

Heart Disease Prediction

Namit Patel,Smit patel,Yogi Patel

April 22, 2024

1 Problem Definition

Heart disease remains one of the leading causes of morbidity and mortality worldwide. Early detection and accurate prediction of heart disease are crucial for effective management and prevention of adverse outcomes. The goal of this project is to develop a predictive model capable of accurately diagnosing the presence and severity of heart disease based on clinical and demographic features. By leveraging machine learning algorithms, the model aims to analyze patterns within a dataset containing patient information such as age, sex, chest pain type, blood pressure, cholesterol levels, and other relevant attributes. The predictive model will classify individuals into different severity levels of heart disease, ranging from no disease to critical conditions. The ultimate objective is to assist healthcare professionals in making timely and informed decisions, leading to improved patient outcomes and reduced cardiovascular morbidity and mortality rates.

2 Dataset Description - Heart Disease Prediction Dataset

2.1 Overview

The Heart Disease Prediction dataset is a collection of clinical data from three different locations: Hungary, Cleveland, and Switzerland and VA long Beach. It contains information about individuals who have been diagnosed with various cardiac conditions, with the aim of facilitating predictive modeling for heart disease. The dataset used has 921 rows and 30 columns.

2.2 Numerical Columns

1. ID
2. Age
3. Trestbps (Resting Blood Pressure)
4. Chol (Serum Cholesterol)
5. Thalch (Maximum Heart Rate Achieved)
6. Oldpeak (ST Depression Induced by Exercise Relative to Rest)
7. CA (Number of Major Vessels Colored by Fluoroscopy)
8. Num (Diagnosis)

2.3 Categorical Columns

1. Sex
2. Dataset
3. CP (Chest Pain Type)
4. FBS (Fasting Blood Sugar)
5. Restecg (Resting Electrocardiographic Results)
6. Exang (Exercise Induced Angina)
7. Slope (Slope of the Peak Exercise ST Segment)
8. Thal (Thallium Stress Test Result)

2.4 Target Column

The target column is "Num," which represents the predicted attribute for heart disease diagnosis. Here's the mapping of the unique values in the "Num" column:

1. 0: No heart disease
2. 1: Mild heart disease
3. 2: Moderate heart disease
4. 3: Severe heart disease
5. 4: Critical heart disease

These values provide the severity levels of heart disease, ranging from no disease to critical condition. This information will be essential for training predictive models to classify the severity of heart disease based on the input features.

3 Attributes

- **ID:** Unique identifier for each patient.
- **Age:** Age of the patient in years.
- **Sex:** Gender of the patient (Male/Female).
- **Dataset:** Location where the data was collected (Hungary, Cleveland, or Switzerland).
- **Chest Pain Type (CP):** Type of chest pain experienced by the patient, categorized as typical angina, atypical angina, non-anginal, or asymptomatic.
- **Resting Blood Pressure (Trestbps):** Resting blood pressure measured in mm Hg.
- **Serum Cholesterol (Chol):** Serum cholesterol level measured in mg/dl.
- **Fasting Blood Sugar (FBS):** Indicates whether the fasting blood sugar level is above 120 mg/dl (TRUE/FALSE).
- **Resting Electrocardiographic Results (Restecg):** Results of the resting electrocardiogram, indicating the presence of left ventricular hypertrophy or normal.

- **Maximum Heart Rate Achieved (Thalch):** Maximum heart rate achieved during exercise.
- **Exercise Induced Angina (Exang):** Indicates whether angina was induced by exercise (TRUE/FALSE).
- **ST Depression Induced by Exercise Relative to Rest (Oldpeak):** ST depression induced by exercise relative to rest.
- **Slope of the Peak Exercise ST Segment (Slope):** Slope of the peak exercise ST segment, categorized as upsloping, flat, or downsloping.
- **Number of Major Vessels Colored by Fluoroscopy (CA):** Number of major vessels colored by fluoroscopy (0-3).
- **Thallium Stress Test Result (Thal):** Result of the thallium stress test, categorized as normal, fixed defect, or reversible defect.
- **Diagnosis (Num):** Diagnosis of heart disease (0: No presence, 1-4: Presence).

4 Purpose

The dataset serves as a valuable resource for predicting the presence or absence of heart disease based on various clinical and demographic factors. By analyzing this dataset, we can develop predictive models to identify individuals at risk of cardiovascular disorders and implement preventive measures accordingly.

