HADOOP

Question 1: Data Ingestion:

 Create a directory in HDFS and transfer the banking dataset from the local system to the HDFS directory.

Answer:

Steps: 1. Access the Hadoop Environment

• 1. Set up the Hadoop environment:

We did this using/installing local hadoop as well as using Hadoop in a Docker environment.

1. Setting up hadoop in local environment:

- 1. Open a terminal on the system where Hadoop is installed. Or you can directly go to hadoop bin using the command line.
- 2. If you're running Hadoop in a Windows environment, you can use **start-all.cmd** to launch the necessary Hadoop services(NameNode, DataNode, ResourceManager, and NodeManager). Or you can use **start-dfs.cmd** and **start-yarn.cmd** to run the yarn demons.

C:\Windows\System32>cd/

C:\>cd hadoop/bin

C:\hadoop\bin>start-all.cmd

2. Create a Directory in HDFS:

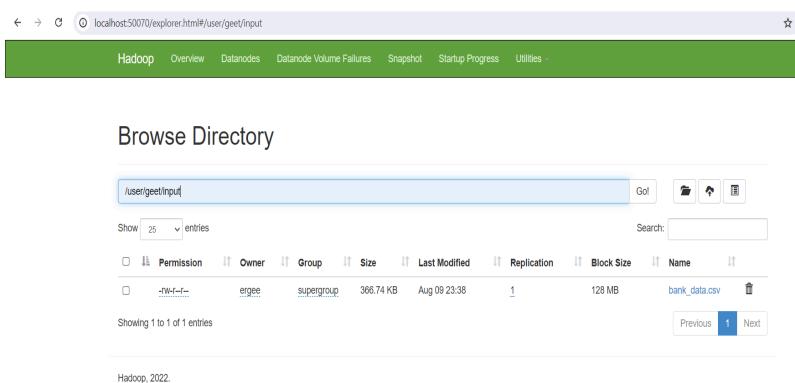
Use the **hadoop fs -mkdir** command to create a directory in HDFS.

C:\hadoop\bin>hadoop fs -mkdir -p /user/geet/input/

3. Transfer the Banking Dataset from Local to HDFS:
Use the hadoop fs -copyFromLocal command to transfer the bank.csv dataset from your local file system to the newly created HDFS directory.

C:\hadoop\bin>hadoop fs -copyFromLocal "C:\Users\ergee\Downloads\bank_data.csv" /user/geet/input/ C:\hadoop\bin>

4. Check if data is loaded to the new directory: We can visualise the data in our localhost by accessing localhost:**50070** in our browser.



2. Setting up Hadoop Using Hadoop docker container:

Steps to create a directory in HDFS and transfer the banking dataset from the local system to the HDFS directory using a Hadoop Docker container:

1. Setting up Hadoop using Docker:

- 1. We used the open source Framework available on git hub to set up Hadoop.
- 2. Navigate to the directory containing your docker file and run the following command to start the Hadoop cluster in detached mode (background):

docker-compose up -d

```
21:59
                              <DIR>
29-06-2020
08-08-2024
08-08-2024
29-06-2020
                                                6 .gitignore
base
                 00:58
                 21:59
21:59
                             <DIR>
                              <DIR>
                                                  datanode
                                          2,522 docker-compose-v3.yml
                 00:58
29-06-2020
29-06-2020
                00:58
00:58
                                          2,507 hadoop.env
08-08-2024
                             <DIR>
                                                  historyserver
                00:58
21:59
                                          1,437 Makefile
29-06-2020
08-08-2024
                              <DIR>
                                                  namenode
08-08-2024
                                                  nginx
08-08-2024
29-06-2020
08-08-2024
                 21:59
                              <DIR>
                                          nodemanager
2,171 README.md
                00:58
                21:59
21:59
                              <DIR>
                                               resourcemanager
                    :59 <DIR> submit
6 File(s) 10,202 bytes
10 Dir(s) 340,735,467,520 bytes free
08-08-2024
C:\docker-hadoop-master\docker-hadoop-master>docker compose up -d time="2024-08-09T17:45:25+05:30" level=warning msg="C:\\docker-hadoop-master\\docker-hadoop-master\\docker-compose.yml:
 version` is obsolete"

√Container resourcemanager Started

  Container nodemanager
Container historyserver
Container namenode
  Container datanode
C:\docker-hadoop-master\docker-hadoop-master>docker exec -it namenode bash
root@ddcb8ff54296:/#
```

Access the Container Shell:

 If you need to run HDFS commands, you can enter the running NameNode container using the following command:

docker exec -it namenode bash

2. Create an input folder: Inside the namenode container use hadoop fs -mkdir command to create a new directory.

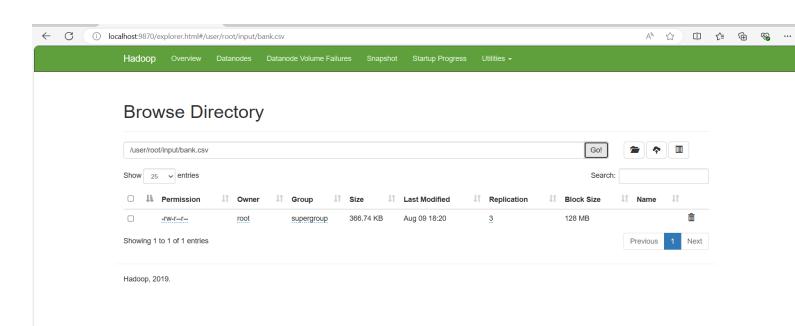
C:\docker-hadoop-master\docker-hadoop-master>docker exec -it namenode bash root@ddcb8ff54296:/# hdfs dfs -mkdir /user/root/input

- 3. Transfer the banking dataset from the local system to the HDFS directory:
 - 1. First move the bank.csv file from local to docker temp using **docker cp bank.csv namenode:/tmp** command.
 - Copy files from temp to the newly created input directory using hdfs dfs -put command.

