# ASSIGNMENT-1

NAME-MADHAV JHA

REGISTRATION NO.-12107418

SECTION-D2112

ROLL. NO.-RD2112A41

COURSE CODE-CAP447

DATE OF SUBMISSION-29-09-2021

# Q1.Explain briefly the main components of RAPID MINER Studio.

Answer:-

RapidMiner Studio is a  where you build and edit analytic processes. RapidMiner Studio and RapidMiner Server connect and interact with each other.RapidMiner Studio is a visual workflow designer for predictive analytics that brings data science and machine learning to everyone on the analytics team.When you're working on a new project of any kind, often the first step will be to go to a whiteboard, where you will plan the workflow and identify the key steps on the way to your goal. If you're a data scientist, the workflow will usually include one or more of the following steps:

- Import data
- Prepare data
- Build a model
- Validate the model
- Apply the model

RapidMiner Studio implements your whiteboard workflow in software, in the *Design View*. The Design View includes numerous *panels*.

- Repository

  Repository: The place where your data, processes, and results are stored, either locally or remotely.

Also known as: folder, workspace, project

When working with RapidMiner Studio, you need a place to save your work. The Repository can be used to store:

- Data
- Processes
- Results
- Operators

Operators: The elements of a Process, each Operator takes input and creates output, depending on the choice of parameters.Also known as: function, formula, node.To use RapidMiner Studio effectively, you have to learn about its Operators. RapidMiner Studio includes hundreds of Operators

but all we not need to know,we have to know only a few operators.

- Parameters

Parameters: Options for configuring the behavior of an Operator.The content of the Parameters Panel is context-dependent. Select any Operator that is displayed in the Process Panel, and the Parameters Panel displays the options for configuring that Operator. Because RapidMiner Studio includes many Operators, each with its own unique functionality, the range of parameters is also quite diverse. By default, RapidMiner Studio will show you only the more commonly used parameters. To see all of the available parameters, click Show advanced parameters.

- Help:-

*Help: Displays a help text for the current Operator.* The content of the *Help Panel* is also context-dependent. Select any Operator that is displayed in the Process Panel, and the Help Panel displays

a help text for that Operator. The Help Panel provides useful background informations.

- Process

Process: A connected set of Operators that help you to transform and analyze your data.Also known as: flow, program, pipeline, diagram .Your goal is to create a finished process, a connected set of Operators that produce a result. For example, your process might read a data set and build a predictive model. When you have connected all your Operators and set their parameters, press the Run Run arrow button at the top of the user interface, and the results will be displayed in the Results View.

Q2(A). Explain any 5 operators along with its usage and snapshots in rapidminer.

Answer:-The five operators along with uses and snapshot are below:-

- Sort:-

This operator sorts the input data set in ascending or descending Order according to several

attributes.The attributes to sort by are specified using the sort by Parameter. For each attribute, sorting is done in ascending or descending order, depending on the setting of the sorting order parameter. The resulting data set is sorted by the first attribute,Then subsets of the same value in the first attribute are sorted by the second attribute.

Figure :-1 ,shows the short operator preview.

Figure :-2 ,shows the short operatoration on sl_no attribute in assending order.



- Join:-

  This Operator joins two ExampleSets using one or more Attributes of the input ExampleSets as key attributes.Identical values of the key attributes indicate matching Examples. An Attribute with id role is selected as key by default but an arbitrary set of one or more Attributes can be chosen as key. Four types of

joins are possible: inner, left, right and outer join. All these types of joins are explained In the parameters section.

Figure :-3 ,shows the join operator preview. which combining two example sets i.e transaction and product.

Figure :-4 ,shows the short operatoration on product id,which combining two example sets i.e transaction and product.



- Replace:-
- Retrieve:-This Operator can access stored information in the Repository and Load them into the Process.The Retrieve Operator loads a RapidMiner Object into the Process. This Object is often an ExampleSet but it can also be a

Collection or a Model. Retrieving data this way also provides the Meta data of the RapidMiner Object.

- Select Attributes:-This Operator selects a subset of Attributes of an ExampleSet and Removes the other Attributes.The Operator provides different filter types to make Attribute selection easy. Possibilities are For example: Direct selection of attributes. Selection by a regular expression or selecting only attributes without missing values. See parameter attribute filter type for a detailed description of the different filter types.

