



# class 07: 2 Machine Learning

Kalodiah Toma (PID: A07606689)

```
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       |

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] 4
```

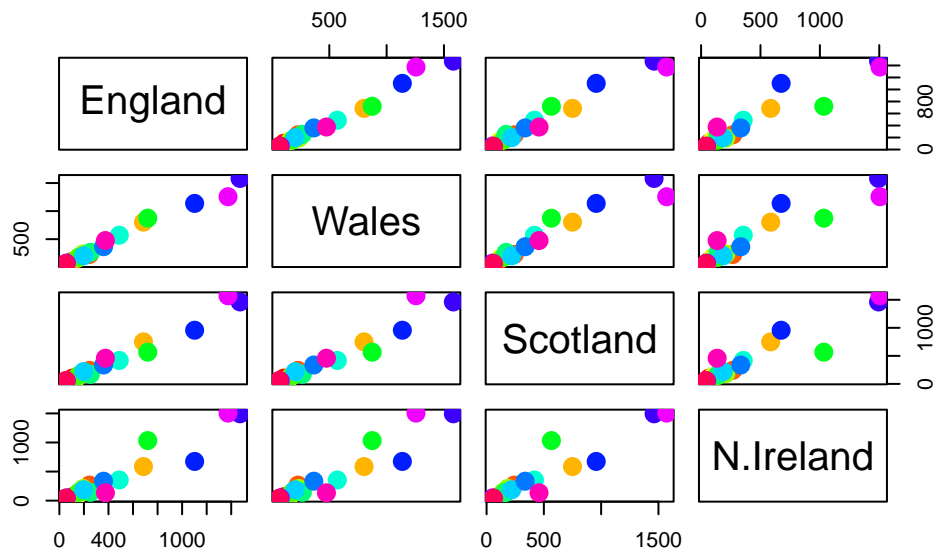
```
head(x)
```

|               | 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       |

```
tail(x)
```

|                  | England | Wales | Scotland | N.Ireland |
|------------------|---------|-------|----------|-----------|
| Fresh_fruit      | 1102    | 1137  | 957      | 674       |
| Cereals          | 1472    | 1582  | 1462     | 1494      |
| Beverages        | 57      | 73    | 53       | 47        |
| Soft_drinks      | 1374    | 1256  | 1572     | 1506      |
| Alcoholic_drinks | 375     | 475   | 458      | 135       |
| Confectionery    | 54      | 64    | 62       | 41        |

```
pairs(x, col=rainbow(17), pch=16, cex=2)
```



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? `x <- read.csv(url, row.names=1)` `head(x)` This is the preferred method as it is not self destructive and will be transferred when rendering.

#PCA to the rescue

Help make sense of this data... The main function for PCA in base R is called `prcomp()`

It wants the transpose (with the `t()`) of our food data for analysis

```
dim(x)
```

```
[1] 17  4
```

```
t(x)
```

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

|           | Fresh_potatoes | Fresh_Veg | Other_Veg | Processed_potatoes |
|-----------|----------------|-----------|-----------|--------------------|
| England   | 720            | 253       | 488       | 198                |
| Wales     | 874            | 265       | 570       | 203                |
| Scotland  | 566            | 171       | 418       | 220                |
| N.Ireland | 1033           | 143       | 355       | 187                |

|           | Processed_Veg | Fresh_fruit | Cereals | Beverages | Soft_drinks |
|-----------|---------------|-------------|---------|-----------|-------------|
| England   | 360           | 1102        | 1472    | 57        | 1374        |
| Wales     | 365           | 1137        | 1582    | 73        | 1256        |
| Scotland  | 337           | 957         | 1462    | 53        | 1572        |
| N.Ireland | 334           | 674         | 1494    | 47        | 1506        |

