{bmatrix}_{n \times d} \quad \mathbf{w} = \begin{bmatrix} 1 \\ \mathbf{w} \\ 1 \end{bmatrix}_{d \times 1}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$\min_{\mathbf{w} \in \mathbb{R}^d} (\mathbf{x}^T \mathbf{w} - \mathbf{y})^T (\mathbf{x}^T \mathbf{w} - \mathbf{y}) \quad \leftarrow \begin{array}{l} \text{unconstrained} \\ \text{quadratic (in w)} \\ \text{optimisation problem} \end{array}$$

Solution: Take derivative (gradient) and set to zero.

$$f(\mathbf{w}) = (\mathbf{x}^T \mathbf{w} - \mathbf{y})^T (\mathbf{x}^T \mathbf{w} - \mathbf{y})$$

$$\nabla f(\mathbf{w}) = 2 (\mathbf{x} \mathbf{x}^T) \mathbf{w} - 2 \mathbf{x} \mathbf{y}$$

Solution satisfies

$$\boxed{(\mathbf{x} \mathbf{x}^T) \mathbf{w}^* = \mathbf{x} \mathbf{y}}$$

OBSERVATION

Like PCA, w^* depends on a "covariance" like matrix. But it also involves y .

$$w^* = (x x^T)^+ x y \quad \text{Pseudo-inverse.}$$

- has Lin reg closed form solution
- GEOMETRIC VIEW?
- COMPUTATIONAL CONSIDERATIONS?
- NON-LINEAR FEATURE \rightarrow LABEL RELATIONSHIP?
- PROBABILISTIC VIEW?

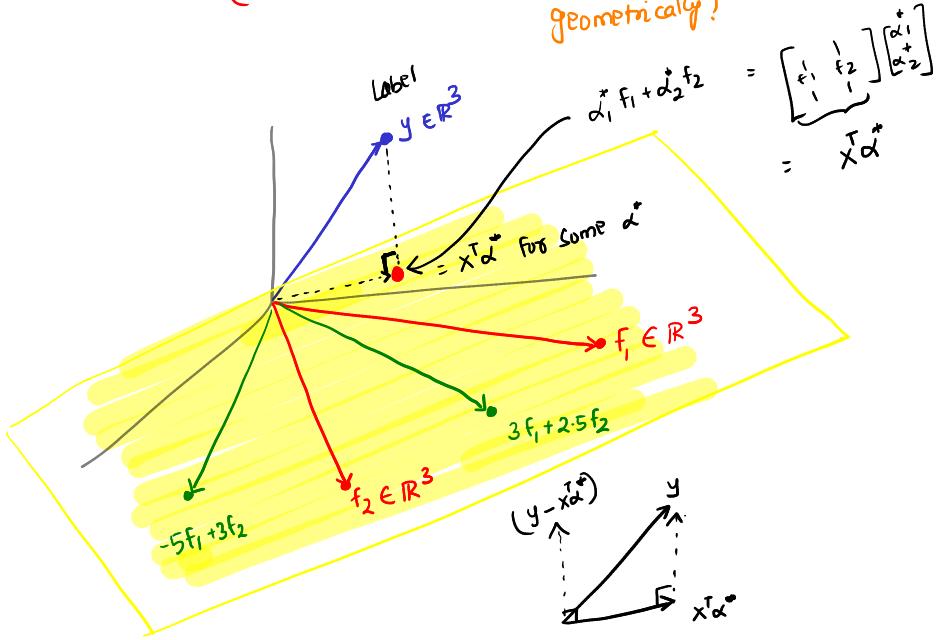
$$w^* = (x x^T)^+ x y \quad \leftarrow \text{How can we interpret this geometrically?}$$

f_1 height
 x_1
 x_2
 x_3

f_2 weight
 $[1]$
 $[1]$
 $[1]$

y_1
 y_2
 y_3

$d=2$
 $n=3$



$$(y - x^T \alpha^*)^T (x^T \alpha^*) = 0$$

$$\underbrace{y^T x^T \alpha^*}_{\text{LHS}} - \underbrace{\alpha^T (x x^T) \alpha^*}_{\text{RHS}} = 0 \quad \text{--- (1)}$$

Recall, $w^* = (x x^T)^+ x y$

Substituting $w^* = \alpha^*$ on L.H.S, we get

$$y^T x ((x x^T)^+ x y) - ((x x^T)^+ x y)^T (x x^T) ((x x^T)^+ x y)$$

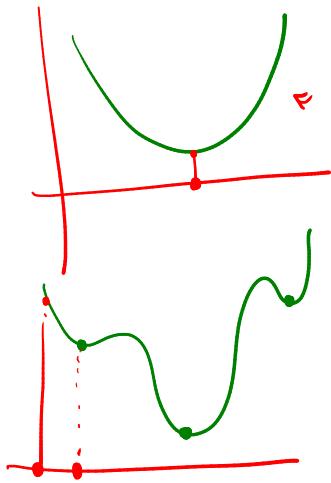
$$= 0 \quad \left[\begin{array}{l} x x^T \text{ is symmetric} \\ \text{use} \end{array} \right]$$

CONCLUSION: $x^T w^*$ is the PROJECTION of the Labels onto
the Subspace spanned by the features.

COMPUTATIONAL CONSIDERATIONS

$$w^* = (x x^T)^+ x y$$

inverse computation
is expensive if d is large
 $O(d^3)$



- We know w^* is the solution of an unconstrained optimization

- we can apply GRADIENT DESCENT.

$$w^{t+1} = w^t - \eta^t \nabla F(w^t)$$

\downarrow \downarrow
Scalar Gradient at w^t
Step size.

$$f(w) = \|xw - y\|^2 = \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\nabla f(w) = 2(x x^T)w - 2xy \quad [\text{verify this}]$$

Gradient descent update for Linear regression

$$w^{t+1} = w^t - \eta^t \left[\underbrace{2(x x^T)w - 2xy}_{\text{}} \right]$$

- what if n is large [hundreds of millions], might want to avoid $x x^T$

- How to adapt gradient descent?

STOCHASTIC GRADIENT DESCENT

for $t=1, \dots, T$

- At each step, sample a bunch of datapoints (R) uniformly at random from the set of all points

- PRETEND this sample is the entire dataset and take a gradient step w.r.t it

$$2(\tilde{x}\tilde{x}^T w^t - \tilde{x}\tilde{y})$$

↳ Manageable because
 $\tilde{x} \in \mathbb{R}^{d \times R}$ ↳ #Sampled points.

end

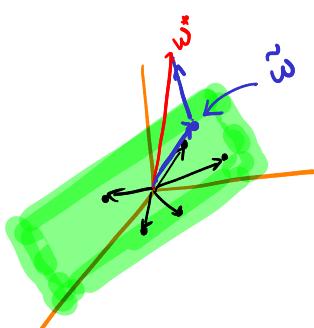
After T rounds, use

$$\underline{w}_{SGD}^T = \frac{1}{T} \sum_{t=1}^T w^t$$

as opposed to w^T as in
Standard Gradient descent

Guaranteed to converge to optima with high probability

NON-LINEAR REGRESSION



$$\underline{w}^* = (\underline{x}\underline{x}^T)^\dagger \underline{x}\underline{y}$$

- \underline{w}^* must lie in the span of data points.

$$\sum_{i=1}^n (\underline{w}^T \underline{x}_i - y_i)^2 = \sum_{i=1}^n (\tilde{\underline{w}}^T \underline{x}_i - y_i)$$

$$+ i \quad \underline{w}^T \underline{x}_i = (\tilde{\underline{w}} + \tilde{\underline{w}}_\perp)^T \underline{x}_i = \tilde{\underline{w}}^T \underline{x}_i + \underbrace{\tilde{\underline{w}}_\perp^T \underline{x}_i}_0$$

$$\omega^* = \boxed{x\alpha^*} \quad \text{for some } \alpha^* \in \mathbb{R}^n$$

$$= (x x^T)^\dagger x y$$

$$x\alpha^* = (x x^T)^\dagger x y$$

$$(x x^T) x \alpha^* = (\cancel{x} \cancel{x^T}) \cancel{(x)} \alpha^* x y$$

$$\cancel{(x x^T)} x \alpha^* = x y$$

$$\cancel{x^T} \cancel{(x x^T)} x \alpha^* = \cancel{x^T} x y$$

$$\underbrace{(x^T x)^2}_{:= K} \alpha^* = (x^T x) y$$

$$\Rightarrow K^2 \alpha^* = Ky$$

$$\alpha^* = \bar{K}^{-1} y \leftarrow \begin{matrix} \\ \in \mathbb{R}^{m \times n} \end{matrix}$$

KERNEL - REGRESSION

PREDICTION

for some $x_{\text{test}} \in \mathbb{R}^d$

$$\omega^* \phi(x_{\text{test}})$$

$$= \left(\sum_{i=1}^n \alpha_i^* \phi(x_i) \right) \phi(x_{\text{test}})$$

$$= \sum_{i=1}^n \alpha_i^* \underbrace{\kappa(x_i, x_{\text{test}})}$$

How important
is i point
towards
 w^*

How
similar
is x_{test}
to x_i

PROBABILISTIC VIEW OF LINEAR REGRESSION

$$x \in \mathbb{R}^d \quad y \in \mathbb{R}$$

$$\left\{ (x_1, y_1), \dots, (x_n, y_n) \right\} \quad \text{Dataset}$$

$$y/x \sim \underbrace{w^T x + \epsilon}_{\substack{\text{Unknown} \\ \text{but fixed}}} \quad \epsilon \sim \underbrace{N(0, \sigma^2)}_{\substack{\text{Noise} \\ \text{Gaussian.}}}$$

- Can view this as an "ESTIMATION" problem

- Solution approach - Maximum Likelihood.

$$\text{Likelihood} \quad L(\omega; \begin{matrix} x_1, \dots, x_n \\ y_1, \dots, y_n \end{matrix}) = \prod_{i=1}^n e^{-\frac{(\omega^T x_i - y_i)^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma}$$

$$\log L(\omega; \begin{matrix} x_1, \dots, x_n \\ y_1, \dots, y_n \end{matrix}) = \sum_{i=1}^n -\frac{(\omega^T x_i - y_i)^2}{2\sigma^2} \cdot \frac{1}{\sqrt{2\pi}\sigma}$$

equivalence

$$\max_{\omega} \sum_{i=1}^n -(\omega^T x_i - y_i)^2$$

$$\equiv \min_{\omega} \sum_{i=1}^n (\omega^T x_i - y_i)^2$$

$$\hat{\omega}_{ML} = \omega^* = (x^T)^+ x y$$

CONCLUSION: Maximum Likelihood estimator assuming ZERO MEAN GAUSSIAN NOISE is same as LINEAR REGRESSION with SQUARED ERROR!

What else have we gained?

- Can study properties of estimators $\hat{\omega}_{ML}$!

$$w^* = \hat{w}_{ML} \leftarrow \text{Max likelihood.} = \boxed{(X^T X)^{-1} X^T y}$$

$$y/x = \underline{w}^T x + \underline{\epsilon} \rightsquigarrow N(0, \underline{\sigma}^2) \quad \begin{cases} (x_1, y_1), \dots, (x_n, y_n) \\ \downarrow \\ w^T x + \epsilon, \end{cases}$$

$$\underline{w} \in \mathbb{R}^d ; \hat{w}_{ML} \in \mathbb{R}^d$$

Want a way to understand how good \hat{w}_{ML} is in estimating w

Mean Squared error $\rightarrow \mathbb{E} [\|\hat{w}_{ML} - w\|^2]$ = $\underline{\sigma}^2 \cdot \underline{\text{trace}}((X^T)^{-1})$

\downarrow
over randomness in y

$$A = \begin{bmatrix} a_1 & a_2 & \dots & a_d \end{bmatrix}$$

$$\text{tr}(A) = \sum_{i=1}^d a_i = \sum_{i=1}^d \lambda_i$$

\uparrow Eigenvalue of A

$$\text{trace}((X^T)^{-1})$$

Let Eigenvalues of (X^T) be $\{\lambda_1, \dots, \lambda_d\}$
 Eigenvalues of $(X^T)^{-1}$ are $\{\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_d}\}$

Mean sq. error (\hat{w}_{ML})

$$\mathbb{E} (\|\hat{w}_{ML} - w\|^2) = \underline{\sigma}^2 \left(\sum_{i=1}^d \frac{1}{\lambda_i} \right)$$

Consider the following estimator:

$$\hat{w}_{\text{new}} = (X^T + \lambda I)^{-1} X^T y$$

$\in \mathbb{R}^d$ $\in \mathbb{R}^{d \times d}$

$$\hat{w}_{ML} = (X^T)^{-1} X^T y$$

For some matrix A , let Eigenvalues be $\{\lambda_1, \dots, \lambda_d\}$

What are Eigenvalues of $A + \lambda I$? $\{\lambda_1 + \lambda, \dots, \lambda_d + \lambda\}$

$$A v_i = \lambda_i v_i$$

$$(A + 2I) v_i = A v_i + 2v_i$$

$$= \lambda_i v_i + 2v_i$$

$$= (\lambda_i + 2)v_i$$

$$\text{trace}((X^T + \lambda I)^{-1}) = \left(\sum_{i=1}^n \frac{1}{\lambda_i + \lambda} \right)$$

EXISTENCE Thm : (Informal)

$\exists \lambda \in \mathbb{R}$ s.t

$\hat{w}_{\text{new}} = (X^T + \lambda I)^{-1} X y$ has lesser mean sq. error
than \hat{w}_{ML}

In practice, find λ by CROSS VALIDATION



- Train on the training set and check for error on validation set.
- Pick λ that gives least error.

K-FOLD CROSS VALIDATION



- Train on folds $\{F_1, \dots, F_{i-1}, F_{i+1}, \dots, F_k\}$
- Validate on F_i

- Pick λ that gives least average error.

LEAVE ONE OUT CROSS VALIDATION

$$\hat{w}_{\text{new}} = (X^T + \lambda I)^{-1} X y$$

IS there an alternate way to understand \hat{w}_{ML} ?

BAYESIAN MODELING

- NEED A PRIOR on w i.e., $P(w) \xrightarrow{\text{R}^d}$

Likelihood

$$y/x \sim N(w^T x, 1)$$

For Simplicity. Can use σ^2 as well.

A CHOICE FOR PRIOR

$$w \sim N(0, \frac{\sigma^2}{\lambda} I)$$

$$\begin{bmatrix} \sigma^2 & & & \\ & \sigma^2 & & \\ & & \ddots & \\ & & & \sigma^2 \end{bmatrix}$$

Covariance Matrix $\mathbb{R}^{d \times d}$

As usual,

$$P(w \mid \{(x_1, y_1), \dots, (x_n, y_n)\}) \propto P(\{(x_1, y_1), \dots, (x_n, y_n)\} \mid w) \cdot P(w)$$

$$\propto \left(\prod_{i=1}^n e^{-\frac{(y_i - w^T x_i)^2}{2}} \right) \cdot \left(\prod_{i=1}^d e^{-\frac{(w_i - 0)^2}{2\sigma^2}} \right)$$

$e^{-\sum_{i=1}^d \frac{w_i^2}{2\sigma^2}}$

$$\propto \left(\prod_{i=1}^n e^{-\frac{(y_i - w^T x_i)^2}{2}} \right) \cdot e^{-\frac{\|w\|^2}{2\sigma^2}}$$

How will the MAP estimate look like?

$$\hat{w}_{MAP} = \arg \max_w \sum_{i=1}^n -\frac{(y_i - w^T x_i)^2}{2} - \frac{\|w\|^2}{2\sigma^2}$$

$$\hat{w}_{MAP} = \arg \min_w \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \frac{1}{2\sigma^2} \|w\|^2$$

$f(w)$

Take gradient, set it to 0 to solve for \hat{w}_{MAP} .

$$\nabla F(\omega) = (X X^T) \omega - X y + \frac{\omega}{\gamma^2} \quad [\text{verify}]$$

$$\hat{\omega}_{MAP} = (X X^T + \frac{1}{\gamma^2} I)^{-1} X y$$

CROSS VALIDATE in PRACTICE

CONCLUSION: MAP ESTIMATION for linear regression
 with a Gaussian prior $N(0, \frac{1}{\gamma^2} I)$ for ω is
 equivalent to "NEW" estimator we used earlier.

LINEAR REGRESSION

$$\hat{\omega}_{ML} = \arg \min_{\omega} \sum_{i=1}^n (\omega^T x_i - y_i)^2$$

RIDGE REGRESSION

$$\hat{\omega}_R = \arg \min_{\omega} \sum_{i=1}^n (\omega^T x_i - y_i)^2 + \lambda \|\omega\|^2$$

LOSS REGULARIZER

f_1 height	f_2 weight	f_3 $2\text{height} + 3\text{weight}$	<u>label</u> : $3\text{height} + 4\text{weight}$
1	1	1	
0	c ₁	c ₂	
3	4	0	

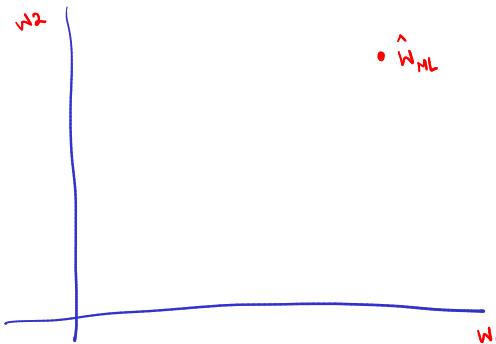
Linear regression / Ridge regression

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2$$

\downarrow

REGULARIZATION

PARAMETER SPACE



$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

Where is

\hat{w}_R → solution of
the ridge regression
problem?

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|^2$$

↑
Ⓐ

is equivalent to

$$\boxed{\begin{aligned} \min_{w \in \mathbb{R}^d} & \sum_{i=1}^n (w^T x_i - y_i)^2 \\ \text{s.t.} & \|w\|^2 \leq \theta \end{aligned}}$$

Ⓑ

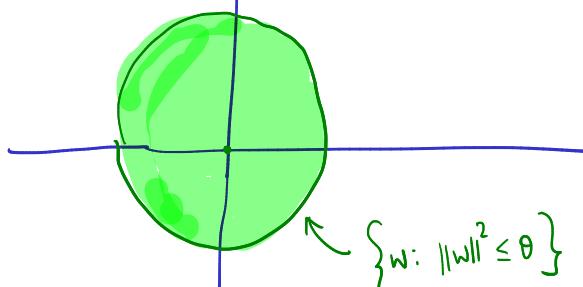
depends on λ

For every choice of $\lambda > 0$, if we set the optimal solutions of problems Ⓐ and Ⓑ coincide.

\hat{w}_{ML}

$$\|w\|^2 \leq \theta$$

$$\Rightarrow w_1^2 + w_2^2 \leq \theta$$



$$\{w : \|w\|^2 \leq \theta\}$$

What is the loss / error / objective function value of linear regression at \hat{w}_{ML}

$$\sum_{i=1}^n (\hat{w}_{ML}^T x_i - y_i)^2 = f(\hat{w}_{ML})$$

Consider the set of all w s.t

$$f(w) = f(\hat{w}_{ML}) + c \quad c \geq 0$$

$$S_c = \left\{ w : f(w) = f(\hat{w}_{ML}) + c \right\}$$

i.e., every $w \in S_c$ satisfies

$$\underbrace{\|x^T w - y\|^2}_{f(w)} = \underbrace{\|x^T \hat{w}_{ML} - y\|^2}_{f(\hat{w}_{ML})} + c$$

On Simplification [Please do this],

one gets

$$(w - \hat{w}_{ML})^T (x^T) (w - \hat{w}_{ML}) = c' \quad \text{Some constant that depends on } c, (x^T), \hat{w}_{ML} \text{ and not on } w$$

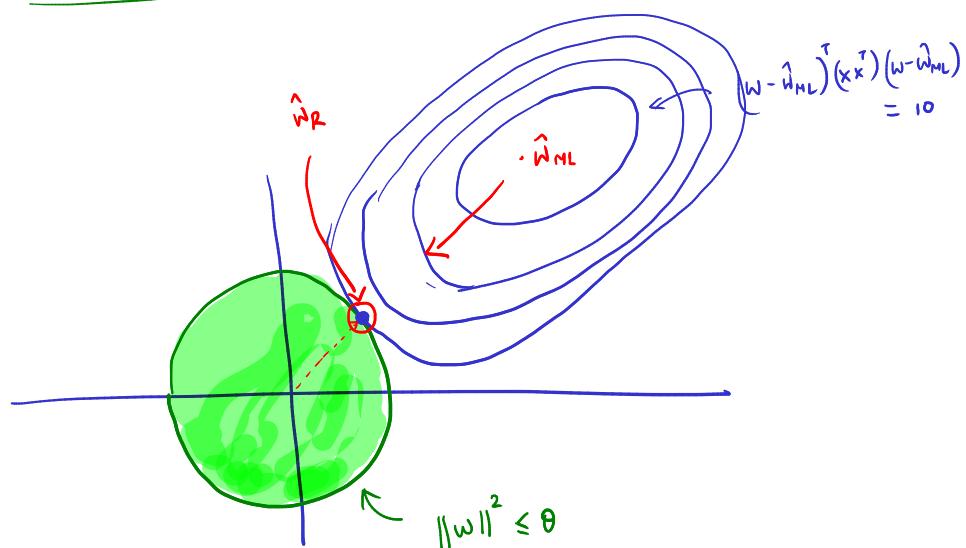
If $x^T = I$

$$(w - \hat{w}_{ML})^T I (w - \hat{w}_{ML}) = c'$$

$$\|w - \hat{w}_{ML}\|^2 = c'$$

$$\|w - \hat{w}_{ML}\|^2 = c'$$

$$\hat{w}_{ML}$$



Conclusions:

- Ridge regression pushes feature values towards 0. But does not necessarily make it 0.

