

CS5525 Final Submission Report

Jennifer Appiah-Kubi, Rebecca DeSipio, Ajinkya Fotedar

12/11/2021

Contents

I. Introduction	1
II. Classification Methods Explored	2
Decision Trees	2
KNN	4
Support Vector Machines (SVM)	4
Logistic Regression	4
III. Analysis	4
Best Model	4
Comparison of Models	4
IV. Conclusion	5
V. Notes and References	5

I. Introduction

The goal of this project was to predict the probability of having a heart attack using 14 variables given in the heart.csv data-set. The classification models chosen to analyze this data-set were: random forest, bagging, support vector machines (SVM), and k-nearest neighbor (KNN). After exploring these various classification methods, we can analyze and interpret the results of each method to determine which might be the “best” classifier. Additionally, the logistic regression methods, lasso and elastic net, were taken into consideration to explore which variables are most important. (add sentence on the importance of this)

Attribute Information:

- **age**
- **sex**
- **cp**: chest pain type (4 values)
- **trestbps**: resting blood pressure
- **chol**: serum cholesterol in mg/dl
- **fbs**: fasting blood sugar > 120 mg/dl

- **restecg**: resting electrocardiograph results (values 0, 1, 2)
- **thalach**: maximum heart rate achieved
- **exang**: exercise induced angina
- **oldpeak**: ST depression induced by exercise relative to rest
- **slope**: the slope of the peak exercise ST segment
- **ca**: number of major vessels (0 - 3) colored by fluoroscope
- **thal**: thalassemia (blood disorder) 0 = normal; 1 = fixed defect; 2 = reversible defect
- **target**: 0 = less chance of heart attack; 1 = more chance of heart attack

II. Classification Methods Explored

Decision Trees

Decision trees are useful in classification to improve prediction accuracy. We explored bagging and random forest models. Figure 1 shows a classification tree used to predict whether someone is more or less likely to have a heart attack based on certain thresholds for each of the 14 predictors.

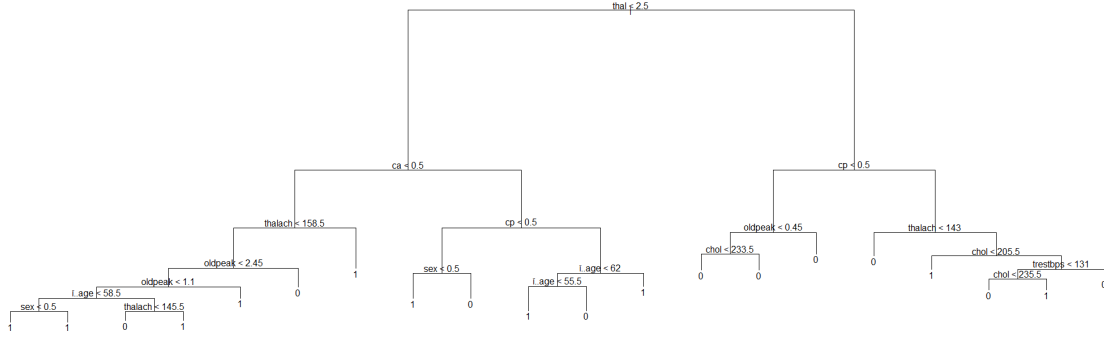


Figure 1: Classification tree for **heart** data-set.

Next, we decided to prune the classification tree, shown in figure 2.

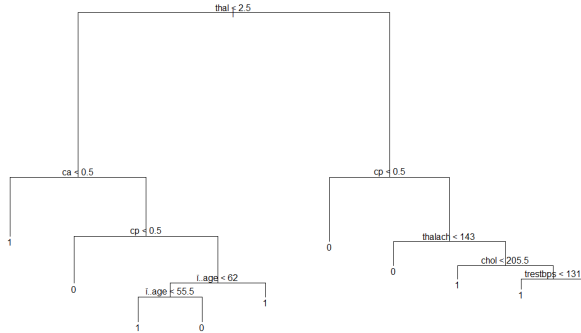


Figure 2: Pruned classification tree.

