

Machine Learning Project Report on Cardiovascular Disease Prediction

Xavier Sánchez Mateus
David Galindo Munté
Maxence Martin

December 17, 2023

1 Introduction

1.1 Subject presentation

This report documents our project's progression in predicting cardiovascular diseases through machine learning techniques. It encompasses our thorough analysis, model training, and exploration of advanced methodologies aimed at enhancing prediction accuracy.

Our journey commenced with an in-depth examination of the dataset, involving meticulous data cleaning and understanding. Subsequently, we delved into training diverse machine learning models, refining them over time. Additionally, we ventured into advanced techniques like ensemble learning, seeking avenues to bolster our predictive capabilities for cardiovascular diseases.

The primary objective of this report is to provide a step-by-step account of our methodologies and findings.

1.2 Presentation of the Data

| Feature Name | Feature Type | Data Type | Description |
|--------------|---------------------|------------------|--|
| age | Objective Feature | int (days) | Age in days |
| height | Objective Feature | int (cm) | Height in centimeters |
| weight | Objective Feature | float (kg) | Weight in kilograms |
| gender | Objective Feature | categorical code | Categorical code for gender (e.g., 1 for male, 2 for female) |
| ap_hi | Examination Feature | int | Systolic blood pressure |
| ap_lo | Examination Feature | int | Diastolic blood pressure |
| cholesterol | Examination Feature | categorical code | 1: normal, 2: above normal, 3: well above normal cholesterol |
| gluc | Examination Feature | categorical code | 1: normal, 2: above normal, 3: well above normal glucose |
| smoke | Subjective Feature | binary | Binary indicator for smoking (1 for yes, 0 for no) |
| alco | Subjective Feature | binary | Binary indicator for alcohol intake (1 for yes, 0 for no) |
| active | Subjective Feature | binary | Binary indicator for physical activity (1 for yes, 0 for no) |
| cardio | Target Variable | binary | Binary indicator for the presence or absence of cardiovascular disease |

Table 1: Presentation of the Data

2 Preliminary Analysis and Data Cleaning

2.1 Preliminary Data Analysis

This section is crucial for understanding our data and, therefore, enhancing the performance of training and

prediction. First, to analyze the data properly, we separated the dataset features into categorical and continuous. This allowed us to create various graphs to easily visualize the data.

Initially, we observed the distribution of data based on categorical features. Through these graphs, we noticed that the categorical target "cardio" is balanced, meaning it has an equal number of sick and healthy individuals. There is a significant imbalance in gender data, with twice as many men as women. This could be explained by the fact that men are more affected by cardiovascular diseases, making it easier to find positive cases. However, within each gender, there are equal numbers of sick and healthy individuals. We learned that 80% of our sample exercises, 95% do not drink alcohol, and 91% do not smoke. Additionally, 85% of our sample has a normal glucose level, and 74.8% have a normal cholesterol level.

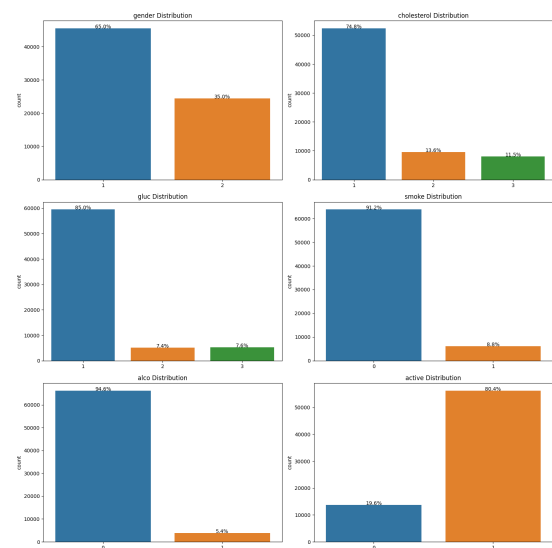


Figure 1: Histogram of Categorical Features

2.2 Cleaning the dataset

The representation of continuous data was more useful. Using boxplots, we identified numerous outliers in the dataset. We adopted two approaches to remove them from the training set. The goal of this removal was to help the model better adapt to the general trend of the data and improve its ability to generalize to new observations.

