

**Name : Sonali Singh**

**Class : B.E(B)**

**Roll No : 12**

**Experiment No. 1**

**Title:**

**For an organization of your choice, choose a set of business processes. Design star / snow flake schemas for analyzing these processes. Create a fact constellation schema by combining them. Extract data from different data sources, apply suitable transformations and load into destination tables using an ETL tool. For Example: Business Origination: Sales, Order, Marketing Process.**

**Objective:**

**Understand the basics of Star/Snowflake/fact Constellation schema & learn the RapidMiner tool for per perform various Operation on in-built or external Datasets.**

**Hardware Requirement:**

**Any CPU with Pentium Processor or similar, 256 MB RAM or more, 1 GB Hard Disk or more**

**Software Requirements:**

**32/64-bit Linux/Windows Operating System, latest RapidMiner Tool**

**Theory:**

**What does ETL mean?**

**ETL stands for Extract, Transform and Load. An ETL tool extracts the data from different RDBMS source systems, transforms the data like applying calculations, concatenate, etc. and then load the data-to-data Warehouse system. The data is loaded in the DW system in the form of dimension and fact tables.**

**Extraction**

**A staging area is required during ETL load. There are various reasons why staging area is required. The source systems are only available for specific period of time to extract data.**

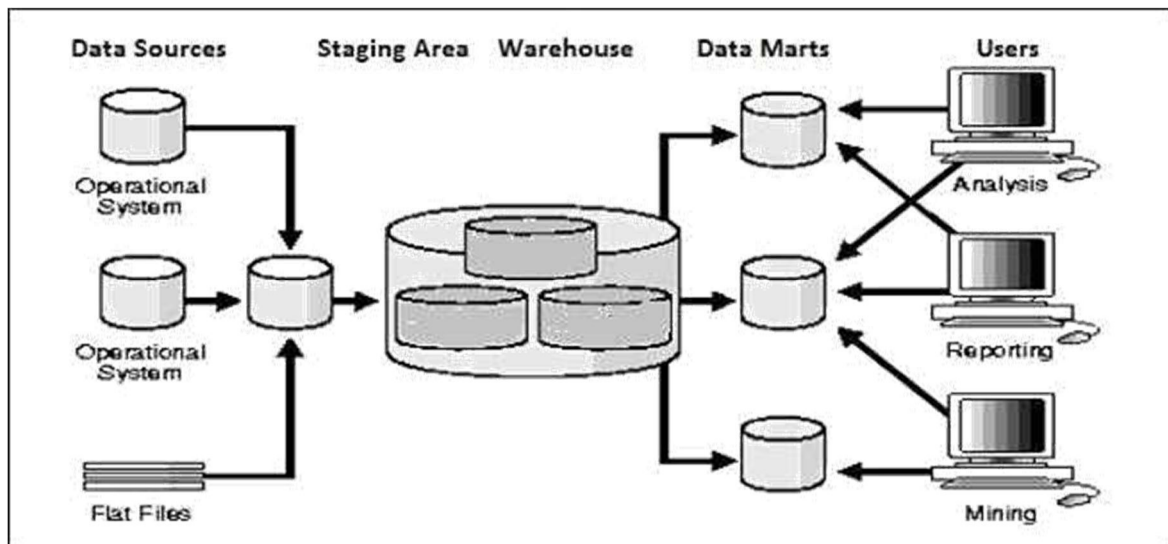
**This period of time is less than the total data-load time. Therefore, staging area allows you to extract the data from the source system and keeps it in the staging area before the time slot ends.**

**Staging area is required when you want to get the data from multiple data sources together or if you want to join two or more systems together. For example, you will not be able to perform a SQL query joining two tables from two physically different databases.**

**Data extractions' time slot for different systems vary as per the time zone and operational hours.**

**Data extracted from source systems can be used in multiple data warehouse system, Operation Data stores, etc.**

**ETL allows you to perform complex transformations and requires extra area to store the data.**



## Transform

In data transformation, you apply a set of functions on extracted data to load it into the target system. Data, which does not require any transformation is known as direct move or pass-through data.

You can apply different transformations on extracted data from the source system. For example, you can perform customized calculations. If you want sum-of-sales revenue and this is not in database, you can apply the SUM formula during transformation and load the data.

For example, if you have the first name and the last name in a table in different columns, you can use concatenate before loading.

## Load

During Load phase, data is loaded into the end-target system and it can be a flat file or a Data Warehouse system.

## Tool for ETL: RAPID MINER

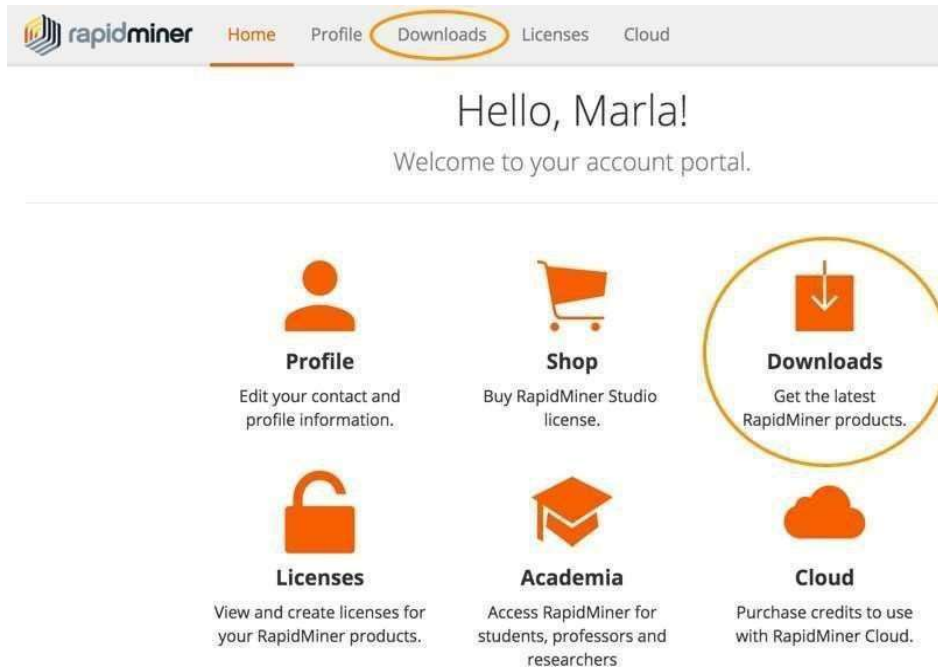
Rapid Miner is a world-leading open-source system for data mining. It is available as a stand-alone application for data analysis and as a data mining engine for the integration into own products. Rapid Miner is now Rapid Miner Studio and Rapid Analytics is now called Rapid Miner Server.

