| Sr.No. | Problem Statement |
|--------|-------------------|
| 1 | Design a distributed application using MapReduce which processes a log file of a system. List out the users who have logged for maximum period on the system. Use simple log file from the Internet and process it using a pseudo distribution mode on Hadoop platform. |
| 2 | Design and develop a distributed application to find frequency of words from sample text data. Use sample text data and process it using MapReduce. |
| 3 | Design a distributed application using MapReduce which processes Music dataset. List out the number of unique listeners and no of times the track was shared with others. Use music dataset and process it using a pseudo distribution mode on Hadoop platform. |
| 4 | Design a distributed application using MapReduce which processes Music dataset. List out the number of times the track was listened on Radio and no of times the track was skipped. Use music dataset and process it using a pseudo distribution mode on Hadoop platform. |
| 5 | Design a distributed application using MapReduce which processes Movie dataset. Recommend the Movie based on the user ratings. Use Movie dataset and process it using a pseudo distribution mode on Hadoop platform. |
| 6 | Write an application using HBase and HiveQL for flight information system which will include<br>a. Create Flight Info Hbase Table(with Flight information, schedule, and delay)<br>b. Demonstrate Creating, Dropping, and altering Database tables in Hbase<br>c. Creating an external Hive table to connect to the HBase for Flight Information Table<br>d. Find the total departure delay in Hive<br>e. Find the average departure delay in Hive<br>f. Create index on Flight information Table |
| 7 | Write an application using HBase and HiveQL for Customer information system which will include<br>a. Creation of –Cutomer_info(Cust-ID,Cust-Name,orderID), order_info(OrderID,ItemID,Quantity), item_info(Item-ID,Item-Name,ItemPrice) tables in Hive<br>b. Load table with data from local storage in Hive.<br>c. Perform Join tables with Hive<br>d. Create Index on Customer information system in Hive.<br>e. Find the total, average sales in Hive<br>f. Find Order details with maximum cost.<br>g. Creating an external Hive table to connect to the HBase for Customer Information System.<br>h. Display records of Customer Information Table in Hbase. |
| 8 | Write an application using HBase and HiveQL for OnlineRetail Dataset which will include<br>i. Create and Load table with Online Retail data in Hive.<br>j. Create Index on Online Retail Table in Hive.<br>k. Find the total, average sales in Hive<br>l. Find Order details with maximum cost.<br>m. Find Customer details with maximum order total. |

| | |
|---|---|
| | n. Find the Country with maximum and minimum sale.<br>o. Creating an external Hive table to connect to the HBase for OnlineRetail.<br>p. Display records of OnlineRetail Table in Hbase. |
| 9 | Perform the following operations using Python on the Facebook metrics data sets<br>a. Create data subsets for type of post<br>b. Merge two subsets<br>c. Sort Data on Page total likes<br>d. Transposing Data<br>e. Melting Data to long format<br>f. Casting data to wide format |
| 10 | Perform the following operations using Python on the Iris data sets<br>g. Create data subsets for different species<br>h. Merge two subsets<br>i. Sort Data Petal Length<br>j. Transposing Data<br>k. Melting Data to long format<br>l. Casting data to wide format |
| 11 | Perform the following operations using Python on Movie data sets<br>m. Create data subsets for different languages(Original Language).<br>n. Merge two subsets<br>o. Sort Data using customer ratings.<br>p. Transposing Data<br>q. Melting Data to long format<br>r. Casting data to wide format |
| 12 | Perform the following operations using Python on census bureau databset(Adult data sets).<br>s. Create data subsets for different Country, Sex, race.<br>t. Merge two subsets<br>u. Sort Data using customer ratings.<br>v. Transposing Data<br>w. Melting Data to long format<br>x. Casting data to wide format |
| 13 | Perform the following operations using Python on Heart Diseases data sets<br>a. Data cleaning(Remove NA, ?, Negative values etc.)<br>b. Error correcting(Outlier detection and removal)<br>c. Data transformation<br>d. Build Data model using regression and kNN methods and compare accuracy of heart disease prediction. |
| 14 | Perform the following operations using Python on Iris data sets<br>e. Data cleaning(Remove NA, ?, Negative values etc.)<br>f. Error correcting(Outlier detection and removal)<br>g. Data transformation |

