



ST4248 Group B2 Final Presentation

Lim Xi Chen Terry (A0199513W)
Chew Bangyao Paul (A0204816J)
Chan Wen Yong (A0201868B)
Low Hon Zheng (A0204803R)

Table of Contents

01

Problem Statement

02

Data Exploration

03

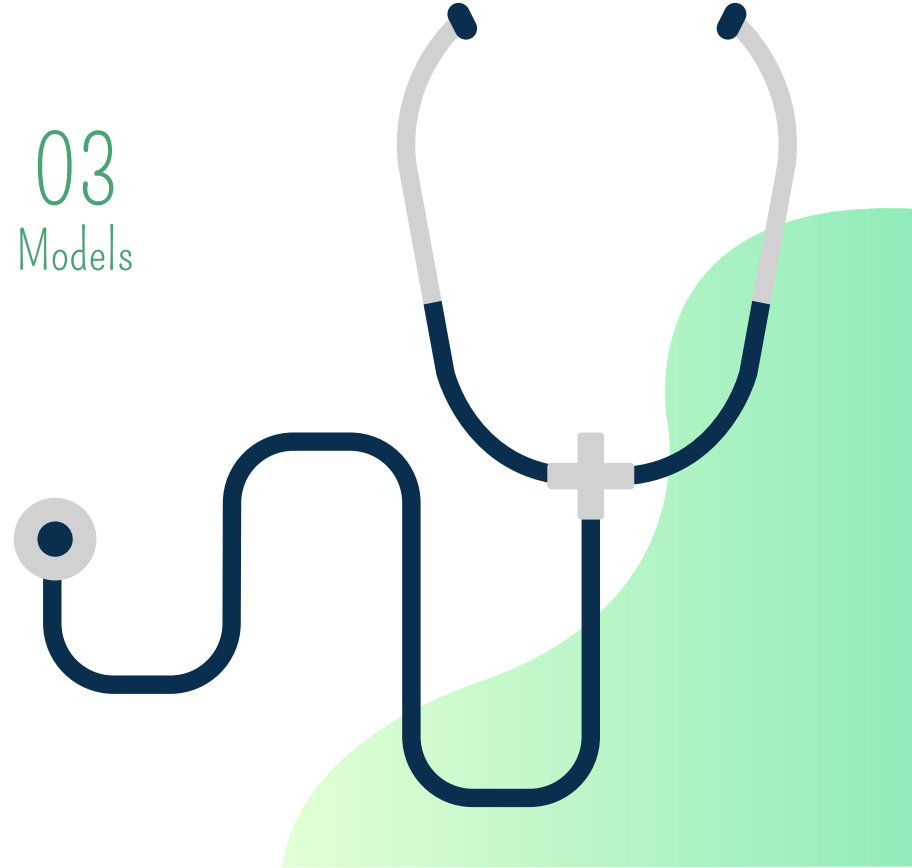
Models

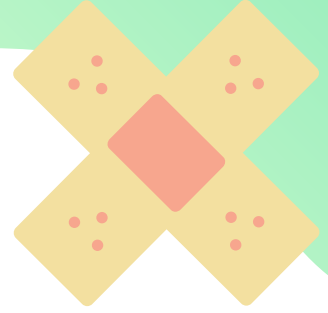
04

Variable Selection

05

Future Works





01

Problem Statement

DIABETES INFOGRAPHICS

Diabetes in numbers

422
MILLION

422 million
people have
diabetes in the
world



that's about 1 out of every 11 people

1 OUT
OF **4**

do not know they
have diabetes

Main types of diabetes



TYPE 1 DIABETES

Body cannot produce
insulin



TYPE 2 DIABETES

Body produce insulin
but can't use it well



GESTATIONAL DIABETES

A temporary condition in
pregnancy (GDM)

Problem Statement

To find out how an individual's profile would determine whether they are likely to be diabetic or not.



