# STATISTICAL METHODS IN A.I

## Reference Books:

> Pattern Classification - Duda Hart and Stork
> Machine Learning - A Probabilistic Perspective by Kevin Murphy
> Neural Networks - Simon Haykin
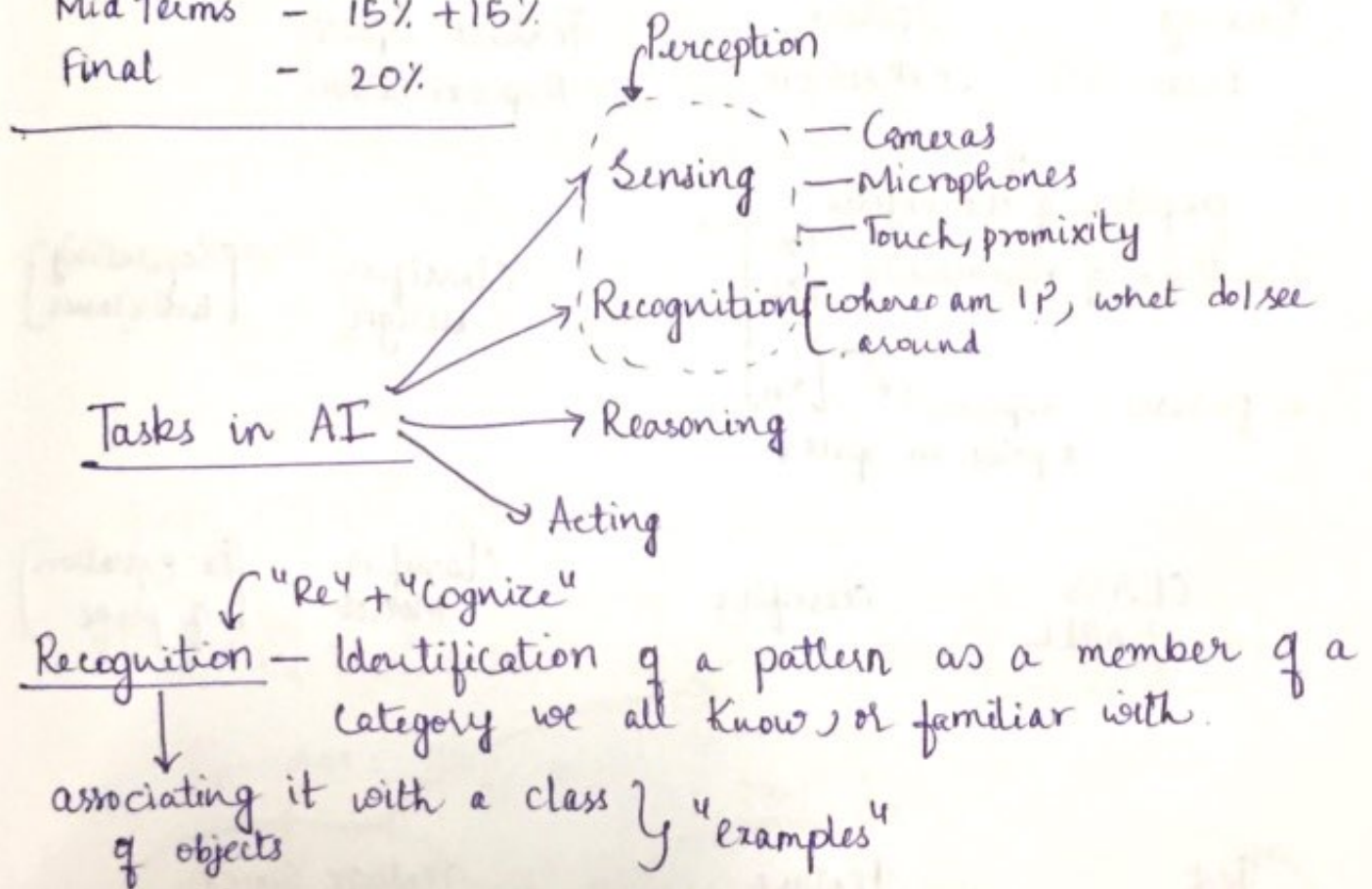> AI - A Modern Approach - Russell and Norwig.

## Grading Schemes -

Mini Project  - 20%
Homeworks    - 30%
Mid Terms    - 15% + 15%
Final        - 20%

Perception

Sensing
 — Cameras
 — Microphones
 — Touch, promixity

Recognition [ wheres am I?, what do I see around

Tasks in AI

→ Reasoning

↳ Acting

Recognition — "Re" + "Cognize"

Recognition — Identification of a pattern as a member of a category we all know, or familiar with.

↓

associating it with a class } "examples"
of objects

Pattern — an entity vaguely defined, that could be given a "name"

→ Same sample might be classified as different classes.

"Class" : collection of "similar" objects
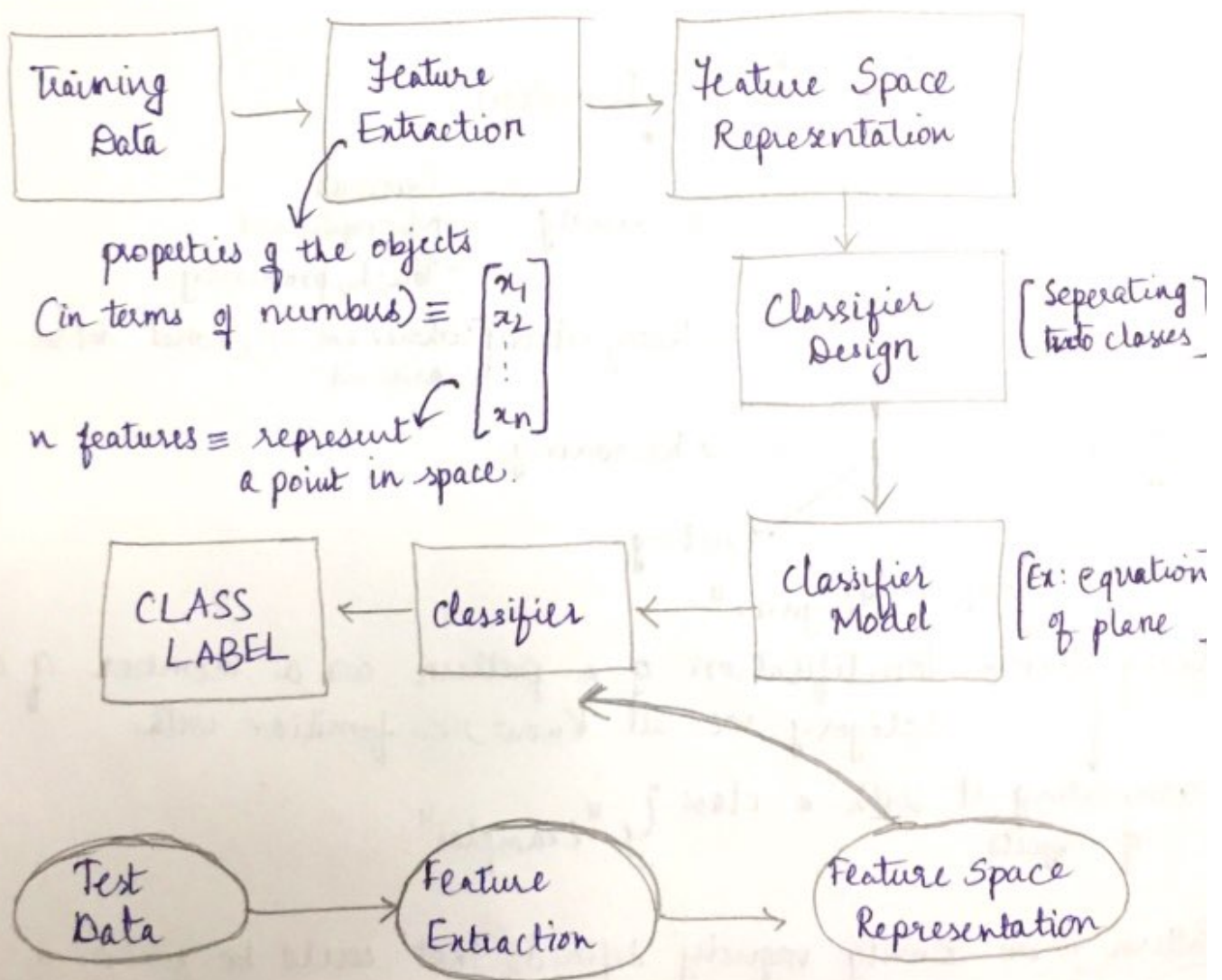 ↳ defined by class samples

Pattern Recognition — inferring a generality from a few examples.

      └→ system "learns" to tell whether or not an object belongs to a class.

* class ⎰ ↗ Inter-class variability
      ⎱ ↘ Intra-class variability

    ↓ represented as $[\omega_1, \omega_2, \ldots \omega_n]$

"train/teach the machine with a "large dataset"."

## Pattern Recognition Process :—

| Training Data | → | Feature Extraction | → | Feature Space Representation |
|---|---|---|---|---|

properties of the objects
(in terms of numbers) $\equiv \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

n features $\equiv$ represent a point in space.

| Classifier Design | [ Seperating two classes ] |
|---|---|

| CLASS LABEL | ← | Classifier | ← | Classifier Model | [ Ex: equation of plane ] |
|---|---|---|---|---|---|

| Test Data | → | Feature Extraction | → | Feature Space Representation |
|---|---|---|---|---|

* class — modeled by a probability density function, $P(x)$

        ↓

class — conditional $\equiv P(x/\omega_i)$

Gaussian Distribution $\equiv$ $P(x) = \dfrac{1}{\sqrt{2\pi}\,\sigma}\, e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \sim$

$$N(\mu, \sigma^2)$$

## Mathematics :- Revision

→ **Linear Algebra :**

$$f(x) \equiv f(\alpha x + \beta y) = \alpha f(x) + \beta f(y) \quad \curvearrowright \text{ for a function to be linear}$$

$$\begin{cases} a_{11}\, x_1 + a_{12}\, x_2 + a_{13}\, x_3 = b_1 \\ a_{21}\, x_1 + a_{22}\, x_2 + a_{23}\, x_3 = b_2 \\ a_{31}\, x_1 + a_{32}\, x_2 + a_{33}\, x_3 = b_3 \end{cases}$$

Basic Representation $\equiv$ matrices and vector form

$$\underbrace{\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}}_{A} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{X \atop (3D\text{-vector})} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}}_{B} \qquad AX = B$$

(Norm/Length)

$\underrightarrow{\underline{\text{Size of Vector}}} \equiv X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$     $|x| \quad \curvearrowleft \text{absolute value}$

$$\|X\|$$

$$\|x\|_2 = \sqrt[2]{x_1^2 + x_2^2 + x_3^2}$$

$$\|x\|_p = \sqrt[p]{|x_1|^p + |x_2|^p + |x_3|^p + \cdots |x_k|^p} \equiv L_p\text{-norm}$$

$$L_p\text{-norm} : \|x\|_p = \sqrt[p]{\left(\sum_{i=1}^{n} |x_i|^p\right)}$$

$\Big\downarrow$ doesn't have any physical significance.

Total distance covered $\equiv L_1$-norm $\Rightarrow \|x\|_1 = \sum_{i=1}^{n} |x_i|$

$$L_0\text{-norm} \equiv \|x\|_0 = \sum_{i=1}^{n} |x_i|^0 \quad \curvearrowright \begin{array}{l}\text{number of non-zero} \\ \text{dimensions / entries in} \\ \text{the vector.}\end{array}$$

$L_\infty$-norm $\equiv \|x\|_\infty = $ | Maximum value of | the entries | $\Big\}$ for Ex: $\begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$L_1 = 3$
$L_2 = 2.24$
$L_3 = 2.08$
$L_4 = 2.03$
⋮
$L_{10} = 2.00019$

$\Rightarrow$ **Span of a set of vectors** —

linearly independent $\equiv$ when any combination of two vectors do not result in the third vector

for Ex: $c = \alpha a + \beta b$ ($c$ is dependent on $a, b$)



$\bar{d} = \bar{a} \alpha + \bar{b} \beta$

If the vectors can be used in expressing the space

are linearly independent, then they can span the whole space.

