14.1

For part 1, I first found which column had the missing values and where the missing values were (column 7). Then, I found the mean of that column for all values that weren't missing. Then i plugged the mean back in for all the values that were missing

```
# Identify missing values
missing_values <- which(data == "?")
num_missing <- colSums(data == "?")
col_with_mv <- names(which(num_missing>0))
print(col_with_mv)

data[,7]

# Calculate mean/mode for column with missing values
replacement_value <- mean(as.numeric(data[,7][data[,7]!="?"], na.rm = TRUE))
# Or use mode function instead of mean function for categorical data

# Replace missing values with the mean/mode value
data_no_mv <- data
data_no_mv <- data
data_no_mv[data_no_mv == "?"] <- replacement_value
which(data_no_mv == "?")</pre>
```

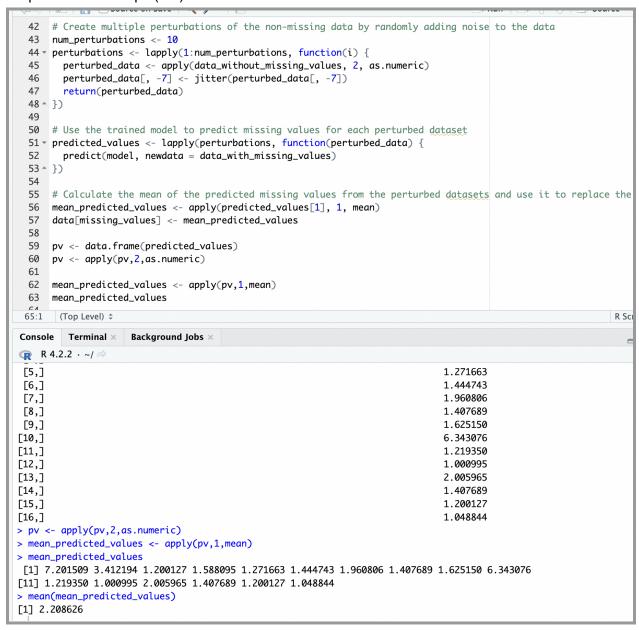
14.2

I first split the data into the rows with complete values and the rows with missing values. I trained a regression model on the rows without missing values. Then, I used that model to predict the values for the dataset with missing values. The result was that all of the predicted values were the same (3.544656). This is likely because the regression model had a poor fit and the coefficient for all of the variables was close to zero. The intercept ended up being the mean of the column with missing values.

```
26
     # Split the dataset into two groups: one with missing values and another without missing values
 27
      data_with_missing_values <- data[which(rowSums(data == "?") > 0), ]
      data_without_missing_values <- data[which(rowSums(data == "?") == 0), ]</pre>
 28
     #train a regression model usining non-missing data to predict missing values
 30
      model \leftarrow lm(V7 \sim V2 + V3 + V4 + V5 + V6 + V8 + V9 + V10 + V11, data = data_without_missing_values
 31
 32
 33 # Use the trained model to predict missing values
 34
      predicted_values <- predict(model, newdata = data_with_missing_values[,-7])</pre>
 35
      predicted_values
 36
 37
     summary(model)
 38
 38:1
      (Top Level) $
                                                                                              R Script $
Console Terminal ×
                    Background Jobs ×
> predicted_values
     24
             41
                      140
                               146
                                        159
                                                 165
                                                          236
                                                                   250
                                                                            276
                                                                                     293
3.544656 3.544656 3.544656 3.544656 3.544656 3.544656 3.544656 3.544656 3.544656
    295
                      316
                               322
                                        412
             298
                                                 618
3.544656 3.544656 3.544656 3.544656 3.544656
> summary(model)
Call:
lm.default(formula = V7 \sim V2 + V3 + V4 + V5 + V6 + V8 + V9 +
   V10 + V11, data = data_without_missing_values)
Residuals:
      Min
                  1Q
                         Median
                                        30
-9.330e-15 -3.250e-15 -1.750e-15 -8.000e-17 1.077e-12
Coefficients:
             Estimate Std. Error
                                    t value Pr(>|t|)
(Intercept) 3.545e+00 6.446e-15 5.499e+14
                                             <2e-16 ***
V2
            8.951e-16 8.244e-16 1.086e+00
                                               0.278
٧3
           -2.165e-16 1.401e-15 -1.550e-01
                                               0.877
٧4
           -4.064e-16 1.363e-15 -2.980e-01
                                               0.766
۷5
           -1.319e-16 8.583e-16 -1.540e-01
                                               0.878
۷6
            4.034e-17 1.149e-15 3.500e-02
                                               0.972
٧8
            8.387e-16 1.109e-15 7.560e-01
                                               0.450
۷9
           -1.188e-16 8.259e-16 -1.440e-01
                                               0.886
            2.160e-17 1.086e-15 2.000e-02
                                               0.984
V10
V11
           -2.654e-15 3.706e-15 -7.160e-01
                                               0.474
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '1
Residual standard error: 4.162e-14 on 673 degrees of freedom
Multiple R-squared: 0.5,
                               Adjusted R-squared: 0.4933
F-statistic: 74.77 on 9 and 673 DF, p-value: < 2.2e-16
```

14.3

For this problem, I first created the 10 perturbations and then wrote a variable that added perturbations to the non-missing data. Then, I used the trained model to predict missing values for each perturbed data set. Finally I calculated the mean, which ended up being different from the previous two steps (2.2).



15.1

In my every day life, I believe optimization would be appropriate to understand the optimal time to do yoga with highly rated yoga instructors, considering i want to do 4-5 classes per week. Some data I would need:

- Instructor rating
- My schedule and obligations

- The location of my work
- The location of my home
- Instructor class schedule
- Location of class

I would want to minimize commute time to get to the desired exercise class while still ensuring I take classes with highly rated instructors and do a minimum of 4-5 classes per week.