

# Loan Approval Prediction

## Import Modules

```
#Importing Required Modules
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import sklearn.metrics as metrics
import warnings
warnings.filterwarnings('ignore')
print("Modules Imported Successfully!")
```

Modules Imported Successfully!

```
#Importing Modules As Required
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, f1_score, confusion_matrix
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from xgboost import XGBClassifier
from lightgbm import LGBMClassifier
print("Success!")
```

Success!

```
#Loading Dataset
df = pd.read_csv('Loan_Train.csv')
print("Dataset Loaded Successfully!")
```

Dataset Loaded Successfully!

```
#Checking the Shape of Dataset
df.shape
```

(614, 13)

There are 614 Rows and 13 Columns in our dataset.

```
#Checking Column Names
df.columns
```

Index(['Loan\_ID', 'Gender', 'Married', 'Dependents', 'Education', 'Self\_Employed', 'ApplicantIncome', 'CoapplicantIncome', 'LoanAmount', 'Loan\_Status', 'Credit\_History', 'Property\_Area', 'Loan\_Status'], dtype='object')

```
#First 5 rows of data
df.head()
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Status
0	LP001002	Male	No	No	Graduate	No	5849	0.0	NaN	360.0
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0
4	LP001008	Male	No	0	Graduate	No	6000	0.0	141.0	360.0

```
#Datatypes of features
df.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
```

```
#Statistical information about data
df.describe()
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.000000	564.000000
mean	5403.459283	1621.245798	146.412162	342.000000	0.842199
std	6109.041673	2926.248369	85.587325	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.000000	0.000000
25%	2877.500000	0.000000	100.000000	360.000000	1.000000
50%	3812.500000	1188.500000	128.000000	360.000000	1.000000
75%	5795.000000	2297.250000	168.000000	360.000000	1.000000
max	81000.000000	41667.000000	700.000000	480.000000	1.000000

```
#Checking unique values in each feature
df.apply(lambda x: x.nunique())
```

Loan\_ID 614  
Gender 3  
Married 3  
Dependents 5  
Education 4  
Self\_Employed 3  
ApplicantIncome 505  
CoapplicantIncome 287  
LoanAmount 264  
Loan\_Amount\_Term 11  
Credit\_History 3  
Property\_Area 3  
Loan\_Status 2  
dtype: int64

```
#Datatypes of features
df.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
```

```
#Checking null values in data, if any
df.isnull().sum()
```

Loan\_ID 0  
Gender 3  
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  
dtype: int64

```
#Loan_ID as index and then dropping it
df.index = df['Loan_ID']
df = df.drop(columns='Loan_ID', axis=1)
df.head()
```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
LP001002	Male	No	0	Graduate	No	5849	0.0	146.412162	360.0	
LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.000000	360.0	
LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.000000	360.0	
LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.000000	360.0	
LP001008	Male	No	0	Graduate	No	6000	0.0	141.000000	360.0	

```
#Checking for Duplicate Rows
df.duplicated().sum()
```

0

```
#Checking data before moving to EDA
df.head()
```

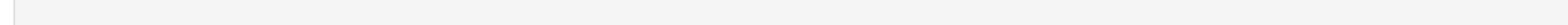
	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
LP001002	Male	No	0	Graduate	No	5849	0.0	146.412162	360.0	
LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.000000	360.0	
LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.000000	360.0	
LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.000000	360.0	
LP001008	Male	No	0	Graduate	No	6000	0.0	141.000000	360.0	

## Exploratory Data Analysis

### Loan Status

```
#Countplot for Loan_Status
ax = sns.countplot(x='Loan_Status', data=df)
ax.bar_label(ax.containers[0])
```

(Text(0, 0, '422'), Text(0, 0, '192'))

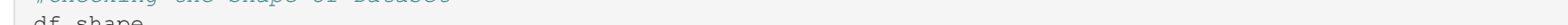


Nearly 68% loans got approved.

### Gender

```
#Countplot for loans on the basis of gender
sns.countplot(x=df['Gender'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Gender', ylabel='count'>



There are more Men than Women. (Approx. 3x)

### Married

```
#Countplot for loans on the basis of marital status
sns.countplot(x=df['Married'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Married', ylabel='count'>

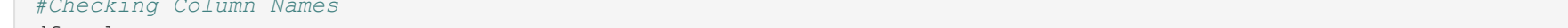


2/3rd of the population in the dataset is Married and Married applicants are more likely to be granted loans.

