

# 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 Norvig

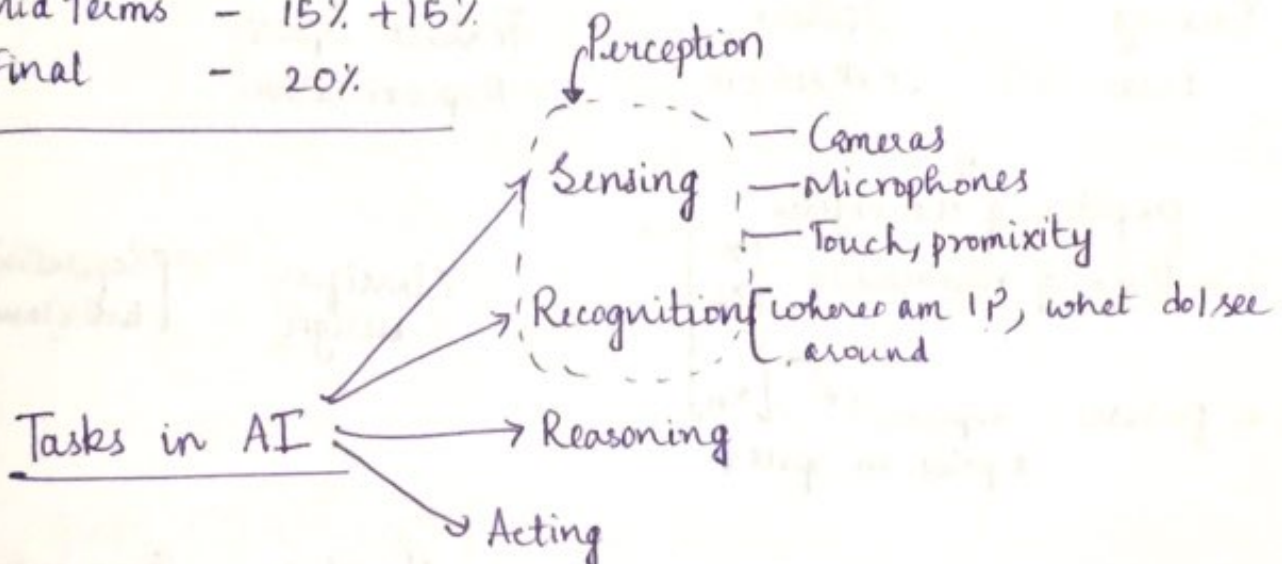
## Grading Schemes-

Mini Project - 20%

Homeworks - 30%

Mid Terms - 15% + 15%

Final - 20%



↓ "Re" + "cognize"  
Recognition - Identification of a pattern as a member of a category we all know, or familiar with.  
↓  
associating it with a class of objects } "examples"

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

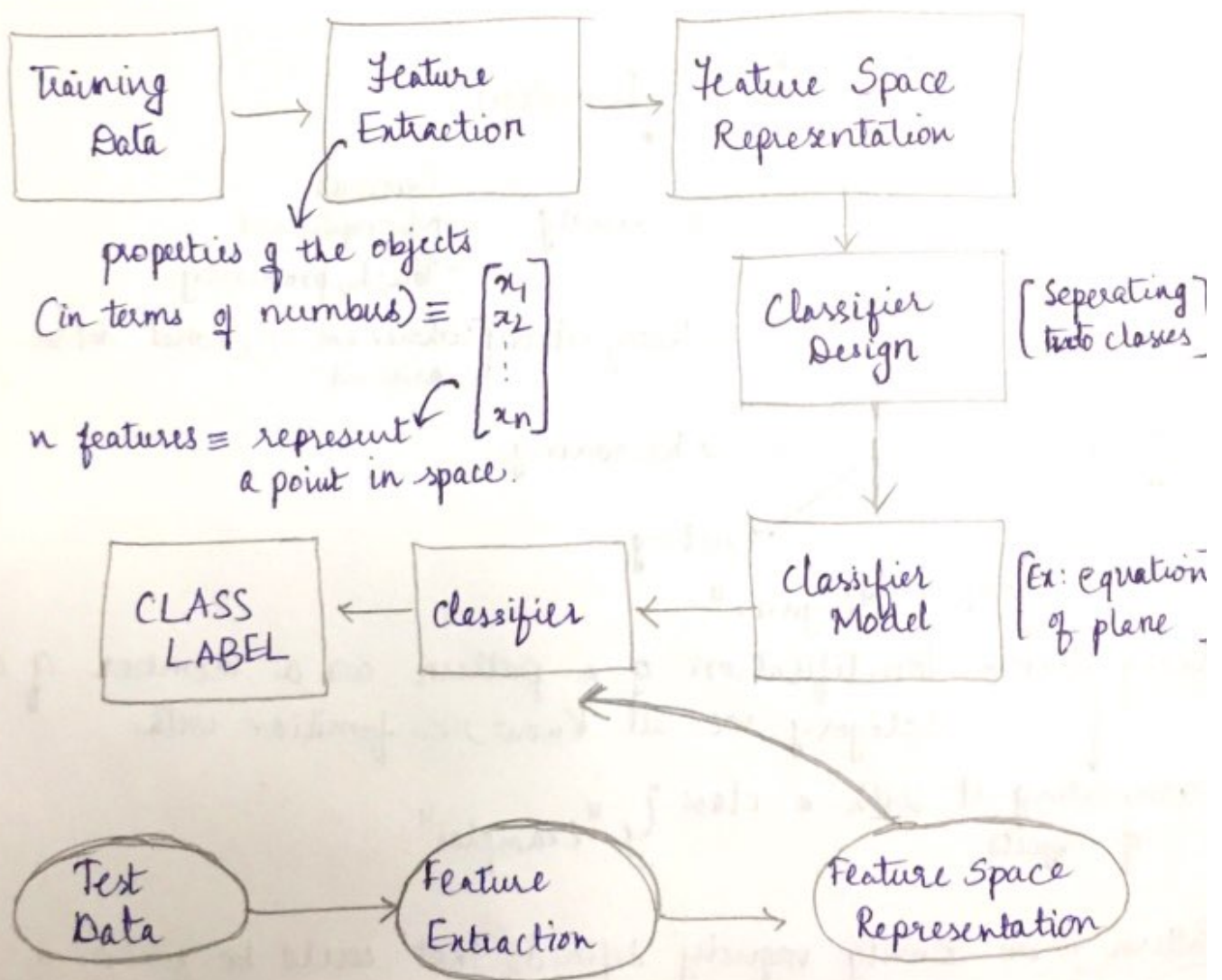
Pattern Recognition — inferring a generality from a few examples

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

\* class  $\begin{cases} \nearrow \text{Inter-class variability} \\ \searrow \text{Intra-class variability} \end{cases}$   
represented as  $[w_1, w_2, \dots, w_n]$

"train/teach the machine with a large dataset"

Pattern Recognition Process :-



\* class — modeled by a probability density function,  $P(x)$   
↓  
class-conditional  $\equiv P(x/w_i)$



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

## Mathematics :- Revision

→ Linear Algebra :

$f(x) \equiv f(\alpha x + \beta y) = \alpha f(x) + \beta f(y)$  ↗ for a function to be linear

$$a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1$$

$$a_{21}x_1 + a_{22}x_2 + a_{23}x_3 = b_2$$

$$a_{31}x_1 + a_{32}x_2 + a_{33}x_3 = b_3$$

Basic Representation  $\equiv$  matrices and vector form

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

$A$ 
 $X$  (3D-vector)
 $B$

$$AX = B$$

(Norm/Length)

Size of Vector  $\equiv$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\|X\|$$

↗ absolute value.

$$\|X\|_2 = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

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

$$L_p\text{-norm: } \|X\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}$$

↙ doesn't have any physical significance.

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

$L_0\text{-norm} \equiv \|X\|_0 = \sum_{i=1}^n |x_i|^0$  ↗ number of non-zero dimensions / entries in the vector.

