

sampling-anomaly-detection-eda-3

November 3, 2024

```
[1]: # Install Kaggle package
!pip install kaggle

# Make a directory for the Kaggle API key
!mkdir -p ~/.kaggle

# Move the Kaggle API key to the correct location
!cp kaggle.json ~/.kaggle/

# Set the permissions for the API key
!chmod 600 ~/.kaggle/kaggle.json
```

```
Requirement already satisfied: kaggle in /usr/local/lib/python3.10/dist-packages
(1.6.17)
Requirement already satisfied: six>=1.10 in /usr/local/lib/python3.10/dist-
packages (from kaggle) (1.16.0)
Requirement already satisfied: certifi>=2023.7.22 in
/usr/local/lib/python3.10/dist-packages (from kaggle) (2024.8.30)
Requirement already satisfied: python-dateutil in
/usr/local/lib/python3.10/dist-packages (from kaggle) (2.8.2)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-
packages (from kaggle) (2.32.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages
(from kaggle) (4.66.6)
Requirement already satisfied: python-slugify in /usr/local/lib/python3.10/dist-
packages (from kaggle) (8.0.4)
Requirement already satisfied: urllib3 in /usr/local/lib/python3.10/dist-
packages (from kaggle) (2.2.3)
Requirement already satisfied: bleach in /usr/local/lib/python3.10/dist-packages
(from kaggle) (6.2.0)
Requirement already satisfied: webencodings in /usr/local/lib/python3.10/dist-
packages (from bleach->kaggle) (0.5.1)
Requirement already satisfied: text-unidecode>=1.3 in
/usr/local/lib/python3.10/dist-packages (from python-slugify->kaggle) (1.3)
Requirement already satisfied: charset-normalizer<4,>=2 in
```

```
/usr/local/lib/python3.10/dist-packages (from requests->kaggle) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-
packages (from requests->kaggle) (3.10)
cp: cannot stat 'kaggle.json': No such file or directory
chmod: cannot access '/root/.kaggle/kaggle.json': No such file or directory
```

```
[2]: # Download the dataset from Kaggle
!kaggle datasets download -d mlg-ulb/creditcardfraud
```

```
Dataset URL: https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud
License(s): DbCL-1.0
Downloading creditcardfraud.zip to /content
 97% 64.0M/66.0M [00:02<00:00, 27.2MB/s]
100% 66.0M/66.0M [00:03<00:00, 22.3MB/s]
```

```
[3]: # Unzip the downloaded file
!unzip creditcardfraud.zip
```

```
Archive: creditcardfraud.zip
  inflating: creditcard.csv
```

```
[4]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.gridspec as gridspec

from sklearn.ensemble import IsolationForest, RandomForestClassifier
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.model_selection import KFold, cross_val_score, StratifiedKFold
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from sklearn.metrics import log_loss, roc_auc_score, f1_score, accuracy_score
from sklearn.preprocessing import StandardScaler

from xgboost import XGBClassifier

from warnings import simplefilter
simplefilter("ignore")
```

```
[5]: # Load the dataset
data = pd.read_csv('creditcard.csv')
```

```
[6]: # Displaying all the columns in the dataset
pd.set_option("display.max_columns", 1000)
```

```
data.head()
```

```
[6]:
```

	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	V10	V11	V12	V13	V14	\
0	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.311169	
1	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.143772	
2	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.165946	
3	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.287924	
4	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.119670	

	V15	V16	V17	V18	V19	V20	V21	\
0	1.468177	-0.470401	0.207971	0.025791	0.403993	0.251412	-0.018307	
1	0.635558	0.463917	-0.114805	-0.183361	-0.145783	-0.069083	-0.225775	
2	2.345865	-2.890083	1.109969	-0.121359	-2.261857	0.524980	0.247998	
3	-0.631418	-1.059647	-0.684093	1.965775	-1.232622	-0.208038	-0.108300	
4	0.175121	-0.451449	-0.237033	-0.038195	0.803487	0.408542	-0.009431	

	V22	V23	V24	V25	V26	V27	V28	\
0	0.277838	-0.110474	0.066928	0.128539	-0.189115	0.133558	-0.021053	
1	-0.638672	0.101288	-0.339846	0.167170	0.125895	-0.008983	0.014724	
2	0.771679	0.909412	-0.689281	-0.327642	-0.139097	-0.055353	-0.059752	
3	0.005274	-0.190321	-1.175575	0.647376	-0.221929	0.062723	0.061458	
4	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153	

	Amount	Class
0	149.62	0
1	2.69	0
2	378.66	0
3	123.50	0
4	69.99	0

- 1 Time-> column represents the time elapsed in seconds between each transaction and the first transaction in the dataset.
- 2 Columns V1-V28-> These are the result of a PCA dimensionality reduction. Their meaning has been made obscure intentionally because of privacy reasons.
- 3 Amount-> Transaction amount.
- 4 Class-> Type of transaction (1 for fraudulent, 0 for legit).

5 DATA EXPLORATION

```
[7]: data.info()
```

```
<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

```

```
[8]: print("rows --->",data.shape[0],"no's")
      print("columns -->",data.shape[1],"no's")
```

```

rows ---> 284807 no's
columns --> 31 no's

```

```
[9]: # Get the number of unique values for each column
      unique_counts = data.nunique()
      print(unique_counts)
```

```

Time      124592
V1         275663
V2         275663
V3         275663
V4         275663
V5         275663
V6         275663
V7         275663
V8         275663
V9         275663
V10        275663
V11        275663
V12        275663
V13        275663
V14        275663
V15        275663
V16        275663
V17        275663
V18        275663
V19        275663
V20        275663
V21        275663
V22        275663
V23        275663
V24        275663
V25        275663
V26        275663
V27        275663
V28        275663
Amount     32767

```

Class 2
dtype: int64

```
[10]: data.describe()  
  
# if cat columns then  
# data.describe(include='object')
```

```
[10]:
```

	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

	V10	V11	V12	V13	V14 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	2.239053e-15	1.673327e-15	-1.247012e-15	8.190001e-16	1.207294e-15
std	1.088850e+00	1.020713e+00	9.992014e-01	9.952742e-01	9.585956e-01
min	-2.458826e+01	-4.797473e+00	-1.868371e+01	-5.791881e+00	-1.921433e+01
25%	-5.354257e-01	-7.624942e-01	-4.055715e-01	-6.485393e-01	-4.255740e-01
50%	-9.291738e-02	-3.275735e-02	1.400326e-01	-1.356806e-02	5.060132e-02
75%	4.539234e-01	7.395934e-01	6.182380e-01	6.625050e-01	4.931498e-01
max	2.374514e+01	1.201891e+01	7.848392e+00	7.126883e+00	1.052677e+01

	V15	V16	V17	V18	V19 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	4.887456e-15	1.437716e-15	-3.772171e-16	9.564149e-16	1.039917e-15
std	9.153160e-01	8.762529e-01	8.493371e-01	8.381762e-01	8.140405e-01
min	-4.498945e+00	-1.412985e+01	-2.516280e+01	-9.498746e+00	-7.213527e+00
25%	-5.828843e-01	-4.680368e-01	-4.837483e-01	-4.988498e-01	-4.562989e-01
50%	4.807155e-02	6.641332e-02	-6.567575e-02	-3.636312e-03	3.734823e-03
75%	6.488208e-01	5.232963e-01	3.996750e-01	5.008067e-01	4.589494e-01
max	8.877742e+00	1.731511e+01	9.253526e+00	5.041069e+00	5.591971e+00

	V20	V21	V22	V23	V24 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	6.406204e-16	1.654067e-16	-3.568593e-16	2.578648e-16	4.473266e-15
std	7.709250e-01	7.345240e-01	7.257016e-01	6.244603e-01	6.056471e-01
min	-5.449772e+01	-3.483038e+01	-1.093314e+01	-4.480774e+01	-2.836627e+00
25%	-2.117214e-01	-2.283949e-01	-5.423504e-01	-1.618463e-01	-3.545861e-01
50%	-6.248109e-02	-2.945017e-02	6.781943e-03	-1.119293e-02	4.097606e-02
75%	1.330408e-01	1.863772e-01	5.285536e-01	1.476421e-01	4.395266e-01
max	3.942090e+01	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

6 V1 to V28: These features have very small mean values (close to zero) and have been scaled and centered around zero. Their standard deviations range between 1 and 10, showing moderate variability.

#Amount: Has a higher standard deviation than its mean, indicating a wide spread in transaction amounts (from very small to very large).

#Class: This categorical feature indicates whether a transaction is regular (0) or fraudulent (1). The mean of Class is close to zero because regular transactions are far more common than fraudulent ones

7 DISTRIBUTION OF V1 ... V28 Features

```
[11]: # Set up the plot grid for 28 features (V1 to V28) with 2 columns
num_features = 28
num_cols = 2
num_rows = num_features // num_cols + (num_features % num_cols > 0) #
    ↳ Calculate number of rows needed

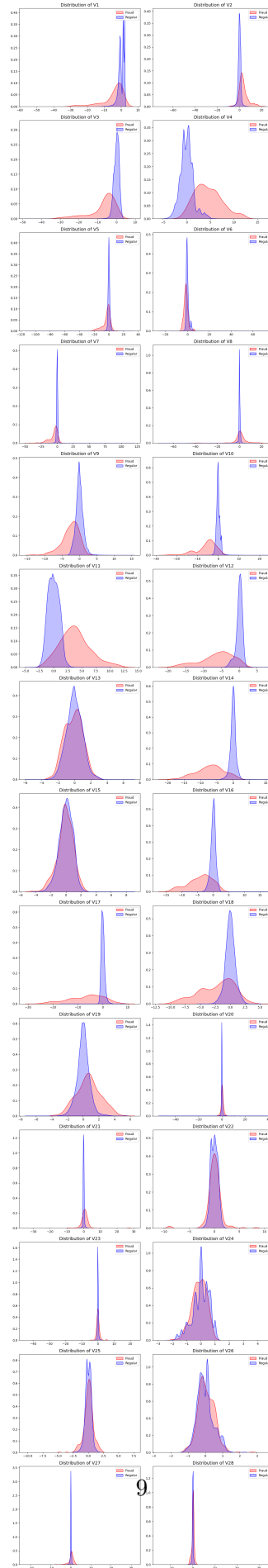
# Create the figure and axes
fig, axes = plt.subplots(num_rows, num_cols, figsize=(12, 5 * num_rows)) #
    ↳ Adjust height based on number of rows

# Flatten the axes array for easy indexing
axes = axes.flatten()

# Plot distributions for each of V1 to V28, separating by fraud and regular
    ↳ transactions
for i in range(1, num_features + 1):
    feature = f'V{i}'
    # This loop iterates over the feature indices from 1 to 28.
    # feature = f'V{i}': This creates a string representing the current feature
        ↳ being plotted (e.g., 'V1', 'V2', ..., 'V28').
        # Plot for fraud transactions
        sns.kdeplot(data[data.Class == 1][feature], ax=axes[i - 1], color='red',
        ↳ label='Fraud', fill=True)
        # Plot for regular transactions
        sns.kdeplot(data[data.Class == 0][feature], ax=axes[i - 1], color='blue',
        ↳ label='Regular', fill=True)

        axes[i - 1].set_title(f'Distribution of {feature}', fontsize=14)
        axes[i - 1].set_xlabel('')
        axes[i - 1].set_ylabel('')
        axes[i - 1].legend()

# Adjust layout to prevent overlap
plt.tight_layout()
plt.show()
```

8 Overview of the Plots

X-axis: Each feature 1 V1 to 28 V28 will have its own plot, representing the values of that feature.

Y-axis: The density of the values, which shows how likely a particular value is to occur.

Classes:

Class 0 (Regular Transactions): This will be represented by one color (e.g., blue).

Class 1 (Fraudulent Transactions): This will be represented by another color (e.g., red).

9 DISTRIBUTION PLOT OF Time Variable

```
[12]: # Set the style for the plot
sns.set(style="whitegrid")

# Create a figure for the distribution plot
plt.figure(figsize=(12, 6))

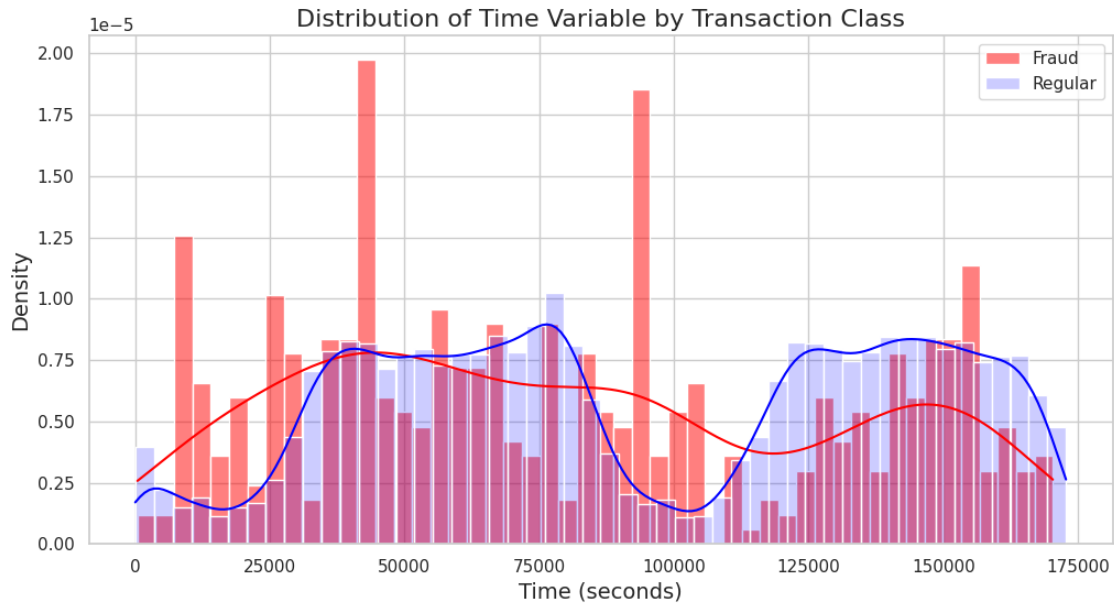
# Plot the distribution of the 'Time' variable for fraud transactions
sns.histplot(data[data.Class == 1]['Time'], bins=50, kde=True, color='red',
             stat='density', alpha=0.5, label='Fraud')

# Plot the distribution of the 'Time' variable for regular transactions
sns.histplot(data[data.Class == 0]['Time'], bins=50, kde=True, color='blue',
             stat='density', alpha=0.2, label='Regular')

# Set title and labels
plt.title('Distribution of Time Variable by Transaction Class', fontsize=16)
plt.xlabel('Time (seconds)', fontsize=14)
plt.ylabel('Density', fontsize=14)

# Add legend
plt.legend()

# Show the plot
plt.show()
```



