

Class 7: Machine Learning

Derek Chang (PID: A16942232)

Today we are going to learn how to apply different machine learning methods, beginning with clustering:

The goal here is to find groups/clusters in input data.

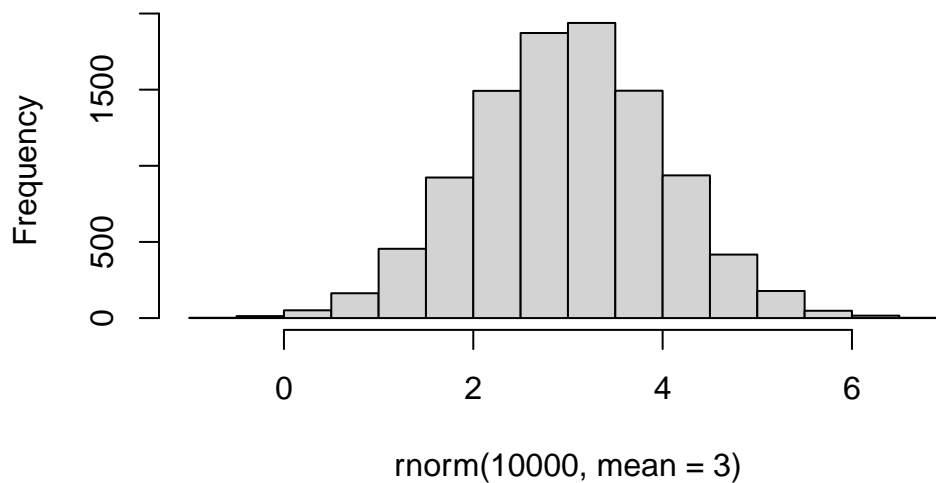
First I will make up some data with clear groups. For this I will use the `rnorm` function.

```
rnorm(10)
```

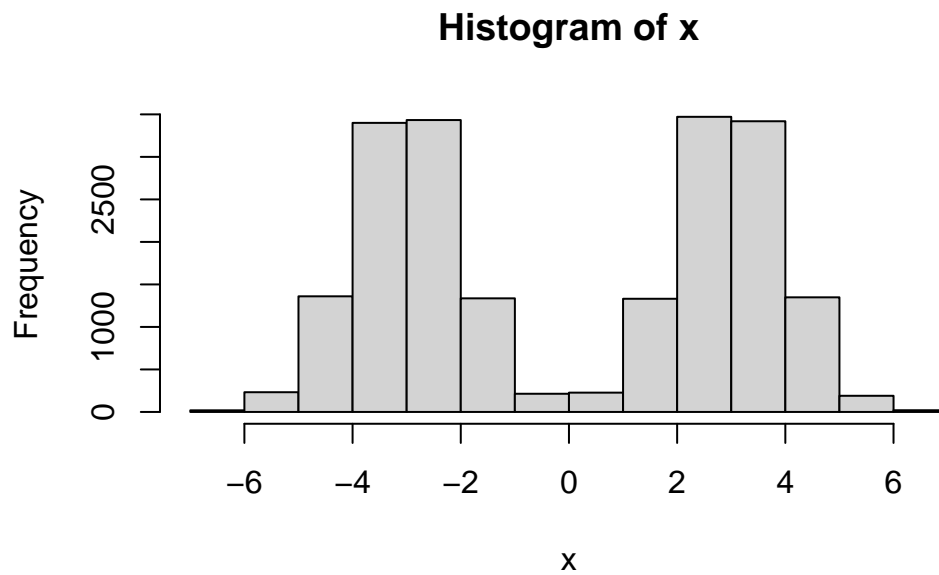
```
[1] -0.55986848  1.20716420 -0.79549189 -0.09963398  0.81469181  0.12504503  
[7] -1.64237939  1.05329969  0.54041293 -1.90624872
```

