# Experiment 5

## Aim:

Experiment to explore Rapid Miner and implement classification models like Decision Tree and Naive Bayes etc.

## Theory:

In machine learning, a decision tree is a predictive model that maps observations about an item to conclusions about its target value. It is a tree-structured model where each internal node represents a test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label. Decision trees can be used for both regression and classification tasks, and they are widely used in various applications, such as finance, medicine, and engineering. The goal of building a decision tree is to create a model that predicts the value of a target variable based on several input variables. The decision tree algorithm uses a recursive process to split the data into smaller subsets based on the input variables until it reaches a point where it can make a prediction.

The basic steps involved in creating a decision tree are as follows:

1. Collect and prepare data: This involves collecting and organizing the data, and then preparing it for analysis. This may include data cleaning, data transformation, and feature selection.
2. Choose an algorithm: There are various algorithms available for building decision trees, including ID3, C4.5, CART, and CHAID. The choice of algorithm will depend on the specific problem and the characteristics of the data.
3. Build the tree: This involves applying the chosen algorithm to the data to create the decision tree. The tree is built by recursively partitioning the data into subsets based on the values of the input features, and selecting the feature that provides the most information gain at each step.
4. Evaluate the tree: Once the tree is built, it needs to be evaluated to determine its accuracy and effectiveness. This may involve using cross-validation or other techniques to estimate the performance of the tree on new data.
5. Use the tree: Finally, the decision tree can be used to make predictions on new data. This involves traversing the tree from the root to a leaf node based on the values of the input features, and outputting the corresponding class label or value.

   **1. Gini Index-**

Gini index is a measure of impurity or inequality used in decision trees to determine the purity of a given split in the data. It ranges from 0 to 1, where 0 represents a completely pure split (all observations belong to the same class) and 1 represents a completely impure split (an equal number of observations belong to each class).

The formula for calculating the Gini index is:

Gini = 1 - (p_1)^2 - (p_2)^2 - ... - (p_k)^2

where k is the number of classes, and p_i is the proportion of observations belonging to class i in the split.

## 2. Information Gain-

Information Gain is a measure used in decision tree algorithms to determine the relevance of a feature to a target variable. It measures the reduction in entropy or degree of disorder in the target variable when a feature is used to split the data into subsets. The formula for Information Gain is:

Information Gain = Entropy(parent) - [Weighted Avg. * Entropy(children)]
where,
Entropy(parent) is the entropy of the target variable for the entire dataset
Weighted Avg. is the weighted average of the entropy for each child node
Entropy(children) is the entropy of the target variable for each child node

