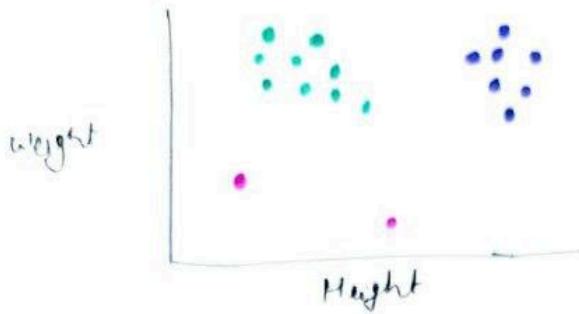


(4)



when we are done with all core points and left with only non-core points that can't be added to any clusters, we are done.

These non-core points is called outliers

- It's 2D, so considering circle, hypersphere else.
- Bottom two aren't core point as doesn't contain point greater than equal to  $m$ .

(i) Find core point

↓  
based on putting circle disc  $m$  points, check if falling as core point, randomly choose any core, find its neighbour

↓  
repeat.

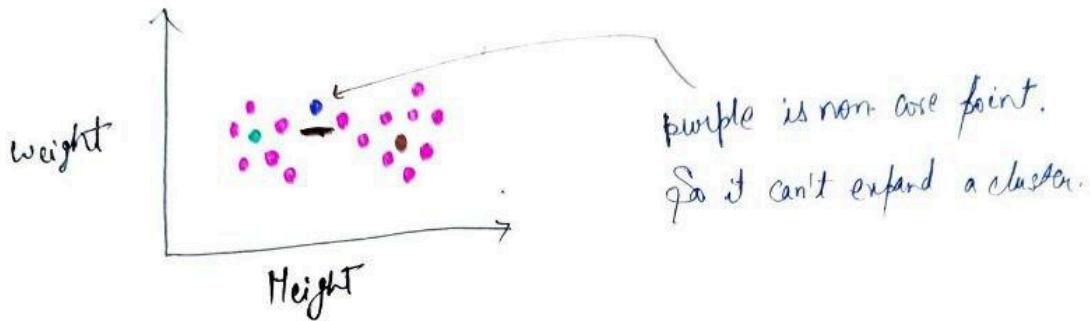
(ii) Repeat with different radius values find which ' $r$ ' considers more point.

(iii) Choose intuitive value of  $r$

There are two parameters  $\sigma$ ,  $m$ ,  $r_L$

## DBSCAN disadvantage

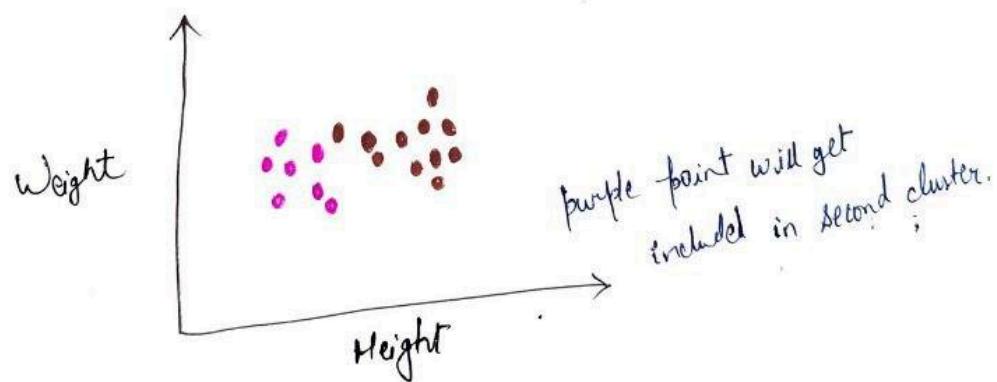
(5)



Case 1 :- If "Green" point is choosed for expansion :—

purple point will be included in first cluster.

Case 2 :- If "Brown" point is choosed for expansion :—

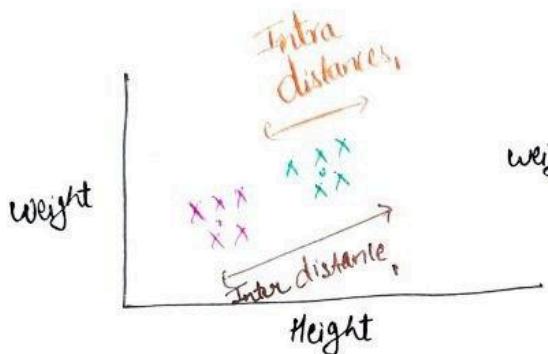


$$\frac{\text{intra distance}}{\text{inter distance}} > 0$$

intra cluster distance  
should be greater than  
not inter class cluster.

## Notion of Cluster Quality

(6)



→ Intra distance<sub>1</sub> > Intra distance<sub>2</sub> (not good)

. Intra distance should be lesser, more compact cluster.

Inter distance<sub>1</sub> < Inter distance<sub>2</sub> (not good)

→ Inter distance b/w cluster should be greater.

DBSCAN + Let  $S = \{x_1, x_2, \dots, x_n\}$  be a set of points to be clustered.

(i) Choose  $\epsilon > 0$  and  $n$

$\epsilon$  = radius,  $n$  = no of points atleast threshold

(ii) Let  $A_i$  be set of points that lies within a disc of radius  $\epsilon$  from  $x_i$ .

Do it for every  $x_i$ .

(iii) If  $A_i < n \rightarrow$  not consider for calculation.

. Points other than this is called core points.

- (iv) Take union of  $A_i$  and  $A_j$  if  $A_i \cap A_j \neq \emptyset$   
• go on doing it till no more union possible.

Disadvantage:-

For a non-core point which falls in mid of 2

clusters.

Clusters depends on from where I start clustering.

Advantage:-

Biggest advantage of DBSCAN is we don't need to know  
no. of cluster previously.

But we need to know  $m, R$ .

What to do with points lying on border?

There might be many points lying on border.

Recheck after converting these point to part of cluster.

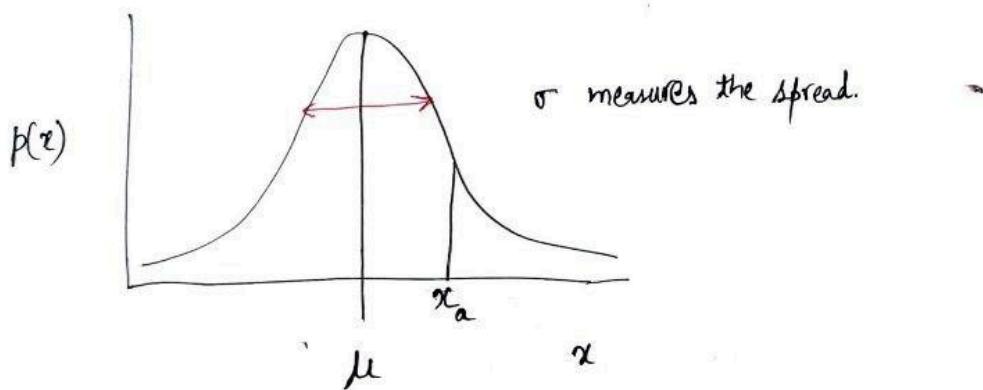
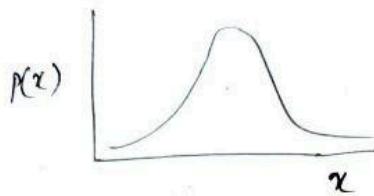
↓  
compare.

$\epsilon$  = radius       $r$  = no of points.

## Gaussian Mixture Model

①

- A gaussian (normal) distribution with mean  $\mu$  and std  $\sigma$



The belonging of point  $x_a$  to normal distribution is

$$P\left(\frac{x_a}{\mu, \sigma}\right) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_a - \mu)^2}{2\sigma^2}}$$

→ If  $x_a$  is closer to  $\mu$ ,

$x_a$  is more likely to belong to the distribution.

$\sigma$  = measure of spread.      ↑  $\sigma$  = distribution is wider

↓  $\sigma$  = distribution is sharp.

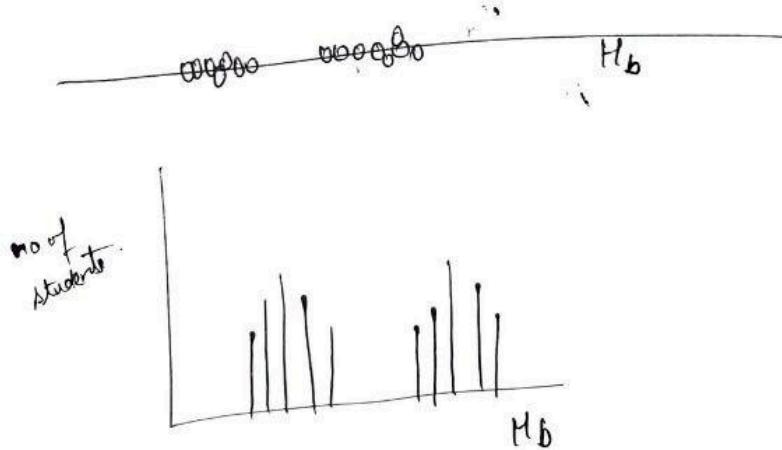
→ If  $x_a$  is closer to mean more likely to belong to distribution. ⑨

\* For multivariate normal distribution:-

For gaussian on multiple random variables, the belongingness of a high-dimensional point  $x_a$  to normal distribution is

$$P(x_a | \mu, \Sigma) = \frac{1}{\sqrt{2\pi\Sigma}} e^{-\frac{(x_a - \mu)^T \Sigma^{-1} (x_a - \mu)}{2}}$$

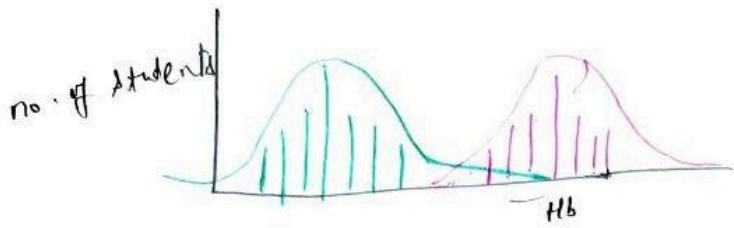
GMM or I. measure Hb of UTJ students:-



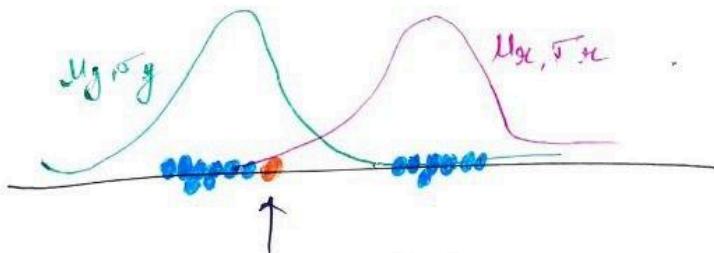
These data points can't be explained by single gaussian.

→ Total distribution of students can be modelled as mixture of two gaussians.

(10)



Inverse Fourier transform decompose sum of sin, cosine waveform.



Consider this student

- Orange point is close to green curve.
- Higher probable to green.

$$e^{-\frac{(x^n - \mu_g)^2}{2\sigma_g^2}}$$

$$p_g(x^n | \mu_g, \sigma_g) = \frac{1}{\sqrt{2\pi}\sigma_g} e^{-\frac{(x^n - \mu_g)^2}{2\sigma_g^2}}$$

$$e^{-\frac{(x^n - \mu_x)^2}{2\sigma_x^2}}$$

$$p_n(x^n | \mu_x, \sigma_x) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(x^n - \mu_x)^2}{2\sigma_x^2}}$$

→ Clustering using GMM :-

(11)

- (i) Randomly initialize  $K$  gaussians.
- (ii) For all data points, check which Gaussian it belongs to.
- (iii) Based on belongingness assign each data point to cluster.
- (iv) Find  $\mu_i$  and  $\Sigma_i$  from each cluster and update Gaussians.
- (v) repeat step (ii)-(iv) unless termination criteria is satisfied.

Termination Criteria :-

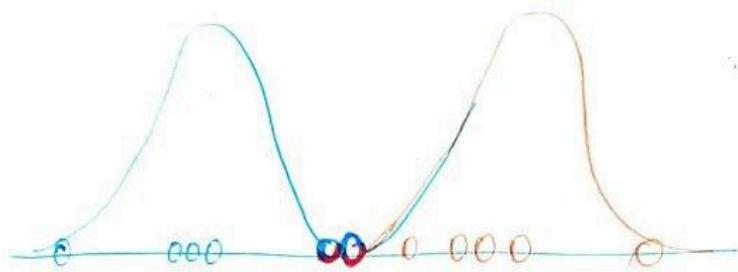
Many termination criteria are possible:-

- (i) we may terminate if  $\mu_i, \Sigma_i$  for each cluster don't change significantly over a few epochs.
- (ii) we may terminate after a certain predefined no of epochs.
- (iii) we may terminate when log likelihood of all data points doesn't change over a few epochs

$$J(\theta) = \sum_{j=1}^N \ln P\left(\frac{x_j}{\theta}\right)$$

(12)

## Soft Clustering using GMM



- Suppose we take  $K$  clusters
- Randomly Initialize  $K$  Gaussians
- The  $i^{\text{th}}$  Gaussian is given by  $\sim N(\mu_i, \Sigma_i)$
- Suppose we take  $K$  clusters  $C_1, C_2, \dots, C_K$ .
- Consider a data point  $x^n$ .
- $p(x^n) = \text{prob. that } x^n \text{ is in } C_1 \text{ or prob. that } x^n \text{ is in } C_2 \text{ or}$   
 $\dots \text{ prob. that } x^n \text{ is in } C_K.$

$$\text{So } p(x^n) = p(x^n, C_1) + p(x^n, C_2) + \dots + p(x^n, C_K) = \sum_{i=1}^K p(x^n, C_i)$$

$$\begin{aligned} \cdot p(x^n) &= \sum_{i=1}^K p(x^n, C_i) = \sum_{i=1}^K p\left(\frac{x^n}{C_i}\right) p(C_i) = \sum_{i=1}^K p\left(\frac{x^n}{C_i}\right) \pi_i \\ \cdot \pi_i &= p(C_i) \rightarrow \text{mixing coefficient.} \end{aligned}$$

Using Baye's theorem we calculate probability that  $x^n$  belongs to cluster  $i$  :- (3)

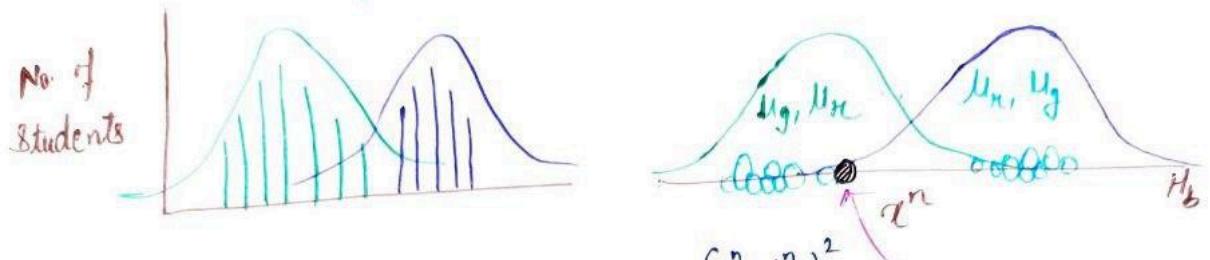
$$P\left(\frac{c_i}{x^n}\right) = \frac{P\left(\frac{x^n}{c_i}\right) \cdot P(c_i)}{P(x^n)} = \frac{\pi_i P\left(\frac{x^n}{c_i}\right)}{\sum_{i=1}^k \left(\pi_i P\left(\frac{x^n}{c_i}\right)\right)}$$

$\pi_i$  = In what proportion this point is mixed with other cluster.

## Lecture-20

### Gaussian Mixture Models

Date: 26-09-2024 ①



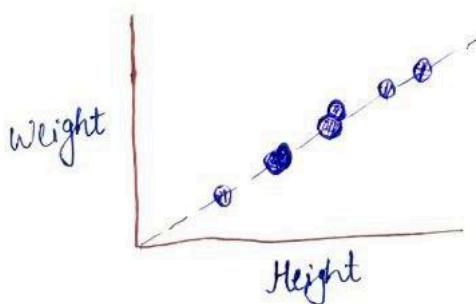
$$p_g(x^n | \mu_g, \sigma_g) = \frac{1}{\sqrt{2\pi} \sigma_g} e^{-\frac{(x^n - \mu_g)^2}{2\sigma_g^2}}$$

$$p_n(x^n | \mu_n, \sigma_n) = \frac{1}{\sqrt{2\pi} \sigma_n} e^{-\frac{(x^n - \mu_n)^2}{2\sigma_n^2}}$$

Consider this student

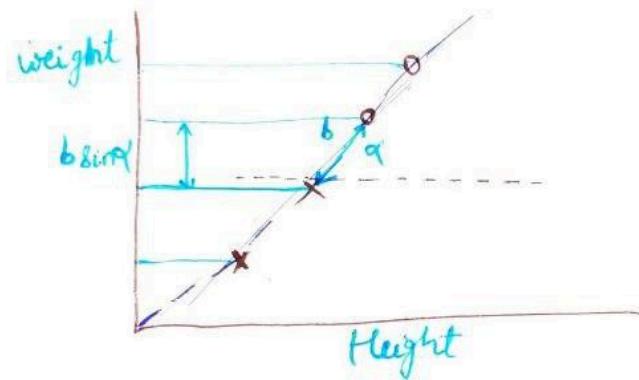
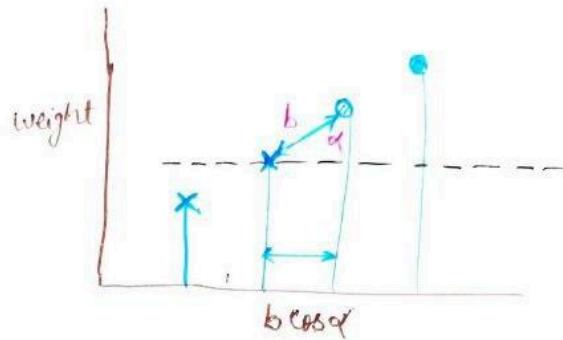
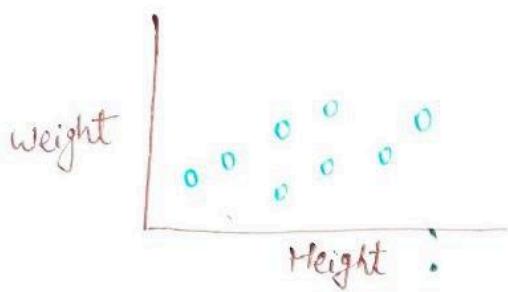
What is belongingness of this student to the Green and Purple Gaussian?

### Classification Problem :-



- I want to project dataset on blue line.

- Variance of data indicate spread of data along blue line.

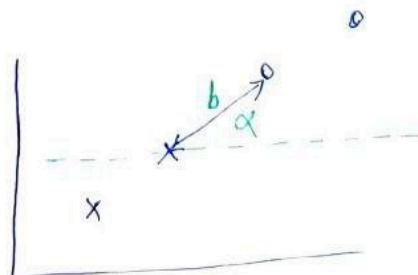
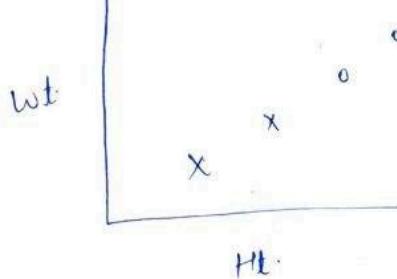


If I try to classify using  
a combination of ht. and wt

→ This gives best decision Boundary

D.B with  $\downarrow$   
highest gap among the samples of different  
classes.

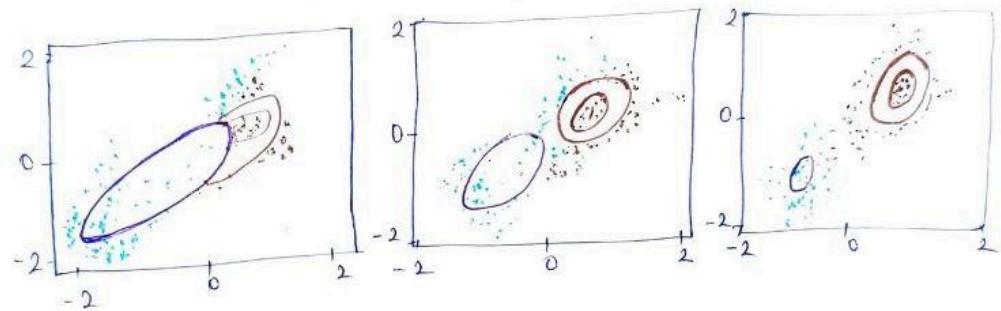
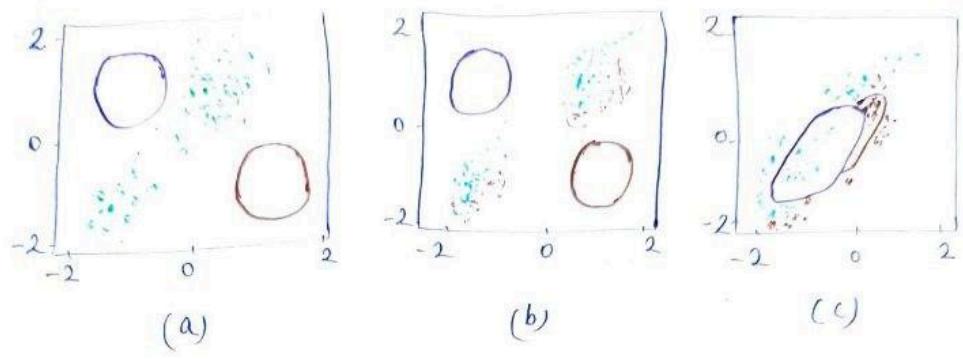
- The diversity of samples are highest along blue (desired) line.



GMM

in  
action

③



### Expectation Minimization :-

- (i) Randomly initialize  $K$  gaussians (randomly initialize  $\mu_i, \Sigma_i, \pi_i$ )  
The  $i^{\text{th}}$  Gaussian is given by  $N(\mu_i, \Sigma_i)$
- (ii) E-step : for each data point, check its belongingness to each of  $K$  Gaussians.
- (iii) M-step : Using belongingness, recalculate  $\mu_i, \Sigma_i, \pi_i$  for each cluster  $i$  to maximize log likelihood.
- (iv) Repeat steps 2-4 unless the termination criteria is classified.

- Soft - Clustering using GMM < Expectation Minimization <sup>(4)</sup>

Calculating  $\mu_i, \Sigma_i$ : It can be shown by taking partial derivative of log-likelihood, that log likelihood will be maximum if:-

$$\mu_i = \frac{\lambda_i^{(1)} x^{(1)} + \lambda_i^{(2)} x^{(2)} + \dots + \lambda_i^{(N)} x^{(N)}}{\lambda_i^1 + \lambda_i^2 + \dots + \lambda_i^N}$$

$$\sum_i = \frac{\lambda_i^{(1)} (x^{(1)} - \mu_i) (x^{(1)} - \mu_i)^T + \lambda_i^{(2)} (x^{(2)} - \mu_i) (x^{(2)} - \mu_i)^T + \dots + \lambda_i^{(N)} (x^{(N)} - \mu_i) (x^{(N)} - \mu_i)^T}{\lambda_i^{(1)} + \lambda_i^{(2)} + \dots + \lambda_i^{(N)}}$$

$$\mu_i = \frac{\sum_{n=1}^N (\lambda_i^n x^n)}{\sum_{n=1}^N \lambda_i^n} \quad \Sigma_i = \frac{\sum_{n=1}^N \left( \lambda_i^n (x^n - \mu_i) (x^n - \mu_i)^T \right)}{\sum_{n=1}^N \lambda_i^n}$$

$$T_i = p(C_i) = \frac{\text{Data points belonging to } C_i}{\text{Total Data Points}} = \frac{\lambda_i^1 + \lambda_i^2 + \dots + \lambda_i^N}{N}$$

$$= \frac{\sum_{n=1}^N \lambda_i^n}{N}$$

5

## Soft Clustering using GMM:

- Suppose we take  $K$  clusters.
- Randomly initialize  $K$  Gaussians. The  $i^{th}$  Gaussian is given by :-  $\mathcal{N}(\mu_i, \Sigma_i)$
- For each data points, check its belongingness to each of  $K$  Gaussians. The belongingness of data point  $x^n$  to  $i^{th}$  Gaussian is :-

$$\lambda_i^{(n)} = \frac{\pi_i \mathcal{N}(x^n | \mu_i, \Sigma_i)}{\sum_{m=1}^K (\pi_m \mathcal{N}(x^n | \mu_m, \Sigma_m))}$$

- Using belongingness, recalculate  $\mu_i, \Sigma_i, \pi_i$  for each cluster  $i$  to maximize log likelihood :-

$$\sum_{j=1}^N \ln \left( \sum_{i=1}^K (\pi_i \mathcal{N}(x^j | \mu_i, \Sigma_i)) \right)$$

$$P(x^n | C_i) = P(x^n | \mu_i, \Sigma_i) = \frac{1}{\sqrt{2\pi} \sigma_i} \exp \left( -\frac{(x^n - \mu_i)^2}{2\sigma_i^2} \right)$$

$$= \mathcal{N}(x^n | \mu_i, \Sigma_i)$$

- Log-likelihood :

(6)

$$\begin{aligned}
 l &= \sum_{j=1}^N \ln P(x^j) = \sum_{j=1}^N \ln \left( \sum_{i=1}^K P(x^j | c_i) \pi_i \right) \\
 &= \sum_{j=1}^N \ln \left( \sum_{i=1}^K P(x^j | c_i) \pi_i \right) \\
 &\quad \downarrow \\
 &\sum_{j=1}^N \ln \left( \sum_{i=1}^K (\pi_i N(x^j | \mu_i, \Sigma_i)) \right)
 \end{aligned}$$

## Soft Clustering using GMM

- Suppose we take K clusters.
- Randomly initialize K Gaussians. The  $i^{th}$  Gaussian is given by :  $N(\mu_i, \Sigma_i)$
- For each data points, check its belongingness to each of K Gaussians.

The belongingness (responsibility) of data point  $x^n$  to the

$i^{th}$  Gaussian is

$$\lambda_i^n = \frac{\pi_i N(x^n | \mu_i, \Sigma_i)}{\sum_{m=1}^K \pi_m N\left(\frac{x^n}{\mu_m, \Sigma_m}\right)}$$

## Belongingness of Data points to Different Clusters:

(7)

Belongingness of  $x^n$  to cluster  $i$  :-

$$\lambda_i^{(n)} = P(C_i | x^n) = \frac{P(x^n | C_i) P(C_i)}{P(x^n)} = \frac{\pi_i P(x^n | C_i)}{\sum_{m=1}^K \pi_m P(x^n | C_m)}$$

$$\lambda_i^{(n)} = P\left(\frac{C_i}{x^n}\right) = \frac{P(x^n | C_i) P(C_i)}{P(x^n)}$$

$$= \frac{\pi_i N(x^n | \mu_i, \Sigma_i)}{\sum_{m=1}^K (\pi_m N(x^n | \mu_m, \Sigma_m))}$$

$$\cdot P(x^n) = P(x^n, c_1) + P(x^n, c_2) + \dots + P(x^n, c_K) = \sum_{i=1}^K P(x^n, c_i)$$

$$\cdot P(x^n) = \sum_{i=1}^K P(x^n, c_i) = \sum_{i=1}^K P\left(\frac{x^n}{c_i}\right) P(c_i) = \sum_{i=1}^K \left(P\left(\frac{x^n}{c_i}\right) \pi_i\right)$$

$$\cdot \pi_i = P(c_i) \rightarrow \text{mixing coefficient}$$

Using Baye's theorem, we calculate probability that  $x^n$  belongs to cluster  $i$  :-

$$P\left(\frac{c_i}{x^n}\right) = \frac{P\left(\frac{x^n}{c_i}\right) P(c_i)}{P(x^n)} = \frac{\pi_i P\left(\frac{x^n}{c_i}\right)}{\sum_{i=1}^K \left(\pi_i P\left(\frac{x^n}{c_i}\right)\right)}$$

$$P\left(\frac{x^n}{c_i}\right) = P(x^n | \mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi} \sigma_i} \exp\left(-\frac{(x^n - \mu_i)^2}{2\sigma_i^2}\right) \quad (8)$$

$$= \mathcal{N}(x^n | \mu_i, \Sigma_i)$$

Belongingness of Data Points to different Clusters :-

- Suppose we take  $K$  clusters  $C_1, C_2, \dots, C_K$
- Consider a data point  $x^n$ .
- $P(x^n) =$  prob. that  $x^n$  is in  $C_1$  or prob. that  $x^n$  is in  $C_2$   
or ... ... ... prob. that  $x^n$  is in  $C_K$ .
- So  $P(x^n) = P(x^n, C_1) + P(x^n, C_2) + \dots + P(x^n, C_K)$   
 $= \sum_{i=1}^K P(x^n, C_i)$

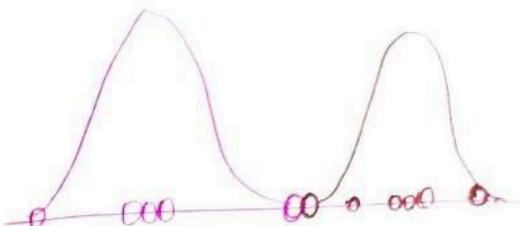
## Soft clustering using GMM

(9)

- we take  $K$  clusters

- Randomly initialize

$K$  gaussians



- The  $i^{th}$  Gaussian is

given by :-  $N(\mu_i, \Sigma_i)$

- For each data points , check its performance belongingness to each of the  $K$  gaussians.

