Task given by 'Priyasha Rathore' from INSAID
Auther Name = Ramdas Balasaheb Yamgar
Project Name = Fraud Transactions

# Importing Libraries

```
In [1]:
 1  import numpy as np # linear algebra
 2  import pandas as pd # Data processing
 3  import matplotlib.pyplot as plt #Data visualization
 4  import seaborn as sns # Data visualization
 5  from sklearn.preprocessing import LabelEncoder # Label Encoding
 6  from sklearn.preprocessing import MinMaxScaler #Data scaling
 7  from imblearn.over_sampling import SMOTE # Balancing data
 8  from sklearn.metrics import confusion_matrix
 9  from sklearn.metrics import accuracy_score , precision_score , recall_
10  from sklearn.metrics import roc_auc_score,roc_curve
11  import warnings
12  warnings.filterwarnings("ignore")
```

```
In [2]:
 1  df = pd.read_csv('Fraud.csv')
```

```
In [3]:
 1  df.shape
```

Out[3]: (6362620, 11)

```
In [4]:
 1  df.head()
```

Out[4]:

| | step | type | amount | nameOrig | oldbalanceOrg | newbalanceOrig | nameDest |
|---|---|---|---|---|---|---|---|
| 0 | 1 | PAYMENT | 9839.64 | C1231006815 | 170136.0 | 160296.36 | M1979787155 |
| 1 | 1 | PAYMENT | 1864.28 | C1666544295 | 21249.0 | 19384.72 | M2044282225 |
| 2 | 1 | TRANSFER | 181.00 | C1305486145 | 181.0 | 0.00 | C553264065 |
| 3 | 1 | CASH_OUT | 181.00 | C840083671 | 181.0 | 0.00 | C38997010 |
| 4 | 1 | PAYMENT | 11668.14 | C2048537720 | 41554.0 | 29885.86 | M1230701703 |

```
In [5]:
 1  df.columns
```

Out[5]: Index(['step', 'type', 'amount', 'nameOrig', 'oldbalanceOrg', 'newbalance
       Orig',
              'nameDest', 'oldbalanceDest', 'newbalanceDest', 'isFraud',
              'isFlaggedFraud'],
            dtype='object')

In [6]: ▶|
```
1  # Checking Nulls
2  df.isna().mean()*100
```

Out[6]:
```
step              0.0
type              0.0
amount            0.0
nameOrig          0.0
oldbalanceOrg     0.0
newbalanceOrig    0.0
nameDest          0.0
oldbalanceDest    0.0
newbalanceDest    0.0
isFraud           0.0
isFlaggedFraud    0.0
dtype: float64
```

In [7]: ▶|
```
1  df.dtypes
```

Out[7]:
```
step                int64
type               object
amount            float64
nameOrig           object
oldbalanceOrg     float64
newbalanceOrig    float64
nameDest           object
oldbalanceDest    float64
newbalanceDest    float64
isFraud             int64
isFlaggedFraud      int64
dtype: object
```

In [8]: ▶|
```
1  pd.set_option('display.float_format', '{:.2f}'.format) # To see actual
2  df.describe()
```

Out[8]:

|       | step       | amount      | oldbalanceOrg | newbalanceOrig | oldbalanceDest | newbalanc |
|-------|------------|-------------|---------------|----------------|----------------|-----------|
| count | 6362620.00 | 6362620.00  | 6362620.00    | 6362620.00     | 6362620.00     | 63626     |
| mean  | 243.40     | 179861.90   | 833883.10     | 855113.67      | 1100701.67     | 12249     |
| std   | 142.33     | 603858.23   | 2888242.67    | 2924048.50     | 3399180.11     | 36741     |
| min   | 1.00       | 0.00        | 0.00          | 0.00           | 0.00           |           |
| 25%   | 156.00     | 13389.57    | 0.00          | 0.00           | 0.00           |           |
| 50%   | 239.00     | 74871.94    | 14208.00      | 0.00           | 132705.66      | 2146      |
| 75%   | 335.00     | 208721.48   | 107315.18     | 144258.41      | 943036.71      | 11119     |
| max   | 743.00     | 92445516.64 | 59585040.37   | 49585040.37    | 356015889.35   | 3561792   |

In [9]: ▶|    `1  df.describe(include= 'object')`

Out[9]:

|         | type     | nameOrig    | nameDest    |
|---------|----------|-------------|-------------|
| count   | 6362620  | 6362620     | 6362620     |
| unique  | 5        | 6353307     | 2722362     |
| top     | CASH_OUT | C1902386530 | C1286084959 |
| freq    | 2237500  | 3           | 113         |