## 3. Entropy-

In the context of machine learning and decision trees, entropy is a measure of impurity or disorder within a set of examples. It is used to quantify the amount of uncertainty or randomness in a set of data. A dataset with low entropy is considered more uniform and has less randomness, whereas a dataset with high entropy is considered more disordered and has more randomness. The entropy of a dataset S with respect to a binary target variable Y can be calculated using the following formula:
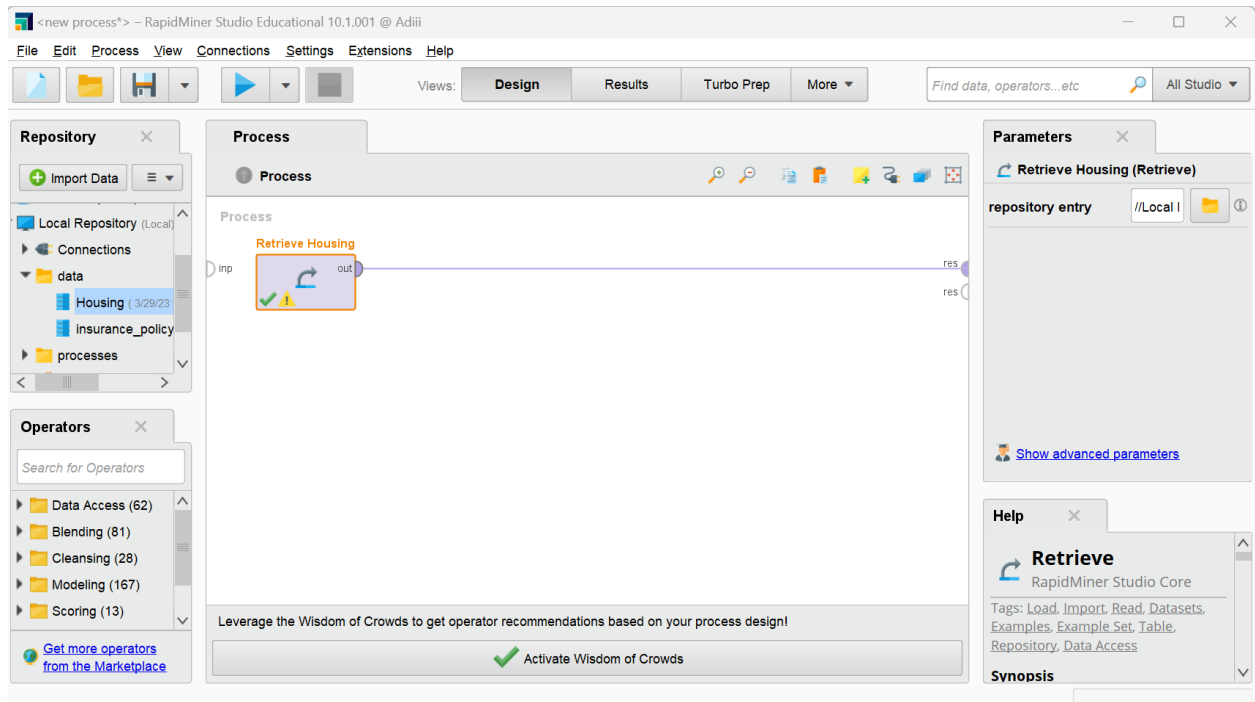
Entropy(S) = -p_1 \log_2 p_1 - p_2 \log_2 p_2$
where p1 is the proportion of examples in S that belong to class 1, and p2 is the proportion of examples in S that belong to class 2. The logarithm base 2 is used to measure the entropy in bits.

# Implementation:

Decision tree

## Create Filters: filters

**Create Filters: filters**
Defines the list of filters to apply.

| hotwaterheating ▼ | contains ▼ | no |

● Match all   ○ Match any   ☑ Preselect comparators   [Add Entry]   [✓ OK]   [✗ Cancel]

---

<new process*> – RapidMiner Studio Educational 10.1.001 @ Adiii

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   **Results**   Turbo Prep   More ▼   Find data, operators...etc   All Studio ▼

Result History | **ExampleSet (Filter Examples)** ✕ | ExampleSet (//Local Repository/data/Housing) ✕ | Repository ✕

Open in   [Turbo Prep]   [Auto Model]   Filter (520 / 520 examples): all

| lrooms | bathrooms | stories | mainroad | guestroom | basement | hotwaterhe... | airconditioni |
|--------|-----------|---------|----------|-----------|----------|---------------|---------------|
| | 4 | 4 | yes | no | no | no | yes |
| | 2 | 2 | yes | no | yes | no | no |
| | 2 | 2 | yes | no | yes | no | yes |
| | 1 | 2 | yes | yes | yes | no | yes |
| | 3 | 1 | yes | no | yes | no | yes |
| | 3 | 4 | yes | no | no | no | yes |
| | 3 | 2 | yes | no | no | no | no |
| | 1 | 2 | yes | yes | yes | no | yes |
| | 2 | 4 | yes | yes | no | no | yes |
| | 1 | 2 | yes | no | yes | no | yes |
| | 2 | 2 | yes | no | no | no | yes |

ExampleSet (520 examples, 0 special attributes, 13 regular attributes)

**Repository**

[⊕ Import Data]   ≡ ▼

▷ 🖥 Training Resources (connected)
▷ 📁 Samples
▷ 👥 Community Samples (connected)
▽ 🖥 Local Repository (Local)
  ▷ 🔌 Connections
  ▽ 📁 data
    📄 Housing ( 3/29/23 9:04 PM – 35 k
    📄 insurance_policy ( 3/29/23 8:50
  ▷ 📁 processes
    ⚙ insurance ( 3/29/23 9:03 PM – 1 kB)
▷ 🖥 DB (Legacy)

Data
Statistics
Visualizations
Annotations

File  Edit  Process  View  Connections  Settings  Extensions  Help

Views:  Design  Results  Turbo Prep  Auto Model

Find data, operators...etc   All Studio ▾

**Repository**

Import Data

▸ Connections
▾ data
   Housing ( 3/29/23 9:04
   insurance_policy ( 3/2
▾ processes
   Decision tree ( 3/30/23
   insurance ( 3/29/23 9:03 P
   naivebayes ( 3/30/23 5:20

**Operators**

performance

▾ Performance (18)
   ▾ Predictive (7)
      Performance (C
      Performance (E
      Performance (F

We found "Model Management" in the Marketplace. Show me!

**Process**

Process ▸

Process

Retrieve Housing    Filter Examples                    Decision Tree

inp    out    exa    exa            tra    mod
                      ori                   exa
                      unm                   wei

Set Role

exa    exa
       ori

Performance

lab    per
per    exa

Apply Model

mod    lab
unl    mod

**Parameters**

Set Role

attribute name    airconditioning ▾

target role    label ▾

set additional roles    Edit List (0)...

Show advanced parameters

✓ Change compatibility (10.1.001)

**Help**

? **Set Role**
RapidMiner Studio Core

Tags: Label, Target, Id, Class, Dependent, Independent, Special, Regular, Inputs, Columns, Attributes, Features, Variables, Types, Deprecated

Synopsis

Leverage the Wisdom of Crowds to get operator recommendations based on your process design!

✓ Activate Wisdom of Crowds

---

PerformanceVector (Performance)    ExampleSet (//Local Repository/data/Housing)

Result History    Tree (Decision Tree)

**Repository**

Import Data

▸ Training Resources (connected)
▸ Samples
▸ Community Samples (connected)
▾ Local Repository (Local)
   ▸ Connections
   ▾ data
      Housing ( 3/29/23 9:04 PM – 35 kB)
      insurance_policy ( 3/29/23 8:50 PM – 3.4 MB)
   ▾ processes
      Decision tree ( 3/30/23 3:47 PM – 3 kB)
      insurance ( 3/29/23 9:03 PM – 1 kB)
      naivebayes ( 3/30/23 5:20 PM – 7 kB)
▸ DB (Legacy)

Zoom

Graph

Tree ▾

Description

☑ Node Labels

☑ Edge Labels

Annotations

price

> 6142500          ≤ 6142500

area                          stories

> 11307.500  ≤ 11307.500      > 3.500    ≤ 3.500

no        bedrooms        price        no

      > 2.500  ≤ 2.500    > 579250  ≤ 5792500

yes        no        no        yes

# Decision Tree

price
> 6142500 — area
≤ 6142500 — stories

area
> 11307.500 — no
≤ 11307.500 — bedrooms

stories
> 3.500 — price
≤ 3.500 — no

bedrooms
> 2.500 — yes
≤ 2.500 — no

price
> 5792500 — no
≤ 5792500 — yes

---

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   Results   Turbo Prep   Auto Model

Find data, operators...etc      All Studio

PerformanceVector (Performance)     ExampleSet (//Local Repository/data/Housing)

Result History     Tree (Decision Tree)

Criterion
accuracy

Performance

Description

Annotations

● Table View   ○ Plot View

accuracy: 77.25%

| | true yes | true no | class precision |
|---|---|---|---|
| pred. yes | 77 | 29 | 72.64% |
| pred. no | 95 | 344 | 78.36% |
| class recall | 44.77% | 92.23% | |

Repository

Import Data

Training Resources (connected)
Samples
Community Samples (connected)
Local Repository (Local)
   Connections
   data
      Housing ( 3/29/23 9:04 PM – 35 kB)
      insurance_policy ( 3/29/23 8:50 PM – 3.4 MB)
   processes
      Decision tree ( 3/30/23 3:47 PM – 3 kB)
      insurance ( 3/29/23 9:03 PM – 1 kB)
      naivebayes ( 3/30/23 5:20 PM – 7 kB)
DB (Legacy)

## Naive Bayes

**ExampleSet (Apply Model) — Data view**

Filter (163 / 163 examples): all

| Row No. | hotwaterhe... | prediction(h... | confidence(... | confidence(... | price | area | bedrooms | ba |
|---|---|---|---|---|---|---|---|---|
| 1 | no | yes | 0.195 | 0.805 | 12215000 | 7500 | 4 | 2 |
| 2 | no | yes | 0.467 | 0.533 | 11410000 | 7420 | 4 | 1 |
| 3 | no | no | 0.770 | 0.230 | 9870000 | 8100 | 4 | 1 |
| 4 | no | no | 0.602 | 0.398 | 9240000 | 7800 | 3 | 2 |
| 5 | no | yes | 0.299 | 0.701 | 9100000 | 6000 | 4 | 1 |
| 6 | no | no | 0.868 | 0.132 | 8890000 | 4600 | 3 | 2 |
| 7 | no | no | 0.965 | 0.035 | 8645000 | 8050 | 3 | 1 |
| 8 | no | no | 0.914 | 0.086 | 8645000 | 4560 | 3 | 2 |
| 9 | no | no | 0.771 | 0.229 | 8575000 | 8800 | 3 | 2 |
| 10 | no | no | 0.990 | 0.010 | 7980000 | 9000 | 4 | 2 |
| 11 | no | no | 0.998 | 0.002 | 7560000 | 6000 | 4 | 2 |
| 12 | no | no | 0.987 | 0.013 | 7420000 | 7440 | 3 | 2 |

ExampleSet (163 examples, 4 special attributes, 12 regular attributes)

**PerformanceVector (Performance)**

Criterion: accuracy

accuracy: 92.02%

| | true no | true yes | class precision |
|---|---|---|---|
| pred. no | 150 | 7 | 95.54% |
| pred. yes | 6 | 0 | 0.00% |
| class recall | 96.15% | 0.00% | |

## Tab bar

SimpleDistribution (Naive Bayes) ✕ | ExampleSet (Split Data) ✕ | ExampleSet (//L...

Result History | ExampleSet (Apply Model) ✕ | % Performanc...

## Performance panel

**Performance** | **Description** | **Annotations**

# PerformanceVector

```
PerformanceVector:
accuracy: 92.02%
ConfusionMatrix:
True:     no        yes
no:       150       7
yes:      6         0
```

## RapidMiner Studio window

`<new process*> – RapidMiner Studio Educational 10.1.001 @ Adiii`

File   Edit   Process   View   Connections   Settings   Extensions   Help

Views:   Design   Results   Turbo Prep   Auto Model

Find data, operators...etc   All Studio ▾

**Auto Model**

Load Data — Select Task — Prepare Target — Select Inputs — Model Types — Results

≪ RESTART   ‹ BACK   ⤓ OPEN PROCESS   ⬇ EXPORT

### Results

- ▾ Comparison
  - Overview
  - ROC Comparison
- ▾ Naive Bayes
  - Model
  - Weights
  - Simulator
  - Performance
  - Lift Chart

SAVE RESULTS

### Overview

Number of Models: **21**

**Accuracy** | **Runtimes (ms)**

| Model | | Accuracy | Standard Deviation | Gains | Total Time | |
|---|---|---|---|---|---|---|
| Naive Bayes | ♀ $ 🏃 | 95.5% | ± 2.9% | 0 | 1 s | |
| Decision Tree | 🏃 | 92.9% | ± 5.3% | -8 | 1 s | |

## Naive bayes

Decision tree



## Conclusion-

We have successfully explored Rapid Miner and implemented classification models like Decision Tree and Naive Bayes.