

class09_mini_project

Hope (PID: A15652616)

10/26/2021

Import data

```
# Save your input data file into your Project directory  
fna.data <- "WisconsinCancer.csv"  
  
# Complete the following code to input the data and store as wisc.df  
wisc.df <- read.csv(fna.data, row.names=1)
```

We don't want to include the analysis given in the file

```
# We can use -1 here to remove the first column  
wisc.data <- wisc.df[,-1]
```

setup a separate new vector called diagnosis that contains the data from the diagnosis column of the original dataset. We will store this as a factor (useful for plotting) and use this later to check our results.

```
# Create diagnosis vector for later  
diagnosis <- as.factor(c(wisc.df[,1]) )
```

Q1. How many observations are in this dataset?

```
dim(wisc.data)
```

```
## [1] 569 30
```

There are 569 patients with 30 different measurements

Q2. How many of the observations have a malignant diagnosis?

```
table(diagnosis)
```

```
## diagnosis  
##    B    M  
## 357 212
```

There were 212 people with a malignant diagnosis.

Q3. How many variables/features in the data are suffixed with `_mean`?

```
?grep
colnames(wisc.df)

## [1] "diagnosis"          "radius_mean"
## [3] "texture_mean"      "perimeter_mean"
## [5] "area_mean"         "smoothness_mean"
## [7] "compactness_mean"  "concavity_mean"
## [9] "concave.points_mean" "symmetry_mean"
## [11] "fractal_dimension_mean" "radius_se"
## [13] "texture_se"        "perimeter_se"
## [15] "area_se"           "smoothness_se"
## [17] "compactness_se"    "concavity_se"
## [19] "concave.points_se" "symmetry_se"
## [21] "fractal_dimension_se" "radius_worst"
## [23] "texture_worst"     "perimeter_worst"
## [25] "area_worst"        "smoothness_worst"
## [27] "compactness_worst" "concavity_worst"
## [29] "concave.points_worst" "symmetry_worst"
## [31] "fractal_dimension_worst"
```

```
length(grep("mean", colnames(wisc.df)))
```

```
## [1] 10
```

There are 10 variables/features with the suffix “mean”

The next step in your analysis is to perform principal component analysis (PCA) on `wisc.data`.

Check the mean and standard deviation of the features (i.e. columns) of the `wisc.data` to determine if the data should be scaled. Use the `colMeans()` and `apply()` functions like you’ve done before.

```
# Check column means and standard deviations
colMeans(wisc.data)

##           radius_mean      texture_mean      perimeter_mean
##      1.412729e+01      1.928965e+01      9.196903e+01
##           area_mean      smoothness_mean      compactness_mean
##      6.548891e+02      9.636028e-02      1.043410e-01
##      concavity_mean      concave.points_mean      symmetry_mean
##      8.879932e-02      4.891915e-02      1.811619e-01
## fractal_dimension_mean      radius_se      texture_se
##      6.279761e-02      4.051721e-01      1.216853e+00
##      perimeter_se      area_se      smoothness_se
##      2.866059e+00      4.033708e+01      7.040979e-03
##      compactness_se      concavity_se      concave.points_se
##      2.547814e-02      3.189372e-02      1.179614e-02
##      symmetry_se      fractal_dimension_se      radius_worst
##      2.054230e-02      3.794904e-03      1.626919e+01
```

```
##          texture_worst      perimeter_worst      area_worst
##          2.567722e+01      1.072612e+02      8.805831e+02
##          smoothness_worst      compactness_worst      concavity_worst
##          1.323686e-01      2.542650e-01      2.721885e-01
##          concave.points_worst      symmetry_worst      fractal_dimension_worst
##          1.146062e-01      2.900756e-01      8.394582e-02
```

```
apply(wisc.data,2,sd)
```

```
##          radius_mean      texture_mean      perimeter_mean
##          3.524049e+00      4.301036e+00      2.429898e+01
##          area_mean      smoothness_mean      compactness_mean
##          3.519141e+02      1.406413e-02      5.281276e-02
##          concavity_mean      concave.points_mean      symmetry_mean
##          7.971981e-02      3.880284e-02      2.741428e-02
##          fractal_dimension_mean      radius_se      texture_se
##          7.060363e-03      2.773127e-01      5.516484e-01
##          perimeter_se      area_se      smoothness_se
##          2.021855e+00      4.549101e+01      3.002518e-03
##          compactness_se      concavity_se      concave.points_se
##          1.790818e-02      3.018606e-02      6.170285e-03
##          symmetry_se      fractal_dimension_se      radius_worst
##          8.266372e-03      2.646071e-03      4.833242e+00
##          texture_worst      perimeter_worst      area_worst
##          6.146258e+00      3.360254e+01      5.693570e+02
##          smoothness_worst      compactness_worst      concavity_worst
##          2.283243e-02      1.573365e-01      2.086243e-01
##          concave.points_worst      symmetry_worst      fractal_dimension_worst
##          6.573234e-02      6.186747e-02      1.806127e-02
```

Execute PCA with the `prcomp()` function on the `wisc.data`, scaling if appropriate, and assign the output model to `wisc.pr`.

```
# Perform PCA on wisc.data by completing the following code
wisc.pr <- prcomp(wisc.data, scale. = TRUE)
```

