

# Loan Prediction

Yifan Chen, Xupeng Tang, Chris  
Yang, Jiajie Yao, Hongyan Zhao





# Background

In the realm of financial services, the accurate prediction of loan approval stands as a cornerstone for responsible lending.

This project delves into the realm of machine learning, specifically leveraging the Loan Prediction Problem Dataset. By scrutinizing applicant data encompassing demographics, financial metrics, and credit history, our goal is to construct a predictive model, to determine whether to approve the loan for the applicant.

This model aims to not only streamline decision-making within financial institutions but also foster a more nuanced understanding of risk factors, contributing to the cultivation of responsible and informed lending practices.



# Data Description

Our dataset comprises several key variables:

- **Loan\_ID**: Unique identifier for each loan application.
- **Gender**: Gender of the applicant (Male, Female).
- **Married**: Marital status of the applicant (Married, Not Married).
- **Dependents**: Number of dependents (0, 1, 2, 3+).
- **Education**: Educational background of the applicant (Graduate, Not Graduate).
- **Self\_Employed**: Whether the applicant is self-employed (Self Employed, Not Self Employed).
- **ApplicantIncome**: Income of the applicant (in dollars).
- **CoapplicantIncome**: Income of the co-applicant (in dollars).
- **LoanAmount**: Amount of the loan requested (in thousands of dollars).
- **Loan\_Amount\_Term**: Term of the loan (in months).
- **Credit\_History**: Credit history of the applicant (1 for good credit history, 0 for otherwise).
- **Property\_Area**: Area where the property is located (Rural, Urban, Semi-Urban).
- **Loan\_Status**: The target variable indicating loan approval status (Y for approved, N for denied).



# Dataset

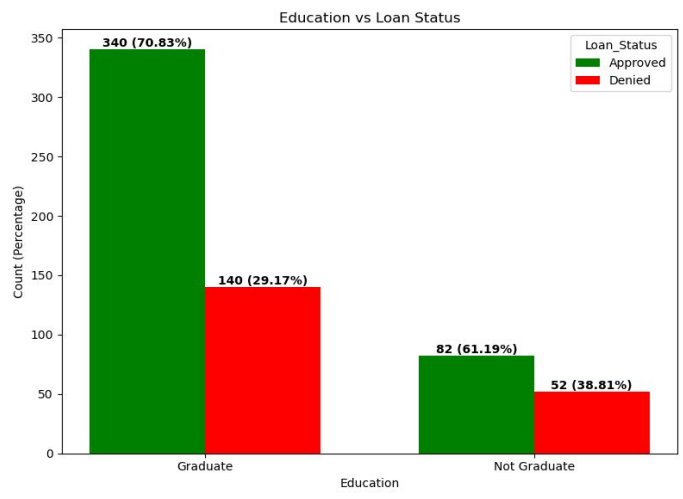
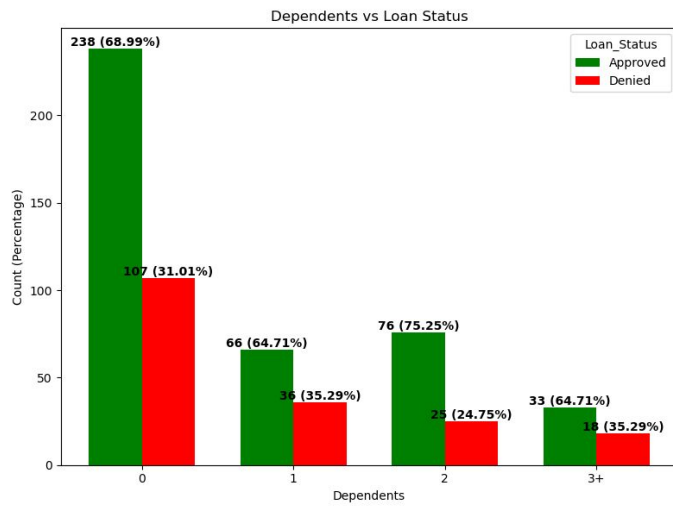
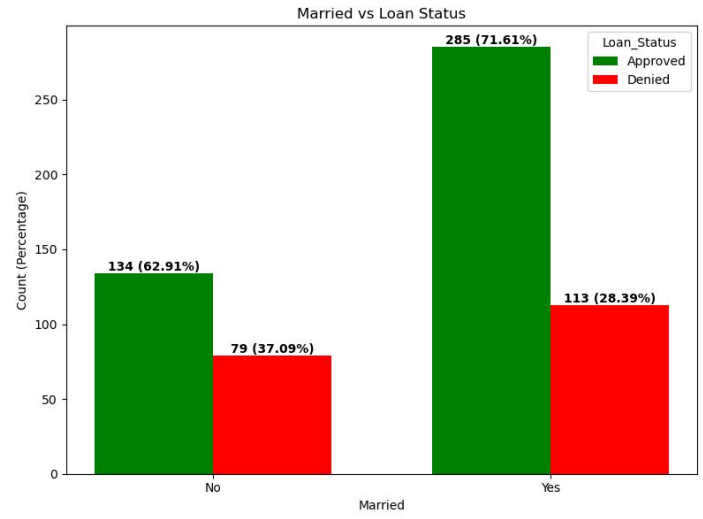
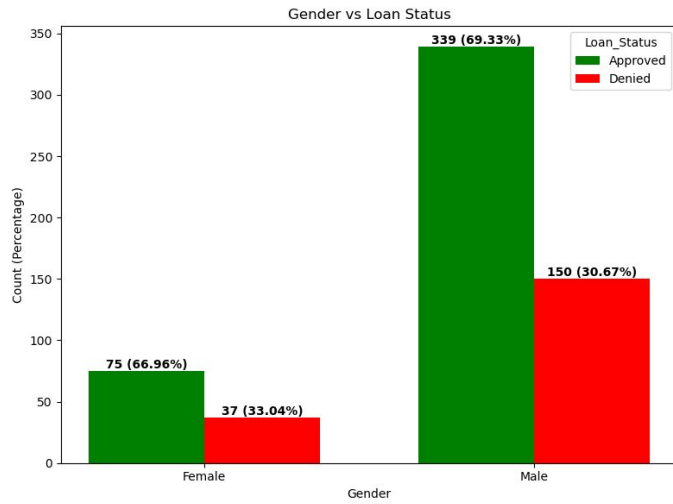
The training data has 614 rows and 13 columns.

Training Data:

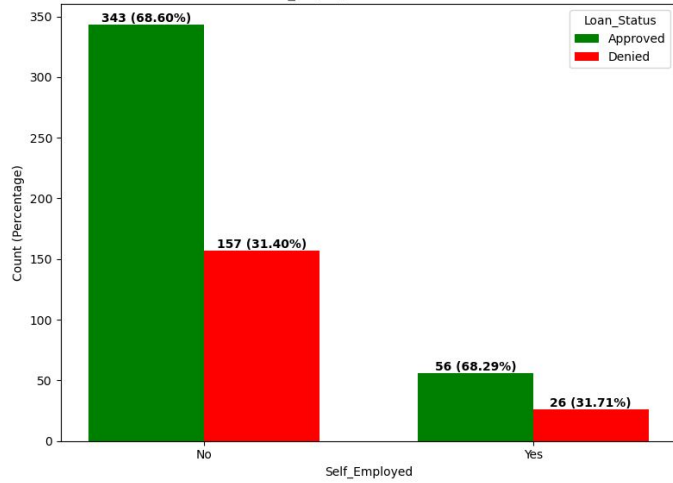
	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
0	LP001002	Male	No	0	Graduate	No	5849	0.0	NaN	360.0	1.0	Urban	Y
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	1.0	Rural	N
2	LP001005	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

The background is a solid orange color. In the top-left corner, there are three vertical bars of varying heights, each composed of several overlapping semi-transparent orange circles. In the bottom-right corner, there are four vertical bars of increasing height from left to right, each also composed of several overlapping semi-transparent orange circles.

