



CREDIT CARD OFFER ACCEPTANCE TRENDS IN BANKING

Group 04

S14853 - PRAMUDI RAJAMANTHRI
S15030 - VIDURA CHATHURANGA
S15091 - CHALANI WIJAMUNIGE

Abstract

This study aims to construct a comprehensive descriptive analysis and fit predictive models with high accuracies for a dataset obtained from the platform ‘Kaggle’ which focuses on identifying the reasons behind customers accepting or rejecting a credit card offers offered by banks. The dataset comprises of records of 18,000 customers with 18 variables containing both categorical and numerical. At the completion of this analysis, the insights gained about credit card offer acceptance trends and possibilities of future acceptance are summarized in this report alongside an accurate predictive model.

Content

| | |
|---|---|
| Abstract..... | 1 |
| List of figures..... | 1 |
| List of Tables..... | 1 |
| 1. Introduction..... | 1 |
| 2. Description of the Problem | 2 |
| 3. Objectives of the Study..... | 2 |
| 4. Description of the data set..... | 2 |
| 5. Data pre-processing..... | 2 |
| 6. Important Results of the Descriptive Analysis..... | 3 |
| 7. Important Results of the Advanced Analysis..... | 5 |
| 8. Issues Encountered and Proposed Solutions..... | 8 |
| 9. Discussion and Conclusion..... | 8 |
| 10. Appendix of the code..... | 9 |

List of Figures

| | |
|---|---|
| Figure 6.1 : Class distribution of the training dataset..... | 3 |
| Figure 6.2 : Multiple bar plot of Reward Type vs Offer Acceptance..... | 3 |
| Figure 6.3 : Stacked bar plot of Bank Accounts Open vs Offer Acceptance..... | 3 |
| Figure 6.4 : Stacked bar plot of Credit Cards Held vs Offer Acceptance..... | 3 |
| Figure 6.5 : Stacked bar plot of Homes Owned vs Offer Acceptance..... | 3 |
| Figure 6.6 : Multiple bar plot of Mailer Type vs Offer Acceptance..... | 4 |
| Figure 6.7 : Multiple bar plot of Credit Rating vs Offer Acceptance..... | 4 |
| Figure 6.8 : Multiple bar plot of Income Level vs Offer Acceptance..... | 4 |
| Figure 6.9 : Box plots of Quarterly Bank Balance vs Offer Acceptance..... | 4 |
| Figure 6.10 : Box plot of Average Balance vs Offer Acceptance..... | 5 |
| Figure 6.11 : Multiple Correspondence Analysis plot of categorical variables..... | 5 |
| Figure 6.12 : Partial Least Square – Discriminant Analysis score plot..... | 5 |
| Figure 6.13 : Spearman correlation matrix of the continuous variables..... | 5 |
| Figure 7.1 : Variable Importance plot of Random Forest model – For SMOTE..... | 7 |
| Figure 7.2 : Variable Importance plot of Random Forest model – For Random Under Sample Technique..... | 7 |
| Figure 7.3 : Partial Dependence plots of the final model..... | 8 |

List of Tables

| | |
|---|---|
| Table 4.1 : Variables of the Dataset..... | 6 |
| Table 7.1 : Evaluation metrics of initial models – without balancing techniques..... | 6 |
| Table 7.2 : Evaluation metrics of initial models – with SMOTE..... | 6 |
| Table 7.3 : Evaluation metrics of initial models – with Random Under Sampling | 6 |
| Table 7.4 : Evaluation metrics of hyper parameter tuned models – with SMOTE | 6 |
| Table 7.5 : Evaluation metrics of hyper parameter tuned models – with Random Under Sampling..... | 6 |
| Table 7.6 : Evaluation metrics of the Voting Classifier – with SMOTE..... | 7 |
| Table 7.7 : Evaluation metrics of initial models applied on reduced dataset – with SMOTE..... | 7 |
| Table 7.8 : Evaluation metrics of hyper parameter tuned models applied on reduced dataset – with SMOTE..... | 7 |
| Table 7.9 : Evaluation metrics of the Voting Classifier applied on reduces dataset – with SMOTE..... | 8 |

1. Introduction

In the everchanging landscape of financial transactions, credit cards emerge as one of the leading methods of payment, shaping the dynamics of modern commerce. To remain competitive, banks strategically offer a variety of credit card promotions, aiming to expand their customer base. However, the pivotal challenge faced by financial institutions lies in understanding the types of offers most likely to be embraced by customers. This data analysis project delves into the comprehensive credit card offer acceptance information of a bank, encompassing records from 18,000 customers and 18 distinct features. By unraveling the complex fabric of customer behavior and identifying the subtle reasons behind offer acceptance or rejection, this endeavor seeks to provide actionable insights into credit card offer acceptance trends. The ultimate objective is to leverage machine learning techniques to predict the likelihood of a customer accepting a credit card offer, offering the financial industry a predictive tool for refining their promotional strategies.

2. Description of the problem

The problem at hand revolves around the challenge faced by banks in identifying credit card offer acceptance trends. With credit cards being a predominant payment method, banks aim to expand their customer base by strategically offering various promotions. The goal of this data analysis project is to uncover the factors influencing customer decisions and employ machine learning techniques to predict whether a customer will accept a credit card offer, providing valuable insights for refining promotional strategies.

3. Objectives of the study

1. Investigate and understand the diverse factors influencing customer behavior in credit card offer acceptance within the dataset of 18,000 customers and 18 distinct features.
2. Uncover actionable insights that reveal the subtle reasons behind customers accepting or rejecting credit card offers, contributing to a comprehensive understanding of offer acceptance trends.
3. Develop and implement machine learning techniques to predict the likelihood of a customer accepting a credit card offer, providing financial institutions with a valuable predictive tool to enhance and refine their promotional strategies.

4. Description of the data set

“Credit Card Offer Acceptance Trends in Banking” dataset contains credit card offer acceptance information of a bank. There are records of 18,000 customers and 18 features (2 unique valued, 8-Qualitative and 8-Quantitative) with the target variable indicating whether an offer is accepted or not.

Source: [Credit Card Offer Acceptance Trends in Banking \(kaggle.com\)](https://www.kaggle.com/datasets/ashishpatel26/credit-card-offer-acceptance-trends-in-banking)

| No | Variable | Description | Type |
|----|----------------------|--|--------------|
| 01 | Index | Index | Unique |
| 02 | Customer Number | Unique identifier for each customer. | Unique |
| 03 | Offer Accepted | Whether the customer accepted the offer. (Yes, No) | Response |
| 04 | Reward | Type of reward offered. (Air Miles, Cash Back, Points) | Qualitative |
| 05 | Mailer Type | Type of mailer used to send the offer. (Letter, Postcard) | Qualitative |
| 06 | Income Level | The customer’s income level. (High, Medium, Low) | Qualitative |
| 07 | Bank Accounts Open | Number of bank accounts the customer has opened. | Quantitative |
| 08 | Overdraft Protection | Whether the customer has overdraft protection on their account. (Yes, No) | Qualitative |
| 09 | Credit Rating | How well the customer’s payment record reflects their ability and willingness to repay the debt under terms accepted by creditors. (High, Medium, Low) | Qualitative |
| 10 | Credit Cards Held | Number of credit cards the customer hold. | Quantitative |
| 11 | Homes Owned | Number of homes the customer has. | Quantitative |
| 12 | Household Size | Size of household that the customer belongs is kept track of. (1,2,3,4,5,6,7,8,9) | Qualitative |
| 13 | Own Your Home | Whether the customer owns their home. (Yes, No) | Qualitative |
| 14 | Average Balance | Average balance across all accounts. | Quantitative |
| 15 | Q1 Balance | The customer’s balance in 1 st quarter. | Quantitative |
| 16 | Q2 Balance | The customer’s balance in 2 nd quarter. | Quantitative |
| 17 | Q3 Balance | The customer’s balance in 3 rd quarter. | Quantitative |
| 18 | Q4 Balance | The customer’s balance in 4 th quarter. | Quantitative |

Table 4.1

5. Data preprocessing

- Index Variables was dropped from the data set as it didn’t provide any meaningful contribution to the study.
- There were no duplicate records were identified in the dataset.
- 24 missing records were detected in each of the Q1 Balance, Q2 Balance, Q3 Balance, and Q4 Balance variables.
- Those missing values were imputed using the mean values of each variable from the training set.
- Initially, the Household Size variable had nine distinct values.
- The Household Size variable was recoded: Categories 1, 2, 3, 4 were combined into a new category named "Small," and categories 5, 6, 7, 8, 9 were combined into a new category named "Big."
- Finally, the dataset was randomly split into training and test datasets which contains 80% and 20% of the entries from original dataset.

6. Important results of the Descriptive Analysis

The pie chart (Figure 6.1) illustrates a notable imbalance in class distribution within the training set, with the representation of the "Yes" category significantly smaller compared to the prevalence of the "No" category. The dataset exhibits a pronounced imbalance, prompting the need for careful consideration and appropriate handling of this skewed distribution during further analysis.

