# Class 8: PCA Mini Project

Woocheol Kim (PID: A16998418)

It is important to consider scalling.

```r
colMeans(mtcars)
```

```
       mpg        cyl       disp         hp       drat         wt       qsec
 20.090625   6.187500 230.721875 146.687500   3.596563   3.217250  17.848750
        vs         am       gear       carb
  0.437500   0.406250   3.687500   2.812500
```

```r
apply(mtcars, 2, sd)
```

```
       mpg        cyl       disp         hp       drat         wt
 6.0269481  1.7859216 123.9386938  68.5628685   0.5346787   0.9784574
      qsec         vs         am       gear       carb
 1.7869432  0.5040161   0.4989909   0.7378041   1.6152000
```

```r
x <- scale(mtcars)
head(x)
```

```
                          mpg        cyl        disp         hp       drat
Mazda RX4           0.1508848 -0.1049878 -0.57061982 -0.5350928   0.5675137
Mazda RX4 Wag       0.1508848 -0.1049878 -0.57061982 -0.5350928   0.5675137
Datsun 710          0.4495434 -1.2248578 -0.99018209 -0.7830405   0.4739996
Hornet 4 Drive      0.2172534 -0.1049878  0.22009369 -0.5350928  -0.9661175
Hornet Sportabout  -0.2307345  1.0148821  1.04308123  0.4129422  -0.8351978
Valiant            -0.3302874 -0.1049878 -0.04616698 -0.6080186  -1.5646078
                           wt       qsec         vs         am       gear
Mazda RX4          -0.610399567 -0.7771651 -0.8680278  1.1899014  0.4235542
Mazda RX4 Wag      -0.349785269 -0.4637808 -0.8680278  1.1899014  0.4235542
Datsun 710         -0.917004624  0.4260068  1.1160357  1.1899014  0.4235542
```

1

```
Hornet 4 Drive      -0.002299538  0.8904872  1.1160357 -0.8141431 -0.9318192
Hornet Sportabout    0.227654255 -0.4637808 -0.8680278 -0.8141431 -0.9318192
Valiant              0.248094592  1.3269868  1.1160357 -0.8141431 -0.9318192
                          carb
Mazda RX4            0.7352031
Mazda RX4 Wag        0.7352031
Datsun 710          -1.1221521
Hornet 4 Drive      -1.1221521
Hornet Sportabout   -0.5030337
Valiant             -1.1221521
```

```r
round(colMeans(x),2)
```

```
 mpg  cyl disp   hp drat   wt qsec   vs   am gear carb
   0    0    0    0    0    0    0    0    0    0    0
```

Key-point: It is usually always a good idea to scale your data before to PCA...

**Preparing the Data**

## Save your input data file into your Project directory

```r
fna.data <- "WisconsinCancer.csv"
```

## Complete the following code to input the data and store as wisc.df

```r
wisc.df <- read.csv(fna.data, row.names=1)
```

```r
View(wisc.df)
```

```r
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
```

| | fractal_dimension_worst | | |
|---|---|---|---|
| 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 |

```r
wisc.data <- wisc.df[,-1]
```

```r
diagnosis <- wisc.df[,1]
table(diagnosis)
```

```
diagnosis
  B   M
357 212
```

Remove this first `diagnosis` column from the dataset as I don't want to pass this to PCA etc It is essentially the expert "answer" that we will compare our analysis results to.

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

| | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean |
|---|---|---|---|---|---|
| 842302 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 |
| 842517 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 |
| 84300903 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 |
| 84348301 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 |
| 84358402 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 |
| 843786 | 12.45 | 15.70 | 82.57 | 477.1 | 0.12780 |

| | compactness_mean | concavity_mean | concave.points_mean | symmetry_mean |
|---|---|---|---|---|
| 842302 | 0.27760 | 0.3001 | 0.14710 | 0.2419 |
| 842517 | 0.07864 | 0.0869 | 0.07017 | 0.1812 |
| 84300903 | 0.15990 | 0.1974 | 0.12790 | 0.2069 |
| 84348301 | 0.28390 | 0.2414 | 0.10520 | 0.2597 |
| 84358402 | 0.13280 | 0.1980 | 0.10430 | 0.1809 |
| 843786 | 0.17000 | 0.1578 | 0.08089 | 0.2087 |

| | fractal_dimension_mean | radius_se | texture_se | perimeter_se | area_se |
|---|---|---|---|---|---|

|          |         |        |        |       |        |
|----------|---------|--------|--------|-------|--------|
| 842302   | 0.07871 | 1.0950 | 0.9053 | 8.589 | 153.40 |
| 842517   | 0.05667 | 0.5435 | 0.7339 | 3.398 | 74.08  |
| 84300903 | 0.05999 | 0.7456 | 0.7869 | 4.585 | 94.03  |
| 84348301 | 0.09744 | 0.4956 | 1.1560 | 3.445 | 27.23  |
| 84358402 | 0.05883 | 0.7572 | 0.7813 | 5.438 | 94.44  |
| 843786   | 0.07613 | 0.3345 | 0.8902 | 2.217 | 27.19  |

