# Project 1
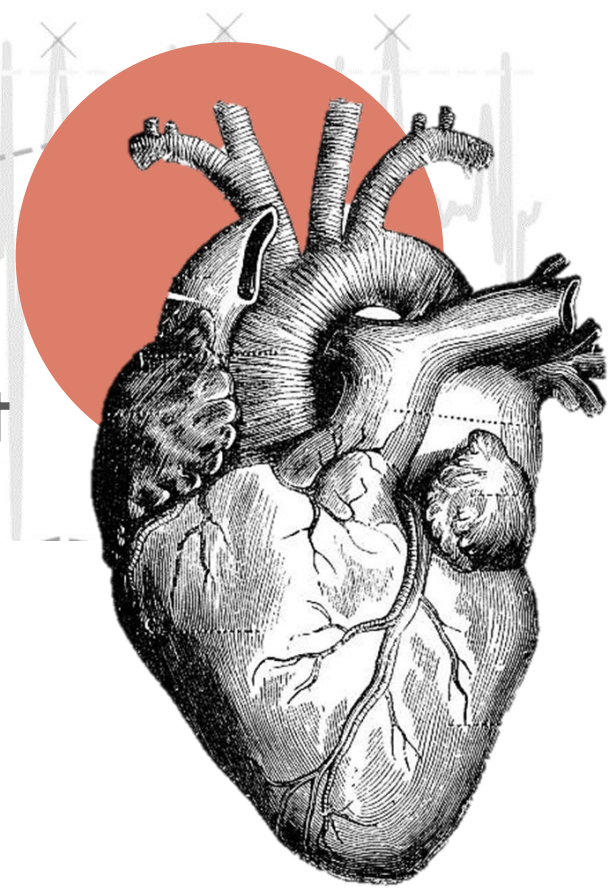
## Data analysis of Cleveland Heart Disease Dataset
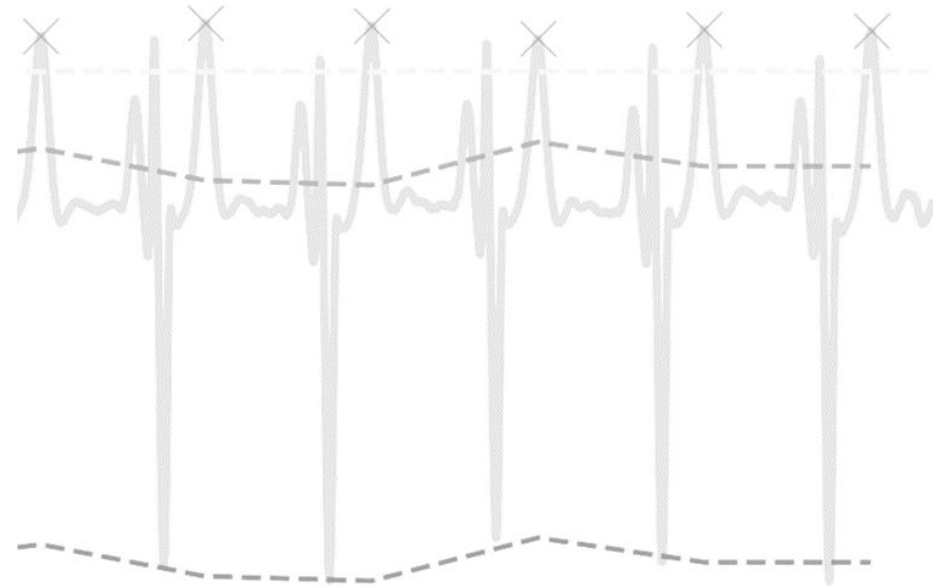
**Group 11:**
Nicolò Retis
Hang Mai Anh Vo (Emmy)
Rinoshan Parameswaran

