

# Data Information

The dataset contains transactions made by credit cards in September 2013 by European cardholders. This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions.

It contains only numerical input variables which are the result of a PCA transformation. Unfortunately, due to confidentiality issues, we cannot provide the original features, and only background information about the data. Features V1, V2, ..., V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. The feature 'Amount' is the transaction Amount, this feature can be used for example-dependant cost-sensitive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.

## Import necessary libraries

```
In [83]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import RobustScaler, StandardScaler
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, f1_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import GradientBoostingClassifier

import warnings
warnings.filterwarnings('ignore')
```

## Load & Explore data

```
In [7]: data=pd.read_csv('D:\Projects\CodeSoft\Data Science Internship\CreditCard Analysis\data.head()'
```

Out[7]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.36378
1	0.0	1.119187	0.266151	0.106480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.25542
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.51465
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.38705
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.81772

5 rows × 11 columns

```
In [8]: # To show basic info about datatype
shape = data.shape
print(f'\n Number of Rows = {shape[0]}\n Number of columns = {shape[1]}')

<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
dtype: object
memory usage: 67.4 MB

Number of Rows = 284807
Number of Columns = 31
```

```
In [9]: # To display stats about the data
data.describe()
```

```
df = data.copy()
temp = data.drop(columns=['Time', 'Amount', 'Class'], axis=1)

fig, ax = plt.subplots(figsize=(20, 40), ncols=4, nrows=7)
index = 0
ax = ax.flatten()

for col in temp.columns:
    sns.distplot(temp[col], ax=ax[index])
    index += 1
```

```
In [10]: # Check for null values
data.isnull().sum()
```

Out[10]:

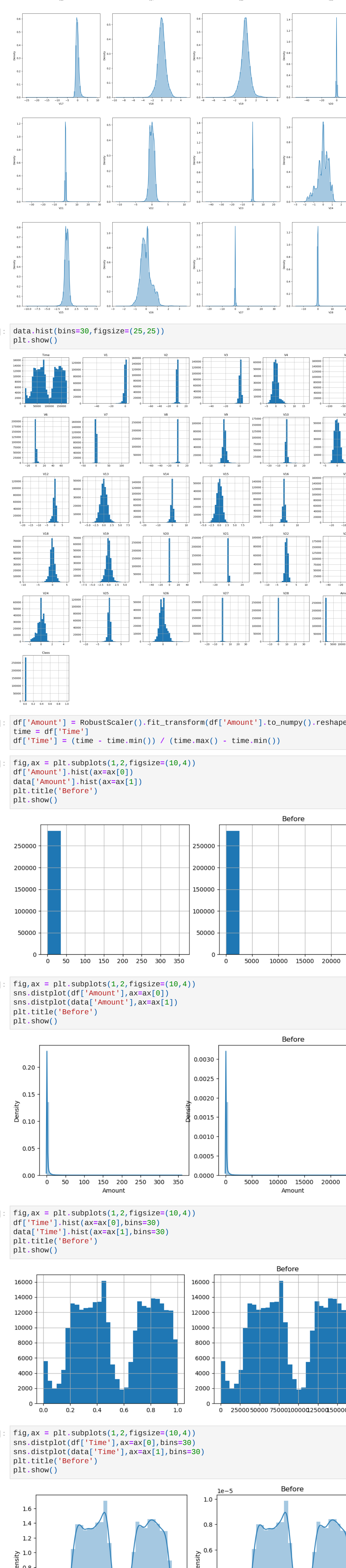
Time	0
V1	0
V2	0
V3	0
V4	0
V5	0
V6	0
V7	0
V8	0
V9	0
V10	0
V11	0
V12	0
V13	0
V14	0
V15	0
V16	0
V17	0
V18	0
V19	0
V20	0
V21	0
V22	0
V23	0
V24	0
V25	0
V26	0
V27	0
V28	0
Amount	0
Class	0
dtype:	int64

