Part A    Assignment 3

Title :-
        write an appliction using HBase and Hive
QL for flight information System which will include
    01) Creating , Dropping and alterning database tables
    2> Creating an external Hive table to connect to
    the HBase for customer information table
    3> Load table with data insert new values and 31S
    field in the table Join tables with hive
    4> Create index on flight information table
    5> Find the average departure delay per day in 2008.

objective :- to apply Hbase and HiveQL on various
            Application
prequisite :-    Java programming


theory :-
        what is the defination of Hive? what is the
present version of hive and explain about
ACID transactions in Hive?
    Hive :      it is a datawarehouse, & database
            where data is typically loaded from batch
    performing for analytical purposes
*    the last version of Hive 3.1.2 released
    on August 26 ,2019
*    use cases that require transaction with ACID
    properties in have
①  Streaming data :- with ACID we can insert
        update ocr Same have pattern without
    affeching performance of table,

(ii) data counceting :- Found Sometimes data Stored may be micovert and Some changes are required which or can easely achev ed by ingend/update delete operations

(iii) Bulk updates using SQL merge statement- with bluk marge, Small dules can be manged into File without affecting read performance

2> what kind of data warehouse is suitable for hive what are types of tables in hive?

* hive is most suitable for data warehouse application because
  - Analya delatively static data
  - has less responsive table
  - does not make rapid changes in data

* there are two types of tables,
  - managed tables
  - external tables

3> is it possible to add 100 nodes when we have 100 nodes allready in hive?

yes we can add by following Steps

* talce new System, erectle usernare and pass
* install SSH and makes node setup SSH connetions
* Add Ssh public id key to authougre key Files
* Add new datanode hostname, ipadres, and other log in to new node
* Start HDFS of newly added Slave and oulpt the sps comd

4> what is metastore in hive? what is difference between local and remote metastore?

metastore is central repository of hive meta data

local Node :- the hive metastore reservive are in the same process but make hive metastore database runs in seperate process as it can be on seprate host in local mode

Remote mode :-
hive metadastore a secources are in its own Jum pwass

conclusion :-
Studied and implemented Hbase and hiveal on various applications.

Part B :- Assignment 1

Title :- Perform the following operations using
        python on Facebook metrics data set a
        a) Create data Subset ⓑ merge data ⓒ Sort data
        ⓓ transposing data ⓔ shape and reshape data

Objective :- To use python concept

Prerequisite :- python programming.

theory :-
        what are key features of python?
    • python is integrated languge it does not need to
        compiled before run
    • python is dynamically typed we don't need to
        state type of variable
    • it object orieanted programming languge
    • writing python code is quick

2) what type of language is python?
            python is capable of scripting but it is
    Considered as general purpose programming
    language

3) what is the name space in python?
            A namespace refers to name which
    is assigned to each object in python the name
    -spaces are maintaned like a dictionary where
    key is namespace and value is adress

4) What are local and global variable in python?

local variable :- any variable inside a function is called local variable

Global variable :-

Variable declared outside a function in global space this can be accessed by any function in programe

5) whate are benifits of using python
- easy do use
- Free and open source
- it can run on any platform
- It provides no
- Python has large library Support
- the data Structure used in python are user frraindly

Conclusion :-

Studied and implemented concepts

Part B :- Assignment No. 2

Title :- perform the following operations using python on the Air quality and heart diseases data sets

    ⓐ data cleaning    ⓑ data integration
    ⓒ data transformation    ⓓ Error correcting
    ⓔ data model building.

Prerequisite :- python programming

theory :-
    what type of language python ? programming or Scripting.
    python is general perpose programming language but Scripting is also possible

2) python an interpreated language Explain
    An interpreted language is any programming language which is not in machine level code before Structure
    python is intrepreted language.

