

Heart Disease Prediction Project

Jarred Priester

11/26/2021

Table of Content

1. Overview

- 1.1 Description of dataset
- 1.2 Goal of project
- 1.3 step to achieve the goal

2. Analysis

- 2.1 Downloading the data
- 2.2 Data cleaning
- 2.3 Data exploration
- 2.4 Visualization
- 2.5 Models

3. Results

- 3.1 Results
- 3.2 Brief thoughts about the results

4. Conclusion

- 4.1 Summary
- 4.2 Limitations
- 4.3 Future Work

1. Overview

1.1 Description of dataset

For this project we will be analyzing the heart disease data set from University of California Irvine machine learning repository. This data set consist fo 14 different features and 303 observations. The description of the features from the website is the following:

- *age*: age in years

- **sex**: sex (1 = male; 0 = female)
- **cp**: chest pain type
 - Value 1: typical angina
 - Value 2: atypical angina
 - Value 3: non-anginal pain
- **trestbps**: resting blood pressure (in mm Hg on admission to the hospital)
- **chol**: serum cholestoral in mg/dl
- **fbs**: (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false)
- **restecg**: resting electrocardiographic results
 - Value 0: normal
 - Value 1: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV)
 - Value 2: showing probable or definite left ventricular hypertrophy by Estes' criteria
- **thalach**: maximum heart rate achieved
- **exang**: exercise induced angina (1 = yes; 0 = no)
- **oldpeak** = ST depression induced by exercise relative to rest
- **slope**: the slope of the peak exercise ST segment
 - Value 1: upsloping
 - Value 2: flat
 - Value 3: downsloping
- **ca**: number of major vessels (0-3) colored by flourosopy
- **thal**: 3 = normal; 6 = fixed defect; 7 = reversable defect
- **num**: diagnosis of heart disease (angiographic disease status)
 - Value 0: < 50% diameter narrowing
 - Value 1: > 50% diameter narrowing

The num feature is the feature we will be trying to predict for this project.

Link to the UCI heart disease data: <https://archive.ics.uci.edu/ml/datasets/heart+disease>

1.2 Goal of project

The goal for this project is to create a model that can predict the patient's heart disease status with an overall accuracy of 85% or higher. The other goals will be to explore the data we have been given and find key insights into heart disease that could be helpful for the medical community going forward.

1.3 Step to acheive the goal

To achieve this goal we will be applying 10 different algorithms and comparing their results. Because the nature of the problem is to determine if a patient is negative or positive, i.e 0 or 1, this is a binary classification problem and we will pick 10 algorithms that work well with binary classification. The algorithms we will be using are the following:

- Logistic Regression
- Linear Discriminant Analysis
- Quadratic Discriminant Analysis
- Loess Model
- K-Nearest Neighbors
- Random Forest

- Tree Models from Genetic Algorithms
- Least Squares Support Vector Machine
- Bayesian Generalized Linear Model
- Neural Network

We will split the data into training and test data. We will be using the K-fold cross variation technique in order to test and validate our models so that we make sure to not over train the models. We will then train the models and apply them to the test set giving us our accuracy.

2. Analysis

2.1 Downloading the data

```
#importing libraries
if(!require(tidyverse)) install.packages("tidyverse")
if(!require(caret)) install.packages("dplyr")
if(!require(dplyr)) install.packages("dplyr")
if(!require(matrixStats)) install.packages("matrixStats")
if(!require(gam)) install.packages("gam")
if(!require(evtree)) install.packages("evtree")

library(tidyverse)
library(caret)
library(dplyr)
library(matrixStats)
library(gam)
library(evtree)

#importing the University of California, Irvine Heart Disease Data set
heart <- read.csv("https://archive.ics.uci.edu/ml/machine-learning-databases/heart-disease/processed.cl
names(heart) <- c( "age", "sex", "cp", "trestbps", "chol", "fbs", "restecg",
                  "thalach", "exang", "oldpeak", "slope", "ca", "thal", "num")
```

2.2 Data Cleaning

First we are going to take a look at the data set and get an idea of what we need to clean and how the data set is structured

```
head(heart)
```

```
##   age sex cp trestbps chol fbs restecg thalach exang oldpeak slope ca thal num
## 1  63  1  1    145   233   1       2    150     0     2.3    3  0    6   0
## 2  67  1  4    160   286   0       2    108     1     1.5    2  3    3   2
## 3  67  1  4    120   229   0       2    129     1     2.6    2  2    7   1
## 4  37  1  3    130   250   0       0    187     0     3.5    3  0    3   0
## 5  41  0  2    130   204   0       2    172     0     1.4    1  0    3   0
## 6  56  1  2    120   236   0       0    178     0     0.8    1  0    3   0
```

```
dim(heart)
```

```
## [1] 303 14
```

```
str(heart)
```

```
## 'data.frame': 303 obs. of 14 variables:
## $ age : num 63 67 67 37 41 56 62 57 63 53 ...
## $ sex : num 1 1 1 1 0 1 0 0 1 1 ...
## $ cp : num 1 4 4 3 2 2 4 4 4 4 ...
## $ trestbps: num 145 160 120 130 130 120 140 120 130 140 ...
## $ chol : num 233 286 229 250 204 236 268 354 254 203 ...
## $ fbs : num 1 0 0 0 0 0 0 0 0 1 ...
## $ restecg : num 2 2 2 0 2 0 2 0 2 2 ...
## $ thalach : num 150 108 129 187 172 178 160 163 147 155 ...
## $ exang : num 0 1 1 0 0 0 0 1 0 1 ...
## $ oldpeak : num 2.3 1.5 2.6 3.5 1.4 0.8 3.6 0.6 1.4 3.1 ...
## $ slope : num 3 2 2 3 1 1 3 1 2 3 ...
## $ ca : num 0 3 2 0 0 0 2 0 1 0 ...
## $ thal : num 6 3 7 3 3 3 3 7 7 ...
## $ num : int 0 2 1 0 0 0 3 0 2 1 ...
```

```
summary(heart)
```

```
##      age      sex      cp      trestbps
## Min.   :29.00 Min.   :0.00000 Min.   :1.000 Min.   : 94.0
## 1st Qu.:48.00 1st Qu.:0.00000 1st Qu.:3.000 1st Qu.:120.0
## Median :56.00 Median :1.00000 Median :3.000 Median :130.0
## Mean   :54.44 Mean   :0.6799 Mean   :3.158 Mean   :131.7
## 3rd Qu.:61.00 3rd Qu.:1.00000 3rd Qu.:4.000 3rd Qu.:140.0
## Max.   :77.00 Max.   :1.00000 Max.   :4.000 Max.   :200.0
##
##      chol      fbs      restecg      thalach
## Min.   :126.0 Min.   :0.00000 Min.   :0.00000 Min.   : 71.0
## 1st Qu.:211.0 1st Qu.:0.00000 1st Qu.:0.00000 1st Qu.:133.5
## Median :241.0 Median :0.00000 Median :1.00000 Median :153.0
## Mean   :246.7 Mean   :0.1485 Mean   :0.9901 Mean   :149.6
## 3rd Qu.:275.0 3rd Qu.:0.00000 3rd Qu.:2.00000 3rd Qu.:166.0
## Max.   :564.0 Max.   :1.00000 Max.   :2.00000 Max.   :202.0
##
##      exang      oldpeak      slope      ca
## Min.   :0.00000 Min.   :0.00 Min.   :1.000 Min.   :0.00000
## 1st Qu.:0.00000 1st Qu.:0.00 1st Qu.:1.000 1st Qu.:0.00000
## Median :0.00000 Median :0.80 Median :2.000 Median :0.00000
## Mean   :0.3267 Mean   :1.04 Mean   :1.601 Mean   :0.6722
## 3rd Qu.:1.00000 3rd Qu.:1.60 3rd Qu.:2.000 3rd Qu.:1.00000
## Max.   :1.00000 Max.   :6.20 Max.   :3.000 Max.   :3.00000
##                                     NA's   :4
##      thal      num
## Min.   :3.000 Min.   :0.00000
## 1st Qu.:3.000 1st Qu.:0.00000
## Median :3.000 Median :0.00000
```

```
## Mean :4.734 Mean :0.9373
## 3rd Qu.:7.000 3rd Qu.:2.0000
## Max. :7.000 Max. :4.0000
## NA's :2
```

The feature num is the angiographic disease status. 0 represents no heart disease while 1-4 represents the extent of heart disease in the patient. So for num, we are going to convert anything greater than 0 to equal 1, leaving us with 0 (no heart disease) and 1 (heart disease)

```
heart$num[heart$num > 0] <- 1
```

check to make sure 1-4 were converted to 1

```
summary(heart$num)
```

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.0000 0.0000 0.0000 0.4587 1.0000 1.0000
```

checking to see if the changes were made correctly

```
heart$sex
```

```
## [1] 1 1 1 1 0 1 0 0 1 1 1 0 1 1 1 1 1 0 1 1 0 1 1 1 0 0 0 1 1 0 1 1 1 1 1 1
## [38] 1 1 1 0 1 0 1 0 1 1 1 0 1 0 1 1 1 1 1 1 1 0 0 1 0 1 1 1 1 1 1 1 1 0 1 1 1
## [75] 1 0 1 0 1 1 1 0 1 1 1 1 1 0 0 0 1 0 1 0 0 1 1 0 1 1 1 1 0 0 1 1 1 1 1 1 0
## [112] 1 1 0 0 1 1 0 1 1 1 0 1 1 1 0 0 1 1 0 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1
## [149] 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 0 1 0 1 1 0 0 1 1 1 1 1 1 1 0 1 1 0
## [186] 0 1 1 1 1 1 1 1 0 0 1 1 0 0 1 0 0 1 0 1 1 1 1 1 0 0 1 1 0 1 1 0 0 0 1 0 0
## [223] 0 1 0 0 1 0 1 1 0 0 1 0 0 1 1 1 0 1 1 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 0 0 0 1
## [260] 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 0 1 1 0 0 1 0 1 1 0 1 1 1 0 1 1 1 1 0 1 1 0 1
## [297] 1 0 1 1 1 0 1
```

check for any NAs