$M = \begin{bmatrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ x_3 & y_3 & z_3 \end{bmatrix}$

Det of $M = |M| \neq 0$, if rows are linearly independent $\Big\}$

form a basis

if Rank = 2, then out of 3, only 2 vectors are linearly independent & the other one is dependent on the first two.

dimensionality of the subspace spanned:

as in feature $\Rightarrow$ typically, vectors are "column vectors".

$\Rightarrow$ **Orthonormal Basis** — Normal + Perpendicular basis

$\downarrow$ |Norm| $= 1$ (linearly independent)

perpendicular ($90°$) [Dot product of orthogonal vectors $= 0$]

$a \perp b \equiv a \cdot b = ab \cos \theta = 0$.

$\Rightarrow$ $\underline{\text{DOT PRODUCT}} \equiv$ $\overset{X}{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}} \overset{Y}{\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}}$ $\Rightarrow$ $\overset{X.Y}{x_1 y_1 + x_2 y_2 + x_3 y_3} = X^T Y$
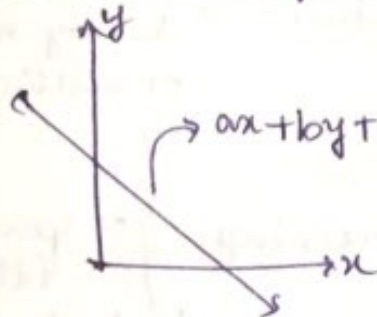
(Inner Product)

$X Y^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} [y_1 \; y_2 \; y_3] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & x_1 y_3 \\ x_2 y_1 & x_2 y_2 & x_2 y_3 \\ x_3 y_1 & x_3 y_2 & x_3 y_3 \end{bmatrix}$ (Outer Product)

$\Rightarrow$ $\underline{\text{CROSS PRODUCT}} \equiv (X, Y) = X \times Y = \begin{vmatrix} i & j & k \\ x_1 & x_2 & x_3 \\ y_1 & y_2 & y_3 \end{vmatrix}$

$\underset{\downarrow}{(a)} = \bar{X} \times \bar{Y}$

perpendicular to both $X$ & $Y$

$\Rightarrow$ $\underline{\text{Representation of Line}}$ :-

is it a vector or a line?
[3D vector]

$\rightarrow ax + by + c = 0$.

$\begin{bmatrix} a \\ b \\ c \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$

$\underset{\omega}{} \quad \underset{x}{}$

$W^T x = 0$  how do you relate both?

Space $\begin{pmatrix} \text{augmented} \\ \text{vector} \end{pmatrix}$

$\rightarrow$ plane perpendicular to $(a, b, c)$

If $\omega$ is scaled, the plane doesn't change and so does the line.

Distance of a point from a plane?

Normalise $\omega$ and take the dot product.

Linear Transformation $\Longleftrightarrow$ Matrix Multiplication
$\downarrow$
by inverting the matrix, we can invert the transformation

$\Rightarrow$ $\underline{\text{Eigen Values and Eigen Vectors}}$

$M = \int u \, \Theta \, v^T \; [\lambda]$ eigen values

$\underset{\text{rows}}{\downarrow} \underset{\text{orthonormal vectors}}{}$

= eigen vectors

↱ physical entity —— represented as a vector?

$X_1 \ldots \qquad X_n \Rightarrow$ vectors $\qquad \qquad \qquad$ with $d \equiv$ dimension

↓

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $d$ —— dimension $\qquad$ where $x_i \in \mathbb{R}^d$

observations.

physical entities are represented differently in different domains.

**Main Problem**

for Ex: $\boxed{x_i \longrightarrow y_i} \in \{0,1\}$ (say, spam or not)

$\{-1, +1\}$

$\begin{cases} y_i \in \mathbb{R} \equiv \text{Real, continuous} \\ \qquad \qquad \text{measure.} \end{cases}$ $\{1, 2, \ldots k\}$ $\left(\begin{array}{l} \text{say, professional,} \\ \text{personal,} \\ \text{etc} \end{array}\right)$

$K$ - class.

It is called as "Regression" $\qquad$ classification $\longrightarrow$ Can be a binary one or multiclass.

representation of an email $\in \mathbb{R}$

↱ spam or not.

Given, $(x_i, y_i)$, $i = 1, \ldots N$.

↓

examples for learning/training $\quad \} \equiv$ Spam filter to be built with "n" examples

[where $\underline{x_i \in \mathbb{R}^d}$, $y_i \in \{0,1\}$]

Find a function, $f: x \to y$

↓

So when given a new email $\equiv f(x) \nearrow^0_{\searrow 1}$ [you should be able to predict whether spam or not]

$x$

$X_1 \qquad X_2 \qquad X_3$

$\begin{bmatrix} \alpha_1 \\ \beta_1 \end{bmatrix}$ $\begin{bmatrix} \alpha_2 \\ \beta_2 \end{bmatrix}$ $\begin{bmatrix} \alpha_3 \\ \beta_3 \end{bmatrix} \to$ 2 features with different values $\quad \}$ can be represented on a plane.

$0 \qquad 1 \qquad 0$ (spam / not spam)

$\longrightarrow f() \equiv f(x) = w^T x + b$

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

$f(x) = w_1 x^1 + w_2 x^2 + b$

Now the problem is reduced to "finding the value of $w_1, w_2, b$".

→ Find $w_1, w_2, b$ ?

↓

OPTIMIZATION problem = Find $w$ that best solves the problem.
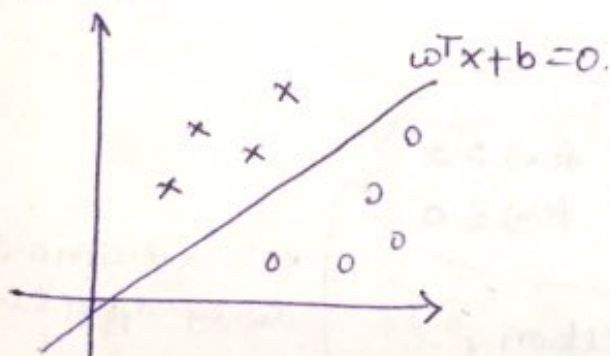
↓

we need an "objective function"

Data $(x_i, y_i)$ pairs

↓

Problem ≡ Find best $w$ ⟶ which minimizes

$$f(x) = w^T x + b.$$

↓

Objective function ≡ Error/Loss

⟶ Regression

⟶ Classification.

changing the loss function can also changes the optimization Algorithm.

↑ because $N \gg$

Representation ≡



$w^T x + b = 0.$

$x$ — Spam
$o$ — Not Spam

**Case ①**

$f(x) = w^T x + b = 0.$

Spam if $f(x) < 0$

Not Spam if $f(x) \geqslant 0$

↓

Discriminant function

**Case ②:**

$f_1(x) = w_1^T x + b_1$

$f_2(x) = w_2^T x + b_2$

Spam if $f_1(x) > f_2(x)$

Not Spam if $f_1(x) \leq f_2(x)$

$f_1 - f_2 > 0.$
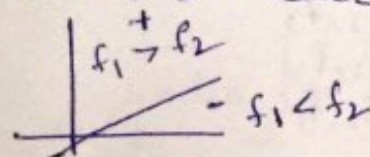
$f_1 - f_2 = (w_1 - w_2)^T x + b_1 - b_2$

$= w_3^T x + b_3.$

In the 1$^{st}$ case ≡

$f(x) \equiv$ function discriminates spam from not spam



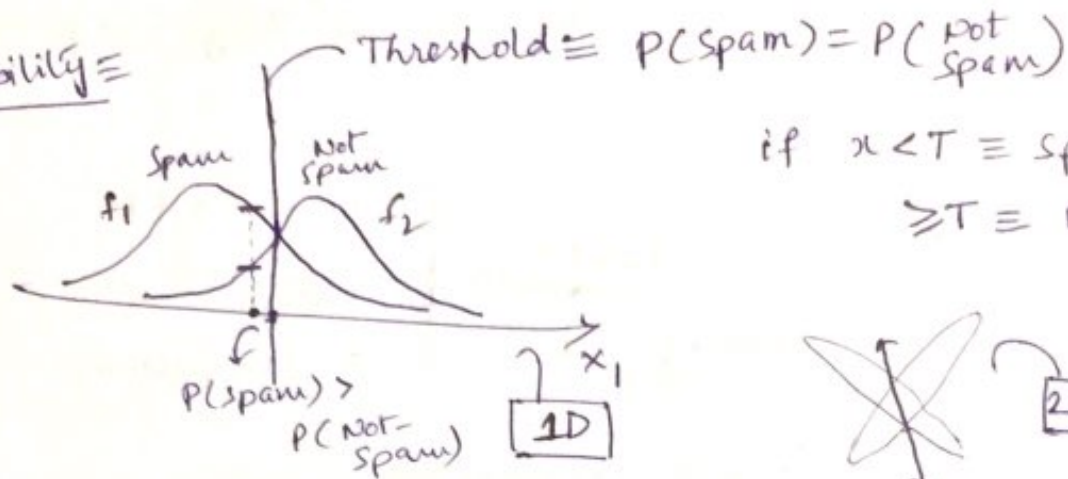In the 2$^{nd}$ case ≡



$f_1 > f_2$

$f_1 < f_2$

? (when to use this approach)

the Second case: Assume $f_1(x) \equiv$ Probability of being spam.

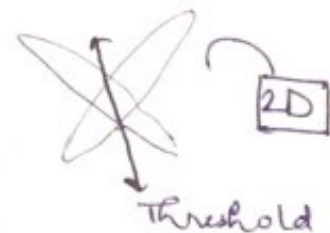$f_2(x) \equiv$ Probability of being not spam

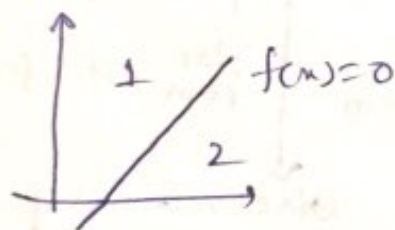both the viewpoints are equivalent in this case.

Probability $\equiv$

Threshold $\equiv P(\text{spam}) = P\left(\begin{smallmatrix}\text{Not}\\\text{Spam}\end{smallmatrix}\right)$



Spam | Not Spam

$f_1$ | $f_2$

$P(\text{spam}) >$
$P(\text{Not-}$
$\text{spam})$

$\boxed{\text{1D}}$

$x_1$

$\boxed{\text{2D}}$

Threshold

So $\equiv$ Spam

if $x < T \equiv$ spam

$\geq T \equiv$ Not Spam.

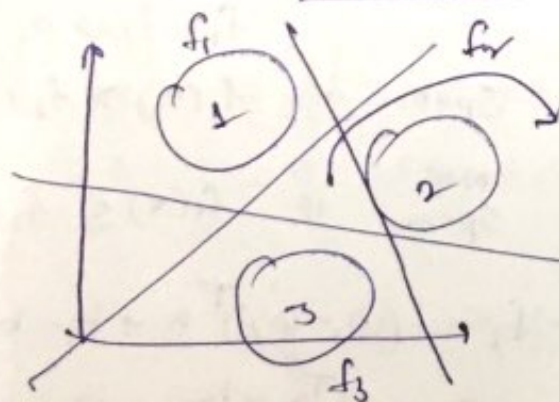These both viewpoints go in parallel. $\checkmark$ (Probabilistic and Non Probabilistic Approach)

2-Class Classification Problem:



1 | $f(x) = 0$
2

1 if $f(x) > 0$
2 if $f(x) \leq 0$

this decision rule doesn't apply here.

what about a 3-class classification?



$f_1$ | $f_2$

1
2
3

$f_3$

what about this sample space.

