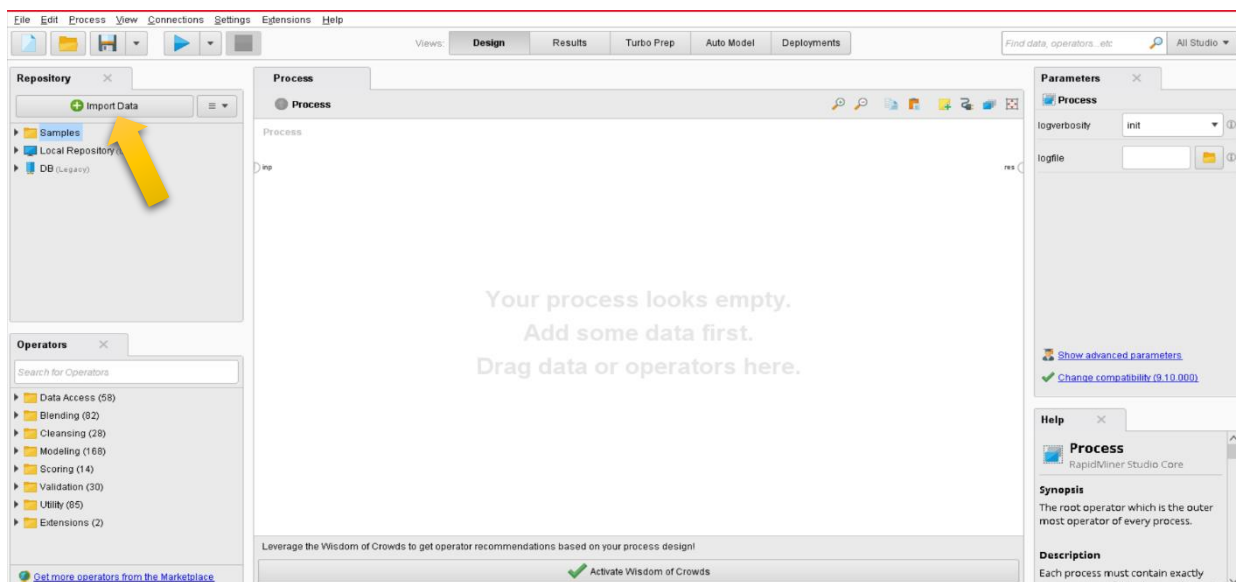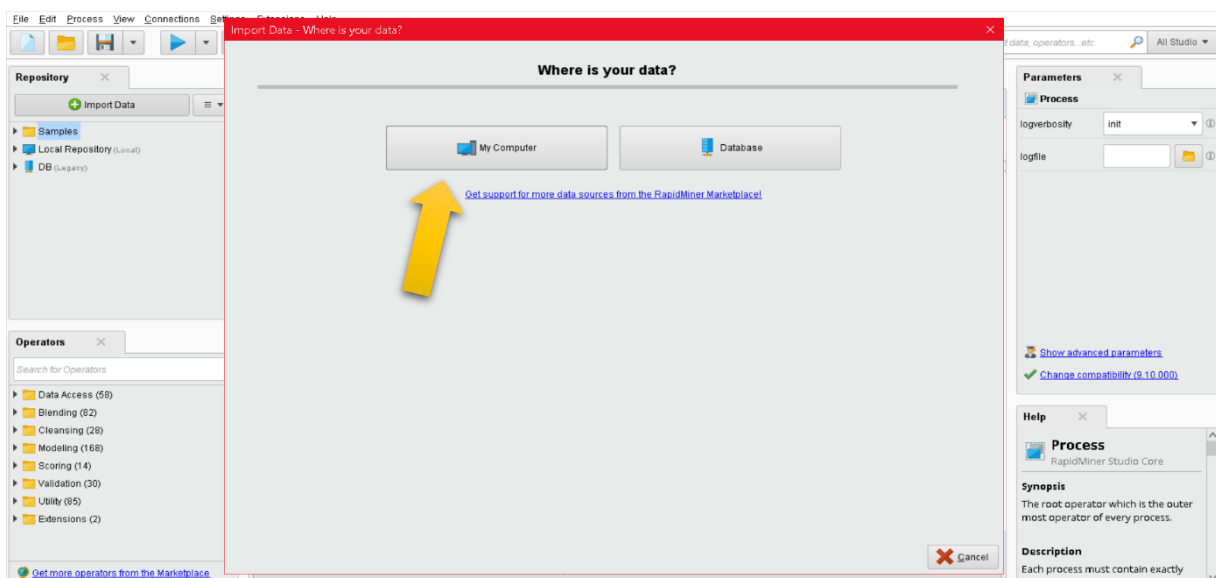## Q1. Perform the steps to import data from your local computer.
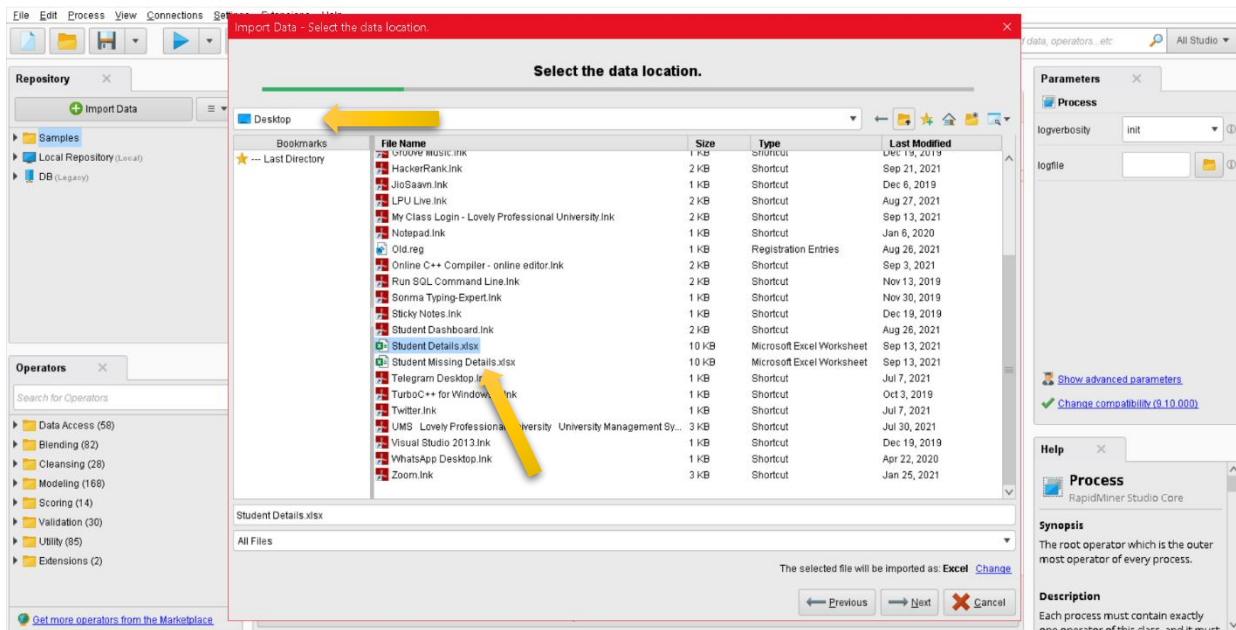
