

1 **Business Analytics Approach to Credit Card Fraud Detection**

2

3 Group 5

4 Chun-Yen Lin 989467253

5 Hsin-Yu Liao 989470611

6 Vishnu Vaibhav Binde 989486228

7 Miteshkumar Dasharathbhai Patel

8 Master Science of Business Analytics, University of the Pacific

9 MSBA 250 Applied Business Analytics

10 Dr. Xun Xu

1 **Business Analytics Approach to Credit Card Fraud Detection**

2 **Abstract**

3 This project investigates the pressing issue of credit card fraud, which poses a significant threat to
4 the security and reliability of financial transactions in the digital era. Utilizing a publicly available
5 dataset of anonymized credit card transactions, the study aims to identify patterns and behavioral
6 indicators commonly associated with fraudulent activity. The analytical process involves
7 comprehensive data preprocessing using Python, followed by the implementation of machine
8 learning models to improve the accuracy of fraud prediction. In parallel, Tableau is employed for
9 interactive data visualization, facilitating intuitive exploration of risk factors across variables such
10 as transaction time, category, and user demographics. The final outcome includes a dynamic fraud
11 monitoring dashboard designed to provide real-time insights and support informed decision-
12 making. Overall, the project demonstrates the effective application of data analytics in combating
13 fraud and strengthening the resilience of financial systems.

14 ***Keywords:*** Credit Card Fraud, Business Analytics, Fraud Detection, Machine Learning,
15 Tableau Visualization

16 **1. Introduction / Background**

17 In an age where digital transactions have become the norm, credit card fraud has emerged as a
18 significant global threat, impacting not just individual consumers but financial institutions and
19 merchants as well. The sophistication of modern fraud schemes has outpaced traditional detection
20 methods, demanding the implementation of advanced analytical tools.

21 According to the Nilson Report (2025), global losses due to credit card fraud have soared to
22 over \$33 billion, further emphasizing the need for data-driven preventive measures. This project,
23 titled Credit Card Fraud Analysis, addresses these challenges by applying business analytics to a

1 dataset sourced from Kaggle. The dataset contains over 55,000 transaction records, complete with
2 customer demographic information, transaction timestamps, locations, and fraud indicators.

3 This study focuses not only on identifying potential fraud patterns, but also on building a
4 comprehensive understanding of customer risk profiles and merchant behaviors. The structure of
5 the report follows a logical progression from background to objectives, data methods, analytics,
6 conclusions, and practical recommendations, offering a complete journey from problem framing
7 to solution implementation.

8 **2. Project Objective**

9 The primary objective of this study is to examine the characteristics of fraudulent credit card
10 transactions and to understand how they differ from legitimate transactions in terms of timing,
11 location, transaction amount, and user behavior. The aim is to identify patterns that allow financial
12 institutions to proactively detect fraud while minimizing false positives that could disrupt genuine
13 customer experiences.

14 To support this objective, the study analyzes trends in customer demographics, including age,
15 gender, and other relevant attributes, and explores their correlation with fraud frequency.
16 Merchant-related variables are also examined to identify categories and regions associated with
17 elevated fraud risk. Furthermore, the project introduces derived features—such as transaction hour
18 and customer age—to uncover behavioral patterns not immediately visible in the raw dataset.

19 One of the key deliverables of this study is an interactive Tableau dashboard that facilitates
20 real-time monitoring of fraud trends, with the ability to assess risk by demographic segment and
21 geographic region. Rather than focusing solely on technical modeling, this study emphasizes the
22 development of tools that provide clear and actionable insights to support operational decision-
23 making. Overall, the project illustrates how business analytics can be effectively applied to real-

1 world challenges, contributing to both improved fraud prevention strategies and enhanced
2 customer confidence.

3 **3. Data / Problem Analytics**

4 **3.1 Data**

5 The dataset used in this study is sourced from the "Credit Card Fraud Prediction" project on
6 the Kaggle platform, comprising approximately 550,000 transaction records. Each transaction
7 includes multiple fields such as transaction time, merchant name, transaction category and amount,
8 as well as basic consumer information including gender, city and state of residence, occupation,
9 and date of birth. The `is_fraud` field indicates whether a given transaction is fraudulent. Among
10 the entire dataset, only about 2,200 transactions are labeled as fraud, accounting for roughly 0.4%
11 of the total sample. This indicates a highly imbalanced classification problem, which poses
12 challenges for model training and prediction accuracy. Therefore, this study also implements
13 corresponding techniques to address the class imbalance issue.

14 **3.2 Methods**

15 Before training the model, this study first conducted data cleaning and transformation using
16 Python to enhance the accuracy of subsequent analysis and modeling. First, the original birthdate
17 field was converted into age, which was used as a model feature to facilitate the exploration of the
18 relationship between different age groups and fraudulent behavior. Second, the transaction time
19 field was parsed to extract the hour of transaction, allowing for analysis of whether fraudulent
20 transactions tend to occur during specific time periods. In addition, the dataset underwent column
21 selection, missing value handling, and format standardization to ensure data consistency and
22 usability, thereby laying a solid foundation for model development.

23 In the feature engineering phase, categorical variables such as gender, transaction type, state,
24 and time period were transformed using LabelEncoder to convert text data into numerical formats

1 acceptable to the model. Numerical variables such as transaction amount, city population, and age
2 were standardized using MinMaxScaler, which scaled all values between 0 and 1 to prevent bias
3 caused by differences in scale. Given that only about 0.4 percent of the original dataset consisted
4 of fraudulent transactions, resulting in a highly imbalanced class distribution, the study applied the
5 Synthetic Minority Over-sampling Technique (SMOTE) to oversample the minority class in the
6 training data. This approach improved the model's ability to learn from fraudulent patterns and
7 enhanced its overall classification performance.

