```
In [1]:  import pandas as pd
         import numpy as np
         import seaborn as sns
         import matplotlib.pyplot as plt
         %matplotlib inline
```

```
In [2]:  train = pd.read_csv('C:/Users/lekshmi/Downloads/train_LE.csv')
         train.head()
```

Out[2]:

|   | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | ApplicantIncome | Co |
|---|---------|--------|---------|------------|-----------|---------------|-----------------|----|
| 0 | LP001002 | Male | No | 0 | Graduate | No | 5849 | |
| 1 | LP001003 | Male | Yes | 1 | Graduate | No | 4583 | |
| 2 | LP001005 | Male | Yes | 0 | Graduate | Yes | 3000 | |
| 3 | LP001006 | Male | Yes | 0 | Not Graduate | No | 2583 | |
| 4 | LP001008 | Male | No | 0 | Graduate | No | 6000 | |

missing values are found

```
In [3]:  train.shape
```

Out[3]:  (614, 13)

```
In [4]:  train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
 #   Column             Non-Null Count  Dtype
---  ------             --------------  -----
 0   Loan_ID            614 non-null    object
 1   Gender             601 non-null    object
 2   Married            611 non-null    object
 3   Dependents         599 non-null    object
 4   Education          614 non-null    object
 5   Self_Employed      582 non-null    object
 6   ApplicantIncome    614 non-null    int64
 7   CoapplicantIncome  614 non-null    float64
 8   LoanAmount         592 non-null    float64
 9   Loan_Amount_Term   600 non-null    float64
 10  Credit_History     564 non-null    float64
 11  Property_Area      614 non-null    object
 12  Loan_Status        614 non-null    object
dtypes: float64(4), int64(1), object(8)
memory usage: 62.5+ KB
```

```
In [5]: train.describe()
```

Out[5]:

|       | ApplicantIncome | CoapplicantIncome | LoanAmount | Loan_Amount_Term | Credit_History |
|-------|-----------------|-------------------|------------|------------------|----------------|
| count | 614.000000      | 614.000000        | 592.000000 | 600.00000        | 564.000000     |
| mean  | 5403.459283     | 1621.245798       | 146.412162 | 342.00000        | 0.842199       |
| std   | 6109.041673     | 2926.248369       | 85.587325  | 65.12041         | 0.364878       |
| min   | 150.000000      | 0.000000          | 9.000000   | 12.00000         | 0.000000       |
| 25%   | 2877.500000     | 0.000000          | 100.000000 | 360.00000        | 1.000000       |
| 50%   | 3812.500000     | 1188.500000       | 128.000000 | 360.00000        | 1.000000       |
| 75%   | 5795.000000     | 2297.250000       | 168.000000 | 360.00000        | 1.000000       |
| max   | 81000.000000    | 41667.000000      | 700.000000 | 480.00000        | 1.000000       |

```
In [6]: pd.crosstab(train['Credit_History'],train['Loan_Status'], margins=True)
```

Out[6]:

| Loan_Status | N | Y | All |
|---|---|---|---|
| Credit_History | | | |
| 0.0 | 82 | 7 | 89 |
| 1.0 | 97 | 378 | 475 |
| All | 179 | 385 | 564 |

higher credit history more eligible

```
In [7]: train.boxplot(column='ApplicantIncome')
```

Out[7]: <Axes: >



ApplicantIncome

outliers are there

```
In [8]: train['ApplicantIncome'].hist(bins=20)
```

Out[8]: <Axes: >

skewed histogram. hence need to normalise

In [9]: `train['CoapplicantIncome'].hist(bins=20)`

Out[9]: <Axes: >

```
In [10]: train.boxplot(column='ApplicantIncome', by = 'Education')
```

```
Out[10]: <Axes: title={'center': 'ApplicantIncome'}, xlabel='Education'>
```



```
In [11]: train.boxplot(column='LoanAmount')
```

```
Out[11]: <Axes: >
```

*Normalisation*

In [12]: 
```
train['LoanAmount_log']=np.log(train['LoanAmount'])
train['LoanAmount_log'].hist(bins=20)
```

