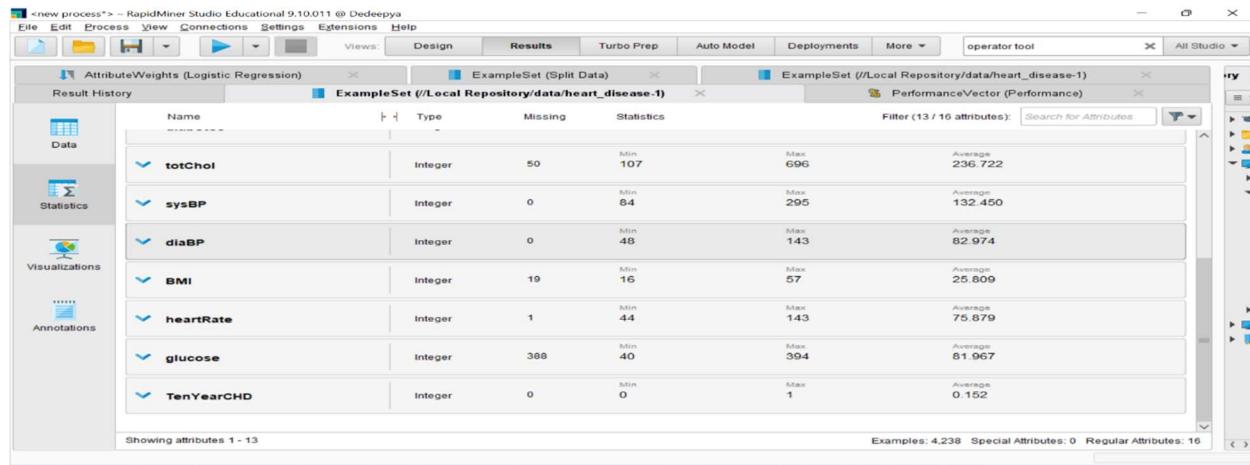


1. Import the dataset "heart\_disease.csv" into Rapid Miner, while importing

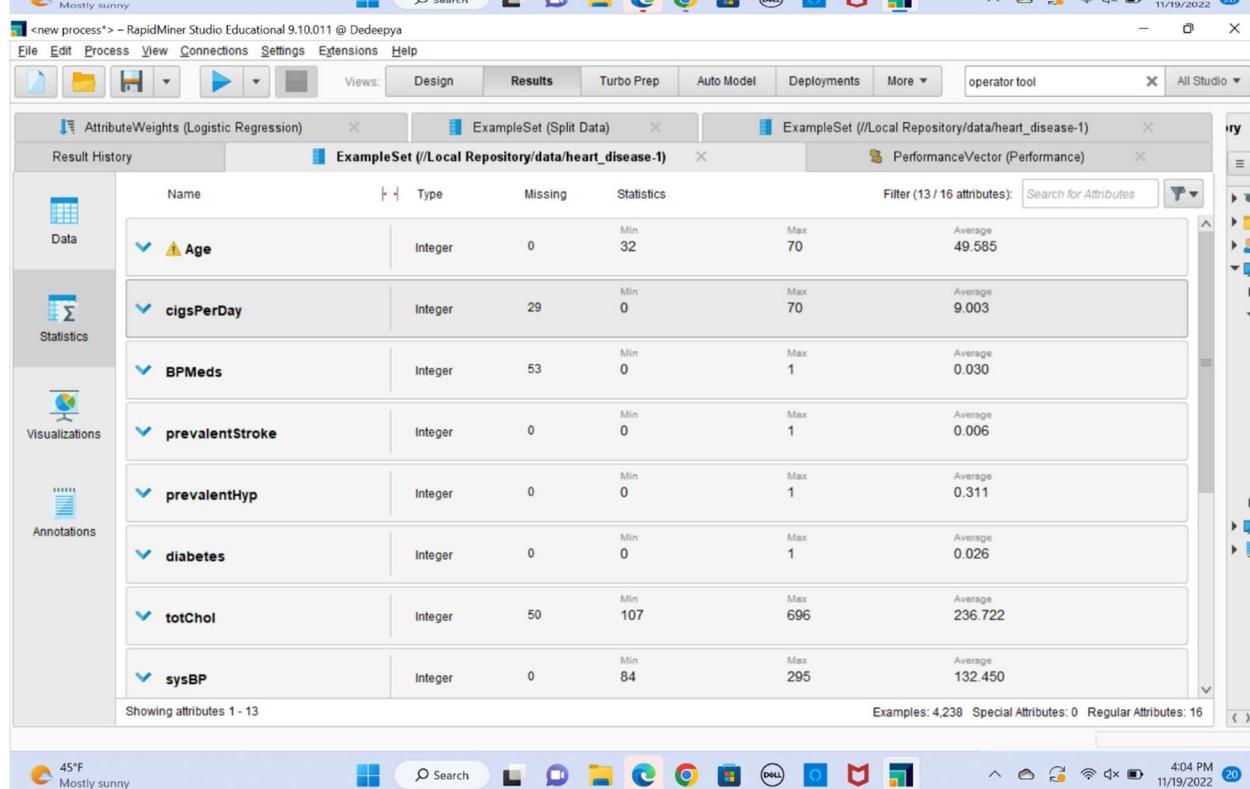
A. Convert the below mentioned columns to integers.

Age	cigsPerDay	BPMed	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP	BMI	heartRate	glucose	TenYearCHD
-----	------------	-------	-----------------	--------------	----------	---------	-------	-------	-----	-----------	---------	------------

B. Click on replace errors with missing values check box to avoid parsing error in cigsPerDay column.

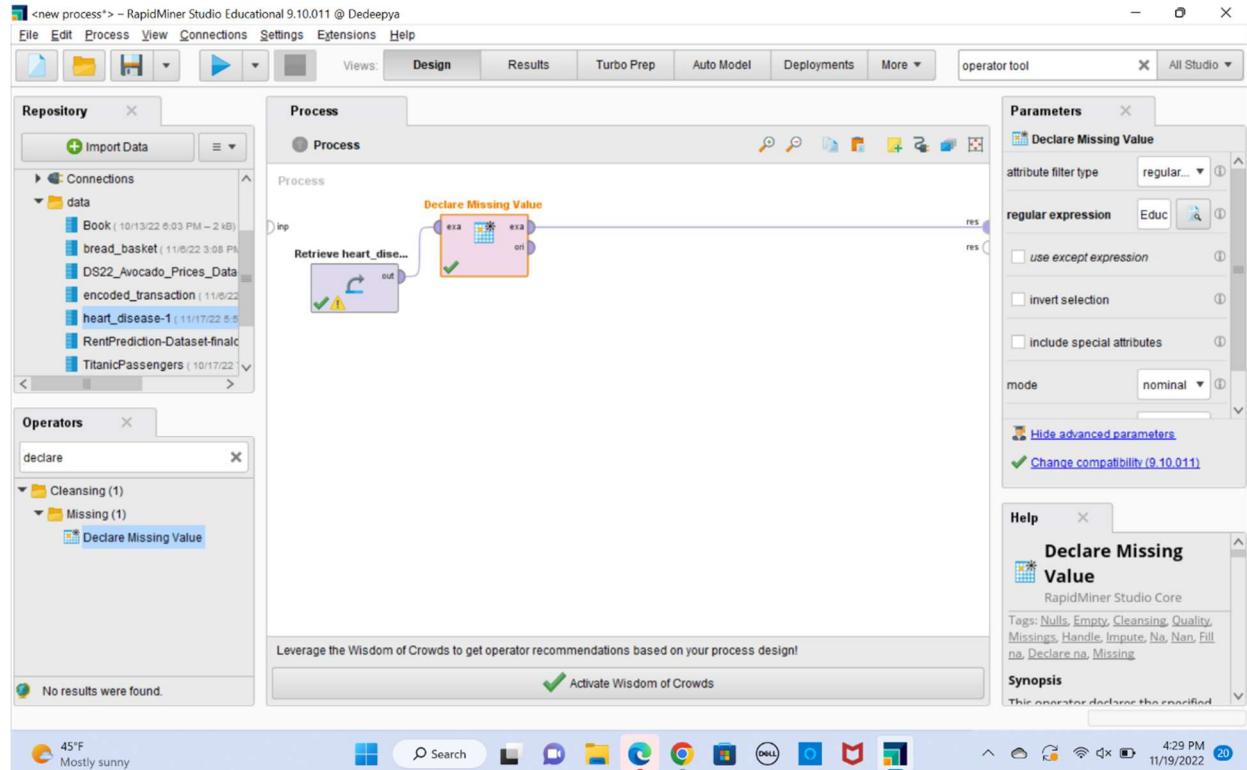


Name	Type	Missing	Statistics
totChol	Integer	50	Min: 107 Max: 696 Average: 236.722
sysBP	Integer	0	Min: 84 Max: 295 Average: 132.450
diaBP	Integer	0	Min: 48 Max: 143 Average: 82.974
BMI	Integer	19	Min: 16 Max: 57 Average: 25.809
heartRate	Integer	1	Min: 44 Max: 143 Average: 75.879
glucose	Integer	388	Min: 40 Max: 394 Average: 81.967
TenYearCHD	Integer	0	Min: 0 Max: 1 Average: 0.152



Name	Type	Missing	Statistics
Age	Integer	0	Min: 32 Max: 70 Average: 49.585
cigsPerDay	Integer	29	Min: 0 Max: 70 Average: 9.003
BPMed	Integer	53	Min: 0 Max: 1 Average: 0.030
prevalentStroke	Integer	0	Min: 0 Max: 1 Average: 0.006
prevalentHyp	Integer	0	Min: 0 Max: 1 Average: 0.311
diabetes	Integer	0	Min: 0 Max: 1 Average: 0.026
totChol	Integer	50	Min: 107 Max: 696 Average: 236.722
sysBP	Integer	0	Min: 84 Max: 295 Average: 132.450

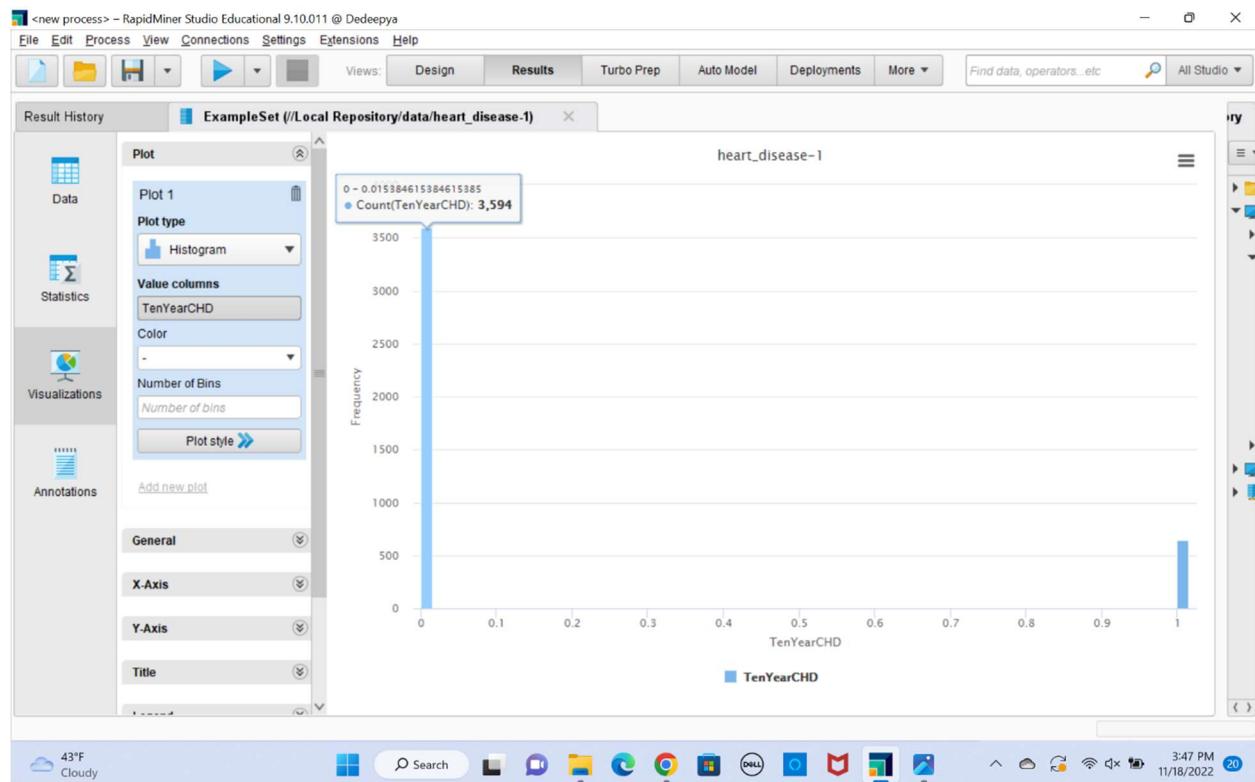
C. Using the “Declare missing values operator” identify “N/A” string in Education level column as Missing values (Use nominal mode and enter “N/A” in expression of declare missing values operator).



