

UNIT 2

Random forest.

- Working
- Advantages
- Disadvantages.

steps

RANDOM → bootstrapping → choosing small amt of data from original all data
(without replacement)

⇒ feature selection.

$x_1, x_2, x_3, x_4 \rightarrow$ columns of a dataset

$\sqrt{4} = 2$ features are chosen.

→ different decision trees

FOREST → collection of trees lead to forest

⇒ many trees may encounter different errors / mistakes and so this will help to learn and know to adjust each of behaviour, so performance improved.

Leo Breiman → proposed Random forest.

i) Random

2 factors.

ii) Bootstrapping → analysis to choose smaller subset of entire dataset

Randomness leads to different decision tree.

each decision tree has different type of mistakes.

③ Forest
collection of different decision trees

Eg: Assume dataset

100 training examples = N

$$\gamma_1, \gamma_2, \gamma_3, \dots, \gamma_n = 100$$

↓
subsets of dataset
30 30.

training examples
because of bootstrapping

WORKING:

① Building Random forest

② It is ensemble learning that is each decision tree is a weak learner but group of weak learning decision trees will give better performance.

③ Assume N training examples from which $\gamma_1, \gamma_2, \dots, \gamma_n$ training examples are chosen randomly.

④ Assume M features out of which m features are chosen

⑤ The decision trees are constructed based on finding the best split node.

⑥ When the new data input arrives it is passed through each of the decision tree in Random forest, then each

tree predicts the label, majority label will be chosen as class of new ^{data} input.

first step is Bootstrapping
second step is Aggregation } together called as
BAGGING

⇒ Regression using Random forest

step (A) remains same

in step (B), it is averaging instead of aggregation.

↓
majority

• Advantages:

① less training time even for larger dataset

② Accuracy is more in cases.

③ large dataset

④ used both for classification and regression

⑤ Noisy data

• Disadvantages:

① memory requirement more. (many DTs).

② blackbox (internal factors of DT, we can't get to know).

③ output is limited to range present in training dataset.

26/11 • Unsupervised learning

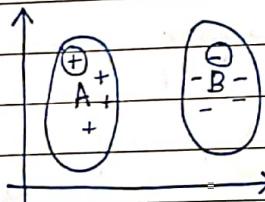
⇒ obtaining the data - planning

⇒ cluster analysis.

- ① no overlap (one & only one cluster)
- ② centroid (every cluster must have a centroid)
(Avg of numeric values in training set)
- ③ similarity within the cluster compared to
data points in other clusters.

→ obtaining useful information from existing data.
→ Principle of unsupervised learning.

cluster analysis Algorithm:



This is a 2 dimensional visualization and is having 2 clusters beyond 3 dimensions. It is difficult for humans to visualize the data and identify the clusters. Therefore cluster analysis algorithm are required.

2 different approaches: of unsupervised learning

- ① clustering
- ② Association (minimal dependency b/w data).

K-Means algorithm

I Input

- ① } Body of algorithm
- ② } (4 substeps)
- ③ }
- ④ }

II Termination

III Output

• INPUT

- 1) set of training egs without labelling
 - 2) no of clusters (no of K values)
- K values → : ① user specified (constant) ② ML specified automatically.
(Software ML generate automatically).

• BODY.

- ① create K clusters, which are initially empty.
calculate the centroid of each cluster
- ② consider a training example x , represented in the form of vector (numerical representation of data is called vector).
calculate the distance of this training eg x from
centroid of each cluster. (clusters)
- Let J denote the nearest cluster to the training example x .
- ③ Move x to J , if x is already in J , do nothing else
separate that x from that cluster to J .
Recalculate the new centroids of clusters.
- ④ repeat steps ② and ③ till termination or stopping

condition is met.

III Every training example is in its meant cluster. no relocation is required. - TERMINATION.

IV OUTPUT.

K clusters with similar examples.

→ Disadvantage: supplying K value is tricky

→ Advantage: simple & easy to understand.

- Euclidean distance - used to calculate the distance b/w training example and centroids of clusters.

