

BEHAVIOURAL FACTORS LEADING TO CARDIAC ISCHAEMIA:

Getting to the Heart of Heart Disease
with Predictive Analytics

Chalamalasetti Sree Vaishnavi

Teo De Xuan Justin

Siah Wee Hung

Lim Zi Xiang

Nagammai Senthil Kumar



PRESENTATION OUTLINE



Introduction



Problem Statement



Project Objective



Data & Methodology



Modelling



Proposed Business Solution



PREVALENCE OF HEART DISEASE

Cardiovascular diseases 2nd leading cause of death in Singapore
– **nearly 1 in 3 deaths** in 2018





PREVALENCE OF HEART DISEASE

Strain on healthcare resources → Overwhelm healthcare system

The screenshot shows a news article from CNA (Channel NewsAsia) dated March 2022. The headline reads: "Longer waiting times at hospitals with some patients told to wait up to 50 hours for a bed". Below the headline is a sub-headline: "Ng Teng Fong and Tan Tock Seng hospitals were among those with wait times of more than half a day, according to patients and health-care workers". A photograph shows several people standing in a hallway, likely a hospital corridor. A caption below the photo states: "Patients waiting outside Ng Teng Fong General Hospital's emergency department on Mar 18, 2022. Photo: CNA". The page includes standard news navigation elements like Top Stories, Latest News, Discover, Singapore, Asia, Comments, Sustainability, CNA Insider, Lifestyle, Health, Letters, and +40 Sections.

(CNA, 2022)

The screenshot shows a news article from THE STRAITS TIMES (Singapore) dated March 2022. The headline reads: "Bed crunch at Singapore hospitals: Some patients are stuck in emergency departments". Below the headline is a sub-headline: "Emergency rooms overflowing as new cases continue to rise, straining already overstuffed facilities". A photograph shows a busy emergency room with medical staff and patients. A sign in the background reads "Emergency". The page includes standard news navigation elements like Top Stories, Latest News, Discover, Singapore, Asia, Comments, Sustainability, CNA Insider, Lifestyle, Health, Letters, and +40 Sections.

(Straits Times, 2022)



PROBLEM STATEMENT



01

Time-consuming
Diagnostic Tests

02

High Costs of Tests

03

Manpower shortage

04

Scant Medical Records



TIME-CONSUMING DIAGNOSTICS TESTS

Ambulatory Electrocardiogram

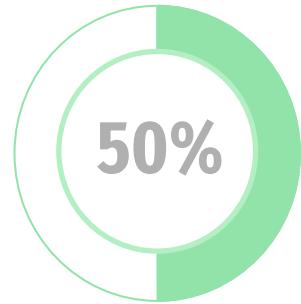


24-48 hours





TIME-CONSUMING DIAGNOSTICS TESTS



Potentially salvageable heart muscle is lost within an hour



Heart attack patients managed to get treatment in time

HIGH COSTS OF TESTS



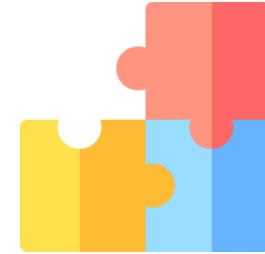
Diagnostic Test	Cost
Consultation	\$180-\$300
Electrocardiogram	\$60-\$80
Echocardiogram	\$450-\$600
Treadmill Test	\$350-\$550

Deters high-risk patients from taking health screenings or seeking treatment →

- Further strain of healthcare resources
- Higher mortality rate



MANPOWER SHORTAGE & SCANT MEDICAL RECORDS



Diversion of critical resources → Could have been better optimised

Initial assessment based on their medical history



GAPS IN EXISTING MODELS



01

High false negative rates of existing models

02

Usefulness of non-medical factors

03

Limitations of use of complex data in existing models



INCOMPLETENESS OF EXISTING MODELS

- ❖ High false negative rates

Naïve Bayes and weighted approach	2 SVM's and XGBoost	SVM and DO	XGBoost
44.5%	47.5%	43.4%	45.7%

- ❖ Life threatening: Heart disease goes undetected





USEFULNESS OF NON-MEDICAL FACTORS

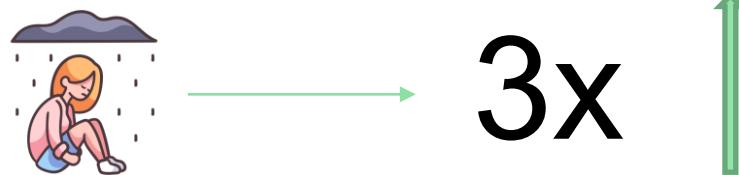
❖ Behavioural Factors

(Global case-study on modifiable behavioural risk factors association with heart disease)



❖ Psychological Factors

(Case-study on depression's association with heart disease)



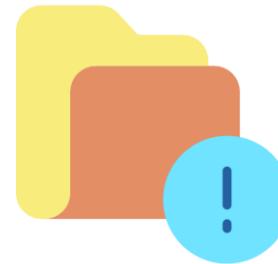


LIMITATIONS OF USE OF COMPLEX DATA IN EXISTING MODELS

- ❖ Complex data needed for existing models for prediction



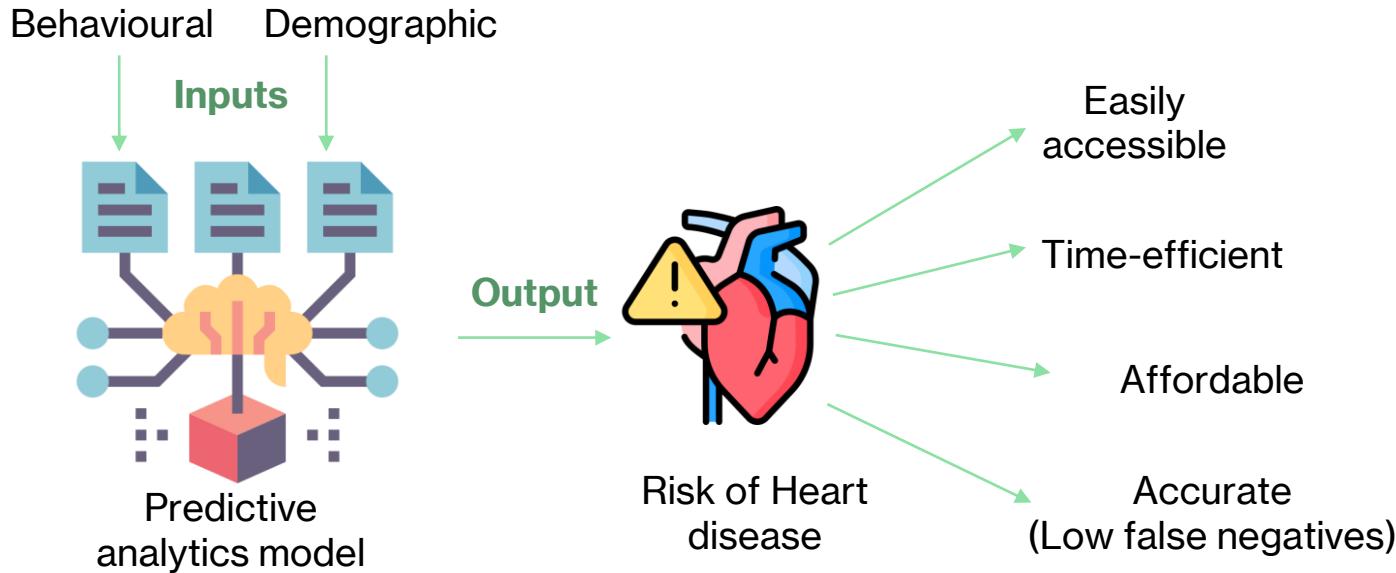
Reluctance



Data unavailable



PROJECT OBJECTIVE



Alleviate stress on NHCS through pre-diagnosis with non-medical data



DATA SOURCING

- ❖ **Aim:** Construct a model that predicts likelihood of CHD based on non-medical data
- ❖ **Need:** Wide range of behavioural, psychological & demographic factors that lead to CHD



- Public health surveys
- Behavioural, Psychological & Demographic factors



DATA CLEANING





INITIAL DATASET

Data

dt 441456 obs. of 330 variables

STATE	FMONTH	IDATE	IMONTH	IDAY	IYEAR	DISPCODE	SEQNO	PSU	CTELENUM	PVTRESID	COLGHOUS	STATERES	CELLFON3	LADULT	HUMADULT	NUMMEN	NUMWOMEN	CTELNUM
2	1	1 b'011202015'	b'01'	b'20'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	3	1.000000e+00	2.000000e+00		
3	1	1 b'02012015'	b'02'	b'01'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
4	1	1 b'01142015'	b'01'	b'14'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
5	1	1 b'01142015'	b'01'	b'14'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	3	1.000000e+00	2.000000e+00		
6	1	1 b'01142015'	b'01'	b'14'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
7	1	1 b'01052015'	b'01'	b'05'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
8	1	1 b'01142015'	b'01'	b'14'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
9	1	1 b'01132015'	b'01'	b'13'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
10	1	1 b'01302015'	b'01'	b'30'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
11	1	1 b'01222015'	b'01'	b'22'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
12	1	1 b'01162015'	b'01'	b'16'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
13	1	1 b'01202015'	b'01'	b'20'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
14	1	1 b'01202015'	b'01'	b'20'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
15	1	1 b'01202015'	b'01'	b'20'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
16	1	1 b'01142015'	b'01'	b'14'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
17	1	1 b'01042015'	b'01'	b'04'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
18	1	1 b'01202015'	b'01'	b'20'	b'2015'	1200	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
19	1	1 b'01272015'	b'01'	b'27'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	5.397605e-79	1.000000e+00		
20	1	1 b'01062015'	b'01'	b'06'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	2	1.000000e+00	1.000000e+00		
21	1	1 b'01142015'	b'01'	b'14'	b'2015'	1100	2.015e+09	2.015e+09	1	N	1	2	NA	1	1.000000e+00	5.397605e-79		

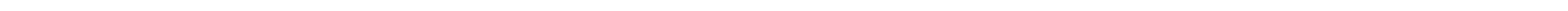


I. RECODE SPECIFIED VALUES



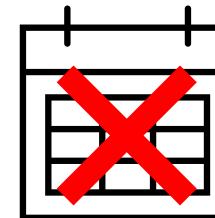
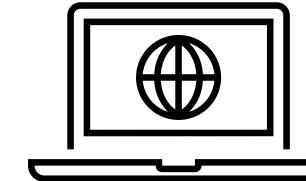
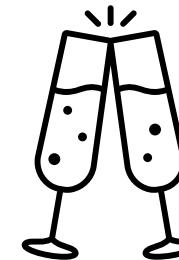
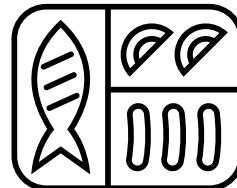
II. REMOVE DUPLICATES

THE DATASET HAD NO DUPLICATES.





