

Project PCA

Chelsey

March 30, 2019

Project PCA

Loading in and exploring the data

```
admissionsData <- read.csv("Admission_Predict_Ver1.1.csv", header = TRUE)
#head(admissionsData)
```

With Response Variable Chance.of.Admit

The variable we are interested in predicting, Chance.of.Admit, is the 9th variable.

Run PCA on the data and remove the response variable (chance of admit) and the unique identifier (serial number)

```
set.seed(43849)
pca.admin <- prcomp(as.matrix(admissionsData[,-c(1,9)]), scale = TRUE)
summary(pca.admin)
```

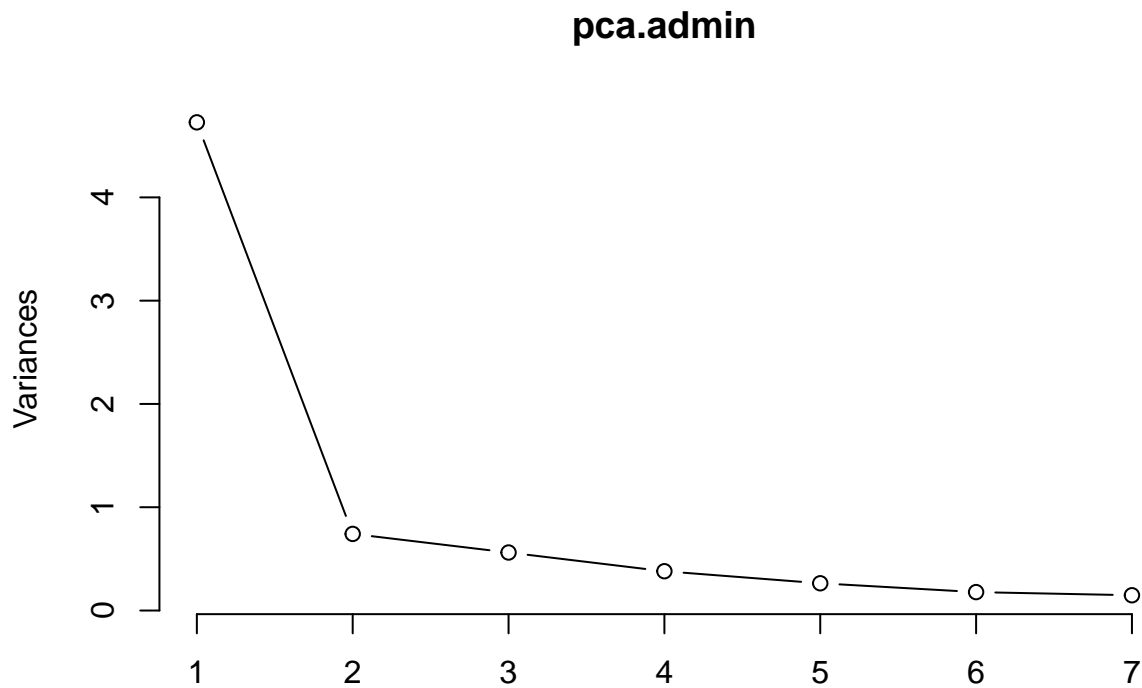
```
## Importance of components:
##              PC1      PC2      PC3      PC4      PC5      PC6
## Standard deviation  2.1740 0.8612 0.74942 0.61674 0.51349 0.42223
## Proportion of Variance 0.6752 0.1060 0.08023 0.05434 0.03767 0.02547
## Cumulative Proportion 0.6752 0.7812 0.86139 0.91573 0.95340 0.97886
##              PC7
## Standard deviation  0.38464
## Proportion of Variance 0.02114
## Cumulative Proportion 1.00000
```

To choose the number of principal components to keep, we can either use the Kaiser criterion, cumulative proportion/percent of variance, or a scree plot.

Using the Kaiser criterion, we keep all principal components with a standard deviation greater than 1 (since the data is scaled). Hence the Kaiser criterion is telling us to keep only the first principal component.

I will now compare this with a scree plot.