Out[12]: <Axes: >



the data looks more normalised now

In [13]: 
```
train.isnull().sum()
```

Out[13]: 
```
Loan_ID               0
Gender               13
Married               3
Dependents           15
Education             0
Self_Employed        32
ApplicantIncome       0
CoapplicantIncome     0
LoanAmount           22
Loan_Amount_Term     14
Credit_History       50
Property_Area         0
Loan_Status           0
LoanAmount_log       22
dtype: int64
```

```
In [14]: train['Gender'].fillna(train['Gender'].mode()[0], inplace=True)
         train['Married'].fillna(train['Married'].mode()[0], inplace=True)
         train['Dependents'].fillna(train['Dependents'].mode()[0], inplace=True)
         train['Self_Employed'].fillna(train['Self_Employed'].mode()[0], inplace=Tru
         train['Credit_History'].fillna(train['Credit_History'].mode()[0], inplace=T
```

```
In [15]: train['Loan_Amount_Term'].fillna(train['Loan_Amount_Term'].mode()[0], inpla
```

```
In [16]: train['LoanAmount']=train['LoanAmount'].fillna(train['LoanAmount'].mean())
```

```
In [17]: train['LoanAmount_log']=train['LoanAmount_log'].fillna(train['LoanAmount_lc
```

```
In [18]: train.isnull().sum()
```

```
Out[18]: Loan_ID              0
         Gender               0
         Married              0
         Dependents           0
         Education            0
         Self_Employed        0
         ApplicantIncome      0
         CoapplicantIncome    0
         LoanAmount           0
         Loan_Amount_Term     0
         Credit_History       0
         Property_Area        0
         Loan_Status          0
         LoanAmount_log       0
         dtype: int64
```

```
In [19]: train['Total_Income']=train['ApplicantIncome']+train['CoapplicantIncome']
         train['Total_Income_log'] = np.log(train['Total_Income'])
```

```
In [20]: train['Total_Income_log'].hist(bins=20)
```

Out[20]: <Axes: >



**division into independent and dependent variables**

```
In [21]: X=train.iloc[:,np.r_[1:5,9:11,13:15]].values
         y=train.iloc[:,12].values
```

```
In [22]: X
```

Out[22]: array([['Male', 'No', '0', ..., 1.0, 4.857444178729352, 5849.0],
               ['Male', 'Yes', '1', ..., 1.0, 4.852030263919617, 6091.0],
               ['Male', 'Yes', '0', ..., 1.0, 4.189654742026425, 3000.0],
               ...,
               ['Male', 'Yes', '1', ..., 1.0, 5.53338948872752, 8312.0],
               ['Male', 'Yes', '2', ..., 1.0, 5.231108616854587, 7583.0],
               ['Female', 'No', '0', ..., 0.0, 4.890349128221754, 4583.0]],
              dtype=object)
```

```
In [23]: y
```

```
Out[23]: array(['Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'Y',
                'N', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'N', 'Y', 'N', 'N', 'N', 'Y',
                'Y', 'Y', 'N', 'Y', 'N', 'N', 'N', 'Y', 'N', 'Y', 'N', 'Y', 'Y',
                'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y',
                'N', 'N', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'N',
                'N', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'N', 'N',
                'N', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
                'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y',
                'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'N',
                'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'N', 'N', 'N', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'N', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'N',
                'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'N', 'Y', 'N', 'N', 'N',
                'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y',
                'N', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y',
                'Y', 'N', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'N',
                'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
                'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'N', 'Y',
                'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'N', 'Y', 'N', 'N', 'N', 'N',
                'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'N', 'Y',
                'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N',
                'N', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'N', 'Y', 'Y', 'Y',
                'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y',
                'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
                'N', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y',
                'N', 'Y', 'N', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'N', 'Y', 'Y', 'Y',
                'Y', 'N', 'N', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y',
                'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'N', 'N', 'Y',
                'Y', 'N', 'Y', 'Y', 'Y', 'N', 'N', 'N', 'Y', 'N', 'Y', 'N', 'Y',
                'N', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y',
                'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'N', 'N', 'Y', 'N', 'Y', 'Y',
                'Y', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y',
                'Y', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y',
                'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y',
                'N', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'N', 'N', 'N',
                'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N',
                'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y',
                'N', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'Y',
                'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'N', 'N', 'Y', 'N',
                'Y', 'N', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'N',
                'N', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'N',
                'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y',
                'Y', 'Y', 'N'], dtype=object)
```

```python
In [24]: from sklearn.model_selection import train_test_split
         X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2,rand
```

```
In [25]: X_train
```

```
Out[25]: array([['Male', 'Yes', '0', ..., 1.0, 4.875197323201151, 5858.0],
                 ['Male', 'No', '1', ..., 1.0, 5.278114659230517, 11250.0],
                 ['Male', 'Yes', '0', ..., 0.0, 5.003946305945459, 5681.0],
                 ...,
                 ['Male', 'Yes', '3+', ..., 1.0, 5.298317366548036, 8334.0],
                 ['Male', 'Yes', '0', ..., 1.0, 5.075173815233827, 6033.0],
                 ['Female', 'Yes', '0', ..., 1.0, 5.204006687076795, 6486.0]],
                 dtype=object)
```

as there are some categorical variables, it should be converted to numerical hence label encoding can be used

