

CREDIT CARD FRAUD DETECTION 2023

End-to-End Data Science
Final Individual Project

GH1019657

22.12.2023

Abstract

In this research, we aim to address the burgeoning issue of credit card fraud in the digital age. Using an extensive set of more than 550,000 credit card transactions made by European cardholders in 2023, we are concentrating on applying cutting-edge AI and machine learning methods to create an effective fraud detection system. The scope of our study includes a thorough examination of the literature, an in-depth investigation of key features in the dataset, the development of relevant research questions, an effective methodology, and an objective evaluation of many machine learning models.

Contents

1. Introduction	3
2. Literature Review	3
• Machine learning approaches	
• Anomaly Detection Techniques	
• Integration of Multiple Techniques	
• Gap Analysis	
3. Dataset Characteristics	4
x	
4. Research Questions	4
5. Methodology	4
5.1 Data Collection	4
5.2 Data Exploration	4 - 5
5.3 Data Cleaning	6
5.4 Outlier Detection	6 - 7
5.5 Box Plot for 'Amount'	7
5.6 Correlation Matrix	8 - 9
5.7 Model Development	9 - 10
5.8 Anomaly Detection	10 - 11
5.9 Real-time Predictions	11 - 12
6. Results	12 - 13
7. Limitations	13
8. Outlooks	13
9. Conclusion	13
10. References	14

1. Introduction

Credit card fraud, an increasingly formidable challenge in today's dynamic digital landscape, poses significant threats to both financial institutions and consumers alike. The growing number of electronic transactions has made it necessary to develop and apply creative solutions that can effectively block the constantly evolving tactics used by fraudsters. The significant financial consequences and imminent risk to customer trust are sufficient to highlight how serious this issue is. Deploying and establishing improved fraud detection methods that can effectively identify and reduce the dangers associated with illegal financial transactions is imperative considering these challenges.

2. Literature Review

The detection of credit card fraud has been a recurring issue in the field of financial security, and researchers have thoroughly investigated several approaches to address this dynamic problem. A thorough analysis of this field of literature indicates the intricacy of credit card theft and the demand for advanced detection technologies.

Machine Learning Approaches:

Numerous studies have delved into the application of machine learning algorithms for fraud detection. Vaishnavi Nath Dornadula et al. (2020) emphasized the efficacy of machine learning in their research on credit card fraud detection. They explored diverse algorithms, highlighting the importance of algorithm selection in achieving high detection accuracy.

Anomaly Detection Techniques:

Anomaly detection has emerged as a pivotal approach in identifying fraudulent transactions. Meenu et al. (2020) conducted research specifically on anomaly detection in credit card transactions using machine learning. Their work shed light on the effectiveness of anomaly detection techniques in capturing irregular patterns indicative of fraud.

Integration of Multiple Techniques:

The landscape of credit card fraud is dynamic, requiring a holistic approach. Research by Btoush et al. (2023) emphasizes the importance of integrating multiple techniques, such as machine learning algorithms, anomaly detection, and deep learning, to create robust and adaptive fraud detection systems.

Gap Analysis:

While existing studies provide valuable insights, there is a recognized gap in the literature. The need for improved precision and adaptability in fraud detection systems is highlighted, aligning with the objectives of the current research.

In conclusion, the assessment of the literature highlights the variety of methods used in credit card fraud detection, from traditional machine learning algorithms to more advanced deep learning approaches. Developing efficient and flexible fraud detection systems that can handle the changing strategies used by fraudsters in the digital era requires integrating different methods.

3. Dataset Features

- 1 - The dataset includes credit card transactions in 2023 carried out by cardholders across Europe.
- 2 - It contains more than 550,000 records with anonymised transaction characteristics, including the time and location of the transaction as well as several features (V1 to V28).
- 3 - A binary label ("Class") indicating whether the transaction is fraudulent (1) or not (0) also has recorded, along with the transaction value.

id	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
0	-0.260648	-0.469648	2.496266	-0.083724	0.129681	0.732898	0.519014	-0.130006	0.727159	...	-0.110552	0.217606	-0.134794	0.165959	0.126280	-0.434824	-0.081230	-0.151045	17982.10	0
1	0.985100	-0.356045	0.558056	-0.429654	0.277140	0.428605	0.406466	-0.133118	0.347452	...	-0.194936	-0.605761	0.079469	-0.577395	0.190090	0.296503	-0.248052	-0.064512	6531.37	0
2	-0.260272	-0.949385	1.728538	-0.457986	0.074062	1.419481	0.743511	-0.095576	-0.261297	...	-0.005020	0.702906	0.945045	-1.154666	-0.605564	-0.312895	-0.300258	-0.244718	2513.54	0
3	-0.152152	-0.508959	1.746840	-1.090178	0.249486	1.143312	0.518269	-0.065130	-0.205698	...	-0.146927	-0.038212	-0.214048	-1.893131	1.003963	-0.515950	-0.165316	0.048424	5384.44	0
4	-0.206820	-0.165280	1.527053	-0.448293	0.106125	0.530549	0.658849	-0.212660	1.049921	...	-0.106984	0.729727	-0.161666	0.312561	-0.414116	1.071126	0.023712	0.419117	14278.97	0