In a few words, Rapid Miner Studio is a "downloadable GUI for machine learning, data mining, text mining, predictive analytics and business analytics". It can also be used (for most purposes) in batch mode (command line mode)

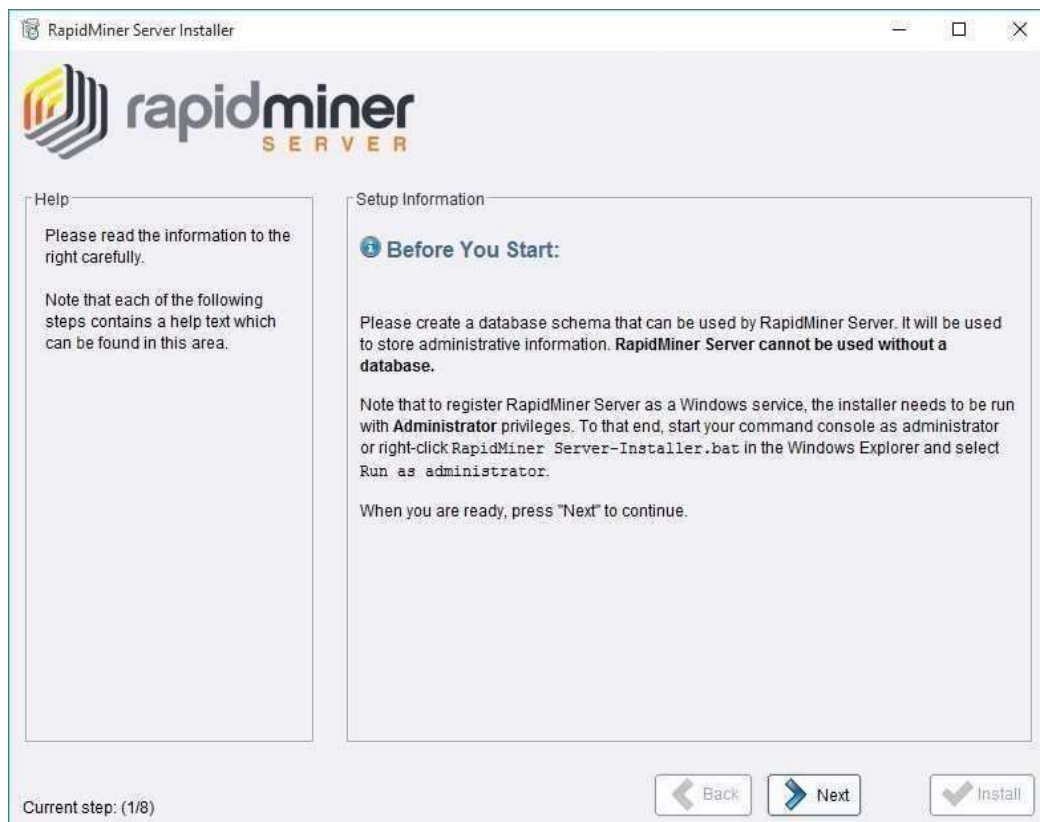
Rapid Miner Support to Nominal, Numerical values, Integers, Real numbers, 2-value nominal, multi-value nominal etc.

## STEPS FOR INSTALLATION:

### 1. Downloading Rapid Miner Server



### 2. Installing Rapid Miner Server



### 3. Configuring Rapid Miner Server settings



The image shows the 'RapidMiner Server Installer' window, specifically the 'Server Settings' tab. The window has a title bar with standard Windows controls. The main area is divided into a 'Help' section on the left and a 'Server Settings' section on the right. The 'Help' section contains text explaining the purpose of this step: specifying a host name and port for clients to connect to. It also mentions memory allocation and the need to specify the Java directory if the JAVA\_HOME environment variable is not set. The 'Server Settings' section contains several input fields and checkboxes. The 'Hostname' field is set to 'rapidminer.example.com'. There is a checkbox for 'Bind to this hostname only'. The 'Port for web interface' is set to '8080' and the 'Internal Port' is set to '5672'. A note states that the server web interface will be available at 'http://rapidminer.example.com:8080'. The 'Server Memory (in MB)' is set to '2048'. The 'Number of bundled Job Containers' is set to '1' and the 'Memory per Job Container (in MB)' is set to '2048'. A note states that the server will allocate memory up to '4,096 MB (System: 20,354 MB)'. There is a checkbox for 'Register as Windows service (needs administrator privileges)' which is checked. The 'Service ID' is set to 'RMS800SVC' and the 'Service Name' is set to 'RapidMiner\_Server\_8\_0\_0'. The 'JAVA\_HOME folder' is set to 'C:\Program Files\Java\jdk1.8.0\_40'. At the bottom, there are buttons for 'Back', 'Next', and 'Install'. The 'Current step' is indicated as '(6/9)' and the 'Version' is '8.0.0'.

**Help**

In this step, you can specify a host name and port under which clients, foremost RapidMiner Studio, will connect to RapidMiner Server. Therefore, you must choose a valid hostname. If you check "Bind to this hostname only", RapidMiner Server will listen only on the respective network interface.

Furthermore, you can assign the amount of memory utilized by RapidMiner Server (in MB) and optionally register it as a Windows service.

If you do not have the JAVA\_HOME Environment variable set, you need to specify your Java directory.

**Server Settings**

Hostname:  ☐ Bind to this hostname only

Port for web interface:  Internal Port:

Server web interface will be available at <http://rapidminer.example.com:8080>

Server Memory (in MB):

Number of bundled Job Containers:  Memory per Job Container (in MB):

RapidMiner Server will allocate memory up to **4,096 MB** (System: 20,354 MB)

☒ Register as Windows service (needs administrator privileges)

Service ID:  Service Name:

JAVA\_HOME folder:  

Current step: (6/9) Version : 8.0.0

### 4. Configuring Rapid Miner Server's database connection



The image shows the 'RapidMiner Server Installer' window, specifically the 'Database Configuration' tab. The window has a title bar with standard Windows controls. The main area is divided into a 'Help' section on the left and a 'Database Configuration' section on the right. The 'Help' section contains text explaining the purpose of this step: configuring the database connection which RapidMiner Server should use. It mentions entering the host or URL, port, and desired DB schema. It also mentions that the username and password can be filled in as needed, and that the user should select the appropriate JDBC driver and choose the driver class via the Dropdown menu. It also mentions that after setting everything up, the user can test the connection to the Database by clicking the Test Connection button. The 'Database Configuration' section contains several input fields and a dropdown menu. The 'Database host' is set to 'localhost' and the 'Database port' is set to '3306'. The 'Database schema' is set to 'rapidminer\_server'. The 'Database username' is set to 'rmUser' and the 'Database password' is masked with dots. There is an information icon and a note stating that the MySQL JDBC driver is not shipped with RapidMiner Server and that the user should click a link for more information. The 'JDBC Driver location' is set to 'C:\Apps\mysql-connector-java-5.1.38\mysql' and there is a folder icon next to it. The 'Database system' is set to 'MySQL'. There is a checkbox for 'Use relative path' which is unchecked. The 'JDBC driver class' is set to 'com.mysql.jdbc.Driver'. At the bottom, there is a 'Test Connection' button. At the bottom of the window, there are buttons for 'Back', 'Next', and 'Install'. The 'Current step' is indicated as '(7/8)'.

**Help**

In this step you can configure your Database connection which RapidMiner Server should use. You will need to enter the host or URL as well as the port and the desired DB schema. Username and Password can be filled in as needed. Then just select the appropriate JDBC driver and choose the driver class via the Dropdown menu. After you have set everything up, you can test the connection to the Database by clicking the Test Connection button.

**Database Configuration**

Database host:  Database port:

Database schema:

Database username:  Database password:

 MySQL JDBC driver is not shipped with RapidMiner Server. Please click [here](#) for more information!

JDBC Driver location:   Database system:

☐ Use relative path

JDBC driver class:



Current step: (7/8)

## 5. Installing Radoop Proxy



## 6. Completing the installation

Once logged in, complete the final installation steps.

1. From the SQL Dialect pull-down, verify that the database type displayed is the one you used to create the Rapid Miner Server database.
2. Verify the setting for the integrated Quartz scheduler, which is enabled by default.
3. Specify the path to the plug-in directory. You can install additional RapidMiner extensions by placing them in, or saving them to, this directory. Note that all extensions bundled with RapidMiner Studio are also bundled with Rapid Miner Server (no installation is necessary). These bundled extensions are stored in a separate directory that is independent of the path specified here. Be sure that you have write permission to the directory.
4. Specify the path to upload directory. This is the directory where RapidMiner Server stores temporary files needed for processes. The installation process creates a local uploads directory in the installation folder. However, if you install Rapid Miner Server on a relatively small hard disk and, for example, use many file objects in processes or if you have large resulting files, consider creating a directory elsewhere in the cluster to store the temporary files. Be sure that you have write permission to the directory.
5. Click Start installation now.
6. Installation gets completed.