```
In [26]: from sklearn.preprocessing import LabelEncoder
         labelencoder_X=LabelEncoder()
```

```
In [27]: for i in range(0,5):
             X_train[:,i]=labelencoder_X.fit_transform(X_train[:,i])
```

```
In [28]: X_train[:,7]=labelencoder_X.fit_transform(X_train[:,7])
```

```
In [29]: X_train
```

```
Out[29]: array([[1, 1, 0, ..., 1.0, 4.875197323201151, 267],
                 [1, 0, 1, ..., 1.0, 5.278114659230517, 407],
                 [1, 1, 0, ..., 0.0, 5.003946305945459, 249],
                 ...,
                 [1, 1, 3, ..., 1.0, 5.298317366548036, 363],
                 [1, 1, 0, ..., 1.0, 5.075173815233827, 273],
                 [0, 1, 0, ..., 1.0, 5.204006687076795, 301]], dtype=object)
```

```
In [30]: labelencoder_y=LabelEncoder()
         y_train=labelencoder_y.fit_transform(y_train)
```

```
In [31]: y_train

Out[31]: array([1, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1,
                0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1,
                1, 0, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0,
                1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1,
                1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0,
                1, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1,
                0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0,
                0, 1, 0, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1,
                0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1,
                0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1,
                1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1,
                1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1,
                1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1,
                1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0,
                1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1,
                1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0,
                1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1,
                1, 1, 1, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 1,
                1, 1, 1, 0, 1, 0, 1])

In [32]: for i in range(0,5):
             X_test[:,i]=labelencoder_X.fit_transform(X_test[:,i])

In [33]: X_test[:,7]=labelencoder_X.fit_transform(X_test[:,7])

In [34]: X_test

Out[34]: array([[1, 0, 0, 0, 5, 1.0, 4.430816798843313, 85],
                [0, 0, 0, 0, 5, 1.0, 4.718498871295094, 28],
                [1, 1, 0, 0, 5, 1.0, 5.780743515792329, 104],
                [1, 1, 0, 0, 5, 1.0, 4.700480365792417, 80],
                [1, 1, 2, 0, 5, 1.0, 4.574710978503383, 22],
                [1, 1, 0, 1, 3, 0.0, 5.10594547390058, 70],
                [1, 1, 3, 0, 3, 1.0, 5.056245805348308, 77],
                [1, 0, 0, 0, 5, 1.0, 6.003887067106539, 114],
                [1, 0, 0, 0, 5, 0.0, 4.820281565605037, 53],
                [1, 1, 0, 0, 5, 1.0, 4.852030263919617, 55],
                [0, 0, 0, 0, 5, 1.0, 4.430816798843313, 4],
                [1, 1, 1, 0, 5, 1.0, 4.553876891600541, 2],
                [0, 0, 0, 0, 5, 1.0, 5.634789603169249, 96],
                [1, 1, 2, 0, 5, 1.0, 5.4638318050256105, 97],
                [1, 1, 0, 0, 5, 1.0, 4.564348191467836, 117],
                [1, 1, 1, 0, 5, 1.0, 4.204692619390966, 22],
                [1, 0, 1, 1, 5, 1.0, 5.247024072160486, 32],
                [1, 0, 0, 1, 5, 1.0, 4.882801922586371, 25],
                [0, 0, 0, 0, 5, 1.0, 4.532599493153256, 1],

In [35]: y_test=labelencoder_y.fit_transform(y_test)
```

```
In [36]: y_test
```

```
Out[36]: array([1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1,
                 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1,
                 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1,
                 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1,
                 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0,
                 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1])
```

it is important to scale data as the varaiables are of different ranges and the algorithm fits better

```
In [37]: from sklearn.preprocessing import StandardScaler
         ss=StandardScaler()
         X_train=ss.fit_transform(X_train)
         X_test=ss.fit_transform(X_test)
```

**Application of Algorithms**

***DECISION TREE***

```
In [38]: from sklearn.tree import DecisionTreeClassifier
         DTClassifier=DecisionTreeClassifier(criterion='entropy',random_state=0)
         DTClassifier.fit(X_train,y_train)
```

