

Smart Lender - Applicant Credibility Prediction For Loan Approval Using Machine Learning

Category: Machine Learning

Project Description:

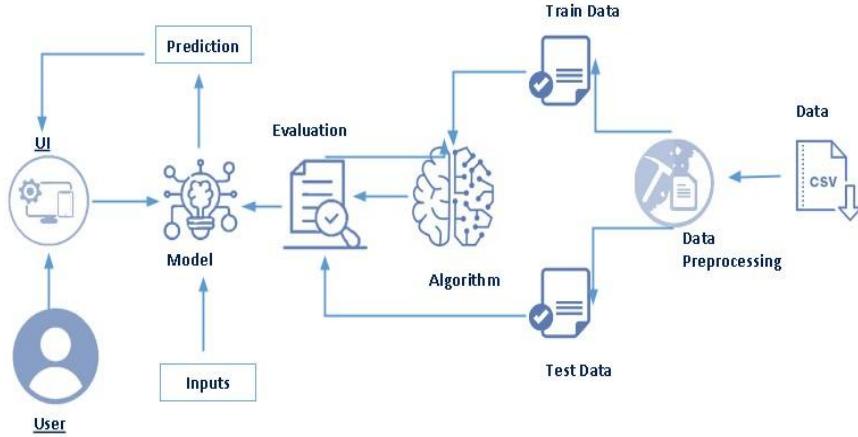
introduction:

One of the most important factors which affect our country's economy and financial condition is the credit system governed by the banks. The process of bank credit risk evaluation is recognized at banks across the globe. "As we know credit risk evaluation is very crucial, there is a variety of techniques are used for risk level calculation. In addition, credit risk is one of the main functions of the banking community.

The prediction of credit defaulters is one of the difficult tasks for any bank. But by forecasting the loan defaulters, the banks definitely may reduce their loss by reducing their non-profit assets, so that recovery of approved loans can take place without any loss and it can play as the contributing parameter of the bank statement. This makes the study of this loan approval prediction important. Machine Learning techniques are very crucial and useful in the prediction of these types of data.

We will be using classification algorithms such as Decision tree, Random forest, KNN, and xgboost. We will train and test the data with these algorithms. From this best model is selected and saved in pkl format. We will be doing flask integration and IBM deployment.

Technical Architecture:



Pre-Requisites:

To complete this project, you must require the following software, concepts, and packages

Python packages:

- o Open anaconda prompt as administrator
- o Type "pip install numpy" and click enter.
- o Type "pip install pandas" and click enter.
- o Type "pip install scikit-learn" and click enter.
- o Type "pip install matplotlib" and click enter.
- o Type "pip install pickle-mixin" and click enter.
- o Type "pip install seaborn" and click enter.
- o Type "pip install Flask" and click enter.

Prior Knowledge

You must have prior knowledge of the following topics to complete this project.

- ML Concepts

Supervised learning

Project Objectives

Write what are all the technical aspects that students would get if they complete this project.

1. Knowledge of Machine Learning Algorithms.
2. Knowledge of Python Language with Machine Learning
3. You'll be able to understand the problem to classify if it is a regression or a classification kind of problem.
4. You will be able to know how to pre-process/clean the data using different data pre-processing techniques.
5. Applying different algorithms according to the dataset and based on visualization.
6. Real-Time Analysis of Project
7. Building ease of User Interface (UI)
8. Navigation of ideas towards other projects(creativity)
9. Knowledge of building ML models.
10. How to build web applications using the Flask framework.

Project Flow

Install Required Libraries.

Data Collection.

- Collect the dataset or Create the dataset

Data Preprocessing.

- Import the Libraries.
- Importing the dataset.
- Understanding Data Type and Summary of features.
- Take care of missing data
- Data Visualization.
- Drop the column from DataFrame & replace the missing value.

- Splitting the Dataset into Dependent and Independent variables
- Splitting Data into Train and Test.

Model Building

- Training and testing the model
- Evaluation of Model
- Saving the Model

Application Building

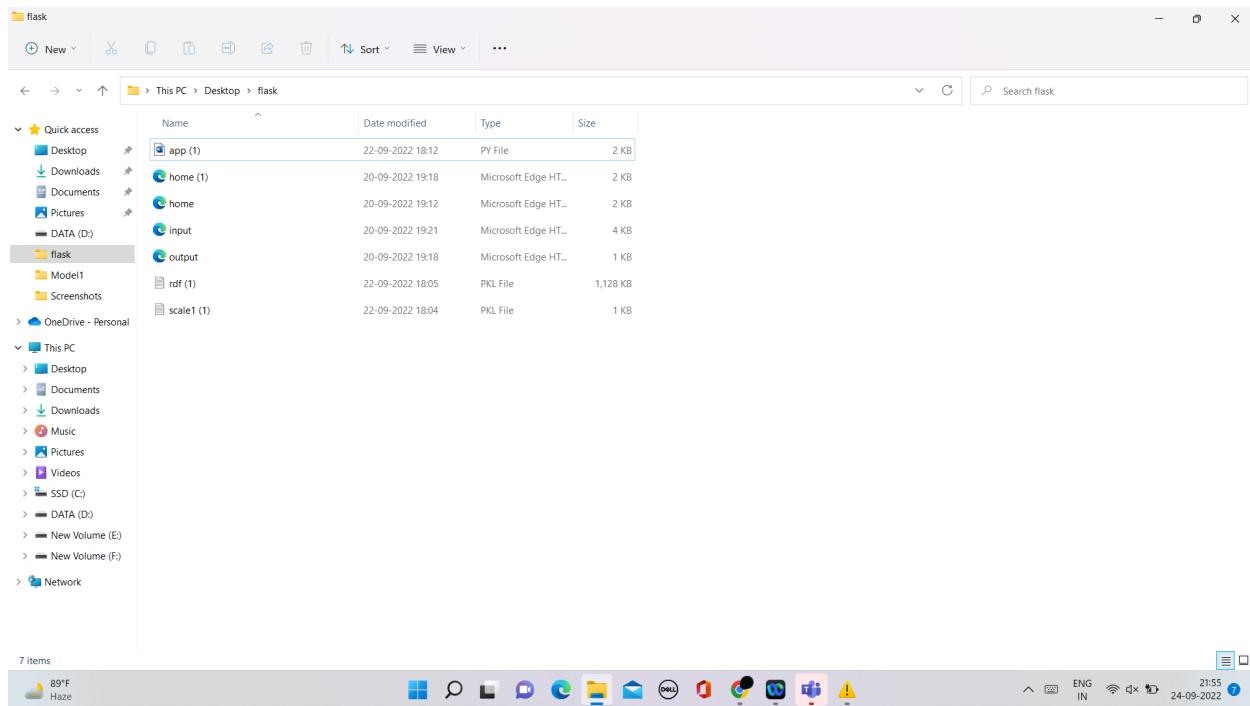
- Create an HTML file
- Build a Python Code