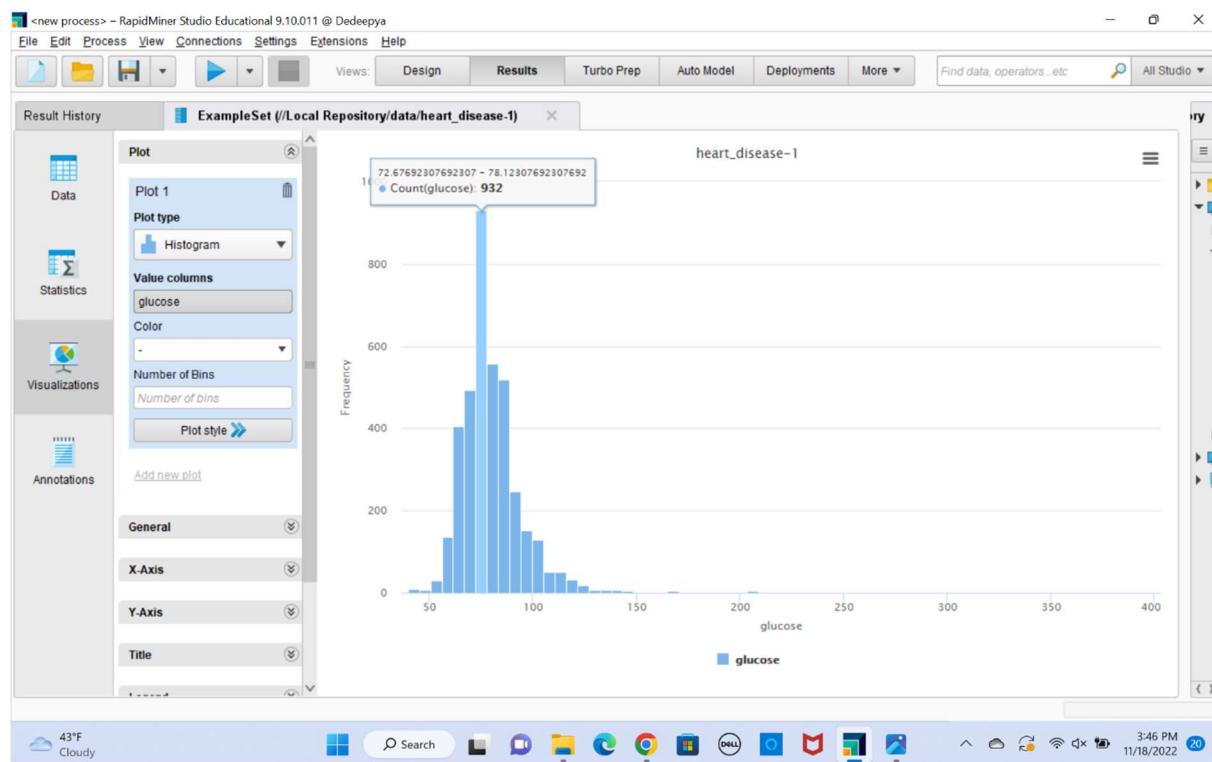
1

- a) Include a histogram of each variable and describe each variable (Use visualization tab after running the dataset in canvas). Need separate histograms for all variables and explain what you see from the histogram plot.
- b) Is there any row that contains missing values on all columns?  
Which column contains missing values or null values?  
How many missing values are there in each column? (Use statistics tab in results view)

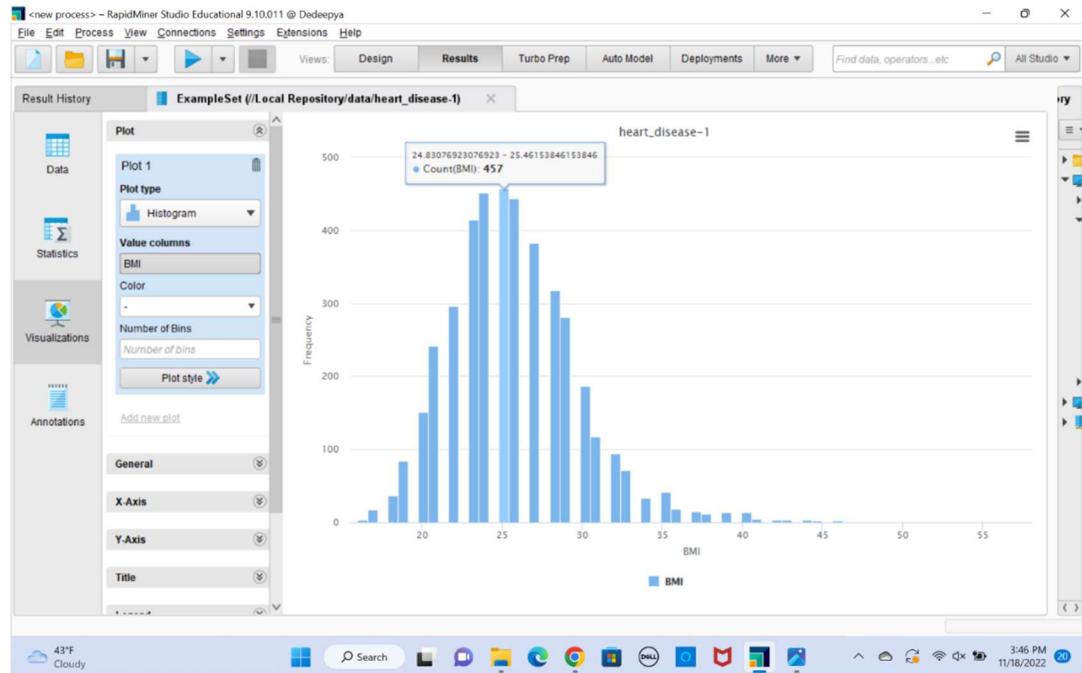
a) TenYearCHD- It has only 2 values 0 and 1 and there are 3594 zeros and 644 ones



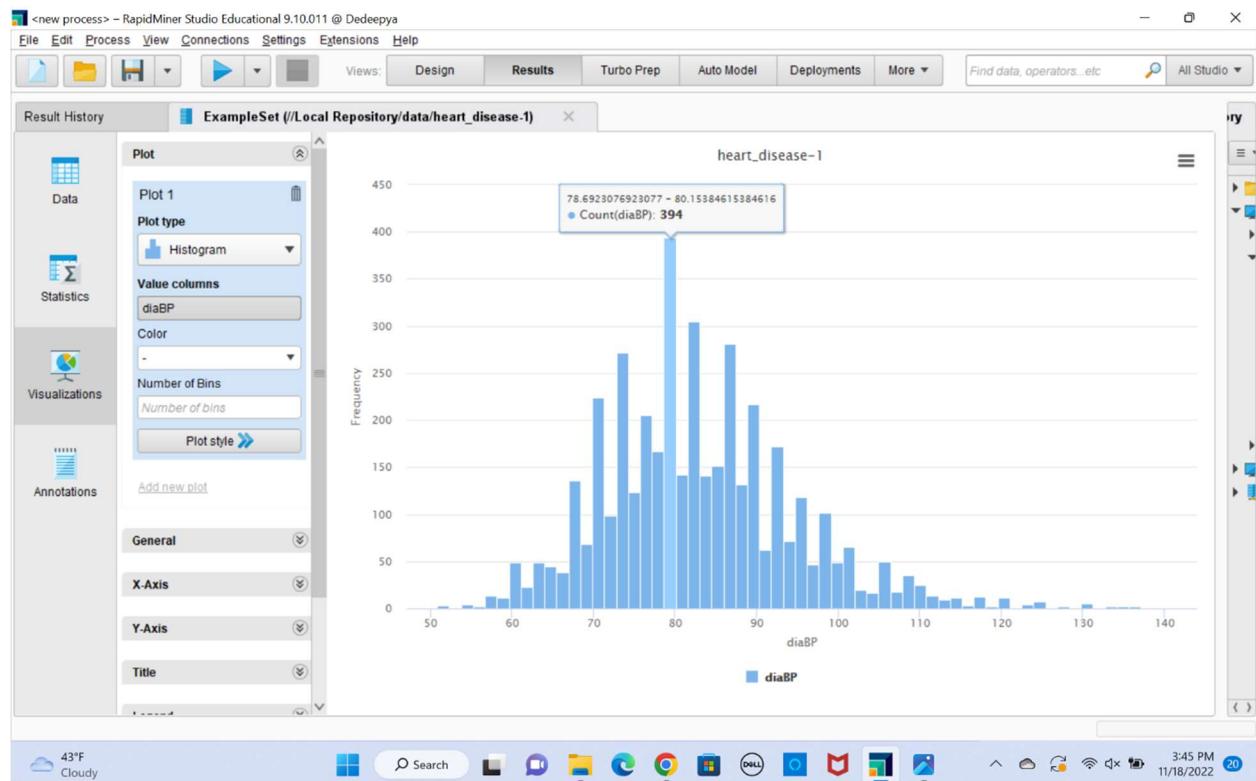
Glucose values ranges from 40 to 208.83, the values between 72.67-78.12 have high frequency of 932.



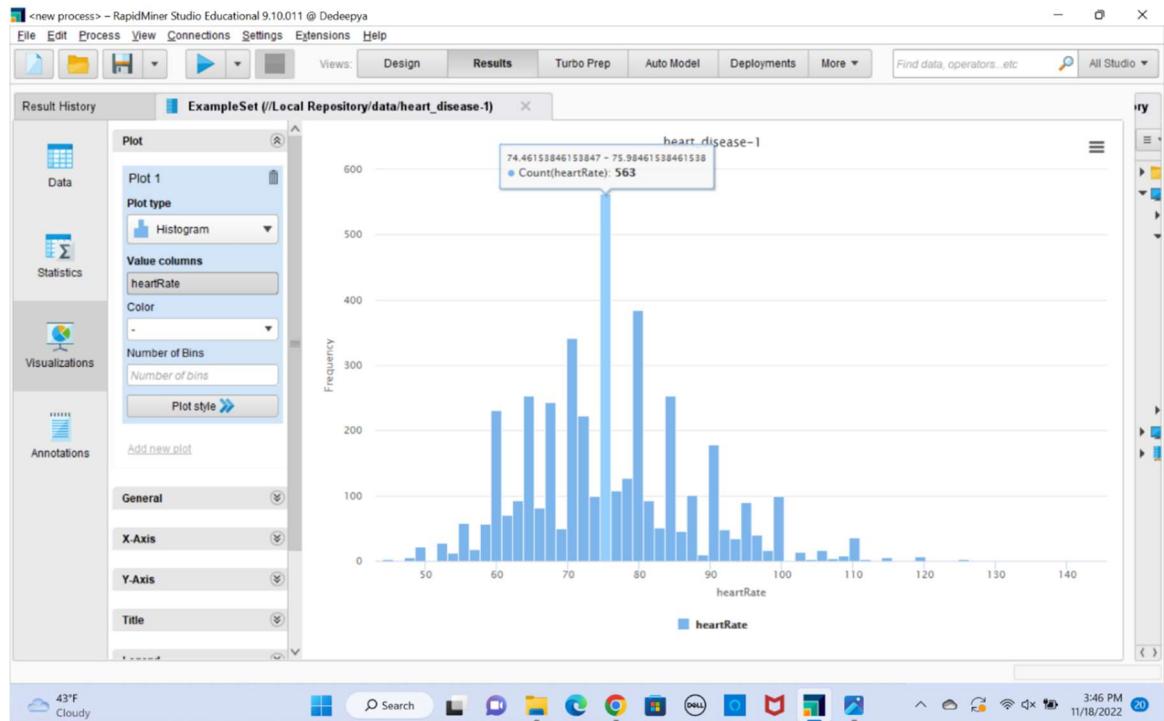
BMI Ranges from 16-46.27 and the values of 24.83-25.46 has high frequency of 457



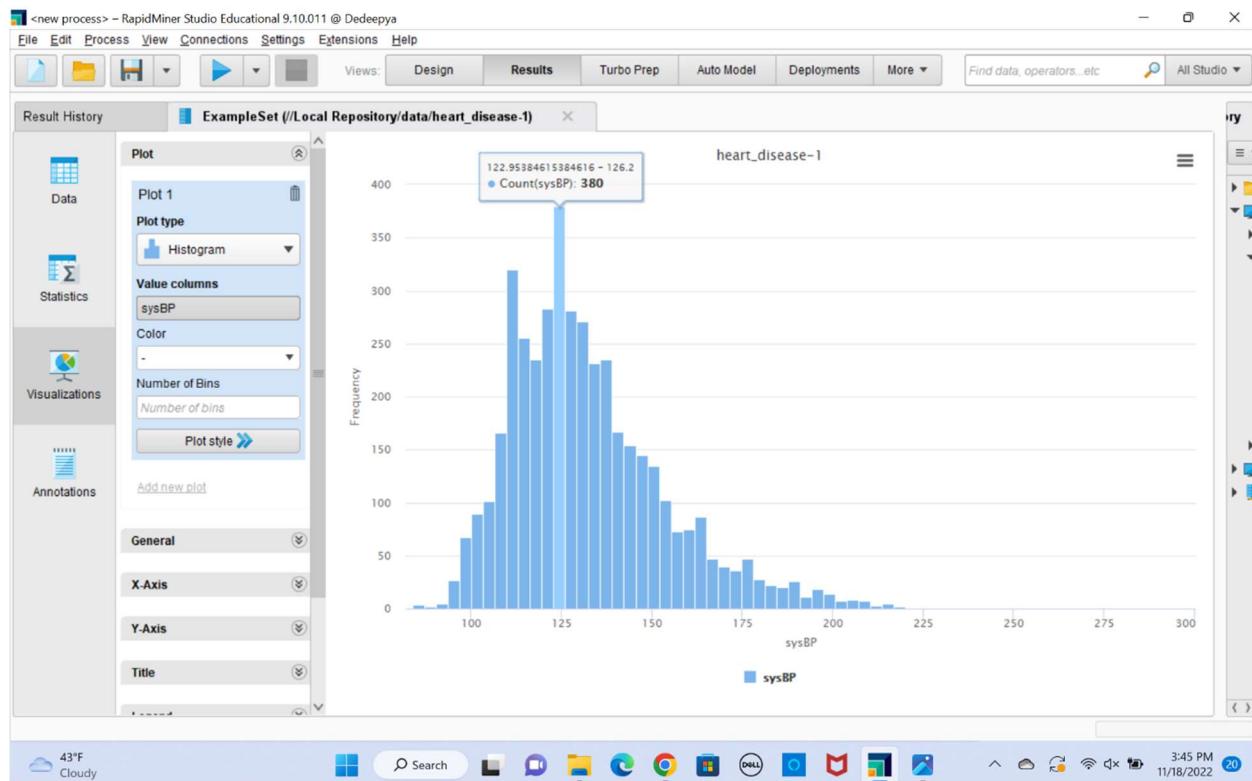
DiaBp- range 50.92-137.15 and, values between 78.69-80.15 has high frequency of 394



