

# Class08

Fan Wu(PID:A15127541)

## Table of contents

|                                                |           |
|------------------------------------------------|-----------|
| Data Import . . . . .                          | 1         |
| Principal Component Analysis (PCA) . . . . .   | 4         |
| Variance Explained . . . . .                   | 9         |
| <b>Communicate PCA Results</b>                 | <b>11</b> |
| Hierarchical Clustering . . . . .              | 12        |
| Selecting number of clusters . . . . .         | 13        |
| Combining methods (PCA & Clustering) . . . . . | 14        |
| 7 Prediction . . . . .                         | 16        |

Data was downloaded from the class website as a CSV file.

## Data Import

```
# 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)

head(wisc.df)
```

|          | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean |
|----------|-----------|-------------|--------------|----------------|-----------|
| 842302   | M         | 17.99       | 10.38        | 122.80         | 1001.0    |
| 842517   | M         | 20.57       | 17.77        | 132.90         | 1326.0    |
| 84300903 | M         | 19.69       | 21.25        | 130.00         | 1203.0    |
| 84348301 | M         | 11.42       | 20.38        | 77.58          | 386.1     |
| 84358402 | M         | 20.29       | 14.34        | 135.10         | 1297.0    |

|          |                 |                        |                  |                     |                   |
|----------|-----------------|------------------------|------------------|---------------------|-------------------|
| 843786   | M               | 12.45                  | 15.70            | 82.57               | 477.1             |
|          | smoothness_mean | compactness_mean       | concavity_mean   | concave.points_mean |                   |
| 842302   | 0.11840         | 0.27760                | 0.3001           | 0.14710             |                   |
| 842517   | 0.08474         | 0.07864                | 0.0869           | 0.07017             |                   |
| 84300903 | 0.10960         | 0.15990                | 0.1974           | 0.12790             |                   |
| 84348301 | 0.14250         | 0.28390                | 0.2414           | 0.10520             |                   |
| 84358402 | 0.10030         | 0.13280                | 0.1980           | 0.10430             |                   |
| 843786   | 0.12780         | 0.17000                | 0.1578           | 0.08089             |                   |
|          | symmetry_mean   | fractal_dimension_mean | radius_se        | texture_se          | perimeter_se      |
| 842302   | 0.2419          | 0.07871                | 1.0950           | 0.9053              | 8.589             |
| 842517   | 0.1812          | 0.05667                | 0.5435           | 0.7339              | 3.398             |
| 84300903 | 0.2069          | 0.05999                | 0.7456           | 0.7869              | 4.585             |
| 84348301 | 0.2597          | 0.09744                | 0.4956           | 1.1560              | 3.445             |
| 84358402 | 0.1809          | 0.05883                | 0.7572           | 0.7813              | 5.438             |
| 843786   | 0.2087          | 0.07613                | 0.3345           | 0.8902              | 2.217             |
|          | area_se         | smoothness_se          | compactness_se   | concavity_se        | concave.points_se |
| 842302   | 153.40          | 0.006399               | 0.04904          | 0.05373             | 0.01587           |
| 842517   | 74.08           | 0.005225               | 0.01308          | 0.01860             | 0.01340           |
| 84300903 | 94.03           | 0.006150               | 0.04006          | 0.03832             | 0.02058           |
| 84348301 | 27.23           | 0.009110               | 0.07458          | 0.05661             | 0.01867           |
| 84358402 | 94.44           | 0.011490               | 0.02461          | 0.05688             | 0.01885           |
| 843786   | 27.19           | 0.007510               | 0.03345          | 0.03672             | 0.01137           |
|          | symmetry_se     | fractal_dimension_se   | radius_worst     | texture_worst       |                   |
| 842302   | 0.03003         | 0.006193               | 25.38            | 17.33               |                   |
| 842517   | 0.01389         | 0.003532               | 24.99            | 23.41               |                   |
| 84300903 | 0.02250         | 0.004571               | 23.57            | 25.53               |                   |
| 84348301 | 0.05963         | 0.009208               | 14.91            | 26.50               |                   |
| 84358402 | 0.01756         | 0.005115               | 22.54            | 16.67               |                   |
| 843786   | 0.02165         | 0.005082               | 15.47            | 23.75               |                   |
|          | perimeter_worst | area_worst             | smoothness_worst | compactness_worst   |                   |
| 842302   | 184.60          | 2019.0                 | 0.1622           | 0.6656              |                   |
| 842517   | 158.80          | 1956.0                 | 0.1238           | 0.1866              |                   |
| 84300903 | 152.50          | 1709.0                 | 0.1444           | 0.4245              |                   |
| 84348301 | 98.87           | 567.7                  | 0.2098           | 0.8663              |                   |
| 84358402 | 152.20          | 1575.0                 | 0.1374           | 0.2050              |                   |
| 843786   | 103.40          | 741.6                  | 0.1791           | 0.5249              |                   |
|          | concavity_worst | concave.points_worst   | symmetry_worst   |                     |                   |
| 842302   | 0.7119          | 0.2654                 | 0.4601           |                     |                   |
| 842517   | 0.2416          | 0.1860                 | 0.2750           |                     |                   |
| 84300903 | 0.4504          | 0.2430                 | 0.3613           |                     |                   |
| 84348301 | 0.6869          | 0.2575                 | 0.6638           |                     |                   |
| 84358402 | 0.4000          | 0.1625                 | 0.2364           |                     |                   |
| 843786   | 0.5355          | 0.1741                 | 0.3985           |                     |                   |

