

```
In [1]: # Import the necessary packages
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler
from matplotlib import gridspec
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_curve, roc_auc_score, auc
```

Loading the Data

```
In [2]: # Load the dataset from the csv file using pandas
# best way is to mount the drive on colab and
# copy the path for the csv file
data = pd.read_csv("creditcard.csv")

In [3]: # Grab a peek at the data
data.head()
```

| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | ... | V21 | V22 | V23 | V24 | V25 | V26 | V27 | V28 | Amount | Class |
|---|------|-----------|-----------|----------|----------|-----------|-----------|-----------|-----------|-----------|-----|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|--------|-------|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.402388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 | -0.110474 | 0.066928 | 0.128539 | -0.189115 | 0.133558 | -0.021053 | 149.62 | 0 |
| 1 | 0.0 | 1.191857 | 0.266151 | 0.166400 | 0.448154 | 0.060318 | -0.082361 | -0.078803 | 0.085102 | -0.255425 | ... | -0.225775 | -0.638672 | 0.101288 | -0.339846 | 0.167170 | 0.125895 | -0.008983 | 0.014724 | 2.69 | 0 |
| 2 | 1.0 | -1.359354 | -1.340163 | 1.773209 | 0.376780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 | -1.514654 | ... | 0.247998 | 0.771679 | -0.909412 | -0.689281 | -0.327642 | -0.139097 | -0.055353 | -0.059752 | 378.66 | 0 |
| 3 | 1.0 | -1.096272 | -0.185226 | 1.792993 | 0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 | -1.387024 | ... | -0.108300 | 0.005744 | -0.909321 | -1.175575 | 0.647376 | -0.221929 | 0.062723 | 0.061458 | 123.90 | 0 |