4. Research Questions

Our exploration is guided by a set of pertinent research questions, steering the investigation towards practical solutions:

- 4.1** - What are the key features or indicators crucial for identifying fraudulent transactions?
- 4.2** - What are the potential implications for customer service and communication when a fraudulent transaction is detected?
- 4.3** - How can innovative approaches like deep learning and anomaly detection enhance the precision of fraud detection?
- 4.4** - What are the most common fraud-related transaction categories, and how can business tactics be modified to counteract these risks?
- 4.5** - How can organizations adapt to emerging fraud tactics over time to maintain the effectiveness of fraud detection models?

5. Methodology

5.1 Data Collection

Over 550,000 records have been collected, all of which conceal the details of credit card transactions made by European cardholders in 2023. The source of this invaluable dataset is attributed to Kaggle.

5.2 Data Exploration

We thoroughly reviewed the dataset, gathering important information and calculating thorough statistics with the help of the `info()` and `describe()` methods. We thoroughly investigated the column data types to ensure that they were compatible with machine learning models. The class distribution was presented visually using a countplot, which let us distinguish between authentic and fraudulent transactions. The dataset's integrity was thoroughly checked for missing values, providing an accurate starting point for further analysis and model building.

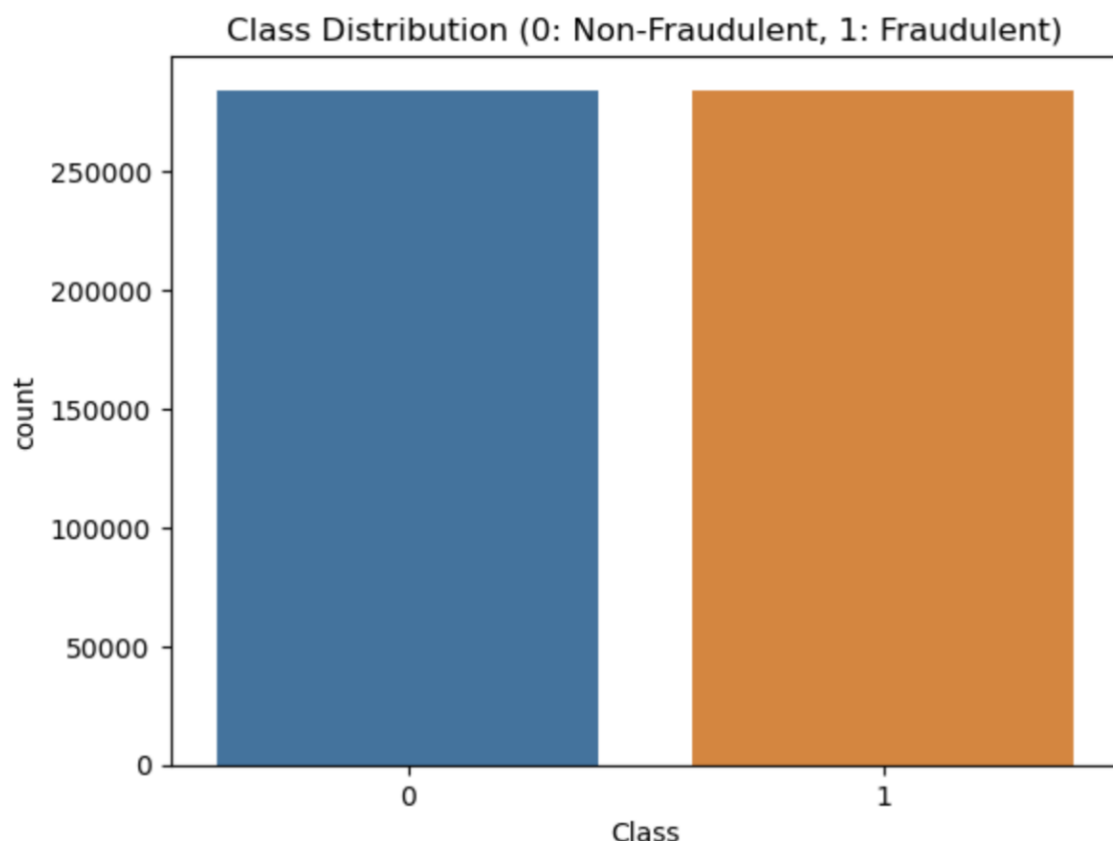
The below figure shows the database possesses 31 columns and 568,630 assets that reflect different credit card transaction information. These characteristics include the target variable ("Class"), which determines whether a transaction is fraudulent (Class=1) or not (Class=0), anonymized features (V1 to V28), and transaction amounts ("Amount"). There are no missing values in any of the columns, suggesting that the data is completely non-null. The dataset takes up around 134.5 MB of RAM. Two columns are of type int64, which is used for discrete integer values, while the remaining columns are of type float64, which represents numerical data. comprehending and becoming ready for additional analysis of the data requires having an in-depth knowledge of the dataset's structure and data types.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 568630 entries, 0 to 568629
Data columns (total 31 columns):
#   Column      Non-Null Count  Dtype
---  -
0   id          568630 non-null  int64
1   V1          568630 non-null  float64
2   V2          568630 non-null  float64
3   V3          568630 non-null  float64
4   V4          568630 non-null  float64
5   V5          568630 non-null  float64
6   V6          568630 non-null  float64
7   V7          568630 non-null  float64
8   V8          568630 non-null  float64
9   V9          568630 non-null  float64
10  V10         568630 non-null  float64
11  V11         568630 non-null  float64
12  V12         568630 non-null  float64
