

Homework Week 6

Chris Messer
2022-09-30

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!))

In the below analysis, I used principal component analysis to reduce the dimensions of the dataset, and then fed them into a linear regression model. I determined that 7 principal components was the optimal number. The resulting linear model is

$$y = 65.22 \cdot PC1 - 70.08 \cdot PC2 + 25.19 \cdot PC3 + 69.45 \cdot PC4 - 229.04 \cdot PC5 - 60.21 \cdot PC6 + 117.26 \cdot PC7 + 905$$

We can then transform those coefficients back into their original (non principal) coefficients, and scale them back and that results in the equation:

$$y = 55.237 \cdot M + 139.757 \cdot So + -6.804 \cdot Ed + 44.586 \cdot Po1 + 46.424 \cdot Po2 + 673.381 \cdot LF + 44.403 \cdot M.F + 0.96 \cdot Pop + 5.685 \cdot NW + -1027.735 \cdot U1 + 24.416 \cdot U2 + 0.029 \cdot Wealth + 12.451 \cdot Ineq + -5170.569 \cdot Prob + -2.215 \cdot Time - 5498.458$$

The R^2 value of the above approach is .5306. Compare this to my best R^2 value using last weeks pure linear regression model with no PCA, which had an R^2 of .69.

It is reasonable that using PCA to reduce dimensionality did not make a better model than using straight linear regression with only a few predictors. That is because PCA takes into account all of the features, even if only using a handful of principal components in the linear model.

See further analysis below.

Analysis

First, let's load in the data.

```
data <- read.csv('uscrime.txt', sep = '\t')
```

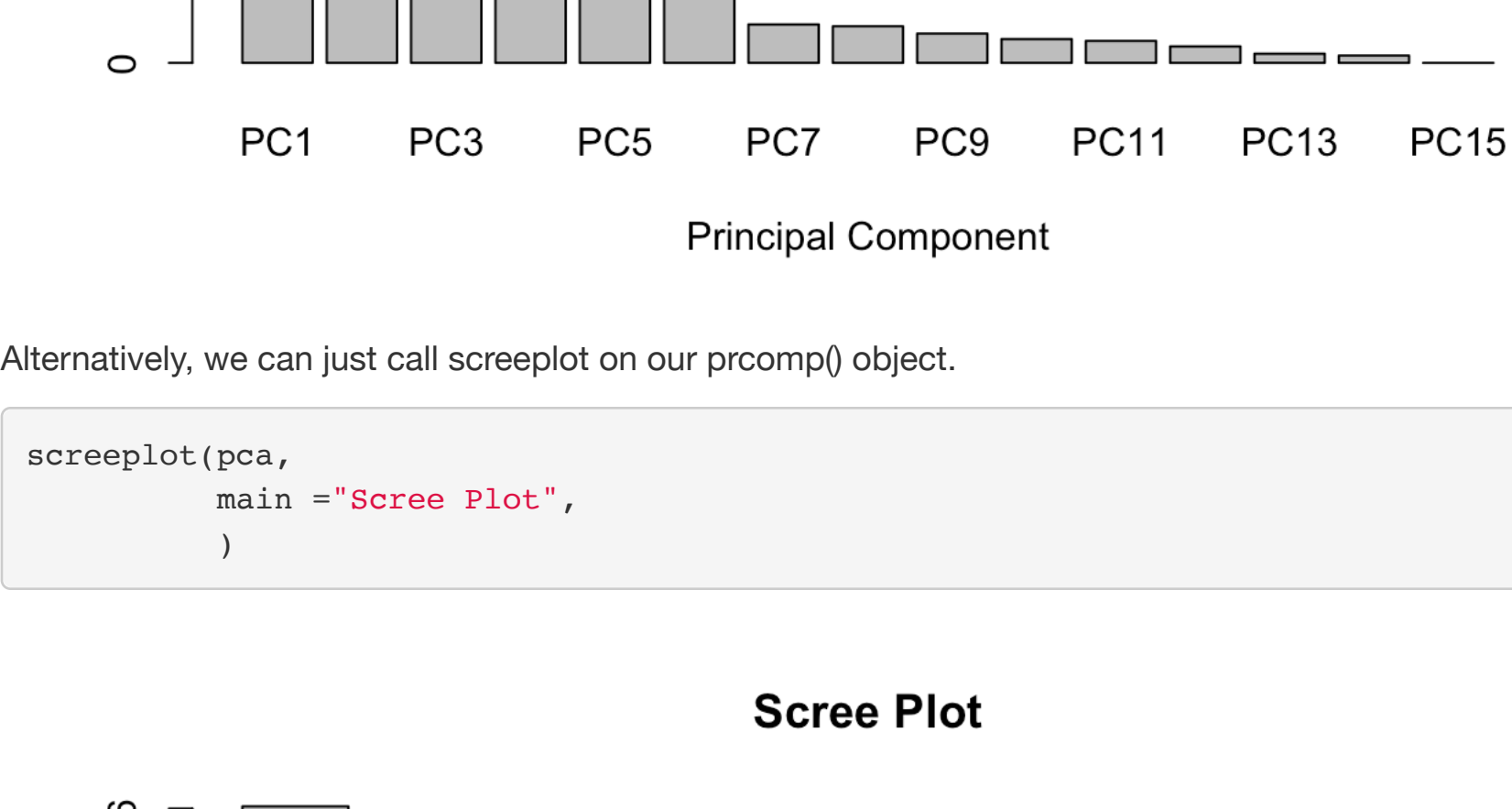
Now, let's call prcomp() to perform principal component analysis.

```
pca <- prcomp(~, data[,1:15], scale = TRUE)
```

```
#These are the eigenvalues, we can manually calculate the percentage variance
pca.var <- pca$sdev^2
pca.var.per <- round(pca.var/sum(pca.var)*100,1)

#or we can get it from prcomp
pca.var.per <- round(summary(pca)$importance[2,]*100,1)

#now plot it out
barplot(pca.var.per,
        main = "Scree Plot",
        xlab = "Principal component",
        ylab = "Percent Variation")
```



Alternatively, we can just call screeplot on our prcomp() object.

