

Final Report

Sujata Avirneni (savirneni@gatech.edu)

INDEX

Introduction & Problem Statement.....	3
Predict Severe Hypoglycemia SH	3
Dataset Details.....	3
DiabScreening_data.....	3
DiabPhysExam_data.....	3
Diab1.....	3
AdvEvent_data.....	4
Diabetes1.....	4
Diabetes2.....	4
Approach	5
Challenge 1: Predicting Severe hypoglycemia.....	5
Feature importance Diab1 dataset	6
Features importance Diabetes2 dataset	11
Corr Plot	6
Exploratory Visualizations.....	11
Logistic Regression	14
Deep Learning Model.....	
Ensemble model.....	17
Comparing models with K-fole cross validation using accuracy with all 3 datasets...	22
Conclusion	24
References.....	25

Introduction & Problem Statement

Diabetes is a common chronic disease. The characteristic of diabetes is that the blood glucose is higher than the normal level, which is caused by defective insulin secretion, Diabetes can lead to chronic damage and dysfunction of various tissues, especially eyes, kidneys, heart, blood vessels and nerves.

Diabetes can be divided into two categories,

Type 1 diabetes (T1D) and Type 2 diabetes (T2D).

My focus is on Type 1 diabetes (T1D) and the complications associated with Type 1 diabetes (T1D). One of the complications is Severe hypoglycemia (SH) is a major acute complication. In severe cases seizure, increased risk for falls, car accidents, unconsciousness or even death can occur.

Predict Severe hypoglycemia (SH)

I built ML models to address two problem statements: predict diabetes and Severe hypoglycemia (SH) complications associated with Type 1 diabetes (T1D).

- 1) Build prediction models focused on FingStkBG-glucose using logistic regression model, deep learning model and support vector machine to predict Severe hypoglycemia (SH).
- 2) Use ensemble methods to improve accuracy.
- 3) Build classification model using decision tree and random forest.
- 4) Evaluate capability of the model validation using hold-out method and k-fold cross model validation method.
- 5) Remove redundant features using feature selection method such as PCA and minimum redundancy maximum relevance (mRMR) to reduce the dimensionality.

Dataset Details

My data source is public study website on diabetics <https://public.jaeb.org/datasets/diabetes> which contains datasets related to diabetics. I used following datasets.

AdvEvent.csv, DiabPhysExam.csv, DiabScreening.csv

DiabScreening Table

This table contains information about patient demographics, viral signs, diagnosis age and date, most recent and last 12 months history of Severe hypoglycemia (SH), insulin mode and type of device, pre-existing condition and eligibility criteria, patient's current medications etc

DiabPhysExam Table

This table contains PatientID and viral signs mostly related to figure stick blood glucose, Heart rate Systolic and diastolic blood pressure.

AdvEvent Table

This table contains the PatientId, medical conditions, adverse events such as disability death etc.

Table 1: Available Features in the Dataset

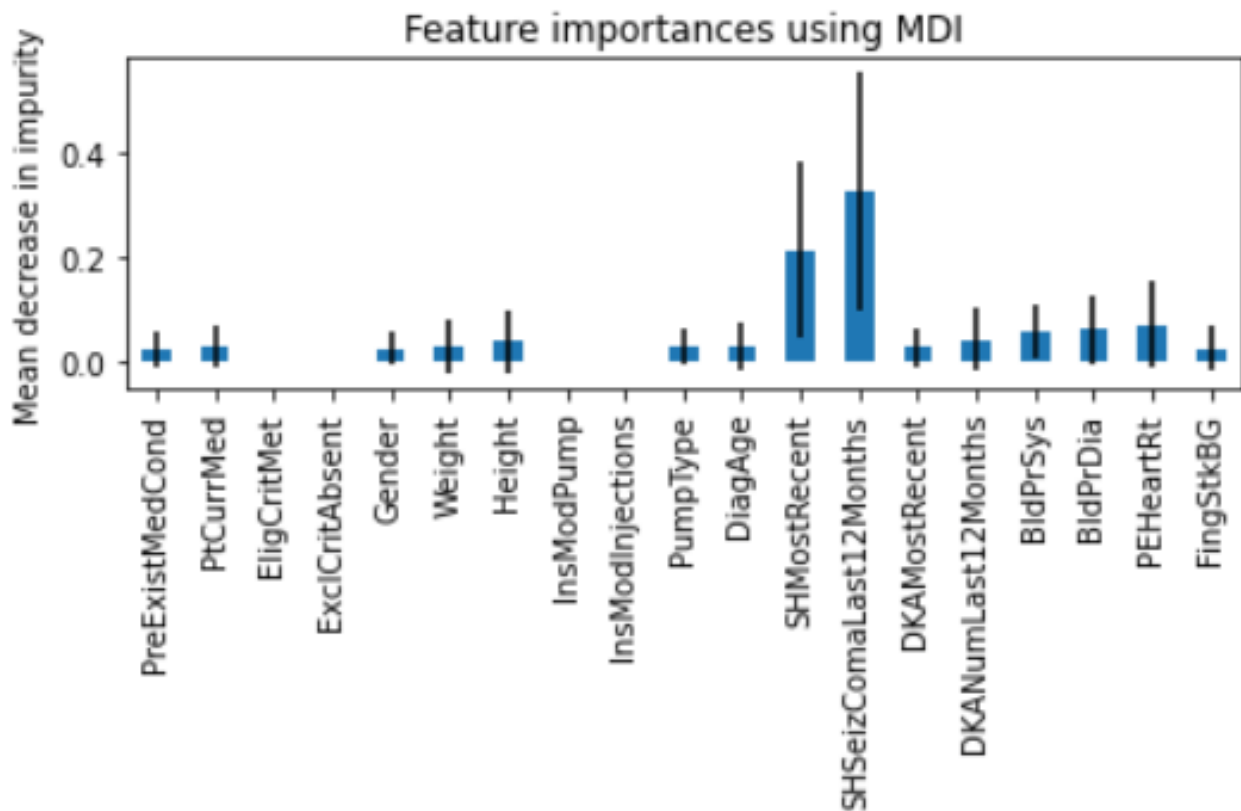
Table	Available predictors
-------	----------------------

DiabScreening_data	['PtID','PreExistMedCond','PtCurrMed','EligCritMet','ExclCritAbsent','Gender','Weight','Height','InsModPump','InsModInjections','PumpType','DiagAge','PEAbnormal','SHMostRecent','SHSeizComa','SHSeizComaLast12Months','DKAMostRecent','DKANumLast12Months']
DiabPhysExam_data	'PtID','BldPrSys','BldPrDia','PEHeartRt','FingStkBG'
Diab1	['PreExistMedCond','PtCurrMed','EligCritMet','ExclCritAbsent','Gender','Weight','Height','InsModPump','InsModInjections','PumpType','DiagAge','SHMostRecent','SHSeizComa','SHSeizComaLast12Months','DKAMostRecent','DKANumLast12Months','BldPrSys','BldPrDia','PEHeartRt','FingStkBG']
AdvEvent_data	[[['PtID','MedicalCondition','MedicalConditionMM','AEPrEnroll','AENotedStdyVisExam','AEIntensity','AERelStdyTrt','AERelStdyProc','AEEffectTrt','AESerious','AETrt','AESurg','AEOthTrt','AEOutcome','AEDeath','AEConAnomaly','AEOthRelHx','AELifeThreat','AEHosp','AEDisability','AEOther','Weight','WeightNotAvail','AERelLabData','AEMedProd','MMAERelStdyTrt','MMAESerious','MMHospDiscRptObtained','AERelStdyDrugDeviceWhich','MMAERelStdyTrtHighLvl']]]
Final dataset after feature importance Diabetes1	['PreExistMedCond','PtCurrMed','EligCritMet','ExclCritAbsent','Gender','Height','InsModPump','InsModInjections','PumpType','DiagAge','SHMostRecent','SHSeizComa','SHSeizComaLast12Months','DKAMostRecent','DKANumLast12Months','BldPrSys','BldPrDia','PEHeartRt','FingStkBG','MedicalCondition','MedicalConditionMM','AEPrEnroll','AENotedStdyVisExam','AEIntensity','AERelStdyTrt','AERelStdyProc','AEEffectTrt','AESerious','AETrt','AESurg','AEOthTrt','AEOutcome','AEConAnomaly','AELifeThreat','AEHosp','AEDisability','AEOther','Weight_y','WeightNotAvail','AERelLabData','AEOthRelHx','AEMedProd','MMAERelStdyTrt','MMAESerious','MMHospDiscRptObtained','AERelStdyDrugDeviceWhich','MMAERelStdyTrtHighLvl']
Diabetes2	'MMAESerious','AEMedProd','AEOthRelHx','AERelLabData','Weight_y','AEOutcome','AEOthTrt','AETrt','AESerious','AEEffectTrt','AEIntensity','AEPrEnroll','MedicalCondition','MedicalConditionMM','FingStkBG','PEHeartRt','BldPrSys','BldPrDia','DKAMostRecent','DKANumLast12Months','SHMostRecent','SHSeizComa','SHSeizComaLast12Months','PumpType','DiagAge','Gender','Height','PtCurrMed','PreExistMedCond'

