



# FINAL REPORT ISE 543, SPRING 2023



Project objective - build a classification model to predict the likelihood of a patient developing Chronic Heart Disease (CHD) in the coming ten years:

Age	BMI
Male	Heart Rate
Education	Glucose
Income	A I c
Current Smoker	Ten Year CHD
Cigarettes per Day	
BP Meds	
Prevalent Stroke	
Prevalent Hyp	
Diabetes	
Total Chol	
Sys BP	
Dia BP	



# **EXPLORATORY DATA ANALYSIS/ DATA PREPARATION**



## Response Variable

- › Ten Year CHD

## Categories

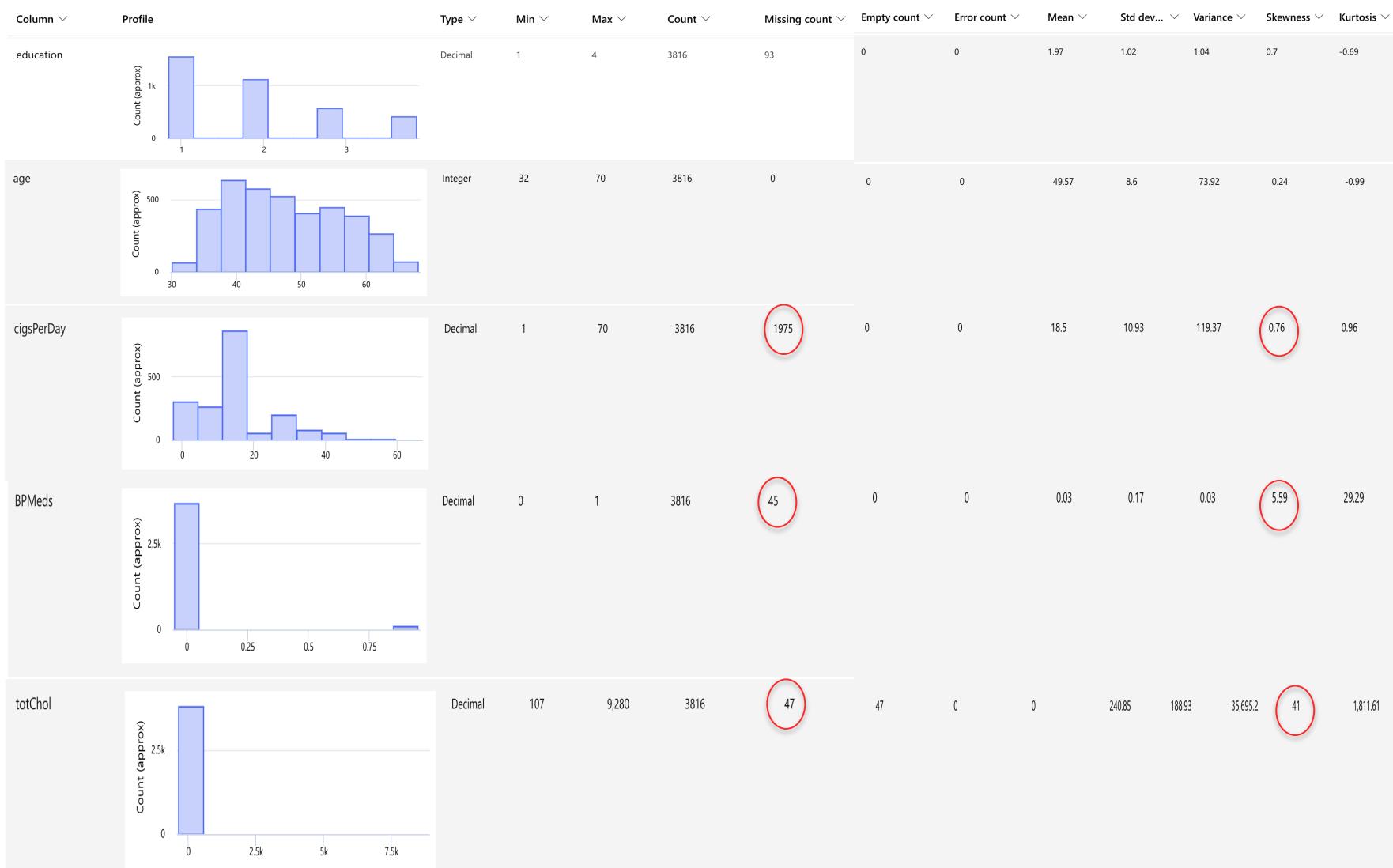
- › Male
- › Education
- › Current Smoker
- › BP Meds
- › Prevalent Stroke
- › Prevalent Hyp
- › Diabetes

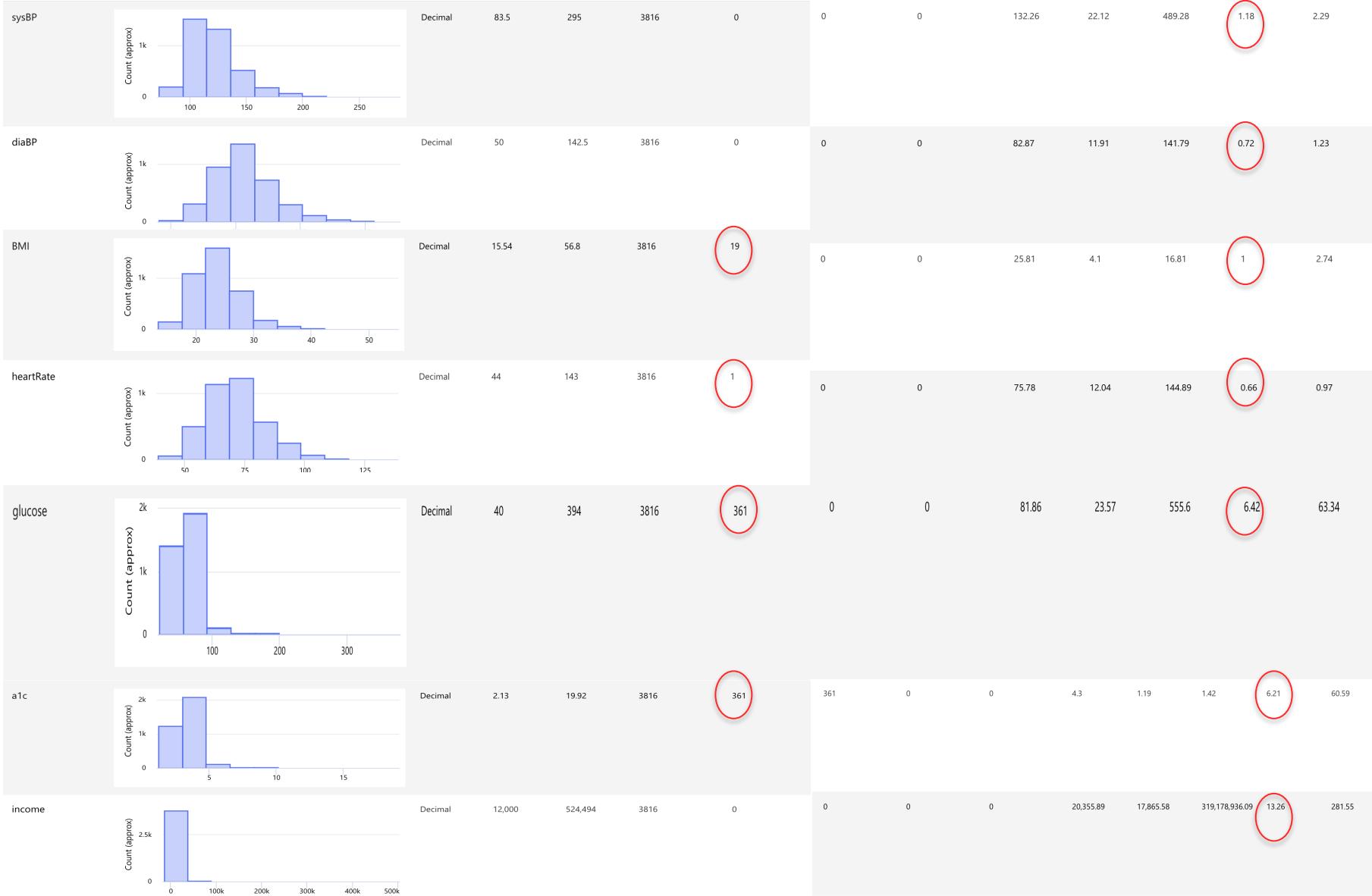
## Measures

- › Income
- › Age
- › Cigarettes per Day
- › Total Chol
- › Sys BP
- › Dia BP
- › BMI
- › Heart Rate
- › Glucose
- › A l c