Summarize

```
# Look at summary of results
summary(wisc.pr)
```

```
## Importance of components:
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation      3.6444 2.3857 1.67867 1.40735 1.28403 1.09880 0.82172
## Proportion of Variance 0.4427 0.1897 0.09393 0.06602 0.05496 0.04025 0.02251
## Cumulative Proportion 0.4427 0.6324 0.72636 0.79239 0.84734 0.88759 0.91010
##          PC8      PC9      PC10      PC11      PC12      PC13      PC14
## Standard deviation      0.69037 0.6457 0.59219 0.5421 0.51104 0.49128 0.39624
## Proportion of Variance 0.01589 0.0139 0.01169 0.0098 0.00871 0.00805 0.00523
## Cumulative Proportion 0.92598 0.9399 0.95157 0.9614 0.97007 0.97812 0.98335
##          PC15      PC16      PC17      PC18      PC19      PC20      PC21
## Standard deviation      0.30681 0.28260 0.24372 0.22939 0.22244 0.17652 0.1731
```

```
## Proportion of Variance 0.00314 0.00266 0.00198 0.00175 0.00165 0.00104 0.0010
## Cumulative Proportion 0.98649 0.98915 0.99113 0.99288 0.99453 0.99557 0.9966
##          PC22      PC23      PC24      PC25      PC26      PC27      PC28
## Standard deviation    0.16565 0.15602 0.1344 0.12442 0.09043 0.08307 0.03987
## Proportion of Variance 0.00091 0.00081 0.0006 0.00052 0.00027 0.00023 0.00005
## Cumulative Proportion 0.99749 0.99830 0.9989 0.99942 0.99969 0.99992 0.99997
##          PC29      PC30
## Standard deviation    0.02736 0.01153
## Proportion of Variance 0.00002 0.00000
## Cumulative Proportion 1.00000 1.00000
```

Q4. From your results, what proportion of the original variance is captured by the first principal components (PC1)?

There is 44% of the original variance captured by PC1.

Q5. How many principal components (PCs) are required to describe at least 70% of the original variance in the data?

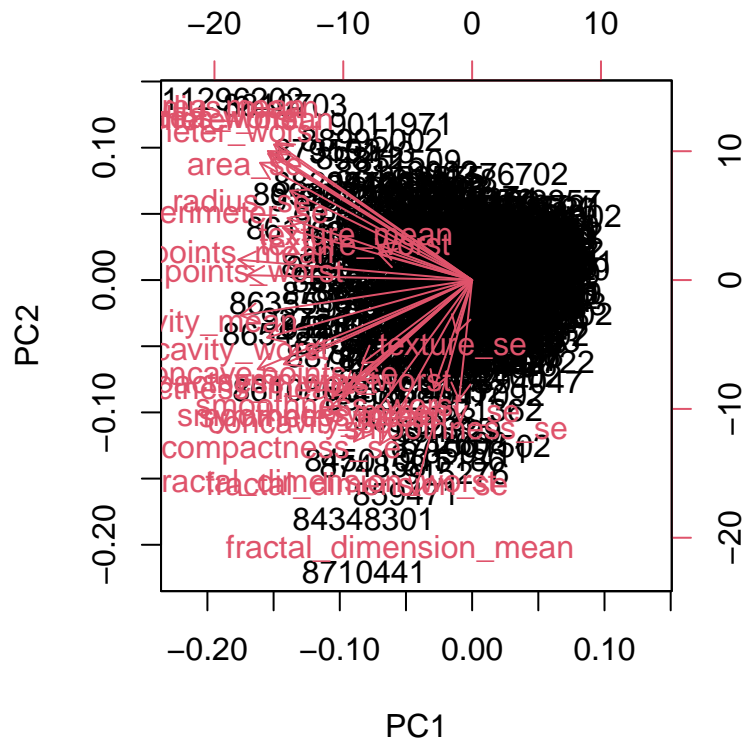
You need 3 principal components to describe at least 70% of the original variance in the data.

Q6. How many principal components (PCs) are required to describe at least 90% of the original variance in the data?

You need 7 principal components to describe at least 90% of the original variance in the data.

Look at biplot

```
biplot(wisc.pr)
```

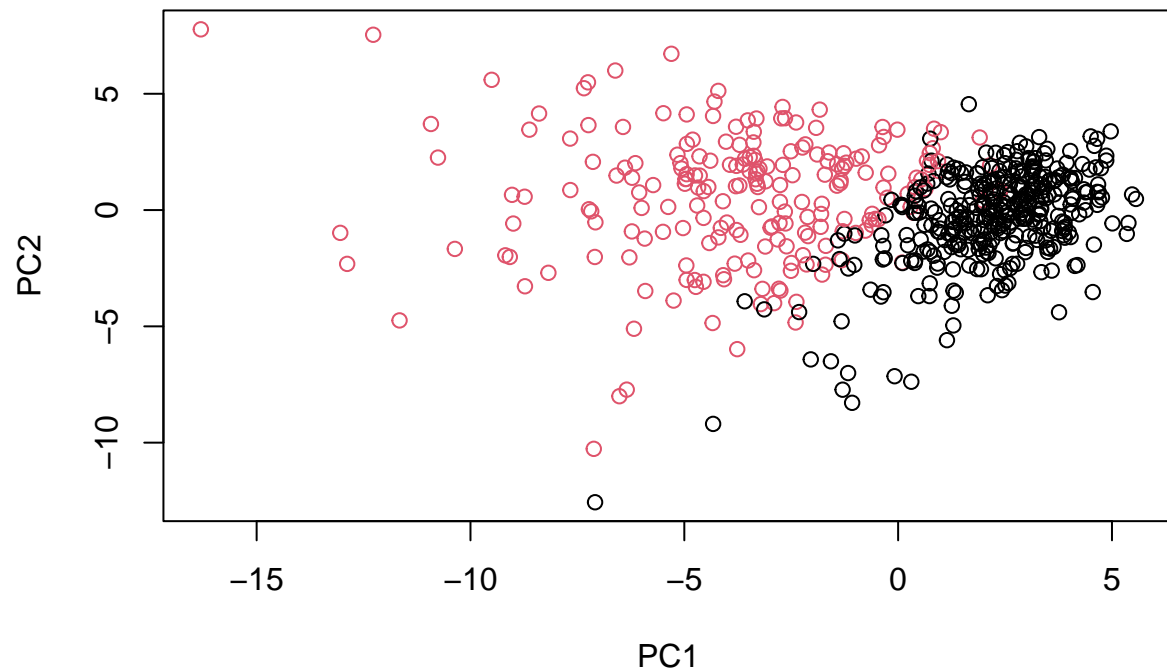


Q7. What stands out to you about this plot? Is it easy or difficult to understand? Why? This plot is difficult to understand and everything is too close to be able to actually read.