Final UI

- Dashboard Of the flask app.

Project Structure

Create a Project folder that contains files as shown below



- We are building a Flask Application that needs HTML pages "home.html", "index.html" stored in the templates folder and a python script app.py for server-side scripting
- The model is built in the notebook floods.ipynb
- We need the model which is saved and the saved model in this content is floods.save and transform.save
- The templates mainly used here are "home.html", "chance.html", "no chance.html",

- “index.html” for showcasing the UI
- The flask app is denoted as app.py
- IBM scoring endpoint consists of templates and app.py

Data Collection

ML depends heavily on data, without data, it is impossible for an “AI” to learn. It is the most crucial aspect that makes algorithm training possible. In Machine Learning projects, we need a training data set. It is the actual data set used to train the model for performing various actions.

Download The Dataset:

Download the dataset from the below link.

- You can collect datasets from different open sources like kaggle.com, data.gov, UCI machine learning repository, etc.
- Here we are using a data set which you can find in the below link and you can download it from the link: [Link](#)

Visualizing And Analyzing The Data

As the dataset is downloaded. Let us read and understand the data properly with the help of some visualization techniques and some analyzing techniques.

Note: There is n number of techniques for understanding the data. But here we have used some of it. In an additional way, you can use multiple techniques.

Importing The Libraries

Import the necessary libraries as shown in the image

Import the required libraries for the model to run. The first step is usually importing the libraries that will be needed in the program.

Import the necessary libraries as shown in the image

Import the required libraries for the model to run. The first step is usually importing the libraries that will be needed in the program.

- **Numpy**- It is an open-source numerical Python library. It contains a multi-dimensional array and matrix data structures. It can be used to perform mathematical operations on arrays such as trigonometric, statistical, and algebraic routines.
- **Pandas**- It is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool, built on top of the Python programming language.
- **Seaborn**- Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
- **Matplotlib**- Visualisation with python. It is a comprehensive library for creating static, animated, and interactive visualizations in Python

Sklearn – which contains all the modules required for model building

```
In [1]: import pandas as pd  
import numpy as np  
import pickle  
import matplotlib.pyplot as plt  
%matplotlib inline  
import seaborn as sns  
import sklearn  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.ensemble import GradientBoostingClassifier, RandomForestClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.model_selection import RandomizedSearchCV  
import imblearn  
from sklearn.model_selection import train_test_split  
from sklearn.preprocessing import StandardScaler  
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, f1_score  
import os
```

Importing the dataset

Reading The Dataset

- Our dataset format might be in .csv, excel files, .txt, .json, etc. We can read the

dataset with the help of pandas.

In pandas, we have a function called `read_csv()` to read the dataset. As a parameter, we have to give the directory of the CSV file

The screenshot shows a Jupyter Notebook interface with three tabs open. The active tab is titled "Importing the dataset". In the code editor, two cells are visible:

```
# importing the dataset which is in csv file
data = pd.read_csv(r"C:\Users\PRAVEEN\Downloads\loan_prediction.csv")
data
```

This cell displays a DataFrame with 614 rows and 13 columns. The columns are: Loan_ID, Gender, Married, Dependents, Education, Self_Employed, ApplicantIncome, CoapplicantIncome, LoanAmount, L, ... (with ellipses). The data includes various demographic and financial information for loan applicants.

```
# dropping the loan id columns because there is no use it for the model building
data.drop(['Loan_ID'], axis=1, inplace=True)
```

The system tray at the bottom shows the date as 24-09-2022, time as 22:08, and battery level as 89%.

Importing The Dataset

- You might have your data in .csv files, .excel files
- Let's load a .csv data file into pandas using `read_csv()` function. We will need to locate the directory of the CSV file at first (it's more efficient to keep the dataset in the same directory as your program).
- If your dataset is in some other location, Then
- **Data=pd.read_csv(r"File_location/datasetname.csv")**
- If the dataset is in the same directory of your program, you can directly read it, without giving raw as r.
- Our Dataset weatherAus.csv contains the following Columns

- Holiday - working day or holiday
- Temp- temperature of the day
- Rain and snow – whether it is raining or snowing on that day or not
- Weather = describes the weather conditions of the day
- Date and time = represents the exact date and time of the day
- Traffic volume – output column

The output column to be predicted is Traffic volume. Based on the input variables we predict the volume of the traffic. The predicted output gives them a fair idea of the count of traffic

Analyse The Data

head() method is used to return top n (5 by default) rows of a DataFrame or series.

Uni-Variate Analysis

In simple words, univariate analysis is understanding the data with single feature. Here we have displayed two different graphs such as distplot and countplot.

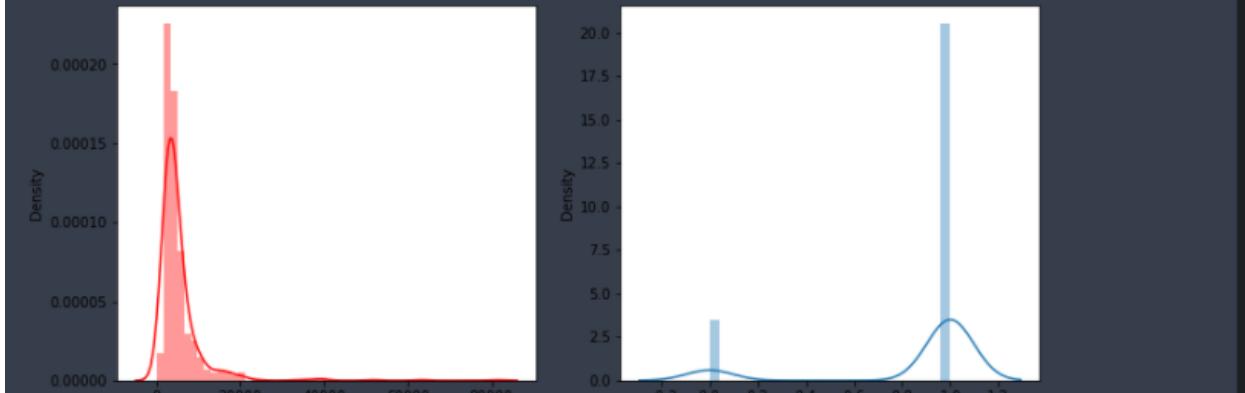
- Seaborn package provides a wonderful function distplot. With the help of distplot, we can find the distribution of the feature. To make multiple graphs in a single plot, we use a subplot.

```

#plotting the using distplot
plt.figure(figsize=(12,5))
plt.subplot(121)
sns.distplot(data['ApplicantIncome'], color='r')
plt.subplot(122)
sns.distplot(data['Credit_History'])
plt.show()

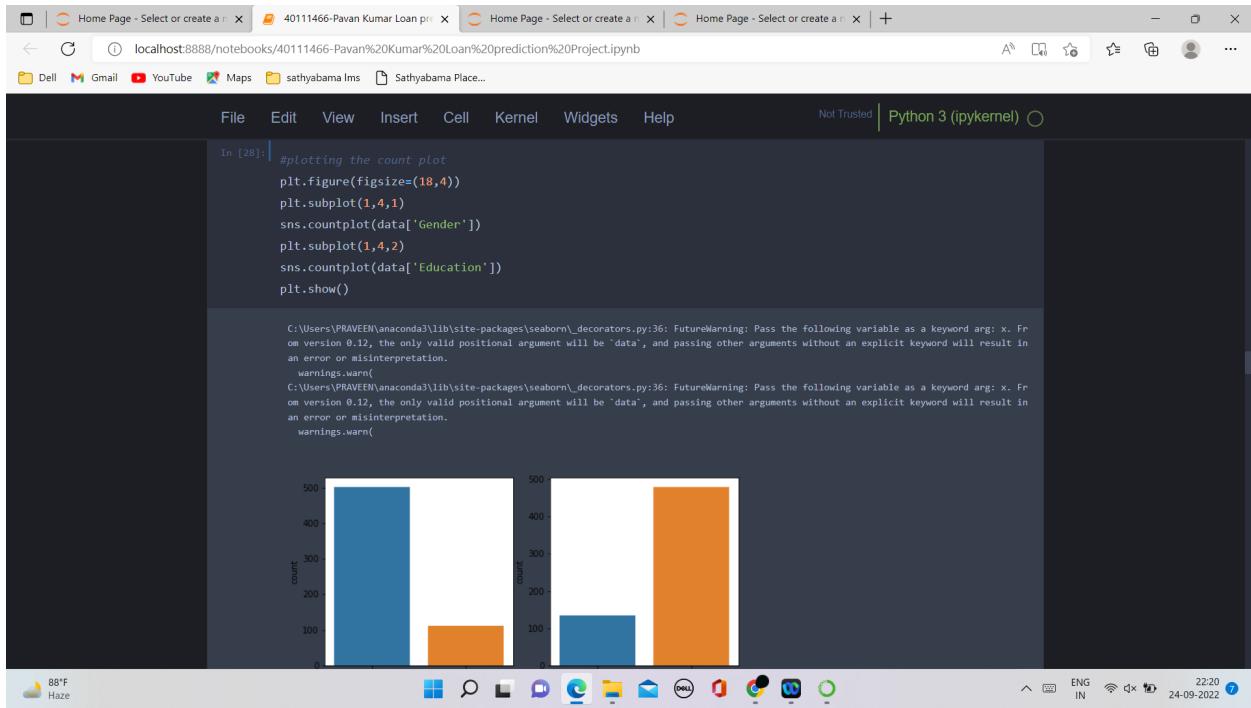
C:\Users\HP\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
C:\Users\HP\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

```



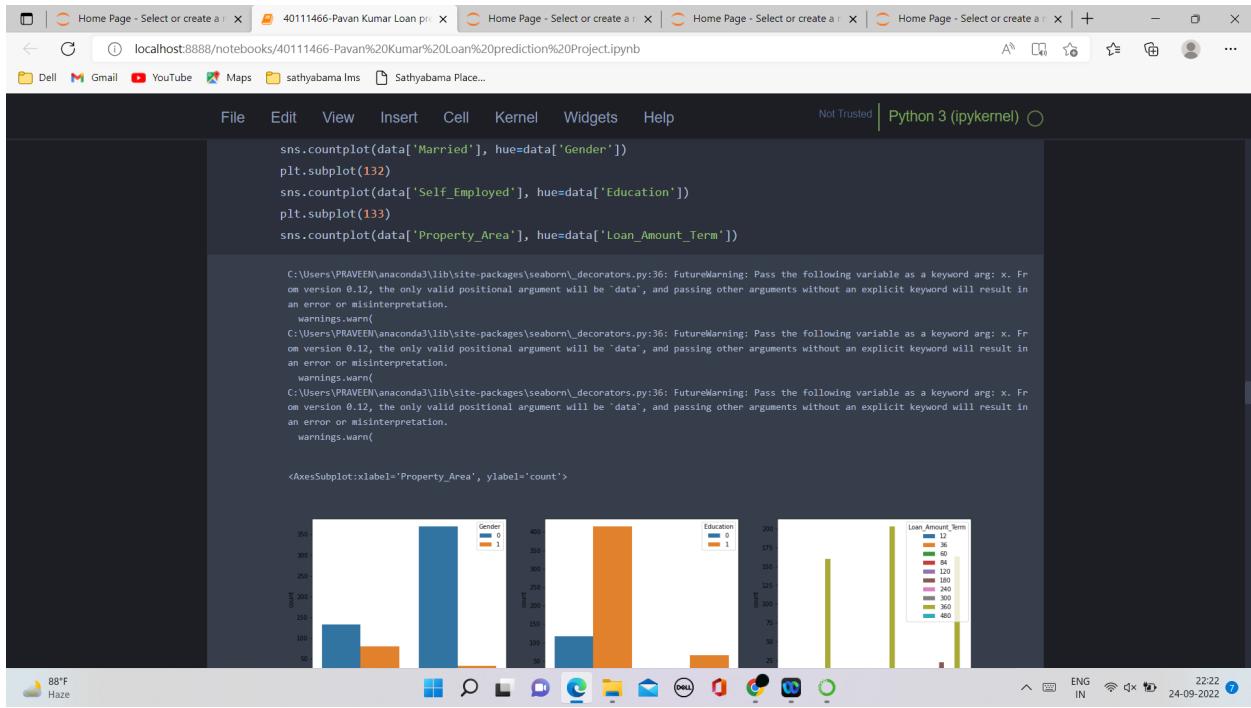
- In our dataset, we have some categorical features. With the count plot function, we are going to count the unique category in those features. We have created a dummy data frame with categorical features. With for loop and subplot, we have plotted the below graph.
- From the plot we came to know, Applicants' income is skewed towards the left side, whereas credit history is categorical with 1.0 and 0.0

Bivariate Analysis



Countplot:-

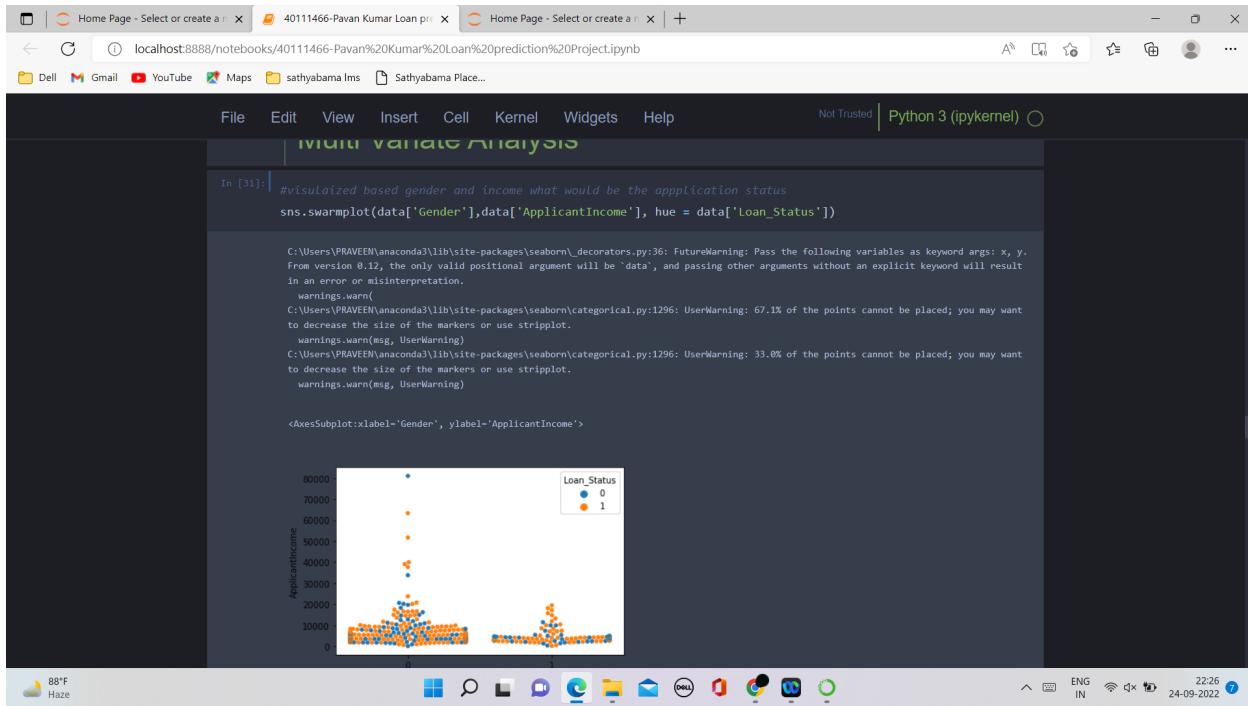
A count plot can be thought of as a histogram across a categorical, instead of a quantitative, variable. The basic API and options are identical to those for barplot() , so you can compare counts across nested variables.



From the above graph, we can infer the analysis such as

- Segmenting the gender column and married column based on bar graphs
- Segmenting the Education and Self-employed based on bar graphs, for drawing insights such as educated people are employed.
- The loan amount term is based on the property area of a person holding

Multivariate Analysis



From the above graph we are plotting the relationship between the Gender, applicants income and loan status of the person

Descriptive Analysis

Descriptive analysis is to study the basic features of data with the statistical process. Here pandas have a worthy function called describe. With this describe function we can understand the unique, top, and frequent values of categorical features. And we can find mean, std, min, max and percentile values of continuous features.

	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	12.00000	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

Data Pre-Processing

As we have understood how the data is let's pre-process the collected data.

The download data set is not suitable for training the machine learning model as it might have so much of randomness so we need to clean the dataset properly in order to fetch good results. This activity includes the following steps.

- Handling missing values
- Handling categorical data
- Handling outliers
- Scaling Techniques
- Splitting dataset into training and test set

Checking For Null Values

- Let's find the shape of our dataset first, To find the shape of our data, df.shape method is used. To find the data type, df.info() function is used.

```

In [24]: #getting bthe total info of the data after perfroming categorical to numericsal and rplacing missng value
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Gender       614 non-null    float64
 1   Married      614 non-null    float64
 2   Dependents   614 non-null    object  
 3   Education    614 non-null    int64  
 4   Self_Employed 614 non-null    float64
 5   ApplicantIncome 614 non-null    int64  
 6   CoapplicantIncome 614 non-null    float64
 7   LoanAmount    614 non-null    float64
 8   Loan_Amount_Term 614 non-null    float64
 9   Credit_History 614 non-null    float64
 10  Property_Area 614 non-null    int64  
 11  Loan_Status   614 non-null    int64  
dtypes: float64(7), int64(4), object(1)
memory usage: 57.7+ KB

In [25]: #changing the datatype of each float column to int

```

- For checking the null values, `df.isnull()` function is used. To sum those null values we use `.sum()` function to it. From the below image we found that there are no null values present in our dataset. So we can skip the handling of the missing values step.

```

In [19]: data['Self_Employed'] = data['Self_Employed'].fillna(data['Self_Employed'].mode()[0])

In [20]: data['LoanAmount'] = data['LoanAmount'].fillna(data['LoanAmount'].mode()[0])

In [21]: data['Loan_Amount_Term'] = data['Loan_Amount_Term'].fillna(data['Loan_Amount_Term'].mode()[0])

In [22]: data['Credit_History'] = data['Credit_History'].fillna(data['Credit_History'].mode()[0])

In [23]: data.isnull().sum()

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

In [24]: #getting bthe total info of the data after perfroming categorical to numericsal and rplacing missng value
data.info()

```

From the above code of analysis, we can infer that columns such as gender, married, dependents, self-employed, loan amount, loan amount tern, and credit history are having the missing values, we

need to treat them in a required way.

The screenshot shows a Jupyter Notebook window with multiple tabs at the top, all titled "Home Page - Select". The main area displays a series of code cells (In [15] through In [23]) used to handle missing values in a dataset named "data". The code includes:

```
In [15]: #replacing + with space for filling the nan values
data['Dependents']=data['Dependents'].str.replace('+','')

C:\Users\PRAVEEN\AppData\Local\Temp\ipykernel_3764\3624576495.py:2: FutureWarning: The default value of regex will change from True to False in a future version. In addition, single character regular expressions will "not" be treated as literal strings when regex=True.
data['Dependents']=data['Dependents'].str.replace('+','')

In [16]: data['Gender'] = data['Gender'].fillna(data['Gender'].mode()[0])

In [17]: data['Married'] = data['Married'].fillna(data['Married'].mode()[0])

In [18]: data['Dependents'] = data['Dependents'].fillna(data['Dependents'].mode()[0])

In [19]: data['Self_Employed'] = data['Self_Employed'].fillna(data['Self_Employed'].mode()[0])

In [20]: data['LoanAmount'] = data['LoanAmount'].fillna(data['LoanAmount'].mode()[0])

In [21]: data['Loan_Amount_Term'] = data['Loan_Amount_Term'].fillna(data['Loan_Amount_Term'].mode()[0])

In [22]: data['Credit_History'] = data['Credit_History'].fillna(data['Credit_History'].mode()[0])

In [23]: data.isnull().sum()
```

The notebook interface includes a toolbar with various icons, a status bar at the bottom showing system information like battery level and date/time, and a menu bar with File, Edit, View, Insert, Cell, Kernel, Widgets, Help, and Python 3 (ipykernel).

We will fill the missing values in numeric data type using the mean value of that particular column and categorical data type using the most repeated value.

Handling Categorical Values

As we can see our dataset has categorical data we must convert the categorical data to integer encoding or binary encoding.

To convert the categorical features into numerical features we use encoding techniques. There are several techniques but in our project, we are using manual encoding with the help of list comprehension.

- In our project, Gender , married,dependents,self-employed,co-applicants income, loan amount , loan amount term, credit history With list comprehension encoding is done.

```

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) ○
Handling categorical Values
In [1]: import jupyterthemes as jt
In [2]: jt --theme onedork
In [12]: data['Gender']=data['Gender'].map({'Female':1,'Male':0})
data['Property_Area']=data['Property_Area'].map({'Urban':2,'Semiurban': 1,'Rural':0})
data['Married']=data['Married'].map({'Yes':1,'No':0})
data['Education']=data['Education'].map({'Graduate':1,'Not Graduate':0})
data['Self_Employed']=data['Self_Employed'].map({'Yes':1,'No':0})
data['Loan_Status']=data['Loan_Status'].map({'Y':1,'N':0})
In [13]: #Handling categorical feature Gender
data.head()

```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoaapplicantIncome	LoanAmount	Loan_Amount
0	0.0	0.0	0	1	0.0	5849	0.0	NaN	360.0
1	0.0	1.0	1	1	0.0	4583	1508.0	128.0	360.0
2	0.0	1.0	0	1	1.0	3000	0.0	66.0	360.0
3	0.0	1.0	0	0	0.0	2583	2358.0	120.0	360.0

converting string datatype variables into integer data type

```

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) ○
In [25]: #changing the datatype of each float column to int
data['Gender']=data['Gender'].astype('int64')
data['Married']=data['Married'].astype('int64')
data['Dependents']=data['Dependents'].astype('int64')
data['Self_Employed']=data['Self_Employed'].astype('int64')
data['CoapplicantIncome']=data['CoapplicantIncome'].astype('int64')
data['LoanAmount']=data['LoanAmount'].astype('int64')
data['Loan_Amount_Term']=data['Loan_Amount_Term'].astype('int64')
data['Credit_History']=data['Credit_History'].astype('int64')
In [26]: data.info()

```

Column	Non-Null Count	Dtype
0 Gender	614 non-null	int64
1 Married	614 non-null	int64
2 Dependents	614 non-null	int64
3 Education	614 non-null	int64
4 Self_Employed	614 non-null	int64
5 ApplicantIncome	614 non-null	int64
6 CoapplicantIncome	614 non-null	int64

Balancing The Dataset

Data Balancing is one of the most important step, which need to be performed for classification models, because when we train our model on imbalanced dataset ,we will get biased results, which means our model is able to predict only one class element

For Balancing the data we are using SMOTE Method.

SMOTE: Synthetic minority over sampling technique, which will create new synthetic data points for under class as per the requirements given by us using KNN method.

The screenshot shows a Jupyter Notebook interface with multiple tabs open. The active tab is titled "Balancing the Dataset". The code in the notebook is as follows:

```
In [32]: #Balancing the dataset by using smote
from imblearn.combine import SMOTETomek

In [33]: smote = SMOTETomek(0.90)

C:\Users\PRAVEEN\anaconda3\lib\site-packages\imblearn\utils\_validation.py:586: FutureWarning: Pass sampling_strategy=0.9 as keyword args. From version 0.9 passing these as positional arguments will result in an error
warnings.warn(
In [34]: #dividing the dataset into dependent and independent y and x respectively
y = data['Loan_Status']
x = data.drop(columns=['loan_Status'],axis=1)

In [35]: #shape of x after seperating from the total data set
x.shape

(614, 11)

In [36]: #shape of y
y.shape

(614,)
```

The system tray at the bottom shows weather information (87°F Haze), taskbar icons (File Explorer, Task View, Start, etc.), and system status (ENG IN, 22:49, 24-09-2022).

The screenshot shows a Jupyter Notebook interface with several tabs open. The active tab is titled 'YARRA PAVAN KUMAR - LOAN P...' and contains Python code for data preprocessing. The code includes:

```
In [35]: #shape of x after separating from the total data set  
x.shape  
  
(614, 11)  
  
In [36]: #shape of y  
y.shape  
  
(614,)  
  
In [37]: #creating a new x and y variables for the balanced set  
x_bal,y_bal = smote.fit_resample(x,y)  
  
In [38]: #printing the values of y before balancing the data and after  
print(y.value_counts())  
print(y_bal.value_counts())  
  
1    422  
0    192  
Name: Loan_Status, dtype: int64  
1    360  
0    317  
Name: Loan_Status, dtype: int64  
  
In [39]: names = x_bal.columns
```

The notebook interface includes a toolbar at the top with File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Not Trusted, and Python 3 (ipykernel). Below the code cells, there's a toolbar with various icons (Windows Start, Search, Task View, File, Mail, Dell logo, Google Sheets, Google Slides, Google Photos) and system status indicators (87°F Haze, ENG IN, 24-09-2022, 22:49).

From the above picture, we can infer that previously our dataset is having 492 class 1, and 192 class items, after applying smote technique on the dataset the size has been changed for minority class.

Scaling The Data

Scaling is one the important process, we have to perform on the dataset, because of data measures in different ranges can leads to mislead in prediction

Models such as KNN, Logistic regression need scaled data, as they follow distance based method and Gradient Descent concept.

The screenshot shows a Jupyter Notebook window titled "Scaling the Dataset". The code in the notebook performs the following steps:

- # performing feature Scaling operation using standard scaler on X part of the dataset because there different type of values in the columns
- sc=StandardScaler()
- x_bal=sc.fit_transform(x_bal)
- x_bal = pd.DataFrame(x_bal,columns=names)
- #splitting the dataset in train and test on balanced dataset
- X_train, X_test, y_train, y_test = train_test_split(x_bal, y_bal, test_size=0.33, random_state=42)
- X_train.shape
- (453, 11)
- X_test.shape
- (224, 11)
- y_train.shape, y_test.shape
- ((453,), (224,))

The notebook is running in Python 3 (ipykernel) and is set to "Not Trusted". The system tray at the bottom shows weather information (87°F Haze), system icons, and a battery level of 22:51.

We will perform scaling only on the input values

Once the dataset is scaled, it will be converted into array and we need to convert it back to dataframe.

Splitting Data Into Train And Test

Now let's split the Dataset into train and test sets

Changes: first split the dataset into x and y and then split the data set

Here x and y variables are created. On the x variable, df is passed by dropping the target variable. And on y target variable is passed. For splitting training and testing data, we are using the `train_test_split()` function from sklearn. As parameters, we are passing x, y, test_size, and random_state.

```
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) ○
sc=StandardScaler()
x_bal=sc.fit_transform(x_bal)

In [41]: x_bal = pd.DataFrame(x_bal,columns=names)

In [42]: #splitting the dataset in train and test on balanced dataset
        X_train, X_test, y_train, y_test = train_test_split(
            x_bal, y_bal, test_size=0.33, random_state=42)

In [43]: X_train.shape

(453, 11)

In [44]: X_test.shape

(224, 11)

In [45]: y_train.shape, y_test.shape

((453,), (224,))

Model building
```

Model Building

Now our data is cleaned and it's time to build the model. We can train our data on different algorithms. for this project, we are applying four classification algorithms. The best model is saved based on its performance.

Decision Tree Model

A function named decision tree is created and train and test data are passed as the parameters. Inside the function, the DecisionTreeClassifier algorithm is initialized and training data is passed to the model with .fit() function. Test data is predicted with the .predict() function and saved in the new variable. For evaluating the model, a confusion

The screenshot shows a Jupyter Notebook interface with several cells of Python code. The first cell (In [42]) contains code to fit a Random Forest model and print accuracy scores. The output shows a score of 1.0 and 0.8609865470852018. The second cell (In [43]) imports a DecisionTreeClassifier and fits it to the same data. The third cell (In [44]) prints the accuracy scores again, showing 1.0 and 0.76694543027619. The notebook is titled 'dileep smartlender project' and is run in Python 3 (ipykernel).

```
model.fit(X_train,y_train)
y_tr = model.predict(X_train)
print(accuracy_score(y_tr,y_train))
yPred = model.predict(X_test)
print(accuracy_score(yPred,y_test))

#printing the train accuracy and test accuracy respectively
RandomForest(X_train,X_test,y_train,y_test)

1.0
0.8609865470852018

#importing and building the Decision tree model
def decisionTree(X_train,X_test,y_train,y_test):
    model = DecisionTreeClassifier()
    model.fit(X_train,y_train)
    y_tr = model.predict(X_train)
    print(accuracy_score(y_tr,y_train))
    yPred = model.predict(X_test)
    print(accuracy_score(yPred,y_test))

#printing the train accuracy and test accuracy respectively
decisionTree(X_train,X_test,y_train,y_test)

1.0
0.76694543027619
```

Random Forest Model

A function named randomForest is created and train and test data are passed as the parameters. Inside the function, the RandomForestClassifier algorithm is initialized and training data is passed to the model with .fit() function. Test data is predicted with .predict() function and saved in a new variable. For evaluating the model, a confusion matrix and classification report are done.

The screenshot shows a Jupyter Notebook interface with several cells of Python code. The first cell (In [44]) prints the shape of X_test, which is (224, 11). The second cell (In [45]) prints the shapes of y_train and y_test, both of which are (453, 11). The third cell (In [46]) imports the RandomForestClassifier and defines a function to build the model. The fourth cell (In [47]) prints the accuracy scores for the train and test sets. The notebook is titled 'dileep smartlender project' and is run in Python 3 (ipykernel).

```
(453, 11)

In [44]: X_test.shape

(224, 11)

In [45]: y_train.shape, y_test.shape

((453,), (224,))

Model building

In [46]: #importing and building the random forest model
def RandomForest(X_train,X_test,y_train,y_test):
    model = RandomForestClassifier()
    model.fit(X_train,y_train)
    y_tr = model.predict(X_train)
    print(accuracy_score(y_tr,y_train))
    yPred = model.predict(X_test)
    print(accuracy_score(yPred,y_test))

In [47]: #printing the train accuracy and test accuracy respectively
RandomForest(X_train,X_test,y_train,y_test)
```

KNN Model

A function named KNN is created and train and test data are passed as the parameters. Inside the function, the KNeighborsClassifier algorithm is initialized and training data is passed to the model with .fit() function. Test data is predicted with .predict() function and saved in a new variable. For evaluating the model, a confusion matrix and classification report is done.

The screenshot shows a Jupyter Notebook interface with three code cells:

- In [48]:**

```
#importing and building the KNN model
def KNN(X_train,X_test,y_train,y_test):
    model = KNeighborsClassifier()
    model.fit(X_train,y_train)
    y_tr = model.predict(X_train)
    print(accuracy_score(y_tr,y_train))
    yPred = model.predict(X_test)
    print(accuracy_score(yPred,y_test))
```
- In [49]:**

```
#printing the train accuracy and test accuracy respectively
KNN(X_train,X_test,y_train,y_test)
```

Output:
0.8454746136865342
0.7633928571428571
- In [50]:**

```
#importing and building the xg boost model
def XGB(X_train,X_test,y_train,y_test):
    model = GradientBoostingClassifier()
    model.fit(X_train,y_train)
    y_tr = model.predict(X_train)
    print(accuracy_score(y_tr,y_train))
    yPred = model.predict(X_test)
```

The notebook is running in Python 3 (ipykernel). The system tray shows battery level (87%), Haze, and a timestamp of 23:07 24-09-2022.

Xgboost Model

A function named xgboost is created and train and test data are passed as the parameters. Inside the function, the GradientBoostingClassifier algorithm is initialized and training data is passed to the model with .fit() function. Test data is predicted with .predict() function and saved in a new variable. For evaluating the model, a confusion matrix and classification report are done.

The screenshot shows a Jupyter Notebook window with several cells of Python code. Cell [50] imports XGBoost and defines a function to build a model. Cell [51] prints accuracy scores for training and testing datasets. Cell [52] initializes a Random Forest classifier. Cell [53] prints parameters for randomized search CV.

```
0.8454746136865342  
0.7633928571428571  
In [50]: #importing and building the Xg boost model  
def XGB(X_train,X_test,y_train,y_test):  
    model = GradientBoostingClassifier()  
    model.fit(X_train,y_train)  
    y_tr = model.predict(X_train)  
    print(accuracy_score(y_tr,y_train))  
    yPred = model.predict(X_test)  
    print(accuracy_score(yPred,y_test))  
  
In [51]: #printing the train accuracy and test accuracy respectively  
XGB(X_train,X_test,y_train,y_test)  
  
0.9381898454746137  
0.7991071428571429  
  
Hyper parameter Tuning  
In [52]: rf = RandomForestClassifier()  
In [53]: # giving some parameters that can be used in randomized search cv
```

Compare The Model

For comparing the above four models compareModel function is defined.

The screenshot shows a Jupyter Notebook window with the same sequence of cells as the first one, demonstrating the execution of the compareModel function.

```
0.8454746136865342  
0.7633928571428571  
In [50]: #importing and building the Xg boost model  
def XGB(X_train,X_test,y_train,y_test):  
    model = GradientBoostingClassifier()  
    model.fit(X_train,y_train)  
    y_tr = model.predict(X_train)  
    print(accuracy_score(y_tr,y_train))  
    yPred = model.predict(X_test)  
    print(accuracy_score(yPred,y_test))  
  
In [51]: #printing the train accuracy and test accuracy respectively  
XGB(X_train,X_test,y_train,y_test)  
  
0.9381898454746137  
0.7991071428571429  
  
Hyper parameter Tuning  
In [52]: rf = RandomForestClassifier()  
In [53]: # giving some parameters that can be used in randomized search cv
```

The screenshot shows a Jupyter Notebook window with several tabs at the top, all titled "Home Page - Select or click". The active tab is "localhost:8889/notebooks/YARRA%20PAVAN%20KUMAR%20-%20LOAN%20PREDICTION%20PROJECT.ipynb". The notebook interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help) and a toolbar with icons for file operations. The code cell In [49] contains Python code for a KNN classifier:

```
model = KNeighborsClassifier()
model.fit(X_train,y_train)
y_tr = model.predict(X_train)
print(accuracy_score(y_tr,y_train))
yPred = model.predict(X_test)
print(accuracy_score(yPred,y_test))

In [49]: #printing the train accuracy and test accuracy respectively
KNN(X_train,X_test,y_train,y_test)
```

