**Problem Statement**

This dataset is about **Human Heart Disease dataset**. The objective is to predict

whether the patient has disease or not.

**About dataset:**

This data set dates from 1988 and consists of four databases: Cleveland, Hungary, Switzerland, and Long Beach V. It contains 76 attributes, including the predicted attribute, but all published experiments refer to using a subset of 14 of them. The "target" field refers to the presence of heart disease in the patient. It is integer valued 0 = no disease and 1 = disease.

**Data Field description:**

age

sex

chest pain type (4 values)

resting blood pressure

serum cholestoral in mg/dl

fasting blood sugar > 120 mg/dl

resting electrocardiographic results (values 0,1,2)

maximum heart rate achieved

exercise induced angina

oldpeak = ST depression induced by exercise relative to rest

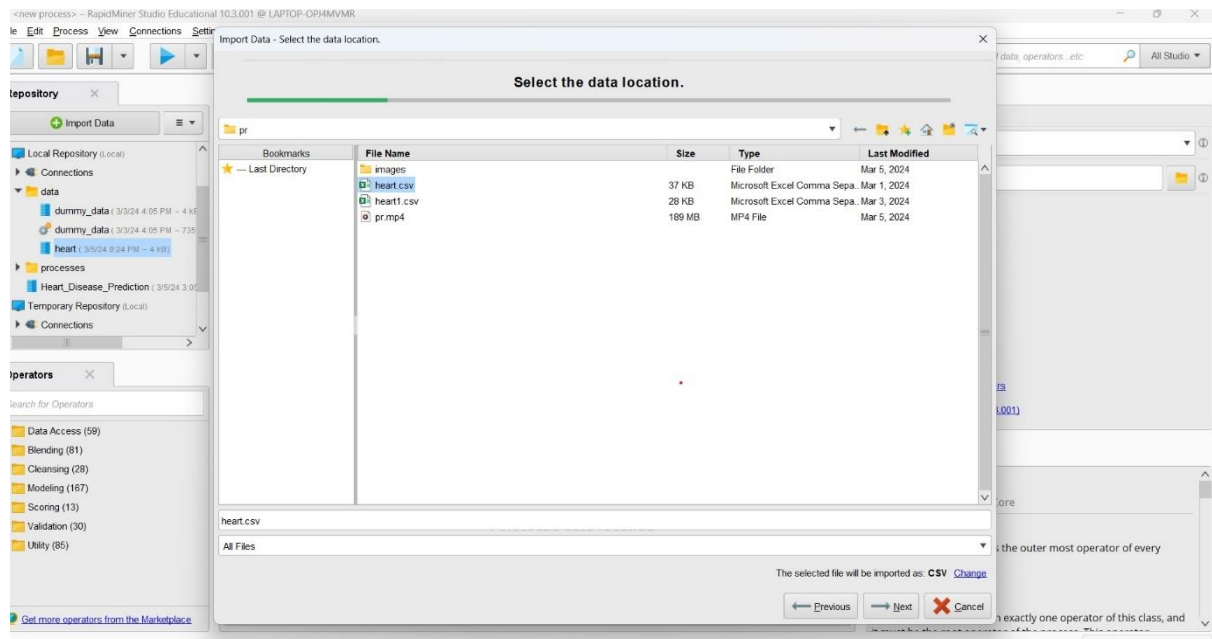the slope of the peak exercise ST segment

number of major vessels (0-3) colored by flourosopy

thal: 0 = normal; 1 = fixed defect; 2 = reversable defect

The names and social security numbers of the patients were recently removed from the database, replaced with dummy values.

**Steps to convert dataset to prediction:**

1. **Load data:** Import dataset from your computer to local repository.
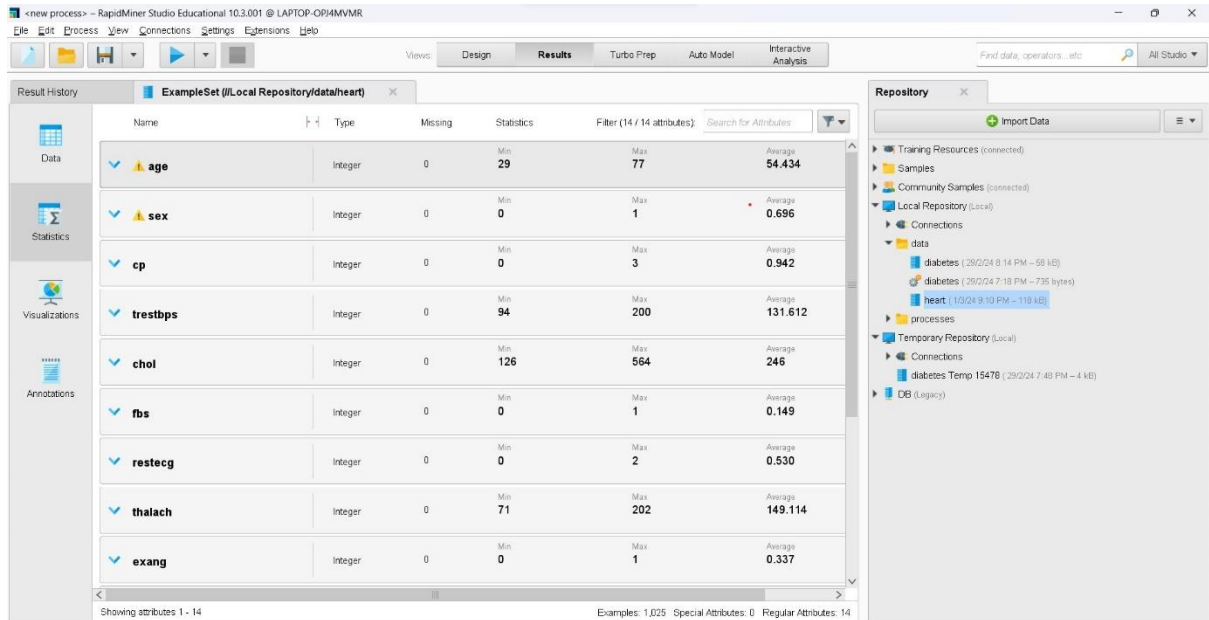
## 2. Exploring Dataset size:



From above figure there are 1025 examples/rows/records and 14 columns/attributes.

Out of 14 columns, 13 columns are predictors and remaining column is target.

### 3. Data cleanup:

a. In this dataset there is no missing values.



b. Drop attributes if required.

We can drop attributes if required. In this dataset there is no need of dropping the attribute.

c. Target selection:

d. Prepare target:

Here target is to predict the target attribute.

e. Select inputs:



f. Select the model:

Out of all models Random Forest has given best performance and best gain.

## 4. Baseline Model:

 Model is more correlated on the oldpeak column and it has given almost equal preference to remaining attributes to predict outcome.

Out of all models Random Forest has given max accuracy about 89.4% with classification error of 10.6%.

Random Forest - Production Model



Random Forest - Lift Chart

**RapidMiner Studio Auto Model — Statistics (target)**

## Results

- Model
- Weights
- Simulator
- Performance
- Lift Chart
- Optimal Parameters
- Predictions
- Production Model
- ▶ Support Vector Machine
- ▼ General
  - Data
  - Statistics
  - Weights by Correlation
  - Correlations

SAVE RESULTS

## Statistics

‹ › **target**

**Summary**

Category

Missing: 0.00%
Infinite: 0.00%
ID-ness: 0.20%
Stability: 51.32%

**Top Values**

**2 Distinct Values:**

| Value | Count | Percentage |
| --- | --- | --- |
| range2 | 526 | 51.32% |
| range1 | 499 | 48.68% |



**RapidMiner Studio Auto Model — Weights by Correlation**

## Weights by Correlation

| Attribute | Weight |
| --- | --- |
| oldpeak | 0.438 |
| exang | 0.438 |
| cp | 0.435 |
| thalach | 0.423 |
| ca | 0.382 |
| slope | 0.346 |
| thal | 0.338 |
| sex | 0.280 |
| age | 0.229 |
| trestbps | 0.139 |
| restecg | 0.134 |
| chol | 0.100 |
| fbs | 0.044 |

## 5. Correlation matrix:



## Output for random input:

**Conclusion:**

Out of all the classifiers tested Random Forest has gave better results in predicting outcome

The conclusion for the Human Heart Disease dataset depends on the specific analysis and modeling performed on it. However, here are some potential conclusions that could be drawn based on the information provided:

Predictive Modeling: Researchers could build predictive models using machine learning algorithms to predict the presence or absence of heart disease based on the provided attributes. These models could be evaluated based on their accuracy, sensitivity, specificity, and other performance metrics.

Feature Selection: Since the dataset contains 76 attributes, but most experiments focus on a subset of 14 attributes, researchers could perform feature selection to identify the most relevant attributes for predicting heart disease. This can help in simplifying models and improving their interpretability.

Risk Factors Identification: By analyzing the relationships between the attributes and the presence of heart disease, researchers can identify potential risk factors or indicators of heart disease. This information can be valuable for healthcare professionals in understanding and managing the disease.

Dataset Bias and Generalization: Researchers should also consider potential biases in the dataset, such as overrepresentation of certain demographics or medical conditions in specific databases. Ensuring that the models generalize well to diverse populations is crucial for their applicability in real-world settings.

Further Research Directions: Depending on the findings from the initial analysis, researchers may identify areas for further investigation, such as exploring interactions between different attributes, investigating the impact of lifestyle factors on heart disease risk, or evaluating the effectiveness of different treatment strategies.

Overall, the conclusion drawn from the dataset would depend on the specific analyses performed and the goals of the research project.