# Be Heart Smart

**The Healthy Healthcare Enthusiasts (Collaborators):**
**(Final-Project Group 7)**

- ❖ **Ayse Ozgun**
- ❖ **Pam Noble**
- ❖ **Subhangi Ghosh**
- ❖ **Krishnakali Sarkar**

# Cardiovascular Disease (CVDs)

Disorders of the heart and blood vessels including coronary heart disease, cerebrovascular disease, rheumatic heart disease and other conditions.
Leading cause of death globally ~ 40% deaths in the US.
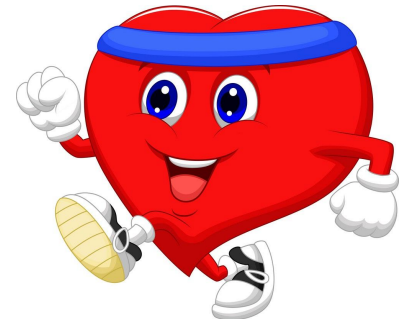
Leading Behavioral Risk Factors :

- Unhealthy diet,
- Physical inactivity
- Tobacco use
- Harmful use of alcohol

Effects of behavioral risk factors :

- Raised blood pressure,
- Raised blood glucose,
- Raised blood lipids,
- Overweight and
- Obesity.

A healthy heart is a happy heart

The purpose of this project is to spread awareness. Embracing a healthy lifestyle at any age can help prevent heart disease, and lower the risks for heart attack or stroke.

MEGAPIXL                                    Download from megapixl.com/30405387

# About the data

Website : [Cardiovascular Disease dataset](#) (Kaggle)

Description :

Three types of input features

➢ Objective

➢ Examination

➢ Subjective

| Objective | Examination | Subjective |
|---|---|---|
| Age (days) | Systolic Blood Pressure | Smoking |
| Height (cm) | Diastolic Blood Pressure | Alcohol Intake |
| Weight (kg) | Cholesterol | Physical Activity |
| Gender | Glucose | |

Target Variable : Presence or Absence of Cardiovascular Disease

# Questions we hope to answer with the data:

★ Is a person at risk of heart disease?

★ What are the potential risk factors for heart disease--smoking, alcohol consumption, obesity, etc?

★ Which factors are the best predictors of heart disease?

<u>Data repository</u>
➔ Database : PostgreSQL

<u>Classification model to predict risk (Yes/No) of heart disease based on different factors</u>
❖ Supervised Machine Learning
  ➢ Logistic Regression
  ➢ Random Forest

❖ Deep Neural Network

# Cleaning, processing, feature engineering

➢ 70,000 observations. Cleaned in PySpark
   ○ Numbers not observed in adult human population removed (systolic BP 16,020)

| | id | BMI | weight_status | obesity_status |
|---|---|---|---|---|
| 35363 | 77629 | 22.0 | normal | no |
| 41697 | 81468 | 29.8 | overweight | yes |
| 37065 | 55211 | 24.2 | normal | no |
| 2697 | 1778 | 19.8 | normal | no |
| 67862 | 73893 | 26.3 | overweight | no |
| 56957 | 13361 | 22.0 | normal | no |
| 16793 | 94697 | 31.2 | obese | yes |

❖ BMI (using information from CDC.gov)
   ➢ BMI between 15-60

❖ Pulse Pressure (Difference between systolic and diastolic blood pressure numbers)
   ➢ Positive and greater than 20

★ New features created - BMI, weight status and obesity, pulse-pressure
★ Total number of observation : 67466

# Interesting insights
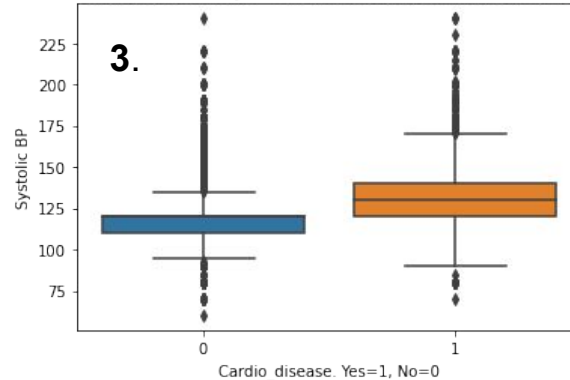
**Orange is positive for CVD**
**Blue is negative for CVD**

