# EXERCISE -9

(i) Implement Simple Linear Regression on an 'Employee' dataset.
(ii) Demonstrate the simple k-Means clustering algorithm on 'iris' dataset.

## 9.1.1 Problem Statement:

Implement Simple Linear Regression on an 'Employee' dataset.

## 9.1.2 Description:

About Dataset used:

Is Employee dataset with eight attributes

MarriedID

MaritalStatusID

GenderID

EmpStatusID      DeptID

PerfScoreID

FromDiversityJobFairID

Salary

Termd

About RapidMiner:

Rapidminer is a comprehensive data science platform with visual workflow design and full automation. It means that we don't have to do the coding for data mining tasks. Rapidminer is one of the most popular data science tools.

This is the graphical user interface of the blank process in rapidminer. It has the repository that holds our dataset. We can import our own datasets. It also offers many public datasets that we can try. We can also work with a database connection.

RapidMiner is an awesome visual workflow designer. The way they present visually is so unique. It helps in speeding and automating the creation of visual models. It helps in creating models in only 5 clicks by automated machine learning.
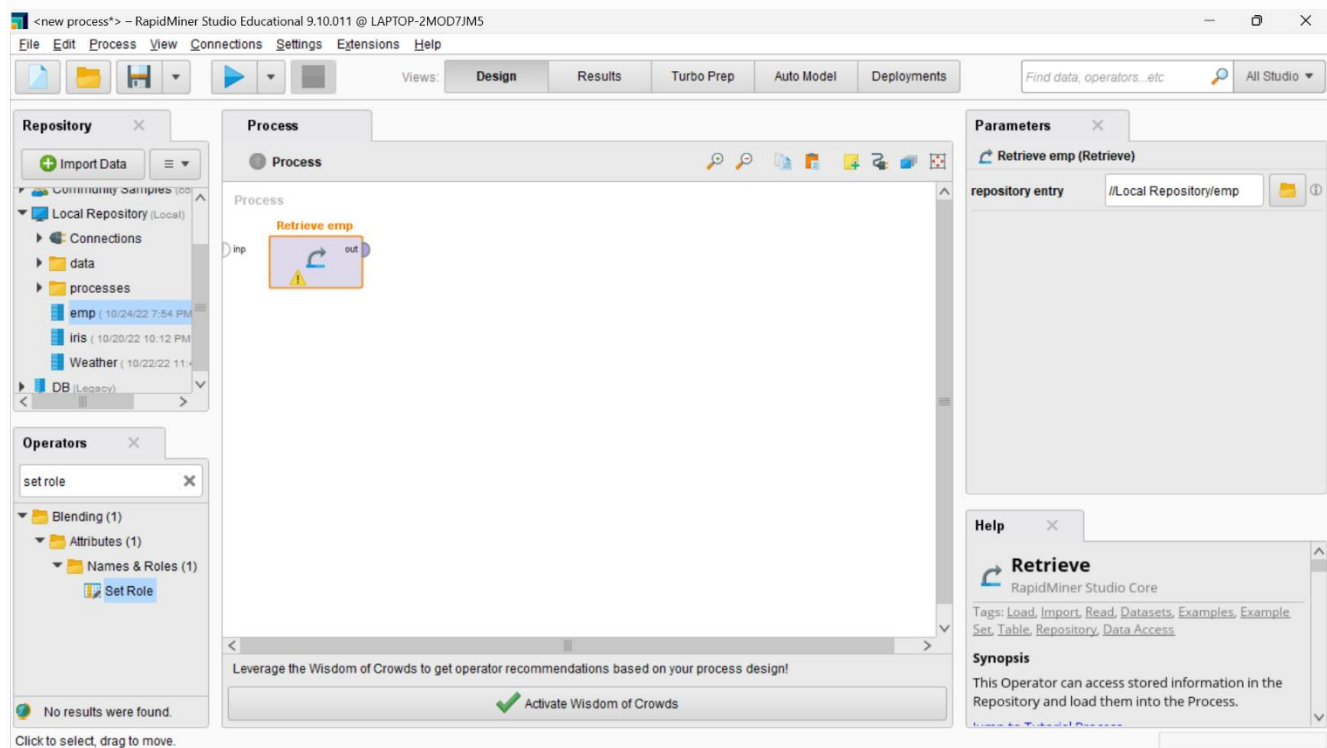
About Linear Regression:

In the simplest words, Linear Regression is the supervised Machine Learning model in which the model finds the best fit linear line between the independent and dependent variable i.e it finds the linear relationship between the dependent and independent variable.
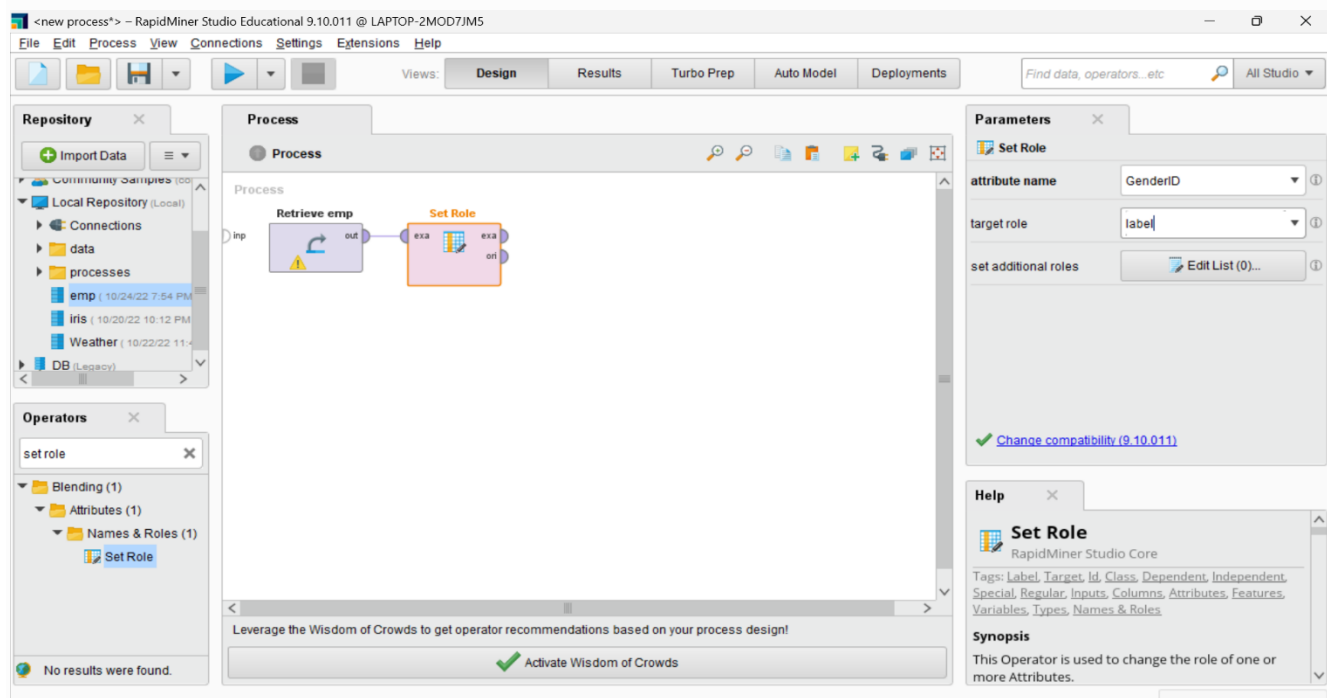
Simple Linear Regression: A linear regression model with one independent and one dependent variable. Multiple Linear Regression: A linear regression model with more than one independent variable and one dependent variable
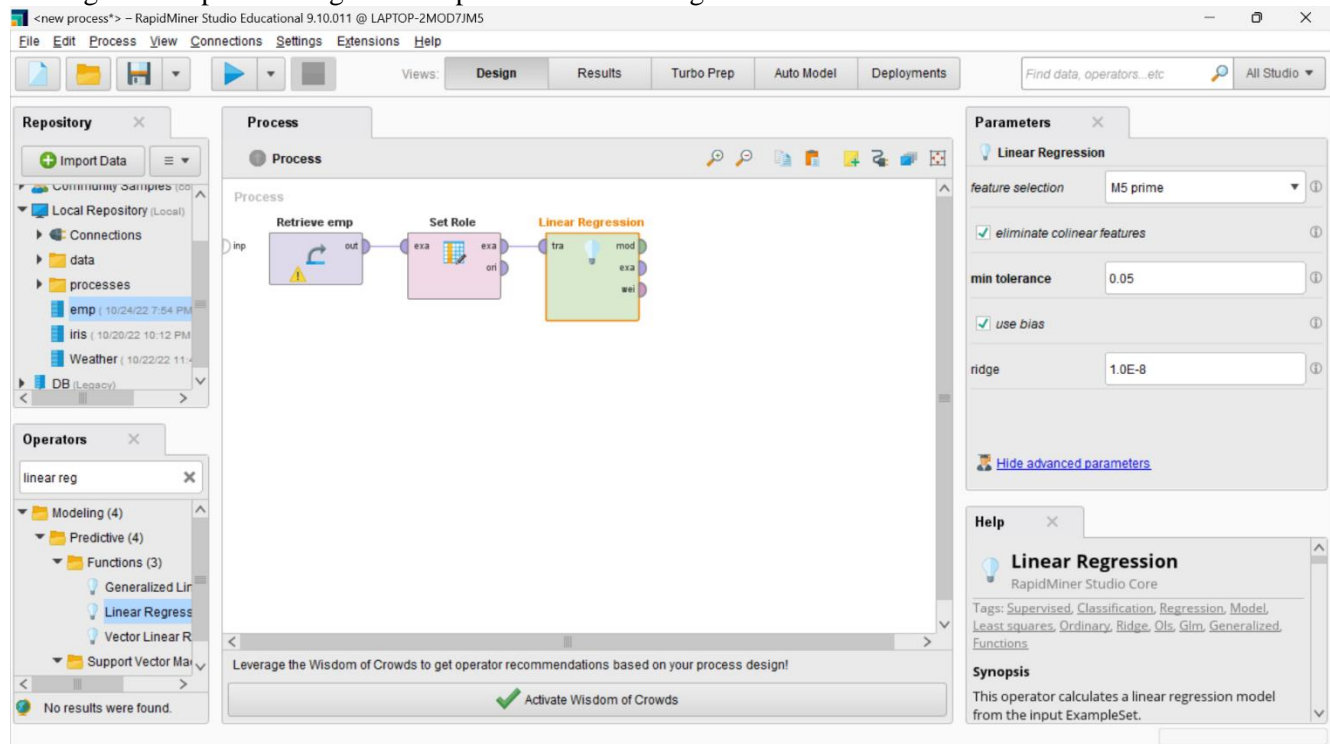
### 9.1.3 Steps to implement simple Linear Regression:

1.To implement the linear regression on Employee dataset using Rapidminer first import the dataset and drop the dataset into the design screen.
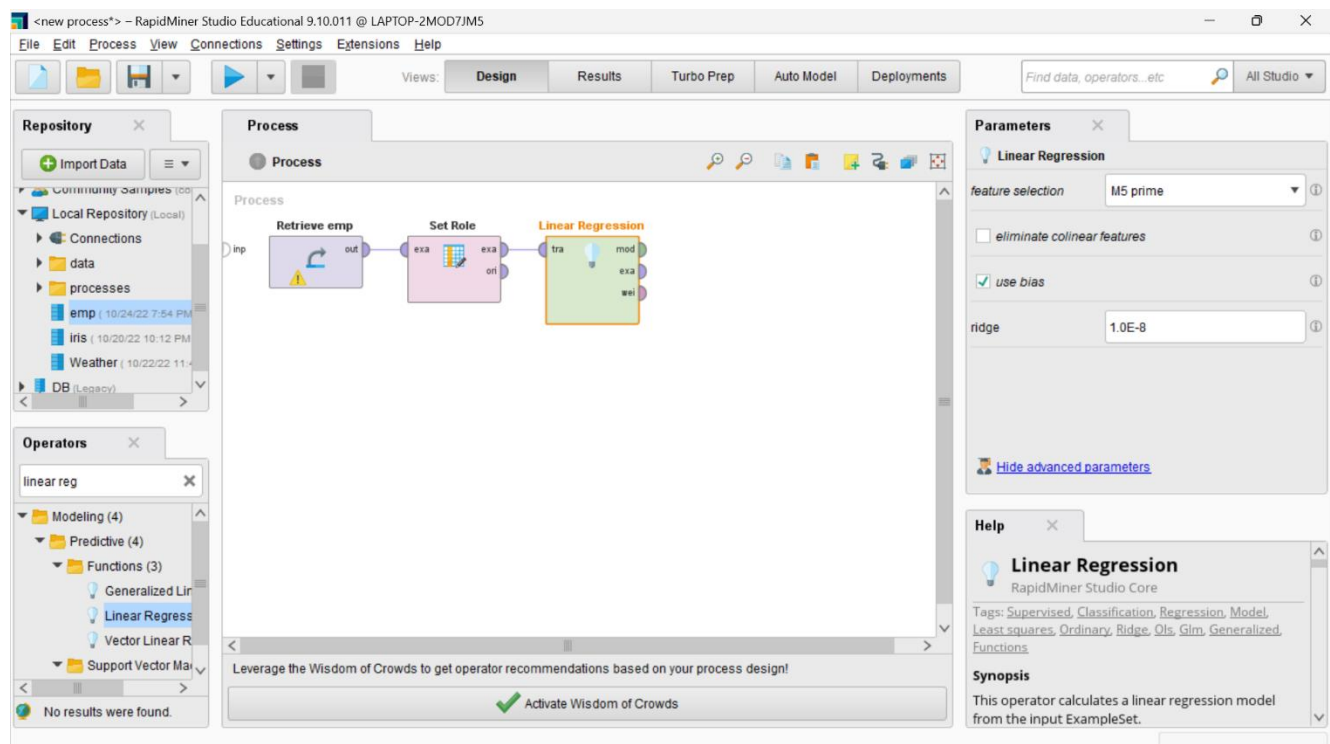


