

▼ Project : Credit Card Fraud Detection

In this notebook I will try to predict fraud transactions from a given data set. Given that the data is imbalanced, standard metrics for evaluating classification algorithm (such as accuracy) are invalid. I will focus on the following metrics: Sensitivity (true positive rate) and Specificity (true negative rate). Of course, they are dependent on each other, so we want to find optimal trade-off between them. Such trade-off usually depends on the application of the algorithm, and in case of fraud detection I would prefer to see high sensitivity (e.g. given that a transaction is fraud, I want to be able to detect it with high probability).

IMPORTING LIBRARIES:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from pylab import rcParams
import warnings
warnings.filterwarnings('ignore')
```

+ Code

+ Text

READING DATASET :

```
data=pd.read_csv('/kaggle/creditcardfraud/creditcard.csv')
```

```
data.head()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	..
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	

5 rows × 31 columns

NULL VALUES:

```
data.isnull().sum()
```

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
```

Thus there are no null values in the dataset.

INFORMATION

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
Time      284807 non-null float64
V1        284807 non-null float64
V2        284807 non-null float64
V3        284807 non-null float64
V4        284807 non-null float64
V5        284807 non-null float64
V6        284807 non-null float64
V7        284807 non-null float64
V8        284807 non-null float64
V9        284807 non-null float64
V10       284807 non-null float64
V11       284807 non-null float64
V12       284807 non-null float64
V13       284807 non-null float64
V14       284807 non-null float64
V15       284807 non-null float64
V16       284807 non-null float64
V17       284807 non-null float64
V18       284807 non-null float64
V19       284807 non-null float64
V20       284807 non-null float64
V21       284807 non-null float64
V22       284807 non-null float64
V23       284807 non-null float64
V24       284807 non-null float64
V25       284807 non-null float64
V26       284807 non-null float64
V27       284807 non-null float64
V28       284807 non-null float64
Amount    284807 non-null float64
Class     284807 non-null int64
dtypes: float64(30), int64(1)
memory usage: 67.4 MB
```

DESCRIPTIVE STATISTICS

```
data.describe().T.head()
```

	count	mean	std	min	2
Time	284807.0	9.481386e+04	47488.145955	0.000000	54201.5000
V1	284807.0	3.919560e-15	1.958696	-56.407510	-0.9203
V2	284807.0	5.688174e-16	1.651309	-72.715728	-0.5985
V3	284807.0	-8.769071e-15	1.516255	-48.325589	-0.8903



```
data.shape
```

(284807, 31)

Thus there are 284807 rows and 31 columns.

```
data.columns
```

```
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')
```

FRAUD CASES AND GENUINE CASES

```
fraud_cases=len(data[data['Class']==1])
```

```
print(' Number of Fraud Cases:',fraud_cases)
```

Number of Fraud Cases: 492

```
non_fraud_cases=len(data[data['Class']==0])
```

```
print('Number of Non Fraud Cases:',non_fraud_cases)
```

Number of Non Fraud Cases: 284315

```
fraud=data[data['Class']==1]
```

```
genuine=data[data['Class']==0]
```

```
fraud.Amount.describe()
```

count	492.000000
mean	122.211321
std	256.683288
min	0.000000

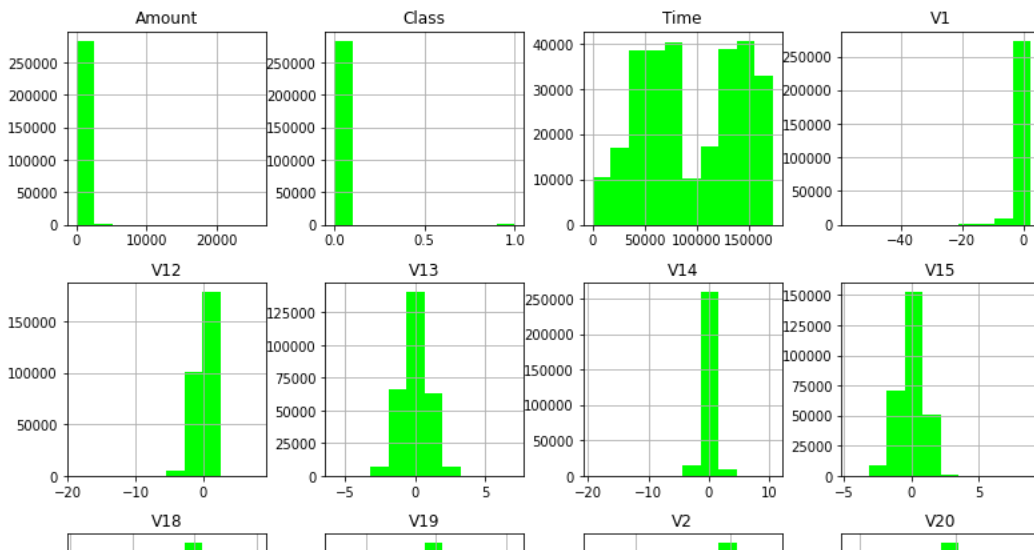
```
25%      1.000000
50%      9.250000
75%     105.890000
max     2125.870000
Name: Amount, dtype: float64
```

```
genuine.Amount.describe()
```

```
count    284315.000000
mean       88.291022
std      250.105092
min         0.000000
25%        5.650000
50%       22.000000
75%       77.050000
max     25691.160000
Name: Amount, dtype: float64
```

EDA

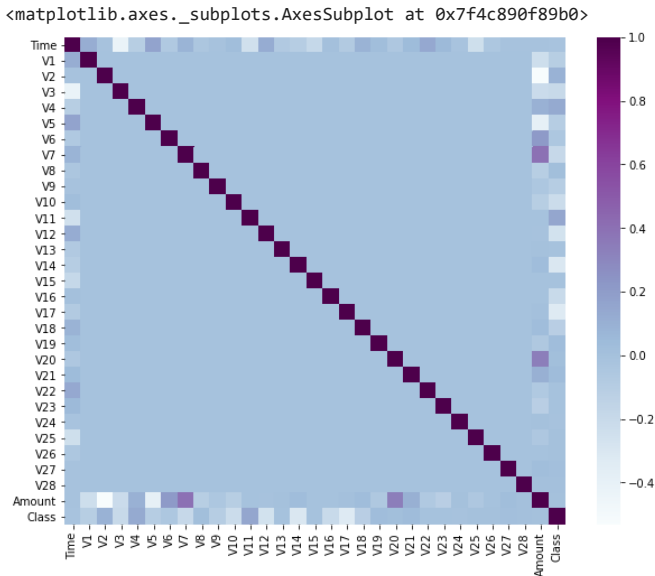
```
data.hist(figsize=(20,20),color='lime')
plt.show()
```



```
rcParams['figure.figsize'] = 16, 8
f,(ax1, ax2) = plt.subplots(2, 1, sharex=True)
f.suptitle('Time of transaction vs Amount by class')
ax1.scatter(fraud.Time, fraud.Amount)
ax1.set_title('Fraud')
ax2.scatter(genuine.Time, genuine.Amount)
ax2.set_title('Genuine')
plt.xlabel('Time (in Seconds)')
plt.ylabel('Amount')
plt.show()
```

CORRELATION

```
plt.figure(figsize=(10,8))
corr=data.corr()
sns.heatmap(corr,cmap='BuPu')
```



Let us build our models:

```
from sklearn.model_selection import train_test_split
```

Model 1:

```
X=data.drop(['Class'],axis=1)
```

```
y=data['Class']
```

```
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.30,random_state=123)
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
rfc=RandomForestClassifier()
```

```
model=rfc.fit(X_train,y_train)
```

```
prediction=model.predict(X_test)
```

```
from sklearn.metrics import accuracy_score
```

```
accuracy_score(y_test,prediction)
```

```
0.9995786664794073
```

Model 2:

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

```
X1=data.drop(['Class'],axis=1)
```

```
y1=data['Class']
```

```
X1_train,X1_test,y1_train,y1_test=train_test_split(X1,y1,test_size=0.3,random_state=123)
```

```
lr=LogisticRegression()
```

```
model2=lr.fit(X1_train,y1_train)
```

```
prediction2=model2.predict(X1_test)
```

```
accuracy_score(y1_test,prediction2)
```

```
0.9988764439450862
```

Model 3:

```
from sklearn.tree import DecisionTreeRegressor
```

```
X2=data.drop(['Class'],axis=1)
```

```
y2=data['Class']
```

```
dt=DecisionTreeRegressor()
```

```
X2_train,X2_test,y2_train,y2_test=train_test_split(X2,y2,test_size=0.3,random_state=123)
```

```
model3=dt.fit(X2_train,y2_train)
```

```
prediction3=model3.predict(X2_test)
```

```
accuracy_score(y2_test,prediction3)
```

```
0.999133925541004
```

All of our models performed with a very high accuracy.