

Class07: Machine Learning 1

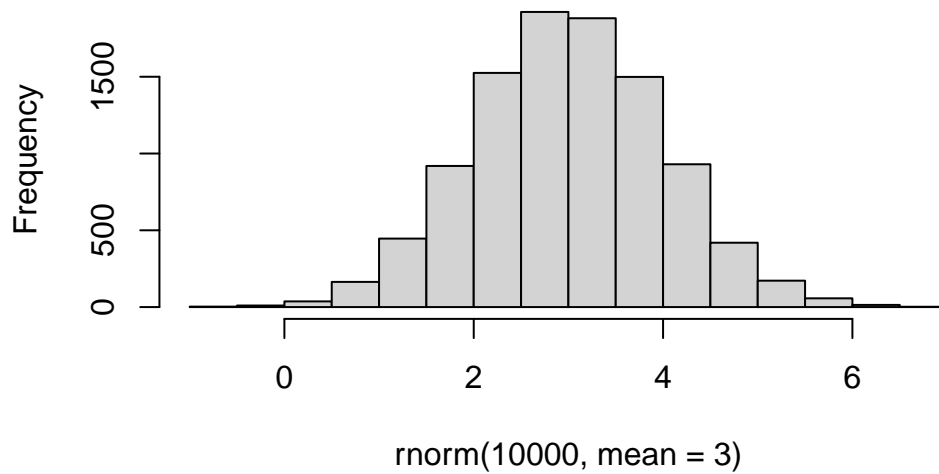
Jiachen Fan (A17662703)

Clustering

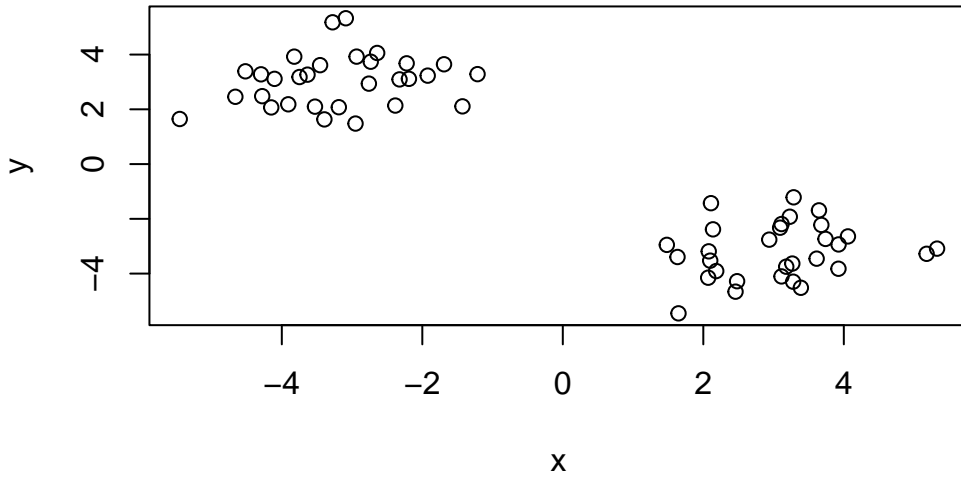
We will start with k-means clustering. To get started let's make some data up:

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

Histogram of rnorm(10000, mean = 3)



```
tmp <- c(rnorm(30,3),rnorm(30, -3))  
x <- cbind(x=tmp, y=rev(tmp))  
plot(x)
```



The main function in R for K-means clustering is called ‘kmeans()’.

```
k <- kmeans(x, centers=2, nstart=20)
k
```

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

Cluster means:

	x	y
1	3.045252	-3.196866
2	-3.196866	3.045252

Clustering vector:

[illegible]

Within cluster sum of squares by cluster:

```
[1] 56.43914 56.43914
(between_SS / total_SS = 91.2 %)
```

Available components:

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

Q1. How many points are in each cluster?

k\$size

[1] 30 30

Q2. The clustering result i.e. membership vector?

```
k$cluster
```

[illegible]