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

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

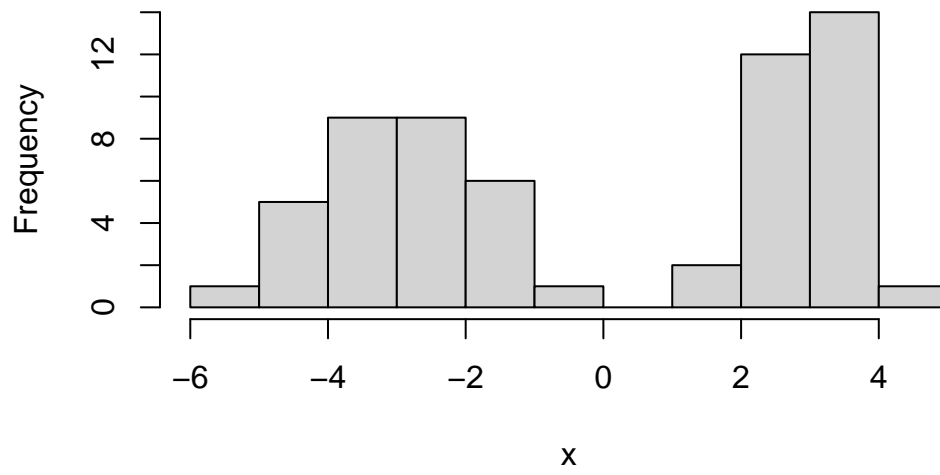


```
n <- 10000  
x <- c(rnorm(n, mean = 3), rnorm(n, mean = -3))  
hist(x)
```



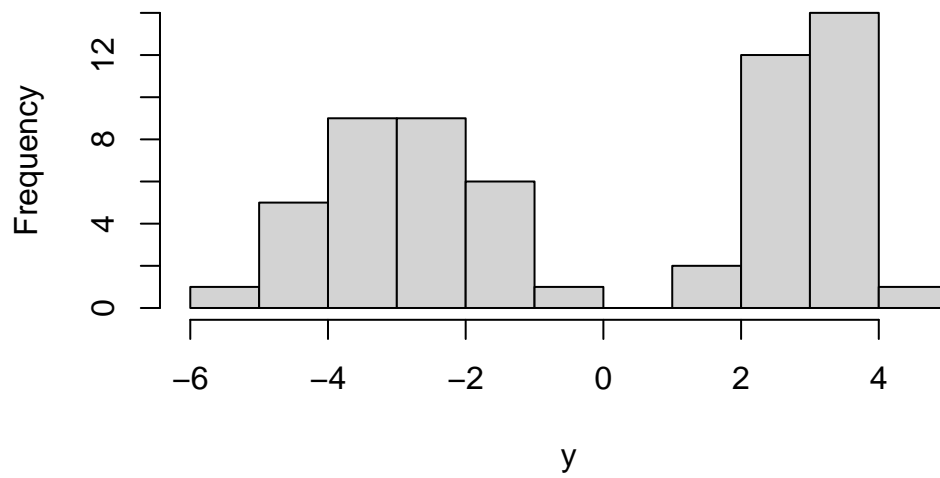
```
n <- 30  
x <- c(rnorm(n, mean = -3), rnorm(n, mean = +3))  
hist(x)
```

Histogram of x



```
y <- rev(x)  
hist(y)
```

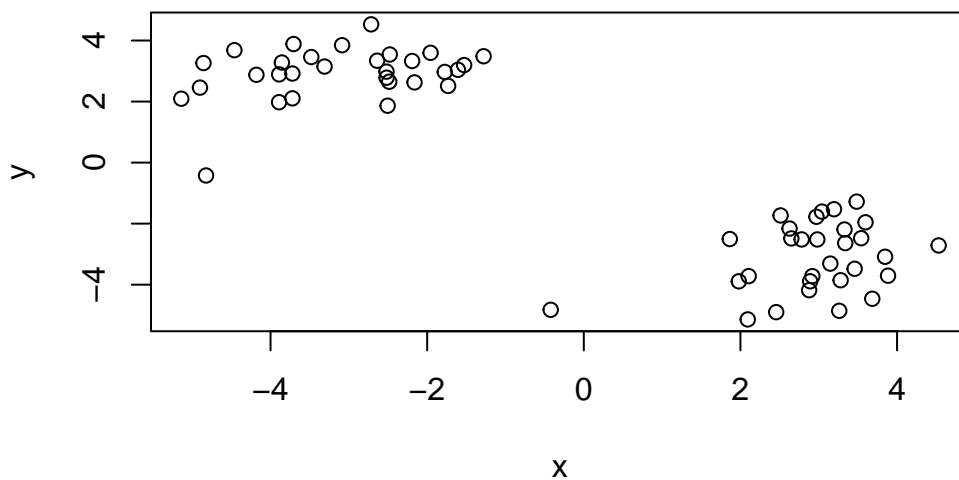
Histogram of y



```
z <- cbind(x,y)
head (z)
```

```
      x      y
[1,] -4.856451 3.261633
[2,] -2.517931 2.781736
[3,] -2.477641 3.541901
[4,] -2.715697 4.530752
[5,] -5.139668 2.093315
[6,] -2.518960 2.981526
```

```
plot(z)
```