$$d(x, y) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

		x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}	x_{19}	x_{20}	x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}	x_{29}	x_{30}	x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}	x_{39}	x_{40}	x_{41}	x_{42}	x_{43}	x_{44}	x_{45}	x_{46}	x_{47}	x_{48}	x_{49}	x_{50}	x_{51}	x_{52}	x_{53}	x_{54}	x_{55}	x_{56}	x_{57}	x_{58}	x_{59}	x_{60}	x_{61}	x_{62}	x_{63}	x_{64}	x_{65}	x_{66}	x_{67}	x_{68}	x_{69}	x_{70}	x_{71}	x_{72}	x_{73}	x_{74}	x_{75}	x_{76}	x_{77}	x_{78}	x_{79}	x_{80}	x_{81}	x_{82}	x_{83}	x_{84}	x_{85}	x_{86}	x_{87}	x_{88}	x_{89}	x_{90}	x_{91}	x_{92}	x_{93}	x_{94}	x_{95}	x_{96}	x_{97}	x_{98}	x_{99}	x_{100}	x_{101}	x_{102}	x_{103}	x_{104}	x_{105}	x_{106}	x_{107}	x_{108}	x_{109}	x_{110}	x_{111}	x_{112}	x_{113}	x_{114}	x_{115}	x_{116}	x_{117}	x_{118}	x_{119}	x_{120}	x_{121}	x_{122}	x_{123}	x_{124}	x_{125}	x_{126}	x_{127}	x_{128}	x_{129}	x_{130}	x_{131}	x_{132}	x_{133}	x_{134}	x_{135}	x_{136}	x_{137}	x_{138}	x_{139}	x_{140}	x_{141}	x_{142}	x_{143}	x_{144}	x_{145}	x_{146}	x_{147}	x_{148}	x_{149}	x_{150}	x_{151}	x_{152}	x_{153}	x_{154}	x_{155}	x_{156}	x_{157}	x_{158}	x_{159}	x_{160}	x_{161}	x_{162}	x_{163}	x_{164}	x_{165}	x_{166}	x_{167}	x_{168}	x_{169}	x_{170}	x_{171}	x_{172}	x_{173}	x_{174}	x_{175}	x_{176}	x_{177}	x_{178}	x_{179}	x_{180}	x_{181}	x_{182}	x_{183}	x_{184}	x_{185}	x_{186}	x_{187}	x_{188}	x_{189}	x_{190}	x_{191}	x_{192}	x_{193}	x_{194}	x_{195}	x_{196}	x_{197}	x_{198}	x_{199}	x_{200}	x_{201}	x_{202}	x_{203}	x_{204}	x_{205}	x_{206}	x_{207}	x_{208}	x_{209}	x_{210}	x_{211}	x_{212}	x_{213}	x_{214}	x_{215}	x_{216}	x_{217}	x_{218}	x_{219}	x_{220}	x_{221}	x_{222}	x_{223}	x_{224}	x_{225}	x_{226}	x_{227}	x_{228}	x_{229}	x_{230}	x_{231}	x_{232}	x_{233}	x_{234}	x_{235}	x_{236}	x_{237}	x_{238}	x_{239}	x_{240}	x_{241}	x_{242}	x_{243}	x_{244}	x_{245}	x_{246}	x_{247}	x_{248}	x_{249}	x_{250}	x_{251}	x_{252}	x_{253}	x_{254}	x_{255}	x_{256}	x_{257}	x_{258}	x_{259}	x_{260}	x_{261}	x_{262}	x_{263}	x_{264}	x_{265}	x_{266}	x_{267}	x_{268}	x_{269}	x_{270}	x_{271}	x_{272}	x_{273}	x_{274}	x_{275}	x_{276}	x_{277}	x_{278}	x_{279}	x_{280}	x_{281}	x_{282}	x_{283}	x_{284}	x_{285}	x_{286}	x_{287}	x_{288}	x_{289}	x_{290}	x_{291}	x_{292}	x_{293}	x_{294}	x_{295}	x_{296}	x_{297}	x_{298}	x_{299}	x_{300}	x_{301}	x_{302}	x_{303}	x_{304}	x_{305}	x_{306}	x_{307}	x_{308}	x_{309}	x_{310}	x_{311}	x_{312}	x_{313}	x_{314}	x_{315}	x_{316}	x_{317}	x_{318}	x_{319}	x_{320}	x_{321}	x_{322}	x_{323}	x_{324}	x_{325}	x_{326}	x_{327}	x_{328}	x_{329}	x_{330}	x_{331}	x_{332}	x_{333}	x_{334}	x_{335}	x_{336}	x_{337}	x_{338}	x_{339}	x_{340}	x_{341}	x_{342}	x_{343}	x_{344}	x_{345}	x_{346}	x_{347}	x_{348}	x_{349}	x_{350}	x_{351}	x_{352}	x_{353}	x_{354}	x_{355}	x_{356}	x_{357}	x_{358}	x_{359}	x_{360}	x_{361}	x_{362}	x_{363}	x_{364}	x_{365}	x_{366}	x_{367}	x_{368}	x_{369}	x_{370}	x_{371}	x_{372}	x_{373}	x_{374}	x_{375}	x_{376}	x_{377}	x_{378}	x_{379}	x_{380}	x_{381}	x_{382}	x_{383}	x_{384}	x_{385}	x_{386}	x_{387}	x_{388}	x_{389}	x_{390}	x_{391}	x_{392}	x_{393}	x_{394}	x_{395}	x_{396}	x_{397}	x_{398}	x_{399}	x_{400}	x_{401}	x_{402}	x_{403}	x_{404}	x_{405}	x_{406}	x_{407}	x_{408}	x_{409}	x_{410}	x_{411}	x_{412}	x_{413}	x_{414}	x_{415}	x_{416}	x_{417}	x_{418}	x_{419}	x_{420}	x_{421}	x_{422}	x_{423}	x_{424}	x_{425}	x_{426}	x_{427}	x_{428}	x_{429}	x_{430}	x_{431}	x_{432}	x_{433}	x_{434}	x_{435}	x_{436}	x_{437}	x_{438}	x_{439}	x_{440}	x_{441}	x_{442}	x_{443}	x_{444}	x_{445}	x_{446}	x_{447}	x_{448}	x_{449}	x_{450}	x_{451}	x_{452}	x_{453}	x_{454}	x_{455}	x_{456}	x_{457}	x_{458}	x_{459}	x_{460}	x_{461}	x_{462}	x_{463}	x_{464}	x_{465}	x_{466}	x_{467}	x_{468}	x_{469}	x_{470}	x_{471}	x_{472}	x_{473}	x_{474}	x_{475}	x_{476}	x_{477}	x_{478}	x_{479}	x_{480}	x_{481}	x_{482}	x_{483}	x_{484}	x_{485}	x_{486}	x_{487}	x_{488}	x_{489}	x_{490}	x_{491}	x_{492}	x_{493}	x_{494}	x_{495}	x_{496}	x_{497}	x_{498}	x_{499}	x_{500}	x_{501}	x_{502}	x_{503}	x_{504}	x_{505}	x_{506}	x_{507}	x_{508}	x_{509}	x_{510}	x_{511}	x_{512}	x_{513}	x_{514}	x_{515}	x_{516}	x_{517}	x_{518}	x_{519}	x_{520}	x_{521}	x_{522}	x_{523}	x_{524}	x_{525}	x_{526}	x_{527}	x_{528}	x_{529}	x_{530}	x_{531}	x_{532}	x_{533}	x_{534}	x_{535}	x_{536}	x_{537}	x_{538}	x_{539}	x_{540}	x_{541}	x_{542}	x_{543}	x_{544}	x_{545}	x_{546}	x_{547}	x_{548}	x_{549}	x_{550}	x_{551}	x_{552}	x_{553}	x_{554}	x_{555}	x_{556}	x_{557}	x_{558}	x_{559}	x_{560}	x_{561}	x_{562}	x_{563}	x_{564}	x_{565}	x_{566}	x_{567}	x_{568}	x_{569}	x_{570}	x_{571}	x_{572}	x_{573}	x_{574}	x_{575}	x_{576}	x_{577}	x_{578}	x_{579}	x_{580}	x_{581}	x_{582}	x_{583}	x_{584}	x_{585}	x_{586}	x_{587}	x_{588}	x_{589}	x_{590}	x_{591}	x_{592}	x_{593}	x_{594}	x_{595}	x_{596}	x_{597}	x_{598}	x_{599}	x_{600}	x_{601}	x_{602}	x_{603}	x_{604}	x_{605}	x_{606}	x_{607}	x_{608}	x_{609}	x_{610}	x_{611}	x_{612}	x_{613}	x_{614}	x_{615}	x_{616}	x_{617}	x_{618}	x_{619}	x_{620}	x_{621}	x_{622}	x_{623}	x_{624}	x_{625}	x_{626}	x_{627}	x_{628}	x_{629}	x_{630}	x_{631}	x_{632}	x_{633}	x_{634}	x_{635}	x_{636}	x_{637}	x_{638}	x_{639}	x_{640}	x_{641}	x_{642}	x_{643}	x_{644}	x_{645}	x_{646}	x_{647}	x_{648}	x_{649}	x_{650}	x_{651}	x_{652}	x_{653}	x_{654}	x_{655}	x_{656}	x_{657}	x_{658}	x_{659}	x_{660}	x_{661}	x_{662}	x_{663}	x_{664}	x_{665}	x_{666}	x_{667}	x_{668}	x_{669}	x_{670}	x_{671}	x_{672}	x_{673}	x_{674}	x_{675}	x_{676}	x_{677}	x_{678}	x_{679}	x_{680}	x_{681}	x_{682}	x_{683}	x_{684}	x_{685}	x_{686}	x_{687}	x_{688}	x_{689}	x_{690}	x_{691}	x_{692}	x_{693}	x_{694}	x_{695}	x_{696}	x_{697}	x_{698}	x_{699}	x_{700}	x_{701}	x_{702}	x_{703}	x_{704}	x_{705}	x_{706}	x_{707}	x_{708}	x_{709}	x_{710}	x_{711}	x_{712}	x_{713}	x_{714}	x_{715}	x_{716}	x_{717}	x_{718}	x_{719}	x_{720}	x_{721}	x_{722}	x_{723}	x_{724}	x_{725}	x_{726}	x_{727}	x_{728}	x_{729}	x_{730}	x_{731}	x_{732}	x_{733}	x_{734}	x_{735}	x_{736}	x_{737}	x_{738}	x_{739}	x_{740}	x_{741}	x_{742}	x_{743}	x_{744}	x_{745}	x_{746}	x_{747}	x_{748}	x_{749}	x_{750}	x_{751}	x_{752}	x_{753}	x_{754}	x_{755}	x_{756}	x_{757}	x_{758}	x_{759}	x_{760}	x_{761}	x_{762}	x_{763}	x_{764}	x_{765}	x_{766}	x_{767}	x_{768}	x_{769}	x_{770}	x_{771}	x_{772}	x_{773}	x_{774}	x_{775}	x_{776}	x_{777}	x_{778}	x_{779}	x_{780}	x_{781}	x_{782}	x_{783}	x_{784}	x_{785}	x_{786}	x_{787}	x_{788}	x_{789}	x_{790}	x_{791}	x_{792}	x_{793}	x_{794}	x_{795}	x_{796}	x_{797}	x_{798}	x_{799}	x_{800}	x_{801}	x_{802}	x_{803}	x_{804}	x_{805}	x_{806}	x_{807}	x_{808}	x_{809}	x_{810}	x_{811}	x_{812}	x_{813}	x_{814}	x_{815}	x_{816}	x_{817}	x_{818}	x_{819}	x_{820}	x_{821}	x_{822}	x_{823}	x_{824}	x_{825}	x_{826}	x_{827}	x_{828}	x_{829}	x_{830}	x_{831}	x_{832}	x_{833}	x_{8