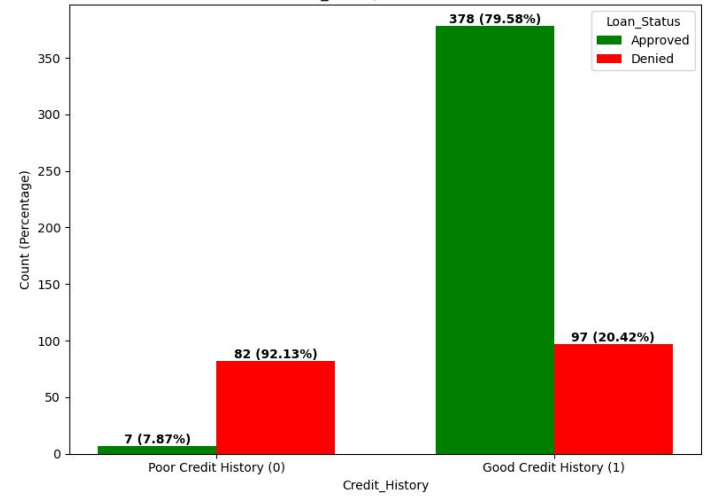
# Exploratory Data Analysis



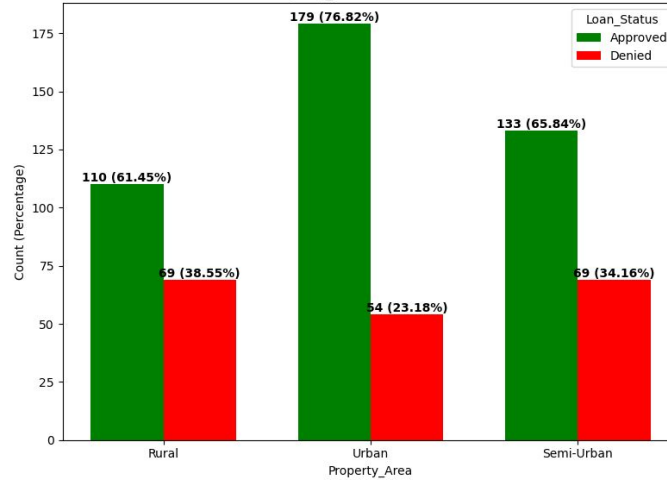
Self\_Employed vs Loan Status

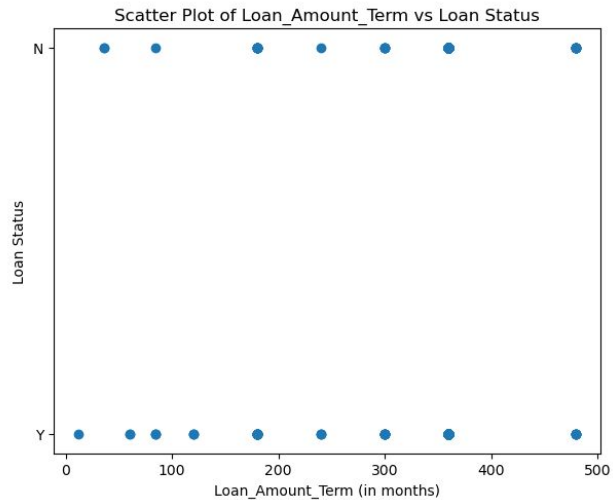
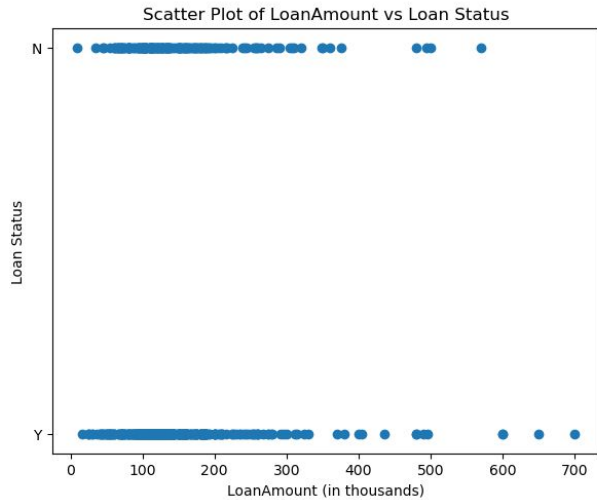
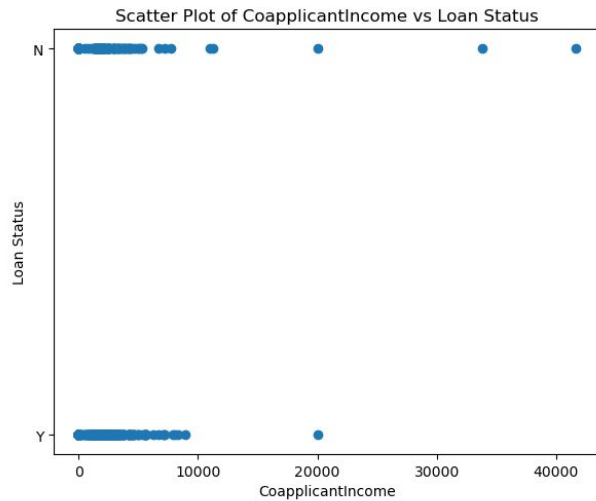
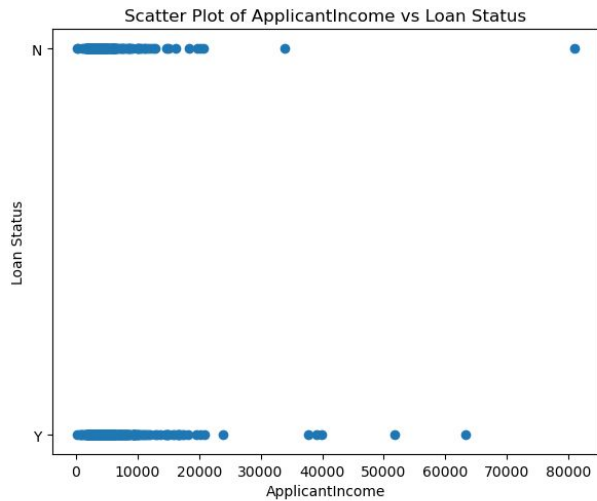