| 4 | 2.0 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 | 0.817739 | ... | -0.009431 | 0.798278 | -0.137458 | 0.141267 | -0.206010 | 0.502292 | 0.219422 | 0.215153 | 69.99 | 0 |

5 rows x 31 columns

Exploitory Data Analysis

```
In [4]: # Print the shape of the data
data = data.sample(frac = 0.1, random_state = 48)
print(data.shape)
print(data.describe())

(284807, 31)
```

| | Time | V1 | V2 | V3 | V4 | ... | V21 | V22 | V23 | V24 | ... | V27 | V28 | Amount | Class |
|-------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| count | 284807.000000 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 284807.000000 |
| mean | 94813.859575 | 1.168375e-15 | 3.416908e-16 | -1.379537e-15 | 2.874095e-15 | 1.654867e-16 | 1.568593e-16 | 2.578648e-16 | 4.473266e-15 | 1.390247e+00 | 1.322771e+00 | 1.237094e+00 | 1.194253e+00 | 1.098632e+00 | 1.415669e+00 |
| std | 47488.145955 | 1.958696e+00 | 1.651309e+00 | 1.516255e+00 | 1.415669e+00 | 1.651309e+00 | 1.516255e+00 | 1.415669e+00 | 1.651309e+00 | 1.516255e+00 | 1.415669e+00 | 1.516255e+00 | 1.415669e+00 | 1.516255e+00 | 1.415669e+00 |
| min | 0.000000 | -5.640751e+01 | -7.271873e+01 | -4.832558e+01 | 5.683771e+00 | -5.640751e+01 | -7.271873e+01 | -4.832558e+01 | 5.683771e+00 | -5.640751e+01 | -7.271873e+01 | -4.832558e+01 | 5.683771e+00 | -5.640751e+01 | -7.271873e+01 |
| 25% | 54201.500000 | -9.203734e-01 | -5.985499e-01 | -8.093648e-01 | -8.486401e-01 | -9.203734e-01 | -5.985499e-01 | -8.093648e-01 | -8.486401e-01 | -9.203734e-01 | -5.985499e-01 | -8.093648e-01 | -8.486401e-01 | -9.203734e-01 | -5.985499e-01 |
