

Date of Submission Sunday, August 27, 2023.

Digital Object Identifier

Exploring Credit Card Approval Prediction with K-Nearest Neighbors (KNN) Classification

PRAJESH SHRESTHA¹, and UJJWAL PAUDEL¹
¹Department of Electronics and Computer Engineering, Thapathali Campus, Institute of Engineering, Tribhuvan University, Kathmandu, Nepal (e-mail: prajesh.762417@thc.tu.edu.np, ujjwal.762417@thc.tu.edu.np)

ABSTRACT K-Nearest Neighbors (KNN) is a machine learning algorithm that classifies data points by examining their proximity to neighboring data points. The study explores the fundamental principles of KNN and employing them to perform an extensive classification analysis on the renowned Credit Card Approval dataset. Our analysis involves a thorough examination of the dataset and the process of selecting the most important features used in training the classification model. The experimental results demonstrate a significant accuracy rate of 78% achieved through the KNN classification approach. The results of this study suggest that the KNN algorithm can be used to predict whether an applicant will be approved for a credit card.

INDEX TERMS Credit Card Approval, Classification, K-Nearest Neighbours (KNN)

I. INTRODUCTION

K-Nearest Neighbors (KNN) is a machine learning technique used for classification tasks. It operates by considering the neighboring data points to classify new instances. This approach makes KNN suitable for pattern recognition and classification tasks. It's particularly effective when the relationship between data points is complex or not easily defined by a simple equation. [1]

The Credit Card Approval dataset is a valuable resource for studying the factors that influence credit card approval. The dataset includes a variety of applicant characteristics, such as personal, financial, and credit-related information. KNN operates by identifying similar applicants based on their attributes, allowing the algorithm to classify new applications by considering the characteristics of their nearest neighbors. This dataset serves as a valuable resource for assessing and predicting the approval status of credit card applications.

II. METHODOLOGY

A. THEORY

The K-Nearest Neighbors (KNN) algorithm is a non-parametric supervised learning algorithm. It is used for both classification and regression tasks. Its primary objective is to predict the class of a given instance by assessing the classes of its neighboring data points. This algorithm capitalizes on the intuitive idea that instances with similar attributes often belong to the same class.

KNN's ability to make predictions based on local patterns within the data makes it particularly useful for tasks involving

complex relationships that might not be easily captured by a formal equation. [2] [3]

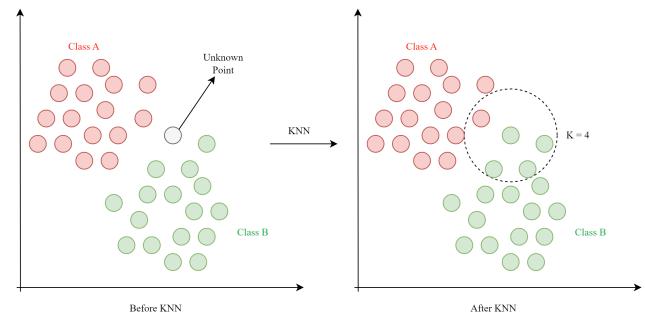


FIGURE 1. K-Nearest Neighbors Visualization

1) Choosing Number of Neighbors (k)

'K' is a hyperparameter that can be tuned to improve the accuracy of the model. A common rule of thumb for the starting value can be an odd number closer to the square root of the number of instances. A smaller 'k' value makes the classification sensitive to noise and outliers, potentially leading to overfitting. Conversely, a larger 'k' value smoothes out decision boundaries but might cause the algorithm to overlook local patterns. Selecting the optimal 'k' can be achieved through techniques like cross-validation or the elbow method, which involves assessing the algorithm's performance for different 'k' values and identifying the point of diminishing returns.

B. ALGORITHM

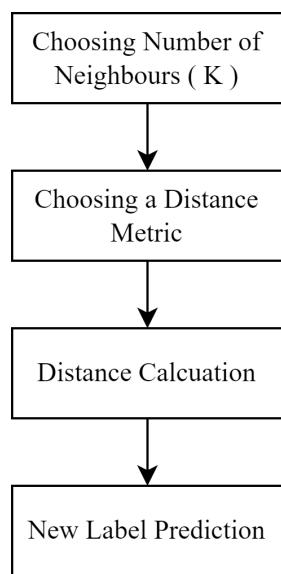


FIGURE 2. KNN Algorithm

1) Choosing a Distance Metric

The distance metric determines how the proximity or similarity between data points is measured. The most common distance metrics are Euclidean distance, Manhattan distance, and Minkowski distance.

2) Distance Calculation

In the "Distance Calculation" step of the K-Nearest Neighbors (KNN) algorithm, distances are computed between the new data point and all instances in the training dataset using the chosen distance metric. Subsequently, these distances are sorted in ascending order to identify the closest neighbors. The next phase involves selecting the K nearest neighbors from the sorted list.

3) New Label Prediction

The class label for the new data point is determined based on the classes of its K nearest neighbors. This involves two common strategies: majority voting and weighted voting.

Majority Voting:

In majority voting, the class label that appears most frequently among the K nearest neighbors is assigned to the new data point. This method assumes that the majority class among the neighbors is likely to be the correct label for the new instance.

$$\text{PredictedClass} = \text{Mode}(\text{ClassLabels of KNearstNeighbors}) \quad (1)$$

Weighted Voting:

In weighted voting, the class labels of the K nearest neighbors are taken into account along with their distances from the new data point. Closer neighbors have a stronger influence on

the prediction. Each neighbor's vote is weighted by a factor inversely proportional to its distance from the new instance.

$$\text{PredictedClass} = \frac{\text{WeightedSumofClassLabels}}{\text{SumofWeights}} \quad (2)$$

Weighted voting gives more importance to closer neighbors and can improve accuracy when some neighbors are more relevant than others.

C. SYSTEM BLOCK DIAGRAM

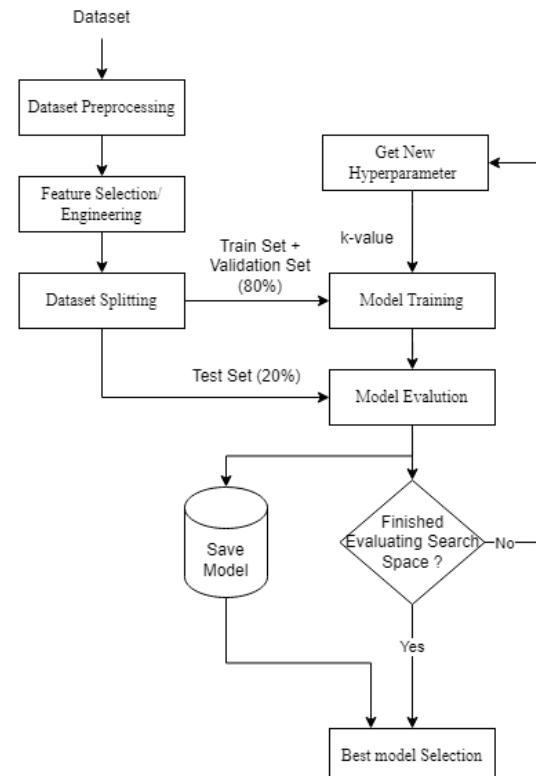


FIGURE 3. System Block Diagram

1) Data Preprocessing

Dataset Preprocessing involves cleaning, transforming, and organizing the raw data to make it suitable for analysis. It includes tasks like handling missing values, scaling features, and encoding categorical variables.

2) Feature Selection/Engineering

In this step, tasks are performed to select relevant features or create new features that can enhance the model's performance. Correlation analysis is used to identify relationships between features and the target variable and keep the relevant features.

3) Dataset Splitting

This block involves dividing the preprocessed data into training, validation, and testing subsets. The training set is used to train the model, the validation set helps tune hyperparameters, and the testing set evaluates the final model's performance.

4) Model Training

This block represents the training process of your machine learning model using the training dataset. It involves selecting an appropriate algorithm, feeding the data, and adjusting model parameters.

5) Model Evaluation

After training, the model's performance needs to be evaluated using the validation and testing datasets. Metrics like accuracy, precision, recall, F1-score, and others can be used to measure how well the model is performing.

6) Hyperparameter Tuning

This block includes optimizing the hyperparameters of the model to improve its performance. Techniques like grid search, random search, or Bayesian optimization can be used to find the best hyperparameter values.

D. INSTRUMENTATION TOOLS

1) Pandas and Numpy

- Pandas is employed for data processing, providing data structures like DataFrames that facilitate data manipulation and analysis.
- NumPy serves as the foundation for linear algebra operations, enabling efficient handling of arrays and matrices.

2) Matplotlib and Plotly

- Matplotlib and Plotly are visualization libraries utilized to create informative and visually appealing graphs, charts, and plots, aiding in data exploration and presentation.
- Matplotlib offers various plotting capabilities, while Plotly provides interactive and dynamic visualizations.

3) Seaborn

- Seaborn complements Matplotlib by providing a higher-level interface for statistical graphics, making it easier to create aesthetically pleasing and informative plots.

4) Scikit-learn

- sklearn.metrics*: This submodule within scikit-learn provides various metrics for evaluating model performance, including accuracy, precision, recall, F1-score, ROC curves, and more.
- sklearn.model_selection.train_test_split*: Used to split the dataset into training and testing subsets, a fundamental step in model evaluation.
- sklearn.neighbors.KNeighborsClassifier*: Implements the K-Nearest Neighbors classifier, a core component of your classification workflow.
- sklearn.preprocessing.LabelEncoder*: Employed for label encoding of categorical columns, transforming categorical data into numeric format suitable for machine learning algorithms.

III. RESULTS AND DISCUSSIONS

A. EXPLORATORY DATA ANALYSIS (EDA)

A credit score plays a pivotal role in assessing and approving credit card applications. Within this dataset, we are presented with two distinct sets of records. While the precise origins of this data collection remain undisclosed, it could either represent fictitious information or a sample from the real world. The overarching aim here is to conduct an insightful analysis to discern the specific features or categories that financial institutions typically record from individuals. These records are integral in their evaluation of whether a prospective client possesses a favorable history of credit card payments.

This dataset comprises two distinct sets of records. The first pertains to application data, encompassing information about individuals applying for credit cards. The second dataset focuses on credit status, providing a historical account of applicants' behavior regarding their repayment of credit loans. Evidently, these two datasets were integrated based on a shared client ID to create a unified dataset, from which valuable insights and patterns were extracted.

At the outset, a thorough examination of each record was conducted to scrutinize their attributes, row counts, and data types. Table 1 provides insights into the applicant data, revealing the presence of eight categorical data types, while the remaining data types are predominantly numerical. Conversely, Table 2 sheds light on the credit status data, comprising only three features. Notably, the categorical attribute 'STATUS' is identified as the target label, encompassing seven distinct notations, each signifying various actions taken by the clients.