Objectives

Prediction



Accurately predicting diabetes using demographics and health – related data

Important Features

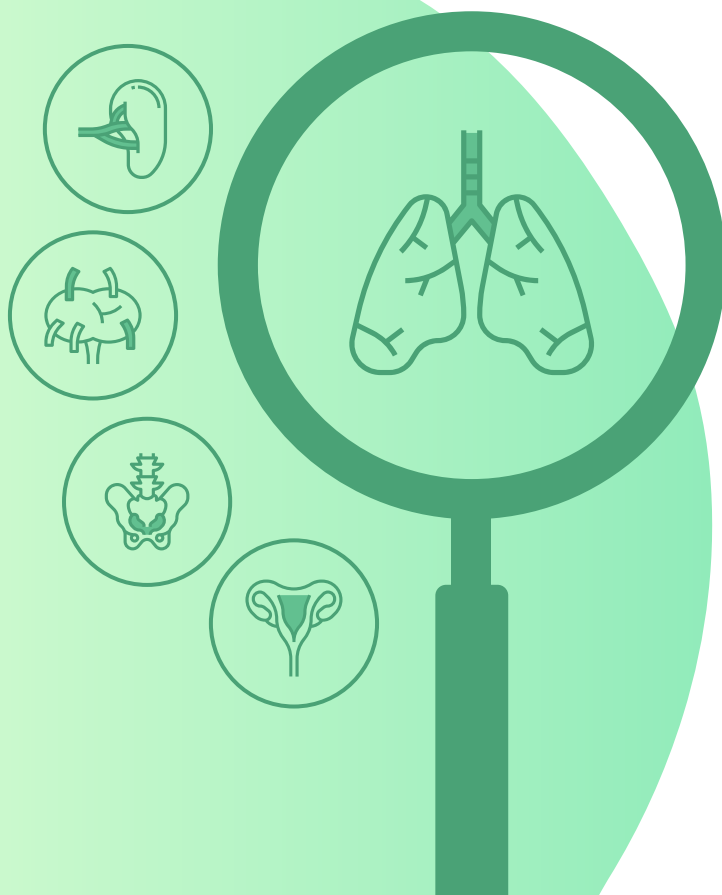


Identifying features that contributes the most to diagnosis of being diabetic

Reduced Feature Model



Using a model trained on important features to accurately predict diabetes



