

## Unit = 3

### Cluster Analysis

#### BASIC CONCEPT of clustering

- Cluster → nothing ↪ just grouping similar data

Clustering → It refer to group the entire data into similar form & other group data will always diff

#### Example

→ E-commerce site decide it entire customer into 3 cluster

1 cluster : "Discount buyers"

2 cluster : "Premium Buyer"

3 cluster : "One-time buyer"

Here site cluster into 3 customer

#### "Measures of similarity"

##### • Euclidean Distance

$$d(P, Q) = \sqrt{\sum (P_i - Q_i)^2}$$

##### • Manhattan Distance

$$d(P, Q) = \sum |P_i - Q_i|$$

##### • Cosine Similarity

$$\text{Similarity}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

##### • Jaccard Similarity

$$\text{Jaccard}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Page No.	
Date:	

## Euclidean distance

Let understand by example

Let a food app;

Customer address point  $(2, 3)$

Restaurant " point  $(7, 8)$

8

6

4

2

$(2, 3)$

$(7, 8)$

2 4 6 8

So by Euclidean distance  
we find ~~for~~ shortest

So the shortest distance,  $d = \sqrt{\sum (x_i - y_i)^2}$   
 $d = \sqrt{(7-2)^2 + (8-3)^2} = \sqrt{50} = 7.07$

Other type ques on Euclidean distance

If we have four diff clusters so one new person come with  $(x, y)$  so which cluster is suitable for new one

Real Example → A Mall divides his customer on the Basis of Age( $x_1$ ) and Income( $y$ ) so new customers come with age 28 & income so tell which cluster is good

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$
 age Inc

clusters

$$\Rightarrow C_{new} = (28, 42)$$

⇒ distance from  $C_1$

$$d = \sqrt{(28-25)^2 + (42-40)^2} = 3.6$$

distance from  $C_2$

$$d = \sqrt{(28-30)^2 + (42-45)^2} = 3.6$$

distance from  $C_3$

$$d = \sqrt{(28-45)^2 + (42-80)^2} = 11.6$$

distance from  $C_4$

$$d = 48.3$$

$C_1$	25	40
$C_2$	30	45
$C_3$	45	80
$C_4$	50	85

so  $C_1$  &  $C_2$

cluster

$C_1$  is

good

for new

Reg No. \_\_\_\_\_  
Date \_\_\_\_\_

Manhattan Distance (यह euclidean distance का straight line distance की जाती है। यह block moves के द्वारा दर्शाया जाता है।)

Let understand by example  
Let suppose you live in grid-type city  
and your location is (2,3) & your friend (7,8)  
So shortest straight distance?

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

$$= |7-2| + (8-3) = 5+5 = 10 \text{ blocks}$$

Actual use :-

Car service actual road distance in Real life  
Kaise use kerti ha C jaha city grid system ho waha

Real life example

Let suppose you are in restaurant origin point  
and you need to go diff. location  
and the city is grid so find all  
distance

Restaurant location R(2,3)

C1(5,7)

C2(8,8)

C3(8,9)

distance b/w Res to C1

$$|5-2| + |7-3| = 7$$

distance b/w Res to C2

$$|8-2| + |8-3| = 11$$

distance b/w Res to C3

$$|8-2| + |9-3| = 12$$

Ques from PyQ :-

$\Rightarrow$  Given two object  $x = (22, 14, 2, 10)$  and  $y = (20, 0, 16, 8)$

(1) Euclidean Distance

$$d = \sqrt{(20-22)^2 + (0-14)^2 + (16-2)^2 + (8-10)^2}$$

$$d = \sqrt{(2)^2 + (14)^2 + (16)^2 + (2)^2} = \sqrt{4+1+36+4} = \sqrt{45} \text{ m}$$

(2) Manhattan Distance

$$d = |20-22| + |0-14| + |16-2| + |8-10|$$

$$= 2 + 1 + 6 + 2 = 11$$

## Type of clusters

### ① Partitioning method

These methods divide the dataset into  $K$  partitions (clusters) where each partition represent a cluster.  
Real life eg:- customers divide into diff group (clusters)

cluster

K-Means:- Assign each point to the nearest cluster centroid and iteratively update centers.