13  V13         568630 non-null  float64
14  V14         568630 non-null  float64
15  V15         568630 non-null  float64
16  V16         568630 non-null  float64
17  V17         568630 non-null  float64
18  V18         568630 non-null  float64
19  V19         568630 non-null  float64
20  V20         568630 non-null  float64
21  V21         568630 non-null  float64
22  V22         568630 non-null  float64
23  V23         568630 non-null  float64
24  V24         568630 non-null  float64
25  V25         568630 non-null  float64
26  V26         568630 non-null  float64
27  V27         568630 non-null  float64
28  V28         568630 non-null  float64
29  Amount      568630 non-null  float64
30  Class       568630 non-null  int64
dtypes: float64(29), int64(2)
memory usage: 134.5 MB
```

The below Figure shows dataset's distribution and attributes are shown by the summary statistics, which show that features V1 to V28 have been standardized and concealed with a mean close to 0 and a standard deviation close to 1. The transaction amounts represented by the "Amount" feature have an estimated mean of 12,041.96 and a standard deviation of 6,919.64. With a mean of 0.5, the "Class" column, which differentiates fraudulent (Class=1) from non-fraudulent (Class=0) transactions, indicates a rather balanced distribution. An in-depth understanding of feature qualities is provided by additional statistics, such as minimum, maximum, and percentile values, which aid in gathering the data for further study.

	id	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	
count	568630.000000	5.686300e+05	5.686300e+05	5.686300e+05	5.686300e+05	5.686300e+05	5.686300e+05	5.686300e+05	5.686300e+05	5.686300e+05	...	5.686300e+05	5.686300e+05	5.686300e+05
mean	284314.500000	-5.638058e-17	-1.323544e-16	-3.518788e-17	-2.879008e-17	7.197521e-18	-3.838678e-17	-3.198898e-17	2.069287e-17	9.116859e-17	...	4.758361e-17	5.398140e-18	5.398140e-18
std	164149.486121	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	1.000001e+00	...	1.000001e+00	1.000001e+00	1.000001e+00
min	0.000000	-3.495584e+00	-4.996657e+01	-3.183760e+00	-4.951222e+00	-9.952786e+00	-2.111111e+01	-4.351839e+00	-1.075634e+01	-3.751919e+00	...	-1.938252e+01	-7.734798e+00	-3.029
25%	142157.250000	-5.652859e-01	-4.866777e-01	-6.492987e-01	-6.560203e-01	-2.934955e-01	-4.458712e-01	-2.835329e-01	-1.922572e-01	-5.687446e-01	...	-1.664408e-01	-4.904892e-01	-2.376
50%	284314.500000	-9.363846e-02	-1.358939e-01	3.528579e-04	-7.376152e-02	8.108788e-02	7.871758e-02	2.333659e-01	-1.145242e-01	9.252647e-02	...	-3.743065e-02	-2.732881e-02	-5.9
75%	426471.750000	8.326582e-01	3.435552e-01	6.285380e-01	7.070047e-01	4.397368e-01	4.977881e-01	5.259548e-01	4.729905e-02	5.592621e-01	...	1.479787e-01	4.638817e-01	1.557
max	568629.000000	2.229046e+00	4.361865e+00	1.412583e+01	3.201536e+00	4.271689e+01	2.616840e+01	2.178730e+02	5.958040e+00	2.027006e+01	...	8.087080e+00	1.263251e+01	3.170

In the below Figure we can see a countplot was used to show the distribution of fraudulent (Class=1) and non-fraudulent (Class=0) transactions. It shows that we have a balanced dataset with an equal number of fraudulent (Class 1) and non-fraudulent (Class 0) transactions, each having 284,315 records.



