## Data Analysis and Management using Hadoop & Hive

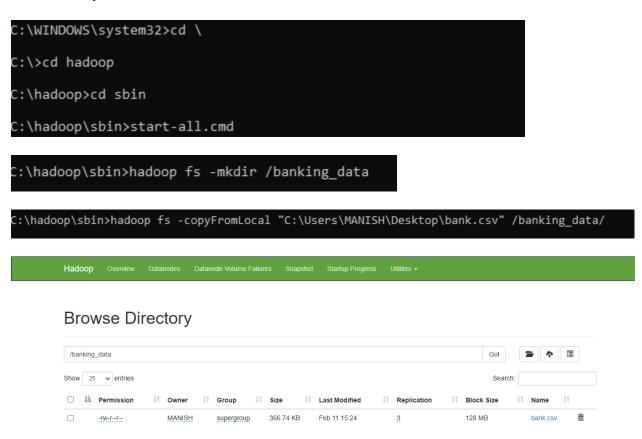
## Hadoop

## 1. Data Ingestion:

Q1: Create a directory in HDFS and transfer the banking dataset from the local system to the HDFS directory.

#### **Steps:**

- 1. Local Setup: Start Hadoop using start-dfs.cmd and start-yarn.cmd.
- 2. HDFS Directory: Create a directory using hadoop fs -mkdir /banking\_data.
- 3. Transfer Data: Use hadoop fs -copyFromLocal bank.csv /banking data/.
- 4. **Verify:** Access localhost: 9870 to check the data.



## 2. Data Transformation with MapReduce:

# Q2.1: Write a MapReduce program in Python that calculates the average account balance for each job type.

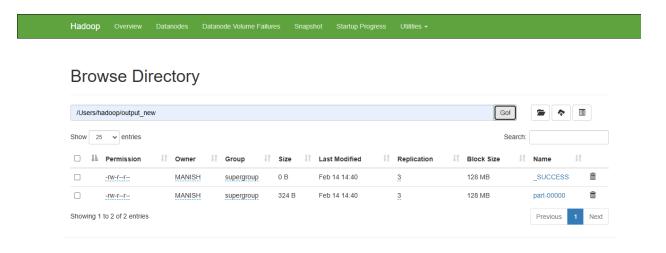
#### **Steps:**

- 1. Write Mapper.py and Reducer.py.
- 2. Upload bank data.csv to HDFS.
- 3. Run the MapReduce job.
- 4. Retrieve results using hadoop fs -cat /output\_new/part-00000.

```
C:\hadoop\bin>hadoop fs -mkdir -p /Users/hadoop/input
C:\hadoop\bin>hadoop fs -copyFromLocal "C:\Users\MANISH\Desktop\Hadoop MapReduce Code\Data transformation MapReduce codes\Question 1\mapper1.py" /Users/hadoop/input
C:\hadoop\bin>hadoop fs -copyFromLocal "C:\Users\MANISH\Desktop\Hadoop MapReduce Code\Data transformation MapReduce codes\Question 1\reducer1.py" /Users/hadoop/input
C:\hadoop\bin>hdfs dfs -get /Users/hadoop/input/mapper1.py C:/Users/hadoop/input/
C:\hadoop\bin>hdfs dfs -get /Users/hadoop/input/reducer1.py C:/Users/hadoop/input/
```

C:\hadoop\bin>hadoop jar "C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar" ^-files "hdfs://localhost:9000/Users/hadoop/input/mapper1.py,hdfs:// localhost:9000/Users/hadoop/input/reducer1.py" ^-mapper "python mapper1.py" ^-reducer "python reducer1.py" ^-input hdfs://localhost:9000/Users/MANISH/input /bank.csv ^-output hdfs://localhost:9000/Users/hadoop/output new

```
C:\hadoop\bin>hadoop fs -cat /Users/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
```



# Q2.2: Write another MapReduce program that counts the number of individuals with and without a housing loan in each education category.

## **Steps:**

• Similar to Q2.1 with modifications in Mapper and Reducer scripts.

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

C:\hadoop\bin>hadoop fs -rm -r /Users/hadoop/output\_new Deleted /Users/hadoop/output\_new

```
C:\hadoop\bin>hadoop fs -rm -r /Users/hadoop/output_new
Deleted /Users/hadoop/output_new
C:\hadoop\bin>hadoop jar "C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar" ^-files "hdfs://localhost:9000/Users/hadoop/input/mapper1.py,hdfs://
localhost:9000/Users/hadoop/input/reducer1.py" ^-mapper "python mapper1.py" ^-reducer "python reducer1.py" ^-input hdfs://localhost:9000/Users/MANISH/input
/bank.csv ^-output hdfs://localhost:9000/Users/hadoop/output_new
```

Q2.3: Perform a MapReduce job to determine the number of clients contacted in each month and their subscription status to term deposits ('y' column).