Data points	I			II			III			cluster	new cluster
	x	y	(2, 10)	(6, 6)	(1.5, 3.5)						
A1	2	10	0	5.656	6.519	1	3	2	2	1	1
A2	2	5	5.	4.12	1.58	3	2	2	2	3	3
A3	8	4	8.485	2.82	6.51	2	2	2	2	2	2
B1	5	8	3.605	2.23	5.700	2	2	2	2	2	2
B2	7	5	7.071	1.414	5.70	2	2	2	2	2	2
B3	6	4	7.211	2	4.527	2	2	2	2	2	2
C1	1	2	8.062	6.40	1.58	3	3	3	3	3	3
C2	4	9	2.936	3.60	6.041	2.1	1	1	1	1	1

$$\Rightarrow A_1, C_2 = \frac{2+4}{2} = 3, \frac{10+9}{2} = 9.5 = A_1 = (3, 9.5).$$

$$2 \Rightarrow A_3, B_1, B_2, B_3 = B_1 (6.5, 5.25)$$

$$C_1 = (1.5, 3.5)$$

	I	II	III	new cluster
	(3, 9.5)	(6.5, 5.25)	(1.5, 3.5)	new cluster
A1	1.118	6.543	6.519	1
A2	4.609	4.506	1.58	3
A3	7.433	1.952	6.51	2
B1	2.5	3.132	5.700	1
B2	6.020	0.559	5.70	2
B3	6.264	1.346	4.527	2
C1	7.762	6.388	1.58	3
C2	1.118	4.506	6.041	1

I	II	III	new cluster
(3.6, 9)	(7, 4.33)	(1.5, 3.5)	new cluster
1.88	7.55	6.519	3
4.308	5.044	1.58	3
6.66	1.05	6.51	2
1.72	4.17	5.700	1
5.25	0.67	5.70	2
5.54	1.05	4.527	2
7.46	6.43	1.58	3
0.4	5.56	6.041	1

Termination condition:
There is no change in cluster, so all data points are in appropriate cluster.

• Hierarchical clustering : ALGORITHM.

Input section : set of unlabelled training examples.

Step 1: each example will form one cluster, so for N training examples there are N clusters initially.
Step 2: calculate the distance b/w 2 clusters. The least cluster to cluster distance will be merged.
 \therefore the clusters reduced from N to $N-1$.

Step 3: the second step is repeated till the termination condition (some threshold specified is reached).

on the no. of clusters.

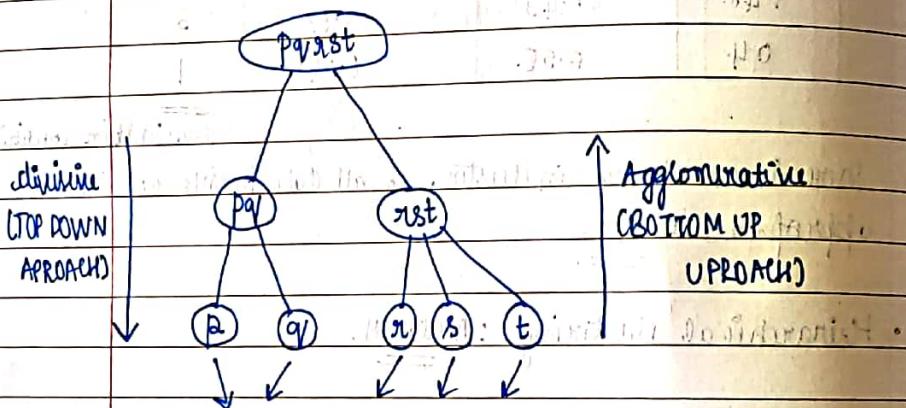
1 April 2024.