- An alternate way to regularize would man be using $\|\cdot\|_1$ norm instead of $\|\cdot\|_2$ norm

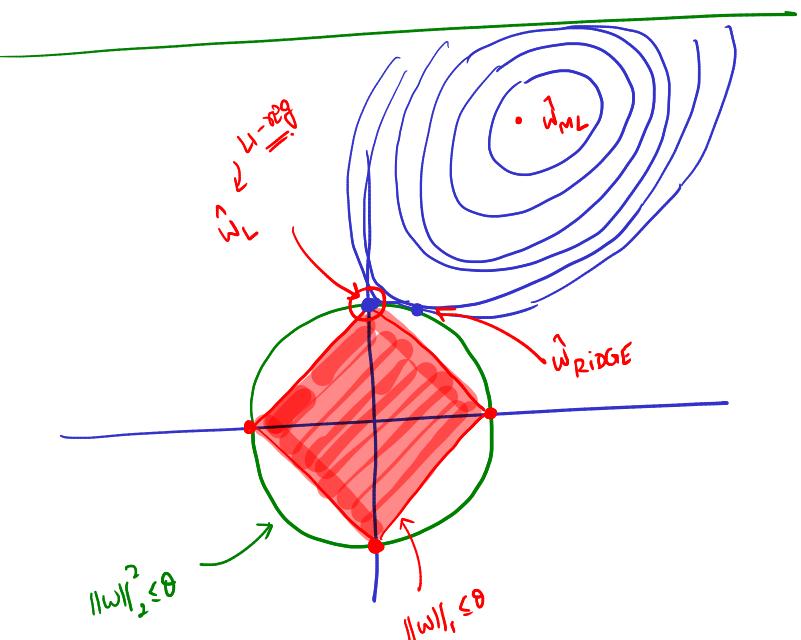
$$\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$$

L1 Regularization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

$$\equiv \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

$$\text{st } \|\mathbf{w}\|_1 \leq \theta$$

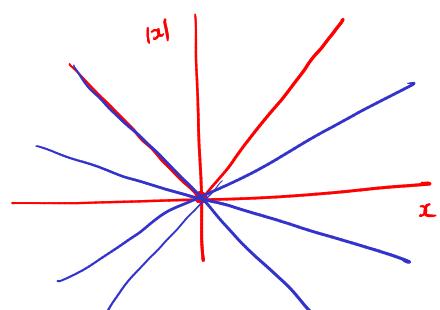
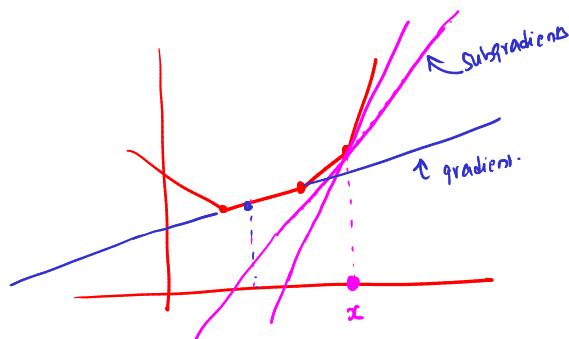


L1 Reg: LASSO - LEAST ABSOLUTE SHRINKAGE and SELECTION OPERATOR.

Points

- LASSO does not have a closed form solution
- Sub-gradient methods are usually used to solve LASSO.

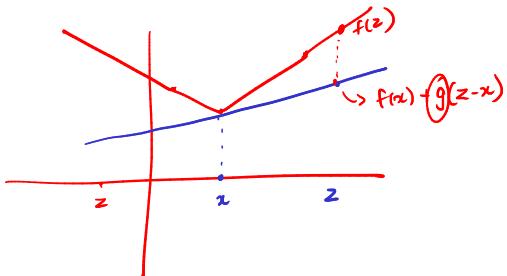
$$\begin{aligned} \text{Subgrad at } 0 \\ = [-1, 1] \end{aligned}$$



Subgradient

A vector $g \in \mathbb{R}^d$ is a sub-gradient of $f: \mathbb{R}^d \rightarrow \mathbb{R}$ at a point $x \in \mathbb{R}^d$ if

$$f(z) \geq f(x) + g^\top(z - x)$$



Why Subgradients?

- If function f to minimize is a convex function, then sub-gradient descent Converges!
- There are other special purpose methods for LASSO \rightarrow (e.g.) IRLS [Iterative Reweighted Least Squares].

SUPERVISED LEARNING

BINARY CLASSIFICATION

$$\left\{ \begin{array}{l} x_1, \dots, x_n \\ y_1, \dots, y_n \end{array} \right\} \quad x_i \in \mathbb{R}^d \quad y_i \in \{0, 1\} / \{-1, +1\}$$

Goal:
 $f: \mathbb{R}^d \rightarrow \{0, 1\}$

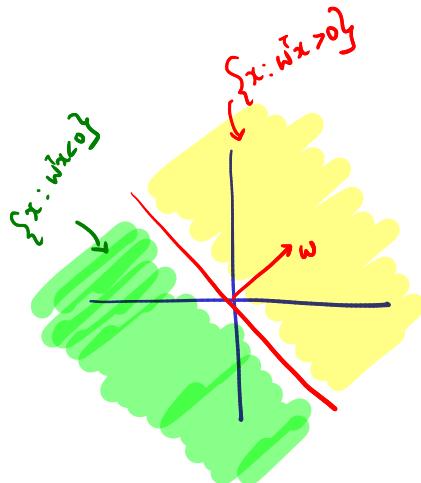
LOSS / ERROR

$$\text{Loss}(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(f(x_i) \neq y_i)$$

↓
0-1-loss

$$\min_{f \in \mathcal{H}_{\text{linear}}} \sum_{i=1}^n \mathbb{1}(f(x_i) \neq y_i)$$

NP-hard problem in general



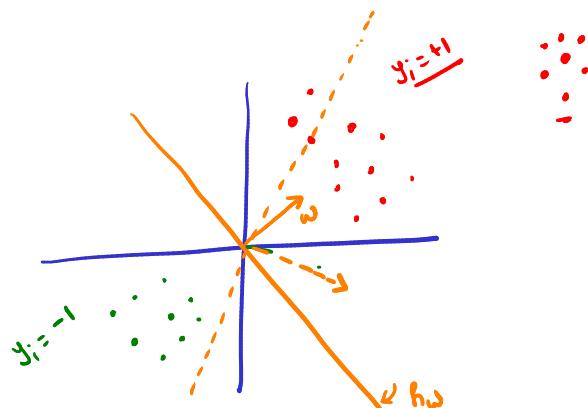
$$\mathbb{1}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{o/w} \end{cases}$$

$$\mathcal{H}_{\text{linear}} = \left\{ f_w: \begin{array}{l} f_w(x) = \text{Sign}(w^T x) \\ w \end{array} \right\}$$

$$\text{Sign}(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{o/w} \end{cases}$$

- Can we use linear regression to solve classification problem?

$$\{(x_1, y_1), \dots, (x_n, y_n)\} \rightarrow \boxed{\text{Lin Reg}} \rightarrow w \in \mathbb{R}^d \rightarrow f_w: \mathbb{R}^d \rightarrow \{0, 1\}$$



Conclusion

Regression is SENSITIVE
 TO location of the
 data points and not just
 the "Side" on which the
 data lies w.r.t separation.

SIMPLE ALGORITHMS FOR CLASSIFICATION

- Given a test point $x_{\text{test}} \in \mathbb{R}^d$, find the closest point x^* to x_{test} in the training set
- Predict $y_{\text{test}} = y^*$

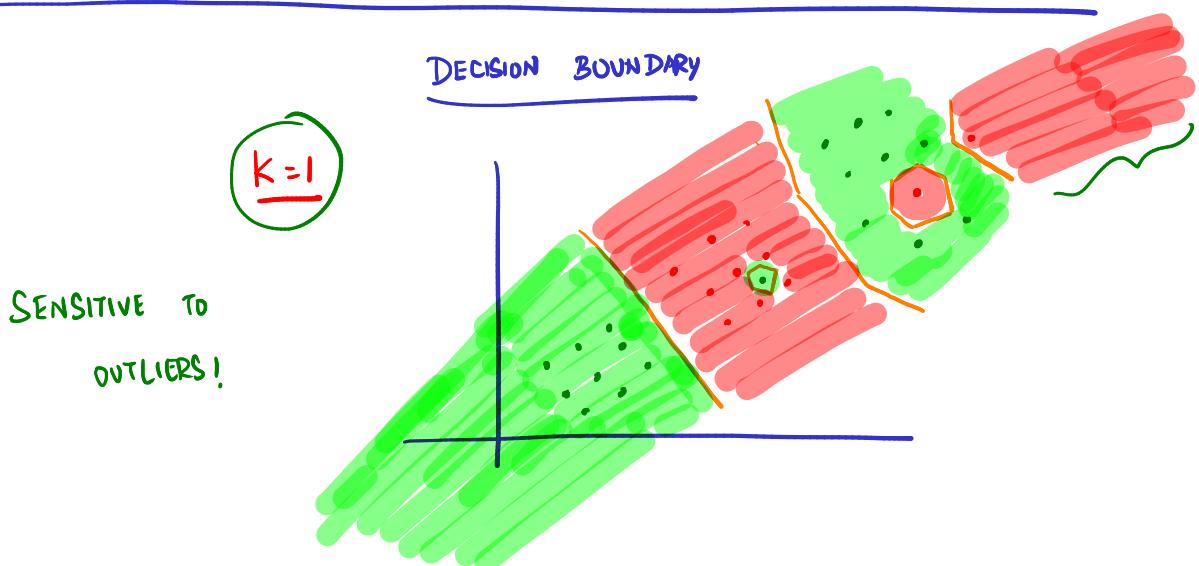
ISSUE: Can get affected by outliers

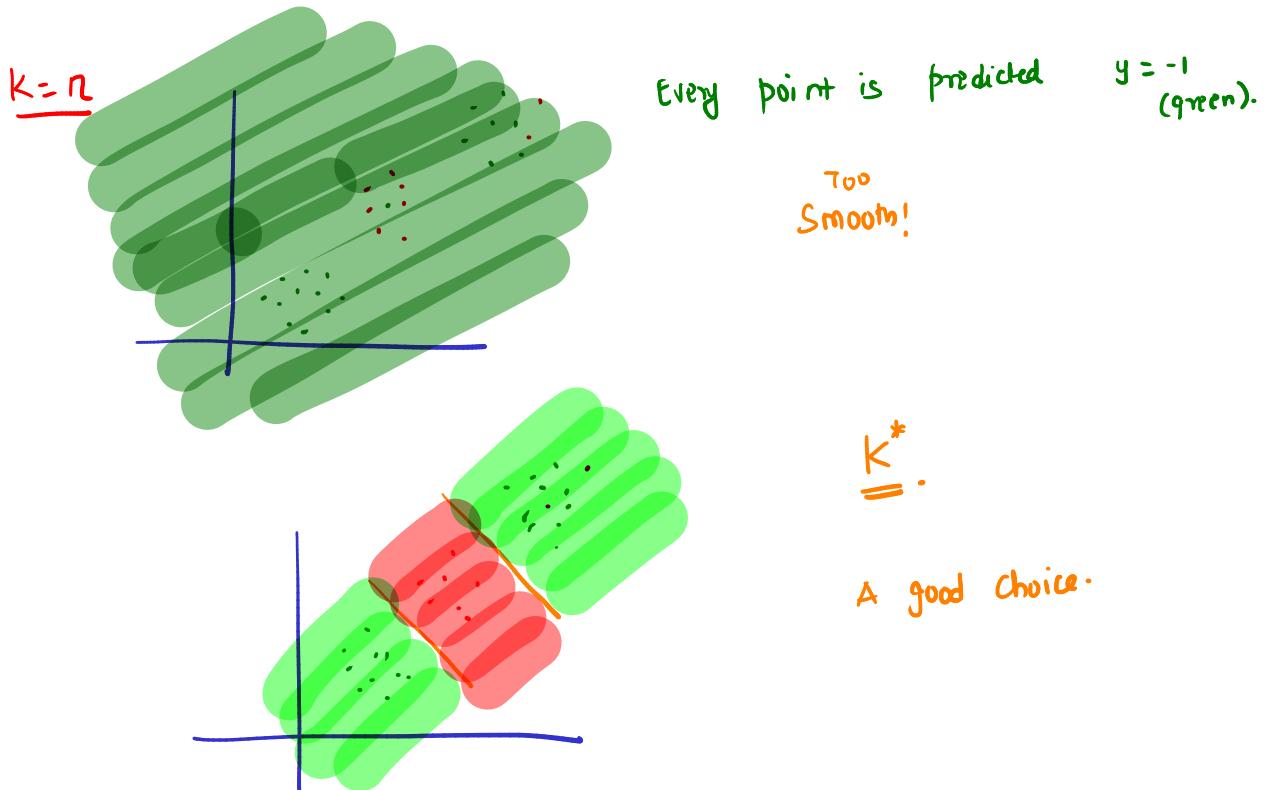
Fix: Ask more neighbours

(K-NN) (K-Nearest Neighbours)

- Given x_{test} , find the k-Closest points in the training set - $(x_1^*, x_2^*, \dots, x_k^*)$

- PREDICT $y_{\text{test}} = \text{majority}(y_1^*, y_2^*, \dots, y_k^*)$





Choosing K

- Can treat as a HYPER PARAMETER
- Smaller the k , complicated the decision boundary
- Soln: CROSS-VALIDATE for k

ISSUES with K-NN

- Choosing a distance function
- PREDICTION is COMPUTATIONALLY EXPENSIVE.
- No MODEL is LEARNT.
 - Cannot throw away data after "LEARNING"

DECISION TREES

$\forall i$
 $x_i \in \mathbb{R}^d$
 $y_i \in \{+1, 0\}$

INPUT:

Dataset $\left\{ (x_1, y_1), \dots, (x_n, y_n) \right\}$

OUTPUT:

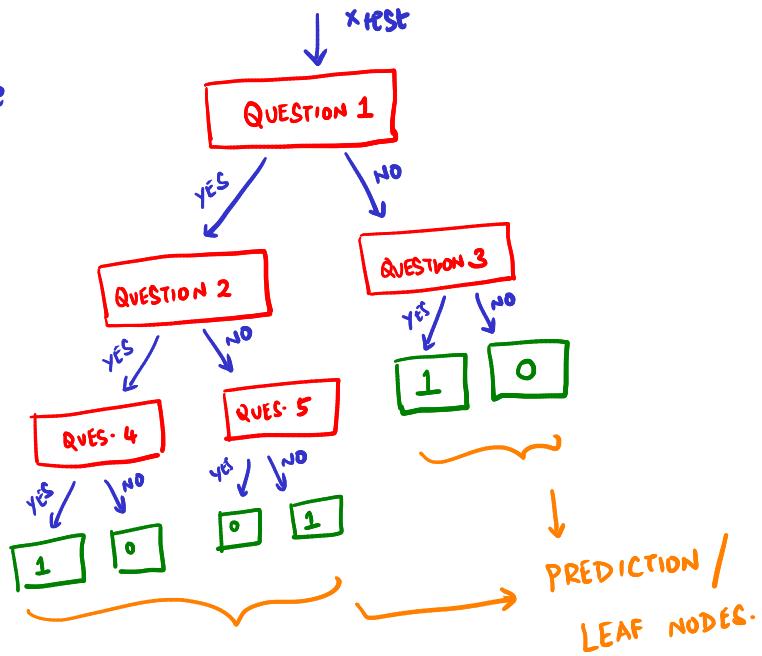
DECISION TREE

PREDICTION: Given x_{test} ,

traverse through the tree to reach a leaf node.

y_{test} = value in leaf node.

DECISION tree



QUESTION: A question is a (feature, value) pair.

Eg: height $\leq 180\text{cm}$?
 $(f_3) \quad \theta$

How to measure "goodness" of a question?

