

Class 7: Machine Learning 1

Jacklyn Jung (PID: A16410001)

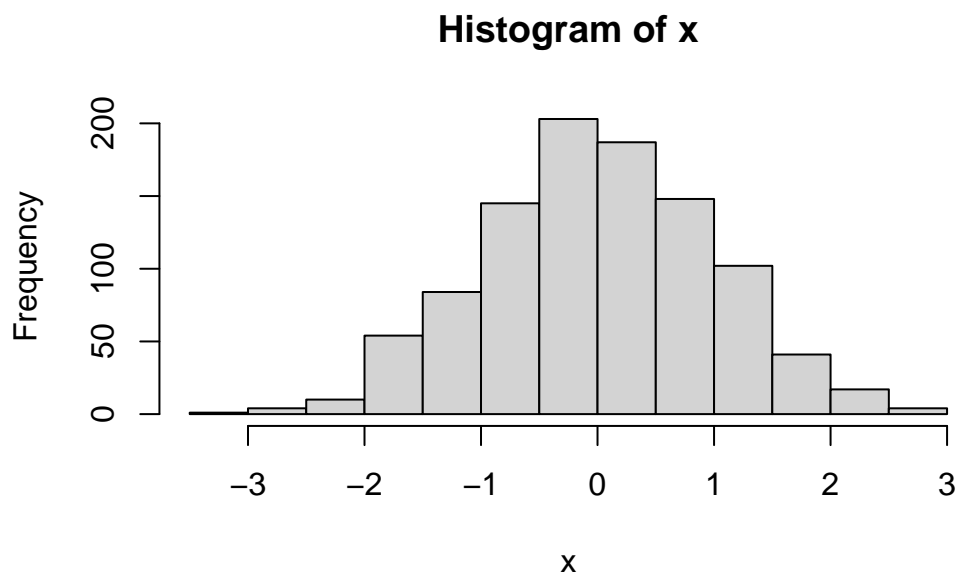
Clustering Methods

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)
```



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

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

```
[1] -3.30148314 -2.23211205 -5.22980332 -5.57722023 -2.84316094 -3.83500133
[7] -0.90352558 -4.11334918 -4.38573154 -2.82409798 -2.49472520 -4.37371447
[13] -4.33350030 -3.28146119 -4.23594252  0.06057436 -2.79877630 -5.16504362
[19] -1.91989112 -2.00076335 -3.24352786 -3.80907680 -3.27851833 -2.87049230
[25] -2.22448866 -3.17069506 -2.47567598 -3.01417182 -2.95624190 -3.24611903
[31]  3.93837087  2.22374823  5.04319823  3.38182938  4.56424599  3.04557143
[37]  2.10468364  3.58145240  1.13732009  2.07057667  2.96269294  5.08661264
[43]  2.85829095  4.20414048  2.44672403  3.13937821  3.49636201  2.96727120
[49]  3.30482983  2.50016454  1.63992143  3.73383871  2.93879456  2.88037212
[55]  3.15305497  4.01693655  3.59277198  3.16423891  2.31631220  3.57881594
```

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

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

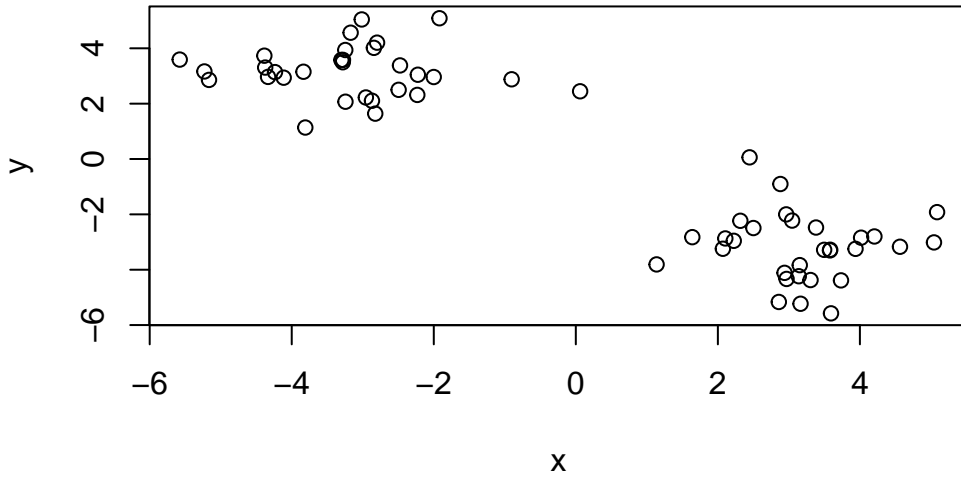
	x	y
[1,]	-3.30148314	3.57881594
[2,]	-2.23211205	2.31631220
[3,]	-5.22980332	3.16423891
[4,]	-5.57722023	3.59277198
[5,]	-2.84316094	4.01693655
[6,]	-3.83500133	3.15305497
[7,]	-0.90352558	2.88037212
[8,]	-4.11334918	2.93879456
[9,]	-4.38573154	3.73383871
[10,]	-2.82409798	1.63992143
[11,]	-2.49472520	2.50016454
[12,]	-4.37371447	3.30482983
[13,]	-4.33350030	2.96727120
[14,]	-3.28146119	3.49636201
[15,]	-4.23594252	3.13937821
[16,]	0.06057436	2.44672403
[17,]	-2.79877630	4.20414048
[18,]	-5.16504362	2.85829095
[19,]	-1.91989112	5.08661264
[20,]	-2.00076335	2.96269294
[21,]	-3.24352786	2.07057667

