```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier

df = pd.read_csv('/content/loan_prediction.csv')
print(df.head())
```

```
⊡      Loan_ID Gender Married Dependents      Education Self_Employed  \
    0  LP001002   Male      No          0      Graduate            No
    1  LP001003   Male     Yes          1      Graduate            No
    2  LP001005   Male     Yes          0      Graduate           Yes
    3  LP001006   Male     Yes          0  Not Graduate            No
    4  LP001008   Male      No          0      Graduate            No

       ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  \
    0             5849                0.0         NaN             360.0
    1             4583             1508.0       128.0             360.0
    2             3000                0.0        66.0             360.0
    3             2583             2358.0       120.0             360.0
    4             6000                0.0       141.0             360.0

       Credit_History Property_Area Loan_Status
    0             1.0         Urban           Y
    1             1.0         Rural           N
    2             1.0         Urban           Y
    3             1.0         Urban           Y
    4             1.0         Urban           Y
```

```python
#I'll drop the loan id column and move further:

df.drop('Loan_ID', axis=1)

df = df.drop('Loan_ID', axis=1)
#Now let's have a look if the data has missing values or not:


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

```
    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
```

```python
#Now let's fill in the missing values. In categorical columns, we can fill in missing values with the mode of each column. The mode represents the v

# Fill missing values in categorical columns with mode

df['Gender'].fillna(df['Gender'].mode()[0], inplace=True)

df['Married'].fillna(df['Married'].mode()[0], inplace=True)

df['Dependents'].fillna(df['Dependents'].mode()[0], inplace=True)

df['Self_Employed'].fillna(df['Self_Employed'].mode()[0], inplace=True)
```

To fill in the missing values of numerical columns, we have to choose appropriate measures:

1.We can fill in the missing values of the loan amount column with the median value. The median is an appropriate measure to fill in missing values when dealing with skewed distributions or when outliers are present in the data;

2.We can fill in the missing values of the loan amount term column with the mode value of the column. Since the term of the loan amount is a discrete value, the mode is an appropriate metric to use;

3.We can fill in the missing values of the credit history column with the mode value. Since credit history is a binary variable (0 or 1), the mode represents the most common value and is an appropriate choice for filling in missing values.

```python
# Fill missing values in LoanAmount with the median
df['LoanAmount'].fillna(df['LoanAmount'].median(), inplace=True)
```

```
# Fill missing values in Loan_Amount_Term with the mode
df['Loan_Amount_Term'].fillna(df['Loan_Amount_Term'].mode()[0], inplace=True)

# Fill missing values in Credit_History with the mode
df['Credit_History'].fillna(df['Credit_History'].mode()[0], inplace=True)
```
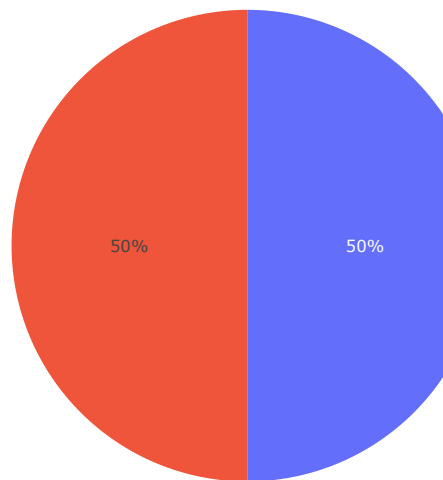
**Exploratory Data Analysis**

Now let's have a look at the distribution of the loan status column:

```
import plotly.express as px

loan_status_count = df['Loan_Status'].value_counts()
fig_loan_status = px.pie(loan_status_count,
                         names=loan_status_count.index,
                         title='Loan Approval Status')
fig_loan_status.show()
```

Loan Approval Status



Now let's have a look at the distribution of the gender column:
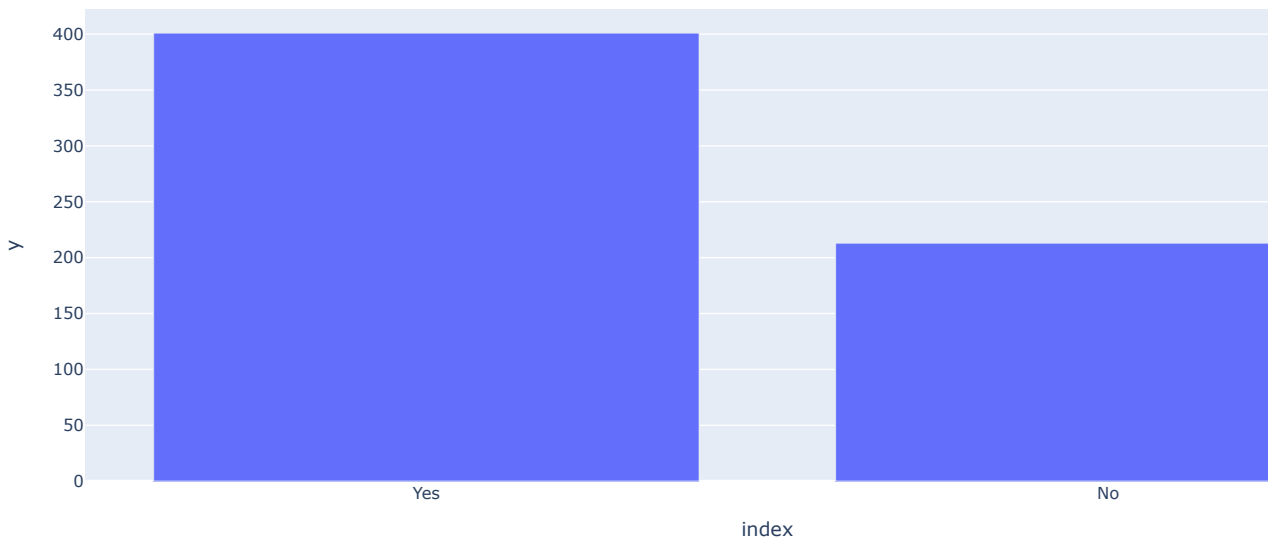
```
gender_count = df['Gender'].value_counts()
fig_gender = px.bar(gender_count,
                    x=gender_count.index,
                    y=gender_count.values,
                    title='Gender Distribution')
fig_gender.show()
```

Now let's have a look at the distribution of the martial status column:

```
married_count = df['Married'].value_counts()
fig_married = px.bar(married_count,
                     x=married_count.index,
                     y=married_count.values,
                     title='Marital Status Distribution')
fig_married.show()
```
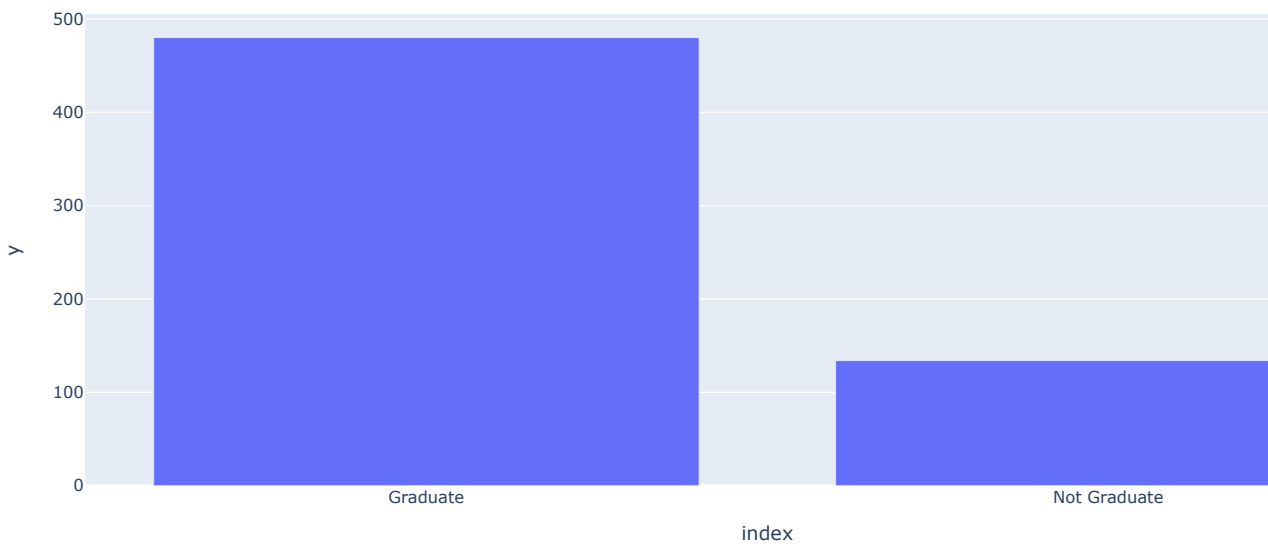
## Marital Status Distribution



Now let's have a look at the distribution of the education column:

```
education_count = df['Education'].value_counts()
fig_education = px.bar(education_count,
                       x=education_count.index,
                       y=education_count.values,
                       title='Education Distribution')
fig_education.show()
```
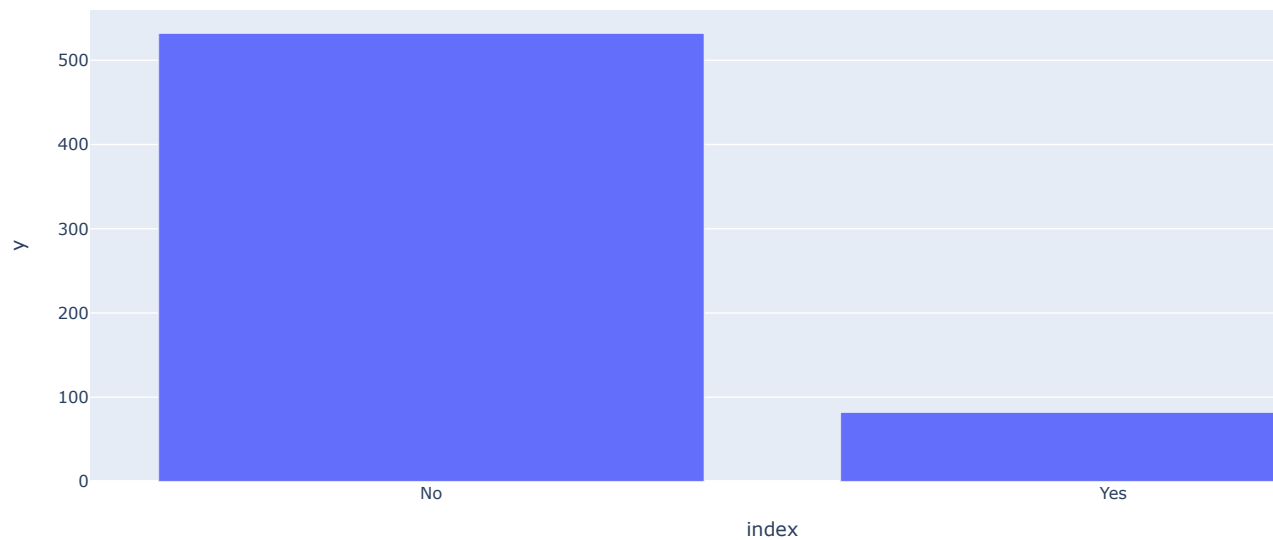
## Education Distribution

Now let's have a look at the distribution of the self-employment column:

```
self_employed_count = df['Self_Employed'].value_counts()
fig_self_employed = px.bar(self_employed_count,
                           x=self_employed_count.index,
                           y=self_employed_count.values,
                           title='Self-Employment Distribution')
fig_self_employed.show()
```
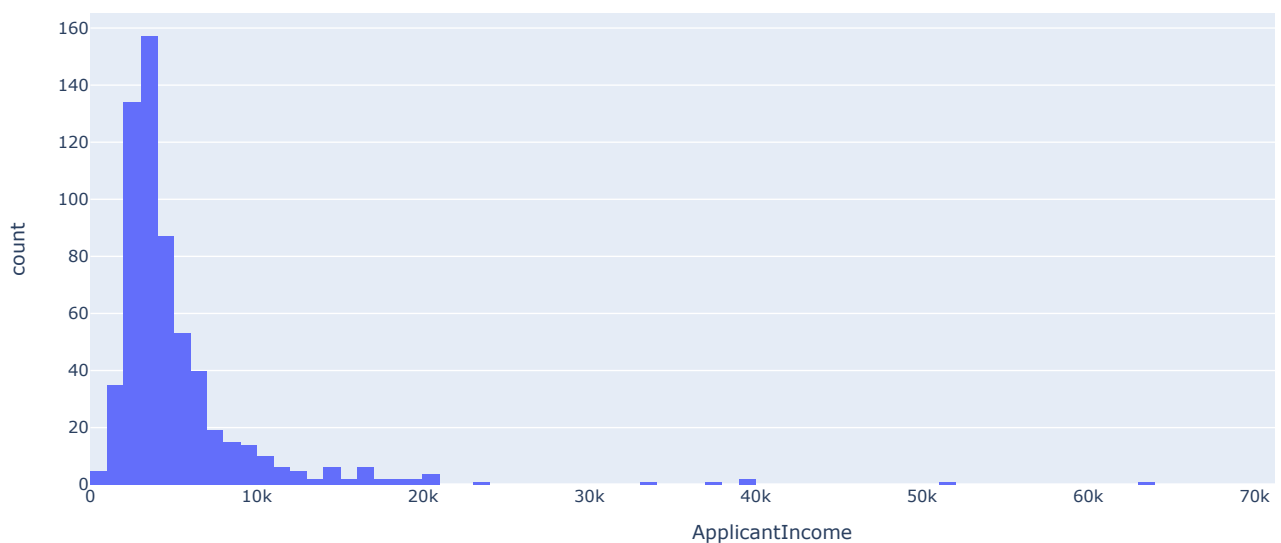
## Self-Employment Distribution



Now let's have a look at the distribution of the Applicant Income column:

```
fig_applicant_income = px.histogram(df, x='ApplicantIncome',
                                    title='Applicant Income Distribution')
fig_applicant_income.show()
```
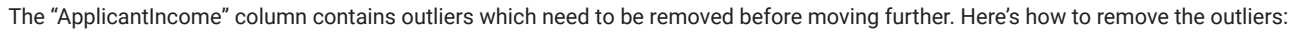
## Applicant Income Distribution



Now let's have a look at the relationship between the income of the loan applicant and the loan status:

```
fig_income = px.box(df, x='Loan_Status',
                y='ApplicantIncome',
                color="Loan_Status",
```

```
                 title='Loan_Status vs ApplicantIncome')
fig_income.show()
```

## Loan_Status vs ApplicantIncome



The "ApplicantIncome" column contains outliers which need to be removed before moving further. Here's how to remove the outliers:

```
# Calculate the IQR
Q1 = df['ApplicantIncome'].quantile(0.25)
Q3 = df['ApplicantIncome'].quantile(0.75)
IQR = Q3 - Q1

# Define the lower and upper bounds for outliers
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Remove outliers
df = df[(df['ApplicantIncome'] >= lower_bound) & (df['ApplicantIncome'] <= upper_bound)]
```

Now let's have a look at the relationship between the income of the loan co-applicant and the loan status:

```
fig_coapplicant_income = px.box(df,
                                x='Loan_Status',
                                y='CoapplicantIncome',
                                color="Loan_Status",
                                title='Loan_Status vs CoapplicantIncome')
fig_coapplicant_income.show()
```

The income of the loan co-applicant also contains outliers. Let's remove the outliers from this column as well:

```python
# Calculate the IQR
Q1 = df['CoapplicantIncome'].quantile(0.25)
Q3 = df['CoapplicantIncome'].quantile(0.75)
IQR = Q3 - Q1

# Define the lower and upper bounds for outliers
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Remove outliers
df = df[(df['CoapplicantIncome'] >= lower_bound) & (df['CoapplicantIncome'] <= upper_bound)]
```
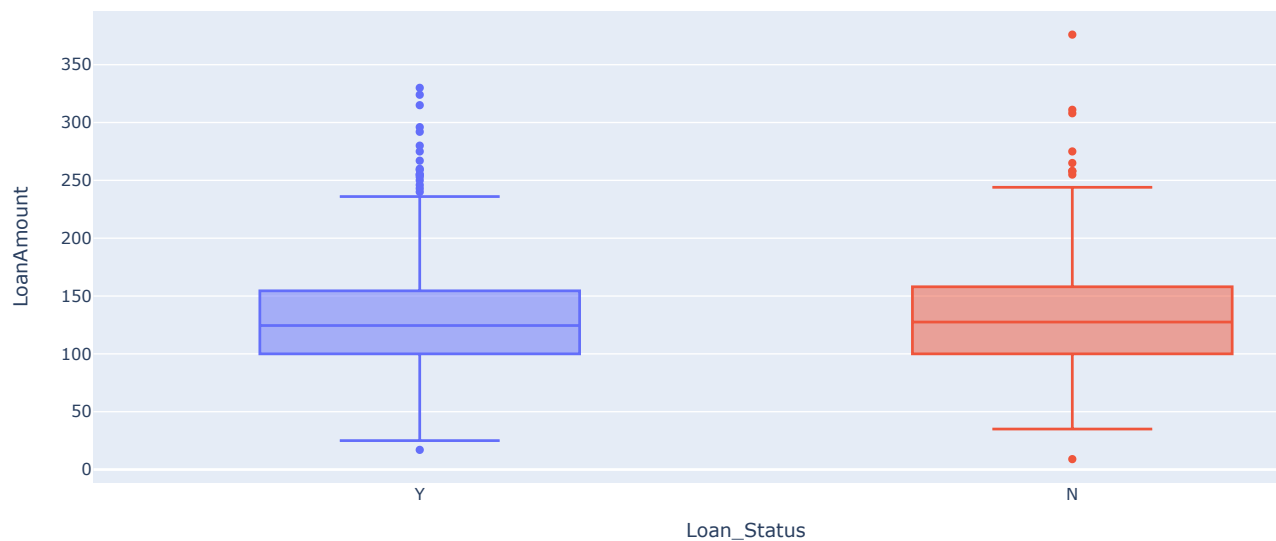
Now let's have a look at the relationship between the loan amount and the loan status

