# Class 7: Machine Learning I

Cameron Jones

In this class we will explore clustering and dimensionality reduction methods.
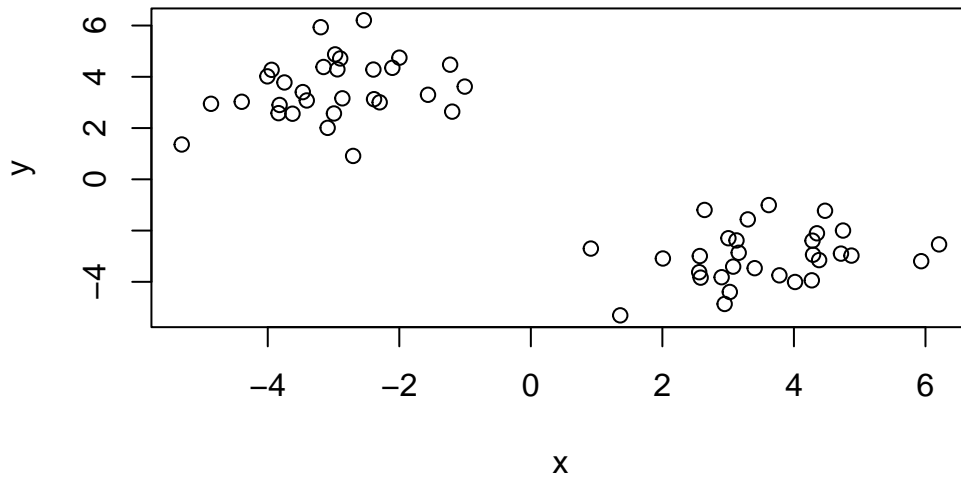
## K-means

Make up some input data where we know what the answer should be.

```
tmp <- c(rnorm(30, -3), rnorm(30, +3))
x <- cbind(x=tmp, y=rev(tmp))
head(x)
```

```
               x         y
[1,] -1.998421 4.7464534
[2,] -1.227297 4.4726691
[3,] -3.836039 2.5817126
[4,] -3.194328 5.9353295
[5,] -1.563857 3.2985113
[6,] -2.702216 0.9123488
```

Quick plot of x to see the two groups at -3,+3 and +3,-3

```
plot(x)
```

Use the 'kmeans()' function setting k to 2 and nstart=20

```
km <- kmeans(x, centers = 2, nstart=20)
km
```

```
K-means clustering with 2 clusters of sizes 30, 30

Cluster means:
          x          y
1 -2.997948  3.550464
2  3.550464 -2.997948

Clustering vector:
 [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2
[39] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2

Within cluster sum of squares by cluster:
[1] 72.43194 72.43194
 (between_SS / total_SS =  89.9 %)

Available components:
```

2

```
[1] "cluster"      "centers"      "totss"      "withinss"      "tot.withinss"
[6] "betweenss"    "size"         "iter"       "ifault"
```

Q. How many points are in each cluster?

```r
km$size
```

```
[1] 30 30
```

Q. What 'component' of your results object details -cluster assignment/membership? -cluster center?

```r
km$cluster
```

```
 [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2
[39] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```
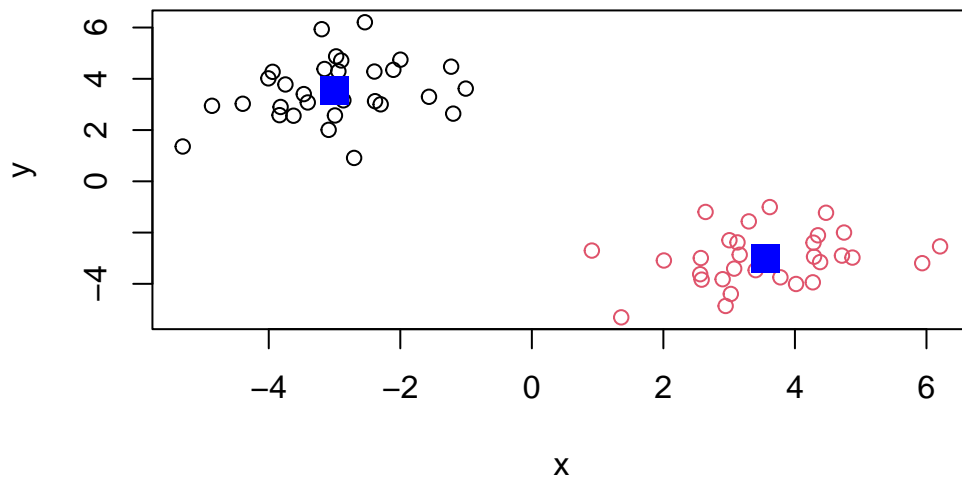
```r
km$centers
```

