HADOOP

1. Data Ingestion:

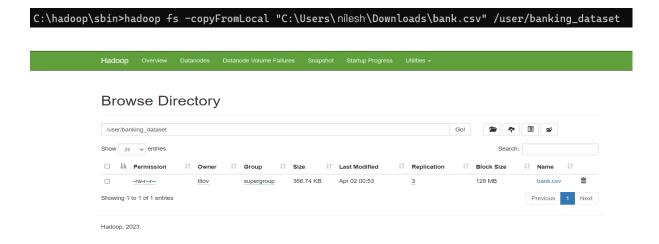
 Create a directory in HDFS and transfer the banking dataset from the local system to the HDFS directory.

Solution:-

Step-1 Create a directory in HDFS



Step-2 Transfer the banking dataset from the local system to the HDFS directory.



2. Data Transformation with MapReduce:

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

```
C:\hadoop\sbin>hadoop fs -copyFromLocal "C:\Users\nilesh\Downloads\Hadoop\mapper.py" /user/nilesh/input/
C:\hadoop\sbin>hadoop fs -copyFromLocal "C:\Users\nilesh\Downloads\Hadoop\reducer.py" /user/nilesh/input/
```

C:\hadoop\sbin>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop\reducer.py"" -mapper "py
thon manger.py" -reducer "nython reducer.py" -innut /user/ -innut

```
C:\hadoop\sbin>hadoop fs -cat /user/hadoop/output new/part-*
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
```

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

```
C:\hadoop\sbin>hadoop fs -copyFromLocal "C:\Users\nilesh\Downloads\Hadoop\New folder\mapper.py" /user/nilesh/input/copyFromLocal: `/user/nilesh/input/mapper.py': File exists
C:\hadoop\sbin>hadoop fs -copyFromLocal "C:\Users\nilesh\Downloads\Hadoop\New folder\reducer.py" /user/nilesh/input/copyFromLocal: `/user/nilesh/input/reducer.py': File exists
```

```
C:\hadoop\shin>hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop\tools\lib\hadoop\streaming-3.2.4.jar -files "file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop\New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop\New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/Hadoop/New_folder/mapper.py", file:///"C:/Users/nilesh/Downloads/New_folder/mapp
```

```
C:\hadoop\sbin>hadoop fs -cat /user/hadoop/output_q2/part-*
primary 94 583
secondary 416 1889
tertiary 173 1176
unknown 7 179
```

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

:\haddop\sbin>hadoop\sbin>hadoop\spin>jar (:\hadoop\shin>hadoop) jar (:\hadoop\shin>hadoop) jar (:\hadoop\shin) jar (:\hadoop\

```
C:\hadoop\sbin>hadoop fs -cat /user/hadoop/output_q3/part-*
        56
                 236
aug
        79
                 553
dec
        8
                 11
feb
        38
                 183
        16
                 131
jan
                 644
jul
        61
jun
        55
                 475
mar
        21
                 27
                 1304
        93
may
                 350
nov
        38
        37
                 42
oct
sep
        16
                 35
```

3. Data Analysis with MapReduce:

Analyze the average duration of contact (in seconds) per campaign outcome ('poutcome').

```
C:\hadoop\sbin>hadoop fs -cat /user/hadoop/output_q5/part-*
failure 254.38367346938776
other 273.83248730964465
success 338.6356589147287
unknown 262.1031039136302
```