Hierarchical clustering :

2 types:

① Agglomeration
(latin word means collecting)

② Division
splitting



• HAC

→ Hierarchical agglomerative clustering

3 types:

① single linkage } minimum distance btw 2 clusters

② complete linkage } (min)

③ Average linkage } more than centroid alt.

① Single linkage

merges the two clusters and finds distance and chooses min

$$(q,s), (t)$$

$$(r,t), (s,t)$$

$$\min_{r,s,t} \dots$$

$$= (a \rightarrow b \rightarrow c) \dots$$

it calculates the minimum distance btwn 2 clusters and min fn gives a single cluster.

2) complete linkage

max fn uses.

usually spherical in shape.

$$(q,s), (t)$$

$$(r,t), (s,t)$$
 max distance btwn 2 options

$$\max(x,y)$$

⇒ convex many eq.

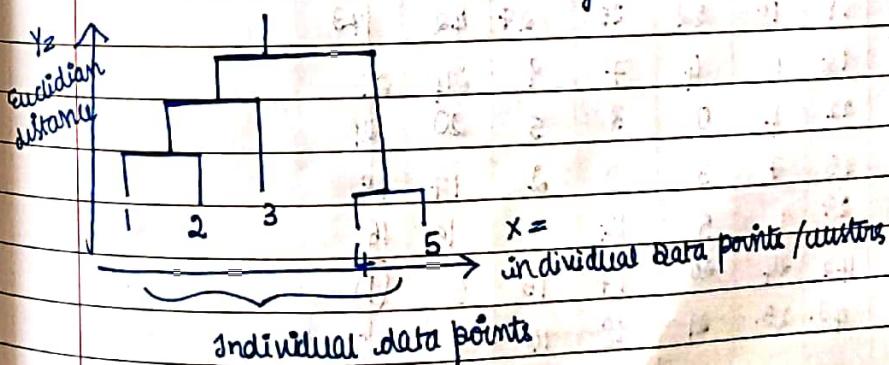
③ Average linkage. (not there).

• Represented by Dendrogram

pictorial representation of hierarchical clustering.

2 greek words: Bundel - Tree

orogram - Drawing



individual data points

→ 1 dimensional points
(x points only given)
(1 coordinate of a point is given) → 2 dimensional points
Both x & y coordinates of point P are given.

• Euclidean = $\sqrt{(x_2 - x_1)^2} = |x_2 - x_1|$
Distance

1. y coordinate is absent

2. square & square root gets cancelled.

→ Applications of HAC:

① fake news detection

② identify the keywords.

③ Appropriate distance will be calculated and will be grouped under one cluster.

i) find the single linkage for the following 10 points.

	18	22	25	27	42	43
18	0	4	7	9	24	25
22	4	0	3	5	20	21
25	7	3	0	2	17	18
27	9	5	2	0	15	16
42	24	20	17	15	0	0
43	25	21	18	16	0	0

① is the minimum zour distance

non-zero
step 2: minimum distance = 1.

42 and 43 are clustered together
and eliminate 43 from row and column.

	18	22	25	27	(42, 43)
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	2	17 ↑
27	9	5	2	0	15
(42, 43)	24	20	17	15	0

step 3:

	18	(22, 25)	(27, 25)	(42, 43)
18	0	4	7	24
(22)	4	0	3	20
(27, 25)	7	3	0	17
(42, 43)	24	20	7	0

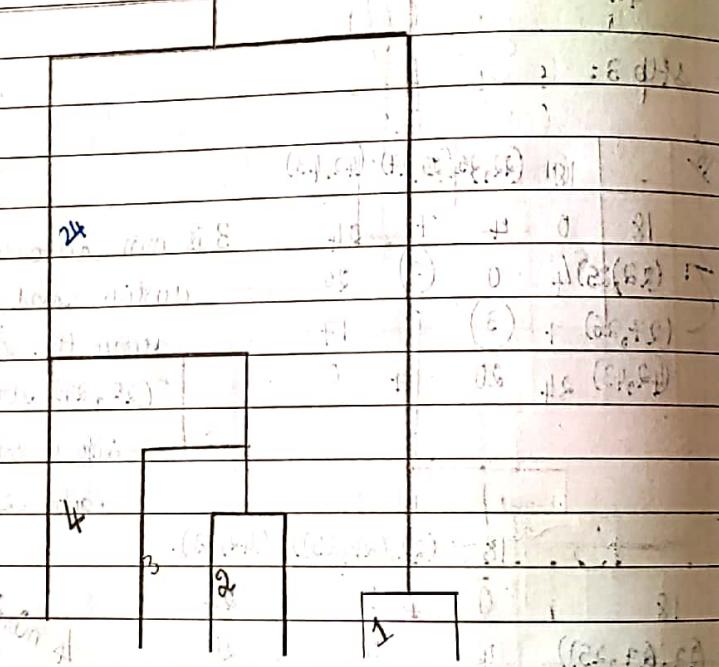
3 is min distance between
clustering and (25, 27)
Hence the cluster (22) and
(25, 27) were merged
into a single cluster
eqt to (22, (25, 27))

	18	(22, (25, 27))	(42, 43)
18	0	4	24
(22, (25, 27))	4	0	20
(42, 43)	24	20	0

4 minimum

IV		(18, (22, 25, 27)))	(42, 43)	
		0	24	0
		(42, 43)	24	0

(18, (22, 25, 27)), (42, 43)



18 22 25 27 42 43 (1, 5)

a) complete linkage

for the following data points : 1, 5, 8, 10, 2.

	1	2	3	4	5	→ column nos
1	1	5	8	10	2	
2	1	0	4	7	9	
3	5	4	0	3	5	3
4	8	7	3	0	2	6
5	10	9	5	2	0	8
6	2	1	3	6	8	0

(1, 5) 2 3 4

(1, 5)	0	4	7	9
2	4	0	3	5
3	7	3	0	2
4	9	5	2	0

(1, 5) = 1.