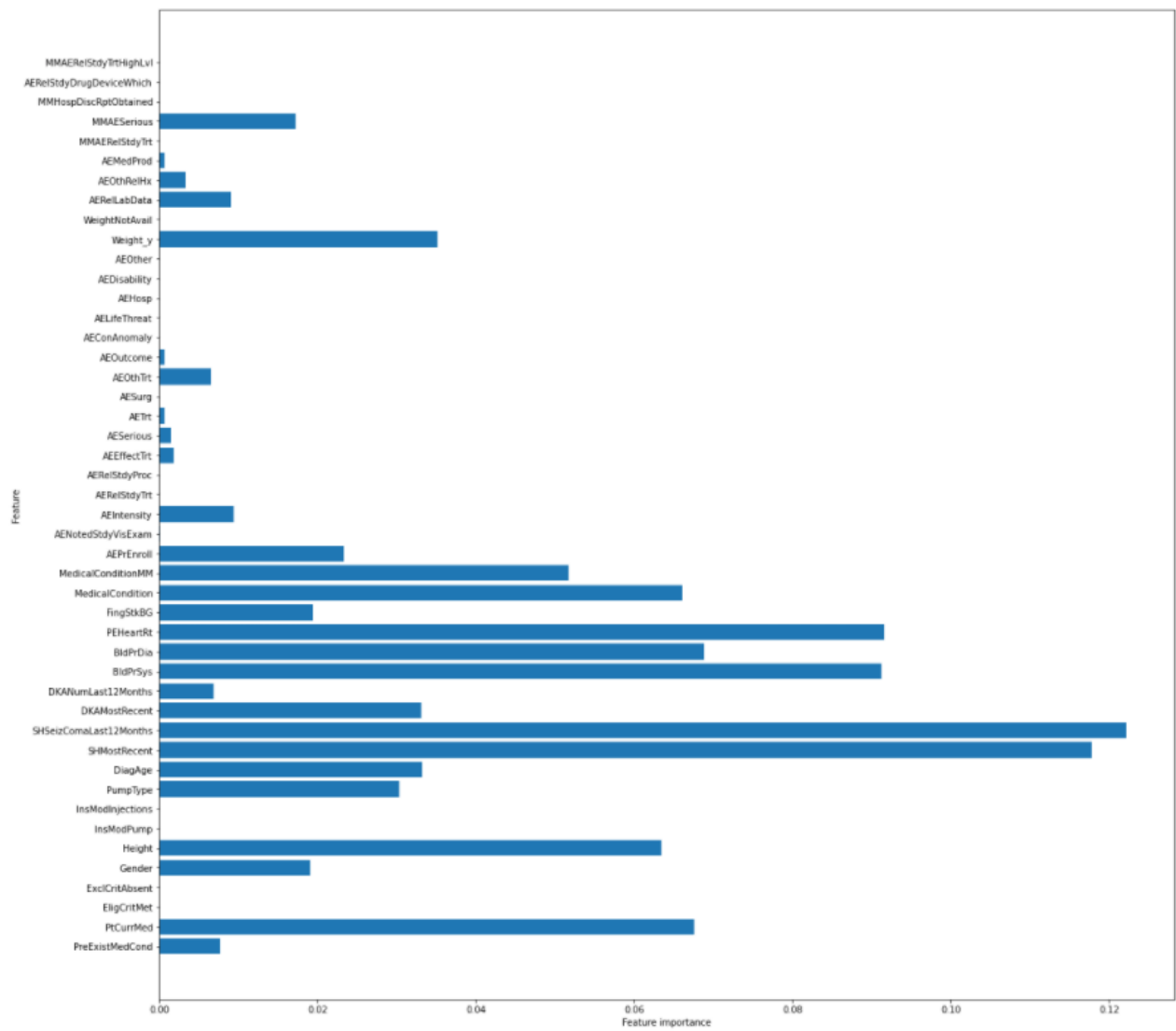
Approach

I dropped features such as PtID, Record Id, race, ethnicity and features which do not have any data recorded. Basically, they were empty columns. Removed features which are correlated with other features and included only features based on feature importance. I used RandomForestClassifier module to use the forest of trees to evaluate the importance of features. The blue bars are the feature importance's of the forest, along with their inter-trees variability represented by the error bars.