} Need to determine how  
to variable

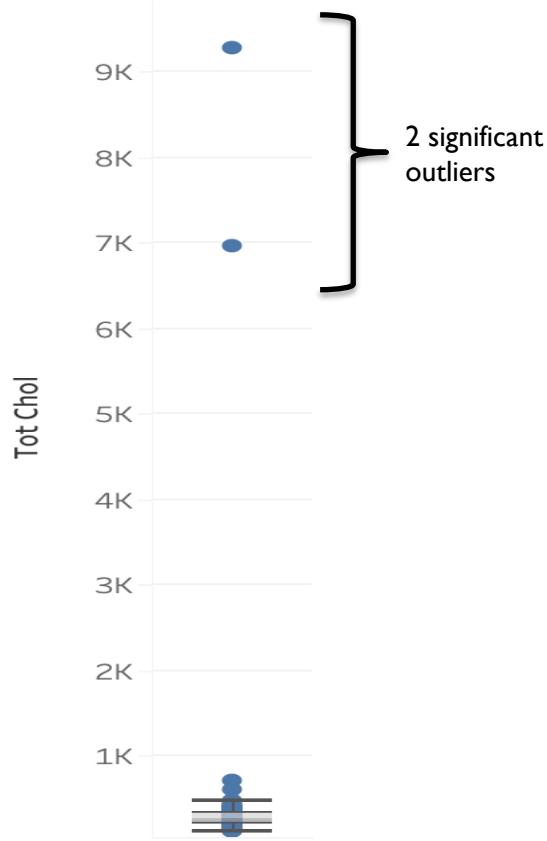
## DATA CLEANSING



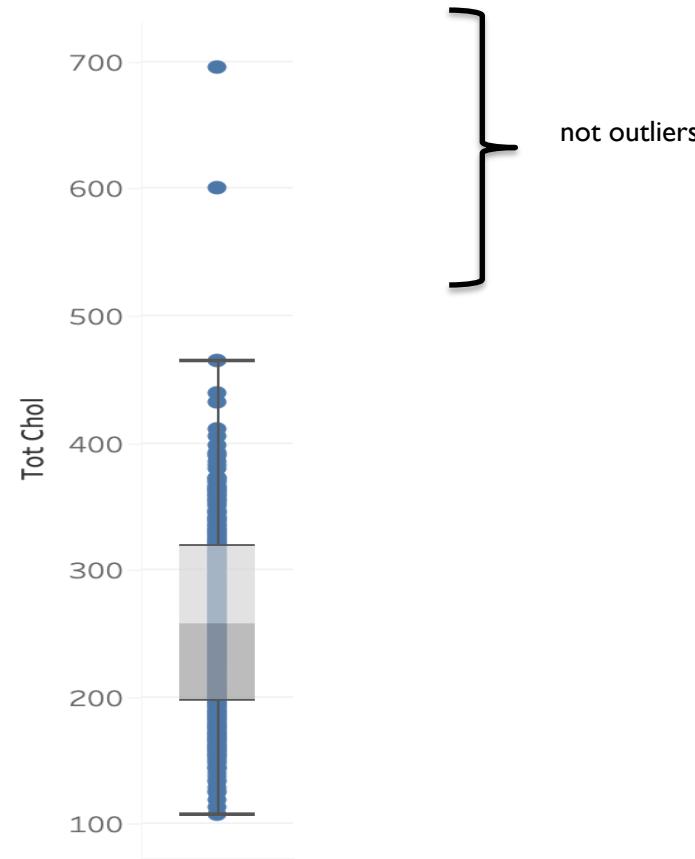




Total Cholesterol



Total Cholesterol



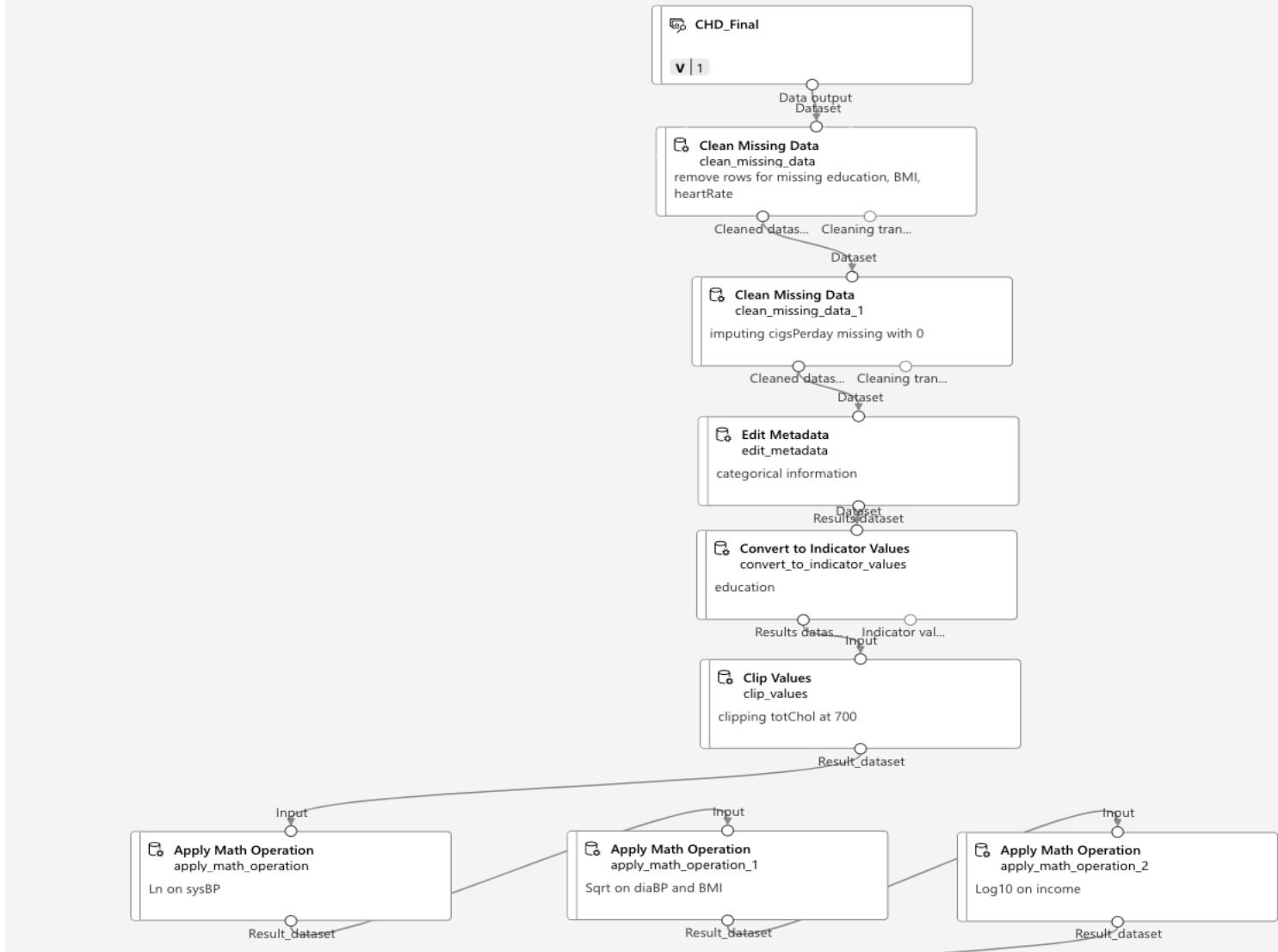


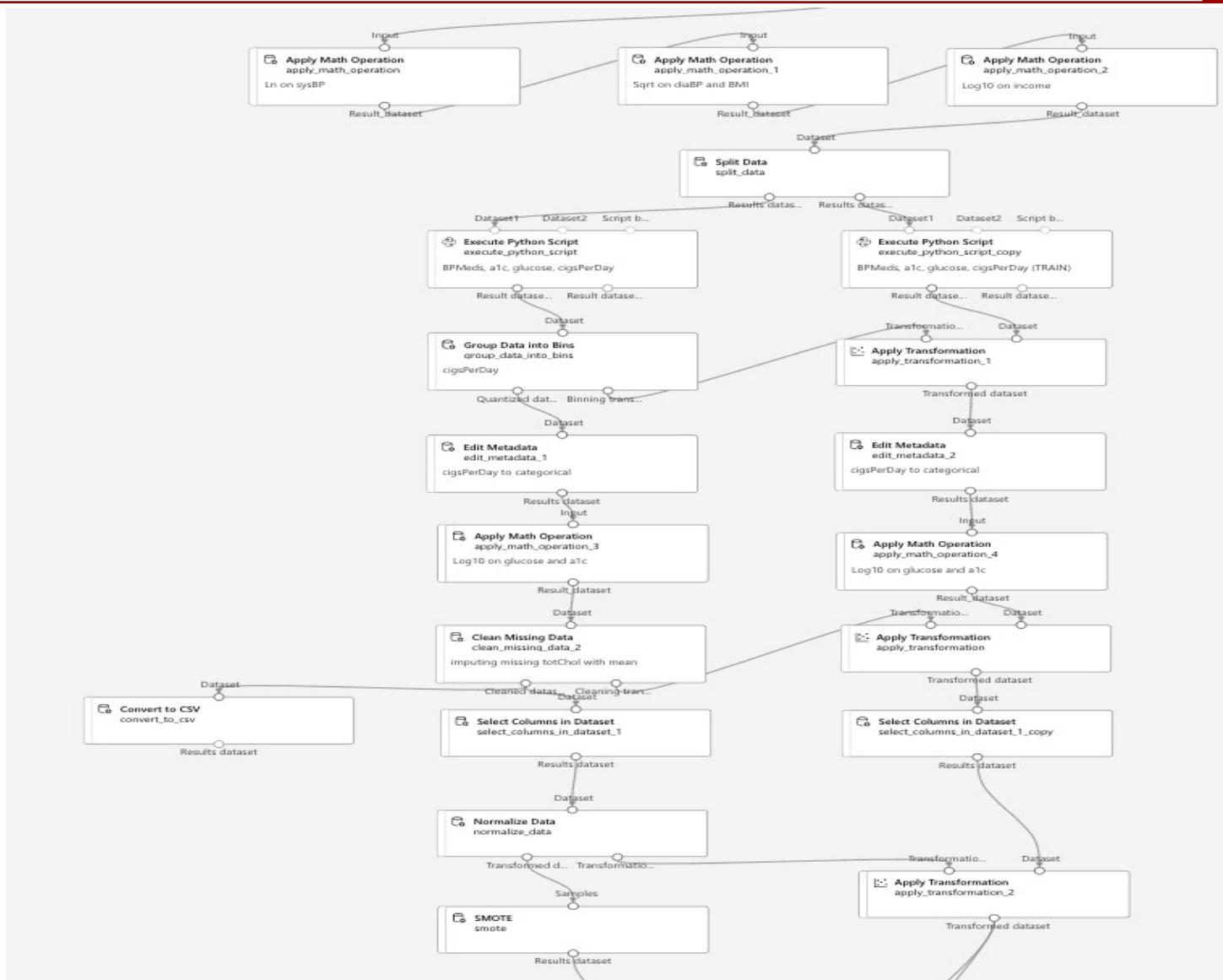
- › Remove rows with missing values for education, BMI and heartrate because which do not have any information which indicates the education or BMI of the patient and imputing summary statistic in these cases can be misleading.
- › Impute 0 for missing values in cigsPerDay variable as most of the missing variables correspond to patients that do are not current smokers.
- › Remove outliers from totChol variable by clipping values greater than 646
- › Impute missing values in BPMeds as 0 if the corresponding sysBP is greater than 140 and 0 otherwise. This is called stage 2 high BP and it is highly likely that they are on medication.
- › Impute missing values in glucose and a1c by taking the average glucose or a1c of the group (diabetic or non-diabetic) that the patient belongs to as both glucose and a1c are related to diabetes



- › For patients that are current smokers and have missing values in cigsPerDay variables, impute with average cigsPerDay
- › Apply binning on cigsPerDay as 0,1,10,31,40,50 and convert it into a categorical variable
- › Impute missing values in totChol with average