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

```
C:\hadoop\sbin>hadoop fs -cat /user/hadoop/output_q6/part-*
19
        393.5
20
        661.3333333333334
21
        1774.2857142857142
22
        1455.3333333333333
23
        2117.95
24
       634.625
25
        1240.0681818181818
26
        788.5584415584416
27
        851.7765957446809
28
        1025.0970873786407
29
        1261.8762886597938
30
        1113.03333333333333
31
        1288.4824120603016
32
        1256.549107142857
33
        1545.4139784946237
        1111.5367965367966
34
35
        1192.827777777778
36
        1226.8936170212767
37
        1463.9192546583852
38
        1718.993710691824
39
        1104.8615384615384
40
       1399.5070422535211
41
       1505.7925925925927
42
       1612.3617021276596
43
       1807.8347826086956
       1836.5523809523809
44
       1187.3660714285713
45
46
       998.7731092436975
47
       1363.0462962962963
       1462.359649122807
48
49
       1591.107142857143
50
       1645.0594059405942
51
       1528.5714285714287
52
        782.2906976744187
53
       1588.3085106382978
54
       1656.661971830986
       1244.94444444443
55
56
       2120.135135135135
57
       1665.6263736263736
58
       1755.0823529411764
59
       1582.4788732394366
60
       2964.574468085106
61
       2407.5
62
       516.1428571428571
63
       2286.375
64
       1103.2857142857142
```

HIVE

1. Data Ingestion and Table Creation:

- Create a Hive database named banking_data.
- Define and create a Hive table client_info with appropriate data types for the bank.csv dataset.

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

hive> SELECT * FROM client_info LIMIT 5;

MITT	job marita	1 education	default	MITT	housing	loan	contact	MITT	month	MITT	MITT	MIII I	MIII I	noutcor	10	V
																y
30	unemployed	married primar	y no	1787	no	no	cellula		19	oct	79	1	-1	0	unknown	no
33	services	married second	ary	no	4789	yes	yes	cellula		11	may	220		339		failure no
35	management	single tertia	ry	no	1350	yes	no	cellula		16	apr	185		330		failure no
30	management	married tertia	ry	no	1476	yes	yes	unknown		jun	199				unknown	no

2. Basic Data Exploration:

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

hive> select count(*) from client_info;

```
OK
2024-05-05 17:27:50,472 INFO ql.Driver: OK
2024-05-05 17:27:50,472 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 17:27:50,478 INFO mapred.FileInputFormat: Total input files to process : 1
2024-05-05 17:27:50,485 INFO exec.ListSinkOperator: RECORDS_OUT_OPERATOR_LIST_SINK_10:1, RECORDS_OUT_INTERMEDIATE:0,
4522
Time taken: 24.905 seconds, Fetched: 1 row(s)
```

Display the first 10 rows of the dataset.

hive> select * from client_info limit 10;

2024-05	-05 17:05:11,599	INFO mapred.File	InputFor	mat: To	tal input	: files :	to proces	ss : 1								
NULL	job marital	education	default	NULL	housing	loan	contact	NULL	month	NULL	NULL	NULL	NULL	poutcom	e	
30	unemployed	married primary	no	1787	no	no	cellular		19	oct	79				unknown	no
33	services	married secondar	y	no	4789	yes	yes	cellular		11	may	220		339		failure no
35	management	single tertiary		no	1350	yes	no	cellular		16	apr	185		330		failure no
30	management	married tertiary		no	1476	yes	yes	unknown		jun	199				unknown	no
59	blue-collar	married secondar	y	no		yes	no	unknown		may	226				unknown	no
35	management	single tertiary		no	747	no	no	cellular			feb	141		176		failure no
36	self-employed	married tertiary		no	307	yes	no	cellular		14	may	341		330		other no
39	technician	married secondar	У	no	147	yes	no	cellula			may	151				unknown no
41	entrepreneur	married tertiary		no	221	yes	no	unknown	14	may	57		-1		unknown	no

3. Data Filtering and Sorting:

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

hive> select * from client info where marital='married' AND loan='yes';

