

# STAT 311Lab 2

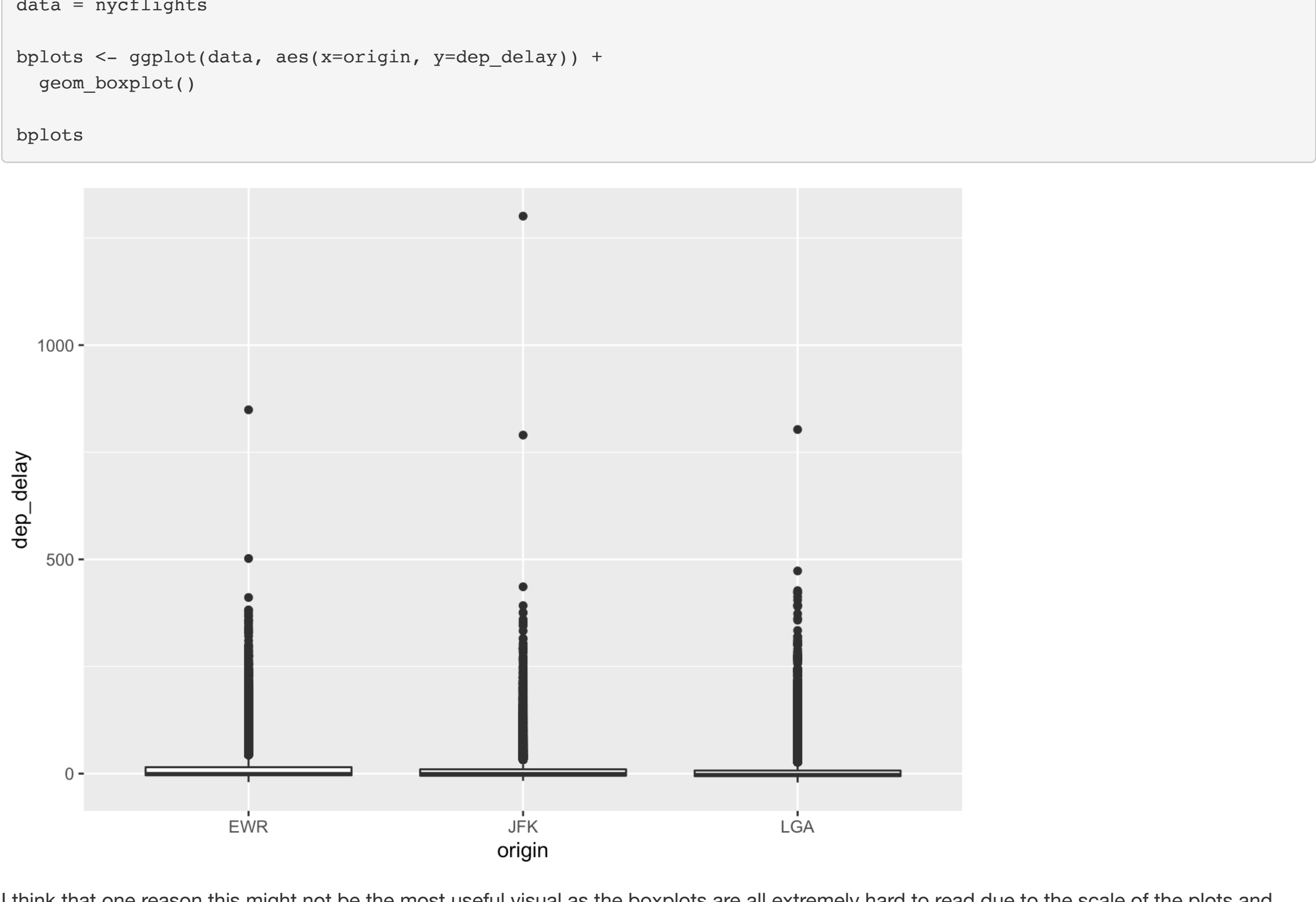
## Group 17

Steven Tran & Owen Cheung

Tuesday, July 26, 2022 @ 11:59pm

## Exercise 1: R Practice

### Part a)



I think that one reason this might not be the most useful visual as the boxplots are all extremely hard to read due to the scale of the plots and because the data points are all clustered together. You really can't the difference between all 3 nor the distributions.