```
sum(is.na(heart) == TRUE)
```

```
## [1] 6
```

we have 6 NAs which is not too many so we will delete those rows

```
heart <- na.omit(heart)
dim(heart)
```

```
## [1] 297 14
```

2.3 Data exploration

Let's look at a summary of the data with just the observations without heart disease

```
heart %>% filter(num == 0) %>% summary()
```

```
##      age      sex      cp      trestbps
##  Min.   :29.00  Min.   :0.0000  Min.   :1.000  Min.   : 94.0
##  1st Qu.:44.75  1st Qu.:0.0000  1st Qu.:2.000  1st Qu.:120.0
##  Median :52.00  Median :1.0000  Median :3.000  Median :130.0
##  Mean   :52.64  Mean   :0.5563  Mean   :2.794  Mean   :129.2
##  3rd Qu.:59.00  3rd Qu.:1.0000  3rd Qu.:3.000  3rd Qu.:140.0
##  Max.   :76.00  Max.   :1.0000  Max.   :4.000  Max.   :180.0
##      chol      fbs      restecg      thalach
##  Min.   :126.0  Min.   :0.0000  Min.   :0.0000  Min.   : 96.0
##  1st Qu.:208.8  1st Qu.:0.0000  1st Qu.:0.0000  1st Qu.:149.0
##  Median :235.5  Median :0.0000  Median :0.0000  Median :161.0
##  Mean   :243.5  Mean   :0.1437  Mean   :0.8438  Mean   :158.6
##  3rd Qu.:268.2  3rd Qu.:0.0000  3rd Qu.:2.0000  3rd Qu.:172.0
##  Max.   :564.0  Max.   :1.0000  Max.   :2.0000  Max.   :202.0
##      exang      oldpeak      slope      ca
##  Min.   :0.0000  Min.   :0.0000  Min.   :1.000  Min.   :0.000
##  1st Qu.:0.0000  1st Qu.:0.0000  1st Qu.:1.000  1st Qu.:0.000
##  Median :0.0000  Median :0.2000  Median :1.000  Median :0.000
##  Mean   :0.1437  Mean   :0.5988  Mean   :1.413  Mean   :0.275
##  3rd Qu.:0.0000  3rd Qu.:1.1000  3rd Qu.:2.000  3rd Qu.:0.000
##  Max.   :1.0000  Max.   :4.2000  Max.   :3.000  Max.   :3.000
##      thal      num
##  Min.   :3.000  Min.   :0
##  1st Qu.:3.000  1st Qu.:0
##  Median :3.000  Median :0
##  Mean   :3.788  Mean   :0
##  3rd Qu.:3.000  3rd Qu.:0
##  Max.   :7.000  Max.   :0
```

Now let's take a look at a summary of the data with just the observations with heart disease

```
heart %>% filter(num == 1) %>% summary()
```

```
##      age      sex      cp      trestbps
##  Min.   :35.00  Min.   :0.0000  Min.   :1.000  Min.   :100.0
##  1st Qu.:53.00  1st Qu.:1.0000  1st Qu.:4.000  1st Qu.:120.0
##  Median :58.00  Median :1.0000  Median :4.000  Median :130.0
##  Mean   :56.76  Mean   :0.8175  Mean   :3.584  Mean   :134.6
##  3rd Qu.:62.00  3rd Qu.:1.0000  3rd Qu.:4.000  3rd Qu.:145.0
##  Max.   :77.00  Max.   :1.0000  Max.   :4.000  Max.   :200.0
##      chol      fbs      restecg      thalach
##  Min.   :131.0  Min.   :0.000  Min.   :0.000  Min.   : 71.0
##  1st Qu.:218.0  1st Qu.:0.000  1st Qu.:0.000  1st Qu.:125.0
##  Median :253.0  Median :0.000  Median :2.000  Median :142.0
##  Mean   :251.9  Mean   :0.146  Mean   :1.175  Mean   :139.1
##  3rd Qu.:284.0  3rd Qu.:0.000  3rd Qu.:2.000  3rd Qu.:157.0
##  Max.   :409.0  Max.   :1.000  Max.   :2.000  Max.   :195.0
##      exang      oldpeak      slope      ca
##  Min.   :0.0000  Min.   :0.000  Min.   :1.000  Min.   :0.000
##  1st Qu.:0.0000  1st Qu.:0.600  1st Qu.:1.000  1st Qu.:0.000
```

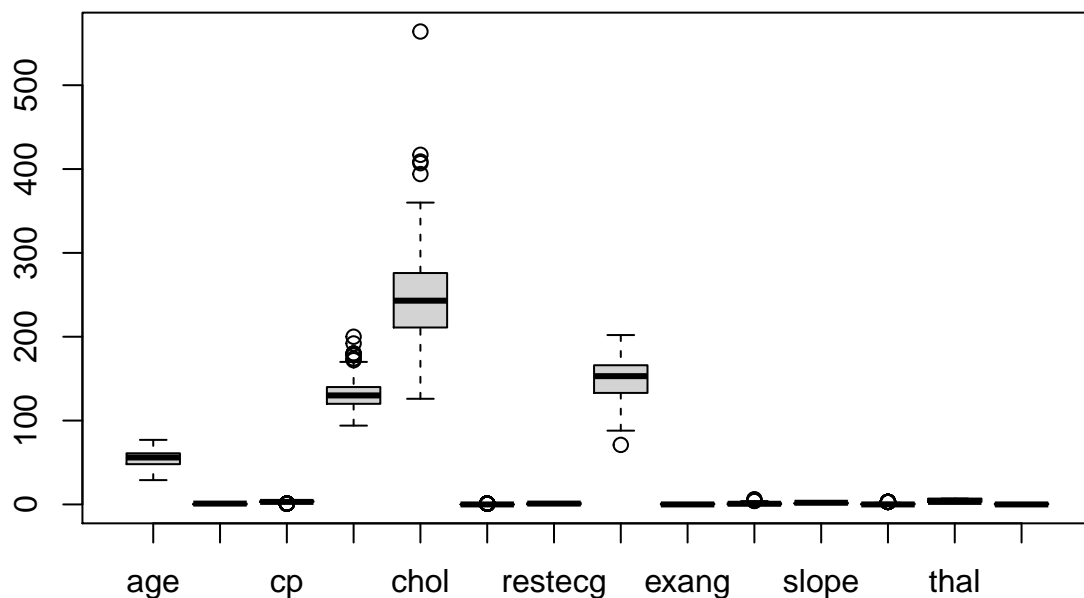
```
## Median :1.0000    Median :1.400    Median :2.000    Median :1.000
## Mean   :0.5401    Mean   :1.589    Mean   :1.825    Mean   :1.146
## 3rd Qu.:1.0000    3rd Qu.:2.500    3rd Qu.:2.000    3rd Qu.:2.000
## Max.   :1.0000    Max.   :6.200    Max.   :3.000    Max.   :3.000
##      thal      num
## Min.   :3.000    Min.   :1
## 1st Qu.:3.000    1st Qu.:1
## Median :7.000    Median :1
## Mean   :5.832    Mean   :1
## 3rd Qu.:7.000    3rd Qu.:1
## Max.   :7.000    Max.   :1
```

There looks to be some differences when you compare the observations for positive heart disease and negative heart disease. For examples, the sex for the positive heart disease leaned more towards male. The age on average was younger for negative heart disease observations. The average maximum heart rate achieved (thalach) was on average higher for the negative heart disease observations. Next let's continue analyzing the data with visualization.

2.4 Visualization

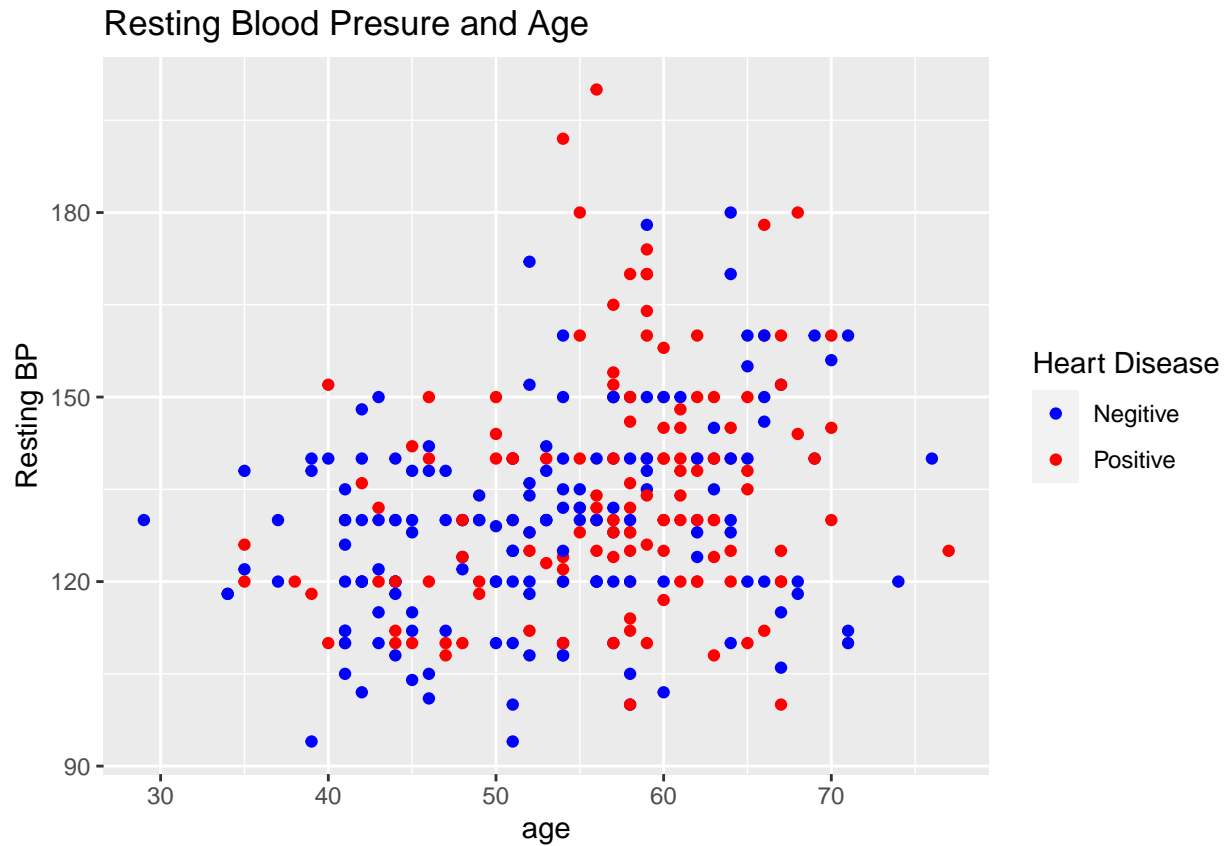
box plot of the data set

