

XP Booster

April 27

For this challenge you will be working in GROUPS OF TWO and using the heart disease dataset (heart.csv). Be sure to clearly write all the code you use.

1. How many observations and variables are in the dataset?
2. What proportion of people have heart disease in the dataset?
3. Out of all males, what proportion have heart disease?
4. Out of all females, what proportion have heart disease?
5. If you randomly select a male and a female, which person has a higher probability of having heart disease?
6. Look at association between heart disease and gender using a relative boxplot.

Solution:

```
1 ### Load Packages and dataset
2 library(ggplot2)
3 library(tidyverse)
4 heart <- read.csv('/data/datasets/heart.csv')
5
6 ### Question 1: How many observations in the dataset, how many variables?
7 glimpse(heart)
8
9 ### Question 2: What is the proportion of people with heart disease
10 sum(heart$HeartDisease=="Yes")/length(heart$HeartDisease) # 0.08559546
11
12 ### Question 3: Proportion of males with heart disease
13 sum(heart$HeartDisease=="Yes" & heart$Gender=="Male")/length(heart$
14   HeartDisease) # 0.05046671
15
16 ### Question 4: Proportion of females with heart disease
17 sum(heart$HeartDisease=="Yes" & heart$Gender=="Female")/length(heart$
18   HeartDisease) # 0.03512875
19
20 ### Question 5: The male proportion appears to be higher.
21
22 ### Question 6: Heart disease Across Gender with relative bar chart
23 heart %>% ggplot(aes(x=Gender, fill=HeartDisease)) +
24   geom_bar(position='fill') +
25   theme_classic()
```