8 This study further developed and compared two supervised classification models: Logistic
9 Regression and Random Forest. Logistic Regression offers strong interpretability, allowing for an
10 understanding of the direction and magnitude of each variable's influence on the likelihood of
11 fraud. In contrast, Random Forest leverages an ensemble of decision trees, providing high
12 tolerance for nonlinear relationships and outliers. After training the models, their performance was
13 evaluated using the test dataset. Evaluation metrics included Accuracy, Precision, Recall, F1-
14 score, and the Confusion Matrix.

15 Finally, the study includes a visual analysis of the model prediction results. Through the use
16 of probability distribution plots and model probability comparison charts, such as scatter plots and
17 ROC curves, the differences and classification stability between the two models were examined.
18 These visualizations provide valuable insights that support subsequent model selection and
19 optimization decisions.

1 **3.3 Data / Problem Analytics**

2 To analyze the underlying trends and characteristics of fraudulent behavior, this study utilized
 3 Tableau for visual analysis to help understand the relationships between various variables and
 4 fraudulent transactions.



5 **Figure 1. Comparison of Fraud Transaction Count and Total Loss Amount by Gender**

6 First, an analysis was conducted based on gender, with two bar charts created: the left chart
 7 shows the number of fraudulent transactions (Count) by gender, while the right chart displays the
 8 total amount of fraud (Total Amount) for each gender.

9 From the figure 1, it can be observed that the number of fraudulent transactions involving
 10 female customers is 1,164, which is higher than the 981 cases involving male customers. However,
 11 in terms of total loss amount, male customers incurred a slightly higher total loss of 575,865
 12 compared to 557,459 for female customers. This result suggests that while fraudulent transactions
 13 occur more frequently among females, the loss per transaction may be higher for males.

14 Next, an analysis was conducted based on transaction category. A bar chart was created to
 15 display the number of fraudulent transactions (Count of Id) across different categories.

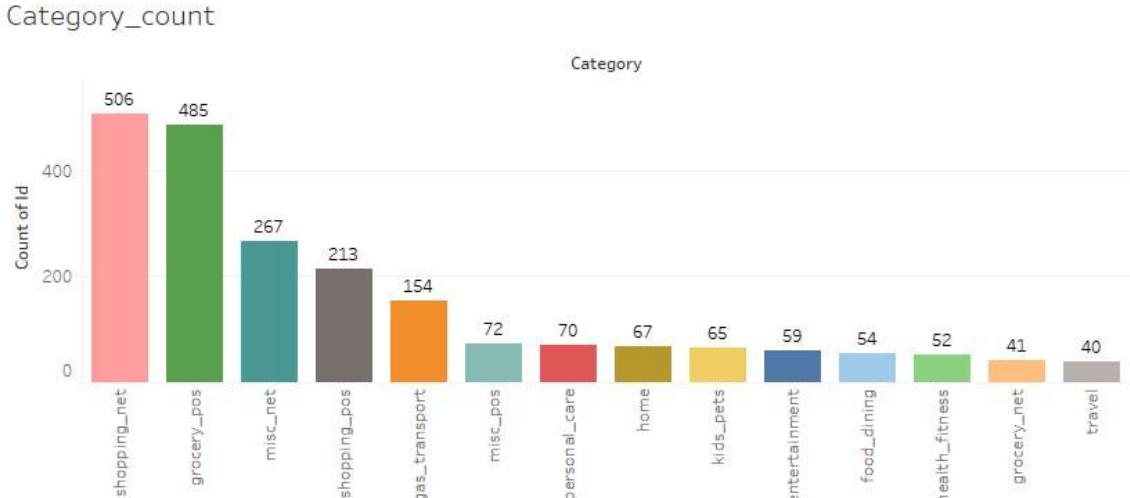


Figure 2. Number of Fraud Transactions by Category

From figure 2, it can be observed that fraudulent transactions occur most frequently in the shopping_net (online shopping) category, with a total of 506 cases. This is followed by grocery_pos (in-store grocery purchases) with 485 cases, and misc_net (other online transactions) with 267 cases. Other notable categories include shopping_pos (213 cases) and gas_transport (154 cases), which also show significant instances of fraud. These results suggest that transaction types related to everyday spending, especially online shopping, are primary targets for fraudsters. The high risk associated with online transactions may be linked to factors such as anonymity and automated payment processes. In-store transactions, such as those in grocery stores, may present opportunities for fraud due to their high frequency. Overall, while fraud is distributed across various categories, the top few with higher concentrations warrant special attention and enhanced security measures.

Next, an analysis was conducted based on transaction location (State), using both a map and a bar chart to present the distribution of fraudulent transaction amounts and counts. The map at the top visualizes the total fraud amount in each state using color intensity, while the bar chart below displays the number of fraudulent transactions (Count of Id) per state.