```
boxplot(heart)
```



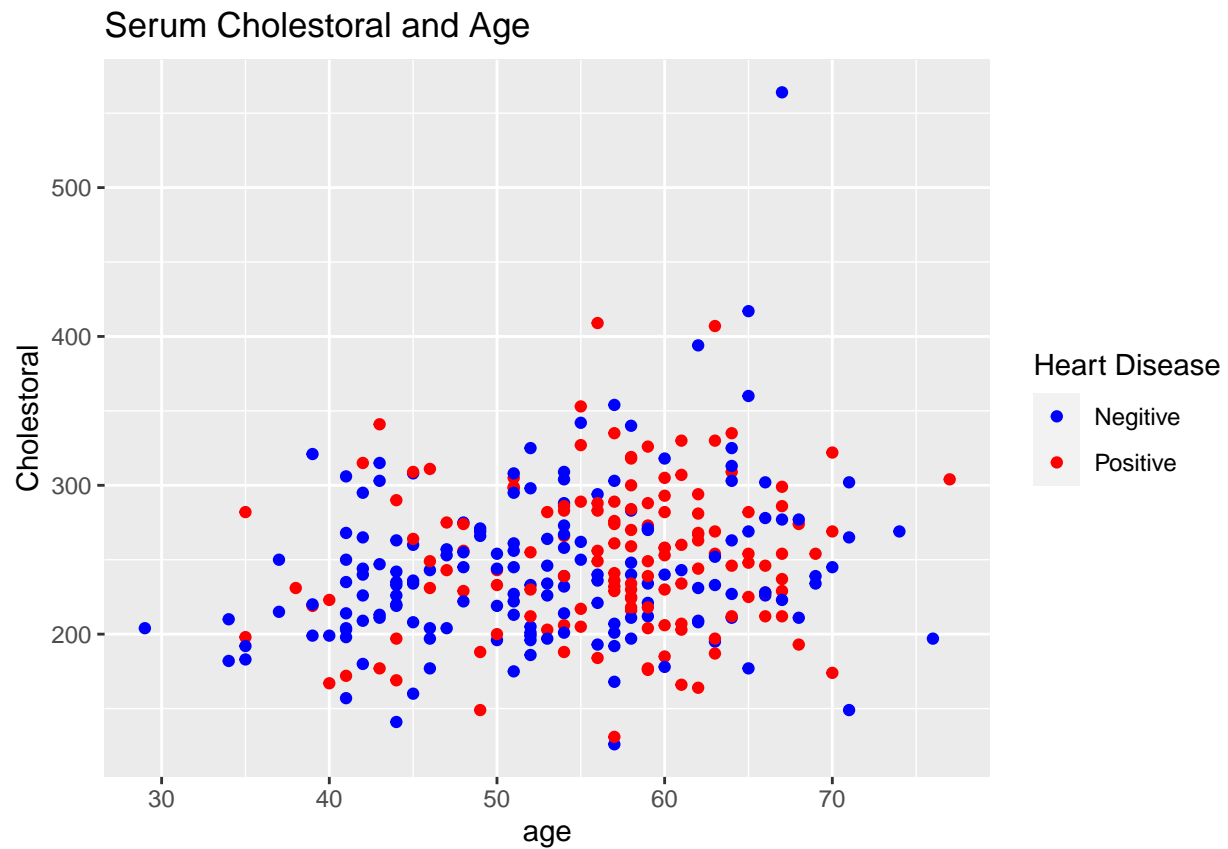
from the box plot, 3 features stand out to me that we will look at in the next few graphs: tresbps, chol, thalach

```
#Scatter plot of resting blood pressure and age
heart %>% ggplot(aes(age,trestbps)) +
  geom_point(aes(color = factor(num))) +
  ggtitle("Resting Blood Presure and Age") +
  ylab("Resting BP") +
  scale_color_manual(name = "Heart Disease",
                     labels = c("Negitive","Positive"),
                     values = c("blue","red"))
```



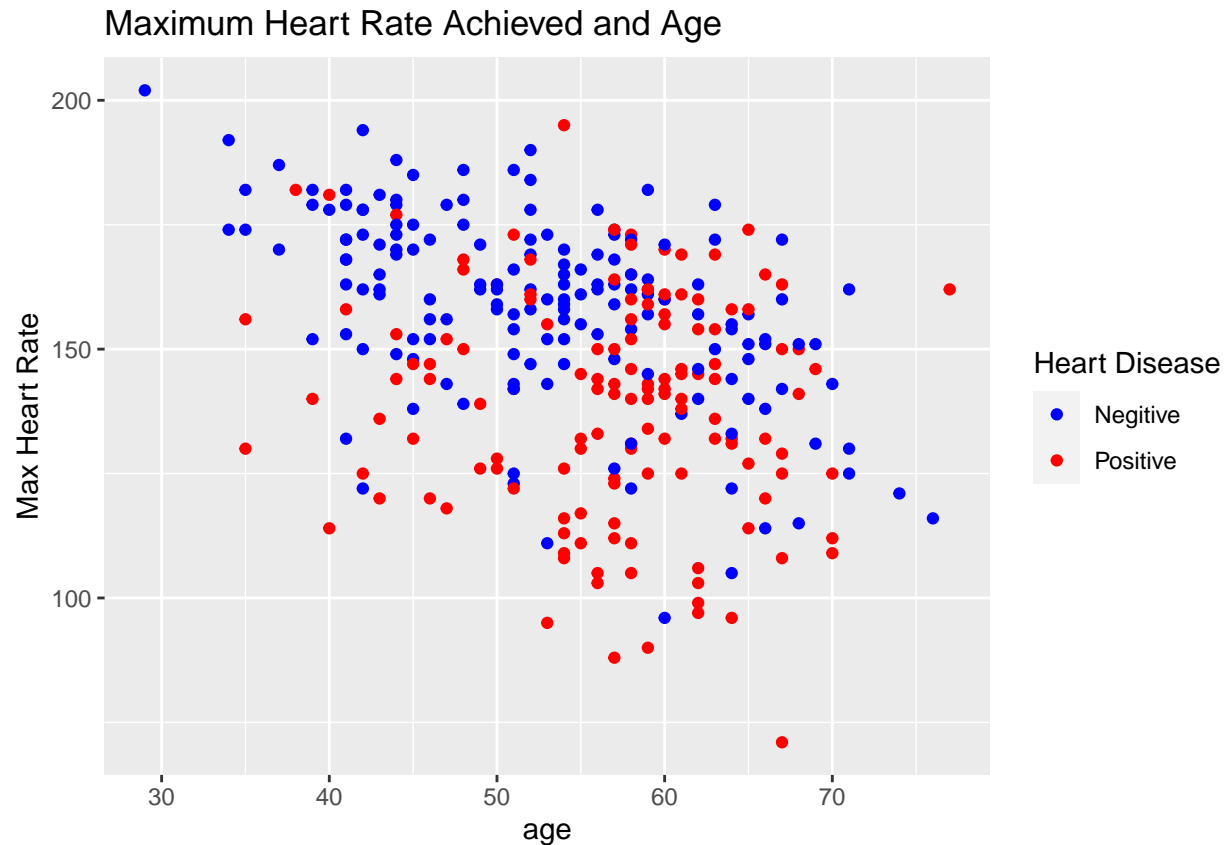
Scatter plot of serum cholestoral and age

