

# Loan Prediction System for Banks

## Abstract

This project delves into the development and implementation of a Loan Prediction System for banks, leveraging the power of machine learning to enhance the efficiency and accuracy of loan approval processes. The primary objective is to address the pivotal questions faced by lending institutions: evaluating the risk associated with borrowers and determining the appropriateness of extending loans to them. By constructing predictive models, this project aims to streamline decision-making, optimize resource utilization, and mitigate the financial risks inherent in the loan approval process.

## Introduction

The proposal method aims to develop a robust and accurate loan prediction system for banks utilizing machine learning algorithms. The system's core functionality involves the analysis of historical customer data, including credit history, income, loan amount, loan term, and employment status, to predict the likelihood of a loan applicant receiving approval. The primary objective is to assist banks in making informed and reliable lending decisions, thereby reducing the risk of financial losses and streamlining the loan approval process.

The successful implementation of the loan prediction system holds the potential to empower banks to make data-driven decisions, minimize risks associated with loan defaults, optimize lending processes, and ultimately enhance their overall financial performance.

In addressing the specific problem presented by Comfort Zone company, the project seeks to automate the loan eligibility process in real-time based on customer details provided during online application submissions. By identifying eligible customer segments, the company aims to target specific customers efficiently. The overarching goal of the project is to predict whether a loan application would be approved or denied.

The lending industry faces two critical questions: How risky is the borrower, and given the borrower's risk, should the loan be granted? In response to these challenges, the project leverages data science teams in banks to build predictive models using machine learning techniques. Loan Prediction, a common real-life problem for retail banks, has the potential to save significant man-hours when approached correctly.

The proposed methodology and analysis plan include key components such as data collection, preprocessing, feature selection, model training, model evaluation, and continuous improvement. These steps are designed to gather historical loan applicant data, clean and prepare the data for analysis, identify relevant features impacting loan approval, train machine learning models, evaluate model performance, and implement mechanisms for continuous improvement.

In conclusion, the implementation of this system holds the promise of reducing the risk of financial losses associated with loan defaults, optimizing the loan approval process, and enhancing the overall operational effectiveness of banks. Through continuous improvement and adaptation to evolving market dynamics, the system aims to ensure its relevance over time.

## Models

The loan prediction system proposed in this project employs several machine learning models to address the specific challenges posed by the lending industry. The models selected for this project are tailored to handle classification problems, specifically binary classification, as the primary goal is to predict whether a loan application will be approved or denied.

## Logistic Regression

Logistic Regression is a fundamental classification algorithm that is well-suited for binary outcomes. In the context of this project, it will be utilized to model the probability of loan approval based on various independent variables such as credit history, income, and loan amount.

$$P(Y = 1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

## Decision Trees

Decision Trees are powerful tools for classification tasks, providing a clear and interpretable decision-making structure. In this project, Decision Trees will be employed to analyze and classify loan applications based on features like credit history, income, and other relevant factors.

$$\text{Decision at node } j : X_i \leq t_j ?$$

## K-Means Clustering

While traditionally used for clustering, K-Means can also be adapted for classification tasks. In this project, K-Means Clustering will be explored to group similar loan applications and identify patterns that may influence the approval or denial of loans.

$$J = \sum_{i=1}^k \sum_{j=1}^{n_i} \|x_{ij} - c_i\|^2$$

## Model Selection Rationale

The choice of these models is motivated by the need for interpretable results, the ability to handle diverse types of data, and the potential for achieving high predictive accuracy. Logistic Regression provides transparency in understanding the impact of each variable, while Decision Trees offer a visual representation of decision-making. K-Means Clustering, although primarily a clustering algorithm, will be explored for its potential in grouping similar loan applications.

These models collectively contribute to the overarching objective of predicting loan approval outcomes with precision, recall, and accuracy, ultimately assisting lending institutions in making well-informed decisions while minimizing the risk of financial losses.

## Data

The dataset used in this project consists of two CSV files: **train** and **test**. These files play a crucial role in training the machine learning models and evaluating their performance. The **train** file is employed for training the models, containing both the independent variables and the target variable. On the other hand, the **test** file includes only the independent variables, and the model's task is to predict the target variable for this test data.

Dataset Link: <https://www.kaggle.com/datasets/altruistdelhite04/loan-prediction-problem-dataset/data>

## Dataset Variables

The dataset comprises the following variables, each with its associated data type:

Variable	Data Type	Description
Loan_ID	Object	Unique Loan ID
Gender	Object	Gender of the applicant
Married	Object	Applicant's marital status
Dependents	Object	Number of dependents
Education	Object	Applicant's education level
Self_Employed	Object	Indicates whether the applicant is self-employed
ApplicantIncome	Int	Income of the applicant
CoapplicantIncome	Float	Income of the coapplicant
LoanAmount	Float	Required loan amount in thousands
Loan_Amount_Term	Float	Term of the loan in months
Credit_History	Float	Credit history of the applicant
Property_Area	Object	Urban/Semi Urban/Rural property area
Loan_Status	Object	Loan approval status (Target Variable)

Table 1: Data description of the variables in the dataset

## Data Structure

The data is structured in rows and columns, with each row representing an individual loan application, and each column representing a specific attribute or feature. The **Loan\_Status** variable in the **train** file serves as the target variable that the machine learning models aim to predict.

The machine learning models will be trained on historical data from the **train** file, learning patterns and relationships between various features and the loan approval outcome. Subsequently, these models will be applied to the **test** file to predict the loan approval status for new, unseen data.

Understanding the nuances of each variable and the relationships within the dataset is essential for the successful development and deployment of the loan prediction system.

## Procedure

The implementation of the loan prediction system involves a systematic procedure. The primary machine learning models employed in this project include Logistic Regression, Decision Trees, and K-Means Clustering. Each model serves a specific purpose in predicting loan approval outcomes based on historical customer data.

## Implementation Steps

The implementation procedure follows key steps:

1. **Data Collection:** Using historical loan applicant data, including both approved and denied loan applications, from the bank's existing database.  
Includes Univariate and Bivariate analysis of the Individual attributes of the dataset as shown through the below plots:

## Univariate Analysis

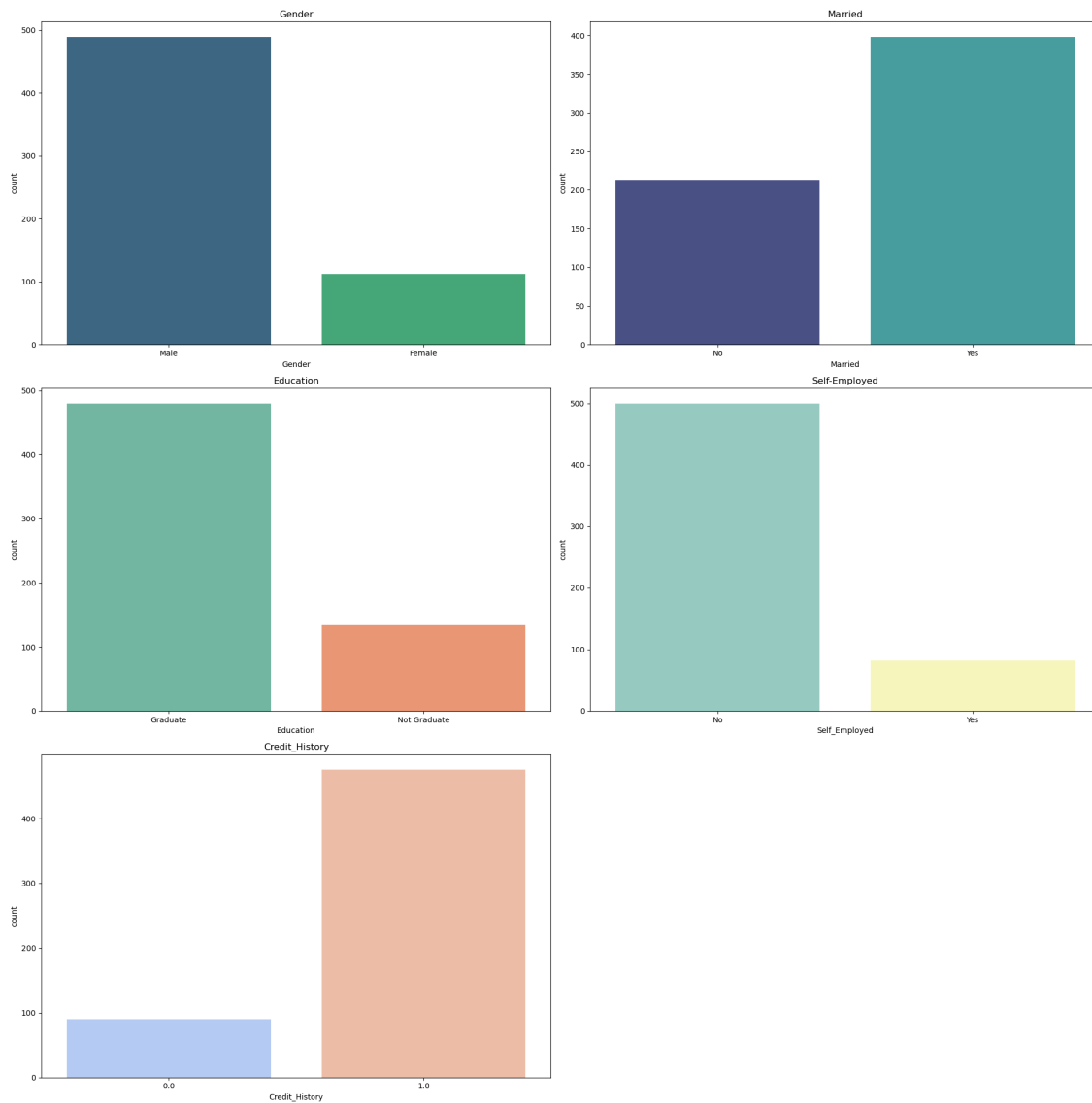


Figure 1: Data Frame plotting different attributes

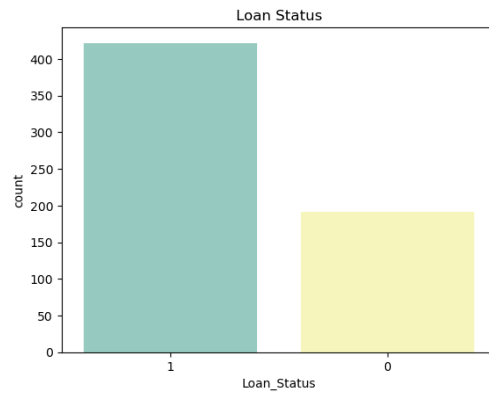


Figure 2: Target Value: Loan Status

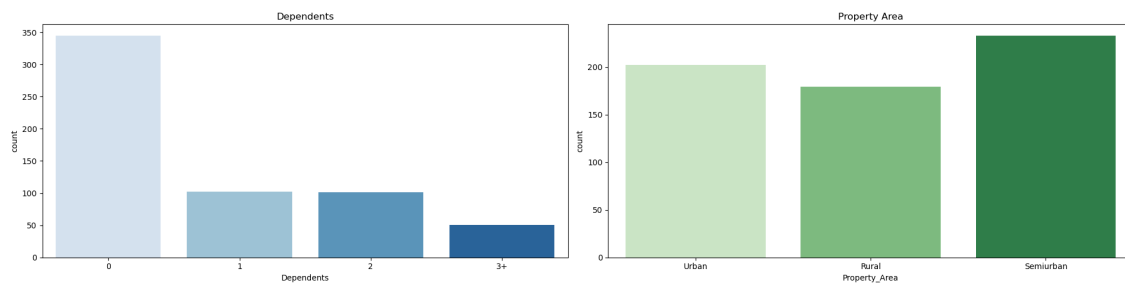


Figure 3: Target Value: Dependents and Property Area Plot

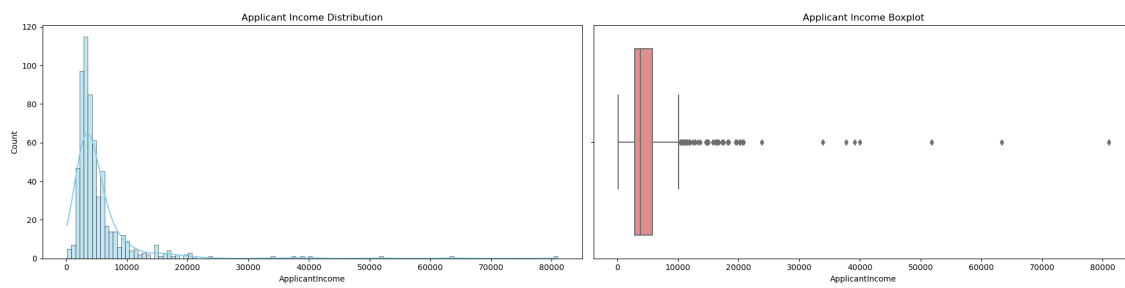


Figure 4: Target Value: Applicant Income

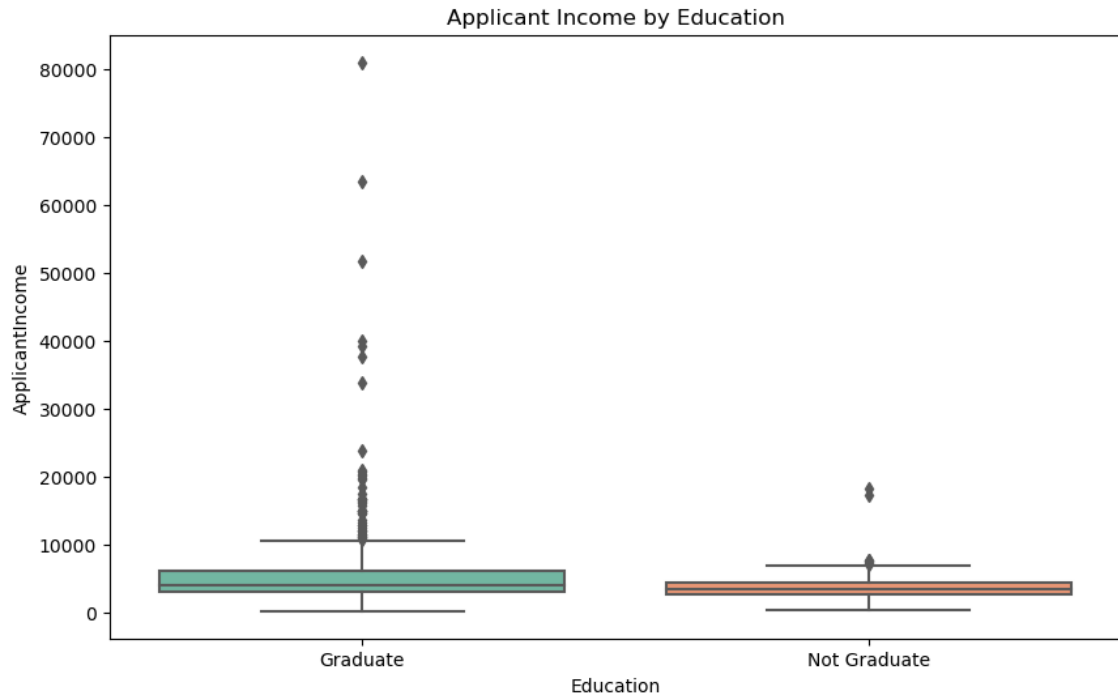


Figure 5: Applicant Income by Education

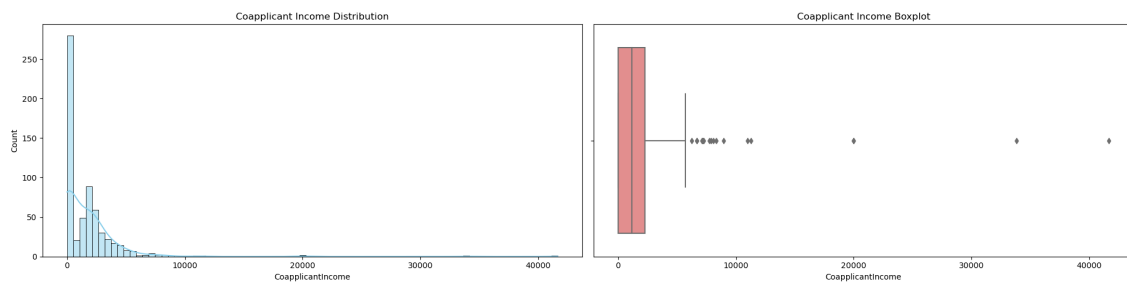


Figure 6: Target Value: Coapplicant Income

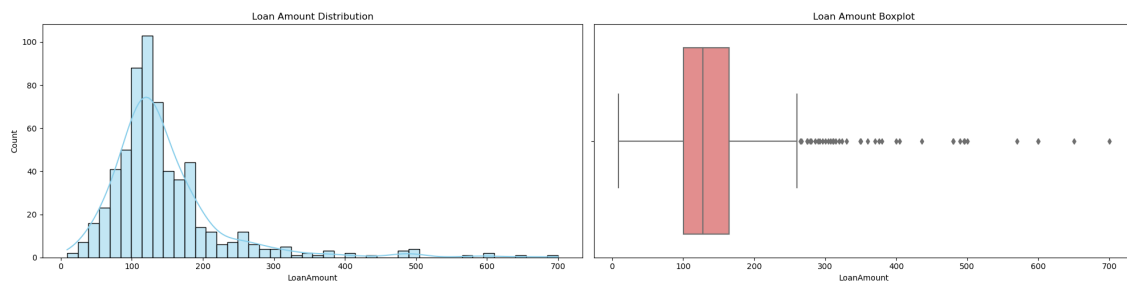
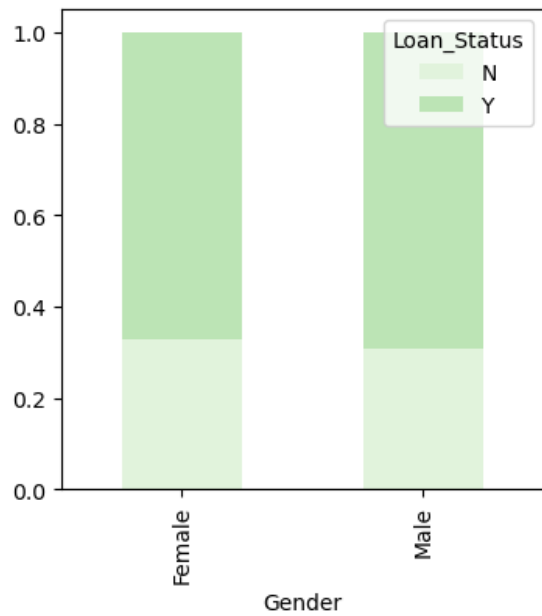
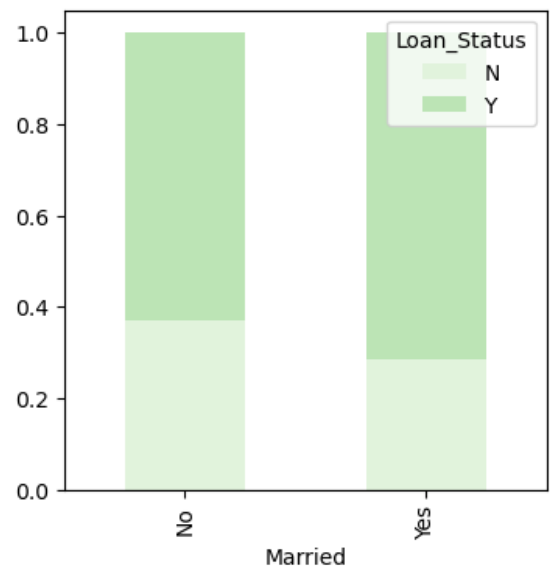


Figure 7: Target Value: Loan Amount

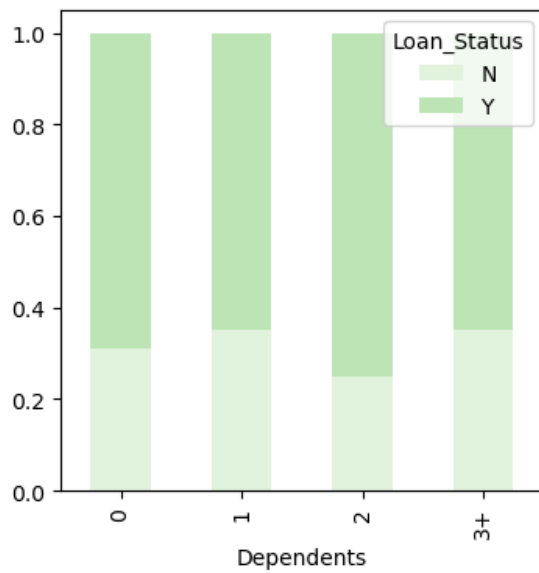
## Bivariate Analysis



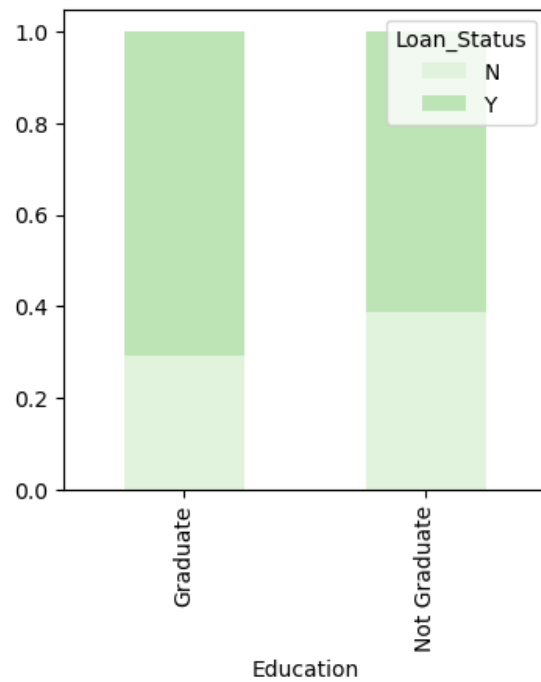
(a) Target Value: Loan Status according to Gender



(b) Target Value: Loan Status according to Marriage

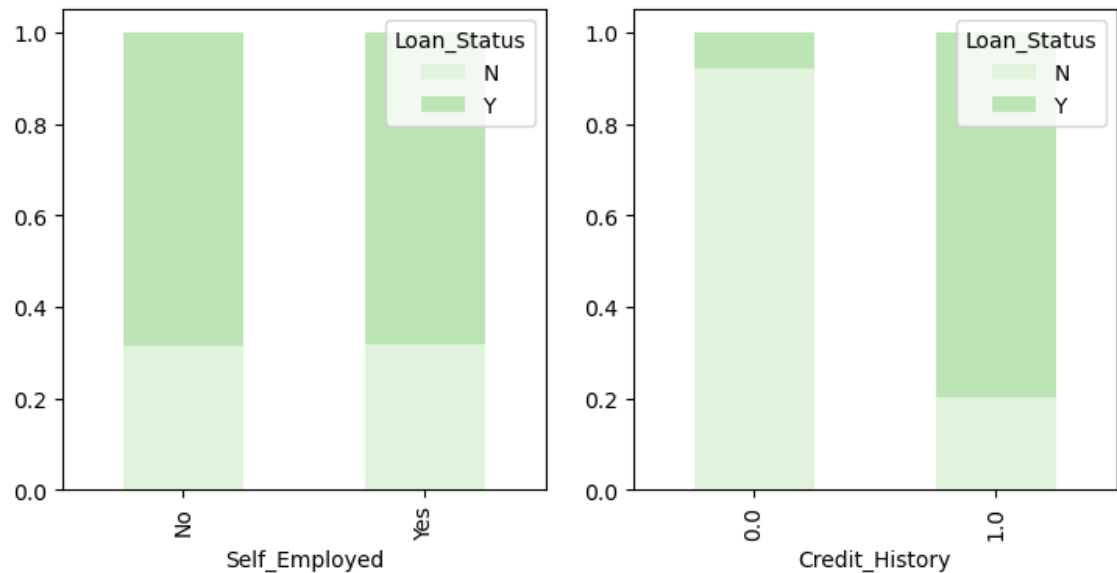


(c) Target Value: Loan Status according to Dependents



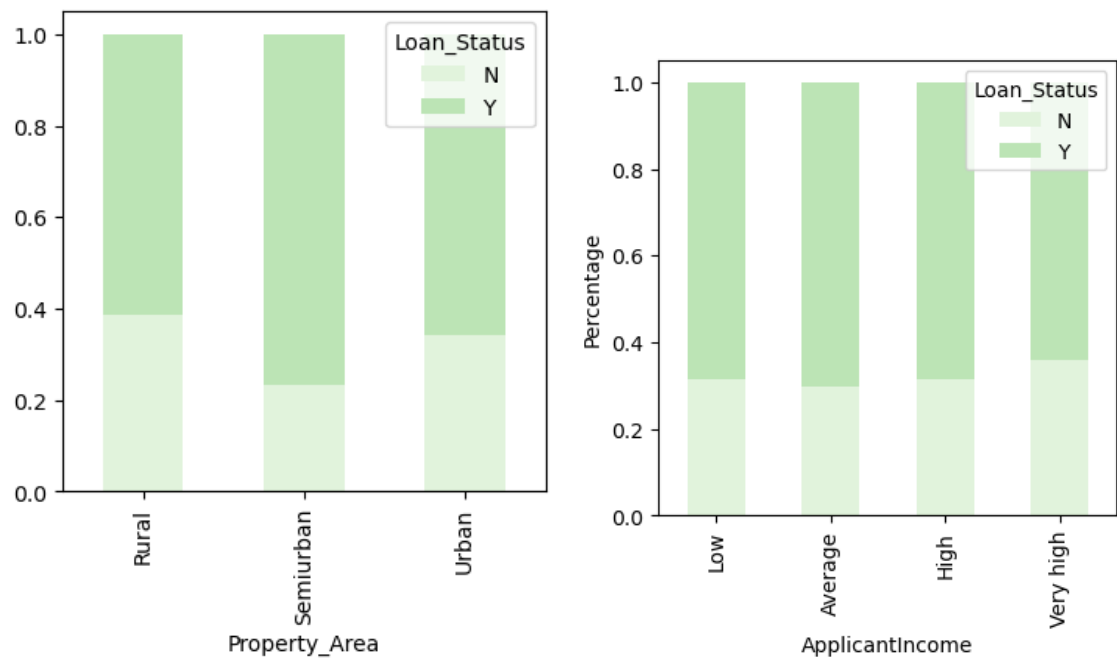
(d) Target Value: Loan Status according to Education

Figure 8: Bivariate Analysis for multiple attributes



(a) Target Value: Loan Status according to Employment (b) Target Value: Loan Status according to Credit History

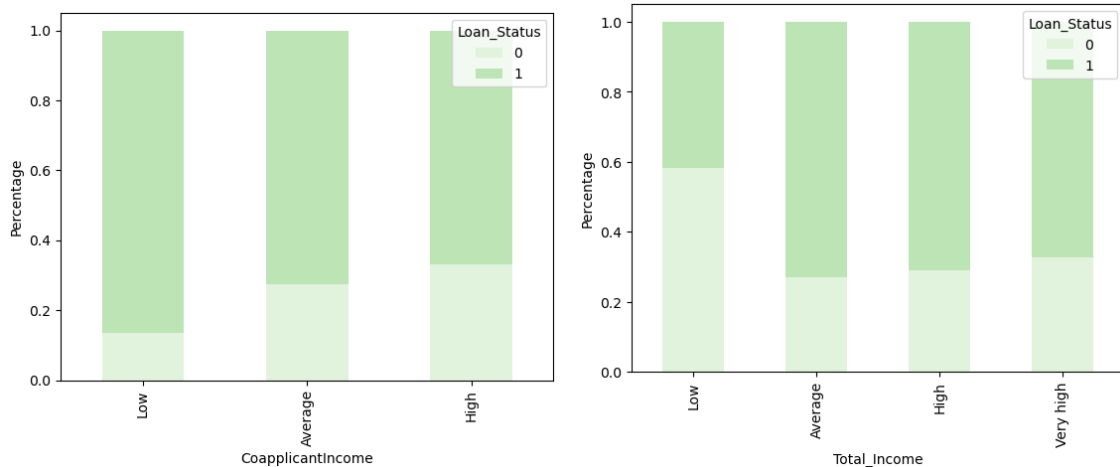
Figure 9: Bivariate Analysis for Employment and Credit History



(a) Target Value: Loan Status according to Property Area (b) Target Value: Loan Status according to Applicant Income

Figure 10: Bivariate Analysis for Property Area and Applicant Income





(a) Target Value: Loan Status according to Coapplicant Income (b) Target Value: Loan Status according to Total Income

Figure 11: Bivariate Analysis for Coapplicant Income and Total Income

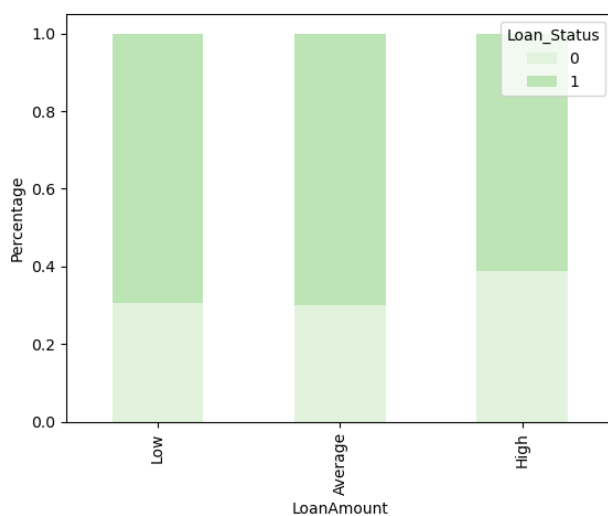


Figure 12: Target Value: Loan Status according to Loan Amount

## Correlation Using HeatMap

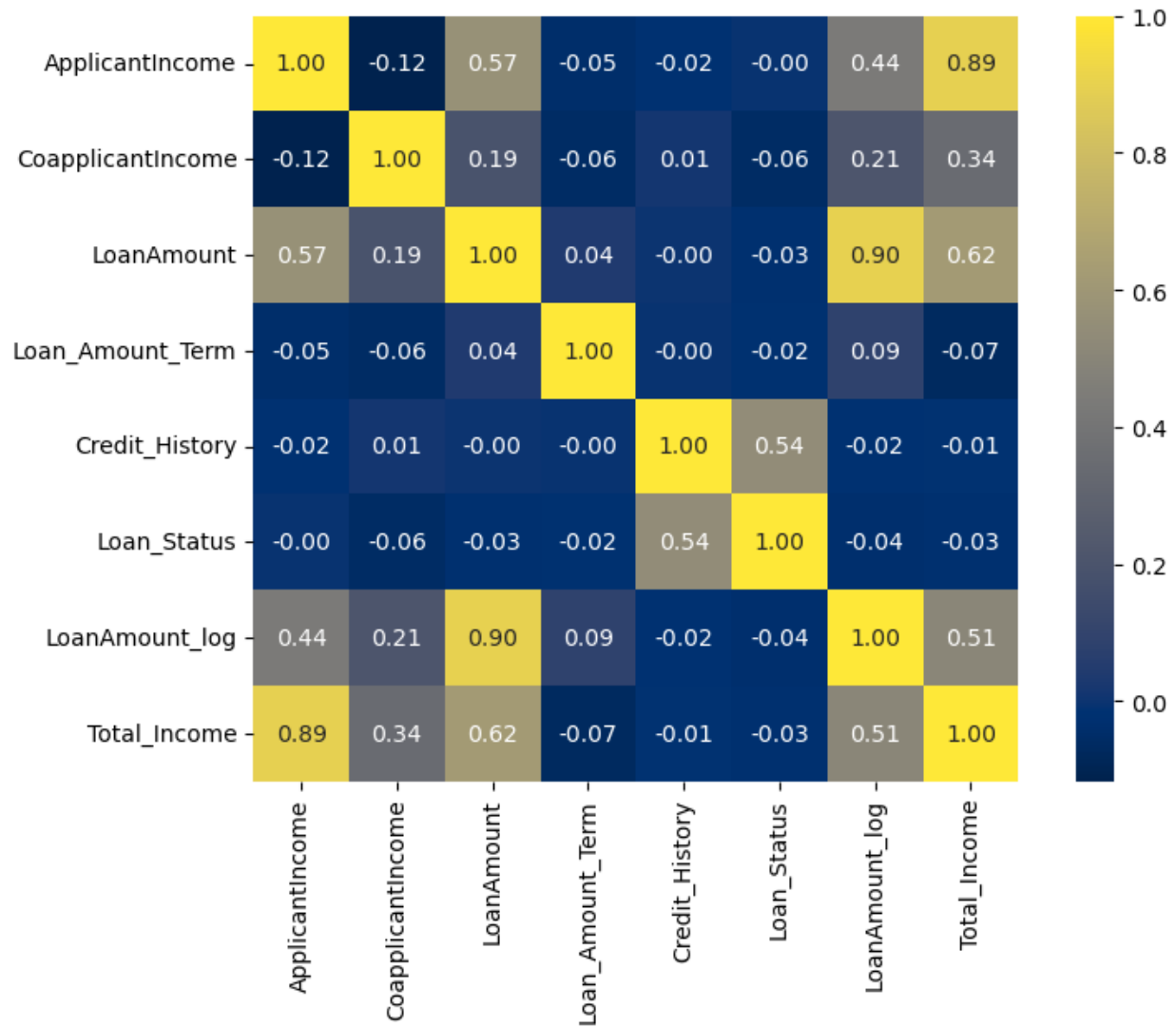


Figure 13: Correlation Matrix for numeric attributes

2. **Data Preprocessing:** Clean and prepare the data for analysis by handling missing values, dropping unwanted columns which will not be used for further analysis.

## Outlier Treatment

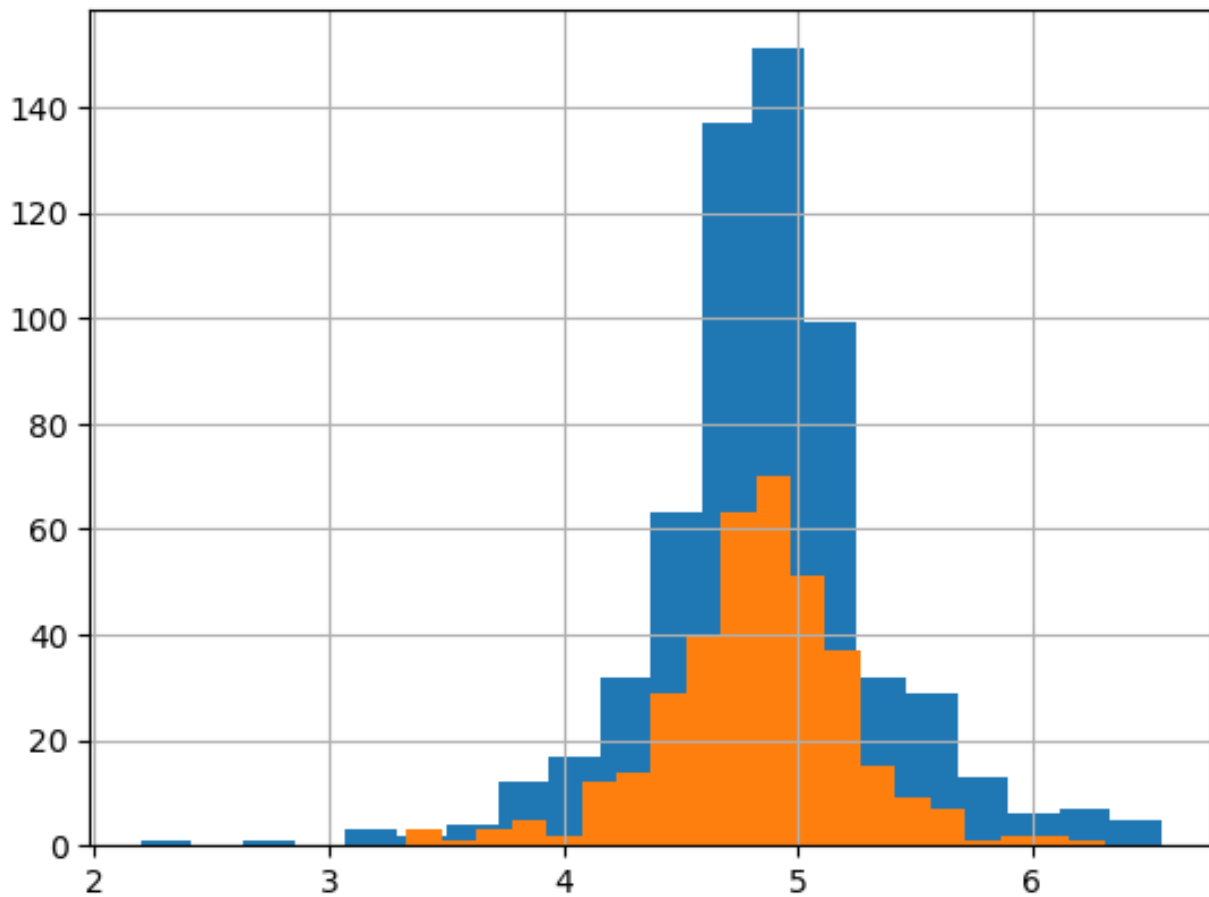


Figure 14: Outlier Treatment for training and testing set

3. **Feature Selection:** Identify the most relevant features that significantly impact loan approval and default rates.
4. **Model Training:** Develop and train machine learning models (Logistic Regression, Decision Trees, K-Means Clustering) using historical data to predict loan approval outcomes.

## Model Building

Loan\_ID will not be used for further analysis. Hence, dropping the Loan\_ID Column

```
In [768]: to_train=to_train.drop("Loan_ID",axis=1)
          to_train.head()
```

```
Out[768]:
```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area
0	Male	No	0	Graduate	No	5849	0.0	128.0	360.0	1.0	U
1	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	1.0	R
2	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0	1.0	U
3	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0	1.0	U
4	Male	No	0	Graduate	No	6000	0.0	141.0	360.0	1.0	U

```
In [769]: to_test=to_test.drop("Loan_ID",axis=1)
          to_test.head()
```

```
Out[769]:
```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area
0	Male	Yes	0	Graduate	No	5720	0	110.0	360.0	1.0	U
1	Male	Yes	1	Graduate	No	3076	1500	126.0	360.0	1.0	U
2	Male	Yes	2	Graduate	No	5000	1800	208.0	360.0	1.0	U
3	Male	Yes	2	Graduate	No	2340	2546	100.0	360.0	1.0	U
4	Male	No	0	Not Graduate	No	3276	0	78.0	360.0	1.0	U

```
In [770]: to_train=to_train.drop("Gender",axis=1)
          to_test=to_test.drop("Gender",axis=1)
```

```
In [771]: to_train=to_train.drop("Dependents",axis=1)
          to_test=to_test.drop("Dependents",axis=1)
```

```
In [772]: to_train=to_train.drop("Self_Employed",axis=1)
          to_test=to_test.drop("Self_Employed",axis=1)
```

Also Dropping the Loan\_Status column and storing it in another variable.

```
In [773]: x=to_train.drop("Loan_Status",axis=1)
          x.head()
```

```
Out[773]:
```

	Married	Education	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	LoanAmount_log
0	No	Graduate	5849	0.0	128.0	360.0	1.0	Urban	4.852030
1	Yes	Graduate	4583	1508.0	128.0	360.0	1.0	Rural	4.852030
2	Yes	Graduate	3000	0.0	66.0	360.0	1.0	Urban	4.189655
3	Yes	Not Graduate	2583	2358.0	120.0	360.0	1.0	Urban	4.787492
4	No	Graduate	6000	0.0	141.0	360.0	1.0	Urban	4.948760

```
In [774]: y=to_train["Loan_Status"]
          y.head()
```

```
Out[774]: 0    1
          1    0
          2    1
          3    1
          4    1
          Name: Loan_Status, dtype: int64
```

Creating Dummy Variable

```
In [775]: x=pd.get_dummies(x)
          to_train=pd.get_dummies(to_train)
          to_test=pd.get_dummies(to_test)
```

```
In [776]: x.head()
```

```
Out[776]:
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	LoanAmount_log	Married_No	Married_Yes	Education_Graduate	Ed
0	5849	0.0	128.0	360.0	1.0	4.852030	True	False	True	
1	4583	1508.0	128.0	360.0	1.0	4.852030	False	True	True	
2	3000	0.0	66.0	360.0	1.0	4.189655	False	True	True	
3	2583	2358.0	120.0	360.0	1.0	4.787492	False	True	False	
4	6000	0.0	141.0	360.0	1.0	4.948760	True	False	True	

5. **Model Evaluation:** Assess the performance of trained models using accuracy metrics to ensure reliability of predictions.

## Applying Logistic Regression

```
In [777]: from sklearn.model_selection import train_test_split
          x_train, x_cv, y_train, y_cv = train_test_split(x,y, train_size = 0.75,random_state=0)
```

```
In [778]: from sklearn.linear_model import LogisticRegression
          from sklearn.metrics import accuracy_score
          model = LogisticRegression()
          model.fit(x_train, y_train)
```

```
Out[778]: LogisticRegression
          LogisticRegression()
```

```
In [779]: pred_cv = model.predict(x_cv)
```

```
In [780]: accuracy_score(y_cv, pred_cv)
```

```
Out[780]: 0.8376623376623377
```

```
In [781]: from sklearn.metrics import confusion_matrix
          c = confusion_matrix(y_cv, pred_cv)
          c
```

```
Out[781]: array([[ 20,  23],
                 [  2, 109]])
```

```
In [782]: to_test.head()
```

Out[782]:

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	LoanAmount_log	Married_No	Married_Yes	Education_Graduate	Ed
0	5720	0	110.0	360.0	1.0	4.700480	False	True	True	
1	3076	1500	126.0	360.0	1.0	4.836282	False	True	True	
2	5000	1800	208.0	360.0	1.0	5.337538	False	True	True	
3	2340	2546	100.0	360.0	1.0	4.605170	False	True	True	
4	3276	0	78.0	360.0	1.0	4.356709	True	False	False	

```
In [783]: pred_test = model.predict(to_test)
```

```
In [784]: submission=pd.read_csv("sample_submission.csv",header=0)
```

```
In [785]: submission["Loan_Status"]=pred_test  
submission["Loan_ID"]=training_original["Loan_ID"]
```

```
In [786]: submission.head()
```

Out[786]:

	Loan_ID	Loan_Status
0	LP001002	1
1	LP001003	1
2	LP001005	1
3	LP001006	1
4	LP001008	1

```
In [787]: submission["Loan_Status"].replace(0, "N",inplace=True)  
submission["Loan_Status"].replace(0, "Y",inplace=True)
```

```
In [788]: pd.DataFrame(submission, columns=["Loan_ID","Loan_Status"]).to_csv("logistic.csv")
```

```
In [789]: submission
```

Out[789]:

	Loan_ID	Loan_Status
0	LP001002	1
1	LP001003	1
2	LP001005	1
3	LP001006	1
4	LP001008	1
...	...	...
362	LP002175	1
363	LP002178	1
364	LP002180	1
365	LP002181	1
366	LP002187	1

367 rows x 2 columns

## Applying Stratified K-means Clustering

```
In [790]: from statistics import mean
```

```
In [791]: from sklearn.model_selection import StratifiedKFold
```

```
In [792]: i = 1
pred_scores=[]
kf = StratifiedKFold(n_splits=5,random_state=1,shuffle=True)
for train_index, test_index in kf.split(x,y):
    print('\n{} of kfold {}'.format(i,kf.n_splits))
    xtr, xvl = x.loc[train_index],x.loc[test_index]
    ytr, yvl = y[train_index],y[test_index]
    model = LogisticRegression(random_state=1)
    model.fit(xtr, ytr)
    pred_test = model.predict(xvl)
    score = accuracy_score(yvl, pred_test)
    print("accuracy_score",score)
    i+=1
    pred_test = model.predict(to_test)
    pred=model.predict_proba(xvl)[: ,1]
    pred_scores.append(score)
print("\n Mean of Accuracy Scores=",mean(pred_scores))
```

```
1 of kfold 5
accuracy_score 0.7967479674796748

2 of kfold 5
accuracy_score 0.8373983739837398

3 of kfold 5
accuracy_score 0.7967479674796748

4 of kfold 5
accuracy_score 0.8130081300813008

5 of kfold 5
accuracy_score 0.7950819672131147

Mean of Accuracy Scores= 0.807796881247501
```

## Decision Tree Algorithm

```
In [793]: from sklearn import tree
```

```
In [794]: i=1
kf = StratifiedKFold(n_splits=5,random_state=1,shuffle=True)
for train_index, test_index in kf.split(x,y):
    print("\n {} of kfold {}".format(i,kf.n_splits))
    xtr,xvl = x.loc[train_index],x.loc[test_index]
    ytr,yvl = y[train_index],y[test_index]
    model = tree.DecisionTreeClassifier(random_state=1)
    model.fit(xtr, ytr)
    pred_test = model.predict(xvl)
    score = accuracy_score(yvl, pred_test)
    print("accuracy_score",score)
    i+=1
    pred_test = model.predict(to_test)
```

```

1 of kfold 5
accuracy_score 0.6585365853658537

2 of kfold 5
accuracy_score 0.7154471544715447

3 of kfold 5
accuracy_score 0.7235772357723578

4 of kfold 5
accuracy_score 0.7154471544715447

5 of kfold 5
accuracy_score 0.6475409836065574

```

```

In [795]: submission["Loan_Status"] = pred_test
          submission["Loan_ID"] = testing_original["Loan_ID"]
          submission.head()

```

Out[795]:

	Loan_ID	Loan_Status
0	LP001015	1
1	LP001022	1
2	LP001031	1
3	LP001035	0
4	LP001051	1

```

In [796]: submission["Loan_Status"].replace(0, "N", inplace=True)
          submission["Loan_Status"].replace(1, "Y", inplace=True)
          pd.DataFrame(submission, columns=["Loan_ID", "Loan_Status"]).to_csv("Decision Tree.csv")

##End

```

## Evaluation

The performance of the loan prediction system is rigorously evaluated using accuracy metrics to ensure the reliability and effectiveness of the implemented machine learning models.

### Evaluation Metrics

The following key evaluation metric is employed to assess the performance of the models:

- **Accuracy:** The ratio of correctly predicted instances to the total instances.

### Results

The performance of the machine learning models — Logistic Regression, Decision Trees, and K-Means Clustering — is evaluated on the test dataset using the aforementioned metrics. The results are summarized below:

Model	Accuracy
Logistic Regression	0.83
Decision Trees	0.64
K-Means Clustering	0.80

Table 2: Performance Metrics for Machine Learning Models



The results indicate that the Logistic Regression model outperforms the other models in terms of accuracy. K-Means Clustering, originally designed for clustering tasks also exhibit reasonable performance, Decision Tree may not be the optimal choice for the problem at hand.

## Conclusion

In conclusion, the Loan Prediction Analysis for Banks project represents a significant stride towards addressing the critical challenges faced by lending institutions in the ever-evolving financial landscape. By harnessing the power of machine learning, we have successfully developed predictive models that contribute to more informed decision-making in the loan approval process.

The project aimed to answer the fundamental questions plaguing the lending industry:

How risky is the borrower, and should a loan be granted based on this risk assessment?

Through extensive exploration of various machine learning algorithms and methodologies, we have crafted a robust predictive system that leverages historical data to accurately evaluate the creditworthiness of loan applicants.

The implementation of the Loan Prediction System not only enhances the accuracy of risk assessment but also streamlines operational processes within banks. By significantly reducing the reliance on manual efforts and optimizing resource utilization, the system serves as a beacon of efficiency in an industry where time and precision are of the essence.

Real-world testing and evaluation have demonstrated the effectiveness of our predictive models, providing a tangible solution to the challenges faced by banks in the loan approval domain. The measured metrics, including accuracy attest to the reliability of the system in making informed predictions.

As we embrace the outcomes of this project, it is essential to acknowledge that the landscape of data and technology is dynamic.

In essence, the Loan Prediction Analysis for Banks project not only marks a successful integration of machine learning into financial decision-making but also serves as a foundation for ongoing advancements in predictive analytics within the banking industry. As we look ahead, the commitment to innovation and adaptability will be crucial in ensuring the sustained effectiveness and relevance of the Loan Prediction System in the face of changing economic landscapes and banking paradigms.