K-Medoids (PAM - Partitioning Around Medoids):

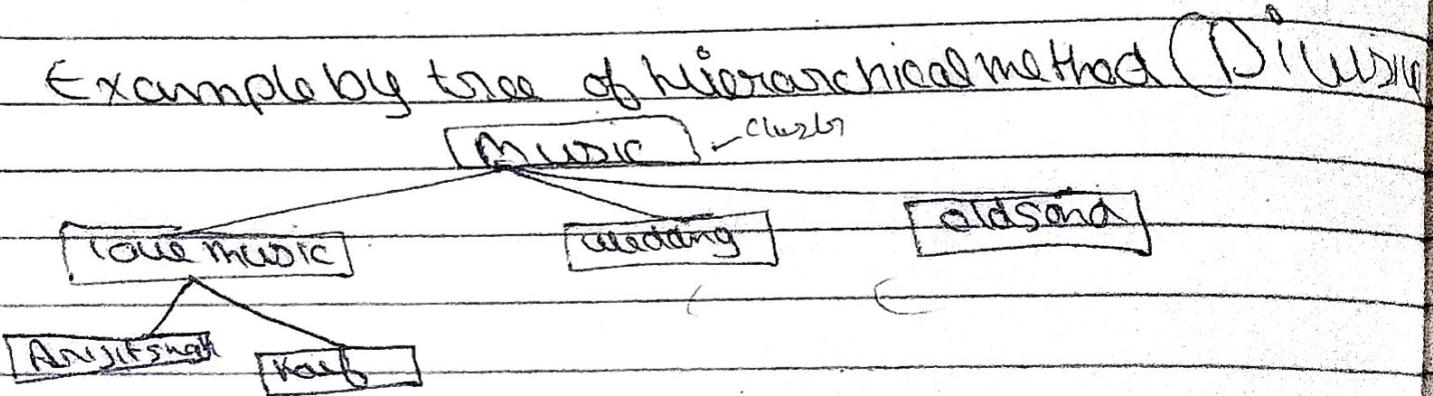
Similar to K-Means but choose actual data point as centroids, making it robust to outliers.

### ② Hierarchical Methods

Hierarchical clustering (unsupervised technique) where entire data are arranged in hierarchical tree structure (dendrogram).

Each node represent a cluster and leaf of node is data.

Example by tree of hierarchical method (Divisive)