|           | Alcoholic_drinks | Confectionery |
|-----------|------------------|---------------|
| England   | 375              | 54            |
| Wales     | 475              | 64            |
| Scotland  | 458              | 62            |
| N.Ireland | 135              | 41            |

```
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 |

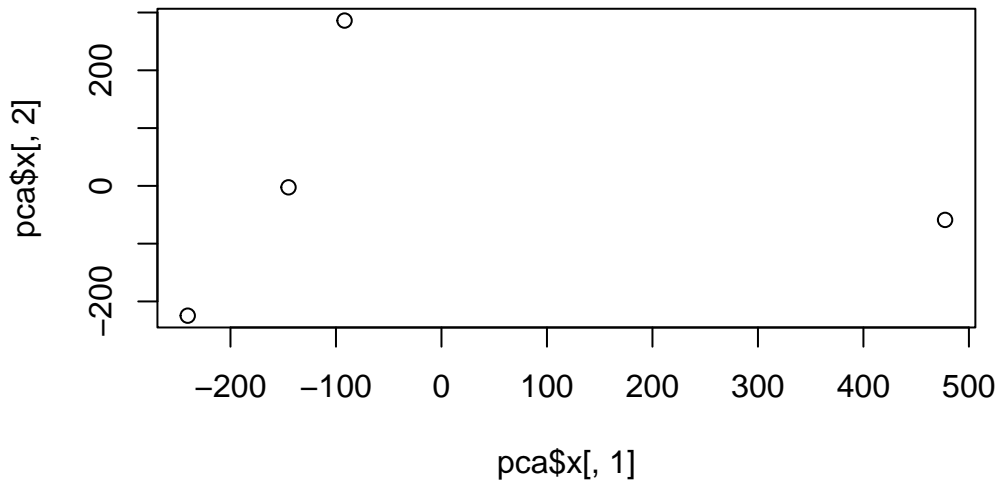
PC1 captures over 67% of data PC2 captures nearly 30% of data so if make a plot of PC1 AND PC1 have 97% of data. Do not need all 17 PC.

One of the main results that folks look for is called the “score plot” aka PC plot, PC1 vs PC2 plot....

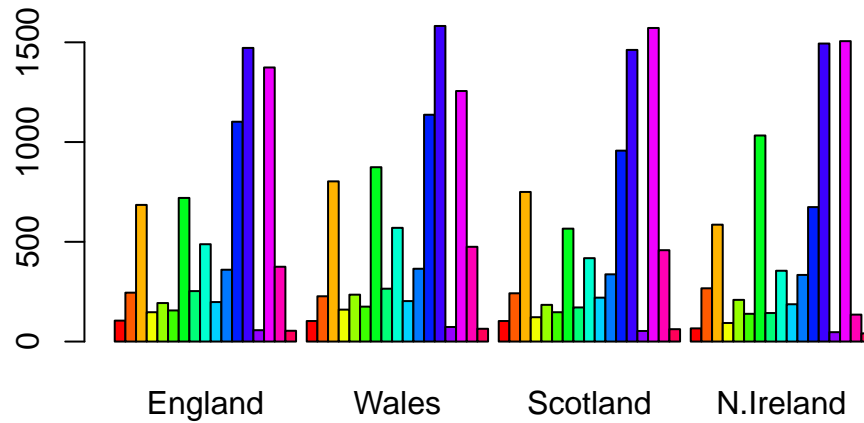
```
pca$x
```

|           | 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  |