| 50% | 84692.000000 | 1.818898e-02 | 5.548556e-02 | 1.799463e-01 | 1.984653e-02 | 1.818898e-02 | 5.548556e-02 | 1.799463e-01 | 1.984653e-02 | 1.818898e-02 | 5.548556e-02 | 1.799463e-01 | 1.984653e-02 | 1.818898e-02 | 5.548556e-02 |
| 59% | 159329.000000 | 1.335642e+00 | 6.837239e-01 | 1.027196e+00 | 7.433433e-01 | 1.335642e+00 | 6.837239e-01 | 1.027196e+00 | 7.433433e-01 | 1.335642e+00 | 6.837239e-01 | 1.027196e+00 | 7.433433e-01 | 1.335642e+00 | 6.837239e-01 |
| max | 172792.000000 | 2.454938e+00 | 2.285773e+01 | 9.382558e+00 | 1.687534e+01 | 2.454938e+00 | 2.285773e+01 | 9.382558e+00 | 1.687534e+01 | 2.454938e+00 | 2.285773e+01 | 9.382558e+00 | 1.687534e+01 | 2.454938e+00 | 2.285773e+01 |

| | V5 | V6 | V7 | V8 | V9 | ... | V21 | V22 | V23 | V24 | ... | V27 | V28 | Amount | Class |
|-------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| count | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 284807.000000 |
| mean | 9.604066e-16 | 1.487313e-15 | -5.556487e-16 | 1.213481e-16 | -2.486331e-15 | 1.487313e-15 | -5.556487e-16 | 1.213481e-16 | -2.486331e-15 | 1.487313e-15 | -5.556487e-16 | 1.213481e-16 | -2.486331e-15 | 1.487313e-15 | -5.556487e-16 |
| std | 1.380247e+00 | 1.322771e+00 | 1.237094e+00 | 1.194253e+00 | 1.098632e+00 | 1.380247e+00 | 1.322771e+00 | 1.237094e+00 | 1.194253e+00 | 1.098632e+00 | 1.380247e+00 | 1.322771e+00 | 1.237094e+00 | 1.194253e+00 | 1.098632e+00 |