# Contents

1. Introduction

2. Exploratory Data Analysis (EDA)

3. Model Selection

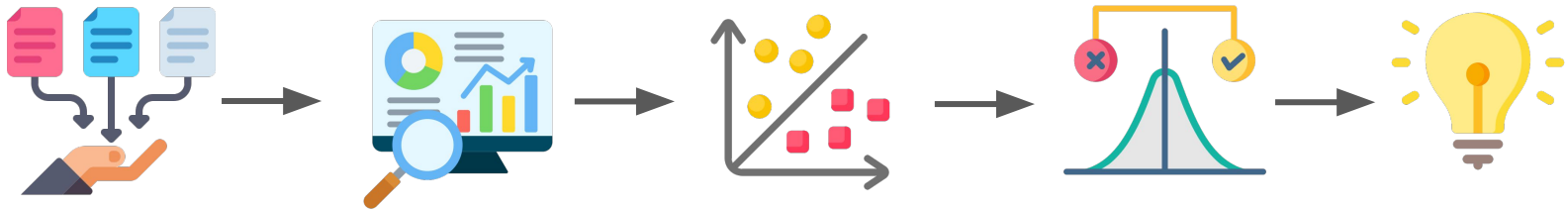4. Model Evaluation and Comparison

5. Key Findings and Conclusions

# 1. Introduction

**Objective**: From Cleveland Heart Disease Datasets

- To see the clinical pictures of heart disease patients
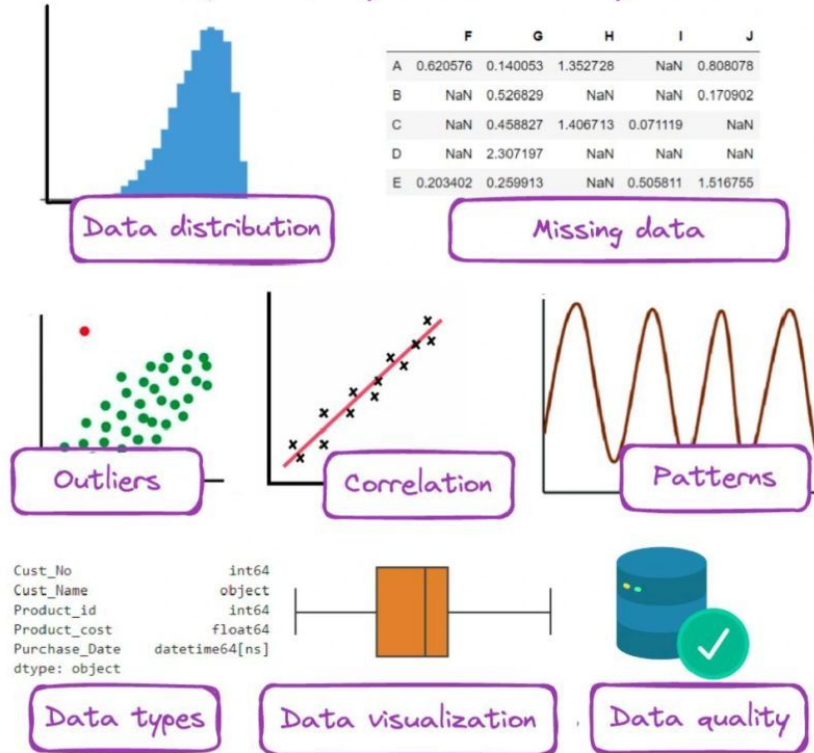- To build prediction models to classify patients with and without heart diseases

**Our work**:

# 2. Exploratory Data Analysis (EDA)

**Objective:** To understand feature distributions and detect patterns associated with heart disease

Source: https://www.markovml.com/blog/exploratory-data-analysis

# What did we do?

- Handle missing data
- Correlation matrix
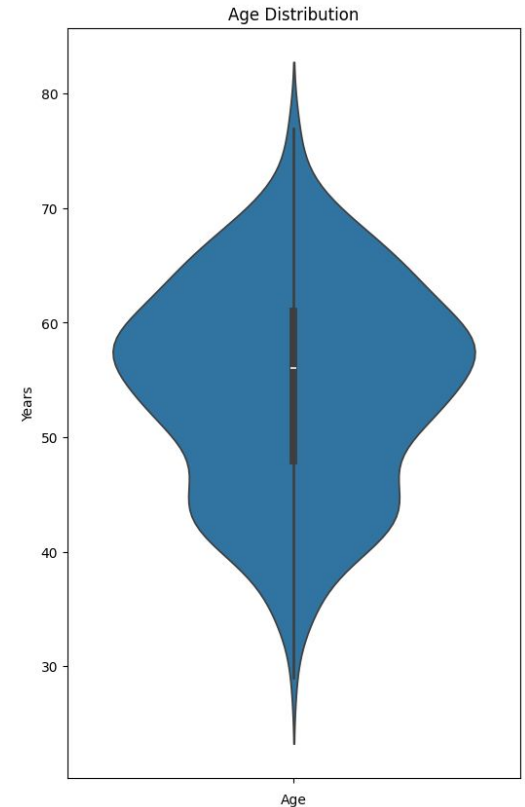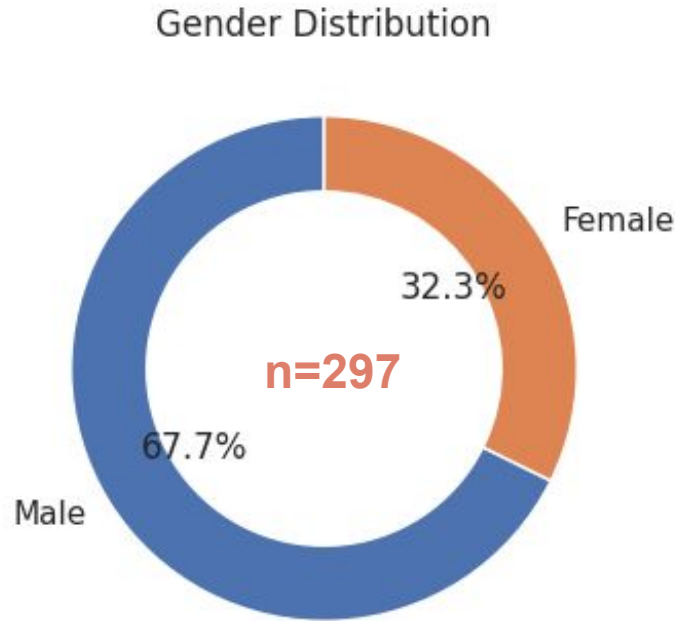- Visualize important features in the dataset

# Dataset overview

- 303 patient records. 6 records has missing data → remove.

| Variable Name | Role | Type | Demographic | Description | Units | Missing Values |
|---|---|---|---|---|---|---|
| age | Feature | Integer | Age | | years | no |
| sex | Feature | Categorical | Sex | | | no |
| cp | Feature | Categorical | | | | no |
| trestbps | Feature | Integer | | resting blood pressure (on admission to the hospital) | mm Hg | no |
| chol | Feature | Integer | | serum cholestoral | mg/dl | no |
| fbs | Feature | Categorical | | fasting blood sugar > 120 mg/dl | | no |
| restecg | Feature | Categorical | | | | no |
| thalach | Feature | Integer | | maximum heart rate achieved | | no |
| exang | Feature | Categorical | | exercise induced angina | | no |
| oldpeak | Feature | Integer | | ST depression induced by exercise relative to rest | | no |
| slope | Feature | Categorical | | | | no |
| ca | Feature | Integer | | number of major vessels (0-3) colored by flourosopy | | yes |
| thal | Feature | Categorical | | | | yes |
| num | Target | Integer | | diagnosis of heart disease | | no |