### Data Warehousing Schemas

1. Star Schema
2. Snowflake Schema
3. Fact Constellation

#### Star Schema

For example, as you can see in the above-given image that fact table is at the center which contains keys to every dimension table like Deal\_ID, Model ID, Date\_ID, Product\_ID, Branch\_ID & other attributes like Units sold and revenue.

## Characteristics of Star Schema:

Every dimension in a star schema is represented with the only one- dimension table. The dimension table should contain the set of attributes.

The dimension table is joined to the fact table using a foreign key The dimension table are not joined to each other

Fact table would contain key and measure

The Star schema is easy to understand and provides optimal disk usage.

The dimension tables are not normalized. For instance, in the above figure, Country\_ID does not have Country lookup table as an OLTP design would have.

The schema is widely supported by BI Tools

## Snowflake Schema

A Snowflake Schema is an extension of a Star Schema, and it adds additional dimensions. It is called snowflake because its diagram resembles a Snowflake.

The dimension tables are normalized which splits data into additional tables. In the following example, Country is further normalized into an individual table.

## Characteristics of Snowflake Schema:

The main benefit of the snowflake schema it uses smaller disk space. Easier to implement a dimension is added to the Schema

Due to multiple tables query performance is reduced

The primary challenge that you will face while using the snowflake Schema is that you need to perform more maintenance efforts because of the more lookup tables.

Star Schema	Snow Flake Schema
Hierarchies for the dimensions are stored in the dimensional table.	Hierarchies are divided into separate tables.
It contains a fact table surrounded by dimension tables.	One fact table surrounded by dimension table which are in turn surrounded by dimension table
In a star schema, only single join creates the relationship between the fact table and any dimension tables.	A snowflake schema requires many joins to fetch the data.
Simple DB Design.	Very Complex DB Design.
Denormalized Data structure and query also run faster.	Normalized Data Structure.

**High level of Data redundancy**

**Very low-level data redundancy**

Single Dimension table contains aggregated data.

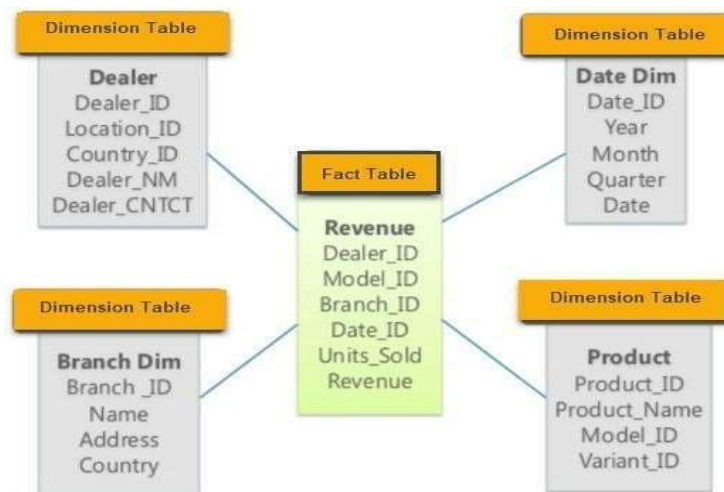
Data Split into different Dimension Tables.

Cube processing is faster.

Cube processing might be slow because of the complex join

Offers higher performing queries using Star Join Query Optimization. Tables may be connected with multiple dimensions.

The Snow Flake Schema is represented by centralized fact table which unlikely connected with multiple dimensions.



Star Schema

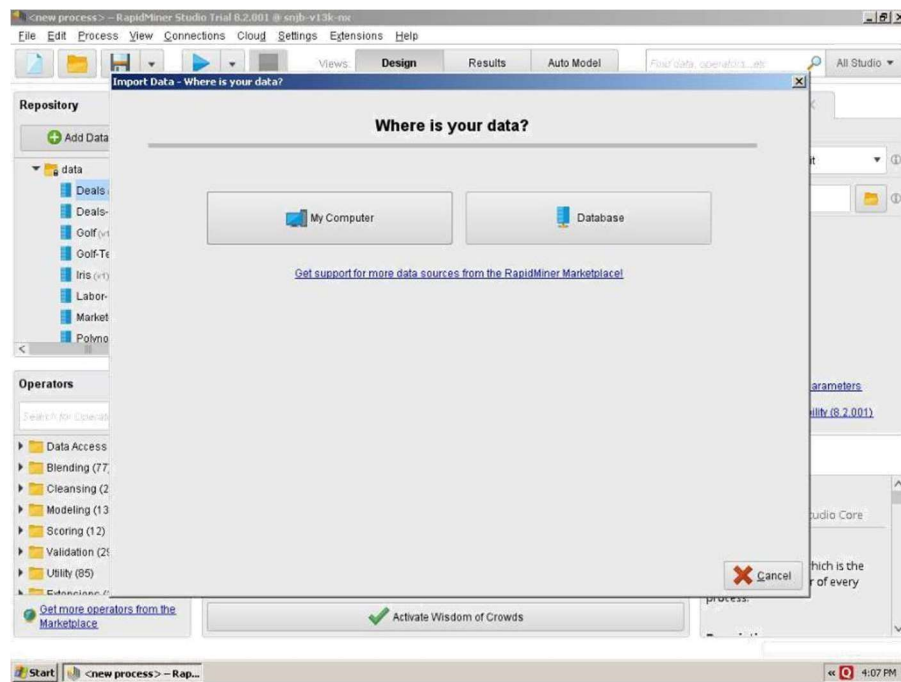
## 1. Design Model

**Step**

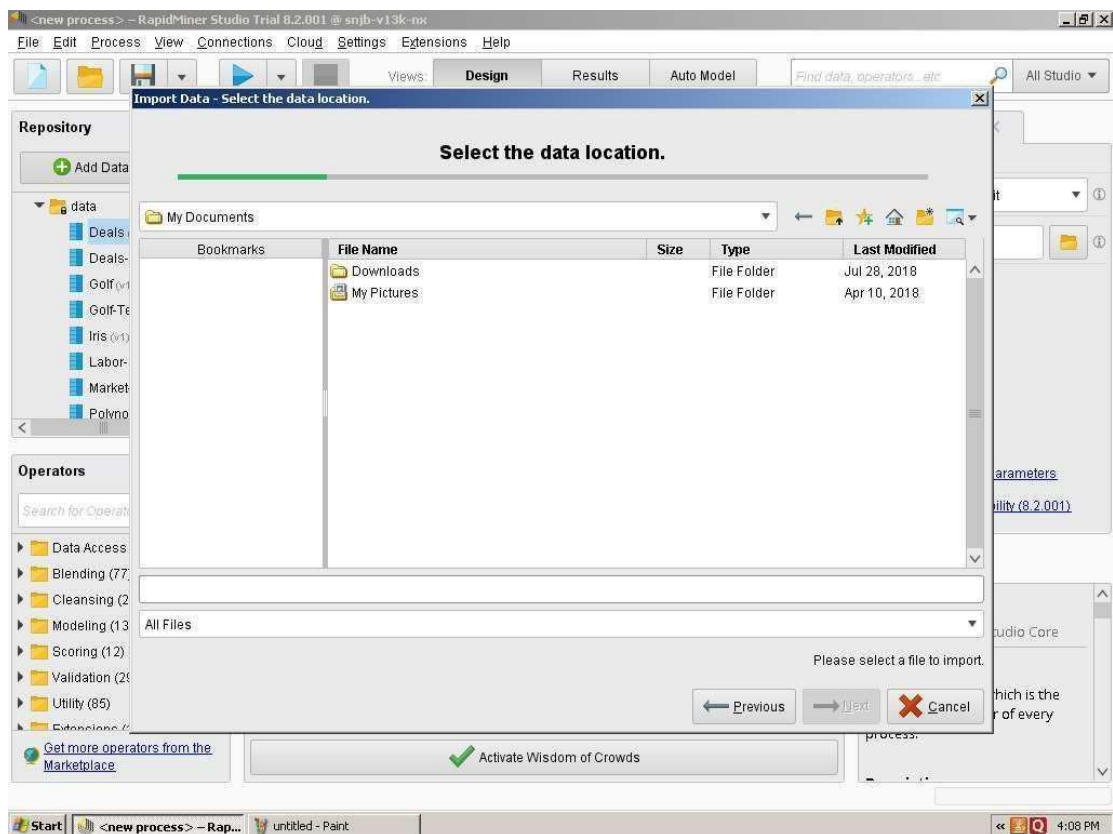
The screenshot shows the RAPID MINER GUI interface. Annotations point to the following components:

