

# PRACTICAL RECORD

Name : Sahil. D. Marbate.

College Name : GCOEN.

Year : 3yr / 6th sem

Subject : Professional Skills -

## PRACTICAL-4

Aim: Correlation and Covariance

- Find the correlation matrix.
- Plot the correlation plot on dataset and visualize giving an overview of relationships among data on iris data.
- Analysis of covariance: variance (ANOVA), if data have categorical variables on iris data.

Theory:

~~Correlation:~~ Correlation measures the strength and direction of the linear relationship between two variables. It ranges from -1 to +1, where:

- +1 indicates a perfect positive linear relationship.
- 1 indicates a perfect negative linear relationship.
- 0 indicates no linear relationship.

~~Covariance:~~ Covariance measures the degree to which two variables change together. If the covariance is positive, it indicates that as one variable increases, the other variable tends to increase. If it is negative, it indicates that as one variable increases, the other variable tends to decrease.

~~Pairplot:~~ A pairplot is a graphical representation of relationships between pairs of variables in a dataset.

~~ANOVA~~ is a statistical method used to compare the means of three or more samples to determine if they are significantly

different from each other.

Some of the methods are: one-way ANOVA, F-statistic and P-value.

Here we are performing ANOVA to test whether there is a significant difference in the 'sepal length' variable among the different species of iris flower ('setosa', 'versicolor', 'virginica').

- Inputs: filename: iris.csv  
having the attributes: sepal-length, sepal-width, petal-length, petal-width, species.
- Libraries used: pandas, numpy, seaborn, matplotlib, scipy, sklearn.
- Conclusion: The program to implement the correlation and covariance, plotting correlation plot on dataset, analysis of covariance (ANOVA) on iris dataset is performed successfully.

Amarkrishna  
8/2/2024 (A)

Aim: REGRESSION MODEL

Import a data from web storage. Name the dataset and now do logistic regression to find out relation between variables that are affecting the admission of a student in an institute based on the his or her GRE score, GPA obtained and rank of the student. Also check the model is fit or not. Dependent (Foreign), independent (MASS).

Theory:

In data analysis and machine learning projects, acquiring data is often the first step. While local datasets are common, sometimes the required data resides on the web.

Logistic regression is a statistical technique used for modelling the probability of a binary outcome based on one or more predictor variables. In this practical, we aim to analyze the factors affecting student admission in an institute using logistic regression. We will import a dataset from web storage, perform logistic regression analysis, and evaluate the model's fit.

It is a regression model used when the dependent variable is binary. It estimates the probability that a given instance belongs to a particular category. In our case, the binary outcome variable is admission.

(1 for admitted, 0 for not admitted)

Steps include importing the Dataset, Data preprocessing, Model fitting, Model evaluation, Interpretation of results.

• Inputs : filename :

Attributes are: admit, gre, gpa, rank

• Libraries used:

(a). pandas

(b). numpy

(c). statsmodels

(d). Scipy

• Conclusion:

Logistic regression is a powerful tool for analyzing the factors influencing binary outcomes, such as student admission. By performing logistic regression analysis, we can identify significant predictors and develop insights into a admission process. This practical demonstrates the institutions using or as an application of logistic regression successfully.

~~Abhayan~~  
15/2/2024 (A)



### Aim: MULTIPLE REGRESSION MODEL

Apply multiple regressions, if data have a continuous Independent Variable. Apply on dataset.

### Theory:

Multiple Linear Regression is a statistical method used to analyze the relationship between two or more independent variables and a single dependent variable. It extends the simple linear regression model, which only considers one independent variable.

- Key Concepts:

- (1) Dependent Variable
- (2) Independent Variables ( $x_1, x_2, \dots, x_n$ )
- (3) Coefficients ( $B_0, B_1, \dots, B_n$ )
- (4) Intercept ( $B_0$ )

- Multiple linear regression relies on several assumptions, including:  
Linearity, Independence, Homoscedasticity, Normality,  
No multicollinearity.

After fitting the regression model, we typically evaluate its performance using metrics such as,

- (1) Mean squared Error (MSE)
- (2) R-squared ( $R^2$ )

- Inputs:

filename : california\_housing.csv.

• Libraries used:

(i) Pandas

(ii) Sklearn

• Algorithm:

- (i). Load the required python packages using import.
- (ii). Load the 'california-housing.csv' dataset using read\_csv() function of panda library.
- (iii). Describe the dataset using describe() and info() function.
- (iv). Print the correlation matrix using corr() function.
- (v). Create the required model and train them.
- (vi). Use the test case to predict using the model.
- (vii). Print the required results and figures.

Conclusion:

We have successfully created and executed a program which applies multiple regression on the given dataset.

