

INFO-I-535 Final Project Report

Name: Vasu Sharma

Email: vasushar@iu.edu

Semester: Fall 2023

Faculty: Inna Kouper

Table of Contents

1.	<u>Introduction</u>	2-3
2.	<u>Background</u>	3
3.	<u>Methodology</u>	
3.1.	<u>Overview</u>	3-4
3.2.	<u>Google Cloud Platform</u>	4-9
3.3.	<u>Pandas, MongoDB, Docker, JetStream</u>	10-15
4.	<u>Results</u>	
4.1.	<u>Visualizations on Google Cloud Platform</u>	16-18
4.2.	<u>Visualizations using Pandas, MongoDB</u>	19-22
5.	<u>Discussion</u>	
5.1.	<u>Visualization discussion</u>	23-24
5.2.	<u>Barrier and Challenges faced</u>	24
6.	<u>Conclusion</u>	24-25
7.	<u>References</u>	25

1.Introduction:

Credit Card is a modern-day currency which is extensively used for various financial transactions. It provides ease of doing business to both consumers and businesses. It plays a vital role in mitigating the need to carry currency and is called plastic money. Credit cards also provide luxury to people to spend money without having it and can be repaid at a later point in time, increasing the spending capacity of a consumer and helping businesses to thrive. Hence credit cards have become an integral part of today's financial system.

One of the major issues that financial institutions are facing in today's rapidly changing environment is understanding and mitigating customer churn in order to sustain a successful business. Churn is basically the phenomenon where the customers discontinue their credit card services. This can have a profound impact on bank's revenue and customer satisfaction.

The different features present in the dataset are as follows:

1. CLIENTNUM
2. Attrition_Flag
3. Customer_Age
4. Gender
5. Dependent_count
6. Education_Level
7. Marital_Status
8. Income_Category
9. Card_Category
10. Months_on_book
11. Total_Relationship_Count
12. Months_Inactive_12_mon
13. Contacts_Count_12_mon
14. Credit_Limit
15. Total_Revolving_Bal
16. Avg_Open_To_Buy
17. Total_Amt_Chng_Q4_Q1
18. Total_Trans_Amt
19. Total_Trans_Ct

20. Total_Ct_Chng_Q4_Q1
21. Avg_Utilization_Ratio
22. Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_1
23. Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_2

2.Background:

We have decided to work on analyzing credit card customers because of my interest and experience in Finance Domain. We wanted to understand the various aspects related to credit card customers and also wanted to understand factors that could potentially cause the credit card customers of a financial institution to drift away. This can help financial institution to know what can be done in order to keep their customers intact. These are some of the reasons because of which we selected this project.

In this project, we have used various methods to analyze and determine the factors that could potentially cause financial institutions such as banks to lose credit card customers. We have tried to analyze the data on various platforms and tried to use various methods required to analyze dataset having information about credit card customers.

3.Methodology:

Overview:

We have analyzed this data using Pandas, MongoDB, Docker, Jetstream, Google Cloud Platform and have used various visualization tools and libraries in order to create visualizations for this dataset.

We have downloaded a credit card customer dataset from Kaggle which has all the features as mentioned above. Once the dataset is downloaded, it is imported on Google Cloud Platform and created a database and a table in which this dataset is imported. Once the database and table is created, we cleaned the dataset using

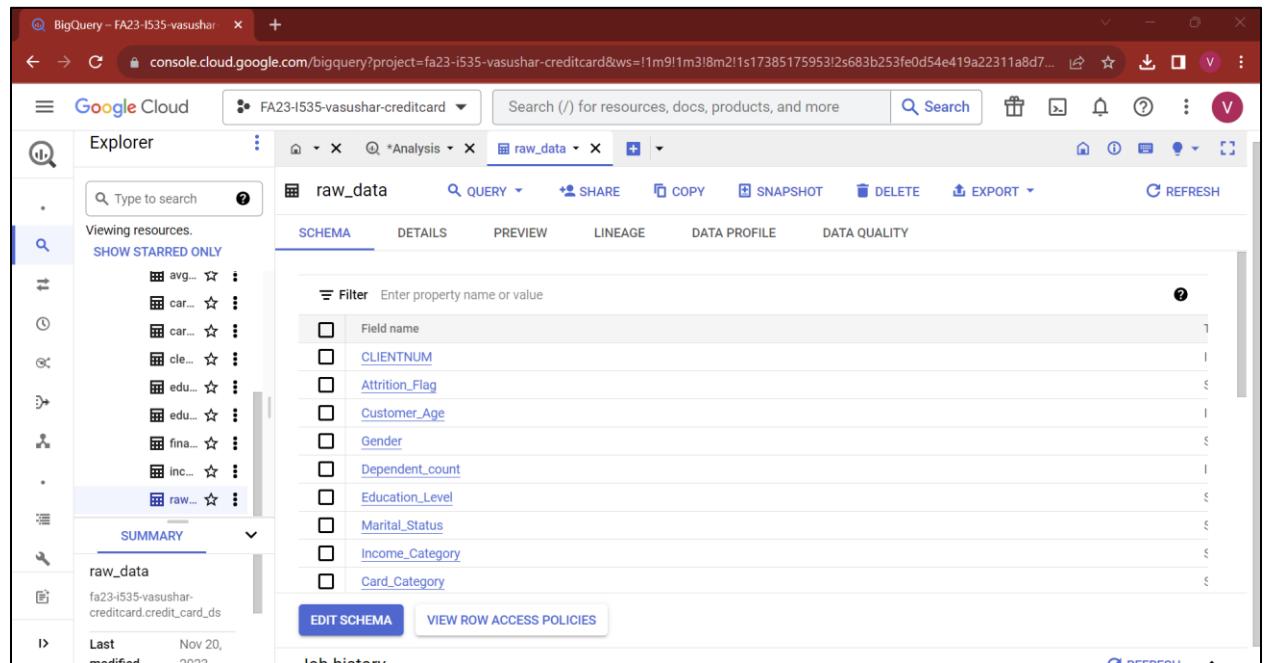
various operations then queried this dataset using Big Query and then tried to derive various conclusions from this dataset. All these conclusions have been visualized on this platform.

We have also created a No SQL database on MongoDB and imported cleaned dataset on MongoDB to create the database. Then we created various visualizations by inserting required data in new collections and then retrieving that information to create visualizations from those collections on MongoDB.

Google Cloud Platform:

I. Selection of Platform and DataSet

- Google Cloud Platform(GCP) is used to clean, analyze and visualize the dataset
- DataSet is selected from Kaggle under category of Credit Card Customer.



The screenshot shows the Google Cloud BigQuery interface. The top navigation bar includes the project name 'FA23-i535-vasushar' and a search bar. Below the navigation is a toolbar with various icons for search, refresh, and data management. The main area is divided into sections: 'Explorer' on the left listing datasets like 'raw_data', 'raw_data', and 'fa23-i535-vasushar-creditcard.credit_card_ds'; 'Analysis' (selected), 'raw_data' (selected); and 'SCHEMA' (selected). The 'SCHEMA' tab displays a list of columns with checkboxes for filtering: Field name, CLIENTNUM, Attrition_Flag, Customer_Age, Gender, Dependent_count, Education_Level, Marital_Status, Income_Category, and Card_Category. At the bottom of the schema view are 'EDIT SCHEMA' and 'VIEW ROW ACCESS POLICIES' buttons.

II. Cleaning of Data:

1. The Dataset consists of 23 features. All of these features are not necessary to analyze the dataset.

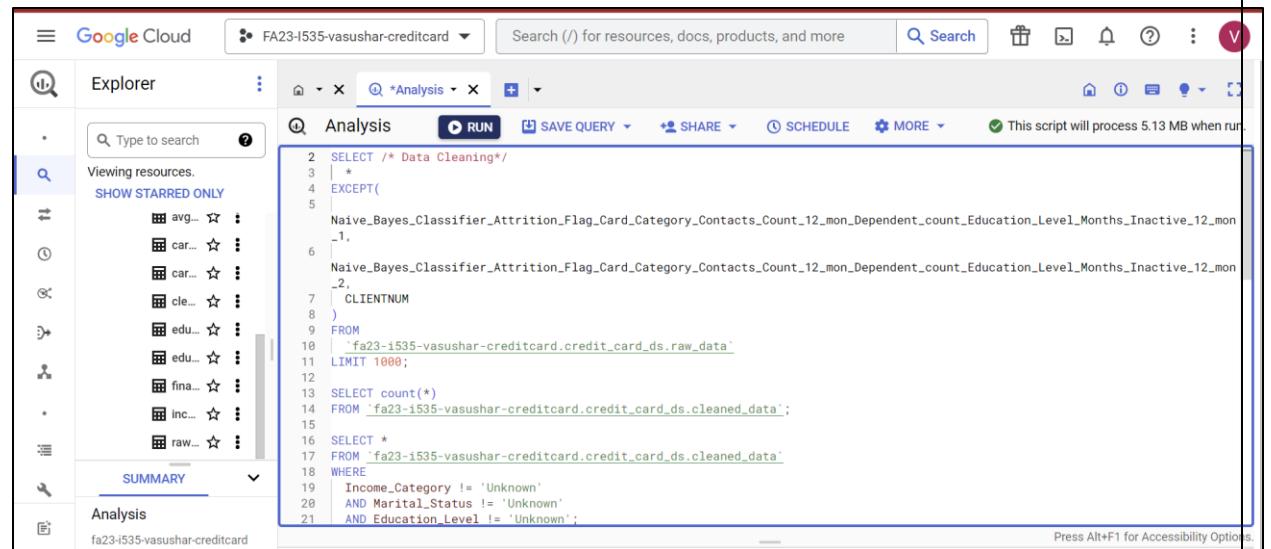
2. Features such as

Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_1,

Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_2,

and CLIENTNUM have been removed and stored in cleaned_data table by using features such as EXCEPT in Big Query on GCP.