Credit\_History vs Loan Status



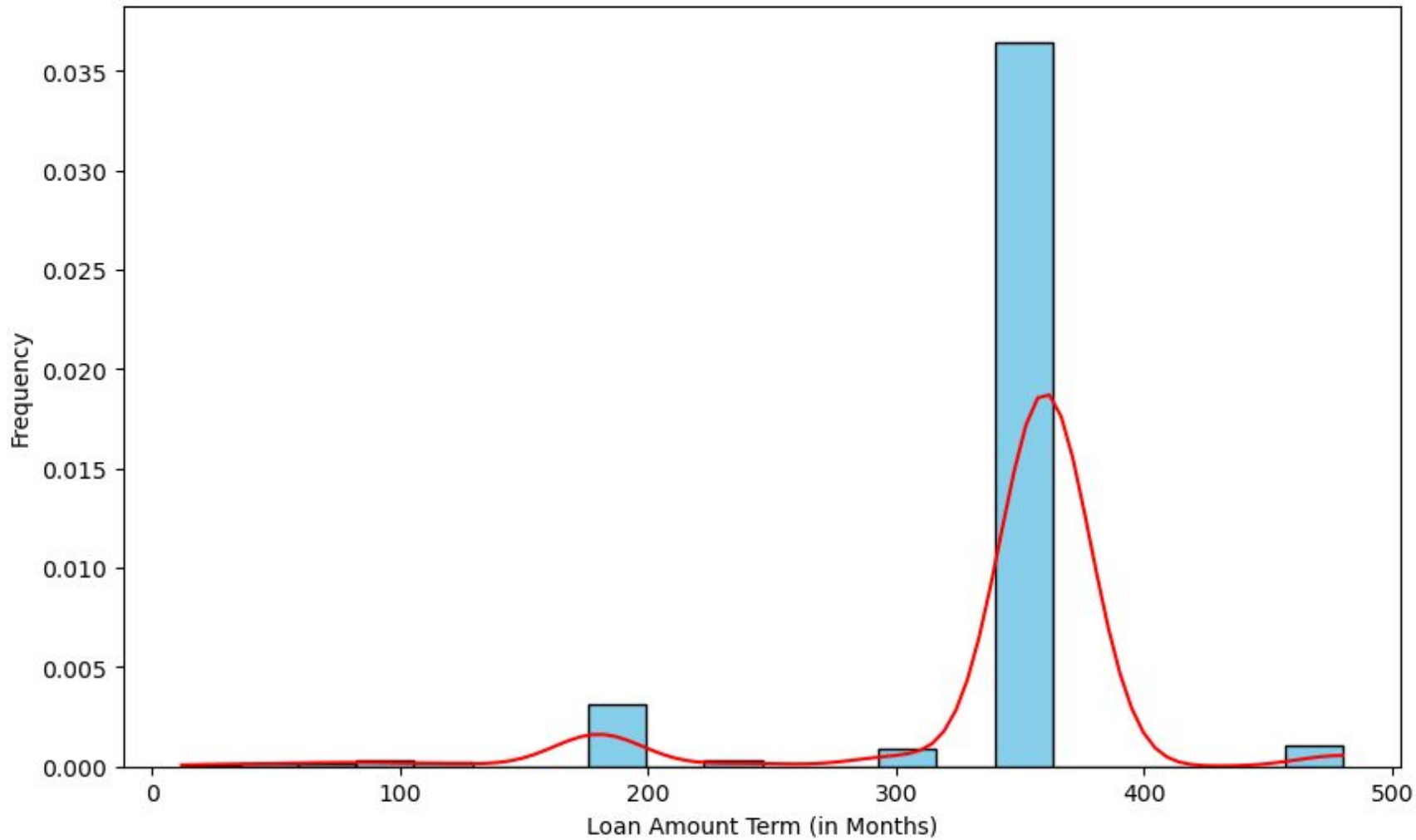
Property\_Area vs Loan Status



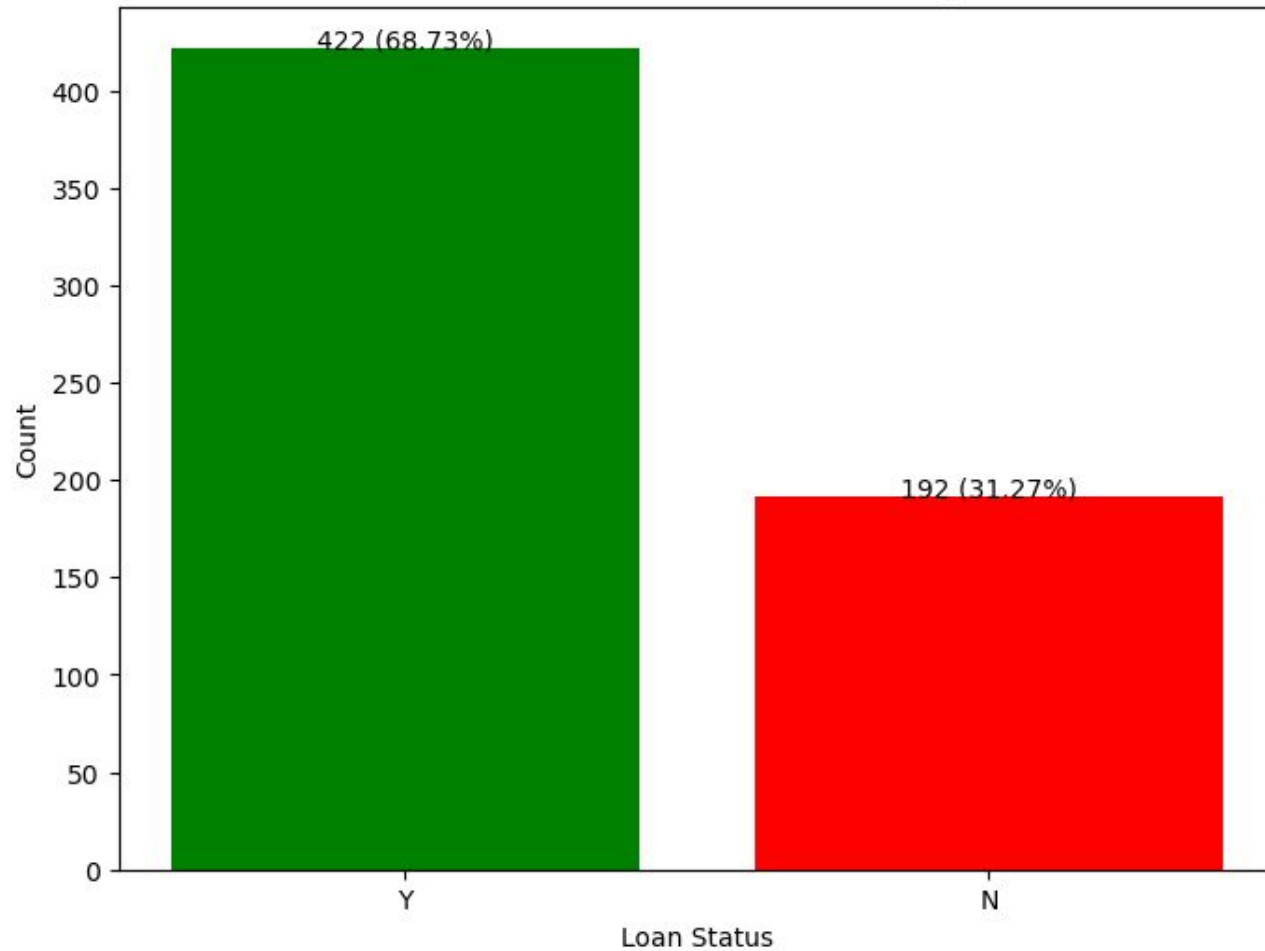




Distribution of Loan Amount Term



Count of Loan Status with Percentages





# Preparation before modeling

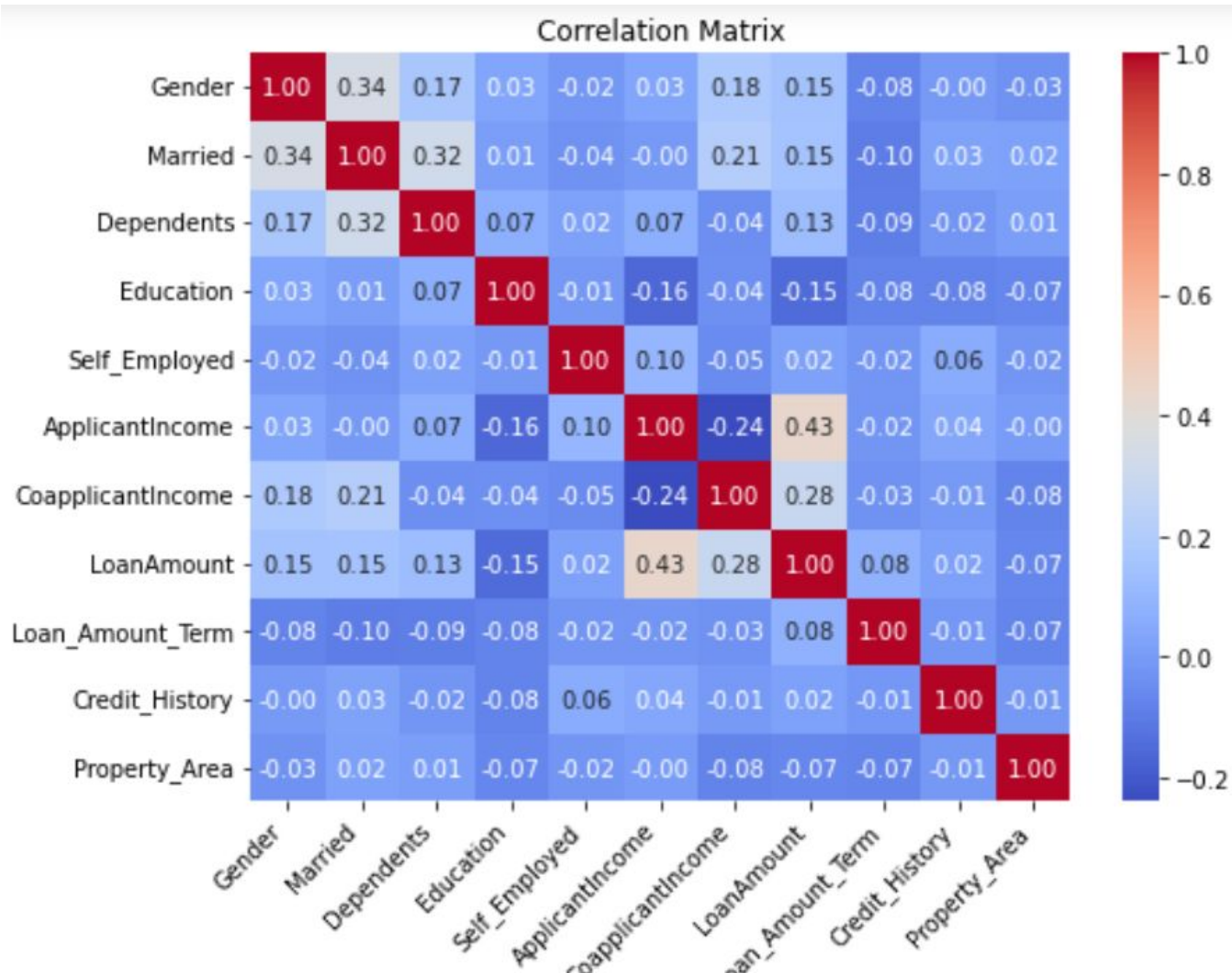


# Preprocessing Data

1. Convert categorical variables and non-standard numerical values to numeric values  
#For simplicity, we use Label Encoding for categorical variables
2. Impute missing values  
# For categorical columns, fill missing values with the mode  
# For numerical columns, fill missing values with the median
3. Remove outliers  
# Here, we use a simple method of removing values that are beyond 3 standard deviations from the mean

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
0	1	0	0.0	0	0	5849	0.0	128.0	360.0	1.0	2	1
1	1	1	1.0	0	0	4583	1508.0	128.0	360.0	1.0	0	0
2	1	1	0.0	0	1	3000	0.0	66.0	360.0	1.0	2	1
3	1	1	0.0	1	0	2583	2358.0	120.0	360.0	1.0	2	1
4	1	0	0.0	0	0	6000	0.0	141.0	360.0	1.0	2	1

# Correlation Matrix





# PCA

Explained variance ratio is set to 95%

## Principal Components

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10
0	-1.048034	0.997088	0.197283	1.117371	-0.411208	-0.202591	0.131674	-0.636012	0.967531	-0.815652