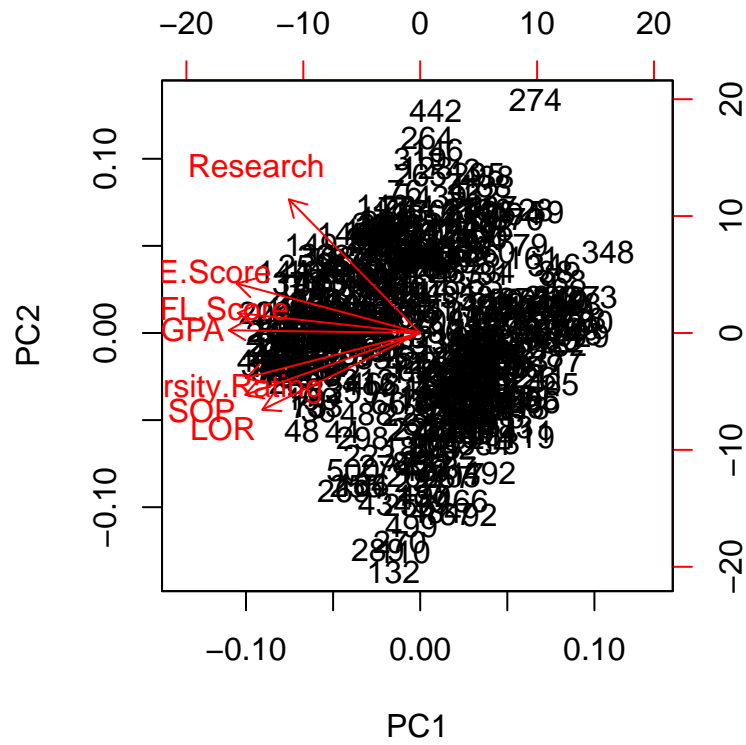
```
plot(pca.admin, type="lines")
```



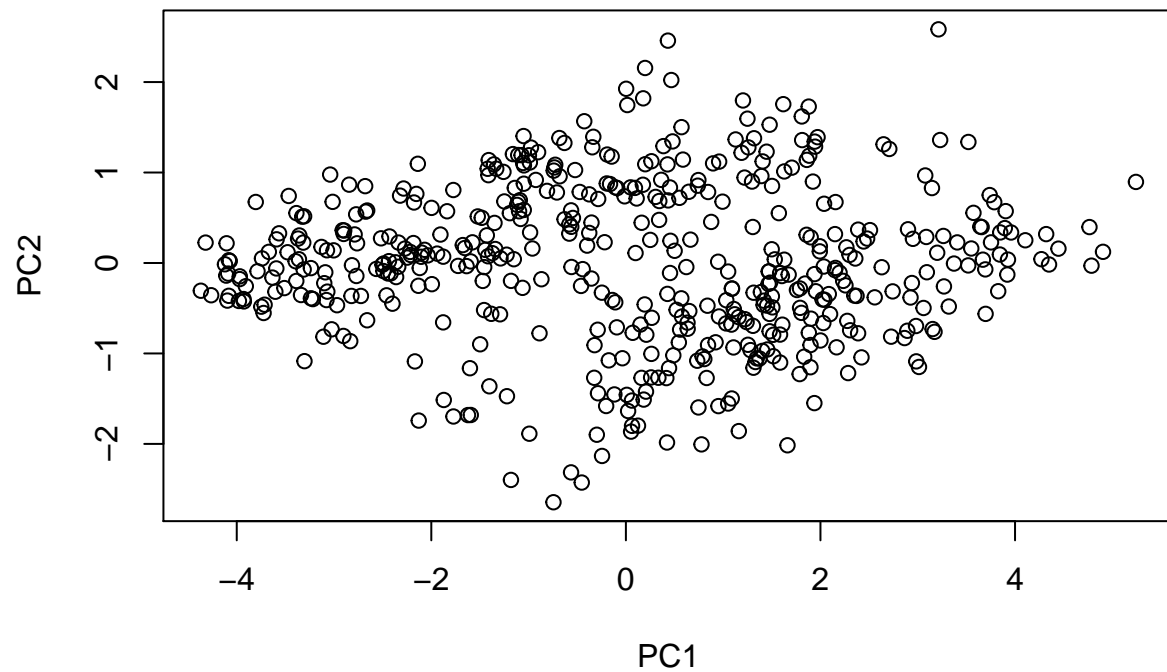
The above scree plot plots the monotonically decreasing eigenvalues and the location of an ‘elbow’ or plateau indicates the number of principal components. The scree plot suggests probably 2 principal components.

The first two principal components that will be retained explain 78% of the variation in the data. We can now view the data projected onto the components using a biplot.

```
biplot(pca.admin)
```