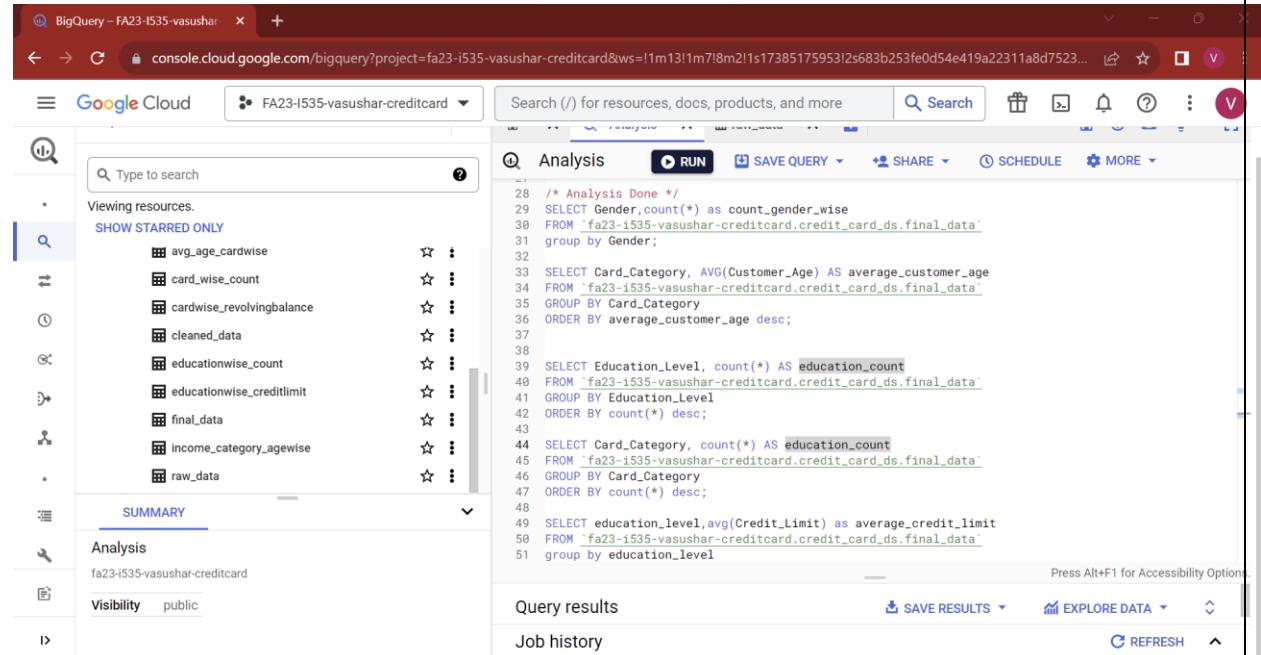
3. Columns such as Income_Category, Marital Status and Education_level contained unknown values. Such rows with unknown values for these columns have been removed.



The screenshot shows the Google Cloud BigQuery interface. The left sidebar has 'Explorer' selected, showing a list of datasets and tables. The main area is titled 'Analysis' with a 'RUN' button. A query is being typed into the editor:

```
2 SELECT /* Data Cleaning*/
3   *
4 EXCEPT(
5   Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon
6   _1,
7   Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon
8   _2,
9   FROM
10   `fa23-i535-vasushar-creditcard.credit_card_ds.raw_data`
11   LIMIT 1000;
12
13   SELECT count(*)
14   FROM `fa23-i535-vasushar-creditcard.credit_card_ds.cleaned_data`;
15
16   SELECT *
17   FROM `fa23-i535-vasushar-creditcard.credit_card_ds.cleaned_data`
18   WHERE
19   Income_Category != 'Unknown'
20   AND Marital_Status != 'Unknown'
21   AND Education_Level != 'Unknown';
```

- Once all the above operations done, data is stored in final_data table.



The screenshot shows the Google Cloud BigQuery interface. On the left, the sidebar displays a list of resources under 'Analysis' and 'Summary'. The main area shows a query titled 'Analysis' with the following SQL code:

```

28 /* Analysis Done */
29 SELECT Gender, count(*) as count_gender_wise
30 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
31 group by Gender;
32
33 SELECT Card_Category, AVG(Customer_Age) AS average_customer_age
34 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
35 GROUP BY Card_Category
36 ORDER BY average_customer_age desc;
37
38
39 SELECT Education_Level, count(*) AS education_count
40 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
41 GROUP BY Education_Level
42 ORDER BY count(*) desc;
43
44 SELECT Card_Category, count(*) AS education_count
45 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
46 GROUP BY Card_Category
47 ORDER BY count(*) desc;
48
49 SELECT education_level,avg(Credit_Limit) as average_credit_limit
50 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
51 group by education_level

```

The interface includes tabs for 'Analysis', 'Query results', and 'Job history'. There are also buttons for 'RUN', 'SAVE QUERY', 'SHARE', 'SCHEDULE', 'MORE', 'SAVE RESULTS', 'EXPLORE DATA', and 'REFRESH'.

III. Calculating Total Revolving balance Card Category Wise

A number of aspects of the dataset have been analyzed and stored in various tables. One of them is as follows:

- We checked the total revolving balance card category wise and renamed the column as sum_total_revolving_balance .
- We also used GROUP BY, ORDER BY methods in order to order the result obtained in descending order.
- Then this result set is stored in a separate table cardwise_revolvingbalance which is then used to create visualizations in further analysis.

The screenshot shows the Google Cloud BigQuery interface. On the left, the Explorer sidebar lists various datasets and tables, including 'raw_data' and 'cardwise_revolvingbalance'. The main area displays an 'Analysis' query:

```

51 group by education_level
52 ORDER BY average_credit_limit desc;
53
54 SELECT
55   Card_Category,
56   SUM(Total_Revolving_Bal) AS sum_total_revolving_balance
57 FROM
58   `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
59 GROUP BY
60   Card_Category
61 ORDER BY sum_total_revolving_balance desc;
62
63 SELECT

```

The 'Query results' section shows the following data:

Card_Category	sum_total_revolving
Blue	7678897
Silver	464701
Gold	108451
Platinum	15030

Process to store this result set is shown as below:

The screenshot shows the Google Cloud BigQuery interface with the 'Export to BigQuery Table' dialog open. The 'Destination' section is configured with:

- Project ***: fa23-i535-vasushar-creditcard
- Dataset ***: credit_card_ds
- Table ***: card_wise_revolvingbal

The 'Advanced options' section is collapsed. The 'Query results' section shows the same data as the previous screenshot:

Card_Category	sum_total_revolving
Blue	7678897
Silver	464701
Gold	108451
Platinum	15030

The screenshot shows the Google Cloud BigQuery interface. In the left sidebar, under 'Explorer', there is a list of resources. One item, 'cardwise_revolvingbalance', is highlighted with a yellow box. The main area displays an 'Analysis' query:

```

51 group by education_level
52 ORDER BY average_credit_limit desc;
53
54 SELECT
55   Card_Category,
56   SUM(Total_Revolving_Bal) AS sum_total_revolving_balance
57   FROM
58   `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
59   GROUP BY
60   Card_Category
61   ORDER BY sum_total_revolving_balance desc;
62
63 SELECT

```

Below the query, the 'Query results' section shows the following data:

Card_Category	sum_total_revolving
Blue	7678897
Silver	464701

IV. Analyzing Existing and attrited customers traits

- We have tried to compare various features of Existing and attrited customers.
- After doing careful analysis, it can be concluded that there is a significant difference between the two categories in average total transaction count, average credit limit and average total revolving balance.

