



JAYAWANT SHIKSHAN PRASARAK MANDAL's
Bhivarabai Sawant Institute of Technology & Research

(Approved by AICTE New Delhi, DTE Mumbai & Affiliated to Savitribai Phule Pune University)

Accredited with B++ Grade by NAAC

Gat No. 719/1 & 2, Wagholi, Pune-Nagar Road, Pune-412207

Ph : 020-067335108, 65217050, 67335100

Telefax : 020-67335100

Website : www.jspm.edu.in / www.bsiotr.org

EN 6311 / CEGP-013100



Prof. Dr. T. J. Sawant
B.E. (Elec.) PGDM, Ph.D
Founder Secretary

Dr. T.K. Nagaraj
ME. (Civil Engg), Ph.D (Civil Engg)
LMISTE, LMIGS, LMIRC
LMISRTT, LMIE
Principal

Department of Information Technology

अहमद एडुकेटल हिल्स जेकराबल गंजिचिप

Lab Maunal

SUBJECT: Labrotary Prattice I
[314448]

TE IT (2019 Course)

Semester v
(A.Y: 2024-25)

Subject In-charge: Prof. Ashvini Pawale

JSPM's
Bhivarabai Sawant Institute of Technology and Research, Wagholi
Department of Information Technology

Vision:

“To be a nucleus nurturing learner to cater current & future digital needs”

Mission:

1. To groom learners for addressing technical challenges by utilizing knowledge and skill sets.
2. To inculcate professional values to develop effective and efficient organization through best practices.

PEOs:

PEO1: Graduate shall have the ability to exhibit excellence in professional career by demonstrating a positive representation of their brand.

PEO2: Graduate shall have the ability to learn latest trends coping present and future needs.

PEO3: Graduate shall have sense of social responsibility by balancing the emotional quotient and strengthening the personal traits.

PSOs:

PSO1: Apply appropriate technologies and employ suitable methodologies by managing the information technology resources of an individual or organization for betterment.

PSO2: Anticipate the ever changing trends in information technology and assess the likely utility of new technologies.

PSO3: Develop IT systems that would resolve issues related to socio-economic development and build the nation through digitization.

Syllabus Of Lab Manual

Savitribai Phule Pune University, Pune		
Third Year Information Technology (2019 Course)		
314448 : Laboratory Practice-I (Machine Learning)		
Teaching Scheme:	Credit Scheme:	Examination Scheme:
Practical (PR) : 4 hrs/week	02 Credits	PR : 25 Marks TW: 25 Marks
Prerequisites: 1. Python programming language		
Course Objectives: 1. The objective of this course is to provide students with the fundamental elements of machine learning for classification, regression, clustering. 2. Design and evaluate the performance of a different machine learning models.		
Course Outcomes: On completion of the course, students will be able to— CO1: Implement different supervised and unsupervised learning algorithms. CO2: Evaluate performance of machine learning algorithms for real-world applications.		
Guidelines for Instructor's Manual		
The faculty member should prepare the laboratory manual for all the experiments and it should be made available to students and laboratory instructor/Assistant.		
Guidelines for Student's Lab Journal		
1. Students should submit term work in the form of a handwritten journal based on a specified list of assignments. 2. Practical Examination will be based on the term work. 3. Students are expected to know the theory involved in the experiment. 4. The practical examination should be conducted if and only if the journal of the candidate is complete in all respects.		
Guidelines for Lab /TW Assessment		
1. Examiners will assess the term work based on performance of students considering the parameters such as timely conduction of practical assignment, methodology adopted for implementation of practical assignment, timely submission of assignment in the form of handwritten write-up along with results of implemented assignment, attendance etc. 2. Examiners will judge the understanding of the practical performed in the examination by asking some questions related to theory & implementation of experiments he/she has carried out. 3. Appropriate knowledge of usage of software and hardware related to respective laboratories should be as a conscious effort and little contribution towards Green IT and environment awareness, attaching printed papers of the program in a journal may be avoided. There must be hand-written write-ups for every assignment in the journal. The DVD/CD containing student programs should be attached to the journal by every student and the same to be maintained by the department/lab In-charge is highly encouraged. For reference one or two journals may be maintained with program prints at Laboratory.		

Guidelines for Laboratory Conduction
<ol style="list-style-type: none">1. All the assignments should be implemented using python programming language2. Implement any 4 assignments out of 63. Assignment clustering with K-Means is compulsory4. The instructor is expected to frame the assignments by understanding the prerequisites, technological aspects, utility and recent trends related to the topic.5. The instructor may frame multiple sets of assignments and distribute them among batches of students.6. All the assignments should be conducted on multicore hardware and 64-bit open-sources software
Guidelines for Practical Examination
<ol style="list-style-type: none">1. Both internal and external examiners should jointly set problem statements for practical examination. During practical assessment, the expert evaluator should give the maximum weightage to the satisfactory implementation of the problem statement.2. The supplementary and relevant questions may be asked at the time of evaluation to judge the student 's understanding of the fundamentals, effective and efficient implementation.3. The evaluation should be done by both external and internal examiners.
List of Laboratory Assignments
Group A

1. Data preparation:

Download heart dataset from following link.

<https://www.kaggle.com/zhaoyingzhu/heartcsv> Perform following operation on given dataset.

- a) Find Shape of Data
- b) Find Missing Values
- c) Find data type of each column
- d) Finding out Zero's
- e) Find Mean age of patients
- f) Now extract only Age, Sex, ChestPain, RestBP, Chol. Randomly divide dataset in training (75%) and testing (25%).

Through the diagnosis test I predicted 100 report as COVID positive, but only 45 of those were actually positive. Total 50 people in my sample were actually COVID positive. I have total 500 samples.

Create confusion matrix based on above data and find

- I. Accuracy
- II. Precision
- III. Recall
- IV. F-1 score

2. Assignment on Regression technique

Download temperature data from below link. <https://www.kaggle.com/venky73/temperatures-of-india?select=temperatures.csv>

This data consists of temperatures of INDIA averaging the temperatures of all places month wise. Temperatures values are recorded in CELSIUS

- a. Apply Linear Regression using suitable library function and predict the Month-wise

temperature.

- b. Assess the performance of regression models using MSE, MAE and R-Square metrics
- c. Visualize simple regression model.

3. Assignment on Classification technique

Every year many students give the GRE exam to get admission in foreign Universities. The data set contains GRE Scores (out of 340), TOEFL Scores (out of 120), University Rating (out of 5), Statement of Purpose strength (out of 5), Letter of Recommendation strength (out of 5), Undergraduate GPA (out of 10), Research Experience (0=no, 1=yes), Admitted (0=no, 1=yes). Admitted is the target variable.

Data Set Available on kaggle (The last column of the dataset needs to be changed to 0 or 1) Data Set : <https://www.kaggle.com/mohansacharya/graduate-admissions>

The counselor of the firm is supposed check whether the student will get an admission or not based on his/her GRE score and Academic Score. So to help the counselor to take appropriate decisions build a machine learning model classifier using Decision tree to predict whether a student will get admission or not.

Apply Data pre-processing (Label Encoding, Data Transformation....) techniques if necessary.

Perform data-preparation (Train-Test Split)

C. Apply Machine Learning Algorithm

D. Evaluate Model.

4. Assignment on Improving Performance of Classifier Models

A SMS unsolicited mail (every now and then known as cell smartphone junk mail) is any junk message brought to a cellular phone as textual content messaging via the Short Message Service (SMS). Use probabilistic approach (Naive Bayes Classifier / Bayesian Network) to implement SMS Spam Filtering system. SMS messages are categorized as SPAM or HAM using features like length of message, word depend, unique keywords etc.

Download Data -Set from : <http://archive.ics.uci.edu/ml/datasets/sms+spam+collection>

This dataset is composed by just one text file, where each line has the correct class followed by the raw message.

- a. Apply Data pre-processing (Label Encoding, Data Transformation....) techniques if necessary
- b. Perform data-preparation (Train-Test Split)
- c. Apply at least two Machine Learning Algorithms and Evaluate Models
- d. Apply Cross-Validation and Evaluate Models and compare performance.
- e. Apply Hyper parameter tuning and evaluate models and compare performance.

5. Assignment on Clustering Techniques

Download the following customer dataset from below link:

Data Set: <https://www.kaggle.com/shwetabh123/mall-customers>

This dataset gives the data of Income and money spent by the customers visiting a Shopping Mall. The data set contains Customer ID, Gender, Age, Annual Income, Spending Score. Therefore, as a mall owner you need to find the group of people who are the profitable customers for the mall owner. Apply at least two clustering algorithms (based on Spending Score) to find the group of customers.

- a. Apply Data pre-processing (Label Encoding , Data Transformation....) techniques if necessary.
- b. Perform data-preparation(Train-Test Split)

- c. Apply Machine Learning Algorithm
- d. Evaluate Model.
- e. Apply Cross-Validation and Evaluate Model

6. Assignment on Association Rule Learning

Download Market Basket Optimization dataset from below link.

Data Set: <https://www.kaggle.com/hemanthkumar05/market-basket-optimization>

This dataset comprises the list of transactions of a retail company over the period of one week. It contains a total of 7501 transaction records where each record consists of the list of items sold in one transaction. Using this record of transactions and items in each transaction, find the association rules between items.

There is no header in the dataset and the first row contains the first transaction, so mentioned header = None here while loading dataset.

- a. Follow following steps :
- b. Data Preprocessing
- c. Generate the list of transactions from the dataset
- d. Train Apriori algorithm on the dataset
- e. Visualize the list of rules

F. Generated rules depend on the values of hyper parameters. By increasing the minimum confidence value and find the rules accordingly

7. Assignment on Multilayer Neural Network Model

Download the dataset of National Institute of Diabetes and Digestive and Kidney Diseases from below link :

Data Set: <https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-indians-diabetes.data.csv>

The dataset is has total 9 attributes where the last attribute is “Class attribute” having values 0 and 1. (1=“Positive for Diabetes”, 0=“Negative”)

- a. Load the dataset in the program. Define the ANN Model with Keras. Define at least two hidden layers. Specify the ReLU function as activation function for the hidden layer and Sigmoid for the output layer.
- b. Compile the model with necessary parameters. Set the number of epochs and batch size and fit the model.
- c. Evaluate the performance of the model for different values of epochs and batch sizes.
- d. Evaluate model performance using different activation functions Visualize the model using ANN Visualizer.

Reference Books:

1. Ethem Alpaydin, Introduction to Machine Learning, PHI 2nd Edition-2013
2. Peter Flach: Machine Learning: The Art and Science of Algorithms that Make Sense of Data, Cambridge University Press, Edition 2012.
3. Hastie, Tibshirani, Friedman: Introduction to Statistical Machine Learning with Applications in R, Springer, 2nd Edition 2012
4. Tom M. Mitchell , Machine Learning, 1997, McGraw-Hill, First EditionC. M. Bishop: Pattern Recognition and Machine Learning, Springer 1st Edition-2013.
5. Ian H Witten, Eibe Frank, Mark A Hall: Data Mining, Practical Machine Learning Tools and Techniques, Elsevier, 3rd Edition
6. Hastie, Tibshirani, Friedman: Introduction to Statistical Machine Learning with Applications in R, Springer, 2nd Edition 2012.

Savitribai Phule Pune University, Pune Third Year Information Technology (2019 Course) 314448 (B) : Laboratory Practice-I (ADBMS)		
Teaching Scheme:	Credit Scheme	Examination Scheme:
Practical (PR) :4 hrs/week	02 Credits	PR : 25 Marks TW : 25 Marks
Prerequisites: 1. Database Management System		
Course Objectives: 1. To learn and understand Database Modeling, Architectures. 2. To learn and understand Advanced Database Programming Frameworks. 3. To learn NoSQL Databases (Open source) such as MongoDB. 4. To design and develop application using NoSQL Database. 5. To design data warehouse schema for given system.		
Course Outcomes: On completion of the course, students will be able to CO1: Apply advanced Database Programming Languages. CO2: Apply the concepts of NoSQL Databases. CO3: Install and configure database systems. CO4: Populate and query a database using MongoDB commands. CO5: Design data warehouse schema of any one real-time: CASE STUDY CO6: Develop small application with NoSQL Database for back-end.		
Guidelines for Instructor's Manual		
The faculty member should prepare the laboratory manual for all the experiments and it should be made available to students and laboratory instructor/Assistant.		
Guidelines for Student's Lab Journal		
1. Student should submit term work in the form of handwritten journal based on specified list of assignments. 2. Practical Examination will be based on all the assignments in the lab manual 3. Candidate is expected to know the theory involved in the experiment. 4. The practical examination should be conducted if and only if the journal of the candidate is complete in all respects.		

Guidelines for Lab /TW Assessment

1. Examiners will assess the student based on performance of students considering the parameters such as timely conduction of practical assignment, methodology adopted for implementation of practical assignment, timely submission of assignment in the form of handwritten write-up along with results of implemented assignment, attendance etc.
2. Appropriate knowledge of usage of software and hardware related to respective laboratory should be checked by the concerned faculty member.
3. As a conscious effort and little contribution towards Green IT and environment awareness, attaching printed papers of the program in journal may be avoided. There must be hand-written write-ups for every assignment in the journal. The DVD/CD containing student's programs should be attached to the journal by every student and same to be maintained by department/lab In- charge is highly encouraged. For reference one or two journals may be maintained with program prints at Laboratory.