We can separate this out a little better.

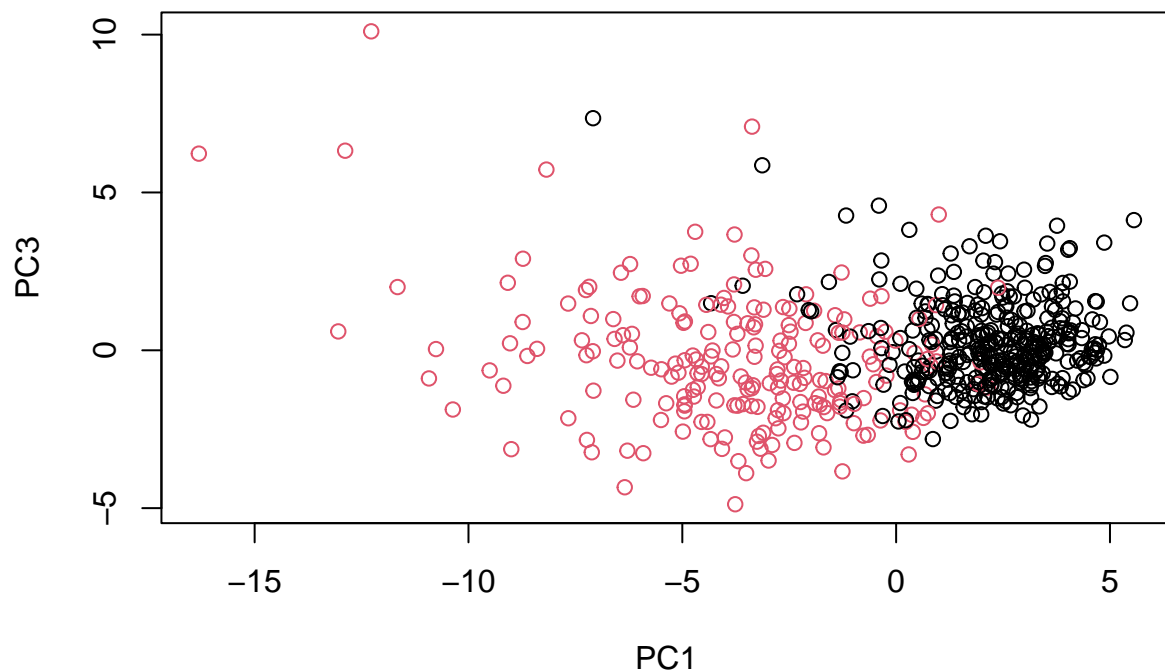
Let's make a score plot

```
# Repeat for components 1 and 3
plot(wisc.pr$x[, 1:2 ], col = as.factor(diagnosis))
```



Q8. Generate a similar plot for principal components 1 and 3. What do you notice about these plots?

```
plot(wisc.pr$x[, c(1,3) ], col = as.factor(diagnosis),  
     xlab = "PC1", ylab = "PC3")
```



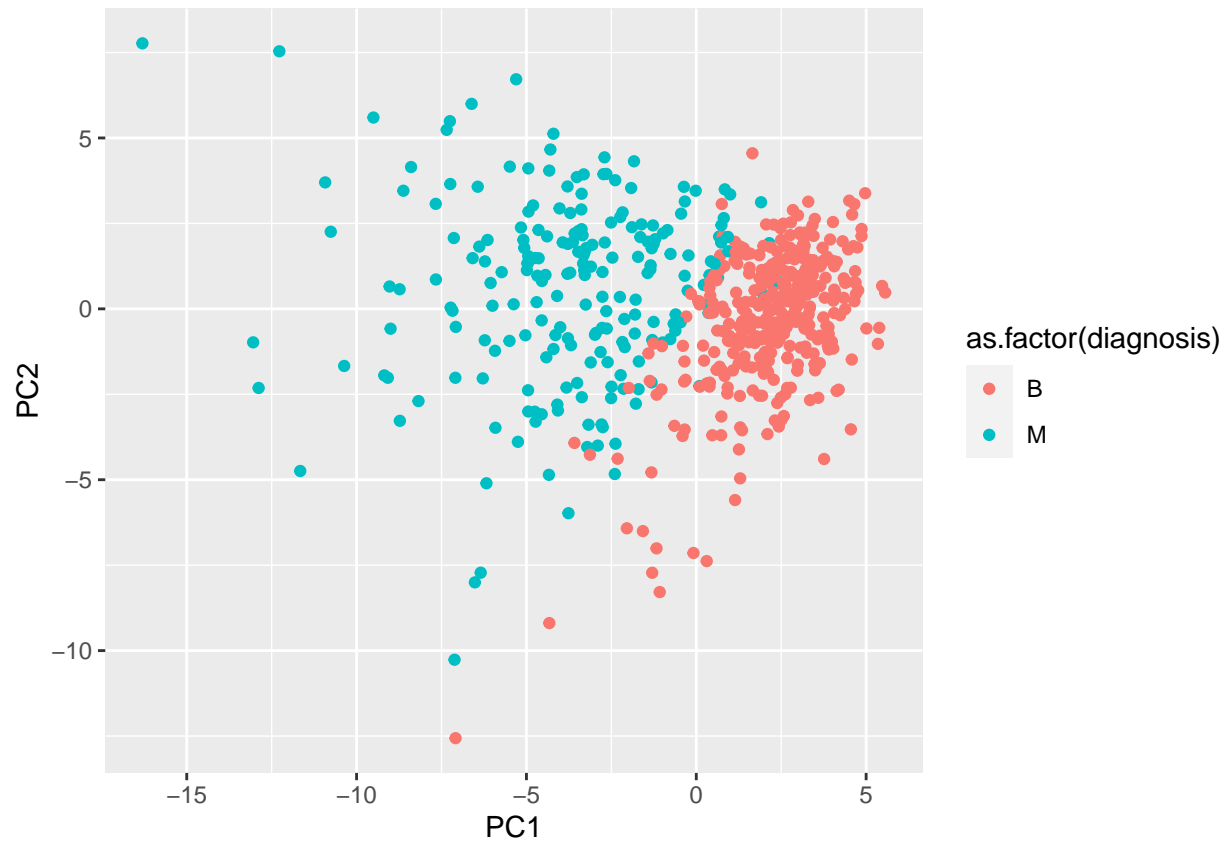
These graphs make it much easier to compare different principal components compared to the biplot seen before. We can see that PC1 explains the most variance compared to both PC2 and PC3. Overall we see there is a bunching of malignant vs benign samples.

