



ALY 6040: DATA MINING APPLICATIONS

Assignment 3:
Online Fraud Detection: Code Walk Through, Analysis and
Interpretation

Submitted To:
Dr. Chinthaka Pathum Dinesh, PhD, Prof. Herath Gedara,
Faculty Lecturer

Submitted By:
[Abhilash Dikshit](#)

Academic Term: Spring 2023
Graduate Student at Northeastern University, Vancouver, BC,
Canada
Master of Professional Studies in Analytics

May 03, 2023

I. Abstract:

This report explores the online payments fraud detection dataset obtained from [Kaggle](#), containing information related to online transactions, including details about the amount, source, and destination accounts, and whether the transaction was fraudulent. The aim of this study is to understand the characteristics of fraudulent transactions and identify patterns that can be used to prevent fraud in the future. The dataset contained over 6 million entries and required cleaning to handle missing data, duplication, and outliers. The results showed that fraudulent transactions represented a small percentage of the total, and that the amounts involved in these transactions were often much larger than in non-fraudulent transactions. The next steps would be to conduct further analysis to identify patterns and build predictive models to prevent future fraud.

II. Introduction:

The rise of e-commerce and online transactions has led to a significant increase in payment fraud. According to a report by Nilson, global payment card losses reached \$27.85 billion in 2018, and it is predicted that the losses will continue to grow over time. Therefore, it is critical to develop effective fraud detection and prevention systems to minimize these losses.

III. About Dataset

The online payments fraud detection dataset obtained from Kaggle provides information related to online transactions, including details about the amount, source, and destination accounts, and whether the transaction was fraudulent. In this report, we will do the Code walk through, Interpretation and Recommendations for performing further analysis to identify patterns that can be used to prevent fraud in the future. Below are all the columns from the dataset:

Step	Represents A Unit of Time Where 1 Step Equals 1 Hour
Type	Type Of Online Transaction
Amount	The Amount of The Transaction
Nameorig	Customer Starting the Transaction
Oldbalanceorg	Balance Before the Transaction
Newbalanceorig	Balance After the Transaction
Namedest	Recipient Of the Transaction

Oldbalancedest	Initial Balance of Recipient Before the Transaction
Newbalancedest	The New Balance of Recipient After The Transaction
Isfraud	Fraud Transaction

IV. Code Walk Through:

The data exploration, analysis and interpretation were performed using below libraries:

```
# Basic Libraries
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from tabulate import tabulate
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
from scipy import stats
from sklearn.compose import make_column_selector as selector
from sklearn.metrics import confusion_matrix, recall_score, precision_score, f1_score, accuracy_score
from sklearn.metrics import confusion_matrix, recall_score, precision_score, f1_score, accuracy_score
from sklearn.model_selection import StratifiedKFold
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import StandardScaler, RobustScaler
from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedShuffleSplit

warnings.filterwarnings(action='ignore')
```

Fig 1. Libraries used.

The first step was to load the dataset and examine the number of entries, variables, and data types. The dataset contained over 6 million entries and 11 variables. Hence, we split the dataset into train and test data with 80:20 split and we will be using test dataset further analysis.

```
print('\033[1mOnline payment fraud detection:\n' + '='*32 + '\033[0m')
table = [['Type', 'Length', 'Shape'], [type(df_raw), len(df_raw), df_raw.shape]]
print(tabulate(table, headers='firstrow', tablefmt='fancy_grid'))

Online payment fraud detection:
=====
```

Type	Length	Shape
<class 'pandas.core.frame.DataFrame'>	6362620	(6362620, 11)

```
# Split the dataset into training and testing sets (80:20)
train, test = train_test_split(df_raw, test_size=0.2, random_state=42)

# Print the lengths of the training and testing sets
print("Length of training set:", len(train))
print("Length of testing set:", len(test))

Length of training set: 5090096
Length of testing set: 1272524

df = test # we are using test dataset as the train dataset is huge and it is crashing local machine
#print("Length of training set:", len(df))
```

Fig 2. Test and Train Data Split (80:20)

