

PREDICTING CREDIT CARD CUSTOMER ATTRITION

BY

FELIX OPARAH

Contents

1. Introduction	3
1.1 Background of Study	3
1.2 Objectives of study.....	3
1.3 Methodology.....	4
1.4 Description of Dataset	4
2. Analytical Design.....	4
3. Data Analytics	5
3.1 Data Importation.....	5
3.2 Data Cleaning	6
3.2.1 Variable type	6
3.2.2 Missing values	6
3.2.3 Summary statistics	7
3.2.4 Variable Dropping	7
3.3 Visualization	7
4. Predictive Machine Learning modeling	13
4.1 Logistic regression.....	14
4.1.1 Interpretation of results.....	15
4.2 Decision Tree.....	16
4.2.1 Interpretation of result	18
4.3 Random Forest.....	18
4.3.1 Interpretation of result	20
5. Comparison of Predictive Accuracy powers of the models.....	21
6. Critical comparison of SAS and Python programing.	21
6.1 Introduction	21
6.2 Overview of SAS.	21
6.3 Overview of Python programing language	21
6.4 Review of SAS and Python programing.....	21
7. Conclusion.....	23
8. References	25
9. Appendix 1 : SAS codes	26
9.1 Data pre-processing and mining	26
9.3: Machine learning Modeling	29
10. Appendix 2: Python codes and results.....	31
10.1 Data pre-processing and mining	31
10.2: Data Visualization	32
10.3: Machine learning Modeling	40

1. Introduction

1.1 Background of Study

Current development in the financial services sector indicates that credit card attrition is one of the major concerns of financial institutions especially banks. This is basically due to its negative impact on profitability. Market dynamics suggests that the banking sector is highly competitive due to the presence of large number of banks, other financial services providers and financial technology companies (Fintechs) competing for the same customers. Hence, each bank tries to retain its customers and maintain their loyalty through the provision of essential products such as credit cards. Customer loyalty is directly related to profit growth; hence, banks should avoid losing customers (Jagadeesan and Indhuja, 2020).

Credit card customer attrition or churn is the loss of customers' usage of a bank's credit card product. Each year, many banks lose their credit card customers and this negatively affects their profit. Since banks seek to maximize profits, it is pertinent to build a system which would proactively predict customers who are likely to stop using this product. Thus appropriate measures could be put in place to forestall such customer churn, probably by offering them better services. Therefore, building a predictive machine learning model becomes pertinent (Parangi, 2022).

This study builds a credit card customer attrition prediction system using relevant historic data. Data on credit card attrition was obtained from the Kaggle website and predictive analytics was implemented on the data using Machine Learning algorithms. This formed a model of predictive analytical system which would predict credit card customer attrition with high accuracy levels.

The major analytical tool adopted for this study was the Statistical Analysis Software SAS programming. However, the Python programming tool was also applied to the same data in order to make comparison between the two programming languages.

1.2 Objectives of study

The general objective of this study is to build a system that would predict credit card customer attrition using the SAS programming language. The specific objectives are to achieve the following:

- a. Analyse the chosen data to understand the characteristics of the data
- b. Perform data visualization using tools such as histogram, charts, boxplots, heat maps, etc so as to discover trends and patterns.
- c. Employ machine learning algorithms to build a predictive model that would predict credit card customer attrition.
- d. Construct a decision making process based on the result.
- e. Adopt the Python programming language to perform the same analysis so as to compare the two languages.

1.3 Methodology

The methodology of this study is a combination of descriptive, analytical and quantitative techniques to create a thorough understanding of the subject, for the purpose of real life implementation. Data mining techniques and Machine Learning (ML) algorithms were employed to extract insights which would aid strategic business decisions. The SAS programming was adopted using the various advanced procedures (PROC) to clean data, visualize data and implement ML algorithms. Several data visualization types such as scatter plots, bar charts, histogram, boxplot, heat maps, etc were used in the study to provide insights from the data.

Three classifier ML algorithms were employed to build the model. These are:

- i. Logistic regression
- ii. Decision tree
- iii. Random forest

Predicting credit card attrition is a binary decision making process, hence a Classification problem. Consequently, these algorithms were selected based on their reliable predictive ability and to provide robustness to the study.

1.4 Description of Dataset

Dataset for this study consists of 10,127 observations and 23 variables showing various customers' information, such as personal detail, demographic details and credit card status details. These include factors like age, gender, income groupings, credit card type, inactive period, utilization ratio, Attrition Flag, etc. The data also shows the spending behaviour of customers such as credit limit and revolving balance.

Exploration of this dataset will enable the understanding of customer- bank relationship dynamics and enable a reliable prediction of customer attrition. Given the set of large data points and multiple variables capture, it is believed customer attrition could be predicted with some high levels of certainty. This dataset was obtained from the Kaggle website however, the original source is Zhyli (2020).

2. Analytical Design

This study adopted the Unified Modeling Language (UML). It is a general purpose language that uses a graphical description which can create an abstract model for visualizing, specifying, constructing, and documenting the artifacts of a software-intensive system. It is also used for business modelling and other non-software systems. It offers a standard way to write a system's blueprints, including conceptual items such as business processes and system functions as well as concrete things such as programming language statements, database schemas, and reusable software components. (Waykar, 2013).

The Activity flow in this study is represented in figure 2 below. At the start of the process, data was obtained and cleaned. It was transformed by removing unwanted variables, visualization and changing the structure so as to prepare it for machine learning modeling. The modeling was done and results obtained. Based on the results, customer retention strategies were developed for implementation. Finally the report was submitted.

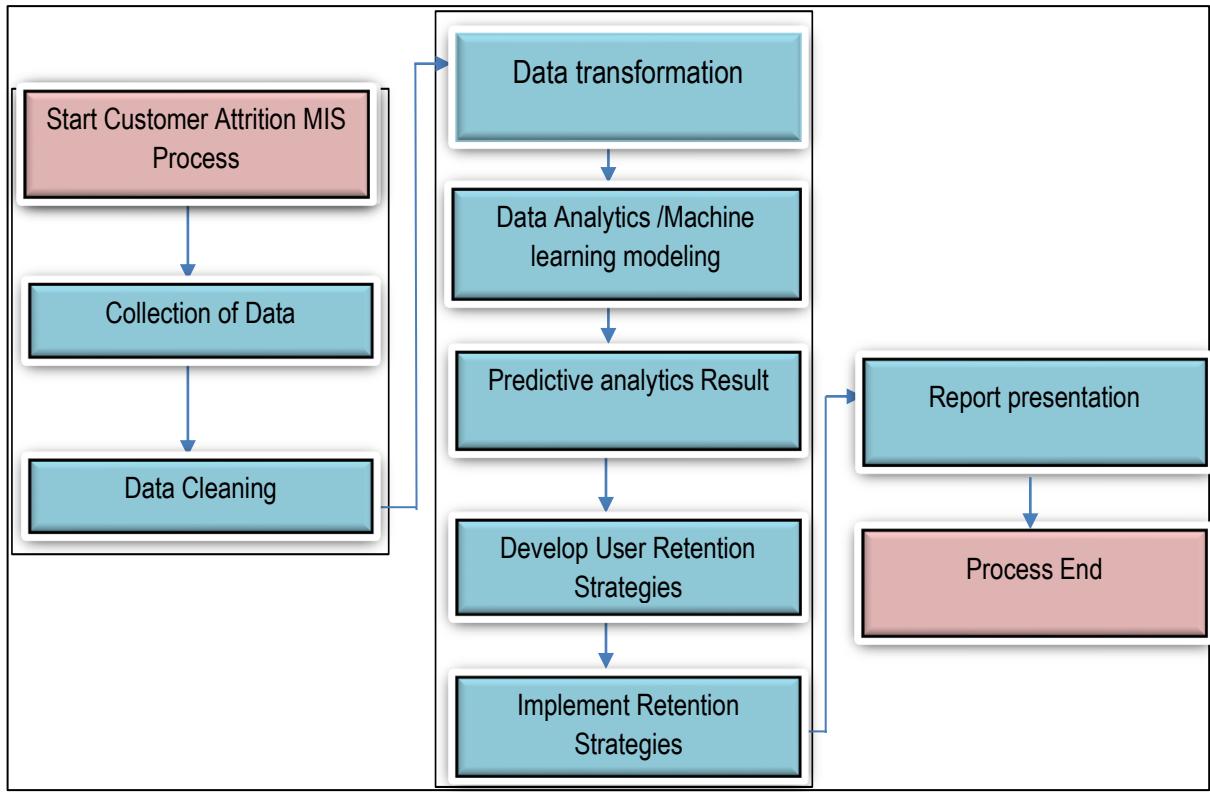


Fig 1: UML Activity design

3. Data Analytics

This section involves mining the data to bring out useful insights. Both SAS and Python were used for all the steps involved in this section. However, only the steps done with SAS are presented in the main body of this analysis, while the steps with Python are contained in the Appendix.

3.1 Data Importation

The dataset is a CSV file and was imported into the SAS Ondemand platform. The first 5 observations are viewed as shown below:

```

1 Libname cust_at '/home/u63398813/';
2 PROC IMPORT DATAFILE='/home/u63398813/BankChurners _1.csv'
3   DBMS=CSV
4   replace
5   OUT=cust_at.Attrition;
6   GETNAMES=YES;
7 RUN;

```

The screenshot shows a table with 5 rows of data. The columns are labeled: Obs, CLIENTNUM, 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, and Avg. The data includes various customer attributes like age, gender, marital status, and income categories, along with financial metrics like credit limit and total revolving balance.

Obs	CLIENTNUM	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.
1	768805383	Existing Customer	45	M	3	High School	Married	\$60K - \$80K	Blue	39	5	1	3	12691	777	
2	818770008	Existing Customer	49	F	5	Graduate	Single	Less than \$40K	Blue	44	6	1	2	8256	864	
3	713982108	Existing Customer	51	M	3	Graduate	Married	\$80K - \$120K	Blue	36	4	1	0	3418	0	
4	769911858	Existing Customer	40	F	4	High School	Unknown	Less than \$40K	Blue	34	3	4	1	3313	2517	
5	709106358	Existing Customer	40	M	3	Uneducated	Married	\$60K - \$80K	Blue	21	5	1	0	4716	0	

Figure 2: Code and output showing first 5 observations in dataset

3.2 Data Cleaning

3.2.1 Variable type

The data was examined to determine the variable types and the number of variables using PROC Contents. The results confirm that data is made of up 22 variables and 10127 observations. The nature of the variables was also shown: numerical and character.

The screenshot shows two tables. The top table, 'The CONTENTS Procedure', provides summary information about the dataset: name (CUST_AT_ATTRITION), member type (DATA), engine (V9), created date (05/12/2023 11:07:26), last modified date (05/12/2023 11:07:26), protection level (NO), and data representation (SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64). The bottom table, 'Alphabetic List of Variables and Attributes', lists 22 variables with their types, lengths, formats, and informats. For example, 'Attrition_Flag' is a character variable with length 17 and informat \$17. Other variables include 'Avg_Open_To_Buy', 'Avg_Utilization_Ratio', and 'Naive_Bayes_Classifier_Attrition'.

The CONTENTS Procedure			
Data Set Name	CUST_AT_ATTRITION	Observations	10127
Member Type	DATA	Variables	22
Engine	V9	Indexes	0
Created	05/12/2023 11:07:26	Observation Length	192
Last Modified	05/12/2023 11:07:26	Deleted Observations	0
Protection		Compressed	NO
Data Set Type		Sorted	NO
Label			
Data Representation	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
Encoding	utf-8 Unicode (UTF-8)		

#	Variable	Type	Len	Format	Informat
1	Attrition_Flag	Char	17	\$17.	\$17.
15	Avg_Open_To_Buy	Num	8	BEST12.	BEST32.
20	Avg_Utilization_Ratio	Num	8	BEST12.	BEST32.
8	Card_Category	Char	6	\$6.	\$6.
12	Contacts_Count_12_mon	Num	8	BEST12.	BEST32.
13	Credit_Limit	Num	8	BEST12.	BEST32.
2	Customer_Age	Num	8	BEST12.	BEST32.
4	Dependent_count	Num	8	BEST12.	BEST32.
5	Education_Level	Char	13	\$13.	\$13.
3	Gender	Char	1	\$1.	\$1.
7	Income_Category	Char	14	\$14.	\$14.
6	Marital_Status	Char	7	\$7.	\$7.
11	Months_Inactive_12_mon	Num	8	BEST12.	BEST32.
9	Months_on_book	Num	8	BEST12.	BEST32.
21	Naive_Bayes_Classifier_Attrition	Num	8	BEST12.	BEST32.
16	Total_Amt_Chng_Q4_Q1	Num	8	BEST12.	BEST32.
19	Total_Ct_Chng_Q4_Q1	Num	8	BEST12.	BEST32.
10	Total_Relationship_Count	Num	8	BEST12.	BEST32.
14	Total_Revolving_Bal	Num	8	BEST12.	BEST32.
17	Total_Trans_Amt	Num	8	BEST12.	BEST32.
18	Total_Trans_Ct	Num	8	BEST12.	BEST32.
22	VAR23	Num	8	BEST12.	BEST32.

Figure 3: Tables showing variable types in dataset

3.2.2 Missing values

Data was also checked for missing values (numeric and categorical). The result shows there are no missing values in the data:

The screenshot shows the 'The MEANS Procedure' output. The first part is a table titled 'Missing Values' showing the count of missing values for each variable. All counts are 0. The second part is a table titled 'Attrition_Flag_miss' showing the count of missing values for specific variables: Attrition_Flag, Gender, Educational_level, Marital_status, Income_Category, and Card_Category. All counts are 0.

Variable	N Miss
Customer_Age	0
Dependent_count	0
Months_on_book	0
Total_Relationship_Count	0
Months_Inactive_12_mon	0
Contacts_Count_12_mon	0
Credit_Limit	0
Total_Revolving_Bal	0
Avg_Open_To_Buy	0
Total_Amt_Chng_Q4_Q1	0
Total_Pt_Chng_Q4_Q1	0
Total_Trans_Amt	0
Total_Trans_Ct	0
Total_Ct_Chng_Q4_Q1	0
Avg_Utilization_Ratio	0
Naive_Bayes_Classifier_Attrition	0
VAR23	0

Attrition_Flag_miss	Gender_miss	Educational_level_miss	Marital_status_miss	Income_Category_miss	Card_Category_miss
0	0	0	0	0	0

Figure 4: Table showing zero missing values

3.2.3 Summary statistics

Summary statistics were obtained using PROC Means as shown in figure 5 below:

```

25 run;
26 /*summary of the numerical variables*/
27 proc means data = cust_at.attrition min q1 median mean q3 max n maxdec=2;
28 run;
29 title "Summary Statistics for numerical variables in Customer Churn Data";
30 run;
31 /*Calculate the standard deviation of the numerical variables*/
32 proc means data = cust_at.attrition Std;
33 run;

```

The MEANS Procedure						
Variable	Minimum	Lower Quartile	Median	Mean	Upper Quartile	Maximum
Customer_Age	26.00	41.00	46.00	48.33	52.00	73.00
Dependent_count	0.00	1.00	2.00	2.50	3.00	8.00
Months_on_book	13.00	31.00	36.00	35.89	40.00	66.00
Total_Relationship_Count	1.00	3.00	4.00	3.81	5.00	8.00
Months_Inactive_12_mon	0.00	2.00	2.00	2.34	3.00	6.00
Contacts_Count_12_mon	0.00	2.00	2.00	2.46	3.00	6.00
Credit_Limit	1436.30	2658.00	4546.00	6955.00	10986.00	34819.00
Total_Revolving_Bal	0.00	357.00	1276.00	1162.81	1784.00	2517.00
Avg_Open_To_Buy	3.00	1324.00	3474.00	7469.14	9861.00	34516.00
Total_Amt_Chang_Q4_Q1	0.00	0.63	0.74	0.76	0.86	3.40
Total_Trans_Amt	\$10.00	2150.00	3899.00	4649.00	4746.00	18469.00
Total_Trans_Ct	10.00	45.00	67.00	64.86	81.00	139.00
Total_Ot_Chang_Q4_Q1	0.00	0.58	0.70	0.71	0.82	3.71
Avg_Utilization_Ratio	0.00	0.02	0.18	0.27	0.50	1.00
Naive_Bayes_Classifier_Attrition	0.00	0.00	0.00	0.16	0.00	1.00
Var23	0.00	1.00	1.00	0.84	1.00	1.00

The MEANS Procedure	
Variable	Std Dev
Customer_Age	8.0168140
Dependent_count	1.2998083
Months_on_book	7.9864163
Total_Relationship_Count	1.5340379
Months_Inactive_12_mon	1.0106224
Contacts_Count_12_mon	1.1062251
Credit_Limit	9488.78
Total_Revolving_Bal	814.9873525
Avg_Open_To_Buy	9090.69
Total_Amt_Chang_Q4_Q1	0.2192068
Total_Trans_Amt	23971.12
Total_Trans_Ct	23.4725704
Total_Ot_Chang_Q4_Q1	0.2380861
Avg_Utilization_Ratio	0.2603169
Naive_Bayes_Classifier_Attrition	0.3653010
Var23	0.3653010

Figure 5: Summary statistics of numerical variables

The summary statistics show that the highest mean value occurred in Credit limit with a value of 11068. This is followed by Avg_Open_To_Buy and Total_Trans_Amt with 9861 and 4741 respectively. The Standard deviation which measures variability in the data shows that the highest variability occurred in Credit Limit with a value of 9088.78. This is followed by Total Revolving Bal with a value of 814.99. This high variability suggests presence of outliers. Most of the variables have zero as their minimum values.

