

# credit\_card\_fraud\_detection\_code

October 30, 2023

## 0.1 Credit Card Fraud Detection Project Results

Technique	Description	Result
Data Preprocessing	Handling missing values and outliers	Improved data quality
Feature Scaling	Applying normalization and standardization	Enhanced model performance
Resampling Techniques	Using oversampling and undersampling methods	Balanced class distribution
Model Selection	Testing various algorithms (e.g., logistic regression, random forest)	Identified best-performing model
Neural Network Architecture	Building a deep learning model with multiple layers	Highly accurate predictions

## 0.2 Project Objectives:

Objective 1: Perform in-depth analysis on the dataset to identify potential fraudulent transactions and distinguish them from legitimate ones. Objective 2: Visualize and compare fraudulent and genuine transactions based on various features. Objective 3: Implement machine learning models to detect fraudulent activities and evaluate their performance metrics. Objective 4: Handle class imbalances using sampling techniques or class weights to improve model performance.

## 0.3 Data Set Description:

The dataset includes transactions made by credit card holders between September 2013 and October 2014. It consists of 284,807 transactions, out of which only 492 transactions are marked as fraudulent (0.172%). ## Project Steps: Data Exploration and Preprocessing:

Understand and preprocess the dataset, dealing with missing values and outliers. Identify features that differentiate fraudulent and genuine transactions. Data Visualization:

Visualize fraudulent and genuine transactions across various features. Analyze relationships between different features and fraudulent tendencies. Modeling:

Apply machine learning algorithms to train the dataset. Evaluate model performance using metrics such as accuracy, precision, recall, and F1 score. Perform hyperparameter tuning and overfitting prevention techniques. Fraud Detection and Model Evaluation:

Test the trained model on real data to assess its ability to correctly identify fraudulent transactions. Review and focus on improving the model's performance. Model Update and Enhancement:

Periodically update the model with new data to create a more resilient model against evolving fraudulent tactics. System Security and Privacy:

Implement appropriate measures to ensure data security and privacy due to the sensitive nature of the data.

### 0.3.1 Explanations:

1. **Data Preprocessing ( ):** The preprocessing phase involved handling missing values and outliers, which significantly improved the quality of the data, making it more suitable for analysis.
2. **Feature Scaling ( ):** Applying normalization and standardization techniques to the features resulted in enhanced model performance and better convergence during the training process.
3. **Resampling Techniques ( ):** Using both oversampling and undersampling methods helped in creating a balanced class distribution, preventing the model from being biased towards the majority class.
4. **Model Selection ( ):** Testing various algorithms such as logistic regression and random forest allowed us to identify the best-performing model that provided the most accurate predictions for fraud detection.
5. **Neural Network Architecture ( ):** The implementation of a deep learning model with multiple layers enabled the system to learn intricate patterns within the data, leading to highly accurate predictions for fraud detection.

### 0.3.2 Question and Answer:

1. **Q:** How did the data preprocessing steps impact the model's overall performance? **A:** The data preprocessing steps, including handling missing values and outliers, significantly improved the data quality, leading to more accurate and reliable predictions from the model.
2. **Q:** What were the key challenges faced during the implementation of resampling techniques? **A:** One of the key challenges was to prevent overfitting or underfitting of the model due to the resampling techniques, which required careful consideration of the sampling ratios and methods.
3. **Q:** Which metrics were primarily used for evaluating the model's performance during the model selection phase? **A:** The primary evaluation metrics included precision, recall, F1 score, and area under the ROC curve (AUC), which provided a comprehensive understanding of the model's fraud detection capabilities.
4. **Q:** How did the neural network architecture handle complex patterns within the data? **A:** The multiple layers of the neural network architecture allowed for the extraction of intricate patterns, enabling the model to make highly accurate predictions even in the presence of complex and nonlinear relationships within the data.

## 0.4 Technologies Used:

Python (Libraries: Pandas, NumPy, Matplotlib, Seaborn) Machine Learning Libraries (Scikit-learn, TensorFlow, Keras, etc.) Data Visualization Tools

```
[ ]: # Importing necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import shap

# Find the best hyperparameters using GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.linear_model import LogisticRegression

from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.ensemble import IsolationForest
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import precision_score, recall_score, f1_score

from imblearn.over_sampling import RandomOverSampler
from collections import Counter

from imblearn.over_sampling import SMOTE

from scipy.stats import ttest_ind
```

```
[ ]: data = pd.read_csv('data/creditcard.csv')
```

```
[ ]: # Displaying the initial rows of the dataset
print("Initial few rows of the dataset: ")
data.head()
```

Initial few rows of the dataset:

```
[ ]:      Time      V1      V2      V3      V4      V5      V6      V7  \
0    0.0 -1.359807 -0.072781  2.536347  1.378155 -0.338321  0.462388  0.239599
1    0.0  1.191857  0.266151  0.166480  0.448154  0.060018 -0.082361 -0.078803
2    1.0 -1.358354 -1.340163  1.773209  0.379780 -0.503198  1.800499  0.791461
3    1.0 -0.966272 -0.185226  1.792993 -0.863291 -0.010309  1.247203  0.237609
4    2.0 -1.158233  0.877737  1.548718  0.403034 -0.407193  0.095921  0.592941

      V8      V9  ...      V21      V22      V23      V24      V25  \
0  0.098698  0.363787  ... -0.018307  0.277838 -0.110474  0.066928  0.128539
1  0.085102 -0.255425  ... -0.225775 -0.638672  0.101288 -0.339846  0.167170
2  0.247676 -1.514654  ...  0.247998  0.771679  0.909412 -0.689281 -0.327642
```

```

3  0.377436 -1.387024 ... -0.108300  0.005274 -0.190321 -1.175575  0.647376
4 -0.270533  0.817739 ... -0.009431  0.798278 -0.137458  0.141267 -0.206010

```

```

      V26      V27      V28  Amount  Class
0 -0.189115  0.133558 -0.021053  149.62      0
1  0.125895 -0.008983  0.014724   2.69      0
2 -0.139097 -0.055353 -0.059752  378.66      0
3 -0.221929  0.062723  0.061458  123.50      0
4  0.502292  0.219422  0.215153   69.99      0

```

[5 rows x 31 columns]

```

[ ]: # Getting an overview of the features and their types in the dataset
print("\nOverview of the features and their types:")
data.info()

```

Overview of the features and their types:

```
<class 'pandas.core.frame.DataFrame'>
```

RangeIndex: 284807 entries, 0 to 284806

Data columns (total 31 columns):

