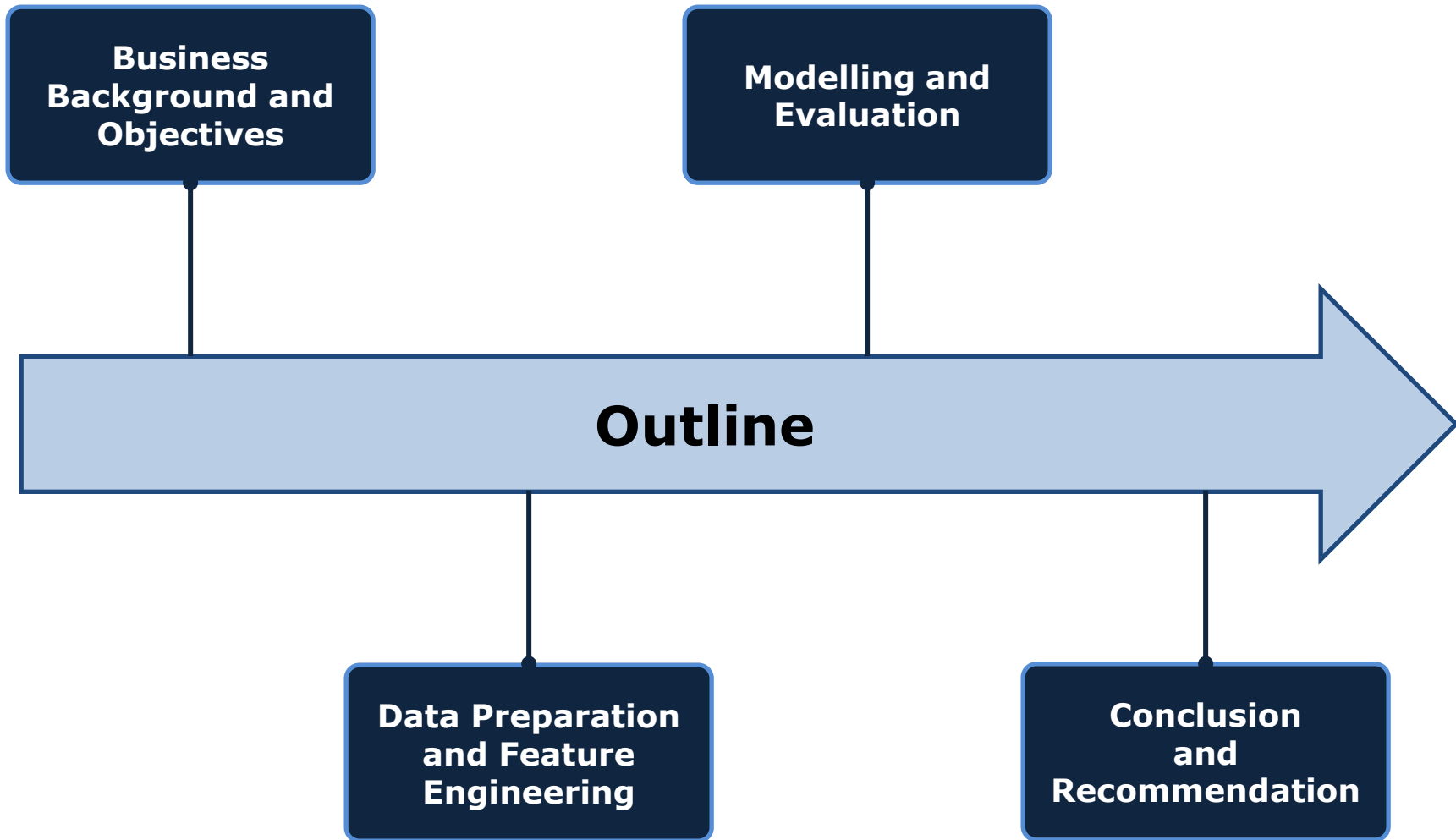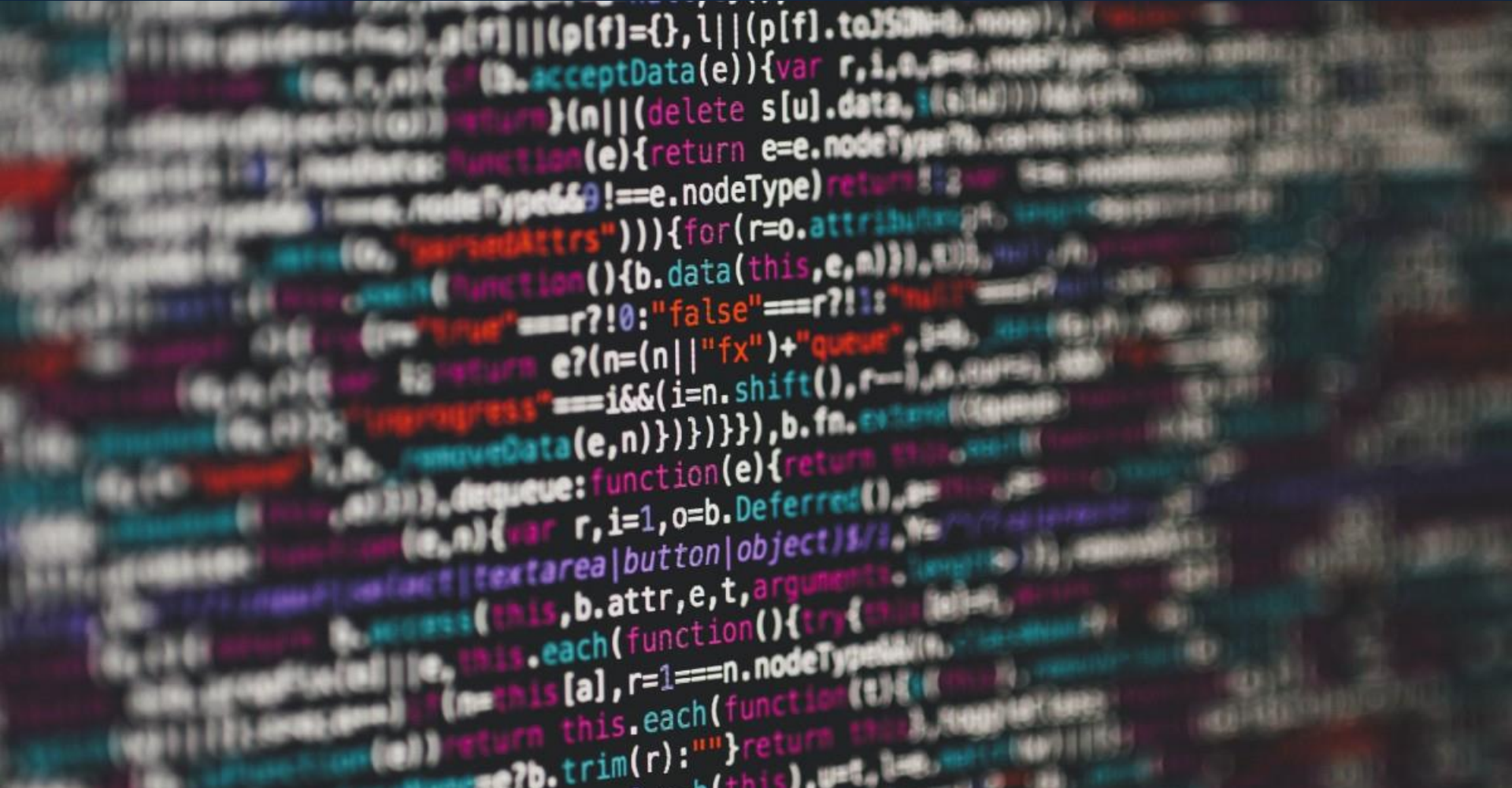# Heart Failure Analysis and Prediction with Logistic Regression