5 Usage

We can utilize this dataset for: Developing machine learning models for heart disease prediction. Investigating the effectiveness of different diagnostic tests and risk factors in predicting heart disease. Understanding geographical variations in cardiac health and disease prevalence. Informing public health policies and interventions aimed at reducing the burden of cardiovascular diseases.

6 Novelty

- **Data Import:** The notebook begins by importing the necessary libraries, such as pandas, numpy, and matplotlib.
- **Data Loading:** The dataset is loaded using pandas' `read_csv` function.
- **Data Cleaning:** The notebook includes code for cleaning the data, such as removing unnecessary columns, handling missing values, and converting categorical data to numerical data.
- **Exploratory Data Analysis (EDA):** The code conducts EDA by examining the dataset's characteristics, such as the age distribution, sex distribution, and unique values in the 'dataset' column. This EDA helps in understanding the data before proceeding with predictive modeling.
- **Gender Analysis:** A unique aspect is the calculation of heart disease percentages based on gender. The code calculates and compares the percentage of males and females with heart disease in the dataset, highlighting the gender disparity in heart disease prevalence.
- **Visualization:** The code includes a histogram plot to visualize the distribution of ages in the dataset. This visual representation enhances the understanding of the age distribution within the dataset.
- **Dataset Analysis:** By exploring the unique values and value counts in the 'dataset' column, the code provides insights into the different datasets included in the analysis, such as 'Cleveland', 'Hungary', 'Switzerland', and 'VA Long Beach'.

- **Age and Gender Analysis:** The code further delves into the dataset by examining the value counts of the 'age' column grouped by 'sex', offering a detailed breakdown of age distribution based on gender.
- **Missing Values and Outliers:** The code utilizes various techniques to handle missing values, such as imputation methods like SimpleImputer, KNNImputer, and IterativeImputer. It also explores the presence of missing values in different columns like 'trestbps', 'chol', 'fbs', 'restecg', 'thalch', 'exang', 'oldpeak', 'slope', 'ca', and 'thal'.
- **Machine Learning Models:** The code incorporates a range of machine learning models for heart disease prediction, including Logistic Regression, K-Nearest Neighbors, Support Vector Machines, Naive Bayes, Decision Trees, Random Forest, AdaBoost, Gradient Boosting, and XGBoost classifiers. It also covers metrics like accuracy score, R2 score, confusion matrix, mean squared error, and mean absolute error to evaluate the performance of these models. Additionally, the notebook utilizes a pipeline approach for model building and includes the MLxtend library for frequent pattern mining using Apriori algorithm.
- **Model Evaluation and Interpretation::** There is a huge importance of model evaluation and interpretation in assessing the predictive model's performance and reliability. We have used the metrics for evaluating model performance, such as accuracy, R-squared coefficient, and model complexity plots, and interpreted their implications in the context of graduate admission prediction.

Overall, the code showcases a comprehensive analysis of the heart disease dataset, incorporating EDA, gender-based analysis, visualization, and dataset-specific insights to gain a deeper understanding of the data and potentially aid in predictive modeling for heart disease.

7 Algorithms

Algorithm 1 Data Cleaning

Input: Data from a CSV file

Output: Cleaned data

```

1: procedure CLEAN DATA(data)
2:   data  $\leftarrow$  data without unnecessary columns
3:   data  $\leftarrow$  data with missing values replaced with NaN
4: end procedure

```

Algorithm 2 Removing Outliers

Input: Data with Outlier**Output:** Data without Outlier

```
1: procedure REMOVE_OUTLIER(data)
2:   Analysis of Boxplot
3:   if data has outlier
4:     if data has outlier then
5:       Compute  $Q1$  and  $Q3$ 
6:        $IQR \leftarrow Q3 - Q1$ 
7:        $min \leftarrow Q1 - 1.5 \times IQR$ 
8:        $max \leftarrow Q3 + 1.5 \times IQR$ 
9:       for value in data do
10:        if value  $\leq$  min then
11:          value  $\leftarrow$  min
12:        end if
13:        if value  $\geq$  max then
14:          value  $\leftarrow$  max
15:        end if
16:      end for
17:    end if
18: end procedure
```

Algorithm 3 Pre-processing the Data

Input: Data**Output:** Processed Data

```
1: procedure PRE-PROCESSING(data)
2:   if feature is categorical then
3:     Impute data by mode of the column
4:     Make a copy of the imputed columns
5:     Drop the original columns
6:     Transform the categorical data using LabelEncoder method
7:   end if
8:   if feature is numeric then
9:     Round off the data to 3 decimal places
10:    Impute the data by mean of the column
11:    Apply PCA to Data
12:   end if
13: end procedure
```