```
heart %>% ggplot(aes(age,chol)) +
  geom_point(aes(color = factor(num))) +
  ggtitle("Serum Cholestoral and Age") +
  ylab("Cholestoral") +
  scale_color_manual(name = "Heart Disease",
                     labels = c("Negitive","Positive"),
                     values = c("blue","red"))
```

Scatter plot of maximum heart rate achieved and age

```
heart %>% ggplot(aes(age,thalach)) +  
  geom_point(aes(color = factor(num)))+  
  ggtitle("Maximum Heart Rate Achieved and Age") +  
  ylab("Max Heart Rate") +  
  scale_color_manual(name = "Heart Disease",  
                     labels = c("Negative","Positive"),  
                     values = c("blue","red"))
```



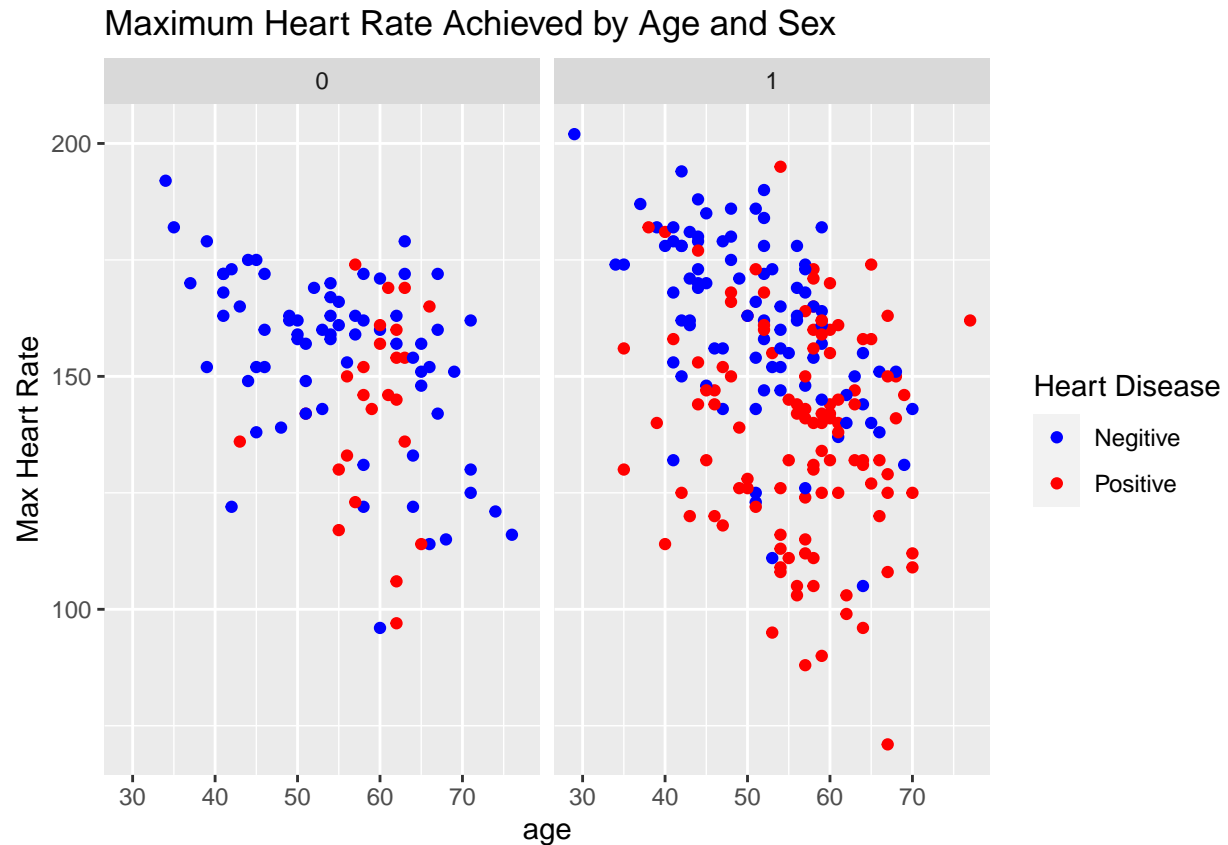
Key insights from the last 3 graphs:

- Max Heart Rate had the largest separation between negative and positive
- Most positive results fall between age 55 - 70
- Very few observations of positive results after the age of 70. This could be due to having a very small sample size. This could also be because people with heart disease do not live past 70 years old very often.

I would like to look more into the max heart rate. This time we will create a scatter plot split by sex.

Scatter plot of maximum heart rate achieved, age and sex 0 = female and 1 = male

```
heart %>% ggplot(aes(age,thalach)) +
  geom_point(aes(color = factor(num))) +
  facet_grid(.~sex) +
  ggtitle("Maximum Heart Rate Achieved by Age and Sex") +
  ylab("Max Heart Rate") +
  scale_color_manual(name = "Heart Disease",
                     labels = c("Negative","Positive"),
                     values = c("blue","red"))
```

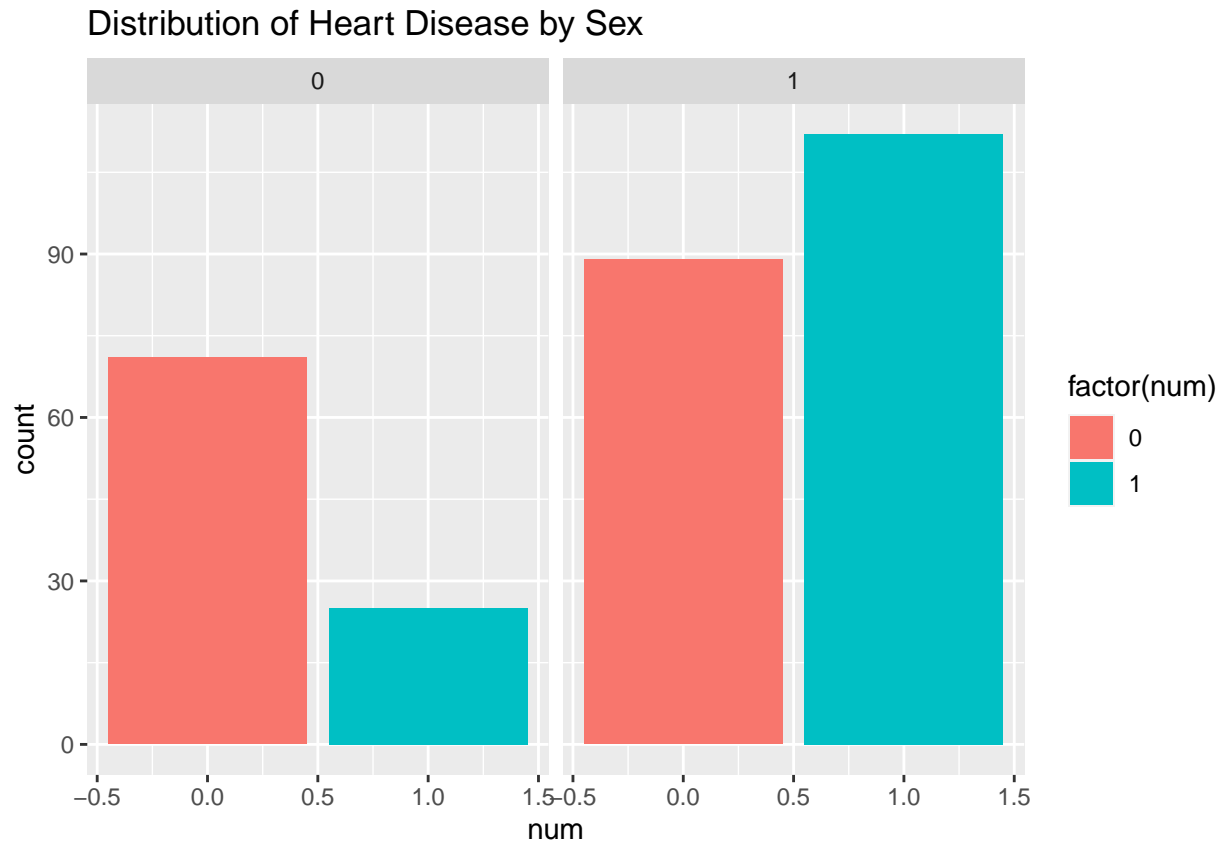


Key insights from this graph:

- Only one female under the age of 50 observed positive
- The majority of male negative observations had max hr over 150
- For the males, the lower the max hr the more positive observations
- For the females, the lower the max hr does not result in more positive observations

bar chart of distribution of heart disease split by sex

