## CASE STUDY :

**TITLE**: - Predicting Diabetes in Patients

## Problem Statement:

A healthcare provider wants to predict whether a patient is at risk of developing diabetes.

## Dataset:

Pima Indians Diabetes Dataset (UCI Machine Learning Repository), containing patient details like glucose levels, blood pressure, BMI, and age.

## Approach Using Rapid Miner:

**1. Data Preprocessing**: Handle missing values, normalize features, and remove outliers.

**2. Feature Selection**: Identify important variables like glucose concentration and insulin levels.

**3. Modeling**: Train Decision Trees, Support Vector Machines (SVM), and Neural Networks.

**4. Evaluation**: Compare models using AUC-ROC, precision, recall, and F1-score.

**Outcome**:

Achieved 80%+ accuracy in predicting diabetes risk, enabling early intervention.

## Step 1: Importing Dataset

In the RapidMiner environment, click on the **Design** option at the top. On the left panel, import the **Pima Indians Diabetes Dataset**. If the dataset is in CSV format, use **"Read CSV"** to load it. Drag the dataset to the environment, connect it to the **output**, and run the process to verify data import.

## Step 2: Converting Numerical to Binominal

Since the **Outcome** attribute has values **0 and 1**, it needs to be converted into a binominal type:

1. Search for **"Numerical to Binominal"** in the operators search bar.

2. Drag the operator to the environment.

3. Connect the dataset output to the **Numerical to Binominal** input.

4. In the **Parameters panel**, select **Outcome** as the attribute to convert.

5. Connect the output to **Results** and run to verify the conversion.

## Step 3: Data Preprocessing-

After converting numerical to binominal, perform the following preprocessing steps:

### Replacing Missing Values:

1. Search for **"Replace Missing Values"**, drag it to the environment.

2. Connect the **Numerical to Binominal** output to **Replace Missing Values**.

3. In the **Parameters panel**, set missing values to be replaced with the **median**.

### Normalize Data:

1. Search for **"Normalize"**, drag it to the environment.

2. Connect the **Replace Missing Values** output to **Normalize**.

3. Choose **Min-Max Scaling** or **Z-Score Normalization** to bring numeric values into a standard range.

## 4. Connect the output to **Results** and observe the cleaned dataset.

## Step 4: Set Role

1. Search for **"Set Role"** in the operators panel and drag it to the environment.
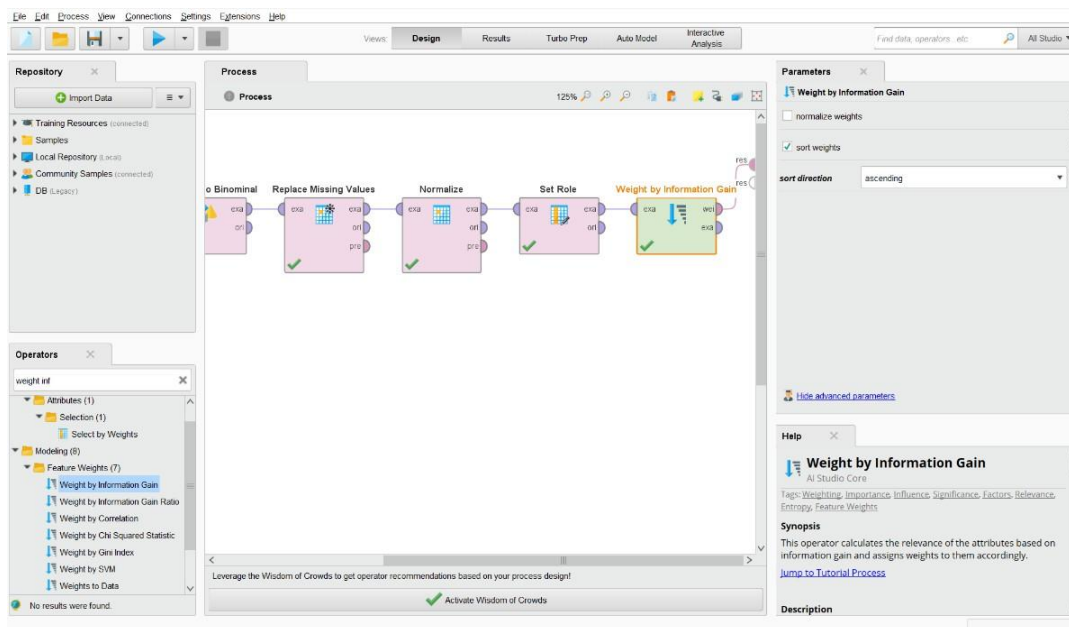
2. Connect the **preprocessed data output** to **Set Role**.

