

## **Modeling Heart Attack Predictions With Patient Data**

Filipp Krasovsky, Rudy Fasano

University of San Diego

Master of Science, Applied Data Science

ADS 503

Section 3

6-28-2021

## Background

Heart attacks occupy a central role in the social consciousness as a commonly known medical condition with incredibly costly medical consequences as well as considerable financial implications for both the victims and their families. Therefore, our team has a powerful social motivation in exploring a data-driven approach to identifying a heart attack in cases of medical ambiguity, as well as predicting the incidence of a heart attack based on known symptoms and medical information about a given patient. We hope that our findings are applied, and improved upon, in medical applications to increase the quality of life.

In this report, we set the context of our findings as being applied to an individual for whom we have easily obtainable medical data – basic information such as age and sex, as well as information that could be determined through short-term medical procedures – resting heart rate, blood sugar, and indications of angina. We further frame the data question as follows:

*“Given some set of available medical information about a patient, are they more likely than not to experience a heart attack?”*

In our investigation, we conclude that a logistic regression model provides the most robust predictive ability with an ROC of approximately 96.6% - this performance was the top out of all other models considered, which include KNN, SVM, NN, GBM, RF, NSC.

## Table of Contents

Background .....	2
Inventory of Resources .....	4
Terminology.....	5
Objectives & Success Criteria .....	5
Exploratory Data Analysis.....	5
Data Description/ Observations .....	8
Data Pre-processing and Splitting.....	16
Modelling Strategy.....	16
Model Performance and Hyperparameter Tuning .....	17
Final Model Selection / Results .....	18
References.....	19

## **Inventory of Resources**

**Data Source:** The heart attack analysis prediction dataset from Kaggle posted by Rashik Rahman in March of 2021.

**Software:** Rstudio version 1.4.1106, R Version 4.04 – “Lost Library Book”

## **Terminology**

1. Angina – a condition characterized by severe chest pain that spreads to the shoulders and arms caused by inadequate blood supply to the heart.
2. Resting Blood Pressure – titularly, refers to blood pressure when not doing any physical activity. Pressure below 120/80mm Hg is considered normal.

## **Objectives & Success Criteria**

The primary goal of this project is a classification problem of predicting the incidence of heart attack where class 0 refers to no incidence and class 1 refers to an incidence of heart attack. In this instance, we focus on optimizing ROC and place equal weight on false positive and false negative results, as both result in medical treatment that may create suboptimal results. While we understand that the consequences of both are incommensurate, we lack the medical foreground to calculate the risks non-arbitrarily.

## **Exploratory Data Analysis**

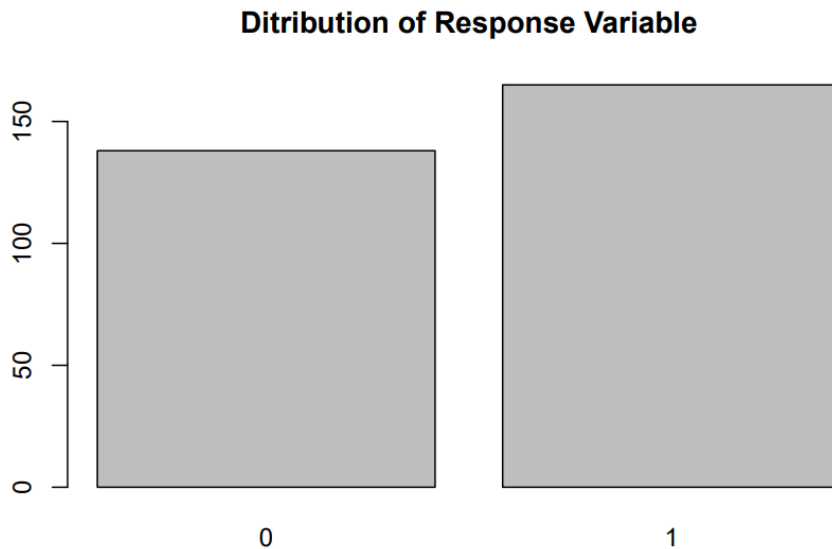
Our dataset consists of  $n=303$  observations with a predictor space of 13 variables stored inside a single CSV file along with a response variable. The following descriptions were salvaged from the context provider by the dataset author:

Field	Description	Type
Output	Response where 1=heart attack	Nominal
Age	Age of patient	Ratio
Sex	0=Female, 1=Male	Nominal
CP	Chest Pain (Angina)  1=typical  2=atypical  3=non-anginal pain  4=asymptomatic	Nominal
Trtbps	Resting Blood Pressure	Ratio
Chol	Cholesterol levels	Ratio
Fbs	1 = fasting blood sugar >  120mg/dl	Nominal
Restecg	Resting electrocardiographic results:  0 = normal  1 = abnormal  2 = hypertrophy	Nominal
Thalachh	Maximum heart rate achieved	Ratio
Exng	1 = exercise induced angina	Nominal
Oldpeak	Previous peak – no context.	Ratio
Slp	Slope – no context.	Interval
Caa	Number of major vessels	Interval

Thal	“thal rate” – no context provided.	Discrete Ratio
------	---------------------------------------	----------------

## Data Description/ Observations

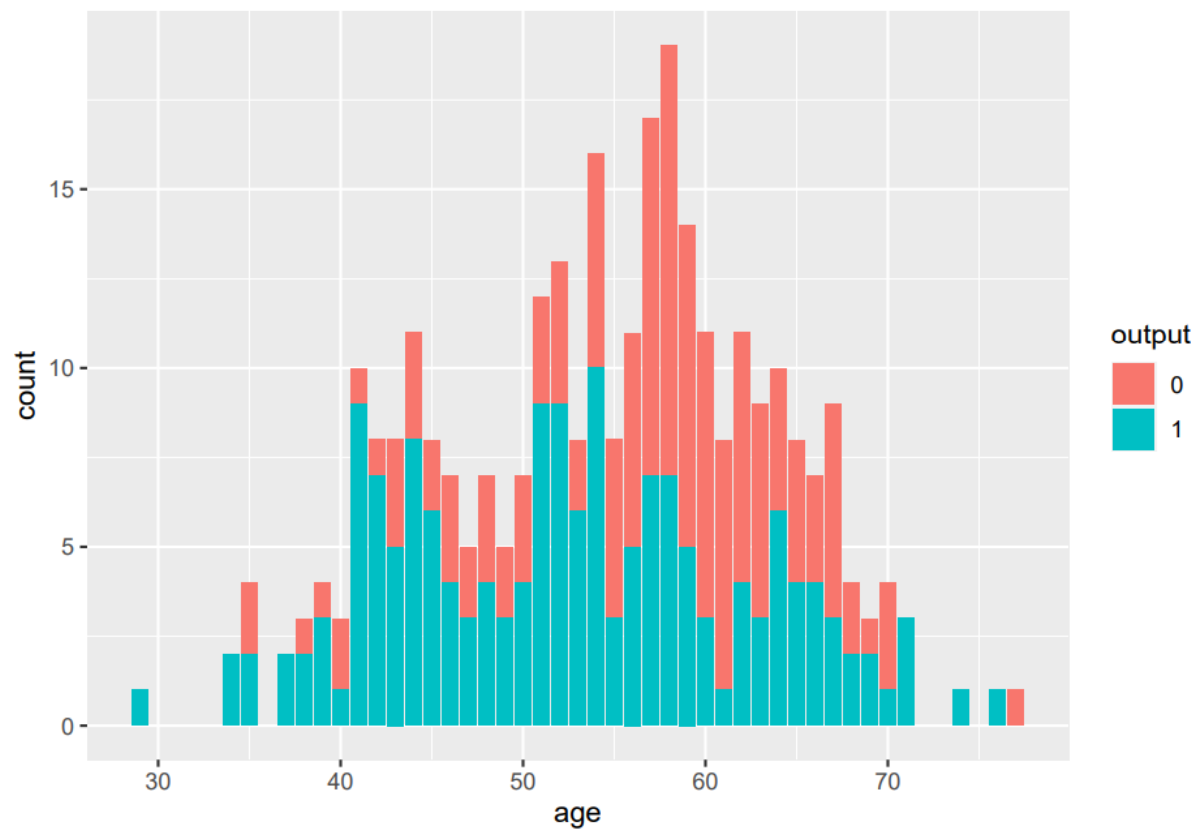
We begin by noting the absence of any missing values in our data, which makes cleaning and pre-processing considerably easier. We also find a near-perfect class balance, which allows us to sample randomly and without stratification:

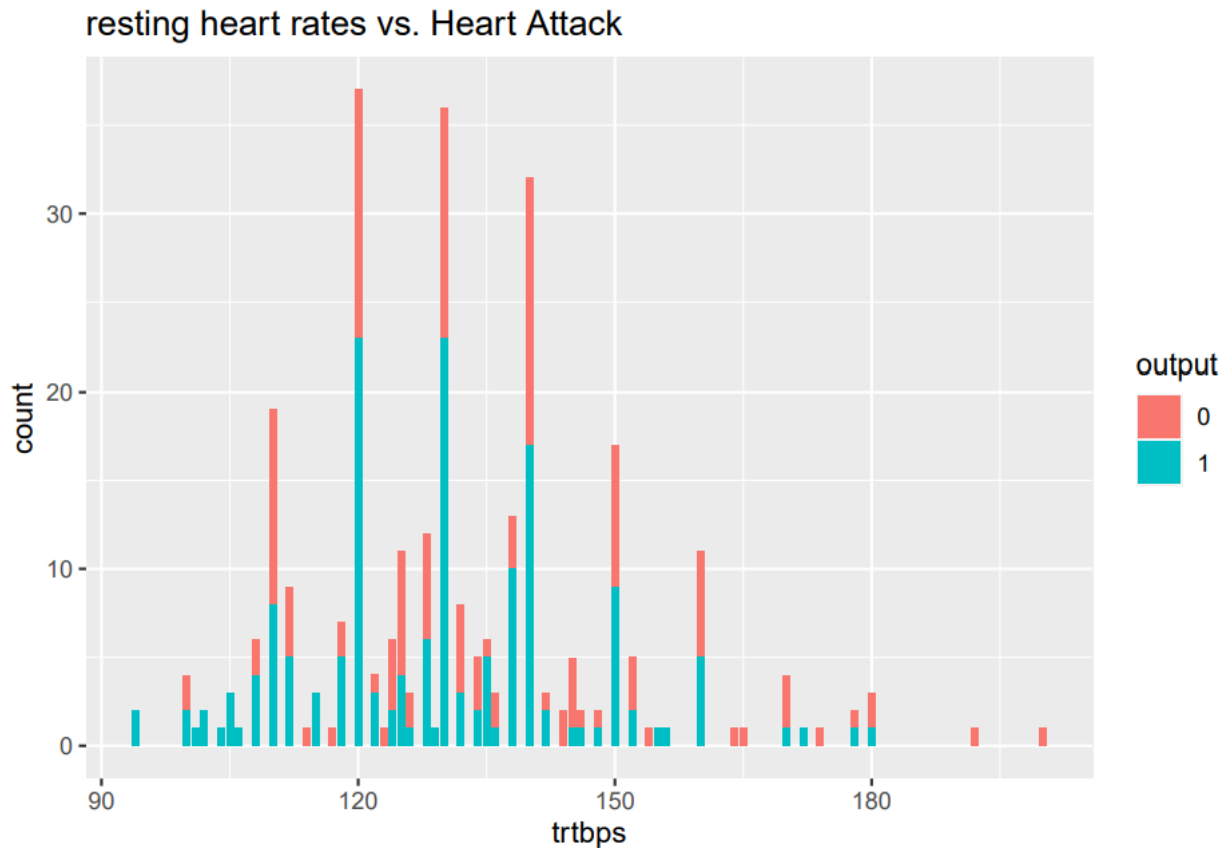


Moving to analysing our predictor space, we find no significant correlations in our continuous space, and can begin observing the relationships between predictors and the response variable.

To begin with, we find a notable pattern in age, wherein lower ages are more disproportionately affected by the incidence of heart attack:







When we observe this relationship for resting heart rates (trtbps), the most obvious face-value observation is that resting heart rate values smaller than 120 bps are disproportionately associated with the risk of heart attack (output=1). This finding confirms the intuitions of our team, which presupposes that lower resting heart rates are linked to a higher incidence of risk.

Although we preliminary suspected that cholesterol levels would have a significant bearing on heart attack risk, we found that there doesn't seem to be a significant difference between response groups:

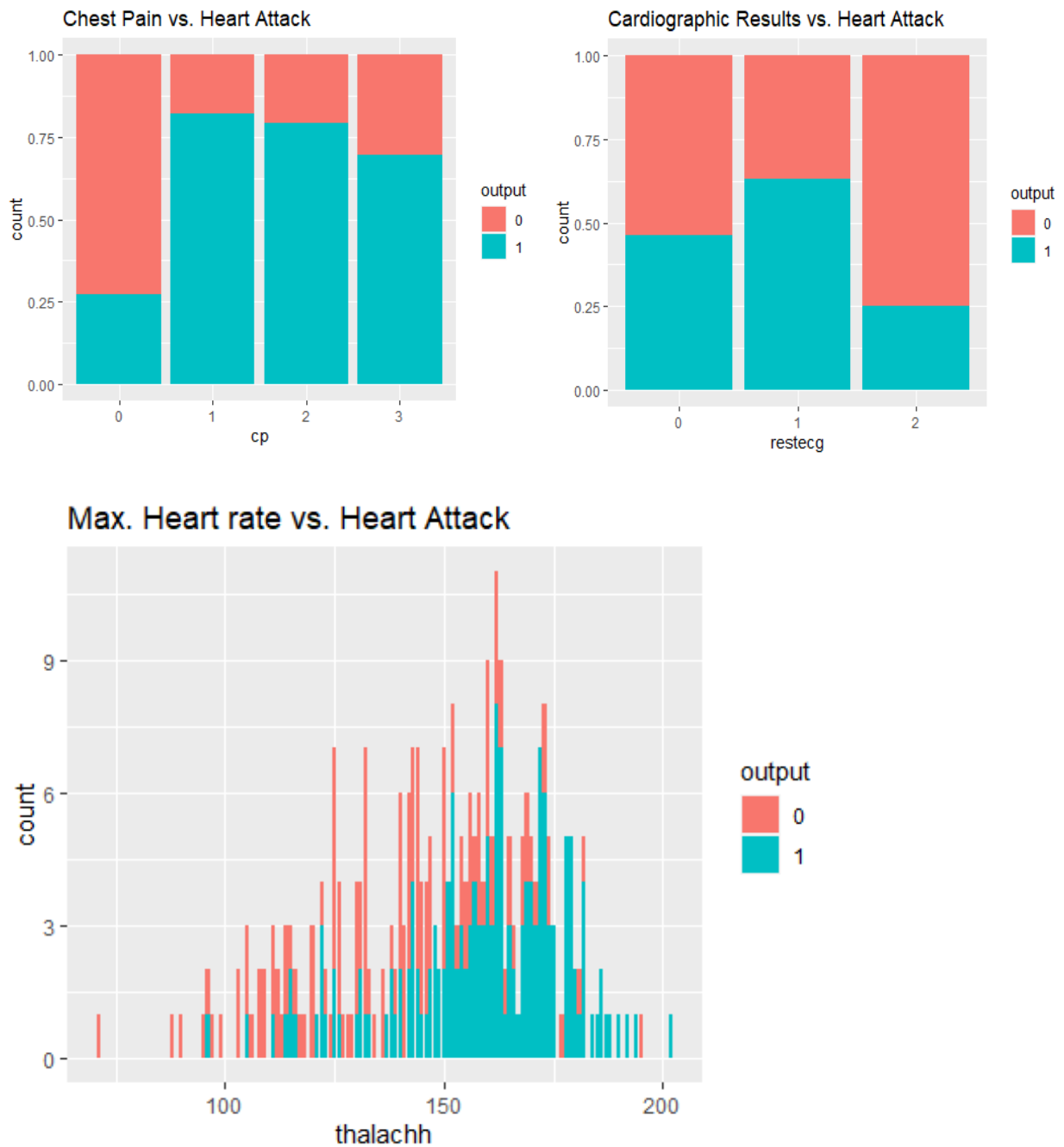


To further fortify the claim that cholesterol is not an informative predictor, a hypothesis test for the difference of means was conducted at an  $\alpha = 0.05$ . Based on these findings, we do not have enough evidence to reject the null hypothesis that the difference of means between the two groups is zero. We did, however, identify a gender disparity with respect to the response variable:

	Heart Attack	No Heart Attack	<b>Total</b>
Male	114	93	207
Female	24	72	96

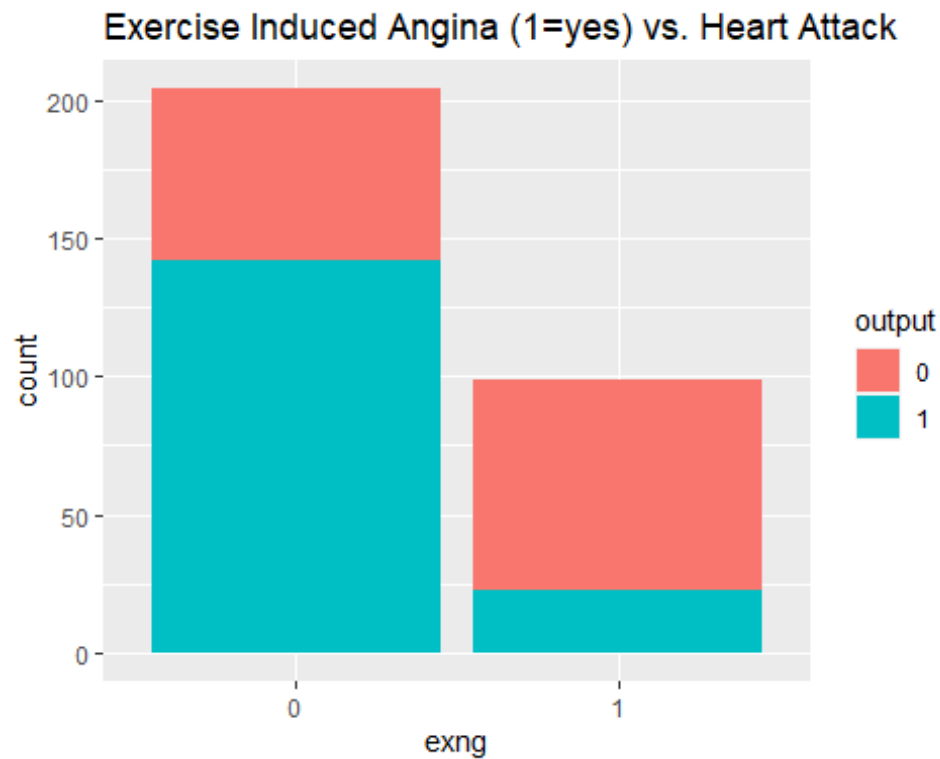
We also note the imbalance in gender between men and women, and recognize that this might pose a risk in representation, and would therefore welcome an opportunity to find more representative data. Other notable predictors that clearly demonstrate a separation between the two response groups include the incidence of chest pain (cp), cardiographic results (restecg), the maximum heart

rate achieved (thalachh), the presence of exercise induced angina (exng), slope and thal, which had no interpretation per se, the number of major vessels, and fasting blood sugar:



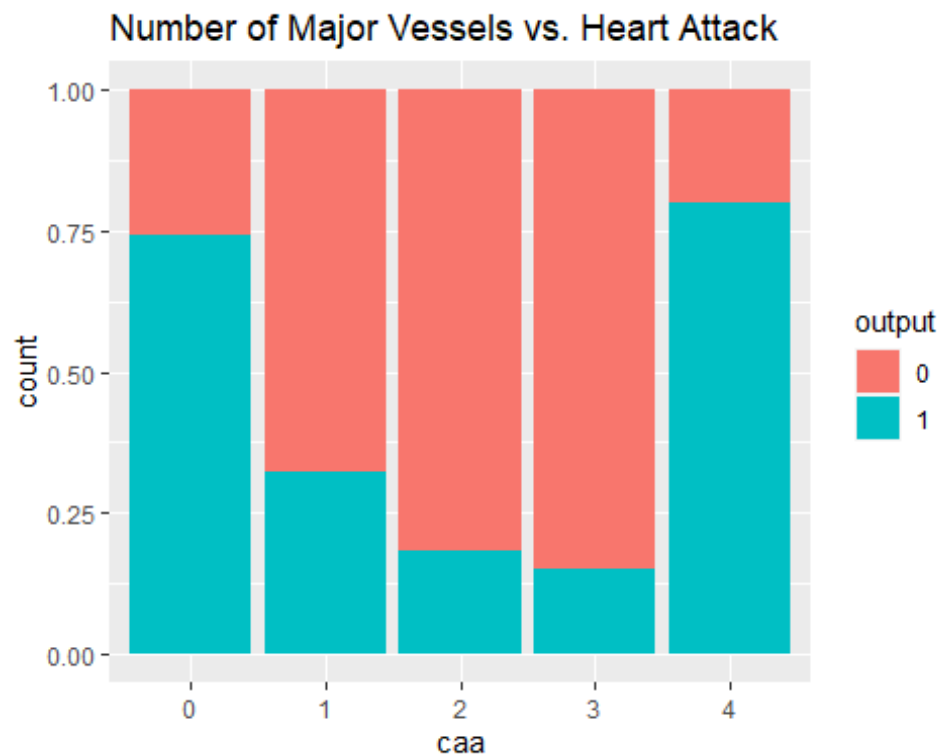
Here, we find the absence of chest pain is strongly associated with the absence of heart attack. Hypertrophy in cardiographic results is also strongly associated with a reduced risk of heart attack. Finally, we find that a higher maximum heart rate is associated with a higher incidence of heart

attack, which was further confirmed by a hypothesis test at  $\alpha=0.05$ , suggesting that we have enough evidence to reject the null hypothesis that the difference of means between the two response groups is zero.



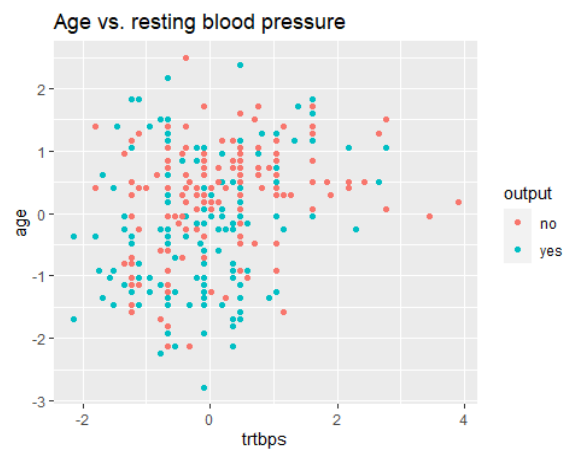
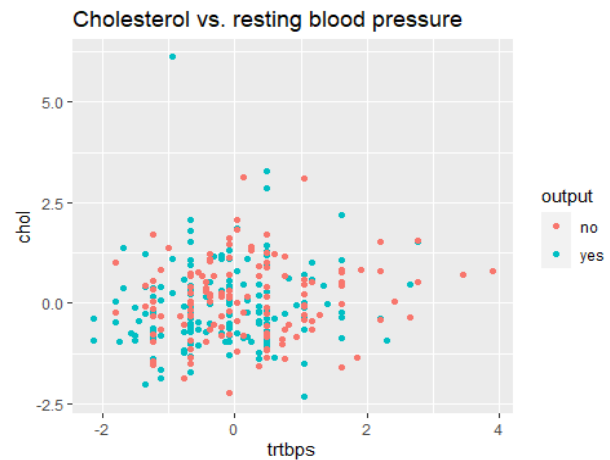
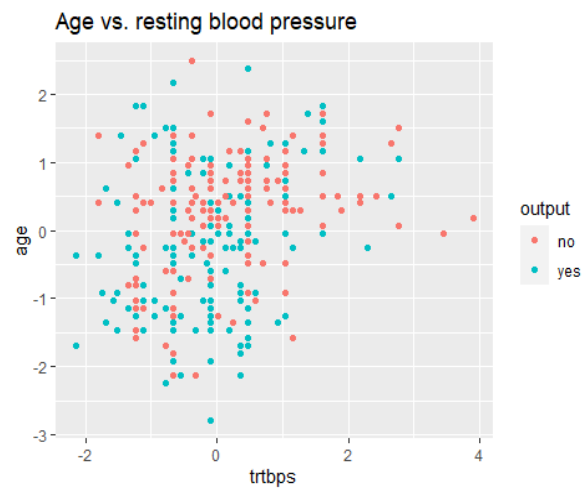
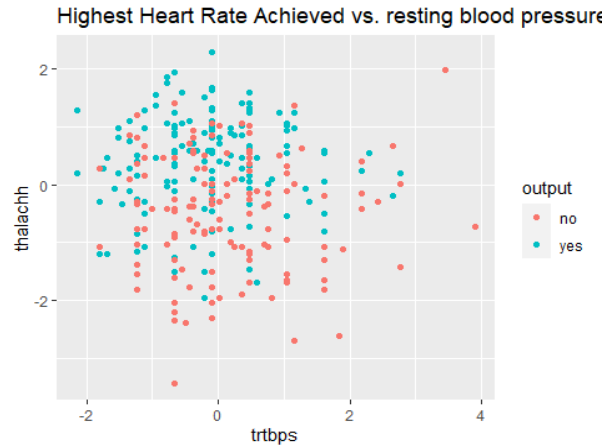
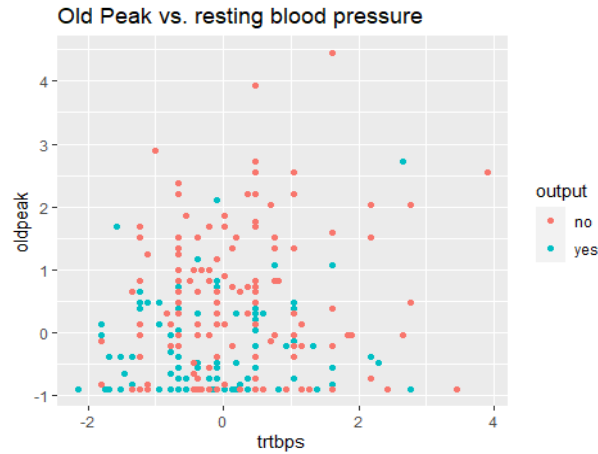
Curiously, we observe that individuals who suffered exercise induced angina are less likely to contract heart disease. From an intuitive point of view, the argument could be made that exercise decreases the risk of heart attack in many instances, and so the angina experienced by patients isn't associated with a heart condition.