Detrano, R., Janosi, A., Steinbrunn, W., Pfisterer, M., Schmid, J. J., Sandhu, S., ... & Froelicher, V. (1989). **International application of a new probability algorithm for the diagnosis of coronary artery disease**. *The American journal of cardiology*, *64*(5), 304-310.
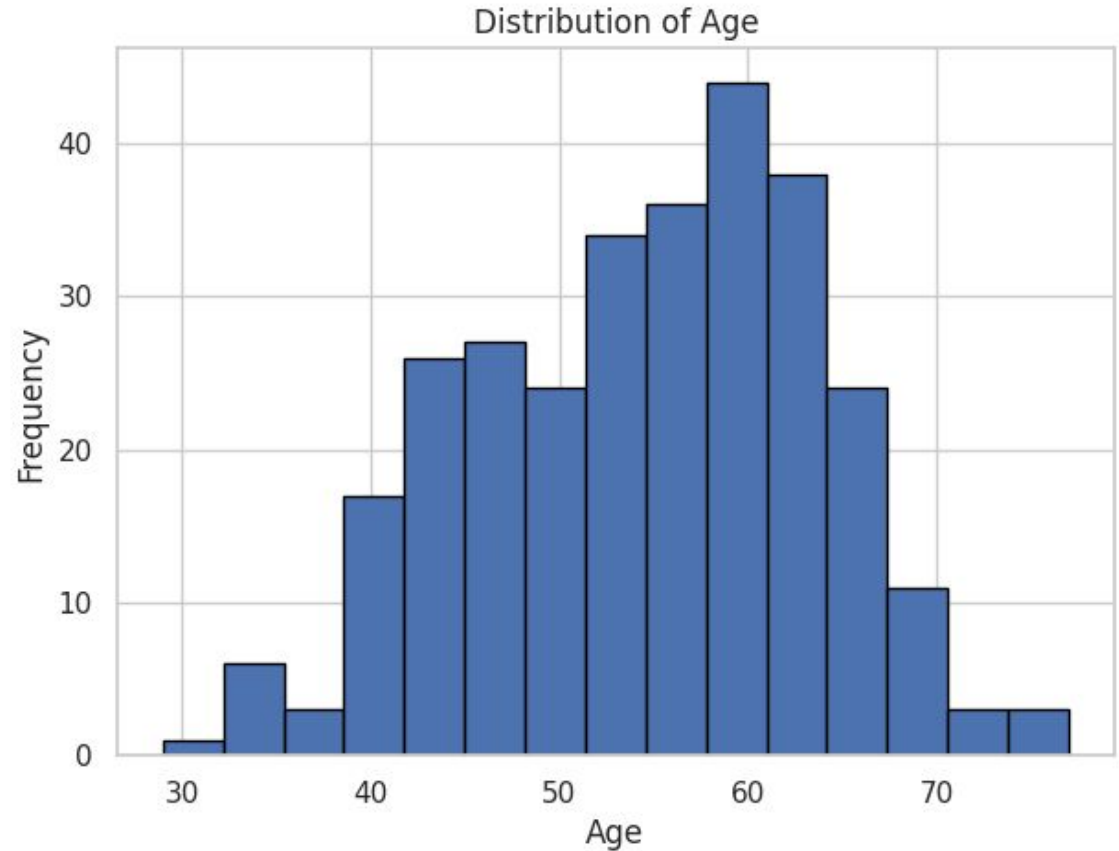
# Gender and age Distribution

- Age [29-77] years.

- **Majority of males.** This shown the same patterns with multiple epidemiological studies. [1]



Gender Distribution

Female 32.3%

n=297

67.7%

Male



Age Distribution

[1] Gao, Z., Chen, Z., Sun, A., & Deng, X. (2019). Gender differences in cardiovascular disease. *Medicine in Novel Technology and Devices*, *4*, 100025.
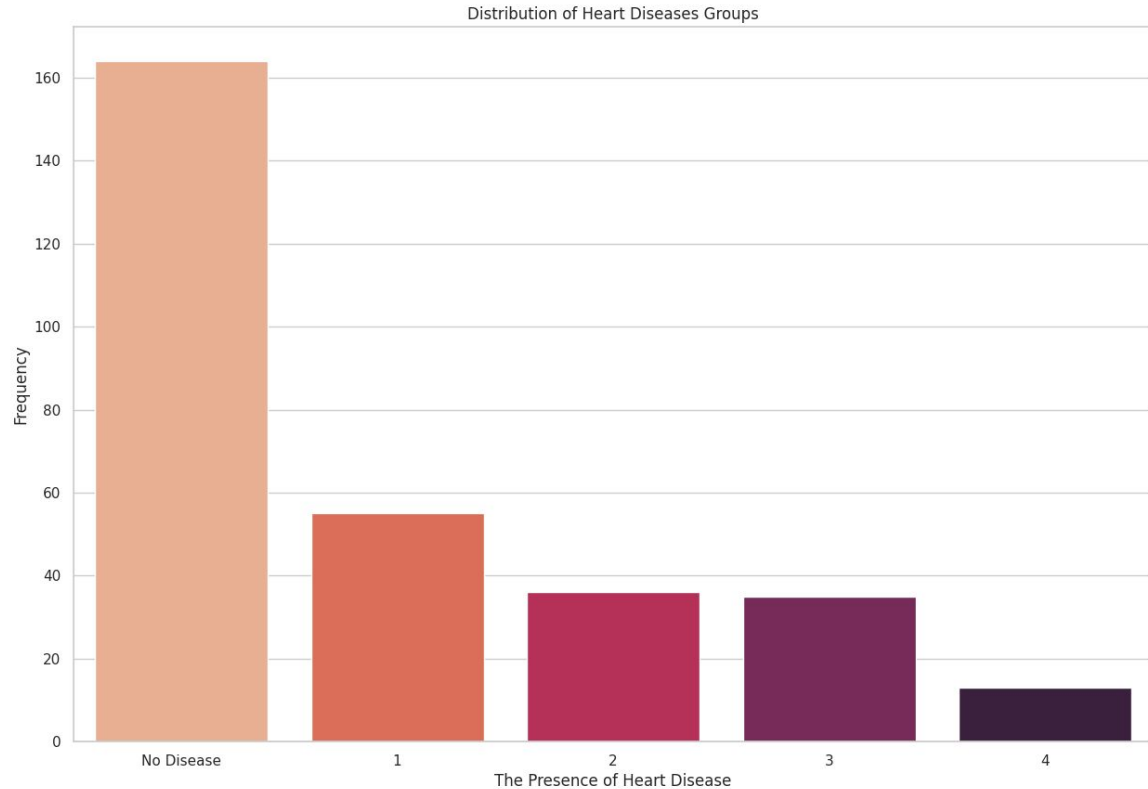
# Age Countplot

- People under 40 and 70+ are under represented.

