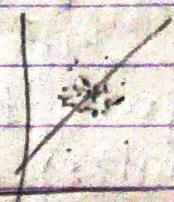


### Unit 3

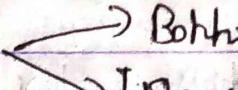
Clustering: It is technique of combining of datasets of similar type in groups.

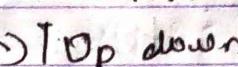
#### Density based

- Clusters are formed at dense regions.
- High accuracy
- Ability to merge two clusters in  e.g DBSCAN

#### Hierarchical based

Clusters are based formed as a tree type structure based on hierarchy

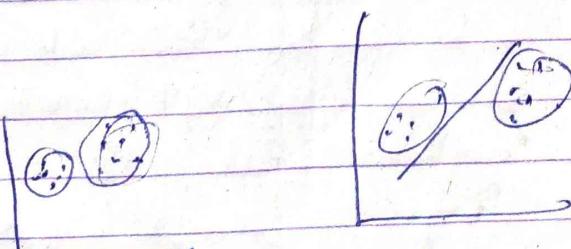
2 types  Bottom up



e.g : CURE (Clustering using representation)

#### Partitioning

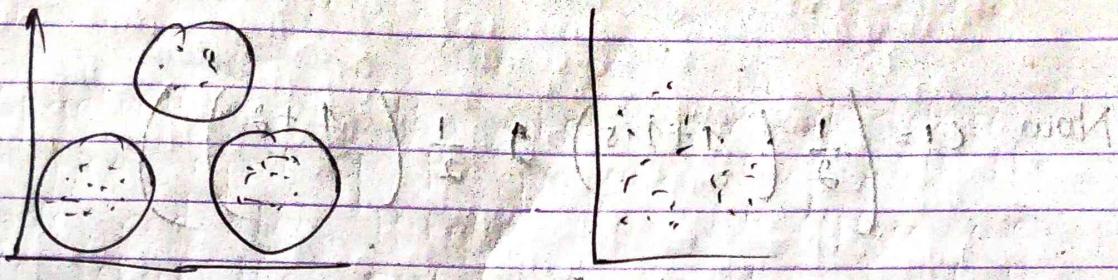
- Clusters are formed by partitioning into k clusters
- No of clusters = No of partitions.  
e.g k-means



K means clustering is an unsupervised learning algorithm that is used to solve clustering problem in ML or data science.

\* defines no of pre defined clusters, that need to be created in process.

Eg  $K = 2$  2 clusters



DS

Sum

$$C_1 = 1.0 \quad 1.0$$

$$C_2 = 3.0 \quad 7.0 + 9.4 + 2.1 + 2.1 + 2.1 + 2.1 + 2.1 + 2.1 + 2.1 + 2.1 + 2.1 + 2.1$$

Student

Var 1.

Var 2.

	x	y
1	4.0	(1.0 4.0 2.0)
2	3.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

$k=1$

Student

	C1
1	0
2	1.11
3	3.60
4	7.21
5	4.71
6	5.31
7	4.30

C2

7.21

6.10

3.60

0

2.5

2.06

2.91

$$\text{Now } C1 = \left( \frac{1}{2} (1+1.5) , \frac{1}{2} (4+2) \right)$$

$$\rightarrow 1.25, 1.5 \rightarrow$$

$$\text{new } C1 = \left( 1.25, 1.5 \right)$$

$$\text{Now } C2 = \left( \frac{1}{5} (3+5+3.5+4.5+3.5) , \frac{1}{5} (4+7+5+5+4.5) \right)$$

$$\rightarrow \left( 3.09, 5.12 \right)$$

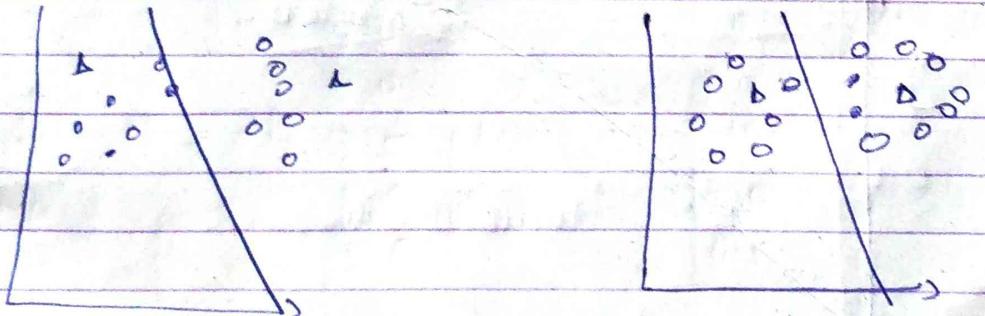
$k=2$

Student

	C1	C2
1	0.55	5.02
2	0.55	3.92
3	3.05	1.42
4	6.65	2.19
5	4.16	0.41
6	4.77	0.60
7	3.75	0.72

Working

1. Select the number of  $k$  to decide no of clusters
2. Select random  $k$  point or centroid
3. Assign each dataset to their closest centroid
4. Apply Euclidean distance to determine grouping of clusters
5. Then calculate the new centroid value by taking avg of values
6. Repeat step 3 then 4 and it goes on till the unchanged  $k$  value
7. The model is ready.