```python
fig_loan_amount = px.box(df, x='Loan_Status',
                         y='LoanAmount',
                         color="Loan_Status",
                         title='Loan_Status vs LoanAmount')
fig_loan_amount.show()
```
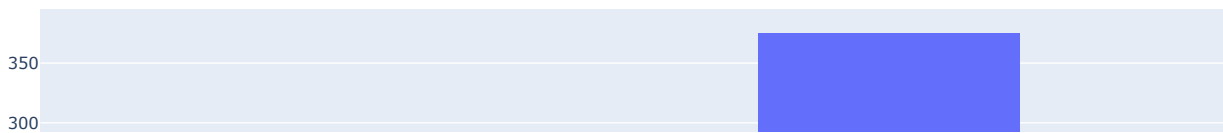
## Loan_Status vs LoanAmount



Now let's have a look at the relationship between credit history and loan status:

```python
fig_credit_history = px.histogram(df, x='Credit_History', color='Loan_Status',
                                  barmode='group',
                                  title='Loan_Status vs Credit_His')
fig_credit_history.show()
```
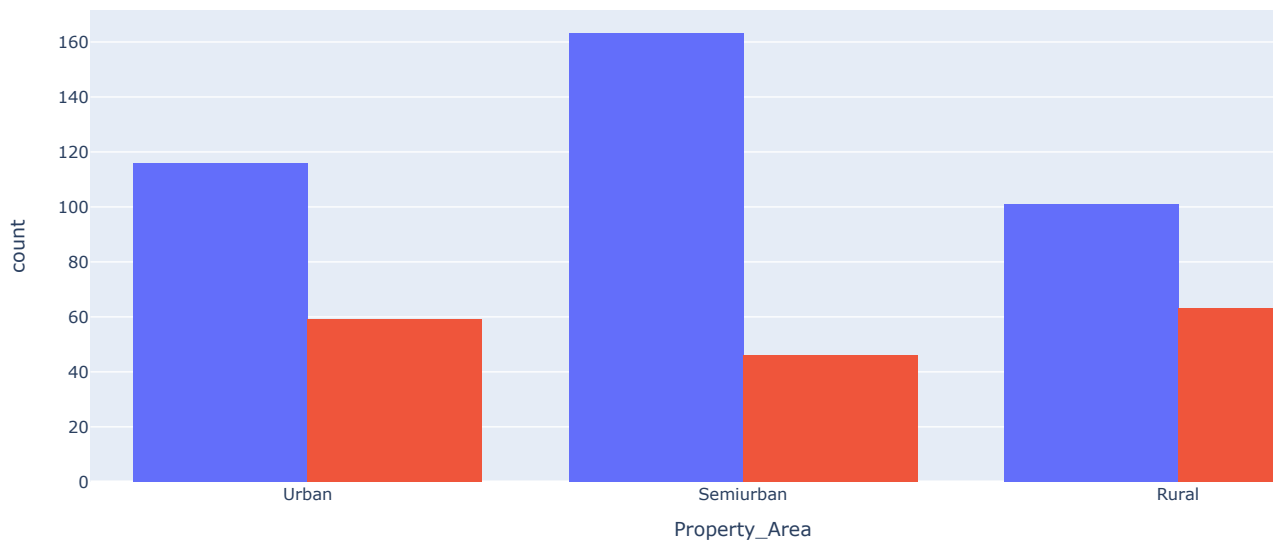
Now let's have a look at the relationship between the property area and the loan status:

```
fig_property_area = px.histogram(df, x='Property_Area', color='Loan_Status',
                                 barmode='group',
                                 title='Loan_Status vs Property_Area')
fig_property_area.show()
```



**Data Preparation and Training Loan Approval Prediction Model**

In this step, we will:

convert categorical columns into numerical ones;

split the data into training and test sets;

scale the numerical features;

train the loan approval prediction model.