| min | -1.137433e+02 | -2.616051e+01 | -4.355724e+01 | -7.321672e+01 | -1.343407e+01 | -1.137433e+02 | -2.616051e+01 | -4.355724e+01 | -7.321672e+01 | -1.343407e+01 | -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 | -6.438976e-01 | -6.915971e-01 | -7.682956e-01 | -5.540759e-01 | -2.086297e-01 | -6.438976e-01 | -6.915971e-01 | -7.682956e-01 | -5.540759e-01 | -2.086297e-01 | -6.438976e-01 |
| 50% | -5.423583e-02 | -2.741871e-01 | -4.018388e-02 | 2.235984e-02 | -5.142879e-02 | -5.423583e-02 | -2.741871e-01 | -4.018388e-02 | 2.235984e-02 | -5.142879e-02 | -5.423583e-02 | -2.741871e-01 | -4.018388e-02 | 2.235984e-02 | -5.142879e-02 |
| 75% | 6.119254e-01 | 3.985649e-01 | 5.794361e-01 | 3.273459e-01 | 5.971390e-01 | 6.119254e-01 | 3.985649e-01 | 5.794361e-01 | 3.273459e-01 | 5.971390e-01 | 6.119254e-01 | 3.985649e-01 | 5.794361e-01 | 3.273459e-01 | 5.971390e-01 |
| max | 3.480167e+01 | 7.330163e+01 | 1.205895e+02 | 2.080721e+01 | 1.559499e+01 | 3.480167e+01 | 7.330163e+01 | 1.205895e+02 | 2.080721e+01 | 1.559499e+01 | 3.480167e+01 | 7.330163e+01 | 1.205895e+02 | 2.080721e+01 | 1.559499e+01 |

| | V21 | V22 | V23 | V24 | ... | V27 | V28 | Amount | Class |
|-------|---------------|---------------|---------------|---------------|--------------|---------------|---------------|---------------|---------------|
| count | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 284807.000000 | 284807.000000 |
| mean | 5.340915e-16 | 1.683437e-15 | -3.668891e-16 | 1.227396e-16 | 88.340915 | 5.340915e-16 | 1.683437e-15 | 1.227396e-16 | 88.340915 |
| std | 5.212781e-01 | 4.822270e-01 | 4.036325e-01 | 3.300833e-01 | 250.120109 | 5.212781e-01 | 4.822270e-01 | 4.036325e-01 | 250.120109 |