2. Now under operators drag and drop the set role operator and edit the parameters like attribute name and target role.
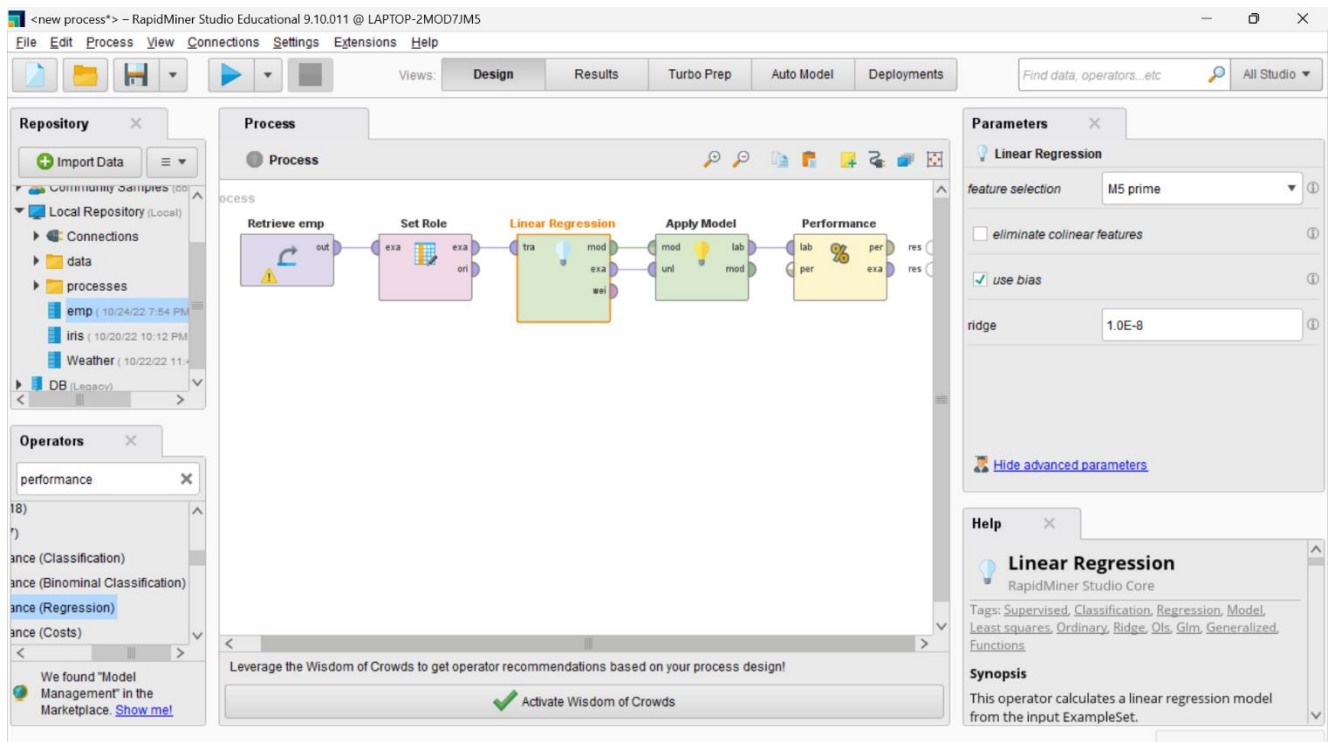
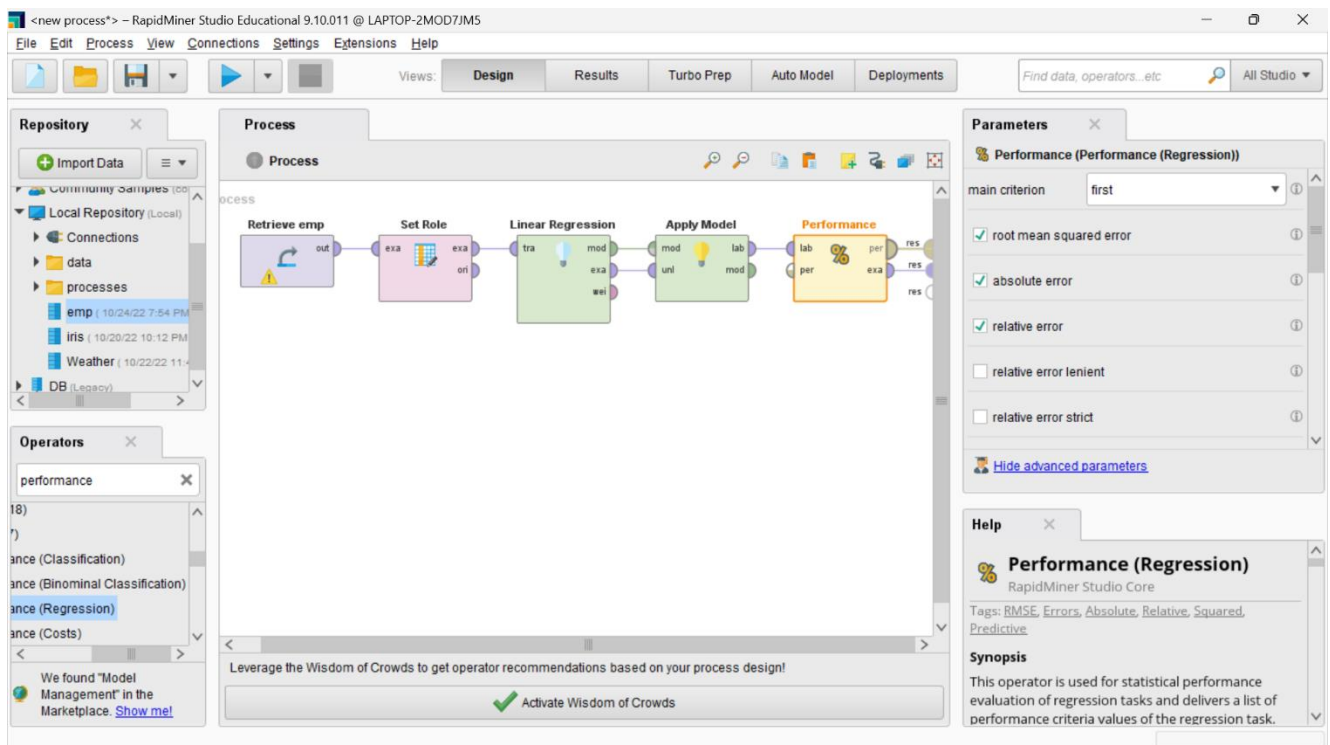3. Drag and Drop Linear regression Operator into the design screen.



