

# CS5525 Final Submission Report

Jennifer Appiah-Kubi, Rebecca DeSipio, Ajinkya Fotedar

12/11/2021

## Contents

<b>I. Introduction</b>	<b>1</b>
<b>II. Classification Methods Explored</b>	<b>2</b>
Decision Trees . . . . .	2
Support Vector Machines . . . . .	4
KNN . . . . .	5
Logistic Regression . . . . .	5
<b>III. Analysis</b>	<b>7</b>
Best Model . . . . .	7
Comparison of Models . . . . .	7
<b>IV. Conclusion</b>	<b>7</b>
<b>V. Notes and References</b>	<b>7</b>

## 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



## 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.

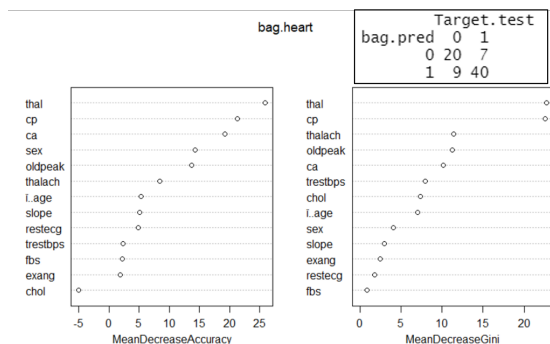


Figure 4: The side-by-side plots show the variables of importance in heart attack predictability. The table in the top-right shows the prediction accuracy for bagging.

## 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.

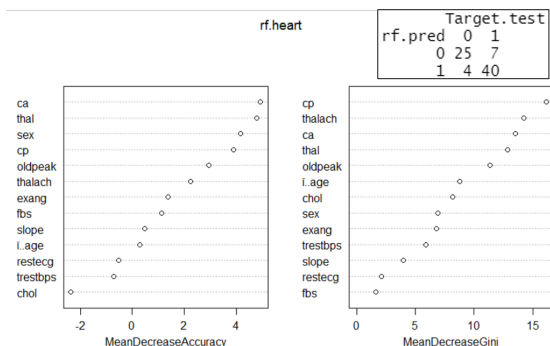


Figure 5: The side-by-side plots show the variables of importance in heart attack predictability. The table in the top-right shows the prediction accuracy for random forest.

## Support Vector Machines

SVM is a supervised classification method that separates data using *hyperplanes*, which act as a decision boundary between the various classes. Here we use a classifier for predicting whether a patient is suffering from any heart disease or not.

### Model Training and Accuracy:

We first train a linear classifier with the help of a *train control* method that uses *repeated cross-validation*, and then calculate prediction accuracy (82.2%) using a *confusion matrix*.

```
## Confusion Matrix and Statistics
##
##
## svm.pred  0  1
##          0 31  4
##          1 12 43
##
##              Accuracy : 0.8222
##              95% CI : (0.7274, 0.8948)
##              No Information Rate : 0.5222
##              P-Value [Acc > NIR] : 2.711e-09
##
##              Kappa : 0.6409
##
##  Mcnemar's Test P-Value : 0.08012
##
##              Sensitivity : 0.7209
##              Specificity : 0.9149
##              Pos Pred Value : 0.8857
##              Neg Pred Value : 0.7818
##              Prevalence : 0.4778
##              Detection Rate : 0.3444
##              Detection Prevalence : 0.3889
##              Balanced Accuracy : 0.8179
##
##              'Positive' Class : 0
##
```

Figure 6: Confusion Matrix with Cost = 1

### Choosing Different Costs:

In order to improve model performance, we play with *Cost* ( $C$ ) values in our classifier. For this we define a grid with specific  $C$  values. We then train our model again using the new  $C$  values. Our prediction accuracy is the most (83.5%) when  $C = 0.01$ , reflected in the plot below.

