

## Association Rule

\* Association rule learning can be divided into  
3 kinds of algorithm

- \* Apriori
- \* Eclat
- \* T-P Growth

\* Association rule learning works on the concept  
of if and else statement

e.g.: If A then B.

where if A is antecedent  
then B is consequent.

Let us discuss about Apriori algorithm

\* A priori algorithm is generally considered as  
unsupervised learning algorithm

problem formula's

For support =  $\frac{\text{No. of items } x \& y \text{ bought together}}{\text{Total no. of transactions}}$

$$\text{support} = \frac{\text{fre}(x,y)}{N}$$

confidence =  $\frac{\text{Total no. of items } x \& y \text{ bought together}}{\text{Total no. of items 'x' bought}}$

$$\text{confidence} = \frac{\text{Freq}(x,y)}{\text{Freq}(x)}$$

$$\text{lift} = \frac{\text{support}(x,y)}{\text{support}(x) + \text{support}(y)}$$

Problem  
 Construct the following transactions apply the association rule mining to get the association rule with minimum support two and confidence of 50%.  
 → Association rule mining is explained using Apriori Algorithm  
 Transactional data for All electronics branch.

T <sub>ID</sub>	List of item IDs
T <sub>100</sub>	T <sub>1</sub> , T <sub>2</sub> , T <sub>5</sub>
T <sub>200</sub>	T <sub>2</sub> , T <sub>4</sub>
T <sub>300</sub>	T <sub>2</sub> , T <sub>3</sub>
T <sub>400</sub>	T <sub>1</sub> , T <sub>2</sub> , T <sub>4</sub>
T <sub>500</sub>	T <sub>1</sub> , T <sub>3</sub>
T <sub>600</sub>	T <sub>2</sub> , T <sub>3</sub>
T <sub>700</sub>	T <sub>1</sub> , T <sub>3</sub>
T <sub>800</sub>	T <sub>1</sub> , T <sub>2</sub> , T <sub>3</sub> , T <sub>5</sub>
T <sub>900</sub>	T <sub>1</sub> , T <sub>2</sub> , T <sub>3</sub>

### Step 1 :

Now scan each itemset and count how many candidate each has.

I-itemset	supcount
$\{I_1\}$	6
$\{I_2\}$	7
$\{I_3\}$	6
$\{I_4\}$	2
$\{I_5\}$	2

### Step 2

Check whether the candidate support count is with minimum of "two".

I-itemset	supcount
$\{I_1\}$	6
$\{I_2\}$	7
$\{I_3\}$	6
$\{I_4\}$	2
$\{I_5\}$	2

Step 3

Now Generate two candidate keys.

2 - itemset

$\{I_1, I_2\}$

$\{I_1, I_3\}$

$\{I_1, I_4\}$

$\{I_1, I_5\}$

$\{I_2, I_3\}$

$\{I_2, I_4\}$

$\{I_2, I_5\}$

$\{I_3, I_4\}$

$\{I_3, I_5\}$

$\{I_4, I_5\}$ .

Step 4 :

Now Scan the support count of candidate keys.

itemset      Supcount.

$\{I_1, I_2\}$       4

$\{I_1, I_3\}$       4

$\{I_1, I_4\}$       1

$\{I_1, I_5\}$       2

$\{I_2, I_3\}$       4

$\{I_2, I_4\}$       2

$\{I_2, I_5\}$       2

$\{I_3, I_4\}$       0

$\{I_3, I_5\}$       1

$\{I_4, I_5\}$       0.