The belongingness of data point  $x^n$  to the  $i^{th}$  gaussian is

$$\lambda_i^n = \frac{\pi_i N(x^n | \mu_i, \Sigma_i)}{\sum_{m=1}^K (\pi_m N(x^n | \mu_m, \Sigma_m))}$$

- Using belongingness , recalculate  $\mu_i, \Sigma_i, \pi_i$  for each cluster  $i$  to maximize log likelihood.

$$\sum_{j=1}^N \ln \left( \sum_{i=1}^K \pi_i N(x^j | \mu_i, \Sigma_i) \right)$$

## Clustering using GMM

(10)

Suppose we take K clusters

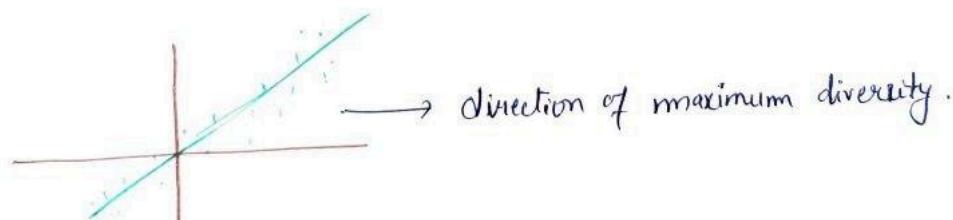
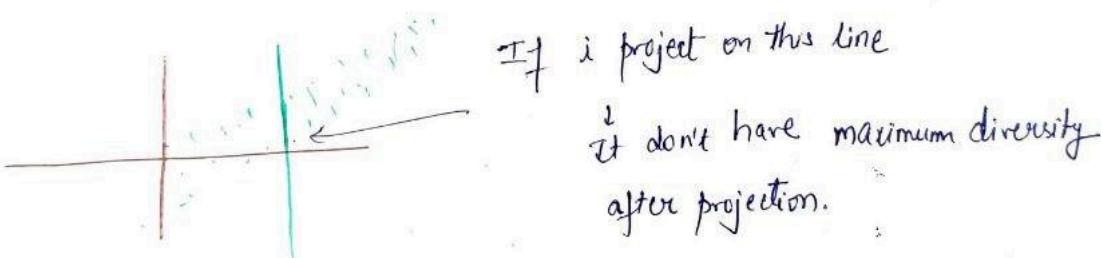
- (i) Randomly initialize K gaussians
- (ii) For all points, check which gaussian it belongs to
- (iii) Based on the belongingness assign each data point to a cluster
- (iv) Find  $\mu_i, \Sigma_i$  from each cluster and update gaussians
- (v) Repeat 2-4 steps unless termination criteria is satisfied.

## Lecture - 21

Date: 27/09/2024

①

- Let's project data on line ( $x$ -axis or  $y$ -axis).
- If projected on line, data has better projection at line.
- I can find direction at which data has maximum possible spread.
- Spread = max. variability for data.
- This projection of data = can do tasks like classification.
- Project data on line to find direction



We calculate covariance matrix by ↗

$$\begin{pmatrix} x_1 & x_2 & x_3 \\ \text{var}x_1 & \text{cov}(x_1, x_2) & \text{cov}(x_1, x_3) \\ \vdots & \vdots & \vdots \\ \text{var}x_2 & & \text{var}x_3 \end{pmatrix}$$

Covariance  
+  
tell spread  
of direction  
of data

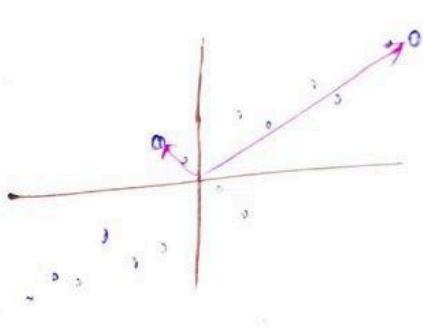
- It is possible to find multiple direction.
- Direction of 2nd max is : 2nd possible Component.

②

Covariance matrix : indicates spread of data points

e.g:- If covariance matrix is :-

$$\begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix}$$



Step 1 :- First take data

Step 2 :- Center it around origin

Step 3 :- Calculate Covariance matrix

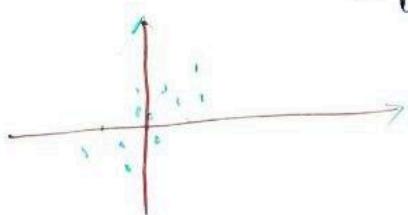
Step 4 :- Calculate eigen vectors of Covariance matrix

Step 5 :- Find corresponding eigen values

Step 6 :- Project all data points into that direction.

Step 2 :- Center it around origin

$$Z_i^{(N)} = \frac{x_i^N - \bar{x}_i}{s_i} \quad \text{standardization}$$



- To fix range when values vary  $\rightarrow$  centering to origin.

②

$\uparrow$

better direction possibility, feature at same scale. (not mandatory)

Step 3:- Calculate covariance matrix :-

Calculation of eigen vector = max. distn of spread of data.

$$\boxed{\lambda \mathbf{z} = \Sigma \mathbf{x}}$$

$$\begin{bmatrix} 9-1 & 4 \\ 4 & 3-1 \end{bmatrix}$$

$\Sigma$  = Covariance matrix

$$(\Sigma - \lambda I) \mathbf{z} = 0 \Rightarrow \text{gives eigen values}$$

$\nearrow$  put  
↓ eigen vector

$$\begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 11 \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$\lambda_1 = 11 \quad \lambda_2 = 1$$

gives 2 vectors

$$\begin{pmatrix} 2 \\ 1 \end{pmatrix} \quad \begin{pmatrix} -1 \\ 2 \end{pmatrix} \text{ direction.}$$

$\Sigma$  = covariance matrix for 2D.

- For d-dimensional data, covariance matrix would be  $d \times d$ .

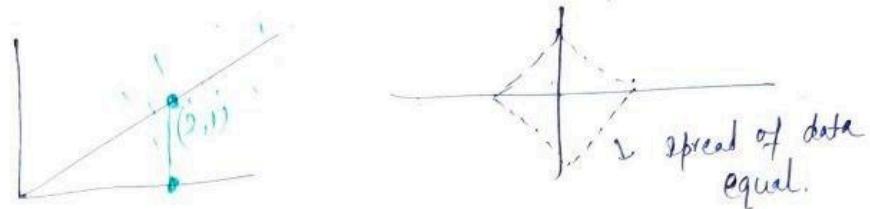
$d \times d = d$  eigen values, vectors

$\downarrow$   
more independent row, column is = d

$$\text{rank} = d$$

(9)

• How to draw line :-



[ Rank is less than 2 when  $\lambda$  are less than 2  
 $\det = 0$  ]

$$\begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \quad \det = 0$$

one of the eigen value vector = 0  
 1 independent column.

Step 4 + calculate eigenvector of covariance matrix

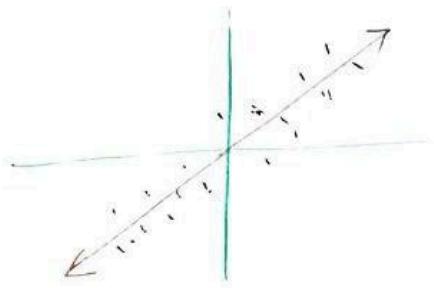
[ Co-linear data points = data in only 1 direction ]

Step 5 :- find corresponding eigen values

Step 6 :- Take dir<sup>n</sup> of maximum eigen value

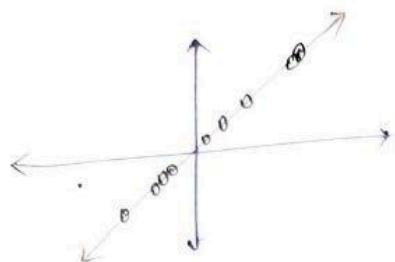
$$\Sigma = \begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix}$$

$\begin{pmatrix} 2 \\ 1 \end{pmatrix}$  eigen vector  
 || eigenvalue (magnitude)



Step 7 Project all points into that direction

(5)

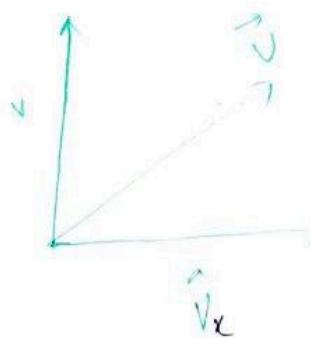


This projected data is called  
first principal Component

- PCA can be used for data reduction.
  - PCA basically gives max direction of spread.
  - 5 dimensional data = 5 max. possible P.C
- ↓  
If we project data in direction of  $\lambda_1, \lambda_2$  :-  
↓ gives 2D plot.  
• we might lose some data.

$5D \rightarrow 2D$  ← Composite features

With PCA Components :- We do our tasks like : compression.



projection still :-  $\sqrt{\cos\theta}$

$$\sqrt{\cos\theta} = \frac{|\vec{v}| |\vec{v}_x| \cos\theta}{|\vec{v}_x|} = \frac{\vec{v} \cdot \vec{v}_x}{|\vec{v}_x|}$$

Depending on data : I'll get P.C  $\Rightarrow$  not its features ④

e.g:- Diff P.C : If we measure height, weight in diff. Country  
e.g+ Germany, India.

- Test data needs to be projected  
 $\uparrow$   
P.C on the direction got from Training data

Composite features = found from 5 features

Eigen stuff

### Principal Components

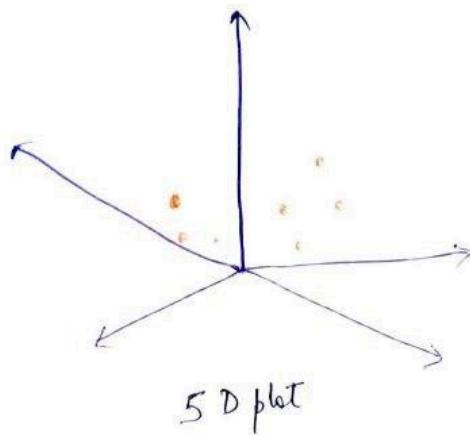
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
1	2	3	4	5
2	3	4	5	6
3	4	5	6	7
4	5	6	7	8

Large Table

$$\begin{pmatrix} \text{Covariance Matrix} \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix} \rightarrow$$

$$\begin{array}{ll} v_1 & \lambda_1 \\ v_2 & \lambda_2 \\ v_3 & \lambda_3 \\ v_4 & \lambda_4 \\ v_5 & \lambda_5 \end{array}$$

Big ↑  
↓ small



If I want to take  
2 pc  $\leftarrow$

①

Eigen stuff

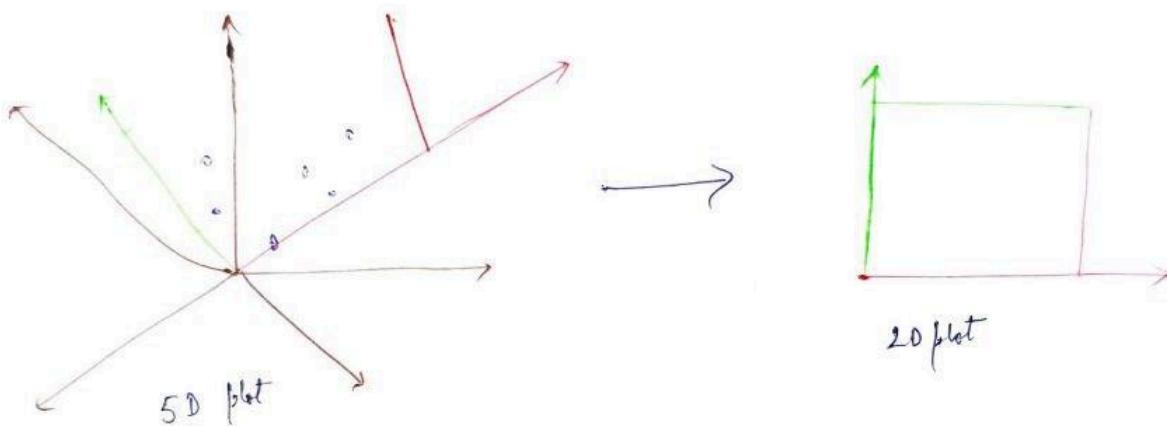
$v_1 \quad \lambda_1$   
 $v_2 \quad \lambda_2$

big  
small

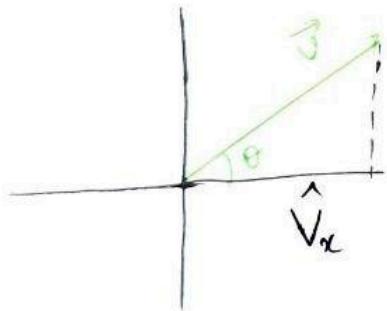
small table

$w_1$	$w_2$
$\vdots$	$\vdots$

Reduced Dimension



$$\vec{v}, \hat{i} = |v| / |\hat{i}| \cos\theta$$

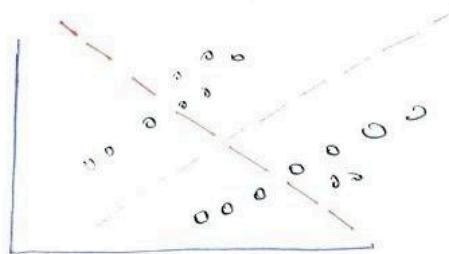


- We'll do stand<sup>n</sup> of train data , also of test data.
- Subtract  $\mu, \sigma$  of train data to test data.  
to get std<sup>n</sup> of test data.
- We'll think test data comes from same dist<sup>n</sup> as of train data.

Date : 27 sept 2024

Lecture : 22

①



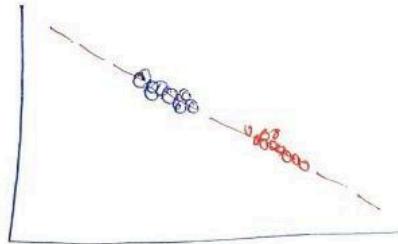
Which line would you like  
to project these data points to?

- A.s per PCA it should be Grey

If these 2 sets belong to different classes.

⇒ Red line will be preferred.

- Projection on red line gives better class separability and dense clusters corresponding to every class.



Goal: Project all points to a direction such that we get best  
projection for classification.

for Labelled data → Classification task

don't only seek for best project, seek:

→ minimum inner cluster variance

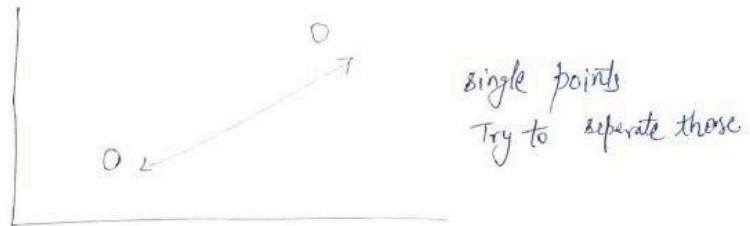
→ max. separability b/w clusters.

- Collapse all blue points of class to one point (2)
- i) collapse all yellow class points to one point.
- ii) [Take these two points as far as possible.]

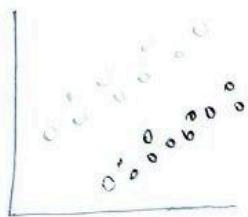
LDA [ $J$  is Linear Discriminant Analysis] It can be used for classification task.

→ we'll try to learn direction.

→  $w$  is unit vector projection.



- $x$  is data point (vector)
- $w$  = unit vector
- If  $x$  is projected on direction of  $w$ .
- Projection is given by  $Z = w^T x$
- Goal: Find  $w$  such that projection of data points towards  $w$  are well separated in terms of their class labels.
- Mean is obtained from vector points, mean is also vector.



If it is 2-class problem with  
classes  $C_1$  and  $C_2$ ,  
mean of classes are:

$$\mu_1 = \frac{1}{N_1} \sum_{x_i \in C_1} x_i, \quad \mu_2 = \frac{1}{N_2} \sum_{x_i \in C_2} x_i$$

$N_1, N_2$ : No. of data points in classes  $C_1, C_2$ .

Mean of classes from the projected data points are :-

$$\mu'_1 = \frac{1}{N_1} \sum_{z_i \in C_1} z_i \quad \mu'_2 = \frac{1}{N_2} \sum_{z_i \in C_2} z_i$$

↓

$$\mu'_1 = \frac{1}{N_1} \sum_{x_i \in C_1} w^T x_i, \quad \mu'_2 = \frac{1}{N_2} \sum_{x_i \in C_2} w^T x_i$$

Goal:- • Maximize distance b/w projected class means

$$\mu'_1 \text{ } \text{and} \text{ } \mu'_2$$

• Minimize within class variance.  
(make denser cluster for each class)

$$U_1' = w^T U_1$$

$$U_2' = w^T U_2$$

$\underbrace{\quad \quad \quad}_{\text{Means after projection}}$

Means after projection

- As points are both +ve and -ve so square of mean needed.

### Covariance of classes before and After projection :

- covariance of classes from original data points are  $\Sigma_1, \Sigma_2$ .
- " " " projected data points  $\Sigma_1', \Sigma_2'$ .

$$\Sigma_1' = w^T \Sigma_1 w \quad \Sigma_2' = w^T \Sigma_2 w$$

$$\text{Distance in 2d} = (\underline{U_1 - U_2})^2$$

As we go for higher dimensions,  $x$  is d dimension so

- Eucledian distance b/w  $U_1'$  and  $U_2'$  is :-

$$\boxed{\beta = w^T (U_1 - U_2) (U_1 - U_2)^T w}$$

$$\beta = (U_1' - U_2')^T (U_1' - U_2')$$

$$\beta = (w^T U_1 - w^T U_2)^T (w^T U_1 - w^T U_2)$$

$$= w^T \underbrace{(U_1 - U_2)(U_1 - U_2)^T}_{{S_B}} w$$

$$\beta = \omega^T (\mu_i - \mu_s) (\mu_i - \mu_s) \omega$$

$$= \mathbf{w}^\top S_B \mathbf{w}$$

$S_B$  = Between class covariance

B: distance b/w

projected mean.

So we have to maximize  $w^T S_B w$

Task 2 :- Minimize within class variance (make denser cluster) for each class

→ Covariance of classes from original data points :  $\Sigma_1, \Sigma_2$

$$\sum_1^1 = w^T \sum_1 w \quad \sum_2^1 = w^T \sum_2 w$$

Minimize :

$$r = \sum_1 + \sum_2$$

$$= \mathbf{w}^T \sum_1 \mathbf{w} + \mathbf{w}^T \sum_2 \mathbf{w}$$

$$= \omega^T (\Sigma_1 + \Sigma_2) \omega = \omega^T S_c \omega$$

$S_c$  = within class variance

$$\text{Maximize} : \quad \beta = w^T S_B w$$

$$\text{Minimize} \quad ; \quad \gamma = w^T S_c w$$

## Achieving Goal 1, 2 +

(6)

Equivalent to maximizing :  $\frac{w^T S_B w}{w^T S_C w}$

- equivalent to maximizing  $w^T S_B w$  such that  $w^T S_C w = 1$  constant

Two ways to maximize :- division, subtract them.

choosed one : division.

Lagrange multiplier is one used for optimization techq.

So numerator is maximized by keeping denominator constant.

- Using lagrange multiplier for optimization problem.  
↳ Solution for  $w$  direction is obtained from eqn

$$S_C^{-1} S_B w = \lambda w$$

- $S_C, S_B$  = matrices  $\lambda$  = scalar
- $w$  = vector

Above eqn can be in form

$$\boxed{A w = \lambda w}$$

$w$  that solve eqn are eigen vector of

$$\boxed{A = S_C^{-1} S_B}$$

$\omega$  = dir<sup>n</sup> where projected data gives max b/w class variance, min within class variance. (7)

How many eigen vector possible if  $A = S_c^{-1} S_B$ ?

It can be shown that only one eigen vector are possible for  $S_c^{-1} S_B$ .

- only 1 eigen vector given direction of maximum spread, minimum covariance.
- one  $\omega$  is obtained we project all data on direction of  $\omega^T$

$$Z = \omega^T X$$

### PCA vs LDA

PCA = unsupervised      LDA : supervised

Linear Transformation ↪

- Transformations (mapping) of points (vectors) from one space to another space.

$$\begin{array}{ccc} & \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} & \begin{pmatrix} 2 & 1 \end{pmatrix} \\ \begin{pmatrix} x \\ y \end{pmatrix} & \xrightarrow{\quad} & \begin{pmatrix} x' \\ y' \end{pmatrix} \end{array}$$

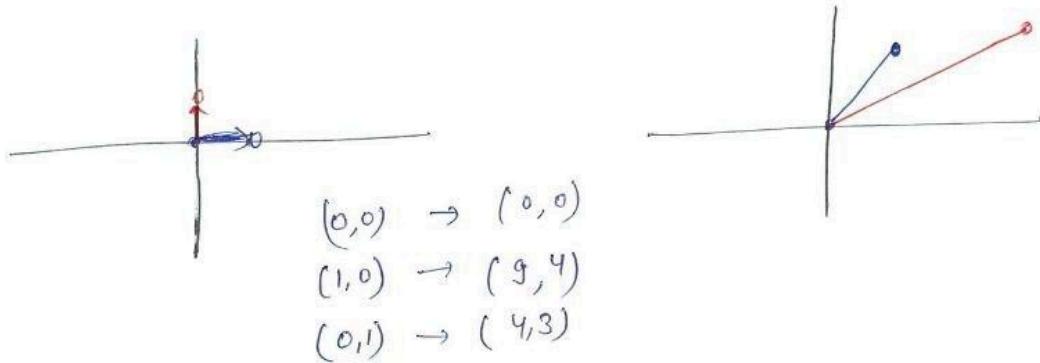
L.T often represented as matrix

(8)

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

$$x' = 9x + 4y$$

$$y' = 4x + 3y$$



For these points direction of vector changes due to linear transformations.

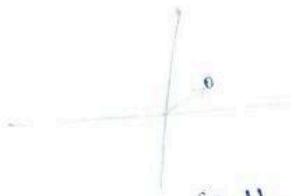
- Linear transformation doesn't cause any direction change.

↳ Direction of eigen vectors of transformation matrix.

→ Corresponding eigenvalue indicates scaling factor.

e.g.  = cylinder rotated around its axis gives  
direction of eigen vector in ↑ vertical direction

- Eigen value : Amount of scaling after Linear Transformation.



$$(2, 1) \longrightarrow (22, 11)$$



$$\begin{aligned}x' &= 9x + 4y \\y' &= 4x + 3y\end{aligned}$$

$\begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix}$  = Transformation matrix

$x$  = point which get linearly transformed.

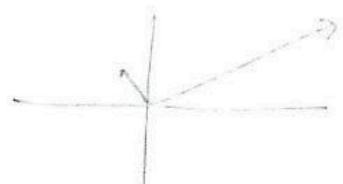
$$\theta = \tan^{-1}\left(\frac{1}{2}\right) \longrightarrow \theta' = \tan^{-1}\frac{11}{22}$$

→ Direction doesn't change due to linear transformation.

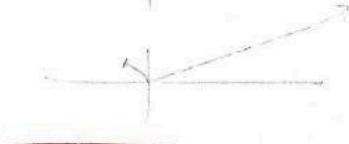


Another set of points, direction doesn't change of this transformation.

$$(-1, 2) \longrightarrow (-1, 2)$$



- For these two directions, linear transformations doesn't change direction.



→ Only scale vectors.

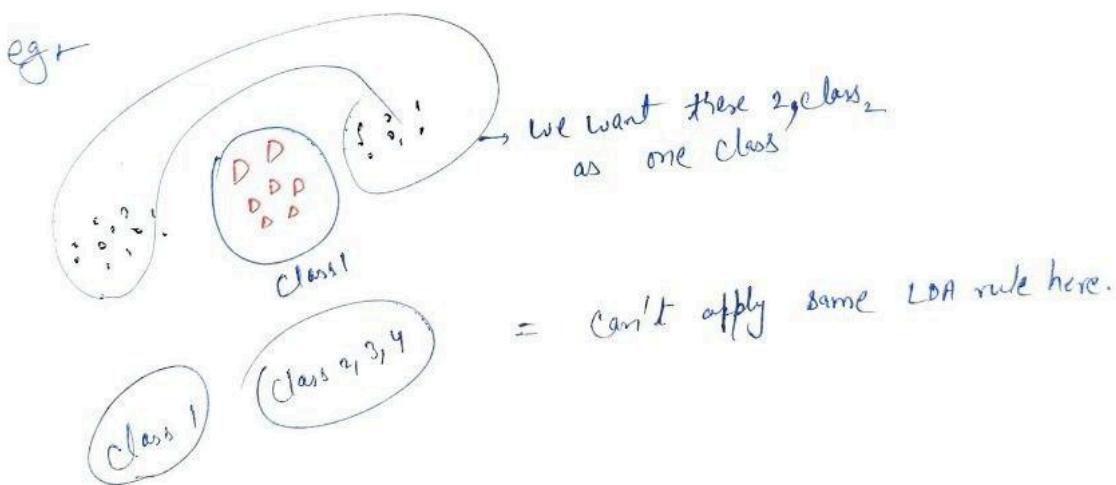
→ Dir<sup>n</sup> of eigen vectors of our transformation matrix.

Q

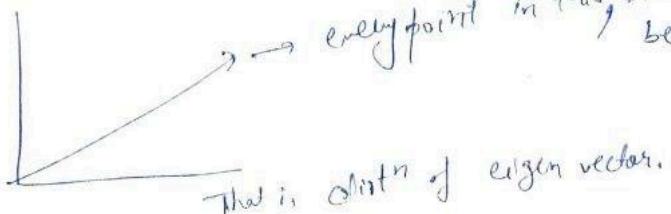
- Corresponding eigenvalues indicates scaling factor

In LDA :-

- Get modified value after projection ↓  
Then train model for it , train for test data
- LDA is applicable for 2 clusters.



- LDA +

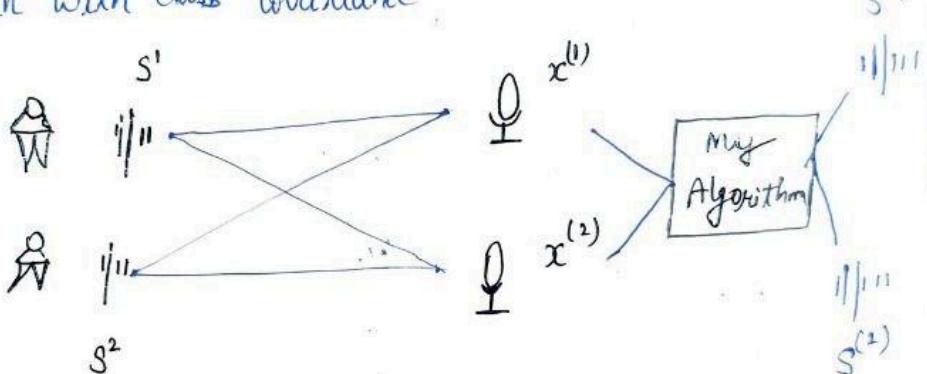


PCA : Unsupervised

LDA : Supervised

- PCA: Finding direction of max variation of data

- LDA: Finding direction of max between class covariance and min with class covariance



## Independent Component Analysis

$$x_1 = a_{11} s_1 + a_{12} s_2$$

$$x_2 = a_{21} s_1 + a_{22} s_2$$

$$\cancel{x = As} \quad x = As$$

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

A = mixing matrix

$\therefore s_i$  is Non-Gaussian

(2)

$$x = As$$

$$s = Wx$$

Let  $p_s(s)$  be pdf of random variable  $s$ .

Let  $p_x(x)$  be pdf of random variable  $x$ .

It can be proved that :-

$$p_x(x) = p_s(s) |W|$$

Signal from each person is independent of others  
 $\uparrow$   
 (random variable)

$s_i$  and  $s_j$  are independent random variables.

$$\begin{aligned} x &= As \\ s &= Wx \end{aligned}$$

$$W = \begin{pmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_N^T \end{pmatrix}$$

So if we have  $N$  persons -

$$\begin{aligned} p_s(s) &= p_s(s_1, s_2, \dots, s_N) = p_s(s_1) p_s(s_2) \dots p_s(s_N) \\ &= \prod_{i=1}^N p_s(s_i) \end{aligned}$$

It can be proved that

$$p_x(x) = p_s(s) |W|$$

If there are  $N$  points :—

(3)

$$P_x(x) = \left( \prod_{i=1}^N P_s(s_i) \right) |w|$$

$$P_x(x) = \left( \prod_{i=1}^N P_s(w_i^T x) \right) |w|$$

$$\begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1M} \\ w_{21} & w_{22} & \cdots & w_{2M} \\ \vdots & \vdots & & \vdots \\ w_{N1} & w_{N2} & & w_{NM} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$$

$$\begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{pmatrix} = \begin{pmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_N^T \end{pmatrix} x$$

$$s_i = w_i^T x$$

Goal To find  $w$  such that  $P_x(x)$  is maximized.