In [10]: ▶|    `1  df.duplicated().sum()`

Out[10]: 0

In [11]: ▶|    `1  sns.countplot(df['type'])`

Out[11]: `<AxesSubplot:xlabel='type', ylabel='count'>`
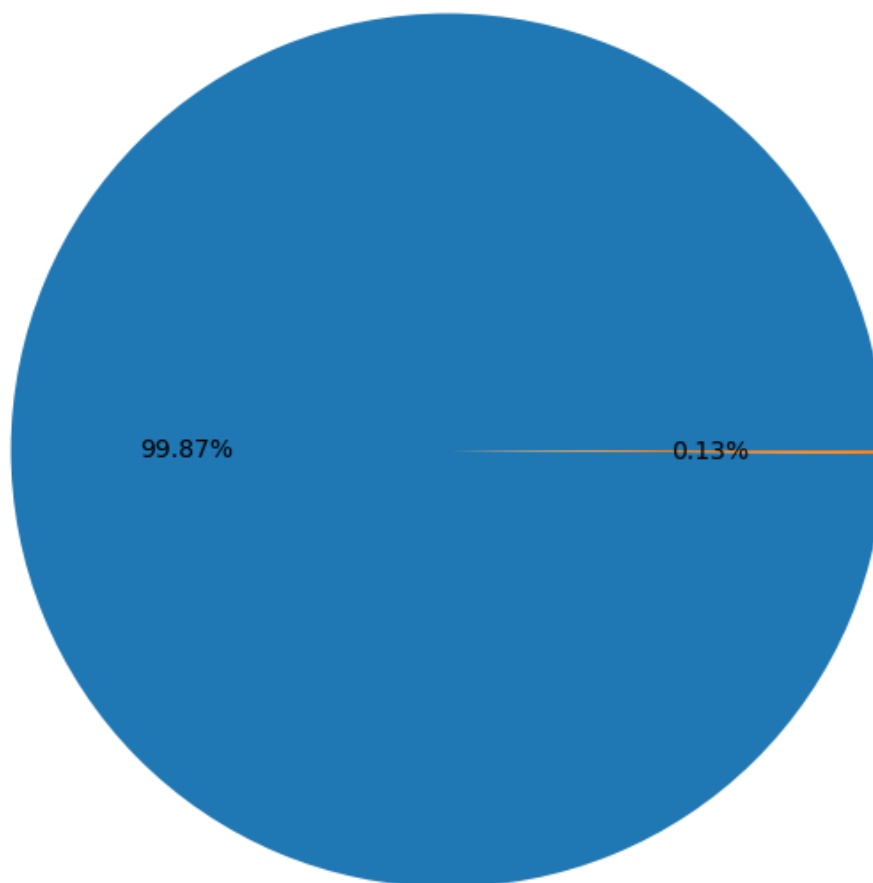


In [12]: ▶|    `1  df['isFraud'].value_counts()`

Out[12]:
```
0    6354407
1       8213
Name: isFraud, dtype: int64
```

```
In [13]:    ▶|   1  plt.figure(figsize=(12,8))
                 2  plt.title('Pie-Chart for Fraud transactions')
                 3  plt.pie(df['isFraud'].value_counts(),  autopct='%.2f%%')
                 4  plt.show()
```