```
plot(pca$x[,1],pca$x[,2])
```

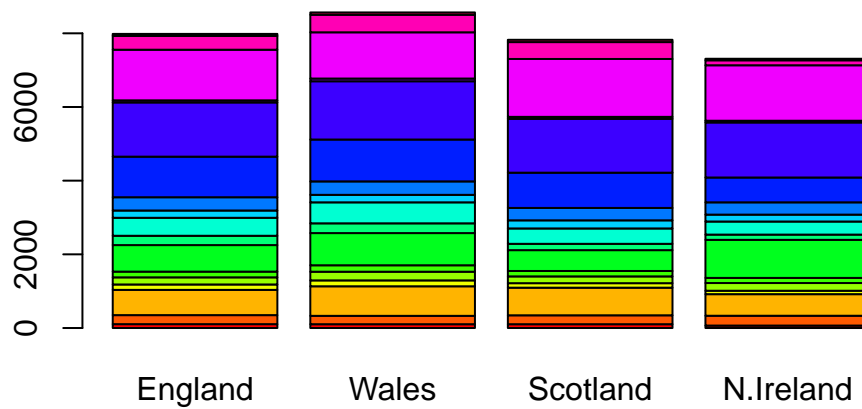


```
barplot(as.matrix(x), beside=T, col=rainbow(nrow(x)))
```



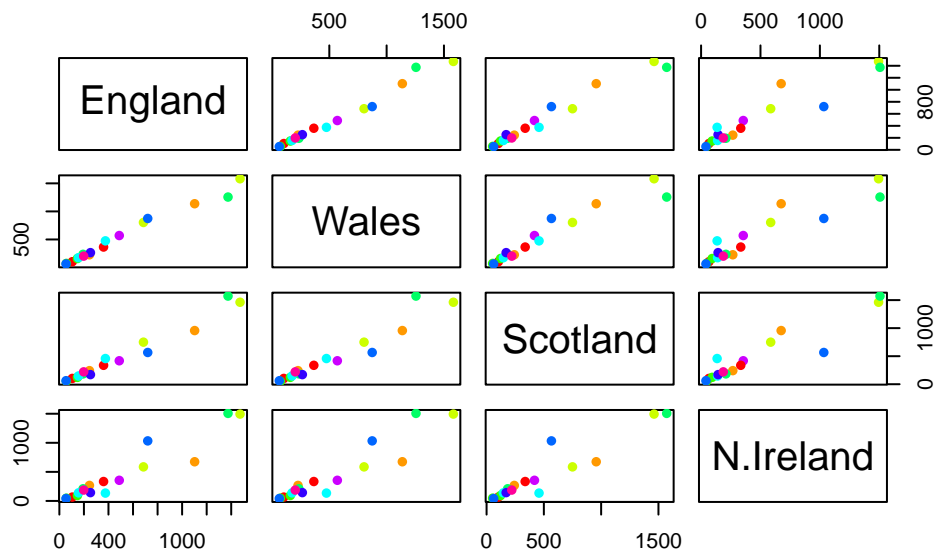
Q3: Changing what optional argument in the above `barplot()` function results in the following plot? leaving `beside` out or setting it to false sets the “value” of `beside` as “0”

```
barplot(as.matrix(x), beside = 0, 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)
```



Q6. What is the main differences between N. Ireland and the other countries of the UK in terms of this data-set? There is one variable (in blue) that is most different and higher in N.Ireland. It is difficult to know what this point is from this type of graph. This is likely alcoholic drinks.

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

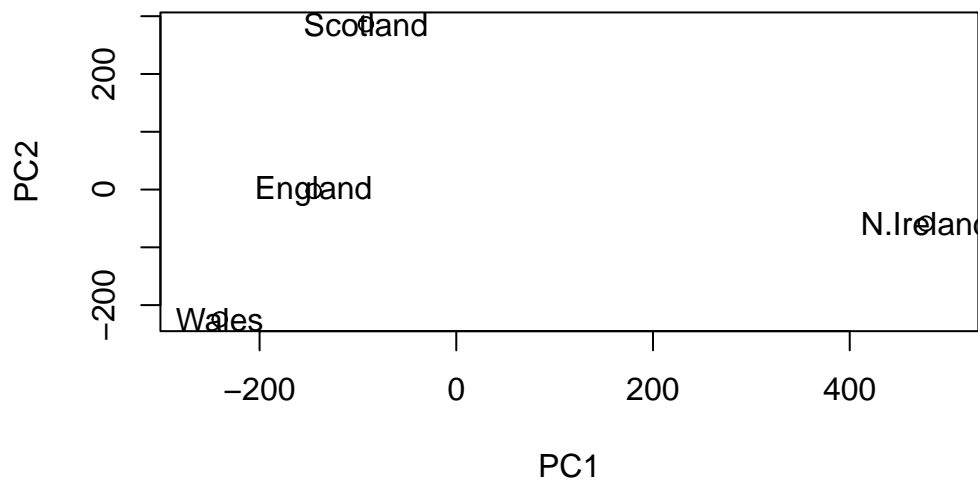
Importance of components:

|                        | 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 |

Q7. Complete the code below to generate a plot of PC1 vs PC2. The second line adds text labels over the data points.

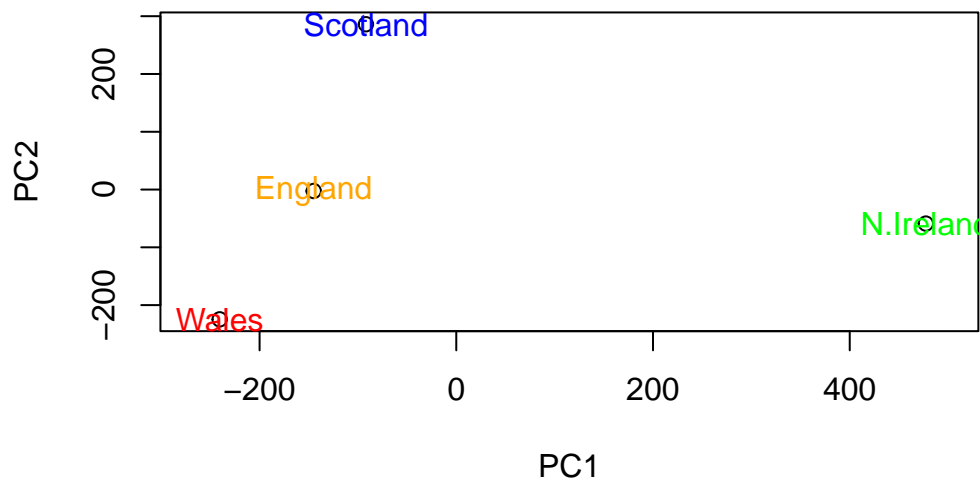
```
plot(pca$x[,1],pca$x[,2], xlab="PC1", ylab="PC2", xlim=c(-270,500))
text(pca$x[,1], pca$x[,2], colnames(x))
```





Q8. Customize your plot so that the colors of the country names match the colors in our UK and Ireland map and table at start of this document.