Maximum likelihood estimation.

If we have  $D$  data points  $x^1, x^2, \dots, x^D$  —

Likelihood —

$$l(w) = \ln \left( \prod_{j=1}^D P_x(x^j) \right)$$

$$\text{putting } p_s(x^j) = \left( \prod_{i=1}^N p_s(w_i^T x^j) \right) |w| \quad (4)$$

↓ in

$$l(w) = \sum_{j=1}^J \ln \left( \left( \prod_{i=1}^N p_s(w_i^T x^j) \right) |w| \right)$$

MLE using gradient descent to maximize:

$$l(w) = \sum_{j=1}^J \ln \left( \left( \prod_{i=1}^N p_s(w_i^T x^j) \right) |w| \right)$$

$$\text{where } q'(s) = \frac{d q(s)}{ds} = \hat{p}_s(s)$$

$\alpha$  = learning rate

$$\text{update rule} \leftarrow w = w + \alpha \nabla_w (l(w))$$

Assuming  $p_s(\cdot)$  to be an uniform distribution  $\in [0,1]$  for data point  $j$ , update rule is

$$\nabla_w (l(w)) = \begin{pmatrix} 1 - 2q(w_1^T x^j) \\ \vdots \\ 1 - 2q(w_N^T x^j) \end{pmatrix} (x^j)^T + (w^T)^{-1}$$

(3)

• Through several gradient update steps, we can find  $w$

• Once  $w$  is found,

we can find  $s^{(j)} = w \cdot x^{(j)}$  for  $j^{\text{th}}$  training date.

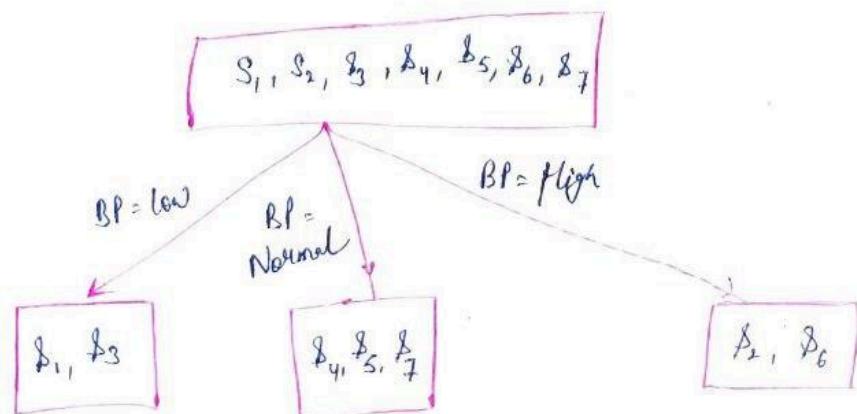
→ So we can find source values  $s_j$  like this for each training measurement  $x_j$

Sample number	Appetite	Weight	BP	Class
S <sub>1</sub>	low	Normal	Low	No Anemia
S <sub>2</sub>	low	Low	High	Anemia
S <sub>3</sub>	Normal	Low	Low	Anemia
S <sub>4</sub>	low	Low	Normal	No Anemia
S <sub>5</sub>	Normal	Low	Normal	Anemia
S <sub>6</sub>	Normal	Normal	High	Anemia
S <sub>7</sub>	Normal	Normal	Normal	No Anemia

Training  
data

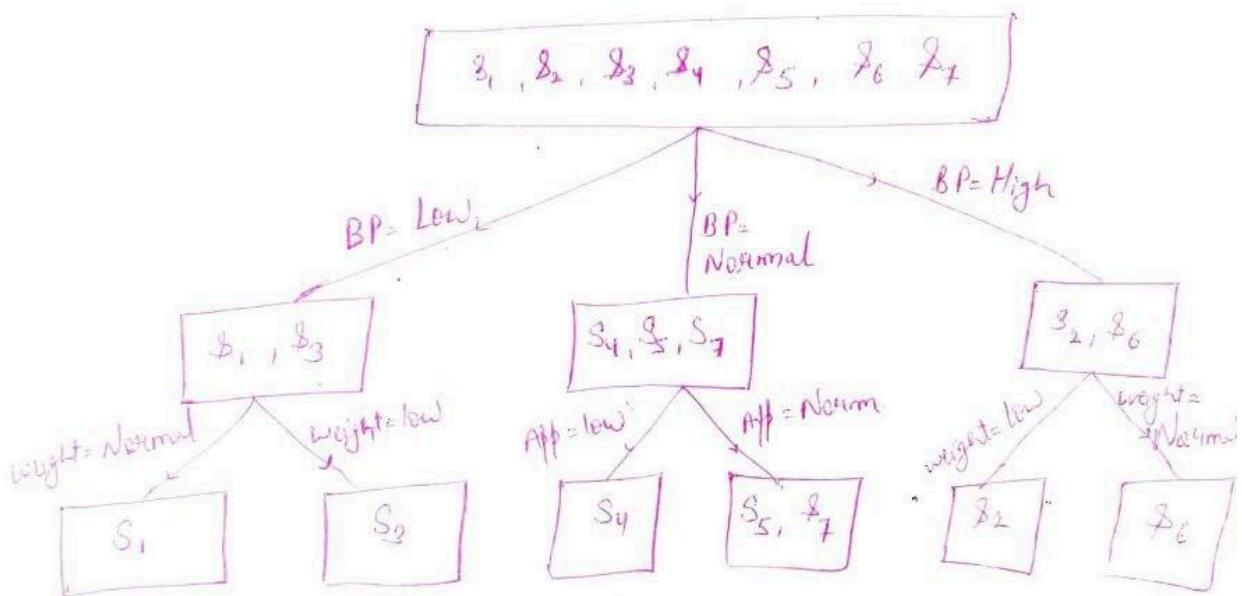
J  
If new sample comes with normal appetite, normal weight,  
low BP. Find if person has Anemia?

- Constructing Decision Tree :- I split the training samples into smaller bags using some criteria based on feature value.



(2) which feature and which value should be chosen?

- splitting : based on some criteria based on feature value.
- Repeat the process in each of smaller bags, may be using some other feature



Q2: How long should we continue splitting?

Goal: For new patient I want to classify patient having anemia or not having Anemia.

Tree satisfied : Training is done.

→ collection of nodes connected by branch = D.T

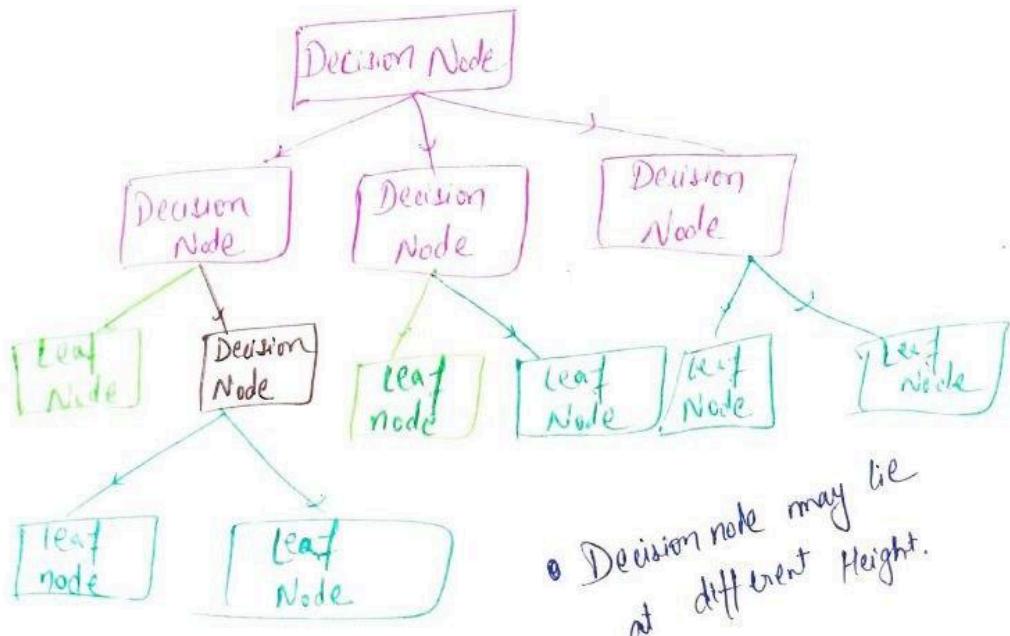
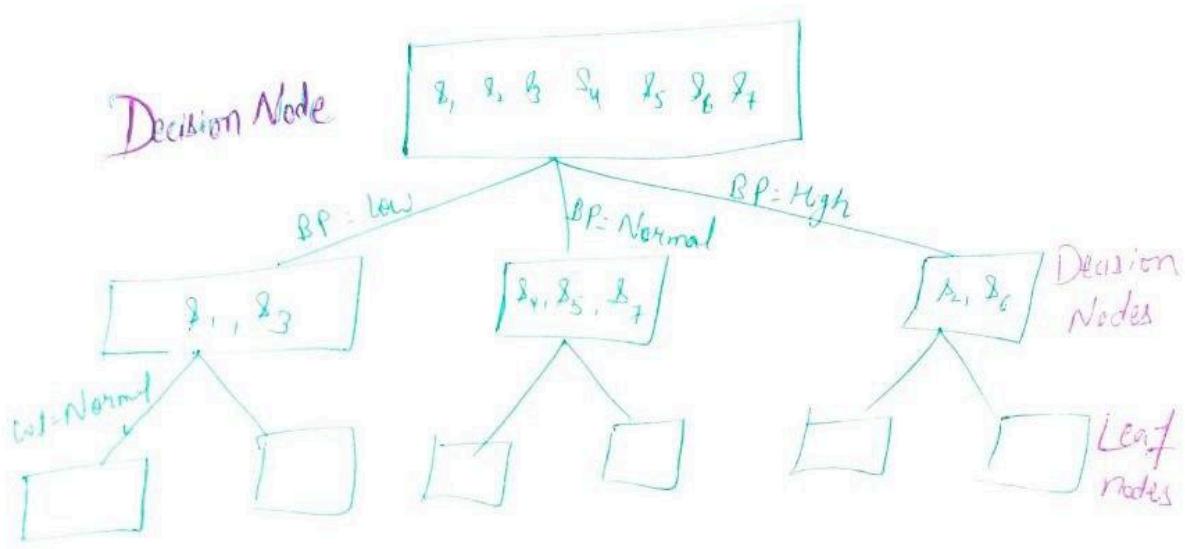
Leaf node may not lie on same height

(3)

Data Reduction is done for :-

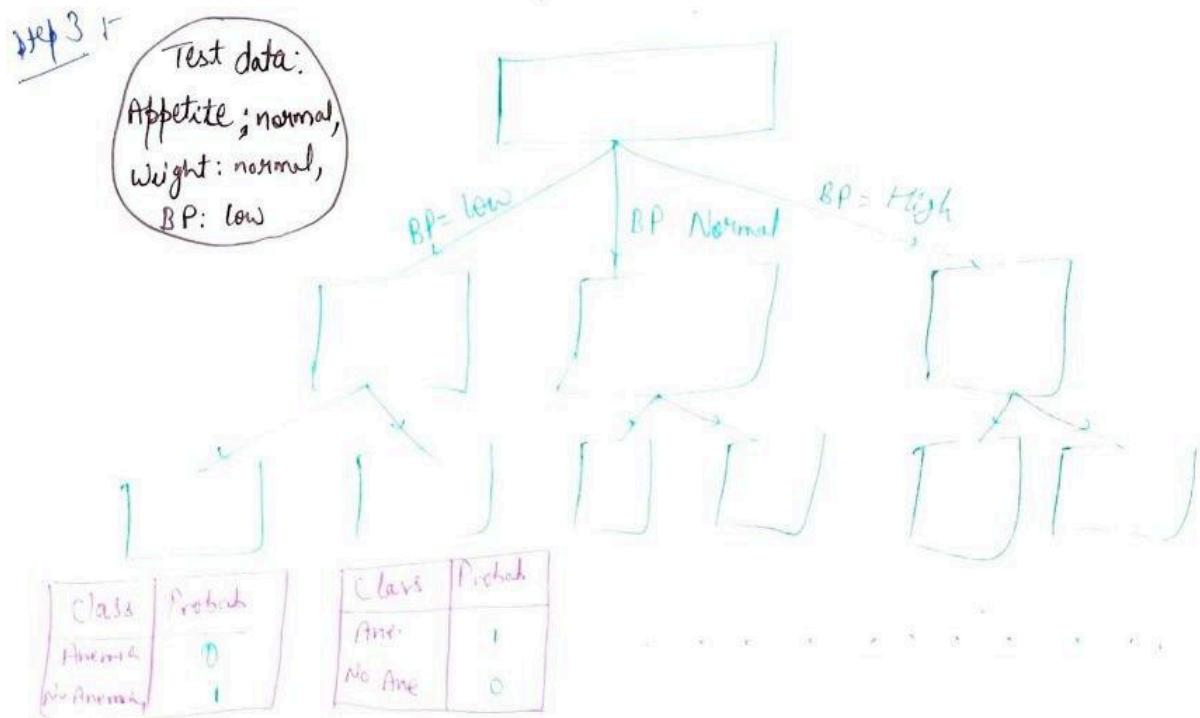
- (i) To compress data , dimensionality reduction
- (ii) D.R also helps to do downstream tasks.

• Leaf node doesn't get splitted further.



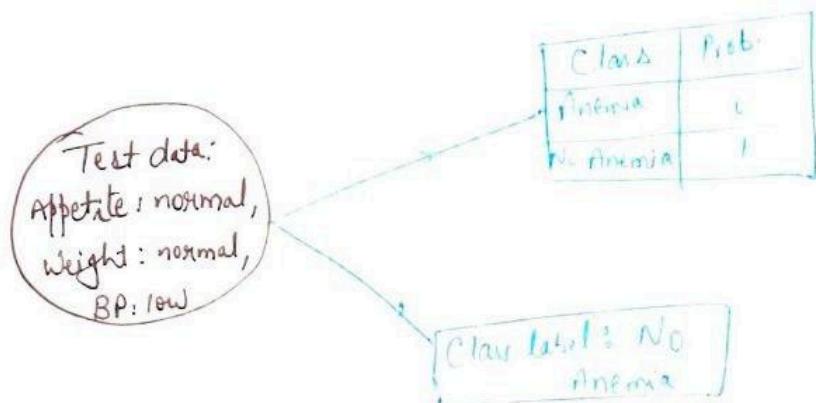
Step 1 :- find out probabilities at each leaf node. ④

Step 2 :- Push down test data through trained tree using the decision rule at each decision node.



Step 4: Check the class probability in leaf node where test data arrives.

## Inference Using a Trained Decision Tree



- for  $s_1$ , probability of Anemia = 0

for  $s_5, s_7$  = 0.5

[ 4<sup>th</sup> node on left can't be assigned Hard Class ]

$\uparrow$   
Soft Class = Probability

### Some Standard Practices :-

- We usually prefer trees with :-

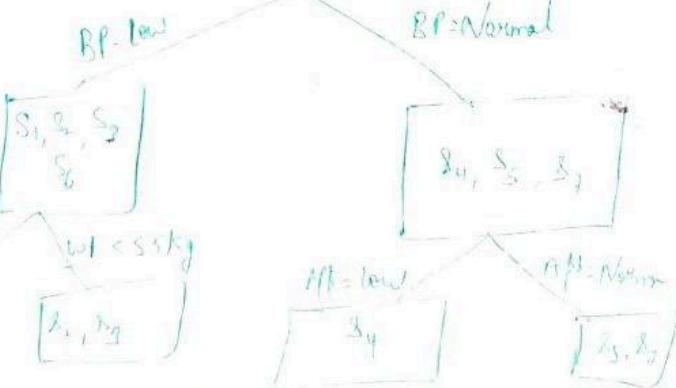
↳ low depth  
 ↳ smaller no of nodes  
 ↳ usually construct binary DT

Inductive Bias of DT

- Each parent node has at most 2 children.

Sample no	Appetite	wt	BP	Class
$s_1$				
:				
$s_7$				

$s_1, s_2, s_3, s_4, s_5, s_6, s_7$



I can either stop

growth of  
tree here

or

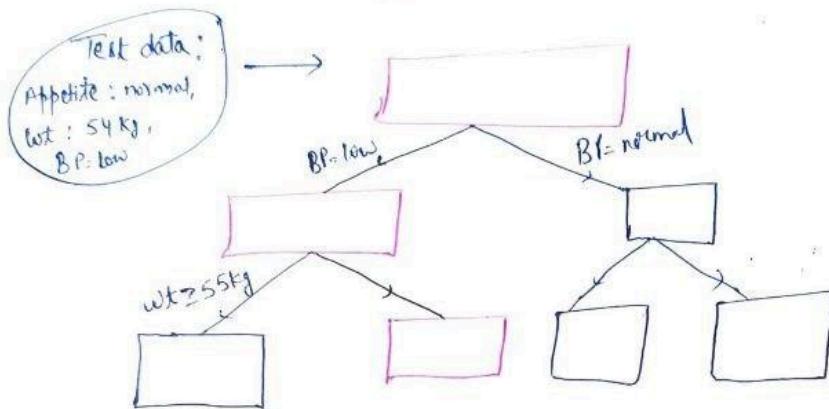
can continue further.

↓  
Calculate class probabilities at  
leaf nodes.

Based on tree we can now test data.

(6)

So if Test data : Appetite : normal,  
Weight : 54 kg, BP : low



- More is depth of Tree : More comparison needed

We want low depth Tree.

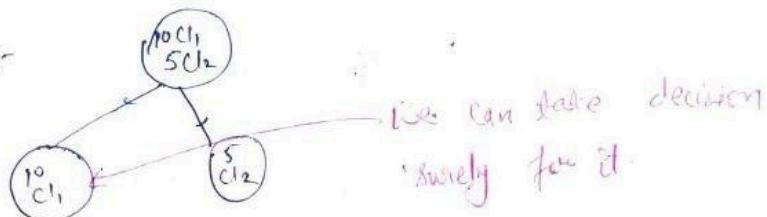
- Inductive Bias of DTF
  - lower depth
  - lower no. of node
  - Binary DT

Binary has criteria of 1 feature

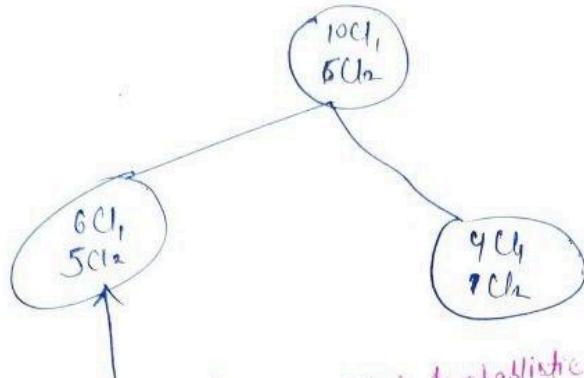
Goal is for Pure Node

Ideally we want child to be of 1 class, while splitting nodes.

e.g :-



(7)

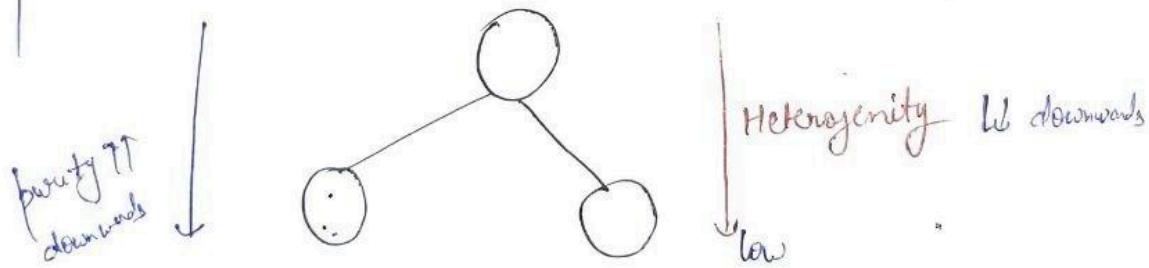


There is confusion, now probabilistic answer

- Ambiguity, slight change in data lead to ambiguous decision for testing.

- In D.T after split we want node to be as Pure as possible.

- Class level at child node should be as Low as possible.



Heterogeneity is fun<sup>n</sup> of class level

- Heterogeneity increases as we move root to child level.
- We should choose feature which is most homogenous.

$$\text{Information} \propto \frac{1}{P}$$

(3)

Entropy : measure of randomness.

→ something obvious is of less interest  $\propto \frac{1}{\text{Probability}}$

→ inversely proportional to probability of it.

If certainty 1 -  $P=1$  info =  $\frac{1}{1} = 1$

$$P = 0.1 \quad \text{Info} = \frac{1}{0.1} = 10 \text{ Info}$$

$$\boxed{P \downarrow \text{Info} \uparrow}$$

→ If  $P_i$  is very low, info increases very high So taking log to mini

event =  $K$

$$I_K = \log \frac{1}{P_K}$$

If random variable =  $X$   
 $\text{Avg} + x P(x) \text{ for discrete} = \sum x P(x)$

Average Info = Entropy

$$= \sum_{K=1}^N P_K I_K$$

(9)

$$\text{Avg Info} = \sum_{k=1}^N p_k \log \frac{1}{p_k}$$

$$\boxed{\text{Avg Info} = - \sum_{k=1}^N p_k \log p_k}$$

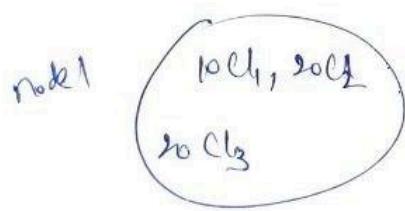
$$p_k = 1 \quad = - \sum p_k \log p_k = 0 = \text{no randomness}$$

Case 2 +  $p_1 = 0.5, p_2 = 0.5$

$$\begin{aligned} & -0.5 \log 0.5 - 0.5 \log 0.5 \\ & = -\log 0.5 = \log 2 \end{aligned}$$

↑↑ Events = more entropy

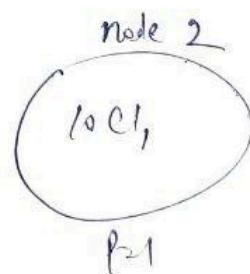
If nodes are like +



$$\text{Probability of Cl}_1 = \frac{10}{50}$$

$$\text{Probability of Cl}_2 = \frac{20}{50}$$

$$\text{Pro. of Cl}_3 = \frac{20}{50}$$



$$\text{entropy} = 0$$

$$\text{Entropy of pure node} = 0$$

more impurity  $\rightarrow$  more entropy

(10)

From Top to Bottom in DT = entropy ↓↓

→ We want to reduce entropy.

We can measure Heterogeneity by +  
Gini Index  
Variance measure.

Entropy  $H$  = Heterogeneity  $H$  = Info loss

Now we have two questions for Decision-tree :-

- (i) which feature to use for splitting?
- (ii) until when to split

- After training we want each leaf node to contain data of only one class.
  - This is called a Completely 'homogeneous' (pure) node.
  - In such a case, confidence in classification is high.
- From root to leaf nodes, we want nodes to become more and more homogeneous.
- For both categorical, continuous type of feature : Same strategy for making tree.

Top most bag = Root node

- For good classification: Leaf node should be as homogeneous as possible.
  - Ideally node should be as pure
  - ↓
    - use some amount of error possible.

↓ child node should become more homogeneous to parent.

Pure node = entropy = 0

- For classification : child is more homogenous than parent when we move downwards.
- Entropy is one of the measure of homogeneity.

$$E = P \ln \frac{1}{P} \quad \text{when } P=1 \quad E = 1 \cdot \ln \frac{1}{1} = 0$$

- Gini Index is also measure of homogeneity.

$$G_I = 1 - \sum_{i=1}^k P_i^2$$

If  $P_j$  is true class for node,  $\perp$

$$G_I = 1 - \sum P_i^2 = 1 - 1 = 0$$

<ul style="list-style-type: none"> <li>Gini is not 0 <math>\Rightarrow</math> Impure node</li> <li>ID3 = Binary DT</li> <li>CART = Any DT</li> </ul>
--

- Goal + In child node entropy is as low as possible
- Calculate entropy based on splitting of all possible feature Binary

D.T :-

$$E = E_1 + E_2$$

③

## Measuring Homogeneity of a Node : Entropy

- Consider a node that contains

- $n_1$  no of data from class 1
- $n_2$  " " " " class 2

- $\vdots$
- $n_K$  no of data from class K.

$$\text{Let } N = n_1 + n_2 + \dots + n_K$$

- In the node, probability

$$\rightarrow \text{of class 1 is } p_1 = \frac{n_1}{N}$$

$$\rightarrow \text{of class 2 is } p_2 = \frac{n_2}{N}$$

$$\rightarrow \text{of class } K \text{ is } p_K = \frac{n_K}{N}$$

### Entropy of Node t

$$E = \sum_{i=1}^K p_i \ln \frac{1}{p_i}$$

If node is pure, data from only one class.

- In pure node, we have data of only one class j.

$$p_j = 1$$

$$p_i = 0 \text{ if } i \neq j.$$

So  $E = 0$   
for pure node, entropy is 0. (4)

As impurity increases, entropy ↑.  
So entropy is measure of impurity.

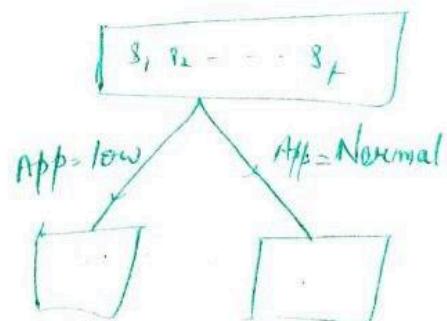
$$G_i = 1 - \sum_{i=1}^k p_i^2$$

$G_i = 0$  for pure node, Gini impurity = 0  
Impurity ↑  $\Rightarrow$  Gini ↑

[After training we want each leaf node to contain data of only 1 class.  
called completely homogeneous (pure) node.  
In such case confidence in classification is high]

Q1:- Which feature and which value should be chosen for splitting a Node?

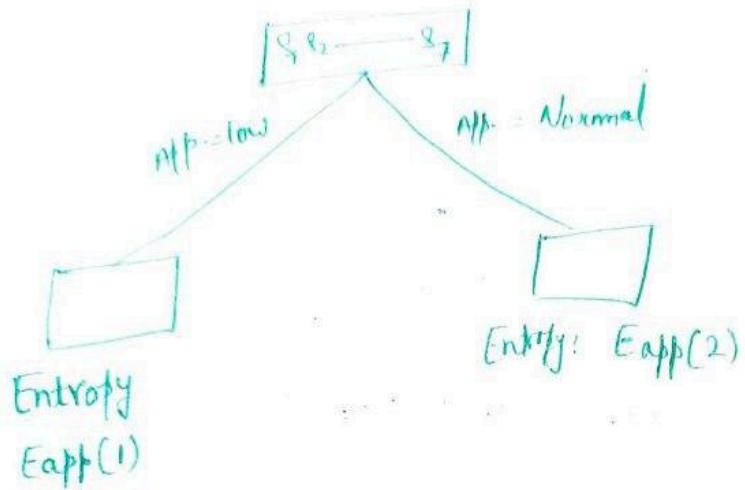
sample no.	App.	ext.	BP	Class
S <sub>1</sub>	:			
S <sub>2</sub>	:			
S <sub>3</sub>	:			
S <sub>4</sub>	:			
S <sub>5</sub>	:			
S <sub>6</sub>	:			
S <sub>7</sub>	:			



first considering feature = Appetite

Possible values: low, Normal

(5)



Hypothetically split the node using these values.

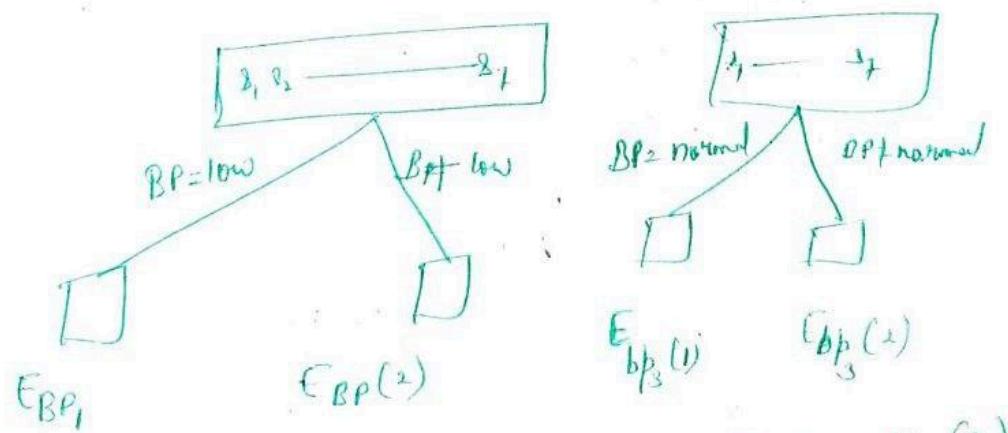
- Calculate total entropy of child nodes

$$\boxed{E_{app} = E_{app}(1) + E_{app}(2)}$$

case 2: Now let's consider next feature BP

Possible values : low, normal, high.