Pie-Chart for Fraud transactions



```
In [14]:    ▶|   1  df.nunique()
```

```
Out[14]:   step               743
           type                 5
           amount         5316900
           nameOrig       6353307
           oldbalanceOrg  1845844
           newbalanceOrig 2682586
           nameDest       2722362
           oldbalanceDest 3614697
           newbalanceDest 3555499
           isFraud              2
           isFlaggedFraud       2
           dtype: int64
```
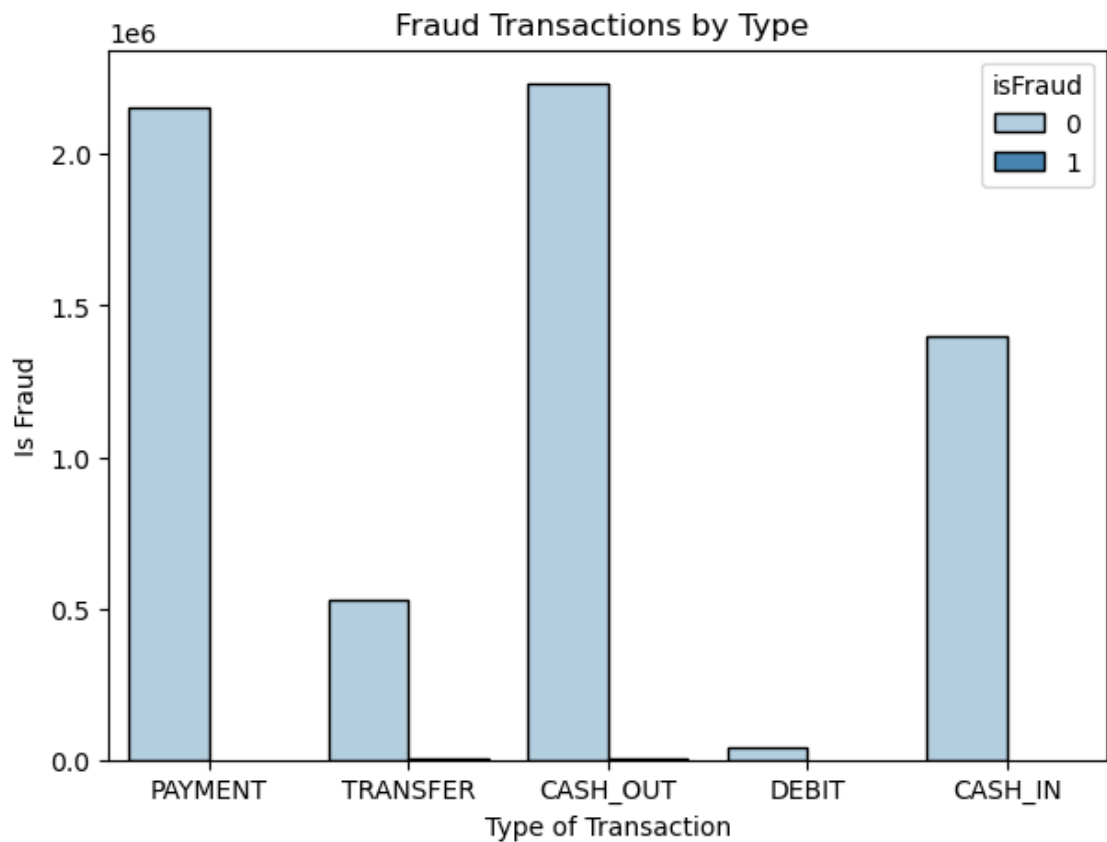
In [15]: ▶| 
```python
1 ((df.loc[df["isFraud"]== 1]).groupby("type").sum())
```

Out[15]:

| type | step | amount | oldbalanceOrg | newbalanceOrig | oldbalanceDest | newb |
|---|---|---|---|---|---|---|
| CASH_OUT | 1513537 | 5989202243.83 | 5984124999.99 | 298767.61 | 4465524469.93 | 1049 |
| TRANSFER | 1512246 | 6067213184.01 | 7564595045.72 | 1579821917.66 | 4397651.53 | 1 |

In [16]: ▶|
```python
1 plt.figure(figsize=(7,5))
2 sns.countplot(x='type',hue='isFraud',palette='Blues',data=df,edgecolor
3 plt.title("Fraud Transactions by Type ")
4 plt.xlabel("Type of Transaction")
5 plt.ylabel("Is Fraud")
```

Out[16]:  Text(0, 0.5, 'Is Fraud')



In [17]: ▶|
```python
1 ((df.loc[df["isFlaggedFraud"]== 1]).groupby("type").sum())
```

Out[17]:

| type | step | amount | oldbalanceOrg | newbalanceOrig | oldbalanceDest | newbalance |
|---|---|---|---|---|---|---|
| TRANSFER | 8601 | 77785563.69 | 125085904.19 | 125085904.19 | 0.00 | |

```
In [18]:    1  plt.figure(figsize=(7,5))
            2  sns.countplot(x='type',hue='isFlaggedFraud',palette='Blues',data=df,ec
            3  plt.title("Flagged Fraud Transactions by Type")
            4  plt.xlabel("Type of Transaction")
            5  plt.ylabel("is Flagged Fraud")
```
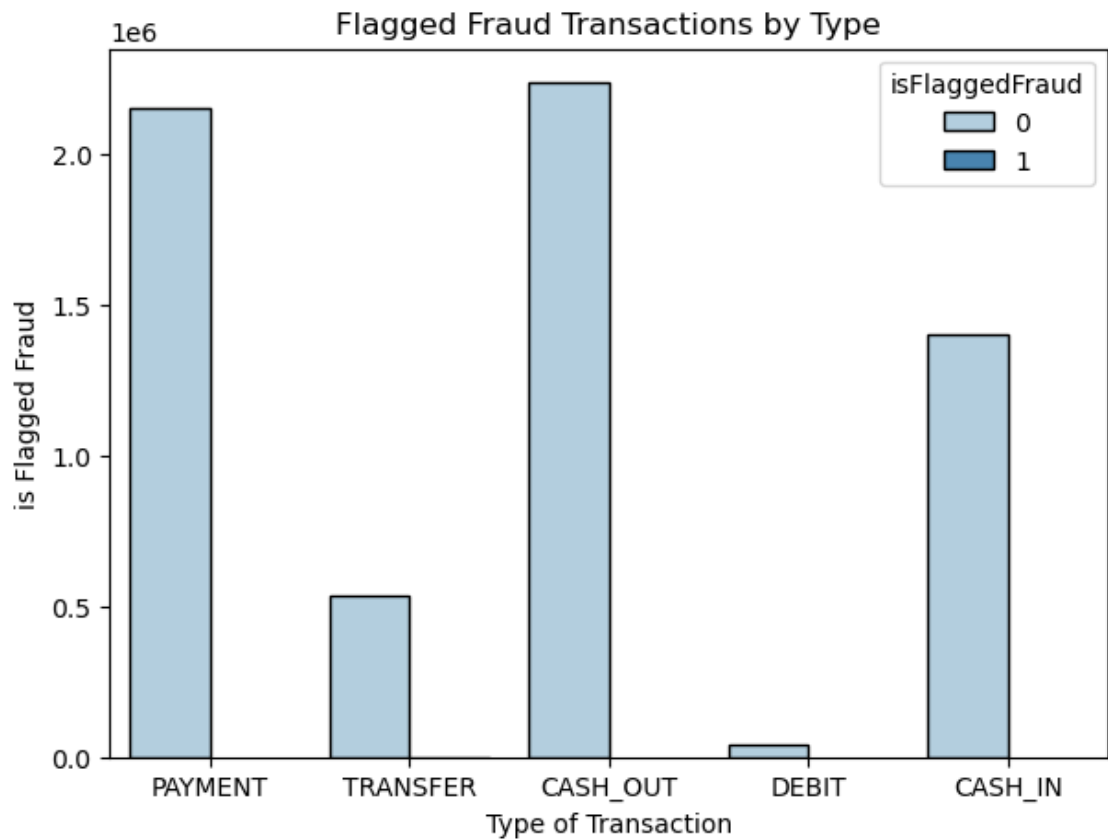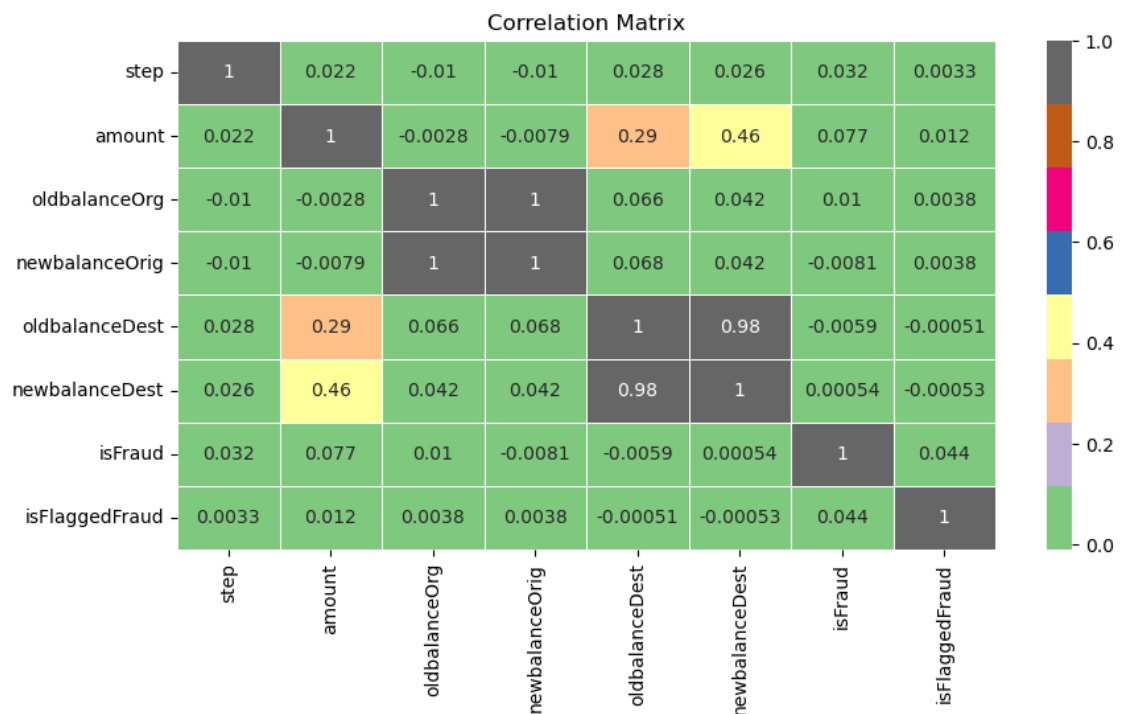
Out[18]:  Text(0, 0.5, 'is Flagged Fraud')