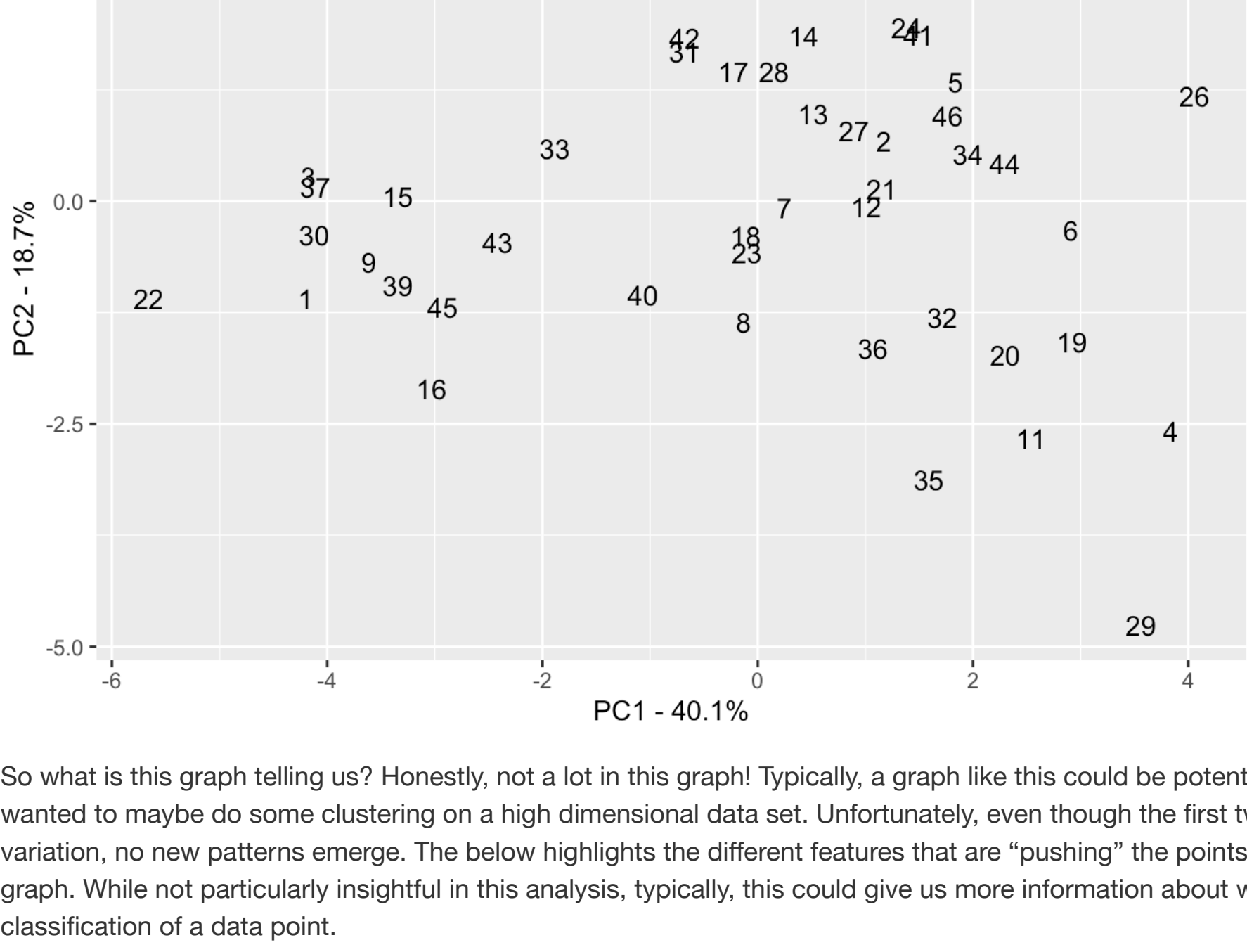
```
screeplot(pca,
          main = "Scree Plot",
          )
```



Now, I'd like to dig in a little further on the first two principal components, as the two of them make up 60% of the variance in the predicted crime rate.

```
library(ggplot2)
pca.data <- data.frame(X=pca$x[,1],
                      Y=pca$x[,2])

ggplot(data = pca.data, aes(x=X, y=Y, label=(1:47))) +
  geom_text() +
  xlab(paste("PC1 - ", pca.var.per[1], "%", sep="")) +
  ylab(paste("PC2 - ", pca.var.per[2], "%", sep="")) +
  ggtitle("PCA Graph")
```



So what is this graph telling us? Honestly, not a lot in this graph! Typically, a graph like this could be potentially helpful in a scenario where we wanted to maybe do some clustering on a high dimensional data set. Unfortunately, even though the first two PC's account for 60% of the variation, no new patterns emerge. The below highlights the different features that are "pushing" the points above to the left or right on the above graph. While not particularly insightful in this analysis, typically, this could give us more information about which features are driving the classification of a data point.