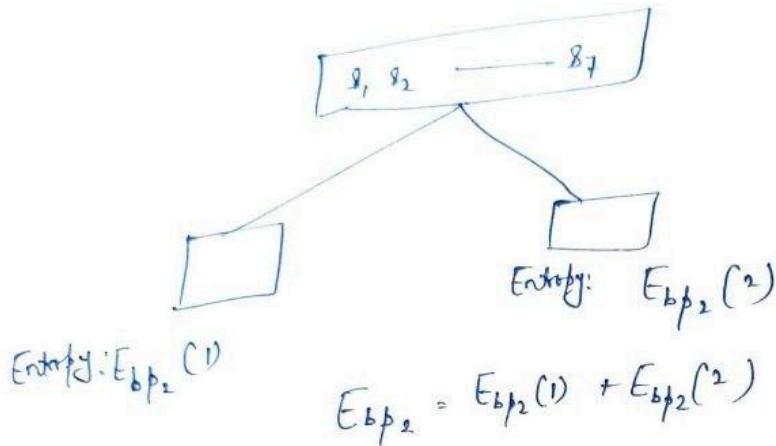
Multiple splits are possible.



$$E_{bp_1} = E_{bp_1}(1) + E_{bp_1}(2)$$

$$E_{bp_3} = E_{bp_3}(1) + E_{bp_3}(2)$$

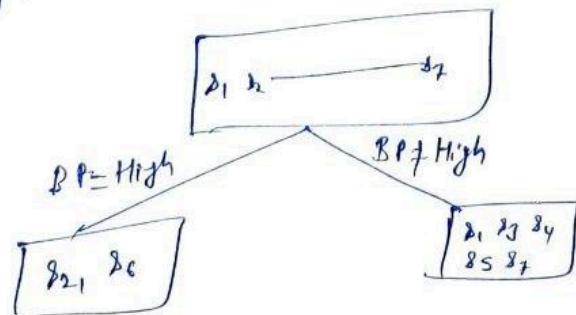
(6)



. we have total entropy of child nodes for various hypothetical

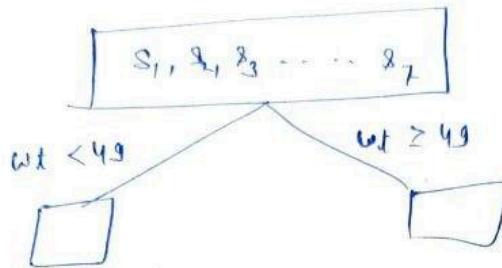
- splits -
- Entropy using feature appetite
- Entropy using weight.
- Entropy using first split point (low vs + low) on feature BP  
similarly  $E_{bp_1}$ ,  $E_{bp_3}$

- find out hypothetical split that leads to smallest total entropy of child nodes.
- suppose  $E_{bp_2}$  is smallest among the above.
- we'll split the node using first split point (High vs + High)  
on feature BP



What happens for Continuous feature?

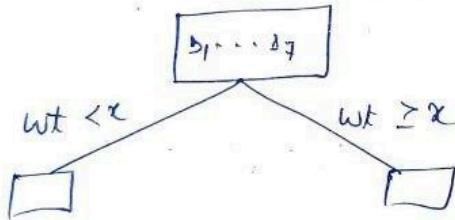
(7)



Total entropy of  
child nodes =  $E_{49}$

For hypothetical split using weight, try hypothetical split with every possible value of weight from dataset.

Possible values of wt = 47, 49, 50, 51, 58, 61, 63



And so on...

- Compare these total entropies with total entropies from other features to find out best split point.

Advantage of DT :-

- Calculations are highly parallelized.
- Less time consumption.

Disadvantage :- Computationally intensive.

- for continuous feature :- Split for all possible weights.
- more split

In some cases DT can be divided into internal based bins, (8)

e.g.  $49 \leq 51 [51 \leftarrow 63] \leq 64 \leq 70$

Binning

- For floating point too approach is same.
- How long should we continue splitting:

There may be diff. criteria:-

- split a node until entropy of node  $< \epsilon$  (a pre-defined value)  
threshold to stop.
- split a node until node contains  $< n$  no of data points  
e.g. if 10 data left after splitting
- split a node until height of tree reaches h  
(a pre-defined value)

Drawback of D.T. A single tree may not be a good classifier  
or regressor.

- If 1000 features, all features might not be used for splitting, not looked into DT splitting for all features.

While doing Regression = we are doing Pattern Identification.

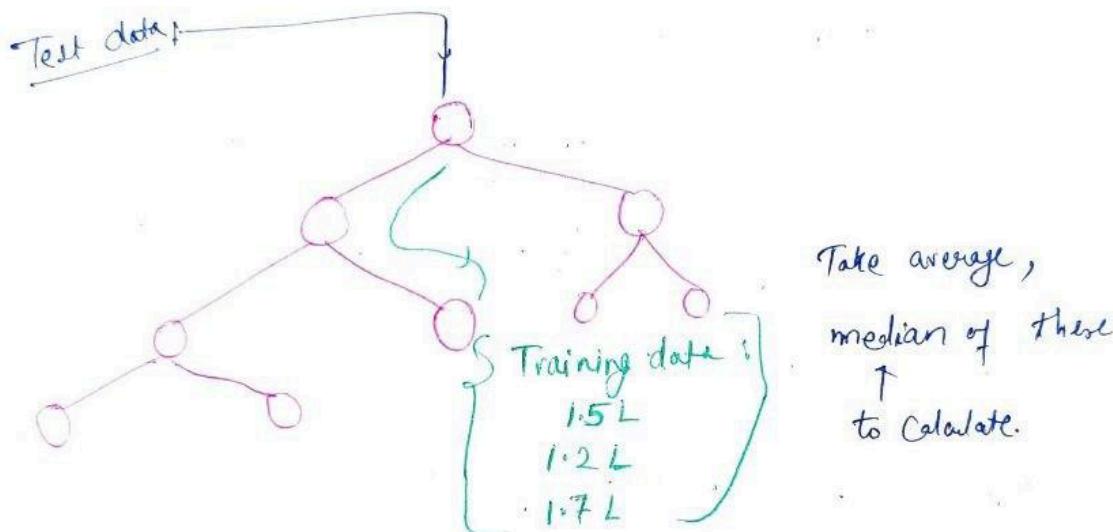
While doing Classification: judging data

②

If node is pure  $\Rightarrow$  variance should be low.

If node is impure  $\Rightarrow$  variance  $\uparrow$

I want to find salary of employ,



$$\text{Info gain} = \text{Entropy of parent} - \sum \text{Total entropy of children}$$

Should be maximized  $\Rightarrow$  less entropy at child should be low.

- A tree may not exploit all imp features of data.
- splitting of nodes may be done using unimportant features.
  - Poor quality of tree

⑩

### Tree to forest +

- Take several D.T construct Forest.
  - voting = taking all possible classification
- Perform voting from ensemble (forest) to make decision about specific classification (or regression)

 In order to get good splitting :- Need to have Diversity.

- We want to grow set of DT diff from each other.  
else biasness may be there.

#### (a) Diversity from perspective of data +

- All the trees must not be grown with same data.
- Each tree should look at diff. training data.

less bias of each tree towards entire training data.

#### (b) Diversity from perspective of features +

- Diff nodes should look at diff features space of data and then split.

- Allowing all features to compete to split nodes.

Lecture - 26 , Date : 18/10/2024

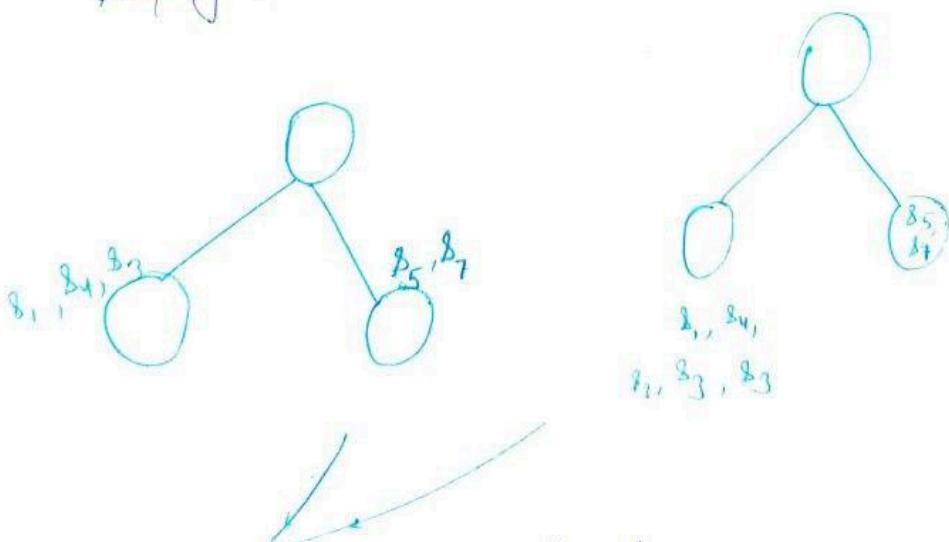
①

- We want heterogeneity as less as possible.
- All tree may grow parallelly.
- If we sample from 100 data items as 10  
↓  
data losses.
- Randomly choosed 70 data = Tree 1  
Randomly choosed 70 data = Tree 2.

Overfitting = less data

Randomly choosing 100 sample with replacement : Bootstrapping

- what if same data is replaced  $\alpha$  no. of time during sampling :-



entropy for both trees are different.  
But info doesn't change

- Feature selection will be affected depending on repeating data or replacement. (2)
- Suppose, there are  $N$  training data points and we want to grow  $T$  no of trees:-
  - Take a random subset of  $N_s$  training data to grow Tree 1
  - Take another random subset  $N_s$  training data to grow Tree 2
  - Take another random subset  $N_s$  training data, to grow Tree  $T$ .

### Diversity from the Prospective of Data

- $N$  training data points and we want to grow  $T$  no of trees.
- Training data for growing diff trees would be different:-
 

$\downarrow$        $\downarrow$

if  $N_s \ll N$  : each tree will be grown with  
small training data,  
more chance of overfitting.

eg If 10 participants, 1 singer is excellent in soul music, else are pop, hip hop etc singer.

- But all should get chance
- every time 1 singer wins : soul music  $\rightarrow$  So

- Dividing data features into small-small pairs then concluding by comparison. (3)
- Reducing feature into smaller parts : More Quality come out.

Bootstrap Sampling : suppose we want to grow T trees with N training data ( $N=8$ ).

Sample no.	$f_1$	$f_2$	$f_3$	Class
$s_1$				
$s_2$				
.				
$s_8$				

- . For growing tree 1, randomly choose N samples with replacement.
- . Let samples :-  $s_5, s_2, s_3, s_6, s_3, s_1, s_2, s_3$

$$D_1 = \{s_5, s_3, s_2, s_3, s_6, s_3, s_1, s_2\}$$

$D_1$  = called bootstrap sample for growing tree 1

Similarly, for Tree 2, bootstrap sample &

$$D_2 = \{s_2, s_1, s_2, s_4, s_4, s_1, s_5, s_8\}$$

Similarly for T trees, we have bootstrap samples

(4)

$$\cdot D_1, D_2, \dots D_T$$

## ② Diversity from Perspective of Features :-

- During growth of forest, trees should exhibit as many features as possible.
- In a D.T., every feature has to compete with every other feature to be chosen for splitting.

$$\begin{cases} \text{if } m = \text{no of feature to choose} \\ \rightarrow n = \text{no of trees} \\ m = \sqrt{n} \quad m = \sqrt{10} = 3 \end{cases}$$

- Quasi-Random : Bootstrapping  
splitting is random

→ Choosing feature :-

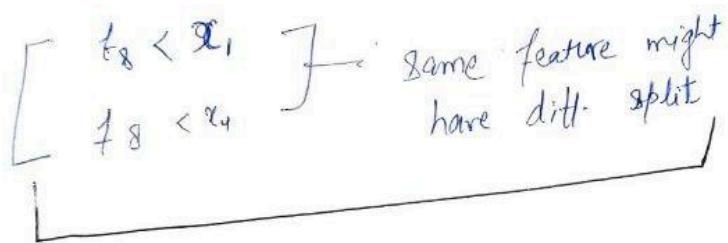
- This becomes tough competition
- If there are a few dominant features, several feature may never get selected for splitting a node.
- such feature would not be utilized in decision making.

→ We want a node splitting process where every feature gets a reasonable chance to be selected for splitting.

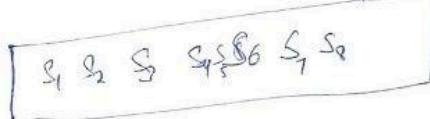
- Random forest mostly used in medical, Bioscience areas.

(3)

- All trees can be grown in parallel.
- Split point different in all trees.
- Particular feature may be best than another.



### Splitting a Node



- Suppose there are features  $f_1, f_2, \dots, f_n$  in my data.

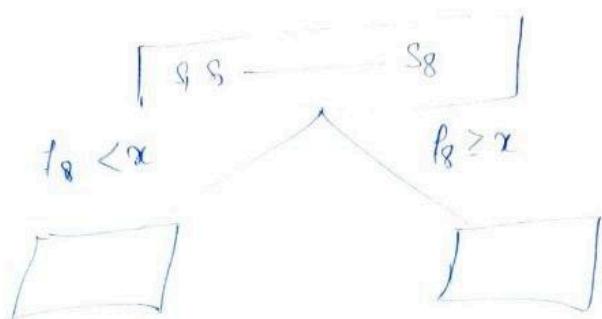
- For splitting node, let randomly choose m no of features (without replacement) out of  $f_1, f_2, \dots, f_n$ .

- Suppose  $m=3$  and I choose  $f_3, f_7, f_8$ .

- Find out best feature out of only  $f_3, f_7, f_8$  using your criteria.  
(Entropy, gini, impurity etc).

- Let split point be  $x$  on feature  $f_8$ .

(6)



Each node effectively look into a subset of features (a feature subspace)

→ Based on subspace, a split point is chosen.

competition among feature reduced.

→ This randomness makes tree diverse.

[ This tree is called Random Decision Tree.

    Growth is continued until a node meets termination criteria.]

## Random Forest

• Decide :-

    → How many trees ( $T$ ) are to be grown

    → How many features one to be selected randomly for splitting a node ( $m$ )

    → Node splitting Criteria (Entropy, Gini Impurity, variance etc)

    → Termination criteria for splitting a node.

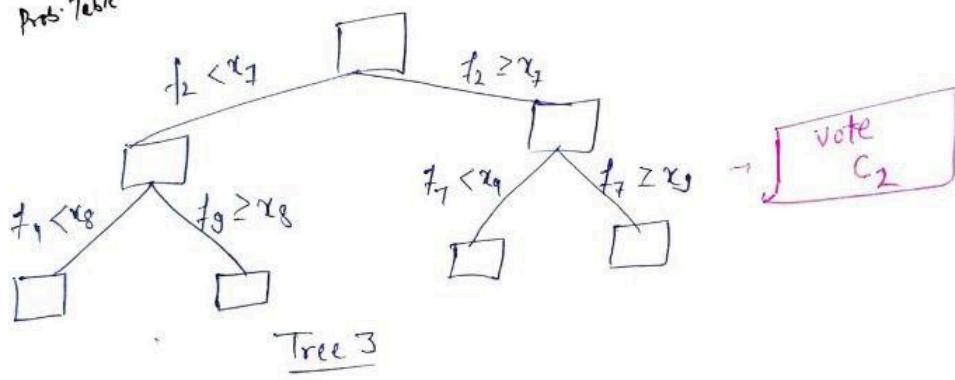
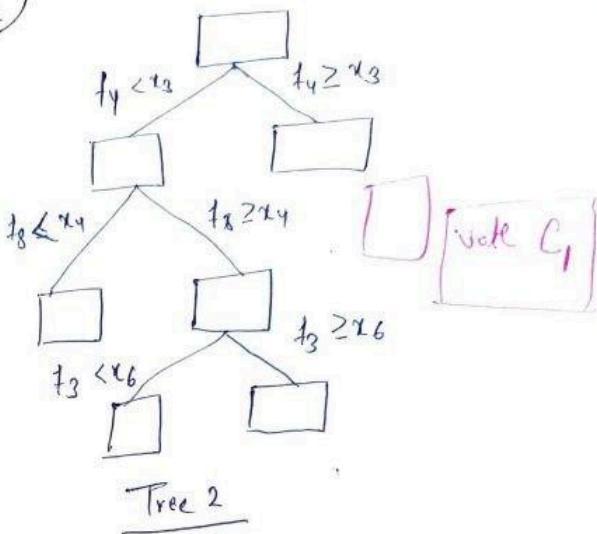
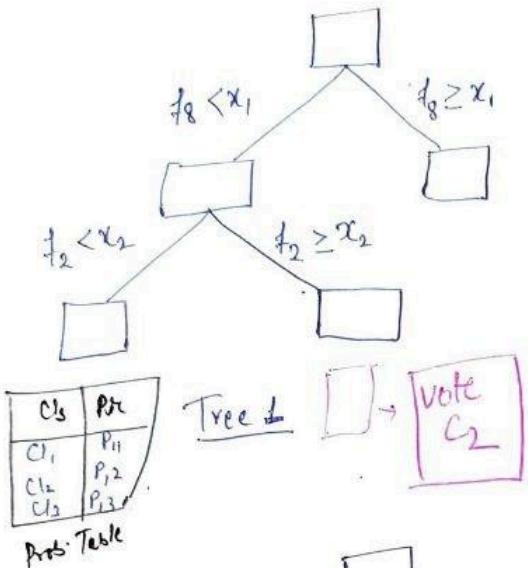
• For  $i = 1:T$

    → Create bootstrap sample for Tree  $i$ .

    → Grow Tree  $i$  using technique mentioned before

    → Calculate class probabilities at each leaf node.

## Random Forest



- Calculate class probabilities at leaf nodes
  - Push test data through all trees
  - At each tree, test data will reach a particular <sup>leaf</sup> nodes.
    - ↓ find the class with highest probability at these nodes.

Class with highest vote is the class for test data.

Class	vote
C <sub>1</sub>	1
C <sub>2</sub>	2
C <sub>3</sub>	0

→ Class of test date is  
Class 2.

At each tree, test data will reach a particular leaf node. (8)  
find class probabilities at those nodes.

Considering all trees, find avg. class probabilities +

Class
$C_1 = \frac{P_{11} + P_{71} + P_{91}}{3}$
$C_2 = \frac{P_{22} + P_{72} + P_{92}}{3}$
$C_3 = \frac{P_{23} + P_{73} + P_{93}}{3}$

→ Class with largest avg. class probability is assigned to test data.  
Class of test data is  $C_3$ .

Random Forest :- Bagging

- Each tree in forest is called a 'weak learner'.
- Each tree is grown independently.
- It can be shown statistically about 66% of training samples will be included in growing trees.
- Other samples are called out of bag sample.

Performance of random forest depends on 2 quantities +

- Strength : indicates the abilities of trees for classification/regression
- Correlation + Similarity among trees, lower better.

R.F is one of the bagging technique.

- Not learning all data some part of data = weak learner.

⑨

. 34% out of bag sample.

• 66% sample used.

→ Error rate in R.F converges.

• Variable importance : Importance of features

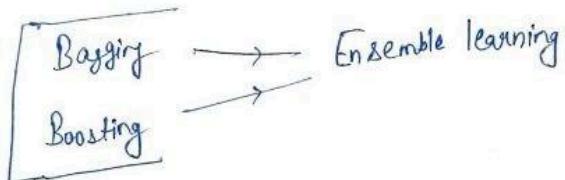
Lecture : 27 , 22/10/2024

①

## Boosting

RF was <sup>not</sup> for all nodes.

- Each tree is creating a hypothesis.
- Each hypothesis gives decision.



- RF + We put equal weightage to each sample
  - not all training data gonna correctly classify.  
↓  
There is impure data : incorrect result.

In Boosting + we put more weightage to incorrect classified data  
DT grow serially in Boosting then we check which is correctly classified, which is not!

- RF all trees growing parallelly.
- In Boosting growing serially → Based on that weightage of P.

↓  
here Stump are independent to each other

↓  
Stumps can grow parallelly.

- Boosting can't grow parallel as whole.

## Boosting +

(2)

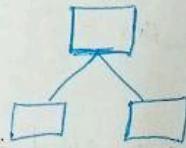
Sample no.	$f_1$	$f_2$	$\dots$	$f_{10}$	Class
$s_1$					$c_2$
$s_2$					$c_1$
$\vdots$					$\vdots$
$s_8$					$c_1$

we have 8 training samples.

→ Assigning equal weights to each sample initially.

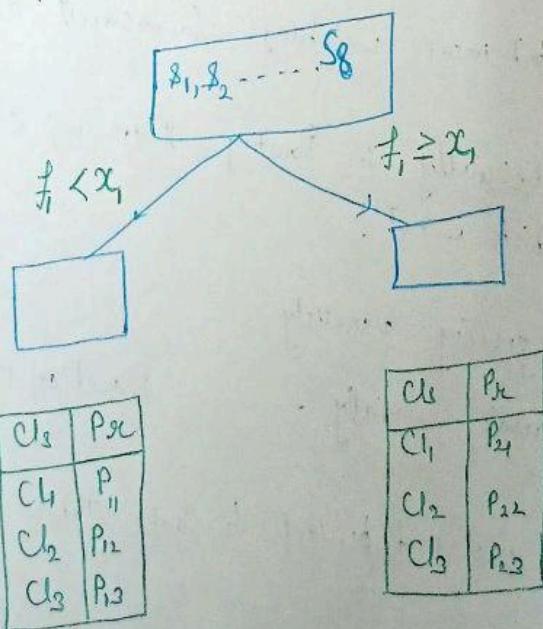
→ with these data we grow stump.

→ Stump is kind of decision tree.  
(with a root, 2 leaves)



Step 1: find out best split point of feature  $f_i$ .

Let split point be  $x_i$ ,



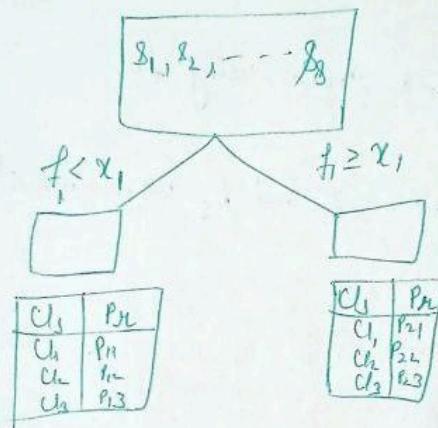
Step 2: Do hypothetical split.

Calculate class probabilities at leaf nodes

Step 3 + Push all training data through the Stump and see how many training data is misclassified. ③

Step 4 + Calculate Gini Index / Entropy of Stump.

Let it be  $g_1$ .



Sample no	$f_1$	$f_2$	...	$f_{10}$	Class	Weight
$S_1$					$C_2$	$\frac{1}{8}$
$S_2$					$C_1$	$\frac{1}{8}$
$S_3$					$C_1$	$\frac{1}{8}$
$S_4$					$C_2$	$\frac{1}{8}$
$S_5$					$C_3$	$\frac{1}{8}$
$S_6$					$C_1$	$\frac{1}{8}$
$S_7$					$C_1$	$\frac{1}{8}$
$S_8$					$C_1$	$\frac{1}{8}$

Step 5 + Calculate Gini Index of stump. Let  $g_2$ .

↑ similarly find best split point of feature  $f_2$ .  
Let split point be  $x_2$

→ do hypothesis split, calculate class probabilities at leaf nodes,  
push training data down, find how many misclassified?

Step 6: So on until  $f_{10}$ .

①

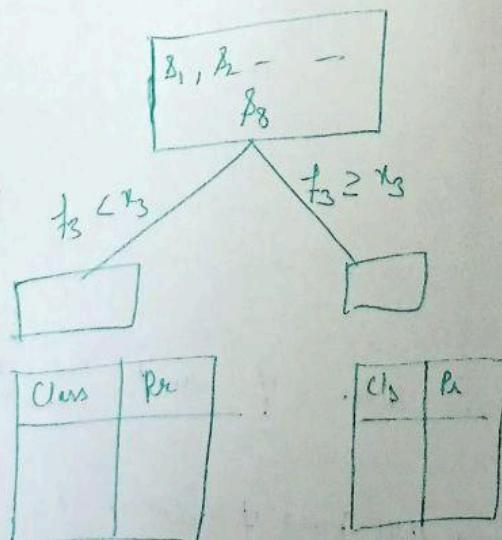
Gini Indices of ten stumps (grown using  $f_1, f_2 \dots f_{10}$ ) are  $g_1, g_2 \dots g_{10}$  respectively.

- Find out stump with lowest Gini Index.

Let it be stump grown using  $f_3$ .

- Assume training samples  $s_3$  and  $s_5$  were misclassified by first Stump

sample no	$f_1$	$f_2$	$\dots f_{10}$	Class	weight
$s_1$					
$s_2$					
$s_3$				$c_1$	$\frac{1}{8}$
$s_4$				$c_2$	$\frac{1}{8}$
$s_5$					
$\vdots \vdots \vdots$					
$s_{10}$					



Total error for Stump = sum of weights of misclassified samples =  $\frac{1}{4}$

$$\frac{1}{8} + \frac{1}{8} = \frac{1}{4}$$

(5)

- Low total error  $\rightarrow$  more imp. Stump.

- Importance of a stump =  $\frac{1}{2} \log \left( \frac{1 - \text{total error}}{\text{total error}} \right)$

$$\text{Importance } II = \text{error } LL$$

$$\text{Imp. } LL = \text{error } II$$

- Incorrectly classified samples by this stump should be treated with more emphasis so that they can be correctly classified by future stumps.  
Boosting.

- Before growing next stump :—
  - $\rightarrow$  increase the weights of incorrectly classified samples
  - $\rightarrow$  decrease the weights of correctly classified samples.

- updated wt. of incorrectly classified samples  $\rightarrow$

$$W_{\text{new}} = W_{\text{old}} * \exp(\text{importance of stump})$$

- updated wt. of correctly classified samples  $\rightarrow$

$$W_{\text{new}} = W_{\text{old}} * \exp(-\text{importance of stump})$$

⑥

If new value is more than 1, we need to normalise

↓  
so divide by sum of all new wts.

→ grow next stump using new wt.

we'll not calculate normal gini index = now weighted gini index

$$\text{earlier} = g = 1 - \sum_{i=1}^N p_i^2$$

$$\text{weighted gini} = 1 - \sum_{i=1}^N w_i p_i^2$$

put more emphasis on incorrectly classified data.

$$\text{entropy (old)} = e = p_i \log \frac{1}{p_i}$$

$$\text{entropy (new)} = e' = w_i p_i \log \frac{1}{w_i p_i}$$

$p_i$  = probability of particular class.

$$\Rightarrow \frac{\text{sum of wt belonging to class}}{\text{sum of wt of all class.}}$$

Step 4 Normalize updated weights

eg:- Normalizing wt weighted more for wrongly classified. (7)

	old	new	<u>Normalize</u> =
$s_1$	$\frac{1}{8}$	$0.17$	$0.17 / \sum \text{all new wts}$
$s_2$	$\frac{1}{8}$	$0.17$	$0.17 / (\sum \text{sum of all wts})$
$s_3$	$\frac{1}{8}$	$0.37$	$0.37 / (\sum \text{sum of all weights})$
$\vdots$	$\vdots$	$0.37$	$\vdots$
$s_8$	$\frac{1}{8}$	$0.17$	

Criterions of Growing boosting :-

- (i) till 50 stump
- (ii) Total error < threshold

[ while growing new stump using ~~same~~ same procedure as that of previous stump ]

[ while growing  $10^{\text{th}}$  stump
 

- we will cover  $9^{\text{th}}$  stump
- we can also consider all earlier 9 stump.