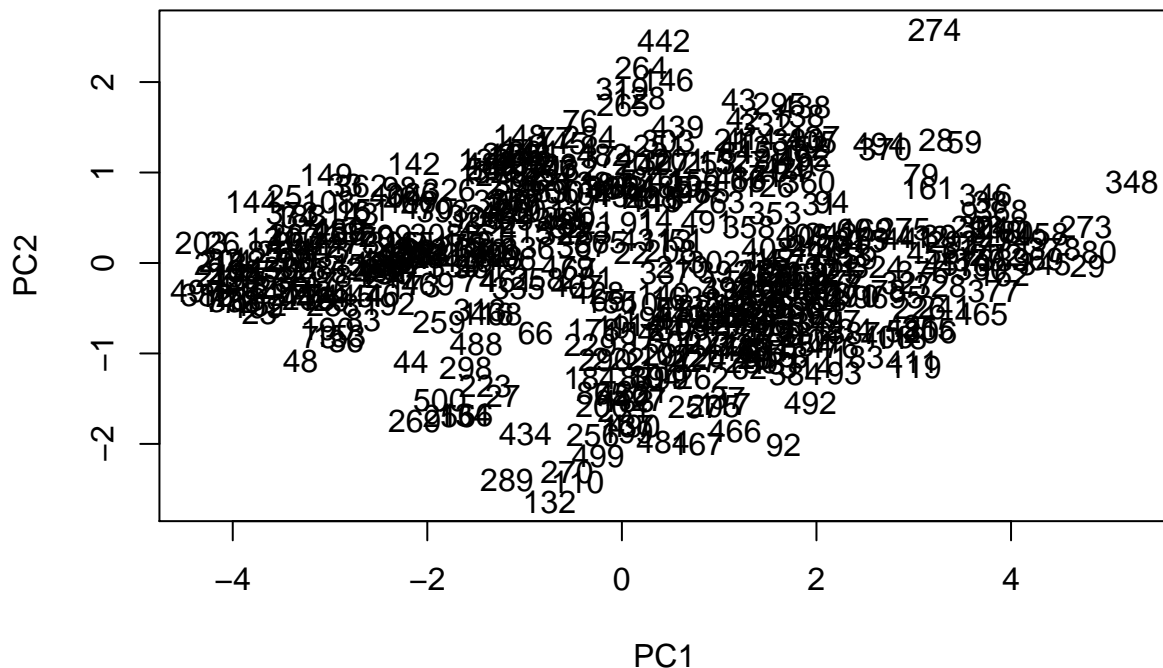


```
plot(pca.admin$x[,1:2])
```



We can put data labels on the biplot by observation number

```
plot(pca.admin$x[,1:2], type = "n")  
text(pca.admin$x[,1:2], labels = 1:nrow(admissionsData))
```



It looks like there are two groups in the above principal component plots.

Take a look at the component loadings (eigenvectors) which provide the coefficients of the original variables, rounded to 2 decimal places.

```
round(pca.admin$rotation[,1:2], 2)
```

```
##           PC1    PC2
## GRE.Score   -0.40  0.27
## TOEFL.Score -0.40  0.11
## University.Rating -0.38 -0.25
## SOP         -0.38 -0.34
## LOR         -0.35 -0.43
## CGPA        -0.42  0.02
## Research    -0.29  0.74
```

These are the coefficients of the original variables. The magnitudes are pretty similar for the first component, perhaps with the exception of research. They are also all containing the same sign. This is a little difficult to interpret, but most likely indicates that the first principal component is equally weighting all predictor variables, with the exception of research.

In the second component, the highest magnitude is the research aspect, along with the letter of recommendation. Perhaps this component indicates previous experience a student has. A reference letter most likely comes from someone you have worked with, conducted research with, volunteered with, or TA'd for. Therefore a good reference letter coupled with research experience could be indicative of research and other activities in both academic and non-academic settings.

We can now look at the four students who scored highest on PC1:

```
admissionsData[order(pca.admin$x[,1], decreasing = TRUE)[1:4], 1:9]
```

```
##      Serial.No. GRE.Score TOEFL.Score University.Rating SOP LOR CGPA
## 348          348      299          94              1 1.0 1.0 7.34
## 80           80      294          93              1 1.5 2.0 7.36
## 29           29      295          93              1 2.0 2.0 7.20
## 273          273      294          95              1 1.5 1.5 7.64
##      Research Chance.of.Admit
## 348          0              0.42
## 80           0              0.46
## 29           0              0.46
## 273          0              0.49
```

It is noted that the four students who performed highest on PC1 all had a low belief of their chance of admit. None of them had research, and all had a similar cumulative GPA. In addition, the universities where all rated low (1 to be exact) and the students had similar GRE and TOEFL scores (well below the average). These students in general seem to be ones who are not performing scoring very well across all predictors.

And the four students who scored highest on PC2:

```
admissionsData[order(pca.admin$x[,2], decreasing = TRUE)[1:4], 1:9]
```

```
##      Serial.No. GRE.Score TOEFL.Score University.Rating SOP LOR CGPA
## 274          274      312          99              1 1.0 1.5 8.01
## 442          442      332          112             1 1.5 3.0 8.66
## 264          264      324          111             3 2.5 1.5 8.79
## 146          146      320          113             2 2.0 2.5 8.64
##      Research Chance.of.Admit
## 274          1              0.52
## 442          1              0.79
## 264          1              0.70
## 146          1              0.81
```

Notice that the four students who performed highest on PC2 all have research experience. In general, these students are scoring better than the students in principal component 1 across the board.

With Response Variable Research

The variable we are interested in predicting, Chance.of.Admit, is the 8th variable.