36	-1 Ø	unknown			070			collula		20		453
36	unemployed 1 183	married 1	secondary failure no	no	-872	yes	yes	cellula		20	nov	153
38	services 6 -1	married 0	secondary unknown no	no	0	no	yes	cellula		16	jul	1473
48	management		tertiary	no	5057	no	yes	cellula		19	nov	37
37	1 -1 management	0 married	unknown no tertiary	no	4039	no	yes	cellula		25	jul	106
	2 -1	0	unknown no									
34	management 1 -1	married 0	tertiary unknown no	no	997	yes	yes	cellula		21	nov	81
31	technician 6 -1	married 0	secondary unknown no	no	0	no	yes	cellula		15	jul	624
32	services	married	secondary	yes	-220	yes	yes	cellula		25	jul	123
58	2 -1 retired married	0 seconda	unknown no ry no	3382	no	yes	cellular	r	28	may	294	2
50	309 2 technician	failure married	no tertiary	no	199	yes	yes	cellula		3	feb	116
	2 253	1	failure no									
35	admin. married	failure		147	yes	yes	cellula	r	29	jan	41	1
54	technician 2 -1	married 0	secondary unknown no	no	2225	no	yes	cellula		13	aug	73
31	housemaid	married	unknown yes	-6	no	yes	telephor	ne	7	jul	94	2
50	-1 0 technician	unknown married	tertiary	no	3337	yes	yes	telepho	ne	31	jul	24
36	14 -1 services	0 married	unknown no secondary	no	895	yes	yes	unknown	4	jun	622	1
	-1 0 services	unknown	no									
37	services 2 -1	married 0	secondary unknown no	no	0	yes	yes	cellula		31	jul	187
58	blue-collar 13 -1	married 0	secondary unknown no	no	-27	no	yes	telepho	ne	31	jul	77
41	blue-collar	married	primary no	293	yes	yes	cellula	r		may	102	1
50	-1 0 management	unknown married	no tertiary	no	19447	yes	yes	cellula		21	nov	166
32	1 -1 admin. married	0 seconda	unknown no ry no	0	yes	yes	cellular	r	17	nov	159	2
	195 2	failure	no									
57	blue-collar 0 unknown	no	primary no	5431	yes	yes	unknown		may	383	1	-1
28	blue-collar -1 0	married unknown	secondary no	no	225	yes	yes	unknown	7	may	866	2
31	blue-collar	married	secondary	no	3653	yes	yes	cellula		21	nov	252
34	1 168 management	married	failure no tertiary	no	436	no	yes	cellula		28	jul	118
30	4 -1 technician	0 married	unknown no tertiary	no	101	yes	yes	cellula			jul	187
	3 -1	0	unknown no									
33	unemployed 1 -1	0	tertiary unknown no	no	302	yes	yes	cellula		16	apr	670
34	management -1 0	married unknown	tertiary no	no	1557	yes	yes	unknown	13	may	213	1
35	services	married	secondary	no	505	yes	yes	unknown	27	may	371	2
37	-1 0 blue-collar		no primary no	190	yes	yes	unknown	8	may	194	1	-1
40	0 unknown management		tertiary	no	-17	yes	yes	cellula		11	may	474
	1 256	1	success yes									
55	self-employed 12 -1		secondary unknown no	no	2678	no	yes	cellula		18	aug	151
42	services 1 -1	married 0	secondary unknown no	no	-91	yes	yes	cellula		5	feb	43
60	self-employed	married	primary no	362	no	yes	cellula	r	29	jul	816	6
42	-1 0 admin. married	unknown unknown		yes	yes	unknown	16	may	509	2	-1	0
57	unknown no self-employed	married	tertiary	yes	-3313	yes	yes	unknown	9	may	153	1
	-1 0	unknown	no	-		,	,					
2024-05	Time taken: 0.219 seconds, Fetched: 453 row(s) 2024-05-05 17:32:25,436 INFO CliDriver: Time taken: 0.219 seconds, Fetched: 453 row(s)											
2024-05-05 17:32:25,436 INFO conf.HiveConf: Using the default value passed in for log id: 14dc1a57-a54e-4b11-98ad-aac325bcacbc 2024-05-05 17:32:25,436 INFO session.SessionState: Resetting thread name to main												
hive	,123,130											