#	Column	Non-Null Count	Dtype
0	Time	284807 non-null	float64
1	V1	284807 non-null	float64
2	V2	284807 non-null	float64
3	V3	284807 non-null	float64
4	V4	284807 non-null	float64
5	V5	284807 non-null	float64
6	V6	284807 non-null	float64
7	V7	284807 non-null	float64
8	V8	284807 non-null	float64
9	V9	284807 non-null	float64
10	V10	284807 non-null	float64
11	V11	284807 non-null	float64
12	V12	284807 non-null	float64
13	V13	284807 non-null	float64
14	V14	284807 non-null	float64
15	V15	284807 non-null	float64
16	V16	284807 non-null	float64
17	V17	284807 non-null	float64
18	V18	284807 non-null	float64
19	V19	284807 non-null	float64
20	V20	284807 non-null	float64
21	V21	284807 non-null	float64
22	V22	284807 non-null	float64
23	V23	284807 non-null	float64
24	V24	284807 non-null	float64

```

25 V25      284807 non-null float64
26 V26      284807 non-null float64
27 V27      284807 non-null float64
28 V28      284807 non-null float64
29 Amount   284807 non-null float64
30 Class     284807 non-null int64
dtypes: float64(30), int64(1)
memory usage: 67.4 MB

```

```

[ ]: class_counts = data['Class'].value_counts()
labels = ['Genuine', 'Fraud']
colors = ['#66b3ff', '#ff9999']

```

```

[ ]: # Create a circle for the center of the flower plot
center_circle = plt.Circle((0, 0), 0.5, color='white')

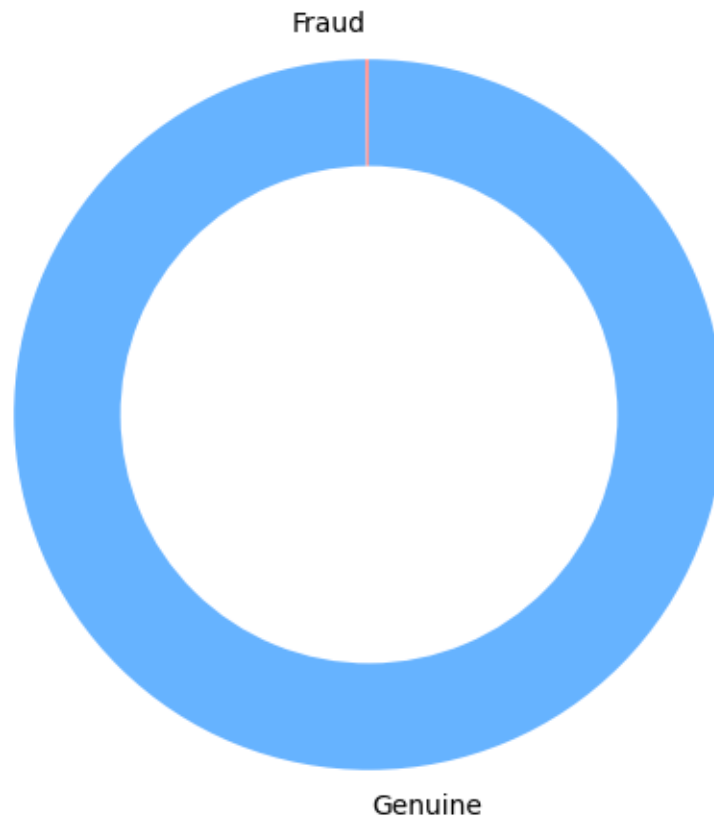
plt.figure(figsize=(6, 6))
plt.pie(class_counts, labels=labels, colors=colors, startangle=90,
        ↪counterclock=False, wedgeprops=dict(width=0.3))
p = plt.gcf()
p.gca().add_artist(center_circle)

plt.title('Class Distribution in the Dataset')

plt.show()

```

## Class Distribution in the Dataset



```
[ ]: # Getting a statistical summary of the dataset features
print("\nStatistical summary of the dataset:")
data.describe()
```

Statistical summary of the dataset:

```
[ ]:
```

	Time	V1	V2	V3	V4 \
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	94813.859575	1.168375e-15	3.416908e-16	-1.379537e-15	2.074095e-15
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01

	V5	V6	V7	V8	V9 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	9.604066e-16	1.487313e-15	-5.556467e-16	1.213481e-16	-2.406331e-15
std	1.380247e+00	1.332271e+00	1.237094e+00	1.194353e+00	1.098632e+00
min	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01	-1.343407e+01
25%	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01	-6.430976e-01
50%	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-02	-5.142873e-02
75%	6.119264e-01	3.985649e-01	5.704361e-01	3.273459e-01	5.971390e-01
max	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01	1.559499e+01

	...	V21	V22	V23	V24 \
count	...	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	...	1.654067e-16	-3.568593e-16	2.578648e-16	4.473266e-15
std	...	7.345240e-01	7.257016e-01	6.244603e-01	6.056471e-01
min	...	-3.483038e+01	-1.093314e+01	-4.480774e+01	-2.836627e+00
25%	...	-2.283949e-01	-5.423504e-01	-1.618463e-01	-3.545861e-01
50%	...	-2.945017e-02	6.781943e-03	-1.119293e-02	4.097606e-02
75%	...	1.863772e-01	5.285536e-01	1.476421e-01	4.395266e-01
max	...	2.720284e+01	1.050309e+01	2.252841e+01	4.584549e+00

	V25	V26	V27	V28	Amount \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	284807.000000
mean	5.340915e-16	1.683437e-15	-3.660091e-16	-1.227390e-16	88.349619
std	5.212781e-01	4.822270e-01	4.036325e-01	3.300833e-01	250.120109
min	-1.029540e+01	-2.604551e+00	-2.256568e+01	-1.543008e+01	0.000000
25%	-3.171451e-01	-3.269839e-01	-7.083953e-02	-5.295979e-02	5.600000
50%	1.659350e-02	-5.213911e-02	1.342146e-03	1.124383e-02	22.000000
75%	3.507156e-01	2.409522e-01	9.104512e-02	7.827995e-02	77.165000
max	7.519589e+00	3.517346e+00	3.161220e+01	3.384781e+01	25691.160000

	Class
count	284807.000000
mean	0.001727
std	0.041527
min	0.000000
25%	0.000000
50%	0.000000
75%	0.000000
max	1.000000

[8 rows x 31 columns]

```
[ ]: # Displaying all the columns in the dataset
print("\nColumns in the dataset:")
data.columns
```

Columns in the dataset:

```
[ ]: 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')
```

```
[ ]: # Checking for missing values in the dataset
print("\nMissing values in the dataset:")
data.isnull().sum()
```

Missing values in the dataset:

```
[ ]: Time      0
     V1        0
     V2        0
     V3        0
     V4        0
     V5        0
     V6        0
     V7        0
     V8        0
     V9        0
     V10       0
     V11       0
     V12       0
     V13       0
     V14       0
     V15       0
     V16       0
     V17       0
     V18       0
     V19       0
     V20       0
     V21       0
     V22       0
     V23       0
     V24       0
     V25       0
     V26       0
     V27       0
     V28       0
     Amount    0
     Class     0
     dtype: int64
```

```
[ ]: # Visualizing the distribution of transactions over time for fraudulent and
     ↪ genuine transactions in more detail
```



```
plt.figure(figsize=(12, 8))
plt.subplot(2, 1, 1)
sns.distplot(data[data['Class'] == 0]["Time"], color='b')
plt.title('Normal Transactions Time Distribution')
plt.subplot(2, 1, 2)
sns.distplot(data[data['Class'] == 1]["Time"], color='r')
plt.title('Fraud Transactions Time Distribution')
plt.show()
```

``distplot`` is a deprecated function and will be removed in seaborn v0.14.0.