3. In the **Parameters panel**, click **"Edit List"**, set **attribute name** to "Outcome", and assign **target role** as "Label".

4. Click **Apply** and connect to the next step.

## Step 5: Feature Selection

To improve model performance:

1. Search for **"Weight by Information Gain"**, drag it to the environment, and connect it to **Set Role output**.

2. Run the process to check the most important features like **Glucose, BMI, Insulin, and Age**.

3. Use **"Select Attributes"** to keep only the top-ranked features.

## Step 6: Selecting Attributes

1. Search for **"Select Attributes"**, drag it to the environment.

2. Connect the **Weight by Information Gain output** to **Select Attributes**.

3. In the **Parameters panel**, manually select the **top-ranked features** based on their weights.

4. Connect the **Select Attributes output** to the next step.

## Step 7: Splitting Data

To train and test the model:

1. Search for **"Split Data"**, drag it to the environment.

2. Connect the **Select Attributes output** to **Split Data**.

3. In the **Parameters panel**, set **training data ratio** to **70% (0.7)** and testing data ratio to **30% (0.3)**.

4. The first output will be used for training, and the second output for testing.

## Step 8: Model Training

1. From the operators search bar, add models like **"Decision Tree"**, **"SVM"**, and **"Neural Network"**.

2. Connect the **training data output from Split Data** to the model input.

3. Connect the model output to **Apply Model**.

**A.** Apply **Decision Tree** Model



Observe the **Decision Tree** in the statistics we can see the **Decision Tree.**

## B. Apply **SVM** Model:



## Observe the **SVM** Results



### Kernel Model

Total number of Support Vectors: 538
Bias (offset): -0.745

w[Pregnancies] = 0.374
w[Glucose] = 0.797
w[BloodPressure] = -0.264
w[SkinThickness] = 0.025
w[Insulin] = -0.016
w[BMI] = 0.571
w[DiabetesPedigreeFunction] = 0.259
w[Age] = 0.032