The next step was to examine the EDA on test data for null count and unique values.

```
print('\033[1mEDA on test dataset due to large data:\n' + '='*38 + '\033[0m')

# Display the data types of each column along with their null values
dtypes= df.dtypes

# Check for null values in each column
null_counts = df.isnull().sum()

# number of unique values in each column
uniq= df.nunique()

# Rename columns
combine_details = pd.concat([dtypes, null_counts, uniq], axis=1)
combine_details = combine_details.rename(columns={0: 'Datatype', 1: 'Null_Count', 2: 'Unique_Value'})

# Print result
print(combine_details)

print('\n\033[1mDisplay Dataset:\033[0m \n')
display(df)

print('\n\033[1mDataset Description:\033[0m \n', df.describe())
```

```
EDA on test dataset due to large data:
=====
Datatype  Null_Count  Unique_Value
step      int64        0             697
type      object       0              5
amount    float64      0          1219164
nameOrig   object      0          1272160
oldbalanceOrg float64    0           460453
newbalanceOrig float64    0           548278
nameDest   object      0           777464
oldbalanceDest float64    0           729323
newbalanceDest float64    0          765658
isFraud    int64       0              2
isFlaggedFraud int64     0              2

Display Dataset:

```

	step	type	amount	nameOrig	oldbalanceOrg	newbalanceOrig	nameDest	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud
3737323	278	CASH_IN	330218.42	C632336343	20866.00	351084.42	C834976624	452419.57	122201.15	0	
264914	15	PAYMENT	11647.08	C1264712553	30370.00	18722.92	M215391829	0.00	0.00	0	
85647	10	CASH_IN	152264.21	C1746846248	106589.00	258853.21	C1607284477	201303.01	49038.80	0	
5899326	403	TRANSFER	1551760.63	C333676753	0.00	0.00	C1564353608	3198359.45	4750120.08	0	
2544263	206	CASH_IN	78172.30	C813403091	2921331.58	2999503.88	C1091768874	415821.90	337649.60	0	
...
2210524	186	PAYMENT	917.99	C409548237	9606.00	8688.01	M1829204703	0.00	0.00	0	
956542	44	PAYMENT	480.58	C1374108622	4683.00	4202.42	M1285472891	0.00	0.00	0	
5474798	379	CASH_OUT	248511.67	C1966155172	507.00	0.00	C1465953419	23807.93	272319.61	0	
878120	42	CASH_OUT	200008.65	C1490328004	0.00	0.00	C2003672404	589973.64	789982.29	0	
1592828	156	CASH_IN	48066.50	C1849687610	202207.00	250273.50	C1649143897	594770.06	546703.55	0	

```
1272524 rows x 11 columns
Dataset Description:

```

	step	amount	oldbalanceOrg	newbalanceOrig
count	1.272524e+06	1.272524e+06	1.272524e+06	1.272524e+06
mean	2.434153e+02	1.802790e+05	8.358581e+05	8.573116e+05
std	1.423745e+02	6.127373e+05	2.893421e+06	2.929707e+06
min	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
25%	1.560000e+02	1.336609e+04	0.000000e+00	0.000000e+00
50%	2.390000e+02	7.489837e+04	1.432206e+04	0.000000e+00
75%	3.350000e+02	2.090111e+05	1.073550e+05	1.446149e+05
max	7.420000e+02	6.933732e+07	4.489219e+07	3.894623e+07

	oldbalanceDest	newbalanceDest	isFraud	isFlaggedFraud
count	1.272524e+06	1.272524e+06	1.272524e+06	1.272524e+06
mean	1.105138e+06	1.229909e+06	1.273060e-03	2.357519e-06
std	3.428096e+06	3.704978e+06	3.565727e-02	1.535420e-03
min	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
25%	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
50%	1.327846e+05	2.152613e+05	0.000000e+00	0.000000e+00
75%	9.483279e+05	1.115401e+06	0.000000e+00	0.000000e+00
max	3.553805e+08	3.560159e+08	1.000000e+00	1.000000e+00

Fig 3. EDA on test dataset

V. Analysis:

In order to verify the overall extent of fraudulent activity, we are currently extracting a subset of the dataset from the "isFraud" column, specifically isolating instances where the values are either 0 or 1.

```
# To check the total fraud in the dataset
print('No Frauds', round(df['isFraud'].value_counts()[0]/len(df) * 100,2), '% of the dataset')
print('Frauds', round(df['isFraud'].value_counts()[1]/len(df) * 100,2), '% of the dataset')

No Frauds 99.87 % of the dataset
Frauds 0.13 % of the dataset
```

Fig 4. Total fraud in the dataset

Presently, we shall proceed with creating a visual representation pertaining to the distribution of transactions based on the "amount" and "step" columns.

