# Mini-project

## Julia Napoli

### 10/27/2021

```
fna.data <- read.csv("WisconsinCancer.csv")
wisc.df <- data.frame(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
```

Let's make sure we don't include the diagnosis column since we won't be needing this for our analysis.

```
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
```

And let's create a diagnosis vector for later. . .

```
diagnosis <- fna.data$diagnosis
# diagnosis <- (data.frame(fna.data, row.names=1))[,1]
diagnosis
```

```
##   [1] "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M"
##  [19] "M" "B" "B" "B" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M"
##  [37] "M" "B" "M" "M" "M" "M" "M" "M" "M" "M" "B" "M" "B" "B" "B" "B" "B" "M"
##  [55] "M" "B" "M" "M" "B" "B" "B" "B" "M" "B" "M" "M" "B" "B" "B" "B" "M" "B"
##  [73] "M" "M" "B" "M" "B" "M" "M" "B" "B" "B" "M" "M" "B" "M" "M" "M" "B" "B"
##  [91] "B" "M" "B" "B" "M" "M" "B" "B" "B" "M" "M" "B" "B" "B" "B" "M" "B" "B"
## [109] "M" "B" "B" "B" "B" "B" "B" "B" "B" "M" "M" "M" "B" "M" "M" "B" "B" "B"
## [127] "M" "M" "B" "M" "B" "M" "M" "B" "M" "M" "B" "B" "M" "B" "B" "M" "B" "B"
## [145] "B" "B" "M" "B" "B" "B" "B" "B" "B" "B" "B" "B" "M" "B" "B" "B" "B" "M"
## [163] "M" "B" "M" "B" "B" "M" "M" "B" "B" "M" "M" "B" "B" "B" "B" "M" "B" "B"
## [181] "M" "M" "M" "B" "M" "B" "M" "B" "B" "B" "M" "B" "B" "M" "M" "B" "M" "M"
## [199] "M" "M" "B" "M" "M" "M" "B" "M" "B" "M" "B" "B" "M" "B" "M" "M" "M" "M"
## [217] "B" "B" "M" "M" "B" "B" "B" "M" "B" "B" "B" "B" "B" "M" "M" "B" "B" "M"
## [235] "B" "B" "M" "M" "B" "M" "B" "B" "B" "B" "M" "B" "B" "B" "B" "B" "M" "B"
## [253] "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "M" "B" "B" "B" "B"
## [271] "B" "B" "M" "B" "M" "B" "B" "M" "B" "B" "M" "B" "M" "M" "B" "B" "B" "B"
## [289] "B" "B" "B" "B" "B" "B" "B" "B" "B" "M" "B" "B" "M" "B" "M" "B" "B" "B"
## [307] "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "M" "B" "B" "B" "M" "B" "M"
```

```
## [325] "B" "B" "B" "B" "M" "M" "M" "B" "B" "B" "B" "M" "B" "M" "B" "M" "B" "B"
## [343] "B" "M" "B" "B" "B" "B" "B" "B" "B" "M" "M" "M" "B" "B" "B" "B" "B" "B"
## [361] "B" "B" "B" "B" "B" "M" "M" "B" "M" "M" "M" "B" "M" "M" "B" "B" "B" "B"
## [379] "B" "M" "B" "B" "B" "B" "B" "M" "B" "B" "B" "M" "B" "B" "M" "M" "B" "B"
## [397] "B" "B" "B" "B" "M" "B" "B" "B" "B" "B" "B" "B" "M" "B" "B" "B" "B" "B"
## [415] "M" "B" "B" "M" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "M" "B"
## [433] "M" "M" "B" "M" "B" "B" "B" "B" "B" "M" "B" "B" "M" "B" "M" "B" "B" "M"
## [451] "B" "M" "B" "B" "B" "B" "B" "B" "B" "B" "M" "M" "B" "B" "B" "B" "B" "B"
## [469] "M" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "M" "B" "B" "B" "B" "B" "B"
## [487] "B" "M" "B" "M" "B" "B" "M" "B" "B" "B" "B" "B" "M" "M" "B" "M" "B" "M"
## [505] "B" "B" "B" "B" "B" "M" "B" "B" "M" "B" "M" "B" "M" "M" "B" "B" "B" "M"
## [523] "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "M" "B" "M" "M" "B" "B" "B"
## [541] "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B" "B"
## [559] "B" "B" "B" "B" "M" "M" "M" "M" "M" "M" "B"
```

```
# To double check that I pulled out a vector, we can check using the is.vector function
# Vectors in R are only horizontal, you cannot have a vertical vector in R !
is.vector(diagnosis)
```

```
## [1] TRUE
```