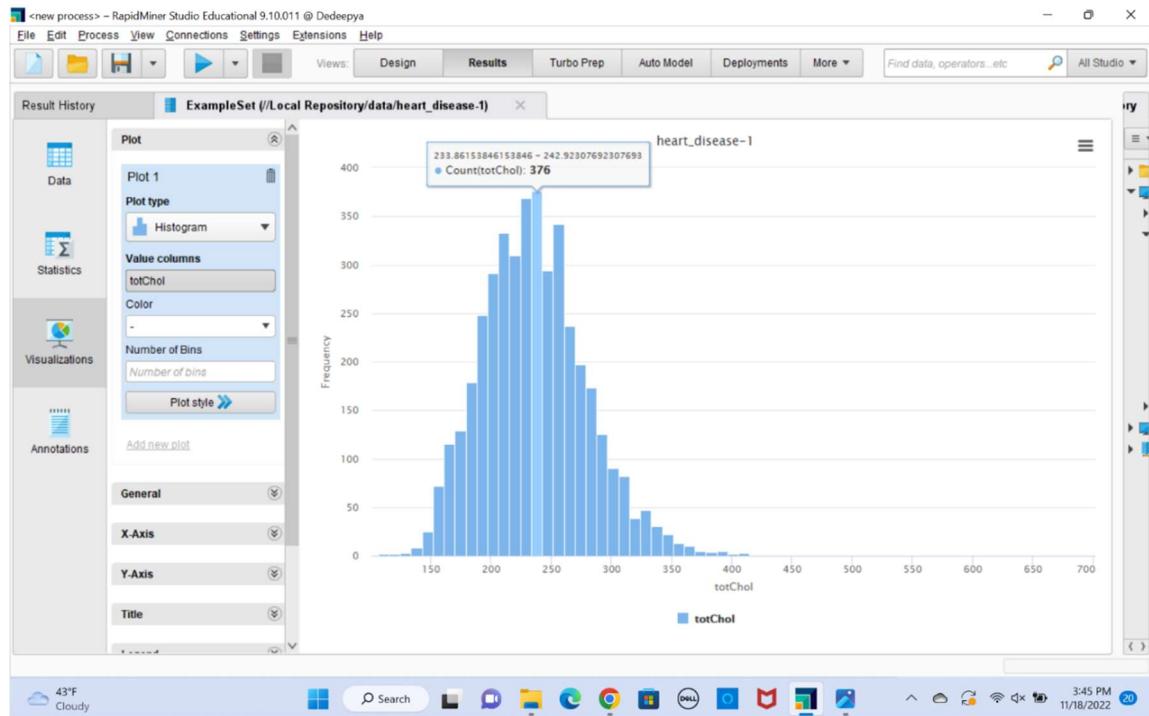
Range of heartrate is 44-126.24, values within range 74.46-75.98 has 563 frequencies



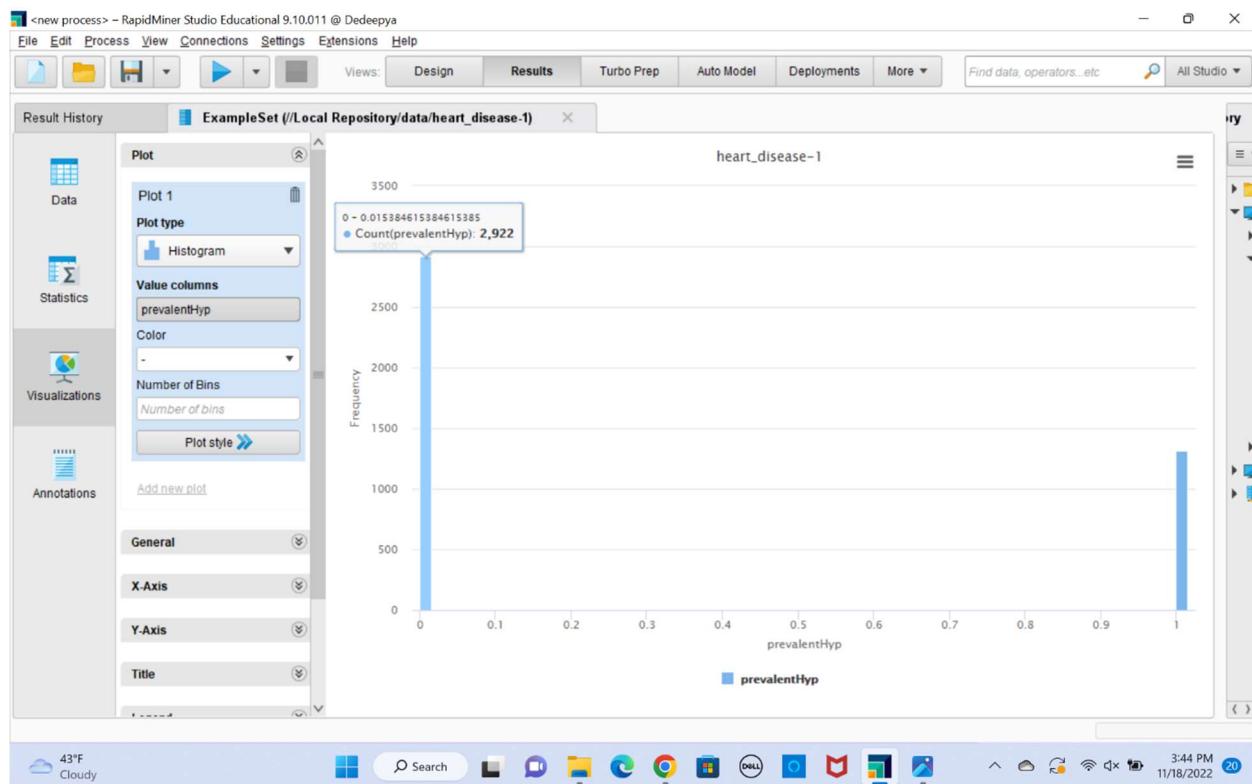
SysBP Range 84-220.33 and values within range 122.95-126.2 has frequency of 380



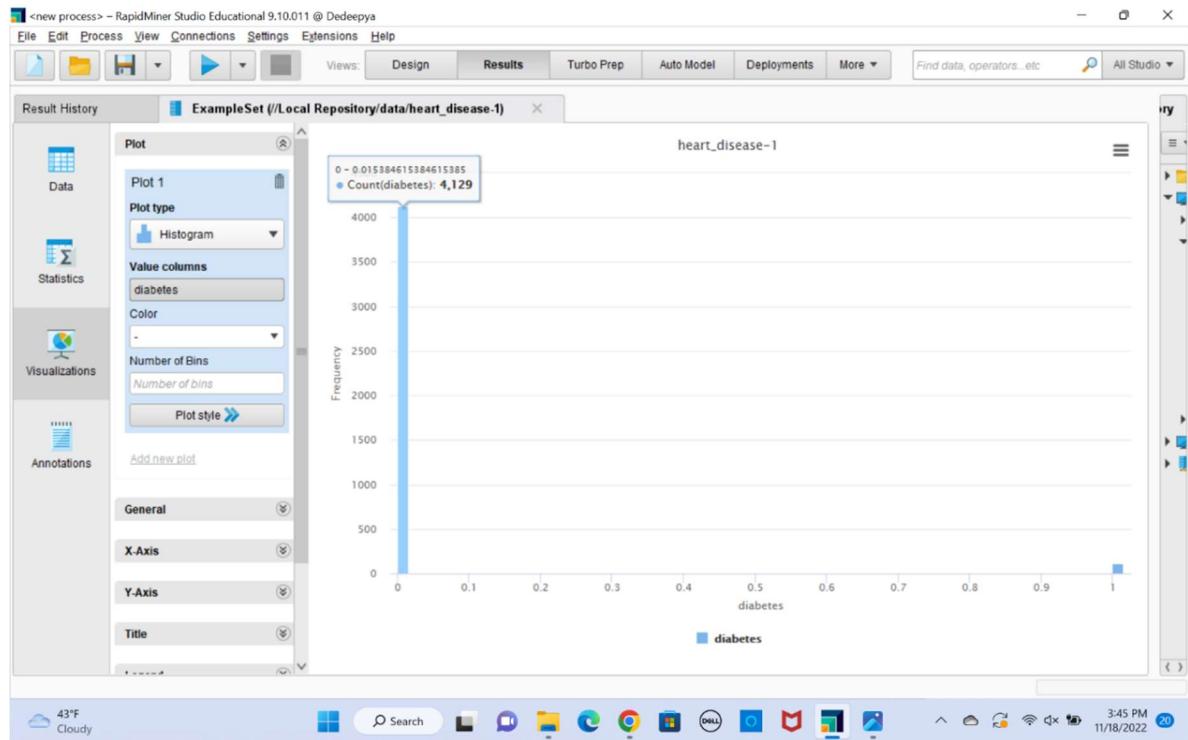
Totchol Range is 107-415, and high frequency of 376 for range 233.86-242.923



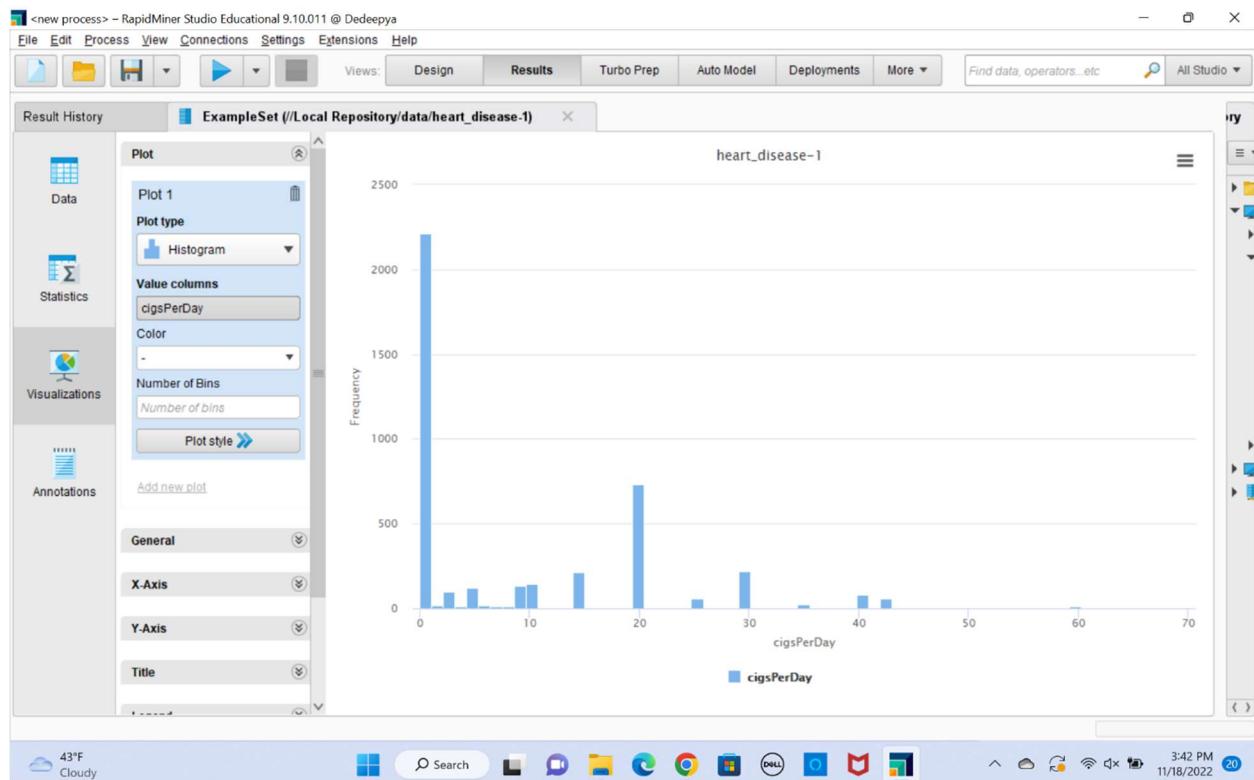
Prevalenthyp has Only 2922 zeros and 1316 ones