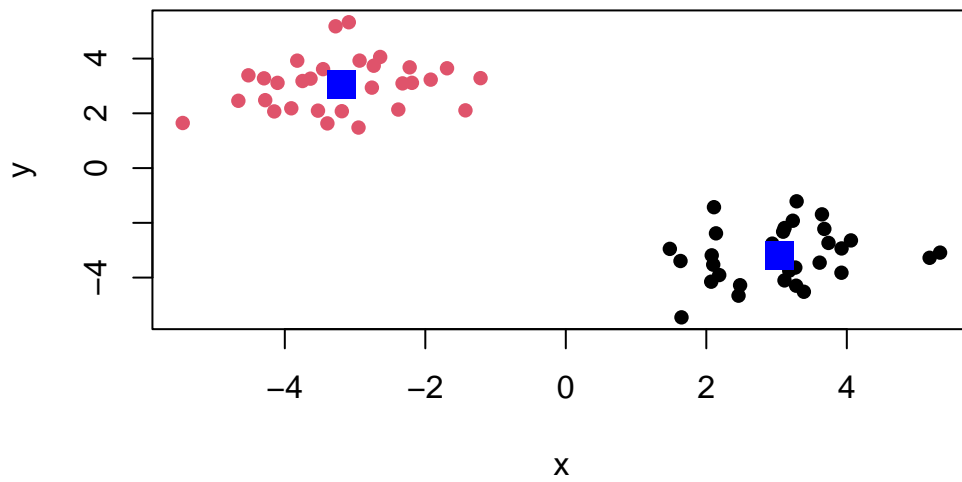
Q3. Cluster centers

k\$centers

	x	y
1	3.045252	-3.196866
2	-3.196866	3.045252

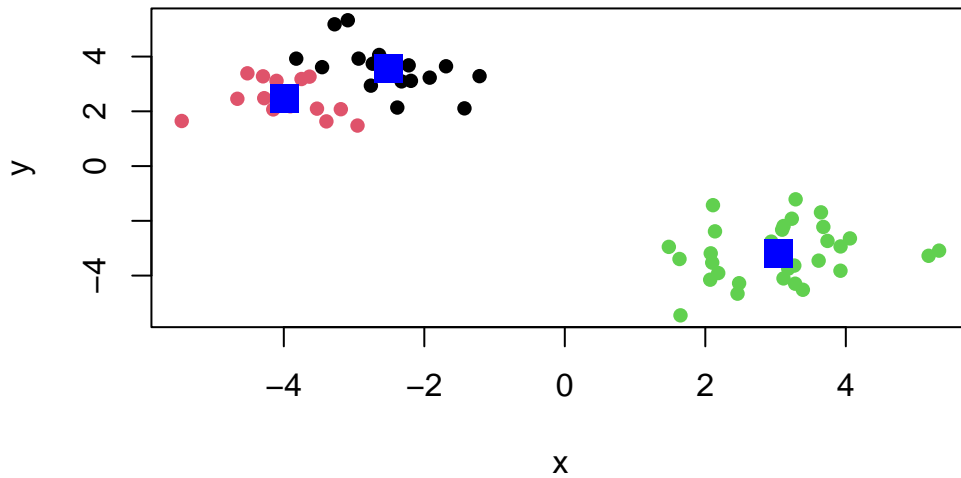
Q4. Make a plot of our data colored by clustering results with optionally the cluster centers shown.

```
plot(x, col=k$cluster, pch=16) # 'pch=16' make points more clear to see
points(k$centers, col='blue', pch=15, cex=2)
```



Q5. Run kmeans again but cluster into 3 groups and plot the results like we did above.

```
k3 <- kmeans(x, centers=3, nstart=20)
plot(x, col=k3$cluster, pch=16)
points(k3$centers, col='blue', pch=15, cex=2)
```



Hierarchical Clustering

It has an advantage in that it can reveal the structure in your data rather than imposing a structure as k-means will.

The main function in ‘base’ R is called ‘hclust()’

It requires a distance matrix as input, not the raw data itself.