C:\docker-hadoop-master\docker-hadoop-master>docker cp bank.csv namenode:/tmp
Successfully copied 377kB to namenode:/tmp

C:\docker-hadoop-master\docker-hadoop-master>docker exec -it namenode bash
root@ddcb8ff54296:/# hdfs dfs -mkdir /user/root/input
root@ddcb8ff54296:/# cd /tmp
root@ddcb8ff54296:/tmp# hdfs dfs -put bank.csv /user/root/input
2024-08-09 12:50:45,452 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false

4. Check if data is loaded to the new directory: We can visualise the data in our localhost by accessing to **localhost:9870** in our browser.



Question 2: Data Transformation with MapReduce:

Question 2.1 Write a MapReduce program in Python that calculates the average account balance for each job type.

Answer:

Step: 1. Write the MapReduce Python script and save as Mapper.py and Reducer.py.

Note: These two code files (Mapper.py and Reducer.py) are submitted with this document in separate files.

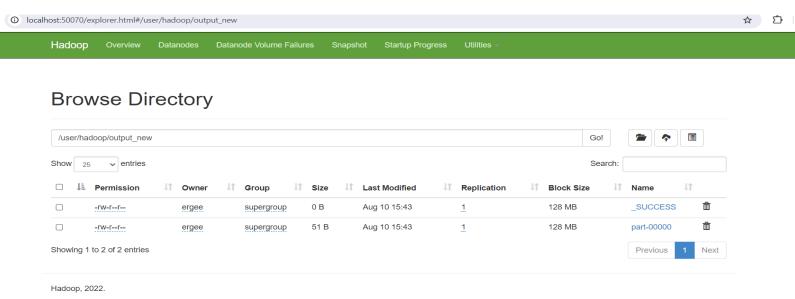
Step: 2. Create a new directory and upload the banking dataset bank_data.csv from local to HDFS.

Step: 3. Run Your MapReduce Job by specifying the input and output paths.



Step: 4. Retrieve the results from the output directory:

Use hadoop fs -cat /user/hadoop/output_new/part-00000 command to visualise and retrieve the output from the output directory.



Output:

```
C:\hadoop\bin>\C:\hadoop\bin>\hadoop fs -cat /user/hadoop/output_new/part-00000
admin. 1226.73640167364
blue-collar 1085.161733615222
entrepreneur 1645.125
housemaid 2083.8035714285716
management 1766.9287925696594
retired 2319.191304347826
self-employed 1392.4098360655737
services 1103.9568345323742
student 1543.8214285714287
technician 1330.99609375
unemployed 1089.421875
unknown 1501.7105263157894
```

Question 2.2. Write another MapReduce program that counts the number of individuals with and without a housing loan in each education category.

Answer:

Step: 1. Write the MapReduce Python script and save as Mapper.py and Reducer.py.

Note: These two code files are submitted with the document in separate files.

Step: 2. Run Your MapReduce Job by specifying the input and output paths.

Note: It is generally a good practice to delete the output directory before running a new MapReduce job in Hadoop to avoid overwriting issues.

*To delete the output directory use the command - hadoop fs -rm -r

C:\hadoop\bin>hadoop fs -rm -r /user/hadoop/output_new Deleted /user/hadoop/output_new

C:\hadoop\bin>hadoop fs -rm -r /user/hadoop/output_new Deleted /user/hadoop/output_new

C:\hadoop\bin>hadoop | gar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-2.10.2.jar -files "file:///C:/hadoop/mapper.py,file:///C:/hadoop/reducer.py" -mapper "python mapper.py" -reducer "python reducer.py" -jout /user/geet/input/bank_data.csv\ -output /user/hadoop/output_new

Output:

Question 2.3 - Perform a MapReduce job to determine the number of clients contacted in each month and their subscription status to term deposits ('y' column).

Answer:

Step 1 Delete the output directory before running a new MapReduce job.

Step: 2. Write the MapReduce Python script and save as Mapper.py and Reducer.py.

Note: These two code files are submitted with the document in separate files.

Step: 3. Run Your MapReduce Job by specifying the input and output paths.

```
C:\hadoop\bin>hadoop fs -rm -r /user/hadoop/output_new
Deleted /user/hadoop/output_new
```

C:\hadoop\bin>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-2.10.2.jar -files "file:///C:/hadoop/mapper.py,file:///C:/hadoop/reducer.py" -mapper "python mapper.py" -reducer "python reducer.py" -i put /user/geet/input/bank data.csv\ -output /user/hadoop/output_new

```
C:\hadoop\bin>
C:\hadoop\bin>hadoop fs -cat
        56
                  236
apr
aug
        79
                  553
dec
        8
                  11
feb
        38
                  183
jan
        16
                  131
jul
        61
                  644
        55
                 475
jun
        20
                  28
mar
        93
                  1304
may
         39
                  349
nov
        37
                  42
oct
                  34
        17
sep
C:\hadoop\bin>
```

3. Data Analysis with MapReduce:

Question 3.1. Analyse the average duration of contact (in seconds) per campaign outcome ('poutcome').

Answer:

Steps:

1. Write the MapReduce script and save as mapper.py and reducer.py

Note: These two code files are submitted with the document in separate files.

- 2. Upload the banking data to HDFS.
- 3. Run the MapReduce job on your Hadoop cluster.

```
C:\hadoop\bin>hadoop fs -rm -r /user/hadoop/output_new
Deleted /user/hadoop/output new
```

C:\hadoop\bin>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-2.10.2.jar -files "file:///C:/hadoop/mapper.py,file:///C:/hadoop/reducer.py" -mapper "python mapper.py" -reducer "python reducer.py" -in put /user/geet/input/bank_data.csv\ -output /user/hadoop/output_new

- 4. Check and interpret the results.
- 5. Clean up the output directory if you need to rerun the job.