|          | fractal_dimension_worst |
|----------|-------------------------|
| 842302   | 0.11890                 |
| 842517   | 0.08902                 |
| 84300903 | 0.08758                 |
| 84348301 | 0.17300                 |
| 84358402 | 0.07678                 |
| 843786   | 0.12440                 |

The first column **diagnosis** is the expert opinion on the sample(i.e. patient FNA).

```
head(wisc.df$diagnosis)
```

```
[1] "M" "M" "M" "M" "M" "M"
```

Remove the diagnosis from data for subsequent analysis

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

```
[1] 569 30
```

Store the diagnosis as a vector for use later when we compare our results to those from experts in the field

```
# Create diagnosis vector for later
diagnosis <- factor(wisc.df$diagnosis)
```

Q1. How many observations are in this dataset?

There are 569 observations/patients in the dataset

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

There are 212 malignant diagnosis

```
#finds out how many Benign/ Malignant there are in the sample
table(diagnosis)
```

```
diagnosis
  B   M
357 212
```

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

```
variables <- colnames(wisc.data)
```

```
grep("_mean", variables)
```

```
[1] 1 2 3 4 5 6 7 8 9 10
```

```
results <- grep("_mean", variables, value = T)
```

```
results
```

```
[1] "radius_mean"      "texture_mean"      "perimeter_mean"
[4] "area_mean"        "smoothness_mean"   "compactness_mean"
[7] "concavity_mean"   "concave.points_mean" "symmetry_mean"
[10] "fractal_dimension_mean"
```

```
total_use_of_mean <- length(results)
```

```
total_use_of_mean
```

```
[1] 10
```

## Principal Component Analysis (PCA)

In general we want to scale and center our data prior to PCA, to ensure that each features contribute equally to the analysis

The `prcomp()` function to do PCA has a `scale = FALSE` default.

We almost always want to set `scale = True` in `prcomp()`, so that certain columns/variables with large standard deviation and mean won't impact when compared to others just because the units of measurement are on different scales.

```
# Check column means and standard deviations
```

```
wisc_mean <- colMeans(wisc.data)
```

```
wisc_sd <- apply(wisc.data, 2, sd)
```

```
max(wisc_mean)
```

```
[1] 880.5831
```

```
min(wisc_mean)
```

```
[1] 0.003794904
```

```
max(wisc_sd)
```

```
[1] 569.357
```

```
min(wisc_sd)
```

```
[1] 0.002646071
```

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

Importance of components:

