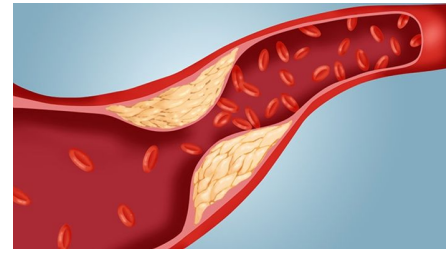# Predicting Cholesterol Without a Clinical Test

**Zehui Lin, Dylan Mendonca, Shimona Narang, Esmond Tang**
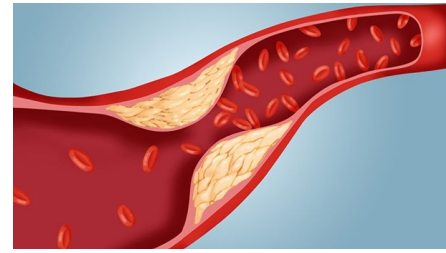
April 6, 2020

# Why Cholesterol?



➢    Cholesterol is a 'waxy' substance produced by the body and found in foods
➢    Levels ≥240 mg/dL are considered borderline/high and increases risk of heart disease, stroke, etc.

https://health.costhelper.com/cholesterol-test.html#extres1
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5485409/

# Why Cholesterol?


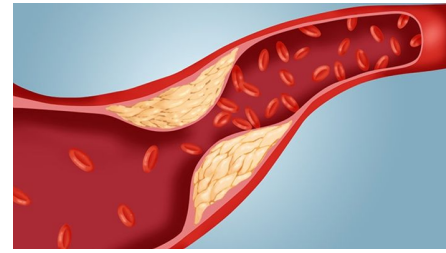
➢ Cholesterol is a 'waxy' substance produced by the body and found in foods
➢ Levels ≥240 mg/dL are considered borderline/high and increase risk of heart disease, stroke, etc.

## $25

Per clinical lab test to check cholesterol level

Expensive for people in developing nations!

https://health.costhelper.com/cholesterol-test.html#extres1
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5485409/

# Why Cholesterol?



➢ Cholesterol is a 'waxy' substance produced by the body and found in foods
➢ Levels ≥240 mg/dL are considered borderline/high and increases risk of heart disease, stroke, etc.

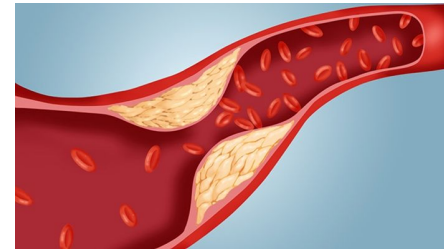## $25

Per clinical lab test to check cholesterol level

Expensive for people in developing nations!

## ~27%

Of urban pop. has high level of cholesterol (>200 mg/dL)

This number is growing!

# Why Cholesterol?



➢ Cholesterol is a 'waxy' substance produced by the body and found in foods
➢ Levels ≥240 mg/dL are considered borderline/high and increases risk of heart disease, stroke, etc.

## $25

Per clinical lab test to check cholesterol level

Expensive for people in developing nations!

## ~27%

Of urban pop. has high level of cholesterol (>200 mg/dL)

This number is growing!

## ~82%

Were **not aware** that they had high cholesterol

That's not good...

*Cholesterol Statistics for the Indian population*

# Supervised ML Algorithms Can Help

**Problem Statement:**

Develop a supervised ML model to predict whether a person has a high total cholesterol level, without having to go to the clinic

# Supervised ML Algorithms Can Help

**Primary Goal:**

Develop a supervised ML model to predict whether a person has a high total cholesterol level, without having to go to the clinic

Inputs

**National Health and Nutrition Survey (NHANES) Data**
- Demographics
- Diet
- Questionnaire
- Examinations

# Supervised ML Algorithms Can Help

**Primary Goal:**

Develop a supervised ML model to predict whether a person has a high total cholesterol level, without having to go to the clinic