**Answer:** For importing data, first we have a data set in our system.
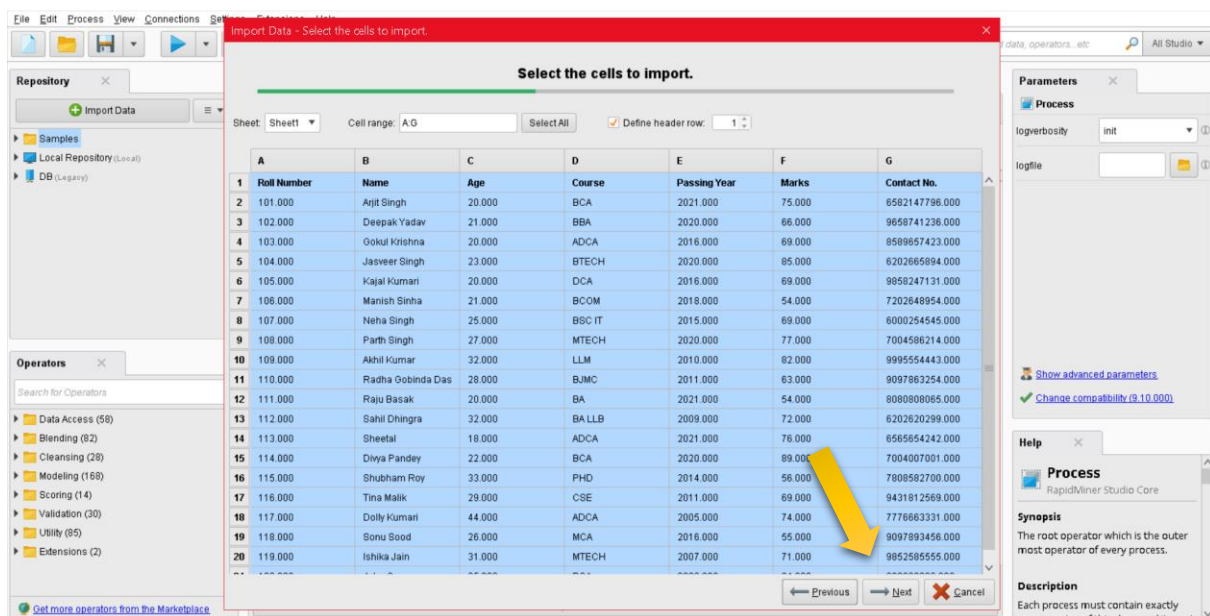**Step 1:** Open rapid miner and click on **Import Data.**
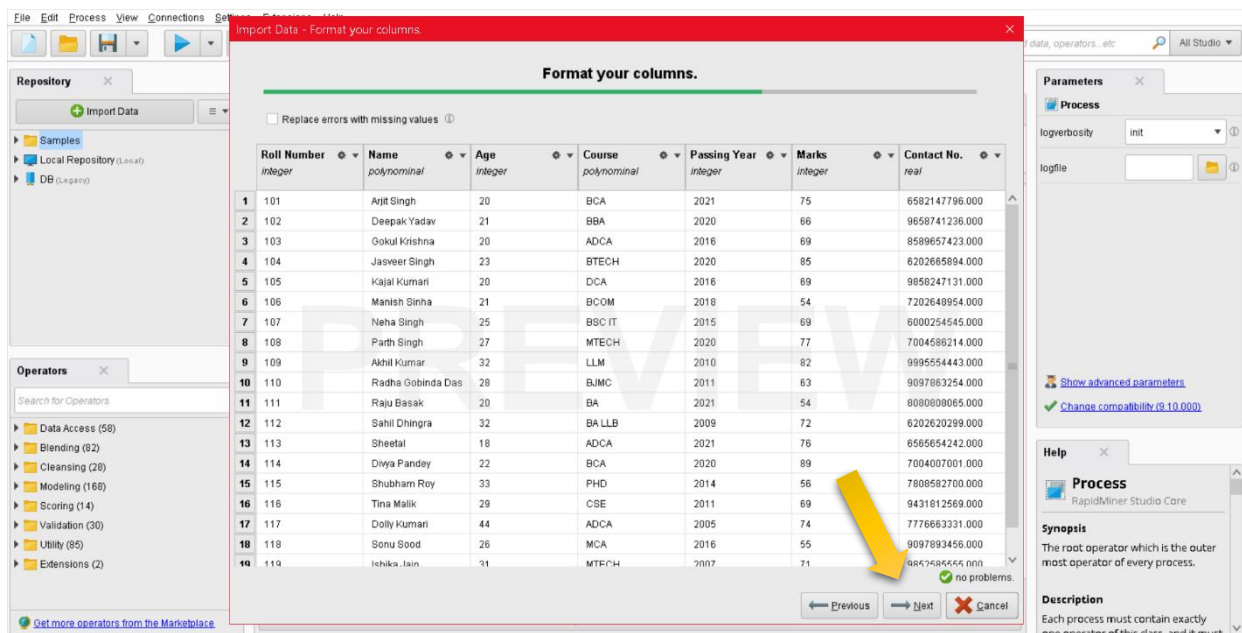