Output:

```
C:\hadoop\bin>hadoop fs -cat /user/hadoop/output_new/part-00000
failure 254.38
other 273.83
success 338.64
unknown 262.10
C:\hadoop\bin>_
```

Summary:

The MapReduce job analyzed the dataset to calculate the average contact duration (in seconds) for each campaign outcome. The results are as follows:

- Failure: Average contact duration is 254.38 seconds.
- Other: Average contact duration is 273.83 seconds.
- Success: Average contact duration is 338.64 seconds.
- Unknown: Average contact duration is 262.10 seconds.

These results show that successful campaign outcomes are associated with the longest average contact duration, while failed outcomes have shorter average durations.

Question 3.2 - Examine the relationship between the age of clients and their balance, and present findings in a summarised form.

Answer:

Summary:

- Purpose: This analysis directly examines how the balance varies with each specific age.
- Implementation: The MapReduce job processes the data by outputting the age and balance in the mapper and then calculates the average balance for each age in the reducer.
- Results: The results will give a detailed view of how balance varies with specific ages, providing a more granular insight into the relationship between age and balance.

This approach avoids grouping by age bins and instead provides a direct average balance for each individual age.

Steps:

1. Write the MapReduce script and save as mapper.py and reducer.py

Note: These two code files are submitted with the document in separate files.

- 2. Upload the banking data to HDFS.
- 3. Run the MapReduce job on your Hadoop cluster.

C:\hadoop\bin>hadoop fs -rm -r /user/hadoop/output_new Deleted /user/hadoop/output_new

C:\hadoop\bin>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-2.10.2.jar -files "file:///C:/hadoop/mapper.py,file:///C:/hadoop/reducer.py" -mapper "python mapper.py" -reducer "python reducer.py" -in put /user/geet/input/bank_data.csv\ -output /user/hadoop/output_new

- 4. Check and interpret the results.
- 5. Clean up the output directory if you need to rerun the job.

Output:

```
C:\hadoop\bin>
C:\hadoop\bin>hadoop fs -ca
19
        393.50
20
        661.33
21
        1774.29
22
        1455.33
23
        2117.95
24
        634.62
25
        1240.07
26
        788.56
27
        851.78
28
        1025.10
29
        1261.88
30
        1113.03
        1288.48
        1256.55
33
        1545.41
34
        1111.54
        1192.83
36
        1226.89
37
        1463.92
38
        1718.99
39
        1104.86
40
        1399.51
41
        1505.79
42
        1612.36
43
        1807.83
44
        1836.55
45
        1187.37
46
        998.77
47
        1363.05
48
        1462.36
49
        1591.11
50
        1645.06
51
        1528.57
52
        782.29
53
        1588.31
54
        1656.66
55
        1244.94
        2120.14
```

```
56
         2120.14
57
         1665.63
         1755.08
58
         1582.48
        2964.57
61
         2407.50
62
        516.14
63
        2286.38
        1103.29
64
        1638.17
         3313.89
        4149.40
67
68
        11753.00
69
         774.33
70
        5084.57
         3787.33
        2526.00
        525.83
74
        1978.33
75
        7046.50
,
76
        1338.00
        2405.17
78
         318.00
79
        4087.75
80
        4183.50
81
        1.00
         380.50
83
         639.00
86
         1503.00
87
         230.00
C:\hadoop\bin>_
```

Summary of Findings:

After running the MapReduce job, the output provides the average account balance for each specific age. Here's a summary of the key points:

- **Age-Specific Averages**: The output shows the average balance corresponding to each client's age. For example, a 23-year-old might have an average balance of 2117.95, while a 25-year-old might have an average balance of 1240.05.
- Trends: If observed across ages, you might notice trends such as:
 - Increase with Age: In some cases, there might be a gradual increase in average balance as age increases.

- Fluctuations: Certain ages might show higher or lower average balances due to specific financial behaviours or life events.
- **Variability:** The average balances might fluctuate significantly across different ages, reflecting the diverse financial situations of clients at various life stages.

Conclusion:

The analysis reveals the average balance associated with each age, helping to identify any patterns or anomalies in financial behaviour across different age groups. This granular insight is valuable for financial institutions to tailor their products and services according to the needs of different age demographics.

HIVE

1. Data Ingestion and Table Creation:

Question 1.1 Create a Hive database named banking_data.

Answer:

Steps to Create a Hive Database:

- 1. Start Hive CLI
- 2. Create a new database named banking_data.

```
Beeline version 2.1.0 by Apache Hive
hive> CREATE DATABASE banking_data;
OK
No rows affected (0.791 seconds)
hive> _
```

3. Switch to the newly created database:

```
OK
No rows affected (0.791 seconds)
hive> USE banking_data;
OK
```

Question 1.2. Define and create a Hive table **client_info** with appropriate data types for the **bank.csv** dataset.

Answer:

• Create the client_info table with appropriate data types based on the columns in the bank.csv file:

```
No rows affected (0.033 seconds)
hive> CREATE TABLE client_info (
 . > age INT,
 . > job STRING,
 . > marital STRING,
 . > education STRING,
 . > default STRING,
 . > balance INT,
 . > housing STRING,
 . > loan STRING,
 . > contact STRING,
 . > day INT,
 . > month STRING,
 . > duration INT,
 . > campaign INT,
 . > pdays INT,
 . > previous INT,
 . > poutcome STRING,
 . > y STRING
 . > )
 . > ROW FORMAT DELIMITED
 . > FIELDS TERMINATED BY ','
 . > STORED AS TEXTFILE;
DΚ
No rows affected (0.492 seconds)
hive> _
```

Question 1.3. Load the data from the bank.csv file into the client_info table.

