

```
#CREDIT CARD FRAUD DETECTION USING MACHINE LEARNING
```

```
#1.IMPORT THE NECESSARY LIBRARIES
```

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
#2.NOW LOAD THE DATA: DATASET-----> FROM KAGGLE
```

```
df=pd.read_csv('creditcard.csv')#data in correct format so no encoding is required
df
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.018307	0.277838	-0.110474
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321
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
...
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.914428	...	0.213454	0.111864	0.057504
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.584800	...	0.214205	0.924384	-0.204233
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.432454	...	0.232045	0.578229	-0.120794
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.392087	...	0.265245	0.800049	-0.063375
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.486180	...	0.261057	0.643078	0.005274

284807 rows × 31 columns

```
#NOW LET US TRY TO UNDERSTAND THE DATA TO GAIN USEFUL INSIGHTS FROM DATASET(EDA))
```

```
df.head(10)#gives first 10 rows from the data
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.018307	0.277838	-0.110474
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321
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
5	2.0	-0.425966	0.960523	1.141109	-0.168252	0.420987	-0.029728	0.476201	0.260314	-0.568671	...	-0.208254	-0.559825	-0.026398
6	4.0	1.229658	0.141004	0.045371	1.202613	0.191881	0.272708	-0.005159	0.081213	0.464960	...	-0.167716	-0.270710	-0.154104
7	7.0	-0.644269	1.417964	1.074380	-0.492199	0.948934	0.428118	1.120631	-3.807864	0.615375	...	1.943465	-1.015455	0.057504
8	7.0	-0.894286	0.286157	-0.113192	-0.271526	2.669599	3.721818	0.370145	0.851084	-0.392048	...	-0.073425	-0.268092	-0.204233
9	9.0	-0.338262	1.119593	1.044367	-0.222187	0.499361	-0.246761	0.651583	0.069539	-0.736727	...	-0.246914	-0.633753	-0.120794

10 rows × 31 columns

```
df.shape#this gives the idea about our data shape that is in the form of matrix of 3973 rows and 31 columns
```

```
(284807, 31)
```

```
df.describe()#the describe method is used to know about the statistics of the data
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8
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
mean	94813.859575	1.168375e-15	3.416908e-16	-1.379537e-15	2.074095e-15	9.604066e-16	1.487313e-15	-5.556467e-16	1.213481e-16
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+00	1.237094e+00	1.194353e+00
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-02
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-01	3.985649e-01	5.704361e-01	3.273459e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01

8 rows × 31 columns

```
df.info()#from this we can conclude that the datatype of every v's are float
```

```
<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
```

```
df.size #this gives us the complete count of the data
```

8829017

```
#now let us find if there are any null values in our dataset
```

```
df.isnull().sum()#from this we can observe that there are some null values
```

```
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

```

df.nunique()#from this we can conclude that our class label/column has only 2 unique values

```

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

```

```

print(df.size)
print(df.shape)

```

```

8829017
(284807, 31)

```

#now let us know about our complete data that is how many valid transactions are being happened and how many of them are fraud

```

fraud = df[df['Class'] == 1]
not_valid = df[df['Class'] == 0]
outlierFraction = len(fraud)/float(len(not_fraud))
print(outlierFraction)
print('Fraud Cases: {}'.format(len(df[df['Class'] == 1])))
print('Valid Transactions: {}'.format(len(df[df['Class'] == 0])))

```

```

0.1239294710327456
Fraud Cases: 492
Valid Transactions: 284315

```

```
#let us now know about the amount details of the fraud and not_fraud transactions
```

```
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%       105.890000
max      2125.870000
Name: Amount, dtype: float64
```

```
print('Amount details of the non_fraudulent transaction')
not_fraud.Amount.describe()
```

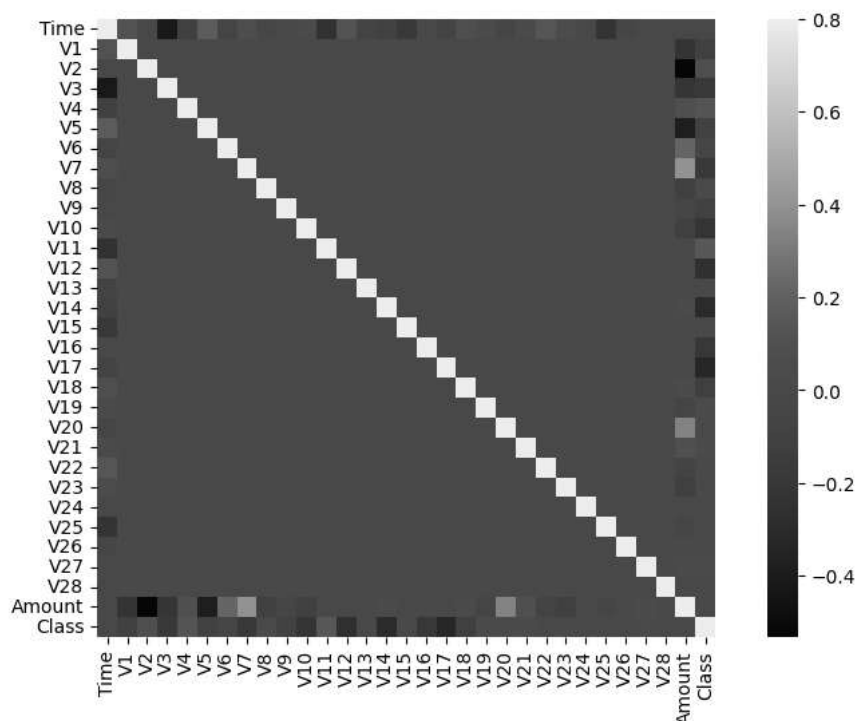
```
Amount details of the non_fraudulent transaction
count   3970.000000
mean      64.899597
std      213.612570
min        0.000000
25%        2.270000
50%       12.990000
75%       54.990000
max      7712.430000
Name: Amount, dtype: float64
```

```
#from the above two descriptions we can conclude that the mean amount of fraud transactions is higher than the non_fraud transactions
```

```
#DATA VISUALISATION
```

```
#NOW LET US FIND THE CORRELATION IN THE DATA SET USING VISUALISATION
```

```
corr = df.corr()
fig = plt.figure(figsize = (10, 6))
sns.heatmap(corr, vmax = .8, square = True)
plt.show()
```



```
#FROM THE ABOVE CORRELATION MATRIX WE CAN CONCLUDE THAT THERE IS NO MUCH CORRELATION AMONG THE FEATURES
```

```
#NOW LET US TRAIN OUR MODEL TO DETECT WHETHER GIVEN IS FRAUDALENT OR NOT
```

```
# dividing the X and Y from the dataset
x = df.iloc[:,0:30]
x
```

	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
...
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006

284807 rows × 30 columns

```
y=df.iloc[:,30]
y
```

0	0
1	0
2	0
3	0
4	0
...	..
284802	0
284803	0
284804	0
284805	0
284806	0

Name: Class, Length: 284807, dtype: int64

```
#let us now perform the logistic regression
```

```
#train and test variables
```

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

```
x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=0)
```

```
#we can perform the normalization technique here if needed but our data is not lengthy so we skip this step for now
```

```
#NOW APPLY SUITABLE CLASSIFIER/REGRESSOR/CLUSTERER
```

```
from sklearn.linear_model import LogisticRegression
model=LogisticRegression()
```

```
#NOW FITTING THE MODEL---->THAT IS GIVING DATA TO LOGISTIC REGRESSION MODEL
```

```
model.fit(x_train,y_train)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: Converger
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = check_optimize_result
```

```
#now let us predict the output
```

```
LogisticRegression()
```

```
y_pred=model.predict(x_test)#through this we can get the predicted values
```

```
y_pred
```

```
array([0, 0, 0, ..., 0, 0, 0])
```

```
#NOW LET US COMPARE THEM WITH ACTUAL VALUES
```

```
y_test
```

```
183484    0
255448    0
244749    0
63919     0
11475     0
..
52247     0
247905    0
78338     0
246056    0
40618     0
Name: Class, Length: 71202, dtype: int64
```

```
# ACCURACY
```

```
from sklearn.metrics import accuracy_score
accuracy_score(y_pred,y_test)*100
```

```
99.89887924496503
```

```
#above model gives the accuracy of 99.89 (accuracy) using logistic regression
```

```
from sklearn.metrics import classification_report, accuracy_score
from sklearn.metrics import precision_score, recall_score
from sklearn.metrics import f1_score, matthews_corrcoef
from sklearn.metrics import confusion_matrix
```

```
n_outliers = len(fraud)
n_errors = (y_pred != y_test).sum()
print("The model used is logistic regression")
```

```
acc = accuracy_score(y_test, y_pred)
print("The accuracy is {}".format(acc))
```

```
prec = precision_score(y_test, y_pred)
print("The precision is {}".format(prec))
```

```
rec = recall_score(y_test, y_pred)
print("The recall is {}".format(rec))
```

```
f1 = f1_score(y_test, y_pred)
print("The F1-Score is {}".format(f1))
```

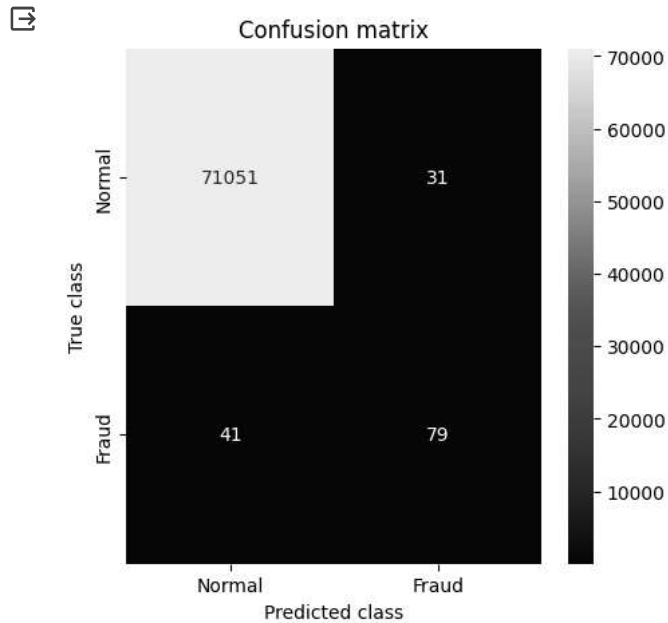
```
MCC = matthews_corrcoef(y_test, y_pred)
print("The Matthews correlation coefficient is{}".format(MCC))
```

```
The model used is logistic regression
The accuracy is 0.9989887924496503
The precision is 0.7181818181818181
The recall is 0.6583333333333333
The F1-Score is 0.6869565217391305
The Matthews correlation coefficient is0.68710290267352
```

```
#THE ABOVE LISTED ARE THE VARIOUS METRICS FOR THE FRAUD DETECTION USING LOGISTIC REGRESSION
```

```
#LET US VISUALIZE THE CONFUSION MATRIX
```

```
LABELS = ['Normal', 'Fraud']  
conf_matrix = confusion_matrix(y_test, y_pred)  
plt.figure(figsize =(5, 5))  
sns.heatmap(conf_matrix, xticklabels = LABELS,  
            yticklabels = LABELS, annot = True, fmt ="d");  
plt.title("Confusion matrix")  
plt.ylabel('True class')  
plt.xlabel('Predicted class')  
plt.show()
```



```
#This is the confusion matrix of the normal and fraud with predicted and true values/classes
```