3.2.4 Variable Dropping

The variables were examined for their suitability in the intended analysis. It was discovered that the variables `CLIENTNUM`, `Naive_Bayes_Classifier_Attrition` and `Var23` were not needed. Therefore they were dropped with the appropriate PROC as shown below:

```

37 run;
38 /*dropping variables not needed in the analysis*/
39 data cust_at.attrition;
40 set cust_at.attrition (drop = CLIENTNUM Naive_Bayes_Classifier_Attrition Var23);
41 run;
42 proc print data = cust_at.attrition(obs=5);
43 run;
44 ods graphics on;
45 /*Creating histogram with proc univariate*/

```

Figure 6: Code showing dropping of variables from dataset

3.3 Visualization

Data was visualized to unravel patterns which would provide insight in the analysis.

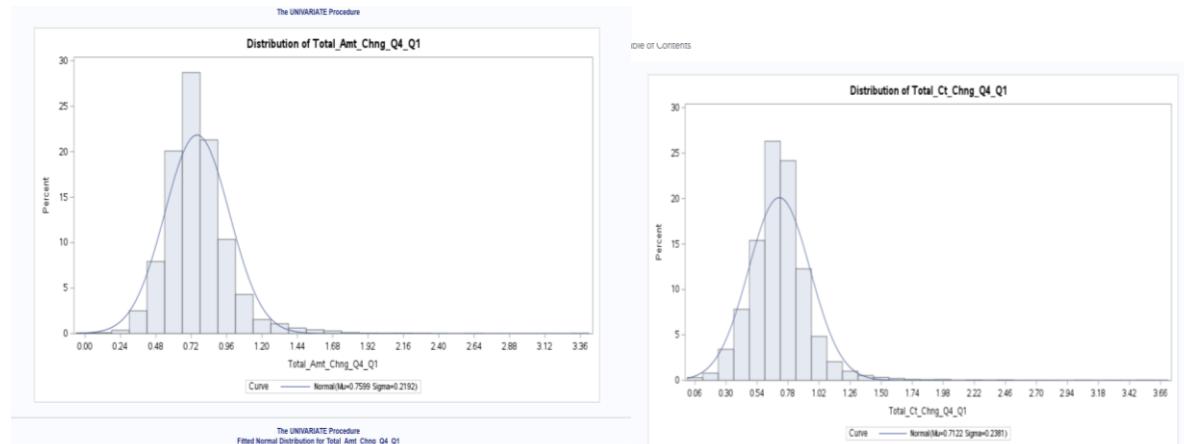


Figure 7: Histogram plots with PROC Univariate

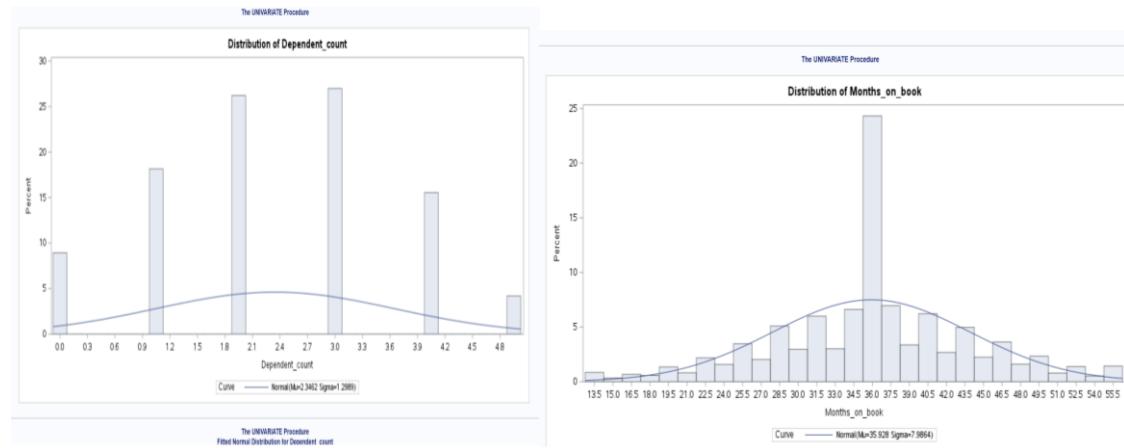


Figure 8: Histogram plots with PROC Univariate

The PROC Univariate was used to obtain histogram for the variables Total_Amt_Chng_Q4_Q1, Total_Ct_Chng_Q1_Q4, Dependent_Count and Months_on_book. The results show that Total_Amt_Chng_Q4_Q1 and Total_Ct_Chng_Q1_Q4 are normally distributed while Dependent_Count and Months_on_book are not normally distributed. This could have resulted from the heavy presence of outliers observed in Dependent_Count and Months_on_book.

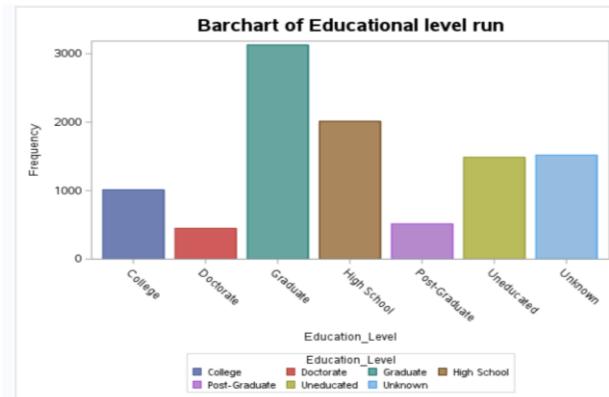


Figure 9: Barchart of Educational Level

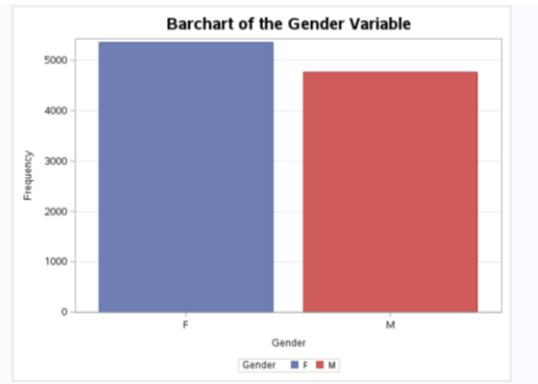


Figure 10: Barchart of Gender

The bar chart in figure 8 shows that graduate education has the highest frequency of 3000 while doctorate degree has the lowest frequency of 500. Similarly, the Gender barchart shows that female customers are more than male customers.

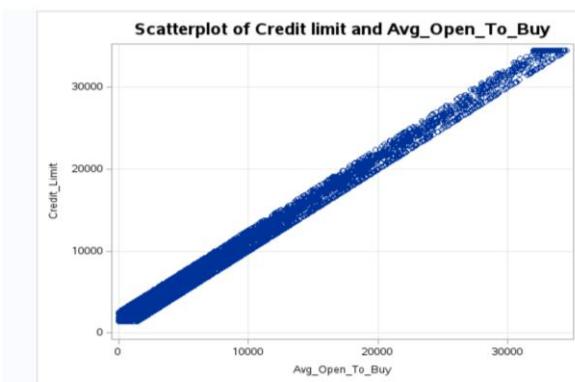


Figure 11: Scatterplot of Credit limit and Avg open to buy

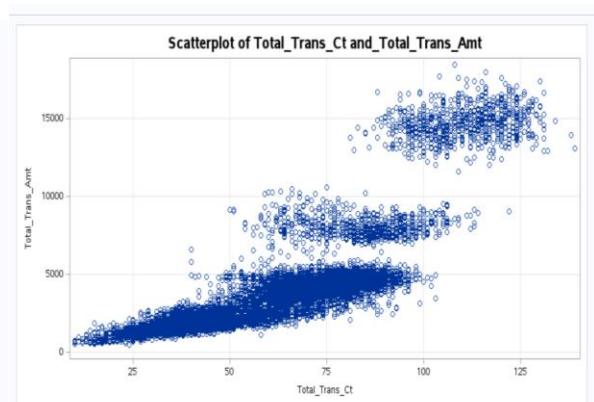


Figure 12: Scatter plot of Total Trans Ct & Total Trans Amt

The Scatterplot in figure 10 shows the points are closely clustered on a straight line and sloped upwards from the origin (zero) to the right hand side of the diagram. This implies a very high positive correlation between Credit limit and Avg open to buy. In other words, the two variables move in the same direction. Similarly, figure 11 shows a positive relationship between Total Trans Ct and Total Trans Amt.

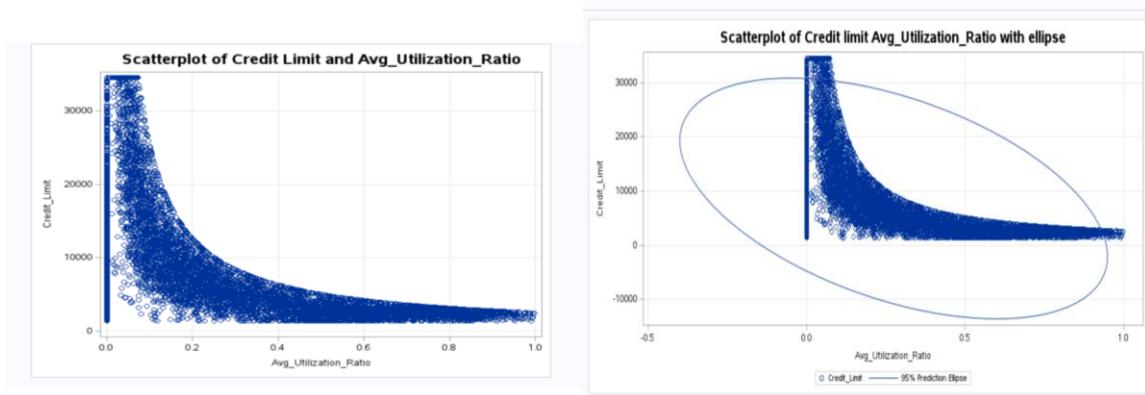


Fig 13 Scatterplots of Credit limit and Ave Utilization Ratio

Scatterplots in figure 11 show an inverse relationship between Credit Limit and Average Utilization Ratio. As Credit limit increases, Average utilization falls, indicating a negative correlation.

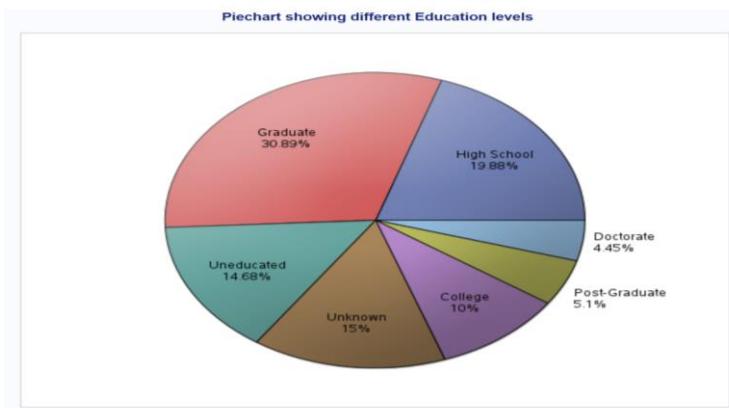


Figure 14: Pie chart of Education levels

The pie chart in figure 14 shows the percentage of the various educational levels. Graduate education with 30.89% is the highest followed by High school with 19.88%. Doctorate education is the lowest with 4.45% while all the other types of education ranges from 5.1% to 15%.

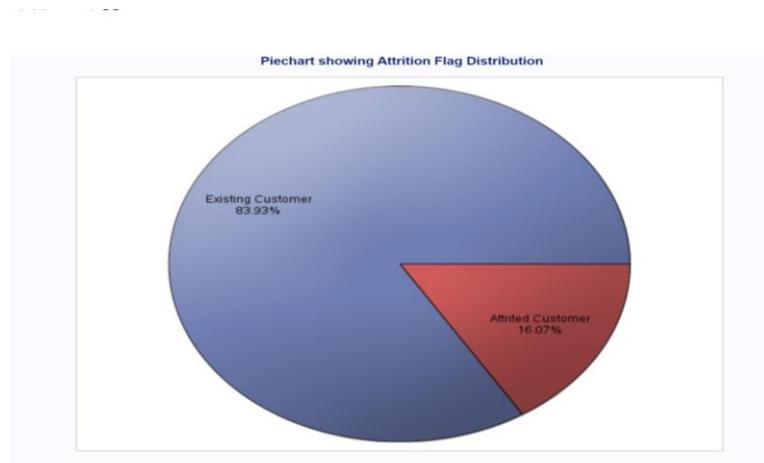


Figure 15: Pie chart of Attrition flag

Figure 15 shows the two types of Attrition flag and their percentage. Existing customers are 83.93% of the total dataset while Attrited Customers are 16.07% or 0.1607 of the entire data base.

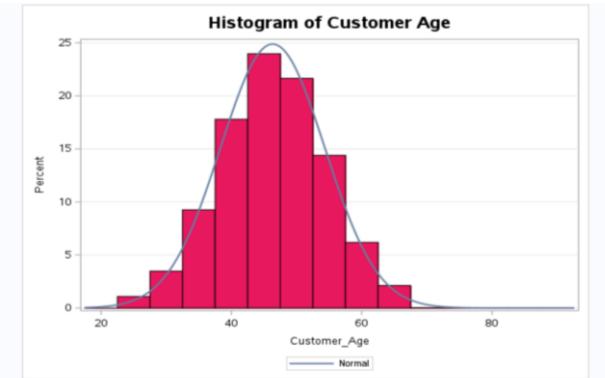


Figure 16: Histogram of Customer Age

The histogram in figure 16 shows the distribution of Customer Age. The bell shaped histogram in this figure suggests normality of the underlying data. It also suggest no outliers.