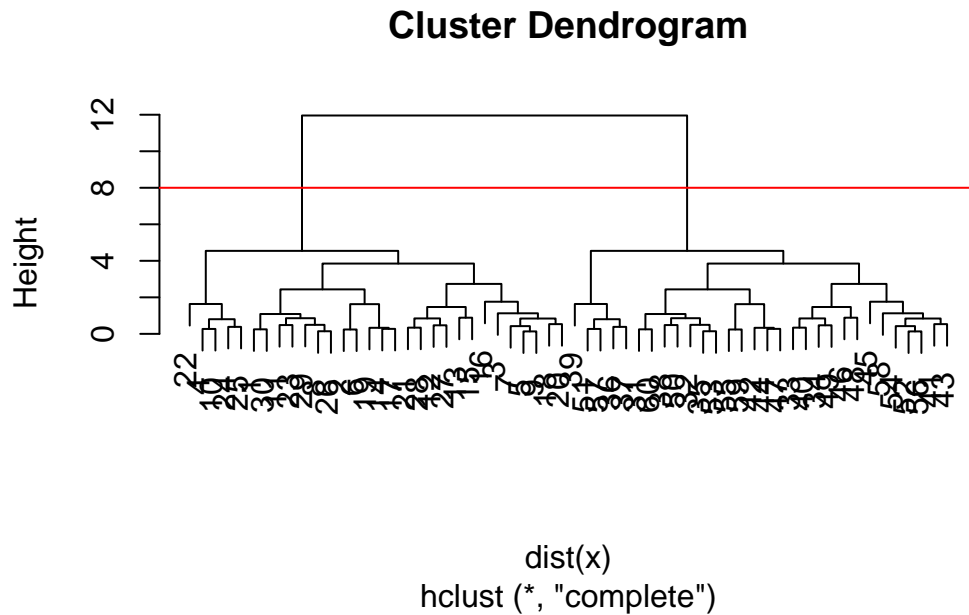
```
hc <- hclust(dist(x))  
hc
```

Call:

```
hclust(d = dist(x))
```

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

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



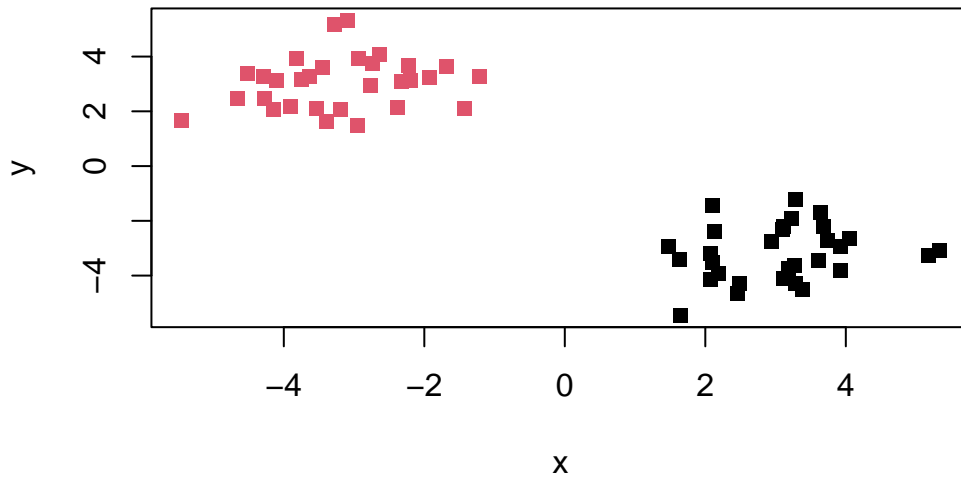
The function to get our clusters from a hclust object is called ‘cutree()’

```
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 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 2 2
```

Q. Plot our hclust results in terms of our data colored by cluster membership.

```
plot(x, col=grps, pch =15)
```



single-link: smallest omplete-link: largest average-link: average

Principal Component Analysis (PCA)

We will work on data from the UK about the strange stuff folks there eat. It has 17 diffenrent foods for 4 countries.

```
url <- 'https://tinyurl.com/UK-foods'
X<- read.csv(url)
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
7	Fresh_potatoes	720	874	566	1033
8	Fresh_Veg	253	265	171	143
9	Other_Veg	488	570	418	355

10	Processed_potatoes	198	203	220	187
11	Processed_Veg	360	365	337	334
12	Fresh_fruit	1102	1137	957	674
13	Cereals	1472	1582	1462	1494
14	Beverages	57	73	53	47
15	Soft_drinks	1374	1256	1572	1506
16	Alcoholic_drinks	375	475	458	135
17	Confectionery	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?

```
dim(X)      # 'dim()'
```

```
[1] 17  5
```

```
# Note how the minus indexing works
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

```
dim(X)
```

```
[1] 17  4
```

There is another way to do it.

