# Data Analysis and Visualization Project Report

## Heart Disease Prediction

Mehreen Junaid

15 December 2023

# 1. Introduction

## 1.1 Background

Cardiovascular diseases represent a pervasive global health challenge, contributing significantly to morbidity and mortality. Timely detection and accurate diagnosis of heart disease are imperative for effective intervention and improved patient outcomes. Advances in data science and machine learning techniques offer unprecedented opportunities to analyze extensive datasets, providing valuable insights into the intricate relationships between various health indicators and the likelihood of cardiovascular events.

## 1.2 Motivation

The motivation behind this project lies in the pursuit of a deeper understanding of the factors influencing heart disease. By leveraging a comprehensive dataset comprising demographic information, clinical metrics, and a target variable indicating the probability of heart disease, we aim to unravel patterns and correlations that can inform preventive strategies and enhance diagnostic precision. The increasing prevalence of heart-related ailments underscores the importance of data-driven approaches in augmenting our ability to identify at-risk individuals and tailor interventions to specific health profiles.

## 1.3 Objectives

This project seeks to conduct a thorough analysis and visualization of key attributes within the heart disease dataset, shedding light on factors that may contribute to or mitigate the risk of heart attacks. By exploring demographic variables, medical readings, and their interactions, our goal is to generate actionable insights that can aid healthcare professionals, policymakers, and researchers in refining existing diagnostic approaches and developing targeted preventive measures.

## 1.4 Significance of Analysis

Understanding the intricate interplay between age, gender, lifestyle factors, and clinical metrics is crucial for developing personalized and effective strategies for cardiovascular health. The utilization of a diverse dataset, with privacy measures in place, ensures the ethical handling of patient information while allowing for robust analyses that can be generalized across populations. The insights derived from this analysis can potentially contribute to the refinement of existing predictive models and the

identification of novel risk factors, ultimately fostering advancements in cardiovascular healthcare.

## 1.5 Structure of the Report

This report encompasses a comprehensive examination of the heart disease dataset, including dataset details, insights from visualizations, correlations among attributes, and a discussion of relevant findings. Additionally, it references a seminal research paper, "International application of a new probability algorithm for the diagnosis of coronary artery disease," providing context for the dataset's origin and its clinical relevance. The report concludes with a discussion of future research directions and a reference section citing the sources utilized in the analysis.

## 2. Dataset Details

## 2.1 Dataset Overview

The dataset encompasses a total of 76 attributes, with a subset of 14 key attributes utilized in published experiments, notably concentrating on the Cleveland database. Of particular interest is the "goal" attribute, quantifying the presence of heart disease and ranging from 0 (indicating no presence) to 4.

It's noteworthy that privacy measures have been implemented, with the names and social security numbers of patients replaced by dummy values. One file, housing the Cleveland database, has undergone processing, while four other unprocessed files coexist in the directory. Additionally, a folder labeled "Costs" contains Test Costs donated by Peter Turney.

The dataset includes diverse attributes such as age, sex, exercise-induced angina, the number of major vessels, chest pain type, resting blood pressure, cholesterol level, fasting blood sugar, resting electrocardiographic results, maximum heart rate achieved, and the target variable signifying the likelihood of a heart attack (0= less chance, 1= more chance).

A discriminant function model, derived from clinical and noninvasive test results of 303 patients at the Cleveland Clinic, has been tested for reliability and clinical utility across three patient test groups in Budapest, Hungary; Long Beach, California; and Zurich and Basel, Switzerland. Notably, the age attribute ranges from 0 to 77, providing essential demographic information. The results of applying the Cleveland algorithm were compared with a Bayesian algorithm (CADENZA), revealing tendencies to over-predict the probability of disease in certain centers.

Clinical utility assessments demonstrated modestly superior performance of the new discriminant function in the Hungarian group and similar performance in the American and Swiss groups compared to CADENZA. The conclusion drawn is that coronary disease probabilities derived from discriminant functions prove reliable and clinically useful, particularly in patients with chest pain syndromes and intermediate disease prevalence.

## 2.2 Data Source

The dataset is sourced from the Cleveland database, with privacy measures implemented for patient confidentiality. Dummy values replace names and social security numbers.

## 2.3 Attribute Information

Key attributes include age, sex, exercise-induced angina, the number of major vessels, chest pain type, resting blood pressure, cholesterol level, fasting blood sugar, resting electrocardiographic results, maximum heart rate achieved, and the target variable indicating the likelihood of a heart attack (0= less chance, 1= more chance).

## 3. Research Paper Reference

## 3.1 Paper Title

"International application of a new probability algorithm for the diagnosis of coronary artery disease."

## 3.2 Authors and Publication Details

Detrano, R., Jánosi, A., et al. (1989). Published in the American Journal of Cardiology. Link to Paper

## 3.3 Summary of Paper

The paper introduces a probability algorithm for diagnosing coronary artery disease and compares its performance with a Bayesian algorithm (CADENZA).

## 4. Detailed Analysis of Results and Visualizations

## 4.1 Exploratory Data Analysis

## 4.1.1 Target Variable Distribution

The dataset contains a subset of 14 key attributes, with the "goal" attribute quantifying the presence of heart disease, ranging from 0 (indicating no presence) to 4. Upon analyzing the distribution of the target variable, we observed a balanced dataset, with approximately 54% indicating the presence

of heart disease. This balance facilitates effective model training without the need for extensive data-balancing techniques.

## 4.1.2 Demographic Insights

### 4.1.2.1 Age Distribution

The age distribution reveals a predominant age range of 50-60 years. This insight provides essential demographic information, emphasizing the representation of individuals in a key age group with a potential risk of heart disease.