**Step 2:** Now click on **My computer**.
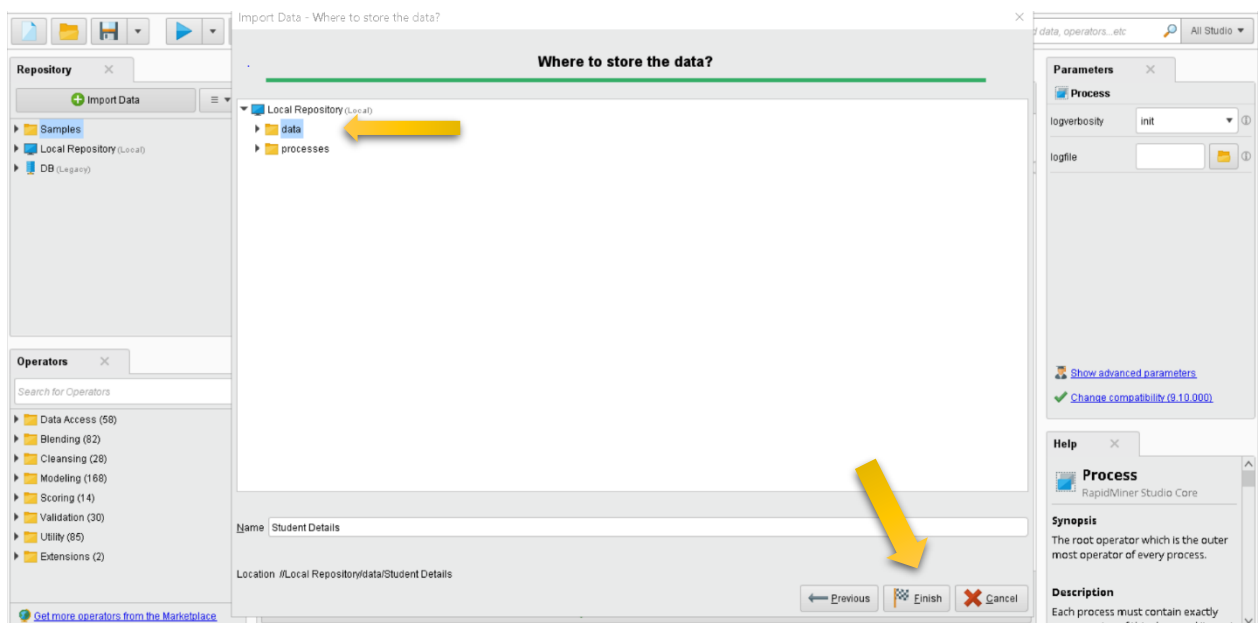
**Step 3:** Select your **Data Location**.