10 Distribution of Amount

```
[13]: # Set the style for the plot
sns.set(style="whitegrid")

# Create a figure for the distribution plot with increased height
plt.figure(figsize=(12, 6)) # Adjust height as needed

# Filter the data to include only amounts below or equal to 5000
amount_filter = data[data['Amount'] <= 2000]

# Plot the distribution of the 'Amount' variable for fraud transactions
sns.histplot(amount_filter[amount_filter.Class == 1]['Amount'], bins=50,
             <-kde=True, color='red', stat='density', alpha=0.5, label='Fraud')

# Plot the distribution of the 'Amount' variable for regular transactions
sns.histplot(amount_filter[amount_filter.Class == 0]['Amount'], bins=50,
             <-kde=True, color='blue', stat='density', alpha=0.5, label='Regular')

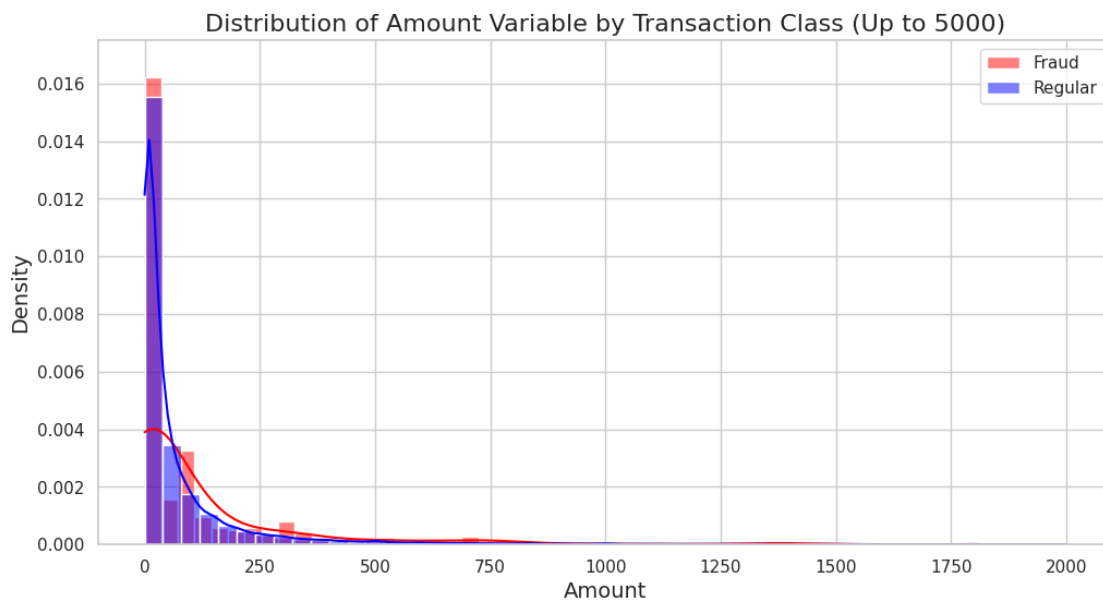
# Set title and labels
plt.title('Distribution of Amount Variable by Transaction Class (Up to 5000)',
         <-fontsize=16)
plt.xlabel('Amount', fontsize=14)
plt.ylabel('Density', fontsize=14)

# Adjust y-axis limits if needed
```

```
plt.ylim(0, 0.0175) # Adjust based on the data

# Add legend
plt.legend()

# Show the plot
plt.show()
```



Red (Fraud Transactions): This color is used to indicate transactions that are classified as fraudulent (where `Class == 1`). These transactions are typically of particular interest when analyzing credit card fraud data because they represent instances of fraudulent activity.

Blue (Regular Transactions): This color represents transactions that are classified as legitimate or regular (where `Class == 0`). These are the majority of transactions in the dataset and serve as a baseline for comparison against fraudulent transactions.

11 we visualized ->The Amount distribution for frauds is peaked on small quantities of money.

12 Duplicate Values

```
[14]: print('There are {} duplicate values in regular transactions out of {}'.  
      ↪format(data[data['Class'] == 0].duplicated().sum(),data[data['Class'] == 0].  
      ↪shape[0]))  
print('There are {} duplicate values in fraudulent transactions out of {}'.  
      ↪format(data[data['Class'] == 1].duplicated().sum(),data[data['Class'] == 1].  
      ↪shape[0]))
```

There are 1062 duplicate values in regular transactions out of 284315

There are 19 duplicate values in fraudulent transactions out of 492

`inplace=True`: This modifies the DataFrame in place. The original DataFrame is updated, and no new DataFrame is returned. As a result, you do not need to assign the result to a new variable or even to the same variable.

`inplace=False` (default behavior): This will return a new DataFrame with the duplicates removed, leaving the original DataFrame unchanged. You would typically assign the result to a new variable or overwrite the existing variable.

```
[15]: print('No. of rows before dropping duplicates: {}'.format(len(data)))  
  
data.drop_duplicates(inplace=True)  
  
print('No. of rows after dropping duplicates: {}'.format(len(data)))
```

No. of rows before dropping duplicates: 284807.

No. of rows after dropping duplicates: 283726.

here few duplicate rows dropped

13 HANDLING OUTLIERS

outlier can be extremely high or low values compared to the other observations and can be caused by measurement errors, natural variations in the data, or even unexpected discoveries.that leads lower performance so we deal with it...

```
[16]: # Create subplots  
fig, (ax1, ax2) = plt.subplots(ncols=2, figsize=(12, 5))  
  
# Boxplot for Time  
sns.boxplot(x=data['Time'], ax=ax1)  
ax1.set_title('Boxplot of Time', fontsize=20)  
ax1.set_facecolor('honeydew')  
  
# Boxplot for Amount  
sns.boxplot(x=data['Amount'], ax=ax2)  
ax2.set_title('Boxplot of Amount', fontsize=20)  
ax2.set_facecolor('honeydew')
```

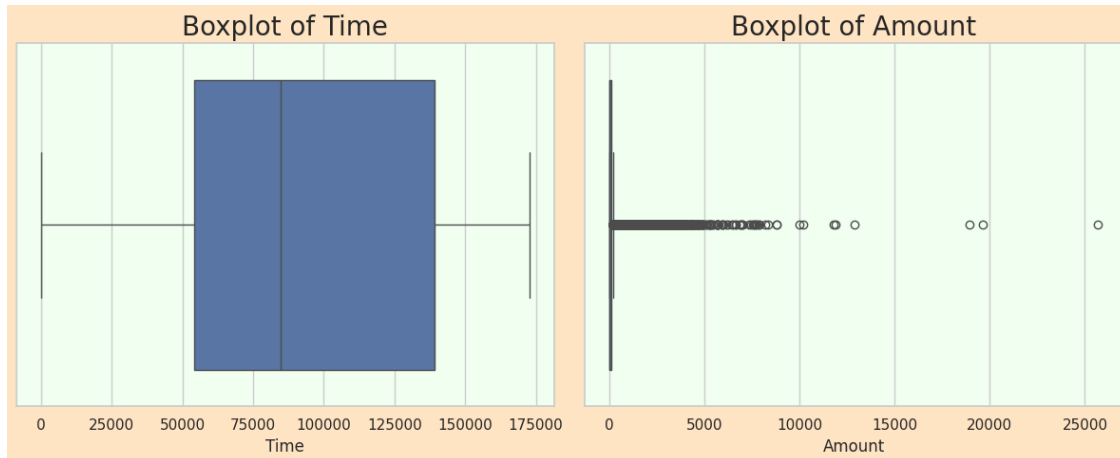
```

# Adjust layout
plt.tight_layout()

# Set figure background color
fig.set_facecolor('bisque')

# Show plot
plt.show()

```



here no outlier in “time” but too many outlier in “amount”

```

[17]: # Select only the V1-V28 features
features = list(data.columns.values)
del features[0]
del features[28]
del features[28]

for i in range(7):#i 0 to 6
    fig,(ax1,ax2,ax3,ax4)=plt.subplots(ncols=4,figsize=(12,5))#wide=12,tall=5
    #plot diagram by 4 rows
    ax1=sns.boxplot(data[features[i*4]],ax=ax1)
    ax1.set_title('Boxplot of '+str(features[i*4]),fontsize=20)
    ax1.set_facecolor('honeydew')

    ax2=sns.boxplot(data[features[i*4+1]],ax=ax2)
    ax2.set_title('Boxplot of '+str(features[i*4+1]),fontsize=20)
    ax2.set_facecolor('honeydew')

    ax3=sns.boxplot(data[features[i*4+2]],ax=ax3)
    ax3.set_title('Boxplot of '+str(features[i*4+2]),fontsize=20)