### Part b)

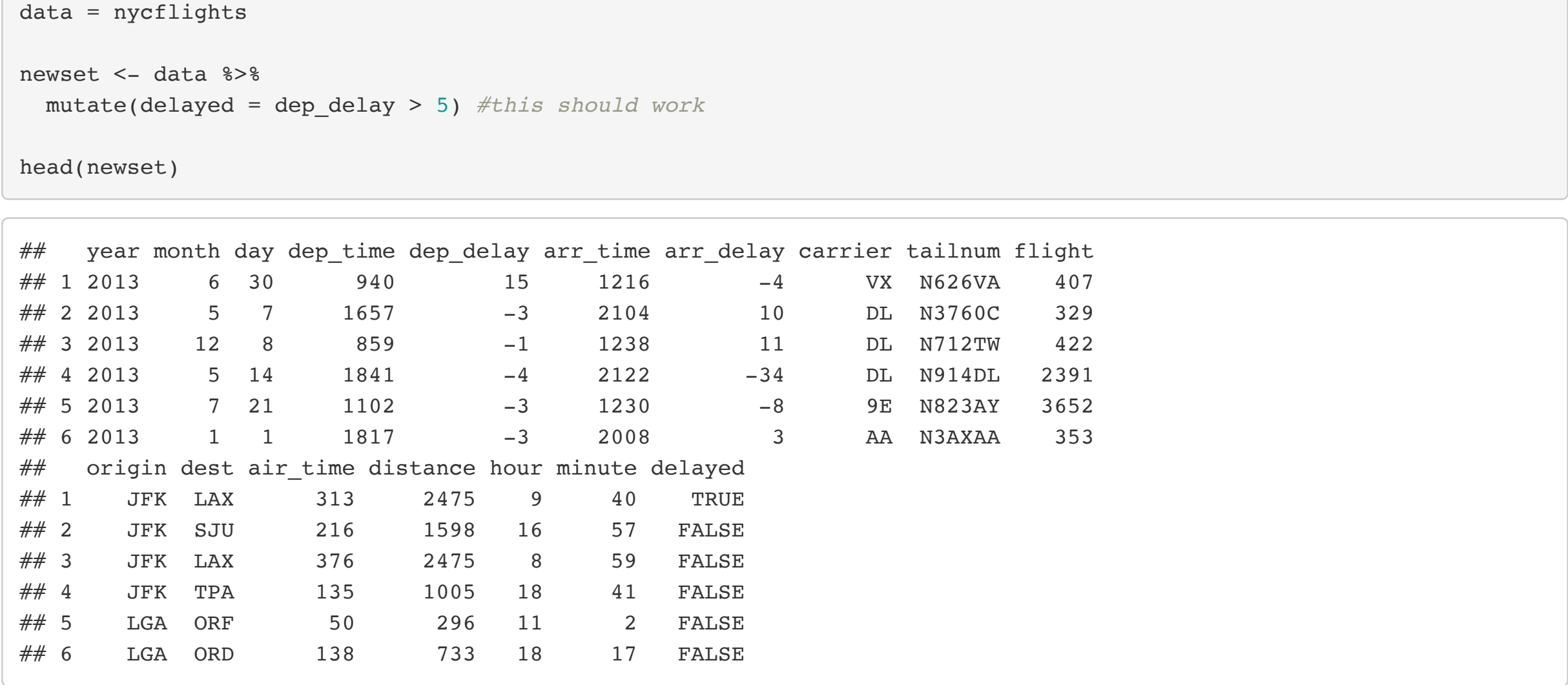


### Part c)

One explanation for the large difference between the means and the medians would be that there are a lot of negative departure delays as well as positive departure delays in the hundreds so if you took a mean/average of that group you would be more likely to get a positive number leaning more towards the positive departure delays while if you took a median of the group its reasonable to get a negative number since there are a large quantity of negative departure delays in the group. I think that I would want to use mean since it might be more representative of the spread of the data

Another explanation could be that there are a lot of outliers in the data. The big difference between the means and the medians could probably be explained by these outliers. Even in the boxplot, there seems to be a large portion of the data clustered around the 0 departure delay time and some outliers as big as 1000. Thus, with this in mind, it's better to use the median as it is less susceptible to large outliers.