(1, 5) 2 (3, 4)

(1, 5) 0 4 9

(2) 0 5

(3, 4) 9 5 0

Distance btwn (2)(1, 5) is
one unit

(2, 1), (2, 5)

max(x, y)

max(dist(x), dist(y))

max(4, 3) =

= 4

Distance (3, (1, 5)) ⇒

max((3, 1) (3, 5))

= (7, 6)

= 7

d(1, 5) (3, 4) Distance max of (4, (1, 5))

= max(1, 5) 3, (1, 5), 4) = (4, 8) (4, 5)

max = max(7, 9) = (9, 8)

= 9

= 9

d(2, (3, 4)) = max(d(2, 3), d(2, 4))

= max(3, 5)

d((3, 4), 2) = (d(3, 2), d(8, 2)) = 9

= (5, 3) max = 5

single linkage - neighbour clustering
complete linkage - farthest neighbour clustering

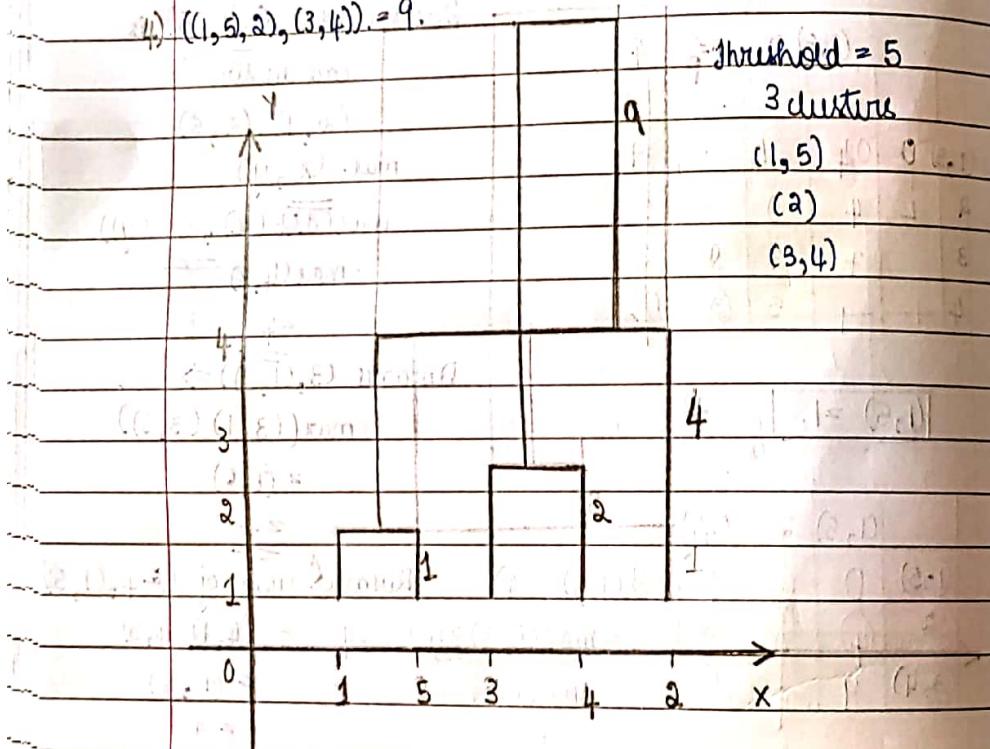
	(1, 5, 2)	(3, 4)	(1, 5), 2, (3, 4).
(1, 5, 2)	0	9	$\max(d((1, 5), 2), 3)$ and $d((1, 5), 2), 4)$
(3, 4)	9	0	= $\max((1, 5) \text{ with } (3, 4) \text{ and } (2) \text{ with } (3, 4))$ = $\max((1, 5)(3, 4) \text{ & } (2)(3, 4))$

$$2) (3, 4) = 2$$

$$= 5, 9 = 9$$

$$3) (1, 5), 2 = 4.$$

$$4) ((1, 5), 2), (3, 4) = 9.$$



→ HAC (Hierarchical Agglomerative clustering)

looks for similarities between 2 points belonging to different clusters.

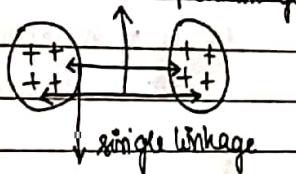
similarly, (HDC)

Divisive hierarchical clustering works for similarities dl
2 points where i pt belongs to cluster A and j pt belongs to cluster B.

→ single linkage / single nearest distance / Neighbour clustering
distance btwn the closest members of 2 clusters A and B.

→ complete linkage / complete farthest distance / farthest neighbour clustering
distance btwn 2 farthest members / pts of 2 diff clusters A & B.

Bendogram:



combining leaf nodes leading to cluster.

→ In HAC, we need to have predefined no of clusters naming k mean clustering.

• Agglomerative consumes more time & space, it ensures nearby points (near) end up in same cluster.

- Identification of criminal activities } Applications
- Document analysis.

8 April

Unit 3

ANN (Artificial neural network).

2 researchers: McCulloch
Pitts

functions of brain: recognizing & discriminating

HUMAN BRAINANN

- ① Neuron switching time approx
 $= 0.001 \mu s$.
- ② Many neural links threshold
 switching units

- ③ No. of neurons approx $= 10^{10}$ or
 10^{11} .
- ④ many

- ⑤ connections per neuron is
 $\approx 10^4$ to 10^5 .
- ⑥ many

- ⑦ sum recognition time approx
 $\approx 10^{-1} \mu s$.
- ⑧ no specific point

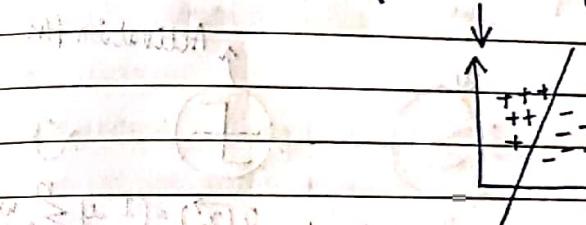
- ⑨ much parallel computation
- ⑩ highly parallel.