### 4.1.2.2 Gender Distribution

The dataset exhibits a predominance of males, underscoring the importance of gender as a demographic factor in heart disease analysis. Males show a higher likelihood of experiencing heart attacks compared to females.

## 4.2 Feature Distribution and Relationships

### 2.2.1 Chest Pain and Heart Attack

A noteworthy finding is the strong correlation between chest pain and the likelihood of a heart attack. Individuals experiencing chest pain, particularly those with higher intensity, are more probable to suffer from heart disease.

### 4.2.2 Blood Pressure and Cholesterol

Distribution plots indicate that individuals with cholesterol levels in the range of 120-250 and blood pressure readings between 110 and 140 are more likely to experience a heart attack. This aligns with established medical knowledge regarding the risk factors associated with cardiovascular diseases.

### 4.2.3 Heart Rate and Heart Attack

The analysis emphasizes the significant correlation between heart rate and heart disease. Individuals with higher heart rates are more likely to suffer from a heart attack, providing a valuable indicator for risk assessment.

### 4.2.4 Exercise-Induced Angina

Individuals without exercise-induced angina are highly likely to suffer from a heart attack. This underscores the potential relevance of exercise-induced angina as a factor influencing the likelihood of a heart attack.

## 4.3 Correlation Heatmap

The correlation heatmap reveals valuable insights into the relationships among different attributes. Positive correlations are observed with chest pain, heart rate, and slope, while negative correlations exist with age, exercise-induced angina, and the number of major vessels.

## 4.4 3D Scatter Plot

The 3D scatter plot provides a visual representation of the relationship between heart attack occurrence, heart rate, and age. It reinforces the observation that individuals with higher heart rates are more likely to experience a heart attack, irrespective of age.

## 4.5 Box Plots and Swarm Plots

Box plots and swarm plots illustrate the correlation of chest pain type, exercise-induced angina, and age with heart rate and heart attack occurrence. These visualizations aid in identifying patterns and outliers, offering nuanced insights into the dataset.

## 5. Discussion of Key Findings

## 5.1 Age and Heart Rate as Critical Indicators

The combined consideration of age and heart rate identifies individuals aged 40-60 with elevated heart rates as a potential high-risk group for heart disease. These findings provide valuable information for targeted preventive strategies.

## 5.2 Chest Pain as a Strong Indicator

Chest pain emerges as a strong indicator or symptom associated with cardiovascular issues. Individuals experiencing chest pain, especially of higher intensity, exhibit a higher probability of suffering from heart disease.

## 5.3 Gender Disparities

Males show a higher likelihood of experiencing heart attacks compared to females. This underscores the importance of gender-specific risk assessments and interventions in cardiovascular healthcare.

## 5.4 Exercise-Induced Angina

The absence of exercise-induced angina is highlighted as a significant factor influencing the likelihood of experiencing a heart attack. This observation emphasizes the potential relevance of exercise-induced angina in risk assessment.

## 6. Conclusion

This comprehensive analysis of the heart disease dataset has unveiled key insights into the factors influencing the likelihood of heart attacks. From demographic patterns to intricate feature correlations, our exploration of the Cleveland database has provided a nuanced understanding of cardiovascular health. The dataset's richness allowed us to uncover associations between

age, gender, chest pain, and various physiological indicators, offering a holistic view of the risk landscape.

## 6.1 Key Findings

In examining demographic trends, we discovered a heightened susceptibility to heart disease among individuals aged 50-60, underlining the significance of age as a crucial risk factor. Gender disparities also surfaced, with males exhibiting a higher propensity for heart attacks than females. Moving beyond demographics, chest pain emerged as a robust indicator, showcasing a strong correlation with the likelihood of a heart attack. Distribution plots for cholesterol, blood pressure, and heart rate further refined our understanding, pinpointing specific ranges associated with elevated probabilities of heart attacks.

## 6.2 Implications and Recommendations

The observed gender disparities call for tailored healthcare strategies, acknowledging the unique risk profiles of males and females. To enhance preventive measures, particularly for those reporting chest pain, swift and accurate diagnosis should be prioritized. The absence of exercise-induced angina emerged as a notable factor influencing the likelihood of a heart attack, warranting further research into its role in cardiovascular risk assessments.

## 6.3 Future Research Directions

As we conclude, avenues for future research emerge. Exploring genetic factors, socioeconomic influences, and lifestyle choices can enrich our understanding of heart disease. Incorporating longitudinal data and leveraging advanced machine learning models may enhance predictive accuracy, contributing to more personalized risk assessments in the realm of cardiology.

## 6.4 Conclusion Statement

In essence, this analysis contributes significantly to the collective knowledge base in cardiology. From demographic insights to predictive models, our findings provide a foundation for informed decision-making in preventive healthcare. As we navigate the complexities of heart disease, the implications drawn from this study are poised to shape strategies, improve patient outcomes, and drive advancements in cardiovascular research.

## 7. References

- **Data Source:** [UCI Machine Learning Repository - Heart Disease Dataset](#)
- **Research Paper:** "International application of a new probability algorithm for the diagnosis of coronary artery disease" - Detrano, J., Janosi, A., Steinbrunn, W. et al.
- **Additional Paper:** [Probability algorithm for diagnosis of coronary artery disease - American Journal of Cardiology](#)
- **Dataset source from kaggle:** https://www.kaggle.com/datasets/rashikrahmanpritom/heart-attack-analysis-prediction-dataset/data?select=heart.csv