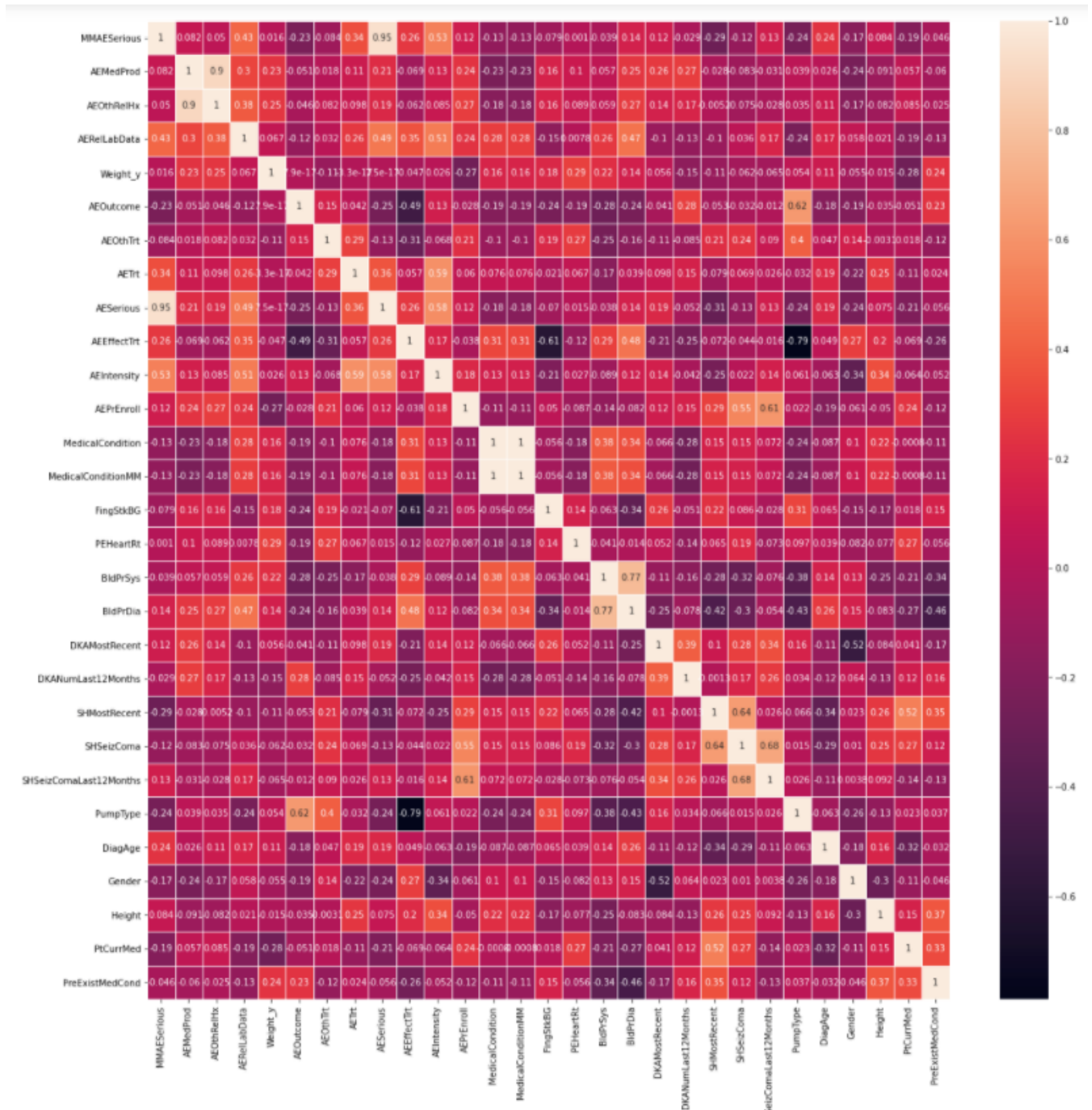
Feature importance Diab1 dataset



Feature importance Diabetes2 dataset



Corr Plot



MMAESerious	1	0.082	0.05	0.43	0.016	-0.23	0.084	0.34	0.95	0.26	0.53	0.12	-0.13	-0.13	0.079	0.001	0.039	0.14	0.12	0.029	-0.29	-0.12	0.13	-0.24	0.24	-0.17	0.084	-0.19	-0.046
AEMedProd	-0.082	1	0.9	0.3	0.23	-0.051	0.018	0.11	0.21	-0.069	0.13	0.24	-0.23	-0.23	0.16	0.1	0.057	0.25	0.26	0.27	-0.028	-0.083	-0.031	0.039	0.026	-0.24	-0.091	0.057	-0.06
AEOTHRelHx	-0.05	0.9	1	0.38	0.25	0.046	0.082	0.098	0.19	0.062	0.085	0.27	-0.18	-0.18	0.16	0.089	0.059	0.27	0.14	0.17	-0.005	-0.075	0.028	0.035	0.11	-0.17	0.082	0.085	0.025
AERelLabData	-0.43	0.3	0.38	1	0.067	-0.12	0.032	0.26	0.49	0.35	0.51	0.24	0.28	0.28	-0.15	0.0078	0.26	0.47	-0.1	-0.13	-0.1	0.036	0.17	-0.24	0.17	0.058	0.021	-0.19	-0.13
Weight_y	-0.016	0.23	0.25	0.067	1	7.9e-17	-0.113	3e-17	5e-17	0.047	0.026	0.27	0.16	0.16	0.18	0.29	0.22	0.14	0.056	-0.15	-0.11	-0.062	0.065	0.054	0.11	-0.055	-0.015	-0.28	0.24
AEOutcome	-0.23	-0.051	0.046	-0.12	7.9e-17	1	0.15	0.042	-0.25	-0.49	0.13	0.026	-0.19	-0.19	-0.24	-0.19	-0.28	-0.24	0.041	0.28	-0.053	0.032	-0.012	0.62	-0.18	-0.19	-0.035	0.051	0.23
AEOTHt	-0.084	0.018	0.082	0.032	-0.11	0.15	1	0.29	-0.13	-0.31	-0.068	0.21	-0.1	-0.1	0.19	0.27	0.25	-0.16	-0.11	0.085	0.21	0.24	0.09	0.4	0.047	0.14	0.003	0.018	-0.12
AETit	0.34	0.11	0.098	0.26	3e-17	0.042	0.29	1	0.36	0.057	0.59	0.06	0.076	0.076	0.021	0.067	-0.17	0.039	0.098	0.15	-0.079	0.069	0.026	-0.032	0.19	-0.22	0.25	-0.11	0.024
AESEious	-0.95	0.21	0.19	0.49	7.5e-17	-0.25	-0.13	0.36	1	0.26	0.58	0.12	-0.18	-0.18	-0.07	0.015	-0.038	0.14	0.19	-0.052	-0.31	-0.13	0.13	-0.24	0.19	-0.24	0.075	-0.21	-0.056
AEEffectTit	0.26	0.069	0.062	0.35	0.047	-0.49	-0.31	0.057	0.26	1	0.17	0.038	0.31	0.31	-0.61	-0.12	0.29	0.48	-0.21	-0.25	-0.072	0.044	0.016	-0.79	0.049	0.27	0.2	-0.069	-0.26
AEIntensity	-0.53	0.13	0.085	0.51	0.026	0.13	-0.065	0.59	0.58	0.17	1	0.18	0.13	0.13	-0.21	0.027	-0.089	0.12	0.14	-0.042	-0.25	0.022	0.14	0.061	-0.063	-0.34	0.34	-0.064	-0.052
AEPrEnroll	-0.12	0.24	0.27	0.24	-0.27	0.028	0.21	0.06	0.12	0.038	0.18	1	-0.11	-0.11	0.05	-0.087	-0.14	-0.082	0.12	0.15	0.29	0.55	0.61	0.022	-0.19	-0.061	0.05	0.24	-0.12
MedicalCondition	-0.13	-0.23	-0.18	0.28	0.16	-0.19	-0.1	0.076	-0.18	0.31	0.13	-0.11	1	1	-0.056	-0.18	0.38	0.34	-0.066	-0.28	0.15	0.15	0.072	-0.24	-0.087	0.1	0.22	-0.000	-0.11
MedicalConditionMM	-0.13	-0.23	-0.18	0.28	0.16	-0.19	-0.1	0.076	-0.18	0.31	0.13	-0.11	1	1	-0.056	-0.18	0.38	0.34	-0.066	-0.28	0.15	0.15	0.072	-0.24	-0.087	0.1	0.22	-0.000	-0.11
FingStkBG	-0.079	0.16	0.16	-0.15	0.18	-0.24	0.19	-0.021	-0.07	-0.61	-0.21	0.05	-0.056	-0.056	1	0.14	-0.063	-0.34	0.26	-0.051	0.22	0.086	-0.028	0.31	0.085	-0.15	-0.17	0.018	0.15
PEHeartRt	-0.001	0.1	0.089	0.0078	0.29	-0.19	0.27	0.067	0.015	-0.12	0.027	0.087	-0.18	-0.18	0.14	1	0.041	0.014	0.052	-0.14	0.065	0.19	-0.073	0.097	0.039	0.082	0.077	0.27	-0.056
BldPrSys	-0.039	0.057	0.059	0.26	0.22	-0.28	-0.25	-0.17	0.038	0.29	-0.089	-0.14	0.38	0.38	0.063	0.041	1	0.77	-0.11	-0.16	-0.28	-0.32	-0.076	-0.38	0.14	0.13	-0.25	-0.21	-0.34
BldPrDia	-0.14	0.25	0.27	0.47	0.14	-0.24	-0.16	0.039	0.14	0.48	0.12	-0.082	0.34	0.34	-0.34	0.014	0.77	1	-0.25	-0.078	-0.42	-0.3	-0.054	-0.43	0.26	0.15	-0.083	-0.27	-0.46
DKAMostRecent	-0.12	0.26	0.14	-0.1	0.056	-0.041	-0.11	0.098	0.19	-0.21	0.14	0.12	-0.066	-0.066	0.26	0.052	-0.11	-0.25	1	0.39	0.1	0.28	0.34	0.16	-0.11	-0.52	0.084	0.041	-0.17
DKANumLast12Months	-0.029	0.27	0.17	-0.13	-0.15	0.28	-0.085	0.15	-0.052	-0.25	-0.042	0.15	-0.28	-0.28	-0.051	-0.14	-0.16	-0.078	0.39	1	-0.0013	0.17	0.26	0.034	-0.12	0.064	-0.13	0.12	0.16
SHMostRecent	-0.29	0.028	0.052	-0.1	-0.11	-0.053	0.21	0.079	-0.31	0.072	-0.25	0.29	0.15	0.15	0.22	0.065	-0.28	-0.42	0.1	-0.0013	1	0.64	0.026	-0.066	-0.34	0.023	0.26	0.52	0.35
SHSeizComa	-0.12	-0.083	-0.075	0.036	-0.062	-0.032	0.24	0.069	-0.13	-0.044	0.022	0.55	0.15	0.15	0.086	0.19	-0.32	-0.3	0.28	0.17	0.64	1	0.68	0.015	-0.29	0.01	0.25	0.27	0.12
SHSeizComaLast12Months	-0.13	-0.031	0.028	0.17	-0.065	0.012	0.09	0.026	0.13	0.016	0.14	0.61	0.072	0.072	0.028	0.073	-0.076	-0.054	0.34	0.26	0.026	0.68	1	0.026	-0.11	0.038	0.092	-0.14	-0.13
PumpType	-0.24	0.039	0.035	-0.24	0.054	0.62	0.4	-0.032	-0.24	-0.79	0.061	0.022	-0.24	-0.24	0.31	0.097	-0.38	-0.43	0.16	0.034	-0.066	0.015	0.026	1	-0.063	-0.26	-0.13	0.023	0.037
DiagAge	-0.24	0.026	0.11	0.17	0.11	-0.18	0.047	0.19	0.19	0.045	-0.063	-0.19	-0.087	-0.087	0.065	0.039	0.14	0.26	-0.11	-0.12	-0.34	-0.29	-0.11	-0.063	1	-0.18	0.16	-0.32	0.032
Gender	-0.17	-0.24	-0.17	0.058	-0.055	-0.19	0.14	-0.22	-0.24	0.27	-0.34	-0.061	0.1	0.1	-0.15	-0.082	0.13	0.15	-0.52	0.064	0.023	0.01	0.038	-0.26	-0.18	1	-0.3	-0.11	-0.046
Height	-0.084	-0.091	-0.082	0.021	0.015	-0.035	0.0031	0.25	0.075	0.2	0.34	-0.05	0.22	0.22	-0.17	-0.077	-0.25	-0.083	-0.084	-0.13	0.26	0.25	0.092	-0.13	0.16	-0.3	1	0.15	0.37
PtCurrMed	-0.19	0.057	0.085	-0.19	-0.28	-0.051	0.018	-0.11	-0.21	-0.069	-0.064	0.24	-0.000	-0.000	0.018	0.27	-0.21	-0.27	0.041	0.12	0.52	0.27	-0.14	0.023	-0.32	-0.11	0.15	1	0.33
PreExistMedCond	-0.046	-0.06	-0.025	-0.13	0.24	0.23	-0.12	0.024	-0.056	-0.26	-0.052	-0.12	-0.11	-0.11	0.15	-0.056	-0.34	-0.46	-0.17	0.16	0.35	0.12	-0.13	0.037	0.032	0.046	0.37	0.33	1
	MMAESerious	AEMedProd	AEOTHRelHx	AERelLabData	Weight_y	AEOutcome	AEOTHt	AETit	AESEious	AEEffectTit	AEIntensity	AEPrEnroll	MedicalCondition	MedicalConditionMM	FingStkBG	PEHeartRt	BldPrSys	BldPrDia	DKAMostRecent	DKANumLast12Months	SHMostRecent	SHSeizComa	SHSeizComaLast12Months	PumpType	DiagAge	Gender	Height	PtCurrMed	PreExistMedCond

Shown below is the ML model to predict SH: Severe hypoglycemia (SH)



Understand the problem statement clearly



Do necessary background study/literature review if available



Analyze and visualize the data to solve the problem



Build models or combination of models on the cleaned different data sets with different predictors



Data cleaning, Feature engineering,



Explore various models, evaluate, compare, and select the model that's most effective in solving the given problem.

Challenge 1: Predict Severe hypoglycemia (SH).

Case: SH event in past 12 months, defined as an event requiring assistance of another person, as a result of altered consciousness or confusion, to administer carbohydrate, glucagon, or other resuscitative actions.

Eligibility and Exclusion Criteria for Cases and Controls:

Eligibility Criteria

- 1) Clinical diagnosis of presumed autoimmune type 1 diabetes
- 2) Age ≥ 60 years old
- 3) Duration of T1D ≥ 20 years
- 4) Insulin presently required
- 5) Has not used a real-time CGM device for the past 2 weeks
- 6) Fluent in English

Exclusion Criteria

- 1) Glomerular filtration rate < 30 (based on available data in medical record from usual care; if no results available, then individual is eligible)
- 2) Diagnosis of dementia that is moderate or advanced
- 3) Serious illnesses where life expectancy is < 1 year
- 4) History of pancreatic transplant

Following tests are performed to record the data. The primary objective of the study is to identify factors associated with the occurrence of severe hypoglycemia in older adults with T1D. It is hoped that the study results can be used to design an intervention study with the goal of reducing the incidence of severe hypoglycemia in older adults.

Testing and Assessments:

1• Cognitive Assessments:

- Montreal Cognitive Assessment
- Symbol Digit Modalities Test
- Trail Making Test
- Hopkins Verbal Learning Test-Revised
- Geriatric Depression Scale
- Diabetes Numeracy Test (shorted version—15 questions)
- Functional Activities Questionnaire
- Grooved Pegboard test
- Vision assessment
- Hypoglycemic Unawareness Assessment 167 –Clarke survey
- Hypoglycemia Fear Survey
- Blood Glucose Attitudes Scale
- Duke Social support scale
- Frailty 10 foot walk
- Blinded continuous glucose monitoring (CGM) for up to 14 days
- Laboratory testing: HbA1c, non-fasting C-peptide, glucose, and creatinine
- Heart activity monitor for up to 14 days for a few participants from selected sites who are willing
- Activity tracker device for up to 14 days for a few participants from selected sites who are willing.

SH frequency increased with longer T1D duration, with the 12-month frequency of one or more events. With respect to HbA1c levels, the frequency of SH was lowest with HbA1c levels. 42 Latino and Hispanic participants had greater frequency of SH than non-Hispanic whites and 298. About 50 pump users fall in group 2-<5 years of range. Injection users had a greater frequency of SH than pump users.

Feature Engineering & Data Cleaning

- First, In DiabScreening_data had string data in few columns PreExistMedCond, PtCurrMed, Gender, PumpType, DKAMostRecent, SHMostRecent, transformed to numeric data using map function to map the column in pandas data frame. I did replace Nan values with mean and 0. DiabScreening_data had 43 features. After, I dropped all empty columns. There are 18 features in this dataset.
- DiabPhysExam contains PtID, Systolic and diastolic blood pressure, Heart rate, Blood glucose I merged DiabPhysExam and DiabScreening_data over PtID and created new dataset called Diabetes_Data
- AdvEvent_data table - This table contains 47 features, PtID, medical conditions, adverse events such as disability death etc. Merged Diabetes_Data and AdvEvent_data table over PtID and created new dataset called Diabetes1. Transformed string data in

Diabetes1 dataset to numeric data using map function to map the column in pandas data frame.

I did replace Nan values with mean and 0.

Dropped following columns based on feature importance and created new dataset called Diabetes2

```
[ 'Weight_x', 'PEAbnormal_x', 'PEAbnormal_y', 'MMUnexpected', 'AERelStdyDrugDevice', 'AERelStdyDrugDeviceUncertain', 'AEDeathCause', 'WeightMeas', 'AEDeathCause', 'AEDeath', 'AEDeathDt', 'AESurgDt', 'AEOnsetDt', 'PtID_x', 'PtID_y', 'RecID', 'PtID', 'ParentLoginVisitID', 'AENotifiedDt', 'AEOnsetDt', 'AEResDt', 'WeightMeas', 'AdverseEventType', 'AERelStdyTrtUncertain', 'MMUnexpected', 'AERelStdyDrugDevice', 'AERelStdyDrugDeviceUncertain', 'AERelStdyTrtHighLvl', 'AERelStdyTrtWhich', 'AERelStdyDrugDeviceHighLvl' ]
```