| min | -1.029546e-01 | -2.694551e+00 | -2.256568e+01 | -1.543088e+01 | 0.000000 | -1.029546e-01 | -2.694551e+00 | -2.256568e+01 | -1.543088e+01 |
| 25% | -3.171451e-01 | 3.269893e-01 | -7.083953e-02 | 5.295979e-02 | 5.600000 | -3.171451e-01 | 3.269893e-01 | -7.083953e-02 | 5.295979e-02 |
| 50% | 1.659350e-02 | -5.213911e-02 | 1.342146e-03 | 1.124383e-02 | 22.000000 | 1.659350e-02 | -5.213911e-02 | 1.342146e-03 | 1.124383e-02 |
| 75% | 3.507158e-01 | 2.489522e-01 | 9.104512e-02 | 7.827995e-02 | 77.105000 | 3.507158e-01 | 2.489522e-01 | 9.104512e-02 | 7.827995e-02 |
| max | 7.019595e+00 | 3.517346e+00 | 3.161228e+01 | 3.384781e+01 | 25691.160000 | 7.019595e+00 | 3.517346e+00 | 3.161228e+01 | 3.384781e+01 |

| | Class |
|-------|---------------|
| count | 284807.000000 |
| mean | 0.991727 |
| std | 0.041527 |
| min | 0.000000 |
| 25% | 0.000000 |
| 50% | 0.000000 |
| 75% | 0.000000 |
| max | 1.000000 |

[8 rows x 31 columns]

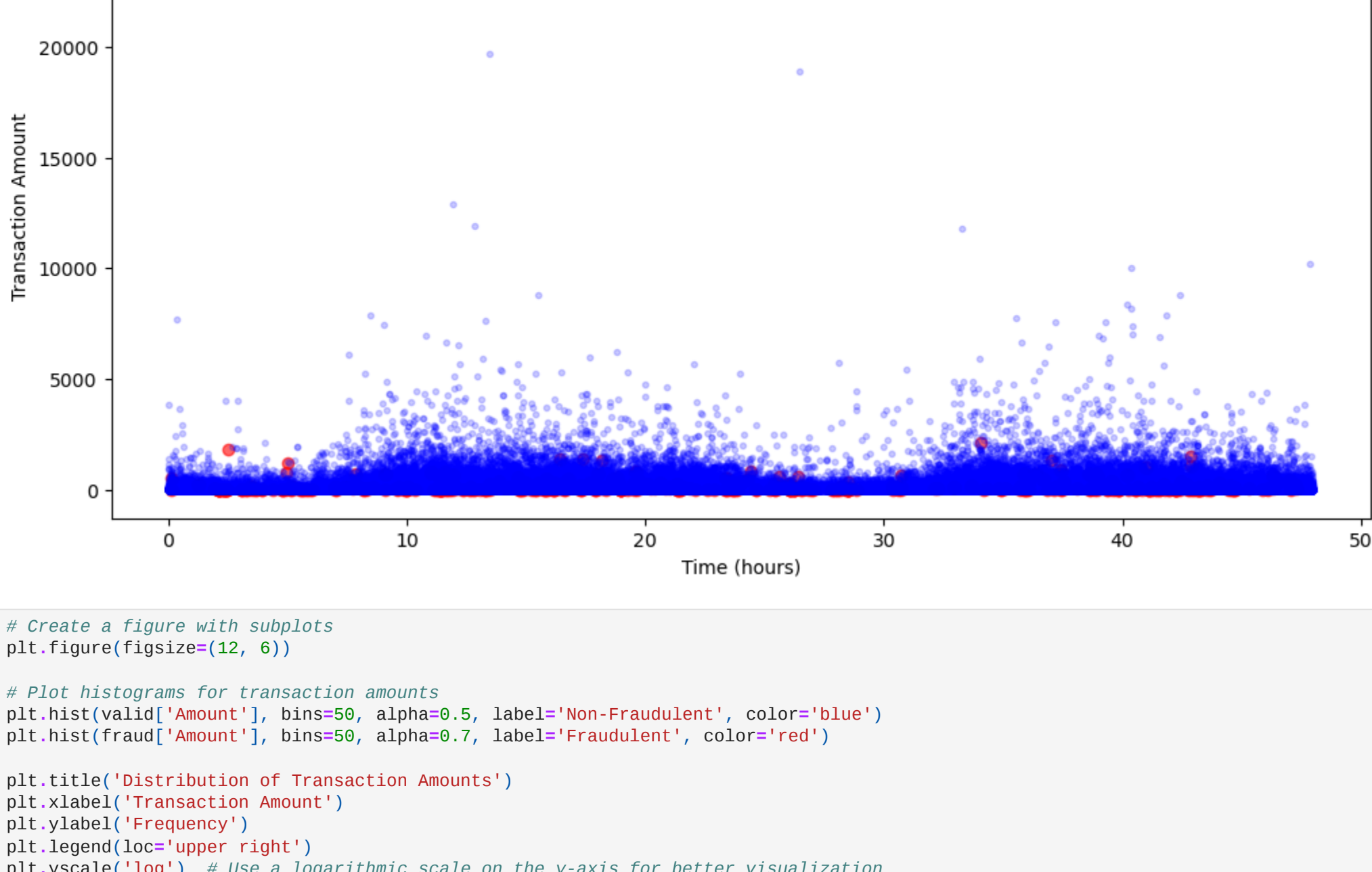
```
In [5]: # Determine number of fraud cases in dataset
data['Time'] = data['Time'] / 3600
fraud = data[data['Class'] == 1]
valid = data[data['Class'] == 0]
outlierFraction = len(fraud)/float(len(valid))
print(outlierFraction)
print('Fraud Cases: {}'.format(len(data[data['Class'] == 1])))
print('Valid Transactions: {}'.format(len(data[data['Class'] == 0])))

0.0017384750013189597
Fraud Cases: 492
Valid Transactions: 284315
```

Visualization

```
In [6]: plt.figure(figsize=(12, 6))
plt.scatter(fraud['Time'], fraud['Amount'], color='red', markers='o', label='Fraud', alpha=0.6)
