# HeartAttack Prediction

Ajanthan

01/05/2021

## Heart Attack Analysis and Predictive Analysis

### Introduction

Heart disease is the leading cause of death for men, women, and people of most racial and ethnic groups. One person dies every 36 seconds in the United States from cardiovascular disease. About 655,000 Americans die from heart disease each year—that's 1 in every 4 deaths. The goal of this project is to find out which factors influences the chances of getting a heart attack. I will be using Heart Attack Analysis & Prediction Dataset, which is a dataset for heart attack classification from kaggle. In addition i'll be using the following classification technique Decision Tree, to provide further analysis of which factors influence heart attacks.

**Description of Dataset**

Age : Age of the patient

Sex : Sex of the patient

exng: exercise induced angina (1 = yes; 0 = no)

caa: number of major vessels (0-3)

cp : Chest Pain type chest pain type

Value 1: typical angina Value 2: atypical angina Value 3: non-anginal pain Value 4: asymptomatic

trtbps : resting blood pressure (in mm Hg)

chol : cholestoral in mg/dl fetched via BMI sensor

fbs : (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false)

rest_ecg : 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

oldpeak: ST depression induced by exercise relative to rest

slp: the slope of the peak exercise ST segment (0 = upsloping; 1 = flat; 2 = downsloping)

thall: 1 = normal; 2 = fixed defect; 3 = reversable defect

target : 0= less chance of heart attack 1= more chance of heart attack

Here is a sample of the first 5 rows of the data

```
##   age sex cp trtbps chol fbs restecg thalachh exng oldpeak slp caa thall output
## 1  63   1  3    145  233   1       0      150    0     2.3   0   0     1      1
## 2  37   1  2    130  250   0       1      187    0     3.5   0   0     2      1
## 3  41   0  1    130  204   0       0      172    0     1.4   2   0     2      1
## 4  56   1  1    120  236   0       1      178    0     0.8   2   0     2      1
## 5  57   0  0    120  354   0       1      163    1     0.6   2   0     2      1
```

Here is the data types of the data set

```
## '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 ...
##  $ trtbps  : 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 0 1 0 ...
##  $ restecg : int  0 1 0 1 1 1 0 1 1 1 ...
##  $ thalachh: int  150 187 172 178 163 148 153 173 162 174 ...
##  $ exng    : int  0 0 0 0 1 0 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 ...
##  $ slp     : int  0 0 2 2 2 1 1 2 2 2 ...
##  $ caa     : int  0 0 0 0 0 0 0 0 0 0 ...
##  $ thall   : int  1 2 2 2 2 1 2 3 3 2 ...
##  $ output  : int  1 1 1 1 1 1 1 1 1 1 ...
```

According to our dataset target which is output is our independant variable. We will use logistic regression
since we are dealing with a categorical variable. We will also build two models one model which is a additative
model and another model which includes interactions. We will evaluate both models performance and choose
the best fitting model for further analysis.

## Cleaning The Dataset

Based on the data set we have qualititave dependent variables and a qualitative independent variable so lets
factor the qualitative variables and relable them so they have meaningful names in the dataset to make it
easier reading the data when analyzing the data.

```
##   age sex               cp trtbps chol   fbs              restecg thalachh
## 1  63   M     asymptomatic    145  233  true               normal      150
## 2  37   M non-anginal pain    130  250 false ST-T wave abnormality      187
## 3  41   F  atypical angina    130  204 false               normal      172
## 4  56   M  atypical angina    120  236 false ST-T wave abnormality      178
## 5  57   F   typical angina    120  354 false ST-T wave abnormality      163
##   exng oldpeak         slp caa        thall      output
## 1   No     2.3    upsloping   0       normal More Chance
## 2   No     3.5    upsloping   0 fixed defect More Chance
## 3   No     1.4 downsloping   0 fixed defect More Chance
## 4   No     0.8 downsloping   0 fixed defect More Chance
## 5  Yes     0.6 downsloping   0 fixed defect More Chance
```

Here is a basic summary of the factored data set

```
## 'data.frame':    303 obs. of  14 variables:
##  $ age     : int  63 37 41 56 57 57 56 44 52 57 ...
##  $ sex     : Factor w/ 2 levels "M","F": 1 1 2 1 2 1 2 1 1 1 ...
##  $ cp      : Factor w/ 4 levels "typical angina",..: 4 3 2 2 1 1 2 2 3 3 ...
##  $ trtbps  : int  145 130 130 120 120 140 140 120 172 150 ...
##  $ chol    : int  233 250 204 236 354 192 294 263 199 168 ...
##  $ fbs     : Factor w/ 2 levels "true","false": 1 2 2 2 2 2 2 2 1 2 ...
##  $ restecg : Factor w/ 3 levels "normal","ST-T wave abnormality",..: 1 2 1 2 2 2 1 2 2 2 ...
##  $ thalachh: int  150 187 172 178 163 148 153 173 162 174 ...
##  $ exng    : Factor w/ 2 levels "Yes","No": 2 2 2 2 1 2 2 2 2 2 ...
##  $ oldpeak : num  2.3 3.5 1.4 0.8 0.6 0.4 1.3 0 0.5 1.6 ...
##  $ slp     : Factor w/ 3 levels "upsloping","flat",..: 1 1 3 3 3 2 2 3 3 3 ...
##  $ caa     : int  0 0 0 0 0 0 0 0 0 0 ...
##  $ thall   : Factor w/ 3 levels "normal","fixed defect",..: 1 2 2 2 2 1 2 3 3 2 ...
##  $ output  : Factor w/ 2 levels "Less Chance",..: 2 2 2 2 2 2 2 2 2 2 ...


##       age          sex                      cp         trtbps          chol
##  Min.   :29.00   M:207    typical angina  :143   Min.   : 94.0   Min.   :126.0
##  1st Qu.:47.50   F: 96    atypical angina : 50   1st Qu.:120.0   1st Qu.:211.0
##  Median :55.00            non-anginal pain: 87   Median :130.0   Median :240.0
##  Mean   :54.37            asymptomatic    : 23   Mean   :131.6   Mean   :246.3
##  3rd Qu.:61.00                                   3rd Qu.:140.0   3rd Qu.:274.5
##  Max.   :77.00                                   Max.   :200.0   Max.   :564.0
##     fbs                          restecg       thalachh        exng
##  true : 45    normal                 :147   Min.   : 71.0   Yes: 99
##  false:258    ST-T wave abnormality  :152   1st Qu.:133.5   No :204
##               left ventricular hypertrophy:  4   Median :153.0
##                                               Mean   :149.6
##                                               3rd Qu.:166.0
##                                               Max.   :202.0
##     oldpeak              slp           caa                    thall
##  Min.   :0.00   upsloping : 21   Min.   :0.0000   normal           : 18
##  1st Qu.:0.00   flat      :140   1st Qu.:0.0000   fixed defect     :166
##  Median :0.80   downsloping:142   Median :0.0000   reversable defect:117
##  Mean   :1.04                    Mean   :0.7294   NA's             :  2
##  3rd Qu.:1.60                    3rd Qu.:1.0000
##  Max.   :6.20                    Max.   :4.0000
##        output
##  Less Chance:138
##  More Chance:165
##
##
##
##
```