Answer:

1. Load the banking data from hadoop directory to the client info table.

```
name. (state=42000,code=40000)
hive> LOAD DATA INPATH '/user/geet/input/bank_data.csv' INTO TABLE client_info;
Loading data to table banking_data.client_info
OK
No rows affected (0.67 seconds)
```

2. Query the client_info table to verify that the data has been loaded correctly.

```
No rows affected (0.67 seconds)
hive> SELECT*FROM client_info LIMIT 10;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
```

Output:

```
job marital education default housing loan contact month poutcome y
30 unemployed married primary no 1787 no no cellular 19 oct 79 1 -1 0 unknown no
33 services married secondary no 4789 yes yes cellular 11 may 220 1 339 4 failure no
35 management single tertiary no 1350 yes no cellular 16 apr 185 1 330 1 failure no
30 management married tertiary no 1476 yes yes unknown 3 jun 199 4 -1 0 unknown no
59 blue-collar married secondary no 0 yes no unknown 5 may 226 1 -1 0 unknown no
35 management single tertiary no 747 no no cellular 23 feb 141 2 176 3 failure no
36 self-employed married tertiary no 307 yes no cellular 14 may 341 1 330 2 other no
39 technician married secondary no 147 yes no cellular 6 may 151 2 -1 0 unknown no
41 entrepreneur married tertiary no 221 yes no unknown 14 may 57 2 -1 0 unknown no
10 rows selected (0.877 seconds)
hive>
```

Here, we can see the asked ten rows in the output. So, it is verified that the data is loaded successfully.

Question 2: Basic Data Exploration:

Question 2.1. Write a HiveQL query to count the total number of clients in the dataset.

Answer:

Steps:

- 1. **FROM client_info**: The query selects data from the client_info table, which contains the records of all clients.
- 2. **COUNT(*) AS total_clients**: The COUNT(*) function counts the total number of rows in the client_info table. The result is labelled as total_clients for clarity.

Query:

```
hive> select count(*) AS total_clients FROM client_info;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
```

```
4522
1 row selected (2.422 seconds)
hive>
```

Summary of the Results:

• **Total Number of Clients**: The query returns a single number representing the total number of clients in the dataset. This number gives you a quick overview of the dataset size, indicating how many client records are available for analysis. So, here we can see that the total number of clients is 4522.

Question: 2.2 Display the first 10 rows of the dataset.

Answer:

Steps:

- 1. **FROM client_info**: The query selects data from the client_info table, which contains all the records in the dataset.
- **SELECT ***: The query selects all columns (*) from the client_info table. This
 means that every piece of information available for each client will be included in the
 result.
- 3. **LIMIT 10**: The query restricts the output to only the first 10 rows of the dataset. This is useful for quickly examining a sample of the data without retrieving the entire dataset.

Query:

```
No rows affected (0.67 seconds)
hive> SELECT*FROM client_info LIMIT 10;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
```

```
job marital education default housing loan contact month poutcome y
30 unemployed married primary no 1787 no no cellular 19 oct 79 1 -1 0 unknown no
33 services married secondary no 4789 yes yes cellular 11 may 220 1 339 4 failure no
35 management single tertiary no 1350 yes no cellular 16 apr 185 1 330 1 failure no
30 management married tertiary no 1476 yes yes unknown 3 jun 199 4 -1 0 unknown no
59 blue-collar married secondary no 0 yes no unknown 5 may 226 1 -1 0 unknown no
35 management single tertiary no 747 no no cellular 23 feb 141 2 176 3 failure no
36 self-employed married tertiary no 307 yes no cellular 14 may 341 1 330 2 other no
39 technician married secondary no 147 yes no cellular 6 may 151 2 -1 0 unknown no
41 entrepreneur married tertiary no 221 yes no unknown 14 may 57 2 -1 0 unknown no
10 rows selected (0.877 seconds)
hive>
```

Summary of the Results:

First 10 Rows of the Dataset: Here, we can see that the output displays all columns
and their values for the first 10 clients in the client_info table. These rows
represent a small sample of the overall dataset, providing a snapshot of the data
structure and contents.

Question 3. Data Filtering and Sorting:

Question 3.1 Retrieve all records of clients who are married and have a personal loan.

Answer:

Steps:

- 1. **FROM client_info**: The query selects data from the client_info table, which includes various details about the clients such as their marital status, loan status, and other attributes.
- 2. **SELECT ***: The query selects all columns (*) from the client_info table. This means that every piece of information available for each client will be included in the result, such as age, job, balance, etc.
- 3. **WHERE marital = 'married' AND loan = 'yes'**: The query filters the data to include only those clients who meet both of the following conditions:
- marital = 'married': The client must be married.
- **loan = 'yes'**: The client must have a personal loan.

Query:

```
1 row selected (2.422 seconds)
hive> select * from client_info where marital ='married'AND loan ='yes';
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
        [-appendToFile <localsrc> ... <dst>]
```

hadoop fs [generic options] -chgrp [-R] GROUP PATH.