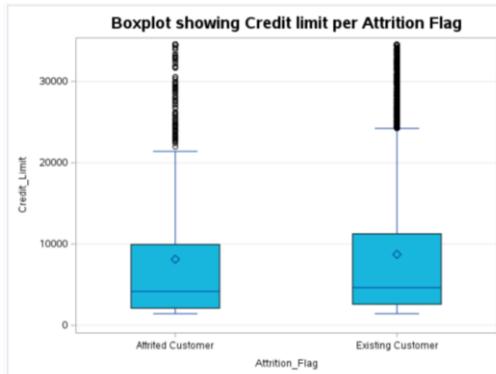


Fig 17: Boxplot of Credit limit & Attrition flag

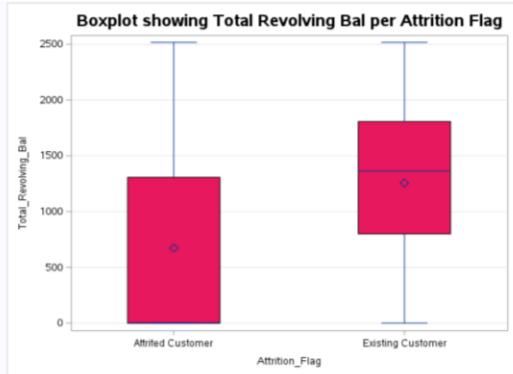


Fig 18: Boxplot of Total Rev bal & Attrition flag

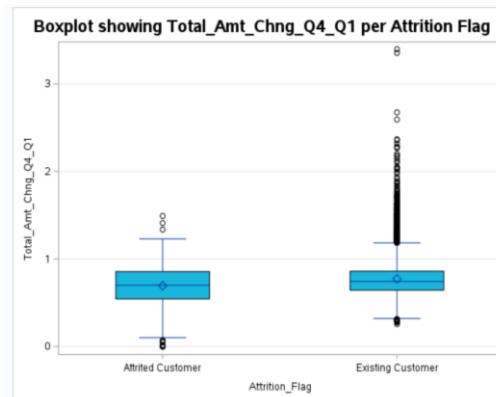


Fig 19: Boxplot of Total Amt Chng & Attrition flag

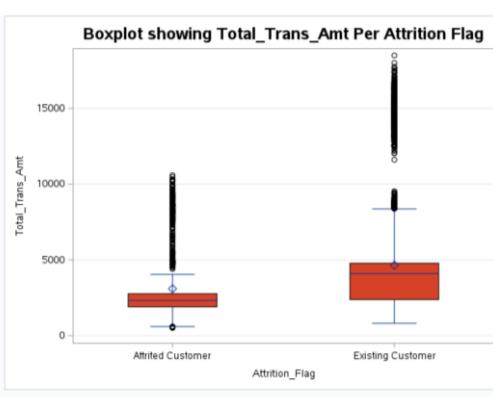


Fig 20: Boxplot of Total Trans Amt & Attrition flag

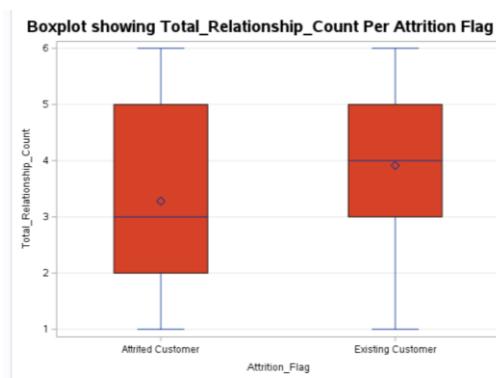


Fig 21: Boxplot Total r/ship count & Attrition flag

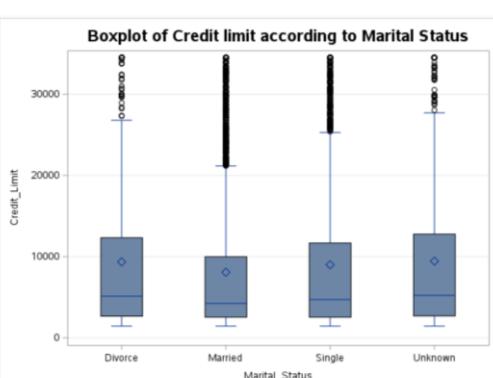


Fig 22: Boxplot of Credit limit & marital status

Figures 17, 18, 19, 20 and 21 are boxplots showing the distribution of Credit limit, Total Revolving bal, Total_Amt_Chng_Q4_Q1, Total Trans Amt and Total_Relationship_Count per Attrition Flag. Similarly, Figure 22 shows boxplot of Credit Limit in relation to marital status. Heavy presence of outliers where

noticed in figures 17, 19, 20 and 22. These outliers impinged on the convergence to normality of these variables as shown in the diagrams.

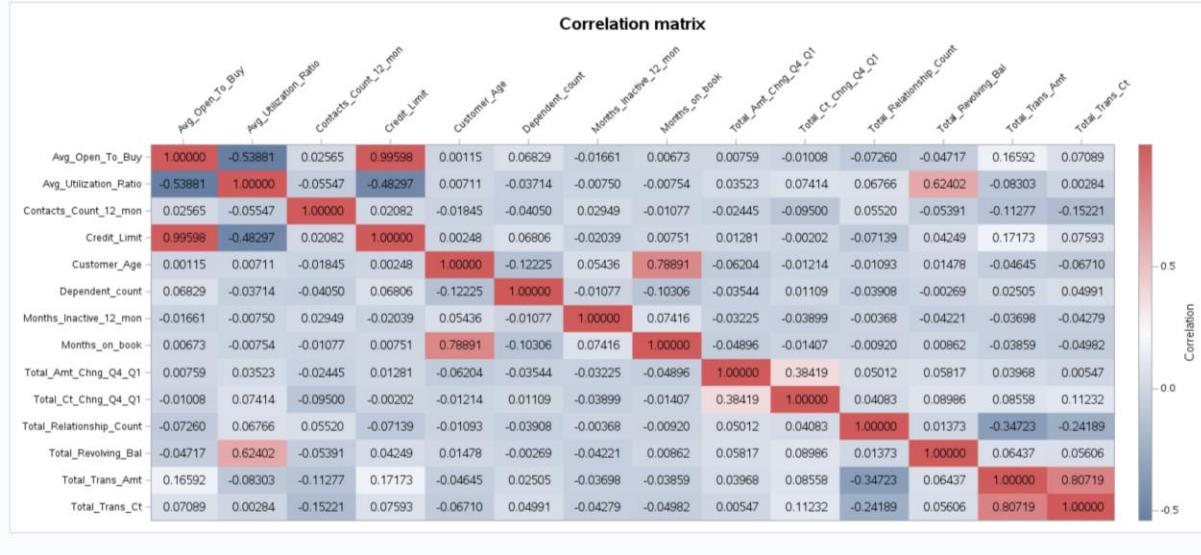


Fig 23: Heatmap of numerical variables

The correlation matrix showing correlation among all the numerical variables was presented in a heatmap as shown in Figure 23. The highest positive correlation is between Credit limit and Avg_Open_To_Buy with a value of 0.996. This is followed by correlation between Total_Trans_Ct and Total_Trans_Amt with a value of 0.807. Most of the variables have very low negative correlation.

4. Predictive Machine Learning modeling

The dataset was evaluated and it was observed the business problem here is a Classification problem. The central problem is to determine how to predict credit card customers that would churn, thereby exposing the business to reduced profitability. Consequently, three machine learning (ML) algorithms were selected and used to model the data accordingly. These ML algorithms are:

- i. Logistic Regression
- ii. Decision Tree
- iii. Random Forest.

Classification algorithm reduces the Attrition_Flag variable to a binary variable of 0 and 1. The data was divided randomly into training and validation (test) set using PROC Surveyselect with a sample rate of 0.7. Hence the training data is 70% of the data while validation data is 30%. The PROC Surveyselect statement is shown below:

```

285 /* Split data into two datasets : 70%- training 30%- validation*/
286 Proc Surveyselect data=CUST_ATT.ATTRITION out=split seed=2345 samprate=.7 outall;
287 Run;
288 Data training validation;
289 Set split;
290 if selected = 1 then output training;
291 else output validation;
292 Run;

```

4.1 Logistic regression

Logistic regression is also called a logit model, is employed to investigate dichotomous outcome variables. In the present case, a binary logistics regression is employed. In the model the log odds of the outcome is demonstrated as a linear combination of the predictor variables.

PROC Logistic was used to implement this as follows:

```

312 /* Logistic regression Model*/;
313 ods graphics on;
314 Proc Logistic Data = training descending;
315 class Attrition_Flag Gender Education_Level Marital_Status Income_Category Card_Category / param = ref;
316 Model Attrition_Flag = Gender Education_Level
317 Marital_Status Income_Category Card_Category Customer_Age Dependent_count
318 Months_on_book Total_Relationship_Count Months_Inactive_12_m
319 Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
320 Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
321 Avg_Utilization_Ratio / selection = stepwise s1stay=0.15 s1entry=0.15 stb;
322 score data=training out = Logit_Training fitstat outroc=troc;
323 score data=validation out = Logit_Validation fitstat outroc=vroc;
324 Run;

```

The result of the logistic regression is presented below



Figure 24: LR model information



Figure 25: Model convergence report (selected variable)

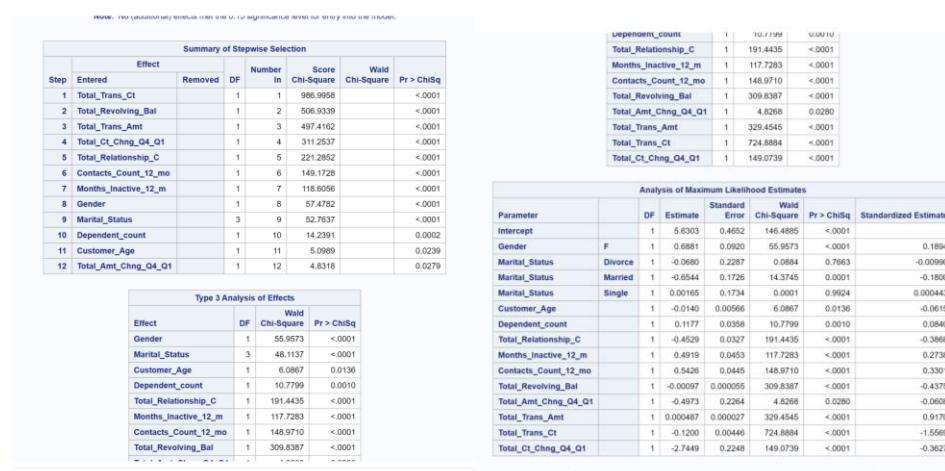


Fig 26: Stepwise selection, Analysis of Effects and Maximum Likelihood estimates

Total_Trans_Amt	1	0.000487	0.000027	329.4545	<.0001	0.9170
Total_Trans_Ct	1	-0.1200	0.00446	724.8884	<.0001	-1.5569
Total_Ct_Chng_Q4_Q1	1	-2.7449	0.2248	149.0739	<.0001	-0.3627
Odds Ratio Estimates						
Effect		Point Estimate		95% Wald Confidence Limits		
Gender F vs M		1.990		1.662 2.383		
Marital_Status Divorce vs Unknown		0.934		0.597 1.463		
Marital_Status Married vs Unknown		0.520		0.371 0.729		
Marital_Status Single vs Unknown		1.002		0.713 1.407		
Customer_Age		0.986		0.975 0.997		
Dependent_count		1.125		1.049 1.207		
Total_Relationship_C		0.636		0.596 0.678		
Months_Inactive_12_m		1.635		1.496 1.787		
Contacts_Count_12_mo		1.720		1.577 1.877		
Total_Revolving_Bal		0.999		0.999 0.999		
Total_Amt_Chng_Q4_Q1		0.608		0.390 0.948		
Total_Trans_Amt		1.000		1.000 1.001		
Total_Trans_Ct		0.887		0.879 0.895		
Total_Ct_Chng_Q4_Q1		0.064		0.041 0.100		
Association of Predicted Probabilities and Observed Responses						
Percent Concordant	92.6	Somers' D	0.853			
Percent Discordant	7.4	Gamma	0.853			
Percent Tied	0.0	Tau-a	0.231			
Score	0.0000	0.0000	0.0000			

Fig 27: Odds Ratio estimates

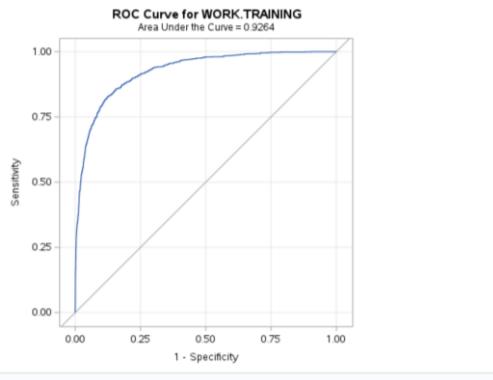


Fig 28: ROC curve for model validation

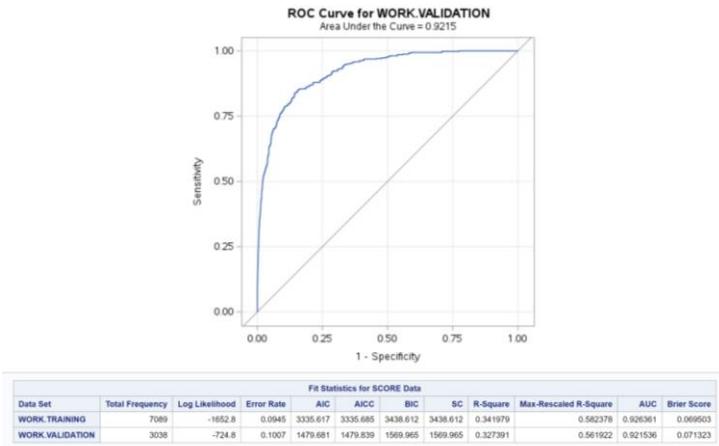


Fig 29: ROC and Model Fit Statistics

4.1.1 Interpretation of results

- i. The model information (fig 24) shows that the response variable is Attrition flag with 2 level of response in a binary Logit model. The observation used for the training data was 7089 made up of 1146 for Attrited customers and 5943 for Existing Customer. The probability modeled is Attrited Customer. This implies the model is built to predict the probability of Customers leaving the credit card scheme.
- ii. Model conversion criteria were satisfied for all the explanatory variables as the model iteratively refines the estimates to arrive at the optimal solution. This is shown in figure 25.
- iii. The “Odds Ratio Estimates” table (figure 27) gives a summary of the significant explanatory variables and shows their associated odds ratios with 95% confidence limits. The coefficients of the model are exponentiated to derive the odds ratio. It could be interpreted as the multiplicative change in the odds for a one unit change in the predictor variable. The odds for Customer Age is 0.986. This means that for 1 year increase in Customer age, the odds of customer attrition will increase by 0.986. Similarly, the odds for Dependent_count is 1.125. This means that for one unit increase in Dependent_count, the odds of Customer Attrition will increase by a factor 1.125. The same explanation hold for other variables.

- iv. The ROC curve (receiver operating characteristic curve) in figure 29 shows the area under the curve for training and validation models are 0.9264 and 0.9215. This implies 92.64% and 92.15% predictive abilities for the training and validation data respectively. This is a very high predictive ability.

4.2 Decision Tree

Decision tree is one of the supervised machine learning algorithms utilized in classification modeling. It employs a hierarchical, tree structure containing root nodes, branches, internal nodes and leaf nodes. The node types perform evaluations to form homogenous subsets represented by leaf nodes, or terminal nodes. The leaf nodes symbolize all the possible outcomes within the dataset. Decision tree also uses a divide and conquer strategy by conducting search to identify the optimal split points within a tree. This process of splitting is then repeated recursively in a top-down fashion until all, or the bulk of records have been classified (Miao and Wang, 2022).