TABLE 1. Description of Applicant Data

S.N.	Column	Non-Null Count	Dtype
1	ID	438,557	int64
2	CODE_GENDER	438,557	object
3	FLAG_OWN_CAR	438,557	object
4	FLAG_OWN_REALTY	438,557	object
5	CNT_CHILDREN	438,557	int64
6	AMT_INCOME_TOTAL	438,557	float64
7	NAME_INCOME_TYPE	438,557	object
8	NAME_EDUCATION_TYPE	438,557	object
9	NAME_FAMILY_STATUS	438,557	object
10	NAME_HOUSING_TYPE	438,557	object
11	DAYS_BIRTH	438,557	int64
12	DAYS_EMPLOYED	438,557	int64
13	FLAG_MOBIL	438,557	int64
14	FLAG_WORK_PHONE	438,557	int64
15	FLAG_PHONE	438,557	int64
16	FLAG_EMAIL	438,557	int64
17	OCCUPATION_TYPE	304,354	object
18	CNT_FAM_MEMBERS	438,557	float64

From Table 1, a noticeable observation emerges when examining the non-null count for the 'OCCUPATION_TYPE'

TABLE 2. Description of Credit Status Data

S.N.	Column	Non-Null Count	Dtype
1	ID	1,048,575	int64
2	MONTHS_BALANCE	1,048,575	int64
3	STATUS	1,048,575	object

feature; it appears to be significantly lower than that of other features. To substantiate this observation, a missing value plot, as shown in Fig. 4, was generated, confirming the presence of a substantial number of missing values within the 'OCCUPATION_TYPE' category. Further investigation revealed that the total count of 'NAN' values in this feature amounted to 134,203, which represents a considerable proportion of the overall dataset. To address this issue effectively, the decision was made to remove the 'OCCUPATION_TYPE' attribute from the dataset.

Each applicant is presumably uniquely identified by their ID number. Under this assumption, it was determined that there were only 47 instances of duplicate IDs within the applicant record data. These duplicates were subsequently removed from the dataset to ensure data integrity and consistency.

Both datasets contain some intriguing features. For instance, the 'DAYS_BIRTH' column represents the integer value of the number of days an individual has been alive, measured as a negative count from a specific reference date backward to their actual birthdate. Consequently, each entry in this column is represented as a negative number. To enhance interpretability, these negative values were transformed into the age of the individual.

For 'DAYS_EMPLOYED' column, which exhibits positive values when the individual is unemployed and negative values when the applicant has been employed for a certain number of days in the past. For simplification, positive values, representing unemployment, were replaced with zero, while negative values were retained to indicate the number of years the individual has been employed.

While the applicant records consist of a total of 18 features to define each data instance, it's important to note that not all of these features carry significant information for data mining and modeling purposes. For instance, the 'FLAG_MOBIL' column, which indicates whether an individual possesses a mobile phone or not, uniformly reveals that every client in the dataset indeed has a mobile phone. Consequently, this column was deemed redundant and was subsequently removed from the dataset.

Similarly, columns like 'FLAG_WORK_PHONE' and 'FLAG_PHONE,' which capture information regarding whether the clients provided their personal or work phone numbers, along with the 'FLAG_EMAIL' column, which signifies whether clients submitted their email addresses, were found to lack meaningful variation. Consequently, these columns were also eliminated from the dataset, streamlining it for more focused and efficient data analysis and modeling.

The issue of granting credit cards to applicants is inherently linked to outlier detection. Detecting instances of bad client behavior, particularly delayed payments on credit loans, is relatively rare within a large dataset. However, it's worth noting that this dataset contains numerical attributes that themselves exhibit outlier values. These outliers must be addressed and removed to facilitate accurate modeling and improve the effectiveness of subsequent analyses, particularly when applying techniques like K-Nearest Neighbors (KNN).

Figure 5 presents a series of box plots illustrating the numerical categories found within the applicant dataset, revealing several notable patterns. First, when examining the number of children, it becomes evident that the majority of applicants either have one or two children or none at all. This aligns with contemporary family sizes, where smaller households are prevalent. In terms of income, most credit card applicants appear to fall within the middle-class income range, with incomes showing relatively low variance. However, there exists a subset of applicants who exhibit significantly higher incomes. This could be attributed to individuals with a deep understanding of financial management, who comfortably handle regular expenses with credit cards while channeling their income into investments, given the minimal interest cost of credit card loans for them.

Regarding the age distribution of applicants, the majority cluster between 33 and 53 years old. This age range typically signifies individuals with stable employment and diverse income sources, providing them with the confidence to manage credit card payments effectively. The box plot for 'YEARS_EMPLOYED' reveals numerous outliers, notably individuals with employment durations exceeding 20 years. This may suggest government employees, who tend to have more secure and extended tenures compared to those in the private sector. Lastly, the 'CNT_FAM_MEMBERS' distribution highlights numerous outliers, indicating that certain applicants have significantly larger families compared to the majority. These insights shed light on the diverse characteristics of credit card applicants and their financial profiles, offering valuable information for subsequent analyses and decision-making.

To address the presence of outlier values within the dataset, all outliers, with the exception of those pertaining to 'AGE_YEARS,' were systematically removed. This removal process was executed by determining the lower and upper bounds for each numerical feature using their respective interquartile range (IQR).

Upon examining the Credit Status dataset, it becomes apparent that the target label, denoted as 'STATUS,' encompasses seven distinct labels. These labels encapsulate the applicant's historical behavior with respect to previous credit card loans. Given the multi-class nature of 'STATUS,' each class holding varying implications, a crucial preprocessing step involves mapping this multiclass problem into a sensible binary classification scheme. In this regard, we have defined 'Good Clients' as those who have paid in the same month ('C' class) or those who have not taken any loan ('X' class), based

on a moral evaluation of their credit behavior. Conversely, 'Bad Clients' are identified as individuals falling within the '2-5 STATUS' classes, indicating a delay in repayment. Additionally, clients in the '0' class, signifying a one-month delay, are also categorized as 'Good Clients,' whereas those in the '1' class, implying a delay of more than one month, are classified as 'Bad Clients.' This binary classification scheme offers a more interpretable and practical approach to modeling and analyzing credit client behavior, allowing for clearer insights into credit risk assessment.

Implementing the aforementioned procedure has led to the creation of the target label, as illustrated in the pie chart displayed in Figure 6. This visualization unmistakably highlights a substantial imbalance within the dataset, with approximately 99% of individuals categorized as 'Good Clients' and only 1% falling into the 'Bad Client' category.

In the context of predictive modeling, the primary objective is to effectively identify potential 'Bad Client' behavior among applicants. However, given the pronounced class imbalance prevalent in this dataset, it is paramount to employ specific techniques to rectify this imbalance, thereby rendering the dataset more amenable for machine learning algorithms and ensuring equitable and precise predictions.

To address this challenge, a strategic approach was implemented. Since each client record encompasses multiple instances of credit card-related activities, these activities were consolidated, and only the aggregate outcome, whether 'Good' or 'Bad,' was considered. This consolidation led to a revised dataset where the pie chart in Figure 7 illustrates a more balanced distribution, with approximately 88% classified as 'Good Clients' and 12% categorized as 'Bad Clients.' This rebalanced dataset better reflects the real-world scenario and enhances the predictive capabilities of the modeling process.

The two datasets were merged based on their shared client 'ID' to create a unified dataset. Given that the 'ID' feature does not hold any significance for the overall dataset, this column was subsequently removed. Furthermore, duplicate records were eliminated from the combined dataset. As a result, the final dataset comprised a total of 11,058 rows and included 12 features along with a class attribute. The distribution of the target variable from this unified dataset is depicted in the pie chart in Fig. 9.

Figure 10 illustrates a correlational heatmap depicting the relationships between all numerical features and the target attribute 'STATUS.' Upon examination, there is no feature that exhibits a high correlation with the target attribute.

Figure 11 presents a pie chart depicting the gender distribution among credit card applicants, while Figure 8 showcases the percentage of approved credit card applications by gender. Notably, females are the dominant group both in terms of credit card application submissions and approvals. This observation indirectly underscores the growing empowerment of women and their significant contribution to economic development, marking a contrast to societal norms of the past where they were predominantly confined to household roles.

The pie charts in Fig. 12 and Fig. 13 provide insights into the financial aspects and property ownership status of the applicants. These two features hold the potential to yield significant insights into the behavior of individuals when it comes to credit card loan repayments.

The bar plot in Fig. 14 provides a compelling insight into the number of children that applicants could potentially have. The data reveals an interesting pattern, with nearly 70% of applicants indicating that they do not have children. This observation can be interpreted as a reflection of the modern societal trend where many individuals or couples prioritize establishing financial stability before starting a family. Consequently, those who are married but have not yet achieved this level of financial stability may be more inclined to apply for credit cards, as they seek to fulfill their needs and manage their expenses through monthly payments.

Figure 15 illustrates the distribution of annual income, clearly indicating that the majority of applicants report an annual income ranging from 1 to 2 lakhs. The pie chart in Fig. 16 provides insights into the sources of income for the applicants. In general, a significant portion of applicants is likely to be employed in various sectors, enabling them to meet their monthly payment obligations. On the other hand, students comprise the smallest percentage, accounting for only 0.03% of the total applicants. This reflects the fact that most students do not engage in employment to repay credit loans.

Figure 17 presents a bar plot depicting the educational attainment of the clients. Upon examination, it becomes evident that the majority have completed secondary education, with higher education being the second most common level of achievement. In Fig. 18, the family status of the clients is illustrated. A significant portion of them are married, which indirectly underscores the notion that marriage often leads to increased expenses and greater financial responsibility to support the family. This dynamic can potentially drive individuals to seek credit card loans as a means to fulfill their desires and financial obligations. Fig. 19 displays the various types of housing arrangements preferred by the applicants, with a majority residing in either apartments or their own houses. Fig. 20, simply shows the age distribution of the clients.

Figure 21 reveals the distribution of the number of years that applicants have been employed. A prominent spike at the 0th value is noticeable, which corresponds to the unemployed category. This spike underscores that individuals without employment, designated as '0,' constitute a significant portion of the applicant pool. This observation implies that the unemployed are often in greater need of financial support, motivating them to apply for credit cards.

Finally, the scatter plot in Fig. 22 provides a compelling and insightful perspective on the dataset. It reveals that applicants with fewer years of employment or those who are unemployed, coupled with lower total income, are the ones whose credit card applications are predominantly rejected. This decision by the financial institution appears to be both

reasonable and understandable

B. FEATURE SELECTION

While irrelevant and unimportant features were already removed during the preprocessing steps, it was essential to compare both categorical and numerical attributes on a common basis to assess their correlation. To facilitate this, all categorical attributes were transformed into numerical attributes using label encoders. Subsequently, a correlation plot depicting the relationships between all features and the target label, 'STATUS,' is presented in the heatmap shown in Fig. 23.

Upon careful observation of the heatmap, it becomes evident that there is no discernible sign of significant correlation between the features and the target labels. In such a scenario, pursuing further feature selection or assigning higher weights to positively correlated features while applying Weighted KNN does not appear to be a suitable strategy.