Maps

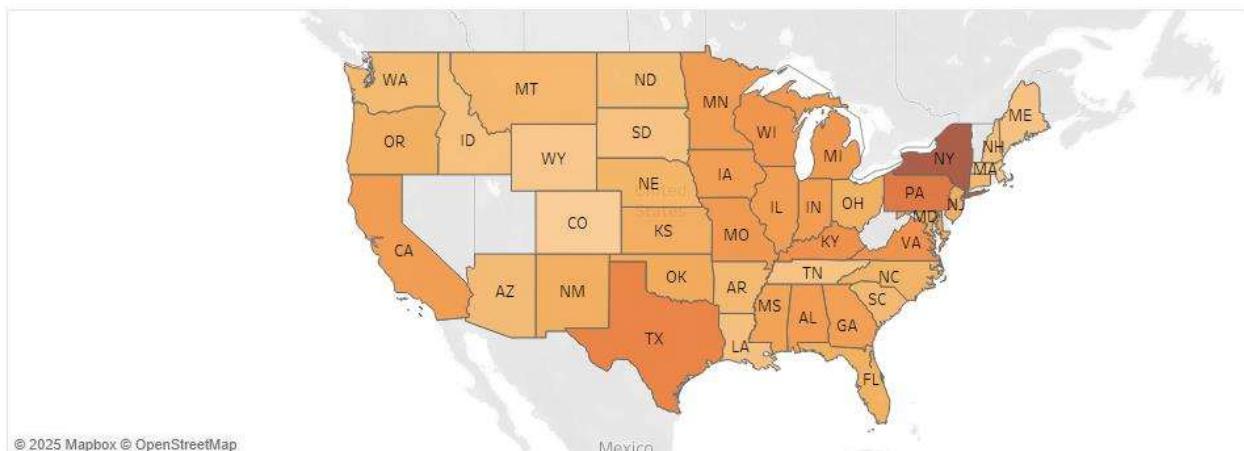


Figure 3 Geographic Distribution of Fraud Losses by U.S. State

State_count

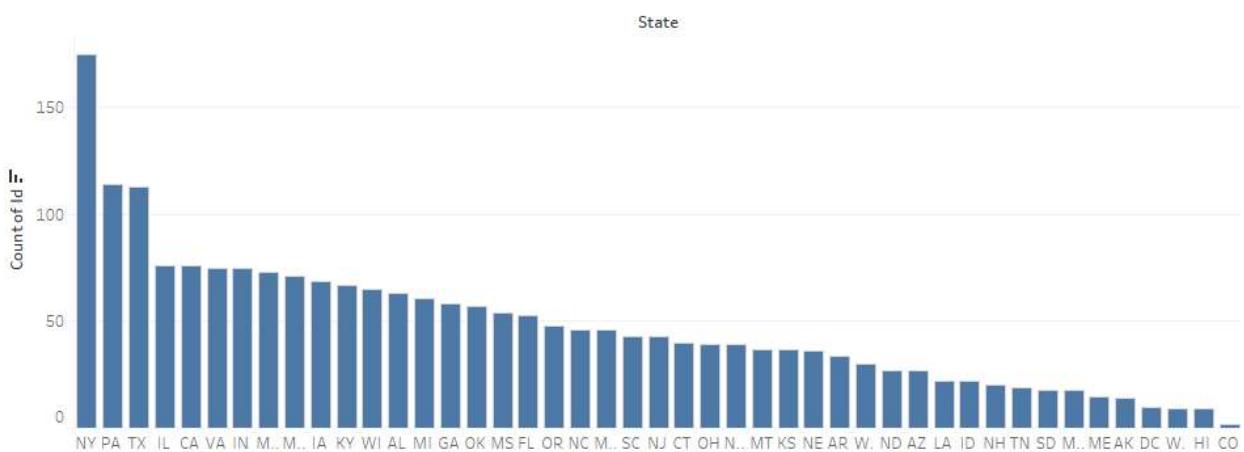


Figure 4 Number of Fraud Cases by U.S. State

From figure 3, it is evident that New York (NY) has the highest number of fraudulent

transactions, totaling 180 cases, followed by Pennsylvania (PA) and Texas (TX), each with over

5 100 cases. These states are among the most populous in the U.S., and the higher volume of

6 transactions may contribute to the increased number of fraud cases. Figure 4 also shows that New

7 York, Pennsylvania, and Texas have relatively high total fraud amounts, indicating that these states

8 not only experience frequent fraud but also significant financial losses. These findings highlight

9 the geographic distribution of fraud cases across the United States, which appears to be related to

- 1 factors such as population density, consumer behavior, and the prevalence of online transactions.
 2 This information serves as a valuable reference for risk assessment and fraud prevention strategies.