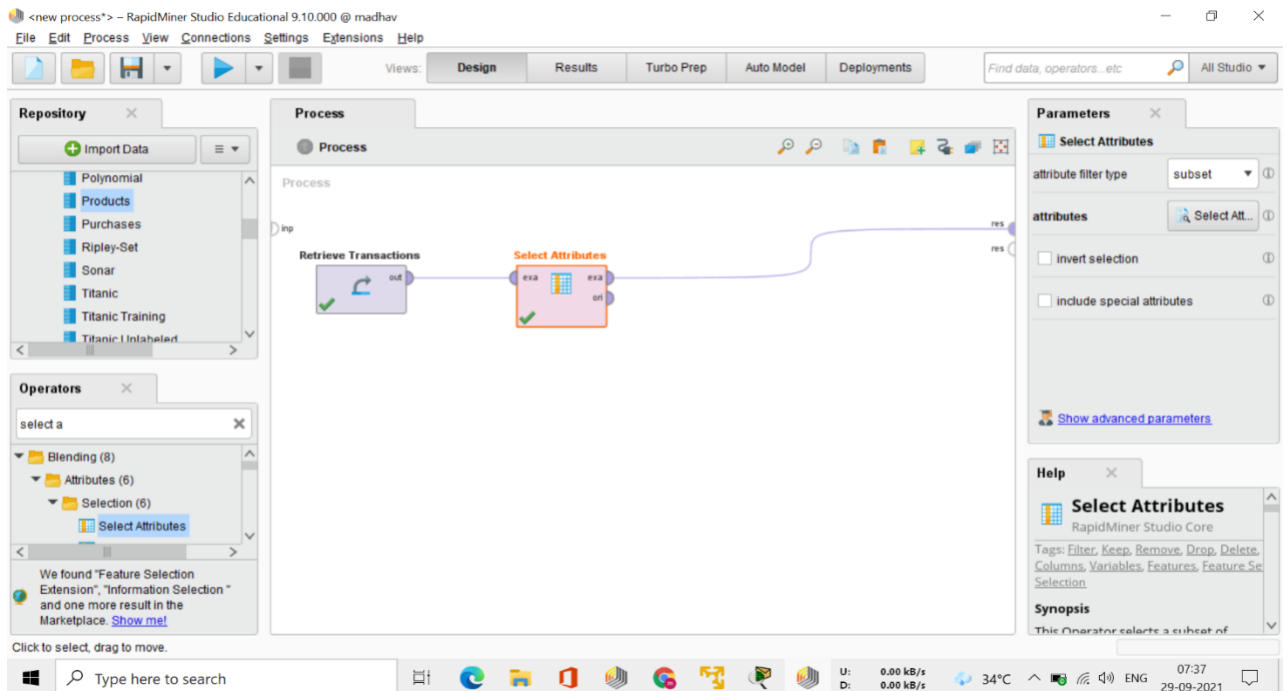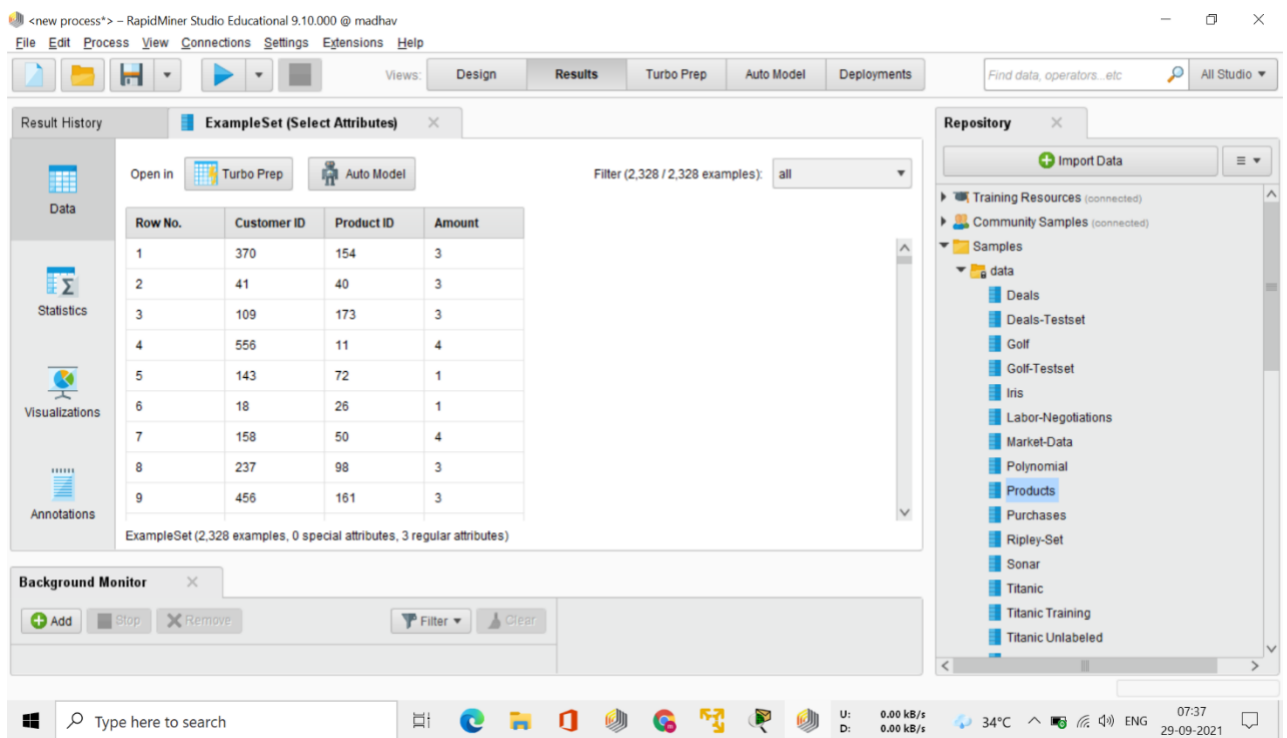Figure :-5 ,shows the select attribute operator preview.