- When to use ANN?
- ⑪ Input is high dimensional
- ⑫ noisy data

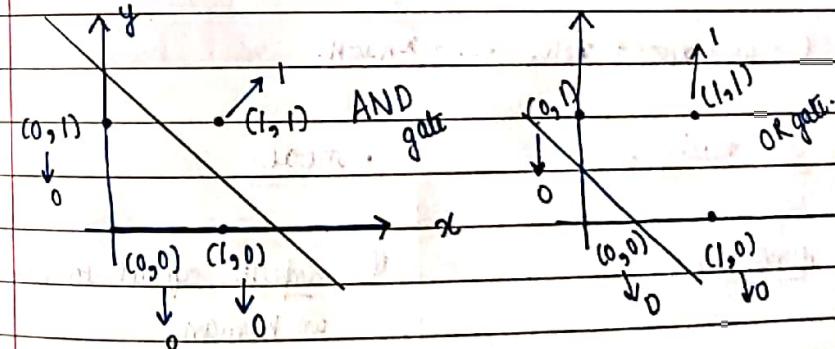
(3) when target fn form is unknown

1st type of ANN is PERCEPTRON (simplest).

- simplest ANN
- proposed by Frank Rosenblatt.
- perceptron - recognition automaton.
- supervised ANN
- used with linear data (eg: linear classification)



Graphical representation for logical gates.



AND and OR are linearly separable but XOR gate is not.

So perceptron can't be used with XOR gate.

Components of perception.① Input layer (X) $\rightarrow x_1, x_2, x_3, \dots, x_n$ ② weight (W) $\rightarrow w_0, w_1, w_2, \dots, w_n$

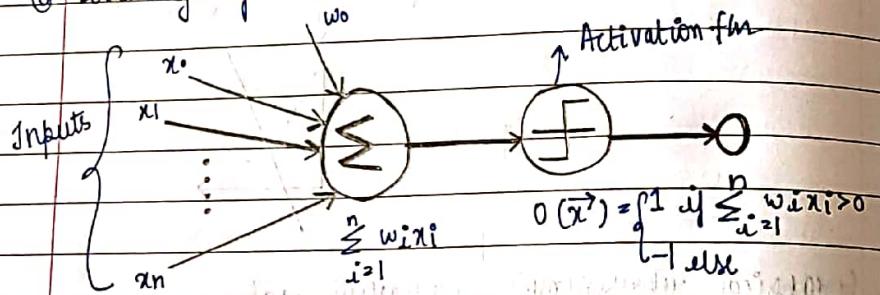
③ Bias.

It adds more flexibility to the network.

④ Activation function

⑤ Output

⑥ Training algorithm.



$$O = w_0 + x_1 w_1 + x_2 w_2 + \dots + x_n w_n$$

• cell nucleus

① dendrite

② synapse.

③ axon

• mode.

① input equal to dendrite in human.

② weights

③ output

• Perception training rule:

$$\Delta w_i \leftarrow w_i + \Delta w_i \quad \text{change in weight}$$

$$\Delta w_i = \eta (t - o) x_i \quad \begin{array}{l} \text{error} \\ \downarrow \\ \text{eta} \end{array} \quad \begin{array}{l} x_i \rightarrow i\text{th training eq. i\th sp.} \\ \text{target value} \end{array}$$

learning rate

lower the better

• Gradient descent.

⇒ Gradient means change, descent means downward.

⇒ it is an optimization algorithm.

(reduce the error).

$$\text{① } E[\vec{w}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

error w.r.t weight

min

$$\text{② } \nabla E[\vec{w}] = \left[\frac{\delta E}{\delta w_0}, \frac{\delta E}{\delta w_n} \right]$$

$$\text{③ } \Delta \vec{w} = -\eta \frac{\delta E}{\delta w_i} \quad (\text{change in weight})$$

$$\frac{\delta}{\delta w_i} \left(\frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \right) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)$$

= $\sum_{d \in D} (t_d - o_d)$ target value is const.

$$= \delta \leq (0 - w_i x_i)$$

$$\frac{\partial \delta}{\partial w_i}$$

$$= -x_i$$

=

- Back-propagation algorithm (9 April 2024).

① multilayer network

② non-linearly separable decisions

INPUTS:

1) training example $\rightarrow \langle \vec{x}, \vec{t} \rangle$

$\vec{x} \rightarrow$ each training example

$\vec{t} \rightarrow$ target value associated with each training example.

2) η :

$\eta \rightarrow$ network in \rightarrow inputs

3) η \rightarrow network
 $\eta_{\text{hidden}} \rightarrow$ units

4) $\eta \rightarrow$ network

$\eta_{\text{out}} \rightarrow$ output units

5) $\eta \rightarrow \text{eta} \rightarrow$ learning rate

All these are inputs of BP algorithm

STEPS:

- I create NN
- II propagate @ forward \rightarrow i/p's
- III update weights
(B) backward - errors

. Back propagation Algorithm of training example.

① Each training example is a pair of the form $\langle \vec{x}, \vec{t} \rangle$ where \vec{x} is a vector of network input values, \vec{t} where is a vector of network output values.

target

steps:

I create a feed forward network with η in inputs, η hidden units, η out output units.

initialise all network weights to some small random value.

(0.5 to 0.5) until the termination condition is met.

② for each training example $\langle \vec{x}, \vec{t} \rangle$, do
propagate the input forward through the network.

II Propagate the input forward through the network.

③ Input the instance \vec{x} to the network and compute the output O_u of every unit u in the network - forward propagation.

④ Propagate the error backward through the network.

for each output network unit k , calculate its error term δ_k .

$$\delta_k \leftarrow \sigma_k(1 - \sigma_k)(t_k - O_k)$$

where $\sigma_k(1 - \sigma_k)$ is the derivative of sigmoid function
 $\downarrow (\sigma(\cdot))$

also called as

sigmoid squashing fn

for each hidden unit h .

$$\delta_h \leftarrow o_h(1-o_h) \sum_{k \in \text{outputs}} w_{hk} \delta_k$$

- derivative of sigmoid activation fn wrt $h \Rightarrow o_h(1-o_h)$
- w_{hk} → weight associated from hidden layer to output layer.

III Update each network weight W_{ji}

$$W_{ji} = W_{ji} + \Delta W_{ji}$$

where

$$\Delta W_{ji} = \eta \delta_j x_{ji}$$

where x_{ji} is input from i th unit to j th unit

and the weight associated with the connection

is W_{ji} .

δ_j → error associated with j th unit

⇒ the stochastic gradient descent version of back propagation algorithm for feed forward networks containing 2 layers of sigmoid units. is as discussed above.

• Stochastic → quick work → guess work

meaning guess or aim

Now, the weights are updated independently of each other and done before the summation (\sum) unlike after \leq in

gradient descent algorithm.