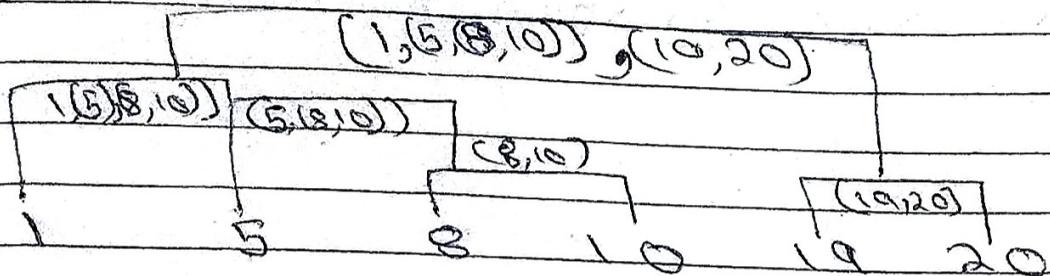


Agglo	
Clust	

## Agglomerative Algorithm

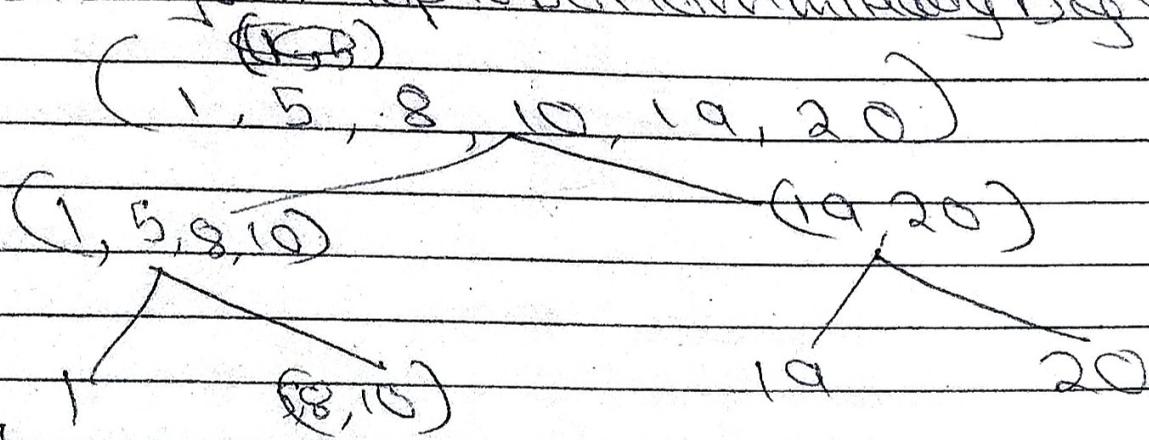
Let understand by example,  
Spotify user create playlist of his song  
so user 1 & Spotify suggest his song into cluster

It start from Bottom to top



divisive

It start from top to bottom initially Big clust



### Similarity Measure

- H.C me<sup>metric loss</sup>

half point decide  
Karna hole

a) single linkage:  
all clusters have  
member k ad w

b) complete linkage



3) Density based method

= Cluster are formed based on the density of data Point

Key algorithm

DBSCAN → Density Based spatial clustering  
of Application with Noise : Finds Cluster  
of arbitrary shapes & identifies noise

Page No.	
Date	

### Ques on K-mean

customers

✓ C1

✓ C2

✓ C3

✓ C4

✓ C5

	Income(X)	Spending Score
C1	2	16
C2	4	12
C3	5	25
C4	10	30
C5	11	35

cluster 1

C1, C2, C3 = C5

cluster 2

C2, C3, C4

let

C1

C1

from C1

$$C_3^1 = \sqrt{(5-2)^2 + (25-16)^2} \\ = \sqrt{(3)^2 + (15)^2} = \sqrt{234} = 15.29 \quad C_3 \text{ from } C_2 = \sqrt{(5-4)^2 + (25-10)^2} \\ = \sqrt{1^2 + (15)^2} = \sqrt{226} = 15.02$$

so C3 from C2 is short so C3 go in C2

now centroid point of cluster 2 =  $\frac{4+5}{2}, \frac{25+10}{2} = 4.5, 18.5$

$$C_4 \text{ from } C_1 = \sqrt{(10-2)^2 + (30-16)^2} ; C_4 \text{ from } C_2 = \sqrt{(10-4.5)^2 + (30-18.5)^2} \\ \Rightarrow \sqrt{(8)^2 + (20)^2} = \sqrt{464} = 21.54 \quad \sqrt{(5.5)^2 + (15)^2} = \sqrt{252.5} = 12.74$$

so C4 from C1 is short so C4 go in C2

now centroid point of cluster 2 =  $\frac{4.5+10}{2}, \frac{18.5+30}{2}$

$$= [7.25, 24.25]$$

C5 from C1

$$(11-2)^2 + (35-16)^2$$

$$= \sqrt{(9)^2 + (19)^2}$$

$$= \sqrt{81 + 361}$$

$$= \sqrt{442} = 26.510 = 29.625$$

C5 from C2

$$\sqrt{(11-7.25)^2 + (35-24.25)^2}$$

$$= \sqrt{(3.75)^2 + (10.75)^2}$$

$$= \sqrt{14.0625} + (115.5625)$$

so C5 from C1 is short so C5 go to C1

now cluster point of C1 =  $\frac{11+2}{2}, \frac{35+10}{2} = 6.5, 22.5$

measure for cluster validation

↳ It means that cluster (group) we make are They really meaningful or not

## Type of Cluster Validation

### (1) Internal Validation

= Cluster ka evaluate karte ha dataset ke ander ke hi info se

#### • Internal Validation check

↳ How much cluster point closely & how much cluster different

↳ → Let suppose we have 100 people dataset where height weight And you make 3-cluster and run K-mean