Q1 How many observations are in this dataset?

```
nrow(wisc.data)
```

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

There are a total of 569 observations in this dataset.

Q2 How many of the observations have a malignant diagnosis?

```
sum(diagnosis == "M")
```

```
## [1] 212
```

A total of 212 of the observations have a malignant diagnosis.

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

```
mean_cols <- grep(pattern = "_mean$", x = colnames(wisc.data), value = TRUE)
# Adding value = TRUE returns the matching elements of the grep functions; value = FALSE (default) simp
mean_cols
```

```
##  [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"
```

```r
length(mean_cols)
```

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

There are a total of 10 variables in the data set that are suffixed with "_mean".

Principal Component Analysis

Check the column means and standard deviations to determine if the data should be scaled.

```r
column_means <- colMeans(wisc.data)
std <- apply(wisc.data,2,sd)
```

```r
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
```

```
wisc.pr <- prcomp(wisc.data, scale = TRUE)
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 What proportion of the original variance is captured by the first principal components (PC1)?

44.27%

Q5 How many principal components (PCs) are required to describe at least 70% of the original variance in the data?
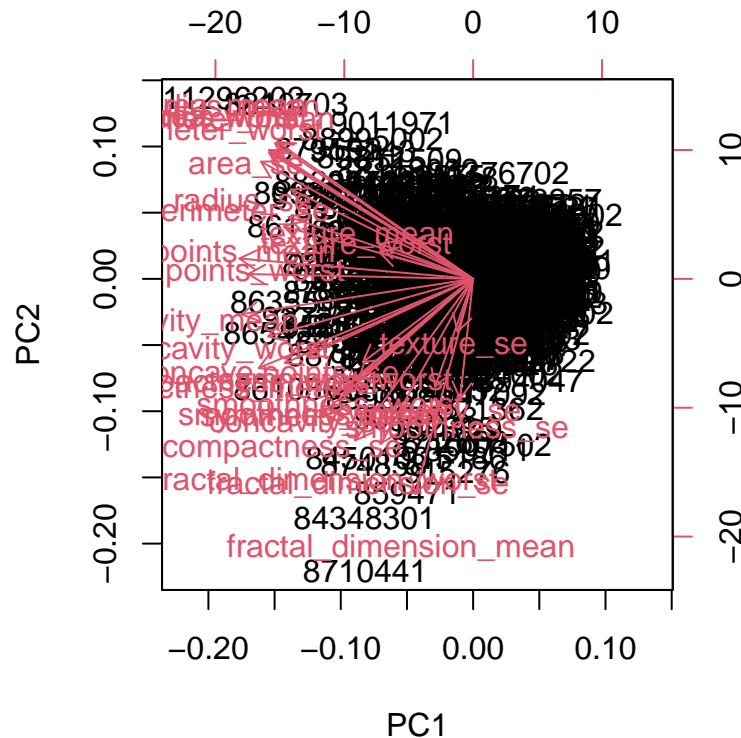
PC1, PC2 & PC3 (3 total components)

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

PC1-PC7 (7 total components)

Interpreting PCA Results

Create a biplot of the `wisc.pr` using the biplot() function
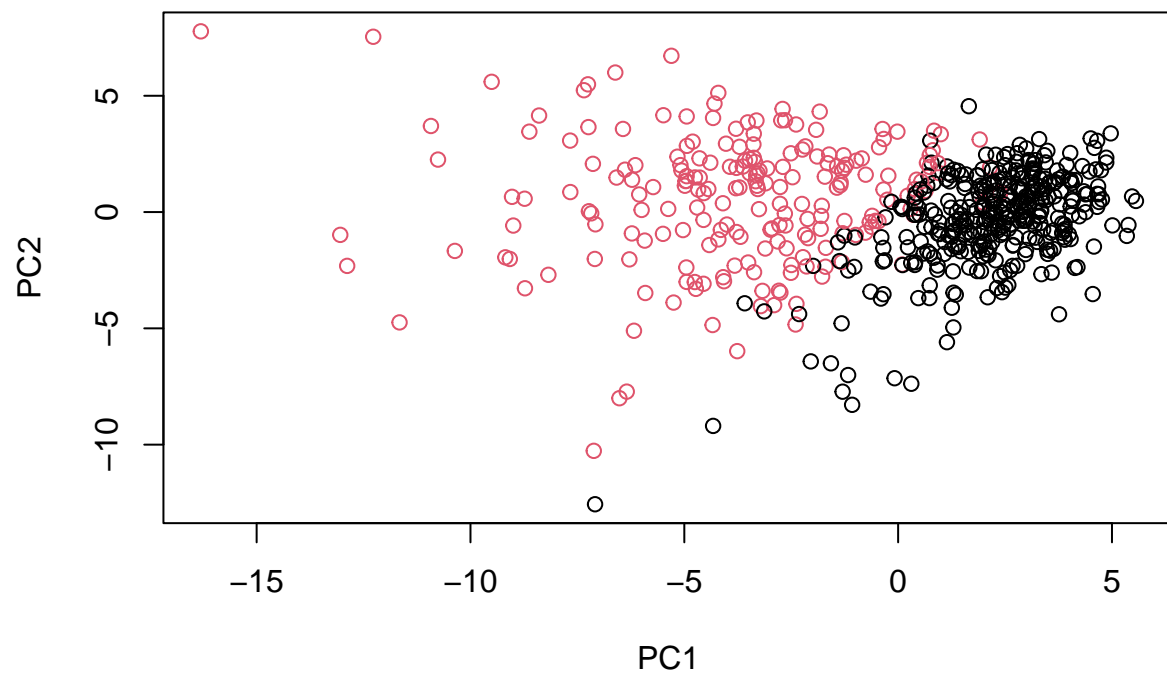
```
biplot(wisc.pr)
```



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

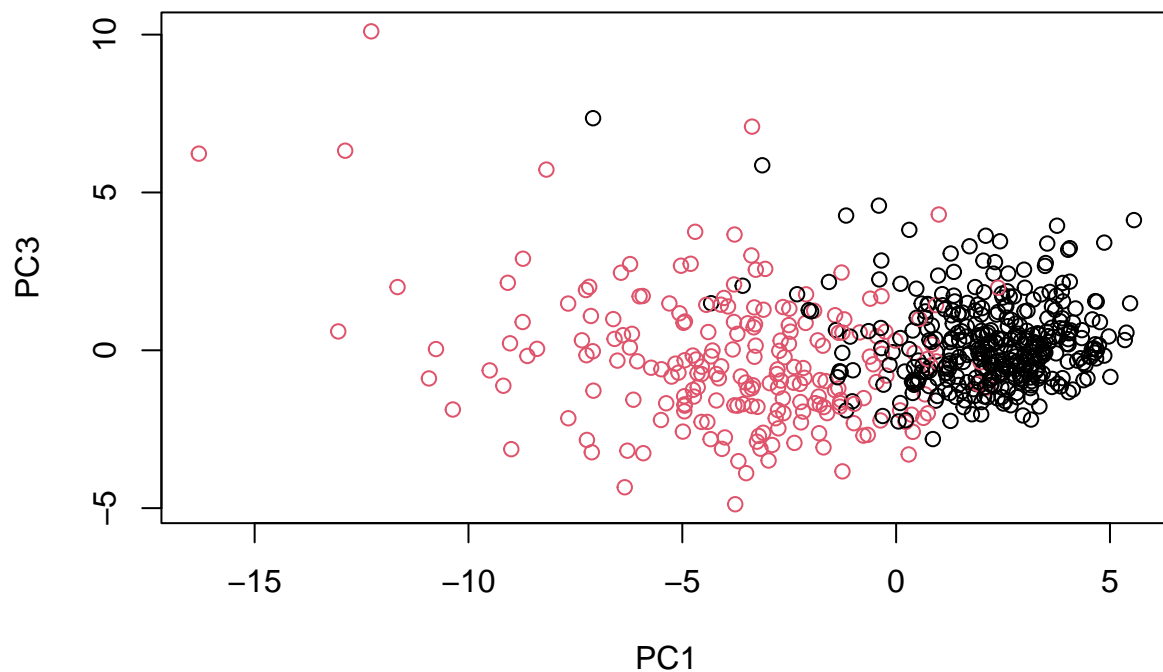It takes a long time to produce and it incredibly difficult to read!

Now let's look at a standard scatter plot of each observation along principal components 1 & 2 and color the points by diagnosis.

