# Class_09_mini_project

## Caitriona Brennan

```r
#read.csv("WisconsinCancer.csv")
```

## 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(fna.data)
```

#Q1. How many observations are in this dataset? Answer - 569 observations

```r
#wisc.df
```

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

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

```
##    [1] M M M M M M M M M M M M M M M M M M M B B B M M M M M M M M M M M M M M M
##   [38] B M M M M M M M M B M B B B B B M M B M M B B B B M B M M B B B B M B M M
##   [75] B M B M M B B B M M B M M M B B B M B B M M B B B M M B B B B M B B M B B
##  [112] B B B B B B M M M B M M B B B M M B M B M M B M M B B M B B M B B B B M B
##  [149] B B B B B B B B M B B B B B M M B M B B M M B B M M B B B B M B B M M M B M
##  [186] B M B B B B M B B M M B M M M M B M M M B M B M B B M B M M M M B B M M B B
##  [223] B M B B B B B M M B B M B B M M B M B B B B M B B B B B M B M M M M M M M M
##  [260] M M M M M M M B B B B B B M B M B B M B B M B B M B M M B B B B B B B B B B
##  [297] B M B M B M B B B B B B B B B B B B B M B B B M B M B B B B B M M M B B
##  [334] B B M B M B M B B B M B B B B B B B M M M B B B B B B M B B B B B M M B M M
##  [371] M B M M B B B B B M B B B B B M B B B M B B M M B B B B B M B B B B B B B B
##  [408] B M B B B B B M B B M B B B B B B B B B B B B B M B M M B M B B B B B M B B
##  [445] M B M B B M B M B B B B B B B M M B B B B B B M B B B B B B B B B B B M B B
##  [482] B B B B B B M B M B B M B B B B B M M B M B M B B B B B M B B B M B M B M M
```

```
## [519] B B B M B B B B B B B B B B B M B M M B B B B B B B B B B B B B B B B B B
## [556] B B B B B B B M M M M M M B
## Levels: B M
```

#Q2. How many of the observations have a malignant diagnosis? #Answer - 212

```
table(diagnosis)
```

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

```
summary(wisc.data)
```

```
##   radius_mean      texture_mean    perimeter_mean     area_mean
## Min.   : 6.981   Min.   : 9.71   Min.   : 43.79   Min.   : 143.5
## 1st Qu.:11.700   1st Qu.:16.17   1st Qu.: 75.17   1st Qu.: 420.3
## Median :13.370   Median :18.84   Median : 86.24   Median : 551.1
## Mean   :14.127   Mean   :19.29   Mean   : 91.97   Mean   : 654.9
## 3rd Qu.:15.780   3rd Qu.:21.80   3rd Qu.:104.10   3rd Qu.: 782.7
## Max.   :28.110   Max.   :39.28   Max.   :188.50   Max.   :2501.0
## smoothness_mean   compactness_mean  concavity_mean   concave.points_mean
## Min.   :0.05263   Min.   :0.01938   Min.   :0.00000   Min.   :0.00000
## 1st Qu.:0.08637   1st Qu.:0.06492   1st Qu.:0.02956   1st Qu.:0.02031
## Median :0.09587   Median :0.09263   Median :0.06154   Median :0.03350
## Mean   :0.09636   Mean   :0.10434   Mean   :0.08880   Mean   :0.04892
## 3rd Qu.:0.10530   3rd Qu.:0.13040   3rd Qu.:0.13070   3rd Qu.:0.07400
## Max.   :0.16340   Max.   :0.34540   Max.   :0.42680   Max.   :0.20120
## symmetry_mean    fractal_dimension_mean   radius_se        texture_se
## Min.   :0.1060   Min.   :0.04996      Min.   :0.1115   Min.   :0.3602
## 1st Qu.:0.1619   1st Qu.:0.05770      1st Qu.:0.2324   1st Qu.:0.8339
## Median :0.1792   Median :0.06154      Median :0.3242   Median :1.1080
## Mean   :0.1812   Mean   :0.06280      Mean   :0.4052   Mean   :1.2169
## 3rd Qu.:0.1957   3rd Qu.:0.06612      3rd Qu.:0.4789   3rd Qu.:1.4740
## Max.   :0.3040   Max.   :0.09744      Max.   :2.8730   Max.   :4.8850
##  perimeter_se       area_se        smoothness_se     compactness_se
## Min.   : 0.757   Min.   :  6.802   Min.   :0.001713   Min.   :0.002252
## 1st Qu.: 1.606   1st Qu.: 17.850   1st Qu.:0.005169   1st Qu.:0.013080
## Median : 2.287   Median : 24.530   Median :0.006380   Median :0.020450
## Mean   : 2.866   Mean   : 40.337   Mean   :0.007041   Mean   :0.025478
## 3rd Qu.: 3.357   3rd Qu.: 45.190   3rd Qu.:0.008146   3rd Qu.:0.032450
## Max.   :21.980   Max.   :542.200   Max.   :0.031130   Max.   :0.135400
##  concavity_se      concave.points_se   symmetry_se      fractal_dimension_se
## Min.   :0.00000   Min.   :0.000000   Min.   :0.007882   Min.   :0.0008948
## 1st Qu.:0.01509   1st Qu.:0.007638   1st Qu.:0.015160   1st Qu.:0.0022480
## Median :0.02589   Median :0.010930   Median :0.018730   Median :0.0031870
## Mean   :0.03189   Mean   :0.011796   Mean   :0.020542   Mean   :0.0037949
## 3rd Qu.:0.04205   3rd Qu.:0.014710   3rd Qu.:0.023480   3rd Qu.:0.0045580
## Max.   :0.39600   Max.   :0.052790   Max.   :0.078950   Max.   :0.0298400
##  radius_worst   texture_worst   perimeter_worst    area_worst
## Min.   : 7.93   Min.   :12.02   Min.   : 50.41   Min.   : 185.2
## 1st Qu.:13.01   1st Qu.:21.08   1st Qu.: 84.11   1st Qu.: 515.3
```

```
##  Median :14.97    Median :25.41    Median : 97.66    Median : 686.5
##  Mean   :16.27    Mean   :25.68    Mean   :107.26    Mean   : 880.6
##  3rd Qu.:18.79    3rd Qu.:29.72    3rd Qu.:125.40    3rd Qu.:1084.0
##  Max.   :36.04    Max.   :49.54    Max.   :251.20    Max.   :4254.0
##  smoothness_worst  compactness_worst concavity_worst  concave.points_worst
##  Min.   :0.07117   Min.   :0.02729   Min.   :0.0000   Min.   :0.00000
##  1st Qu.:0.11660   1st Qu.:0.14720   1st Qu.:0.1145   1st Qu.:0.06493
##  Median :0.13130   Median :0.21190   Median :0.2267   Median :0.09993
##  Mean   :0.13237   Mean   :0.25427   Mean   :0.2722   Mean   :0.11461
##  3rd Qu.:0.14600   3rd Qu.:0.33910   3rd Qu.:0.3829   3rd Qu.:0.16140
##  Max.   :0.22260   Max.   :1.05800   Max.   :1.2520   Max.   :0.29100
##  symmetry_worst    fractal_dimension_worst
##  Min.   :0.1565    Min.   :0.05504
##  1st Qu.:0.2504    1st Qu.:0.07146
##  Median :0.2822    Median :0.08004
##  Mean   :0.2901    Mean   :0.08395
##  3rd Qu.:0.3179    3rd Qu.:0.09208
##  Max.   :0.6638    Max.   :0.20750
```

#Q3. How many variables/features in the data are suffixed with _mean? #Answer - 10

```r
#needed to take just the col names
#Then find a pattern in col names that ends with _mean
#then use length to count these instances
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"
```

```r
length(grep("_mean", colnames(wisc.df)))
```

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

# Check column means and standard deviations

```r
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
```

```r
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
```

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

```r
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
```
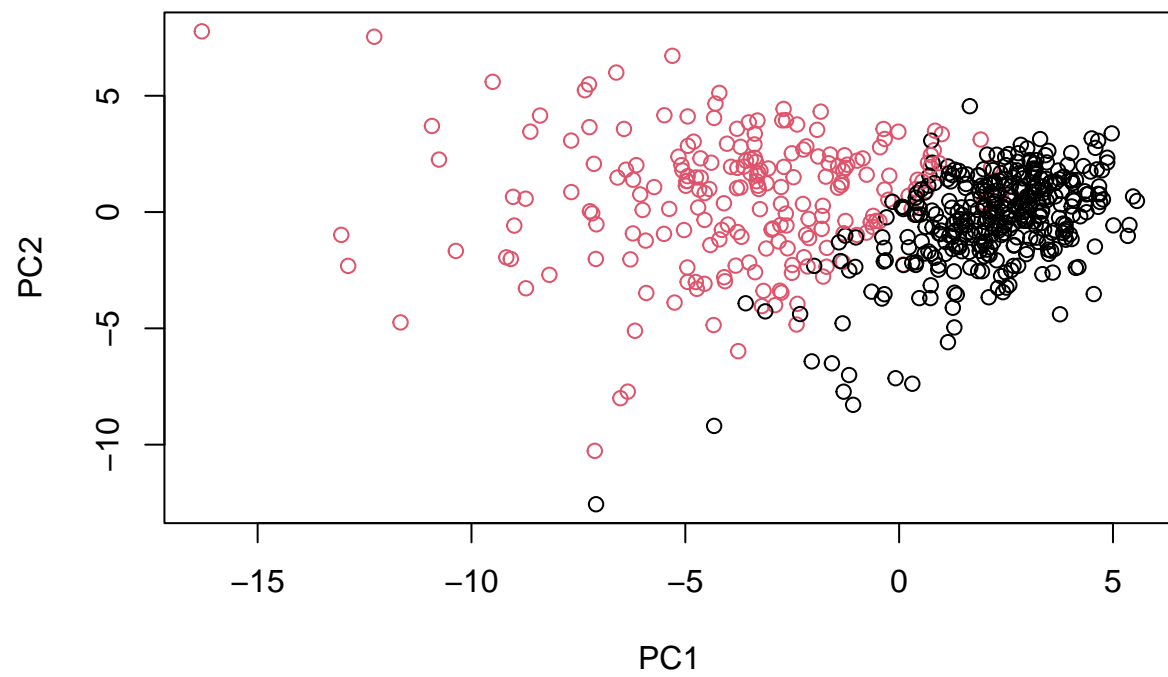
```
#biplot(wisc.pr)
```

#Q4. From your results, what proportion of the original variance is captured by the first principal components (PC1)? # Answer - 44.27%
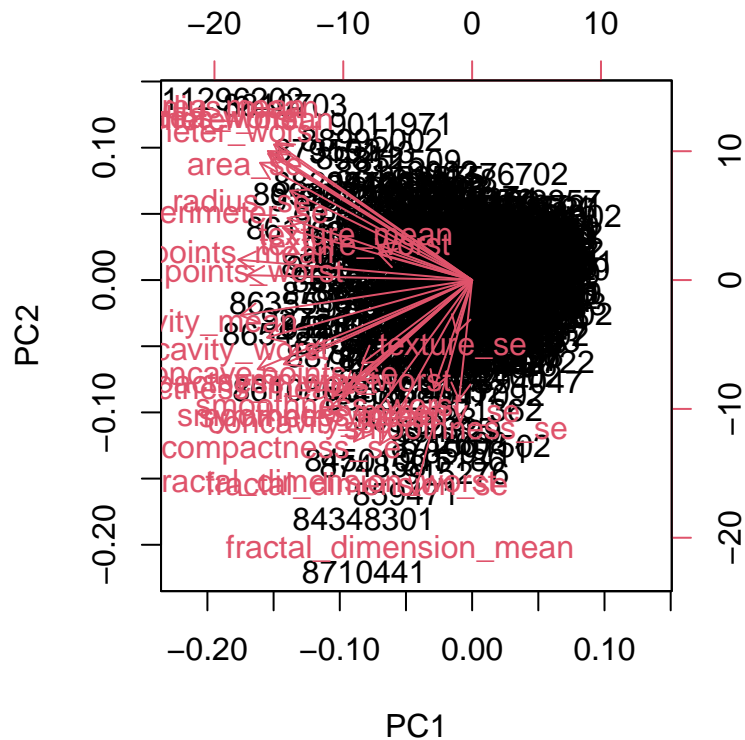
#Q5. How many principal components (PCs) are required to describe at least 70% of the original variance in the data? #answer - 3

#Q6. How many principal components (PCs) are required to describe at least 90% of the original variance in the data? #Answer - 7
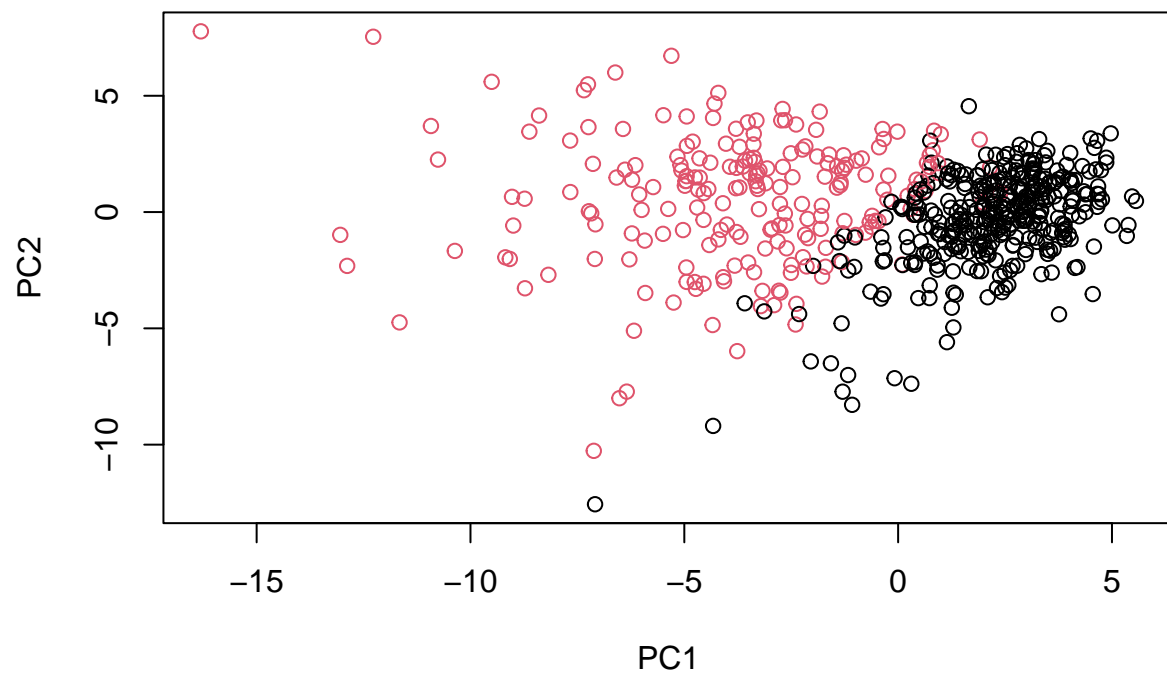
```
plot(wisc.pr$x[,1:2], col=diagnosis)
```
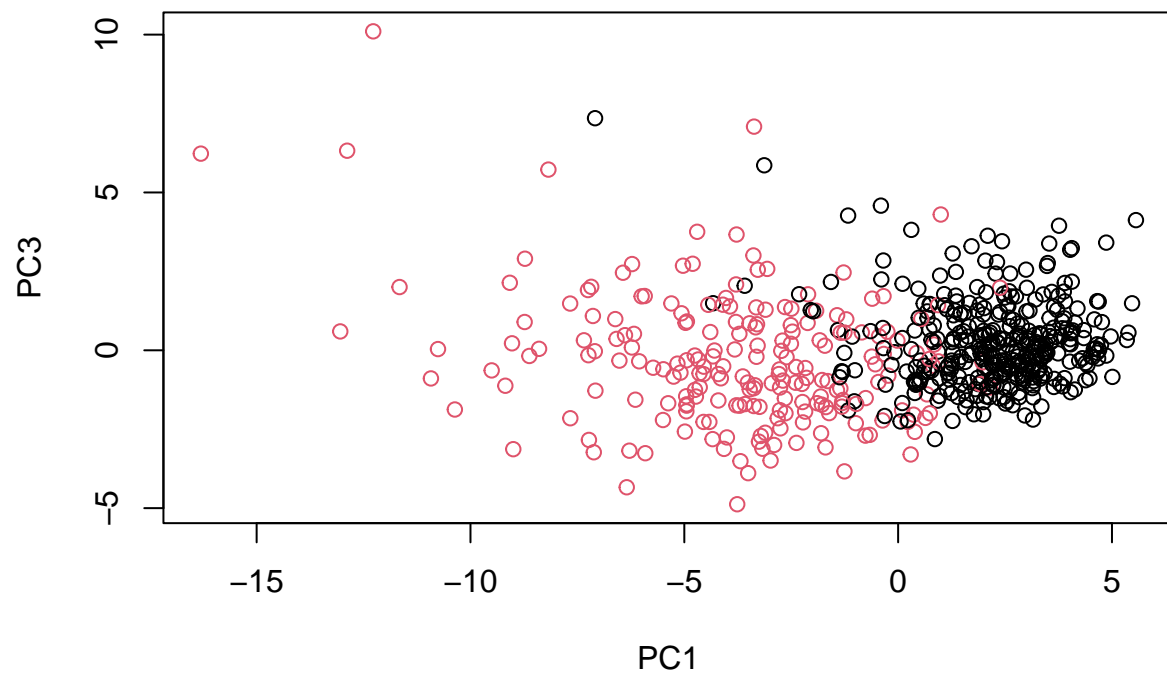
```
biplot(wisc.pr)
```

#Q7. What stands out to you about this plot? Is it easy or difficult to understand? Why? #Answer - difficult to understand, it is too compressed and saturated with data.

```r
# Scatter plot observations by components 1 and 2
plot( wisc.pr$x[,1:2] , col = 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$x[,c(1,3)] , col = diagnosis ,
      xlab = "PC1", ylab = "PC3")
```
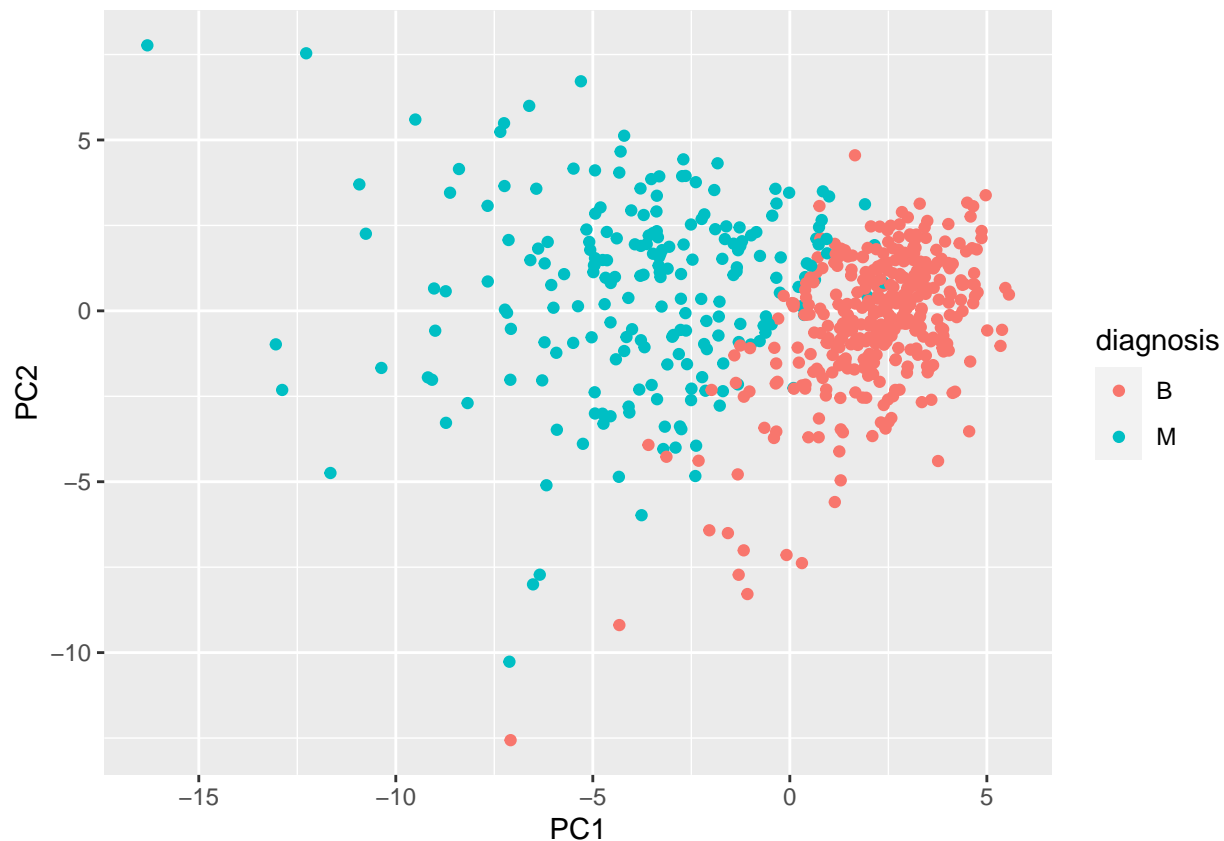
```
# Create a data.frame for ggplot
df <- as.data.frame(wisc.pr$x)
df$diagnosis <- diagnosis

# Load the ggplot2 package
library(ggplot2)

# Make a scatter plot colored by diagnosis
ggplot(df) +
  aes(PC1, PC2, col=diagnosis) +
  geom_point()
```
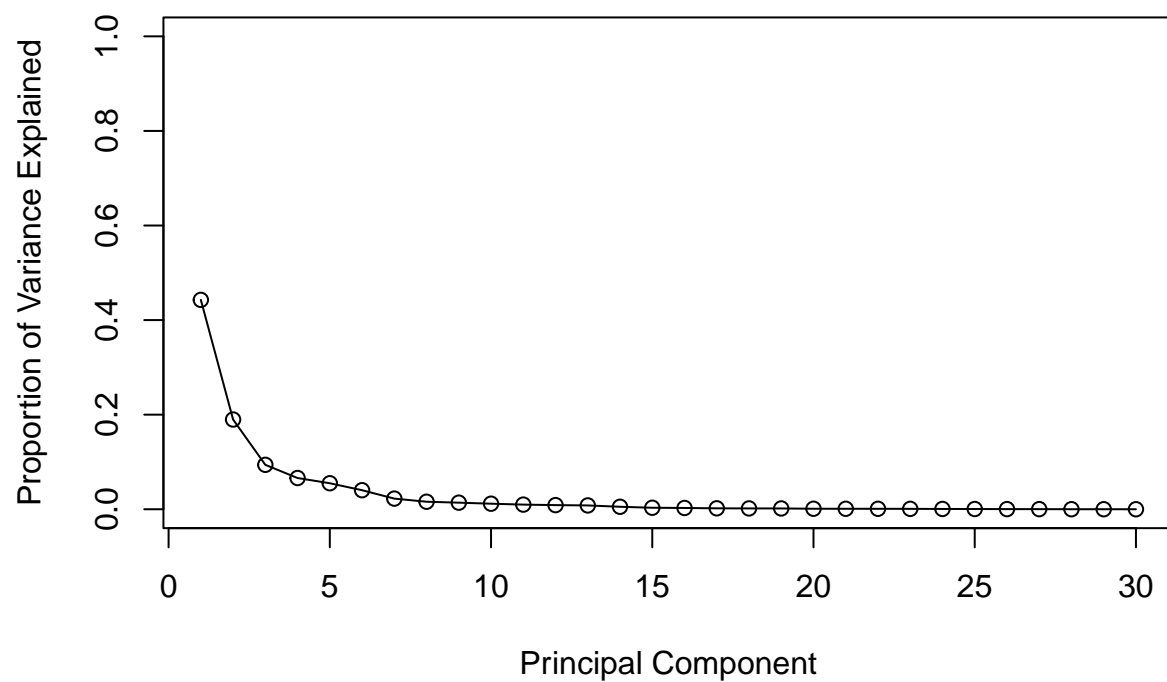
```r
# Calculate variance of each component
pr.var <- wisc.pr$sdev^2
head(pr.var)
```

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

```r
# Variance explained by each principal component: pve
total.var <- sum(pr.var)
pve <- pr.var / total.var
pve
```
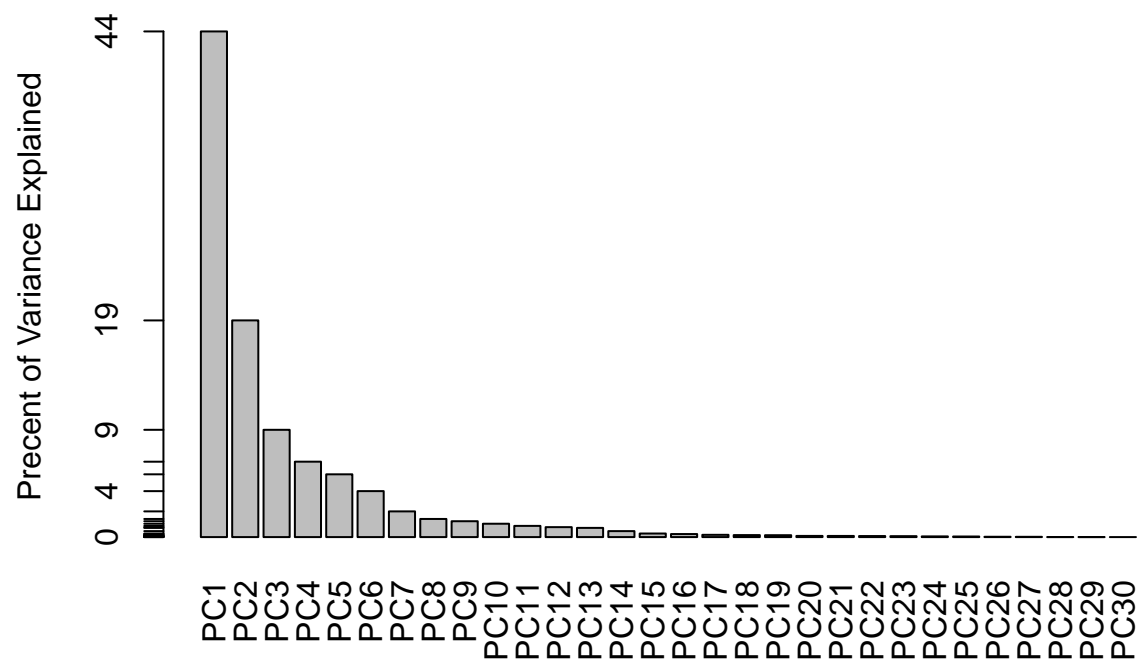
```
##  [1] 4.427203e-01 1.897118e-01 9.393163e-02 6.602135e-02 5.495768e-02
##  [6] 4.024522e-02 2.250734e-02 1.588724e-02 1.389649e-02 1.168978e-02
## [11] 9.797190e-03 8.705379e-03 8.045250e-03 5.233657e-03 3.137832e-03
## [16] 2.662093e-03 1.979968e-03 1.753959e-03 1.649253e-03 1.038647e-03
## [21] 9.990965e-04 9.146468e-04 8.113613e-04 6.018336e-04 5.160424e-04
## [26] 2.725880e-04 2.300155e-04 5.297793e-05 2.496010e-05 4.434827e-06
```

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

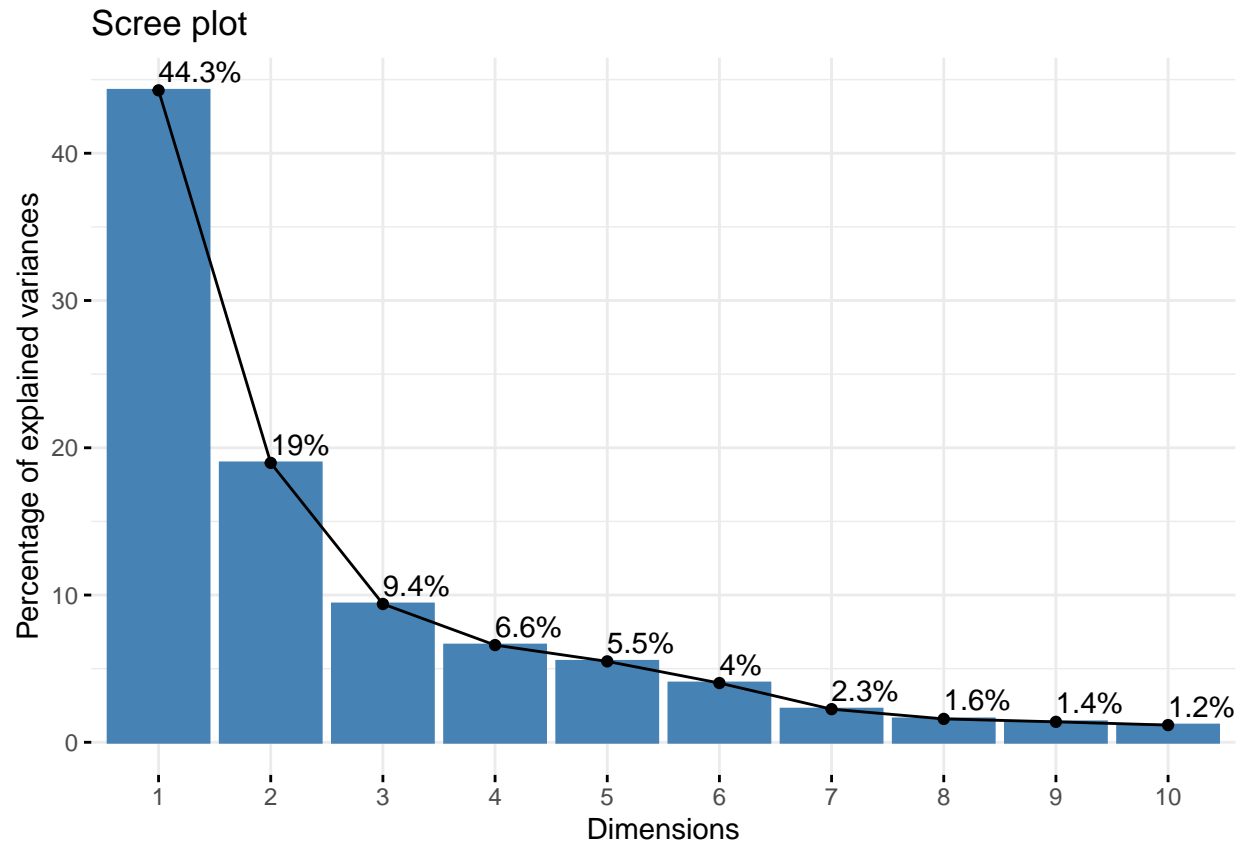**Alternative scree plot of the same data, note data driven y-axis**

```r
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 )
```

```
library(factoextra)
```

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

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

## Scree plot



#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? #Answer - -0.26085376

```
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? #Answer - 5

```
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
```

```
#need to call summary of our PCA and then assign it to a variable
#this summary has importance table. we get the sum of the 3 row of this table and see when it is less t
#this returns 4 but we would beed at least 5 PCs to get over 0.8.
var <- summary(wisc.pr)
sum(var$importance[3,] < 0.8)
```
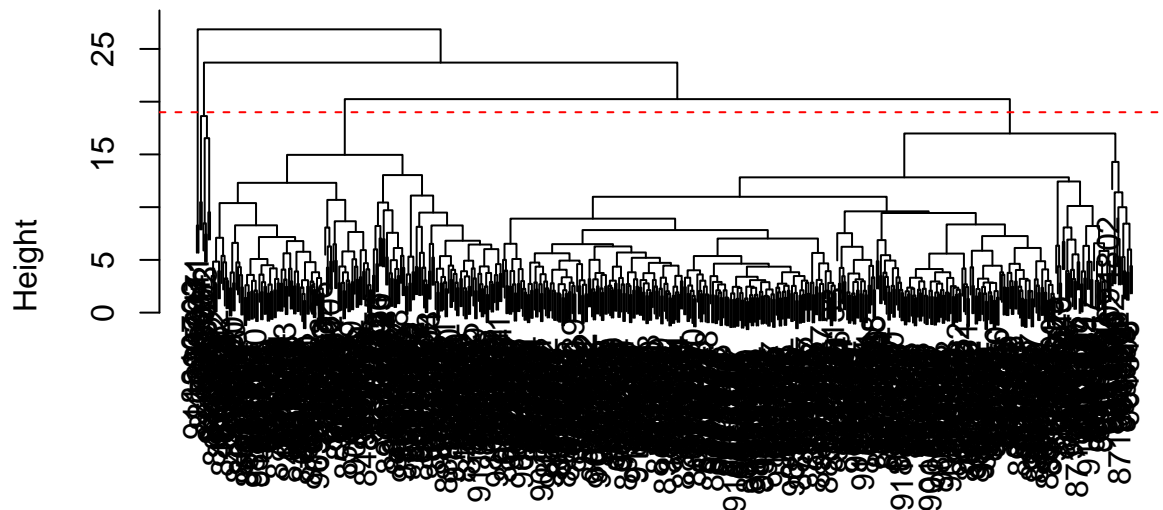
```
## [1] 4
```

```
# 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" )
```

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

## Cluster Dendrogram



data.dist
hclust (*, "complete")

#Q11. Using the plot() and abline() functions, what is the height at which the clustering model has 4 clusters? #Answer - 19

```
wisc.hclust.clusters <- cutree(wisc.hclust, k = 2, h = 19)
```

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

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

```
##                      diagnosis
## wisc.hclust.clusters   B   M
##                    1 357 210
##                    2   0   2
```

#Section 5 Here we aim to combine our PCA results with clustering. Essentially, we are going to cluster in "PC space", that is cluster on the 'wisc.pr$x' 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
```
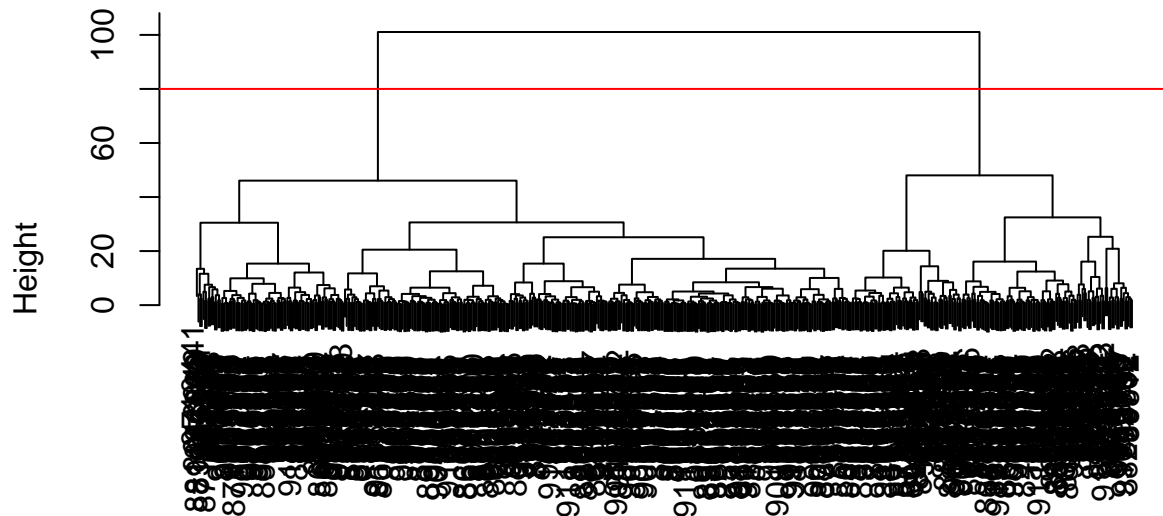
15

```
##                          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
```

##CLuster my PCA results I will use 4 PCs and 'hclust()' and 'dist()' as an input.

```
wisc.pr.hlust <- hclust( dist(wisc.pr$x[,1:4]), method = "ward.D2")
```

```
plot(wisc.pr.hlust)
abline(h=80, col = "red")
```



**Cluster Dendrogram**

dist(wisc.pr$x[, 1:4])
hclust (*, "ward.D2")

let's find our cluster membership vector by cutting this tree into k=2 groups

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

```
## grps
##   1   2
## 171 398
```

Now let's compare to the expert M and B vector

```
table(diagnosis)
```

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

we can do a cross-table by giving the 'table()' function two inputs #the B (6) below can be called false neg as they are in that clsuter that has mostly M

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

```
##      diagnosis
## grps   B   M
##    1   6 165
##    2 351  47
```

**Accuracy** essentially how many did we get correct? #nrow to get the total number of patients We got 90% accuracy

```
(165+351)/nrow(wisc.data)
```

```
## [1] 0.9068541
```

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

**Sensitivity** refers to a test's ability to correctly detect ill patients who do have the condition. In our example here the sensitivity is the total number of samples in the cluster identified as predominantly malignant (cancerous) divided by the total number of known malignant samples. In other words: TP/(TP+FN).

```
(165)/((165)+(6))
```

```
## [1] 0.9649123
```

**Specificity** relates to a test's ability to correctly reject healthy patients without a condition. In our example specificity is the proportion of benign (not cancerous) samples in the cluster identified as predominantly benign that are known to be benign. In other words: TN/(TN+FN).

```
351/(351+6)
```

```
## [1] 0.9831933
```

# PREDICTION

We will use the predict() function that will take our PCA model from before and new cancer cell data and project that data onto our PCA space.

```
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
```

#Now add these new samples to our PCA plot #pch = making points solid fill cex - makes these points bigger size 3 text here labels the points 1 and 2 and makes these numbers white

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