


## Importing the Dependencies

```
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

# loading the dataset to a Pandas DataFrame
credit_card_data = pd.read_csv('/content/creditcard.csv')


# first 5 rows of the dataset
credit_card_data.head()
```



	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	0.277861
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638651
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771651
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005271
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.009431	0.798271

5 rows × 31 columns


```
credit_card_data.tail()
```



	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22
284802	172786.0	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.914428	...	0.213451	0.213451
284803	172787.0	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.584800	...	0.214201	0.214201
284804	172788.0	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.432454	...	0.232041	0.232041
284805	172788.0	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.392087	...	0.265241	0.265241
284806	172792.0	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.486180	...	0.261051	0.261051

5 rows × 31 columns

```
# dataset informations
credit_card_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

```

```

# checking the number of missing values in each column
credit_card_data.isnull().sum()

```



	count
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

```

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

```



count	
Class	
0	284315
1	492



This Dataset is highly unblanced

0 -> Normal Transaction

1 -> fraudulent transaction

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

```
print(legit.shape)
print(fraud.shape)
```



```
(284315, 31)
(492, 31)
```

```
# statistical measures of the data
legit.Amount.describe()
```



Amount	
count	284315.000000
mean	88.291022
std	250.105092
min	0.000000
25%	5.650000
50%	22.000000
75%	77.050000
max	25691.160000



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



Amount	
count	492.000000
mean	122.211321
std	256.683288
min	0.000000
25%	1.000000
50%	9.250000
75%	105.890000
max	2125.870000



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



	V6	V7	V8	V9	...
12419	0.009637	-0.000987	0.004467	...	-0.000
17737	-5.568731	0.570636	-2.581123	...	0.372



## Under-Sampling

Build a sample dataset containing similar distribution of normal transactions and Fraudulent Transactions

Number of Fraudulent Transactions -> 492

```
legit_sample = legit.sample(n=492)
```

Concatenating two DataFrames

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

```
new_dataset.head()
```



	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21
61552	49886.0	0.982837	-0.368623	1.126019	1.105429	-1.150071	-0.140456	-0.677621	0.275431	0.587823	...	0.197364
61650	49924.0	-0.718149	-0.182675	0.784151	0.458554	0.440625	-0.540078	0.546506	-0.633099	0.208116	...	0.032462
224193	143688.0	-1.665787	-0.503365	2.137617	-0.818380	0.139817	-0.877459	0.483466	0.142268	0.008394	...	0.464638
140658	83852.0	-1.733278	1.594123	0.544248	0.616726	-0.004772	-0.066265	0.312356	-0.875175	-0.358318	...	1.016659
224375	143765.0	-1.172259	0.964566	0.026250	-2.020997	2.259373	4.466028	-1.004763	-1.210856	0.160314	...	2.128421

5 rows × 31 columns



```
new_dataset.tail()
```



	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21
279863	169142.0	-1.927883	1.125653	-4.518331	1.749293	-1.566487	-2.010494	-0.882850	0.697211	-2.064945	...	0.778584
280143	169347.0	1.378559	1.289381	-5.004247	1.411850	0.442581	-1.326536	-1.413170	0.248525	-1.127396	...	0.370612
280149	169351.0	-0.676143	1.126366	-2.213700	0.468308	-1.120541	-0.003346	-2.234739	1.210158	-0.652250	...	0.751826
281144	169966.0	-3.113832	0.585864	-5.399730	1.817092	-0.840618	-2.943548	-2.208002	1.058733	-1.632333	...	0.583276
281674	170348.0	1.991976	0.158476	-2.583441	0.408670	1.151147	-0.096695	0.223050	-0.068384	0.577829	...	-0.164350

5 rows × 31 columns



```
new_dataset['Class'].value_counts()
```



	count
Class	
0	492
1	492



```
new_dataset.groupby('Class').mean()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...
Class											
0	99687.058943	0.028644	-0.045205	-0.084760	-0.034788	-0.014881	-0.026232	-0.059963	0.028633	0.014886	... -0.070
1	80746.806911	-4.771948	3.623778	-7.033281	4.542029	-3.151225	-1.397737	-5.568731	0.570636	-2.581123	... 0.372

2 rows x 30 columns

## Splitting the data into Features & Targets