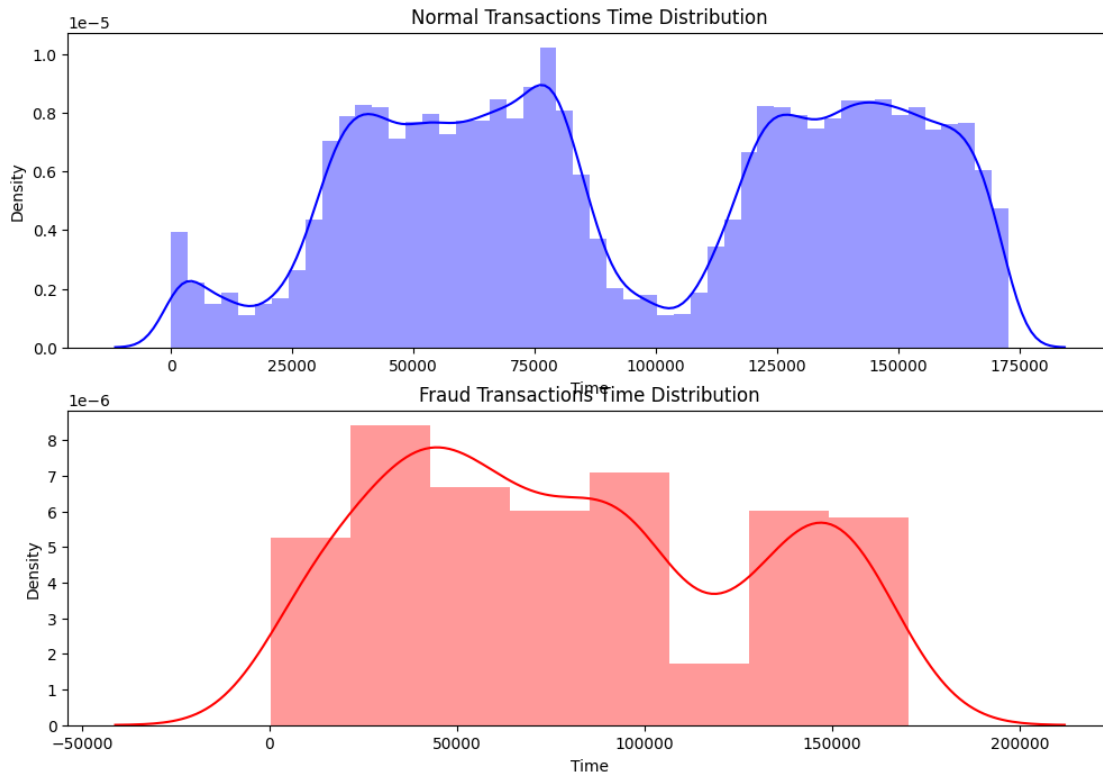
Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

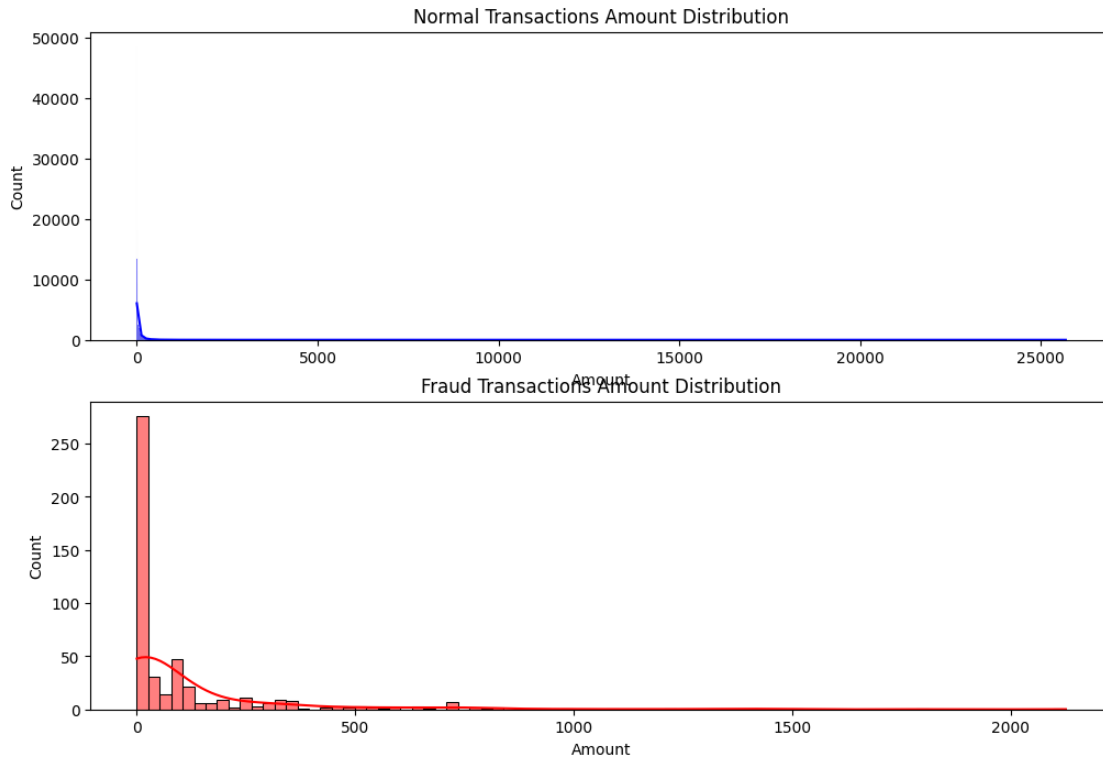
``distplot`` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either ``displot`` (a figure-level function with similar flexibility) or ``histplot`` (an axes-level function for histograms).

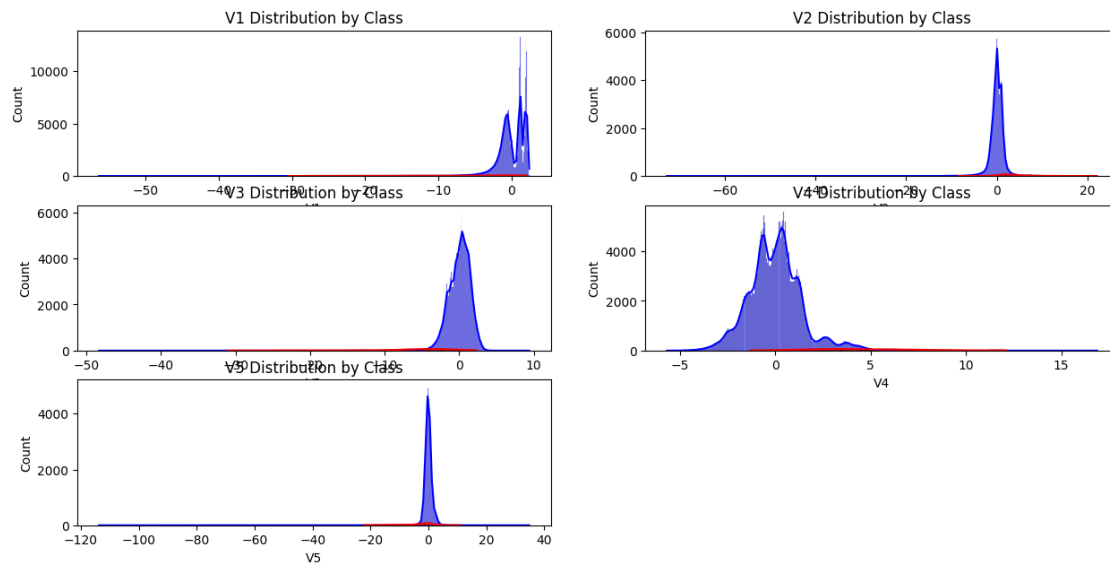
For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>