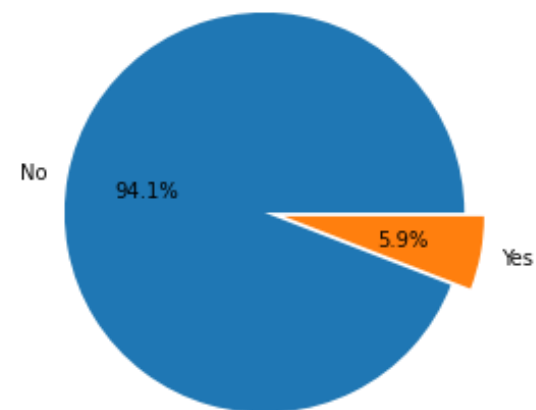


Figure 6.1

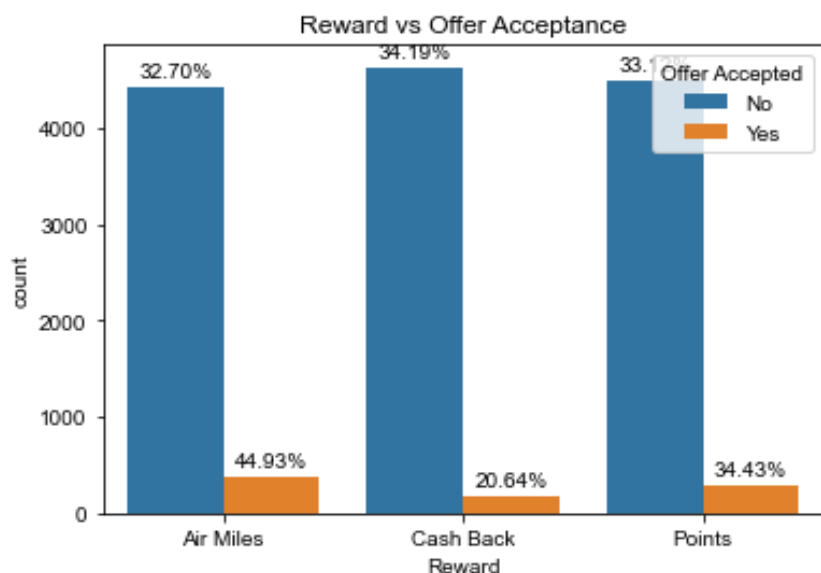


Figure 6.2

The Figure 6.2 plot is used to analyze the relationship between the type of reward received and the response variable - Offer Acceptance. Upon observation, it is apparent that the distribution of not accepting the offer is relatively consistent across all categories of the 'Reward' variable. However, in the case of the 'Accepting' category, there is a slight variation among the three reward categories. This suggests that rewards may play a significant role in influencing credit card offer acceptance.

These findings align with insights from the American Bankers Association, as mentioned in their articles. According to them, the analysis, coupled with other studies, indicates that card rewards are accessible, valuable, and well-understood by consumers across all income levels. Credit card offers, as revealed by this analysis, indeed exert a noteworthy influence on the acceptance or rejection of credit card offers. (<https://www.aba.com/news-research/analysis-guides/the-benefits-of-credit-card-rewards>)

This visualization (Figure 6.3) illustrates the relationship between the response variable 'Offer Accepted' and the 'Number of Bank Accounts Opened' variable. However, the overall distribution of acceptance and rejection of credit card offers appears similar across the three categories of the 'Number of Bank Accounts Opened' variable. This suggests that the number of bank accounts alone may not be a decisive factor in determining credit card offer acceptance, as the distribution remains consistent across different account counts.



Figure 6.3

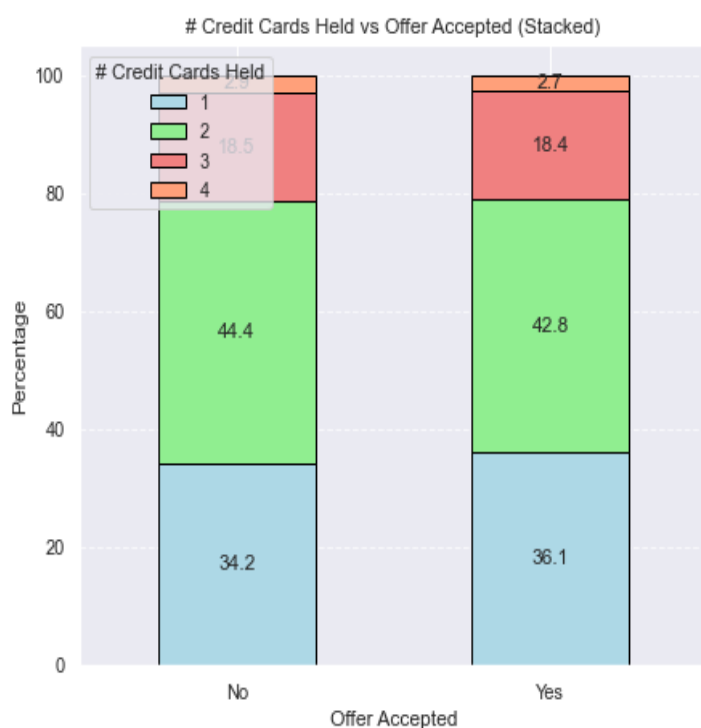


Figure 6.4

The presented visualization (Figure 6.4) showcases the correlation between the response variable 'Offer Accepted' and the 'Number of Credit cards held' variable. Notably, the distribution patterns for both acceptance and rejection of credit card offers seem consistent across the four categories of the 'Number of Credit cards held' variable. This implies that the singular consideration of the number of credit cards held may not be a determining factor in credit card offer acceptance, as the distribution remains uniform across varying counts of credit cards held.

The depicted plot (Figure 6.5) visualizes the relationship between the variables "No of Houses Owned" and "Offer Accepted." In this portrayal, the distribution among the categories of the "No of Homes Owned" variable appears relatively consistent for both categories of the "Offer Accepted" variable. This observation suggests that the number of homes owned may not exert a significant influence on the acceptance of credit cards, as the distribution remains comparable across different counts of homes owned.



Figure 6.5

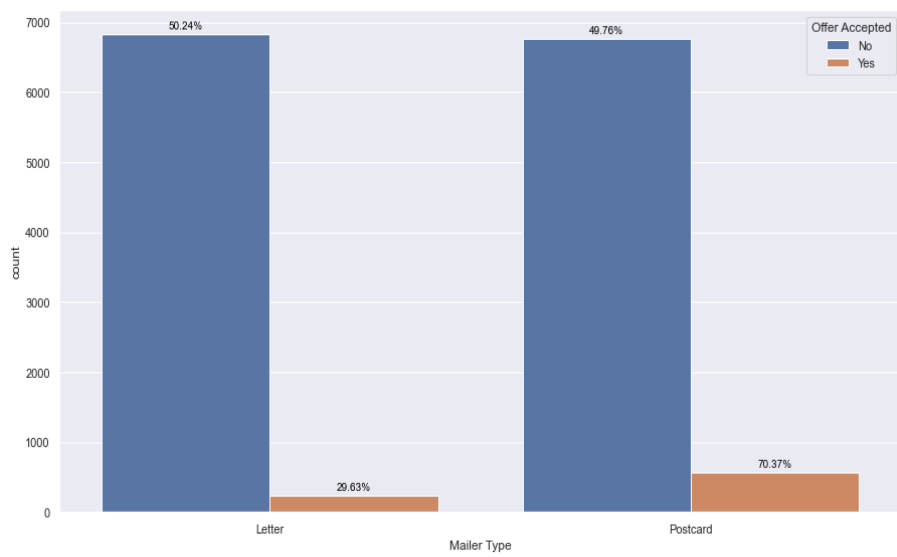


Figure 6.6

"In our examination of the Credit Rating variable's impact on credit card acceptance, a notable and substantial behavioral difference becomes evident between accepted and rejected incidents. This pronounced contrast strongly indicates a significant association between these variables, emphasizing the pivotal role of Credit Rating in determining credit card acceptance. The discernible impact implies that a deeper analysis and careful consideration of Credit Rating are essential for gaining comprehensive insights into the dynamics of customer responses to credit card offers.

The importance of credit scores is reiterated in an article by Investopedia

(https://www.investopedia.com/articles/pf/07/credit_card_rating.asp#:~:text=Your%20credit%20score%20is%20an,you%20are%20able%20to%20borrow.), where it is highlighted that your credit score plays a crucial role in your ability to obtain a credit card or borrow money. The ideas presented in this reference align with and validate the results obtained in our analysis."

Regarding the Mailer type (Figure 6.6), while the distribution among the "Not accepted" category remains similar, there is a noticeable distinction between the "Accepted" category concerning the two mailer types, "Letter" and "Postcard." This discrepancy suggests a meaningful impact of the Mailer type on the acceptance of credit cards by customers. The observed significant difference in distribution implies that the method of communication, whether through letters or postcards, plays a substantial role in influencing customer decisions regarding credit card offers.