The screenshot shows the Google Cloud BigQuery interface. In the left sidebar, under 'Explorer', there is a list of resources. The 'Analysis' section contains a query:

```

81   AVG(Credit_Limit) AS avg_total_Credit_Limit,
82   AVG(Total_Revolving.Bal) AS avg_total_revolving_balance,
83   AVG(Months_Inactive_12_mon) AS avg_inactive_month,
84   AVG(Months_on_book) AS month_on_books,
85   AVG(Total_Relationship_Count) AS avg_relationship_count,
86   AVG(Avg_Open_To_Buy) AS avg_open_to_buy,
87   FROM
88   `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
89   GROUP BY
90   Attrition_Flag;
91
92
93 SELECT
94   Attrition_Flag

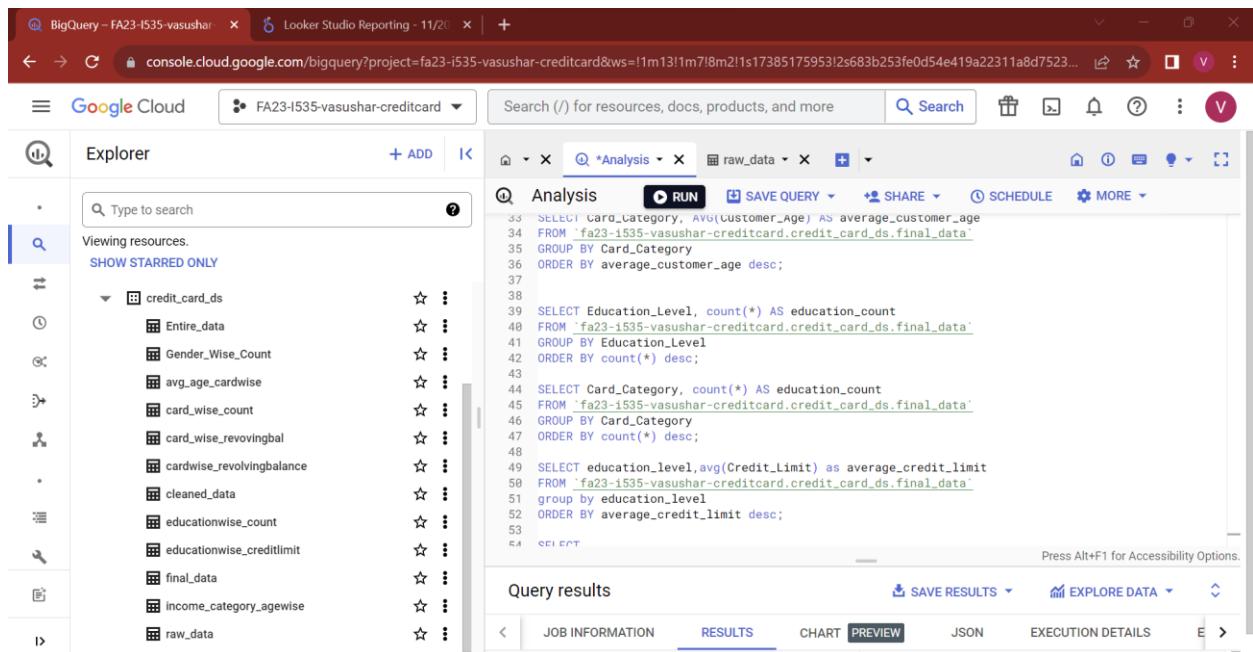
```

The 'Query results' section shows the following data:

Attrition_Flag	Avg_Age	avg_total_tran_cnt	avg_total_Credit_Lim	avg_total_revolving	avg_inactive_month	month_on_books	avg_relationship_cou
Existing Customer	46.31735924932...	68.17644101876...	8555.099195710...	1260.589979892...	2.276977211796...	35.94939678284...	3.916554959785...
Attrited Customer	46.51033243486...	44.80772686433...	8158.579964061...	668.3539982030...	2.694519317160...	36.15274034141...	3.298292902066...

V. Extracting various other insights and creating sub tables for these query results

- In a similar way as described above, we have tried to analyze and correlate various other features of the dataset and tried to extract relationships among those features.
- We have used various functions provided Google Cloud platform Big Query such as aggregate functions such as count, avg, sum and GROUP BY, ORDER BY clauses in order to create various sub tables
- Later on, visualizations are created on these sub table only for better visualizations.



The screenshot shows the Google Cloud BigQuery interface. The top navigation bar includes tabs for 'BigQuery - FA23-I535-vasushar-' and 'Looker Studio Reporting - 11/20'. Below the navigation is a search bar and a 'Search' button. The main area is divided into two sections: 'Explorer' on the left and 'Analysis' on the right.

Explorer: Shows a tree view of datasets and tables under 'credit_card_ds'. Nodes include 'Entire_data', 'Gender_Wise_Count', 'avg_age_cardwise', 'card_wise_count', 'card_wise_revovingbal', 'cardwise_revolvingbalance', 'cleaned_data', 'educationwise_count', 'educationwise_creditlimit', 'final_data', 'income_category_agewise', and 'raw_data'. Each node has a star icon and a more options menu.

Analysis: The current tab is 'Analysis'. It contains a query editor with the following SQL code:

```
33 SELECT Card_Category, AVG(Customer_Age) AS average_customer_age
34 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
35 GROUP BY Card_Category
36 ORDER BY average_customer_age desc;
37
38
39 SELECT Education_Level, count(*) AS education_count
40 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
41 GROUP BY Education_Level
42 ORDER BY count(*) desc;
43
44 SELECT Card_Category, count(*) AS education_count
45 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
46 GROUP BY Card_Category
47 ORDER BY count(*) desc;
48
49 SELECT education_level, avg(Credit_Limit) as average_credit_limit
50 FROM `fa23-i535-vasushar-creditcard.credit_card_ds.final_data`
51 group by education_level
52 ORDER BY average_credit_limit desc;
53
54
```

The 'Query results' section below the code is currently empty. Navigation buttons at the bottom include 'JOB INFORMATION', 'RESULTS' (which is selected), 'CHART', 'PREVIEW', 'JSON', and 'EXECUTION DETAILS'.

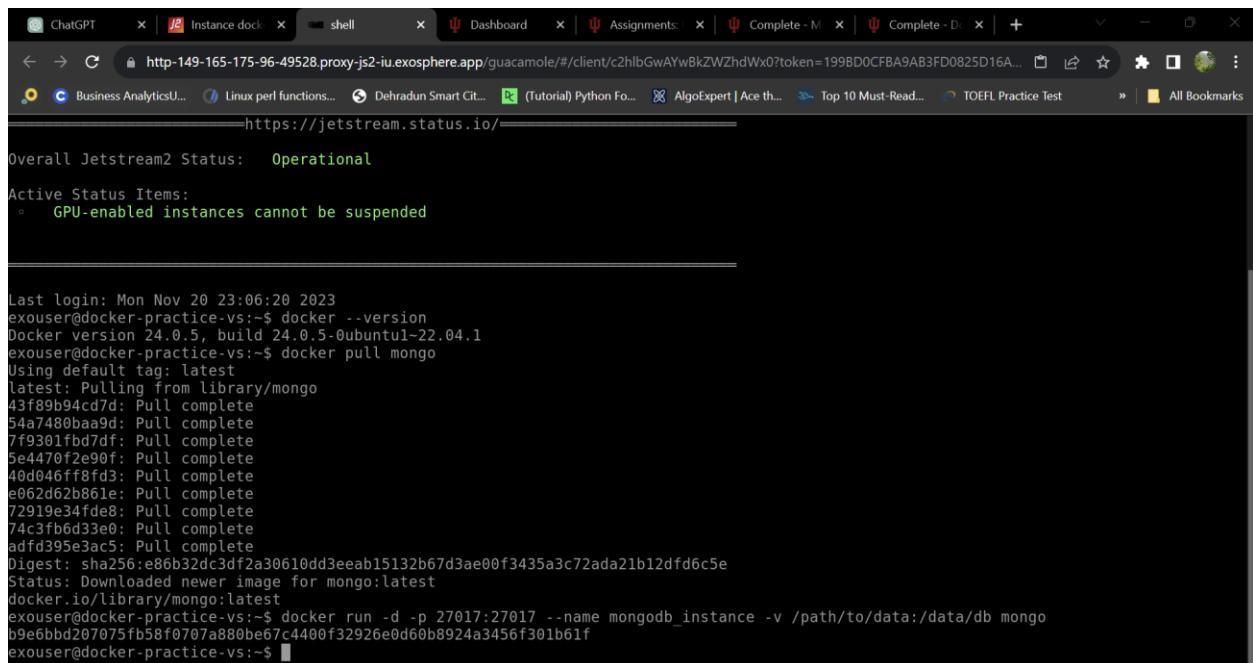
JetStream, Docker, Pandas and MongoDB:

Step 1: Creating a JetStream Docker instance with MonogoDb installed

First we created a docker instance on Jetstream. Jetstream enables higher scalability and cost effectiveness. We opted for Docker and MongoDB for managing this dataset due to their scalability and cost-effectiveness. These are key attributes essential for establishing a robust and sustainable solution in the long term. We have pulled latest MongoDB image on docker to run MongoDB on this instance.

The screenshot shows the Jetstream2 Instances page. At the top, it displays resource usage statistics: 85 of 100 total instances used, 225 of 1,164 total cores used, and 755.7 GB of 4.6 TB RAM used. Below this, a search bar shows a filter for 'Created by me' and a result count of '2 instances filtered from 85 total'. Two instances are listed: 'Spark_Practice_VS' (m3.small flavor) and 'docker-practice-vs' (m3.medium flavor). Each instance entry includes a 'Connect to' button, an IP address (149.165.175.142 for the first and 149.165.175.96 for the second), and a delete icon.