**Step 4:** Now your data set is showing, select the cells that you want to import and then click on **Next**.
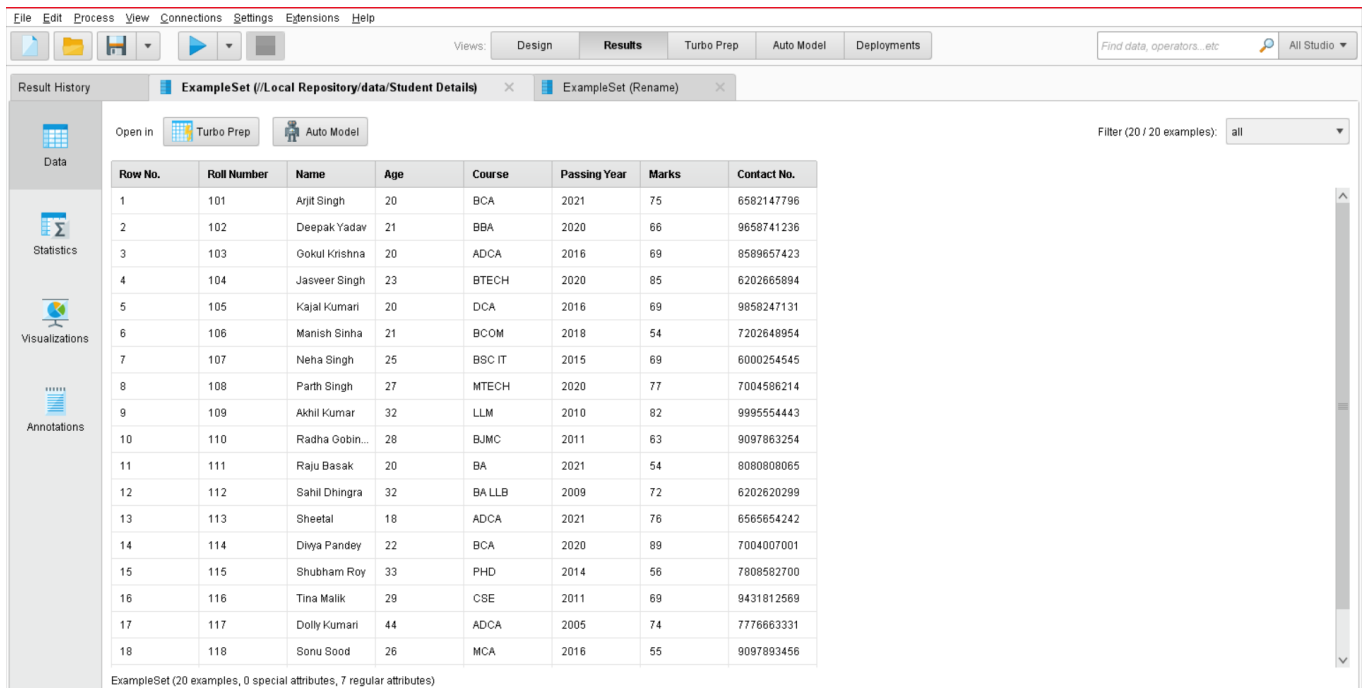
**Step 5:** Your selected data set has been loaded, now click on **Next**.



