



Experiment Assignment No. 03

Title of the Assignment:

Descriptive Statistics - Measures of Central Tendency & variability.

Perform the following operations on any open source dataset (eg. data.csv)

1. Provide summary statistics (mean, median, minimum, maximum, standard deviation) for a dataset (age, income, etc.) with numeric variables grouped by one of the qualitative (categorical) variables. For example, if your categorical variable is age groups & quantitative variable is income, then provide summary statistics of income grouped by the age groups. Create a list that contains a numeric value for each response to the categorical variable.
2. Write a python program to display some basic statistical details like percentile, mean, standard deviation etc. of the species of 'Iris-setosa', 'Iris-versicolor' & 'Iris-virginica' of Iris.csv dataset.

Provide the codes with o/p & explain everything that you do in this step.



Objective of the Assignment →

Students should be able to perform the statistical operations using python on any open source dataset.

Prerequisite :→

1. Basic of python programming
2. Concept of statistics such as mean, median, minimum, maximum, standard deviation etc.



Experiment No. 04

Title of the Assignment :-

Create a Linear Regression Model using python / R to predict home prices using Boston Housing Dataset (<http://www-kaggle.com/c/boston-housing>).

The Boston housing dataset contains information about various houses in Boston through different parameters. There are 506 samples & 14 feature variables in this dataset.

The objective is to predict the value of prices, of the house using the given features.

Objectives of the Assignment :-

Students should be able to do data analysis using linear regression using python for any open source dataset.

Prerequisite :-

1. Basic of python programming.
2. Concept of Regression.

Theory :-

(1) Linear Regression :-

- It is a mlc learning algo. based on supervised learning. It targets predicts values on the basis of independent variables.



- It is preferred to find out the relationship between forecasting & variables.
- A linear relationship between a dependent variable (y) & continuous; while relationship should be available to independent variable (x). relationship may be continuous. A linear relationship should be available to "predictor" & target variable & so known as linear regression.
- If it is shown as an eqn of line like:

$$Y = mX + b + e$$

where,

$b \Rightarrow$ Intercepted

$m \Rightarrow$ Slope of line

$e \Rightarrow$ error term

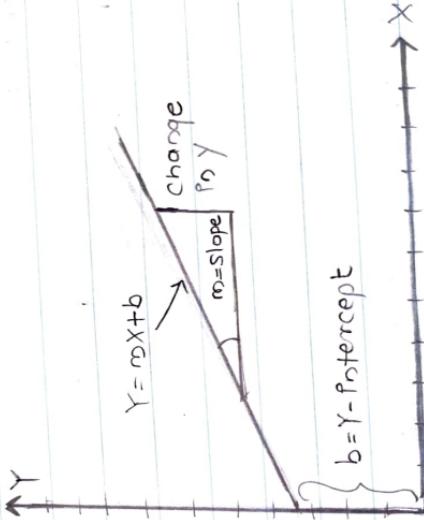


Fig. 01 : Geometry of linear regression



- This eqⁿ can be used to predict the value of target variable Y based on given predictor variables(X). as shown in Fig.01.
- Fig.02 shown below is about the relation betⁿ weight (in kg) & height (in cm), a linear relation. It is an approach of studying in a statistical manner to sum upse of all the relationships among continuous (quantitative) variables.
- Here a variable, denoted by 'x' is considered as predictor, explanatory or independent variable. & Y is considered as the response, outcome or dependent variable.

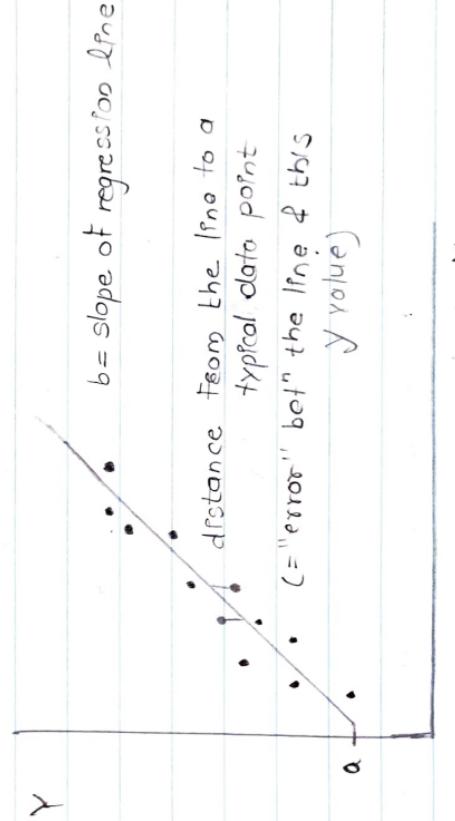


Fig.02: Relation betⁿ weight (in kg) of height (in cm)



Multivariate Regression?