Figure :-6 ,shows the selection of attributes operatoration on transaction data set.

- Replace Missing Values:-

This Operator replaces missing values in Examples of selected Attributes by a specified replacement.Missing values can be replaced by the minimum, maximum or average value of that Attribute.zero can also be used to replace missing values. Any replenishment value can also be specified as a replacement of missing values.
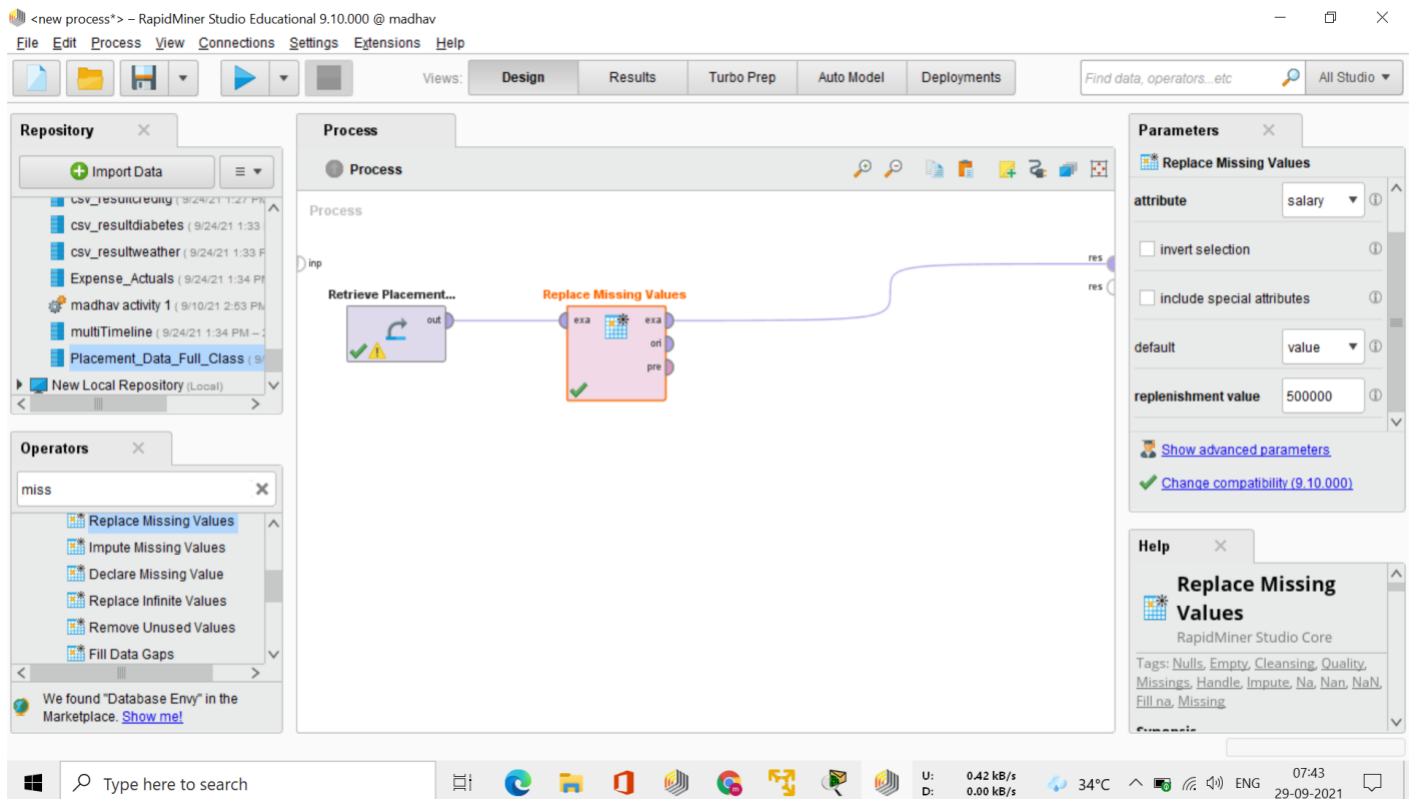
Figure :-7 ,shows the replace missing values operator preview.