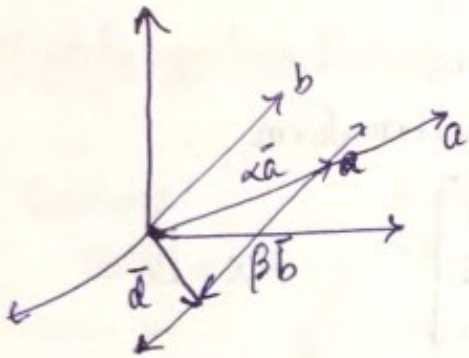
$L_\infty$ -norm  $\equiv ||x||_\infty =$  Maximum value of the entries for ex!  $\begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$\Rightarrow$  Span of a set of vectors -

linearly independent  $\equiv$  when any combination of two vectors don't result in the third vector

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

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



$$\vec{c} = \vec{a}\alpha + \vec{b}\beta$$

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

form a basis

if Rank  $\geq 2$ , 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 (linearly independent)  
 $\downarrow$   $|w_{\text{norm}}| = 1$

perpendicular ( $90^\circ$ ) [Dot product of orthogonal vectors  $= 0$ ]

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



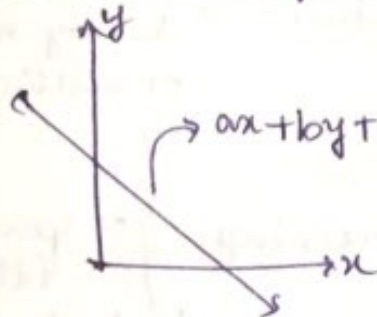
$$\Rightarrow \text{DOT PRODUCT} \equiv \begin{matrix} X & Y \\ \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} & \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \end{matrix} \Rightarrow x_1 y_1 + x_2 y_2 + x_3 y_3 = X^T \cdot Y \quad (\text{Inner Product})$$

$$XY^T = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} [y_1 \ y_2 \ y_3] = \begin{bmatrix} x_1 y_1 & x_1 y_2 & x_1 y_3 \\ x_2 y_1 & x_2 y_2 & x_2 y_3 \\ x_3 y_1 & x_3 y_2 & x_3 y_3 \end{bmatrix} \quad (\text{Outer Product})$$

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

$\hat{a} = \overline{X} \times \overline{Y}$   
 perpendicular to both  $X$  &  $Y$

Representation of line:

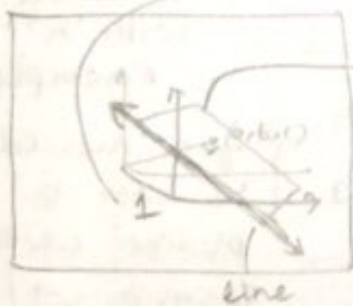


$$\begin{bmatrix} a \\ b \\ c \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$w$        $x$

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

how do you relate both?  
 $w^T x = 0$   
Space (augmented vector)



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

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

Distance of a point from a plane?  
Normalise  $w$  and take the dot product.

Linear Transformation  $\Leftrightarrow$  Matrix Multiplication

by inverting the matrix, we can invert the transformation

Eigen Values and Eigen Vectors

$$M = \underbrace{U}_{\text{rows = eigen vectors}} \underbrace{D}_{\text{orthonormal vectors}} \underbrace{V^T}_{\text{eigen values}}$$

physical entity — represented as a vector  
 $x_1, \dots, x_n \Rightarrow$  vectors  
 with  $d \equiv$  dimension  
 where  $x_i \in \mathbb{R}^d$

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$   $d$  — dimension  
 observations

physical entities are represented differently in different domains.

For Ex:

$$x_i \rightarrow y_i$$

$\in \{0, 1\}$  (say, spam or not)

$\{-1, +1\}$

$\{1, 2, \dots, k\}$

(say, professional, person, etc)  
 $k$ -class.

$y_i \in \mathbb{R} \equiv$  Real, continuous measure.

It is called as "Regression"

representation of an email  $\in \mathbb{R}$   
 spam or not.

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

examples for learning/training

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

$\equiv$  Spam filter to be built with "n" examples

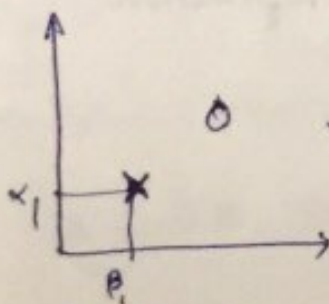
Find a function,  $f: x \rightarrow y$

So when given a new email,  $\equiv f(x)$

you should be able to predict whether spam or not

$x_1, x_2, x_3$   
 $\begin{bmatrix} x_1 \\ \beta_1 \end{bmatrix}, \begin{bmatrix} x_2 \\ \beta_2 \end{bmatrix}, \begin{bmatrix} x_3 \\ \beta_3 \end{bmatrix} \rightarrow$  2 features with different values  
 0 1 0 (spam / not spam)

can be represented on a plane.



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

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

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



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

↓  
 OPTIMIZATION problem  $\equiv$  Find  $w$  that best solves the problem.

↓  
 we need an "objective function"

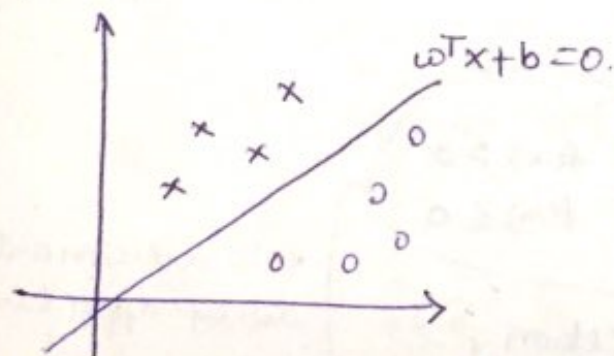
Data  $(x_i, y_i)$  pairs

↓  
 Problem  $\equiv$  Find best  $w$  which minimizes  
 ↓  
 Objective function  $\equiv$  Error / Loss

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

changing the loss function can also changes the optimization Algorithm. because  $N \gg$

Representation  $\equiv$



x - Spam  
 o - Not Spam

Case ①

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

Spam if  $f(x) < 0$

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

↓  
 Discriminant function

Case ②:

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

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

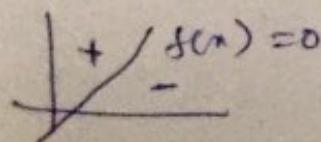
$$f_1 - f_2 > 0$$

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

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

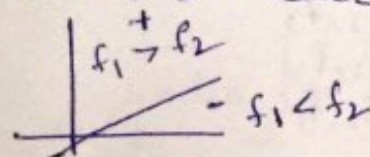
In the 1st case  $\equiv$

$f(x) \equiv$  function discriminates spam from notspam



$$f_1 - f_2 = (w_1 - w_2)^T x + b_1 - b_2 = w_3^T x + b_3$$

In the 2nd case  $\equiv$



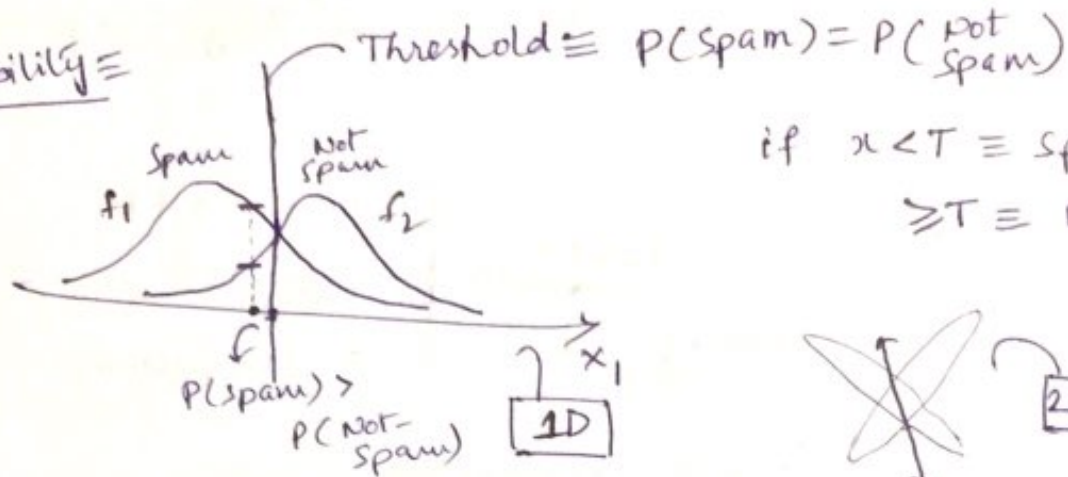
? (when to use this approach)

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

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

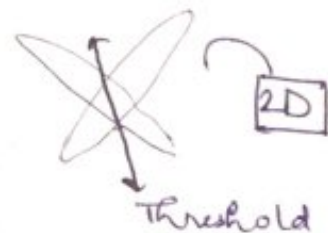
both the viewpoints are equivalent in this case.