- It concerns the study of two or more predictor variables.
- A simple linear model $y = a + bx$ is transformed into original feature will be transformed into polynomial feature is transformed & further a linear regression applied to fit & if will be something like

$$y = a + bx + cx^2$$

(2) Least Square Method for Linear Regression

- A simple linear model is the one which involves only one dependent & one independent variable. Regression Model are usually denoted in matrix notations.
- However for a simple univariate linear model it can be denoted by the regression eq?

$$\hat{y} = \beta_0 + \beta_1 x \quad \text{--- (1)}$$

where,

$\hat{y} \Rightarrow$ dependent or response variable

$x \Rightarrow$ Independent or r/p variable

$\beta_0 \Rightarrow$ value of y when $x=0$ or y intercept.

$\beta_1 \Rightarrow$ value of slope of line & is the error or the noise.



- This regression eqⁿ represents a line also known as the regression line!
- This technique estimate parameters β_0 & β_1 by trying to minimise the square of errors at all the points in the sample set. The error is the deviation of the actual sample data point from the regression line.
- The technique can be represented by the eqⁿ.

$$\min \sum_{i=0}^n (\hat{y}_i - y_i)^2 \quad \text{--- (2)}$$

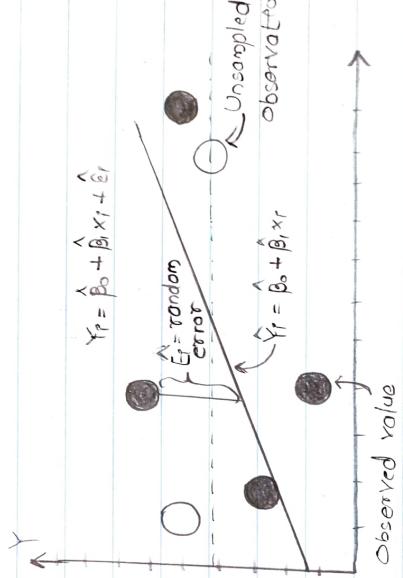
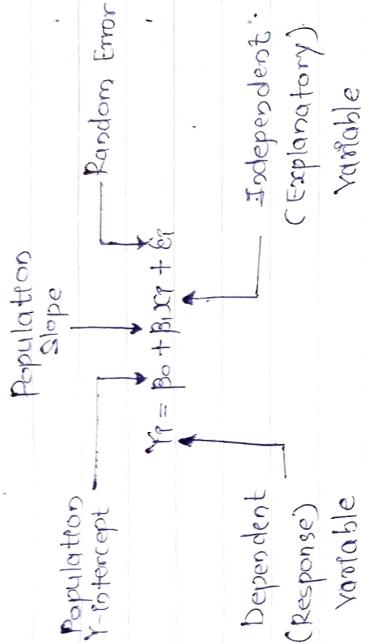


Fig. 09



Using differential calculus on eqⁿ if we can find the values of β_0 that the sum of squares (that is equation 2) is minimum.

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{--- (3)}$$

$$\beta_0 = \bar{Y} - \beta_1 \bar{x} \quad \text{--- (4)}$$

(8) Measuring Performance of linear regression :

Mean Square Error : (MSE)

It represents the error of the estimation or predictive model created based on the given set of observations in the sample.