```
# In order to use diagnosis as a color we must change it from a character vector to a factor vector!
plot(wisc.pr$x[,1:2], 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?

```
plot(wisc.pr$x[,1],wisc.pr$x[,3], col = as.factor(diagnosis), xlab = "PC1", ylab = "PC3")
```
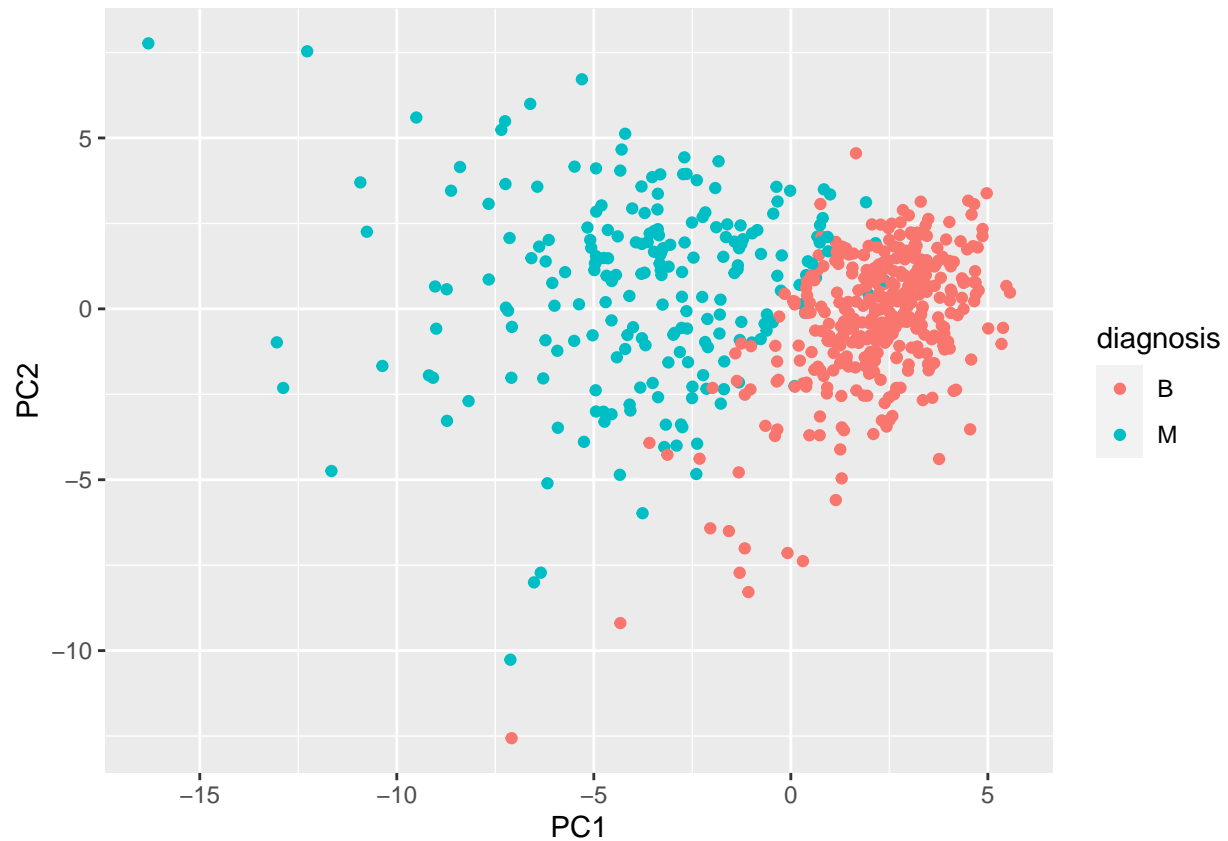
The first plot has a more clear differentiation between the two clusters than the second plot.