## Exploratory Analysis

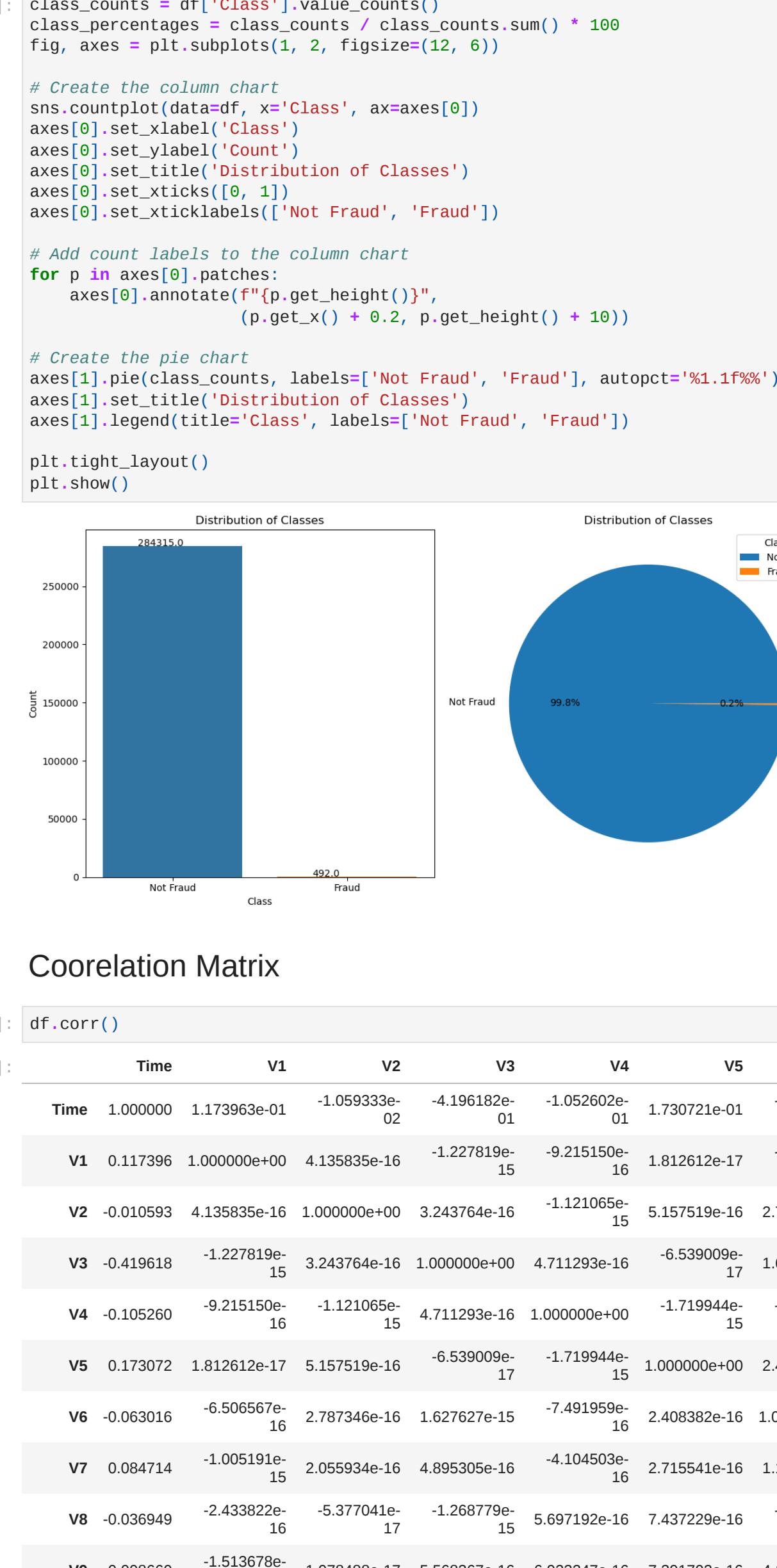
```
In [11]: df=data.copy()
temp=data.drop(columns=['Time', 'Amount', 'Class'],axis=1)
```

```
In [12]: fig,ax=plt.subplots(figsize=(20,40),ncols=4,nrows=7)
index=0
ax=ax.flatten()
for col in temp.columns:
    sns.distplot(temp[col],ax=ax[index])
    index+=1

plt.tight_layout(pad=0.5,w_pad=0.5,h_pad=5)
```