Use the `kmeans()` function setting `K` to 2 and `nstart = 20`

Inspect/print the results

Q. How many points are in each cluster?: 30 in each cluster

Q. What 'component' of your result object details - cluster size? - cluster assignment/membership? - Cluster center? Q. plot x colored by the kmeans cluster assignment and add cluster centers as blue points

```
km <- kmeans(z, 2)
km
```

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

Cluster means:

| | x | y |
|---|-----------|-----------|
| 1 | -3.103776 | 2.930182 |
| 2 | 2.930182 | -3.103776 |

Clustering vector:

[illegible]

Within cluster sum of squares by cluster:

```
[1] 58.67352 58.67352
(between_SS / total_SS = 90.3 %)
```

Available components:

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

Results in kmeans object `km`

```
attributes(km$centers)
```

```
$dim
[1] 2 2
```

```
$dimnames
$dimnames[[1]]
[1] "1" "2"
```

```
$dimnames[[2]]
[1] "x" "y"
```

Cluster size?

```
km$size
```

[1] 30 30

```
km$centers
```

| | x | y |
|---|-----------|-----------|
| 1 | -3.103776 | 2.930182 |
| 2 | 2.930182 | -3.103776 |

```
km$cluster
```

[illegible]

Cluster membership?

```
km$cluster
```

[illegible]