We can also graph this using ggplot to make it look a little better.

```
# Create a data.frame for ggplot
df <- as.data.frame(wisc.pr$x)
df$diagnosis <- diagnosis

# Load the ggplot2 package
library(ggplot2)

# Make a scatter plot colored by diagnosis
ggplot(df) +
  aes(PC1, PC2, col=as.factor(diagnosis)) +
  geom_point()
```



We now want to explain the variance

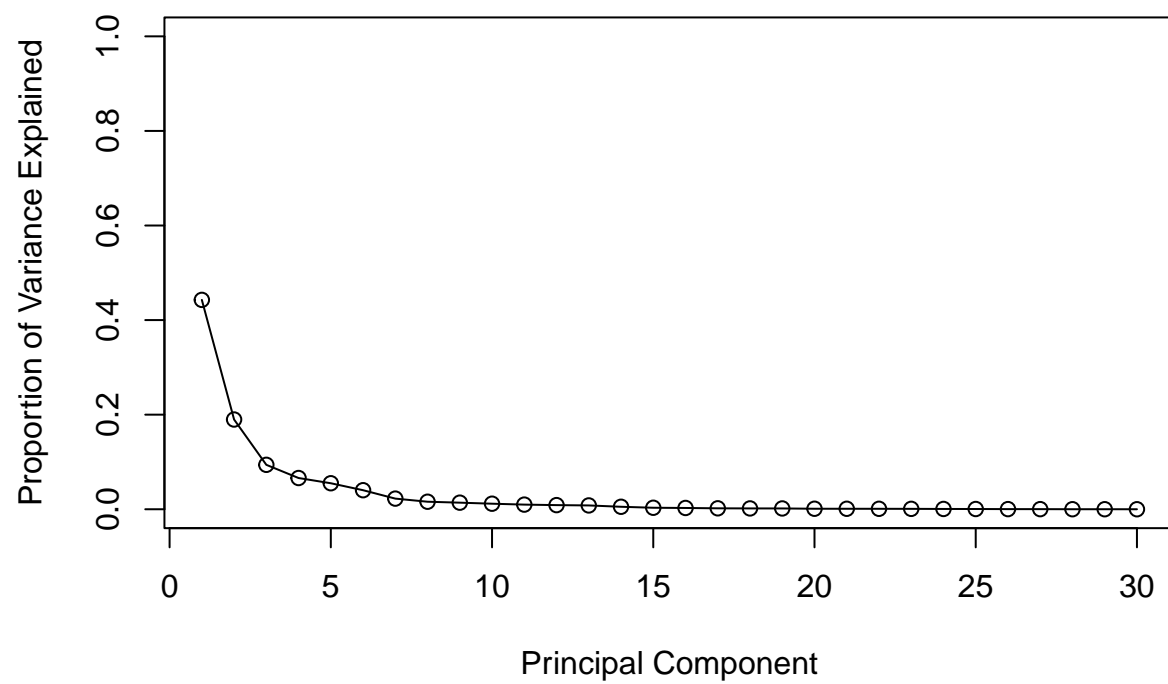
```
# Calculate variance of each component
pr.var <- wisc.pr$sdev^2
head(pr.var)
```

```
## [1] 13.281608  5.691355  2.817949  1.980640  1.648731  1.207357
```

Calculate the variance explained by each principal component by dividing by the total variance explained of all principal components. Assign this to a variable called pve and create a plot of variance explained for each principal component.

