

Loan Prediction

Binary Classification using Logistic Regression

Importing Libraries

In [161]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

Importing & Loading the dataset

In []:

In [123]:

```
df = pd.read_csv('train.csv')
df.head()
```

Out[123]:

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	Coa
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	

Dataset Info:

In [124]:

```
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
```

Dataset Shape:

In [125]:

```
df.shape
```

Out[125]:

```
(614, 13)
```

Data Cleaning

Checking the Missing Values

In [126]:

```
df.isnull().sum()
```

Out[126]:

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

dtype: int64

First we will fill the Missing Values in "LoanAmount" & "Credit_History" by the 'Mean' & 'Median' of the respective variables.

In [127]:

```
df['LoanAmount'] = df['LoanAmount'].fillna(df['LoanAmount'].mean())
```

In [128]:

```
df['Credit_History'] = df['Credit_History'].fillna(df['Credit_History'].median())
```

Let's confirm if there are any missing values in 'LoanAmount' & 'Credit_History'

In [129]:

```
df.isnull().sum()
```

Out[129]:

```
Loan_ID          0
Gender           13
Married          3
Dependents       15
Education         0
Self_Employed    32
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount        0
Loan_Amount_Term 14
Credit_History   0
Property_Area     0
Loan_Status       0
dtype: int64
```

Now, Let's drop all the missing values remaining.

In [130]:

```
df.dropna(inplace=True)
```

Let's check the Missing values for the final time!

In [131]:

```
df.isnull().sum()
```

Out[131]:

```
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
dtype: int64
```

Here, we have dropped all the missing values to avoid disturbances in the model. The Loan Prediction requires all the details to work efficiently and thus the missing values are dropped.

Now, Let's check the final Dataset Shape

In [132]:

```
df.shape
```

Out[132]:

```
(542, 13)
```

Exploratory Data Analysis

Comparison between Parameters in getting the Loan:

In [133]:

```
plt.figure(figsize=(100,50))
sns.set(font_scale=5)

plt.subplot(331)
sns.countplot(x = 'Gender', hue = 'Loan_Status', data = df)

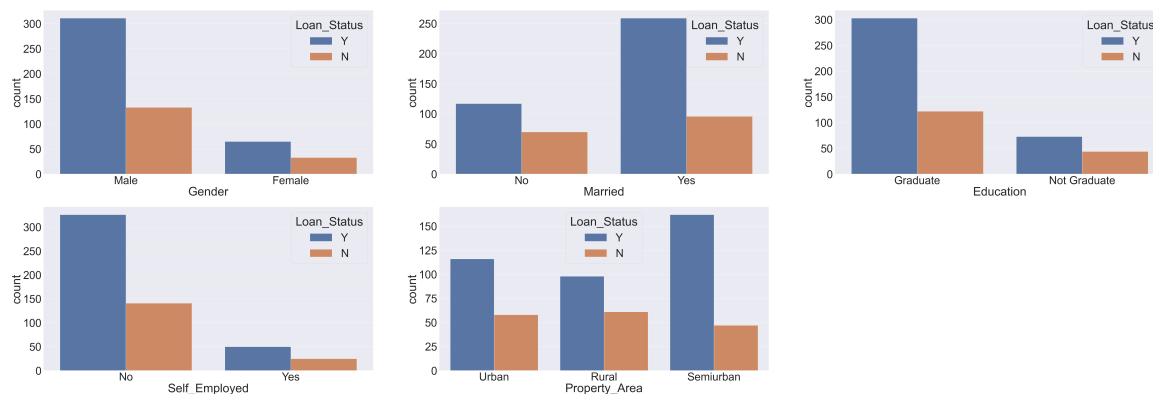
plt.subplot(332)
sns.countplot(x = 'Married', hue = 'Loan_Status', data = df)

plt.subplot(333)
sns.countplot(x = 'Education', hue = 'Loan_Status', data = df)

plt.subplot(334)
sns.countplot(x = 'Self_Employed', hue = 'Loan_Status', data = df)

plt.subplot(335)
sns.countplot(x = 'Property_Area', hue = 'Loan_Status', data = df)

plt.show()
```



Let's replace the Variable values to Numerical form & display the Value Counts

The data in Numerical form avoids disturbances in building the model.