Guidelines for Laboratory Conduction

1. Group A assignments are compulsory and should be performed by individual student.
2. Group B case study may be performed in group of 3/4.
3. Mini project of Group C can be implemented using any suitable front-end. But back-end must be MongoDB.

Guidelines for Practical Examination

1. Practical Examination will be based on the all topics covered.
2. Examiners will judge the understanding of the practical performed in the examination by asking some questions related to theory & implementation of experiments he/she has carried out.

List of Laboratory Assignments

Group A : MongoDB

1. Create a database with suitable example using MongoDB and implement
 - Inserting and saving document (batch insert, insert validation)
 - Removing document
 - Updating document (document replacement, using modifiers, up inserts, updating multipledocuments, returning updated documents)
 - Execute at least 10 queries on any suitable MongoDB database that demonstrates following:
 - a. Find and find One (specific values)
 - b. Query criteria (Query conditionals, OR queries, \$not, Conditional semantics) Type-specific queries (Null, Regular expression, Querying arrays)
 - c. \$ where queries
 - d. Cursors (Limit, skip, sort, advanced query options)

2. Implement Map-reduce and aggregation, indexing with suitable example in MongoDB. Demonstrate the following:

- Aggregation framework
- Create and drop different types of indexes and explain () to show the advantage of the indexes.

3. **Case Study:** Design conceptual model using Star and Snowflake schema for any one database.

4. Mini Project

Pre-requisite: Build the mini project based on the requirement document and design prepared as a part of Database Management Lab in second year.

1. Form teams of around 3 to 4 people.

2. Develop the application:

Build a suitable GUI by using forms and placing the controls on it for any application. Proper data entry validations are expected.

Add the database connection with front end. Implement the basic CRUD operations.

3. Prepare and submit report to include: Title of the Project, Abstract, List the hardware and software requirements at the backend and at the front end, Source Code , Graphical User Interface, Conclusion.

Reference Books:

1. Silberschatz A., Korth H., Sudarshan S., "Database System Concepts", 6th Edition, McGraw Hill Publishers, ISBN 0-07-120413-X.
2. Kristina Chodorow, MongoDB The definitive guide, O'Reilly Publications, ISBN:978-93-5110-269-4, 2nd Edition.
3. Jiawei Han, Micheline Kamber, Jian Pei "Data Mining: concepts and techniques", 2nd Edition, Publisher: Elsevier/Morgan Kaufmann.
4. <http://nosql-database.org/>.

PART A:

Machine Learning

Assignment No -1

Title: Regression technique

Problem Statement:

This data consists of temperatures of INDIA averaging the temperatures of all place's month wise. Temperatures values are recorded in CELSIUS

- a) Apply Linear Regression using suitable library function and predict the Month-wise temperature.
- b) Assess the performance of regression models using MSE, MAE and R-Square metrics
- c) Visualize simple regression model.

Objective: This assignment will help the students to realize how the Linear Regression can be used and predictions using the same can be performed.

S/W Packages and H/W apparatus used:

Linux OS: Ubuntu/Windows, Jupyter notebook.
PC with the configuration as Pentium IV 1.7 GHz. 128M.B RAM, 40 G.B HDD,
15''Color Monitor, Keyboard, Mouse

Theory:

Definition of Linear Regression

In layman terms, we can define linear regression as **it is used for learning the linear relationship between the target and one or more forecasters**, and it is probably one of the most popular and well inferential algorithms in statistics. Linear regression endeavours to demonstrate the connection between two variables by fitting a linear equation to observed information. One variable is viewed as an explanatory variable, and the other is viewed as a dependent variable.

Normally, linear regression is divided into two types: Multiple linear regression and Simple linear regression.

1. Multiple Linear Regression

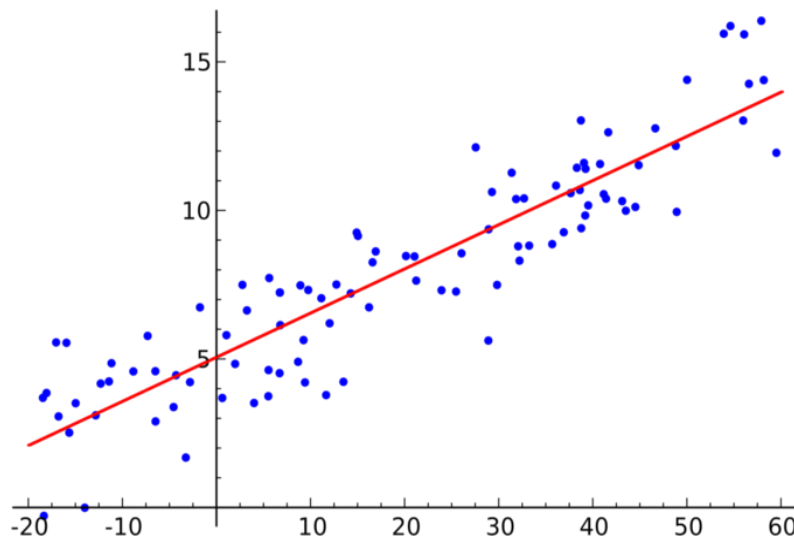
In this type of linear regression, we always attempt to discover the relationship between two or more independent variables or inputs and the corresponding dependent variable or output and the independent variables can be either continuous or categorical.

This linear regression analysis is very helpful in several ways like it helps in foreseeing trends, future values, and moreover predict the impacts of changes.

2. Simple Linear Regression

In simple linear regression, we aim to reveal the relationship between a single independent variable or you can say input, and a corresponding dependent variable or output. We can discuss this in a simple line as $y = \beta_0 + \beta_1 x + \epsilon$

Here, Y speaks to the output or dependent variable, β_0 and β_1 are two obscure constants that speak to the intercept and coefficient that is slope separately, and the error term is ϵ Epsilon. We can also discuss this in the form of a graph and here is a sample simple linear regression model graph.



Simple Linear Regression graph

What Actually is Simple Linear Regression?

It can be described as a method of statistical analysis that can be used to study the relationship between two quantitative variables.

Primarily, there are two things which can be found out by using the method of simple linear regression:

1. **Strength of the relationship between the given duo of variables.** (For example, the relationship between global warming and the melting of glaciers)
2. **How much the value of the dependent variable is at a given value of the independent variable.** (For example, the amount of melting of a glacier at a certain level of global warming or temperature)

Regression models are used for the elaborated explanation of the relationship between two given variables. There are certain types of regression models like logistic regression models, nonlinear regression models, and linear regression models. The linear regression model fits a straight line into the summarized data to establish the relationship between two variables.

Assumptions of Linear Regression

To conduct a simple linear regression, one has to make certain assumptions about the data. This is because it is a parametric test. The assumptions used while performing a simple linear regression are as follows:

- **Homogeneity of variance (homoscedasticity)**- One of the main predictions in a simple linear regression method is that the size of the error stays constant. This simply means that in the value of the independent variable, the error size never changes significantly.
- **Independence of observations**- All the relationships between the observations are transparent, which means that nothing is hidden, and only valid sampling methods are used during the collection of data.
- **Normality**- There is a normal rate of flow in the data.

These three are the assumptions of regression methods.

However, there is one additional assumption that has to be taken into consideration while specifically conducting a linear regression.

- **The line is always a straight line**- There is no curve or grouping factor during the conduction of a linear regression. There is a linear relationship between the variables (dependent variable and independent variable). If the data fails the assumptions of

homoscedasticity or normality, a nonparametric test might be used. (For example, the Spearman rank test)

Example of data that fails to meet the assumptions: One may think that cured meat consumption and the incidence of colorectal cancer in the U.S have a linear relationship. But later on, it comes to the knowledge that there is a very high range difference between the collection of data of both the variables. Since the homoscedasticity assumption is being violated here, there can be no linear regression test. However, a Spearman rank test can be performed to know about the relationship between the given variables.

Applications of Simple Linear Regression

1. **Marks scored by students based on number of hours studied (ideally)-** Here marks scored in exams are dependent and the number of hours studied is independent.
2. **Predicting crop yields based on the amount of rainfall-** Yield is a dependent variable while the measure of precipitation is an independent variable.
3. **Predicting the Salary of a person based on years of experience-** Therefore, Experience becomes the independent while Salary turns into the dependent variable.

Limitations of Simple Linear Regression

Indeed, even the best information doesn't recount a total story. Regression investigation is ordinarily utilized in examination to set up that a relationship exists between variables. However, correlation isn't equivalent to causation: a connection between two variables doesn't mean one causes the other to occur. Indeed, even a line in a simple linear regression that fits the information focuses well may not ensure a circumstances and logical results relationship.

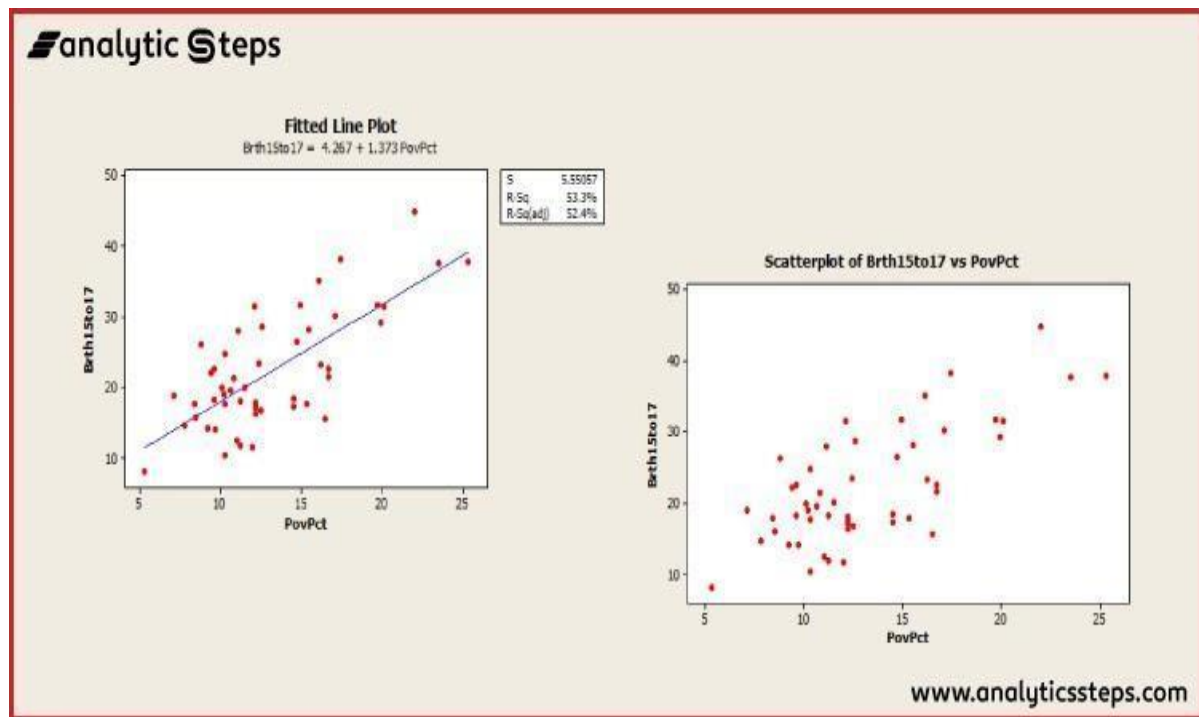
Utilizing a linear regression model will permit you to find whether a connection between variables exists by any means. To see precisely what that relationship is and whether one variable cause another, you will require extra examination and statistical analysis.

Examples of Simple Linear Regression

Now, let's move towards understanding simple linear regression with the help of an example. We will take an example of teen birth rate and poverty level data.

This dataset of size $n = 51$ is for the 50 states and the District of Columbia in the United States. The variables are y = year 2002 birth rate for each 1000 females 15 to 17 years of age and x = destitution rate, which is the percent of the state's populace living in families with wages underneath the governmentally characterized neediness level. (Information source: Mind On Statistics, 3rd version, Utts and Heckard).

Below is the graph (right image) in which you can see the (birth rate on the vertical) is indicating a normally linear relationship, on average, with a positive slope. As the poverty level builds, the birth rate for 15 to 17-year-old females will in general increment too.



Example graph of simple linear regression

Here is another graph (left graph) which is showing a regression line superimposed on the data.

The condition of the fitted regression line is given close to the highest point of the plot. The condition should express that it is for the "average" birth rate (or "anticipated" birth rate would be alright as well) as a regression condition portrays the normal estimation of y as a component of at least one x -variables. In statistical documentation, the condition could be composed $\hat{y} = 4.267 + 1.373x$.

- The interpretation of the slope (value = 1.373) is that the 15 to 17-year-old birth rate increases 1.373 units, on average, for each one unit (one per cent) increase in the poverty rate.