Now internal validation say

cluster 1 people height weight variation (good)  
if cluster 1 & cluster 2 belong all diverse  
alag hoga to (good & separate)

#### • External Validation

→ If you have student dataset of percentage you make 3 cluster (Mark, Avg, Stng) but you have actual label that is achieve you now compare the actual label with cluster label

Rand index → tell how much cluster label is close the actual label

#### Relative Validation

→ do make  $K=3$  and  $K=4$  and now check which is more useful

# DATA MINING

[Unit  $\Rightarrow$  4]

Association Rule mining: Transaction data-set, frequent itemset, support measure, rule generation, confidence of association rule, Apriori algorithm, Apriori principle

- Transaction data-set  $\Rightarrow$  A Transaction data-set is a collection of transaction where each transaction represent a set of item bought by a customer at a specific time.

How transaction Data-set look like

Transaction ID	Item Bought	Each row $\Rightarrow$ one customer shopping basket
T <sub>1</sub>	Bread, Milk	
T <sub>2</sub>	Bread, Diaper, Beer, Egg	Each "Item Bought" =
T <sub>3</sub>	Milk, Diaper, Beer, Cola	The set of products purchased together
T <sub>4</sub>	Bread, Milk, Diaper, Beer	
T <sub>5</sub>	Bread, Milk, Diaper, Cola	

This entire table  $\Rightarrow$  Transaction Data-set  
Real life example/use

- 2) This transaction Data Set have hidden feature which are useful for mall company find best use full pattern from transaction data set

Frequent itemset  $\rightarrow$  is a group of items that appear together frequently in a transaction dataset.

In short items are bought together many times, they form a frequent itemset.

Support  $\rightarrow$  means how frequently an item or items appear in the entire dataset  
 $\rightarrow$  out of all transactions, in how many transactions does this item appear.

Support =  $\frac{\text{no. of transaction containing X}}{\text{no. of total transaction}}$

From our previous dataset we find support

Item	Appearing	Support (Count/5)	Support %
Bread	4	4/5	80%
Milk	4	4/5	80%
Diaper	4	4/5	80%
Bear	3	3/5	60%
Cola	2	2/5	40%
Egg	1	1/5	20%

What is Support Threshold (minimum support)?

When we analyzed big data, we can't consider every small item we only want frequent items that appear many times so we decide a minimum support threshold [if support  $> 60\%$  we call it frequent]. So after applying minimum support frequent the frequent is Bread, Milk, Diaper, Bear,  $> 60\%$ .

Page No.	
Date	

Support for itemset (Two item Together)

Itemset	Apriori	Support	Support %	
{Bread, Milk}	T1, T4, T5	$\frac{3}{5} = 0.6$	60%	Very we not
{Milk, Diaper}	T3, T4, T5	$\frac{3}{5} = 0.6$	60%	make Cola, Egg
{Diaper, Bread}	T2, T3, T4	$\frac{3}{5} = 0.6$	60%	like comb. coz they are already more than 60%.
{Bread, Diaper}	T2, T4, T5	$\frac{3}{5} = 0.6$	60%	
{Bread, Beary}	T2, T4	$\frac{2}{5} = 0.4$	40%	

Interpretation  $\Rightarrow$  all Two item set appear  
more together 60% of all Transaction except  
one combination that {Bread, Beary}

### Rule Generation

$\Rightarrow$  Rule Generation mean creating if-then rule  
from the frequent itemsets we discovered  
to show: when one item is related to another  
item other frequent itemsets:

• {Bread, Milk}

• {Milk, Diaper}

• {Diaper, Beary}

• {Bread, Diaper}

for itemset {Bread, Milk} we can generate

the following rule

(1) Bread  $\rightarrow$  Milk

If a person buys Bread; they are  
likely to buy Milk.