from above we conlcude that

1) Only in CASH_OUT and TRANSFER type "Fraud" is happen
2) Only 16 times "Flagged Fraud" happen which is also in TRANSFER type.

In [19]: ▶|
```python
1  #check for correlation
2  ig ,ax = plt.subplots(figsize=(10,5))
3  sns.heatmap(data=df.corr(),cmap="Accent",annot=True,linewidths=.5,ax=a
4  plt.show()
```

Correlation Matrix

|                  | step | amount | oldbalanceOrg | newbalanceOrig | oldbalanceDest | newbalanceDest | isFraud | isFlaggedFraud |
|------------------|------|--------|---------------|----------------|----------------|----------------|---------|----------------|
| step             | 1    | 0.022  | -0.01         | -0.01          | 0.028          | 0.026          | 0.032   | 0.0033         |
| amount           | 0.022| 1      | -0.0028       | -0.0079        | 0.29           | 0.46           | 0.077   | 0.012          |
| oldbalanceOrg    | -0.01| -0.0028| 1             | 1              | 0.066          | 0.042          | 0.01    | 0.0038         |
| newbalanceOrig   | -0.01| -0.0079| 1             | 1              | 0.068          | 0.042          | -0.0081 | 0.0038         |
| oldbalanceDest   | 0.028| 0.29   | 0.066         | 0.068          | 1              | 0.98           | -0.0059 | -0.00051       |
| newbalanceDest   | 0.026| 0.46   | 0.042         | 0.042          | 0.98           | 1              | 0.00054 | -0.00053       |
| isFraud          | 0.032| 0.077  | 0.01          | -0.0081        | -0.0059        | 0.00054        | 1       | 0.044          |
| isFlaggedFraud   | 0.0033| 0.012 | 0.0038        | 0.0038         | -0.00051       | -0.00053       | 0.044   | 1              |

So multicollinearity exists between "oldbalanceOrg" and "newbalanceOrg" ; "oldbalanceDest" and "newbalanceDest".

Let us try to build the model without removing features.

# Feature engineering

In [20]: ▶|
```
1 print("Shape : ",(df.loc[(df["isFraud"] == 1 ) & (df["isFlaggedFraud"]
2 df.loc[(df["isFraud"] == 1 ) & (df["isFlaggedFraud"] == 1 )]
```

Shape :  (16, 11)

Out[20]:

| | step | type | amount | nameOrig | oldbalanceOrg | newbalanceOrig | |
|---|---|---|---|---|---|---|---|
| 2736446 | 212 | TRANSFER | 4953893.08 | C728984460 | 4953893.08 | 4953893.08 | C6 |
| 3247297 | 250 | TRANSFER | 1343002.08 | C1100582606 | 1343002.08 | 1343002.08 | C11 |
| 3760288 | 279 | TRANSFER | 536624.41 | C1035541766 | 536624.41 | 536624.41 | C11 |
| 5563713 | 387 | TRANSFER | 4892193.09 | C908544136 | 4892193.09 | 4892193.09 | C8 |
| 5996407 | 425 | TRANSFER | 10000000.00 | C689608084 | 19585040.37 | 19585040.37 | C13 |
| 5996409 | 425 | TRANSFER | 9585040.37 | C452586515 | 19585040.37 | 19585040.37 | C11 |
| 6168499 | 554 | TRANSFER | 3576297.10 | C193696150 | 3576297.10 | 3576297.10 | C4 |
| 6205439 | 586 | TRANSFER | 353874.22 | C1684585475 | 353874.22 | 353874.22 | C17 |
| 6266413 | 617 | TRANSFER | 2542664.27 | C786455622 | 2542664.27 | 2542664.27 | C6 |
| 6281482 | 646 | TRANSFER | 10000000.00 | C19004745 | 10399045.08 | 10399045.08 | C18 |
| 6281484 | 646 | TRANSFER | 399045.08 | C724693370 | 10399045.08 | 10399045.08 | C19 |
| 6296014 | 671 | TRANSFER | 3441041.46 | C917414431 | 3441041.46 | 3441041.46 | C10 |
| 6351225 | 702 | TRANSFER | 3171085.59 | C1892216157 | 3171085.59 | 3171085.59 | C13 |
| 6362460 | 730 | TRANSFER | 10000000.00 | C2140038573 | 17316255.05 | 17316255.05 | C13 |
| 6362462 | 730 | TRANSFER | 7316255.05 | C1869569059 | 17316255.05 | 17316255.05 | C18 |
| 6362584 | 741 | TRANSFER | 5674547.89 | C992223106 | 5674547.89 | 5674547.89 | C13 |

Conclusion: When "Flagged Fraud" is happen also "Fraud" is happen, So, Don't need of "Flagged Fraud" feature that's why we drop it

# Lebel Encoding

In [21]: ▶|
```
1 df["type"].unique()
```

Out[21]: array(['PAYMENT', 'TRANSFER', 'CASH_OUT', 'DEBIT', 'CASH_IN'],
              dtype=object)

In [22]: ▶|
```
1 # checking Object Columns
2 df.columns[df.dtypes == 'object']
```

Out[22]: Index(['type', 'nameOrig', 'nameDest'], dtype='object')

```python
In [23]:   1  # for Label Encoder
           2
           3  le = {}
           4  for i in df.select_dtypes('object').columns:
           5      le[i] = LabelEncoder()
           6      df[i] = le[i].fit_transform(df[i])
```

```python
In [24]:   1  df.head(2)
```

Out[24]:

| | step | type | amount | nameOrig | oldbalanceOrg | newbalanceOrig | nameDest | oldbalanceDe |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 3 | 9839.64 | 757869 | 170136.00 | 160296.36 | 1662094 | 0.0 |
| **1** | 1 | 3 | 1864.28 | 2188998 | 21249.00 | 19384.72 | 1733924 | 0.0 |

# Balancing Dependent Variable

```python
In [25]:   1  # cheacking dependent variable is balanced or not?
           2  df["isFraud"].value_counts()
```

Out[25]:  0    6354407
          1       8213
          Name: isFraud, dtype: int64

```python
In [26]:   1  X = df.drop(columns=['isFraud'])
           2  y = df['isFraud']
```

```python
In [27]:   1  over_sample = SMOTE()
           2  X,y = over_sample.fit_resample(X,y)
```

```python
In [28]:   1  y.value_counts() #resampled
```

Out[28]:  0    6354407
          1    6354407
          Name: isFraud, dtype: int64

# Train Test Split

```python
In [29]:   1  #import required libraries and split the data into train anad test
           2
           3  from sklearn.model_selection import train_test_split
           4
           5  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
```

```
In [30]:  ▶  1  print(X_train.shape)
              2  print(X_test.shape)
              3  print(y_train.shape)
              4  print(y_test.shape)
```

```
(10167051, 10)
(2541763, 10)
(10167051,)
(2541763,)
```

# Scaling values

```
In [31]:  ▶  1  from sklearn.preprocessing import MinMaxScaler
              2  scaler = MinMaxScaler()
              3
              4  X_train = scaler.fit_transform(X_train)
              5
              6  X_test = scaler.transform(X_test)
```