**Steps:** Write the MapReduce Python script and save as Mapper.py and Reduce

C:\hadoop\bin>hadoop jar "C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar" ^-files "hdfs://localhost:9000/Users/hadoop/input/mapper1.py,hdfs:// localhost:9000/Users/hadoop/input/reducer1.py" ^-mapper "python mapper1.py" ^-reducer "python reducer1.py" ^-input hdfs://localhost:9000/Users/MANISH/input /bank.csv ^-output hdfs://localhost:9000/Users/hadoop/output new

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

## 3. Data Analysis with MapReduce:

## Q3.1: Calculate the average duration of contact per campaign outcome.

#### **Summary:**

• Successful campaigns have the longest average contact duration.

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

C:\hadoop\bin>hadoop jar "C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar" ^-files "hdfs://localhost:9000/Users/hadoop/input/mapper1.py,hdfs://
localhost:9000/Users/hadoop/input/reducer1.py" ^-mapper "python mapper1.py" ^-reducer "python reducer1.py" ^-input hdfs://localhost:9000/Users/MANISH/input
/bank.csv ^-output hdfs://localhost:9000/Users/hadoop/output_new
```

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

# Q3.2: Examine the relationship between the age of clients and their balance, and present findings in a summarised form.

#### **Summary:**

• Shows how balance varies across different age groups.

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

C:\hadoop\bin>hadoop jar "C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.2.4.jar" ^-files "hdfs://localhost:9000/Users/hadoop/input/mapper1.py,hdfs://
localhost:9000/Users/hadoop/input/reducer1.py" ^-mapper "python mapper1.py" ^-reducer "python reducer1.py" ^-input hdfs://localhost:9000/Users/MANISH/input
/bank.csv ^-output hdfs://localhost:9000/Users/hadoop/output_new
```

```
:\hadoop\bin>
:\hadoop\bin>hadoop fs -cat
19
        393.50
20
        661.33
21
        1774.29
        1455.33
22
23
        2117.95
24
        634.62
25
        1240.07
26
        788.56
27
        851.78
                                         1665.63
28
        1025.10
                                         1755.08
29
                              59
60
61
63
64
65
66
67
        1261.88
                                         1582.48
30
        1113.03
                                         2964.57
31
                                         2407.50
        1288.48
                                         516.14
32
        1256.55
                                         2286.38
33
        1545.41
                                         1103.29
34
        1111.54
                                         1638.17
35
        1192.83
                                         3313.89
36
        1226.89
                                         4149.40
37
        1463.92
                                         11753.00
38
        1718.99
                              69
                                         774.33
39
        1104.86
                              70
71
72
73
74
75
76
                                         5084.57
40
        1399.51
                                         3787.33
41
        1505.79
                                         2526.00
42
                                         525.83
        1612.36
43
                                         1978.33
        1807.83
44
                                         7046.50
        1836.55
45
                                         1338.00
        1187.37
                              77
78
79
                                         2405.17
46
        998.77
                                         318.00
47
        1363.05
                                         4087.75
48
        1462.36
                              80
                                         4183.50
49
        1591.11
                              81
                                         1.00
50
        1645.06
                               83
                                         380.50
51
        1528.57
                              84
                                         639.00
52
        782.29
                                         1503.00
                              86
53
        1588.31
                              87
                                         230.00
        1656.66
                              C:\hadoop\bin>_
        1244.94
        2120.14
```

## **Summary of Findings:**

After executing the **MapReduce** job, we obtained the **average account balance for each specific age** from the dataset. Below are the key observations:

#### • Age-Specific Averages:

- o The output provides the **average balance** associated with each age group.
- 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 Observed in the Data:

- Increase with Age: In many cases, there is a gradual increase in average balance as individuals grow older, potentially due to higher salaries, career growth, or accumulated savings.
- o **Fluctuations**: Some age groups show **higher or lower** average balances, possibly due to factors such as student loans, mortgages, or retirement planning.

#### • Variability Across Age Groups:

- o The **average balance varies significantly** across different age groups, reflecting different **financial habits** and life circumstances.
- Some anomalies may indicate specific financial behaviors, such as a sudden dip
  in balances due to major life expenses or a rise in balances for retirees with large
  savings.