Features available in our final Diabetes2 data frame

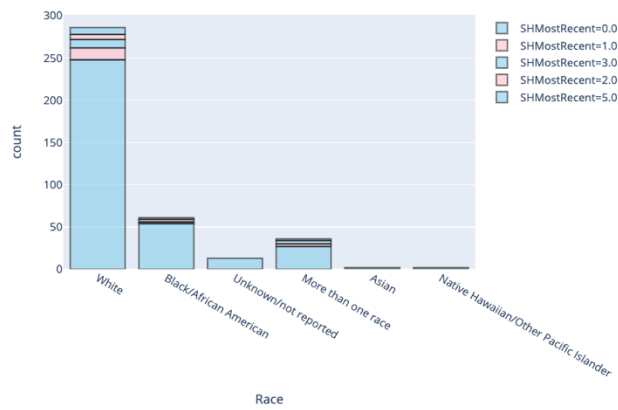
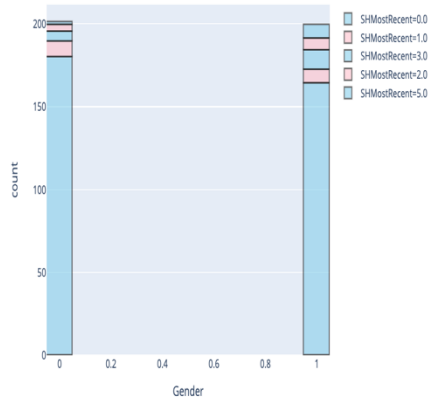
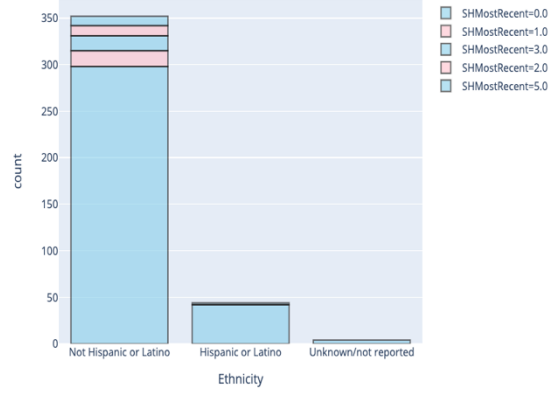
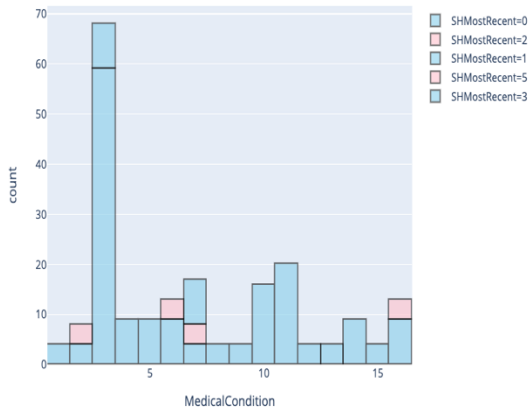
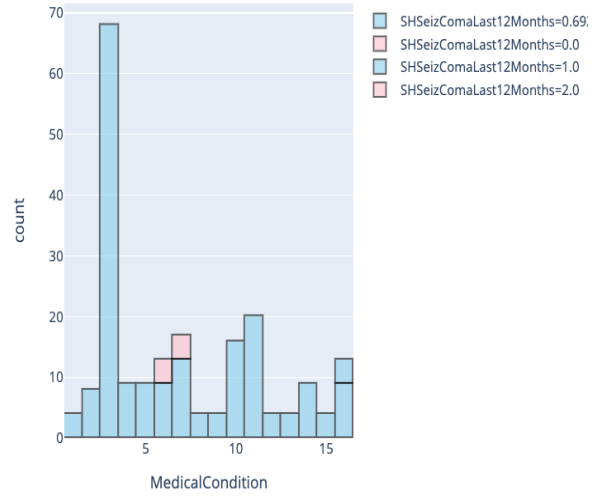
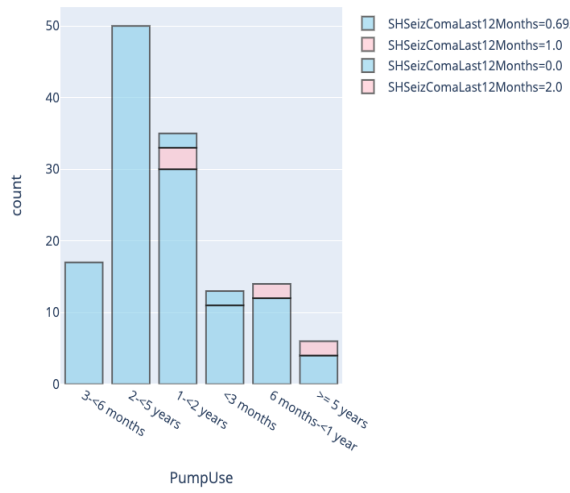
	MMAESerious	AEMedProd	AEOTHRelHx	AERelLabData	Weight_y	AEOutcome	AEOTHTrt	AETrt	AESerious	AEEffectTrt	...	DKANumLast12Months	SHM...
0	0.0	0.0	0.0	0.0	27.063492	1	1.0	1	0	1	...	1	1
1	0.0	0.0	0.0	0.0	27.063492	1	1.0	1	0	1	...	1	1
2	1.0	1.0	1.0	1.0	15.000000	0	1.0	1	1	2	...	0	0
3	1.0	1.0	1.0	1.0	15.000000	0	1.0	1	1	2	...	0	0
4	0.0	1.0	1.0	1.0	25.000000	0	0.0	1	1	2	...	0	0
...
201	1.0	0.0	0.0	0.0	27.063492	0	0.0	1	1	2	...	1	1
202	1.0	0.0	0.0	0.0	27.063492	0	0.0	1	1	2	...	1	1
203	0.0	0.0	0.0	0.0	27.063492	0	1.0	1	0	2	...	0	0
204	0.0	0.0	0.0	0.0	27.063492	0	1.0	1	0	2	...	0	0
205	0.0	0.0	0.0	0.0	27.063492	0	1.0	1	0	2	...	0	0

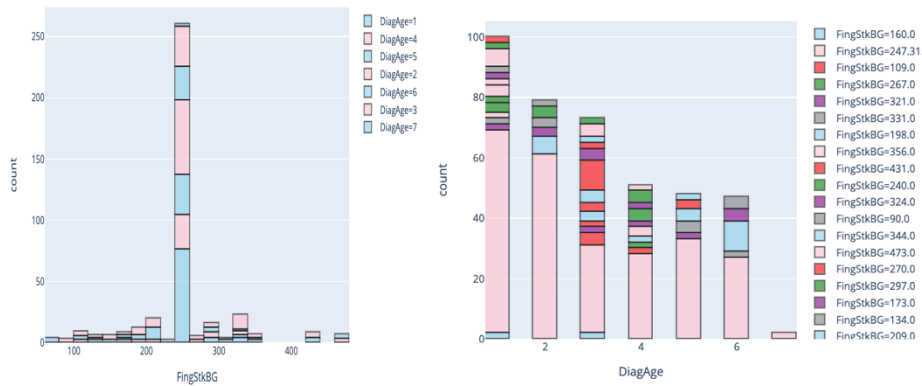
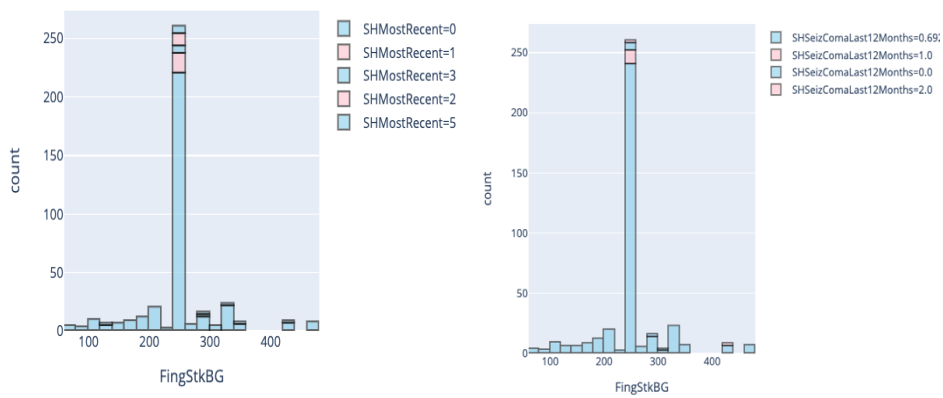
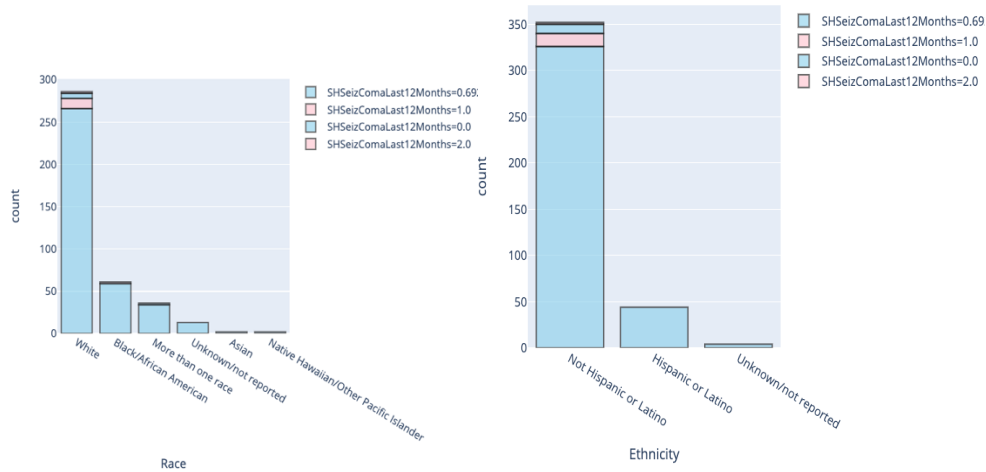
206 rows × 29 columns

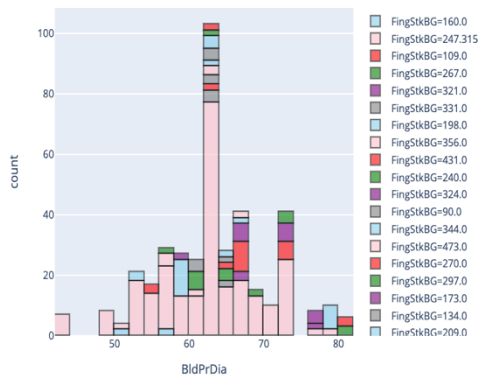
In our final Diabetes2 data frame has 29 features. I created Xm, Ym dataframe with predictors and response variables and used train_test_split function to split 80% data training set and 20% data to testset.

Exploratory Dataset Visualizations

Exploratory visualizations on our final data frame used to predict SH Severe hypoglycemia (SH).