(2) Milk  $\rightarrow$  Bread

If a person buys Milk  
they are likely to buy bread

Kuchnaia Rule Generation  
we bus home combination  
so if there no  
wana has

confidence measure / confidence of association

confidence tell us how often the rule has been found true

e.g. → when a customer buys A good so how much time customer buy B good

$$\text{Confidence } (A \rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}$$

Let take a example

$$\text{e.g. } \{ \text{Bread} \rightarrow \{ \text{Milk} \}$$

meaning "if a customer buy bread, they are also likely to buy milk"

$$\begin{array}{l} A = \text{Bread} \\ B = \text{Milk} \end{array}$$

$$\text{Confidence } (A \rightarrow B) = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}$$

Calculate support value

$\text{Support}(\text{Bread}) \rightarrow$  How many transaction have bread

$$T_1, T_2, T_3, T_4, T_5 \rightarrow 4 \text{ transactions}$$

$$= \frac{4}{5} = 0.8$$

$\text{Support}(\text{Bread} \cup \text{Milk}) \rightarrow$  How many transaction

$$T_1, T_2, T_4, T_5 = 4$$

$$= \frac{4}{8} = 0.5$$

$$\text{Confidence} = (\text{Bread} \rightarrow \text{Milk}) = \frac{0.8}{0.5} = 1.6$$

$$\text{confidence } (\text{Milk} \rightarrow \text{Dialy}) = \frac{\text{Support}(\text{Milk} \cup \text{Dialy})}{\text{Support}(\text{Milk})}$$

$$= \frac{0.30}{0.40} = 0.75$$

110+31 score

more to no abn)

Anten

$$P_{ab} = 3^d - 2^{d+1} + 1$$

Design No.	
Date	

## A Priori Algorithm

→ A Priori Algorithm is used to find frequent itemsets (group of item that often appear together in transaction)

A shopkeeper have 5 transaction record

$T_1 = \text{Tea, Sugar, Milk}$ ;  $T_2 = \text{Tea, Sugar}$ ;  $T_3 = \text{Tea, Milk}$   
 $T_4 = \text{Sugar, Milk}$ ;  $T_5 = \text{Tea, Sugar, Milk}$

Goal - Find out which item are frequently bought together and make strong associations like if a customer buy tea they are likely to buy sugar

ROUGH  
DRAFT

## Step by step of A Priori Algorithm

Let take minimum support = 60%.

That mean

Minimum count  $\Rightarrow 60\% \cdot \frac{1}{5} = 3$  transaction

So, any item or combination that appear in at least 3 transaction is consider frequent

### Step-2 find frequent itemset (L1)

Item	Count	Sup. per	Frequent	
Tea	4	80%	✓	{ all above the minimum }
sugar	4	80%	✓	
milk	4	80%	✓	

$$L_1 = \{\text{Tea}, \text{sugar}, \text{milk}\}$$

### Step - 3 Generate candidate 2-item

Candidate	count	support	Frequent
{Tea, sugar}	3	60%	—
{Tea, milky}	3	60%	—
{Sugar, milky}	3	6%	—

All 3 pair are frequent since each appears in 7/ 3 transaction

$$L_2 = \{ \text{Tea, sugar} \cup \text{Tea, milky} \cup \text{Sugar, milky} \}$$

### Step 4 Generate candidate 3-item

Candidate	count	support	Frequent
{Tea, sugar, milky}	2	40%	X

### Step 5 Final frequent itemset

$$L_1 \Rightarrow \{ \text{Tea} \}; \{ \text{Sugar} \}; \{ \text{Milk} \}$$

$$L_2 \Rightarrow \{ \text{Tea, sugar} \}; \{ \text{Tea, milky} \}; \{ \text{Sugar, milky} \}$$

### Step 6 Generate association rule

- 1)  $\{\text{Tea}\} \rightarrow \{\text{Sugar}\}$ ;  $\{\text{Sugar}\} \rightarrow \{\text{Tea}\}$
- (2)  $\{\text{Tea}\} \rightarrow \{\text{milk}\}$  (4)  $\{\text{milk}\} \rightarrow \{\text{Tea}\}$
- (5)  $\{\text{Sugar}\} \rightarrow \{\text{milk}\}$  (6)  $\{\text{milk}\} \rightarrow \{\text{Sugar}\}$