plt.scatter(valid['Time'], valid['Amount'], color='blue', markers='.', label='Non-Fraud', alpha=0.2)

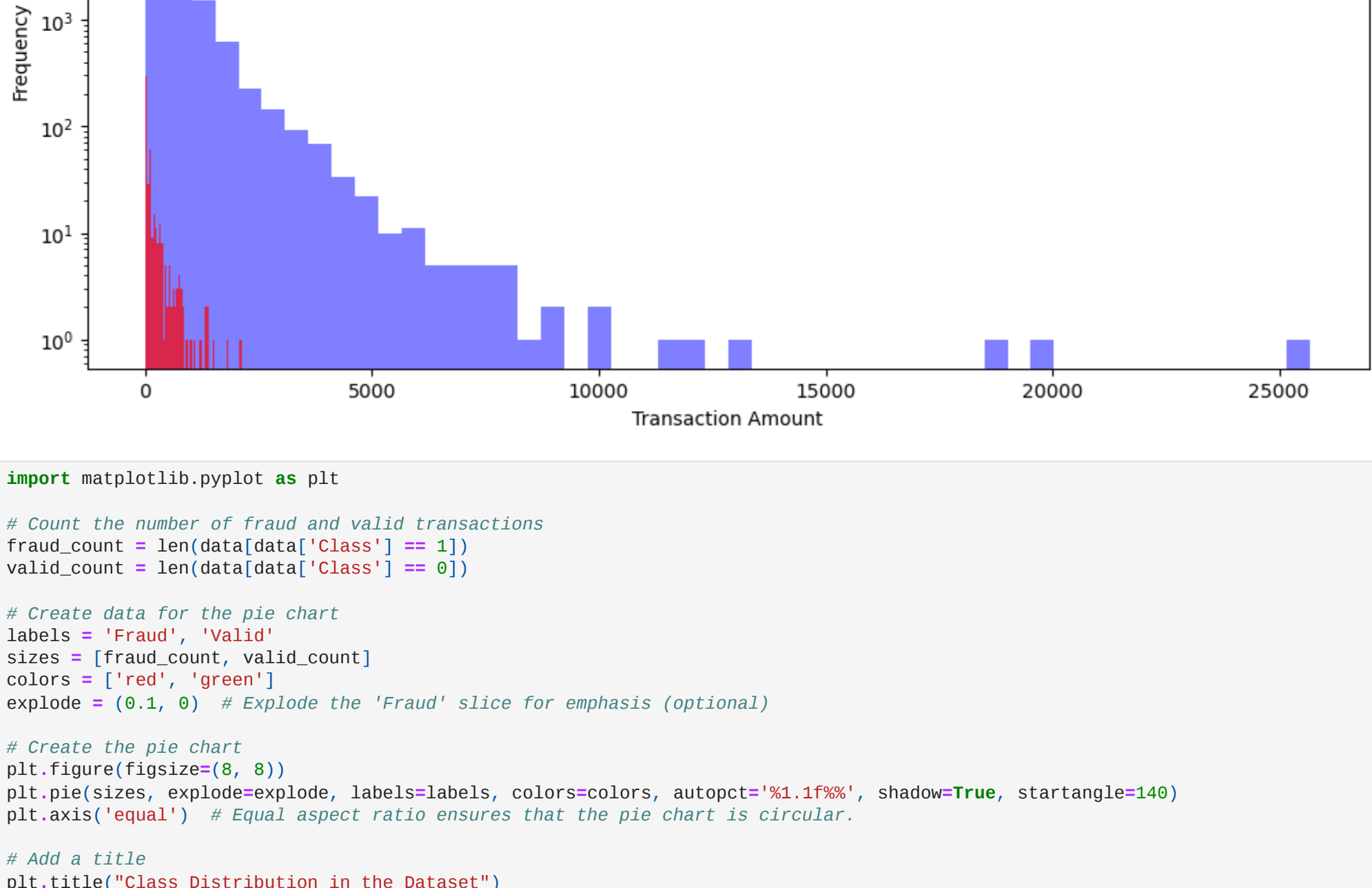
plt.title('Time Series of Fraud vs. Non-Fraud Transactions')
plt.xlabel('Time (hours)')
plt.ylabel('Transaction Amount')
plt.legend(loc='upper right')
plt.show()
```



```
In [7]: # Create a Figure with subplots
plt.figure(figsize=(12, 6))

# Plot histograms for transaction amounts
plt.hist(valid['Amount'], bins=50, alpha=0.5, label='Non-Fraudulent', color='blue')
plt.hist(fraud['Amount'], bins=50, alpha=0.7, label='Fraudulent', color='red')

plt.title('Distribution of Transaction Amounts')
plt.xlabel('Transaction Amount')
plt.ylabel('Frequency')
plt.legend(loc='upper right')
plt.yscale('log') # Use a logarithmic scale on the y-axis for better visualization
plt.show()
```



```
In [8]: import matplotlib.pyplot as plt

# Count the number of fraud and valid transactions
fraud_count = len(data[data['Class'] == 1])
valid_count = len(data[data['Class'] == 0])

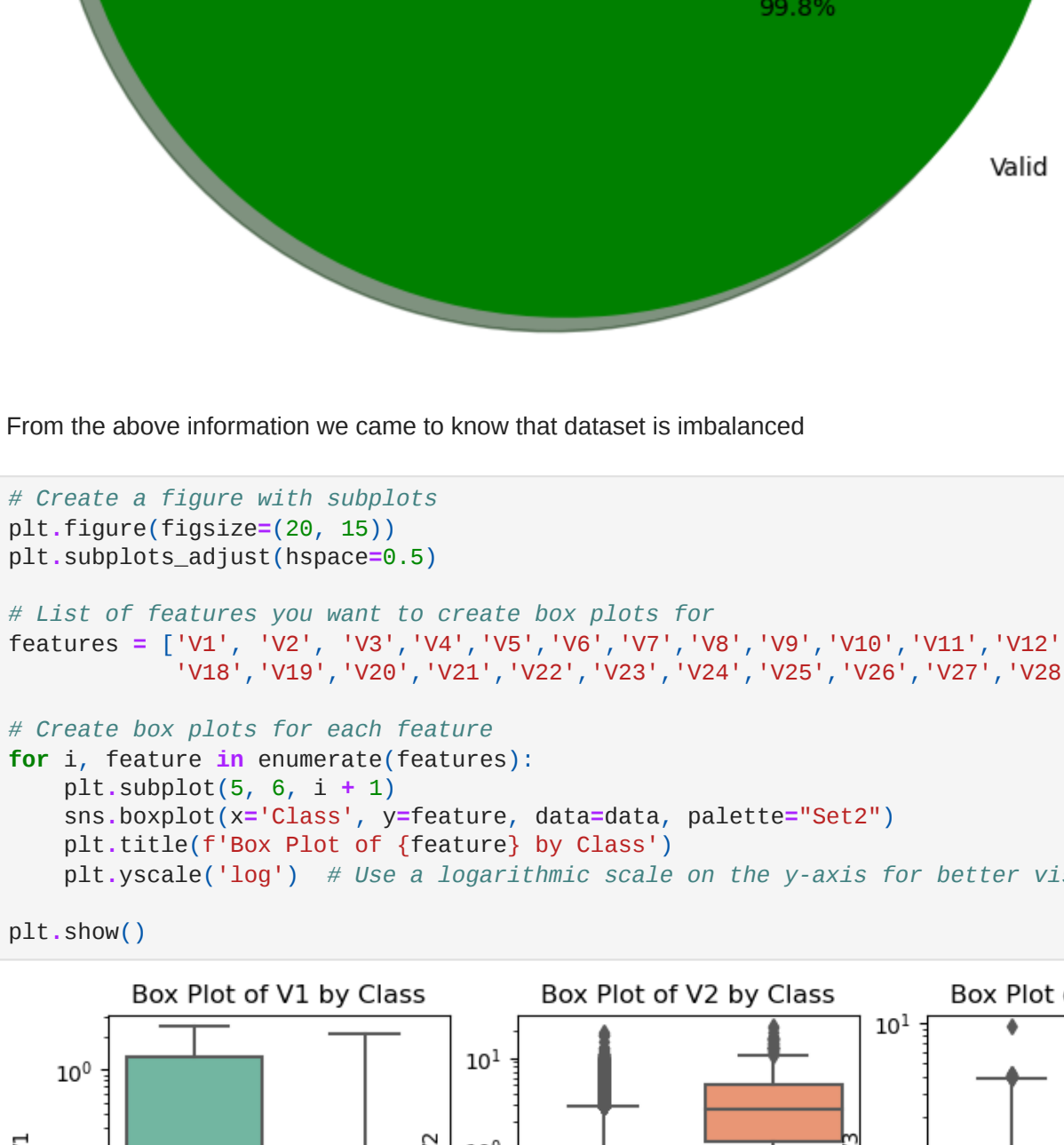
# Create data for the pie chart
labels = 'Fraud', 'Valid'
sizes = [fraud_count, valid_count]
colors = ['red', 'green']
explode = (0.1, 0) # Explode the 'Fraud' slice for emphasis (optional)

# Create the pie chart
plt.figure(figsize=(8, 8))
plt.pie(sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%', shadow=True, startangle=140)
plt.axis('equal') # Equal aspect ratio ensures that the pie chart is circular.

# Add a title
plt.title('Class Distribution in the Dataset')

# Display the pie chart
plt.show()
```