| | |
|---|---|
| | h. Build Data model using regression and Naïve Bayes methods and compare accuracy of Iris Species Prediction. |
| 15 | Perform the following operations using Python on Breast Cancer data sets<br>i. Data cleaning(Remove NA, ?, Negative values etc.)<br>j. Error correcting(Outlier detection and removal)<br>k. Data transformation<br>l. Build Data model using regression and Naïve Bayes methods and compare accuracy of benign and malignant tumors in Breast Cancer Dataset. |
| 16 | Perform the following operations using Python on census bureau databset(Adult data sets)<br>m. Data cleaning(Remove NA, ?, Negative values etc.)<br>n. Error correcting(Outlier detection and removal)<br>o. Data transformation<br>p. Build Data model using regression and Naïve Bayes methods for prediction of income category (>=50k or <=50k) and compare accuracy Prediction. |
| 17 | Visualize the Heart disease dataset by plotting the following graphs using Python. (Define objective for every graph)<br>    a. Histograms<br>    b. Dot Plots<br>    c. Bar Plots<br>    d. Line Charts<br>    e. Add Histogram and Scatter plot to box plot. |
| 18 | Visualize the Heart disease dataset by plotting the following graphs using Python. (Define objective for every graph)<br>    a. Histograms<br>    b. Pie Charts<br>    c. Box Plots<br>    d. Scatter Plots<br>    e. `Add boxplots to a scatterplot` |
| 19 | Perform the data visualization operations using Tableau to get answers to various business questions on Retail dataset.<br>    a. Find and Plot top 10 products based on total sale<br>    b. Find and Plot product contribution to total sale<br>    c. Find and Plot the month wise sales in year 2010 in descending order<br>    d. Find and Plot most loyal customers based on purchase order<br>    e. Find and Plot yearly sales comparison<br>    f. Find and Plot country wise total sales price and show on Geospatial graph |
| 20 | Perform the data visualization operations using Tableau to get answers to various business questions on Retail dataset.<br>    a. Find and Plot country wise popular product<br>    b. Find and Plot bottom 10 products based on total sale<br>    c. Find and Plot top 5 purchase order<br>    d. Find and Plot most popular products based on sales<br>    e. Find and Plot half yearly sales for the year 2011 |

| | |
|---|---|
| | f. Find and Plot country wise total sales quantity and show on Geospatial graph |
| 21 | Visualize the census bureau databset(Adult data sets)by plotting the following graphs using Python. (Define objective for every graph)<br>  f. Histograms<br>  g. Dot Plots<br>  h. Bar Plots<br>  i. Line Charts<br>  j. Add Histogram and Scatter plot to box plot. |
| 22 | Visualize the census bureau databset(Adult data sets)by plotting the following graphs using Python. (Define objective for every graph)<br>  a. Histograms<br>  b. Pie Charts<br>  c. Box Plots<br>  d. Scatter Plots<br>  e. Add boxplots to a scatterplot |
| 23 | Perform the data visualization operations using Tableau to get answers to various questions on the census bureau databset(Adult data sets).<br>  a. Find and Plot Income class of People whose education is master's and doctorate.<br>  b. Find and Plot Income class of people who have private jobs.<br>  c. Find and Plot yearly sales comparison<br>  d. Find and Plot country wise statistics on Geospatial graph<br>  e. Plot agewise- education vs salary statistics.<br>  f. Plot Countrywise male female ratio.<br>  g. Plot Income class based on workclass(Government and other) |
| 24 | Perform the following operations using Python on ForestFires Dataset.<br>  a. Create data subsets by making classes for amount of region affected.(e.g. NotAffected, Partially affected, Mostlyaffected).<br>  b. Merge two subsets<br>  c. Sort Data using Temperature, wind and area.<br>  d. Transposing Data<br>  e. Melting Data to long format<br>  f. Casting data to wide format |
| 25 | Perform the following operations using Python on Hepatitis Dataset.<br>  a. Create data subsets for different sex.<br>  b. Merge two subsets<br>  c. Sort Data using age, SGOT, PROTIME.<br>  d. Transposing Data<br>  e. Melting Data to long format<br>  f. Casting data to wide format |

| 26 | Perform the following operations using Python on Hepatitis dataset. |
|---|---|
| | q. Data cleaning(Remove NA, ?, Negative values etc.) |
| | r. Error correcting(Outlier detection and removal) |
| | s. Data transformation |
| | t. Build Data model using regression and Naïve Bayes methods for prediction class DIE, LIVE and compare accuracy Prediction. |