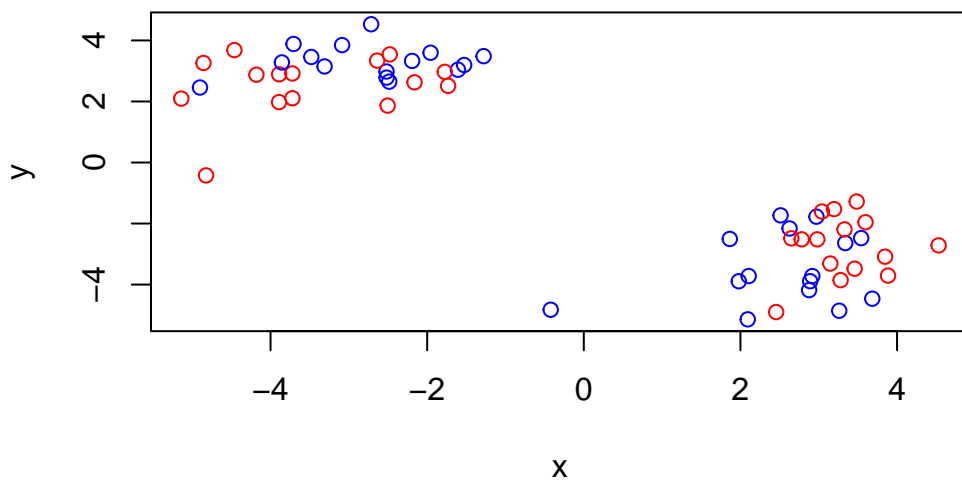
Cluster center?

```
km$centers
```

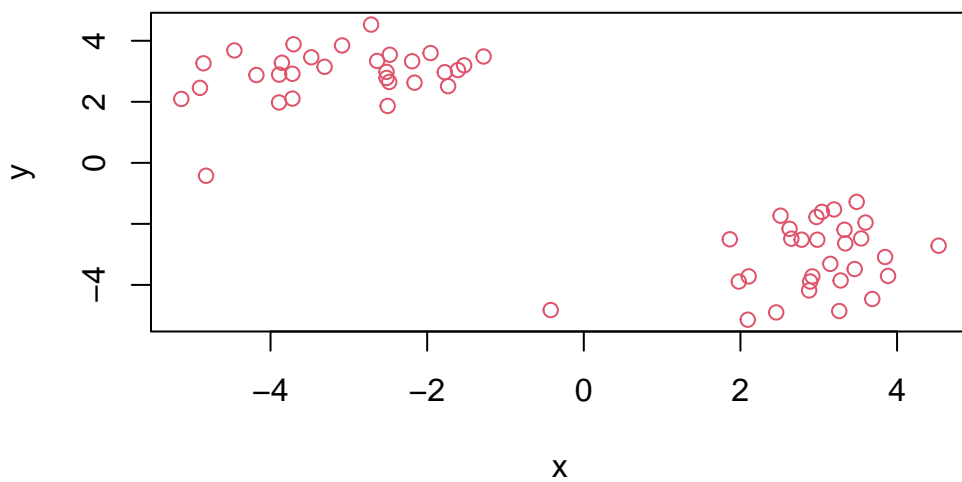
| | x | y |
|---|-----------|-----------|
| 1 | -3.103776 | 2.930182 |
| 2 | 2.930182 | -3.103776 |

R will re-cycle the shorter color vector to be the same length as the longer (number of data points) in z

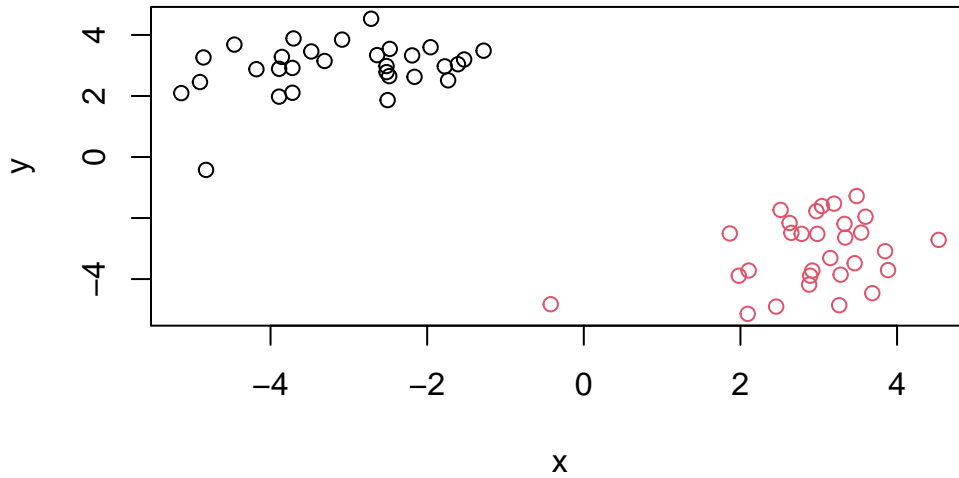
```
plot(z, col = c("red", "blue"))
```



```
plot(z, col = 2)
```