| Inputs | Supervised ML Algorithms |
|---|---|

**National Health and Nutrition Survey (NHANES) Data**
- Demographics
- Diet
- Questionnaire
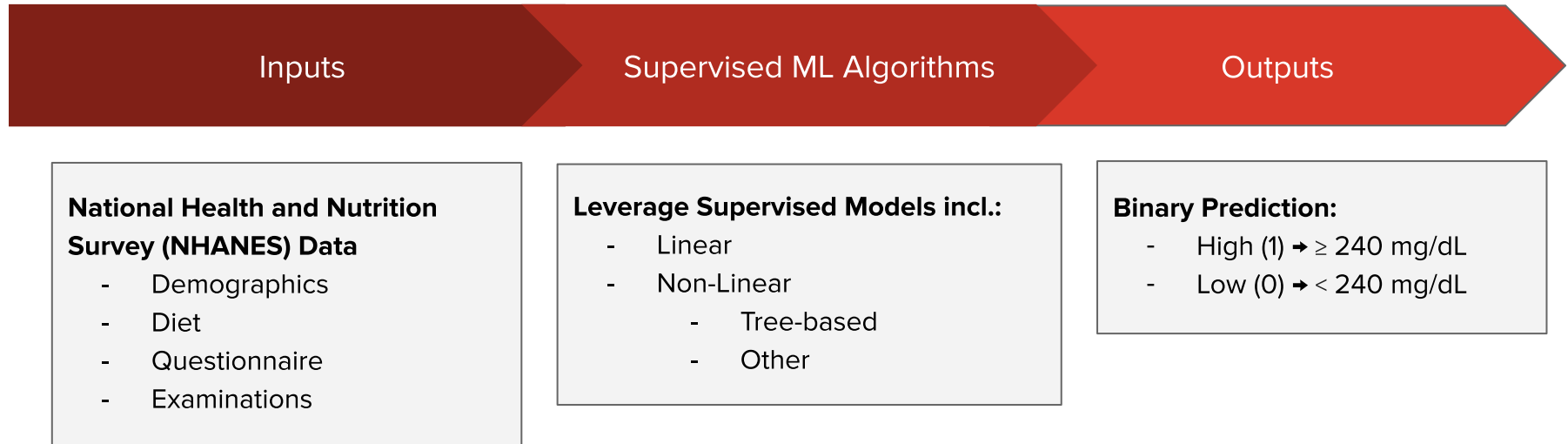- Examinations

**Leverage Supervised Models incl.:**
- Linear
- Non-Linear
  - Tree-based
  - Other

# Supervised ML Algorithms Can Help

**Primary Goal:**

Develop a supervised ML model to predict whether a person has a high total cholesterol level, without having to go to the clinic
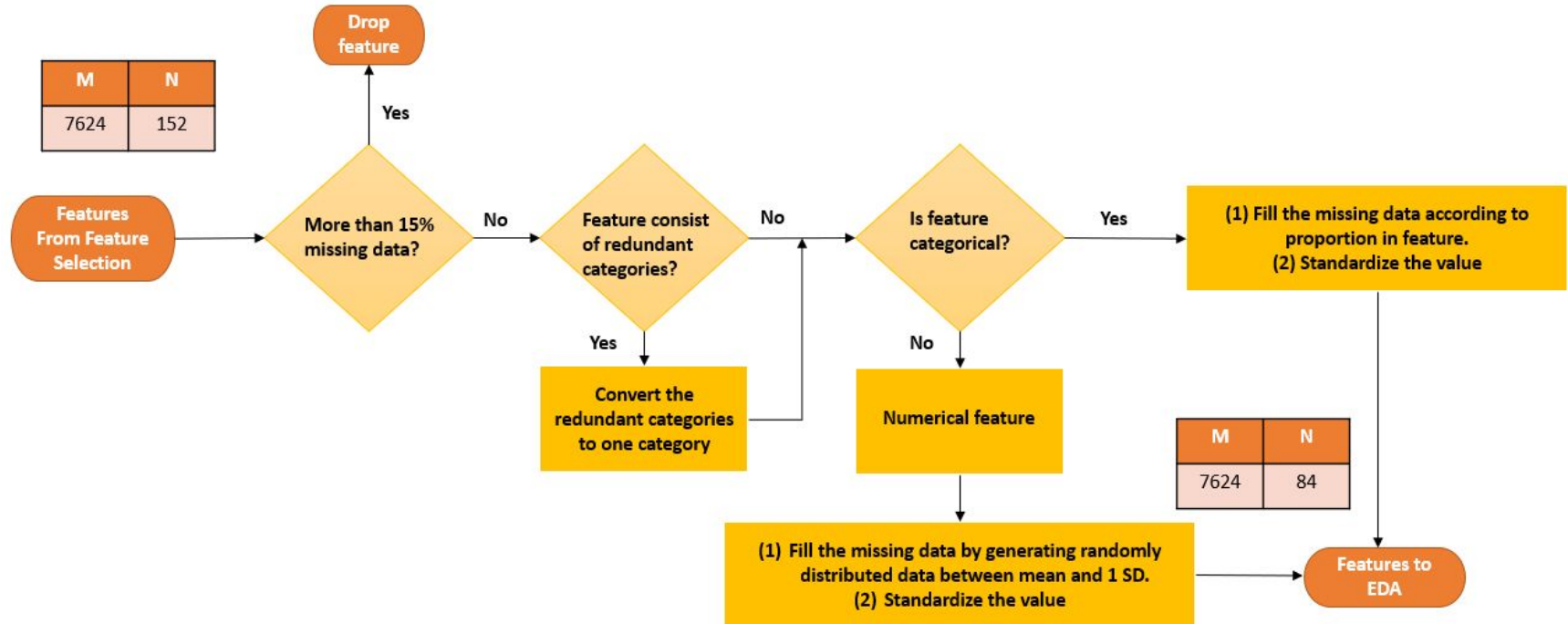
| Inputs | Supervised ML Algorithms | Outputs |
|---|---|---|