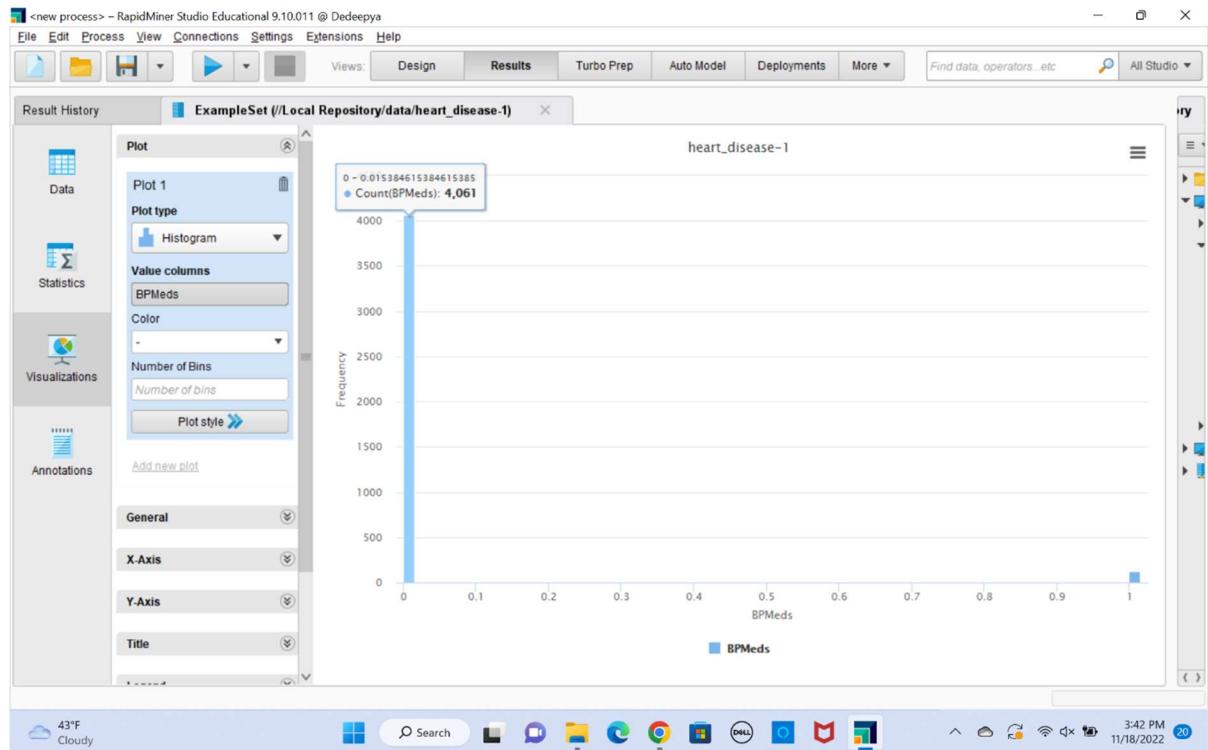
Diabetes has 4129 zeros, 109 ones



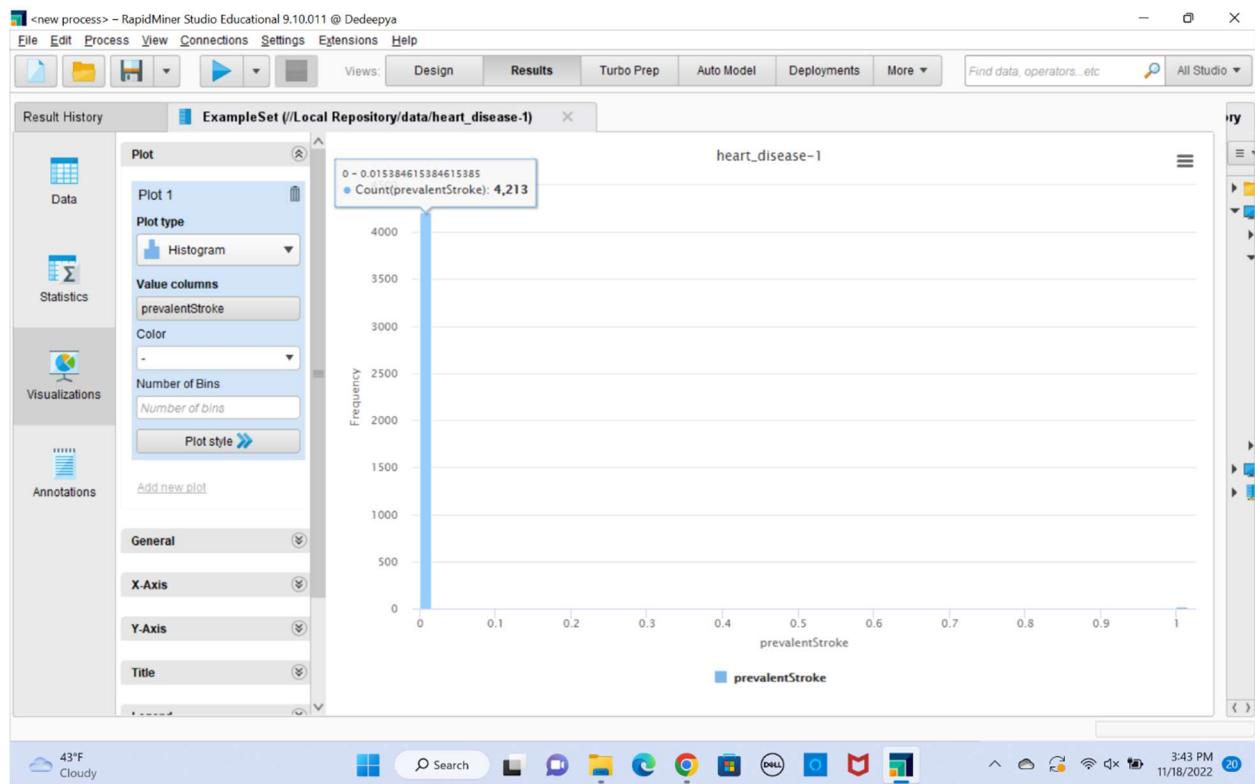
Cigsperday has high frequency of 2211 for range 0-1.076 and values ranges from 0-60.30



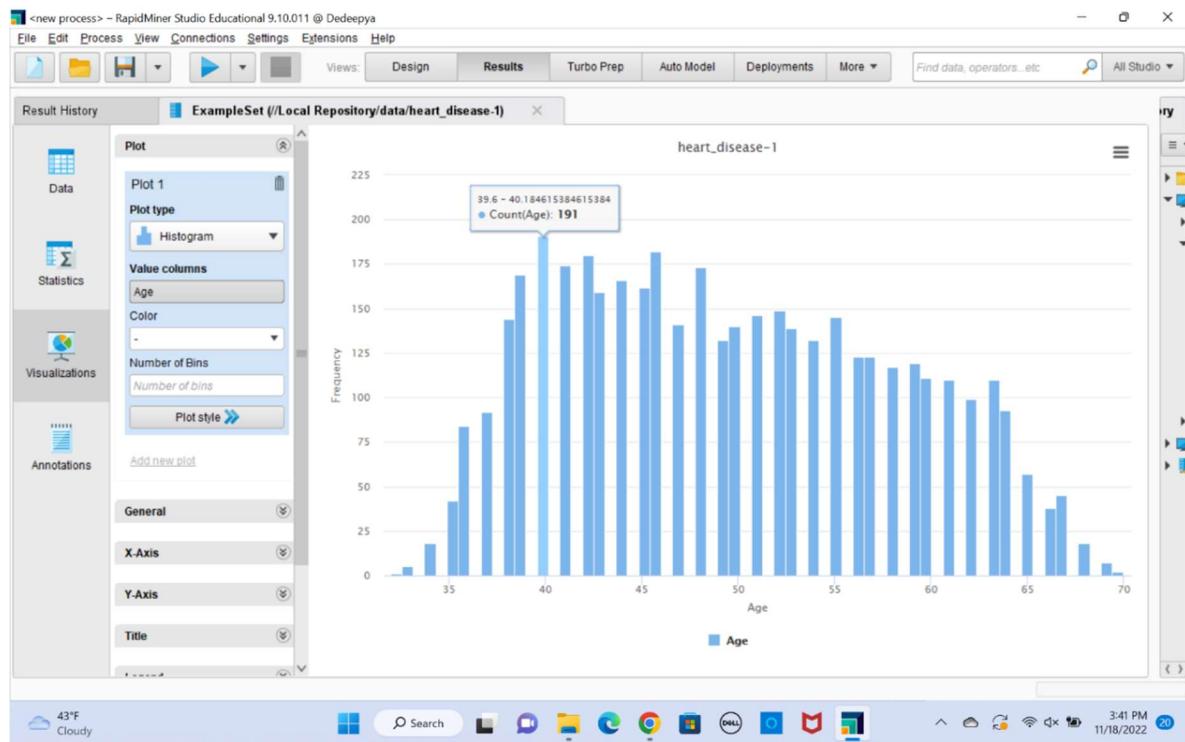
BPmeds has Freq-4061 for 0, and for 1, we have freq-124.