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

The square of difference b/w actual & predicted



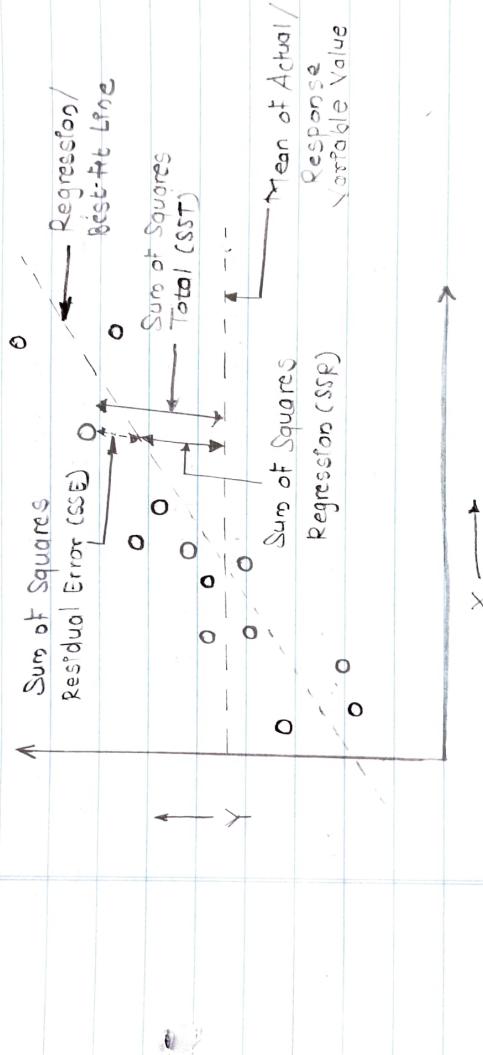
An MSE of zero(0) represents the fact that the predictor is a perfect predictor.

RMSE (Root Mean Squared Error) :

This is a method that calculate the least-squares error & take root of summed values

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}$$

R-Squared :



R-Squared is the ratio of sum of squares regression (SSR) & the sum of squares total (SST).



$$SST = \sum_{p=1}^n (\gamma_p - \bar{\gamma})^2$$

$$SSR = \sum_{p=1}^n (\hat{\gamma}_p - \bar{\gamma})^2$$

$$SSE = \sum_{p=1}^n (\gamma_p - \hat{\gamma}_p)^2$$

$$R^2 = \frac{SSR}{SST} = \frac{\sum (\hat{\gamma}_p - \bar{\gamma})^2}{\sum (\gamma_p - \bar{\gamma})^2}$$

A value of R-squared closer to 1 could mean that the regression model covers most part of the variance of the values of the response variable & can be termed as a good model.

(4) Example of Linear Regression :-
Consider following data for 5 student.

Student	Score in X standard (X)	Score in XII standard (Y)
1	95	85
2	85	95
3	80	70
4	70	65
5	60	70



x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
95	85	19	8	289	1936
85	95	9	18	81	126
80	70	2	-17	4	-14
70	65	-8	-12	64	96
60	70	-18	-17	324	126
$\bar{x} = 78$		$\bar{y} = 91$		$\sum (x - \bar{x})^2 = 132$	$\sum (x - \bar{x})(y - \bar{y}) = 470$

⑨ Linear regression eqⁿ that best predicts standard ZTth score.

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\beta_1 = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\beta_1 = 470 / 132 = 0.644$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\beta_0 = 91 - (0.644 * 78) = 26.768$$

$$\hat{y} = 26.76 + 0.644 x$$

⑩ Interpretation of regression line
Interpretation I:

For an increase in value of x by 0.644 units there is an increase in value of y in one unit.
Interpretation 2

Even if $x=0$ value of independent variable y is 26.768.
It is expected that value of y is 26.768.



If a student's score is 65 in std X, then its expected score in XII standard is 78.288

For $x = 80$ the y value will be

$$\hat{y} = 26.76 + 0.644 \cdot 65 = 68.38$$

(5) Training data set & Testing data set

- M/c learning algo. has two phases -
- 1. Training
- 2. Testing

Training of
learning algo.

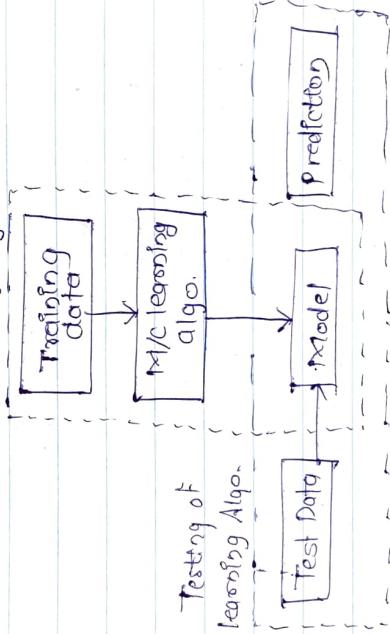


Fig. Training & testing phase in m/c learning.



⑤ Training Phase ↗

- Training dataset is provided as input to this phase.
- Training dataset is a dataset having attributes of class labels & used for training mlc learning algo. to prepare models.

⑥ Testing phase ↗

- Testing dataset is provided as input to this phase.
- Test dataset is a dataset for which class label is unknown. It is tested using model.
- A test dataset used for assessment of the finally chosen model.