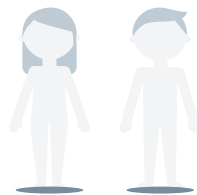
02

Data Exploration

BEHAVIOURAL RISK FACTOR SURVEILLANCE SYSTEMS 2020



America

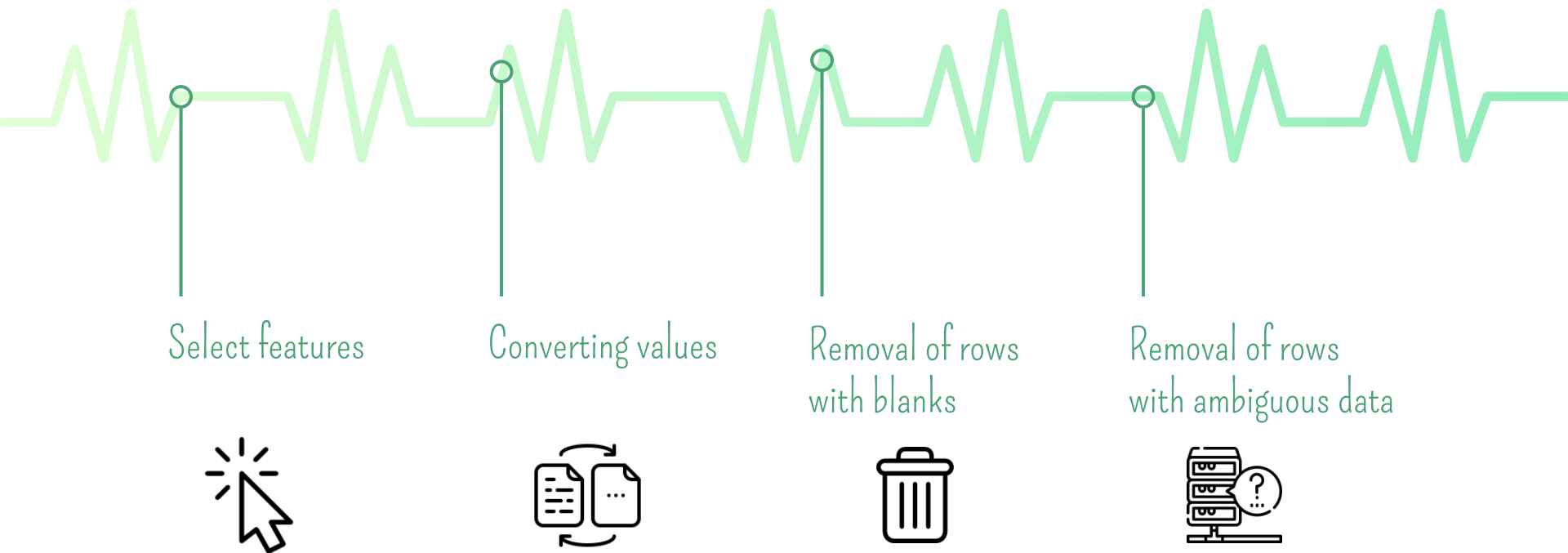


401,958
respondents



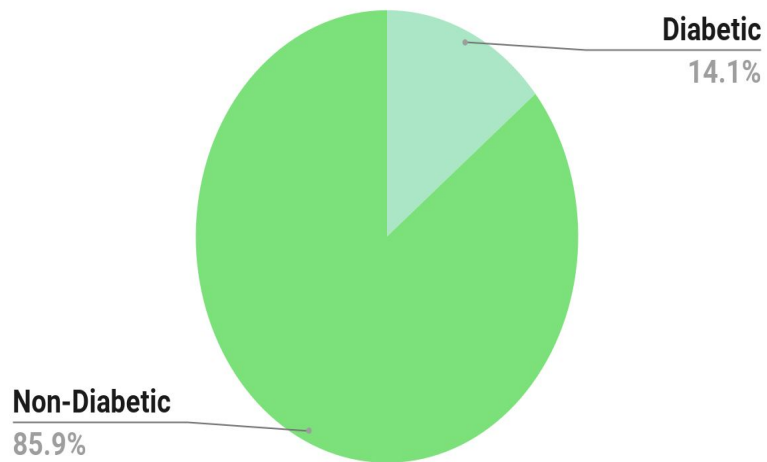
279
questions

Data Cleaning