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

```
2024-05-05 17:50:42,035 INFO ql.Driver: OK
2024-05-05 17:50:42,036 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 17:50:42,036 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 17:50:42,046 INFO mapred.FileInputFormat: Total input files to process: 1
retired married 71188
entrepreneur married 42045
technician single 27733
management married 27359
technician married 27669
housemaid single 26965
retired married 26452
services married 26394
management divorced 26304
retired single 25824
2024-05-05 17:50:42,066 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_8:10,
Time taken: 40.449 seconds, Fetched: 10 row(s)
2024-05-05 17:50:42,077 INFO CliDriver: Time taken: 40.449 seconds, Fetched: 10 row(s)
2024-05-05 17:50:42,078 INFO conf.HiveConf: Using the default value passed in for log id: 14dc1a57-a54e-4b11-98ad-aac325bcacbc
2024-05-05 17:50:42,079 INFO session.SessionState: Resetting thread name to main
```

4. Data Aggregation and Grouping:

Calculate the average age of clients for each job category.

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

```
2024-05-05 19:55:16,731 INFO mapred.FileInputFormat: Total input files to process: 1
2024-05-05 19:55:16,743 INFO exec.ListSinkOperator: RECORDS_OUT_OPERATOR_LIST_SINK_10:13, RECORDS_OUT_INTERMEDIATE:0, admin. 39.68200836820084 blue-collar 40.15644820295983 entrepreneur 42.01190476190476 housemaid 47.339285714285715 job NULL management 40.54076367389061 retired 61.869565247391305 self-employed 41.45355191256831 services 38.57074340527578 student 26.821428571428573 technician 39.470052083333336 unemployed 49.90625 unknown 48.10526315789474 Time taken: 51.054 seconds, Fetched: 13 row(s) 2024-05-05 19:55:16,765 INFO Coliriver: Time taken: 51.054 seconds, Fetched: 13 row(s) 2024-05-05 19:55:16,765 INFO Conf. HiveConf: Using the default value passed in for log id: 14dc1a57-a54e-4b11-98ad-aac325bcacbc 2024-05-05 19:55:16,766 INFO session.SessionState: Resetting thread name to main
```

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

```
OK
2024-05-05 20:00:06,421 INFO ql.Driver: OK
2024-05-05 20:00:06,421 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 20:00:06,421 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 20:00:06,422 INFO mapred.FileInputFormat: Total input files to process: 1
2024-05-05 20:00:06,432 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_13:4,
primary yes 10
secondary yes 46
tertiary yes 17
unknown yes 3
Time taken: 24.233 seconds, Fetched: 4 row(s)
2024-05-05 20:00:06,441 INFO CliDriver: Time taken: 24.233 seconds, Fetched: 4 row(s)
2024-05-05 20:00:06,441 INFO conf.HiveConf: Using the default value passed in for log id: 14dc1a57-a54e-4b11-98ad-aac325bcacbc
2024-05-05 20:00:06,442 INFO session.SessionState: Resetting thread name to main
```

5. Complex Queries for Insights:

 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.

We directly calculate the average balance (AVG(ci.balance)) for each job category.

We also calculate the total number of clients (COUNT(*)) and the number of clients who have subscribed to a term deposit (SUM(CASE WHEN ci.y = 'yes' THEN 1 ELSE 0 END)).

We calculate the subscription rate as the percentage of subscribed clients out of the total number of clients.

The results are grouped by job category and sorted in descending order of average balance.

We limit the output to the top 5 job categories with the highest average balance using the LIMIT clause.

This query will give you 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.

