



FINAL PROJECT

ISE 543, SPRING 2023

XIAOYI WANG
2403234885



- › For this exam, you are to create a predictive model in Azure ML Studio for the attached dataset and turn in a report as specified in the following pages. You should use whichever data preparation, modeling, and model assessment techniques that were covered in this portion of the class that you believe result in the best model.
- › You will be performing an Exploratory Data Analysis, Model Development and Training, and Model Deployment activities and preparing a report in PowerPoint form
- › See the sample report that is part of this assignment for a template and example
- › When you are complete, save this file as a PDF and upload it to Gradescope
- › It is due at 11:00PM on Monday, May 8
- › As a reminder, the work that you submit must be done individually. Unlike the homework assignments, working together is not permitted and the graders will be looking for identical solutions.



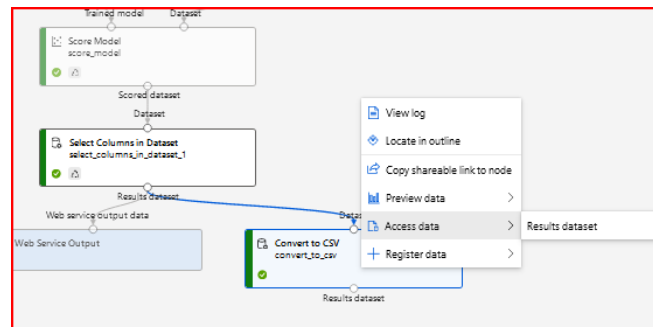
For this exam, you will use Azure ML Studio Designer to build a classification model to predict the likelihood of a patient developing Chronic Heart Disease (CHD) in the coming ten years. The dataset you will be using has been distributed with this exam and consists of the variables on the following page



Variable	Description
Age	age of the participant at the time of examination
Male	gender of the participant (male =1, female = 0)
Education	Educational level of the patient (1 = less than high school, 2 = completed high school or equivalent, 3 = some college, 4= completed college or higher)
Income	Income of the patient
Current Smoker	whether the participant is currently a smoker (yes or no)
Cigarettes per Day	the average number of cigarettes smoked per day by current smokers
BP Meds	whether the participant is taking blood pressure medication (yes or no)
Prevalent Stroke	whether the participant has a history of stroke (yes or no)
Prevalent Hyp	whether the participant has a history of hypertension (yes or no)
Diabetes	whether the participant has diabetes (yes or no)
Total Chol	total cholesterol level in milligrams per deciliter
Sys BP	systolic blood pressure in millimeters of mercury
Dia BP	diastolic blood pressure in millimeters of mercury
BMI	body mass index in kilograms per square meter
Heart Rate	resting heart rate in beats per minute
Glucose	Blood glucose level in milligrams per deciliter
A1c	Hemoglobin A1c (%)
Ten Year CHD	whether the participant developed coronary heart disease (CHD) within 10 years of the examination (yes or no)



- › When complete, do a final run of your inference pipeline by copying test data from the file “Final Project Evaluation Dataset.csv” into the “Enter Data Manually” component (see Sample Final Report for an example)
 - › Do this by opening the file in WordPad or a similar text editor and then copying the data from there
- › Include a “Convert to CSV” component at the end of your inference pipeline (see Sample Final Report for an example)
- › Download the CSV file containing your labels (predictions) by right-clicking the Convert to CSV component, selecting Access data / Results Dataset





- › This will take you to an Azure data folder that will contain a file called data.csv. Download this file (right-click, “Download”) to your laptop and rename it to be “Final scored dataset – xxxxxx.csv” where xxxxxx is your student ID number
- › Upload this file to this Google folder:
 - » https://drive.google.com/drive/folders/1rlrbiYMkQHoWIA3Km_moODAnsFHlksQS?usp=sharing
- › The TAs will use this file to calculate the AUC for your model (by comparing it to the actual values of the response variable)



Please follow the provided template/example and structure your final report into the following three sections:

- › Exploratory Data Analysis
- › Model Development
- › Model Deployment



Report contents: 80%

- › Attribute summary
- › Data cleansing – summary of decisions made
- › Data cleansing pipeline (portion of your overall pipeline)
- › Univariate analysis
- › Bivariate analysis (each variable vs the response variable)
- › Feature section/engineering decisions
- › Model pipeline screenshot
- › Model evaluation results screenshot
- › Inference pipeline screenshot
- › REST Endpoint URL and authentication key (in PPT and in Google drive spreadsheet)
- › Screenshot of scored test dataset

Model performance: 20%

- › Based on TAs calling your endpoint with test data



Response Variable

- › Ten Year CHD

Categories

- › PatientID (UID)
- › Male
- › Education
- › Current Smoker
- › BP Meds
- › Prevalent Stroke
- › Pravalent Hyp
- › Diabetes

Measures

- › Age
- › Income
- › Cigarettes per Day
- › Total Chol
- › Sys BP
- › Dia BP
- › BMI
- › Heart Rate
- › Glucose
- › A1c



Rows 19 Columns 23

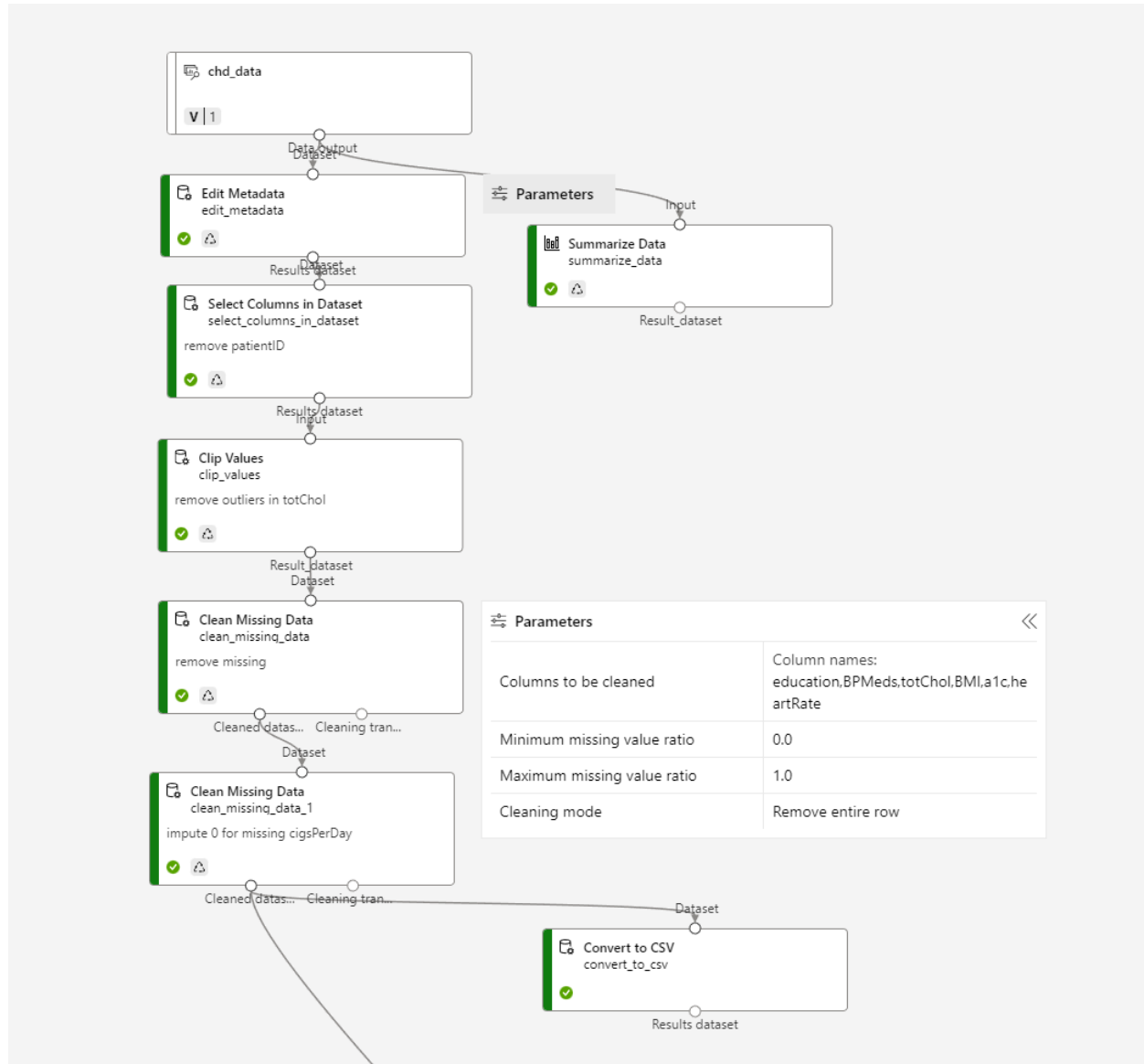
Feature	Count	Unique Value Count	Missing Value Count	Min	Max	Mean	Mean Deviation	1st quantile	Median	3rd quantile	Mode	Range	Sample Variance	Std Dev
patientID	3816	3816	0	100002	999826	554019.06499	222606.258181	336251.75	555421.5	772812.25	{100002, 100012, 100130, 100214, 100754, 101056, 101294, 101524, 102052, 102439, 102712, 102873, 102935, 103840, 103955, 103992, 104276, 104985, 104993, 105054,	899824	66073719998.68935	258181
male	3816	2	0	0	1	0.427673	0.489538	0	0	1	0	1	0.244833	0.494838
age	3816	39	0	32	70	49.567348	7.350087	42	49	56	40	38	73.918923	8.598961
education	3723	4	93	1.0	4.0	1.974483	0.812462	1	2	3	1.0	3.0	1.037769	1.019608
currentSmoker	3816	2	0	0	1	0.489518	0.49978	0	0	1	0	1	0.249956	0.49978
cigsPerDay	1841	31	1975	1.0	70.0	18.500272	8.12733	10	20	20	20.0	69.0	119.365353	10.924838
BPMeds	3771	2	45	0.0	1.0	0.02917	0.056638	0	0	0	0.0	1.0	0.028327	0.168327
prevalentStroke	3816	2	0	0	1	0.006027	0.011982	0	0	0	0	1	0.005992	0.077509
prevalentHyp	3816	2	0	0	1	0.306604	0.425196	0	0	1	0	1	0.212654	0.461265
diabetes	3816	2	0	0	1	0.024895	0.048551	0	0	0	0	1	0.024282	0.155623
totChol	3769	246	47	107.0	9280.0	240.852746	39.378735	205	234	263	240.0	9173.0	35695.202567	188.66823
sysBP	3816	232	0	83.5	295.0	132.260089	16.890304	117	128	143.5	130.0	211.5	489.279059	222.60089
diaBP	3816	142	0	50.0	142.5	82.874214	9.158874	75	82	89.5	80.0	92.5	141.791251	11.58874
BMI	3797	1319	19	15.54	56.8	25.814791	3.113207	23.07	25.4	28.04	22.91	41.26	16.807305	4.113207
heartRate	3815	73	1	44.0	143.0	75.775098	9.288791	68	75	82	75.0	99.0	144.885431	12.250119
glucose	3455	134	361	40.0	394.0	81.856151	12.250119	71	78	87	75.0	354.0	555.598004	23.540307
TenYearCHD	3816	2	0	0	1	0.15173	0.257415	0	0	0	0	1	0.128741	0.353366
a1c	3455	3455	361	2.134768768113766	19.917371285750395	4.296312	0.631685	3.738947	4.126325	4.564732	{2.134768768113766, 2.3023530613653507, 2.349277968047848, 2.355203600535006, 2.3861003033933663, 2.38903897942708, 2.4155486171620075, 2.4374983635432623, 2.4946768795983005, 2.51440307446826,	17.78260251763663	1.424309	1.019608
income	3816	3282	0	12000.0	524494.0	20355.886792	7743.282613	13562.5	16055	21395.5	14623.0	512494.0	319178936.087836	17893.6087836