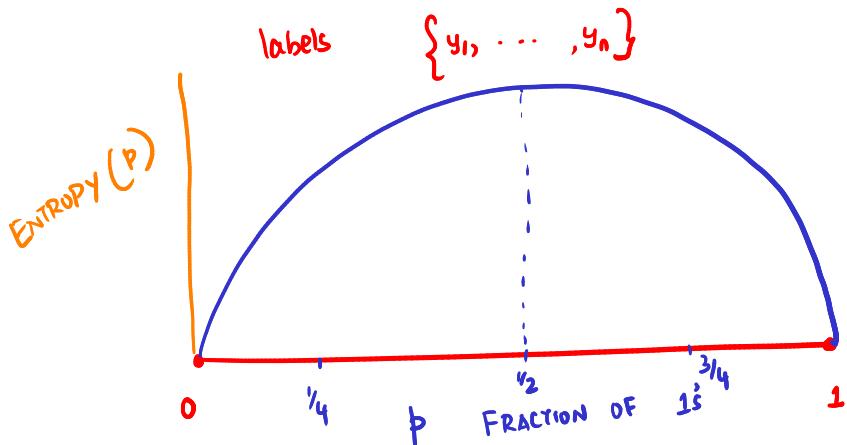
DATASET
 $D = \left\{ (x_1, y_1), \dots, (x_n, y_n) \right\}$

$f_k \leq \theta ?$

YES NO

$D_{yes} = \left\{ (x_1, y_1), (x_{10}, y_{10}), \dots \right\}$ $D_{no} = \left\{ (x_2, y_2), (x_3, y_3), \dots \right\}$

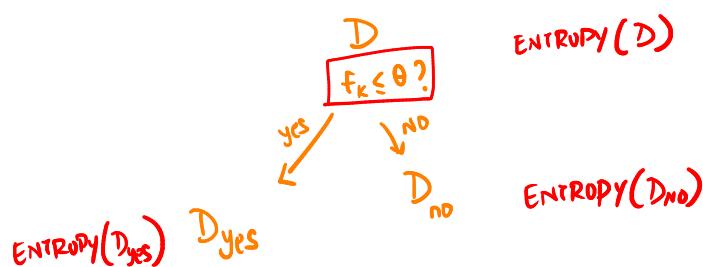
• Need is a measure of "Impurity" for a set of



$$\text{ENTROPY}(\{y_1, \dots, y_n\}) = \text{ENTROPY}(p)$$

$$= -\left(p \log p + (1-p) \log(1-p) \right)$$

[convention
 $\log(0)=0$]



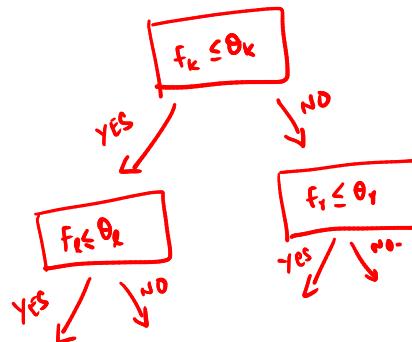
INFORMATION GAIN (feature, value) =

$$\text{ENTROPY}(D) - \left[\gamma \text{ ENTROPY}(D_{\text{yes}}) + (1-\gamma) \text{ ENTROPY}(D_{\text{no}}) \right]$$

$$\gamma = \frac{|D_{\text{yes}}|}{|D|}$$

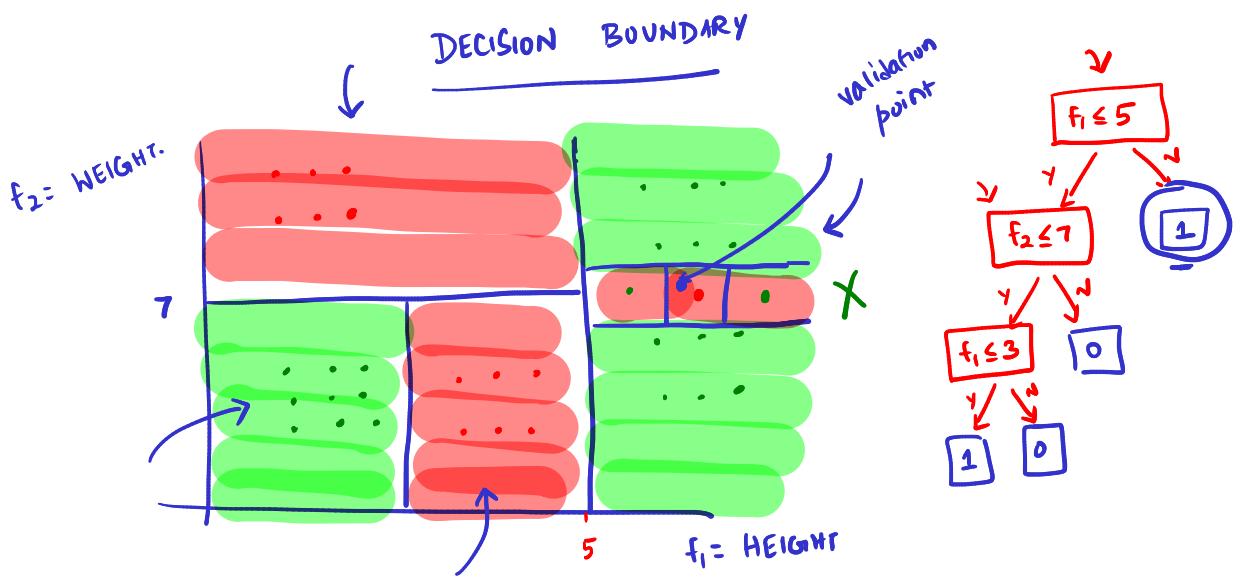
ALGORITHM - DECISION TREE

- DISCRETIZE each feature in $[\min, \max]$ range
- Pick the Question that has the largest Information gain.
- Repeat the procedure for D_{yes} , D_{no}

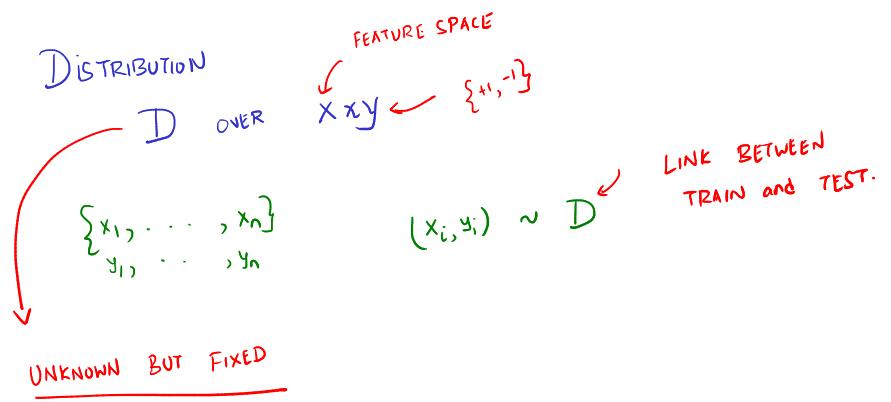


POINTS

- Can stop growing a tree if a node becomes "SUFFICIENTLY" pure.
- DEPTH of the tree is a hyperparameter
- There are alternate measures for "goodness" of a question
 - GINI INDEX



TYPES OF MODELING



CLASSIFICATION

- GENERATIVE MODEL
- DISCRIMINATIVE MODEL

GENERATIVE MODEL

- MODEL $P(x, y)$

↳ NEXT.

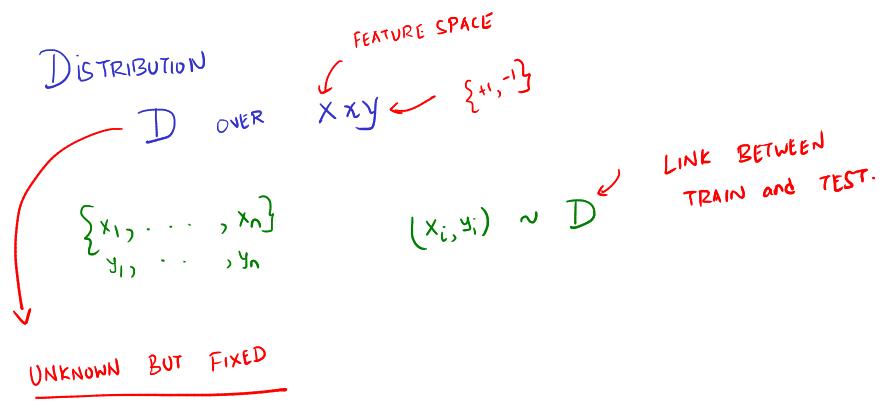
DISCRIMINATIVE MODEL

- MODEL $P(y|x)$

EG: K-NN DECISION TREES

$P(y=1|x) = 1$ if decision tree for x says 1.
 $P(y=1|x) = 1$ if majority of neighbours say 1.
 $= 0$ otherwise

TYPES OF MODELING



CLASSIFICATION

- GENERATIVE MODEL
- DISCRIMINATIVE MODEL

GENERATIVE MODEL

- MODEL $P(x, y)$

↳ NEXT.

DISCRIMINATIVE MODEL

- MODEL $P(y|x)$

EG: K-NN DECISION TREES

$P(y=1|x) = 1$ if decision tree for x says 1.
 $P(y=1|x) = 1$ if majority of neighbours say 1.
 $= 0$ otherwise

$$\text{DATA} = \left\{ (x_1, y_1), \dots, (x_n, y_n) \right\}$$

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

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

EXAMPLE : SPAM CLASSIFICATION

$$x_i \in \{0, 1\}^d \quad \# \text{ words in dictionary}$$

Eg: "Hello, how are you?"

$$\rightarrow \begin{bmatrix} \text{About} & \dots & \text{ARE} & \dots & \text{HELLO} & \dots & \text{HOW} & \dots & \text{you} \\ 0 & 0 & \dots & 1 & \dots & 1 & \dots & 1 & \dots & 0 \end{bmatrix} \quad \text{ZEBRA}$$

$$P(x, y) \quad ("Hello, how are you?", \text{ spam}) = 0.01$$

$$P(x, y) = \frac{P(x) \cdot P(y|x)}{P(\text{email}) P(\text{spam}/\text{email})} = \frac{\underbrace{P(y)}_{P(\text{spam})} \cdot \underbrace{P(x|y)}_{P(\text{email}/\text{spam})}}{P(\text{spam})}$$

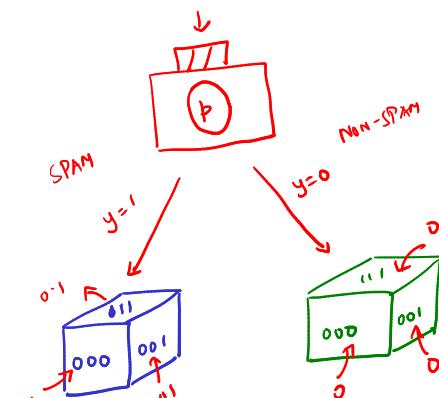
GENERATIVE STORY

STEP 1: DECIDE the Label by tossing a coin

$$P(y_i = 1) = p$$

STEP 2: DECIDE features using the label in step 1
by $P(x_i | y_i)$

	w_1	w_2	w_3	
0 0 0	0	0	0	0.1
0 0 1	0	0	1	0.01
0 1 0	0	1	0	0.02
0 1 1	0	1	1	0.03
1 0 0	1	0	0	0.04
1 0 1	1	0	1	0.5
1 1 0	1	1	0	0.2
1 1 1	1	1	1	0.1



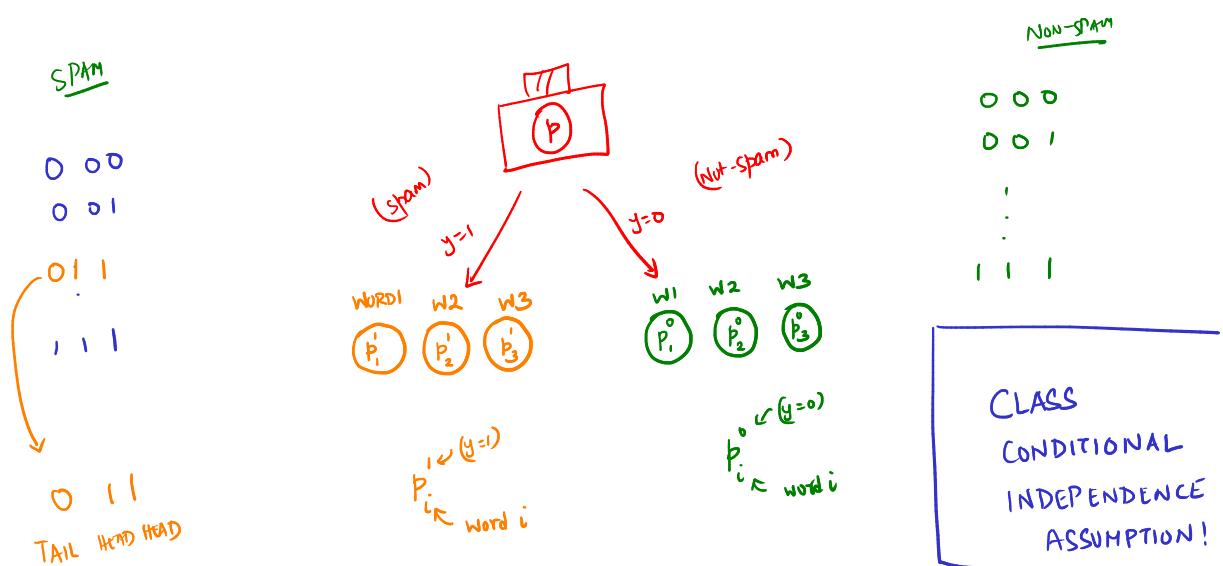
w_1	w_2	w_3	
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	0.5
1	0	0	0.5
1	0	1	0
1	1	0	0
1	1	1	0

$= 1$

$$\begin{aligned} \# \text{ parameters in } \text{the} & \quad = 1 + \underbrace{\binom{2^d-1}{2}}_{P(x|y=1)} + \underbrace{\binom{d}{2-1}}_{P(x|y=0)} \\ \text{model} & \quad \uparrow \quad \uparrow \\ & \quad \text{to decide} \\ & \quad \text{label} \\ & = 1 + 2(2^d-1) \\ & = 2^{d+1}-1 \end{aligned}$$

ISSUE

- Too many parameters!
- Not a reasonable story.



$$\# \text{ parameters} = 1 + d + d = \underbrace{2d+1}_{\text{MANAGEABLE!}}$$

$$\text{STEP 1: } P(y=1) = p$$

$$P(x = [f_1 \ f_2 \ \dots \ f_d] | y=1) \\ = \\ (1-p_1) \cdot (p_2) \cdot (1-p_3)$$

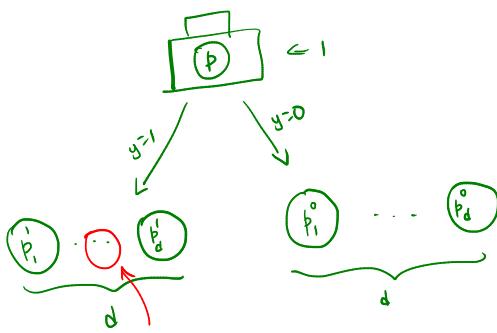
STEP 2:

$$P(x = [f_1 \ f_2 \ \dots \ f_d] | y)$$

$$= \prod_{i=1}^d (p_i^{y_i} (1-p_i)^{1-y_i})$$

FEATURES ARE
CONDITIONALLY INDEPENDENT
GIVEN LABEL

- How to estimate the parameters
- SOLUTION: MAXIMUM LIKELIHOOD!



CLASS CONDITIONAL INDEPENDENCE
2d + 1

PARAMETER ESTIMATION

$$p, \{p'_1, \dots, p'_d\}, \{p_1, \dots, p_d\}$$

MAX. LIKELIHOOD ESTIMATES

$$1 \rightarrow \hat{p} = \frac{1}{n} \sum_{i=1}^n y_i \rightarrow \{\text{Fraction of Spam emails in the dataset}\}$$

$$\begin{aligned} \forall j \in \{1, \dots, d\} \quad \hat{p}_j^y &= \frac{\sum_{i=1}^n \mathbb{1}(f_j^i = 1, y_i = y)}{\sum_{i=1}^n \mathbb{1}(y_i = y)} \\ &\quad \leftarrow \text{Number of emails with label } y. \end{aligned}$$

↳ Fraction of y-labelled emails that contain the jth word.

PREDICTION

Given $x^{\text{test}} \in \{0, 1\}^d$, what is \hat{y}^{test} ?

$$P(y^{\text{test}} = 1 | x^{\text{test}}) > P(y^{\text{test}} = 0 | x^{\text{test}})$$

$$\Rightarrow \hat{y}^{\text{test}} = 1$$

= 0 otherwise

How to obtain $P(y/x)$ from $P(y)$ and $P(x/y)$?

BAYES RULE!

$$P(y^{\text{test}}=1 | x^{\text{test}}) = \frac{P(x^{\text{test}} | y^{\text{test}}=1) \cdot P(y^{\text{test}}=1)}{P(x^{\text{test}})}$$

$$P(y^{\text{test}}=0 | x^{\text{test}}) = \frac{P(x^{\text{test}} | y^{\text{test}}=0) \cdot P(y^{\text{test}}=0)}{P(x^{\text{test}})}$$

$$P(x^{\text{test}} | y^{\text{test}}=1) \cdot P(y^{\text{test}}=1)$$

$$= P(x^{\text{test}} = [f_1 f_2 \dots f_d] | y^{\text{test}}=1) \cdot P(y^{\text{test}}=1)$$

$$= \left(\prod_{j=1}^d \left(\hat{p}_j^{f_j} (1-\hat{p}_j)^{1-f_j} \right) \right) \cdot \hat{p}$$

$$x^{\text{test}} = \begin{bmatrix} f_1 & f_2 & f_3 & f_4 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\left(\prod_{j=1}^d \left(\hat{p}_j^{f_j} (1-\hat{p}_j)^{1-f_j} \right) \cdot \hat{p} \right) > \left(\prod_{j=1}^d \left(\hat{p}_j^0 (1-\hat{p}_j)^{1-f_j} \right) \cdot (1-\hat{p}) \right)$$

⇒ PREDICT $\hat{y}^{\text{test}} = 1$
else $\hat{y}^{\text{test}} = 0$.

MODEL uses 2 main things

NAIVE
BAYES
ALGORITHM.

- CLASS CONDITIONAL INDEPENDENCE
- BAYES THEOREM

- may not hold in practice
- NAIVE ASSUMPTION
- Still works well in practice.

PITFALLS IN NAIVE BAYES TO
WATCH OUT FOR.

- IF a word does not appear in the train set but appears in a test datapoint,

$$\hat{p}_j^1 = 0 \quad \hat{p}_j^0 = 0$$

$$P(y_{\text{test}}=1 | x_{\text{test}}=[f_1, f_2, \dots, f_d]) \propto \left(\prod_{i=1}^d \underbrace{\left(\hat{p}_i^1\right)^{1_{f_i}}}_{=1} \underbrace{\left(1 - \hat{p}_i^1\right)^{1_{(1-f_i)}}}_{=1} \right) \hat{p}$$

$$P(y_{\text{test}}=0 | x_{\text{test}}=[f_1, f_2, \dots, f_d]) \propto \left(\prod_{i=1}^d \underbrace{\left(\hat{p}_i^0\right)^{1_{f_i}}}_{=0} \underbrace{\left(1 - \hat{p}_i^0\right)^{1_{(1-f_i)}}}_{=1} \right) (1 - \hat{p})$$

Possible Fix

- Can add two "pseudo" emails with all words present - one email has label 0 and another has label 1

LAPLACE
SMOOTHING

$$\begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & 1 & 1 & \dots & & 1 \end{bmatrix} \quad \begin{array}{l} \text{pseudo} \\ y_1 \\ 2 \\ 0 \end{array}$$

DECISION FUNCTION
OF NAIVE BAYES

Given x_{test} : $y_{\text{test}} = 1$ if $\frac{P(y_{\text{test}}=1/x_{\text{test}})}{P(y_{\text{test}}=0/x_{\text{test}})} \geq 1$

$$\log \left(\frac{P(y_{\text{test}}=1|x_{\text{test}})}{P(y_{\text{test}}=0|x_{\text{test}})} \right) \geq 0$$

$$\log \left(\frac{\underbrace{P(x_{\text{test}}|y_{\text{test}}=1)}_{P(x_{\text{test}})} \cdot P(y_{\text{test}}=1)}{\underbrace{P(x_{\text{test}}|y_{\text{test}}=0)}_{P(x_{\text{test}})} \cdot P(y_{\text{test}}=0)} \right) \geq 0$$

$$x_{\text{test}} = [f_1, f_2, \dots, f_d]$$

$$\log \left(\prod_{i=1}^d \frac{\hat{p}_i^{f_i} (1-\hat{p}_i)^{1-f_i}}{\hat{p}_i^0 (1-\hat{p}_i)^{1-f_i}} \cdot \frac{\hat{p}}{1-\hat{p}} \right) \geq 0$$

$$\therefore \log \left(\prod_{i=1}^d \left(\frac{\hat{p}_i^1}{\hat{p}_i^0} \right)^{f_i} \left(\frac{1-\hat{p}_i^1}{1-\hat{p}_i^0} \right)^{1-f_i} \cdot \frac{\hat{p}}{1-\hat{p}} \right) \geq 0$$

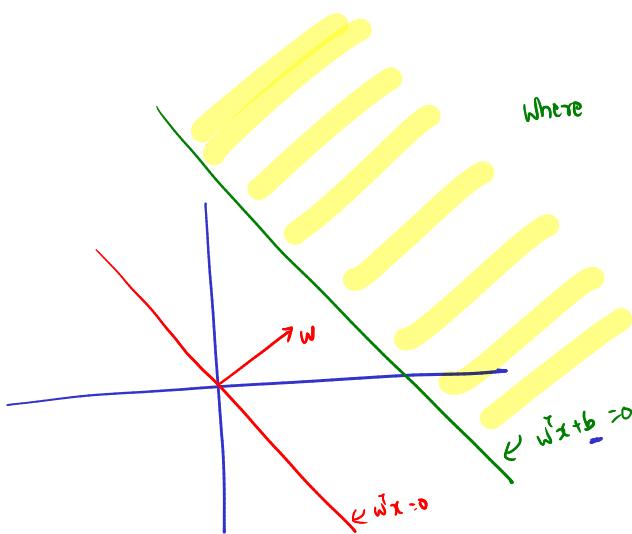
$$\therefore \sum_{i=1}^d \left(f_i \log \left(\frac{\hat{p}_i^1}{\hat{p}_i^0} \right) + (1-f_i) \log \left(\frac{1-\hat{p}_i^1}{1-\hat{p}_i^0} \right) + \log \left(\frac{\hat{p}}{1-\hat{p}} \right) \right) \geq 0$$

$$\therefore \sum_{i=1}^d f_i \left(\log \left(\frac{\hat{p}_i^1 (1-\hat{p}_i^0)}{\hat{p}_i^0 (1-\hat{p}_i^1)} \right) \right) + \log \left(\frac{(1-\hat{p}_i^1)}{(1-\hat{p}_i^0)} \right) + \log \left(\frac{\hat{p}}{1-\hat{p}} \right) \geq 0$$

$$x_{\text{test}} = [f_1, \dots, f_d]$$

DECISION FUNCTION is of the form

Predict $y_{\text{test}} = 1$ if $w^T x_{\text{test}} + b \geq 0$



where $w_i = \log \left(\frac{\hat{p}_i^1 (1-\hat{p}_i^0)}{\hat{p}_i^0 (1-\hat{p}_i^1)} \right)$, $b =$

CONCLUSION:

→ DECISION FUNCTION OF NAIVE BAYES IS LINEAR!