### Step 7 Calculate Confidence

$$\{\text{Tea}\} \rightarrow \{\text{Sugar}\}$$

$$\text{confidence}(\text{Tea} \rightarrow \text{Sugar}) = \frac{\text{SCT}(\text{US})}{\text{SCT}(\text{S})} = 0.6 / 0.75$$

$$\text{Confidence} = 0.75 / 0.8$$

mean  $\rightarrow$  whenever customer buy Tea there's a 75% chance they also buy Sugar

Project	
Date	
Page No.	
Subject	
Topic	

## Step 8: Real-life Business Use (Optional)

Shopkeepers learn from this:

- Customers who buy Tea also buy sugar (5%)
- Customers who buy Milk also buy sugar (6%)

So he can:

- keep Tea, Sugar, Milk together on the shelf
- offer combo discounts:

"Buy Tea & get 25% off on sugar"

This increases sale and convenience

## Apriori Property

All subsets of frequent itemset must also be frequent

Ex → If {Tea, Sugar, Milk} was frequent  
then {Tea, Sugar}; {Sugar, Milk} and {Tea, Milk}  
must also be frequent

# DATA-MINING

Unit  $\Rightarrow$  5

## Classification

Classification  $\Rightarrow$  Naive Bayes classifier,  
Nearest Neighbour classifier, decision tree,  
overfitting, confusion matrix, evaluation  
metrics and model evaluation

$\Rightarrow$  Classification  $\Rightarrow$  divide the entire data  
into different categories, based on some  
feature using trained model

### ONE Real life example

trained model make a model for spam  
where two cat. spam and not spam  
so we give some keyword like if email  
contain win a phone now / congrats you won 10000  
this type message & so it will go into spam  
else non-spam

### In short

$\hookrightarrow$  Classification  $\Rightarrow$  predicting a  
category label (notano) using data  
pattern

Regression no numeric output at all  
eg (house price)

Classification no categorical output  
at all eg predicting house type

## Naive Bayes classifier

→ It's based on Bayes' Theorem

$$P(\text{Class} / \text{Data}) = \frac{P(\text{Data}) \times P(\text{Class})}{P(\text{DATA})}$$

If  $P(\text{Spam Email}) > P(\text{Not Spam Email})$ , so email is spam

What "Naive" mean

→ completely independent, which mean  
two data / email didn't similar

Real-life Example

Problem → Predict which email is spam  
which is not spam  
on the basis of two word → ["offer" & "win"]

Training Data

Email	Offer	Win	Class [Spam   nonspam]
E1	Yes -	Yes	Spam
E2	Yes -	No	Spam
E3	No	Yes	Spam
E4	Yes -	No	Not Spam
E5	No	No	Not Spam

$$P(\text{Spam}) = \frac{3}{5} \rightarrow P(\text{not Spam}) = \frac{2}{5}$$

Step 2

## Conditional Probabilities

From Data:

Condition

$$P(\text{Offer} = \text{Yes} \mid \text{Spam}) = \frac{2}{3}$$

$$P(\text{Win} = \text{Yes} \mid \text{Spam}) = \frac{2}{3}$$

$$P(\text{Offer} = \text{Yes} \mid \text{NotSpam}) = \frac{1}{2}$$

$$P(\text{Win} = \text{Yes} \mid \text{NotSpam}) = \frac{0}{2} = 0$$

Step-3 new email com your offer and click that won't stop killing both so

for spam

$$P(\text{Spam} \mid \text{Offer, Win}) = P(\text{Offer})_{\text{spam}} \times P(\text{Win})_{\text{spam}} \times P(\text{Spam})$$

$$\Rightarrow \frac{2}{3} \times \frac{2}{3} \times \frac{3}{5} = \frac{4}{15} = 0.266$$