Now in the parameters of Linear regression enable the eliminate colinear features.
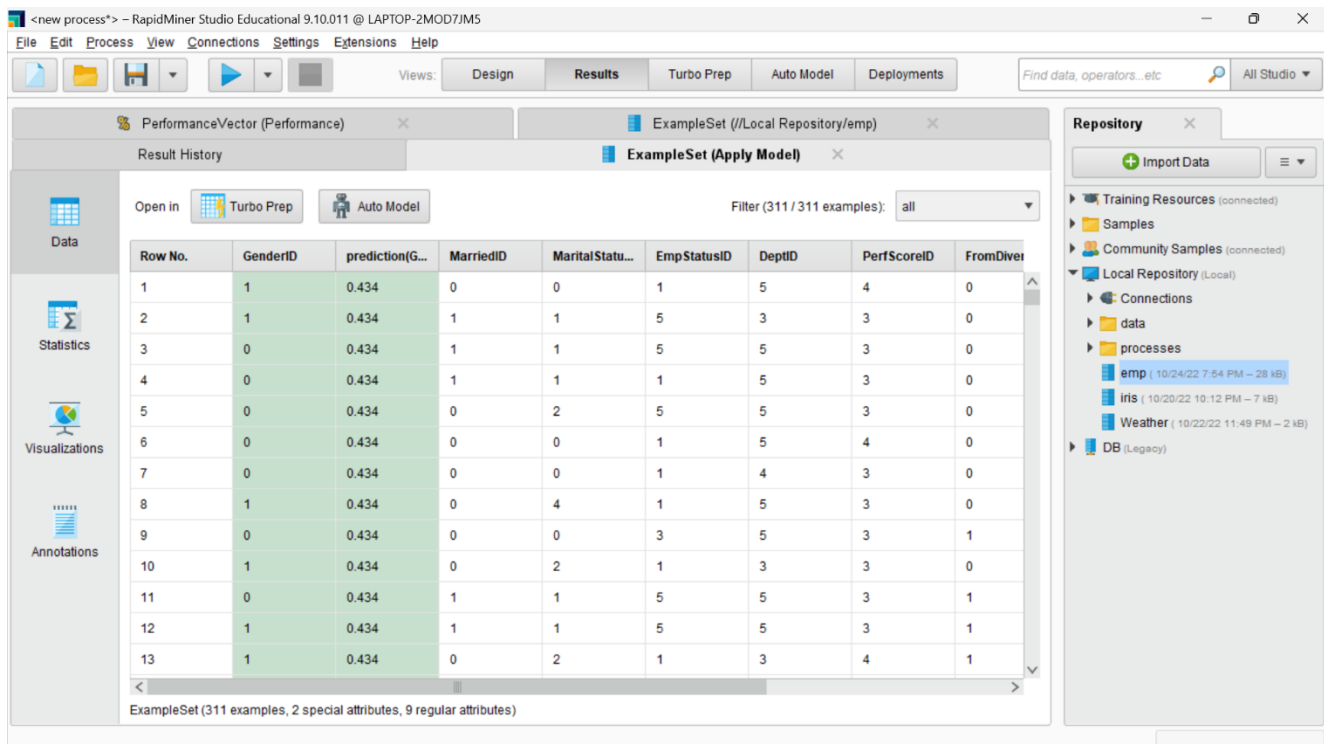