```

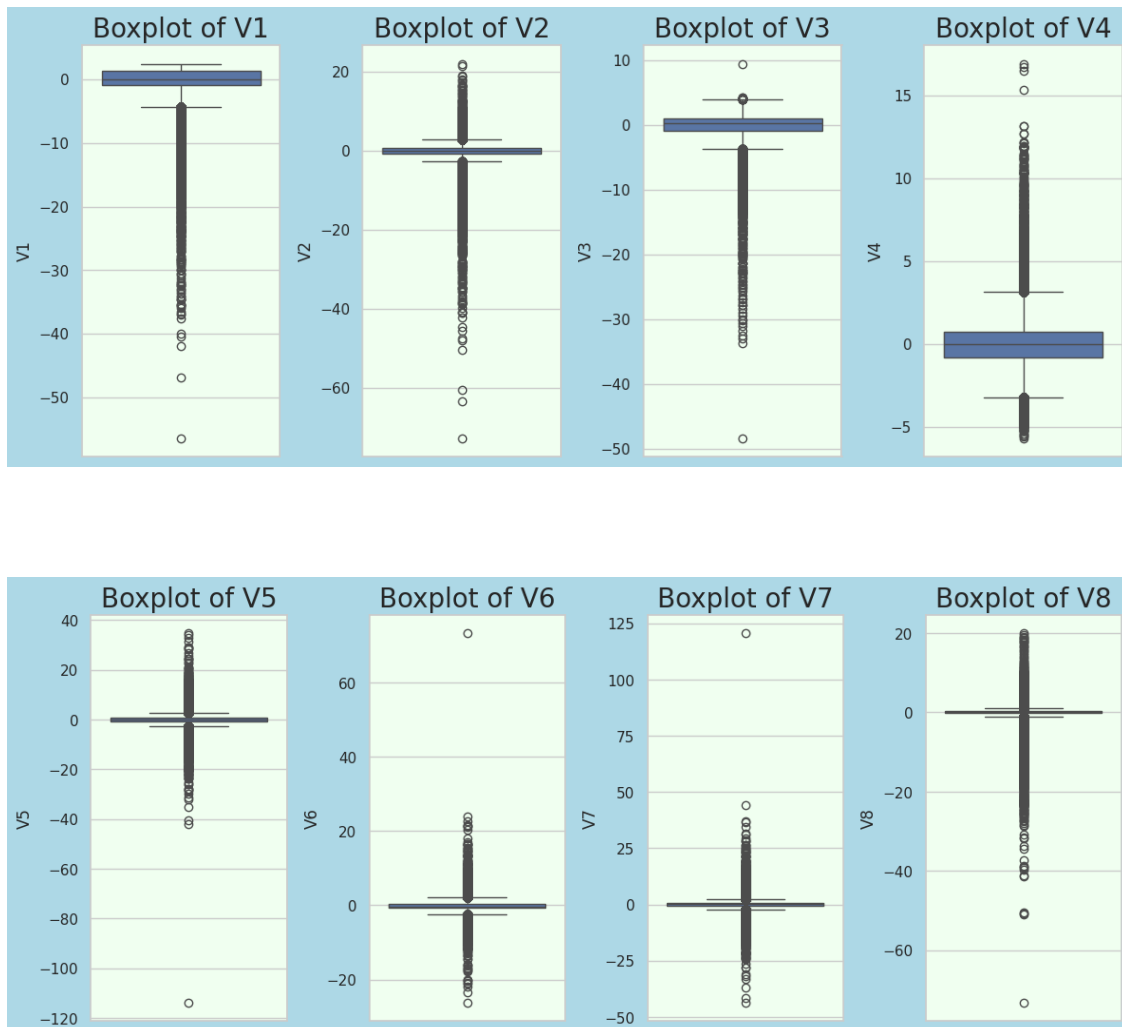
```

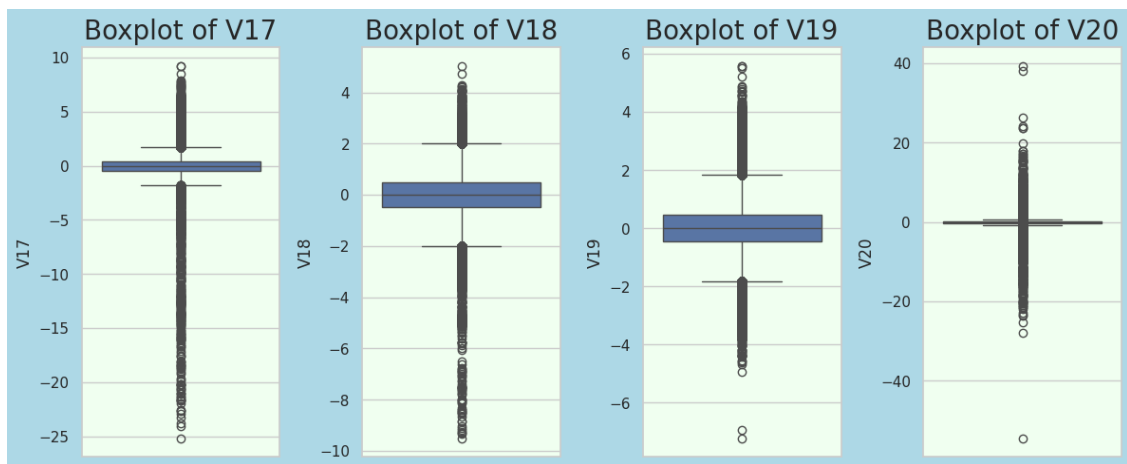
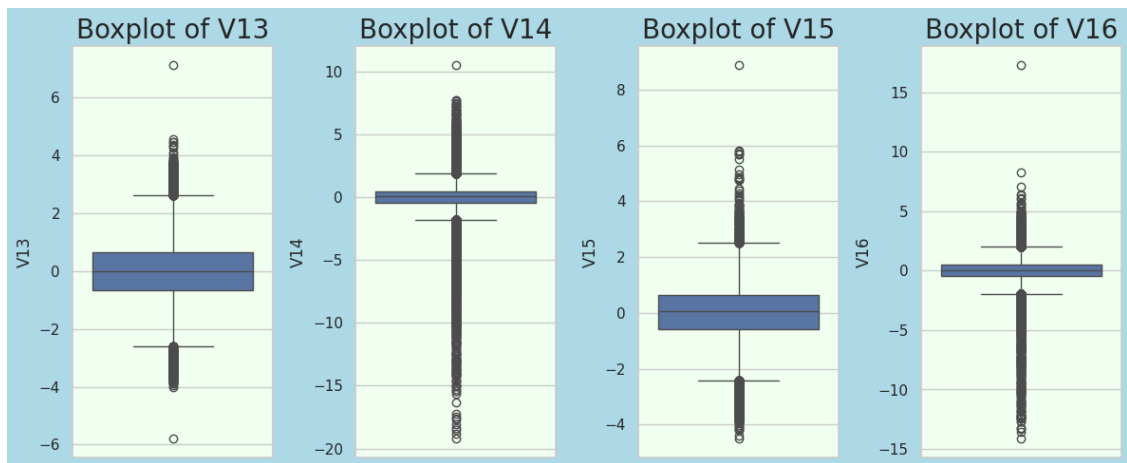
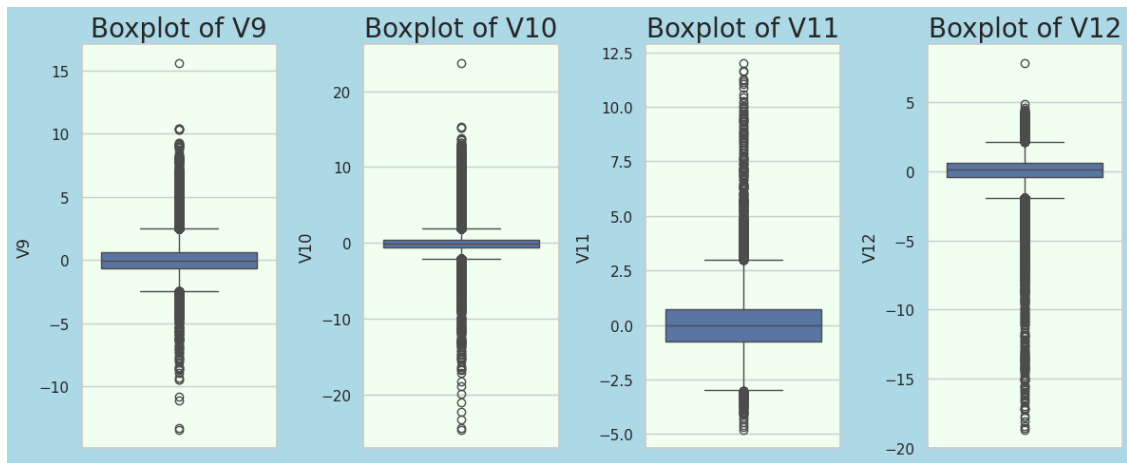
ax3.set_facecolor('honeydew')

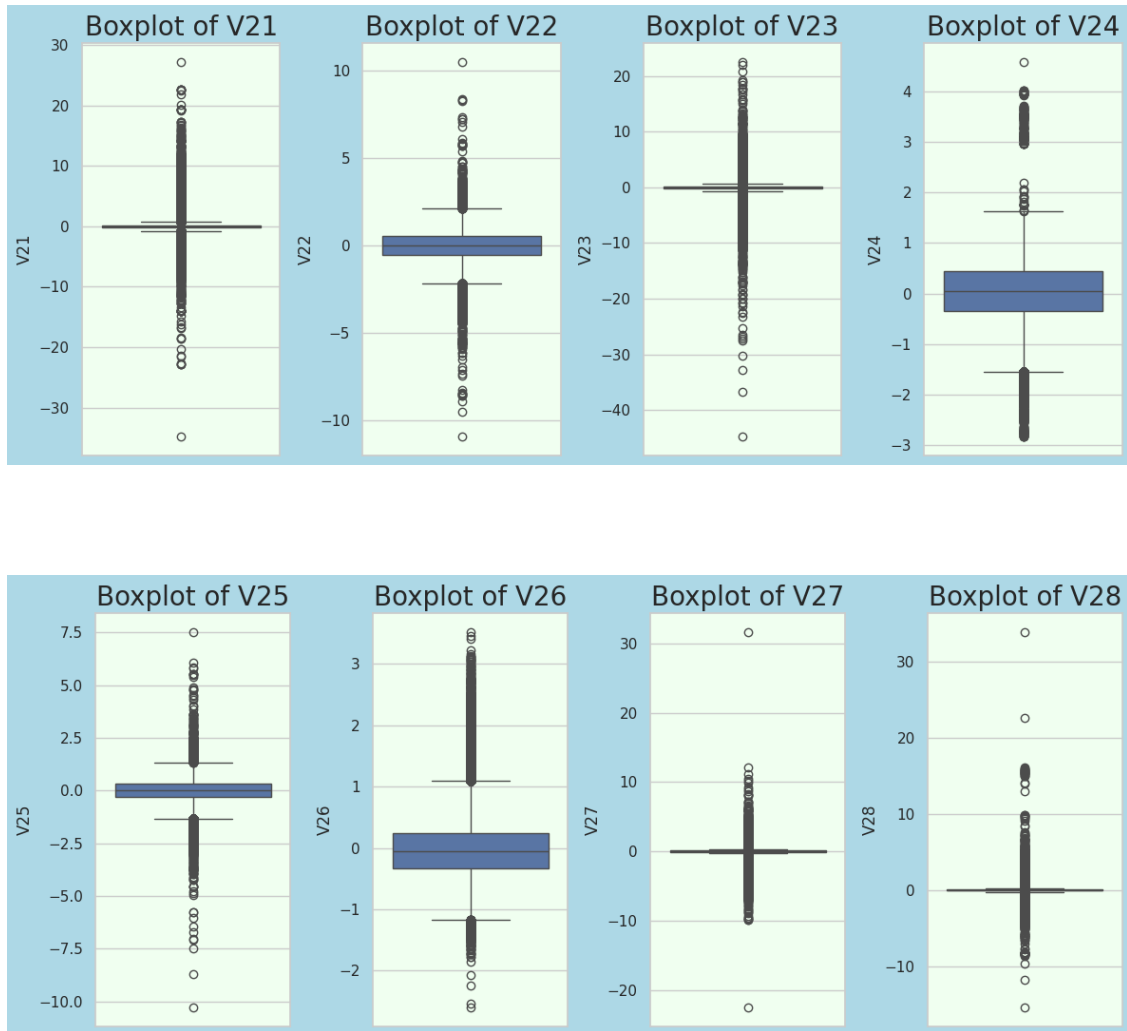
ax4=sns.boxplot(data[features[i*4+3]],ax=ax4)
ax4.set_title('Boxplot of '+str(features[i*4+3]),fontsize=20)
ax4.set_facecolor('honeydew')

plt.tight_layout()
fig.set_facecolor('#ADD8E6')

```







```
[18]: for column in features:
      #check how much skew
      print(f"Skewness for {column}: {data[column].skew()}")
```

```
Skewness for V1: -3.273271248440309
Skewness for V2: -4.6951619005404694
Skewness for V3: -2.1519839570997124
Skewness for V4: 0.6715041706728241
Skewness for V5: -2.414079246966253
Skewness for V6: 1.829880383771521
Skewness for V7: 2.890271192715498
Skewness for V8: -8.310970330052545
Skewness for V9: 0.5376630534496958
Skewness for V10: 1.2529670787468168
Skewness for V11: 0.34407419325686267
Skewness for V12: -2.1990082816149954
```

```

Skewness for V13: 0.06429340464018111
Skewness for V14: -1.9188037137586451
Skewness for V15: -0.3096590822936595
Skewness for V16: -1.0511614715174662
Skewness for V17: -3.690497194148406
Skewness for V18: -0.24866145737243997
Skewness for V19: 0.1083118109324772
Skewness for V20: -2.0431210560273323
Skewness for V21: 2.820033113572543
Skewness for V22: -0.18232972797521269
Skewness for V23: -5.867220791006341
Skewness for V24: -0.5521292366718961
Skewness for V25: -0.41574386205469593
Skewness for V26: 0.5802923172348093
Skewness for V27: -0.7538039138186547
Skewness for V28: 11.555115084196773

```

-ve indicates left skew

+ve indicates right skew

#CAPING OUTLIER

```

[19]: def outlier_imputer(data,features):

    data_out = data.copy()

    for column in features:

        # First define the first and third quartiles
        Q1 = data_out[column].quantile(0.25)
        Q3 = data_out[column].quantile(0.75)
        # Define the inter-quartile range
        IQR = Q3 - Q1
        # ... and the lower/higher threshold values
        lowerL = Q1 - 1.5 * IQR
        higherL = Q3 + 1.5 * IQR

        # Impute 'left' outliers
        data_out.loc[data_out[column] < lowerL,column] = lowerL
        #which rows value lower than lowerL that capped
        # Impute 'right' outliers
        data_out.loc[data_out[column] > higherL,column] = higherL
        #data_out.loc[mask, column] locates the rows in column where
        →data_out[column] > higherL is True.-->[rows,columns]

    return data_out

#only made because of extracting features

```

```
data2 = data.drop('Class',axis=1)

feats = list(data2.columns.values) #list of column names

capped_data = outlier_imputer(data,feats)
```

```
[20]: pd.set_option("display.max_columns",1000)

capped_data.head()
```

```
[20]:
```

	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	V10	V11	V12	V13	V14	\
0	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.311169	
1	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.143772	
2	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.165946	
3	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.287924	
4	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.119670	

	V15	V16	V17	V18	V19	V20	V21	\
0	1.468177	-0.470401	0.207971	0.025791	0.403993	0.251412	-0.018307	
1	0.635558	0.463917	-0.114805	-0.183361	-0.145783	-0.069083	-0.225775	
2	2.345865	-1.952416	1.109969	-0.121359	-1.828484	0.524980	0.247998	
3	-0.631418	-1.059647	-0.684093	1.965775	-1.232622	-0.208038	-0.108300	
4	0.175121	-0.451449	-0.237033	-0.038195	0.803487	0.408542	-0.009431	

	V22	V23	V24	V25	V26	V27	V28	\
0	0.277838	-0.110474	0.066928	0.128539	-0.189115	0.133558	-0.021053	
1	-0.638672	0.101288	-0.339846	0.167170	0.125895	-0.008983	0.014724	
2	0.771679	0.611926	-0.689281	-0.327642	-0.139097	-0.055353	-0.059752	
3	0.005274	-0.190321	-1.175575	0.647376	-0.221929	0.062723	0.061458	
4	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153	

	Amount	Class
0	149.620	0
1	2.690	0
2	185.375	0
3	123.500	0
4	69.990	0

the outliers have been imputed and the extreme values have been set to either $Q_1 - 1.5 * IQR$ (lower threshold) or to $Q_3 + 1.5 * IQR$ (higher threshold), where Q_1

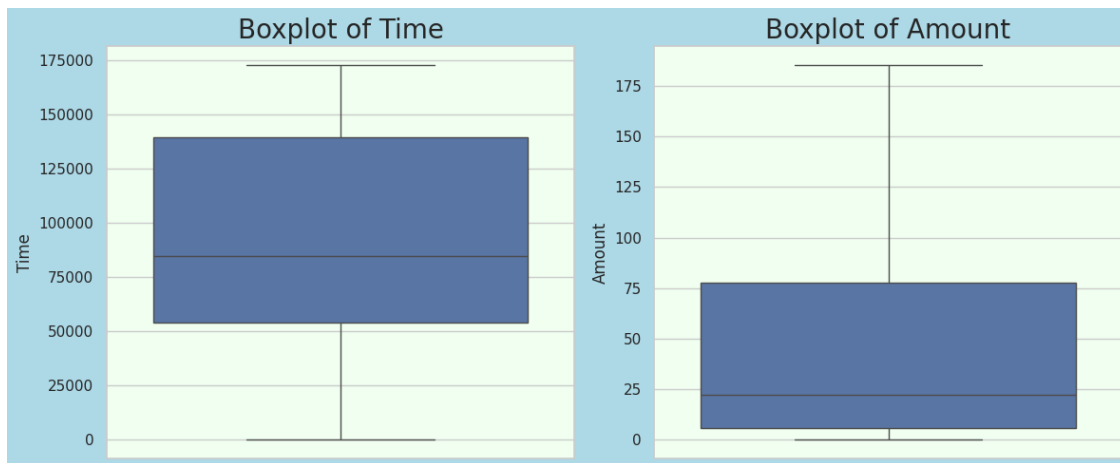
and Q 3 are the first and third quartiles and I Q R is the so called interquartile range