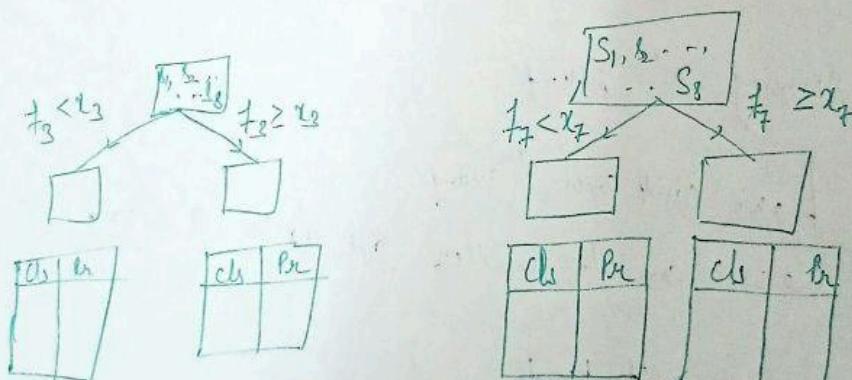
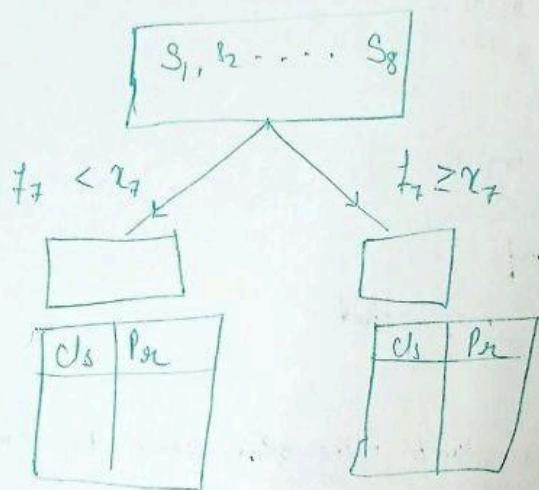
]

finding best split point  
 ↓  
 use weighted Gini importance

$$g = 1 - \sum_{i=1}^N w_i p_i^2$$

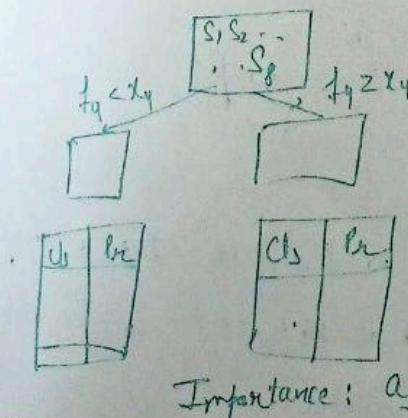
. using best feature and split point, grow the stump :-  
 Let feature be  $f_7$  and split point to be  $x_7$

- Recalculate weight of samples and normalize. (8)
- Grow next stump with normalized updated weights.
- Continue growing more stumps until termination criteria is reached.



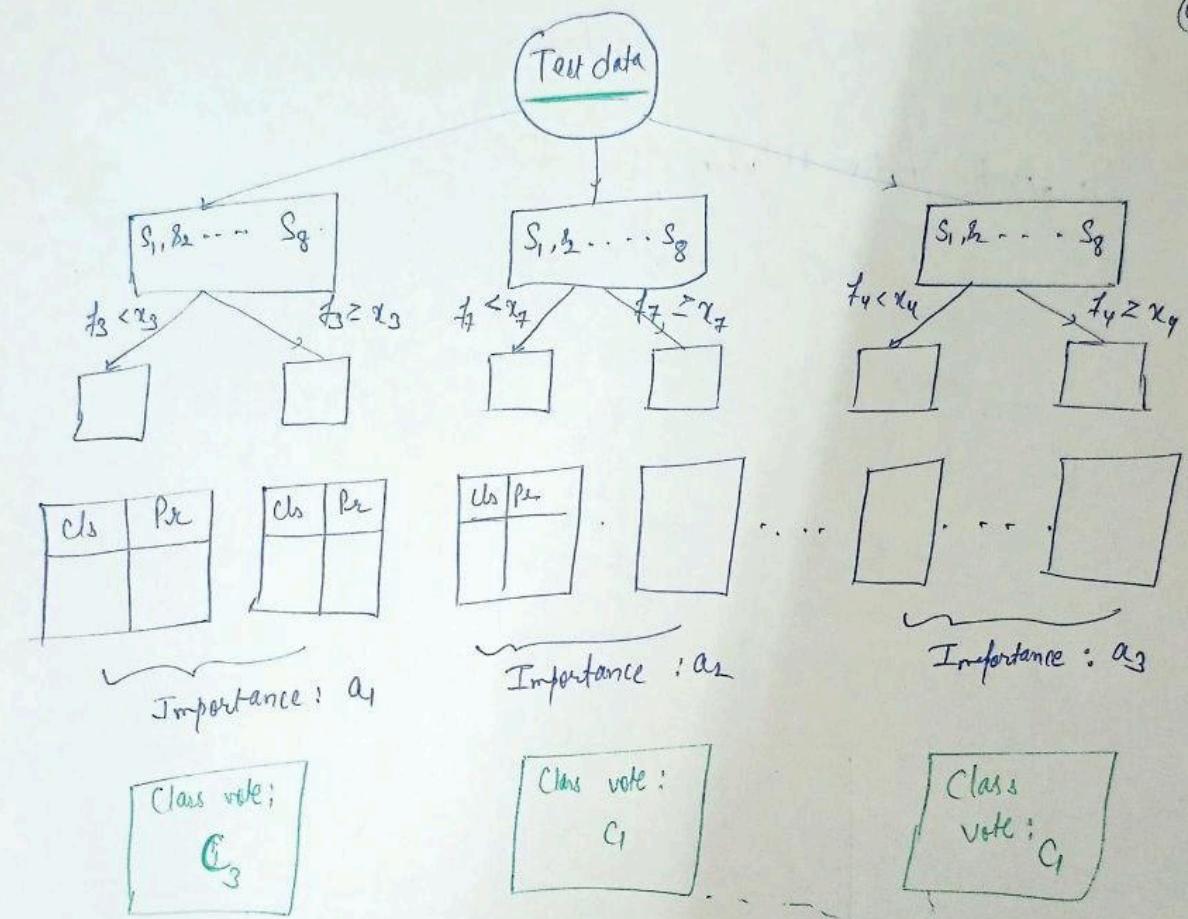
Importance:  $a_1$

Importance:  $a_2$



Importance:  $a_3$

(Q)



Calculate total Importance score of each class! —

Class	Total Importance Score
\$C_1\$	\$a_2 + a_3\$
\$C_2\$	0
\$C_3\$	\$a_1\$

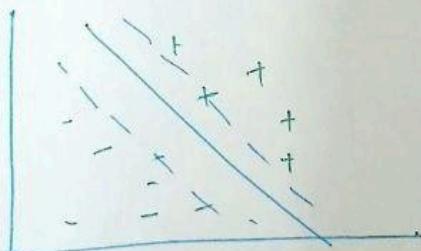
- variants of AdaBoost : XGBoost

(10)

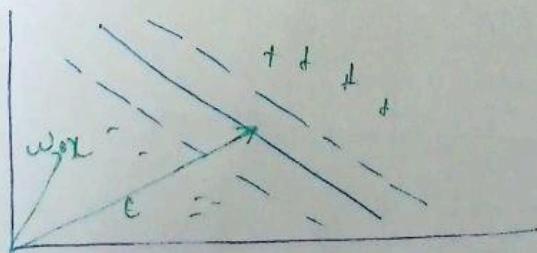
Support Vector Machine :-  
Let's say 2 class classification problem :-

- Linear classification.
- Find a line decision boundary among these class points.
  - Infinite decision boundary are possible.

S.t. line : Linear Decision Boundary.



Find DB that minimizes error.



find D.B whose margin max ( $\gamma_1, \gamma_2$ )

- we want DB to be in mid of close boundary / Margin.

•  $x$  is point in any class : vector.

→ Take projection of  $x$  in direction of  $w$ .

$$\vec{x} \cdot \vec{w} = d = \text{Scalar quantity} = \text{distance from origin}$$

Date - 24/10/2024      Lecture - 28

Goal :- How to draw Linear Boundary

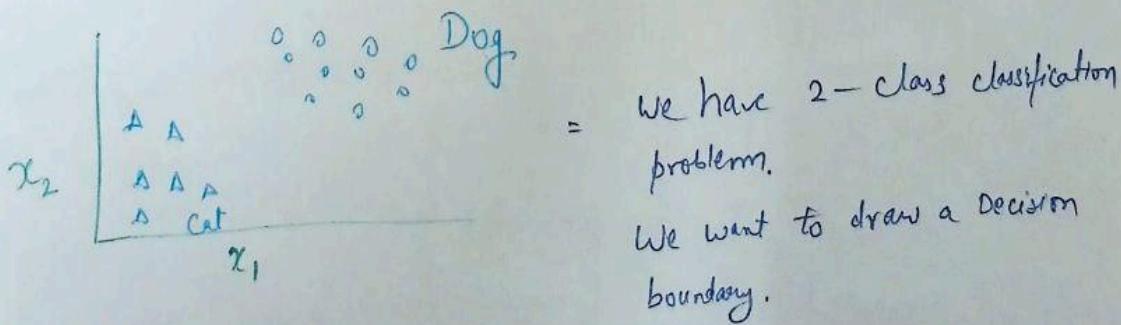
Initially in SVM, trying to solve linear decision boundary.

→ Linear Classification.

• we want good separation of points so that they are far apart from D.B.

• Slight change in feature level ( $H_t, w$  here) should not change D.B. classification.

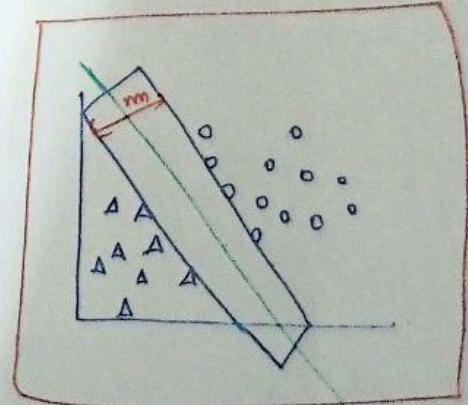
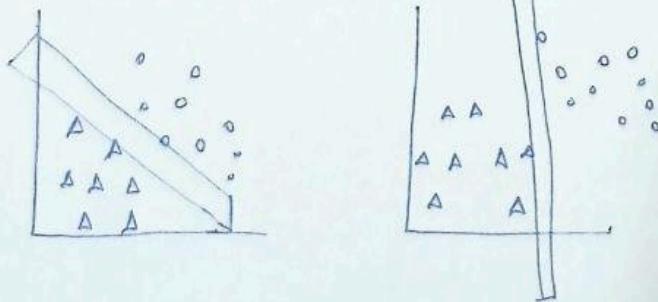
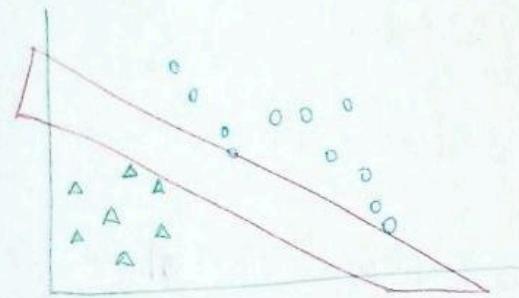
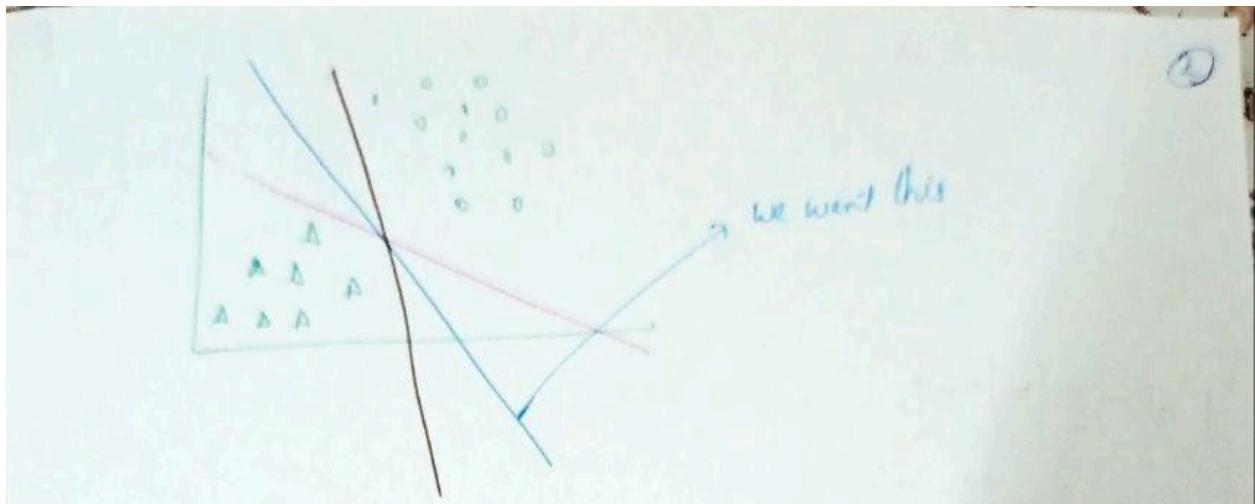
• There may be many hypothesis possible, from the Hypothesis space  
↓  
Choose hypothesis with good separation for class points.



We want to draw a Decision boundary.

$x_1$  : Height

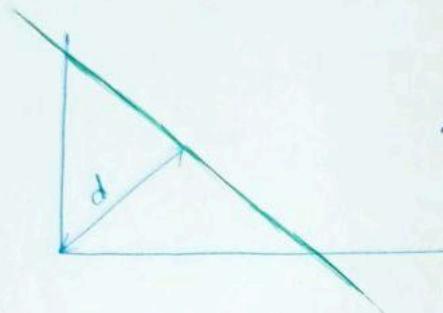
$x_2$  : weight



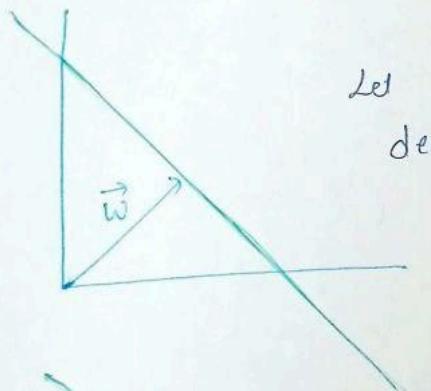
→ The best decision boundary is line through the middle of margin of 2 classes.

→ ' $m'$ ' is maximum margin b/w 2 classes.

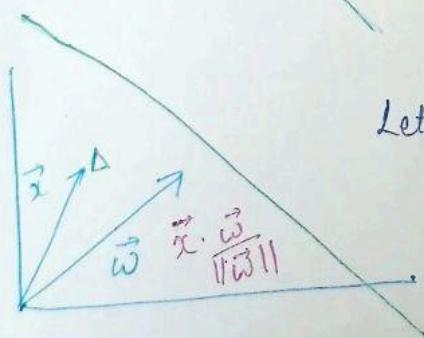
'm' : represents maximum separation b/w two-class boundaries. (3)



Let 'd' be distance of DB from origin.



Let  $\vec{w}$  be a vector perpendicular to decision boundary.



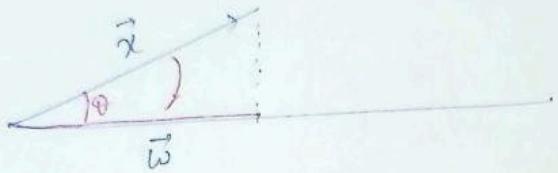
Let  $\vec{x}$  be another vector representing a data point.

Projection of  $\vec{x}$  on direction of  $\vec{w}$  is :-

$$\vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

Data point  $x$  is vector = 2d point, at some angle.

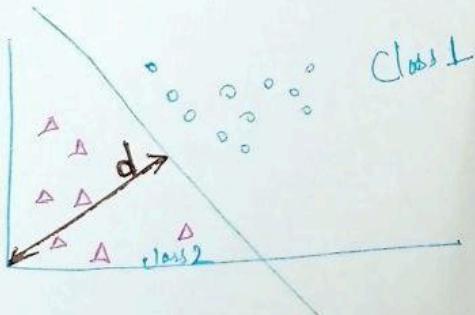
(4)



$$\vec{x} \cdot \vec{w} \cos \theta$$

If projection  $\frac{\vec{x} \cdot \vec{w}}{\|\vec{w}\|} \geq d$

↓  
data point  $\vec{x}$  belongs to class 1      Use it belongs to class 2.



So if  $\vec{x}$  belongs to class 1 if

$$\vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|} \geq d$$

$$\vec{x} \cdot \vec{w} \geq d \|\vec{w}\|$$

$$\vec{x} \cdot \vec{w} + b \geq 0 \quad \text{where } b = -d \|\vec{w}\|$$

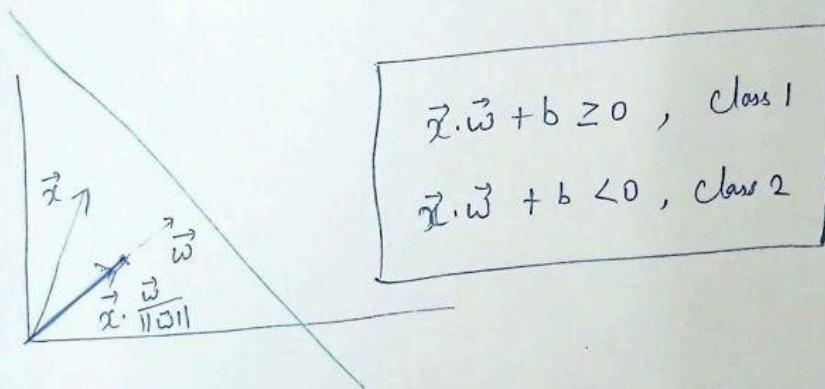
So if we can say  $\vec{x}$  belongs to class 2 if

(2)

$$\frac{\vec{x} \cdot \vec{w}}{\|\vec{w}\|} < d$$

$$\vec{x} \cdot \vec{w} < d \|\vec{w}\|$$

$$\vec{x} \cdot \vec{w} + b < 0 \quad \text{where } b = -d \|\vec{w}\|$$

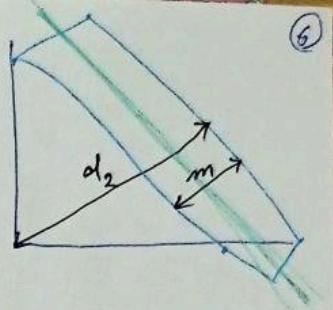
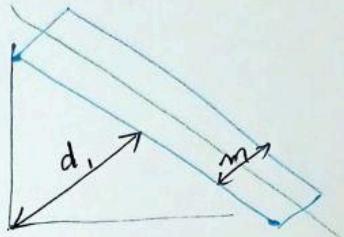
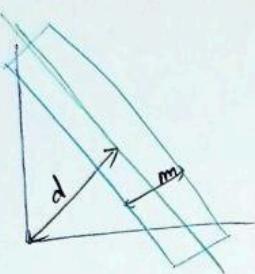


- if length of projection of  $x$  towards  $d$  is :-

$\rightarrow \downarrow d$  point  $x$  belongs to class i.

$\rightarrow \uparrow d$  belongs to class j.

• Slight change in position lead to different class.  
So we want Margin.



$$m = d_2 - d_1$$

Case 1 Suppose, we want  $\vec{x}$  belongs to class 1 if :-

$$\vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|} \geq d_2$$

$$\vec{x} \cdot \vec{w} \geq d_2 \|\vec{w}\|$$

$$\Rightarrow \vec{x} \cdot \vec{w} + b \geq 1 \quad \text{where } b = 1 - d_2 \|\vec{w}\|$$

Case 2 if  $\vec{x}$  belongs to class 2 +

$$\vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|} \leq d_1$$

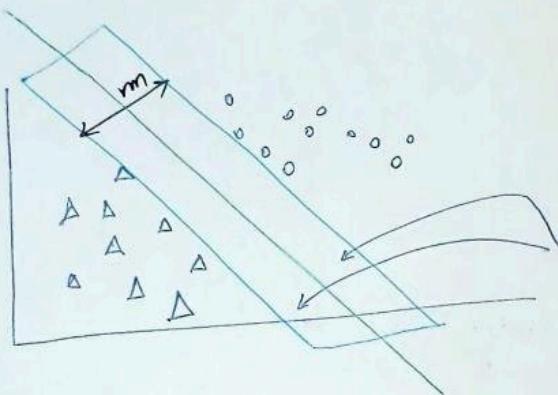
$$\vec{x} \cdot \vec{w} \leq d_1 \|\vec{w}\|$$

$$\vec{x} \cdot \vec{w} + b \leq -1, \quad \text{where}$$

$$b = -1 - d_1 \|\vec{w}\|$$

From training data we want to  
find  $w, b$

(7)



we want sample to fall  
in these regions.

$$\vec{x} \cdot \vec{w} + b \geq 1 \text{ then class 1.}$$

$$\vec{x} \cdot \vec{w} + b \leq -1 \text{ then class 2.}$$

. For classification, support vector is important.

Support vector help in making classification decision.

For  $i^{\text{th}}$  sample  $\vec{x}_i$  with class label  $y_i$  →

$$\vec{x}_i \cdot \vec{w} + b \geq 1, \text{ then class 1 } (80, y_i = 1)$$

$$\vec{x}_i \cdot \vec{w} + b \leq -1 \text{ then class 2 } (80, y_i = -1)$$



our decision boundary rule is →

$$y_i (\vec{x}_i \cdot \vec{w} + b) \geq 1$$

(8)

$$\vec{x}_i \cdot \vec{w} + b \leq -1 \quad \text{if } y_i = -1$$

$$y_i (\vec{x}_i \cdot \vec{w} + b) \geq 1 \quad (-1 * -1 = +1)$$

So we need to classify Test data :-

Data  $\vec{x}$ ,

if  $\vec{x} \cdot \vec{w} + b \geq 1$ , then  $\vec{x}$  belongs to class 1.

if  $\vec{x} \cdot \vec{w} + b \leq -1$ , then  $\vec{x}$  belongs to class 2.

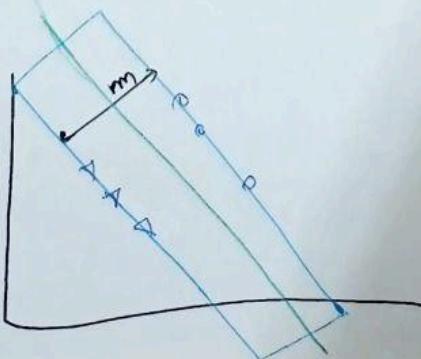
So test data be  $\vec{x}_t$ , so

[ if  $\vec{x}_t \cdot \vec{w} + b \geq 1$ , then  $\vec{x}_t$  belongs to class 1 ]

[ if  $\vec{x}_t \cdot \vec{w} + b \leq -1$ , then  $\vec{x}_t$  belongs to class 2 ]

To make this decision, we need to find out  $\vec{w}, b$ .

The goal of SVM is to find out  $\vec{w}, b$  from training data.

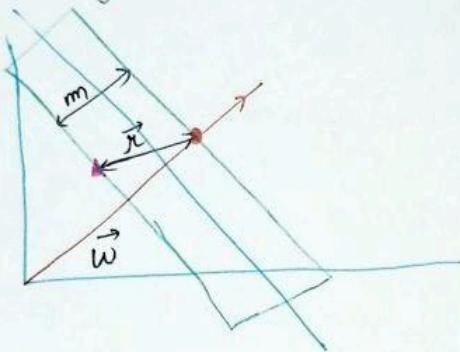


$\vec{x}_i \cdot \vec{w} + b = 1$ , class 1
$y_i = 1$
$\vec{x}_i \cdot \vec{w} + b = -1$ , class 2
$y_i = -1$

$$y_i (\vec{x}_i \cdot \vec{w} + b) = 1$$

⑦

These training data points are called Support vectors.



Let  $\vec{x}$  be vector b/w samples on two class-boundaries.

Width  $m$  = projection of  $\vec{x}$  on perpendicular direction of boundary

$m = \vec{x} \cdot (\text{unit vector on perpendicular direction of boundary})$

$$\boxed{m = \vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|}}$$

width  $m = \vec{x} \cdot \frac{\vec{w}}{\|\vec{w}\|}$

$$x = (\vec{x}_2 - \vec{x}_1)$$

Width  
 $m = (\vec{x}_2 - \vec{x}_1) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{(\vec{x}_2 \cdot \vec{w} - \vec{x}_1 \cdot \vec{w})}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$

- We want to find such a  $\vec{w}$ ,  $b$  so that width is maximized, <sup>(10)</sup>  
subject to constraint of our decision rule

$$y_i(\vec{x}_i \cdot \vec{w} + b) \geq 1$$

To solve constrained optimization problem :-

Langrangian Function.

- This is equivalent to minimizing  $\frac{1}{2} \|\vec{w}\|^2$  subject to

$$y_i(\vec{x}_i \cdot \vec{w} + b) \geq 1$$

$m \uparrow\uparrow$	$w \downarrow\downarrow$	$m = \frac{2}{\ w\ }$
$m \downarrow\downarrow$	$w \uparrow\uparrow$	

$$L = \text{take item we want to minimize} - \sum_{i=1}^N \alpha_i [y_i(\vec{x}_i \cdot \vec{w} + b) - 1]$$

add or subtract other constraint.

$\alpha_i$  = value will diff. for each training data.

$L = \frac{1}{2} \ \vec{w}\ ^2 - \sum_{i=1}^N [\alpha_i (y_i(\vec{x}_i \cdot \vec{w} + b) - 1)]$
--

$\alpha_i$  = Lagrange multiplier

So, considering N training data points, the lagrangian fun is - ⑪

$$L = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^N [\alpha_i y_i (\vec{x}_i \cdot \vec{w} + b) - 1]$$

→ I want to minimize L wrt  $\vec{w}$  and  $b$ , so -

$$\frac{\partial L}{\partial w} = 0 \quad \text{and} \quad \frac{\partial L}{\partial b} = 0$$

Note :- The minus (-) sign in the Lagrangian fun indicates that I want to maximize L wrt  $\alpha$ .

setting  $\frac{\partial L}{\partial w} = 0$ , we get :-

$$\vec{w} = \sum_{i=1}^N (\alpha_i y_i \vec{x}_i)$$

setting  $\frac{\partial L}{\partial b} = 0$ , we get :-

$$\sum_{i=1}^N \alpha_i y_i = 0$$

After putting value and simplifying, we get :-

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) + \sum_{i=1}^N \alpha_i$$

subject to  $\alpha_i \geq 0, \forall i \in [1, N]$

$$\sum_{i=1}^N (\alpha_i y_i) = 0$$

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_j \alpha_j y_j (\vec{x}_i \cdot \vec{x}_j) + \sum_{i=1}^N \alpha_i$$

. Optimization depends on dot product of pair of training samples  
 $(\vec{x}_i \cdot \vec{x}_j)$

→ For all Support vector :-

$$\boxed{\alpha_i (y_i (\vec{w}_i \cdot \vec{x} + b) - 1) = 0}$$

Then either

$$\boxed{\alpha_i = 0}$$

or

$$\boxed{y_i (\vec{w}_i \cdot \vec{x} + b) - 1 = 0}$$

(i) for support vector :-

$$y_i (\vec{w}_i \cdot \vec{x} + b) = 1$$

(ii) for all not support vector  $\alpha_i = 0$

It can be shown that solution must satisfy :—

(13)

$$\alpha_i \geq 0, \forall i \in [1, N]$$

$$y_i(\vec{x}_i \cdot \vec{w} + b) - 1 \geq 0$$

and

$$\alpha_i(y_i(\vec{x}_i \cdot \vec{w} + b) - 1) = 0$$

Last condition requires—

- either  $y_i(\vec{x}_i \cdot \vec{w} + b) = 1$  (which is true for support vectors)

or

$\alpha_i = 0$  (i.e. for training data points, other than support vectors,  $\alpha_i = 0$ )

$\alpha_i = 0$  = non-support vectors

Can absorb  $\alpha_i = 0$  for support vector.

- The Dual form :-

It can be shown that solution must satisfy :-

$$\alpha_i \geq 0, \forall i \in [1, N]$$

$$y_i(\vec{x}_i \cdot \vec{w} + b) - 1 \geq 0$$

and

$$\alpha_i (y_i(\vec{x}_i \cdot \vec{w} + b) - 1) = 0$$

Last condition requires :-

either  $y_i(\vec{x}_i \cdot \vec{w} + b) = 1$  { which is true for support vectors }

$\alpha_i = 0$  { i.e. for training data points, other than support vectors,  $\alpha_i = 0$  }

[Lagrange multiplier of S.V may/may not be 0]

By maximization of :-

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) + \sum_{i=1}^N \alpha_i$$

→ Maximization of L wrt  $\alpha$ , we can find values of  $\alpha_i$

$\forall i \in [1, N]$  (i.e., all training data)

We can find :-

$$\vec{w} = \sum_{i=1}^N (\alpha_i y_i \vec{x}_i) \quad (2)$$

and, it can be shown that :-

$$b = \frac{1}{N_s} \sum_{i=1}^{N_s} \left( y_i - \sum_{j=1}^N \alpha_j y_j (\vec{x}_i \cdot \vec{x}_j) \right)$$

$N_s$  : no. of support vectors

If not support vector  $\alpha_j = 0$   $\Downarrow$   $w, b = 0$  non-support vectors  
don't contribute to  
calculation)

St. line only depends on support vectors.

↓  
Ignore other points than s.v not useful in decision making.

Started with ..  
 $x_t \cdot w + b$   
 replaced

Testing :-

Suppose I want to find class label of test data  $x_t$ .

I'll calculate :-

$$q_t = (\vec{x}_t \cdot \vec{w} + b) = \vec{x}_t \left( \sum_{i=1}^N (\alpha_i y_i \vec{x}_i) \right) + b = \sum_{i=1}^N (\alpha_i y_i \vec{x}_i \cdot \vec{x}_t) + b$$

If  $q_t \geq 1$ , test data  $x_t$  belongs to class  $y_t = 1$ . ③

If  $q_t \leq -1$ , test data  $x_t$  belongs to class  $y_t = -1$ .