### Part d)



### Part e)

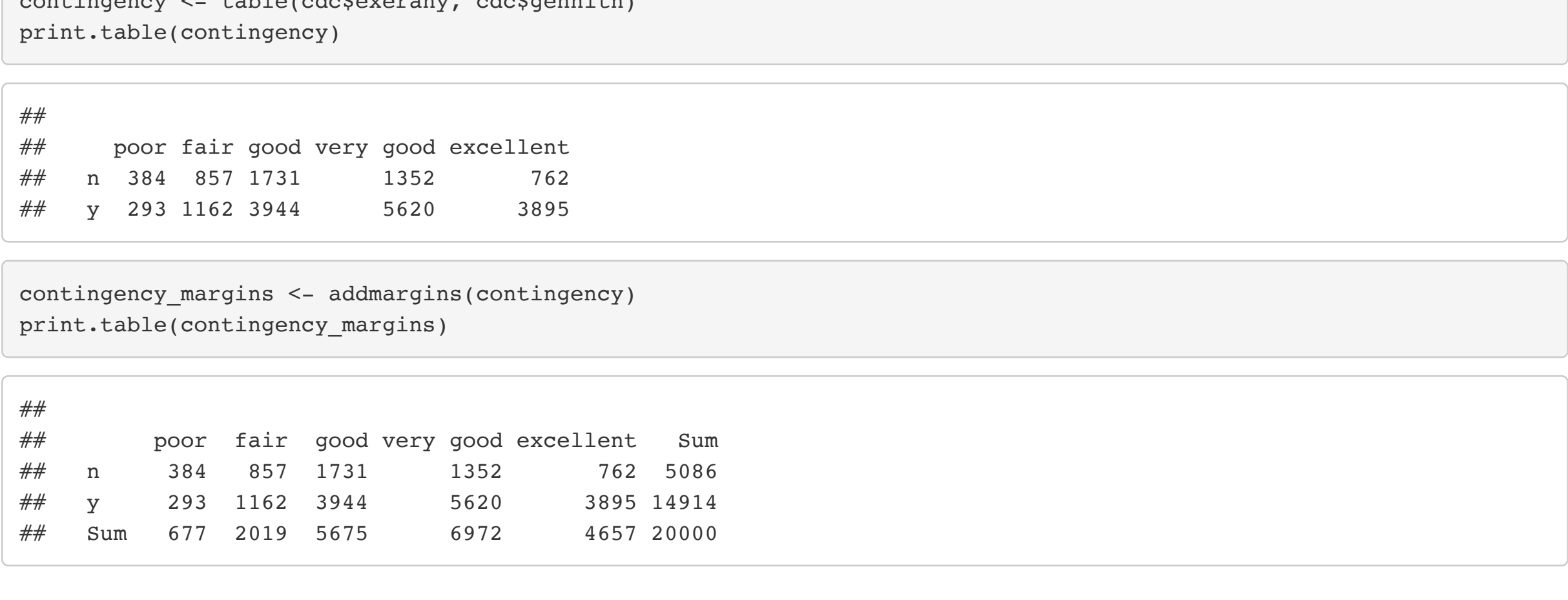


### Part f)

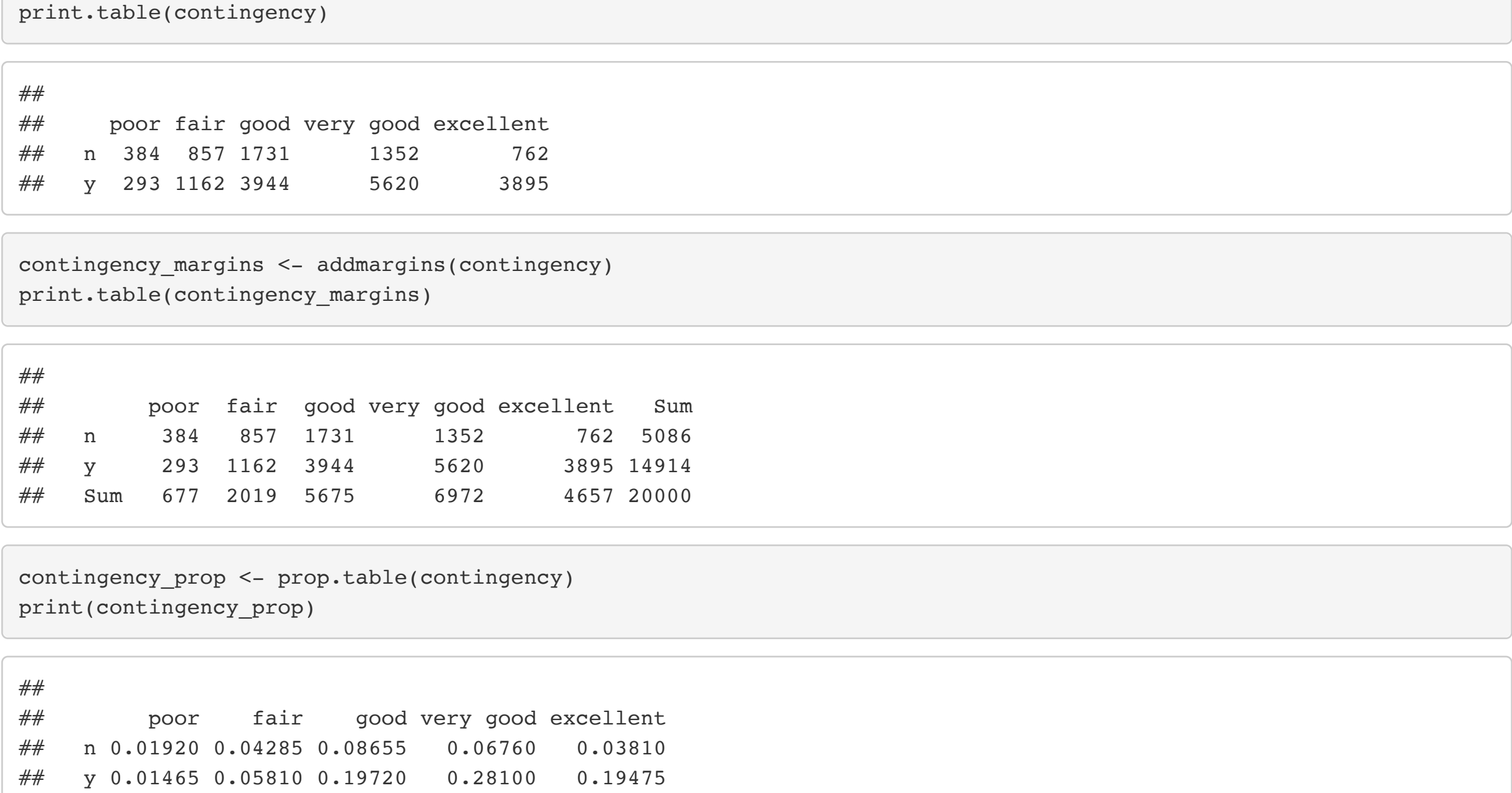
EWR had the highest percentage of delayed flights with 35% compared to JFK with 29% and LGA with 26%