Identify the top 5 job categories with the highest average balance:

```
2024-05-05 22:37:47,922 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_18:5, retired 2319.191304347826 230 54 23.47826086956522 housemaid 2083.8035714285716 112 14 12.5 management 1766.928792596594 969 131 13.519091847265221 entrepreneur 1645.125 168 15 8.928571428571429 student 1543.8214285714287 84 19 22.61904761904762 time taken: 50.115 seconds, Fetched: 5 row(s) 2024-05-05 22:37:47,929 INFO CliDriver: Time taken: 50.115 seconds, Fetched: 5 row(s)
```

 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). We directly calculate the success rate for each month by counting the total number of contacts and the number of clients who subscribed to a term deposit within each month.

We use the GROUP BY clause to group the data by month.

The results are sorted in descending order of the number of contacts, and we limit the output to only the first row, which represents the month with the highest number of contacts.

This query will give you the month with the highest number of contacts and the success rate of the campaign in that month.

```
2024-05-05 22:29:00,692 INFO ql.Driver: OK
2024-05-05 22:29:00,692 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 22:29:00,692 INFO mapred.FileInputFormat: Total input files to process : 1
2024-05-05 22:29:00,708 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_18:1,
may 1398 93 6.652360515021459
Time taken: 46.87 seconds, Fetched: 1 row(s)
2024-05-05 22:29:00,708 INFO CliDriver: Time taken: 46.87 seconds, Fetched: 1 row(s)
```

6. Correlation Analysis:

Calculate the correlation between age and balance for the clients.

In this simplified query:

- We calculate the numerator of the correlation formula: the sum of the products of age and balance minus the count of records times the average of age times the average of balance.
- 2. We calculate the denominator of the correlation formula: the square root of the difference between the sum of the squares of age and the count times the square of the average of age, multiplied by the square root of the difference between the sum of the squares of balance and the count times the square of the average of balance.
- 3. We divide the numerator by the denominator to get the correlation coefficient.

This query should provide you with the correlation coefficient between age and balance for the clients in your dataset.

```
hive> SELECT (

> SUM(age*balance)-COUNT(*)*AVG(age)*AVG(balance)

> )/(

> SQRT(SUM(POW(age,2))-COUNT(*)*POW(AVG(age),2))*

> SQRT(SUM(POW(balance,2))-COUNT(*)*POW(AVG(balance),2))

> )AS correlation_age_balance

> FROM client_info;

2024-05-05 21:03:10,887 INFO conf.HiveConf: Using the default value passed in for log id: 2fe3dd08-4d00-4eec-a267-7f5cbbf8629c
2024-05-05 21:03:10,887 INFO session.SessionState: Updating thread name to 2fe3dd08-4d00-4eec-a267-7f5cbbf8629c main
2024-05-05 21:03:10,889 INFO ql.Driver: Compiling command(queryId=titov_20240505210310_5d5c2c37-75d7-4d3b-a6bf-f37b1c95a93c): SELECT (
SUM(age*balance)-COUNT(*)*AVG(age)*AVG(balance)

)/(
SQRT(SUM(POW(age,2))-COUNT(*)*POW(AVG(age),2))*
SQRT(SUM(POW(balance,2))-COUNT(*)*POW(AVG(balance),2))

AS correlation_age_balance

FROM client_info
```

```
2024-05-05 21:03:56,534 INFO ql.Driver: OK
2024-05-05 21:03:56,537 INFO ql.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 21:03:56,550 INFO mapred.FileInputFormat: Total input files to process : 1
2024-05-05 21:03:56,637 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_11:1,
0.0835552815119586
Time taken: 45.649 seconds, Fetched: 1 row(s)
2024-05-05 21:03:56,733 INFO CliDriver: Time taken: 45.649 seconds, Fetched: 1 row(s)
```

7. Trend Analysis:(TO BE DONE)

• Analyze the year-over-year trend in the number of clients contacted.

In this query:

We use the SUBSTRING function to extract the year from the month column.

Assuming the month column is in the format "YYYY-MM", we extract the first four characters to get the year.

We count the number of clients contacted (COUNT (*)) for each year.

We group the results by the extracted year.

Finally, we order the results by year to see the trend over time.