⑦ Generalization ↗

- Generalization is the prediction of the future based on the past system.
- It needs to generalize beyond the training data to some future data that it might not have seen yet.
- The ultimate aim of the mlc learning model is to minimize the generalization error.
- The generalization error is essentially the average error for data the model has never seen.



④ Training Phase ↗

- Training dataset is provided as i/p to this phase
- Training dataset is a dataset having attributes & class labels & used for training machine learning algo. to prepare models.

⑤ Testing phase ↗

- Testing dataset is provided as input to this phase.
- Test dataset is a dataset for which class label is unknown. It is tested using model.
- A test dataset used for assessment of the finally chosen model.

⑥ Generalization ↗

- Generalization is the prediction of the future based on the past system.
- It needs to generalize beyond the training data to some future data that it might not have seen yet.
- The ultimate aim of the machine learning model is to minimize the generalization error.
 - The generalization error is essentially the average error for data the model has never seen).



Algorithm (Synthesis Dataset):

- Step-01: Import libraries & create alias for pandas, Numpy & matplotlib.
- Step-02: Create a Dataframe with Dependent Variable (y) & independent variable x .
- Step-03: Create Linear Regression Model using polyfit fun.
- Step-04: Observe the coeff. of the model.
- Step-05: Predict the y value for x & observe the o/p.
- Step-06: Predict the y pred for all values of x .
- Step-07: Evaluate the performance of model (R-Square)
- Step-08: Plotting the linear regression model

Algorithm (Boston Dataset):

- Step-01: Import libraries & create alias for pandas, numpy & matplotlib.



Algorithm (Synthesis Dataset):

- Step-01: Import libraries & create alias for Pandas, Numpy & matplotlib.
- Step-02: Create a Dataframe with Dependent Variable (x) & independent variable y.
- Step-03: Create Linear Regression Model using polyfit fun.

Step-04: Observe the coeff. of the model.

Step-05: Predict the Y value for x & observe the o/p.

Step-06: Predict the Y-pred for all values of x

Step-07: Evaluate the performance of model (R-Square)

Step-08: Plotting the linear regression model

Algorithm (Boston Dataset):

- Step-01: Import libraries & create alias for pandas, numpy & matplotlib.



Algorithm (Synthesizes Dataset):

Step-01: Import libraries & create alias for pandas, Numpy & matplotlib.

Step-02: Create a DataFrame with Dependent Variable (y) & independent variable x.

Step-03: Create Linear Regression Model using 'polyfit' fn.

Step-04: Observe the coeff. of the model.

Step-05: Predict the Y value for x & observe the o/p.

Step-06: Predict the Y-pred for all values of x

Step-07: Evaluate the performance of model (R-Square)

Step-08: Plotting the linear regression model

Algorithm (Boston Dataset):

Step-01:

Import libraries & create alias for pandas, Numpy & matplotlib.



Algorithm (Synthesise Dataset):

Step-01: Import libraries & create alias for pandas, NumPy & matplotlib.

Step-02: Create a DataFrame with Dependent Variable (y) & independent variable x.

Step-03: Create Linear Regression Model using polyfit fun

Step-04: Observe the coeff. of the model.

Step-05: Predict the Y value for x & observe the op.

Step-06: Predict the Y-pred for all values of x

Step-07: Evaluate the performance of model (R-Square)

Step-08: Plotting the linear regression model

Algorithm (Boston Dataset)

Step-01:

Import libraries & create alias for pandas, NumPy & matplotlib.



Step-02 :

Import the Boston Housing dataset.

Step-03 :

Initialize the datframe

Step-04 :

Add the feature names to datframe

Step-05 :

Adding target variable to datframe.

Step-06 :

Perform Data Preprocessing (check for missing values)

Step-07 :

Split dependent variable & independent variable.

Step-08 :

Splitting data to training & testing dataset.

Step-09 :

Use linear regression (train the model) to create model.

Step-10 :

Predict the y-pred for all values of train_x & train_y.



Step-11:

Evaluate the performance of Model for training & testing

Step-12:

Calculate Mean Square error for training & testing

→

Step-13:

Plotting the Linear regression model.

Conclusion:

In this way we have done data analysis using linear regression for Boston Dataset to predict the price of houses using the feature of the Boston Dataset.