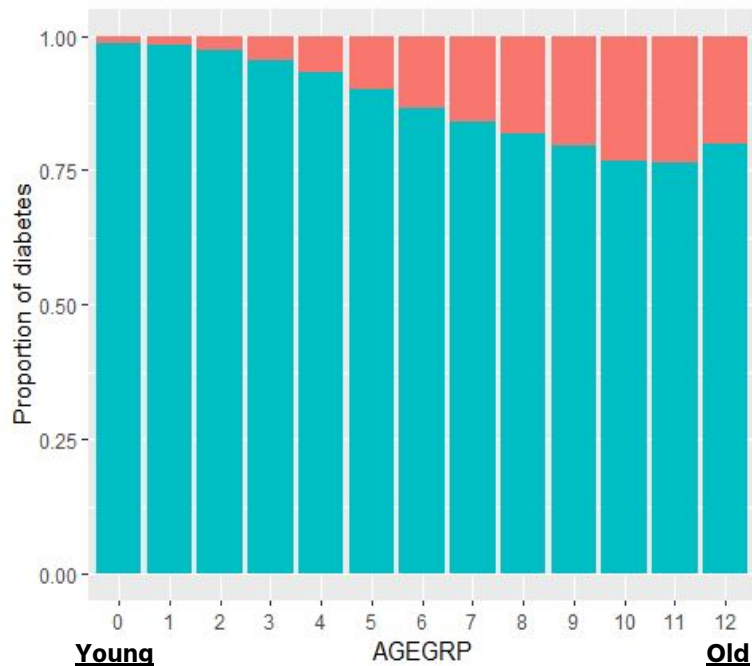
Imbalanced Data

% of respondents



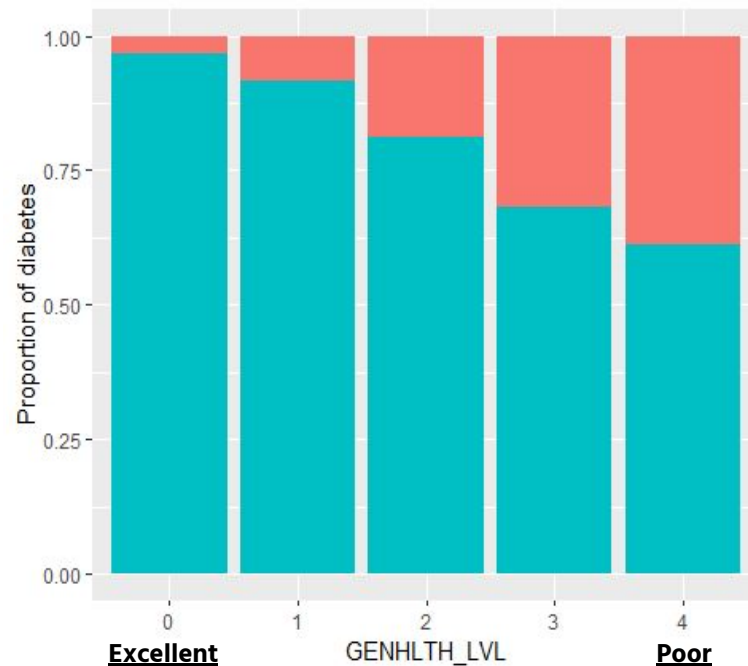
EDA

Age Group



Diabetes Status: ■ = Diabetes ■ = No Diabetes

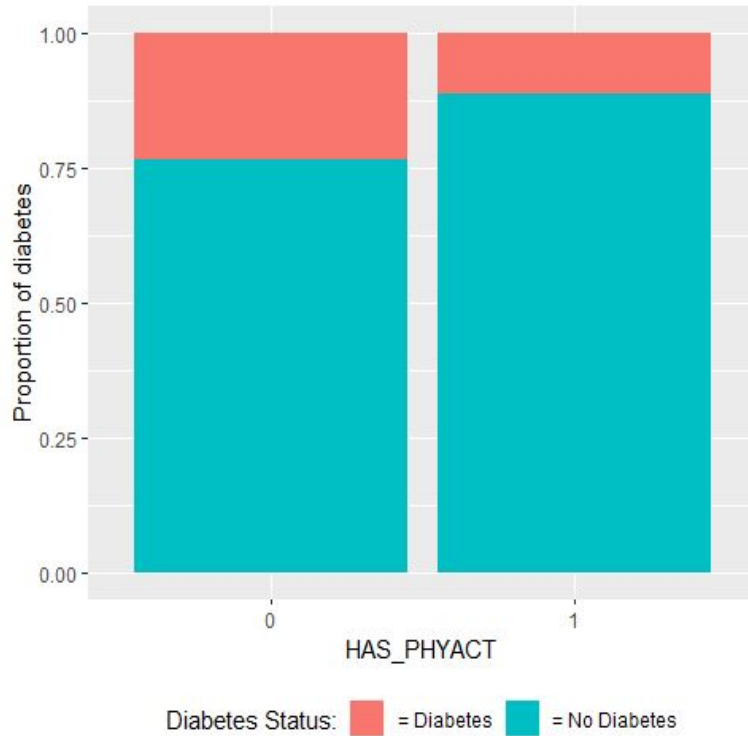
General health



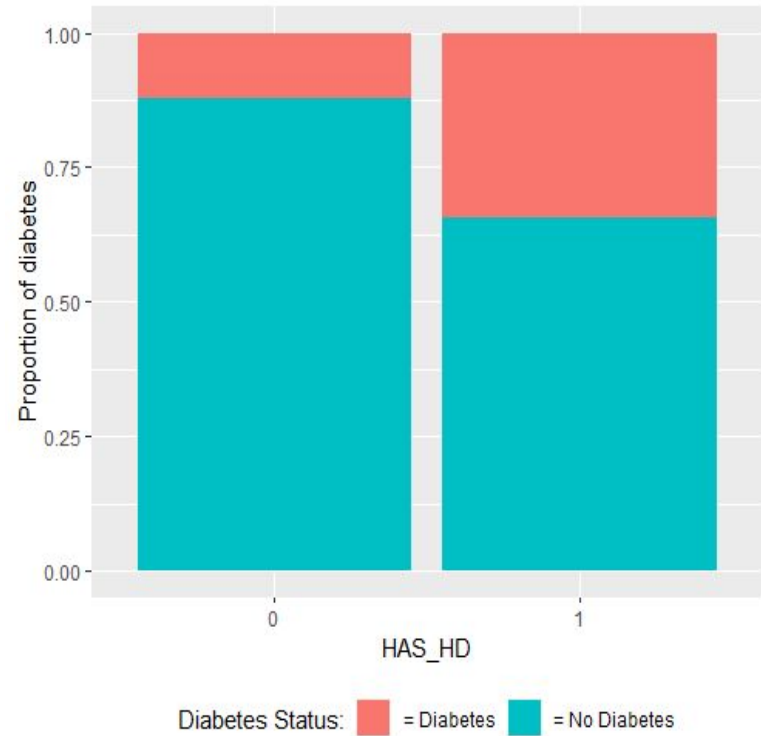
Diabetes Status: ■ = Diabetes ■ = No Diabetes