**National Health and Nutrition Survey (NHANES) Data**
- Demographics
- Diet
- Questionnaire
- Examinations

**Leverage Supervised Models incl.:**
- Linear
- Non-Linear
    - Tree-based
    - Other

**Binary Prediction:**
- High (1) ➜ ≥ 240 mg/dL
- Low (0) ➜ < 240 mg/dL

# Explaining the Dataset

**Feature Selection**

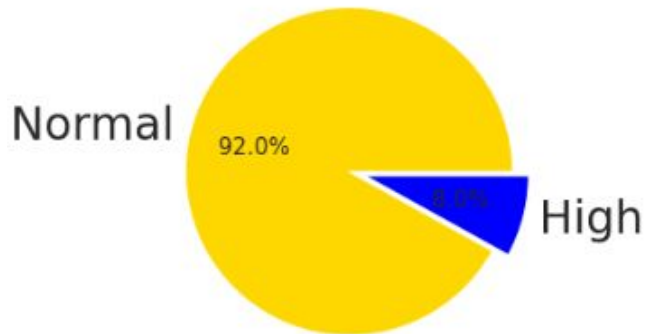| Steps | Observation (M) | Number of Features (N) |
|---|---|---|
| (1)  Original dataset | 9813 | 1390 |
| (2)  After feature selection | 9813 | 152 |
| (3)  Dropping the observations with no target output | 7624 | 152 |

# Explaining the Dataset

## <u>Data Imputation Flow Chart</u>

# Exploratory Data Analysis

Classification on basis of cholestrol levels



Distribution of feature: BMXWT

Distribution of feature: DR1TMAGN
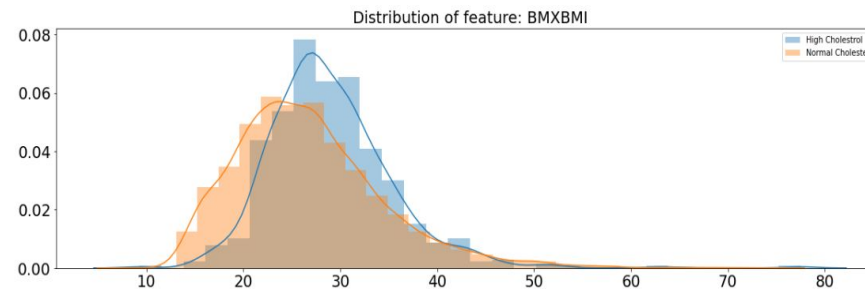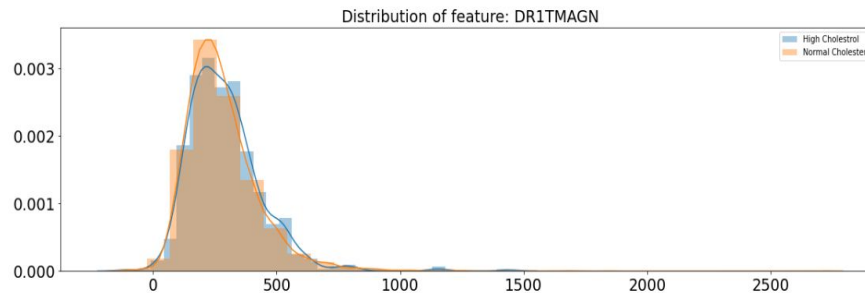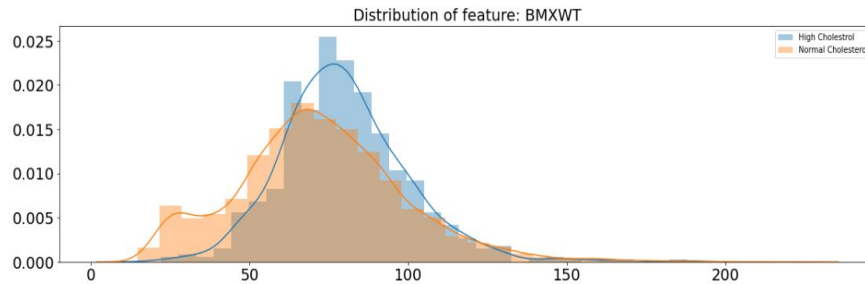
Distribution of feature: BMXBMI

Findings from Pie Chart:
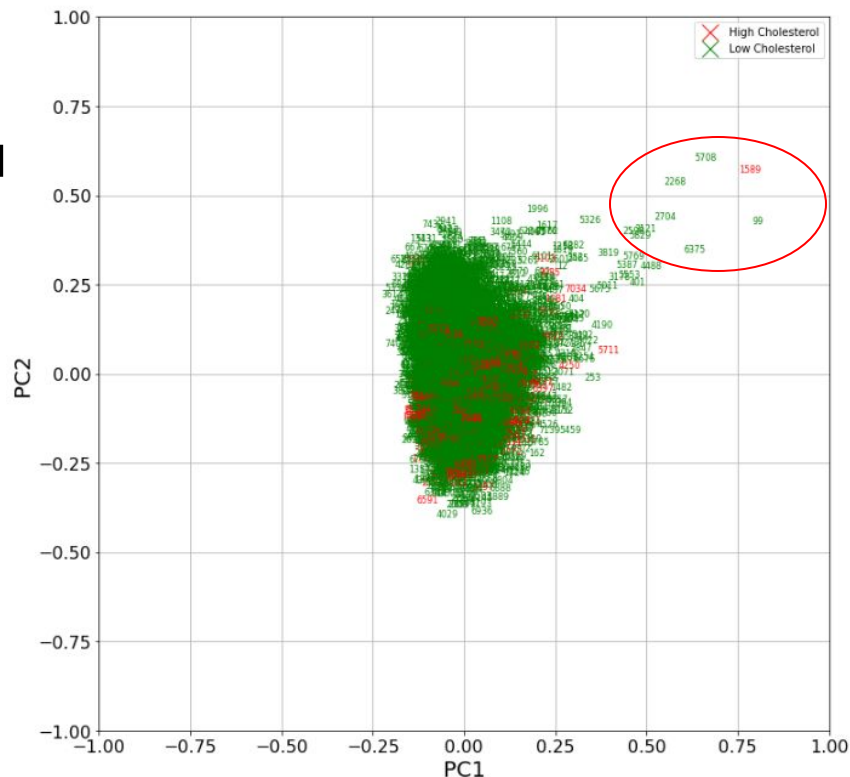➢ Created a balanced train set and an imbalanced test set

