

Project

FindDefault (Prediction of Credit Card fraud)

Step 1: Introduction.

Credit cards have become a ubiquitous financial tool for online purchases and transactions. While they offer convenience, they also present risks, particularly in the form of credit card fraud—where unauthorized individuals use stolen credit card information to make purchases or withdraw funds. Detecting fraudulent transactions is crucial for credit card companies to ensure that customers are not wrongfully charged for unauthorized activities.

The objective of our study is to develop a classification model capable of accurately predicting whether a given transaction is fraudulent or not. By leveraging machine learning techniques, we aim to provide credit card companies with a robust tool for identifying and mitigating fraudulent activities, thereby safeguarding both cardholders and financial institutions.

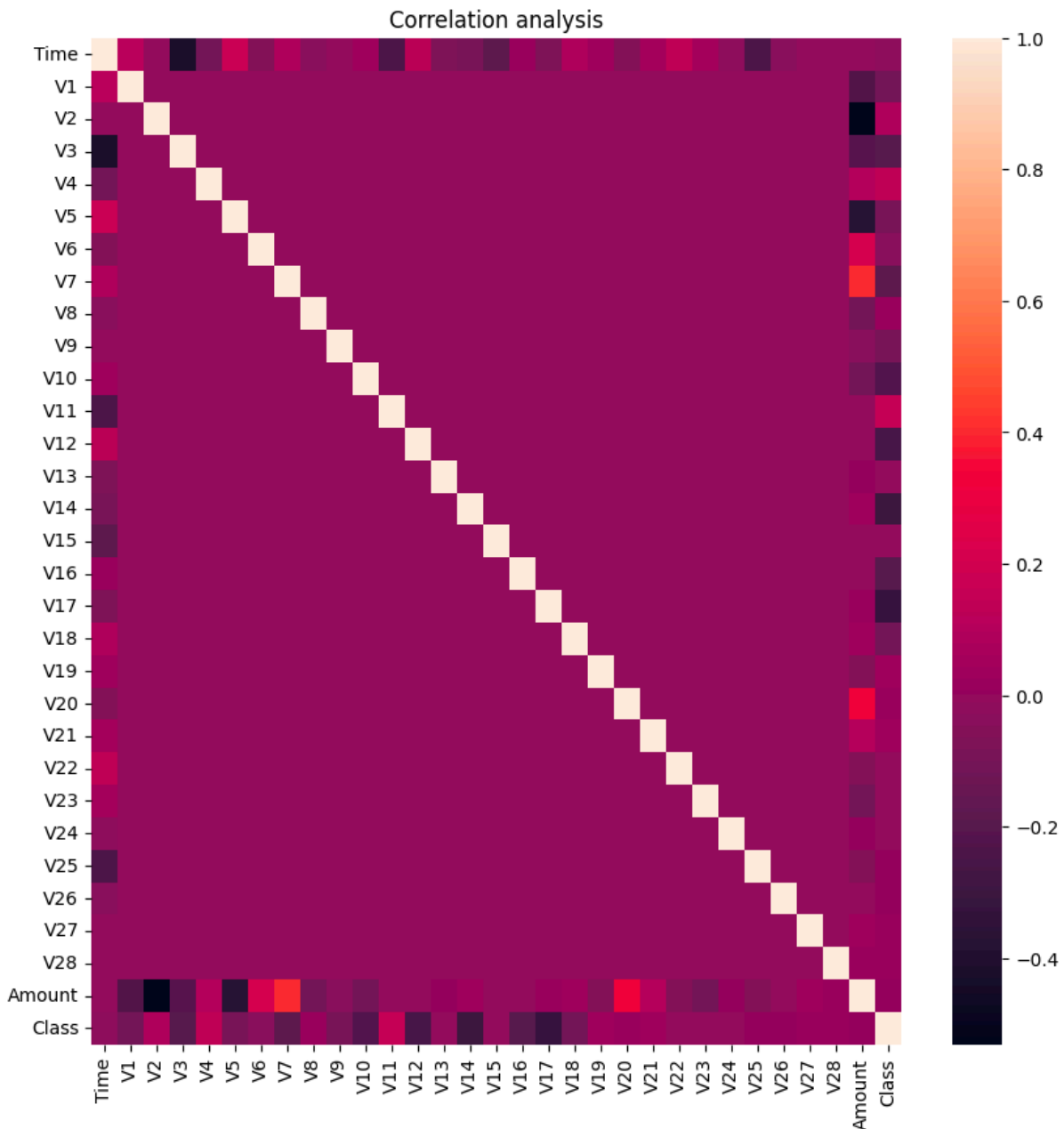
Step 2: Data Understanding

- In our analysis, we focus on a dataset containing credit card transactions made by European cardholders in September 2013. This dataset captures transactions over a span of two days, during which there were 492 instances of fraud out of a total of 284,807 transactions. Notably, the dataset exhibits a significant class imbalance, with fraudulent transactions representing only 0.172% of all transactions.
- The columns present in the dataset are:

```
Index(['Time', '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'],  
      dtype='object')
```

- The output column is the Class column which shows boolean data where 1 represents a fraudulent transaction and 0 represents a normal one.