```
plot(pca$x[,1],pca$x[,2], xlab="PC1", ylab="PC2", xlim=c(-270,500))
text(pca$x[,1], pca$x[,2], colnames(x), col=c("orange","red", "blue", "green"))
```



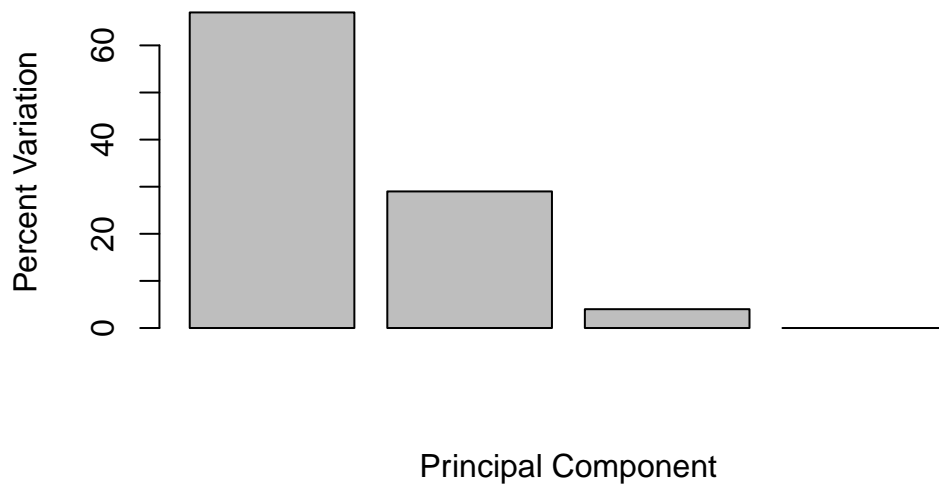
```
v <- round( pca$sdev^2/sum(pca$sdev^2) * 100 )
v
```