Run PCA on the data and remove the response variable (research) and the unique identifier (serial number)

```
set.seed(43849)
pca.admin2 <- prcomp(as.matrix(admissionsData[, -c(1,8)]), scale = TRUE)
summary(pca.admin2)
```

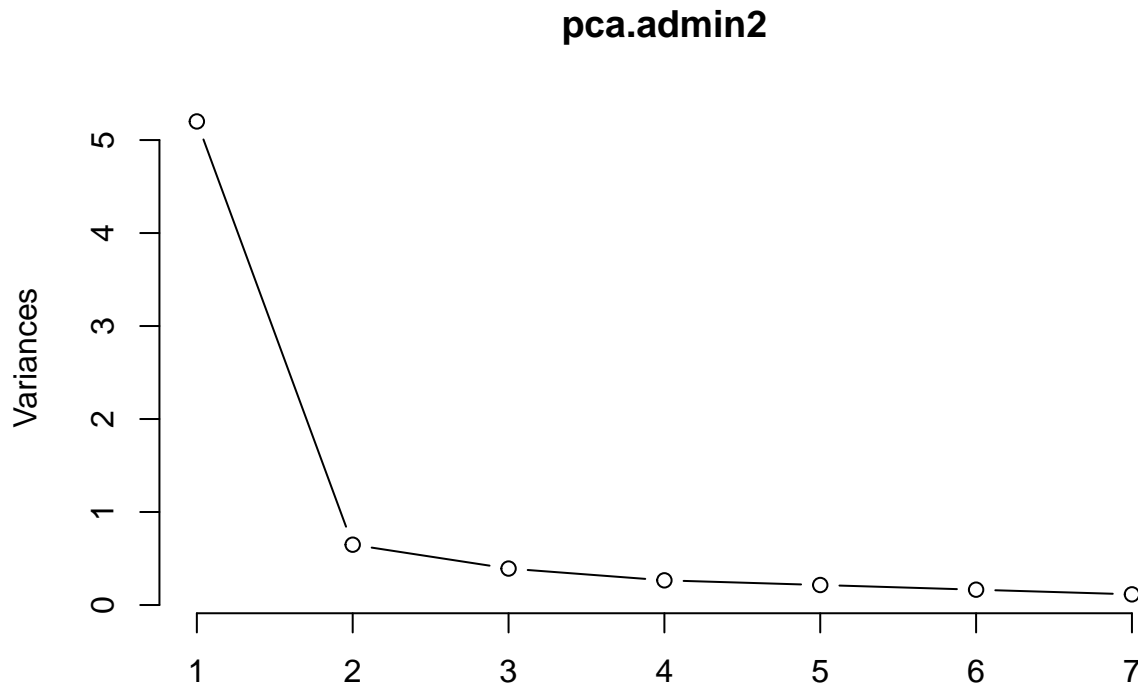
```
## Importance of components:
##              PC1      PC2      PC3      PC4      PC5      PC6
## Standard deviation  2.2803 0.80529 0.62599 0.5150 0.46369 0.40586
## Proportion of Variance 0.7429 0.09264 0.05598 0.0379 0.03071 0.02353
## Cumulative Proportion 0.7429 0.83549 0.89147 0.9294 0.96008 0.98361
##              PC7
## Standard deviation  0.33868
## Proportion of Variance 0.01639
## Cumulative Proportion 1.00000
```

To choose the number of principal components to keep, we can either use the Kaiser criterion, cumulative proportion/percent of variance, or a scree plot.

Using the Kaiser criterion, we keep all principal components with a standard deviation greater than 1 (since the data is scaled). Hence the Kaiser criterion is telling us to keep the first principal component.

I will now compare this with a scree plot.