Age

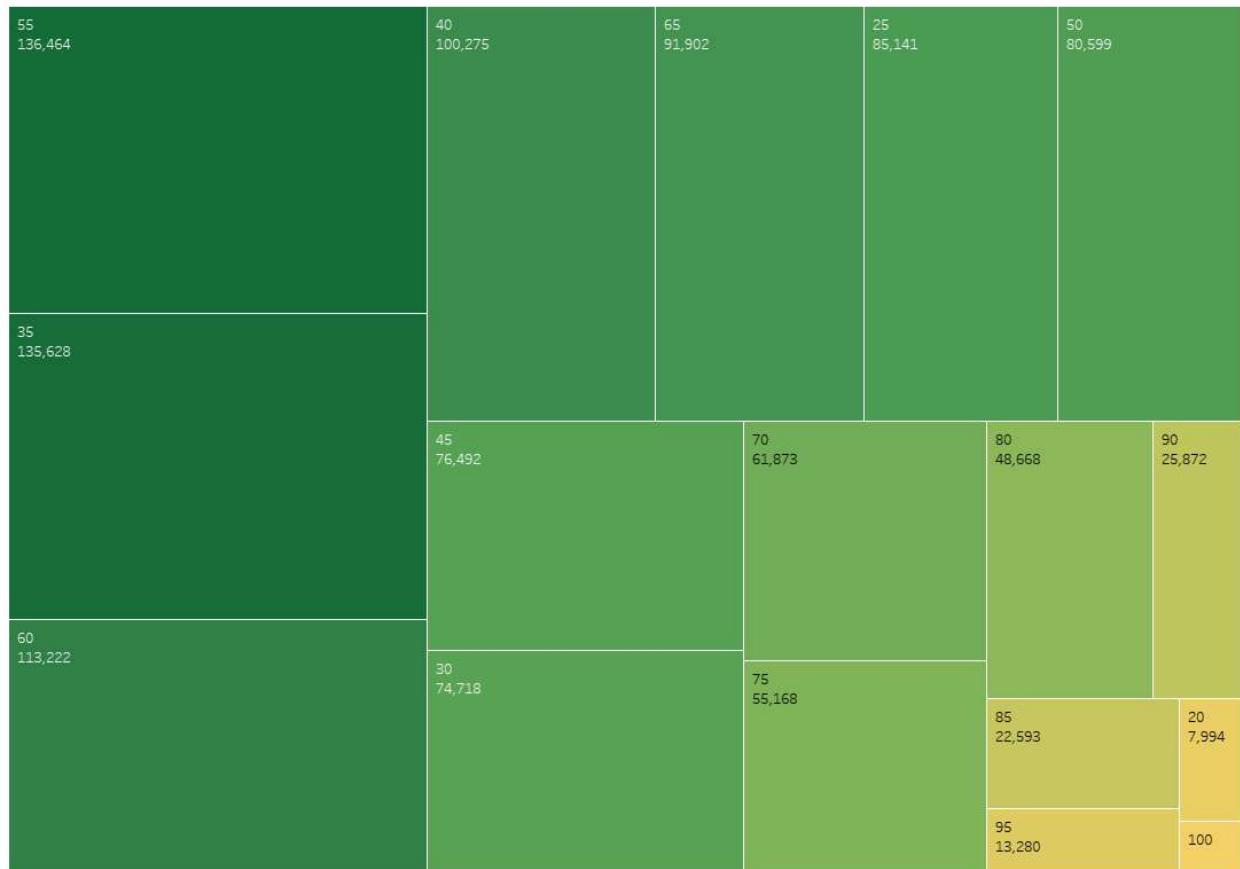
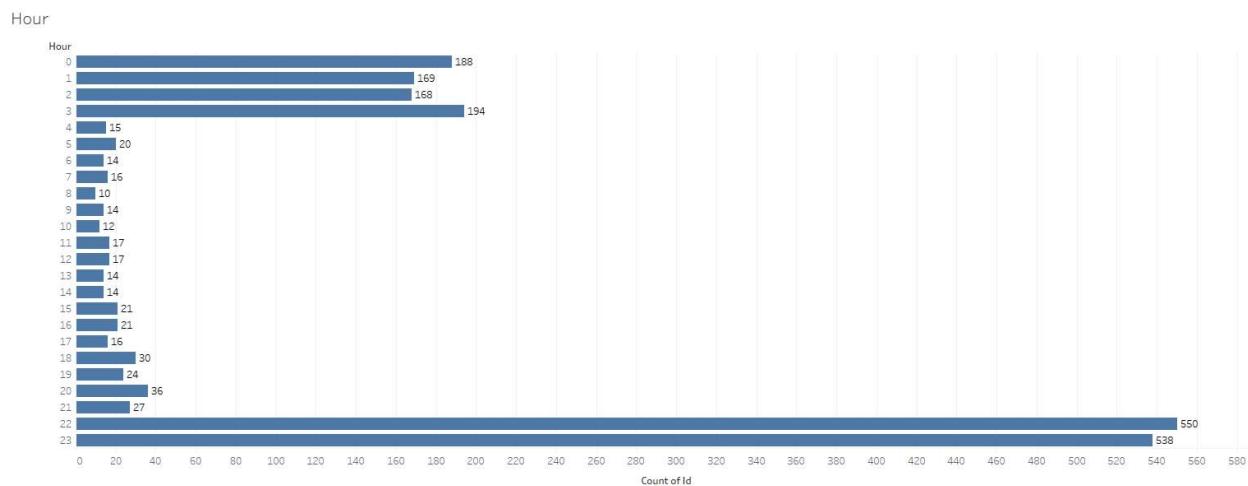


Figure 5 Treemap of Fraud Losses by Customer Age

- 3
 4 Next, an analysis was conducted based on age, with a treemap created to display the total fraud
 5 amount (Total Amount) accumulated by different age groups in fraudulent transactions. Figure 5
 6 shows that individuals aged 55 experienced the highest total fraud loss, reaching 136,464, followed
 7 by those aged 35 (135,628) and 60 (113,222), indicating that these three age groups account for
 8 the largest portions of fraudulent transaction amounts. In addition, the age groups of 40, 65, and
 9 25 also recorded losses exceeding 85,000, marking them as high-risk segments. In contrast,
 10 younger individuals (such as those aged 20) and older adults (such as those aged 85, 90, and 95)
 11 showed relatively lower fraud amounts, ranging between 10,000 and 25,000. These results suggest
 12 that middle-aged to early senior individuals (approximately ages 35 to 65) not only appear more

1 frequently in fraud cases but also incur higher average losses. This may be attributed to their higher
 2 level of economic activity and greater asset holding. Therefore, enhanced financial security
 3 education and fraud prevention efforts should be prioritized for this demographic group.