✓  
Amar  
7/3/2024

18

## PRACTICAL No. 7

Aim: CLASSIFICATION MODEL

- Install relevant package for classification.
- Choose classifier for classification problem
- Evaluate the performance of classifier.

Theory:

Classification is a supervised learning task where the goal is to predict the where the goal is to predict the categorical class labels of new instances based on past observations. In other words, it involves assigning instances to one of several predefined classes or categories.

Key concepts:

- Classifier
- Features ( $x$ )
- Target variable ( $y$ )
- Training and Testing Data.

Once the classifier is trained and tested, its performance is evaluated using various metrics including:

- Accuracy.
- Confusion matrix
- Precision, Recall and F-1 score

Support Vector machine (SVM) is a powerful supervised learning algorithm used for classification tasks.

It finds the optimal hyper plane that best separates the classes in the features space.

- Inputs:

File name : iris.csv

- Libraries used:

(1). Pandas

(2). sklearn. datasets

(3). sklearn. svm.

(4). sklearn. metrics

(5). sklearn. linear\_model

(6). sklearn. model\_selection

### Conclusion:

We have successfully created and executed a program which installs relevant packages for classification, chooses classifier for classification problem and evaluates the performance of classifier.

Akhare  
11/3/2024 (X)

Aim: CLUSTERING MODEL

- Clustering algorithms for unsupervised classification.
- Plot the cluster using matplotlib visualizations.

Theory:

Clustering Algorithms for unsupervised classification: Clustering algorithms are a fundamental part of unsupervised machine learning used for grouping similar data points together. One of the most commonly used clustering algorithms is k-means. K-means is a partition-based clustering algorithm that divides the dataset into 'k' distinct, non-overlapping clusters.

## Plotting Cluster Data using matplotlib:

Matplotlib is a popular python library used for creating static, interactive and animated visualizations in Python. In the context of clustering, matplotlib is often used to visualize the clustered data points and the cluster centers.

~~Practical Implementation:~~

- Data Generation.
- k-means Clustering
- Fitting the Model
- Visualizations
- Display

Libraries Used:

- matplotlib
- sklearn.cluster
- sklearn.datasets



Inputs: dataset make\_blobs

- This dataset is imported from the 'sklearn.datasets' module.
- Isotropic gaussian blobs for clustering.

Algorithm:

- (1) Choose the number of clusters  $k$ . Randomly initialize  $k$  cluster centroids  $c_1, c_2, \dots, c_k$ .
- (2) for each data point  $x_i$ , calculate distance to each centroid  $c_j$  using distance matrix.
- (3) Assign  $x_i$  to cluster with nearest centroid  $\arg \min_j \text{distance}(x_i, c_j)$ .
- (4) For each cluster calculate new centroid  $c_j$  as the mean of all data points assigned to cluster  $j$ .
- (5) Check if the centroids have converged. i) The centroids no longer change significantly between iterations. ii). Maximum number of iterations is reached.
- (6) Once convergence is achieved, each data point is assigned to a cluster, and the centroids represent the final cluster centers.

Conclusion: The combination of k-means clustering and matplotlib visualization allows us to effectively group data points into cluster and visually inspect the clustering results. This technique is valuable for various applications including outcome segmentation, anomaly detection and pattern recognition among others.

~~21/3/2021  
Monday~~

## PRACTICAL No. 9



Aim: Write a program to implement k-Nearest Neighbour algorithm to classify the iris dataset. Print both correct and incorrect predictions.

Theory:

## (1). Iris Dataset:

The Iris dataset is a classic dataset in machine learning and statistics, containing measurements for 150 iris flowers from three different species: setosa, Versicolor and Virginica. Each sample consists of the four features: sepal length, sepal width, petal length and petal width.

## (2). k-Nearest neighbors (k-NN) algorithm:

The k-nearest neighbor algorithm is a simple yet effective supervised learning algorithm used for classification and regression tasks.

For classification, the algorithm predicts the class of a new data point by considering the majority class among its k nearest neighbors, where k is a hyperparameter set by the user.

## (3). Practical Implementation.

(a). Loading the Dataset

(b). Splitting the Data

(c). Creating the k-NN classifier

(d). Training the model

(e). Making Predictions

(f). Evaluating Predictions



Inputs: filename: iris.csv.