5:

Check whether the candidate support count is with minimum of 'two', if not remove the candidate key. (from step 4)

2-itemset	supcount
$\{I_1, I_2\}$	4
$\{I_1, I_3\}$	4
$\{I_1, I_5\}$	2
$\{I_2, I_3\}$	4
$\{I_2, I_4\}$	2
$\{I_2, I_5\}$	2

Step 6:

Now Generate three candidate keys. with sup count

3-itemset	supcount	
$\{I_1, I_2, I_3\}$	2	$I_1, I_3, I_5$ 1
$\{I_1, I_2, I_5\}$	2	$I_1, I_4, I_5$ 0
$\{I_1, I_2, I_4\}$	1	$I_2, I_3, I_4$ 0
$\{I_1, I_3, I_4\}$	0	$I_2, I_3, I_5$ 1
		$I_2, I_4, I_5$ 0

Step 7:

check whether the candidate support count is with minimum of 'two'

3-itemset	supcount
$\{I_1, I_2, I_3\}$	2
$\{I_1, I_2, I_5\}$	2.

Step 8:

Generate four candidate keys.

itemset

freq count

$\{I_1, I_2, I_3, I_5\}$

1

So not possible, because it has only 1 freq count

\* Non-empty subsets of frequency sets are three item sets

$\{I_1, I_2, I_3\}$   $\{I_1, I_2, I_5\}$ .

$\rightarrow \{ (I_1), (I_2), (I_3), (I_1, I_2), (I_1, I_3), (I_2, I_3) \}$

$\rightarrow \{ (I_1), (I_2), (I_5), (I_1, I_2), (I_1, I_5), (I_2, I_5) \}$ .

\* Association rule for every non-empty subsets

$S \Rightarrow (I - S)$

\* Let us consider, Association rule mining subset creation between three items sets.

$I \Rightarrow \{I_1, I_2, I_3\}$ .

Non-empty subsets

$\{ (I_1) (I_2) (I_3) (I_1, I_2) (I_1, I_3) (I_2, I_3) \}$

Rule 1 :

$$\begin{array}{ccc} \{1, 2\} & \rightarrow & \{1, 2, 3\} \\ \downarrow & & \downarrow \\ S & & T-S \end{array}$$

$$\text{Support} = \frac{2}{9} = 22.22\%$$

$$\text{Confidence} = \frac{\text{Support}(1, 2, 3)}{\text{Support}(1)} = \frac{2/9}{6/9} = \frac{1}{3} = 33.33\%$$

R ∵ The condition is invalid ( $< 50\%$ )

Rule 2

$$\{2, 3\} = \{1, 3\}$$

$$\text{Support} = \frac{2}{9} = 22.22\%$$

$$\text{Confidence} = \frac{\text{Support}(1, 2, 3)}{\text{Support}(2)} = \frac{2/9}{7/9} = \frac{2}{7} = 28.57\%$$

∴ The condition is invalid ( $< 50\%$ )

Rule 3

$$\{3\} = \{1, 2\}$$

$$\text{Support} = \frac{2}{9} = 22.22\%$$

$$\text{Confidence} = \frac{\text{Support}(1, 2, 3)}{\text{Support}(3)} = \frac{2/9}{6/9} = \frac{2}{6} = \frac{1}{3} = 33.33\%$$

∴ The condition is invalid ( $< 50\%$ )

Rule 4

$$\{1, 2\} = \{3\}$$

$$\text{Support} = \frac{2}{9} = 22.22\%$$

$$\text{Confidence} = \frac{\text{Support}(1,2,3)}{\text{Support}(1,2)} = \frac{2/9}{4/9} = \frac{2}{4} = 50\%.$$

$\therefore$  The above condition is valid ( $\geq 50\%$ ).

Ex 5

$$\{1,3\} \subseteq \{2,3\}$$

$$\text{Support} = \frac{2}{9} = 22.22\%.$$

$$\text{Confidence} = \frac{\text{Support}(1,2,3)}{\text{Support}(1,3)} = \frac{2/9}{4/9} = \frac{2}{4} = 50\%.$$

$\therefore$  The above condition is valid ( $\geq 50\%$ ).

Ex 6

$$\{2,3\} \subseteq \{1,3\}$$

$$\text{Support} = \frac{2}{9} = 22.22\%.$$

$$\text{Confidence} = \frac{\text{Support}(1,2,3)}{\text{Support}(2,3)} = \frac{2/9}{7/9} = \frac{2}{7} = 28.57\%.$$

$\therefore$  The above condition is invalid ( $< 50\%$ ).