**2.**



Relationship between Obesity and Cardiovascular Disease

2. Obesity increases incidence of CVD
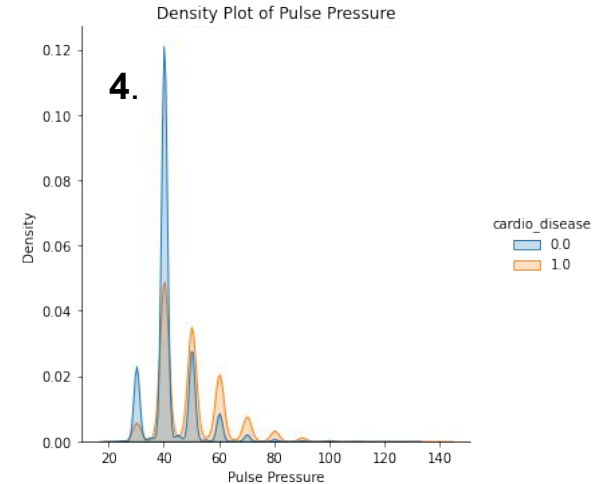
**1.**



Density Plot of Age

1 .Shift towards positive for CVD with increasing age

**3.**



Distribution of Systolic Blood Pressure Measurements
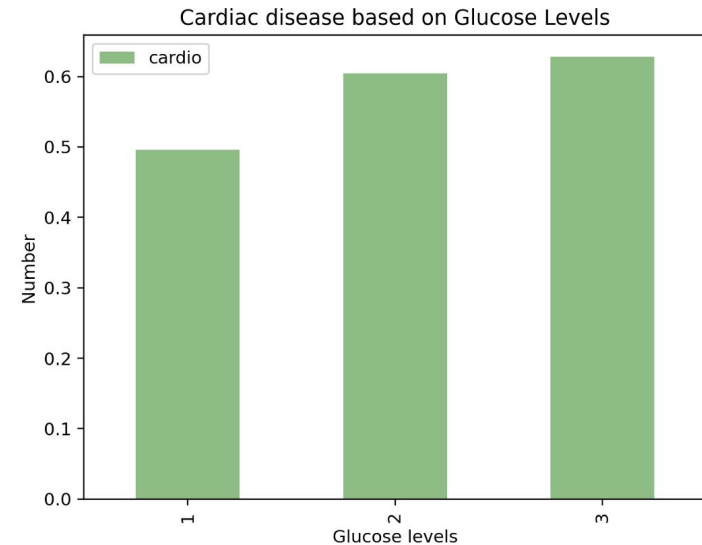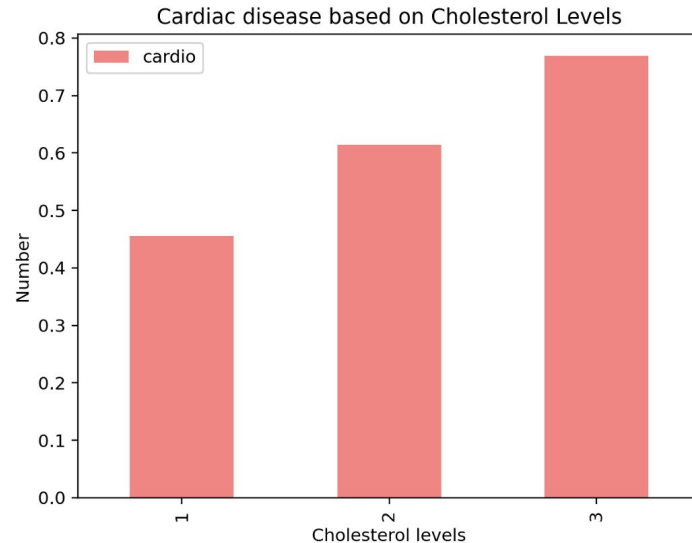
3. Patients with CVDs have higher blood pressure numbers