#### **Conclusion:**

• The MapReduce job successfully calculated the **average balance for each age group**, revealing important **financial trends**. These insights can be used by **banks and financial institutions** to create targeted **financial products**, such as **age-specific loan offers**, savings plans, and investment strategies.

## **Hive**

## 1. Data Ingestion and Table Creation:

#### Q1.1: Create a Hive database named banking\_data.

```
CREATE DATABASE banking_data;
USE banking_data;
```

```
NIVEY CREATE DATABASE banking data;
2025-02-21T18:38:07,980 INFO [10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.ql.session.SessionState - METASTORE_FILTER_HOOK will be ignored, since hive.security.authorization.manager is set to instance of HiveAuthorizerFactory.

Note taken: 1.719 seconds
2025-02-21T18:38:09,839 INFO [10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.ql.session.SessionState - Resetting thread name to main

hive> USE banking data;
2025-02-21T18:39:00,796 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main

OK
Time taken: 0.094 seconds
2025-02-21T18:39:00,796 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b main] org.apache.hadoop.hive.conf.HiveConf -
```

## Q1.2: Define and create a Hive table client\_info with appropriate data types for the bank.csv dataset.

```
CREATE TABLE client_info (
    age INT, job STRING, balance FLOAT, ...
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;
```

```
NIVEY CREATE TABLE client_info(

age INT,

job STRING,

marital STRING,

default STRING,

default STRING,

housing STRING,

housing STRING,

loan STRING,

contact STRING,

contact STRING,

month STRING,

month STRING,

duration INT,

month STRING,

duration INT,

previous INT,

previous INT,

previous INT,

previous INT,

previous INT,

previous FTRING,

)

STRING AS TEXTFILE;

205-202-21T18:47:27,940 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: 10a63bfa-a255-4d01-8f85-991187dace5b

205-02-21T18:47:27,940 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to 10a63bfa-a255-4d01-8f85-991187dace5b

main String As accounts

Time taken: 2.437 seconds

String As accounts

Time taken: 2.437 seconds

String As accounts

Time taken: 2.437 seconds

Time taken: 2.437 seconds

String As accounts

Time taken: 2.437 seconds

Time taken: 2.437 seconds

String As accounts

Time taken: 2.437 seconds

Time taken: 2.437 seconds

String As accounts

Time taken: 2.437 seconds

Time taken: 2.438 seconds

Time taken: 2.438 seconds

Time taken: 2.439 seconds

Time taken: 2.437 seconds

Time taken: 2.438 seco
```

## Q1.3: Load the data from the bank.csv file into the client\_info table.

```
ive> LOAD DATA LOCAL INPATH 'C:/Users/MANISH/Desktop/bank.csv' INTO TABLE client info;
1025-02-22T11:08:16,105 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d
025-02-22T11:08:16,105 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
oading data to table banking_data.client_info
.
ime taken: 3.469 seconds
025-02-22T11:08:19,580 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log i
  25-02-2711:08:19,580 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.ql.session.SessionState - Resetting thread name to
 TIME JOINT CLAIM 100, 1100 CLAIM 100, 2025-02-22111:88:57,899 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 825-02-22111:88:57,900 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main | Org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs://localhost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main | Org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs://localhost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive_2025-02-22_11-08-57_925_7904764322295344929-1/-mm-10001/.hive-staging_hive_2025-02-22
 11-08-57 925_7904764322295340499-1
025-02-22T11:09:01,635 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.conf.Configuration.deprecation - mapred.task.is.map is deprecated. In tead, use mapreduce.task.ismap
                                                      education default NULL
married primary no 1787
married secondary no
single tertiary no
married tertiary no
                                                                                                                                  housing loan
no no
4789 yes
1350 yes
                                 marital education
                                                                                                                                                                          contact NULL month
cellular 19
                                                                                                                                                                         contact f
cellular
yes
no
yes
no
no
no
                                                                                                                                  no
4789
1350
1476
                                                                                                                                                                                                                                                                                                                 0
339
330
                                                                                                                                                                                                                                                      79
may
apr
199
226
feb
may
may
57
                                                                                                                                                                                            cellular
cellular
unknown 3
unknown 5
                                                                                                                                                                                                                                                                                                                                                      failure no
failure no
        management married tertiary no 1476 yes yes unknown 3 jun 199 4 -1 0 unknown no blue-collar married secondary no 0 yes no unknown 5 may 226 1 -1 0 unknown no management single tertiary no 747 no no cellular 23 feb 141 2 176 3 failure no self-employed married tertiary no 307 yes no cellular 14 may 341 1 330 2 other no technician married secondary no 147 yes no cellular 6 may 151 2 -1 0 unknown no entrepreneur married tertiary no 221 yes no unknown 14 may 57 2 -1 0 unknown no taken: 3.735 seconds, Fetched: 10 row(s) -62-22711:99:02,072 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed/inclored id-
          02-22111:09:02,072 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.ql.session.SessionState - Resetting Chicago Walk to a chiagin Windo
```