This query will give you the year-over-year trend in the number of clients contacted based on the data in your client info table.

```
2024-05-05 21:12:45,360 INFO mapred.FileInputFormat: Total input files to process : 1
2024-05-05 21:12:45,377 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_16:13, 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
Time taken: 100.264 seconds, Fetched: 13 row(s)
```

8. Anomaly Detection:

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

We calculate the average yearly balance for each education level. We also calculate the overall average yearly balance for each year using the window function AVG(AVG(balance)) OVER (PARTITION BY SUBSTRING (month, 1, 4)).

We calculate the z-score for each average yearly balance within each education level by subtracting the overall average yearly balance and dividing by the standard deviation, both calculated over the partition of years.

By examining the z-scores, you can identify any unusual patterns in the average yearly balance across different education levels. A z-score significantly higher or lower than zero indicates that the average yearly balance for a particular education level is unusually high or low compared to the overall average yearly balance.

```
hive> 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,

> (AVG(balance)-AVG(AVG(balance)) OVER (PARTITION BY SUBSTRING(month,1,4)))/STDDEV(balance) OVER(PARTITION BY SUBSTRING(month,1,4)) AS z_score

> FROM client_info

> GROUP BY SUBSTRING(month,1,4),education;
```

9. Advanced Analysis:

 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.

In this query:

- 1. We group the data by the poutcome column to analyze the impact of previous campaign outcomes.
- 2. We count the total number of clients (total_clients) and the number of clients who subscribed to term deposits (subscribed_clients) for each poutcome category.
- 3. We calculate the subscription rate (subscription_rate) as the percentage of clients who subscribed to term deposits out of the total number of clients for each poutcome category.

This query will provide you with the subscription rate for each poutcome category, allowing you to analyze the impact of previous campaign outcomes on the current campaign's success.

```
2024-05-05 21:20:32,943 INFO q1.Driver: OK
2024-05-05 21:20:32,943 INFO q1.Driver: Concurrency mode is disabled, not creating a lock manager
2024-05-05 21:20:32,948 INFO mapred.FileInputFormat: Total input files to process : 1
2024-05-05 21:20:32,956 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_11:5, failure 490 63 12.857142857142856
other 197 38 19.289340101522843
poutcome 1 0 0.0
success 129 83 64.34108527131784
unknown 3705 337 9.095816464237517
Time taken: 24.245 seconds, Fetched: 5 row(s)
2024-05-05 21:20:32,967 INFO CliDriver: Time taken: 24.245 seconds, Fetched: 5 row(s)
```

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

In this query:

- 1. We group the data by the subscription status ($_{Y}$ column), which indicates whether the client subscribed to a term deposit.
- 2. We calculate the average contact duration (avg_contact_duration) for each group separately.

This query will provide you with the average contact duration for clients who subscribed and who did not subscribe to a term deposit, allowing you to compare the contact durations between the two groups.

```
nive> SELECT y AS sub_status,
> AVG(duration) AS avg_contact_duration
> FROM
> client_info
> GROUP BY y
> ORDER BY AVG(duration) DESC;
```

```
2024-05-05 21:26:25,506 INFO mapred.FileInputFormat: Total input files to process: 1
2024-05-05 21:26:25,511 INFO exec.ListSinkOperator: RECORDS_OUT_INTERMEDIATE:0, RECORDS_OUT_OPERATOR_LIST_SINK_16:3,
yes 552.7428023032629
no 226.3475
y NULL
Time taken: 54.453 seconds, Fetched: 3 row(s)
2024-05-05 21:26:25,517 INFO CliDriver: Time taken: 54.453 seconds, Fetched: 3 row(s)
2024-05-05 21:26:25,517 INFO conf.HiveConf: Using the default value passed in for log id: 2fe3dd08-4d00-4eec-a267-7f5cbbf8629c
2024-05-05 21:26:25,518 INFO session.SessionState: Resetting thread name to main
```

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.