**CLASSIFICATION:**

- \* Classification is a process of categorizing a given set of data into classes.
- \* The classes are often referred to as target, label or categories.
- \* Classification is defined as the process of recognition, understanding, and grouping of objects and ideas to proceed categories.

In classification algorithms are used to predict the output of the categorical data.

**CLASSIFIER:**

The algorithm which implement the classification on a dataset is known as "classifier."

**NAIVE BAYESIAN CLASSIFICATION:**

- \* Invented by Reverend Thomas Bayes in 1761.
- \* Naive Bayes assumes that all attributes are equally important, independent of one another in a given class.

**BAYES THEOREM:**

Bayes' Theorem is one of the most popular machine learning concept that helps to calculate the probability of occurring one event with uncertain knowledge while others are already occurred.

$$P(\text{class}|\text{data}) = \frac{(P(\text{data|class}) * P(\text{class}))}{P(\text{data})}$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$P(A|B)$  - The probability of 'A' being true given that 'B' is true.

$P(A)$  - The probability of 'A' being true.

7/2/23

P.T.W.S

Dataset: Weather condition.

Problem: If the weather is sunny, then the player should play or not?

	outlook	play
0	Rainy	yes
1	Sunny	yes
2	overcast	yes
3	overcast	yes
4	sunny	yes
5	Rainy	yes
6	Sunny	yes
7	overcast	yes
8	Rainy	No
9	Sunny	No
10	Sunny	No
11	Rainy	No
12	overcast	yes
13	overcast	yes

FREQUENCY TABLE:

Weather	Yes	No
overcast	5	0
Rainy	2	2
sunny	3	2

Likelihood Table:

Weather	Yes	No	
overcast	5	0	$\frac{5}{14} = 0.35$
Rainy	2	2	$\frac{2}{14} = 0.29$
sunny	3	2	$\frac{5}{14} = 0.35$

$$\text{All } \quad 10/14 \quad 4/14 = 0.29 \\ - 0.71$$

Applying Bayes' Theorem: —

$$P(\text{player}/\text{sunny}) = \frac{P(\text{sunny}|\text{yes}) * P(\text{yes})}{P(\text{sunny})}$$

$$P(\text{sunny}|\text{yes}) = 3/10 = 0.3$$

$$P(\text{sunny}) = 0.35$$

$$P(\text{yes}) = 0.71$$

$$\therefore P(\text{player}/\text{sunny}) = \frac{0.3 \times 0.71}{0.35} = 0.60$$

$$P(\text{no}/\text{sunny}) = P(\text{sunny}/\text{no}) * P(\text{no}) \\ P(\text{sunny})$$

$$P(\text{sunny}/\text{no}) = \frac{9}{4} = 0.5$$

$$P(\text{no}) = 0.29$$

$$P(\text{sunny}) = 0.35$$

$$\therefore P(\text{no}/\text{sunny}) = \frac{0.5 \times 0.29}{0.35} = 0.41$$

$$\therefore P(\text{player}/\text{sunny}) > P(\text{no}/\text{sunny})$$

Hence on sunny days player can play the game.

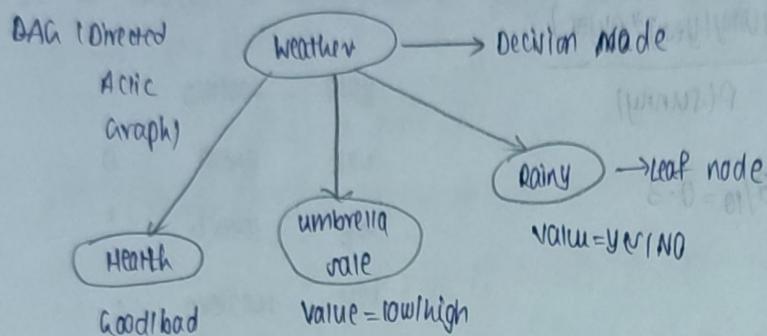
13/2/23

BAYERIAN NETWORK:

\* Bayesian networks are the type of probabilistic graphical model that can be used to build models from data.