1	0.712545	-0.113204	-0.507286	-0.475809	0.336084	-0.954566	0.339778	0.958273	0.291930	0.259092
2	-0.734810	-0.195042	0.841875	1.113072	0.963517	0.008335	-1.693967	0.074033	1.056185	0.736767
3	-0.018778	-1.580996	-0.121421	0.434776	0.042295	-0.132564	-0.145632	-2.056172	0.185487	0.616348
4	-0.925248	1.126658	0.155562	1.090603	-0.446349	-0.163390	0.170565	-0.708914	0.918788	-0.826128



# PCA

## Explained Variance Ratio

	<b>PC</b>	<b>Explained Variance Ratio</b>
<b>0</b>	PC1	0.166351
<b>1</b>	PC2	0.135786
<b>2</b>	PC3	0.113989
<b>3</b>	PC4	0.097313
<b>4</b>	PC5	0.094600
<b>5</b>	PC6	0.085422
<b>6</b>	PC7	0.081580
<b>7</b>	PC8	0.071628
<b>8</b>	PC9	0.070197
<b>9</b>	PC10	0.052805



# Models







# SVM

Classifier	Parameters
SVM	{'C': [0.01, 1, 100]}
KNN	{'n_neighbors': [1, 2, 3, 4]}
Logistic Regression	{'C': [0.01, 1, 100]}

Using Grid search to find the best model.



# SVM

## Model Accuracies

	Classifier	Parameters	Accuracy	Model
0	SVM	{'C': 0.01}	0.669	SVC(C=0.01, random_state=42)
1	KNN	{'n_neighbors': 3}	0.585	KNeighborsClassifier(n_neighbors=3)
2	Logistic Regression	{'C': 1}	0.763	LogisticRegression(C=1, max_iter=5000)



# KNN

Classifier	Parameters
SVM	{'C': [0.01, 1, 100]}
KNN	{'n_neighbors': [1, 2, 3, 4]}
Logistic Regression	{'C': [0.01, 1, 100]}

Using Grid search to find the best model.