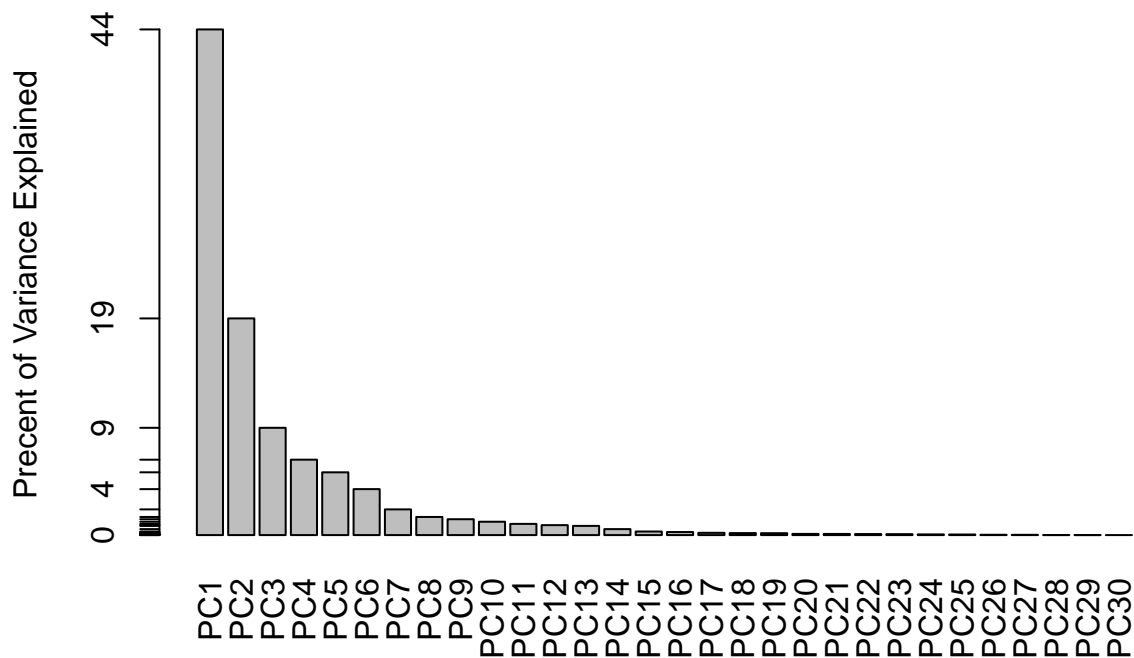
```
# Variance explained by each principal component: pve
pve <- pr.var / sum(pr.var)

# Plot variance explained for each principal component
plot(pve, xlab = "Principal Component",
     ylab = "Proportion of Variance Explained",
     ylim = c(0, 1), type = "o")
```

Look at a scree plot

```
# Alternative scree plot of the same data, note data driven y-axis
barplot(pve, ylab = "Precent of Variance Explained",
        names.arg=paste0("PC",1:length(pve)), las=2, axes = FALSE)
axis(2, at=pve, labels=round(pve,2)*100 )
```



Q9. For the first principal component, what is the component of the loading vector (i.e. `wisc.pr$rotation[,1]`) for the feature `concave.points_mean`?

```
wisc.pr$rotation["concave.points_mean",1]
```

```
## [1] -0.2608538
```

Q10. What is the minimum number of principal components required to explain 80% of the variance of the data?

```
summary(wisc.pr)
```

```
## Importance of components:
##              PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation    3.6444  2.3857  1.67867  1.40735  1.28403  1.09880  0.82172
## Proportion of Variance 0.4427  0.1897  0.09393  0.06602  0.05496  0.04025  0.02251
## Cumulative Proportion 0.4427  0.6324  0.72636  0.79239  0.84734  0.88759  0.91010
```

##	PC8	PC9	PC10	PC11	PC12	PC13	PC14
## Standard deviation	0.69037	0.6457	0.59219	0.5421	0.51104	0.49128	0.39624
## Proportion of Variance	0.01589	0.0139	0.01169	0.0098	0.00871	0.00805	0.00523
## Cumulative Proportion	0.92598	0.9399	0.95157	0.9614	0.97007	0.97812	0.98335
##	PC15	PC16	PC17	PC18	PC19	PC20	PC21
## Standard deviation	0.30681	0.28260	0.24372	0.22939	0.22244	0.17652	0.1731
## Proportion of Variance	0.00314	0.00266	0.00198	0.00175	0.00165	0.00104	0.0010
## Cumulative Proportion	0.98649	0.98915	0.99113	0.99288	0.99453	0.99557	0.9966
##	PC22	PC23	PC24	PC25	PC26	PC27	PC28
## Standard deviation	0.16565	0.15602	0.1344	0.12442	0.09043	0.08307	0.03987
## Proportion of Variance	0.00091	0.00081	0.0006	0.00052	0.00027	0.00023	0.00005
## Cumulative Proportion	0.99749	0.99830	0.9989	0.99942	0.99969	0.99992	0.99997
##	PC29	PC30					
## Standard deviation	0.02736	0.01153					
## Proportion of Variance	0.00002	0.00000					
## Cumulative Proportion	1.00000	1.00000					

You need 5 PC's to explain 80% of the variance

Let's start looking at hierarchical clustering

Scale the wisc.data data using the "scale()" function

```
data.scaled <- scale(wisc.data)
```

Calculate the (Euclidean) distances between all pairs of observations in the new scaled dataset and assign the result to data.dist.

```
data.dist <- dist(data.scaled)
```

Create a hierarchical clustering model using complete linkage. Manually specify the method argument to hclust() and assign the results to wisc.hclust.

```
wisc.hclust <- hclust(data.dist)
```