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

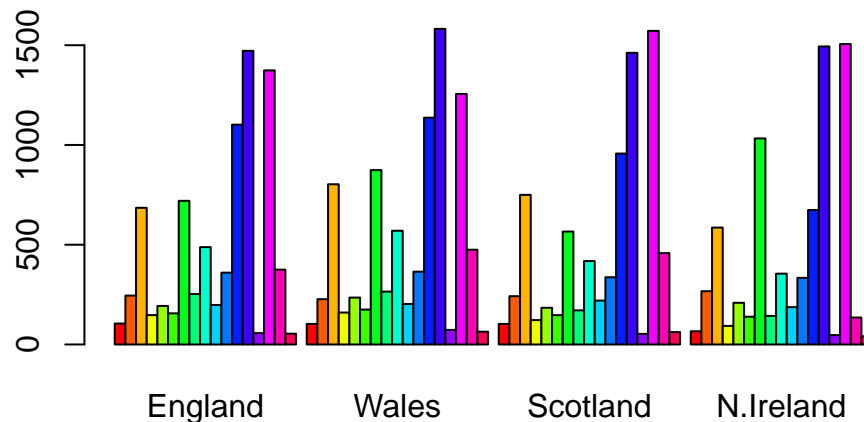

	England	Wales	Scotland	N.Ireland
Cheese	105	103	103	66
Carcass_meat	245	227	242	267
Other_meat	685	803	750	586
Fish	147	160	122	93
Fats_and_oils	193	235	184	209
Sugars	156	175	147	139

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?

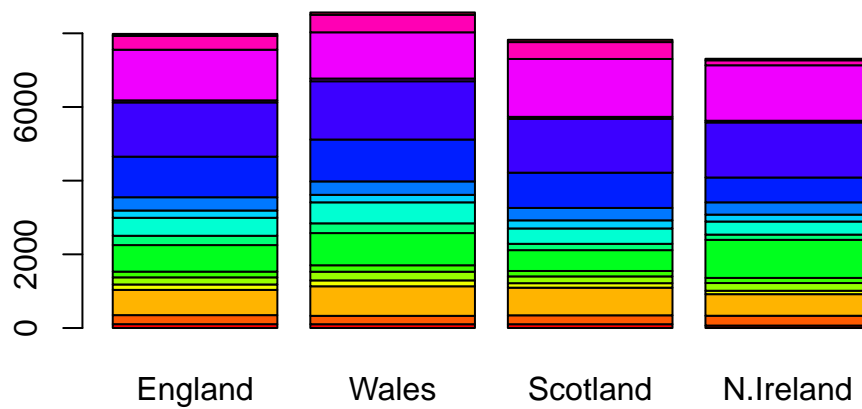
I love the second one since it is concise and make no change to the raw data. If I repeat the first one, it will drop one more column.

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

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



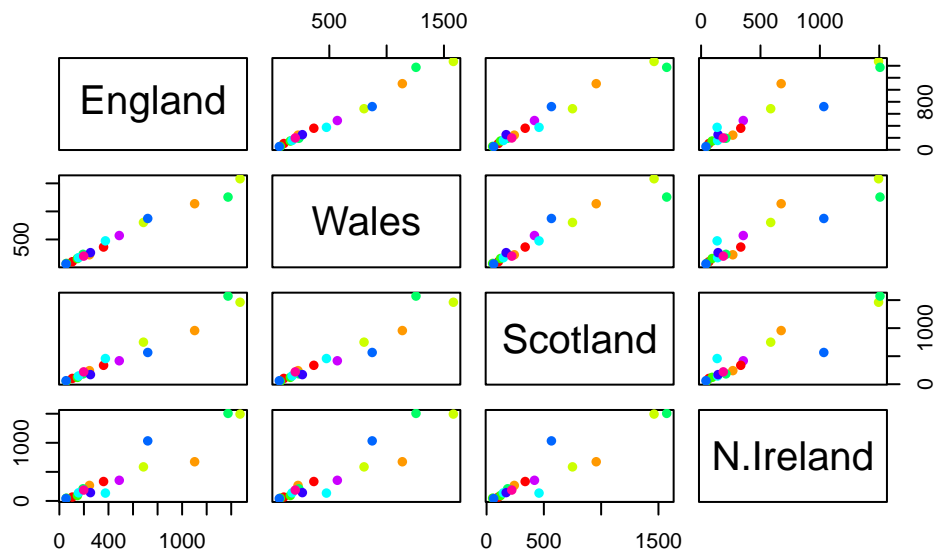
```
barplot(as.matrix(x), beside=FALSE, col=rainbow(nrow(X)))
```



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?