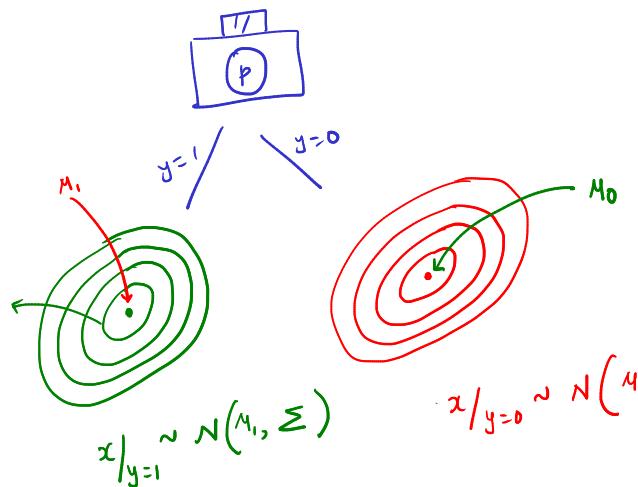
$$\text{DATA: } \left\{ (x_1, y_1), \dots, (x_n, y_n) \right\}$$

$$x_i \in \mathbb{R}^d \quad y_i \in \{0, 1\}$$

A GENERATIVE STORY

PARAMETERS

- β
- μ_0, μ_1
- Σ



NOTE: In this model, covariances are assumed to be same

MAXIMUM LIKELIHOOD ESTIMATES

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i}{n} \quad \leftarrow \quad \begin{matrix} \text{FRACTION OF points} \\ \text{labelled 1.} \end{matrix}$$

$$\hat{\mu}_1 = \frac{\sum_{i=1}^n \mathbb{1}(y_i=1) \cdot x_i}{\sum_{i=1}^n \mathbb{1}(y_i=1)} \quad \leftarrow \quad \begin{matrix} \text{Sample mean of} \\ \text{data points labelled 1.} \end{matrix}$$

$$\hat{\mu}_0 = \frac{\sum_{i=1}^n \mathbb{1}(y_i=0) \cdot x_i}{\sum_{i=1}^n \mathbb{1}(y_i=0)} \quad \leftarrow \quad \begin{matrix} \text{Sample mean of} \\ \text{data points labelled 0.} \end{matrix}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{M}_{y_i}) (x_i - \hat{M}_{y_i})^T$$

PREDICTION ? Bayes rule:

$$P(y_{\text{test}} | x_{\text{test}}) \propto \underbrace{P(x_{\text{test}} | y_{\text{test}})}_{f(x_{\text{test}}; \hat{M}_{y_{\text{test}}}, \hat{\Sigma})} \cdot \underbrace{P(y_{\text{test}})}_{\hat{p}}$$

Predict $y_{\text{test}} = 1$ if

$$f(x_{\text{test}}; \hat{M}_1, \hat{\Sigma}) \cdot \hat{p} \geq f(x_{\text{test}}; \hat{M}_0, \hat{\Sigma}) \cdot (1-\hat{p})$$

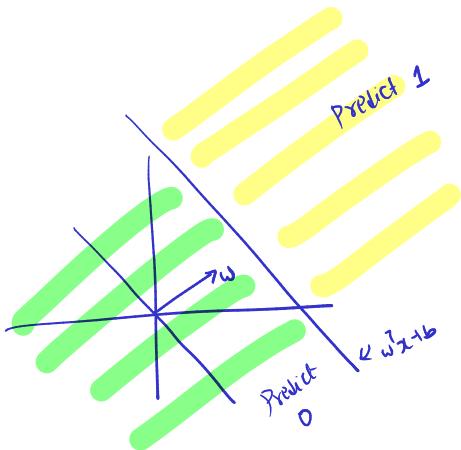
$$\frac{-(x_{\text{test}} - \hat{M}_1)^T \hat{\Sigma}^{-1} (x_{\text{test}} - \hat{M}_1)}{\hat{p}} \geq \frac{-(x_{\text{test}} - \hat{M}_0)^T \hat{\Sigma}^{-1} (x_{\text{test}} - \hat{M}_0)}{e^{(1-\hat{p})}}$$

On simplification

[take log]

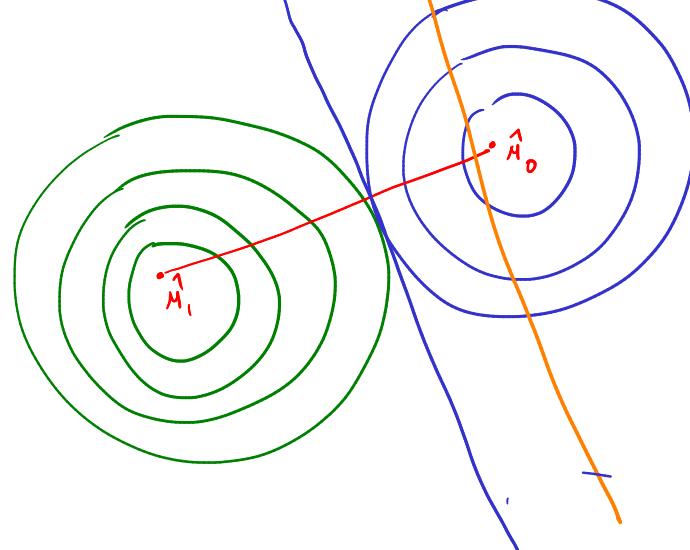
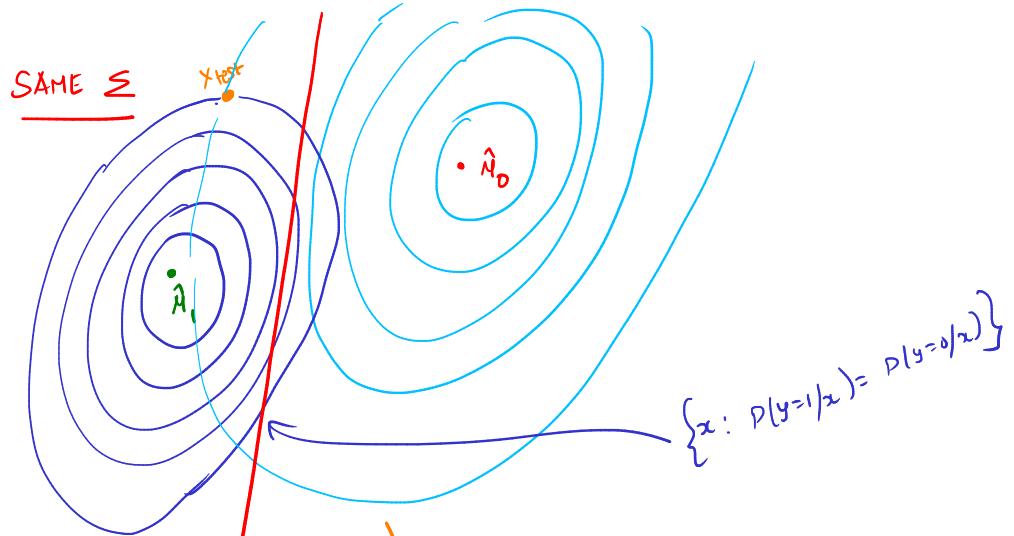
Predict 1 if

$$\left(\underbrace{(\hat{M}_1 - \hat{M}_0)^T \hat{\Sigma}^{-1}}_{w} \right) x_{\text{test}} + \underbrace{\hat{M}_0^T \hat{\Sigma}^{-1} \hat{M}_0 - \hat{M}_1^T \hat{\Sigma}^{-1} \hat{M}_1}_{+ \log(\frac{1-\hat{p}}{\hat{p}})} \geq 0$$

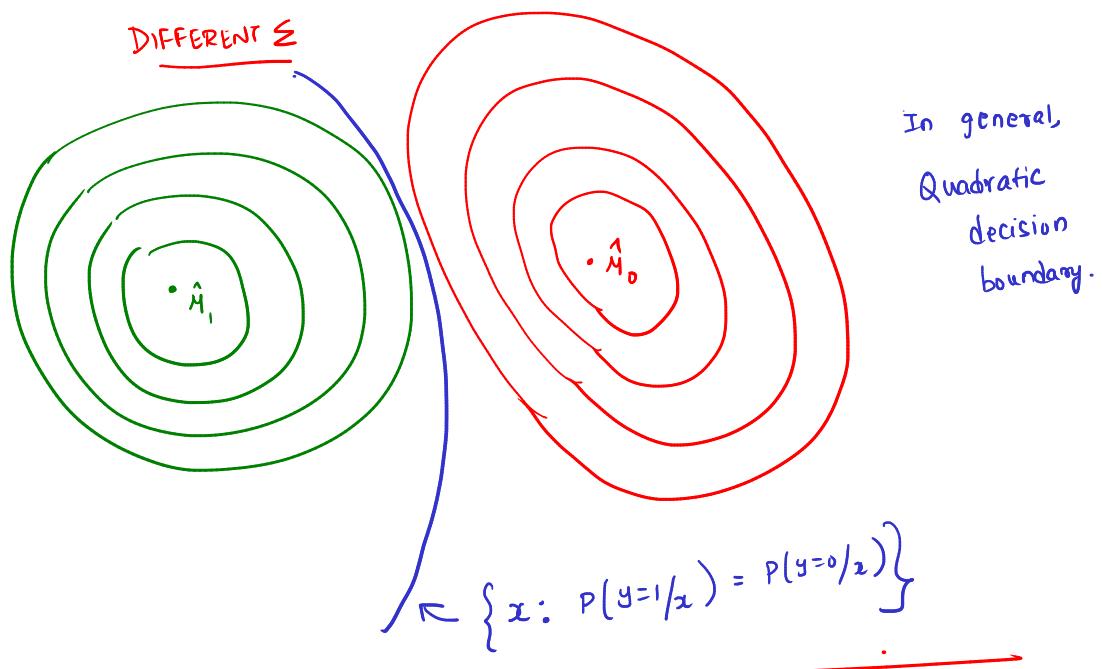


DECISION FUNCTION is LINEAR!

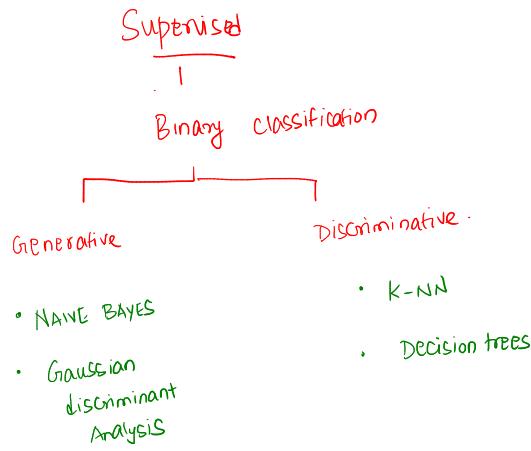
↳ Σ is same for both classes.



→



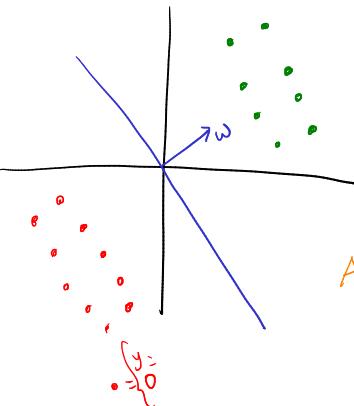
→ GAUSSIAN NAIVE BAYES



Discriminative models for classification

How to model $P(y=1/x)$?

$$\left\{ \begin{array}{l} y=1 \\ y=0 \end{array} \right.$$



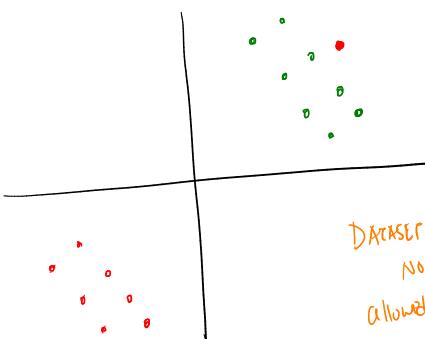
Allowed under our model!

SIMPLEST ASSUMPTION

$$P(y=1/x) = 1 \quad \text{if } \underline{w^T x \geq 0} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

0 otherwise.

LINEAR SEPARABILITY ASSUMPTION



Dataset is NOT allowed under our model.

Goal.

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

NP-HARD for a general dataset even if \mathbb{H} is just linear hypotheses.

- How about with extra "Linear separability" assumption.

LINEAR SEPARABILITY ASSUMPTION

$$\mathbf{w} \in \mathbb{R}^d \text{ s.t. } \text{Sign}(\mathbf{w}^T \mathbf{x}_i) = y_i \quad \forall i \in [n]$$

PERCEP TRON

[Rosenblatt, 1950's]

Input: $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$\begin{aligned} x_i &\in \mathbb{R}^d \\ y_i &\in \{+1, -1\} \end{aligned}$$

$$\overset{\text{ITERATION}}{\downarrow} \quad w^0 = 0 \in \mathbb{R}^d$$

$$[0 \ 0 \ \dots \ 0]$$

until convergence- Pick (x_i, y_i) pair from the dataset- IF $\text{sign}(w^t x_i) = y_i$

do nothing

ELSE

$$w^{t+1} = w^t + x_i y_i \quad \leftarrow \text{UPDATE RULE.}$$

$\downarrow \downarrow$
 $\mathbb{R}^d \quad \{\pm\}$

end.

end.

UPDATE RULE

if mistake

$$w^{t+1} = w^t + x_i y_i$$

MISTAKE TYPE 1

$$\text{Predicted} = 1 \leftarrow$$

$$\text{Actual} = -1$$

$$\boxed{(w^t x_i) \geq 0}$$

$y_i = -1$

MISTAKE TYPE 2

$$\begin{aligned} \text{Pred} &= -1 \quad \leftarrow (w^t x_i) < 0 \\ \text{Act} &= +1 \quad \leftarrow y_i = +1 \end{aligned}$$

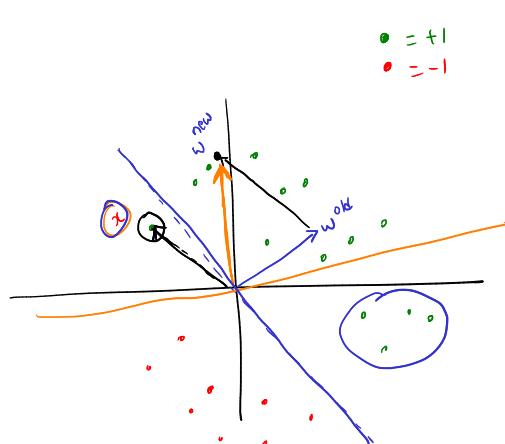
$$w^{t+1} = w^t + x_i y_i$$

$$(w^{t+1})^\top x_i = (w^t + x_i y_i)^\top x_i$$

$$\begin{aligned} &= w^t x_i + y_i \|x_i\|^2 \\ &\stackrel{\geq 0}{=} -1 \stackrel{\downarrow}{>} 0 \\ &\quad \text{Negative.} \end{aligned}$$

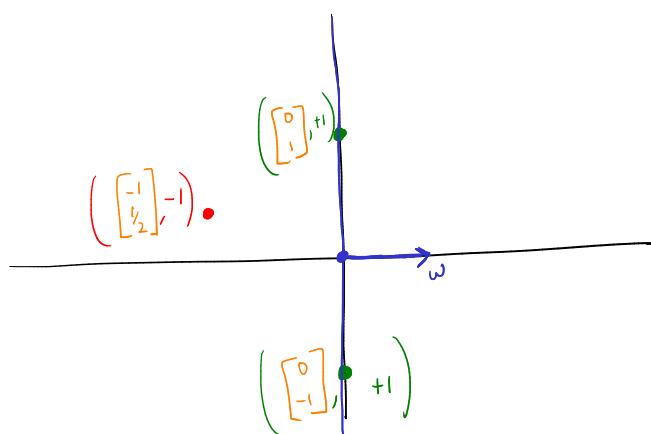
$$(w^{t+1})^\top x_i = \underbrace{w^t x_i}_{\geq 0} + \underbrace{y_i \|x_i\|^2}_{+1} \stackrel{< 0}{\rightarrow} \stackrel{> 0}{\rightarrow}$$

\Rightarrow update rule pushes w in the "right" direction for x_i



$$w^{\text{new}} = w^{\text{old}} + \alpha \cdot \frac{y}{\|y\|} \cdot x$$

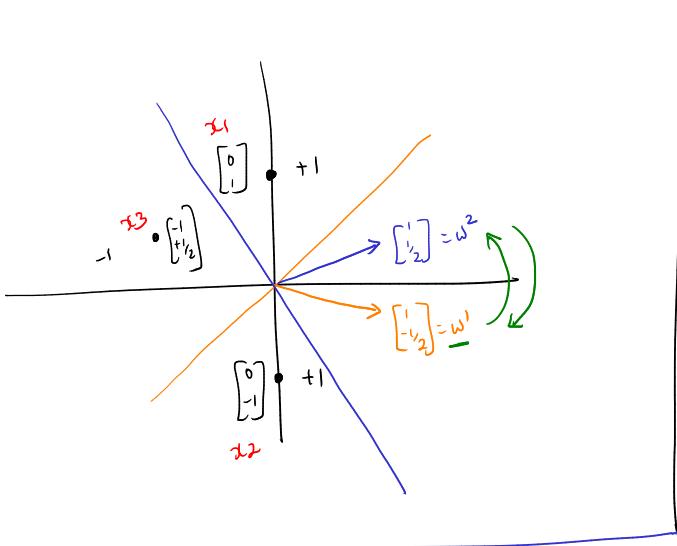
- Fixing w for one x might affect decision for other data points.
- So need more careful argument for convergence.