+ They are also commonly referred as Bayes' Network, Belief Net.

Bayesian networks are probabilistic because they are built from probability distribution.



- Probability of Good health based on weather.

- Probability of umbrella sale based on weather.

- Random variables (weather, health) can give a hypothesis.

Types of probabilities:

1. Joint probability

2. Conditional probability.

DECISION TREE:

\* Decision Tree are some of the most used machine learning algorithm in supervised learning.

\* Decision Trees are used in both classification & Regression.

\* They can be used for both linear and non-linear data but they are mainly used in non-linear data.

YET is a tree structured classifier.

1/12/23

Predict the probability of a class variable using Bayesian classification.

Training dataset: All electronics customer database

Variable (X) = (age = youth, income = Medium, student = yes, credit = rating = fair)

ID	Age	Income	Student	Credit Rating	buys-computer
1.	youth	high	no	fair	no
2.	youth	high	no	Excellent	no
3.	Middle age	high	no	fair	yes
4.	senior	medium	no	fair	yes
5.	senior	low	yes	fair	yes
6.	senior	low	yes	Excellent	no
7.	Middle age	low	yes	Excellent	yes
8.	youth	Medium	No	fair	no
9.	youth	Low	yes	fair	yes
10.	senior	Medium	yes	fair	yes
11.	youth	Medium	yes	Excellent	yes
12.	Medium Age	Medium	no	Excellent	yes
13.	MA	High	yes	fair	yes
14.	senior	Medium	no	Excellent	no

Class:  $C_1$  : buys-computer = yes

$C_2$  : buys-computer = no

probability:

$$P(C_1) = P(\text{buys-computer} = \text{Yes}) = \frac{9}{14} = 0.643$$

$$P(C_2) = P(\text{buys-computer} = \text{No}) = \frac{5}{14} = 0.357$$

Compute  $P(X|C_i)$  for each class:

$$P(\text{age} = \text{young} | \text{buys-computer} = \text{yes}) = \frac{2}{9} = 0.222$$

$$P(\text{age} = \text{medium} | \text{buys-computer} = \text{yes}) = \frac{4}{9} = 0.444$$

$$P(\text{student} = \text{yes} | \text{buys-computer} = \text{yes}) = \frac{6}{9} = 0.667$$

$$P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{yes}) = \frac{6}{9} = 0.667$$

my

$$P(\text{age} = \text{young} | \text{buys-computer} = \text{no}) = \frac{3}{5} = 0.6$$

$$P(\text{income} = \text{medium} | \text{buys-computer} = \text{no}) = \frac{2}{5} = 0.4$$

$$P(\text{student} = \text{yes} | \text{buys-computer} = \text{no}) = \frac{1}{5} = 0.2$$

$$P(\text{credit-rating} = \text{fair} | \text{buys-computer} = \text{no}) = \frac{2}{5} = 0.4$$

$$\therefore \text{so } P(X|C_1) = P(\text{age} = \text{young} | \text{buys-computer} = \text{yes})$$

$$= 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.004$$

$$P(X | \text{buys-computer} = \text{no}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

$$P(X|C_1) \times P(C_1)$$

$$P(X | \text{buys-computer} = \text{yes}) \times P(\text{buys-computer} = \text{yes})$$

$$= 0.004 \times 0.643 = 0.00257$$

$$P(X | \text{buys-computer} = \text{no}) \times P(\text{buys-computer} = \text{no})$$

$$= 0.019 \times 0.357 = 0.007$$

RESULT:

$P(X | \text{buys-computer} = \text{yes}) > P(X | \text{buys-computer} = \text{no})$  greater than  $P(X | \text{buys-computer} = \text{yes})$

$$P(\text{buys-computer} = \text{no})$$

The variable X (age = youth, income = medium, student = yes, credit - ranking = fair),

Have the probability of buying a computer.

15/2/23

### ASSOCIATION RULE:

\* Association rule learning is a type of supervised learning in machine learning.

\* It checks the dependency of one data item on another data item and maps accordingly.