### Algorithm:

- (1). Load the dataset containing features ( $x$ ) and target label ( $y$ )
- (2). Split the data into training and testing sets, typically, 80% of the data.
- (3). Choose the number of neighbors ( $k$ ) for the kNN algorithm.
- (4). Initialize the kNN classifier with chosen  $k$  value, then the model using the training data ( $x\text{-train}, y\text{-train}$ ) using the 'fit' method.
- (5). Calculate the distance between data point and all training data points.
- (6). Select the  $k$  nearest neighbors based on calculated distance.
- (7). Determine the class label of data point based on majority class  $k$  nearest neighbors.
- (8). Store predicted class label for the data point.
- (9). Compare the predicted labels with the actual labels from the set ( $y\text{-test}$ ).
- (10). Calculate the accuracy of the model by dividing the number of correct predictions by total number of predictions.
- (11). Print the accuracy, both correct and incorrect predictions, along with the corresponding actual and predicted labels.



### Conclusion:

The practical demonstrate the implementation of the k-NN algorithm for classification task using the Iris dataset. Successfully, the k-NN algorithm performs well on the Iris dataset achieving high accuracy in predicting the species of the Iris flower based on the their features measurements.

(P)

Prakash  
07/01/2024

PRACTICAL No. 10

Aim: Implement k-mean algorithm to classify the given data set.

Theory:K-means Algorithm:

k-Means is an iterative clustering algorithm that partitions data into k clusters. It aims to minimize the within-cluster sum of squares (WCSS) by iteratively assigning data points to the nearest cluster centroid and updating the centroids until convergence.

1. Initialization:

k initial centroids are randomly chosen from the data points or generated using a specific initialization technique like k-means++.

(2). ~~Assignment step:~~ Each data point is assigned to the nearest centroid based on a distance metric, commonly the Euclidean distance.

(3). ~~Update step:~~ After all data points have been assigned, the centroids are recalculated as the mean of all data points assigned to each cluster.

(4). ~~Convergence:~~ steps 2 and 3 are repeated until the centroids no longer change significantly or a maximum number of iterations is reached.

• Inputs: Iris.csv

• Libraries used:

- (1). sklearn (scikit-learn)
- (2). matplotlib

### Algorithm:

- (1). Importing libraries: The necessary libraries are imported, including pandas, train-split for data splitting, etc.
- (2). The iris dataset is loaded using 'pd.read\_csv' from the specified file path.
- (3). The features ( $x$ ) and target ( $y$ ) are separated from the dataset. In this case,  $x$  contains all columns except the last one (features), and  $y$  contains the last column.
- (4). The dataset is split into training and testing sets using 'train test split'.
- (5). A k-means clustering model is initialized with 3 clusters ( $k=3$ ) and fitted the training data using 'fit'.
- (6). The model predicts cluster labels for the test data using 'predict'.
- (7). The accuracy of the clustering model is calculated using 'accuracy\_score' by comparing the predicted labels with the actual labels from the test set.
- (8). The accuracy score, dataframe created to display the actual and predicted labels for each sample in the test set are displayed.



### Conclusion:

The k-means algorithm to classify the iris dataset is implemented successfully. The k-means algorithm provides a powerful tool for clustering and exploratory data analysis. We can uncover patterns and relationships among iris species based on their morphological characteristics.

\*Shakal  
10/1/2024

# GOVERNMENT COLLEGE OF ENGINEERING



Sector- 27, Mihan Rehabilitation Colony Khapri, Nagpur-441 108 (Maharashtra)

Phone No. 07103&295226(P) 295220(O) Website : [www.gcoen.ac.in](http://www.gcoen.ac.in)  
E-mail : principal.gcoenagpur@dtmh.maharashtra.gov.in

## Subject: Professional Skills-II Lab

### List of Practical's

Sr. No.	Aim of the Practical
1	Perform <b>DESCRIPTIVE STATISTICS</b> using summary, str, quartile function on mtcars & cars datasets.
2	Implement <b>READING AND WRITING DIFFERENT TYPES OF DATASETS</b> (.txt, .csv) from Web and specific disk location
3	Implement <b>VISUALIZATIONS</b> of data distributions using box, scatter plot and outliers using histogram, bar chart and pie chart plot.
4	Find <b>CORRELATION</b> matrix, plot <b>COVARIANCE</b> and Analysis of covariance: variance (ANOVA) on iris dataset.
5	Implement <b>REGRESSION MODEL</b> using Logistic Regression.
6	Implement <b>MULTIPLE REGRESSION MODEL</b> using multiple regressions.
7	Implement <b>CLASSIFICATION MODEL</b> and evaluate the performance of classifier.
8	Implement <b>CLUSTERING MODEL</b> and plot the cluster.
9	Implement <b>k-Nearest Neighbour</b> algorithm to classify the iris data set.
10	Implement <b>k-means</b> algorithm to classify the iris data set.

Roma Goel &  
Monali Thakre

Subject Co-ordinator

Dr.Latesh Malik  
Head of Department