PROC hpsplit was utilized to perform decision tree modeling in this study. This is shown below:

```

325 /* Decision Tree Model*/
326 proc hpsplit data= cust_at.Attrition;
327 class Attrition_Flag Gender Education_Level Marital_Status Income_Category Card_Category _CHARACTER_;
328 Model Attrition_Flag = Gender Education_Level
329   Marital_Status Income_Category Card_Category Customer_Age Dependent_count
330   Months_on_book Total_Relationship_Count Months_Inactive_12_mon
331   Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
332   Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct Chng_Q4_Q1
333   Avg_Utilization_Ratio ;
334 PARTITION FRACTION(validate=0.3 SEED=12345);
335 grow entropy;
336 prune costcomplexity;
337 run;
```

The results are presented in the diagrams below:

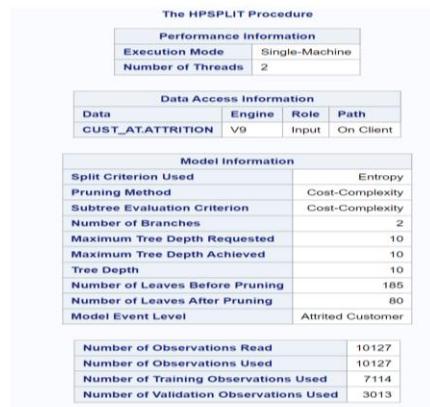


Fig. 30: Model information

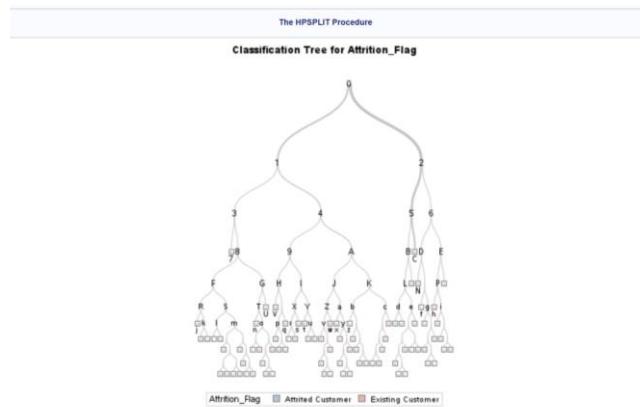


Fig. 31: Classification tree

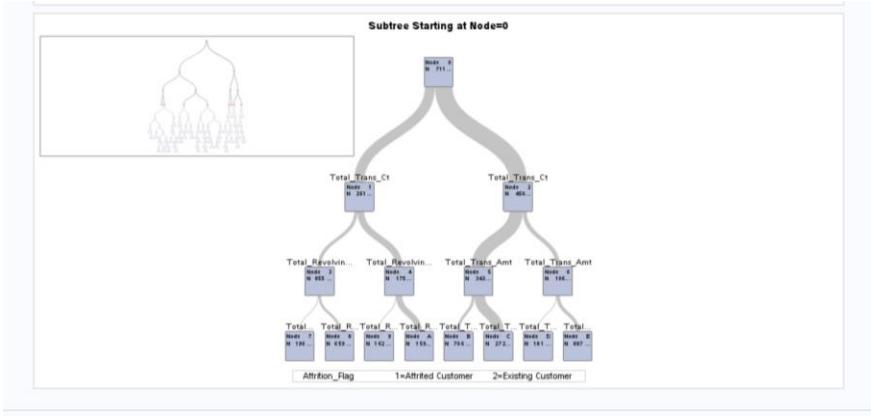


Fig. 32: Sub tree

The HPSPPLIT Procedure						
		Confusion Matrices				
		Predicted				
Training	Actual	Attrited Customer	Existing Customer	Error Rate		
	Attrited Customer	992	144	0.1266		
Validation	Attrited Customer	73	5905	0.0122		
	Existing Customer	382	109	0.2220		
Fit Statistics for Selected Tree						
	N	Leaves	A SE	Mis-class	Sensitivity	Specificity
Training	80	0.0252	0.0305	0.8732	0.9878	0.1381
Validation	80	0.0449	0.0534	0.7780	0.9794	0.1791
					RSS	AUC
Training	80				358.0	0.9850
Validation	80				270.8	0.9432

Fig. 33: Confusion Matrix & Fit Stats

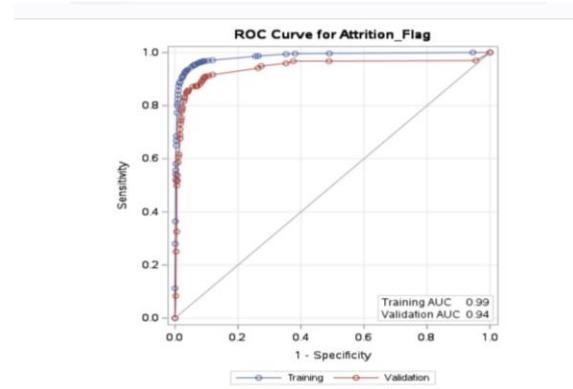


Fig 34: ROC Curve

Variable	Variable Importance					
	Training		Validation		Relative Ratio	Count
	Relative	Importance	Relative	Importance		
Total_Trans_Ct	1.0000	23.1944	1.0000	14.9171	1.0000	10
Total_Revolving_Bal	0.7883	18.2840	0.7949	11.6583	1.0084	5
Total_Trans_Amt	0.5611	13.0143	0.5948	8.8727	1.0601	19
Total_Relationship_Count	0.5525	12.8148	0.5407	8.0655	0.9786	4
Total_Ct_Chng_Q4_Q1	0.4901	11.3682	0.4869	7.2632	0.9934	8
Total_Amt_Chng_Q4_Q1	0.2832	6.5687	0.2015	3.0064	0.7117	7
Months_on_book	0.2172	5.0389	0.1832	2.7328	0.8433	3
Months_Inactive_12_mon	0.3119	7.2349	0.1814	2.7059	0.5815	4
Contacts_Count_12_mon	0.1468	3.4042	0.1298	1.9332	0.8830	3
Customer_Age	0.1921	4.4557	0.1120	1.6711	0.5831	3
Income_Category	0.1892	4.3881	0.0942	1.4055	0.4980	4
Avg_Open_To_Buy	0.1726	4.0031	0.0867	1.2934	0.5024	2
Credit_Limit	0.1722	3.9942	0.0000	0	0.0000	3
Dependent_count	0.1105	2.5640	0.0000	0	0.0000	2
Education_Level	0.1048	2.4300	0.0000	0	0.0000	2

Fig 35: Variable importance

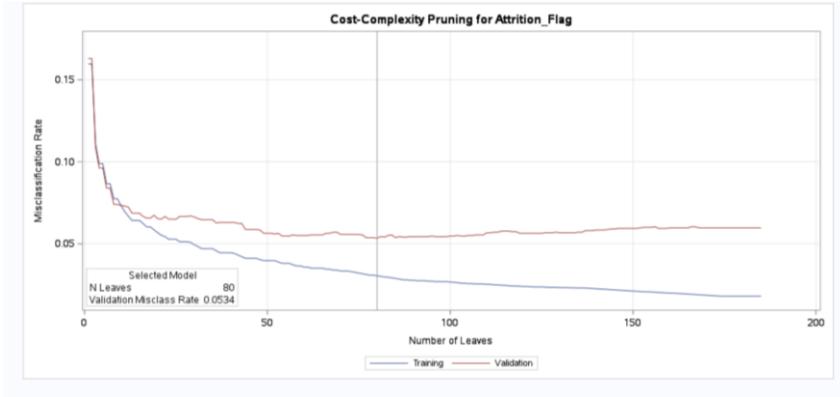


Fig 36: Cost complexity Pruning

4.2.1 Interpretation of result

- i. The data was randomly split into two (training and validation data) at the ratio of 70:30. The decision tree was developed on the bank_train data set.
- ii. The training data consist of 7114 observations while validation data consist of 3013 observations.
- iii. The tree has two branches with initial leaves of 185. However, this was pruned down to 80 leaves as shown by 30.
- iv. Figure 33 shows the confusion matrix and fit statistics. The final tree consists of 80 leaves. The training Entropy is 0.1381 while the validation Entropy is 0.1791. The model misclassification rate 0.0305 and 0.0534 for training and validation models respectively.
- v. Figure 35 shows the variable importance arranged in descending order of magnitude. The 5 most important variables are Total_trans_ct, Total_Revolving_bal, Total_trans_amt, Total_Relationship_count and Total_ct_chng_Q4_Q1.
- vi. The ROC curve show 99% area under the curve for training model and 94% for validation model. This implies a very high predictive ability of the model.

4.3 Random Forest

Random forest is a Supervised ML algorithm used extensively in Classification. It pools the results of multiple decision trees to arrive at a single result. Random forest is an extension of the bagging method as it employs both bagging and feature randomness to generate an uncorrelated forest of decision trees. It has 3 key hyperparameters, which are usually set up before training. The hyperparameters comprise the number of trees, node size, and the number of sample features.

PROC hpforest was used to implement this algorithm as follows:

```

330 /* RANDOM FOREST MODEL */
331 proc hforest data = cust_at.Attrition maxtrees = 100 seed = 12345 trainfraction=0.7 leafsize=5
332   vars_to_try=4 alpha= 0.1 maxdepth=50;
333   target Attrition_Flag / level = Nominal;
334   input Gender Education_Level
335     Marital_Status Income_Category Card_Category /level=nominal;
336   input Customer_Age Dependent_count Months_on_book Total_Relationship_Count Months_Inactive_12_mon
337     Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
338     Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
339     Avg_Utilization_Ratio /level=interval;
340   ods output FitStatistics = fit_at_runtime;
341   ods output VariableImportance = Variable_Importance;
342   ods output Baseline = Baseline;
343   save file = "/home/u63398813/model_fit.bin";
344   run;
345   title "The Average Square Error";
346   run;
347   /* visualizing the error/accuracy of the Model*/
348   proc sgplot data = fit_at_runtime;
349   series x=NTrees y=PredAll/legendlabel='Train Error';
350   series x=NTrees y=PredOOB/legendlabel='OOB Error';
351   xaxis values=(0 to 50 by 1);
352   yaxis values=(0 to 0.3 by 0.05) label='Average Square Error';
353   run;
354   /* visualizing the importance level of each variable included in the Model*/
355   proc sgplot data = Variable_Importance;
356   vbar Variable /response=Gini groupdisplay = cluster categoryorder=respdesc;
357   title "Variable Importance (Gini)";
358   run;

```

The results are presented in the diagrams below:

The HFOREST Procedure			
Performance Information			
Execution Mode	Single-Machine		
Number of Threads	2		
Data Access Information			
Data	Engine	Role	Path
CUST_ATATTRITION	V9	Input	On Client
Model Information			
Parameter	Value		
Variables to Try	4		
Maximum Trees	100		
Actual Trees	100		
Inbug Fraction	0.7		
Prune Fraction	0 (Default)		
Prune Threshold	0.1 (Default)		
Leaf Fraction	0.00001 (Default)		
Leaf Size Setting	5		
Leaf Size Used	5		
Category Bins	30 (Default)		
Interval Bins	100		
Minimum Category Size	5 (Default)		
Node Size	100000 (Default)		
Maximum Depth	50		
Alpha	0.1		
Exhaustive	5000 (Default)		
Rows of Sequence to Skip	5 (Default)		

Fig 37: Model Information

Baseline Fit Statistics						
Statistic	Value					
Average Square Error	0.135					
Misclassification Rate	0.161					
Log Loss	0.441					

Fit Statistics						
Number of Trees	Number of Leaves	Average Square Error (Train)	Average Square Error (OOB)	Misclassification Rate (Train)	Misclassification Rate (OOB)	Log Loss (Train)
1	102	0.0396	0.0501	0.0522	0.0562	0.162 0.313
2	203	0.0385	0.0602	0.0486	0.0774	0.139 0.342
3	286	0.0342	0.0556	0.0437	0.0993	0.127 0.294
4	382	0.0348	0.0567	0.0408	0.0719	0.130 0.272
5	443	0.0365	0.0578	0.0424	0.0751	0.140 0.244
6	557	0.0368	0.0574	0.0399	0.0758	0.141 0.247
7	645	0.0365	0.0566	0.0411	0.0728	0.140 0.225
8	745	0.0356	0.0544	0.0406	0.0704	0.139 0.213
9	854	0.0352	0.0525	0.0406	0.0674	0.137 0.196
10	957	0.0357	0.0530	0.0419	0.0662	0.138 0.200
11	1026	0.0351	0.0511	0.0411	0.0652	0.136 0.192
12	1151	0.0351	0.0506	0.0415	0.0646	0.136 0.191
13	1243	0.0349	0.0506	0.0407	0.0642	0.136 0.189
14	1345	0.0345	0.0492	0.0406	0.0620	0.135 0.180
15	1433	0.0346	0.0486	0.0403	0.0617	0.136 0.175
16	1542	0.0351	0.0490	0.0407	0.0634	0.138 0.177
17	1642	0.0349	0.0485	0.0405	0.0621	0.137 0.175

Fig 38: Fit Statistics

Loss Reduction Variable Importance					
Variable	Number of Rulers	Gini	OOB Gini	Margin	OOB Margin
Total_Trans_Ct	1098	0.040640	0.03638	0.081279	0.076630
Total_Trans_Amt	1576	0.032739	0.02737	0.065478	0.059833
Total_Revolving_Bal	798	0.022798	0.02024	0.045578	0.042690
Total_Ct_Chng_Q4_Q1	798	0.022014	0.01883	0.044023	0.040590
Avg_Utilization_Ratio	421	0.012608	0.01112	0.025217	0.023593
Total_Relationship_Count	669	0.012390	0.01018	0.024780	0.022554
Total_Amt_Chng_Q4_Q1	745	0.008414	0.00445	0.016829	0.012822
Contacts_Count_12_mon	440	0.004952	0.00324	0.009194	0.007824
Months_Inactive_12_mon	531	0.004266	0.00287	0.008532	0.007105
Gender	388	0.002500	0.00171	0.005000	0.004239
Customer_Age	411	0.002376	0.00087	0.004742	0.003197
Dependent_count	312	0.001652	0.00036	0.003304	0.001967
Marital_Status	211	0.000789	0.00015	0.001578	0.001172
Credit_Limit	284	0.001174	0.00000	0.002348	0.001238
Avg_Open_To_Buy	306	0.001233	0.00001	0.002466	0.001249
Card_Category	111	0.000139	-0.00002	0.000277	0.000250
Dependent_count	191	0.000632	-0.00021	0.001285	0.000437
Education_Level	129	0.000253	-0.00035	0.000505	0.000065
Income_Category	158	0.000500	-0.00035	0.001000	0.000219

Fig 39: Loss reduction variable importance

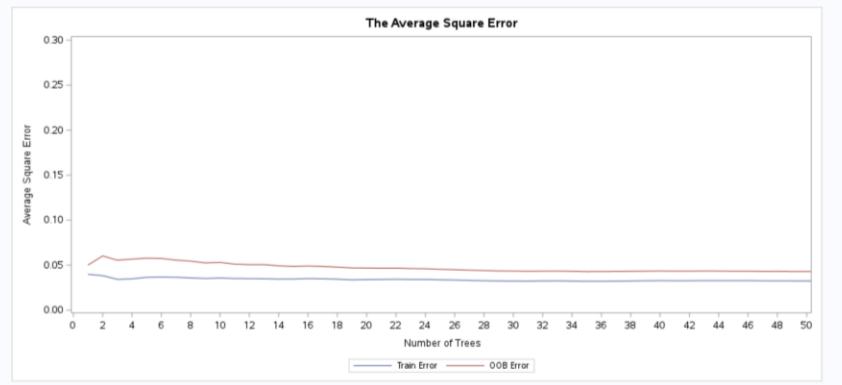


Fig 40: The Average Square Error (ASE)

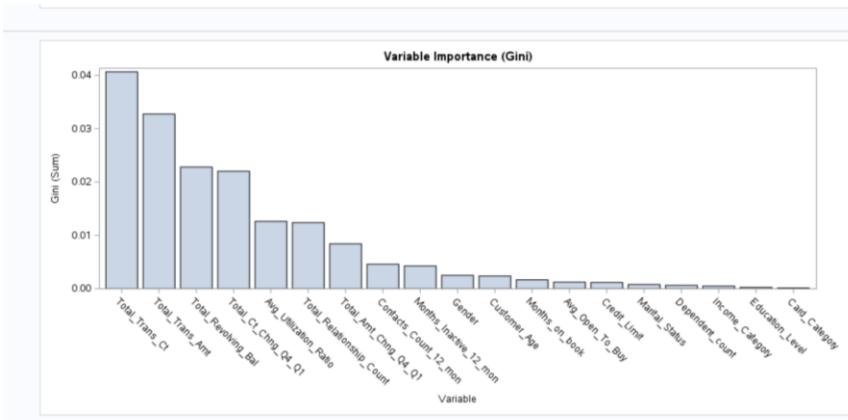


Fig 41: Variable importance

4.3.1 Interpretation of result

- i. From Fig 37, there are 4 variables selected for the possible split of tree or to test each node and the maximum number of tree was 100.
- ii. From the Baseline fit statistics (fig 38), misclassification rate is 0.161 or 16.1%. This means that 83.9% was correctly classified.
- iii. The model also shows the Average Square error for the train and OOB (out-of-bag-evaluation) data (fig 40). The results show that model has misclassification reached its lowest point when the tree number is 100.
- iv. From the Variable importance diagram (fig 41), the contribution of each variable to the classification of Attrition flag has been shown in order of importance using a bar chart. The chart shows that Total_trans_ct is of highest importance. This is followed by Total_trans_Amt while Card category is of least importance.