\* It is mostly used in questions, decision making and predict behaviour.

\* It is a rule based machine learning method for discovering interesting relations between variables in large database.

### ASSOCIATION RULE - APRIORI ALGORITHM:

\* Apriori algorithm is used to find the frequent itemset in a dataset for boolean association rule.

\* Apriori means prior knowledge of frequent item set properties.

\* Apriori is an algorithm used for association rule mining.

\* Types of Association are multi-relational, quantitative and generalized.

### APRIORI ALGORITHM MEASUREMENTS:

There are 2 measured mostly used in association rule.

1. support.

$$\text{support} = \frac{\text{frequency}(x,y)}{N}$$

2. confidence

$$\text{confidence} = \frac{\text{frequency}(x,y)}{\text{frequency}(x)}$$

Example: Apriori Algorithm.

The following are the set of items transacted in a market. Find the support and confidence?

Trans ID	
101	milk, bread, eggs
102	milk, juice
103	Juice, Butter
104	milk, bread, eggs
105	coffee, eggs
106	coffee
107	coffee, juice
108	milk, bread, cookies, eggs
109	cookies, butter
110	milk, bread,

Index the item sets

Item	Numbers
Milk	1
Bread	2
Eggs	3
Juice	4
Butter	5
Coffee	6
Cookies	7

Find the overall support (Total item = 10 = 100%)

Itemset	support	%
Milk	5	50%
Bread	4	40%
Eggs	4	40%
Juice	3	30%
Butter	2	20%
Coffee	3	30%
cookies	2	20%

find the support for 2 item set from Index.

Itemset	support	%	
{1,2}	4	40%	
{1,3}	3	30%	
{1,4}	1	10%	
{1,5}	0	-	
{1,6}	0	-	
{1,7}	1	10%	
{2,3}	3	30%	
{2,4}	0	-	
{2,5}	0	-	
{2,6}	0	-	$\frac{8}{8} = \frac{(6,2,1) \text{ freq}}{(6,1) \text{ freq}}$
{2,7}	1	10%	
{3,4}	0	-	$\frac{8}{8} = \frac{(8,1,1) \text{ freq}}{(8,1) \text{ freq}}$
{3,5}	0	-	
{3,6}	1	10%	$\frac{8}{8} = \frac{(8,1,1) \text{ freq}}{(8,1) \text{ freq}}$
{3,7}	1	10%	
{4,5}	1	10%	
{4,6}	1	10%	
{4,7}	0	-	
{5,6}	0	10%	
{5,7}	1	10%	$\frac{8}{8} = \frac{(8,1,1) \text{ freq}}{(8,1) \text{ freq}}$
{6,7}	0	-	

Itemset	support
{1,2}	4
{1,3}	3
{2,3}	3

Itemset	support
{1,2,3}	3

use 20% for the minimum support

CONFIDENCE: 70% or greater in the item set

$$\frac{\text{support}(1,2,3)}{\text{support}(1)} = \frac{3}{5} = 60\% < 70\% \quad \times$$

$$\frac{\text{support}(1,2,3)}{\text{support}(2)} = \frac{3}{4} = 75\% > 70\% \quad \checkmark$$

$$\frac{\text{support}(1,2,3)}{\text{support}(3)} = \frac{3}{9} = 33\% > 70\% \quad \checkmark$$

confidence  $\Rightarrow$ :

R

	confidence
Bread $\rightarrow \{milk, eggs\}$	75%
Eggs $\rightarrow \{milk, bread\}$	75%

$$\frac{\text{support}(1,2,3)}{\text{support}(1,2)} = 75$$

$$\frac{\text{support}(1,2,3)}{\text{support}(1,3)} = 100\%$$

$$\frac{\text{support}(1,2,3)}{\text{support}(2,3)} = 100\%$$

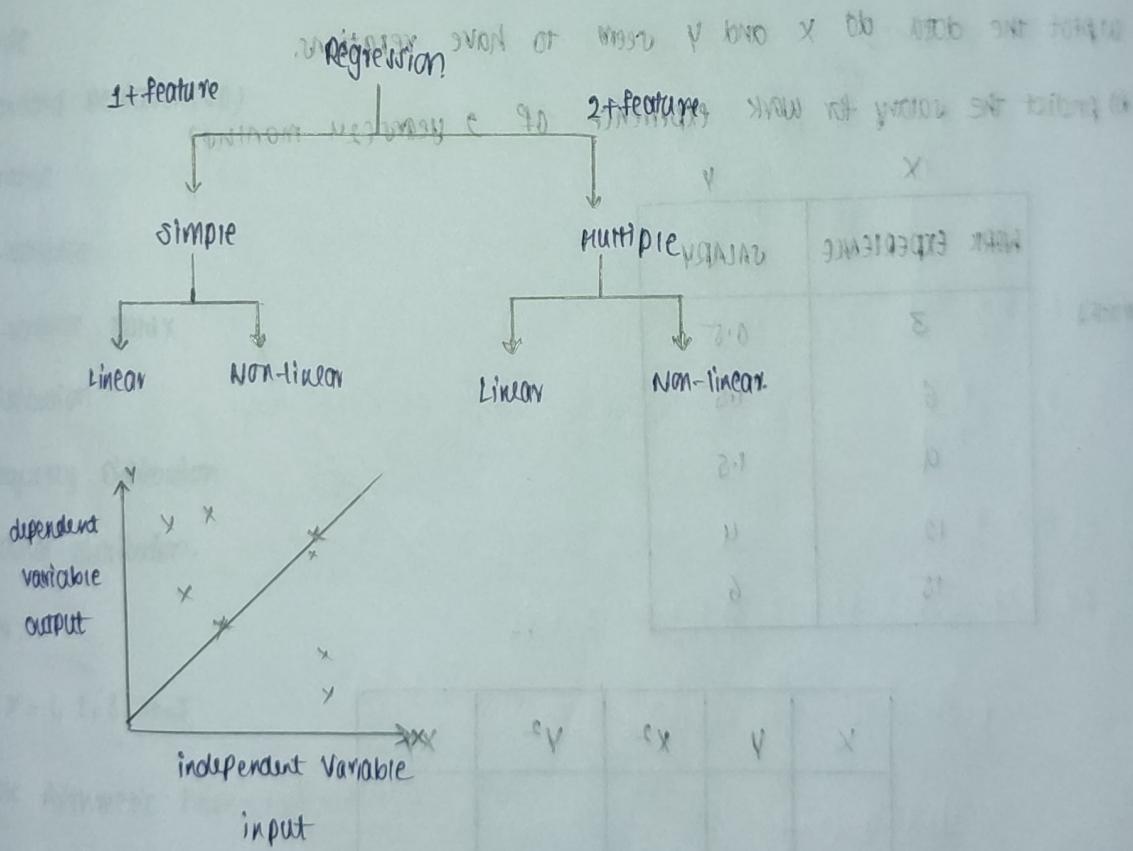
$S(2,3)$

16/1/23

## REGRESSION:

- \* Regression is a statistical method to determine the strength and character of the relationship b/w one dependent variable (usually denoted by  $y$ ) and a series of the other variables (known as independent variable  $x$ )
- \* Regression is a supervised Machine learning algorithm.
- \* Regression analysis is the process of estimating the relationship b/w a dependent variable and two independent variables.
- \* Regression will predict a continuous outcome ( $y$ ) based on the value of one or more predictive variables ( $x$ )

## REGRESSION MODELS:



## SIMPLE LINEAR REGRESSION:

Simple linear regression is a technique which plots a straight line within the data points. It is one of the most simple and basic type of machine learning regression.

## MULTIPLE LINE REGRESSION:

It is a technique used when more than one independent variable is used. Polynomial is an example of a multilinear regression.