**Step 6:** Now select where you want to store your data set, you have two options **data/processes,** here I selected **data** and then click on **Finish**.
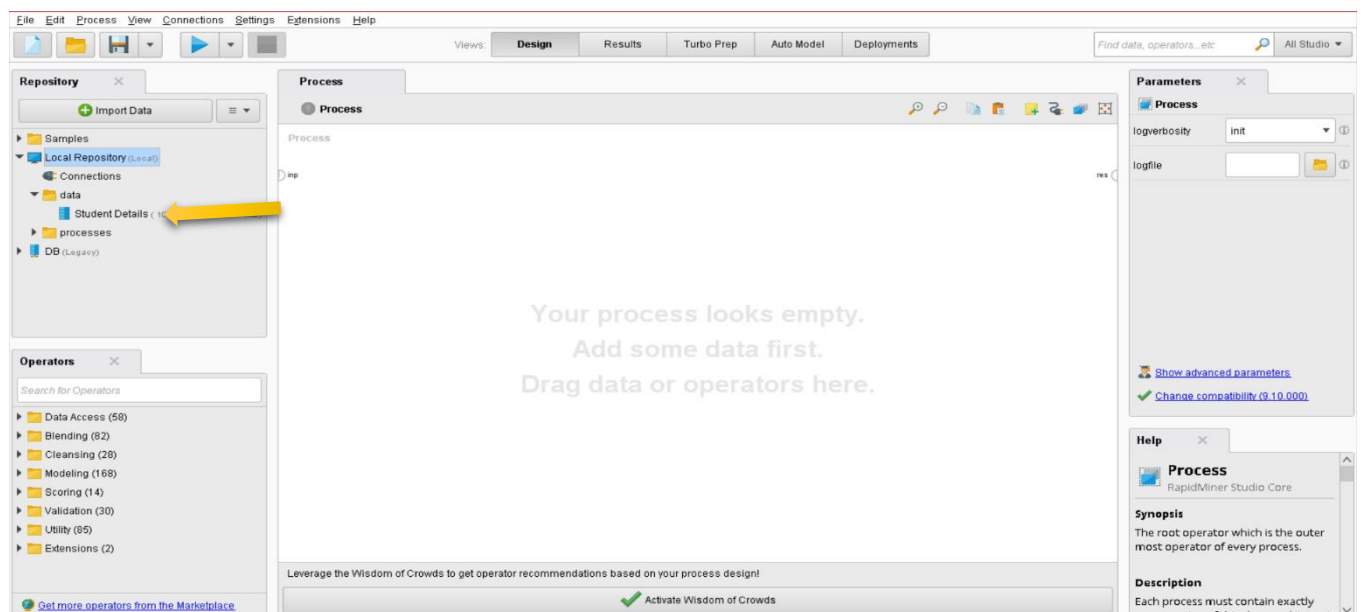
**Step 7:** Now you can see your data set has been **Imported Successfully.**