## 2. Basic Data Exploration:

#### Q2.1: Write a HiveQL query to count the total number of clients in the dataset.

SELECT COUNT(\*) AS total clients FROM client info;

```
hive> select count(*) AS total_clients FROM client_info;
2025-02-22T11:13:11,970 INFO [main] org.apache.hadoop.hive.conf.HiveConf -
2025-02-22T11:13:11,971 INFO [main] org.apache.hadoop.hive.ql.session.Sess:
2025-02-22T11:13:12,905 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] or
ost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive_2025-02-
76205251589225147-1
Query ID = MANISH_20250222111311_06e0cea0-5ae0-4deb-b6ab-be8152e4b487
Total jobs = 1
Launching Job 1 out of 1
```

```
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 9.371 sec HDFS Read: 390574 HDFS Write: 104 SUCCESS
Total MapReduce CPU Time Spent: 9 seconds 371 msec
OK
4522
Time taken: 60.881 seconds, Fetched: 1 row(s)
2025-02-22T11:14:13,019 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.conf.HiveCo
dd8-44d7-bc75-b4465eb86d38
2025-02-22T11:14:13,020 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.ql.session.
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.

#### Q2.2: Display first 10 rows.

```
SELECT * FROM client info LIMIT 10;
```

```
hive> SELECT * FROM client_info LIMIT 10;_
2025-02-2211:18:28,112 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-2211:18:28,112 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
2025-02-2211:18:28,478 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs:
//localhost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive_2025-02-22_11-18-28_138_3095909105960493795-1/-mr-10001/.hive-staging_hive_2025-02-22
//l-18-28_138_3095909105960493795-1
```

```
OK

NULL job marital education default NULL housing loan contact NULL month NULL NULL NULL poutcome y

unemployed married primary no 1787 no no cellular 19 oct 79 1 -1 0 unknown no

services married secondary no 4789 yes yes cellular 11 may 220 1 339 4 failure no

so management single tertiary no 1350 yes no cellular 16 apr 185 1 330 1 failure no

so management married tertiary no 1476 yes yes unknown 3 jun 199 4 -1 0 unknown no

so blue-collar married secondary no 0 yes no unknown 5 may 226 1 -1 0 unknown no

management single tertiary no 747 no no cellular 23 feb 141 2 176 3 failure no

self-employed married tertiary no 307 yes no cellular 14 may 341 1 330 2 other no

self-employed married secondary no 147 yes no cellular 14 may 341 1 330 2 other no

self-employed married tertiary no 221 yes no cellular 6 may 151 2 -1 0 unknown no

tentrepreneur married tertiary no 221 yes no unknown 14 may 57 2 -1 0 unknown no

time taken: 0.406 seconds, Fetched: 10 row(s)

2025-09-22111:18:28,709 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.conf.HiveConf - Using the default value pression of the content of the cont
```

#### **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.

## 3. Data Filtering and Sorting

#### Q3.1: Retrieve all records of clients who are married and have a personal loan.