By: Aldimeola Alfarisy

# Business Background and Objectives

# Business Background and  Objectives

## Background

      Cardiovascular diseases (CVDs) are the leading cause of death globally with taking an estimated 17.9 million lives each year. More than four out of five CVD deaths are due to heart attacks and one third of these deaths occur prematurely in people under 70 years of age. Heart failure is a common event caused by CVDs. early detection and management wherein a machine learning model can be of great help.

## Objectives

- What factors affect heart failure?
- How accurate a machine learning model can predict to heart disease early detection?

Image source:
https://news.harvard.edu/gazette/story/ 2022/04/infertility-history-linked-with- increased-risk-of-heart-failure/

Data Preparation and Feature Engineering

# Data Preparation and Feature Engineering

## Dataset Information

### **918** rows

### **11** feature

**Numerical Feature**

- **Age**
- **RestingBP**
- **Cholesterol**
- **FastingBS**
- **MaxHR**
- **Oldpeak**

**Categorical Feature**

- **Sex**
- **ChestPainType**
- **RestingECG**
- **ExerciseAngina**
- **ST_Slope**

### **1** target

**Heart Disease**

Source dataset: https://www.kaggle.com/datasets/fedesoriano/heart-failure-prediction

# Data Preparation and Feature Engineering

## Dataset Attribute Information

| Column name | Description |
|---|---|
| Age | Age of the patient [years] |
| Sex | Sex of the patient [M: Male, F: Female] |
| ChestPainType | Chest pain type [TA: Typical Angina, ATA: Atypical Angina, NAP: Non-Anginal Pain, ASY: Asymptomatic] |
| RestingBP | Resting blood pressure [mm Hg] |
| Cholesterol | Serum cholesterol [mm/dl] |
| FastingBS | Fasting blood sugar [1: if FastingBS > 120 mg/dl, 0: otherwise] |
| RestingECG | Resting electrocardiogram results [Normal: Normal, ST: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV), LVH: showing probable or definite left ventricular hypertrophy by Estes' criteria] |
| MaxHR | Maximum heart rate achieved [Numeric value between 60 and 202] |
| ExerciseAngina | Exercise-induced angina [Y: Yes, N: No] |
| Oldpeak | ST [Numeric value measured in depression] |
| ST_Slope | The slope of the peak exercise ST segment [Up: upsloping, Flat: flat, Down: downsloping] |
| HeartDisease | Output class [1: heart disease, 0: Normal] |

# Data Preparation and Feature Engineering

## General Info

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column          Non-Null Count   Dtype
---  ------          --------------   -----
 0   Age             918 non-null     int64
 1   Sex             918 non-null     object
 2   ChestPainType   918 non-null     object
 3   RestingBP       918 non-null     int64
 4   Cholesterol     918 non-null     int64
 5   FastingBS       918 non-null     int64
 6   RestingECG      918 non-null     object
 7   MaxHR           918 non-null     int64
 8   ExerciseAngina  918 non-null     object
 9   Oldpeak         918 non-null     float64
 10  ST_Slope        918 non-null     object
 11  HeartDisease    918 non-null     int64
dtypes: float64(1), int64(6), object(5)
memory usage: 86.2+ KB
```

```
[9]  data.duplicated().sum()

      0
```

### Result

- **No Missing Value**

- **No Duplicated Value**

# Data Preparation and Feature Engineering

## Boxplot for Numerical Data



Outliers
Cholesterol

# 19,93%

Trim data

Outliers
Oldpeak

# 1,74%

Drop data

**After filtering row becomes 902 (before 918)**

# Data Preparation and Feature Engineering

## Exploratory Data Analysis (EDA) Insights

Numerical Data



- There is no redundant feature(s)
- Oldpeak has highest correlation to heart disease (41%)
- Meanwhile MaxHR has high correlation to heart disease, but in negative way (-40%)
- The higher Oldpeak and the lower MaxHR tend have heart disease and heart failure

# Data Preparation and Feature Engineering

## Exploratory Data Analysis (EDA) Insights

### Categorical Data



- The highest category that risky have heart disease:

  - Sex                   : Male
  - ChestPainType     : ASY (Asymptomatic)
  - RestingECG          : Normal
  - ExcerciseAngina    : Yes
  - ST_Slope            : Flat

- Those condition needs to be proven by further check other aspects (ex. Lifestyle)

# Data Preparation and Feature Engineering

## Feature Engineering (Label Encoding)

| Features | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| Sex | M | F | | |
| ChestPainType | ASY | NAP | ATA | TA |
| RestingECG | Normal | ST | LVH | |
| ExerciseAngina | N | Y | | |
| ST_Slope | Up | Flat | Down | |

## Feature Engineering (Standardization)

# Modelling and Evaluation

## Split Data

# Modelling and Evaluation

## Logistic Regression



Training Data

Test Data

**Model from test data not far from the model training data and has high accuracy and precision**

# Conclusion and Recommendation

## Conclusions

- All features in the dataset are used to analysing (no redundant features)
- Individual's Old peak is the highest factor that cause heart disease and affects heart failure
- Logistic regression model is capable to used due to high accuracy and precision

## Recommendations

Adding more features about lifestyle. Example:

- Smoker status
- Daily food

# THANK YOU

Email : aldimeolaalfarisy@yahoo.com
LinkedIn : aldimeolaalfarisy
Github : aldimeolaalfarisy