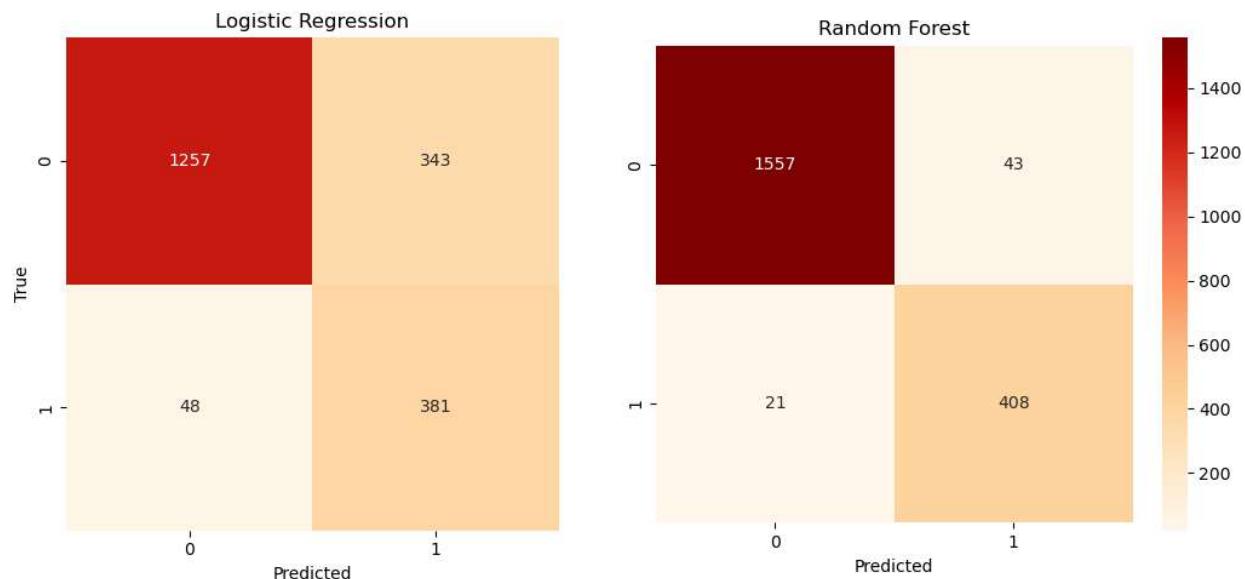
4 Finally, an analysis was conducted based on transaction time (Hour), with a bar chart created
 5 to display the number of fraudulent transactions (Count of Id) across different time periods within
 6 a 24-hour span.



7 **Figure 6 Distribution of Fraudulent Transactions by Hour**

8 Figure 6 clearly shows that fraudulent transactions occur most frequently at 10 PM (550 cases)
 9 and 11 PM (538 cases), followed by 3 AM (194 cases), 12 AM (188 cases), and 1 AM (169 cases),
 10 all of which fall within the late-night to early-morning hours. Overall, there is a clear trend of fraud
 11 being concentrated during night-time and early-morning periods, while daytime transactions (from
 12 6 AM to 5 PM) are significantly lower, typically ranging between 10 and 30 cases per hour. This
 13 pattern may be related to reduced user alertness at night, fraudsters intentionally avoiding peak
 14 monitoring hours, or increased mobile device usage by victims during nighttime. These findings
 15 suggest that the time of transaction is a critical factor in fraud risk assessment, and that enhanced
 16 real-time detection mechanisms and transaction risk controls are especially necessary during late-
 17 night hours.

1 In addition to conducting preliminary exploratory data analysis using Tableau, this study also
 2 developed fraud prediction models using two classification algorithms: Logistic Regression and
 3 Random Forest. The performance of these models was compared using various evaluation metrics
 4 and visualizations.



5 **Figure 7. Confusion Matrix Comparison between Logistic Regression and Random Forest**

6 First, the confusion matrix (Figure 7) provides insight into the models' performance in actual
 7 predictions. For instance, with the Random Forest model, out of 2,029 test samples, the model
 8 successfully identified 408 fraudulent cases (True Positives) and missed only 21 (False Negatives),
 9 demonstrating high accuracy. In contrast, while the Logistic Regression model also achieved a
 10 respectable recall of 0.89, it exhibited a higher false positive rate, misclassifying 343 legitimate
 11 transactions as fraudulent (False Positives), which could lead to excessive disruptions.