C. K-NEAREST NEIGHBORS MODELING

The dataset was divided into a training set and a testing set in an 80:20 ratio. Subsequently, a basic KNN model was employed with a chosen k-value of 7. The results of this classification task are presented in Table 3. Notably, the accuracy achieved through this modeling effort peaked at 77%.

TABLE 3. Classification Report: KNN Model with K = 7

Class	Precision	Recall	F1-Score	Support
Good Client	78.85%	96.85%	86.93%	1744
Bad Client	21.43%	3.21%	5.58%	468
Accuracy			77.03%	2212
Macro Avg.	50.14%	50.03%	46.25%	2212
Weighted Avg.	66.70%	77.03%	69.72%	2212

The support counts for the 'Good Client' and 'Bad Client' labels underscore the dataset's imbalance, a factor that significantly impacts model performance. The corresponding confusion matrix for this modeling effort is displayed in Fig. 24.

Furthermore, it's worth noting that although weighted-KNN may not be the ideal choice for this dataset, an experiment was conducted wherein weighted-KNN was applied with random initialization and assigned high values to certain features. Surprisingly, this experiment yielded the same accuracy of 78.84%. However, the confusion matrix from this experiment, shown in Fig. 25, reveals an interesting outcome. The model predominantly predicts 'Good Clients,' with no predictions made for 'Bad Clients,' indicating that applying weighted-KNN to features that exhibit zero correlation with the target label does not yield meaningful results.

To determine the optimal number of nearest neighbors (k) for the KNN model, a grid search cross-validation was conducted, employing stratified k-folding with k values ranging from 2 to 20 neighbors. This rigorous process revealed that the highest accuracy was achieved when k was set to 10. In a 5-

fold stratified cross-validation, the average accuracy reached 78%.

However, when the experiment was replicated using a basic KNN model with k equal to 10, the accuracy dropped significantly to just 63.88%, as detailed in the classification report presented in Table 4. This observation highlights the nuanced nature of model performance, where the selected hyperparameter, k, can impact the outcomes differently depending on the specific context and dataset.

TABLE 4. Classification Report: GridSearchCV Applied KNN Model

Class	Precision	Recall	F1-Score	Support
Good Client	77.78%	75.86%	76.81%	1744
Bad Client	17.61%	19.23%	18.39%	468
Accuracy			63.88%	2212
Macro Avg.	47.70%	47.55%	47.60%	2212
Weighted Avg.	65.05%	63.88%	64.45%	2212

The confusion matrix for this specific experiment is depicted in Fig. 26. The disparity in performance between cross-validation and the simple experiment could likely be attributed to the quality of the data partitioning into training and validation sets. This explanation gains credence, particularly considering the highly imbalanced nature of the dataset.

Indeed, the dataset's significant class imbalance, with a limited number of 'Bad Client' instances compared to 'Good Client' instances, plays a pivotal role. As evident from the classification report, the support count for 'Bad Client' is a mere 468, while 'Good Client' instances amount to 1744. This imbalance can skew model learning and influence performance, especially in situations where random partitioning may lead to an insufficient representation of 'Bad Client' cases in the training and testing sets.

In an effort to address the class imbalance issue, Synthetic Minority Over-sampling Technique (SMOTE) was applied to create synthetic instances of the minority class ('Bad Client'). However, the outcomes from this strategy did not yield significant improvements in the KNN model's performance. The classification report for this experiment is detailed in Table 5.

TABLE 5. Classification Report: KNN Model on Balanced Dataset using SMOTE

Class	Precision	Recall	F1-Score	Support
Good Client	60.49%	75.86%	67.31%	1744
Bad Client	38.99%	23.74%	29.51%	1133
Accuracy			55.34%	2877
Macro Avg.	49.74%	49.80%	48.41%	2877
Weighted Avg.	52.02%	55.34%	52.43%	2877

Several factors could contribute to the limited success of SMOTE in improving the KNN model's performance in this context. Firstly, while SMOTE is effective in generating synthetic instances to balance class distributions, its success can be influenced by the complexity and nature of the dataset. In

situations where the decision boundaries between classes are intricate and not easily discernible, as may be the case in this dataset, generating synthetic instances might not effectively capture the underlying patterns.

Secondly, the choice of k-value in the KNN algorithm becomes critical, especially after applying SMOTE. Selecting an appropriate k-value is inherently challenging, and an incorrect choice can lead to the model overfitting or underperforming. In this specific experiment, despite optimizing k via cross-validation, the model may still be sensitive to the intricate class distribution, leading to suboptimal results.

Additionally, SMOTE introduces synthetic data points that can potentially introduce noise into the dataset, further complicating the model's ability to distinguish between genuine and synthetic instances. This noise can impact the overall model performance.

In conclusion, the suitability of the dataset for KNN classification appears limited. Despite various attempts, it becomes evident that relying solely on the KNN model, or any method within its purview, may not yield promising results for this particular dataset.

IV. CONCLUSION

In conclusion, we explored the practical implementation of the K-Nearest Neighbors (KNN) algorithm for classification tasks, using the Credit Card Approval Dataset. Through a systematic approach, we undertook data analysis, feature enhancement, dataset balancing, model training, hyperparameter tuning, and rigorous evaluation. As an outcome, an accuracy rate of 78% achieved by the KNN model, thus demonstrating how KNN can be a valuable tool for making predictions in real-world situations.

REFERENCES

- [1] J. H. Friedman, J. L. Bentley, and R. A. Finkel, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, 1967.
- [2] D. Angluin and L. J. Valiant, "k-Nearest Neighbors in the 21st Century," in *Proceedings of the Twenty-Ninth Annual ACM Symposium on Theory of Computing*, 1997. DOI: 10.1145/258533.258660.
- [3] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient k-means clustering algorithm: Analysis and implementation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2002. DOI: 10.1109/34.677568.



UJJWAL P. is a dedicated undergraduate student currently pursuing Computer Engineering at Thapathali College, Institute of Engineering (IOE), Nepal. With a keen interest in Machine Learning and Artificial Intelligence, he embraces the role of a passionate learner, constantly exploring and experimenting with various ML/AI models. Driven by a passion to make meaningful contributions to the AI field, he is committed to continuous growth and learning in this exciting and dynamic domain.

For any inquiries or potential collaborations, you can reach out to Ujjwal at ujjwal.762416@thc.tu.edu.np.



PRAJESH S. is an enthusiastic undergraduate student pursuing Computer Engineering at Thapathali College, Institute of Engineering (IOE), Nepal. Curious about the vast world of Machine Learning and Artificial Intelligence, Prajesh enjoys immersing himself in learning and experimenting with diverse ML/AI models. He is fascinated by the potential applications of these technologies and aspires to utilize them for addressing real-world challenges. For any queries or potential collaborations, you can contact Prajesh at prajesh.762417@thc.tu.edu.np.