It plots food consumption of every country versus other countries. So it has 12 figures in total.

```
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?

N. Ireland eating habit are quite different from other countries because there are some points out of the line. We can see they eat more 'blue' representing food compared to others.

PCA to the rescue

The main function for PCA is called 'prcomp()'

It wants the transpose (with the 't()') of our food data for analysis

```
# Use the prcomp() PCA function
pca <- prcomp(t(x))
summary(pca)
```

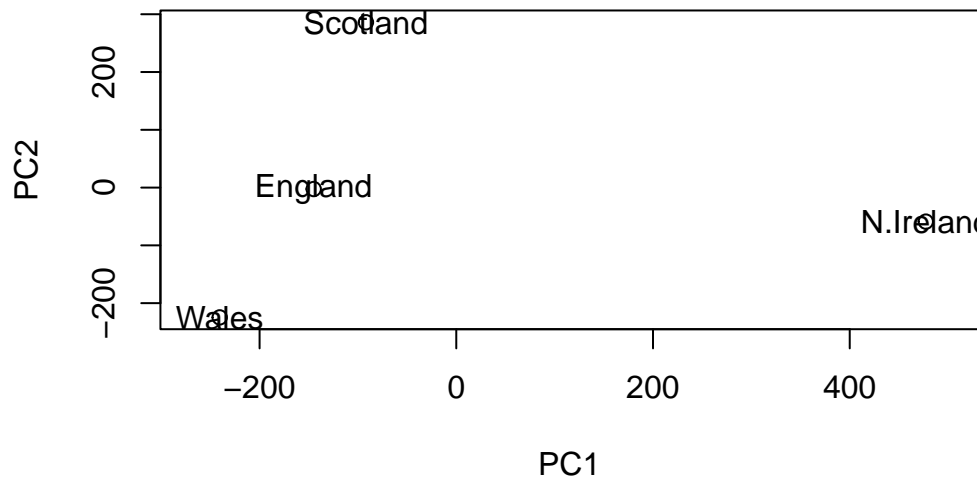
Importance of components:

	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

One of the main results that folks look for is called the ‘score plot’ a.k.a. PC plot, PC1 vs PC2 plot...

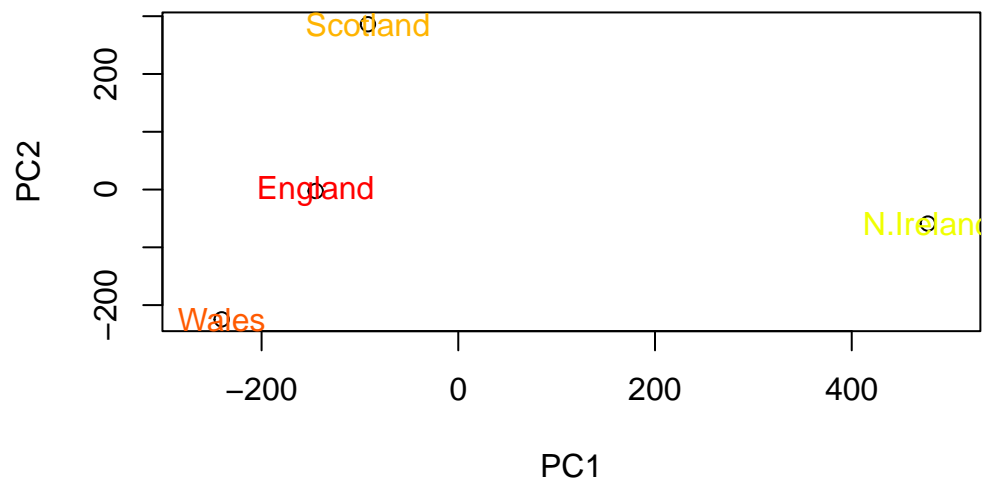
Q7. Complete the code below to generate a plot of PC1 vs PC2. The second line adds text labels over the data points.

```
# Plot PC1 vs PC2
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.