## DATA CLEANSING PIPELINE

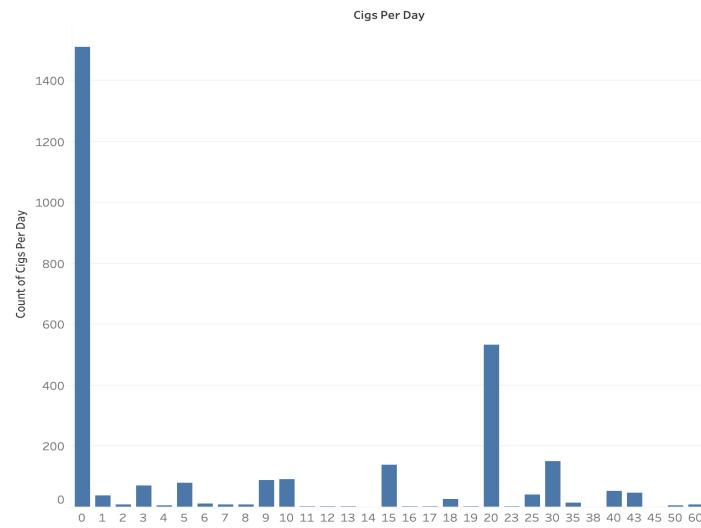




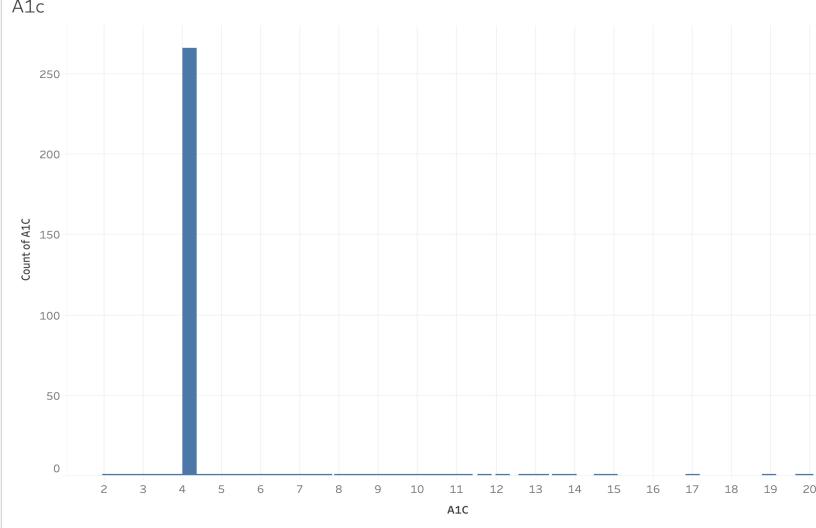
# UNIVARIATE ANALYSIS MEASURES



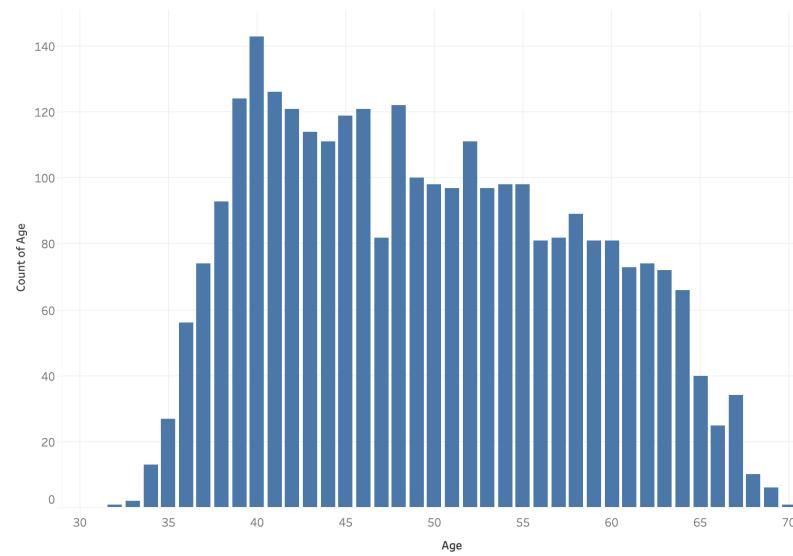
CigsPerDay



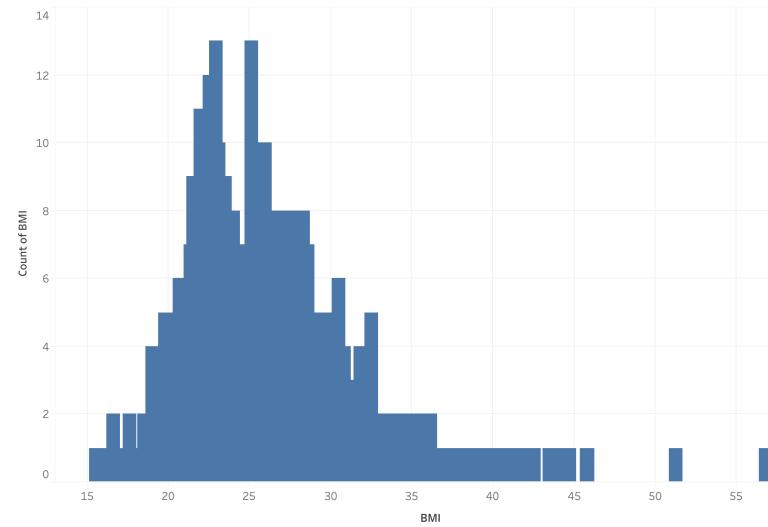
A1c



Age



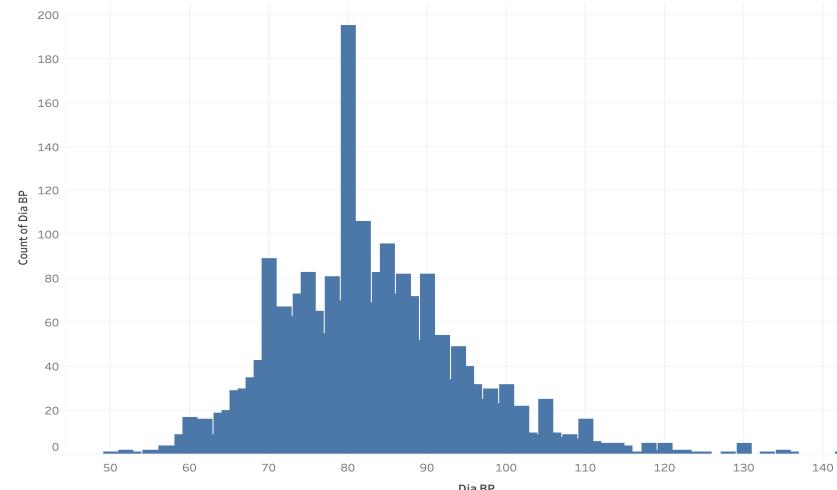
BMI



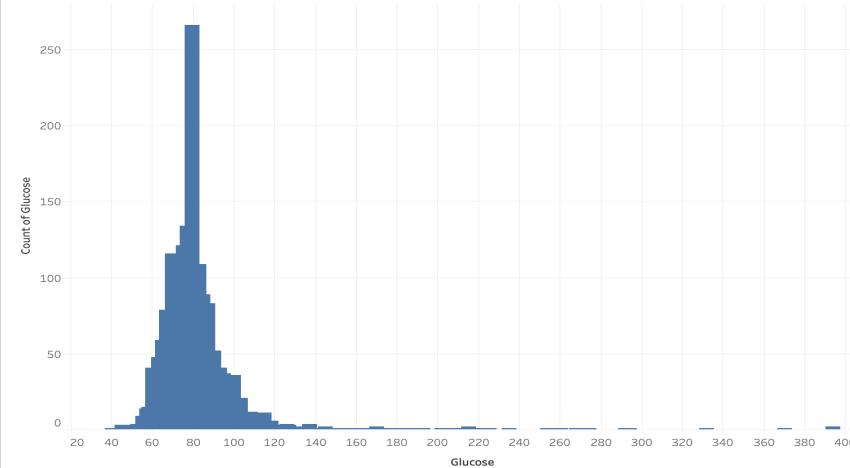
# UNIVARIATE ANALYSIS MEASURES



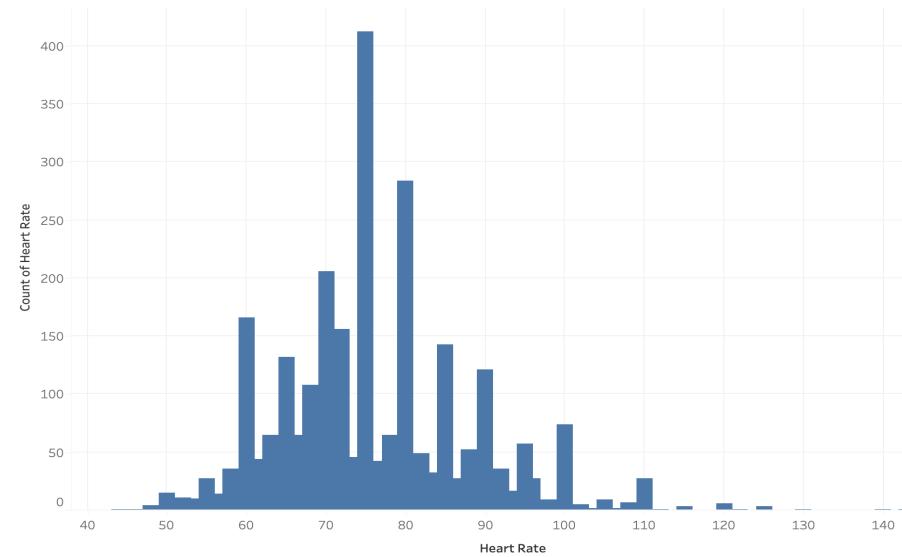
DiaBP



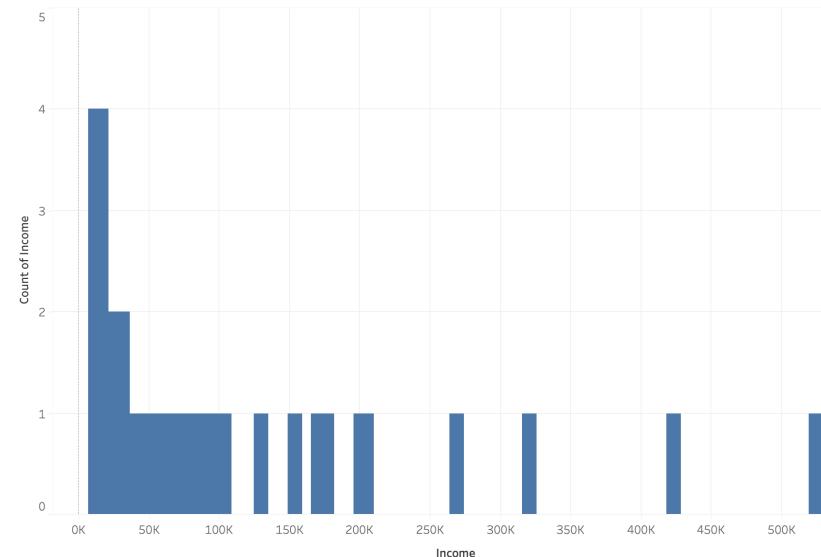
Glucose



Heart Rate

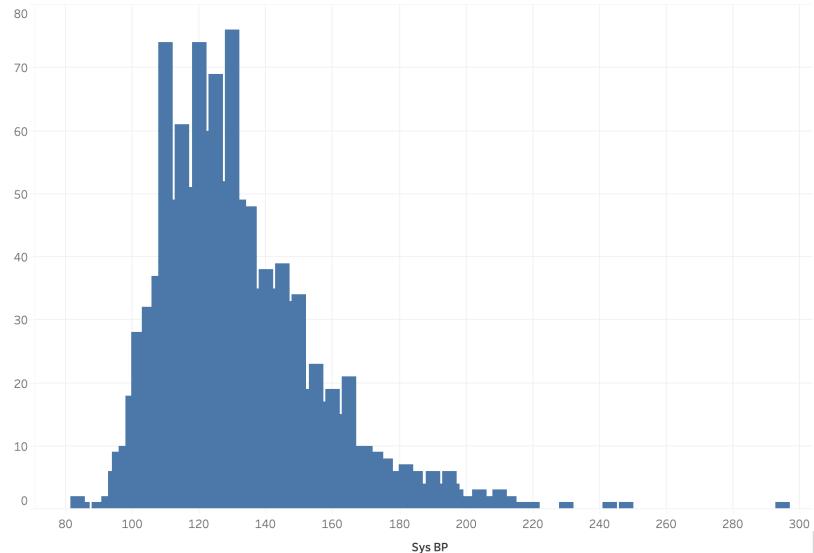


Income

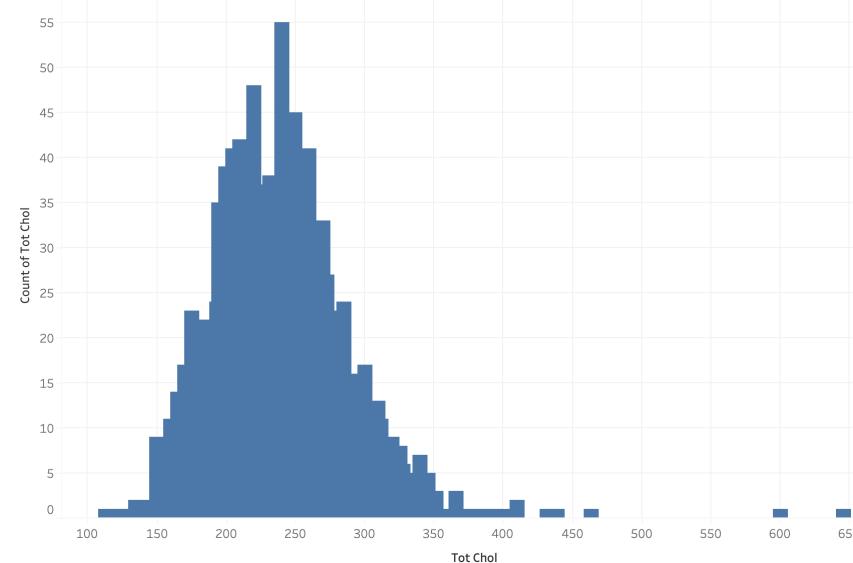




Sys BP

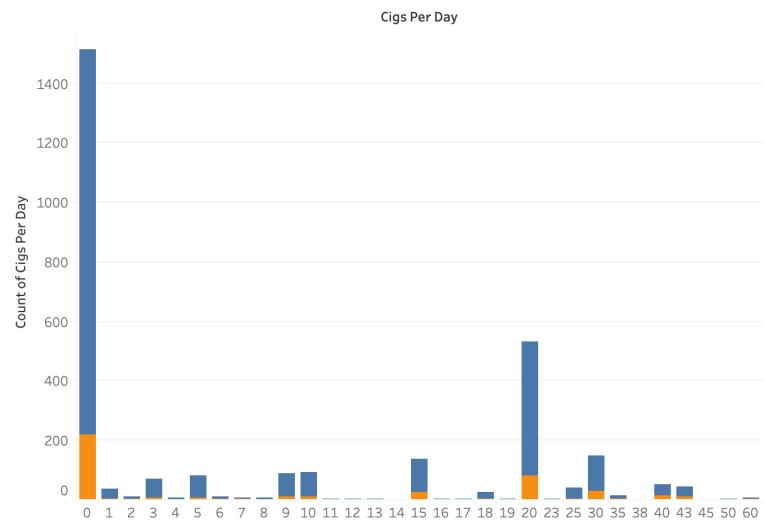


Tot Chol

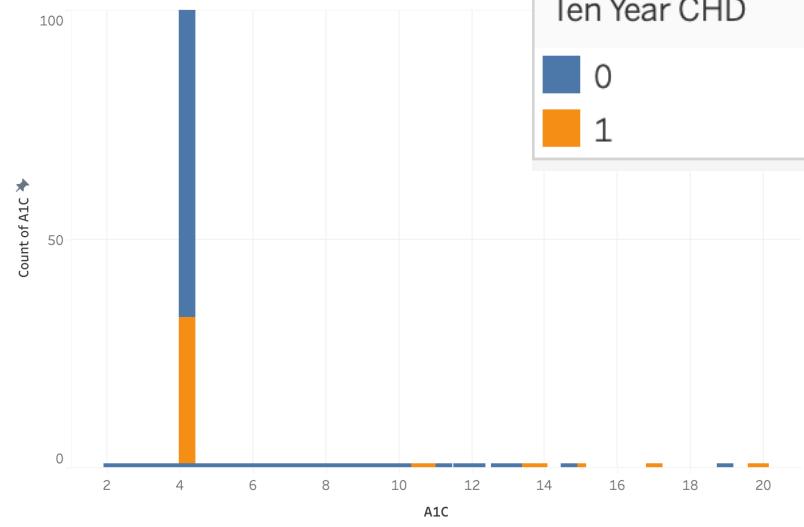




CigsPerDay vs CHD



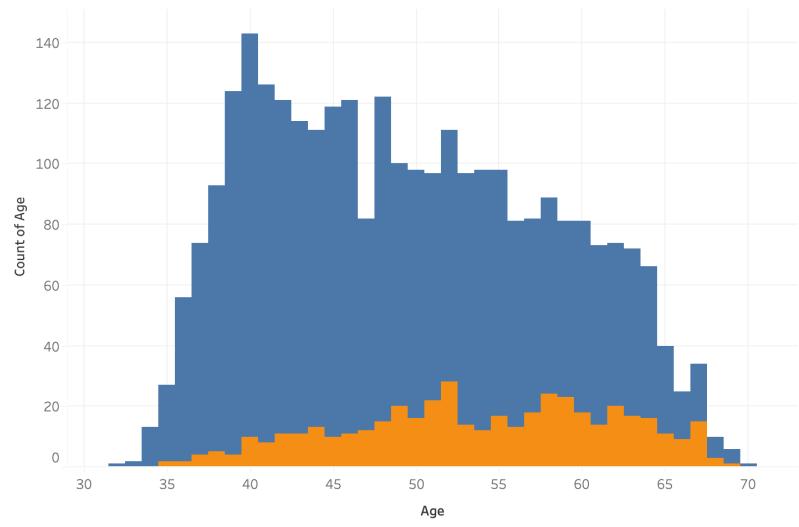