4. Drag and drop Apply Model and performance (Regression) operator into the design screen.
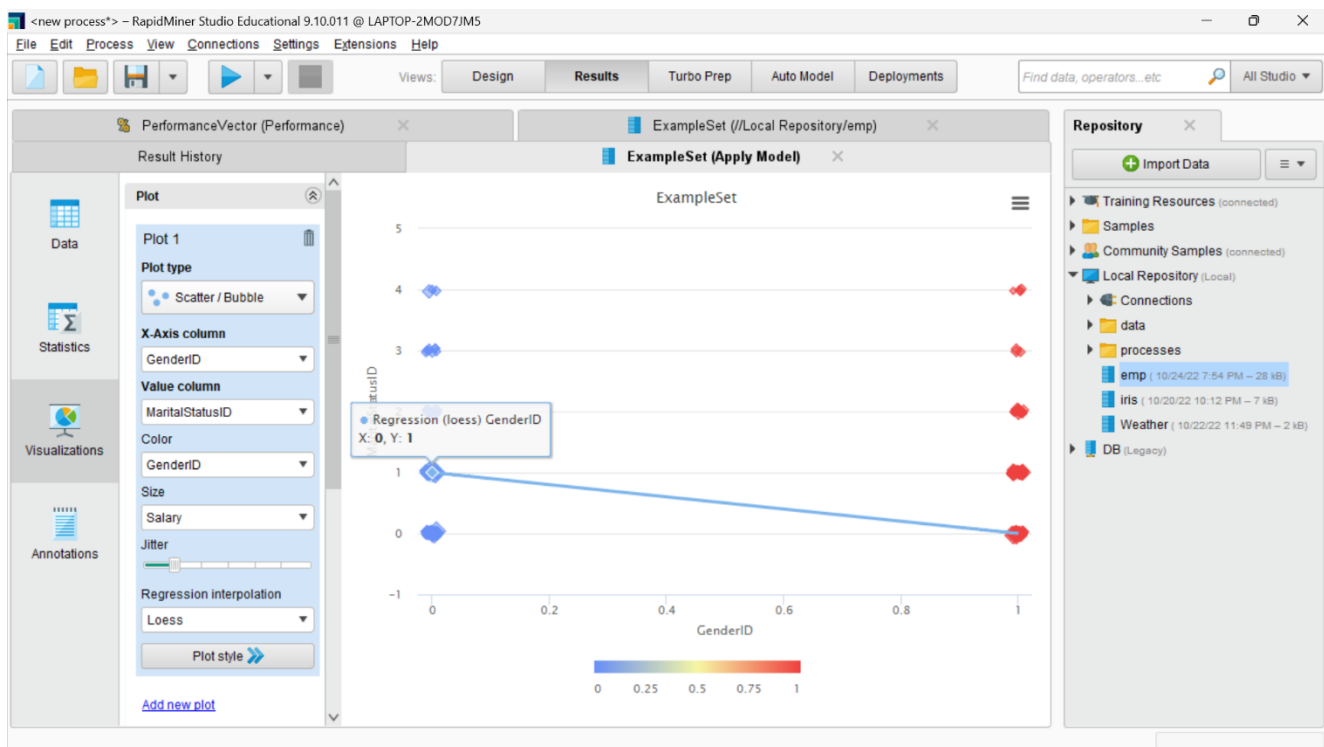
5. In Performance select the parameters that has to display result when the design is executed and give the flow of connections as shown in the below figure, and click on run button to see the result.
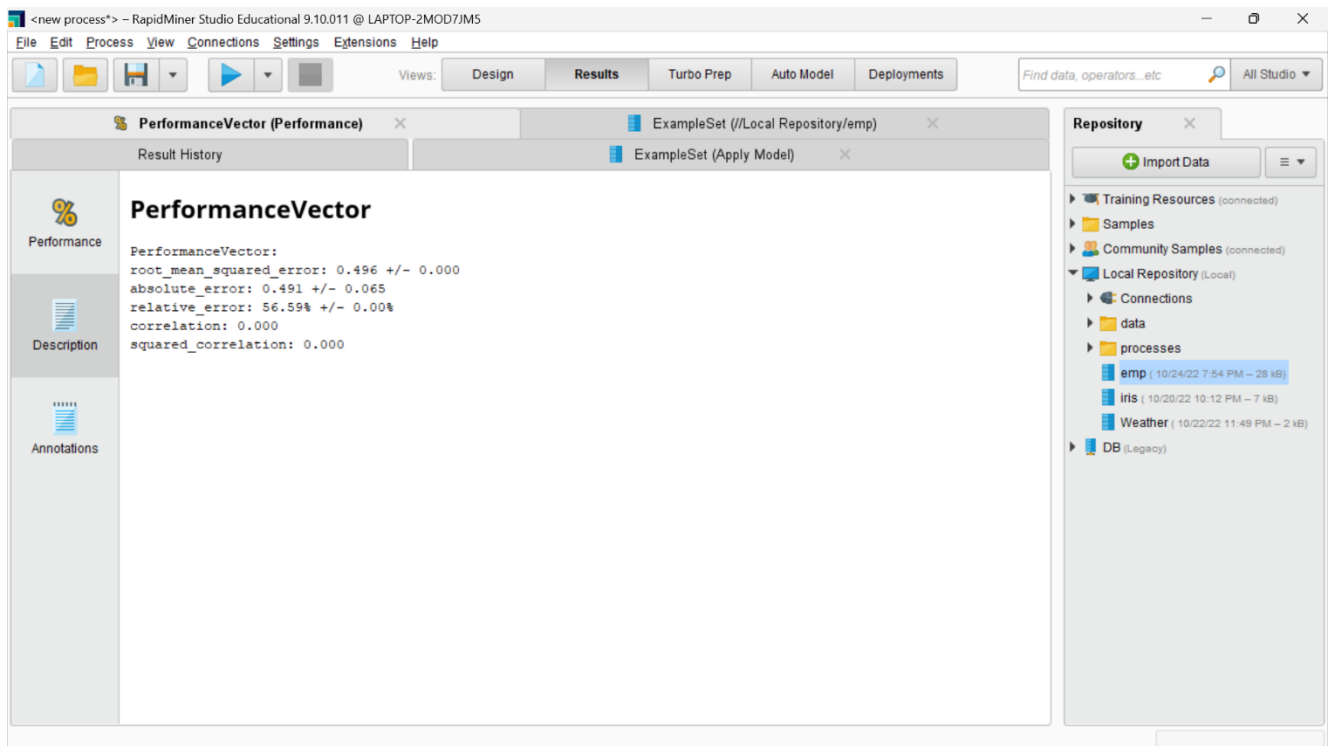


6. Now result of the design is shown in many ways like Data, Statistics, Visualizations, Performance vector.

7. In Visualization we can plot the graph according to our modifications.

## 9.1.4 Results and Discussion:

Therefore, Implementation of Simple Linear Regression on an 'Employee' dataset using RapidMiner is executed successfully.

### 9.2.1 Problem Statement:

Demonstrate the simple k-Means clustering algorithm on 'iris' dataset.

### 9.2.2 Description:

About Dataset used:

The attributes are:

1. sepal length in cm

2. sepal width in cm

3. petal length in cm

4. petal width in cm

5. class

6. Number of samples of each species of iris flowers

7. Predicted attribute: class of iris plant

8. Missing Attribute Values: None

The Iris Dataset contains information of three species of Iris flowers (Iris setosa, Iris virginica and Iris versicolor). The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are NOT linearly separable from each other.

Data Set Characteristics: Multivariate

Area: Life Sciences

Number of samples (or instances) in the dataset: 150

Number of attributes (or features): 05

Attribute Information:

-- Iris Setosa

-- Iris Versicolour

-- Iris Virginica Class Distribution: 33.3% for each of 3 classes

About RapidMiner:

Rapidminer is a comprehensive data science platform with visual workflow design and full automation. It means that we don't have to do the coding for data mining tasks. Rapidminer is one of the most popular data science tools.

This is the graphical user interface of the blank process in rapidminer. It has the repository that holds our dataset. We can import our own datasets. It also offers many public datasets that we can try. We can also work with a database connection.

RapidMiner is an awesome visual workflow designer. The way they present visually is so unique. It helps in speeding and automating the creation of visual models. It helps in creating models in only 5 clicks by automated machine learning.