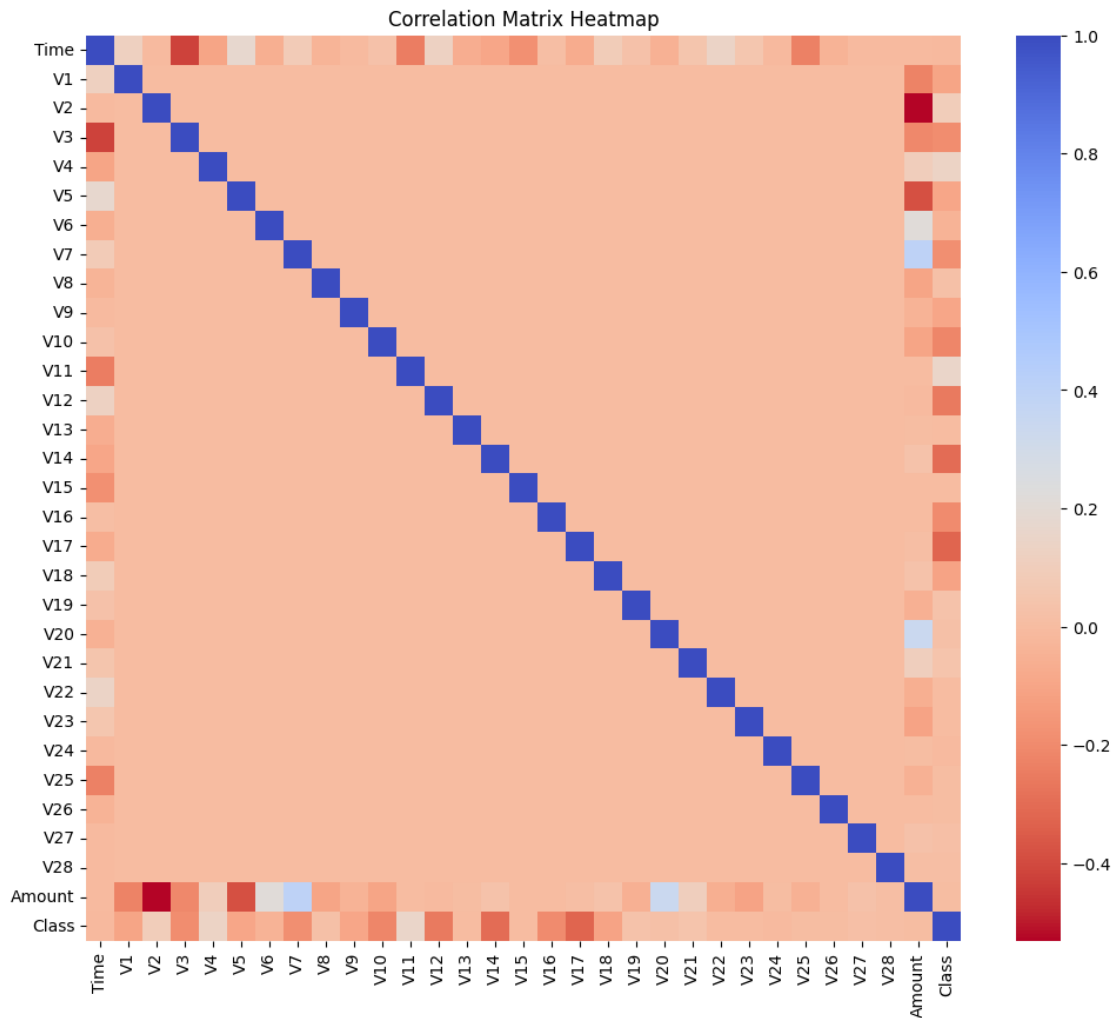
```
[ ]: # Visualizing the distribution of transaction amounts for fraudulent and
      ↪ genuine transactions in more detail
plt.figure(figsize=(12, 8))
plt.subplot(2, 1, 1)
sns.histplot(data[data['Class'] == 0]["Amount"], color='b', kde=True)
plt.title('Normal Transactions Amount Distribution')
plt.subplot(2, 1, 2)
sns.histplot(data[data['Class'] == 1]["Amount"], color='r', kde=True)
plt.title('Fraud Transactions Amount Distribution')
plt.show()
```



```
[ ]: # Analyzing the distribution of other features for fraudulent and genuine
      ↪ transactions (e.g., V1, V2, etc.)
features = ['V1', 'V2', 'V3', 'V4', 'V5']
plt.figure(figsize=(15, 35))
for i, feature in enumerate(features, 1):
    plt.subplot(14, 2, i)
    sns.histplot(data[data['Class'] == 0][feature], color='b', kde=True)
    sns.histplot(data[data['Class'] == 1][feature], color='r', kde=True)
    plt.title(f'{feature} Distribution by Class')
plt.show()
```



```
[ ]: # Analyzing the correlation between features using a heatmap
plt.figure(figsize=(12, 10))
corr = data.corr()
sns.heatmap(corr, cmap='coolwarm_r', annot_kws={'size': 10})
plt.title('Correlation Matrix Heatmap')
plt.show()
```



Density-Based Plots:

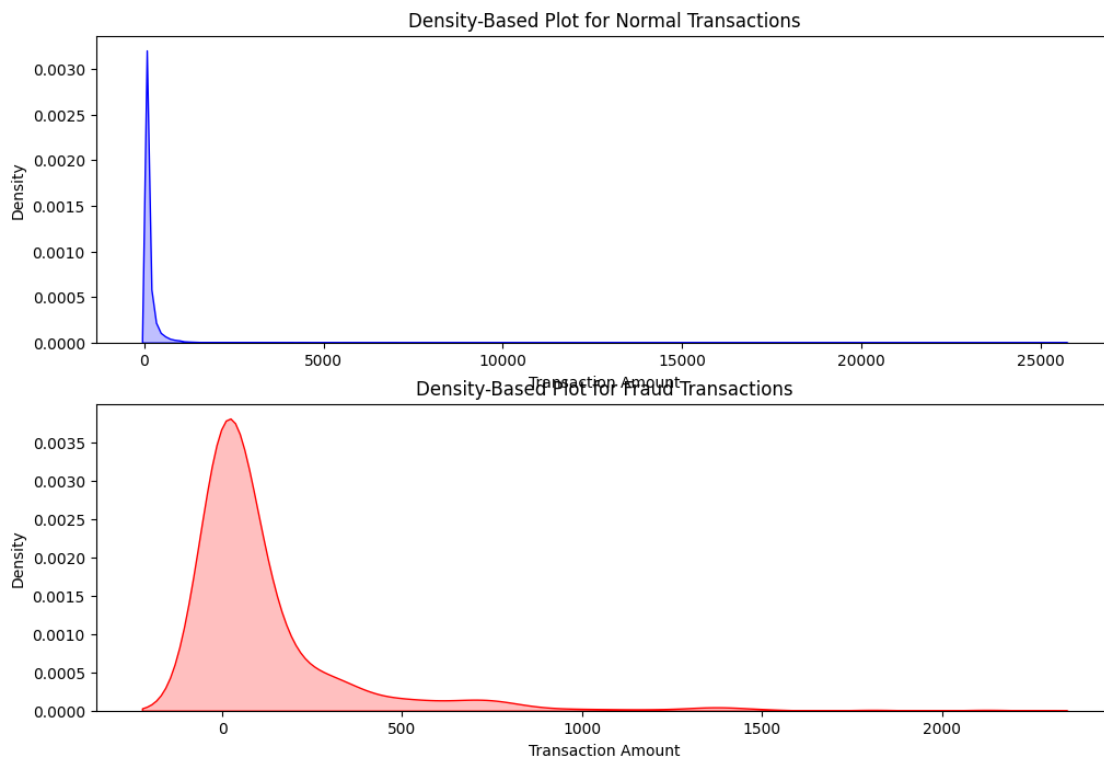
Density-based visual analysis of fraud and genuine transactions can help you understand transaction densities and trends more effectively.

```
[ ]: # Density-based plots for fraud and genuine transactions
plt.figure(figsize=(12, 8))
plt.subplot(2, 1, 1)
sns.kdeplot(data[data['Class'] == 0]["Amount"], shade=True, color='b',
            label='Normal Transactions')
plt.title('Density-Based Plot for Normal Transactions')
plt.xlabel('Transaction Amount')
plt.ylabel('Density')
plt.subplot(2, 1, 2)
sns.kdeplot(data[data['Class'] == 1]["Amount"], shade=True, color='r',
            label='Fraud Transactions')
```

```
plt.title('Density-Based Plot for Fraud Transactions')
plt.xlabel('Transaction Amount')
plt.ylabel('Density')
plt.show()
```

`shade` is now deprecated in favor of `fill`; setting `fill=True`.  
This will become an error in seaborn v0.14.0; please update your code.

`shade` is now deprecated in favor of `fill`; setting `fill=True`.  
This will become an error in seaborn v0.14.0; please update your code.

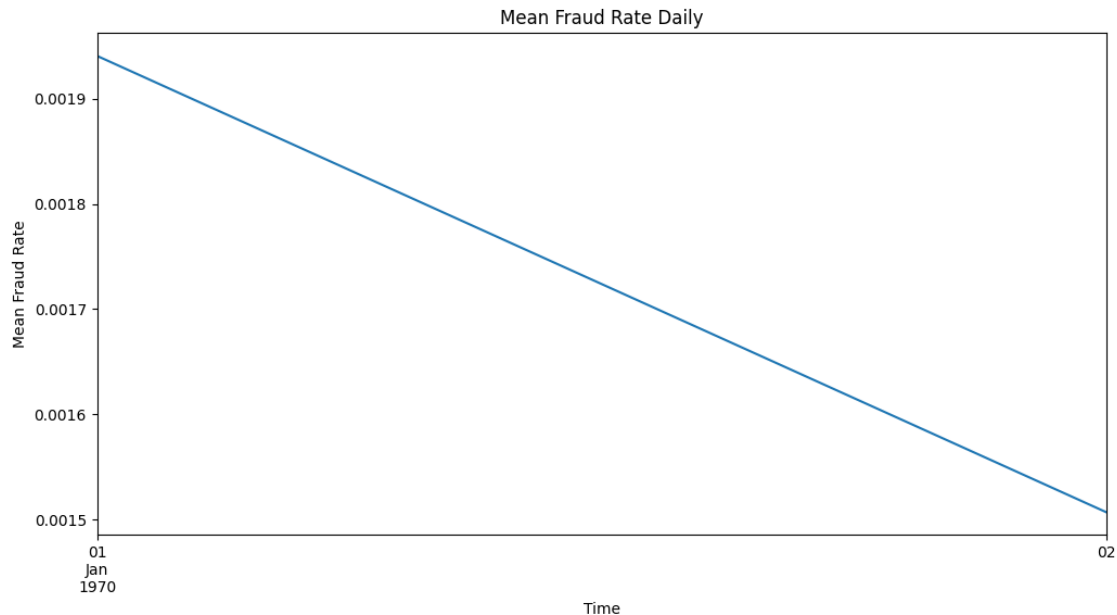


Time Series Analysis:

Conduct time series analysis to understand the trends of fraud cases over time.

```
[ ]: # Time series analysis for fraud cases
plt.figure(figsize=(12, 6))
data['Time'] = pd.to_datetime(data['Time'], unit='s')
data.set_index('Time', inplace=True)
```

```
data['Class'].resample('D').mean().plot()
plt.title('Mean Fraud Rate Daily')
plt.xlabel('Time')
plt.ylabel('Mean Fraud Rate')
plt.show()
```



### Statistical Tests:

Perform statistical tests to determine if there are statistically significant differences between fraud and normal transactions.

```
[ ]: # Performing t-test for transaction amounts between fraud and normal
      ↪ transactions
normal_transactions = data[data['Class'] == 0]['Amount']
fraud_transactions = data[data['Class'] == 1]['Amount']
t_stat, p_val = ttest_ind(normal_transactions, fraud_transactions)
print(f"T-statistic: {t_stat}, P-value: {p_val}")
```

T-statistic: -3.00555231397141, P-value: 0.002651220649191683

### Anomaly Detection Models:

Develop more advanced anomaly detection models using machine learning for fraud detection.

```
[ ]: # Implementing Isolation Forest for anomaly detection
X = data.drop('Class', axis=1)
y = data['Class']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
      ↪ random_state=42)
```

```
[ ]: # Training the model
model = IsolationForest(contamination=0.01, random_state=42)
model.fit(X_train)
```

```
[ ]: IsolationForest(contamination=0.01, random_state=42)
```

```
[ ]: # Predicting on the test set
y_pred = model.predict(X_test)
```

```
[ ]: # Generating classification report
print("Classification Report for Anomaly Detection Model:")
print(classification_report(y_test, y_pred))
```

```
Classification Report for Anomaly Detection Model:
              precision    recall  f1-score   support

     -1         0.00         0.00         0.00         0
         0         0.00         0.00         0.00    56864
         1         0.00         0.51         0.00         98

 accuracy                   0.00    56962
 macro avg              0.00    0.17    0.00    56962
 weighted avg           0.00    0.00    0.00    56962
```

Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero\_division` parameter to control this behavior.  
Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero\_division` parameter to control this behavior.  
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero\_division` parameter to control this behavior.  
Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero\_division` parameter to control this behavior.  
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero\_division` parameter to control this behavior.  
Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use `zero\_division` parameter to control this behavior.