|          | smoothness_se | compactness_se | concavity_se | concave.points_se |
|----------|---------------|----------------|--------------|-------------------|
| 842302   | 0.006399      | 0.04904        | 0.05373      | 0.01587           |
| 842517   | 0.005225      | 0.01308        | 0.01860      | 0.01340           |
| 84300903 | 0.006150      | 0.04006        | 0.03832      | 0.02058           |
| 84348301 | 0.009110      | 0.07458        | 0.05661      | 0.01867           |
| 84358402 | 0.011490      | 0.02461        | 0.05688      | 0.01885           |
| 843786   | 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                 |

```
dim(wisc.df)
```

```
[1] 569  31
```

```
nrow(wisc.df)
```

```
[1] 569
```

```
length( grep("_mean", colnames(wisc.data)) )
```

```
[1] 10
```

## Exploratory data analysis

Q1. How many observations are in this dataset?

569

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

212

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

10

```
wisc.pr <- prcomp( wisc.data, scale=T )
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
```

Main "PC score plot", "PC1 vs PC2 plot"

See what is in our PCA result object:

```
attributes(wisc.pr)
```

```
$names
[1] "sdev"     "rotation" "center"   "scale"    "x"

$class
[1] "prcomp"
```

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

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

0.4427

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

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

Q7. What stands out to you about this plot? Is it easy or difficult to understand? Why?

There is so much on this plot as the it is very hard to read and understand what is going on. There is so much datasets that there is all the information in this one plot when we do not need it.

```
# Scatter plot observations by components 1 and 2
plot( wisc.pr$x[,1], wisc.pr$x[,2] , col = as.factor(diagnosis) ,
      xlab = "PC1", ylab = "PC2")
```

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

From these plots, I noticed that PC1 roughly remains the same, and the two plots are pretty similar. PC3 does go up to 10 though but mostly stays the same.

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

```r
x <- summary(wisc.pr)
plot(x$importance[2,], typ="b")
```

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?

concave.points_mean is -0.25840048

```
wisc.pr$rotation[,1]
```

| radius_mean | texture_mean | perimeter_mean |
|---|---|---|
| -0.21890244 | -0.10372458 | -0.22753729 |
| area_mean | smoothness_mean | compactness_mean |
| -0.22099499 | -0.14258969 | -0.23928535 |
| concavity_mean | concave.points_mean | symmetry_mean |
| -0.25840048 | -0.26085376 | -0.13816696 |
| fractal_dimension_mean | radius_se | texture_se |
| -0.06436335 | -0.20597878 | -0.01742803 |
| perimeter_se | area_se | smoothness_se |
| -0.21132592 | -0.20286964 | -0.01453145 |
| compactness_se | concavity_se | concave.points_se |
| -0.17039345 | -0.15358979 | -0.18341740 |
| symmetry_se | fractal_dimension_se | radius_worst |
| -0.04249842 | -0.10256832 | -0.22799663 |
| texture_worst | perimeter_worst | area_worst |
| -0.10446933 | -0.23663968 | -0.22487053 |
| smoothness_worst | compactness_worst | concavity_worst |
| -0.12795256 | -0.21009588 | -0.22876753 |
| concave.points_worst | symmetry_worst | fractal_dimension_worst |
| -0.25088597 | -0.12290456 | -0.13178394 |

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

The minumum number of principal components required to explain 80% of the variance of the data is PC 5.

## Combine PCA and clustering

Our PCA results were in `wisc.pr$x`

```
# distance matrix from PCA result
d <- dist( wisc.pr$x[,1:3] )
hc <- hclust(d, method = "ward.D2")
plot(hc)
```
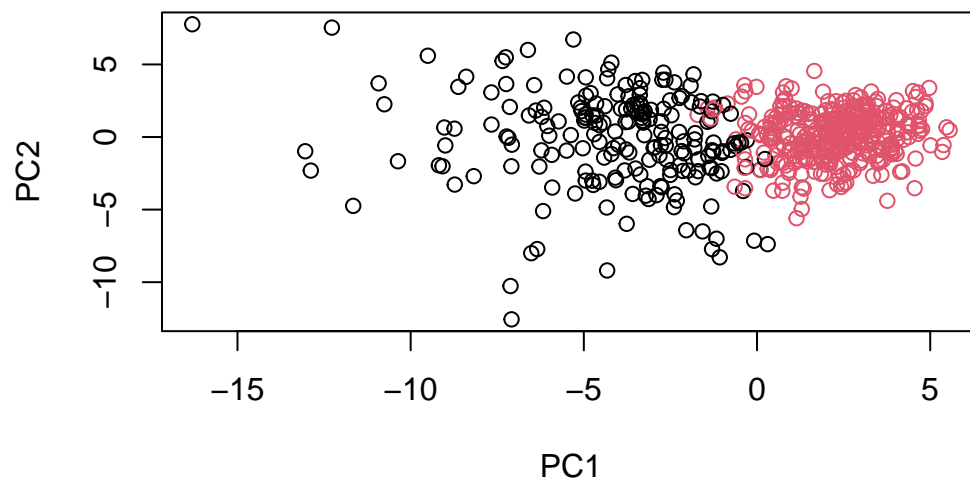
## Cluster Dendrogram