```
X = new_dataset.drop(columns='Class', axis=1)
Y = new_dataset['Class']
```

```
print(X)
```

	Time	V1	V2	V3	V4	V5	V6	\
61552	49886.0	0.982837	-0.368623	1.126019	1.105429	-1.150071	-0.140456	
61650	49924.0	-0.718149	-0.182675	0.784151	0.458554	0.440625	-0.540078	
224193	143688.0	-1.665787	-0.503365	2.137617	-0.818380	0.139817	-0.877459	
140658	83852.0	-1.733278	1.594123	0.544248	0.616726	-0.004772	-0.066265	
224375	143765.0	-1.172259	0.964566	0.026250	-2.020997	2.259373	4.466028	
...	...	...	...	...	...	...	...	
279863	169142.0	-1.927883	1.125653	-4.518331	1.749293	-1.566487	-2.010494	
280143	169347.0	1.378559	1.289381	-5.004247	1.411850	0.442581	-1.326536	
280149	169351.0	-0.676143	1.126366	-2.213700	0.468308	-1.120541	-0.003346	
281144	169966.0	-3.113832	0.585864	-5.399730	1.817092	-0.840618	-2.943548	
281674	170348.0	1.991976	0.158476	-2.583441	0.408670	1.151147	-0.096695	
	V7	V8	V9	...	V20	V21	V22	\
61552	-0.677621	0.275431	0.587823	...	-0.071244	0.197364	0.303175	
61650	0.546506	-0.633099	0.208116	...	-0.305341	0.032462	0.925323	
224193	0.483466	0.142268	0.008394	...	0.365694	0.464638	0.841699	
140658	0.312356	-0.875175	-0.358318	...	-0.409179	1.016659	0.313676	
224375	-1.004763	-1.210856	0.160314	...	-0.511513	2.128421	-1.350318	
...	...	...	...	...	...	...	...	
279863	-0.882850	0.697211	-2.064945	...	1.252967	0.778584	-0.319189	
280143	-1.413170	0.248525	-1.127396	...	0.226138	0.370612	0.028234	
280149	-2.234739	1.210158	-0.652250	...	0.247968	0.751826	0.834108	
281144	-2.208002	1.058733	-1.632333	...	0.306271	0.583276	-0.269209	
281674	0.223050	-0.068384	0.577829	...	-0.017652	-0.164350	-0.295135	
	V23	V24	V25	V26	V27	V28	Amount	
61552	-0.073249	0.282617	0.222791	-0.382392	0.033521	0.040084	85.00	
61650	0.580951	0.068247	-1.318343	0.822810	-0.474397	-0.249802	8.02	
224193	0.014607	0.527027	0.795392	-0.581914	-0.014156	0.085019	150.00	
140658	-0.063209	0.010602	-0.294514	-0.441019	-1.063089	0.008448	28.33	
224375	0.328534	0.653189	-0.374396	0.321437	0.195329	0.147984	1.00	
...	...	...	...	...	...	...	...	
279863	0.639419	-0.294885	0.537503	0.788395	0.292680	0.147968	390.00	
280143	-0.145640	-0.081049	0.521875	0.739467	0.389152	0.186637	0.76	
280149	0.190944	0.032070	-0.739695	0.471111	0.385107	0.194361	77.89	
281144	-0.456108	-0.183659	-0.328168	0.606116	0.884876	-0.253700	245.00	
281674	-0.072173	-0.450261	0.313267	-0.289617	0.002988	-0.015309	42.53	

[984 rows x 30 columns]

```
print(Y)
```

61552	0
61650	0
224193	0
140658	0
224375	0
..	
279863	1
280143	1
280149	1
281144	1
281674	1

Name: Class, Length: 984, dtype: int64

## Split the data into Training data & Testing Data

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, stratify=Y, random_state=2)
```

```
print(X.shape, X_train.shape, X_test.shape)
```

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

## Model Training

### Logistic Regression

```
model = LogisticRegression()
```

```
# training the Logistic Regression Model with Training Data
model.fit(X_train, Y_train)
```

```
➦ /usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to converge
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](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

```
n_iter_i = _check_optimize_result(
```

```
    ▾ LogisticRegression ⓘ ?
```

```
    LogisticRegression())
```

## Model Evaluation

### Accuracy Score

```
# accuracy on training data
```

```
X_train_prediction = model.predict(X_train)
```

```
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)
```

```
print('Accuracy on Training data : ', training_data_accuracy)
```

```
➦ Accuracy on Training data : 0.9415501905972046
```

```
# accuracy on test data
```

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

```
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)
```

```
print('Accuracy score on Test Data : ', test_data_accuracy)
```

```
➦ Accuracy score on Test Data : 0.9035532994923858
```