# Class 7: Machine Learning 1

Yu (Ericsson) Cao (PID: A16421048)
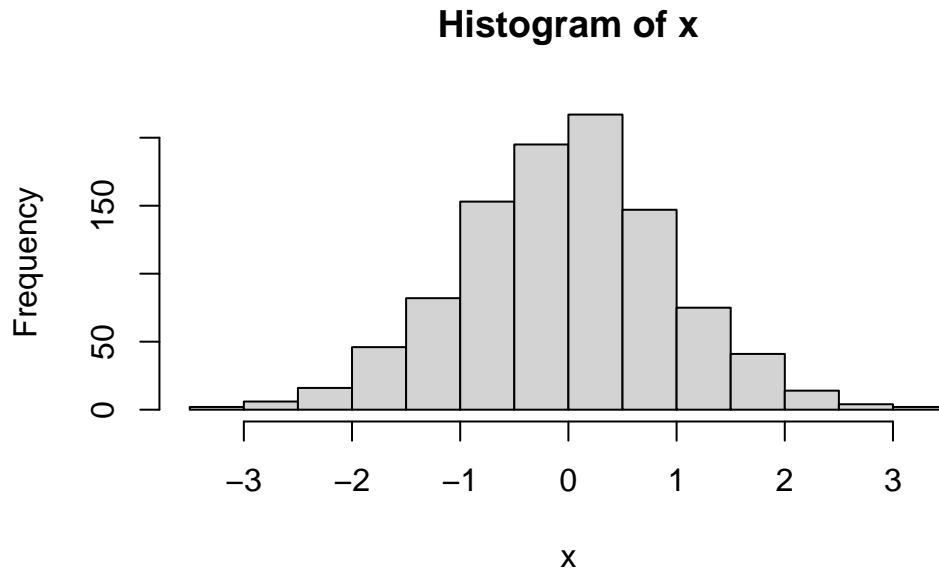
#Clustering Methods

The broad goal here is to find groupings (clusters) in your input data.

## Kmeans

First, let's make up some data to cluster.

```
x <- rnorm(1000)
hist(x)
```
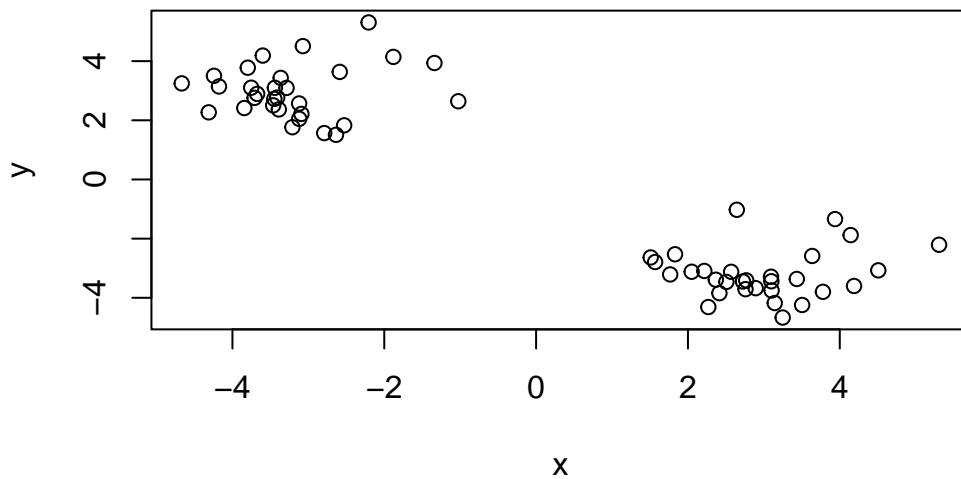
**Histogram of x**

Make a vector of length 60 with 30 points centered at -3 and 30 points centered at +3

```r
tmp <- c(rnorm(30, mean=-3), rnorm(30, mean=3))
tmp
```

```
 [1] -3.287247 -2.586324 -2.206540 -3.463605 -3.410626 -3.441128 -1.339530
 [8] -3.708213 -3.448137 -3.389156 -4.668121 -1.880439 -3.091712 -3.752532
[15] -4.243483 -3.364603 -2.791006 -2.528414 -3.121589 -3.798048 -4.177534
[22] -4.312539 -1.024811 -3.600514 -3.210222 -3.120242 -3.072386 -2.635961
[29] -3.674701 -3.843721  2.414547  2.892438  1.509369  4.507604  2.049495
[36]  1.766157  4.187772  2.643593  2.268707  3.143261  3.778334  2.569274
[43]  1.830291  1.567639  3.434061  3.504261  3.102103  2.215567  4.143188
[50]  3.248103  2.367188  2.724765  2.758572  3.937583  3.099056  2.767694
[57]  2.506343  5.308139  3.636785  3.094533
```

I will now make a wee x and y dataset with 2 groups of points.

```r
x <- cbind(x=tmp, y=rev(tmp))
plot(x)
```

```
k <- kmeans(x, centers=2)
k
```

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

Cluster means:
          x          y
1  2.965881 -3.206436
2 -3.206436  2.965881

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

Within cluster sum of squares by cluster:
[1] 43.18066 43.18066
 (between_SS / total_SS =  93.0 %)

Available components:

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

Q. From your result object **k** how many points are in each cluster?

```
k$size
```

```
[1] 30 30
```

Q. What "component" of your result object details the cluster membership?

```
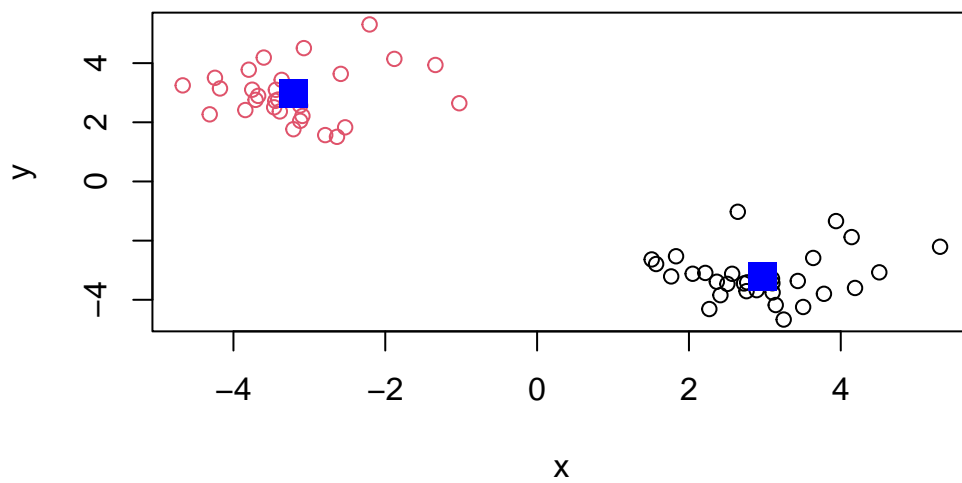k$cluster
```

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

Q. Cluster centers?

```
k$centers
```

```
           x          y
1   2.965881 -3.206436
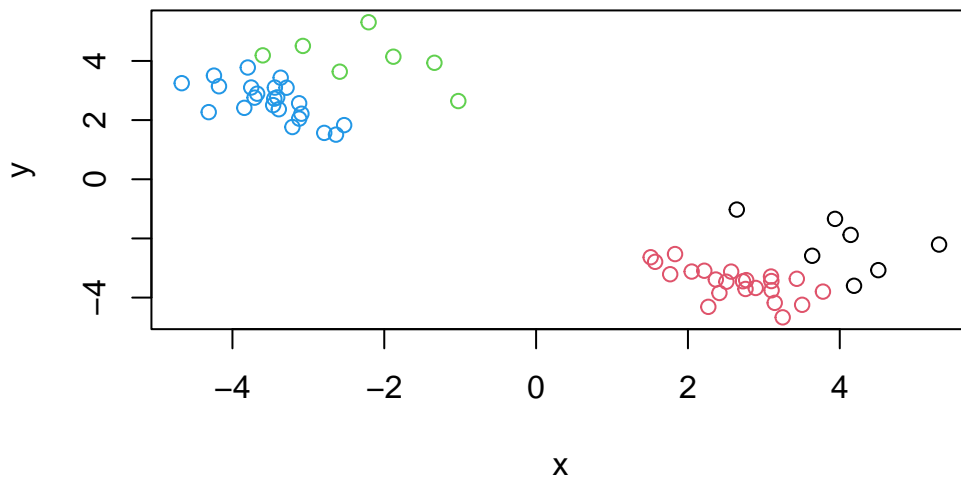2  -3.206436  2.965881
```

Q. Plot of our clustering results.

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



We can cluster into 4 groups

```r
# kmeans
k4 <- kmeans(x, centers=4)
# plot results
plot(x, col=k4$cluster)
```

A big limitation of kmeans is that it does what you ask even if you ask for silly clusters.

## Hierarchial Clustering

The main base R function for Hierarchial Clustering is `hclust()`. Unlike `kmeans()` you can not just pass it your data as input. You first need to calculate a distance matrix.

```
d <- dist(x)
hc <- hclust(d)
hc
```

```
Call:
hclust(d = d)

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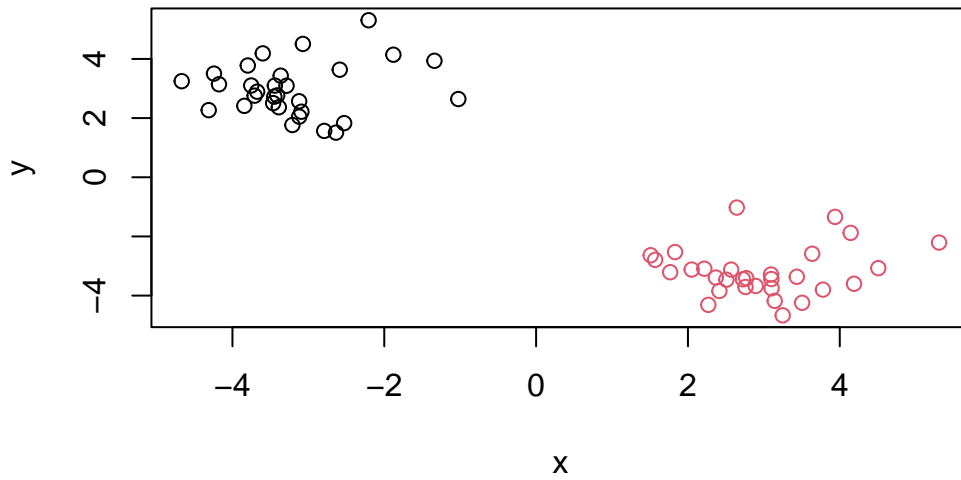```

Use plot() to view results

```r
plot(hc)
abline(h=10, col="red")
```

**Cluster Dendrogram**



d
hclust (*, "complete")

To make the "cut" and get our cluster membership vector we can use the `cutree()` function.

```r
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
```

Make a plot of our data colored by hclust results

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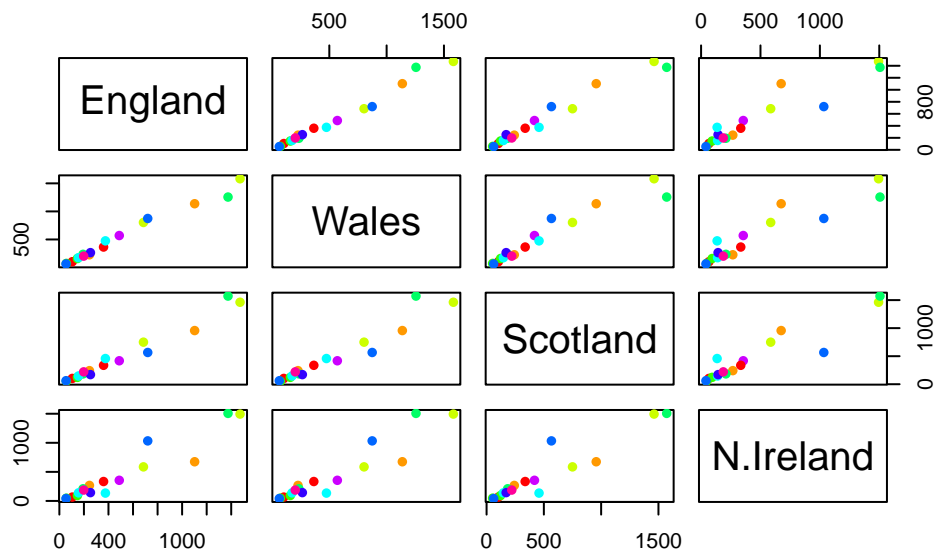plot(x, col=grps)
```

## Principal Component Analysis (PCA)

Here we will do Principal Component Analysis (PCA) on some food data from the UK.

```
url <- "https://tinyurl.com/UK-foods"
x <- read.csv(url, row.names=1)
```

```
#rownames(x) <- x[,-1]
#x <- x[,-1]
#x
```

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

## PCA to the rescue

The main "base" R function for PCA is called `prcomp()`. Here we need to take the transpose of our input as we want the countries in the rows and food as the column.

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

```
Importance of components:
                          PC1      PC2      PC3       PC4
Standard deviation    324.1502 212.7478 73.87622 2.921e-14
Proportion of Variance  0.6744   0.2905  0.03503 0.000e+00
Cumulative Proportion   0.6744   0.9650  1.00000 1.000e+00
```

Q. How much variance is captured in 2 PCs

96.5%

To make our main "PC score plot" (a.k.a. or "PC1 vs PC2 plot", or "PC plot" or "ordination plot").

```
attributes(pca)
```

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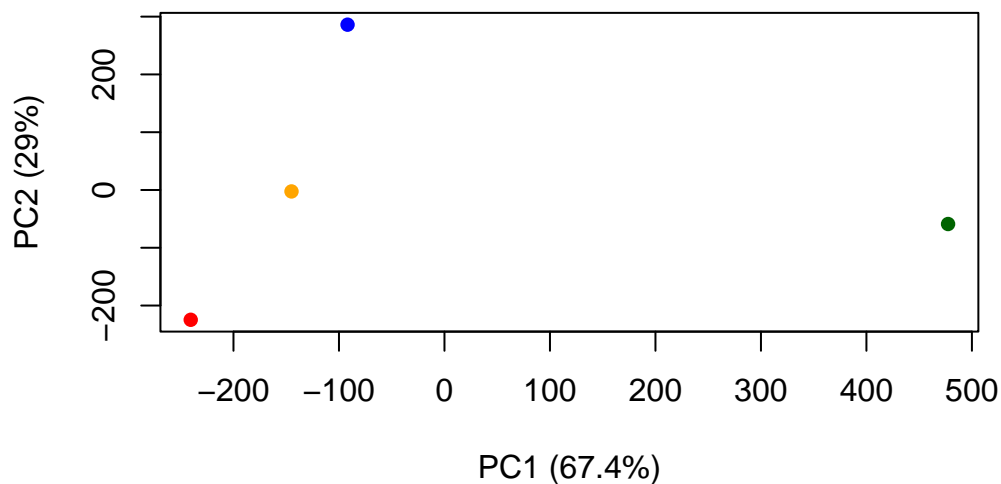

$class
[1] "prcomp"
```

We are after the `pca$x` result component to make our main PCA plot.

```
pca$x
```

```
                 PC1            PC2          PC3           PC4
England    -144.99315    -2.532999 105.768945 -9.152022e-15
Wales      -240.52915 -224.646925 -56.475555  5.560040e-13
Scotland    -91.86934  286.081786 -44.415495 -6.638419e-13
N.Ireland   477.39164  -58.901862  -4.877895  1.329771e-13
```

```
mycols <- c("orange", "red", "blue", "darkgreen")
plot(pca$x[,1], pca$x[,2], col=mycols, pch=16, xlab="PC1 (67.4%)", ylab="PC2 (29%)")
```



Another important result from PCA is how the original variables (in this case the foods) contribute to the PCs.

This is contained in the `pca$rotation` object - folks often call this the "loadings" or "contributions" to the PCs

```
pca$rotation
```

```
                          PC1           PC2          PC3           PC4
Cheese           -0.056955380   0.016012850   0.02394295  -0.409382587
Carcass_meat      0.047927628   0.013915823   0.06367111   0.729481922
Other_meat       -0.258916658  -0.015331138  -0.55384854   0.331001134
Fish             -0.084414983  -0.050754947   0.03906481   0.022375878
Fats_and_oils    -0.005193623  -0.095388656  -0.12522257   0.034512161
Sugars           -0.037620983  -0.043021699  -0.03605745   0.024943337
Fresh_potatoes    0.401402060  -0.715017078  -0.20668248   0.021396007
Fresh_Veg        -0.151849942  -0.144900268   0.21382237   0.001606882
Other_Veg        -0.243593729  -0.225450923  -0.05332841   0.031153231
Processed_potatoes -0.026886233   0.042850761  -0.07364902  -0.017379680
Processed_Veg    -0.036488269  -0.045451802   0.05289191   0.021250980
Fresh_fruit      -0.632640898  -0.177740743   0.40012865   0.227657348
Cereals          -0.047702858  -0.212599678  -0.35884921   0.100043319
Beverages        -0.026187756  -0.030560542  -0.04135860  -0.018382072
Soft_drinks       0.232244140   0.555124311  -0.16942648   0.222319484
Alcoholic_drinks -0.463968168   0.113536523  -0.49858320  -0.273126013
Confectionery    -0.029650201   0.005949921  -0.05232164   0.001890737
```

We can make a plot along PC1.

```
library(ggplot2)

contrib <- as.data.frame(pca$rotation)

ggplot(contrib) +
  aes(PC1, rownames(contrib)) +
  geom_col()
```