

# LOAN ELIGIBILITY PREDICTION

## PROBLEM STATEMENT

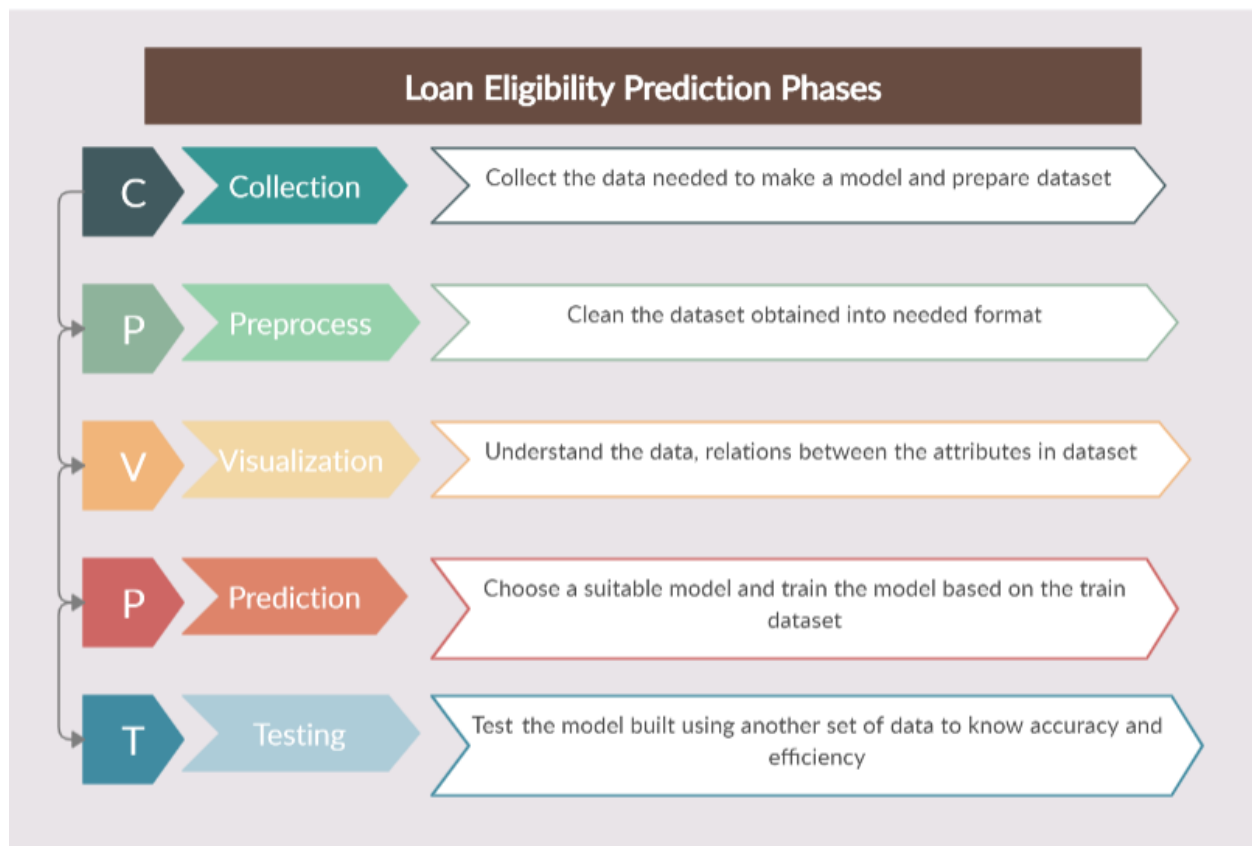
Loans are the core business of banks. The main profit comes directly from the loan's interest. The loan companies grant a loan after an intensive process of verification and validation. However, they still don't have assurance if the applicant is able to repay the loan with no difficulties.

## SOLUTION PROPOSED

Build a predictive model to predict using machine learning if an applicant can repay the lending company or not.

## SOLUTION USING MACHINE LEARNING

### FLOW OF PROJECT



## STEP 1: Analysing data

After collecting dataset, analyse the dataset. Some of the details to be noticed in the dataset:

- Shape of dataset (number of columns and rows present in dataset)
- Understand columns present in the dataset.
- Analyse the structure of the dataset.
- Check if missing values are present in the dataset.

Structure of train_data.csv															
Index	Column Name	Column Type													
0	Loan_ID	STRING													
1	Gender	STRING													
2	Married	STRING													
3	Dependents	STRING													
4	Education	STRING													
5	Self_Employed	STRING													
6	ApplicantIncome	INTEGER													
7	CoapplicantIncome	DOUBLE													
8	LoanAmount	INTEGER													
9	Loan_Amount_Term	INTEGER													
10	Credit_History	INTEGER													
11	Property_Area	STRING													
12	Loan_Status	BOOLEAN													

Summary of data															
Summary	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount		LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
Count	614	614	614	614	614	614	614	614	614	4	614	614	614	614	614
Unique	614	3	3	5	2	3								3	1
Top Freq	LP002888	Male	Yes	0	Graduate	No								Semiurban	
Top Freq	1	489	398	345	480	500	3317724	995444.91998864		4				233	
sum							5403.4592833876195	1621.2457980270997		7					
Mean							150	0		9					
Min							81000	41667	700	7	9	12	0		
Max							80050	41667	691	7	700	480	1		
Range							37320390.167181246	8562029.518387228		8		468	1		
Variance							6109.041673387181	2926.2483692241894		4					
Std. Dev															192
false															422
true															

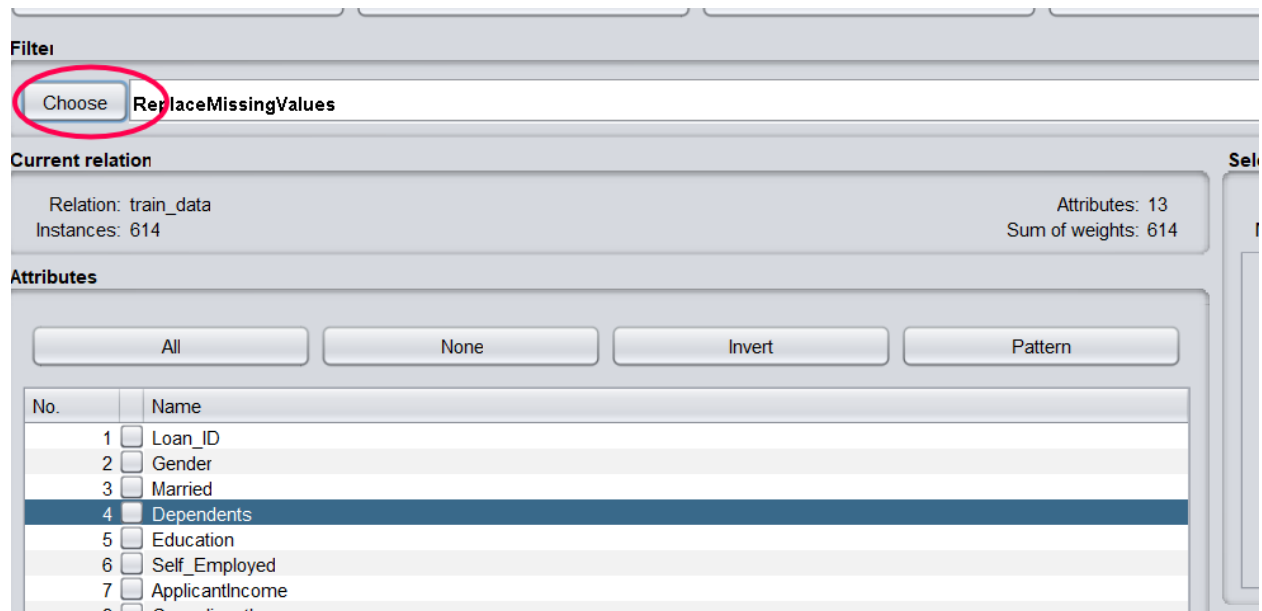
## STEP 2: Preprocess data

It is important to remove/fill missing values, reducing noisy data, handling outliers, and removing duplicates to get accurate prediction results. If data is not preprocessed, the prediction results will be deviated and corrupted.

- Can be done using Java or on Weka.

Using Weka:

- Open the file that requires to be preprocessed.
- Choose the filter to be applied such as "ReplaceMissingValues" to replace all the missing values with mode and mean of that columns.
- To remove duplicates use "RemoveDuplicates" filter.
- If you feel any attribute is not needed during prediction or deviates the result, use "RemoveAttribute" filter and delete those columns.
- Save the edited file to use it for prediction.

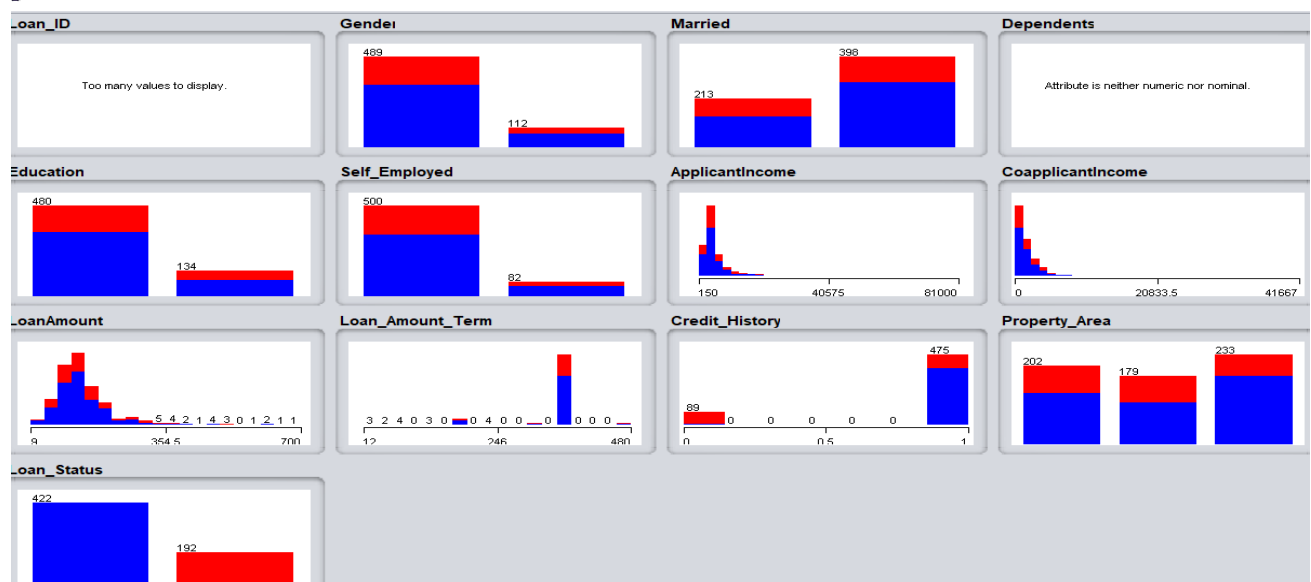


### STEP 3: Visualization

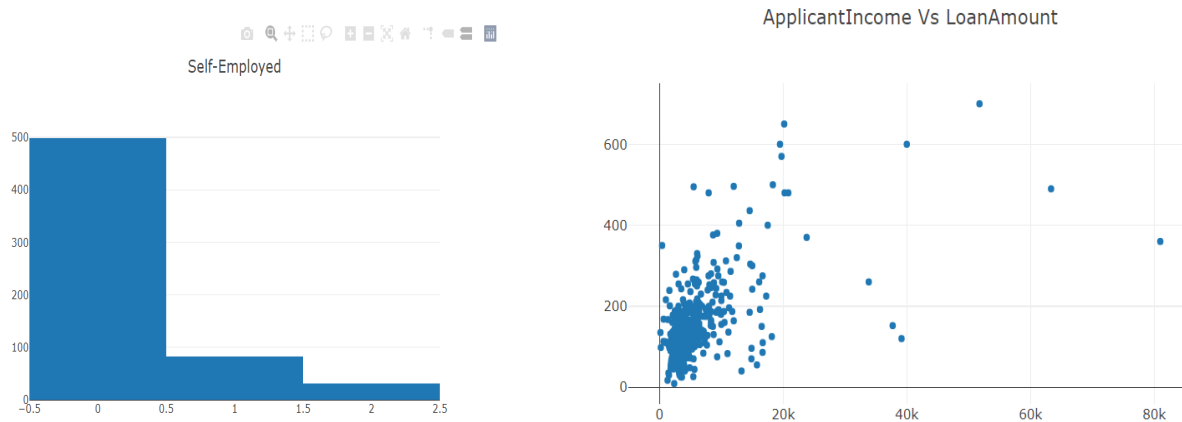
Understand the relationships between different attributes or columns.

- Can be done using Java or on Weka.
- Plot any graph (scatter plot, bar plot, histogram etc.) and understand how data is classified.

Using Weka:



## Using Java:



## STEP 4: Prediction and Testing

- Choose a model that you think will give accurate results.
- Construct and train the model using training data.
- Pass the testing data to the trained model to check prediction accuracy.

Prediction using Weka: (i took logistic regression)

Choose **Logistic -R 1.0E-8 -M -1 -num-decimal-places 4**

**Test options**

- ☐ Use training set
- ☐ Supplied test set
- ☒ Cross-validation Folds
- ☐ Percentage split %

om) Loan\_Status

**Result list (right-click for options)**

- 14:00:48 - rules.ZeroR
- 14:02:10 - trees.DecisionStump
- 14:02:51 - functions.Logistic

**Classifier output**

Credit\_History 50.2967  
Property\_Area=Urban 0.8084  
Property\_Area=Rural 0.6987  
Property\_Area=Semiurban 1.6717

Time taken to build model: 0.12 seconds

=== Stratified cross-validation ===  
=== Summary ===

Correctly Classified Instances	496	80.7818 %
Incorrectly Classified Instances	118	19.2182 %
Kappa statistic	0.4772	
Mean absolute error	0.296	
Root mean squared error	0.3892	
Relative absolute error	68.8135 %	
Root relative squared error	83.9458 %	
Total Number of Instances	614	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	FRC Area	Class
	0.981	0.573	0.790	0.981	0.875	0.535	0.759	0.837	true
	0.427	0.019	0.911	0.427	0.582	0.535	0.759	0.677	false
Weighted Avg.	0.808	0.400	0.828	0.808	0.783	0.535	0.759	0.787	

=== Confusion Matrix ===

```
a  b  <-- classified as
414  8 | a = true
110  82 | b = false
```

## Using Java: (i took logistic regression)

---

```
614
Logistic Regression on loan prediction data:
summary :
Correctly Classified Instances      496           80.7818 %
Incorrectly Classified Instances    118           19.2182 %
Kappa statistic                     0.4772
Mean absolute error                 0.2911
Root mean squared error             0.3815
Relative absolute error             67.6867 %
Root relative squared error         82.3008 %
Total Number of Instances          614

Confusion matrix:
[414.0, 8.0]
[110.0, 82.0]

Area under the curve
0.7898030213270142
0.7898030213270142
[Correct, Incorrect, Kappa, Total cost, Average cost, KB relative, KB information, Correlation, Complexity 0, Complexity scheme, Complexity improvement, MAE, RMSE, R

Recall :
0.43
0.43

Precision:
0.91
0.91

F1 score:0.58
0.58

Accuracy:0.81
0.81
```

## RESULT AND CONCLUSION

Successfully created a logistic model to predict loan eligibility given an instance with accuracy of 81% and 91% precision with the help of machine learning through Java and Weka.