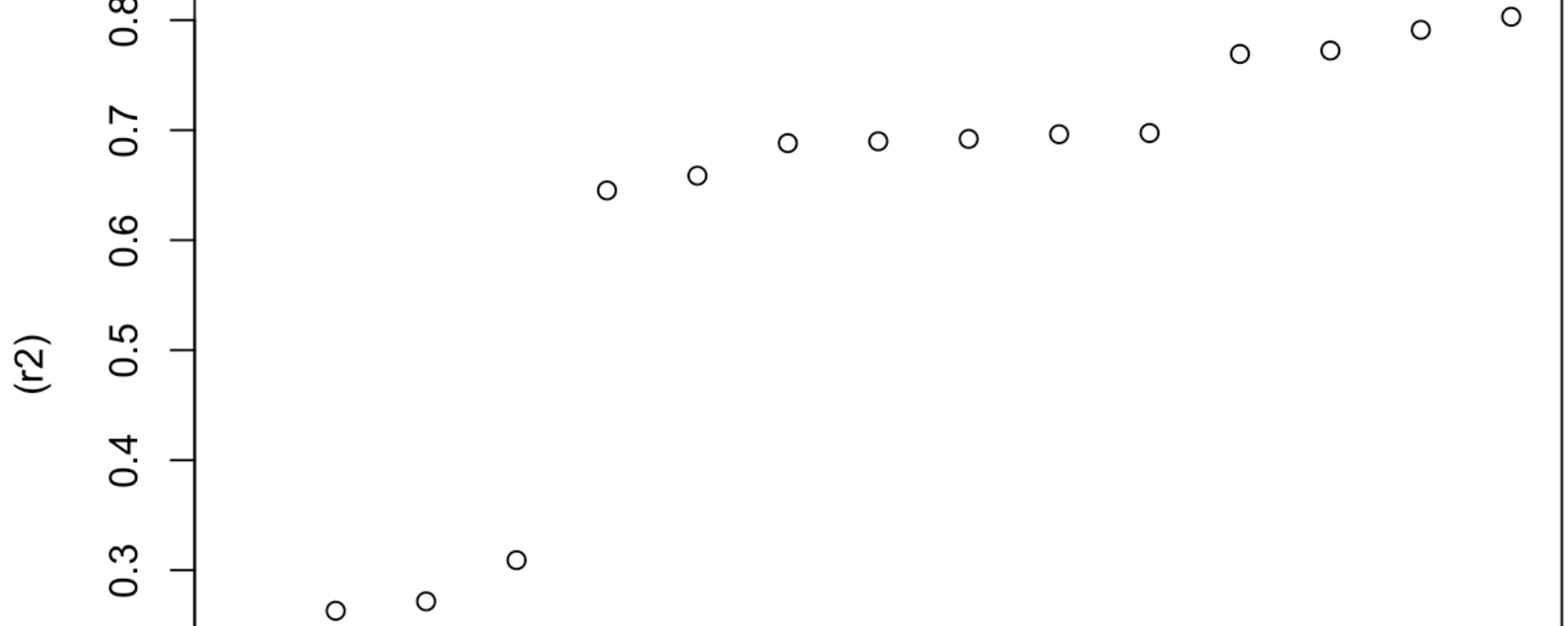
```
loading_scores <- abs(loading_scores)
feature_score_ranked <- sort(feature_score[,1], decreasing = T)
top <- names(feature_score_ranked)
pca$rotation[top,1]
```

```
##      Wealth      Ineq      Ed      So      Po2      Po1
## 0.37970331 -0.36579778 0.33962148 -0.33088129 0.31099285 0.30863412
##      M      NW      Prob      LF      M.F      Pop
## -0.30371194 -0.29358647 -0.25888661 0.17617757 0.11638221 0.11397836
##      U1      U2      Ineq
## 0.04050137 -0.02062867 0.01812228
```

Now, let's go ahead and build a linear model, and loop over all of the PC combinations and look at the r^2 for each.

```
r2 <- c()
for (i in 1:15){
  linear_data <- as.data.frame(cbind(pca$x[,1:i],Crime = data[,16]))
  pca.model <- lm(Crime~., data = linear_data)
  r2[i] <- summary(pca.model)$r.squared
}

plot(r2)
```



Looking at our Scree Plot from above, it looks our optimal number of principal components to look at is probably around 4 or 5. So why does our model just get better and better the more PC's we use in the model? Well, if you remember from last week, there was an issue with over fitting when using all 15 features to predict. And when using all 15 features of PCA in a linear regression model, we are effectively using all 15 features from last week's analysis! In fact, we get the exact same r^2 values.

So, a better measure would be to cross validate the linear regression model that uses the principal components.