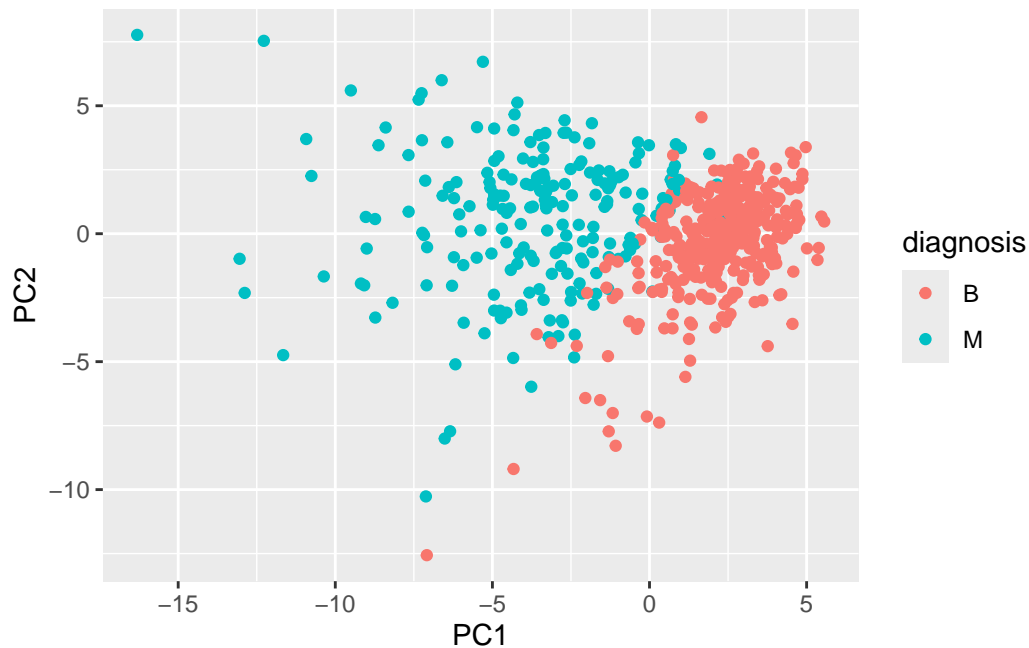
|                        | PC1     | PC2     | PC3     | PC4     | PC5     | PC6     | PC7     |
|------------------------|---------|---------|---------|---------|---------|---------|---------|
| Standard deviation     | 3.6444  | 2.3857  | 1.67867 | 1.40735 | 1.28403 | 1.09880 | 0.82172 |
| Proportion 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 |         |         |         |         |         |

The main PC result figure is called a “score plot” or “PC plot” or “ordination plot”...

```
library(ggplot2)

# wisc.pr$x

ggplot(wisc.pr$x)+ aes(x = PC1, y = PC2, col = diagnosis) + geom_point()
```



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

44.27%

```
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 |         |         |         |         |         |

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

3 PCs

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

7 PCs

Create a biplot of the wisc.pr using the biplot() function. > Q7. What stands out to you about this plot? Is it easy or difficult to understand? Why?

Very packed; It's hard to understand, because there are so many dimensions involved.

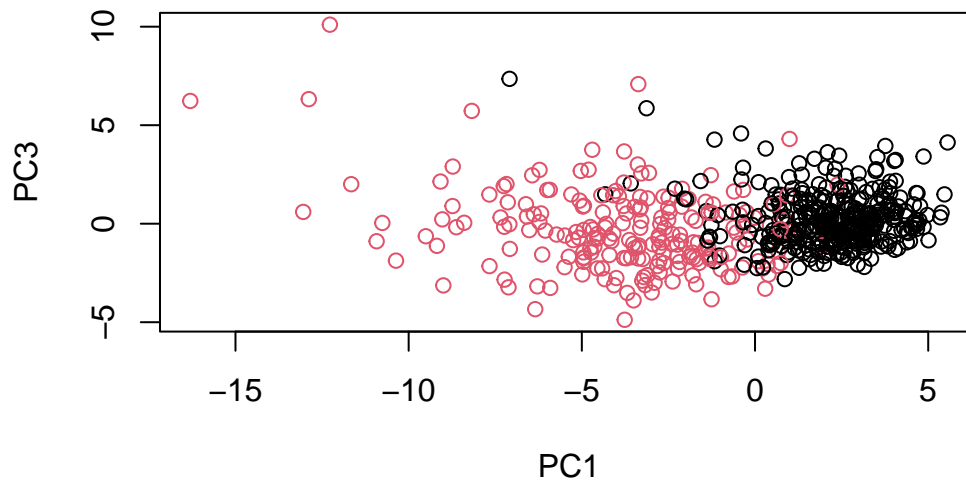
```
biplot(wisc.pr)
```





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