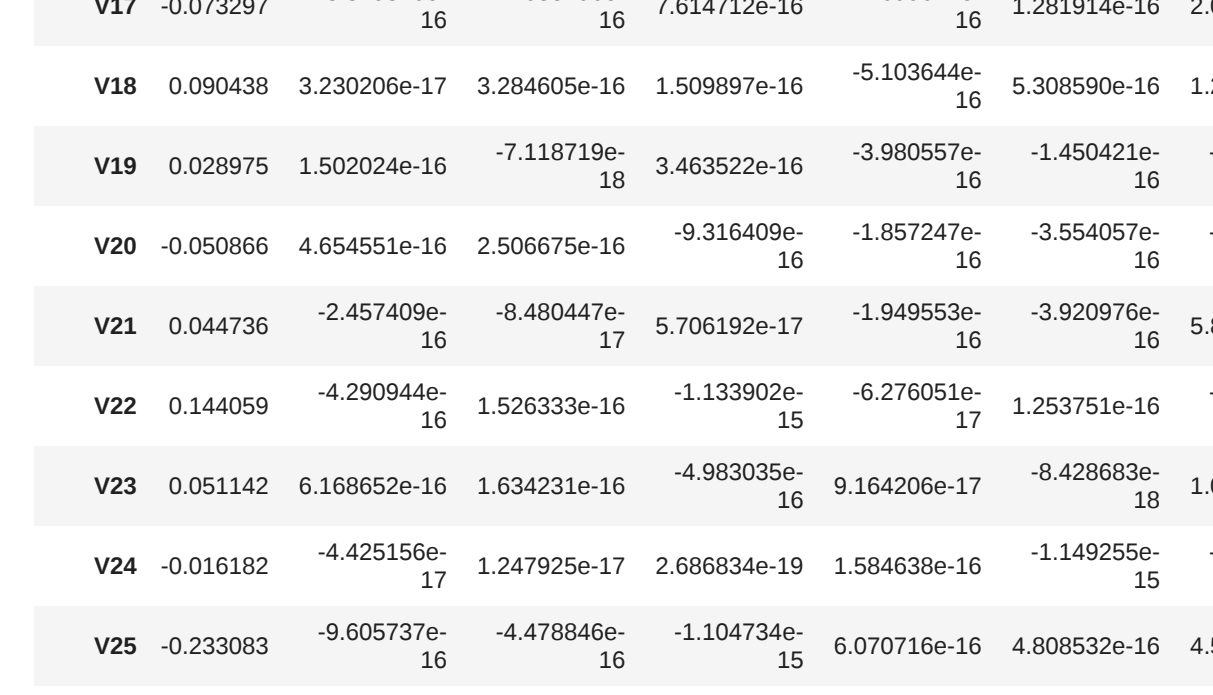


```
In [13]: data.hist(bins=30,figsize=(25,25))
plt.show()
```

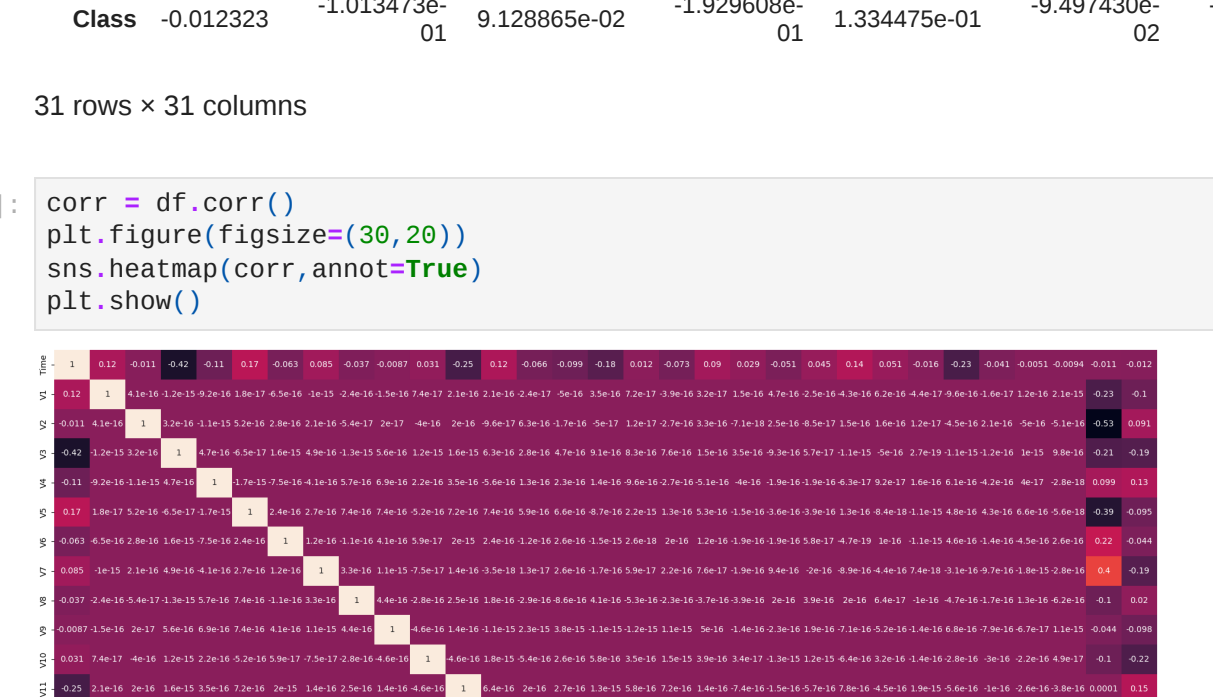


```
In [15]: df['Amount']=RobustScaler().fit_transform(df['Amount']).to_numpy().reshape(-1,1)
df['Time']=df['Time']-(time.time.min())/(time.max()-time.min())
```

```
In [16]: fig,ax=plt.subplots(1,2,figsize=(10,4))
df['Time'].hist(ax=ax[0],bins=30)
data['Amount'].hist(ax=ax[1])
plt.title('Before')
plt.show()
```



```
In [18]: fig,ax=plt.subplots(1,2,figsize=(10,4))
df['Time'].hist(ax=ax[0],bins=30)
data['Amount'].hist(ax=ax[1],bins=30)
plt.title('Before')
plt.show()
```



```
In [25]: fig,ax=plt.subplots(1,2,figsize=(10,4))
sns.distplot(df['Time'],ax=ax[0],bins=30)
sns.distplot(data['Amount'],ax=ax[1],bins=30)
plt.title('Before')
plt.show()
```