- Navigate repositories**: Points to the 'Repository' pane on the left, which lists 'Samples', 'DB', 'Local Repository (connected)', 'Personal (connected)', and 'Cloud Repository (disconnected)'.
- Available operators**: Points to the 'Operators' pane on the left, which lists various operators like 'Data Access (46)', 'Blending (77)', 'Cleaning (26)', 'Modeling (325)', 'Scoring (10)', 'Validation (30)', and 'Utility (85)'.
- Log of activities, including errors. If this is missing, add from View/Show Panel**: Points to the 'Log' pane at the bottom left.
- Process design window**: Points to the central 'Process' pane, which displays the message: 'Your process looks empty. Add some data first. Drag data or operators here.'
- Explanation of the selected operator**: Points to the 'Help' pane on the right, which shows the 'Process' operator's synopsis: 'The root operator which is the outer most operator of every contract.'
- Parameter settings for selected**: Points to the 'Parameters' pane on the right, which shows settings for the 'Process' operator, such as 'logverbosity' (set to 'info'), 'logfile', 'resultfile', 'random seed' (2001), and 'send mail' (set to 'never').

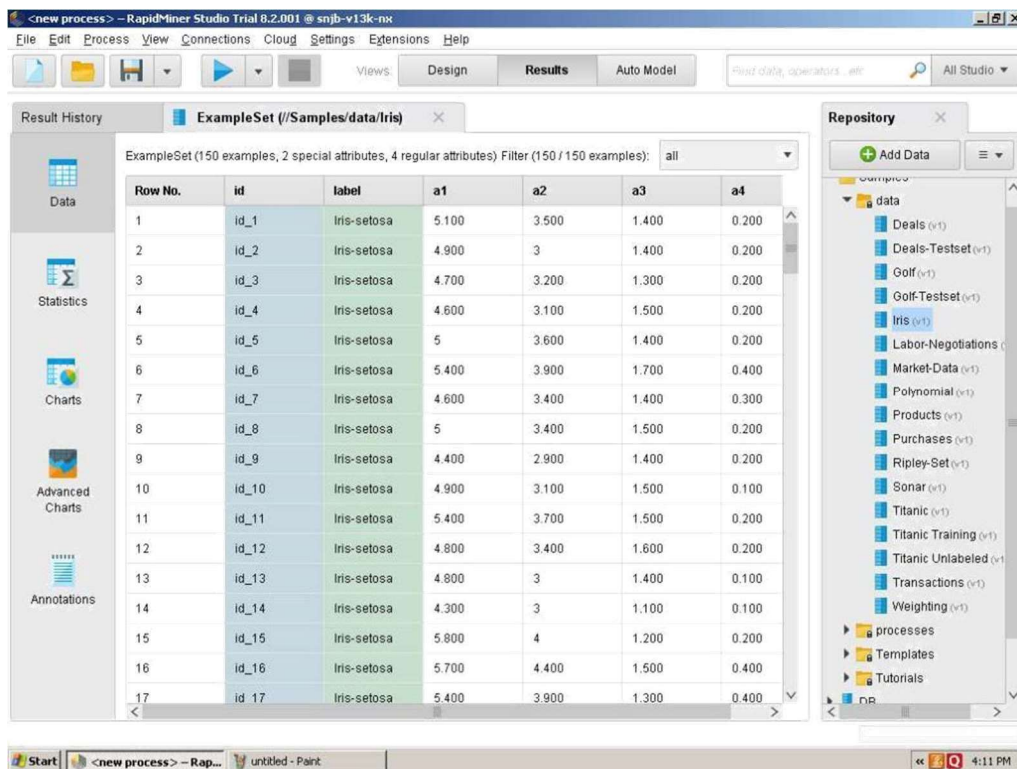




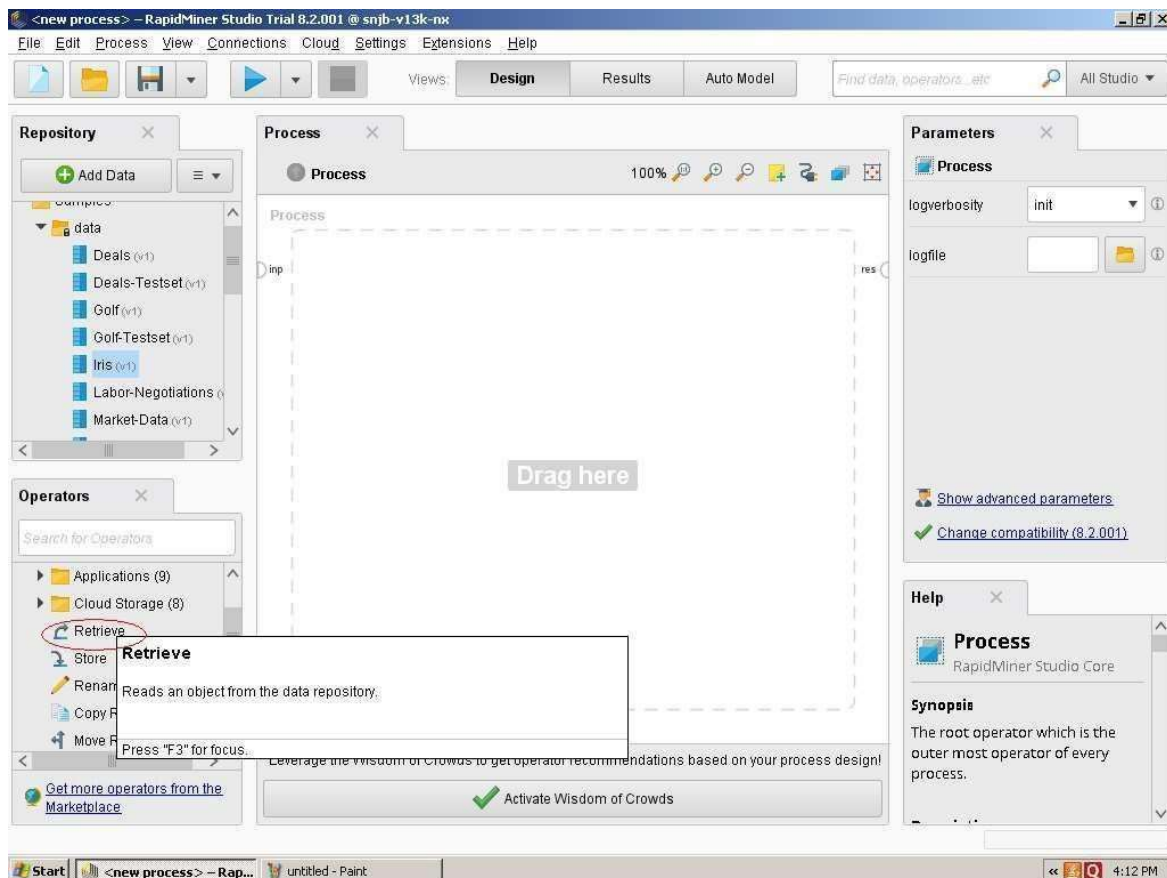
## Step-2 Select Data Location



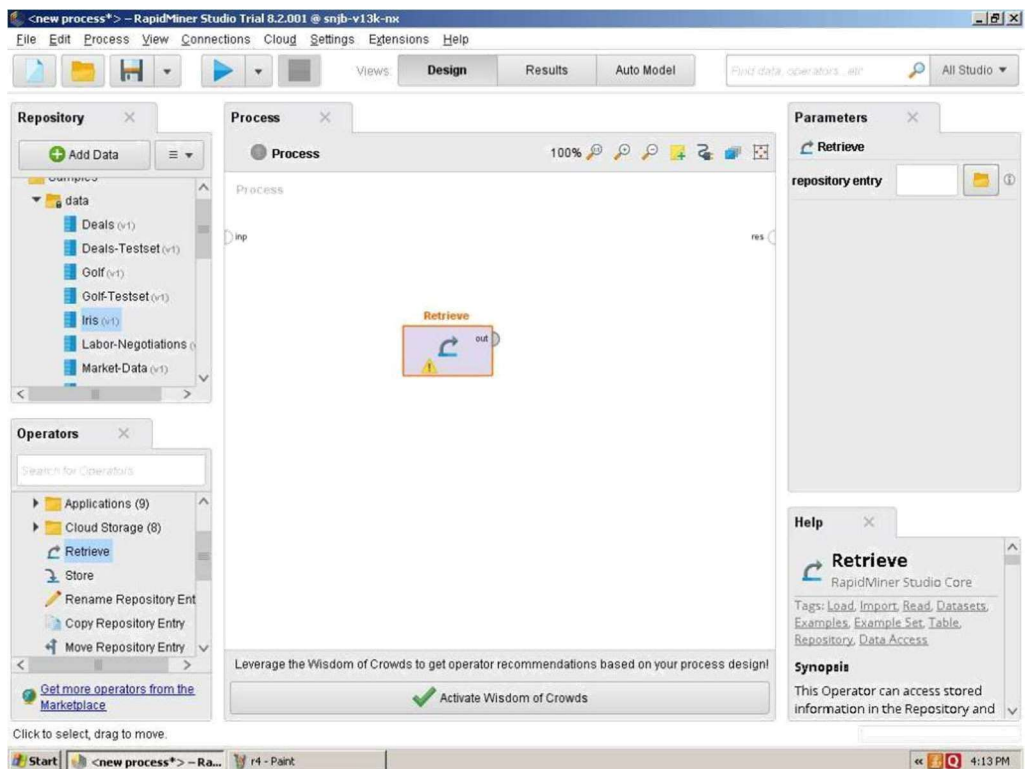
## Step-3 Open Sample Data Set e.g. Iris dataset available inbuilt with tool



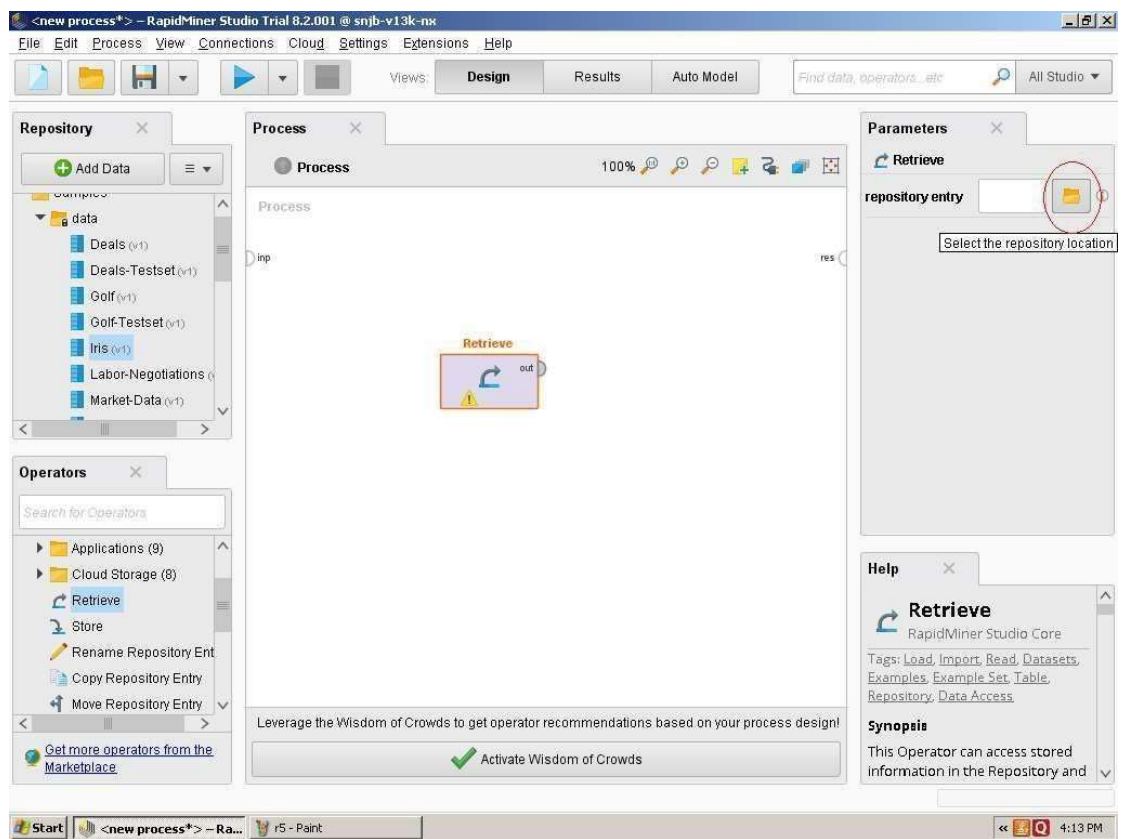
#### Step-4 Click on retrieve Operator Drag in Process View