14 AFTER CAPPING

```
[21]: fig,(ax1,ax2) = plt.subplots(ncols=2,figsize=(12,5))
```

```
ax1 = sns.boxplot(capped_data['Time'],ax=ax1)
ax1.set_title('Boxplot of Time',fontsize=20)
ax1.set_facecolor('honeydew')
ax2 = sns.boxplot(capped_data['Amount'],ax=ax2)
ax2.set_title('Boxplot of Amount',fontsize=20)
ax2.set_facecolor('honeydew')
```

```
plt.tight_layout()
fig.set_facecolor('#ADD8E6')
```



```
[22]: # Select only the V1-V28 features
features_c = list(capped_data.columns.values)
del features_c[0]
del features_c[28]
del features_c[28]

for i in range(7):#i 0 to 6
    fig,(ax1,ax2,ax3,ax4)=plt.subplots(ncols=4,figsize=(12,5))#wide=12,tall=5
    #plot diagram by 4 rows
    ax1=sns.boxplot(capped_data[features_c[i*4]],ax=ax1)
    ax1.set_title('Boxplot of '+str(features_c[i*4]),fontsize=20)
    ax1.set_facecolor('honeydew')
```

```

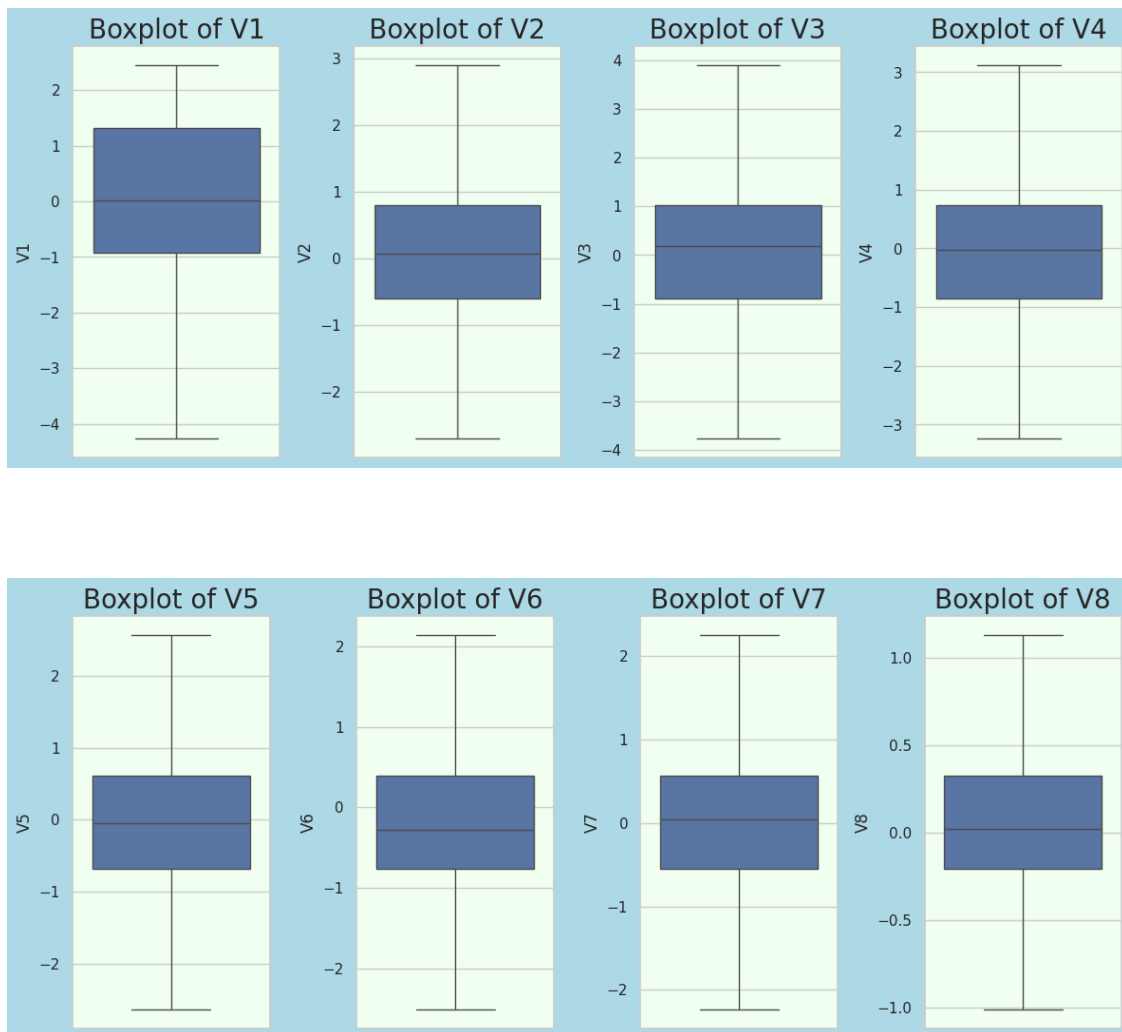
ax2=sns.boxplot(capped_data[features_c[i*4+1]],ax=ax2)
ax2.set_title('Boxplot of '+str(features_c[i*4+1]),fontsize=20)
ax2.set_facecolor('honeydew')

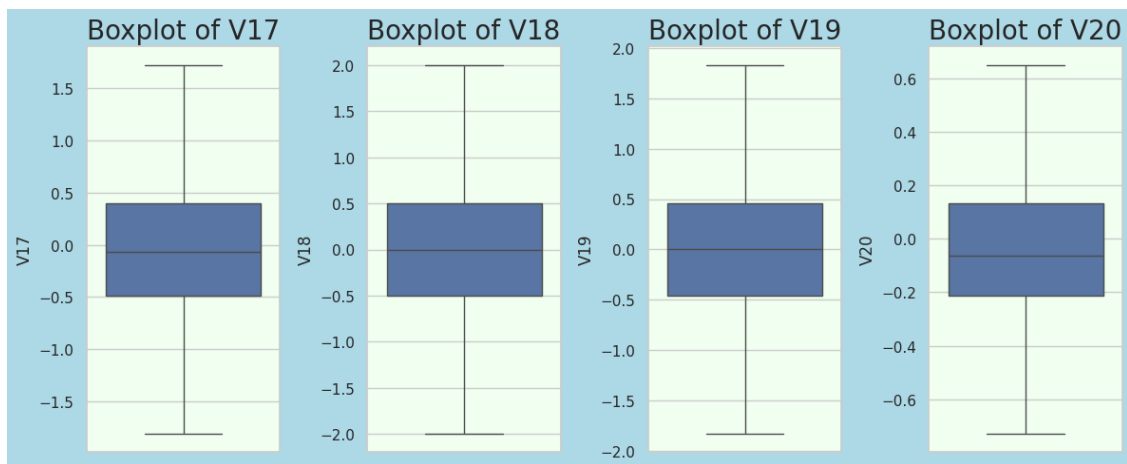
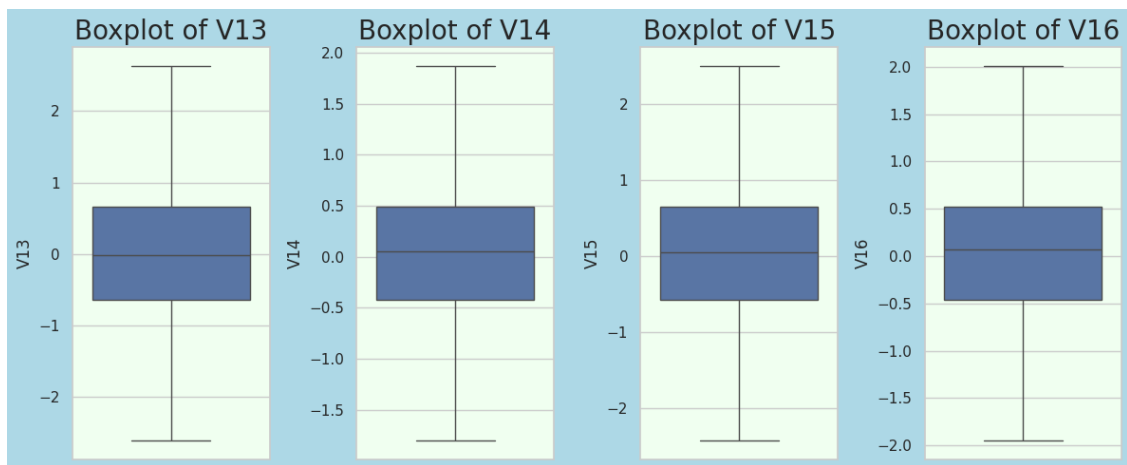
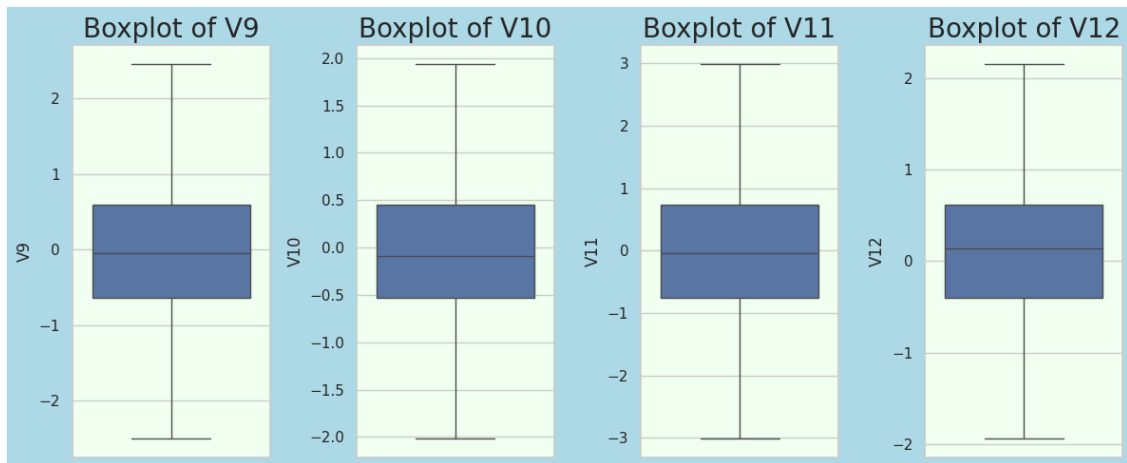
ax3=sns.boxplot(capped_data[features_c[i*4+2]],ax=ax3)
ax3.set_title('Boxplot of '+str(features_c[i*4+2]),fontsize=20)
ax3.set_facecolor('honeydew')

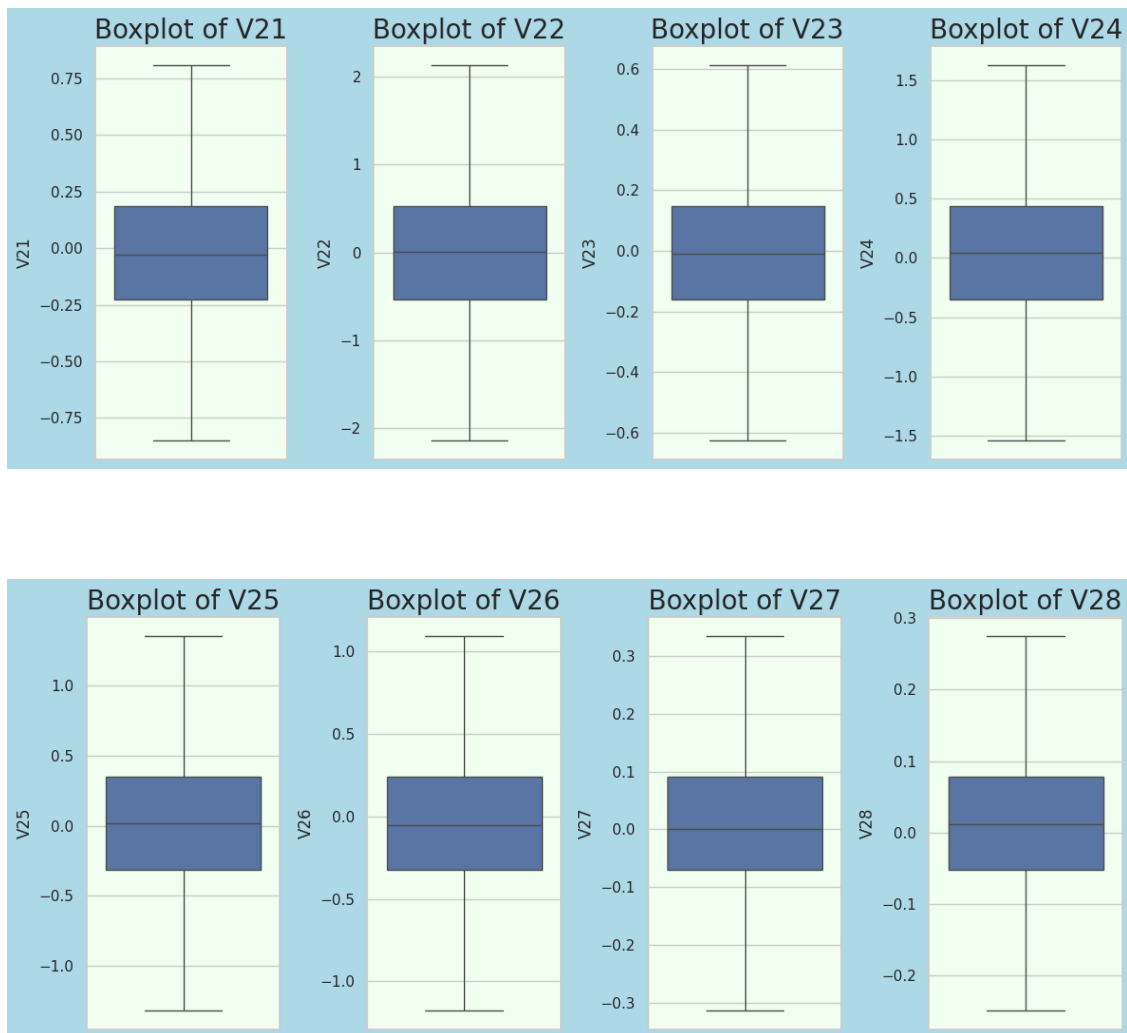
ax4=sns.boxplot(capped_data[features_c[i*4+3]],ax=ax4)
ax4.set_title('Boxplot of '+str(features_c[i*4+3]),fontsize=20)
ax4.set_facecolor('honeydew')

plt.tight_layout()
fig.set_facecolor('#ADD8E6')

```







NOTE : after capping outliers, the data points are restricted to a defined range, which can potentially lead to a situation where many values are concentrated around the capped limits. This is where using anomaly detection techniques like Isolation Forest becomes especially important.

15 Isolation Forest & Dropping Outliers

```
[23]: data3 = data.copy()
      data3 = data3.drop('Class',axis=1)
```

```
[24]: model = IsolationForest(n_estimators=150,max_samples='auto',contamination=float(0.
      ↪1),max_features=1.0)
```

```
model.fit(data3)
```

```
[24]: IsolationForest(contamination=0.1, n_estimators=150)
```

Then, I am adding the score values and also an anomaly column to the dataframe. The negative score value of -1 indicates the presence of an anomaly. The value of 1 for the anomaly represents normal data.