**You can also see that your data set has been stored in the data section of the local repository as well.**

| Course Code: **CAP 447** | Course Title: **Data Warehousing and Data Mining Lab** |
|---|---|
| Course Instructor: **Dr. Geeta Sharma** | |
| Student's Roll no: **RD2110B79**      Student's Reg. no: **12102801** | |
| Name: **Atul Kumar** | |
| Question No. **02**      Page No. **05**      Total Pages. **17** | |

## Q2. Perform following Transformation operations on your data you imported in Q1:

a) Sorting
b) Filter numerical data
c) Filter String data
d) Remove attribute

**Answer:** For performing all the following transformation in the Imported data set first we have to click on **Turbo Prep** and then click on **Load Data** and select the imported data set. (Imported data sets are in the local repository.) Now click on **TRANSFORM**.

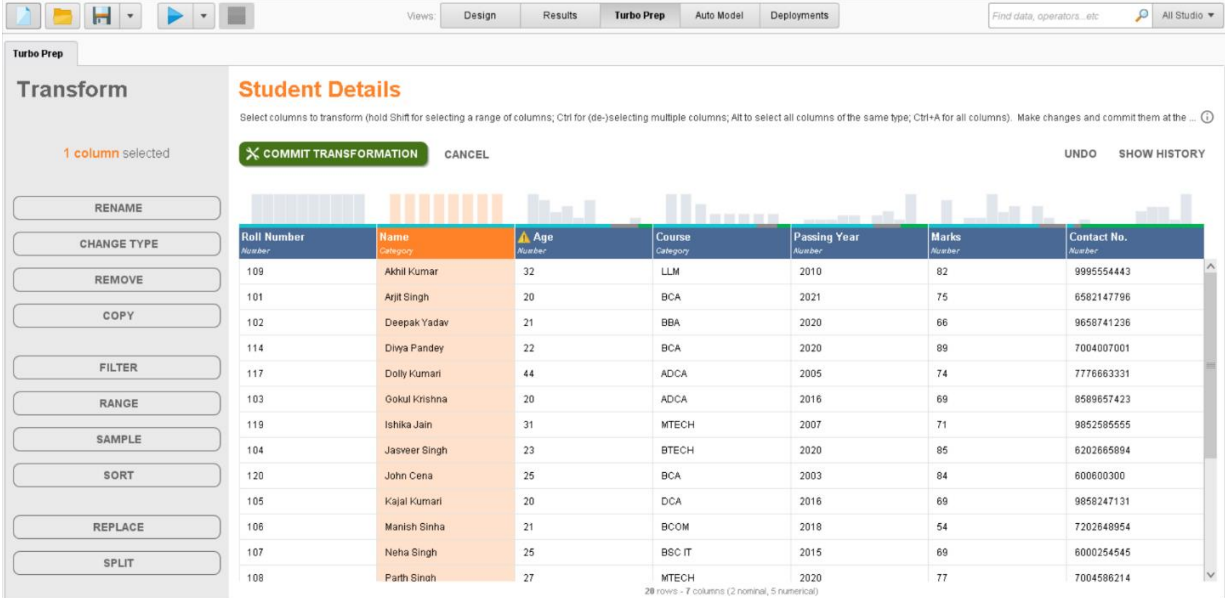| Course Code: **CAP 447** | Course Title: **Data Warehousing and Data Mining Lab** |
|---|---|
| Course Instructor: **Dr. Geeta Sharma** | |
| Student's Roll no: **RD2110B79** | Student's Reg. no: **12102801** |
| Name: **Atul Kumar** | |
| Question No. **02** | Page No. **06** | Total Pages. **17** |

This interface will come up after clicking on **TRANSFORM**.



### a) Sorting

**Step 1:** First you have to select the attribute where you want to use the sorting then click on **Sort**. Now we have two options sorted in **Ascending/Descending.** Here I selected **Ascending.**

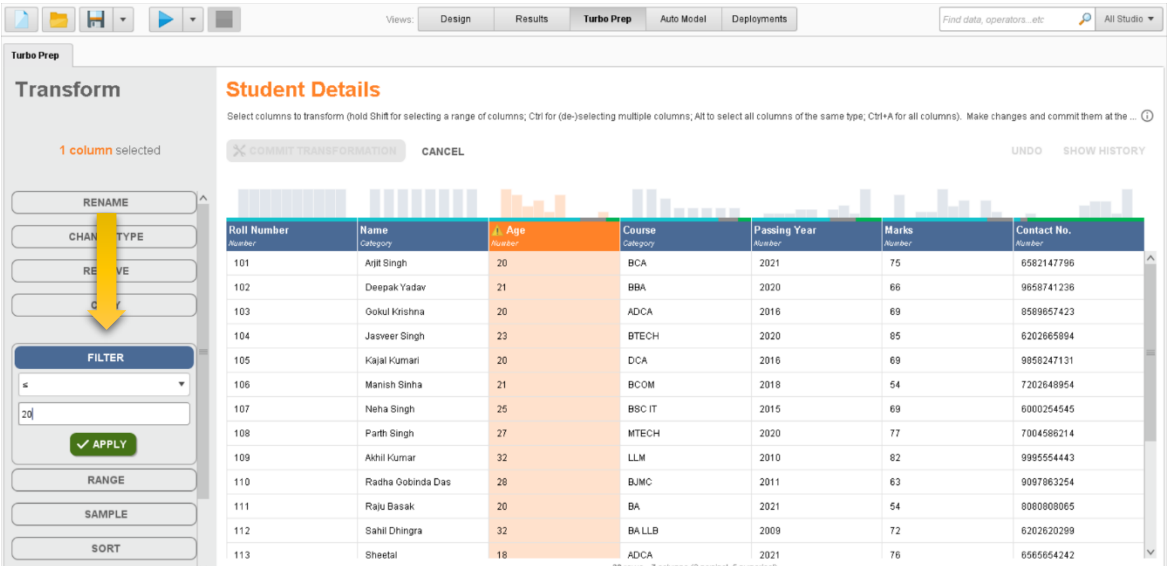**Step 2:** Now you can see the name attribute has been sorted in ascending order.