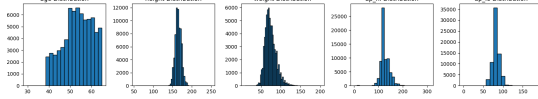


Figure 2: Distribution of Numerical Features

Firstly, we eliminated values that were deemed impossible based on different sources and scientific articles. This was done by setting lower and upper bounds for each category and removing all rows with values outside these bounds. After this, we generated the first dataset, which we would use later for modeling. Before setting it aside, we plotted the histogram of the target feature to check if, after removing rows with impossible values, the dataset remained balanced. In our case, this is particularly concerning because it is likely that abnormal values come from unhealthy individuals, and removing these rows could significantly reduce the number of people with cardiovascular diseases. However, after analyzing the graph, we could see that the dataset remained balanced.

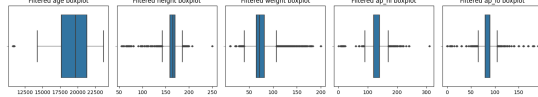


Figure 3: All Boxplots after filtering impossible values

Secondly, we used the dataset without impossible values to build a second dataset, this time without extreme values. For this dataset, we removed outliers in a more conventional way. Specifically, an outlier is considered low if it is below ($Q1 - 1.5 * \text{interquartile range}$) and high if it is above ($Q3 + 1.5 * \text{interquartile range}$). With this manipulation, we obtained our final dataset.

2.3 Principal Component Analysis

Finally, we attempted to reduce the number of features in our dataset by performing Principal Component Analysis (PCA). The idea was to simplify the input features, making the algorithm training faster and

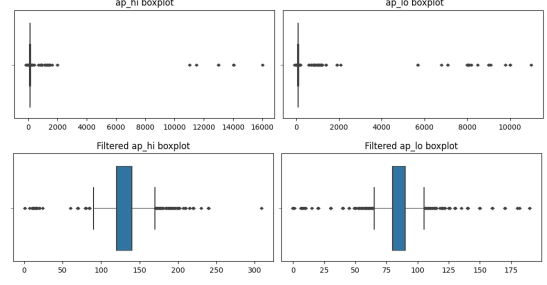


Figure 4: Boxplot of Systolic and Diastolic Pressure Before and After Filtering Impossible Values

allowing us to try higher-degree polynomial transformations. However, after conducting PCA, we found that there was no benefit in performing PCA on our data, as the variance explained by each axis was very low. Moreover, the cumulative explained variance exceeded 97% with 10 components out of 11. Thus, the trade-off between reducing the number of features and preserving information was not worthwhile. In general, the cumulative explained variance graph is nearly linear, without a clear elbow, which is characteristic of a dataset whose dimensionality cannot be reduced without significant information loss.

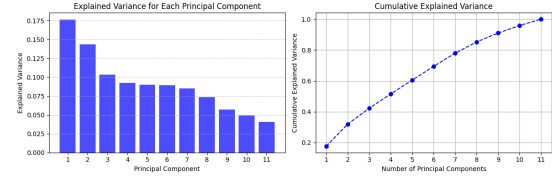


Figure 5: Explained Variance of Each Component and Cumulative Explained Variance

We will train using three datasets: the original dataset serves as a baseline, the second excludes impossible values, and the third removes impossible and extreme values. This approach aims to evaluate the impact of removing outliers on model performance and determine whether excluding outliers improves the results.