In the initial phase of data exploration, we performed a comprehensive check for missing values within the dataset. Our analysis revealed that there are no missing values in any of the columns, confirming the completeness and integrity of the dataset. This is a pivotal finding as it ensures that the dataset is suitable for further analysis and model development.

The next stage is to determine which columns are necessary for our study and perhaps eliminate those that are not. We can better concentrate on appropriate features for our machine learning models by simplifying this approach.

5.3 Data Cleaning

Given the dataset's basic cleanliness, the only contain found was the 'id' column, which was considered unnecessary for further analysis.

5.4 Outlier Detection

Using the Z-score approach, outliers in the dataset were identified. Every data point's Z-score was computed using the `detect_outliers_zscore` function, and outliers were identified using a threshold of 3. Outlier-filled rows have been identified and shown. These anomalies will be taken into account in the stages of our study that follow if they have an influence on the analysis.

In the below figure it shows the DataFrame indicates that no rows with outliers were found in this study. This implies that there are no significant numerical outliers in the dataset.

Rows with outliers:

Empty DataFrame

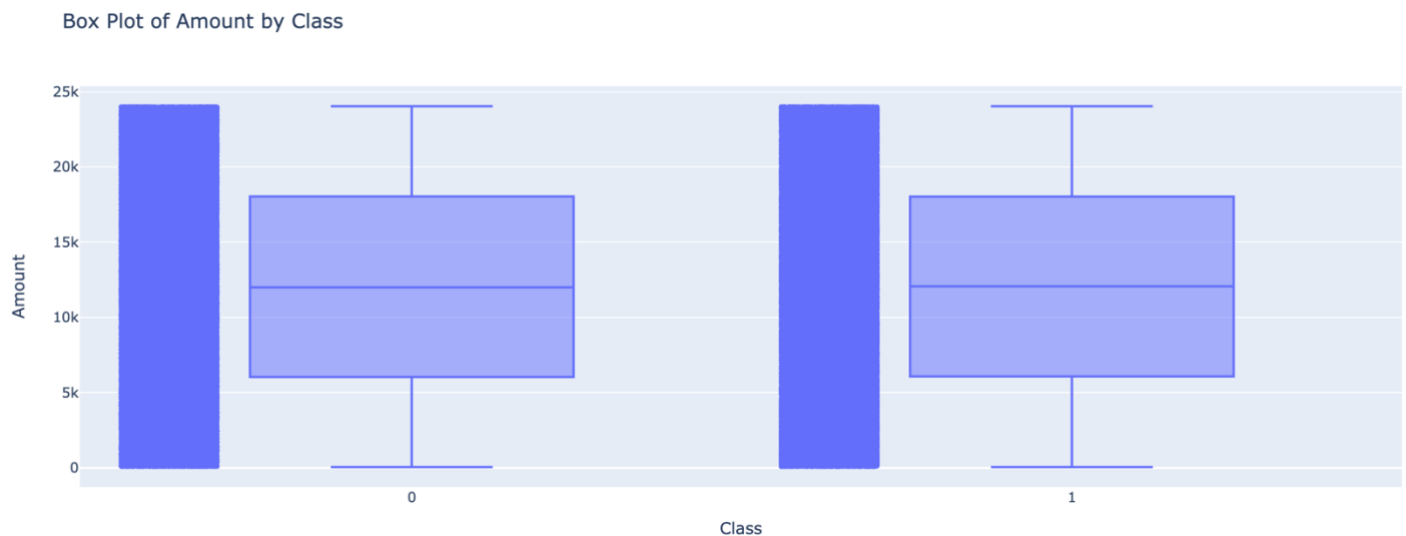
Columns: [V1, V2, V3, V4, V5, V6, V7, V8, V9, V10, V11, V12, V13, V14, V15, V16, V17, V18, V19, V20, V21, V22, V23, V24, V25, V26, V27, V28, Amount, Class]

Index: []

[0 rows x 30 columns]

5.5 Box Plot for Amount

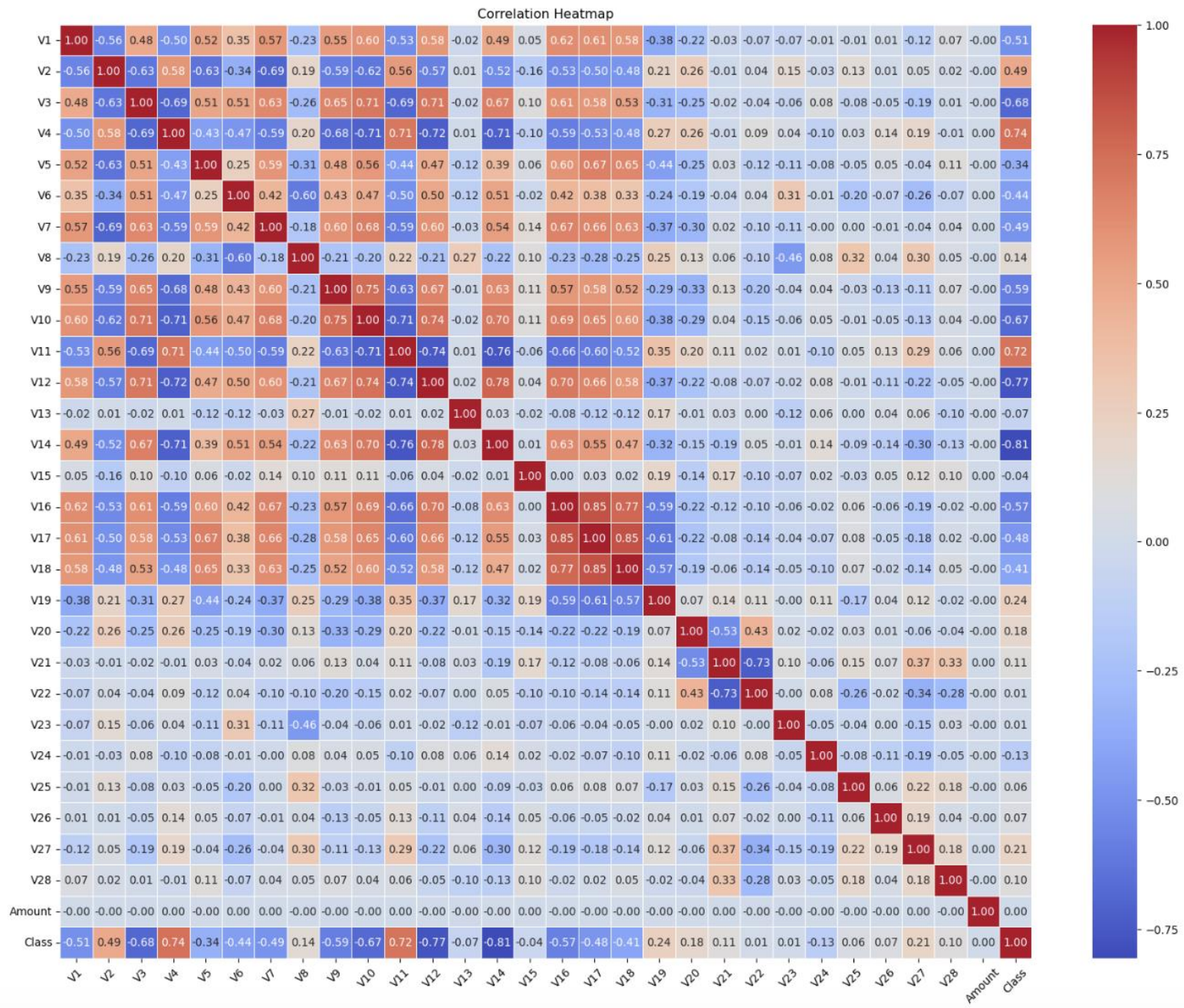
Plotly Express was used to create a box plot that showed the distribution of transaction amounts by class (fraudulent and non-fraudulent). For every class, the plot displays the distribution, the central tendency, and any possible outliers. Box plots are used to display the values for the 'Amount' for both classes. The interquartile range (IQR) for each class is displayed in a box, and a line inside the box represents the median. Individual data points outside of this range are regarded as outliers. The whiskers extend to the lowest and greatest values within a specified range. To differentiate between transactions that are fraudulent (Class=1) and those that are not (Class=0), the 'Class' variable is utilized as the x-axis. Plotting gives the distribution of transaction amounts across the two groups a visual comparison, which is essential for spotting variations that may help in the identification of fraud.



In the above figure, the box plot analysis shows that although the range and variability of transaction amounts are identical, there are slight variations in the central tendency of transaction amounts between fraudulent and non-fraudulent transactions. This suggests that additional variables or analysis are required to improve fraud detection accuracy and that transaction quantity alone may not be an accurate indicator of fraud.