Is this a
Linearly Separable
dataset?

$$\left\{ \begin{array}{l} \begin{bmatrix} 0 \\ 1 \end{bmatrix}, +1 \\ \begin{bmatrix} 0 \\ -1 \end{bmatrix}, +1 \\ \begin{bmatrix} -1 \\ -1 \end{bmatrix}, -1 \end{array} \right\}$$

Is there a $w \in \mathbb{R}^2$ s.t. $w^T x_i \geq 0 \Rightarrow y_i = +1$
 $w^T x_i < 0 \Rightarrow y_i = -1$



PERCEPTRON

$$\vec{w}^0 = [0 \ 0]$$

$$\vec{w}^0 \cdot x_1 = 0 ; \vec{w}^0 \cdot x_2 = 0 ; \vec{w}^0 \cdot x_3 = 0$$

$$\hat{y}_1 = +1 \quad \hat{y}_2 = +1 \quad \hat{y}_3 = +1$$

$$\vec{w}^1 = \vec{w}^0 + x_3 \cdot y_3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -1 \\ 1/2 \end{bmatrix} \times -1$$

$$= \begin{bmatrix} 1 \\ -1/2 \end{bmatrix}$$

$$\vec{w}^2 = \vec{w}^1 + x_1 \cdot y_1 = \begin{bmatrix} 1 \\ -1/2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \cdot 1$$

$$\vec{w}^2 = \begin{bmatrix} 1 \\ 1/2 \end{bmatrix}$$

$$\vec{w}^3 = \vec{w}^2 + x_2 \cdot y_2 = \begin{bmatrix} 1 \\ 1/2 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} \times 1$$

$$\vec{w}^3 = \begin{bmatrix} 1 \\ -1/2 \end{bmatrix} = \vec{w}^1$$

ASSUMPTION

(1)

LINEAR SEPARABILITY with γ -MARGIN

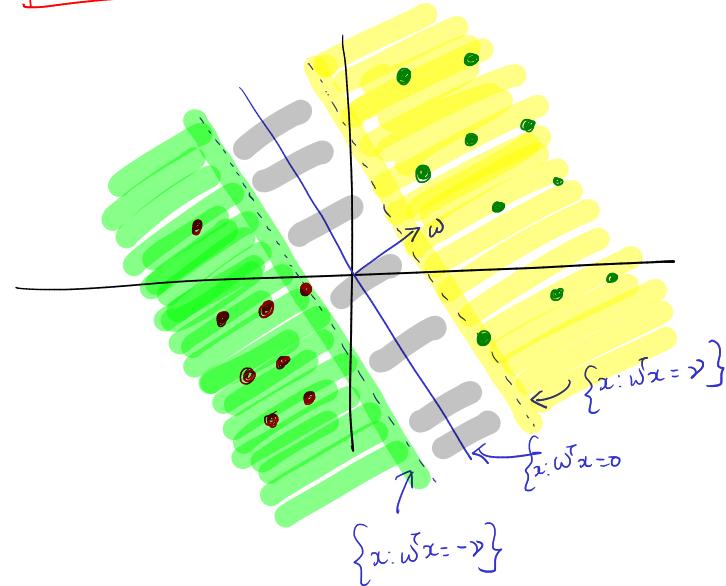
A Dataset
 $\{(x_1, y_1), \dots, (x_n, y_n)\}$

is Linearly Separable
 with γ -margin

if $\exists w^* \in \mathbb{R}^d$ st

$$(w^{*T} x_i) y_i \geq \gamma \quad \forall i$$

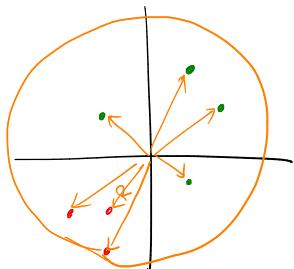
for some $\gamma > 0$



(2)

RADIUS ASSUMPTION

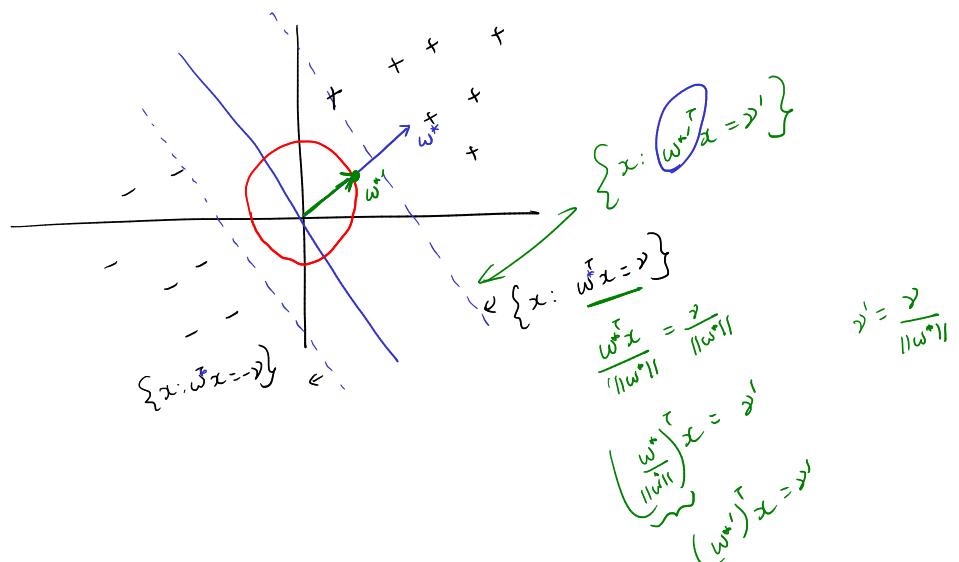
$$\forall i \in D \quad \|x_i\|_2 \leq R \quad \text{for some } R > 0.$$



(3)

Without loss of generality,

assume $\|w^*\| = 1$



ANALYSIS OF "MISTAKES" OF
PERCEPTRON

- Observe that an update happens only when a mistake occurs.
- Say w^e is the current guess and a mistake happens w.r.t (x, y)

$$w^{e+1} = w^e + x \cdot y$$

$$\|w^{e+1}\|^2 = \|w^e + x \cdot y\|^2$$

$$= (w^e + x \cdot y)^T (w^e + x \cdot y)$$

$$= \|w^e\|^2 + \underbrace{2 \cdot (w^e \cdot x) y}_{\leq 0} + \underbrace{\|x\|^2 \cdot y^2}_{\leq R^2}$$

because
mistake.

$$\|w^{e+1}\|^2 \leq \|w^e\|^2 + R^2$$

$$\|w^{e+1}\|^2 \leq (\|w^e\|^2 + R^2) + R^2$$

$$\|w^{e+1}\|^2 \leq \|w^e\|^2 + 2R^2$$

$\Rightarrow \boxed{\|w^{e+1}\|^2 \leq 2R^2}$

①

$$(w^{e+1})^T w^* = (w^e + x \cdot y)^T w^*$$

$$= w^e \cdot w^* + \underbrace{(w^e \cdot x) y}_{\geq 0}$$

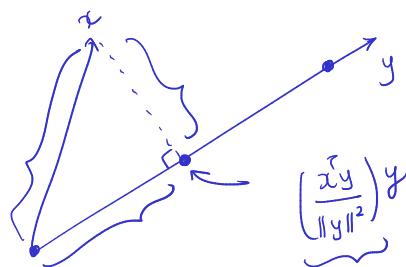
$$(w^{e+1})^T w^* \geq \underbrace{w^e \cdot w^*}_{\geq 0} + \gamma$$

$$\geq (\omega^{l+1}^\top \omega^*) + \gamma$$

$$(\omega^{l+1})^\top \omega^* \geq \underbrace{\omega^0 \omega^*}_{0} + \lambda \gamma$$

$$\Rightarrow (\omega^{l+1})^\top \omega^* \geq \lambda \gamma \quad \text{--- (2)}$$

For any x, y



$$\left\| \frac{x^T y}{\|y\|^2} y \right\|^2 \leq \|x\|^2$$

[Pythagoras]

Cauchy-Schwarz

$$\frac{(x^T y)^2}{\|y\|^2} \leq \|x\|^2$$

$$(x^T y)^2 \leq \|x\|^2 \|y\|^2 \quad \Rightarrow \quad \boxed{(x^T y)^2 \leq \|x\|^2 \|y\|^2}$$

From (2)

$$\lambda \gamma \leq (\omega^{l+1})^\top \omega^*$$

$$\Rightarrow \lambda^2 \gamma^2 \leq ((\underbrace{\omega^{l+1}}_{+ y} \underbrace{\omega^*}_{+ y})^\top)^2 \leq \|\omega^{l+1}\|^2 \underbrace{\|\omega^*\|^2}_{1} \quad \text{[From C.S.]}$$

$$\Rightarrow \boxed{\|\omega^{l+1}\|^2 \geq \lambda^2 \gamma^2} \quad \text{--- (3)}$$

$$\lambda^2 \gamma^2 \leq \|\omega^{l+1}\|^2 \leq \lambda R^2$$

↑
From (3) ↑
From (1)

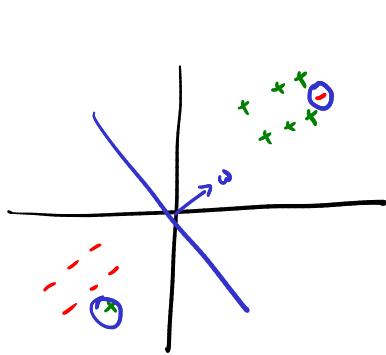
$$\Rightarrow \lambda^2 \gamma^2 \leq \lambda R^2 \quad \lambda \leftarrow \# \text{ mistakes}$$

$$\Rightarrow \boxed{\lambda \leq R^2 / \gamma^2} \quad \leftarrow \text{RADIUS MARGIN BOUND}$$

\Rightarrow # mistakes is bounded [because $\gamma > 0$].

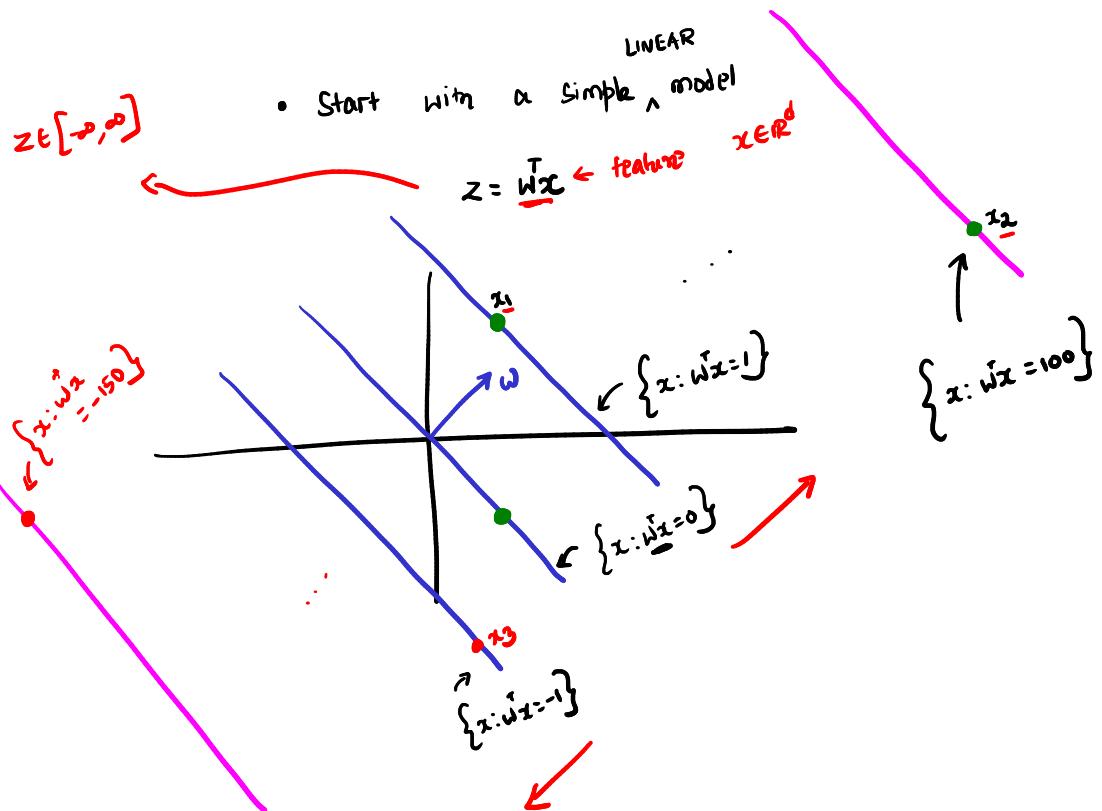
\Rightarrow PERCEPTRON CONVERGES!

Perception mistakes $\leq \frac{R^2}{y^2} \leftarrow \text{Radius}$
 $\leftarrow \text{margin}$



$$\left\{ \begin{array}{l} P(y=1|x) = 1 \quad \text{if } w^T x \geq 0 \\ = 0 \quad \text{otherwise.} \end{array} \right.$$

Can we model probabilities differently?



- Larger the score ($z = w^T x$), more the probability of being +1

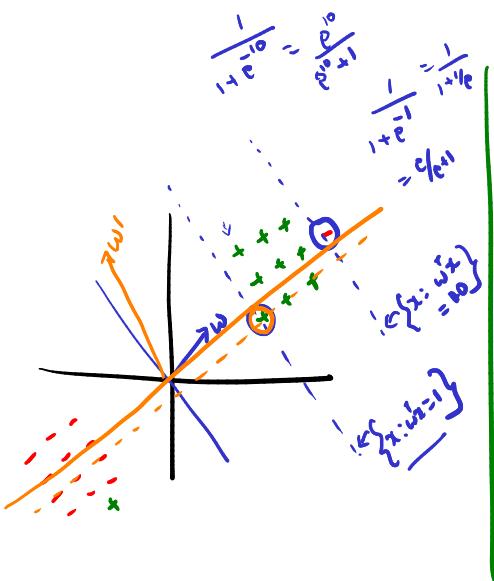
^{Score}

$$g(z) = 0.5 \quad \text{if } z=0$$

 LINK FUNCTION

$$g(z) \rightarrow 1 \quad \text{as } z \rightarrow \infty$$

$$g(z) \rightarrow 0 \quad \text{as } z \rightarrow -\infty$$



one
popular
choice

$$g(z) = \frac{1}{1 + e^{-z}}$$

SIGMOID /
LOGISTIC
FUNCTION

MODEL: LOGISTIC REGRESSION

$$P(y=1/x) = \frac{1}{1 + e^{-w^T x}} = g(w^T x)$$

Dataset: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ $y_i \in \{0, 1\}$.

How to find w : Maximum Likelihood.

$$\mathcal{L}(w; \text{Data}) = \prod_{i=1}^n \left(g(w^T x_i) \right)^{y_i} \left(1 - g(w^T x_i) \right)^{(1-y_i)}$$

$$\begin{aligned} \log \mathcal{L}(w; \text{Data}) &= \sum_{i=1}^n y_i \log(g(w^T x_i)) + (1-y_i) \log(1 - g(w^T x_i)) \\ &= \sum_{i=1}^n y_i \log\left(\frac{1}{1 + e^{-w^T x_i}}\right) + (1-y_i) \log\left(\frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}}\right) \end{aligned}$$

$$= \sum_{i=1}^n [(1-y_i)(-w^T x_i) - \log(1 + e^{-w^T x_i})]$$

Goal:

$$\boxed{\underset{w}{\operatorname{max}} \sum_{i=1}^n [(1-y_i)(w^T x_i) - \log(1 + e^{-w^T x_i})]} \quad \log L(w)$$

- No closed form expression
- Can perform Gradient descent. [ascent]

$$\nabla \log L(\omega) = \sum_{i=1}^n \left[(y_i)(-x_i) - \left(\frac{e^{-\omega^T x_i}}{1+e^{-\omega^T x_i}} \right) (-x_i) \right]$$

$$= \sum_{i=1}^n \left[-x_i + y_i x_i + x_i \left(\frac{e^{-\omega^T x_i}}{1+e^{-\omega^T x_i}} \right) \right]$$

$$= \sum_{i=1}^n y_i x_i - x_i \left(\frac{1}{1+e^{-\omega^T x_i}} \right)$$

$$= \sum_{i=1}^n x_i \left(y_i - \frac{1}{1+e^{-\omega^T x_i}} \right)$$

Gradient update rule

$$\begin{aligned} w_{t+1} &= w_t + \eta_t \nabla \log L(\omega) \\ &= w_t + \eta_t \left(\sum_{i=1}^n x_i \left(y_i - \frac{1}{1+e^{-\omega^T x_i}} \right) \right) \end{aligned}$$

$\underbrace{\{0,1\}}_{g(w_t^T x_i)}$ $\underbrace{\frac{1}{1+e^{-\omega^T x_i}}}_{\theta_i}$

$x_{\text{test}} \in \mathbb{R}^d$

$\hat{y}_{\text{test}} = \text{Sign}(\hat{w}^T x_{\text{test}})$

KERNEL VERSION

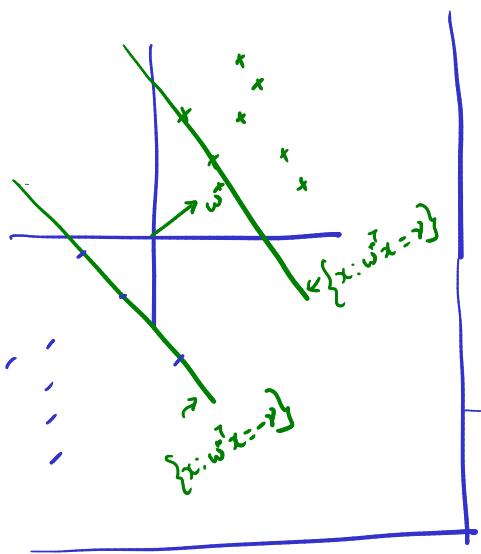
- Can argue $w^* = \sum_{i=1}^n \alpha_i x_i$

[Formal theorem
is called the
Representer
theorem.]

REGULARIZED VERSION

$$\min_w \sum_{i=1}^n \left[\log \left(1 + e^{\omega^T x_i} \right) + \omega^T x_i (1 - y_i) \right] + \frac{\lambda}{2} \|\omega\|^2$$

CROSS VALIDATE Hyper parameters
Regularized



PERCEPTRON

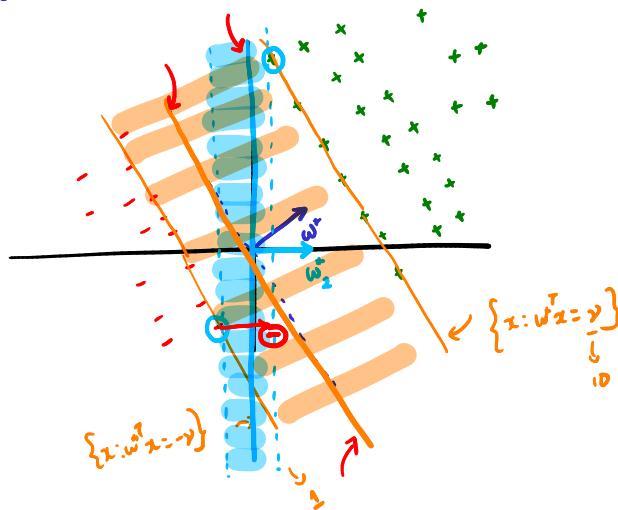
$$\# \text{ mistakes} \leq \frac{R^2}{\gamma^2}$$

$$\|x_i\|^2 \leq R^2$$

Dataset - L.S with margin γ

$$(\bar{w}^T x_i) y_i \geq \gamma + i \quad \gamma > 0$$

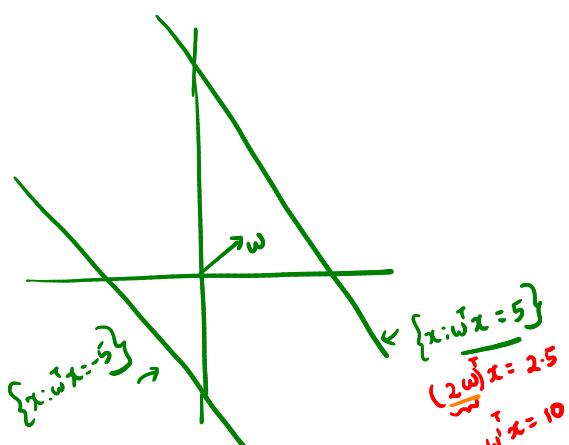
"QUALITY" OF FINAL SOLUTION



Observation

- ① # mistakes depends on the best possible w's margin.
- ② w_{perc} need not necessarily be w. It could be w₂ also (blue line)

Goal: To come up with a formulation that maximizes "margin"