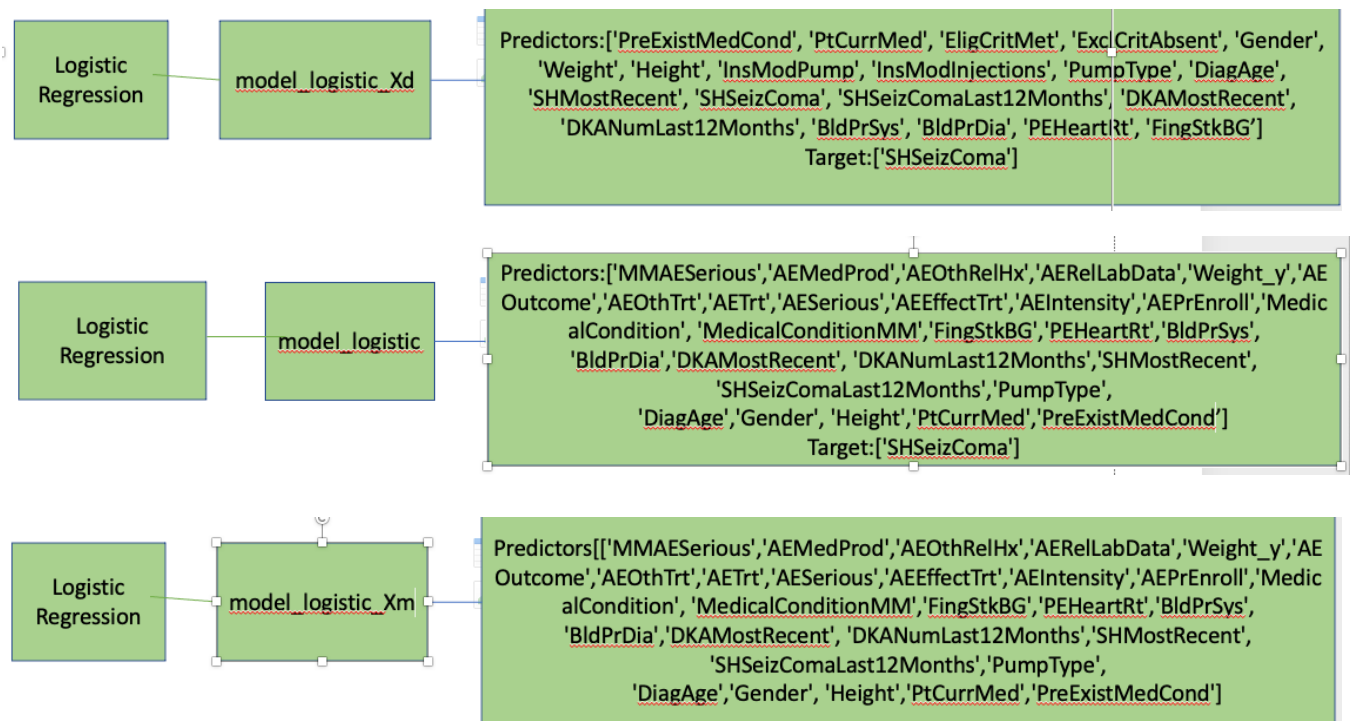






Logistic Regression

In our logistic regression model the target variable is 'SHSeizComa'. The end goal of this model is to predict the SH Severe hypoglycemia (SH). Built 3 different models with different features from the 3 data sets generated Diab1, Diabetes1 and Diabetes2. Dataset details are in feature engineering section.



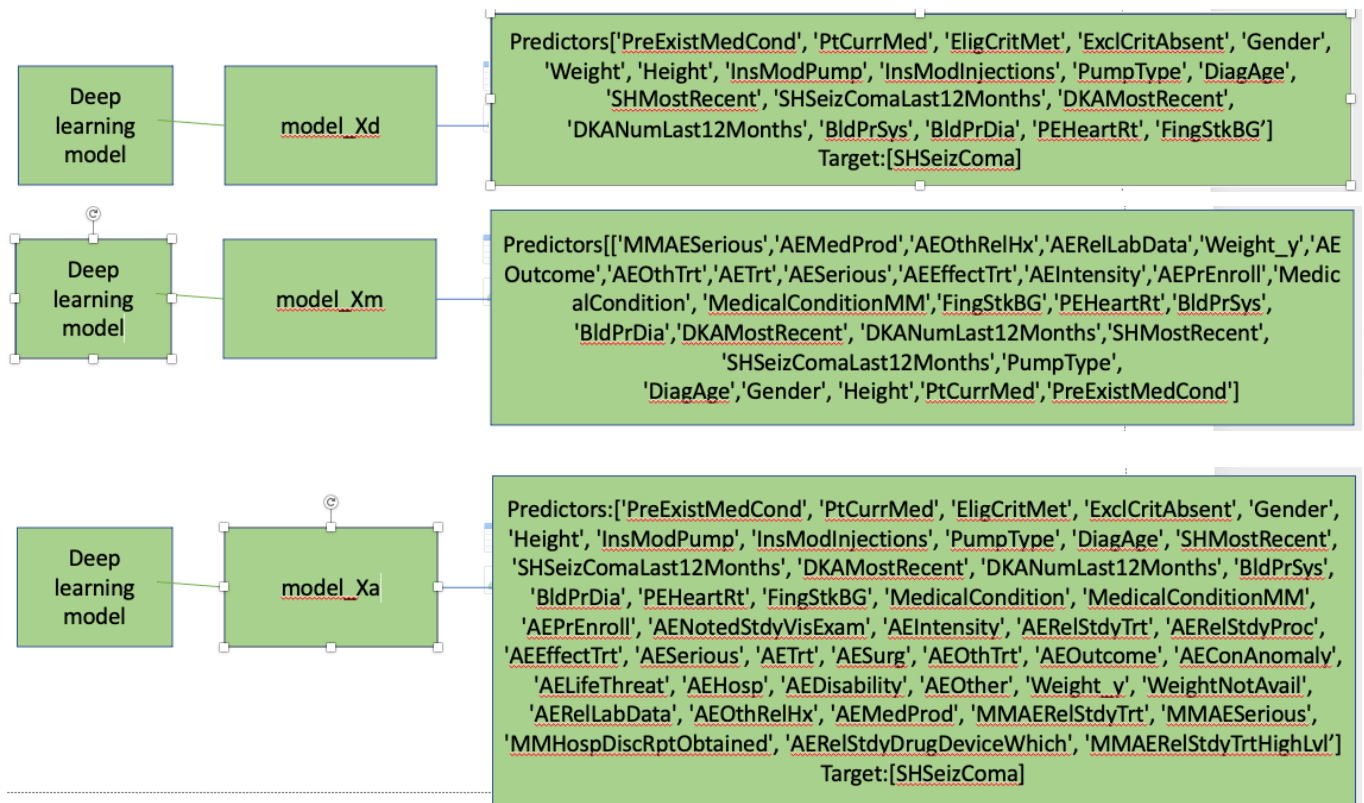
Results of Logistic Regression Models:

Model_logistic_Xm had slightly better accuracy since I included only variables based on feature importance.

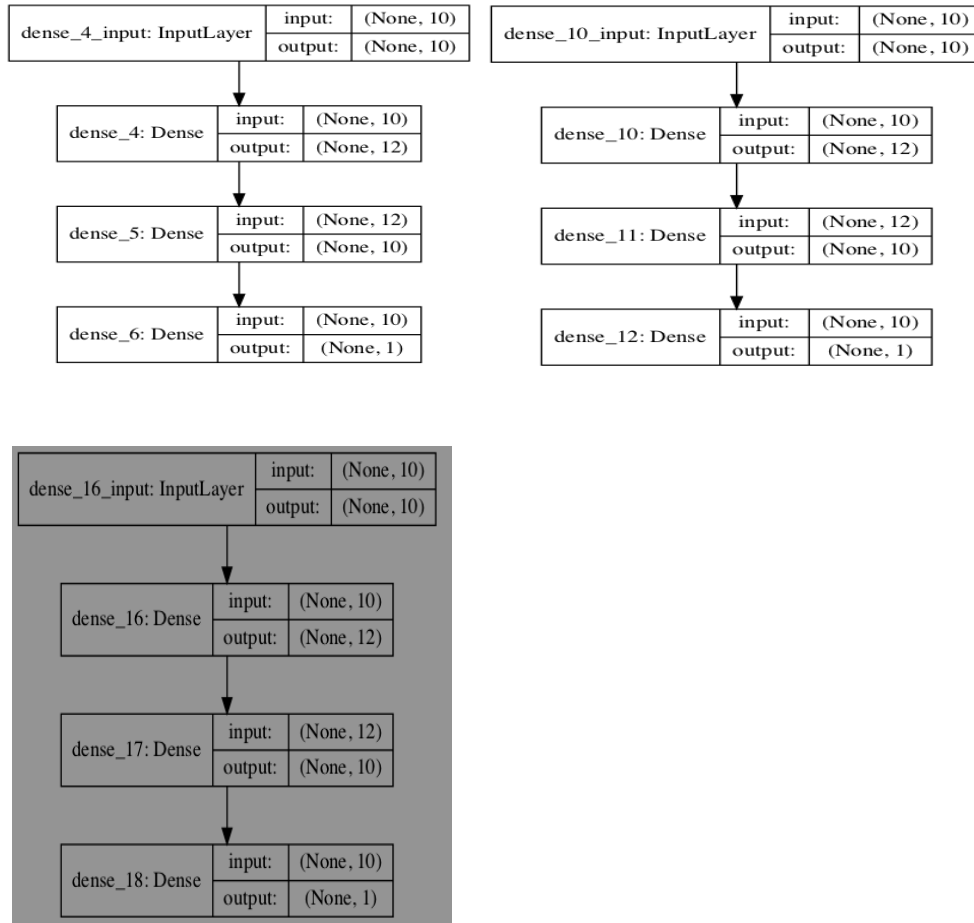
Model_name	Test Accuracy
model_logistic_Xd	0.95238
model_logistic_Xa	0.95
model_logistic_Xm	0.9761

Deep learning Model

Following deep learning classification models were explored in this study:



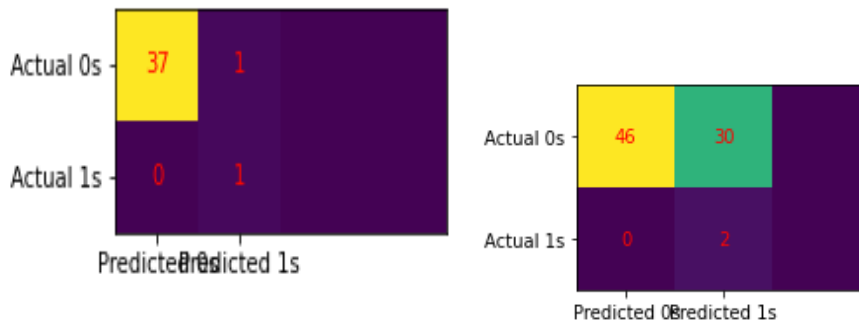
Deep learning model has a simple 3-layer neural net architecture as shown below:
Neural Net architecture use in Model_Xd, Model_Xa, Model_Xm



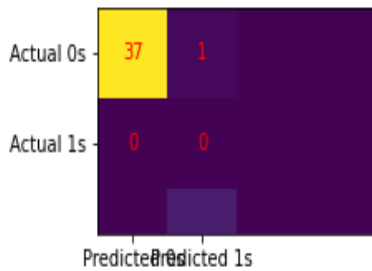
Results of Deep Learning models:

Model name	Test Accuracy
model_Xm	0.88
model_Xd	0.60
model_Xa	0.90

Confusion matrix generated by model_Xa model_Xd



model_Xm



In this case model_Xa has better accuracy compared to other models. Predictors from AdvEvent dataset helped improve the accuracy. Unlike in logistic regression model, including AdvEvent predictors improved the accuracy to 90% of the Neural Net model, seems like the NN model is able to better learn from the adverse event predictors. Overall, of all the models we explored, Model_Xm (logistic regression model after removing features based on feature importance turned out to be better model with 97% score.