```
SELECT * FROM client info WHERE marital = 'married' AND loan = 'yes';
```

```
hive> select * from client_info where marital ='married' ANd loan='yes';
2025-02-22T11:21:37,043 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:21:37,044 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
2025-02-22T11:21:38,043 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs://locallost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive_2025-02-22_11-21-37_067_3878738514662863417-1/-mr-10001/.hive-staging_hive_2025-02-22_11-21-37_067_38
```

services	married secondary	no	4789	yes	yes	cellular	11	may	220		339		failure	no
management	married tertiary	no	1476	yes	yes	unknown 3	jun	199				unknowr		
services	married primary no	-88	yes	yes	cellula		apr	313		147		failure		
blue-collar	married secondary	no	360	yes	yes	cellular	29	jan	89		241		failure	
management	married tertiary		194	no	yes	cellular	29	aug	189				unknown	no no
self-employed	married secondary		784		yes	cellular	30	jul	149				unknown	no no
admin. married	secondary no	105		yes	cellula	ır 21	aug	74				unknowr	n no	
management	married secondary	no	82	no	yes	telephone		feb	140				unknown	no no
blue-collar	married primary no	-516	no	yes	telepho		jul	554		-1	0	unknowr	n no	
management	married secondary	no	0	no	yes	cellular		jul	630			0	unknown	no no
blue-collar	married secondary	no	427	yes	ves	unknown 9	jun	371		-1	0	unknowr	n no	
self-employed	married secondary	no	217	yes	ves	cellular	15	jul	317	5	-1	0	unknown	no no
blue-collar	married secondary	no	-231	no	ves	cellular	15	jul	779	2	-1	0	unknown	
admin. married		323	yes	yes	unknown		280	2	-1	ø	unknowr			
management	married tertiary	no	106	no	yes	cellular	11	aug	588	2	-1	е	unknown	no.
	primary no 1906	no	ves	unknown		jun 45	9	-1	0	unknowr			dilicitoni	
convices	primary no 1906 married secondary	no	978	ves	ves	unknown 26	may	82	2	-1	0	unknowr	no.	
admin. married	secondary no	-465		yes	cellula	ir 23	jul	166	1	-1	9	unknown		
admin. married		5181	yes	ves	cellula		iul	18	7	-1	9	unknowr		
			yes											
blue-collar	married primary no	0	yes	yes	telepho		jul	97	6 152		0	unknowr a		
services	married secondary	no	1	no	yes	cellular	21	nov		2	-1		unknown	
technician	married secondary	no	2030	yes	yes	cellular		jul	196		-1		unknown	no
blue-collar	married primary no	305	yes	yes	telepho		jul	834	10		0	unknowr		
technician	married tertiary	no		yes	yes	cellular		sep	112		62		other	no
admin. married		-247	yes	yes	unknown		344				unknowr			
technician	married secondary	no		no	yes	cellular		may			172		failure	no
blue-collar	married secondary	no	989	yes	yes	unknown 23	may	246				unknowr		
management	married tertiary	no	415	no	yes	cellular		jul	361				unknown	no no
housemaid	married secondary	no	209	yes	yes	cellular		jul					unknown	no no
entrepreneur	married tertiary	no	624	no	yes	cellular		jul	180				unknown	
self-employed	married secondary	no	1516	yes	yes	unknown 23	may	373				unknowr	n no	
technician	married tertiary	no	-988	yes	yes	cellular	15	jul					unknown	no
admin. married		69	no	yes	cellula		aug	120	3	1	6	success		
entrepreneur	married secondary	no	593	yes	yes	cellular	24	jul	1484	24	-1	0	unknown	ves
management	married unknown no	353	no	ves	cellula		jul	171	2	-1	0	unknowr		,
technician	married secondary	no	205	no	yes	cellular	23	jul	442	2	-1	0	unknown	no .
unemployed	married secondary married primary no	1147	ves	yes	unknown		249	5	-1	9	unknowr		amenowi	
blue-collar	married secondary	no	8545	yes	yes	cellular	6	may	199	7			failuse	no
management	married secondary	no	2	no		cellular	20		472	2	-1	A <b>ģ</b> tivate	: үүрдас	WS
admin mannied	primary no 276	no		unknown	yes 17	jun 641	6	aug -1	0	unknowr		Go to Sett	ings to act	tivate Winc
blue coller	primary no 276 married primary no		yes					1	-1				mgs to at	arace wille
pide-collar	married primary no married secondary	214	yes 1760	yes no	unknown	9 jun cellular	168 19			0 1	unknowr -1	n no 0	lee	