```

[22,] -3.80907680  1.13732009
[23,] -3.27851833  3.58145240
[24,] -2.87049230  2.10468364
[25,] -2.22448866  3.04557143
[26,] -3.17069506  4.56424599
[27,] -2.47567598  3.38182938
[28,] -3.01417182  5.04319823
[29,] -2.95624190  2.22374823
[30,] -3.24611903  3.93837087
[31,]  3.93837087 -3.24611903
[32,]  2.22374823 -2.95624190
[33,]  5.04319823 -3.01417182
[34,]  3.38182938 -2.47567598
[35,]  4.56424599 -3.17069506
[36,]  3.04557143 -2.22448866
[37,]  2.10468364 -2.87049230
[38,]  3.58145240 -3.27851833
[39,]  1.13732009 -3.80907680
[40,]  2.07057667 -3.24352786
[41,]  2.96269294 -2.00076335
[42,]  5.08661264 -1.91989112
[43,]  2.85829095 -5.16504362
[44,]  4.20414048 -2.79877630
[45,]  2.44672403  0.06057436
[46,]  3.13937821 -4.23594252
[47,]  3.49636201 -3.28146119
[48,]  2.96727120 -4.33350030
[49,]  3.30482983 -4.37371447
[50,]  2.50016454 -2.49472520
[51,]  1.63992143 -2.82409798
[52,]  3.73383871 -4.38573154
[53,]  2.93879456 -4.11334918
[54,]  2.88037212 -0.90352558
[55,]  3.15305497 -3.83500133
[56,]  4.01693655 -2.84316094
[57,]  3.59277198 -5.57722023
[58,]  3.16423891 -5.22980332
[59,]  2.31631220 -2.23211205
[60,]  3.57881594 -3.30148314

```

```
plot(x)
```



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

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

Cluster means:

	x	y
1	-3.202591	3.169084
2	3.169084	-3.202591

Clustering vector:

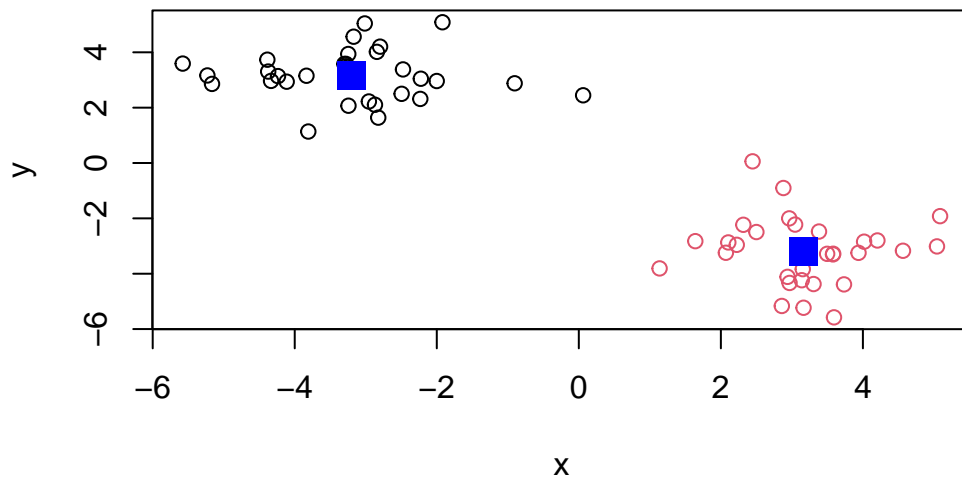
[illegible]

Within cluster sum of squares by cluster:

```
[1] 67.18885 67.18885
      (between_SS / total_SS =  90.1 %)
```

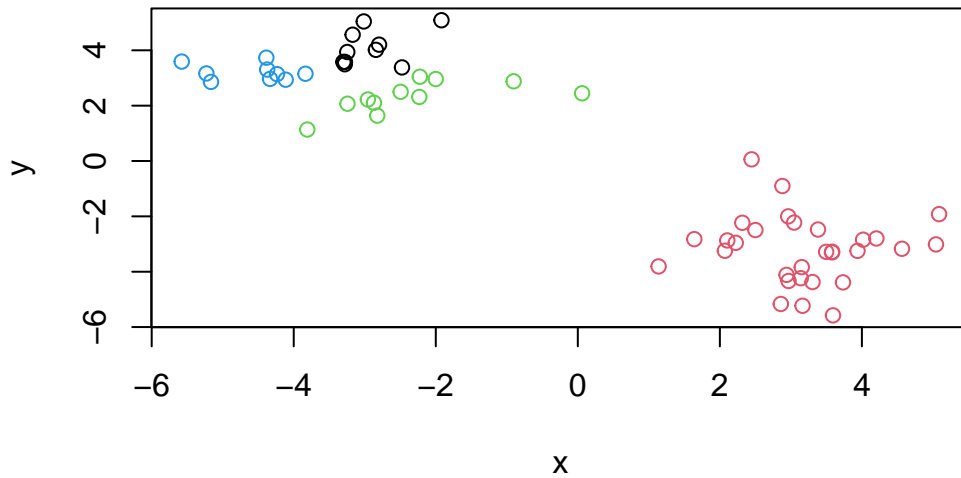
Available components:

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

We can cluster into 4 groups

```
#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.

Hierarchical Clustering

The main base R function for Hierarchical 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
```

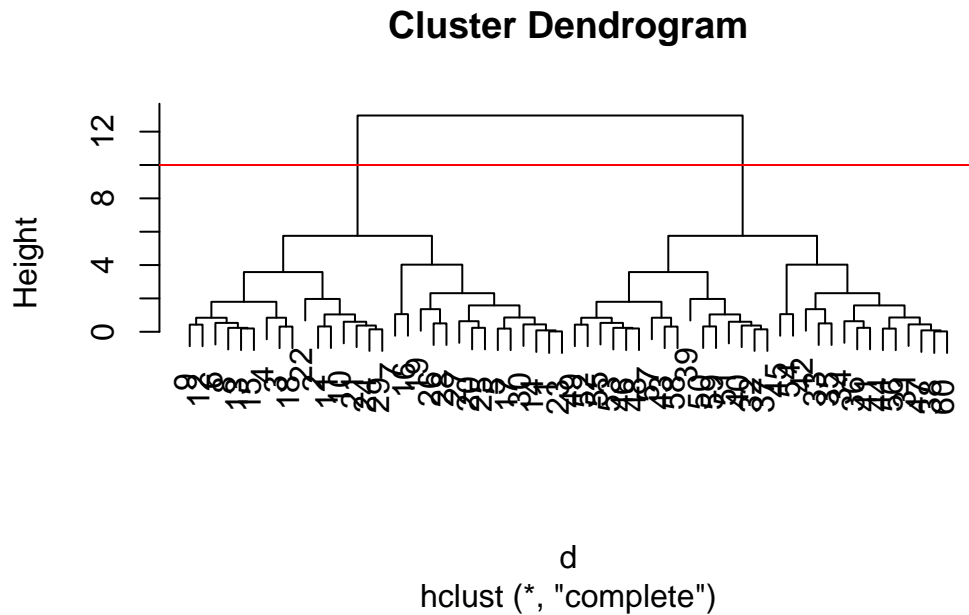
Call:

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

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

Use 'plot()' to see results.

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



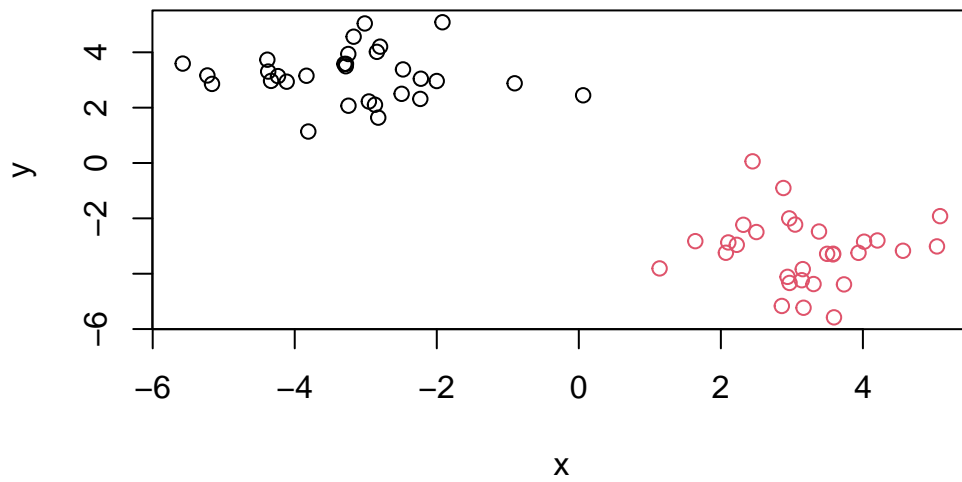
To make the “cut” and get our cluster membership vector we can use the ‘cutree()’ function.