"No-man"

$\rightarrow$ "multi-class classification problem"

Decision Rule:
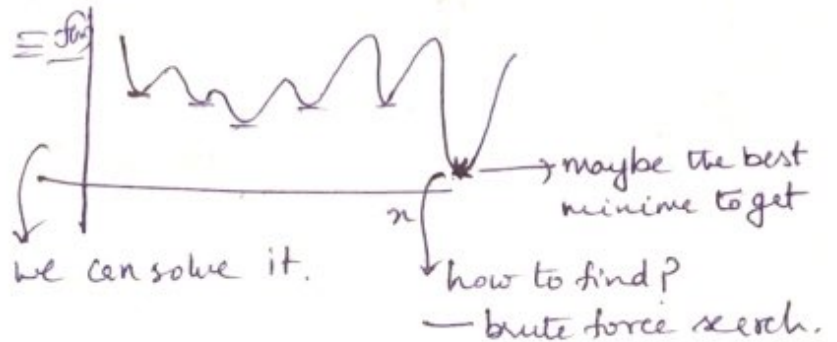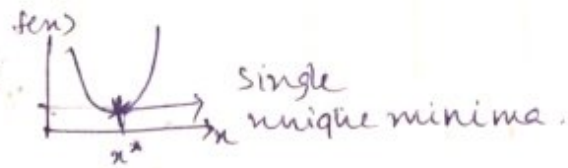
1: $f_1 > f_2, f_3$
2: $f_2 > f_1, f_3$
3: $f_3 > f_2, f_1$

Output — Discrete (Not Continuous)
          ↓ discrete
becomes an optimization
              problem.

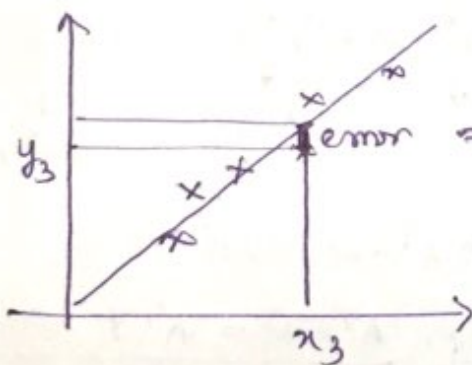Loss function ⌐
          ↓ Optimization
             Problem

Optimization
          ↗ Convex ≡ [f(x) graph] Single unique minima. $x^*$

          ↘ Non - convex

≡ f(x) [graph]

→ maybe the best minima to get

we can solve it.

how to find?
— brute force search.

## Regression ≡

$(x_1, y_1) \cdots (x_n, y_n)$    where $y_i \in R$.

$$f(x) = w^T x + b$$

[graph with $y_3$, $x_3$, error]

error ⟹ $\varepsilon_3 = w^T x_3 + b - y_3$

Error ≡ $\Sigma = \sum\limits_{i=1}^{N} \left[ y_i - (w^T x_i^\circ + b) \right]^2$

‖                predicted value by
min Σ              the function.
w, b

⇓

$\min\limits_{w} \left( \sum\limits_{i=1}^{N} (y_i - w^T x_i)^2 \right)$

can be +ve / -ve
[above or below the] line

$w^T x$
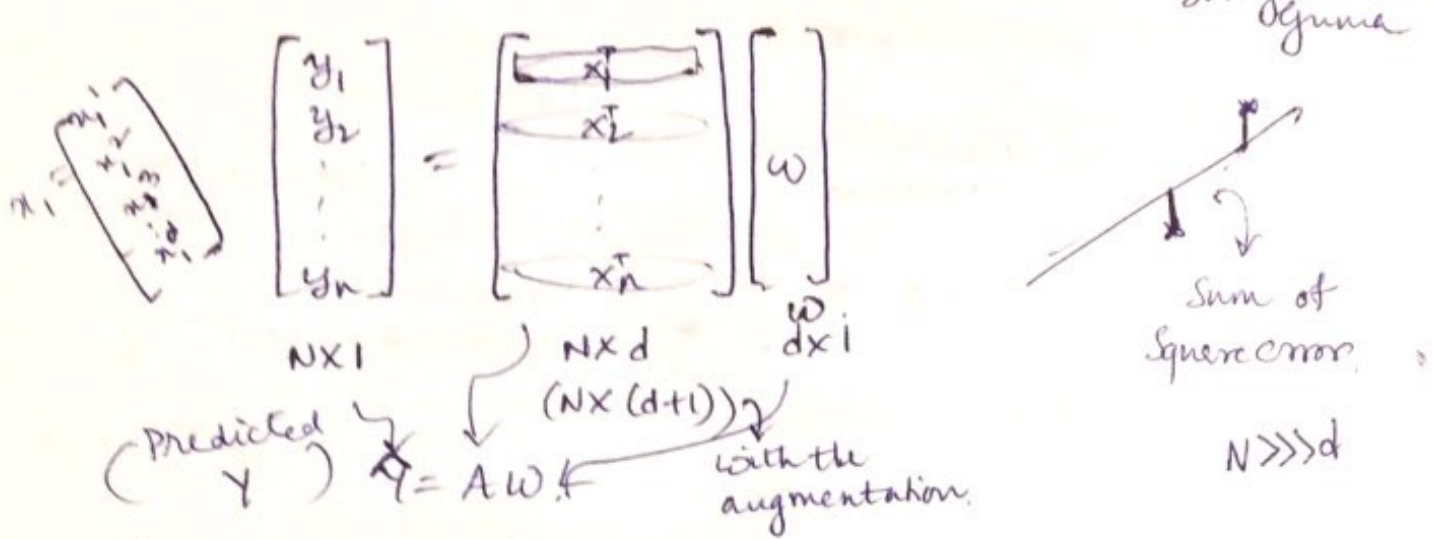
$x$ is augmented with 1.

$w$    $x$
$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + b$
↓
Crow to remove b
⇓

$w' \| x'$
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ b \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

⇓

Minimum Square Error (MSE)

$$x_i = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \overline{\quad x_1^T \quad} \\ \overline{\quad x_2^T \quad} \\ \vdots \\ \overline{\quad x_n^T \quad} \end{bmatrix} \begin{bmatrix} w \end{bmatrix}$$

$N \times 1$  $\qquad N \times d$  $\qquad d \times i$

(Predicted $Y$) $\hat{Y} = AW$ $\leftarrow$ $(N \times (d+1))$ with the augmentation.



Sum of Square error.

$N \ggg d$

$$\sum_{i=1}^{N} (y_i^\circ - w^T x_i^\circ)^2 \equiv \begin{pmatrix} \text{in} \\ \text{matrix} \\ \text{form} \end{pmatrix} \equiv [Y - AW]^T [Y - AW]$$

$(N \times 1)^T$  $(N \times 1)$

$= (1 \times 1)$

Sigma is gone

To minimize this $\Longleftarrow Y^T Y + (AW)^T AW - 2(AW)^T Y$
equation

$\begin{bmatrix} \text{Differentiate} \\ \text{and equate} = 0 \end{bmatrix} \Longrightarrow \dfrac{\partial}{\partial w} \left[ Y^T Y + (AW)^T AW - 2(AW)^T Y \right]$

$= \dfrac{\partial}{\partial w} \left[ Y^T Y + w^T A^T A w - 2 w^T A^T Y \right] = 0$

$\qquad \qquad \underset{0}{\downarrow} \qquad \underset{2A^TAW}{\downarrow} \qquad \underset{2A^TY}{\downarrow}$
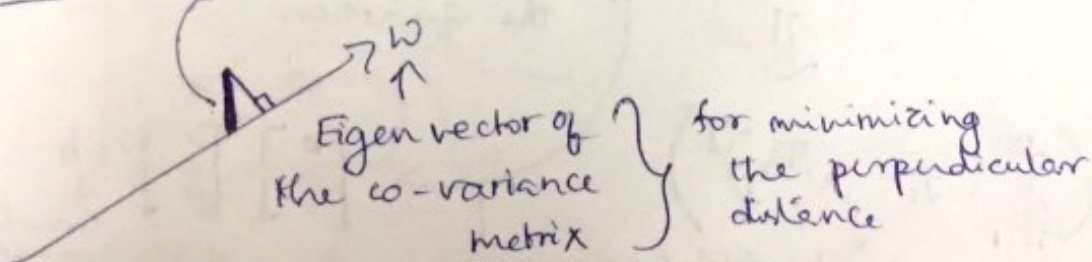
$2 A^T A W - 2 A^T Y = 0$

$A^T A W = A^T Y$

$\Longrightarrow \boxed{W = (A^T A)^{-1} A^T Y}$

Tom Minka ≡ Matrix Differentiation.

Issues with

**Error:**

we took this as error



Eigen vector of the co-variance matrix $\Big\}$ for minimizing the perpendicular distance

> A Size $\equiv N \times d$ $\longrightarrow$ could be extremely large $\equiv$ "Gradient Descent"

Examples of Regression ≡



rain ↑ ... → wind

$y = ax + b$

$\Rightarrow x = \begin{bmatrix} x \\ 1 \end{bmatrix}$ (augmented vector)

$\begin{bmatrix} a \\ b \end{bmatrix}^T \begin{bmatrix} x \\ 1 \end{bmatrix}$

$y = w_1 x_1 + w_2 x_2 + w_3$

$\Rightarrow x = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$   $\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$ → this even works for non-linear

Representation of $x$ ∈

$Y = f(x) = w^T z \cap \begin{bmatrix} x \\ 1 \end{bmatrix}$

$\begin{bmatrix} x^2 \\ x \\ 1 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$

— not restrictive  $w_1 x^2 + w_2 x + w_3 = 0$
— due to augmentation (not passing thru origin ✓)
— New vectors ≡ Non linear functions embedded in $w^T x = 0$.

↓ line fitting

Can be represented as ≡

$x = \begin{bmatrix} x^2 \\ x \\ 1 \end{bmatrix}$   $w = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$

$\hookrightarrow w^T x = 0$.

$ax + b = 0$
$ax^2 + bx + c = 0$.
Nonlinear ≡ reality.

* Main goal:- Find "$\underline{w}$", for a function of the form ≡ $\underline{w^T x}$

Not restrictive

$\Rightarrow Y = A w \Rightarrow$ To find $w ≡$
$\quad | \quad | dx1$
$N×1 \quad N×d$

$w = A^{-1} Y$ (doesn't work.
as long as A is not singular and square matrix (N=d).

→ $x → z$ can be mapped

$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \Rightarrow \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \\ x_1 x_3 \\ \vdots \end{bmatrix} z$   these are not independent.

$w = \boxed{(A^T A)^{-1} A^T} Y$
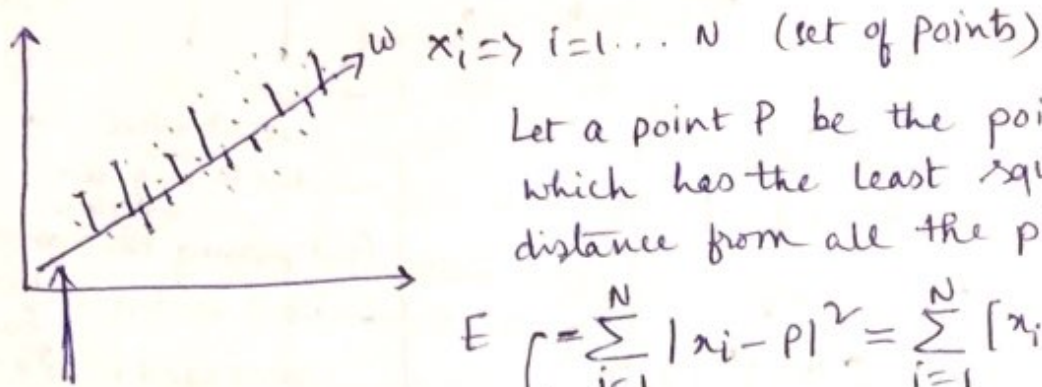
(pseudo inverse of A.)

↓

approach ≡ Minimum Square Error

→ Doesn't minimize perpendicular (orthogonal) dist.
* $(A^T A)^{-1} ≡$ computational storage
— Data comes incrementally (data is always not offline)

$\omega$ — needs to be incrementally updated.

$\Rightarrow$ **Orthogonal Distance Minimizing** $\quad \to$ Find $P$, a point
$\quad \to$ Find $\omega$, the direction.