Figure :-8 ,shows the select missing values operatoration on placement data set,in which i had some missing values in salary column and replacing missing values with 50000.

- Filter examples:-

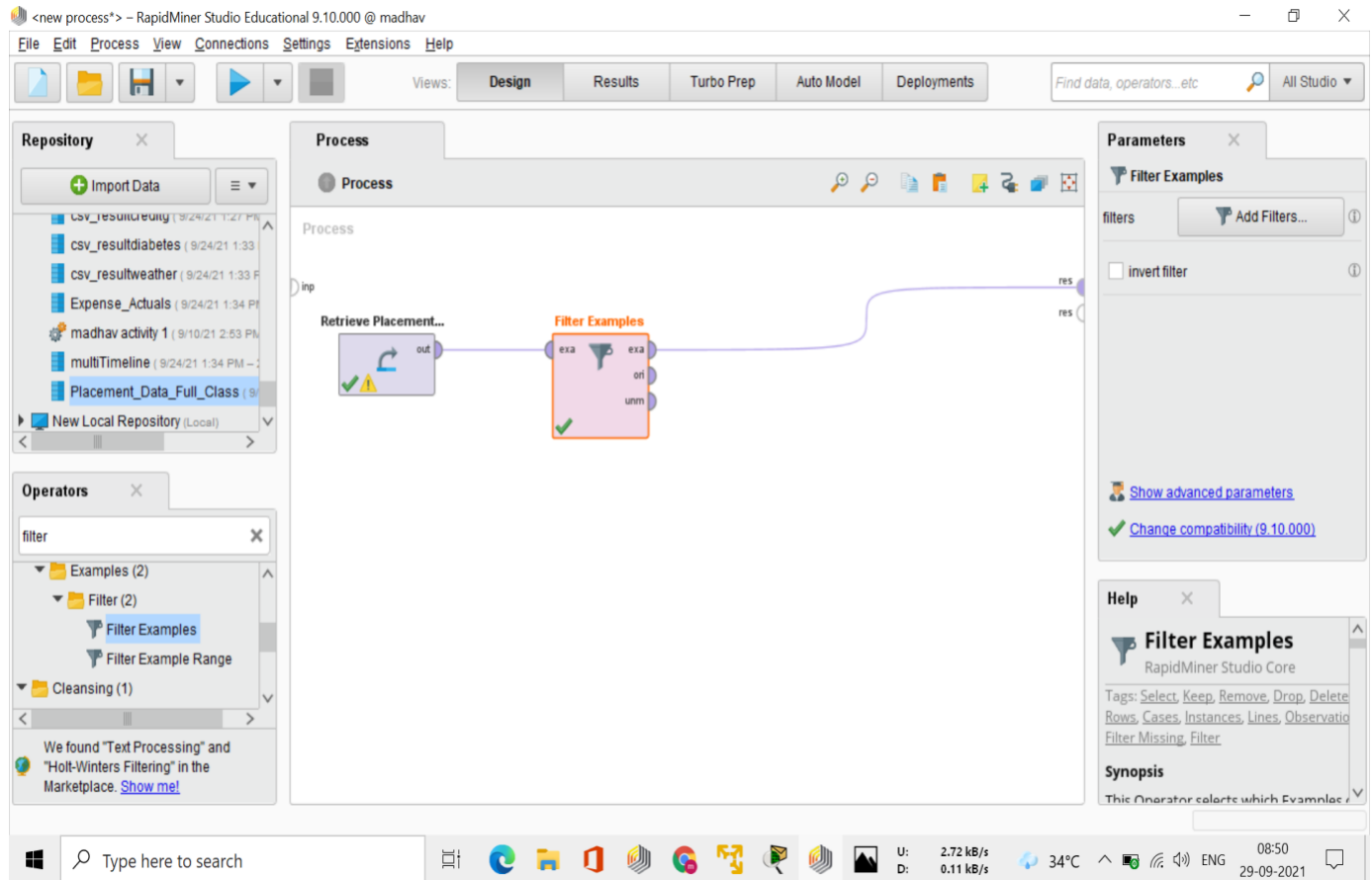Figure :-9 ,shows the filter examples operator preview.

