

```

# Project Name: Loan Prediction System

# Project Description
# Loading & Exploring the Data
# Working With Missing Values
# Dropping Unnecessary Columns
# Visualization For Data Presentation
# Encoding the Categorical Data
# Model Deployment:
# Data Segregation
# Using GaussianNB
# Loss Function
# Using SVC with grid search CV
# XGBoost Classifier
# Decision tree using randomized search
# Random forest using randomized search
# saving the model
# coding a web app that will take user input and predict chances of user getting loan

# Hello Folks, My name is Nirbhay Tiwari I am Data Scientist at NibbusPost, In this Project We will see How we Deal with Classification Problem using Supervised Algorithm

In [99]:
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib inline

In [100]:
train_data = pd.read_csv('loan_data_vst.csv')

In [101]:
train_data.head()

Out[101]:
   Loan_ID  Gender  Married  Dependents  Education  Self_Employed  ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  Credit_History  Property_Area  Loan_Status
0  LP001002   Male    No         0         Graduate    No          5849              0.0         NaN          360.0          1.0         Urban          Y
1  LP001003   Male    Yes        1         Graduate    No          4583             1508.0        128.0          360.0          1.0         Rural          N
2  LP001005   Male    Yes        0         Graduate    Yes          3000              0.0         66.0          360.0          1.0         Urban          Y
3  LP001006   Male    Yes        0         Not Graduate    No          2583             2358.0       120.0          360.0          1.0         Urban          Y
4  LP001008   Male    No         0         Graduate    No          6000              0.0        141.0          360.0          1.0         Urban          Y

In [102]:
# Lets check the shape of the data
train_data.shape

Out[102]:
(614, 13)

In [103]:
train_data.describe(include='all')

Out[103]:
   Loan_ID  Gender  Married  Dependents  Education  Self_Employed  ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  Credit_History  Property_Area  Loan_Status
count      614      601      611      599      614      592      614.000000      614.000000      592.000000      600.000000      564.000000      614      614
unique      614      2      2      4      2      2      NaN      NaN      NaN      NaN      NaN      NaN      3      2
top  LP001002   Male    Yes        0         Graduate    No          NaN      NaN      NaN      NaN      NaN      NaN      Semrurban      Y
freq         1      489      398      345      480      500      NaN      NaN      NaN      NaN      NaN      NaN      233      422
mean      NaN      NaN      NaN      NaN      NaN      NaN      NaN      5403.459283      1621.245798      146.412162      342.000000      0.842199      NaN      NaN
std       NaN      NaN      NaN      NaN      NaN      NaN      NaN      6109.014173      2926.248369      85.587325      65.12041      0.364878      NaN      NaN
min       NaN      NaN      NaN      NaN      NaN      NaN      NaN      150.000000      0.000000      9.000000      12.000000      0.000000      NaN      NaN
25%      NaN      NaN      NaN      NaN      NaN      NaN      NaN      2877.500000      0.000000      100.000000      360.000000      1.000000      NaN      NaN
50%      NaN      NaN      NaN      NaN      NaN      NaN      NaN      3812.500000      1188.500000      128.000000      360.000000      1.000000      NaN      NaN
75%      NaN      NaN      NaN      NaN      NaN      NaN      NaN      5795.000000      2297.250000      168.000000      360.000000      1.000000      NaN      NaN
max       NaN      NaN      NaN      NaN      NaN      NaN      NaN      81000.000000      41667.000000      700.000000      480.000000      1.000000      NaN      NaN

In [104]:
train_data.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  592 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 [105]:
def missing_values(df):
    n = df.isnull().sum()
    return n

In [106]:
missing_values(train_data)

Out[106]:
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

In [107]:
train_data.drop(["Loan_ID", "Dependents"], axis=1, inplace=True)

In [108]:
train_data

Out[108]:
   Gender  Married  Education  Self_Employed  ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  Credit_History  Property_Area  Loan_Status
0      Male    No      Graduate    No          5849              0.0         NaN          360.0          1.0         Urban          Y
1      Male    Yes     Graduate    No          4583             1508.0        128.0          360.0          1.0         Rural          N
2      Male    Yes     Graduate    Yes          3000              0.0         66.0          360.0          1.0         Urban          Y
3      Male    Yes  Not Graduate    No          2583             2358.0       120.0          360.0          1.0         Urban          Y
4      Male    No      Graduate    No          6000              0.0        141.0          360.0          1.0         Urban          Y
...