Probability  $\equiv$



if  $x < T \equiv$  Spam

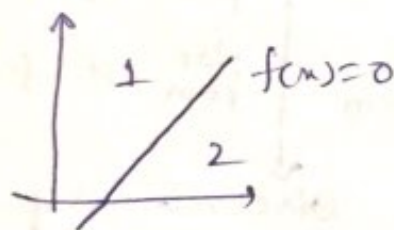
$\geq T \equiv$  Not Spam.



$S_0 \equiv$  Spam

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

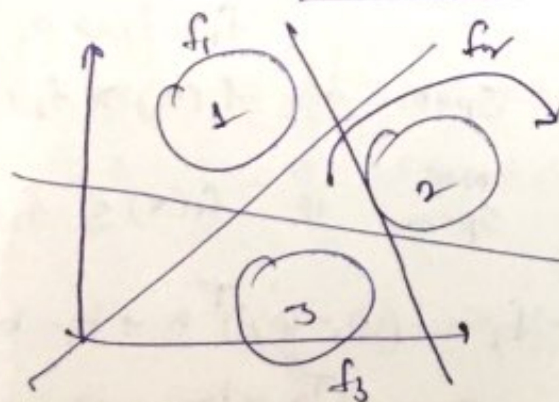
2-Class Classification Problem:



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

this decision rule doesn't apply here.

What about a 3-class classification?



What about this sample space.

"No-man"

"multi-class classification problem"

Decision Rule:

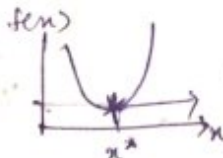
- 1:  $f_1 > f_2, f_3$
- 2:  $f_2 > f_1, f_3$
- 3:  $f_3 > f_1, f_2$




Output - Discrete (Not Continuous)

↓ discrete  
becomes an optimization problem.

Loss function ↓  
Optimization Problem

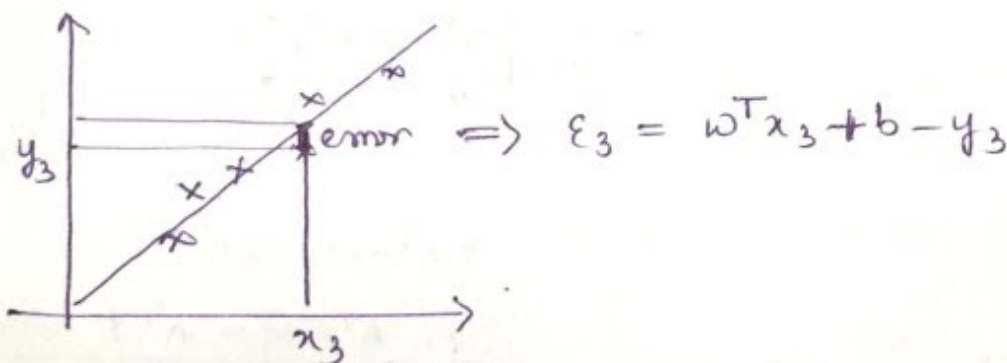
Optimization → Convex  $\equiv$   single unique minima.  
→ Non-convex

$\equiv f(x)$    
we can solve it. → maybe the best minima to get  
how to find? — brute force search.

Regression  $\equiv$

$(x_1, y_1) \dots (x_n, y_n)$  where  $y_i \in \mathbb{R}$ .

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

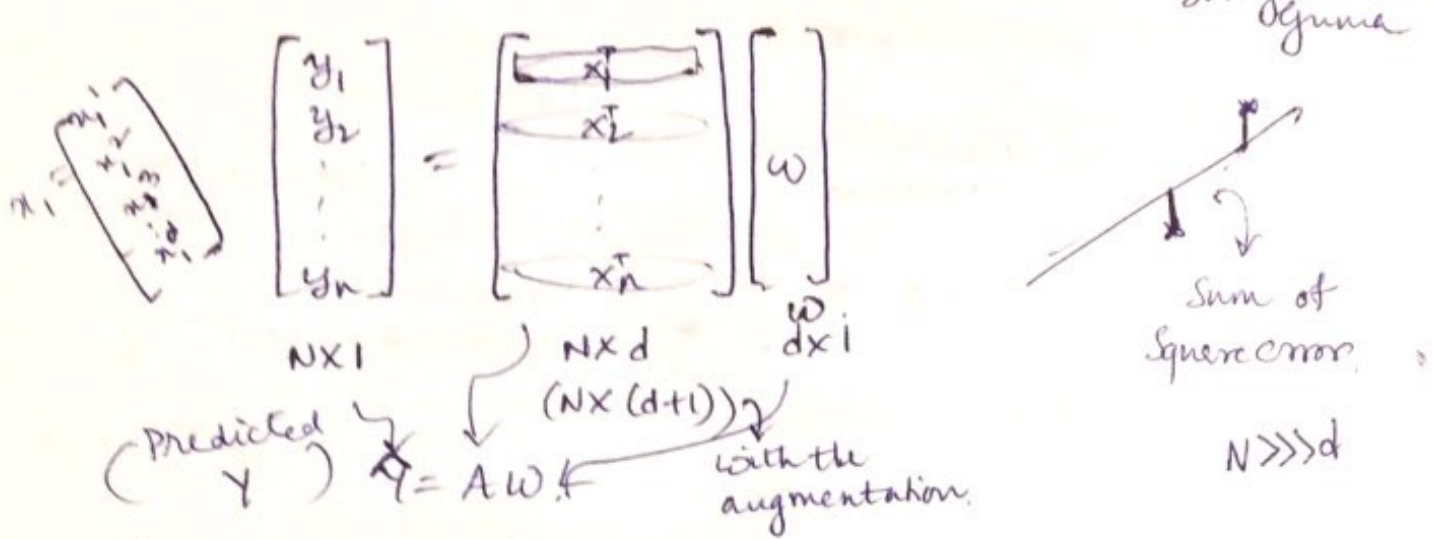


Error  $\equiv \sum_{i=1}^N [y_i - (w^T x_i + b)]^2$   
 $\Downarrow$   
 $\min_{w, b} \sum$

$\Downarrow$   
 predicted value by the function.

$\min_w \left( \sum_{i=1}^N (y_i - w^T x_i)^2 \right)$   
 can be +ve / -ve  
 [above or below the line]  
 Minimum Square Error (MSE)

$w^T x$   
 $\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + b$   
 $\downarrow$  Crawl to remove b  
 $\Downarrow$   
 $\begin{bmatrix} 0 \\ 0 \\ 0 \\ b \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$   
 $x$  is augmented with 1



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

(N x 1)<sup>T</sup> (N x 1) = (1 x 1)

sigma is gone

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

$$[\text{Differentiate and equate to 0}] \Rightarrow \frac{\partial}{\partial w} [Y^T Y + (AW)^T AW - 2(AW)^T Y]$$

$$= \frac{\partial}{\partial w} [Y^T Y + W^T A^T A W - 2W^T A^T Y] = 0$$

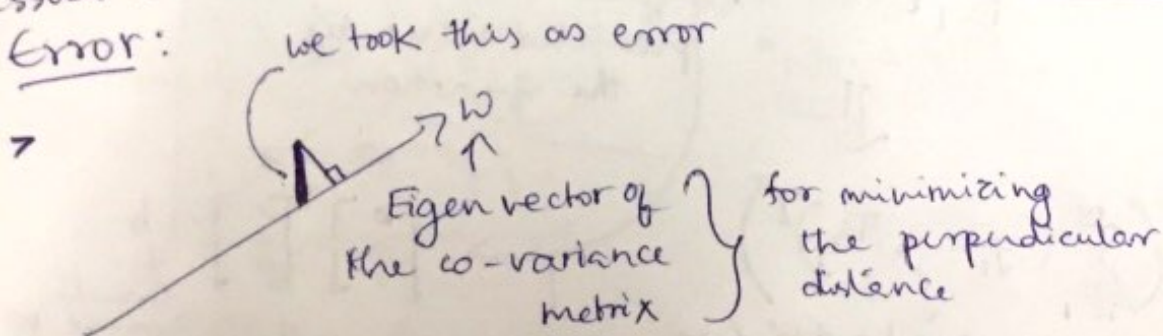
0      2A<sup>T</sup>A      2A<sup>T</sup>Y

$$2A^T A W - 2A^T Y = 0$$

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

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

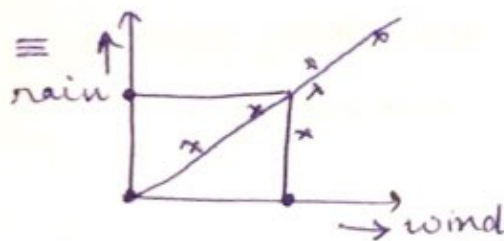
Tom Minka = Matrix Differentiation.  
Issues with Error:



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



# Examples of Regression



$$Y = f(x) = w^T Z$$

$$\begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$Y = ax + b$$

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

(augmented vector)

$$y = w_1 x_1 + w_2 x_2 + w_3$$

$$\Rightarrow x = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix} \quad \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

— not restrictive  
— due to augmentation (not passing thru origin  $v$ )  
New vectors = Non linear functions embedded in  $w^T x = 0$ .

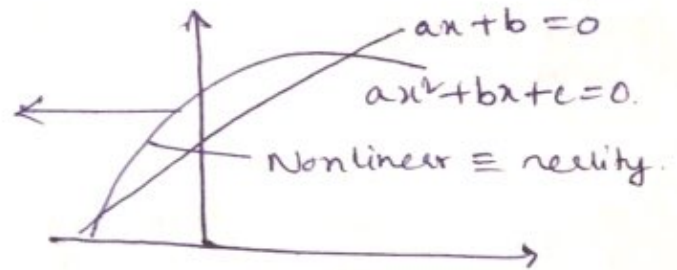
— this even works for non-linear

line fitting

can be represented as

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

$$\hookrightarrow w^T x = 0$$



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

$\Rightarrow Y = A w$   
 $\begin{matrix} 1 & | & d \times 1 \\ N \times 1 & & N \times d \end{matrix}$   
 To find  $w \equiv w = A^{-1} Y$   
 (doesn't work)  
 as long as A is not singular and square matrix ( $N=d$ ).

Not restrictive

$x \rightarrow Z$  can be mapped

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_x \Rightarrow \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \\ x_1 x_3 \\ \vdots \end{bmatrix}_Z$$

these are not independent.

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

pseudo inverse of A.

approach  $\equiv$  Minimum Square Error

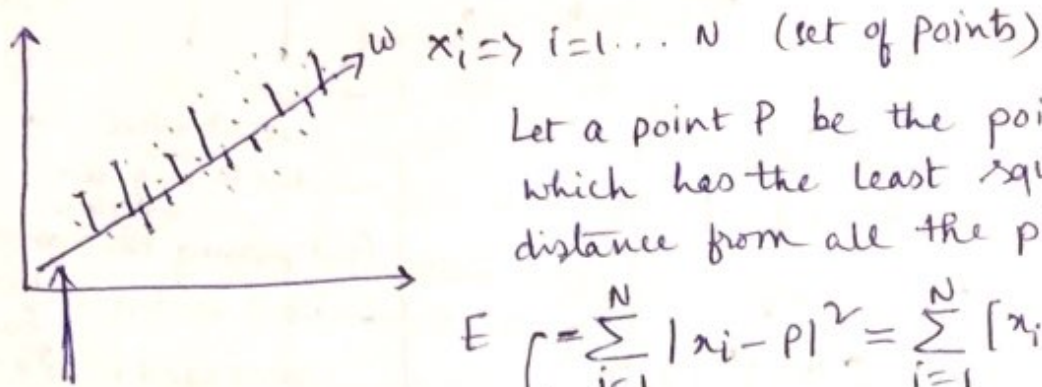
$\rightarrow$  Doesn't minimize perpendicular (orthogonal) dist.

\*  $(A^T A)^{-1} \equiv$  computational storage

— Data comes incrementally (data is always not offline)

$w$  - needs to be incrementally updated.

⇒ Orthogonal Distance Minimizing → find  $P$ , a point  
→ find  $w$ , the direction.



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

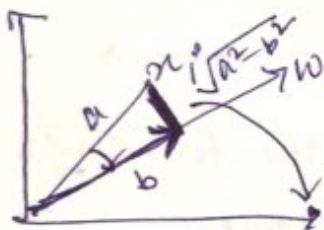
$$E = \sum_{i=1}^N |x_i - P|^2 = \sum_{i=1}^N [x_i - P]^T [x_i - P]$$

eigenvector for minimum, from the covariance matrix.  
(passing thru the mean)

$$\frac{\partial E}{\partial P} = 0$$

$$\frac{\partial}{\partial P} \sum_{i=1}^N \left( \underbrace{x_i^T x_i}_0 + \underbrace{P^T P}_{2P} - 2 \underbrace{P^T x_i}_{2x_i} \right)$$

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



Normalized.

$$w^T w = 1$$

Norm of  $x$ 's are fixed.

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

independent of  $w$ .

$$\sum_{i=1}^N P = \sum_{i=1}^N x_i^2$$

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

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

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

$$\downarrow$$

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

$$(xw)^T xw$$

$$w^T x^T x w$$

$$\begin{pmatrix} (w^T x)^T & (w^T x) \\ x^T w & w^T x \end{pmatrix}$$

$$\text{Max: } w^T x^T x w$$

wrong.

look at regression.



Maximize  $\equiv w^T X^T X w$  st  $w^T w = 1$

$\downarrow$   $w$  is eigen vector of (largest) eigen value  
 $\Downarrow$  why?

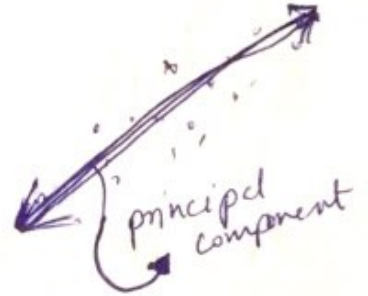
$(X^T X) \equiv$  PCA  
 $\uparrow$  covariance matrix

$Ax = \lambda x$

Max  $(w^T X^T X w)$   
 $\uparrow \lambda w$   
 $w^T \lambda w$

To Maximize this quantity  $w^T \lambda w$ , we have to take largest value of  $\lambda$

largest eigen value



max  $(w^T A w)$

$w^T \lambda w \equiv \lambda w^T w = \lambda$   
 $\downarrow$  eigen value  $\downarrow 1$

Direction is given by  $w$ , where  $w$  is the eigen vector corresponding to the largest value of  $X^T X$