```
heart %>% ggplot(aes(num)) +
  geom_bar(aes(fill=factor(num))) +
  facet_grid(.~sex) +
  ggtitle("Distribution of Heart Disease by Sex") +
  scale_color_manual(name = "Heart Disease",
    labels = c("Negative", "Positive"),
    values = c("blue", "red"))
```



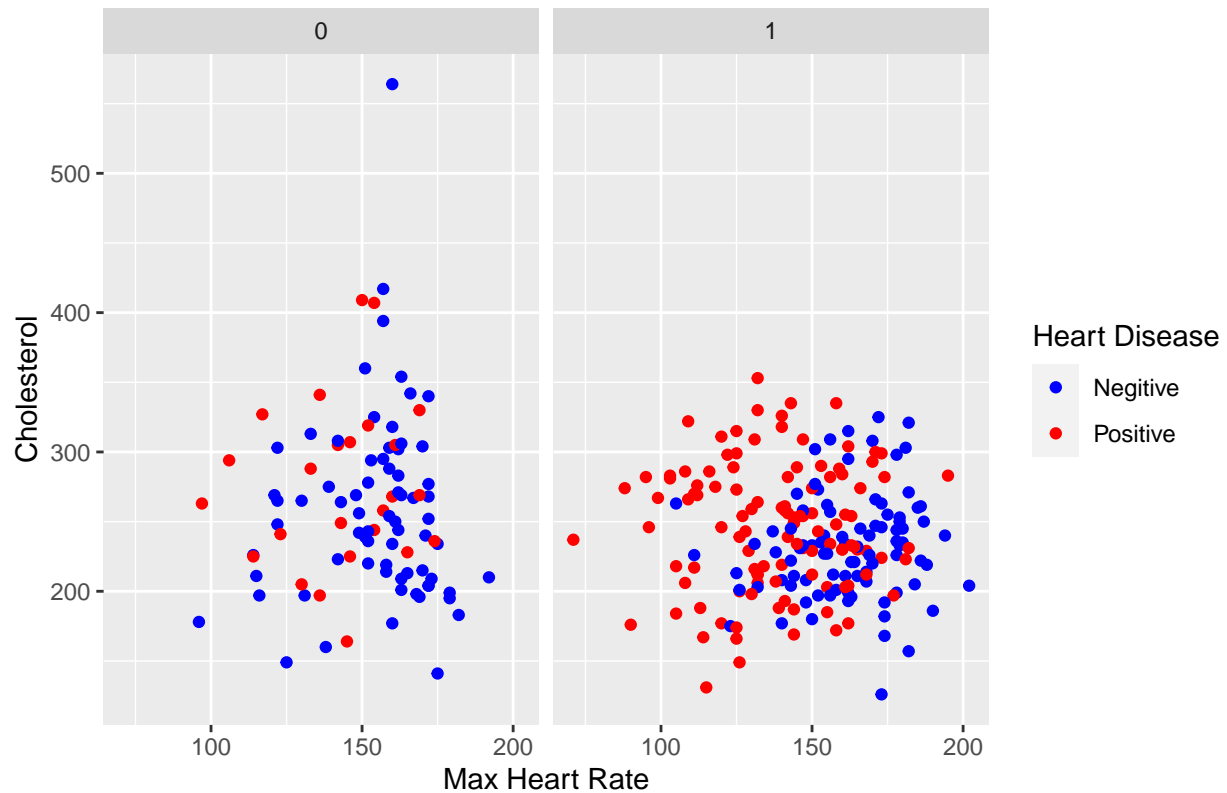
Key insight from the graph:

- The large majority of positive observations in this data set are male

scatter plot of maximum heart rate achieved, cholesterol, sex

```
heart %>% ggplot(aes(thalach, chol)) +
  geom_point(aes(col=factor(num))) +
  facet_grid(.~sex) +
  labs(title = "Maximum Heart Rate Achieved, Cholesterol and Sex",
       x = "Max Heart Rate", y = "Cholesterol") +
  scale_color_manual(name = "Heart Disease",
                    labels = c("Negative", "Positive"),
                    values = c("blue", "red"))
```

Maximum Heart Rate Achieved, Cholesterol and Sex



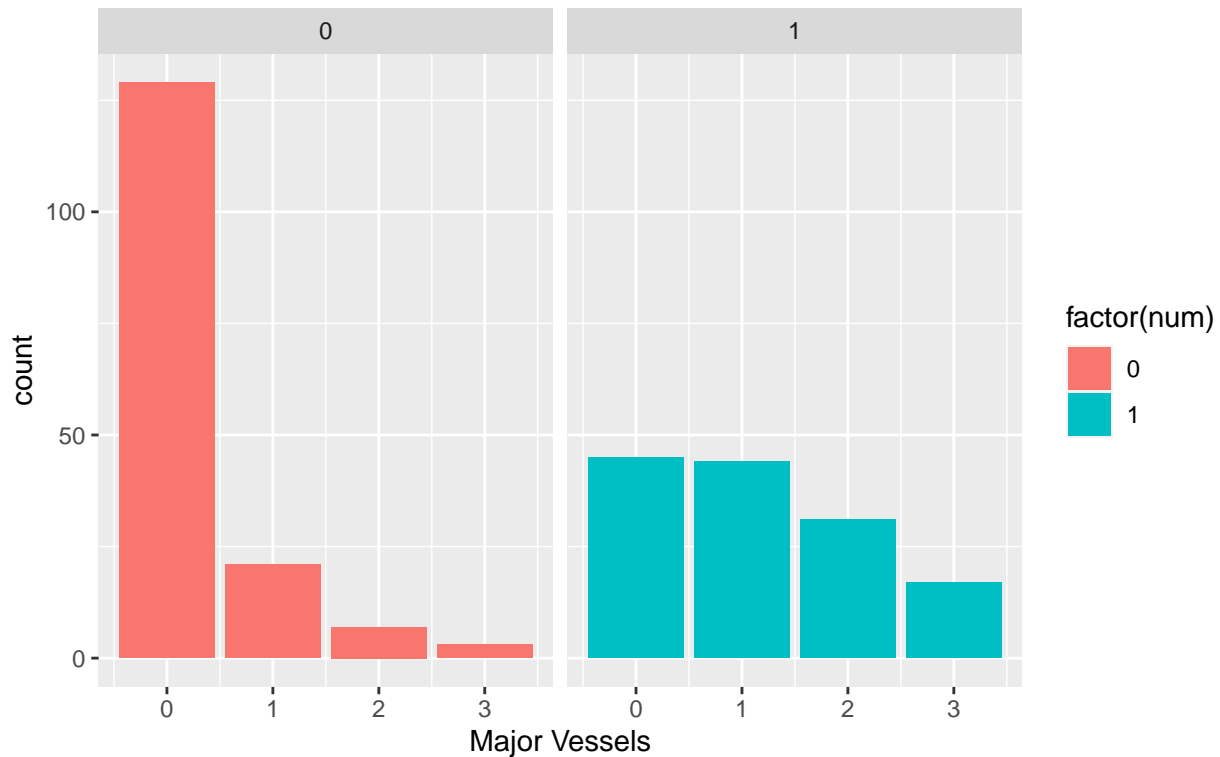
Key insights from this graph:

- For both males and females, high cholesterol and low max hr is more likely to have a positive observation
- For both males and females, low cholesterol and high max hr is more likely to have a negative observation

bar chart of number of major vessels (0-3) colored by flourosopy during exam split by Heart Disease diagnosis

```
heart %>% ggplot(aes(ca)) +
  geom_bar(aes(fill=factor(num))) +
  facet_grid(.~num) +
  ggtitle("Number of Major Vessels (0-3) Colored by Flourosopy During Exam
Split by Heart Disease Diagnosis") +
  xlab("Major Vessels") +
  scale_color_manual(name = "Heart Disease",
                     labels = c("Negitive", "Positive"),
                     values = c("blue", "red"))
```

Number of Major Vessels (0–3) Colored by Flourosopy During Exam
Split by Heart Disease Diagnosis



Key insight from this graph:

- Majority of positive observations have all three major vessels functioning properly

2.5 Models

creating variable x which will consist of the data set except the feature we are trying to predict

```
x <- heart[, -14]
```

creating variable y which will consist the feature we are trying to predict

```
y <- heart$num
```

We will split these two up into training and test

```
set.seed(10, sample.kind = "Rounding")
test_index <- createDataPartition(y, times = 1, p = .2, list = FALSE)
test_x <- x[test_index,]
test_y <- y[test_index]
train_x <- x[-test_index,]
train_y <- y[-test_index]
```

checking to see if the proportions are the same for the train and test set

```
mean(test_y == 0)
```

```
## [1] 0.5666667
```

```
mean(train_y == 0)
```

```
## [1] 0.5316456
```

Will be using k-fold cross validation on all the algorithms creating the k-fold parameters, k is 10

```
control <- trainControl(method = "cv", number = 10, p = .9)
```

Logistic regression model

```
#training the model using train set
train_glm <- train(train_x, as.factor(train_y), method = "glm",
                  family = "binomial",
                  trControl = control)

#creating the predictions
glm_preds <- predict(train_glm, test_x)

#getting the overall accuracy
logistic_regression <- confusionMatrix(glm_preds, as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy results
logistic_regression
```

```
## [1] 0.85
```

Linear discriminant analysis

```
#training the model using the train set
train_lda <- train(train_x, as.factor(train_y), method = "lda",
                  trControl = control)

#creating the predictions
lda_preds <- predict(train_lda, test_x)

#getting the overall accuracy
LDA <- confusionMatrix(lda_preds, as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy results
LDA
```

```
## [1] 0.8666667
```

Quadratic discriminant analysis

```
#training the model using the train set
train_qda <- train(train_x, as.factor(train_y), method = "qda",
                  trControl = control)

#creating the predictions
qda_preds <- predict(train_qda, test_x)

#getting the overall accuracy
QDA <- confusionMatrix(qda_preds, as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy results
QDA
```