```
# Repeat for components 1 and 3
plot(wisc.pr$x[,c(1,3)], col = diagnosis,
     xlab = "PC1", ylab = "PC3")
```



## Variance Explained

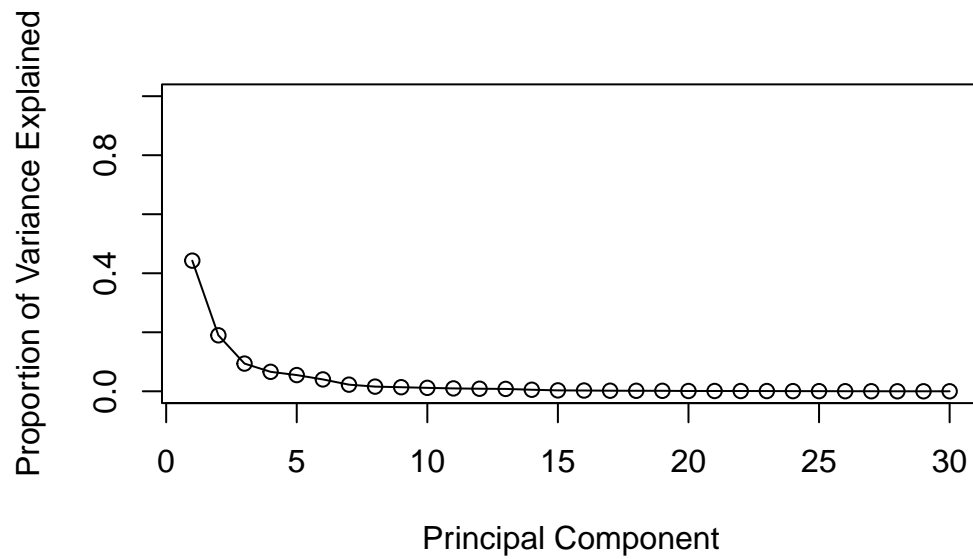
```
# Calculate variance of each component
pr.var <- wisc.pr$sdev^2
head(pr.var)
```

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

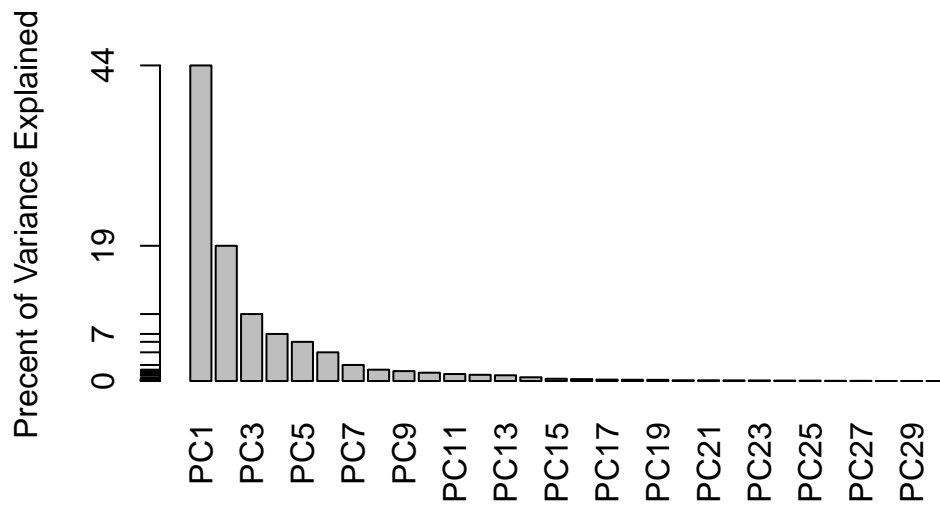
```
tot_variance = sum(pr.var)
# Variance explained by each principal component: pve
pve <- pr.var / tot_variance