We can also look at the prediction accuracy (??) for both the original (see figure 3) and pruned trees and compare the outputted results.

	Target.test	
tree.pred	0	1
	0	21 10
	1	8 37

	Target.test	
prune.pred	0	1
	0	22 6
	1	7 41

Figure 3: Prediction accuracy of the original tree (left) compared to the pruned tree (right).

Bagging

Next, we used bagging, which is useful to reduce the variance and improve the prediction accuracy. The results of bagging are shown in figure 4. Given these results, we can conclude that thalassemia and chest pain type of the two most important attributes.

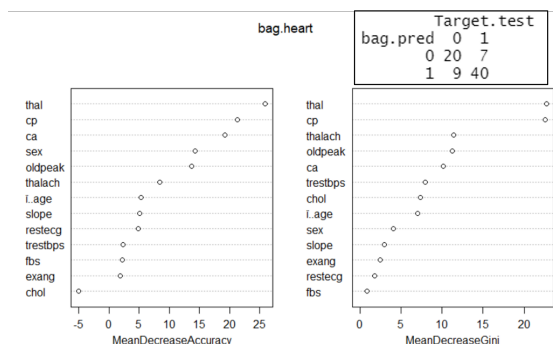


Figure 4: The side-by-side plots show the variables of importance in heart attack predictability. The table in the top-right shows the prediction accuracy for bagging.

Random Forest

Finally, random forest was used to improve upon bagging and further improve prediction accuracy, shown in figure 5. We can see that after performing random forest, the variables of importance have changed. While thalassemia and chest pain type are still important, the random forest analysis shows that maximum heart rate achieved and number of major vessels are also of importance.

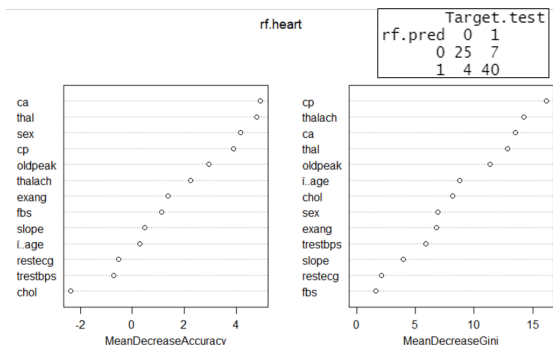


Figure 5: The side-by-side plots show the variables of importance in heart attack predictability. The table in the top-right shows the prediction accuracy for random forest.

KNN

Support Vector Machines (SVM)

Logistic Regression

Lasso

Elastic Net

III. Analysis

Best Model

Random Forest Plots

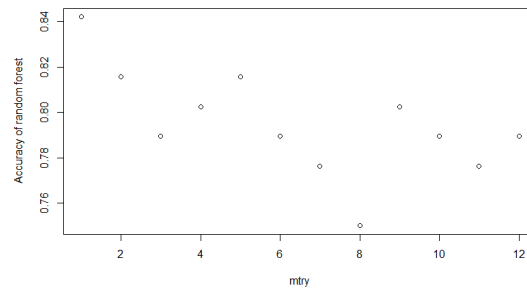


Figure 6: Accuracy plot.

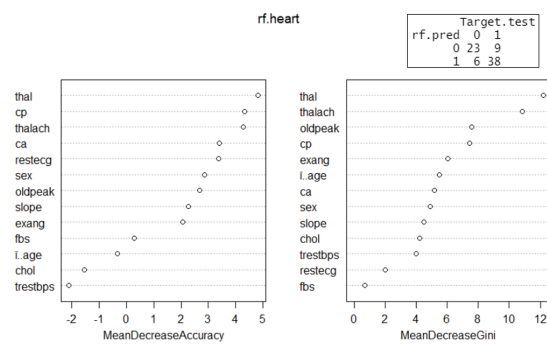


Figure 7: Random Forest variables of importance.

{A couple paragraphs on determining a best model for each classification method - should have 6 total}

Comparison of Models

{Determine then the best model by looking at prediction accuracy and MSE (for logistic regression)}

IV. Conclusion

V. Notes and References