- The translation of the intercept (value=4.267) is that if there were states with a population rate = 0, the anticipated normal for the 15 to 17-year-old birth rate would be 4.267 for those states. Since there are no states with a poverty rate = 0 this understanding of the catch isn't basically significant for this model.
- In the chart with a regression line present, we additionally observe the data that $s = 5.55057$ and $r^2 = 53.3\%$.
- The estimation of s discloses to us generally the standard deviation of the contrasts between the y -estimations of individual perceptions and expectations of y dependent on the regression line. The estimation of r^2 can be deciphered to imply that destitution rates "clarify" 53.3% of the noticed variety in the 15 to 17-year-old normal birth rates of the states.

The R^2 (adj) value (52.4%) is a change in accordance with R^2 dependent on the number of x -variables in the model (just one here) and the example size. With just a single x -variable, the charged R^2 isn't significant.

Implementation:

```
#Import Libraries import
numpy as np
import matplotlib.pyplot as plt import pandas
as pd
```

Importing the csv file

```
from google.colab import files uploaded =
files.upload()
```

```
#Read Student Grades .csv file and divide the data into dependent and inde pendent variables.
data = pd.read_csv('Student_Grades_Data.csv') data.head()
```

```
data.shape (50,
```

```
2)
```

```
X = data.iloc[:, :-1].values y =
data.iloc[:, 1].values
```

```
X
```

```
array([[ 1],
       [ 5],
       [ 7],
       [ 3],
```

```
}SPM,BSIOTR,WAGHOLI,PUNE
```

[2],
[9],
[6],
[12],
[11],
[2],
[4],
[8],
[13],
[9],
[14],
[10],
[6],
[12],
[1],
[4],
[14],
[10],
[11],
[4],
[5],
[8],
[1],
[2],
[3],
[7],
[8],
[14],
[7],
[8],
[1],
[2],
[3],
[4],
[5],
[6],
[7],
[8],
[9],
[10],
[11],
[12],
[13],
[14],
[8],
[2]])

Y

array([1.5, 2.7, 3.1, 2.1, 1.8, 3.9, 2.9, 4.5, 4.3, 1.8, 2.4, 3.5, 4.8,
3.9, 5. , 4.1, 2.9, 4.5, 1.5, 2.4, 5. , 4.1, 4.3, 2.4, 2.7, 3.5,
1.5, 1.8, 2.1, 3.1, 3.5, 5. , 3.1, 3.5, 1.5, 1.8, 2.1, 2.4, 2.7,

}SPM,BSIOTR,WAGHOLI,PUNE

2.9, 3.1, 3.5, 3.9, 4.1, 4.3, 4.5, 4.8, 5. , 3.5, 1.8])

```
#Split the data into training and test datasets from
sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)
```

```
y_test
```

```
array([2.1, 3.5, 2.4, 3.5, 3.1, 1.8, 2.7, 5. , 4.3, 1.8, 3.5, 1.8, 1.5,
       1.5, 1.5])
```

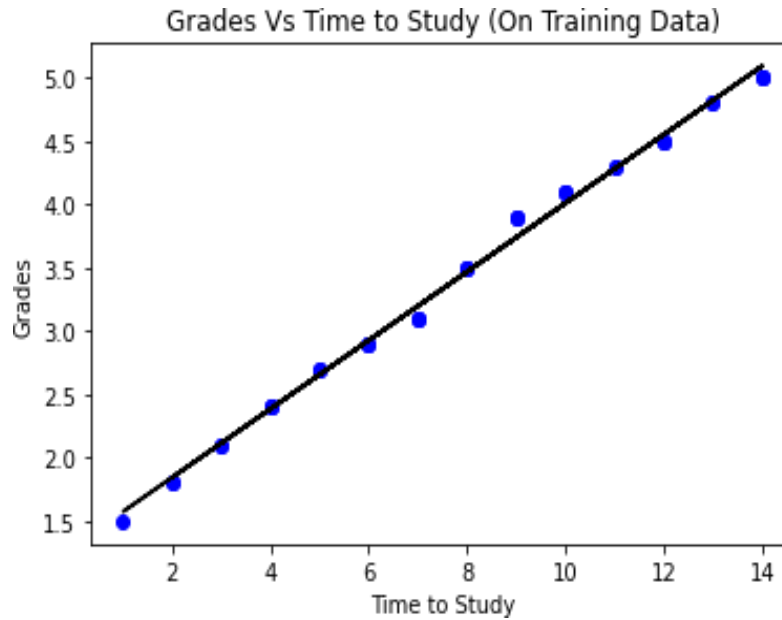
```
#Fit the Simple Linear Regression Model
from sklearn.linear_model import LinearRegression LinReg =
LinearRegression()
LinReg.fit(X_train, y_train)
```

```
#Print the
print(f'a0 = {LinReg.intercept_}') print(f'a1 =
{LinReg.coef_}')
```

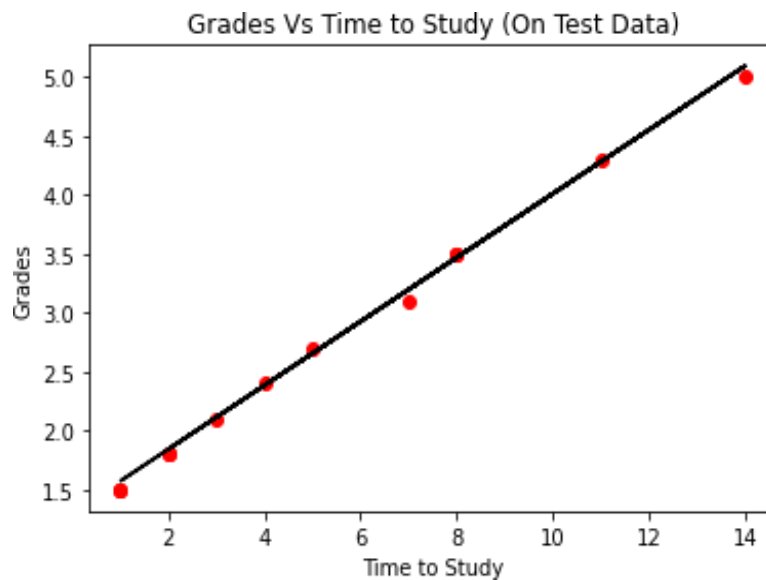
```
#Predicted grade scores from test dataset y_predict =
LinReg.predict(X_test) y_predict
```

```
#Actual grade scores from test dataset y_test
```

```
#Grades Vs Time to Study visualization on Training Data plt.scatter(X_train,
y_train, color='Blue') plt.plot(X_train, LinReg.predict(X_train), color='Black')
plt.title('Grades Vs Time to Study (On Training Data)') plt.xlabel('Time to
Study')
plt.ylabel('Grades') plt.show()
```



```
#Grades Vs Time to Study visualization on Test Data plt.scatter(X_test,
y_test, color='Red') plt.plot(X_train, LinReg.predict(X_train),
color='Black') plt.title('Grades Vs Time to Study (On Test Data)')
plt.xlabel('Time to Study')
plt.ylabel('Grades') plt.show()
```



```
#Predicting Grade of a student when he studied for 10 Hrs. Example of how to pass an external
value,
#Independent of Test or Training Dataset
```

```
Predict_Grade = LinReg.predict([[10]]) Predict_Grade
```

```
}SPM,BSIOTR,WAGHOLI,PUNE
```

```
#Model Evaluation using R-Square from
sklearn import metrics
r_square = metrics.r2_score(y_test, y_predict) print('R-
Square Error:', r_square)

#Model Evaluation using Mean Square Error (MSE)
print('Mean Squared Error:', metrics.mean_squared_error(y_test, y_predict)
)

#Model Evaluation using Root Mean Square Error (RMSE)
print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, y_predict)))

#Model Evaluation using Mean Absolute Error (MAE)
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_predict))
```

Conclusion

Simple linear regression is a regression model that figures out the relationship between one independent variable and one dependent variable using a straight line.

References:

- [1] Riya Kumari, <https://www.analyticssteps.com/blogs/simple-linear-regression-applications-limitations-examples>.
- [2] Ethem Alpaydin, Introduction to Machine Learning, PHI 2nd Edition, 2013.
- [3] Peter Flach: Machine Learning: The Art and Science of Algorithms that Make Sense of Data, Cambridge University Press, Edition 2012.

Assignment No -2

Title: Classification using Machine Learning

Problem Statement:

Perform following operations on given dataset:

- a) Apply Data pre-processing (Label Encoding, Data Transformation)
techniques if necessary.
- b) Perform data-preparation (Train-Test Split)
- c) Apply Decision tree classification Algorithm
- d) Evaluate Model.

Objective:

This assignment will help the students to realize how the decision tree classifier can be used and predictions using the same can be performed.

S/W Packages and H/W apparatus used:

Linux OS: Ubuntu/Windows , Jupyter notebook.

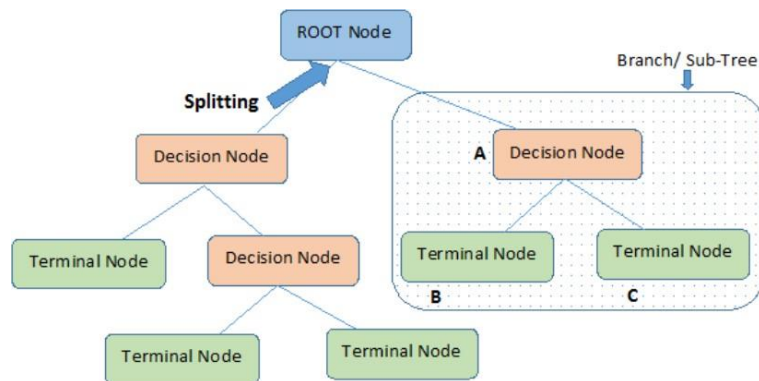
PC with the configuration as Pentium IV 1.7 GHz. 128M.B RAM, 40 G.B HDD,
15’’Color Monitor, Keyboard, Mouse

Theory:**Classification:**

Classification is a **process of categorizing a given set of data into classes**. It can be performed on both structured or unstructured data. The process starts with predicting the class of given data points. The classes are often referred to as target, label or categories.

What is a Decision Tree?

It uses a flowchart like a tree structure to show the predictions that result from a series of feature-based splits. It starts with a root node and ends with a decision made by leaves.[1]



Root Nodes – It is the node present at the beginning of a decision tree. from this node the population starts dividing according to various features.

Decision Nodes – the nodes we get after splitting the root nodes are called Decision Node

Leaf Nodes – the nodes where further splitting is not possible are called leaf nodes or terminal nodes

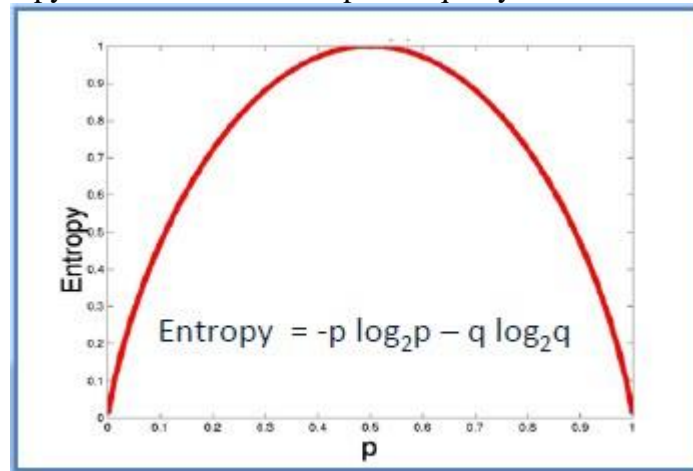
Sub-tree – just like a small portion of a graph is called sub-graph similarly a sub-section of this decision tree is called sub-tree.

Pruning – It is cutting down some nodes to stop overfitting.



Entropy:

Entropy is used to calculate the homogeneity of a sample. If the sample is completely homogeneous the entropy is zero and if the sample is equally divided it has entropy of one.



$$\text{Entropy} = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$$

a) Entropy using the frequency table of one attribute:

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Play Golf	
Yes	No
9	5



$$\begin{aligned} \text{Entropy}(\text{PlayGolf}) &= \text{Entropy}(5,9) \\ &= \text{Entropy}(0.36, 0.64) \\ &= -(0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\ &= 0.94 \end{aligned}$$

b) Entropy using the frequency table of two attributes:

$$E(T, X) = \sum_{c \in X} P(c)E(c)$$

		Play Golf		
		Yes	No	
Outlook	Sunny	3	2	5
	Overcast	4	0	4
	Rainy	2	3	5
				14



$$\begin{aligned}
 E(\text{PlayGolf}, \text{Outlook}) &= P(\text{Sunny}) * E(3,2) + P(\text{Overcast}) * E(4,0) + P(\text{Rainy}) * E(2,3) \\
 &= (5/14) * 0.971 + (4/14) * 0.0 + (5/14) * 0.971 \\
 &= 0.693
 \end{aligned}$$

Information Gain

The information gain is based on the decrease in entropy after a dataset is split on an attribute. Constructing a decision tree is all about finding attributes that return the highest information gain (i.e., the most homogeneous branches).[2]

Step 1: Calculate entropy of the target.

$$\begin{aligned}
 \text{Entropy}(\text{PlayGolf}) &= \text{Entropy}(5,9) \\
 &= \text{Entropy}(0.36, 0.64) \\
 &= - (0.36 \log_2 0.36) - (0.64 \log_2 0.64) \\
 &= 0.94
 \end{aligned}$$

Step 2: The dataset is then split on the different attributes. The entropy for each branch is calculated.