The screenshot shows the Jetstream2 Instance details page for 'docker-practice-vs'. The top navigation bar includes 'Home', 'Project CIS230181', 'Instances', and 'Instance docker-practice-vs'. It also features 'Remove Allocation' and 'Create' buttons. The main section displays the instance's status as 'Ready' with a green button. Below this, the 'Info' card provides details: created 2 months ago, by user vsharma4@access-ci.org, from image 9ff7bc85...5f605dfee8be, and flavor m3.medium. It also shows a burn rate of 8.00 SU/hour. A 'Resource Usage' section at the bottom shows CPU, RAM, Root Disk, and Total Disk usage.



Overall Jetstream2 Status: Operational

Active Status Items:

- GPU-enabled instances cannot be suspended

```
Last login: Mon Nov 20 23:06:20 2023
exouser@docker-practice-vs:~$ docker --version
Docker version 24.0.5, build 24.0.5-0ubuntu1-22.04.1
exouser@docker-practice-vs:~$ docker pull mongo
Using default tag: latest
latest: Pulling from library/mongo
43f89b94cd7d: Pull complete
54a7480baa9d: Pull complete
7f9301fb7df: Pull complete
5e4470f2e90f: Pull complete
40d046ff8fd3: Pull complete
e062d62b861e: Pull complete
72919e34fde8: Pull complete
74c3fb6bd33e0: Pull complete
adfd395e3ac5: Pull complete
Digest: sha256:e86b32dc3df2a30610dd3eeab15132b67d3ae00f3435a3c72ada21b12dfd6c5e
Status: Downloaded newer image for mongo:latest
docker.io/library/mongo:latest
exouser@docker-practice-vs:~$ docker run -d -p 27017:27017 --name mongodb_instance -v /path/to/data:/data/db mongo
b9e6bbd207075fb58f0707a880be67c4400f32926e0d60b8924a3456f301b61f
exouser@docker-practice-vs:~$
```

Step 2: Reading the dataset

- Now we read the raw data into pandas and extracted various information regarding details about the dataset. Once we have all the information, we proceed to do data cleaning of the dataset.
- We have used various functions such as head(), info(), describe() to get more information about the dataset.
- All this information helps to do data cleaning in a better and efficient manner.

The screenshot shows a Jupyter Notebook interface in VS Code. The code cell at the top imports pandas, numpy, matplotlib.pyplot, seaborn, and MongoClient. The next cell reads a CSV file named 'BankChurners.csv' into a DataFrame named credit_df. The third cell displays the first five rows of the DataFrame.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from pymongo import MongoClient

credit_df = pd.read_csv('BankChurners.csv')

credit_df.head(5)
```

The screenshot shows a Jupyter Notebook interface in VS Code. The code cell at the top prints the list of columns in the credit_df DataFrame. The second cell shows the shape of the DataFrame. The third cell displays the summary information for the DataFrame.

```
credit_df.columns

credit_df.shape

credit_df.info()
```

Step 3: Data Cleaning

- All the data cleaning steps that have been on GCP, are repeated in Pandas.
- This helps to analyze and make visualizations and derive conclusions that are more insightful and meaningful.

The screenshot shows a Jupyter Notebook interface with the following details:

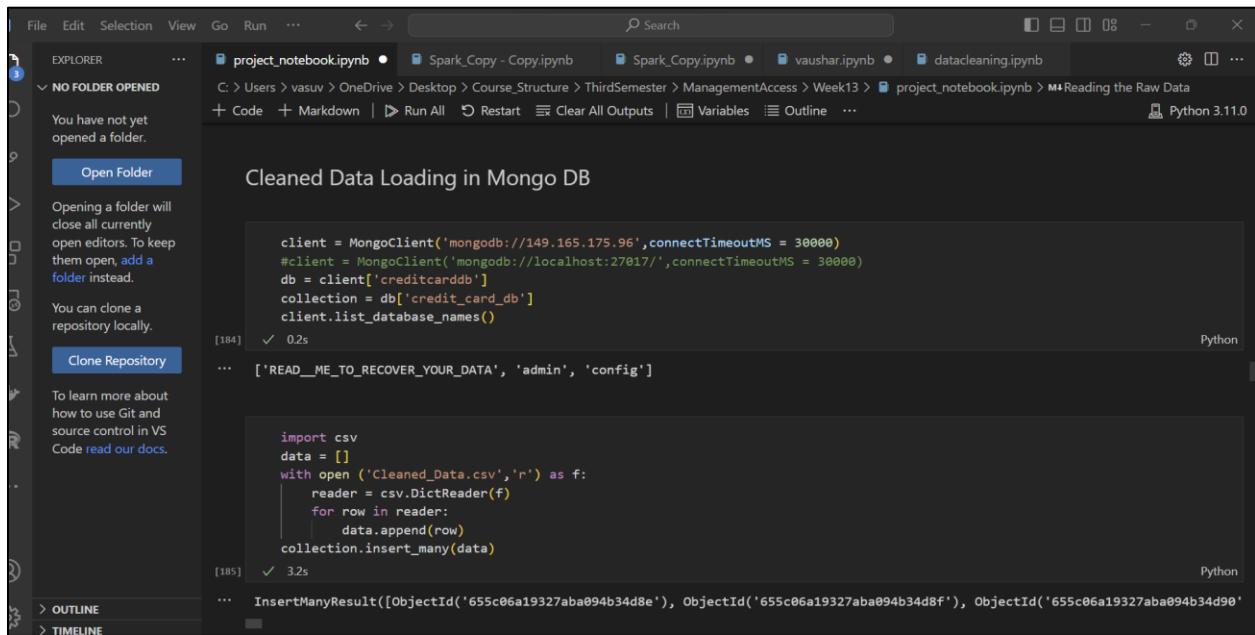
- File Bar:** File, Edit, Selection, View, Go, Run, ...
- Search Bar:** Search
- Toolbar:** Explorer, Project Notebook.ipynb, Spark_Copy - Copy.ipynb, Spark_Copy.ipynb, vaushar.ipynb, datacleaning.ipynb, Python 3.11.0
- Left Sidebar:**
 - Explorer: NO FOLDER OPENED
 - Open Folder
 - Opening a folder will close all currently open editors. To keep them open, add a folder instead.
 - Clone Repository
 - To learn more about how to use Git and source control in VS Code [read our docs](#).
- Code Cell 1:** credit_df.drop('Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Mont', axis=1,inplace = True)
credit_df.drop('Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Mont', axis=1,inplace = True)
credit_df.drop('CLIENTNUM',axis = 1,inplace = True)
#credit_df.drop('_id',axis = 1,inplace = True)
- Code Cell 2:** [159] ✓ 0.0s credit_df.columns
- Code Cell 3:** [160] ✓ 0.0s ... Index(['Attrition_Flag', 'Customer_Age', 'Gender', 'Dependent_count', 'Education_Level', 'Marital_Status', 'Income_Category', 'Card_Category', 'Months_on_book', 'Total_Relationship_Count', 'Months_Inactive_12_mon', 'Contacts_Count_12_mon', 'Credit_Limit', 'Total_Revolving_Bal', 'Avg_Open_To_Buy', 'Total_Amt_Chng_Q4_Q1', 'Total_Trans_Amt', 'Total_Trans_Ct', 'Total_Ct_Chng_Q4_Q1', 'Avg_Utilization_Ratio'], dtype='object')
- Code Cell 4:** [161] ✓ 0.0s print(credit_df.isna().sum())
print(credit_df.isnull().sum())
- Bottom Navigation:** OUTLINE, TIMELINE

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, ...
- Search Bar:** Search
- Toolbar:** Explorer, Project Notebook.ipynb, Spark_Copy - Copy.ipynb, Spark_Copy.ipynb, vaushar.ipynb, datacleaning.ipynb, Python 3.11.0
- Left Sidebar:**
 - Explorer: NO FOLDER OPENED
 - Open Folder
 - Opening a folder will close all currently open editors. To keep them open, add a folder instead.
 - Clone Repository
 - To learn more about how to use Git and source control in VS Code [read our docs](#).
- Code Cell 1:** income_unknown = credit_df['Income_Category'].eq('Unknown').sum()
marital_status_unknown = credit_df['Marital_Status'].eq('Unknown').sum()
education_level_unknown = credit_df['Education_Level'].eq('Unknown').sum()
- Code Cell 2:** [162] ✓ 0.0s print(f'Income_category with Unknown status = {income_unknown}, Marital_Status with Unknown status = {marital_status_unknown} , Education_Level with Unknown status = {education_level_unknown}')
- Code Cell 3:** [163] ✓ 0.0s credit_df.shape
- Code Cell 4:** [164] ✓ 0.0s ... (7081, 20)
- Code Cell 5:** [165] ✓ 0.0s credit_df.columns
- Code Cell 6:** ... Index(['Attrition_Flag', 'Customer_Age', 'Gender', 'Dependent_count', 'Education_Level', 'Marital_Status', 'Income_Category', 'Card_Category', 'Months_on_book', 'Total_Relationship_Count', 'Months_Inactive_12_mon', 'Contacts_Count_12_mon', 'Credit_Limit', 'Total_Revolving_Bal', 'Avg_Open_To_Buy', 'Total_Amt_Chng_Q4_Q1', 'Total_Trans_Amt', 'Total_Trans_Ct', 'Total_Ct_Chng_Q4_Q1', 'Avg_Utilization_Ratio'],
- Bottom Navigation:** OUTLINE, TIMELINE

Step 4: Data Loading

- Once the dataset is cleaned, it is loaded into MongoDB hosted in Docker on JetStream in order to create a NO SQL document based database.
- Later, we have created various collections within this database to create meaningful visualizations.



The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, ...
- Toolbar:** Search, Python 3.11.0
- File List:** project_notebook.ipynb, Spark_Copy - Copy.ipynb, Spark_Copy.ipynb, vaushar.ipynb, datacleaning.ipynb
- Path:** C:\Users\vasuv\OneDrive\Desktop\Course Structure\ThirdSemester\ManagementAccess\Week13>project_notebook.ipynb>M4:Reading the Raw Data
- Code Cell 1 (Index 184):** Python code to connect to MongoDB and insert data from a CSV file.