**4.**



Density Plot of Pulse Pressure

4. Higher pulse pressure shows increasing density for CVD

# Data Processing and Exploratory Data Analysis
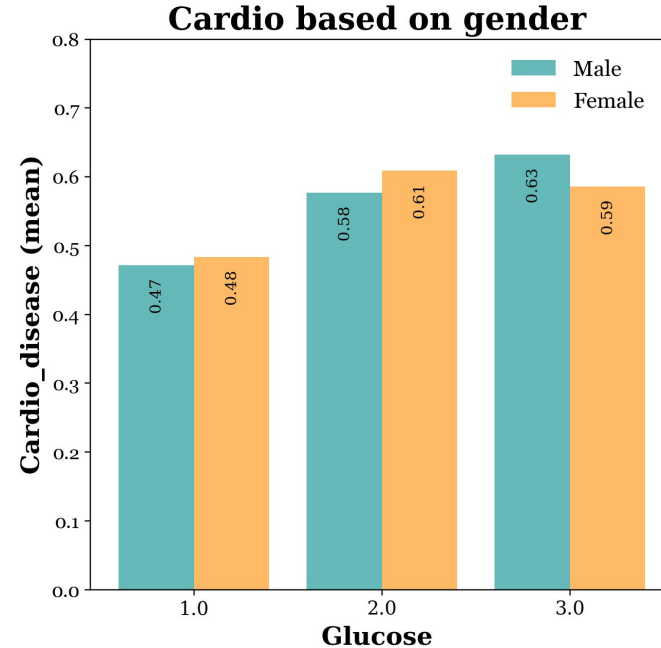
## Exploratory Data-Analysis:

Performed on the initial trial pre-processed data on Excel

# The effect of Cholesterol and Glucose on cardiac disease based on Gender



- Men with high Cholesterol have a higher chance of developing cardiac disease.

- Men with high Glucose levels have a higher chance of developing cardiac diseases.

# Supervised Machine Learning :   Logistic Regression

**Purpose :** Given a set of health and lifestyle conditions, the algorithm will be able to predict if the user <u>has or does not have</u> cardiovascular disease (CVD) -- a binary classification

Final_Table

| age | gender_M | height | weight | BMI | underweight | overweight | obese | is_obese | systolic_bp | diastolic_bp | pulse_pressure | cholesterol_moderate | cholesterol_high | glucose_moderate | glucose_high | smoker | alcohol_intake | active | cardio_disease |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 62.0 | 0 | 143.0 | 34.0 | 16.6 | 1 | 0 | 0 | 0 | 100.0 | 70.0 | 30.0 | 0 | 0 | 0 | 0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 43.0 | 0 | 143.0 | 36.0 | 17.6 | 1 | 0 | 0 | 0 | 90.0 | 60.0 | 30.0 | 0 | 0 | 0 | 0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 61.0 | 0 | 145.0 | 36.0 | 17.1 | 1 | 0 | 0 | 0 | 120.0 | 80.0 | 40.0 | 0 | 0 | 0 | 0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 56.0 | 0 | 144.0 | 36.0 | 17.4 | 1 | 0 | 0 | 0 | 100.0 | 70.0 | 30.0 | 0 | 0 | 0 | 0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 58.0 | 0 | 152.0 | 38.0 | 16.4 | 1 | 0 | 0 | 0 | 110.0 | 80.0 | 30.0 | 0 | 0 | 0 | 0 | 0.0 | 0.0 | 1.0 | 0.0 |

★   Data divided into Train, Validation, and Test sets (60 : 20 : 20 %). Test set is unseen

★   Scaled using standard scaler. Fit on the Train set, and transformed Train, Validation, and Test

★   K Fold Cross Validation with k=10 (scoring on recall) :

```
Cross-Validation Performance on Recall:

[0.67162698 0.68154762 0.67162698 0.66815476 0.65195835 0.657412
 0.67278136 0.65989093 0.6817055  0.66418651]

Mean Recall Score : 66.8089099794603
```
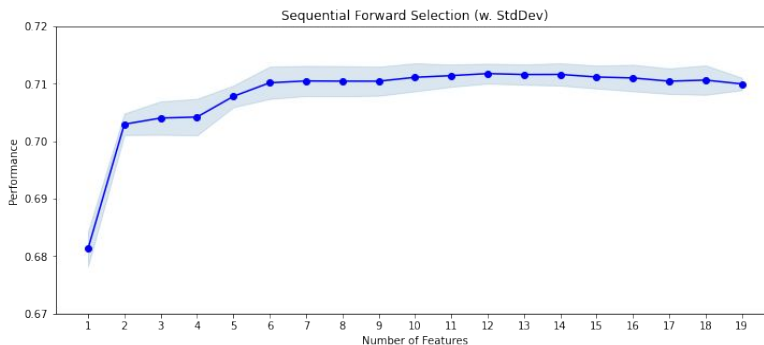
# Logistic Regression :

## Feature selection :
(Sequential Feature Selector was used)

**(On the Validation Set)**

Accuracy : 72.7 %,    Recall : 67 %

Precision : 75 %,    F1-score : 71 %



Features:
1. age
2. underweight
3. is_obese
4. systolic_bp
5. pulse_pressure
6. cholesterol_high
7. active

## It is important to not missing patients with cardiovascular disease, therefore recall is maximized



Threshold : 0.4 (on Test Set)

| | |
|---|---|
| Accuracy | 70.9 % |
| **Recall** | **80 %** |
| Precision | 67.2 % |
| F1-score | 73 % |

### Key Takeaways

➔ LR model was hypertuned to detect cardiovascular risk, even at the expense of including some false positives