You might set a threshold where scores less than 0 are considered anomalies, while scores equal to or greater than 0 are normal.

```
[25]: scores = model.decision_function(data3)
      scores
```

```
[25]: array([ 0.06763218,  0.09511447, -0.01269494, ...,  0.04583759,
            0.01526097,  0.06730932])
```

Calculating Scores: `scores = model.decision_function(data3)` computes the anomaly score for each data point in `data3`. The score reflects how isolated a data point is:

Points with lower scores are more likely to be anomalies.

scores that are more negative (e.g., -0.5, -1.0). These scores suggest that the point is more isolated, meaning it does not fit the typical pattern of the majority of the data.

```
[26]: anomaly = model.predict(data3)
      anomaly
```

```
[26]: array([ 1,  1, -1, ...,  1,  1,  1])
```

generates binary predictions for the data points: The method returns -1 for points identified as anomalies and 1 for normal points.

```
[27]: #adding scores and anomalies
      data3['score'] = scores
      data3['anomaly'] = anomaly

      data3.head()
```

```
[27]:
```

	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	V10	V11	V12	V13	V14	\
0	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.311169	
1	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.143772	
2	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.165946	


```

3  0.377436 -1.387024 -0.054952 -0.226487  0.178228  0.507757 -0.287924
4 -0.270533  0.817739  0.753074 -0.822843  0.538196  1.345852 -1.119670

```

```

          V15          V16          V17          V18          V19          V20          V21  \
0  1.468177 -0.470401  0.207971  0.025791  0.403993  0.251412 -0.018307
1  0.635558  0.463917 -0.114805 -0.183361 -0.145783 -0.069083 -0.225775
2  2.345865 -2.890083  1.109969 -0.121359 -2.261857  0.524980  0.247998
3 -0.631418 -1.059647 -0.684093  1.965775 -1.232622 -0.208038 -0.108300
4  0.175121 -0.451449 -0.237033 -0.038195  0.803487  0.408542 -0.009431

```

```

          V22          V23          V24          V25          V26          V27          V28  \
0  0.277838 -0.110474  0.066928  0.128539 -0.189115  0.133558 -0.021053
1 -0.638672  0.101288 -0.339846  0.167170  0.125895 -0.008983  0.014724
2  0.771679  0.909412 -0.689281 -0.327642 -0.139097 -0.055353 -0.059752
3  0.005274 -0.190321 -1.175575  0.647376 -0.221929  0.062723  0.061458
4  0.798278 -0.137458  0.141267 -0.206010  0.502292  0.219422  0.215153

```

```

      Amount      score  anomaly
0  149.62  0.067632         1
1    2.69  0.095114         1
2  378.66 -0.012695        -1
3  123.50  0.048105         1
4   69.99  0.059486         1

```

```

[28]: anomaly = data3.loc[data3['anomaly'] == -1]
      print('The total number of outliers is {} out of {}'.format(len(anomaly),len(data)))

```

The total number of outliers is 28373 out of 283726.

dropping the outliers.

```

[29]: anomaly_index = list(anomaly.index)
      forest_data = data.drop(anomaly_index,axis=0).reset_index(drop=True)

```

16 FINAL DISTRIBUTION

```

[30]: for i in range(14):
      fig,(ax1,ax2) = plt.subplots(ncols=2,figsize=(12,5))
      ax1 = sns.distplot(forest_data[forest_data.Class == 1][features[i*2]],ax=ax1,hist=False)
      ax1 = sns.distplot(forest_data[forest_data.Class == 0][features[i*2]],ax=ax1,hist=False)
      ax1.set_title('Distribution of '+str(features[i*2]),fontsize=20)
      ax1.set_facecolor('honeydew')

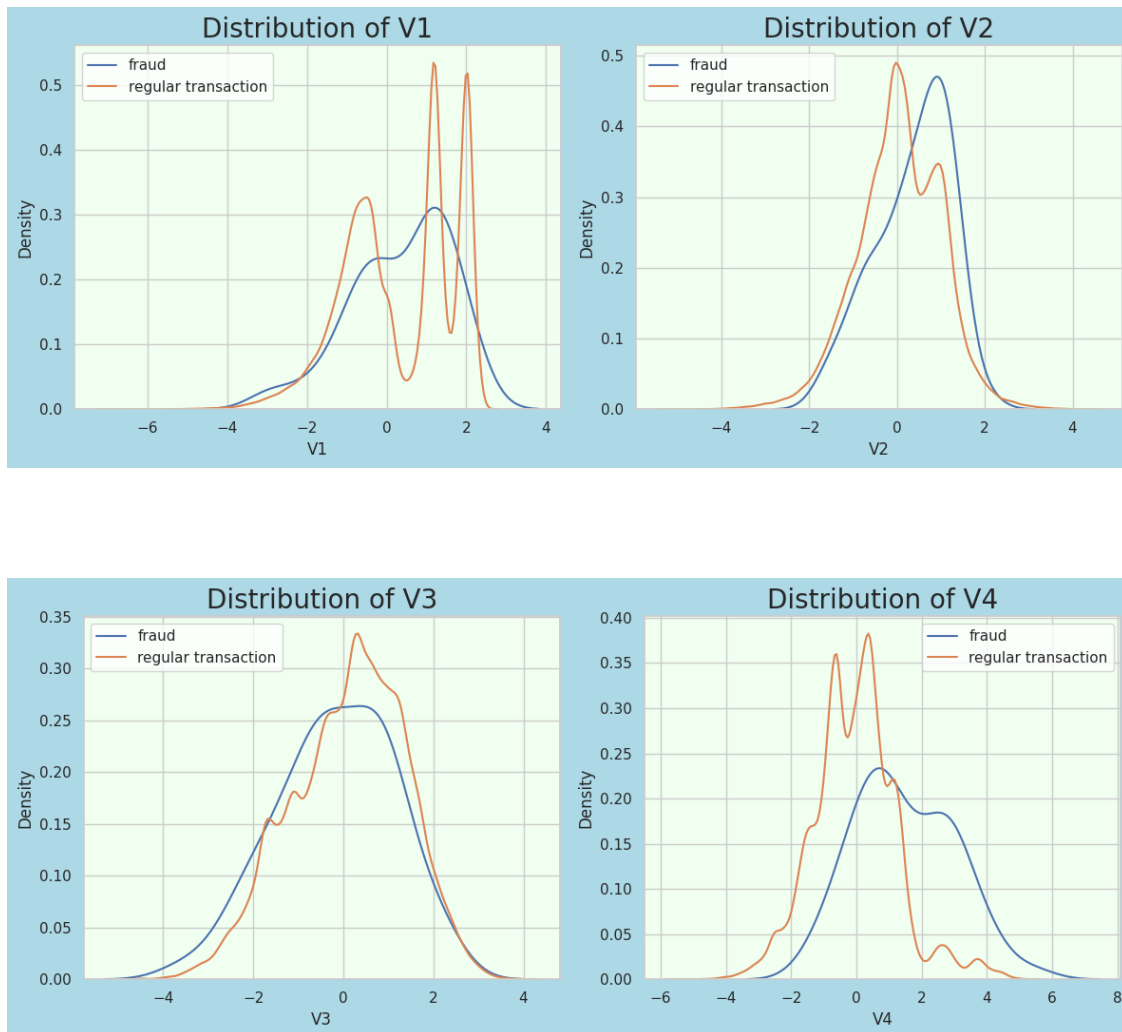
```