$x_i \Rightarrow i = 1 \ldots N$ (set of points)

Let a point $P$ be the point ~~such~~ which has the least square error distance from all the points.

$$E = \sum_{i=1}^{N} |x_i - P|^2 = \sum_{i=1}^{N} [x_i - P]^T [x_i - P]$$

eigen vector for minimum, $\dfrac{\partial E}{\partial P} = 0$.
from the
covariance
matrix.
$\begin{pmatrix}\text{passing thru the}\\ \text{mean}\end{pmatrix}$

$$\frac{\partial}{\partial P} \sum_{i=1}^{N} \left( \underset{0}{\underbrace{x^T x}} + \underset{2P}{\underbrace{P^T P}} - \underset{2x_i}{\underbrace{2P^T x_i}} \right)$$

$$= \sum_{i=1}^{N} 2P - 2x_1 = 0.$$



Normalized.

$$\boxed{\omega^T \omega = 1}$$

Norm of $x$'s are fixed.

$$\sum_{i=1}^{N} \left( \|x_i\|^2 - (\omega^T x_i)^2 \right).$$

independent of $\omega$.

$$\sum_{i=1}^{N} P = \sum_{i=1}^{N} x_i$$

$$P \cdot N = \sum_{i=1}^{N} x_i$$

$$P = \frac{\sum_{i=1}^{N} x_i}{N}.$$

Minimize $\equiv \sum_{i=1}^{N} - (\omega^T x_i)^2$ s.t $\omega^T \omega = 1$

$\downarrow$

Maximize $\equiv \sum_{i=1}^{N} (\omega^T x_i)^2$

$(x\omega)^T x\omega$

$\omega^T x^T x \omega$

$\boxed{\begin{array}{l}(\omega^T x)^T (\omega^T x) \\ x^T \omega \omega^T x \quad ?\end{array}}$ wrong.

$\underleftarrow{\quad}$ look at regression.

Max: $\quad \omega^T x^T x \omega$.

Maximize $\equiv$ $\boxed{w^T x^T x w}$    st    $w^T w = 1$

$\downarrow$ w is eigen vector of $\boxed{x^T x}$ $\equiv$ PCA
(largest) eigen value

$\downarrow$ covariance matrix

$Ax = \lambda x$

Max ( $\cancel{\ast}$ $w^T x^T x w$ ) $\left. \substack{\lambda w \\ \\} \right\}$ $\downarrow$ To Maximize this quantity $w^T \lambda w$, we have to take largest value of $\lambda$

$w^T \lambda w$.

$\underline{\hspace{2cm}}$ largest eigenvalue

man ( $w^T A w$ )
‖
$w^T \lambda w$ $\equiv$ $\lambda w^T w$ = $\boxed{\lambda}$
$\downarrow$ eigen value.     $\downarrow$ 1.

$\downarrow$

Mainly used in PCA.

$\underline{\hspace{2cm}}$ Direction is given by w, where w is the eigen vector corresponding to the largest value of $x^T x$.

## Gradient Descent $\Rightarrow$

$(x_i, y_i)$, $i = 1 \dots N$.

$\sum (y_i - f(x_i))$

loss function = need to be minimized



Loss function.

negative gradient

$\rightarrow$ * Randomly initialize $\underline{w}$

gradient=0

* Move on the negative grad

$w^{n+1} \leftarrow w^n$
$-\eta \nabla J$

$\{$ learning rate $\}$

* Repeat steps until you reach turning$^n$.

* Goal
  — objective
* Data (x, y)
* Optimization
  — closed form Algo
  — Iterative Algo $\longrightarrow$ gradient Descent

Given $(x_i, y_i)$, $i = 1, \dots n$.
Find $y = f(x)$ such that it fits the existing (and future) data
or minimize a loss/error function.

$\rightarrow y_i \sim f(x_i)$ $\forall_i^o$ $\qquad f(x) = w^T x$ } Assumed it to be a linear function.

$\swarrow \searrow$

$R \qquad \{0,1\}$
$\qquad C$

$$\begin{bmatrix} x^1 \\ x^2 \\ x^3 \\ x^4 \end{bmatrix} = w_1 x^1 + w_2 x^2 + w_3 x^3 + w_4 x^4$$

$w^T x \equiv$ Not restrictive

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \searrow \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 \\ x_2 \\ 1 \end{bmatrix} \rightarrow \text{Find the } w \text{ suiting these dimensions/ features.}$$

$2d$

$5d \longrightarrow$ In five dimension, it is linear.

## Gradient Descent:

> Make a random guess.
$\downarrow$
you should either
increase w/decrease w
to reach optimal $\omega$.

> check the gradient
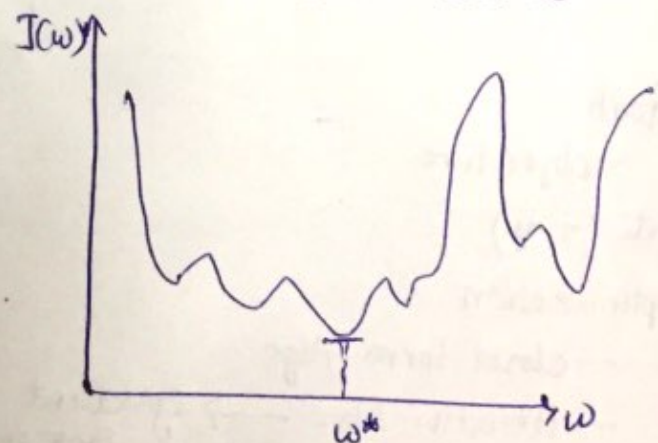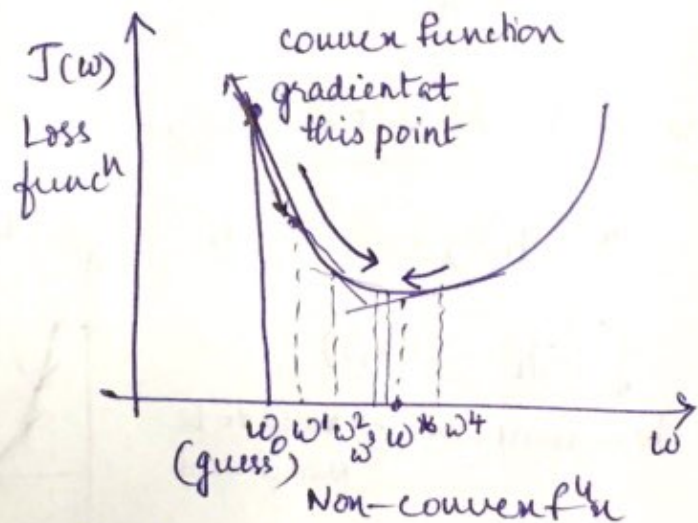descent.
and calculate the next
move $\equiv$



$J(\omega)$
Loss func$^n$

convex function
gradient at
this point

$\omega_0 \ w^1 w^2_3 \ w^* w^4$
(guess) $\qquad \qquad \omega$
Non-convex fn

$J(\omega)$

$\omega = \omega^o - \eta \nabla J(\omega^o)$
$\qquad \underbrace{\qquad}$
amount of movement.

At the minimum;
$\Delta J(\omega) = 0$ and
$\omega,$ doesnt change $\rightarrow$ Algorithm converges

$\omega^*$ $\qquad \omega$

may not reach $\nabla J(w^0) = 0$, we can put an error constant threshold.

In nonconvex function; the minimum (local) $\nabla J(w_0) < 0.001$ that we get maynot be the best one. by using the same equation

changing it doesnt give much of an advantage

What is $\eta$ (eta)? —— Adaptive over time / iterations based on the behavior of algo.

Learning Rate

$w^{n+1} \leftarrow w^n - \textcircled{$\eta$} \nabla J(w^n)$

→ Given is data $(x_i, y_i)$ not the loss function.

when we make a move into the negative gradient

— Assume the function to be a line at the initial stage. and compute the slope.



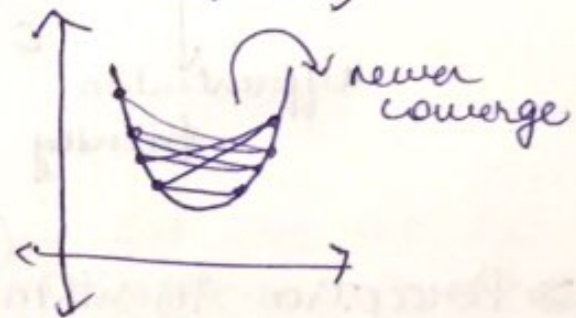This is why, we take $\eta$ and go step by step.

←— (if function nature is known, it would have been easy).

⇒ Start with a very small $\eta$

Reason : (with large $\eta$, you take huge leaps ( to and fro into positive and negative gradient regions))

⇒ $\eta$ is increased, if direction of $w$ in the prev step and present step are same.

⇒ and if direction changes between steps, $\eta \equiv$ decreased



newer converge.

for Ex:  $\eta \leftarrow 1.5 \times \eta$, if prev dr = present dr

$\eta \leftarrow 0.5 \times \eta$, if direcn change.

Taylor series expansion ≡          $x = w^{n+1}$
          $a = w^n$

$$f(x) = f(a) + f'(a) \cdot (x-a) + f''(a)(x-a)^2 + \cdots$$

$$f(w^{n+1}) = f(w^n) + [w^{n+1} - w^n] \cdot f'(w^n)$$

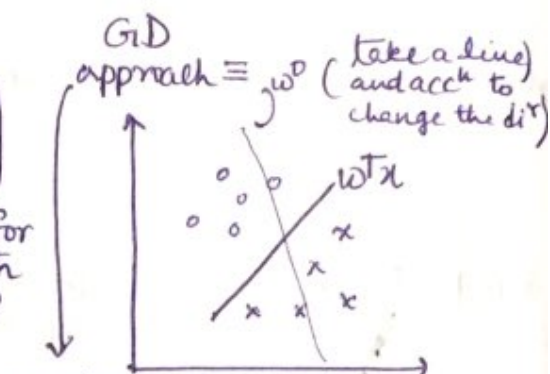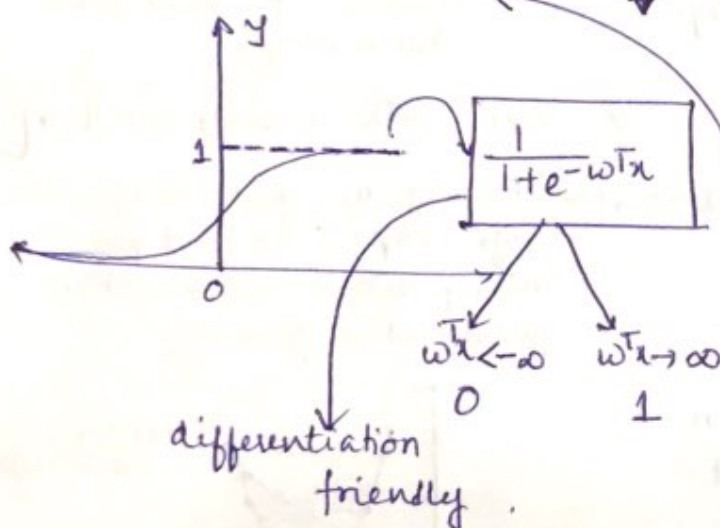negative quantity.   $w^{n+1} \leftarrow w^n - \eta \Delta J(w^n)$

→ **GD for Classification:**   New function value = Old function value − (Positive quantity).
(less than the old one).

⟫ **Logistic Regression** ≡

⟫ **Perceptron Algorithm** ≡

(instead of step) → search for smooth step

GD approach ≡ $w^0$ (take a line and acc$^n$ to change the dir$^n$)

$w^T x$

minimum error. ≡   not differentiatio friendly to eliminate the errors

Linear classifier.
$$f(x) = 1, \ w^T x > 0$$
$$0, \ w^T x \le 0.$$

But instead;

1 ─────

$\dfrac{1}{1 + e^{-w^T x}}$

$w^T x \to -\infty$   $w^T x \to \infty$
    0              1

differentiation friendly.

(created a new objective func$^n$)

⟫ **Perceptron Algorithm** ≡

⟫    $J = \sum w^T x$
        $x \in$ Misclassified Samples

} Minimize the no. of such (error) classifications.

Sum over misclassified samples. (instead of over all the samples)

when, Misclassified samples = $\{\phi\}$ ↝ we found the "line".

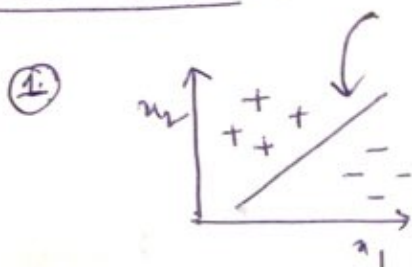Logistic Regression $\equiv$ step $\rightarrow$ smooth step

Perceptron Algorithm $\equiv$ creating a new objective function.