## SIMPLE LINEAR REGRESSION:

- \* One variable ( $x$ ) is called independent variable or predictor
- \* The other variable ( $y$ ) is known as dependent variable or outcome, and the simple linear regression equation is

$$Y = b_0 + b_1 x$$

What is linear regression? The following table shows the work experience and salary obtained from employees.

a) Plot the data do  $x$  and  $y$  seem to have relations.

b) predict the salary for work experience of 2 years(24 months)

X	Y
WORK EXPERIENCE	SALARY
3	0.5
6	1.5
9	1.5
12	4
15	6

X	Y	$X^2$	$Y^2$	$XY$
$\Sigma X = 45$	13.5	162	49.5	56.75

$$y = b_0 + b_1 x$$

$$b_0 = \frac{\sum xy - \bar{x}\bar{y}}{n \cdot \bar{x}^2 - (\bar{x})^2}$$

$$b_1 = \frac{n \cdot \bar{x}\bar{y} - \sum xy}{n \cdot \bar{x}^2 - (\bar{x})^2}$$

$$y = -1.35 + 0.45x$$

$$= 0.45$$

20/2/23

## STATISTICS:

1. simple Arithmetic mean.

2. Median

3. Mode

4. Standard Deviation (SD)

5. Variance

6. Co-Variance

7. Co-Variance Matrix

8. Distribution

9. Frequency Distribution

10. Normal distribution.

## MEAN:

$$x = 1, 2, 3, 4, 5$$

$$\text{Simple Arithmetic Mean} = \frac{1+2+3+4+5}{5}$$

$$= \frac{15}{5} = 3$$

## MEDIAN:

50<sup>th</sup> percentile position

$$\text{Position} = N+1/2$$

$$Y = 15, 10, 8, 12, 14$$

Arrange in highest to lowest 15, 14, 12, 10, 8

$$N=5 \quad \frac{N+1}{2} = \frac{6}{2} = 3^{\text{rd}} \text{ position}$$

$$= 12$$

MODE: No. of occurring of data.

$$3, 5, 7, 10 - \text{No mode}$$

$$3, 5, 3, 7, 3, 10, 3 \quad \text{Mode} = 3$$

POPULATION VARIANCE:

$$\sigma^2 = \frac{\sum (X - \bar{X})^2}{n} \rightarrow \begin{array}{l} \text{variable} \\ \text{Mean} \\ \text{No. of score} \end{array}$$

SAMPLE VARIANCE:

$$\sigma^2 = \frac{\sum (X - \bar{X})^2}{n-1}$$

$$\text{mean } X = 3$$

$$PV = (1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2 / 5$$

$$= (-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2 / 5$$

$$= \frac{4+1+0+1+4}{5}$$

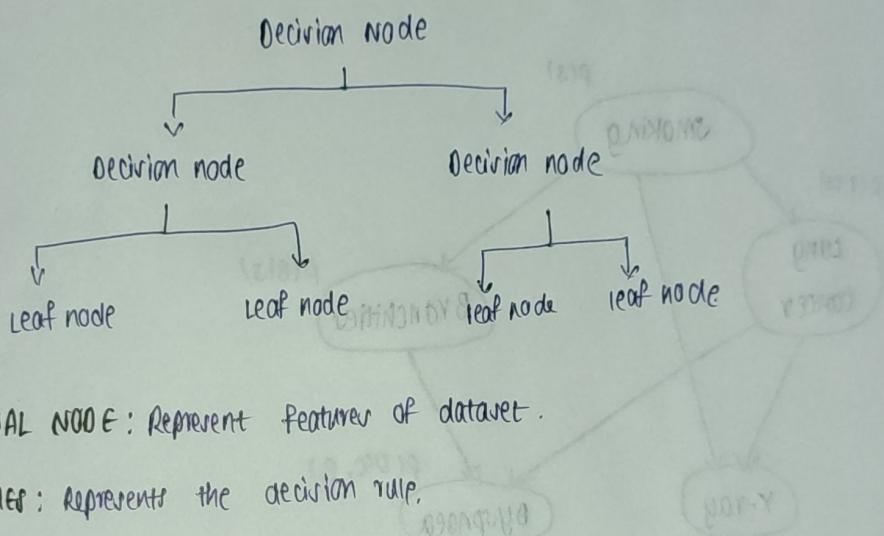
$$\text{P.C.V} = \frac{(X - \bar{X})^2 + (Y - \bar{Y})^2}{n}$$

Standard deviation: Square root of variance.