```
client = MongoClient('mongodb://149.165.175.96',connectTimeoutMS = 30000)
#client = MongoClient('mongodb://localhost:27017/',connectTimeoutMS = 30000)
db = client['creditcarddb']
collection = db['credit_card_db']
client.list_database_names()

[184]  ✓  0.2s
...  ['READ__ME_TO_RECOVER_YOUR_DATA', 'admin', 'config']
```
- Code Cell 2 (Index 185):** Python code to read a CSV file and insert its contents into the MongoDB collection.

```
import csv
data = []
with open ('Cleaned_Data.csv','r') as f:
    reader = csv.DictReader(f)
    for row in reader:
        data.append(row)
collection.insert_many(data)

[185]  ✓  3.2s
...  InsertManyResult([ObjectId('655c06a19327aba094b34d8e'), ObjectId('655c06a19327aba094b34d8f'), ObjectId('655c06a19327aba094b34d90')])
```

The database created in MongoDB is can be seen below:

Database: creditcard

Collection: credit_card_db

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, ...
- Search Bar:** Search
- Toolbar:** Open Folder, Clone Repository, Python 3.11.0
- Left Sidebar:**
 - EXPLORER:** NO FOLDER OPENED, Open Folder, Clone Repository.
 - Opening a folder will close all currently open editors. To keep them open, add a folder instead.
 - You can clone a repository locally.
 - To learn more about how to use Git and source control in VS Code read our docs.
- Code Cell 1:** [186] collection_names = db.list_collection_names()

```
# Print the list of collection names
print(collection_names)
```

Output: [186] ✓ 0.0s

```
... ['credit_card_db']
```
- Code Cell 2:** [187] collection = db['credit_card_db']
result = collection.find()
for document in result:
 print(document)

mongo_document = collection.find()
credit_df = pd.DataFrame(list(mongo_document))
- Output 187:** ✓ 1.6s

```
... {'_id': ObjectId('655c06a19327aba094b34d8e'), ':': '0', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '45', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d8f'), ':': '1', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '49', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34d90'), ':': '2', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '51', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d91'), ':': '4', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '40', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d92'), ':': '5', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '44', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d93'), ':': '8', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '37', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d94'), ':': '9', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '48', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d95'), ':': '12', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '56', 'Gender': 'M'}
```

The database created in Mongo DB is shown below:

```
... {'_id': ObjectId('655c06a19327aba094b34d8e'), ':': '0', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '45', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d8f'), ':': '1', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '49', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34d90'), ':': '2', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '51', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d91'), ':': '4', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '40', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d92'), ':': '5', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '44', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d93'), ':': '8', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '37', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d94'), ':': '9', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '48', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d95'), ':': '12', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '56', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d96'), ':': '14', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '57', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34d97'), ':': '16', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '48', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d98'), ':': '18', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '61', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d99'), ':': '20', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '47', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d9a'), ':': '21', 'Attrition_Flag': 'Attrited Customer', 'Customer_Age': '62', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34d9b'), ':': '22', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '41', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d9c'), ':': '25', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '41', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34d9d'), ':': '29', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '47', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d9e'), ':': '31', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '53', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34d9f'), ':': '32', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '41', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34da0'), ':': '33', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '53', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34da1'), ':': '34', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '58', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34da2'), ':': '35', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '55', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34da3'), ':': '36', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '55', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34da4'), ':': '37', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '42', 'Gender': 'F'
{'_id': ObjectId('655c06a19327aba094b34da5'), ':': '40', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '45', 'Gender': 'M'
{'_id': ObjectId('655c06a19327aba094b34da6'), ':': '42', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '50', 'Gender': 'F'
...
{'_id': ObjectId('655c06a19327aba094b36933'), ':': '10121', 'Attrition_Flag': 'Existing Customer', 'Customer_Age': '56', 'Gender': 'M'}
```

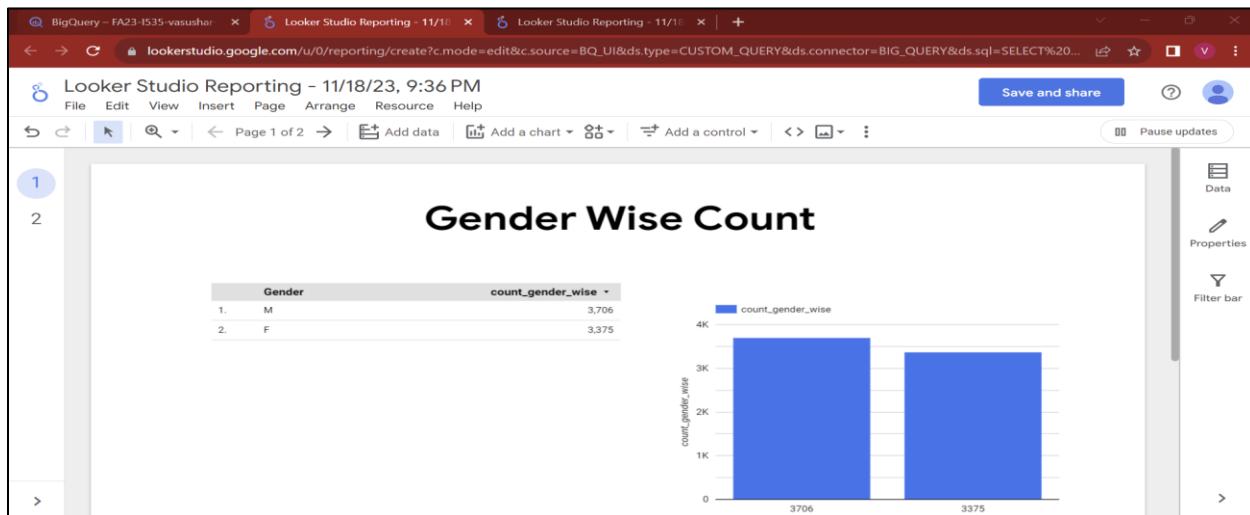
4. Results:

Google Cloud Platform

Visualizations created are as follows:

Gender Wise Count in dataset:

- Number of males credit card customer are higher as compared to females credit card customers in the dataset.



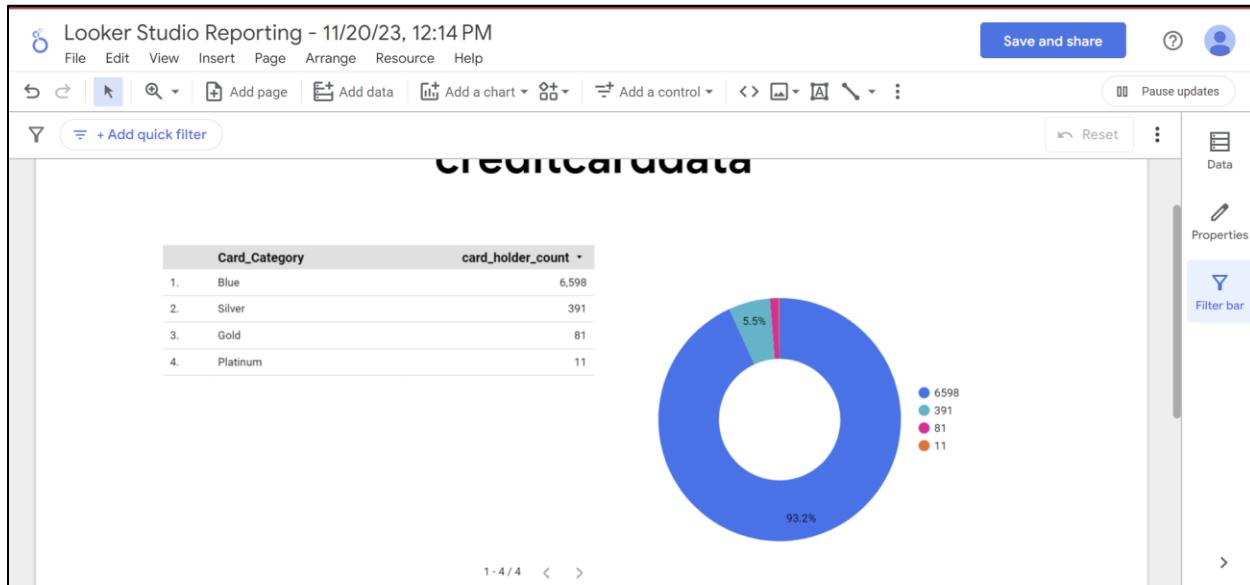
Card Category Wise Total revolving Balance in dataset:

- As can be seen from the pie chart, Blue credit card category has the highest total revolving balance as compared to Platinum Card Category with least revolving balance.



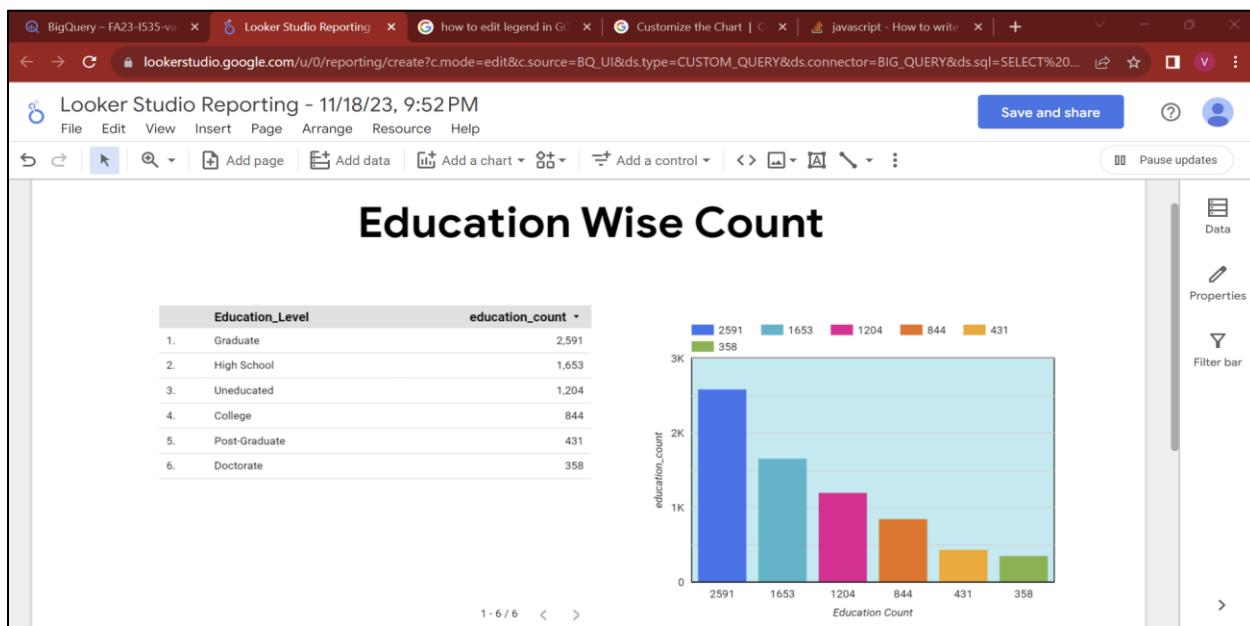
Credit Card Category Wise count:

- People holding Blue category credit card are the highest as compared to Platinum cards holders with lowest in number



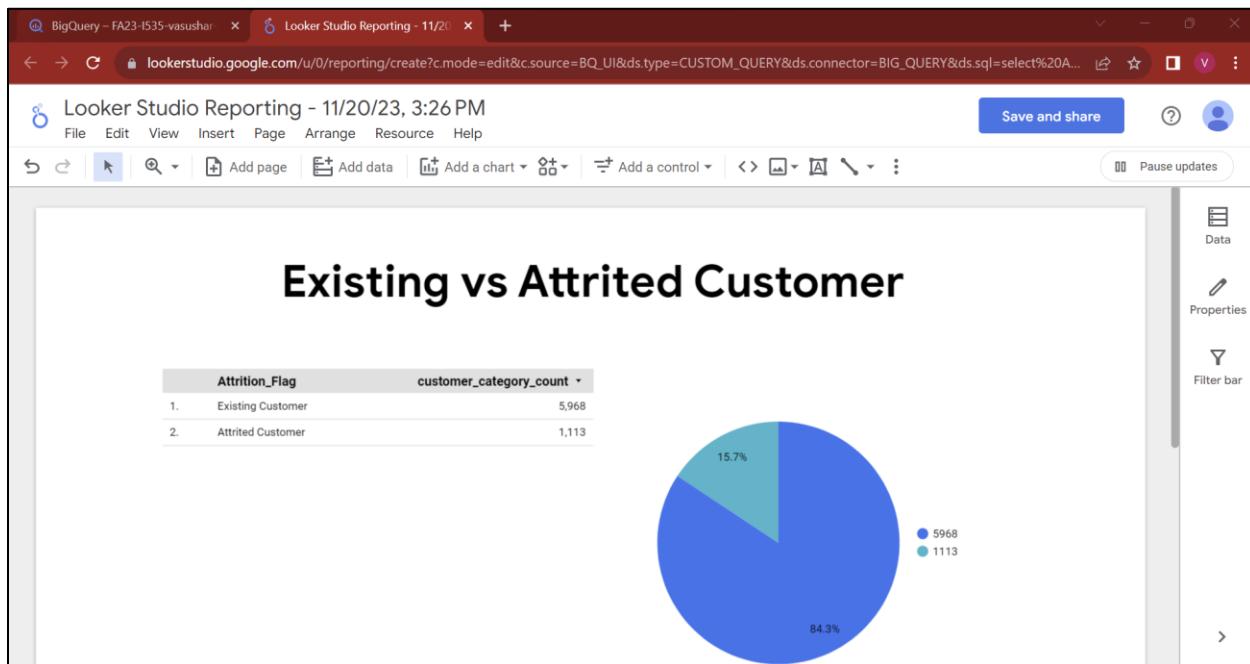
Education Wise Credit Card holder Count:

- We can see that Graduates holds the card in highest ratio as compared to Doctorate with lowest count in Credit Card holders.



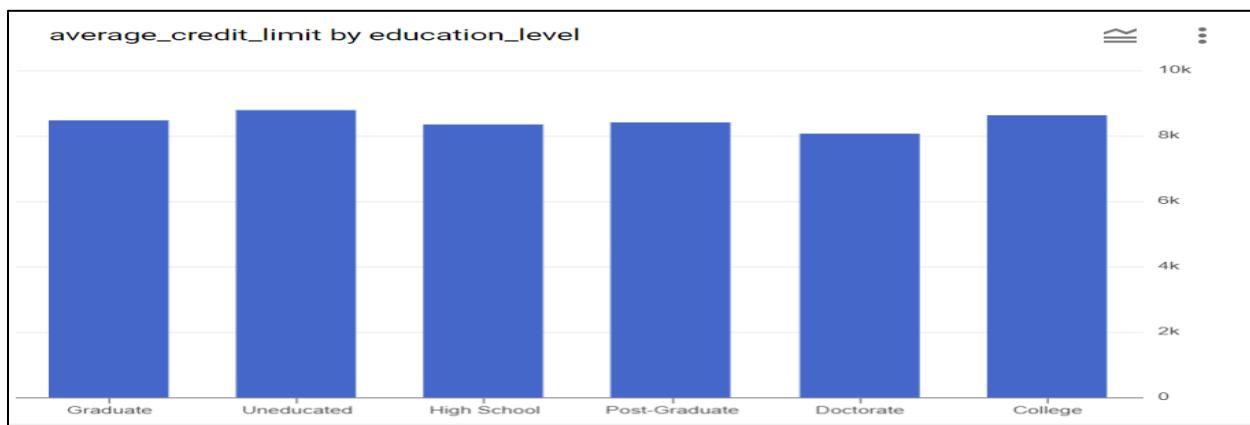
Existing Customer vs Attrited Customer:

- Following graph have shown the number of customers that bank has retained vs the number of customers who have left the institution



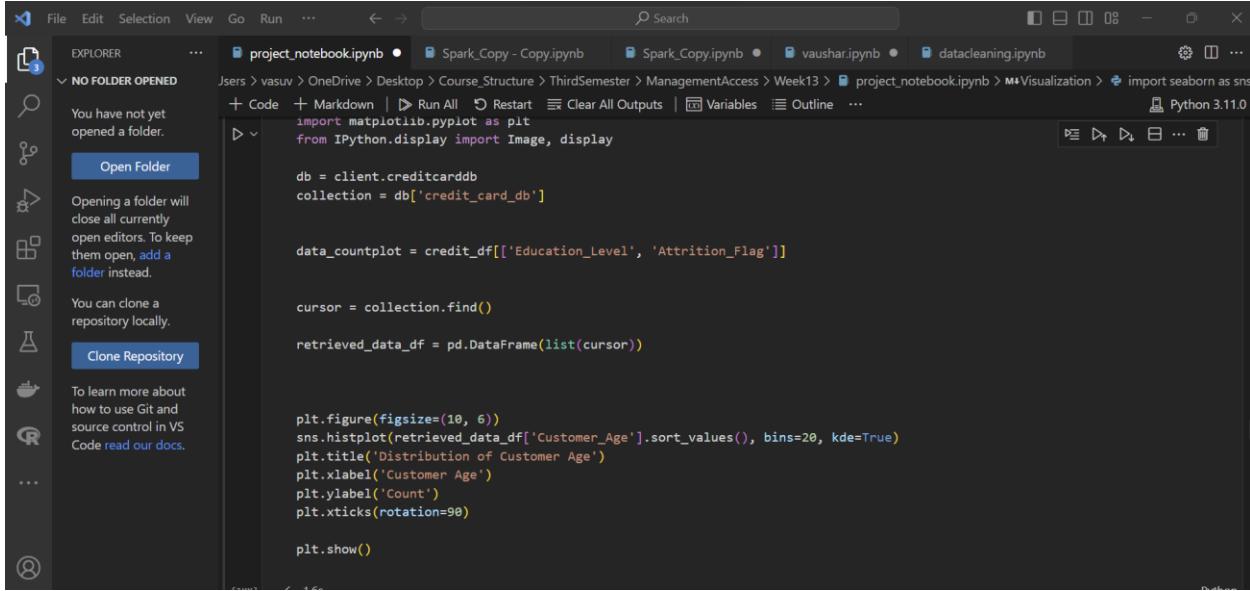
Average Credit Limit per Education level:

- As we can see there is not much difference in average credit limit as per education level and all of them are mostly the same.



Visualizations created using Pandas and MongoDB on JetStream:

The following graph shows the age distribution of existing and attrited credit card customers of the bank.



A screenshot of a Jupyter Notebook interface. The left sidebar shows the 'EXPLORER' tab with a message: 'You have not yet opened a folder.' Below it are buttons for 'Open Folder', 'Clone Repository', and a link to learn about Git and source control. The main area displays the following Python code:

```
File Edit Selection View Go Run ... Search
project_notebook.ipynb • Spark_Copy - Copy.ipynb • Spark_Copy.ipynb • vaushar.ipynb • datacleaning.ipynb
Jusers > vasuv > OneDrive > Desktop > Course_Structure > ThirdSemester > ManagementAccess > Week13 > project_notebook.ipynb > M>Visualization > import seaborn as sns
+ Markdown | ▶ Run All ⌂ Restart ⌂ Clear All Outputs | ⌂ Variables ⌂ Outline ...
import matplotlib.pyplot as plt
from IPython.display import Image, display

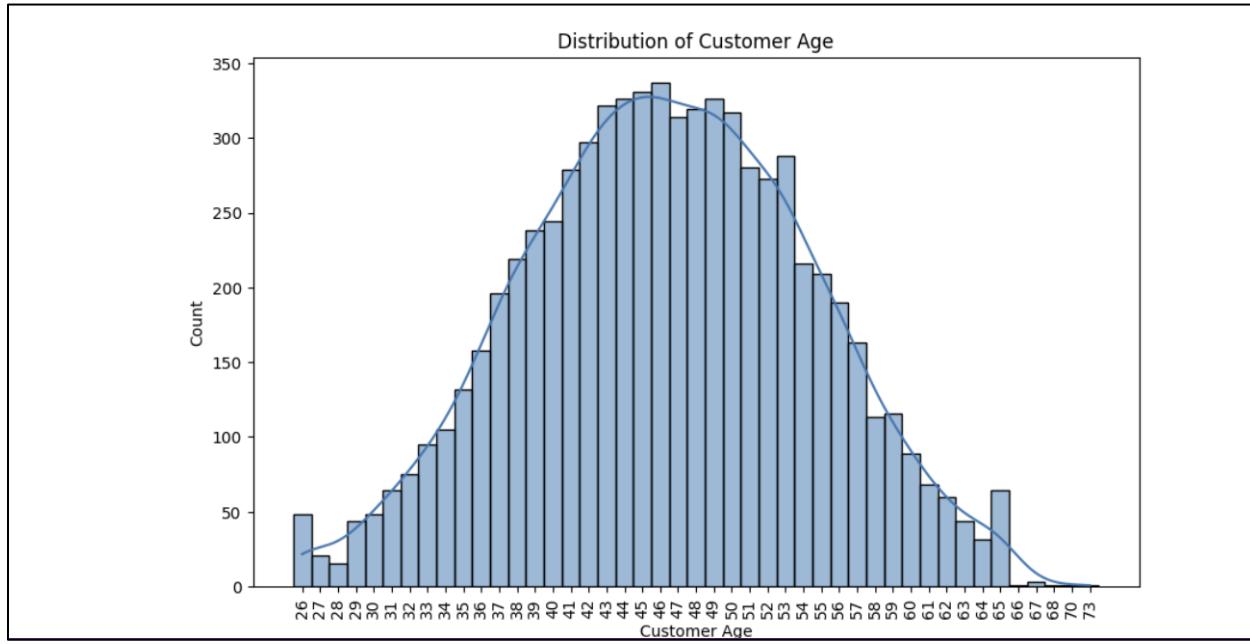
db = client.creditcarddb
collection = db['credit_card_db']

data_countplot = credit_df[['Education_Level', 'Attrition_Flag']]

cursor = collection.find()
retrieved_data_df = pd.DataFrame(list(cursor))

plt.figure(figsize=(10, 6))
sns.histplot(retrieved_data_df['Customer_Age'].sort_values(), bins=20, kde=True)
plt.title('Distribution of Customer Age')
plt.xlabel('Customer Age')
plt.ylabel('Count')
plt.xticks(rotation=90)

plt.show()
```



Education Level by Attrition Flag:

- The education level of various customers is shown.
- Graduate students seems to be highest in number for both existing and attrited customers.

```
import seaborn as sns
import matplotlib.pyplot as plt
from IPython.display import Image, display

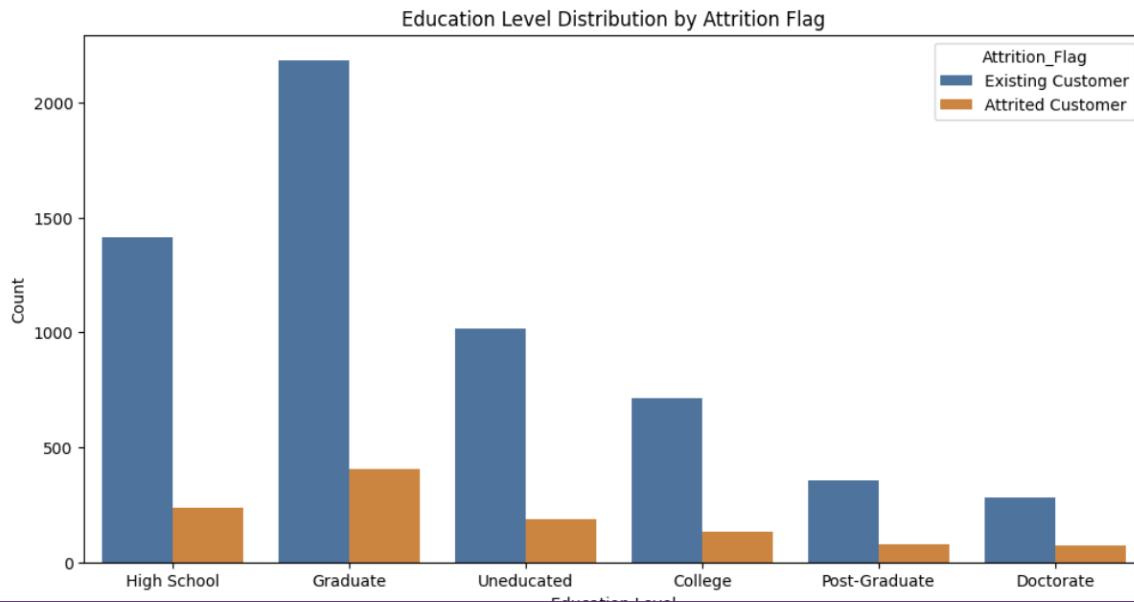
db = client.creditcarddb
education_level_collection = db['education_level_collection']

data_countplot = credit_df[['Education_Level', 'Attrition_Flag']]

data_countplot_dict = data_countplot.to_dict(orient='records')
education_level_collection.insert_many(data_countplot_dict)

cursor = education_level_collection.find()
retrieved_data_df = pd.DataFrame(list(cursor))