Let's now use ggplot2 to make a more fancy figure of the results!

```
diagnosis <- as.factor(diagnosis)
df <- as.data.frame(wisc.pr$x)
df$diagnosis <- diagnosis

# Load ggplot2 package
library(ggplot2)

# Create a scatter plot
ggplot(df) + aes(PC1,PC2, col = diagnosis) + geom_point()
```

## Variance explained

Calculate the variance of each principal 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 the total variance explained of all principal components.

```
pve <- pr.var / sum(pr.var)
plot(pve, xlab = "Principal Component", ylab = "Proportion of Variance Explained", ylim = c(0,1), type =
```

Let's make an alternative scree plot of the same data...

```r
barplot(pve, ylab = "Percent of Variance Explained", names.arg = paste0("PC",1:length(pve)),las = 2, ax
axis(2, at = pve, labels =round(pve,2)*100)
```
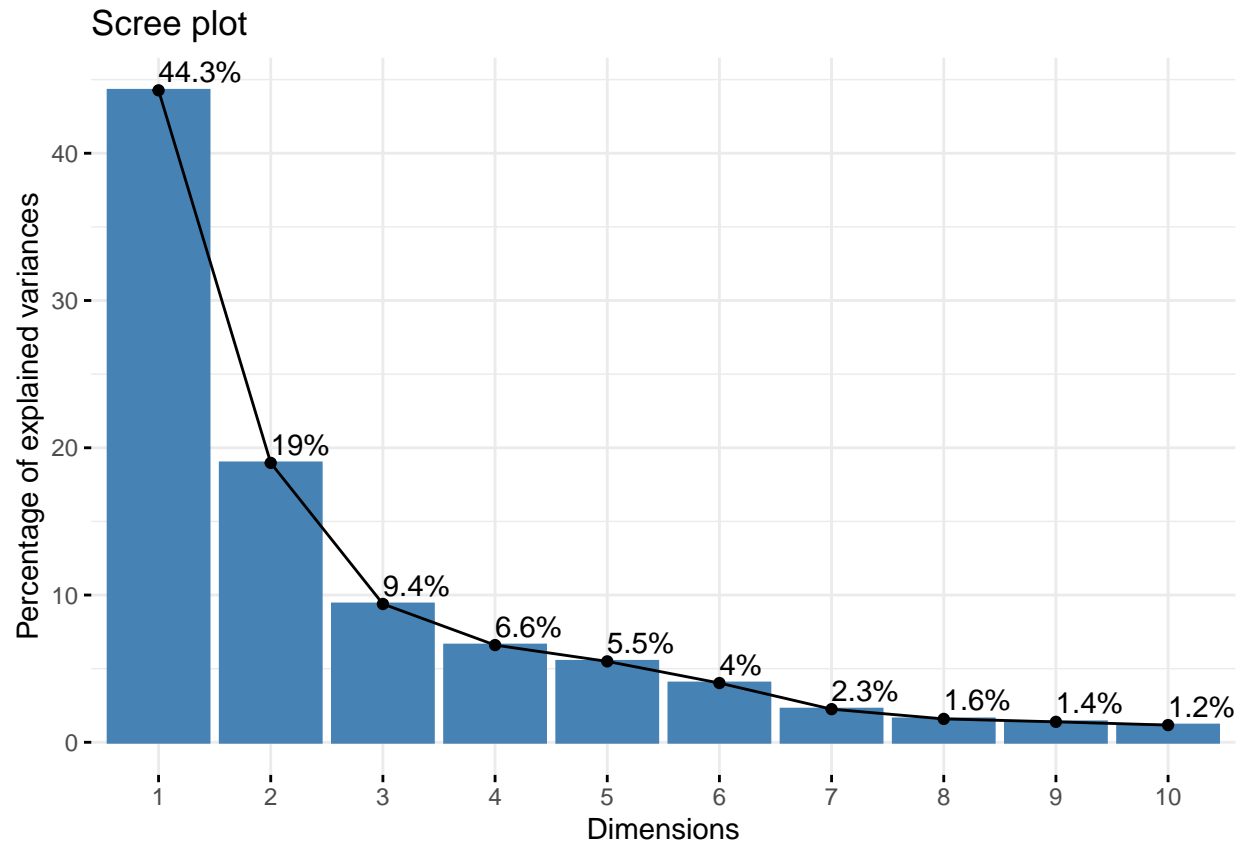
Optional! Checking out the factoextra package from CRAN.

```r
# install.packages("factoextra")
library(factoextra)
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```r
fviz_eig(wisc.pr, addlabels = TRUE)
```