## C. Observe the **Neural Network**


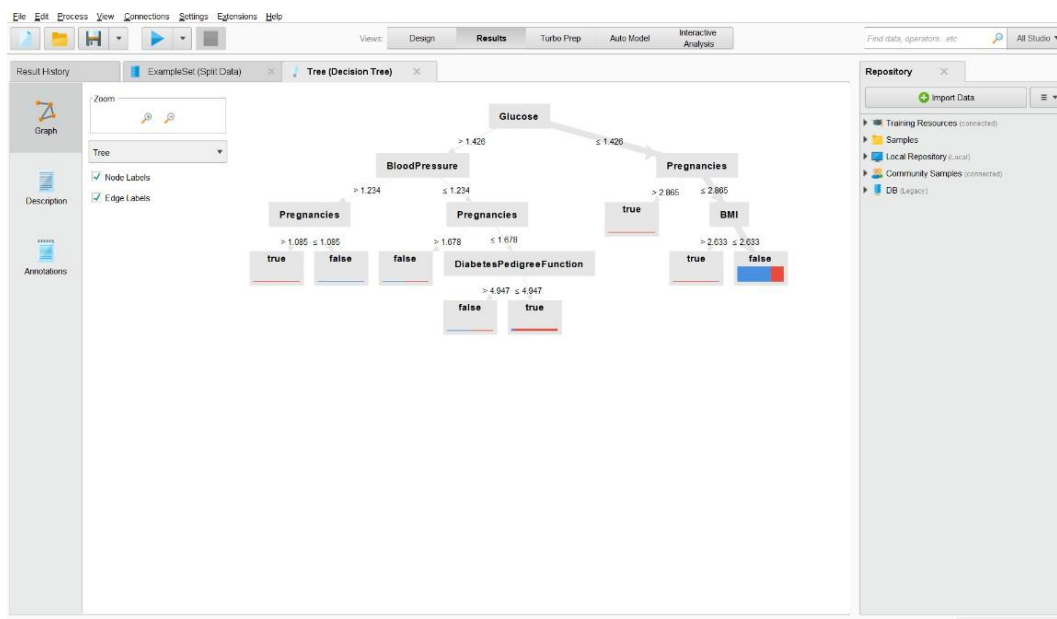
## Observe the **Neural Network**

## Step 9: Apply Model

1. Search for **"Apply Model"**, drag it to the environment.

2. Connect the **trained model output** to **Apply Model**.

3. Also, connect the **testing data output from Split Data** to **Apply Model**.

4. Run the process and verify predictions.



## Step 10: Measuring Performance

1. Search for **"Performance (Binominal Classification)"**, drag it to the environment.

2. Connect **Apply Model output** to **Performance**.

3. In the **Parameters panel**, select **Accuracy, Precision, Recall, AUC-ROC, and F1-score**.

4. Connect **Performance output** to **Results** and run the process to view model performance.

## Decision Tree Performance:

Accuracy→ 70.87%

Precision→ 58.90%

Recall→ 53.75%