Findings from Histograms:
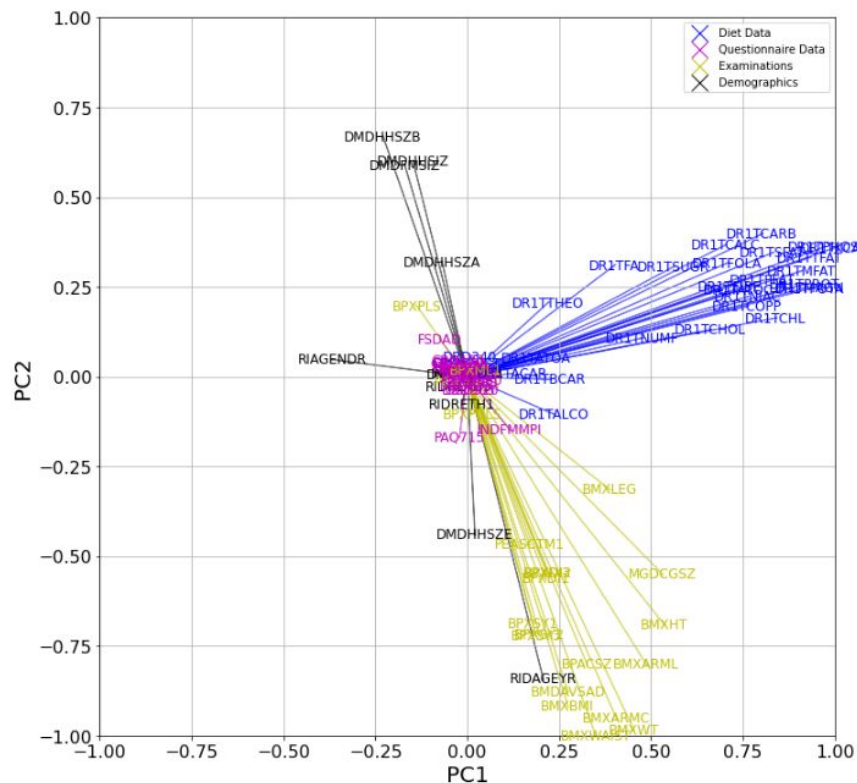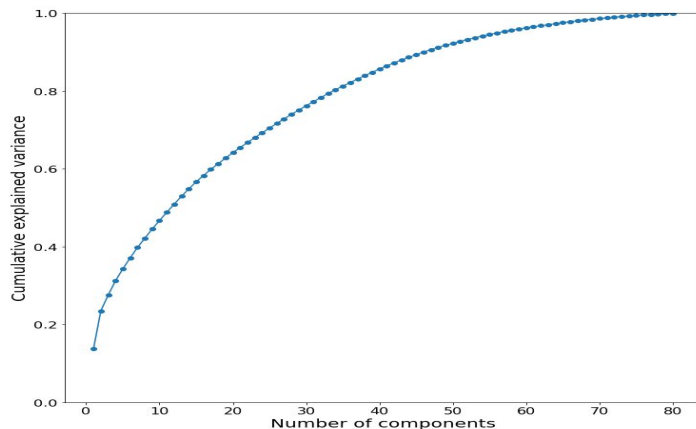➢ Most of the numerical features have overlapping classes

12

# Exploratory Data Analysis, contd.

➢ Our data doesn't seem to be linearly separable
➢ Outliers are mostly low cholesterol and have really high values for diet-related variables