# KNN

## Model Accuracies

	Classifier	Parameters	Accuracy	Model
0	SVM	{'C': 0.01}	0.669	SVC(C=0.01, random_state=42)
1	KNN	{'n_neighbors': 3}	0.585	KNeighborsClassifier(n_neighbors=3)
2	Logistic Regression	{'C': 1}	0.763	LogisticRegression(C=1, max_iter=5000)



# Logistic Regression

Classifier	Parameters
SVM	{'C': [0.01, 1, 100]}
KNN	{'n_neighbors': [1, 2, 3, 4]}
Logistic Regression	{'C': [0.01, 1, 100]}

Using Grid search to find the best model.



# Logistic Regression

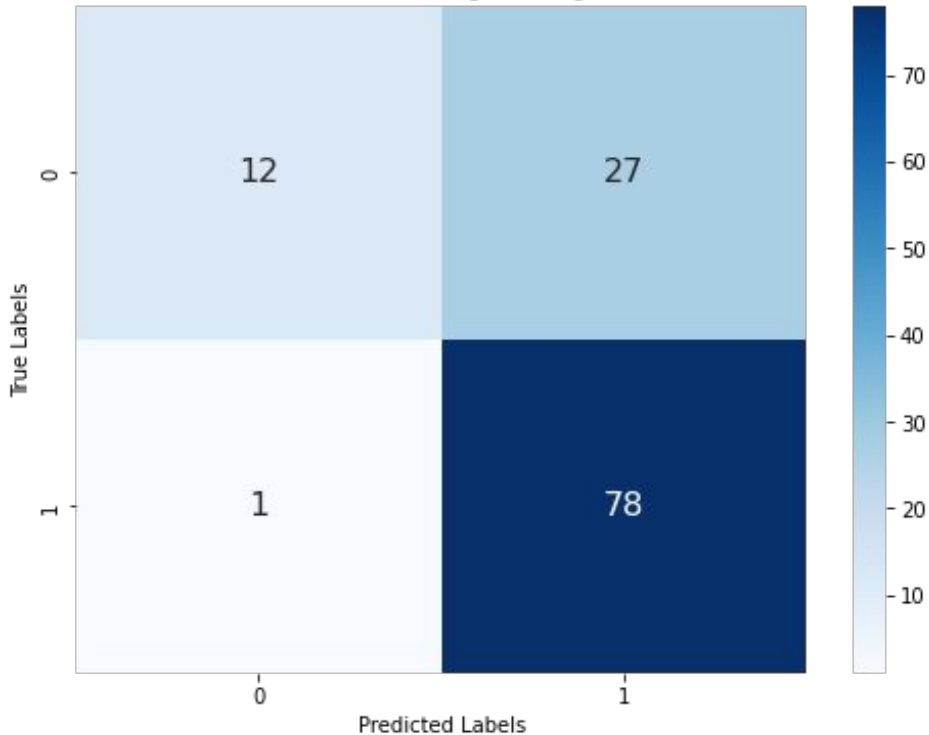
## Model Accuracies

	Classifier	Parameters	Accuracy	Model
0	SVM	{'C': 0.01}	0.669	SVC(C=0.01, random_state=42)
1	KNN	{'n_neighbors': 3}	0.585	KNeighborsClassifier(n_neighbors=3)
2	Logistic Regression	{'C': 1}	0.763	LogisticRegression(C=1, max_iter=5000)

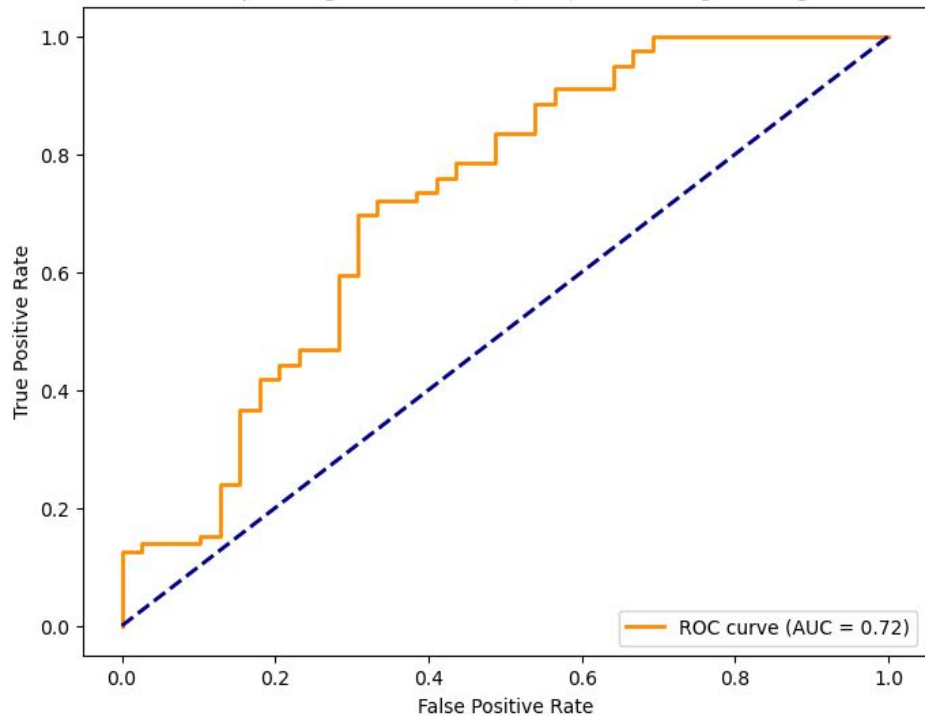


# Logistic Regression

Confusion Matrix - Logistic Regression



Receiver Operating Characteristic (ROC) Curve - Logistic Regression





# Logistic Regression

Logistic Regression Coefficients and Odds Ratios

	Feature	Coefficient	Odds Ratio
0	Gender	-0.5559	0.5735
1	Married	0.5258	1.6918
2	Dependents	0.1056	1.1113
3	Education	-0.6439	0.5253
4	Self_Employed	-0.1953	0.8226
5	ApplicantIncome	0.0000	1.0000
6	CoapplicantIncome	0.0001	1.0001
7	LoanAmount	-0.0058	0.9942
8	Loan_Amount_Term	-0.0020	0.9980
9	Credit_History	3.1788	24.0175
10	Property_Area	0.0850	1.0887

Compare Odds Ratio with 1

Education, Married, and Credit\_history are most significant variables.





