Question 9.1

Using the same crime data set uscrime.txt as in Question 8.2, apply Principal Component Analysis and then create a regression model using the first few principal components. Specify your new model in terms of the original variables (not the principal components), and compare its quality to that of your solution to Question 8.2. You can use the R function prcomp for PCA. (Note that to first scale the data, you can include scale. = TRUE to scale as part of the PCA function. Don't forget that, to make a prediction for the new city, you'll need to unscale the coefficients (i.e., do the scaling calculation in reverse)!)

ANSWER:

I first ran a PCA analysis on the data using prcomp

```
myPCA <- prcomp(crime[,1:15], scale = TRUE)
myPCA
summary(myPCA)</pre>
```

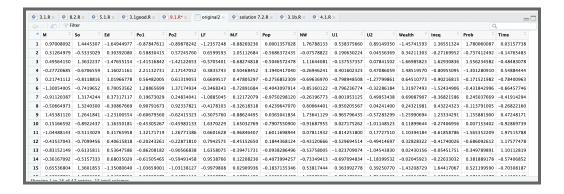
I used the summary of the PCA analysis to determine which principal components were most important in explaining the data. I determined that PC1 - PC9 were useful since it captured 95% of the variance

```
> summary(myPCA)
Importance of components:
                                PC2
                                       PC3
                                               PC4
                                                       PC5
                         PC1
                                                               PC6
Standard deviation
                      2.4534 1.6739 1.4160 1.07806 0.97893 0.74377
Proportion of Variance 0.4013 0.1868 0.1337 0.07748 0.06389 0.03688
Cumulative Proportion 0.4013 0.5880 0.7217 0.79920 0.86308 0.89996
                          PC7
                                  PC8
                                          PC9
                                                 PC10
                                                         PC11
                                                                 PC12
                      0.56729 0.55444 0.48493 0.44708 0.41915 0.35804
Standard deviation
Proportion of Variance 0.02145 0.02049 0.01568 0.01333 0.01171 0.00855
Cumulative Proportion 0.92142 0.94191 0.95759 0.97091 0.98263 0.99117
                         PC13 PC14
                                        PC15
Standard deviation
                      0.26333 0.2418 0.06793
Proportion of Variance 0.00462 0.0039 0.00031
Cumulative Proportion 0.99579 0.9997 1.00000
```

I then used the rotation and x of the PCA model to restate the PCA in terms of the original variables.

```
pcaroation2 <- myPCA$rotation[,1:9]
standardized2 <- myPCA$x[,1:9]
original2 <- t(pcaroation2 %*% t(standardized2))</pre>
```

The output of original2 looks like this:



We now have our original variables and not the principal components. However, this data is still scaled, so I had to use the center output to unscale the data.

```
unscaled_data <- t(t(original2)+ myPCA$center)
```

Now i had my PCA 1-9 unscaled and expressed in the original variables. I was ready to run the linear regression on this data.

There were a few issues with this output, mainly that the summary showed that many of the variables could not be defined because of singularities. I researched this online and learned that it means that the predictor variables are highly correlated (which also seemed to be stated in the office hours).

```
Coefficients: (6 not defined because of singularities)
           Estimate Std. Error t value Pr(>|t|)
(Intercept) -5185.56
                       8634.35 -0.601
                                        0.5518
             156.17
                         67.67 2.308
                                        0.0267 *
So
             231.32
                         92.98
                                2.488
                                        0.0175 *
Ed
            -116.09
                        106.83 -1.087
                                        0.2842
Po1
            3563.54
                       1788.64
                                1.992
                                        0.0538 .
Po2
           -3114.45
                       1751.55 -1.778
                                        0.0836 .
LF
              75.51
                         69.68
                               1.084
                                        0.2855
M.F
              39.65
                         76.14
                                        0.6056
                                0.521
Pop
             -35.34
                         70.98 -0.498
                                        0.6215
                                        0.0458 *
NW
            -284.19
                        137.50 -2.067
U1
                 NA
                            NA
                                   NA
                                            NA
U2
                 NA
                            NA
                                   NA
                                            NA
Wealth
                 NA
                            NA
                                   NA
                                            NA
Ineq
                 NA
                            NA
                                   NA
                                            NA
Prob
                 NA
                            NA
                                   NA
                                            NA
Time
                 NA
                                   NA
                                            NA
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' '1
Residual standard error: 239.3 on 37 degrees of freedom
Multiple R-squared: 0.692,
                              Adjusted R-squared: 0.6171
F-statistic: 9.239 on 9 and 37 DF, p-value: 3.588e-07
```

Ultimately, the output of the model for the test data point made no sense as it was a negative number (-11,000).