III. COLUMN PRESERVATION





IV. CHOOSE 'Y' – MICHD

Ever had CHD or MI

Calculated Variables: 6.1 Calculated Variables

Type: Num

Column: 1899

SAS Variable Name: _MICHD

Prologue:

Description: Respondents that have ever reported having coronary heart disease (CHD) or myocardial infarction (MI)

Value	Value Label	Frequency	Percentage	Weighted Percentage
1	Reported having MI or CHD Notes: CVDINFR4=1 OR CVDCRHD4=1	38,633	8.83	6.42
2	Did not report having MI or CHD Notes: CVDINFR4=2 AND CVDCRHD4=2	398,881	91.17	93.58
BLANK	Not asked or Missing Notes: CVDINFR4=7, 9 OR MISSING OR CVDCRHD4=7, 9, OR MISSING	3,942		



V. CORRECT NUMERIC DATA

Computed Weight in Kilograms		Type: Num		
Calculated Variables:	Column: 1983-1987			
Prologue:				
Description:	Reported weight in kilograms	SAS Variable Name: WTKG3		
Value Value Label Frequency Percentage Weighted Percentage				
2300 - 29500	Weight in kilograms [2 implied decimal places] Notes: 0001 <= WEIGHT2 <= 650 or 9023 <= WEIGHT2 <= 9295 (non-metric WEIGHT2 value divided by 2.2046)	410,535	93.00	93.29
99999	Don't know/Refused/Missing Notes: WEIGHT2 = 7777 or 9999 or Missing or not in accepted values	30,921	7.00	6.71

WTKG3, BMI5, GRENDAY, FRUITDA1, BEANDAY : * 0.01



VI. REMOVAL OF MISSING DATA

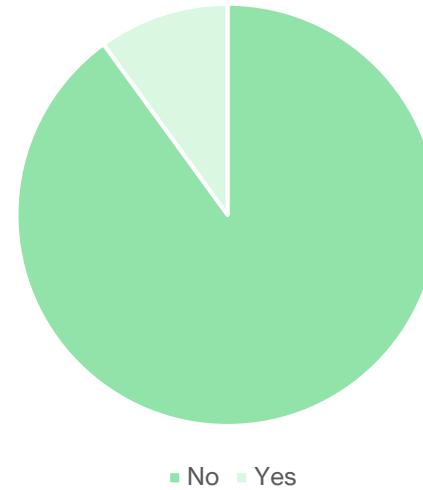
- Removed rows with **missing 'Y'**
- Removed rows with missing values in columns with **>20%** data missing (SCNTWRK1, PA1VIGM, PA1MIN)
 - No need to resort to imputation – enough data (441,456 obs.)
- Removed rows with **>20%** missing values



VII. BALANCE DATA

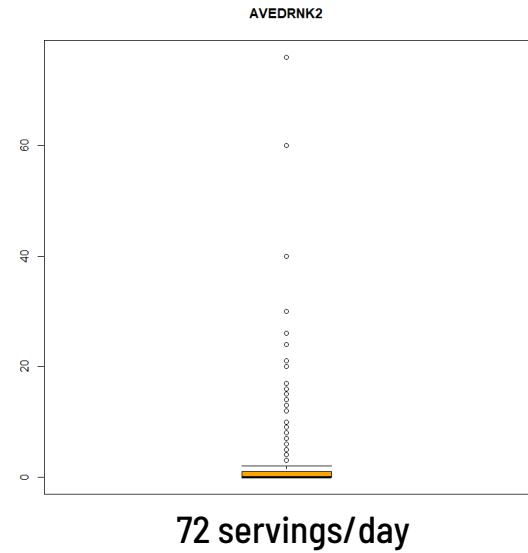
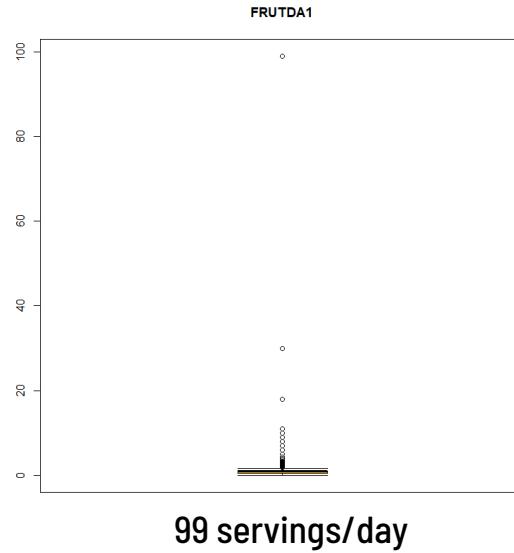
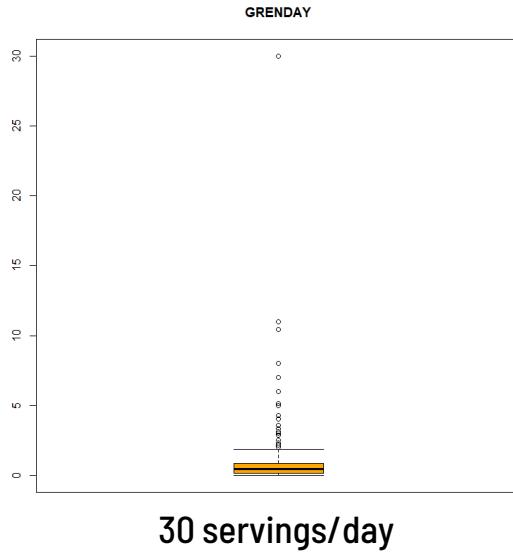
- 1:9 [Disease : No Disease] – bias towards predicting no disease
- **Undersampling**
- $1:9 \rightarrow 1:1$

Occurrence of Heart Disease (Y)



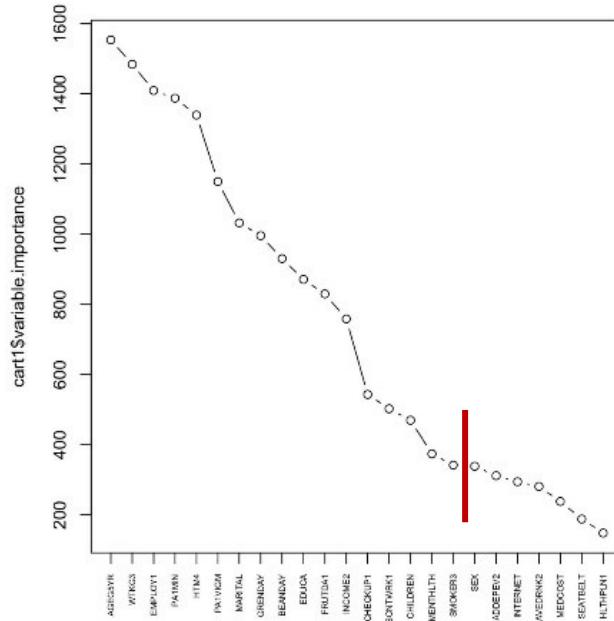


VIII. REMOVAL OF OUTLIERS





IX. REDUCE DIMENSION WITH CART



variable.importance : 27 → 18 Columns



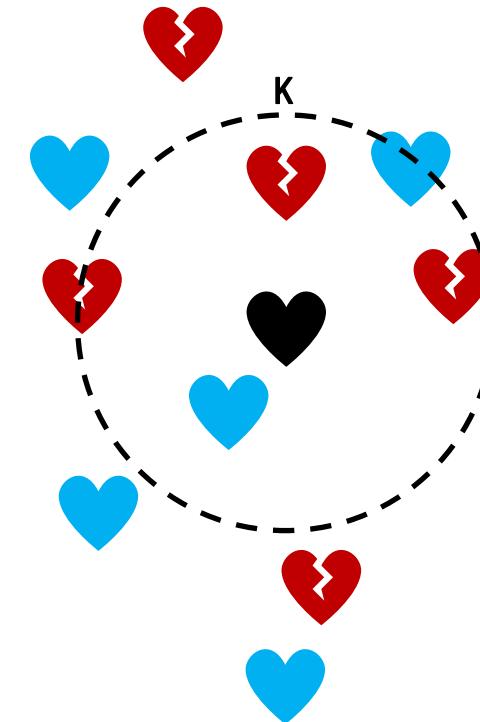
POST-CLEANING

- Dimensions: **21,214 rows, 18 columns**
- Impute data to fill missing values



DATA IMPUTATION

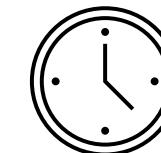
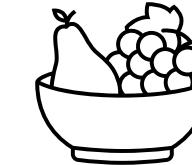
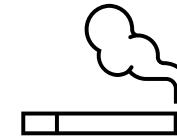
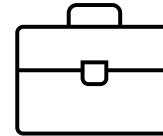
- **K-nearest neighbours (kNN)** algorithm
- More robust than regression models
 - Can impute **both numerical, categorical** variables
 - Distance-based algorithm – can impute even with missing values by using alternative columns





FINAL DATASET

- **21,214 rows, 18 columns**
- Data includes:

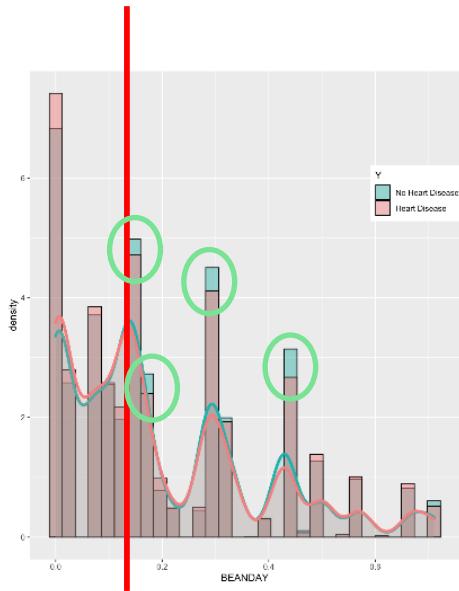


DATA EXPLORATION

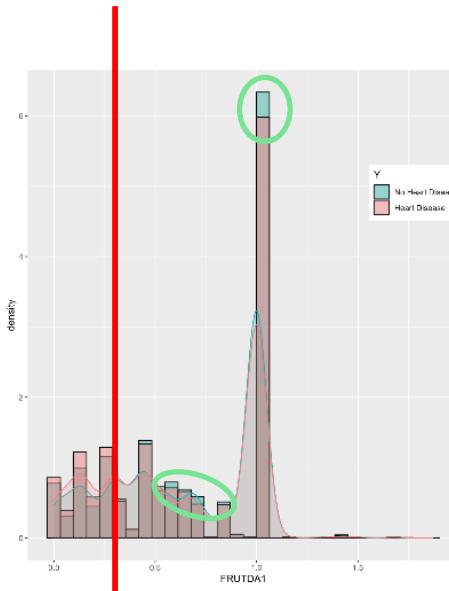




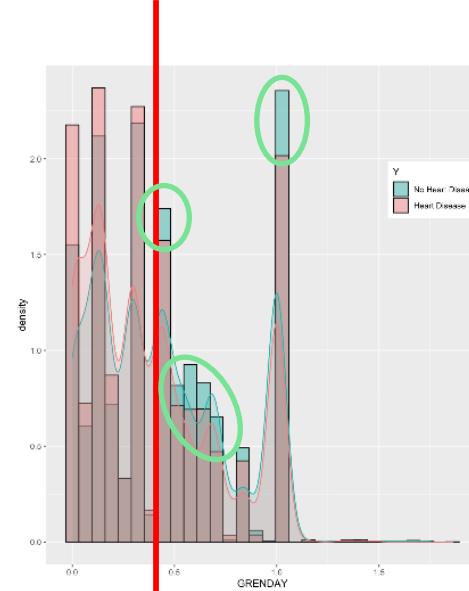
NUTRITION



DAILY BEAN INTAKE



DAILY FRUIT INTAKE

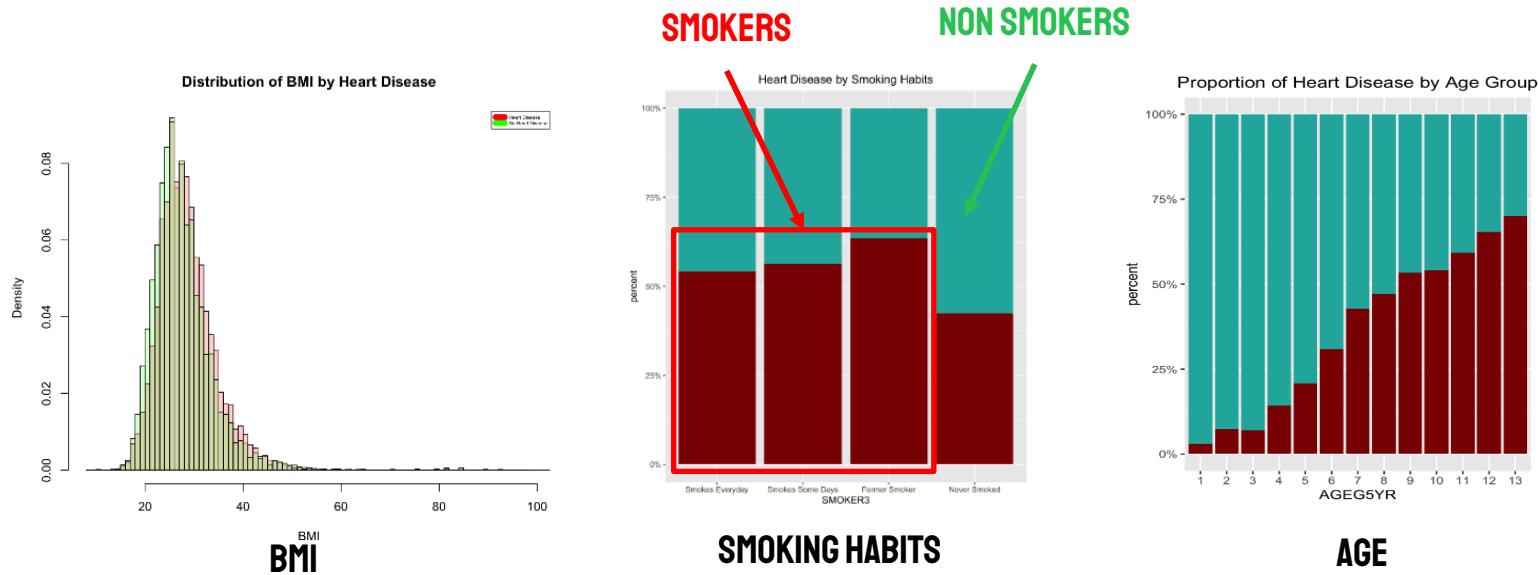


DAILY VEGETABLE INTAKE

HIGHER NUTRITION OFFERS SOME PROTECTIVE EFFECT

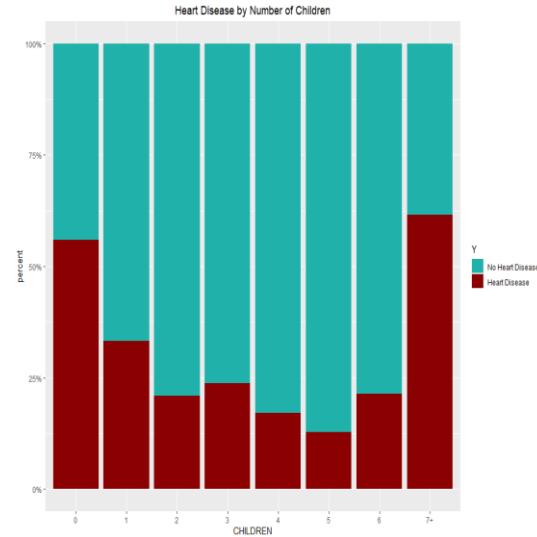


COMMON HEALTH FACTORS





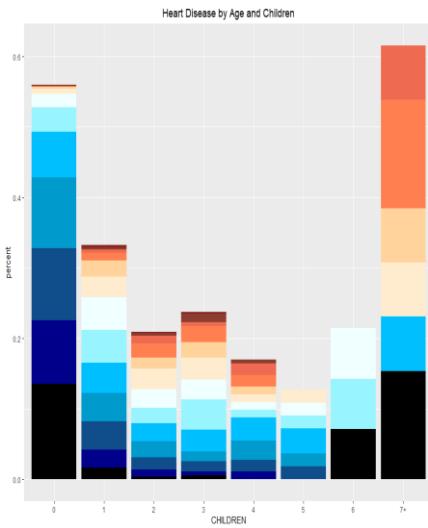
NO. CHILDREN



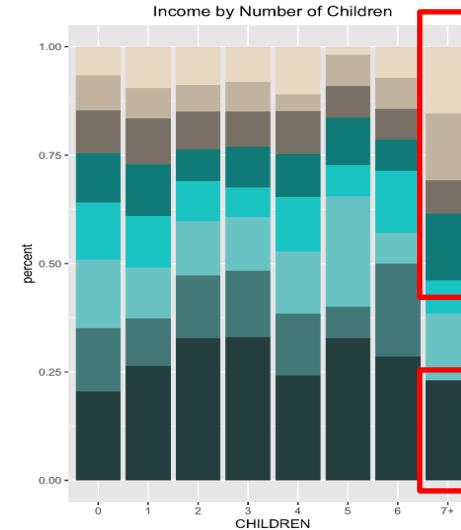
NO. CHILDREN
VS HEART DISEASE



NO. CHILDREN



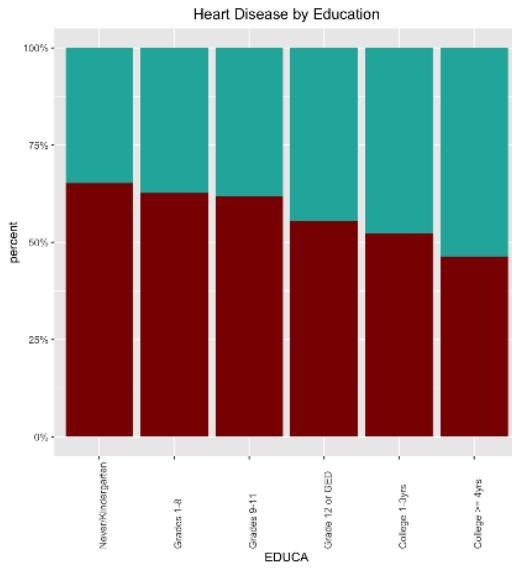
**NO. CHILDREN
VS HEART DISEASE**



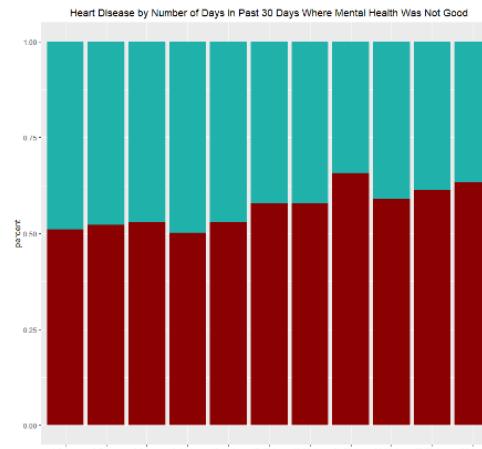
LOW INCOME GROUPS

HIGH INCOME GROUP

EDUCATION & MENTAL HEALTH

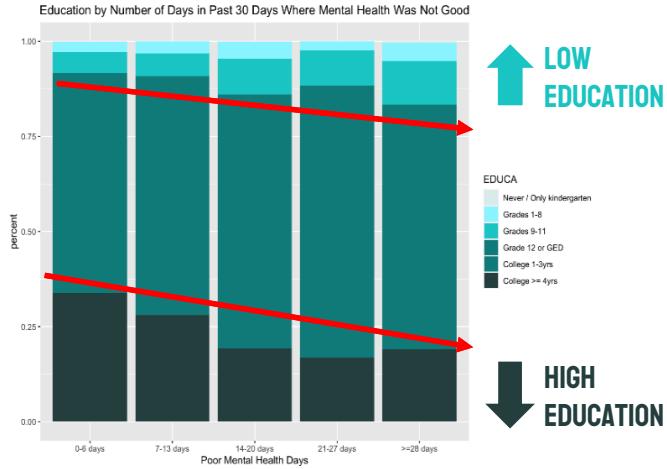


EDUCATION

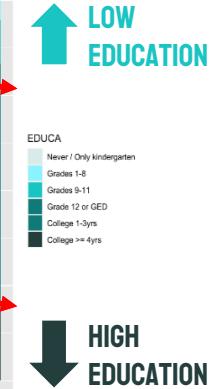


MENTAL HEALTH

"NUMBER OF DAYS IN PAST 30 DAYS WHERE MENTAL HEALTH IS NOT GOOD"

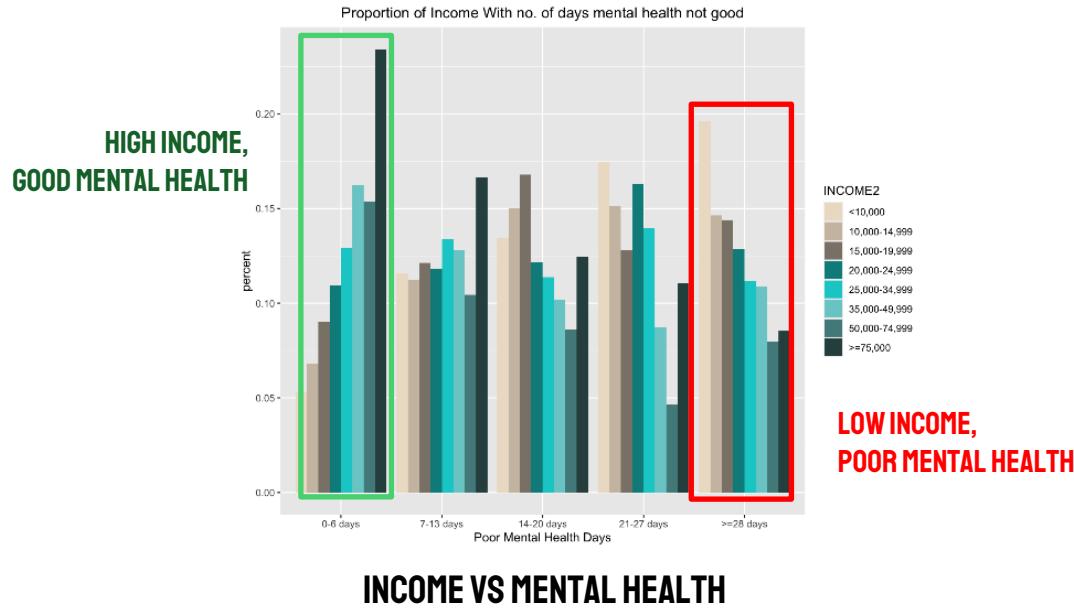


EDUCATION VS MENTAL HEALTH



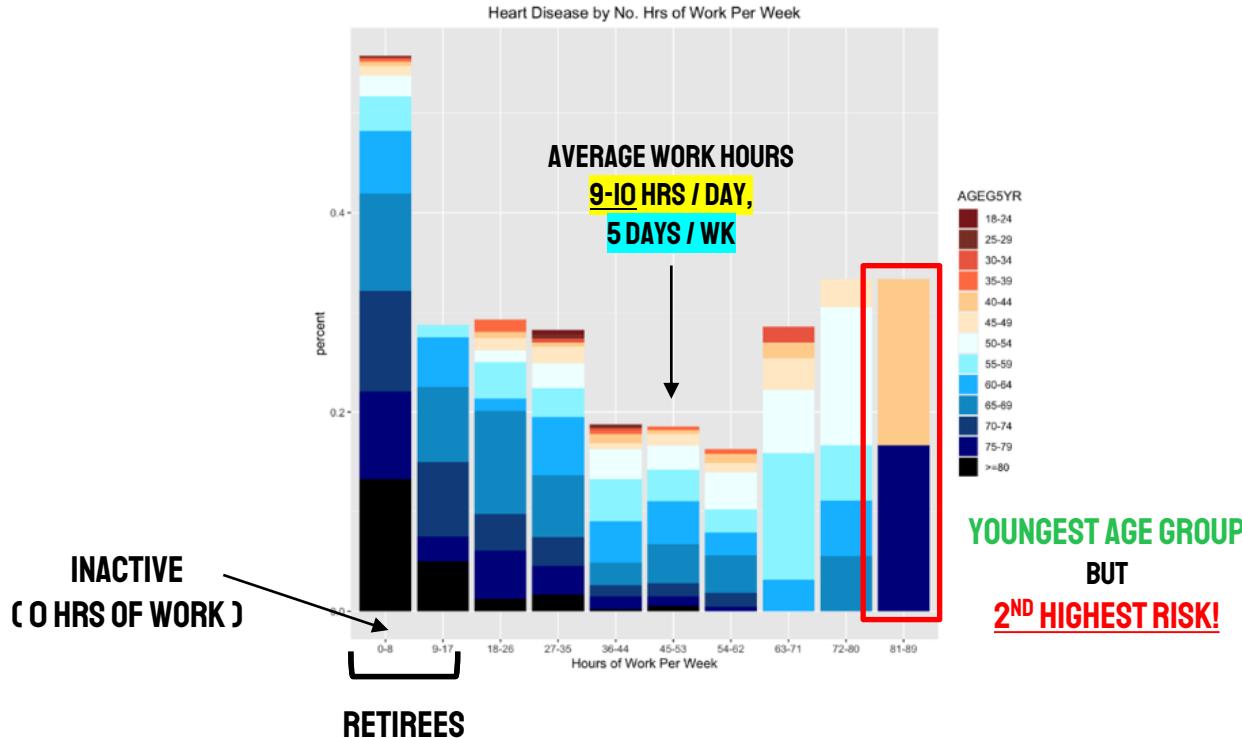


INCOME & MENTAL HEALTH



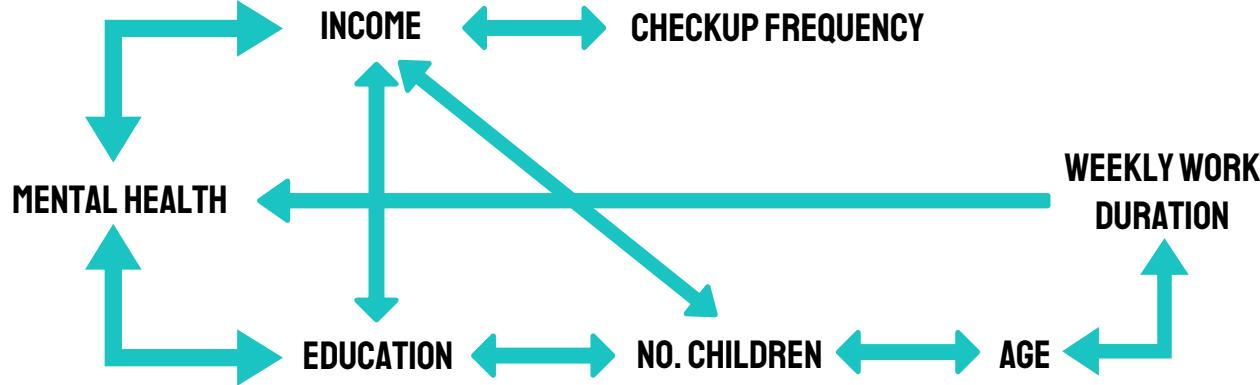


WEEKLY HOURS OF WORK





INTERCONNECTEDNESS OF VARIABLES



1

CORRELATION ≠ CAUSALITY

2

INTERDEPENDENT VARIABLES
= MULTICOLLINEARITY?