unemployed	married secondary	no	1700	110	yes	Cellulai.	19	nov	162	1	-1	U	unknown	i iio
solf omployed	manniad cosendary	no	900	no	VOC	cellular	9	503	170	1	1	0	unknown	no
self-employed	married secondary	no	800	no				jul	170	1			unknown	110
blue-collar	married primary no	506	yes 210	yes	unknown		122	1		0	unknown a			
services	married secondary	no		yes		unknown 20	may	201				unknown		
unemployed	married secondary	no	-872	yes	yes	cellular	20	nov	153		183		failure	
services	married secondary	no	0	no		cellular	16	jul	1473				unknown	
management	married tertiary	no	5057	no		cellular	19	nov	37				unknown	
management	married tertiary	no	4039	no	yes	cellular	25	jul	106				unknown	
management	married tertiary	no	997	yes		cellular	21	nov	81				unknown	
technician	married secondary					cellular		jul	624				unknown	
services	married secondary	yes	-220	yes		cellular	25	jul	123				unknown	no
retired married	secondary no	3382	no	yes	cellular		may	294		309		failure		
technician	married tertiary	no	199	yes		cellular		feb	116				failure	no
admin. married	secondary no	147	yes	yes	cellular		jan			184		failure		
technician	married secondary	no	2225	no		cellular	13	aug					unknown	no
housemaid	married unknown yes		no	yes	telephor	ne 7	jul	94				unknown	no	
technician	married tertiary	no	3337	yes		telephone	31	jul	24	14			unknown	no
services	married secondary	no	895	yes		unknown 4	jun	622				unknown		
services	married secondary	no	0	yes	yes	cellular	31	jul	187				unknown	no
blue-collar	married secondary	no	-27	no	yes	telephone	31	jul	77	13	-1		unknown	
blue-collar	married primary no	293	ves	yes	cellular		may	102	1	-1	0	unknown		
management	married tertiary	no	19447	yes		cellular	21	nov	166	1	-1		unknown	no
admin. married	secondary no	9	yes	yes	cellular		nov	159	2	195	2	failure		
blue-collar	married primary no	5431	yes	yes	unknown		383	1	-1	0	unknown			
blue-collar	married secondary	no	225	yes		unknown 7	may	866	2	-1		unknown	no	
blue-collar	married secondary	no	3653	yes		cellular	111ay 21	nov	252	1			failure	no
management	married secondary married tertiary	no	436	no	yes	cellular	28	jul	118	4			unknown	
						cellular	28 7			3			unknown	
technician	married tertiary	no	101 302	yes		cellular	16	jul	187 670				unknown	
unemployed	married tertiary	no		yes				apr						IIO
management	married tertiary	no	1557	yes	yes	unknown 13	may	213	1	-1	0	unknown		
services	married secondary	no	505	yes		unknown 27	may	371		-1		unknown	по	
blue-collar	married primary no	190	yes	yes	unknown		194		-1	0	unknown			
management	married tertiary	no	-17	yes		cellular	11	may	474		256		success	
self-employed	married secondary	no	2678	no		cellular	18	aug	151	12			unknown	
services	married secondary		-91	yes		cellular		feb					unknown	no
self-employed	married primary no	362	no	yes	cellular		jul	816				unknown	yes	
admin. married	unknown no 642	yes	yes	unknown		may 509	2			unknown				
admin. married self-employed	unknown no 642 married tertiary	yes yes	yes -3313			may 509 unknown 9	2 may	-1 153	1	unknown -1		unknown	no	

# Q3.2: List the top 10 clients with the highest balance, displaying their job, marital status, and balance.

taken: 1.199 seconds, Fetched: 453 row(s)
-02-22T11:21:38,686 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-dd7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-dd7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.ql.session.SessionState - Resetting threatD frame(n) for magnivate Windows.

SELECT job, marital, balance FROM client\_info ORDER BY balance DESC LIMIT 10;

2025-02-22711:28:00,238 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 2025-02-22711:28:00,238 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 2025-02-22711:28:00,239 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main 2025-02-22711:28:00,595 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs: //localhost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-28-00\_262\_7415283650087693896-1/-mr-10001/.hive-stagin