APPENDIX A PLOTS AND FIGURES

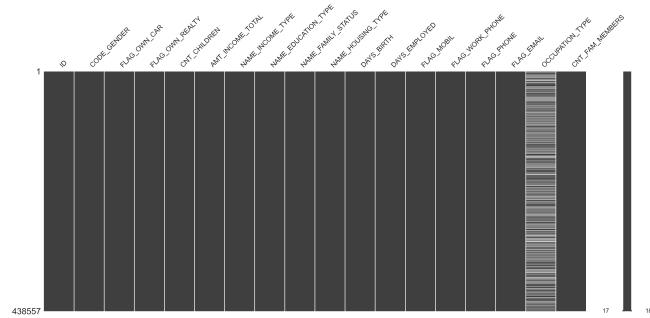


FIGURE 4. Missing Values Visualization Graph

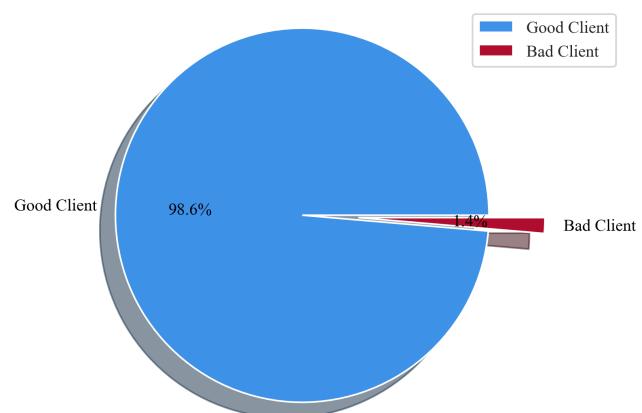


FIGURE 6. Pie Chart: Initial Target Distribution

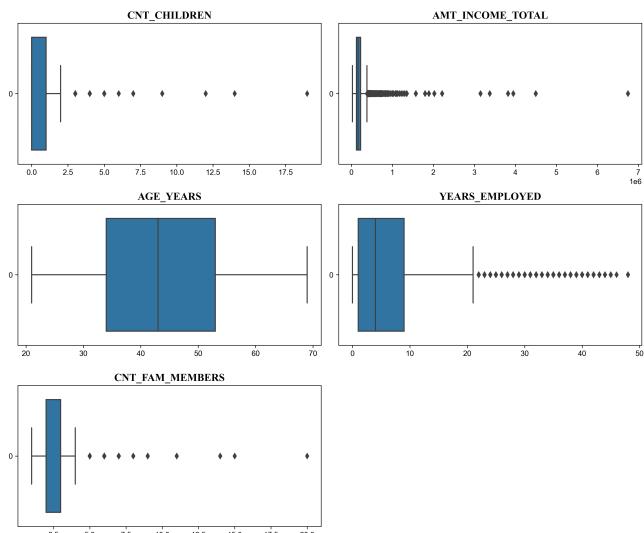


FIGURE 5. Box Plot: Outliers in each Numerical Attributes

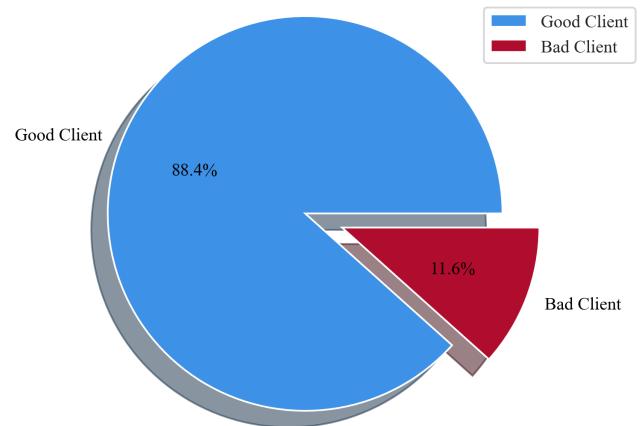
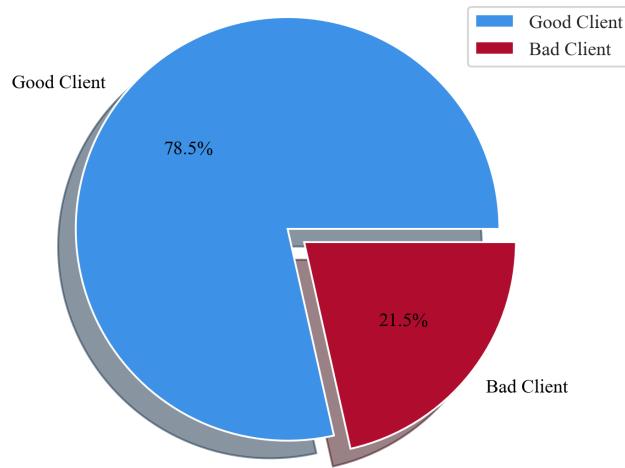
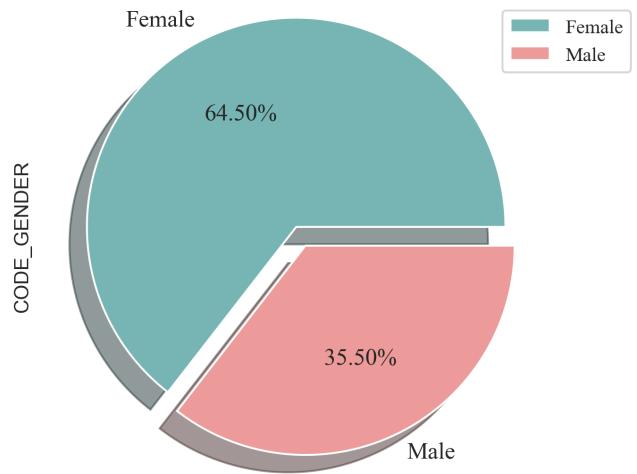
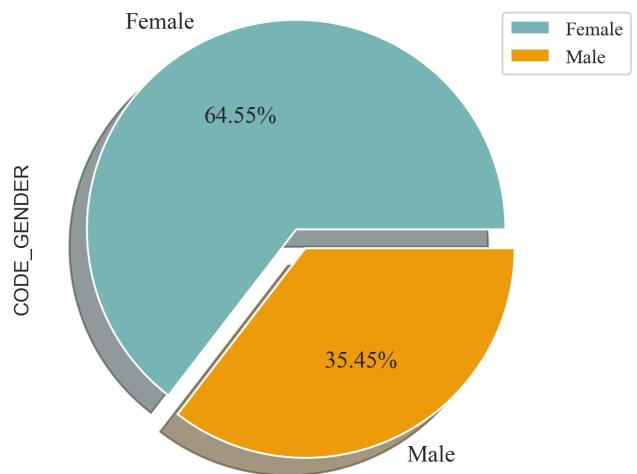
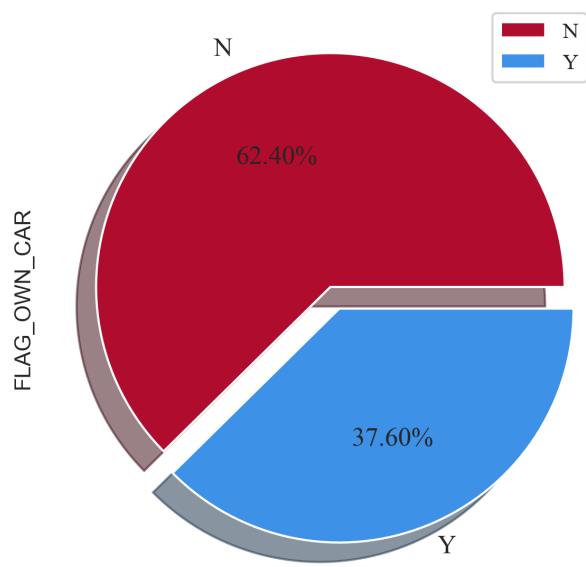
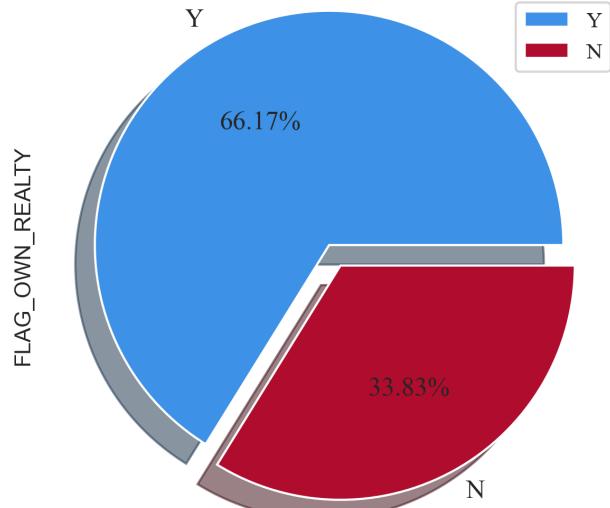
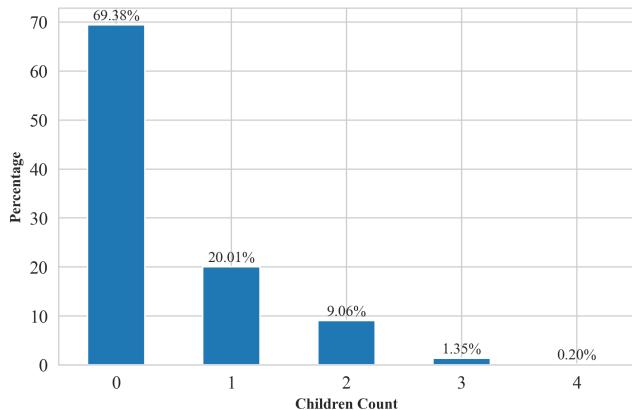
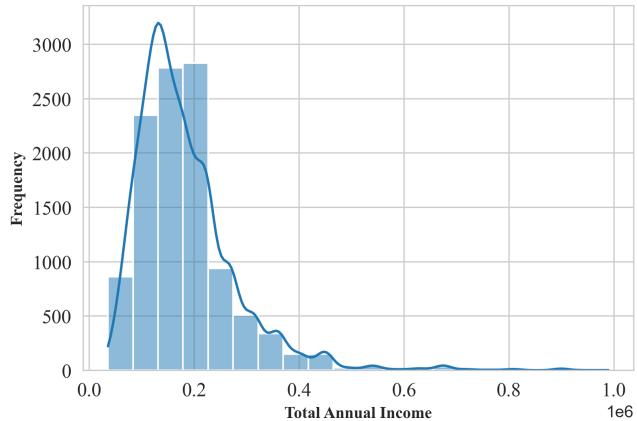
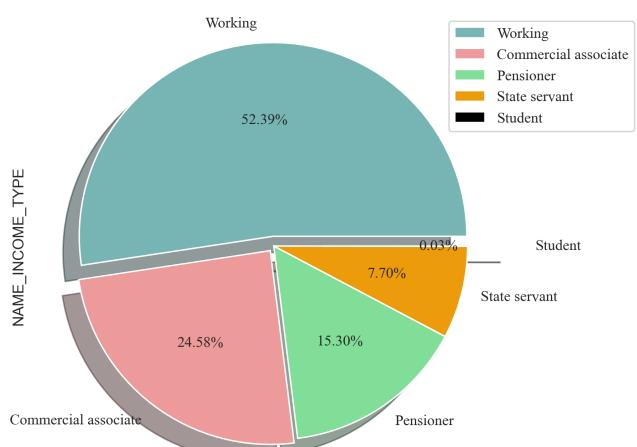


FIGURE 7. Pie Chart: Aggregated Target Distribution

**FIGURE 8.** Pie Chart: Final Target Distribution**FIGURE 10.** Pie Chart: Percentage of Application Submission based on Gender**FIGURE 9.** Correlational Heat Map for Numerical Attributes**FIGURE 11.** Pie Chart: Percentage of Application Approved based on Gender

**FIGURE 12.** Pie Chart: Percentage of Applicants who owns Car**FIGURE 13.** Pie Chart: Percentage of Applicants who has Real Estate Property**FIGURE 14.** Children Counts Bar Plot**FIGURE 15.** Histogram Plot: Total Annual Income**FIGURE 16.** Pie Chart: Percentage of Applicants based on Income Types

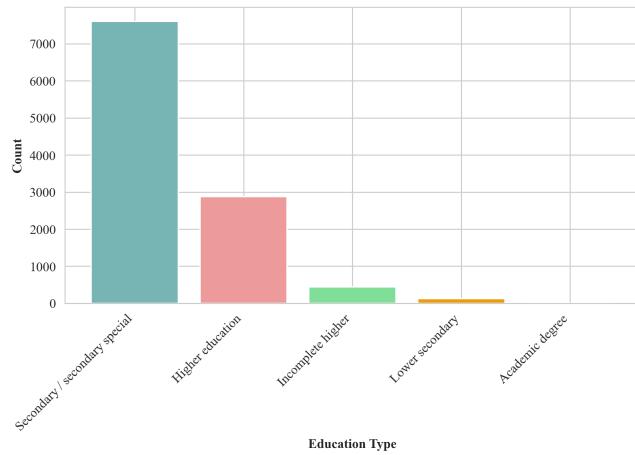


FIGURE 17. Pie Chart: Percentage of Applicants based on Education Acquired

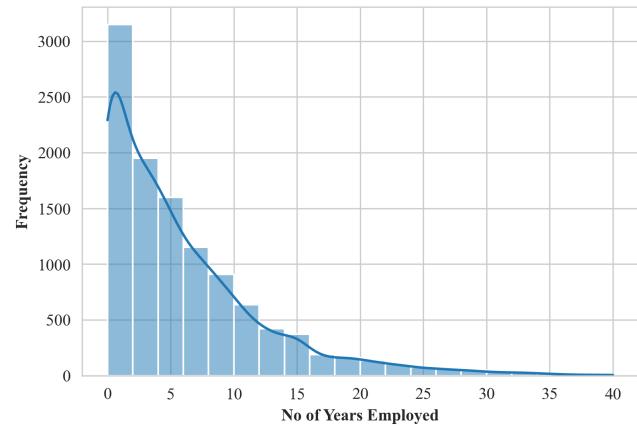


FIGURE 20. Histogram Plot: Applicants Years of Employment

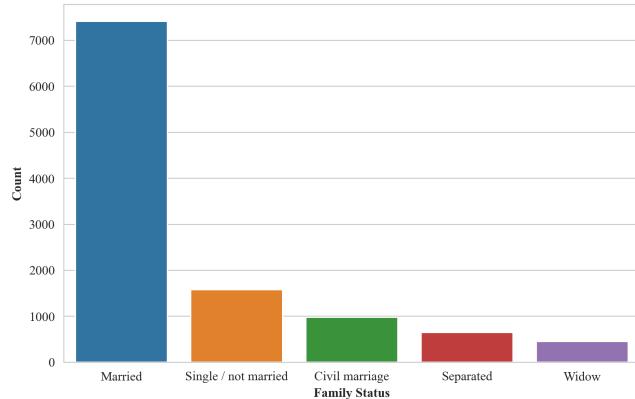


FIGURE 18. Pie Chart: Percentage of Applicants based on Family Status Types

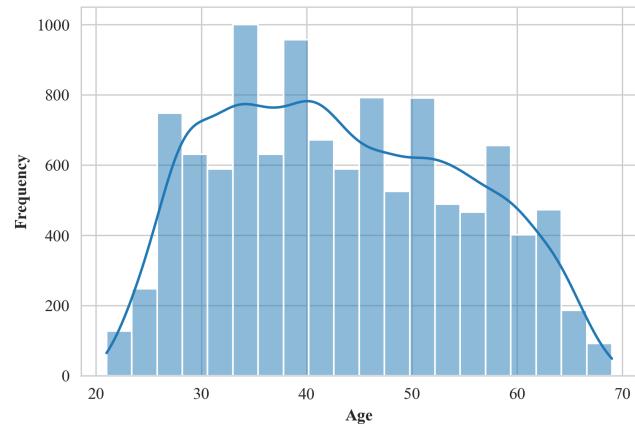


FIGURE 19. Histogram Plot: Applicant Age Distribution

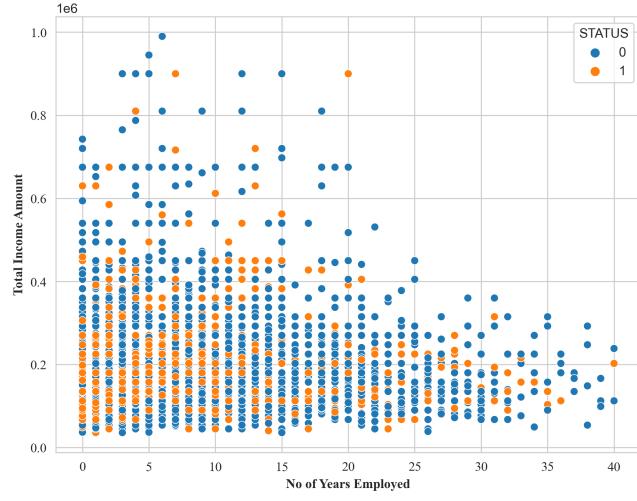
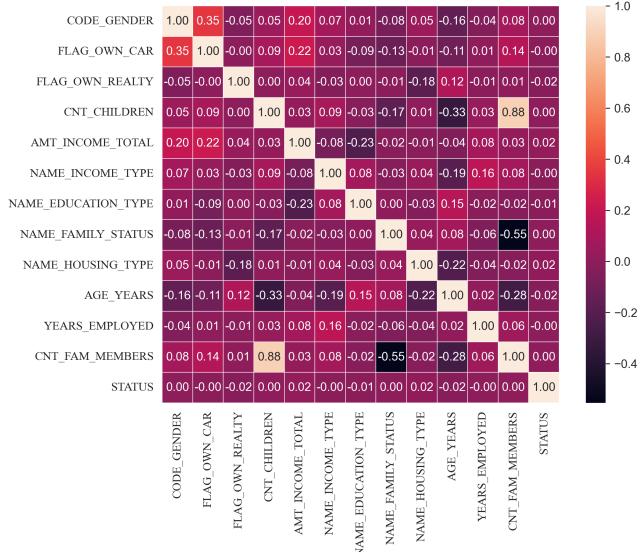
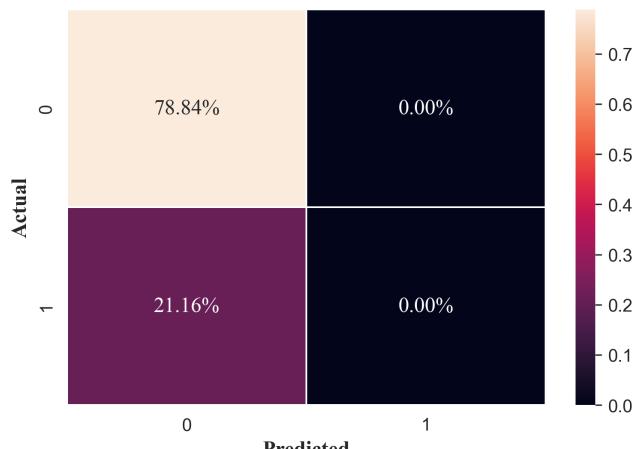
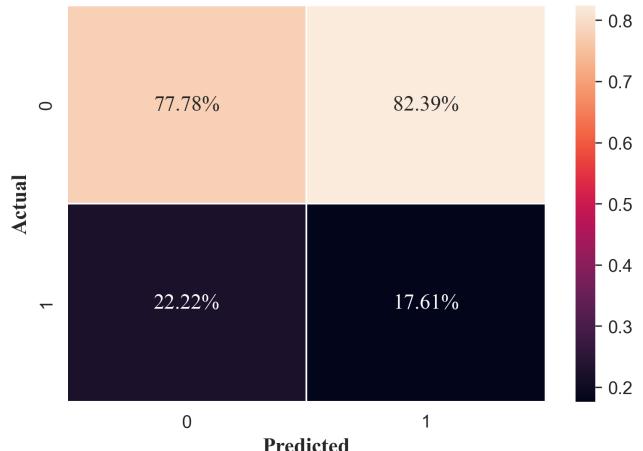
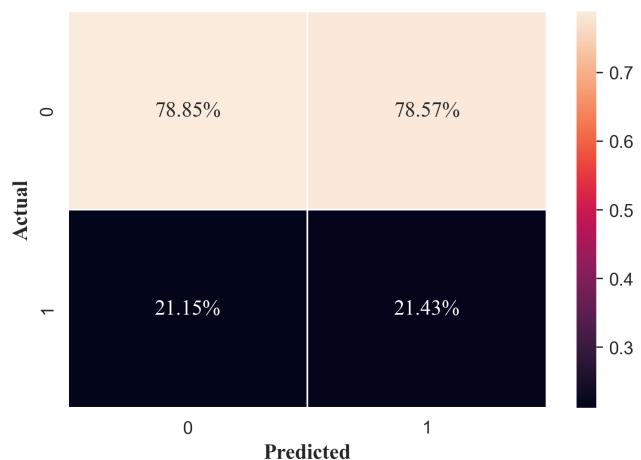
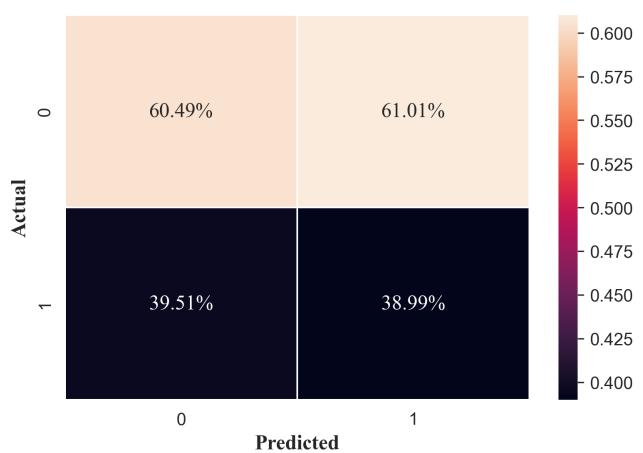


FIGURE 21. Scatter Plot: Years Employed vs Total Annual Income of Applicants

**FIGURE 22.** Correlation Heat Map between all Features and Target Label**FIGURE 24.** Confusion Matrix: Weighted KNN Model**FIGURE 25.** Confusion Matrix: KNN Model with K = 10**FIGURE 23.** Confusion Matrix: Basic KNN Model with K = 7**FIGURE 26.** Confusion Matrix: KNN Model on Balanced Dataset using SMOTE

APPENDIX A

CODE

1) Imported Modules and Functions

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from pprint import pprint
6 from tqdm import tqdm
7 from matplotlib.font_manager import
     FontProperties
8
9 from sklearn.metrics import
     classification_report, accuracy_score
     , confusion_matrix
10 from sklearn.model_selection import
     train_test_split
11
12 import matplotlib.style as style
13 import matplotlib.pyplot as plt
14 from matplotlib import colors
15 from mpl_toolkits.mplot3d import Axes3D
16 from matplotlib.colors import
     ListedColormap
17 from IPython.display import Image
18 import plotly.express as px
19 import plotly.graph_objs as go
20 import plotly.figure_factory as ff
21 from plotly.subplots import
     make_subplots
22 import plotly.offline as pyo
23 from plotly import tools
24 import seaborn as sns
25
26 import missingno as msno #to visualize
     missing data
27
28 from imblearn.over_sampling import SMOTE
29 import itertools
30
31 from sklearn.model_selection import
     train_test_split
32 from sklearn.metrics import
     accuracy_score, confusion_matrix,
     precision_score, recall_score,
     roc_auc_score, f1_score, roc_curve
33 from sklearn.neighbors import
     KNeighborsClassifier
34
35 from sklearn.preprocessing import
     LabelEncoder #label encoding for
     categorical columns
36
37 import warnings
38 warnings.filterwarnings('ignore')

```

2) Dataset Exploration, Dataset Cleaning

```

1 pd.set_option('display.max_rows', 10)
2 pd.set_option('display.max_columns', 6)
3 app_df = pd.read_csv("./
     application_record.csv")
4 app_df.head(10)
5 credit_df = pd.read_csv("./credit_record
     .csv")
6 credit_df.head(8)
7 print("Application Record Shape: ",
     app_df.shape)
8 print("Credit Record Shape: ", credit_df
     .shape)
9 app_df.info()
10 credit_df.info()
11 app_df.describe()
12 app_df.isnull().sum()
13
14 plt_missing_1 = msno.matrix(app_df)
15 plt.show()
16
17 print("TOTAL NAN in 'OCCUPATION_TYPE': "
     , len(app_df[app_df['OCCUPATION_TYPE'
     ].isna()].OCCUPATION_TYPE))
18 # dropping occupation type which has
     many null values
19 app_df.drop('OCCUPATION_TYPE', axis=1,
     inplace=True)
20 # Checking duplicates in 'ID' column
21 len(app_df['ID']) - len(app_df['ID'].unique())
22 # Dropping duplicate entries from ID
     column
23 app_df = app_df.drop_duplicates('ID',
     keep='last')
24 # Checking Non-Numerical Columns
25 cat_columns = app_df.columns[(app_df.
     dtypes == 'object').values].tolist()
26 print("Non-Numerical Columns: ")
27 cat_columns
28 # Checking Numerical Columns
29 print("Numerical Columns")
30 app_df.columns[(app_df.dtypes != 'object'
     ).values].tolist()
31 # Checking unique values from
     Categorical Columns
32 for i in app_df.columns[(app_df.dtypes
     == 'object').values].tolist():
33     print(i, '\n')
34     print(app_df[i].value_counts())
35     print('
     ')
-----)
36 app_df['CNT_CHILDREN'].value_counts()

```

```

37 print('Min DAYS_BIRTH :', app_df['
    DAYS_BIRTH'].min(), '\nMax DAYS_BIRTH
    :', app_df['DAYS_BIRTH'].max())
38
39 # Converting 'DAYS_BIRTH' values from
# Day to Years
40 app_df['DAYS_BIRTH'] = round(app_df['
    DAYS_BIRTH']/-365, 0)
41 app_df.rename(columns={'DAYS_BIRTH':'
    AGE_YEARS'}, inplace=True)
42 app_df.head(1)
43
44 # Checking unique values greater than 0
45 app_df[app_df['DAYS_EMPLOYED']>0]['
    DAYS_EMPLOYED'].unique()
46
47 # As mentioned in document, if '
# DAYS_EMPLOYED' is positive no, it
# means person currently unemployed,
# hence replacing it with 0
48 app_df['DAYS_EMPLOYED'].replace(365243,
    0, inplace=True)
49 app_df['DAYS_EMPLOYED']
50
51 # Converting 'DAYS_EMPLOYED' values from
# Day to Years
52 app_df['DAYS_EMPLOYED'] = abs(round(
    app_df['DAYS_EMPLOYED']/-365, 0))
53 app_df.rename(columns={'DAYS_EMPLOYED':'
    YEARS_EMPLOYED'}, inplace=True)
54
55 app_df['FLAG_MOBIL'].value_counts()
56
57 # As all the values in column are 1,
# hence dropping column
58 app_df.drop('FLAG_MOBIL', axis=1,
    inplace=True)
59 app_df['FLAG_WORK_PHONE'].value_counts()
60
61 # This column only contains 0 & 1 values
# for Mobile no submitted, hence
# dropping column
62 app_df.drop('FLAG_WORK_PHONE', axis=1,
    inplace=True)
63 app_df['FLAG_PHONE'].value_counts()
64
65 # This column only contains 0 & 1 values
# for Phone no submitted, hence
# dropping column
66 app_df.drop('FLAG_PHONE', axis=1,
    inplace=True)
67 app_df['FLAG_EMAIL'].value_counts()
68
69 # This column only contains 0 & 1 values
# for Email submitted, hence dropping
# column

```

```

70 app_df.drop('FLAG_EMAIL', axis=1,
    inplace=True)
71
72 app_df['CNT_FAM_MEMBERS'].value_counts()

3) Dataset Visualization
1 fig, axes = plt.subplots(3, 2, figsize
    =(12, 10))
2 font = FontProperties(family='Times New
    Roman', size = 10)
3 font_label = FontProperties(family='
    Times New Roman', size = 14, weight =
    "bold")
4 font_xticks = FontProperties(family='
    Times New Roman', size = 12)
5 font_yticks = FontProperties(family='
    Times New Roman', size = 12)
6 sns.set_style('whitegrid')
7
8 sns.boxplot(app_df['CNT_CHILDREN'],
    orient='h', ax=axes[0][0])
9 sns.boxplot(app_df['AMT_INCOME_TOTAL'],
    orient='h', ax=axes[0][1])
10 sns.boxplot(app_df['AGE_YEARS'], orient=
    'h', ax=axes[1][0])
11 sns.boxplot(app_df['YEARS_EMPLOYED'],
    orient='h', ax=axes[1][1])
12 sns.boxplot(app_df['CNT_FAM_MEMBERS'],
    orient='h', ax=axes[2][0])
13
14 axes[0][0].set_title('CNT_CHILDREN',
    font = font_label)
15 axes[0][1].set_title('AMT_INCOME_TOTAL',
    font = font_label)
16 axes[1][0].set_title('AGE_YEARS', font =
    font_label)
17 axes[1][1].set_title('YEARS_EMPLOYED',
    font = font_label)
18 axes[2][0].set_title('CNT_FAM_MEMBERS',
    font = font_label)
19
20 axes[2][1].axis('off')
21 plt.tight_layout()
22 plt.show()
23
24 high_bound = app_df['CNT_CHILDREN'].quantile(0.999)
25 print('high_bound :', high_bound)
26 low_bound = app_df['CNT_CHILDREN'].quantile(0.001)
27 print('low_bound :', low_bound)
28 app_df = app_df[(app_df['CNT_CHILDREN']
    >= low_bound) & (app_df['CNT_CHILDREN'
    ] <= high_bound)]

```

```

29 high_bound = app_df['AMT_INCOME_TOTAL'].quantile(0.999)
30 print('high_bound :', high_bound)
31 low_bound = app_df['AMT_INCOME_TOTAL'].quantile(0.001)
32 print('low_bound :', low_bound)
33 app_df = app_df[(app_df['AMT_INCOME_TOTAL'] >= low_bound) & (app_df['AMT_INCOME_TOTAL'] <= high_bound)]
34 high_bound = app_df['YEARS_EMPLOYED'].quantile(0.999)
35 print('high_bound :', high_bound)
36 low_bound = app_df['YEARS_EMPLOYED'].quantile(0.001)
37 print('low_bound :', low_bound)
38 app_df = app_df[(app_df['YEARS_EMPLOYED'] >= low_bound) & (app_df['YEARS_EMPLOYED'] <= high_bound)]
39 high_bound = app_df['CNT_FAM_MEMBERS'].quantile(0.999)
40 print('high_bound :', high_bound)
41 low_bound = app_df['CNT_FAM_MEMBERS'].quantile(0.001)
42 print('low_bound :', low_bound)
43 app_df = app_df[(app_df['CNT_FAM_MEMBERS'] >= low_bound) & (app_df['CNT_FAM_MEMBERS'] <= high_bound)]
44
45 # categorizing 'STATUS' column to binary classification 0 : Good Client and 1 : bad client
46 credit_df['STATUS'].replace(['C', 'X'], 0, inplace=True)
47 credit_df['STATUS'].replace(['2', '3', '4', '5'], 1, inplace=True)
48 credit_df['STATUS'] = credit_df['STATUS'].astype('int')
49
50 # Get the value counts as percentages
51 percentage_counts = credit_df['STATUS'].value_counts(normalize=True) * 100
52
53 # Create labels for the pie chart
54 labels = ['Good Client', 'Bad Client']
55
56 # Create a pie chart
57 plt.figure(figsize=(5, 5))
58 plt.pie(percentage_counts, explode=(0.1, 0.2), labels=labels, autopct='%1.1f%%', shadow=True, colors=[(0.24, 0.57, 0.90), (0.69, 0.05, 0.18)], textprops = {'family':'Times New Roman', 'size': 10, 'color':'black'})
59 plt.legend(loc = 'best', bbox_to_anchor = (1.15, 0.95), prop = font)
60 plt.show()
61
62 credit_df_trans = credit_df.groupby('ID').agg(max).reset_index()
63 credit_df_trans.drop('MONTHS_BALANCE', axis=1, inplace=True)
64 credit_df_trans.head()
65
66 credit_df_trans['STATUS'].value_counts(normalize=True)*100
67 # Get the value counts as percentages
68 percentage_counts = credit_df_trans['STATUS'].value_counts(normalize=True) * 100
69
70 # Create labels for the pie chart
71 labels = ['Good Client', 'Bad Client']
72
73 # Create a pie chart
74 plt.figure(figsize=(5, 5))
75 plt.pie(percentage_counts, explode=(0.1, 0.1), labels=labels, autopct='%1.1f%%', shadow=True, colors=[(0.24, 0.57, 0.90), (0.69, 0.05, 0.18)], textprops = {'family':'Times New Roman', 'size': 10, 'color':'black'})
76 plt.legend(loc = 'best', bbox_to_anchor = (1.15, 0.95), prop = font)
77 plt.show()
78
79 # merging the two datasets based on 'ID'
80 final_df = pd.merge(app_df, credit_df_trans, on='ID', how='inner')
81 final_df.head()
82
83 # dropping 'ID' column as it is having only unique values (not required for ML Model)
84 final_df.drop('ID', axis=1, inplace=True)
85
86 # checking if there are still duplicate rows in Final Dataframe
87 len(final_df) - len(final_df.drop_duplicates())
88
89 # Dropping duplicate records
90 final_df = final_df.drop_duplicates()
91 final_df.reset_index(drop=True, inplace=True)
92
93 final_df.isnull().sum()
94 final_df['STATUS'].value_counts(normalize = True)*100
95

```

```

96 percentage_counts = final_df['STATUS'].value_counts(normalize = True)*100
97
98 labels = ['Good Client', 'Bad Client']
99
100 plt.figure(figsize=(5, 5))
101 plt.pie(percentage_counts, explode = (0.05, 0.05), labels=labels, autopct='%.1f%%', shadow=True, colors =[(0.24, 0.57, 0.90), (0.69, 0.05, 0.18)], textprops = {'family':'Times New Roman', 'size': 10, 'color':'black'})
102 plt.legend(loc = 'best', bbox_to_anchor = (1.15, 0.95), prop = font)
103 plt.show()
104
105 plt.figure(figsize = (8, 6))
106 font = FontProperties(family='Times New Roman', size = 10)
107 font_label = FontProperties(family='Times New Roman', size = 14, weight = "bold")
108 font_xticks = FontProperties(family='Times New Roman', size = 10)
109 font_yticks = FontProperties(family='Times New Roman', size = 10)
110 sns.heatmap(final_df.corr(), annot=True, square = True)
111 plt.xticks(font = font_xticks)
112 plt.yticks(font = font_yticks)
113 plt.show()
114
115 # Create a 2x2 grid of subplots
116 fig, axes = plt.subplots(2, 2, figsize =(10, 6))
117
118 # Plot the first pie chart in the upper-left subplot
119 axes[0, 0].pie(final_df['CODE_GENDER'].value_counts(), labels=['Female', 'Male'], autopct='%.1f%%')
120 axes[0, 0].set_title('% of Applications \nsubmitted based on Gender')
121
122 # Plot the second pie chart in the upper-right subplot
123 axes[0, 1].pie(final_df[final_df['STATUS'] == 0]['CODE_GENDER'].value_counts(), labels=['Female', 'Male'], autopct='%.1f%%')
124 axes[0, 1].set_title('% of Applications \nApproved based on Gender')
125
126 # Plot the third pie chart in the lower-left subplot
127 axes[1, 0].pie(final_df['FLAG_OWN_CAR'].value_counts(), labels=['No', 'Yes'], autopct='%.1f%%')
128 axes[1, 0].set_title('% of Applications \nsubmitted based on owning a Car')
129
130 # Plot the fourth pie chart in the lower-right subplot
131 axes[1, 1].pie(final_df['FLAG_OWN_REALTY'].value_counts(), labels=['Yes', 'No'], autopct='%.1f%%')
132 axes[1, 1].set_title('% of Applications \nsubmitted based on owning Real Estate')
133 plt.subplots_adjust(wspace = 0.4, hspace = 10)
134
135 # Add borders between subplots
136 for ax in axes.flatten():
137     ax.spines['top'].set_visible(True)
138     ax.spines['right'].set_visible(True)
139
140 # Adjust layout
141 plt.tight_layout()
142
143 # Show the plot
144 plt.show()
145
146 # Plot the first pie chart (Gender)
147 g1 = final_df['CODE_GENDER'].value_counts().plot.pie(explode =[0.05, 0.05], labels=['Female', 'Male'], autopct='%.1f%%', shadow=True, colors=["#76B5B3", "#EC9B9A"], textprops={'fontsize': 12, 'family':'Times New Roman'})
148 # g1.set_title("Percentage of Applications \nSubmitted based on Gender", font = font_label)
149 plt.legend(loc = 'best', prop = font, bbox_to_anchor = (1.15, 0.95))
150 plt.show()
151
152 # Plot the first pie chart (Gender)
153 g2 = final_df[final_df['STATUS'] == 0]['CODE_GENDER'].value_counts().plot.pie(explode=[0.05, 0.05], labels=['Female', 'Male'], autopct='%.1f%%', shadow=True, colors=["#76B5B3", "#EC9B0B"], textprops={'fontsize': 12, 'family': 'Times New Roman'})
154 # g2.set_title("Percentage of Applications \nApproved based on Gender", font = font_label)
155 plt.legend(loc = 'best', prop = font, bbox_to_anchor = (1.15, 0.95))

```