## b) Filter Numeric Data

**Step 1:** First we have to select numeric data attribute, here I selected **Age**. Now click on **Filter** and select filter operator like **=, >, <, ≠, ≤, ≥, is missing, is not missing.** Here I selected **≤20.**

| Course Code: **CAP 447** | Course Title: **Data Warehousing and Data Mining Lab** |
|---|---|
| Course Instructor: **Dr. Geeta Sharma** | |

| Student's Roll no: **RD2110B79** | Student's Reg. no: **12102801** |
|---|---|
| Name: **Atul Kumar** | |

| Question No. **02** | Page No. **08** | Total Pages. **17** |
|---|---|---|

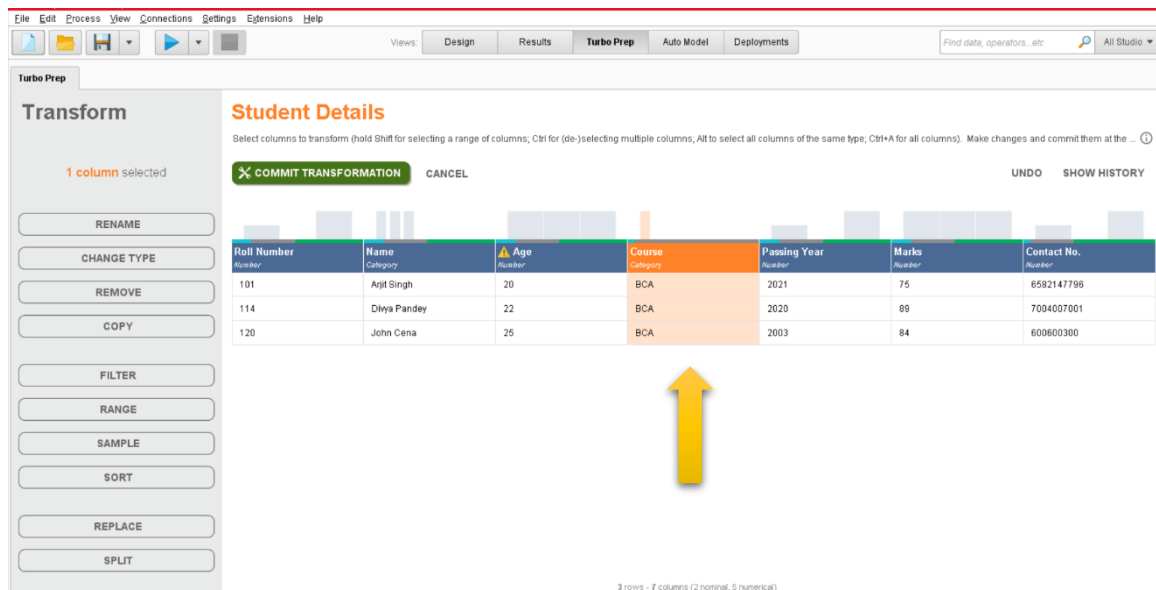**Step 2:** Here we can see our Age attribute is **Filtered** as per the given operation. **(≤20)**



## c) Filter String Data

**Step 1:** First we have to select string data attribute, here I selected **Course**. Now click on **Filter** and select filter operator like **equals, does not equal, is in, is not in, contains, does not contain, start with etc.** Here I selected **start with (BCA).**

**Step 2:** Here we can see our Course attribute is **Filtered** as per the given operation. (**start with BCA**)



### d) Remove Attribute

**Step 1:** First we have to select the attribute which we want to remove from the data set. Here I selected **Contact No.** attribute.

Now click on **Remove** and then click on **Apply.**

**Step 2:** Here we can see Contact No. attribute has been **Removed**.



**Q3.** **Perform following Cleanse operations on your data.**

   **a) Auto Cleanse**
   **b) Normalisation**
   **c) Discretization**

**Answer:** For performing all the following transformation in the Imported data set first we have to click on **Turbo Prep** and then click on **Load Data** and select the imported data set. (Imported data sets are in the local repository.) Now click on **CLEANSE**.