→ as update is individual, performance is better.

sample No.	X	y
P1	0.40	0.53
P2	0.22	0.38
P3	0.35	0.32
P4	0.26	0.19
P5	0.08	0.41
P6	0.45	0.20

①	P1	P2	P3	P4	P5	P6
(0.40, 0.53)	0	0.234	0.215	0.367	0.3417	0.23
(0.22, 0.38)	0.234	0	0.1431	0.1974	0.1431	0.145
(0.35, 0.32)	0.215	0.1431	0	0.1581	0.2846	0.10
(0.26, 0.19)	0.367	0.1974	0.1581	0	0.284	0.29
(0.08, 0.41)	0.3417	0.1431	0.2846	0.284	0	0.386
(0.45, 0.20)	0.235	0.243	0.10	0.211	0.386	0

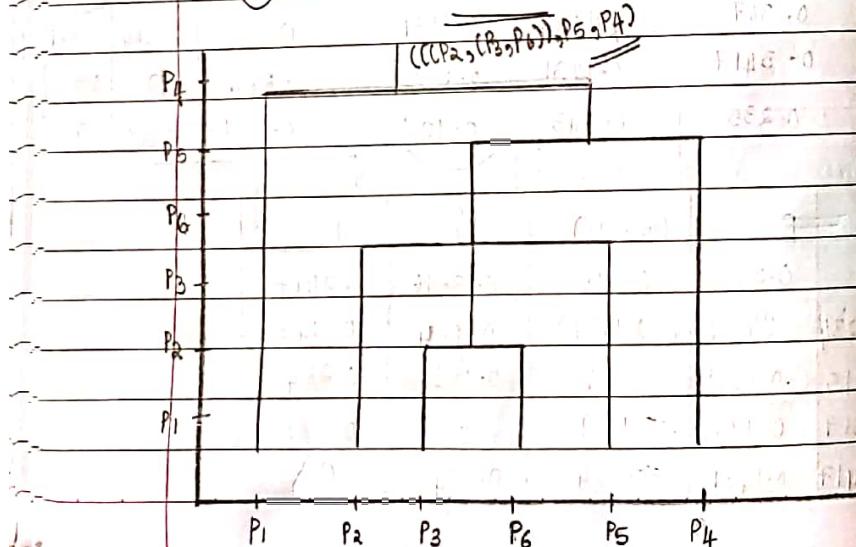
②	P1	P2	(P3, P6)	P4	P5
P1	0	0.23	0.215	0.3676	0.3417
P2	0.234	0	0.1431	0.1974	0.1431
(P3, P6)	0.215	0.1431	0	0.1581	0.284
P4	0.367	0.1974	0.1581	0	0.284
P5	0.3417	0.1431	0.284	0.284	0

(3) P ₁	$(P_2, (P_3, P_6))$	P ₄	P ₅
P ₁	0	0.234	0.3417
$(P_2, (P_3, P_6))$	0.234	0	0.1431
P ₄	0.367	0.1974	0
P ₅	0.3417	0.1431	0

(4) P ₁	$(P_2, (P_3, P_6), P_5)$	P ₄	
P ₁	0	0.234	0.3676
$(P_2, (P_3, P_6), P_5)$	0.234	0	0.1974
P ₄	0.3676	0.1974	0

(5) P ₁	$((P_2, (P_3, P_6)), P_5, P_4)$		
P ₁	0	0.234	
$((P_2, (P_3, P_6)), P_5, P_4)$	0.234	0	

(P₁) and $((P_2, (P_3, P_6)), P_5, P_4)$ are the 2 clusters



Naïve Bayes Derivation

→ Bayes' theorem

→ find

→ $P(A|B)$: known/assumed

→ more informed decision

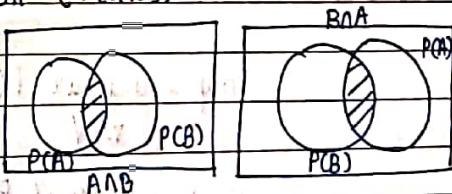
$P(A) \rightarrow$ general

conditional probability

→ Based on Bayes' theorem.

→ conditional probability: w. the help of which we can make more informed decision ($P(A|B)$)

→ general probability:



$$\rightarrow P(h|D) = P(D|h) P(h) \quad \text{---(1)}$$

$$\rightarrow P(D) \quad \text{---(2)}$$

$$\text{hmap} = \underset{h \in H}{\operatorname{argmax}} P(h|D) \quad \text{---(3)}$$

Put (1) in (3)

$$= \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h) \cdot P(h)}{P(D)} \quad \text{---(3)}$$

in ④ $P(D)$ is constant

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(D|h) P(h) \rightarrow ④$$

where $P(D)$ is constant

But the equation ④ cannot be used for real world examples

\Rightarrow ∴ modified equation is:

v_{MAP} = most probable target value.

Given the features $(a_1, a_2 \dots a_n)$

which describe the training example

$$v_{MAP} = \underset{v \in V}{\operatorname{argmax}} \underset{h \in D}{P(v_j | a_1, a_2 \dots a_n)} \quad ③$$

$$v_{MAP} = \underset{v \in V}{\operatorname{argmax}} \frac{P(a_1, a_2 \dots a_n | v_j) \cdot P(v_j)}{P(a_1, a_2 \dots a_n)} \quad - ③$$

where

$P(a_1, a_2 \dots a_n)$ is constant

$$v_{MAP} = \underset{v \in V}{\operatorname{argmax}} P(a_1, a_2 \dots a_n | v_j) \cdot P(v_j) \quad - ④$$

$$v_{MAP} = \underset{v \in V}{\operatorname{argmax}} P(v_j) \prod_i P(a_i | v_j)$$

Evaluation metrics:

confusion matrix

		Predicted	
		No	Yes
Actual	No	TN = 45	FP = 5
	Yes	FN = 5	TP = 95

① Accuracy

Total Actual Yes = 100

Total Actual No = 50

$$\frac{\text{No. of correctly classified instances}}{\text{No. of All instances}} = \text{Accuracy}$$

$$\frac{TN + TP}{all} = \frac{140}{150} = 0.93 \times 100 = 93\%$$

② Precision

When the machine predicts yes, how often it is correct

$$\text{Precision} = \frac{\text{No. of correctly labelled +ve instances (TP)}}{\text{All +ve labelled instances}}$$

$$= \frac{TP}{TP + FP} = \frac{95}{100} = 95\%$$

3) Recall (also called as sensitivity or true positive rate)