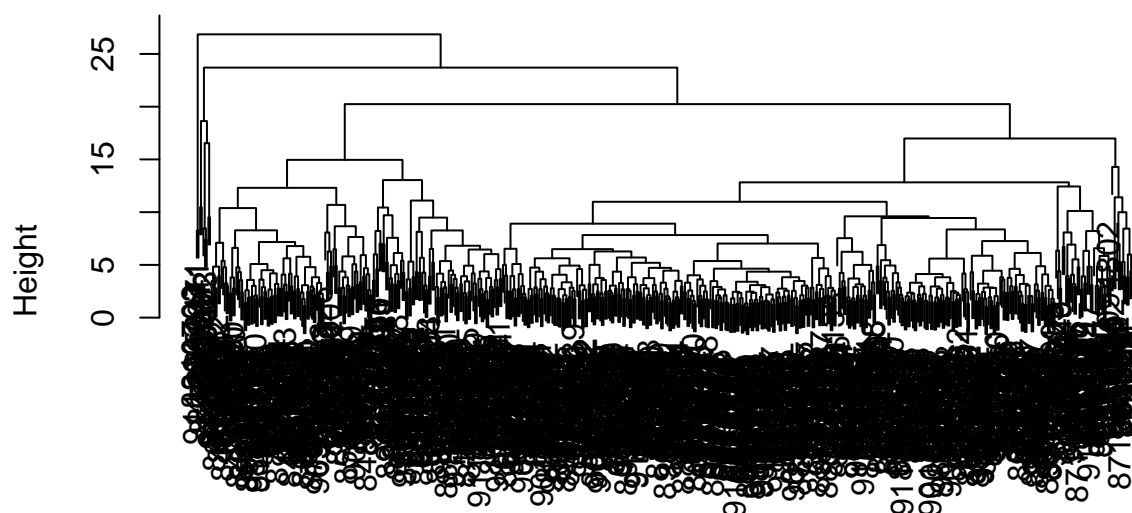
Results of hierarchical clustering

Let's use the hierarchical clustering model you just created to determine a height (or distance between clusters) where a certain number of clusters exists.

Q11. Using the plot() and abline() functions, what is the height at which the clustering model has 4 clusters?

```
plot(wisc.hclust)
```

Cluster Dendrogram



```
data.dist
hclust (*, "complete")
```

The height is 21

Selecting number of clusters

This exercise will help you determine if, in this case, hierarchical clustering provides a promising new feature. Use `cutree()` to cut the tree so that it has 4 clusters. Assign the output to the variable `wisc.hclust.clusters`.

```
wisc.hclust.clusters <- cutree(wisc.hclust, k=4)
```

We can use the `table()` function to compare the cluster membership to the actual diagnoses.

```
table(wisc.hclust.clusters, diagnosis)
```

```
##           diagnosis
## wisc.hclust.clusters  B  M
##           1  12 165
##           2   2   5
##           3 343  40
##           4   0   2
```

Q12. Can you find a better cluster vs diagnoses match by cutting into a different number of clusters between 2 and 10?

The best cluster match is 2.

Q13. Which method gives your favorite results for the same data.dist dataset? Explain your reasoning.

```
wisc.pc.hclust <- hclust(dist(wisc.pr$x[,1:3]),  
                        method="ward.D2")
```

```
wisc.pc.hclust <- hclust(dist(wisc.pr$x[,1:3]),  
                        method="single")
```

ward.d2 was the best method because there was a logical flow to the clustering that was easier for me to understand.

Clustering on PCA results

In this final section, you will put together several steps you used earlier and, in doing so, you will experience some of the creativity and open endedness that is typical in unsupervised learning.

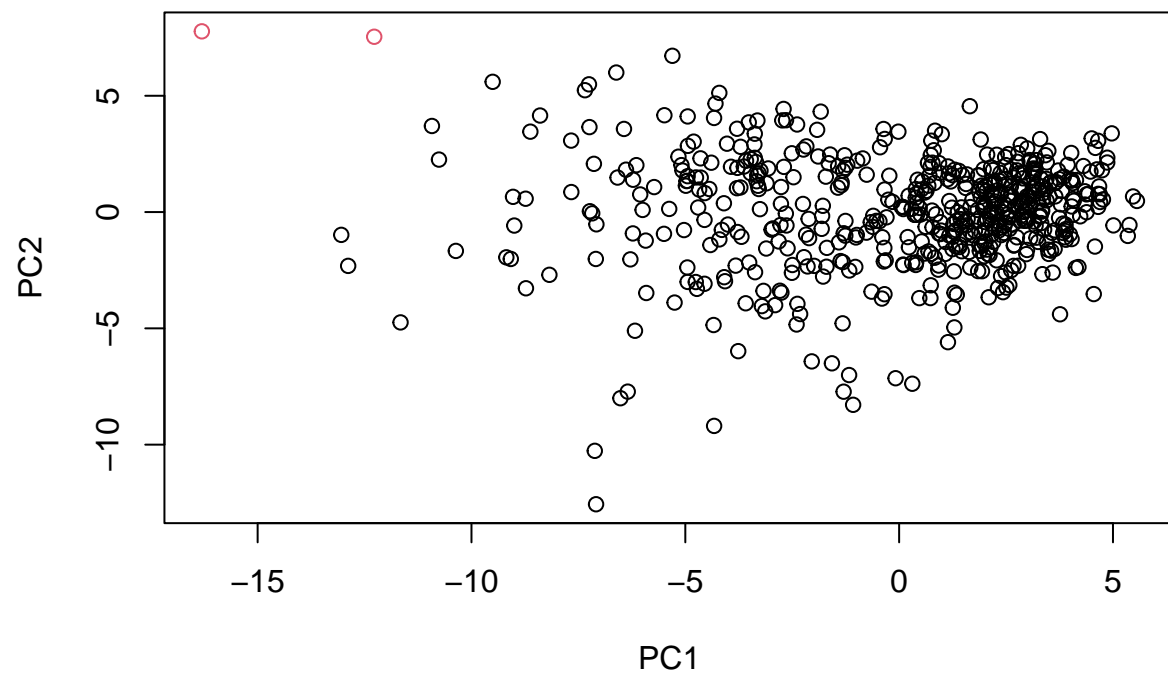
```
grps <- cutree(wisc.pc.hclust, k=2)  
table(grps)
```

```
## grps  
##    1    2  
## 567    2
```