```
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 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 2 2 2 2 2 2 2 2 2 2
```

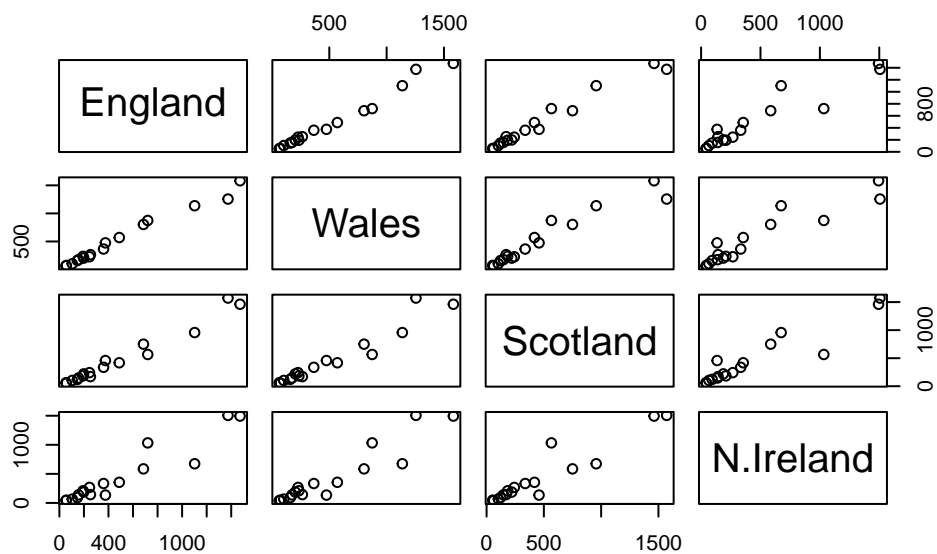
Make a plot of our data colored by hclust results

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

PCA of UK Food Data

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



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

	Wales	Scotland	N.Ireland
105	103	103	66
245	227	242	267
685	803	750	586
147	160	122	93
193	235	184	209
156	175	147	139
720	874	566	1033
253	265	171	143
488	570	418	355
198	203	220	187
360	365	337	334
1102	1137	957	674
1472	1582	1462	1494
57	73	53	47
1374	1256	1572	1506
375	475	458	135
54	64	62	41

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