# Model Building

## Logistic Regression

```
In [32]:  ▶  1  from sklearn.linear_model import LogisticRegression
              2  logreg = LogisticRegression()
              3
              4  logreg.fit(X_train, y_train)
              5
              6  pred_logreg = logreg.predict(X_test)
```
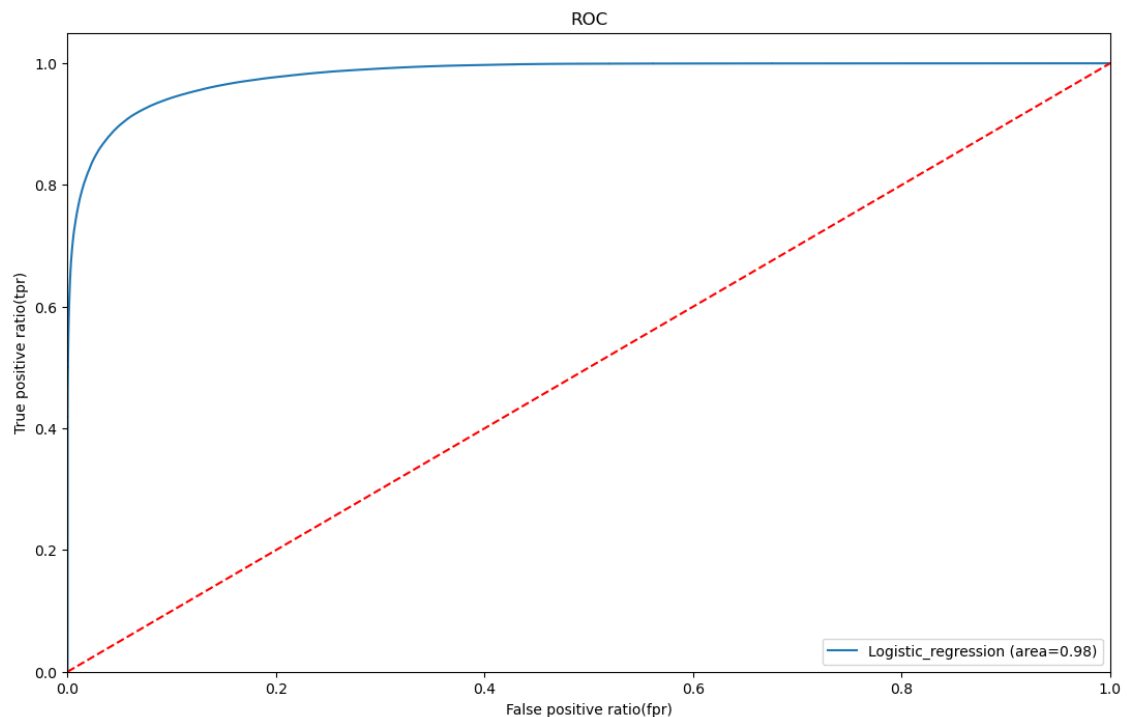
In [33]:

```python
1  conf_logreg = confusion_matrix(y_test , pred_logreg)
2  print("Confusion Matrix :")
3  print(conf_logreg)
4  print(" ")
5  print("Accuracy     :", accuracy_score( y_test, pred_logreg ) * 100)
6  print("Precision    :", precision_score( y_test, pred_logreg ) * 100)
7  print("TPR/ Recall :", recall_score( y_test, pred_logreg ) * 100)
8  print("FPR          :", conf_logreg[0][1] / ( conf_logreg[0][1] * conf_
9  print("F1_Ratio     :", f1_score( y_test, pred_logreg ) * 100)
```

```
Confusion Matrix :
[[1199565   73043]
 [ 116268 1152887]]

Accuracy     : 92.55198065279886
Precision    : 94.04182946824044
TPR/ Recall : 90.83894402180978
FPR          : 8.33635526211585e-05
F1_Ratio     : 92.41264325664257
```

```
In [34]:   ▶  1  #plotting the roc curve and getting the value of auroc
               2
               3  logit_roc_auc=roc_auc_score(y_test,logreg.predict_proba(X_test)[:,1])
               4  fpr,tpr,thresholds=roc_curve(y_test,logreg.predict_proba(X_test)[:,1])
               5  plt.figure(figsize=(13,8))
               6  plt.plot(fpr,tpr,label="Logistic_regression (area=%0.2f)"% logit_roc_a
               7  plt.plot([0,1],[0,1],"r--")
               8  plt.xlim([0.0,1.0])
               9  plt.ylim([0.0,1.05])
              10  plt.xlabel("False positive ratio(fpr)")
              11  plt.ylabel("True positive ratio(tpr)")
              12  plt.title("ROC")
              13  plt.legend(loc="lower right")
              14  plt.savefig("Log_ROC")
              15  plt.show()
```



# Checking our Model is Overfit or Not

```
In [35]:   ▶  1  from sklearn.model_selection import cross_val_score
```

```
In [36]:   ▶  1  cvs = cross_val_score(logreg,X,y,cv=3)
               2  print(cvs)
```

```
[0.94134843 0.95013728 0.90772687]
```

```
In [37]:   ▶  1  cvs.mean()
```

Out[37]:  0.9330708586751232

Cross validation score is less than accuracy. So this is a case of underfitting.