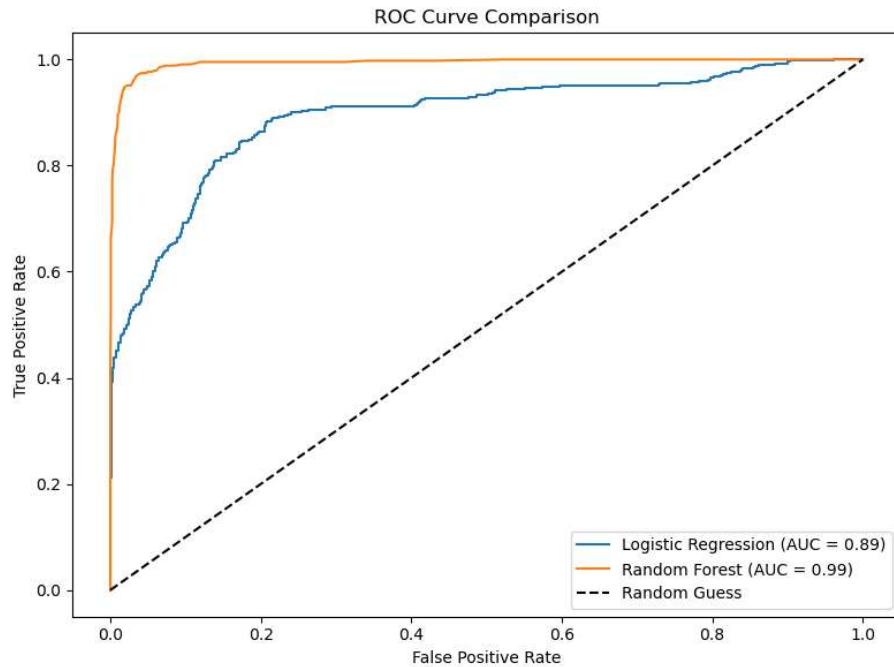


Figure 8. ROC Curve Comparison Between Logistic Regression and Random Forest

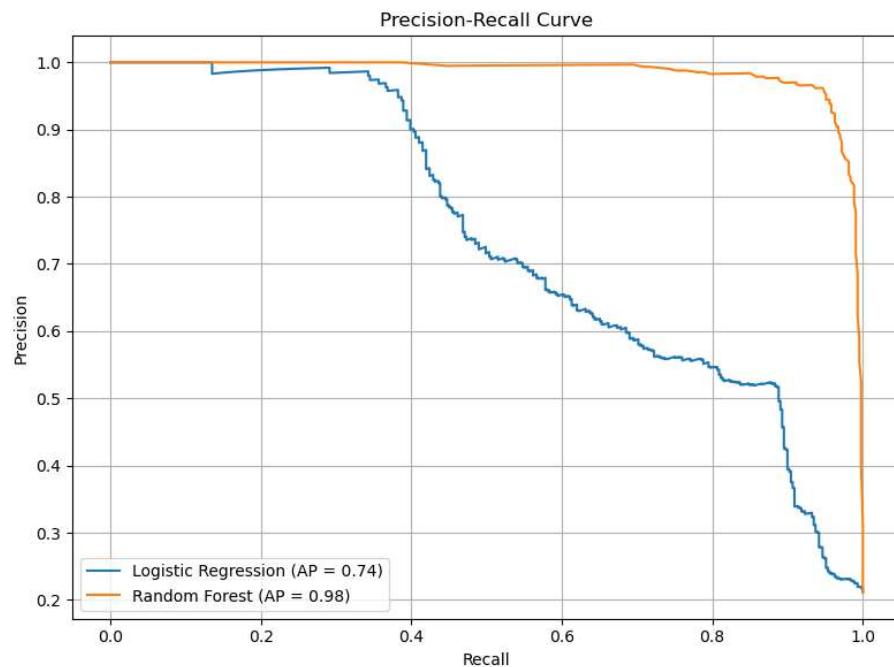


Figure 9. Precision-Recall Curve Comparison Between Logistic Regression and Random Forest

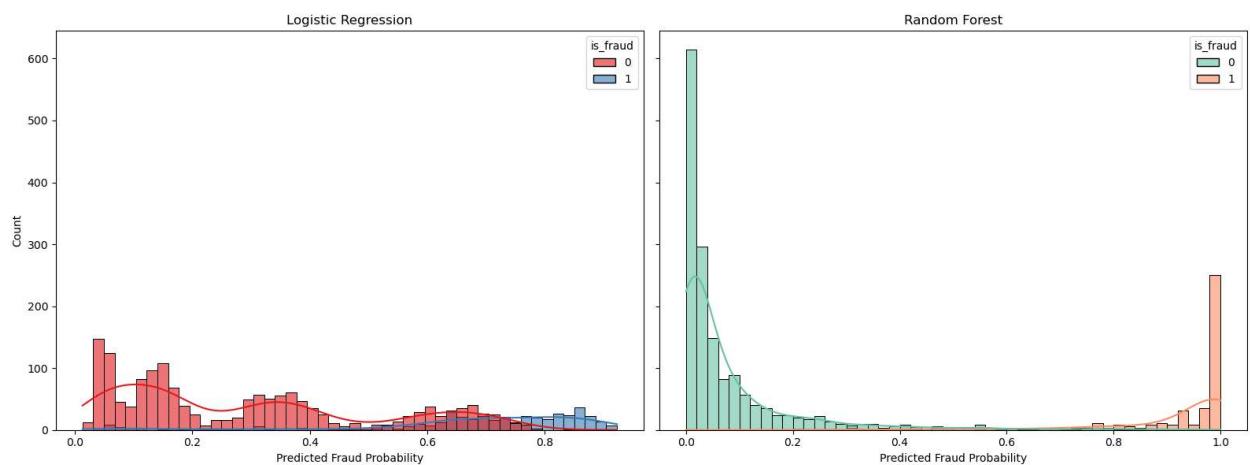
Further analysis using the ROC curve (Figure 8) and Precision-Recall curve (Figure 9)

reveals that the Random Forest model outperforms Logistic Regression in overall performance.

The Random Forest achieved an AUC of 0.99, significantly higher than Logistic Regression's

1 0.89. In terms of average precision (AP), Random Forest also demonstrated superior results,
 2 scoring 0.98 compared to 0.74 for Logistic Regression. These findings indicate that Random Forest
 3 is more effective in detecting fraudulent transactions within highly imbalanced datasets.

4 In addition, a similar trend was observed in the Precision-Recall curve. The Random Forest
 5 model achieved an average precision (AP) of 0.98, significantly outperforming Logistic
 6 Regression, which scored 0.74. This indicates that Random Forest demonstrates more robust
 7 performance when handling highly imbalanced data.