MODELLING



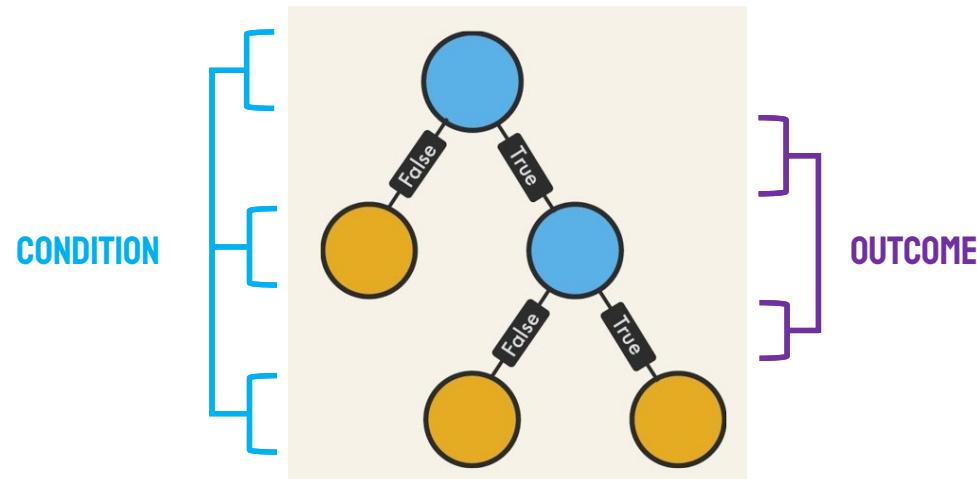
CLASSIFICATION AND
REGRESSION TREE

GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE





MODELLING



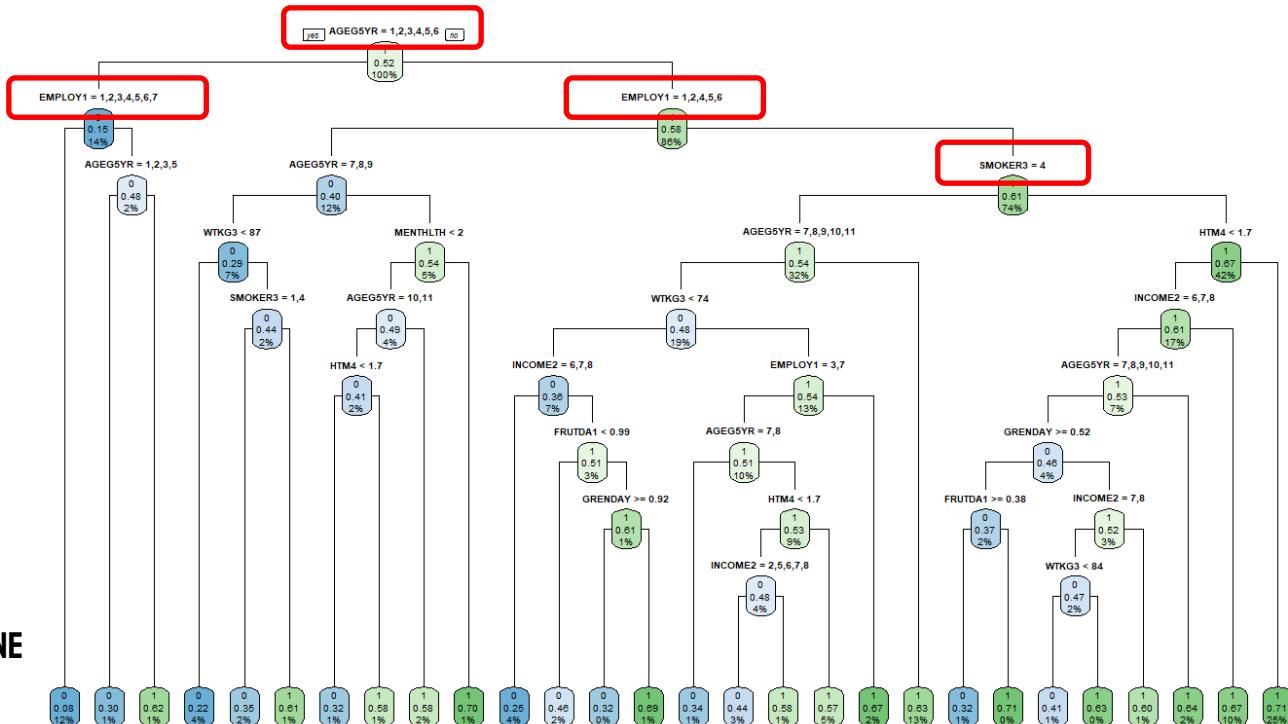
**CLASSIFICATION AND
REGRESSION TREE**

**GRADIENT BOOSTING
MACHINE**

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE





MODELLING



CLASSIFICATION AND
REGRESSION TREE



GRADIENT BOOSTING
MACHINE



LOGISTIC REGRESSION



NEURAL NETWORK



SUPPORT VECTOR MACHINE



INPUT
DATA



PREDICTIONS
(y)

RESIDUALS₁



INPUT
DATA



PREDICTIONS
(RESIDUALS₁)

RESIDUALS₂

⋮

⋮

⋮



COMBINE WEAK LEARNERS



NEW TREE CORRECTS THE ERRORS OF PREVIOUS TREE



MODELLING



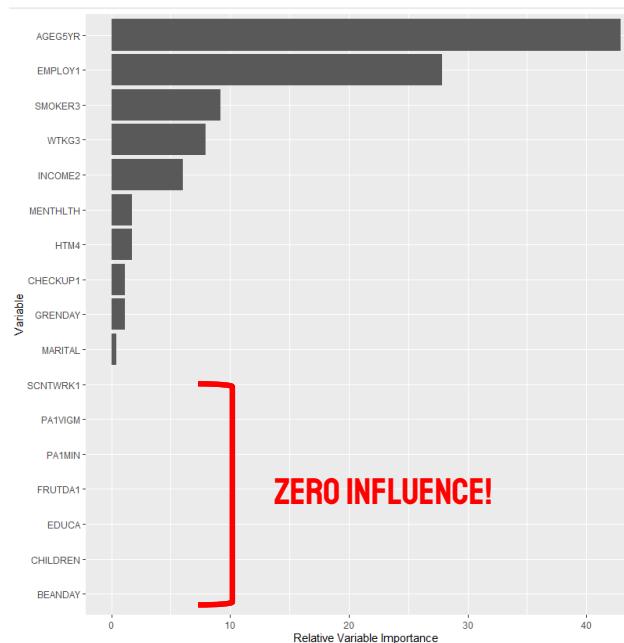
CLASSIFICATION AND
REGRESSION TREE

GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE



ZERO INFLUENCE!



MODELLING



CLASSIFICATION AND
REGRESSION TREE

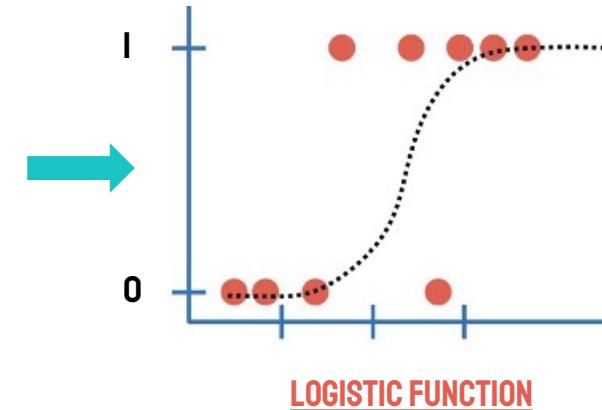
GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE

$$Y_i = \beta_0 + \beta_1 X_i$$





MODELLING



CLASSIFICATION AND
REGRESSION TREE

GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE

Backward Stepwise Logistic Regression

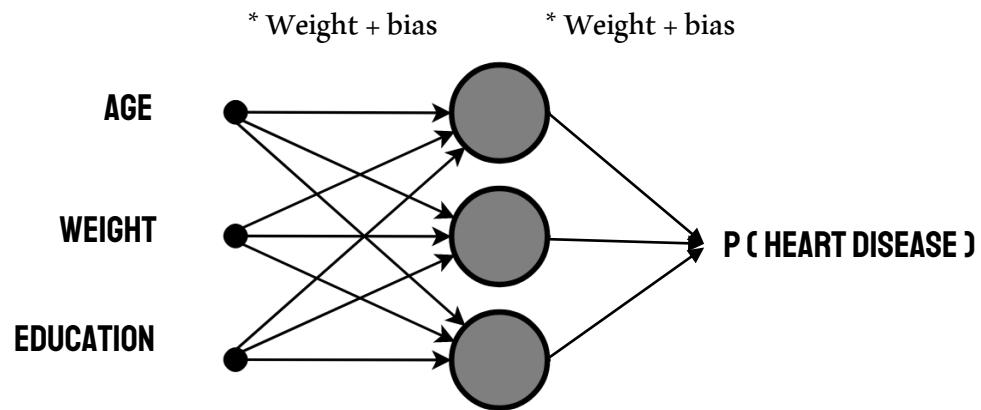
- Variables removed from the original model step-by-step
- Prevent overfitting
- Reduce multicollinearity

WTKG3
HTM4
AGEG5YR
EMPLOY1
GRENDAY
MARITAL
INCOME2
PA1VIGM
CHECKUP1
MENTHLTH
SMOKER3

PA1MIN
BEANDAY
FRUTDA1
EDUCA
CHILDREN
SCNTWRK1



MODELLING





MODELLING



CLASSIFICATION AND
REGRESSION TREE

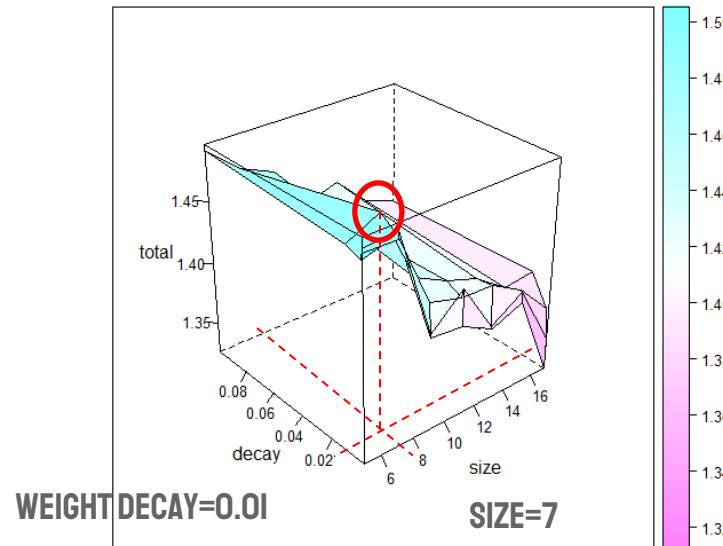
GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE

GRID SEARCH





MODELLING



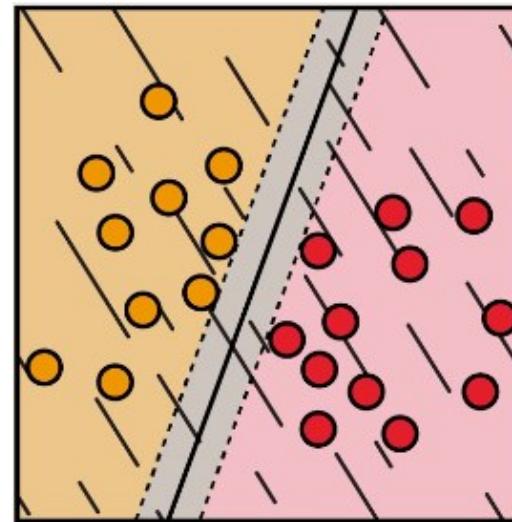
CLASSIFICATION AND
REGRESSION TREE

GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE





MODELLING



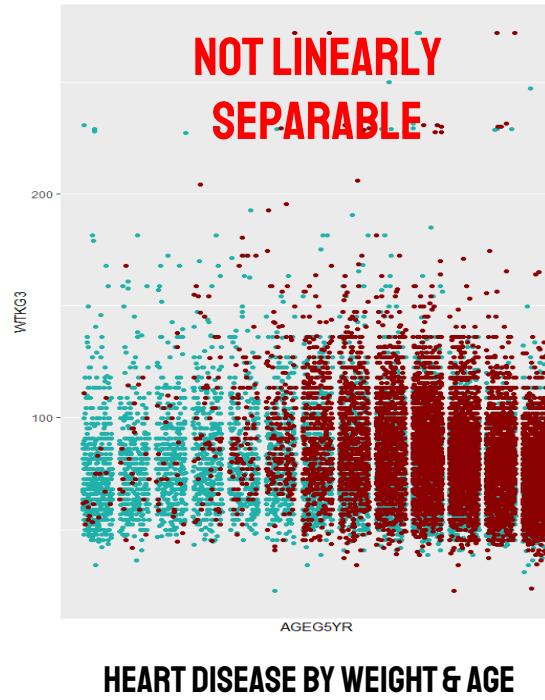
CLASSIFICATION AND
REGRESSION TREE

GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE





MODELLING



CLASSIFICATION AND
REGRESSION TREE

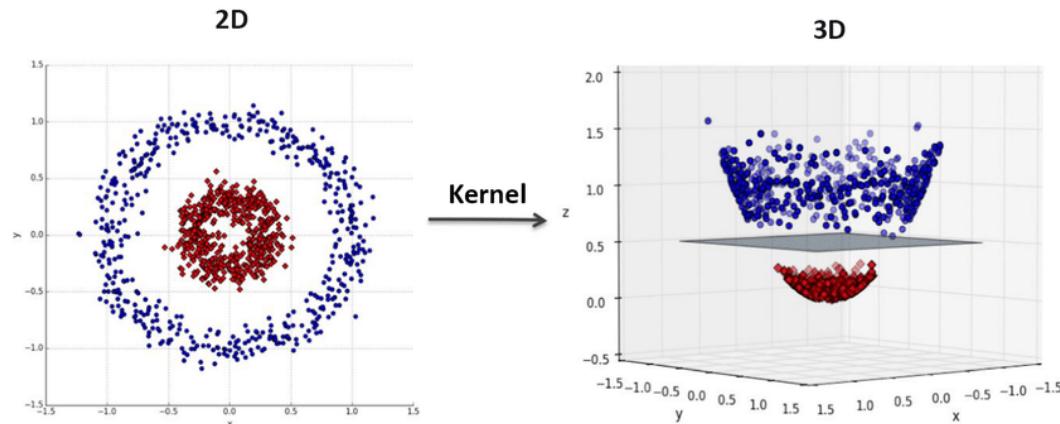
GRADIENT BOOSTING
MACHINE

LOGISTIC REGRESSION

NEURAL NETWORK

SUPPORT VECTOR MACHINE

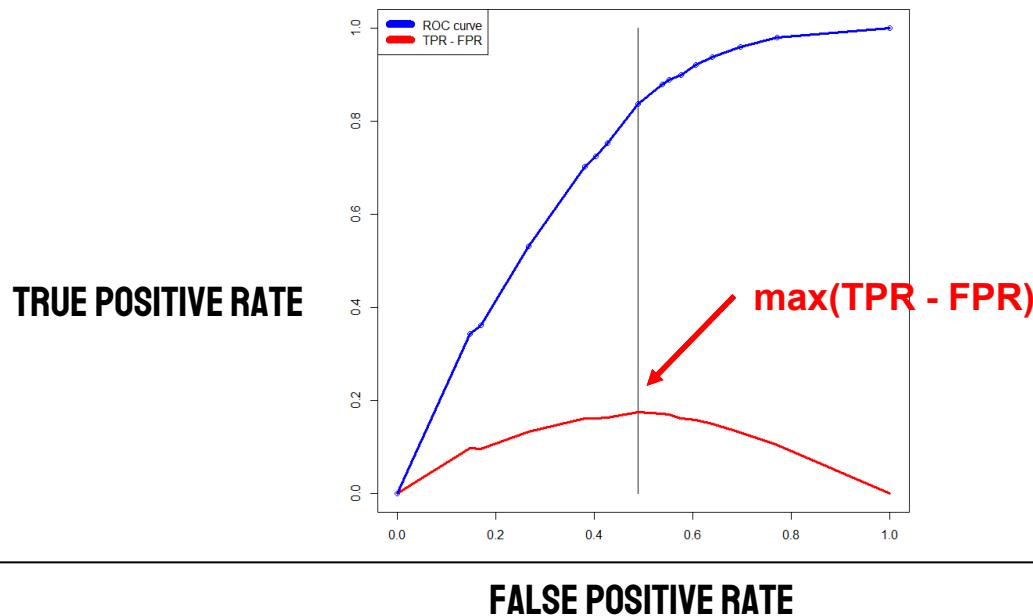
KERNEL TRICK





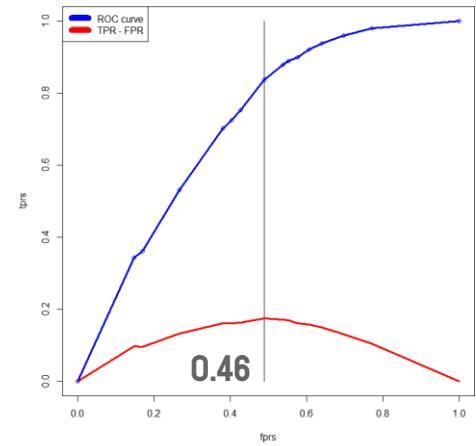
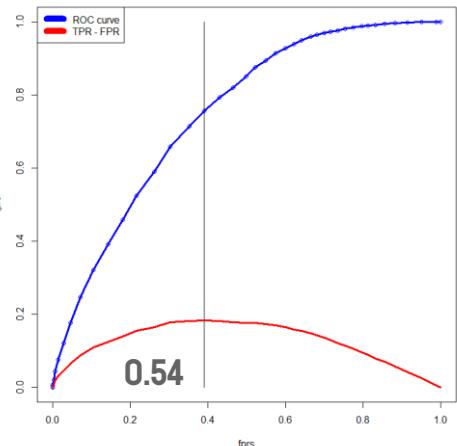
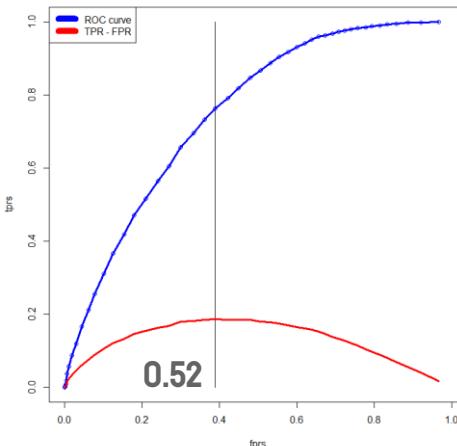
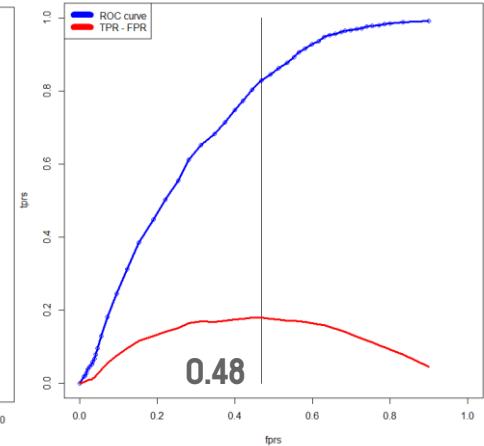
THRESHOLD MOVING

Receiver Operating Characteristic (ROC) Curve Method





THRESHOLD MOVING

CART**GBM****LOGISTIC REGRESSION****NEURAL NETWORK**



MODEL EVALUATION

Validation Accuracy



TRUE POSITIVES
ACTUAL POSITIVES

Recall



TRUE POSITIVES
POSITIVE PREDICTIONS

Precision

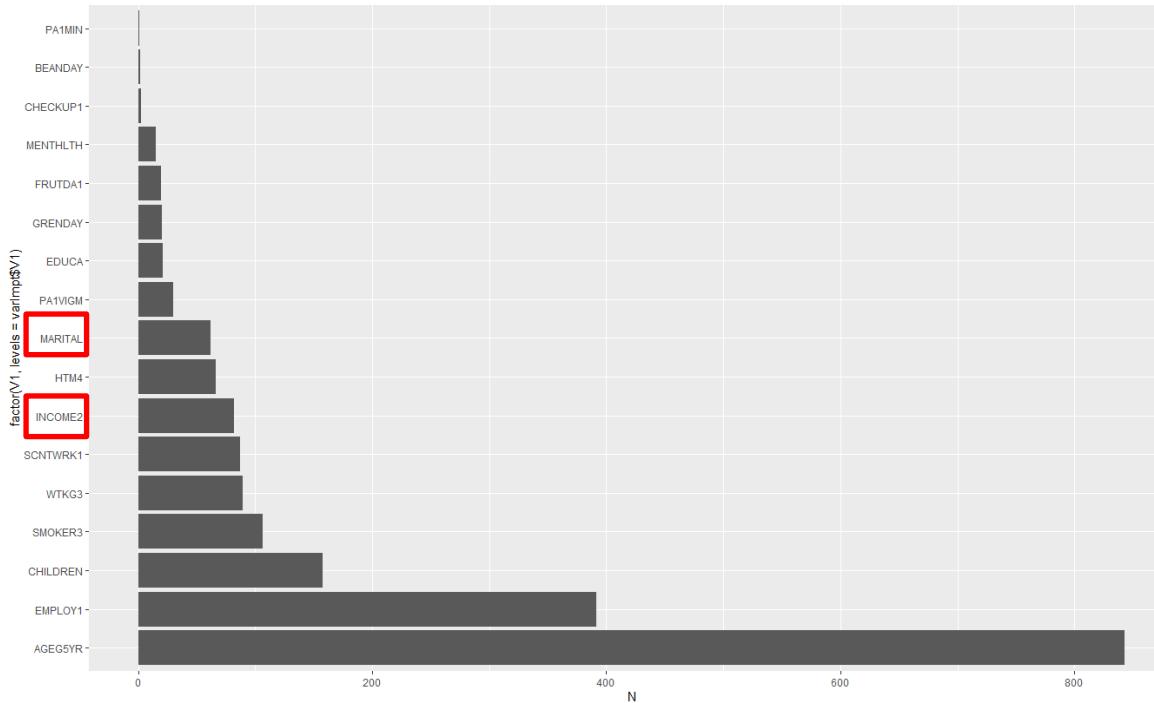
HARMONIC MEAN OF
RECALL & PRECISION

F1 Score



MODEL INSIGHTS

CART'S VARIABLE IMPORTANCE





MODEL INSIGHTS

CHI-SQUARE HYPOTHESIS

Chi-square test of categorical association

Variables: **Y, MARITAL**

Hypotheses:

null: variables are independent of one another
alternative: some contingency exists between variables

Observed contingency table:

MARITAL

Y	1	2	3	4	5	6
0	5231	1320	1679	153	1519	224
1	5469	1712	2888	229	670	120

Expected contingency table under the null hypothesis:

MARITAL

Y	1	2	3	4	5	6
0	5107	1447	2180	182	1045	164
1	5593	1585	2387	200	1144	180

Test results:

X-squared statistic: 709.708

degrees of freedom: 5

p-value: <.001

Other information:

estimated effect size (Cramer's v): 0.183

Chi-square test of categorical association

Variables: **Y, INCOME2**

Hypotheses:

null: variables are independent of one another
alternative: some contingency exists between variables

Observed contingency table:

INCOME2

Y	1	2	3	4	5	6	7	8
0	660	643	821	974	1181	1517	1516	2814
1	805	1033	1270	1400	1531	1755	1518	1776

Expected contingency table under the null hypothesis:

INCOME2

Y	1	2	3	4	5	6	7	8
0	699	800	998	1133	1295	1562	1448	2191
1	766	876	1093	1241	1417	1710	1586	2399

Test results:

X-squared statistic: 532.651

degrees of freedom: 7

p-value: <.001

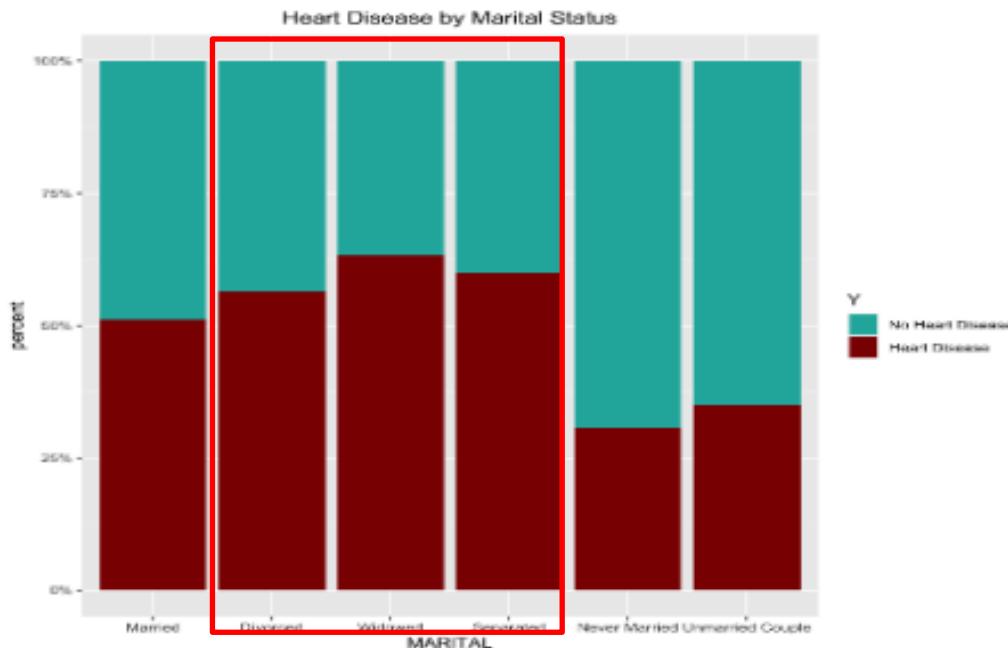
Other information:

estimated effect size (Cramer's v): 0.158



MODEL INSIGHTS

MARITAL

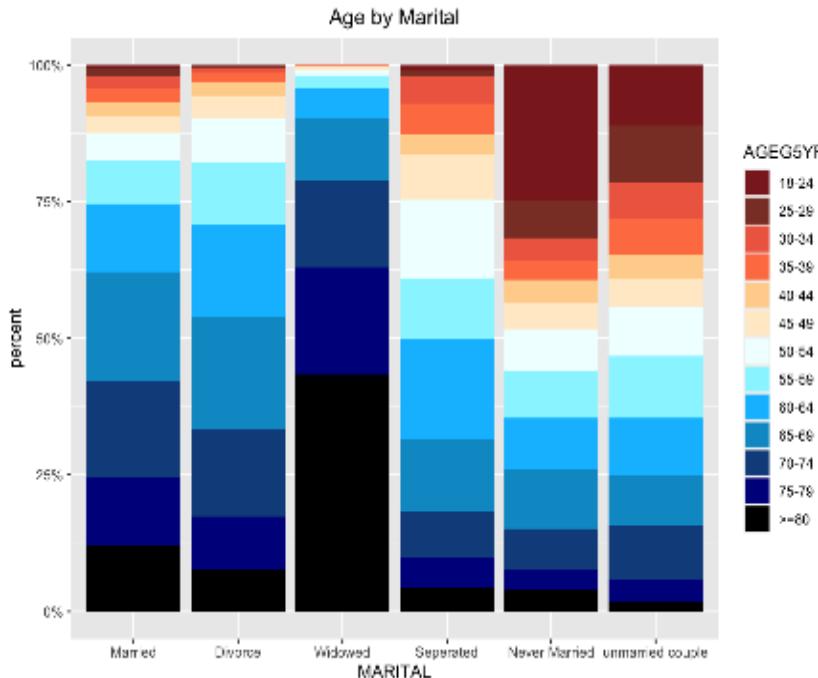


- Higher percentage of **"Divorced"**, **"Widowed"**, **"Separated"** who contracted heart disease



MODEL INSIGHTS

MARITAL BY AGE

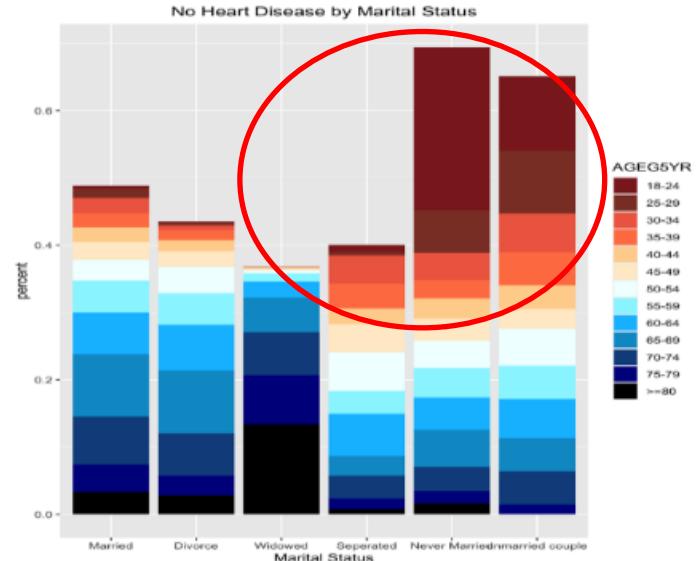
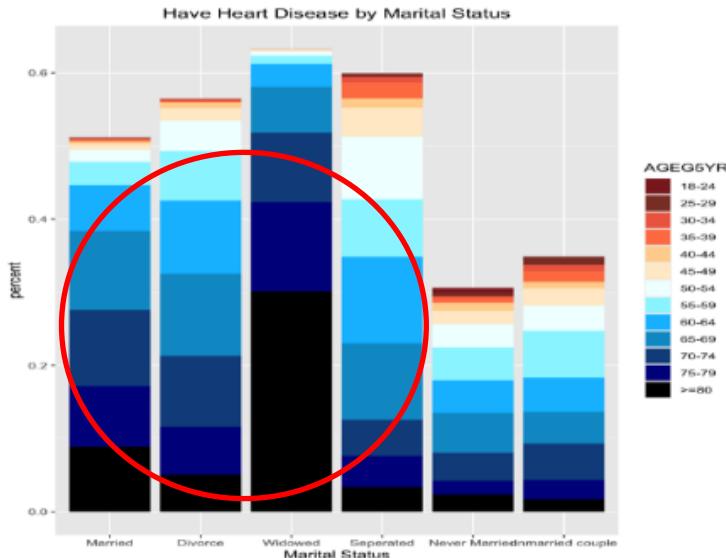


- **Respondents under 40** tend to fall under “non-married” and “unmarried couple”
- **Widowers** have the largest percentage of seniors



MODEL INSIGHTS

MARITAL BY AGE

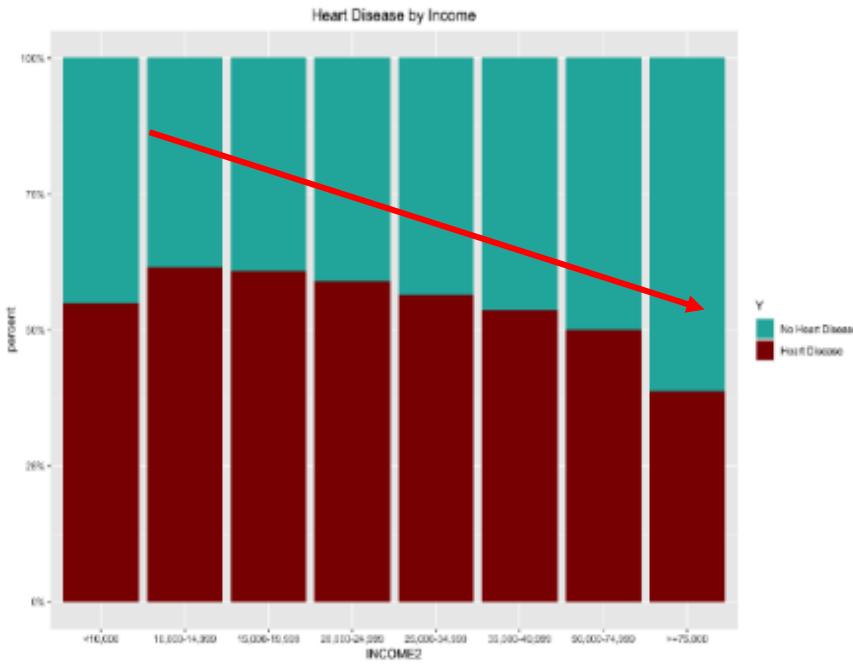


- **Age** factor in Marital Status may be the underlying factor in the prediction of heart disease



MODEL INSIGHTS

INCOME

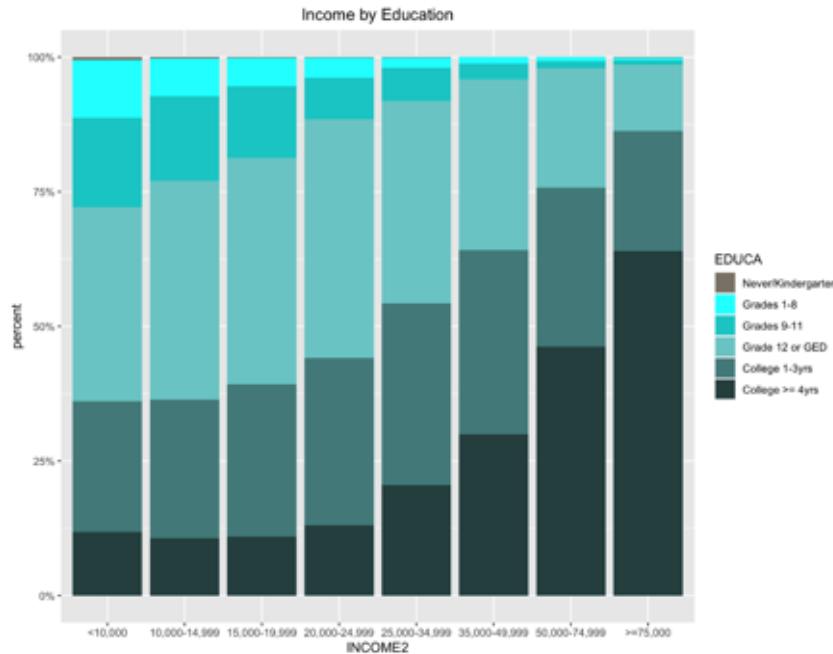


- As income increases, the risk of contracting heart disease decreases



MODEL INSIGHTS

INCOME BY EDUCATION

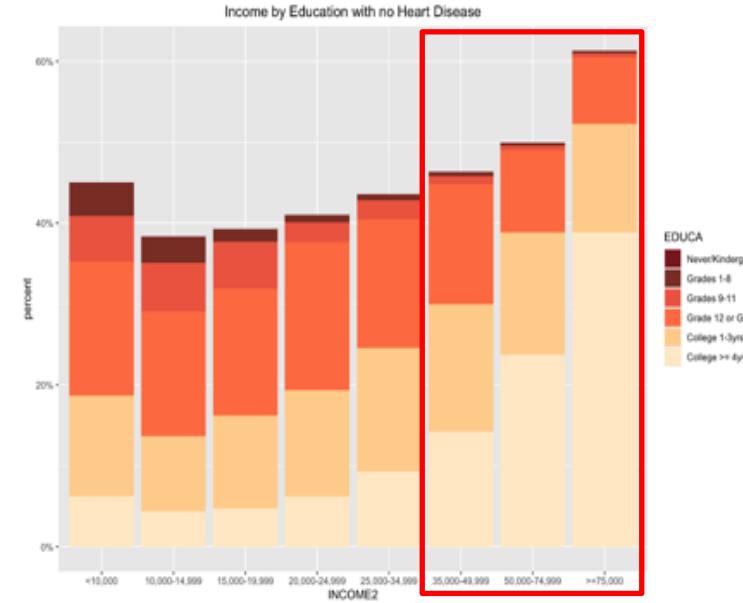
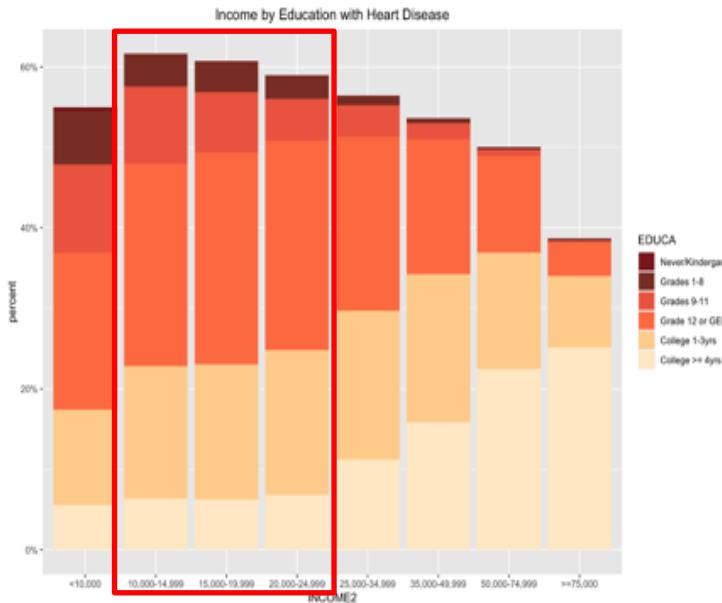