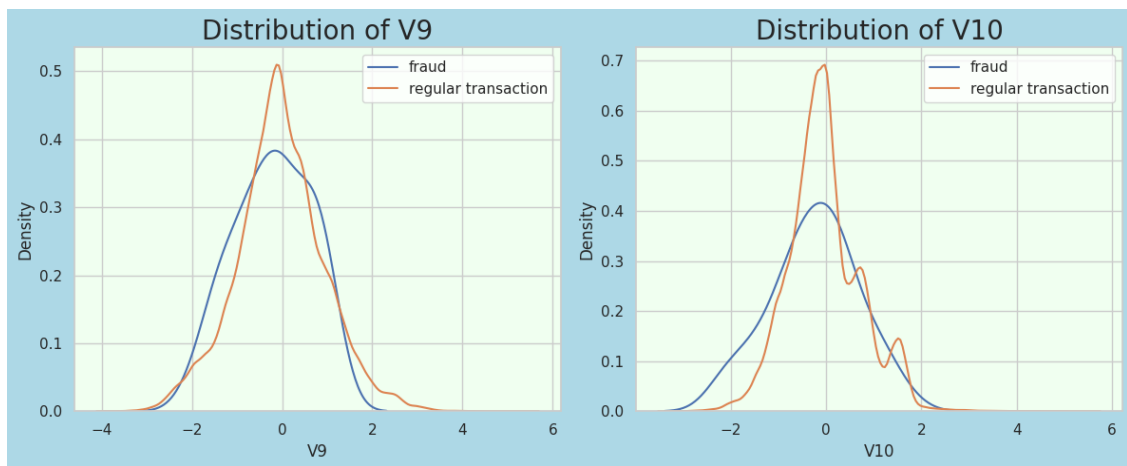
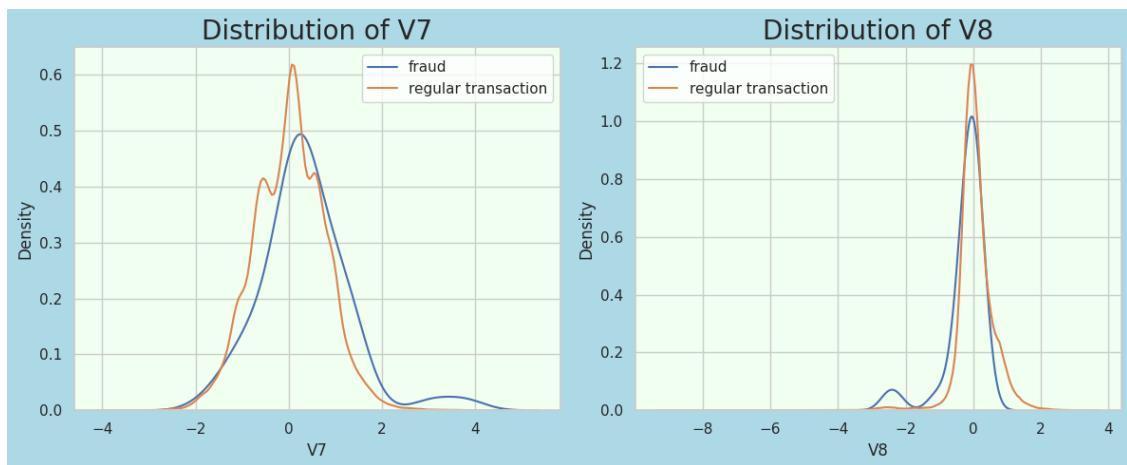
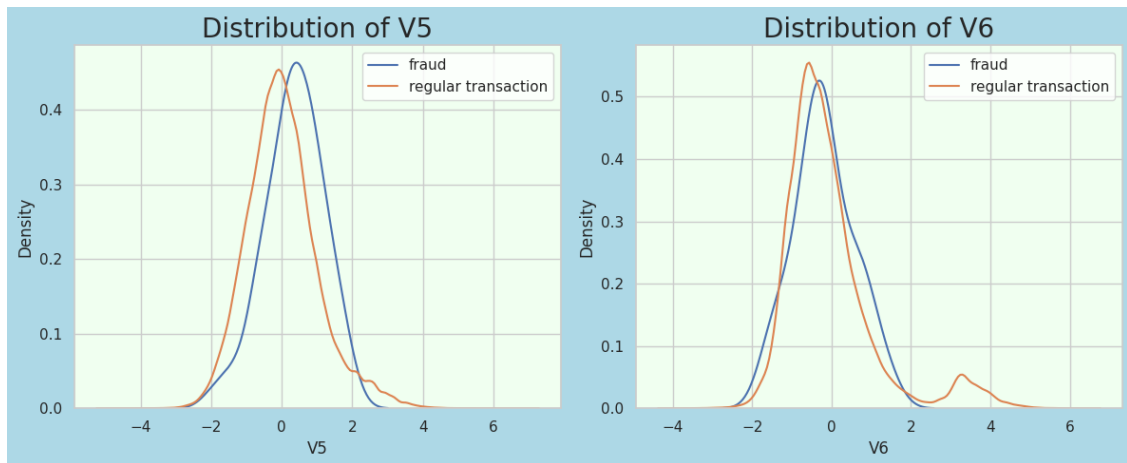
```

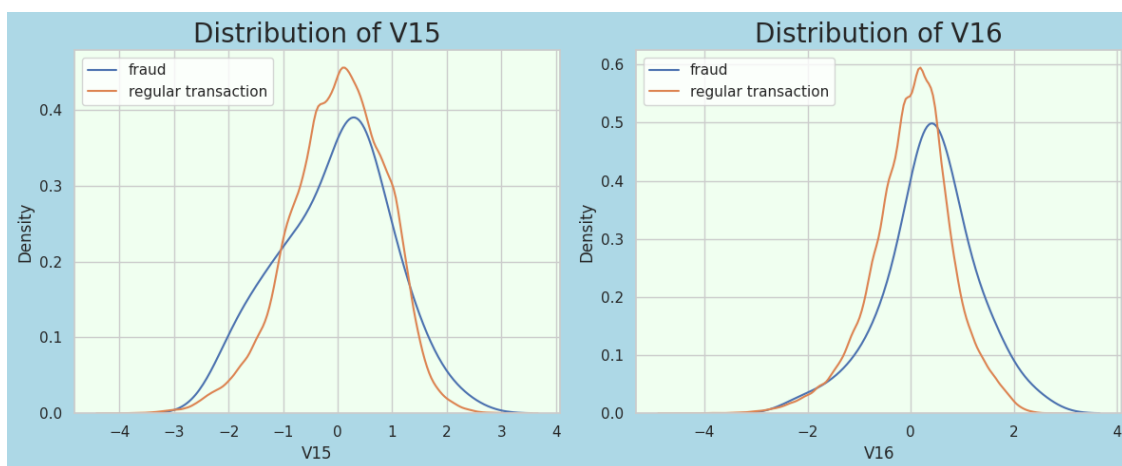
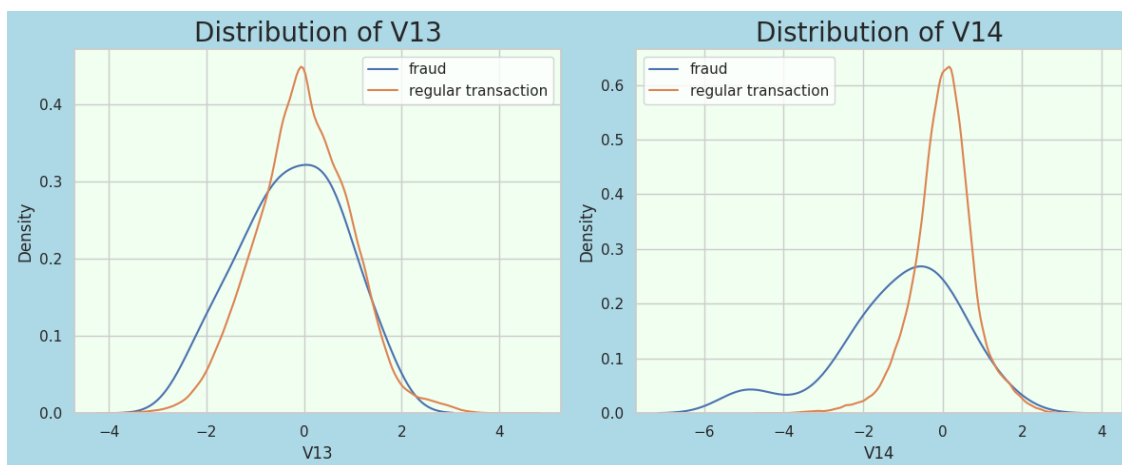
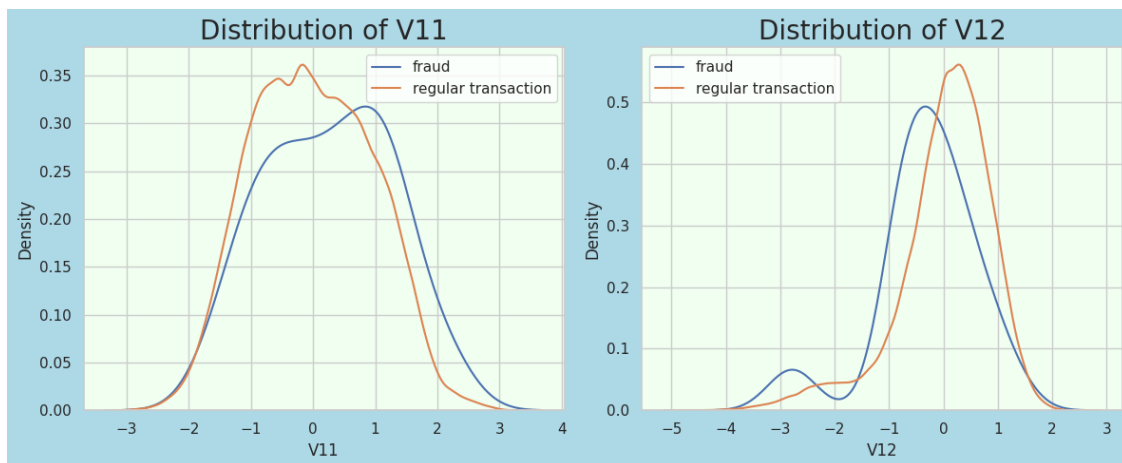
ax1.legend(labels=['fraud', 'regular transaction'])
ax2 = sns.distplot(forest_data[forest_data.Class == 1]
    ↪ [features[i*2+1]], ax=ax2, hist=False)
ax2 = sns.distplot(forest_data[forest_data.Class == 0]
    ↪ [features[i*2+1]], ax=ax2, hist=False)
ax2.set_title('Distribution of '+str(features[i*2+1]), fontsize=20)
ax2.set_facecolor('honeydew')
ax2.legend(labels=['fraud', 'regular transaction'])

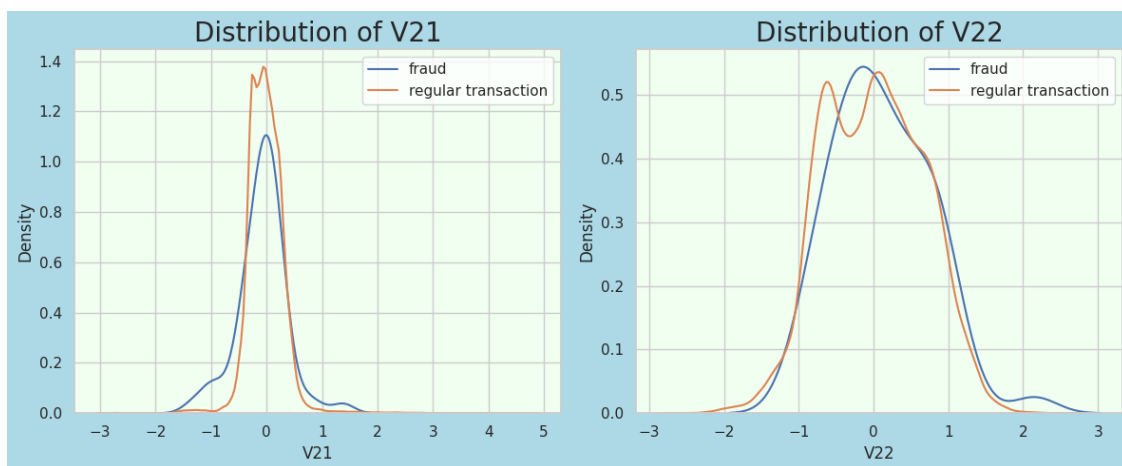
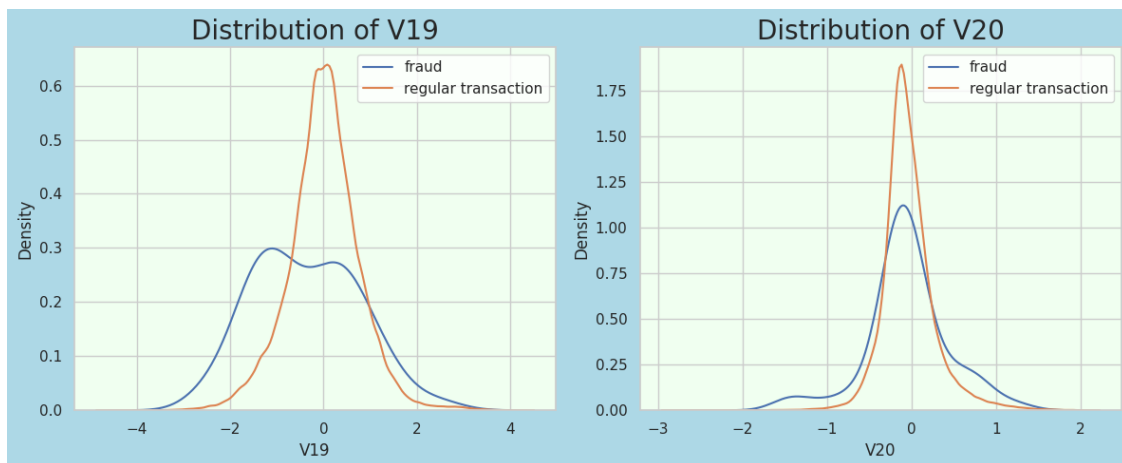
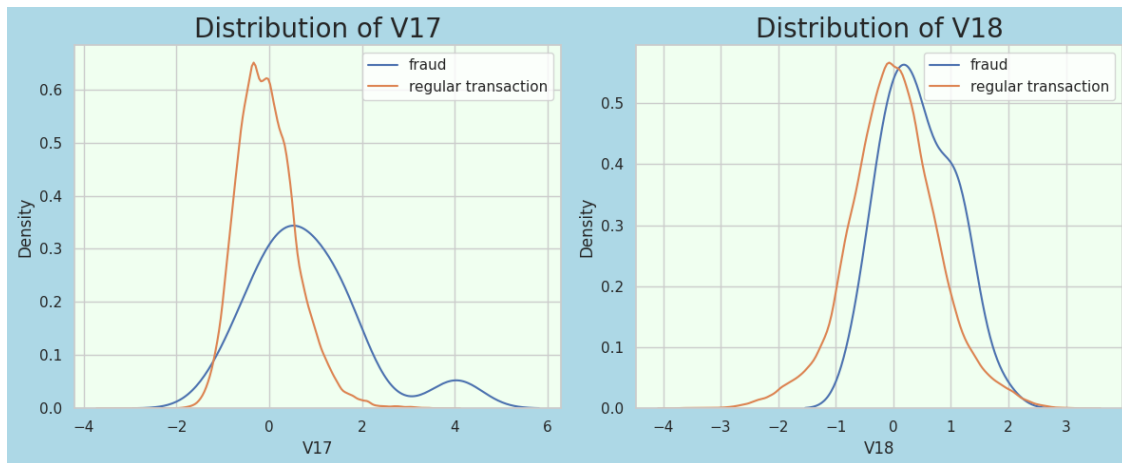
plt.tight_layout()
fig.set_facecolor('#ADD8E6')

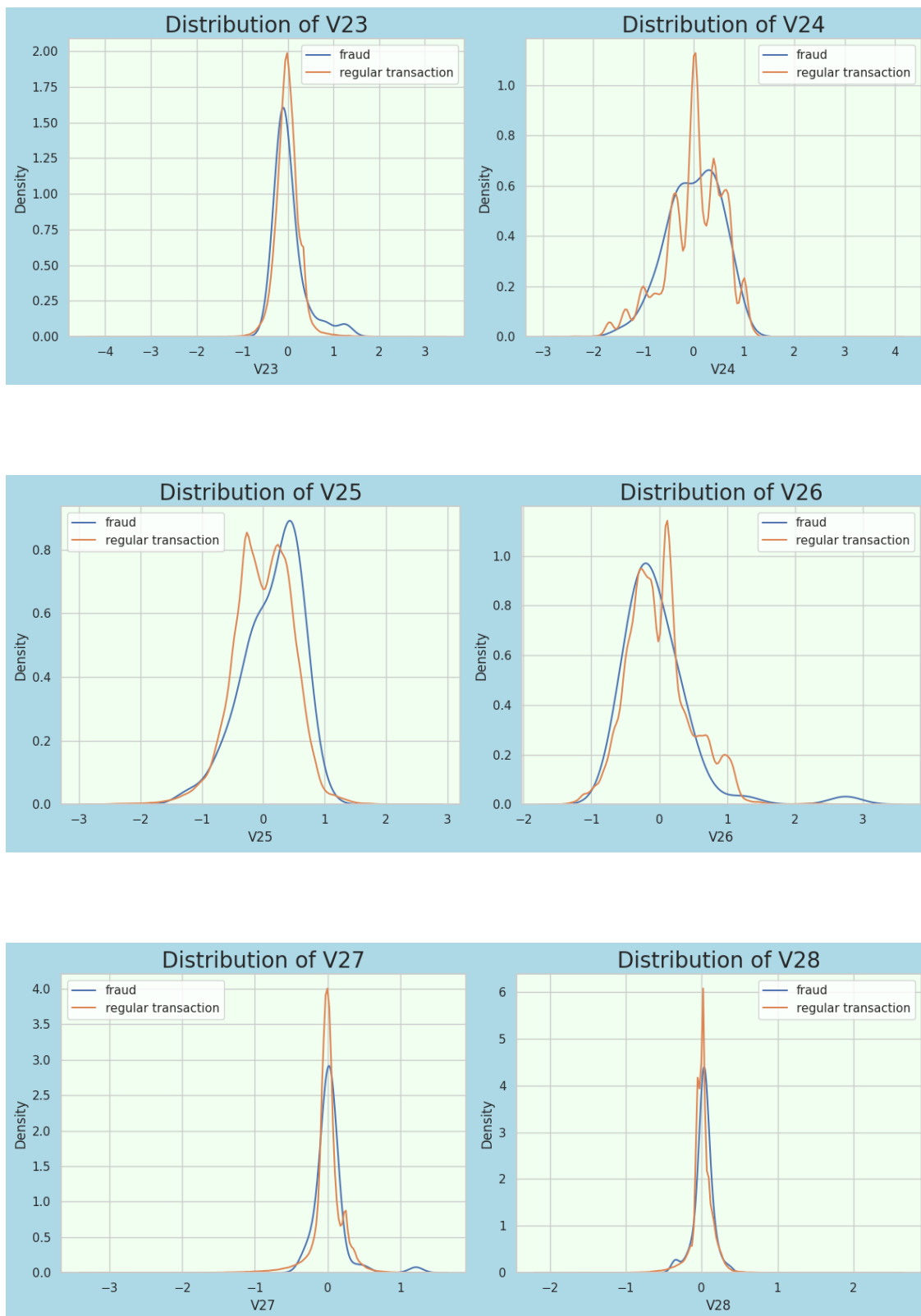
```











COMPARING

```
[31]: fig,(ax1,ax2) = plt.subplots(ncols=2,figsize=(12,5))

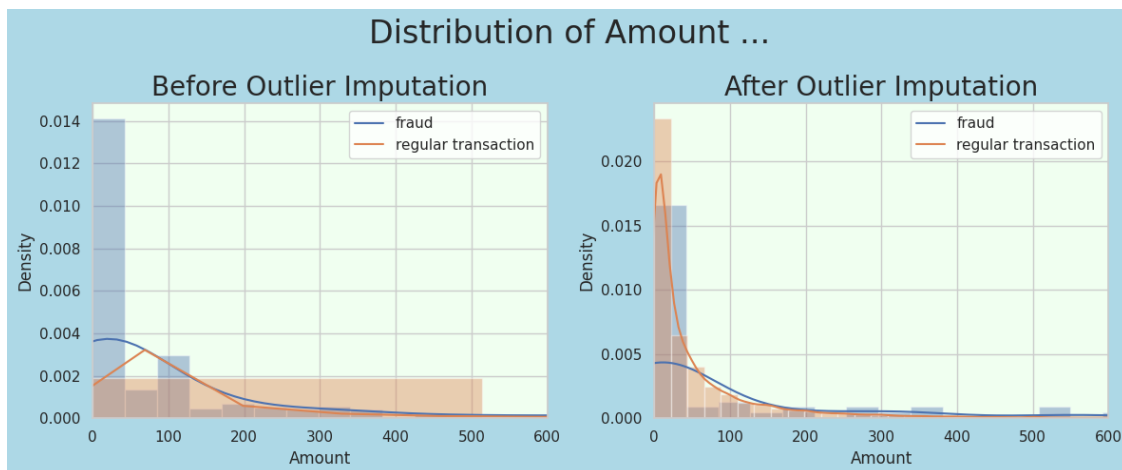
#BEFORE
ax1 = sns.distplot(data['Amount'][data.Class == 1],ax=ax1)
ax1 = sns.distplot(data['Amount'][data.Class == 0],ax=ax1)
ax1.set_xlim(0,600)
ax1.set_title('Before Outlier Imputation',fontsize=20)
ax1.legend(labels=['fraud','regular transaction'])
ax1.set_facecolor('honeydew')

#AFTER
ax2 = sns.distplot(forest_data['Amount'][forest_data.Class == 1],ax=ax2)
ax2 = sns.distplot(forest_data['Amount'][forest_data.Class == 0],ax=ax2)
ax2.set_xlim(0,600)
ax2.set_title('After Outlier Imputation',fontsize=20)
ax2.legend(labels=['fraud','regular transaction'])
ax2.set_facecolor('honeydew')

#PLOT

fig.suptitle("Distribution of Amount ...",fontsize=24)

plt.tight_layout()
fig.set_facecolor('#ADD8E6')
```



Before Outlier Imputation:

In the left plot, some transactions have very high amounts, which makes the data look stretched out to the right. Regular transactions usually have smaller amounts, while fraud transactions

sometimes go higher.

After Outlier Imputation:

In the right plot, those high-amount transactions have been reduced or adjusted, so the data is now focused more on smaller amounts. This makes the graph more compact, with fewer large transactions.

In short: MACHINE FOCOUS ON HIGH VALUE SO ,

We reduced the impact of big, unusual transaction amounts, so now the data focuses more on typical, smaller amounts. This can help make patterns in the data clearer.