5. Comparison of Predictive Accuracy powers of the models.

The Accuracy levels of all the ML algorithms are evaluated in table 1 below:

Predictive Accuracy powers of the models		
S/N	Model	Accuracy level
1	Logistic regression	92.15%
2	Decision tree	94%
3	Random forest	83.90%

Table 1: Model Prediction power

The predictive powers of the models in table 1 above show that Decision tree with 94% has the highest predictive power. It is therefore recommended that Decision tree model should be given preference over the other models in the prediction of credit card customer churn.

6. Critical comparison of SAS and Python programing.

6.1 Introduction

In this section of the study, the SAS (including SAS Ondemand for Academics) and Python programing tools were critically investigated. Both were used to analyze the dataset hence a comparison of their efficiency levels is made based on defined criteria.

6.2 Overview of SAS.

SAS stands for Statistical Analysis Software. SAS stands for Statistical Analysis Software. It was created in the year 1960 by the SAS Institute. From 1st January 1960, SAS was used for data management, business intelligence, Predictive Analysis, Descriptive and Prescriptive Analysis etc. It is employed to explore, visualize and analyze data at high levels of data granularity, so as to extract insight for informed decision making. Unlike open source data analytics tools such as R, Python, Apache Spark, etc, SAS is a proprietary or commercial software. However, it offers some open source platforms such as the SAS OnDemand for Academics which comes at no cost.

6.3 Overview of Python programing language

Python is a powerful open source and easy to learn programming language. Python is a general-purpose language, meaning it can be used to create a variety of different programs and everyday tasks. It could be used for developing websites and software, task automation, data analytics, and data visualization. It has a well-designed syntax and together with its interpreted nature (Jupiter Notebook), makes it a superlative language for scripting and application development in many areas. Python can be used as a scripting language or can be compiled to byte-code for building large applications. It can also be easily integrated with C, C++, COM, ActiveX, CORBA, and Java.

6.4 Review of SAS and Python programing

- i. **Installation and setup:** For SAS the right to set up would be bought from the company however, the SAS Ondemand version requires only creating an online account. Python requires downloading the Jupiter Notebook program from <https://www.python.org/>.

- ii. **Navigation and User experience.** Both involve understanding the syntax or a set of rules associated with each. The SAS uses PROC statements while Python uses physical lines as well as logical lines or statements. Both require some form of learning and mastering. However, the SAS has some form of point and click features which make it a bit easier for those who find the codes more difficult to learn.
- iii. **Data mining:** Python provides a robust platform for data mining. Data importation, cleaning, manipulation and visualization could be handled with simple codes. It also produces powerful visuals that make for enriching and insightful reports. On the other hand, the process of data importation in SAS is a bit more complex. File conversion and library creation make it somewhat cumbersome. Besides, certain PROC cannot run on the SAS Ondemand for Academics .
- iv. **Scalability:** Python provides a better structure and support for large programs than shell scripting. SAS Ondemand for Academics has a smaller processing capacity than Python. For SAS Ondemand, any result greater than 3MB will result in some form of malfunction while Python can produce results far more than 3MB without any issues. The main SAS can run higher capacity processing.
- v. **Statistical Analysis and Machine learning.** Both can handle statistical analysis however, Python have an edge over SAS as it could handle more complex statistical models and machine learning modeling with ease
- vi. **Database management:** Python supports a wide range of other third party platforms and databases. It is easily integrated and can run on a wide variety of hardware platforms maintaining the same same interface on all platforms. SAS is more restrictive due to proprietary rights.

vii. **Output.**

The output quality of SAS and Python are compared in the table below. The highest score is 5 while the lowest is 1.

Output scoring		
Criteria	SAS Ondemand for Academics	Python
Installation and setup		
Ease of installation	3	4
Platform Dependence	4	5
Navigation and user experience		
User Friendliness	4	4
Easy of navigation	5	4
Efficiency	4	5
Data mining		
Data importation	3	5
Data cleaning	4	5
Data combination	4	4
Data sub-setting	4	4
Data transformation	4	4
Other criteria		
Scalability	3	5
Speed	3	5
Statistical analysis and machine learning modeling	4	5
Output	4	5
Data visualization	4	5
Reporting	5	4
Integration with other platforms	4	5
Storage, retrieval and sharing	4	4
Data base management system adaptability	4	5
Total	72	82

Table 2. SAS Ondemand for Academics and Python scoring

Table 2 indicates that overall Python scored more than SAS Ondemand for Academics. However, a good knowledge of the two is very relevant so that one could serve as a backup. For businesses, such decision may depend on the strategic goals of the company.

7. Conclusion

This study investigated customer attrition in a credit card scheme with a view to providing a predictive system that would correctly predict customer attrition. Such prediction will provide financial managers pragmatic decision making basis in eliminating or reducing customer attrition.

Three Machine Learning models Logistic regression, Decision tree and Random forest were deployed to investigate the phenomenon. The result shows that Decision tree performed better than the other models

in making the required prediction. Decision tree provided 94% predictive accuracy hence should be the preferred model. Logistic regression and Random forest provided 92.15% and 83.9 predictive accuracy respectively. It is believed that the adoption of this result will offer the bank (financial institution) opportunity to put measure in place that will improve their services so as to effectively retain their customers. This will ultimately improve the profitability and market share of the bank.

8. References

Jagadeesan A.P and Indhuja. R, 2020. Bank Customer Retention Prediction And Customer Ranking Based On Deep Neural Networks. *International Journal of Scientific Development and Research (IJSR) Volume 5 Issue 9*

Kumar, D.A.; Ravi, V.(2008) Predicting credit card customer churn in banks using data mining. *Int. J. Data Anal. Tech. Strateg.* **2008**, 1, 4–28

Miao, X.; Wang, H.(2022) Customer churn prediction on credit card services using random forest method. In Proceedings of the 2022 7th International Conference on Financial Innovation and Economic Development (ICFIED 2022), Online, 14–16 January 2022; Atlantis Press: Paris, France; pp. 649–656.

Parangi, A. (2022). Predicting Credit Card Customer Churn. Accessed 11/05/2023.
<https://www.akkio.com/post/credit-card-churn-prediction>

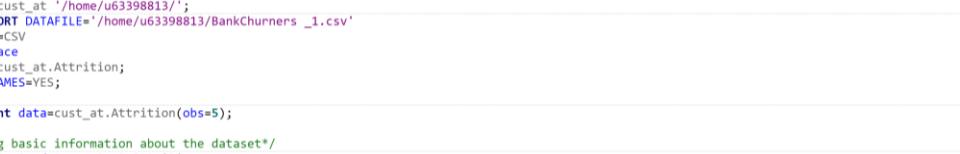
Waykar, Y. (2013) A Study of Importance of UML diagrams: With Special Reference to Very Large-sized Projects. Conference paper at International Conference on Reinventing Thinking beyond boundaries to Excel At: FARIDABAD, INDIA

Zhyli. (2020). Prediction of Churning Credit Card Customers [Data set]. Zenodo.
<https://doi.org/10.5281/zenodo.4322342>

<https://olvy.co>blog>customer.churn.analysis>

9. Appendix 1 : SAS codes

9.1 Data pre-processing and mining



The screenshot shows a SAS Enterprise Guide interface with the following details:

- Program 1.sas** is the active file.
- Import Data 1.ctl** is the current step.
- CODE** tab is selected.
- LOG**, **RESULTS**, and **OUTPUT DATA** tabs are also present.
- Toolbar icons** include: Run, Stop, Save, Print, Copy, Paste, Find, Replace, Undo, Redo, and Help.
- Line #** column header is visible.
- Code content:**

```
1 |Libname cust_at '/home/u63398813/';
2 PROC IMPORT DATAFILE='/home/u63398813/BankChurners _1.csv'
3   DBMS=CSV
4   replace
5   OUT=cust_at.Attrition;
6   GETNAMES=YES;
7 RUN;
8 proc print data=cust_at.Attrition(obs=5);
9 run;
10 /*Getting basic information about the dataset*/
11 Proc Contents data=cust_at.Attrition;
12 run;
13 /*checking for missing numbers in numerical variables*/
14 proc means data=cust_at.attrition
15   NMISS;
16 run;
17 /*checking for missing numbers in character variables*/
18 proc sql;
19   select nmiss(Attrition_Flag) as Attrition_Flag_miss, nmiss(Gender) as Gender_miss, nmiss(Education_Level) as Educational_level_miss, nmiss(Marital_Status)as Marital_Status_miss
20   from cust_at.attrition;
21 quit;
22 /*checking for the number of churned customers and number of existing customers*/
23 PROC FREQ DATA=cust_at.attrition;
24 Table Attrition_Flag;
25 run;
26 /*summary of the numerical variables*/
27 proc means data = cust_at.attrition min q1 median mean q3 max n maxdec=2;
```

The screenshot shows the SAS Enterprise Guide interface with the 'Program 1.sas' file open. The code editor displays the following SAS script:

```
27 proc means data = cust_at.attrition min q1 median mean q3 max n maxdec=2;
28 run;
29 title "Summary Statistics for numerical variables in Customer Churn Data";
30 run;
31
32 proc means data = cust_at.attrition Std;
33 run;
34 /*summary of the categorical variables*/
35 proc freq data=cust_at_attrition;
36 tables Attrition_Flag Gender Education_Level Marital_Status Income_Category Card_Category;
37 run;
38 /*dropping a variable not needed in the analysis*/
39 data cust_at_attrition;
40   set cust_at_attrition (drop = CLIENTNUM);
41 run;
42 proc print data = cust_at.attrition(obs=5);
43 run;
44 ods graphics on;
45 /*Creating histogram with proc univariate*/
46 proc univariate data=cust_at_attrition noprint;
47 var Total_Amt_Chng_Q4_Q1 Total_Ct_Chng_Q4_Q1 ;
48 histogram/normal(color=CXe7185d);
49 run;
50 proc univariate data=cust_at_attrition noprint;
51 var Dependent_count Months_on_book ;
52 histogram/normal(color=CXe7185d);
53 run;
```

9.2: Data Visualization

Program 1.sas Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```
50 proc univariate data=cust_at_attrition noprint;
51 var Dependent_count Months_on_book ;
52 histogram/normal(color=xCe7185d);
53 run;
54
55 /*Creating a Bchart of the Education level variable*/
56 proc sgplot data= cust_at.attrition;
57 vbar Education_Level/group=Education_Level groupdisplay=cluster;
58 title height=14pt "Bchart of Educational level";
59 run;
60
61 /*Creating a Scatter plot of Credit Limit and Avg_Open_To_Buy*/
62 proc sgplot data= cust_at.attrition;
63 scatter y=Credit_Limit x=Avg_Open_To_Buy;
64 xaxis grid;
65 yaxis grid;
66 title height=14pt "Scatterplot of Credit limit and Avg_Open_To_Buy";
67 run;
68 /*Creating a Scatter plot of Credit Limit and Avg_Utilization_Ratio*/
69 proc sgplot data= cust_at.attrition;
70 scatter y=Credit_Limit x=Avg_Utilization_Ratio;
71 xaxis grid;
72 yaxis grid;
73 title height=14pt "Scatterplot of Credit Limit and Avg_Utilization_Ratio";
74 run;
75 /*Creating a Scatter plot of Credit Limit and Avg_Utilization_Ratio with ellipse*/
76 proc sgplot data= cust_at.attrition;
77
```

/home/u63398813/Program 1.sas

Line 18, Column 1

Program 1.sas Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```
78 run;
79 /*Creating a Scatter plot of Credit Limit and Avg_Utilization_Ratio with ellipse*/
80 proc sgplot data= cust_at.attrition;
81 scatter y=Credit_Limit x=Avg_Utilization_Ratio;
82 xaxis grid;
83 yaxis grid;
84 ellipse y=Credit_Limit x=Avg_Utilization_Ratio;
85 title height=14pt "Scatterplot of Credit limit Avg_Utilization_Ratio with ellipse";
86 run;
87
88 /*Creating a Scatter plot of Total_Trans_Amt and Total_Trans_Ct*/
89 proc sgplot data= cust_at.attrition;
90 scatter y=Total_Trans_Amt x=Total_Trans_Ct;
91 xaxis grid;
92 yaxis grid;
93 title height=14pt "Scatterplot of Total_Trans_Ct and Total_Trans_Amt";
94 run;
95
96 /*Creating a Scatter plot of Customer_Age and Months_on_book*/
97 proc sgplot data= cust_at.attrition;
98 scatter y=Customer_Age x=Months_on_book;
99 xaxis grid;
100 yaxis grid;
101 title height=14pt "Scatterplot of Customer_Age and Months_on_book with ellipse";
102 ellipse y=Customer_Age x=Months_on_book;
103 run;
104
```

/home/u63398813/Program 1.sas

Line 18, Co

Program 1.sas Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```

244 run;
245
246 /*calculate correlation matrix for the data*/;
247 ods graphics / reset width=10.4in height=4.8in imagemap;
248 ods output PearsonCorr=Corr_P;
249 PROC CORR DATA=CUST_AT_ATTRITION;
250   VAR Customer_Age Dependent_count
251     Months_on_book Total_Relationship_Count Months_Inactive_12_mon
252     Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
253     Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
254     Avg_Utilization_Ratio ;
255 RUN;
256 /*sort for transposes;
257 proc sort data=Corr_P;
258   by VARIABLE;
259 run;
260 /*restructure data so that it's in a long format for graphing;
261 *need to transpose the correlation and p-values separately;
262 proc transpose data=Corr_P out=CorrLong(rename=(COL1=Correlation)
263   name=CorrelationID;
264   var Customer_Age
265     Months_on_book Total_Relationship_Count Months_Inactive_12_mon
266     Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
267     Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
268     Avg_Utilization_Ratio Dependent_count;
269   by Variable;
270 run;
271
```

/home/u63398813/Program 1.sas

9.3: Machine learning Modeling

Program 1.sas Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```

285
286 /*sort for graphing;
287 proc sort data=CorrLong;
288   by VARIABLE CorrelationID;
289 run;
290
291 /*create a heat map graph with correlation values in the squares;
292 proc sgplot data=CorrLong noautolegend;
293 heatmap x=Variable y=CorrelationID / colorresponse=Correlation name="nope1" discretex discretey xaxis colormodel=ThreeColorRamp; /*Colorresponse allows discrete squares
294 text x=Variable y=CorrelationID text=p_value / textattr=(size=10pt) xaxis name='nope2';
295 yaxis reverse display=(nolabel);
296 xaxis display=(nolabel);
297 gradlegend;
298 title height=14pt "Correlation matrix";
299 run;
300
301
302 /*Model Building*/
303 /*Splitting data into training and test for model building*/
304 /* Split data into two datasets : 70%- training 30%- validation*/
305 Proc Surveyselect data=CUST_AT_ATTRITION out=split seed=2345 samprate=.7 outall;
306 Run;
307 Data training validation;
308 Set split;
309 if selected = 1 then output training;
310 else output validation;
311 Run;
```