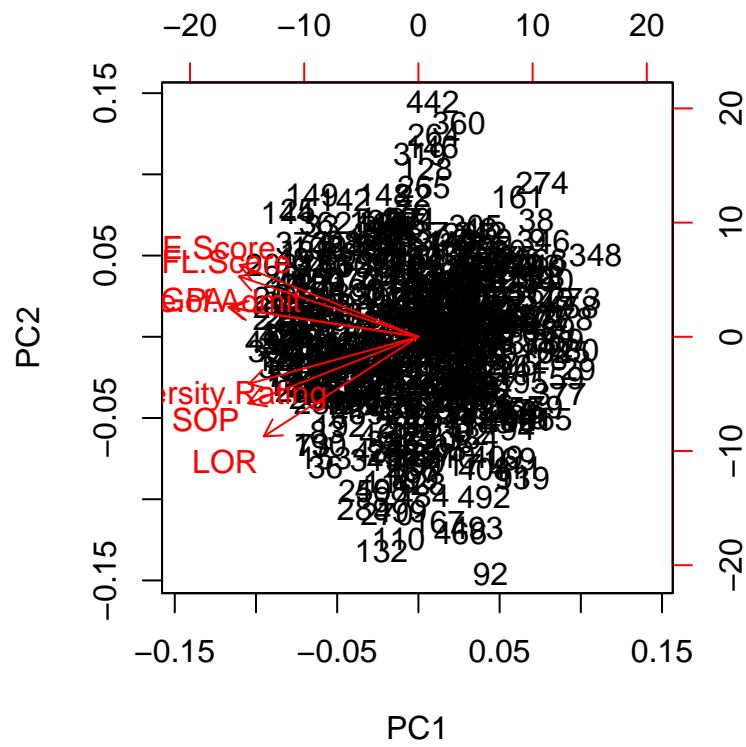
```
plot(pca.admin2, type="lines")
```



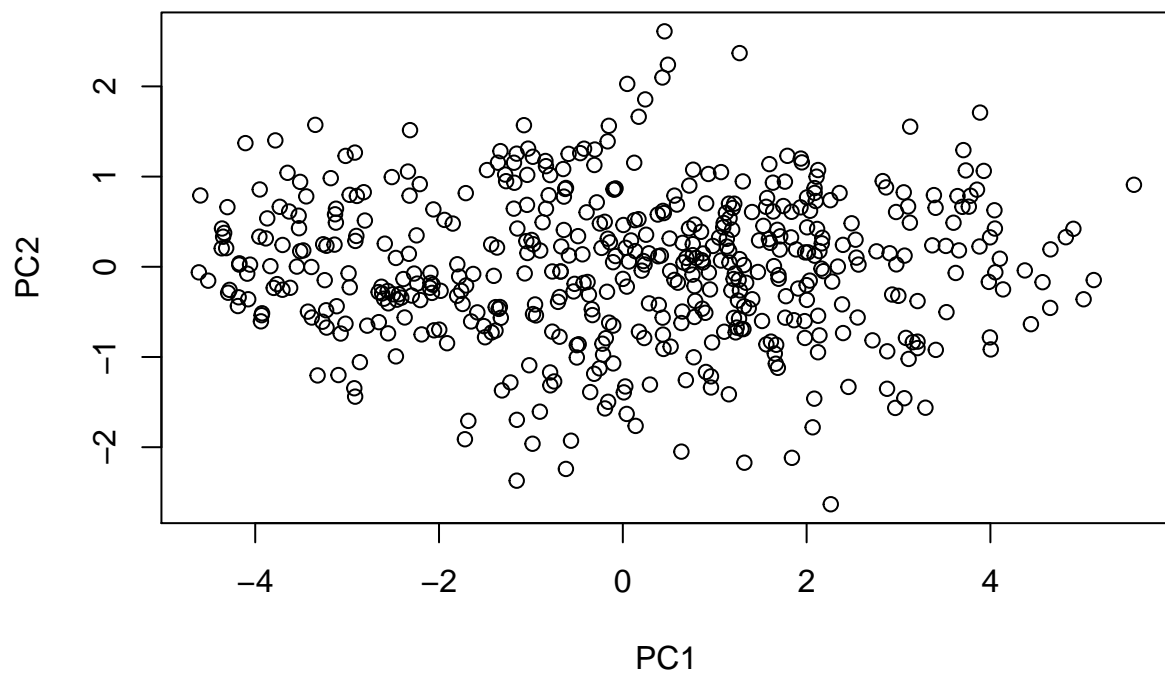
The above scree plot plots the monotonically decreasing eigenvalues and the location of an ‘elbow’ or plateau indicates the number of principal components. The scree plot suggests probably 2 principal components.

The first two principal components that will be retained explain 84% of the variation in the data. We can now view the data projected onto the components using a biplot.

```
biplot(pca.admin2)
```