17 FEATURE IMPORTANCE

identify which features are most influential in making predictions.

```
[32]: X = forest_data.drop('Class',axis=1)
      y = forest_data['Class']
```

I can look at the relative importance of the features by means of a random forest classifier.

```
[33]: #USING RANDOM FOREST FIND FEATURE IMPORTANCE
      # Random Forest Model
      random_forest = RandomForestClassifier(random_state=1,max_depth=4)
      random_forest.fit(X,y)
```

```
[33]: RandomForestClassifier(max_depth=4, random_state=1)
```

```
[34]: FI=random_forest.feature_importances_
      importances = pd.DataFrame({'feature':X.columns,'importance':np.round(FI,3)})
      importances = importances.sort_values('importance',ascending=False).
      ↪set_index('feature')
      importances.head(30)
```

```
[34]:
```

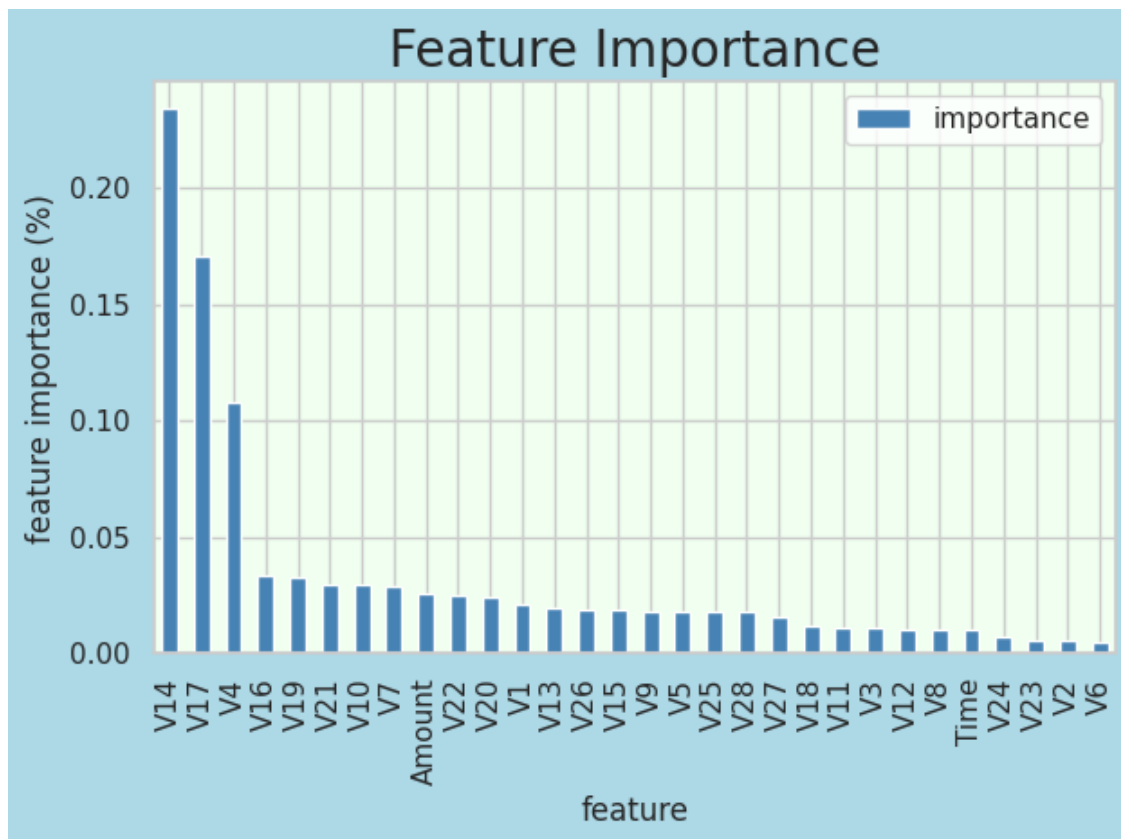
	importance
feature	
V14	0.234
V17	0.171
V4	0.108
V16	0.034
V19	0.033
V21	0.030
V10	0.030
V7	0.029
Amount	0.026
V22	0.025
V20	0.024
V1	0.021

V13	0.020
V26	0.019
V15	0.019
V9	0.018
V5	0.018
V25	0.018
V28	0.018
V27	0.016
V18	0.012
V11	0.011
V3	0.011
V12	0.010
V8	0.010
Time	0.010
V24	0.007
V23	0.006
V2	0.006
V6	0.005

```
[35]: importances.plot.bar(color='steelblue')

plt.ylabel('feature importance (%)',fontsize=12)
plt.title('Feature Importance',fontsize=20)

plt.tight_layout()
plt.gcf().patch.set_facecolor('#ADD8E6')
plt.gca().set_facecolor('honeydew')
plt.show()
```



TO AVOID OVERFITTING :Apparently, there are three dominant features: V14, V17 and V4. All the others are much less important.SO WE DROP SOME COLUMNS

#experiment :

without drop any column random_forest accuracy -0.9997 DRP -
 'V1','V2','V3','V5','V6','V7','V8','V8','V9','V11','V13','V14','V15','V16','V18','V19','V20','V21',
 'V22','V23','V24','V25','V26','V27','V28' random_forest accuracy —>0.9997

so no DIFFERENCE

[36]: *# Drop all of the features that have very similar distributions between frauds and non-frauds.*

```
X = X.drop(['V2', 'V5', 'V6', 'V7', 'V8', 'V11', 'V15', 'V16', 'V18',
            'V22', 'V23', 'V25', 'V26', 'V27', 'V28'],axis=1)
```

[37]: *# Train-test split*

```
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.
            3,random_state=0)
```

[38]: random_forest = RandomForestClassifier(class_weight='balanced')

```
random_forest.fit(X_train,y_train)
```

```
[38]: RandomForestClassifier(class_weight='balanced')
```

```
[39]: def get_test_scores(model_name:str,preds,y_test_data):

    accuracy = accuracy_score(y_test_data,preds)
    precision = precision_score(y_test_data,preds,average='macro')
    recall    = recall_score(y_test_data,preds,average='macro')
    f1        = f1_score(y_test_data,preds,average='macro')

    table = pd.DataFrame({'model': [model_name],'precision': [
↪precision],'recall': [recall],
                        'F1': [f1],'accuracy': [accuracy]})

    return table
```

```
[40]: # Use the model to predict on test data
rf_test_preds = random_forest.predict(X_test)

rf_test_results = get_test_scores('RF (test)',rf_test_preds,y_test)
rf_test_results
```

```
[40]:      model  precision  recall      F1  accuracy
0  RF (test)    0.499902     0.5  0.499951  0.999804
```

#Accuracy gives a general idea of how well the model performs overall.

#Precision is important when the cost of false positives is high (e.g., in spam detection).

#Recall is crucial when the cost of false negatives is high (e.g., in disease detection).

#F1 Score is a balanced measure that is useful when you need a balance between precision and recall, especially when the class distribution is imbalanced.

18 Resampling and Cross Validation

Undersampling

```
[41]: from imblearn.under_sampling import RandomUnderSampler

# Create a RandomUnderSampler object
rus = RandomUnderSampler(random_state=42,sampling_strategy='majority')

# Balancing the data
X_resampled,y_resampled = rus.fit_resample(X_train,y_train)# sample we used
```

Random Forest Classifier

```
[42]: random_forest.fit(X_resampled,y_resampled)
```

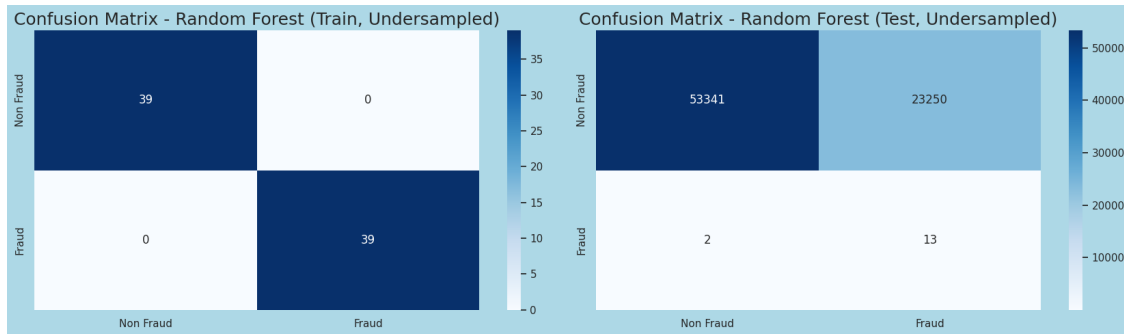
```
# Use the model to predict on train data  
rf_train_resampled_preds = random_forest.predict(X_resampled)
```

```
[43]: # Use the model to predict on test data  
rf_test_preds = random_forest.predict(X_test)
```

```
rf_test_results = get_test_scores('RF (test)',rf_test_preds,y_test)  
rf_test_results
```

```
[43]:      model  precision    recall      F1  accuracy  
0  RF (test)   0.500261   0.781553   0.411082   0.696473
```

```
[44]: fig, (ax1, ax2) = plt.subplots(ncols=2, figsize=(17, 5))  
# Generate array of values for confusion matrix  
cm1 =   
    ↪confusion_matrix(y_resampled,rf_train_resampled_preds,labels=random_forest.  
    ↪classes_)  
# Plot confusion matrix for the training set  
sns.heatmap(cm1, annot=True, ax=ax1, fmt='d', cmap='Blues')  
ax1.xaxis.set_ticklabels(['Non Fraud', 'Fraud'])  
ax1.yaxis.set_ticklabels(['Non Fraud', 'Fraud'])  
ax1.set_title('Confusion Matrix - Random Forest (Train, Undersampled)',  
    ↪fontsize=18)  
ax1.set_facecolor('honeydew')  
  
# Generate array of values for confusion matrix  
cm2 = confusion_matrix(y_test,rf_test_preds,labels=random_forest.classes_)  
# Plot confusion matrix for the test set  
sns.heatmap(cm2, annot=True, ax=ax2, fmt='d', cmap='Blues')  
ax2.xaxis.set_ticklabels(['Non Fraud', 'Fraud'])  
ax2.yaxis.set_ticklabels(['Non Fraud', 'Fraud'])  
ax2.set_title('Confusion Matrix - Random Forest (Test, Undersampled)',  
    ↪fontsize=18)  
ax2.set_facecolor('honeydew')  
  
#plot  
# Set figure background color  
plt.gcf().patch.set_facecolor('#ADD8E6')  
  
# Adjust layout  
plt.tight_layout()  
plt.show()
```



XGBoost Classifier

```
[45]: # Instantiate the XGBoost classifier
xgb = XGBClassifier(objective='binary:logistic',random_state=42)

xgb.fit(X_resampled,y_resampled)

# Use the model to predict on train data
xgb_train_resampled_preds = xgb.predict(X_resampled)
```

```
[46]: # Use the model to predict on test data
xgb_test_preds = xgb.predict(X_test)

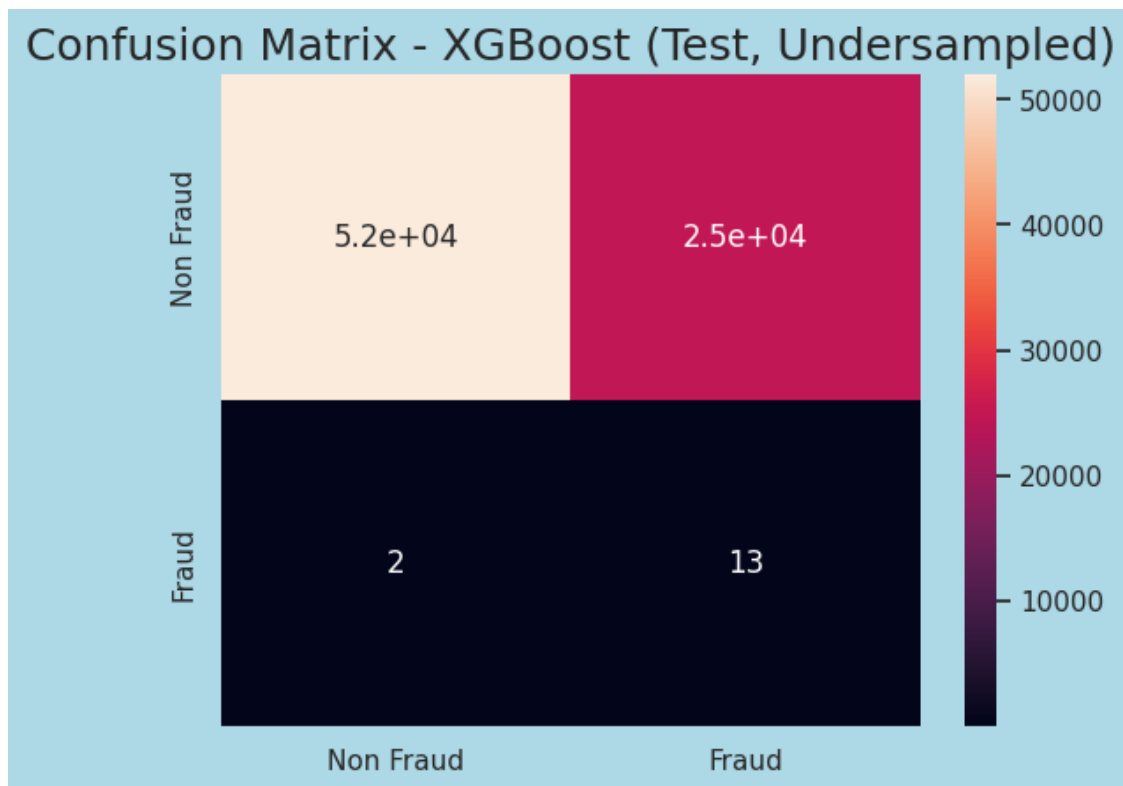
xgb_test_results = get_test_scores('XGB (test)',xgb_test_preds,y_test)
xgb_test_results
```

```
[46]:      model  precision    recall  F1   accuracy
0  XGB (test)   0.500245   0.772786  0.404895   0.678942
```

```
[47]: # Generate array of values for confusion matrix
cm = confusion_matrix(y_test,xgb_test_preds,labels=xgb.classes_)

ax = sns.heatmap(cm,annot=True)
ax.xaxis.set_ticklabels(['Non Fraud','Fraud'])
ax.yaxis.set_ticklabels(['Non Fraud','Fraud'])
ax.set_title('Confusion Matrix - XGBoost (Test, Undersampled)',fontsize=18)

plt.gcf().patch.set_facecolor('#ADD8E6')
```



Oversampling

```
[48]: from imblearn.over_sampling import SMOTE

smote = SMOTE()

# Balancing the data
X_oversampled,y_oversampled = smote.fit_resample(X_train,y_train)
```