Then it is added proportionally, to get total entropy for the split. The resulting entropy is subtracted from the entropy before the split. The result is the Information Gain, or decrease in entropy.

		Play Golf	
		Yes	No
Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3
		Gain = 0.247	

		Play Golf	
		Yes	No
Temp.	Hot	2	2
	Mild	4	2
	Cool	3	1
		Gain = 0.029	

		Play Golf	
		Yes	No
Humidity	High	3	4
	Normal	6	1
		Gain = 0.152	

		Play Golf	
		Yes	No
Windy	False	6	2
	True	3	3
		Gain = 0.048	

$$Gain(T, X) = Entropy(T) - Entropy(T, X)$$

$$G(\text{PlayGolf}, \text{Outlook}) = E(\text{PlayGolf}) - E(\text{PlayGolf}, \text{Outlook})$$

$$= 0.940 - 0.693 = 0.247$$

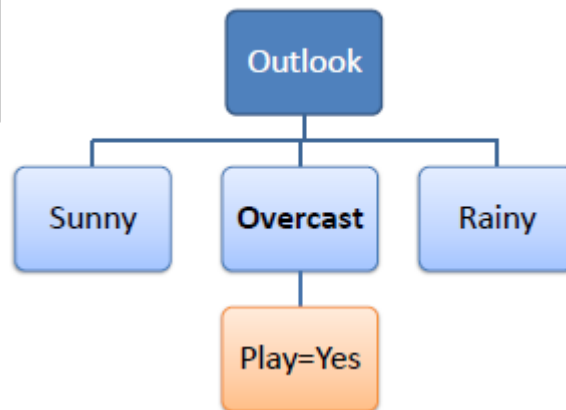
Step 3: Choose attribute with the largest information gain as the decision node, divide the dataset by its branches and repeat the same process on every branch.

		Play Golf	
		Yes	No
★ Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3
		Gain = 0.247	

		Outlook	Temp	Humidity	Windy	Play Golf
Outlook	Sunny	Sunny	Mild	High	FALSE	Yes
		Sunny	Cool	Normal	FALSE	Yes
		Sunny	Cool	Normal	TRUE	No
		Sunny	Mild	Normal	FALSE	Yes
		Sunny	Mild	High	TRUE	No
	Overcast	Overcast	Hot	High	FALSE	Yes
		Overcast	Cool	Normal	TRUE	Yes
		Overcast	Mild	High	TRUE	Yes
		Overcast	Hot	Normal	FALSE	Yes
	Rainy	Rainy	Hot	High	FALSE	No
		Rainy	Hot	High	TRUE	No
		Rainy	Mild	High	FALSE	No
		Rainy	Cool	Normal	FALSE	Yes
		Rainy	Mild	Normal	TRUE	Yes

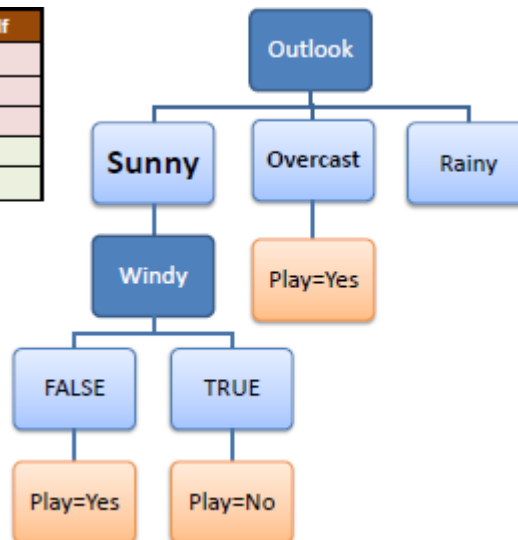
Step 4a: A branch with entropy of 0 is a leaf node.

Temp	Humidity	Windy	Play Golf
Hot	High	FALSE	Yes
Cool	Normal	TRUE	Yes
Mild	High	TRUE	Yes
Hot	Normal	FALSE	Yes



Step 4b: A branch with entropy more than 0 needs further splitting.

Temp	Humidity	Windy	Play Golf
Mild	High	FALSE	Yes
Cool	Normal	FALSE	Yes
Mild	Normal	FALSE	Yes
Cool	Normal	TRUE	No
Mild	High	TRUE	No



Step 5: The ID3 algorithm is run recursively on the non-leaf branches, until all data is classified.

Decision Tree to Decision Rules

A decision tree can easily be transformed to a set of rules by mapping from the root node to the leaf nodes one by one.

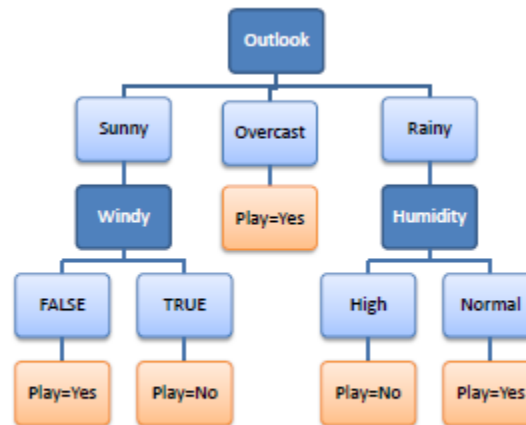
R_1 : IF (Outlook=Sunny) AND (Windy=FALSE) THEN Play=Yes

R_2 : IF (Outlook=Sunny) AND (Windy=TRUE) THEN Play=No

R_3 : IF (Outlook=Overcast) THEN Play=Yes

R_4 : IF (Outlook=Rainy) AND (Humidity=High) THEN Play=No

R_5 : IF (Outlook=Rain) AND (Humidity=Normal) THEN Play=Yes



Pruning:

It is another method that can help us avoid overfitting. It helps in improving the performance of the tree by cutting the nodes or sub-nodes which are not significant. It removes the branches which have very low importance.

There are mainly 2 ways for pruning:

- (i) **Pre-pruning** – we can stop growing the tree earlier, which means we can prune/remove/cut a node if it has low importance **while growing** the tree.
- (ii) **Post-pruning** – once our **tree is built to its depth**, we can start pruning the nodes based on their significance.

Implementation:

Importing all the necessary libraries

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt
```

Importing the csv file

```
df = pd.read_csv('../input/Admission_Predict.csv')
```

```
# Check null values in the dataset
```

```
df.isnull().sum()
```

```
Out[ ]:
```

```
Serial No.      0
GRE Score       0
TOEFL Score     0
```

```
}SPM,BSIOTR,WAGHOLI,PUNE
```

University Rating 0

```
LOR          0
CGPA         0
Research      0
Chance of Admit  0
dtype: int64
```

```
df.columns
```

```
Index(['Serial No.', 'GRE Score', 'TOEFL Score', 'University Rating', 'SOP',
      'LOR ', 'CGPA', 'Research', 'Chance of Admit '],
      dtype='object')
```

Updating chance of admission column

```
# if chance >= 80% CHANCE = 1
# if chance < 80% CHANCE = 0
```

```
dataset.loc[dataset['Chance of Admit '] < 0.8, 'Chance of Admit '] = 0
dataset.loc[dataset['Chance of Admit '] >= 0.8, 'Chance of Admit '] = 1
```

Initializing the variables

```
X = df.drop(['Chance of Admit ', 'Serial No.'], axis=1)
y = df['Chance of Admit ']
```

Split the data into training and testing set

```
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test =
train_test_split(X, y, test_size=0.25, random_state=123)
```

```
# importing required libraries
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics
```

```
# Creating Decision Tree classifier object
clf = DecisionTreeClassifier()
```

```
# Training Decision Tree Classifier
clf = clf.fit(X_train, Y_train)
```

```
#Predicting for the test data
y_pred = clf.predict(X_test)
```

Confusion matrix:

```
print("confusion matrix:\n")
```

```
print(metrics.confusion_matrix(Y_test, y_pred))
```

confusion matrix:


```
[[79 3]  
 [ 4 39]]
```

```
print("1. Accuracy Score:", metrics.accuracy_score(Y_test, y_pred))  
print("2. Precision Score:", metrics.precision_score(Y_test, y_pred))  
print("3. Recall Score:", metrics.recall_score(Y_test, y_pred))  
print("4. f1 Score:", metrics.f1_score(Y_test, y_pred))
```

1. Accuracy Score: 0.944
2. Precision Score: 0.9285714285714286
3. Recall Score: 0.9069767441860465
4. f1 Score: 0.9176470588235294

Application:

Helpful in solving classification problems.

References:

- [1] Anshul Saini ,Analytics Vidhya,Decision Tree Algorithm – A Complete Guide
- [2] Dr. Saed Sayad,https://www.saedsayad.com/decision_tree.htm

Assignment No -3

Title: K Means Clustering

Problem Statement:

Perform following operations on given dataset:

- a) Apply Data pre-processing (Label Encoding, Data Transformation.....)
techniques if necessary.
- b) Perform data-preparation (Train-Test Split)
- c) Apply Machine Learning Algorithm
- d) Evaluate Model.
- e) Apply Cross-Validation and Evaluate Model

Objective:

This assignment will help the students to realize how to do Clustering using K Means Clustering algorithm.

S/W Packages and H/W apparatus used:

Linux OS: Ubuntu/Windows , Jupyter notebook.

PC with the configuration as Pentium IV 1.7 GHz. 128M.B RAM, 40 G.B HDD,

15’’Color Monitor, Keyboard, Mouse

Theory:**Introduction to K-means Clustering**

K-means clustering is a type of unsupervised learning, which is used when you have unlabeled data (i.e., data without defined categories or groups). The goal of this algorithm is to find groups in the data, with the number of groups represented by the variable K. [1]

The algorithm works iteratively to assign each data point to one of K groups based on the features that are provided. Data points are clustered based on feature similarity.

The results of the K-means clustering algorithm are:

1. The centroids of the K clusters, which can be used to label new data
2. Labels for the training data (each data point is assigned to a single cluster) Rather than defining groups before looking at the data, clustering allows you to find and analyze the groups that have formed organically.

The "Choosing K" section below describes how the number of groups can be determined. Each centroid of a cluster is a collection of feature values which define the resulting groups. Examining the centroid feature weights can be used to qualitatively interpret what kind of group each cluster represents.

This introduction to the K-means clustering algorithm covers:

- Common business cases where K-means is used
- The steps involved in running the algorithm

Some examples of use cases are:

Behavioral segmentation:

- o Segment by purchase history
- o Segment by activities on application, website, or platform.
- o Define personas based on interests
- o Create profiles based on activity monitoring

Inventory categorization:

- o Group inventory by sales activity
- o Group inventory by manufacturing metrics

Sorting sensor measurements:

- o Detect activity types in motion sensors
- o Group images
- o Separate audio
- o Identify groups in health monitoring

Detecting bots or anomalies:

- o Separate valid activity groups from bots
 - o Group valid activity to clean up outlier detection
- In addition, monitoring if a tracked data point switches between groups over time can be used to detect meaningful changes in the data.

Algorithm:

The K-means clustering algorithm uses iterative refinement to produce a final result. The

algorithm inputs are the number of clusters K and the data set. The data set is a collection of features for each data point. The algorithms start with initial estimates for the K centroids, which can either be randomly generated or randomly selected from the data set. The algorithm then iterates between two steps:

1. **Data assignment step:**

Each centroid defines one of the clusters. In this step, each data point is assigned to its nearest centroid, based on the squared Euclidean distance. More formally, if c_i is the collection of centroids in set C , then each data point x is assigned to a cluster based on

$$\operatorname{argmin}_{c_i \in C} \operatorname{dist}(c_i, x)^2$$

where $\operatorname{dist}(\cdot)$ is the standard (L2) Euclidean distance. Let the set of data point assignments for each i th cluster centroid be S_i .

2. **Centroid update step:**

In this step, the centroids are recomputed. This is done by taking the mean of all data points assigned to that centroid's cluster.

$$c_i = \frac{1}{|S_i|} \sum_{x_i \in S_i} x_i$$

The algorithm iterates between steps one and two until a stopping criteria is met (i.e., no data points change clusters, the sum of the distances is minimized, or some maximum number of iterations is reached).

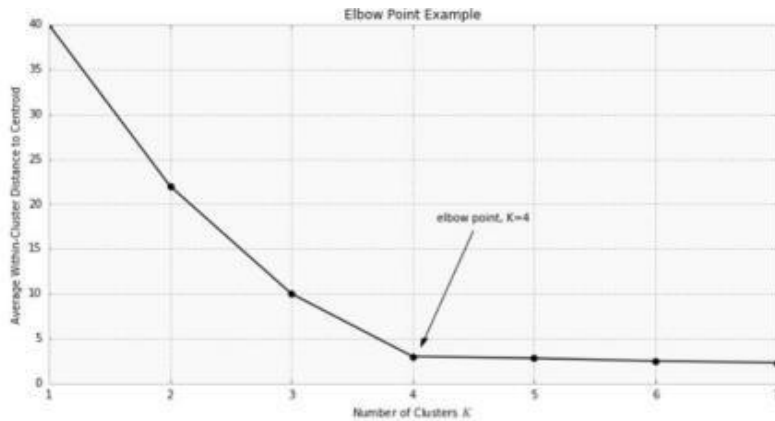
This algorithm is guaranteed to converge to a result. The result may be a local optimum (i.e. not necessarily the best possible outcome), meaning that assessing more than one run of the algorithm with randomized starting centroids may give a better outcome.

Choosing K

The algorithm described above finds the clusters and data set labels for a particular pre-chosen K . To find the number of clusters in the data, the user needs to run the K -means clustering algorithm for a range of K values and compare the results. In general, there is no method for determining exact value of K , but an accurate estimate can be obtained using the following techniques. One of the metrics that is commonly used to compare results across different values of K is the mean distance between data points and their cluster centroid. Since increasing the number of clusters will always reduce the distance to data points, increasing K will always decrease this metric, to the extreme of reaching zero when K is the same as the number of data points. Thus, this metric cannot be used as the sole target. Instead, mean distance to the centroid as a function of K is plotted and the "elbow point," where the rate of decrease sharply shifts, can be used to roughly determine K .

A number of other techniques exist for validating K , including cross-validation, information criteria, the information theoretic jump method, the silhouette method, and the G -means algorithm. In addition, monitoring the distribution of data points across groups provides

insight into how the algorithm is splitting the data for each K.



Implementation:

Importing Dataset