3) what are common build in data types in python
- Numbers :- they include int, flot and complex
- list : an ordered sequence of iteam
- tuple :- an ordered immutable sequence
- string : sequence of character
- Set :- collection of unique items
- dictornory :- Story value in key and value pair
- Boolean :- true / false

## Conclusion :-

Studied and Implemented python concepts

PART B :- Assignment No. 3

Title :- integrate python and Hadoop and perform the following operations on Forest fire dataset (a) data analysis using the map Reduce in pyHadoop (b) Data mining in Hive

Objective :- To use python Concept

Prerequisite :- Python programming

theory :-
① Explain what is Name Node in Hadoop ?
NameNode works as master in Hadoop cluster the main function are
✱ states metadata of actual data
✱ manages File system namespaces
✱ regulate client access request for actual file data

2> what is Jobtracker in hadoop & what are active followed by hadoop ?
In hadoop for submmiting and tracking map reduce roles, Jobtracker is used Jobtracker run on its own Jvm process
Following actions perfound by Hadoop:-
• Client application submit Job to Jobtracker
• Jobtracker communccate to NameNode to determine data location
• Near location Jobtracker locates track trolcer nodes
• an chosen node it submit work

3) mention what are most common input format defined by Hadoop?

the most common input formats is Hadoop are

- Sequencefile input Format
- text input format
* keyvalue input format

Conclusion :-

Studed and applied concepts of python

Part B :- Assignment 4

Title :- visualize the data using python libraries matplotlib, seaborn by plotting the graphs for assignment no 2 and 3 (Group B)

Objective :- to use python concepts

theory :-

① How do create multiple subplots in matplotlib in python ?

to create multiple plots use matplotlib pyplot subplot method which structures the figure along with axes object or array of axes object n row in code attribute of subplot(s) method determine the m. of source end columns subplot grid

.by default it returns figure with single subplot

How it works :-

① when we call the subplots() method by stating only in on direction it returns array of axes objects i·e subplots
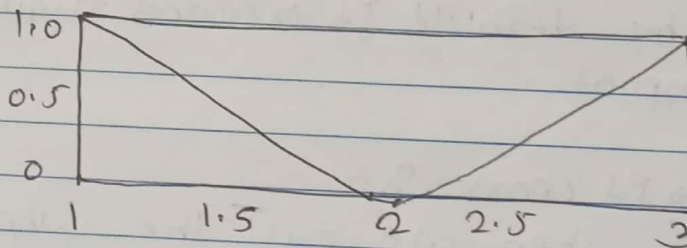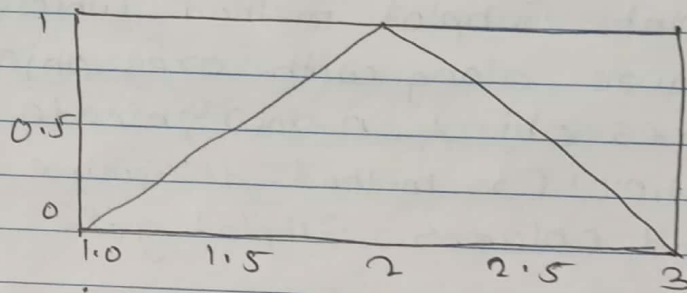
② we can access these axes object wing ndicess do create specific subplot call matplotlib·pyplot.plot() on corresponding index of axes

example:-

　　1-D array of subplot
　　import matplotlib.pyplot as pt
　　x = [1,2,3]
　　y = [0,1,0]
　　z = [1,0,1]
　　fig, ax = plt.subplot(2)
　　ax[0].plot(x,y)
　　ax[1].plot(x,y]
　　output

Conclusion:-
　　Studied and implemented concepts of python

Part B :- Assignment 5

Title :- perform the following data visualisation operations using tableau on Adult and Iris datasets @ 1D (linear) data visualization ⓑ 2D (planner) Data Visualization © 3D visualizati -on Ⓕ tree/hiarechical data Visualization ⓖ Network data Visualization

objective :- To use python Concept

theory :-
    What is TABLEAU
    TABLEAU is a visual analyties platform transform ating the way we use data to solve problems empouring people and organization to make most of them data

    What is data Visualization ?
        it is graphical representation of information and data by using visual elements like chart graphs and maps provides acceptble way do see undeastand trands outless patterens in data

    3> list out TABLEAU file extensions :-
        workbooks (.tood) - tubes workbook hold one move worksheet
        Bookwaks (.tbm) - it is easy way to quidxy share your worc

- package workbooks (twbx) :-

it is single zip file that contains a workbook along with any supporting file data.

- extract (.hyper or .tde) -

extract file are local way of subset or entire data set

4> what are file size limitations with tableau?

A site convers with 100 GB storage capacity workbooks published data source individual workbook pubished to your site can maximum size of 15GB

Conclusion :-

Studied and implemented concepts of python programing