- Income has a positive correlation to the **level of education** one attains



MODEL INSIGHTS

INCOME BY EDUCATION



- Correlation between Education -> Income -> Heart Disease



MODEL INSIGHTS

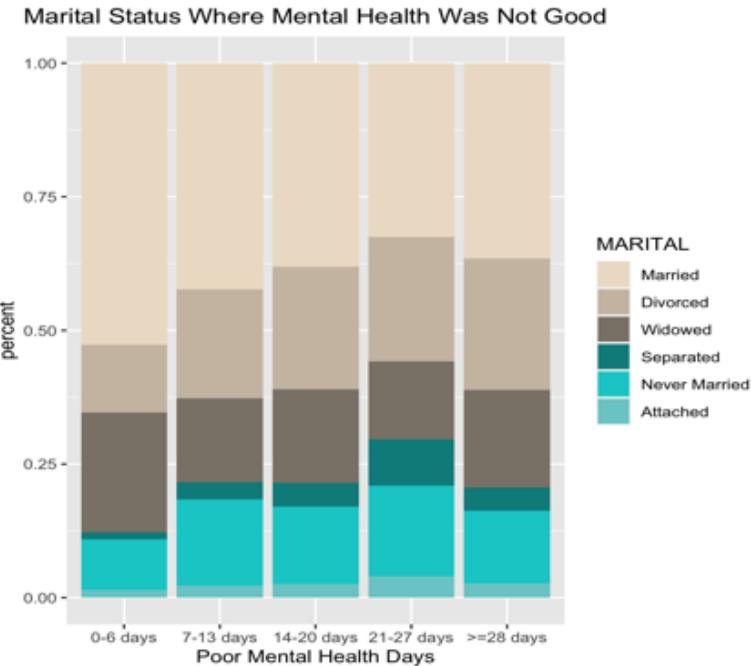
MENTAL HEALTH

- Relatedness exist between marital status and heart disease through a syndrome named "**broken heart**" (AHA, 2022)
- Studies have shown that **low income and work stress** are researched to be linked to heart disease (ECS, 2019)



MODEL INSIGHTS

MARITAL BY MENTAL HEALTH

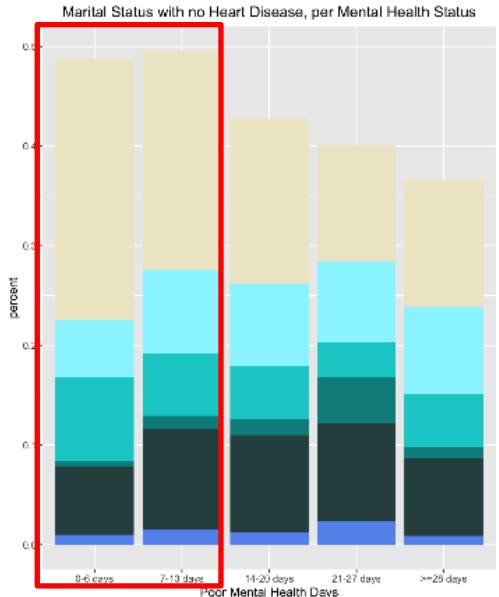
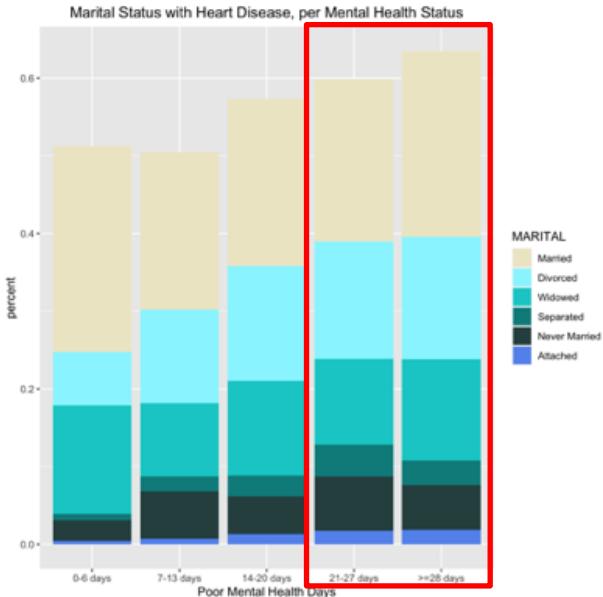


- As no. of poor mental health days decrease, ratio of married to divorce and separate increase



MODEL INSIGHTS

MARITAL BY MENTAL HEALTH

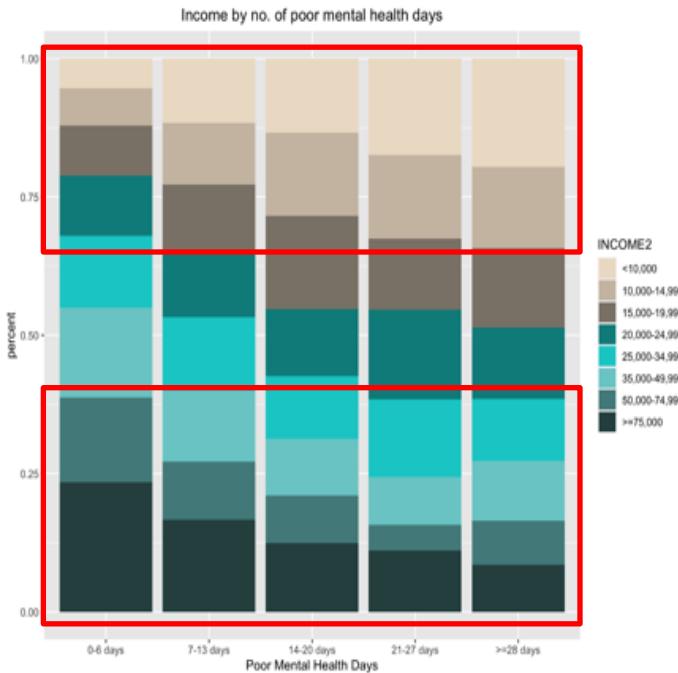


- Correlation between Marital Status → Mental Health → Heart Disease



MODEL INSIGHTS

INCOME BY MENTAL HEALTH

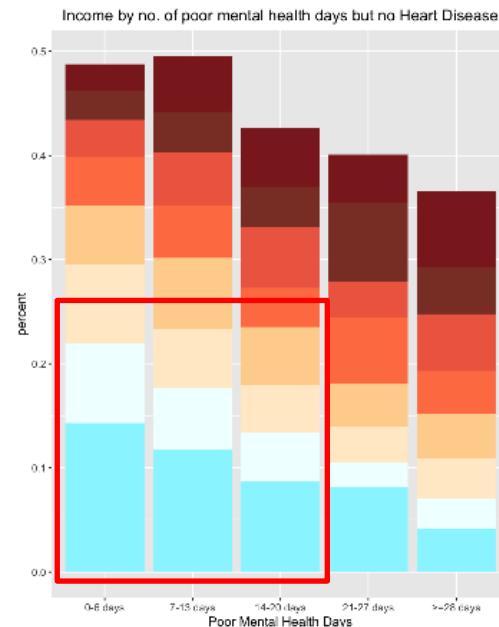
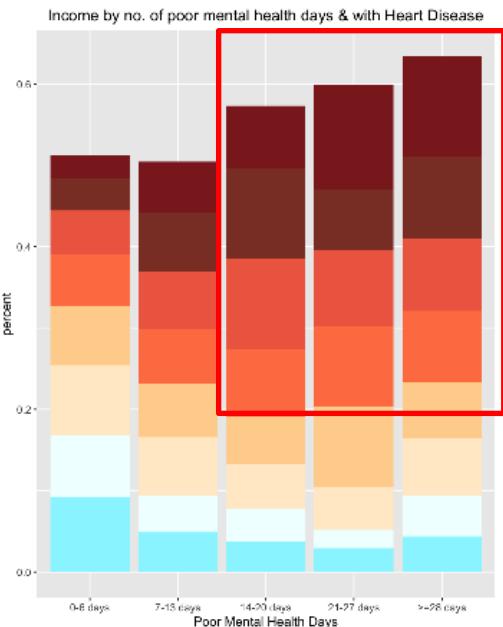


- As income increases, the no. of poor mental health days drop



MODEL INSIGHTS

INCOME BY MENTAL HEALTH



- Correlation between Income-> Mental Health -> Heart Disease



MODEL INSIGHTS

MENTAL HEALTH

- Correlation exists between
 - Income/Marital status
 - Mental Health
 - Heart disease
- More research has to be done
- More emphasis must be placed on mental health



HUMAN TENDENCY

4 COMMON BELIEFS



"I am healthy, so health screening is unnecessary"



"Going for health screening is a hassle."



"I have done my health screening 5 years ago, why must I do it again?"



"Health screening is expensive."



PROPOSED BUSINESS SOLUTION

Prevention is **KEY** as first-time events may be fatal





IS THIS APP A REPLACEMENT TO THE EXISTING MEDICAL TESTS?

NO

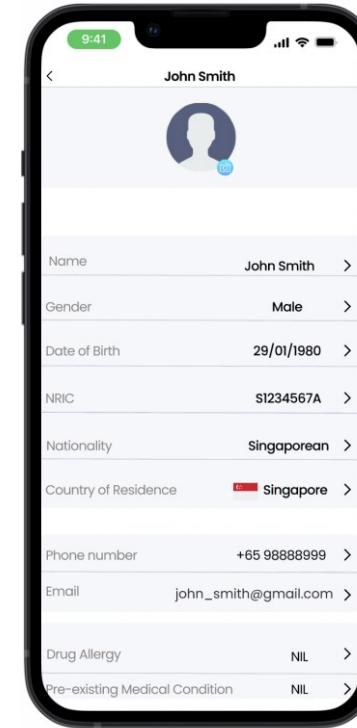
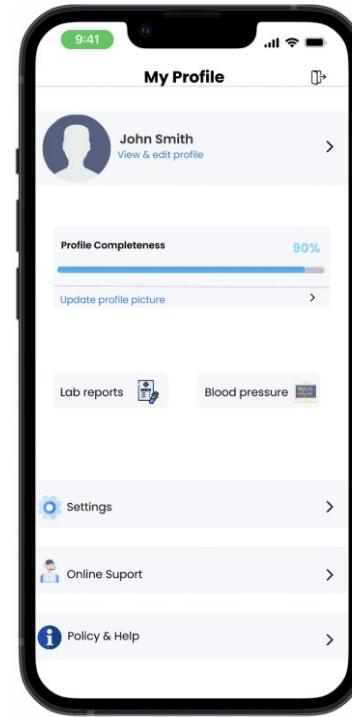
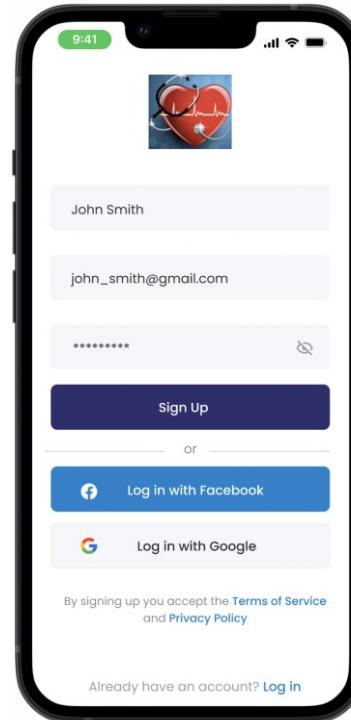
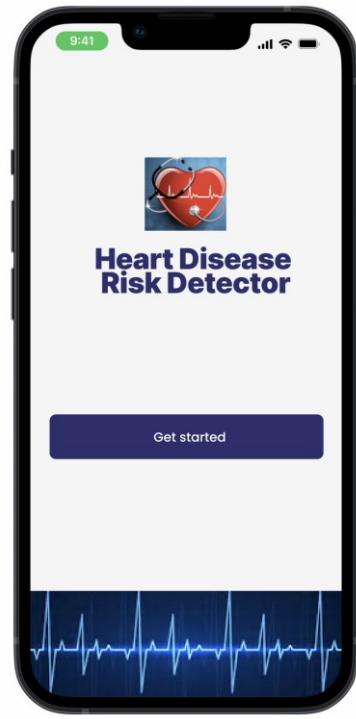
The heart disease Risk Detector app is a supplement to the existing medical tests