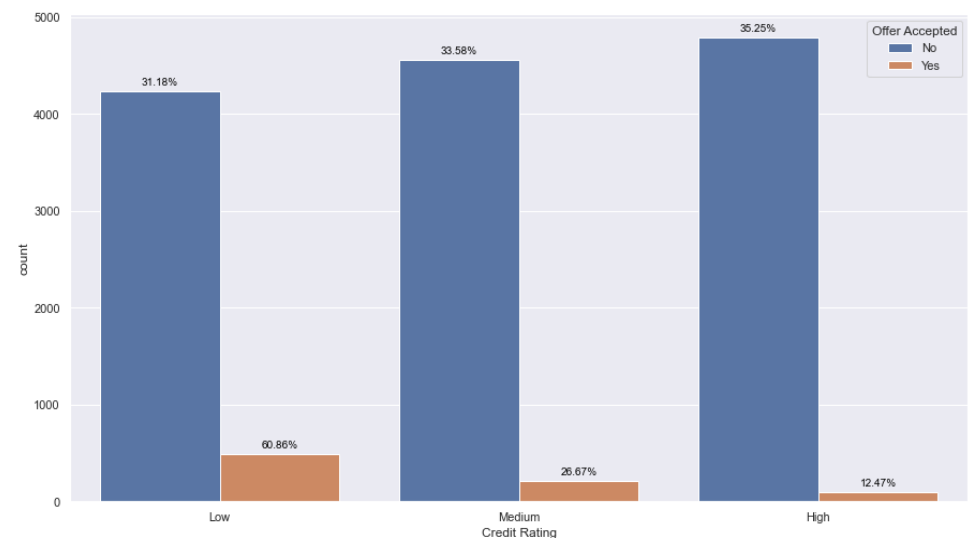


Figure 6.7

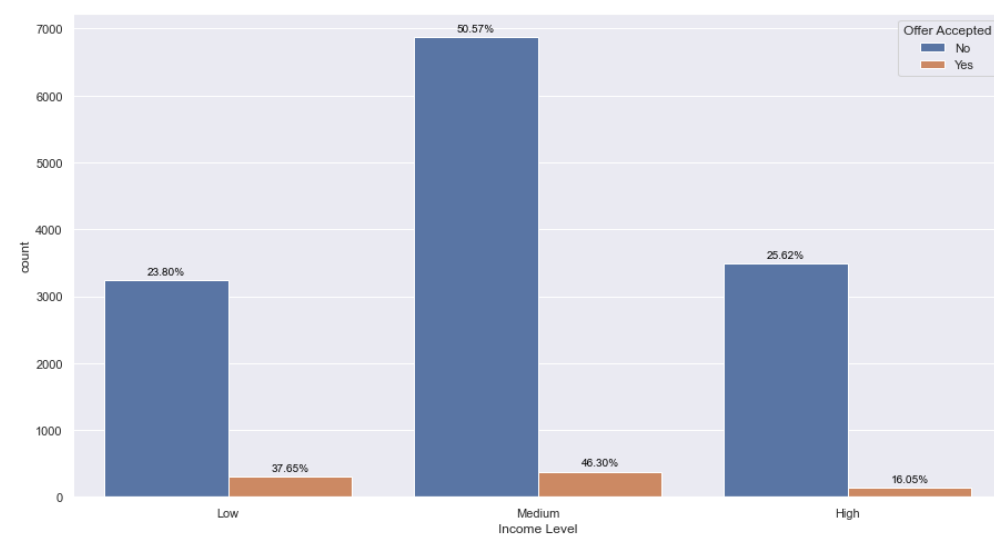


Figure 6.8

highlighting the intrinsic connection between income and credit card approval decisions. They underscore the imperative nature of factoring in diverse income levels in the assessment of credit card acceptance trends.

This analysis underscores the pivotal role of income level as a decisive factor influencing credit card acceptance. The distributions and proportions of acceptance and rejection manifest notable variations across the distinct categories of the 'Income Level' variable as shown in figure 6.8.

This observation aligns cohesively with insights discussed in articles from CNBC (<https://www.cnbc.com/select/how-does-salary-and-income-impact-your-credit-score/>) and IndusInd Bank (<https://www.indusind.com/iblogs/categories/manage-your-finance/the-crucial-role-of-income-in-credit-card-approval/>)

In the latter, it is emphasized that credit card approval transcends merely having a good credit score, with income being a paramount factor considered by banks. These references fortify our findings,

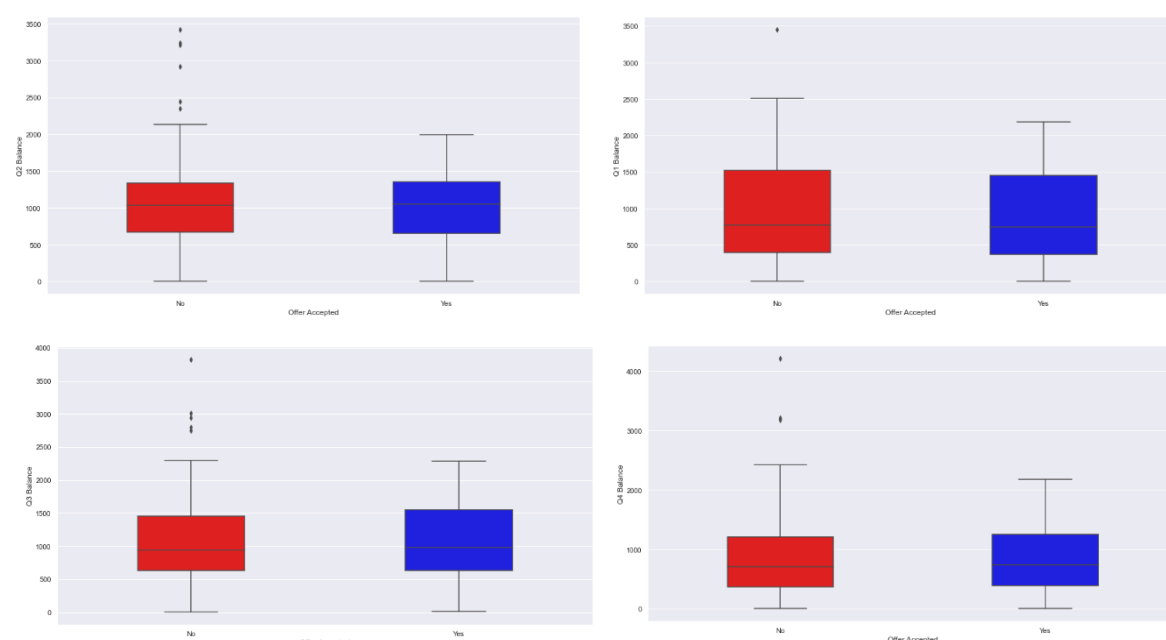


Figure 6.9

Next, we need to assess the relationship between quantitative variables and the response variable, Offer Acceptance. Upon analyzing the impact of customer balances in each quarter of the year, we have uncovered significant insights into their behavior (Figure 6.9).

In each quarter, the distribution of balances for the "Offer Not Accepted" category appears considerably wider compared to the "Accepted" category. Moreover, numerous outliers are observed specifically within the "Not Accepted" category. The distributions of balances in each quarter exhibit noticeable differences between the two acceptance categories. These findings suggest that there is a substantial impact of balances in Q1, Q2, Q3, and Q4 on the acceptance of credit cards by customers.

Turning our attention to the numerical variable of Average Balance, insights derived from the boxplots unveil notable patterns. The distribution of Average Balance associated with the "Offer Not Accepted" category appears broader compared to the distribution for the "Offer Accepted" category. Moreover, a higher prevalence of outliers is evident in the "Not Accepted" category. These distinct distribution patterns strongly indicate a substantial impact of Average Balance on offer acceptance. The observed differences in distribution strongly suggest that Average Balance plays a pivotal role in determining the likelihood of offer acceptance. Coincidentally, The Nest discusses a related notion in their report titled "Do Credit Cards Check Your Bank Balance Before Issuing You a Card?" (source: <https://budgeting.thenest.com/credit-cards-check-bank-balance-before-issuing-card-24697.html>), aligning closely with our findings and emphasizing the significance of Average Balance in credit card acceptance dynamics.

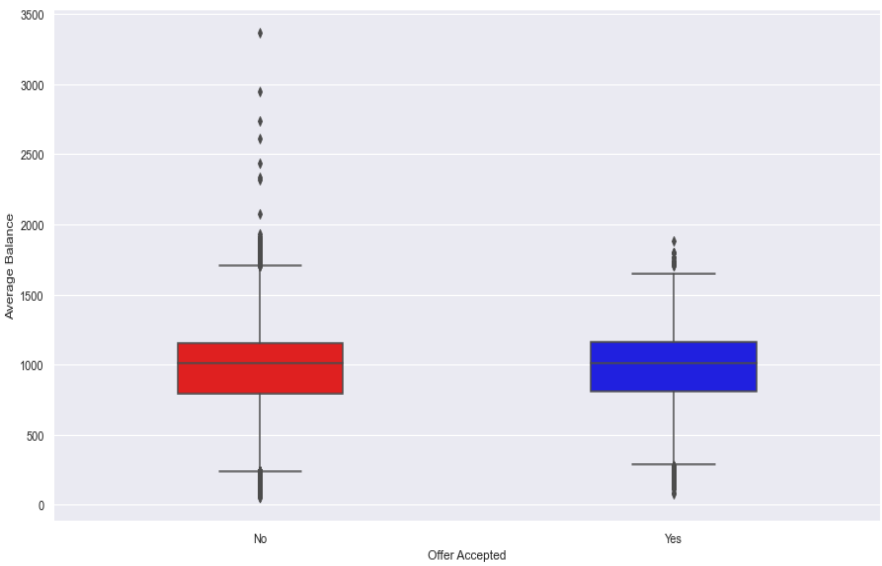


Figure 6.10

Multiple Correspondence Analysis for Categorical Variables

In the Multiple Correspondence Analysis conducted for the categorical variables, it was observed that the first 2 components explain 20.41% of the entire variance. Considering the association of each variable with the response variables, the information uncovered on the surface analysis are further clarified as the variables Reward, Mailer Type, Income Level, Credit Rating shows a high association with the response variable. But those are scattered around the category ‘Offer_Accepted_No’, which suggests that those are more related towards that category while ‘Offer_Accepted_Yes’ category shows distant relationships. The reason behind this can be the imbalanced nature of the dataset, as when there are plenty of datapoints, that category captures most of the variance.