# Plot variance explained for each principal component
```

```
plot(pve, xlab = "Principal Component",
     ylab = "Proportion of Variance Explained",
     ylim = c(0, 1), type = "o")
```



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



## Communicate PCA Results

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`?

```
PC1_loading <- wisc.pr$rotation[,1]
PC1_loading["concave.points_mean"]
```

```
concave.points_mean
-0.2608538
```

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

5Pcs

```
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 |         |         |         |         |         |

## Hierarchical Clustering

The goal of this section is to do hierarchical clustering of the original data. Recall from class that this type of clustering does not assume in advance the number of natural groups that exist in the data.

As part of the preparation for hierarchical clustering, the distance between all pairs of observations are computed. Furthermore, there are different ways to link clusters together, with single, complete, and average being the most common linkage methods.

First scale the `wisc.data` data and assign the result to `data.scaled`.

```
# 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`.

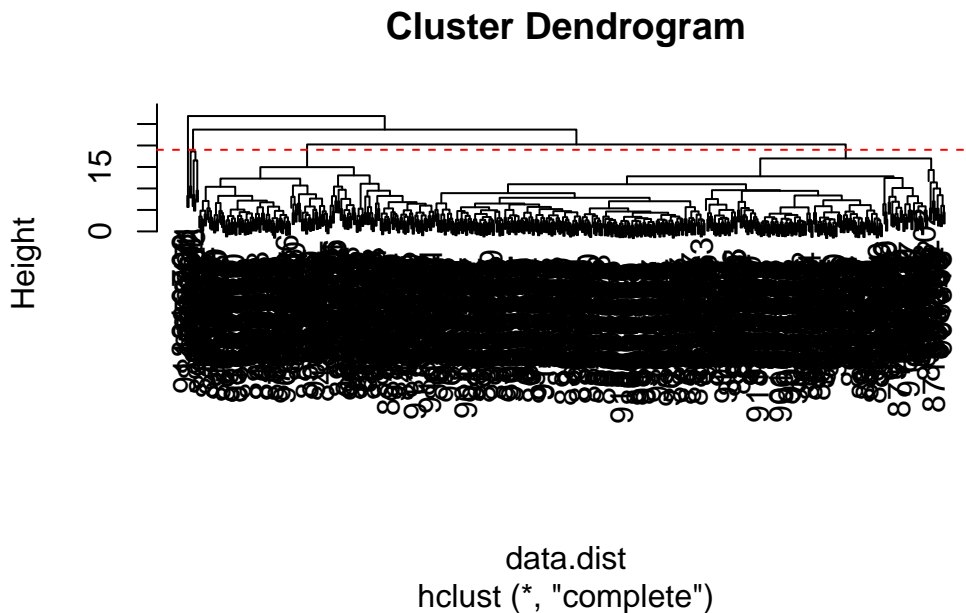
```
data.dist <- dist(data.scaled)
```

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, method = "complete")
```

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

```
plot(wisc.hclust)  
abline(h = 19, col="red", lty=2)
```



### Selecting number of clusters

In this section, you will compare the outputs from your hierarchical clustering model to the actual diagnoses. Normally when performing unsupervised learning like this, a target variable (i.e. known answer or labels) isn't available. We do have it with this dataset, however, so it can be used to check the performance of the clustering model.

When performing supervised learning - that is, when you're trying to predict some target variable of interest and that target variable is available in the original data - using clustering to create new features may or may not improve the performance of the final model.