- Only 21 patients belonging to these groups.



Distribution of Age

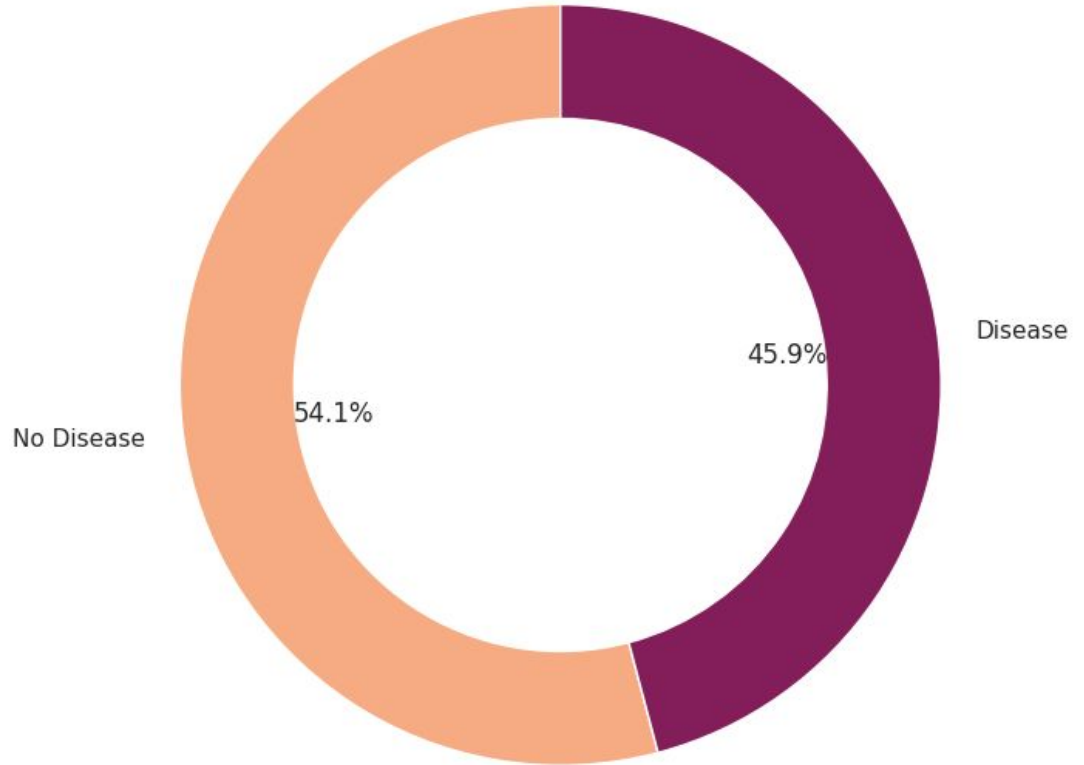# Heart disease distribution

- The majority of the patients don't have any heart disease.

- Number from 1 to 4 are used to rank the severe of heart disease.



Distribution of Heart Diseases Groups

# Grouping approach

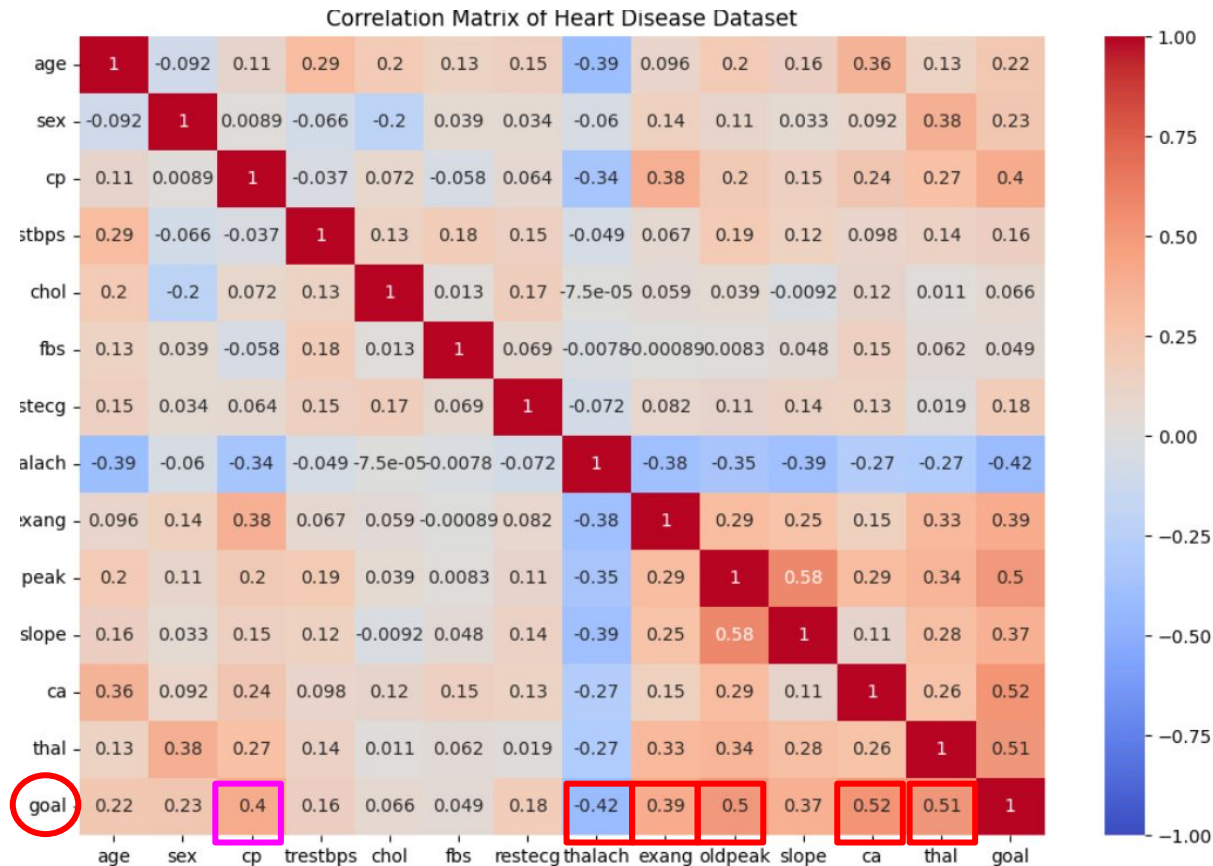Patients were splitted into 2 groups, based on whether they do or don't have heart disease:

- More robust data inference
- Better understanding of the data
- Balancing the data

## Distribution of Heart Disease Groups



Disease 45.9%

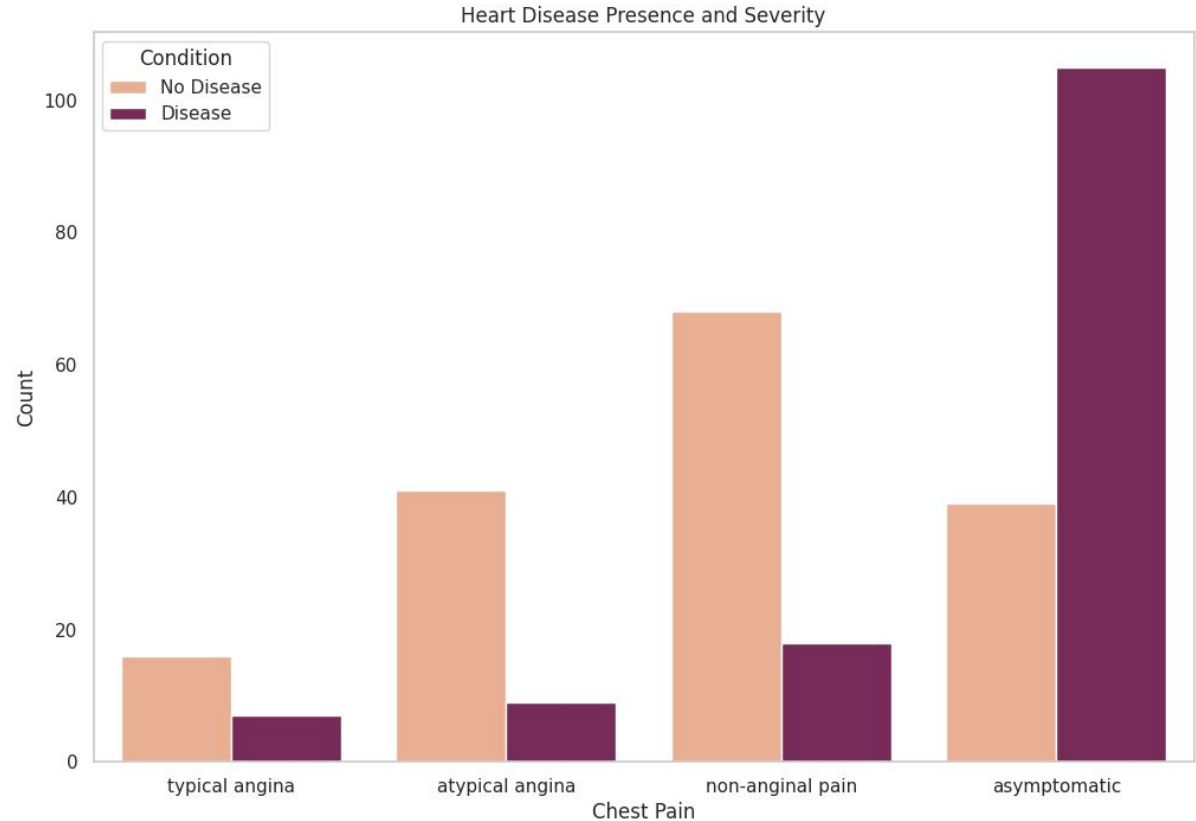No Disease 54.1%

# Correlation Matrix

Features with high corr.

- Chest Pain Type (**cp**)
- ST Depression (**oldpeak**)
- Exercise-Induced Angina (**exang**)
- Number of Blocked Vessels (**ca**)
- Thalassemia (**thal**)
- Maximum Heart Rate Achieved (**thalach**)



Correlation Matrix of Heart Disease Dataset

# Chest Pain Distribution

- Chest Pain correlates with disease.
- Patient with disease tend not to manifest chest pain.



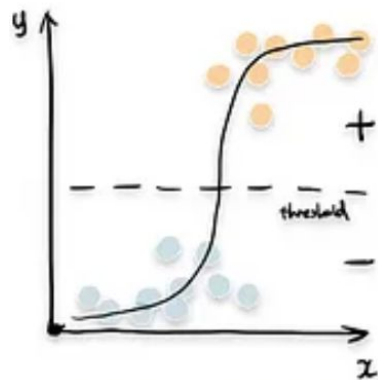Heart Disease Presence and Severity

# 3. Model selection
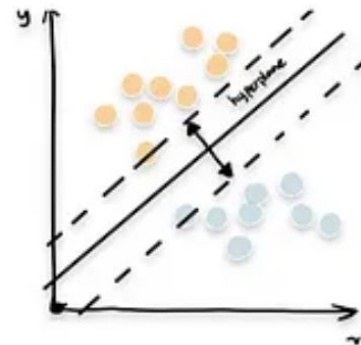
Supervised machine learning models

**Objective:** To find a good model to predict if a patient has heart diseases or not
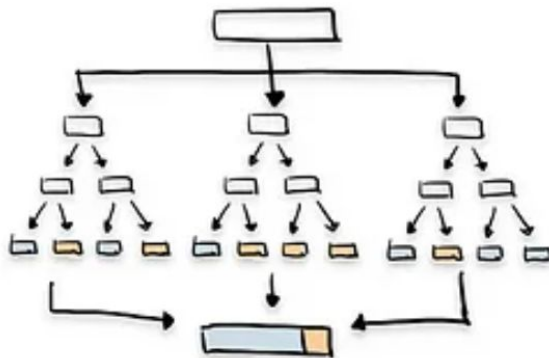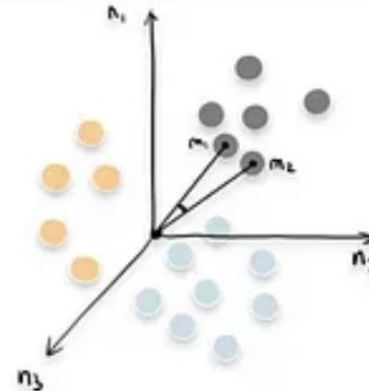
# Models