```
Out[38]:                      DecisionTreeClassifier
         ▼

         DecisionTreeClassifier(criterion='entropy', random_state=0)
```

```
In [39]: y_pred=DTClassifier.predict(X_test)
         y_pred
```

```
Out[39]: array([0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
                 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1,
                 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1,
                 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1,
                 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1,
                 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1])
```

```
In [40]: from sklearn import metrics
         metrics.accuracy_score(y_pred,y_test)
```

```
Out[40]: 0.7073170731707317
```

as the accuracy score is low, another algorithm can be applied

***NAIVE BAYES***

```
In [41]: from sklearn.naive_bayes import GaussianNB
         NBClassifier=GaussianNB()
         NBClassifier.fit(X_train,y_train)
```

Out[41]:  ▾ GaussianNB

         GaussianNB()

```
In [42]: y_pred=NBClassifier.predict(X_test)
         y_pred
```

Out[42]: array([1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1,
                1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1])

In [43]: metrics.accuracy_score(y_pred,y_test)
```

Out[43]: 0.8292682926829268

## IMPLEMENTATION ON TEST DATA

```
In [44]: test = pd.read_csv('C:/Users/lekshmi/Downloads/test_LE.csv')

In [45]: test.head()
```

Out[45]:

|   | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | ApplicantIncome | Co... |
|---|---------|--------|---------|------------|-----------|---------------|-----------------|-------|
| 0 | LP001015 | Male | Yes | 0 | Graduate | No | 5720 | |
| 1 | LP001022 | Male | Yes | 1 | Graduate | No | 3076 | |
| 2 | LP001031 | Male | Yes | 2 | Graduate | No | 5000 | |
| 3 | LP001035 | Male | Yes | 2 | Graduate | No | 2340 | |
| 4 | LP001051 | Male | No | 0 | Not Graduate | No | 3276 | |

```
In [46]:  test.info()

          <class 'pandas.core.frame.DataFrame'>
          RangeIndex: 367 entries, 0 to 366
          Data columns (total 12 columns):
           #   Column             Non-Null Count  Dtype
          ---  ------             --------------  -----
           0   Loan_ID            367 non-null    object
           1   Gender             356 non-null    object
           2   Married            367 non-null    object
           3   Dependents         357 non-null    object
           4   Education          367 non-null    object
           5   Self_Employed      344 non-null    object
           6   ApplicantIncome    367 non-null    int64
           7   CoapplicantIncome  367 non-null    int64
           8   LoanAmount         362 non-null    float64
           9   Loan_Amount_Term   361 non-null    float64
           10  Credit_History     338 non-null    float64
           11  Property_Area      367 non-null    object
          dtypes: float64(3), int64(2), object(7)
          memory usage: 34.5+ KB

In [47]:  test.isnull().sum()

Out[47]:  Loan_ID               0
          Gender               11
          Married               0
          Dependents           10
          Education             0
          Self_Employed        23
          ApplicantIncome       0
          CoapplicantIncome     0
          LoanAmount            5
          Loan_Amount_Term      6
          Credit_History       29
          Property_Area         0
          dtype: int64

In [48]:  test['Gender'].fillna(test['Gender'].mode()[0], inplace=True)
          test['Dependents'].fillna(test['Dependents'].mode()[0], inplace=True)
          test['Self_Employed'].fillna(test['Self_Employed'].mode()[0], inplace=True)
          test['Credit_History'].fillna(test['Credit_History'].mode()[0], inplace=Tru
          test['Loan_Amount_Term'].fillna(test['Loan_Amount_Term'].mode()[0], inplace
          test['LoanAmount'].fillna(test['LoanAmount'].mean(), inplace=True)

In [49]:  test['LoanAmount_log']=np.log(test['LoanAmount'])
```
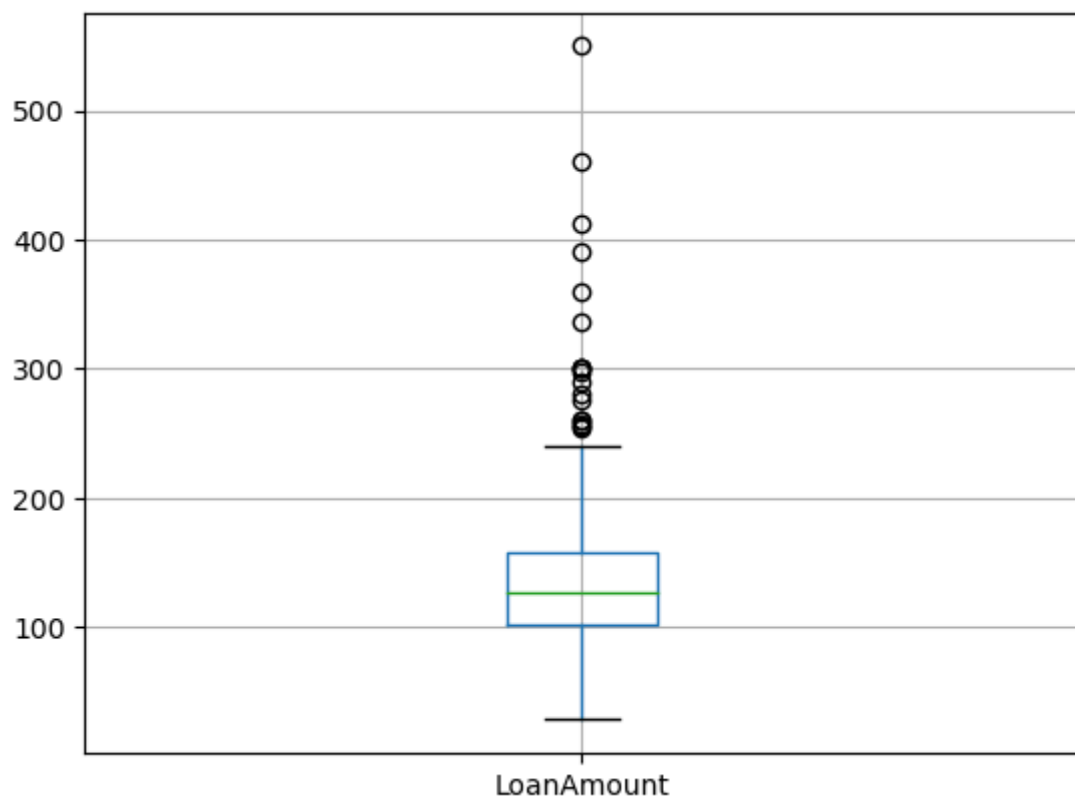
```
In [50]: test.isnull().sum()
```

```
Out[50]: Loan_ID              0
         Gender               0
         Married              0
         Dependents           0
         Education            0
         Self_Employed        0
         ApplicantIncome      0
         CoapplicantIncome    0
         LoanAmount           0
         Loan_Amount_Term     0
         Credit_History       0
         Property_Area        0
         LoanAmount_log       0
         dtype: int64
```

```
In [51]: test.boxplot(column='LoanAmount')
```

```
Out[51]: <Axes: >
```



```
In [52]: test['Total_Income']=test['ApplicantIncome']+test['CoapplicantIncome']
         test['Total_Income_log'] = np.log(test['Total_Income'])
```

```
In [53]: test.head()
```