## Scree plot



## Communicating 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`?

```
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
```

-0.26085376

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
```

```
summary_pcr <- summary(wisc.pr)
sum(summary_pcr$importance[3,] <= 0.8)
```

```
## [1] 4
```

PC1-PC5 (cumulative 84.7%), so a total of 5 principal components. [Note: PC1-PC4 covers a cumulative 79.2% variance].

Using code to pull out the answer gives us 4 the answer of 4 principal components under 80% of variance, rounding to the tenth decimal point.

## Heirarchical Clustering

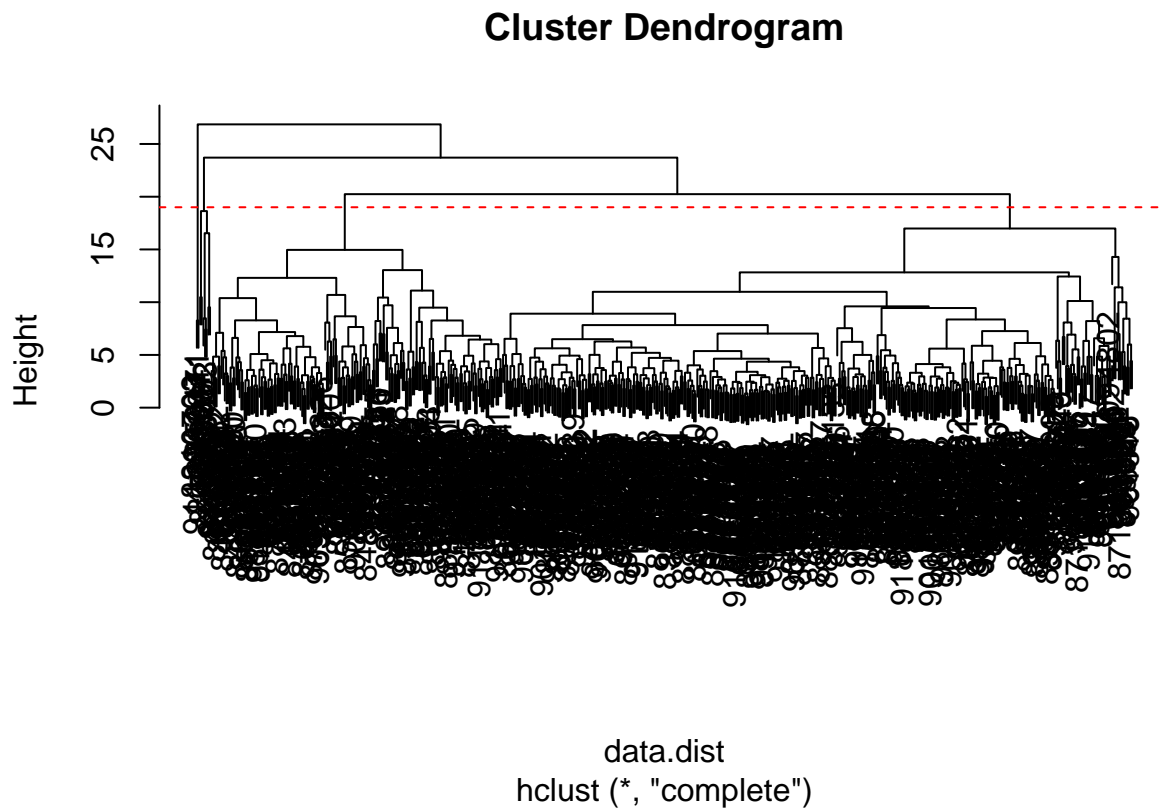Scale the wisc.data using the scale() function