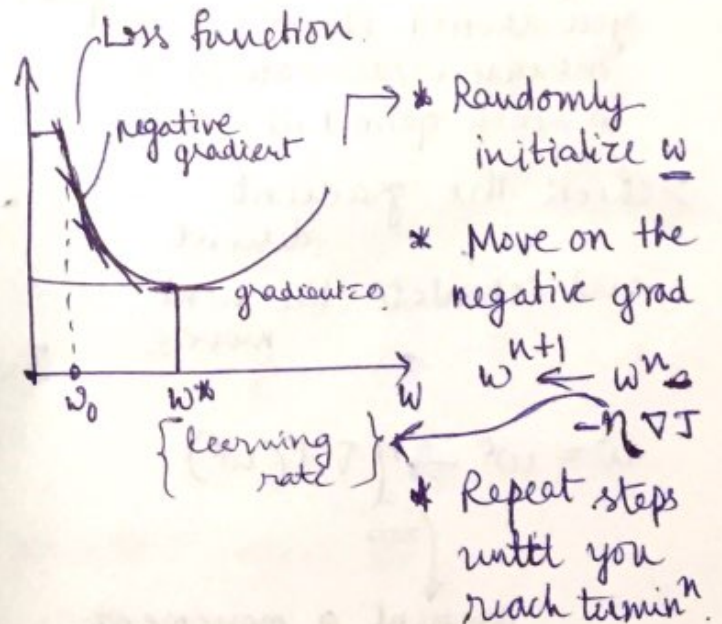
$\downarrow$   
 Mainly used in PCA

## Gradient Descent $\Rightarrow$

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

$\downarrow (y_i - f(x_i))$   
 loss function = need to be minimized

- \* Goal
  - objective
- \* Data  $(x, y)$
- \* Optimization
  - closed form Algo
  - Iterative Algo  $\rightarrow$  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.

$$\rightarrow y_i \sim f(x_i) \quad \forall_i$$

$\swarrow \searrow$   
 $R \quad \{0,1\}$   
 $C$

$f(x) = w^T x$  Assumed it to be a linear function.

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

$w^T x \equiv$  Not restrictive

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

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

## Gradient Descent:

> Make a random guess.

$\downarrow$   
 you should either increase  $w$ /decrease  $w$  to reach optimal  $w$ .

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

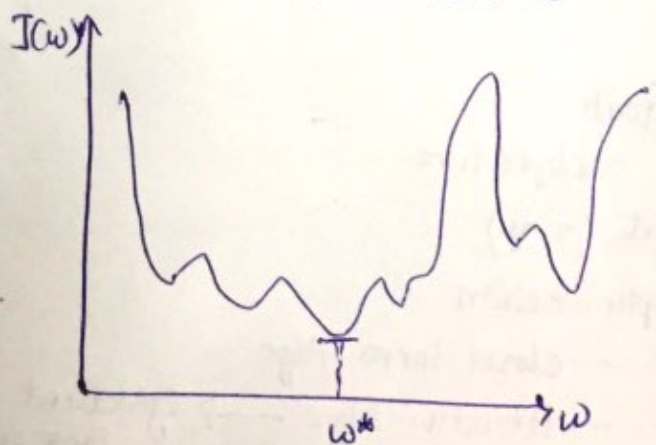
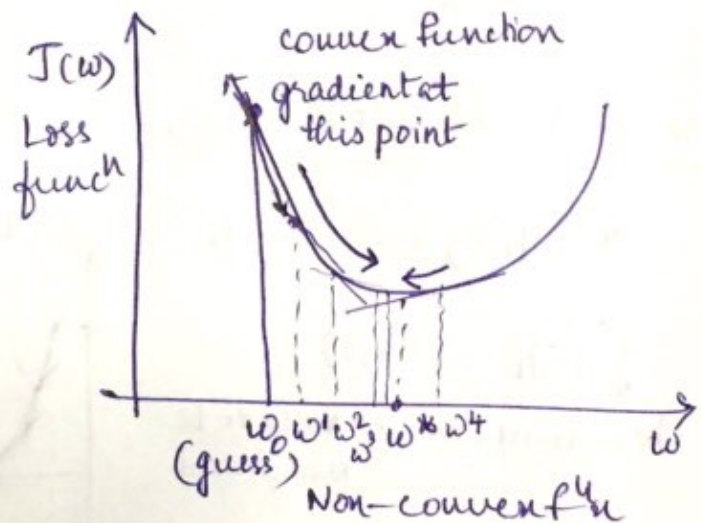
$$w = w^0 - \eta \nabla J(w^0)$$

$\eta$   $\rightarrow$  amount of movement.

At the minimum;

$$\Delta J(w) = 0 \text{ and}$$

$w$  doesn't change  $\rightarrow$  Algorithm converges





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

In nonconvex function; the minimum <sup>(local)</sup> that we get may not be the best one by using the same equation.

$$\nabla J(w^0) < 0.001$$

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

Learning Rate

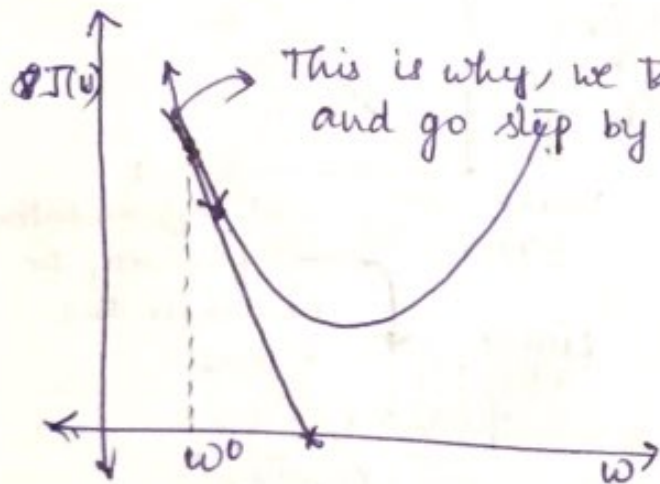
$$w^{n+1} \leftarrow w^n - (\eta) \nabla J(w^n)$$

changing it doesn't give much of an advantage

when we make a move into the negative gradient

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

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



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

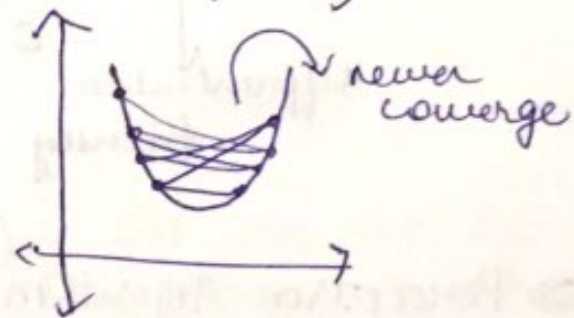
— (if function nature is known, it would have been easy).

⇒ 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~~ <sup>prev</sup> step and ~~next~~ <sup>present</sup> step are same.

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



for  $\eta$  :

$$\eta \leftarrow 1.5 \times \eta, \text{ if prev dir} = \text{present dir}$$

$$\eta \leftarrow 0.5 \times \eta, \text{ if dir chenge.}$$

Taylor series expansion  $\equiv$

$$x = w^{n+1}$$

$$a = w^n$$

$$f(x) = f(a) + f'(a) \cdot (x-a) + f''(a) \frac{(x-a)^2}{2} + \dots$$

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

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

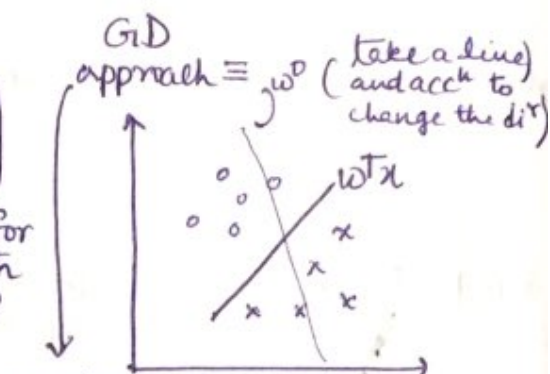
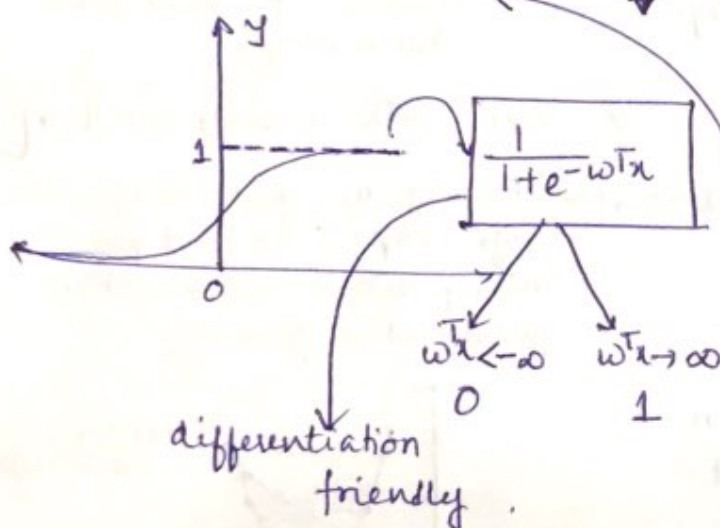
GD for New function value = Old function value - (Positive quantity).