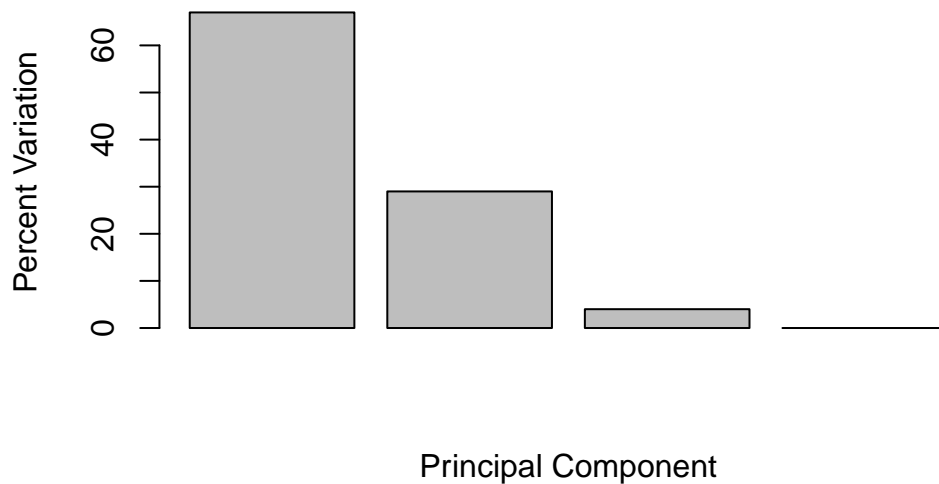
```
# Plot PC1 vs PC2
plot(pca$x[,1], pca$x[,2], xlab="PC1", ylab="PC2", xlim=c(-270,500))
text(pca$x[,1], pca$x[,2], colnames(x), col = rainbow(nrow(x)))
```



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

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

```
barplot(v, xlab="Principal Component", ylab="Percent Variation")
```