Prevalent stroke has 4213 frequencies for 0 and frequency of 25 for 1

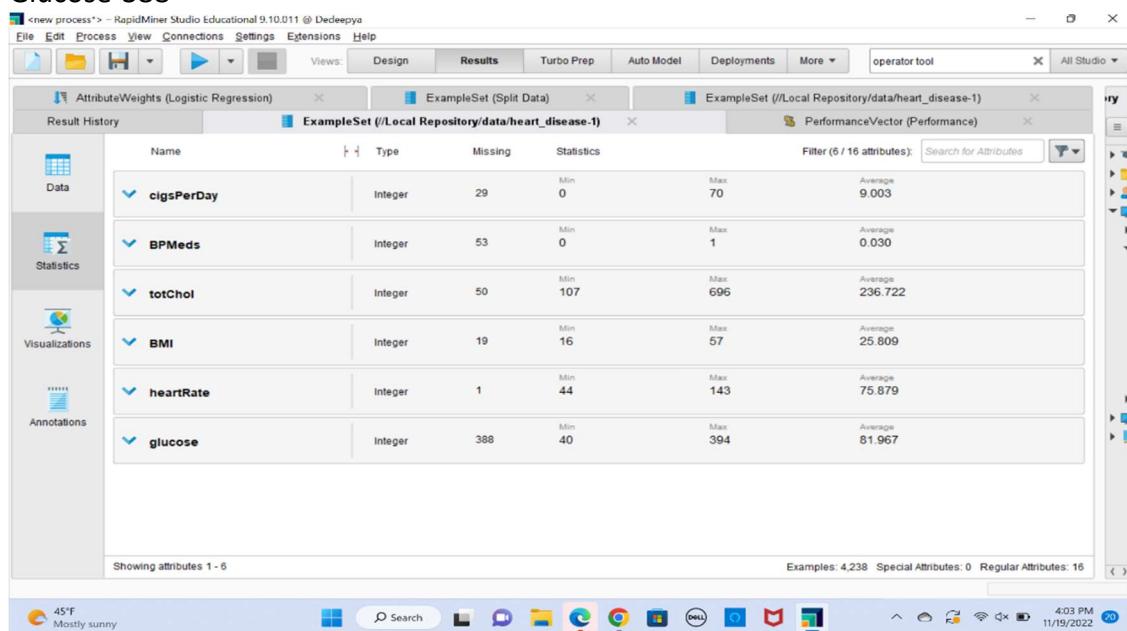


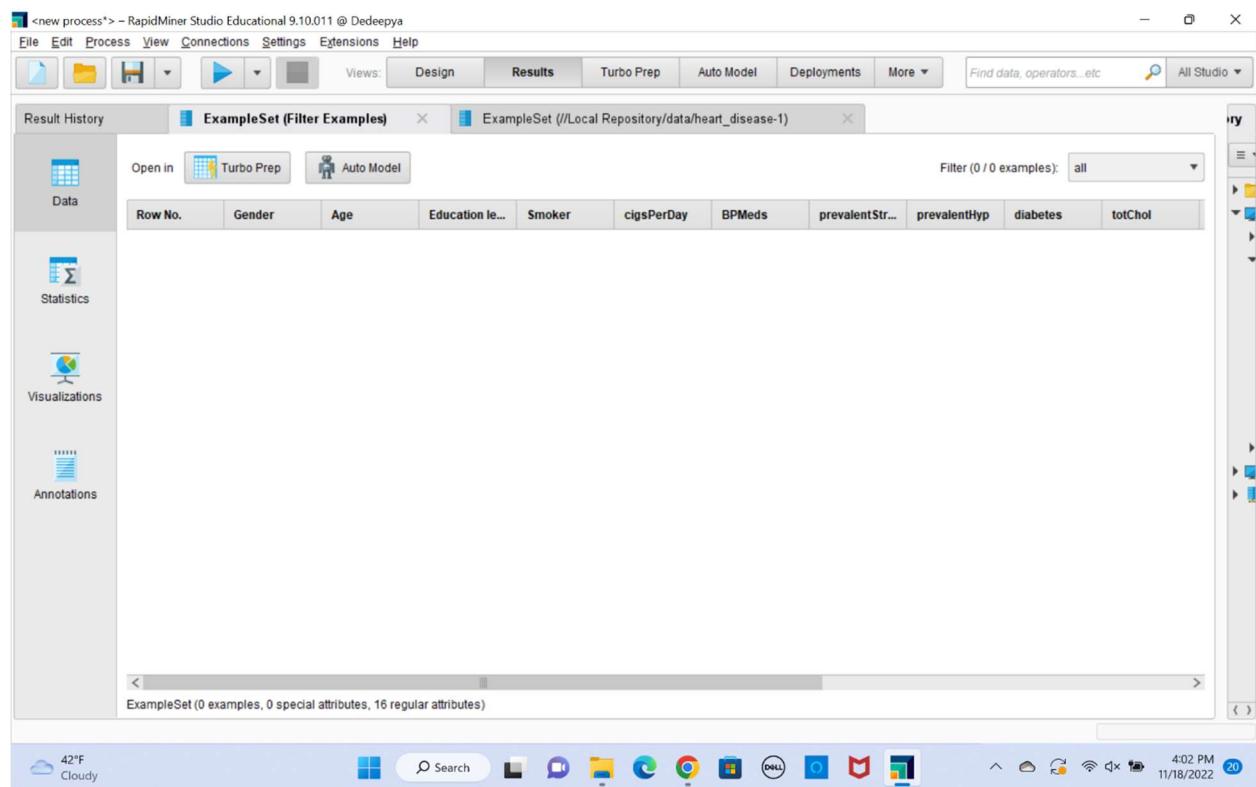
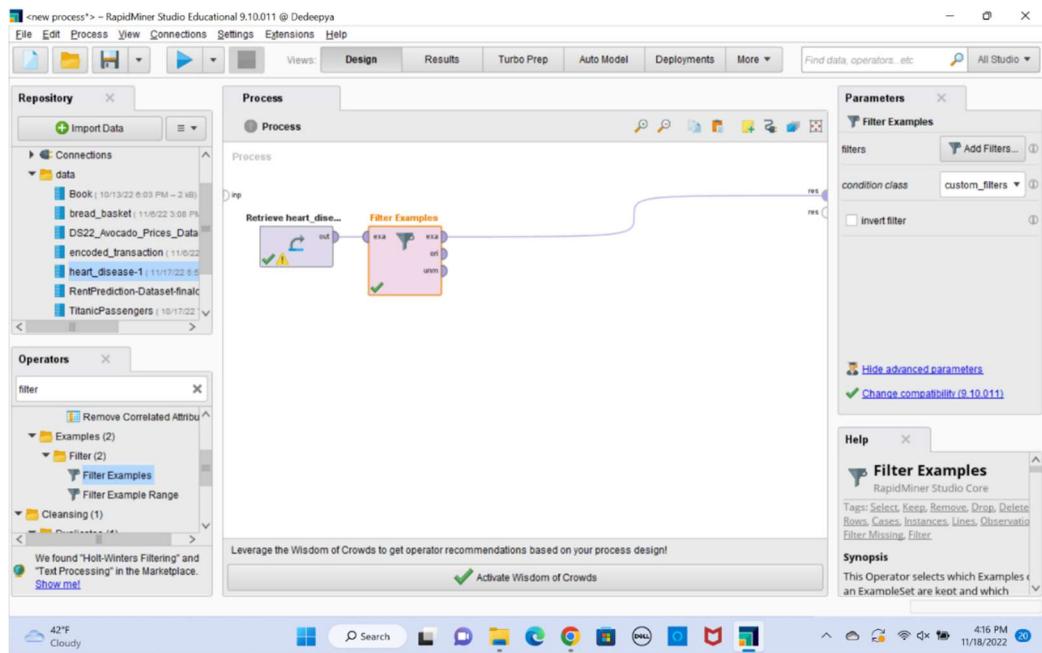
Age has range 32-70 and high frequency of 191 for values between 39.6-40.18



b) There is no record with missing values in all columns.

Cigsperday-29  
 BPMeds-53  
 Totchol-50  
 BMI-19  
 Heartrate-1  
 Glucose-388





- 
2. For each record that has a missing value in “cigsPerDay” and is identified as a smoker, replaces the missing values with the average of “cigsPerDay” column (use filter examples for filtering the values who are identified as smokers)

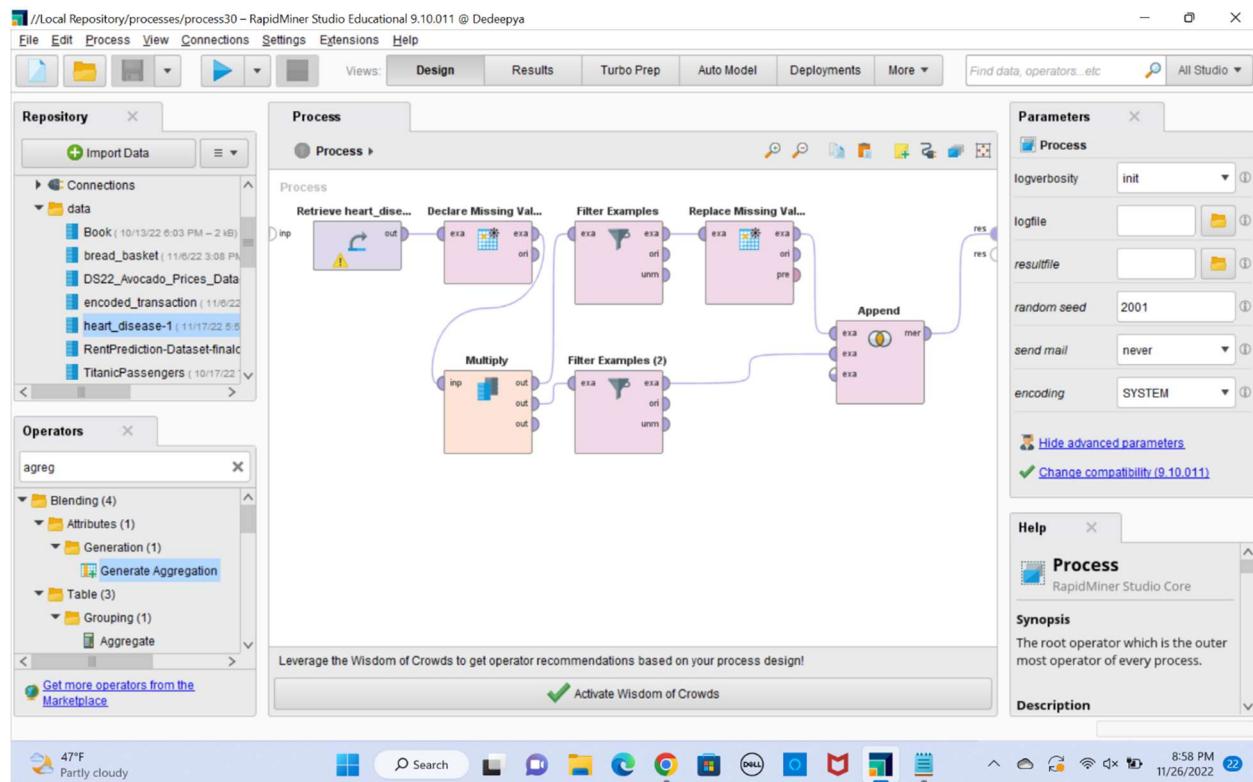
Use replaces missing values operator and). Answer the following questions

Use append operator to combine the non-smoker data to missing values replaced data.)

Include a screenshot of Result set that shows number of data on top right corner of result view using the all filter.

Answer the following questions

- A. Out of all the records that contain missing values in “cigsPerDay”, how many of them are male, and how many of them are female?
  - B. Find the age range of the Male smokers that have missing values in “cigsPerDay”. Do the same thing for Female smokers with missing values in “cigsPerDay”.
  - C. What is the average “cigsPerDay” among the Male smokers, whose age is between 30 and 39 years old; 40 and 49 years old; 50 and 59 years old? (Round to the nearest integer)
  - D. What is the average “cigsPerDay” among the Female smokers, whose age is between 30 and 39 years old; 40 and 49 years old; 50 and 59 years old; 60 and 69 years old? (Round to the nearest integer)
-



File Edit Process View Connections Settings Extensions Help

Views: Design Results Turbo Prep Auto Model Deployments More Find data, operators...etc All Studio

**Result History**

**ExampleSet (Append)**

ExampleSet (/Local Repository/data/heart\_disease-1)

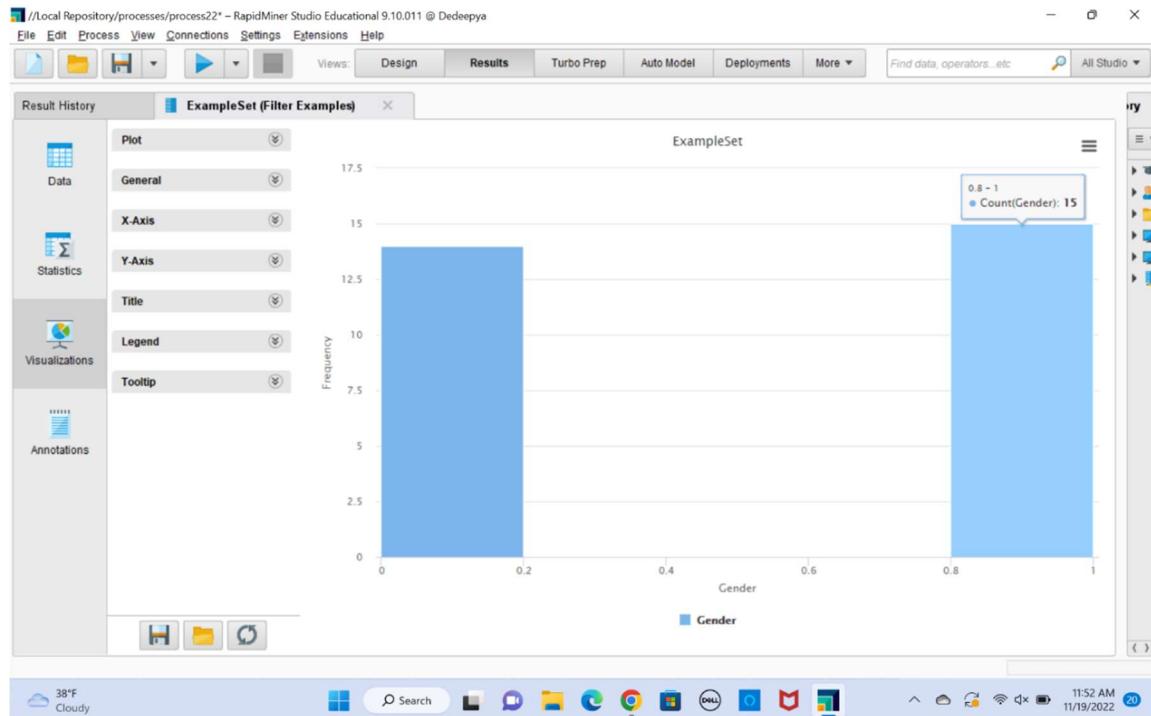
Data

Row No.	cigsPerDay	Gender	Age	Education le...	Smoker	BPMeds	prevalentStr...	prevalentHyp	diabetes	totChol
1	20	Male	48	middle school	Yes	0	0	0	0	245
2	30	Female	61	bachelor	Yes	0	0	1	0	225
3	23	Female	46	bachelor	Yes	0	0	0	0	285
4	20	Female	45	high school	Yes	0	0	0	0	313
5	30	Male	43	middle school	Yes	0	0	1	0	225
6	15	Male	46	middle school	Yes	0	0	1	0	294
7	9	Female	39	high school	Yes	0	0	0	0	226
8	20	Female	38	high school	Yes	0	0	1	0	221
9	10	Male	48	bachelor	Yes	0	0	1	0	232
10	20	Female	46	high school	Yes	0	0	0	0	291
11	5	Female	38	high school	Yes	0	0	0	0	195
12	30	Female	42	high school	Yes	0	0	0	0	190
13	20	Female	52	bachelor	Yes	0	0	0	0	215
14	30	Male	44	high school	Yes	0	0	1	0	270
15	20	Male	47	higher educat...	Yes	0	0	0	0	294

ExampleSet (4,238 examples, 0 special attributes, 16 regular attributes)