```
# Convert categorical columns to numerical using one-hot encoding
cat_cols = ['Gender', 'Married', 'Dependents', 'Education', 'Self_Employed', 'Property_Area']
df = pd.get_dummies(df, columns=cat_cols)

# Split the dataset into features (X) and target (y)
X = df.drop('Loan_Status', axis=1)
y = df['Loan_Status']

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Scale the numerical columns using StandardScaler
scaler = StandardScaler()
numerical_cols = ['ApplicantIncome', 'CoapplicantIncome', 'LoanAmount', 'Loan_Amount_Term', 'Credit_History']
X_train[numerical_cols] = scaler.fit_transform(X_train[numerical_cols])
X_test[numerical_cols] = scaler.transform(X_test[numerical_cols])

from sklearn.svm import SVC
model = SVC(random_state=42)
model.fit(X_train, y_train)
```

Now let's make predictions on the test set:

```
y_pred = model.predict(X_test)
print(y_pred)
```

```
['Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y'
 'Y' 'N' 'Y' 'Y' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y'
 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'N' 'N' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y'
 'Y' 'N' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y' 'Y'
 'Y' 'Y' 'Y' 'N' 'Y' 'N' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'Y'
 'Y' 'N' 'Y' 'Y' 'N' 'Y' 'Y' 'N' 'Y' 'Y' 'Y' 'Y' 'N' 'Y' 'Y' 'N' 'Y' 'Y'
 'Y' 'Y']
```

Now let's add the predicted loan approval values to the testing set as a new column in a DataFrame called X_test_df and show the predicted loan approval values alongside the original features:

```
# Convert X_test to a DataFrame
X_test_df = pd.DataFrame(X_test, columns=X_test.columns)

# Add the predicted values to X_test_df
X_test_df['Loan_Status_Predicted'] = y_pred
print(X_test_df.head())
```

```
     ApplicantIncome  CoapplicantIncome  LoanAmount  Loan_Amount_Term  \
277        -0.544528          -0.037922   -0.983772          0.305159
84         -0.067325          -0.931554   -1.571353         -1.430680
275        -0.734870           0.334654   -0.298262          0.305159
392        -0.824919           0.522317   -0.200332          0.305159
537        -0.267373          -0.931554   -0.454950          0.305159

     Credit_History  Gender_Female  Gender_Male  Married_No  Married_Yes  \
277        0.402248              0            1           0            1
84         0.402248              0            1           0            1
275        0.402248              0            1           0            1
392        0.402248              0            1           0            1
537        0.402248              0            1           1            0

     Dependents_0  ...  Dependents_2  Dependents_3+  Education_Graduate  \
277             1  ...             0              0                   1
84              0  ...             0              0                   1
275             0  ...             0              0                   1
392             1  ...             0              0                   1
537             0  ...             1              0                   1

     Education_Not Graduate  Self_Employed_No  Self_Employed_Yes  \
277                       0                 1                  0
84                        0                 1                  0
275                       0                 1                  0
392                       0                 1                  0
537                       0                 1                  0

     Property_Area_Rural  Property_Area_Semiurban  Property_Area_Urban  \
277                    0                        0                    1
84                     0                        0                    1
275                    0                        1                    0
392                    0                        0                    1
537                    0                        1                    0

     Loan_Status_Predicted
277                      Y
84                       Y
275                      Y
392                      Y
537                      Y

[5 rows x 21 columns]
```

So this is how you can train a Machine Learning model to predict loan approval using Python.

## ▾ Summary

Loan approval prediction involves the analysis of various factors, such as the applicant's financial history, income, credit rating, employment status, and other relevant attributes. By leveraging historical loan data and applying machine learning algorithms, businesses can build models to determine loan approvals for new applicants. I hope you liked this article on Loan Approval Prediction with Machine Learning using Python. Feel free to ask valuable questions in the comments section below.

Double-click (or enter) to edit

0s    completed at 3:55 PM