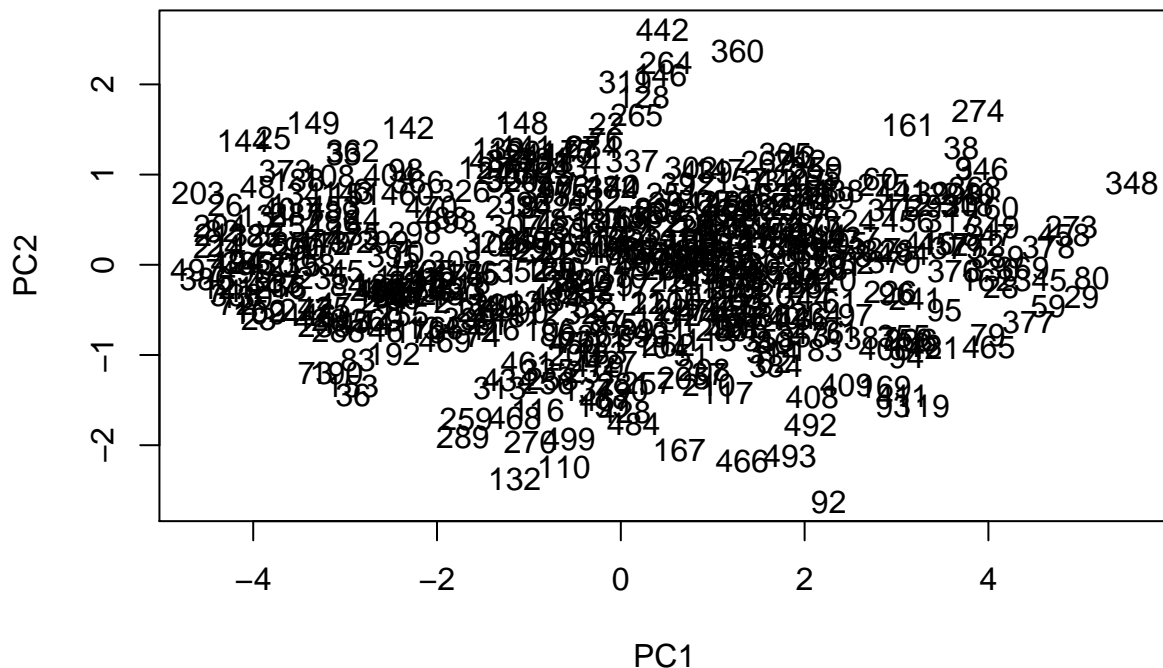


```
plot(pca.admin2$x[,1:2])
```

We can put data labels on the biplot by observation number

```
plot(pca.admin2$x[,1:2], type = "n")  
text(pca.admin2$x[,1:2], labels = 1:nrow(admissionsData))
```



It looks like there are two groups in the above principal component plots.

Take a look at the component loadings (eigenvectors) which provide the coefficients of the original variables, rounded to 2 decimal places.

```
round(pca.admin2$rotation[,1:2], 2)
```

```
##          PC1    PC2
## GRE.Score   -0.38  0.44
## TOEFL.Score -0.39  0.37
## University.Rating -0.36 -0.29
## SOP         -0.37 -0.40
## LOR         -0.33 -0.61
## CGPA        -0.41  0.18
## Chance.of.Admit -0.40  0.17
```

These are the coefficients of the original variables. The magnitudes are extremely similar for the first component. They are also all containing the same sign. This is a little difficult to interpret again, but most likely indicates that the first principal component is equally weighting all predictor variables.

In the second component, the highest magnitude is the letter of recommendation which has a negative sign. Other variables with the same sign include the SOP score and the university rating. Variables of opposite sign with higher magnitude include GRE Score, TOEFL Score, as well as CGPA and Chance of Admit having a lower magnitude. Students who score high on this principal component, likely scored high on their standardized tests.

We can now look at the four students who scored highest on PC1:

```
admissionsData[order(pca.admin2$x[,1], decreasing = TRUE)[1:4],1:9]
```

```
##      Serial.No. GRE.Score TOEFL.Score University.Rating SOP LOR CGPA
## 348      348      299      94              1 1.0 1.0 7.34
## 80       80      294      93              1 1.5 2.0 7.36
## 29       29      295      93              1 2.0 2.0 7.20
## 273      273      294      95              1 1.5 1.5 7.64
##      Research Chance.of.Admit
## 348      0      0.42
## 80       0      0.46
## 29       0      0.46
## 273      0      0.49
```

The top four students in this first principal component are the same as the first four students in the previous PC1 (compared using Serial.No.). Even when looking at the loadings, this principal component is very similar to the principal component in the previous section.

And the four students who scored highest on PC2:

```
admissionsData[order(pca.admin2$x[,2], decreasing = TRUE)[1:4], 1:9]
```

```
##      Serial.No. GRE.Score TOEFL.Score University.Rating SOP LOR CGPA
## 442      442      332      112              1 1.5 3.0 8.66
## 360      360      321      107              2 2.0 1.5 8.44
## 264      264      324      111              3 2.5 1.5 8.79
## 146      146      320      113              2 2.0 2.5 8.64
##      Research Chance.of.Admit
## 442      1      0.79
## 360      0      0.81
## 264      1      0.70
## 146      1      0.81
```