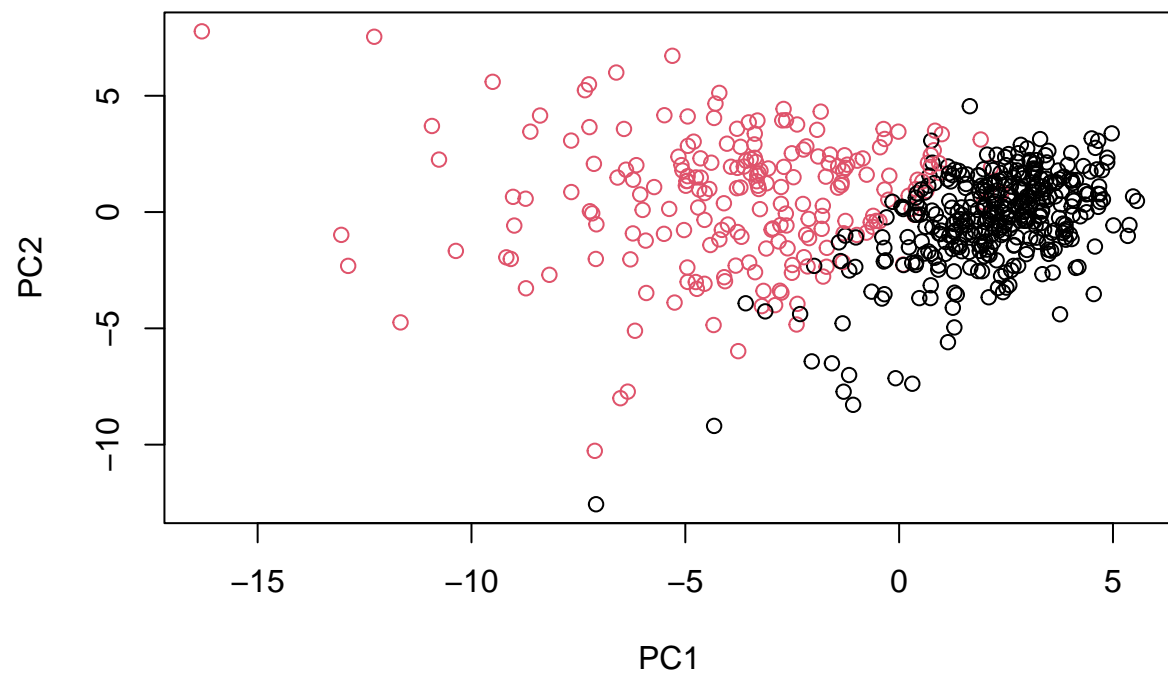
```
table(grps, diagnosis)
```

```
##      diagnosis  
## grps    B    M  
##    1 357 210  
##    2   0   2
```

```
plot(wisc.pr$x[,1:2], col=grps)
```



```
plot(wisc.pr$x[,1:2], col=diagnosis)
```



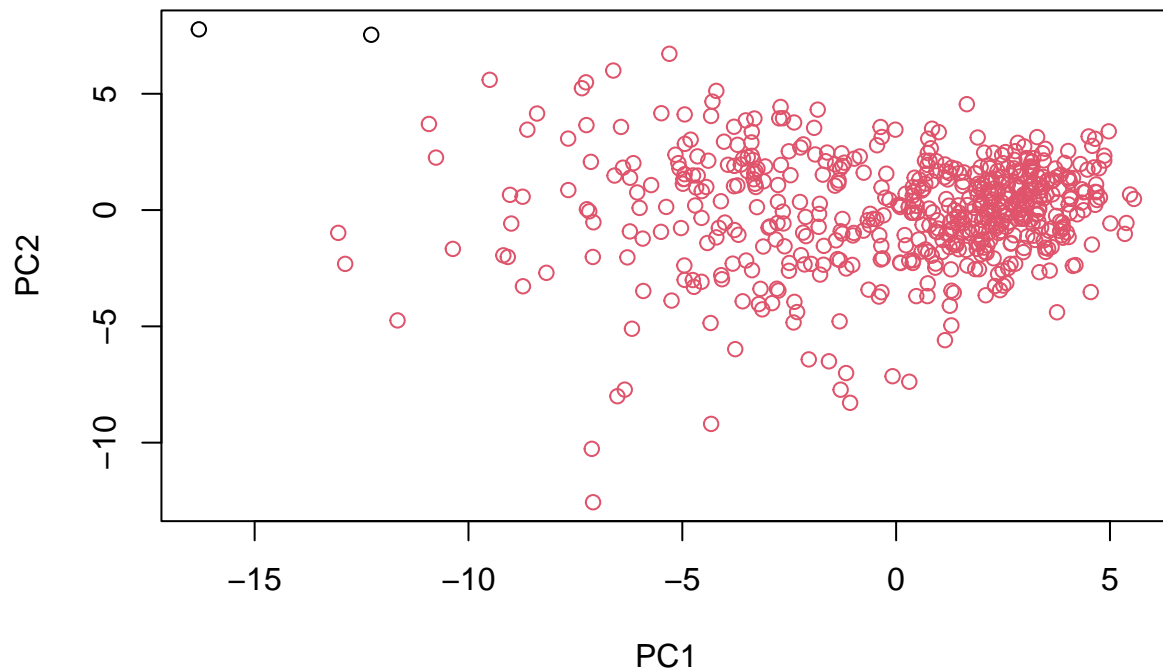
```
g <- as.factor(grps)
levels(g)
```

```
## [1] "1" "2"
```

```
g <- relevel(g,2)
levels(g)
```

```
## [1] "2" "1"
```

```
# Plot using our re-ordered factor
plot(wisc.pr$x[,1:2], col=g)
```



```
## Use the distance along the first 7 PCs for clustering i.e. wisc.pr$x[, 1:7]
wisc.pr.hclust <- hclust(dist(wisc.pr$x[, 1:7]), method="ward.D2")
```

```
wisc.pr.hclust.clusters <- cutree(wisc.pr.hclust, k=2)
```

Q15. How well does the newly created model with four clusters separate out the two diagnoses?

It separates the two diagnoses out well.