# Exploratory Data Analysis, contd.

➢ Loadings show that features coming from the same dataset match

➢ Explained variance makes dimensionality reduction less "appealing"
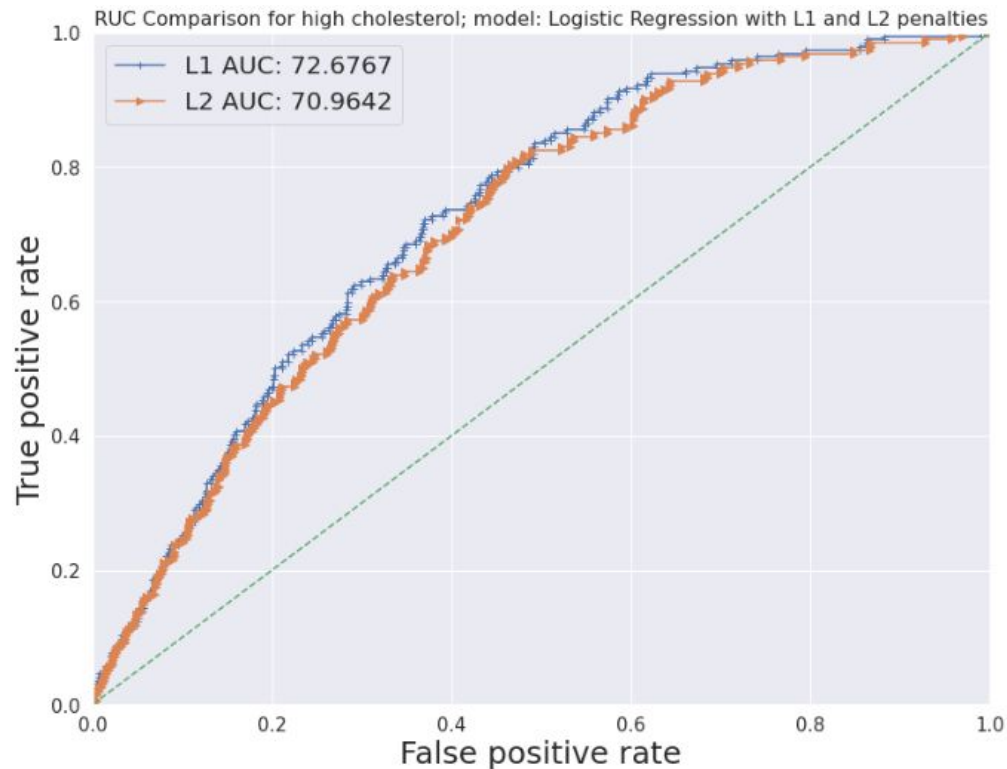
# Approach for Model Building

1. Training a model on balanced dataset and tested on imbalanced dataset
2. Fine-tuning hyperparameters using GridSearchCV
   - Score: Recall (exception: XGBoost with Area under precision and recall curve)
   - CV folds: 5
3. Comparing the models based on recall.

Recall measures the proportion of actual high cholesterol cases that are correctly identified as 'high' by the model.

Therefore, higher the recall, better the model

# Linear methods for classification

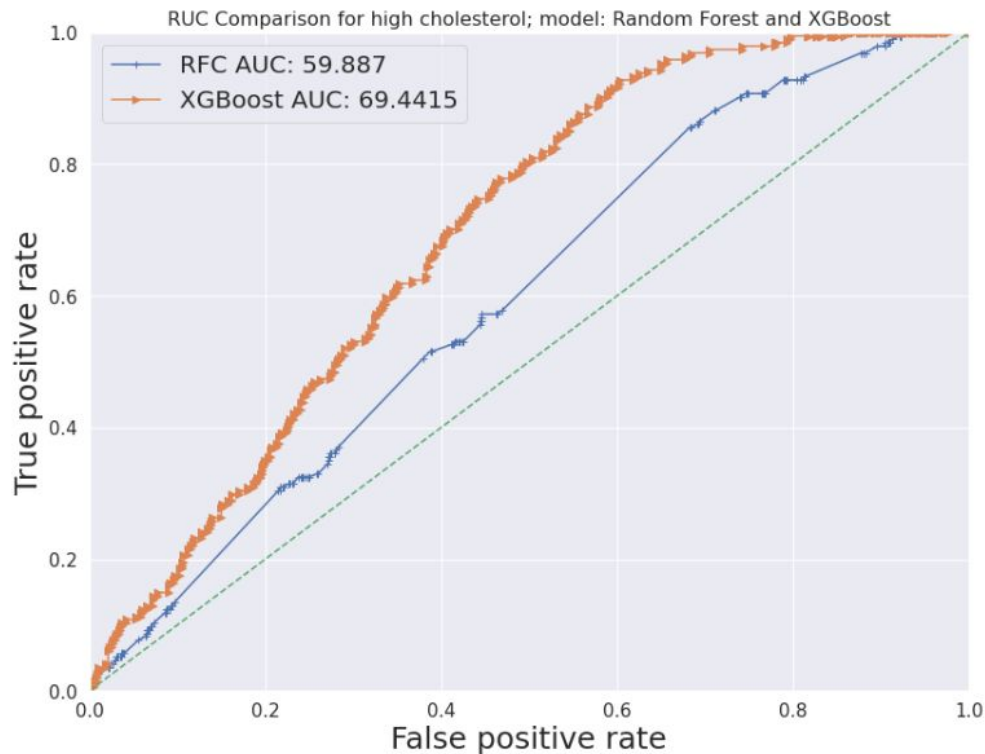Logistic Regression with L1 and L2 regularization



RUC Comparison for high cholesterol; model: Logistic Regression with L1 and L2 penalties

L1 AUC: 72.6767
L2 AUC: 70.9642

➢ Parameters
   ○ Solver: 'saga'
   ○ Regularizer strength C (Fine Tuned)

|  | Model Type | |
| --- | --- | --- |
| Metrics | Penalty L1 | Penalty L2 |
| Recall | 0.74 | 0.70 |
| Accuracy | 0.61 | 0.61 |