## Decision Tree

Training Data Set Accuracy: 1.0

Training Data F1 Score 1.0

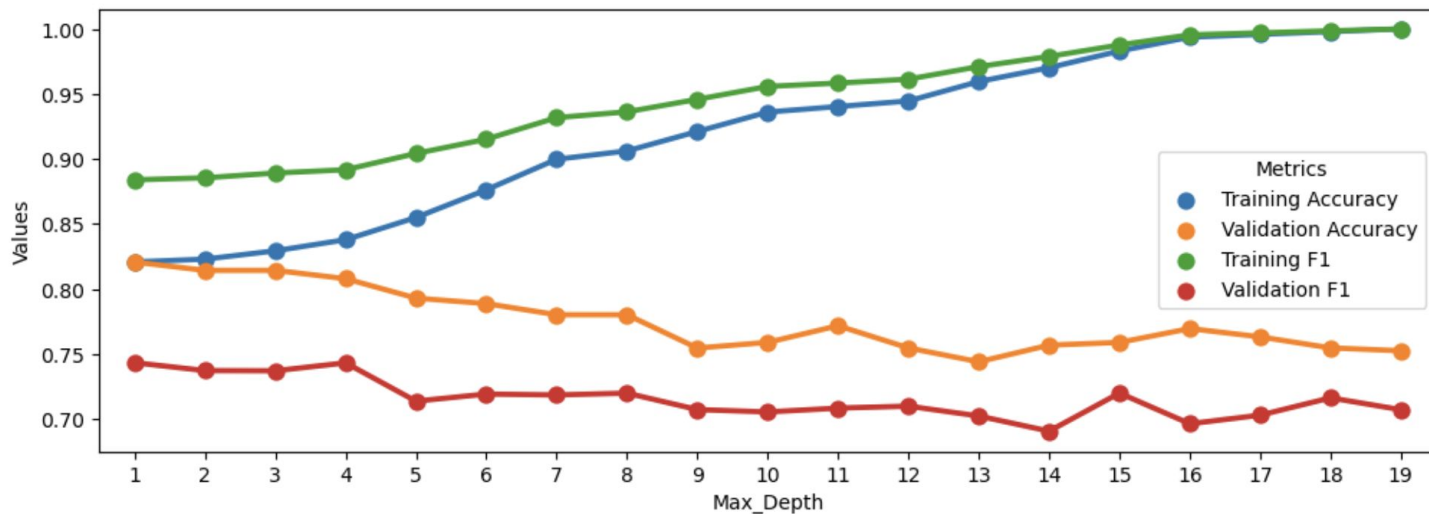
Validation Mean F1 Score: 0.704

Validation Mean Accuracy: 0.739



# Decision Tree

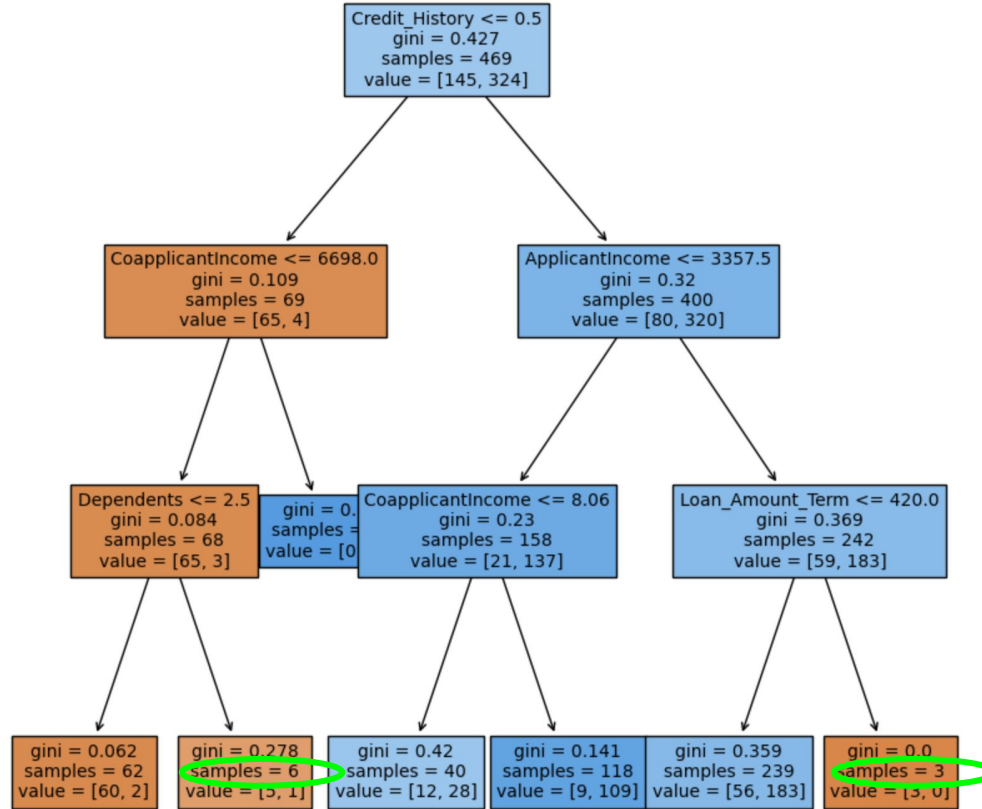
Tuning Max\_Depth



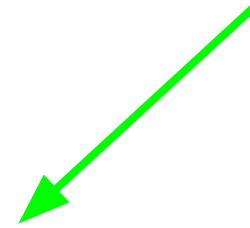
Test Accuracy: 0.779

Test F1 Score: 0.858

# Decision Tree



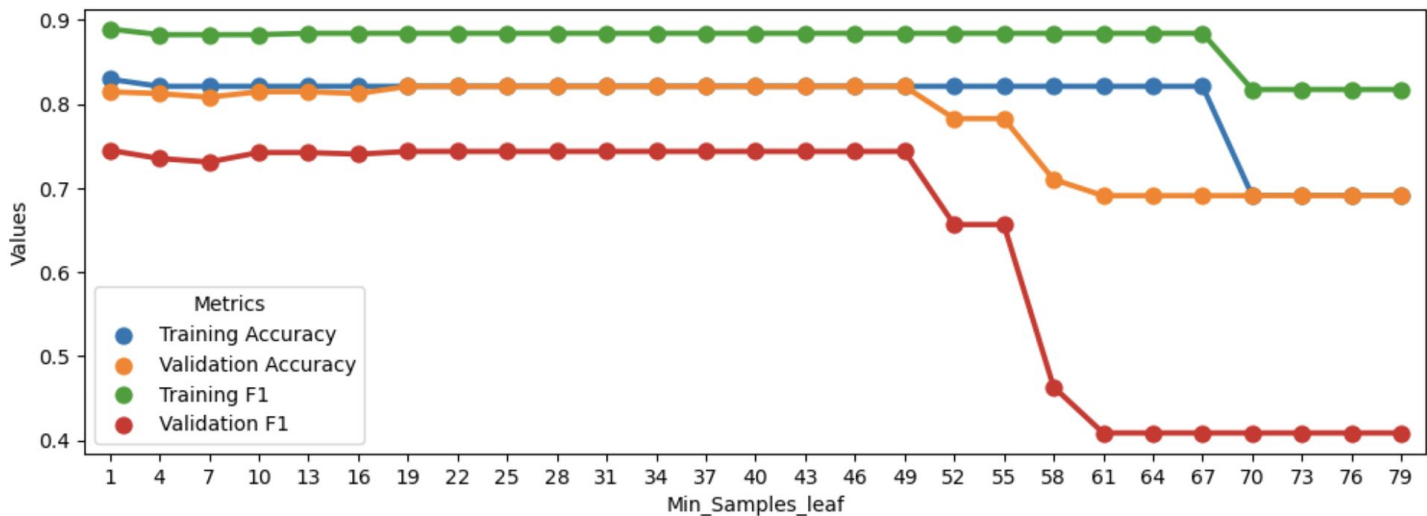
Overfitting





# Decision Tree

Tuning Min\_Samples\_leaf



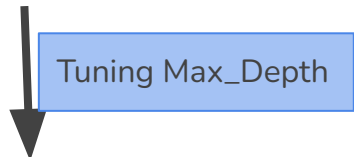
Test Accuracy: 0.796

Test F1 Score: 0.869



# Decision Tree

Test Accuracy: 0.677  
Test F1 Score: 0.771



Test Accuracy: 0.779  
Test F1 Score: 0.858



Test Accuracy: 0.796  
Test F1 Score: 0.869



## Random Forest

```
'n_estimators': [20, 30, 50],  
'max_depth': [None, 10, 20],  
'min_samples_split': [2, 5, 10],  
'min_samples_leaf': [10, 20, 35]
```

Using Grid  
search to find  
the best model.



