

(2)

Mean:-

arithmetic average of a dataset
• adding no. in dataset divided by obser

$$\bar{x} = \frac{\sum x}{N} \rightarrow \text{observation}$$

mean $\frac{\sum}{N} \rightarrow$ no. of observation

(3) Mode:-

Most frequently occurring observation or value

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

f_m = freq. possessed by modal class
 f_1 = class before
 f_2 = class after modal class
 h = width

(3) Median

Middle number in a dataset
when no listed either ascending
and descending order.

for odd obs.

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ observation}$$

for even

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

Smoothing: →

Process that is used to remove noise from the dataset using some algorithms. It helps in predicting the patterns.

Aggregation: →

is nothing but data collection. It 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, 30, 40) 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: sq. of variance

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

\bar{x} - each value

\bar{x} - Mean

n = no. of value

①

Linear Regression:-

Machine learning algorithm based on supervised learning

It targets prediction value 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 the error of the estimator or predictive model created based on the given set of observation in the sample

$$MSE = \frac{1}{n} \sum (y - \hat{y})^2$$

Standard deviation

Actual - predicted

an and data

ly 4 common for two

king binary

har q

r not

Using likelihood (MLE)

⑥
⑤
④

⑤ Logistic regression

Logistic Regression:-

classification techniques are an essential part of machine learning and data mining application.
• is one of the most simple & commonly used machine learning algorithm for two class classification.

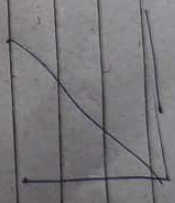
Statistical method for predicting binary class

linear logistic

continuous o/p constant o/p

eg. House price cancer survival

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



Accuracy:- Calculated as 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 classified instance divided by the total number of instance

$$error = 1 - acc$$

precision:-

Calculated as the no. of correctly classified true instance divided by the total no. of instance which are predicted positive. it is also called Confidence

$$precision = \frac{TP}{TP+FP}$$

Recall:-

$$recall = \frac{TP}{TP+FN}$$

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.

actual \ predicted	TP	FN
	FP	TN

Confusion Matrix

(6)

Naive Bayes:- Used for classification of categorical data

Prior probabilities:-

Calculate for some event based on no other information

$$\text{Bayes Rule}:- P\left(\frac{A}{B}\right) = \frac{P(B|A)P(A)}{P(B)}$$

(7)

Text Analysis

is deriving & extending high quality information from text

high quality information from text

Pie-processing

Method in Text analysis

1. Tokenization

The process of breaking text paragraph into small chunk or sentence