Logistic regression
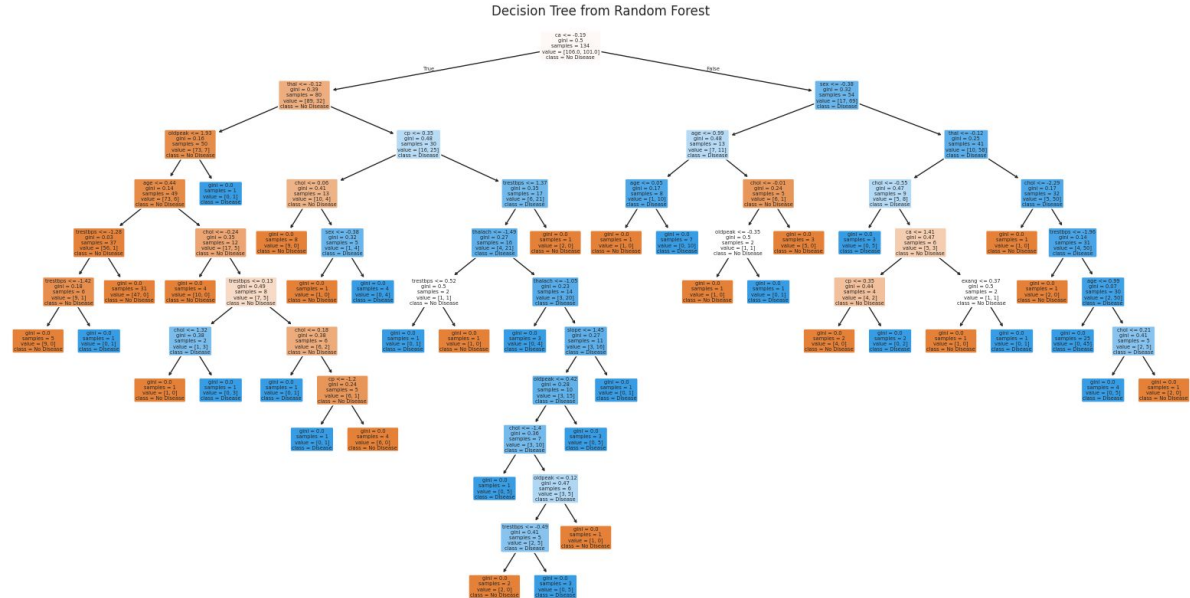
Support Vector
Machine
(SVM)

Random Forest

K Nearest
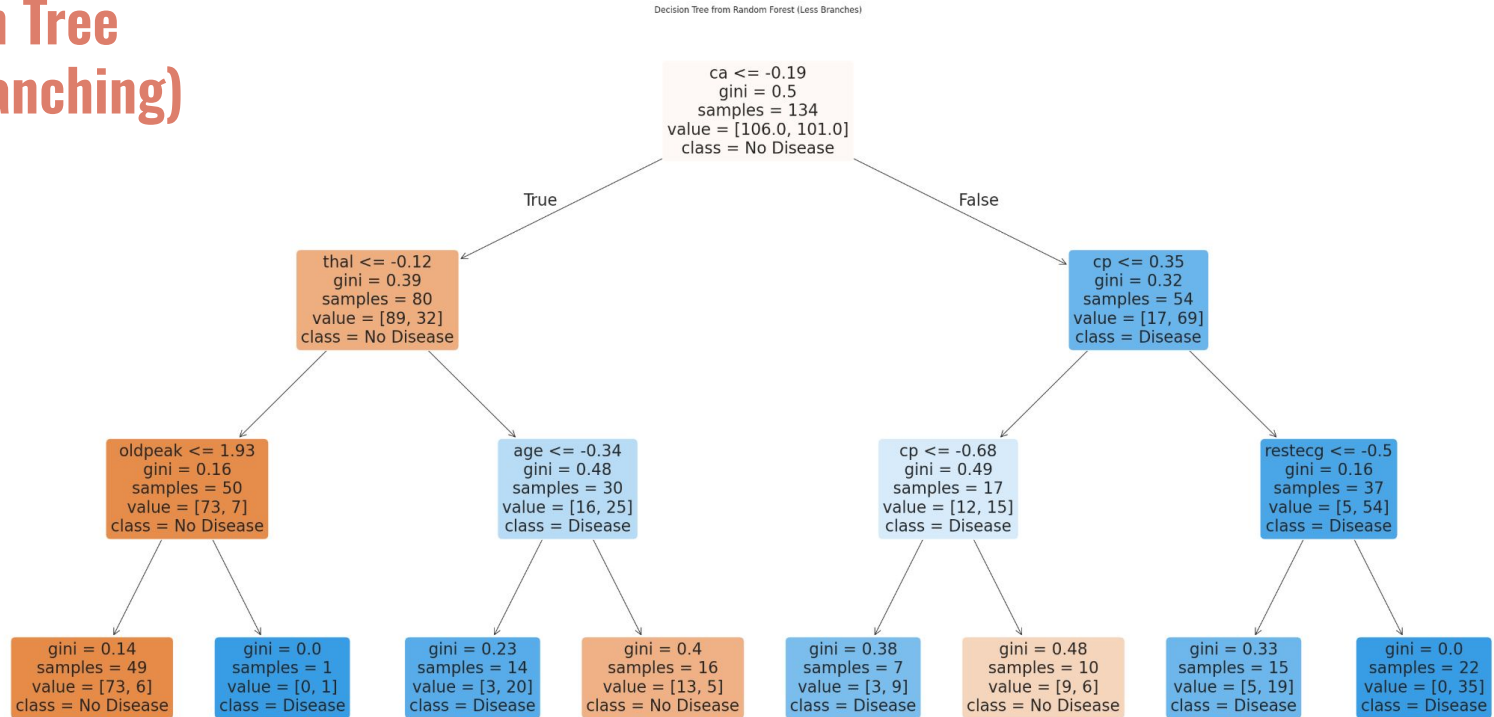Neighbor
(KNN)