This interface will come up after clicking on **CLEANSE.**

## a) Auto Cleanse
### Step 1: Click on AUTO CLEANSING.



### Step 2: Select the attribute, here I selected Marks.

| Course Code: **CAP 447** | Course Title: **Data Warehousing and Data Mining Lab** |
|---|---|
| Course Instructor: **Dr. Geeta Sharma** | |
| Student's Roll no: **RD2110B79**      Name: **Atul Kumar** | Student's Reg. no: **12102801** |
| Question No. **03**     Page No. **13** | Total Pages. **17** |

**Step 3:** Rapid miner automatically select some attributes whose have very **low quality** and **remove** them at the last.



**Step 4:** If you want to change the type of selected attribute then Select change otherwise select keep original. Here I selected **Keep Original.**

**Step 5:** Now we have two options **Perform PCA/Perform normalization.** Here I selected **Perform normalization.**



**Step 6:** We can see a green tick sign means all the operations have been completed. Now click on **APPLY AUTO CLEANSING.**

**Step 7:** Now we can see the **Final Result** of auto cleansing.



### b) Normalization

**Step 1:** Select the attribute where we want to perform normalization. Here I selected **Marks.** Now **Define range.**
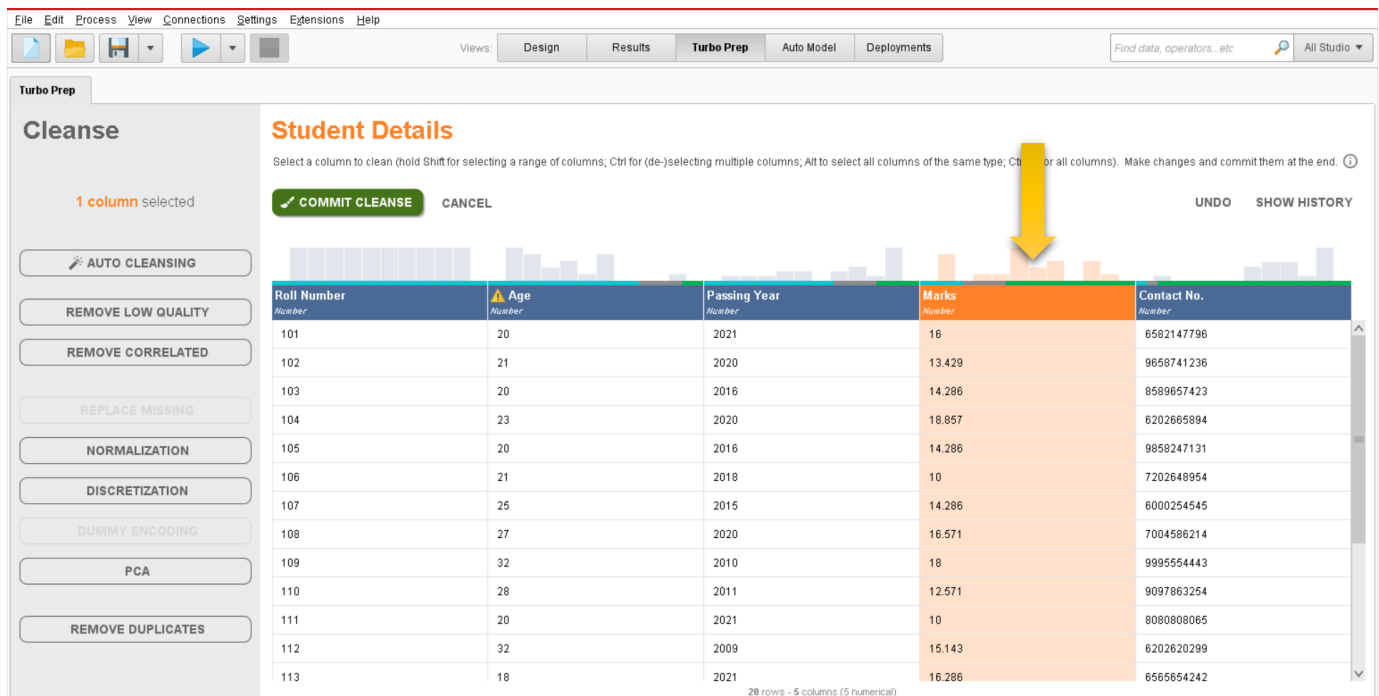
**Step 2:** Now we can see marks are **Normalized** between 10-20.



## c) Discretization

**Step 1:** Here I selected **Roll Number** for discretization.
Now set the number of bins, here I write **5** and then click on **Apply.**

**Step 2:** Here we can see the data of Roll Number attributes has been changed to **5 equal bins (range1 to range5).**



# Thank You