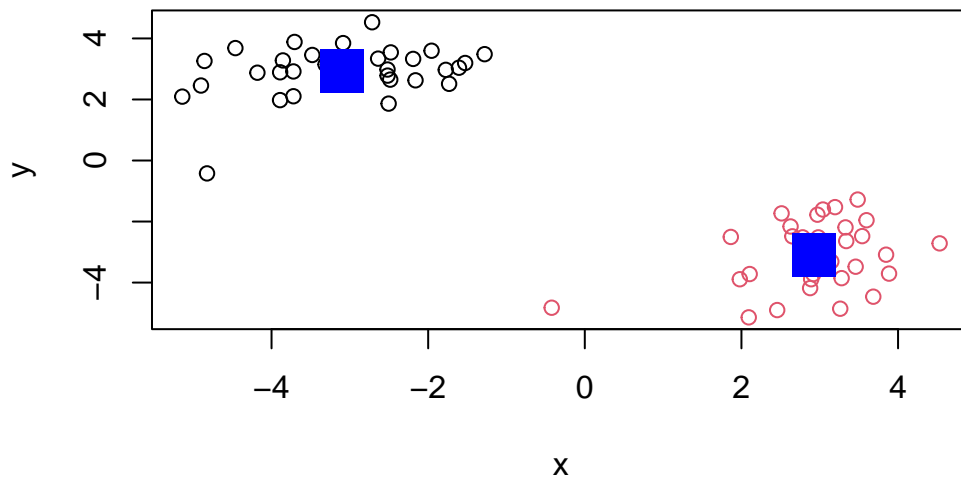


```
plot(z, col = km$cluster)
```



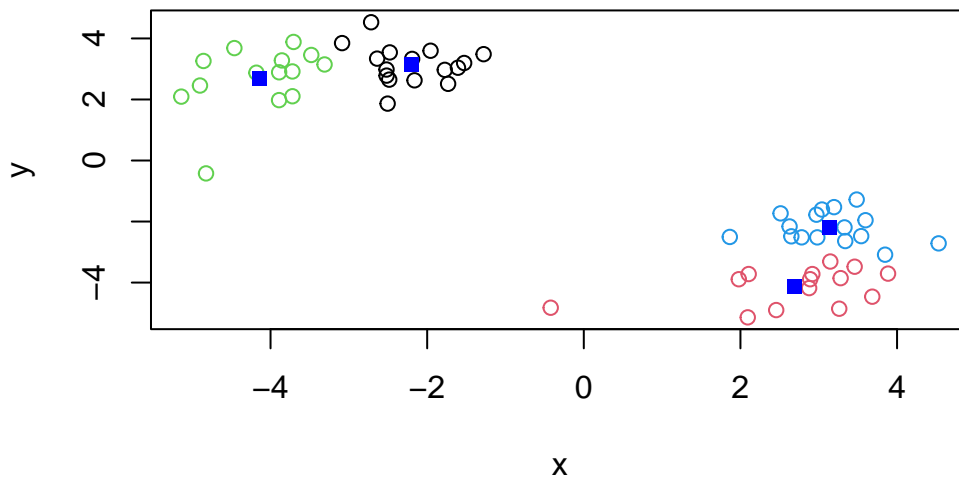
We can use the `points()` function to add new points to an existing plot, like the cluster centers.

```
plot(z, col = km$cluster)
points(km$centers, col = "blue", pch = 15, cex = 3)
```

Q. Can you run kmeans and ask for 4 clutsters and plot the results.

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



Hierarchical Clustering

Let's take our same made-up data `z` and see how `hclust` works.

First we need a distance matrix of our data to be clustered.