609     Female  No      Graduate    No          2900              0.0         71.0          360.0          1.0         Rural          Y
610     Male    Yes     Graduate    No          4106              0.0         40.0          180.0          1.0         Rural          Y
611     Male    Yes     Graduate    No          8072             240.0       253.0          360.0          1.0         Urban          Y
612     Male    Yes     Graduate    No          7583              0.0       187.0          360.0          1.0         Urban          Y
613     Female  No      Graduate    Yes          4583              0.0       133.0          360.0          0.0         Semrurban  N
614 rows x 11 columns

In [109]:
# Dealing with Categorical Missing Values Replacing All the missing category with most frequent value using mode function
cols = train_data[["Gender", "Married", "Self_Employed"]]
for i in cols:
    train_data[i].fillna(train_data[i].mode()[0], inplace=True)

In [110]:
train_data.isnull().sum()

Out[110]:
Gender      0
Married     0
Self_Employed 0
ApplicantIncome 0
CoapplicantIncome 0
LoanAmount   22
Loan_Amount_Term 14
Credit_History 50
Property_Area 0
Loan_Status  0
dtype: int64

In [111]:
# Dealing with Numerical Missing Values Replacing all missing values with mean value
n_cols = train_data[["LoanAmount", "Loan_Amount_Term", "Credit_History"]]
for i in n_cols:
    train_data[i].fillna(train_data[i].mean(axis=0), inplace=True)

In [112]:
# Now we will create a function that will take parameter and show bar chart
def bar_chart(col):
    Approved = train_data[train_data["Loan_Status"]=="Y"][col].value_counts()
    Disapproved = train_data[train_data["Loan_Status"]=="N"][col].value_counts()
    df1 = pd.DataFrame([Approved, Disapproved])
    df1.index = ["Approved", "Disapproved"]
    df1.plot(kind='bar')

In [113]:
bar_chart("Gender"), bar_chart("Married"), bar_chart("Education"), bar_chart("Self_Employed")

Out[113]:
(None, None, None, None)

In [114]:
# Lets Encode the categorical columns value by assigning a numeric identifier to each of them
from sklearn.preprocessing import OrdinalEncoder
ord_enc = OrdinalEncoder()
train_data[["Gender", "Married", "Education", "Self_Employed", "Property_Area", "Loan_Status"]] = ord_enc.fit_transform(train_data[["Gender", "Married", "Education", "Self_Employed", "Property_Area", "Loan_Status"]])

In [115]:
train_data.head()

Out[115]:
   Gender  Married  Education  Self_Employed  ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  Credit_History  Property_Area  Loan_Status
0      1.0      0.0      0.0      0.0          5849              0.0      146.412162      360.0          1.0      2.0      1.0
1      1.0      1.0      0.0      0.0          4583             1508.0       128.000000      360.0          1.0      0.0      0.0
2      1.0      1.0      1.0      1.0          3000              0.0      66.000000      360.0          1.0      2.0      1.0
3      1.0      1.0      1.0      0.0          2583             2358.0      120.000000      360.0          1.0      2.0      1.0
4      1.0      0.0      0.0      0.0          6000              0.0     141.000000      360.0          1.0      2.0      1.0

In [116]:
# Now Lets change the decimal value with whole number by converting the datatype as integer
train_data[["Gender", "Married", "Education", "Self_Employed", "Property_Area", "Loan_Status"]] = train_data[["Gender", "Married", "Education", "Self_Employed", "Property_Area", "Loan_Status"]].astype('int64')

In [117]:
train_data

Out[117]:
   Gender  Married  Education  Self_Employed  ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  Credit_History  Property_Area  Loan_Status
0      1      0      0      0          5849              0.0      146.412162      360.0          1.0      2      1
1      1      1      0      0          4583             1508.0     128.000000      360.0          1.0      0      0
2      1      1      1      1          3000              0.0      66.000000      360.0          1.0      2      1
3      1      1      1      0          2583             2358.0     120.000000      360.0          1.0      2      1
4      1      0      0      0          6000              0.0     141.000000      360.0          1.0      2      1
...
609      0      0      0      0          2900              0.0      71.000000      360.0          1.0      0      1
610      1      1      0      0          4106              0.0      40.000000      180.0          1.0      0      1
611      1      1      0      0          8072             240.0     253.000000      360.0          1.0      2      1
612      1      1      0      0          7583              0.0     187.000000      360.0          1.0      2      1
613      0      0      0      1          4583              0.0     133.000000      360.0          0.0      1      0
614 rows x 11 columns

In [118]:
from sklearn.model_selection import train_test_split
X = train_data.drop("Loan_Status", axis=
```