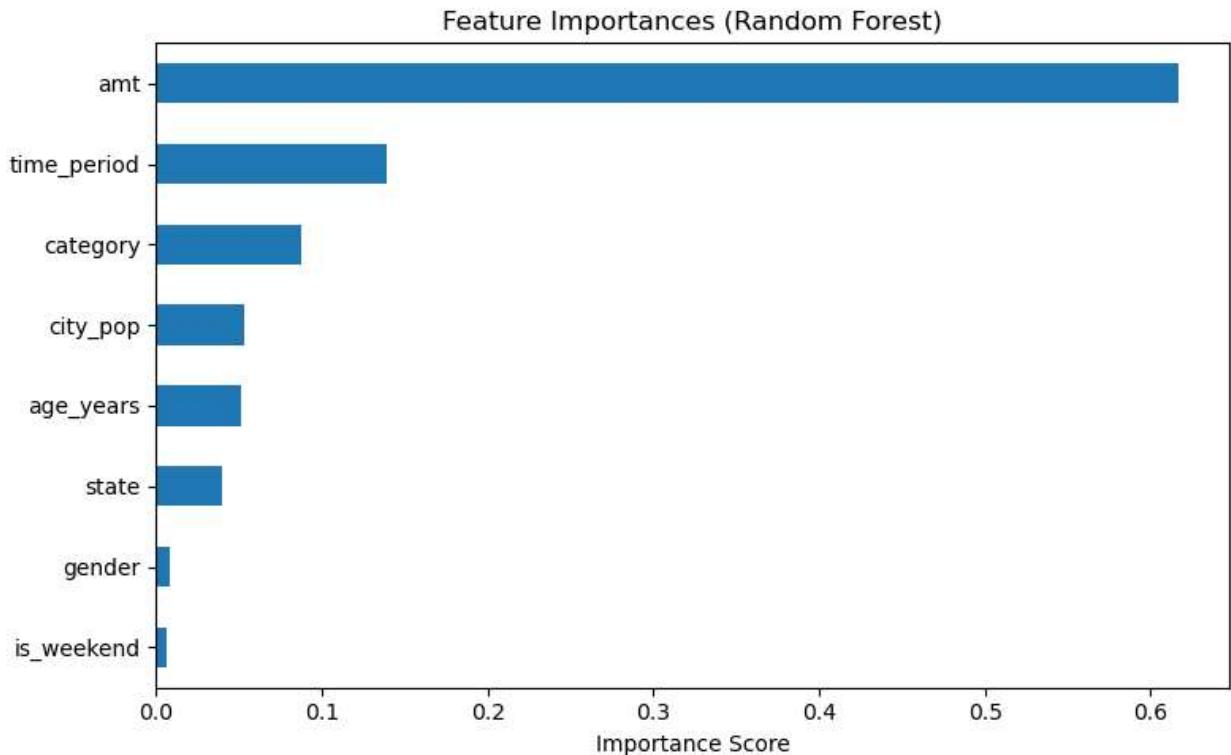


8 **Figure 10. Probability Distribution Comparison between Logistic Regression & Random Forest**

9 To further examine the distribution of model predictions, this study also show the predicted
 10 fraud probability distributions for both models (Figure 10). The results show that the Logistic
 11 Regression model produces a more dispersed distribution, with significant overlap between
 12 fraudulent and non-fraudulent samples, making the model more prone to confusion. In contrast,
 13 the Random Forest model is able to concentrate most fraudulent samples in the high-probability
 14 range, demonstrating greater prediction confidence.

15 Finally, to understand the contribution of each feature to the model's predictions, this study
 16 utilized the Feature Importance (Figure 11) provided by the Random Forest model to assess
 17 variable significance. The results indicate that transaction amount (amt) is the most critical factor,

1 followed by time period and transaction category. These variables demonstrate strong explanatory
 2 power in identifying fraudulent behavior.



3 **Figure 11.** Feature Importances from Random Forest Model

4 The detailed classification report further reveals the overall performance of both models in
 5 terms of Precision, Recall, and F1-score. For specific figures, please refer to Table 1.

	Logistic Regression			Random Forest			Support
	Precision	Recall	F1-Score	Precision	Recall	F1-Score	
0	0.96	0.79	0.87	0.99	0.97	0.98	1600
1	0.53	0.89	0.66	0.90	0.95	0.93	429
Accuracy	0.81			0.97			2066
Marco Avg	0.74	0.84	0.76	0.95	0.96	0.95	2029
Weighted Avg	0.87	0.81	0.82	0.97	0.97	0.97	2029

6 **Table 1.** Classification Report of Logistic Regression and Random Forest Models

1 4. **Findings / Conclusions**

2 Through the visual analysis of fraudulent transaction data in this study, several key findings
3 were identified:

4 First, regarding gender, although the number of fraudulent transactions involving female
5 customers is higher than that of male customers, the total loss amount is slightly greater for males.
6 This suggests that male victims may experience higher losses per transaction, while females may
7 be more susceptible to fraud in terms of frequency.

8 Second, in terms of transaction categories, fraud occurs most frequently in those related to
9 everyday purchases, particularly shopping_net (online shopping) and grocery_pos (in-store
10 grocery shopping). This indicates that fraudsters tend to target high-frequency, widely used
11 transaction types, with online channels being the primary focus due to their anonymity and ease
12 of exploitation.

13 In terms of geographic distribution, fraud cases and amounts are highly concentrated in large
14 states such as New York, Pennsylvania, and Texas. This suggests a potential positive correlation
15 between fraud risk and factors like population density and economic activity. These regions may
16 require enhanced fraud alerts and risk control mechanisms from local financial institutions and
17 platforms.

18 The age group analysis reveals that individuals between the ages of 35 and 60—classified as
19 middle-aged and working adults—experience the highest total fraud losses. Notably, those aged
20 55 and 35 each recorded losses exceeding 130,000. This may be due to their stronger financial
21 capacity and higher engagement in online or financial transactions, making them prime targets for
22 fraudsters.

23 Lastly, the analysis of transaction time indicates that fraudulent activities are concentrated
24 during the late evening and early morning hours, with 10 PM and 11 PM being peak times—

1 significantly higher than other periods. This suggests that fraudsters may be exploiting reduced
2 vigilance among victims during nighttime hours to carry out their attacks.