removing the NA's from the dataset

```
##       age          sex                      cp         trtbps          chol
##  Min.   :29.00   M:206    typical angina  :142   Min.   : 94.0   Min.   :126.0
##  1st Qu.:47.00   F: 95    atypical angina : 50   1st Qu.:120.0   1st Qu.:211.0
##  Median :56.00            non-anginal pain: 86   Median :130.0   Median :241.0
##  Mean   :54.38            asymptomatic    : 23   Mean   :131.6   Mean   :246.5
##  3rd Qu.:61.00                                   3rd Qu.:140.0   3rd Qu.:275.0
```

```
## Max.    :77.00                                       Max.    :200.0   Max.    :564.0
##    fbs                              restecg         thalachh       exng
## true : 44    normal                     :146   Min.    : 71.0   Yes: 98
## false:257    ST-T wave abnormality      :151   1st Qu.:134.0   No :203
##              left ventricular hypertrophy:  4   Median :153.0
##                                                 Mean    :149.7
##                                                 3rd Qu.:166.0
##                                                 Max.    :202.0
##    oldpeak            slp            caa                      thall
## Min.    :0.000   upsloping  : 21   Min.    :0.0000   normal            : 18
## 1st Qu.:0.000   flat       :139   1st Qu.:0.0000   fixed defect      :166
## Median :0.800   downsloping:141   Median :0.0000   reversable defect:117
## Mean    :1.043                     Mean    :0.7342
## 3rd Qu.:1.600                     3rd Qu.:1.0000
## Max.    :6.200                     Max.    :4.0000
##        output
## Less Chance:137
## More Chance:164
##
##
##
##
```
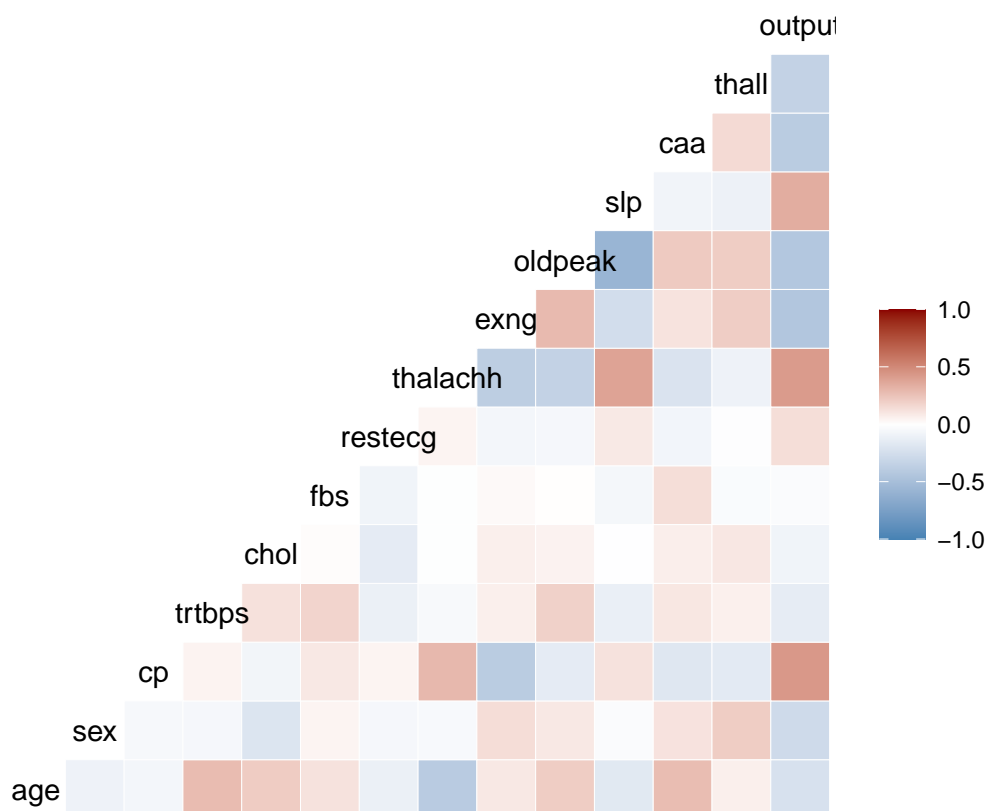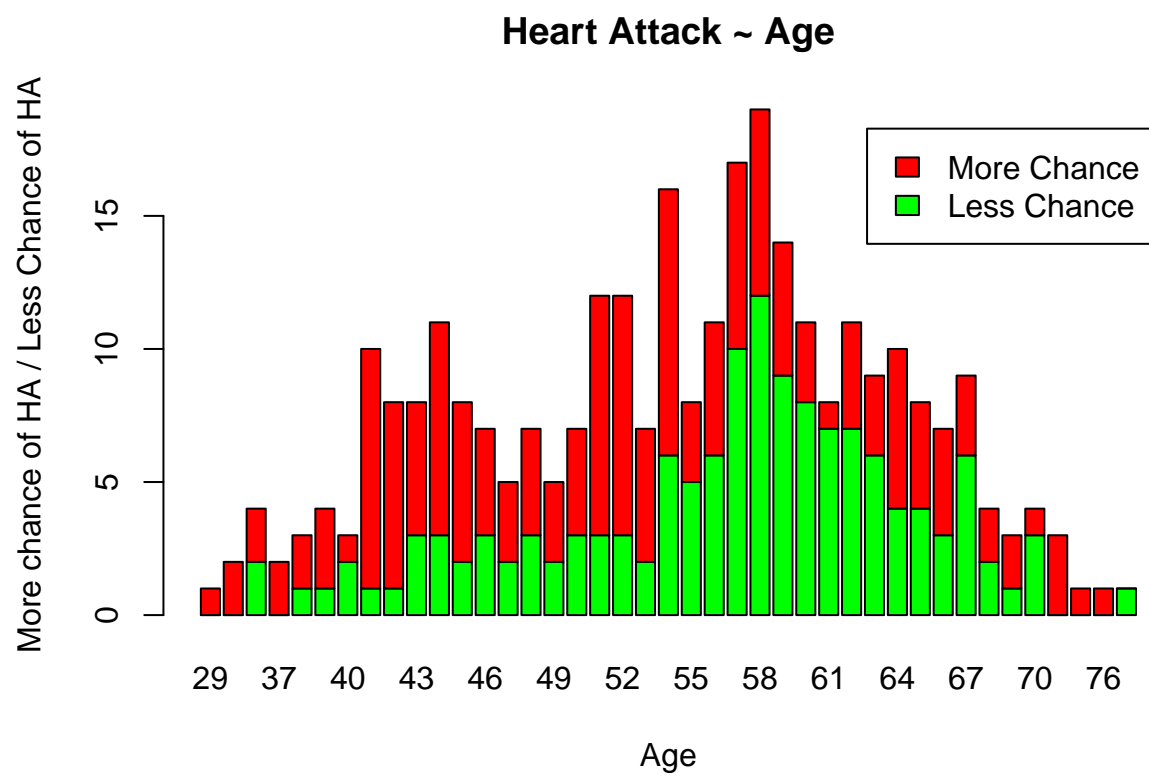
# Data Preperation

## Correlation

Lets see if there is any form of multicollinearity in the dataset by looking at the correlation plot of the data
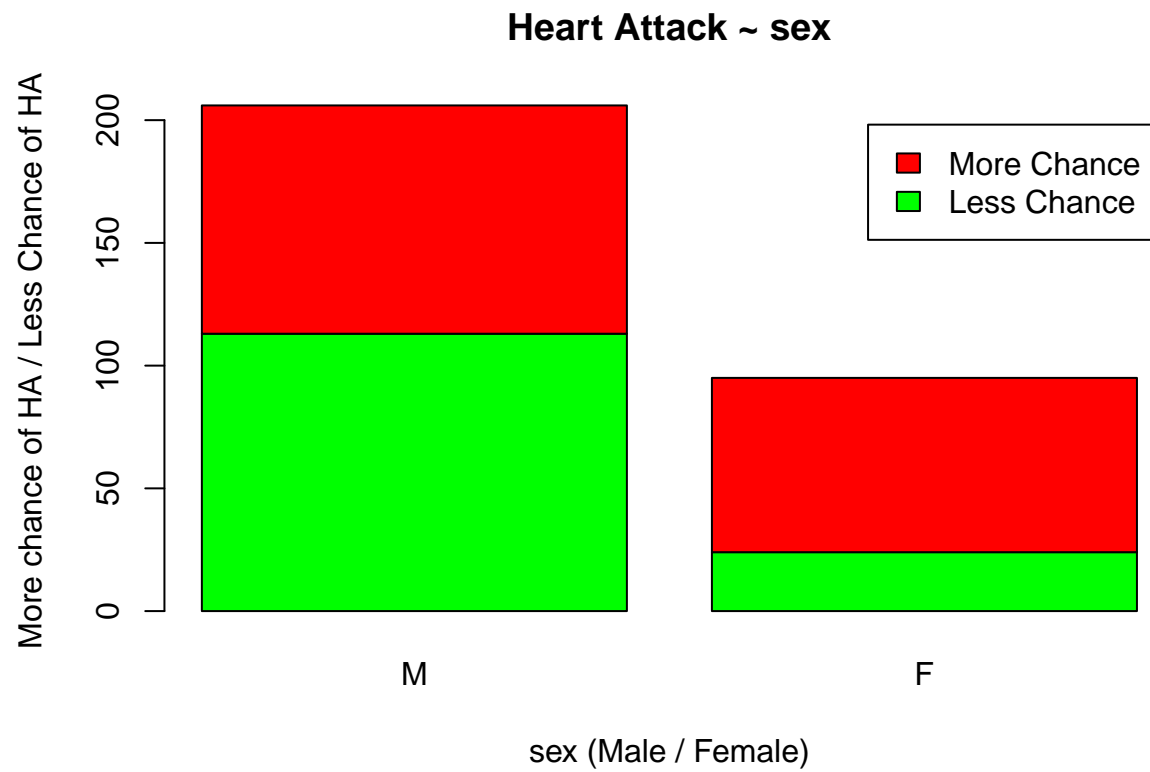
By looking at the following figure we can see that there is high correlation between the predictors cp, thalach and slp to tell whether or not a patient has a higher chance of a heart attack or less chance of a heart attack. This makes sense because cp is the type of chest pain, thalach the maximum heart rate achieved and slp the slope of the heart rate. These are major factors that help distinguish whether a patient is more likely or less likely to have a hearth attack. We can also see that there is some aspects of multicollinearity in the model for there shows high significants amongst other predictors such as age and trtbps, cp and thalachh, thalach and slp, etc. These factors must be taken into account when building the model.