```
retired married 71188
entrepreneur
             married 42045
technician
             single 27733
             married 27359
management
technician
             married 27069
housemaid
             single 26965
retired married 26452
services
             married 26394
management
             divorced
                           26306
retired single 25824
Time taken: 47.462 seconds, Fetched: 10 row(s)
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:28:47,830 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38
hive>
```

## 4. Data Aggregation and Grouping

#### Q4.1: Average age per job category.

```
SELECT job, AVG(age) AS average_age FROM client_info GROUP BY job;
```

hive> select job,AVG(age) AS average\_age from client\_info group by job;\_\_
2025-02-22711:30:32,696 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22711:30:32,697 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
2025-02-22711:30:33,040 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs:
//localhost:9000/tmp/hive/NANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_722\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_5607926747041170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_72\_560792674704170607-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-30-32\_

```
admin. 39.68200836820084
blue-collar 40.15644820295983
entrepreneur
               42.01190476190476
housemaid
               47.339285714285715
job
        NULL
management
               40.54076367389061
retired 61.869565217391305
self-employed 41.45355191256831
services
                38.57074340527578
student 26.821428571428573
technician
               39.470052083333336
               40.90625
unemployed
unknown 48.10526315789474
Time taken: 46.533 seconds, Fetched: 13 row(s)
2025-02-22T11:31:19,353 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:31:19,354 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.
```

Q4.2: Find the total number of clients for each education level who have defaulted on credit.

SELECT education, COUNT(\*) AS total\_defaulted\_clients FROM client\_info WHERE
default = 'yes' GROUP BY education,default;

hive> select education,default,count(\*) AS total\_defaulted\_client from client\_info where default = 'yes' group by education,default;

2025-02-22711:33:59,681 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38

2025-02-22711:33:59,681 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main

2025-02-22711:33:59,681 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs:

//localhost:9000/tmp/hive/MANTSh/f49f5dee-fdd8-44d7-bc75-b4465eb86d38/hive\_2025-02-22\_11-33-59\_703\_5822997283630260050-1/-mr-10001/.hive-staging\_hive\_2025-02-22

\_11-33-59\_703\_5822997283630260050-1

```
OK
primary yes 10
secondary yes 46
tertiary yes 17
unknown yes 3
Time taken: 49.307 seconds, Fetched: 4 row(s)
2025-02-22T11:34:49,107 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 ma
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:34:49,107 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 ma
```

## 5. Complex Queries for Insights

Q5.1: 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.

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;

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 I else 0 end) as subscribed\_clients from client\_info group by job order by avg\_balance desc limit 5)sc; 2025-02-22111:41:03,384 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 2025-02-22111:41:03,384 INFO [main] org.apache.hadoop.hive.cons.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main 2025-02-22111:41:03,873 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: hdfs://localhost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_11-41-03\_415\_1402845492661504678-1/-mr-10001/.hive-staging\_

```
retired 2319.191304347826
                           23.47826086956522
housemaid 2083.8035714285716
                                    12.5
              1766.9287925696594
                                     13.519091847265221
management
entrepreneur 1645.125
                              8.928571428571429
student 1543.8214285714287
                              22.61904761904762
Time taken: 90.571 seconds, Fetched: 5 row(s)
2025-02-22T11:42:34,072 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apac
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:42:34,073 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apac
hive>
```

Q5.2: 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).

Select month, total\_contacts, (successful\_contacts/total\_contacts) \*100 as success\_rate from (select month,count(\*) as total\_contacts,sum(case when y='yes' then 1 else 0 end) as successful\_contacts from client\_info group by month order by total contacts DESC limit 1) as top\_month;

hive> select month,total\_contacts,(successful\_contacts/total\_contacts) \* 100 as success\_nate\_from (select month,count(\*) as total\_contacts,sum(case when y='yes' then 1 else 0 end) as successful\_contacts from client\_info group by month order by total\_contacts DESC limit 1) as top month;
2025-02-22T11:50:01,320 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
2025-02-22T11:50:01,321 INFO [main] org.apache.hadoop.hive.ql.session.State - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main [076.apache.hadoop.hive.common.fileUtils - Creating directory if it doesn't exist: hdfs:
//localhost:9000/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive\_2025-02-22\_11-50-01\_346\_8132715559582893855-1/-mr-10001/.hive-staging\_hive\_2025-02-22\_
11-50-01\_346\_8132715559582893855-1

```
OK
may 1398 6.652360515021459
Time taken: 89.291 seconds, Fetched: 1 row(s)
2025-02-22T11:51:30,729 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:51:30,730 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
hive>
```

## 6. Correlation Analysis

Q6: Calculate the correlation between age and balance for the clients.

SELECT CORR(age, balance) as age balance correlation from client info;

```
hive's select CORR(age,balance) as age_balance_correlation from client_info;
2025-02-22T11153:28,720 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:53:28,720 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main
2025-02-22T11:53:29,088 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FleUtils - Creating directory if it doesn't exist: hdfs:
//localhosti9900/tmp/hive/MANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive_2025-02-22_11-53-28_742_622968810028760442-1/-mr-10001/.hive-staging_hive_2025-02-22
```

```
OK
0.08382014224477742
Time taken: 47.859 seconds, Fetched: 1 row(s)
2025-02-22T11:54:16,703 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 ma
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:54:16,704 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 ma
hive>
```

## 7. Trend Analysis

Q7: 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).

SELECT SUBSTRING(month, 1, 4) AS year, COUNT(\*) AS num\_clients\_contacted from client info group by substring(month,1,4) order by year;