Experiment No. 05



Title of the Assignment \Rightarrow

1. Implement logistic regression using python/R to perform classification on Social_Network_Ads.csv dataset.

2. Compute Confusion matrix to find T_p , F_p , T_N , F_N , Accuracy, Error rate, precision, recall on the given dataset.

Objective of the Assignment \Rightarrow

Students should be able to do analysis using logistic regression using python for any specific source dataset.

Prerequisite \Rightarrow

1. Basic of python programming
2. Concept of regression.

Theory \Rightarrow

(1) Logistic Regression \Rightarrow

- There are lots of classification problems that are available, but logistic regression is common & is a useful regression method for solving the binary classification problem.

- Logistic regression can be used for various classification problem such as spam detection.
- Logistic Regression is one of the most simple & commonly used machine learning algorithms for two-class classification.
- Logistic regression is a statistical method for predicting binary classes.
- In a Linear Regression Eq:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

where, $y \Rightarrow$ dependent variable
 $x_1, x_2, \dots, x_n \Rightarrow$ explanatory variables.

Sigmoid Eq \Rightarrow

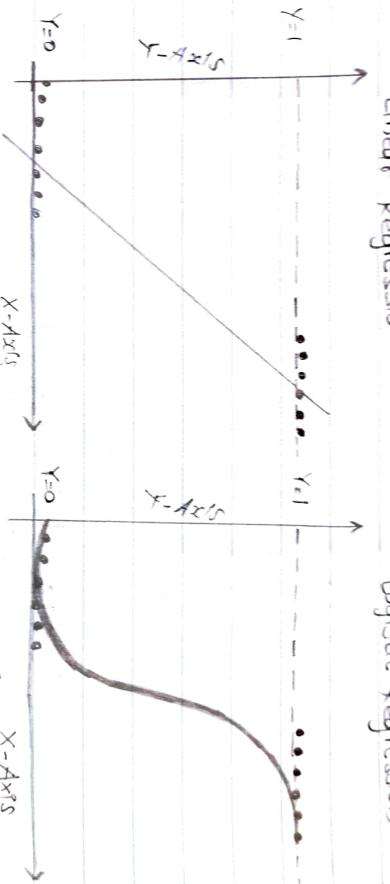
$$p = \frac{1}{1+e^{-y}}$$

Apply Sigmoid fun on Linear regression:

$$p = \frac{1}{1+e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$



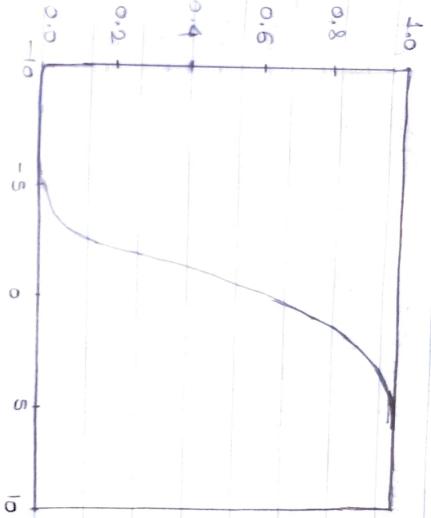
- (2) Differentiate b/w Linear & Logistic regression.
- Linear regression gives you a continuous op.,
but logistic regression provides a constant op.
 - Linear regression is estimated using Ordinary Least Squares (OLS) while logistic regression is estimated using Maximum Likelihood Estimation (MLE) approach.



(3) Sigmoid fun \Rightarrow

- The sigmoid fun, also called logistic fun, gives an 'S' shaped curve that can take any real-valued number & map a value b/w 0 & 1.
- If curve goes to +ve infinity, y predicted will become 1 & if the curve goes to -ve infinity, y predicted will become 0.

$$f(x) = \frac{1}{1 + e^{-x}}$$



(4) Types of Logistic Regression \Rightarrow

- Binary Logistic Regression
- Multinomial Logistic Regression
- Ordinal Logistic Regression

(5) Confusion Matrix Evaluation Metrics

- Contingency table or confusion matrix is often used to measure the performance of classifiers.



The following table shows the confusion matrix for two class classifier.

predicted

		actual		P	
		TP	FN		
N	FP	TN	N		

Confusion Matrix

Some important measures derived from confusion matrix are :

- Number of positive (Pos) :

Total number instances which are labelled as +ve in a given dataset.

- Number of negative (Neg) :

Total Number instances which are labelled as -ve in a given dataset.

- No. of True +ve (TP) :

No. of instances which are actually labelled as positive and the predicted class by classifier is also +ve.