Class Distribution in the Dataset



From the above information we came to know that dataset is imbalanced

```
In [18]: # Create a Figure with subplots
plt.figure(figsize=(20, 15))
plt.subplots_adjust(hspace=0.5)

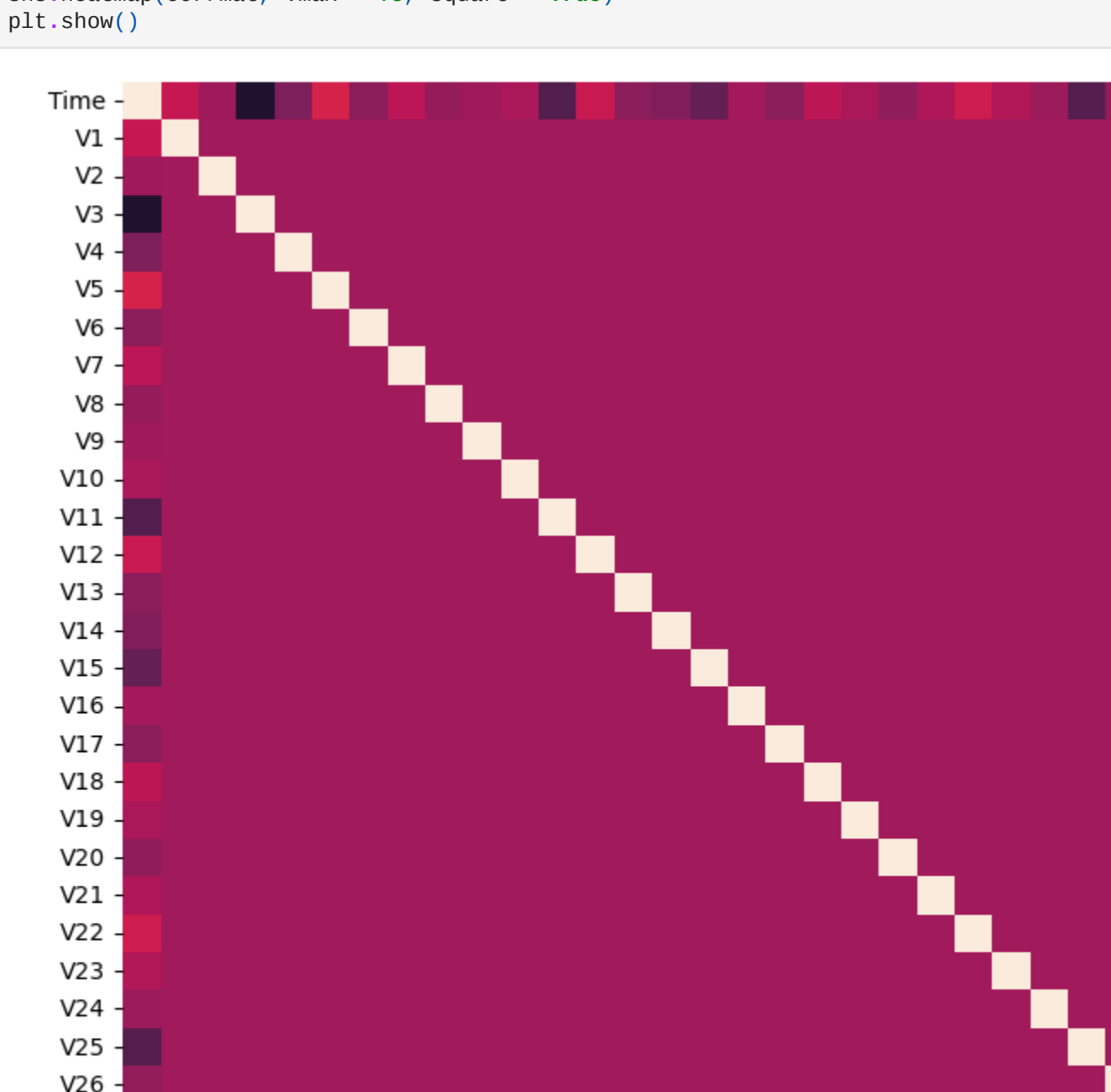
# List of features you want to create box plots for
features = ['V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20', 'V21', 'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount'] # Add more Features as needed

# Create box plots for each feature
for f, feature in enumerate(features):
    plt.subplot(5, 6, 1 + f)
    sns.boxplot(x='Class', y=feature, data=data, palette='Set2')
    plt.title('Box Plot of {} by Class'.format(feature))
    plt.yscale('log') # Use a logarithmic scale on the y-axis for better visualization

plt.show()
```



```
In [19]: # Create a boxplot for the 'Amount' column with outliers displayed
plt.figure(figsize=(8, 6))
sns.boxplot(x=data['Amount'], showfliers=True)
plt.title('Box Plot of Amount with Outliers')
plt.ylabel('Amount')
plt.show()
```



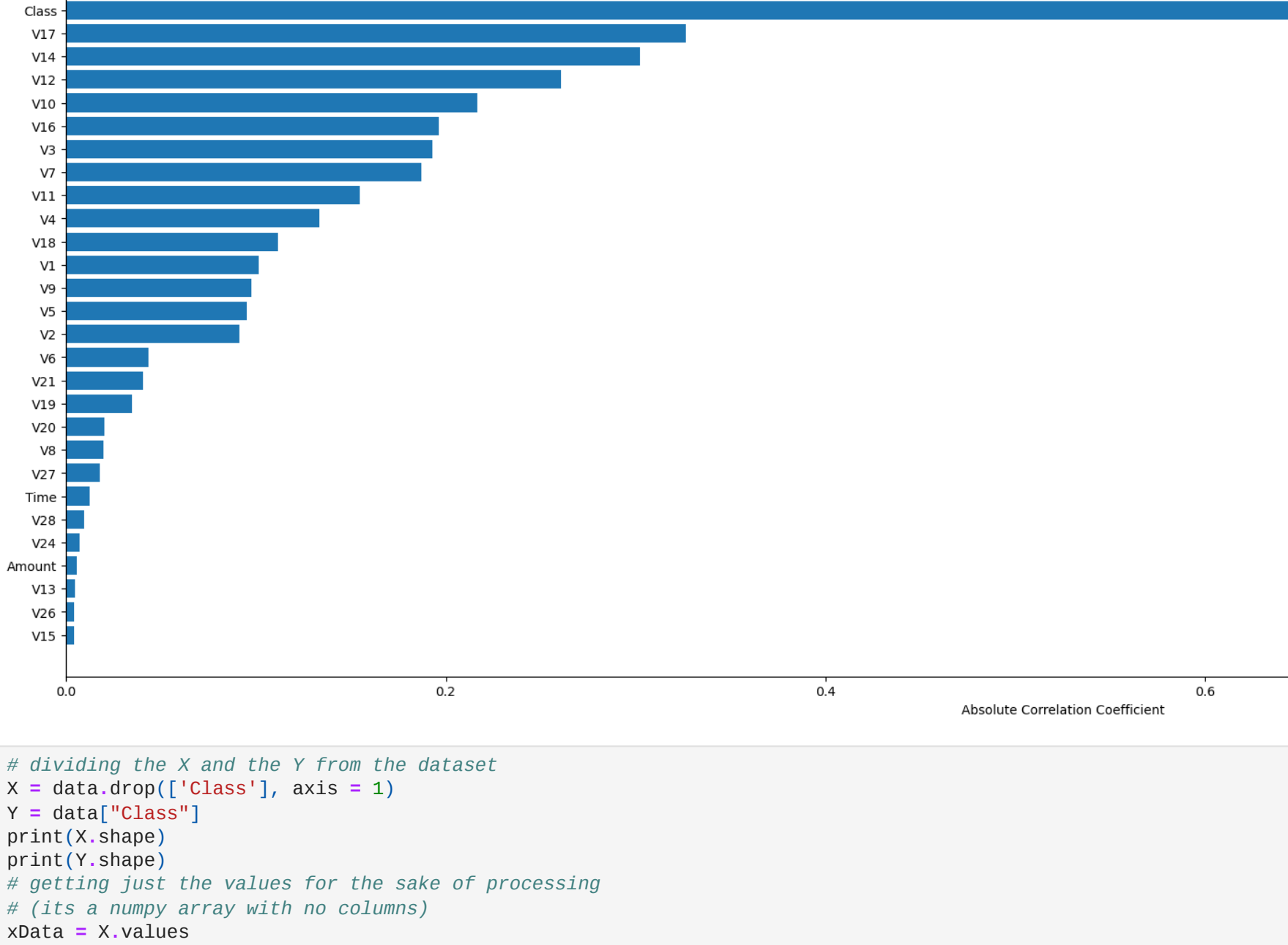
```
In [11]: #Fraudulent transaction
print("Amount details of the fraudulent transaction")
fraud.Amount.describe()

Amount details of the fraudulent transaction
count    492.000000
mean     122.211321
std      256.683288
min        0.000000
25%        1.000000
50%         9.250000
75%       105.899000
max      2125.870000
Name: Amount, dtype: float64

In [12]: #Normal Transaction
print("details of valid transaction")
valid.Amount.describe()

details of valid transaction
count    284315.000000