Classification: (less than the old one)

$\Rightarrow$  Logistic Regression  $\equiv$  (instead of step  $\downarrow$  search for smooth step)

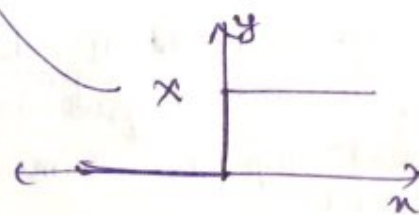
$\Rightarrow$  Perceptron Algorithm  $\equiv$

But instead;



minimum error.  $\equiv$  not differentiable friendly to eliminate the errors

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



$\Rightarrow$  Perceptron Algorithm  $\equiv$  (created a new objective func<sup>n</sup>)

$$J = \sum_{x \in \text{Misclassified Samples}} w^T x$$

Sum over misclassified samples. (instead of <sup>over</sup> all the samples)

Minimize the no. of such (error) classifications.

when, misclassified samples =  $\{\emptyset\}$   $\sim$  we found the "line".



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

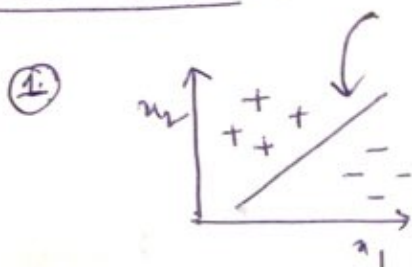
Perceptron Algorithm  $\equiv$  creating a new objective function.

every sample that gets misclassified is added to it

$$w^{n+1} \leftarrow w^n + \eta \cdot x \cdot y$$

If the samples are "linearly separable",  $\exists$   $f(w)$  which can separate +ve and -ve examples) perceptron will find that  $w$

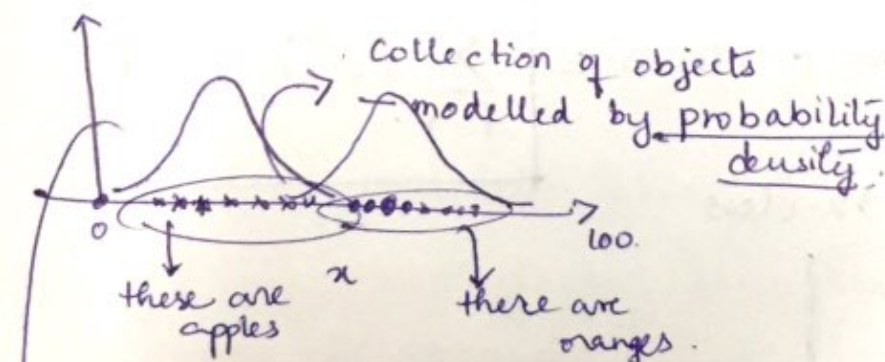
Classification  $\equiv$  Model of classifier is [Given this model, we cannot generate an apple]  $\equiv$  Line.



we model the line in the feature space  
Boundary (line between class A and class B.)  
 $\rightarrow$  Discriminative classifier/model

② Generative Model (we do not model the boundary)  $\equiv$  reduced representation  
For Ex: given a fruit, how apple-like or orange-like is it? instead model the data [as in "a point"]  
[Ex: model oranges and apples independently]

Model the classes?



[if I see 10.5, how likely is to classify it as apple]

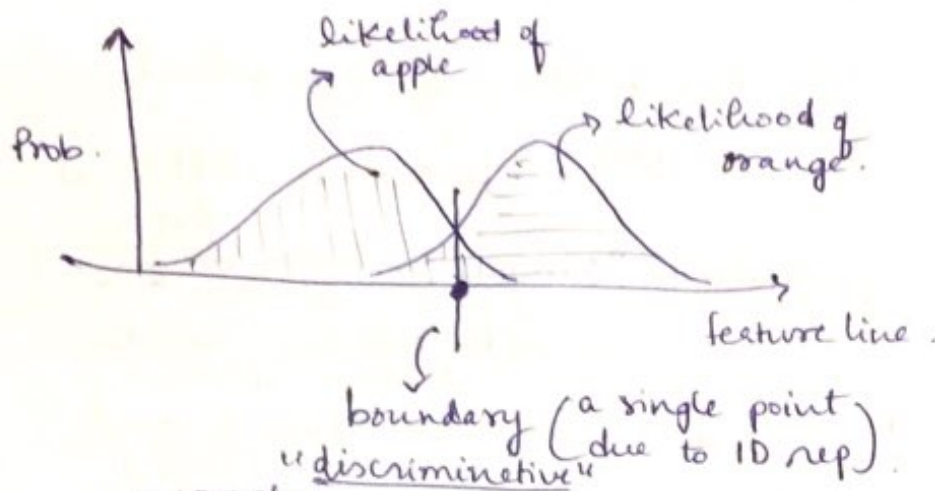
likelihood functions = Independently Learning about apples/oranges.

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

Discriminative classifier  $\equiv$  we cannot generate samples as it can lie anywhere over the apple side of the line.

Not rich enough to capture the details  $\rightarrow$  cannot explicitly model apples.

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



$w \equiv$  class

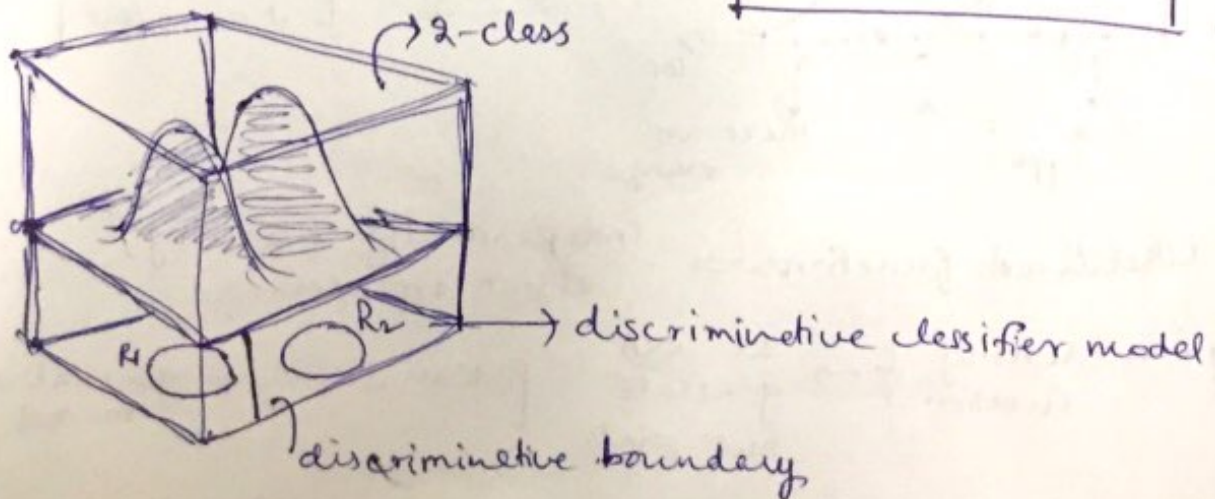
$P(w_i/x) \equiv$  "Probability" that given feature vector,  $x$  belongs to the  $w_i$  class.  $[0 < P < 1]$

feature vector

$P(w_1/x), P(w_2/x)$  can be compared to determine the class.