# Non-Linear methods for classification
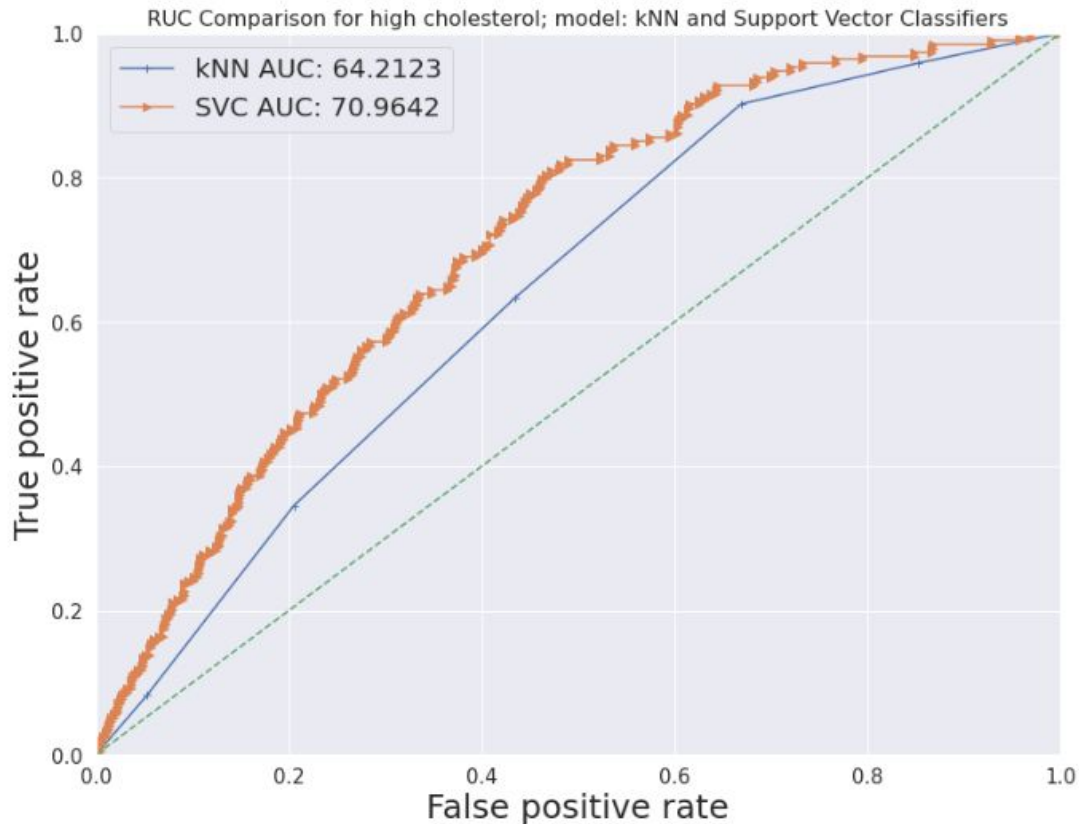
Random Forest Classifier and XGBoost



RUC Comparison for high cholesterol; model: Random Forest and XGBoost

- RFC AUC: 59.887
- XGBoost AUC: 69.4415

➢ Parameters to train & fine tune Random Forest
  ○ n_estimators
  ○ Max_depth

➢ Parameters to train & fine tune XGBoost
  ○ Eta (Learning rate)
  ○ Max_depth
  ○ Gamma
  ○ subsample

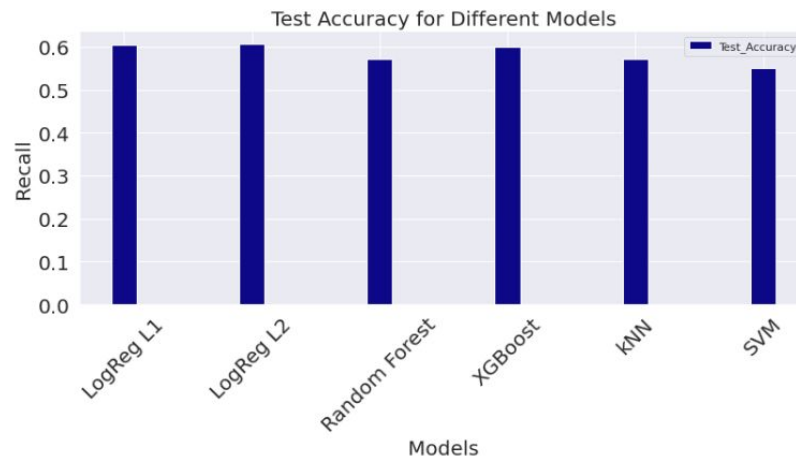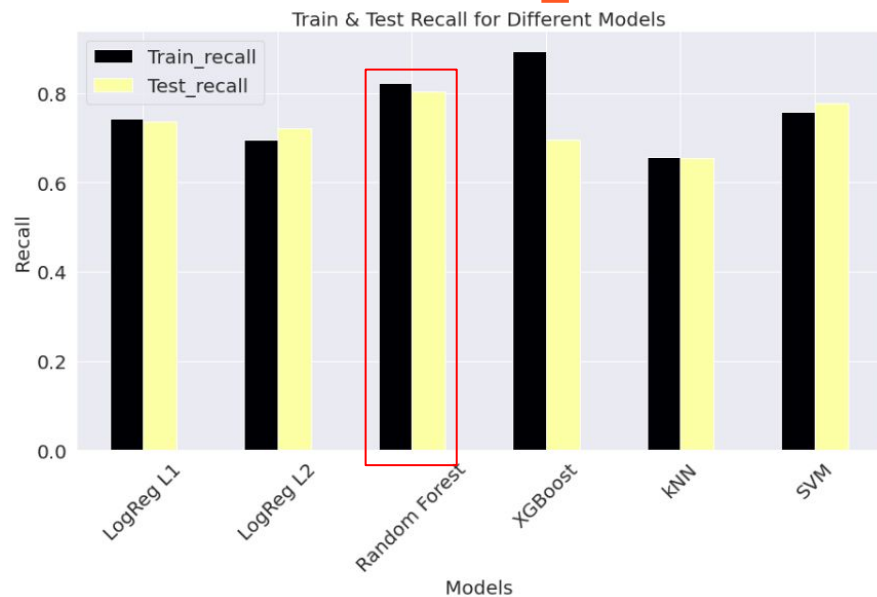| | Model Type | |
|---|---|---|
| **Metrics** | **Random Forest** | **XGBoost** |
| **Recall** | 0.80 | 0.70 |
| **Accuracy** | 0.56 | 0.60 |

