# Assignment 1

Name – Akash Raj                                          Subject code – CAP 447

Roll no – RD2112B78                                    Section – D2112

Reg no – 12108382                                      Submitted to – Poonam Rattan Mam
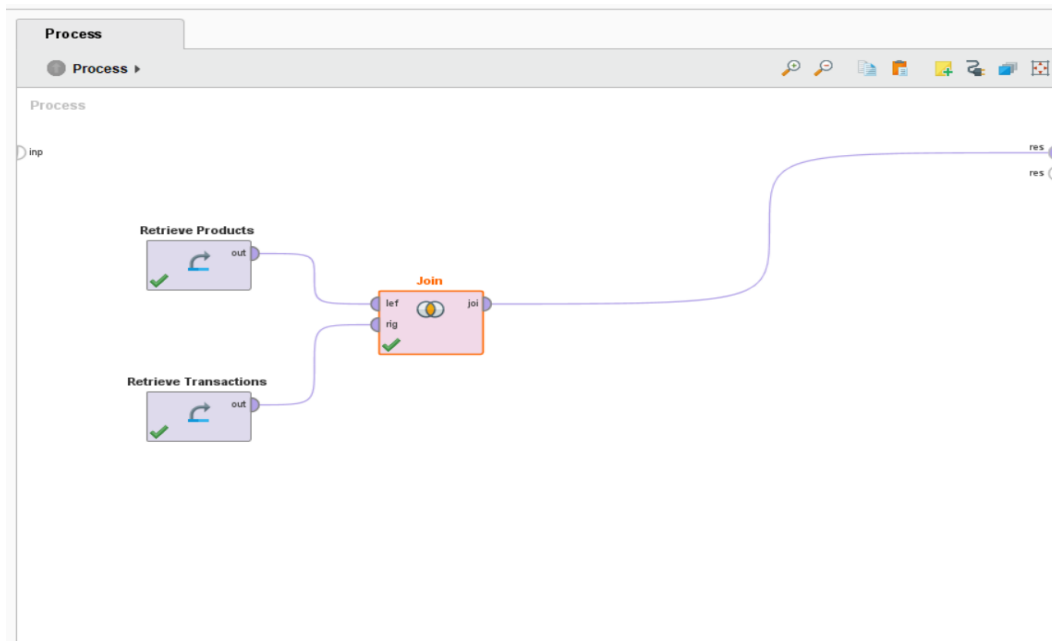
**1. Explain briefly the main components of Rapid miner.**

**Ans**- The main components of rapid minor are –

i) Repository- A repository is simply a folder that holds all of your RapidMiner data sets (we call them "ExampleSets), processes, and other file objects that you will create using RapidMiner Studio. This folder can be stored locally on your computer, or on a RapidMiner Server.

ii) Operators - The building blocks, grouped by function, used to create RapidMiner processes. An *operator* has input and output ports; the action performed on the input ultimately leads to what is supplied to the output. Operators parameter control those actions. There are more than 1500 operators available in RapidMiner. Operators, in the Operator panel of the Design view, are both browsable and searchable.

iii) Parameters - The setting(s) whose value(s) determine the characteristics or behavior of an operator. RapidMiner presents parameters in the Parameter panel of the Design view. There are regular parameters and expert parameters. The expert parameters are indicated by italic names and are displayed or hidden by clicking the Show/Hide advanced parameters link at the bottom of the panel.

iv) Help – It gives the definition and description of the operators or the attributes which are being used.

v) Views – This section consists of few important elements which are used in rapid miner. It consist of  design, result, turboprop, auto model and deployment. The design and result options are used to switch between the visualization part and the result part.

**2. a) Explain any 5 operators with its usage with its snapshot in rapid miner.**

Ans- Five operators with its usage are –

1) **Join Operator –** The join operator is used to combine two or more data sets into a single using one or more attributes of the given data set known as key attributes.
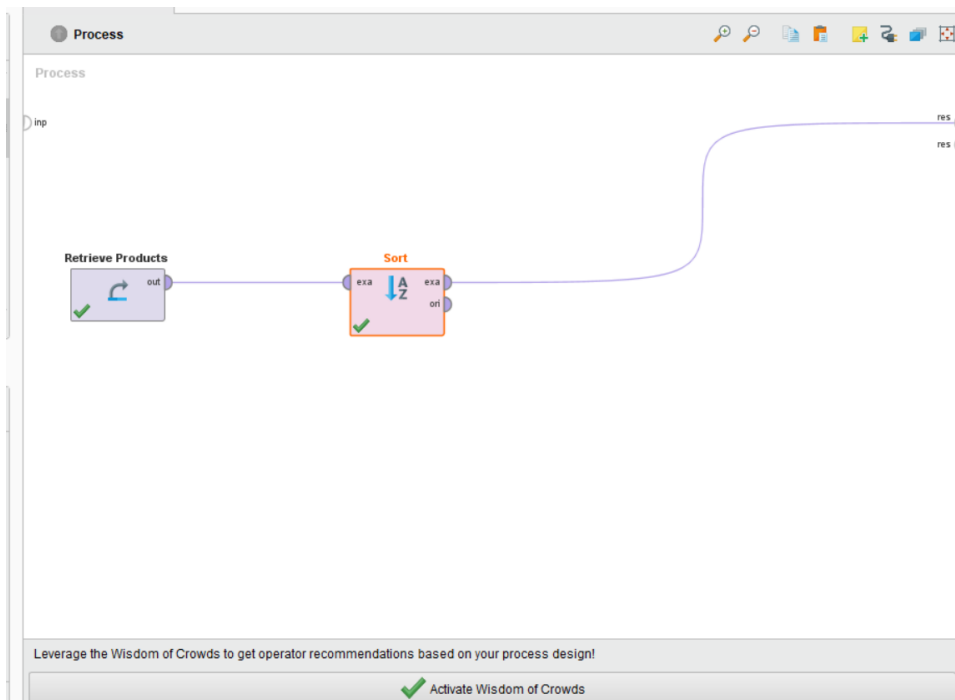
<u>Product and transaction are combined using join operator</u>