```
          x          y
1 -2.997948   3.550464
2  3.550464  -2.997948
```
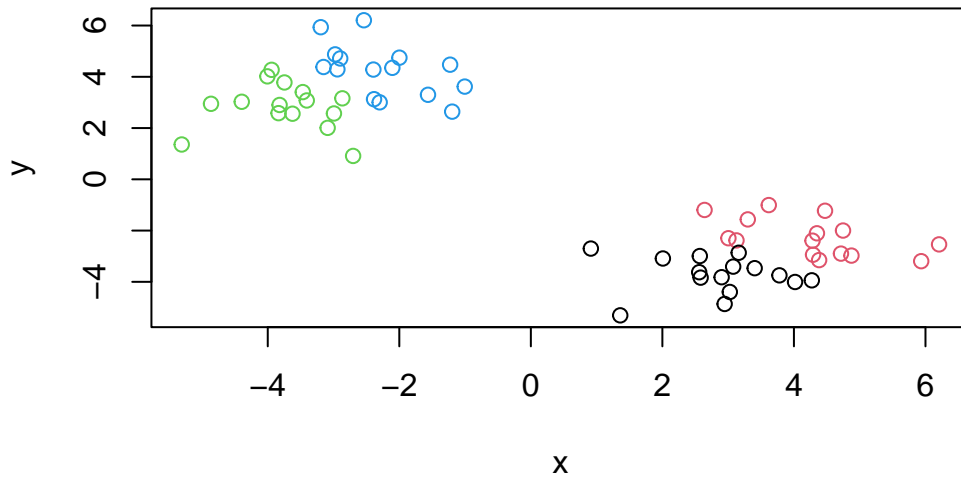
Q. Plot x colored by the kmeans cluster assignment and add cluster centers as blue points

```r
plot(x, col=km$cluster)
points(km$centers, col="blue", pch=15, cex=2)
```

Play with kmeans and ask for different number of clusters

```
km <- kmeans(x, centers = 4, nstart=20)
plot(x, col=km$cluster)
points(km$clusters, col="blue", pch=16, cex=2)
```

#Hierarchical Clustering

This is another very useful and widely employed clustering method which has the advantage over k-means in that it an help reveal the something of the true grouping in your data.

The 'hclust()' function wants a distance matrix as input. We can get this from the 'dist()' function.
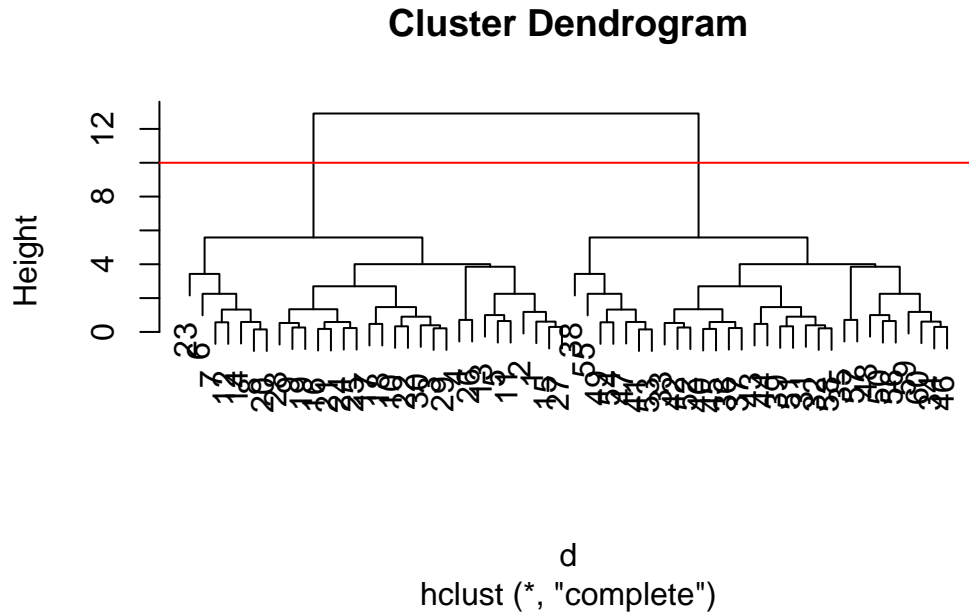
```
d <- dist(x)
hc <- hclust(d)
hc
```