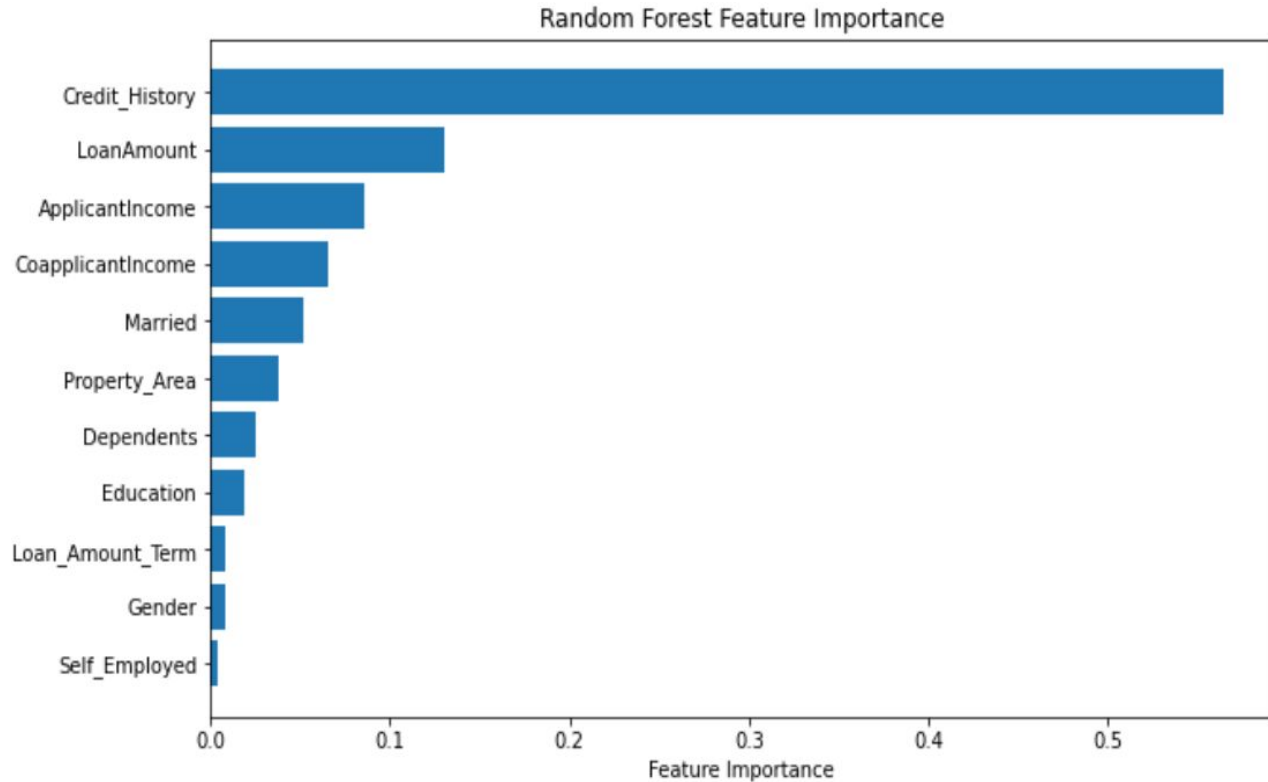
# Random Forest

Best Parameters: {'max\_depth': None, 'min\_samples\_leaf': 10, 'min\_samples\_split': 2, 'n\_estimators': 20}