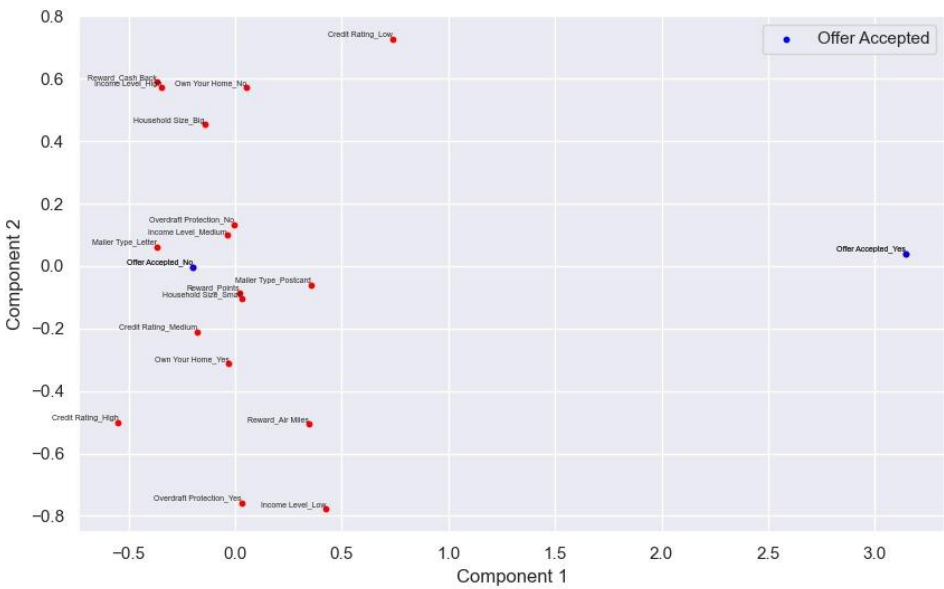


Figure 6.11

Partial Least Square – Discriminant Analysis

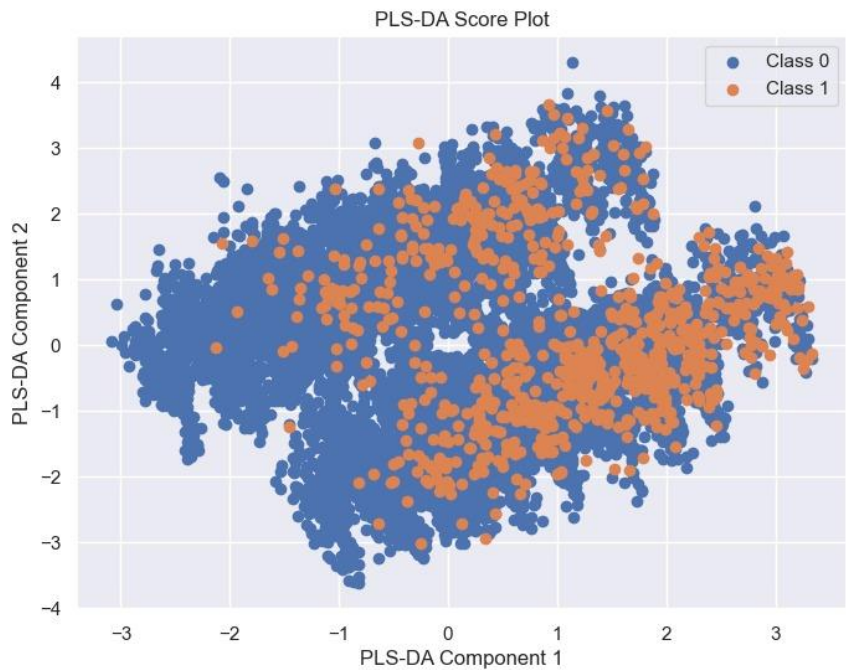


Figure 6.12

The variance explained by the first 2 components of PLS – DA is approximately added up to 18.13% which is considerably low. Therefore, the findings of this analysis cannot be clearly interpreted and according to the score plot where there are no distinct clusters to be observed. But a certain pattern is quite visible in the scattering of scores in the plot which can be addressed as a highlighting fact.

Spearman Correlation of Continuous Variables.

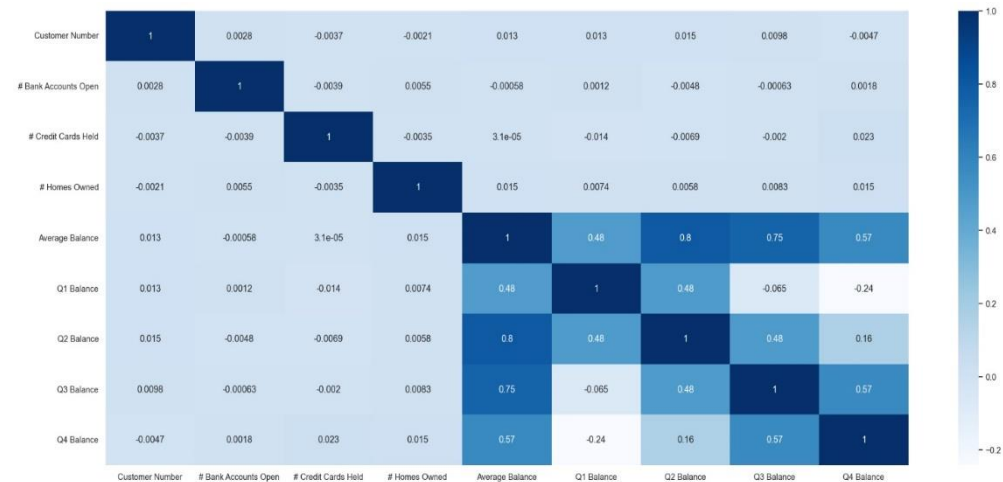


Figure 6.13

Considering the correlations between continuous variables of the dataset, only the Average Balance and Quarterly Balances shows some correlation and the highest correlation is shown between Average Balance and Q2 Balance which is 0.8. Therefore, a significant multicollinearity is not observed. Due to this, its possible to start model fitting, keeping logistic regression as a baseline.

7. Important results of the Advanced Analysis

In the advanced analysis phase, the main objective is to delve deep into the classifying the acceptance of credit card offers based on customer details. It is aimed to evaluate performance of a range of predictive models in their capacity to accurately classify instances based on the ‘Offer Accepted’ variable. Since, no significant multicollinearity was observed in the descriptive study, models will be chosen keeping logistic regression as a baseline. To initiate the model fitting, Logistic Regression, Logistic Ridge, Logistic Lasso, K-Nearest Neighbor, Gaussian Naïve Bays, Support Vector Machine, Random Forest, and XGBoost models were fitted and evaluated.

| | Without Applying Balancing Techniques | | | | | | | | | | | | | | | |
|-----------|---------------------------------------|----------------|----------------|------|----------------------|------|---------------|---------|---------------------|----------------|----------------|------|----------------------|------|---------------|---------|
| | Training Set | | | | | | | | Testing Set | | | | | | | |
| | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost |
| Accuracy | 0.95 | 0.95 | 0.95 | 0.94 | 0.95 | 0.95 | 1.0 | 0.97 | 0.94 | 0.94 | 0.94 | 0.95 | 0.94 | 0.94 | 0.95 | 0.95 |
| Precision | 0.48 | 0.48 | 0.48 | 0.54 | 0.48 | 0.48 | 1.0 | 0.98 | 0.47 | 0.47 | 0.47 | 0.75 | 0.47 | 0.47 | 0.48 | 0.58 |
| Recall | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 1.0 | 0.72 | 0.5 | 0.5 | 0.5 | 0.52 | 0.5 | 0.5 | 0.5 | 0.5 |
| F1 score | 0.49 | 0.49 | 0.49 | 0.5 | 0.49 | 0.49 | 1.0 | 0.79 | 0.48 | 0.48 | 0.48 | 0.52 | 0.48 | 0.48 | 0.49 | 0.5 |

Table 7.1

According to the results given in the table, almost every model has not been successful in predicting the target class which is the acceptance of the offer as the evaluation metrics precision, recall and F1 score show values ranging from 45% to 60%. The reason behind these low metric values can be again the extreme imbalanced nature of the dataset as the number of offer accepted cases is significantly low compared to offer rejected cases.

To overcome this, 2 balancing techniques which are Synthetic Minority Oversampling Technique (SMOTE) and Random Under Sampling Technique, were introduced in order to improve precision and recall which represents the accuracy in identifying the positive cases (accepted cases) as that’s where the main focus of this study is relied on.

| | With SMOTE | | | | | | | | | | | | | | | |
|-----------|---------------------|----------------|----------------|------|----------------------|------|---------------|---------|---------------------|----------------|----------------|------|----------------------|------|---------------|---------|
| | Training Set | | | | | | | | Testing Set | | | | | | | |
| | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost |
| Accuracy | 0.83 | 0.83 | 0.83 | 0.99 | 0.83 | 0.81 | 1.0 | 0.93 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.88 | 0.83 |
| Precision | 0.55 | 0.55 | 0.55 | 0.99 | 0.55 | 0.55 | 1.0 | 0.71 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 | 0.8 | 0.65 |
| Recall | 0.66 | 0.66 | 0.66 | 0.89 | 0.61 | 0.66 | 1.0 | 0.82 | 0.63 | 0.63 | 0.63 | 0.66 | 0.66 | 0.64 | 0.81 | 0.66 |
| F1 score | 0.56 | 0.56 | 0.56 | 0.94 | 0.55 | 0.56 | 1.0 | 0.75 | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 | 0.8 | 0.66 |