Run a cross validated linear regression model for all combinations of principal components

```
library("DAAG")

nComp = 15
r2cross <- c()
r2cross2 <- c()
for (Comp in seq(nComp)){

  principal_components <- as.data.frame(cbind(pca$x[,1:Comp],Crime = data[,16]))
  pca.model.lm <- lm(Crime~., data = principal_components)
  pca.model.cv <- cv.lm(data = principal_components, pca.model.lm, m = 5, printit=F, plotit=F)

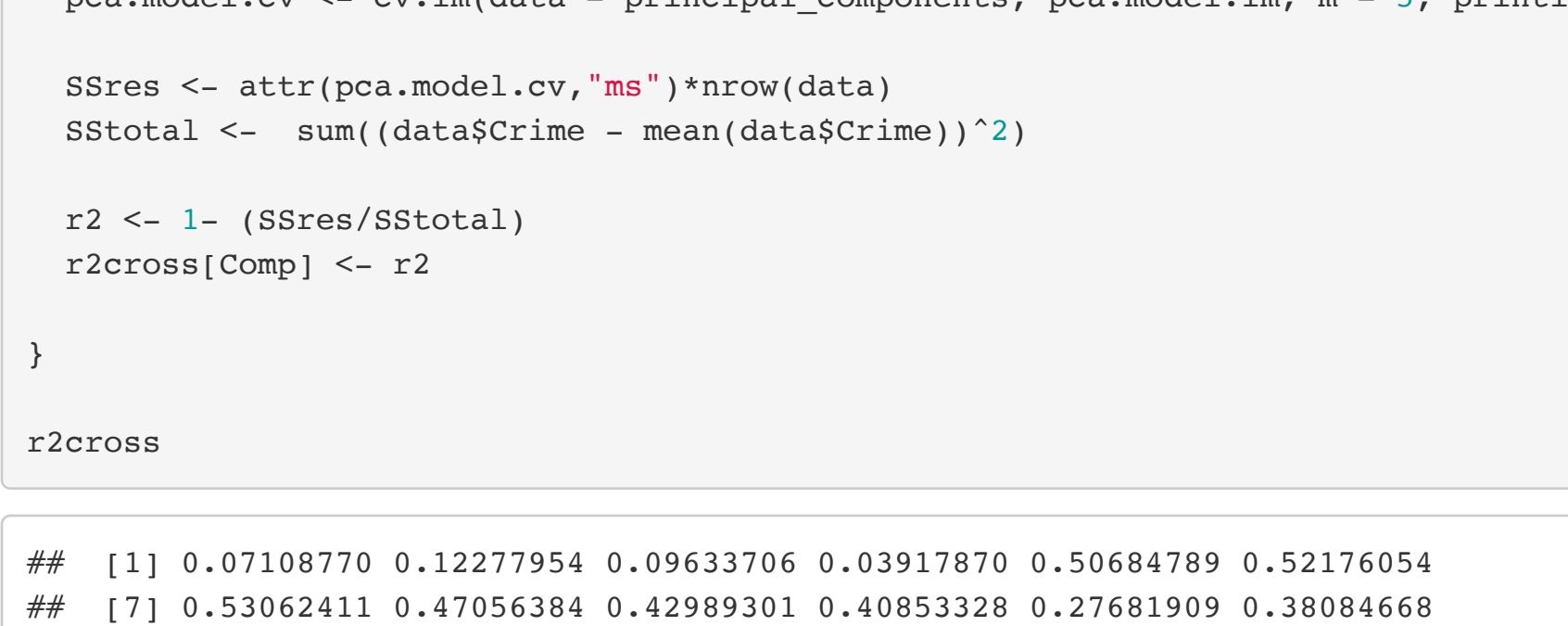
  SRes <- attr(pca.model.cv, "ma")$nrow(data)
  SStotal <- sum((data$Crime - mean(data$Crime))^2)

  r2 <- 1 - (SRes/SStotal)
  r2cross[Comp] <- r2
}

r2cross
```

```
## [1] 0.07108779 0.12277954 0.09633766 0.03917870 0.50684769 0.52176054
## [7] 0.53062411 0.47056384 0.42989301 0.40853328 0.27681909 0.38084660
## [13] 0.34605680 0.41716859 0.41975895
```

```
plot(r2cross)