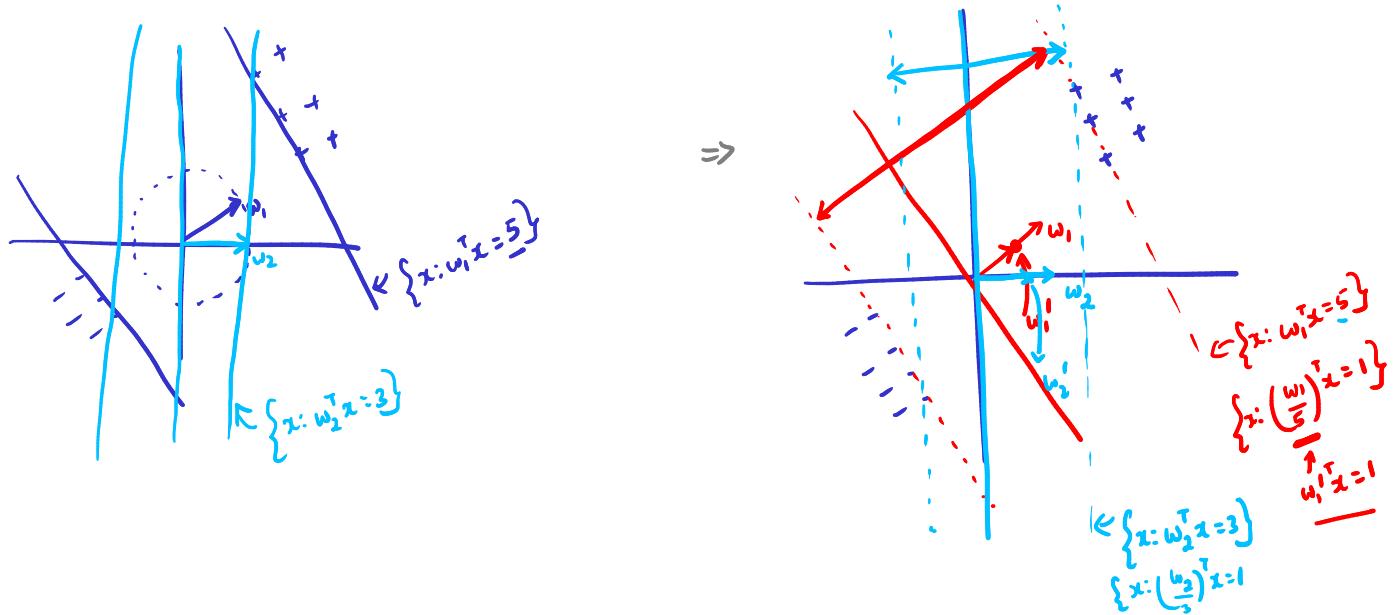
$$\max_{w, \gamma} \gamma$$

such that

$$(\bar{w}^T x_i) y_i \geq \gamma + i$$

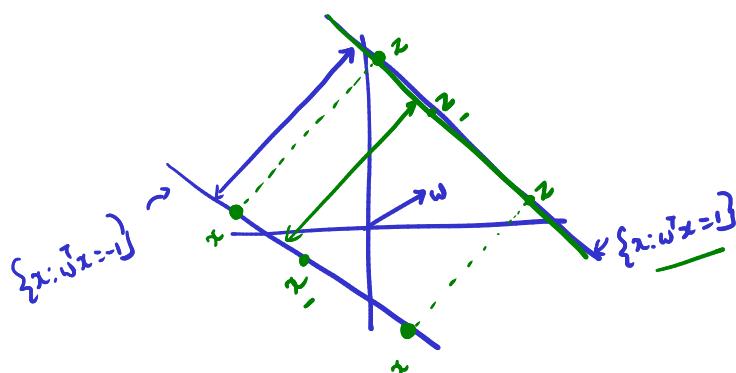
Issue: Can scale w arbitrarily.

$$\boxed{\begin{aligned} \max_{w, \gamma} & \gamma \\ \text{such that} & (\bar{w}^T x_i) y_i \geq \gamma + i \\ & \|w\|^2 = 1 \end{aligned}}$$



$$\boxed{\begin{array}{ll} \max_w & \underline{\text{width}(w)} \\ \text{s.t.} & (\omega^T x_i) y_i \geq 1 \end{array}}$$

What is width(w) ?



$$\min_z \frac{1}{2} \|z - z'\|^2$$

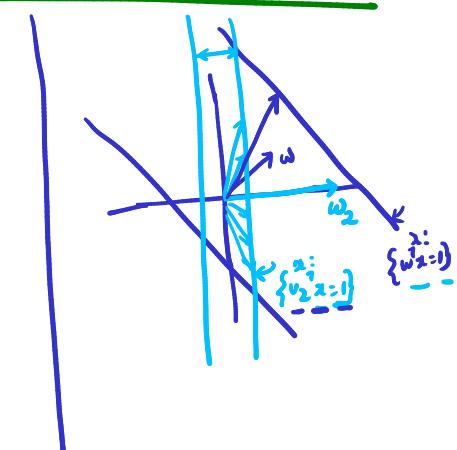
s.t.

$$\omega^T z = +1$$

$$\omega^T z = -1$$

Solution: $\text{width}(w) = \boxed{\frac{2}{\|w\|^2}}$

$$\boxed{\begin{array}{ll} \max_w & \frac{2}{\|w\|^2} \\ \text{s.t.} & (\omega^T x_i) y_i \geq 1 \end{array}}$$



Equivalently

$$\min_w \frac{1}{2} \|w\|^2$$

$$\text{s.t. } (w^T x_i) y_i \geq 1 \quad \forall i$$

$$\begin{array}{ll} \min_{w \in \mathbb{R}^d} & \frac{1}{2} \|w\|^2 \\ \text{s.t.} & \underline{x_i} (w^T x_i) y_i \geq 1 \end{array}$$

- A

DETOUR

$$\begin{array}{ll} \min_w & f(w) \\ \text{s.t.} & g(w) \leq 0 \end{array}$$

$$L(w, \alpha) = f(w) + \alpha g(w)$$

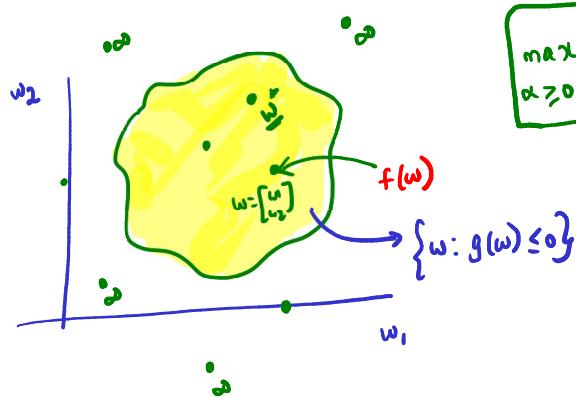
Fix any w .

Consider

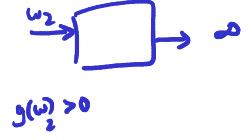
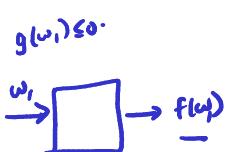
$$\max_{\alpha \geq 0} L(w, \alpha) = \max_{\alpha \geq 0} \underbrace{f(w)}_{\alpha} + \underbrace{\alpha g(w)}_{\geq 0}$$

w	$f(w)$	$g(w)$	$L(w, \alpha)$
$[1 2 3 4]$	-100	5	$\boxed{100}$
	$\max_{\alpha \geq 0} -100 + \alpha 5$		$\boxed{100}$
	$\alpha=1$	-95	$\boxed{95}$
	$\alpha=10$	-60	$\boxed{50}$
	$\alpha=100$	400	$\boxed{100}$

$$\left\{ \begin{array}{ll} \infty & g(w) > 0 \\ f(w) & g(w) \leq 0 \end{array} \right.$$



$$\max_{\alpha \geq 0} f(w) + \alpha g(w)$$



$$\min_w f(w) = \min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right]$$

- Can we swap min and max in (B)?
 - In general, No! But if f and g are "nice" functions [convex functions], then yes!
↳ [Quadratic / Linear]

$$\min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right] = \max_{\alpha \geq 0} \left[\min_w f(w) + \alpha g(w) \right]$$

Convex f and g.

For convex f and g .

For multiple constraints

$$\begin{array}{ll}
 \boxed{\min_w f(w)} \\
 \text{s.t. } g_i(w) \leq 0 \quad i=1 \dots k
 \end{array} \quad = \quad \min_w \left[\max_{\substack{d_1, \dots, d_k \\ \geq 0}} \geq 0 \cdot \geq 0 \right] f(w) + \alpha_1 g_1(w) + \alpha_2 g_2(w) + \dots + \alpha_k g_k(w)$$

$$\min_w \frac{1}{2} \|w\|^2$$

← Quadratic in w

s.t.

$$(\vec{w}^\top x_i) y_i \geq 1 \quad \forall i = 1, \dots, n$$

← Linear in w

\equiv

$$\underbrace{1 - (\vec{w}^\top x_i) y_i}_{g_i(w)} \leq 0 \quad \forall i = 1, \dots, n$$

$$L(\omega, \alpha) = \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \alpha_i (1 - (\omega^T x_i) y_i)$$

$$\min_w \max_{\alpha \geq 0} \left[\frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - (w^T x_i) y_i) \right] = \underbrace{\max_{\alpha \geq 0}}_{\frac{1}{w}} \left[\min_w \left[\frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - (w^T x_i) y_i) \right] \right]$$

Fix Some $\alpha \geq 0$. $\alpha = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix}$

$$\min_w \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - \underbrace{\langle w^T x_i \rangle}_{y_i})$$

$$w_d^* + \sum_{i=1}^n \alpha_i (-x_i y_i) = 0$$

$$w_d^* = \sum_{i=1}^n \alpha_i x_i y_i$$

≥ 0 $\{x_i, y_i\}$

→ Substitute back value of w_d^* in the objective

$$w_d^* = \underline{X Y \alpha}$$

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}_{d \times n} \begin{bmatrix} y_1 & \dots & 0 \\ 0 & \dots & \vdots \\ y_n & \dots & 0 \end{bmatrix}_{n \times n} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

$$\frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - \langle w^T x_i \rangle y_i)$$

Substitute $w_d^* = \underline{X Y \alpha}$ into

on Simplification

$$= \alpha^T 1 - \frac{1}{2} (X Y \alpha)^T (X Y \alpha)$$

PRIMAL

$$\min_w \frac{1}{2} \|w\|^2$$

$$\text{s.t. } (I - X)w \leq 0$$

DUAL PROBLEM

$$\max_{\alpha \geq 0} \alpha^T 1 - \frac{1}{2} \alpha^T Y^T X^T X Y \alpha$$

$\in \mathbb{R}^{n \times n}$

Kernel K

What have we gained?

- Dual variable dimension is \mathbb{R}^n while primal problem dimension is \mathbb{R}^d

- Dual constraints are "easier"

- More importantly dual depends on $x^T x$ and so

can be "KERNELIZED"!

$$w^*_{\alpha^*} = \sum_{i=1}^n \alpha_i^* x_i y_i$$

→ This says optimal w^* is a linear combination of the data points where importance of a datapoint is given by α_i^* (for i^{th} point)

→ Question: Where are the "IMPORTANT" points? (i.e., points for which $\alpha_i^* > 0$)

REVISITING THE LAGRANGIAN

$$\min_w \left[\max_{\alpha \geq 0} f(w) + \alpha g(w) \right] = \max_{\alpha \geq 0} \left[\min_w f(w) + \alpha g(w) \right]$$

w^* is the primal solution

α^* is the dual solution

$$\max f(w^*) + \alpha^* g(w^*) = \min_w f(w) + \alpha^* g(w)$$

$$f(\omega^*) = \min_{\omega} f(\omega) + \alpha^* g(\omega)$$

$$\leq f(\omega^*) + \alpha^* g(\omega^*)$$

$$\Rightarrow f(\omega^*) \leq f(\omega^*) + \alpha^* g(\omega^*)$$

$$\Rightarrow \boxed{\alpha^* g(\omega^*) \geq 0} \quad -①$$

But we already know $\underline{\alpha^* \geq 0}$ & $\underline{g(\omega^*) \leq 0}$

$$\Rightarrow \boxed{\alpha^* g(\omega^*) \leq 0} \quad -②$$

$$① \text{ } \& \text{ } ② \Rightarrow \boxed{\alpha^* g(\omega^*) = 0} \rightarrow \text{COMPLEMENTARY SLACKNESS}$$

For multiple constraints

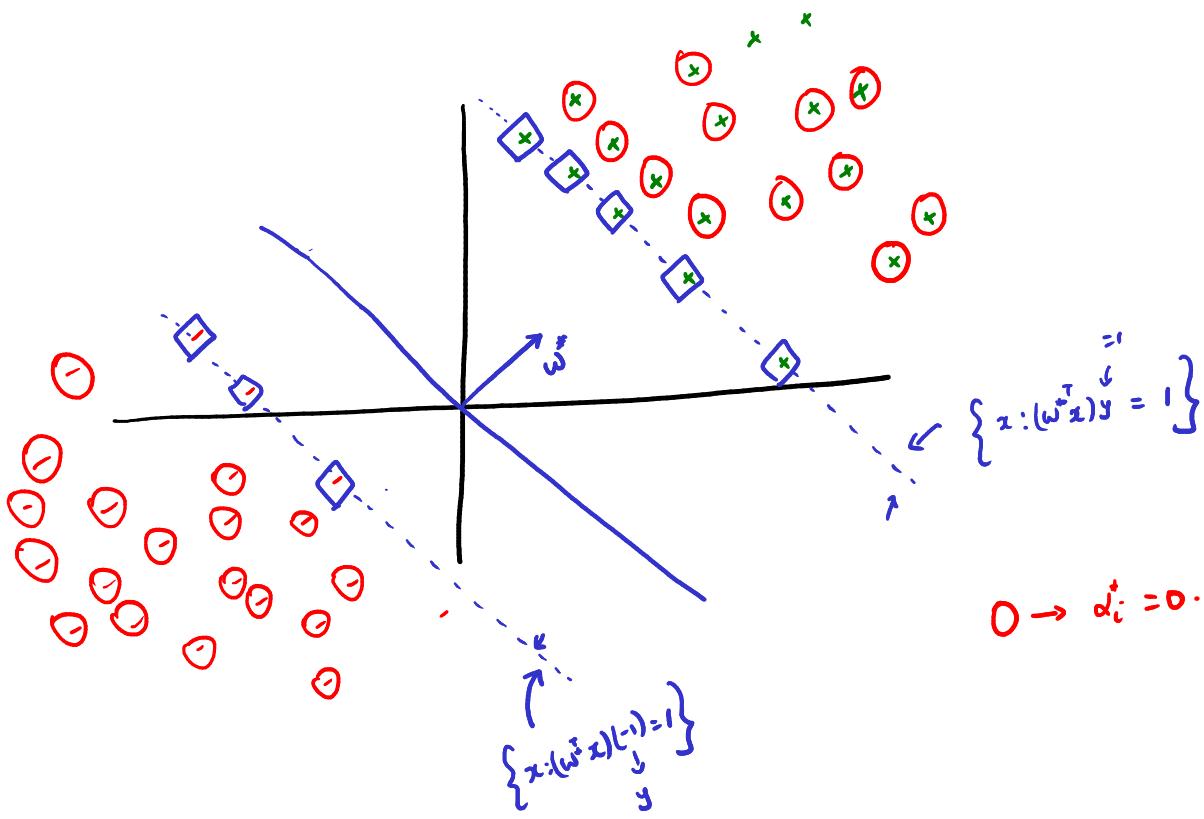
$$\boxed{\alpha_i^* g_i(\omega^*) = 0 \quad \forall i}$$

In our problem

$$\underline{\alpha_i^*} \left(1 - \frac{(\omega^T x_i) y_i}{g_i(\omega^*)} \right) = 0 \quad \forall i \quad \left[\begin{array}{l} \text{by} \\ \text{complementary} \\ \text{slackness} \end{array} \right]$$

$$\Rightarrow \text{If } \underline{\alpha_i^*} > 0 \quad \stackrel{\text{C.S}}{\Rightarrow} \quad 1 - (\omega^T x_i) y_i = 0$$

$$\boxed{(\omega^T x_i) y_i = 1}$$



- ▶ only the points that are on the "SUPPORTING" hyperplane can contribute to w^*
- ▶ These special points are called "SUPPORT VECTORS"
- ▶ ALGORITHM → SUPPORT VECTOR MACHINE (SVM)
 - [Vapnik et.al]
- ▶ w^* is a sparse linear combination of the data points.

Given

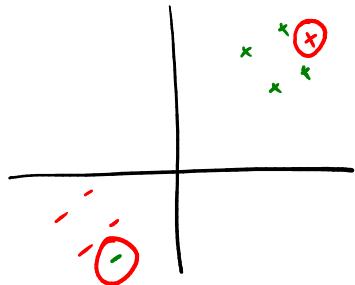
$$x_{\text{test}} = w^T x_{\text{test}} = \left(\sum_{i=1}^n \alpha_i^+ y_i \right)^T x_{\text{test}}$$

$$= \sum_{i=1}^n \alpha_i^+ y_i (x_i^T x_{\text{test}})$$

$$x_{\text{test}}^T w^T \phi(x_{\text{test}}) = \sum_{i=1}^n \alpha_i^+ y_i k(x_i, x_{\text{test}})$$

QUESTIONS

- How to adapt the SVM algorithms when data has outliers.



- KERNELS can help but is not the right way to solve this!

$$\begin{aligned} \min_w & \frac{1}{2} \|w\|^2 \\ \text{st } & (\underline{w^T z_i}) y_i \geq 1 \quad \forall i \end{aligned}$$

Insight: Make every w feasible.

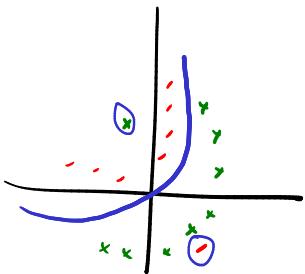
- Fix any w . w classifies some points correctly and misclassifies some points
- The incorrectly classified points "pay bribe" to go to the "correct" side!

SOFT MARGIN PRIMAL FORMULATION

$$\begin{aligned} \text{MODIFIED FORMULATION} & \min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \quad \xrightarrow{\geq 0} \text{HYPER PARAMETER.} \\ \text{st} & (\underline{w^T z_i}) y_i + \xi_i \geq 1 \quad \forall i \\ & \xi_i \geq 0 \quad \forall i \end{aligned}$$

$C = 0$
 \Rightarrow Bribes don't cost
 $\Rightarrow w = 0 \in \mathbb{R}^d$ is the solution

$C = \infty$
 \Rightarrow Linear separable case.



$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

st

$$\rightarrow w^T x_i y_i + \xi_i \geq 1 \quad \leftarrow 1 - w^T x_i y_i - \xi_i \leq 0$$

$$\rightarrow \xi_i \geq 0 \quad \leftarrow -\xi_i \leq 0$$

$$L(w, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - w^T x_i y_i - \xi_i)$$

$$+ \sum_{i=1}^n \beta_i (-\xi_i)$$

DUAL PROBLEM

$$\min_{w, \xi} \left[\max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - w^T x_i y_i - \xi_i) + \sum_{i=1}^n \beta_i (-\xi_i) \right]$$

III Duality

$$\max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \left[\min_{w, \xi} \left[\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - w^T x_i y_i - \xi_i) + \sum_{i=1}^n \beta_i (-\xi_i) \right] \right]$$

Fix $\underline{\alpha, \beta}$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow \boxed{w_{\alpha, \beta}^* = \sum_{i=1}^n \alpha_i x_i y_i} \quad -①$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow C + \alpha_i (-1) + \beta_i (-1) = 0$$

$$\Rightarrow \boxed{\alpha_i + \beta_i = C} \quad -②$$

Back Substituting $w_{\alpha, \beta}^*$ into the Lagrangian

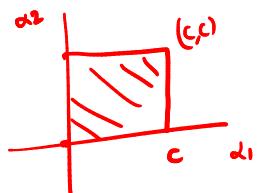
$$w_{\alpha, \beta}^* = \underline{x y \alpha}$$

Dual problem

$$\begin{aligned} \max_{\alpha \geq 0} \quad & \frac{\alpha^T - \frac{1}{2} \alpha^T y^T X^T X \alpha}{\alpha^T y} \\ \text{subject to} \quad & \alpha^T y = 1 \\ & \alpha \geq 0 \\ & \alpha^T + \beta^T = C \\ & \alpha_i + \beta_i = c \cdot x_i \end{aligned}$$

↑ no β term

III



Box Constraints

$$\max_{0 \leq \alpha \leq C} \alpha^T - \frac{1}{2} \alpha^T y^T X^T X \alpha$$

Kernelizable

$$\text{if } C = 0 \Rightarrow \alpha^* = 0 \in \mathbb{R}^n \Rightarrow w^* = \sum \alpha_i^* x_i y_i \Rightarrow w^* = 0$$

$C = \infty \Rightarrow$ Hard-margin

What do the COMPLEMENTARY SLACKNESS (C-S)
say about the SOFT-MARGIN SVM?

Let (w^*, ξ^*) be the primal optimal solutions.

Let (α^*, β^*) be the dual optimal solutions.