Table 7.2

| | With Random Under Sampling | | | | | | | | | | | | | | | |
|-----------|----------------------------|----------------|----------------|------|----------------------|------|---------------|---------|---------------------|----------------|----------------|------|----------------------|------|---------------|---------|
| | Training Set | | | | | | | | Testing Set | | | | | | | |
| | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost |
| Accuracy | 0.83 | 0.83 | 0.83 | 0.78 | 0.81 | 0.83 | 0.97 | 0.89 | 0.87 | 0.87 | 0.83 | 0.77 | 0.82 | 0.83 | 0.9 | 0.76 |
| Precision | 0.56 | 0.56 | 0.55 | 0.54 | 0.55 | 0.56 | 0.83 | 0.66 | 0.56 | 0.56 | 0.55 | 0.52 | 0.55 | 0.55 | 0.56 | 0.63 |
| Recall | 0.65 | 0.65 | 0.65 | 0.89 | 0.63 | 0.65 | 0.89 | 0.86 | 0.63 | 0.63 | 0.66 | 0.59 | 0.64 | 0.66 | 0.58 | 0.62 |
| F1 score | 0.57 | 0.57 | 0.56 | 0.94 | 0.55 | 0.57 | 0.86 | 0.71 | 0.57 | 0.57 | 0.56 | 0.51 | 0.55 | 0.56 | 0.56 | 0.61 |

Table 7.3

After balancing the dataset using both the balancing techniques, the performances of all models have improved significantly. By carefully examining, we identified that Random Forest and XGBoost in both scenarios perform better in term of all the evaluation metrics and also, they show considerably less overfitting. Therefore, to further improve their performance, hyper parameter tuning was applied.

| | With SMOTE | | | |
|-----------|---------------|----------|--------------|----------|
| | Random Forest | | XGBoost | |
| | Training Set | Test Set | Training Set | Test Set |
| Accuracy | 1.0 | 0.89 | 0.97 | 0.87 |
| Precision | 1.0 | 0.85 | 0.94 | 0.74 |
| Recall | 1.0 | 0.81 | 0.95 | 0.79 |
| F1 score | 1.0 | 0.83 | 0.88 | 0.75 |

Table 7.4

Best Parameters (Random Forest) : {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100}

Best Parameters (XGBoost) : {'learning_rate': 0.2, 'max_depth': 7, 'min_child_weight': 1, 'n_estimators': 200, 'subsample': 0.8}

| | With Random Under Sampling | | | |
|-----------|----------------------------|----------|--------------|----------|
| | Random Forest | | XGBoost | |
| | Training Set | Test Set | Training Set | Test Set |
| Accuracy | 0.77 | 0.73 | 0.81 | 0.75 |
| Precision | 0.77 | 0.64 | 0.75 | 0.54 |
| Recall | 0.81 | 0.5 | 0.81 | 0.61 |
| F1 score | 0.72 | 0.57 | 0.82 | 0.55 |

Table 7.5

Best Parameters (Random Forest) : {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 10, 'n_estimators': 200}
Best Parameters (XGBoost) : {'learning_rate': 0.01, 'max_depth': 5, 'min_child_weight': 3, 'n_estimators': 100, 'subsample': 0.8}

As observed in these 2 tables, after applying the tuned parameters, the improvement in precision and recall is quite evident, By comparing the model performance under 2 balancing techniques, it is visible that the models fitted to the dataset balanced using SMOTE performs better than the models fitted to the under sampled dataset. It also should be noted that there is a possible information loss when using the random under sampling technique compared to the oversampling technique. Therefore, the study will be carried further on the dataset balanced with SMOTE.

Considering the performance of Random Forest and XGBoost with SMOTE, as both models show better performance, both were combined to form an ensemble model using the voting classifier using soft voting.

| | Voting Classifier | |
|-----------|-------------------|-------------|
| | Training Set | Testing Set |
| Accuracy | 0.99 | 0.92 |
| F1-Score | 0.96 | 0.83 |
| Precision | 0.99 | 0.88 |
| Recall | 0.97 | 0.82 |

Table 7.6

As observed in the table, the voting classifier shows a significant improvement in performance as the training accuracies are almost 100% in terms of all evaluation metrics.

Now to further reduce the complexity of the models as there are 17 predictor variables, the dimensionality of the dataset was reduced. According to the variable importance plot obtained from Random Forest, all the

initial models were refitted considering the most important 15, 13 and 10 variables sequentially.

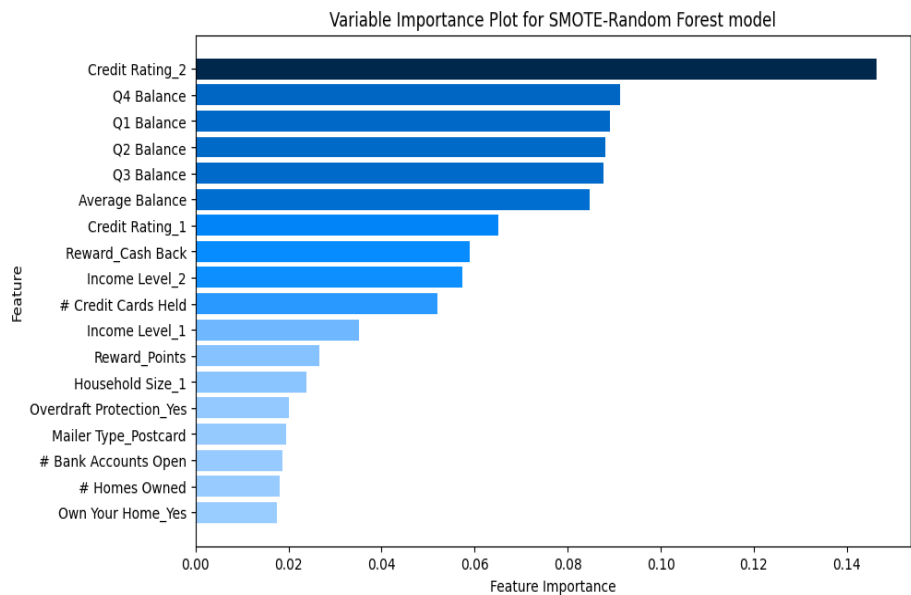


Figure 7.1

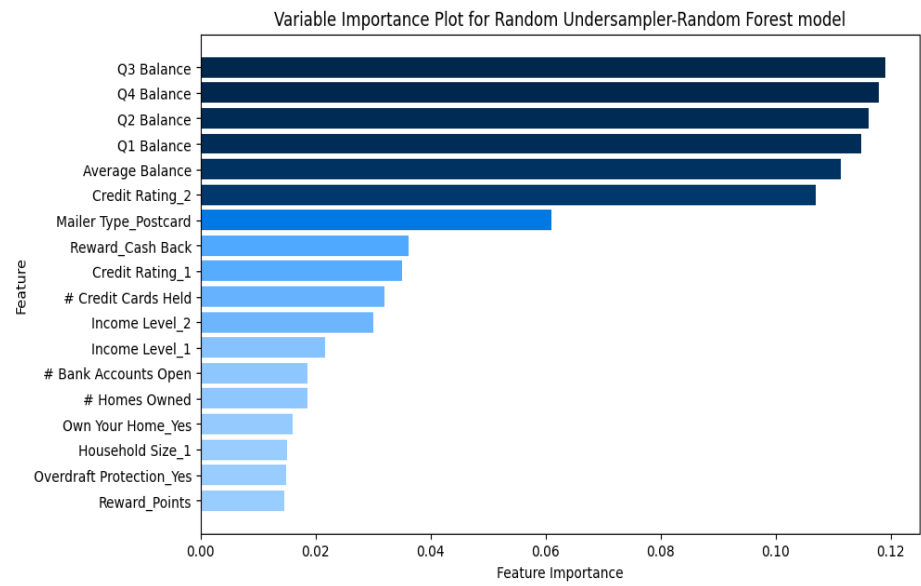


Figure 7.2

According to the performances of the reduced models, the model with top 10 most important features which are 'Q1 Balance', 'Q2 Balance', 'Q3 Balance', 'Q4 Balance', 'Average Balance', 'Credit Rating_2', 'Mailer Type_Postcard', 'Reward_Cash Back', 'Credit Rating_1', '# Credit Cards Held', gave better and justifiable results in term of all the evaluation metrics as given in the following table. Here the models are fitted on the dataset balanced using SMOTE.

| | Training Set | | | | | | | | Testing Set | | | | | | | |
|-----------|---------------------|----------------|----------------|------|----------------------|------|---------------|---------|---------------------|----------------|----------------|------|----------------------|------|---------------|---------|
| | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost | Logistic Regression | Logistic Ridge | Logistic Lasso | KNN | Gaussian Naïve Bayes | SVM | Random Forest | XGBoost |
| Accuracy | 0.95 | 0.95 | 0.95 | 0.99 | 0.95 | 0.95 | 1.0 | 0.93 | 0.86 | 0.86 | 0.86 | 0.83 | 0.86 | 0.83 | 0.87 | 0.84 |
| Precision | 0.57 | 0.57 | 0.55 | 1.0 | 0.75 | 0.55 | 1.0 | 0.73 | 0.56 | 0.56 | 0.56 | 0.83 | 0.56 | 0.48 | 0.83 | 0.72 |
| Recall | 0.61 | 0.61 | 0.6 | 0.83 | 0.71 | 0.64 | 1.0 | 0.77 | 0.51 | 0.51 | 0.51 | 0.82 | 0.61 | 0.4 | 0.81 | 0.66 |
| F1 score | 0.57 | 0.57 | 0.57 | 0.86 | 0.72 | 0.55 | 1.0 | 0.75 | 0.5 | 0.5 | 0.5 | 0.82 | 0.57 | 0.49 | 0.82 | 0.72 |

Table 7.7

According to the accuracies shown in the table, it is identified that models K-Nearest Neighbor, Random Forest and XGBoost show the best performance out of all the fitted models under the reduced variable scenario. Further to improve the performance and to ensure minimum overfitting of Random Forest and XGBoost, those 2 models were hyper parameter tuned.

| | Random Forest | | XGBoost | |
|-----------|---------------|----------|--------------|----------|
| | Training Set | Test Set | Training Set | Test Set |
| Accuracy | 1.0 | 0.85 | 0.93 | 0.86 |
| Precision | 1.0 | 0.83 | 0.84 | 0.75 |
| Recall | 1.0 | 0.87 | 0.94 | 0.71 |
| F1 score | 1.0 | 0.83 | 0.79 | 0.75 |

Table 7.8

Since all the 3 models, which are the Reduced KNN model, Reduced and hyper parameter tuned Random Forest and XGBoost models still shows high accuracies, to collectively capture their intelligence, an ensemble model was fitted using the voting classifier with 'soft' voting. The evaluation metrics of the voting classifier is as follows.

| | Voting Classifier | |
|-----------|-------------------|-------------|
| | Training Set | Testing Set |
| Accuracy | 0.97 | 0.92 |
| F1-Score | 0.92 | 0.82 |
| Precision | 0.85 | 0.83 |
| Recall | 0.97 | 0.78 |

It is evident and clearly observable that the ensemble model capturing the KNN, hyper parameter tuned Random Forest and XGBoost models, which is the voting classifier build upon the dataset balanced using SMOTE with reduced dimensionality outperforms all the models which were fitted so far in terms of simplicity, high accuracies, and less overfitting.