Accuracy on Test Set: 0.771



# Random Forest

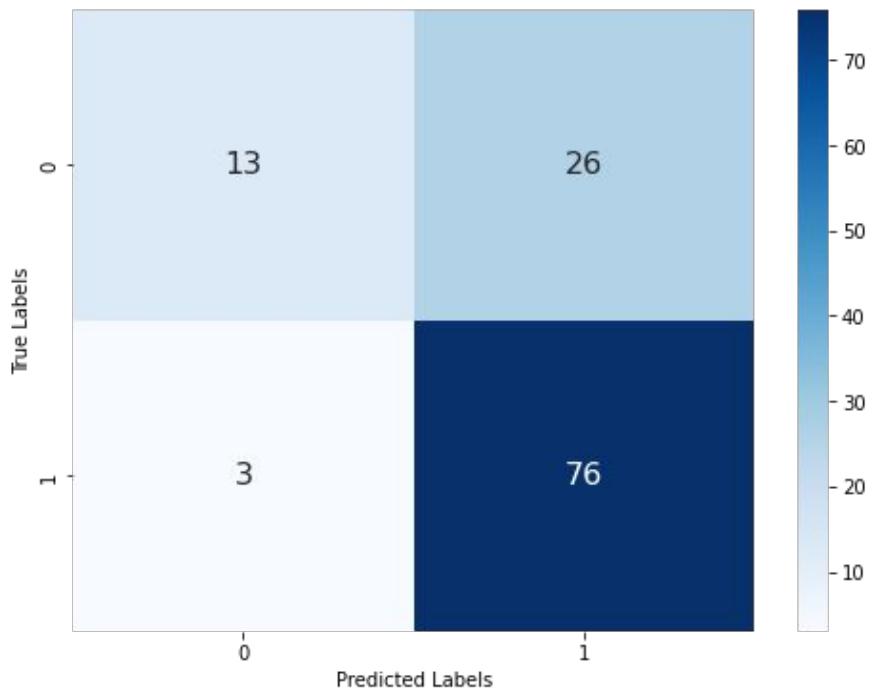




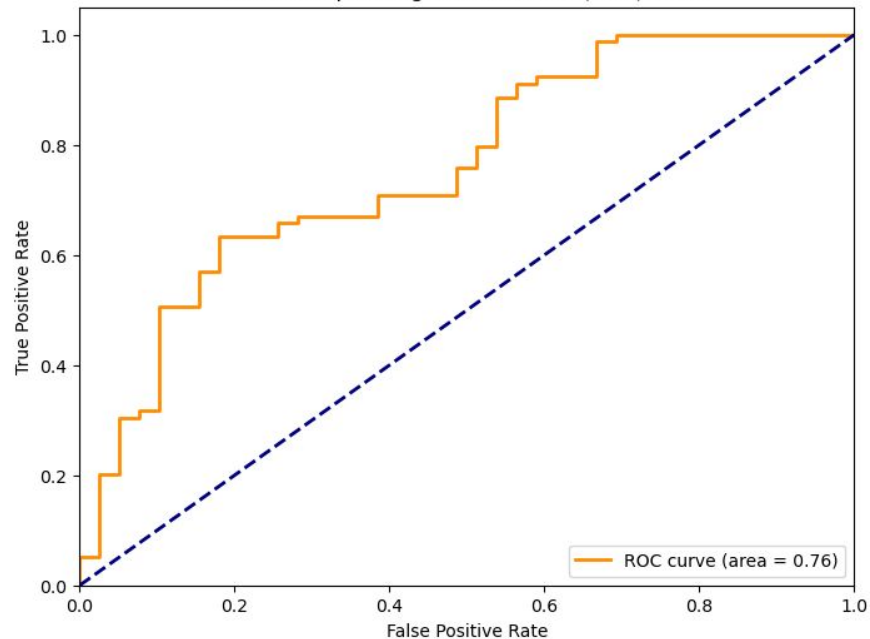


# Random Forest

Confusion Matrix

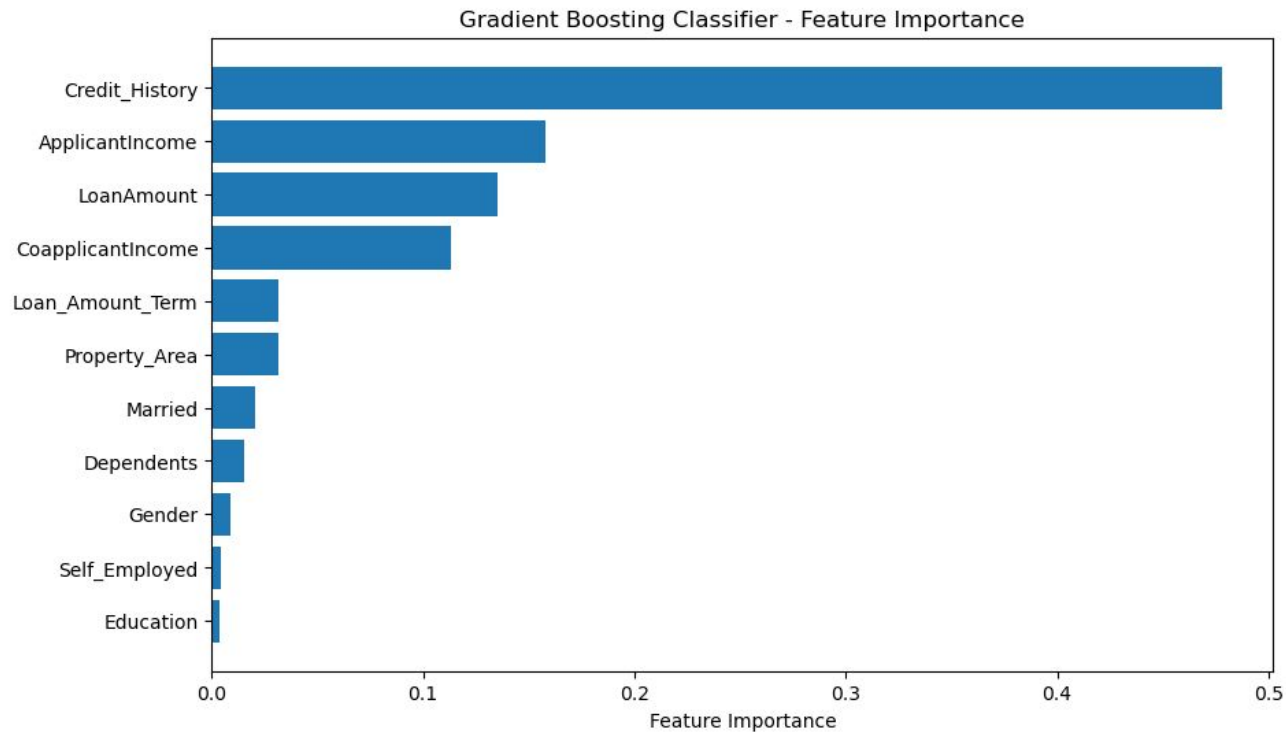


Receiver Operating Characteristic (ROC) Curve



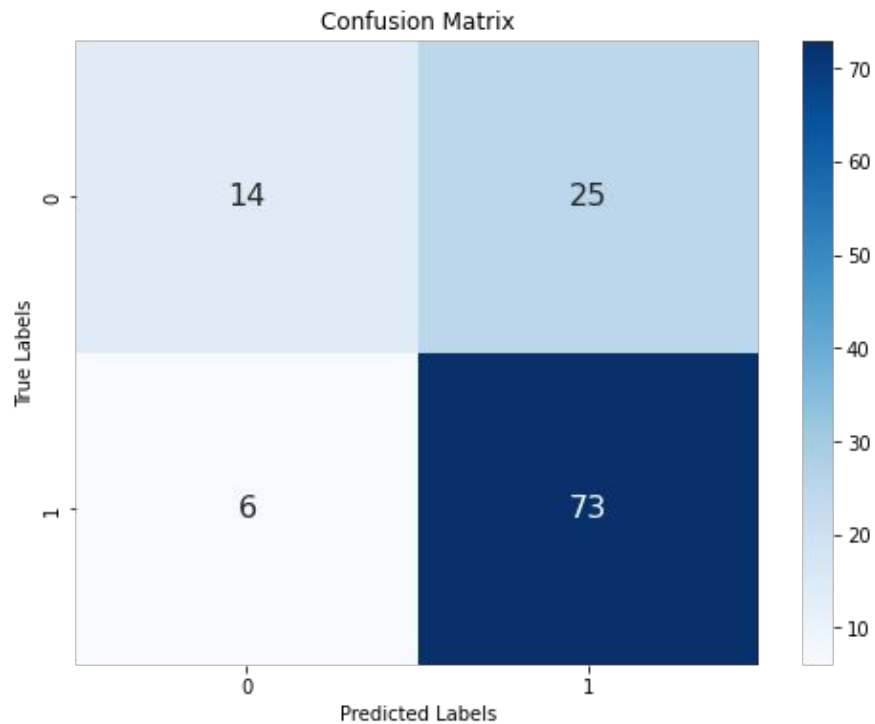


# Gradient Boosting





# Gradient Boosting





# Gradient Boosting

Accuracy of Gradient Boosting on Test Set: 0.737



# Best Model

Model	Accuracy on test dataset
SVM	0.669
KNN	0.585
Logistic Regression	0.763
Decision Tree	0.796
Random Forest	0.771
Gradient Boosting	0.737



## Summary & Conclusion

- Decision tree performs the best when predicting the approval or denial of a loan application
- Important features differ between models:
  - Random Forest & Gradient Boosting: Credit History, Loan Amount, Applicant & Co Applicant Income
  - Logistic Regression: Credit History, Education, Marital Status
- Credit History matters



## **Future work based on weakness and limitations**

### Enhancing Data Quality:

Data Enrichment: Incorporate more diverse and comprehensive data sources to capture a wider range of variables that might affect loan approval decisions.

Handling Missing Data: Explore advanced imputation techniques or data augmentation methods to better handle missing data without introducing significant bias.



## **Future work based on weakness and limitations**

Advanced Modeling Techniques:

Ensemble Methods: Further explore ensemble methods that combine the predictions from multiple models to improve accuracy and robustness.





## **Future work based on weakness and limitations**

Model Validation and Testing:

Cross-Validation with Diverse Datasets: Test the model on various datasets to ensure its generalizability and robustness.

Real-World Testing: Pilot the model in a real-world setting to observe its performance and gather feedback.

**Thank you for  
listening.**

Yifan Chen, Xupeng Tang, Chris  
Yang, Jiajie Yao, Hongyan Zhao