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

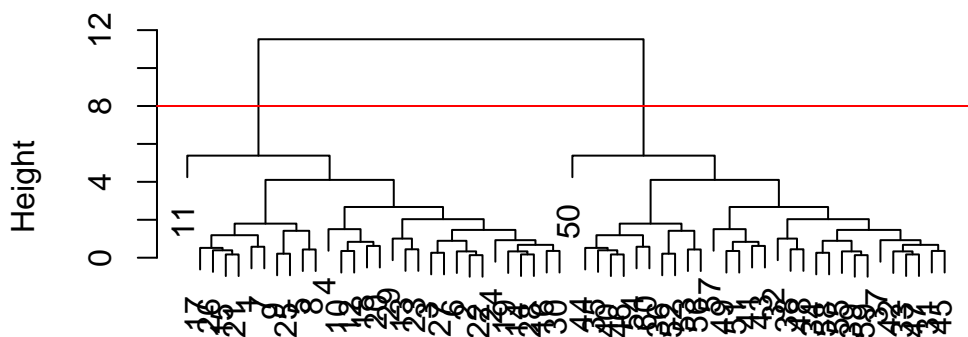
Call:

```
hclust(d = d)
```

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

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

Cluster Dendrogram



```
hclust (*, "complete")
```

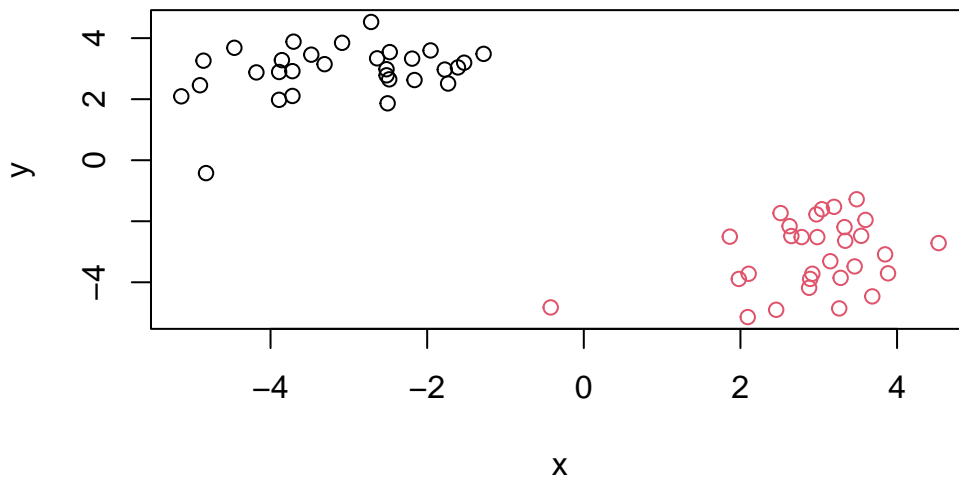
I can get my cluster membership by cutting the tree with the `cutree()` function like so:

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

[illegible]

can you plot `z` colored by our hclust results:

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



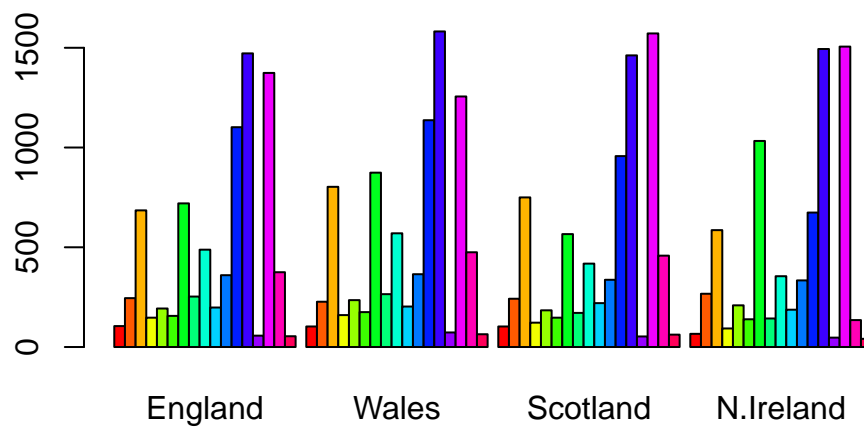
PCA of UK Food Data

Read data from the UK on food consumption in different parts of the UK.

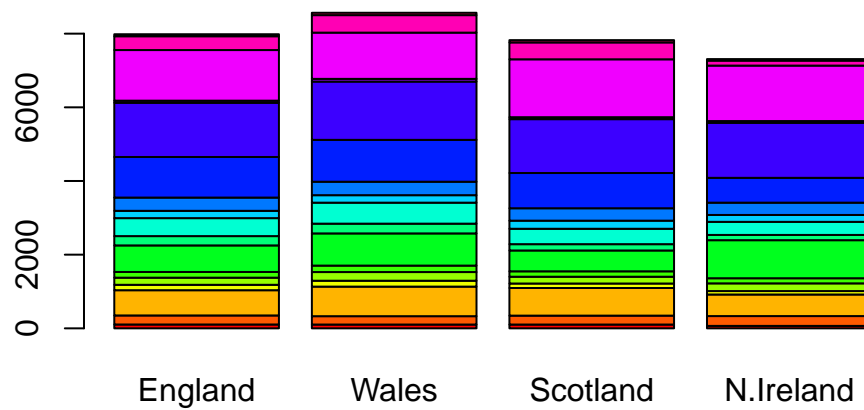
```
url <- "https://tinyurl.com/UK-foods"
x <- read.csv(url, row.names = 1)
head(x)
```