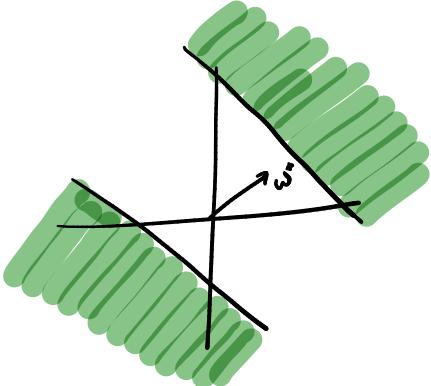
$$\begin{aligned} \stackrel{\text{CS}}{\underline{1}} \quad & \forall i \quad \underline{\alpha_i^*} \left(1 - w^T x_i y_i - \xi_i^* \right) = 0 \quad \leftarrow \\ \stackrel{\text{CS}}{\underline{2}} \quad & \forall i \quad \underline{\beta_i^*} (\xi_i^*) = 0 \end{aligned}$$

Various cases possible

$$\underline{1} \quad \underline{\alpha_i^*} = 0 \Rightarrow \underline{\beta_i^*} = C \quad [\underline{\alpha_i^* + \beta_i^* = c}]$$

↓ CS ②

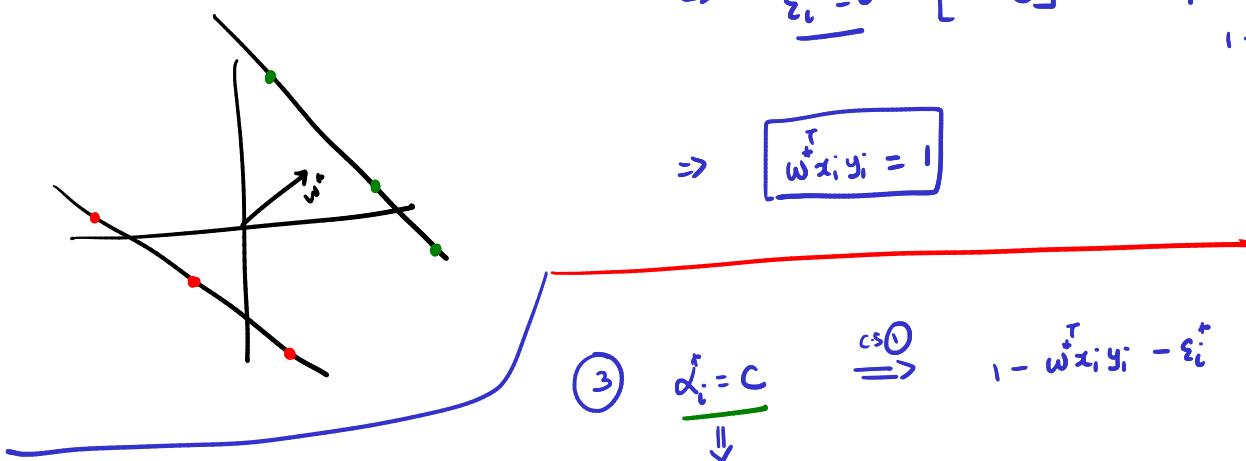
$$\xi_i^* = 0$$



We know $w^T x_i y_i + \hat{\epsilon}_i \geq 1$
 $\Rightarrow w^T x_i y_i \geq 1 \Rightarrow w^* \text{ classifies } (x_i, y_i) \text{ correctly!}$

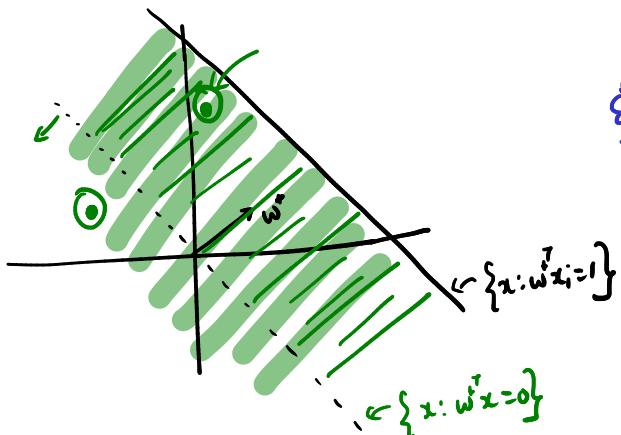
② $\hat{\epsilon}_i \in (0, c)$ $0 < \hat{\epsilon}_i < c \Rightarrow [c.s. ①]$

 $\Rightarrow \hat{p}_i \in (0, c) \quad 0 < \hat{p}_i < c \quad |_{1 - w^T x_i y_i - \hat{\epsilon}_i = 0}$
 $\Rightarrow \hat{\epsilon}_i = 0 \quad [c.s. ②] \quad |_{1 - w^T x_i y_i = 0}$



③ $\hat{\epsilon}_i = c \stackrel{c.s. ①}{\Rightarrow} 1 - w^T x_i y_i - \hat{\epsilon}_i = 0$

$\hat{p}_i = 0 \quad \hat{\epsilon}_i = 1 - w^T x_i y_i \geq 0$
 $\hat{\epsilon}_i \geq 0 \quad \Rightarrow \boxed{w^T x_i y_i \leq 1}$



Points where either
 x_i is incorrectly classified
 by w^* or
 correctly classified but with
 margin ≤ 1

Let's see this from the primal point of view.

Case ①

$\underline{w^T x_i y_i < 1}$

$w^T x_i y_i + \hat{\epsilon}_i \geq 1 \Rightarrow \hat{\epsilon}_i \geq 1 - \underline{w^T x_i y_i}$
 $\Rightarrow \hat{\epsilon}_i > 0 \stackrel{c.s. ②}{\Rightarrow} \hat{p}_i = 0 \Rightarrow \underline{\hat{\epsilon}_i = c}$

Case ②

$$\omega^T x_i y_i = 1$$

$$\xi_i^* \geq 1 - \omega^T x_i y_i$$

$$\xi_i^* \geq 0 \Rightarrow \alpha_i^* \in [0, C]$$

Case ③

$$\boxed{\omega^T x_i y_i > 1}$$

$$\Rightarrow 1 - \underbrace{\omega^T x_i y_i} - \underbrace{\xi_i^*} < 0$$

$$\stackrel{\text{as } ①}{\Rightarrow} \boxed{\alpha_i^* = 0}$$

SUMMARY

$$\alpha_i^* = 0 \Rightarrow \omega^T x_i y_i \geq 1$$

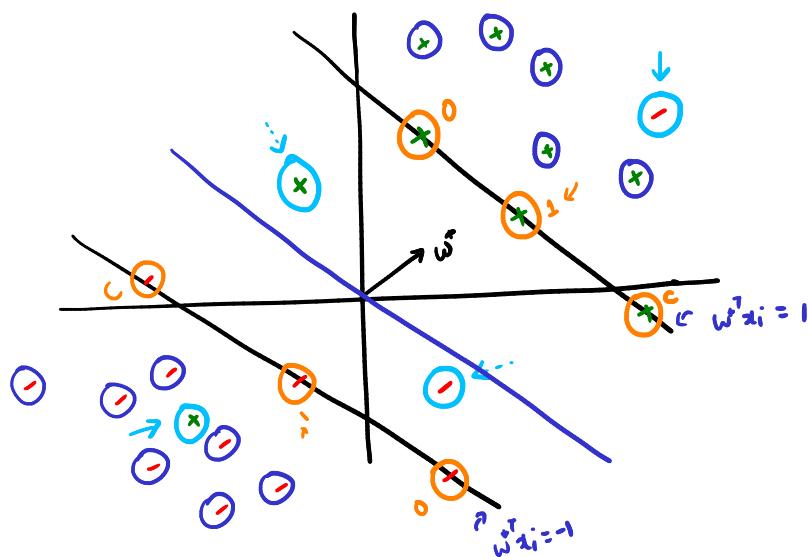
$$0 < \alpha_i^* < C \Rightarrow \omega^T x_i y_i = 1$$

$$\alpha_i^* = C \Rightarrow \omega^T x_i y_i \leq 1$$

$$\omega^T x_i y_i < 1 \Rightarrow \alpha_i^* = C$$

$$\omega^T x_i y_i = 1 \Rightarrow 0 \leq \alpha_i^* \leq C$$

$$\omega^T x_i y_i > 1 \Rightarrow \alpha_i^* = 0$$



$$0 \rightarrow \alpha_i^* = 0$$

$$0 \rightarrow \alpha_i^* = C$$

$$0 \rightarrow \alpha_i^* \in [0, C]$$

So-far

Generative

- Naive-Bayes
- G-D-A/G-N-B

Discriminative

- K-NN
- Decision trees
- Logistic regression
- Perceptron
- Support-vector machines

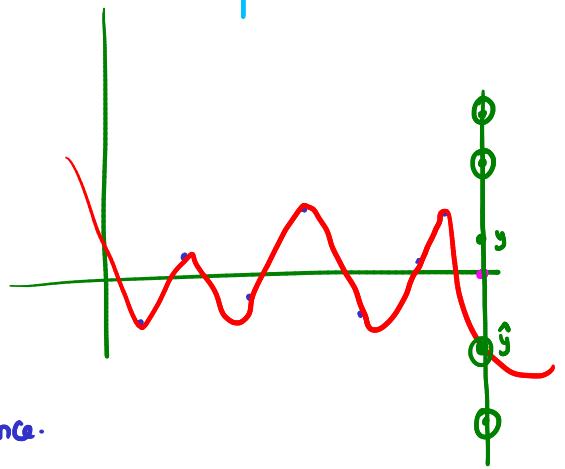
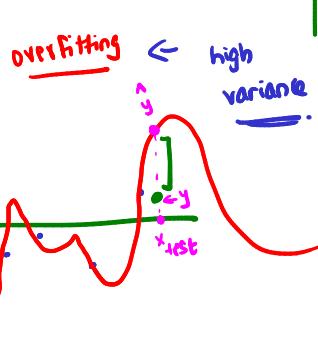
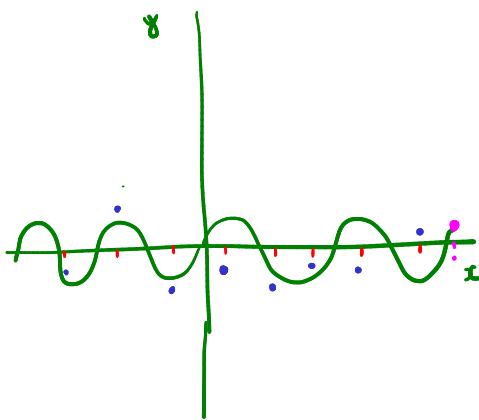
Goal :- Meta classifiers (or) Ensemble classifiers.

WEAK LEARNERS → STRONG LEARNERS.

(better than random)

overfitting - fit noise

underfitting - missing out on structure
thinking it is noise.



$$\text{Error} = \text{bias} + \text{variance}.$$

BAGGING

- BOOTSTRAP AGGREGATION

$$\{x_1, x_2, \dots, x_n\} \sim N(\mu, 1)$$

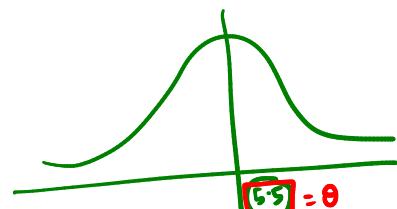
$$\begin{aligned}\hat{\theta}_1 &= x_1 \\ \hat{\theta}_2 &= x_2 \\ &\vdots \\ \hat{\theta}_n &= x_n\end{aligned}$$

$$\hat{\theta}_{\text{ML}} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$E[\hat{\theta}_{\text{ML}}] = \mu$$

$$E[\hat{\theta}] = \theta \leftarrow \text{unbiased estimator}$$

$$E[\hat{\theta}_i] = E[x_i] = \mu$$



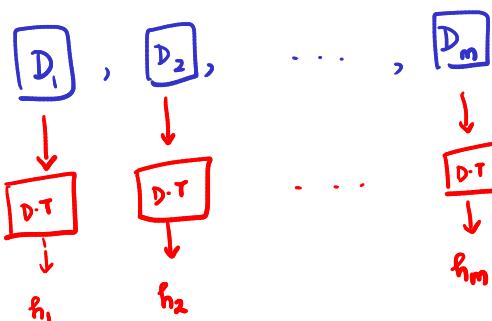
$$\{5, 7.3, 6.7, 8.2, \dots, 9\}$$

$$\begin{aligned}\hat{\theta}_1 &= 5 \\ \hat{\theta}_2 &= 7.3 \\ &\vdots \\ \hat{\theta}_n &= 9\end{aligned}$$

"Averaging reduces variance"

Bagging

$$h_i : \mathbb{R}^d \rightarrow \{\pm 1\}$$



Each D_i has n datapoints.

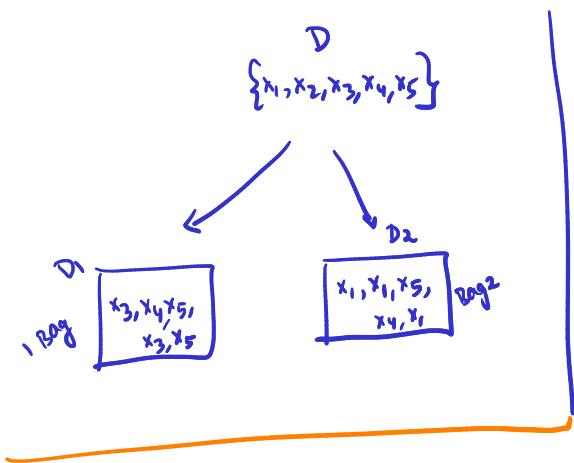
$$\begin{aligned}D_i &= \{x'_1, x'_2, \dots, x'_n\} \\ &\quad y'_1, y'_2, \dots, y'_n \\ x'_i &\in \mathbb{R}^d, y'_i \in \{\pm 1\}\end{aligned}$$

$$f_n^*(x) = \text{Sign}\left(\frac{1}{m} \sum_{i=1}^m h_i(x)\right)$$

$$\begin{aligned}\text{Sign}(z) &= +1 \text{ if } z \geq 0 \\ &= -1 \text{ otherwise}\end{aligned}$$

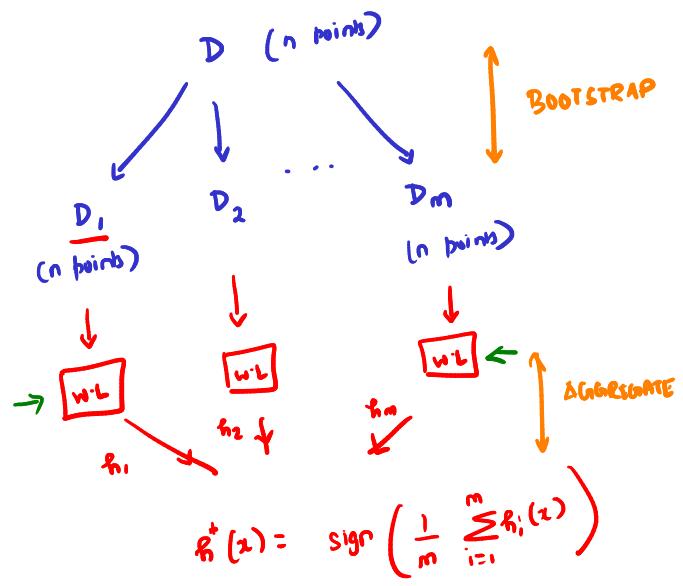
$$\text{Input} \quad D = \{ (x_i, y_i), \dots, (x_n, y_n) \} \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \{\pm 1\} \end{array}$$

Bootstrapping \rightarrow Sampling with replacement

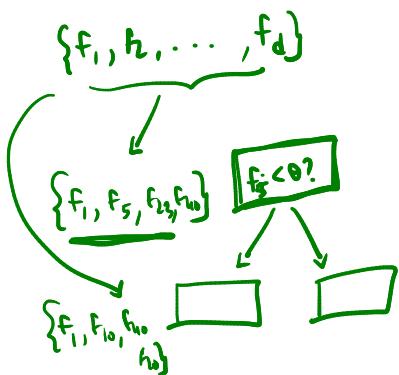


$$P_i(x_j \in D_i) = 1 - \left(1 - \frac{1}{n}\right)^n \approx 1 - \frac{1}{e}$$

$\approx 67\% \text{ for large } n$



Bagging reduces variance!



Random Forest

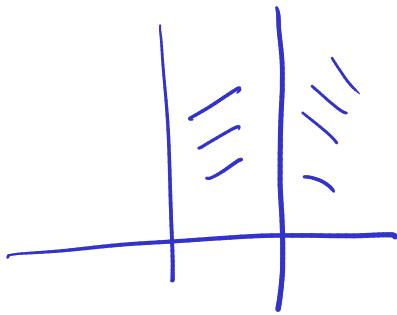
- Bags Decision trees (typically overfit trees)
- Feature bagging (\sqrt{d})

- Bootstrapping - Sampling uniform with replacement
- Bag - Averaging

BOOSTING

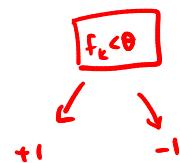
BOOSTING [ADA BOOST]

L FREUND & SCHAPIRA



Weak Learner → Strong Learner

Decision Stumps



1-level or
2-level

ADA-BOOST ALGORITHM

↳ Adaptive.