```
from pandas import read_csv
A=read_csv("E:/DS1/Mall_Customers.csv")
```

Dropping the irrelevant columns

```
B=A.drop(["Customer_ID"],axis=1)
```

Label encoding

```
# Import label encoder
from sklearn import preprocessing
```

```
# label_encoder object knows how to understand word labels.
label_encoder = preprocessing.LabelEncoder()
```

```
# Encode labels in column 'species'.
B['Genre']= label_encoder.fit_transform(B['Genre'])
```

```
B['Genre'].unique()
array([1, 0], dtype=int64)
```

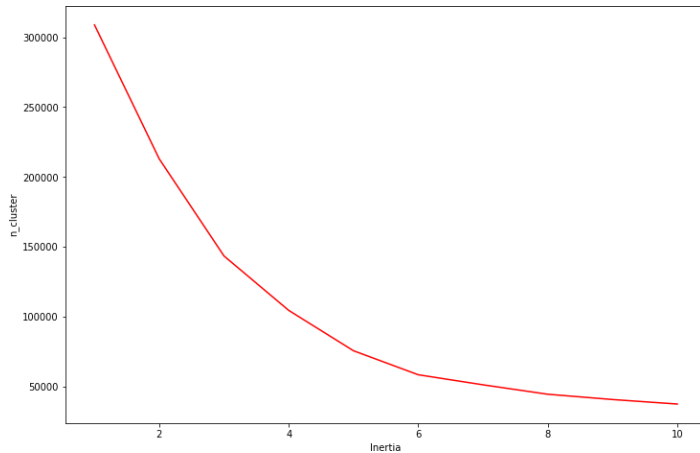
Finding K

```
from sklearn.cluster import KMeans
cluster = []
for k in range (1, 11):
    kmean = KMeans(n_clusters=k).fit(B)
    cluster.append(kmean.inertia_)
```

```
import matplotlib.pyplot as plt
plt.figure(figsize=(12, 8))
```

```
}SPM,BSIOTR,WAGHOLI,PUNE
```

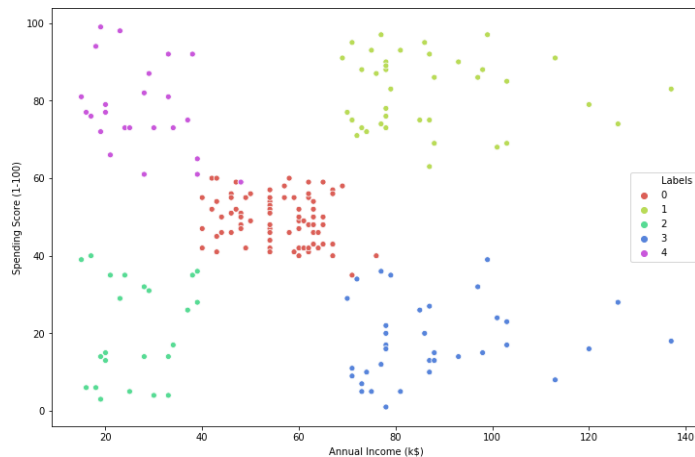
```
plt.plot(range(1, 11), cluster, 'r-')
plt.xlabel('Inertia')
plt.ylabel('n_cluster')
plt.show()
```



With above value of K, Create K means clustering model

```
km = KMeans(n_clusters=5).fit(B)
B['Labels'] = km.labels_
```

```
import seaborn as sns
plt.figure(figsize=(12, 8))
sns.scatterplot(B['Annual Income (k$)'], B['Spending Score (1-100)'], hue=B['Labels'],
               palette=sns.color_palette('hls', 5))
plt.show()
```



References:

[1]Andrea Trevino,Introduction to K-means Clustering ,Oracle AI & Data Science Blog

Assignment No -4

Title: Association Rule Learning

Problem Statement:

Download Market Basket Optimization dataset from below link.

Data Set: <https://www.kaggle.com/hemanthkumar05/market-basket-optimization>

This dataset comprises the list of transactions of a retail company over the period of one week. It contains a total of 7501 transaction records where each record consists of the list of items sold in one transaction. Using this record of transactions and items in each transaction, find the association rules between items.

There is no header in the dataset and the first row contains the first transaction, so mentioned header = None here while loading dataset.

- a. Follow following steps :
- b. Data Preprocessing
- c. Generate the list of transactions from the dataset
- d. Train Apriori algorithm on the dataset
- e. Visualize the list of rules

Generated rules depend on the values of hyper parameters. By increasing the minimum confidence value and find the rules accordingly

Objective:

This assignment will help the students to understand and implement Apriori Algorithm.

S/W Packages and H/W apparatus used:

Linux OS: Ubuntu/Windows , Jupyter notebook.

PC with the configuration as Pentium IV 1.7 GHz. 128M.B RAM, 40 G.B HDD, 15''Color Monitor, Keyboard, Mouse

References:

1. Ethem Alpaydin, Introduction to Machine Learning, PHI 2nd Edition, 2013.
2. Peter Flach: Machine Learning: The Art and Science of Algorithms that Make Sense of Data, Cambridge University Press, Edition 2012.

Theory:

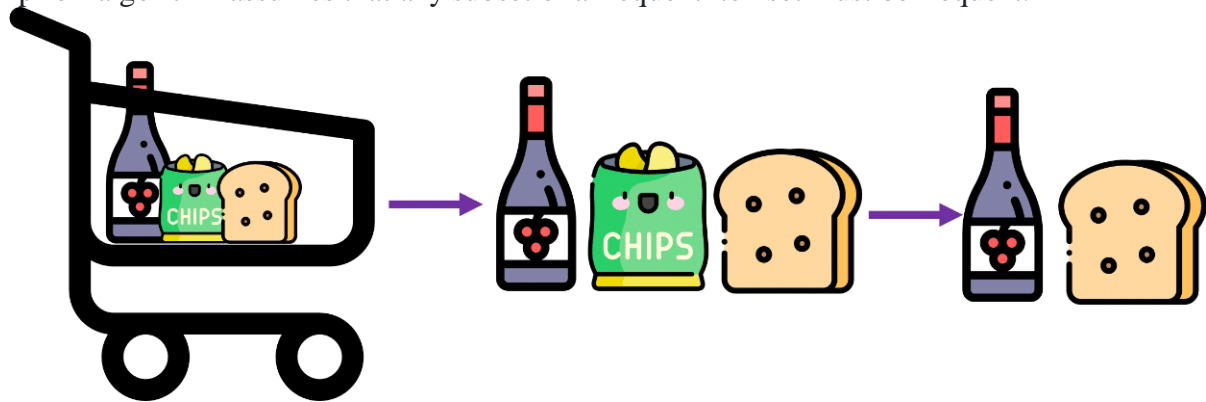
the Apriori algorithm is used for the purpose of [association rule mining](#). Association rule mining is a technique to identify frequent patterns and associations among a set of items.

For example, understanding customer buying habits. By finding correlations and associations between different items that customers place in their ‘shopping basket,’ recurring patterns can be derived.

This process of identifying an association between products/items is called association rule mining. To implement association rule mining, many algorithms have been developed. Apriori algorithm is one of the most popular and arguably the most efficient algorithms among them.

Apriori Algorithm

Apriori algorithm assumes that any subset of a frequent itemset must be frequent.



Say, a transaction containing {wine, chips, bread} also contains {wine, bread}. So, according to the principle of Apriori, if {wine, chips, bread} is frequent, then {wine, bread} must also be frequent.

The key concept in the Apriori algorithm is that it assumes all subsets of a frequent itemset to be frequent. Similarly, for any infrequent itemset, all its supersets must also be infrequent.

Here is a dataset consisting of six transactions in an hour. Each transaction is a combination of 0s and 1s, where 0 represents the absence of an item and 1 represents the presence of it.

Here is a dataset consisting of six transactions in an hour. Each transaction is a combination of 0s and 1s, where 0 represents the absence of an item and 1 represents the presence of it.

Transaction ID	Wine	Chips	Bread	Milk

1	1	1	1	1
2	1	0	1	1
3	0	0	1	1
4	0	1	0	0
5	1	1	1	1
6	1	1	0	1

We can find multiple rules from this scenario. For example, in a transaction of wine, chips, and bread, if wine and chips are bought, then customers also buy bread.

`{ wine, chips } => { bread }`

In order to select the interesting rules out of multiple possible rules from this small business scenario, we will be using the following measures:

- Support
- Confidence
- Lift
- Conviction

Support

}SPM,BSIOTR,WAGHOLI,PUNE

Support of item x is nothing but the ratio of the number of transactions in which item x appears to the total number of transactions.

i.e.,

$$\text{Support(wine)} = \frac{\text{Number of transactions in which the item wine appears}}{\text{Total number of transactions}}$$

$$\text{Support(wine)} = 4/6 = 0.66667$$

Confidence

Confidence ($x \Rightarrow y$) signifies the likelihood of the item y being purchased when item x is purchased.

This method takes into account the popularity of item x .

i.e.,

$$\text{Conf}(\{\text{wine, chips}\} \Rightarrow \{\text{bread}\}) = \frac{\text{support(wine,chips,bread)}}{\text{support(wine,chips)}}$$

$$\frac{2}{6} = \frac{1}{3} :$$

$$\text{Conf}(\{\text{wine, chips}\} \Rightarrow \{\text{bread}\}) = 0.667$$

Lift

Lift ($x \Rightarrow y$) is nothing but the ‘interestingness’ or the likelihood of the item y being purchased when item x is sold. Unlike confidence ($x \Rightarrow y$), this method takes into account the popularity of the item y .

i.e.,

$$\text{lift}(\{ \text{wine, chips} \} \Rightarrow \{ \text{bread} \}) = \frac{\text{support}(\text{wine, chips, bread})}{\text{support}(\text{wine, chips})}$$

$$\text{lift}(\{ \text{wine, chips} \} \Rightarrow \{ \text{bread} \}) = \frac{\frac{2}{6}}{\frac{3}{6} * \frac{4}{6}} = 1$$

- Lift $(x \Rightarrow y) = 1$ means that there is no correlation within the itemset.
- Lift $(x \Rightarrow y) > 1$ means that there is a positive correlation within the itemset, i.e., products in the itemset, x and y , are more likely to be bought together.
- Lift $(x \Rightarrow y) < 1$ means that there is a negative correlation within the itemset, i.e., products in itemset, x and y , are unlikely to be bought together.

Dataset

Below is the transaction data from Day 1. This dataset contains 6 items and 22 transaction records.

	A	B	C	D	E	F
1	Wine	Chips	Bread	Butter	Milk	Apple
2	Wine		Bread	Butter	Milk	
3			Bread	Butter	Milk	
4		Chips				Apple
5	Wine	Chips	Bread	Butter	Milk	Apple
6	Wine	Chips			Milk	
7	Wine	Chips	Bread	Butter		Apple
8	Wine	Chips			Milk	
9	Wine		Bread			Apple
10	Wine		Bread	Butter	Milk	
11		Chips	Bread	Butter		Apple
12	Wine			Butter	Milk	Apple
13	Wine	Chips	Bread	Butter	Milk	
14	Wine		Bread		Milk	Apple
15	Wine		Bread	Butter	Milk	Apple
16	Wine	Chips	Bread	Butter	Milk	Apple
17		Chips	Bread	Butter	Milk	Apple
18		Chips		Butter	Milk	Apple
19	Wine	Chips	Bread	Butter	Milk	Apple
20	Wine		Bread	Butter	Milk	Apple
21	Wine	Chips	Bread		Milk	Apple
22		Chips				

Environment Setup:

Before we move forward, we need to install the ‘apriori’ package first.