Images from: https://towardsdatascience.com/top-machine-learning-algorithms-for-classification-2197870ff501

# Decision Tree

- Can be interpreted by physician.
- Although reducing the branching may be necessary for the understanding and to limit overfitting.
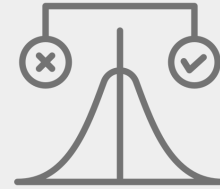


Decision Tree from Random Forest

# Decision Tree (less branching)



Decision Tree from Random Forest (Less Branches)

ca <= -0.19
gini = 0.5
samples = 134
value = [106.0, 101.0]
class = No Disease

True

False

thal <= -0.12
gini = 0.39
samples = 80
value = [89, 32]
class = No Disease

cp <= 0.35
gini = 0.32
samples = 54
value = [17, 69]
class = Disease

oldpeak <= 1.93
gini = 0.16
samples = 50
value = [73, 7]
class = No Disease

age <= -0.34
gini = 0.48
samples = 30
value = [16, 25]
class = Disease

cp <= -0.68
gini = 0.49
samples = 17
value = [12, 15]
class = Disease

restecg <= -0.5
gini = 0.16
samples = 37
value = [5, 54]
class = Disease

gini = 0.14
samples = 49
value = [73, 6]
class = No Disease

gini = 0.0
samples = 1
value = [0, 1]
class = Disease

gini = 0.23
samples = 14
value = [3, 20]
class = Disease

gini = 0.4
samples = 16
value = [13, 5]
class = No Disease

gini = 0.38
samples = 7
value = [3, 9]
class = Disease

gini = 0.48
samples = 10
value = [9, 6]
class = No Disease

gini = 0.33
samples = 15
value = [5, 19]
class = Disease
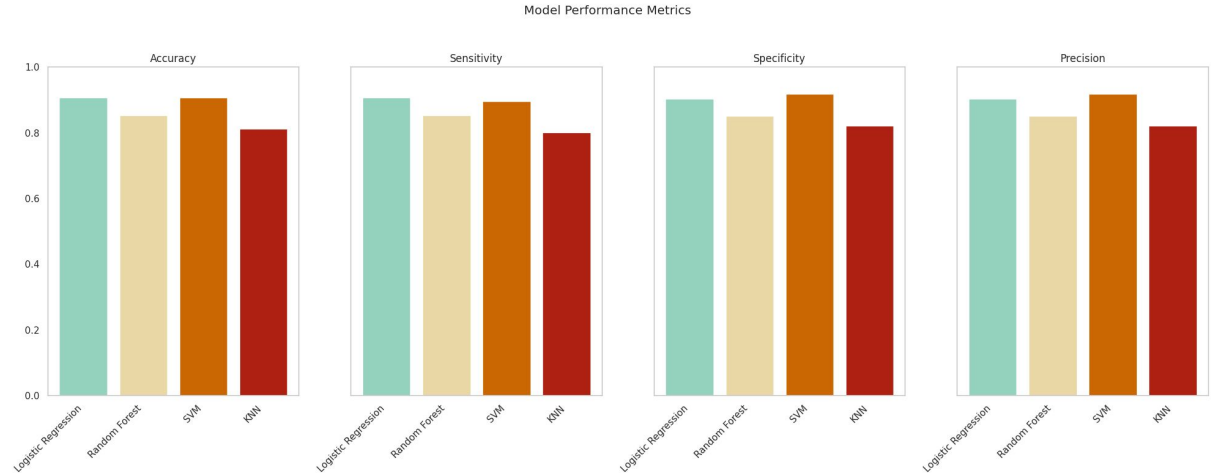
gini = 0.0
samples = 22
value = [0, 35]
class = Disease

# 4. Model evaluation

**Objective:** To understand the performance of each selected models with the provided dataset
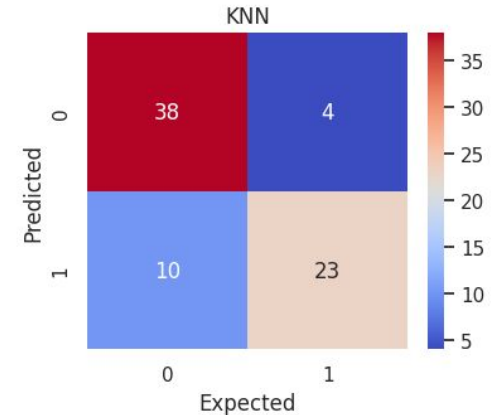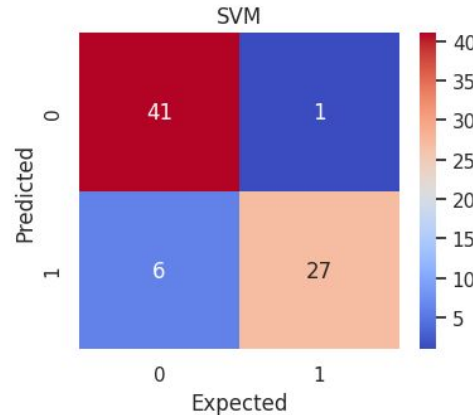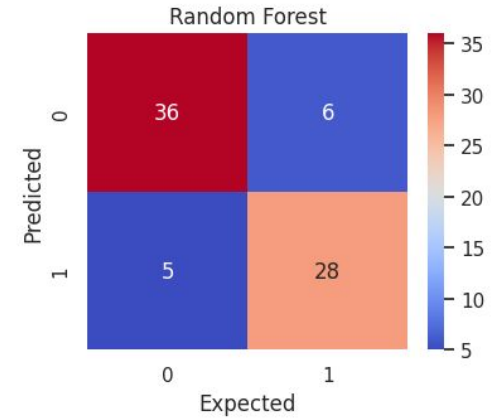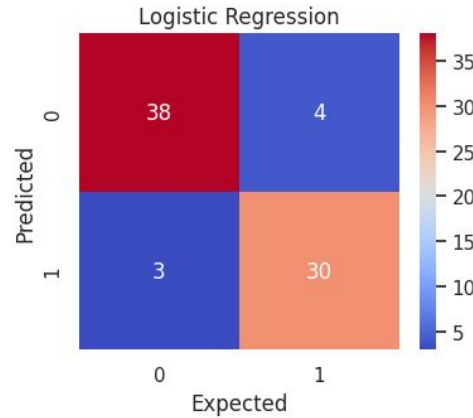
# Model comparison (model metrics)