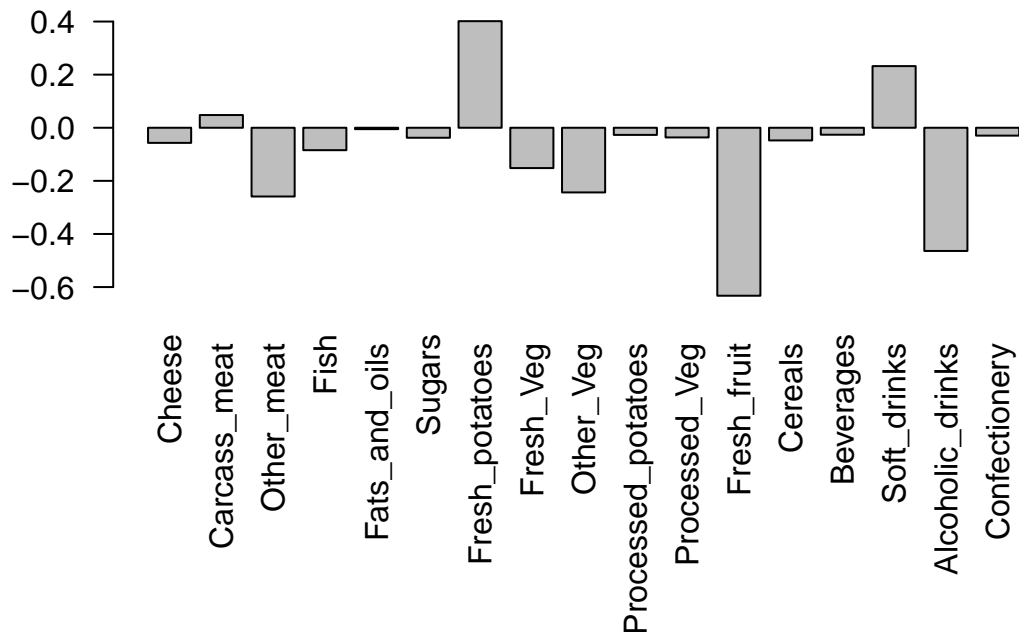
Digging deeper (variable loadings)

```
# Let's focus on PC1 as it accounts for > 90% of variance
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

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

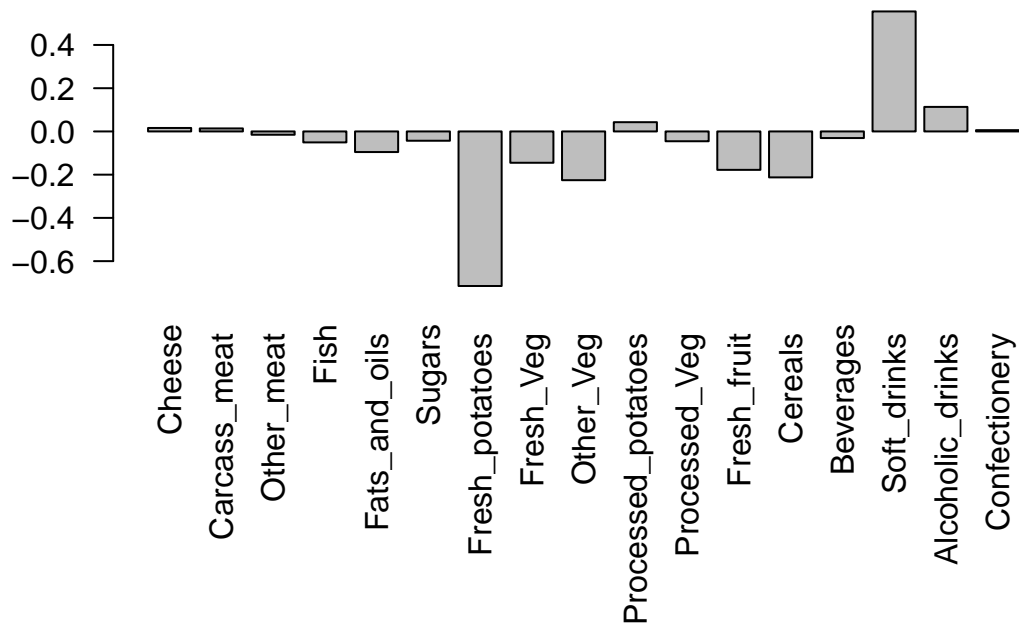


Q9: Generate a similar 'loadings plot' for PC2. What two food groups feature prominently and what does PC2 mainly tell us about?

Fresh_potatoes and Soft_drinks are the two groups contribute most to the variance in PC2.

PC2 mainly tells us the variance between different countries are from original variables: Fresh_potatoes and Soft_drinks. These differences contribute to the variance in countries.

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

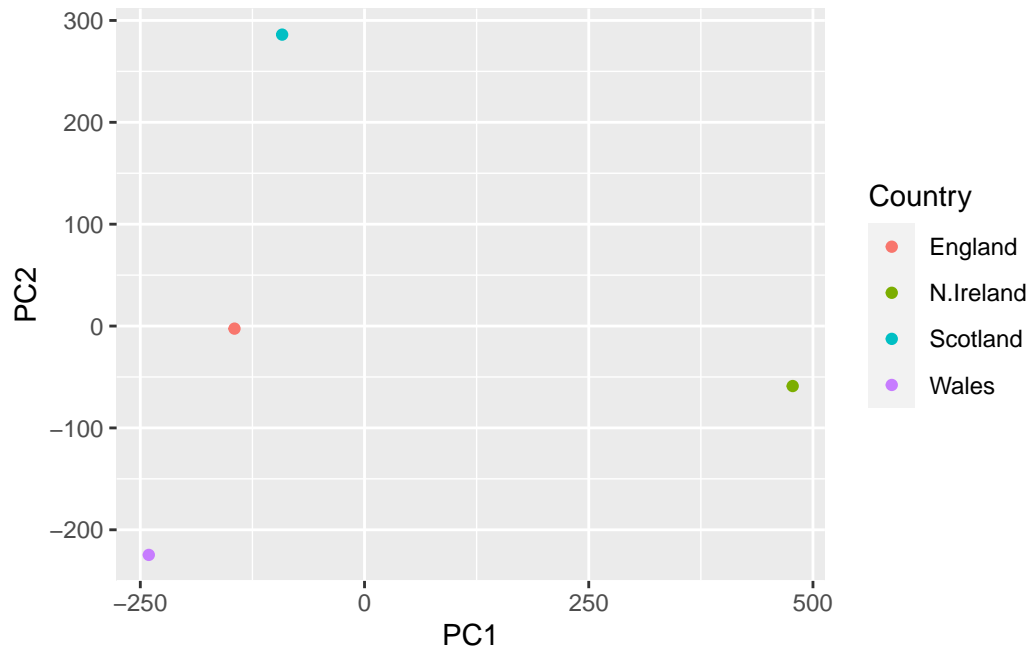


Using ggplot for these figures

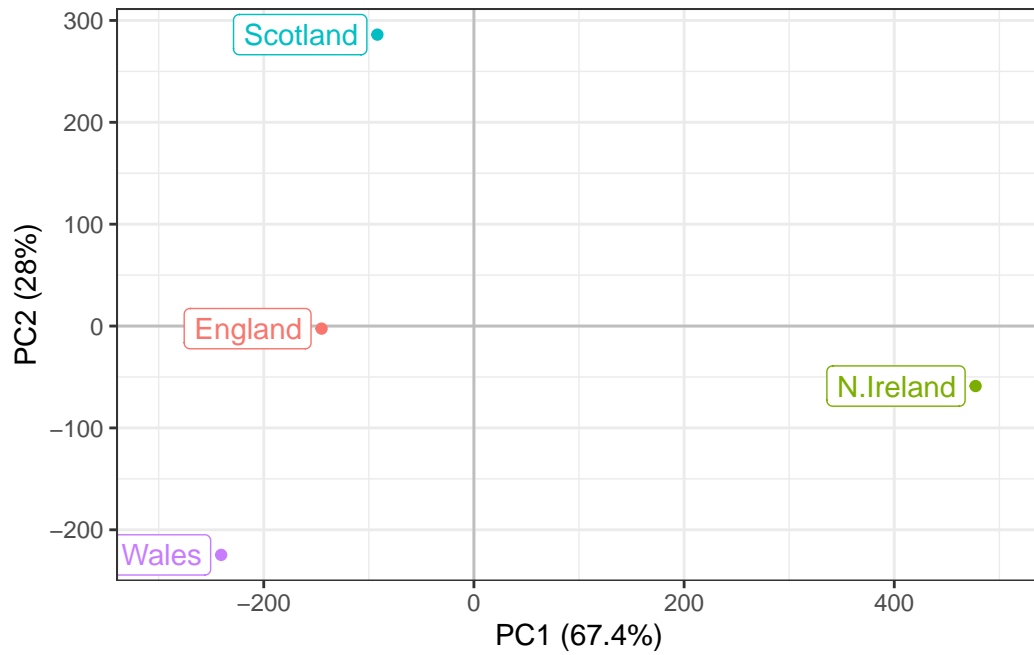
```
library(ggplot2)