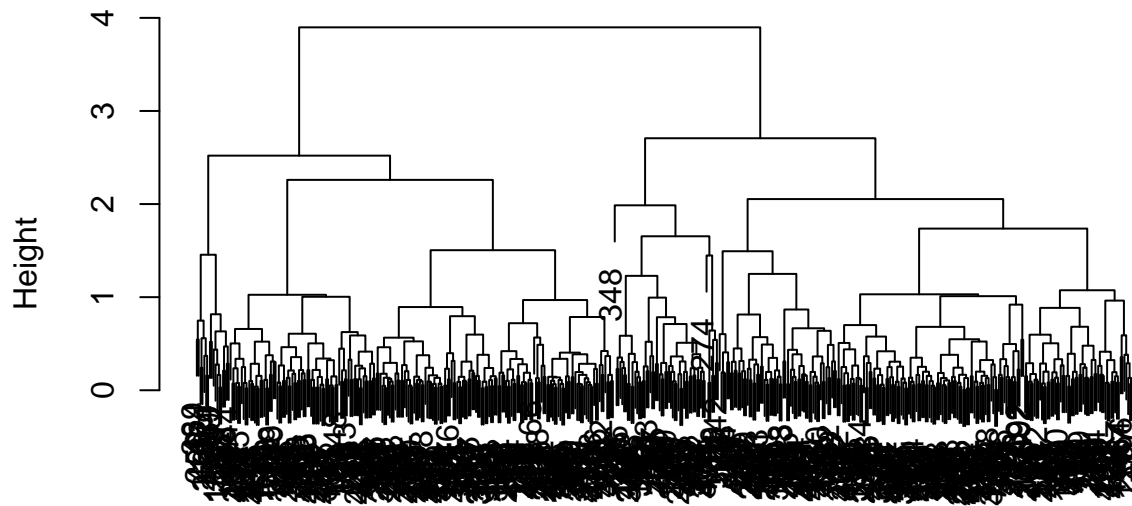
As hypothesized above, the first four students in PC2 are scoring higher on their standardized tests (GRE.Score and TOEFL.Score). These students are performing the at, or above average on these standardized tests. However, they all have a below average score on SOP, and LOR. The CGPA of the students scoring high on PC2 hovers fairly close to the mean. This proves the initial hypothesis that standardized testing is most important for PC2.

Trying to Cluster on the First Two Principal Components

It appeared that in the first PCA analysis, with the predictor chance.of.Admit removed, there were two groups in the remaining principal component plots. Here we will perform hierarchical clustering to try to find these groups.

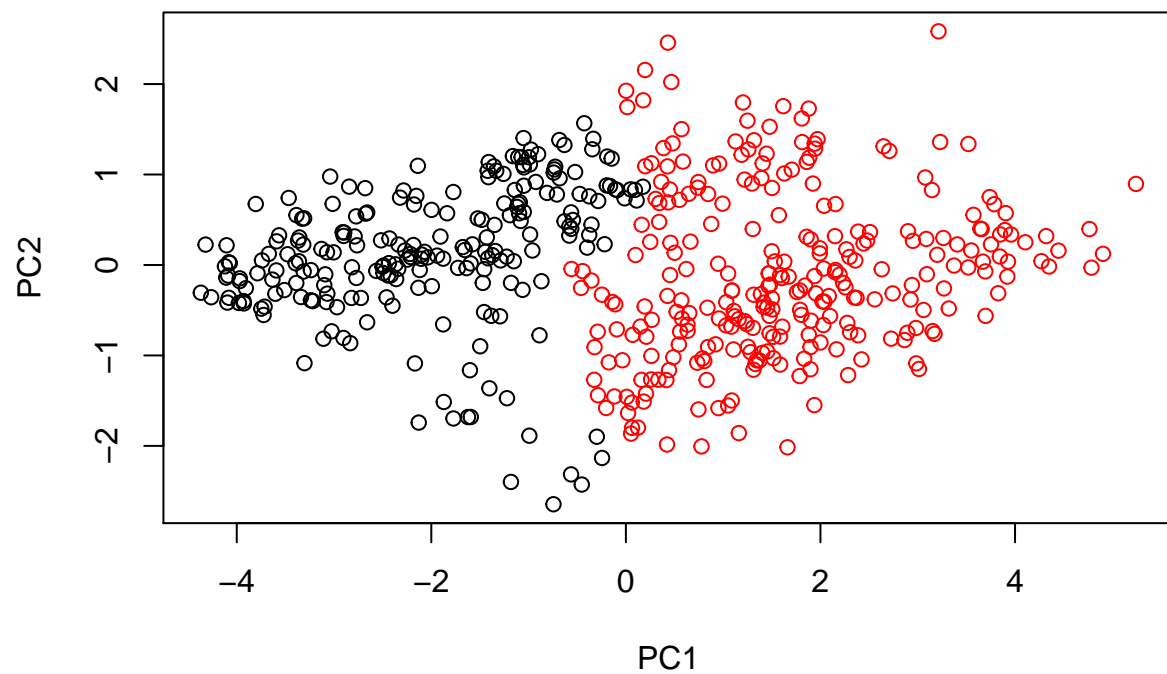
```
set.seed(574847)
clusts <- hclust(dist(pca.admin$x[,1:2]), method="average")
plot(clusts)
```

