# Class 7: Machine Learning 1

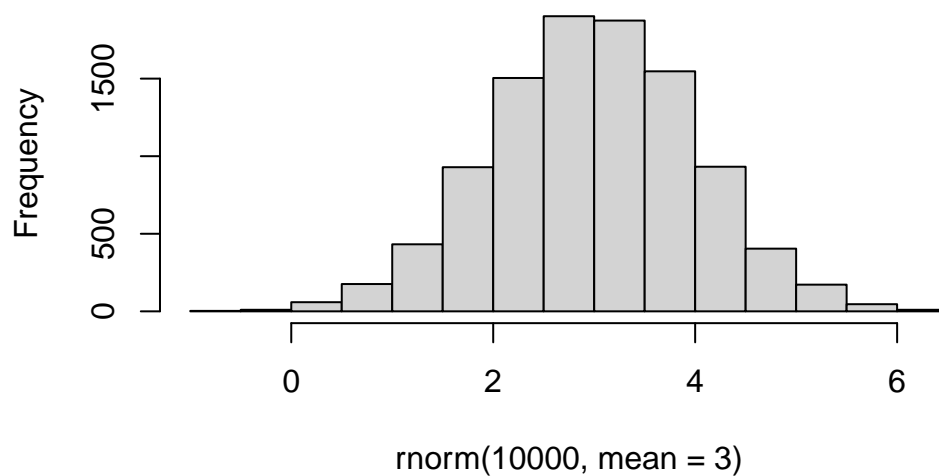Barry

## Table of contents

Today we will explore unsupervised machine learning methods starting with clustering and dimensionality reduction.

## Clustering

To start let's make up some data to cluster where we know what the answer should be. The `rnorm()` function will help us here.

```
hist( rnorm(10000, mean=3) )
```

## Histogram of rnorm(10000, mean = 3)



Return 30 numbers centred on -3

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

x <- cbind(x=tmp, y=rev(tmp))

x
```
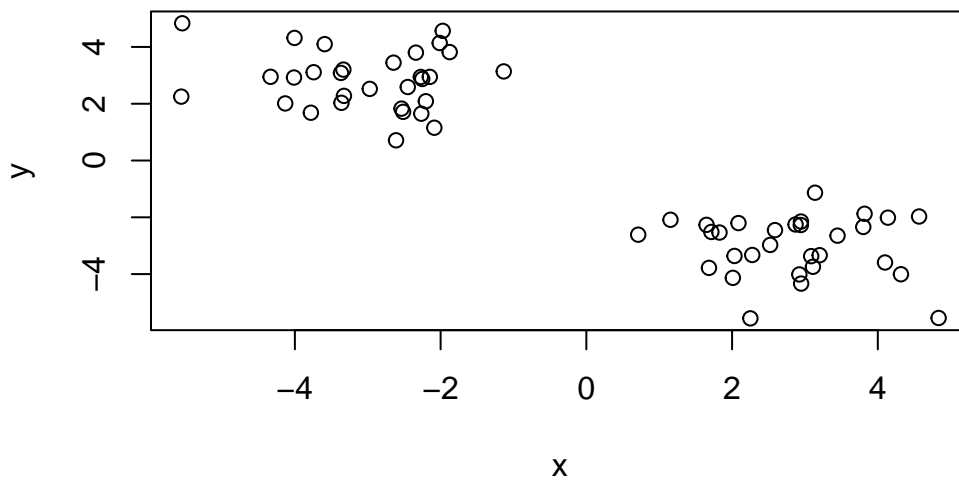
```
             x           y
 [1,] -4.0052498   4.3175912
 [2,] -2.5365722   1.8276559
 [3,] -3.5911897   4.0986464
 [4,] -2.4508472   2.5887270
 [5,] -4.1325932   2.0106003
 [6,] -3.7806618   1.6829283
 [7,] -3.3607000   2.0327000
 [8,] -2.6114337   0.7115766
 [9,] -2.0864165   1.1543562
[10,] -5.5458617   4.8342991
[11,] -1.1335161   3.1381491
[12,] -2.2518117   2.8710214
[13,] -2.1482412   2.9451070
```

```
[14,] -2.2651313  1.6493587
[15,] -2.5145047  1.7154397
[16,] -1.8756006  3.8184801
[17,] -3.7421402  3.1095511
[18,] -3.3345238  3.2019636
[19,] -2.2717764  2.9436002
[20,] -4.0122721  2.9230563
[21,] -2.3392715  3.8003348
[22,] -3.3265262  2.2767546
[23,] -2.9726308  2.5225962
[24,] -1.9712435  4.5665173
[25,] -2.0126858  4.1387376
[26,] -4.3339338  2.9478143
[27,] -3.3677454  3.0854448
[28,] -2.2030091  2.0877410
[29,] -5.5590850  2.2508092
[30,] -2.6464799  3.4467805
[31,]  3.4467805 -2.6464799
[32,]  2.2508092 -5.5590850
[33,]  2.0877410 -2.2030091
[34,]  3.0854448 -3.3677454
[35,]  2.9478143 -4.3339338
[36,]  4.1387376 -2.0126858
[37,]  4.5665173 -1.9712435
[38,]  2.5225962 -2.9726308
[39,]  2.2767546 -3.3265262
[40,]  3.8003348 -2.3392715
[41,]  2.9230563 -4.0122721
[42,]  2.9436002 -2.2717764
[43,]  3.2019636 -3.3345238
[44,]  3.1095511 -3.7421402
[45,]  3.8184801 -1.8756006
[46,]  1.7154397 -2.5145047
[47,]  1.6493587 -2.2651313
[48,]  2.9451070 -2.1482412
[49,]  2.8710214 -2.2518117
[50,]  3.1381491 -1.1335161
[51,]  4.8342991 -5.5458617
[52,]  1.1543562 -2.0864165
[53,]  0.7115766 -2.6114337
[54,]  2.0327000 -3.3607000
[55,]  1.6829283 -3.7806618
[56,]  2.0106003 -4.1325932
```

```
[57,]   2.5887270 -2.4508472
[58,]   4.0986464 -3.5911897
[59,]   1.8276559 -2.5365722
[60,]   4.3175912 -4.0052498
```

Make a plot of x

```
plot(x)
```



**K-means**

The main function in "base" R for K-means clustering is called `kmeans()`:

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

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

Cluster means:
          x        y
1 -3.012788  2.823278
```
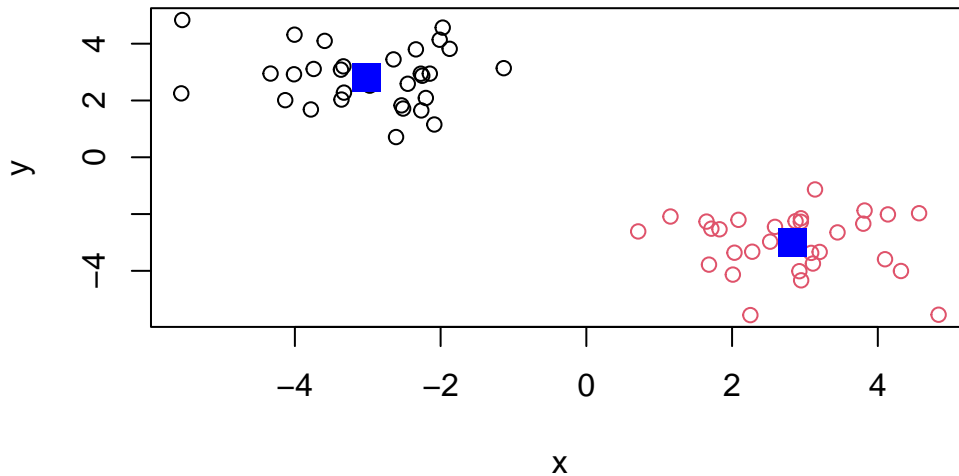
```
2   2.823278 -3.012788

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] 62.38541 62.38541
 (between_SS / total_SS =  89.1 %)

Available components:

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

The `kmeans()` function return a "list" with 9 components. You can see the named components of any list with the `attributes()` function.

```
attributes(km)
```

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

$class
[1] "kmeans"
```

    Q. How many points are in each cluster?

```
km$size
```

```
[1] 30 30
```

    Q. Cluster assignment/membership vector?

```
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
```

    Q. Cluster centers?

```
km$centers
```

```
          x          y
1 -3.012788  2.823278
2  2.823278 -3.012788
```
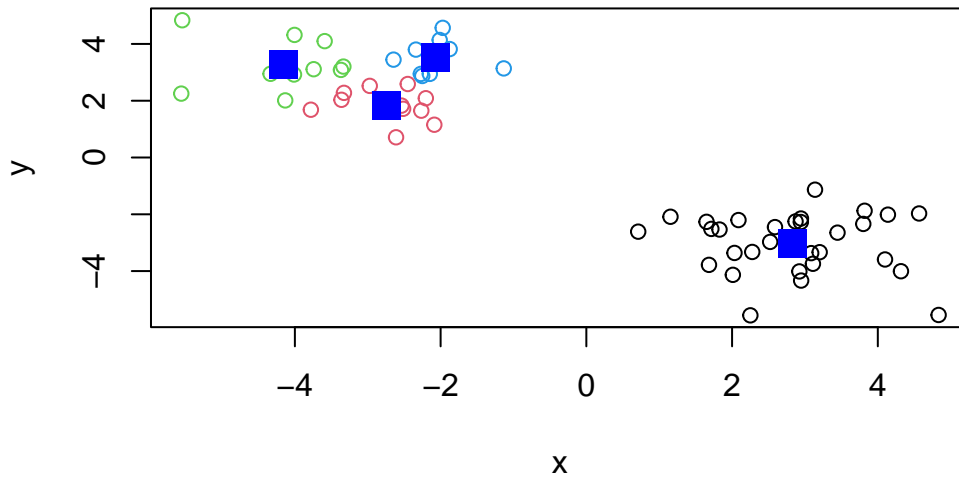
Q. Make a plot of our `kmeans()` results showing cluster assignment using different colors for each cluster/group of points and cluster centers in blue.

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



Q. Run `kmeans()` again on `x` and this cluster into 4 groups/clusters and plot the same result figure as above.

```
km4 <- kmeans(x, centers = 4)
plot(x, col=km4$cluster )
points(km4$centers, col="blue", pch=15, cex=2)
```

**key-point**: K-means clustering is supper popular but can be miss-used. One big limmitation is that it can impose a clustering pattern on your data even if clear natural grouping don't exist - i.e. it does what you tell it to do in therms of `centers`.

### Hierarchical Clustering

The main function in "base" R for hierarchical clustering is called `hclust()`.

You can't just pass our dataset as is into `hclust()` you must give "distance matrix" as input. We can get this from the `dist()` function in R.
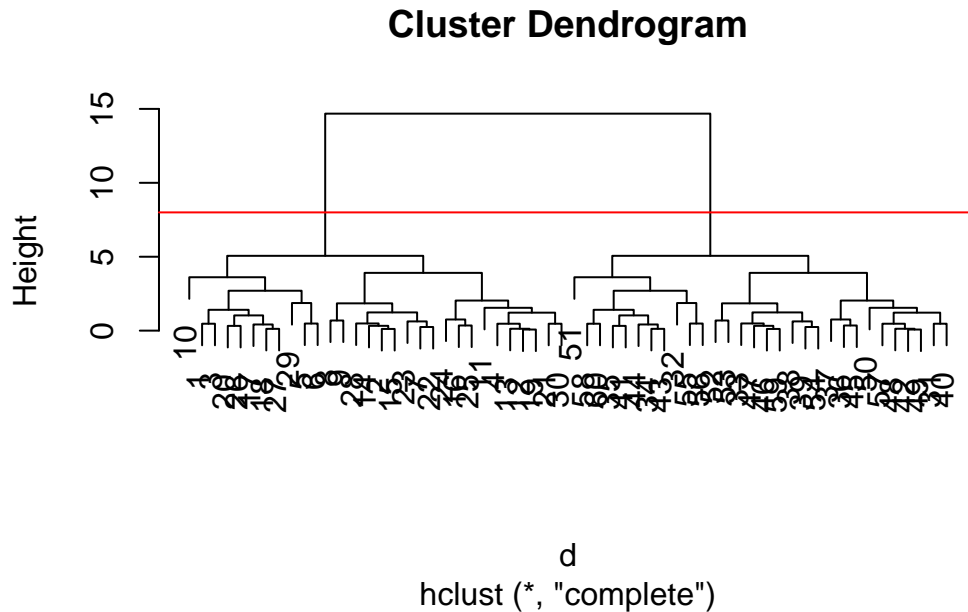
```
d <- dist(x)
hc <- hclust(d)
hc
```

```
Call:
hclust(d = d)

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

The results of `hclust()` don't have a useful `print()` method but do have a special `plot()` method.

```
plot(hc)
abline(h=8, col="red")
```

**Cluster Dendrogram**



d
hclust (*, "complete")

To get our main cluster assignment (membership vector) we need to "cut" the tree at the big goal posts…

```
grps <- cutree(hc, h=8)
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
```
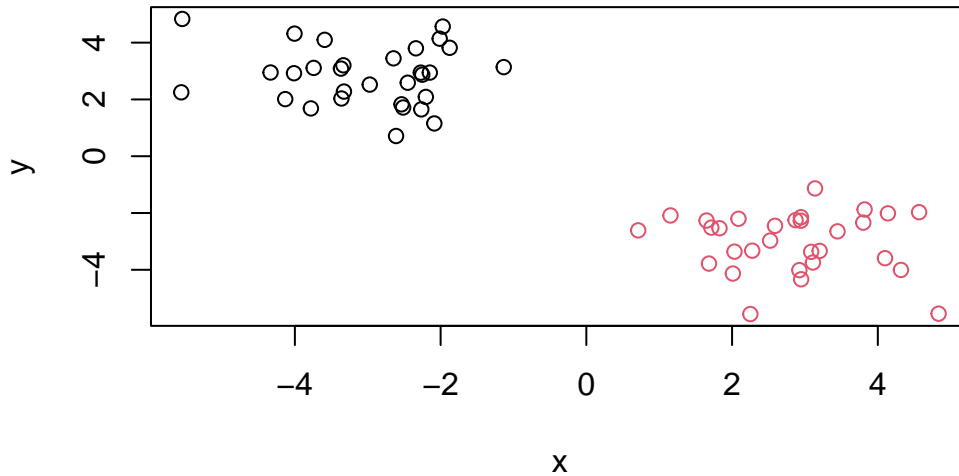
```
table(grps)
```

```
grps
 1  2
30 30
```

8

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



Hierarchical clustering is distinct in that the dendrogram (tree figure) can reveal the potential grouping in your data (unlike K-means)

## Principal Component Analysis (PCA)

PCA is a common and highly useful dimensionality reduction technique used in many fields - particullary bioinformatics.

Here we will analyze some data from the UK on food consumption.

### Data import

```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
```

```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
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
```

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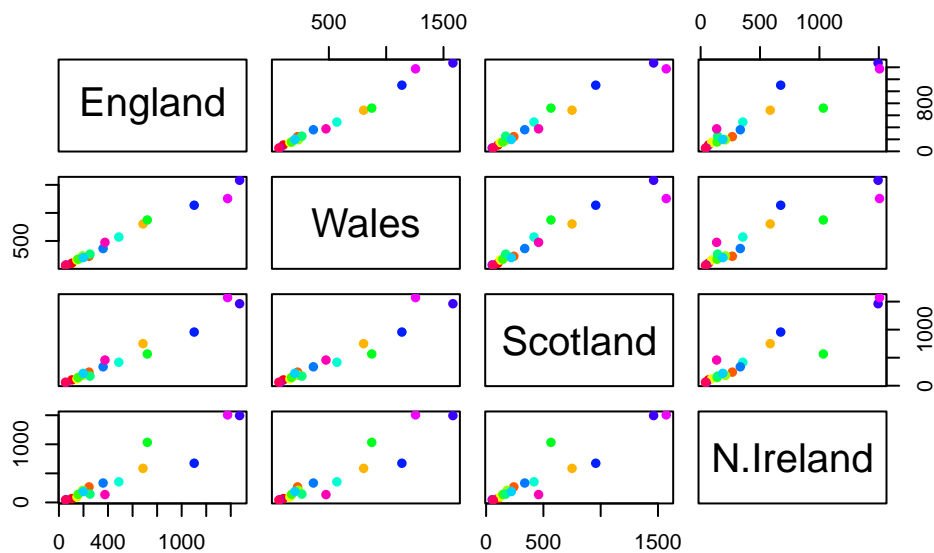
```
barplot(as.matrix(x), beside=F, col=rainbow(nrow(x)))
```



One conventional plot that can be useful is called a "paris" plot.

11

```r
pairs(x, col=rainbow(nrow(x)), pch=16)
```



**PCA to the rescue**

The main function in base R for PCA is called `prcomp()`.

```r
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
```

**Interperting PCA Results**

The `prcomp()` fucntion returns a list object of our results with five attributes/components

```
attributes(pca)
```

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

$class
[1] "prcomp"
```

The two main "results" in here are `pca$x` and `pca$rotation`. The first of these (`pca$x`) contains the scores of the data on the new PC axis - we use these to make our "PCA plot".
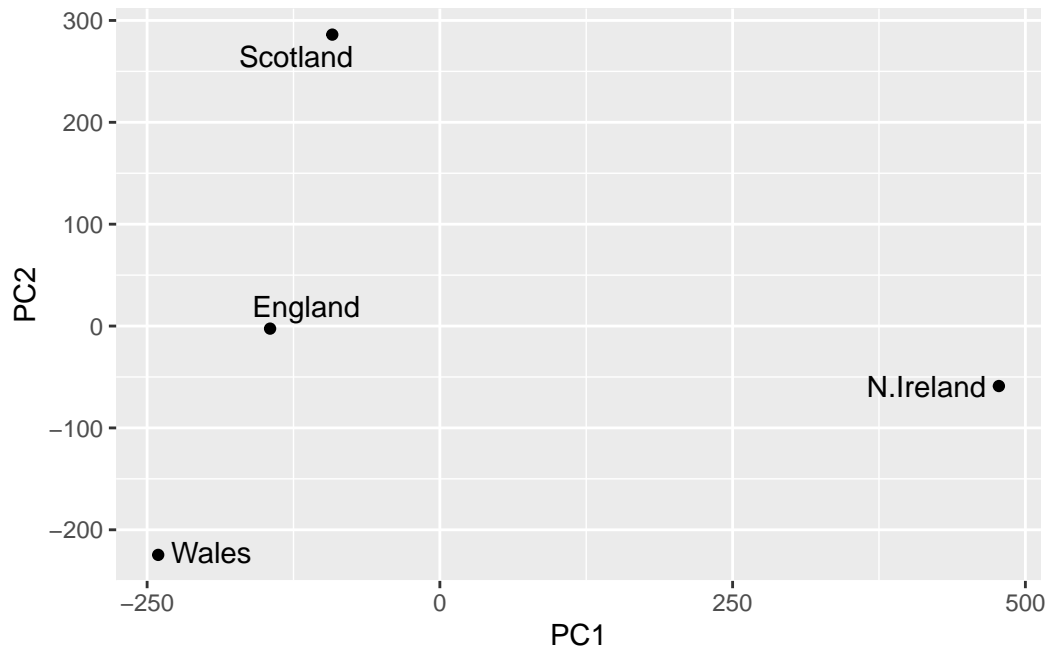
```
pca$x
```

```
                 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
```

```r
library(ggplot2)
library(ggrepel)

# Make a plot of pca$x with PC1 vs PC2
ggplot(pca$x) +
  aes(PC1, PC2, label=rownames(pca$x)) +
  geom_point() +
  geom_text_repel()
```

The second major result is contained in the `pca$rotation` object or component. Let's plot this to see what PCA is picking up...

```
ggplot(pca$rotation) +
  aes(PC1, rownames(pca$rotation)) +
  geom_col()
```