## Exercise 2: Exercise and General Health

### Part a)

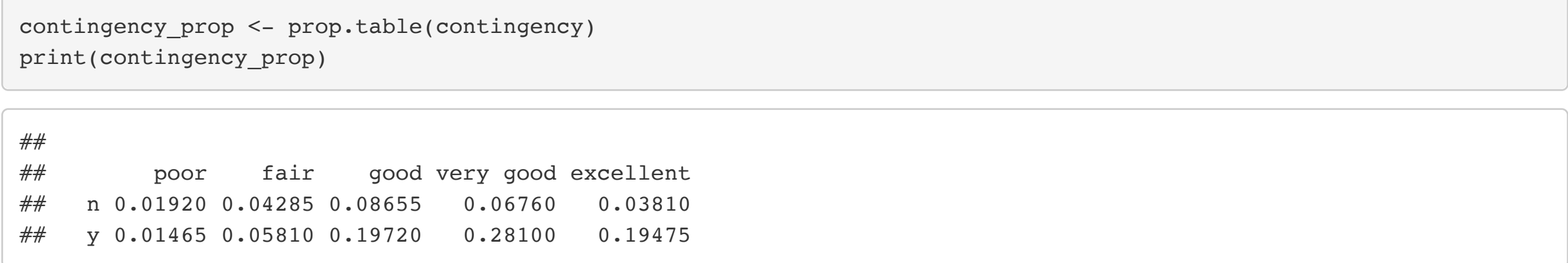


### Part b)



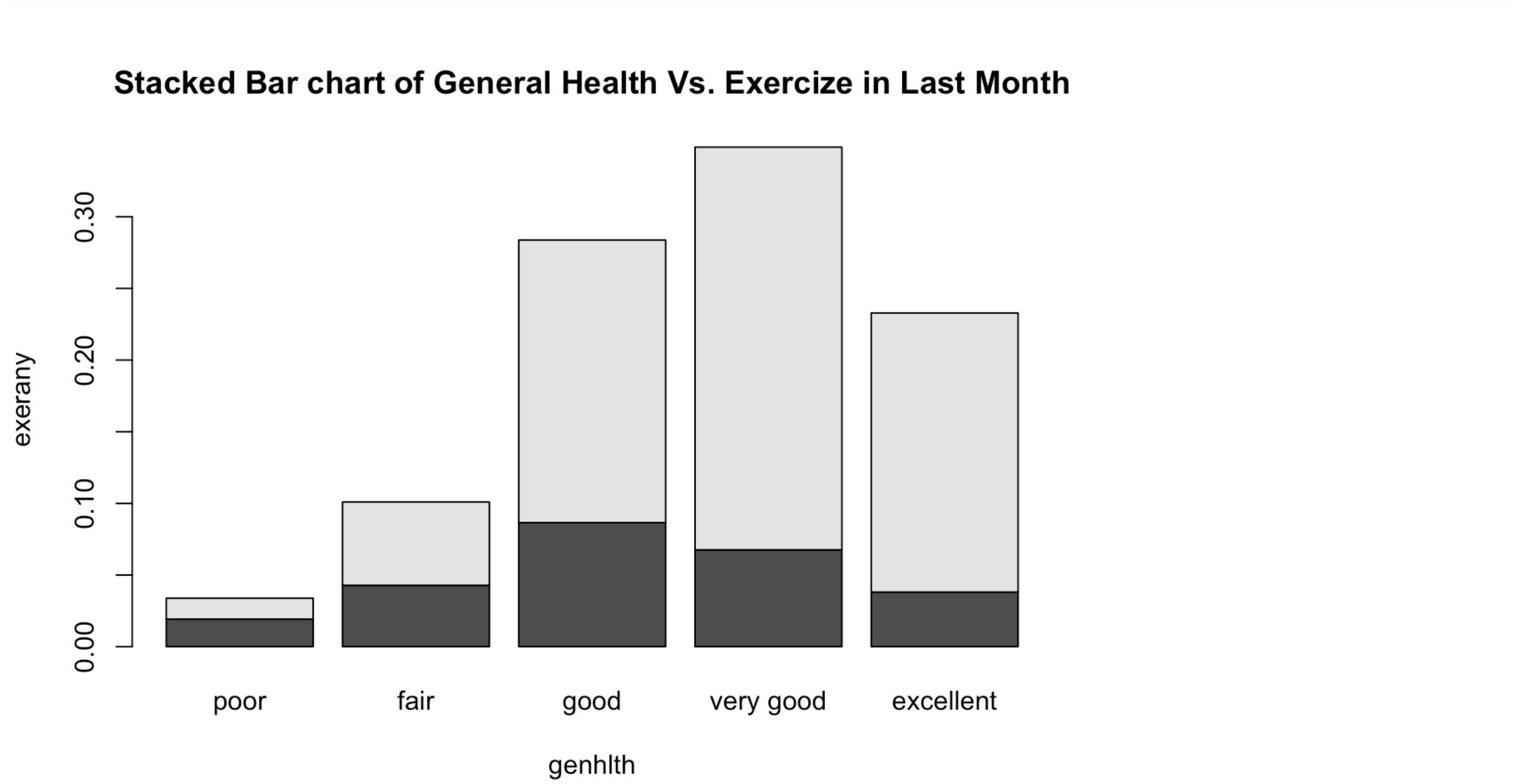
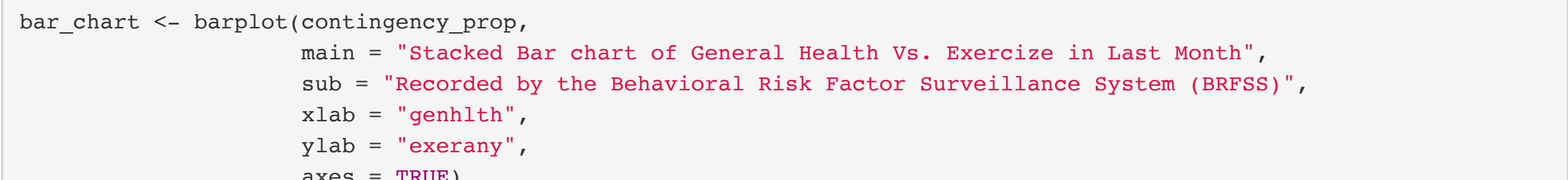
The proportion of those who have exercised in the past month which is the yes to the variable exerany is 14914/20000, which is 0.7457. The proportion of the sample reporting excellent health is 4657/20000 or 0.23285. These numbers are supported by both contingency\_prop and contingency\_margins .