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

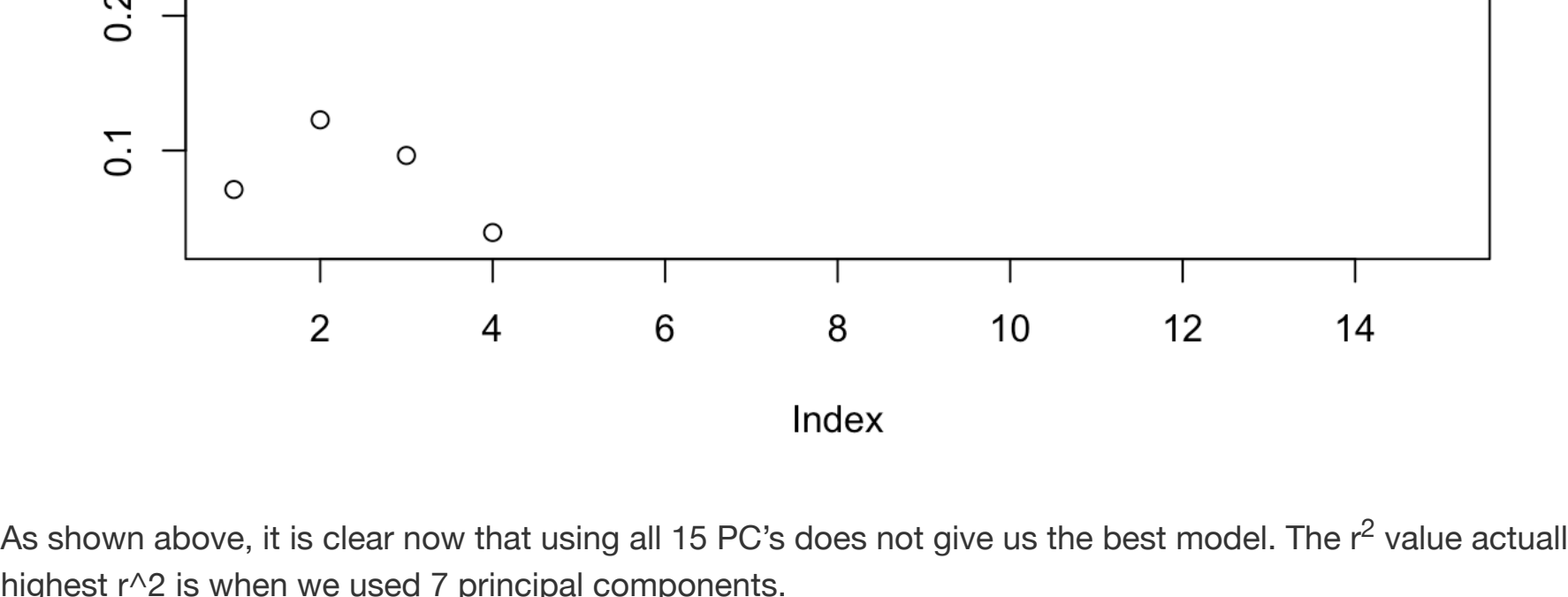
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

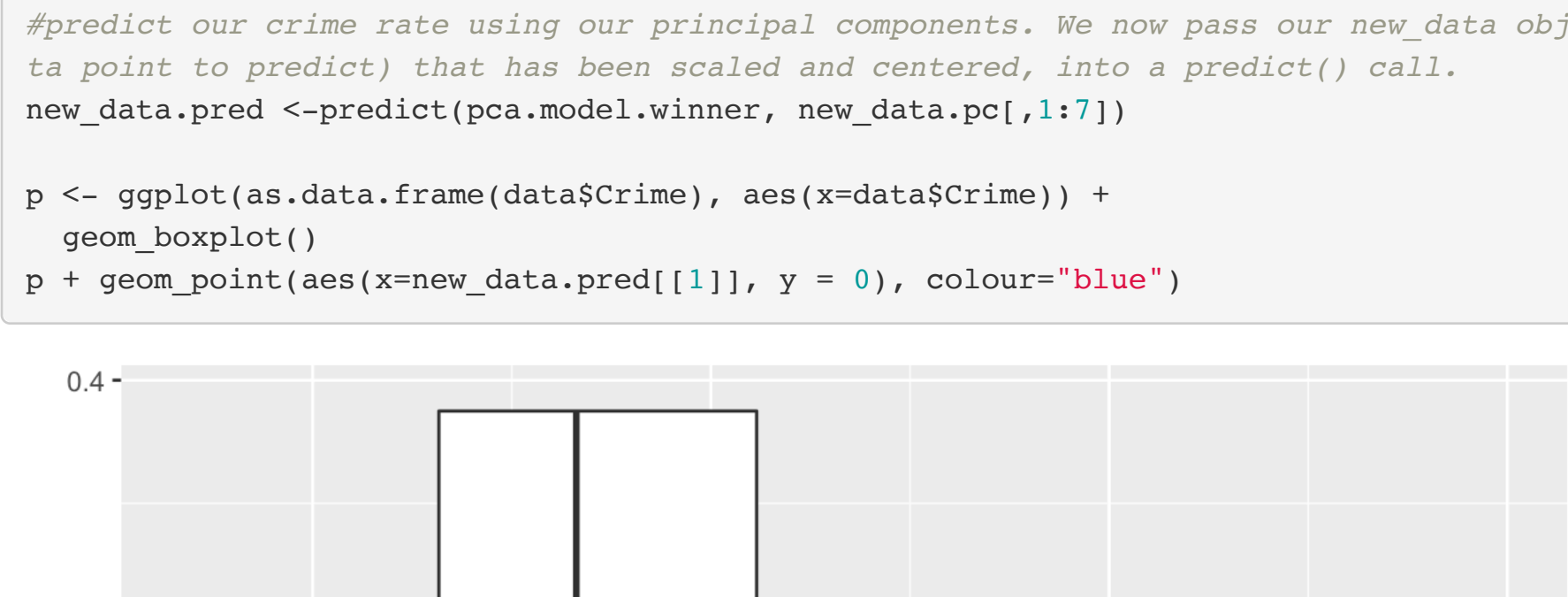
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

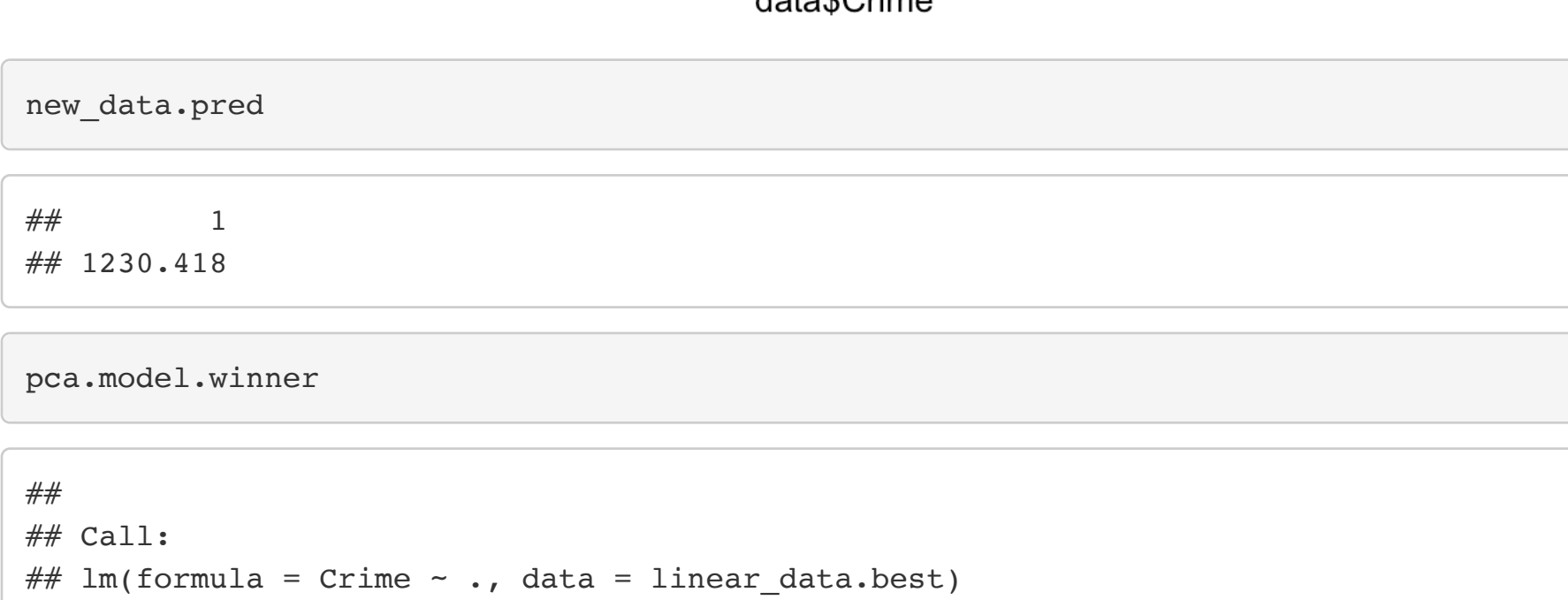
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

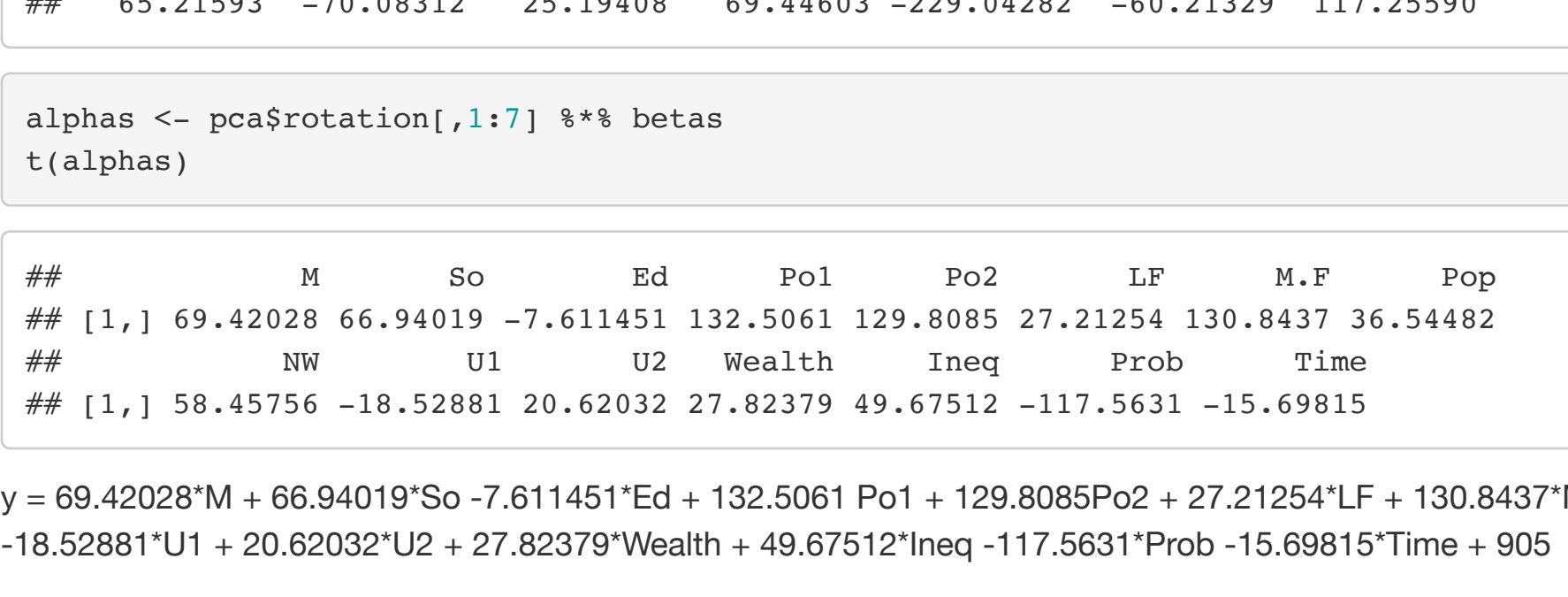
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