Figure :-10 ,shows the filter examples operatoration on salary attribute whose values are missing and gender is equal to female.
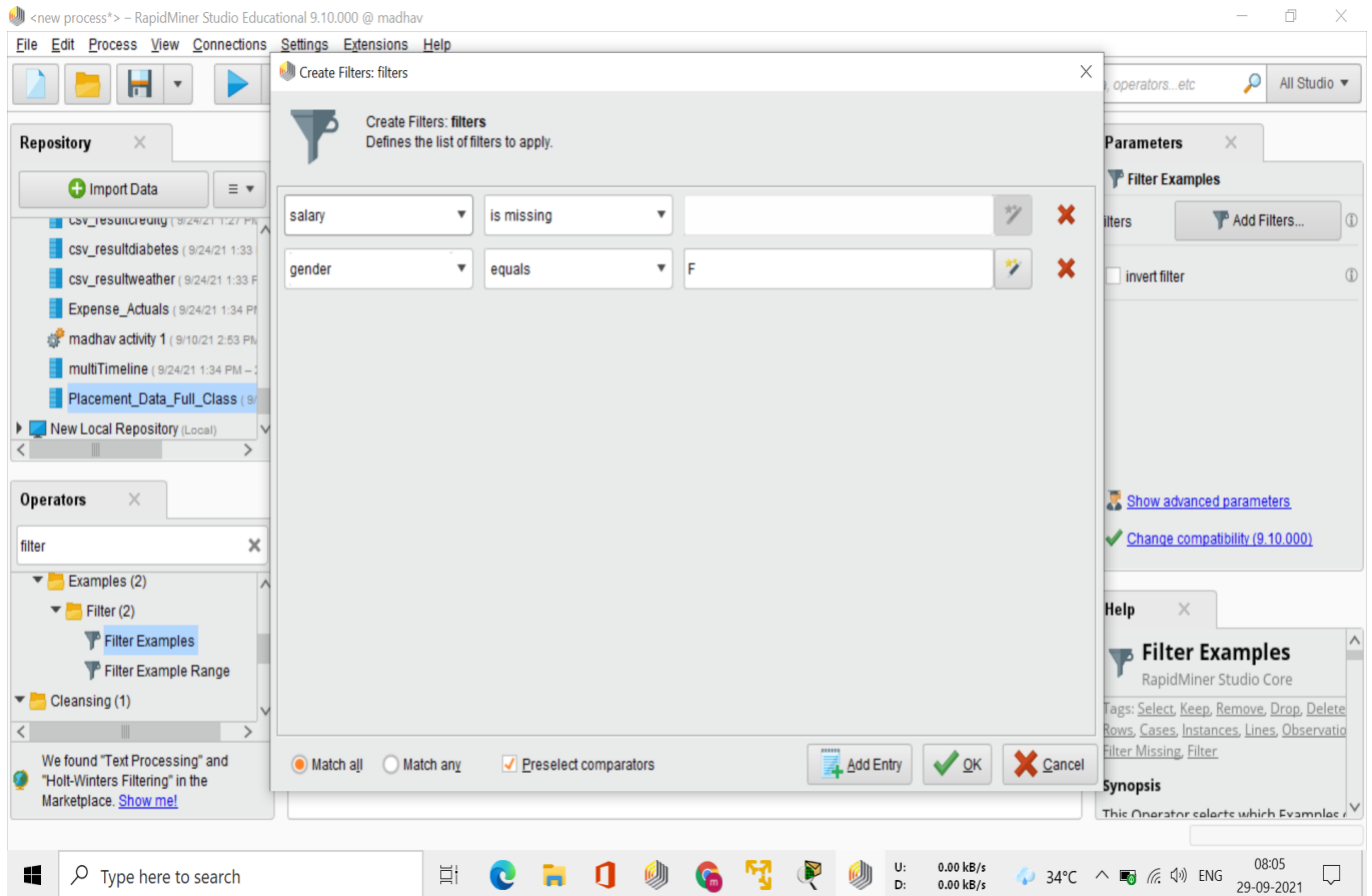
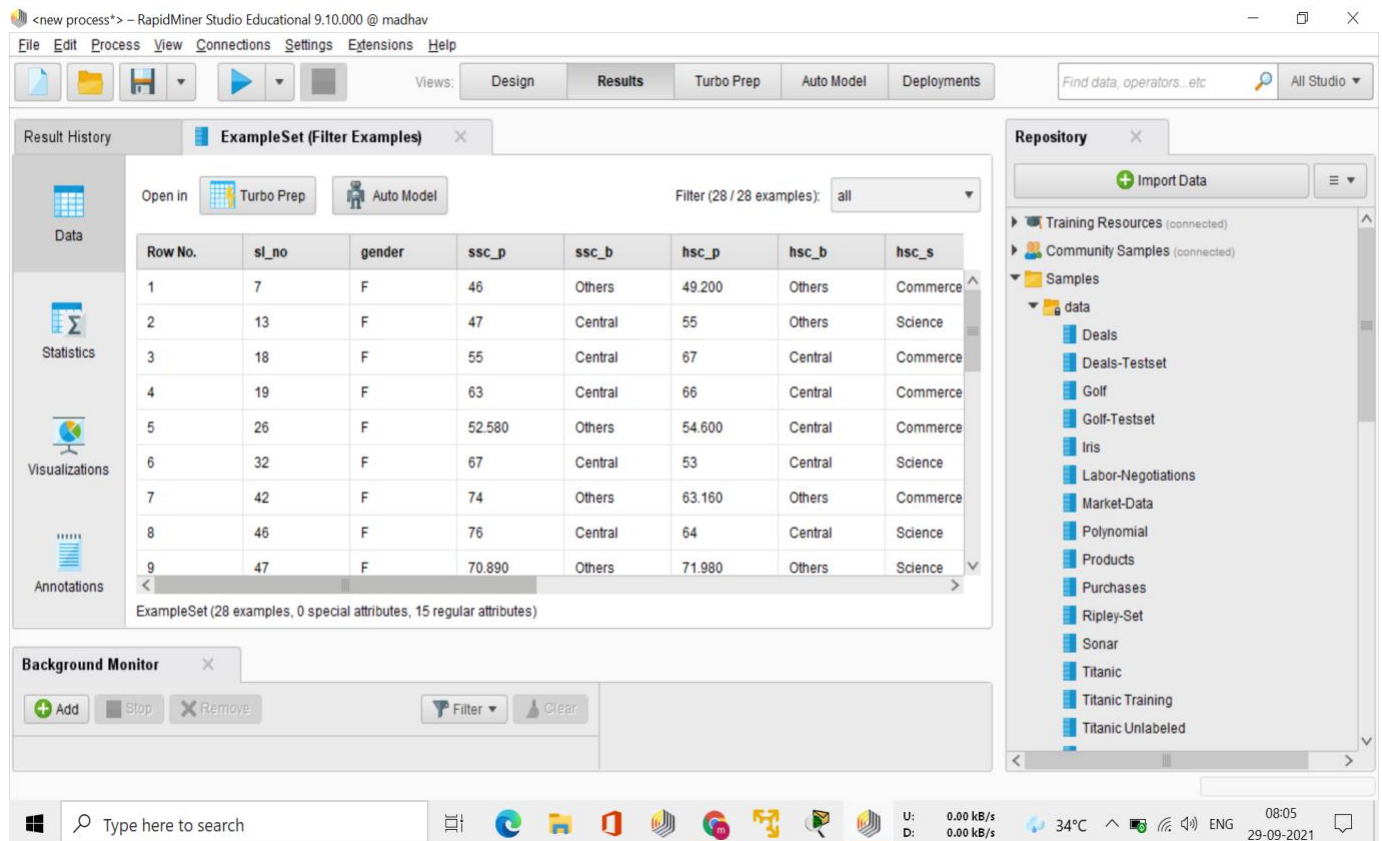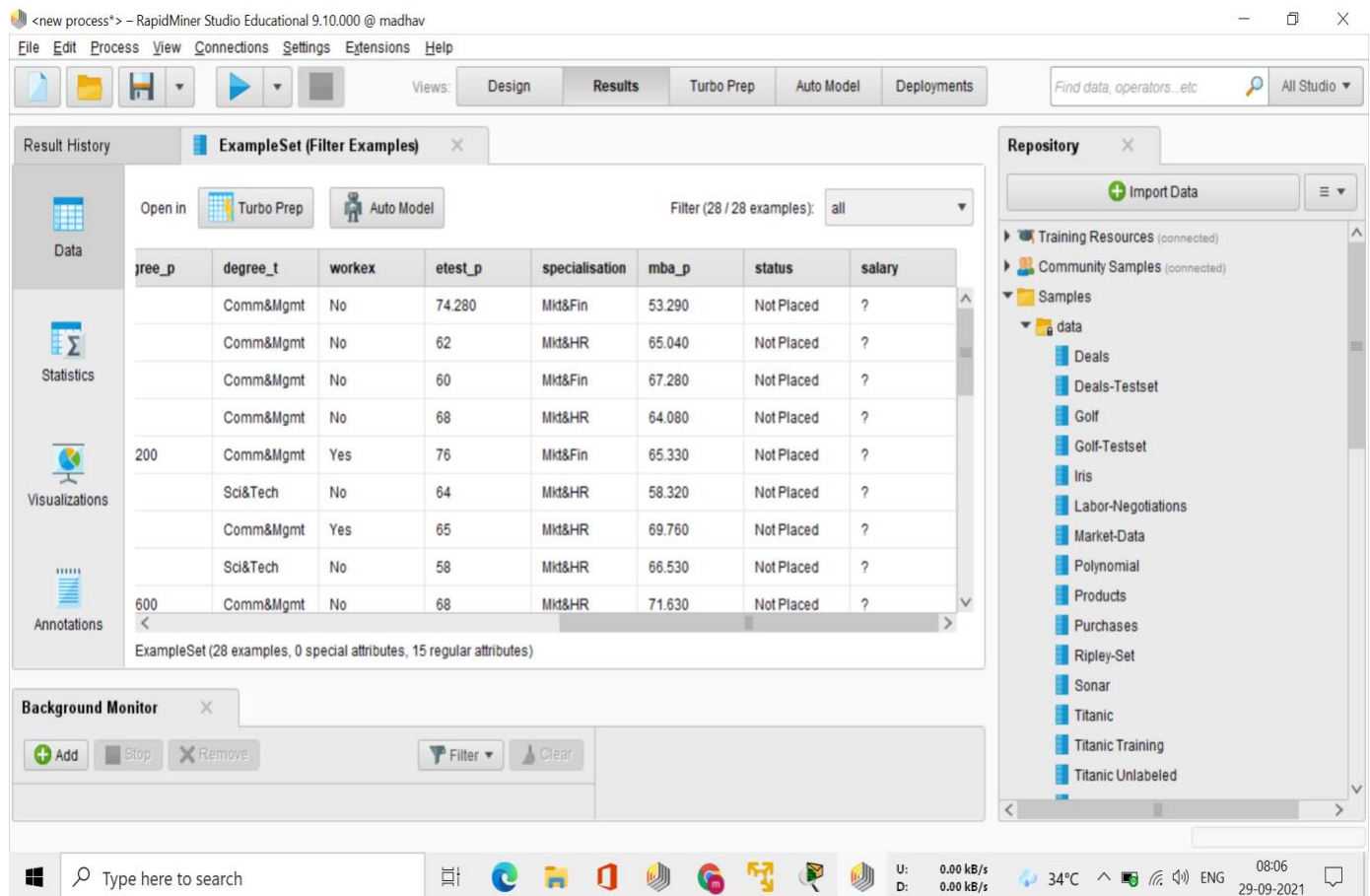Figure :-11 ,shows the filter examples operatoration on gender values which is equal to female.