In [134]:

```
df['Loan_Status'].replace('Y',1,inplace=True)  
df['Loan_Status'].replace('N',0,inplace=True)
```

In [135]:

```
df['Loan_Status'].value_counts()
```

Out[135]:

```
1    376  
0    166  
Name: Loan_Status, dtype: int64
```

In [136]:

```
df.Gender=df.Gender.map({'Male':1,'Female':0})  
df['Gender'].value_counts()
```

Out[136]:

```
1    444  
0     98  
Name: Gender, dtype: int64
```

In [137]:

```
df.Married = df.Married.map({'Yes':1, 'No': 0})  
df['Married'].value_counts()
```

Out[137]:

```
1    355  
0    187  
Name: Married, dtype: int64
```

In [138]:

```
df.Dependents=df.Dependents.map({'0':0, '1':1, '2':2, '3+':3})  
df['Dependents'].value_counts()
```

Out[138]:

```
0    309  
1     94  
2     94  
3     45  
Name: Dependents, dtype: int64
```

In [139]:

```
df.Education = df.Education.map({'Graduate':1, 'Not Graduate':0})  
df['Education'].value_counts()
```

Out[139]:

```
1    425  
0    117  
Name: Education, dtype: int64
```

In [140]:

```
df.Self_Employed = df.Self_Employed.map({'Yes':1, 'No':0})  
df['Self_Employed'].value_counts()
```

Out[140]:

```
0    467  
1     75  
Name: Self_Employed, dtype: int64
```

In [141]:

```
df.Property_Area = df.Property_Area.map({'Urban':2, 'Rural':0, 'Semiurban':1})  
df['Property_Area'].value_counts()
```

Out[141]:

```
1    209  
2    174  
0    159  
Name: Property_Area, dtype: int64
```

In [142]:

```
df['LoanAmount'].value_counts()
```

Out[142]:

```
146.412162    19  
120.000000     15  
100.000000     14  
110.000000     13  
187.000000     12  
..  
280.000000      1  
240.000000      1  
214.000000      1  
59.000000       1  
253.000000      1  
Name: LoanAmount, Length: 195, dtype: int64
```

In [143]:

```
df['Loan_Amount_Term'].value_counts()
```

Out[143]:

```
360.0    464  
180.0     38  
480.0     13  
300.0     12  
84.0       4  
120.0       3  
240.0       3  
60.0        2  
36.0         2  
12.0         1  
Name: Loan_Amount_Term, dtype: int64
```

In [144]:

```
df['Credit_History'].value_counts()
```

Out[144]:

```
1.0    468
0.0     74
Name: Credit_History, dtype: int64
```

From the above figure, we can see that **Credit_History** (Independent Variable) has the maximum correlation with **Loan_Status** (Dependent Variable). Which denotes that the Loan_Status is heavily dependent on the Credit_History.

Final DataFrame

In [145]:

```
df.head()
```

Out[145]:

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome
0	LP001002	1	0	0	1	0	5849	
1	LP001003	1	1	1	1	0	4583	
2	LP001005	1	1	0	1	1	3000	
3	LP001006	1	1	0	0	0	2583	
4	LP001008	1	0	0	1	0	6000	

Importing Packages for Classification algorithms

In [148]:

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn import metrics
```

Splitting the data into Train and Test set

In [149]:

```
X = df.iloc[1:542,1:12].values
y = df.iloc[1:542,12].values
```

In [152]:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)
```


Logistic Regression (LR)

Logistic regression is a supervised learning classification algorithm used to predict the probability of a target variable.

Mathematically, a logistic regression model predicts $P(Y=1)$ as a function of X . It is one of the simplest ML algorithms that can be used for various classification problems such as spam detection, Diabetes prediction, cancer detection etc.

Sigmoid Function

In [156]:

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

lr_prediction = model.predict(X_test)
print('Logistic Regression accuracy = ', metrics.accuracy_score(lr_prediction,y_test))

Logistic Regression accuracy = 0.7914110429447853
```

In [165]:

```
print("y_predicted",lr_prediction)
print("y_test",y_test)
```

```
y_predicted [1 0 0 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
0 0 0 1 1 0
1 1 1 1 0 0 1 1 1 1 1 1 0 1 1 0 1 1 1 0 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 0 1 0 1 1 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 0 1 1 1 1 1 1 1 1 0 1 1 1 1 0 1 1 0 1
1 0 1 0 1 1 1 1 1 1 1 1 1 1 0]
y_test [0 0 0 0 0 1 0 1 1 0 1 1 1 1 0 0 1 1 1 0 1 0 1 1 1 1 1 0 1 1 0 0
0 0 1 0
1 1 1 1 0 1 0 1 1 1 1 1 0 1 1 0 1 0 1 0 1 1 1 1 0 1 0 1 0 1 0 1 1
1 1 1 1 1 1 0 0 0 1 0 0 0 1 0 1 0 1 0 1 1 1 1 0 1 0 1 1 0 1 1 1 1
1 1 1 0 1 1 1 1 1 1 1 0 1 1 0 1 1 1 1 1 0 1 1 1 1 0 0 1 1 0 0
1 0 0 0 0 1 0 1 0 1 1 1 1 1 0]
```

CONCLUSION:

1. The Loan Status is heavily dependent on the Credit History for Predictions.
2. The Logistic Regression algorithm gives us the maximum Accuracy (79% approx) compared to the other 3 Machine Learning Classification Algorithms.

In []: