# DMBI Lab

## EXPERIMENT NO. 5

**Name :** Kaushik Kotian      **Roll no.** 30      **Div :** D15B      **Batch:**B


**AIM :** To explore Rapid Miner and implement classification models like Decision Tree and Naive Bayes.

1. Preprocess data. Split data into train and test set

2. Build Classification model using Rapid miner on training data (use decision tree and naive Bayes method)

3. Calculate metrics based on test data

4. Compare the models based on metrics and find which model is best suited for the dataset.

**Link to the dataset :**
https://www.kaggle.com/datasets/bhavikjikadara/loan-status-prediction?resource=downl oad

**THEORY :**

**RapidMiner** is a comprehensive open-source platform for data science and machine learning, providing a user-friendly interface for building, deploying, and maintaining predictive models. It offers a visual workflow environment where users can design and execute data analysis tasks without needing extensive programming knowledge. RapidMiner supports various data preprocessing techniques, feature engineering, model training using diverse algorithms such as decision trees and neural networks, and model evaluation with metrics like accuracy and F1-score. It also enables model deployment for making predictions on new data or integrating models into production systems. With its intuitive interface and robust functionality, RapidMiner empowers data scientists and analysts to explore data, create powerful machine learning models, and derive valuable insights from complex datasets.

A **classification model** is a type of machine learning model that is used to categorize input data into predefined classes or categories. The goal of a classification model is to

learn the mapping between input features and the corresponding output labels, enabling it to accurately predict the class of new, unseen data points.

Key components of a classification model include:

**1. Input Data:** The data used to train and test the classification model typically consists of features or attributes that describe the input, along with the corresponding output labels or classes.
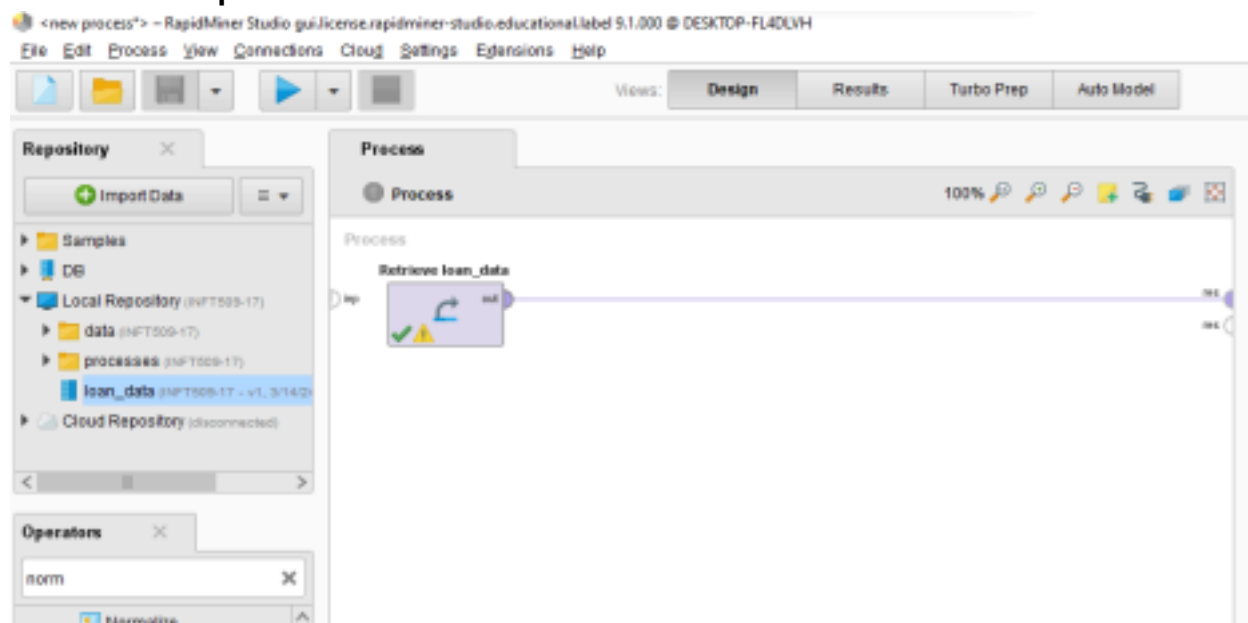
**2. Training Phase:** During the training phase, the classification model learns from the input data by adjusting its internal parameters based on the provided input-output pairs. Different algorithms, such as decision trees, logistic regression, support vector machines (SVM), k-nearest neighbors (KNN), and neural networks, can be used for training classification models.

**3. Prediction:** Once the model is trained, it can be used to predict the class labels of new, unseen data points. The model applies the learned mapping to the input features of the new data and outputs the predicted class or probability distribution over classes.

**4. Evaluation:** The performance of a classification model is evaluated using metrics such as accuracy, precision, recall, F1-score, and area under the receiver operating characteristic curve (ROC-AUC). These metrics assess how well the model generalizes to unseen data and its ability to correctly classify instances into their respective classes.

**Steps :**

**1. Data Preparation**





ExampleSet (301 examples, 0 special attributes, 13 regular attributes)

**2. Filter the data to eliminate missing values**

**3.**

**Check the statistics tab in results, and eliminate missing values from entire dataset**



**After removing all missing values:**

## 4. Building decision tree :
### Set label attribute



### Set criterion (information gain/ gini index) :

## Numeric description of the decision tree

Graph

Description

Annotations

# Tree

```
Credit_History = ?
|   ApplicantIncome > 1937
|   |   Dependents = 0: Y {N=0, Y=17}
|   |   Dependents = 1: Y {N=1, Y=3}
|   |   Dependents = 2: Y {N=0, Y=5}
|   |   Dependents = 3+: N {N=1, Y=1}
|   ApplicantIncome ≤ 1937: N {N=2, Y=0}
Credit_History > 0.500
|   Property_Area = Rural
|   |   CoapplicantIncome > 1563.500
|   |   |   LoanAmount > 113.500
|   |   |   |   ApplicantIncome > 1905.500
|   |   |   |   |   LoanAmount > 133
|   |   |   |   |   |   LoanAmount > 137: Y {N=0, Y=5}
|   |   |   |   |   |   LoanAmount ≤ 137: N {N=1, Y=1}
|   |   |   |   |   LoanAmount ≤ 133: Y {N=0, Y=12}
|   |   |   |   ApplicantIncome ≤ 1905.500: N {N=2, Y=0}
|   |   |   LoanAmount ≤ 113.500: Y {N=0, Y=14}
|   |   CoapplicantIncome ≤ 1563.500
|   |   |   LoanAmount > 135.500: Y {N=0, Y=5}
|   |   |   LoanAmount ≤ 135.500
|   |   |   |   LoanAmount > 132.500: N {N=3, Y=0}
|   |   |   |   LoanAmount ≤ 132.500
|   |   |   |   |   ApplicantIncome > 4788.500: Y {N=1, Y=7}
|   |   |   |   |   ApplicantIncome ≤ 4788.500
|   |   |   |   |   |   CoapplicantIncome > 1248.500: N {N=4, Y=0}
```

**Conclusion :** Hence we have understood and implemented a classification model using Rapid Miner software and performed data pre processing , data cleaning, and built the decision tree.