About k-Means Clustering algorithm:

K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if K=2, there will be two clusters, and for K=3, there will be three clusters, and so on. It allows us to cluster the data

into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.
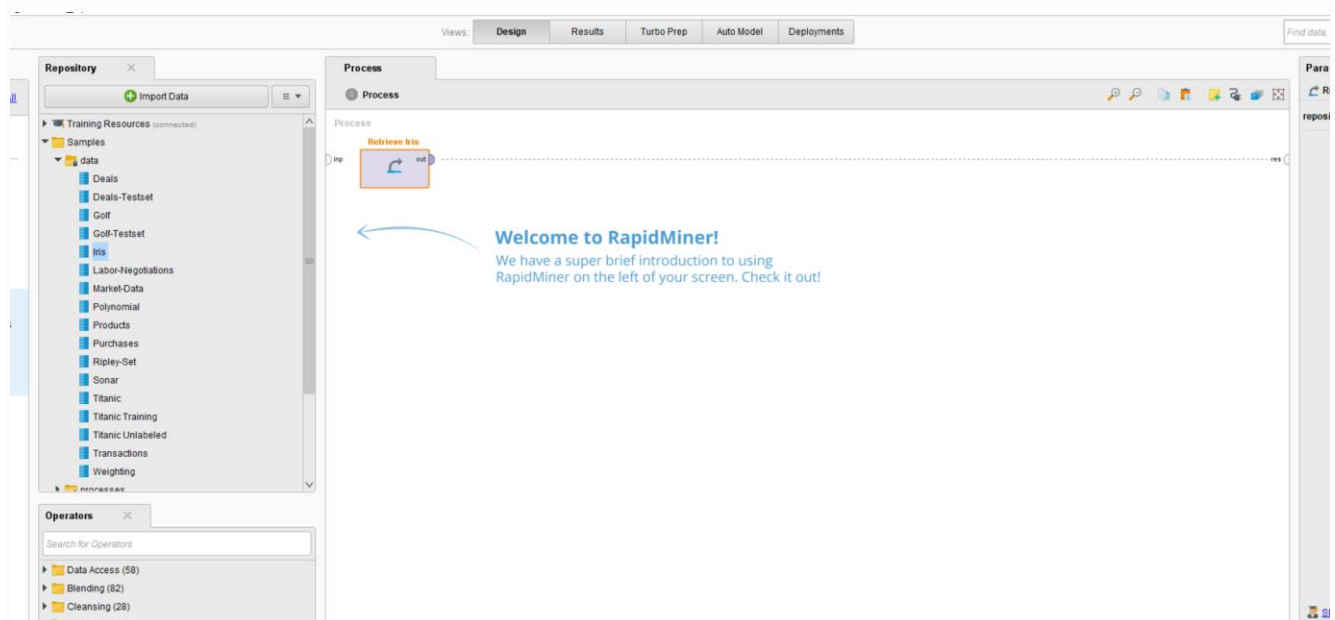
The k-means clustering algorithm mainly performs two tasks:

- o Determines the best value for K center points or centroids by an iterative process.

- o Assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.
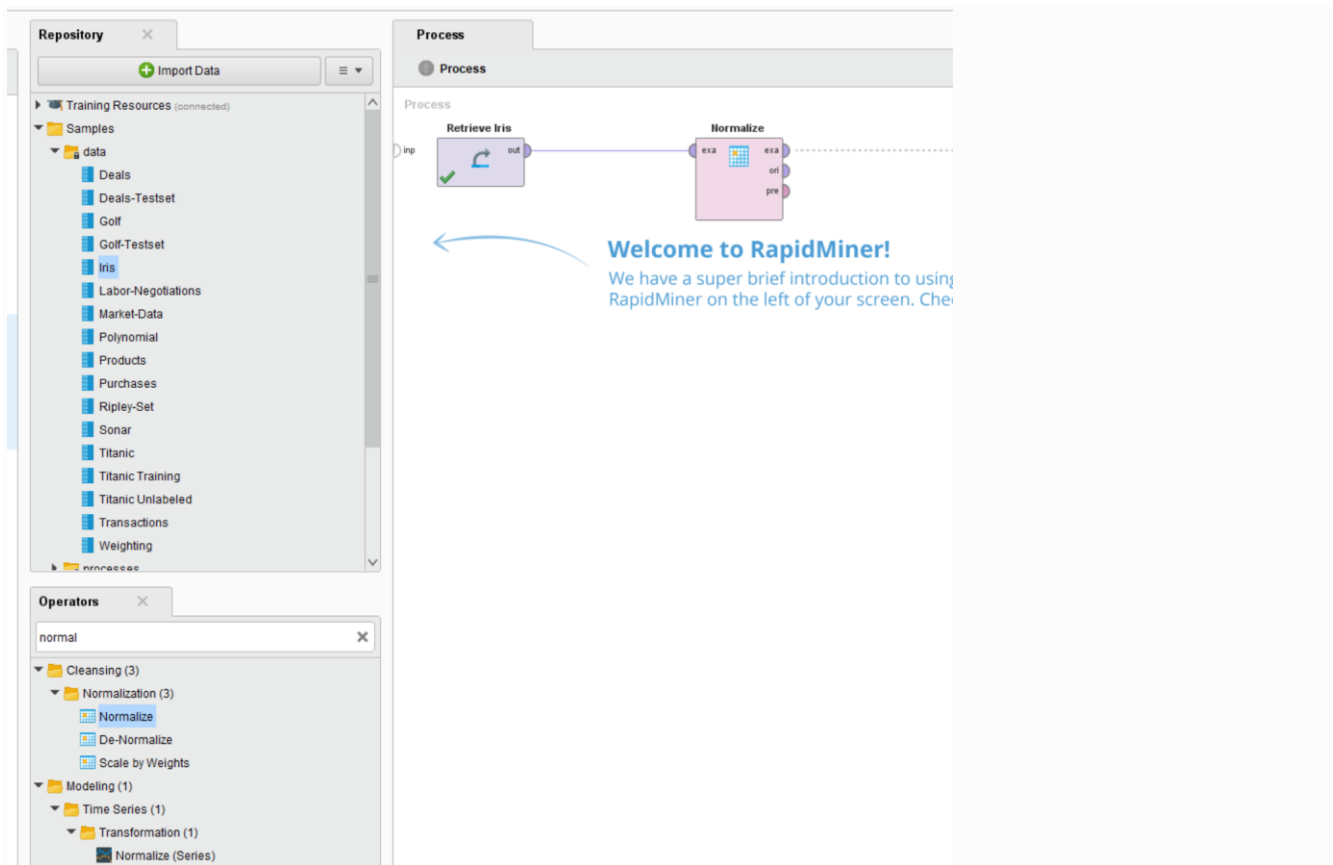
Hence each cluster has datapoints with some commonalities, and it is away from other clusters.