EDA

Physical Activity

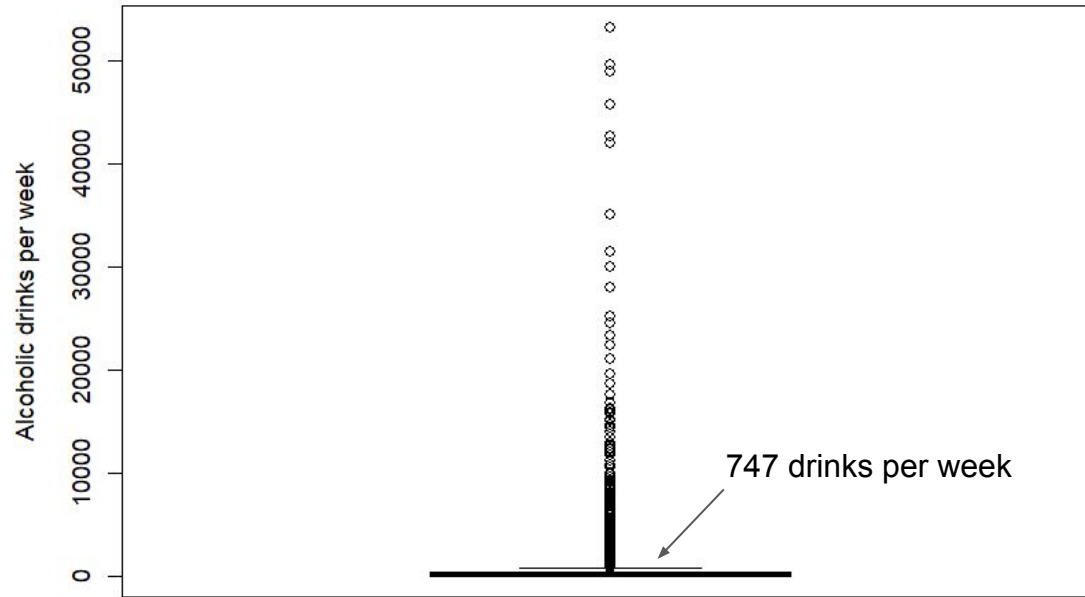


Heart Disease



Outlier removal

Boxplot of Alcoholic drinks per week



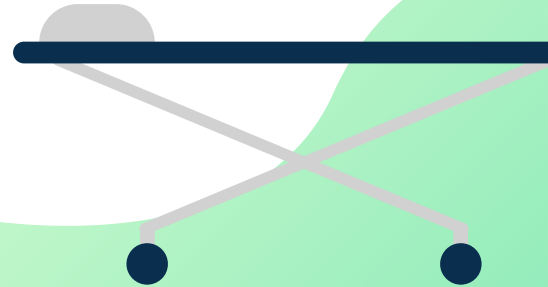
Correlation

Comparison	Measure	Analysis	Result
<u>Continuous</u> variables	Pearson Correlation Coefficient	Moderate correlation found between: <ul style="list-style-type: none">- NUM_POORHLTH and NUM_POORPHYHLTH (0.611)- NUMPOORHLTH and NUM_POORMENTHLTH (0.451)	Remove NUMPOORHLTH
<u>Binary Categorical</u> variables	Cramér's V	All values < 0.5 and are close to 0.	No correlation found