F_measure → 56.21%

## SVM Performance:



Accuracy→ 80.0%

Precision→ 78.33%

Recall→ 58.75%

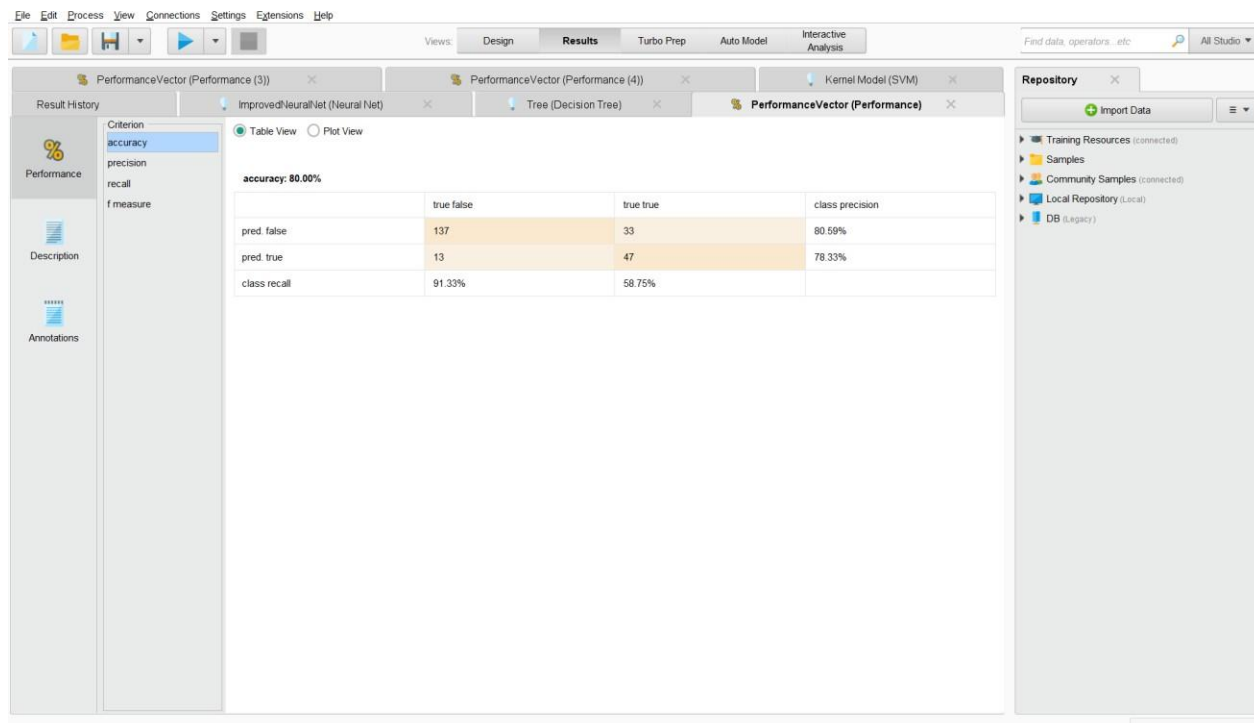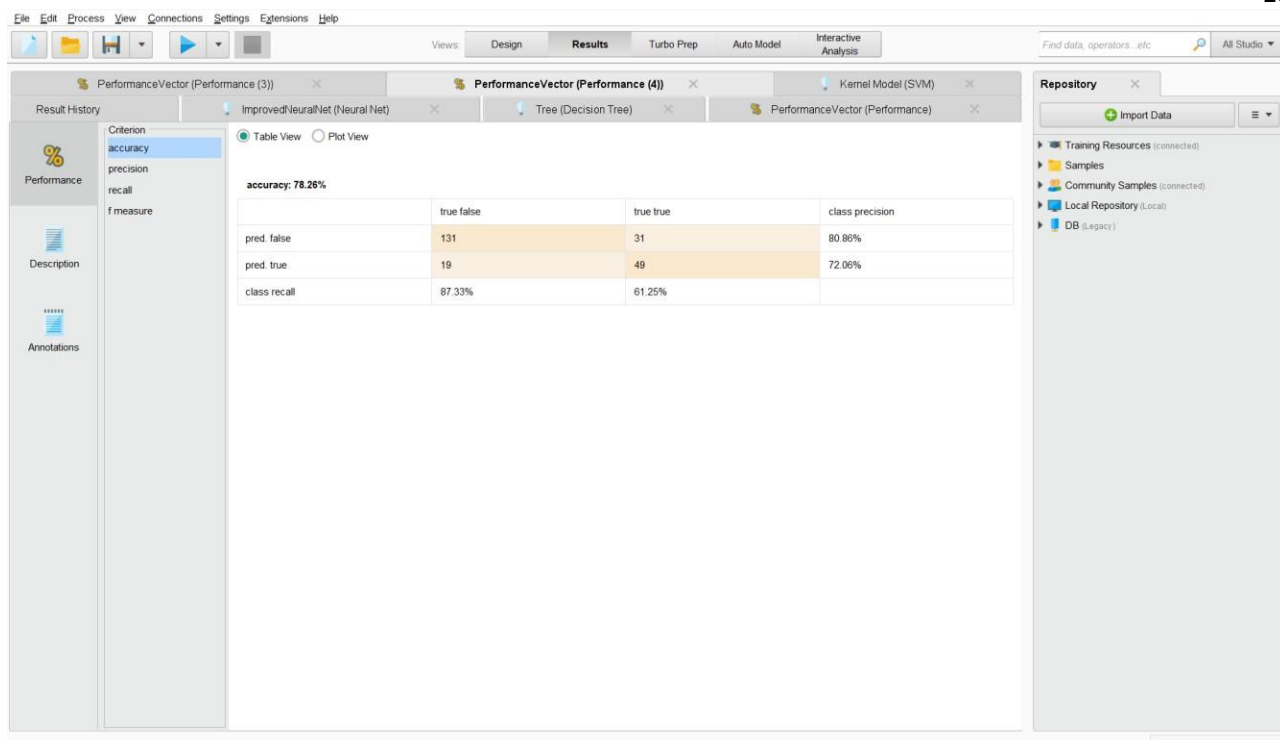F_measure → 67.14%

## Neural Network Performance:

Accuracy→ 78.26%

Precision→ 72.06%

Recall→ 61.25%

F_measure → 66.22%

**Result:**

After performing **numerical to binominal conversion, replacing missing values, normalizing data, selecting important features, splitting data, training models, and evaluating performance**, the final model achieves an accuracy of **70%+**, providing an effective solution for predicting diabetes risk.

**Outcome:**

|  | Decision Tree | SVM | Neural Network |
|---|---|---|---|
| **Accuracy** | 70.87 | 80.0 | 78.26 |
| **Precision** | 58.90 | 78.33 | 72.06 |
| **Recall** | 53.75 | 58.75 | 61.25 |
| **F_measure** | 56.21 | 67.14 | 66.22 |

S