➔ Out of the 19 original features, the model returned 7 as key determinants

# Supervised Machine Learning

## Random Forest Classifier

- Split the data into 'Train: Validation: Test' sets → 70: 20: 10
- Perform 10-fold Cross-validation
- Predict Random Forest classifier model--with default hyperparameters

|  | Initial model |
|---|---|
| **Accuracy** | 69.1% |
| **Recall** | 69.3% |
| **F1** | 69.2% |

# Supervised Machine Learning

## Random Forest Classifier

- Perform feature selection using sklearn's feature_importances_
- Create a RF classifier with selected features only
- Predict the model



| | Model with feature selection |
|---|---|
| **Accuracy** | 66.19% |
| **Recall** | 64.0% |
| **F1** | 66.0% |

# Supervised Machine Learning

## Random Forest Classifier

- Optimize our model with hyperparameter tuning
- Search for best parameters using scikit_learn's GridSearchCV function

**--Top 5 Hyperparameters--**

```
n_estimators = [300, 800]
max_depth = [5, 8]
min_samples_split = [2, 5]
min_samples_leaf = [1, 5]
```

| | param_max_depth | param_n_estimators | param_min_samples_leaf | param_min_samples_split | mean_test_score |
|---|---|---|---|---|---|
| 10 | 8 | 300 | 1 | 5 | 0.730785 |
| 12 | 8 | 300 | 5 | 2 | 0.730638 |
| 14 | 8 | 300 | 5 | 5 | 0.730638 |
| 9 | 8 | 800 | 1 | 2 | 0.730511 |
| 8 | 8 | 300 | 1 | 2 | 0.730469 |

- Plug in the best hyperparameters into our Random Forest Classifier, retrain our model and predict it on the unseen test set.

```
{'max_depth': 8,
 'min_samples_leaf': 1,
 'min_samples_split': 5,
 'n_estimators': 300}
```

# Supervised Machine Learning

## Random Forest Classifier

### Classification Report

```
            precision   recall   f1-score   support

         0     0.71      0.80       0.75       3455
         1     0.76      0.66       0.71       3292

  accuracy                          0.73       6747
 macro avg     0.74      0.73       0.73       6747
weighted avg   0.74      0.73       0.73       6747
```

### Confusion Matrix

```
[[2763  692]
 [1110 2182]]
```

| | Final model |
|---|---|
| **Accuracy** | 73.3% |
| **Recall** | 66.0% |
| **Precision** | 76.0% |
| **F1** | 71.0% |

# Neural network Model

❖ The deep neural network model was run on the final merged dataset.

❖ Activation function for input: Relu

❖ Output function: Sigmoid

❖ No of hidden layers: 5

❖ The loss function:binary_crossentrophy

❖ Optimizer: **rmsprop**

❖ The accuracy of this model is 73%

Be Heart Smart Neural Network Model

Input Layer (+3)

(+690)

(+490)

(+190)

(+90)

(+40)

Output Layer

# Comparison of Machine Learning

|  | Logistic Regression | Random Forest | Deep Neural network |
|---|---|---|---|
| **Accuracy** | 70% | 73.3% | 73.0% |
| **Precision** | 67% | 76.0% | - |
| **Recall** | 80% | 66.0% | - |
| **F1** | 73% | 71.0% | - |

**We want to maximize recall without compromising accuracy**

# Thank you

We would like to thank

- ❖ Klaus
- ❖ Artem
- ❖ Jacob
- ❖ Trent
- ❖ Jackson
- ❖ Geoff
- ❖ Gael
- ❖ Tutors
- ❖ All our classmates
- ❖ Last but not the least the Amazing Group 7 members .

# Tableau Dashboard

# Be Heart Smart

Use the interactive charts below to explore the dataset

## Patients ID No:

0

### Demographic Info

ID: 0
AGE: 50
GENDER: 2
HEIGHT: 168.0
WEIGHT: 62.0
SYSTOLIC_BP: 110.0
DIASTOLIC_BP: 80.0
CHOLESTEROL: 1
GLUCOSE: 1
SMOKER: 0
ALCOHOL_INTAKE: 0
ACTIVE: 1
CARDIO_DISEASE: 0
BMI: 22.0
WEIGHT_STATUS: normal
OBESITY_STATUS: no

## Systolic BP

| Factor | Value |
|--------|-------|
| Age | 50 |
| BMI | 22.0 |
| Height | 168.0 |
| Weight | 62.0 |
| Systolic | 110.0 |
| Diastolic | 80.0 |

## Cardiac Disease Indicator

0

## Factors that cause Heart Disease

- Cholesterol
- Glucose
- Smoker
- Alcohol Intake
- Active

Count / Factors

# Dashboard

❤️ Click Me

# Dashboard :
# Web Application to  Predict Cardiovascular Disease