

Author: Mahendra Kumar Indian Institute of Technology Jammu

For Model Training - Logistic Regression

For Model Evaluation- accuracy_score compared

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

```
[2] from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
[3] credit_card_data = pd.read_csv('/content/drive/MyDrive/Credit card farAUD/creditcard.csv')
```

```
[4] credit_card_data.head()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28
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	0.066928	0.128539	-0.189115	0.133558	-0.021053
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	-0.339846	0.167170	0.125895	-0.008983	0.014724
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	-0.689281	-0.327642	-0.139097	-0.055353	-0.059752
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	-1.175575	0.647376	-0.221929	0.062723	0.061458
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.009431	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153

5 rows × 31 columns

```
<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
```

```
# checking the number of missing values in each column  
credit_card_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
```

✓
0s

```
# distribution of legit transactions & fraudulent transactions  
credit_card_data['Class'].value_counts()
```

```
0    284315  
1      492  
Name: Class, dtype: int64
```

✓
0s

```
[8] # separating the data for analysis  
legit = credit_card_data[credit_card_data.Class == 0]  
fraud = credit_card_data[credit_card_data.Class == 1]
```

✓
0s

```
[9] # statistical measures of the data  
legit.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
```

✓
0s

```
[10] 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
```

0s



```
# compare the values for both transactions
credit_card_data.groupby('Class').mean()
```



	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V20	V21	V22	V23	V
Class																
0	94838.202258	0.008258	-0.006271	0.012171	-0.007860	0.005453	0.002419	0.009637	-0.000987	0.004467	...	-0.000644	-0.001235	-0.000024	0.000070	0.000...
1	80746.806911	-4.771948	3.623778	-7.033281	4.542029	-3.151225	-1.397737	-5.568731	0.570636	-2.581123	...	0.372319	0.713588	0.014049	-0.040308	-0.105...

2 rows × 30 columns



0s

```
[12] legit_sample = legit.sample(n=492)
```

0s

```
[13] new_dataset = pd.concat([legit_sample, fraud], axis=0)
```

0s

```
[14] new_dataset.head()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25
161198	113971.0	1.923030	-0.089917	-2.411952	0.533412	0.649872	-1.002331	0.567675	-0.378595	0.315880	...	0.230511	0.672861	-0.160096	0.607241	0.372230
27515	34581.0	1.224690	0.465365	0.120432	1.118882	-0.033035	-1.035069	0.488549	-0.362718	-0.326955	...	-0.004920	0.054008	-0.123612	0.429692	0.764971
43249	41432.0	-4.348409	3.528210	-2.179063	0.888274	-1.578307	-0.419320	-1.694995	2.973809	-1.403018	...	0.234748	0.159508	-0.042368	-0.278676	1.059277
221996	142780.0	1.905575	-1.382459	-0.554888	-0.915659	-1.154882	-0.207503	-1.022227	-0.002865	-0.193405	...	0.510261	1.255992	-0.019834	-0.338258	-0.230694
203571	134873.0	1.920715	-0.587574	0.164317	0.304182	-0.906992	0.329663	-1.242625	0.363511	1.139973	...	0.295114	0.855471	0.297016	0.707362	-0.600389

5 rows × 31 columns

✓ [15] new_dataset['Class'].value_counts()

0s

```
0    492
1    492
Name: Class, dtype: int64
```

✓ [16] new_dataset.groupby('Class').mean()

0s

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V20	V21	V22	V23	V24
Class																
0	95764.756098	0.059783	0.037790	-0.183396	0.090261	0.042669	-0.014866	0.027514	0.085745	-0.016563	...	-0.011256	0.010688	-0.042701	0.001503	-0.047461
1	80746.806911	-4.771948	3.623778	-7.033281	4.542029	-3.151225	-1.397737	-5.568731	0.570636	-2.581123	...	0.372319	0.713588	0.014049	-0.040308	-0.105130

2 rows × 30 columns



✓ [17] X = new_dataset.drop(columns='Class', axis=1)

0s

```
Y = new_dataset['Class']
```

✓ [18] X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, stratify=Y, random_state=2)

0s

✓ [19] print(X.shape, X_train.shape, X_test.shape)

0s

```
(984, 30) (787, 30) (197, 30)
```

Model Training

```
✓ [20] model = LogisticRegression()  
0s
```

```
✓ [22] # training the Logistic Regression Model with Training Data  
0s model.fit(X_train, Y_train)  
  
LogisticRegression()
```

Model Evaluation

```
✓ [23] # accuracy on training data  
0s X_train_prediction = model.predict(X_train)  
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)
```

```
✓ [24] print('Accuracy on Training data : ', training_data_accuracy)  
0s
```

Accuracy on Training data : 0.9161372299872935

```
✓ [27] # accuracy on test data  
0s X_test_prediction = model.predict(X_test)  
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)
```

```
✓ [28] print('Accuracy score on Test Data : ', test_data_accuracy)  
0s
```

Accuracy score on Test Data : 0.8934010152284264

Double-click (or enter) to edit