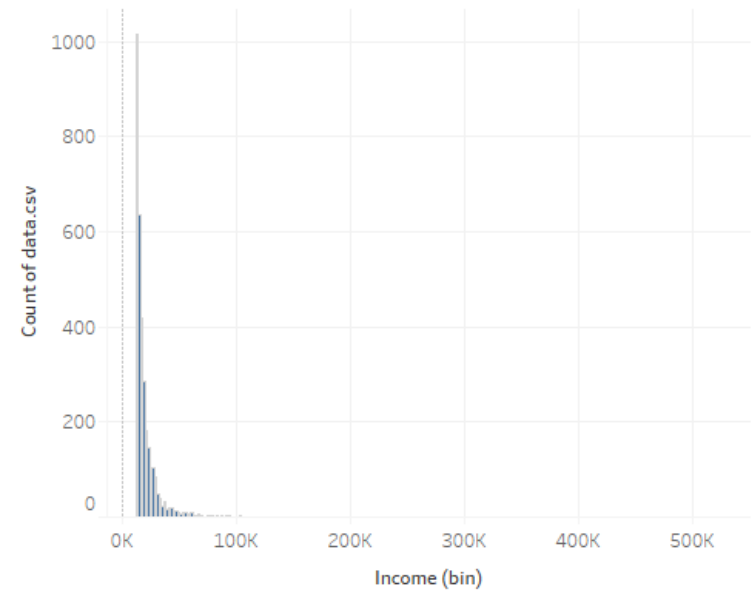
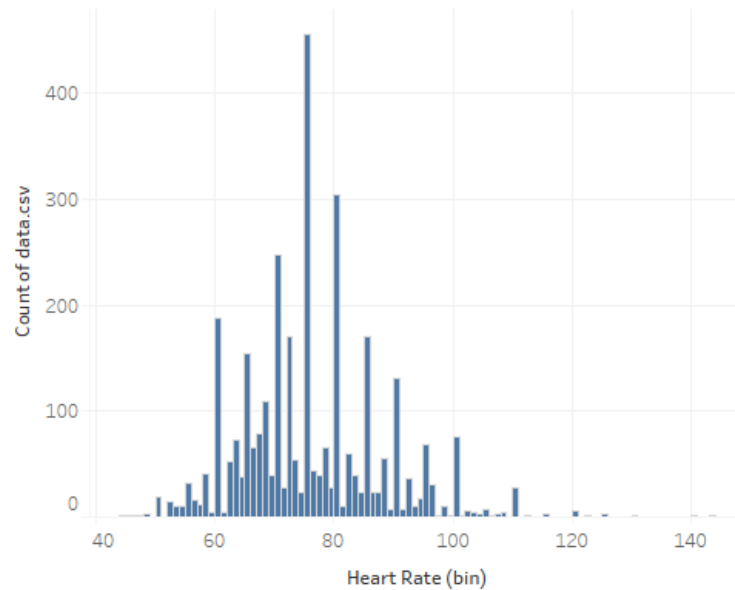
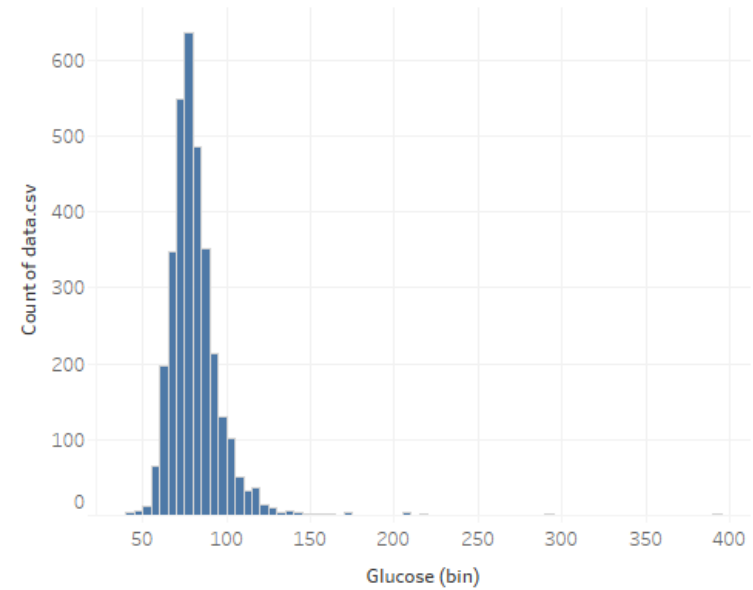
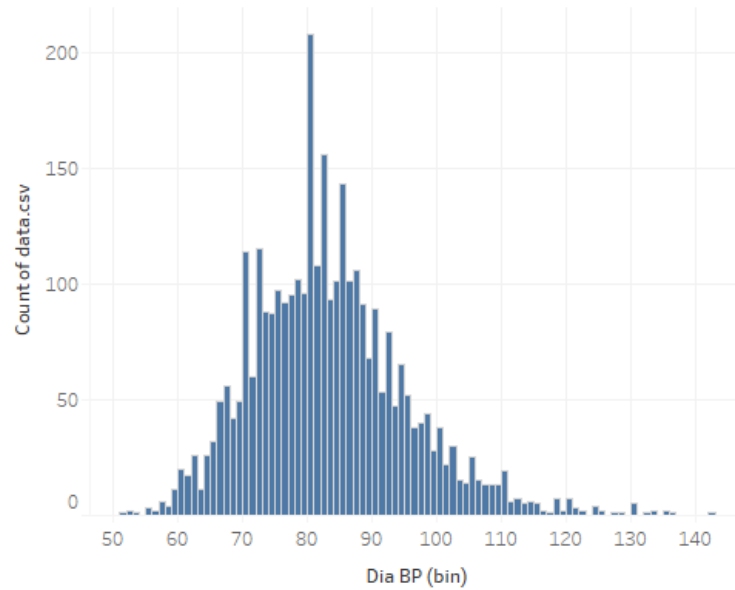
many missing values

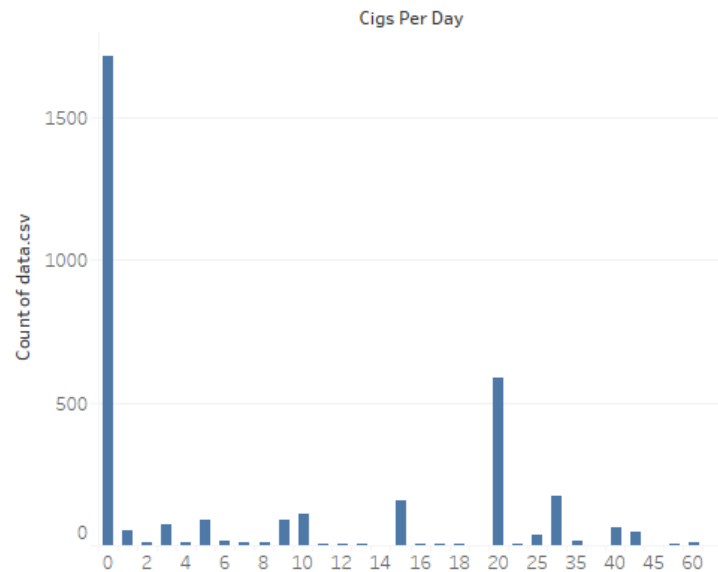
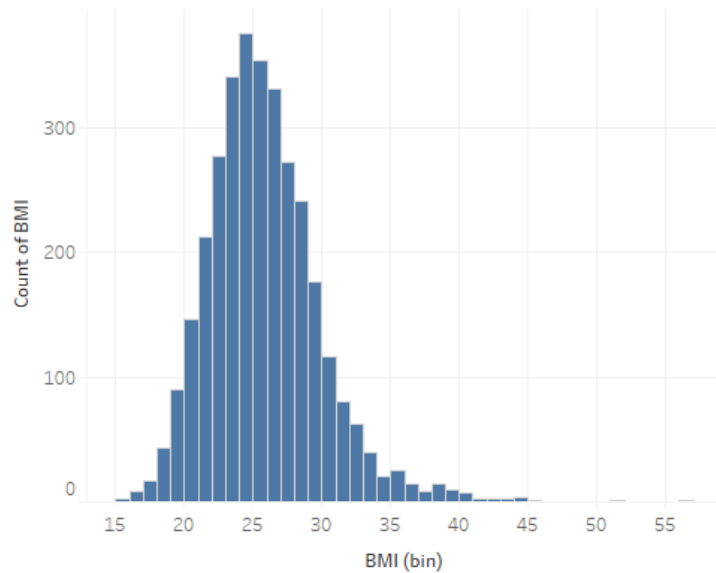
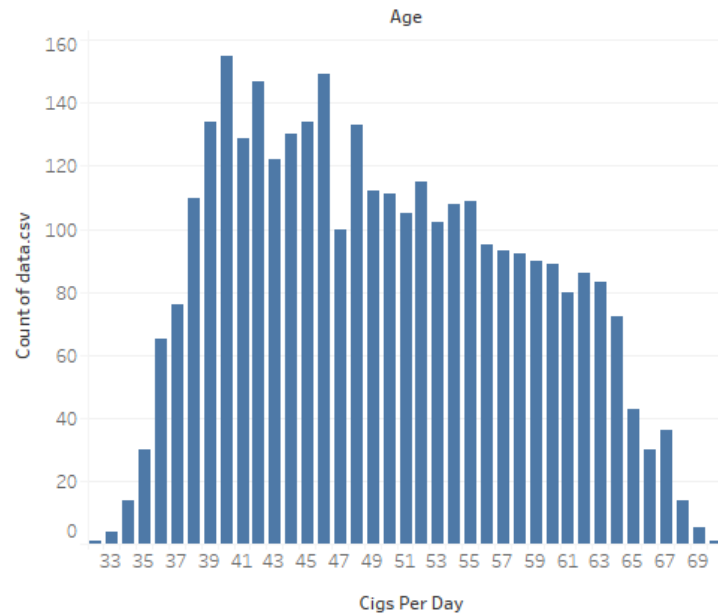
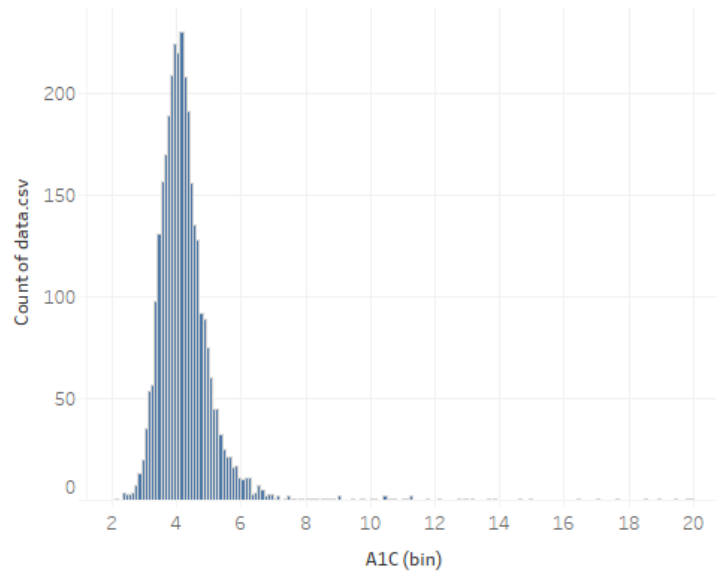
large outliers