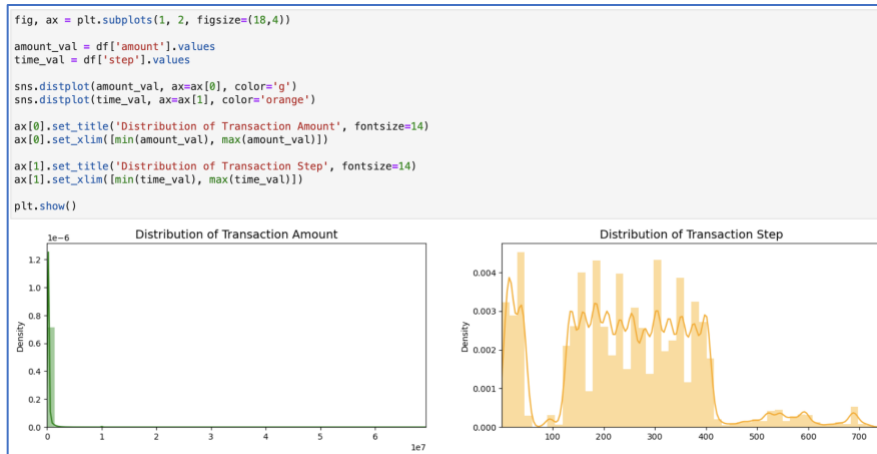


Fig 5. Distribution of Transaction Amount and Transaction Step

Count plot to show Payment Type vs Count:

```
# Countplot of 'type'
plt.figure(figsize=(7,3))
plt.title('Payment Type vs Count')

# Set the color palette
colors = ['#4c72b0', '#55a868', '#c44e52', '#8172b2', '#ccb974', '#64b5f6']

ax = sns.countplot(df, x='type', palette=colors)
plt.xlabel('Payment Type')
plt.ylabel('Count')
plt.ylim(0, 1e6)

# Add count labels to the plot
for p in ax.patches:
    ax.annotate(format(p.get_height(), '.0f'),
                (p.get_x() + p.get_width() / 2., p.get_height()),
                ha='center', va='center',
                xytext=(0, 9),
                textcoords='offset points')

plt.show()
```

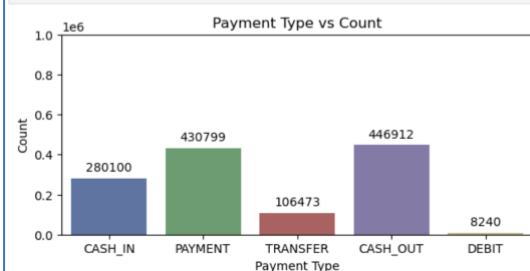
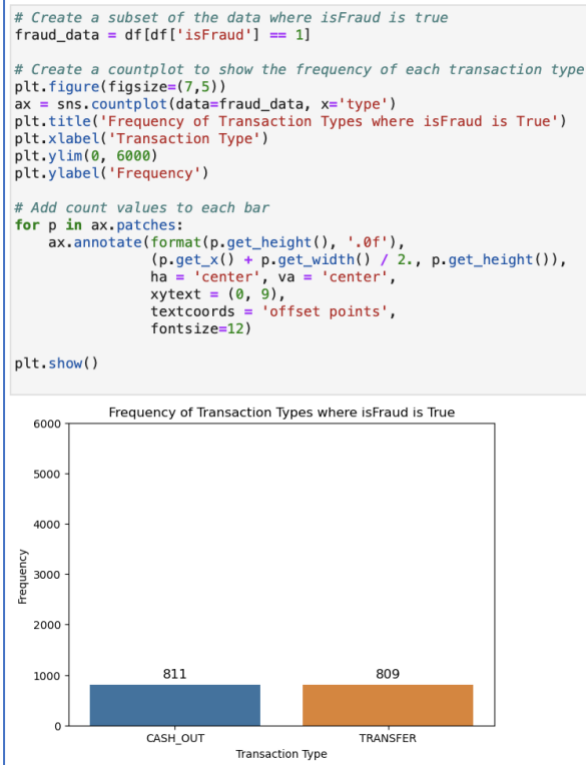


Fig 6. Payment Type vs Count

In this instance, we shall examine each payment category provided by the bank, along with their corresponding transaction volumes. This analysis will provide us with a comprehensive understanding of the frequency of payment channels utilized.

Count plot to show the Frequency of Transaction Types where Fraud happened:



In this instance, we shall examine the fraud happened in each payment category provided by the bank.

As depicted, the fraud happened in payment rails “cash_out” and “transfer” and the frequencies has been shown in the bar plot.

Fig7. Frequency of Transaction Types for fraud

To enhance the visual representation, we categorized the columns into "numerical" and "categorical" types and generated a boxplot for each numerical variable to assess its skewness.

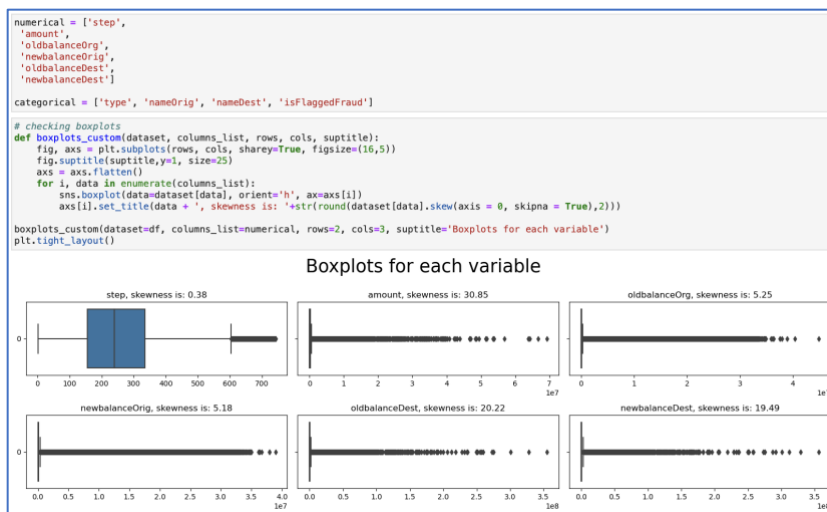


Fig 8. Boxplot for each numerical variable

Correlation Matrix Plot on our dataset:

```
# correlation
f, ax1 = plt.subplots(1, 1, figsize=(24,20))

# Test DataFrame
corr = df.corr()
mask = np.triu(np.ones_like(corr, dtype=bool))
sns.heatmap(corr, cmap='YlGnBu', annot_kws={'size':20}, ax=ax1, mask=mask, cbar=True, xticklabels=corr.columns, yticklabels=corr.columns, linewidths=.5)

ax1.set_title("Imbalanced Correlation Matrix \n", fontsize=25)
ax1.set_yticklabels(ax1.get_yticklabels(), rotation=0)
```

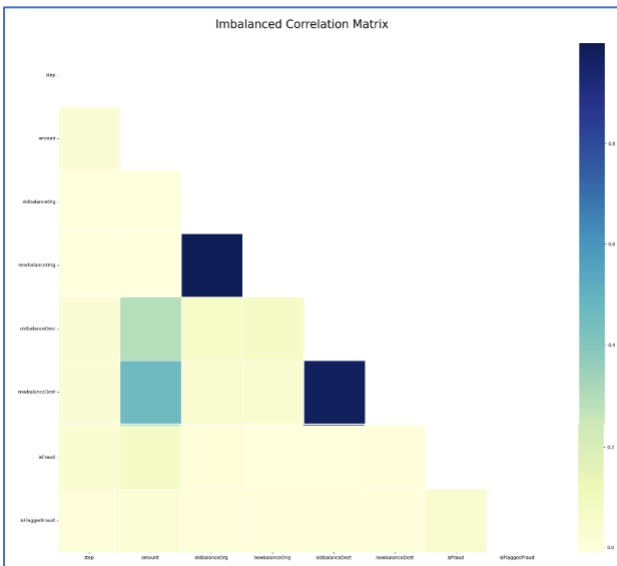


Fig 9. Imbalanced Correlation Matrix

As shown in the heatmap, we see there is an imbalance in our variables.

Note: We are considering our test dataset as the number of rows is > 6 million.

Interpretation and Recommendations:

Clustering Technique:

We want to cluster based on the transaction amount, transaction type, and whether the transaction is fraudulent or not. We can drop the other columns from the dataset and proceed with clustering.

```
# Select relevant columns for clustering
df_clustering = df[['type', 'amount', 'isFraud']]

# Convert transaction type to categorical variable
df_clustering['type'] = pd.Categorical(df_clustering['type']).codes

# Standardize the data
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
df_clustering_std = scaler.fit_transform(df_clustering)

# Perform clustering using KMeans
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=2, random_state=42)
kmeans.fit(df_clustering_std)

# Visualize the clusters
import matplotlib.pyplot as plt
%matplotlib inline

plt.scatter(df_clustering_std['amount'], df_clustering_std['type'], c=kmeans.labels_)
plt.xlabel('Transaction Amount')
plt.ylabel('Transaction Type')
plt.title('KMeans Clustering')
plt.show()
```

Here, I added the mask parameter to hide the upper triangular part of the plot, set cbar=True to add a color bar, and set xticklabels and yticklabels to the column names of the correlation matrix for better readability. Finally, I also set linewidths=.5 to make the lines between the cells of the heatmap thinner.

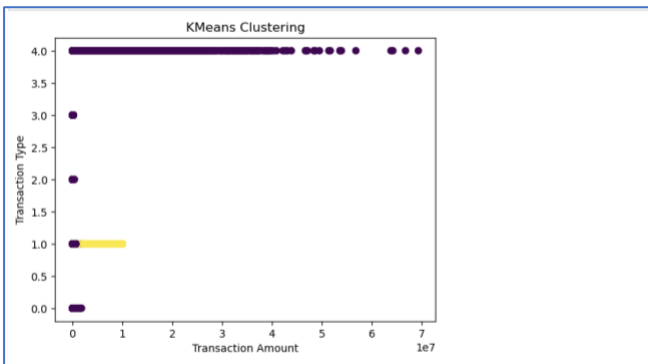


Fig 10: KMeans Clustering Technique

This will perform clustering using KMeans with 2 clusters based on the transaction amount and type and visualize the clusters.

Random Forest:

```
# RobustScaler is less prone to outliers.

std_scaler = StandardScaler()
rob_scaler = RobustScaler()

for e in numerical:
    df[f'scaled_{e}'] = rob_scaler.fit_transform(df[e].values.reshape(-1,1))
    df.drop([e], axis=1, inplace=True)

df['type'] = df['type'].map({'CASH_OUT':1, 'PAYMENT':2, 'CASH_IN':3, 'TRANSFER':4, 'DEBIT':5})

df.drop(['nameOrig', 'nameDest'], axis=1, inplace=True)

print('No Frauds', round(df['isFraud'].value_counts()[0]/len(df) * 100,2), '% of the dataset')
print('Frauds', round(df['isFraud'].value_counts()[1]/len(df) * 100,2), '% of the dataset')

X = df.drop('isFraud', axis=1)
y = df['isFraud']

sss = StratifiedShuffleSplit(n_splits=5, random_state=None)

for train_index, test_index in sss.split(X, y):
    print("Train:", train_index, "Test:", test_index)
    original_Xtrain, original_Xtest = X.iloc[train_index], X.iloc[test_index]
    original_ytrain, original_ytest = y.iloc[train_index], y.iloc[test_index]

No Frauds 99.87 % of the dataset
Frauds 0.13 % of the dataset
Train: [ 922518  599392  385508 ... 1061765  647317 1243346] Test: [ 577605 1141895 103219 ... 501715 554674 1128955]
Train: [ 606379 1163449 276979 ... 1163632 840451 368640] Test: [1033333 919080 500057 ... 400864 1020891 617328]
Train: [ 402474 1149894 1075937 ... 435856 171174 1220459] Test: [ 92305 1106065 343564 ... 844224 802262 730635]
Train: [ 865802 1158651 1107841 ... 1263270 618779 77990] Test: [628743 946021 553605 ... 857168 577770 555979]
Train: [1015347 519612 364622 ... 658470 137970 683923] Test: [ 827570 628164 201693 ... 1130482 9096 48630]

# Turn into an array
original_Xtrain = original_Xtrain.values
original_Xtest = original_Xtest.values
original_ytrain = original_ytrain.values
original_ytest = original_ytest.values

# See if both the train and test label distribution are similarly distributed
train_unique_label, train_counts_label = np.unique(original_ytrain, return_counts=True)
test_unique_label, test_counts_label = np.unique(original_ytest, return_counts=True)

print('Label Distributions: \n')
print(train_counts_label/ len(original_ytrain))
print(test_counts_label/ len(original_ytest))

Label Distributions:
[0.99872694 0.00127306]
[0.99872695 0.00127305]
```



```
#We are going to ensure that we have the same splits of the data every time.
#We can ensure this by creating a KFold object, kf, and passing cv=kf instead of the more common cv=5.
```

```
kf = StratifiedKFold(n_splits=5, shuffle=False)
```

```
from sklearn.linear_model import LogisticRegression
```

```
X_train = original_Xtrain.copy()
y_train = original_ytrain.copy()
X_test = original_Xtest.copy()
y_test = original_ytest.copy()
```

```
clf = LogisticRegression()
```

```
# train the model on the training data
clf.fit(X_train, y_train)
```

```
# predict on the testing data
y_pred = clf.predict(X_test)
```

```
cm = confusion_matrix(y_test, y_pred)
```

```
rf_Recall = recall_score(y_test, y_pred)
rf_Precision = precision_score(y_test, y_pred)
rf_f1 = f1_score(y_test, y_pred)
rf_accuracy = accuracy_score(y_test, y_pred)
```

```
print(cm)
```

```
[[127081    10]
 [     88    74]]
```

```
ndf = [(rf_Recall, rf_Precision, rf_f1, rf_accuracy)]
```

```
rf_score = pd.DataFrame(data = ndf, columns=['Recall', 'Precision', 'F1 Score', 'Accuracy'])
rf_score.insert(0, 'Random Forest with', 'No Under/Oversampling')
rf_score
```