for not spam

$$P(\text{notSpam} \mid \text{Offer, Win}) = P(\text{Offer})_{\text{notspam}} \times P(\text{Win})_{\text{notspam}} \times P(\text{notSpam})$$

$$\Rightarrow \frac{1}{2} \times 0 \times \frac{2}{5} = 0$$

$$P(\text{Spam}) > P(\text{notSpam})$$

so mail is spam

another new email we offer you a car

for spam =

$$P(\text{Spam} \mid \text{Offer}) = P(\text{Offer})_{\text{spam}} \times P(\text{Spam}) = \frac{2}{3} \times \frac{3}{5} = \frac{6}{15}$$

$$P(\text{notSpam} \mid \text{Offer}) = P(\text{Offer})_{\text{notspam}} \times P(\text{NotSpam})$$

$$\Rightarrow \frac{1}{2} \times \frac{2}{3} = 20\% \text{ or } 40\%$$

so mail is not spam

## Nearest Neighbour classifier (K-NN)

Take the decision by watch neighbour data

e.g.: If 10 people and your closure people take a cricket ban (3-4 people) so by (K-NN) you will also say i am cricket fan.

Real life example

Customer like Shoes like loans Class

	like Shoes	like loans	
A	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	Fashion lover
B	<input checked="" type="checkbox"/>	<input type="checkbox"/>	Fashion lover
C	<input type="checkbox"/>	<input checked="" type="checkbox"/>	Fashion lover
D	<input type="checkbox"/>	<input type="checkbox"/>	Non Fashion lover

New customer come E which are like Shoes & not like loans so E is nearest to B so E will blur into Class Fashion lover where B

In numeric data some used Euclidean

Customer	Income	Credit Score	Loan Applied	
A	25	300	No	new customer
B	50	600	Yes	E
C	45	580	Yes	Income = 45
D	20	200	No	Score credit = 590

Instance calculate for E (40, 590)

$$\text{Distance}(E, B) = 14.14$$

$$\text{Distance}(E, C) = 11.18 - \text{most nearest so E come in C group}$$

$$\text{Distance}(E, A) = 2.90$$

$$\text{Distance}(E, D) = 3.91$$

## decision tree over fitting

Decision Tree classifier

= DecisionTree is a supervised learning algorithm which do both classification as well as prediction

How Tree structure look?

- Root node → Started point
- Branches → Conditions
- Leaf Nodes → Final decision/outlet

Let make a loan approval system for bank

we have Income (High/low), Credit score (Good/bad), Age (old/young)

Credit Score

Good

bad

/ Income

Reject

High

Low

Appeal

Reject

## Confusion Matrices

- Confusion Matrix is a table where we see how much prediction are correct and how much is wrong
- They provide deep insight. They tell us which type of mistake model done

### Real life example

We make a machine learning model who predict which customer is good for loan approval

- So, possible result
  - Yes
  - No

Model prediction & actual result both have 4 possible combination

### Confusion matrix Structure

Predicted Yes      Predicted No

Actual Yes	True Positive	False Negative
Actual No	False Positive	True Negative

### Real life meaning of TP, FP, FN, TN

meaning

loan Example

True Positive	Model no correctly "Yes" bala	Actual loan - worthy person Ko model no approve Karo
False Negative	Model no correctly "no" bala	Actual risky person Ko model no select Karo
False Positive	Model no galat "Yes" bala	Risky person Ko model no galat approve Karo
False negative	Model no galat "no" bala	loan worthy person Ko model no select Karo

Actual 100

45 (TP: Post) 15 (False Neg)  
5 (False Pos) 85 (TN: -ve)

60 (Yes)  
40 (No)

### Example

We have 100 people data where 60 are actually loan-worthy & 40 are non-worth

Model me Predict Kya:

- 50 Kya approve Kya BOKO & geet

	Predict Yes	Predict No
Actual Yes (60)	45 (TP)	15 (FP)
Actual No (40)	10 (FN)	30 (TN)