```
pip install apyori
```

Market Basket Analysis Implementation within Python

With the help of the apyori package, we will be implementing the Apriori algorithm in order to help the manager in market basket analysis.

Step 1: Import the libraries

```
In [1]: #Importing the required datasets
import numpy as np
import pandas as pd
from apyori import apriori
```

Step 2: Load the dataset

```
In [2]: #Loading the dataset
store_data = pd.read_csv('Day1.csv', header=None)
```

Step 3: Have a glance at the records

```
In [3]: #Having a glance at the records
store_data
```

```
Out[3]:
```

	0	1	2	3	4	5
0	Wine	Chips	Bread	Butter	Milk	Apple
1	Wine	NaN	Bread	Butter	Milk	NaN
2	NaN	NaN	Bread	Butter	Milk	NaN
3	NaN	Chips	NaN	NaN	NaN	Apple
4	Wine	Chips	Bread	Butter	Milk	Apple
5	Wine	Chips	NaN	NaN	Milk	NaN
6	Wine	Chips	Bread	Butter	NaN	Apple
7	Wine	Chips	NaN	NaN	Milk	NaN
8	Wine	NaN	Bread	NaN	NaN	Apple
9	Wine	NaN	Bread	Butter	Milk	NaN
10	NaN	Chips	Bread	Butter	NaN	Apple
11	Wine	NaN	NaN	Butter	Milk	Apple
12	Wine	Chips	Bread	Butter	Milk	NaN
13	Wine	NaN	Bread	NaN	Milk	Apple
14	Wine	NaN	Bread	Butter	Milk	Apple
15	Wine	Chips	Bread	Butter	Milk	Apple
16	NaN	Chips	Bread	Butter	Milk	Apple
17	NaN	Chips	NaN	Butter	Milk	Apple
18	Wine	Chips	Bread	Butter	Milk	Apple
19	Wine	NaN	Bread	Butter	Milk	Apple
20	Wine	Chips	Bread	NaN	Milk	Apple
21	NaN	Chips	NaN	NaN	NaN	NaN

Step 4 :Convert Pandas DataFrame into a list of lists

```
In [5]: #Converting the pandas dataframe into a list of lists
records = []
for i in range(0, 22):
    records.append([str(store_data.values[i,j]) for j in range(0, 6)])
```

Step 5: Build the Apriori model

```
In [7]: #Building the first apriori model
association_rules = apriori(records, min_support=0.50, min_confidence=0.7, min_lift=1.2, min_length=2)
association_results = list(association_rules)
```

Step 6: Print out the number of rules

```
In [8]: #Getting the number of rules
print(len(association_results))
```

Step 7: Have a glance at the rule

```
In [10]: #Glancing at the first rule
print(association_results)
```

```
[RelationRecord(items=frozenset({'Milk', 'Butter', 'Bread'}), support=0.5, ordered_statistics=[OrderedStatistic(items_base=frozenset({'Milk', 'Bread'}), items_add=frozenset({'Butter'}), confidence=0.8461538461538461, lift=1.241025641025641)])]
```

The support value for the first rule is 0.5. This number is calculated by dividing the number of transactions containing ‘Milk,’ ‘Bread,’ and ‘Butter’ by the total number of transactions.

The confidence level for the rule is 0.846, which shows that out of all the transactions that contain both “Milk” and “Bread”, 84.6 % contain ‘Butter’ too.

The lift of 1.241 tells us that 'Butter' is 1.241 times more likely to be bought by the customers who buy both 'Milk' and 'Butter' compared to the default likelihood sale of 'Butter.'

Conclusion:

Apriori algorithm is implemented

PART B:

Advanced Database Management Systems

Assignment: 1

AIM:

Create a database with suitable example using MongoDB and implement

1. Inserting and saving document (batch insert, insert validation)
2. Removing document
3. Updating document (document replacement, using modifiers, up inserts, updating multiple documents, returning updated documents)
4. Execute at least 10 queries on any suitable MongoDB database that demonstrates following:
 - Find and find One (specific values)
 - Query criteria (Query conditionals, OR queries, \$not, Conditional semantics)
 - Type-specific queries (Null, Regular expression, Querying arrays) where queries
 - Cursors (Limit, skip, sort, advanced query options)

PROBLEM STATEMENT /DEFINITION

Create a database with suitable example using MongoDB and implement

- Inserting and saving document (batch insert, insert validation)
- Removing document
- Updating document (document replacement, using modifiers, upserts, updating multiple documents, returning updated documents)

OBJECTIVE:

To understand MongoDB basic commands

To implement the concept of document oriented databases.

To understand MongoDB retrieval commands

THEORY:

SQL Vs MongoDB

SQL Concepts	MongoDB Concepts
Database	Database
Table	Collection

Row	Document or BSON Document
Column	Field
Index	Index
Table Join	Embedded Documents & Linking
Primary Key	Primary Key
Specify Any Unique Column Or Column Combination As Primary Key.	In MongoDB, The Primary Key Is Automatically Set To The <u>_id</u> Field.
Aggregation (E.G. Group By)	Aggregation Pipeline

1. Create a collection in mongodb

```
db.createCollection("Teacher_info")
```

2. Create a capped collection in mongodb

```
>db.createCollection("audit", {capped:true, size:20480})
{ "ok" : 1 }
```

3. Insert a document into collection

```
db.Teacher_info.insert( { Teacher_id: "Pic001", Teacher_Name: "Ravi",Dept_Name:
"IT", Sal:30000, status: "A" } )
```

```
db.Teacher_info.insert( { Teacher_id: "Pic002", Teacher_Name: "Ravi",Dept_Name:
"IT", Sal:20000, status: "A" } )
```

```
db.Teacher_info.insert( { Teacher_id: "Pic003", Teacher_Name: "Akshay",Dept_Name:
"Comp", Sal:25000, status: "N" } )
```

4. Update a document into collection

```
db. Teacher_info.update( { sal: { $gt: 25000 } }, { $set: { Dept_name: "ETC" } }, { multi: true } )
```

```
db. Teacher_info.update( { status: "A" } , { $inc: { sal: 10000 } }, { multi: true } )
```

5. Remove a document from collection

```
db.Teacher_info.remove({ Teacher_id: "pic001"}); db. Teacher_info.remove({})
```

6. Alter a field into a mongodb document

```
db.Teacher_info.update( { }, { $set: { join_date: new Date() } }, { multi: true } )
```

7. To drop a particular collection

```
db.Teacher_info.drop()
```

Retrieval From Database:-

When we retrieve a document from mongodb collection it always add a `_id` field in the every document which contain unique `_id` field.

ObjectId(<hexadecimal>)

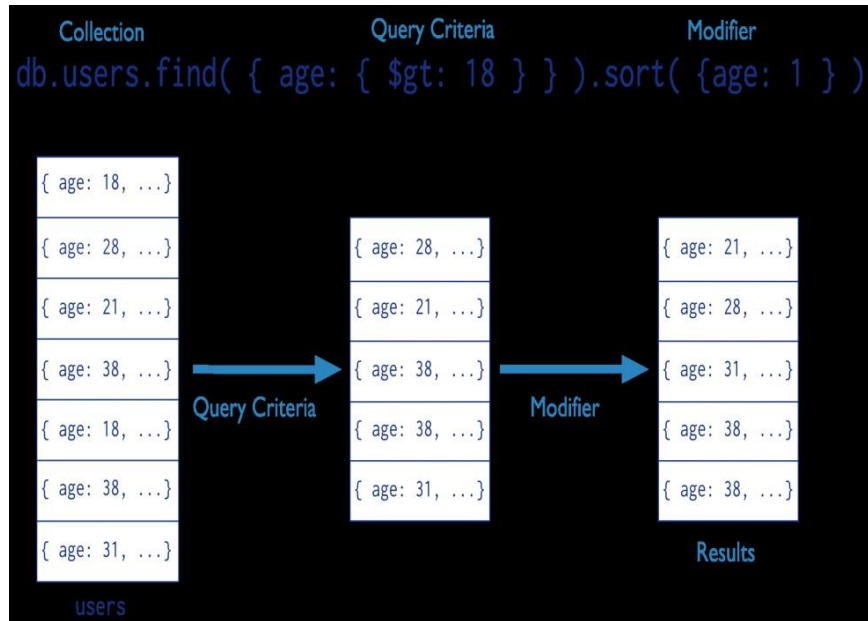
Returns a new ObjectId value. The 12-byte ObjectId value consists of: 4-byte

value representing the seconds since the Unix epoch,

3-byte machine identifier,

2-byte process id, and

3-byte counter, starting with a random value.



1. Retrieve a collection in mongodb using Find command

`db.Teacher.find()`

```
{ "_id" : 101, "Name" : "Dev",
  "Address" : [ { "City" : "Pune", "Pin" : 444043 } ],
  "Department" : [ { "Dept_id" : 111, "Dept_name" : "IT" } ], "Salary" : 78000 }

{ "_id" : 135, "Name" : "Jennifer", "Address" : [ { "City" :
  "Mumbai", "Pin" : 444111 } ], "Department" : [ { "Dept_id" : 112, "Dept_name" :
  "COMP" } ], "Salary" : 65000 }
{ "_id" : 126, "Name" : "Gaurav", "Address" : [ { "City" : "Nashik",
  "Pin" : 444198 } ], "Department" : [ { "Dept_id" : 112, "Dept_name"
  : "COMP" } ], "Salary" : 90000 }
{ "_id" : 175, "Name" : "Shree", "Address" : [ { "City" : "Nagpur",
  "Pin" : 444158 } ], "Department" : [ { "Dept_id" : 113, "Dept_name"
  : "ENTC" } ], "Salary" : 42000 }
{ "_id" : 587, "Name" : "Raman", "Address" : [ { "City" : "Bangalore",
  "Pin" : 445754 } ], "Department" : [ { "Dept_id" : 113, "Dept_name"
  : "ENTC" } ], "Salary" : 79000 }
{ "_id" : 674, "Name" : "Mandar", "Address" : [ { "City" : "Jalgaon",
  "Pin" : 465487 } ], "Department" : [ { "Dept_id" : 111, "Dept_name"
  : "IT" } ], "Salary" : 88000 }
{ "_id" : 573, "Name" : "Manish", "Address" : [ { "City" : "Washim",
  "Pin" : 547353 } ], "Department" : [ { "Dept_id" : 112, "Dept_name"
  : "COMP" } ], "Salary" : 65000 }
```

2. Retrieve a document from collection in mongodb using Find command using condition

`{SPM,BSIOTR,WAGHOLI,PUNE`

```
>db.Teacher_info.find({sal: 25000})
```

3. Retrieve a document from collection in mongodb using Find command using or operator

```
>db.Teacher_info.find( { $or: [ { status: "A" } , { sal:50000 } ] } )
```

4. Retrieve a document from collection in mongodb using Find command using greater than , less than, greater than and equal to ,less than and equal to operator

```
>db. Teacher_info.find( { sal: { $gt: 40000 } } )
```

```
>db.media.find( { Released : { $gt : 2000 } }, { "Cast" : 0 } )
```

```
{ "_id" : ObjectId("4c4369a3c603000000007ed3"), "Type" : "DVD", "Title" : "Toy Story 3", "Released" : 2010 }
```

```
>db.media.find ( { Released : { $gte : 1999 } }, { "Cast" : 0 } )
```

```
{ "_id" : ObjectId("4c43694bc603000000007ed1"), "Type" : "DVD", "Title" : "Matrix, The", "Released" : 1999 }
```

```
{ "_id" : ObjectId("4c4369a3c603000000007ed3"), "Type" : "DVD", "Title" : "Toy Story 3", "Released" : 2010 }
```

```
>db.media.find ( { Released : { $lt : 1999 } }, { "Cast" : 0 } )
```

```
{ "_id" : ObjectId("4c436969c603000000007ed2"), "Type" : "DVD", "Title" : "Blade Runner", "Released" : 1982 }
```

```
>db.media.find( { Released : { $lte: 1999 } }, { "Cast" : 0 } )
```

```
{ "_id" : ObjectId("4c43694bc603000000007ed1"), "Type" : "DVD", "Title" : "Matrix, The", "Released" : 1999 }
```

```
{ "_id" : ObjectId("4c436969c603000000007ed2"), "Type" : "DVD", "Title" : "Blade Runner", "Released" : 1982 }
```

```
>db.media.find( { Released : { $gte: 1990, $lt : 2010 } }, { "Cast" : 0 } )
```

```
{ "_id" : ObjectId("4c43694bc603000000007ed1"), "Type" : "DVD", "Title" : "Matrix, The", "Released" : 1999 }
```

Retrieval a value from document which contain array field Exact Match

on an Array

```
db.inventory.find( { tags: [ 'fruit', 'food', 'citrus' ] } )
```

Match an Array Element db.inventory.find({

tags: 'fruit' }) **Match a Specific Element of an**

Array db.inventory.find({ 'tags.0' : 'fruit' })

6. MongoDB provides a db.collection.findOne() method as a special case of find() that returns a single document.

7. Exclude One Field from a Result Set

```
>db.records.Find( { "user_id": { $lt: 42 } }, { history: 0 } )
```

8. Return Two fields and the _id Field

```
>db.records.find( { "user_id": { $lt: 42 } }, { "name": 1, "email": 1 } )
```

9. Return Two Fields and Exclude _id