```
data.scaled <- scale(wisc.data)
# Calculate the Euclidean distances between all pairs of observations
```

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

# Create a heirarchical clustering model
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)
```

## Cluster Dendrogram



data.dist
hclust (*, "complete")

## Select number of 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?

```
wisc.hclust.clusters.test <- cutree(wisc.hclust, k=5)
table(wisc.hclust.clusters.test,diagnosis)
```

```
##                          diagnosis
## wisc.hclust.clusters.test   B   M
##                         1  12 165
##                         2   0   5
##                         3 343  40
##                         4   2   0
##                         5   0   2
```
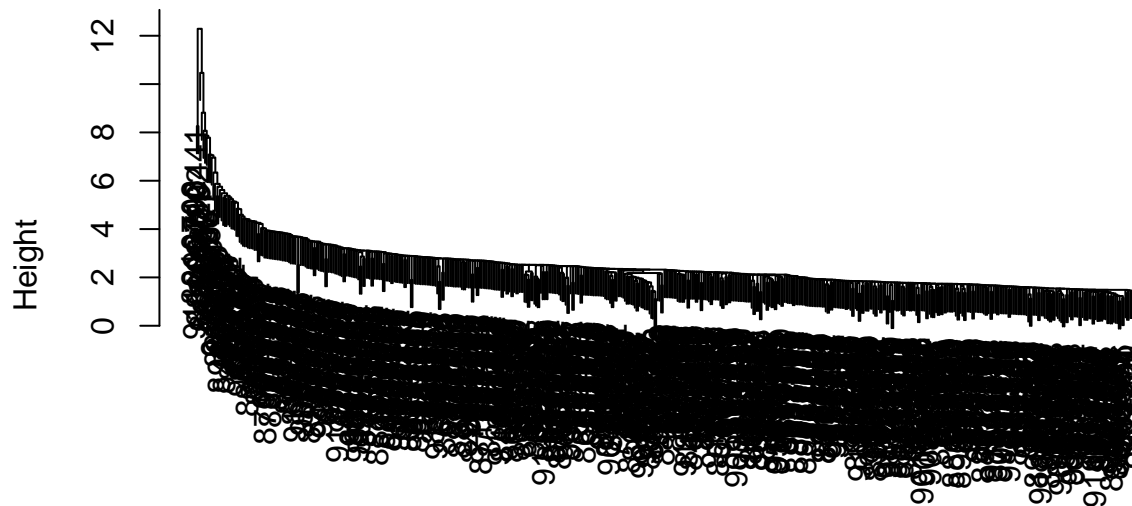
Both k = 4 and k = 5 are good options, because the clustering results are enough to split up the malignant v. benign tumors into their own clusters, while at the same time there are not too many extra clusters being added that aren't really accounting for anything else in the data (k > 5 clusters). Because there isn't a huge difference between k =4 and k=5, I would choose k = 4 to be the most ideal clustering due to its simplicity.