6

Program 1.sas | Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```

311 /* Logistic regression Model*/
312 /* Logistic regression Model*/
313 ods graphics on;
314 Proc Logistic Data = training descending;
315 class Attrition_Flag Gender Education_Level Marital_Status Income_Category Card_Category/ param = ref;
316 Model Attrition_Flag = Gender Education_Level
317 Marital_Status Income_Category Card_Category Customer_Age Dependent_count
318 Months_on_book Total_Relationship_Count Months_Inactive_12_mon
319 Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
320 Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
321 Avg_Utilization_Ratio / selection = stepwise slstay=<0.15 slentry=<0.15 stb;
322 score data=training out = Logit_Training fitstat outroc=roc;
323 score data=validation out = Logit_Validation fitstat outroc=roc;
324 Run;
325 /* Decision Tree Model*/
326 proc hpsplit data= cust_at Attrition;
327 class Attrition_Flag Gender Education_Level Marital_Status Income_Category Card_Category _CHARACTER_;
328 Model Attrition_Flag = Gender Education_Level
329 Marital_Status Income_Category Card_Category Customer_Age Dependent_count
330 Months_on_book Total_Relationship_Count Months_Inactive_12_mon
331 Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
332 Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
333 Avg_Utilization_Ratio ;
334 PARTITION FRACTION(VALIDATE=0.3 SEED=12345);
335 grow entropy;
336 prune costcomplexity;
337 run;
338 
```

5

Program 1.sas | Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```

335 /* Random Forest Model*/
336 PARTITION FRACTION(VALIDATE=0.3 SEED=12345);
337 grow entropy;
338 prune costcomplexity;
339 run;
340 /* Random Forest Model*/
341 proc hforest data = cust_at Attrition maxtrees = 100 seed = 12345 trainfraction=0.7 leafsize=5
342 vars_to_try=4 alpha= 0.1 maxdepth=50;
343 target Attrition_Flag / level = Nominal;
344 input Gender Education_Level
345 Marital_Status Income_Category Card_Category /level=nominal;
346 input Customer_Age Dependent_count Months_on_book Total_Relationship_Count Months_Inactive_12_mon
347 Contacts_Count_12_mon Credit_Limit Total_Revolving_Bal Avg_Open_To_Buy
348 Total_Amt_Chng_Q4_Q1 Total_Trans_Amt Total_Trans_Ct Total_Ct_Chng_Q4_Q1
349 Avg_Utilization_Ratio /level=interval;
350 ods output FitStatistics = fit_at_runtime;
351 ods output VariableImportance = Variable_Importance;
352 ods output Baseline = Baseline;
353 save file = "/home/u63398813/model_fit.bin";
354 run;
355 title "The Average Square Error";
356 /* visualizing the error/accuracy of the Model*/
357 proc sgplot data = fit_at_runtime;
358 series x=NTrees y=PredAll/legendlable='Train Error';
359 series x=NTrees y=PredOOB/legendlable='OOB Error';
360 xaxis values=(0 to 50 by 1);
361 yaxis values=(0 to 0.3 by 0.05) label='Average Square Error';
362 
```

6

Program 1.sas | Import Data 1.ctl

CODE LOG RESULTS OUTPUT DATA

```

347 Avg_Utilization_Ratio /level=interval;
348 ods output FitStatistics = fit_at_runtime;
349 ods output VariableImportance = Variable_Importance;
350 ods output Baseline = Baseline;
351 save file = "/home/u63398813/model_fit.bin";
352 run;
353 title "The Average Square Error";
354 /* visualizing the error/accuracy of the Model*/
355 proc sgplot data = fit_at_runtime;
356 series x=NTrees y=PredAll/legendlable='Train Error';
357 series x=NTrees y=PredOOB/legendlable='OOB Error';
358 xaxis values=(0 to 50 by 1);
359 yaxis values=(0 to 0.3 by 0.05) label='Average Square Error';
360 run;
361 /* visualizing the importance level of each variable included in the Model*/
362
363 proc sgplot data = Variable_Importance;
364 vbar Variable /response=Gini groupdisplay = cluster categoryorder=respdesc;
365 title "Variable Importance (Gini)";
366 run;
367 /* Making prediction with the Model*/
368 proc hp4score data=validation;
369 score file= "/home/u63398813/model_fit.bin"
370 out=Predictions;
371 run;
372 proc print data=Predictions (obs=50);
373 run;
374 
```

10. Appendix 2: Python codes and results

10.1 Data pre-processing and mining

The screenshot shows a Jupyter Notebook interface with two code cells and one data preview cell.

In [4]:

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

import warnings
warnings.simplefilter("ignore")
```

In [6]:

```
import warnings
warnings.filterwarnings("ignore")
from datetime import datetime, timedelta
```

In [15]:

```
df=pd.read_csv("BankChurners _1.csv")
```

In [16]:

```
df.head()
```

Out[16]:

	CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book
0	768805383	Existing Customer	45	M	3	High School	Married	60K-80K	Blue	39
1	818770008	Existing Customer	49	F	5	Graduate	Single	Less than \$40K	Blue	44
2	713982108	Existing Customer	51	M	3	Graduate	Married	80K-120K	Blue	36
3	709911858	Existing Customer	40	F	4	High School	Unknown	Less than \$40K	Blue	34
4	715116494	Existing	46	M	3	Undergraduated	Married	60K-80K	Blue	21

In [29]:

```
df1.info()
```

Out[29]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10127 entries, 0 to 10126
Data columns (total 20 columns):
 #   Column           Non-Null Count  Dtype  
 --- 
 0   Attrition_Flag    10127 non-null   object 
 1   Customer_Age      10127 non-null   int64  
 2   Gender            10127 non-null   object 
 3   Dependent_Count   10127 non-null   int64  
 4   Education_Level   10127 non-null   object 
 5   Marital_Status     10127 non-null   object 
 6   Income_Category    10127 non-null   object 
 7   Card_Category      10127 non-null   object 
 8   Months_on_book    10127 non-null   int64  
 9   Total_Relationship_Count 10127 non-null   int64  
 10  Months_Inactive_12_mon 10127 non-null   int64  
 11  Contacts_Count_12_mon 10127 non-null   int64  
 12  Credit_Limit       10127 non-null   float64 
 13  Total_Revolving_Bal 10127 non-null   int64  
 14  Avg_Open_To_Buy    10127 non-null   float64 
 15  Total_Amt_Chng_Q4_Q1 10127 non-null   float64 
 16  Total_Trans_Amt    10127 non-null   int64  
 17  Total_Trans_Ct     10127 non-null   int64  
 18  Total_Ct_Chng_Q4_Q1 10127 non-null   float64 
 19  Avg_Utilization_Ratio 10127 non-null   float64 
dtypes: float64(5), int64(9), object(6)
memory usage: 1.5+ MB
```

In []:

```
#checking for missing values
```

In [30]:

```
df1.isnull().sum()
```

Out[30]:

```
Attrition_Flag      0
Customer_Age        0
Dependent_Count    0
Education_Level     0
Gender              0
Income_Category     0
Marital_Status      0
Card_Category        0
Months_on_book      0
Total_Relationship_Count 0
Months_Inactive_12_mon 0
Contacts_Count_12_mon 0
Credit_Limit         0
Total_Revolving_Bal 0
Avg_Open_To_Buy     0
Total_Amt_Chng_Q4_Q1 0
Total_Trans_Amt      0
Total_Trans_Ct        0
Total_Ct_Chng_Q4_Q1 0
Avg_Utilization_Ratio 0
```

lhost:8888/notebooks/Untitled8-Copy1.ipynb#

jupyter Untitled8-Copy1 Last Checkpoint: 16 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

memory usage: 1.5+ MB

```
In [ ]: #checking for missing values
```

```
In [30]: df1.isnull().sum()
```

```
Out[30]: Attrition_Flag      0
Customer_Age          0
Gender                0
Dependent_count       0
Education_Level        0
Marital_Status         0
Income_Category        0
Card_Category          0
Months_on_book         0
Total_Relationship_Count 0
Months_Inactive_12_mon 0
Contacts_Count_12_mon 0
Credit_Limit            0
Total_Revolving_Bal    0
Avg_Open_To_Buy         0
Total_Amt_Chng_Q4_Q1    0
Total_Trans_Amt          0
Total_Trans_Ct           0
Total_Ct_Chng_Q4_Q1     0
Avg_Utilization_Ratio   0
dtype: int64
```

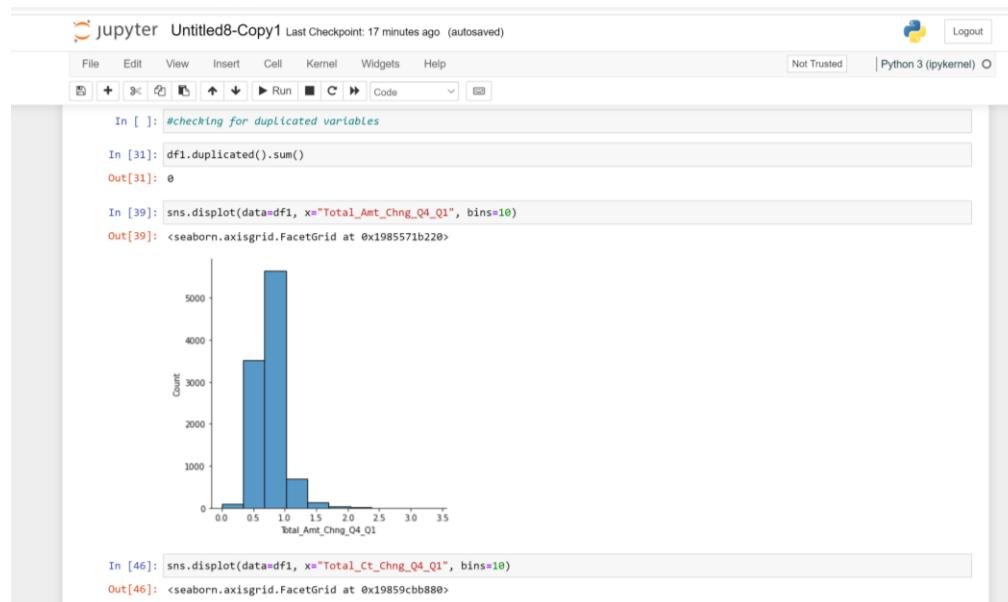
```
In [ ]: #checking for duplicated variables
```

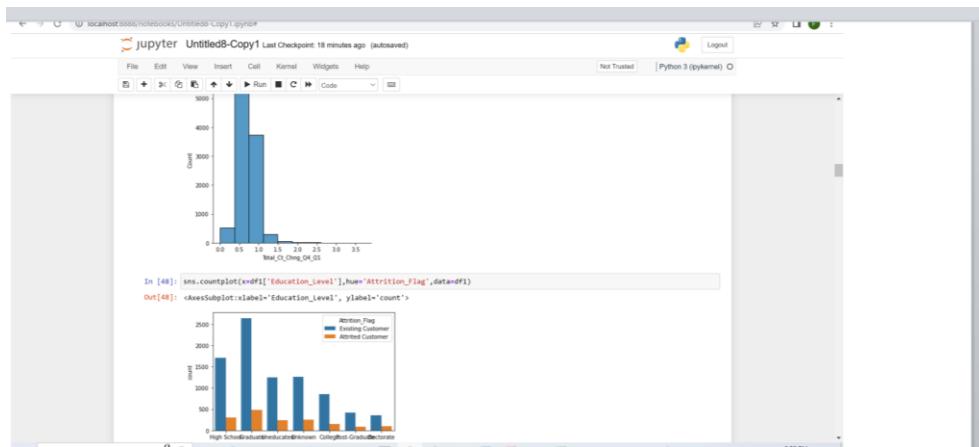
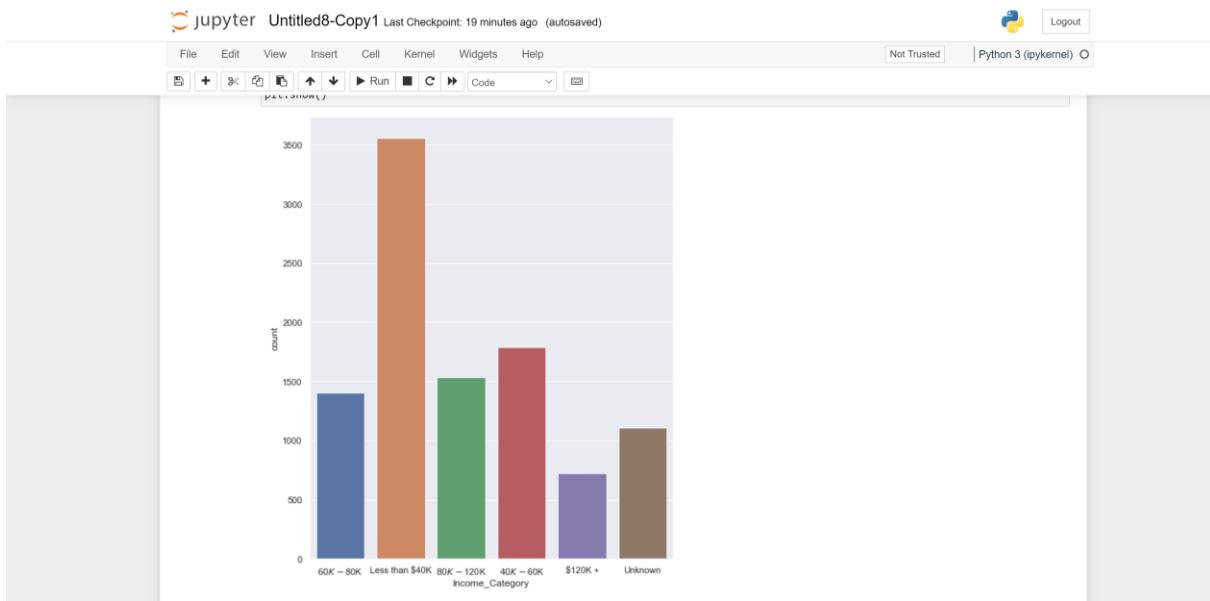
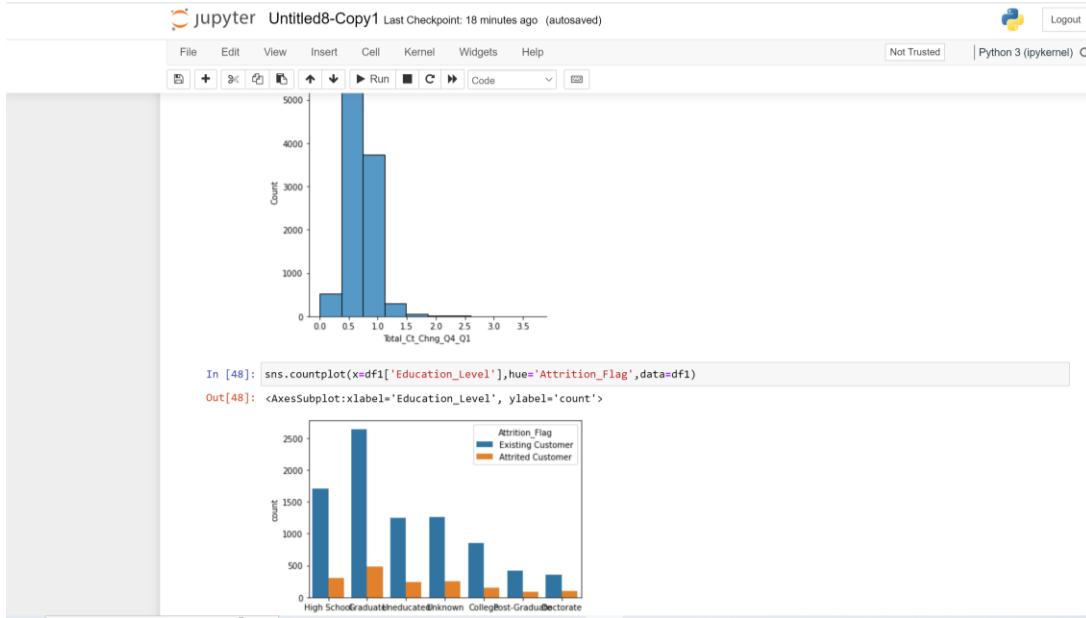
```
In [31]: df1.duplicated().sum()
```