```

156 plt.show()
157
158 # Plot the second pie chart (Car
    Ownership)
159 g3 = final_df['FLAG_OWN_CAR'].value_counts().plot.pie(explode=[0.05, 0.05], autopct='%1.2f%%', shadow=True, colors=[(0.69, 0.05, 0.18), (0.24, 0.57, 0.90)], textprops={'fontsize': 12, 'family': 'Times New Roman'})
160
161 # g3.set_title("Percentage of Applications \nSubmitted based on owning a Car", font = font_label)
162 plt.legend(loc = 'best', prop = font, bbox_to_anchor = (0.75, 0.85))
163 plt.show()
164
165 # Plot the third pie chart (Realty
    Ownership)
166 g4 = final_df['FLAG_OWN_REALTY'].value_counts().plot.pie(explode=[0.05, 0.05], autopct='%1.2f%%', shadow=True, colors=[(0.24, 0.57, 0.90), (0.69, 0.05, 0.18)], textprops={'fontsize': 12, 'family': 'Times New Roman'})
167 # g4.set_title("Percentage of Applications \nSubmitted based on owning Real Estate")
168 plt.legend(loc = 'best', prop = font, bbox_to_anchor = (0.99, 0.95))
169 plt.show()
170
171 # Create a 1x3 grid of subplots
172 fig, axes = plt.subplots(2, 2, figsize=(8, 8))
173
174
175 # Plot the first pie chart (Gender)
176 g1 = final_df['CODE_GENDER'].value_counts().plot.pie(explode=[0.1, 0.1], labels=['Female', 'Male'], autopct='%1.1f%%', shadow=True, colors=["#76B5B3", "#EC9B9A"], textprops={'fontsize': 10}, ax=axes[0][0])
177 g1.set_title("Percentage of Applications \nSubmitted based on Gender", font = font_label)
178
179
180 # Plot the first pie chart (Gender)
181 g2 = final_df[final_df['STATUS'] == 0]['CODE_GENDER'].value_counts().plot.pie(explode=[0.1, 0.1], labels=['Female', 'Male'], autopct='%1.1f%%', shadow=True, colors=["#76B5B3", "#EC9B9A"], textprops={'fontsize': 10}, ax=axes[0][1])
182 g2.set_title("Percentage of Applications \nApproved based on Gender", font = font_label)
183
184 # Plot the second pie chart (Car
    Ownership)
185 g3 = final_df['FLAG_OWN_CAR'].value_counts().plot.pie(explode=[0.1, 0.1], labels=['Female', 'Male'], autopct='%1.1f%%', shadow=True, colors=["#80DE99", "#00CECB"], textprops={'fontsize': 10}, ax=axes[1][0])
186 g3.set_title("Percentage of Applications \nSubmitted based on owning a Car", font = font_label)
187
188 # Plot the third pie chart (Realty
    Ownership)
189 g4 = final_df['FLAG_OWN_REALTY'].value_counts().plot.pie(explode=[0.1, 0.1], autopct='%1.1f%%', shadow=True, colors=["#76B5B3", "#00CECB"], textprops={'fontsize': 10}, ax=axes[1][1])
190 g4.set_title("Percentage of Applications \nSubmitted based on owning Real Estate")
191
192 # Adjust layout
193 plt.tight_layout()
194
195 # Show the plot
196 plt.show()
197
198 # Plot the first pie chart (Gender)
199 g1 = final_df['CNT_CHILDREN'].value_counts().plot.pie(explode = [0.1, 0, 0, 0.1, 0.25], labels=final_df['CNT_CHILDREN'].value_counts().index, autopct='%1.1f%%', shadow=True, colors=["#76B5B3", "#EC9B9A", "#80DE99", "#EC9B0B", "#000000"], textprops={'fontsize': 10})
200 g1.set_xlabel("% of Applications submitted based on Children count")
201 # Adjust layout
202 plt.tight_layout()
203 plt.legend(loc = 'right', bbox_to_anchor = (1.05, 0.85))
204 # Show the plot
205 plt.show()

```

```

206
207 font = FontProperties(family='Times New
   Roman', size = 10)
208 font_label = FontProperties(family='
   Times New Roman', size = 10, weight =
   "bold")
209 font_xticks = FontProperties(family='
   Times New Roman', size = 12)
210 font_yticks = FontProperties(family='
   Times New Roman', size = 12)
211
212 # Calculate the percentage distribution
   of 'CNT_CHILDREN'
213 percentage_distribution = final_df['
   CNT_CHILDREN'].value_counts(normalize
   =True) * 100
214
215 # Create a bar plot
216 plt.figure(figsize=(6, 4))
217 percentage_distribution.plot(kind='bar')
218
219 # Add labels and percentages on top of
   the bars
220 for x, y in enumerate(
   percentage_distribution):
221     plt.text(x, y + 1, f'{y:.2f}%', ha='
       center', fontsize = 10,
       fontfamily = 'Times New Roman')
222
223 # Set labels and title
224 plt.xlabel('Children Count', font =
   font_label)
225 plt.ylabel('Percentage', font =
   font_label)
226
227 # Show the plot
228 plt.xticks(rotation=0, font =
   font_xticks)
229 plt.yticks(font = font_yticks)
230 plt.tight_layout()
231 plt.savefig('./Figures/10
   _children_count_bar_plot.png', dpi =
   300, bbox_inches = 'tight')
232 plt.show()
233
234 # This graph shows that, majority of
   applicatant's income lies between 1
   to 3 lakh
235 font = FontProperties(family='Times New
   Roman', size = 10)
236 font_label = FontProperties(family='
   Times New Roman', size = 10, weight =
   "bold")
237 font_xticks = FontProperties(family='
   Times New Roman', size = 12)
238 font_yticks = FontProperties(family='
   Times New Roman', size = 12)
239 plt.figure(figsize=(6, 4))
240 sns.histplot(final_df['AMT_INCOME_TOTAL']
   ], bins=20, kde=True)
241 plt.xlabel('Total Annual Income', font =
   font_label)
242 plt.ylabel('Frequency', font =
   font_label)
243 plt.xticks(font = font_xticks)
244 plt.yticks(font = font_yticks)
245 plt.savefig("./Figures/11
   _total_annual_income_dist.png", dpi =
   300, bbox_inches = 'tight')
246 plt.show()
247
248 # Plot the first pie chart (Gender)
249 g1 = final_df['NAME_INCOME_TYPE'].value_counts().plot.pie(explode =
   (0.05, 0.03, 0, 0, 0.25), labels=
   final_df['NAME_INCOME_TYPE'].value_counts().index, autopct='%.2f
   %%', shadow=True, colors=["#76B5B3",
   "#EC9B9A", "#80DE99", "#EC9B0B", "
   #000000"], textprops={'fontsize': 10,
   'family': 'Times New Roman'})
250 # g1.set_xlabel("% of Applications
   submitted based on Income Type")
251 # Adjust layout
252 plt.tight_layout()
253 plt.legend(loc = 'right', bbox_to_anchor
   = (1.26, 0.85), prop = font)
254 plt.savefig('./Figures/12_income_types.
   png', dpi = 300, bbox_inches = 'tight
   ')
255 plt.show()
256
257 font = FontProperties(family='Times New
   Roman', size = 10)
258 font_label = FontProperties(family='
   Times New Roman', size = 10, weight =
   "bold")
259 font_xticks = FontProperties(family='
   Times New Roman', size = 10)
260 font_yticks = FontProperties(family='
   Times New Roman', size = 10)
261
262 # Calculate the counts of each education
   type
263 education_counts = final_df['
   NAME_EDUCATION_TYPE'].value_counts()
264
265 # Create a bar graph
266 plt.figure(figsize=(7, 5)) # Adjust the
   figure size as needed

```

```

267 plt.bar(education_counts.index,
           education_counts.values, color=["#76B5B3", "#EC9B9A", "#80DE99", "#EC9B0B", "#000000"])
268
269 # Add labels and title
270 plt.xlabel("Education Type", font = font_label)
271 plt.ylabel("Count", font = font_label)
272
273 # Rotate x-axis labels if needed for better readability
274 plt.xticks(rotation = 45, ha='right', font = font_xticks)
275 plt.yticks(font = font_yticks)
276
277 # Show the plot
278 plt.tight_layout()
279 plt.savefig('./Figures/13_education_types.png', dpi = 300, bbox_inches = 'tight')
280 plt.show()
281
282 # Plot the first pie chart (Gender)
283 g1 = final_df['NAME_EDUCATION_TYPE'].value_counts().plot.pie(explode = (0.1, 0, 0.25, 0, 0.35), labels = final_df['NAME_EDUCATION_TYPE'].value_counts().index, autopct='%1.1f %%', shadow=True, colors=[ "#76B5B3", "#EC9B9A", "#80DE99", "#EC9B0B", "#000000"], textprops={'fontsize': 10, 'family': 'Times New Roman'})
284 # g1.set_xlabel("% of Applications submitted based on Education")
285 # Adjust layout
286 plt.tight_layout()
287 plt.legend(loc = 'right', bbox_to_anchor = (1.34, 0.85), prop = font)
288 # Show the plot
289 plt.show()
290
291 # This graph shows that, majority of applicatant's are married
292 # Sample data
293 family_status_counts = final_df['NAME_FAMILY_STATUS'].value_counts()
294
295 sns.set_style('whitegrid')
296
297 # Create a barplot using Seaborn
298 plt.figure(figsize=(8, 5))
299 sns.barplot(x=family_status_counts.index, y=family_status_counts.values)
300 plt.xticks(rotation = 0) # Rotate x-axis labels for better visibility if needed
301 plt.xlabel('Family Status', font = font_label)
302 plt.ylabel('Count', font = font_label)
303 # plt.title('Applications submitted based on Family Status')
304 plt.xticks(font = font_xticks)
305 plt.yticks(font = font_yticks)
306 plt.savefig('./Figures/14_family_status.png', dpi = 300, bbox_inches = 'tight')
307 plt.show()
308
309 # This graph shows that, majority of applicatant's lives in House / Apartment
310 family_status_counts = final_df['NAME_HOUSING_TYPE'].value_counts()
311
312 # Create a barplot using Seaborn
313 plt.figure(figsize=(6, 4))
314 sns.barplot(x=family_status_counts.index, y=family_status_counts.values)
315 plt.xticks(rotation = 45) # Rotate x-axis labels for better visibility if needed
316 plt.xlabel('Housing Types', font = font_label)
317 plt.ylabel('Count', font = font_label)
318 plt.xticks(font = font_xticks)
319 plt.yticks(font = font_yticks)
320 plt.savefig('./Figures/15_housing_types.png', dpi = 300, bbox_inches = 'tight')
321 plt.show()
322
323 plt.figure(figsize=(6, 4))
324 sns.histplot(final_df['AGE_YEARS'], bins = 20, kde=True)
325 plt.xlabel('Age', font = font_label)
326 plt.ylabel("Frequency", font = font_label)
327 plt.xticks(font = font_xticks)
328 plt.yticks(font = font_yticks)
329 plt.savefig('./Figures/15_age_distribution.png', dpi = 300, bbox_inches = 'tight')
330 plt.show()
331
332 # This graph shows that, majority of applicatant's are Employed for 0 to 7 years
333 plt.figure(figsize=(6, 4))
334 sns.histplot(final_df['YEARS_EMPLOYED'], bins=20, kde=True)

```