Now calculate Metric

Metric	Formula	Meaning
Accuracy	$(TP+TN)/\text{Total}$	Model overall

$$\frac{75}{100} = 75\% \quad \text{Soh kisi b}$$

-ve brediti

Precision

$$\frac{TP}{TP+FP} = \frac{45}{45+10} = 81.8\%$$

Jiske Ko app

Dale unme

se kitne ko

sachme

loan worthy

Jiske sachme

loan worthy

th unme se

kitne ko

model ne ko

approve ki

F1 Score

$$\frac{2 \times \text{Precision} \times \text{Recall}}{2 \times \text{Precision} + \text{Recall}}$$

Balanced

b/w

$$= 2 \times (81.8\% \times 75\%)$$

$$= 81.8 \times 75\%$$

Precision  
+ Recall

phi

angry

$$\Rightarrow 78\%$$

Consider the class we have

200 [80 screen] are up for no. 6c

00 are not in

no. more

Example you build a machine model to  
predict whether person has diabetes or not  
yes  $\rightarrow$  diabetes Person  
no  $\rightarrow$  non-diabetes  
100 patient

80 application Actual 50  $\rightarrow$  loan-worthy

30  $\rightarrow$  non loan-worthy

ML = 55 loan-worthy

25 loan non-worthy

confusion matrix

	ML Yes	Prediction Yes	Prediction No	Actual Yes
Actual Yes	45	5	5	50
Actual No	10	20	20	30

Kuch Brisbane se RF

$$\frac{2 \times P \times R}{P + R} \rightarrow 2 \times Tn \quad \text{RF} \rightarrow$$

RR un

How

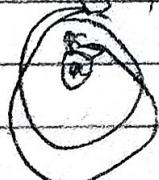
$$Bf = t - R \quad \text{Balap}$$

$$2 - 2 = 0 \quad 6\% \cos \cdot D \text{ at } 0^\circ$$

B-T Help

bz

AC 1 ✓



my KBST v logn

left ne R me

Small

Bad

AC = 100

1 PM = 60

(20)

-10, +1 ✓ Awt

$Bf = \text{height of hST} - \text{height of RST}$

Begin  
Re

0, I & no und

~~Evaluation Matrices~~

Ques → A bank want to predict whether a customer should get a loan (Yes/No) using an ML model.

Actual Yes → 12 Prediction Yes = 11

Actual No → 8 Prediction No = 9

Confusion Matrix

	Pred. Yes " 11 "	Pred. No " 9 "
A. Yes 12	10 TP	2 FN
A. No 8	1 FP	7 TN

Evaluation Matrix

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{10 + 7}{20} = \frac{17}{20} = 85\%$$

SML Kitne se kitne karte hain?

$$\text{Precision, TP} = \frac{10}{10 + 1} = 90.90\%$$

$$\text{Yes before } \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{10}{10 + 1}$$

me no del kitna sahiha

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{10}{10 + 2} = \frac{10}{12} = 83.33\%$$

mt no actual yes

me se kitne karte hain?

Seh me yes karte hain

F1-Score

$$\frac{2 \times \text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})} = \frac{2 \times 90.90\% \times 83.33\%}{90.90\% + 83.33\%}$$

$$1.650894 = 0.90710$$

$$1.66347 - 90.710$$

Q=1 Comparison b/w K mean and DB scan

→ K-mean is centroid based algorithm  
in K-mean Data are divided into K no. clusters

DB Scan Density Based algorithm

DBScan have two thing [parameter]

$\epsilon$  (Epsilon) → maximum distance to consider neighbors & Minpoint minimum no. of neighbors to form a dense region

→ where dense areas so cluster

Q=2 How to determine optimal number of clusters

→ Aim → To identify how much K are naturally present if  $K \rightarrow$  less no. of clusters  
vague & mixed if K more so undergo splitting

most common method

→ Elbow Method ✓

→ Here using diff. value of K e.g. K=100  
then calculate Within Cluster Summation

(Within K cluster) Square Sum

$WCSS = \sum_i (of K=100) j^2$

then make ~~WCSS~~ WCSS + K Graph

K	WCSS
1	1000
2	500
3	290
4	230
5	220