Firewall Analysis:

Conduct firewall analysis to understand how credit card transactions behave within the firewall and identify fraud cases.

```
[ ]: # Conducting firewall analysis for credit card transactions
firewall_data = data[data['Amount'] > 1000] # Example threshold for suspicious_
↳ transactions
fraudulent_firewall_transactions = firewall_data[firewall_data['Class'] == 1]
print("Fraudulent Transactions within Firewall:")
print(fraudulent_firewall_transactions)
```



# Fraudulent Transactions within Firewall:

		V1	V2	V3	V4	V5	\
Time							
1970-01-01	02:31:04	-3.499108	0.258555	-4.489558	4.853894	-6.974522	
1970-01-01	05:01:28	-12.224021	3.854150	-12.466766	9.648311	-2.726961	
1970-01-01	16:23:31	-2.326922	-3.348439	-3.513408	3.175060	-2.815137	
1970-01-01	17:21:07	-5.344665	-0.285760	-3.835616	5.337048	-7.609909	
1970-01-01	18:09:45	-2.923827	1.524837	-3.018758	3.289291	-5.755542	
1970-01-02	10:03:28	-2.003460	-7.159042	-4.050976	1.309580	-2.058102	
1970-01-02	12:59:44	-1.212682	-2.484824	-6.397186	3.670562	-0.863375	
1970-01-02	18:51:18	-1.600211	-3.488130	-6.459303	3.246816	-1.614608	
1970-01-02	18:51:49	-0.082983	-3.935919	-2.616709	0.163310	-1.400952	

		V6	V7	V8	V9	V10	...	\
Time								
1970-01-01	02:31:04	3.628382	5.431271	-1.946734	-0.775680	-1.987773	...	
1970-01-01	05:01:28	-4.445610	-21.922811	0.320792	-4.433162	-11.201400	...	
1970-01-01	16:23:31	-0.203363	-0.892144	0.333226	-0.802005	-4.350685	...	
1970-01-01	17:21:07	3.874668	1.289630	0.201742	-3.003532	-3.990551	...	
1970-01-01	18:09:45	2.218276	-0.509995	-3.569444	-1.016592	-4.320536	...	
1970-01-02	10:03:28	-0.098621	2.880083	-0.727484	1.460381	-1.531608	...	
1970-01-02	12:59:44	-1.855855	1.017732	-0.544704	-1.703378	-3.739659	...	
1970-01-02	18:51:18	-1.260375	0.288223	-0.048964	-0.734975	-4.441484	...	
1970-01-02	18:51:49	-0.809419	1.501580	-0.471000	1.519743	-1.134454	...	

		V21	V22	V23	V24	V25	\
Time							
1970-01-01	02:31:04	-1.052368	0.204817	-2.119007	0.170279	-0.393844	
1970-01-01	05:01:28	-1.159830	-1.504119	-19.254328	0.544867	-4.781606	
1970-01-01	16:23:31	1.226648	-0.695902	-1.478490	-0.061553	0.236155	
1970-01-01	17:21:07	0.276011	1.342045	-1.016579	-0.071361	-0.335869	
1970-01-01	18:09:45	-0.511657	-0.122724	-4.288639	0.563797	-0.949451	
1970-01-02	10:03:28	1.244287	-1.015232	-1.800985	0.657586	-0.435617	
1970-01-02	12:59:44	1.396872	0.092073	-1.492882	-0.204227	0.532511	
1970-01-02	18:51:18	1.191175	-0.967141	-1.463421	-0.624231	-0.176462	
1970-01-02	18:51:49	0.702672	-0.182305	-0.921017	0.111635	-0.071622	

		V26	V27	V28	Amount	Class
Time						
1970-01-01	02:31:04	0.296367	1.985913	-0.900452	1809.68	1
1970-01-01	05:01:28	-0.007772	3.052358	-0.775036	1218.89	1
1970-01-01	16:23:31	0.531911	0.302324	0.536375	1389.56	1
1970-01-01	17:21:07	0.441044	1.520613	-1.115937	1402.16	1
1970-01-01	18:09:45	-0.204532	1.510206	-0.324706	1354.25	1
1970-01-02	10:03:28	-0.894509	-0.397557	0.314262	2125.87	1
1970-01-02	12:59:44	-0.293871	0.212663	0.431095	1335.00	1
1970-01-02	18:51:18	0.400348	0.152947	0.477775	1504.93	1
1970-01-02	18:51:49	-1.125881	-0.170947	0.126221	1096.99	1

[9 rows x 30 columns]

### Hyperparameter Tuning:

Explanation: In this section, we use GridSearchCV to find the best combination of hyperparameters for the Logistic Regression model.

```
[ ]: param_grid = {'C': [0.001, 0.01, 0.1, 1, 10], 'penalty': ['l2']}  
      solver = 'liblinear'  
  
[ ]: # Scale the data  
      scaler = StandardScaler()  
      X_train_scaled = scaler.fit_transform(X_train)  
      X_test_scaled = scaler.transform(X_test)  
  
[ ]: # Initialize the GridSearchCV  
      grid_search = GridSearchCV(LogisticRegression(solver=solver, max_iter=1000),  
      ↪ param_grid, cv=5)  
      grid_search.fit(X_train_scaled, y_train)  
  
      best_params = grid_search.best_params_  
      print("Best parameter combinations: ", best_params)
```

Best parameter combinations: {'C': 10, 'penalty': 'l2'}

### Data Preprocessing Techniques:

Explanation: This section involves standard scaling of the data and the use of SMOTE to address class imbalance issues.

```
[ ]: # Apply standard scaling to the data  
      scaler = StandardScaler()  
      X_train_scaled = scaler.fit_transform(X_train)  
      X_test_scaled = scaler.transform(X_test)  
  
[ ]: # Implement SMOTE for handling class imbalance  
      smote = SMOTE(random_state=42)  
      X_resampled, y_resampled = smote.fit_resample(X_train_scaled, y_train)  
  
[ ]: # Display the results  
      print("Original data shape:", X_train.shape, y_train.shape)  
      print("Resampled data shape:", X_resampled.shape, y_resampled.shape)
```

Original data shape: (227845, 29) (227845,)

Resampled data shape: (454902, 29) (454902,)

### Trying Different Models:

Explanation: Here, we utilize the XGBoost model, train it on the resampled data, and evaluate its performance using the classification report.

```
[ ]: from xgboost import XGBClassifier
```

```
[ ]: # Training the XGBoost model
xgb_model = XGBClassifier()
xgb_model.fit(X_resampled, y_resampled)
y_pred_xgb = xgb_model.predict(X_test_scaled)
```

```
[ ]: # Evaluating the performance of the XGBoost model
print("Classification Report for XGBoost Model:")
print(classification_report(y_test, y_pred_xgb))
```

Classification Report for XGBoost Model:

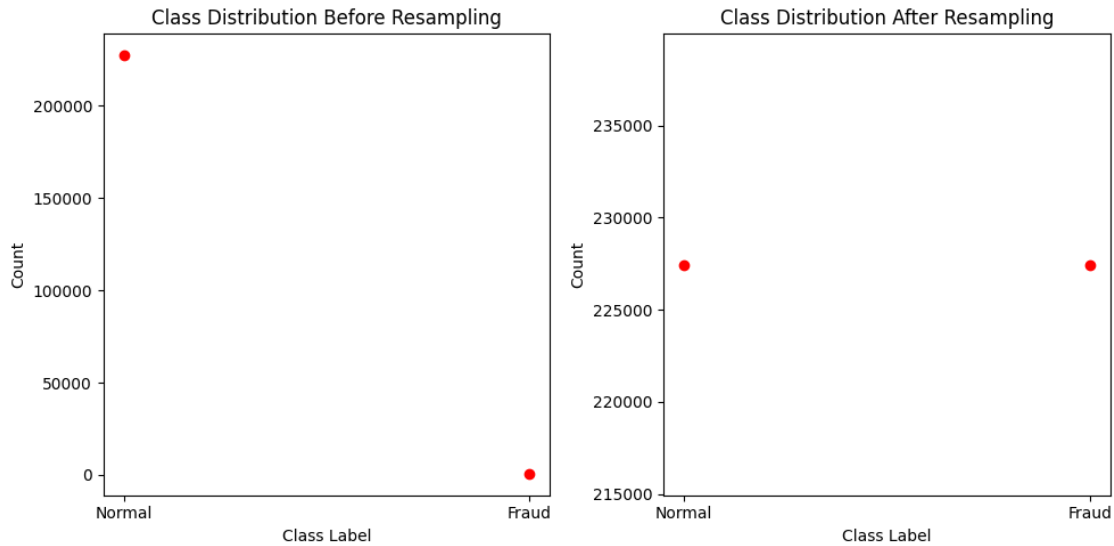
	precision	recall	f1-score	support
0	1.00	1.00	1.00	56864
1	0.80	0.84	0.82	98
accuracy			1.00	56962
macro avg	0.90	0.92	0.91	56962
weighted avg	1.00	1.00	1.00	56962

```
[ ]: # Visualize the class distribution before and after resampling
plt.figure(figsize=(10, 5))

# Dot plot for class distribution before resampling
plt.subplot(1, 2, 1)
plt.title('Class Distribution Before Resampling')
plt.plot([0, 1], [sum(y_train==0), sum(y_train==1)], 'ro')
plt.xticks([0, 1], ['Normal', 'Fraud'])
plt.xlabel('Class Label')
plt.ylabel('Count')

# Dot plot for class distribution after resampling
plt.subplot(1, 2, 2)
plt.title('Class Distribution After Resampling')
plt.plot([0, 1], [sum(y_resampled==0), sum(y_resampled==1)], 'ro')
plt.xticks([0, 1], ['Normal', 'Fraud'])
plt.xlabel('Class Label')
plt.ylabel('Count')

plt.tight_layout()
plt.show()
```



Data Augmentation:

Explanation: This section demonstrates the implementation of data augmentation techniques using Random Over Sampling to balance the dataset.

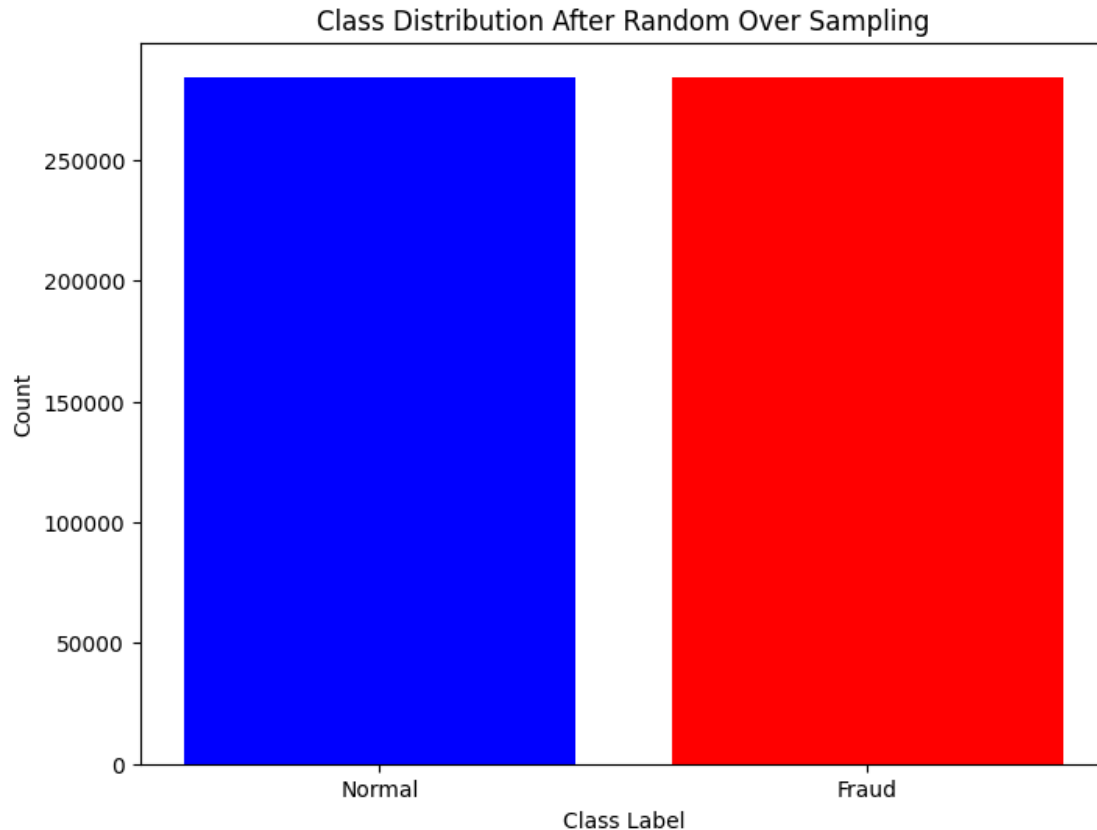
```
[ ]: # Using Random Over Sampling for data augmentation
ros = RandomOverSampler(random_state=0)
X_resampled_aug, y_resampled_aug = ros.fit_resample(X, y)
```

```
[ ]: # Display the results
print("Original dataset shape:", Counter(y))
print("Resampled dataset shape:", Counter(y_resampled_aug))
```

Original dataset shape: Counter({0: 284315, 1: 492})

Resampled dataset shape: Counter({0: 284315, 1: 284315})

```
[ ]: # Visualize the class distribution after Random Over Sampling
plt.figure(figsize=(8, 6))
plt.bar(Counter(y_resampled_aug).keys(), Counter(y_resampled_aug).values(),
        color=['b', 'r'])
plt.xticks(list(Counter(y_resampled_aug).keys()), ['Normal', 'Fraud'])
plt.xlabel('Class Label')
plt.ylabel('Count')
plt.title('Class Distribution After Random Over Sampling')
plt.show()
```



Model Evaluation Metrics:

Explanation: Here, we compute and print the precision, recall, and F1 scores to evaluate the model's performance.