mean      85.231022
std       250.105992
min         0.000000
25%         5.650000
50%        22.000000
75%        77.050000
max     25691.160000
Name: Amount, dtype: float64
```

```
In [13]: # Correlation matrix
corrmat = data.corr()
fig = plt.figure(figsize=(12, 9))
sns.heatmap(corrmat, vmax = .8, square = True)
plt.show()
```



```
In [14]: # Calculate the absolute Pearson correlation coefficient between features and the target variable
correlations = data.corrwith(data['Class']).abs()

# Sort feature importances in descending order
sorted_correlations = correlations.sort_values(ascending=False)

# Get the top n important features (e.g., top 10)
top_correlations = sorted_correlations[:n]
top_feature_names = top_correlations.index

# Create a bar chart to visualize feature importance
plt.figure(figsize=(28, 10))
plt.barh(range(n), top_correlations, align='center')
plt.yticks(range(n), top_feature_names)
plt.xlabel('Absolute Correlation Coefficient')
plt.title('Top Feature Correlations with Target (Class)')
plt.gca().invert_yaxis() # Invert the y-axis for better visualization
plt.show()
```



```
In [15]: # dividing the X and the Y from the dataset
X = data.drop('Class', axis = 1)
Y = data['Class']
print(X.shape)
print(Y.shape)

# getting just the values for the sake of processing
# (its a numpy array with no columns)
xdata = X.values
ydata = Y.values

(284807, 30)
(284807,)
```

```
In [16]: # Split the data into training and testing sets
xTrain, xTest, yTrain, yTest = train_test_split(xdata, ydata, test_size = 0.2, random_state = 42)
```

```
In [ ]:
```