# Non-Linear methods for classification

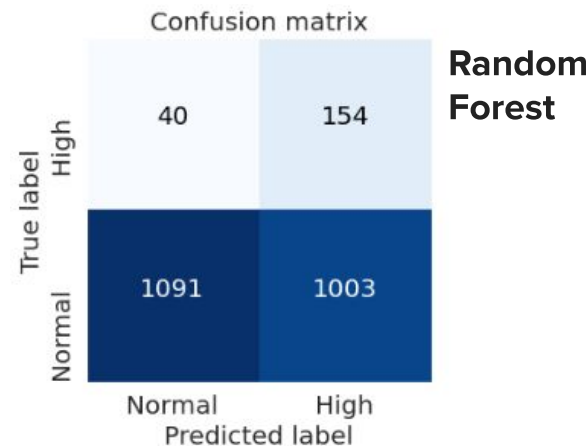K Nearest Neighbors (kNN) and Support Vector Classifier (SVC)



RUC Comparison for high cholesterol; model: kNN and Support Vector Classifiers

kNN AUC: 64.2123
SVC AUC: 70.9642

➢ Parameters to train & fine tune kNN
   ○ n_neighbors
   ○ Distance metric: euclidean
➢ Parameters to train & fine tune SVC
   ○ Kernel: RBF kernel
   ○ Regularization parameter C

| | Model Type | |
|---|---|---|
| **Metrics** | **kNN** | **SVC** |
| **Recall** | **0.63** | **0.79** |
| **Accuracy** | **0.57** | **0.55** |

# Model Comparison



Train & Test Recall for Different Models



Test Accuracy for Different Models



Confusion matrix — **Random Forest**

➢ Random Forest predicts high cholesterol more efficiently
➢ Random Forest > SVM > XGBoost > Logistic Regression L1 > Logistic Regression L2 > kNN
➢ Normal cholesterol cases are not classified accurately; this may be due to overlapping of classes
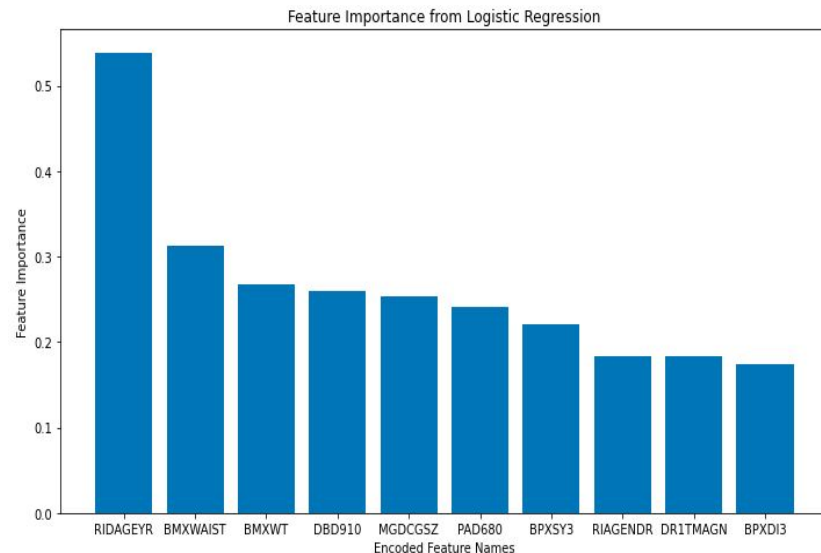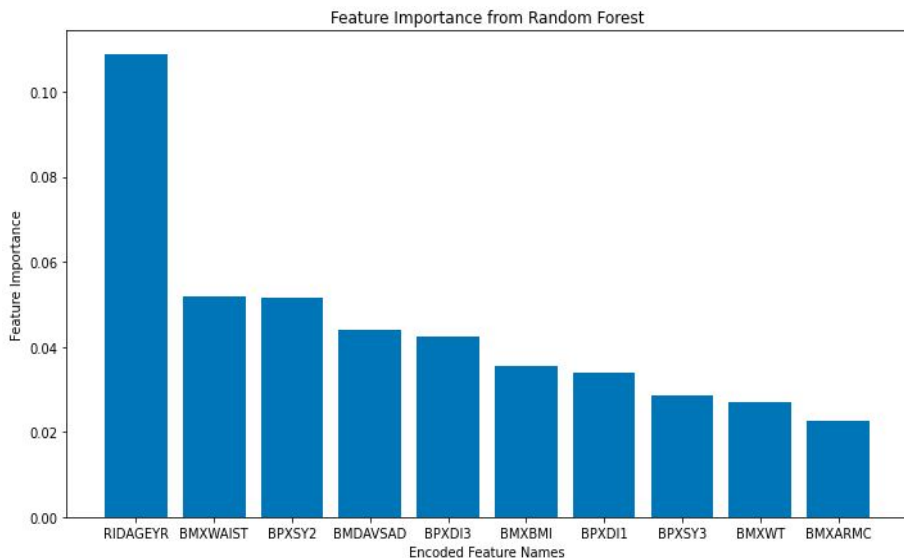
# Building for Deployment