- Logistic regression and SVM perform similarly.

- Random Forest metrics are slightly better than KNN classifier.
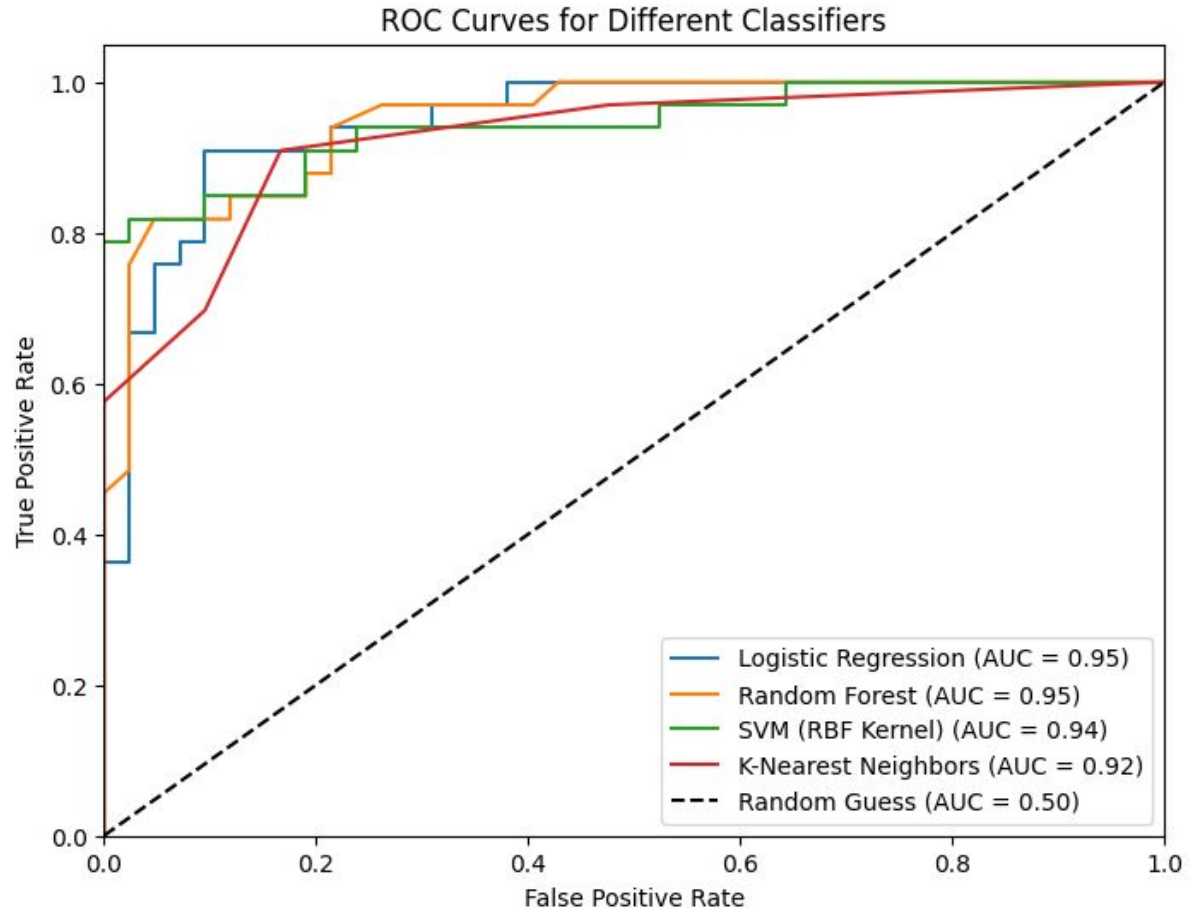


Model Performance Metrics

# Model comparison (confusion matrix)

- Logistic regression and SVM perform similarly.

- Random Forest metrics are slightly better than KNN classifier.

# Model comparison (ROC)

- Logistic regression and SVM perform similarly.

- Random Forest metrics are slightly better than KNN classifier.



ROC Curves for Different Classifiers

# 5. Insights and Conclusion

**Predictive Features**: Blocked vessels (`ca`), ST depression (`oldpeak`), maximum heart rate (`thalach`) are strong indicators of heart disease (corr. value >= 0.4)

Chest pain type (`cp`), thalassemia (`thal`) is not even it has high corr. value because the number not reflect the severity, but just for categorizing.

**Binary Classification:** Converting to a binary target improves performance metrics, focusing on diagnosing disease presence rather than categorizing severity.

**Model Evaluation:** Random Forest, Logistic Regression, and SVM provide reliable predictions.

High AUC scores from Random Forest and SVM suggesting these models may well-suited for real-world clinical applications.

**Clinical Application:** These insights could aid in early identification of at-risk patients, helping healthcare providers prioritize testing and intervention.

# Thank you

- This is the end -

# Chest Pain

Patients divided in:

- 1: Typical angina
- 2: Atypical Anginal Pain.
- 3: Non Anginal Pain.
- 4: Asymptomatic.