<u>Non-Binary Categorical</u> variables	Phi Coefficient	All values are either < 0.5 or > -0.5 and are close to 0.	
<u>Continuous</u> & <u>Binary Categorical</u> variables	Point-Biserial Correlation Coefficient	All values are either < 0.5 or > -0.5 and are close to 0.	



03

Models



Models



Logistic Regression (Baseline)

No Down-sampling

No scaling

Parametric method



Comparison

Logistic Regression

Accuracy 0.862

Precision 0.540

Recall 0.149

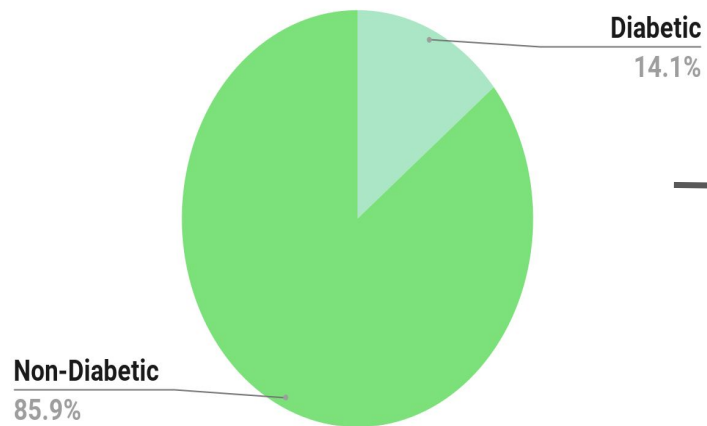
F1 score 0.233

AUC 0.564

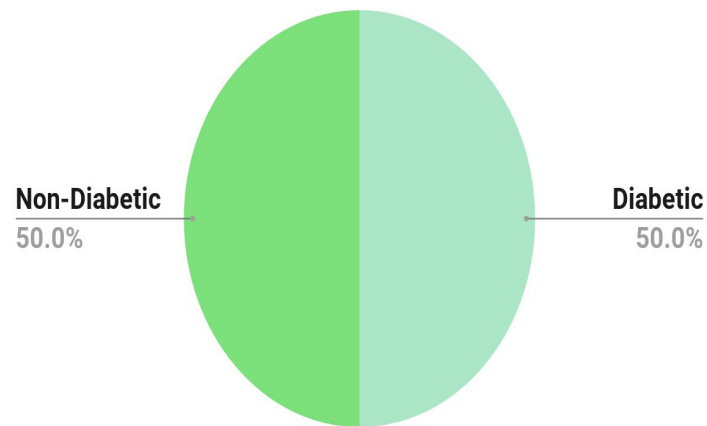


Imbalanced Data

% of respondents



% of train data



Models



Logistic Regression
(Baseline)