2) **Sort Operator –** The sort operator is used to sort the data set provided at the input port. The complete data set is sorted according to a single or more attributes. The sorting is done in either ascending or descending order.

The price of the product dataset has been sorted using sort operator in ascending order.

3) **Aggregate operator –** The aggregate operator performs the aggregation function known from SQL. This operator provides a lot of functionalities in the same format as provided by the SQL aggregation function. It performs mathematical operations like average, count, aggregate, Max, Min and Sum.

_The Count and mode functions are used from the Aggregate function_

4) **Discretize By Binning** - This operator discretizes the selected numerical attributes to nominal attributes. The number of bins parameter is used to specify the required number of bins. This discretization is performed by simple binning. The range of numerical values is partitioned into segments of equal size.

The price has been discretized using the Discretize by binning method

5) **Replace Missing Values –** This operator is use to replace missing value in example of selected attributes by a specified replacement. The missing values can be replaced by minimum, maximum or average value of that attribute. Zero can also be used to replace the missing values. Any replenishment value can also be specified as the replacement of the value.

Replacing missing values of price

**b. What do you understand by graphical representation and statistics in rapid miner. Attach any 5 different types of graphs and their interpretation.**

**Ans-** Graphical representation is a method to show and represent values, increases, decreases, comparisons. It is the visual representation of data through charts and
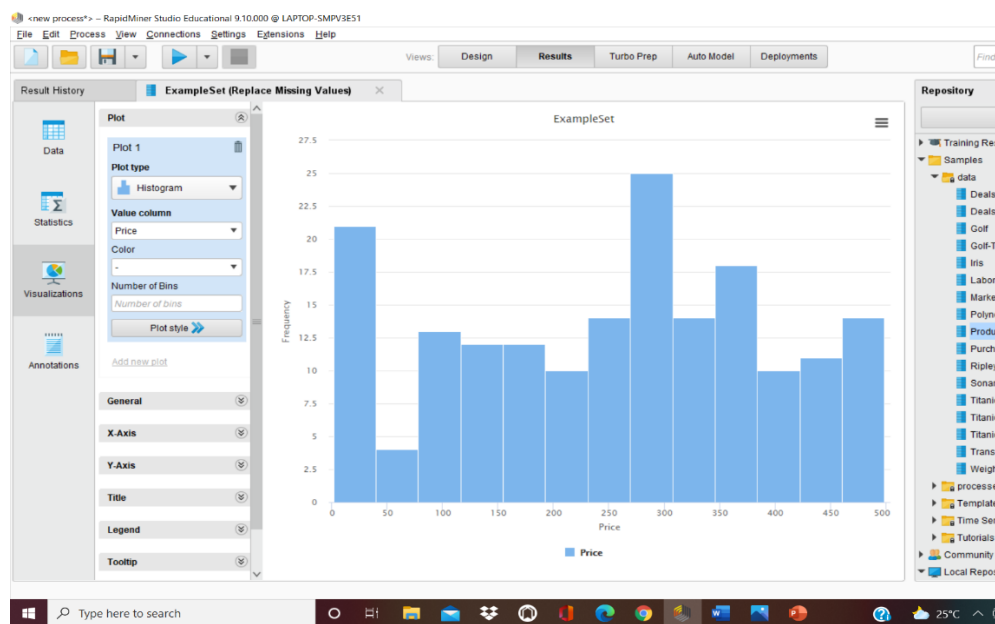
graphs. It helps to make prediction or show report of how certain situation was yesterday and now.

Statistics provides meaning to the data represented by the graphs.

Five different types of graphs are-

i)  **Histogram –** Histogram is a non-cumulative frequency graph. It is drawn on a natural scale in which the representative frequencies of the different class of values are represented through vertical rectangles drawn closed to each other. Measure of central tendency, mode can be easily determined with the help of this graph.



A Histogram

ii)  **Bar-Chart –** A bar chart or bar graph is a chart or a graph that presents categorical data with rectangular bars with height or lengths proportional to the values that they represent. The bars can be plotted vertically or horizontally. A vertical bar graph is sometimes called a line graph.

