

heart_disease_project_doc

November 21, 2025

0.1 DSCI 522 – Milestone 1 - Group 25

0.1.1 Project Name: Heart Disease Prediction Model

0.1.2 Team Members:

Johnson Chuang | Eduardo Sanches | Azadeh Ramesh | Jose Davila

Data analysis and workflow project for DSCI 522 (Data Science Workflows), a course in the Master of Data Science program at the University of British Columbia.

git hub link: <https://github.com/stoyq/heart-disease-predictor>

0.1.3 Summary

Heart disease is one of the leading causes of death globally, and early detection is critical for prevention and treatment. In this project, we use the UCI Heart Disease dataset to build a machine-learning model that predicts whether a patient is likely to have heart disease based on clinical and physiological attributes. We load the dataset directly from the web, clean and wrangle the data, perform exploratory data analysis (EDA), and train a classification model (Decision Tree) to identify important predictors of heart disease. Our results highlight key risk indicators that align with well-known medical knowledge, demonstrating how machine learning can support early screening and clinical decision-making.

0.1.4 Introduction

The objective of this project is to develop a predictive model that determines whether a patient is at risk of heart disease using a set of clinical measurements. Heart disease diagnoses often rely on many interacting factors such as chest pain symptoms, blood pressure, cholesterol levels, and exercise response. Machine-learning models can help uncover patterns in these variables and support early identification of high-risk patients.

Our research question is:

“Given a patient’s clinical and physiological attributes, can we accurately predict whether they have heart disease?”

To answer this question, we use the publicly available Heart Disease dataset from the UCI Machine Learning Repository. This dataset contains multiple medically relevant variables, making it suitable for a classification model such as a Decision Tree.

0.1.5 Discussion

The Decision Tree model was able to identify meaningful patterns to predict heart disease based on the data, with a test score of 0.61 and train score of 0.78. Based on these results, it might indicate that there was some overfitting based on the large difference between training and test results.

From the EDA, we see that various features such as age, sex, chol and more have clear differences in their distribution between disease and no disease which will help the model to predict between the two. For a better predictor, we may want to incorporate additional features given the complexity of heart disease.

0.1.6 Methods & Results:

We build a machine-learning classification model using the UCI Heart Disease dataset:

1. Load data from the original source on the web: <https://archive.ics.uci.edu/dataset/45/heart+disease>
2. Wrangle and clean the data
 - Replace missing values
 - Assign meaningful column names
 - Convert categorical variables to numeric where needed
 - Ensure that the target variable is binary (0 = no heart disease, 1 = heart disease)
3. Perform exploratory data analysis (EDA)
 - Summary statistics for continuous variables
 - Count plots for categorical variables
 - Histograms and boxplots to understand feature distributions
4. Create visualizations relevant to the classification task
 - Pairplots to explore relationships between key features
 - Distribution of target classes
 - Feature correlation matrix
5. Build a classification model
 - A Decision Tree Classifier is trained to predict heart disease.
 - We split the dataset into training and testing subsets and evaluate model accuracy.
6. Visualize the model results
 - Plot of the trained Decision Tree
 - Feature importance bar chart

0.1.7 Dataset Description

We use the Heart Disease dataset from the UCI Machine Learning Repository, a widely used benchmark dataset for medical prediction tasks. The dataset includes the following 14 attributes: - Age - Sex - Chest Pain Type (cp) - Resting Blood Pressure (trestbps) - Cholesterol (chol) - Fasting Blood Sugar (fbs) - Resting ECG results (restecg) - Maximum heart rate achieved (thalach) - Exercise induced angina (exang) - ST depression (oldpeak) - Slope of ST segment (slope) - Number of major vessels (ca) - Thalassemia result (thal) - num (Target: the predicted attribute (0 = no heart disease, 1 = heart disease))

These variables include both continuous and categorical measurements commonly used in clinical diagnostics.

0.1.8 Results and Conclusion

Our analysis shows that several clinical features differ noticeably between patients with and without heart disease. As seen in the EDA histograms, patients with heart disease tend to have higher resting blood pressure (trestbps), higher ST-depression values (oldpeak), and lower maximum heart rate achieved (thalach) compared to individuals without disease. After preprocessing the dataset using scaling for numerical variables and one-hot encoding for categorical variables, we trained a Support Vector Classifier (SVC) model. Cross-validation results indicate an average test accuracy of 0.61, with a higher training accuracy of 0.78, suggesting some overfitting. When evaluating predictions on the unseen test set, the model correctly identified many cases but also showed several misclassifications, especially where the model predicted “0” (no disease) but the true label was “1” or “2.” Overall, while the model captures meaningful patterns in the dataset, its moderate predictive performance suggests that further tuning, alternative models, or feature engineering may be needed to improve accuracy and reduce classification bias.

0.1.9 References

- UCI Machine Learning Repository. Heart Disease Dataset:
<https://archive.ics.uci.edu/dataset/45/heart+disease>
- International application of a new probability algorithm for the diagnosis of coronary artery disease. By R. Detrano, A. Jánosi, W. Steinbrunn, M. Pfisterer, J. Schmid, S. Sandhu, K. Guppy, S. Lee, V. Froelicher. 1989 Published in American Journal of Cardiology