47°F Partly cloudy 8:58 PM 11/26/2022

2A Male-14, female-15

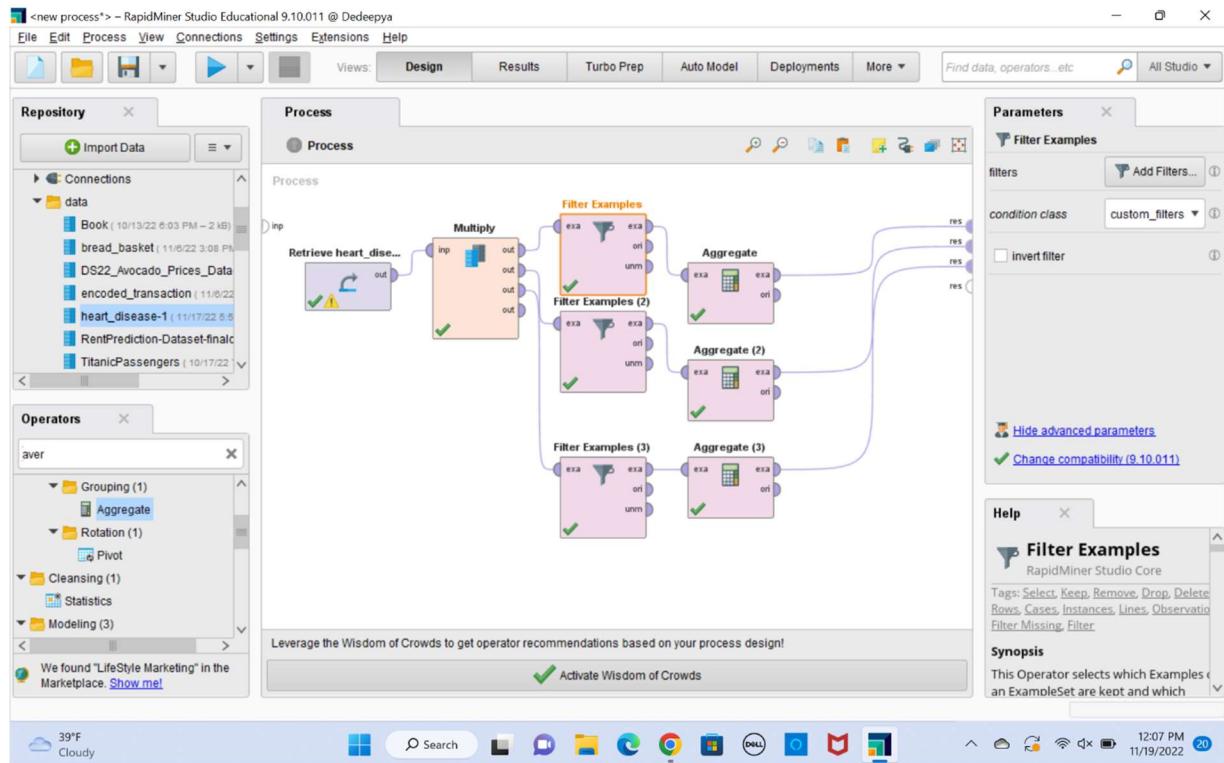


2b

Female age range 37-58, Male Age range - 39-61

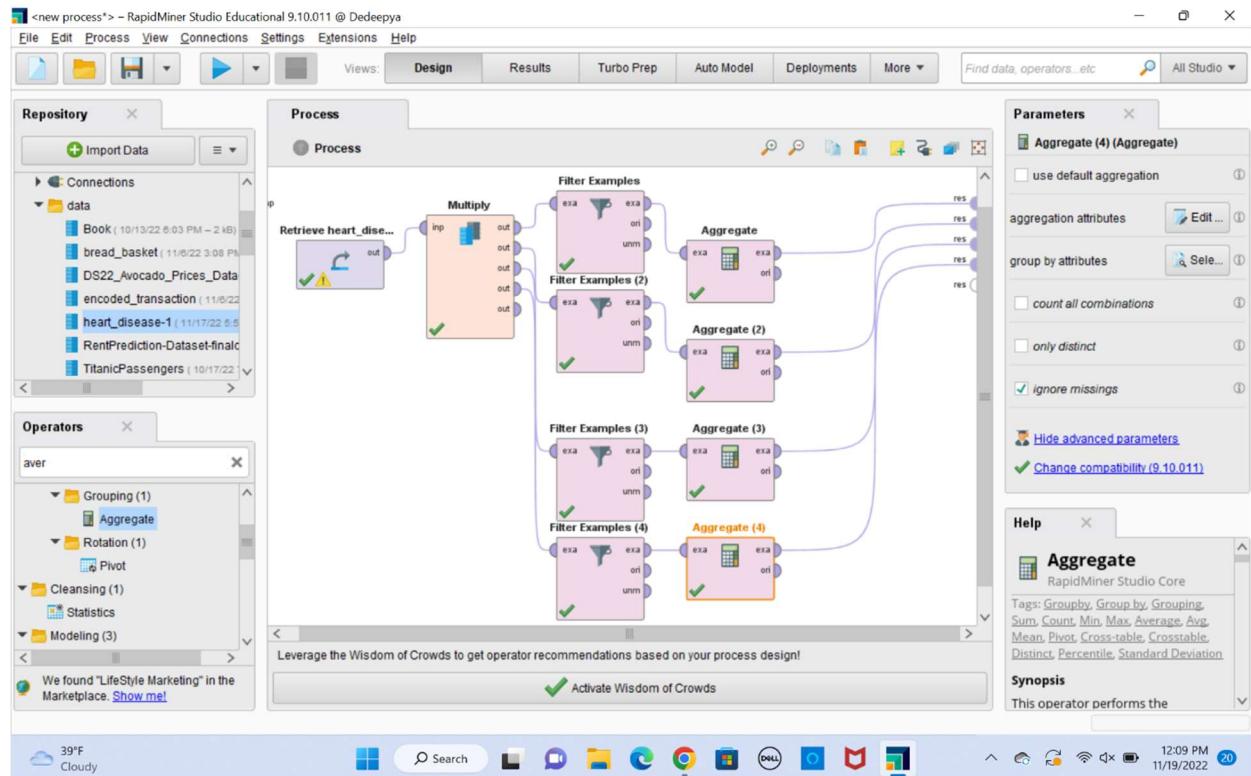
2c

Male Smoker age range	Average <u>Cigerperday</u>
30-39	22
40-49	23
50-59	22



2d.

Female Smoker age range	Average Cigerperday
30-39	15
40-49	14
50-59	14
60-69	14



3. Fill in the missing values in the other columns and show proofs that you have completed these steps by including the screenshots of before and after replacing of missing values count in statistics tab: (For all the below questions use the replace missing values operator just like we used in question 2).
- A. For each record that has a missing value in “glucose”, replace the missing values with the average “glucose”.
- B. For each record that has a missing value in “totChol”, replace the missing values with the average “totChol”.
- C. Replace the missing values in “BPMeds” with 0.
- D. Replace the missing values in “BMI” with the average BMI.
- E. Replace the missing values in “Education level” with the value “Unknown”.
- F. Replace the missing values in “heartRate” with the average heart rate value.

### 3. Before

Name	Type	Missing	Statistics
cigsPerDay	Integer	29	Min: 0 Max: 70 Average: 9.003
BPMeds	Integer	53	Min: 0 Max: 1 Average: 0.030
totChol	Integer	50	Min: 107 Max: 696 Average: 236.722
BMI	Integer	19	Min: 16 Max: 57 Average: 25.809
heartRate	Integer	1	Min: 44 Max: 143 Average: 75.879
glucose	Integer	388	Min: 40 Max: 394 Average: 81.967

## After

Result History    ExampleSet (Replace Missing Values)    ExampleSet (//Local Repository/data/heart\_disease-1)

Name	Type	Missing	Statistics	Filter (16 / 16 attributes):	Search for Attributes
diabetes	Integer	0	0	1	0.026
totChol	Integer	0	Min 107 Max 696	236.725	Average
sysBP	Integer	0	Min 84 Max 295	132.450	Average
diaBP	Integer	0	Min 48 Max 143	82.974	Average
BMI	Integer	0	Min 16 Max 57	25.810	Average
heartRate	Integer	0	Min 44 Max 143	75.879	Average
glucose	Integer	0	Min 40 Max 394	81.970	Average
TenYearCHD	Integer	0	Min 0 Max 1	0.152	Average

Showing attributes 1 - 16    Examples: 4,238    Special Attributes: 0    Regular Attributes: 16

Cloudy 39°F    Search    Windows Icons    12:18 PM 11/19/2022

<new process\*> - RapidMiner Studio Educational 9.10.011 @ Dedeepya

File Edit Process View Connections Settings Extensions Help

Views: Design Results Turbo Prep Auto Model Deployments More Find data, operators...etc All Studio

Result History ExampleSet (Replace Missing Values) ExampleSet (/Local Repository/data/heart\_disease-1)

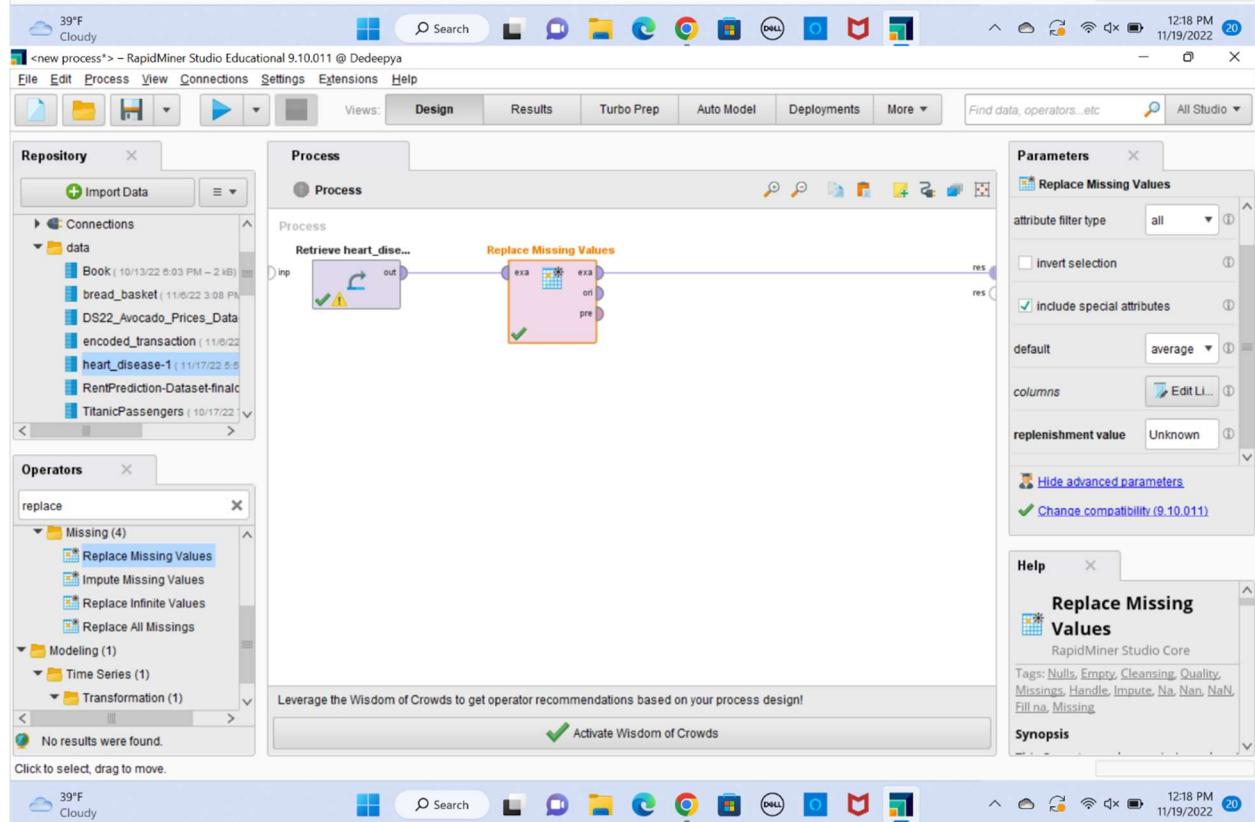
Filter (16 / 16 attributes): Search for Attributes

Data Statistics Visualizations Annotations

Attributes:

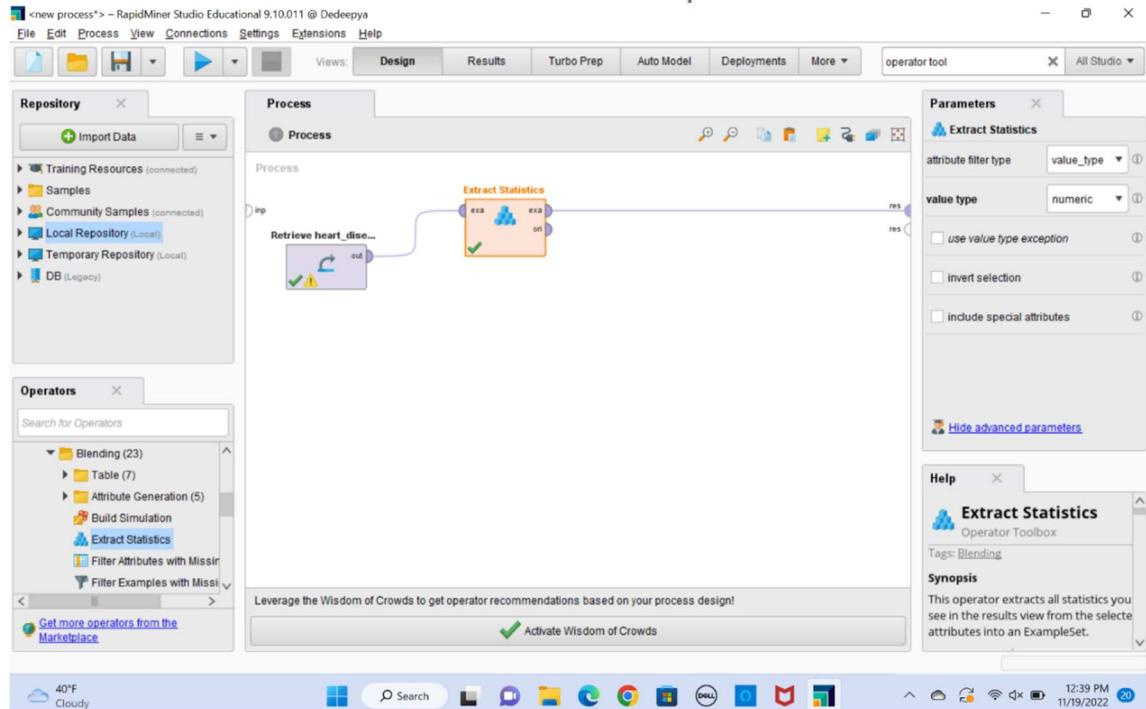
Name	Type	Missing	Statistics	Filter (16 / 16 attributes):
Gender	Polynomial	0	Least Male (1819) Most Female (2419) Values Female (2419), Male (1819)	Search for Attributes
Age	Integer	0	Min 32 Max 70 Average 49.585	
Education level	Polynomial	0	Least Unknown (0) Most middle school (1720) Values middle school (1720), high school (1253)	
Smoker	Polynomial	0	Least Yes (2094) Most No (2144) Values No (2144), Yes (2094)	
cigsPerDay	Integer	0	Min 0 Max 70 Average 9.003	
BPMeds	Integer	0	Min 0 Max 1 Average 0.029	
prevalentStroke	Integer	0	Min 0 Max 1 Average 0.006	
prevalentHyp	Integer	0	Min 0 Max 1 Average 0.311	

Showing attributes 1 - 16 Examples: 4,238 Special Attributes: 0 Regular Attributes: 16