A1c vs CHD



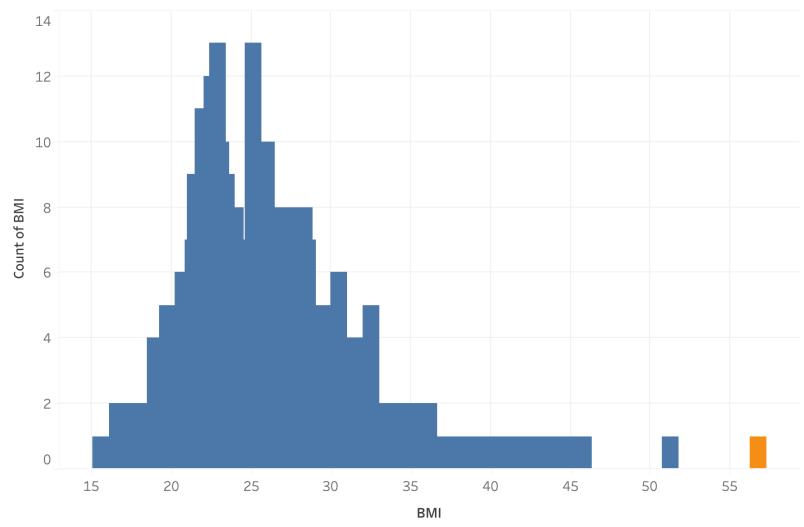
Ten Year CHD

0  
1

Age vs CHD

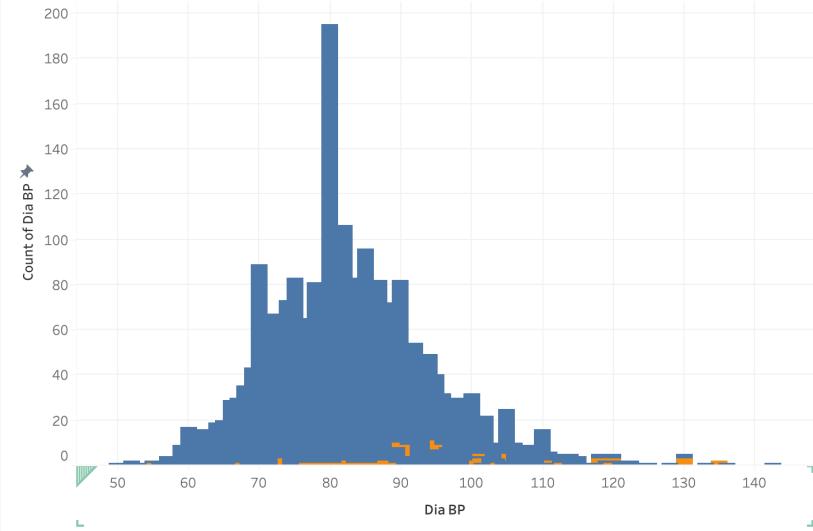


BMI vs CHD

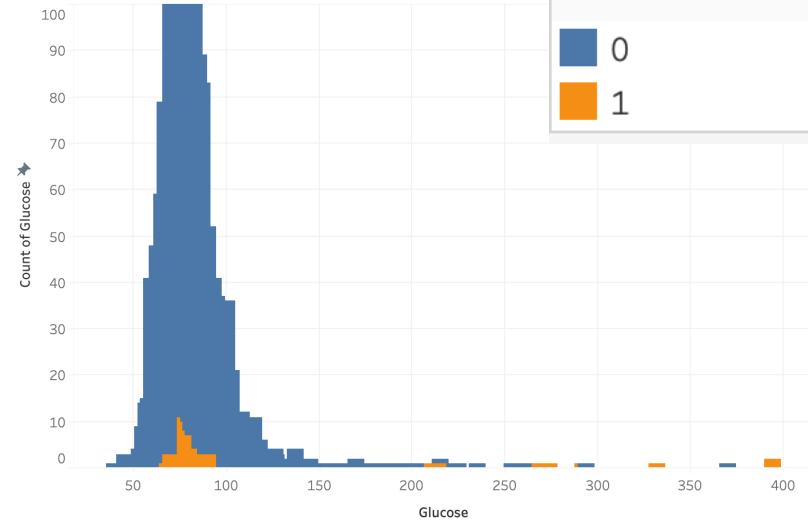




DiaBP vs CHD



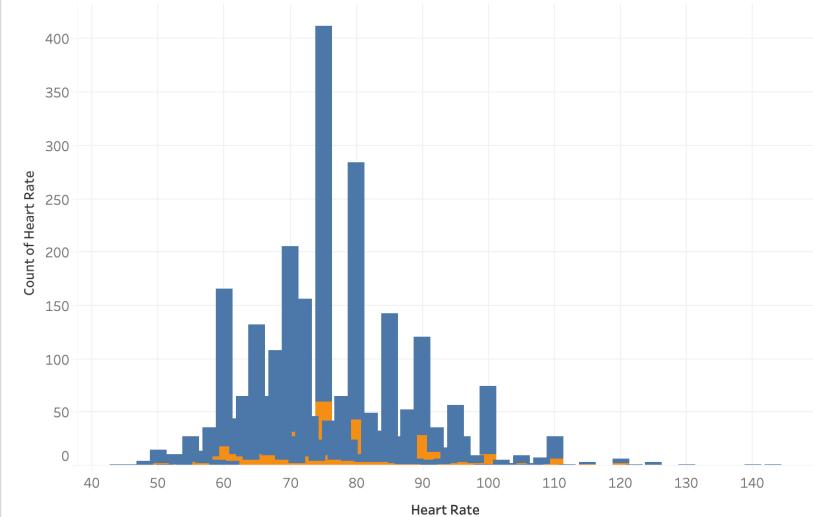
Glucose vs CHD



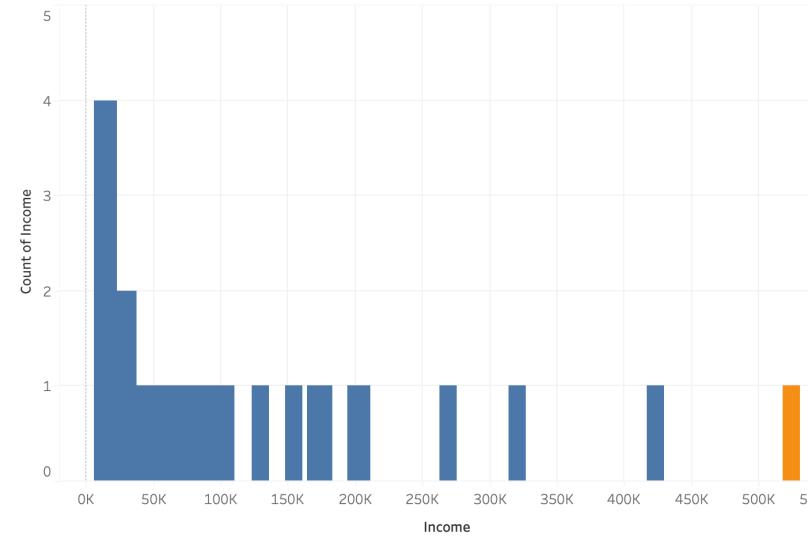
Ten Year CHD

0  
1

Heart Rate vs CHD

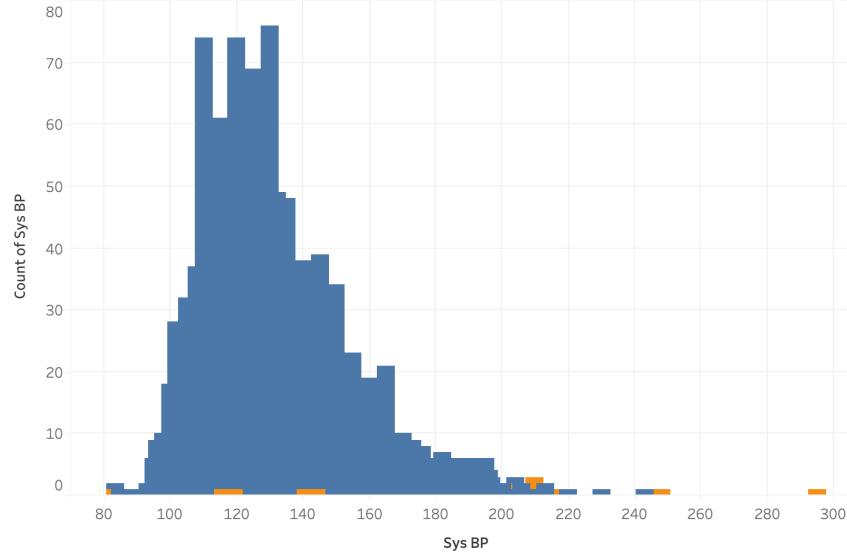


Income vs CHD





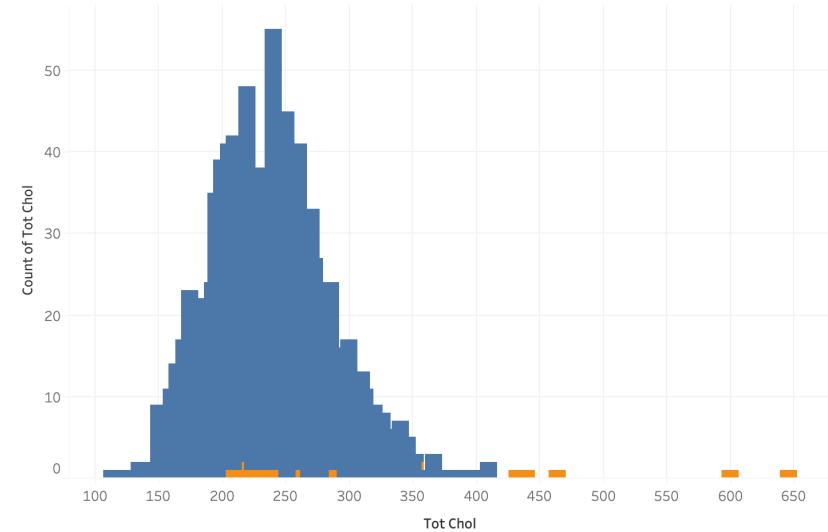
Sys BP vs CHD



Ten Year CHD

- 0
- 1

Tot Chol vs CHD

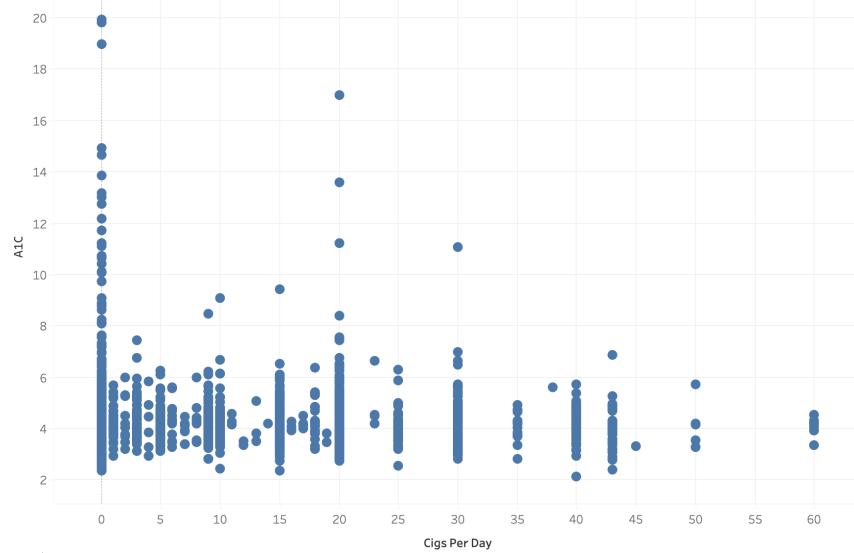


Determined that outliers are not to be deleted considering the relationship to response variable and business knowledge (nature of the variable) i.e. it is not uncommon to see such large/extreme values

