③

Mean :-
arithmetic average of a dataset
• adding no.an dataset divided by obser

$$\bar{x} = \frac{\Sigma x}{N}$$

mean     $\Sigma x \longrightarrow$ observation
$N \longrightarrow$ no. of observation

② mode :-
Most frequently occuring observation or value.

$$mode = 1 + \left[ \frac{fm - f_1}{2fm - f_1 + f_2} \right] \times h$$

fm = freq. possereb by modal class
f₁ = class before
f₂ = class after modal class
h = width

3) Median
Middle number in a data set
When no listed either assending and desending order.

for odd obs.
$$Media = \left(\frac{n+1}{2}\right)^{th} observation$$

for even
$$Median = \frac{(n/2)^{th} obs + (n/2 + 1)^{th} obs}{2}$$

## smoothing :→

Process that is used to remove noise from the dataset using some Algorithms

It helps in predicting the Patterns

## Aggregation :→

is a nothing but data collection is the method of storing and presenting data in a summary format

## Generalization :-

It converts low level data attributes to high-level attribute using concept hierarchy.

eg. Age in Numerical form (20, 21) converted into categorical value (young, old)

## Normalization :-

Data Normalization involves converting all data variables into a given range.

## Min - Max normalization

Transfer original data linearly

## Z-score normalization :-

In Z score Normalization the value of an attribute (A), are normalized based on the mean of A and its S.D.

Standard deviation: sqre of variance

$$s = \sqrt{\frac{\varepsilon(\bar{x} - x)^2}{n-1}}$$

$\sigma^2$

$x$ — each value

$\bar{x}$ — Mean

$n =$ no. of value

### (1) Linear Regression:-

Machine learning algorithm based on supervised learning. It targets prediction vaguer on the basis of independent variables.

### MultiVariate Regression:-

it concerns the study of two or more predictor variable

$$y = a + bx + cx_2$$

### Mean square Error:-

represent r the error of the estimator or predictive model created based on the given set of observation in the sample

$$MSE = \frac{1}{n} \varepsilon \left( y - \bar{y} \right)^2$$

sq of diff betn
actual & predicted

⑤ Logistic regression

## Logistic Regression :-

classification techniques are an essential part of Machine learning and data mining application.

- is one of the most simply & commonly used machine learning algorithm for two class classification.

statistical method for predicting binary class
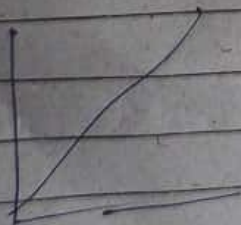
| linear | logistic |

continuous o/p                    constant o/p

eg. House price                eg. patient has a cancer or not

estimated using            estimated using
Ordinary Least           Maximum Likelihood
Square (OLS)             Estimation (MLE)

Accuracy :- calculated or the number of correctly instance divided by total number of instance.

$$acc = \frac{Tp + TN}{Tp + Fp + TN + FN} = \frac{Tp + TN}{Pos + Neg}$$

error Rate :- calculated as the no of incorrectly clarrified instance divided by the total Number of instance

$$error = 1 - acc.$$

Precision :-

calculated as the no. of correctly clarrified tve instance divided by the total no. of instance Which are predicted positive. it is also called Confidence

$$precision = \frac{Tp}{Tp + Fp}$$

Recall :-

$$= \frac{Tp}{Tp + Fn}$$

**Q.** Types of Logistic Regression

Binary Logistic Regression:-
The target variable has only two possible outcomes such as spam or Not spam, Cancer or No Cancer.

Multinomial Logistic Regression:-
The target variable has three or more nominal categories such as predicting the type of wine.

Ordinal logistic Regression:-
The target variable has three or more ordinal categories such as restaurant or product rating from 1 to 5

Confusion Matrix
It contains information about actual and predicted classification done by Classification system.

|  | Predicted | |
|---|---|---|
| actual | TP | FN |
| | FP | TN |

Confusion Matrix

(6) Naive Bayes:- used for classification of categorical data

prior probabilities:-
Calculate for some event base on no other information

Bayes Rule :- $P\left(\dfrac{A}{B}\right) = \dfrac{P(B/A)\, P(A)}{P(B)}$

(7) Text Analysis
is a deriving & extending high quality information from text

Pre-processing Method in Text analysis
1. Tokenization.
The process of breaking text paragraph into small chunk or sentence