- No. of True -ve (TN) :

No. of instances which are actually labelled as -ve & the predicted class by classifier is also -ve.



- No. of False +ve (FP) \Rightarrow
No. of instances which are actually labelled
as -ve & predicted class by classifier is +ve.
- No. of False -ve (FN) \Rightarrow
No. of instances which are actually labelled
as +ve & predicted class by classifier is -ve.

- Accuracy \Rightarrow

$$\text{acc} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} = \frac{\text{TP} + \text{TN}}{\text{Pos} + \text{Neg}}$$

- Error Rate:

$$\text{err} = \frac{\text{FP} + \text{FN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} = \frac{\text{FP} + \text{FN}}{\text{Pos} + \text{Neg}} \quad \text{or}$$

$$\text{err} = 1 - \text{acc}$$

- Precision:

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

- Recall:

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



Algorithm (Boston Dataset):

- Step-01: Import libraries & create alias for pandas, Numpy & Matplotlib.
- Step-02: Import the Social-Media-Adv Dataset.
- Step-03: Initialize the DataFrame
- Step-04: Perform Data preprocessing.
 - Convert Categorical to Numerical Values if applicable
 - Check for Null values.
 - Covariance Matrix to select the most promising features.
- Divide the dataset into independent (X) Dependent (y) variables.
- Split the dataset into training & testing datasets
- Scale the features if necessary.
- Step-05: Use Logistic regression (Train the ML) to create model
- Step-06: Predict the Ypred for all values of train-X & test-Y.
- Step-07: Evaluate the performance of Model for train-Y & test-Y.
- Step-08: Calculate the required evaluation parameters.



Conclusion ↗

In this way we have done data analysis using logistic regression for Social Media Adv. & evaluate the performance of model.



Experiment No. 06

Title of the Assignment :-

1. Implement Simple Naive Bayes classification algo. using python / R on Iris.csv dataset.
2. Compute confusion matrix to find TP, FP, TN, FN, Accuracy, Error rate, Precision, Recall on the given dataset.

Objective :-

Students should be able to do data analysis using Naive Bayes classification algo. using python for any open source dataset.

Prerequisite :-

1. Basic of python programming.
2. Concept of Joint & Marginal Probability.

Theory :-

1) Concept used is Naive Bayes classifier.

- Naive Bayes classifier can be used for classification of categorical data.

- Let there a 'j' number of classes, $C = \{1, 2, \dots, j\}$
- Let, p/p observation is specified by 'p' feature.

Therefore p/p observation x is given, $x = \{x_1, x_2, \dots, x_p\}$

- Naive Bayes classifier depends on Bayes rules from probability theory.



- Prior probabilities: Probabilities which are calculated for some event based on no other information are called prior probabilities.
- For e.g. $P(A)$, $P(B)$, $P(C)$ are prior probabilities because while calculating $P(A)$, occurrence of event B or C are not concerned i.e. no information about occurrence of any other event is used.

Conditional probabilities:

$$P\left(\frac{A}{B}\right) = \frac{P(A \cap B)}{P(B)} \cdot \text{ if } P(B) \neq 0 \rightarrow \textcircled{1}$$

$$P\left(\frac{B}{A}\right) = \frac{P(B \cap A)}{P(A)} \leftarrow \textcircled{2}$$

From eqn \textcircled{1} & \textcircled{2}

$$P(A \cap B) = P\left(\frac{A}{B}\right) \cdot P(B) = P\left(\frac{B}{A}\right) \cdot P(A)$$

$$\therefore P\left(\frac{A}{B}\right) = \frac{P\left(\frac{B}{A}\right) \cdot P(A)}{P(B)}$$

is called Bayes Rule.



Conditional Probability \Rightarrow

$$P(C_k | X) = \frac{P(X | C_k) * P(C_k)}{P(X)}$$

Now we have two classes & 4 features, so if we write this formula for class C_1 , it will be,

$$P(C_1 | x_1 \cap x_2 \cap x_3 \cap x_4) = \frac{P(x_1 \cap x_2 \cap x_3 \cap x_4 | C_1) * P(C_1)}{P(x_1 \cap x_2 \cap x_3 \cap x_4)}$$

Here we replaced C_k with C_1 & x with the intersection x_1, x_2, x_3, x_4 .

$$P(C_1 | x_1 \cap x_2 \cap x_3 \cap x_4) = \frac{P(x_1 | C_1) * P(x_2 | C_1) * P(x_3 | C_1) * P(x_4 | C_1) * P(C_1)}{P(x_1) * P(x_2) * P(x_3) * P(x_4)}$$

Algorithm (Iris Dataset):

Step-01: Import libraries & create alias for Pandas, Numpy & matplotlib.