A Vertical Bar Graph

iii) **Line Graph** – A line graph is a type of chart which displays information as a series of data points called 'markers' connected by a straight line segments.



A Line Graph

iv) **Pie-Chart** - A pie chart is a circular statistical graph, which is divided into slices to illustrate numerical proportion. In this, the arc length of each slice, is proportional to the quantity it represents.

A Pie Chart

v) **Scatter plot -** A scatter plot uses dots to represent values for two different numeric variables. The position of each dot on the horizontal and vertical axis indicates values for an individual data point. Scatter plots are used to observe relationships between variables.
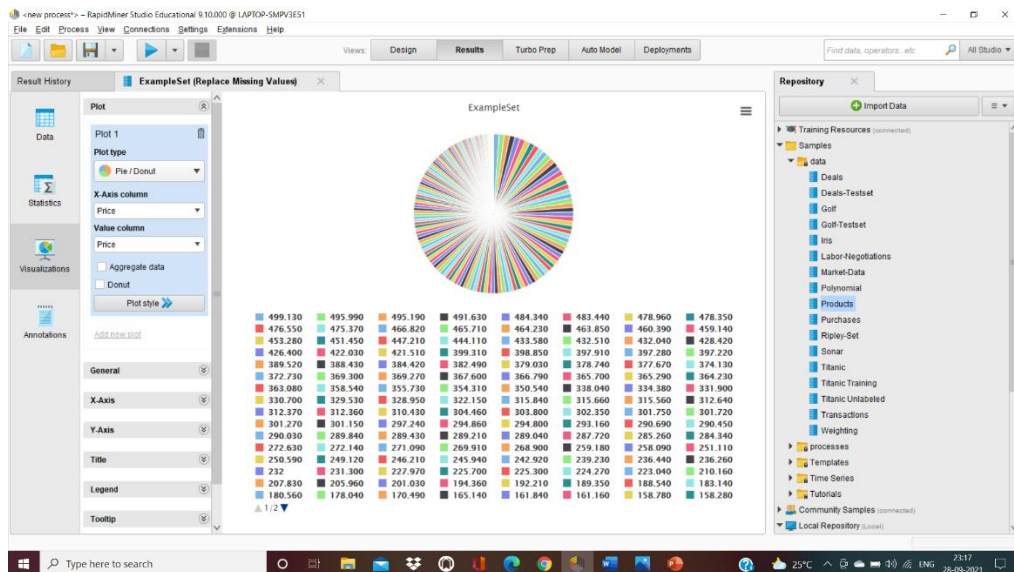


A Scatter Plot

**Ques 3 – What do you mean by data preprocessing in Data Mining? Explain any 5 operations to handle missing values and attach screenshot of same.**

**Ans –** Data preprocessing is use to improve the quality of data in the data warehouse. Using this technique we remove the noisy data. It also helps in correcting the inconsistency in data

as well as the incomplete data, due to which the efficiency of data is increased as well as it is easy to mine the data.

**Steps Involved in Data Preprocessing:**
1. **Data Cleaning:** The data can have many irrelevant and missing parts. To handle this part, data cleaning is done. It involves handling of missing data, noisy data etc.
2. **Data Transformation:**
   This step is taken in order to transform the data in appropriate forms suitable for mining process. This involves following ways:
   1. **Normalization:**
      It is done in order to scale the data values in a specified range (-1.0 to 1.0 or 0.0 to 1.0)

   2. **Attribute Selection:**
      In this strategy, new attributes are constructed from the given set of attributes to help the mining process.

   3. **Discretization:**
      This is done to replace the raw values of numeric attribute by interval levels or conceptual levels.

   4. **Concept Hierarchy Generation:**
      Here attributes are converted from lower level to higher level in hierarchy. For Example-The attribute "city" can be converted to "country".

3. **Data Reduction:** Since data mining is a technique that is used to handle huge amount of data. While working with huge volume of data, analysis became harder in such cases. In order to get rid of this, we use data reduction technique. It aims to increase the storage efficiency and reduce data storage and analysis costs.

   Five operations to handle missing values are-
1) **Ignore the Data Row-** This is usually done when the class label is missing (assuming your data mining goal is classification), or many attributes are missing from the row (not just one). However, you'll obviously get poor performance if the percentage of such rows is high.
2) **Use a global constant to fill in for missing values -** Decide on a new global constant value, like "unknown", "*N/A*" or minus infinity, that will be used to fill all the missing values. This technique is used because sometimes it just doesn't make sense to try and predict the missing value.
3) **Use attribute mean –** Replacing missing values of an attribute with the mean (or median if its discrete) value for that attribute in the database.
4) **Use attribute mean for all samples belonging to the same class -** Instead of using the mean (or median) of a certain attribute calculated by looking at all the rows in a database, we can limit the calculations to the relevant class to make the value more relevant to the row we're looking at.

5) **Using a data mining algorithm to predict most probable value -** The value can be determined using regression, inference based tools using Bayesian formalism, decision trees, clustering algorithms (K-Mean\Median etc.).