## Data visualization

Lets use bar graphs to visualize the effects of the predictors on the independent variable output.

**Heart Attack ~ Age**

As the age increases we can see the chances of having a heart attack increases.

## Heart Attack ~ sex



We can see most of the patients in the data set is male compared to females. The ratio of having a heart attack for males for having a higher chance of having a heart attack to a lower chance if having a heart attack is less than the ratio for females.

**Heart Attack ~ exng**

Exercise included agnia

We can see that for patients who experience chest pain during exercise has a less chance of having a heart attack for patients who experienced chest pain when they were not exercising.

**Heart Attack ~ Chest Pain**

We can see patients who experienced non-anginal pain in there chest have more chance of a heart attack than patients who experienced typical chest pains. Therefore non-anginal pain in the chest is a significant factor to having a higher chance of heart attack.

## Heart Attack ~ Fasting blood sugar



We can see that if the fasting blood sugar is less than 120 mg/dl, there is more of a chance of having a heart attack compared to someone with fasting blood sugar greater than 120 mg/dl.

**Heart Attack ~ resting electrocardiographic**

We can see patients who's resting electrocardiographic is showing forms of ST-T wave abnormality have a higher chance of having a heart attack compared to patients who show left venticular hypertrophy. Having a normal reading doesn't give enough information of whether the patient has a higher chance of having a heart to less of chance because the ratio of more chance to less chance of having a heart attack is one to one.

# Heart Attack ~ slope



We can see if the slope of the cardiactric machine is downwards slopping there it shows that there is more of a chance to have a heartattack compared to seeing a upward slope or a flat slope on the machine.

## Heart Attack ~ defects



We can see that if it is a fixed defect there is more chance of a heart attack to occur compared to a reversable defect or a normal defect.

##Model Selection ### Additative Model Based on the data set lets first evaluate the additative model

```
##
## Call:
## glm(formula = output ~ age + sex + cp + trtbps + chol + fbs +
##     restecg + thalachh + exng + oldpeak + slp + caa + thall,
##     family = "binomial", data = newheart)
##
## Deviance Residuals:
##     Min      1Q   Median      3Q      Max
## -2.7668  -0.3527   0.1548   0.5312   2.5923
##
## Coefficients:
##                              Estimate Std. Error z value Pr(>|z|)
## (Intercept)                  7.071e-01  2.695e+00   0.262  0.79308
## age                         -2.299e-05  2.353e-02  -0.001  0.99922
## sexF                         1.477e+00  5.231e-01   2.823  0.00475 **
## cpatypical angina            9.769e-01  5.636e-01   1.733  0.08301 .
## cpnon-anginal pain           1.909e+00  4.794e-01   3.982 6.83e-05 ***
## cpasymptomatic               1.985e+00  6.508e-01   3.050  0.00229 **
## trtbps                      -1.712e-02  1.068e-02  -1.603  0.10890
## chol                        -4.303e-03  3.877e-03  -1.110  0.26705
## fbsfalse                    -2.134e-01  5.692e-01  -0.375  0.70774
## restecgST-T wave abnormality 5.885e-01  3.755e-01   1.567  0.11704
```

```
## restecgleft ventricular hypertrophy -2.395e-01  2.246e+00  -0.107  0.91507
## thalachh                              1.782e-02  1.079e-02   1.652  0.09857 .
## exngNo                                7.468e-01  4.260e-01   1.753  0.07958 .
## oldpeak                              -4.861e-01  2.254e-01  -2.157  0.03103 *
## slpflat                              -7.007e-01  8.622e-01  -0.813  0.41637
## slpdownsloping                        1.875e-01  9.372e-01   0.200  0.84144
## caa                                  -8.302e-01  2.036e-01  -4.077 4.55e-05 ***
## thallfixed defect                     6.695e-02  7.714e-01   0.087  0.93084
## thallreversable defect               -1.326e+00  7.572e-01  -1.751  0.08003 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 414.85  on 300  degrees of freedom
## Residual deviance: 201.19  on 282  degrees of freedom
## AIC: 239.19
##
## Number of Fisher Scoring iterations: 6


## [1] "R-squared: 0.515021967180605"
```

## Wald Test Evaluation

This is the Wald test for additative model. A wald test is conducted in logistic regression is to check whether a predictor is significant or not. From looking at the wald test we see that some of the p-values are high which indicates that there are some varaibles that are not significant to the model. Doing a quick glance we see that age, fbs, and restec at the left verticular hypertrophy has high p-values showing that they are insignificant to the model. This doesnt mean we should drop them right away, we shall do further analysis on the model to get a clearer idea on whether or not these variables are significant to showing the patient has a lower or higher chance of getting a heart attack.

```
##   numDF denDF F.value p.value
##      19   282 4.34943 <.00001
##                                          Estimate  Std.Error DF   t-value   p-value
## (Intercept)                              0.707054  2.695471  282   0.262312 0.79327
## age                                     -0.000023  0.023533  282  -0.000977 0.99922
## sexF                                     1.476875  0.523100  282   2.823315 0.00509
## cpatypical angina                        0.976943  0.563577  282   1.733468 0.08411
## cpnon-anginal pain                       1.908910  0.479388  282   3.981973 0.00009
## cpasymptomatic                           1.985185  0.650795  282   3.050401 0.00250
## trtbps                                  -0.017119  0.010678  282  -1.603164 0.11002
## chol                                    -0.004303  0.003877  282  -1.109892 0.26799
## fbsfalse                                -0.213385  0.569199  282  -0.374886 0.70803
## restecgST-T wave abnormality             0.588548  0.375510  282   1.567329 0.11816
## restecgleft ventricular hypertrophy     -0.239541  2.246060  282  -0.106649 0.91514
## thalachh                                 0.017820  0.010788  282   1.651829 0.09968
## exngNo                                   0.746836  0.426008  282   1.753105 0.08067
## oldpeak                                 -0.486090  0.225392  282  -2.156649 0.03188
## slpflat                                 -0.700728  0.862186  282  -0.812734 0.41706
## slpdownsloping                           0.187488  0.937188  282   0.200054 0.84158
## caa                                     -0.830169  0.203604  282  -4.077363 0.00006
## thallfixed defect                        0.066946  0.771424  282   0.086783 0.93091
```

```
## thallreversable defect                   -1.325557 0.757231   282 -1.750530 0.08111
##                                           Lower 0.95 Upper 0.95
## (Intercept)                               -4.598742   6.012851
## age                                       -0.046346   0.046300
## sexF                                        0.447200   2.506551
## cpatypical angina                         -0.132409   2.086296
## cpnon-anginal pain                          0.965277   2.852543
## cpasymptomatic                              0.704153   3.266217
## trtbps                                    -0.038138   0.003900
## chol                                      -0.011934   0.003328
## fbsfalse                                  -1.333803   0.907033
## restecgST-T wave abnormality              -0.150611   1.327708
## restecgleft ventricular hypertrophy -4.660711   4.181630
## thalachh                                  -0.003415   0.039055
## exngNo                                    -0.091722   1.585394
## oldpeak                                   -0.929754  -0.042427
## slpflat                                   -2.397864   0.996409
## slpdownsloping                            -1.657284   2.032260
## caa                                       -1.230946  -0.429392
## thallfixed defect                         -1.451534   1.585426
## thallreversable defect                    -2.816100   0.164987
```

## Stepwise Regression

Stepwise regression is a modification of the forward selection so that after each step in which a variable was added, all candidate variables in the model are checked to see if their significance has been reduced below the specified tolerance level. If a nonsignificant variable is found, it is removed from the model. Applying stepwise regression both backward and forward leaves us with output ~ sex + cp + trtbps + chol + thalachh + exng + oldpeak + slp + caa + thall

```
## Start:  AIC=239.19
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall
##
##             Df Deviance    AIC
## - age        1    201.19 237.19
## - fbs        1    201.33 237.33
## - restecg    2    203.75 237.75
## - chol       1    202.41 238.41
## <none>            201.19 239.19
## - slp        2    205.37 239.37
## - trtbps     1    203.81 239.81
## - thalachh   1    204.02 240.02
## - exng       1    204.23 240.23
## - oldpeak    1    206.14 242.14
## - sex        1    209.80 245.80
## - thall      2    214.12 248.12
## - cp         3    222.08 254.08
## - caa        1    219.17 255.17
##
## Step:  AIC=237.19
## output ~ sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall
```

```
##
##             Df Deviance    AIC
## - fbs        1   201.34 235.34
## - restecg    2   203.77 235.77
## - chol       1   202.45 236.45
## <none>           201.19 237.19
## - slp        2   205.37 237.37
## - trtbps     1   203.94 237.94
## - exng       1   204.24 238.24
## - thalachh   1   204.49 238.49
## + age        1   201.19 239.19
## - oldpeak    1   206.14 240.14
## - sex        1   209.94 243.94
## - thall      2   214.16 246.16
## - cp         3   222.08 252.08
## - caa        1   219.92 253.92
##
## Step:  AIC=235.34
## output ~ sex + cp + trtbps + chol + restecg + thalachh + exng +
##     oldpeak + slp + caa + thall
##
##             Df Deviance    AIC
## - restecg    2   203.88 233.88
## - chol       1   202.55 234.55
## <none>           201.34 235.34
## - slp        2   205.46 235.46
## - trtbps     1   203.97 235.97
## - exng       1   204.30 236.30
## - thalachh   1   204.64 236.64
## + fbs        1   201.19 237.19
## + age        1   201.33 237.33
## - oldpeak    1   206.46 238.46
## - sex        1   209.99 241.99
## - thall      2   214.57 244.57
## - cp         3   223.57 251.57
## - caa        1   219.99 251.99
##
## Step:  AIC=233.88
## output ~ sex + cp + trtbps + chol + thalachh + exng + oldpeak +
##     slp + caa + thall
##
##             Df Deviance    AIC
## <none>           203.88 233.88
## - chol       1   205.93 233.93
## - slp        2   208.63 234.63
## - exng       1   206.75 234.75
## - trtbps     1   206.96 234.96
## + restecg    2   201.34 235.34
## - thalachh   1   207.44 235.44
## + fbs        1   203.77 235.77
## + age        1   203.86 235.86
## - oldpeak    1   209.12 237.12
## - sex        1   213.78 241.78
## - thall      2   216.18 242.18
```

```
## - cp         3    226.06 250.06
## - caa        1    222.65 250.65


##
## Call:  glm(formula = output ~ sex + cp + trtbps + chol + thalachh +
##     exng + oldpeak + slp + caa + thall, family = "binomial",
##     data = newheart)
##
## Coefficients:
##          (Intercept)                    sexF        cpatypical angina
##             1.025922                1.535349                 1.000189
##     cpnon-anginal pain         cpasymptomatic                   trtbps
##             1.943932                1.965324                -0.017488
##                 chol               thalachh                   exngNo
##            -0.005343                0.018627                 0.722047
##              oldpeak                slpflat            slpdownsloping
##            -0.484829               -0.711536                 0.228519
##                  caa       thallfixed defect  thallreversable defect
##            -0.817229               -0.032660                -1.345274
##
## Degrees of Freedom: 300 Total (i.e. Null);   286 Residual
## Null Deviance:         414.8
## Residual Deviance: 203.9      AIC: 233.9
```

## Model with interactions

```
##
## Call:
## glm(formula = output ~ age + sex + cp + trtbps + chol + fbs +
##     restecg + thalachh + exng + oldpeak + slp + caa + thall +
##     age * sex + age * cp + trtbps * age + chol * age + fbs *
##     age + exng * age + sex * chol + chol * cp + chol * fbs +
##     chol * exng, family = "binomial", data = newheart)
##
## Deviance Residuals:
##     Min       1Q   Median       3Q       Max
## -2.8295  -0.3641   0.1254   0.4979   2.6450
##
## Coefficients:
##                                      Estimate Std. Error z value Pr(>|z|)
## (Intercept)                          7.129e+00  1.417e+01    0.503 0.614769
## age                                 -3.171e-02  2.490e-01   -0.127 0.898648
## sexF                                 2.733e-01  4.001e+00    0.068 0.945546
## cpatypical angina                    6.375e+00  4.942e+00    1.290 0.197078
## cpnon-anginal pain                   4.903e+00  4.084e+00    1.201 0.229902
## cpasymptomatic                       6.266e+00  6.493e+00    0.965 0.334525
## trtbps                               1.619e-02  7.928e-02    0.204 0.838210
## chol                                -2.187e-02  3.155e-02   -0.693 0.488213
## fbsfalse                            -9.888e+00  6.110e+00   -1.618 0.105587
## restecgST-T wave abnormality         6.448e-01  4.029e-01    1.600 0.109506
## restecgleft ventricular hypertrophy  5.637e-01  1.868e+00    0.302 0.762825
## thalachh                             1.676e-02  1.117e-02    1.501 0.133293
## exngNo                              -3.486e+00  3.680e+00   -0.947 0.343556
```

```
## oldpeak                             -5.458e-01  2.451e-01  -2.227 0.025962 *
## slpflat                             -3.739e-01  9.492e-01  -0.394 0.693661
## slpdownsloping                       4.920e-01  1.055e+00   0.467 0.640817
## caa                                 -8.771e-01  2.349e-01  -3.734 0.000189 ***
## thallfixed defect                    3.013e-01  8.280e-01   0.364 0.715938
## thallreversable defect              -1.107e+00  7.948e-01  -1.393 0.163617
## age:sexF                            -1.492e-02  5.244e-02  -0.285 0.775986
## age:cpatypical angina               -2.248e-02  6.621e-02  -0.339 0.734235
## age:cpnon-anginal pain              -6.404e-02  6.200e-02  -1.033 0.301641
## age:cpasymptomatic                   2.177e-02  7.251e-02   0.300 0.764009
## age:trtbps                          -5.562e-04  1.391e-03  -0.400 0.689337
## age:chol                            -4.015e-05  5.378e-04  -0.075 0.940496
## age:fbsfalse                         1.082e-01  8.910e-02   1.215 0.224524
## age:exngNo                           4.046e-02  5.202e-02   0.778 0.436661
## sexF:chol                            8.117e-03  1.078e-02   0.753 0.451553
## cpatypical angina:chol              -1.602e-02  1.571e-02  -1.020 0.307925
## cpnon-anginal pain:chol              2.073e-03  1.190e-02   0.174 0.861672
## cpasymptomatic:chol                 -2.238e-02  2.176e-02  -1.028 0.303742
## chol:fbsfalse                        1.280e-02  1.241e-02   1.032 0.302217
## chol:exngNo                          8.572e-03  1.114e-02   0.770 0.441566
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 414.85  on 300  degrees of freedom
## Residual deviance: 191.63  on 268  degrees of freedom
## AIC: 257.63
##
## Number of Fisher Scoring iterations: 6


## [1] "R-squared: 0.538066697274319"
```

**Wald test for interactive model**

```
##   numDF denDF F.value p.value
##     33   268 2.39469   7e-05
##                                    Estimate  Std.Error DF   t-value   p-value
## (Intercept)                        7.129284  14.165714 268  0.503277  0.61518
## age                               -0.031709   0.248955 268 -0.127370  0.89874
## sexF                               0.273257   4.000752 268  0.068301  0.94560
## cpatypical angina                  6.375048   4.942201 268  1.289921  0.19819
## cpnon-anginal pain                 4.903159   4.083889 268  1.200610  0.23096
## cpasymptomatic                     6.265785   6.492777 268  0.965039  0.33540
## trtbps                             0.016188   0.079281 268  0.204184  0.83837
## chol                              -0.021866   0.031545 268 -0.693154  0.48881
## fbsfalse                          -9.887685   6.109733 268 -1.618350  0.10676
## restecgST-T wave abnormality       0.644823   0.402909 268  1.600418  0.11068
## restecgleft ventricular hypertrophy 0.563728  1.868053 268  0.301773  0.76306
## thalachh                           0.016764   0.011166 268  1.501241  0.13447
## exngNo                            -3.485743   3.680194 268 -0.947163  0.34441
## oldpeak                           -0.545767   0.245092 268 -2.226781  0.02679
## slpflat                           -0.373894   0.949231 268 -0.393891  0.69397
```

```
## slpdownsloping                          0.492027  1.054593 268   0.466557 0.64120
## caa                                     -0.877077  0.234917 268  -3.733559 0.00023
## thallfixed defect                        0.301297  0.827982 268   0.363893 0.71622
## thallreversable defect                  -1.107233  0.794850 268  -1.393009 0.16477
## age:sexF                                -0.014921  0.052437 268  -0.284554 0.77621
## age:cpatypical angina                   -0.022477  0.066208 268  -0.339498 0.73450
## age:cpnon-anginal pain                  -0.064039  0.061998 268  -1.032920 0.30257
## age:cpasymptomatic                       0.021770  0.072515 268   0.300220 0.76424
## age:trtbps                              -0.000556  0.001391 268  -0.399755 0.68966
## age:chol                                -0.000040  0.000538 268  -0.074647 0.94055
## age:fbsfalse                             0.108222  0.089102 268   1.214587 0.22559
## age:exngNo                               0.040463  0.052019 268   0.777843 0.43735
## sexF:chol                                0.008117  0.010782 268   0.752829 0.45221
## cpatypical angina:chol                  -0.016020  0.015713 268  -1.019587 0.30884
## cpnon-anginal pain:chol                  0.002073  0.011897 268   0.174246 0.86180
## cpasymptomatic:chol                     -0.022379  0.021760 268  -1.028442 0.30467
## chol:fbsfalse                            0.012805  0.012412 268   1.031691 0.30315
## chol:exngNo                              0.008572  0.011140 268   0.769552 0.44224
##                                         Lower 0.95 Upper 0.95
## (Intercept)                            -20.760956 35.019525
## age                                     -0.521866  0.458447
## sexF                                    -7.603644  8.150158
## cpatypical angina                       -3.355430 16.105527
## cpnon-anginal pain                      -3.137427 12.943745
## cpasymptomatic                          -6.517553 19.049123
## trtbps                                  -0.139904  0.172280
## chol                                    -0.083974  0.040242
## fbsfalse                               -21.916863  2.141494
## restecgST-T wave abnormality            -0.148447  1.438092
## restecgleft ventricular hypertrophy     -3.114198  4.241655
## thalachh                                -0.005222  0.038749
## exngNo                                 -10.731513  3.760027
## oldpeak                                 -1.028318 -0.063216
## slpflat                                 -2.242791  1.495004
## slpdownsloping                          -1.584313  2.568368
## caa                                     -1.339595 -0.414559
## thallfixed defect                       -1.328880  1.931474
## thallreversable defect                  -2.672177  0.457711
## age:sexF                                -0.118161  0.088319
## age:cpatypical angina                   -0.152830  0.107876
## age:cpnon-anginal pain                  -0.186104  0.058026
## age:cpasymptomatic                      -0.121000  0.164541
## age:trtbps                              -0.003296  0.002183
## age:chol                                -0.001099  0.001019
## age:fbsfalse                            -0.067207  0.283652
## age:exngNo                              -0.061956  0.142882
## sexF:chol                               -0.013111  0.029346
## cpatypical angina:chol                  -0.046956  0.014916
## cpnon-anginal pain:chol                 -0.021351  0.025497
## cpasymptomatic:chol                     -0.065221  0.020463
## chol:fbsfalse                           -0.011632  0.037242
## chol:exngNo                             -0.013360  0.030505
```