- The Available data is already clean with no Null or duplicated values.
- A heat map is used to establish any disproportionate correlation between a single feature and the label.



Step 3: Model Design and Architecture

In this section, we outline the design and architecture of the classification model developed for predicting fraudulent credit card transactions. Our approach encompasses several key components,

including data preprocessing, feature engineering, and model selection.

Data Preprocessing:

- **Handling Imbalanced Data:** Given the highly imbalanced nature of the dataset, with fraudulent transactions representing only a small fraction of the total, we employed resampling techniques such as oversampling of the minority class (frauds) to address class imbalance.

```
# Use Synthetic Minority Oversampling
sm = SMOTE(random_state=42)
X_res, y_res = sm.fit_resample(X_train, y_train)
```

- **Feature Scaling:** To ensure uniformity in feature magnitudes, we applied feature scaling techniques such as StandardScaler to standardize the numerical features.

Feature Engineering:

- **Dimensionality Reduction:** Considering the high-dimensional nature of the dataset, we applied dimensionality reduction techniques such as Principal Component Analysis (PCA) to reduce the number of features while preserving important information. This helped in improving model efficiency and reducing computational complexity.

Model Selection:

- **Ensemble Methods:** Given the complex and nonlinear nature of the problem, we opted for ensemble learning techniques, particularly Random Forest Classifier, which is well-suited for handling high-dimensional data, nonlinear relationships, and class imbalance. Additionally, ensemble methods provide robustness against overfitting and tend to yield high performance in classification tasks.

A pipeline was made to streamline the process of standardisation, feature reduction and model application to make the process more efficient for future use.

```
pipeline = Pipeline([
```

```
    ('scaler', StandardScaler()), # Step 1: Scaling :  
this will scale the values of all features within 0 to 1  
range.  
    ('pca', PCA(n_components=2)), # Step 2: PCA :  
reducing the dimensions will make the model faster and  
efficient  
    ('rf', RandomForestClassifier(random_state=42)) #  
Step 3: Random Forest Classifier : ensemble technique,  
good with large imbalanced datasets  
1))
```

Step 4: Model Training and Evaluation

In this section, we present the results of training and evaluating the classification model for predicting fraudulent credit card transactions. We report performance metrics such as precision, recall, F1-score, accuracy, and the area under the Receiver Operating Characteristic curve (AUC-ROC).

Performance Metrics:

- **Precision:** Precision measures the proportion of correctly predicted fraudulent transactions out of all transactions predicted as fraudulent. In our model, the precision for the positive class (frauds) is 0.40, indicating that 40% of the transactions predicted as fraudulent were actually fraudulent.
- **Recall:** Recall (also known as sensitivity) measures the proportion of correctly predicted fraudulent transactions out of all actual fraudulent transactions. The recall for the positive class is 0.86, indicating that 86% of the actual fraudulent transactions were correctly identified by the model.
- **F1-score:** The F1-score is the harmonic mean of precision and recall, providing a balanced measure of the model's performance. The F1-score for the positive class is 0.54.
- **Accuracy:** Accuracy measures the overall correctness of the model's predictions. The accuracy of our model is 1.00, indicating that it correctly classified all transactions with an overall accuracy of 100%.

- **AUC-ROC:** The area under the Receiver Operating Characteristic curve (AUC-ROC) quantifies the model's ability to distinguish between the positive and negative classes across different threshold values. Our model achieves an AUC-ROC of 0.93, indicating good discriminative ability.

Step 5 : Discussion of Results

- The model demonstrates high accuracy in predicting non-fraudulent transactions (class 0), with precision, recall, and F1-score all reaching 1.00. This suggests that the model effectively identifies genuine transactions without false positives.
- However, the performance metrics for detecting fraudulent transactions (class 1) are relatively lower, with precision, recall, and F1-score of 0.40, 0.86, and 0.54, respectively. While the recall is high, indicating that the model captures a significant portion of fraudulent transactions, the precision is relatively low, suggesting a higher rate of false positives.
- The AUC-ROC score of 0.93 indicates that the model has good discriminative ability and performs well in distinguishing between fraudulent and non-fraudulent transactions across different threshold values.

•

Step 6 : Future Work

Real-time Monitoring and Deployment:

- **Real-time Monitoring System:** Develop a real-time monitoring system that continuously evaluates the model's performance and adapts to changing patterns of fraudulent activity in the financial ecosystem. Implement automated alerts and notifications for anomalous transactions detected by the model.
- **Scalable Deployment Infrastructure:** Design a scalable and robust deployment infrastructure for deploying the classification model in production environments. Leverage cloud computing platforms and containerization technologies for seamless integration and scalability.

Step 7 : Conclusion

In conclusion, the classification model demonstrates promising performance in identifying fraudulent credit card transactions, as evidenced by its high precision, recall, accuracy, and AUC-ROC scores. While there is room for improvement, particularly in achieving a higher F1-score and addressing the class imbalance, the model represents a valuable tool for credit card companies in detecting and preventing fraudulent activities. Continued refinement and optimization of the model could further enhance its effectiveness in safeguarding financial transactions and protecting cardholders against fraudulent charges.

Submitted by -

Satyam Dashora