Recall = % of correctly labelled +ve instances.
All +ve instances (or actual Yes).

$$= \frac{TP}{TP+FPN} = \frac{95}{100} = 95\%$$

4) $\frac{2 \times Precision \times Recall}{Precision + Recall}$ = F-Score.

$$= \frac{2 \times 95\% \times 95\%}{95\% + 95\%} = \frac{2 \times 95 \times 95}{100 \times 100} = \frac{190}{10000} = 0.95 = 95\%$$

1) find if the patient is having cough or not from the following :

$$P(\text{cough}) = 0.008$$

Report $\rightarrow +$ -ve

$$P(+|\text{cough}) = 0.98$$

$$P(-|\text{cough}) = \underline{\underline{0.02}}$$

$$P(+|\text{noncough}) = 0.03$$

$$P(-|\text{noncough}) = 0.97$$

② Verify if the patient is having cough or not given report is +ve.

$$P(\text{cough} | +ve) = \frac{P(+ve | \text{cough}) \times P(\text{cough})}{P(+ve)}$$

$$= 0.0078$$

$$P(\text{noncough} | +ve) = \frac{P(+ve | \text{noncough}) \times P(\text{noncough})}{P(+ve)}$$

$$= 0.03 \times 0.992$$

$$= 0.02976$$

Vmap = No cough.

$$\text{Jacc: } P(\text{cough} | -ve) = \frac{P(-ve | \text{cough}) \times P(\text{cough})}{P(-ve)}$$

$$= 0.02 \times 0.008$$

$$= 0.00016$$

$$P(\text{cough} | -ve) = \frac{P(-ve | \text{noncough}) \times P(\text{noncough})}{P(-ve)}$$

$$= 0.992 \times 0.97$$

$$= 0.96224$$

Vmap = No cough (highest value).

PCA : principal component analysis

It is a dimensionality reduction technique.

Algorithm has 5 steps:

- compute mean of all features from the dataset having

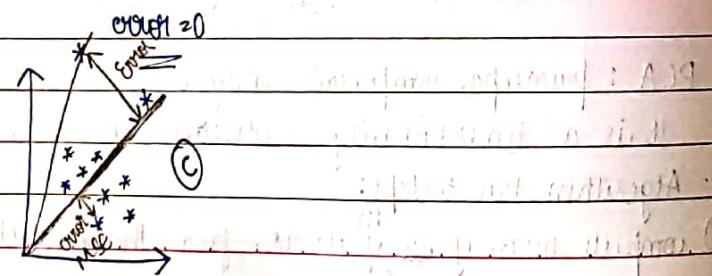
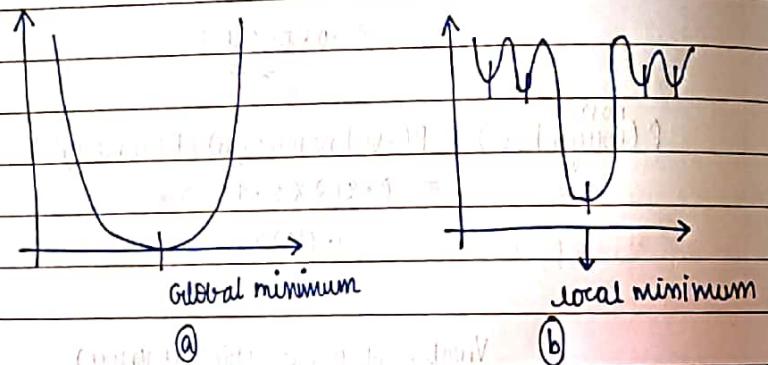
n features and N examples.

- ② calculate the covariance matrix
- ③ calculate eigen values & eigen vectors of covariance matrix
- ④ Normalize the eigen vectors.
- ⑤ Create new dataset by arranging eigen values in descending order. unit eigen vector corresponding to the largest eigen value = PCA

18 April • Mean squared error. $(y_i - \hat{y}_i) = \text{error}$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

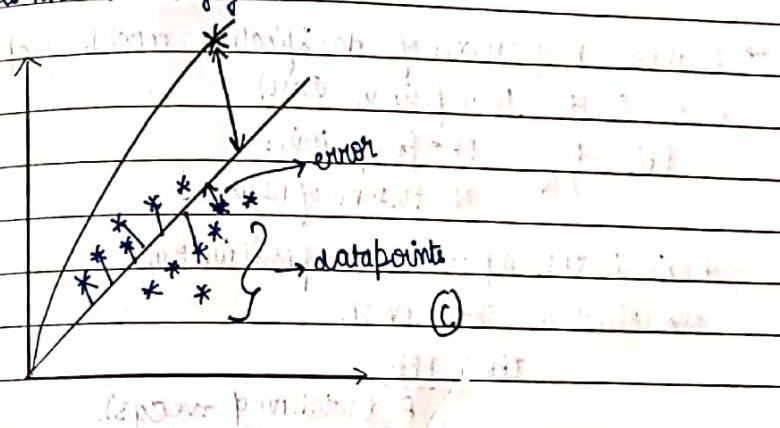
→ RMSE → square root of MSE



mean squared error:

fig @ is differentiable at each point and hence slope can be calculated at any point. It has only one global minimum, ∴ converges/convergence is faster.

disadvantage: it can / is not robust in outliers as shown in fig (c).



→ RMSE : standard to way to measure error.

$$\text{RMSE} = \sqrt{\text{MSE}}$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

only used for regression,

also used for hyper parameter tuning

- Hyper parameter : values which are initialized before the process. (High level parameter). Eg: weights.
- parameter : values which are initialized during the process

lowest RMSE = 0.

⇒ Error rate = 1 - Accuracy.

↓
how many miss classification done.

- ROC (Receiving operating characteristic curve).

→ It shows performance of classification model graphically.

x axis is FPR (False positive rate)

$$FPR = \frac{FP}{N} \quad FP = \text{false positives}$$

$N = \text{total no of negatives}$

y axis is TPR is plotted (true positive rate)
also called as SENSITIVITY.

$$TPR = \frac{TP}{P} \quad P = \text{(total no of +ve obs.)}$$

→ Area under curve (AUC).

Higher the value of AUC, good is the performance of model.

AUC = 0 → Bad

= 1 → Good.

→ It is a global measure of ability of a test to determine if a specific condition is present or not.