Compare to actual diagnoses

```
table(wisc.pr.hclust.clusters, diagnosis)
```

```
##              diagnosis
## wisc.pr.hclust.clusters  B  M
##              1  28 188
##              2 329  24
```


Q16. How well do the k-means and hierarchical clustering models you created in previous sections (i.e. before PCA) do in terms of separating the diagnoses? Again, use the `table()` function to compare the output of each model (`wisc.km$cluster` and `wisc.hclust.clusters`) with the vector containing the actual diagnoses.

```
table(wisc.hclust.clusters, diagnosis)
```

```
##               diagnosis
## wisc.hclust.clusters  B  M
##                   1 12 165
##                   2  2  5
##                   3 343 40
##                   4  0  2
```

```
table(diagnosis)
```

```
## diagnosis
##    B    M
## 357 212
```

The k-means and hierarchical clustering models separate the diagnosis very well. This resulted in the exact same results that the original data concluded.

Sensitivity/Specificity

Q17. Which of your analysis procedures resulted in a clustering model with the best specificity? How about sensitivity?

Sensitivity: $TP/(TP+FN)$

```
179/(179+33)
```

```
## [1] 0.8443396
```

Specificity: $TN/(TN+FP)$

```
333/(333+24)
```

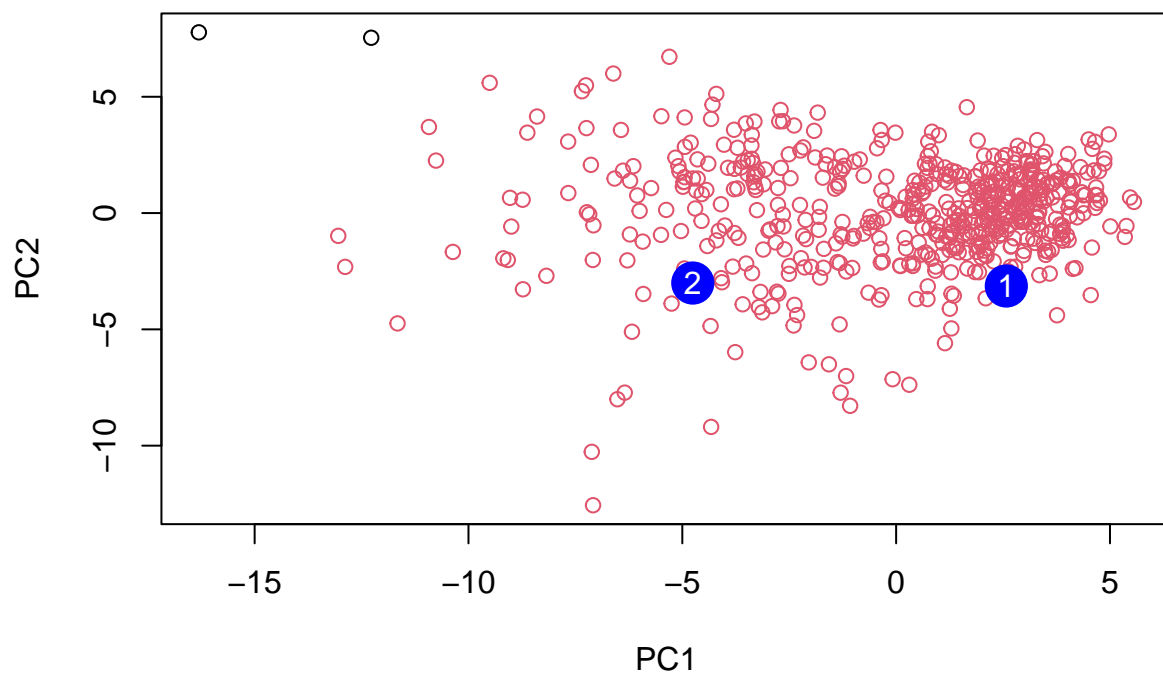
```
## [1] 0.9327731
```

Prediction

```
#url <- "new_samples.csv"
url <- "https://tinyurl.com/new-samples-CSV"
new <- read.csv(url)
npc <- predict(wisc.pr, newdata=new)
npc
```

```
##          PC1          PC2          PC3          PC4          PC5          PC6          PC7
## [1,]  2.576616 -3.135913  1.3990492 -0.7631950  2.781648 -0.8150185 -0.3959098
## [2,] -4.754928 -3.009033 -0.1660946 -0.6052952 -1.140698 -1.2189945  0.8193031
##          PC8          PC9          PC10          PC11          PC12          PC13          PC14
## [1,] -0.2307350 0.1029569 -0.9272861 0.3411457  0.375921 0.1610764 1.187882
## [2,] -0.3307423 0.5281896 -0.4855301 0.7173233 -1.185917 0.5893856 0.303029
##          PC15          PC16          PC17          PC18          PC19          PC20
## [1,] 0.3216974 -0.1743616 -0.07875393 -0.11207028 -0.08802955 -0.2495216
## [2,] 0.1299153 0.1448061 -0.40509706 0.06565549 0.25591230 -0.4289500
##          PC21          PC22          PC23          PC24          PC25          PC26
## [1,] 0.1228233 0.09358453 0.08347651 0.1223396 0.02124121 0.078884581
## [2,] -0.1224776 0.01732146 0.06316631 -0.2338618 -0.20755948 -0.009833238
##          PC27          PC28          PC29          PC30
## [1,] 0.220199544 -0.02946023 -0.015620933 0.005269029
## [2,] -0.001134152 0.09638361 0.002795349 -0.019015820
```

```
plot(wisc.pr$x[,1:2], col=g)
points(npc[,1], npc[,2], col="blue", pch=16, cex=3)
text(npc[,1], npc[,2], c(1,2), col="white")
```



Q18. Which of these new patients should we prioritize for follow up based on your results?

Patient 2 should be prioritized based on these results.