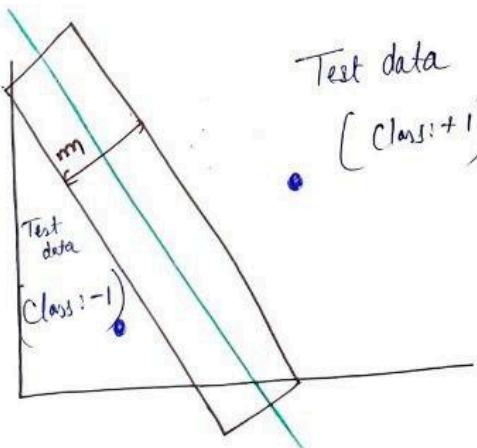
here 
$$\boxed{q_t = \sum_{i=1}^N \alpha_i y_i (x_i \cdot x_t) + b}$$

dot product of training with data of test shows :—

→ All training data contribute to decision making for testing data.

here also  $\alpha_i = 0$  for non-support vector.

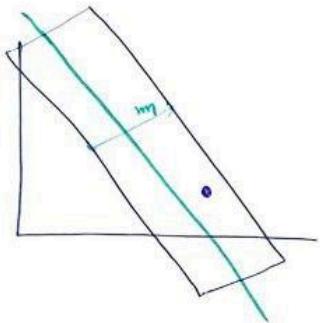
only those data lying on boundary of class.  
They decide for decision of test data.



$$q_t = \sum_{i=1}^N (\alpha_i y_i \vec{x}_i \cdot \vec{x}_t) + b$$

Therefore, data points which are not support vector  
don't play role in determining class label of test  
data.

(4)

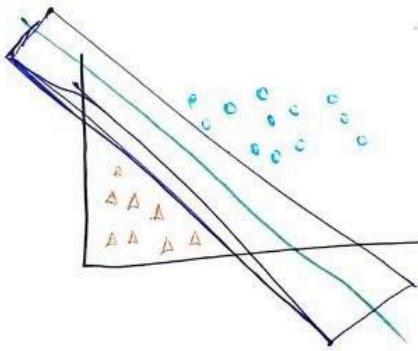


For finding out class label of test data  $x_t$ :

what happens if

$$1 > q_t > -1 ?$$

### Maximum Margin Classifier:



Since we are trying to maximize margin during training,  
it is called as Maximum Margin Classifier.

$\ell_v = \delta h^o = \text{Slack variable.}$

### SVM with Training data within Margin $\rightarrow$ Soft Margin SVM

Decision rule for  $i^{th}$  sample  $x_i$  with class label  $y_i$

$$\vec{x}_i \cdot \vec{w} + b \geq 1 - \ell_{v_i}, \text{ then class 1 } (so, y_i = 1)$$

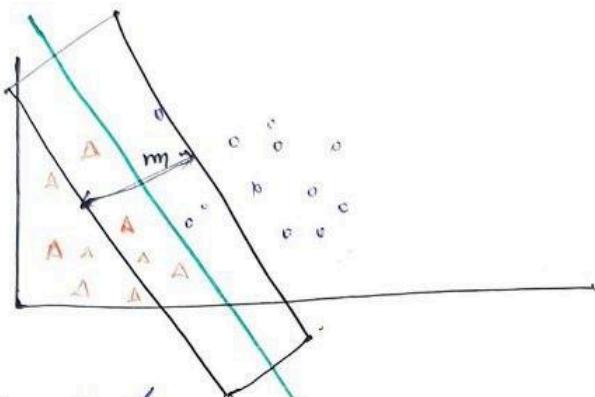
$$\vec{x}_i \cdot \vec{w} + b \leq -(1 - \ell_{v_i}) \text{ then class 2 } (so, y_i = -1)$$

Considering both conditions

our decision rule +

$$\boxed{y_i(\vec{x}_i \cdot \vec{w} + b) \geq 1 - \ell_{v_i}}$$

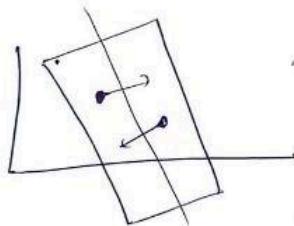
(5)



Characteristics of  $\epsilon_i$

$\epsilon_i \geq 0$  is called slack variable.

if  $\epsilon_i > 1$ , data is misclassified.



If  $\epsilon_i = 1$ , data lies on decision boundary / b.c

$\epsilon_i > 1$  = data crossed decision boundary.

$\epsilon_i = 1$  exactly lying on D.B

$0 < \epsilon_i < 1$  : Soft Margin SVM

For Hard Margin Classifier

$$y_i(\vec{x}_i \cdot \vec{w} + b) \geq 1 - \epsilon_i$$

$\uparrow$   
0

$$y_i(\vec{x}_i \cdot \vec{w} + b) \geq 1$$

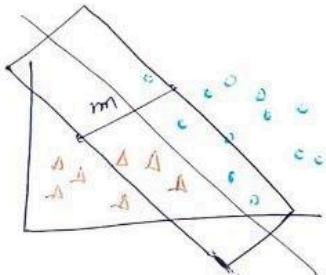
Ideally we want  $\ell_w$  to be as less as possible. ⑥

or  $\ell_w = 0$  (want)

Soft Margin SVM:-

$$\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^N \ell_{\nu_i}$$

Subject to -  $y_i(\vec{x}_i \cdot \vec{w} + b) \geq 1 - \ell_{\nu_i}$  and  $\ell_{\nu_i} \geq 0$



$C$  : determines weightage on penalty on slack variable.

{ for training data, we need to know  $w, b, \ell_{\nu_i}$  }

$C$  = constant, which is how much penalty we apply

if  $\ell_w > 0$

→ Second term in the ftn to be minimized is penalty on points crossing class boundary

→ Minimization of 2nd term ensures that many training data points don't cross boundary by a large amount.

⑦

## Soft Margin SVM: Lagrangian function :-

Minimize :-

$$\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^N \epsilon_i$$

Subject to :-

$$y_i (\vec{x}_i \cdot \vec{w} + b) \geq 1 - \epsilon_i \quad \text{and} \quad \epsilon_i \geq 0$$

The Lagrangian funt is :-

$$L = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^N \epsilon_i - \sum_{i=1}^N [\alpha_i \{y_i (\vec{x}_i \cdot \vec{w} + b) - 1 + \epsilon_i\}] - \sum_{i=1}^N \mu_i \epsilon_i$$

$$L = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^N \epsilon_i - \sum_{i=1}^N [\alpha_i \{y_i (\vec{x}_i \cdot \vec{w} + b) - 1 + \epsilon_i\}] - \sum_{i=1}^N \mu_i \epsilon_i$$

$\alpha_i \geq 0, \mu_i \geq 0$  : Lagrange Multiplier

From this we can solve for  $\vec{w}$ ,  $b$ , and  $\epsilon_i$ .

Minimize :-

$$\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^N \epsilon_i$$

Subject to

$$y_i (\vec{x}_i \cdot \vec{w} + b) \geq 1 - \epsilon_i \quad \text{and} \quad \epsilon_i \geq 0$$

⑧

$C$  is a hyperparameter of method.

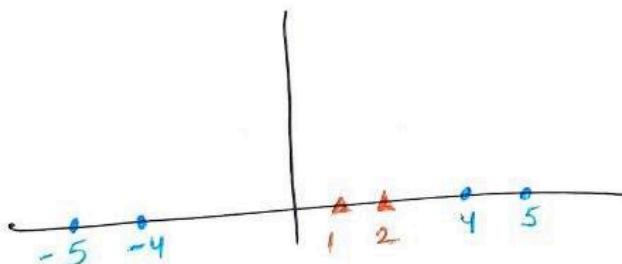
$$\begin{aligned} (\vec{x}_i \cdot \vec{w} + b) &\geq 1 - \epsilon_i \rightarrow \text{class 1 (+1)} \\ &\leq -(1 - \epsilon_i) \rightarrow \text{class 2 (-1)} \end{aligned}$$

$w, b$  = we'll minimize

for  $\ell_C$  = we'll maximize;

as  $-\ell_C$  = already -ve sign there

Let's consider :-



Considering I have one-dimensional dataset with feature  $x_i$ .

There are 2 classes.

Can we classify data using our SVM designed so far?

No, because data points are not linearly separable.

Now if we still want to classify linearly:-

Making new dimension :-  $x_i^2$

⑨

If we still want to classify:

new dimension  $x_1^2$ .

$$x_1 \xrightarrow{\phi(x)} (x_1, x_1^2)$$

$$1 \rightarrow (1, 1)$$

$$-4 \rightarrow (-4, 16)$$

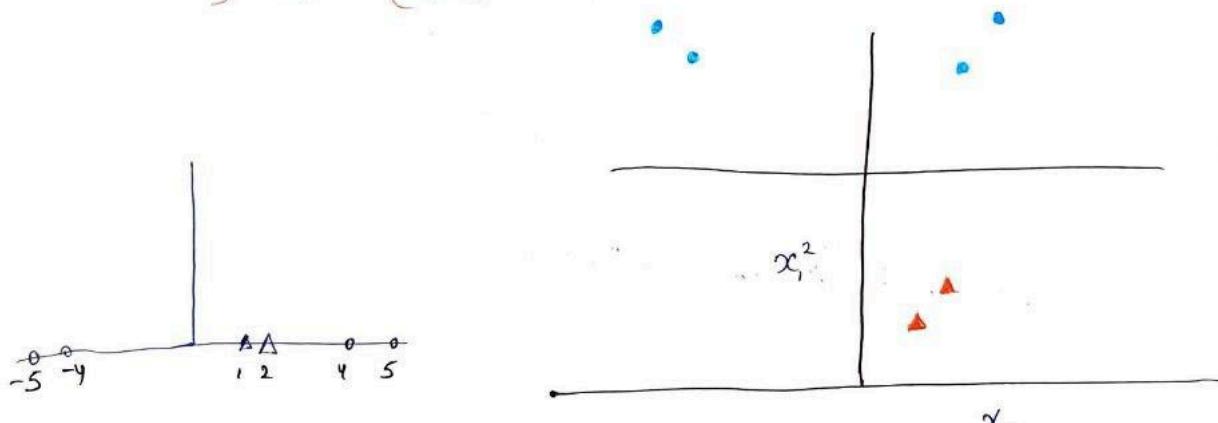
$$2 \rightarrow (2, 4)$$

$$-5 \rightarrow (-5, 25)$$

$$3 \rightarrow (3, 9)$$

$$4 \rightarrow (4, 16)$$

$$5 \rightarrow (5, 25)$$



$$\phi(x)$$

transformed data is linearly separable.

So SVM can be used to classify them.

$\phi(x)$ : takes a vector  $x$  and transforms it into a higher-dimensional vector  $\phi(x)$ .

(10)

## SVM with Transformed Data: Testing

SVM with original SVM:-

Suppose I want to find out the class label of test data  $\vec{x}_t$ .

I'll calculate :-

$$\begin{aligned} q_t &= (\vec{x}_t \cdot \vec{w} + b) \\ &= \vec{x}_t \left( \sum_{i=1}^N (\alpha_i y_i \vec{x}_i) \right) + b \\ &= \sum_{i=1}^N (\alpha_i y_i \cdot \vec{x}_i \cdot \vec{x}_t) + b \end{aligned}$$

## SVM with Transformed Data

Suppose I want to find out class label of test data  $\vec{x}_t$ .

I'll calculate :-

$$\begin{aligned} q_t &= (\phi(\vec{x}_t) \cdot \vec{w}_t + b) \\ &= \vec{x}_t \cdot \left( \sum_{i=1}^N (\alpha_i y_i \phi(\vec{x}_i)) \right) + b \\ &= \sum_{i=1}^N (\alpha_i y_i \phi(\vec{x}_i) \cdot \phi(\vec{x}_t)) + b \end{aligned}$$

(1)

- For every data point, finding transformation  $\phi(x)$ .  
Taking dot products is computationally challenging.

The Dual Form :-

After putting values and simplifying, we get :-

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\vec{\phi}(x_i) \cdot \vec{\phi}(x_j)) + \sum_{i=1}^N \alpha_i$$

We want to maximize  $L$  w.r.t  $\alpha$ 

Subject to :-

$$\alpha_i \geq 0, \forall i \in [1, N]$$

$$\sum_{i=1}^N (\alpha_i y_i) = 0$$

Testing → To find out class label of test data  $x_t$  :-

$$q_t = (\vec{\phi}(x_t) \cdot \vec{w} + b)$$

$$q_t = \sum_{i=1}^N (\alpha_i y_i \vec{\phi}(x_i) \cdot \vec{\phi}(x_t)) + b$$

If we take data to higher dimension space we can often  
not always classify linearly data.

now what we are doing :- Transform data → Take Dot Product

What if we can do! —

(12)

The Dual Form :-

After putting values and simplifying :-

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (\vec{\phi}(x_i) \cdot \vec{\phi}(x_j)) + \sum_{i=1}^N \alpha_i$$

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j) + \sum_{i=1}^N \alpha_i$$

Why Kernel :-  
if 100 million data  $\xrightarrow{\text{Higher Dimen}}$   $(100 \text{ million})^2$

what if 100 million data with 200 dimension ?

It is difficult to compute  $\vec{\phi}$  for all dimension with square.

what if infinite dimension : Hard

Practically not feasible for transforming data into that much

Higher D

Testing :- Suppose I want to find out class label of test data  $x_t$

$$q_t = (\vec{\phi}(x_t), \vec{w} + b)$$

$$= \sum_{i=1}^N (\alpha_i y_i \underbrace{K(\vec{x}_i, \vec{x}_t)}_{}) + b$$

$$K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

(3)

Suppose  $x$  in 1D,  $\phi$  transforms  $x$  to 2D then dot product

Kernel function → It helps us getting dot product of transformed data without requiring calculation of transformation of every data point.

Kernel function → Converts dot product after transformation to original data product and (square if  $1D \rightarrow 2D$ )

Example → We have 2D points →

one data point is  $x: (x_1, x_2)$

Another data point is  $z: (z_1, z_2)$

using transformation →

$$\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$\phi(z) = (z_1^2, \sqrt{2}z_1z_2, z_2^2)$$

$$\text{So } \phi(\vec{x}) \cdot \phi(\vec{z}) = (x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2)$$

$$= (x_1 z_1 + x_2 z_2)^2$$

$$\Rightarrow (\vec{x} \cdot \vec{z})^2$$

no need to go for  
Higher dimension of

$$= K(\vec{x}, \vec{z})$$

(dot product)<sup>2</sup>  
kernel

(14)

$$K(\vec{x}, \vec{z}) = (\vec{x} \cdot \vec{z})^2$$

By using kernel, get rid of need to calculate  $\phi(x)$  and  $\phi(z)$

$$\text{eg: } A = (a_1, a_2)$$

$$B = (b_1, b_2)$$

$$A \cdot B = (a_1 b_1 + a_2 b_2)$$

(Dot product is = Multiplication & addition)

we get rid of higher computational overhead to changing data point to higher dimension to dot product and square in same dimension.

Kernel function : If I have 2D data points

one data point is  $x = (x_1, x_2)$

Another data point  $z = (z_1, z_2)$

↓  
Using Transformation :

$$\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$\phi(z) = (z_1^2, \sqrt{2}z_1z_2, z_2^2)$$

$$\text{So } \vec{\phi}(x) \cdot \vec{\phi}(z) = (x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2) \\ = (x_1 z_1 + x_2 z_2)^2$$

$$= K(\vec{x}, \vec{z})$$

$$K(\vec{x}, \vec{z}) = (\vec{x} \cdot \vec{z})^2$$

Using Kernel, I get rid of  
need to calculate  $\phi(x), \phi(z)$

$$\begin{array}{ccc} 2D & \longrightarrow & 3D \\ (x_1, x_2)(z_1, z_2) & & x_1^2, \sqrt{2}x_1x_2, x_2^2 \\ \text{Dot product } \phi(x) \cdot \phi(z) = (\vec{x} \cdot \vec{z})^2 & & \end{array}$$

no need to  
go in higher  
dimension

"Kernel fun" helps to get rid of moving/transforming to higher D. ①

e.g. PCA need to find dot product to find PCA components.

Kernel fun can be applied to PCA

### Kernel PCA

RBF Kernel (Radial Basis funtn | Gaussian Kernel)

$$K(\vec{x}, \vec{z}) = \exp\left(-\frac{\|\vec{x} - \vec{z}\|^2}{2\sigma^2}\right)$$

can also be represented as:

$$K(\vec{x}, \vec{z}) = \exp\left(-\gamma \|\vec{x} - \vec{z}\|^2\right)$$

$$= \exp\left(-\gamma (\|\vec{x}\|^2 + \|\vec{z}\|^2 - 2\|\vec{x}\|\|\vec{z}\|)\right)$$

$$\exp\left(-\gamma (\|\vec{x}\|^2 + \|\vec{z}\|^2)\right) \exp\left(2\gamma \|\vec{x}\|\|\vec{z}\|\right)$$

$$= \exp\left(-\gamma (\|\vec{x}\|^2 + \|\vec{z}\|^2)\right) \left[ 1 + \underbrace{\frac{2\gamma \|\vec{x}\|\|\vec{z}\|}{1!}}_{+} + \underbrace{\frac{(2\gamma \|\vec{x}\|\|\vec{z}\|)^2}{2!}}_{+} + \dots \right]$$

$$= \gamma^2 \left[ 1 + \underbrace{\frac{\alpha^2 \|\vec{x}\|\|\vec{z}\|}{1!}}_{+} + \underbrace{\frac{(\alpha^2 \|\vec{x}\|\|\vec{z}\|)^2}{2!}}_{+} + \dots \right]$$

where  $\gamma^2 = \exp\left(-\gamma (\|\vec{x}\|^2 + \|\vec{z}\|^2)\right)$ ,  $\alpha^2 = 2\gamma$

RBF: Circular type of Bias fun<sup>n</sup>

②

$$\gamma = \frac{1}{2\sigma^2}$$

for support vector user need to input  $\gamma$ .

when using kernel if one need to go in infinite space

- infinite term in  $\phi(\vec{x})$
- " " "  $\phi(\vec{z})$

$\phi(\vec{x}) \cdot \phi(\vec{z})$  can take data to  $\infty$  term / space

Normally can't take data from finite to infinite dimension space.  
Kernel fun<sup>n</sup> do.

- RBF Kernel is used widely.
- RBF kernel practically utilizes an  $\infty$  dimensional transformation.

Polynomial Kernel

$$K(\vec{x}, \vec{z}) = (\vec{x} \cdot \vec{z} + r)^d$$

e.g. if  $r = \frac{1}{2}$ ,  $d = 2$

$$K(\vec{x}, \vec{z}) = (\vec{x} \cdot \vec{z} + \frac{1}{2})^2$$

$$= \|\vec{x}\| \|\vec{z}\| + \|\vec{x}\| \cdot \|\vec{z}\|^2 + \frac{1}{4}$$

$$= \left( \|\vec{x}\|, \|\vec{z}\|, \frac{1}{2} \right) \cdot \left( \|\vec{x}\|, \|\vec{z}\|^2, \frac{1}{2} \right)$$

$$= \phi(\vec{x}) \cdot \phi(\vec{z})$$

$$\phi(\vec{x}) = (\|\vec{x}\|, \|\vec{x}\|^2, \frac{1}{2}) \quad (4)$$

$$\phi(\vec{z}) = (\|\vec{z}\|, \|\vec{z}\|^2, \frac{1}{2})$$

- every fun is not Kernel function.
- Need to follow some criteria

(Gram Matrix)

for a fun, I have to prove Gram matrix to be  
semi definite.

$$\phi(\vec{x}) \cdot \phi(\vec{z}) \text{ should be able to} = K(\vec{x} \cdot \vec{z})$$

Neural Networks S.V need training data label, supervised.

→ Random forest need to find probability can't be applied. (supervised)

- NN can be applied to unsupervised data too.
- If data is image, video convert data into feature vector in supervised ML.  
↓ NN can directly used image, video in input.

$w_j$  = How much importance to  $x_1, x_2, \dots$  Layer linear.

$y$  = linear, every layer produces linear o/p

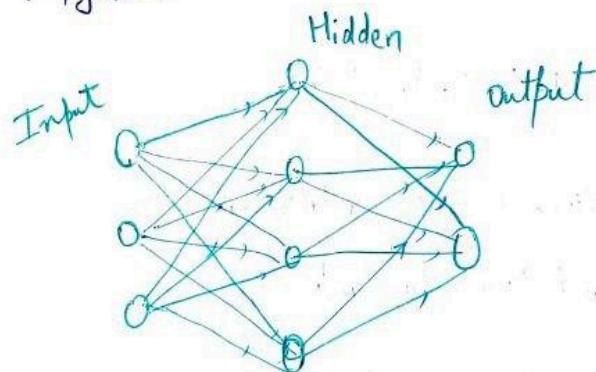
Neural Netw r  
why NN:- It can be used for different type of tasks → ⑤

- Classification      • Segmentation      many more...
- Regression      • Translation
- Detection      • Feature extraction

Why now:- Availability of hardware  
 . CPU      . Storage      . Cloud services

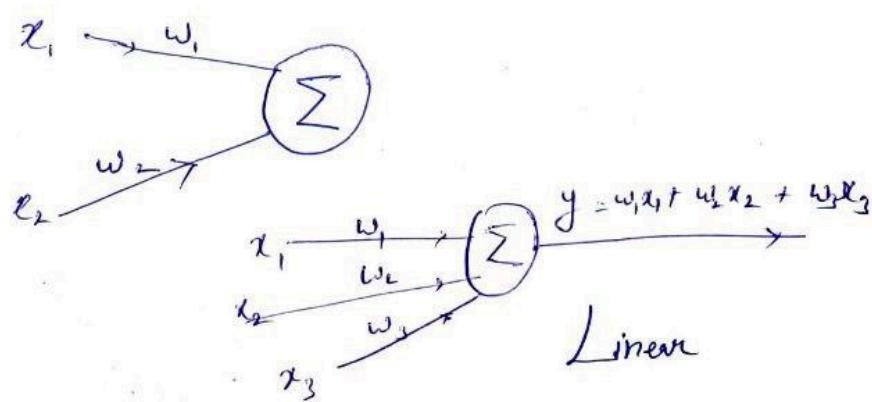
- Availability of large datasets.

→ Availability of implementation platforms:-  
 . PyTorch      . Tensorflow

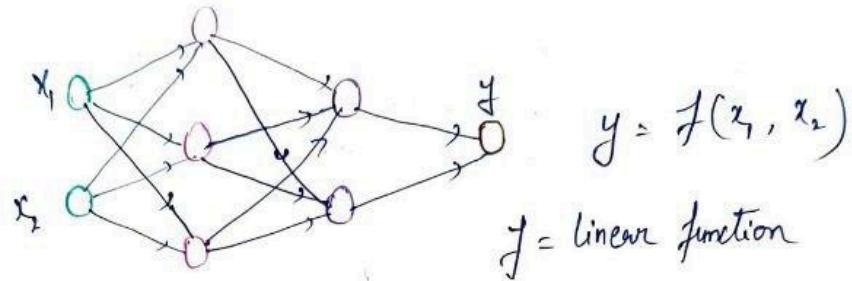


Neural netw came from biological context.

Neuron :-



(6)



$x_1$ : No. of cars per sq. km

$x_2$ : No. of industries per sq. km

$y$ : Predicted AQI (air quality index)

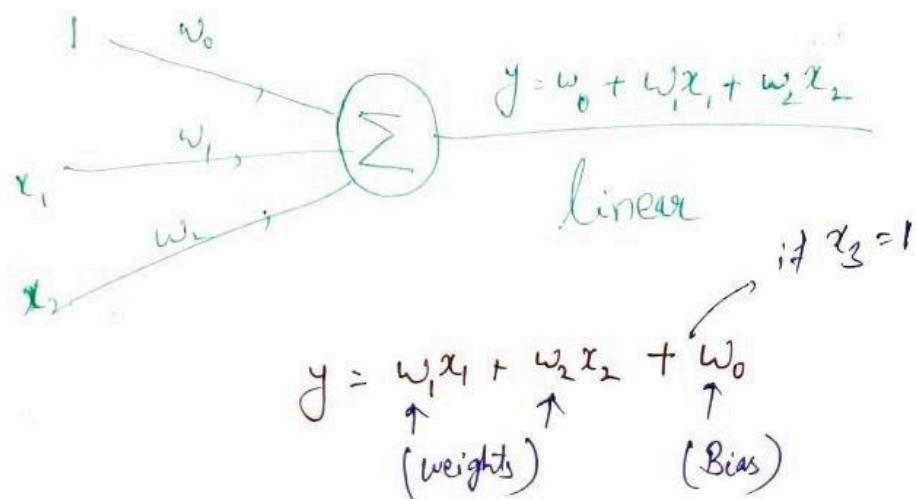
$$y = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3$$

goal is to find  
 $\alpha_1, \alpha_2, \alpha_3$

Assumptions: linear relationship

→ we don't know constants  $\alpha_i$

N.N can learn each  $\alpha_i$



(7)

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\begin{aligned} y &= w_1 x_1 + w_2 x_2 + b \\ &\downarrow \\ w^T x + b &= w^T x' \end{aligned}$$

$$w' = \begin{bmatrix} w_1 \\ w_2 \\ w_0 \end{bmatrix} \quad x' = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

$$\text{Also if } y(1) = w_1 x_1(1) + w_2 x_2(1) + b$$

$$\downarrow \\ w^T x(1) + b$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad x(1) = \begin{bmatrix} x_1(1) \\ x_2(1) \end{bmatrix}$$

↑  
for one data point (its features)

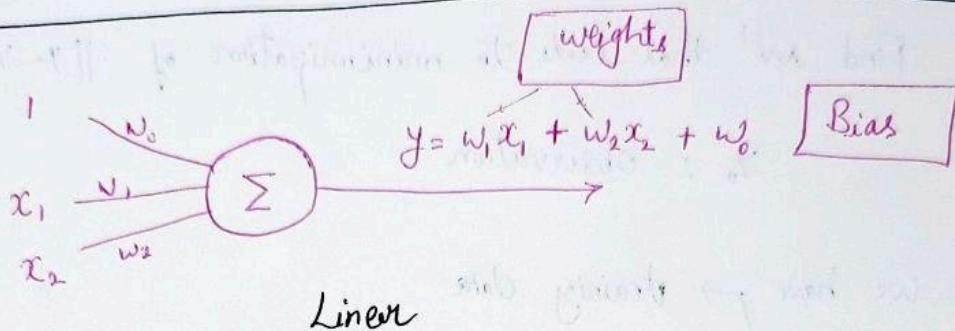
Suppose n such data points :-

$x(1), x(2), \dots, x(n)$  then -

$$X = \begin{bmatrix} x(1) & x(2) & \dots & x(n) \end{bmatrix} = \begin{bmatrix} x_1(1) & \dots & x_1(n) \\ x_2(1) & \dots & x_2(n) \\ x_3(1) & \vdots & \vdots \end{bmatrix}$$

$x(1)$  represent 1st data.