```
nrow(x)
```

```
[1] 17
```

```
ncol(x)
```

```
[1] 3
```

```
dim(x)
```

```
[1] 17 3
```

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?

- I prefer the second method since the first method can give errors if we run it multiple times.

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

	Scotland	N.Ireland
103	103	66
227	242	267
803	750	586
160	122	93
235	184	209
175	147	139

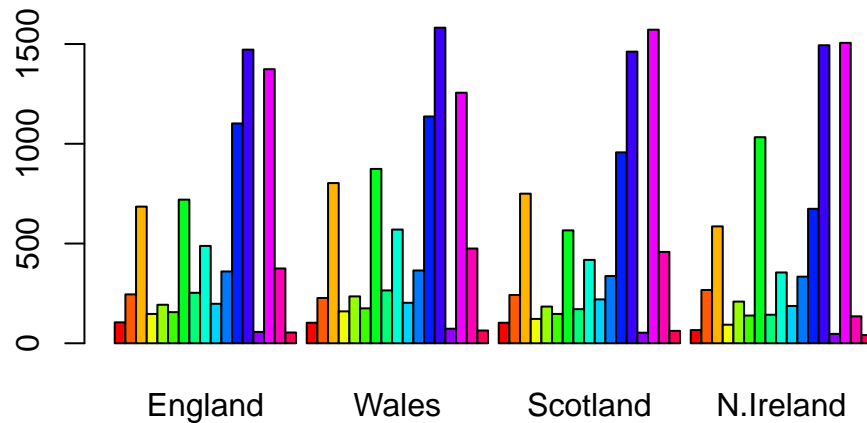
```
x <- read.csv(url, row.names=1)  
head(x)
```

	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

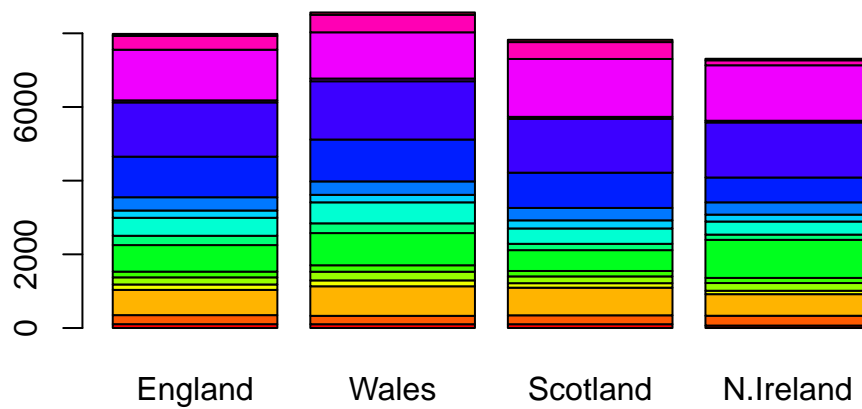
Q3: Changing what optional argument in the above `barplot()` function results in the following plot?

- `beside=T` to `beside=F`
- if `FALSE`, the columns of height are portrayed as stacked bars, and if `TRUE` the columns are portrayed as juxtaposed bars.

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

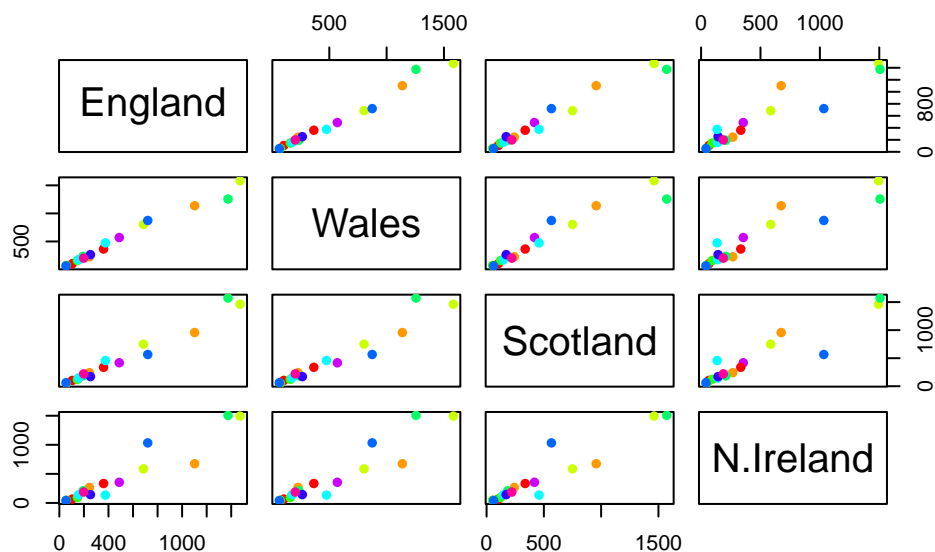


```
barplot(as.matrix(x), beside=F, col=rainbow(nrow(x)))
```



Q5: Generating all pairwise plots may help somewhat. Can you make sense of the following code and resulting figure? What does it mean if a given point lies on the diagonal for a given plot?