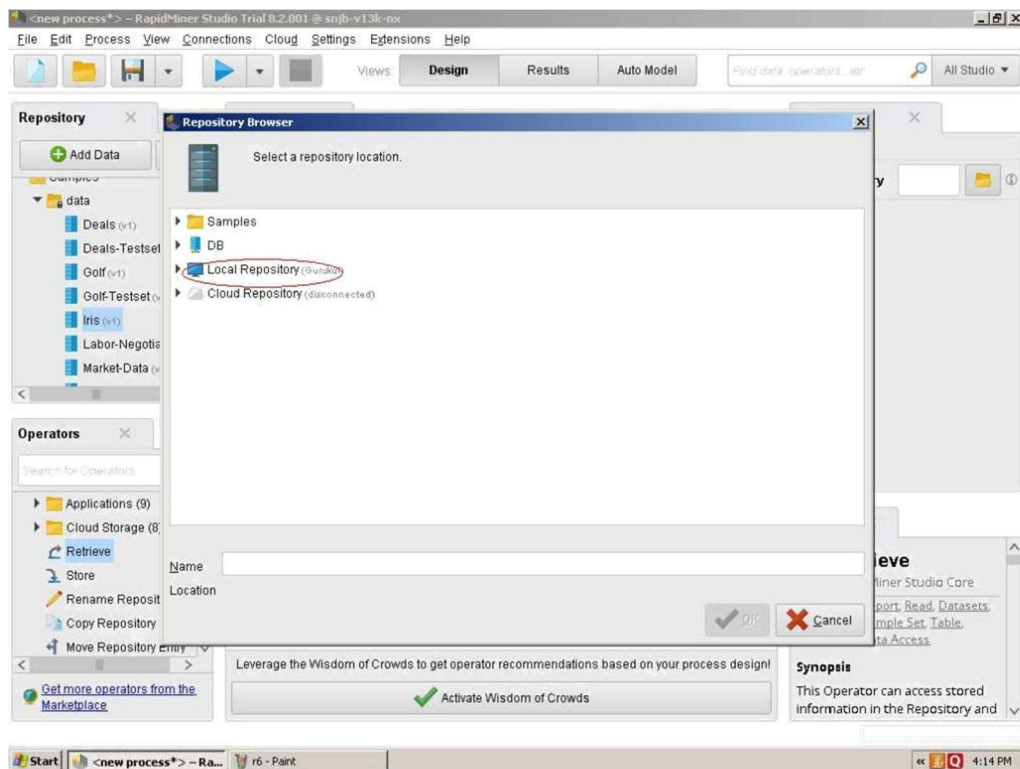
#### Step-5 Retrieve icon shows in Process View it has input and out Operator