Ensemble - Decision Tree/Random Forest

Built Ensemble model on all datasets Diab1, Diabetes1, Diabetes2 using Bagging method, BaggingClassifier, meta-estimator, taking as input a user-specified base estimator along with parameters specifying the strategy to draw random subsets. `max_samples` and `max_features` control the size of the subsets. `bootstrap` and `bootstrap_features` control whether samples and features are drawn with or without replacement. When using a subset of the available samples the generalization accuracy can be estimated with the out-of-bag samples by setting `oob_score=True`. Controls the randomness of the estimator. The features are always randomly permuted at each split, even if `splitter` is set to "best". When `max_features < n_features`, the algorithm will select `max_features` at random at each split before finding the best split among them. But the best found split may vary across different runs, even if `max_features=n_features`

Model name	Test Accuracy
model_Xm	0.83
model_Xd	0.91
model_Xa	1.0

Decision Tree

I built model again on all 3 data sets Diab1, Diabetes1, Diabetes2, with default parameters of the method so decision tree is big and shows that it is overfitting. To validate, I made predictions by using training data and the accuracy was around 90%. Model model_Xd had better accuracy compared to all other models. I did calculate miss-classification rate for all models and plot is shown below for one model. Also, Evaluated decision tree performance on train and test sets with different tree depths with all 3 datasets.

Some of the model parameters are explained below.

`criterion{"gini", "entropy"}, default="gini"`

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain. When the node is pure, value of gini-index or cross-

entropy is small and close to zero. Decision tree algorithm splits nodes as long as this value decreases till it reaches zero or there is no other parameter to stop it.

splitter{*“best”, “random”*}, *default=“best”*

The strategy used to choose the split at each node. Supported strategies are “best” to choose the best split and “random” to choose the best random split.

max_depth*int, default=None*

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

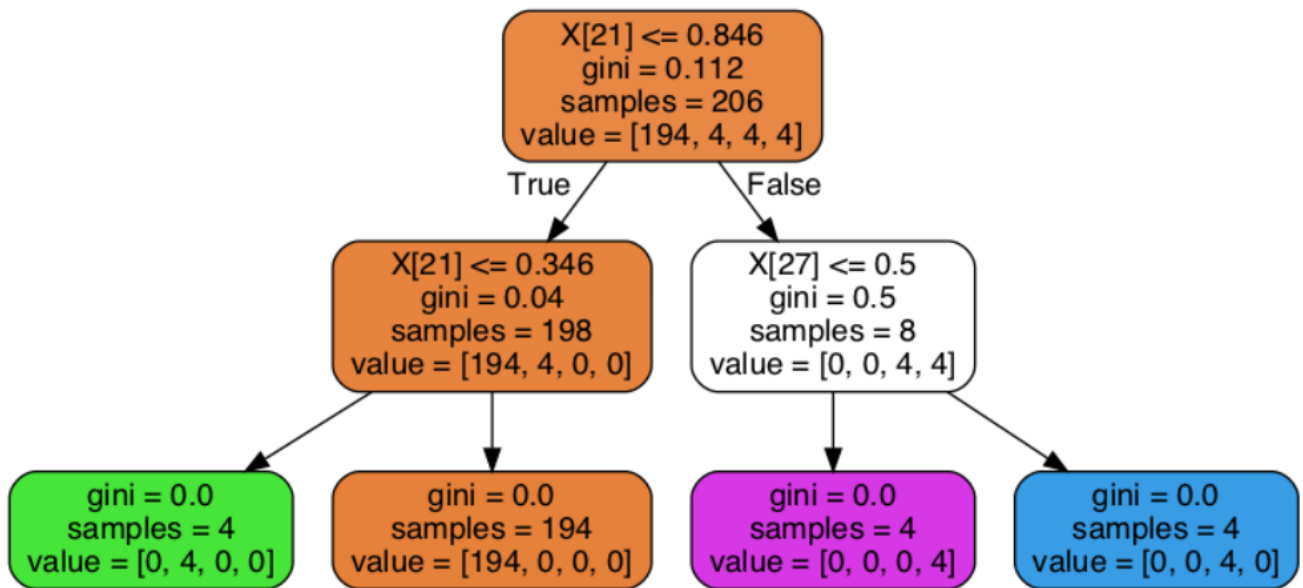
min_samples_split*int or float, default=2*

The minimum number of samples required to split an internal node:

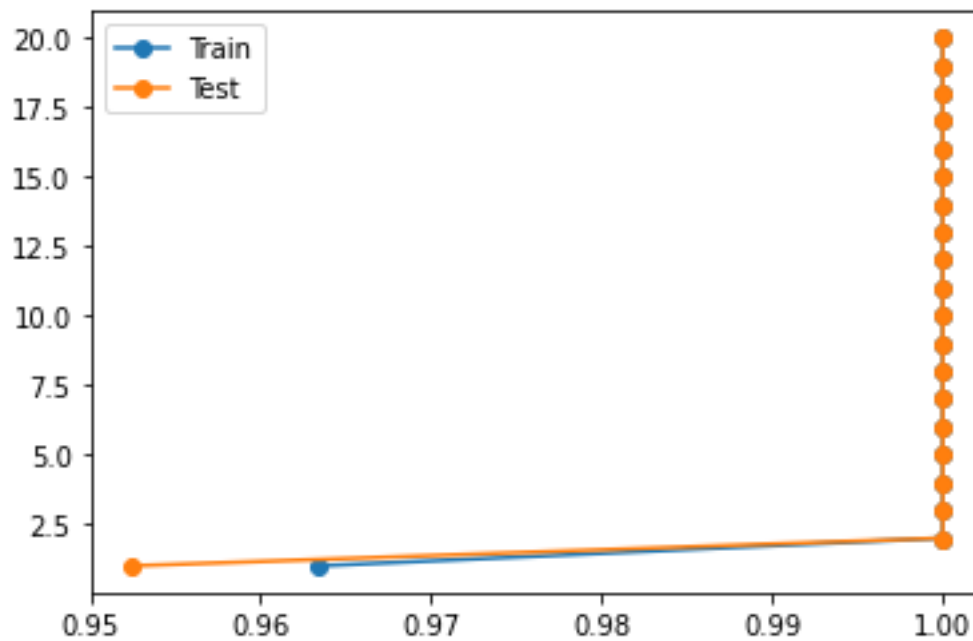
- If int, then consider min_samples_split as the minimum number.
- If float, then min_samples_split is a fraction and $\text{ceil}(\text{min_samples_split} * n_samples)$ are the minimum number of samples for each split.