management married tertiary no 176 yes yes unknown 1 jun 8 1 687 2 failure no unemployed married primary no 180 yes yes unknown 27 may 62 2 -1 0 unknown no entrepreneur married tertiary no 37 yes yes cellular 6 may 207 1 363 2 failure no technician married secondary no 410 no yes cellular 21 nov 162 5 -1 0 unknown no 45 unemployed married primary no 180 yes yes unknown 27 may 62 2 - 1 0 unknown no 66 entrepreneur married tertiary no 37 yes yes cellular 6 may 207 1 363 2 failure no 46 technician married secondary no 410 no yes cellular 21 nov 162 5 - 1 0 unknown no 30 technician married primary no 177 yes yes telephone 21 jul 742 4 - 1 0 unknown no 30 technician married secondary no -163 yes yes cellular 16 jul 435 1 - 1 0 unknown no 29 blue-collar married secondary no -163 yes yes cellular 8 aug 502 2 - 1 0 unknown no 30 services married secondary no 227 no yes unknown 12 may 241 2 - 1 0 unknown no 30 services married secondary no 280 no yes cellular 9 jul 170 1 - 1 0 unknown no 56 self-employed married secondary no 800 no yes cellular 9 jul 170 1 - 1 0 unknown no 31 services married secondary no 270 yes yes unknown 12 jun 122 1 - 1 0 unknown no 32 services married secondary no 270 yes yes unknown 20 may 201 1 - 1 0 unknown no 36 unemployed married secondary no -872 yes yes cellular 20 nov 153 1 183 1 failure no 38 services married secondary no -872 yes yes cellular 20 nov 153 1 183 1 failure no 38 services married secondary no 9 no yes cellular 16 jul 1473 6 - 1 0 unknown no 37 management married tertiary no 997 yes yes cellular 20 nov 153 1 183 1 failure no 38 management married tertiary no 997 yes yes cellular 25 jul 160 2 - 1 0 unknown no 32 services married secondary no 8 no yes cellular 15 jul 1624 6 - 1 0 unknown no 32 services married secondary no 8 no yes cellular 25 jul 132 2 - 10 unknown no 32 services married secondary no 8 no yes cellular 15 jul 624 6 - 1 0 unknown no 32 services married secondary no 3382 no yes cellular 25 jul 132 2 - 10 unknown no 38 cervices married secondary no 3387 no yes cellular 15 jul 624 6 - 1 0 unknown no 50 technician married tertiary no 3387 yes yes telephone 31 jul 74 1 1 4 - 1 0 unknown no 51 technician married secondary no 80 yes yes cellular 25 jul 123 2 - 1 0 unknown no 52 damin married secondary no 80 yes yes unknown 4 jun 622 1 - 1 0 unknown no 53 tous services married second

Summary of result:

- This query filters the data to retrieve records of clients who are married and have a personal loan.
- **Output**: We can see the output is a list of all married clients who have taken a personal loan, including all columns of data.

Question 3.2. List the top 10 clients with the highest balance, displaying their job, marital status, and balance.

Answer:

Steps:

- 1. **FROM client_info**: The query selects data from the client_info table, which contains information about clients, including their job, marital status, and account balance.
- 2. **SELECT job, marital, balance**: The query specifies that it wants to retrieve the job, marital, and balance columns from the client_info table.
- 3. **ORDER BY balance DESC**: The query orders the results by the balance column in descending order (DESC), meaning that the clients with the highest balances will appear first.
- 4. **LIMIT 10**: The query limits the results to the top 10 records. This means only the 10 clients with the highest balances will be shown.

Query:

hive> select job,marital,balance from client_info order by balance desc limit 10;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group

Result:

```
retired married 71188
entrepreneur married 42045
technician single 27733
management married 27359
technician married 27069
housemaid single 26965
retired married 26452
services married 26394
management divorced 26306
retired single 25824
10 rows selected (1.677 seconds)
hive>
```

Summary of the Results:

Top 10 Clients by Balance: The output of this query will display the job, marital status, and balance of the 10 clients who have the highest account balances. So, here we see that **retired clients** are generally married and have the highest balance.

Question 4.: Data Aggregation and Grouping:

Question 4.1 Calculate the average age of clients for each job category.

Steps:

The query retrieves data from the client_info table, which contains information about clients, including their job and age.

GROUP BY job: The query groups the data by the job column. This means that the data will be aggregated separately for each unique job category.

AVG(age) AS average_age: For each job category, the query calculates the average age of the clients using the AVG(age) function. The result is stored in a column named average_age.

SELECT job, AVG(age) AS average_age: Finally, the query selects and displays the job category (job) alongside the calculated average age (average_age) for each group.

Query:

hive> select job, AVG(age) AS average_age from client_info group by job;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group

Result:

```
WARNING: Hive-on-MR is deprecated in Hive 2 and may
admin. 39.68200836820084
blue-collar 40.15644820295983
entrepreneur 42.01190476190476
housemaid 47.339285714285715
iob
management 40.54076367389061
retired 61.869565217391305
self-employed 41.45355191256831
services 38.57074340527578
student 26.821428571428573
technician 39.470052083333336
unemployed 40.90625
unknown 48.10526315789474
13 rows selected (1.643 seconds)
hive>
```

Summary of the Results:

 Average Age by Job Category: The query outputs the average age of clients for each job category. This provides insight into the typical age of clients in different professions. Here, we can see that the average age of the majority of clients for different job categories is between 35 to 45.

Question. 4.2. Find the total number of clients for each education level who have defaulted on credit.

Answer:

Steps:

FROM client_info: The query starts by selecting data from the client_info table, which contains information about the clients, including their education level and whether they have defaulted on credit.

WHERE default = 'yes': The query filters the data to include only those clients who have defaulted on credit. The default column is checked, and only records where default = 'yes' are selected. This ensures that the query is only counting clients who have actually defaulted.

GROUP BY education and default: The query then groups the filtered data by the education column and default column.

COUNT(*) AS total_defaulted_clients: For each education level, the query counts the number of clients who have defaulted on credit using COUNT(*). The result is stored in a column named total_defaulted_clients.

SELECT education, default, COUNT(*) AS total_defaulted_clients: Finally, the query selects and displays the education level alongside the total number of clients who have defaulted in that education level.

Query:

```
13 rows selected (1.643 seconds)
```

hive> select education, default, count(*) AS total_defaulted_client from client_info where default = 'yes' group by education,default; -chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group

Usage: hadoop fs [generic options]

```
WARNING: Hive-on-MR is deprecated primary yes 10 secondary yes 46 tertiary yes 17 unknown yes 3 4 rows selected (1.726 seconds) hive>
```

Summary of result:

 This query finds the total number of clients for each education level who have defaulted on credit. Here, we can see that the majority of clients who defaulted have a secondary education level.

Question: 5 Complex Queries for Insights:

Question 5.1. Identify the top 5 job categories with the highest average balance and the percentage of clients in each of these job categories who have subscribed to a term deposit.