# CORRELATION ANALYSIS



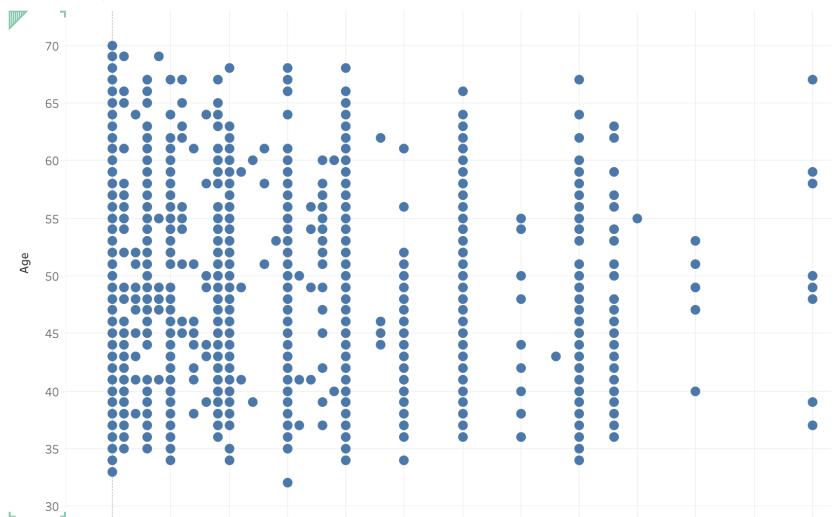
CigsPerDay vs a1c



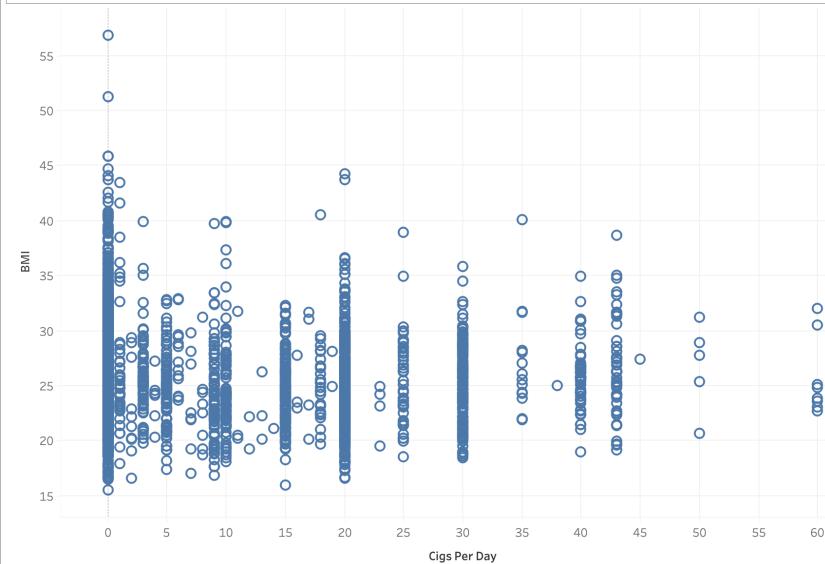
BMI vs Age



cigsPerDay vs Age

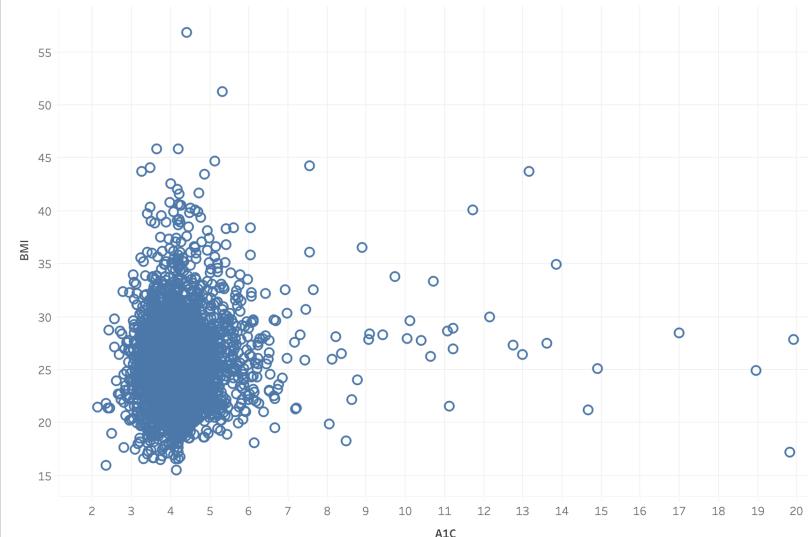


CigsPerDay vs BMI

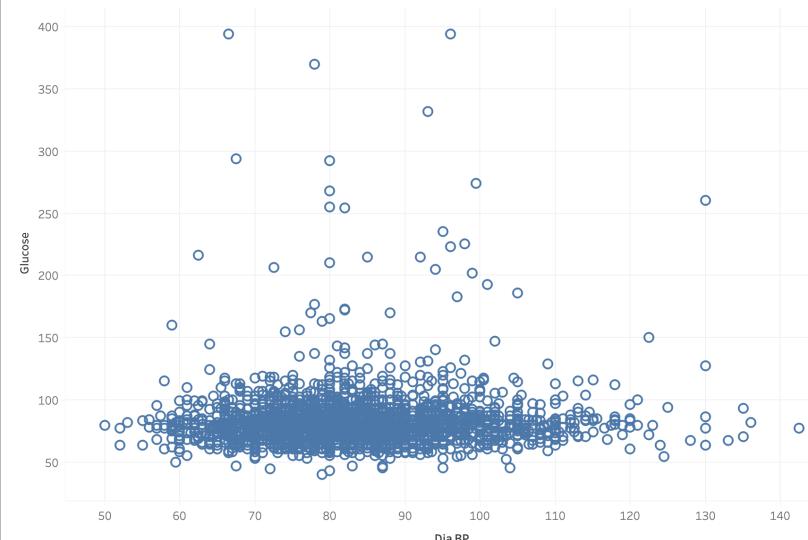




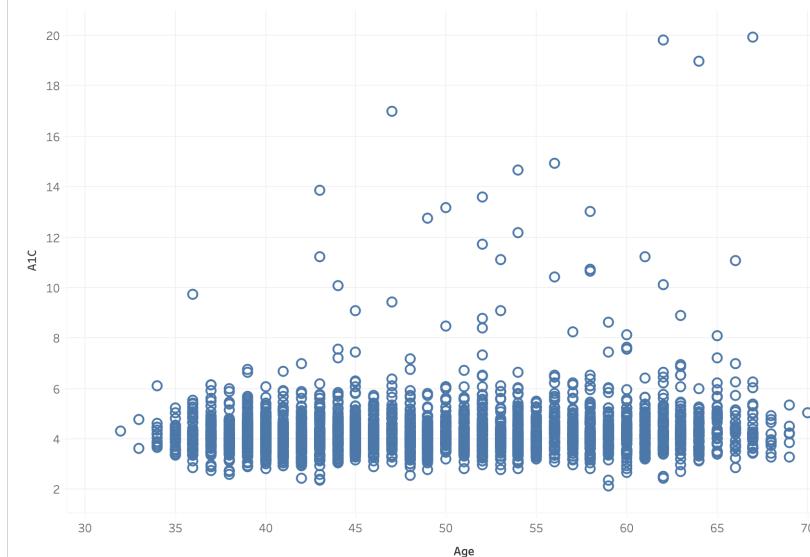
a1c vs BMI



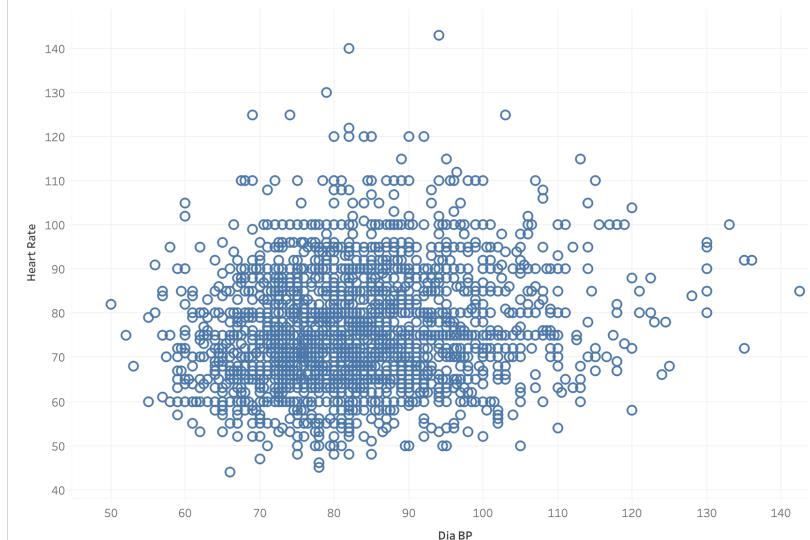
DiaBP vs Glucose



a1c vs age

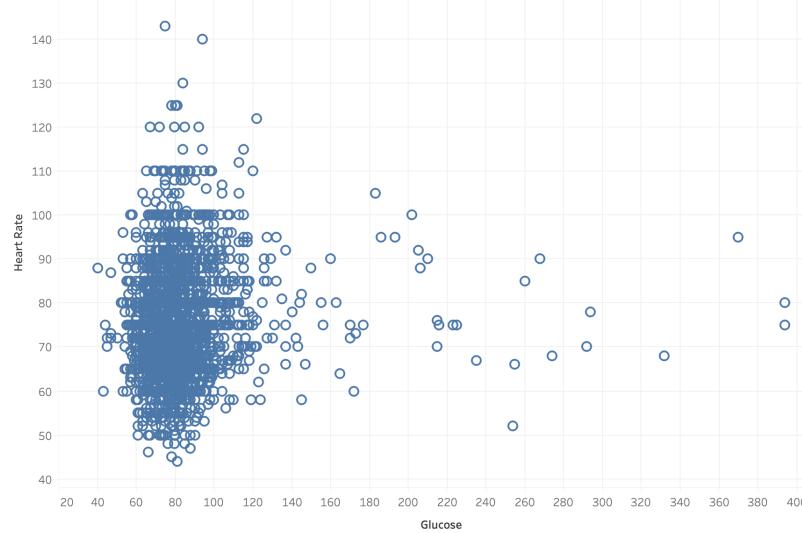


DiaBP vs Heart Rate



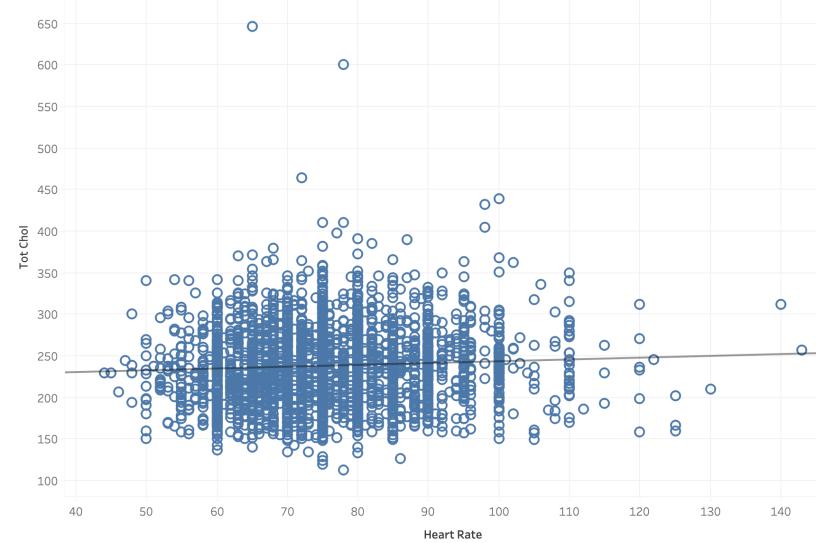


glucose vs heart rate



No Significant correlations found

heart rate vs tot Chol





## Features engineering

- › Convert education into a categorical variable
- › Apply Ln function on sysBP
- › Apply Sqrt function on diaBP and BMI
- › Apply Log10 function on income, glucose a1c
- › Bin cigsPerDay variavble into 6 bins after looking at bivariate analysis and convert them into categories