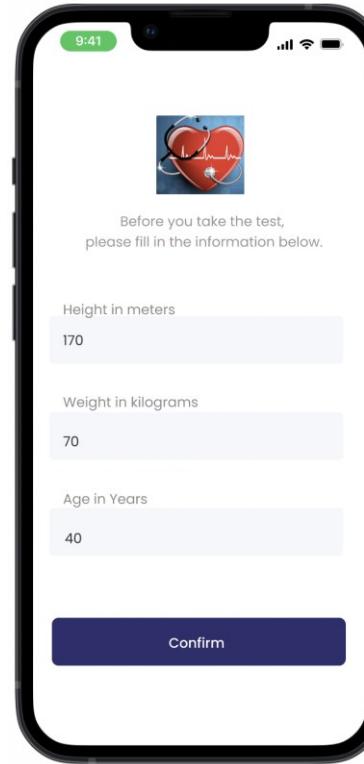
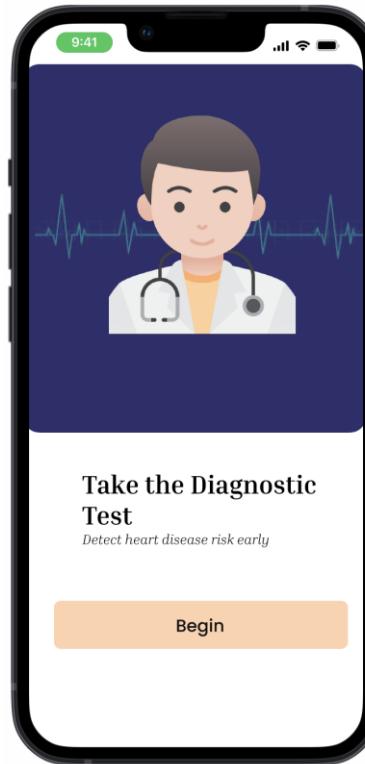
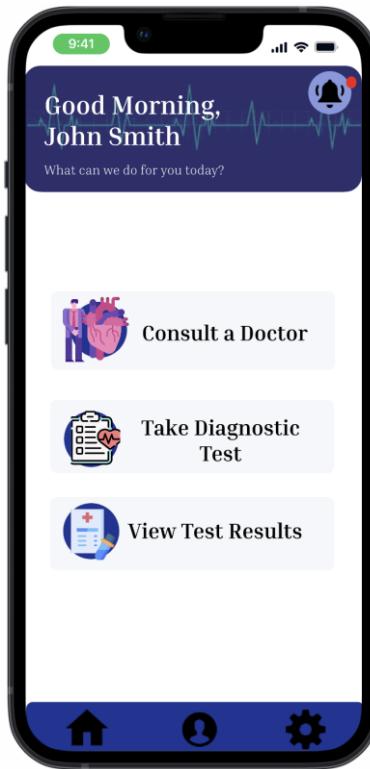
HEART DISEASE RISK DETECTOR APP



First, users create an account



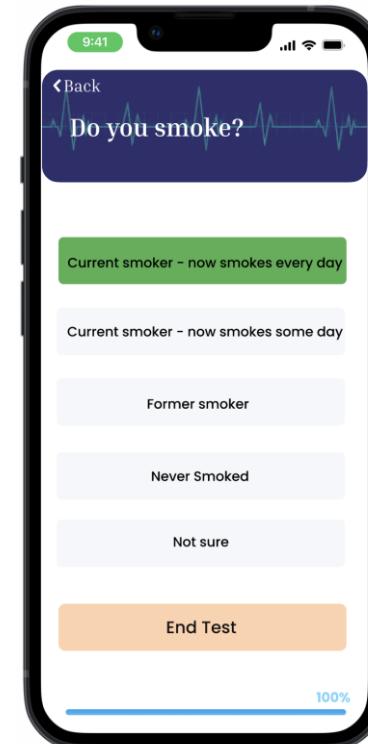
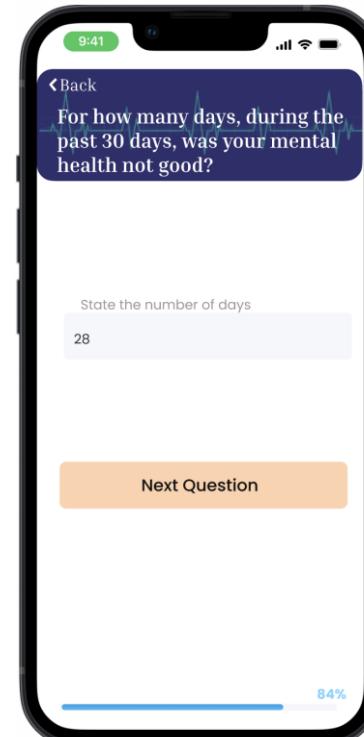
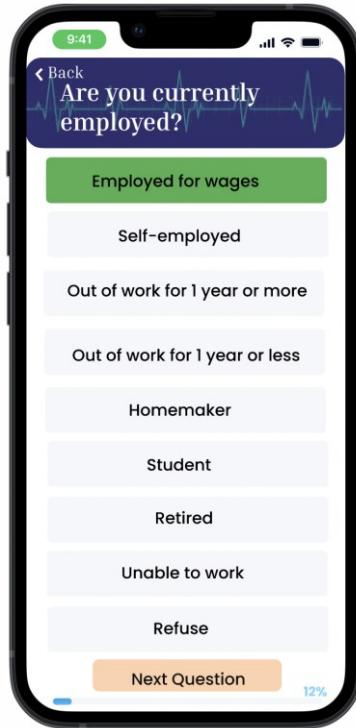
HEART DISEASE RISK DETECTOR APP



HEART DISEASE RISK DETECTOR APP



Users take the given test, providing key non-medical data

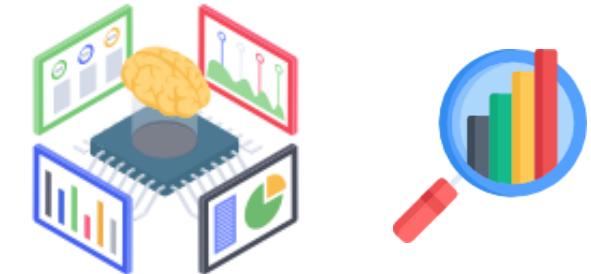




HEART DISEASE RISK DETECTOR APP



Data collected from the users is run through the CART model





HEART DISEASE RISK DETECTOR APP



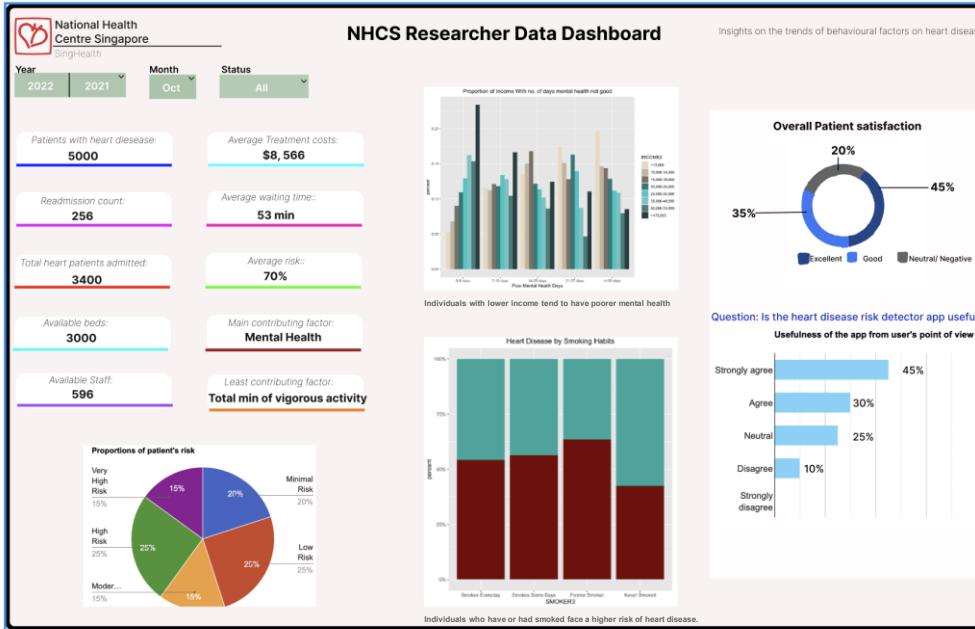
0 to 15% : Minimal risk
16 to 30%: Low risk
31 to 50%: Moderate risk
51 to 75%: High risk
76 to 100%: Very high risk



Results of the test will be displayed based on the CART model and patients should take the necessary actions accordingly.



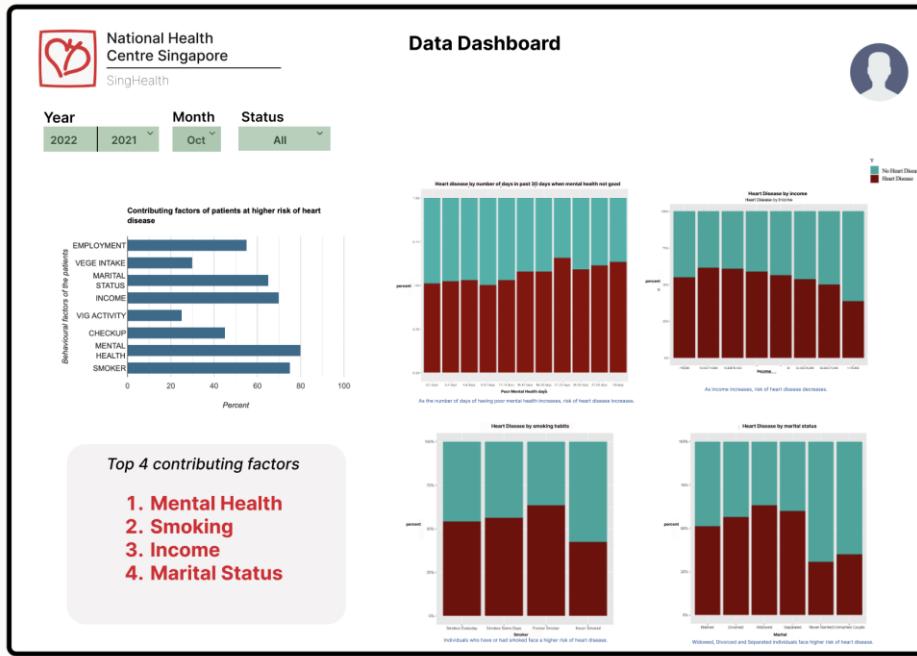
NHCS DASHBOARD MOCK-UP



NHCS Researcher's Data Dashboard



APP USER'S DASHBOARD MOCK-UP



User interface



VALUE PROPOSITION

Valued benefits over existing models:

1. Efficiency and cost effective
2. Comprehensible and user friendly
3. Alleviates stress on NHCS and the medical professionals



VALUE PROPOSITION



1. Efficiency and cost effectiveness

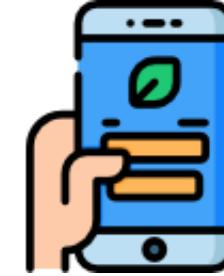


ALL NEW
HEART
DISEASE RISK
DETECTOR

VALUE PROPOSITION



2. Comprehensible and user-friendly



**ALL NEW
HEART
DISEASE RISK
DETECTOR**

VALUE PROPOSITION



3. Alleviates stress on NHCS and doctors



**ALL NEW
HEART
DISEASE RISK
DETECTOR**

LIMITATIONS



1. Success of our app
2. Presence of missing values in the dataset
3. Our app may pose difficulty in terms of usage for elderlies



**ALL NEW
HEART
DISEASE RISK
DETECTOR**

Increased predictive accuracy by up to 50% by taking into account behavioural factors

Prevention is KEY as first-time events may be fatal

CONCLUSION



In order to make well-informed and timely decisions, NHCS could implement the Heart Disease Risk Detector app which allows users to self-check their risk of heart disease.



THANK YOU!