Out[53]:

| | Loan_ID | Gender | Married | Dependents | Education | Self_Employed | ApplicantIncome | Co |
|---|---------|--------|---------|------------|-----------|---------------|-----------------|----|
| 0 | LP001015 | Male | Yes | 0 | Graduate | No | 5720 | |
| 1 | LP001022 | Male | Yes | 1 | Graduate | No | 3076 | |
| 2 | LP001031 | Male | Yes | 2 | Graduate | No | 5000 | |
| 3 | LP001035 | Male | Yes | 2 | Graduate | No | 2340 | |
| 4 | LP001051 | Male | No | 0 | Not Graduate | No | 3276 | |

```
In [54]: testX=test.iloc[:,np.r_[1:5,9:11,13:15]].values
```

```
In [55]: for i in range(0,5):
             testX[:,i]=labelencoder_X.fit_transform(testX[:,i])
```

```
In [56]: testX[:,7]=labelencoder_X.fit_transform(testX[:,7])
```

```
In [57]: testX
```

```
Out[57]: array([[1, 1, 0, ..., 1.0, 5720, 207],
                [1, 1, 1, ..., 1.0, 4576, 124],
                [1, 1, 2, ..., 1.0, 6800, 251],
                ...,
                [1, 0, 0, ..., 1.0, 5243, 174],
                [1, 1, 0, ..., 1.0, 7393, 268],
                [1, 0, 0, ..., 1.0, 9200, 311]], dtype=object)
```

```
In [58]: testX=ss.fit_transform(testX)
```

```
In [59]: pred=NBClassifier.predict(testX)
```

```
In [60]: pred
```

```
Out[60]: array([1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1,
                0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0,
                1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1,
                1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1,
                0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0,
                1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
                1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
                1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1,
                1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1])
```

```python
In [ ]:
```

```python
In [ ]:
```