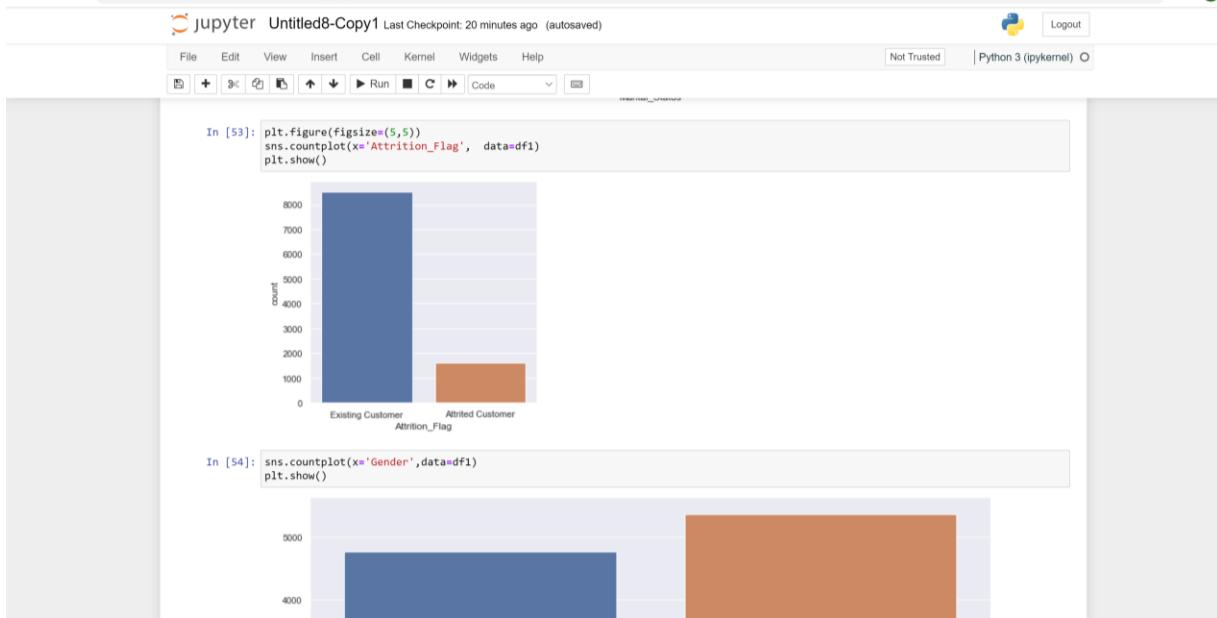
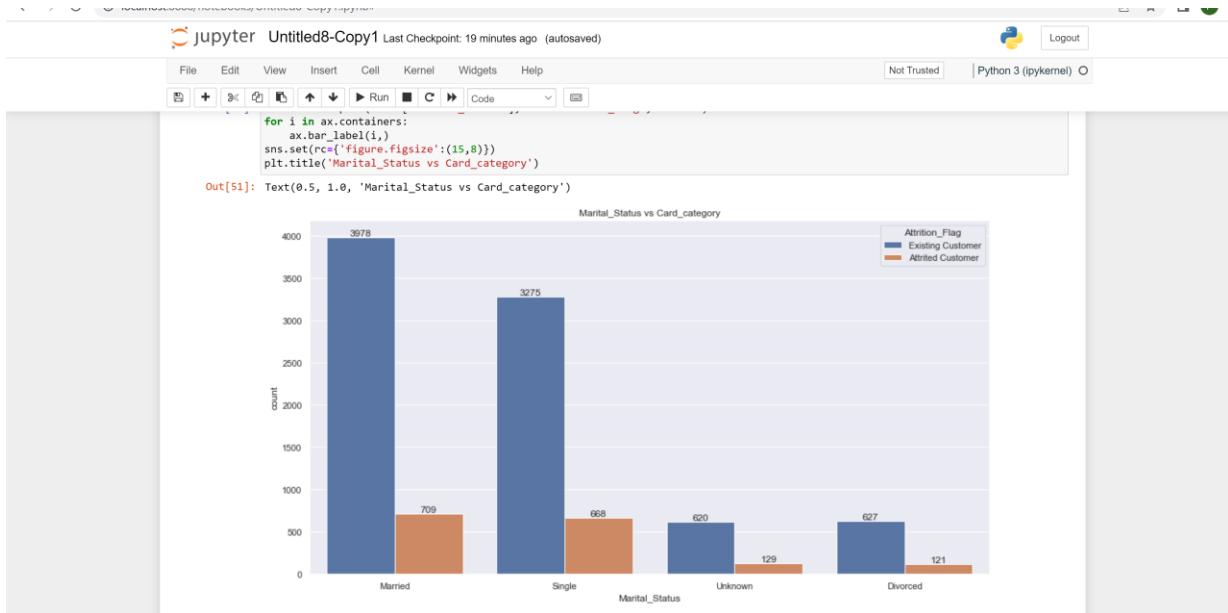
```
Out[31]: 0
```

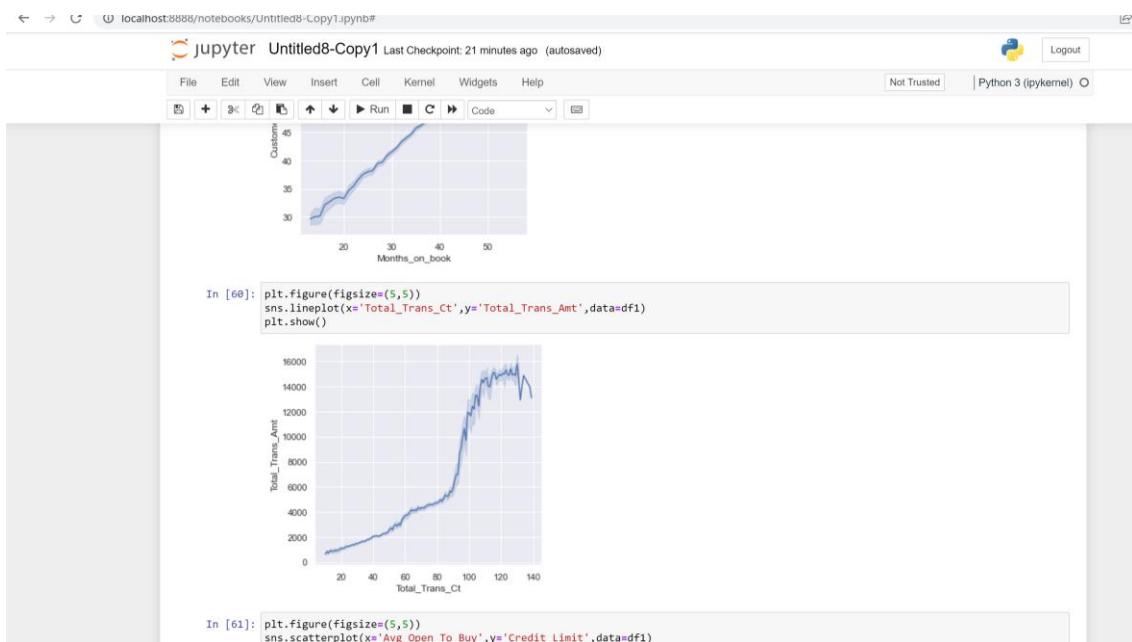
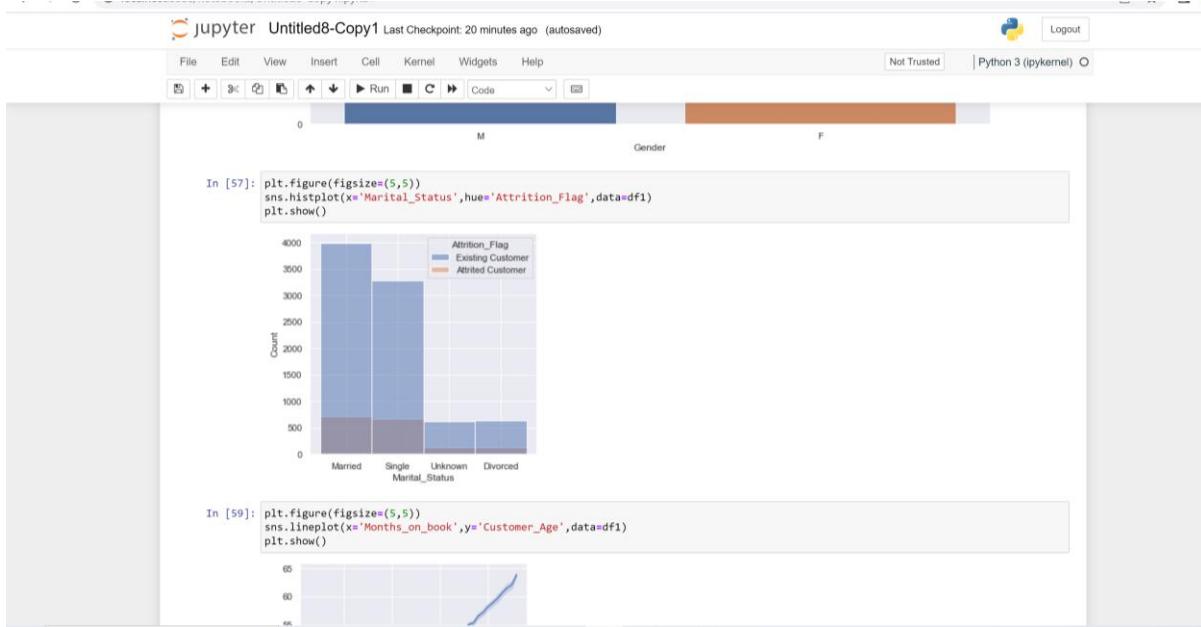
```
In [39]: sns.displot(data=df1, x="Total_Amt_Chng_Q4_Q1", bins=10)
```