2/12/23

## DECISION TREE:

- \* Decision Tree are one of the most used machine learning algorithm.
- \* Decision Tree is a supervised learning algorithm.
- \* Decision Tree can be used in both classification and regression techniques.
- \* Decision Tree can be used for both linear (continuous) and non-linear (discrete) data but mainly used in non-linear data.



**INTERNAL NODE:** Represent features of dataset.

**BRANCHES:** Represents the decision rule.

**LEAF NODE:** Represents the outcome.

## BAYESIAN NETWORK:

Bayesian Networks are a type of probabilistic graphical model that can be used to build models from data. They are also commonly referred to as Bayesian Network, Belief Network, or BN model.

\* Bayesian Network are probabilistic because they are built from probability distribution.

2 types of probability

1. Joint probability

2. conditional

28/12/23

DATA MODELLING

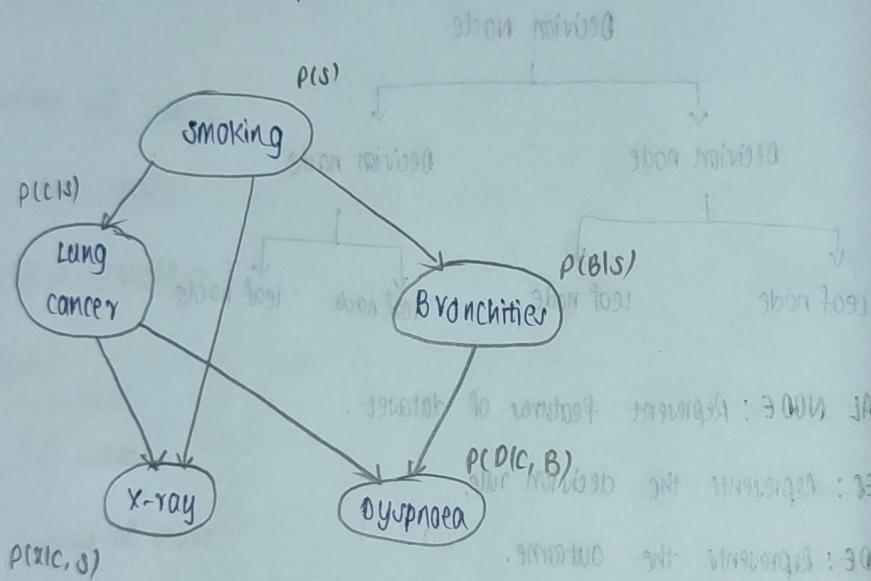
## BAYESIAN NETWORK:

Structured Graphical representation of probabilistic relationship b/w several random variables.

### PROBABILITIES:

1. Joint probability

2. Conditional probability distribution (cpd)



BAYESIAN NETWORK (BN =  $G, \theta$ )

$\rightarrow$  Directed Acyclic Graph (DAG) nodes - random variables  
edges - direct dependencies

$\theta$ -set of parameters in all conditional probability Distribution (cpd)

cpd:-

C	B	$D=0$	$D=1$
---	---	-------	-------

0	0	0.1	0.9
---	---	-----	-----

0 - False, 1 - True

0	1	0.7	0.3
---	---	-----	-----

1	0	0.8	0.2
---	---	-----	-----

1	1	0.9	0.1
---	---	-----	-----

$$\text{CPD of node } X = p(X | \text{parents}(X))$$

X - random variable (smoking, lung cancer, Bronchitis, X-ray, dyspnoea)

compact representation of Joint Distribution in a product form

$$p(s, c, B, X, O) = p(s) p(c|s) p(B|s) p(X|c, s) p(O|c, B)$$

$$= 1 + 2 + 2 + 4 + 4 = 13 \text{ parameters instead of } 2^5 = 32$$