```
## [1] 0.8
```

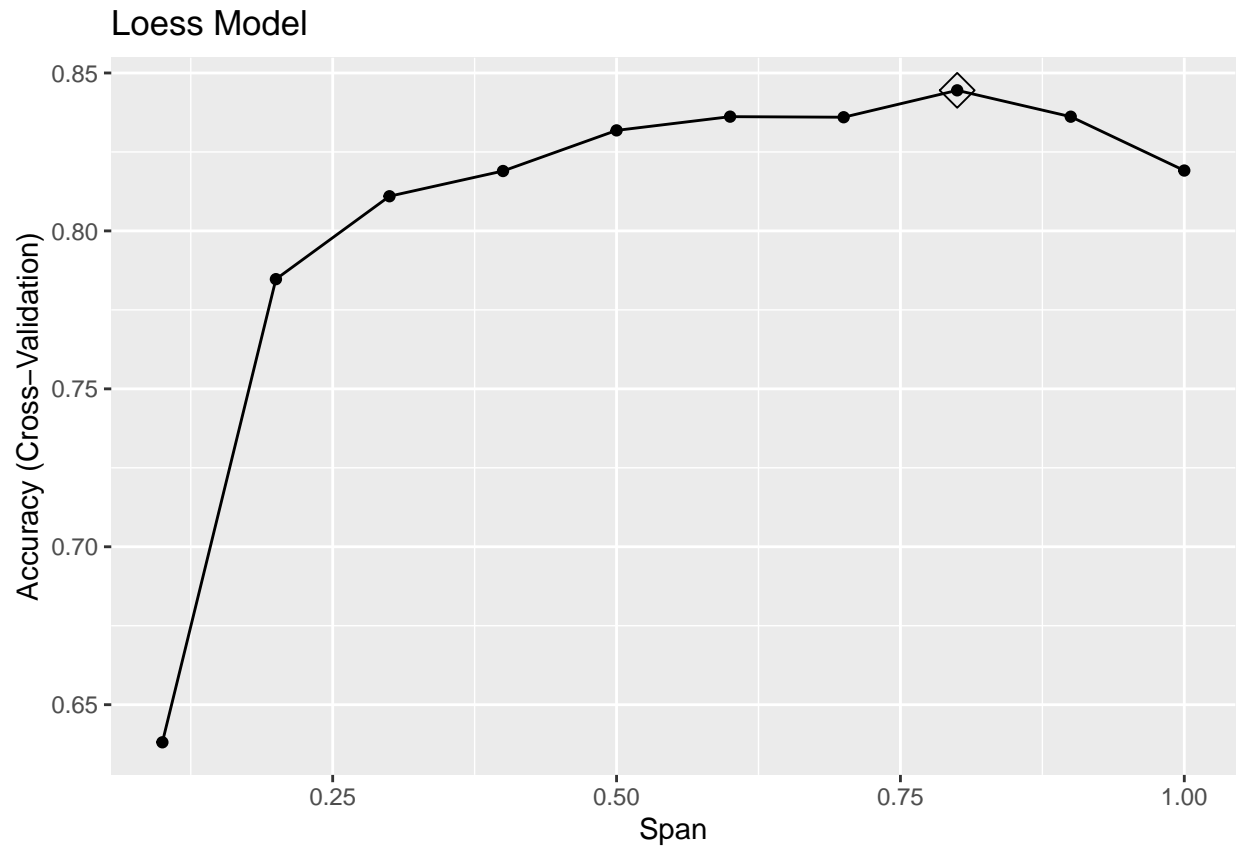
Loess model

```
#creating grid for the two parameter: span, degree
grid <- expand.grid(span = seq(.1, 1, len = 10), degree = 1)

#setting the seed
set.seed(5, sample.kind = "Rounding")

#training the model using the training set
train_loess <- train(train_x, as.factor(train_y), method = "gamLoess",
                   trControl = control,
                   tuneGrid = grid)

#creating graph of the tuning result
ggplot(train_loess, highlight = TRUE) +
  ggtitle("Loess Model")
```

```
#creating the predictions
loess_preds <- predict(train_loess, test_x)

#getting the accuracy
Loess <- confusionMatrix(loess_preds, as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy result
Loess
```

```
## [1] 0.85
```

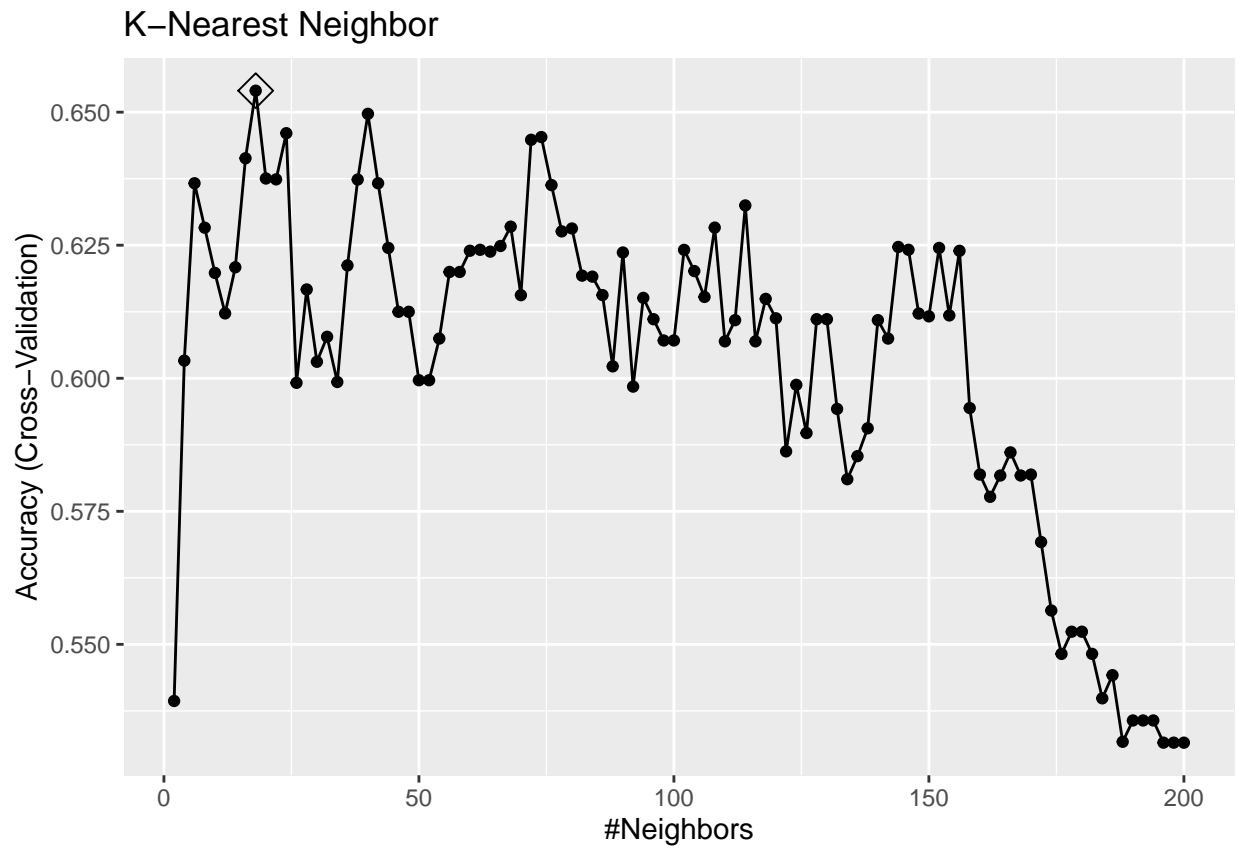
K-nearest neighbors model

```
#setting the seed
set.seed(7, sample.kind = "Rounding")

#creating tuning parameter for K
tuning <- data.frame(k = seq(2, 200, 2))

#training the model with the training set
train_knn <- train(train_x, as.factor(train_y), method = "knn",
                  trControl = control,
                  tuneGrid = tuning,)
```

```
#creating graph of tuning result
ggplot(train_knn, highlight = TRUE) +
  ggtitle("K-Nearest Neighbor")
```



```
#finding best tuning
train_knn$bestTune
```

```
##      k
## 9 18
```

```
#creating prediction
knn_preds <- predict(train_knn, test_x)
```

```
#getting accuracy
Knearest_neighbors <- confusionMatrix(knn_preds, as.factor(test_y))$overall[["Accuracy"]]
```