G. What is the mean, median, min, max, and standard deviation of all numerical variables?

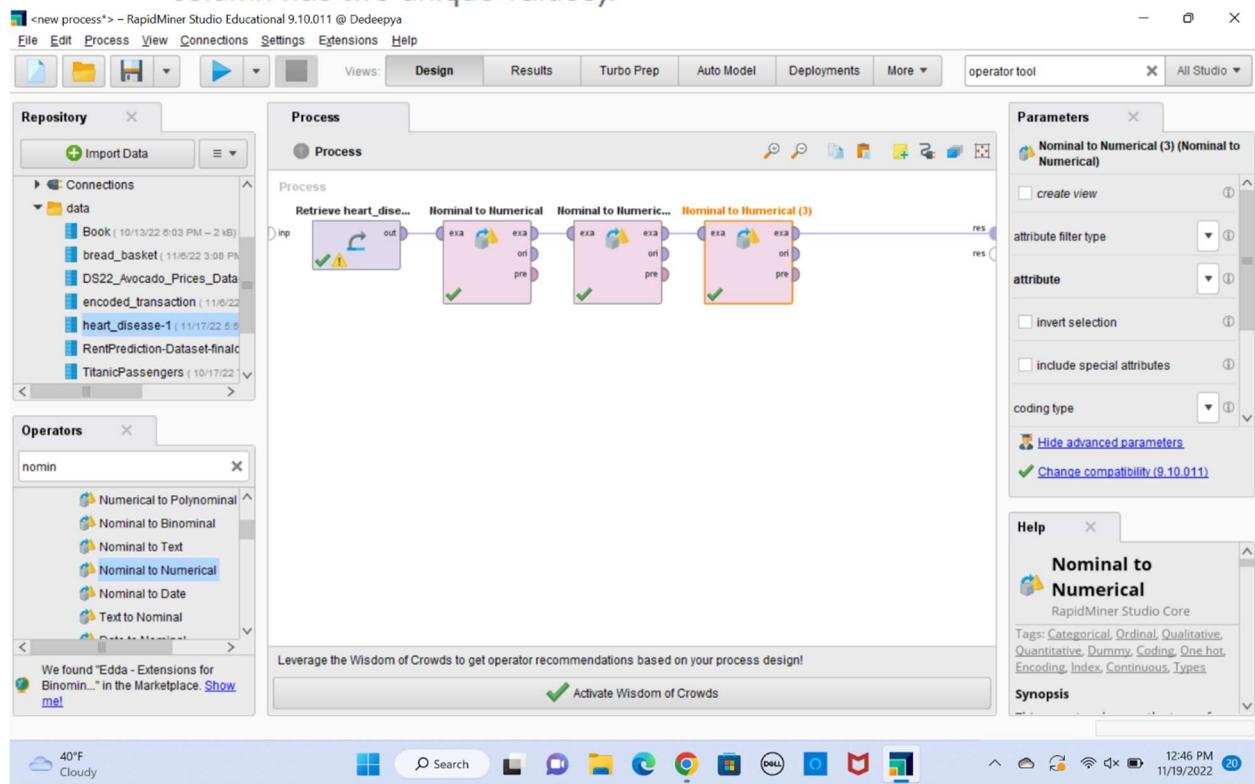
Install Operator toolbox extension by click on Extensions tab on top toolbar in rapid miner. This extension contains operator name “Extract Statistics” use that to answer the above question.



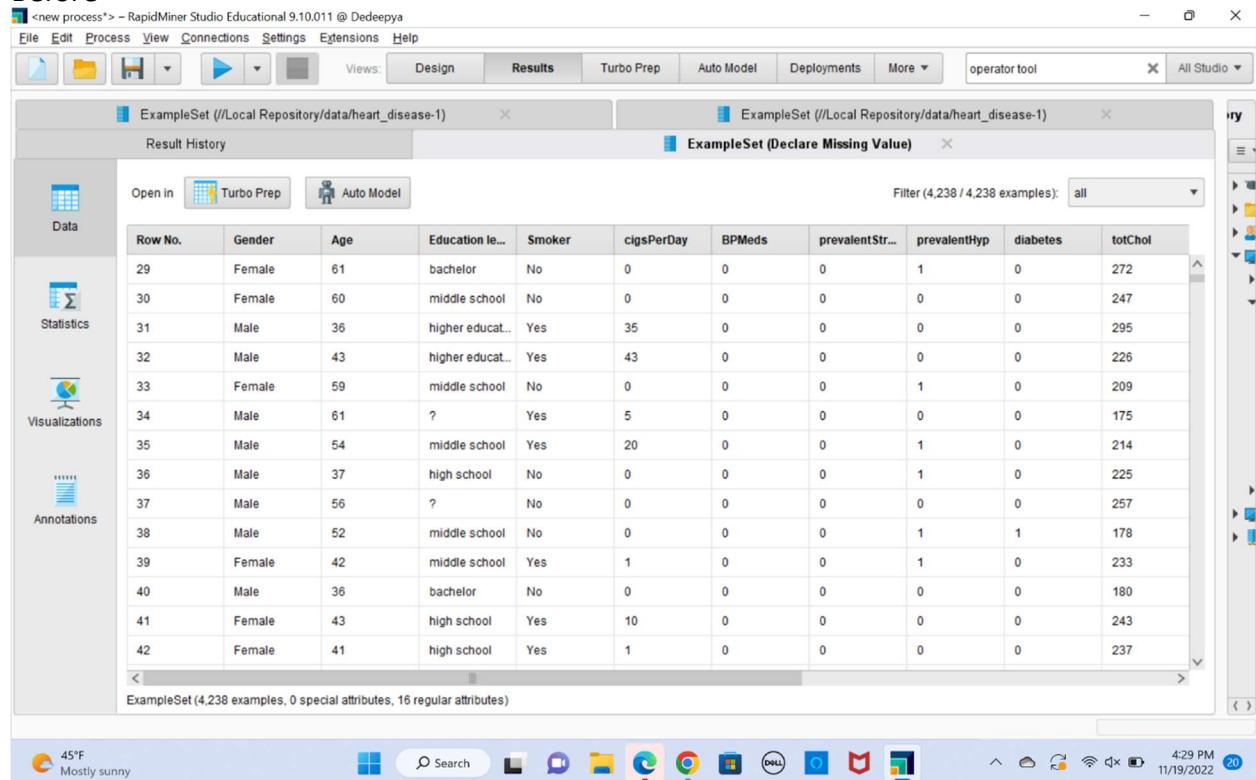
The screenshot shows the 'ExampleSet (Extract Statistics)' results table in RapidMiner Studio. The table lists 13 examples with various numerical attributes like Age, BPMed, and glucose. The 'Annotations' sidebar on the left is visible.

Row No.	Name	Role	Type	Missing	Minimum	Maximum	Average	Deviation	Least	Most	Values	Earliest ...	Latest ...	Duration
1	Age	regular	integer	0	32	70	49.585	8.572	?	?	?	?	?	?
2	cigsPer...	regular	integer	29	0	70	9.003	11.920	?	?	?	?	?	?
3	BPMed	regular	integer	53	0	1	0.030	0.170	?	?	?	?	?	?
4	prevalent...	regular	integer	0	0	1	0.006	0.077	?	?	?	?	?	?
5	prevalent...	regular	integer	0	0	1	0.311	0.463	?	?	?	?	?	?
6	diabetes	regular	integer	0	0	1	0.026	0.158	?	?	?	?	?	?
7	totChol	regular	integer	50	107	696	236.722	44.590	?	?	?	?	?	?
8	sysBP	regular	integer	0	84	295	132.450	22.037	?	?	?	?	?	?
9	diaBP	regular	integer	0	48	143	82.974	11.907	?	?	?	?	?	?
10	BMI	regular	integer	19	16	57	25.809	4.092	?	?	?	?	?	?
11	heartRate	regular	integer	1	44	143	75.879	12.027	?	?	?	?	?	?
12	glucose	regular	integer	388	40	394	81.967	23.960	?	?	?	?	?	?
13	TenYear...	regular	integer	0	0	1	0.152	0.359	?	?	?	?	?	?

4. Encode categorical variables and include the before and after screenshots of each categorical variable.
  - A. For a categorical variable with more than 2 unique values, create a dummy variable for each unique value. (Use nominal to Numerical operator, this will create one column for each type of value in columns)
  - B. For a binary categorical variable, encode it using 0 and 1. (Use the Nominal to numerical operator to replace binary categorical variables to numerical, use Unique integers coding type, it will just convert to integer by replacing one value with 0 and one value with 1, use this only if the column has two unique values).



## Before-



ExampleSet (//Local Repository/data/heart\_disease-1) ExampleSet (//Local Repository/data/heart\_disease-1) ExampleSet (Declare Missing Value)

Result History

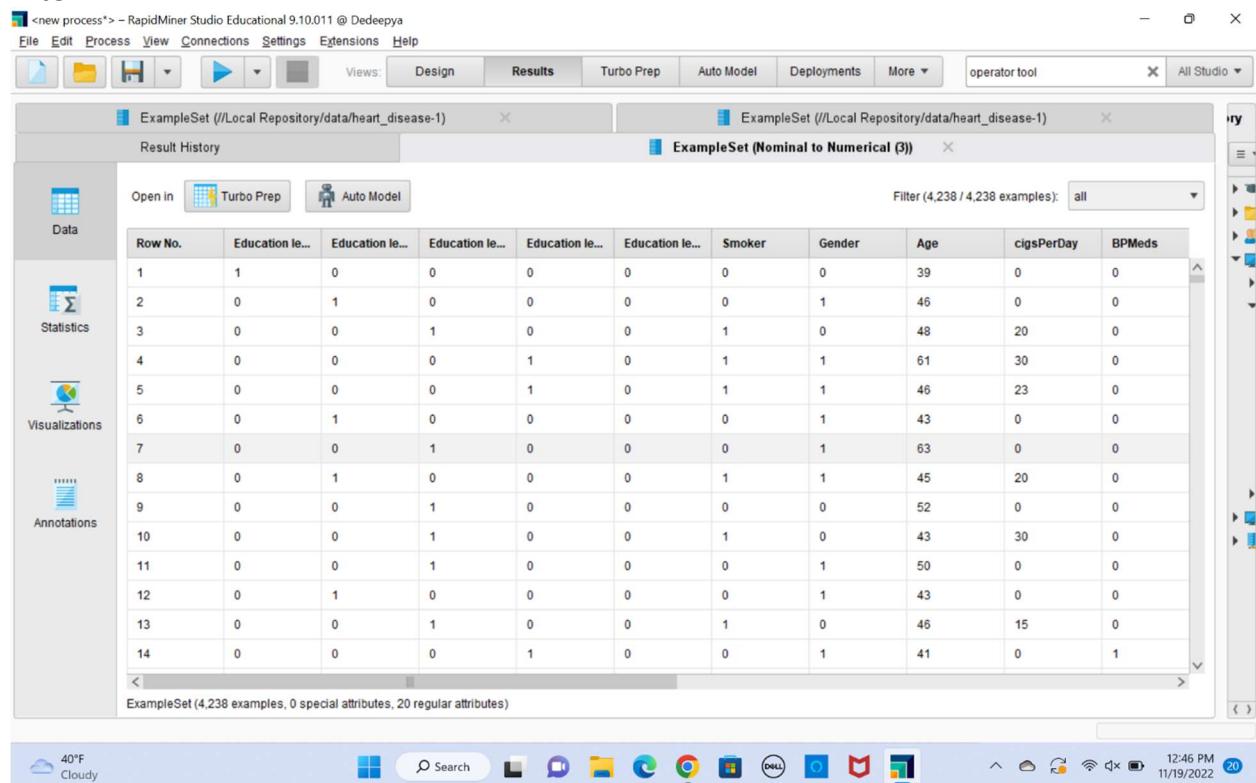
Open in Turbo Prep Auto Model

Filter (4,238 / 4,238 examples): all

Row No.	Gender	Age	Education level	Smoker	cigsPerDay	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol
29	Female	61	bachelor	No	0	0	0	1	0	272
30	Female	60	middle school	No	0	0	0	0	0	247
31	Male	36	higher education	Yes	35	0	0	0	0	295
32	Male	43	higher education	Yes	43	0	0	0	0	226
33	Female	59	middle school	No	0	0	0	1	0	209
34	Male	61	?	Yes	5	0	0	0	0	175
35	Male	54	middle school	Yes	20	0	0	1	0	214
36	Male	37	high school	No	0	0	0	1	0	225
37	Male	56	?	No	0	0	0	0	0	257
38	Male	52	middle school	No	0	0	0	1	1	178
39	Female	42	middle school	Yes	1	0	0	1	0	233
40	Male	36	bachelor	No	0	0	0	0	0	180
41	Female	43	high school	Yes	10	0	0	0	0	243
42	Female	41	high school	Yes	1	0	0	0	0	237

ExampleSet (4,238 examples, 0 special attributes, 16 regular attributes)

## After-



ExampleSet (//Local Repository/data/heart\_disease-1) ExampleSet (//Local Repository/data/heart\_disease-1) ExampleSet (Nominal to Numerical (3))