5.6 Correlation Matrix

The below heatmap shows the relationships between the dataset's different features. The correlations are displayed by colors that indicate their strength and direction (positive or negative), and the correlation coefficients are given precisely by the numerical values in each cell



The relationship between characteristics and the "Class" column is particularly relevant when it comes to fraud detection. The following are some significant findings from the correlation matrix:

- The "Class" column displays significantly large positive correlations with features V17, V14, V12, V10, V11, V4, V2, V7, and V19, suggesting that these features may be more important for identifying fraudulent transactions.

- The "Class" column shows significant negative correlations with features V3, V16, V1, V6, V9, V18, V5, and V21, indicating an inverse relationship between these parameters and fraudulent transactions.
- The 'Amount' feature shows a very low correlation with the 'Class' feature, this indicated that it might not be a strong indicator of fraud on an individual basis.

5.7 Model Development

For the machine learning classification task, One-Class SVM for anomaly detection, Random Forest Classifier, and Logistic Regression were used to prepare the data.

For further binary classification procedure that identifies fraudulent detection:

1 - Data Preparation: First, the dataset is split into the target variable (y) and its features (X). 'Class,' the target variable, which specifies whether a transaction is fraudulent (1) or not (0).

2 - Data Split: The dataset was then split into training and testing using the 'train_test_split' function from scikit learn. Here, the dataset is split into 80% and 20% for training and testing respectively.

3 - Model Initialization:

There are three initialised machine learning models:

lr_model: A popular linear classification approach is called logistic regression.

A decision tree-based ensemble learning technique is called Random Forest (**rf_model**).

OneClassSVM (svm_model): Novelty detection using a support vector machine model. Here, it's applied to detect abnormalities or outliers in the data.

4 - Model Training: The training i.e X_train and y_train is used to train on each models. The lr_model, rf_model, svm_model were used to train the data respectively. Then the train data were saved using joblib as this allows for easy reuse of the models without the need to retrain them in the future.

5 - Model Predictions: The three trained models (lr_model, rf_model, and svm_model) are used to make predictions on the test data (X_test).

6 - Evaluation Metrics: Based on the test data and the model's predictions, the algorithm calculates and generates the evaluation metrics listed below for each model:

- **Accuracy:** Calculates the amount of accurate forecasts.
- **Precision:** Evaluate how well the model can anticipate favorable outcomes.
- **F1 Score:** An equilibrium between recall and accuracy that is particularly helpful in handling unbalanced datasets.