```
#viewing accuracy result
Knearest_neighbors
```

```
## [1] 0.65
```

Random Forest

```

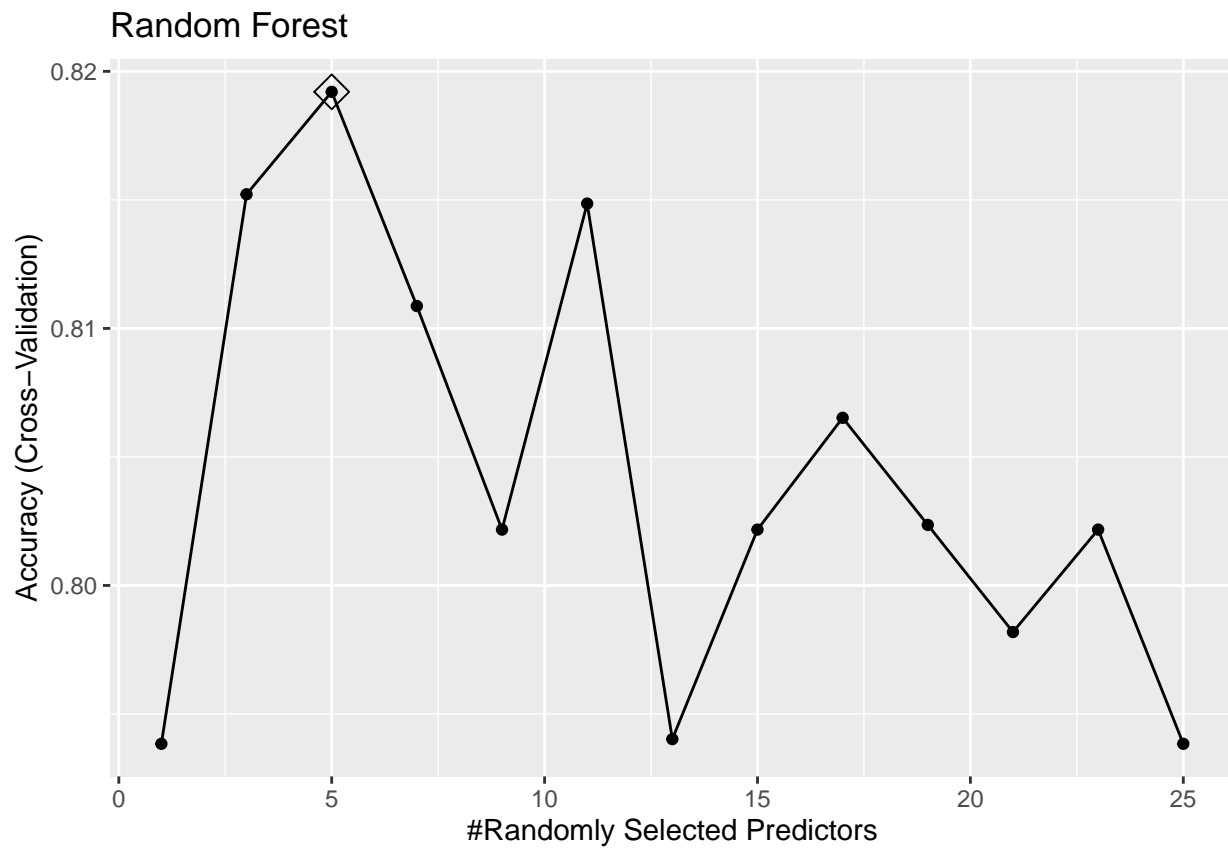
#setting the seed
set.seed(9, sample.kind = "Rounding")

#setting the tuning parameters
tuning <- data.frame(mtry = seq(1,25,2))

#training the model using the train set
train_rf <- train(train_x, as.factor(train_y),method = "rf",
                  tuneGrid = tuning,
                  trControl = control,
                  importance = TRUE)

#creating graph of tuning results
ggplot(train_rf, highlight = TRUE) +
  ggtitle("Random Forest")

```



```

#finding the best tuning result
train_rf$bestTune

```

```

##      mtry
## 3      5

```

```

#creating predictions
rf_preds <- predict(train_rf, test_x)

#getting accuracy
Random_Forest <- confusionMatrix(rf_preds,as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy result
Random_Forest

```

```
## [1] 0.8333333
```

```

#viewing importance
varImp(train_rf)

```

```

## rf variable importance
##
##      Importance
## thal      100.00
## ca        99.83
## cp        84.80
## oldpeak   75.11
## sex       52.99
## thalach   43.50
## exang     38.74
## slope     35.54
## age       34.25
## restecg   33.88
## trestbps  30.07
## fbs       15.11
## chol      0.00

```

Tree Models from Genetic Algorithms

```

#setting the seed
set.seed(7, sample.kind = "Rounding")

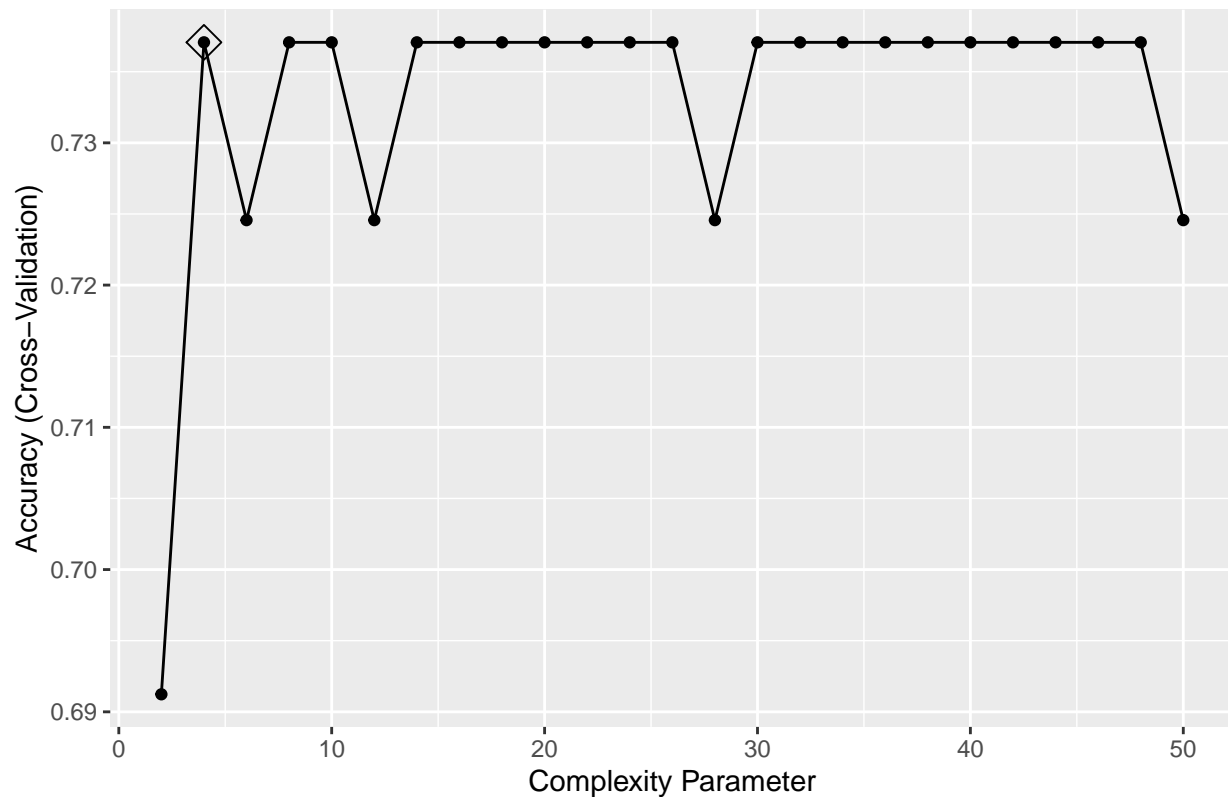
#setting the tuning parameter alpha
tuning <- data.frame(alpha = seq(2,50,2))

#training the model on the train set
train_tree <- train(train_x, as.factor(train_y),method = "evtree",
                    tuneGrid = tuning,
                    trControl = control)

#creating a graph for the tuning results
ggplot(train_tree, highlight = TRUE) +
  ggtitle("Tree Models From Genetic Algorithms")

```

Tree Models From Genetic Algorithms



```
#finding best tune
train_tree$bestTune
```

```
##      alpha
## 2      4
```

```
#creating predictions
tree_preds <- predict(train_tree, test_x)
```

```
#getting accuracy results
tree_model <- confusionMatrix(tree_preds,as.factor(test_y))$overall[["Accuracy"]]
```

```
#viewing accuracy results
tree_model
```

```
## [1] 0.7833333
```

Least Squares Support Vector Machine

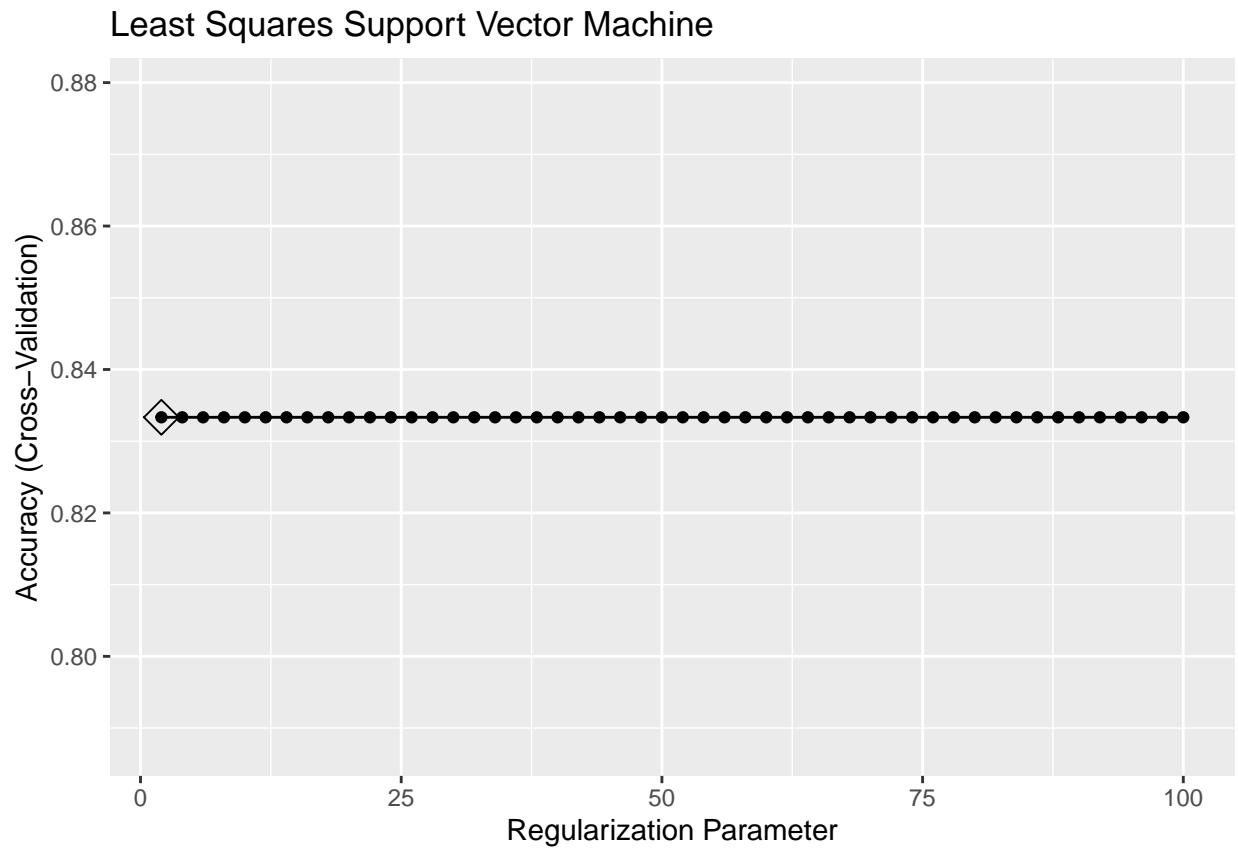
```
#setting the seed
set.seed(7, sample.kind = "Rounding")
```