Table 7.9

Partial Dependence plots

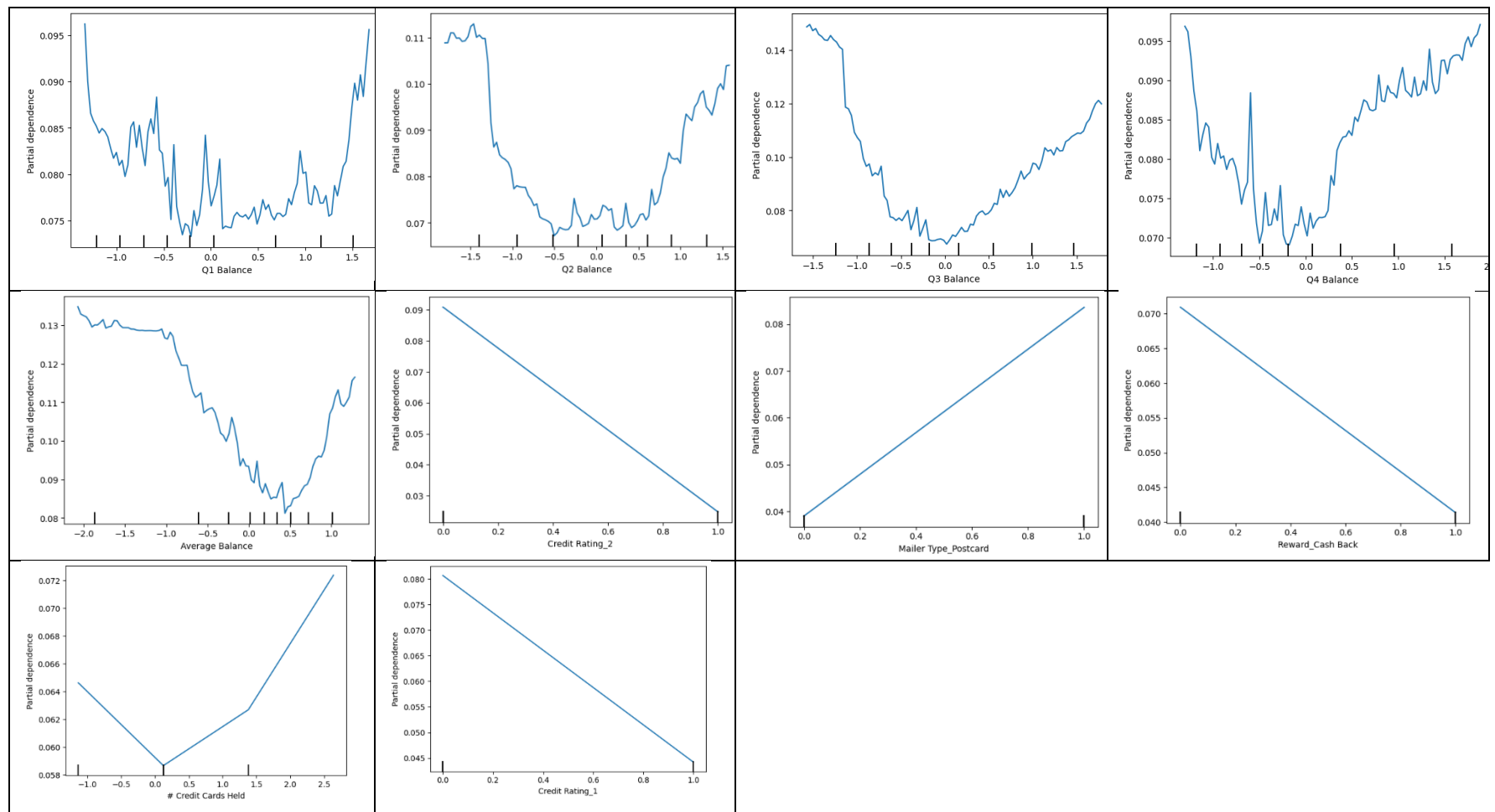


Figure 7.3

Observing the above partial dependence plots considering the selected variables utilized when fitting the ensemble model based on soft voting which we selected as the best model in this prediction scenario depict that variables namely 'Q1 Balance', 'Q2 Balance', 'Q3 Balance', 'Q4 Balance', 'Average Balance' effect on the acceptance of the credit card offers at higher values while lower values causes non-acceptance of the offers. Moreover, the acceptance of the credit card offer is highly probable when the number of credit cards are at a high value which is quite justifiable as customers who maintain several credit cards due to their high financial potential.

8. Issues Encountered and Proposed Solutions

- The 'Household Size' categorical variable had 9 levels in it which is a considerably large number of categories to have and may further increase dimensionality in model fitting. To reduce the number of levels, the 9 levels were divided in to 2 categories capturing level 1,2,3 and 4 to one category and the rest of the levels to another category labeled as 'small' and 'big' respectively.
- The data distribution over the 2 categories of the response variable was highly imbalanced and this caused the majority class to highly influence the evaluation metrics of the fitted models where the precision and F1 score turned out to be significantly low when compared to accuracy due to the lower number of accepted offer (true positive) instances. To overcome this, the dataset was balance using Synthetic Minority Over Sampling Technique (SMOTE) and Random under sampling technique.
- In the initially fitted models to the original training dataset, none of the models predicted 'Offer_Accepted_Yes', in other words the true positive cases accurately. This disrupts the entire study as the main focus of this study is to predict offers that are getting accepted. Dataset balancing techniques were applied to balance the dataset to overcome this issue as well.
- The initial dataset had 17 predictor variables and after applying one-hot encoding on categorical variables, the dimensionality further increased making the models more complex, while inhibiting its performance. To reduce dimensionality, a feature importance plot was obtained and the most important 10 features were extracted as it gave the optimum performance and models were refitted.
- Some of the fitted models showed high performance in terms of all evaluation metrics, so selecting a single model which outperforms rest of the models was difficult, therefore to capture the performance of all the high performing models, an ensemble model was fitted using a voting classifier and that outperformed rest of the models.

9. Discussion and Conclusion

- At the surface analysis it was discovered that features such as Quarterly Balances, Average Balance, Credit Rating, Reward, Mailer Type, Credit Cards Held and Income Level are the most significant features which determine the acceptance of a credit card offer and that was further proved in variable importance plot of Random Forest model.
- Due to the imbalanced nature of the dataset, the models fitted on the dataset which was balanced using the oversampling technique SMOTE and Random Under Sampling Technique performed better.

- Due to the possible information loss in Random Under Sampling Technique, the accuracies of models fitted on data balanced using SMOTE showed higher values and better performance.
- Out of all the initial models, Random Forest and XGBoost models performed the best and their performance was further improved while reducing overfitting by applying hyper parameter tuning.
- The best performing models in the initial stage was combined together in an ensemble model using the Voting Classifier and that showed the best accuracies so far.
- The performance of models was further improved by reducing the dimensionality where the dataset with 10 features with the highest feature importance were extracted and models were refitted.
- Out of all models fitted after dimensionality reduction, KNN, hyper parameter tuned Random Forest and XGBoost models showed the highest accuracies and they collectively formed an ensemble model which is the Voting Classifier and that outperformed all the predictive models fitted so far showing the best accuracies in term of precision, F1 score and recall.
- Therefore, Voting Classifier fitted on the reduced dataset with balanced data (applying SMOTE) was identified as the best model with the most accurate predictions and less information loss.