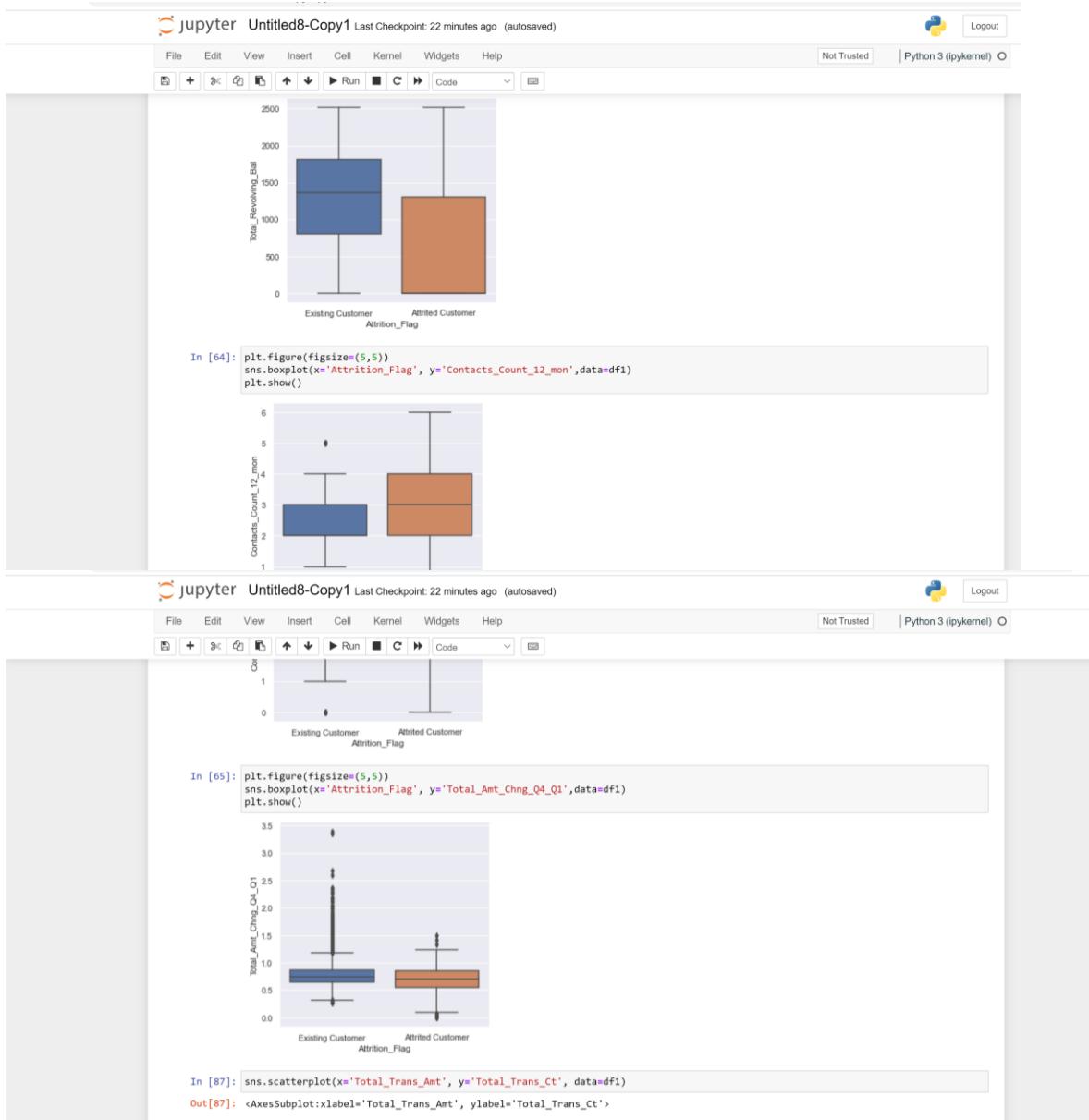
10.2: Data Visualization

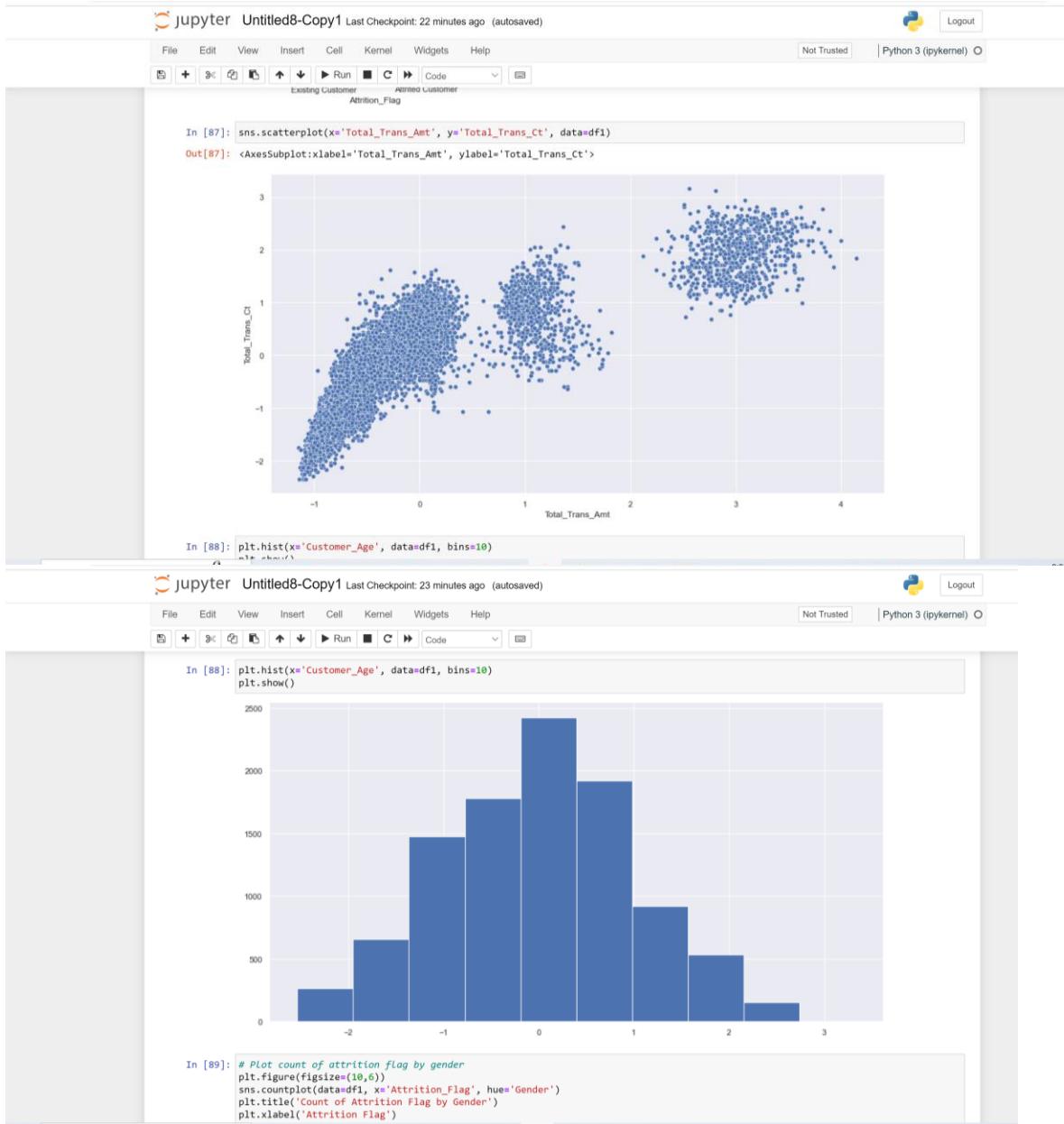


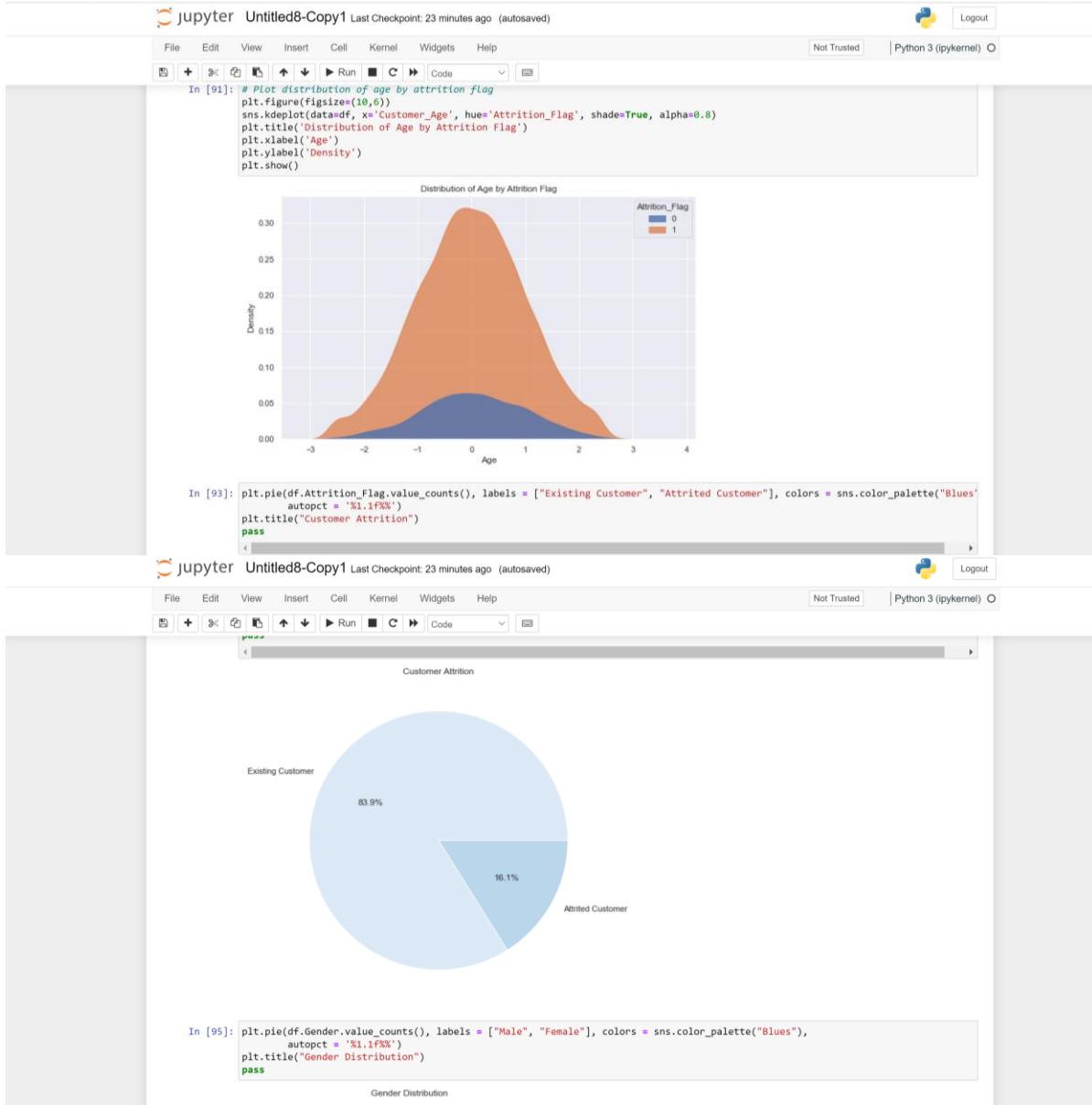


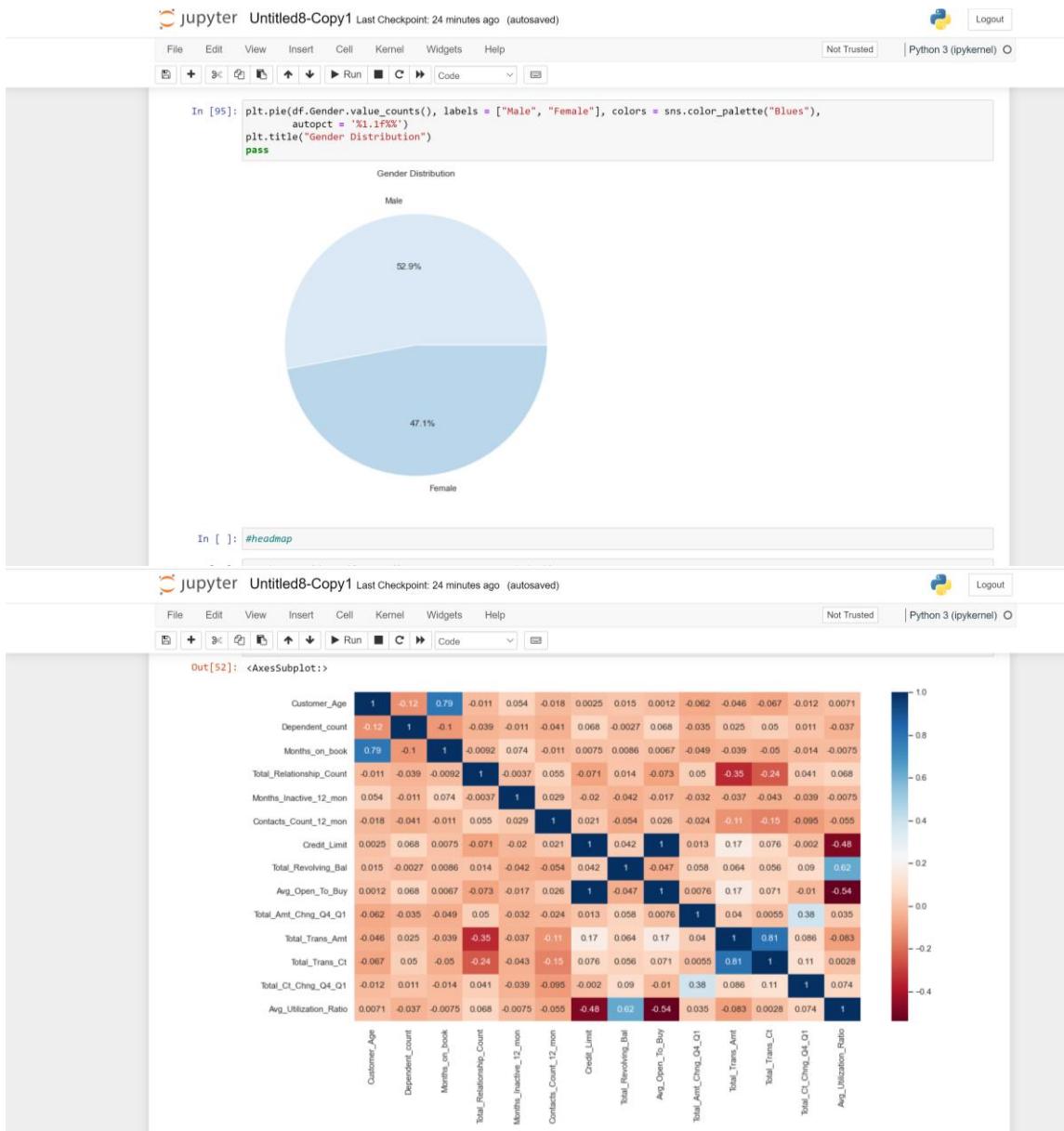












10.3: Machine learning Modeling

The screenshot shows two Jupyter Notebook sessions. The top session (In [68] to Out[73]) demonstrates standardizing numerical data types. The bottom session (In [74] to Out[76]) shows encoding categorical columns using LabelEncoder.

Session 1 (Top): Standardizing Data

```
In [68]: df.select_dtypes(['int64']).columns
Out[68]: Index(['CLIENTNUM', 'Customer_Age', 'Dependent_count', 'Months_on_book',
       'Total_Relationship_Count', 'Months_Inactive_12_mon',
       'Contacts_Count_12_mon', 'Total_Revolving_Bal', 'Total_Trans_Amt',
       'Total_Trans_Ct'], dtype='object')

In [71]: from sklearn.preprocessing import StandardScaler
def standardize_int_columns(data):
    int_columns = df1.select_dtypes(['int64']).columns
    sd = StandardScaler()
    for col in int_columns:
        data[col] = sd.fit_transform(data[col].values.reshape(-1,1))
    return data

In [72]: df=standardize_int_columns(df1)

In [73]: df.select_dtypes(include='object')
```

Session 2 (Bottom): Encoding Category Columns

```
In [74]: # encoding category columns
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
df1['Attrition_Flag']=le.fit_transform(df1['Attrition_Flag'])
df1['Gender']=le.fit_transform(df1['Gender'])
df1['Education_Level']=le.fit_transform(df1['Education_Level'])
df1['Marital_Status']=le.fit_transform(df1['Marital_Status'])
df1['Income_Category']=le.fit_transform(df1['Income_Category'])
df1['Card_Category']=le.fit_transform(df1['Card_Category'])

In [75]: df1.head()
```

	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book	Total_Relations
0	1	-0.165406	1	0.503368	3	1	2	0	0.384621	
1	1	0.333570	0	2.043199	2	2	4	0	1.010715	
2	1	0.583058	1	0.503368	2	1	3	0	0.008965	
3	1	-0.789126	0	1.273283	3	3	4	0	-0.241473	
4	1	-0.789126	1	0.503368	5	1	2	0	-1.869317	

```
In [76]: df1.columns
Out[76]: 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')
```

jupyter Untitled8-Copy1 Last Checkpoint: 26 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In [78]: `y=df1['Attrition_Flag']`

In [79]: `from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=42,test_size=0.3)`

In [80]: `y_train.value_counts()`

Out[80]:

1	5957
0	1131
Name: Attrition_Flag, dtype: int64	

In [81]: `from sklearn.linear_model import LogisticRegression
lr=LogisticRegression()
lr.fit(x_train,y_train)
lr_pred=lr.predict(x_test)`

In [82]: `from sklearn.metrics import confusion_matrix,classification_report
print('Classification Report:\n',classification_report(y_test,lr_pred))
print('Confusion Matrix:\n',confusion_matrix(y_test,lr_pred))`

Classification Report:

	precision	recall	f1-score	support
0	0.71	0.28	0.40	496
1	0.87	0.98	0.92	2543
accuracy			0.86	3039
macro avg	0.79	0.63	0.66	3039
weighted avg	0.85	0.86	0.84	3039

Confusion Matrix:

[[146 356]
[57 2486]]

In [83]: `from sklearn.tree import DecisionTreeClassifier
ds=DecisionTreeClassifier()`

jupyter Untitled8-Copy1 Last Checkpoint: 26 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In [83]: `from sklearn.tree import DecisionTreeClassifier
ds=DecisionTreeClassifier()
ds.fit(x_train,y_train)
ds_pred=ds.predict(x_test)`

In [84]: `from sklearn.metrics import confusion_matrix,classification_report
print('Classification Report:\n',classification_report(y_test,ds_pred))
print('Confusion Matrix:\n',confusion_matrix(y_test,ds_pred))`

Classification Report:

	precision	recall	f1-score	support
0	0.78	0.73	0.75	496
1	0.95	0.96	0.95	2543
accuracy			0.92	3039
macro avg	0.86	0.84	0.85	3039
weighted avg	0.92	0.92	0.92	3039

Confusion Matrix:

[[360 136]
[101 2442]]

In [85]: `from sklearn.ensemble import RandomForestClassifier
rf=RandomForestClassifier()
rf.fit(x_train,y_train)
rf_pred=rf.predict(x_test)`

In [86]: `from sklearn.metrics import confusion_matrix,classification_report
print('Classification Report:\n',classification_report(y_test,rf_pred))
print('Confusion Matrix:\n',confusion_matrix(y_test,rf_pred))`

Classification Report:

	precision	recall	f1-score	support
0	0.90	0.71	0.79	496

jupyter Untitled8-Copy1 Last Checkpoint: 27 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) Logout

In [85]:

```
from sklearn.ensemble import RandomForestClassifier
rf=RandomForestClassifier()
rf.fit(x_train,y_train)
rf_pred=rf.predict(x_test)
```

In [86]:

```
from sklearn.metrics import confusion_matrix,classification_report
print('Classification Report:\n',classification_report(y_test,rf_pred))
print('Confusion Matrix:\n',confusion_matrix(y_test,rf_pred))

Classification Report:
precision    recall   f1-score   support
          0       0.90      0.71      0.79      496
          1       0.95      0.99      0.96     2543

   accuracy        0.94      3039
  macro avg       0.92      0.85      0.88      3039
weighted avg       0.94      0.94      0.94      3039

Confusion Matrix:
[[ 351 145]
 [ 37 2506]]
```

In []: