# Homework 5 Julia Fish

**Task 1: Conceptual Questions**

- Question 1: **What is the purpose of using cross-validation when fitting a random forest model?**
  Cross-validation helps estimate a model's performance with "new" data. When fitting a random forest model, cross-validation ensures that the model is not overfitting by evaluating its performance across multiple data splits (i.e. leaving one fold out at a time). Because of this, it gives a more reliable estimate of prediction error than a single train-test split.

- Question 2: **Describe the bagged tree algorithm.**
  The baged tree algorithm involves generating multiple training sets using bootstrapping techniques. A decision tree is trained on each of the bootstrap resamples, and predictions are averaged over. This reduces variance and increases model stability compared to fitting a single tree model.

- Question 3: **What is meant by a general linear model?**
  A general linear model is a linear model that can have different forms of response variable than just all real numbers. For example, binomial, poisson, gamma, etc. can be used in general linear models to help modify the response range as well as oher logistical factors in the model fit.

- Question 4: **When fitting a multiple linear regression model, what does adding an interaction term do?**
  Adding an interaction term to a MLR model allows the effect of one variable on the response to depend on the level of another variable. That is, the variables are allowed to communicate with one another to predict a response value instead of only having just additive effects.

- Question 5: **Why do we split our data into a training and test set?**
  Using a training and test set allows us to see model performance on "unseen" data with known response values. That way we can get a better idea of whether or not our model performs well or has just been overfit to the data it was given to train the model.

## Task 2: Data Prep

### Packages and Data

First, we will load in the packages we need for this task. We will also read in our `heart` data set as a tibble.

```
library(tidyverse)
library(tidymodels)
library(caret)
library(yardstick)


heart <- as_tibble(read_csv("https://www4.stat.ncsu.edu/~online/datasets/heart.csv"))
```

### Question 1

Now, we will run and discuss the `summary()` of this tibble.

```
summary(heart)
```

```
      Age              Sex            ChestPainType         RestingBP
 Min.   :28.00   Length:918         Length:918         Min.   :  0.0
 1st Qu.:47.00   Class :character   Class :character   1st Qu.:120.0
 Median :54.00   Mode  :character   Mode  :character   Median :130.0
 Mean   :53.51                                         Mean   :132.4
 3rd Qu.:60.00                                         3rd Qu.:140.0
 Max.   :77.00                                         Max.   :200.0
  Cholesterol      FastingBS       RestingECG            MaxHR
 Min.   :  0.0   Min.   :0.0000   Length:918         Min.   : 60.0
 1st Qu.:173.2   1st Qu.:0.0000   Class :character   1st Qu.:120.0
 Median :223.0   Median :0.0000   Mode  :character   Median :138.0
 Mean   :198.8   Mean   :0.2331                      Mean   :136.8
 3rd Qu.:267.0   3rd Qu.:0.0000                      3rd Qu.:156.0
 Max.   :603.0   Max.   :1.0000                      Max.   :202.0
 ExerciseAngina       Oldpeak          ST_Slope           HeartDisease
 Length:918         Min.   :-2.6000   Length:918         Min.   :0.0000
 Class :character   1st Qu.: 0.0000   Class :character   1st Qu.:0.0000
 Mode  :character   Median : 0.6000   Mode  :character   Median :1.0000
                    Mean   : 0.8874                      Mean   :0.5534
                    3rd Qu.: 1.5000                      3rd Qu.:1.0000
                    Max.   : 6.2000                      Max.   :1.0000
```

Above, we see the full summary of this data set. The `HeartDisease` variable seems to be a quantitative variable, which does not make sense in this case since it is supposed to represent the presence or absence of heart disease. That should be categorical.

### Question 2

Now, we will create a categorical`HeartDisease` variable named `HD`. In the same pipeline, we will also remove `ST_Slope` and the original `HeartDisease` variable. Lastly, we will save this new version of the data under `new_heart`. We will do this below:

```
new_heart <- heart |>
  mutate(HD = factor(HeartDisease)) |>
  select(-ST_Slope, -HeartDisease)
```
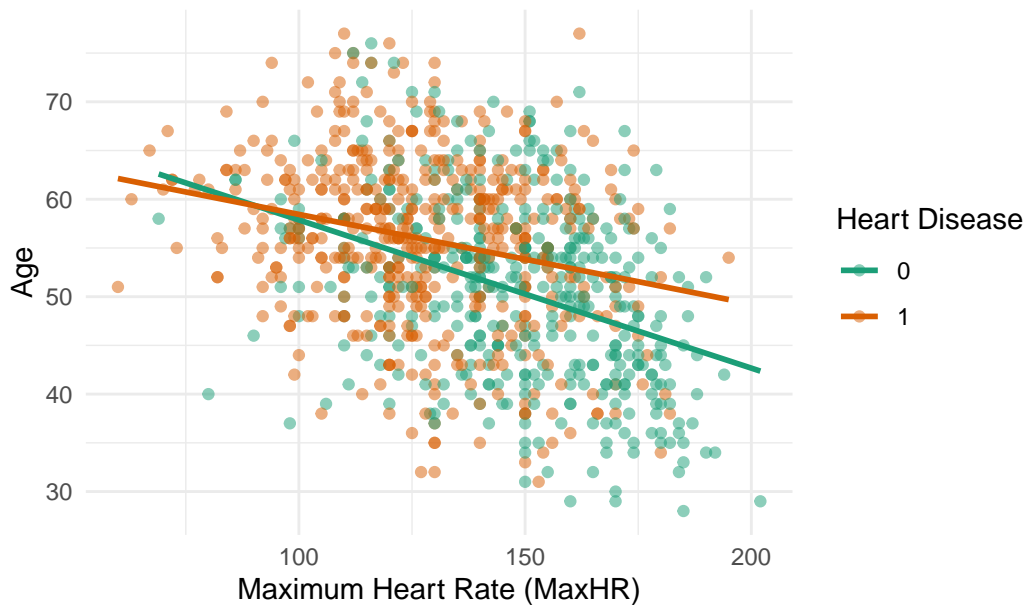
## Task 3: EDA

### Question 1

Here, we hope to make a scatterplot to visualize age vs. a function of heart disease and their max heart rate. We will do this by making the x-axis contain max heart rate, the y-axis containing age, and the points being colored by whether or not the person has heart disease with a colorblind friendly pallet. Then, the trend lines will be put on the plot as well.

```
library(ggplot2)
library(scales)

ggplot(new_heart, aes(x = MaxHR, y = Age, color = HD)) +
  geom_point(alpha = 0.5) +
  geom_smooth(method = "lm", se = FALSE) +
  scale_color_brewer(palette = "Dark2") +
  labs(
    title = "Age vs. Max Heart Rate by Heart Disease Presence",
    x = "Maximum Heart Rate (MaxHR)",
    y = "Age",
    color = "Heart Disease"
  ) +
  theme_minimal()
```

```
`geom_smooth()` using formula = 'y ~ x'
```

## Age vs. Max Heart Rate by Heart Disease Presence



**Question 2**

Based on this plot, I think an interaction plot would be better than an additive model. That is because the lines are not parallel (or almost parallel). That means that the slope changes with the heart disease status, which points to an interaction model being a more appropriate fit.

### Task 4: Testing and Training

Now, we will split our `new_heart` into an 80-20 spli for training and testing. This will be done below (with a seed of 101 set):

```
set.seed(101)

heart_split <- initial_split(new_heart, prop = 0.8)

train <- training(heart_split)
test <- testing(heart_split)
```