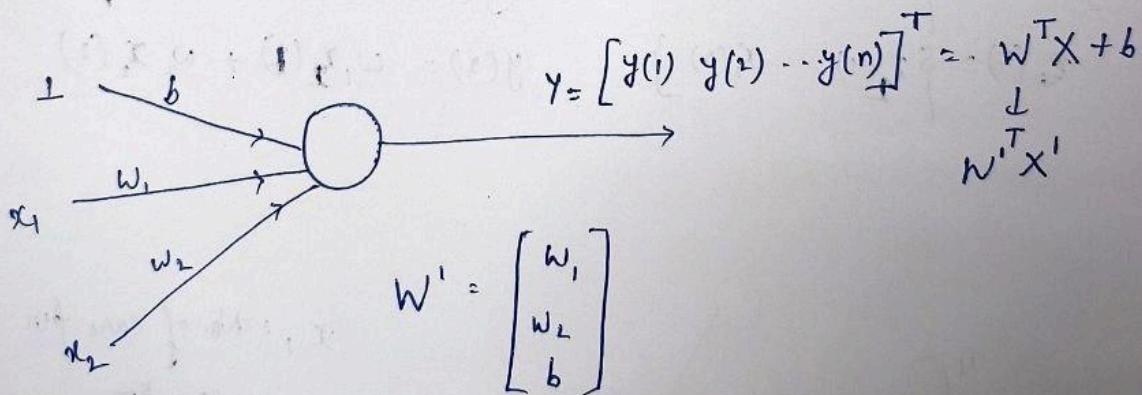
n = no. of data points, each with some specific dimension.



$$w^I = \begin{bmatrix} w_1 \\ w_2 \\ w_0 \end{bmatrix} \quad x^I = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

$$x^{(1)} = [x_{1(1)} \quad x_{1(2)} \dots]^T$$

$$x^{(2)} = [x_{2(1)} \quad x_{2(2)} \dots]^T$$



$$X^I = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \end{bmatrix} = \begin{bmatrix} x_{1(1)} & \dots & x_{1(n)} \\ x_{2(1)} & \dots & x_{2(n)} \\ \vdots & \ddots & \vdots \end{bmatrix}$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \end{bmatrix} = A \otimes I$$

$$Y_o = \begin{bmatrix} y_o^{(1)} & y_o^{(2)} & \dots & y_o^{(n)} \end{bmatrix} \text{ and } NN \text{ outputs}$$

$$y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(n)} \end{bmatrix}$$

How to find  $w'$ , i.e.  $w_1, w_2, b$  analytically? ②

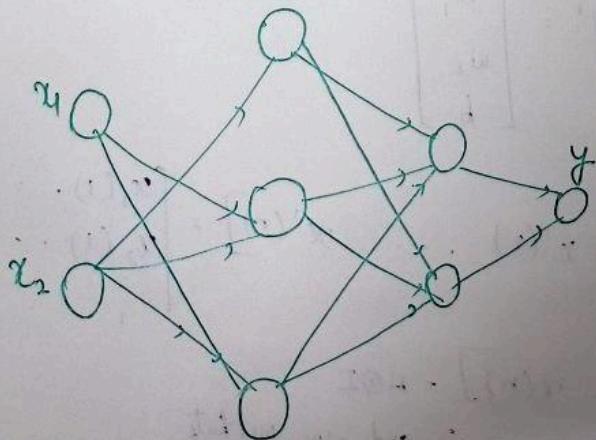
Find  $w'$  that leads to minimization of  $\|y - y_0\|^2$

$y_0$  : observation

- We have
  - training data
  - Actual fb value
- Closed form solution: Analytical approach
  - ↓ iteratively updating eqn  $w, b$ .
  - equating  $w, b = 0$  then updating eqn.

$$x(1) = \begin{bmatrix} x_1(1) & x_2(1) \end{bmatrix} \quad y(1) = w_1 x_1(1) + w_2 x_2(1)$$

$$x(2) = \begin{bmatrix} x_1(2) & \dots & x_n(2) \end{bmatrix} \quad y(2) = w_1 x_1(2) + w_2 x_2(2).$$



$x_1$ : No. of cars per sq. Km

$x_2$ : No. of industries per sq. Km

Assumption : Linear Relationship

But we don't know constants  $\alpha_i$

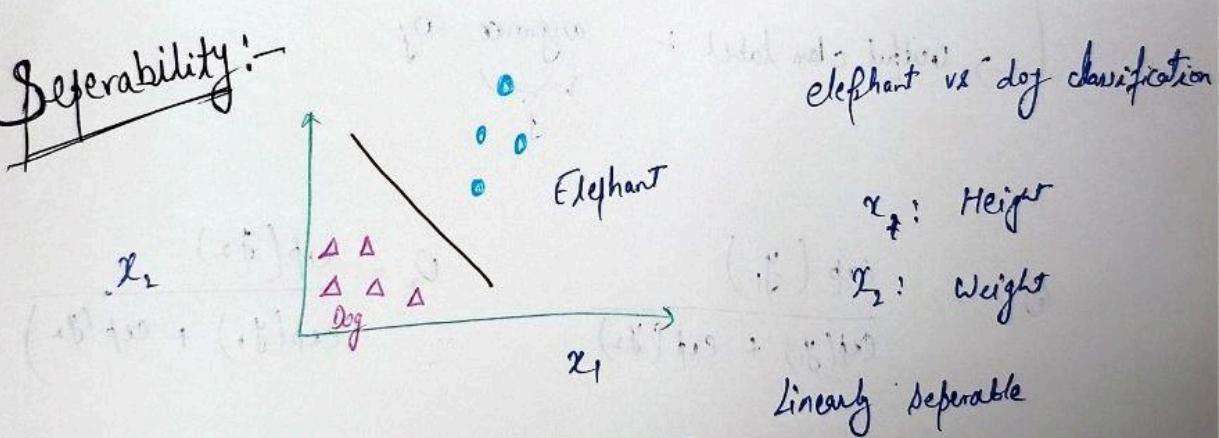
NN can learn each  $\alpha_i$ .

$$y = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3$$

Solution for  $w'$  will give values of  $\alpha_i$

→ One particular dimension will be 1, else all 0.

Separability:-



Class Name	Class Number	Encoding
Elephant	1	1 0
Dog	2	0 1

if n no. of classes

Class 1 :- 1 0 0 0 ... 0

Class 2 :- 0 1 0 0 ... 0

Class n :- 0 0 0 ... 1

①

$$y_1 = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3$$

$$y_2 = \beta_1 x_1 + \beta_2 x_2 + \beta_3$$

There is not guarantee that  $y_i$  will be 0 or 1

So Softmax

$$\rightarrow o_i = \frac{\exp(y_i)}{\sum_j \exp(y_j)}$$

$$\rightarrow \text{output class label} := \underset{j}{\operatorname{argmax}} o_j$$

$$o_1 = \frac{\exp(y_1)}{\exp(y_1) + \exp(y_2)}$$

$$o_2 = \frac{\exp(y_2)}{\exp(y_1) + \exp(y_2)}$$

So on

$o_1, \dots, o_n$  will be range of  $(0,1)$

$\exp$  of anything =  $\frac{\text{true value}}{\text{large term}}$

range  $(0,1)$

\* Each neuron Corresponding O/P  $\rightarrow$  Probability O/P

eg +  $y_1 = 0, 0.3 \rightarrow 0$  (5)



$$\rightarrow 0, y_2 \rightarrow 0, 0.5 \rightarrow 1$$

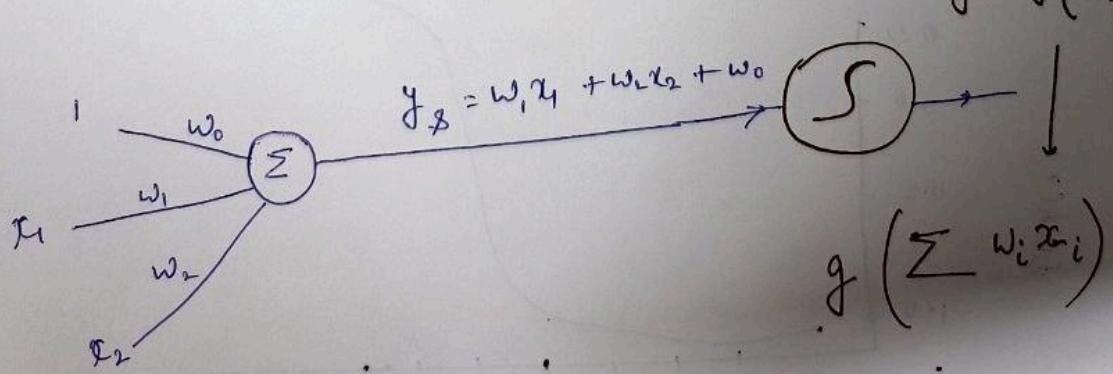
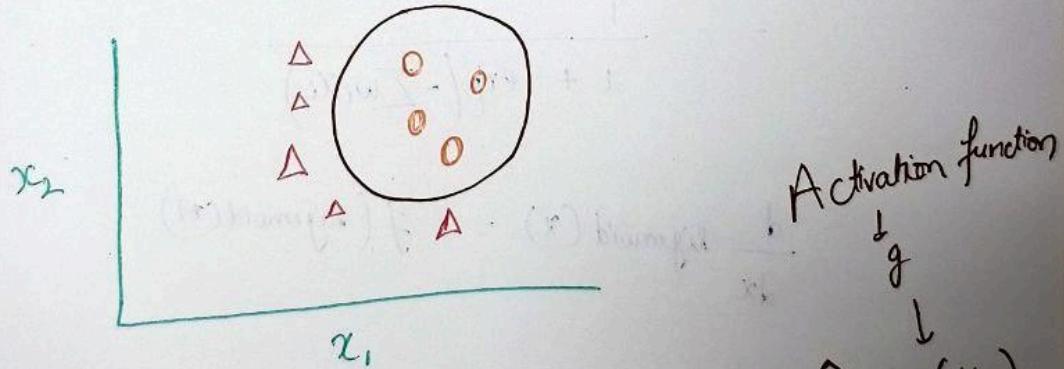
$$\rightarrow 0, y_3 \rightarrow 0, 0.2 \rightarrow 0$$

So 0th class = Class 2.

→ for getting 0th  $\Rightarrow$  either Softmax

$\rightarrow$  or 1 hot encoding

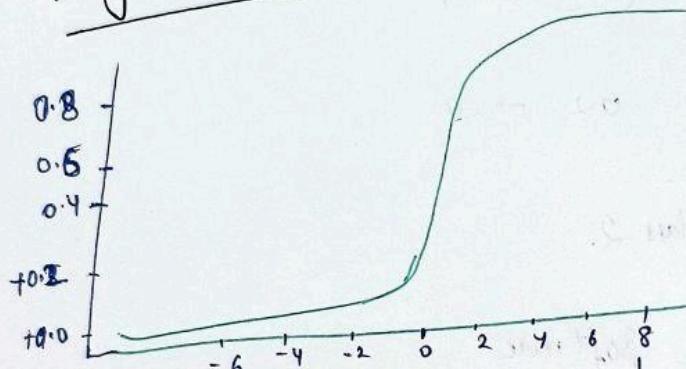
But if data is not linearly separable?



⑥

Another Non-linear fun<sup>n</sup>

Sigmoid



$$\text{Sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

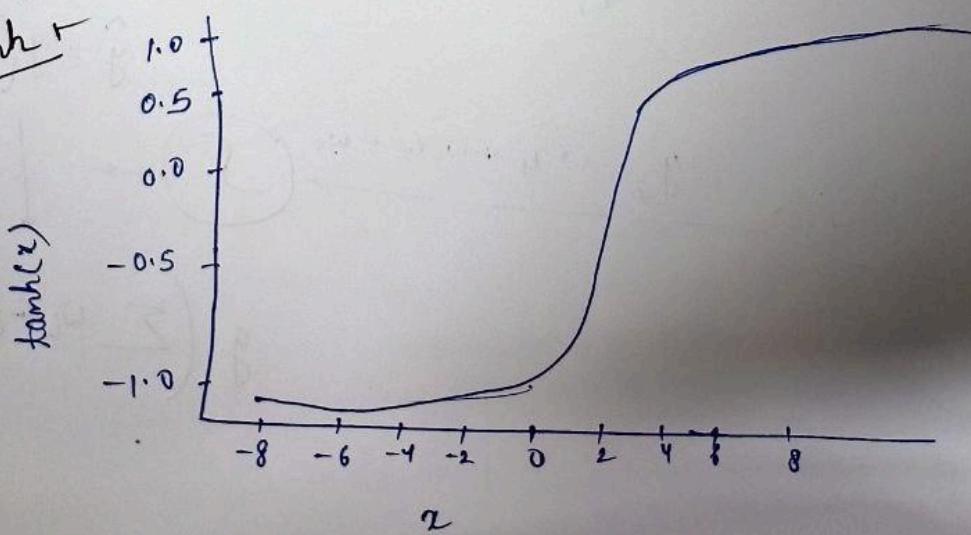
derivative of Sigmoid = also in terms of Sigmoid

$$y_0 = w_1 x_1 + w_2 x_2 + w_0 = \sum_{i=0}^2 w_i x_i$$

$$= \frac{1}{1 + \exp(-\sum w_i x_i)}$$

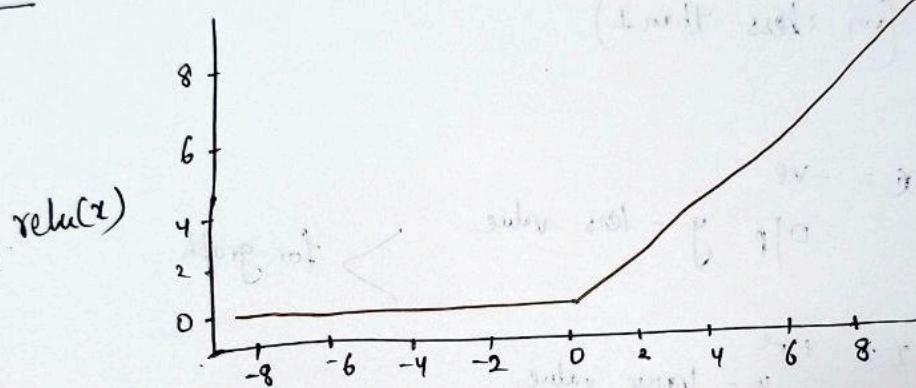
$$\frac{d}{dx} \text{sigmoid}(x) = f(\text{sigmoid}(x))$$

Tanh



$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)} \quad (7)$$

ReLU :-

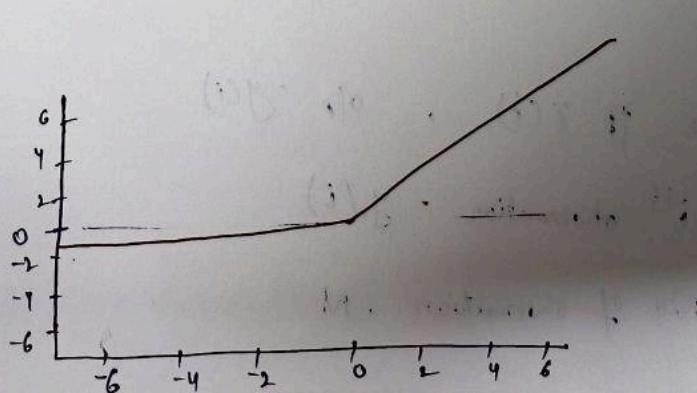


$$\text{ReLU}(x) = \max(0, x)$$

$$\begin{aligned} \text{ReLU}(x) &= x \text{ if } x > 0 \\ &= 0 \text{ otherwise} \end{aligned}$$

$$\text{ReLU}\left(\sum w_i x_i\right) = \begin{cases} \sum w_i x_i & \text{if } \sum w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

Caty ReLU



$$\text{Leaky ReLU}(x) = \max(mx, x) \quad \text{(r)ad.} \quad \textcircled{8}$$

$m$  = gradient of line

$m$  = very small value e.g.  $m = 0.001$

( $m$  less than 1)

if  $x = -ve$

o/p  $y = \text{less value}$

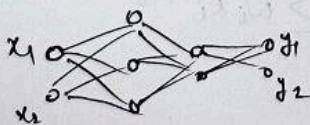
> for graph

if  $x = +ve$

$y = \text{large value}$

By observation we see Leaky ReLU works better than ReLU.

How to measure performance of NN?



Loss

Dif. b/w o/p and

reality (observation)

for ip  $x(i)$ , o/p  $= y(i)$

$i^{th}$  observation  $= y_0(i)$

No. of observation : N

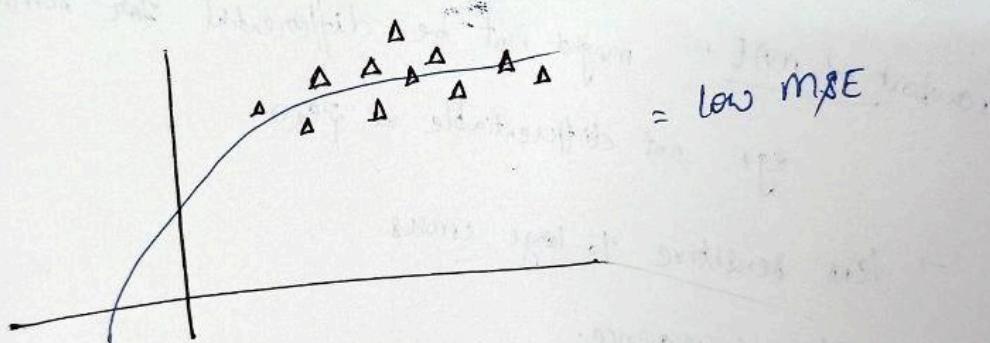
$$MSE = \frac{1}{N} \sum_{i=1}^N (y_{\text{obs}}(i) - y_{\text{pred}}(i))^2$$

(1)

why square?  $\Rightarrow$  pos, -ve loss may cancel out.

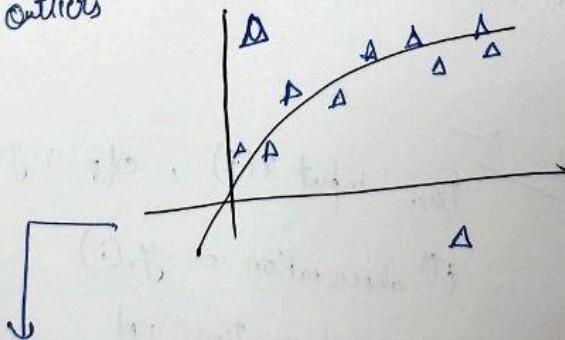
MSE loss is mostly used in Regression.

Case 1: if graph is



= low MSE

if 2 outliers



Overall MSE ↑ even by 2 outliers

high MSE.

Mean Absolute Error

(10)

For i<sup>th</sup> input  $x(i)$ , output =  $y(i)$

i<sup>th</sup> observation =  $y_o(i)$

No. of observations = N

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_o(i) - y(i)|$$

Drawback of MAE → might not be differentiable for some cases.  
e.g. not differentiable at zero.

→ less sensitive to large errors

→ slow convergence.

→ less interpretable for variance

Huber Loss

For input  $x(i)$ , output :  $y(i)$

i<sup>th</sup> observation =  $y_o(i)$

No. of observations : N

$$L_\delta(a) = \begin{cases} \frac{1}{2} a^2 & \text{if } |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta) & \text{if } |a| > \delta \end{cases}$$

$a = y - f(x)$  is residual,  $\delta$  = threshold parameter.

Why huber loss  $\rightarrow$

(1)

- (i) MSE penalizes large errors more heavily due to  
(a) squared term, so sensitive to outliers.
- (b) MAE is less sensitive to outliers as it only considers absolute error.  
Sometime it reduces emphasis on large errors that may need correction.

(ii) Huber loss or it introduces a threshold parameter (delta).

(a) error less than delta,  $\downarrow$  MSE  
penalizes error sharply

(b) error  $>$  threshold  
 $\downarrow$   
MAE

$\rightarrow$  limits influence of outliers.

(iii) (Huber loss) is differentiable everywhere.

(iv) allows model to converge faster than MAE by penalizing error more similar to MSE.

(v) Prevents overfitting of noisy data

- smooth gradients, sensitivity to smaller error (like MSE)
- Robustness against outliers (like MAE)

Cross Entropy Loss

(12)

$$\text{Observed Label} - y_0 = [y_{01}, y_{02}]$$

$$\text{Predicted Label} - y = [y_1, y_2]$$

$$L = - \sum_{i=1}^C y_{0i} \log y_i$$

CEL :- How much similarity b/w predicted and actual value.

e.g.  $y_0 = [y_{01}, y_{02}]$        $y = [y_1, y_2]$

$$\begin{aligned} L &= - (1 \log 0.9 + 0 \log 0.1) \\ &= -\log 0.9 \\ &= -(-0.04) = +ve \text{ value} \end{aligned}$$

confidence in prediction ↑↑  $\Rightarrow$  loss ↓↓

if  $y_1$   $y_2$

$$[0.6 \quad 0.4]$$

$$\begin{aligned} L &= - (1 \log 0.6 + 0.4 \log 0.4) \\ &\approx -\log 0.6 = 0.2 \end{aligned}$$

Confidence in prediction  $\downarrow \downarrow$  = loss  $\uparrow \uparrow$  ③

if  $y_0 = \begin{bmatrix} y_{01} & y_{02} \\ 1 & 0 \end{bmatrix}$   $y = \begin{bmatrix} y_1 & y_2 \\ 0.999 & 0.0001 \end{bmatrix}$

$$L = - \left( 1 \log 0.999 + 0 \log 0.0001 \right)$$
$$= 0.0004$$

if  $y_0 = \begin{bmatrix} 1 & 0 \end{bmatrix}$   $y = \begin{bmatrix} 1 & 0 \end{bmatrix}$

$$L = - \left( 1 \log 1 + 0 \log 0 \right)$$
$$= 0$$

absolute value  $y=1$ ,  $\log=0$

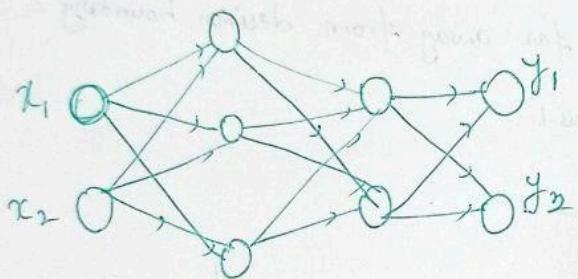
Binary cross entropy = if 2 classes

Multi cross entropy = if more than 2 classes.

Cross-Entropy loss :-

$$\text{Observed label : } y_0 = [y_{01} \ y_{02}]$$

$$\text{Predicted label : } \hat{y} = [\hat{y}_1 \ \hat{y}_2]$$

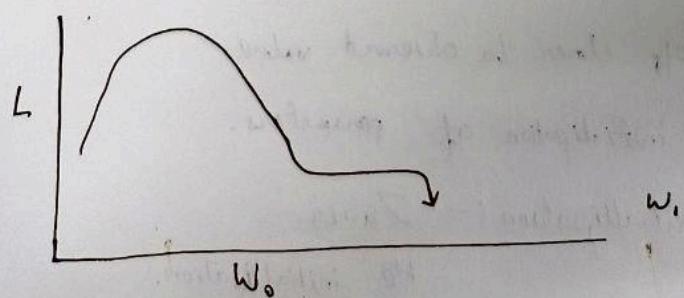


$$L = - \sum_{i=1}^c y_{0i} \log \hat{y}_i = f(\omega)$$

Loss is a function of weights and biases (parameters).

Training of NN :- We can't have an analytical solution of NN with Non-linearity.

→ Goal of training NN :- get  $w, b$  that minimizes loss.



(a) SVM :-  
Can SVM be used with multiclass classification? (2)

Two strategies :-

(a) One versus rest / all classification

(b) One versus group of classes.

e.g. 1 vs (3) group of class

If test data is far away from decision boundary = Class 1.

Less distant not class 1.

Another such cases may be :-

2 | (13)                    3 | (12)

whichever probability is higher → assign.

(b) One Versus One

Probability 1 & arg. cross all Decision boundary.

Probability 2 & arg. cross all D.B.

We want  $\alpha_1$  closer to observed value.

Randomly initialization of parameters.

Technique to initialization :- Xavier

He initialization.

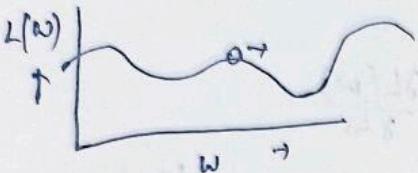
(3)

### How To reach Global Minima :-

if moving to Right :-

- increasing value of  $w$  reduces value of  $L(w)$ .

$$\frac{\partial L(w)}{\partial w} < 0$$

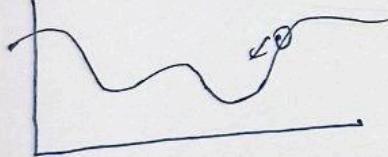


To increase value of  $w$  :-

$$w \leftarrow w - \eta \frac{\partial L(w)}{\partial w}$$

$\eta$ : positive constant

if moving to left :-



$$\frac{\partial L(w)}{\partial w} > 0$$

To decrease value of  $w$  :-

$$w \leftarrow w - \eta \frac{\partial L(w)}{\partial w}$$

$\eta$  = pre. constant

When to stop updating parameter?

a) If decreasing of loss becomes constants

Loss Saturates  $\leftarrow$  Stop at that point

b) Pre-decide no of epoch to run loop for updating loss.

c) Try and error (on validation data first).

## Gradient Descent

⑨

- (i) initialize parameters
- (ii) loop until convergence

(a) Compute gradient  $\frac{\partial L(w)}{\partial w}$

(b) update parameters  $w \leftarrow w - \eta \frac{\partial L(w)}{\partial w}$

(iii) Return parameters

This is called (Hyper Parameter tuning)

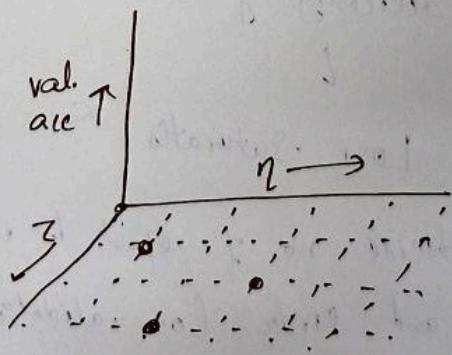
Best possible technique is Grid Search

$\eta$ , learning rate,  $T$  epoch.

$$\eta = \begin{cases} 0.001 \\ 50 \end{cases} \rightarrow 0.0001 \text{ range of values} \quad | \quad 100$$

$$\begin{array}{l} \eta_1 = 0.0001 \\ T_1 = 50 \end{array} \quad | \quad \begin{array}{l} \eta_2 = 0.0002 \\ T_2 = 55 \end{array} \quad | \quad \begin{array}{l} \eta_3 = 0.00019 \\ T_3 = 55 \end{array}$$

try diff. diff config of parameter to find best one from grid.

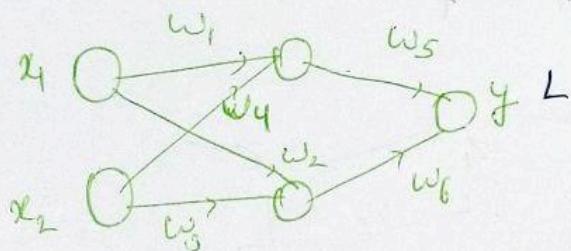


• can do grid search on even 6 dimensional data. (5)

Parameter : which N.N learns by gradients

Hyperparameter : user set by diff mechanism.

### Computing Gradients : Backpropagation +



partial derivative here:

$$\frac{\text{Loss}}{\text{diff } w} = \text{partial}$$

$$w_6 \leftarrow \frac{w_6}{\text{old value}} - \alpha \cdot \frac{\partial L}{\partial w_6}$$

To update  $w_5 \leftarrow$

$$w_5 \leftarrow w_5 - \alpha \cdot \frac{\partial L}{\partial w_5}$$

$$L = f(y, y_0) + \text{observed value constant}$$

$$L = f(y)$$

$$y = g(z_1 w_5 + z_2 w_6)$$

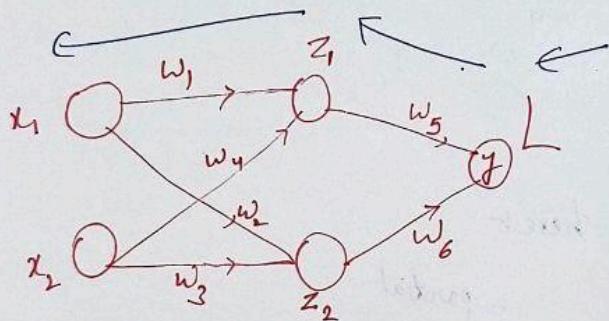
From last layer I can calculate  $w_5, w_6$  last layers  
Connecting parameters. ⑥

- Can't directly

$$z_1 = (x, w_1, w_4, x_2)$$

$$z_2 = (w_2, w_3, x_1, x_2)$$

$z_1, z_2$  are sum of above parameters



Backpropagation :- To update  $w_1$ , we'll do :-

$$w_1 \leftarrow w_1 - \eta \frac{\partial L}{\partial w_1}$$

$$w_1 \leftarrow w_1 - \eta \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

chain rule

loss function and output of nodes needs to be differentiable.

It is not correct

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_1} \frac{\partial z_1}{\partial w_1} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial w_1}$$

$\downarrow$   
no connection.

(2)

$$\frac{\partial L}{\partial y} \quad \frac{\partial y}{\partial z_1} \quad \frac{\partial z_1}{\partial w_1}$$

Backpropagation :- It automatic calculation looks like propagating.

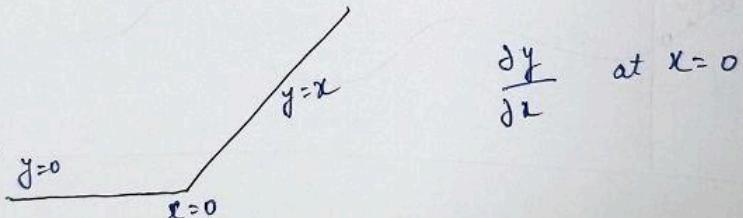
only forward propagation there.

At each epoch :- every parameter  $w, b$  update.

- each of these "fun" has to be differentiable.

- Loss "fun" and o/p of nodes need to be differentiable.

ReLU is not differentiable at 0



$$\frac{\partial y}{\partial x} \text{ at } x=0$$

$$\frac{\partial y}{\partial x} = 1, \frac{\partial y}{\partial x} = 0 \text{ at -ve}, \quad \frac{1+0}{2} = \frac{1}{2} \text{ possibility approximated}$$

at +ve

$$f(x+\delta)(x-\delta) \leftarrow \text{differential now}$$

Optimizing loss :- may be challenging because of

- local minima
- Vanishing gradients
- Saddle points

- local minima how to find Optimum local minima.