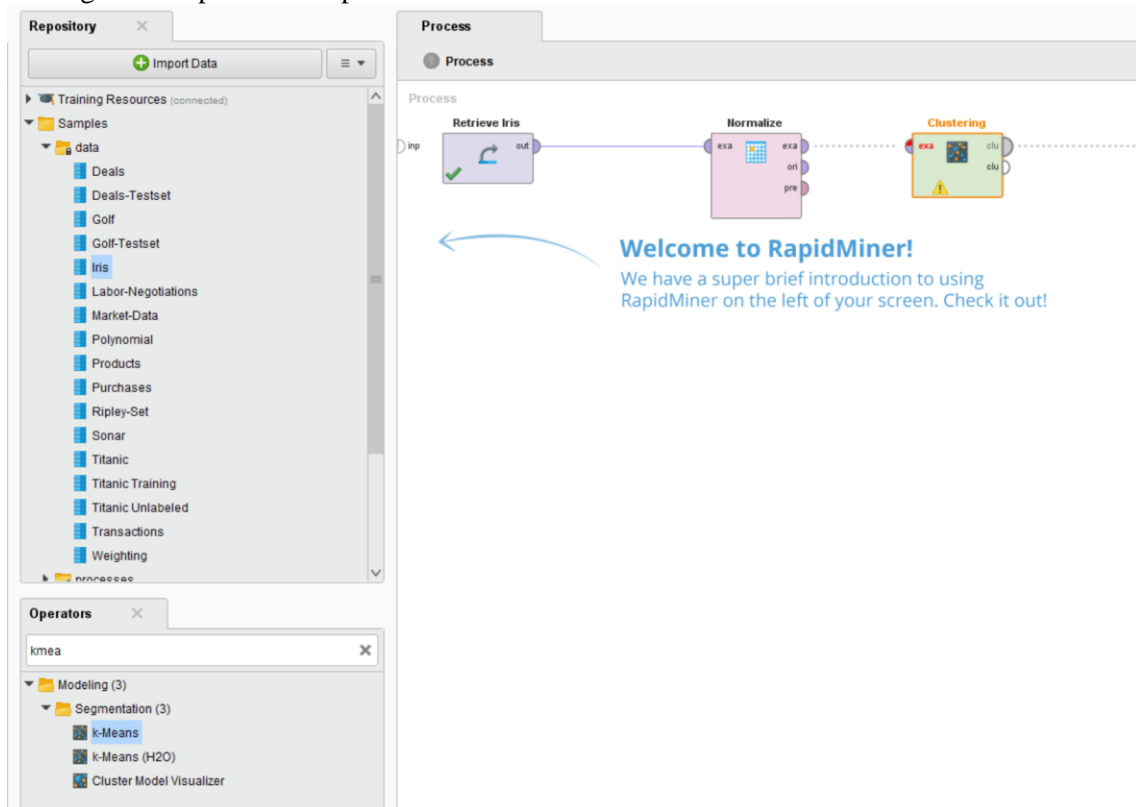### 9.2.3 Steps to implement simple K-Means clustering algorithm:

1. To implement the K-Means clustering on Iris dataset using Rapidminer, first import the dataset and drop the dataset into the design screen.



2. Now under operators drag and drop the Normalize operator and connect the dataset to it.

3. Drag and Drop k-means operator and connect it.



4. Set the k value to 3 to create three clusters and run it and see the clusters.

**Cluster Model**

```
Cluster 0: 50 items
Cluster 1: 44 items
Cluster 2: 56 items
Total number of items: 150
```

5. Drag and drop the Cluster Distance Performance from the operators and see the performance of our k-means algorithm.



6. Now the performance vector is seen.

**PerformanceVector**

```
PerformanceVector:
Avg. within centroid distance: -0.935
Avg. within centroid distance_cluster_0: -0.963
Avg. within centroid distance_cluster_1: -0.988
Avg. within centroid distance_cluster_2: -0.867
Davies Bouldin: -0.834
```

### 9.2.4 Results and Discussion:

 Therefore, Implementation of Simple k-means clustering algorithm on 'iris' dataset using RapidMiner is executed successfully.