Answer:

Step 1: Inner Query - Identifying Top 5 Job Categories by Average Balance

- 1. **FROM client_info**: The query selects data from the client_info table.
- 2. **GROUP BY job**: It groups the data by the job column, meaning it will aggregate data for each job category.
- 3. **AVG(balance) AS avg_balance**: For each job category, the query calculates the average balance (AVG(balance)) and stores it as avg_balance.
- 4. **COUNT(*) AS total_clients**: It counts the total number of clients in each job category using COUNT(*) and stores this value as total_clients.
- 5. SUM(CASE WHEN y = 'yes' THEN 1 ELSE 0 END) AS subscribed_clients: It counts the number of clients who subscribed to a term deposit (y = 'yes') within each job category. This is done using a CASE statement, where 1 is added for each subscription and 0 for non-subscriptions. The result is stored as subscribed_clients.
- 6. **ORDER BY avg_balance DESC**: The job categories are ordered by their average balance in descending order, so the job categories with the highest average balance come first.
- 7. **LIMIT 5**: The query limits the result to the top 5 job categories with the highest average balance.

Step 2: Outer Query - Calculating Subscription Percentage

- 1. **FROM (...) sc**: The outer query selects data from the result of the inner query, which is aliased as sc.
- 2. **SELECT sc.job, sc.avg_balance**: The outer query directly selects the job and avg_balance columns from the inner query's results.
- 3. (sc.subscribed_clients / sc.total_clients) * 100 AS subscription_percentage: It calculates the subscription percentage by dividing the number of subscribed clients (subscribed_clients) by the total number of clients (total_clients) for each job category, then multiplying by 100 to express it as a percentage. This value is stored as subscription_percentage.

Query:

```
hive> select sc.job, sc.avg_balance, (sc.subscribed_clients / sc.total_clients) * 100 as subscription_percentage from (select job, AVG(balance) as avg_balance, count(*) as total_clients, sum(case when y ='yes' then 1 else 0 end) as subscribed_clients from client_info group by job order by avg_balance desc limit 5) sc;

-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group

Usage: hadoop fs [generic options]

[-appendToFile <localsrc> ... <dst>]

[-cat [.ignore(rc] <src> ]

[-cat [.ignore(rc] <src> ]
```

Output:

```
retired 2319.191304347826 23.47826086956522
housemaid 2083.8035714285716 12.5
management 1766.9287925696594 13.519091847265221
entrepreneur 1645.125 8.928571428571429
student 1543.8214285714287 22.61904761904762
5 rows selected (2.922 seconds)
hive>
```

Summary of the Results:

- Top 5 Job Categories by Average Balance: The query identifies the five job
 categories with the highest average balances. These are the clients with the most
 significant average account balances across all job types. So, here retired and
 housemaids have the highest average balance.
- Subscription Percentage: The query then calculates how successful the current campaign was in converting clients in these high-balance job categories into term deposit subscribers. The subscription_percentage indicates the effectiveness of the campaign for each job category. E.g. Here it was not very successful for entrepreneurs because the subscription_rate was only 8.928 percent.

Question 5.2. Determine the month with the highest number of contacts and the success rate of the campaign in that month (percentage of clients who subscribed to a term deposit).

Α	n	C	۱۸	Δ	r.
$\overline{}$		J	٧v	·	ι.

Steps:

Inner query:

- 1. **FROM client_info**: The query starts by selecting data from the client_info table.
- 2. **GROUP BY month**: It groups the data by the month column to calculate statistics for each month.
- 3. **COUNT(*) AS total_contacts**: For each month, it counts the total number of records (or contacts) using COUNT(*). This value is stored as total_contacts.
- 4. **SUM(CASE WHEN y = 'yes' THEN 1 ELSE 0 END) AS successful_contacts**: For each month, the query counts the number of records where y = 'yes' (i.e., the client subscribed to a term deposit). This is done using a CASE statement that returns 1 for successful contacts and 0 otherwise. The sum of these values gives the number of successful contacts, stored as successful_contacts.
- 5. **ORDER BY total_contacts DESC**: The query orders the results by total_contacts in descending order, meaning the month with the highest number of contacts will appear first.
- 6. **LIMIT 1**: Finally, it limits the results to just the top month (the month with the highest number of contacts).

Outer Query:

- 1. FROM (...) AS top_month: The outer query selects data from the result of the inner query, which has been aliased as top_month.
- 2. SELECT month, total_contacts: The outer query selects the month and total_contacts fields directly from the result of the inner query.
- 3. (successful_contacts / total_contacts) * 100 AS success_rate: It calculates the success_rate by dividing the number of successful_contacts by total_contacts and then multiplying by 100 to get a percentage.

Query:

Output:

```
WARNING: Hive-on-MR is deprecated i
may 1398 6.652360515021459
1 row selected (2.877 seconds)
hive>
```

Summary:

- This query determines the month with the highest number of client contacts and calculates the success rate of the campaign in that month (the percentage of clients who subscribed to a term deposit).
- **Output**: May, is the month with the highest contacts, the total number of contacts are 1398, and the success rate (percentage of clients who subscribed to a term deposit) during that month is 6.652.

Question .6. Correlation Analysis:

Question .6. Calculate the correlation between age and balance for the clients.

Answer:

The CORR(age, balance) function calculates the Pearson correlation coefficient between the age and balance columns.

This coefficient will range from -1 to 1, where:

- 1 indicates a perfect positive correlation,
- -1 indicates a perfect negative correlation, and
- **0** indicates no correlation.

Query:

```
1 row selected (2.877 seconds)
hive> select CORR(age, balance) as age_balance_correlation
. . > from client_info;
```

```
0.08382014224477739
1 row selected (1.585 seconds)
hive>
```

Summary Of Result:

This result gives you an idea of how strongly the age of clients is related to their account balance. Here, Since the output is 0.0838 i.e. between 0 and 1 so, we can say that age is slightly related to balance.