(C)

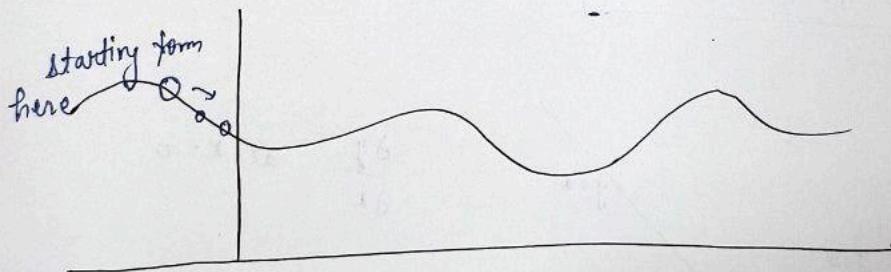
Gradient may 0 at local minima.

→ value of gradient = 0 at saddle point. → Training stops

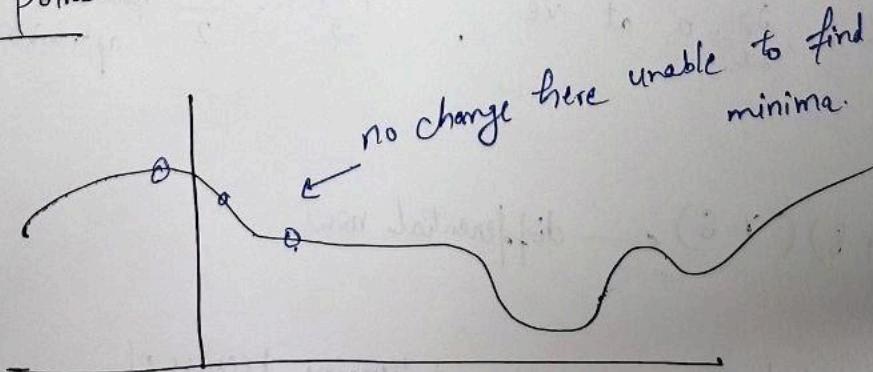
[ Most NN saddle point predominant than  
local minima. ]

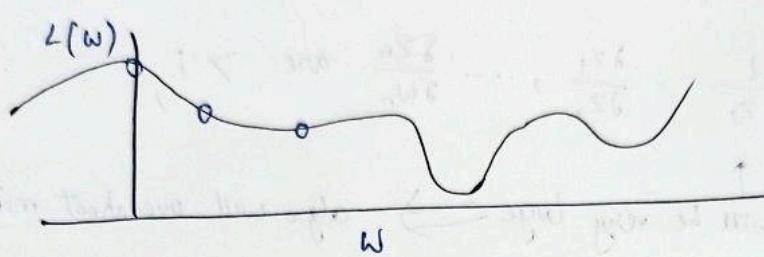
→ Mostly we try to get rid of saddle point.

Local Minima ↗



Saddle Point





$$\frac{\partial L}{\partial w_i} = 0 \text{ its minima}$$

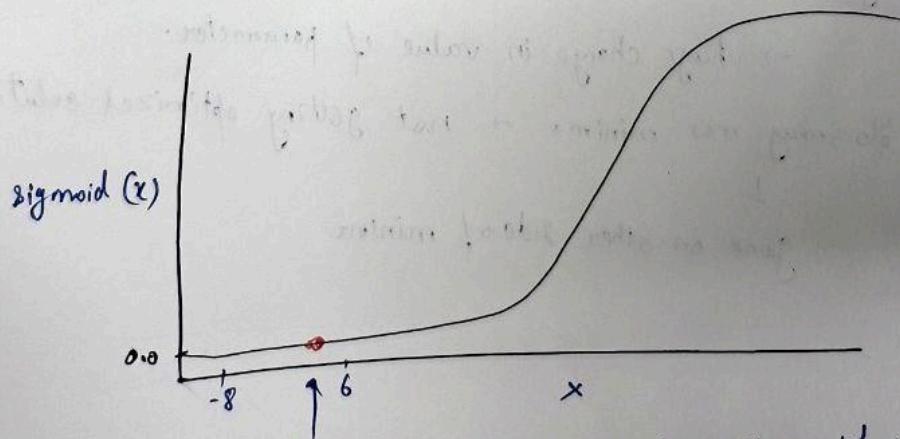
$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i}$$

w<sub>i</sub>, doesn't change further, same situation happens for other parameters.

Vanishing Gr. D +  $w_i \leftarrow w_i - \eta \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial w_i}$

for  $w_i$  to change  $\frac{\partial L}{\partial z_i}$  and  $\frac{\partial L}{\partial w_i}$  must be non-zero

$$w_n \leftarrow w_n - \eta \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial z_2} \dots \frac{\partial z_n}{\partial w_n}$$



for any of nodes related to derivatives, if we are here, derivative term will be zero.

for  $w_i$  to change, each of  $\frac{\partial L}{\partial z_1}, \frac{\partial z_1}{\partial z_2}, \dots, \frac{\partial z_n}{\partial w_n}$  must be non-zero. ②

Q) If each  $\frac{\partial L}{\partial z_1}, \frac{\partial z_1}{\partial z_2}, \dots, \frac{\partial z_n}{\partial w_n}$  are  $> 1$ ,  
product will be very large  $\Rightarrow$  algo will overshoot minima.

b) for vanishing Cr.D = product  $< 1$  product smaller.

c) If any of DNN one term so  
overall update  $= 0$  = vanishing Crd.

To prevent vanishing Cr.D +

(a) residual connection

(b) skip connection.

Exploding Cr.D + each differential product has very large value.

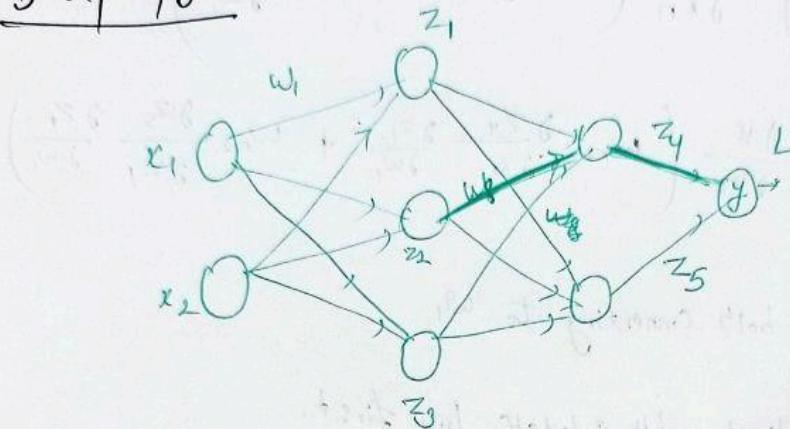
$\rightarrow$  huge change in value of parameter.

So may miss minima  $\rightarrow$  not getting optimized solution.

$\downarrow$   
gone on other side of minima.

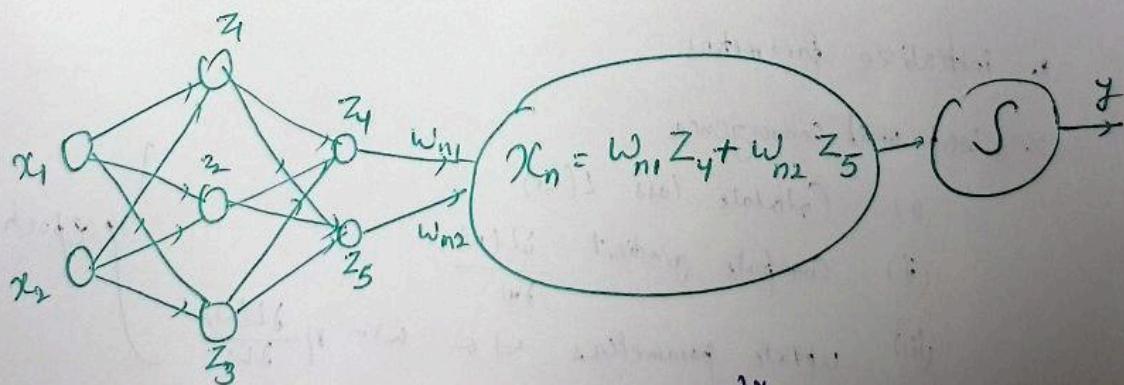
③

### Backpropagation



$$\frac{\partial L}{\partial w_8} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_4} \frac{\partial z_4}{\partial w_8}$$

$z_1$  impacting  $z_4, z_5$  both considering path.  
 $\rightarrow$  every node has same probability no vanishing gradient here.



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x_n} \frac{\partial x_n}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x_n} \frac{\partial}{\partial w_1} \left( \frac{w_{n1}z_4 + w_{n2}z_5}{\partial w_1} \right)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x_n} \left( w_{n1} \frac{\partial z_4}{\partial w_1} + w_{n2} \frac{\partial z_5}{\partial w_1} \right) \quad \textcircled{1}$$

$$\Rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x_n} \left( w_{n1} \frac{\partial z_4}{\partial z_1} \frac{\partial z_1}{\partial w_1} + w_{n2} \frac{\partial z_5}{\partial z_1} \frac{\partial z_1}{\partial w_1} \right)$$

here  $z_4, z_5$  both connecting to  $w_1$

- To calculate gradient we'll calculate loss first.

### Gradient Descent: challenges

Suppose 10 training data  
 [ if 10 million training data → to update every parameter  
 to pass through nn is very hard ]

1. initialize parameters

2. loop until convergence

- Calculate loss  $L(w)$
  - Compute gradient  $\frac{\partial L(w)}{\partial w}$
  - update parameters  $w \leftarrow w - \eta \frac{\partial L(w)}{\partial w}$
- } epoch

3. Return parameters.

In Batch GD with  $N$  samples (data points)

$$L(w) = \frac{1}{N} \sum_{i=1}^N L_i(w)$$

$L_i(w)$ : loss for  $i$ th data point (3)

$$L(w) = E(L_i(w)) = \frac{1}{N} \sum_{i=1}^N L_i(w)$$

So if we calculate loss for each sample  
↓  
and take expected value  $\Rightarrow$  Loss for batch

I calculate loss for one sample =  $L_j(w)$

[ do parameter update

[ do it again for another sample = continue doing for each individual samples

same impact on parameters

→ Each update would require only one gradient computation

(less time, less storage for gradient values)

better than this

### Stochastic G.D.

1. initialize parameters

2. loop until convergence

{ i) Randomly shuffle samples in training dataset

ii) for training sample  $i = 1, 2, 3, \dots, N$

(a) compute gradient  $\frac{\partial L_i(w)}{\partial w}$

(b) update parameters  $w \leftarrow w - \eta \frac{\partial L_i(w)}{\partial w}$

iii) Return parameters

epoch.

BGD ↗ not data efficient when data is similar  
↗ high computational time for each step

SGD + [ not computationally efficient  
  ↗ noisy result since at a time, one data point  
    is considered.  
    ↓ solution

### MiniBatch GD +

1. Initialise parameters
2. Loop until convergence
  - (i) Randomly shuffle samples in training dataset
  - (ii) From training dataset, Create b minibatches of size m
  - (iii) for each minibatch  $i = 1, 2, \dots, b$ 
    - (a) Compute gradient  $\frac{\partial L_i(w)}{\partial w}$
    - (b) Update parameters  
 $w \leftarrow w - \eta \frac{\partial L_i(w)}{\partial w}$
3. Return parameters

Adaptive learning Rate + L.R :- How much gradient update at every step.

Do we want complete or shorter step update.

→ Larger step might miss Minima.

→ Smaller takes longer period of time.

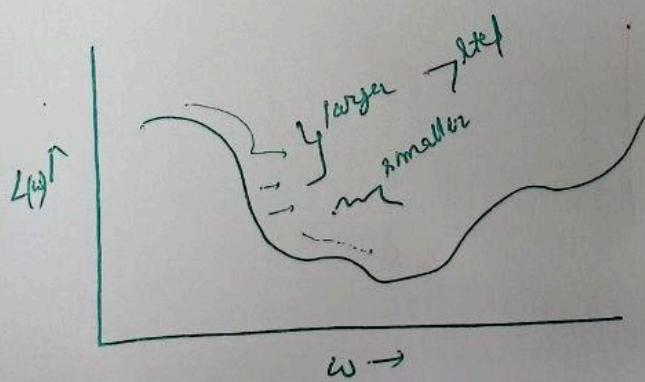
There should be some trade-off for learning rate.

Time dependent learning rate :-

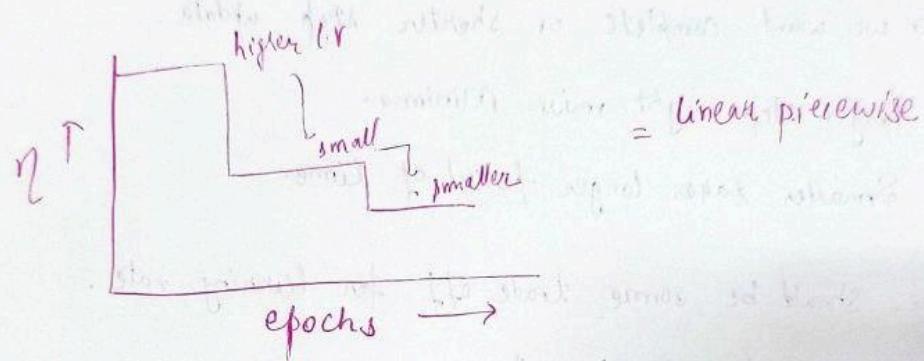
$$\begin{cases} \eta(t) = \eta_i & \text{if } t_i \leq t \leq t_{i+1} & \text{piecewise constant} \\ \eta(t) = \eta_0 \cdot e^{-\lambda t} & \text{exponential decay} \\ \eta(t) = \eta_0 \cdot (\beta t + 1)^{-\alpha} & \text{polynomial decay} \end{cases}$$

initially  $\downarrow r \rightarrow$  larger step

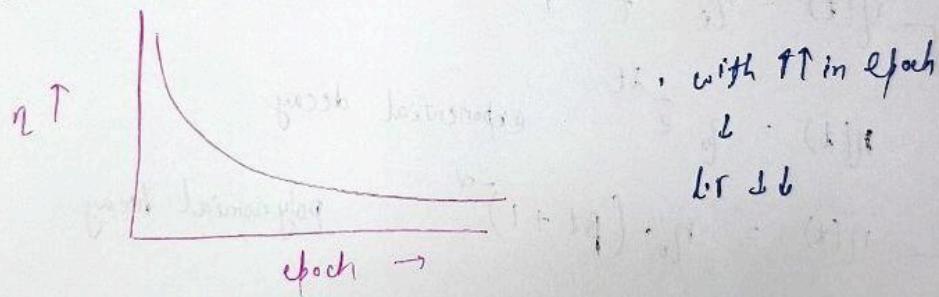
near minima  $\rightarrow$  smaller step.



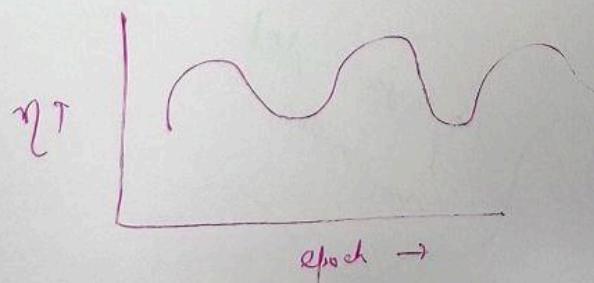
(a) Interval wise !  
 egt  $10^3, 10, 1, \dots, 10^{-2}$   
 such kind of lr in intervals.



(b) exponential



(c) some times Cosine LR



Why NN :- It can be applied to all fields / problems. ③

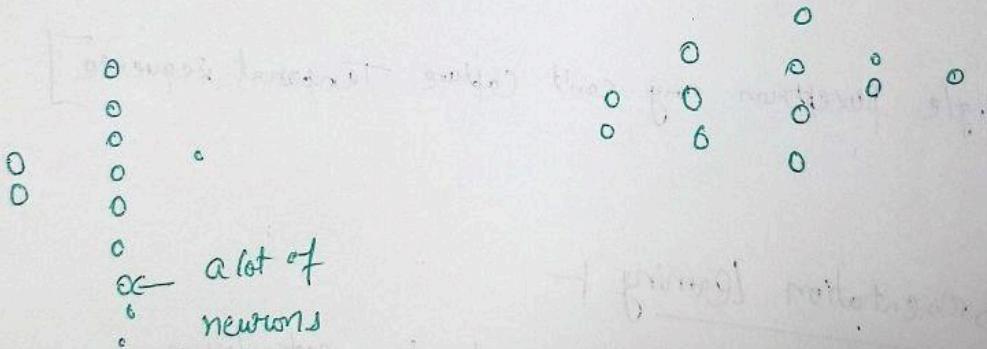
→ something applicable and gives good result across different problem.

- In ML we mostly try to find curve.

### Universal Approximation Theorem

→ Even with single layer of NN, NN can approximate any function with arbitrary precision.

Single layer → a lot of neuron



Multilayer Perception (MLP)  $\xrightarrow{\text{In}}$  node of previous layer connected to all future.

(I/p layer hidden layer o/p layer)

• Single layer may need huge no. of nodes.

having that much wider n/w is very hard.

• Better to take deeper n/w with less wider n/w.

UAT & A feedforward n/w with linear o/p layer. (Q)  
and atleast 1 hidden layer with any "squashing"  
activation fun" i.e logistic sigmoid activation fun.

↓

can approximate any "fun"

↓

from one finite dimensional space to another

↓

with any desired non-zero amount of error

↓

grauded n/w is given enough hidden units.

[single perception can't capture Temporal sequence]

Representation Learning &

FNN data is represented in vector form.

→ Pixel wise operation might not give good result of image.

Case Based on prior knowledge

↓  
Extracting handcrafted feature is hard.

Hand crafted feature are useful in ML if correct  
features are extracted.

It's often difficult to know correct features, (5)

(specially images, audios, videos etc)

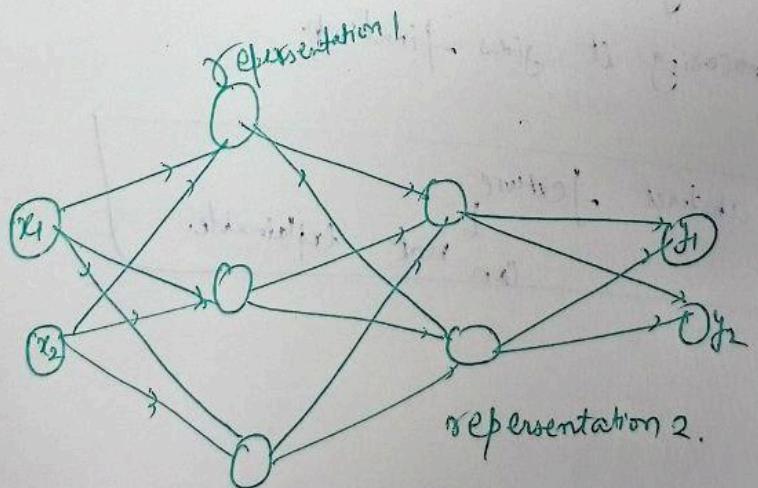
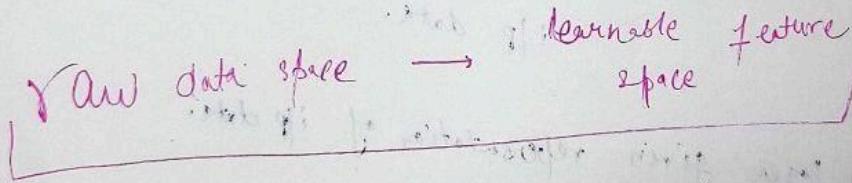
→ Various NN (RNN, CNN etc) not only map i/p data to o/p

↳ but also learn a representation of data in various layers.

Representation learning & don't have to rely on prior knowledge.

apply raw data, model itself will extract feature

→ Instructing meaningful data



$x_1, x_2$  = 2 features.

Linearly Combining feat<sup>n</sup> → Activation fun

- Initial layers learn low-level explainable features  
(e.g. edge, corner, texture etc)
  - Deeper layers learn high-level abstract features  
(by combining suitable low-level features)
- Top node of layer = getting some form of modified / derived feature

layer 2 → getting more derived feature

layer 3 → getting more

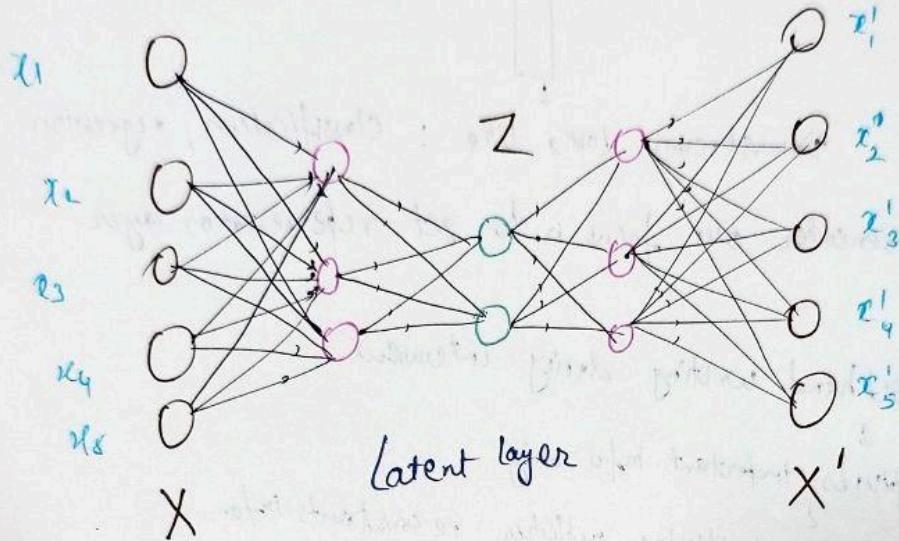
[ Desired feature vector + Modified representation of raw ip data. ]

→ Each layer given representation of ip data.  
Processing it gives final op.

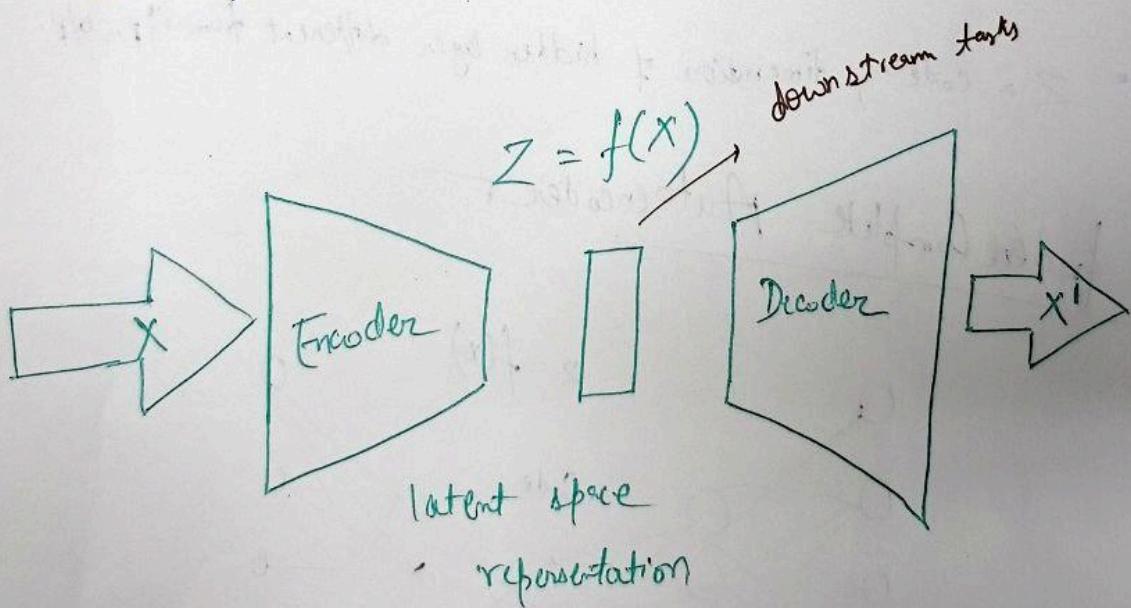
abstract feature  
↓  
can not explainable.

Auto encoders ↪

$$x' = g(z) \quad \text{⑦}$$



$Z = f(x)$  : Latent space representation / code



Encoder creates code  $Z$  from  $X$

Decoder reconstructs  $X$  from  $Z$ .. so  $X' = X$  ideally

(8)

Representation space  $\leftarrow z = f(x)$



Downstream tasks like : classification, regression

In Autoencoder our focus is to get representation layer.

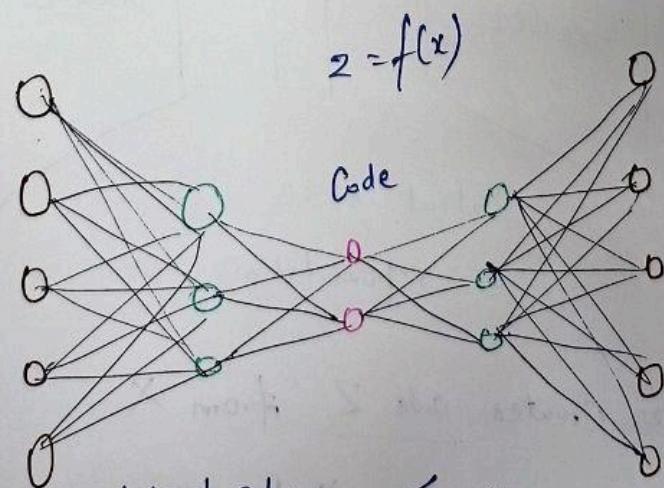
e.g. shorthand writing during interview

captures important info only  
during publishing publisher reconstructs info.

Having good  $Z$  only  $\rightarrow$  can get  $x'$

- $Z = \text{Code}$ , dimension of hidden layer different from i/p, o/p.

UnderComplete Autoencoder



Dimension of latent space  $<$  dimension of i/p/o/p space.

User  $\downarrow$  may be used for dimensionality reduction / compression of data. ①

(ii) Extraction of salient features

$$\boxed{\text{Loss} = \text{MSE } (x, x')}$$

- Minimization of this loss requires  $x'$  to be as similar as to  $x$ .
- In order to do,  $Z$  must contain salient features.

Extraction of Salient features

e.g. To diff among zebra and dog we need combination of features alone 1 feature can't differentiate them <sup>may be linear</sup>

$x_1$ : weight

$x_3$ : Major color

$x_2$ : Height

$x_4$ : Major color 2

$x_5$ : Repetition of color pattern (yes/no)

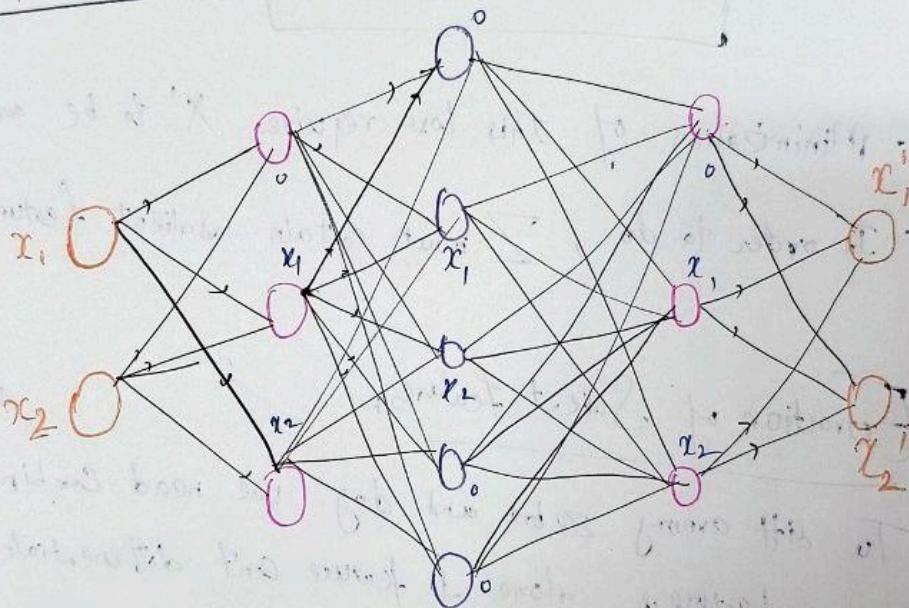
→ We are not using label here.

Self Supervised learning - data itself tackling/ data teaching

- The blending of features is created by AE in latent space.
- latent space creates a salient representation of data.

### OverComplete Autoencoders

$$\text{Code } z = f(x)$$



Dimension of latent space  $>$  dimension of i/p, o/p space

latent space ↑↑

- It might copy data
- might not learn anything new      then why we need it?

(1)

(i) Better feature learning

→ more neurons in hidden layer than i/p forces to learn  
more hierarchical, abstract or sparse representation of data.

→ Rather than directly copying data

↓

It learns to represent data in more informative,  
compressed form that captures underlying structure.

(ii) Regularization, Sparsity +

By L<sub>1</sub> regularization, sparsity autoencoders where only a few  
units in hidden layer are active at any given time.

→ It helps in reducing overfitting, forces to learn more compact,  
efficient features.

(iii) Non-linearities and complex representations +

It can learn complex, non linear mapping from i/p to hidden  
layer

→ It allows to create a higher-dimensional representation of data,  
making it more effective, generalized.

(iv) Improved generative modeling

(v) Generalization to unseen data