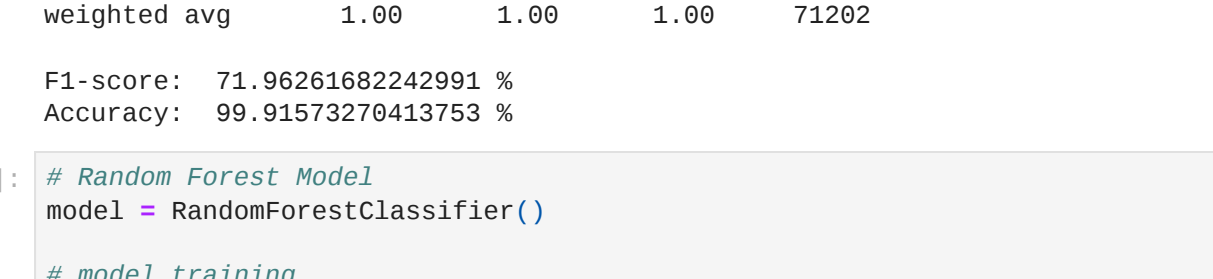


```
In [34]: class_counts = df['Class'].value_counts()
class_percentages = class_counts / class_counts.sum() * 100
fig, axes = plt.subplots(1, 2, figsize=(12, 6))
```

```
# Create the column chart
sns.countplot(data=df, x='Class', ax=axes[0])
axes[0].set_xlabel('Class')
axes[0].set_ylabel('Count')
axes[0].set_title('Distribution of Classes')
axes[0].set_xticks([0, 1])
axes[0].set_xticklabels(['Not Fraud', 'Fraud'])
```

```
# Add count labels to the column chart
for p in axes[0].patches:
    axes[0].annotate(f"({p.get_height()}",
                    (p.get_x() + 0.2, p.get_height() + 10))

# Create the pie chart
axes[1].pie(class_counts, labels=['Not Fraud', 'Fraud'], autopct='%1.1f%%')
axes[1].set_title('Distribution of Classes')
axes[1].legend(title='Class', labels=['Not Fraud', 'Fraud'])
```