Model name	Test Accuracy
model_Xm	0.90
model_Xd	0.95
model_Xa	0.90

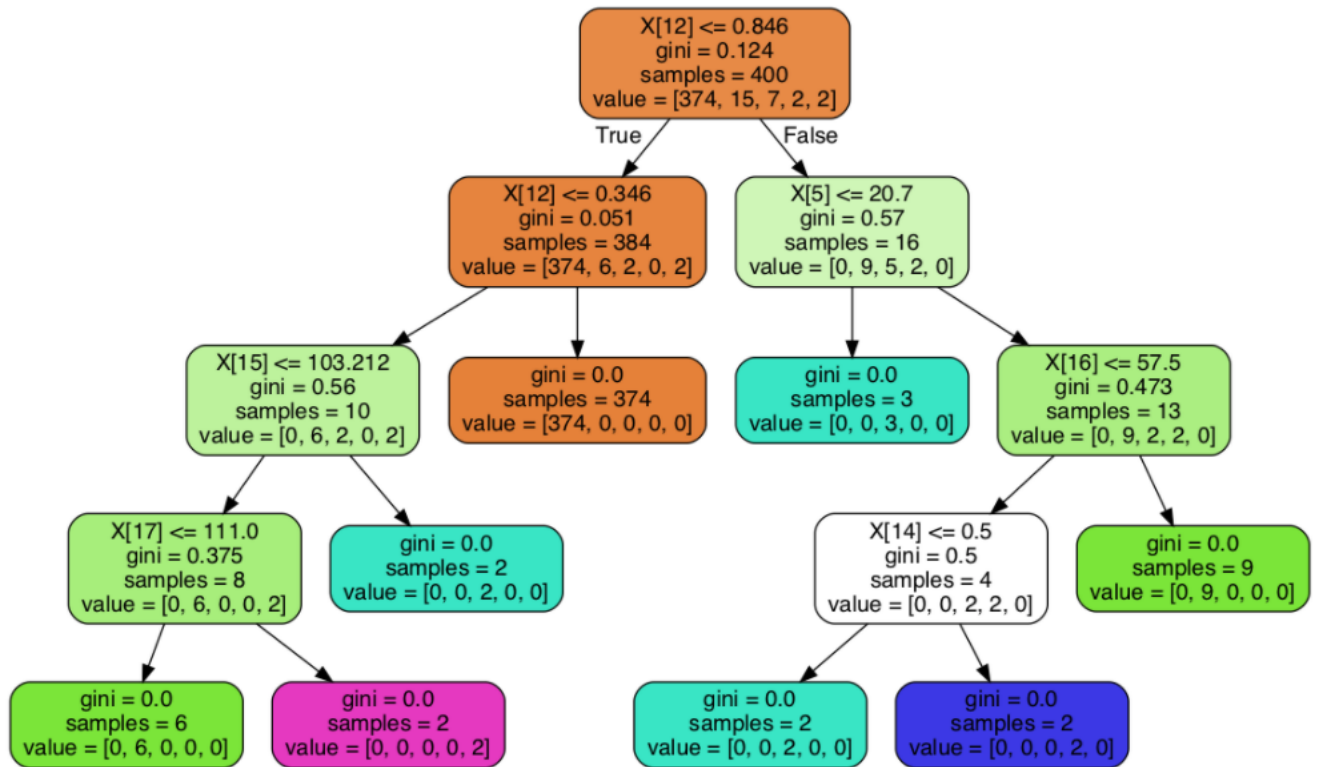
Model with Xa dataset



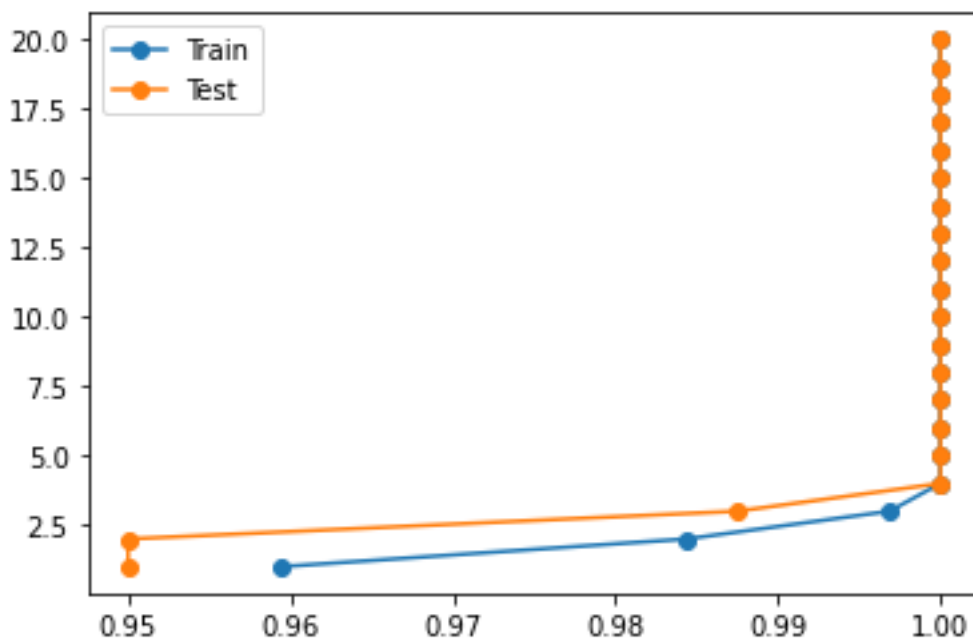
Evaluate decision tree performance on train and test sets with different tree depths with Xa dataset



Model with Xd dataset

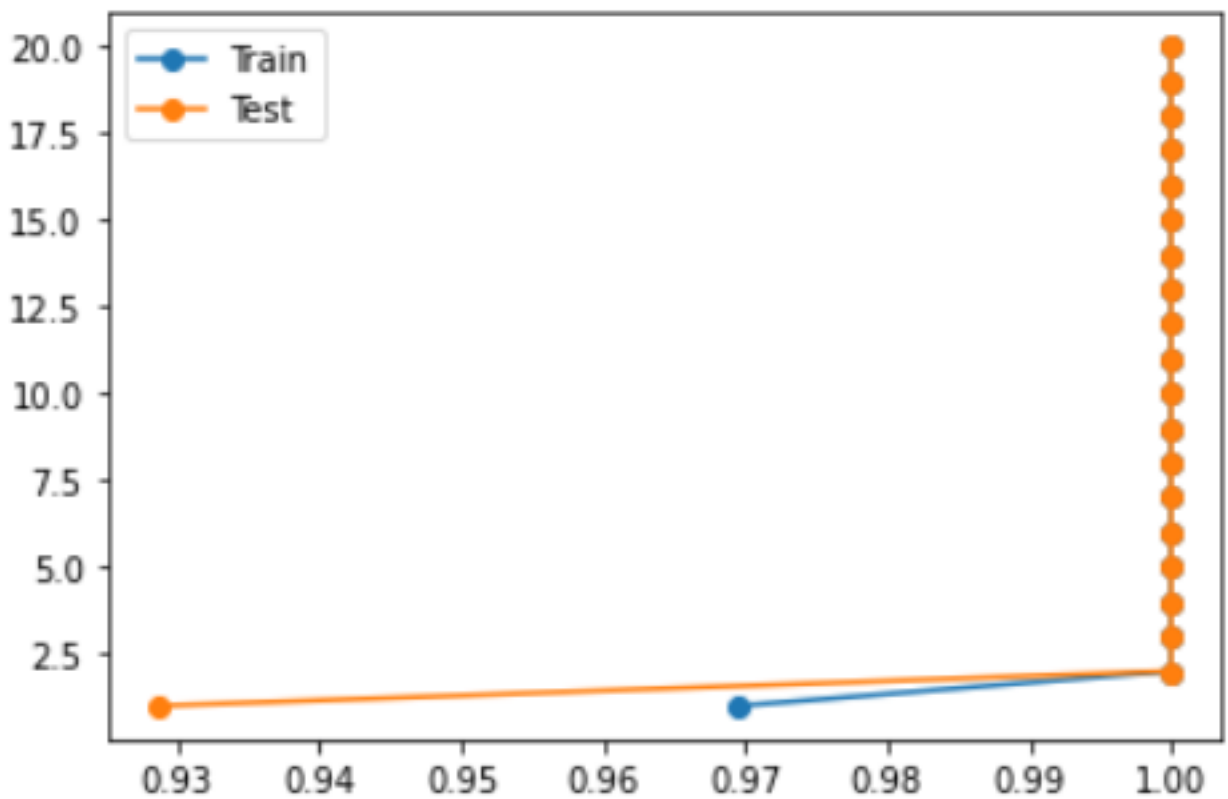
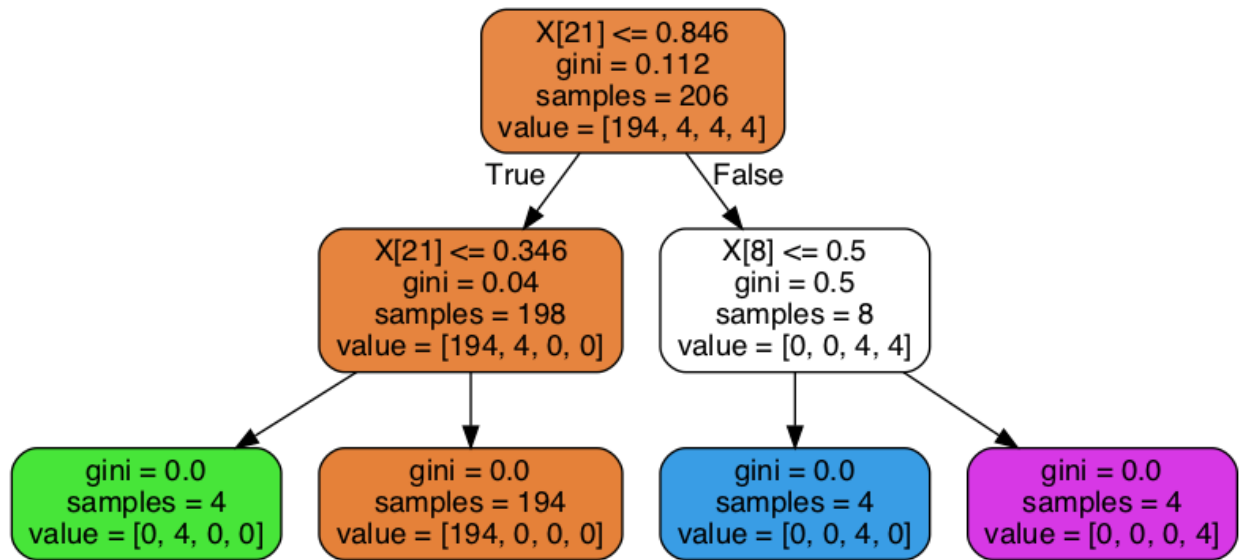


Evaluate decision tree performance on train and test sets with different tree depths with Xd dataset

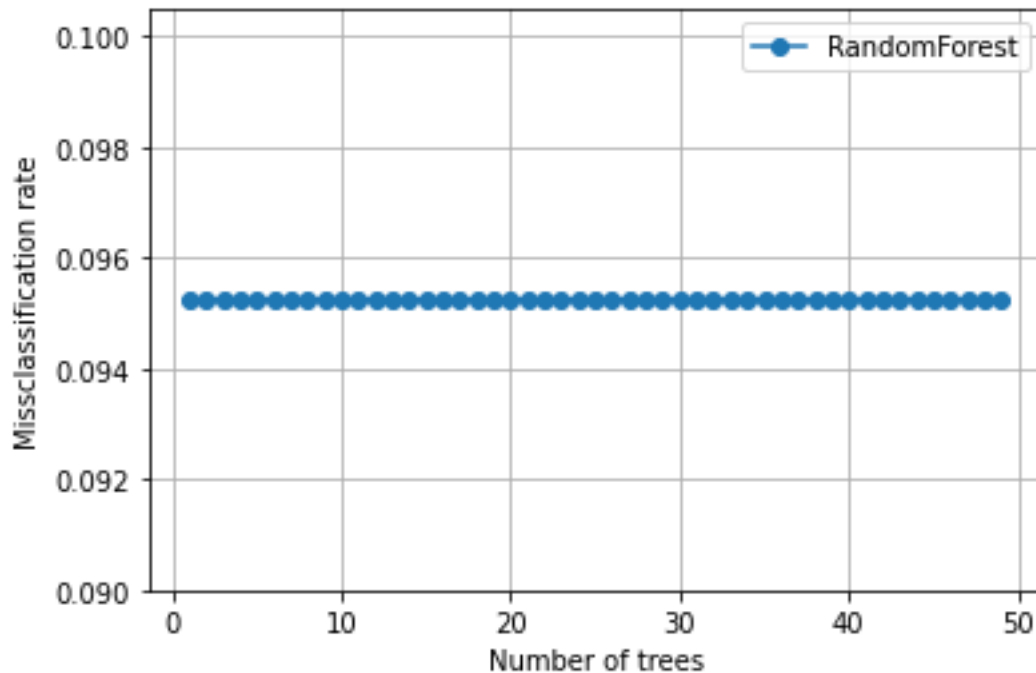


Model with Xm

18]:



Miss classification rate plot with RandomForestClassifier with Xm dataset



Comparing models with K-fold cross validation using accuracy with all 3 datasets

Model scores with Xa dataset Model scores with Xd dataset

PCA
 Logistic regression 0.9876543209876544
 Gaussian Naive Bayes 1.0
 4-Neighbors 0.9512345679012344
 Decisiontree Gini 0.9938271604938271
 Decisiontree Information gain 1.0
 Support vector machine 0.9512345679012344

BFE
 Logistic regression 1.0
 Gaussian Naive Bayes 1.0
 4-Neighbors 0.9419191919191918
 Decisiontree Gini 1.0
 Decisiontree Information gain 1.0
 Support vector machine 0.9419191919191918

All Data
 Logistic regression 1.0
 Gaussian Naive Bayes 1.0
 4-Neighbors 0.9512345679012344
 Decisiontree Gini 1.0
 Decisiontree Information gain 1.0
 Support vector machine 0.9512345679012344

PCA
 Logistic regression 0.9308278867102396
 Gaussian Naive Bayes 0.6844589687726942
 4-Neighbors 0.9244734931009442
 Decisiontree Gini 0.9531590413943355
 Decisiontree Information gain 0.9219317356572257
 Support vector machine 0.9308278867102396