Figure :-12 ,shows the filter examples operatoration on missing values of salary.

**Q2(B).What do you understand by graphical representation and statistics in rapidminer. Attach any 5 different types of graph and their interpretation.**

Answer:-

A statistical graph or chart is defined as the pictorial representation of statistical data in

graphical form. The statistical graphs are used to represent a set of data to make it easier to understand and interpret statistical information.

Graphical representation and statical based analysis helps us to view data in different view.

# Q3.What do you mean by preprocessing in Data Mining. Explain any 5 operations to handle missing values and attach screenshots of same.

Answer:-

Data preprocessing involves transforming raw data to well-formed data sets so that data mining analytics can be applied. Raw data is often incomplete and has inconsistent formatting .
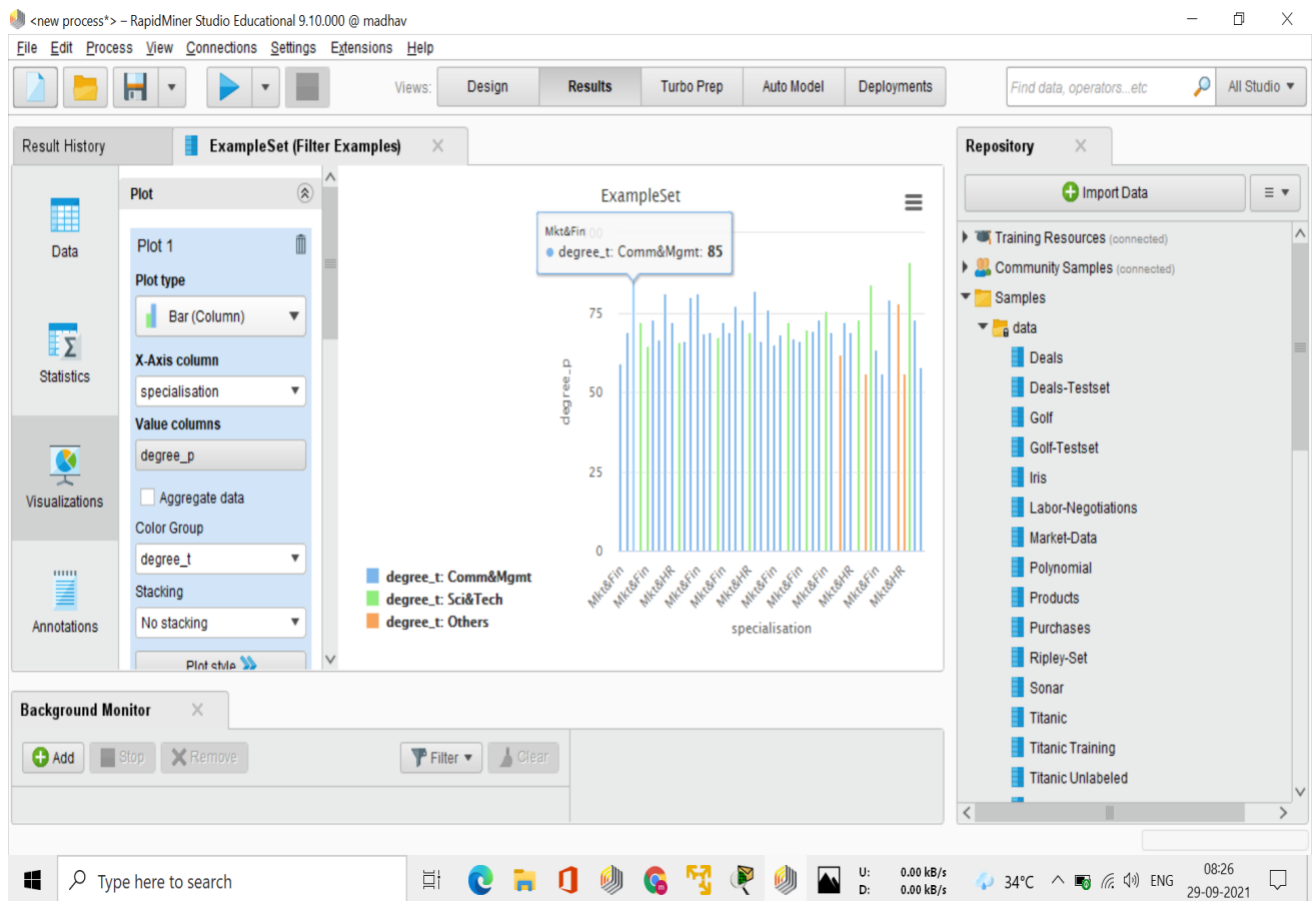Preprocessing involves both data validation and data imputation. The goal of data validation is to assess whether the data in question is both complete and accurate. The goal of data imputation is to correct errors and input missing values , either manually or automatically.

- There are few steps to preprocess:-
- Data cleaning
- Data integration
- Data transformation
- Data reduction
- Data DDiscretization

>Data cleaning:-

In data cleaning process we fill the missing,identify outliers and remove them,smooth out noisy data,and correct inconsistent data sets.

>Data integration: combines data from multiple sources into a coherent store.Careful integration can help reduce and avoid redundancies and inconsistencies in resulting data set.This can help improve the accuracy and speed of the subsequent data mining process.This is referred to as the entity identification problem.

>Data Transformation:-
Data is normalized and generalized. Normalization is a process that ensures that no data is redundant, it is all stored in a single place, and all the dependencies are logical.

>Data Reduction:-
When the volume of data is huge, databases can become slower, costly to access, and challenging to properly store. Data reduction aims to present a reduced representation of the data in a data warehouse.
Data reduction strategies:-
1)Dimensionality reduction
2)Numerosity reduction
3)Data compression

>Discretization :-

Reduce the number of values for a given continuous attribute by dividing the range of the attribute into intervals. Interval labels can then be used to replace actual data values.

>Five operations to handle missing values:-