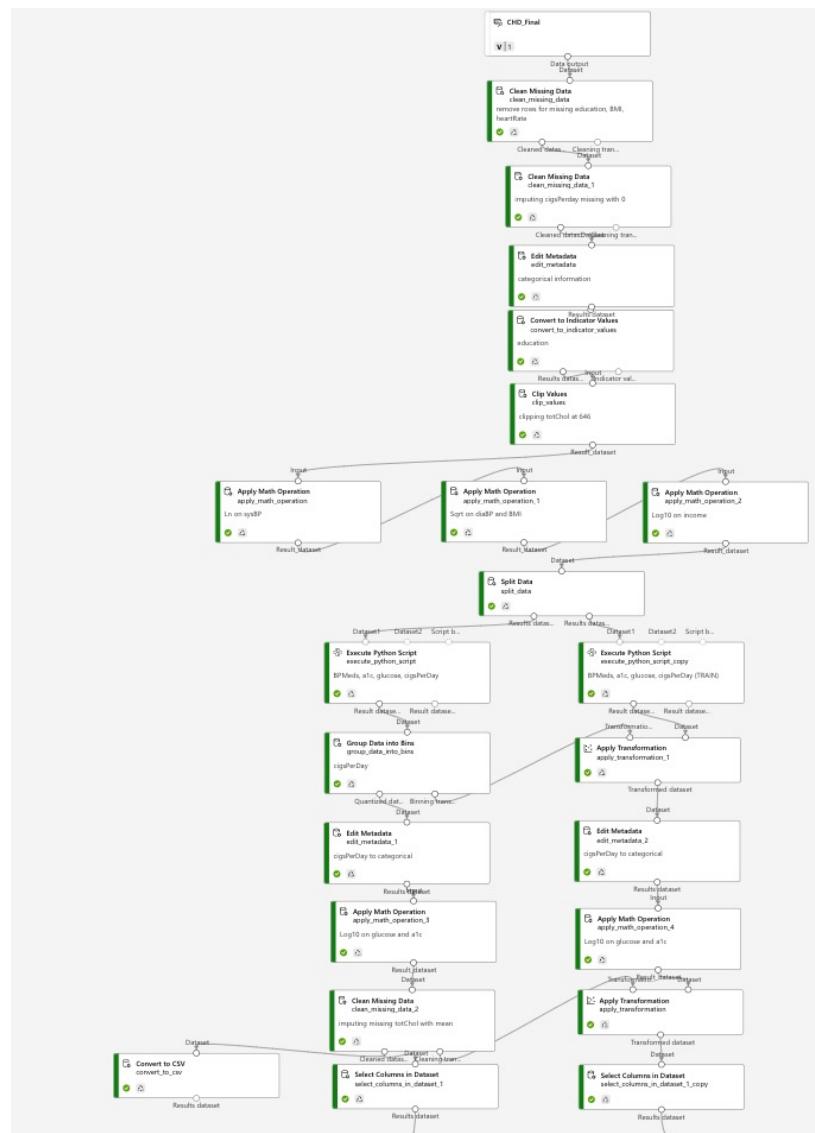
## Other Decisions

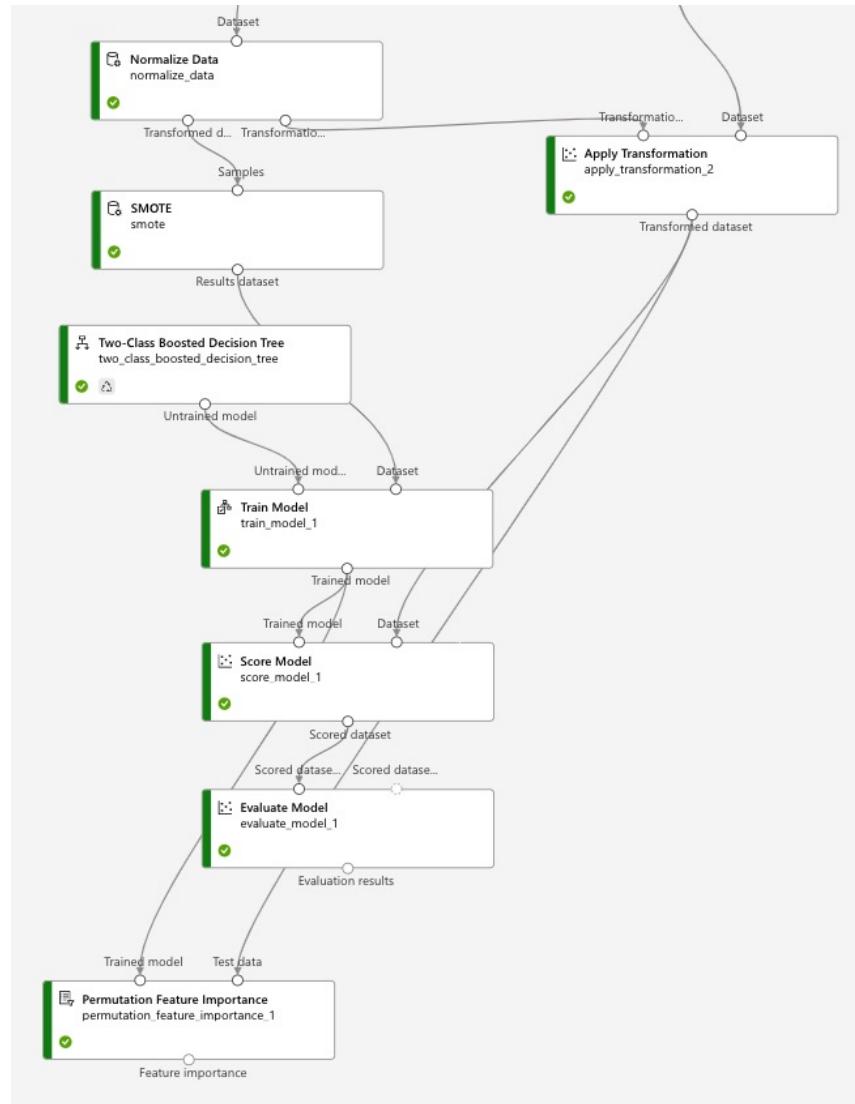
- › Split data into 80-20 and prevent information leakage
- › Normalize all measures (except patientID) using Zscore transformation
- › Use SMOTE, with 200% and 3 nearest neighbors, as positive obs in response variable are only 15%