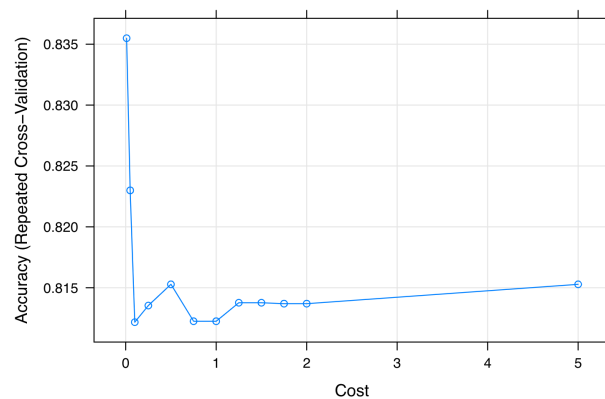


Figure 7: Accuracy Plot with Varying Costs

## Tuned Model:

Finally, we test the model for the same  $C$  values. We do this by using *predict* over the *tuned* training model and the testing data-set, and checking the accuracy using the *confusion matrix*. We note that this ends up giving a higher accuracy rate (83.3%).

```
## Confusion Matrix and Statistics
##
##
## svm.pred.grid 0 1
##               0 30 2
##               1 13 45
##
##               Accuracy : 0.8333
##               95% CI : (0.74, 0.9036)
##               No Information Rate : 0.5222
##               P-Value [Acc > NIR] : 6.187e-10
##
##               Kappa : 0.6623
##
## Mcnemar's Test P-Value : 0.009823
##
##               Sensitivity : 0.6977
##               Specificity : 0.9574
##               Pos Pred Value : 0.9375
##               Neg Pred Value : 0.7759
##               Prevalence : 0.4778
##               Detection Rate : 0.3333
##               Detection Prevalence : 0.3556
##               Balanced Accuracy : 0.8276
##
##               'Positive' Class : 0
##
```

Figure 8: Confusion Matrix of the Tuned Model

## KNN

### Logistic Regression

LR predicts whether something is *True* or *False* instead of predicting something continuous. Here, we try to predict the probability that a person will get a heart disease or not.

### Pre-processing Data:

For accurate predictions, we first process the raw data and convert a significant chunk of the variables into factors.

```
## 'data.frame': 303 obs. of 14 variables:
## $ age : int 63 37 41 56 57 57 56 44 52 57 ...
## $ sex : int 1 1 0 1 0 1 0 1 1 1 ...
## $ cp : int 3 2 1 1 0 0 1 1 2 2 ...
## $ trestbps: int 145 130 130 120 120 140 140 120 172 150 ...
## $ chol : int 233 250 204 236 354 192 294 263 199 168 ...
## $ fbs : int 1 0 0 0 0 0 0 1 0 ...
## $ restecg: int 0 1 0 1 1 1 0 1 1 1 ...
## $ thalach: int 150 187 172 178 163 148 153 173 162 174 ...
## $ exang : int 0 0 0 1 0 0 0 0 ...
## $ oldpeak: num 2.3 3.5 1.4 0.8 0.6 0.4 1.3 0 0.5 1.6 ...
## $ slope : int 0 0 2 2 1 1 2 2 2 ...
## $ ca : int 0 0 0 0 0 0 0 0 ...
## $ thal : int 1 2 2 2 2 1 2 3 3 2 ...
## $ target : int 1 1 1 1 1 1 1 1 1 ...

## 'data.frame': 303 obs. of 14 variables:
## $ age : num 63 37 41 56 57 57 56 44 52 57 ...
## $ sex : Factor w/ 2 levels "F","M": 2 2 1 2 1 2 1 2 2 ...
## $ cp : Factor w/ 4 levels "0","1","2","3": 4 3 2 2 1 1 2 2 3 ...
## $ trestbps: num 145 130 130 120 120 140 140 120 172 150 ...
## $ chol : num 233 250 204 236 354 192 294 263 199 168 ...
## $ fbs : Factor w/ 2 levels "0","1": 2 1 1 1 1 1 1 2 1 ...
## $ restecg: Factor w/ 3 levels "0","1","2": 1 2 1 2 2 2 1 2 2 ...
## $ thalach: num 150 187 172 178 163 148 153 173 162 174 ...
## $ exang : Factor w/ 2 levels "0","1": 1 1 1 1 2 1 1 1 1 ...
## $ oldpeak: num 2.3 3.5 1.4 0.8 0.6 0.4 1.3 0 0.5 1.6 ...
## $ slope : Factor w/ 3 levels "0","1","2": 1 1 3 3 2 2 3 3 3 ...
## $ ca : Factor w/ 5 levels "0","1","2","3": 1 1 1 1 1 1 1 1 ...
## $ thal : Factor w/ 4 levels "0","1","2","3": 2 3 3 3 3 2 3 4 4 3 ...
## $ target : Factor w/ 2 levels "Healthy","Unhealthy": 2 2 2 2 2 2 2 2 ...
```