Step-02: Import the Iris dataset by calling URL

Step-03: Initialize the DataFrame

Step-04: Perform Data preprocessing

Step-05: Use Naïve Bayes algorithm (Train the m/c) to create Model



Conditional Probability :-

$$P(c_k|x) = \frac{P(x|c_k)*P(c_k)}{P(x)}$$

Now we have two classes of 4 features, so if we write this formula for class c_1 , it will be,

$$P(c_1|x_1 \cap x_2 \cap x_3 \cap x_4) = \frac{P(x_1 \cap x_2 \cap x_3 \cap x_4 | c_1)^* P(c_1)}{P(x_1 \cap x_2 \cap x_3 \cap x_4)}$$

Here we replaced c_k with c_1 & x with the intersection x_1, x_2, x_3, x_4 .

$$P(c_1|x_1 \cap x_2 \cap x_3 \cap x_4) = \frac{P(x_1|c_1)^* P(x_2|c_1)^* P(x_3|c_1)^* P(x_4|c_1)^* P(c_1)}{P(x_1)^* P(x_2)^* P(x_3)^* P(x_4)}$$

Algorithm (Iris Dataset):

Step-01: Import libraries & create alias for Pandas, Numpy & Matplotlib.

Step-02: Import the Iris dataset by calling URL

Step-03: Initialize the DataFrame

Step-04: Perform Data preprocessing

Step-05: Use Naive Bayes algorithm (Train the m/c) to create Model



Step-06: Predict the y-pred for all values of train.x & test.x.

Step-07: Evaluate the performance of model for train.y & test.y.

Step-08: Calculate the required evaluation parameter

Conclusion :-

In this way we have done data analysis using Naive Bayes Algo. for Iris dataset & evaluate the performance of model.



Experiment No. 07

Title :-

1. Extract sample document & apply following document preprocessing method:
Tokenization, POS Tagging, stop words removal, stemming & Lemmatization.
2. Create representation of document by calculating Term Frequency & Inverse Document Frequency.

Objective :-

Students should be able to perform Text Analysis using TFIDF Algo.

Prerequisite :-

1. Basic of Python Programming.
2. Basic of English language.

Theory :-

(1) Basic concepts of Text Analysis :-

- Text mining is also referred to as text analytics.
- Text mining is a process of exploring sizable textual data & finding patterns.
- Text mining processes the text itself, while NLP processes with the underlying metadata.
- Finding frequency count of words, length of the sentence, presence / absence of specific words is known as text mining.



- Text mining is preprocessed data for text analytics.
- In text analytics, statistical & ml learning algo. are used to classify information.

(2) Text Analysis Operations using natural language toolkit :

- NLTK is a leading platform, for building python programs to work with human language data.
- It provides easy-to-use interfaces of lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing & semantic reasoning & many more.

(1) Tokenization :

Tokenization is the first step in text analysis. The process of breaking down a text paragraph into smaller chunks such as words or sentences is called tokenization.

- Sentence tokenization : split a paragraph into list of sentences using `sent_tokenize()` method.
- Word tokenization : split a sentence into list of words using `word_tokenize()` method.



② Stop words removal :-

Stopwords considered as noise in the text.

Text may contain stop words such as it is, am, are, this, a, an, the etc. In NLTK for removing stopwords, you need to create a list of stopwords & filter out your list of tokens from these words.

③ Stemming & Lemmatization :-

Stemming is a normalization technique where lists of tokenized words are converted into shortened root words to remove redundancy.

Lemmatization in NLTK is the algo. process of finding the lemma of a word depending on its meaning & context.

④ POS Tagging :-

POS (Parts of Speech) tell us about grammatical information of words of the sentence by assigning specific token (Determiner, noun, adjective, adverb, verb, personal Pronoun etc) as tag (DT, NN, JJ, RB, VB, PRP, etc) to each words.

② Stop words removal :-

Stopwords considered as noise in text.

Text may contain stop words like,

am, are, this, a, an, the.

In NLTK for removing stopwords, you need to create a list of stopwords & filter out your list of tokens from these words.

③ Stemming & Lemmatization :-

Stemming is a normalization technique where lists of tokenized words are converted into shortened root words to remove redundancy.