Finally, we find that individuals with 0 and 4 major vessels are considerably more associated with the risk of heart attack than all other values in between:



For this report, we choose not to visualize slope or  $\theta_{\text{al}}$ , as we lack the interpretation for either of these variables, but still choose to include them during the modelling process.

The last portion of EDA includes a discussion of whether or not this problem is linearly separable- as a preliminary attempt to understand the data environment, we visualize our continuous variables and attempt to find face-value evidence that the problem might be separated by a line in a two dimensional space:



Based on a preliminary assessment of the data, it's unclear that there's an opportunity to create a linear separation in the continuous predictor space, and so we're unlikely to use LDA or any other linear separators.

## **Data Pre-processing and Splitting**

For the modelling portion of our investigation, predictors were scaled and centered. For the Neural Net portion of the modelling process, a spatial sign transformation was applied to continuous predictors where applicable. Since no missing values exist in the set, interpolation was not required. Data was split into a 90%-10% training/testing set due to the fact that n vastly outpaces the size of the predictor space. Of 303 samples, ~270 were used for the training set.

## **Modelling Strategy**

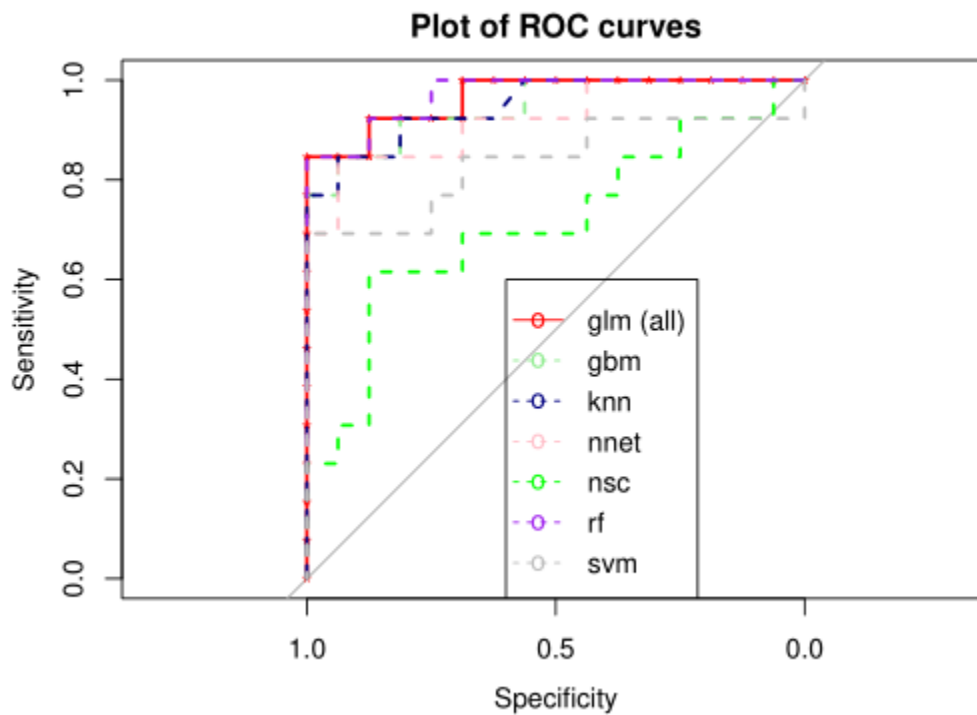
Given that we have a combination of categorical and continuous data, one of the challenges we encountered during the process of creating a model building strategy was that some R packages were resistant to training a combination of both variables. For instance, while logistic regression / glm did end up outperforming all other models, the **train()** function from the caret package was harder to apply during the model building process than the regular glm() function. Because of this, we fielded a combination of tree-based models such as random forest and gbm, regression models like glm, and models of varying complexity such as SVM, Neural Networks, Nearest Shrunken Centroids, and K-Nearest Neighbours. By offering a broad range of classification models that all take fundamentally different approaches to the data problem, we hope to find a balance between ROC maximization and parsimony.



## Model Performance and Hyperparameter Tuning

Of the seven models that were tested (KNN, NSC, SVM, NN, GLM, GBM, and RF), our neural net model performed the best, followed closely by random forest and gbm. Glm, svm, nsc, and knn all had below-90% ROC, of which KNN performed the best.

Model	ROC	Params	Model	ROC	Params
GLM	0.9663	NA	KNN	0.9495	K=40
GBM	0.9471	Trees=1000 depth =1, shrinkage=0.01, minobs=30	NNET	0.9231	Size=5, decay=1
NSC	0.7115	Thresh=0	RF	0.9712	Mtry=2
SVM	0.8365	Sigma=.054 C=.0625			



### Final Model Selection / Results

In this instance, we find that Random Forest provides a robust ROC performance (+97%), which is comparable to the performance of the most parsimonious model, the logistic regression model.

We observe the following variable importance values for RF below:

Variable	Overall
age	58.92492
sex	5.67629
cp	100
trtbps	44.00965
ochol	49.39511
fbs	10.1194
restecg	0
thalachh	90.27633
exng	22.08506
oldpeak	82.97035
slp	27.32064
caa	93.07395
thall	93.66576

Although we would normally consider parsimony a driving factor (ie model complexity) in deciding which model to move forward with, we make a special caveat for the given application, which is determining the incidence of a heart attack. In this particular instance, maximizing ROC is paramount because of the high social cost of misclassification in either direction – failing to identify a heart attack could result in potential death, while misdiagnosing a heart attack and applying defensive treatment might result in incredible financial expenses for the patient or possibly other secondary health detriments as a result of emergency treatment. Because of this, we choose to move forward with the random forest model.

## References

Rashik Rahman, (March 2021). *Heart Attack Analysis & Prediction Dataset*, Kaggle, retrieved June, 2021 from <https://www.kaggle.com/rashikrahmanpritom/heart-attack-analysis-prediction-dataset>