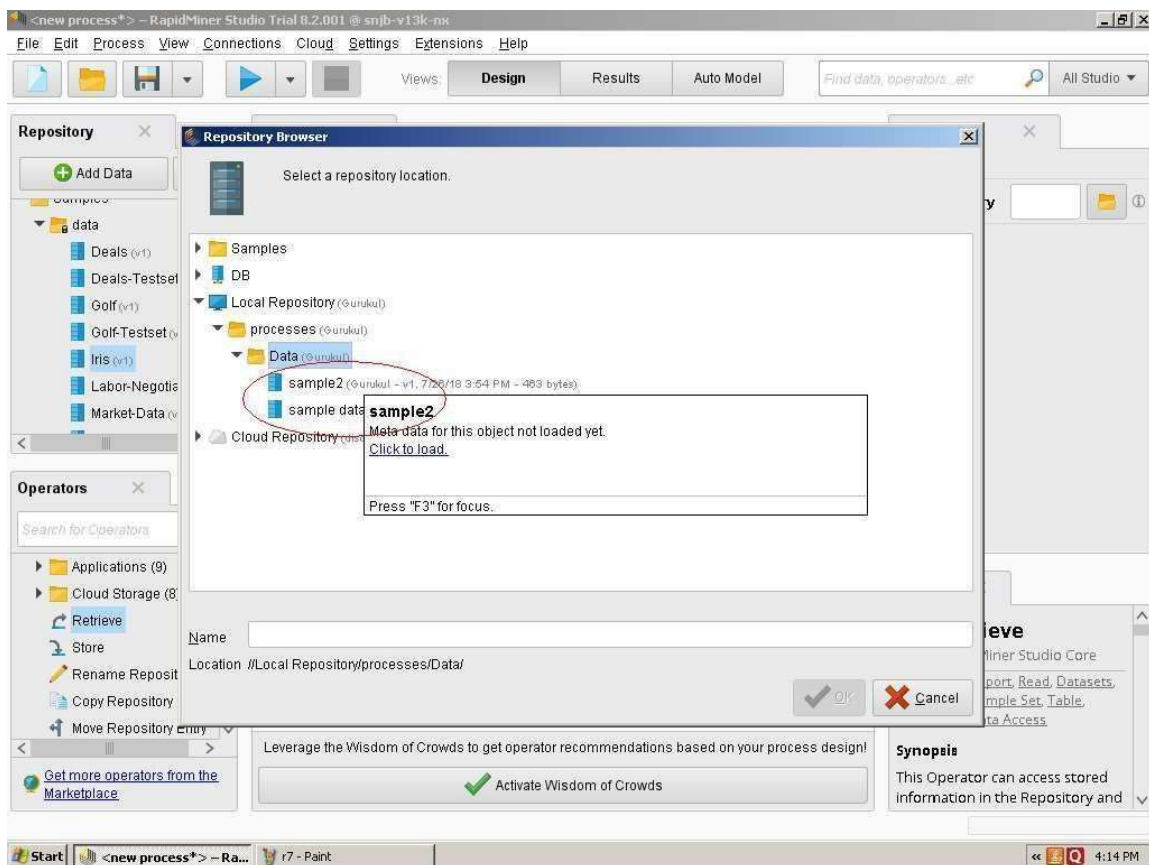
## Step-6 Click on repository entry



## Step-7 Select Local Repository

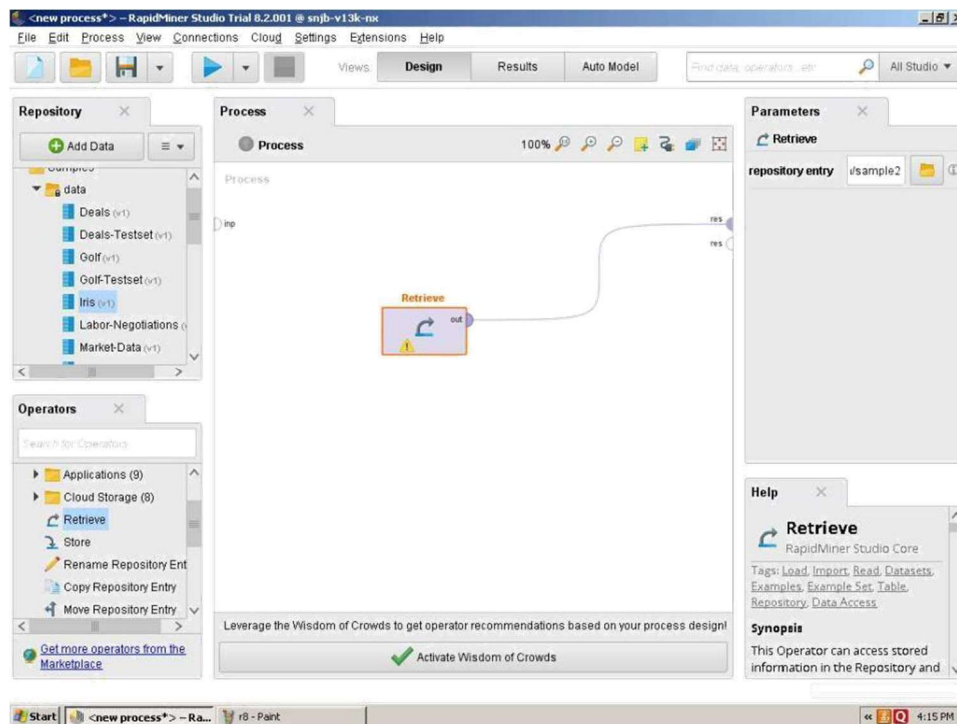


## Step-8 Select Sample file

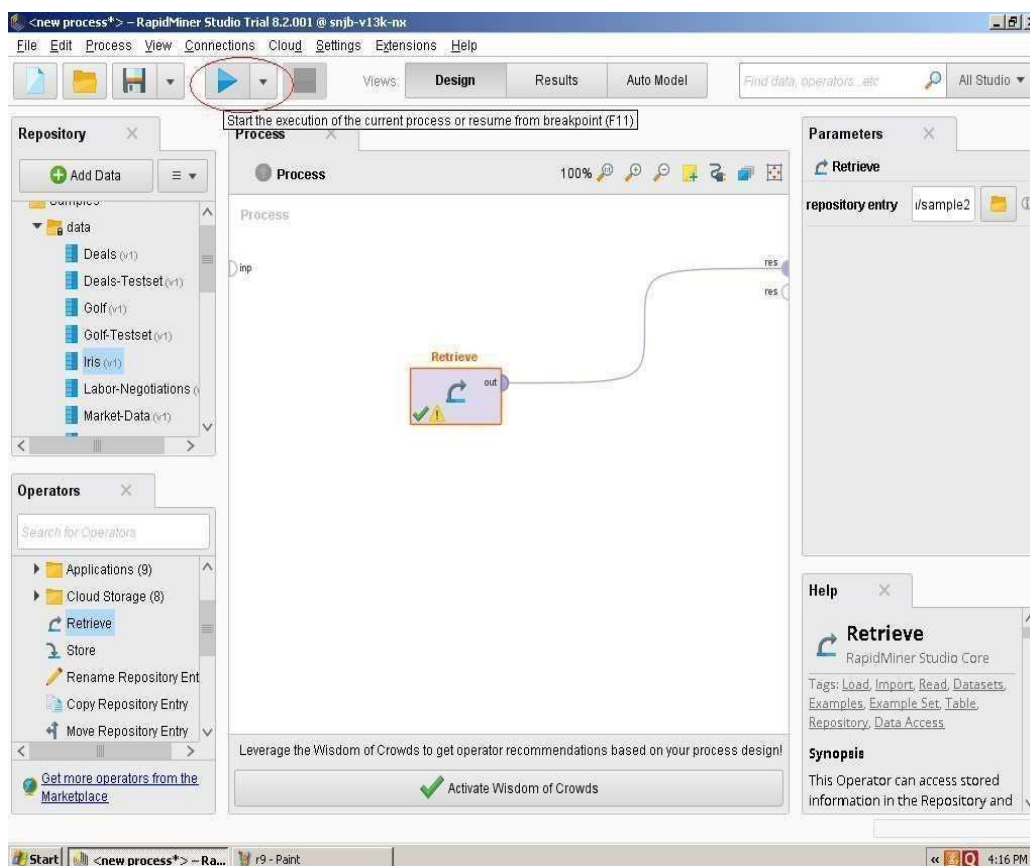