10. Appendix of the Code

| | |
|---|--|
| <pre> Final Project - Credit Card Offer Acceptance in Banking In []: import os import numpy as np import pandas as pd import seaborn as sns import matplotlib.pyplot as plt import warnings warnings.filterwarnings("ignore") In []: df = pd.read_csv('creditcardmarketing-bm.csv') In []: df.head() In []: df.columns In []: df.shape In []: duplicate_rows = df[df.duplicated()] len(duplicate_rows) In []: df['Offer Accepted'].value_counts() In []: df.info() In []: df.isna().sum() In []: for column in df.columns: unique_values = df[column].unique() print("Unique values in column '{column}':") print(unique_values) print() In []: df.dtypes In []: for column in df.columns: value_counts = df[column].value_counts() print("Value counts in column '{column}':") print(value_counts) print() In []: def categorize_size(size): if size in [1, 2, 3, 4]: return "Small" else: return "Big" df['Household Size'] = df['Household Size'].apply(categorize_size) In []: df = df.drop('index',axis=1) In []: df['Offer Accepted'] = pd.Categorical(df['Offer Accepted'], categories=['No', 'Yes']) df['Reward'] = pd.Categorical(df['Reward'], categories=['Air Miles', 'Cash Back', 'Points']) df['Mailer Type'] = pd.Categorical(df['Mailer Type'], categories=['Letter', 'Postcard']) df['Overdraft Protection'] = pd.Categorical(df['Overdraft Protection'], categories=['No', 'Yes']) df['Credit Rating'] = pd.Categorical(df['Credit Rating'], categories=['Low', 'Medium', 'High'],ordered=True) df['Household Size'] = pd.Categorical(df['Household Size'], categories=['Small', 'Big'],ordered=True) df['Own Your Home'] = pd.Categorical(df['Own Your Home'], categories=['No', 'Yes']) df['Income Level'] = pd.Categorical(df['Income Level'], categories=['Low', 'Medium', 'High'],ordered=True) In []: from sklearn.model_selection import train_test_split trainset, testset = train_test_split(df, test_size=0.2, random_state=100) print("Training set shape:", trainset.shape) print("Test set shape:", testset.shape) In []: trainset['Offer Accepted'].value_counts() In []: testset['Offer Accepted'].value_counts() In []: trainset.dtypes In []: for column in trainset.columns: value_counts = trainset[column].value_counts() print("Value counts in column '{column}':") print(value_counts) print() In []: for column in testset.columns: value_counts = testset[column].value_counts() print("Value counts in column '{column}':") print(value_counts) print() In []: trainset.isna().sum() In []: testset.isna().sum() </pre> | <pre> In []: #Handling missing values Average_Balance_mean=trainset['Average Balance'].mean() Q1_Balance_mean=trainset['Q1 Balance'].mean() Q2_Balance_mean=trainset['Q2 Balance'].mean() Q3_Balance_mean=trainset['Q3 Balance'].mean() Q4_Balance_mean=trainset['Q4 Balance'].mean() trainset['Average Balance'].fillna(Average_Balance_mean, inplace=True) trainset['Q1 Balance'].fillna(Q1_Balance_mean, inplace=True) trainset['Q2 Balance'].fillna(Q2_Balance_mean, inplace=True) trainset['Q3 Balance'].fillna(Q3_Balance_mean, inplace=True) trainset['Q4 Balance'].fillna(Q4_Balance_mean, inplace=True) testset['Average Balance'].fillna(Average_Balance_mean, inplace=True) testset['Q1 Balance'].fillna(Q1_Balance_mean, inplace=True) testset['Q2 Balance'].fillna(Q2_Balance_mean, inplace=True) testset['Q3 Balance'].fillna(Q3_Balance_mean, inplace=True) testset['Q4 Balance'].fillna(Q4_Balance_mean, inplace=True) In []: trainset.isna().sum() Descriptive Analysis using the training set In []: target=trainset['Offer Accepted'].value_counts() fig, ax1 = plt.subplots() ax1.plot(target, labels = target.index, autopct = '%1.1f%%', shadow = False, explode = [0.1, 0]) ax1.axis('equal') plt.show() print('Total number of customers in the training set:', trainset['Offer Accepted'].count()) print(trainset['Offer Accepted'].value_counts()) In []: trainset.columns In []: ax = sns.countplot(x='Reward', hue='Offer Accepted', data=trainset) total_counts = trainset['Reward'].value_counts() for container in ax.containers: for bar in container: for bar in container: percentage = f'({bar.get_height() / total * 100:.2f})%' ax.annotate(percentage, xybar.get_x() + bar.get_width() / 2, bar.get_height(), xytext=(0, 3), # 3 points vertical offset textcoords='offset points', ha='center', va='bottom', fontsize=8, color='black') sns.set(rc={'figure.figsize': (15, 8)}) plt.title('Reward vs Offer Acceptance') plt.show() In []: categorical_cols = ['Bank Accounts Open', 'Credit Cards Held', 'Home Owned', 'Reward', 'Mailer Type', 'Income Level', 'Overdraft Protection', 'Credit Rating', 'Household Size', 'Own Your Home'] for column in categorical_cols: ax = sns.countplot(x=column, hue='Offer Accepted', data=trainset) total_counts = trainset[column].value_counts() for container in ax.containers: total = sum(bar.get_height() for bar in container) for bar in container: percentage = f'({bar.get_height() / total * 100:.2f})%' ax.annotate(percentage, xybar.get_x() + bar.get_width() / 2, bar.get_height(), xytext=(0, 3), textcoords='offset points', ha='center', va='bottom', fontsize=8, color='black') sns.set(rc={'figure.figsize': (15, 8)}) plt.show() In []: col = ['Reward', 'Mailer Type', 'Income Level', 'Overdraft Protection', 'Credit Rating', 'Household Size', 'Own Your Home'] sns.set(rc={'figure.figsize': (15, 8)}) for column in col: ax = sns.countplot(x=column, hue='Offer Accepted', data=trainset) for container in ax.containers: total_heights = [bar.get_height() for bar in container] sns_total = total_heights[1] if len(total_heights) > 1 else 0 no_total = total_heights[0] if len(total_heights) > 1 else 0 for i, bar in enumerate(container): total = sns_total if i % 2 == 1 else no_total ax.annotate(total, xybar.get_x() + bar.get_width() / 2, bar.get_height(), xytext=(0, 3), textcoords='offset points', ha='center', va='bottom', fontsize=8, color='black') plt.title(f'({column}) vs Offer Acceptance') plt.show() In []: #Descriptive statistics </pre> |
| <pre> In []: #Descriptive statistics In []: numerical_cols = ['Average Balance', 'Q1 Balance', 'Q2 Balance', 'Q3 Balance', 'Q4 Balance'] for column in numerical_cols: ax = sns.boxplot(x='Offer Accepted', y=column, data=trainset, palette=['red', 'blue'], width=0.4) plt.show() In []: numerical_cols = ['Bank Accounts Open', 'Credit Cards Held', 'Home Owned'] response_categories = ['No', 'Yes'] sns.set(rc={'figure.figsize': (15, 8)}) for column in numerical_cols: ax = sns.countplot(x=column, hue='Offer Accepted', data=trainset, hue_order=response_categories) for container in ax.containers: for i, bar in enumerate(container): total = bar.get_height() ax.annotate(total, xybar.get_x() + bar.get_width() / 2, bar.get_height(), xytext=(0, 3), textcoords='offset points', ha='center', va='bottom', fontsize=8, color='black') plt.title(f'({column}) vs Offer Acceptance') plt.show() In []: columns_to_plot = ['Bank Accounts Open', 'Credit Cards Held', 'Home Owned', 'Average Balance', 'Q1 Balance', 'Q2 Balance', 'Q3 Balance', 'Q4 Balance'] sns.pairplot(trainset[columns_to_plot]) plt.show() MCA In []: from prince import MCA In []: mca_cols = trainset.select_dtypes(['category']).columns print(len(mca_cols), 'features used for MCA are', mca_cols.tolist()) In []: df_encoded = pd.get_dummies(df, mca_cols) In []: df_encoded = df_encoded.drop(columns=['Customer Number']) In []: mca = MCA() mca_data=trainset[mca_cols] mca_data.head() mcal = mca.fit(mca_data) In []: mca.eigenvalues_summary In []: row_coordinates=mca.row_coordinates(mca_data) row_coordinates In []: column_coordinates=mca.column_coordinates(mca_data) column_coordinates In []: plt.figure(figsize=(10, 6)) plt.scatter(column_coordinates[0], column_coordinates[1], markers='o', s=10, color='red') for label, x, y in zip(column_coordinates.index, column_coordinates[0], column_coordinates[1]): plt.text(x, y, label, fontsize=5, ha='right', va='bottom') classes_to_highlight = ['Offer Accepted_No', 'Offer Accepted_Yes'] class_coordinates_to_highlight = column_coordinates.loc[classes_to_highlight] plt.scatter(class_coordinates_to_highlight[0], class_coordinates_to_highlight[1], markers='o', s=10, color='blue', labels='Offer Accepted') for label, x, y in zip(classes_to_highlight, class_coordinates_to_highlight[0], class_coordinates_to_highlight[1]): plt.text(x, y, label, fontsize=5, ha='right', va='bottom') plt.xlabel('Component 1') plt.ylabel('Component 2') plt.legend() plt.show() PLS - DA In []: from sklearn.cross_decomposition import PLSRegression from sklearn.model_selection import train_test_split from sklearn.metrics import accuracy_score, classification_report In []: X_train = trainset.drop('Offer Accepted', axis=1) y_train = trainset['Offer Accepted'] X_test = testset.drop('Offer Accepted', axis=1) y_test = testset['Offer Accepted'] In []: from sklearn.preprocessing import LabelEncoder label_encoder = LabelEncoder() X_train_encoded = pd.get_dummies(X_train) y_train_encoded = label_encoder.fit_transform(y_train) </pre> | <pre> X_test_encoded = pd.get_dummies(X_test) X_test_encoded = label_encoder.fit_transform(y_test) In []: n_components = 2 plsda = PLSRegression(n_components=n_components) In []: plsda.fit(X_train_encoded, y_train_encoded) X_train_plsda = plsda.transform(X_train_encoded) X_test_plsda = plsda.transform(X_test_encoded) In []: scores = plsda.fit(X_train_encoded, y_train_encoded).x_scores[:, :2] total_variance = np.var(X_train_encoded, axis=0).sum() scores_variance = np.var(scores, axis=0).sum() scores_variance variance_explained = scores_variance / total_variance print("Variance explained by the first two components: (variance_explained * 100)%") In []: import matplotlib.pyplot as plt unique_classes = np.unique(y_train_encoded) plt.figure(figsize=(8, 6)) for cls in unique_classes: mask = y_train_encoded == cls plt.scatter(X_train_plsda[mask, 0], X_train_plsda[mask, 1], label=f'Class {cls}') plt.title('PLS-DA Score Plot') plt.xlabel('PLS-DA Component 1') plt.ylabel('PLS-DA Component 2') plt.legend() plt.grid(True) plt.show() Correlation In []: corr_matrix = trainset.corr(method='pearson') print(corr_matrix) In []: from scipy.stats import spearmanr In []: trainset.corr(method='spearmanr') In []: plt.figure(figsize=(25,10)) sns.heatmap(trainset.corr(numeric_only=True), annot=True, cmap='Blues') Advanced Analysis In []: trainset['Offer Accepted'] = trainset['Offer Accepted'].replace(['No', '0', 'Yes', '1']) testset['Offer Accepted'] = testset['Offer Accepted'].replace(['No', '0', 'Yes', '1']) In []: trainset.head() In []: X_train = trainset.drop('Offer Accepted', axis=1) y_train = trainset['Offer Accepted'] X_test = testset.drop('Offer Accepted', axis=1) y_test = testset['Offer Accepted'] In []: ordinal_mapping = { 'Credit Rating': ['Low', 'Medium', 'High'], 'Household Size': ['Small', 'Big'], 'Income Level': ['Low', 'Medium', 'High']} for column, categories in ordinal_mapping.items(): X_train[column] = pd.Categorical(X_train[column], categories=categories, ordered=True).codes X_test[column] = pd.Categorical(X_test[column], categories=categories, ordered=True).codes In []: nominal_columns = ['Reward', 'Mailer Type', 'Income Level', 'Overdraft Protection', 'Credit Rating', 'Household Size', 'Own Your Home'] X_train = pd.get_dummies(X_train, columns=nominal_columns, drop_first=True) X_test = pd.get_dummies(X_test, columns=nominal_columns, drop_first=True) In []: X_train.head() In []: X_train = X_train.drop(columns=['Customer Number']) X_test = X_test.drop(columns=['Customer Number']) In []: y_train.head() In []: trainset.to_csv('trainset.csv', index=False) testset.to_csv('testset.csv', index=False) X_train.to_csv('X_train.csv', index=False) X_test.to_csv('X_test.csv', index=False) In []: X_train.dtypes </pre> |


```

print("Classification Report for the tuned Random Underampler model:")
print(report_underampler)

In [ ]:
from xgboost import XGBClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import classification_report
from sklearn.over_sampling import SMOTE
from sklearn.under_sampling import RandomUnderSampler
from sklearn.model_selection import StratifiedFold

param_grid_xgb = {
    "n_estimators": [50, 100, 200],
    "max_depth": [3, 5, 7],
    "learning_rate": [0.01, 0.1, 0.2],
    "subsample": [0.5, 1.0],
    "min_child_weight": [1, 5, 5],
}

smote = SMOTE(random_state=2)
X_resampled_smote, y_resampled_smote = smote.fit_resample(X_train_scaled, y_train)

xgb_classifier_smote = XGBClassifier(random_state=2)
grid_search_smote_xgb = GridSearchCV(xgb_classifier_smote, param_grid_xgb, cv=StratifiedFold(n_splits=5), scoring='f1', n_jobs=-1)
grid_search_smote_xgb.fit(X_resampled_smote, y_resampled_smote)

best_model_smote_xgb = grid_search_smote_xgb.best_estimator_

print("Best parameters for SMOTE XGBoost model:", grid_search_smote_xgb.best_params_)

y_pred_smote_xgb = best_model_smote_xgb.predict(X_test_scaled)

report_smote_xgb = classification_report(y_test, y_pred_smote_xgb)
print("Classification Report for the tuned SMOTE XGBoost model:")
print(report_smote_xgb)

undersampler = RandomUnderSampler(random_state=2)
X_resampled_undersampler, y_resampled_undersampler = undersampler.fit_resample(X_train_scaled, y_train)

xgb_classifier_undersampler = XGBClassifier(random_state=2)
grid_search_undersampler_xgb = GridSearchCV(xgb_classifier_undersampler, param_grid_xgb, cv=StratifiedFold(n_splits=5), scoring='f1', n_jobs=-1)
grid_search_undersampler_xgb.fit(X_resampled_undersampler, y_resampled_undersampler)

best_model_undersampler_xgb = grid_search_undersampler_xgb.best_estimator_

print("Best parameters for Random Underampler XGBoost model:", grid_search_undersampler_xgb.best_params_)

y_pred_undersampler_xgb = best_model_undersampler_xgb.predict(X_test_scaled)

report_undersampler_xgb = classification_report(y_test, y_pred_undersampler_xgb)
print("Classification Report for the tuned Random Underampler XGBoost model:")
print(report_undersampler_xgb)

In [ ]:
y_pred_smote = best_model_smote_xgb.predict(X_train_scaled)

report_smote = classification_report(y_train, y_pred_smote)
print("Classification Report for the tuned SMOTE model:")
print(report_smote)

y_pred_undersampler = best_model_undersampler_xgb.predict(X_train_scaled)

report_undersampler = classification_report(y_train, y_pred_undersampler)
print("Classification Report for the tuned Random Underampler model:")
print(report_undersampler)

In [ ]:
import matplotlib.pyplot as plt
import numpy as np
from matplotlib.colors import LinearSegmentedColormap

def generate_color_gradient(importances):
    norm_importances = (importances - min(importances)) / (max(importances) - min(importances))

    # Define a custom colormap starting from light blue and degrading towards dark blue
    cmap = LinearSegmentedColormap.from_list('custom_gradient', ['#99CCFF', '#66B2FF', '#0080FF', '#0071D7', '#0055A0', '#004B7A', '#00274C'])

    colors = cmap(norm_importances)
    return colors

feature_importances_smote = best_model_smote_xgb.feature_importances_
sorted_ids_smote = feature_importances_smote.argsort()

colors_smote = generate_color_gradient(feature_importances_smote)

plt.figure(figsize=(10, 6))
plt.barh(range(X_train_scaled.shape[1]), feature_importances_smote[sorted_ids_smote], color=colors_smote[sorted_ids_smote])
plt.xticks(range(X_train_scaled.shape[1]), X_train.columns[sorted_ids_smote])
plt.xlabel('Feature Importance')
plt.ylabel('Feature')
plt.title('Variable Importance Plot for SMOTE-Random Forest model')
plt.show()

feature_importances_undersampler = best_model_undersampler_xgb.feature_importances_
sorted_ids_undersampler = feature_importances_undersampler.argsort()

colors_undersampler = generate_color_gradient(feature_importances_undersampler)

plt.figure(figsize=(10, 6))
plt.barh(range(X_train_scaled.shape[1]), feature_importances_undersampler[sorted_ids_undersampler], color=colors_undersampler[sorted_ids_undersampler])