BFE
 Logistic regression 0.9348924022837066
 Gaussian Naive Bayes 0.8975625823451909
 4-Neighbors 0.9250109793588055
 Decisiontree Gini 1.0
 Decisiontree Information gain 0.9873737373737375
 Support vector machine 0.9348924022837066

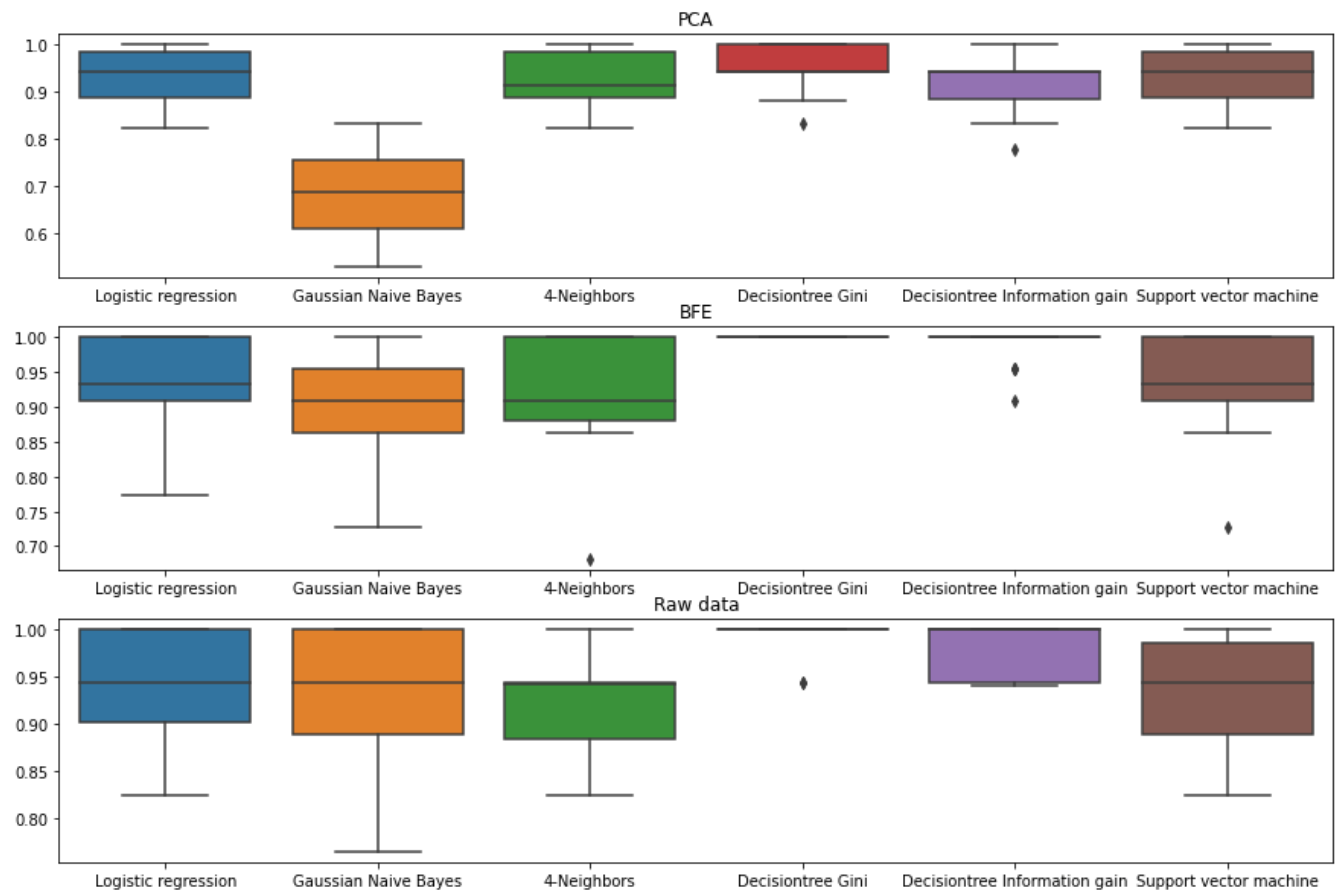
All Data
 Logistic regression 0.9400871459694989
 Gaussian Naive Bayes 0.9221132897603486
 4-Neighbors 0.9213870733478576
 Decisiontree Gini 0.9938271604938271
 Decisiontree Information gain 0.9812999273783587
 Support vector machine 0.9308278867102396

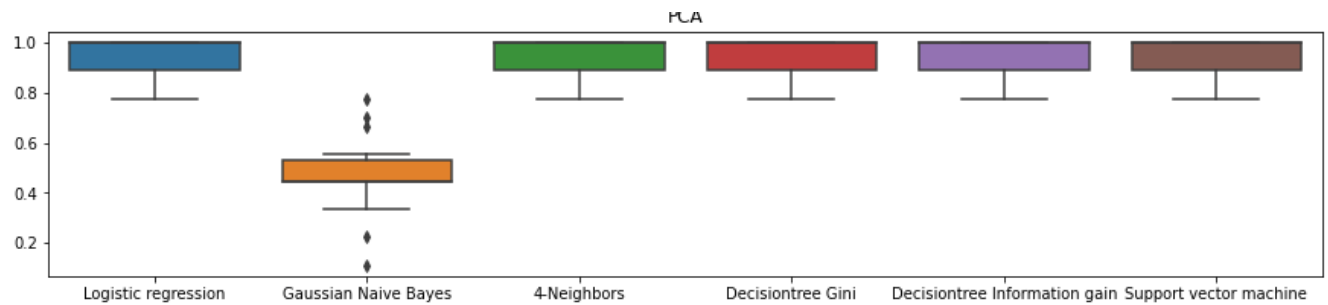
PCA
 Logistic regression 0.9444444444444444
 Gaussian Naive Bayes 0.46234567901234563
 4-Neighbors 0.9506172839506173
 Decisiontree Gini 0.9444444444444444
 Decisiontree Information gain 0.9444444444444444
 Support vector machine 0.9506172839506173

BFE
 Logistic regression 1.0
 Gaussian Naive Bayes 1.0
 4-Neighbors 0.9419191919191918
 Decisiontree Gini 1.0
 Decisiontree Information gain 1.0
 Support vector machine 0.9419191919191918

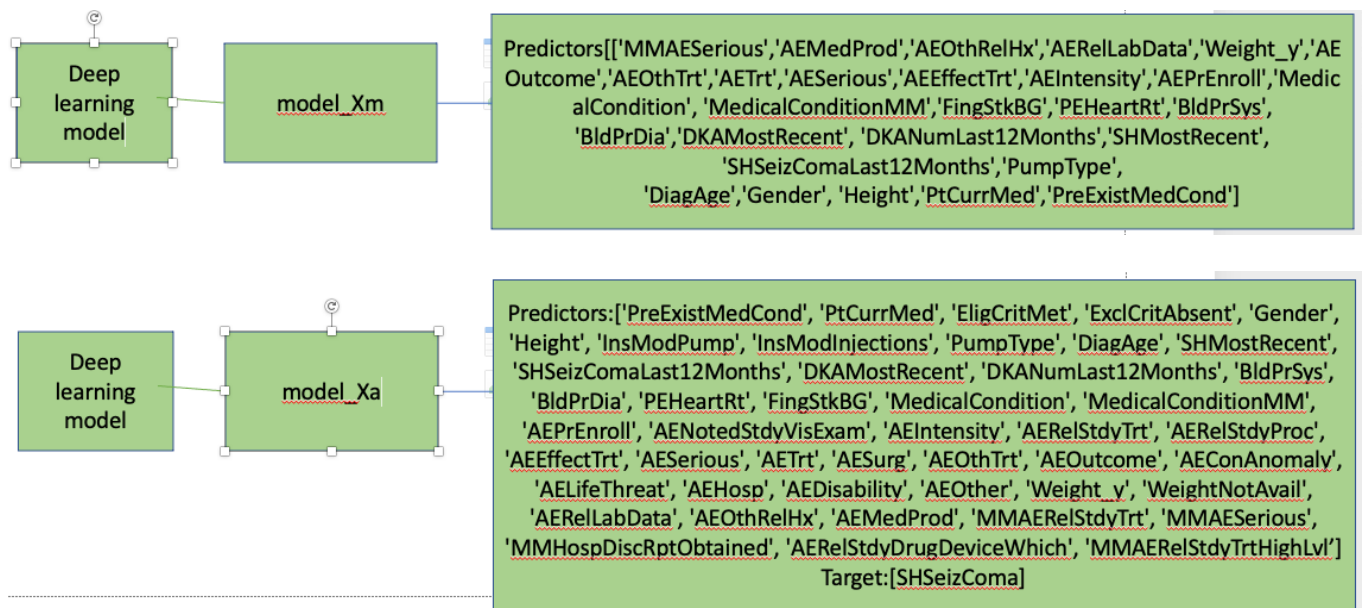
All Data
 Logistic regression 1.0
 Gaussian Naive Bayes 1.0
 4-Neighbors 0.9506172839506173
 Decisiontree Gini 1.0
 Decisiontree Information gain 1.0
 Support vector machine 0.9506172839506173

Boxplots with Xd dataset – Cross validation models with Xd dataset with better scores.





Conclusion: For predicting Severe Hypoglycemia SH Logistic regression model with predictors after removing features based on feature importance seems to be a better model in predicting Severe Hypoglycemia SH with 97%. Deep learning model performed better with prediction score of 90%.



Decision tree model	model_Xd	0.95
Ensemble model	Model_Xa	1.0

References

1. <https://public.jaeb.org/datasets/diabetes>

2. <https://scikit-learn.org/stable/>
3. <https://www.ncbi.nlm.nih.gov/pmc/articles/>
1. Beck RW, Tamborlane WV, Bergenstal RM, Miller KM, Dubose SN, Hall CA. The T1D Exchange Clinic Registry. *J Clin Endocrinol Metab.* 2012.
2. Nordin C. The case for hypoglycaemia as a proarrhythmic event: basic and clinical evidence. *Diabetologia.* 2010; 53:1552-61.
3. McCoy RG, Van Houten HK, Ziegenfuss JY, Shah ND, Wermers RA, Smith SA. Increased mortality of patients with diabetes reporting severe hypoglycemia. *Diabetes Care.* 2012; 35:1897-901.
4. Snell-Bergeon JK, Wadwa RP. Hypoglycemia, diabetes, and cardiovascular disease. *Diabetes Technol Ther.* 2012; 14 Suppl 1:S51-8.
5. Gill GV, Woodward A, Casson IF, Weston PJ. Cardiac arrhythmia and nocturnal hypoglycaemia in type 1 diabetes--the 'dead in bed' syndrome revisited. *Diabetologia.* 2009; 52:42-5.
6. Rothenbuhler A, Bibal CP, Le Fur S, Bougneres P. Effects of a controlled hypoglycaemia
1. 533 test on QTc in adolescents with Type 1 diabetes. *Diabet Med.* 2008; 25:1483-5.