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

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

```
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?

maybe 6

## Combining methods (PCA & Clustering)

- Most important

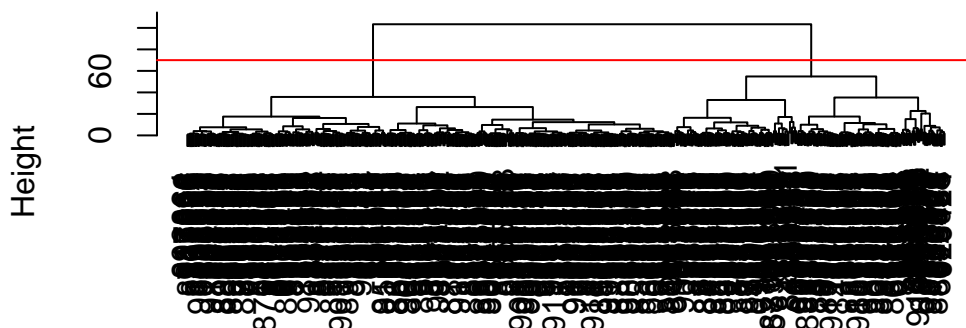
Clustering the original data was not very productive. The PCA results look promising. Here we combine these methods by clustering from our PCA results. In other words, “clustering in PC space”..

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

View the tree..

```
plot(wisc.pr.hclust)
abline(h=70, col = "red")
```

## Cluster Dendrogram



```
dist.pc
hclust (*, "ward.D2")
```

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

```
## 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)
```

```
# Compare to actual diagnoses
table(wisc.pr.hclust.clusters, diagnosis)
```

```

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

To get our clustering membership vector (i.e. our main clustering result) we “cut the tree at a desired height or to yield a desired number of”k” groups.

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

```
grps
  1  2
216 353
```

How does this clustering grps compare to the expert diagnosis

```
# Compare to actual diagnoses
table(grps, diagnosis)
```

```
      diagnosis
grps   B    M
  1  28 188
  2 329  24
```

Sensitivity :  $TP / (TP + FN)$  Specificity:  $TN / (TN + FP)$

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.

K-means clustering is skipped as an optional section

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

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

## 7 Prediction

We can use our PCA model for prediction with new input patient samples

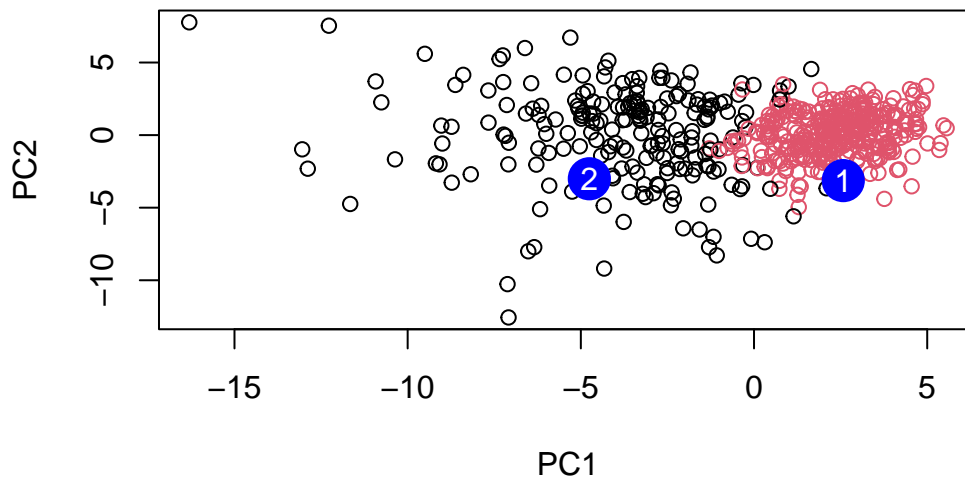
```
#url <- "new_samples.csv"
url <- "https://tinyurl.com/new-samples-CSV"
new <- read.csv(url)
npc <- predict(wisc.pr, newdata=new)
head(npc)
```



|      | 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 |             |              |            |

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

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?

Group 1

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

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

table(grps, diagnosis)
```

|      | diagnosis |     |
|------|-----------|-----|
| grps | B         | M   |
| 1    | 28        | 188 |
| 2    | 329       | 24  |