```
Logistic Regression:
Accuracy: 0.9601586268751209
Precision: 0.9783120223263958
F1 Score: 0.9594719094088497
-----
Random Forest:
Accuracy: 0.9998768971035648
Precision: 0.9997543428671697
F1 Score: 0.9998771563448748
-----
SVM:
Accuracy: 0.49803035365703535
Precision: 0.49032258064516127
F1 Score: 0.08968123614676851
```

The above figure is the output of using models like Logistic Regression, Random Forest and SVM:

- The Random Forest model has a very high F1 score, accuracy, and precision, and works very effectively. It indicates that fraudulent transaction detection is an attribute of the Random Forest model.
- High accuracy, precision, and F1 scores are also achieved using the Logistic Regression model, which also performs well. It's a reliable choice for this study.
- In contrast, the SVM model (OneClassSVM) has inadequate performance. Its considerably lower F1 score, accuracy, and precision imply that it would not be appropriate for this particular fraud detection task.

Based on the evaluation metrics provided, the Random Forest model appears to be the best choice for the task of fraud detection in this study. It has demonstrated high accuracy, precision, and F1 score, indicating its ability to effectively detect fraudulent transactions.

5.8 Anomaly Detection

In the pursuit of anomaly detection, a Random Forest classifier was employed, a process delineated below:

1 - Importing necessary libraries:

- Utilization of RandomForestClassifier from scikit-learn's ensemble module to construct the Random Forest model.
- Application of train_test_split to bifurcate the dataset into training and testing sets.
- Employment of precision_score, recall_score, and f1_score from scikit-learn's metrics module to assess model performance.

2 - Initialize and train Random Forest Model: A RandomForestClassifier object (rf_model) with 100 trees is generated. Fit is then used to train the model using the training set of data.

3 - Make Predictions: Predict is used to make predictions based on the testing results. The predictions in this instance are binary (0 for true transactions and 1 for fraudulent ones).

4 - Set a threshold: A threshold value (in this case, 0.5) is used to identify anomalies. The probability at which a transaction can be considered fraudulent is defined by this threshold.

5 - Evaluate the model: The accuracy, recall, and F1 score of the anomaly detection model are computed to evaluate its performance. These metrics assist in evaluating the model's overall efficacy, accuracy, and completeness in detecting fraudulent transactions.

Anomaly Detection (Random Forest):

Precision: 0.9997894293535481

Recall: 1.0

F1 Score: 0.9998947035906075

To summarise, the Random Forest model shows remarkable performance in anomaly identification, as demonstrated by its high accuracy, recall, and F1 score. This model is a great option for anomaly detection since it is very good at detecting fraudulent transactions while reducing false alarms.

5.9 Real-time Predictions

In order to make predictions in real-time, a Random Forest model that has been trained beforehand is used to classify obtained transaction data as real or fake. The steps involved are as follows:

Model Import and Application:

Importing a pre-trained Random Forest model and using it to handle incoming transaction data for fraud categorization is the purpose of the code. The model classifies the given data in binary after being trained on a variety of features (V1 to V28 and Amount).

Interpretation of the Output:

```
# Load the trained Random Forest model
loaded_model = joblib.load('fraud_detection_model.pkl')

new_data = pd.DataFrame({
    'V1': [0.5],
    'V2': [-0.3],
    'V3': [1.2],
    'V4': [-0.7],
    'V5': [0.4],
    'V6': [0.8],
    'V7': [0.6],
    'V8': [-0.2],
    'V9': [0.5],
    'V10': [0.7],
    'V11': [0.1],
    'V12': [0.2],
    'V13': [0.3],
    'V14': [0.4],
    'V15': [0.5],
    'V16': [0.6],
    'V17': [0.7],
    'V18': [0.8],
    'V19': [0.9],
    'V20': [0.10],
    'V21': [0.11],
    'V22': [0.12],
    'V23': [0.13],
    'V24': [0.14],
    'V25': [0.15],
    'V26': [0.16],
    'V27': [0.17],
    'V28': [0.18],
    'Amount': [100.0]
})

# Make predictions on the new data using the trained model
predictions = loaded_model.predict(new_data)

# Define a threshold for classifying anomalies
threshold = 0.5

# Classify transactions based on the threshold
classified = (predictions > threshold).astype(int)

# Depending on your business logic, you can take various actions based on the classification results:
if classified == 1:
    # Take actions for fraudulent transactions (e.g., alert, block, or investigate)
    print("Fraudulent Transaction Detected!")
else:
    # Process normal transactions
    print("Normal Transaction")
```

Normal Transaction

After the code has run, the “Normal Transaction” output is seen. This label indicates that the transaction data is considered normal by the loaded Random Forest model. Based on the input attributes and the predetermined threshold, it is implied that the model did not flag the transaction as fraudulent when there was no “Fraudulent Transaction” result

In essence, the real-time prediction is a key tool for evaluating incoming transactions and offers an instantaneous categorization that facilitates the detection of any fraudulent activity. The model’s capacity to identify typical transactions highlights how well it can discriminate between legitimate and questionable financial activity in real-time situations.

6. Results

In our comprehensive analysis of the credit card fraud detection dataset, various aspects were explored, leading to the development and evaluation of machine learning models for fraud detection.

Exploring Data Analysis:

- **Dataset Summary:** The dataset, comprising 568,630 entries, exhibited no missing values, providing a solid foundation for subsequent analyses.
- **Visualization:** Utilizing a box plot, we observed higher amounts in fraudulent transactions, emphasizing their distinctive nature.
- **Statistical Analysis:** Correlation matrix analysis identified features with strong correlations, particularly 'V14' and 'V17,' indicating their significance in fraud identification.
- **Outlier Detection:** Although an outlier detection method based on Z-scores was applied, no outliers were found.

Model Building and Evaluation:

- **Data Split:** The dataset was divided into 80% training and 20% testing sets.
Models Trained: One-Class SVM, Random Forest, and Logistic Regression were the three models that were trained.
- **Model Evaluation:** Compared to the other models, the Random Forest model performed better, obtaining a high F1 score, accuracy, and precision. The model that was found to be most appropriate for detecting fraud was accepted.
- **Anomaly detection:** The Random Forest model was used to discover anomalies, and a threshold for classifying them was determined. The model had a strong ability to detect fraudulent transactions, as indicated by its excellent accuracy, recall, and F1 score.

Answers to Research Questions:

5.1 - What are the key features or indicators that can be used in the identification of fraudulent transactions?

- The dataset analysis revealed that certain features, such as 'V1' through 'V28' and 'Amount,' are crucial for identifying fraudulent transactions. Notably, features having the strongest negative correlations with the 'Class' variable are 'V14' and 'V17,' which can be valuable indications. Another way to increase feature relevance is to study feature engineering and selection methods.

5.2 - What are the potential implications for customer service and communication when a fraudulent transaction is detected?

- Better customer service can result from detecting fraudulent transactions as they stop unauthorized transactions and guarantee safety for customers. It's critical to notify customers as soon as a fraudulent transaction becomes apparent in order to explain the circumstances, walk them through the next steps, and offer help. Reducing possible interruptions and upholding confidence are two benefits of having an efficient customer communication strategy.