$\downarrow$

every sample theta gets misclassified is added to it

$w^{n+1} \leftarrow w^n \underbrace{x}_{u}$ ?

If the samples are "linearly seperable", $\exists \, w$ which (can seperate the and $-ve$ examples) perceptron will find that $\underline{\underline{w}}$

Classification $\equiv$     Model of classifier is [Given this model, we $\underline{\underline{cannot}}$ generate an apple]
                   $\equiv$ Line.
                                 we model the line in the feature space

①    Boundary/line between class A and class B.

$\rightarrow$ Discriminative classifier/model

② Generative Model ( we donot model the boundary ) $\equiv$ reduced representation

       $\downarrow$                       instead model the data [ as in "a point" ]

for Ex : given a fruit, how apple-like or orange-like is it?     [Ex: model oranges and apples independently.]

$\downarrow$

Model the classes?

Collection of objects modelled by probability density.     [if I see 10.5, how likely is to classify it as apple]

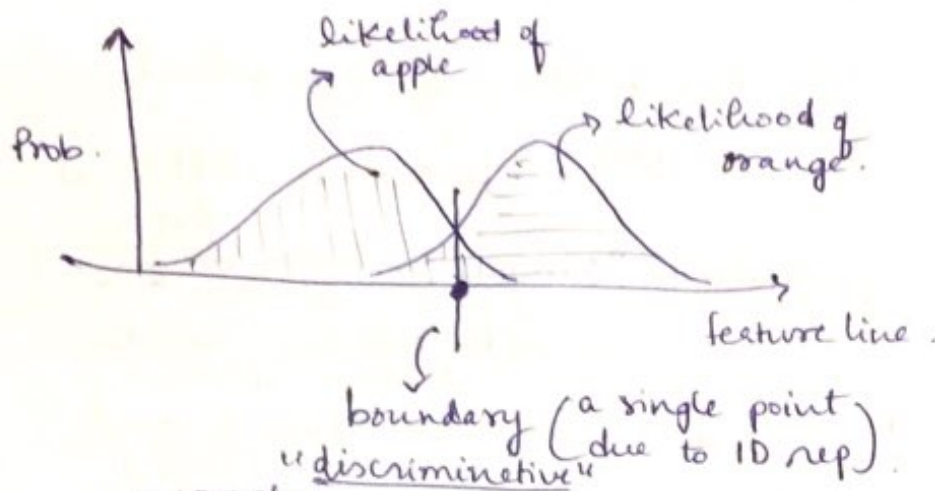these are apples          there are oranges.

likelihood functions $\equiv$ Independently Learning about apples/oranges.

$\Rightarrow$ given density function $\} \rightarrow$ we can generate new points    [that is why $\equiv$ generative model]

Discriminative  
Classifier $\Big\}$ ≡ $\begin{cases} \text{we cannot generate samples as it can} \\ \text{lie anywhere over the apple side of} \\ \qquad\qquad\qquad\text{the line.} \end{cases}$

↓

Not rich enough to capture the details

Cannot explicitly model apples.

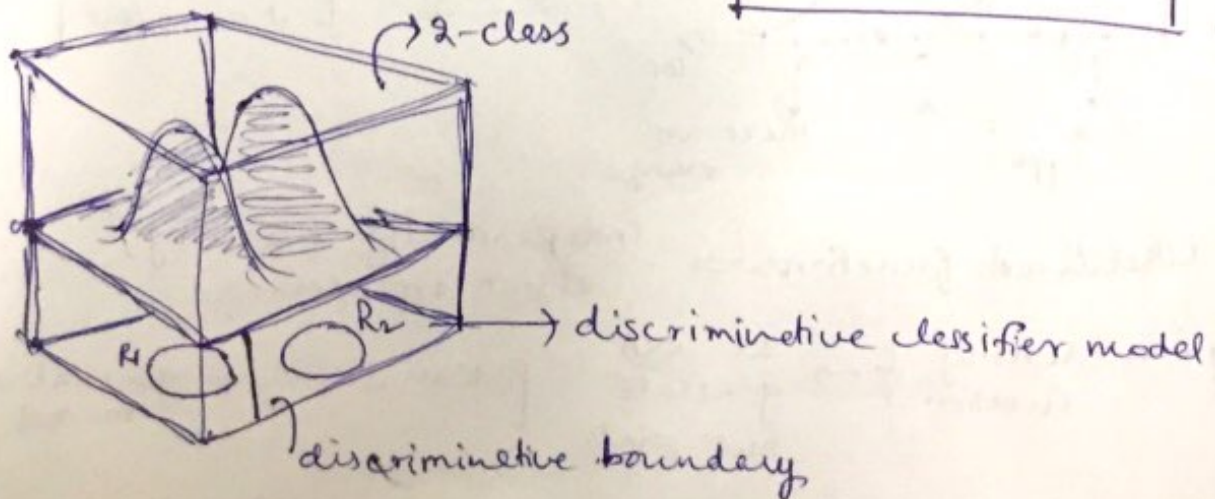[can only say that anything that lies on one side is apple]



likelihood of apple

Prob.

likelihood of orange.

feature line.

boundary $\begin{pmatrix} \text{a single point} \\ \text{due to 1D rep} \end{pmatrix}$.

"discriminative"

$\lceil$ ω ≡ class

$P\left(\dfrac{\omega_1}{x}\right)$ ≡ "Probability" that given feature vector, $x$ belongs to the $\omega_1$ class. $\qquad [0 < P < 1]$

↓

feature vector

$P\left(\dfrac{\omega_1}{x}\right)$, $P\left(\dfrac{\omega_2}{x}\right)$. can be compared to determine the class.

$\Big[$in the case of 2-class problems $\Big\} ≡ P\left(\dfrac{\omega_1}{x}\right) + P\left(\dfrac{\omega_2}{x}\right) = 1\Big\}$.

Whereas, if there are $c$ classes ≡ $\boxed{\displaystyle\sum_{i=1}^{c} P\left(\dfrac{\omega_i}{x}\right) = 1}$



→ 2-class

→ discriminative classifier model

discriminative boundary

$P$ [Probability] ;    $P(^{\omega_1}/_x)$, $P(^{\omega_2}/_x)$.

$p$ [pdf]

$P(\omega_1, x) = P(^{\omega_1}/_x) \cdot P(x)$  (or)

$\uparrow$
Joint Probability            $P(^x/_{\omega_1}) \cdot P(\omega_1)$.

$$\therefore P\left(^{\omega_1}/_x\right) = \frac{p(^x/_{\omega_1}) \cdot P(\omega_1)}{p(x)}$$

Posterior prob.

$\overline{P(A,b) = P(^A/_B) \cdot P(b)}$

(independent events)
$= P(A) \cdot P(B)$

$\downarrow$
Bayes Theorem.

$\Rightarrow P(\omega_1) \equiv$ Prob. of a class irrespective of given sample.

$\downarrow$ Prior probability / belief

$\Rightarrow P(^x/_{\omega_i}) \Longrightarrow$ likelihood.

$\Rightarrow p(x) \longrightarrow$ how likely is that we'll observe $x$.
"evidence".

Posterior prob $\equiv$ Post observing $x$, what is the probability
$P(^{\omega_1}/_x)$

Ex: Disease and Test

$P(D) = 0.001$

$P(+/_D) = 0.99$  ($\sim$ If you have the disease, tests will
say + ,99%  (Accuracy of test)

$P(D/+) = \dfrac{0.99 \times 0.001}{P(+)}$

Prob of having disease, given
test came out +ve

$\rightsquigarrow$ Evidence = likelihood that the
test come out +ve.
$\Big($ irrespective of
disease or not$\Big)$
on random person

$P(+) = P(+/_D) \cdot P(D) + P(+/_{ND}) \cdot P(ND)$
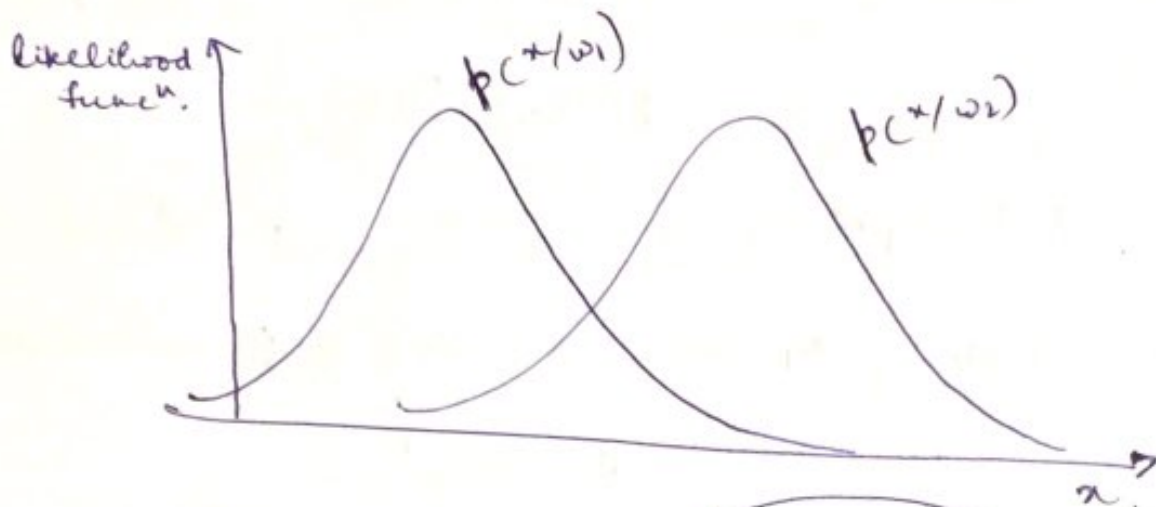$($                        $\uparrow$
person with D        person with ND
gets +ve result

$= 0.99 \times 0.001 + 0.01 \times 0.999$
$= 0.01 (0.099 + 0.999) = 1.098 \times 0.01$

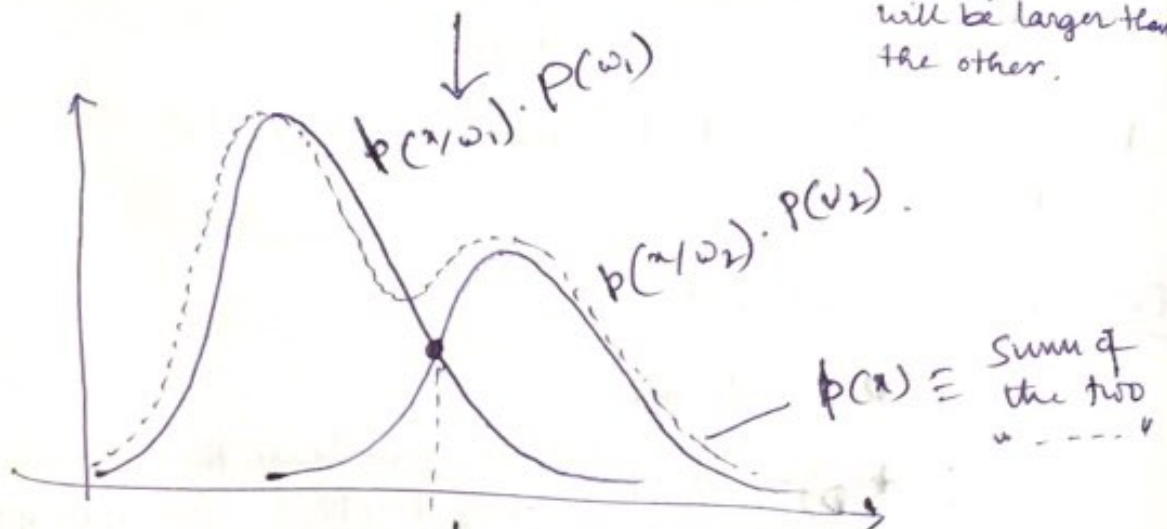$$P(^D/_+) = \frac{0.99 \times 0.001}{0.01098} = 0.09 \ (9\%).$$
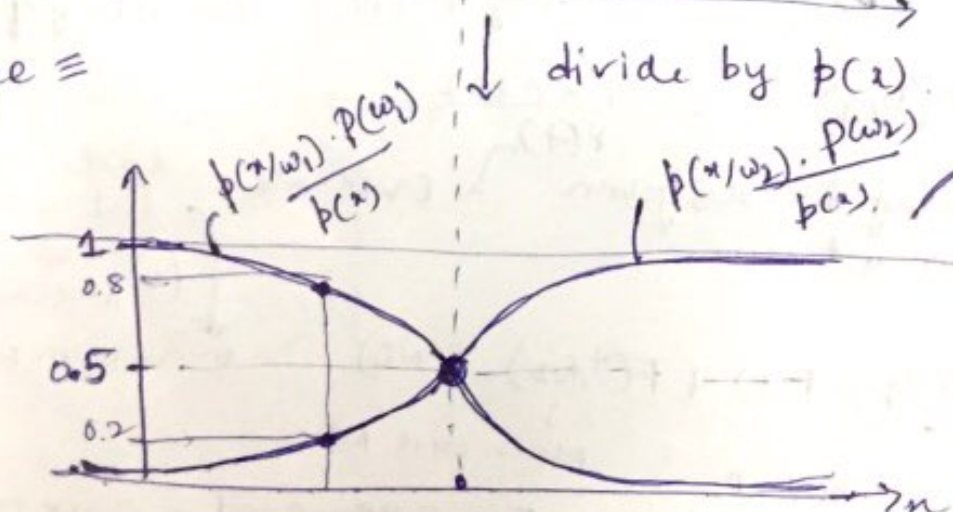
$\uparrow$ chance of having the disease.

likelihood funcⁿ.