```
Call:
hclust(d = d)

Cluster method   : complete
Distance         : euclidean
Number of objects: 60
```

There is a plot method for hclust results:
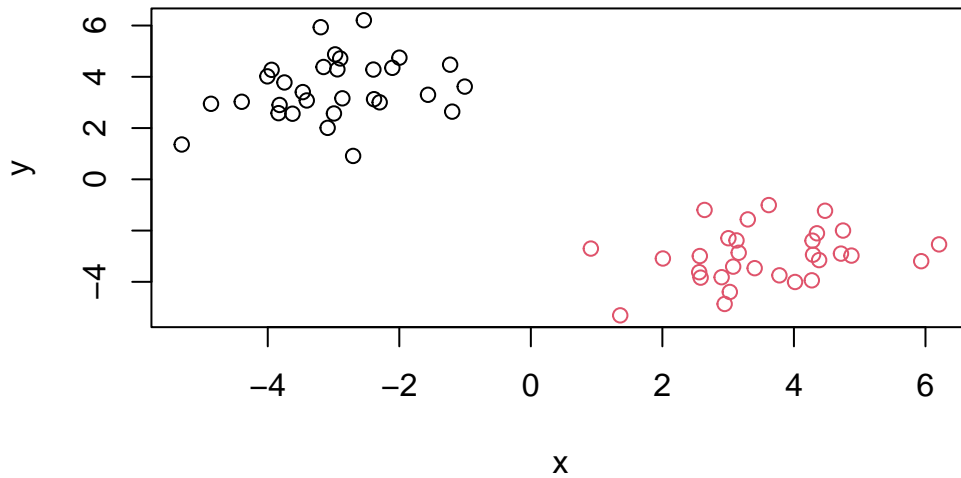
```
plot(hc)
abline(h=10, col="red")
```

## Cluster Dendrogram



d
hclust (*, "complete")

To get my cluster membership vector, I need to "cut" my tree to yield sub-trees or branches with all the members of a given cluster residing on the same cut branch. The function to do this is called 'cutree()'

```
grps <- cutree(hc, h=10)
grps
```

```
 [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2
[39] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

```
plot(x, col=grps)
```

It is often helpful to use the 'k=' argument to cutree rather than the height 'h=' of cutting with 'cutree()'. This will cut the tree to yield the number of clusters you want.

```r
cutree(hc, k=2)
```

```
 [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2
[39] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
```

## Principal Component Analysis (PCA)

The base R function for PCA is called 'prcomp()' Let's play with some 17D data

##PCA of UK food Data

Import the data

```r
url <- "https://tinyurl.com/UK-foods"
x <- read.csv(url)
head(x)
```

```
          X England Wales Scotland N.Ireland
1        Cheese     105    103      103        66
2  Carcass_meat     245    227      242       267
3    Other_meat     685    803      750       586
4          Fish     147    160      122        93
5 Fats_and_oils     193    235      184       209
6        Sugars     156    175      147       139
```

Q1. How many rows and columns are in your new data frame named x? What R functions could you use to answer this questions?

```r
dim(x)
```

```
[1] 17  5
```

```r
nrow(x)
```

```
[1] 17
```

```r
ncol(x)
```

```
[1] 5
```

```r
rownames(x) <- x[,1]
x <- x[,-1]
head(x)
```

```
              England Wales Scotland N.Ireland
Cheese            105   103      103        66
Carcass_meat      245   227      242       267
Other_meat        685   803      750       586
Fish              147   160      122        93
Fats_and_oils     193   235      184       209
Sugars            156   175      147       139
```
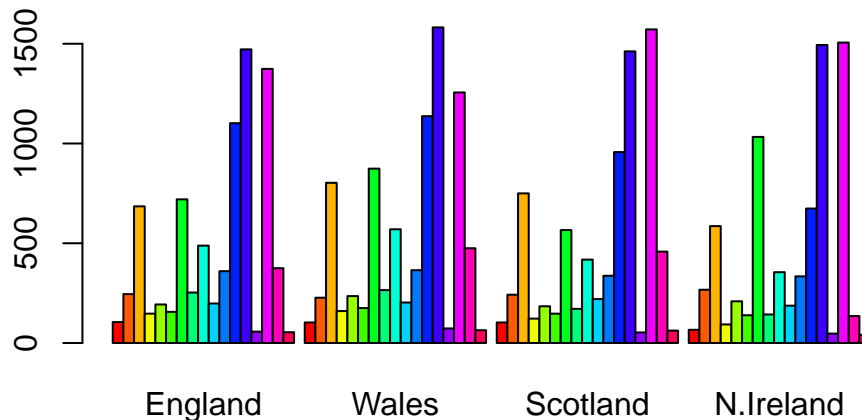
```r
dim(x)
```

```
[1] 17  4
```

```
x <- read.csv(url, row.names=1)
head(x)
```

```
              England Wales Scotland N.Ireland
Cheese            105   103      103        66
Carcass_meat      245   227      242       267
Other_meat        685   803      750       586
Fish              147   160      122        93
Fats_and_oils     193   235      184       209
Sugars            156   175      147       139
```
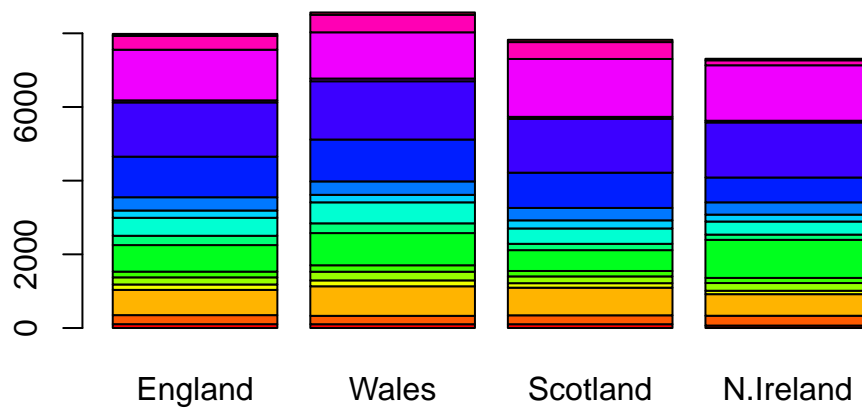
Q2. Which approach to solving the 'row-names problem' mentioned above do you prefer and why? Is one approach more robust than another under certain circumstances?

The second approach is better, as the first approach, if repeated, will start deleting columns!
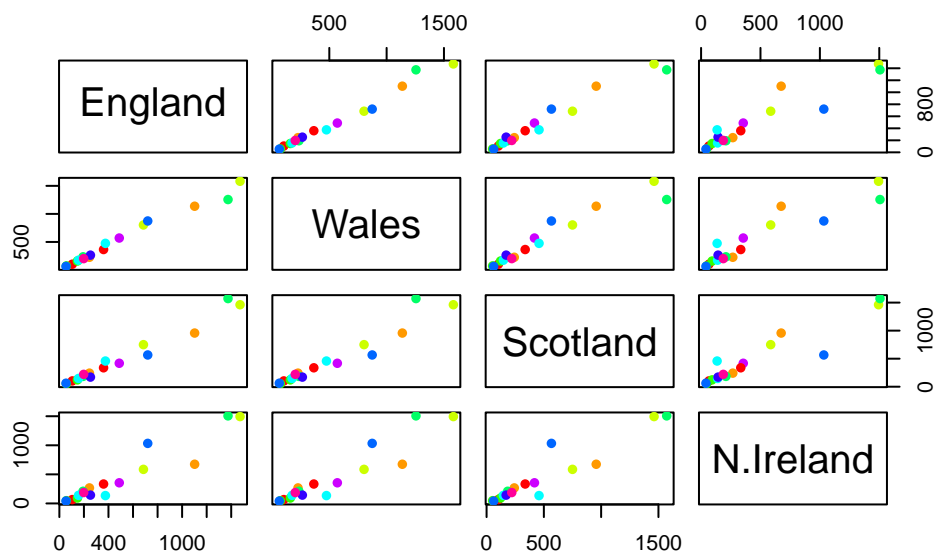
```
barplot(as.matrix(x), beside=T, col=rainbow(nrow(x)))
```



```
barplot(as.matrix(x), beside=F, col=rainbow(nrow(x)))
```

```
pairs(x, col=rainbow(10), pch=16)
```



