# BIOS 707 Final Project: Preliminary Analysis

*Angel Chen*

*12/3/2018*

**load data**

```r
setwd("~/Downloads")
raw <- read.csv("data.csv")
#for some reason an extra column was read in as NA values
#remove last column
raw <- raw[,1:32]
str(raw)
```

```
## 'data.frame':    569 obs. of  32 variables:
##  $ id                     : int  842302 842517 84300903 84348301 84358402 843786 844359 84458202 844
##  $ diagnosis              : Factor w/ 2 levels "B","M": 2 2 2 2 2 2 2 2 2 2 ...
##  $ radius_mean            : num  18 20.6 19.7 11.4 20.3 ...
##  $ texture_mean           : num  10.4 17.8 21.2 20.4 14.3 ...
##  $ perimeter_mean         : num  122.8 132.9 130 77.6 135.1 ...
##  $ area_mean              : num  1001 1326 1203 386 1297 ...
##  $ smoothness_mean        : num  0.1184 0.0847 0.1096 0.1425 0.1003 ...
##  $ compactness_mean       : num  0.2776 0.0786 0.1599 0.2839 0.1328 ...
##  $ concavity_mean         : num  0.3001 0.0869 0.1974 0.2414 0.198 ...
##  $ concave.points_mean    : num  0.1471 0.0702 0.1279 0.1052 0.1043 ...
##  $ symmetry_mean          : num  0.242 0.181 0.207 0.26 0.181 ...
##  $ fractal_dimension_mean : num  0.0787 0.0567 0.06 0.0974 0.0588 ...
##  $ radius_se              : num  1.095 0.543 0.746 0.496 0.757 ...
##  $ texture_se             : num  0.905 0.734 0.787 1.156 0.781 ...
##  $ perimeter_se           : num  8.59 3.4 4.58 3.44 5.44 ...
##  $ area_se                : num  153.4 74.1 94 27.2 94.4 ...
##  $ smoothness_se          : num  0.0064 0.00522 0.00615 0.00911 0.01149 ...
##  $ compactness_se         : num  0.049 0.0131 0.0401 0.0746 0.0246 ...
##  $ concavity_se           : num  0.0537 0.0186 0.0383 0.0566 0.0569 ...
##  $ concave.points_se      : num  0.0159 0.0134 0.0206 0.0187 0.0188 ...
##  $ symmetry_se            : num  0.03 0.0139 0.0225 0.0596 0.0176 ...
##  $ fractal_dimension_se   : num  0.00619 0.00353 0.00457 0.00921 0.00511 ...
##  $ radius_worst           : num  25.4 25 23.6 14.9 22.5 ...
##  $ texture_worst          : num  17.3 23.4 25.5 26.5 16.7 ...
##  $ perimeter_worst        : num  184.6 158.8 152.5 98.9 152.2 ...
##  $ area_worst             : num  2019 1956 1709 568 1575 ...
##  $ smoothness_worst       : num  0.162 0.124 0.144 0.21 0.137 ...
##  $ compactness_worst      : num  0.666 0.187 0.424 0.866 0.205 ...
##  $ concavity_worst        : num  0.712 0.242 0.45 0.687 0.4 ...
##  $ concave.points_worst   : num  0.265 0.186 0.243 0.258 0.163 ...
##  $ symmetry_worst         : num  0.46 0.275 0.361 0.664 0.236 ...
##  $ fractal_dimension_worst: num  0.1189 0.089 0.0876 0.173 0.0768 ...
```

```r
summary(raw)
```

```
##       id            diagnosis radius_mean    texture_mean
## Min.   :     8670  B:357    Min.   : 6.981   Min.   : 9.71
## 1st Qu.:   869218  M:212    1st Qu.:11.700   1st Qu.:16.17
## Median :   906024           Median :13.370   Median :18.84
## Mean   : 30371831           Mean   :14.127   Mean   :19.29
## 3rd Qu.:  8813129           3rd Qu.:15.780   3rd Qu.:21.80
## Max.   :911320502           Max.   :28.110   Max.   :39.28
## perimeter_mean     area_mean      smoothness_mean   compactness_mean
## Min.   : 43.79  Min.   : 143.5  Min.   :0.05263  Min.   :0.01938
## 1st Qu.: 75.17  1st Qu.: 420.3  1st Qu.:0.08637  1st Qu.:0.06492
## Median : 86.24  Median : 551.1  Median :0.09587  Median :0.09263
## Mean   : 91.97  Mean   : 654.9  Mean   :0.09636  Mean   :0.10434
## 3rd Qu.:104.10  3rd Qu.: 782.7  3rd Qu.:0.10530  3rd Qu.:0.13040
## Max.   :188.50  Max.   :2501.0  Max.   :0.16340  Max.   :0.34540
## concavity_mean    concave.points_mean symmetry_mean
## Min.   :0.00000  Min.   :0.00000     Min.   :0.1060
## 1st Qu.:0.02956  1st Qu.:0.02031     1st Qu.:0.1619
## Median :0.06154  Median :0.03350     Median :0.1792
## Mean   :0.08880  Mean   :0.04892     Mean   :0.1812
## 3rd Qu.:0.13070  3rd Qu.:0.07400     3rd Qu.:0.1957
## Max.   :0.42680  Max.   :0.20120     Max.   :0.3040
## fractal_dimension_mean   radius_se       texture_se      perimeter_se
## Min.   :0.04996        Min.   :0.1115  Min.   :0.3602  Min.   : 0.757
## 1st Qu.:0.05770        1st Qu.:0.2324  1st Qu.:0.8339  1st Qu.: 1.606
## Median :0.06154        Median :0.3242  Median :1.1080  Median : 2.287
## Mean   :0.06280        Mean   :0.4052  Mean   :1.2169  Mean   : 2.866
## 3rd Qu.:0.06612        3rd Qu.:0.4789  3rd Qu.:1.4740  3rd Qu.: 3.357
## Max.   :0.09744        Max.   :2.8730  Max.   :4.8850  Max.   :21.980
##    area_se        smoothness_se     compactness_se     concavity_se
## Min.   :  6.802  Min.   :0.001713  Min.   :0.002252  Min.   :0.00000
## 1st Qu.: 17.850  1st Qu.:0.005169  1st Qu.:0.013080  1st Qu.:0.01509
## Median : 24.530  Median :0.006380  Median :0.020450  Median :0.02589
## Mean   : 40.337  Mean   :0.007041  Mean   :0.025478  Mean   :0.03189
## 3rd Qu.: 45.190  3rd Qu.:0.008146  3rd Qu.:0.032450  3rd Qu.:0.04205
## Max.   :542.200  Max.   :0.031130  Max.   :0.135400  Max.   :0.39600
## concave.points_se    symmetry_se       fractal_dimension_se
## Min.   :0.000000  Min.   :0.007882  Min.   :0.0008948
## 1st Qu.:0.007638  1st Qu.:0.015160  1st Qu.:0.0022480
## Median :0.010930  Median :0.018730  Median :0.0031870
## Mean   :0.011796  Mean   :0.020542  Mean   :0.0037949
## 3rd Qu.:0.014710  3rd Qu.:0.023480  3rd Qu.:0.0045580
## 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
```

```r
#there are no NA values in the entire dataset, HOORAY!
bcdat <- raw
```
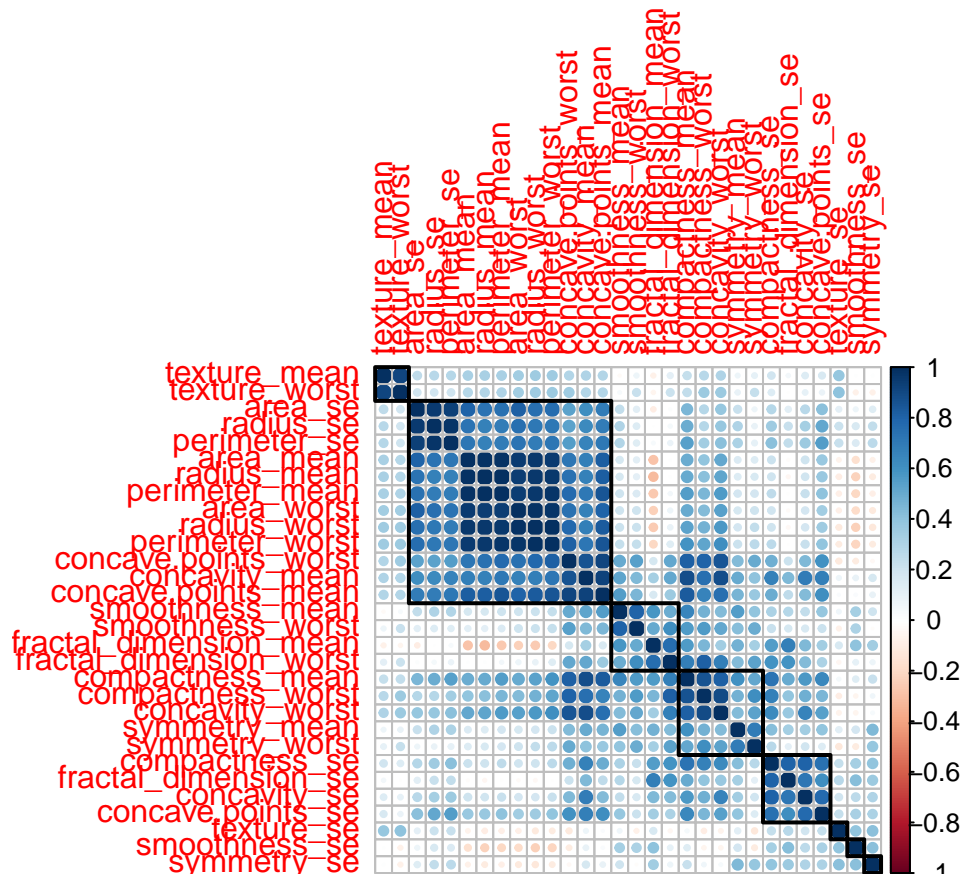
**data cleaning**

I suspect high correlation exists in this data because we have variables such as radius, area, perimeter and shape as a linear or some form of combination of each other

```r
correlations <- cor(bcdat[,-c(1:2)])
library(corrplot)
corrplot(correlations, order = "hclust", tl.cex=1, addrect = 8)
#find vars that are highly correlated
library(caret)
```

```
## Loading required package: lattice
```

```
## Loading required package: ggplot2
```

```
highCorr <- findCorrelation(correlations, cutoff = .85)
colnames(bcdat)[highCorr]
```

```
##  [1] "smoothness_mean"    "compactness_mean"   "area_mean"
##  [4] "compactness_worst"  "smoothness_worst"   "radius_worst"
##  [7] "symmetry_se"        "radius_mean"        "texture_worst"
## [10] "id"                 "radius_se"          "texture_se"
## [13] "diagnosis"
```

Should we remove some of these correlated values?

**PCA analysis (eliminated those that are highly correlated)**

```
bcdat_new <- bcdat[, -highCorr]
pca <- prcomp(bcdat_new[,-c(1:2)], scale = TRUE, center = TRUE)
summary(pca)
```

```
## Importance of components:
##                          PC1    PC2    PC3     PC4     PC5     PC6     PC7
## Standard deviation     2.868 1.8740 1.3566 1.01197 0.80043 0.70610 0.55982
## Proportion of Variance 0.484 0.2066 0.1082 0.06024 0.03769 0.02933 0.01844
## Cumulative Proportion  0.484 0.6906 0.7988 0.85908 0.89677 0.92610 0.94453
##                          PC8    PC9   PC10    PC11    PC12    PC13
```
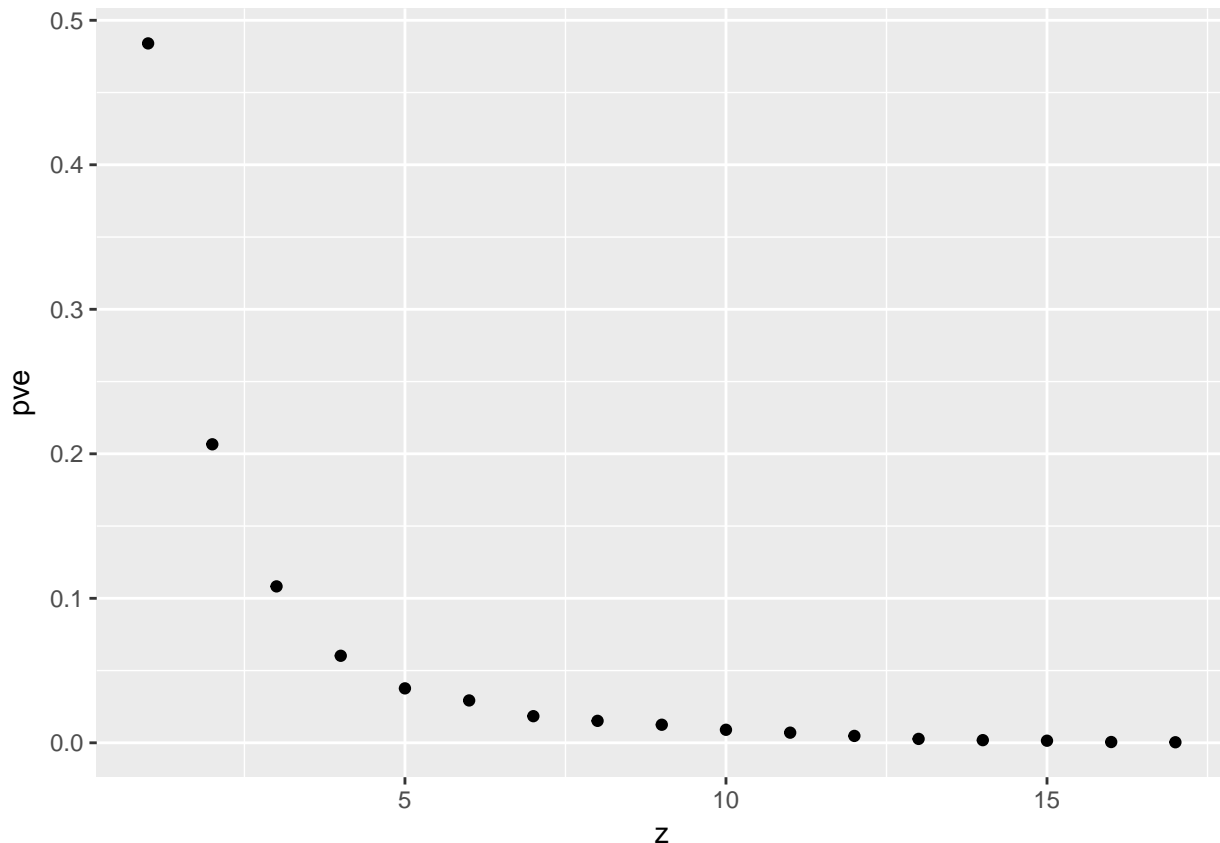
4

```
## Standard deviation      0.50842 0.4610 0.39145 0.34576 0.28519 0.2144
## Proportion of Variance 0.01521 0.0125 0.00901 0.00703 0.00478 0.0027
## Cumulative Proportion  0.95974 0.9722 0.98125 0.98829 0.99307 0.9958
##                              PC14    PC15    PC16    PC17
## Standard deviation      0.17672 0.15693 0.09750 0.08058
## Proportion of Variance 0.00184 0.00145 0.00056 0.00038
## Cumulative Proportion  0.99761 0.99906 0.99962 1.00000
```

```r
#calculate the standard deviation
pca.var = pca$sdev^2

#Calculate proportion of variance explained
pve = pca.var/sum(pca.var)
z = seq(1,17)

#Calculate cummulative PVE
cumpve = cumsum(pve)
pve.table = as.data.frame(cbind(z,pve, cumpve))

#plot variables against proportion variance explained
ggplot(pve.table, aes(x=z,y=pve))+geom_point()
```



```r
#plot variables against cumulative varinace explained


#to get 95% of the information in our data we need about 8 PCs
```

```
#explore the first 3 PCs
#install.packages("GGally")
library(GGally)

PCs <- data.frame(pca$x)
PCs$diagnosis <- bcdat$diagnosis
ggpairs(data=PCs,columns = 1:3, ggplot2::aes(color=diagnosis))
```