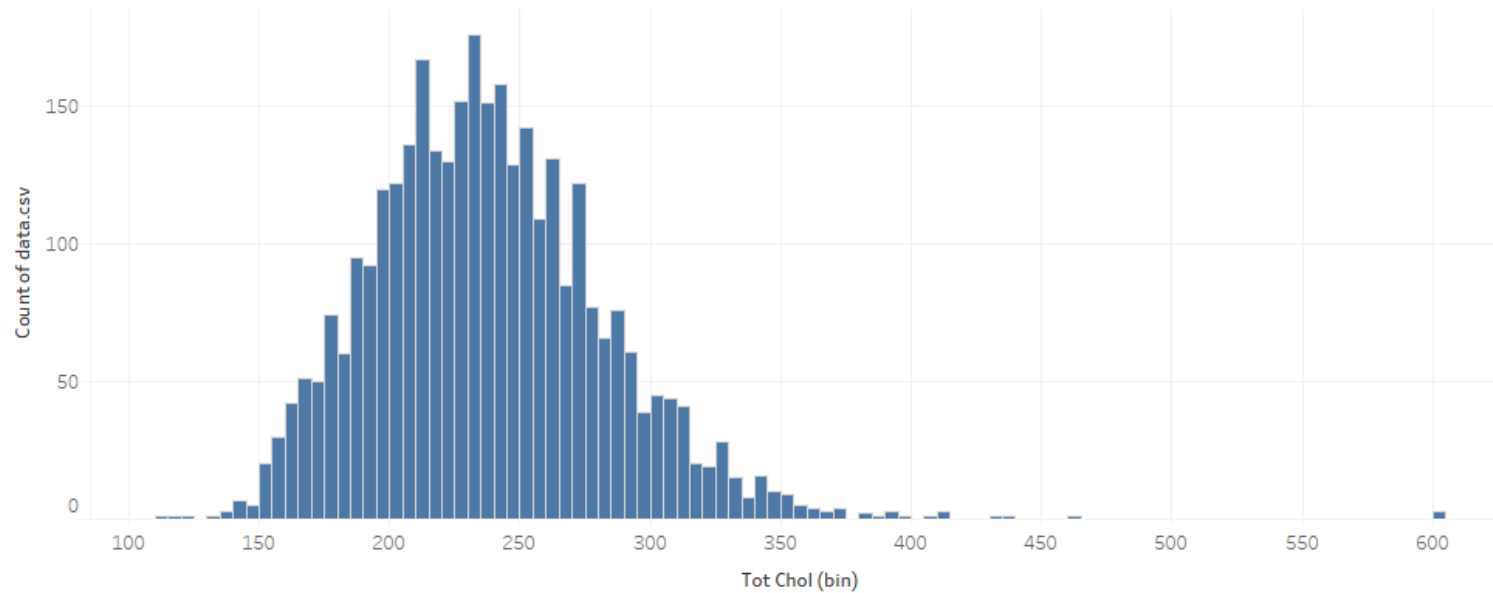
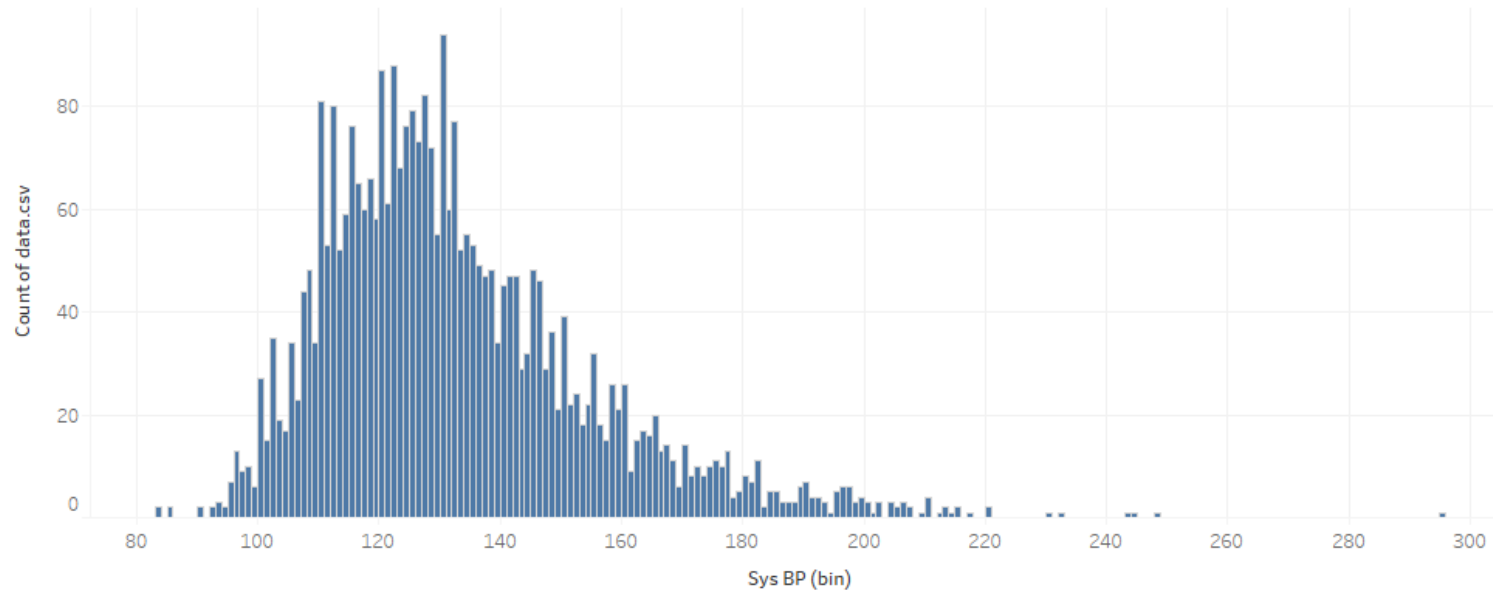


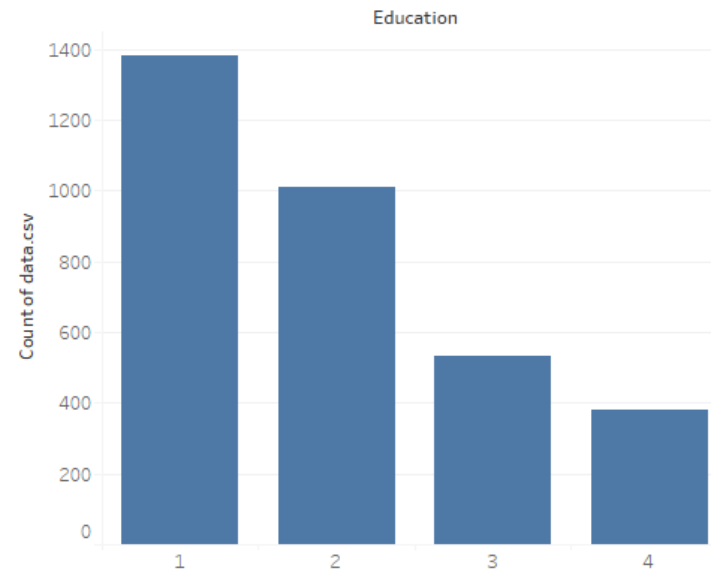
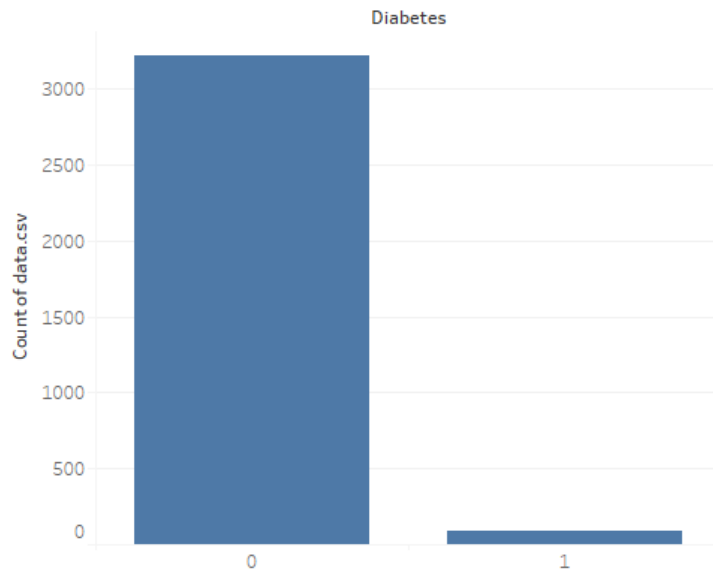
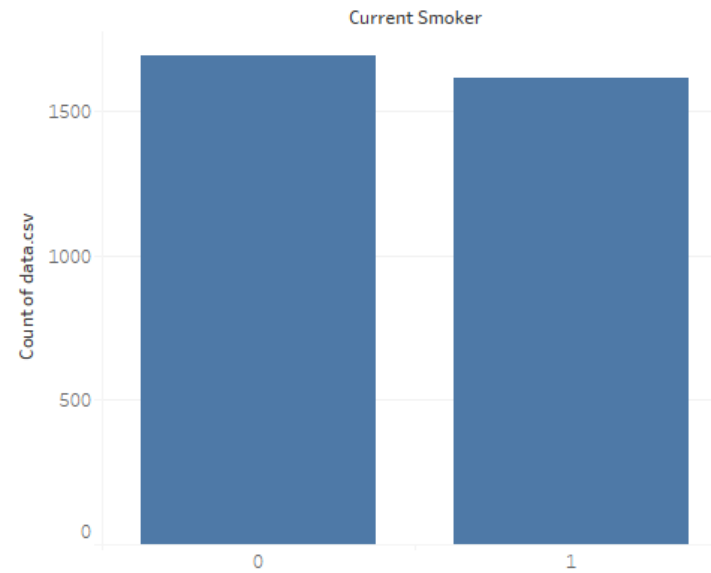
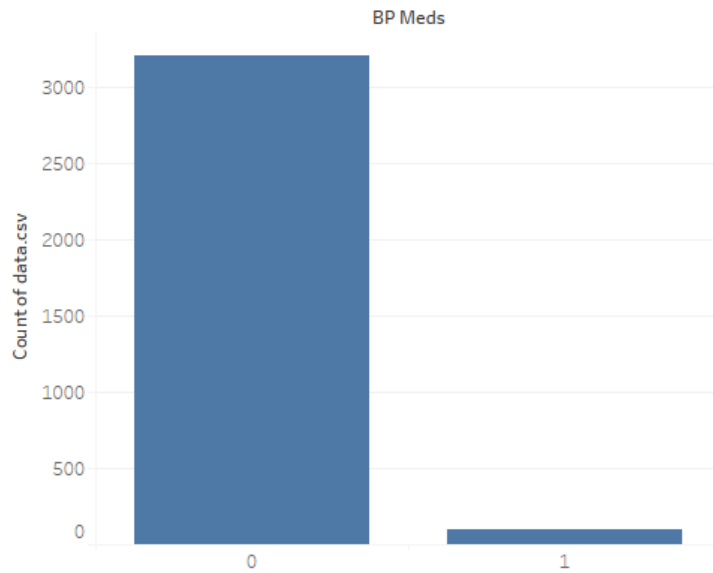
- › change response variable TenYearCHD to categorical
- › outliers in totChol: seem to be erroneous data. clip outliers in totChol at a threshold of 600
- › cigarettes per day: impute missing values with 0 because these non-responses are mostly due to non current smokers
- › remove rows with missing values: education, BPMeds, totChol, BMI, heartRate, glucose, a1c
- › (Note: normalization and transformation are done in feature engineering)

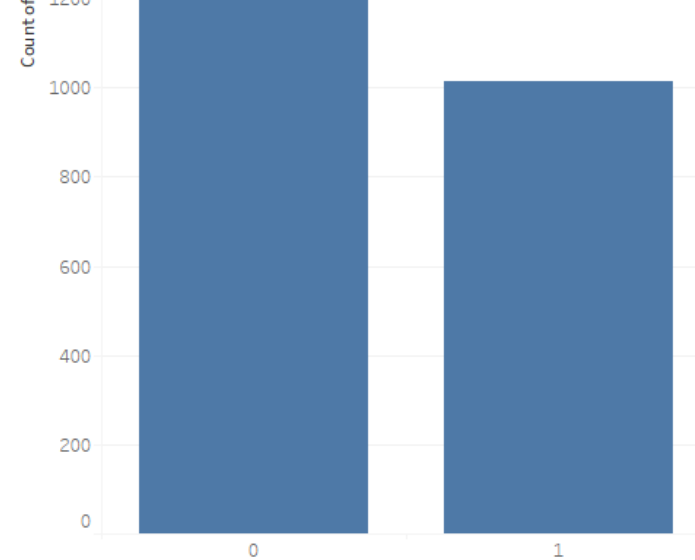
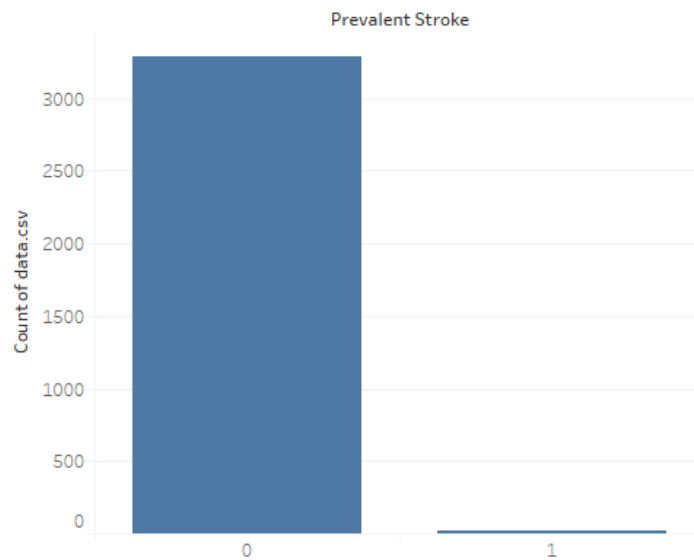
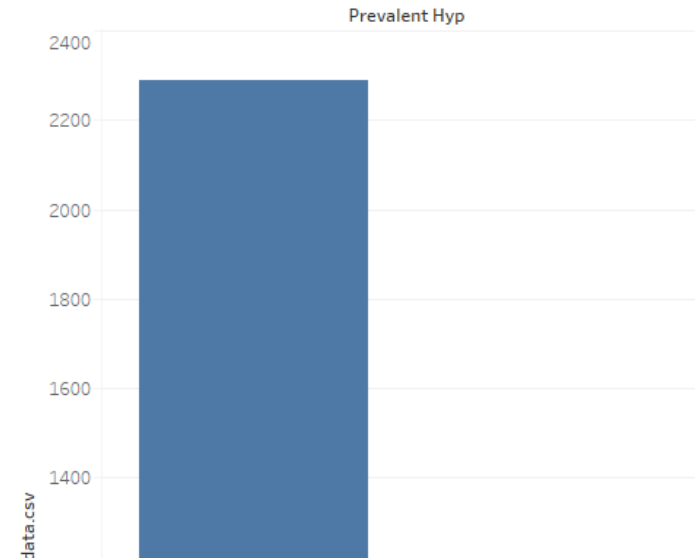
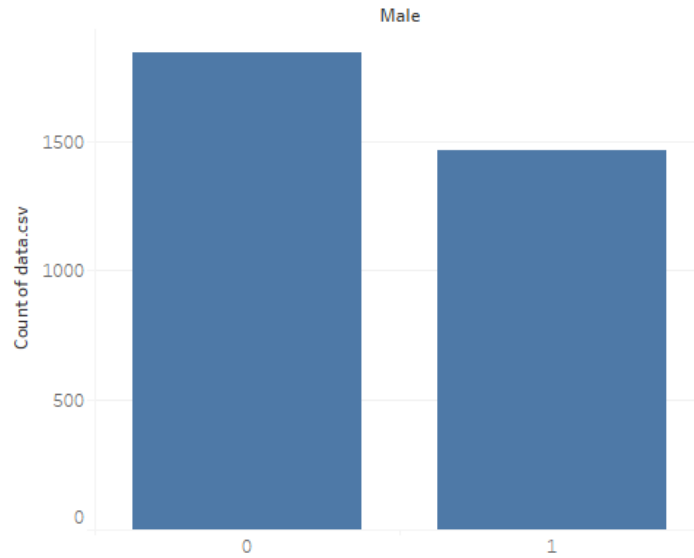


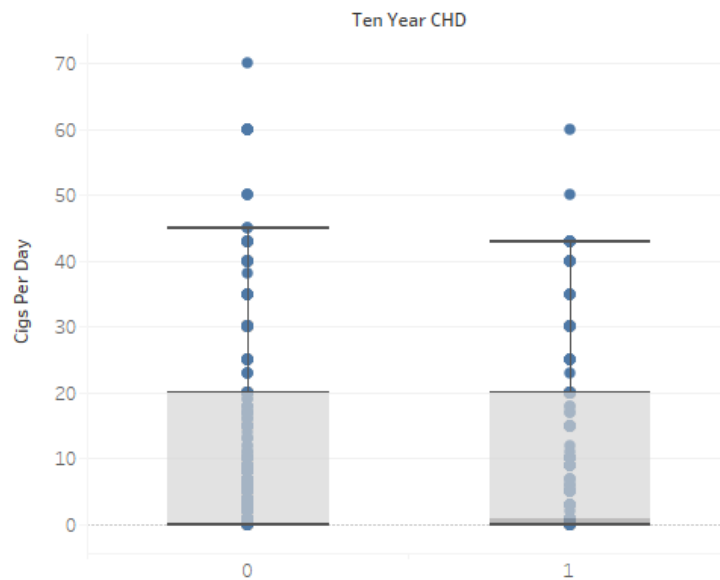
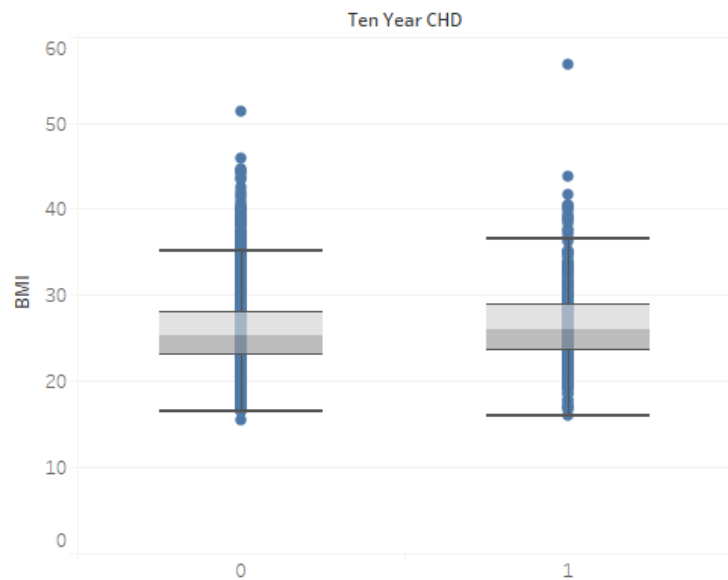
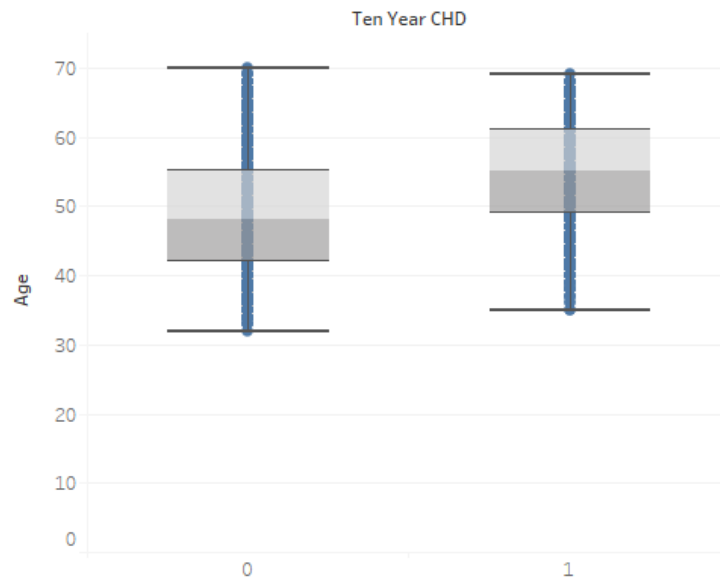
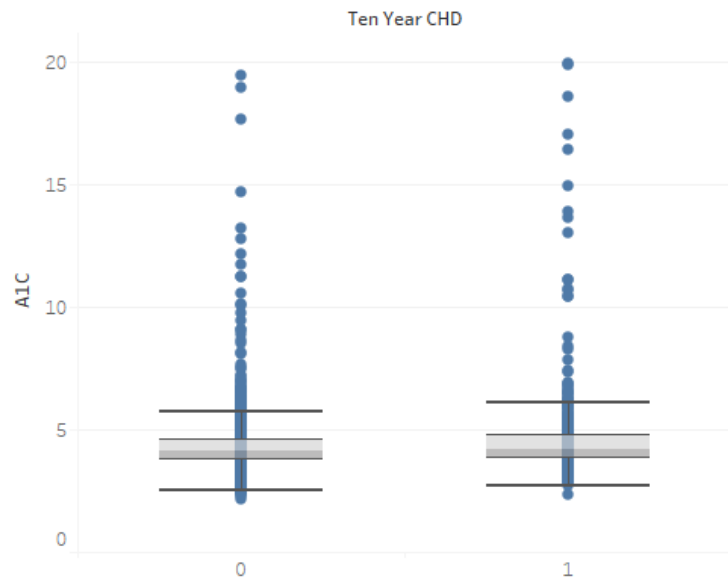


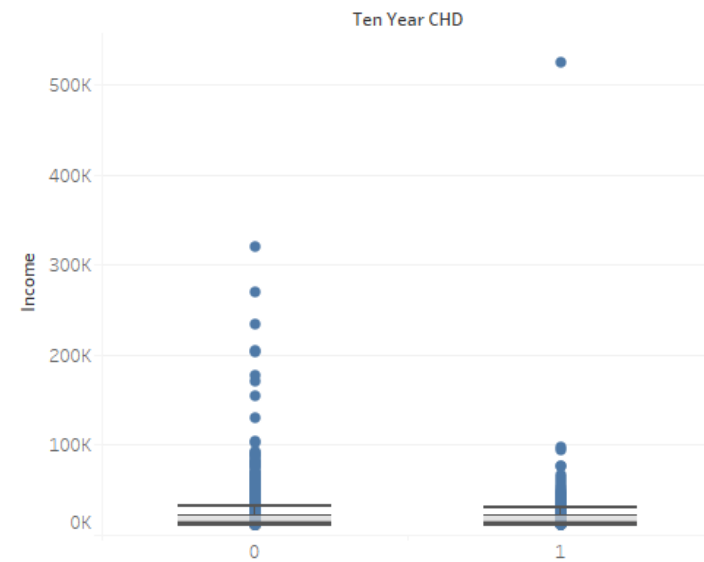
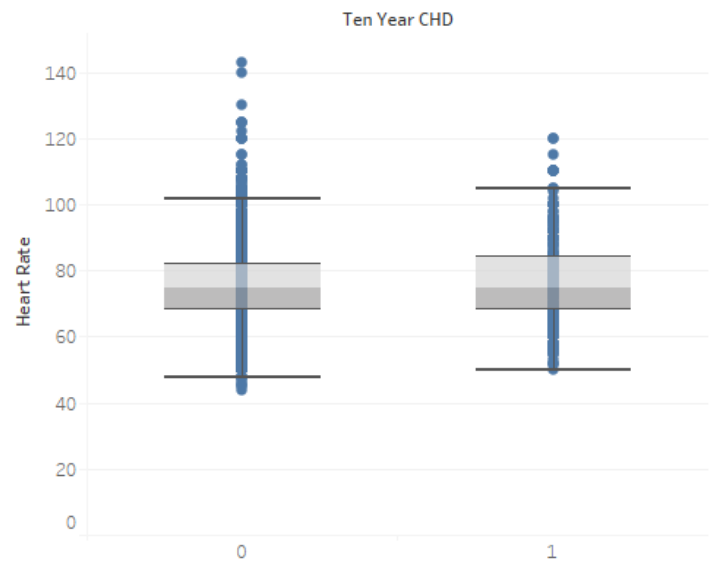
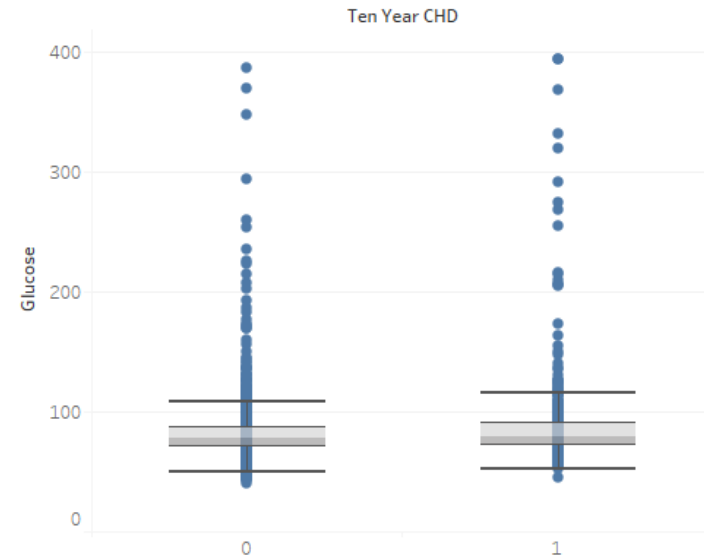
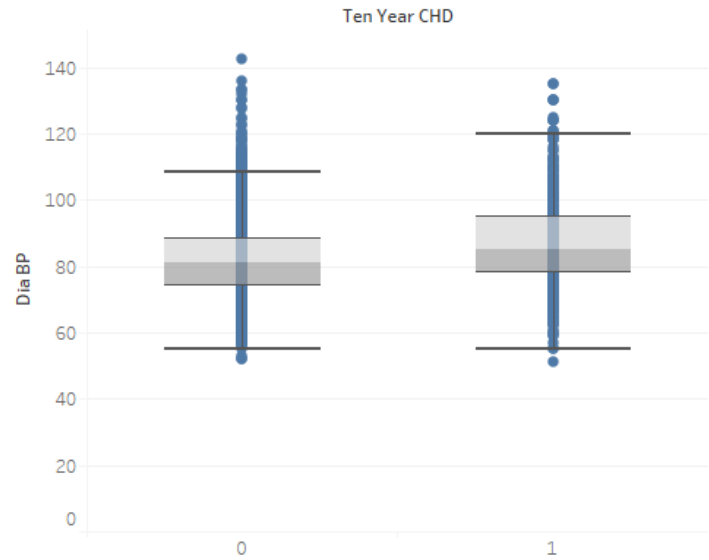