```

```

print("Best parameters for SMOTE XGBoost model:", grid_search_smote_xgb.best_params_)

y_pred_smote_xgb = best_model_smote_xgb.predict(X_test_scaled_selected)

report_smote_xgb = classification_report(y_test, y_pred_smote_xgb)
print("Classification Report for the tuned SMOTE XGBoost model:")
print(report_smote_xgb)

In [ ]:
y_pred_smote = best_model_smote_xgb.predict(X_train_scaled_selected)

report_smote = classification_report(y_train, y_pred_smote)
print("Classification Report for the tuned SMOTE model:")
print(report_smote)

y_pred_smote_xgb = best_model_smote_xgb.predict(X_train_scaled_selected)

report_smote_xgb = classification_report(y_train, y_pred_smote_xgb)
print("Classification Report for the tuned SMOTE XGBoost model:")
print(report_smote_xgb)

Voting Classifier

In [ ]:
from sklearn.ensemble import VotingClassifier
from sklearn.metrics import classification_report
import joblib

# Load the pre-trained models
best_reduced_model_knn = joblib.load("smote_reduced_knn_model.joblib")
xgb_reduced_smote = best_reduced_model_smote_xgb
rf_reduced_smotes = best_reduced_model_smote

voting_classifier = VotingClassifier(
    estimators=[
        ("xgb_smote", xgb_reduced_smote),
        ("rf_smote", rf_reduced_smote),
        ("knn", best_reduced_model_knn)
    ],
    voting='soft'
)

voting_classifier.fit(X_resampled_smote, y_resampled_smote)

y_pred_ensemble = voting_classifier.predict(X_train_scaled_selected)

report_ensemble = classification_report(y_train, y_pred_ensemble)
print("Classification Report for the Ensemble Model:")
print(report_ensemble)

y_pred_ensemble = voting_classifier.predict(X_test_scaled_selected)

report_ensemble = classification_report(y_test, y_pred_ensemble)
print("Classification Report for the Ensemble Model:")
print(report_ensemble)

Partial Dependence plots

In [ ]:
from sklearn.inspection import PartialDependenceDisplay

import matplotlib.pyplot as plt
features_to_plot = [
    'Q1 Balance', 'Q2 Balance', 'Q3 Balance', 'Q4 Balance',
    'Average Balance', 'Credit Rating_2', 'Mailer Type_Postcard',
    'Household_Cash Back', 'Credit Rating_1', 'A Credit Cards Held']

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Q1 Balance'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Q2 Balance'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Q3 Balance'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Q4 Balance'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Average Balance'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Credit Rating_2'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Mailer Type_Postcard'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Household_Cash Back'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'A Credit Cards Held'])
plt.show()

disp = PartialDependenceDisplay.from_estimator(best_reduced_model_smote, X_train_scaled_selected,
                                              features_to_plot=[
                                                  'Credit Rating_1'])
plt.show()

```