## Coorelation Matrix

```
In [62]: df.corr()
```

```

accuracy      1.00      71202
macro avg     0.98      0.89      0.93      71202
weighted avg  1.00      1.00      1.00      71202

F1-score:  85.71428571428571 %
Accuracy:  99.95085744220668 %

# Gradient boost Model
model = GradientBoostingClassifier()

# model training
model.fit(x_train, y_train)

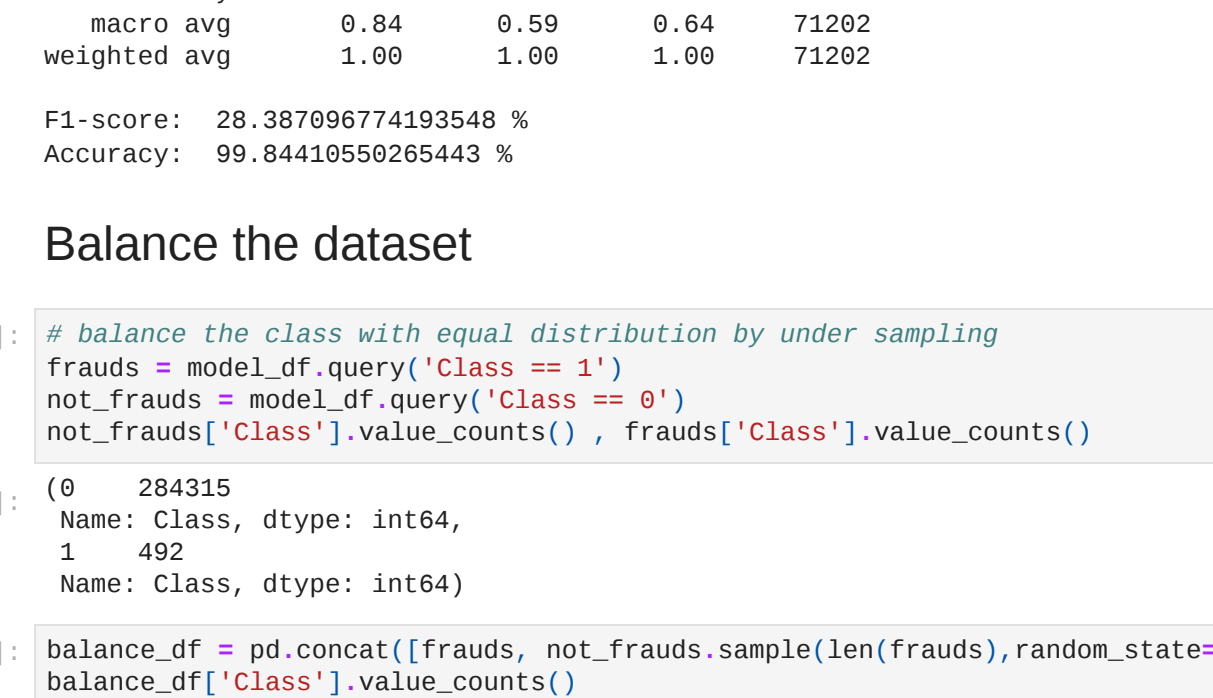
# model test
y_pred = model.predict(x_test)

# print classification report get
print(classification_report(y_test, y_pred, target_names=['Not Fraud', 'Fraud']))

# print matrix to get F1 score

```

```
In [35]: corr = df.corr()
plt.figure(figsize=(30,20))
sns.heatmap(corr,annot=True)
```



## Model Training

```
In [72]: model_df = df.copy()
# Input split
x = model_df.drop(columns=['Class'], axis=1)
y = model_df['Class']
```

```
In [38]: # Standard scaling
sc = StandardScaler()
x_scaler = sc.fit_transform(x)
```

```
In [51]: # Train test split
x_train, x_test, y_train, y_test = train_test_split(x_scaler, y, test_size=0.25,
```

```
In [52]: # Logistic Regression Model
model = LogisticRegression()

# model training
model.fit(x_train, y_train)

# model test
y_pred = model.predict(x_test)

# print classification report get
print(classification_report(y_test, y_pred, target_names=['Not Fraud', 'Fraud']))

# print matrix to get F1-score
print('F1-score: ', f1_score(y_test, y_pred)*100,'%')

# print matrix to get performance
print('Accuracy: ', model.score(x_test, y_test)*100,'%')
```

	precision	recall	f1-score	support
Not Fraud	1.00	1.00	1.00	71879
Fraud	0.85	0.63	0.72	123
accuracy			1.00	71282
macro avg	0.92	0.81	0.86	71282
weighted avg	1.00	1.00	1.00	71282

F1-score: 71.96261682242991 %  
Accuracy: 99.91573278413753 %

```
In [63]: # Random Forest Model
model = RandomForestClassifier()

# model training
model.fit(x_train, y_train)

# model test
y_pred = model.predict(x_test)

# print classification report get
print(classification_report(y_test, y_pred, target_names=['Not Fraud', 'Fraud']))

# print matrix to get F1-score
print('F1-score: ', f1_score(y_test, y_pred)*100,'%')

# print matrix to get performance
print('Accuracy: ', model.score(x_test, y_test)*100,'%')
```

	precision	recall	f1-score	support
Not Fraud	1.00	1.00	1.00	71079
Fraud	0.95	0.78	0.86	123
accuracy			1.00	71282
macro avg	0.98	0.89	0.94	71282
weighted avg	1.00	1.00	1.00	71282