3 Training and Model Performance Evaluation

3.1 Preparation of Training and Test Data

Before training our models, we preprocessed our data to facilitate learning for our models. We separated cat-

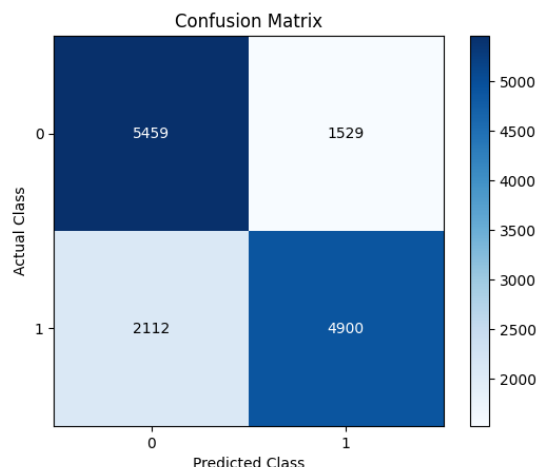


Figure 6: Confusion Matrix of the Best Classifier

threshold. A high AUC indicates superior discriminatory ability, while a value close to 0.5 suggests performance equivalent to random chance. These combined metrics allow for a thorough evaluation of the quality of our models.

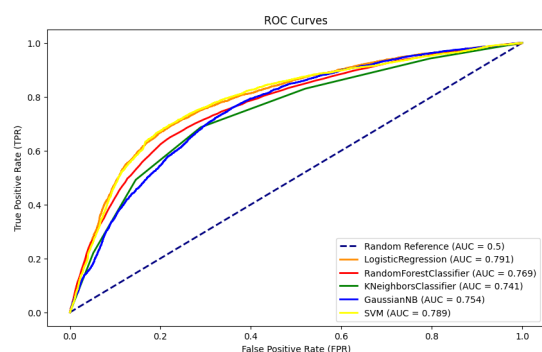


Figure 7: ROC curves for classifiers with default settings, degree 1, and training dataset without extreme values

3.4 Exploring Ensemble Learning Algorithms: Boosting and Voting

In order to try to improve the performance of our previous models we explored the possibilities offered by ensemble learning. So we trained several AdaBoost meta-classifiers with some of our previous base classifiers. We also trained several gradient boosting and voting classifiers with one “hard” vote and one “soft” vote. All results are available in Excel ‘results_ensemble_learning.xlsx’. It appears from all these experiments that we still can-

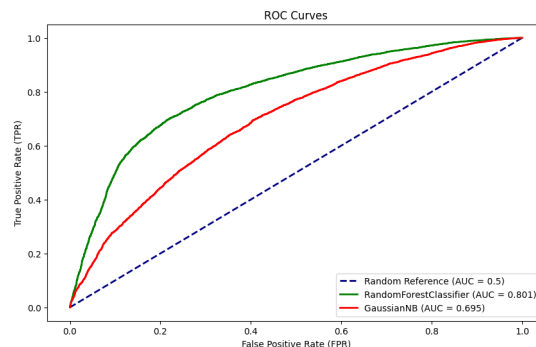


Figure 8: ROC curves of our best and worst classifiers to visualize the difference in performance

not exceed the fateful threshold of 74% with the boosting and voting methods.

4 Conclusion

Despite the efforts invested in this project, the accuracy settled at a somewhat disappointing plateau of 74%. However, upon exhaustive exploration of the Kaggle platform, the primary source of our dataset, it became evident that this performance level might represent the ceiling attainable within the constraints of this specific dataset.

Our project not hitting the expected accuracy shows how tough it is to predict complicated health problems. It really shows how tricky it is to make predictions in healthcare. We found lots of complicated things and uncertainties while trying to predict health issues using the dataset we had. This made it hard to get very precise predictions.

Furthermore, it shows how important good data is for making accurate predictions. The quality of our predictions depends a lot on how good, complete, and relevant the dataset we use is. This teaches us how much better our predictions can be if we pay close attention to getting the right data for healthcare predictions.

In essence, our project journey, though not hitting the accuracy we wanted, shows how hard it is to predict tricky health illnesses. It tells us how important good data is for making predictions that work well. This experience pushes us to keep digging, getting better, and trying new things to make healthcare predictions more accurate.