plt.figure(figsize=(12, 6))
sns.countplot(x='Education_Level', data=retrieved_data_df, hue='Attrition_Flag')
plt.title('Education Level Distribution by Attrition Flag')
plt.xlabel('Education Level')
plt.ylabel('Count')
plt.show()
```



Credit Limit vs Income Category:

- The credit limit of customers based on their income is shown
- People with highest income category have highest credit limit.

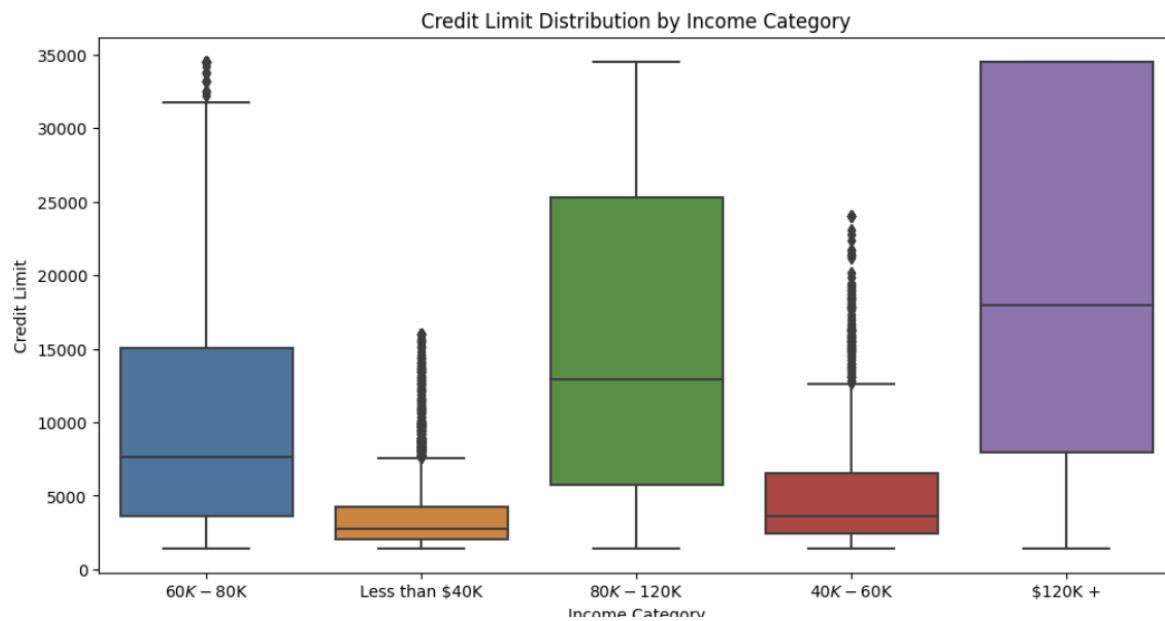
```
db = client.creditcarddb
income_category_collection = db['income_category_collection']

data_for_boxplot = credit_df[['Income_Category', 'Credit_Limit']]

data_for_boxplot_dict = data_for_boxplot.to_dict(orient='records')
income_category_collection.insert_many(data_for_boxplot_dict)

cursor = income_category_collection.find()
retrieved_data_df = pd.DataFrame(list(cursor))

plt.figure(figsize=(12, 6))
sns.boxplot(x='Income_Category', y='Credit_Limit', data=retrieved_data_df)
plt.title('Credit Limit Distribution by Income Category')
plt.xlabel('Income Category')
plt.ylabel('Credit Limit')
plt.show()
```



Pairplot:

- The pairplot for various attributes of the customers has been plotted as below

```
# Assuming 'pairplot_collection' is the collection name and 'creditcarddb' is the database name
db = client.creditcarddb
pairplot_collection = db['pairplot_collection']

pairplot_collection_df = credit_df[['Customer_Age', 'Credit_Limit', 'Total_Trans_Amt', 'Avg_Utilization_Ratio']]

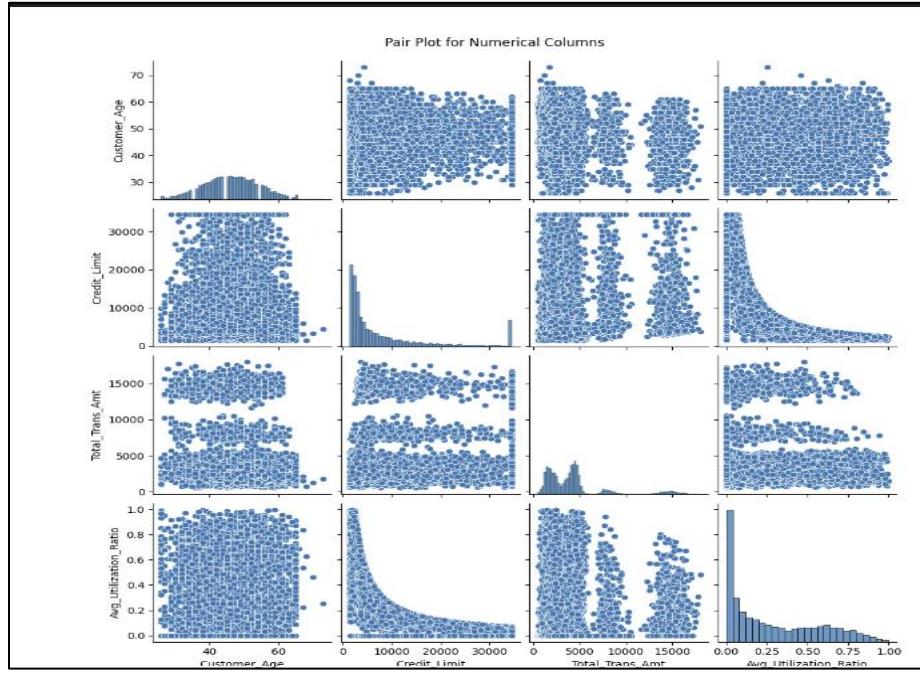
# Insert the prepared data into MongoDB
pairplot_collection_dict = pairplot_collection_df.to_dict(orient='records')
pairplot_collection.insert_many(pairplot_collection_dict)

# Retrieve data from MongoDB
cursor = pairplot_collection.find()
retrieved_data_df = pd.DataFrame(list(cursor))

retrieved_data_df['Credit_Limit'] = pd.to_numeric(retrieved_data_df['Credit_Limit'], errors='coerce')
retrieved_data_df['Customer_Age'] = pd.to_numeric(retrieved_data_df['Customer_Age'], errors='coerce')
retrieved_data_df['Total_Trans_Amt'] = pd.to_numeric(retrieved_data_df['Total_Trans_Amt'], errors='coerce')
retrieved_data_df['Avg_Utilization_Ratio'] = pd.to_numeric(retrieved_data_df['Avg_Utilization_Ratio'], errors='coerce')

# Check for missing values and handle them
retrieved_data_df = retrieved_data_df.dropna()

# Create pairplot
sns.pairplot(retrieved_data_df[['Customer_Age', 'Credit_Limit', 'Total_Trans_Amt', 'Avg_Utilization_Ratio']])
plt.suptitle('Pair Plot for Numerical Columns', y=1.02)
plt.show()
```



5.Discussions:

After doing analysis and making visualizations on various platforms, one can conclude the following:

- As we have compared various characteristics of Existing and Attrited Customer, Existing customers have much higher number of total transaction count as compared to Attrited customer. This means that existing customers uses credit cards much frequently as compared to attrited customers.
- Number of Males Credit Card Customers are higher as compared to number of female Credit Card Customers.
- Blue Credit card category people sums up the highest total revolving balance as compared to other card categories. One of the reasons could be that the highest number of people have subscribed to Blue Card category as compared to other card categories.
- As mentioned earlier, Blue category credit card holder have higher count as compared to all other credit card category. Blue credit card makes almost 93% of all credit card that have been issued.
- Customers whose education level is graduate are the highest number of customers and credit card holders as per this dataset
- Currently as per the dataset present, bank has more existing customers as compared to Attrite credit card customers which is a positive sign for financial growth of the institution.
- People between the ages of 40-50 are the highest credit card holders as compared to people in other age brackets.
- People in highest income category are having the highest credit limit.
- People with graduate degree are the highest number in existing customer as well as attrited employee.

After completing this analysis, we could say that this course has helped a me lot in order to employ various techniques and strategies required to analyze the data.

The new tools that I have learnt and used in this project are MongoDB, Docker, JetStream and Google Cloud platform. These platforms are great to work with and provides a lot of flexibility in terms of analyzing the dataset. Also though this project, I came to know that one can create informative visualizations on MongoDB itself and Google Cloud Platform itself.

Barriers and Challenges:

The barrier that I faced while building this project was that we found it quite difficult to use MongoDB on Docker on Jetstream platform. We weren't aware of how to integrate MongoDB with Docker and then how to use MongoDB using pandas. It took me sometime to figure all this out and to create visualizations using this platform.

We also tried to integrate Spark in this tech stack of Jetstream, Docker, MongoDB, Pandas but was repeatedly facing lag issue. Hence even after creating Spark instance and queries, we didn't use it in the final project.

Also it was quite difficult to edit graphs and customize graphs created on Google Cloud Platform. The options to edit graphs and to create graphs with desirable

Measures were quite difficult to find. It seems that GCP can work a lot on LookerStudio.

6.Conclusions:

The conclusion that is drawn from the above analysis is that there can be a number of factors influencing people holding credit card such as their marital status, age, income etc. One can not depend only on one parameter to assess what parameters actually influences the credit card holders and their capacity to spend money.

Also although one can say that attrition rate of credit card holders doesn't solely depend on factors such as gender, education, income etc one of the factors that could determine if a person will continue a credit card or not is total number of transaction count done on a credit card. As mentioned earlier, since people who

have left credit card service have much lower transaction count as compared to people who are continuing credit card service of a financial institution. Hence number of transaction done by a customer can help financial institutions to assess if a person will continue their credit card service or not.

7. References:

<https://www.kaggle.com/search?q=bank+product>

<https://cloud.google.com/docs>

<https://www.mongodb.com/docs/>

<https://spark.apache.org/docs/latest/api/python/index.html>