| | 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 |

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

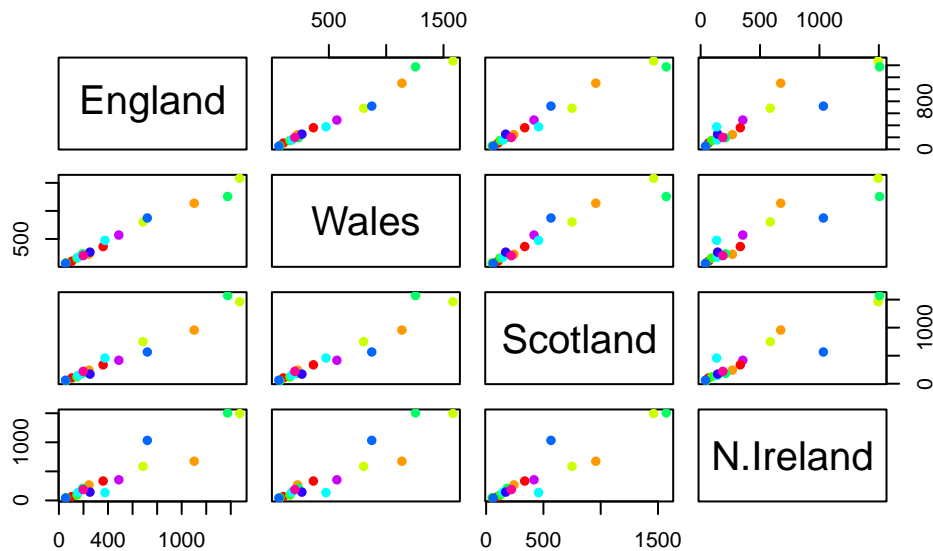


```
barplot(as.matrix(x), beside=F, col=rainbow(nrow(x)))
```



A so-called “Pairs” plot can be useful for small datasets like this

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



It is hard to see structure and trends even in this small data-set. How will we ever do this when we have big datasets with 1000s or 10s of thousands of measured things

PCA to the rescue

Let's see how PCA deals with this dataset. So main function in base R to do PCA is called `prcomp()`

```
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 |

Let's See what is inside this `pca` object that we created from running `prcomp(t(x))`

```
attributes(pca)
```

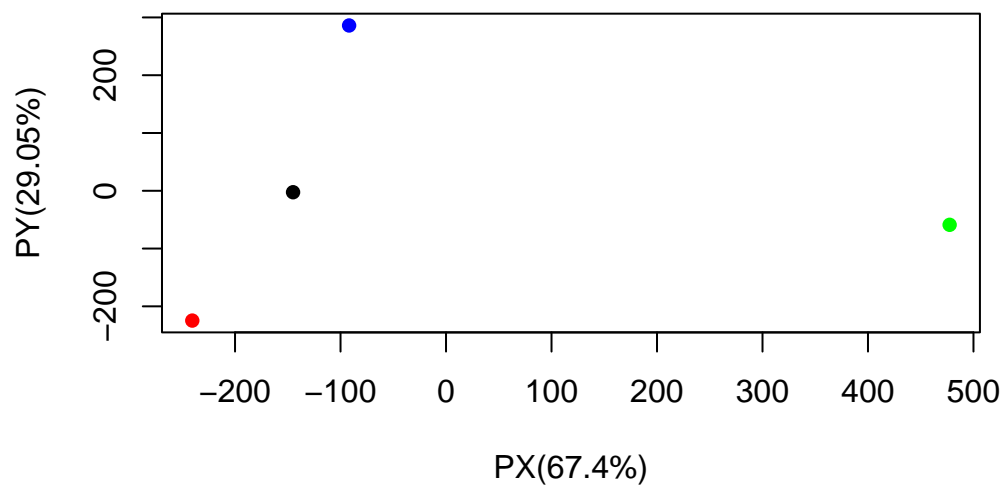
```
$names  
[1] "sdev"      "rotation" "center"    "scale"     "x"
```

```
$class  
[1] "prcomp"
```

```
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 |

```
plot(pca$x[,1], pca$x[,2], col = c("black","red","blue","green"), pch = 16, xlab = "PX(67.4%")
```



Loadings plot

```
par(mar=c(10, 3, 0.35, 0))
barplot( pca$rotation[,1], las=2 )
```