```
>db.records.find( { "user_id": { $lt: 42 } }, { "_id": 0, "name": 1, "email": 1 } )
```

10. Retrieve a collection in mongodb using Find command and pretty appearance

```
>db.<collection>.find().pretty()
```

db.users.find(collection
{age:{\$gt:18}},	query
criteria	
{name :1,address:1}	projection

Retrieve a document in ascending or descending order using 1 for ascending and -1 for descending from collection in mongodb

```
>db.Teacher_info.find( { status: "A" } ).sort( { sal: -1 } )
```

```
>db.audit.find().sort( { $natural: -1 } ).limit ( 10 )
```

```
>db.Employee.find().sort({_id:-1})
```

```
{ "_id" : 106, "Name" : "RAJ", "Address" : [ { "City" : "NASIK", "Pin" :
411002 } ], "Department" : [ { "Dept_id" : 113, "Dept_name" : "ACCOUNTING"
} ], "Salary" : 50000 }
{ "_id" : 105, "Name" : "ASHOK", "Address" : [ { "City" : "NASIK", "Pin" :
411002 } ], "Department" : [ { "Dept_id" : 113, "Dept_name" : "ACCOUNTING"
} ], "Salary" : 40000 }
{ "_id" : 104, "Name" : "JOY", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 20000 }
{ "_id" : 103, "Name" : "RAM", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 10000 }
{ "_id" : 102, "Name" : "AKASH", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 111, "Dept_name" : "HR" } ], "Salary" : 80000 }
{ "_id" : 101, "Name" : "Dev", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 111, "Dept_name" : "HR" } ], "Salary" : 78000 }
```

```
>db.Employee.find().sort({_id:1})
```

```
{ "_id" : 101, "Name" : "Dev", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 111, "Dept_name" : "HR" } ], "Salary" : 78000 }
{ "_id" : 102, "Name" : "AKASH", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 111, "Dept_name" : "HR" } ], "Salary" : 80000 }
```



```
{ "_id" : 103, "Name" : "RAM", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 10000 }
{ "_id" : 104, "Name" : "JOY", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 20000 }
{ "_id" : 105, "Name" : "ASHOK", "Address" : [ { "City" : "NASIK", "Pin" :
411002 } ], "Department" : [ { "Dept_id" : 113, "Dept_name" : "ACCOUNTING"
} ], "Salary" : 40000 }
{ "_id" : 106, "Name" : "RAJ", "Address" : [ { "City" : "NASIK", "Pin" :
411002 } ], "Department" : [ { "Dept_id" : 113, "Dept_name" : "ACCOUNTING"
} ], "Salary" : 50000 }
>db.Employee.find().sort({$natural:-1}).limit(2)
{ "_id" : 106, "Name" : "RAJ", "Address" : [ { "City" : "NASIK", "Pin" :
411002 } ], "Department" : [ { "Dept_id" : 113, "Dept_name" : "ACCOUNTING"
} ], "Salary" : 50000 }
{ "_id" : 105, "Name" : "ASHOK", "Address" : [ { "City" : "NASIK", "Pin" :
411002 } ], "Department" : [ { "Dept_id" : 113, "Dept_name" : "ACCOUNTING"
} ], "Salary" : 40000 }
```

>db.Employee.find().sort({\$natural:1}).limit(2)

```
{ "_id" : 101, "Name" : "Dev", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 111, "Dept_name" : "HR" } ], "Salary" : 78000 }
{ "_id" : 102, "Name" : "AKASH", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 111, "Dept_name" : "HR" } ], "Salary" : 80000 }
>db.Employee.find({Salary:{$in:[10000,30000]}})
{ "_id" : 103, "Name" : "RAM", "Address" : [ { "City" : "Pune", "Pin" :
444043 } ], "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 10000 }
>db.Employee.update({"Name":"RAM"},{ $set :{Address:{City: "Nasik"}}}) WriteResult({
  "nMatched" : 1, "nUpserted" : 0, "nModified" : 1 })
>db.Employee.find({"Name":"RAM"})
{ "_id" : 103, "Name" : "RAM", "Address" : { "City" : "Nasik" },
  "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 10000 }
>db.Employee.update({"Name":"RAM"},{$inc :{"Salary": 10000 } }) WriteResult({
  "nMatched" : 1, "nUpserted" : 0, "nModified" : 1 })
>db.Employee.find({"Name":"RAM"})
{ "_id" : 103, "Name" : "RAM", "Address" : { "City" : "Nasik" },
  "Department" : [ { "Dept_id" : 112, "Dept_name" : "SALES" } ], "Salary" : 20000 }
```

Retrieve document with a particular from collection in mongodb

>db.Employee.find().limit(2).pretty()

```
{
  "_id" : 101,
```

```
"Name" : "Dev",
"Address" : [
  {
    "City" : "Pune", "Pin" :
    444043
  }
],
"Department" : [
  {
    "Dept_id" : 111,
    "Dept_name" : "HR"
  }
],
"Salary" : 78000
}
{
  "_id" : 102, "Name" :
  "AKASH",
  "Address" : [
    {
      "City" : "Pune", "Pin" :
      444043
    }
  ],
  "Department" : [
    {
      "Dept_id" : 111,
      "Dept_name" : "HR"
    }
  ],
  "Salary" : 80000
}
```

Retrieve document skipping some documents from collection in mongodb

```
>db.Employee.find().skip(3).pretty()
```

```
{
  "_id" : 104, "Name" :
  "JOY",
  "Address" : [
    {
      "City" : "Pune", "Pin" :
      444043
    }
  ],
  "Department" : [
    {
      "Dept_id" : 112, "Dept_name" :
      "SALES"
    }
  ]
}
```

```
}SPM,BSIOTR,WAGHOLI,PUNE
```

```
}
{
  "_id" : 105, "Name" :
  "ASHOK",
  "Address" : [
    {
      "City" : "NASIK",
      "Pin" : 411002
    }
  ],
  "Department" : [
    {
      "Dept_id" : 113, "Dept_name" :
      "ACCOUNTING"
    }
  ],
  "Salary" : 40000
}
{
  "_id" : 106, "Name" :
  "RAJ",
  "Address" : [
    {
      "City" : "NASIK",
      "Pin" : 411002
    }
  ],
  "Department" : [
    {
      "Dept_id" : 113, "Dept_name" :
      "ACCOUNTING"
    }
  ],
  "Salary" : 50000
}
```

REFERENCE BOOK:

Kristina Chodorow, MongoDB The definitive guide, O'Reilly Publications, ISBN:978-93-5110-269-4, 2nd Edition.

CONCLUSION:

Understand to implement data from mongodb database with the help of statement and operators.

Assignment: 2

AIM: Implement Map reduces operation with suitable example on above MongoDB database

- Aggregation framework
- Create and drop different types of indexes and explain () to show the advantage of the indexes.

PROBLEM STATEMENT /DEFINITION

Implement Map reduces operation with suitable example on above MongoDB database

- Aggregation framework
- Create and drop different types of indexes and explain () to show the advantage of the indexes.

OBJECTIVE:

To understand the concept of Mapreduce in mongodb.

To understand the concept of Aggregation in mongodb.

To implement the concept of document oriented databases.

THEORY:

- Implements the MapReduce model for processing large data sets.
- Can choose from one of several output options (inline, new collection, merge, replace, reduce)
- MapReduce functions are written in JavaScript.
- Supports non-sharded and sharded input collections.
- Can be used for incremental aggregation over large collections.
- MongoDB 2.2 implements much better support for sharded map reduce output.
- New feature in the Mongodb2.2.0 production release (August, 2012).
- Designed with specific goals of **improving performance** and **usability**.
- Returns result set inline.
- Supports **non-sharded and sharded** input collections.

- Uses a "**pipeline**" approach where objects are transformed as they pass through a series of pipeline operators such as matching, projecting, sorting, and grouping.
- **Pipeline operators need not produce one output document for every input document:** operators may also generate new documents or filter out documents.
- Map/Reduce involves two steps:
- first, map the data from the collection specified;
- second, reduce the results.

```
>db.createCollection("Order")
```

- { "ok" : 1 }

```
>db.order.insert({cust_id:"A123",amount:500,status:"A"})
```

- WriteResult({ "nInserted" : 1 })

```
>db.order.insert({cust_id:"A123",amount:250,status:"A"})
```

- WriteResult({ "nInserted" : 1 })

```
>db.order.insert({cust_id:"B212",amount:200,status:"A"})
```

- WriteResult({ "nInserted" : 1 })

```
>db.order.insert({cust_id:"A123",amount:300,status:"d"})
```

- WriteResult({ "nInserted" : 1 })

- **Map Function**

- var mapFunction1 = function()
- { emit(this.cust_id, this.amount);};

- **Reduce Function**

- var reduceFunction1 = function(key, values)
- {return Array.sum(values);};

db.order.mapReduce

(mapFunction1, reduceFunction1, {query: {status: "A" },
out: "order_totals"});

```
"result" : "order_totals",
"timeMillis" : 28, "counts" : {
  "input" : 3,
  "emit" : 3,
  "reduce" : 1,
  "output" : 2
},
"ok" : 1,}

>db.order.mapReduce(
Map Function ->    function() { emit( this.cust_id, this.amount);},
Reduce   Function -> function( key, values ) { return Array.sum ( values )}, Query à
    {query: { status:"A"}},
Output collection à    out: "order_    totals"})
{
  "result" : "order_totals",
  "timeMillis" : 27, "counts" : {
    "input" : 3,
    "emit" : 3,
    "reduce" : 1,
    "output" : 2
  },
  "ok" : 1,
}
```

To display result of mapReduce function use collection created in OUT.

Db.<collection name>.find();

```
db.order_totals.find();

{ "_id" : "A123", "value" : 750 }

{ "_id" : "B212", "value" : 200 }
```

Implementation of Aggregation:-

```
> use Teacher
switched to db Teacher
>db.Teacher.find()
{ "_id" : 101, "Name" : "Dev", "Address" : [ { "City" : "Pune", "Pin" : 444043 } ],
"Department" : [ { "Dept_id" : 111, "Dept_name" : "IT" } ], "Salary" : 78000 }
{ "_id" : 135, "Name" : "Jennifer", "Address" : [ { "City" : "Mumbai", "Pin" : 444111 } ],
"Department" : [ { "Dept_id" : 112, "Dept_name" : "COMP" } ], "Salary" : 65000 }
{ "_id" : 126, "Name" : "Gaurav", "Address" : [ { "City" : "Nashik", "Pin" : 444198 } ],
"Department" : [ { "Dept_id" : 112, "Dept_name" : "COMP" } ], "Salary" : 90000 }
{ "_id" : 175, "Name" : "Shree", "Address" : [ { "City" : "Nagpur", "Pin" : 444158 } ],
"Department" : [ { "Dept_id" : 113, "Dept_name" : "ENTC" } ], "Salary" : 42000 }
{ "_id" : 587, "Name" : "Raman", "Address" : [ { "City" : "Banglore", "Pin" : 445754 } ],
"Department" : [ { "Dept_id" : 113, "Dept_name" : "ENTC" } ], "Salary" : 79000 }
{ "_id" : 674, "Name" : "Mandar", "Address" : [ { "City" : "Jalgaon", "Pin" : 465487 } ],
"Department" : [ { "Dept_id" : 111, "Dept_name" : "IT" } ], "Salary" : 88000 }
{ "_id" : 573, "Name" : "Manish", "Address" : [ { "City" : "Washim", "Pin" : 547353 } ],
"Department" : [ { "Dept_id" : 112, "Dept_name" : "COMP" } ], "Salary" : 65000 }
```

```
>db.Teacher.aggregate([
... {$group:{$_id:"$Department",totalsalary:{$sum:"$Salary"}}}
... ])
{
  "result" : [
    {
      "_id" : [
        {
          "Dept_id" : 113,
          "Dept_name" : "ENTC"
        }
      ],
      "totalsalary" : 121000
    },
    {
      "_id" : [
        {
          "Dept_id" : 112,
          "Dept_name" : "COMP"
        }
      ]
    }
  ]
}
```

```

        ],
        "totalsalary" : 220000
    },
    {
        "_id" : [
            {
                "Dept_id" : 111,
                "Dept_name" : "IT"
            }
        ],
        "totalsalary" : 166000
    }
],
"ok" : 1
}
>db.Teacher.aggregate([
{$group:{_id:"$Department",totalsalary:{$sum:"$Salary"}},{ $group:{_id:"$_id.Department",AvgSal:{$sum:"$totalsalary"}}})
{ "result" : [ { "_id" : [ ], "AvgSal" : 507000 } ], "ok" : 1 }
>db.Teacher.aggregate([
{$group:{_id:"$Department",totalsalary:{$sum:"$Salary"}},{ $match:{totalsalary:{$gte:200000}}})
{
    "result" : [
        {
            "_id" : [
                {
                    "Dept_id" : 112,
                    "Dept_name" : "COMP"
                }
            ],
            "totalsalary" : 220000
        }
    ],
    "ok" : 1
}
>db.Teacher.aggregate([ {$group:{_id:"$Department",totalsalary:{$sum:"$Salary"}}, {
$sort:{totalsalary:1}}])
{
    "result" : [
        {
            "_id" : [
                {
                    "Dept_id" : 113,
                    "Dept_name" : "ENTC"
                }
            ]
        }
    ]
}

```