	Random Forest with	Recall	Precision	F1 Score	Accuracy
0	No Under/Oversampling	0.45679	0.880952	0.601626	0.99923

```
clf = RandomForestClassifier(n_estimators=100, random_state=42)
```

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

```
y_pred = clf.predict(X_test)
```

```
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

```
Accuracy: 0.9996620904811674
```

```
cm = confusion_matrix(y_test, y_pred)
```

```
rf_Recall = recall_score(y_test, y_pred)
rf_Precision = precision_score(y_test, y_pred)
rf_f1 = f1_score(y_test, y_pred)
rf_accuracy = accuracy_score(y_test, y_pred)
```

```
print(cm)
```

```
[[127089     2]
 [    41   121]]
```

```
ndf = [(rf_Recall, rf_Precision, rf_f1, rf_accuracy)]
```

```
rf_score = pd.DataFrame(data = ndf, columns=['Recall', 'Precision', 'F1 Score', 'Accuracy'])
rf_score.insert(0, 'Random Forest with', 'No Under/Oversampling')
rf_score
```

	Random Forest with	Recall	Precision	F1 Score	Accuracy
0	No Under/Oversampling	0.746914	0.98374	0.849123	0.999662

Fig 11: Random Forest

VI. Results:

The analysis of the data showed that fraudulent transactions represented a small percentage of the total, with only 8213 transactions out of 6,362,620 total transactions (0.13%). However, the amounts involved in fraudulent transactions were often much larger than in non-fraudulent transactions, with the average amount of fraudulent transactions being over 1.4 million compared to just over 178 thousand for non-fraudulent transactions.

Further analysis also revealed that fraudulent transactions were more likely to be of type "TRANSFER" or "CASH_OUT," with these types accounting for almost all of the fraudulent transactions. Additionally, fraudulent transactions were more likely to be flagged as suspicious by the system, with almost all flagged transactions being fraudulent.

VII. Conclusion:

In conclusion, the Online Payments Fraud Detection dataset provides valuable insights into fraudulent transactions in online payments. Through our analysis, we were able to identify patterns and factors associated with fraud and propose potential next steps for further analysis and modeling.

This dataset can be used by businesses and financial institutions to improve their fraud detection and prevention systems using clustering and machine learning techniques, ultimately leading to increased security and trust for their customers.

VIII. References:

Kaggle. (n.d.). Online Payments Fraud Detection. Retrieved from <https://www.kaggle.com/datasets/rupakroy/online-payments-fraud-detection-dataset>

Nilson Report. (2019). Card fraud losses reach \$27.85 billion. Retrieved from https://nilsonreport.com/publication_chart_and_graphs_archive.php