$p(^x/\omega_1)$    $p(^x/\omega_2)$

$x$.

$$P(^{\omega_1}/_x) = \boxed{\frac{p(^x/\omega_1) \cdot P(\omega_1)}{p(x)}} \leftarrow \text{constant}$$

$\searrow$ one of the curves will be larger than the other.

$p(^x/\omega_1) \cdot P(\omega_1)$

$p(^x/\omega_2) \cdot P(\omega_2)$.

$p(x) \equiv$ Sum of the two "-----"

Normalize $\equiv$     $\downarrow$ divide by $p(x)$.

$\frac{p(^x/\omega_1) \cdot P(\omega_1)}{p(x)}$     $\frac{p(^x/\omega_2) \cdot P(\omega_2)}{p(x)}$

add upto 1 always for any $x$

1
0.8
0.5
0.2

$n$

# Bayes Theorem ≡

$$P(\omega_j / x) = \frac{P(x/\omega_j) \cdot P(\omega_j)}{P(x)}$$

↑
after observing $x$,
state of nature ≡ $\omega_j$

probability that state of nature, $\omega_j$ is true.

↗ Pick $\omega_j$ that has maximum prob.

Are all errors equally costly? __No__

so we define a loss function ≡ $\lambda(\alpha_i / \omega_j)$

↳ action, $\alpha_i$

↳ state of nature ≡ class

Loss is not symmetrical,
$\xi_1$ if class A is misclassified as class B or
$\xi_2$ if class B is misclassified as class A.
$$\xi_1 \neq \xi_2$$

→ there can be many actions

Risk, $R(\alpha_i / x)$ → observe $x$, risk of taking $\alpha_i$ action.

→ Prob that $\omega_j$ is true

(of taking action $\alpha_i$) $$R(\alpha_i / x) = \sum_{j=1}^{C} \lambda(\alpha_i / \omega_j) \cdot P(\omega_j / x)$$

↳ loss incurred because of taking $\alpha_i$ in $\omega_j$ state of nature

↓

c classes.

Strategy to pick the best action = based on action plan

Evaluate the action plan?

$\alpha(x)$ — if i see this observa'n, i'll do this.

↓ given an $x$, what action needs to be taken.

$$R(\alpha(x) / x) \cdot p(x)$$

↗ prob of observing $x$

↑
Risk of following action plan, $\alpha(x)$ on observing $x$.

" Action plans can be eval. based on risks"

— Integrate over all possible $x$,

Overall Risk ≡ $$R_{\alpha(x)} = \int_x R(\alpha(x)/x) \cdot p(x) \, dx.$$

To Minimize overall risk, by minimize $R(\alpha(x)/x)$ for each $x$.

choosing appropriate/
best "$\alpha(x)$"

From Bayes —
$$R(\alpha_i/x) = \sum_{j=1}^{c} \lambda(\alpha_i/w_j) \cdot P(w_j/x)$$

$R^* = $ Bayes Risk.

Bayes Action Plan.

$$R(\alpha(x)/x) = \underset{\alpha_i}{argmin}\left(R(\alpha_i/x)\right)$$

Risk $\propto P(error)$.

**✳**

$P(w_j) \rightarrow$ can change
(current assumptions)

change in $\underset{Risk}{=}$ Linear

vs
ground truth



$P(error)$

$P(w_j)$ 0

( Mix of the     Mix of
  classes $\equiv$ Ex: Kashmir
  and Washington
  apples
  in the market )

"Neyman-Pearson
criterion"

if the market
fluctuates, the
maximum risk
is minimized
choose such
action plan.

$$P(w_i/x) = \frac{p(x/w_j) \cdot P(w_j)}{p(x)} \longrightarrow \text{need not be computed to decide/compare } \begin{bmatrix} \text{independent} \\ \text{of } w_j \end{bmatrix}$$

$$g_i(x) = p(x/w_i) \cdot P(w_i)$$

(this is different) $g_i(x) = \ln p(x/w_i) + \ln P(w_i)$
$g$

To classify, you need to compute $P(w_i/x)$, you can
either compute $g_i(x)$ or $g_i(x)$ (latter one)

Equation of Normal density $\Rightarrow p(x) = \dfrac{1}{(2\pi)^{d/2}|\Sigma_i|^{1/2}} e^{-\frac{1}{2}(x-\mu_i)^t \Sigma_i^{-1}(x-\mu_i)} = P(x/\omega_i)$

$\left(\begin{array}{c}\text{can be}\\ \text{ignored}\end{array}\right)$ constant

$g_i(x) = -\dfrac{d}{2}\ln 2\pi - \dfrac{1}{2}\ln|\Sigma_i| - \dfrac{1}{2}(x-\mu_i)^t \Sigma_i^{-1}(x-\mu_i) + \ln P(\omega_i)$

$g_1(x), g_2(x)$ can be computed for the normal densities

$*$ $\quad g_i(x) = -\dfrac{1}{2}\ln|\Sigma_i| - \dfrac{1}{2}(x-\mu_i)^t \Sigma_i^{-1}(x-\mu_i) + \ln P(\omega_i)$

**Case 2:** if $\Sigma_i = \Sigma$ , if Kashmir and Washington, variance (both the classes) may be same, same covariance matrix

$g_i(x) = -\dfrac{1}{2}(x-\mu_i)^t \Sigma_i^{-1}(x-\mu_i) + \ln P(\omega_i)$

$= -\dfrac{1}{2}(\underbrace{x^t\Sigma^{-1}x} - 2\Sigma^{-1}\mu_i^t x + \mu_i^t \Sigma^{-1}\mu_i) + \ln(P(\omega_i))$

$\underset{\substack{\text{independent}\\ \text{of i, class}}}{\uparrow}$

same stretch and tilt but just shifted

linear boundary

$\downarrow$ tends to be "linear"

$g_i(x) = \Sigma^{-1}\mu_i^t x - \dfrac{1}{2}\mu_i^t \Sigma^{-1}\mu_i + \ln(P(\omega_i))$

Discriminant function.

$= w_{i1}^T x + w_{i0}\quad \}$ linear function.

In 2-D, decision boundary, would be a linear function (a line)

In 3-D, it'll be a plane

$*$ $\Sigma_i \neq \Sigma$

$\Rightarrow$ if covariance matrix is not the same, then the decision boundary will no more be a linear function but a complex function. (it could be an ellipse, hyperbola ...)

**Case 3:** $\Sigma_i = \Sigma_i$

$g_i(x) = -\dfrac{1}{2}(x^t\Sigma^{-1}x - 2\Sigma^{-1}\mu_i^t x + \mu_i^t \Sigma^{-1}\mu_i) + \ln P(\omega_i)$

## Loss :

In Regression — $R \equiv J \equiv \sum\limits_{i=1}^{N} (y_i - w^T x_i)^2$

In Classification — $C \equiv 0\%$ misclassifications.

$g(z) = +1 \; ; \; z \geqslant 0$
$\quad\quad -1 \; ; \; z < 0$

Ex:    $y_i * g(w^T x)$.


Classification

1. $y_i = +1$  ✓ correct $= 1 = y_i (g(w^T x))$
   $w^T x = 10$.

2. $y_i = -1$  ✓ correct $= 1 = y_i (g(w^T x))$
   $w^T x = -10$.

3. $y_i = -1$  ✗ incorrect
   $w^T x = 10$    (misclassified).

for this, instead of using



its better to use



$* \quad C \equiv J = \dfrac{1}{2N} \sum\limits_{i=1}^{N} (1 - y_i \cdot g(w^T x_i))$

— Good
— Not optimisation friendly.

this is
( Not Logistic
  Regression )

optimize it by
changing the 〔g〕 ⟶


[Logistic Regression]

(or)



## Gradient Descent:

1. start with $w^0$, $n = 0$.

2. $w^{n+1} \leftarrow w^n - n \nabla J$
   ⤷ changes, when loss function changes
   ⟶ varies over time (iteration) [ usually constant ]

3. $n \leftarrow n + 1$

4. Repeat 2-4 until some convergence criteria is met.
   ?
   — change in $w$ is very very small.
   — change in loss is very small.

→ Loss function is not linear but the discriminant function is linear.

$$f(y) = f(x) + (y-x)f'(x) + \frac{(y-x)^2}{2}f''(x)$$

approximated $\equiv f(Y) = f(x) + (y-x)^T \nabla f(x) + \frac{1}{2}(y-x)^T H (y-x) = \boxed{H \quad \text{Hessian} \quad \text{matrix}}$ — operator, $d \times d$ matrix

$$\boxed{f(w^{n+1}) = f(w^n) + [w^{n+1} - w^n]^T \nabla f(w^n)} + \frac{1}{2}(w^{n+1} - w^n)^T H [w^{n+1} - w^n]$$

↓

$$w^{n+1} \leftarrow w^n - \eta \nabla f$$

$$w^{n+1} - w^n = -\eta \nabla f \qquad \boxed{\nabla f^T \cdot \nabla f}$$

only proves that

$$f(w^{n+1}) = f(w^n) + (w^{n+1} - w^n)^T \nabla f$$

$$f(w^{n+1}) = f(w^n) - \eta (\text{+ve quantity})$$

⟹ GD reduces the loss in each iteration if

$\eta$ is +ve and small

but can we go faster?

$$f(w^{n+1}) = f(w^n) - \eta \|\nabla f\|^2 + \frac{\eta^2}{2} \nabla f^T H \{\nabla f\}$$

Since we want minimum, diff wrt $\eta$

$f(w^{n+1})$

$$\frac{\partial}{\partial \eta} f(w^{n+1}) = 0 - \|\nabla f\|^2 + \frac{2\eta}{2} \cdot \nabla f^T H \{\nabla f\} = 0$$

$$\eta \cdot \nabla f^T H (\nabla f) = \|\nabla f\|^2$$

$$\eta = \frac{\nabla f \cdot \nabla f^T}{\nabla f^T \cdot H\{\nabla f\}} = \frac{\nabla f^T \nabla f}{\nabla f^T H \nabla f}$$

$$w^{n+1} - w^n \equiv S$$

Better update rule than this: $w^{n+1} \leftarrow w^n - \eta \nabla f$ ?

$$f(w^{n+1}) = f(w^n) + S^T \nabla f + \frac{1}{2} S^T H\{S\}$$
$$= f(w^n) + S^T \nabla f + \frac{1}{2} S^T H \cdot S$$

H - matrix
S - vector

$$\frac{\partial}{\partial S} = 0 \Rightarrow \qquad \nabla f + HS = 0$$

$$S = -H^{-1} \nabla f$$
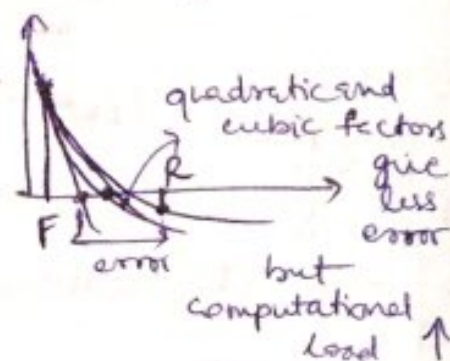
↓

Better update rule, — Newton iteration. (but not used)
Descent

↳ Not used → maybe H might be singular. ✓
(inverse doesn't exist)

if |H| is very less ↘ inverse can be very
sensitive

Not computationally attractive

First of all,
How is this better?
than the first one? GD
→ In the first case, approximating to the first order
[GD] gives a smaller value (higher error)



quadratic and
cubic factors
→ give
less
error
but
computational
load ↑

For now, constant $\eta$. ✓ but we have to
handle the loss, J

optimizing
wrong objective
func$^n$

**Perceptron** ≡  Sum the loss only over the
misclassified samples.
objective function.

$w^T x \geqslant 0, \; y_i = +1$
$w^T x < 0, \; y_i = -1$

$$\sum_{x_i \in \mathcal{E}} -y_i \, w^T x_i$$

↳ misclassified.