Input: $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$$x_i \in \mathbb{R}^d \\ y_i \in \{\pm 1\}$$

Initialize $D_0(i) = \frac{1}{n}$

iteration

for $t = 1, \dots, T$

- $h_t = \text{Input}(S, D_t)$ to a weak learner
to get h_t

$$h_t: \mathbb{R}^d \rightarrow \{\pm 1\}$$

$$\tilde{D}_{t+1}(i) = \begin{cases} D_t(i) \cdot \frac{\alpha_t}{e} & \text{if } h_t(x_i) \neq y_i \\ D_t(i) \cdot \frac{-\alpha_t}{e} & \text{if } h_t(x_i) = y_i \end{cases}$$

$$D_{t+1}(i) = \frac{\tilde{D}_{t+1}(i)}{\sum_j \tilde{D}_{t+1}(j)}$$

Do $\begin{bmatrix} x_1 & x_2 & x_3 \\ 0.33 & 0.33 & 0.33 \end{bmatrix} \leftarrow$
 $D_1 \begin{bmatrix} 0.3 & 0.5 & 0.2 \end{bmatrix} \rightarrow h_1$

$$D_2 \begin{bmatrix} \frac{0.3 \times e^{a_1}}{z} & \frac{0.5 \times e^{-a_1}}{z} & \frac{0.2 \times e^{a_1}}{z} \end{bmatrix}$$

$$z = 0.3e^{a_1} + 0.5e^{-a_1} + 0.2e^{a_1}$$

end

h_1, h_2, \dots, h_T

$$h_T^*(x) = \text{sgn} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$$

$$\alpha_t = \ln \sqrt{\frac{1 - \epsilon_t(h_t)}{\epsilon_t(h_t)}}$$

One can prove

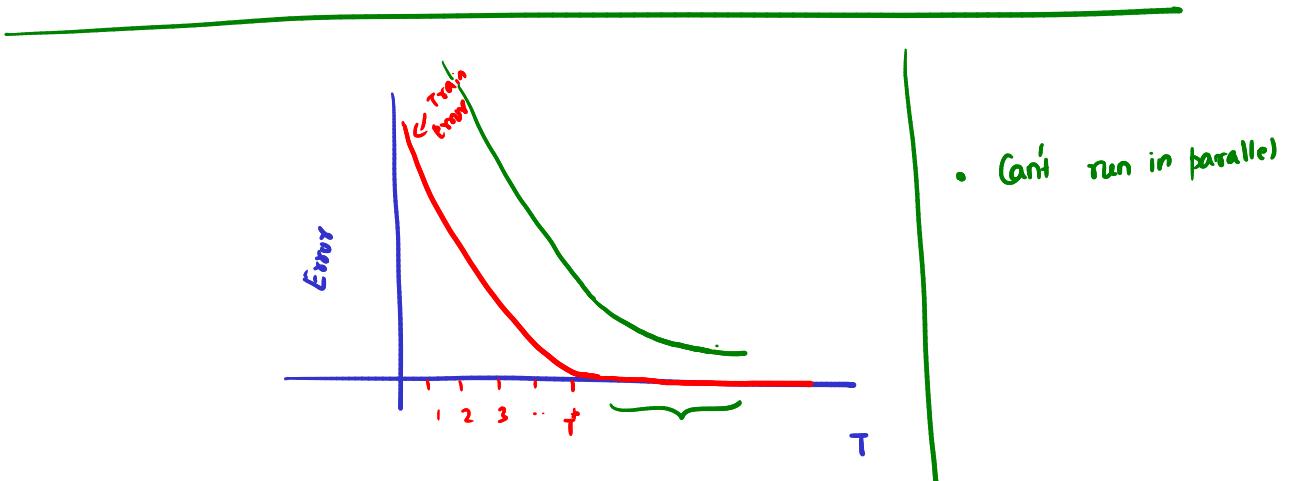
$$\text{If } T \geq \frac{1}{2\gamma^2} \ln(2n)$$

, then

Training error = 0

$$\frac{1}{2\gamma^2}$$

How good is
my weak learner?



$$\text{Dataset} = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad x_i \in \mathbb{R}^d \quad y_i \in \{+1, -1\}$$

$$\text{Goal: } h: \mathbb{R}^d \rightarrow \{\pm 1\}$$

Performance measure

$$\sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

$$\mathbb{1}(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$$\min_{h \in \mathcal{H}_{\text{linear}}} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

$$= \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \mathbb{1}\left(\text{sign}(w^T x_i) \neq y_i\right)$$

NP-HARD

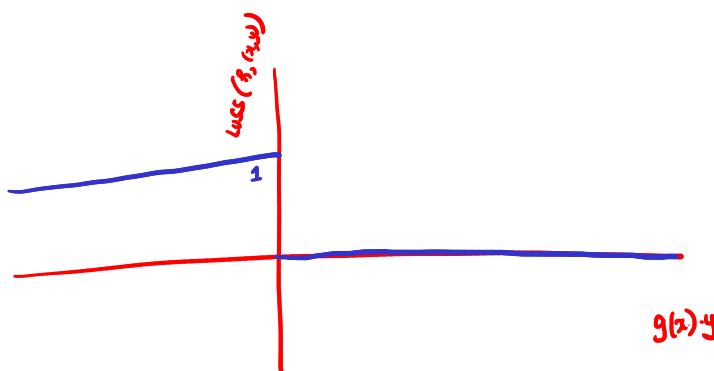
$$\begin{aligned} x &\in \mathbb{R}^d \\ x &= [x_1 \ 1] \in \mathbb{R}^{d+1} \\ w &= [w_1 \ \dots \ w_d \ b] \\ w^T x &= w^T x + b \end{aligned}$$

LOSS-FUNCTION VIEW

$$\frac{\mathbb{1}(x, y)}{\mathbb{R}^d}, \quad h: \mathbb{R}^d \rightarrow \{\pm 1\}$$

$$h(x) = \text{Sign}(w^T x)$$

$$\mathbb{1}(h(x) \neq y) = \frac{\mathbb{1}((w^T x) \cdot y < 0)}{\leq 0 \ -1 \\ > 0 \ -1 \\ < 0 \ +1}$$



$$\sum_{i=1}^n \frac{\mathbb{1}(g(x_i) \cdot y_i < 0)}{1}$$

Alg 1:

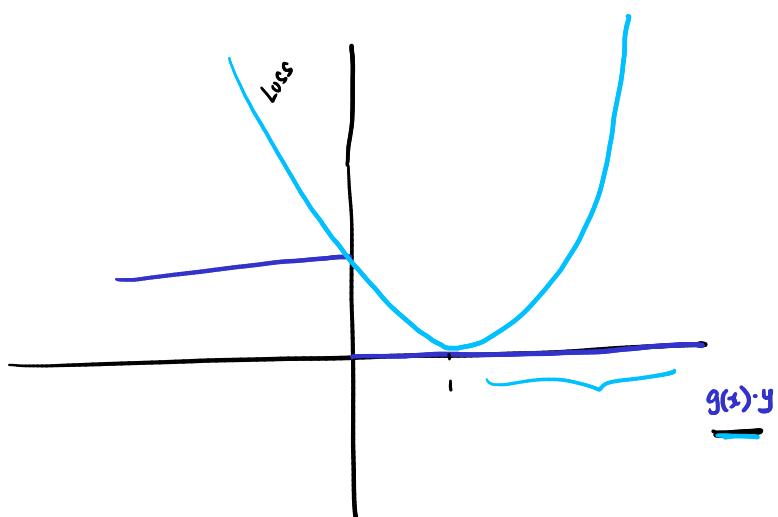
Using Regression for Classification

$$h(x) = \text{Sign}(g(x))$$

$$\begin{aligned} \text{Loss}(g, (x, y)) &= \frac{(g(x) - y)^2}{\frac{g(x) - y}{2}} = (g(x))^2 + y^2 - 2g(x) \cdot y \\ &= \frac{(g(x) - 1)^2}{(g(x) - 1)^2} = (g(x))^2 + 1 - 2g(x) \cdot y - 1 \\ &= (g(x) \cdot y)^2 + 1 - 2g(x) \cdot y \end{aligned}$$

$$= (g(x))^2 + 1 - 2g(x)y - 2$$

$$\textcircled{1} = \textcircled{2}$$



- \rightarrow 0-1 loss
- $\rightarrow (g(x) \cdot y - 1)^2 \rightarrow$ SQUARED LOSS

SUPPORT VECTOR MACHINES

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } (w^T x_i) y_i + \xi_i \geq 1 \\ \xi_i \geq 0$$

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } \begin{cases} \xi_i \geq 1 - (w^T x_i) y_i \\ \xi_i \geq 0 \end{cases}$$

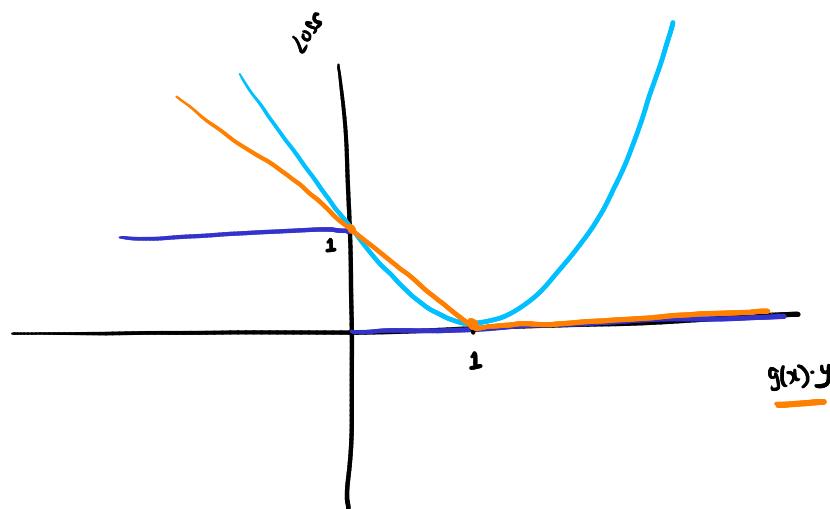
$$\xi_i \geq \max(0, 1 - (w^T x_i) y_i)$$

$$\min_{w \in \mathbb{R}^d} \underbrace{\frac{1}{2} \|w\|^2}_{\text{Regularized}} + C \underbrace{\sum_{i=1}^n \max(0, 1 - (w^T x_i) y_i)}_{\text{Data dependent Loss}} \rightarrow g(x_i)$$

Model dependent

REGULARIZATION +

Data dependent
Loss



- 0-1 loss
- squared loss $(g(x) - 1)^2$
- $\max(0, 1 - g(x) \cdot y)$ HINGE LOSS

LOGISTIC REGRESSION

$$\sigma(a) = \frac{1}{1+e^{-a}}$$

$$\max_w \prod_{i=1}^n \left(\sigma(w^T x_i) \right)^{z_i} \left(1 - \sigma(w^T x_i) \right)^{(1-z_i)}$$

$$\begin{aligned} z_i &= 1 && \text{if } y_i = +1 \\ z_i &= 0 && \text{if } y_i = -1 \end{aligned}$$

$$\max_w \sum_{i=1}^n z_i \log(\sigma(w^T x_i)) + (1-z_i) \log(1 - \sigma(w^T x_i))$$

$$= \min_w \sum_{i=1}^n [-z_i \log(\sigma(w^T x_i)) + (1-z_i) \log(1 - \sigma(w^T x_i))]$$

Loss for a single point $z_i = 1$ ($y_i = +1$)

$$\begin{aligned} &= -\log(\sigma(w^T x_i)) = -\log\left(\frac{1}{1+e^{-w^T x_i}}\right) \\ &= \log\left(1 + e^{-w^T x_i}\right) = \boxed{\log\left(1 + e^{-w^T x_i y_i}\right)} \end{aligned}$$

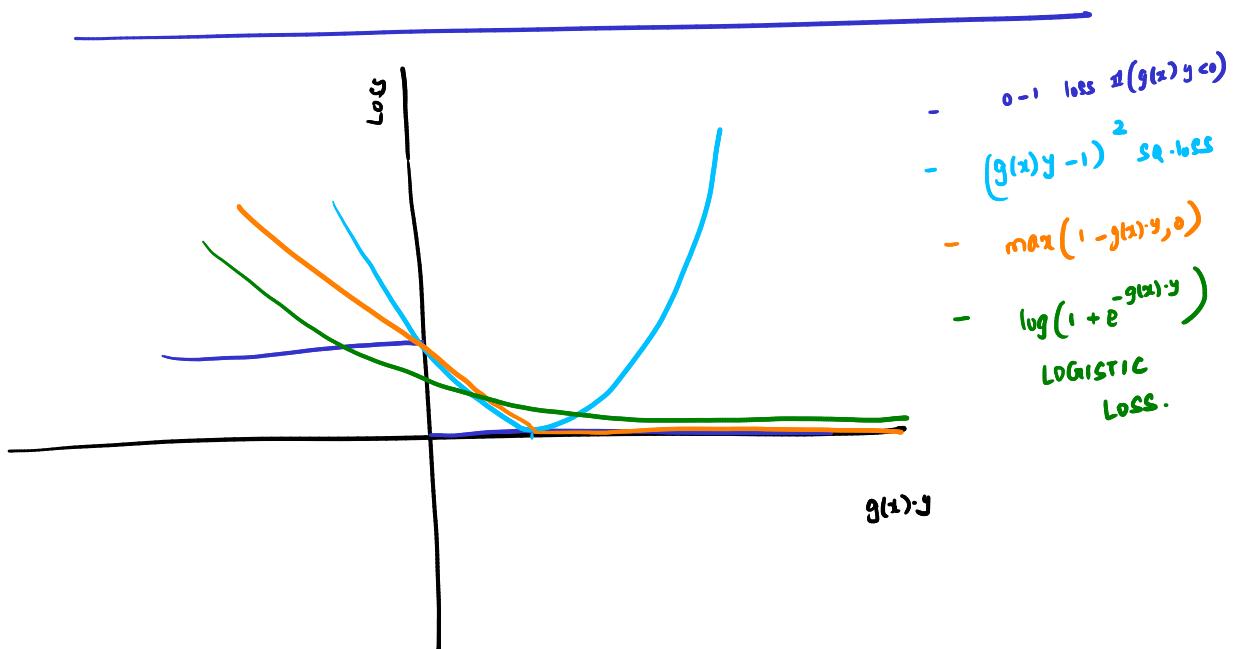
Loss for a single point $z_i = 0$ ($y_i = -1$)

$$\begin{aligned} &= -\log(1 - \sigma(w^T x_i)) \\ &= -\log\left(1 - \frac{1}{1+e^{w^T x_i}}\right) \end{aligned}$$

$$\begin{aligned} &= -\log\left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}}\right) \\ &= -\log\left(\frac{1}{e^{w^T x_i} + 1}\right) \end{aligned}$$

$$\begin{aligned} &= \log\left(1 + e^{-w^T x_i}\right) \\ &= \boxed{\log\left(1 + e^{-w^T x_i y_i}\right)} \end{aligned}$$

$$\equiv \min_{\mathbf{w}} \sum_{i=1}^n \log \left(1 + e^{-\mathbf{w}^T \mathbf{x}_i y_i} \right)$$



CONCLUSIONS

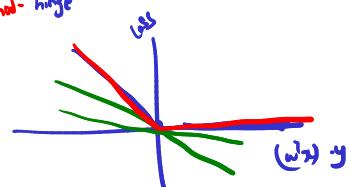
- 0-1 loss is NP-hard to minimize
- Different algorithms use different "surrogate" loss
- Surrogates are convex and hence easy to minimize.

PERCEPTRON

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \underline{x_t y_t}$$

HINGE LOSS

$$l(\mathbf{w}, (\mathbf{x}, y)) = \max(0, -(\mathbf{w}^T \mathbf{x})y)$$



$$\nabla_w \ell_{\text{hinge}} = \begin{cases} -x y & (\underline{\omega^T x}) y < 0 \\ 0 & (\underline{\omega^T x}) y > 0 \\ [-1, 0] x y & (\underline{\omega^T x}) y = 0 \end{cases}$$

↳ chooses $\underline{-x y}$ when

mistake,
 $0 \parallel w$.

$$\omega_{t+1} = \omega_t - \eta \nabla_w \ell_{\text{hinge}}(\omega_t)$$

$$= \omega_t - (-x_t y_t)$$

- Perceptron can be interpreted as S.G.D with modified hinge loss with step size = 1

BOOSTING

$$\bullet \text{Loss}(h, (x, y)) = \frac{-y h(x)}{e} \rightarrow \text{Exponential loss.}$$

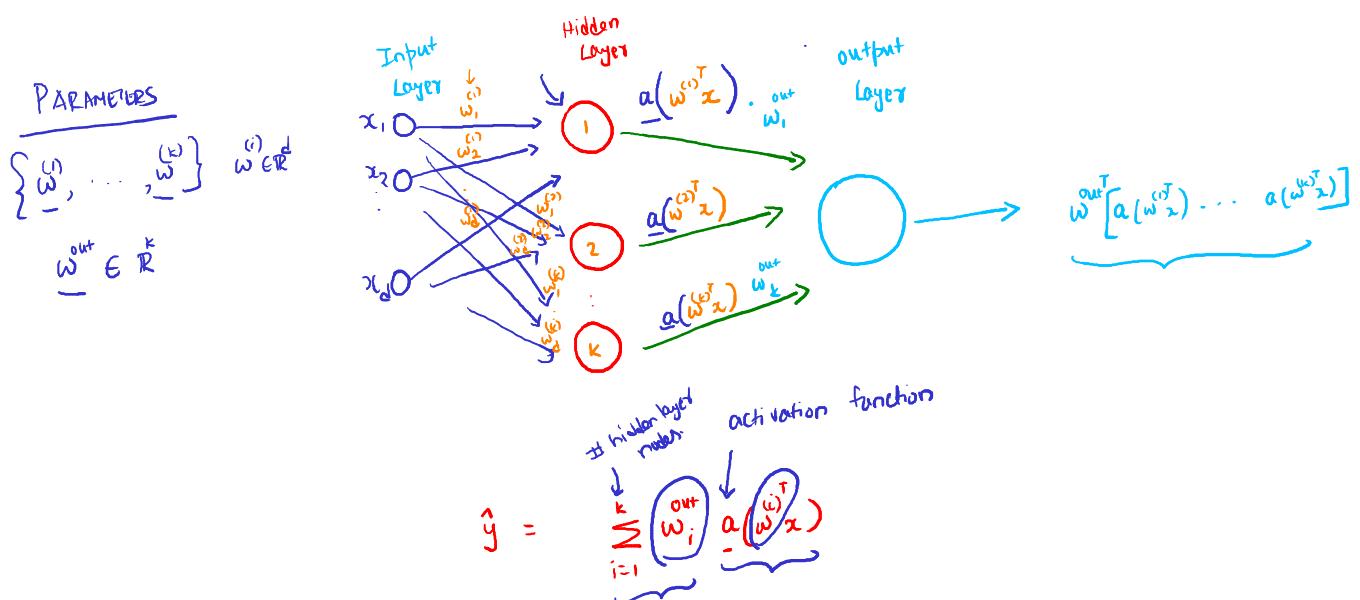
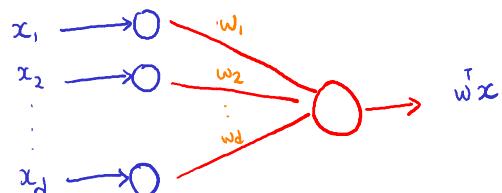
$$\min_w \underbrace{\sum_{i=1}^n L(w^T x_i, y_i)}_{\text{Loss}} + R(w) \quad \hookrightarrow \text{Regularizer}$$

NEURAL NETWORKS

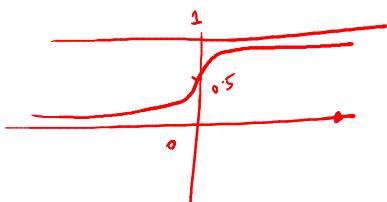
$$x \in \mathbb{R}^d \quad \text{Sign}(w^T x)$$

\downarrow

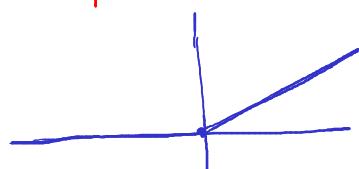
$$[x_1, x_2, \dots, x_d]$$



Examples of activation functions / non-linearities



- $a(z) = \frac{1}{1 + e^{-z}}$ [SIGMOID]



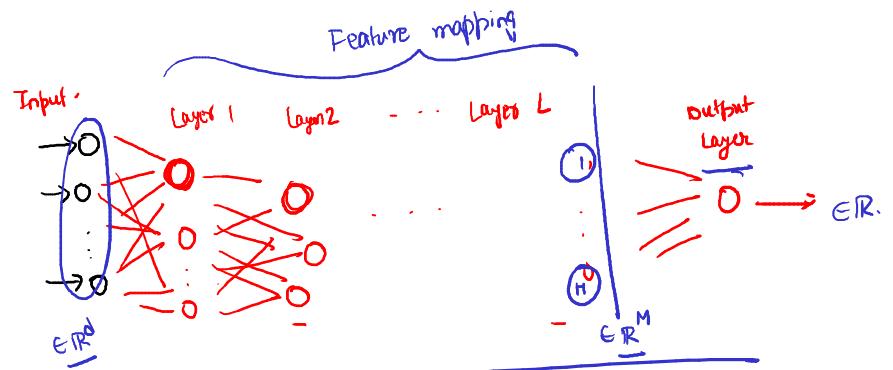
- $a(z) = \max(0, z)$ [Rectified Linear unit]

Regression

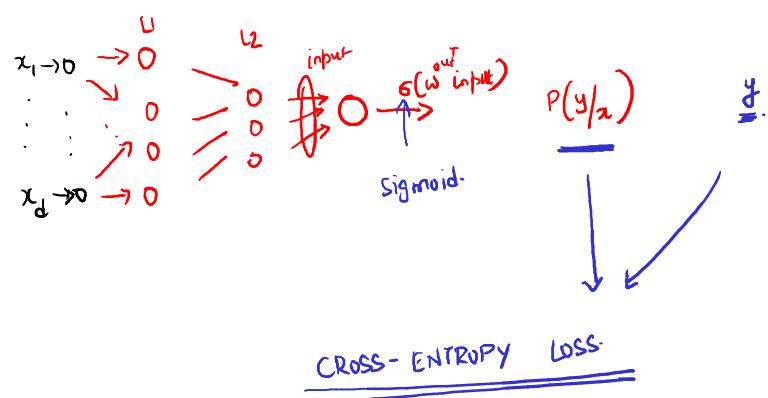
$$L(\text{NN}(x_i; \theta), y_i)$$

$$= \sum_{i=1}^n (\underbrace{\text{NN}(x_i; \theta)}_{\equiv w^T x_i} - y_i)^2$$

Learn θ^* using Gradient descent



- Gradient computed taking advantage of Chain rule \rightarrow BACK-PROPAGATION
- Converges to local minima!



Conclusion

- | | |
|---|--|
| CNN, RNN,
LSTM, Attention
Transformers, | <ul style="list-style-type: none"> • Learn's local minima of non-concave functions • Typically works very well in <u>practical</u> especially for unstructured data. |
|---|--|