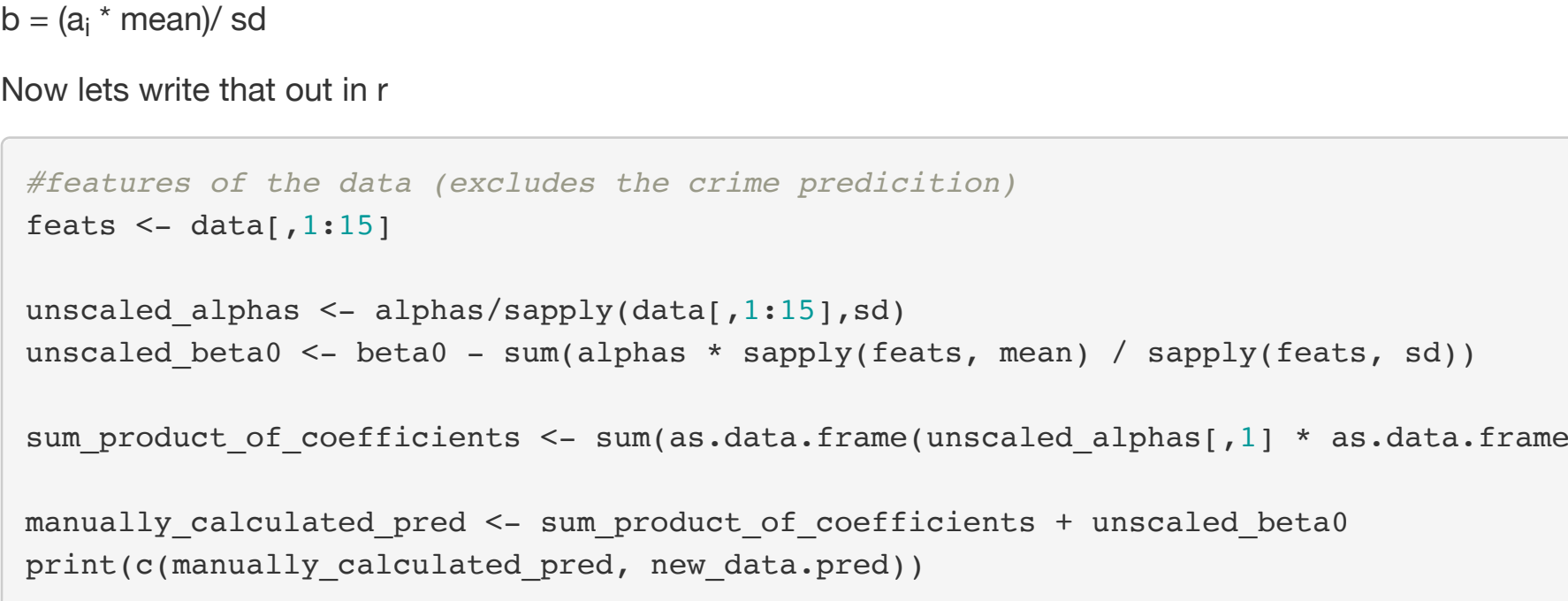
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

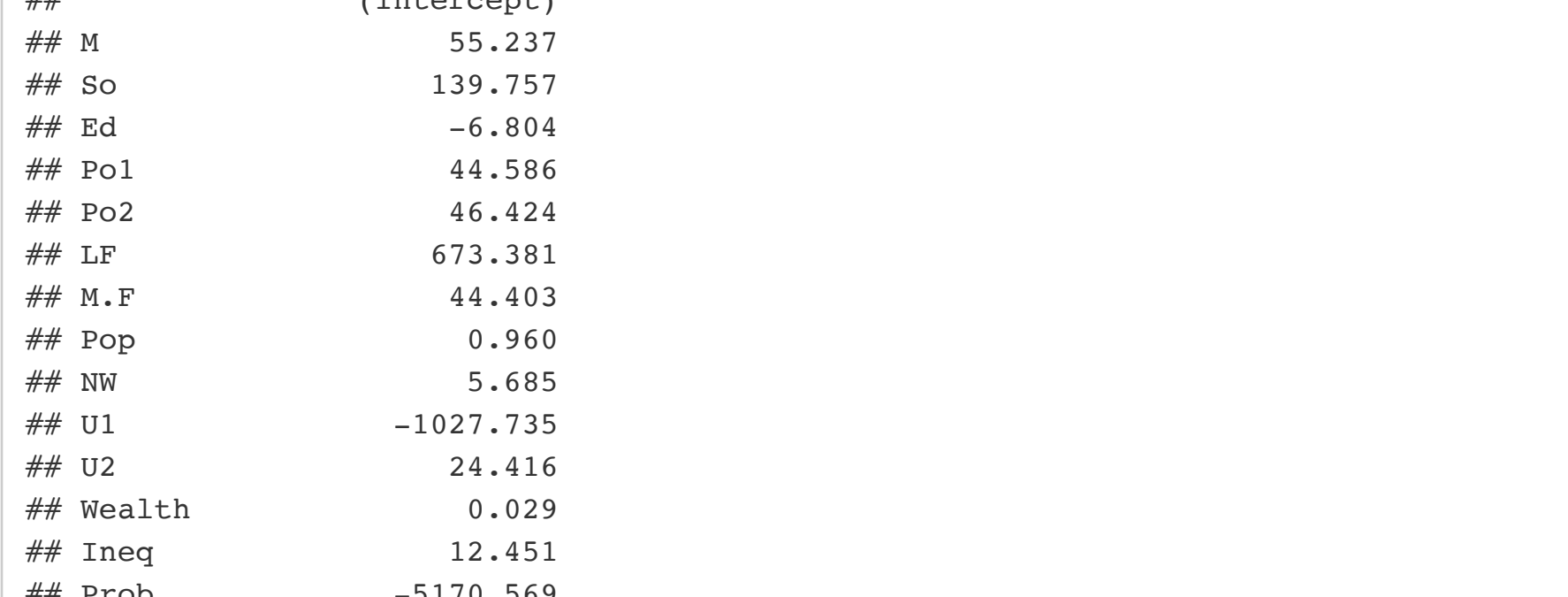
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

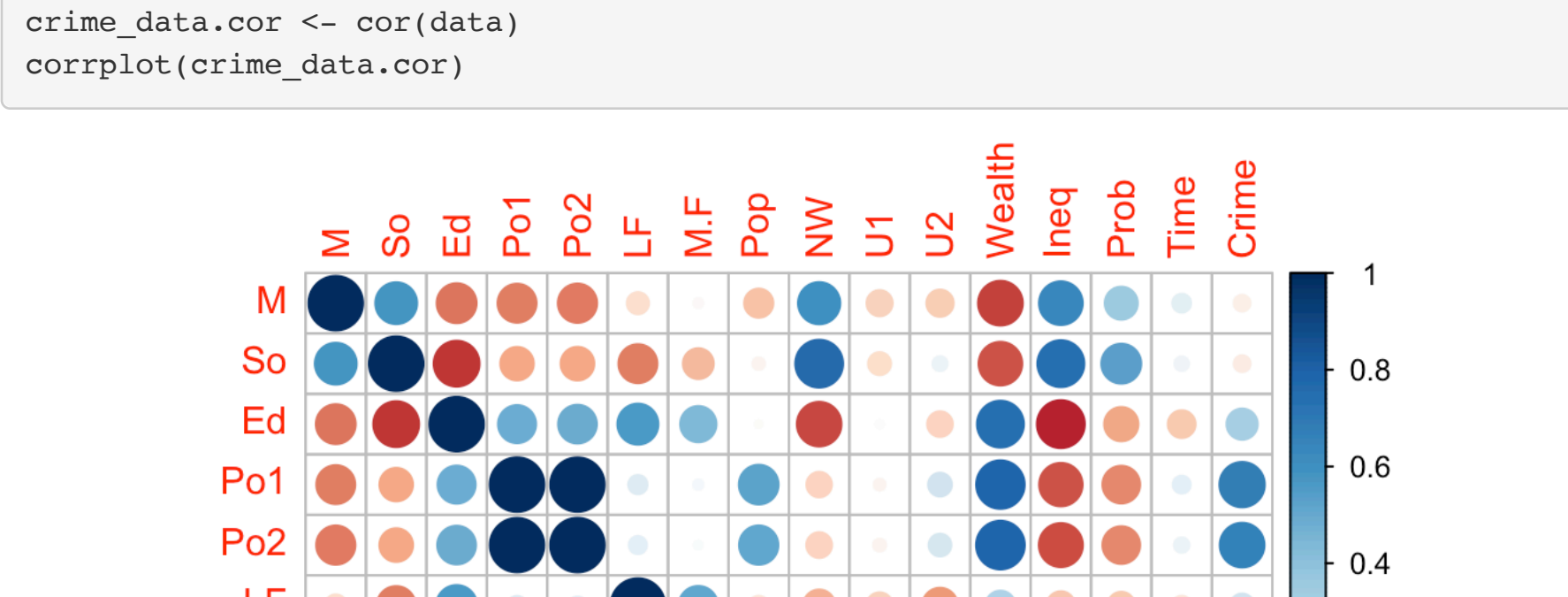
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

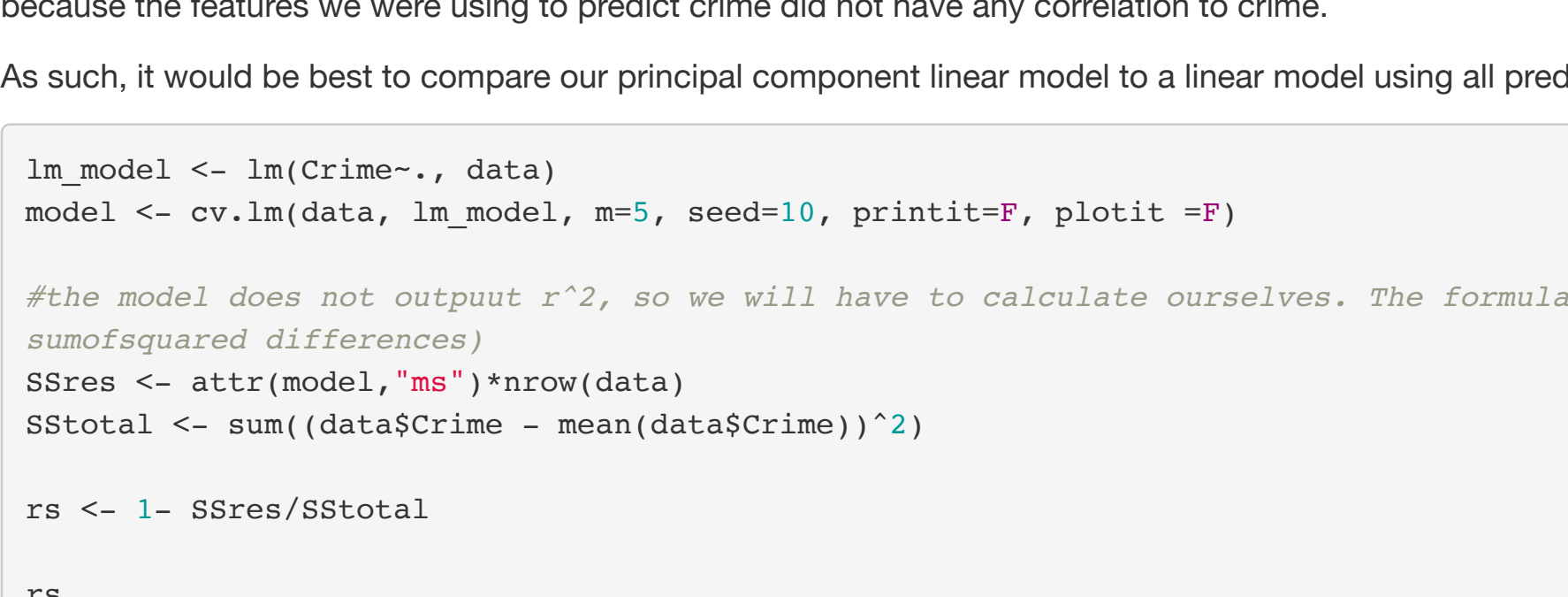
```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```