### Dependents

```
#Countplot for loans on the basis of Dependents
sns.countplot(x=df['Dependents'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Dependents', ylabel='count'>

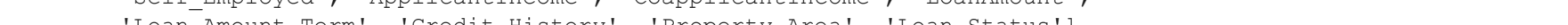


Majority of the population have 0 dependents and are also likely to be accepted for loan.

### Education

```
#Countplot for loans on the basis of Education
sns.countplot(x=df['Education'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Education', ylabel='count'>

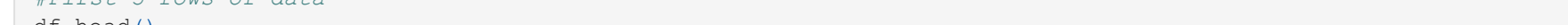


Nearly 5/6th population is graduate and are more likely to be approved for loan

### Self Employed

```
#Countplot for loans on the basis of Employment
sns.countplot(x=df['Self_Employed'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Self\_Employed', ylabel='count'>



5/6th of the population is not self-employed.

### Loan Amount Term

```
#Countplot for loans on the basis of Term
sns.countplot(x=df['Loan_Amount_Term'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Loan\_Amount\_Term', ylabel='count'>



Majority of the loans are taken for 360 Months.(30 Years)

### Credit History

```
#Countplot for loans on the basis of Credit History
sns.countplot(x=df['Credit_History'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Credit\_History', ylabel='count'>



Applicants with credit history are more likely to be approved.

### Property Area

```
#Countplot for loans on the basis of Area
sns.countplot(x=df['Property_Area'], data=df, hue='Loan_Status')
```

<AxesSubplot:xlabel='Property\_Area', ylabel='count'>

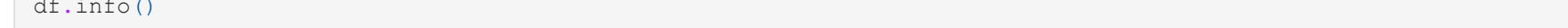


More applicants from Semi-urban and also more likely to be granted loans

### Applicant Income

```
#Histogram for Applicant Income
sns.distplot(df['ApplicantIncome'])
```

<AxesSubplot:xlabel='ApplicantIncome', ylabel='Density'>



### Coapplicant Income

```
#Histogram for Coapplicant Income
sns.distplot(df['CoapplicantIncome'])
```

<AxesSubplot:xlabel='CoapplicantIncome', ylabel='Density'>



### Loan Amount

```
#Histogram for Loan Amount
sns.distplot(df['LoanAmount'])
```

<AxesSubplot:xlabel='LoanAmount', ylabel='Density'>



Other than the skewedness of numerical data, there is nothing much to correlate the data in numerical features.

## Preprocessing

```
#Encoding categorical values to numerical
df = pd.get_dummies(df, drop_first=True)
```

```
#Splitting the data into Feature and Target Variable
X = df.drop(columns = 'Loan_Status_Y')
y = df['Loan_Status_Y']
```

```
#Splitting the data in Training and Test set
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,stratify=y,random_state =42)
```

## Training the models

### Logistic Regression

```
LR = LogisticRegression()
LR.fit(X_train,y_train)
```

LogisticRegression()

y\_pred = LR.predict(X\_test)

```
acc_LR = accuracy_score(y_test,y_pred)
f1_LR = f1_score(y_test,y_pred)
print("Accuracy: ", acc_LR)
print("F1 Score: ", f1_LR)
```

Accuracy: 0.8048780487804879  
F1 Score: 0.8666666666666667

### Gaussian Naive Bayes

```
GNB = GaussianNB()
GNB.fit(X_train,y_train)
```

GaussianNB()

y\_pred = GNB.predict(X\_test)

```
acc_GNB = accuracy_score(y_test,y_pred)
f1_GNB = f1_score(y_test,y_pred)
print("Accuracy: ", acc_GNB)
print("F1 Score: ", f1_GNB)
```

Accuracy: 0.7642276422764228  
F1 Score: 0.8432652652652652

### Decision Tree

```
DT = DecisionTreeClassifier()
DT.fit(X_train,y_train)
```

DecisionTreeClassifier()

y\_pred = DT.predict(X\_test)

```
acc_DT = accuracy_score(y_test,y_pred)
f1_DT = f1_score(y_test,y_pred)
print("Accuracy: ", acc_DT)
print("F1 Score: ", f1_DT)
```

Accuracy: 0.6829268292682927  
F1 Score: 0.7636363636363637

### Random Forest

```
RF = RandomForestClassifier()
RF.fit(X_train,y_train)
```

RandomForestClassifier()

y\_pred = RF.predict(X\_test)

```
acc_RF = accuracy_score(y_test,y_pred)
f1_RF = f1_score(y_test,y_pred)
print("Accuracy: ", acc_RF)
print("F1 Score: ", f1_RF)
```

Accuracy: 0.7804878048780488  
F1 Score: 0.8432652652652652

### Support Vector Machine

```
SVM = SVC()
SVM.fit(X_train,y_train)
```

SVC()

y\_pred = SVM.predict(X\_test)