Trying to improve accuracy by removing multi-collinearity and selecting the best features

# Improving accuracy

Select the best 8 features

In [38]: ▶ | 1 df

Out[38]:

|  | step | type | amount | nameOrig | oldbalanceOrg | newbalanceOrig | nameDest | old |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 3 | 9839.64 | 757869 | 170136.00 | 160296.36 | 1662094 | |
| 1 | 1 | 3 | 1864.28 | 2188998 | 21249.00 | 19384.72 | 1733924 | |
| 2 | 1 | 4 | 181.00 | 1002156 | 181.00 | 0.00 | 439685 | |
| 3 | 1 | 1 | 181.00 | 5828262 | 181.00 | 0.00 | 391696 | |
| 4 | 1 | 3 | 11668.14 | 3445981 | 41554.00 | 29885.86 | 828919 | |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 6362615 | 743 | 1 | 339682.13 | 5651847 | 339682.13 | 0.00 | 505863 | |
| 6362616 | 743 | 4 | 6311409.28 | 1737278 | 6311409.28 | 0.00 | 260949 | |
| 6362617 | 743 | 1 | 6311409.28 | 533958 | 6311409.28 | 0.00 | 108224 | |
| 6362618 | 743 | 4 | 850002.52 | 2252932 | 850002.52 | 0.00 | 319713 | |
| 6362619 | 743 | 1 | 850002.52 | 919229 | 850002.52 | 0.00 | 534595 | |

6362620 rows × 11 columns

Now we can build a new logistic regression model using only these 8 features

In [39]: ▶ | 1 df.drop(["nameOrig","nameDest"],inplace=True,axis=1)

In [40]: ▶ | 1 df.tail()

Out[40]:

|  | step | type | amount | oldbalanceOrg | newbalanceOrig | oldbalanceDest | newbalan |
|---|---|---|---|---|---|---|---|
| 6362615 | 743 | 1 | 339682.13 | 339682.13 | 0.00 | 0.00 | 339 |
| 6362616 | 743 | 4 | 6311409.28 | 6311409.28 | 0.00 | 0.00 | |
| 6362617 | 743 | 1 | 6311409.28 | 6311409.28 | 0.00 | 68488.84 | 6379 |
| 6362618 | 743 | 4 | 850002.52 | 850002.52 | 0.00 | 0.00 | |
| 6362619 | 743 | 1 | 850002.52 | 850002.52 | 0.00 | 6510099.11 | 7360 |

Now we can build a new logistic regression model using only these 8 features

```
In [41]:  ▶   1  X_new = df.iloc[:,:-1]
              2  y_new = df.iloc[:,-1]
```

```
In [42]:  ▶   1  X_new_train, X_new_test, y_new_train, y_new_test = train_test_split(X_
```

```
In [43]:  ▶   1  #  Scaling values
              2  X_new_train = scaler.fit_transform(X_new_train)
              3
              4  X_new_test = scaler.transform(X_new_test)
```

```
In [44]:  ▶   1  logreg_new = LogisticRegression()
              2
              3  logreg_new.fit(X_new_train, y_new_train)
              4
              5  pred_new_logreg = logreg_new.predict(X_new_test)
```

```
In [45]:  ▶   1  conf_new_logreg = confusion_matrix(y_new_test , pred_new_logreg)
              2  print("Confusion Matrix :")
              3  print(conf_logreg)
              4  print(" ")
              5  print("Accuracy     :", accuracy_score(y_new_test , pred_new_logreg) *
              6  print("Precision    :", precision_score(y_new_test , pred_new_logreg) *
              7  print("TPR/ Recall :", recall_score(y_new_test , pred_new_logreg) * 10
              8  print("FPR          :", conf_new_logreg[0][1] / ( conf_new_logreg[0][1]
              9  print("F1_Ratio    :", f1_score(y_new_test , pred_new_logreg) * 100)
```

```
Confusion Matrix :
[[1199565    73043]
 [ 116268 1152887]]

Accuracy     : 99.99984283204088
Precision    : 0.0
TPR/ Recall : 0.0
FPR          : nan
F1_Ratio    : 0.0
```

# So final accuracy turns out to be 99.99%, which is a big improvement from previous case

```
In [ ]:  ▶   1
```

```
In [ ]:  ▶   1
```