### Part c)



Among the people who exercised in the past month, a proportion of about 0.19475 of them reported excellent health. Among the people who didn't exercise in the past month, only a proportion of 0.03810 of them reported excellent health.

### Part d)



It seems as though most people would rate their health as very good, though those who exercise more make up more of the excellent, very good, good, and fair categories. I also see some response bias here where those who responded tended to have better health, which may perhaps mean that they had strong feelings for this topic. There doesn't seem to be a large 'no exercise group'.

### Part e)



Based on the plots, both the stacked bar chart and the mosaic plot, it seems as though those who exercised in the past month rated their general health higher than the population that didn't. It also seems like the people who didn't exercise may have been subjected to wording bias, maybe they felt bad for having not exercised. There is also response and non-response bias where there is a lot more responses from the people who exercised, and those who didn't exercise perhaps wouldn't have worse health than shown here on the graph.

### Part f)

No, it doesn't seem like the two variables exerany and genhealth are independent. The people who said they exercised in the past month tend to indicate a higher level of general health. However, all we can say is that there is a correlation, not a causation, because there may be many other confounding variables that may impact a person's general health.

## Exercise 3: More Research Questions

### Research Question 1

#### Proposed Question

What is the relationship between mother's age and lengths of pregnancy in weeks?

#### Proposed Statistical Method

Since there are 2 numerical variables, we can use a histogram to analyze the data. I chose this method because we feel as though it works the best with 2 numerical variables. Also, we wanted to know if the mother's age had an influence in the development of the child, or if older mothers would produce more premie babies.

We could also use a scatter plot with x being mother's age and y being lengths of pregnancy. It just depends on which method shows the correlation between the 2 variables best.

### Research Question 2

#### Proposed Question

Is there a relationship between the maturity status of the mother and the premie status of the baby?

#### Proposed Statistical Method

Since there are 2 categorical variables, in the maturity status of the mother and the premie status of the baby, we can use a chi-squared test. Using a chi-square test, we can test if there is a correlation between the status of the mother and premie status of the baby, or if it is due to random chance.

### Research Question 3

#### Proposed Question

Is there a trend between the weight gained by the baby measured in pounds and the smoking status of the mother?

#### Proposed Statistical Method

In this question, there is 1 categorical and 1 numerical variable. As such, we can run a t-test to see if the average weight gained by the baby when the mother smokes is greater or less than when the mother doesn't smoke.