Stepwise regression for model with interactions

```
## Start:  AIC=257.63
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age * sex + age * cp +
##     trtbps * age + chol * age + fbs * age + exng * age + sex *
##     chol + chol * cp + chol * fbs + chol * exng
##
##              Df Deviance    AIC
## - age:cp      3   193.11 253.11
## - cp:chol     3   194.15 254.15
## - age:chol    1   191.64 255.64
## - age:sex     1   191.71 255.71
## - age:trtbps  1   191.79 255.79
## - sex:chol    1   192.21 256.21
## - chol:exng   1   192.23 256.23
## - age:exng    1   192.24 256.24
## - restecg     2   194.25 256.25
## - slp         2   194.70 256.70
## - chol:fbs    1   192.71 256.71
## - age:fbs     1   193.08 257.08
## <none>            191.63 257.63
## - thalachh    1   193.94 257.94
## - oldpeak     1   196.97 260.97
## - thall       2   201.97 263.97
## - caa         1   207.81 271.81
##
## Step:  AIC=253.11
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:sex + age:trtbps +
##     age:chol + age:fbs + age:exng + sex:chol + cp:chol + chol:fbs +
##     chol:exng
##
##              Df Deviance    AIC
## - cp:chol     3   195.40 249.40
## - age:chol    1   193.12 251.12
## - age:sex     1   193.35 251.35
## - age:trtbps  1   193.38 251.38
## - age:exng    1   193.40 251.40
## - chol:exng   1   193.71 251.71
## - sex:chol    1   193.79 251.79
## - slp         2   196.04 252.04
## - chol:fbs    1   194.15 252.15
## - restecg     2   196.35 252.35
## - age:fbs     1   194.49 252.49
## <none>            193.11 253.11
## - thalachh    1   195.46 253.46
## - oldpeak     1   198.97 256.97
## + age:cp      3   191.63 257.63
## - thall       2   205.82 261.82
## - caa         1   208.37 266.37
##
## Step:  AIC=249.4
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:sex + age:trtbps +
##     age:chol + age:fbs + age:exng + sex:chol + chol:fbs + chol:exng
```

```
##
##                 Df Deviance    AIC
## - age:chol     1    195.41 247.41
## - chol:exng    1    195.70 247.70
## - age:sex      1    195.73 247.73
## - age:trtbps   1    195.73 247.73
## - age:exng     1    195.82 247.82
## - restecg      2    198.48 248.48
## - slp          2    198.52 248.52
## - age:fbs      1    196.54 248.54
## - sex:chol     1    196.91 248.91
## - chol:fbs     1    197.01 249.01
## <none>              195.40 249.40
## - thalachh     1    198.18 250.18
## + cp:chol      3    193.11 253.11
## - oldpeak      1    201.40 253.40
## + age:cp       3    194.15 254.15
## - thall        2    209.82 259.82
## - cp           3    215.74 263.74
## - caa          1    211.75 263.75
##
## Step:  AIC=247.41
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:sex + age:trtbps +
##     age:fbs + age:exng + sex:chol + chol:fbs + chol:exng
##
##                 Df Deviance    AIC
## - chol:exng    1    195.70 245.70
## - age:trtbps   1    195.74 245.74
## - age:sex      1    195.74 245.74
## - age:exng     1    195.84 245.84
## - restecg      2    198.49 246.49
## - slp          2    198.53 246.53
## - age:fbs      1    196.59 246.59
## - chol:fbs     1    197.02 247.02
## - sex:chol     1    197.35 247.35
## <none>              195.41 247.41
## - thalachh     1    198.18 248.18
## + age:chol     1    195.40 249.40
## + cp:chol      3    193.12 251.12
## - oldpeak      1    201.44 251.44
## + age:cp       3    194.15 252.15
## - thall        2    209.83 257.83
## - cp           3    215.91 261.91
## - caa          1    212.28 262.29
##
## Step:  AIC=245.7
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:sex + age:trtbps +
##     age:fbs + age:exng + sex:chol + chol:fbs
##
##                 Df Deviance    AIC
## - age:sex      1    196.00 244.00
## - age:trtbps   1    196.05 244.05
```

```
## - age:exng     1    196.24 244.24
## - restecg      2    198.74 244.74
## - age:fbs      1    196.80 244.80
## - slp          2    198.87 244.87
## - chol:fbs     1    197.23 245.23
## <none>              195.70 245.70
## - sex:chol     1    197.93 245.93
## - thalachh     1    198.51 246.51
## + chol:exng    1    195.41 247.41
## + age:chol     1    195.70 247.70
## + cp:chol      3    193.71 249.71
## - oldpeak      1    201.82 249.82
## + age:cp       3    194.51 250.51
## - thall        2    210.42 256.42
## - cp           3    216.37 260.37
## - caa          1    212.95 260.95
##
## Step:  AIC=244
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:trtbps + age:fbs +
##     age:exng + sex:chol + chol:fbs
##
##               Df Deviance    AIC
## - age:trtbps   1    196.33 242.33
## - age:exng     1    196.55 242.55
## - restecg      2    198.80 242.80
## - slp          2    198.97 242.97
## - age:fbs      1    197.09 243.09
## - chol:fbs     1    197.51 243.51
## <none>              196.00 244.00
## - sex:chol     1    198.05 244.05
## - thalachh     1    198.78 244.78
## + age:sex      1    195.70 245.70
## + chol:exng    1    195.74 245.74
## + age:chol     1    195.99 245.99
## + cp:chol      3    193.96 247.96
## - oldpeak      1    202.26 248.26
## + age:cp       3    194.62 248.62
## - thall        2    210.53 254.53
## - cp           3    216.44 258.44
## - caa          1    213.43 259.43
##
## Step:  AIC=242.33
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:fbs + age:exng +
##     sex:chol + chol:fbs
##
##               Df Deviance    AIC
## - age:exng     1    196.84 240.84
## - restecg      2    199.21 241.21
## - slp          2    199.40 241.40
## - age:fbs      1    197.63 241.63
## - chol:fbs     1    197.78 241.78
## - sex:chol     1    198.23 242.23
```

```
## <none>                     196.33 242.33
## - trtbps      1    198.57 242.57
## - thalachh    1    199.29 243.29
## + age:trtbps  1    196.00 244.00
## + chol:exng   1    196.05 244.05
## + age:sex     1    196.05 244.05
## + age:chol    1    196.32 244.32
## + cp:chol     3    194.23 246.23
## - oldpeak     1    202.42 246.42
## + age:cp      3    194.84 246.84
## - thall       2    211.22 253.22
## - cp          3    216.90 256.90
## - caa         1    213.53 257.53
##
## Step:  AIC=240.84
## output ~ age + sex + cp + trtbps + chol + fbs + restecg + thalachh +
##     exng + oldpeak + slp + caa + thall + age:fbs + sex:chol +
##     chol:fbs
##
##               Df Deviance    AIC
## - restecg     2    199.47 239.47
## - slp         2    200.00 240.00
## - age:fbs     1    198.09 240.09
## - chol:fbs    1    198.20 240.20
## - sex:chol    1    198.72 240.72
## <none>                     196.84 240.84
## - trtbps      1    199.15 241.15
## - thalachh    1    199.68 241.68
## - exng        1    199.84 241.84
## + age:exng    1    196.33 242.33
## + chol:exng   1    196.46 242.46
## + age:sex     1    196.54 242.54
## + age:trtbps  1    196.55 242.55
## + age:chol    1    196.83 242.83
## - oldpeak     1    202.57 244.57
## + cp:chol     3    194.66 244.66
## + age:cp      3    195.77 245.77
## - thall       2    211.25 251.25
## - cp          3    217.75 255.75
## - caa         1    214.58 256.58
##
## Step:  AIC=239.47
## output ~ age + sex + cp + trtbps + chol + fbs + thalachh + exng +
##     oldpeak + slp + caa + thall + age:fbs + sex:chol + chol:fbs
##
##               Df Deviance    AIC
## - age:fbs     1    200.68 238.68
## - chol:fbs    1    201.00 239.00
## - sex:chol    1    201.10 239.10
## - slp         2    203.16 239.16
## <none>                     199.47 239.47
## - trtbps      1    202.14 240.14
## - thalachh    1    202.35 240.35
## - exng        1    202.47 240.47
```

```
## + restecg     2   196.84 240.84
## + age:trtbps  1   199.11 241.11
## + chol:exng   1   199.13 241.13
## + age:exng    1   199.21 241.21
## + age:sex     1   199.41 241.41
## + age:chol    1   199.47 241.47
## - oldpeak     1   205.22 243.22
## + cp:chol     3   197.38 243.38
## + age:cp      3   197.94 243.94
## - thall       2   212.82 248.82
## - cp          3   220.30 254.30
## - caa         1   216.99 254.99
##
## Step:  AIC=238.68
## output ~ age + sex + cp + trtbps + chol + fbs + thalachh + exng +
##     oldpeak + slp + caa + thall + sex:chol + chol:fbs
##
##              Df Deviance    AIC
## - age         1   200.71 236.71
## - sex:chol    1   202.15 238.15
## - chol:fbs    1   202.33 238.33
## - slp         2   204.65 238.65
## <none>            200.68 238.68
## + age:fbs     1   199.47 239.47
## - exng        1   203.53 239.53
## - trtbps      1   203.71 239.71
## - thalachh    1   203.74 239.74
## + restecg     2   198.09 240.09
## + age:trtbps  1   200.11 240.11
## + chol:exng   1   200.43 240.43
## + age:exng    1   200.44 240.44
## + age:sex     1   200.62 240.62
## + age:chol    1   200.63 240.63
## - oldpeak     1   206.13 242.13
## + cp:chol     3   198.88 242.88
## + age:cp      3   199.03 243.03
## - thall       2   213.14 247.14
## - cp          3   221.07 253.07
## - caa         1   217.76 253.76
##
## Step:  AIC=236.7
## output ~ sex + cp + trtbps + chol + fbs + thalachh + exng + oldpeak +
##     slp + caa + thall + sex:chol + chol:fbs
##
##              Df Deviance    AIC
## - sex:chol    1   202.19 236.19
## - chol:fbs    1   202.34 236.34
## - slp         2   204.68 236.68
## <none>            200.71 236.71
## - exng        1   203.54 237.54
## - trtbps      1   203.96 237.96
## + restecg     2   198.09 238.09
## + chol:exng   1   200.46 238.46
## - thalachh    1   204.49 238.49
```