10

```
pca <- prcomp(t(x))
summary(pca)
```

```
Importance of components:
                         PC1      PC2      PC3      PC4
Standard deviation   324.1502 212.7478 73.87622 4.189e-14
Proportion of Variance  0.6744   0.2905  0.03503 0.000e+00
Cumulative Proportion   0.6744   0.9650  1.00000 1.000e+00
```
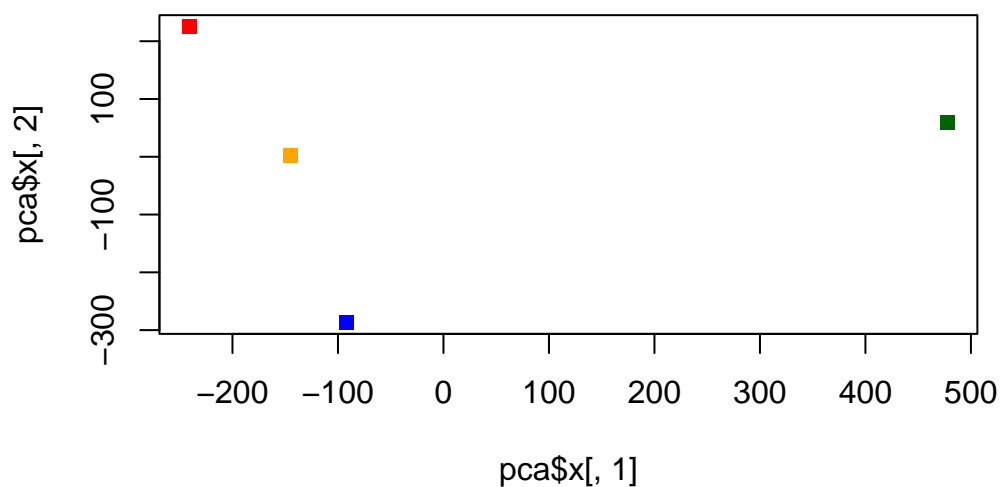
A "PCA" plot (aka "Score plot", PC1vsPC2 plot, etc.)

```
pca$x
```

```
                 PC1         PC2          PC3          PC4
England    -144.99315    2.532999 -105.768945  2.842865e-14
Wales      -240.52915  224.646925   56.475555  7.804382e-13
Scotland    -91.86934 -286.081786   44.415495 -9.614462e-13
N.Ireland   477.39164   58.901862    4.877895  1.448078e-13
```

```
plot(pca$x[,1], pca$x[,2],
     col=c("orange", "red", "blue", "darkgreen"), pch=15)
```
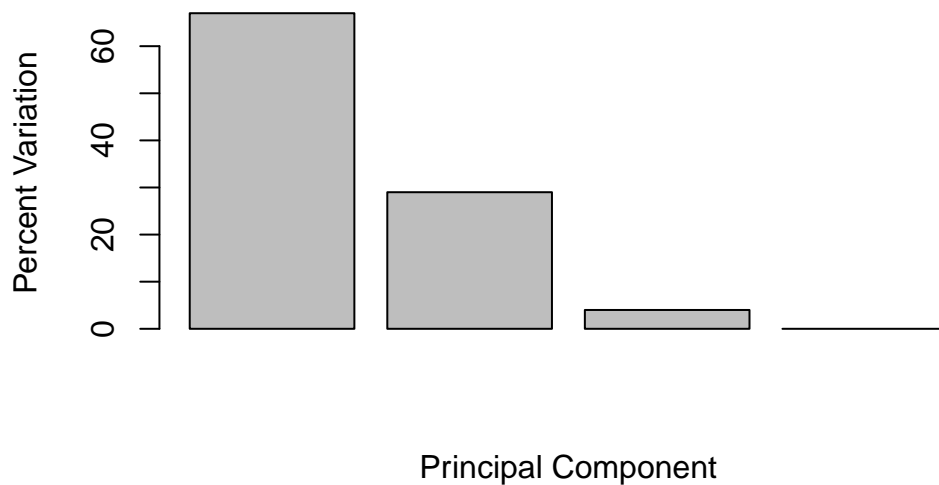
```
v <- round( pca$sdev^2/sum(pca$sdev^2) * 100 )
v
```

[1] 67 29  4  0

```
z <- summary(pca)
z$importance
```

|  | PC1 | PC2 | PC3 | PC4 |
|---|---|---|---|---|
| Standard deviation | 324.15019 | 212.74780 | 73.87622 | 4.188568e-14 |
| Proportion of Variance | 0.67444 | 0.29052 | 0.03503 | 0.000000e+00 |
| Cumulative Proportion | 0.67444 | 0.96497 | 1.00000 | 1.000000e+00 |

```
barplot(v, xlab="Principal Component", ylab="Percent Variation")
```



```
## Lets focus on PC1 as it accounts for > 90% of variance
par(mar=c(10, 3, 0.35, 0))
barplot( pca$rotation[,1], las=2 )
```