```
acc_SVM = accuracy_score(y_test,y_pred)
f1_SVM = f1_score(y_test,y_pred)
print("Accuracy: ", acc_SVM)
print("F1 Score: ", f1_SVM)
```

Accuracy: 0.6910569105691057  
F1 Score: 0.8173076923076924

### K-Nearest Neighbors

```
KNN = KNeighborsClassifier()
KNN.fit(X_train,y_train)
```

KNeighborsClassifier()

y\_pred = KNN.predict(X\_test)

```
acc_KNN = accuracy_score(y_test,y_pred)
f1_KNN = f1_score(y_test,y_pred)
print("Accuracy: ", acc_KNN)
print("F1 Score: ", f1_KNN)
```

Accuracy: 0.8432652652652652  
F1 Score: 0.762430932265194

### XGBoost

```
XGB = XGBClassifier(eval_metric='mlogloss')
XGB.fit(X_train,y_train)
```

XGBClassifier(base\_score=0.5, booster='gbtree', colsample\_bylevel=1, colsample\_byrow=1, colsample\_bynode=1, colsample\_bytree=1, eval\_metric='mlogloss', gamma=0, gpu\_id=-1, importance\_type='gain', interaction\_constraints='', learning\_rate=0.300000012, max\_delta\_step=0, max\_depth=6, min\_child\_weight=1, missing=nan, monotone\_constraints=(), n\_estimators=100, n\_jobs=4, num\_parallel\_tree=1, random\_state=0, reg\_alpha=0, reg\_lambda=1, scale\_pos\_weight=1, subsample=1, tree\_method='exact', validate\_parameters=1, verbosity=None)

y\_pred = XGB.predict(X\_test)

```
acc_XGB = accuracy_score(y_test,y_pred)
f1_XGB = f1_score(y_test,y_pred)
print("Accuracy: ", acc_XGB)
print("F1 Score: ", f1_XGB)
```

Accuracy: 0.7317073170731707  
F1 Score: 0.807017543156492

### Light GBM

```
LGBM = LGBMClassifier()
LGBM.fit(X_train,y_train)
```

LGBMClassifier()

y\_pred = LGBM.predict(X\_test)

```
acc_LGBM = accuracy_score(y_test,y_pred)
f1_LGBM = f1_score(y_test,y_pred)
print("Accuracy: ", acc_LGBM)
print("F1 Score: ", f1_LGBM)
```

Accuracy: 0.7479674796747967  
F1 Score: 0.820809248594133

## Evaluation

### Accuracy Scores

```
models = pd.DataFrame({
    'Model': ['Logistic Regression', 'Naive Bayes', 'Decision Tree', 'Random Forest', 'Support Vector Machines',
              'K - Nearest Neighbors', 'XGBost Classifier', 'Light Gradient Boosting Machine'],
    'Score': [acc_LR, acc_GNB, acc_DT, acc_RF, acc_SVM, acc_KNN, acc_XGB, acc_LGBM]})
models.sort_values(by='Score', ascending=False)
```

Model Score

0 Logistic Regression 0.804878

3 Random Forest 0.780488

1 Naive Bayes 0.764228

7 Light Gradient Boosting Machine 0.747967

6 XGBoost Classifier 0.731707

4 Support Vector Machines 0.691057

2 Decision Tree 0.682927

5 K - Nearest Neighbors 0.650407

### F1 Scores

```
models = pd.DataFrame({
    'Model': ['Logistic Regression', 'Naive Bayes', 'Decision Tree', 'Random Forest', 'Support Vector Machines',
              'K - Nearest Neighbors', 'XGBost Classifier', 'Light Gradient Boosting Machine'],
    'Score': [f1_LR, f1_GNB, f1_DT, f1_RF, f1_SVM, f1_KNN, f1_XGB, f1_LGBM]})
models.sort_values(by='Score', ascending=False)
```

Model Score

0 Logistic Regression 0.866667

3 Random Forest 0.843931

1 Naive Bayes 0.843243

7 Light Gradient Boosting Machine 0.820809

4 Support Vector Machines 0.817308

6 XGBoost Classifier 0.807018

2 Decision Tree 0.763636

5 K - Nearest Neighbors 0.762431

### Confusion Matrix

```
#Confusion Matrix using Logistic Regression
y_pred = LR.predict(X_test)
cmf = confusion_matrix(y_test, y_pred)
cmf
```

array([[21, 17],
 [ 7, 78]], dtype=int64)

```
print("Confusion Matrix on Test Data")
sns.heatmap(pd.DataFrame(cmf), annot=True, cmap='YlGnBu')
plt.show()
```

Confusion Matrix on Test Data