No Down-sampling

No scaling

Parametric method



L1 Regularised
Logistic Regression

Down-sampling

No scaling

Parametric method

Comparison



Logistic
Regression

L1 Regularised
Logistic Regression

Accuracy

0.862

0.731

Precision

0.540

0.313

Recall

0.149

0.755



F1 score

0.233

0.442

AUC

0.564

0.741

Models



Logistic Regression
(Baseline)

No Down-sampling

No scaling

Parametric method



L1 Regularised
Logistic Regression

Down-sampling

No scaling

Parametric method



XGBoost

Down-sampling

No scaling

Non parametric method

Comparison



	Logistic Regression	L1 Regularised Logistic Regression	XGBoost
Accuracy	0.862	0.731	0.747
Precision	0.540	0.313	0.326
Recall	0.149	0.755	0.741
F1 score	0.233	0.442	0.452
AUC	0.564	0.741	0.744



Models



Logistic Regression
(Baseline)

No Down-sampling

No scaling

Parametric method



L1 Regularised
Logistic Regression

Down-sampling

No scaling

Parametric method

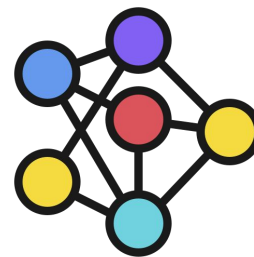


XGBoost

Down-sampling

No scaling

Non parametric method



NN (3 hidden layers)

Down-sampling

Standard scaling

Parametric method

Comparison



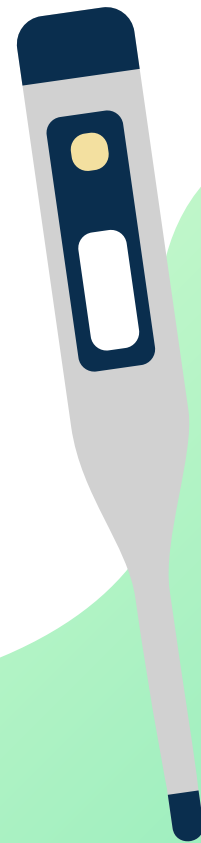
	Logistic Regression	L1 Regularised Logistic Regression	XGBoost	NN
Accuracy	0.862	0.731	0.747	0.747
Precision	0.540	0.313	0.326	0.325
Recall	0.149	0.755	0.741	0.734
F1 score	0.233	0.442	0.452	0.451
AUC	0.564	0.741	0.744	0.742





04

Variable Selection



Reduced feature list

CATEGORICAL VARIABLES

- HAS STROKE
- HAS HEART DISEASE
- HAS PHYSICAL ACTIVITY
- HAS FINANCIAL ISSUES
- RACE
- AGE GROUP
- SEX
- DAYS SINCE LAST CHECK UP
- GENERAL LEVEL HEALTH

CONTINUOUS VARIABLES

- BMI
- NUMBER OF DRINKS DRANK PER WEEK

Reduced feature model

Logistic Regression with
reduced features

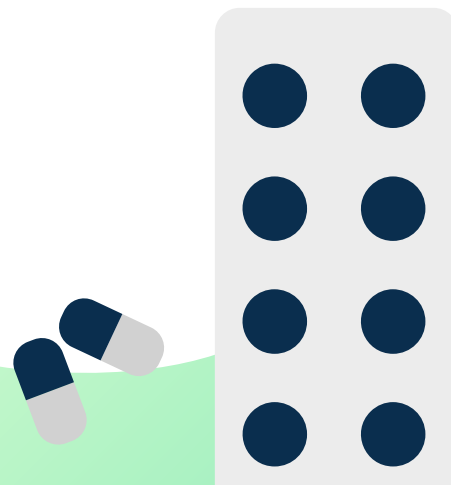
L1 Regularised
Logistic Regression

Accuracy	Precision	Recall	F1 score	AUC
0.723	0.309	0.780	0.443	0.747
0.731	0.313	0.755	0.442	0.741



05

Future Works



Additional Data Collection



Blood pressure

Cholesterol

Diet



THANK YOU!