# MODEL DEVELOPMENT

# BUILD MODEL PIPELINE







Threshold  0.5

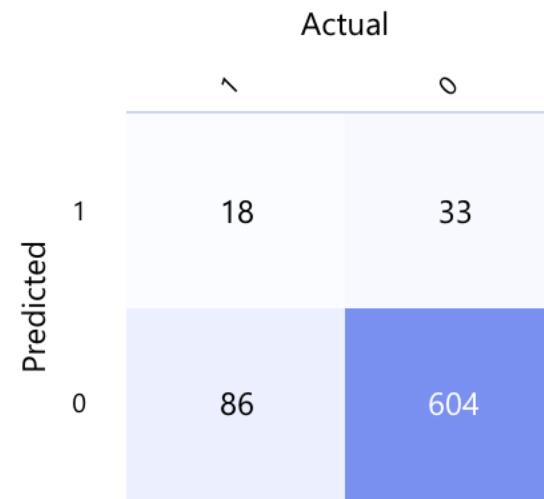
Accuracy 0.839

Precision 0.353

Recall 0.173

F1 Score 0.232

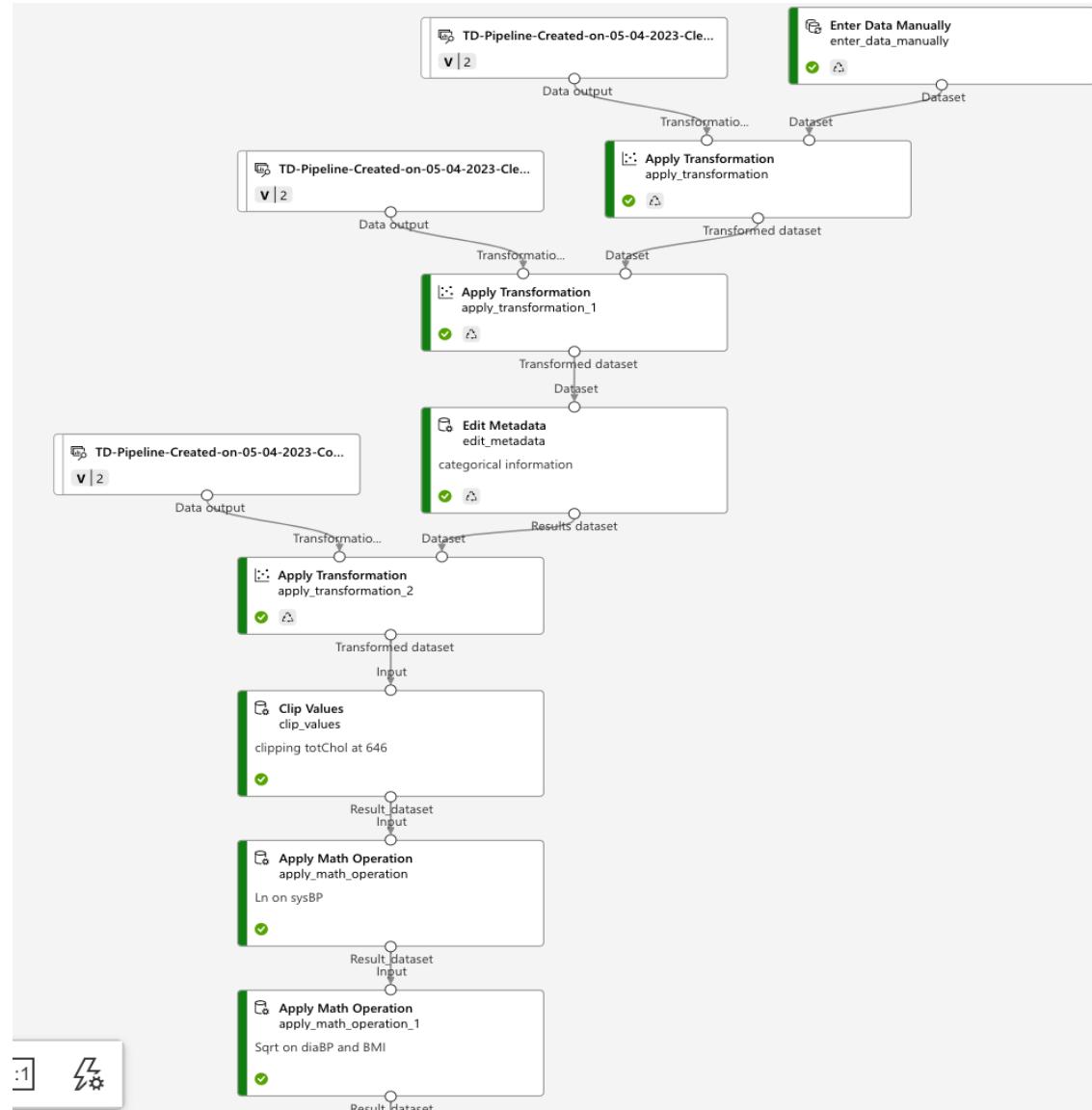
AUC 0.662

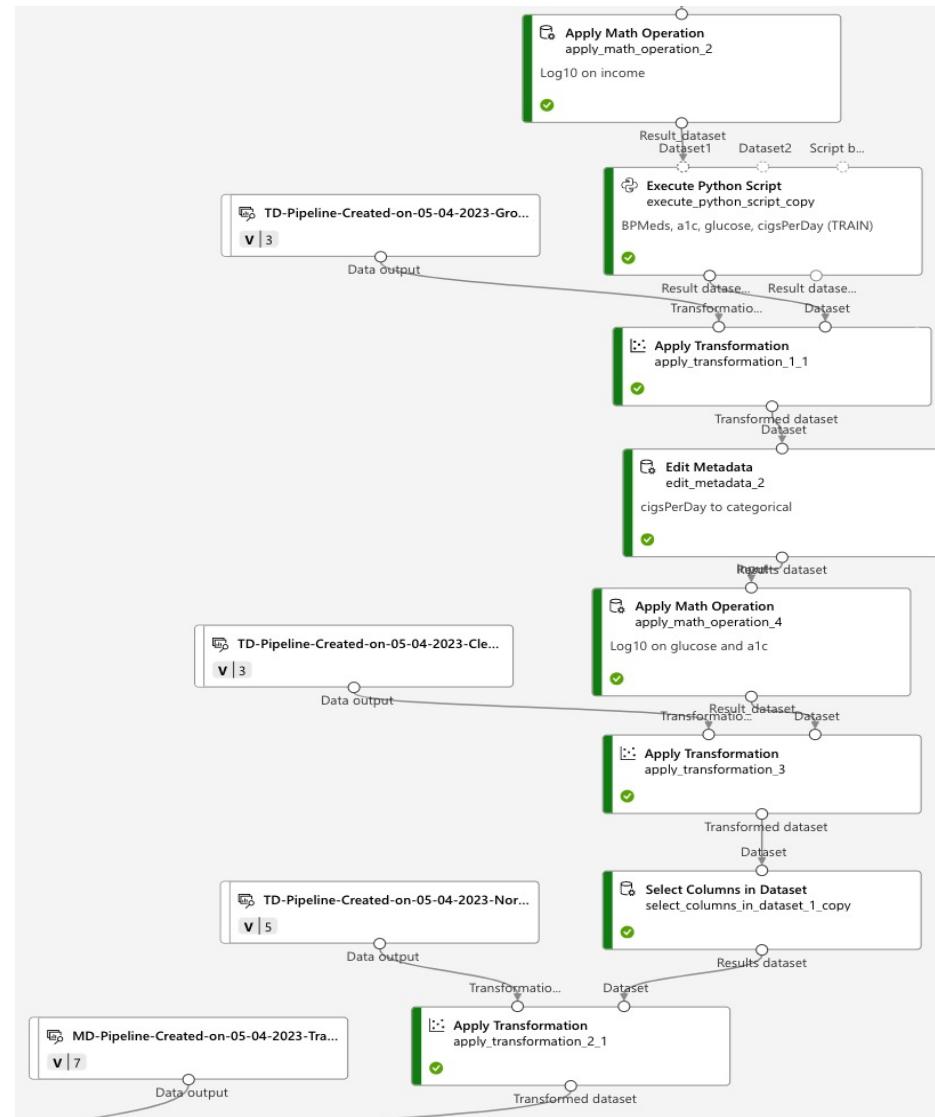


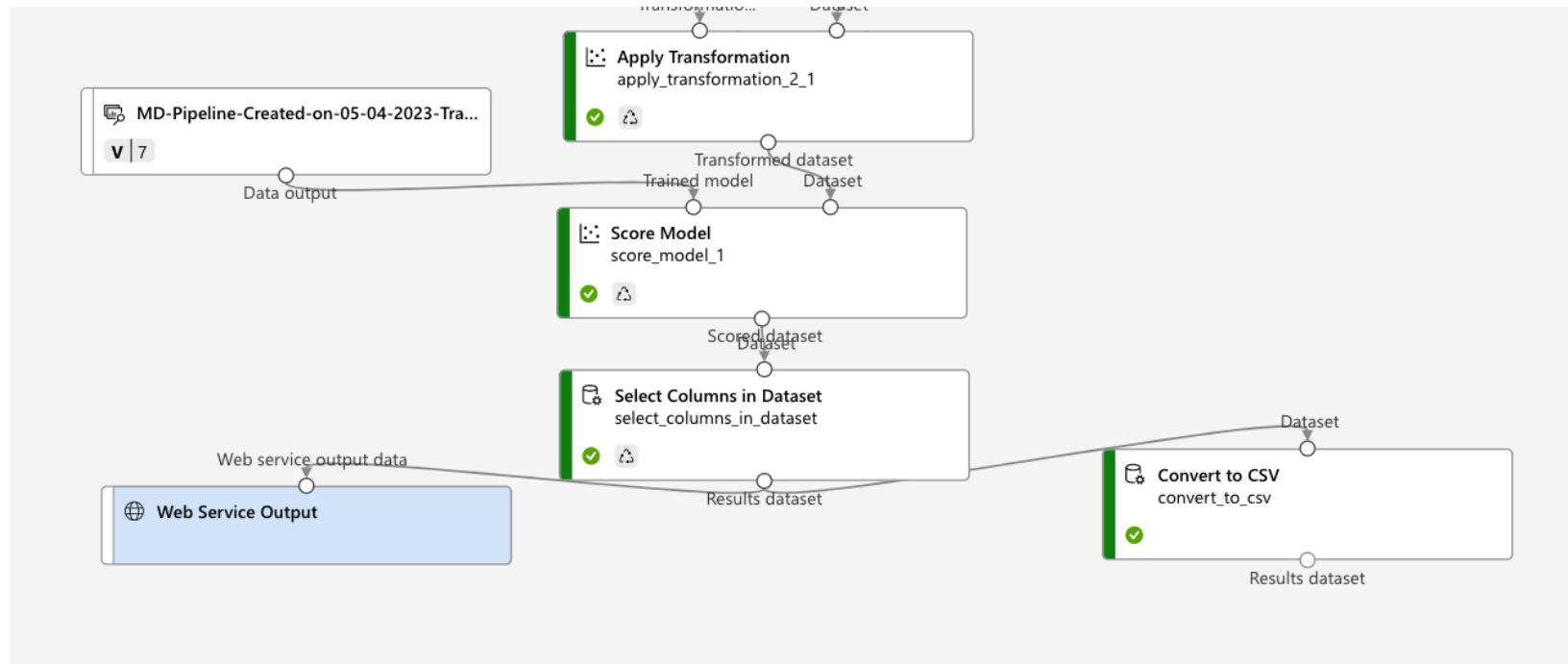


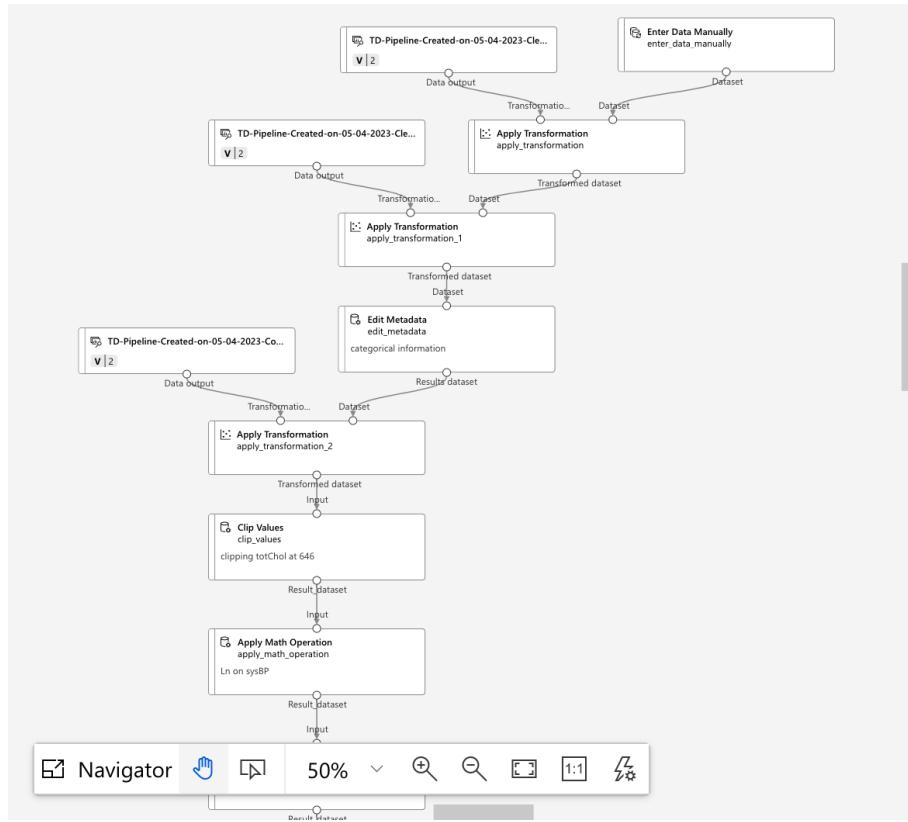
# MODEL DEPLOYMENT

# MODEL DEPLOYMENT









## Enter Data Manually

Data format i \*

CSV

Has header i \*

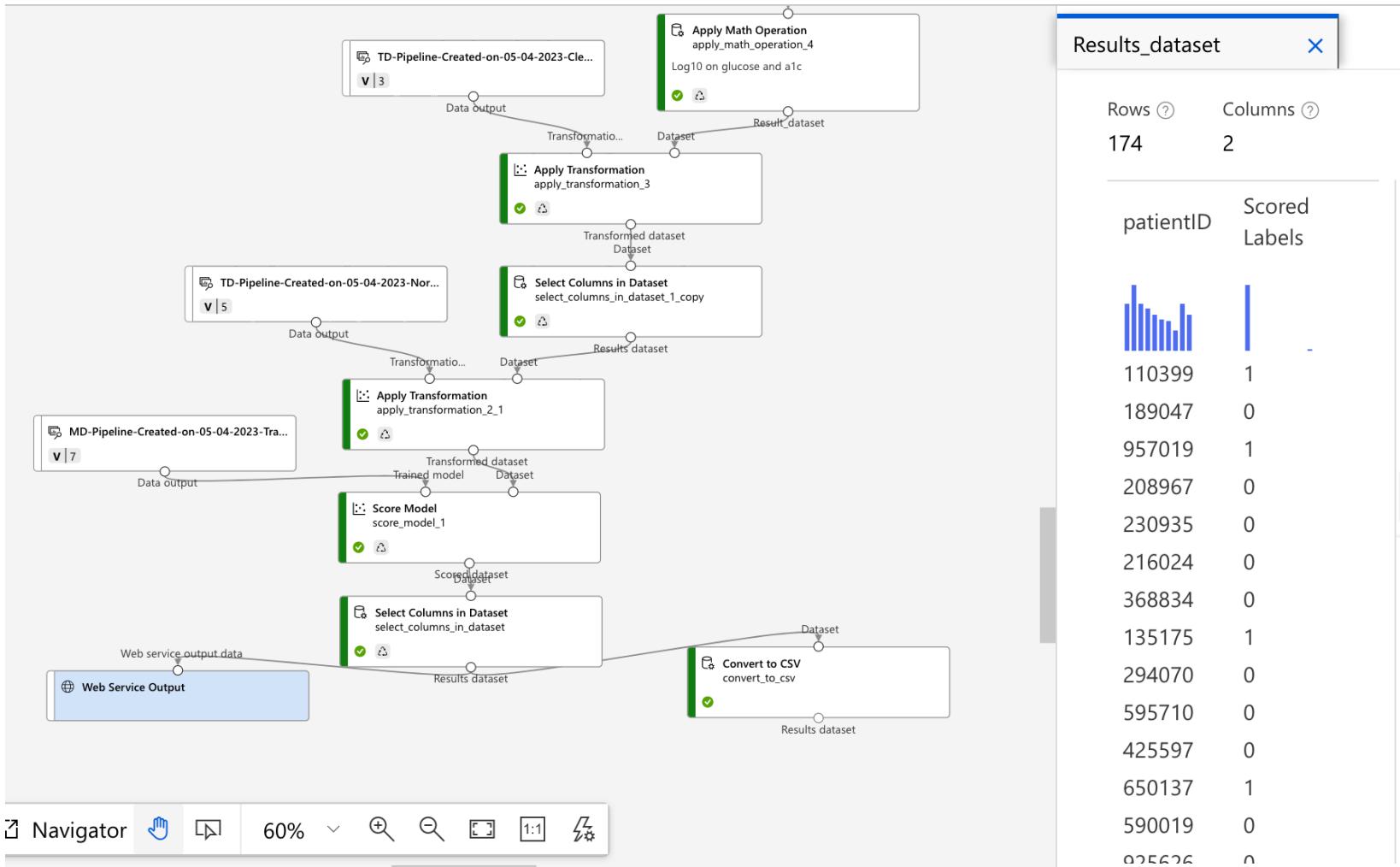
True

Data i \*

```

1 patientID,male,age,education,currentSmoker,cigsPer
2 110399,1,48,3,1,10,0,0,1,0,232,138,90,22.37,64,72,
3 189047,1,41,2,0,,0,0,0,0,195,139,88,26.88,85,65,3,
4 957019,1,54,1,1,20,0,0,1,0,214,147,74,24.71,96,87,
5 208967,1,37,2,0,,0,0,1,0,225,124.5,92.5,38.53,95,8
6 230935,0,63,1,1,3,0,0,1,0,267,156.5,92.5,27.1,60,7