Total Error (**Amount** of misclassification)
Not number, because
we are taking $w^T x$.
instead of $g(w^T x)$.

update rule changes to ≡  $w^{n+1} \leftarrow w^n + \eta \sum_{x_i \in \mathcal{E}} y_i \, x_i$

Changes with every step.
so

Not
computationally ← alternative ≡
different

$$w^{n+1} \leftarrow w^n + \eta \sum_{i=1}^{N} (t_i - o_i) \, x_i$$

target   output

if target = 1, output = 1
$t - o = 0$ (correctly classified)    doesn't need to
go over misclassified set

Convergence criteria $\Rightarrow$ until no class is misclassified.

## How does this work? (Intuitive Approach)

① Assume, $t_i = +1$, $o_i = -1$ and $x_i = +ve$.

$\underbrace{\qquad}_{\text{Error is present.}}$

$w^T x \geq 0 \; , \; +1$

$w^T x < 0 \; , \; -1$

initially $\underline{\underline{w < 0}}$ as $x > 0$ & $w^T x < 0$

with these values, $w$'s going to $\underline{increase}$. $\Rightarrow$ increase $w$

$\qquad \qquad \hookrightarrow w$ becomes $+ve$. $\underline{\underline{\phantom{x}}}$

② $\qquad t_i = +1, \; o_i = -1$ and $x_i = -ve$

$w^T x < 0$

$\wedge$

$x < 0 \quad w > 0$

$(t - o)x \Rightarrow -ve \longrightarrow w \; \underline{decreases}.$

$\qquad \qquad \hookrightarrow w$ becomes $-ve$ $\underline{\underline{\phantom{x}}}$

$\text{ll}^{y}$ for diff values of $t$ & $o$.

$\longrightarrow$ If a sample is misclassified, add/subtract it. $\left( \begin{array}{c} \text{based on values} \\ \text{of } t \text{ and } o \end{array} \right)$

---

Binary (2-class) classification :- $\underline{\text{Loss/obj/Error}}$

$\overbrace{\qquad \qquad \qquad \qquad}$

optimised using gradient descent

* Perceptron — $\underset{w}{argmin} \sum\limits_{x \in \mathcal{E}} - y_i \, w^T x . \equiv \begin{array}{c} OBJ \\ Func^n \end{array}$

* Logistic Regression $\qquad \qquad$ (misclassified samples)

$w^{k+1} \xleftarrow{} w^k - \eta \boxed{\nabla J}$

Algo depends on $\nabla J$. $\overset{\curvearrowleft}{J}$

## Algo :-

* Initialize $w^0$, $k = 0$. $\qquad \overset{\frown}{\qquad} \checkmark$ change is done per sample?

$w^{k+1} \xleftarrow{} w^k + \eta \, Y_k X_k \; , \; \text{if } X_k \in \Sigma$

$k \xleftarrow{} k + 1$

$\Big($ Repeat until $\Sigma$ is empty.

$\downarrow$ $\qquad \qquad \text{sol}^n \qquad \left( \begin{array}{c} \text{update for} \\ \text{every sample} \end{array} \right)$

3 variations $= \longrightarrow$ changes with every sample (every sample)

* online (imp: data streaming (continuously flowing samples)).
* Batch $\longrightarrow$ "compute derivative over all the samples"
* Stochastic

Online   (samples one at a time).
— Start with the first sample.
— If error, update. else go to next sample
— Repeat step 2.

⇒ *If Batch goes into oscillations, you never come out of it

Stochastic   (Not deterministic approach)
↓    ↓randomly picked up.
Why? , because there might be patterns leading
to oscillations.

Mini Batch

— Instead of taking all the samples, consider a set of samples!
for En: out of 100 samples, sample a set of 10.

| mini batch stochastic batch gradient |   * mini SGD

→ online is a form of stochastic.

sample has to be replaced back

* Pseudocodes of all the algos   *
      online   miniSGD.
      stochastic

Perceptron (Recap) ≡

$$w^{k+1} \leftarrow w^k + \eta \sum_{i=1}^{N} (t_i - o_i) x_i$$

$$w^{k+1} = w^k + \eta \sum y_i x_i$$

(this misclassification.

| $t_i$ | $o_i$ | $x_i$ | $w$ | (to remove (this error) |
|-------|-------|-------|-----|-------|
| +1 | −1 | +ve | ⇒ −ve | w has to ↑ |
| −1 | +1 | +ve | ⇒ +ve | w has to ↓ |

\* **Separable** — $\exists \vec{w}$ such that data can be correctly classified.

Cannot have a linear separating hyperplane for every problem

→ Problem is linearly separable, there exists a solution weight vector.

Termination criteria — $\xi$ should be empty of perceptron.

↓ misclassification sample set.

\* If data is linearly separable; samples/

$$\boxed{y_i \cdot w^T x_i \geq \gamma_0}$$ ⟹ this means, they are correctly classified

↓ represents the Classification

→ small positive quantity (margin)

\* **Bounded** ⟹ $\|x_i\| \leq R$.

$\eta = 1 \Rightarrow$   $w^{k+1} = w^k + y_i x_i$

Let $w^*$ be the separating plane. (solution ≡ optima).

$$w^{*T} \cdot w^{k+1} = \underbrace{w^{*T} w^k + y_i w^{*T} x_i}_{y_i} \quad (> \gamma)$$

$$\underbrace{w^{*T} \cdot w^{k+1}}_{\text{previous term}} > w^{*T} w^k + \gamma$$

$$(\because \text{let } w^0 = 0)$$

$$w^{*T} w^k > w^{*T} w^{k-1} + \gamma$$

$$\vdots$$

$$w^{*T} w^1 > w^{*T} w^0 + \gamma$$

$$\Rightarrow \quad w^{*T} w^{k+1} > (k+1)\gamma$$

$$\downarrow \quad w^{k+1} = w^k + y_i \, w_i^T x_i$$

Squaring on both the sides. $\equiv$

$$\| w^{k+1} \|^2 = \| w^k + y_i \, w_i^T x_i \|^2$$

$$\underbrace{\| x_i \|^2}_{} \longrightarrow \text{Bounded}$$
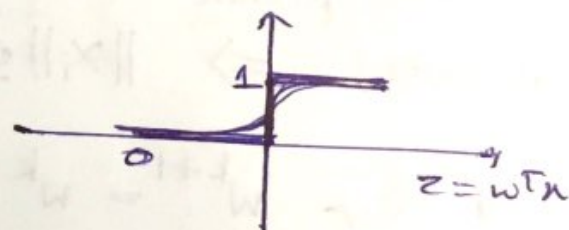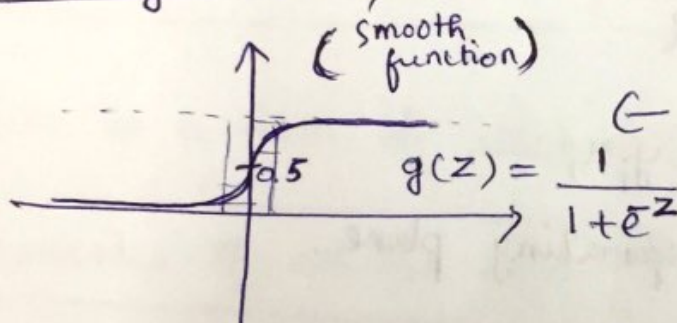
$$\downarrow$$

$$k^2 \gamma^2 < k R^2$$

→ (Perceptron algo guarantees that it converges in **finite** no. of steps)

$$K < \frac{R^2}{\gamma^2}$$

"(Converge in these many finite no. of steps)"

$\overset{*}{\equiv} \gamma \rightarrow$ large ; $K \rightarrow$ converge faster

(data sets are nicely separated)  (upper bound on number of iterations)

## Logistic Regression $\Rightarrow$

(smooth function)

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

$$z = w^T x$$

can be $\dfrac{1}{1 + e^{-kz}}$ (gets more closer to the axis)

$$\downarrow \quad g(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

$P(y = 1/x)$

If $w^T x >>$ large, belongs to class-1
else, belongs to class-0

$$g(w^T x) = P\left(y = \frac{1}{x}\right)$$

$$1 - g(w^T x) = P(y = 0/x)$$

$$\underbrace{\qquad\qquad}_{}$$

$$\downarrow$$

$\Rightarrow$ Classify as 1, if $P(y = 1/x) > 0.5$

Classify as 0,  " $\leq 0.5$

though both the curves (smooth & step) infer the same thing $\Leftarrow$ {

$w^T x = 0$ is the decision boundary.

$$J = \sum_{i=1}^{N} \left( y_i - g(w^T x_i) \right)^2 \quad \times \text{ Not a good function} \quad \left( \begin{array}{c} \text{due to} \\ \text{local} \\ \text{minima} \end{array} \right)$$

either 1/0          smooth curve $= \dfrac{1}{1+e^{-z}}$



convex

differentiable but it is not <u>convex</u>?



non convex.

but linear regression $J = \sum_{i=1}^{N} (y_i - w^T x_i)^2 \checkmark$

$$J = \sum_{i=1}^{N} \log \left( 1 + e^{-y_i\, g(w^T x_i)} \right)$$

$$= \sum_{i=1}^{N} y_i \log (w^T x_i) + (1 - y_i) \log (1 - g(w^T x_i))$$

→ ① * J = objective function + $\boxed{\text{Regularizer}}$ ?      — usually a function of $w$

diff (J) = diff (obj f$^n$) + diff (Regularizer)      additional constraints.
algo to.                                              beyond objective function.

→ ② Extending to multi-class (k-classes)

* [combining many binary to make a multi-class]

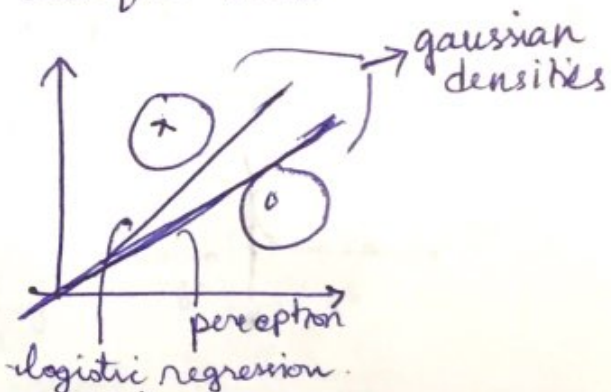# Perceptron :- (Recap) $\longrightarrow$ Finds the classifier line

$$J = \sum_{x_i \in E} (\text{Loss})$$

$$w^{k+1} \leftarrow w^k + \eta\, y_i\, x_i$$

$$w^{k+1} \leftarrow w^k + \eta(t_i - o_i)\, x_i$$

gaussian densities

perceptron
logistic regression.

may give one line which may be very close to the boundary
( but not optimal in bayesian sense)

Bayesian optimized

Perceptron may give this

$w^*$ (optimal $w$)

Perceptron should work for future data also.
↓                    (all the samples in the future)
"generalize"
        (requirement).

* Computation was done on existing training data but the model has to be extended for future values also. — ?

$$y_i w^T x_i > \gamma \quad \text{[classified samples]}$$

↓

at least $\gamma$ away from the decision boundary.