24

```
## + age          1   200.68 238.68
## - oldpeak      1   206.16 240.16
## + cp:chol      3   198.93 240.93
## - thall        2   213.23 245.23
## - cp           3   221.08 251.08
## - caa          1   218.80 252.80
##
## Step:  AIC=236.19
## output ~ sex + cp + trtbps + chol + fbs + thalachh + exng + oldpeak +
##     slp + caa + thall + chol:fbs
##
##                Df Deviance    AIC
## - chol:fbs     1   203.77 235.77
## <none>             202.19 236.19
## + sex:chol     1   200.71 236.71
## - exng         1   205.00 237.00
## - slp          2   207.10 237.10
## - trtbps       1   205.57 237.57
## + chol:exng    1   201.74 237.74
## + restecg      2   199.80 237.80
## - thalachh     1   205.80 237.80
## + age          1   202.15 238.15
## - oldpeak      1   207.14 239.14
## + cp:chol      3   199.58 239.58
## - thall        2   214.32 244.32
## - sex          1   212.73 244.73
## - cp           3   223.47 251.47
## - caa          1   220.29 252.29
##
## Step:  AIC=235.77
## output ~ sex + cp + trtbps + chol + fbs + thalachh + exng + oldpeak +
##     slp + caa + thall
##
##                Df Deviance    AIC
## - fbs          1   203.88 233.88
## <none>             203.77 235.77
## - chol         1   205.86 235.86
## + chol:fbs     1   202.19 236.19
## + sex:chol     1   202.34 236.34
## - slp          2   208.59 236.59
## - exng         1   206.71 236.71
## - trtbps       1   206.95 236.95
## + restecg      2   201.19 237.19
## - thalachh     1   207.31 237.31
## + chol:exng    1   203.45 237.45
## + age          1   203.75 237.75
## + cp:chol      3   200.74 238.74
## - oldpeak      1   208.83 238.83
## - sex          1   213.77 243.77
## - thall        2   215.83 243.83
## - cp           3   224.60 250.60
## - caa          1   222.60 252.60
##
## Step:  AIC=233.88
```

```
## output ~ sex + cp + trtbps + chol + thalachh + exng + oldpeak +
##     slp + caa + thall
##
##               Df Deviance    AIC
## <none>           203.88 233.88
## - chol        1   205.93 233.93
## + sex:chol    1   202.46 234.46
## - slp         2   208.63 234.63
## - exng        1   206.75 234.75
## - trtbps      1   206.96 234.96
## + restecg     2   201.34 235.34
## - thalachh    1   207.44 235.44
## + chol:exng   1   203.59 235.59
## + fbs         1   203.77 235.77
## + age         1   203.86 235.86
## + cp:chol     3   201.08 237.08
## - oldpeak     1   209.12 237.12
## - sex         1   213.78 241.78
## - thall       2   216.18 242.18
## - cp          3   226.06 250.06
## - caa         1   222.65 250.65


##
## Call:  glm(formula = output ~ sex + cp + trtbps + chol + thalachh +
##     exng + oldpeak + slp + caa + thall, family = "binomial",
##     data = newheart)
##
## Coefficients:
##           (Intercept)                       sexF        cpatypical angina
##              1.025922                   1.535349                 1.000189
##     cpnon-anginal pain           cpasymptomatic                   trtbps
##              1.943932                   1.965324                -0.017488
##                  chol                   thalachh                   exngNo
##             -0.005343                   0.018627                 0.722047
##               oldpeak                    slpflat             slpdownsloping
##             -0.484829                  -0.711536                 0.228519
##                   caa          thallfixed defect  thallreversable defect
##             -0.817229                  -0.032660                -1.345274
##
## Degrees of Freedom: 300 Total (i.e. Null);  286 Residual
## Null Deviance:        414.8
## Residual Deviance: 203.9      AIC: 233.9
```

### Model Comparision

I noticed when adding interactions to the model the significance of the variables in the model became less effective. We can see this by looking at the both the p-values in the wald test. All the p-values for the model with interactions are significantly more higher than the additative model. We can also see that When running step wise regression on the model with interactions. It eliminated all the interaction terms because they were insigificant to the model and gave us back the additative model. This is telling us that the additative model will be more viable to use for the analysis. Thus after evaluating both models we will stick with the additative model for further analysis.

```
## [1] "R-squared: 0.508545330225056"
```

#Model evaluation

## VIF

Since our selected model is the additative model we can now use the model to answer the question. Which predictor influences the factors of having a heart attack. As metioned before during the correlation test we notice there was some aspects of multicollinearity in the dataset. Lets check the VIF of the additative model. In statistics, the variance inflation factor is the quotient of the variance in a model with multiple terms by the variance of a model with one term alone. It quantifies the severity of multicollinearity in an ordinary least squares regression analysis

```
##              GVIF Df GVIF^(1/(2*Df))
## sex      1.491120  1        1.221114
## cp       1.569147  3        1.077980
## trtbps   1.094189  1        1.046035
## chol     1.187864  1        1.089892
## thalachh 1.318093  1        1.148082
## exng     1.176216  1        1.084535
## oldpeak  1.502516  1        1.225772
## slp      1.805862  2        1.159234
## caa      1.150145  1        1.072448
## thall    1.312593  2        1.070367
```

Since the VIF score of our predictors are less than 10. There is no aspects of multicollinearity in our model.
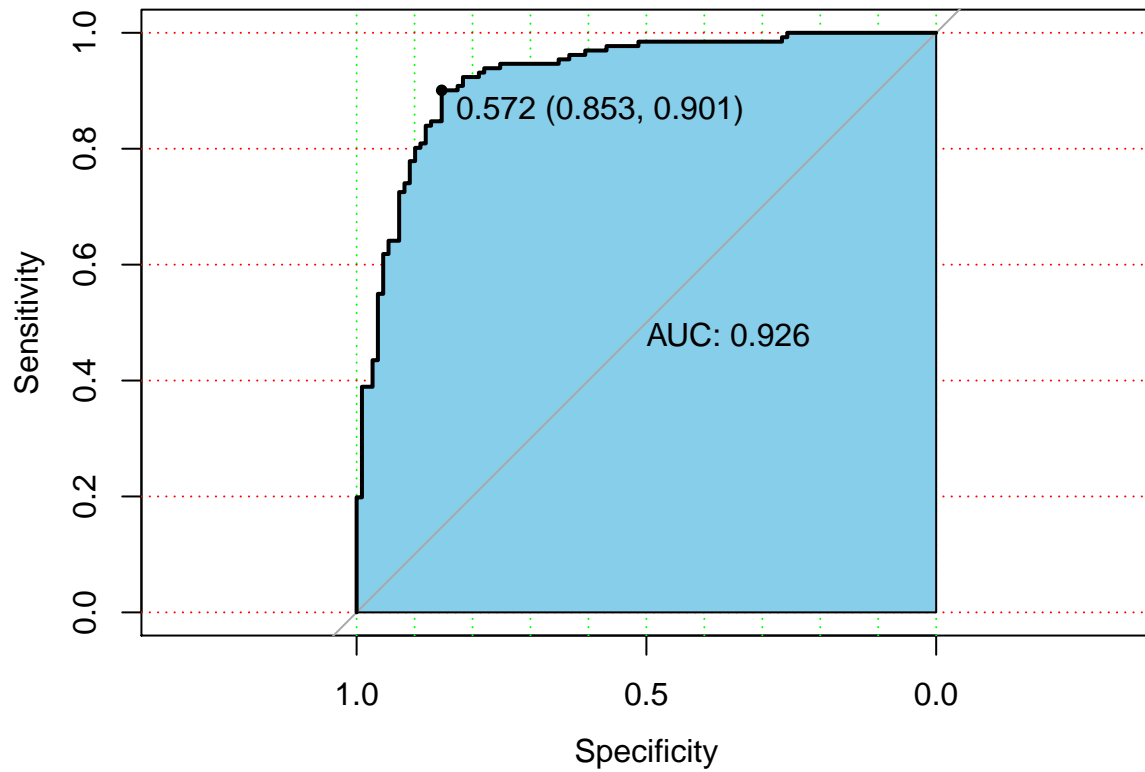
#Model Analysis