Cluster Dendrogram

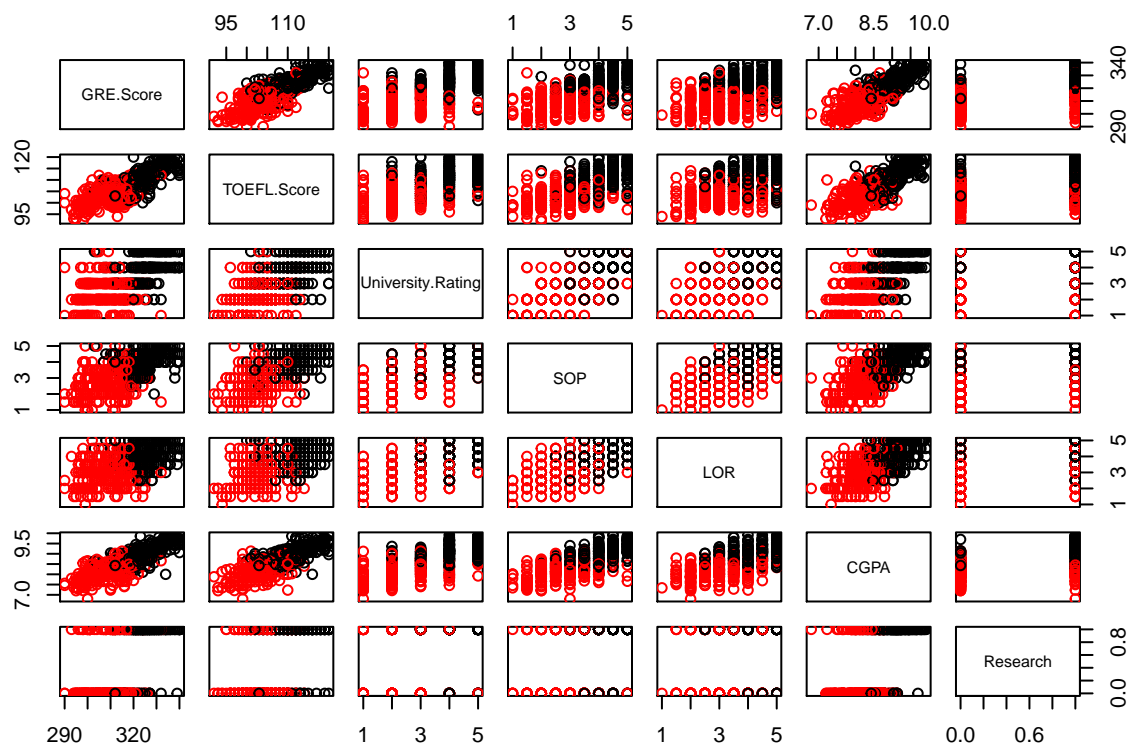


```
dist(pca.admin$x[, 1:2])  
hclust (*, "average")
```

```
plot(pca.admin$x[, 1:2], col=cutree(clusts, 2))
```



```
pairs(admissionsData[, -c(1,9)], col=cutree(clusts,2))
```



These are not quiet the groups we noticed by eye. Let's try clustering with a mixture model.

```
#install.packages("mclust")
#install.packages("teigen")
library(mclust)
library(teigen)

mPCA <- Mclust(dist(pca.admin$x[,1:2]), G=1:5, scale = TRUE)
summary(mPCA)

mPCA2 <- Mclust(dist(pca.admin$x[,1:2]), G = 2)
summary(mPCA2)
plot(mPCA2)

plot(pca.admin$x[,1:2], col=mPCA2$classification)

set.seed(2521)
tfaith <- teigen(as.matrix(dist(pca.admin$x[,1:2])), models = "CCCC", Gs = 1:4, verbose = FALSE)
plot(tfaith, what = "uncertainty", cex = 1.5, uncmult = 1.5)
plot(tfaith, what = "contour")
plot(tfaith, ymarg = NULL, lwd = 2)

tPCA <- teigen(as.matrix(dist(pca.admin$x[,1:2])), Gs=1:9, models="all", scale= FALSE, verbose = TRUE,
plot(tPCA, what = "uncertainty")
plot(tPCA, what = "contour")
```

```

tfaith <- teigen(faithful, Gs=1:5, model="UUUU", scale=FALSE, verbose = FALSE)
plot(tfaith, what="uncertainty")
plot(tfaith, what="contour")
library(gclus)
data(wine)
twine <- teigen(wine[, -1], Gs=1:5, model="UCCU", scale=FALSE, verbose=FALSE)
plot(twine, xmarg=1, ymarg=7, what="contour")
plot(twine, xmarg=1, ymarg=10, what="uncertainty")
table(wine[, 1], twine$classification)
plot(twine$allbic, type="l")
points(twine$allbic)

```