```
[1] 67 29 4 0
```

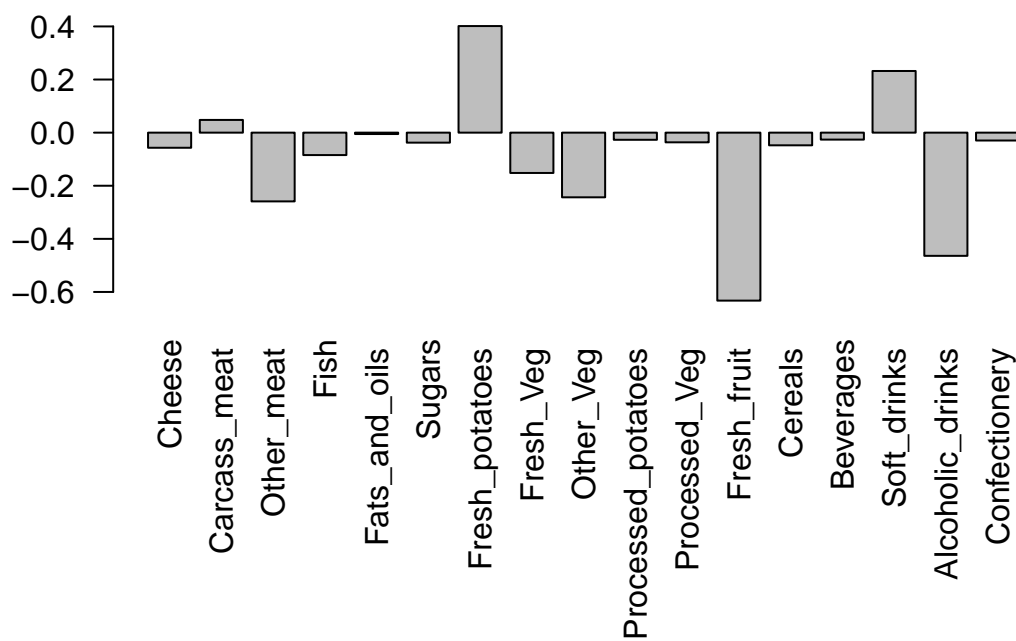
```
z <- summary(pca)
z$importance
```

|                        | PC1       | PC2       | PC3      | PC4          |
|------------------------|-----------|-----------|----------|--------------|
| Standard deviation     | 324.15019 | 212.74780 | 73.87622 | 3.175833e-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")
```

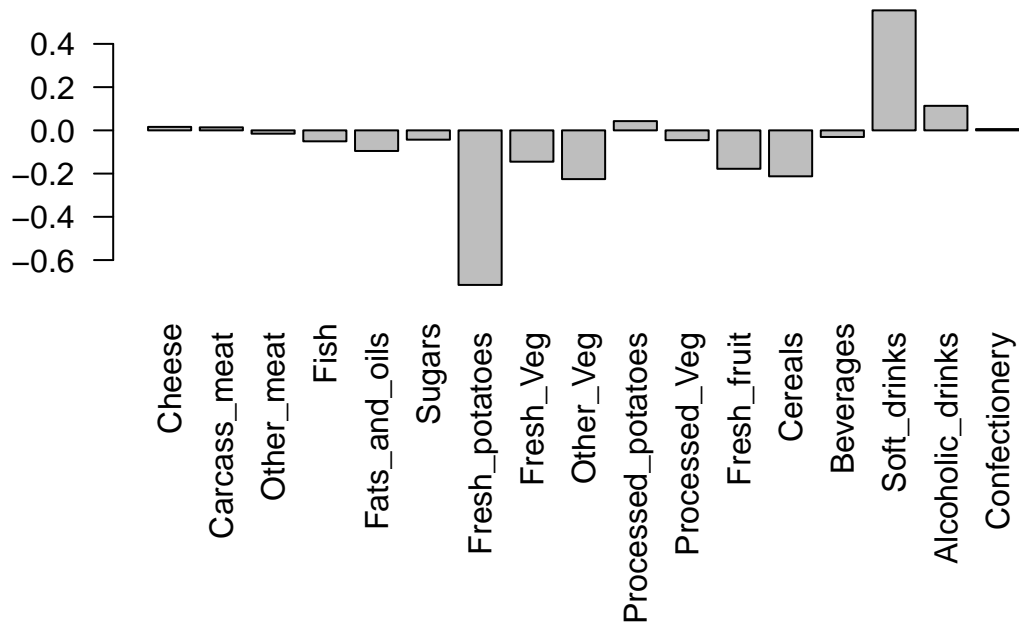


```
par(mar=c(10, 3, 0.35, 0))
barplot( pca$rotation[,1], las=2 )
```



Q9: Generate a similar ‘loadings plot’ for PC2. What two food groups feature prominently and what does PC2 mainly tell us about?

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



PC2 mainly tells us the variance between Wales and Scotland. The two food groups that feature prominently are “fresh potatoes” and “soft-drinks”.