― Maximize "$\gamma$", ⟨margin⟩

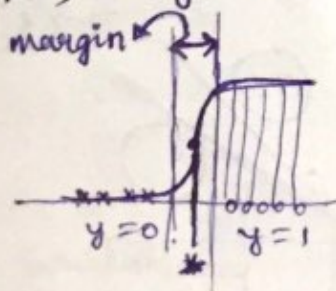↓ why? make sure it is less prone to error in the case of future values

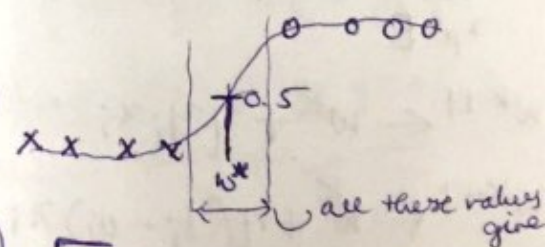Logistic Regression ⟹

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

If $g(w^T x) > 0.5$, $\hat{y} = 1 \equiv w^T x > 0$,
else        ($y = 0 \equiv w^T x < 0$) ⟶ again the same classification.

$P(y = 1/x)$ ― given $x$, probability that it belongs to class 1

margin

$P(y = 1/x)$
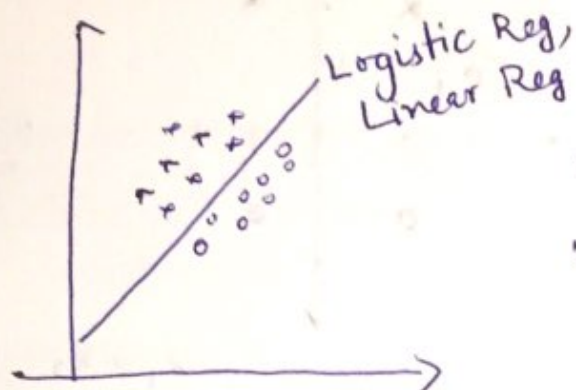
(Not the best) gives the first solution and stop

$\prod P(y = 1/x_i)$ and also for $y = 0$.

to maximize the product [probabilities]

gives an unique solution

all these values give

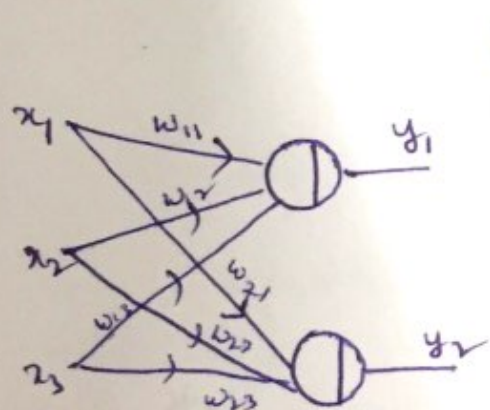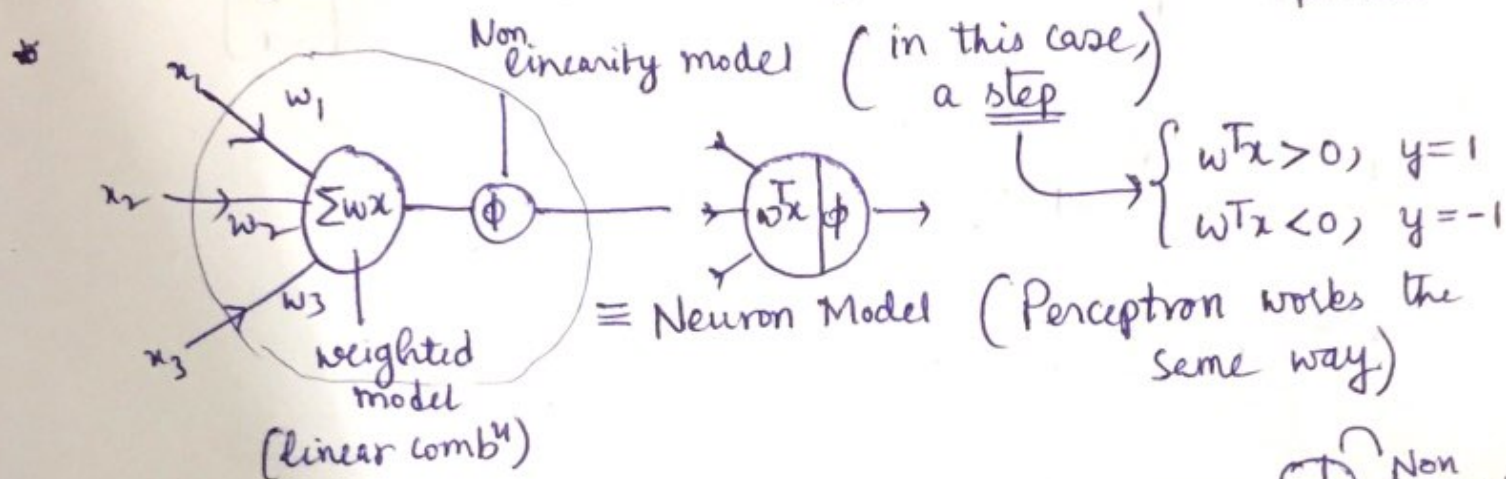→ slide this
(there exists a region)

Asymmetric costs ⅂ — costs of misclassifying samples.

Logistic Reg,
Linear Reg

when
⟶
samples
are
added

Logistic Reg
Linear Reg

Linear Regression changes but logistic regression gives the optimal.

Non linearity model ( in this case, a step )

$\Sigma wx$    $\phi$

weighted model
(linear comb⁴)

$\overline{w}^T x$   $\phi$ ⟶

$$\begin{cases} w^T x > 0, & y = 1 \\ w^T x < 0, & y = -1 \end{cases}$$

≡ Neuron Model (Perceptron works the same way.)

Non linearity is present

ForEx: sigmoid, $\left(\dfrac{1}{1+e^{-w^T x}}\right)$ step..

Perceptron

$\Rightarrow \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \phi\left( \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right)$

$2 \times 3$    $3 \times 1$

⟱

$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \phi \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

where let $\phi$ = sigmoid

$y_1 = \dfrac{1}{1+e^{-z_1}}$

$y_2 = \dfrac{1}{1+e^{-z_2}}$

$\text{Samples} = \{ (1,2) +1, (-1,2) +1,$
$\qquad\qquad (2,1) -1, (-1,-1) -1,$
$\qquad\qquad (1,-1) +1 \}$



Line $\equiv -x_1 + x_2 - 1 = 0$

$\begin{pmatrix} w_2 \\ w_1 \\ w_0 \end{pmatrix} w_0 = \begin{bmatrix} +1 \\ -1 \\ -1 \end{bmatrix}$

$\boxed{\eta = 0.1}$

$\qquad\qquad w_2 \quad w_1 \quad w_0$
$P \equiv \begin{bmatrix} 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} +2 \\ -2 \\ 1 \end{bmatrix} \begin{matrix} x_2 \\ x_1 \\ x_0 \end{matrix}$
$(-2,1)$

$0 = w^T x \equiv 1 + 2 - 1 = 2 > 0 \implies +1$
$\qquad\qquad t = -1$

$P \text{ (misclassified)}$

$\qquad\qquad\qquad \eta \, y_i \, x_i$
$w' \leftarrow w^0 + 0.1 \times (-1)(x_p)$

$w' = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} - 0.1 \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.9 \\ -0.8 \\ -1.1 \end{bmatrix}$

$> \dfrac{1}{1 + e^{-w^T x}} = \dfrac{e^{w^T x}}{1 + e^{w^T x}}$

**Multiclass :**
$\qquad\qquad\qquad P(y = c / x)$
$\qquad\qquad\qquad\qquad w_c^T x$

$* \; 2 \text{ class } - "w"$
$\quad k \text{ class } - ?$

$> \dfrac{e^{w_c^T x}}{e^{w_1^T x} + e^{w_2^T x} + \dots e^{w_k^T x}} \quad \left(\text{probability of being in } c^{th} \text{ class}\right)$

$> \dfrac{e^{w_c^T x}}{\sum\limits_{i=1}^{k} e^{w_i^T x}} \quad (\text{Soft Max})$

In the case of 2 classes;

$\dfrac{e^{w_1^T x}}{e^{w_1^T x} + e^{w_2^T x}} = \dfrac{e^{(w_1 - w_2)^T x}}{1 + e^{(w_1 - w_2)^T x}} \doteq \boxed{\dfrac{e^{w^T x}}{1 + e^{w^T x}}}$
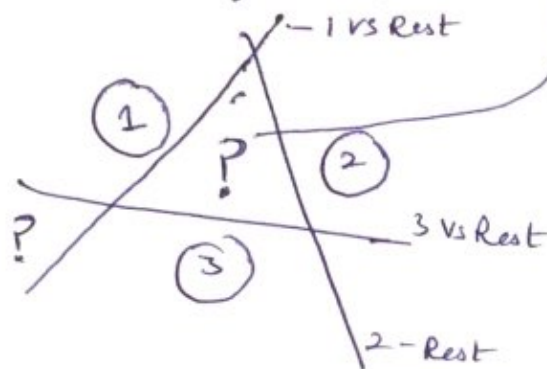
Problems — multiclass
Solutions — binary.

uses soft max type of classifier

DNN —C1, —C2, —C3, ... } K classes

last layer

K-classes

(1). ≡ K 1vs Rest classifier

with logistic regression kind of loss.

(one class is positive whereas all the others are negative)

(2). ≡ $KC_2$ 1vs 1 classifier

test for all K classifiers.

In, K 1vs Rest classifiers → what if all the classes say negative ("Rest")?

— 1 vs Rest

(1)

?

(2)

?

3 vs Rest

(3)

2 — Rest

gives ambiguous results.

$KC_2$ 1vs1 classifiers ⟹ (majority voting)

↓
more computational complexity

(for fx: 1000 class classification).

vs

DAG ↘
Evaluation Complexity ≡ O(k)

DAG — Directed Acyclic Graph.

1vs2

1vs3     2vs3
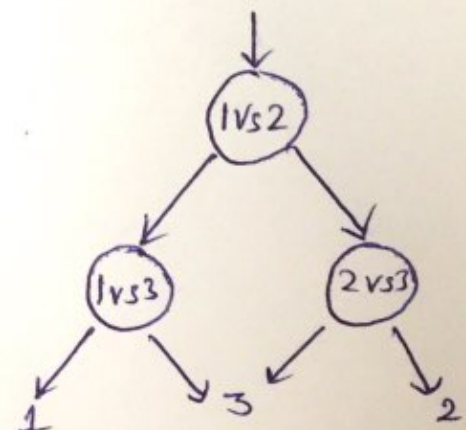
1    3    2

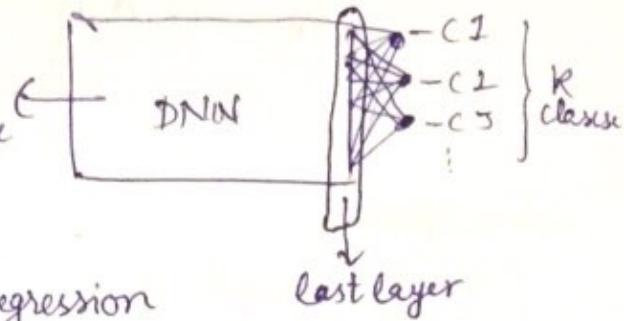Regularize Objective Function :—

↓
define a new objective function "regularized"

$$I_R = J + \lambda \|w\|^2 \longrightarrow \text{small}.$$

[Ex: perceptron, logistic Regression, SVM]

$$J_R = J + \lambda ||w||^2 + 1000 w_3^2 + 16 w_4^2$$

$\downarrow$ small.    $\downarrow$ 0    $\downarrow$ 0

(to remove a lot of dimensions / features)

— computationally good at the test time.

$W \rightarrow$ sparse

(lot of elements = 0, supress $\uparrow$ features)

$\rightarrow$ Find $w$ as sparse as possible.

● can use $L_1$ norm.

✓ $L_2$ norm (easier) $\Rightarrow$ this to make it small

$\left( ||w||^2 \right)$