## Step-9 Join Out Operator to result Operator

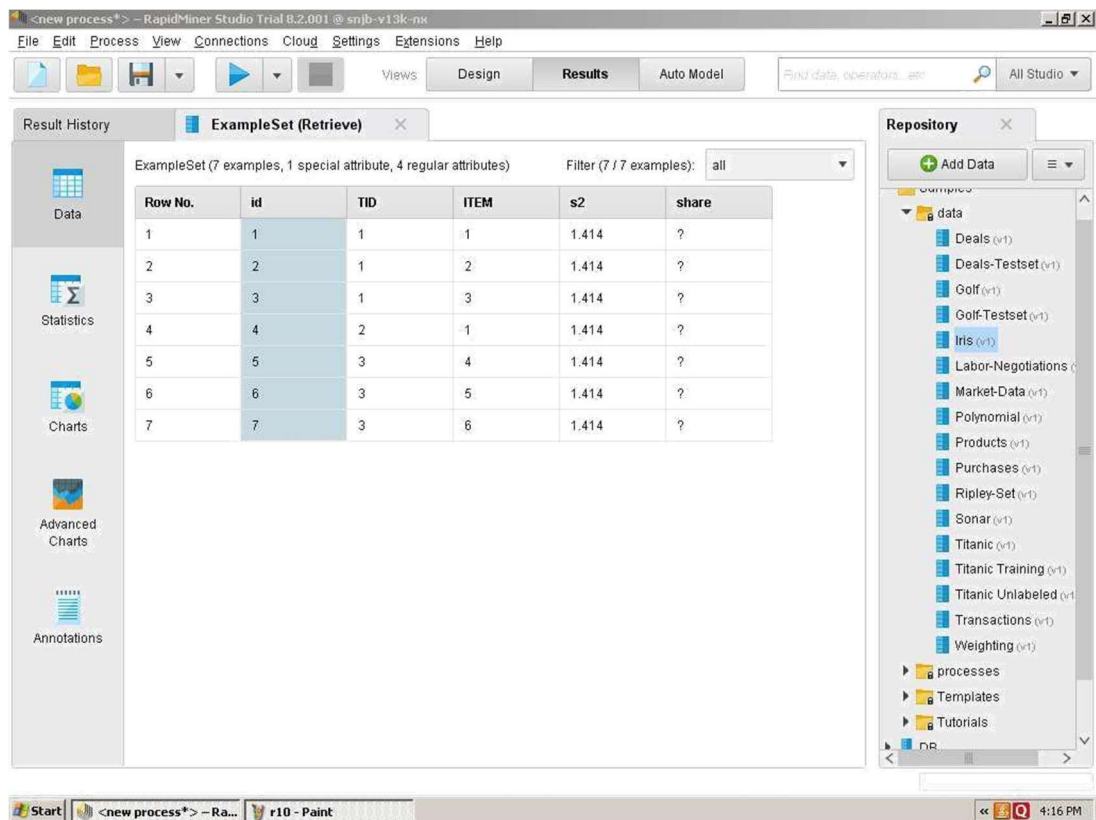




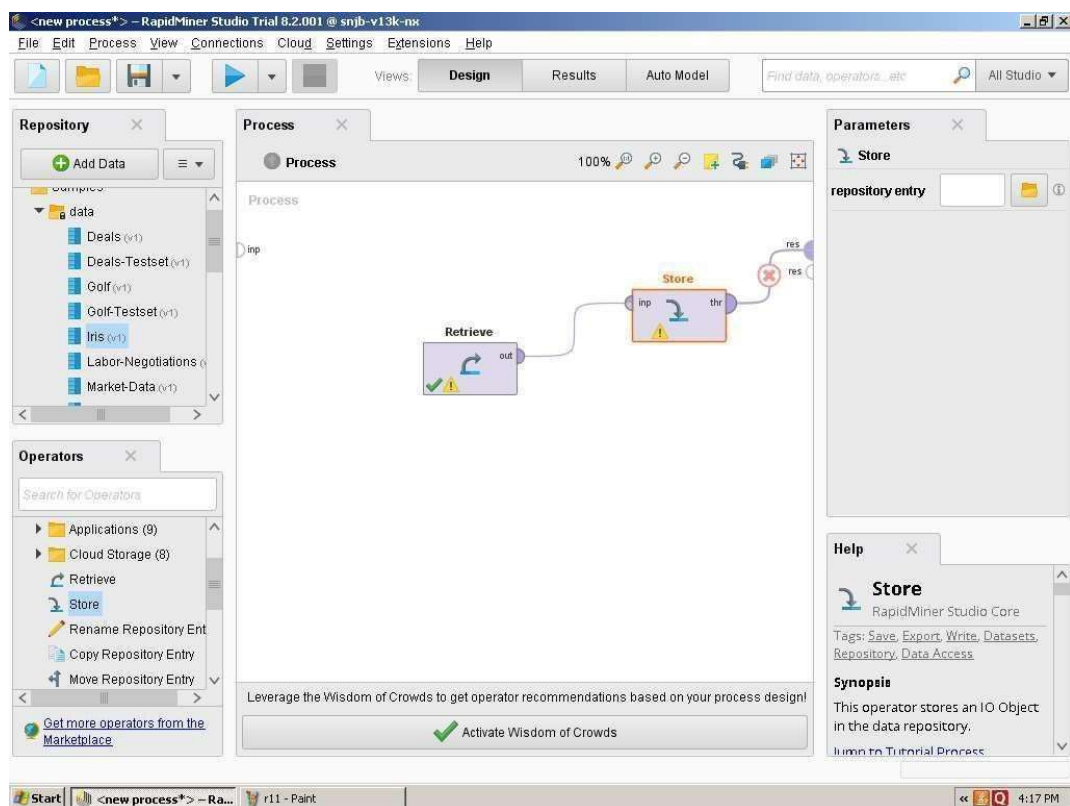
## Step-10 Start Execution of Current Process