Lemmatization in NLTK is the algo. process of finding the lemma of a word depending on its meaning & context.

④ PoS Tagging :-

PoS (Parts of Speech) tell us about grammatical information of words of the sentence by assigning specific token (Determiner, noun, adjective, adverb, verb, Personal Pronoun etc) as tag (DT, NN, JJ, RB, VB, PRP, etc) to each words.



(8) Text Analysis Model using TF-IDF

Term Frequency - Inverse document

Frequency (TFIDF), is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.

- Term Frequency (TF)

It is a measure of the frequency of a word (w) in a document (d).

$$TF(w, d) = \frac{\text{occurrences of } w \text{ in document } d}{\text{total no. of words in document } d}$$

- Inverse Document Frequency (IDF)

It is the measure of the importance of a word.

$$IDF(w, D) = \log \left(\frac{\text{Total No. of documents (CN) in corpus } D}{\text{No. of documents containing } w} \right)$$

- Term Frequency - Inverse Document Frequency (TFIDF)

It is the product of TF & IDF

$$TFIDF(w, d, D) = TF(w, d) * IDF(w, D)$$



(4) Bag of Words (BoW)

- Machine learning algo. cannot work with raw text directly. Rather, the text must be converted into vectors of no.
- In natural language processing, a common technique for extracting features from text is to place all of the words that occur in the text in a bucket. This approach is called a bag of words model or BoW.

Algorithm for Tokenization, PoS Tagging, stop words removal, stemming & lemmatization:

Step-01: Download the required packages.

Step-02: Initialize the text.

Step-03: Perform Tokenization

Step-04: Perform Tokenization Removing Punctuations & stop word.

Step-05: Perform stemming

Step-06: Perform Lemmatization

Step-07: Apply PoS Tagging to text.

Algorithm for Create representation of document by calculating TFIDF.

Step-01: Import necessary libraries.

Step-02: Initialize documents

Step-03: Create Bag of words (BoW) for document A & B.



Step-04: Create collection of unique words from Document A & B.

Step-05: Create a dictionary of words & their occurrence for each document in the corpus.

Step-06: Compute term frequency for each of our documents

Step-07: Compute the term Inverse Document Frequency.

Step-08: Compute the term TF-IDF for all words.

Conclusion →

- In this way we have done text data analysis using TF-IDF algo.



Step-4: Create collection of unique words.
from Document A & B.

Step-5: Create a dictionary of counts of their
occurrence for each document in
the corpus.

Step-6: Compute term frequency for each
of our documents

Step-7: Compute the term Inverse Document
Frequency.

Step-8: Compute the term TF-IDF for all words.

Conclusion →

In this way we have done text data
analysis using TF-IDF algo.



Experiment No. 08

Title :-

Data Visualization - I

1. Use the prebuilt dataset 'titanic'. The dataset contains 891 rows & contains information about the passengers who boarded the unfortunate Titanic ship. Use the Seaborn library to see if we can find any patterns in the data.
2. Write a code to check how the price of the ticket (column name : 'Fare') for each passengers distributed by plotting a histogram.

Objective :-

Students should be able to perform the data visualization operation using python on any open source dataset.

Prerequisite :

1. Basic of python programming
2. Seaborn Library, Concept of Data Visualization



Experiment No. 09

Title :-

Data Visualization II

1. Use the inbuilt dataset 'titanic' as used in the above problem. Plot a box plot for distribution of age with respect to each gender along with the information about whether whether they survived or not. (columns name : 'sex' & 'age')
2. Write observations on the inference from the above statistics.

Objective :-

Students should be able to perform the data visualization operations using python on any open source dataset.

Prerequisite :-

1. Basics of python programming
2. Seaborn Library, Concept of Data Visualization.



Experiment No. 10

Title ↗

Data Visualization III

Download the Iris flower dataset or any dataset

from a Dataframe (e.g. <http://archive.ics.uci.edu/ml/datasets/Iris>). Scan the dataset & give the inference as:

1. List down the features & their types (e.g. numerical, nominal) available in the dataset.
2. Create a histogram for each feature in the dataset to illustrate the feature distributions.
3. Create a box plot for each feature in the dataset.
4. Compare distributions & identify outliers

Objective ↗

Students should be able to perform the data visualization operation using python on any open source dataset

Prerequisite ↗

1. Basic of python programming
2. seaborn Library, Concept of Data Visualization
3. Types of variables.