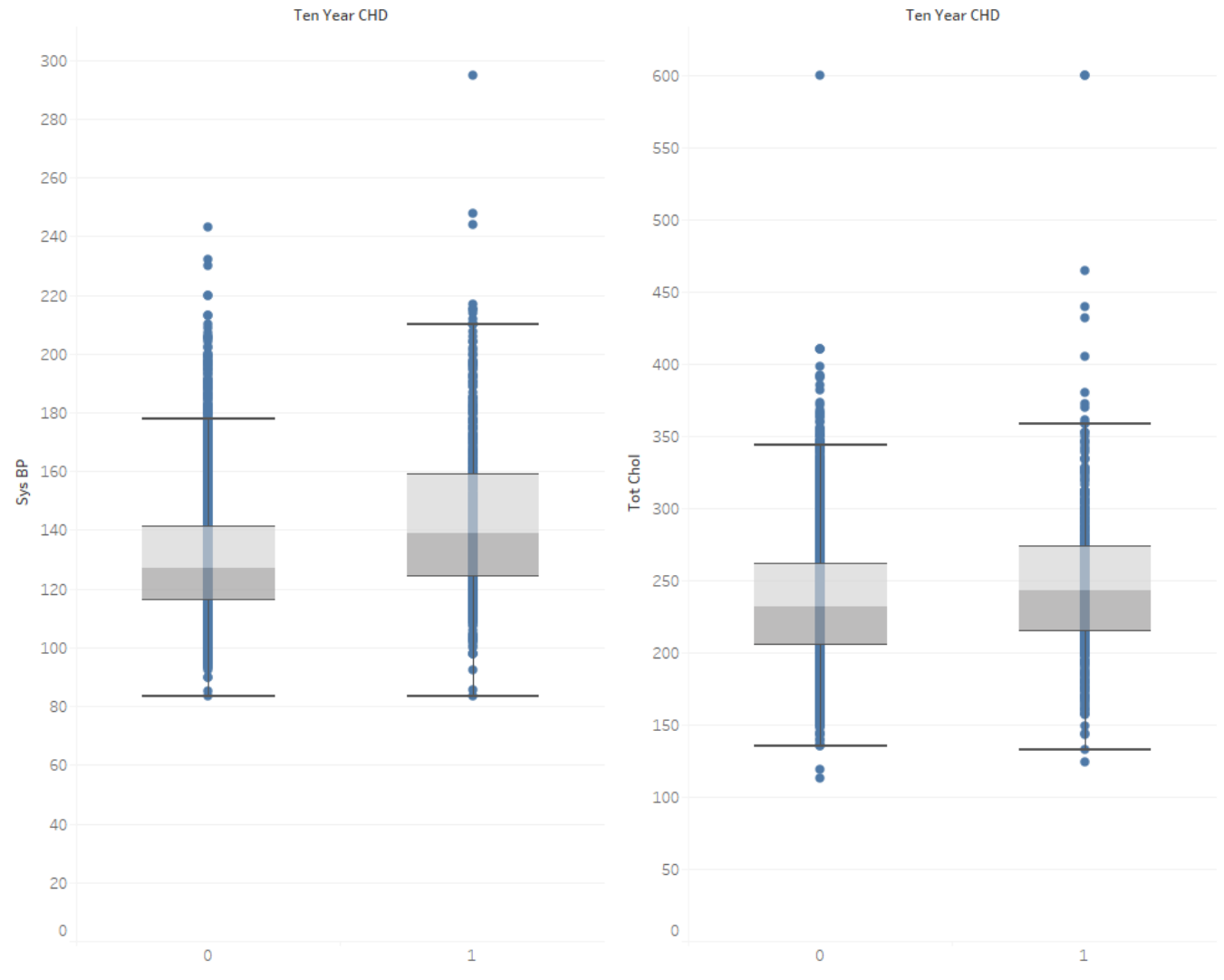


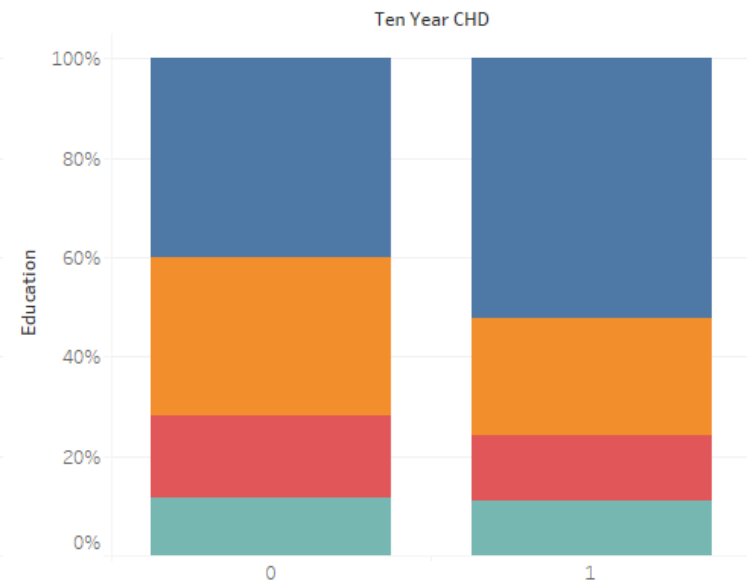
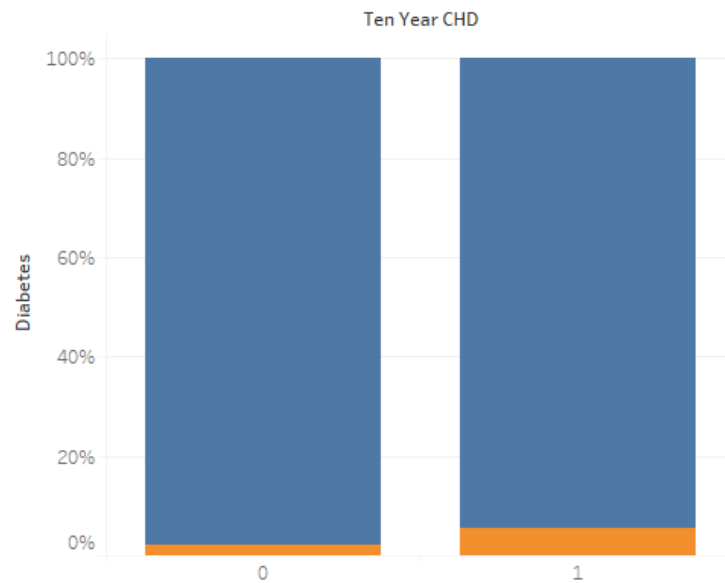
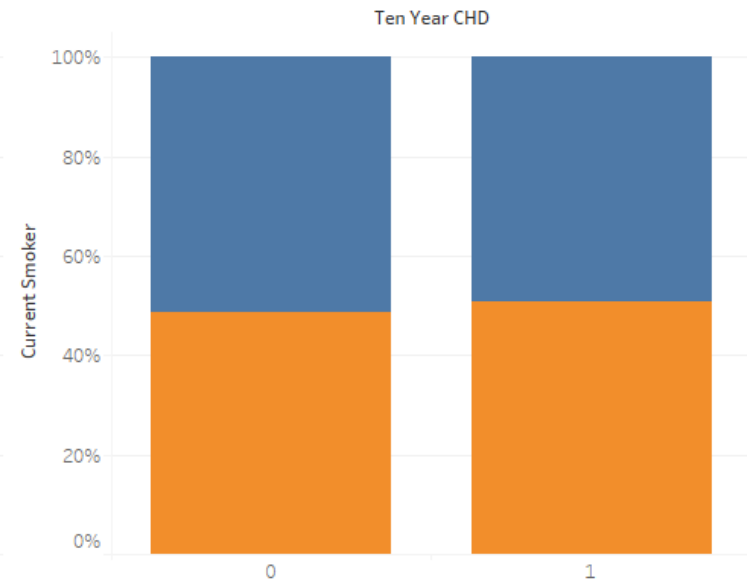
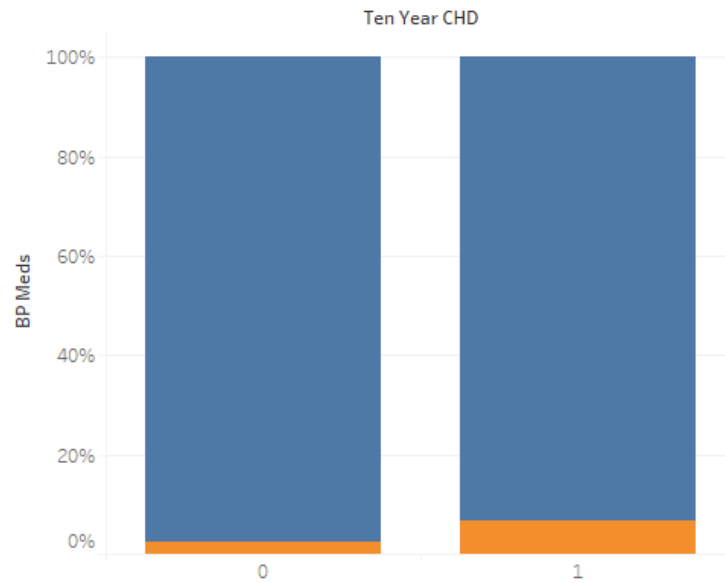


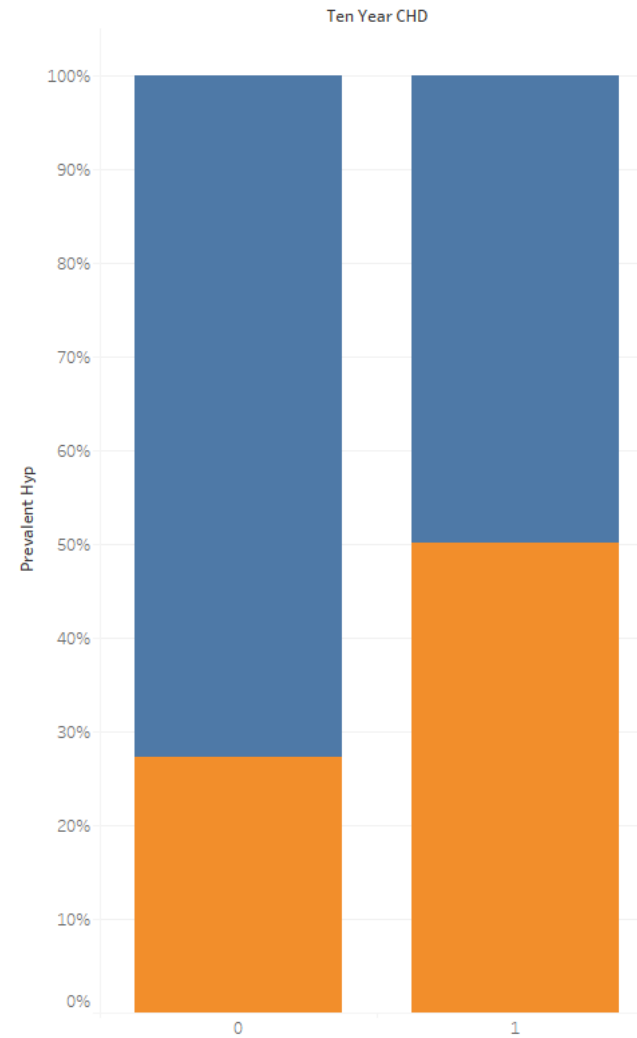
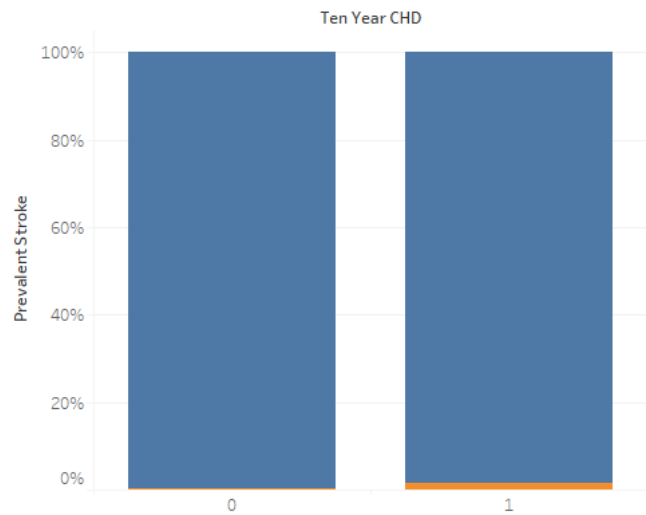
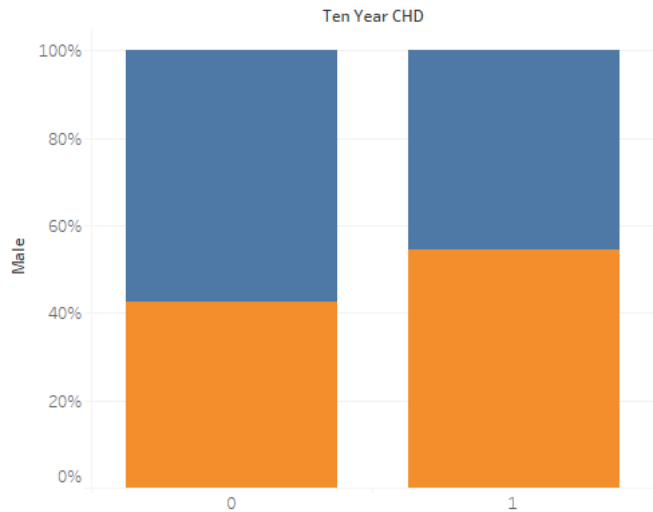






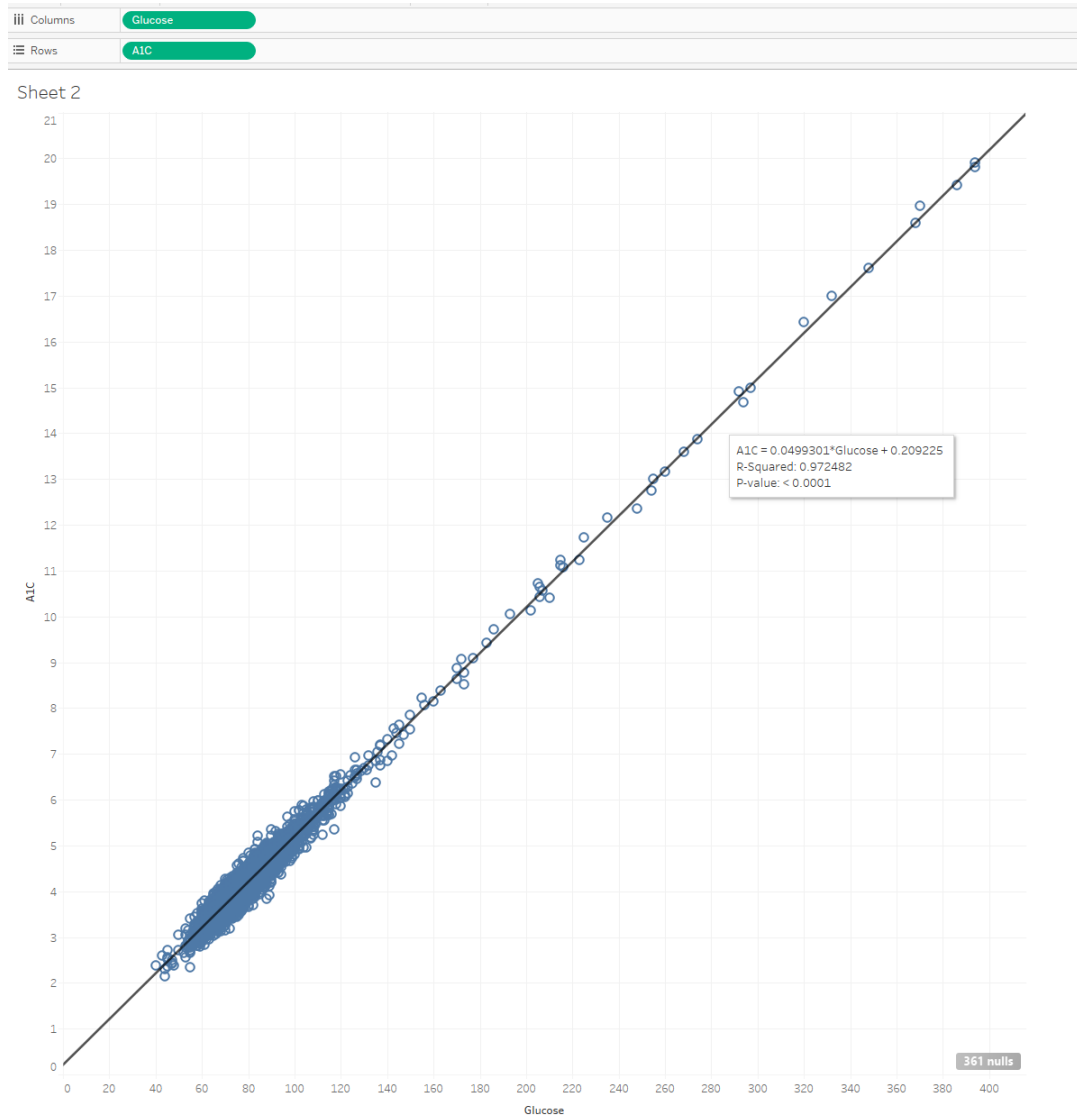








- › glucose vs. A1C: very strong correlation





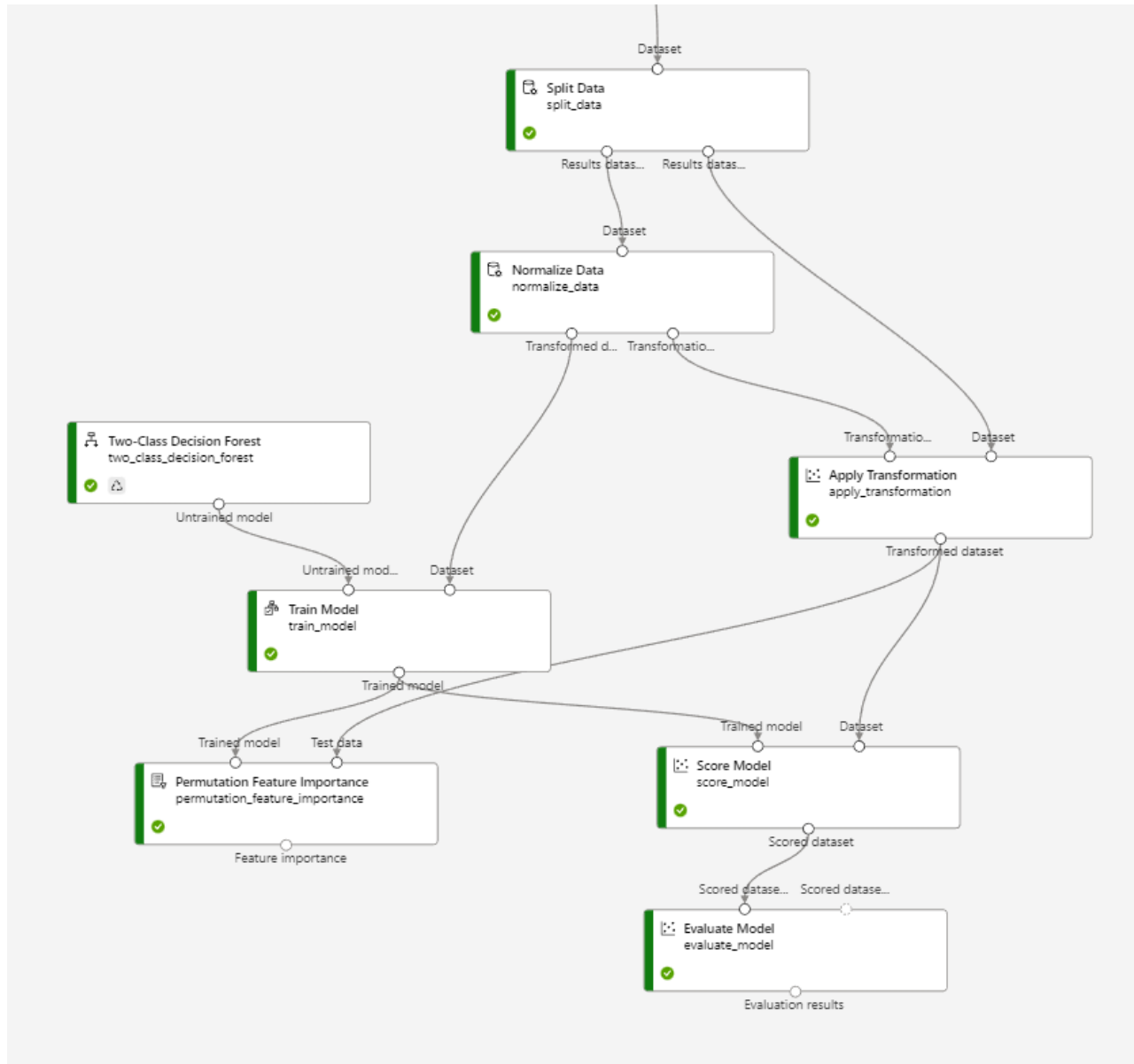
› Dia BP vs. Sys BP: strong correlation





- › Features removed
 - » patient ID – irrelevant for prediction
 - » highly correlated variables: despite the high correlation between two pairs of predictors (glucose vs. a1c, dia bp vs. sys bp), I decide not to remove those two predictors because doing so will significantly decrease the model performance
- › Feature engineering
 - » create square of sysBP as a new predictor
 - » normalize income
 - » log transformation:
 - » very strong skewness (Log10): glucose, a1c
 - » moderate skewness (LnPlus1): diaBP, sysBP, BMI, heartRate

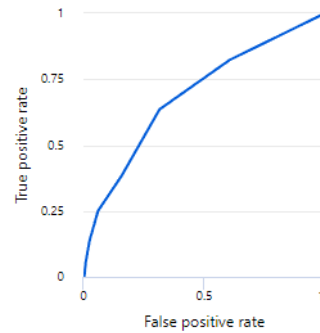




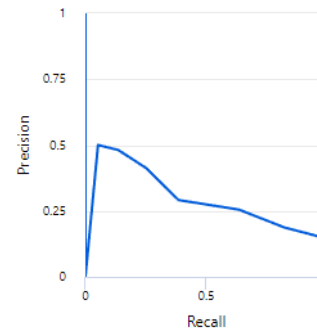


● Scored dataset (left port)

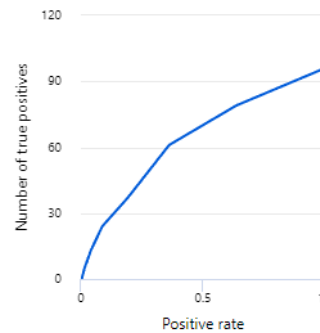
ROC curve



Precision-recall curve



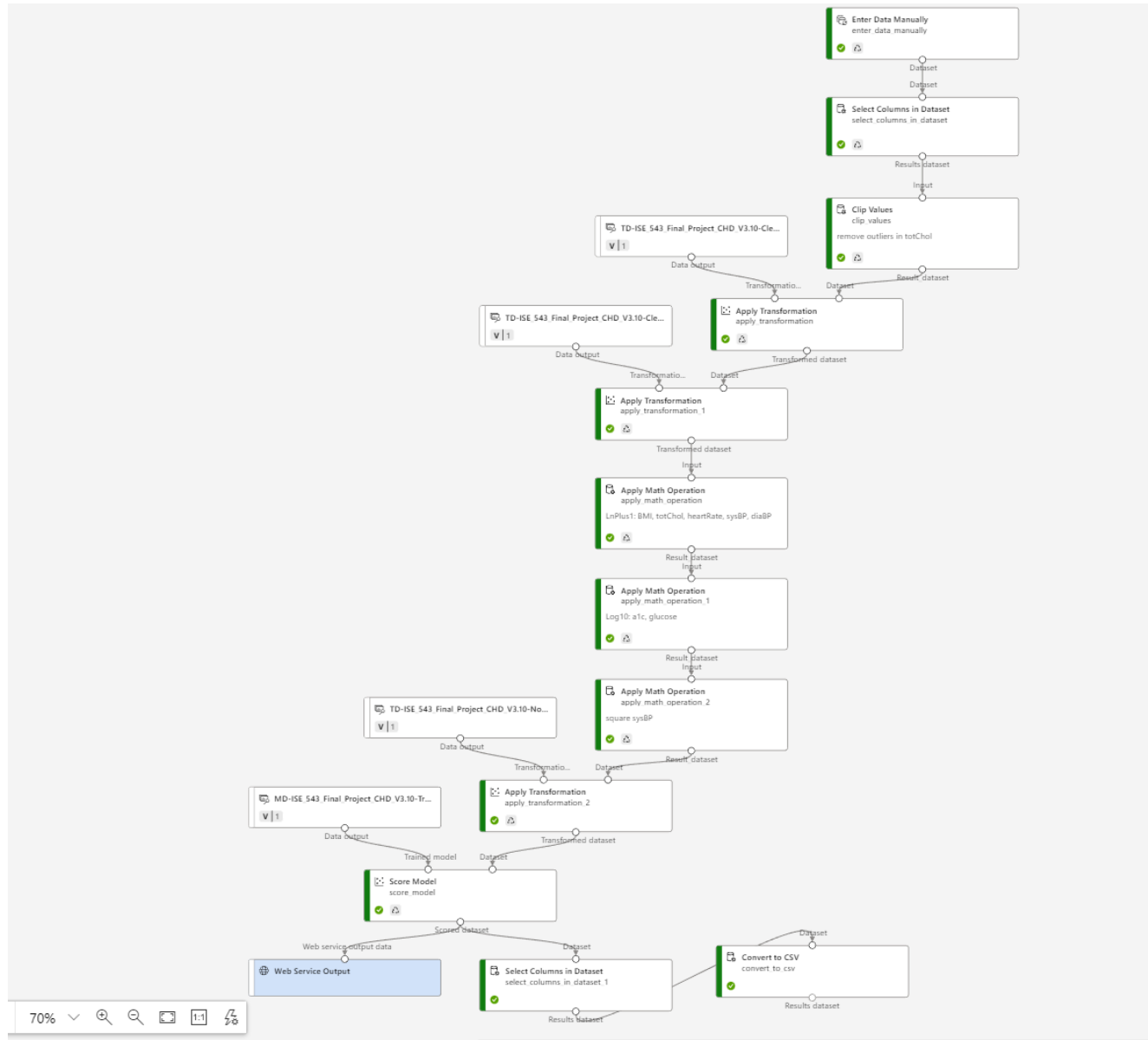
Lift curve



Threshold 0.5

Accuracy 0.84
Precision 0.414
Recall 0.25
F1 Score 0.312
AUC 0.691

	Actual	
	1	0
Predicted	24	34
	72	531





ISE 543 Final Project CHD V3.10-real time inference

Save Settings

Enter Data Manually
enter_data_manually

Dataset

Dataset

Select Columns in Dataset
select_columns_in_dataset
exclude label

Results dataset

Input

Clip Values
clip_values
remove outliers in totChol

Result dataset

Dataset

Transformed dataset

Apply Transformation
apply_transformation

Transformed dataset

Apply Transformation
apply_transformation_1

Transformed dataset

Input

Apply Math Operation
apply_math_operation
LnPlus1: BMI, totChol, heartRate, sysBP, diaBP

Result dataset

Input

Apply Math Operation
apply_math_operation_1
Log10: a1c, glucose

Navigator

100%

Enter Data Manually

Data format

CSV

Has header

True

Data

1	patientID,male,age,education,currentSmoker,cigsPerDay,BPMeds,preval
2	110399,1,48,3,1,10,0,0,1,0,232,138,90,22.37,64,72,4.050698516,13026
3	189047,1,41,2,0,,0,0,0,0,195,139,88,26.88,85,65,3.789559338,18858
4	957019,1,54,1,1,20,0,0,1,0,214,147,74,24.71,96,87,4.571277531,15439
5	208967,1,37,2,0,,0,0,1,0,225,124.5,92.5,38.53,95,83,4.24288858,1580
6	230935,0,63,1,1,3,0,0,1,0,267,156.5,92.5,27.1,60,79,4.370722316,197
7	216024,1,57,1,0,,0,0,0,0,220,136,84,26.84,75,64,3.046747617,13403
8	368834,0,56,1,0,,0,0,1,0,296,180,90,23.72,75,120,6.22075543,17803
9	135175,0,48,1,0,,0,0,1,0,265,145,77,24.23,74,64,3.914695094,37822
10	294070,1,66,3,0,,0,0,0,0,288,109,71,29.29,80,80,4.610732928,14508
11	595710,1,38,4,0,,0,0,0,0,235,118,77,25.87,60,82,4.988670391,20519
12	425597,1,53,1,1,30,0,0,0,0,244,106,67.5,21.84,88,65,3.341763577,262
13	650137,0,59,1,1,1,0,0,1,0,259,141,86,25.97,70,86,4.57160478,16288
14	590019,0,38,3,1,3,1,0,1,0,,125,80,22.79,98,,13340
15	925626,0,36,3,1,20,0,0,0,0,159,121.5,73,20.41,72,75,3.859034015,270
16	276518,1,41,4,1,20,1,0,1,0,244,139,86,30.77,60,67,3.54930057,22263
17	342284,0,37,1,0,,0,0,0,0,300,112,60,23.67,81,75,3.77075934,28637
18	469306,0,45,4,1,15,0,0,0,0,224,117,74.5,16.75,68,87,4.632244066,162
19	197764,0,55,1,0,,0,0,0,0,245,144.5,83.5,28.96,72,65,3.381004969,155
20	416488,1,38,2,1,20,0,0,1,0,253,133,92,28.82,80,63,4.424617023,27512
21	208652,0,55,2,1,9,0,0,0,0,248,157,82.5,22.91,89,83,4.32518232,15356
22	562216,0,38,2,0,,0,0,0,0,171,111,68,18.76,90,83,4.479458239,13711
23	115448,1,52,2,1,15,0,0,0,0,240,94,66.5,22.93,70,88,4.82557568,15047
24	224178,0,64,1,0,,0,0,1,0,273,155,86,27.53,100,91,4.847750302,16831
25	271781,1,53,2,0,,0,0,0,0,193,142,89,29.56,70,78,4.265167455,14462
26	887721,1,65,1,1,15,0,0,1,0,219,148,90,29.35,77,97,4.868143542,13917
27	338781,1,46,1,1,30,,0,1,0,253,147,85,30.62,100,75,3.920203843,12875
28	262782,1,61,1,0,,0,0,0,0,224,124,74,21.9,55,75,3.94315094,14924
29	583133,0,46,1,0,,0,0,0,0,219,150,81,25.43,69,93,4.763662967,12061
30	196173,0,45,1,30,0,0,0,0,203,131,85,23.47,94,70,3.632547241,14348
31	779064,0,48,3,0,,0,0,0,0,230,129,84.5,24.73,78,,12181
32	676839,0,54,1,0,,0,0,0,0,219,143.5,89,28.47,73,96,4.912113375,15665
33	141292,0,61,3,0,,0,0,1,0,290,170,98,26.98,80,84,4.101478176,26960
34	881246,0,54,1,0,,0,0,0,0,245,117,76,26.64,65,76,3.97022846,17582

Edit code

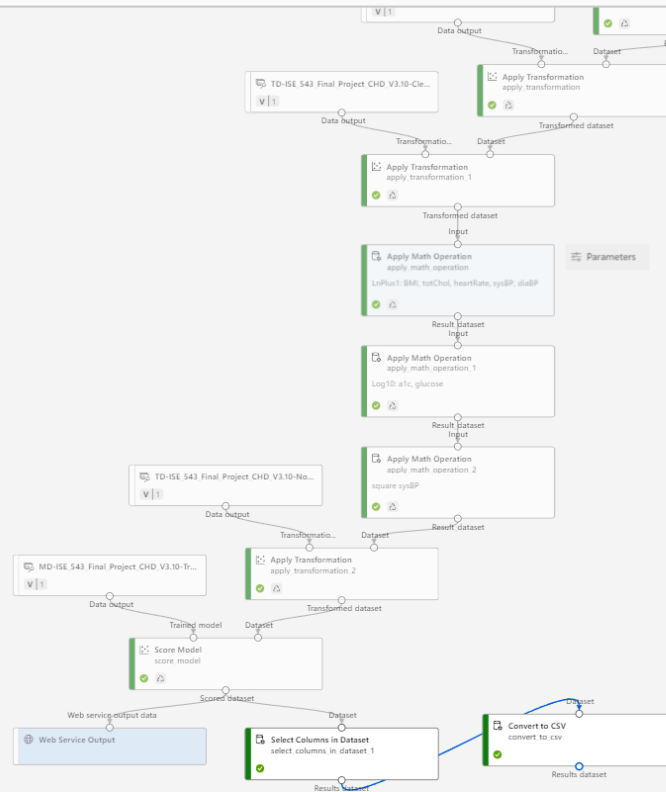
Output settings

Input settings

Run settings



ISE 543 Final Project CHD V3.10-real time inference ✓ Completed



Results_dataset	
Rows	Columns
162	2
patientID	Scored Labels
110399	0
189047	0
957019	0
208967	0
230935	0
216024	0
368834	0
135175	0
294070	0
595710	0
425597	0
650137	0
925626	0
276518	0
342284	0
469306	0
197764	0
416488	0
208652	0
562216	0
115448	0
224178	0
271781	0
887721	0
262782	0
583133	0
676839	0
141292	0
881246	0
982834	0
752601	0
349422	0
530203	0
382977	0
926016	0
757693	0