Figure 9: Raw (left) and Processed (right) Data

### Comparing Models:

We now compared the  $R^2$  and BIC values to two models - one with only age as the independent variable and another with all the variables. The second model is the better model since it has a higher  $R^2$  and a lower

BIC, which is an indicator of a better fit. We note that since the median age in our data-set is 55, it makes sense why it *is not* a statistically significant variable in our complex model (that covers all variables).

<pre>R_sq_1 &lt;- 1 - logistic\$deviance / logistic\$null.deviance R_sq_1  ## [1] 0.05947945 BIC_1 &lt;- logistic\$deviance + 2 * log(dim(data)[1]) BIC_1  ## [1] 404.2246</pre>	<pre>R_sq_2 &lt;- 1 - logistic\$deviance / logistic\$null.deviance R_sq_2  ## [1] 0.569889 BIC_2 &lt;- logistic\$deviance + 14 * log(dim(data)[1]) BIC_2  ## [1] 259.623</pre>
--	--

Figure 10: Simple (left) and Complex (right) Models

### Predicting Probability of Heart Disease:

Finally, we plot the probability of predicting whether a person in our data-set has a heart disease. The upper-right portion of the logistic curve (cyan) shows the predicted probability a person will get a heart disease. The bottom-right portion of the logistic curve (light orange) shows the predicted probability a person will not get a heart disease.

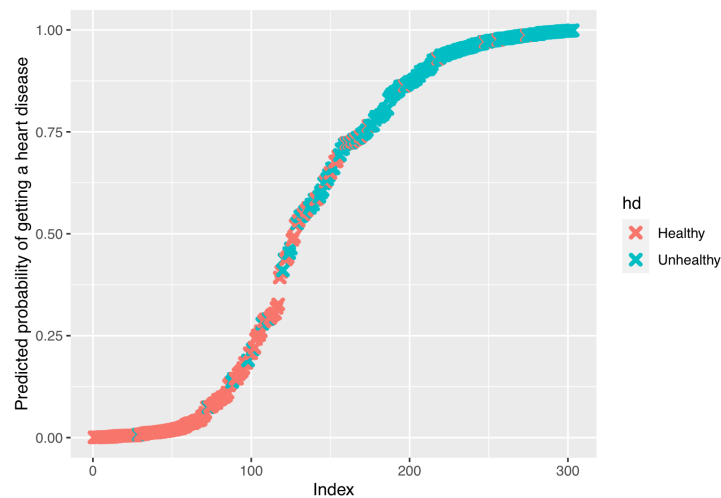


Figure 11: Predicted Heart Disease Probability

### III. Analysis

#### Best Model

Random Forest Plots

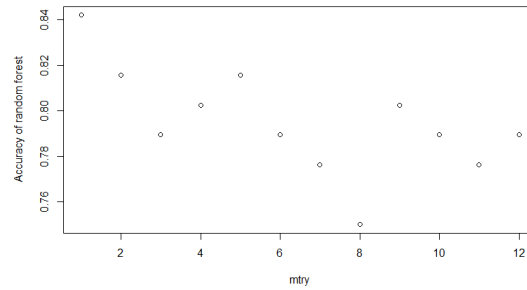


Figure 12: Accuracy plot.

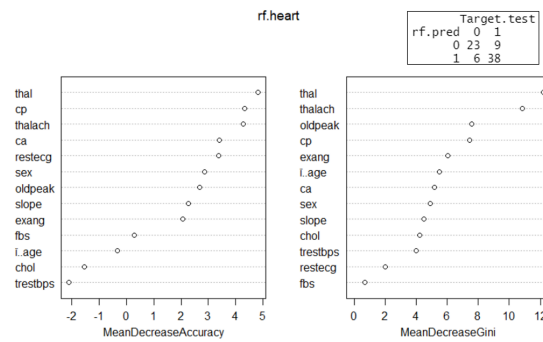


Figure 13: Random Forest variables of importance.

{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