Question. 7. Trend Analysis:

Analyse the year-over-year trend in the number of clients contacted:

*There is no data in the bank_data.csv dataset which represents the year. But let's say the first four characters from the month column represent the year (e.g., 2023 from 2023-Jan).

Answer:

Steps:

- 1. Extract Year: The SUBSTRING(month, 1, 4) function extracts the first four characters from the month column, assuming they represent the year (e.g., 2023 from 2023–Jan).
- 2. Count Contacts by Year: We group the data by the extracted year and count the number of clients contacted in each year.

Query:

```
1 row selected (1.585 seconds)
hive> SELECT SUBSTRING(month,1,4) AS year,
. . > COUNT(*) AS num_clients_contacted
. . > FROM
. . > CLIENT_INFO
. . > GROUP BY SUBSTRING(month,1,4)
. . > ORDER BY year;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
```

```
apr 293
aug 633
dec 20
feb 222
jan 148
jul 706
jun 531
mar 49
may 1398
mont 1
nov 389
oct 80
sep 52
13 rows selected
```

Summary of result: This output will give us a clear view of the year-over-year trend in the number of client contacts. Here, we did not have year information in the bank_data.csv file. So, we could not get the desired result.

Question.8.Anomaly Detection:

Identify any unusual patterns in the average yearly balance across different education levels.

Answer:

Inner Query (Subquery):

- SUBSTRING(month, 1, 4) AS year: Extracts the year from the month column, assuming month is in a format that includes the year (e.g., 2024-08).
- **education**: Selects the education column from the client_info table.
- AVG(balance) AS avg_yearly_balance: Calculates the average balance for clients grouped by year and education.
- AVG(AVG(balance)) OVER (PARTITION BY SUBSTRING(month, 1, 4))
 AS overall_avg_balance: Calculates the overall average balance for each year across all education levels using a window function.
- STDDEV(AVG(balance)) OVER (PARTITION BY SUBSTRING(month, 1, 4)) AS stddev_balance: Calculates the standard deviation of the average balance for each year across all education levels using a window function.
- **GROUP BY SUBSTRING(month, 1, 4), education**: Groups the data by year and education to calculate the metrics above.

Outer Query:

(avg_yearly_balance - overall_avg_balance) / stddev_balance AS
 z_score: Calculates the z-score for the average yearly balance. The z-score
 measures how many standard deviations a data point is from the mean. In this
 context, it tells us how much the average balance for a specific year and education
 level deviates from the overall average balance for that year.

Query:

```
hive> select year, education, (avg_yearly_balance - overall_avg_balance) / stddev_balance AS z_score
. . > FROM (
. . > SELECT
. . > SUBSTRING(month, 1, 4) AS year,
. . > education,
. . > AVG(balance) AS avg_yearly_balance,
. . > AVG(balance) OVER (PARTITION BY SUBSTRING(month, 1, 4)) AS overall_avg_balance,
. . > STDDEV(AVG(balance)) OVER (PARTITION BY SUBSTRING(month, 1, 4)) AS stddev_balance
. . > FROM client_info
. . > GROUP BY SUBSTRING(month, 1, 4), education
. . > ) AS subquery;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
```

```
OK
WARNING: Hive-on-MR is deprecated in Hi
apr primary 1.267794378462834
apr secondary -0.7235770026264813
apr secondary -0.7235770026264813
apr tertiary 0.6555354947097891
apr unknown -1.1997528705461413
aug primary -0.8113812295586758
aug secondary -0.5775732890112949
aug tertiary -0.31633238372712924
aug unknown 1.7052869022970998
dec primary -0.8088635843923103
dec secondary -0.397568378390461
dec tertiary -0.5058241544774784
dec unknown 1.7122561172602497
feb primary -1.0701923004119014
feb secondary -0.9216091477916096
feb tertiary 0.8917649610011248
feb unknown 1.1000364872023862
jan primary 1.1404156384061253
jan secondary -0.4369029752604172
jan tertiary 0.7115176049750513
jan unknown -1.41503026812076
jul secondary -0.4258456259680814
jul tertiary 0.9014112560610348
jul unknown -1.4404947331772213
jun primary -0.15682779176634196
jun secondary -1.0365919570822748
jun tertiary 1.6428352994501707
jun unknown -0.44941555060155386
mar primary -1.0722876073529977
mar secondary -0.24628224810074777
mar tertiary 1.6391398006892093
mar unknown -0.3205697123287344
may primary -1.0048248602462826
may secondary -0.9457054758622359
may tertiary 0.6640708393343474
may primary -1.086497054758622359
may tertiary 0.6640708393343474
may unknown 1.2864595047741691
mont education
nov primary -0.8137802862217836
nov secondary -0.9457054758622359
may tertiary 1.5167439470393023
nov unknown -0.98157898775337383
oct primary 1.7097889755337383
oct primary -1.097889755337383
oct primary -0.8154396810172013
sep primary -1.3086820436917062
```

```
dec unknown 1.7122561172602497
feb primary -1.0701923004119014
feb secondary -0.9216091477916096
feb tertiary 0.8917649610011248
feb unknown 1.1000364872023862
jan primary 1.1404156384061253
jan secondary -0.4369029752604172
jan tertiary 0.7115176049750513
ian unknown -1.41503026812076
jul primary 0.9649291030842655
jul secondary -0.4258456259680814
jul tertiary 0.9014112560610348
jul unknown -1.4404947331772213
jun primary -0.15682779176634196
jun secondary -1.0365919570822748
jun tertiary 1.6428352994501707
iun unknown -0.44941555060155386
mar primary -1.0722876073529977
mar secondary -0.2462824810074777
mar tertiary 1.6391398006892093
mar unknown -0.3205697123287344
may primary -1.0048248682462826
may secondary -0.9457054758622359
may tertiary 0.6640708393343474
mav unknown 1.2864595047741691
mont education
nov primary -0.8137802862217836
nov secondary 0.2770762379753206
nov tertiary 1.5167439470393023
nov unknown -0.9800398987928378
oct primary 1.7097889755337383
oct secondary -0.37052787477325194
oct tertiary -0.5238214197432853
oct unknown -0.8154396810172013
sep primary -1.3086820436917062
sep_secondary -0.4707154583719365
sep tertiary 0.3984500713422833
sep unknown 1.38094743072136
49 rows selected (1.701 seconds)
hive> _
```

Summary of the Results:

- year: The year extracted from the month column.
- education: The education level of the clients.
- z_score: The z-score indicates how the average yearly balance for a specific education level compares to the overall average for that year. A positive z-score means the average balance is above the overall average, while a negative z-score indicates it is below the overall average.