[in the case of 2-class problems]  $\equiv P(w_1/x) + P(w_2/x) = 1$

Whereas, if there are  $c$  classes  $\equiv \sum_{i=1}^c P(w_i/x) = 1$





$P$  [Probability] ;  $P(w_1/x)$ ,  $P(w_2/x)$ .

$p$  [pdf]

$$P(w_1, x) = P(w_1/x) \cdot P(x) \text{ (or)}$$

Joint Probability

$$P(x/w_1) \cdot P(w_1)$$

$$P(A, B) = P(A/B) \cdot P(B)$$

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

$$\therefore P(w_1/x) = \frac{P(x/w_1) \cdot P(w_1)}{P(x)}$$

Posterior prob.

Bayes Theorem

$\Rightarrow P(w_1) \equiv$  Prob. of ~~any~~ class irrespective of given sample.

$\hookrightarrow$  Prior probability / belief

$\Rightarrow P(x/w_i) \Rightarrow$  likelihood.

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

Posterior prob  $\equiv$  Post observing  $x$ , what is the probability  
 $P(w_1/x)$

Ex: Disease and Test

$$P(D) = 0.001$$

$$P(+/D) = 0.99$$

(If you have the disease, tests will say +, 99% (Accuracy of test))

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

Prob of having disease, given test came out +ve

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

$$P(+) = P(+/D) \cdot P(D) + P(+/ND) \cdot P(ND)$$

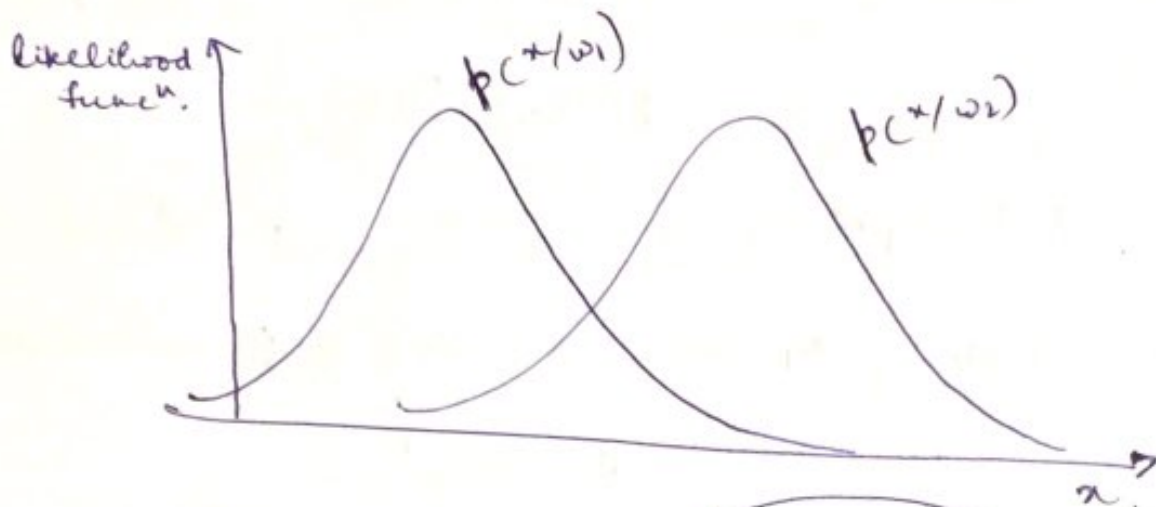
(person with D gets +ve result)      (person with ND)

$$= 0.99 \times 0.001 + 0.01 \times 0.999$$

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

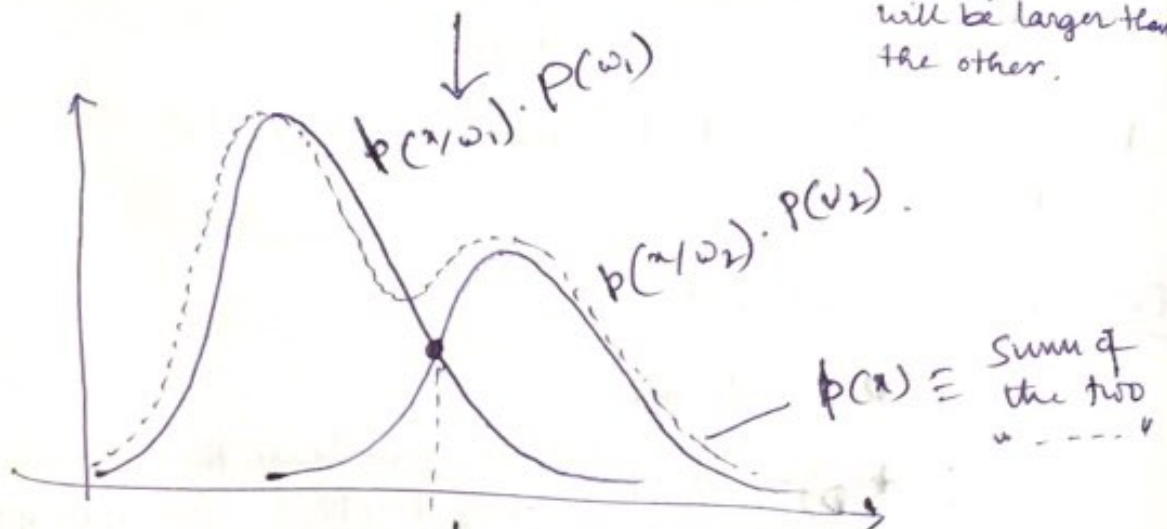
$$P(D/+) = \frac{0.99 \times 0.001}{0.01098} = 0.09 \text{ (9\%)}$$

↑  
chance of having the disease.

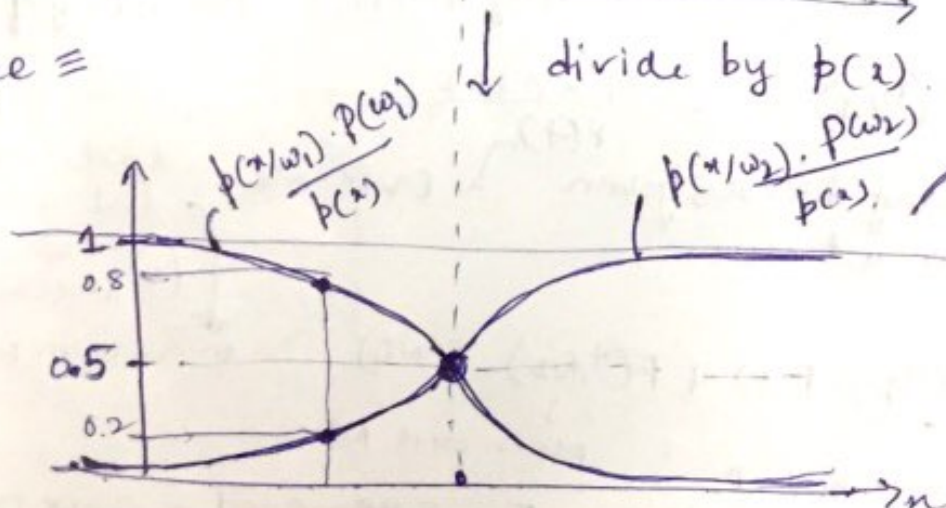


$$P(w_1/x) = \frac{p(x/w_1) \cdot P(w_1)}{p(x)}$$

constant  
↓  
one of the curves will be larger than the other.



Normalize  $\equiv$





# Bayes Theorem $\equiv$

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

↑  
after observing  $x$ ,  
state of nature  $\equiv w_j$

probability that state of nature,  $w_j$  is true. Pick  $w_j$  that has maximum prob.

Are all errors equally costly? No

so we define a loss function  $\equiv$

$$\lambda(\alpha_i/w_j)$$

↑ action,  $\alpha_i$   
↑ state of nature  $\equiv$  class

Loss is not symmetrical,

$E_1$  if class A is misclassified as class B or

$E_2$  if class B is misclassified as class A.

$$E_1 \neq E_2$$

→ there can be many actions

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

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

↓  
c classes.