```

335 plt.xlabel('No of Years Employed', font
            = font_label)
336 plt.ylabel("Frequency", font =
            font_label)
337 plt.xticks(font = font_xticks)
338 plt.yticks(font = font_yticks)
339 plt.savefig('./Figures/16_years_employed
            .png', dpi = 300, bbox_inches =
            'tight')
340 plt.show()
341
342 # This graph shows that, majority of
    applications are rejected if Total
    income & years of Employment is less
343 plt.figure(figsize=(8, 6))
344 sns.scatterplot(data=final_df, x='
    YEARS_EMPLOYED', y='AMT_INCOME_TOTAL'
    , hue='STATUS')
345 plt.xlabel('No of Years Employed', font
            = font_label)
346 plt.ylabel("Total Income Amount", font =
            font_label)
347 plt.xticks(font = font_xticks)
348 plt.yticks(font = font_yticks)
349 plt.savefig('./Figures/17
            _scatter_plot_employed_vs_income.png'
    , dpi = 300, bbox_inches = 'tight')
350 plt.show()

```

4) Feature Selection

```

1 # Taking only the categorical Columns
2 cat_columns = final_df.columns[(final_df
        .dtypes =='object').values].tolist()
3 cat_columns
4
5 #Converting all Non-Numerical Columns to
    Numerical
6 from sklearn.preprocessing import
    LabelEncoder
7
8 for col in cat_columns:
9     globals()['LE_{}'.format(col)] =
        LabelEncoder()
10    final_df[col] = globals()['LE_{}'.format(
        col)].fit_transform(
        final_df[col])
11 final_df.head()
12
13 for col in cat_columns:
14     print(col , " : ", globals()['LE_{}'.format(
        col)].classes_)
15
16 correlation_matrix = final_df.corr()
17
18 plt.figure(figsize=(8, 6))

```

```

19 plt.xticks(font = font_xticks)
20 plt.yticks(font = font_yticks)
21 sns.heatmap(correlation_matrix, annot=
            True, fmt=".2f", linewidths = 0.01,
            square = True)
22 plt.savefig("./Figures/18
            _Correlation_all_heat_map.png", dpi =
            300, bbox_inches = 'tight')
23 plt.show()

```

5) Data Modeling: KNN Models

```

1 features = final_df.drop(['STATUS'],
            axis=1)
2 label = final_df['STATUS']
3
4 from sklearn.model_selection import
    train_test_split
5 x_train, x_test, y_train, y_test =
    train_test_split(features, label,
    test_size=0.2, random_state = 10)
6
7 # K Nearest Neighbor classification
8 from sklearn.neighbors import
    KNeighborsClassifier
9
10 knn_model = KNeighborsClassifier(
    n_neighbors = 7)
11
12 knn_model.fit(x_train, y_train)
13
14 print('KNN Model Accuracy : ', knn_model
    .score(x_test, y_test)*100, '%')
15
16 prediction = knn_model.predict(x_test)
17 print('\nConfusion matrix :')
18 print(confusion_matrix(y_test,
    prediction))
19
20 print('\nClassification report :')
21 print(classification_report(y_test,
    prediction))
22
23 font = FontProperties(family='Times New
    Roman', size = 12, weight = 'bold')
24 plt.figure(figsize=(6, 4))
25 sns.heatmap(confusion_matrix(y_test,
    prediction, normalize = 'pred'), fmt
    = '.2%', annot=True, linewidths=0.5,
    annot_kws={"fontsize": 12, 'family':
    'Times New Roman'})
26
27 # Add labels and title
28 plt.xlabel('Predicted', font = font)
29 plt.ylabel('Actual', font = font)

```

```

30 plt.savefig("./Figures/19_CM_simple_KNN.png", dpi = 300, bbox_inches = 'tight')
31 plt.show()
32
33 # scaling all features
34 from sklearn.preprocessing import
    MinMaxScaler
35 MMS = MinMaxScaler()
36 x_train_scaled = pd.DataFrame(MMS.
    fit_transform(x_train), columns=
    x_train.columns)
37 x_test_scaled = pd.DataFrame(MMS.
    transform(x_test), columns=x_test.
    columns)
38
39 # adding samples to minority class using
    SMOTE
40 from imblearn.over_sampling import SMOTE
41 oversample = SMOTE(sampling_strategy =
    0.65, random_state = 69)
42
43 x_train_oversam, y_train_oversam =
    oversample.fit_resample(x_train,
    y_train)
44 x_test_oversam, y_test_oversam =
    oversample.fit_resample(x_test,
    y_test)
45
46 # after using SMOTE
47 y_train_oversam.value_counts(normalize =
    True) * 100
48
49 # Original majority and minority class
50 y_train.value_counts(normalize = True) *
    100
51
52 # K Nearest Neighbor classification
53 knn_model = KNeighborsClassifier(
    n_neighbors = 10)
54 knn_model.fit(x_train_oversam,
    y_train_oversam)
55
56 print('KNN Model Accuracy : ', knn_model
    .score(x_test_oversam, y_test_oversam
    )*100, '%')
57
58 prediction = knn_model.predict(
    x_test_oversam)
59 print('\nConfusion matrix :')
60 print(confusion_matrix(y_test_oversam,
    prediction))
61
62 print('\nClassification report :')
63 print(classification_report(
    y_test_oversam, prediction))
64
65 font = FontProperties(family='Times New
    Roman', size = 12, weight = 'bold')
66 plt.figure(figsize=(6, 4))
67 sns.heatmap(confusion_matrix(
        y_test_oversam, prediction, normalize
        = 'pred'), fmt = '.2%', annot=True,
        linewidths=0.5, annot_kws={"fontsize":
        12, 'family':'Times New Roman'})
68
69 # Add labels and title
70 plt.xlabel('Predicted', font = font)
71 plt.ylabel('Actual', font = font)
72 plt.savefig("./Figures/21
    _CM_SMOTE_balanced.png", dpi = 300,
    bbox_inches = 'tight')
73 plt.show()
74
75 import random
76 weights_list = []
77 for _ in range(12):
78     weights_list.append(random.random())
79 len(weights_list)
80 print(weights_list)
81
82 clf_KNN = KNeighborsClassifier(
    metric_params = {'w': weights_list})
83 clf_KNN.fit(x_train, y_train)
84
85 print('KNN Model Accuracy : ', clf_KNN.
    score(x_test, y_test)*100, '%')
86
87 prediction = knn_model.predict(x_test)
88 print('\nConfusion matrix :')
89 print(confusion_matrix(y_test,
    prediction))
90
91 print('\nClassification report :')
92 print(classification_report(y_test,
    prediction))
93
94 font = FontProperties(family='Times New
    Roman', size = 12, weight = 'bold')
95 plt.figure(figsize=(6, 4))
96 sns.heatmap(confusion_matrix(y_test,
    prediction, normalize = 'pred'), fmt
    = '.2%', annot=True, linewidths=0.5,
    annot_kws={"fontsize": 12, 'family':
    'Times New Roman'})
97
98 # Add labels and title
99 plt.xlabel('Predicted', font = font)
100 plt.ylabel('Actual', font = font)
101 plt.savefig("./Figures/20
    _CM_weighted_KNN.png", dpi = 300,
    bbox_inches = 'tight')

```

```
102 plt.show()
103
104 features = ['CODE_GENDER', 'FLAG_OWN_CAR',
105     'FLAG_OWN_REALTY', 'CNT_CHILDREN',
106     'AMT_INCOME_TOTAL', 'NAME_INCOME_TYPE',
107     'NAME_EDUCATION_TYPE',
108     'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE',
109     'AGE_YEARS',
110     'YEARS_EMPLOYED',
111     'CNT_FAM_MEMBERS']
112
113 weighted = [0.2775000393397501,
114     0.34514382952446754,
115     0.03019638152758386,
116     0.5569195160369883, 4,
117     0.9607734156146648,
118     0.08218285220614585,
119     0.0731939869231566, 4,
120     0.09920635007728695,
121     0.2678344226623953,
122     0.774548695366705]
123
124 clf_KNN = KNeighborsClassifier(
125     metric_params = {'w': weighted})
126 clf_KNN.fit(x_train, y_train)
127
128 print('KNN Model Accuracy : ', clf_KNN.
129     score(x_test, y_test)*100, '%')
130
131 prediction = knn_model.predict(x_test)
132 print('\nConfusion matrix :')
133 print(confusion_matrix(y_test,
134     prediction))
135
136 # Model accuracy score on Train set
137 # : ", model_entropy.best_score_)
138 y_pred = model_entropy.predict(x_test)
139 print(f"Accuracy after GridSearch &
140     Stratified K-fold: ", accuracy_score(
141     y_true = y_test, y_pred = y_pred))
142
143 clf_KNN = KNeighborsClassifier(
144     n_neighbors = 10)
145 clf_KNN.fit(x_train, y_train)
146
147 print('KNN Model Accuracy : ', clf_KNN.
148     score(x_test, y_test)*100, '%')
149
150 prediction = knn_model.predict(x_test)
151
152 sns.heatmap(confusion_matrix(y_test,
153     prediction, normalize = 'pred'), fmt =
154     '.2%', annot=True, linewidths=0.5,
155     annot_kws={"fontsize": 12, 'family':
156     'Times New Roman'})
157
158 # Add labels and title
159 plt.xlabel('Predicted', font = font)
160 plt.ylabel('Actual', font = font)
161 plt.savefig("./Figures/21_CM_best_KNN.
162     png", dpi = 300, bbox_inches = 'tight')
163
164 plt.show()
```