Random Forest *Classifier*

```
[49]: random_forest.fit(X_oversampled,y_oversampled)

# Use the model to predict on train data
rf_train_oversampled_preds = random_forest.predict(X_oversampled)

rf_train_oversampled_results = get_test_scores('RF (train,␣
↪oversampled)',rf_train_oversampled_preds,y_oversampled)
rf_train_oversampled_results
```

```
[49]:
```

	model	precision	recall	F1	accuracy
0	RF (train, oversampled)	1.0	1.0	1.0	1.0

```
[50]: # Use the model to predict on test data
rf_test_preds = random_forest.predict(X_test)

rf_test_results = get_test_scores('RF (test)',rf_test_preds,y_test)
rf_test_results
```

```
[50]:          model  precision    recall      F1  accuracy
0  RF (test)    0.571337  0.533294  0.545389  0.999739
```

XGBoost Classifier

```
[51]: # 1. Instantiate the XGBoost classifier
xgb = XGBClassifier(objective='binary:logistic',random_state=42)

xgb.fit(X_oversampled,y_oversampled)
```

```
[51]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                  colsample_bylevel=None, colsample_bynode=None,
                  colsample_bytree=None, device=None, early_stopping_rounds=None,
                  enable_categorical=False, eval_metric=None, feature_types=None,
                  gamma=None, grow_policy=None, importance_type=None,
                  interaction_constraints=None, learning_rate=None, max_bin=None,
                  max_cat_threshold=None, max_cat_to_onehot=None,
                  max_delta_step=None, max_depth=None, max_leaves=None,
                  min_child_weight=None, missing=nan, monotone_constraints=None,
                  multi_strategy=None, n_estimators=None, n_jobs=None,
                  num_parallel_tree=None, random_state=42, ...)
```

```
[52]: # Use the model to predict on test data
xgb_test_preds = xgb.predict(X_test)

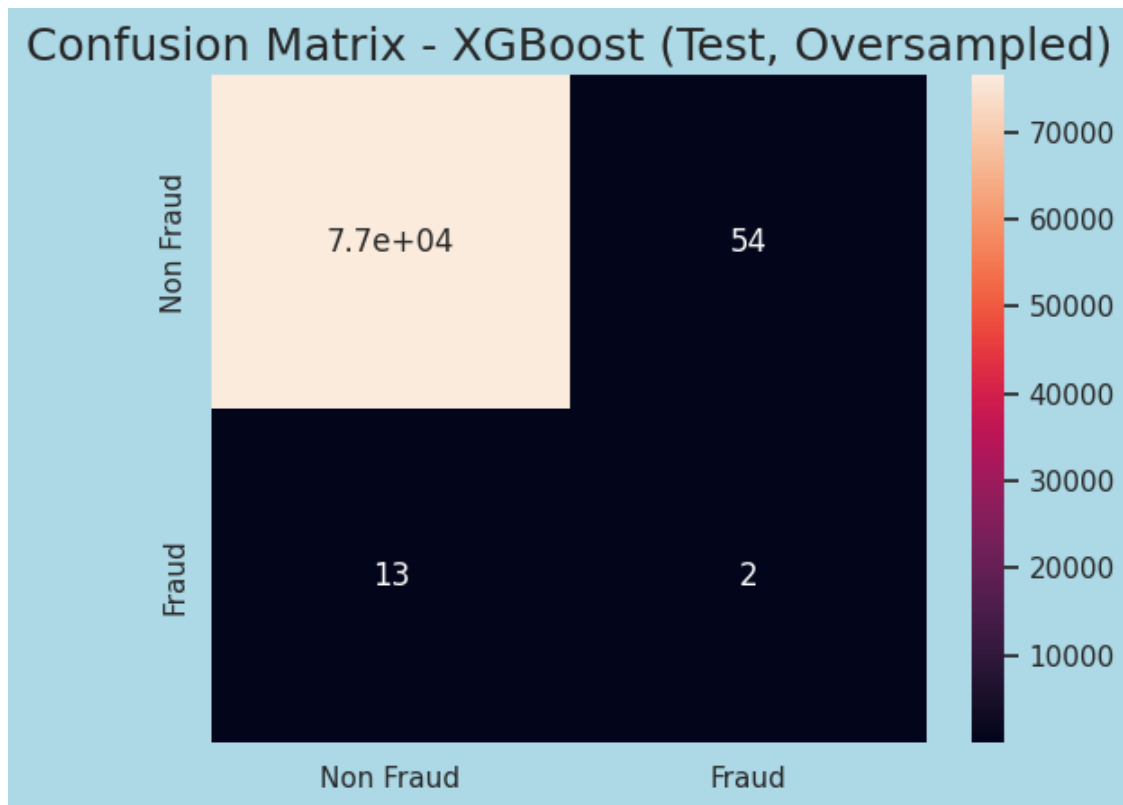
xgb_test_results = get_test_scores('XGB (test)',xgb_test_preds,y_test)
xgb_test_results
```

```
[52]:          model  precision    recall      F1  accuracy
0  XGB (test)    0.517772  0.566314  0.52795  0.999125
```

```
[53]: # Generate array of values for confusion matrix
cm = confusion_matrix(y_test,xgb_test_preds,labels=xgb.classes_)

ax = sns.heatmap(cm,annot=True)
ax.xaxis.set_ticklabels(['Non Fraud','Fraud'])
ax.yaxis.set_ticklabels(['Non Fraud','Fraud'])
ax.set_title('Confusion Matrix - XGBoost (Test, Oversampled)',fontsize=18)

plt.gcf().patch.set_facecolor('#ADD8E6')
```



It seems that undersampling improves the classification of the minority class (frauds), while over-sampling tends to do better on the majority class (non frauds, i.e. regular transactions).

In light of this, I will keep the undersampling strategy and try to improve its results by also implementing other strategies.

since undersampling random forest give best recall score so we cross validate to this and check score

```
[54]: # K-fold stratified cross validation
kfold = StratifiedKFold(n_splits=10)
kfold
```

```
[54]: StratifiedKFold(n_splits=10, random_state=None, shuffle=False)
```

```
[55]: #GRID SEARCH CV TAKE TOO MUCH TIM SO I COMMENT OUT IF U HAVE 6 CORE PROCESSOR
      ↳ THEN IT WILL BE DONE IN 20 MIN. ACCURACY IS NEAR TO -->0.77

# RFC Parameters tuning
# RFC = RandomForestClassifier(random_state=42)

# rf_param_grid = {
#     'max_depth': [2,3,4,5,None],
```



```

#     'max_features': [1.0],
#     'max_samples': [1.0],
#     'min_samples_leaf': [2,3,4],
#     'min_samples_split': [2,2,4],
#     'n_estimators': [200,300,400]
# }

# gsRFC = GridSearchCV(RFC,param_grid=rf_param_grid,cv=kfold,scoring="f1")

# gsRFC.fit(X_resampled,y_resampled)

#FOR TIME TAKING I USE RANDOM SEARCH CV
from sklearn.model_selection import RandomizedSearchCV
RFC = RandomForestClassifier(random_state=42)
# Define the Stratified K-Folds cross-validator
kfold = StratifiedKFold(n_splits=10)

# Define the parameter grid for RandomizedSearchCV
rf_param_dist = {
    'max_depth': [2, 3, 4, 5, None],
    'max_features': [1.0],
    'max_samples': [1.0],
    'min_samples_leaf': [2, 3, 4],
    'min_samples_split': [2, 4],
    'n_estimators': [200, 300, 400]
}

# Initialize RandomizedSearchCV
rsRFC = RandomizedSearchCV(RFC, param_distributions=rf_param_dist, n_iter=100,
    cv=kfold, scoring="f1", random_state=42)

# Fit the model
rsRFC.fit(X_resampled, y_resampled)

# Optionally, you can retrieve the best parameters
best_params = rsRFC.best_params_
print("Best parameters found: ", best_params)

```

Best parameters found: {'n_estimators': 200, 'min_samples_split': 2, 'min_samples_leaf': 2, 'max_samples': 1.0, 'max_features': 1.0, 'max_depth': 4}

```

[56]: # Use the model to predict on test data
rf_test_preds = rsRFC.predict(X_test)

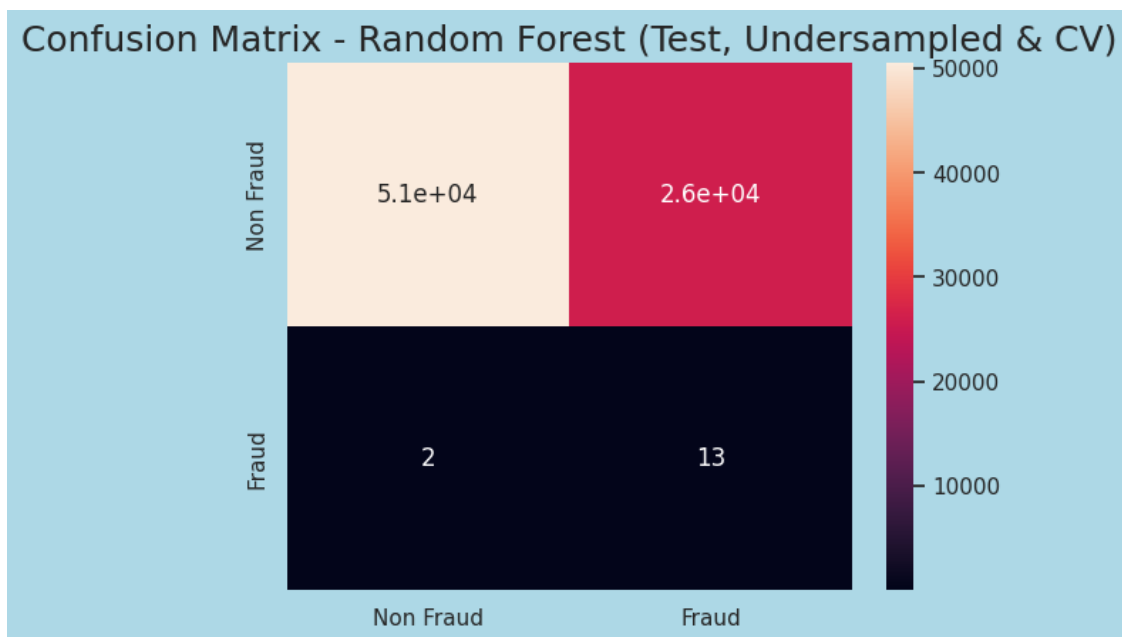
rf_test_results = get_test_scores('RF (test)',rf_test_preds,y_test)
rf_test_results

```

```
[56]:      model  precision    recall  F1   accuracy
0  RF (test)    0.50023  0.763366  0.398106  0.660105
```

```
[57]: # Generate array of values for confusion matrix
cm = confusion_matrix(y_test, rf_test_preds, labels=rsRFC.classes_)

ax = sns.heatmap(cm, annot=True)
ax.xaxis.set_ticklabels(['Non Fraud', 'Fraud'])
ax.yaxis.set_ticklabels(['Non Fraud', 'Fraud'])
ax.set_title('Confusion Matrix - Random Forest (Test, Undersampled & CV)')
plt.gcf().patch.set_facecolor('#ADD8E6')
```



```
[58]: !pip install catboost

from sklearn.ensemble import VotingClassifier
from sklearn.ensemble import RandomForestClassifier
from lightgbm import LGBMClassifier
from catboost import CatBoostClassifier
from xgboost import XGBClassifier

# Instantiate individual classifiers
rf = RandomForestClassifier(n_estimators=100, random_state=42)
lgbm = LGBMClassifier(objective='binary', random_state=42)
xgb = XGBClassifier(objective='binary:logistic', random_state=42)
```

```

catboost = CatBoostClassifier(iterations=1000, learning_rate=0.1, depth=6,
    ↪random_seed=42, silent=True)

# Instantiate the Voting Classifier
voting_clf = VotingClassifier(
    estimators=[
        ('rf', rf),
        ('lgbm', lgbm),
        ('xgb', xgb),
        ('catboost', catboost)
    ],
    voting='hard' # Use 'soft' for soft voting
)

# Fit the Voting Classifier on the resampled training data
voting_clf.fit(X_resampled, y_resampled)

# Use the model to predict on train data
voting_train_resampled_preds = voting_clf.predict(X_resampled)

# Use the model to predict on test data
voting_test_preds = voting_clf.predict(X_test)

# Get test results using your scoring function
voting_test_results = get_test_scores('Voting Classifier (test)',
    ↪voting_test_preds, y_test)
voting_test_results

```

Collecting catboost

Downloading catboost-1.2.7-cp310-cp310-manylinux2014_x86_64.whl.metadata (1.2 kB)

Requirement already satisfied: graphviz in /usr/local/lib/python3.10/dist-packages (from catboost) (0.20.3)

Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from catboost) (3.8.0)

Requirement already satisfied: numpy<2.0,>=1.16.0 in /usr/local/lib/python3.10/dist-packages (from catboost) (1.26.4)

Requirement already satisfied: pandas>=0.24 in /usr/local/lib/python3.10/dist-packages (from catboost) (2.2.2)

Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from catboost) (1.13.1)

Requirement already satisfied: plotly in /usr/local/lib/python3.10/dist-packages (from catboost) (5.24.1)

Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from catboost) (1.16.0)

Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24->catboost) (2.8.2)

Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24->catboost) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.24->catboost) (2024.2)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (1.3.0)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (4.54.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (1.4.7)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (24.1)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (10.4.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib->catboost) (3.2.0)
Requirement already satisfied: tenacity>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from plotly->catboost) (9.0.0)
Downloading catboost-1.2.7-cp310-cp310-manylinux2014_x86_64.whl (98.7 MB)
98.7/98.7 MB

6.3 MB/s eta 0:00:00

Installing collected packages: catboost

Successfully installed catboost-1.2.7

[LightGBM] [Info] Number of positive: 39, number of negative: 39

[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000929 seconds.

You can set `force_col_wise=true` to remove the overhead.

[LightGBM] [Info] Total Bins 411

[LightGBM] [Info] Number of data points in the train set: 78, number of used features: 15

[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.500000 -> initscore=0.000000

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf

[illegible]

[illegible]

[58]:		model	precision	recall	F1	accuracy
0	Voting Classifier	(test)	0.500247	0.739589	0.427763	0.745842

- 19 Till now Best Recall Score is \rightarrow under sampling , random forest \rightarrow recall score of 0.78 for your credit card fraud detection model using Random Forest with undersampling! A high recall is crucial in fraud detection, as it indicates that your model is effectively identifying fraudulent transactions.
- 20 Accuracy \rightarrow 0.999804

[58] :