→ Prob that  $w_j$  is true

loss incurred because of taking  $\alpha_i$  in  $w_j$  state of nature

Strategy to pick the best action = based on action plan

Evaluate the action plan?

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

prob of observing  $x$

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

$\alpha(x)$  — if I see this observation, I'll do this.

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

"Action plans can be eval. based on risks"

— Integrate over all possible  $x$ ,

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

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

From Bayes

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

$R^* = \text{Bayes Risk}$

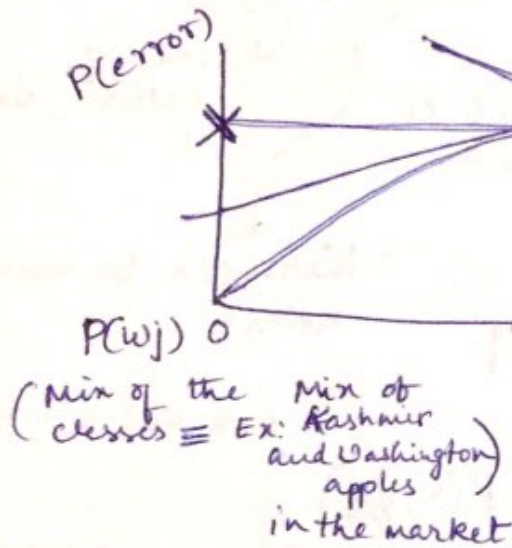
Bayes Action Plan.

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

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

Risk  $\propto$  P(error).

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



change in Risk = Linear

vs ground truth

if the market fluctuates, the maximum risk

is minimized  
choose such action plan.

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

$p(x) \rightarrow$  need not be computed to  
decide/compare [independent  
of  $w_j$ ]

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

(this is different g)

$$g_i(x) = \ln p(x/w_i) + \ln P(w_i)$$

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  $\Rightarrow p(x) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)}$

(can be ignored) constant  $= P(x/w_i)$

$$g_i(x) = -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) + \ln P(w_i)$$

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

\*  $g_i(x) = -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) + \ln P(w_i)$

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

$$g_i(x) = -\frac{1}{2} (x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i) + \ln P(w_i)$$

$$= -\frac{1}{2} (x^T \Sigma^{-1} x - 2 \Sigma^{-1} \mu_i^T x + \mu_i^T \Sigma^{-1} \mu_i) + \ln(P(w_i))$$

independent of i, class

linear boundary

same stretch and tilt but just shifted

$$g_i(x) = \underbrace{\Sigma^{-1} \mu_i^T x}_{\text{independent of i, class}} - \underbrace{\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i}_{\text{constant}} + \ln(P(w_i))$$

tends to be "linear"

Discriminant function.

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

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, hyperbole...)

Case 3:  $\Sigma_i = \Sigma$

$$g_i(x) = -\frac{1}{2} (x^T \Sigma^{-1} x - 2 \Sigma^{-1} \mu_i^T x + \mu_i^T \Sigma^{-1} \mu_i) + \ln P(w_i)$$

Loss :

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

$$g(z) = \begin{cases} +1 & ; z \geq 0 \\ -1 & ; z < 0 \end{cases}$$

In Classification -  $C \equiv$  % misclassifications.

Ex:  $y_i \neq g(w^T x_i)$ .

1.  $y_i = +1$   
 $w^T x_i = 10$ .  $\checkmark$  correct  $= 1 = y_i (g(w^T x_i))$

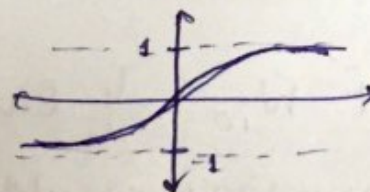
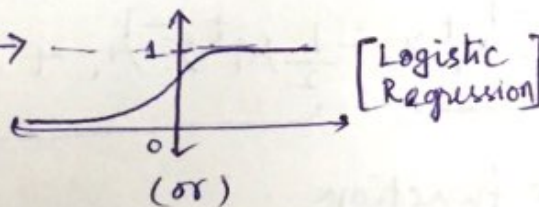
2.  $y_i = -1$   
 $w^T x_i = -10$ .  $\checkmark$  correct  $= 1 = y_i (g(w^T x_i))$

3.  $y_i = -1$   
 $w^T x_i = 10$   $\times$  incorrect  
(misclassified).

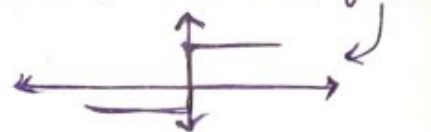
$$C \equiv J = \frac{1}{2N} \sum_{i=1}^N (1 - y_i \cdot g(w^T x_i))$$

— Good  
— Not optimisation friendly.

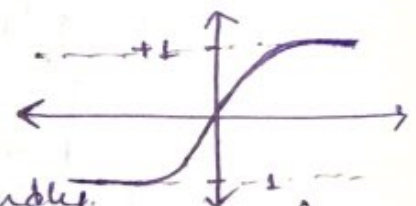
optimize it by  
changing the g



For this, instead of using



its better to use



this is  $\downarrow$   
(Not Logistic  
Regression)

Gradient Descent:

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

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

3.  $n \leftarrow n+1$   $\rightarrow$  changes, when loss function changes  
 $\rightarrow$  varies over time (iteration) [usually constant]

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  $\equiv f(y) = f(x) + (y-x)^T \nabla f(x) + \frac{1}{2}(y-x)^T H(y-x) \equiv$  operator  
Hessian matrix

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

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

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

$$\nabla f^T \cdot \nabla f$$

only proves that

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

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

$\Rightarrow$  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  $\eta$

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

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

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

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

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

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

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

H-matrix  
s-vector

$$\frac{\partial}{\partial s} = 0 \Rightarrow \nabla f + Hs = 0$$

$$s = -H^{-1} \nabla f$$

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



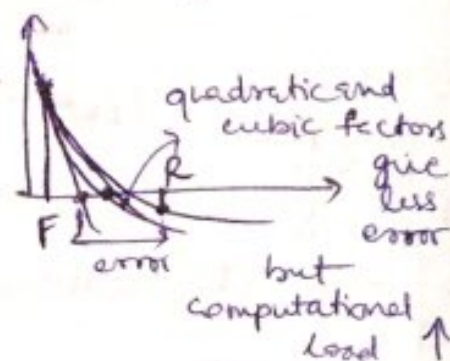
Not used  $\rightarrow$  maybe  $H$  might be singular.  $\checkmark$   
(inverse doesn't exist)

if  $|H|$  is very less  $\rightarrow$  inverse can be very sensitive

Not computationally attractive

First of all,  
how is this better  
than the first one?  
GD

$\rightarrow$  In the first case, approximating to the first order  
[GD] gives a smaller value (higher error)



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

Perceptron  $\equiv$  sum the loss only over the misclassified samples.  
objective function.

$$\sum_{x_i \in \Sigma} -y_i w^T x_i$$

$\hookrightarrow$  misclassified.

$$w^T x \geq 0, y_i = +1$$

$$w^T x < 0, y_i = -1$$

optimizing wrong objective func.

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

Update rule changes to  $\equiv$

$$w^{n+1} \leftarrow w^n + \eta \sum_{x_i \in \Sigma} y_i x_i$$

Changes with every step.

Not computationally different  $\leftarrow$  alternative  $\equiv$

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

if target = 1, output = 1  
 $1 - 0 = 0$  (correctly classified)

doesn't need to ~~be~~ misclassified set.  
go over



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

How does this work? (Intuitive Approach)

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

Error is present.

With these values,  $w$ 's going to increase.  $\Rightarrow$  increase  $w$   
 $\rightarrow w$  becomes +ve.

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

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

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

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

$(t - o)x \Rightarrow -ve \rightarrow w$  decreases.

$\rightarrow w$  becomes -ve.

$$w^T x < 0$$

$$x < 0 \quad \underline{\underline{w > 0}}$$

11<sup>th</sup> for diff values of  $t \& o$ .

$\rightarrow$  if a sample is misclassified, add/subtract it. (based on values of  $t$  and  $o$ )