```
#setting the tuning parameter alpha
```

```
tuning <- data.frame(tau = seq(2,100,2))

#training the model on the train set
train_SVM <- train(train_x, as.factor(train_y),method = "lssvmLinear",
                  tuneGrid = tuning,
                  trControl = control)

#creating a graph for the tuning results
ggplot(train_SVM, highlight = TRUE) +
  ggtitle("Least Squares Support Vector Machine")
```



```
#finding best tune
train_SVM$bestTune
```

```
##    tau
## 1    2
```

```
#creating predictions
SVM_preds <- predict(train_SVM, test_x)

#getting accuracy results
SVM <- confusionMatrix(SVM_preds,as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy results
SVM
```

```
## [1] 0.8666667
```

Bayesian Generalized Linear Model

```
#setting the seed
set.seed(7, sample.kind = "Rounding")

#training the model on the train set
train_nb <- train(train_x, as.factor(train_y), method = "bayesglm",
                  trControl = control)

#creating predictions
nb_preds <- predict(train_nb, test_x)

#getting accuracy results
nb <- confusionMatrix(nb_preds, as.factor(test_y))$overall[["Accuracy"]]

#viewing accuracy results
nb
```

```
## [1] 0.85
```

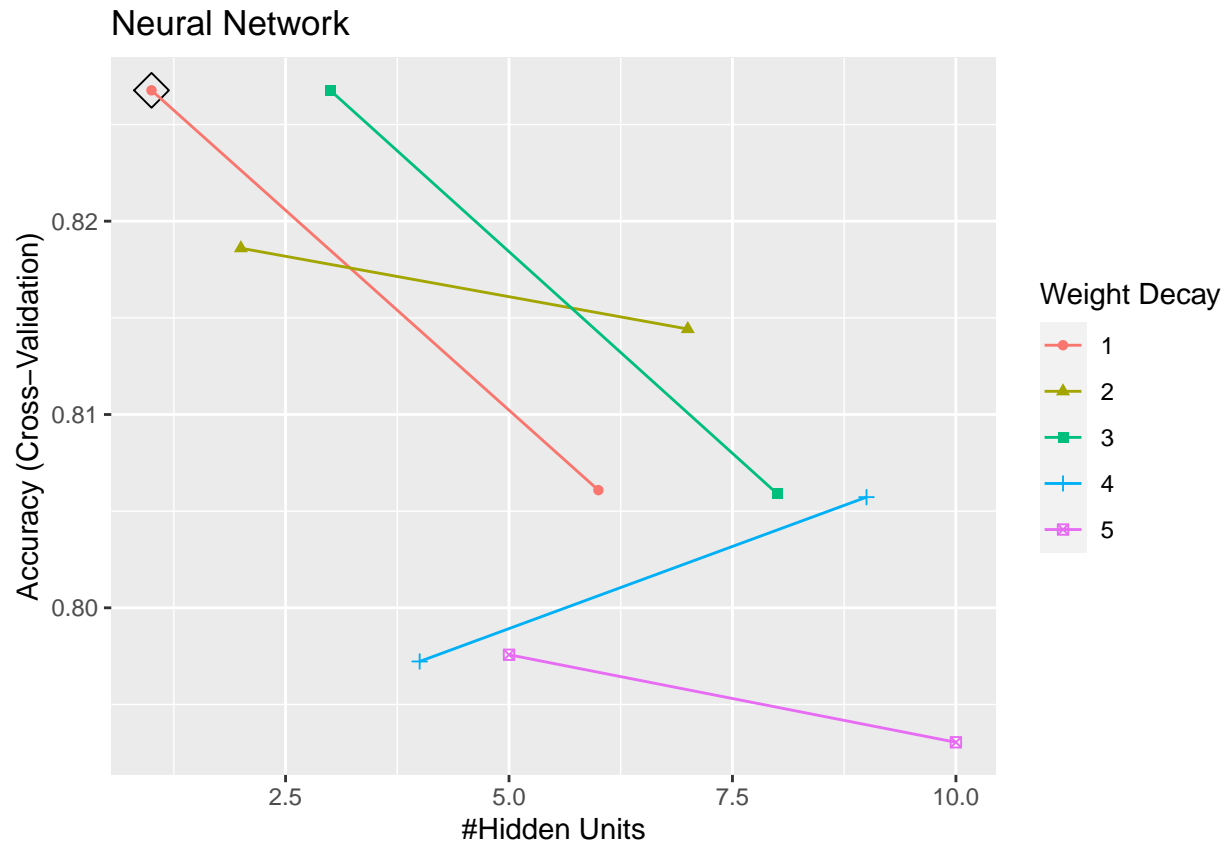
Neural Network

```
#setting the seed
set.seed(7, sample.kind = "Rounding")

#setting the tuning parameter alpha
tuning <- data.frame(size = seq(10), decay = seq(1, 5, 1))

#training the model on the train set
train_nn <- train(train_x, as.factor(train_y), method = "nnet",
                  tuneGrid = tuning,
                  trControl = control)

#creating a graph for the tuning results
ggplot(train_nn, highlight = TRUE) +
  ggtitle("Neural Network")
```



```
#finding best tune
train_nn$bestTune

#creating predictions
nn_preds <- predict(train_nn, test_x)

#getting accuracy results
nn <- confusionMatrix(nn_preds, as.factor(test_y))$overall[["Accuracy"]]
```

```
#viewing accuracy results
nn
```

```
## [1] 0.85
```

3. Results

3.1 Results

creating a results table

```
result_table <- data.frame(Algorithm = c("Logistic Regression",
    "Linear Discriminant Analysis",
    "Quadratic Discriminant Analysis",
```



```

"Loess Model",
"K-Nearest Neighbors",
"Random Forest",
"Tree Models from Genetic Algorithms",
"Least Squares Support Vector Machine",
"Bayesian Generalized Linear Model",
"Neural Network"),
Result = c(round(logistic_regression,2),
round(LDA,2),
round(QDA,2),
round(Loess,2),
round(Knearest_neighbors,2),
round(Random_Forest,2),
round(tree_model,2),
round(SVM,2),
round(nb,2),
round(nn,2)))

```

```

#viewing final results
result_table

```

##	Algorithm	Result
## 1	Logistic Regression	0.85
## 2	Linear Discriminant Analysis	0.87
## 3	Ouadratic Discriminant Analysis	0.80
## 4	Loess Model	0.85
## 5	K-Nearest Neighbors	0.65
## 6	Random Forest	0.83
## 7	Tree Models from Genetic Algorithms	0.78
## 8	Least Squares Support Vector Machine	0.87
## 9	Bayesian Generalized Linear Model	0.85
## 10	Neural Network	0.85

3.2 Brief thoughts about the results

The best overall accuracy came from both Linear Discriminant Analysis and Least Squares Support Vector Machine at 87%. We set out with a goal of training a model that could predict the correct diagnosis with an accuracy of 85%. Not only did we achieve this goal but 6 out of the 10 algorithms were able to predict at least 85%

4. Conclusion

4.1 Summary

We set out to use the UCI data set on Heart Disease to create a model that could correctly predict Heart Disease diagnoses. We set a goal of achieving an overall accuracy of 85%. We started by downloading the UCI data set on Heart Disease. We then cleaned the data set and prepared it for analysis. We split the data set into training and test sets. We trained 10 algorithms using the train set and applying the k-fold cross validation technique with a k of 10. We achieved our goal with 6 algorithms that had an overall accuracy of at least 85% with two algorithms achieving 87%.

4.2 Limitations

For me the biggest limitation in this project is the size of the data set. With only 303 observations this is a very small sample size. The other limitation is the data within the data set. 14 features is enough to achieve a high prediction accuracy, as we proved, but I think with more features we could achieve an overall accuracy over 90%.

4.3 Future Work

For the future, I would be curious to see how these algorithms perform on a much larger data set, say 10 million plus observation data set. Along with a data set that has more features such as: height, weight, if parents had heart disease, use of drugs and alcohol, exercise amount, etc. I would be curious to see which algorithms perform better and if any perform worse. One final thing that I would include is adding more algorithms to this project. These 10 algorithms are not the only algorithms that work well with classification and they may produce a higher overall accuracy, along with an ensemble of these algorithms too.