Algorithm 4 Exploratory Data Analysis (EDA)

```
1: procedure EDA(data)
2:   1. Summary Statistics
3:   Compute basic summary statistics (mean, median, mode, standard deviation, min, max, quartiles, etc.) for each variable.
4:   2. Data Visualization
5:   Plot histograms for each variable to visualize the distribution.
6:   Plot boxplots for each variable to identify outliers and compare distributions.
7:   Create scatter plots for pairs of variables to explore potential relationships.
8:   Create correlation heatmaps to examine correlations between variables.
9:   3. Missing Values Analysis
10:  Identify missing values in the data.
11:  Compute the percentage of missing values for each variable.
12:  4. Data Cleaning
13:  Handle missing values (e.g., imputation, deletion).
14:  Handle outliers (e.g., capping, transformation, removal).
15:  5. Data Transformation
16:  Normalize or standardize numerical variables, if necessary.
17:  Encode categorical variables (e.g., one-hot encoding, label encoding).
18:  6. Feature Engineering
19:  Generate new features from existing ones, if applicable.
20:  7. Preliminary Insights
21:  Identify key trends, patterns, or anomalies in the data.
22:  Document key findings and insights for further analysis.
23: end procedure
```

Algorithm 5 Algorithm for Feature Importance

Input: Processed Data

Output: Reduced Selected Data

```
1: procedure FEATURE_IMPORTANCE(data)
2:   Split the data into train and test sets.
3:   Initialize a Random Forest Tree Classifier model.
4:   Train the model on the training set.
5:   Evaluate the model's performance on the test set.
6:   Retrieve the importance of features from the trained model.
7:   Plot a graph of the feature importance.
8: end procedure
```

Algorithm 6 Impute Categorical Missing Data

```
1: procedure IMPUTE_CATEGORICAL_MISSING_DATA(passed_col)
2:   1. Split the DataFrame into two parts: df_null containing rows with missing values in passed_col, and df_not_null
   containing rows without missing values in passed_col
3:   2. Define X as the features from df_not_null without passed_col, and y as the target column passed_col
4:   3. Define other_missing_cols as a list of columns with missing values except passed_col
5:   4. Initialize label_encoder as a LabelEncoder instance
6:   5. Encode categorical features in X using label_encoder
7:   6. Instantiate an IterativeImputer with RandomForestRegressor as estimator
8:   7. Impute missing values in other_missing_cols in X using IterativeImputer
9:   8. Split X and y into training and testing sets using train_test_split
10:  9. Instantiate RandomForestClassifier as rf_classifier
11:  10. Fit rf_classifier on X_train and y_train
12:  11. Predict y_test using rf_classifier
13:  12. Calculate accuracy_score between y_test and y_pred
14:  13. Print the accuracy_score
15:  14. Replace missing values in df_null[passed_col] with predictions from rf_classifier
16:  15. Concatenate df_not_null and df_null to get df_combined
17:  16. Return df_combined[passed_col]
18: end procedure
```

Algorithm 7 Heart Disease Prediction: Model Training and Evaluation

```
1: procedure MODEL TRAINING AND EVALUATION(Data)
2:   1. Train Models
3:   Train Logistic Regression
4:   Train K-Nearest Neighbors (KNN)
5:   Train Support Vector Machine (SVM)
6:   Train Naive Bayes
7:   Train Decision Tree
8:   Train Random Forest
9:   Train AdaBoost model
10:  Train Gradient Boosting model
11:  Train XGBoost model
12:  2. Evaluate Models
13:  Evaluate each model using the testing data and the following metrics:
14:    Accuracy score
15:    Confusion matrix
16:    R2 score
17:    Mean squared error (MSE)
18:    Mean absolute error (MAE)
19:  3. Model Comparison and Selection
20:  Compare the performance of each model using evaluation metrics.
21:  Select the model with the best performance for deployment.
22:  4. Frequent Pattern Mining (using MLxtend)
23:  Utilize the Apriori algorithm from the MLxtend library to mine frequent patterns from the data.
24:  Analyze the patterns to identify associations between different features.
25:  5. Document Findings
26:  Summarize the performance of different models and chosen model.
27:  Document frequent patterns and potential insights from the mining process.
28: end procedure
```