```
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 columns.

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

Importance of components:

	PC1	PC2	PC3	PC4
Standard deviation	324.1502	212.7478	73.87622	3.176e-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” or (a.k.a.”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.

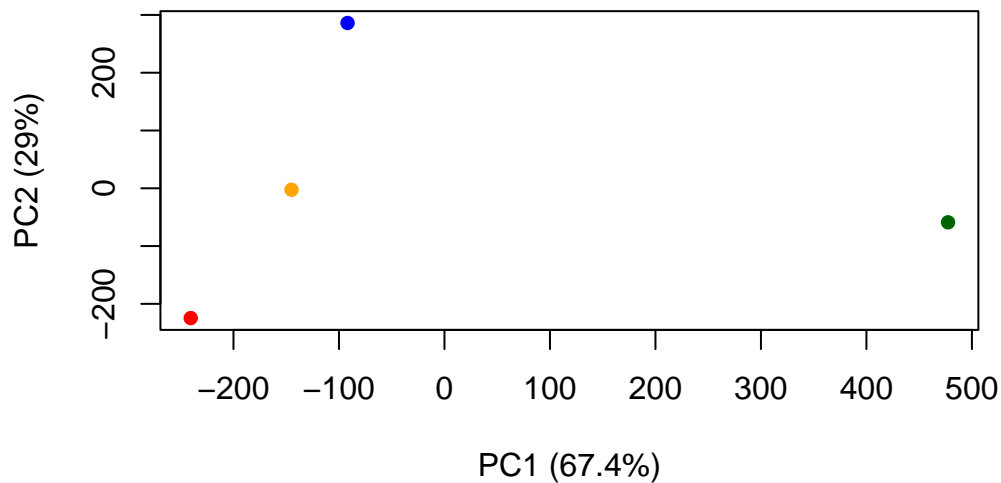
```
pca$x
```

	PC1	PC2	PC3	PC4
England	-144.99315	-2.532999	105.768945	-4.894696e-14
Wales	-240.52915	-224.646925	-56.475555	5.700024e-13
Scotland	-91.86934	286.081786	-44.415495	-7.460785e-13
N.Ireland	477.39164	-58.901862	-4.877895	2.321303e-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	-0.694538519
Carcass_meat	0.047927628	0.013915823	0.06367111	0.489884628
Other_meat	-0.258916658	-0.015331138	-0.55384854	0.279023718
Fish	-0.084414983	-0.050754947	0.03906481	-0.008483145
Fats_and_oils	-0.005193623	-0.095388656	-0.12522257	0.076097502
Sugars	-0.037620983	-0.043021699	-0.03605745	0.034101334
Fresh_potatoes	0.401402060	-0.715017078	-0.20668248	-0.090972715
Fresh_Veg	-0.151849942	-0.144900268	0.21382237	-0.039901917
Other_Veg	-0.243593729	-0.225450923	-0.05332841	0.016719075
Processed_potatoes	-0.026886233	0.042850761	-0.07364902	0.030125166
Processed_Veg	-0.036488269	-0.045451802	0.05289191	-0.013969507
Fresh_fruit	-0.632640898	-0.177740743	0.40012865	0.184072217
Cereals	-0.047702858	-0.212599678	-0.35884921	0.191926714
Beverages	-0.026187756	-0.030560542	-0.04135860	0.004831876
Soft_drinks	0.232244140	0.555124311	-0.16942648	0.103508492
Alcoholic_drinks	-0.463968168	0.113536523	-0.49858320	-0.316290619
Confectionery	-0.029650201	0.005949921	-0.05232164	0.001847469

We can make a plot along PC1.

```
library(ggplot2)
contrib <- as.data.frame(pca$rotation)
ggplot(contrib) +
  aes(PC1, rownames(contrib)) +
  geom_col()
```