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

```
plot(hclust(data.dist, method = "single"))
```
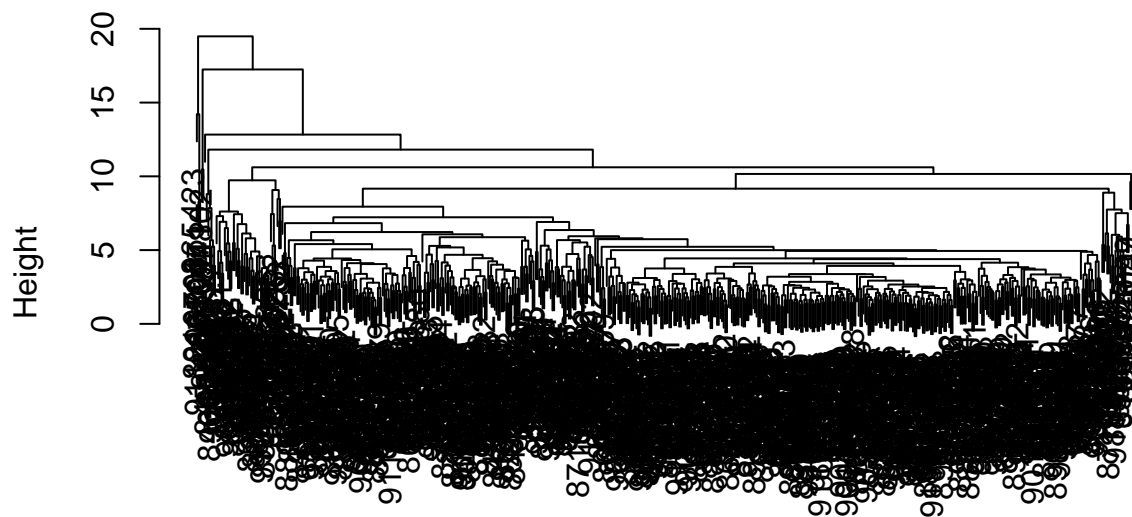
## Cluster Dendrogram



data.dist
hclust (*, "single")

```
plot(hclust(data.dist, method = "average"))
```
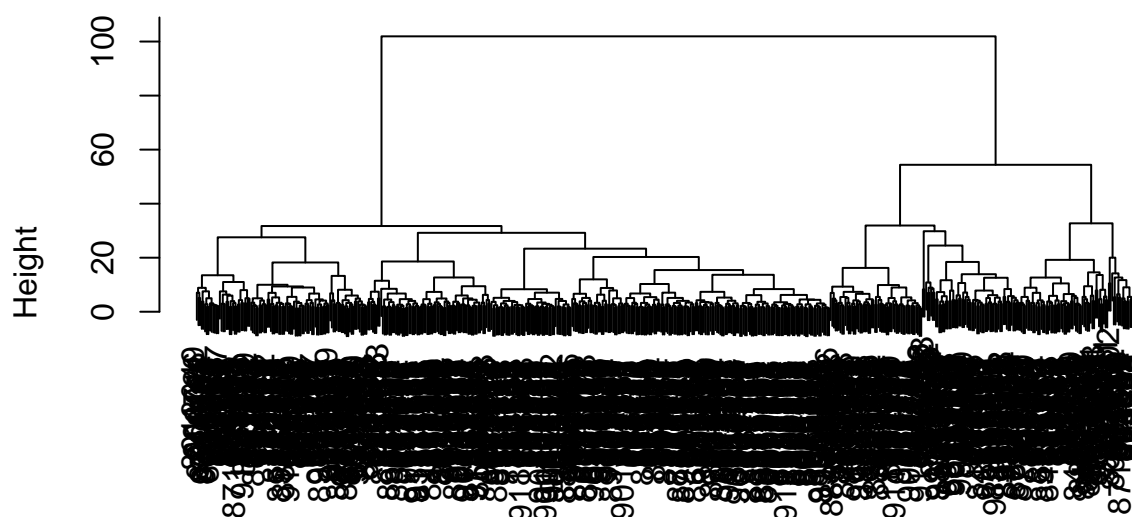
## Cluster Dendrogram



Height

data.dist
hclust (*, "average")

```
plot(hclust(data.dist, method = "ward.D2"))
```

## Cluster Dendrogram



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

Personally, I'm a big fan of the "complete" method. The different methods tell R the different ways to plot the dendrogram. The single one is the worst, because it branches everything off the first singular cluster. Average did fine, but split the clusters out a bit more. Ward. D2. splits everything in two right off the bat, and then goes from there, which really only seems ideal if you for sure have two clusters. Overall, I think complete is the best way to visualize the clusters in the dataset.