Result History

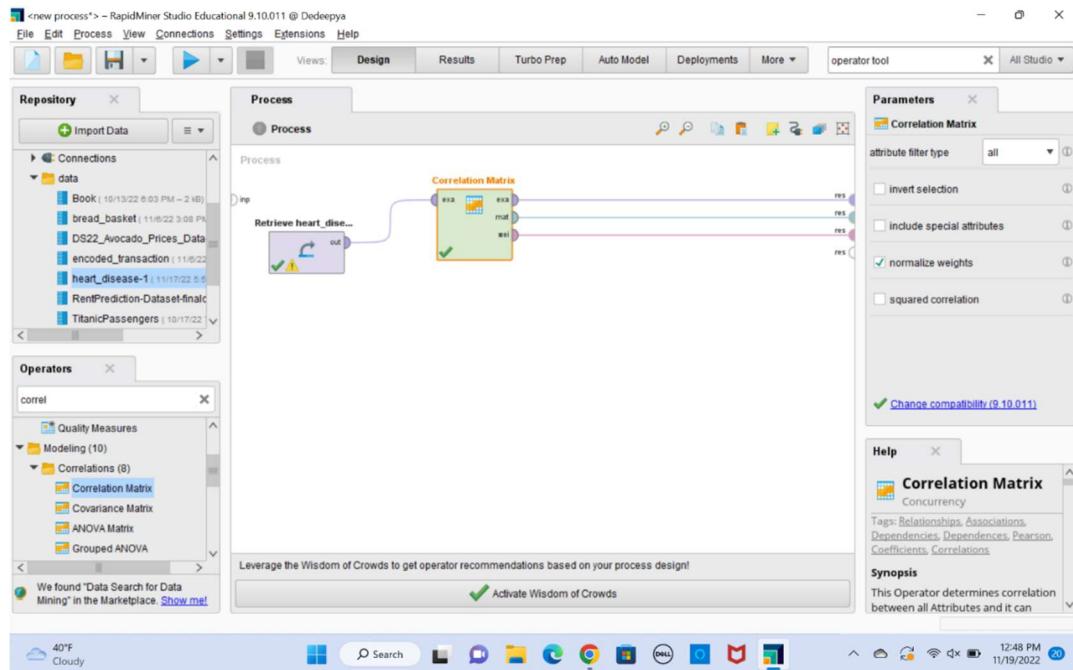
Open in Turbo Prep Auto Model

Filter (4,238 / 4,238 examples): all

Row No.	Education level	Smoker	Gender	Age	cigsPerDay	BPMeds					
1	1	0	0	0	0	0	0	0	39	0	0
2	0	1	0	0	0	0	0	1	46	0	0
3	0	0	1	0	0	1	0	0	48	20	0
4	0	0	0	1	0	1	1	1	61	30	0
5	0	0	0	1	0	1	1	1	46	23	0
6	0	1	0	0	0	0	1	1	43	0	0
7	0	0	1	0	0	0	1	1	63	0	0
8	0	1	0	0	0	1	1	1	45	20	0
9	0	0	1	0	0	0	0	0	52	0	0
10	0	0	1	0	0	1	0	0	43	30	0
11	0	0	1	0	0	0	1	1	50	0	0
12	0	1	0	0	0	0	1	1	43	0	0
13	0	0	1	0	0	1	0	0	46	15	0
14	0	0	0	1	0	0	1	1	41	0	1

ExampleSet (4,238 examples, 0 special attributes, 20 regular attributes)

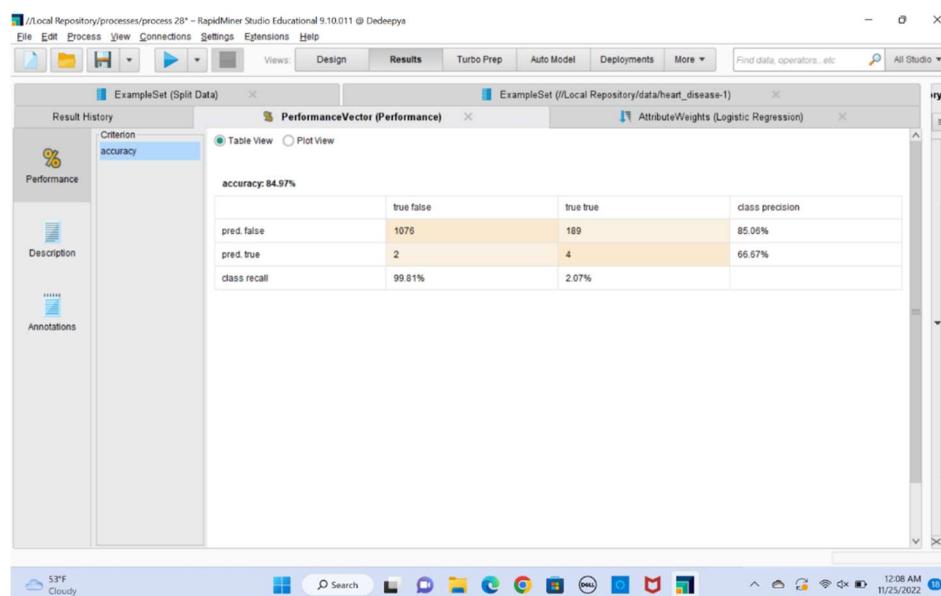
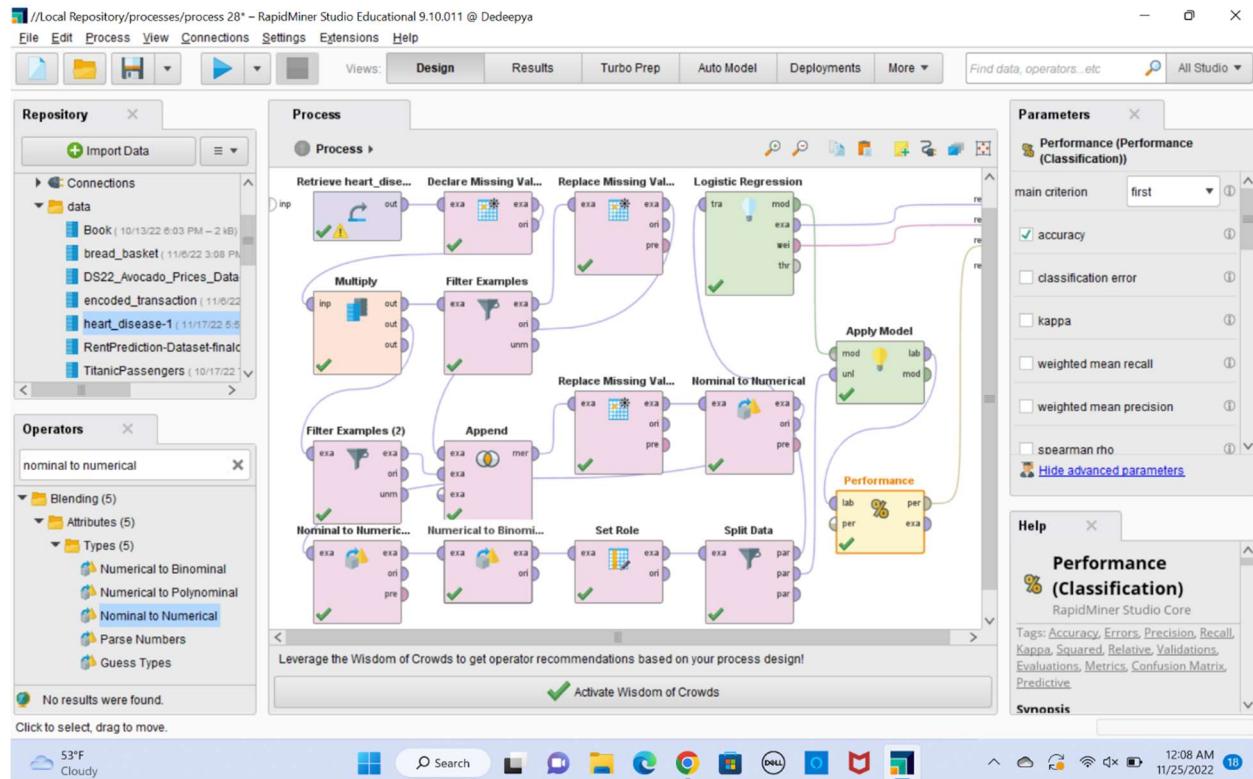
5. Create a correlation matrix and analyze it. Round the correlation values to the 2<sup>nd</sup> nearest decimal place. List the variables whose correlation with "TenYearCHD" is at least 0.1. List the variables whose correlation with "TenYearCHD" is between -0.1 and -0.01.



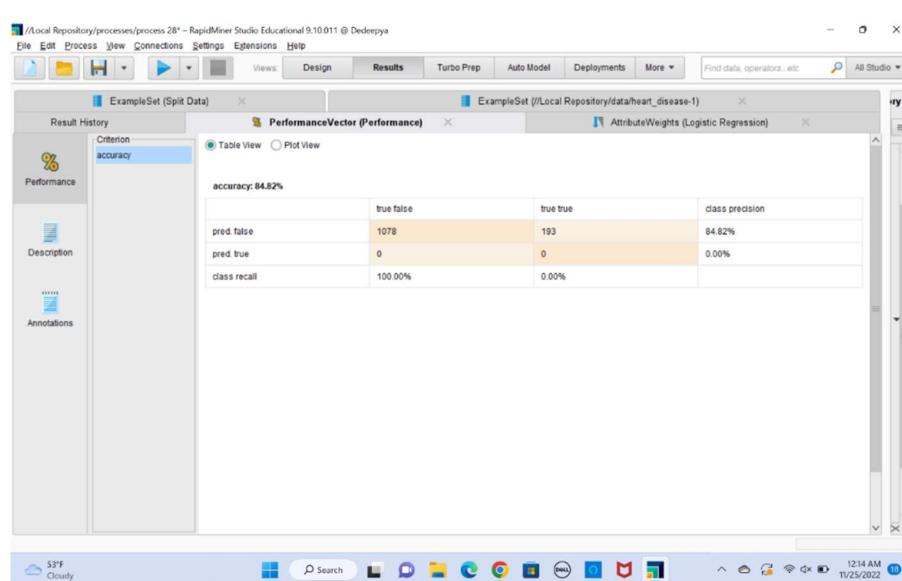
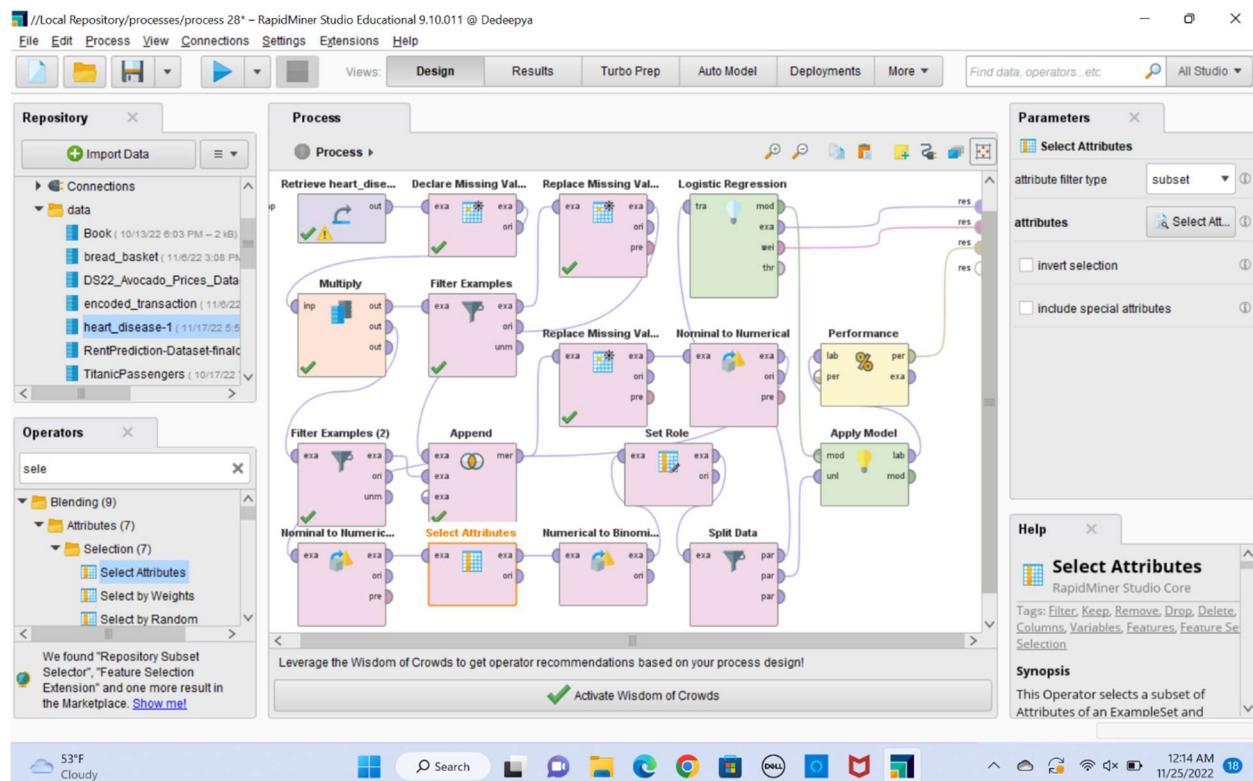
Columns with values of correlation At least 0.1 are - age, prevalenthyp, sysbp, diabp, glucose.

Columns with values of correlation Between -0.1 to -0.01 are – gender,

6. Split the data into training and test datasets with a ratio of 70:30 using a random state of 58. Build a Logistic Regression model using all independent variables. What is the test accuracy? How many records were misclassified?



7. Using the same training and test sets above, build a Logistic Regression model using the variables identified in Question 5 and a random state of 58, and evaluate it with the test set. What is the test accuracy? How many records were misclassified? Compared to the model built in Question 6, which model is better? Why?



From both the models we can see that there is high accuracy for model build with all the features instead of models with less features. As we reduced number of features by using correlation matrix the performance of second model decreased.