3 Based on the above analysis, it is recommended that businesses and institutions enhance their
4 fraud prevention mechanisms by focusing on high-risk user groups, such as middle-aged
5 individuals, high-risk time periods, particularly nighttime hours, and high-risk transaction
6 categories, such as online shopping. The implementation of real-time monitoring and alert systems
7 can further assist in reducing potential risks and strengthening overall fraud detection efforts.

8 Additionally, when comparing the performance of the two classification models—Logistic
9 Regression and Random Forest—the Random Forest model clearly outperforms Logistic
10 Regression. It achieved an overall accuracy of 97%, compared to just 81% for Logistic Regression.
11 In terms of precision and recall, Random Forest also demonstrated better balance, particularly in
12 detecting the minority class (fraudulent transactions), with a recall rate as high as 95%,
13 significantly improving the model's ability to identify anomalous cases.

14 The confusion matrix further highlights the performance gap between the two models in real-
15 world predictions. For Logistic Regression, although it accurately predicts the majority of normal
16 transactions, it shows a relatively high false positive rate when identifying fraudulent cases. In
17 contrast, the Random Forest model effectively detects most fraudulent transactions with fewer
18 misclassifications, demonstrating greater feasibility and reliability in practical applications.

19 In the feature importance analysis, both models identified transaction amount (amt) as the most
20 significant variable influencing prediction outcomes, followed by time period (time_period) and
21 transaction category (category). This indicates that fraudulent transactions often occur within
22 specific amount ranges and time windows, and are closely associated with certain transaction
23 types—findings that align with the earlier visual exploration results.

1 Further examining model stability and performance through ROC and Precision-Recall curves,
2 the Random Forest model achieved an AUC of 0.99 and an Average Precision (AP) of 0.98—
3 significantly outperforming Logistic Regression, which recorded an AUC of 0.89 and an AP of
4 0.74. These results once again confirm the superiority of Random Forest in fraud detection tasks.

5 Based on the evaluation of machine learning models and visual analysis, this study
6 recommends the use of ensemble learning methods such as Random Forest for handling highly
7 imbalanced fraud datasets. When combined with oversampling techniques like SMOTE, these
8 methods can significantly improve the model's ability to identify minority classes. Looking ahead,
9 integrating deep learning models with real-time anomaly detection mechanisms holds great
10 promise for developing more timely and accurate fraud prevention systems in practical
11 applications.

12 **5. Managerial Implications**

13 The findings of this study offer actionable insights that managers can leverage to strengthen
14 fraud prevention strategies and enhance overall operational performance. These analytics-driven
15 results can inform strategic decision-making across various organizational functions.

16 For instance, the discovery that fraudulent transactions peak between 10 PM and 11 PM
17 highlights the need for enhanced real-time monitoring during late-night hours. Additionally,
18 identifying middle-aged consumers (ages 35–60) as the most vulnerable group calls for targeted
19 financial education campaigns and increased security measures tailored to this demographic.
20 Geographic insights also pinpoint high-risk states such as New York, Texas, and Pennsylvania,
21 suggesting that resource allocation for fraud detection systems should be prioritized in these areas
22 to reduce financial losses and boost consumer trust.

23 By implementing these findings, organizations can improve internal control mechanisms,
24 safeguard customer assets, and reinforce their reputation in an increasingly complex digital

1 environment. Ultimately, applying the results of this project enables firms to proactively address
2 fraud risk while driving smarter, data-informed management decisions.

3 **6. Idea Sharing**

4 This project provided a valuable opportunity to apply business analytics concepts in a real-
5 world context, revealing both the technical challenges and practical potential of data-driven
6 problem solving. One key learning outcome was the importance of managing imbalanced data,
7 which is often encountered in fraud detection. The use of oversampling techniques such as SMOTE
8 significantly improved the model's ability to detect rare but critical events.

9 Another important takeaway was the role of data visualization in translating complex analytical
10 findings into clear, actionable insights. Tools like Tableau enabled the team to discover patterns
11 that might have been missed through raw numerical analysis alone, reinforcing the importance of
12 visual storytelling in communicating with diverse stakeholders.

13 The comparison between Random Forest and Logistic Regression models emphasized that
14 model selection should consider not only predictive performance, but also the practical context in
15 which the model will be deployed. This reinforced the idea that business analytics is not just about
16 technical methods, but about applying analytical thinking to support strategic decision-making.

17 Ultimately, the project deepened our understanding of how analytics can bridge the gap
18 between data and decisions. It reaffirmed the value of thoughtful analysis, careful model
19 evaluation, and clear communication in delivering insights that can truly support business
20 performance and customer protection.

21

1 **7. References**

2 Kelue, K. (2024, March 11). *Credit Card Fraud Prediction*. Kaggle.

3 <https://www.kaggle.com/datasets/kelvinkelue/credit-card-fraud-prediction>

4 Report, N. (2025, January 6). Payment Card Fraud Losses Approach \$34 Billion.

5 *GlobeNewswire News Room*. <https://www.globenewswire.com/news>

6 [release/2025/01/06/3004931/0/en/Payment-Card-Fraud-Losses-Approach-34-Billion.html](https://www.globenewswire.com/news-release/2025/01/06/3004931/0/en/Payment-Card-Fraud-Losses-Approach-34-Billion.html)

7 **8. Appendix**

Time	Contents
Saturday January 25, 2025	Project Group Forms
Friday February 24, 2025	Defineing the issue and Topic
Monday February 24, 2025	Project Proposal Presentation
Monday March 24, 2025	Data preprocessing
Monday April 7, 2025	Training the model
Friday April 11, 2025	Visual Exploration of the Data
Monday April 21, 2025	Final Project Paper Submission and Final Project Presentation Slides Submission