The output of this cell shows the accuracy scores:

```
0.8454746136865342
0.7633928571428571
```

The code cell In [50] contains Python code for an XGBoost classifier:

```
#importing and building the Xg boost model
def XGB(X_train,X_test,y_train,y_test):
    model = GradientBoostingClassifier()
    model.fit(X_train,y_train)
    y_tr = model.predict(X_train)
    print(accuracy_score(y_tr,y_train))
    yPred = model.predict(X_test)
    print(accuracy_score(yPred,y_test))

In [50]:
```

Evaluating Performance Of The Model And Saving The Model

From sklearn, cross_val_score is used to evaluate the score of the model. On the parameters, we have given rf (model name), x, y, cv (as 5 folds). Our model is performing well. So, we are saving the model by pickle.dump().

Note: To understand cross-validation, refer this [link](#)

```
from sklearn.model_selection import cross_val_score

# Random forest model is selected

rf = RandomForestClassifier()
rf.fit(x_train,y_train)
yPred = rf.predict(x_test)

f1_score(yPred,y_test,average='weighted')
0.9679166666666668

cv = cross_val_score(rf,x,y,cv=5)

np.mean(cv)
0.985

#saving the model by using pickle function
pickle.dump(model,open('rdf.pkl','wb'))
```

Application Building

In this section, we will be building a web application that is integrated to the model we built. A UI is provided for the user where he has to enter the values for predictions. The entered values are given to the saved model and prediction is showcased on the UI.

This section has the following tasks

- Building HTML Pages
- Building serverside script

Building Html Pages

For this project create three HTML files namely

- home.html
- predict.html
- submit.html

and save them in the templates folder.

Let's see how our home.html page looks like:

Welcome to Loan Prediction

Loan Approval is based on lot of this,rather than going to a bank and getting rejected. We made it simple that you can get your loan approval prediction by our machine learning model, for to predict we need some of your information

[Predict](#)



Now when you click on predict button from top right corner you will get redirected to predict.html

Lets look how our predict.html file looks like:

Enter your Details for Loan Approval Prediction	
Gender	<input type="text" value="Male"/>
Married	<input type="text" value="Yes"/>
Dependents	<input type="text" value="No of Dependents on you....."/>
Education	<input type="text" value="Graduate"/>
Self Employed	<input type="text" value="Yes"/>
Applicant Income	<input type="text" value="Your Income...."/>
Co Applicant Income	<input type="text" value="Your Co Applicant Income..."/>
Loan Amount	<input type="text" value="Enter the Loan Amount ..."/>
Loan Amount Term	<input type="text" value="Enter the Term Loan Amount in days ..."/>
Credit History	<input type="text" value="Enter the Your Previous Credit History ..."/>
Property Area	<input type="text" value="Urban"/>
<input type="button" value="Submit"/> <	

Now when you click on submit button from left bottom corner you will get redirected to submit.html

Lets look how our submit.html file looks like:



Build Python Code

```
y x
from flask import Flask, render_template, request
import numpy as np
import pickle
```

Load the saved model. Importing the flask module in the project is mandatory. An object of Flask class is our WSGI application. Flask constructor takes the name of the current module (`_name_`) as an argument.

```
app = Flask(__name__)
model = pickle.load(open(r'rdf.pkl', 'rb'))
scale = pickle.load(open(r'scale1.pkl', 'rb'))
```

Render HTML page:

```
@app.route('/')
def home():
    return render_template('home.html')
```

Here we will be using declared constructor to route to the HTML page which we have created earlier.

In the above example, '/' URL is bound with home.html function. Hence, when the home page of the web server is opened in browser, the html page will be rendered. Whenever you enter the values from the html page the values can be retrieved using POST Method.

Retrieves the value from UI:

```
@app.route('/submit',methods=["POST","GET"])# route to show the predictions in a web UI
def submit():
    # reading the inputs given by the user
    input_feature=[int(x) for x in request.form.values() ]
    #input_feature = np.transpose(input_feature)
    input_feature=[np.array(input_feature)]
    print(input_feature)
    names = ['Gender', 'Married', 'Dependents', 'Education', 'Self_Employed', 'ApplicantIncome',
             'CoapplicantIncome', 'LoanAmount', 'Loan_Amount_Term', 'Credit_History', 'Property_Area']
    data = pandas.DataFrame(input_feature,columns=names)
    print(data)

    #data_scaled = scale.fit_transform(data)
    #data = pandas.DataFrame(),columns=names)

    # predictions using the loaded model file
    prediction=model.predict(data)
    print(prediction)
    prediction = int(prediction)
    print(type(prediction))

    if (prediction == 0):
        return render_template("output.html",result ="Loan will Not be Approved")
    else:
        return render_template("output.html",result = "Loan will be Approved")
    # showing the prediction results in a UI
```

Here we are routing our app to predict() function. This function retrieves all the values from the HTML page using a Post request. That is stored in an array. This array is passed to the model.predict() function. This function returns the prediction. And this prediction value will be rendered to the text that we have mentioned in the submit.html page earlier.

Main Function:

```
if __name__=="__main__":
    # app.run(host='0.0.0.0', port=8000,debug=True)      # running the app
    port=int(os.environ.get('PORT',5000))
    app.run(debug=False)
```

Run The Application

Run the application

- Open the anaconda prompt from the start menu

- Navigate to the folder where your python script is.
- Now type “python app.py” command
- Navigate to the localhost where you can view your web page.

Click on the predict button from the top left corner, enter the inputs, click on the submit button, and see the result/prediction on the web

```
base) D:\TheSmartBridge\Projects\2. DrugClassification\Drug classification> python app.py
* Serving Flask app "app" (lazy loading)
* Environment: production
WARNING: This is a development server. Do not use it in a production setting.
Use a production WSGI server instead.
* Debug mode: off
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
```