5.3 - Are there any innovative approaches or technologies, such as deep learning or anomaly detection, that could improve the precision of fraud detection in the future?

- Indeed, cutting-edge technologies with the potential to improve fraud detection precision include deep learning and anomaly detection. Neural networks and other deep learning algorithms are able to recognize intricate patterns in data and adjust to changing fraud strategies. Through the reduction of false positives and the identification of new fraud trends, anomaly detection techniques—especially when integrated with dynamic thresholding—can improve detection precision.

5.4 - What are the most regular fraud-related transaction categories in the dataset, and how can the business modify its fraud protection tactics to counteract these particular risks?

- Transaction types are not specifically classified in the dataset. Transaction data would need to be analyzed or classified in order to determine the most common fraud-related categories. This might reveal which kinds of transactions are more frequently the target of fraud. Changing fraud prevention strategies would entail focusing on the weaknesses connected to these kinds of high-risk transactions.

5.5 - In order to maintain the effectiveness of the fraud detection models over time, how can the organization adapt to emerging fraud tactics and techniques?

- The following actions should be taken by the organization to respond to new fraud strategies and techniques:
 - a) Use dynamic thresholding to continually change model sensitivity.
 - b) Retrain and update models on a regular basis with new data to take changing strategies into consideration.
 - c) Explore adding new features and cutting-edge methods to expand the feature set.
 - d) To remain ahead of new risks and study advanced fraud detection technology, make research and development investments.

In the context of evolving fraud methods, these tactics will assist the organisation in continuing to detect fraud effectively.

7. Limitations

Recognizing the limits is essential even if our study yielded insightful results. Specific fraud-related categories are difficult to identify in the dataset due to the absence of clear transaction type classification. Furthermore, it highlights the dynamic nature of fraud detection and the need for constant effort to adapt to new fraud strategies.

8. Outlooks

In order to improve the accuracy of fraud detection, further studies may involve methods based on deep learning, feature engineering, and dynamic thresholding. It should continue to be a priority to update models often and implement preventative measures to counter new fraud techniques.

9. Conclusion

Our study aims to avoid fraud with credit cards by developing and evaluating machine learning models. The best-performing model was the Random Forest model, highlighting the significance of model selection. Model implementation, dynamic thresholding, continuous improvements, feature improvement, and proactive communication with consumers are among the suggestions. Organizations need to be on the lookout for emerging fraud strategies and make continuous investments in research and development.

References

N.E., 2023, Credit Card Fraud Detection Dataset 2023, Kaggle. Available at: <https://www.kaggle.com/datasets/nelgiriyeewithana/credit-card-fraud-detection-dataset-2023> (Accessed: 05 November 2023).

Learn, no date, scikit. Available at: <https://scikit-learn.org/stable/index.html> (Accessed: 05 November 2023).

Joblib, no date, Running python functions as pipeline jobs. Available at: <https://joblib.readthedocs.io/en/latest/> (Accessed: 05 November 2023).

Vaishnavi Nath Dornadula et al., 2020, Credit card fraud detection using machine learning algorithms, Procedia Computer Science. Available at: <https://www.sciencedirect.com/science/article/pii/S187705092030065X> (Accessed: 05 November 2023).

Meenu et al., 2020, Anomaly detection in credit card transactions using Machine Learning, SSRN. Available at: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3670230 (Accessed: 05 November 2023).

Marazqah Btoush, E.A.L. et al., 2023, A systematic review of literature on credit card cyber fraud detection using machine and Deep Learning, PeerJ. Computer science. Available at: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10280638/> (Accessed: 05 November 2023).

Btoush, E.A.L. et al., 2023, "A systematic review of literature on credit card cyber fraud detection using machine and Deep Learning," PeerJ. Computer science. Available at: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10280638/> (Accessed: 05 November 2023).