Height

d
hclust (*, "ward.D2")

Cut tree into two groups/branches/clusters...

```r
grps <- cutree(hc, k=2)
```

```r
plot(wisc.pr$x, col=grps)
```

Compare my clustering result (my `grps`) to the expert `diagnosis`

```
table(diagnosis)
```

```
diagnosis
  B   M
357 212
```

```
table(grps)
```

```
grps
  1   2
203 366
```

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

```
         grps
diagnosis   1   2
        B  24 333
        M 179  33
```

```
# Scale the wisc.data data using the "scale()" function
data.scaled <- scale(wisc.data)
```

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

```
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?

Height is 19.

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

## Cluster Dendrogram



data.dist
hclust (*, "complete")

```
wisc.hclust.clusters <- cutree(wisc.hclust, k=10)
table(wisc.hclust.clusters, diagnosis)
```

```
                     diagnosis
wisc.hclust.clusters   B   M
                  1   12  86
```

```
 2    0  59
 3    0   3
 4  331  39
 5    0  20
 6    2   0
 7   12   0
 8    0   2
 9    0   2
10    0   1
```

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

A better cluster vs diagnoses match by cutting into a different number of clusters would be clust 2 because it has the highest count of M found.

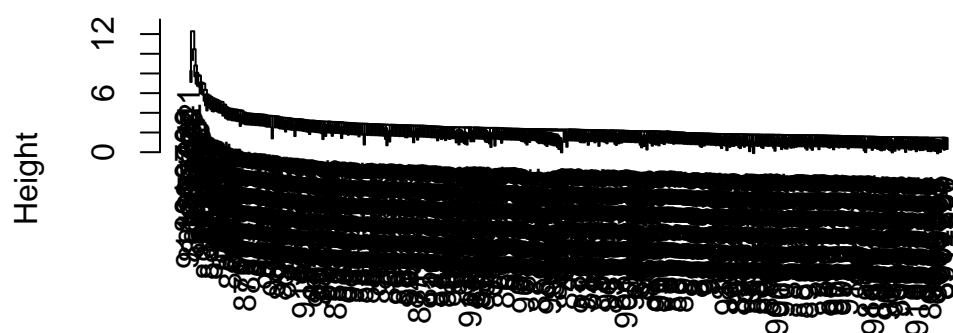> Q13. Which method gives your favorite results for the same data.dist dataset? Explain your reasoning.

"Ward.D2" gives my favorite results for the same data.dist dataset because the branches is the most neat and clearly can see the separation between the points.

```r
wisc.hclust <- hclust(data.dist, method = "single")

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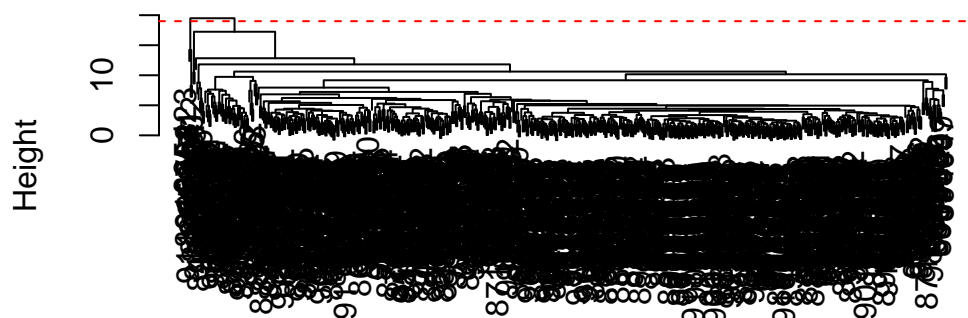
**Cluster Dendrogram**



data.dist
hclust (*, "single")

```r
wisc.hclust <- hclust(data.dist, method = "average")

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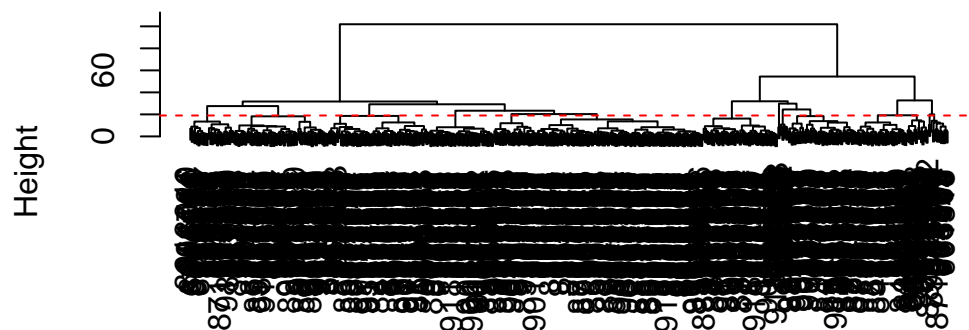
# Cluster Dendrogram



Height

data.dist
hclust (*, "average")

```r
wisc.hclust <- hclust(data.dist, method = "ward.D2")

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

# Cluster Dendrogram



data.dist
hclust (*, "ward.D2")