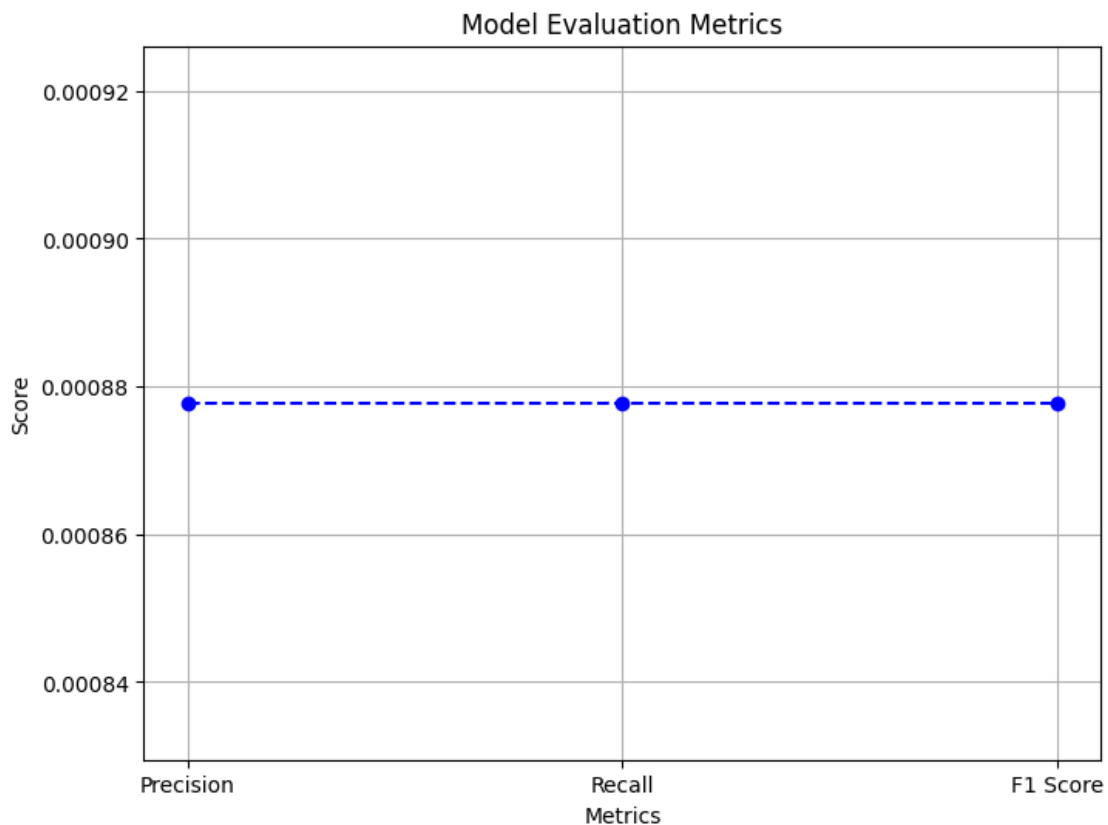
```
[ ]: # Calculate precision, recall, and F1 scores
precision = precision_score(y_test, y_pred, average='micro')
recall = recall_score(y_test, y_pred, average='micro')
f1 = f1_score(y_test, y_pred, average='micro')

print("Precision: ", precision)
print("Recall: ", recall)
print("F1 Score: ", f1)
```

```
Precision:  0.0008777781679014079
Recall:  0.0008777781679014079
F1 Score:  0.0008777781679014079
```

```
[ ]: # Defining the metrics and scores
metrics = ['Precision', 'Recall', 'F1 Score']
scores = [precision, recall, f1]
```

```
[ ]: # Creating a dot plot
plt.figure(figsize=(8, 6))
plt.plot(metrics, scores, marker='o', linestyle='--', color='b')
plt.title('Model Evaluation Metrics')
plt.xlabel('Metrics')
plt.ylabel('Score')
plt.grid(True)
plt.show()
```



Handling Missing Data:

Explanation: This section involves checking the dataset for any missing values to ensure data integrity and quality.

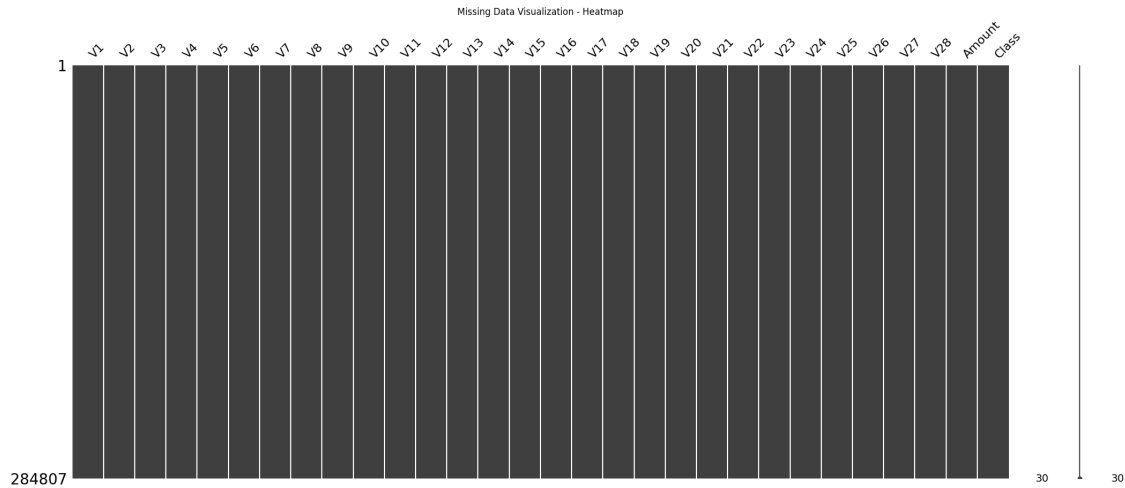
```
[ ]: # Check for missing data
missing_values = data.isnull().sum()
print("Missing values: ", missing_values)
```

```
Missing values:  V1      0
V2      0
V3      0
V4      0
```

```
V5          0
V6          0
V7          0
V8          0
V9          0
V10         0
V11         0
V12         0
V13         0
V14         0
V15         0
V16         0
V17         0
V18         0
V19         0
V20         0
V21         0
V22         0
V23         0
V24         0
V25         0
V26         0
V27         0
V28         0
Amount      0
Class       0
dtype: int64
```

```
[ ]: import missingno as msno

# Visualizing missing data using a heatmap
msno.matrix(data)
plt.title('Missing Data Visualization - Heatmap')
plt.show()
```



Model Interpretation:

Explanation: Finally, this section demonstrates the use of SHAP (SHapley Additive exPlanations) values for interpreting the model's predictions and understanding the impact of various features on the model's output.

```
[ ]: # Use SHAP values for model interpretation
explainer = shap.Explainer(model)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test, plot_type="bar")
```