## Linear regression.

- It is machine learning algo based on supervised learning.
- Regression is mainly used to find the relationship between dependent & independent variables.
- Linear regression performs task to predict a dependent variable value ( $y$ ) based on given independent variable ( $x$ ).

Eg  $y \rightarrow$  dependent  $x \rightarrow$  Salary.  
(Single)

$$y = \alpha_0 + \alpha_1 x_1$$

$\alpha_0 \rightarrow$  Intercept.

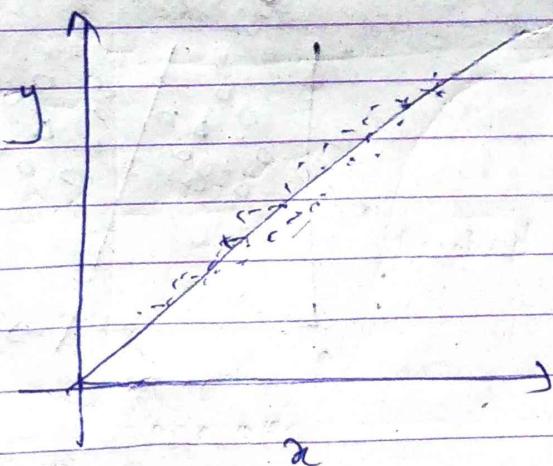
$\alpha_1 \rightarrow$  coefficient of  $x_1$ .

(multiple)

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$$

$\alpha_0 \rightarrow$  Intercept

$\alpha_i \rightarrow$  coefficient of  $x_i$



Higher value of  $\alpha$  more weightage is given in data training.

	$x$	$y$	$xy$	$x^2$	Predicted value	Brown
1	3	3	9	1	2.9	0.2
2	4	8	32	16	4.1	0.1
3	5	15	75	25	5.4	0.4
4	7	28	196	49	6.7	0.3
	$\sum x = 10$	$\sum y = 19$	$\sum xy = 54$	$\sum x^2 = 30$		

$$\text{Find eqn } y = bx + a$$

To find  $b$  and  $a$ .

$$a = \frac{(\sum y)(\sum x^2) - (\sum x)(\sum xy)}{n(\sum x^2) - (\sum x)^2}$$

$$b = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

$$a = \frac{(19)(30) - (10)(54)}{4(30) - (10)^2}$$

$$\Rightarrow \frac{570 - 540}{120 - 100} = \frac{30}{20} = 1.5$$

$$b = \frac{4(54) - (10)(19)}{4(30) - (10)^2}$$

$$\Rightarrow \frac{\cancel{54} - 190}{120 - 100}$$

$$= \frac{2b}{20} = \frac{13}{10} = 1.3$$

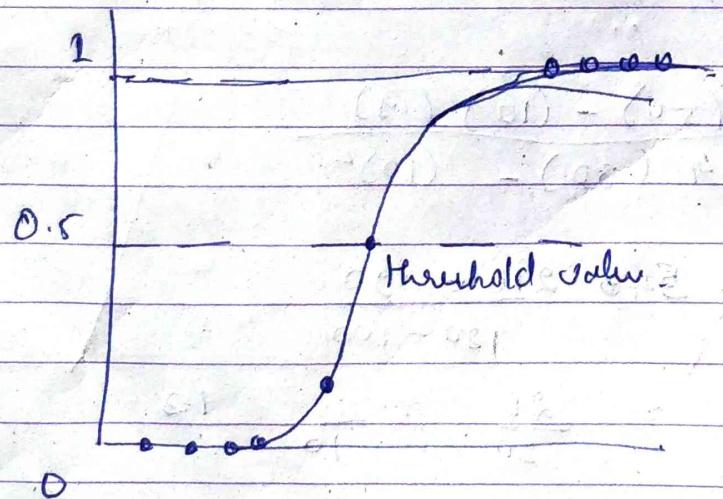
$$y = 1.3x + 1.5$$

What next?

Sub given value of  $x$  & find  $\phi$  or  
or find error in predicted & actual value

### Logistic Regression

- It is a supervised algorithm.
- It is used for predicting the categorical dependent variable using given set of independent variables.
- As it just predicts the output of categorical dependent variable. Therefore the outcome must be categorical or discrete.
- Like it can be either 0 or 1, true or false etc.
- It basically gives probabilistic values which lie in b/w 0 & 1.
- Logistic regression is similar to linear regression except the type of data used.
- The graph or curve is fit in 'S' shape if predicts 2 values (0 or 1).



## Sigmoid function

It is a function which helps us to map the values and get the outputs <sup>in probability</sup>  $[0 \text{ or } 1]$ .

$$\rightarrow Y = \frac{1}{1 + e^{-x}}$$

It is used for

- fraud detection
- disease diagnosis
- spam no spam
- emergency detection

## Measures

- The dataset should be free of missing values
- Absence of multicollinearity
- Linearity in variables
- Irrelevant variable bias
- functional form
- Structural Break