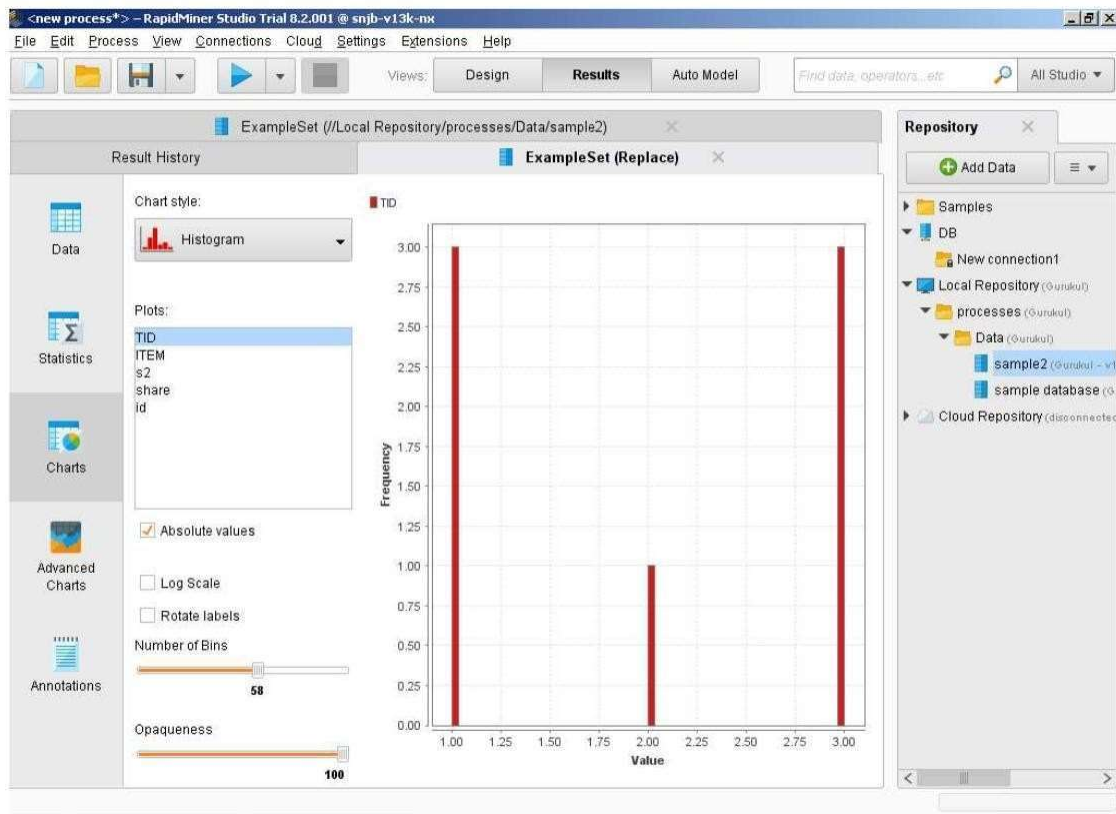
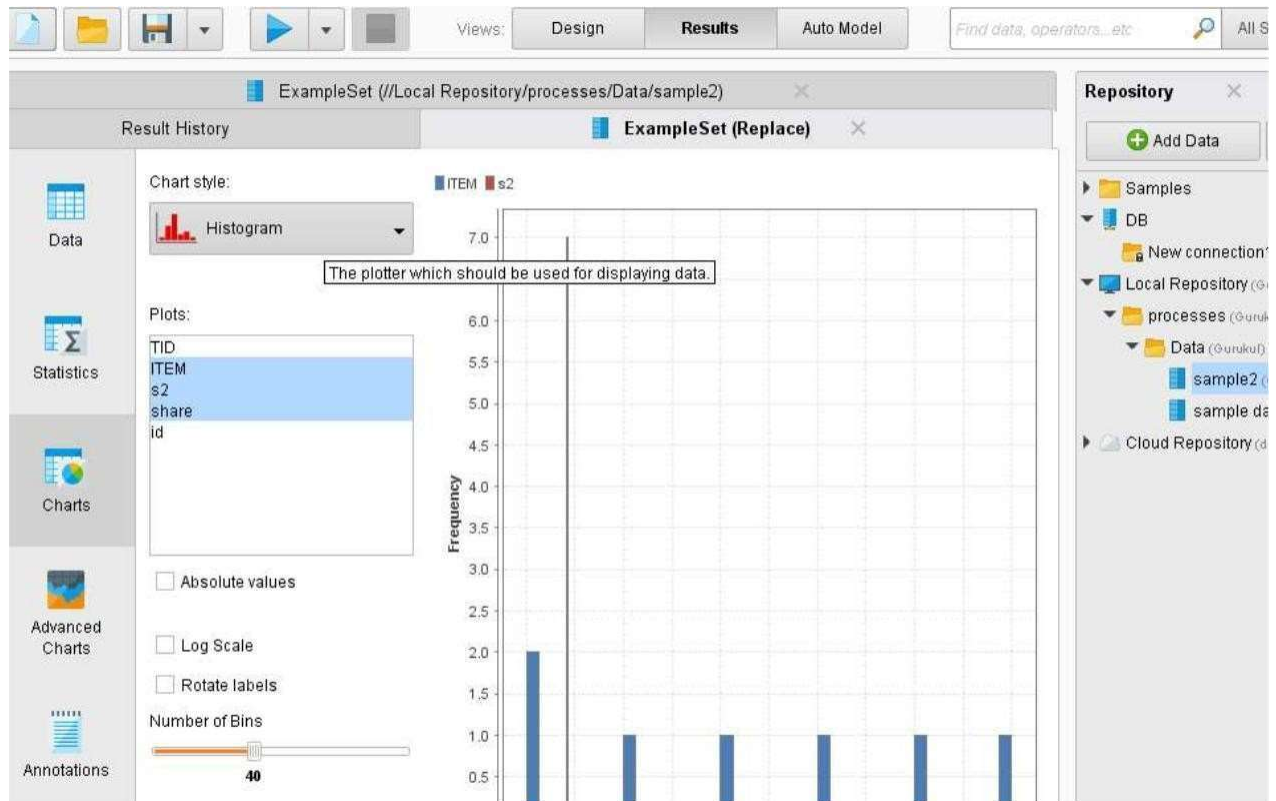
## Step-11 Output Result Generated after Execution of Current Process



**Step-12 Now you can add Store Operator and connect to result operator**



**Step-13 You can also plot Charts of Sample Data set**



A nice functionality for data preparation, called RapidMiner Turbo Prep, is where you simply drag and drop

**data to create amazing interfaces.**

**Conclusion:**

**Hence, we are able to study RapidMiner Tools us can Perform ETL operations on SampleData sets and can perform analysis on sample data sets.**