F1-score: 85.71428571428571 %  
Accuracy: 99.95585744228668 %

```
In [64]: # Gradient boost Model
model = GradientBoostingClassifier()

# model training
model.fit(x_train, y_train)

# model test
y_pred = model.predict(x_test)

# print classification report get
print(classification_report(y_test, y_pred, target_names=['Not Fraud', 'Fraud']))

# print matrix to get F1-score
print('F1-score: ', f1_score(y_test, y_pred)*100,'%')

# print matrix to get performance
print('Accuracy: ', model.score(x_test, y_test)*100,'%')
```

	precision	recall	f1-score	support
Not Fraud	1.00	1.00	1.00	71079
Fraud	0.69	0.18	0.28	123
accuracy			1.00	71282
macro avg	0.84	0.59	0.64	71282
weighted avg	1.00	1.00	1.00	71282

F1-score: 28.387966774193548 %  
Accuracy: 99.84418556265443 %

## Balance the dataset

```
In [74]: # balance the class with equal distribution by under sampling
frauds = model_df.query('Class == 1')
not_frauds = model_df.query('Class == 0')
not_frauds['class'].value_counts()
frauds['class'].value_counts()
```

Out[74]:

0	284315
1	492
Name:	Class, dtype: int64

```
In [76]: balance_df = pd.concat([frauds, not_frauds.sample(len(frauds), random_state=1)])
balance_df['class'].value_counts()
```

Out[76]:

1	492
0	492
Name:	Class, dtype: int64

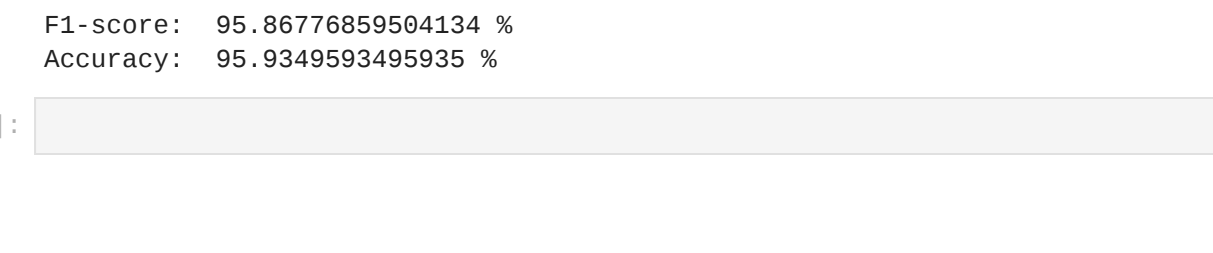
```
In [77]: class_counts = balance_df['class'].value_counts()
class_percentages = class_counts / class_counts.sum() * 100
fig, axes = plt.subplots(1, 2, figsize=(12, 6))

# Create the column chart
sns.countplot(data=balance_df, x='Class', ax=axes[0])
axes[0].set_xlabel('Class')
axes[0].set_ylabel('Count')
axes[0].set_title('Distribution of Classes')
axes[0].set_xticks([0, 1])
axes[0].set_xticklabels(['Not Fraud', 'Fraud'])

# Add count labels to the column chart
for p in axes[0].patches:
    axes[0].annotate(f"({p.get_height()}",
                    (p.get_x() + 0.2, p.get_height() + 10))

# Create the pie chart
axes[1].pie(class_counts, labels=['Not Fraud', 'Fraud'], autopct='%1.1f%%')
axes[1].set_title('Distribution of Classes')
axes[1].legend(title='Class', labels=['Not Fraud', 'Fraud'])

plt.tight_layout()
plt.show()
```



```
In [78]: # retrain the model with balanced dataset
x = balance_df.drop(columns=['Class'], axis=1)
y = balance_df['Class']
```

```
In [79]: # Standard scaling
sc = StandardScaler()
x_scaler = sc.fit_transform(x)
```

```
In [80]: # Train test split
x_train, x_test, y_train, y_test = train_test_split(x_scaler, y, test_size=0.25,
```