```

        ],
        "totalsalary" : 121000
    },
    {
        "_id" : [
            {
                "Dept_id" : 111,
                "Dept_name" : "IT"
            }
        ],
        "totalsalary" : 166000
    },
    {
        "_id" : [
            {
                "Dept_id" : 112,
                "Dept_name" : "COMP"
            }
        ],
        "totalsalary" : 220000
    }
],
"ok" : 1
}

>db.Teacher.aggregate([ {$group:{$_id:"$Department",totalsalary:{$sum:"$Salary"}}}, {
$group: { $_id:"$_id.Department", big: { $last: "$_id.Dept_name" }, bigsalary: {
$last:"$totalsalary"}, small: { $first:"$_id.Dept_name"}, smallsalary: {
$first:"$totalsalary"} } } ])
{
    "result" : [
        {
            "_id" : [],
            "big" : [
                "IT"
            ],
            "bigsalary" : 166000,
            "small" : [
                "ENTC"
            ],
            "smallsalary" : 121000
        }
    ],
    "ok" : 1
}

```

REFERENCE BOOK:

Kristina Chodorow, MongoDB The definitive guide, O'Reilly Publications, ISBN:978-93-5110-269-4, 2nd Edition.

CONCLUSION:

Understand to mapreduce operation in mongodb

Assignment 3

Aim : Case Study: Design conceptual model using Star and Snowflake schema for any one

database **Problem statement:** Design conceptual model using Star and Snowflake schema for

any one database **OBJECTIVE:**

1. To understand concepts of multidimensional data.
2. To understand the the relational implementation of the multidimensional data model is typically a star schema, or a snowflake schema.

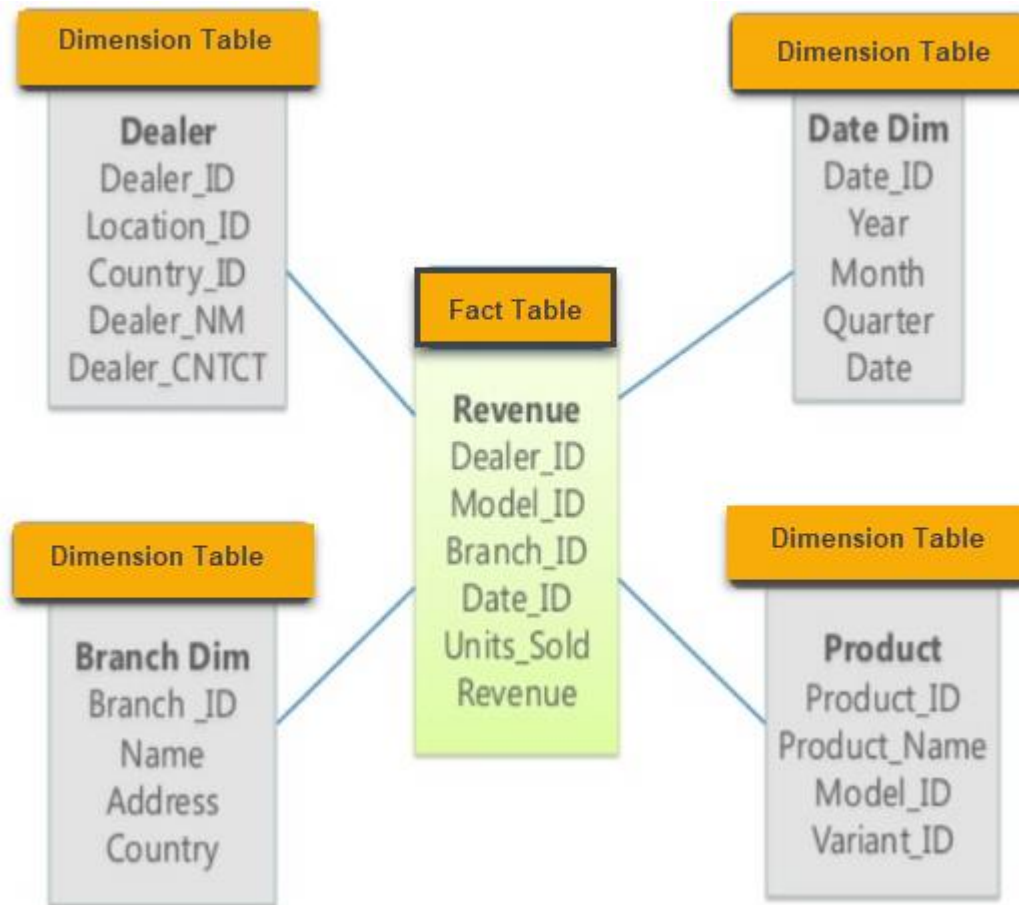
Theory:

The data warehouses are considered modern ancient techniques, since the early days for the relational databases, the idea of the keeping a historical data for reference when it needed has been originated, and the idea was primitive to create archives for the historical data to save these data, despite of the usage of a special techniques for the recovery of these data from the different storage modes. This research applied of structured databases for a trading company operating across the continents, has a set of branches each one has its own stores and showrooms, and the company branch's group of sections with specific activities, such as stores management, showrooms management, accounting management, contracts and other departments. It also assumes that the company center exported software to manage databases for all branches to ensure the safety performance, standardization of processors and prevent the possible errors and bottlenecks problems. Also the research provides this methods the best requirements have been used for the applied of the data warehouse (DW), the information that managed by such an applied must be with high accuracy. It must be emphasized to ensure compatibility information and hedge its security, in schemes domain, been applied to a comparison between the two schemes (Star and Snowflake Schemas) with the concepts of multidimensional database. It turns out that Star Schema is better than Snowflake Schema in (Query complexity, Query performance, Foreign Key Joins), And finally it has been concluded that Star Schema center fact and change, while Snowflake Schema center fact and not change.

Example:

1. *Star Schema*

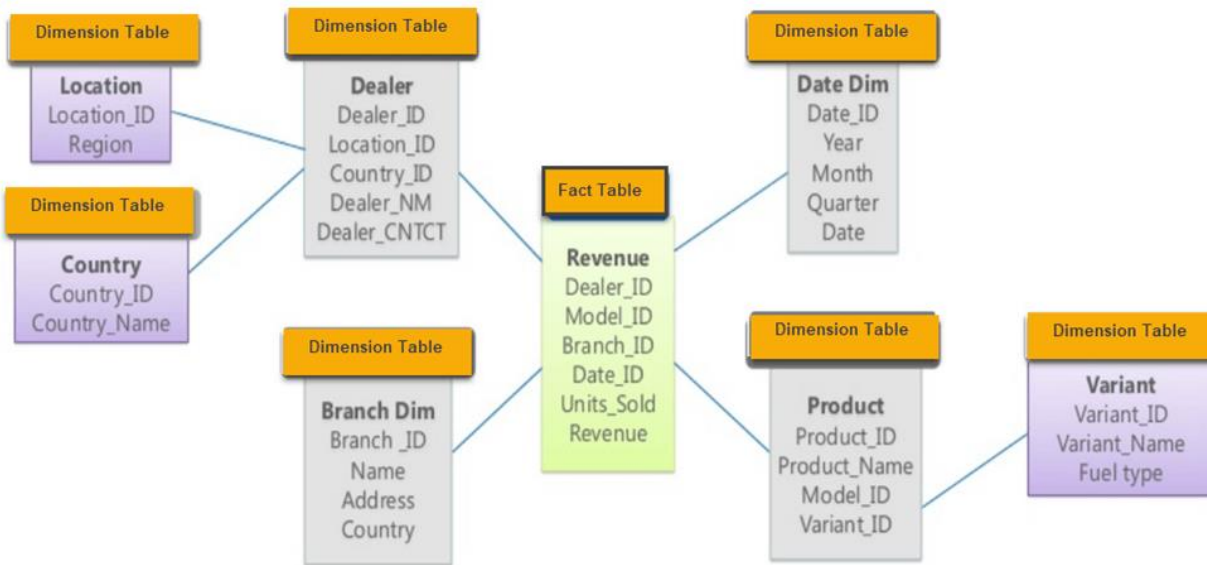
Star Schema in data warehouse, in which the center of the star can have one fact table and a number of associated dimension tables. It is known as star schema as its structure resembles a star. The Star Schema data model is the simplest type of Data Warehouse schema. It is also known as Star Join Schema and is optimized for querying large data sets



2. Snowflake Schema?

Snowflake Schema in data warehouse is a logical arrangement of tables in a multidimensional database such that the [ER diagram](#) resembles a snowflake shape. A Snowflake Schema is an extension of a Star Schema, and it adds additional dimensions. The dimension tables are normalized which splits data into additional tables.

In the following Snowflake Schema example, Country is further normalized into an individual table.



- Multidimensional schema is especially designed to model data warehouse systems
- The star schema is the simplest type of Data Warehouse schema. It is known as star schema as its structure resembles a star.
- Comparing Snowflake vs Star schema, a Snowflake Schema is an extension of a Star Schema, and it adds additional dimensions. It is called snowflake because its diagram resembles a Snowflake.
- In a star schema, only single join defines the relationship between the fact table and any dimension tables.
- Star schema contains a fact table surrounded by dimension tables.
- Snowflake schema is surrounded by dimension table which are in turn surrounded by dimension table
- A snowflake schema requires many joins to fetch the data.
- Comparing Star vs Snowflake schema, Start schema has simple DB design, while Snowflake schema has very complex DB design.

REFERENCE BOOK:

Jiawei Han, Micheline Kamber, Jian Pei “Data Mining: concepts and techniques”, 2nd Edition,
 Publisher: Elsevier/Morgan Kaufmann.

CONCLUSION:

Understand to concept of multidimensional data.

PART C:

MINI PROJECT

Mini Project

AIM: Build the mini project based on the relevant applicable concepts of Machine Learning / DAA / ADBMS by forming teams of around 3 to 4 students.

A] Sample ML mini-project Format:

PROBLEM STATEMENT /DEFINITION

Design and Implement any Data science Application using Python. Obtain Data, Scrub data (Data cleaning), Explore data, Prepare and validate data model and Interpret data (Data Visualization). Visualize data using any visualization tool like Matplotlib, ggplot, Tableau etc. Prepare Project Report.

OBJECTIVE:

1. To explore the Data science project life cycle.
2. To identify need of project and define problem statement.
3. To extract and process data.
4. To interpret and analyze results using data visualization.

Mini Project Report Format:

Abstract

Acknowledgement

List of Tables & Figures

Contents

1. Introduction
 - 1.1 Purpose, Problem statement
 - 1.2 Scope, Objective
 - 1.3 Definition, Acronym, and Abbreviations
 - 1.4 References
2. Literature Survey
 - 2.1 Introduction
 - 2.2 Detail Literature survey

2.4 Findings of Literature survey

3. System Architecture and Design

3.1 Detail Architecture

3.2 Dataset Description

3.3 Detail Phases

3.4 Algorithms

4. Experimentation and Results

4.1 Phase-wise results

4.2 Explanation with example

4.3 Comparison of result with standard

4.4 Accuracy

4.5 Visualization

4.6 Tools used

5. Conclusion and Future scope

5.1 Conclusion

5.2 Future scope

References

Annexure:

A. GUIs / Screen Snapshot of the System Developed

B. Implementation /code

B] Sample ADBMS mini-project Format:**PROBLEM STATEMENT /DEFINITION**

Build the mini project based on the requirement document and design prepared as a part of Database Management Lab in second year.

Form teams of around 3 to 4 people.

- A. Develop the application: Build a suitable GUI by using forms and placing the controls on it for any application. Proper data entry validations are expected.
- B. Add the database connection with front end. Implement the basic CRUD operations.
- C. Prepare and submit report to include Title of the Project, Abstract, List the hardware and software requirements at the backend and at the front end, Source Code, Graphical User Interface, Conclusion.

OBJECTIVE:

- 3. To understand applications of document-oriented database by implementing mini project.
- 4. To learn effective UI designs.
- 5. To learn to design & implement database system for specific domain.
- 6. To learn to design system architectural & flow diagram.

Mini Project Report Format:

Abstract

Acknowledgement

List of Tables & Figures

Contents

- 1. Introduction
 - 1.1 Purpose
 - 1.2 Scope
 - 1.3 Definition, Acronym, and Abbreviations
 - 1.4 References
 - 1.5 Developers' Responsibilities: An Overview
- 2. General Description
 - 2.1 Product Function Perspective

2.2 User Characteristics.

2.4 General Constraints

2.5 Assumptions and Dependencies

3. Specific Requirements

3.1 Inputs and Outputs

3.2 Functional Requirements

3.3 Functional Interface Requirements

3.3 Performance Constraints

3.4 Design Constraints

3.6 Acceptance criteria

4. System Design

5. System Implementation

5.1 Hardware and Software Platform description

5.2 Tools used

5.3 System Verification and Testing (Test Case Execution)

5.4 Future work / Extension

5.5 Conclusion

References

Annexure:

A. GUIs / Screen Snapshot of the System Developed

THANK YOU!