**Challenge:**

➢ Determine cholesterol level by filling out online surveys

➢ 84 questions are too overwhelming for an online survey

**Solution:**

➢ Select the 10 most important features from logistic regression model

➢ Design a survey based on these 10 features

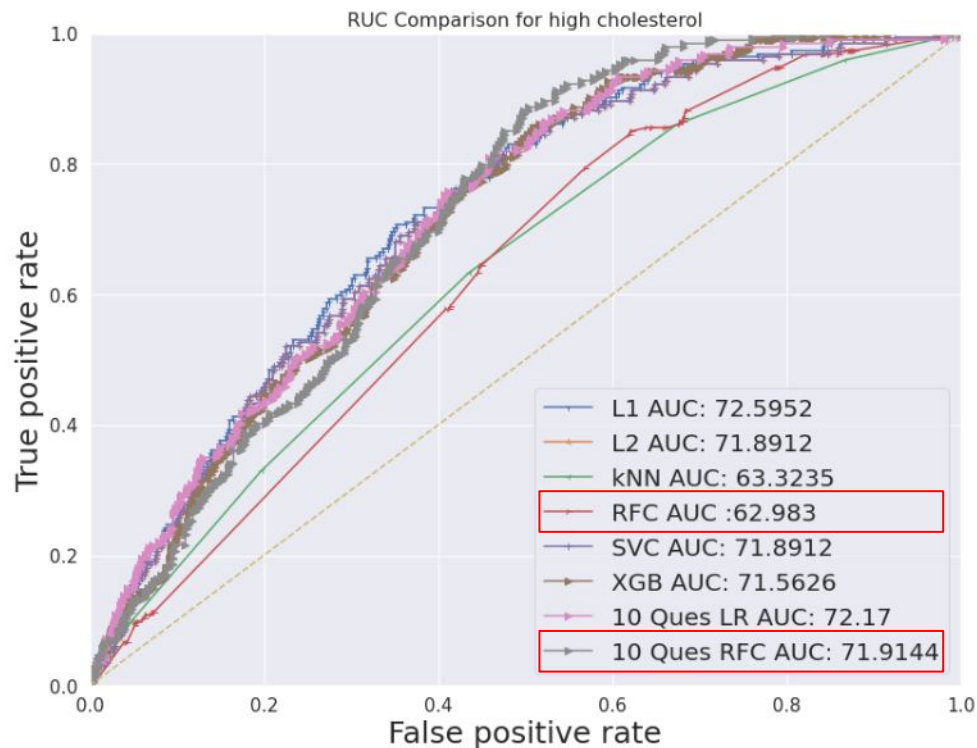➢ Non-linear models such as Random Forest shows promising results

# Building for Deployment Cont'd.

**Ten Most Important Features for Random Forest (Left) and Logistic Regression (Right):**



- ➢ Age and waist circumference are ranked as the 2 most important features by both methods
- ➢ Features related to blood pressure are important as well
- ➢ Majority of the features reflect the participants' diet habit and lifestyle

# Building for Deployment Cont'd.



RUC Comparison for high cholesterol

Legend:
- L1 AUC: 72.5952
- L2 AUC: 71.8912
- kNN AUC: 63.3235
- RFC AUC :62.983
- SVC AUC: 71.8912
- XGB AUC: 71.5626
- 10 Ques LR AUC: 72.17
- 10 Ques RFC AUC: 71.9144

| | Model Type | |
|---|---|---|
| **Metrics** | **Logistic Regression** | **Random Forest** |
| **Recall** | 0.77 | 0.84 |
| **Accuracy** | 0.61 | 0.56 |

# Conclusion

➢ Random forest algorithm can build the best performance model

➢ Questionnaire with only 10 questions can be used to classify the patient's cholesterol level

➢ Dataset Limitations: Only represent a population from USA

# Thank You!

Questions?