Output Interpretation:

- Positive Z-Scores: Indicate that the clients with a specific education level have a higher-than-average balance compared to others in the same year.
- Negative Z-Scores: Indicate that the clients with a specific education level have a lower-than-average balance compared to others in the same year.

Question . 9. Advanced Analysis:

Question 9.1. Analyze the impact of previous campaign outcomes (**poutcome**) on the current campaign's success. Calculate the subscription rate (to term deposits) for each **poutcome** category.

Answer:

Select Columns:

- **poutcome**: This column indicates the outcome of the previous marketing campaign (e.g., 'success', 'failure', etc.).
- **count(*) AS total_clients**: Counts the total number of clients for each poutcome category.
- sum(case when y = 'yes' then 1 else 0 end) as
 subscribed_clients: Counts the number of clients who subscribed to the term deposit (y = 'yes') for each poutcome category.
- ROUND(SUM(case when y = 'yes' then 1 else 0 end) * 100.0 /
 COUNT(*), 2) AS subscription_rate: Calculates the subscription rate as a
 percentage and rounds it to two decimal places. This rate is the proportion of clients
 who subscribed to the term deposit out of the total clients in each poutcome
 category.

From:

• The guery uses the client_info table, which contains the relevant data.

Group By:

• Groups the results by the poutcome column, meaning that the calculations (total clients, subscribed clients, and subscription rate) are done separately for each distinct value of poutcome.

Order By:

 Orders the results by subscription_rate in descending order, showing the outcomes with the highest subscription rates first.

Query:

```
hive> select
.. > poutcome,
.. > count(*) AS total_clients,
.. > sum(case when y = 'yes' then 1 else 0 end) as subscribed_clients,
.. > ROUND(SUM(case when y = 'yes' then 1 else 0 end) * 100.0 / COUNT(*), 2) AS subscription_rate
.. > FROM
.. > client_info
.. > GROUP BY
.. > poutcome
.. > ORDER BY
.. > subscription_rate DESC;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
```

Result:

```
success 129 83 64.34
other 197 38 19.29
failure 490 63 12.86
unknown 3705 337 9.1
poutcome 1 0 0.0
5 rows selected (2.943 seconds)
hive>
```

Summary of the Results:

- **poutcome**: Represents different outcomes of the previous marketing campaigns, such as success or failure.
- total_clients: The total number of clients associated with each outcome.
- **subscribed_clients**: The number of clients who subscribed to the term deposit within each poutcome category.
- **subscription_rate**: The percentage of clients who subscribed to the term deposit, calculated per poutcome category.

The results of this query help to understand the effectiveness of previous marketing outcomes (poutcome) by showing how many clients subscribed to the term deposit in each category. By sorting the results by subscription_rate, it highlights which previous outcomes were most successful in leading to subscriptions.

Interpretation:

- High Subscription Rate: A high subscription rate in poutcome category (e.g., success) indicates that clients who had a positive outcome in the previous campaign are more likely to subscribe again in the current campaign.
- Low Subscription Rate: A low subscription rate (e.g., failure or unknown) suggests that clients with unsuccessful or unknown outcomes in the previous campaign are less likely to subscribe in the current campaign.

Here, we can see that the subscription_rate is very high i.e. 64.34 for success poutcome. So, we can say that clients who had a positive outcome in the previous campaign are more likely to subscribe again in the current campaign.

Question 9.2. Compare the average contact duration for clients who subscribed and who did not subscribe to a term deposit.

Answer:

Select Columns:

- y as subscription_status: The y column is renamed as subscription_status. This column indicates whether the client subscribed to a term deposit (yes or no).
- AVG(duration) AS avg_contact_duration: Calculates the average duration of the contact in seconds for each subscription_status. The duration column represents the last contact duration with the client.

From:

• The data is retrieved from the client_info table.

Group By:

• The results are grouped by the subscription_status, meaning the average contact duration is calculated separately for clients who subscribed (y = 'yes') and those who did not (y = 'no').

```
Query: 5 rows selected (2.943 seconds)
hive> select
    . . > y as subscription_status,
    . . > AVG(duration) AS avg_contact_duration
    . . > FROM
    . . > client_info
    . . > group by
    . . > y;
-chgrp: 'GREAT_SUCCESS\ergee' does not match expected pattern for group
Usage: hadoop fs [generic options]
```

Result:

```
no 226.3475
y
yes 552.7428023032629
3 rows selected (1.637 seconds)
hive> _
```

Summary of the Results:

- **subscription_status**: This indicates whether the client subscribed to the term deposit (yes or no).
- avg_contact_duration: This is the average duration (in seconds) of contact with the client, grouped by whether they subscribed or not.

The result provides insight into whether there is a difference in the average contact duration between clients who subscribed to the term deposit and those who did not. For instance, here higher average contact duration for the yes group suggests that longer interactions are more effective in convincing clients to subscribe.

Submission Guidelines:

- Make a copy of this doc file.
- Perform the analysis in your local system using Hadoop and Hive and provide screenshots of both the **code** and the **output** under each question.
- Upload the doc file with other files and submit it in the submission dashboard.