```
apr
      293
aug
      633
dec
      20
feb
      222
jan
      148
jul
      706
jun
      531
      49
mar
      1398
may
mont
      1
      389
nov
oct
      80
      52
sep
Time taken: 87.188 seconds, Fetched: 13 row(s)
9f5dee-fdd8-44d7-bc75-b4465eb86d38
2025-02-22T11:59:02,604 INFO [fa9f5dee-fdd8-44d7-bc
hive>
```

## 8. Anomaly Detection

Q8: Detect unusual patterns in average yearly balance across education levels.

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

```
nive> 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

> A Subquery;

2025-02.22T12:15:04,927 INFO [main] org.apache.hadoop.hive.conf.HiveConf - Using the default value passed in for log id: fa9f5dee-fdd8-44d7-bc75-b4465eb86d3

2025-02.22T12:15:04,927 INFO [main] org.apache.hadoop.hive.ql.session.SessionState - Updating thread name to fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main

2025-02.22T12:15:05,304 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.apache.hadoop.hive.common.FileUtils - Creating directory if it doesn't exist: If //localnost:9000/tmp/hive/NANISH/fa9f5dee-fdd8-44d7-bc75-b4465eb86d38/hive_2025-02-22_12-15-04_950_6068760278293976634-1/-mr-10001/.hive.stagging_hive_2025-04.
```

```
OK
apr
        primary 1.267794378462834
                        -0.7235770026264813
apr
        secondary
apr
        tertiary
                        0.6555354947097891
        unknown -1.1997528705461413
apr
        primary -0.8113812295586758
aug
        secondary
                        -0.5775732890112949
aug
aug
        tertiary
                        -0.31633238372712924
        unknown 1.7052869022970998
aug
dec
        primary -0.8088635843923103
dec
        secondary
                        -0.397568378390461
dec
        tertiary
                        -0.5058241544774784
        unknown 1.7122561172602497
dec
feb
        primary -1.0701923004119014
feb
        secondary
                        -0.9216091477916096
feb
        tertiary
                        0.8917649610011248
        unknown 1.1000364872023862
feb
        primary 1.1404156384061253
ian
jan
        secondary
                        -0.4369029752604172
        tertiary
                        0.7115176049750513
ian
        unknown -1.41503026812076
jan
iul
        primary 0.9649291030842655
jul
        secondary
                        -0.4258456259680814
jul
        tertiary
                        0.9014112560610348
iul
        unknown -1.4404947331772213
iun
        primary -0.15682779176634196
jun
        secondary
                        -1.0365919570822748
jun
        tertiary
                        1.6428352994501707
iun
        unknown -0.44941555060155386
        primary -1.0722876073529977
mar
mar
        secondary
                        -0.2462824810074777
mar
        tertiary
                        1.6391398006892093
        unknown -0.3205697123287344
mar
        primary -1.0048248682462826
may
        secondary
                        -0.9457054758622359
may
may
        tertiary
                        0.6640708393343474
may
        unknown 1.2864595047741691
        education
mont
                        NULL
        primary -0.8137802862217836
nov
        secondary
                    0.2770762379753206
nov
nov
        tertiary
                       1.5167439470393023
```

```
1.5167439470393023
nov
        tertiary
nov
       unknown -0.9800398987928378
       primary 1.7097889755337383
oct
oct
       secondary
                        -0.37052787477325194
       tertiary
                        -0.5238214197432853
oct
       unknown -0.8154396810172013
oct
       primary -1.3086820436917062
sep
       secondary
                        -0.4707154583719365
sep
                        0.3984500713422833
sep
       tertiary
sep
       unknown 1.38094743072136
Time taken: 45.104 seconds, Fetched: 49 row(s)
2025-02-22T12:15:50,234 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.
dd8-44d7-bc75-b4465eb86d38
2025-02-22T12:15:50,235 INFO [fa9f5dee-fdd8-44d7-bc75-b4465eb86d38 main] org.
hive>
```

## 9. Advanced Analysis

Q9.1: 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.

```
OK
success 129
         83
              64.34
         38
other
    197
              19.29
failure 490
         63
              12.86
unknown 3705
         337
              9.10
poutcome
                   0.00
         1
              О
Time taken: 91.906 seconds, Fetched: 5 row(s)
9f5dee-fdd8-44d7-bc75-b4465eb86d38
hive>
```

# Q9.2: Compare the average contact duration for clients who subscribed and who did not subscribe to a term deposit.

SELECT y AS subscription\_status, AVG(duration) AS avg\_contact\_duration FROM client info GROUP BY y;

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