As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```


As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```


As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```


As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.
new_data.pc <- as.data.frame(predict(pca, new_data))

#Create a new model using the best cross validated model from above
linear_data.best <- as.data.frame(cbind(pca$x[,1:7],Crime = data[,16]))
pca.model.winner <- lm(Crime~., data = linear_data.best)

#Predict our crime rate using our principal components. We now pass our new_data object (which is our original data point to predict) that has been scaled and centered, into a predict() call.
new_data.pred <- predict(pca.model.winner, new_data.pc[,1:7])

p <- ggplot(as.data.frame(data$Crime), aes(x=data$Crime)) +
  geom_boxplot()
p + geom_point(aes(x=new_data.pred[[1]], y = 0), colour="blue")
```


As shown above, it is clear now that using all 15 PC's does not give us the best model. The r^2 value actually went from .5 to around .4, and our highest r^2 is when we used 7 principal components.

Now that we have our best model that uses PCA, we will use make a prediction using our model.

```
#Store the point to predict into a dataframe
new_data <- data.frame(M = 14.0, So = 0, Ed = 10.0, Po1 = 12.0, Po2 = 15.5, LF = 0.640, M.F = 94.0, Pop = 150, NW = 1.1, U1 = 0.120, U2 = 3.6, Wealth = 3200, Ineq = 20.1, Prob = 0.040, Time = 39.0)

#Transform our datapoint into it's principal components, so that we can multiply it by the coefficients from our linear model. Also note, we DO NOT need to scale our data, because when we call predict in the below line of code and pass in our pca object and the data point to convert to principal components, it scales and centers the data for us.

```