## ROC Curve for Logistic Regression

The ROC curve shows the trade-off between sensitivity (or TPR) and specificity (1 – FPR). Classifiers that give curves closer to the top-left corner indicate a better performance. As a baseline, a random classifier is expected to give points lying along the diagonal (FPR = TPR). The closer the curve comes to the 45-degree diagonal of the ROC space, the less accurate the test. To compare different classifiers, it can be useful to summarize the performance of each classifier into a single measure. One common approach is to calculate the area under the ROC curve, which is abbreviated to AUC

```
## Setting levels: control = Less Chance, case = More Chance
```
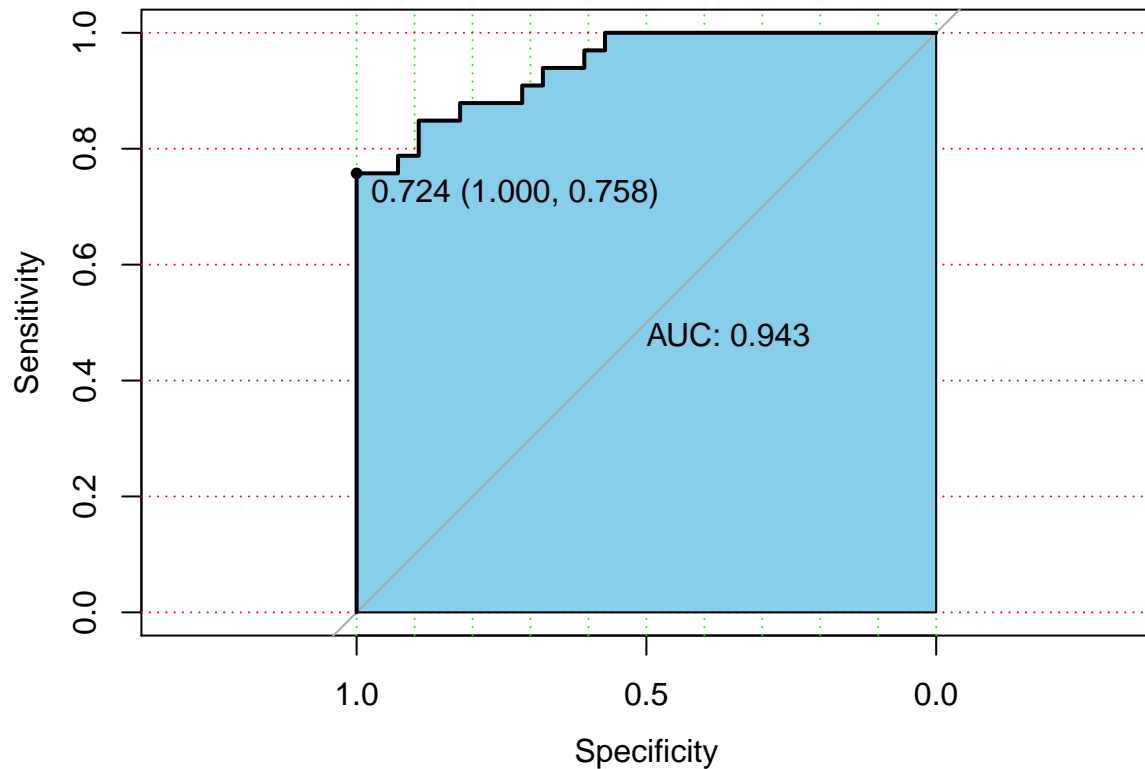
```
## Setting direction: controls < cases
```

The AUC on training set is 0.918 which is approximately 92%, This is an indication that the model prediction performance is good.

```
## Setting levels: control = Less Chance, case = More Chance
```

```
## Setting direction: controls < cases
```
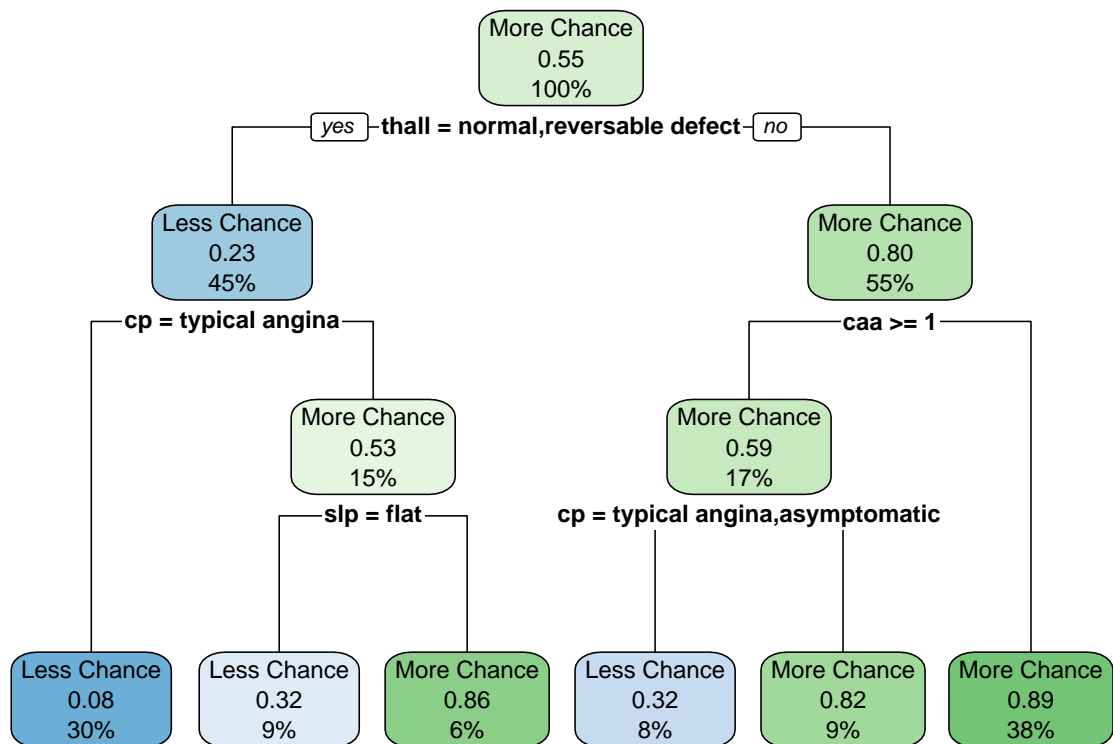
The AUC on testing set is 0.974 which is approximetly 97%. This is an indication that the model prediction performance is good.

Overall comparing the model performance on both the testing and training set we get similar performances which tells us that the model we have is able to classify whether the patient is highly in risk of having a heart attack or less in risk of having a heart attack, based on logistic regression.

## Decision Tree Classification

We will use our aditative model for decision tree classification also. The Decision tree builds classification or regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. the Decision trees can handle both categorical and numerical data

More Chance
0.55
100%

yes ⊣ **thall = normal,reversable defect** ⊢ no

Less Chance
0.23
45%

More Chance
0.80
55%

**cp = typical angina**

**caa >= 1**

More Chance
0.53
15%

More Chance
0.59
17%

**slp = flat**

**cp = typical angina,asymptomatic**

Less Chance
0.08
30%

Less Chance
0.32
9%

More Chance
0.86
6%

Less Chance
0.32
8%

More Chance
0.82
9%

More Chance
0.89
38%

At the top is the overal percentage of patients getting heart attacks. It shows the porportion of patients who have a higher chance of recieving a heart attack to patients having a less chance of recieving a heart attack. 53% of the patients have a higher chance of recieving a heart attack. You can keep going down the nodes to understand what features impact the chances of having a higher chance of a heart attack to having a less chance of having a heart attack. For example, if its a fixed defect and if the chest pain type is asymptomatic then there is a 91% chance of having a heart attack.

## Making A Prediction

Accuracy of test set

```
##              predict_unseen
##             Less Chance More Chance
##   Less Chance          24           4
##   More Chance           8          25
```

Confusion matrix of the patients who have more chance of a heart attack and the patients who have a less chance of a heart attack using the testing data set.

```
## [1] "Accuracy for test 0.80327868852459"
```

The accuracy of the testing set is approximatly 85%. Which is good because it tells us that 82% of the predictions made are correct based on the testing set.

Accuracy of training set

```
##                 predict_unseen1
##               Less Chance More Chance
##   Less Chance           93         16
##   More Chance           19        112
```

```
## [1] "Accuracy for train 0.854166666666667"
```

The accuracy of the training set is approximatly 84%. This is good because it tells us that 84% of the predictions made are correct based on the training set.

#Conclusion

Comparing the performance of our model through logistic regression and using the decision tree method we can see that the model performs well to predict whether the patient is classified in having a higher chance of getting a heart attack to having a less chance of having a heart attack. We can also see that Logistic regression seems to be more viable since it has a higher chance of prediciting more or less chance of having a heart attack compared to the decision tree method. Finally, we see that the predictors that show the most influence to our independant variable 'output' is sex, cp, trtbps, chol, thalachh, exng, oldpeak, slp, caa and the thall. Theoretically speaking this makes sense. Factors such as chest pain, cholesterol levels, etc, are related to heart attacks. Some recommendations that will probably improve the performance of the model is having more meaningful predictors in the data set for detecting heart attacks.