7 216024,1,57,1,0,,0,0,0,0,220,136,84,26.84,75,64,3,
8 368834,0,56,1,0,,0,0,1,0,296,180,90,23.72,75,120,6
9 135175,0,48,1,0,,0,0,1,0,265,145,77,24.23,74,64,3,
10 294070,1,66,3,0,,0,0,0,0,288,109,71,29.29,80,80,4,
11 595710,1,38,4,0,,0,0,0,0,235,118,77,25.87,60,82,4,
12 425597,1,53,1,1,30,0,0,0,0,244,106,67.5,21.84,88,6
13 650137,0,59,1,1,1,0,0,1,0,259,141,86,25.97,70,86,4
14 590019,0,38,3,1,3,1,0,1,0,,125,80,22.79,98,,,13340
15 925626,0,36,3,1,20,0,0,0,0,159,121.5,73,20.41,72,7
16 276518,1,41,4,1,20,1,0,1,0,244,139,86,30.77,60,67,
17 342284,0,37,1,0,,0,0,0,0,300,112,60,23.67,81,75,3.
--
```

# TEST DATA SCORING



Results\_dataset

Rows ② 174  
Columns ② 2

patientID Scored Labels

patientID	Scored Labels
110399	1
189047	0
957019	1
208967	0
230935	0
216024	0
368834	0
135175	1
294070	0
595710	0
425597	0
650137	1
590019	0
025626	0