df <- as.data.frame(pca$x)
df_lab <- tibble::rownames_to_column(df, "Country")

# Our first basic plot
ggplot(df_lab) +
  aes(PC1, PC2, col=Country) +
  geom_point()
```

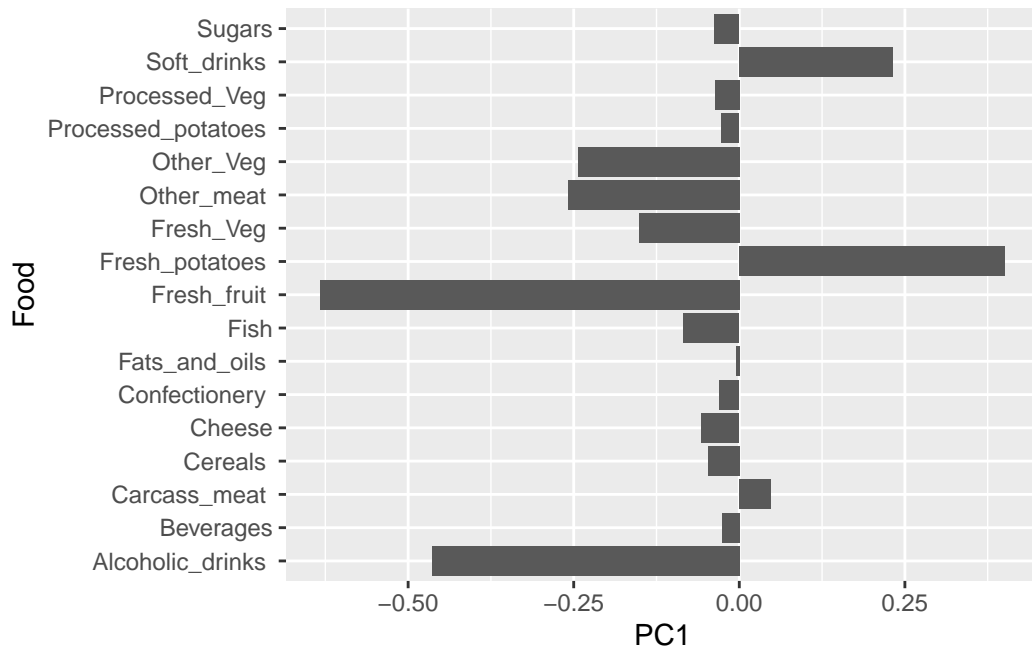



```
ggplot(df_lab) +  
  aes(PC1, PC2, col=Country, label=Country) +  
  geom_hline(yintercept = 0, col="gray") +  
  geom_vline(xintercept = 0, col="gray") +  
  geom_point(show.legend = FALSE) +  
  geom_label(hjust=1, nudge_x = -10, show.legend = FALSE) +  
  expand_limits(x = c(-300,500)) +  
  xlab("PC1 (67.4%)") +  
  ylab("PC2 (28%)") +  
  theme_bw()
```

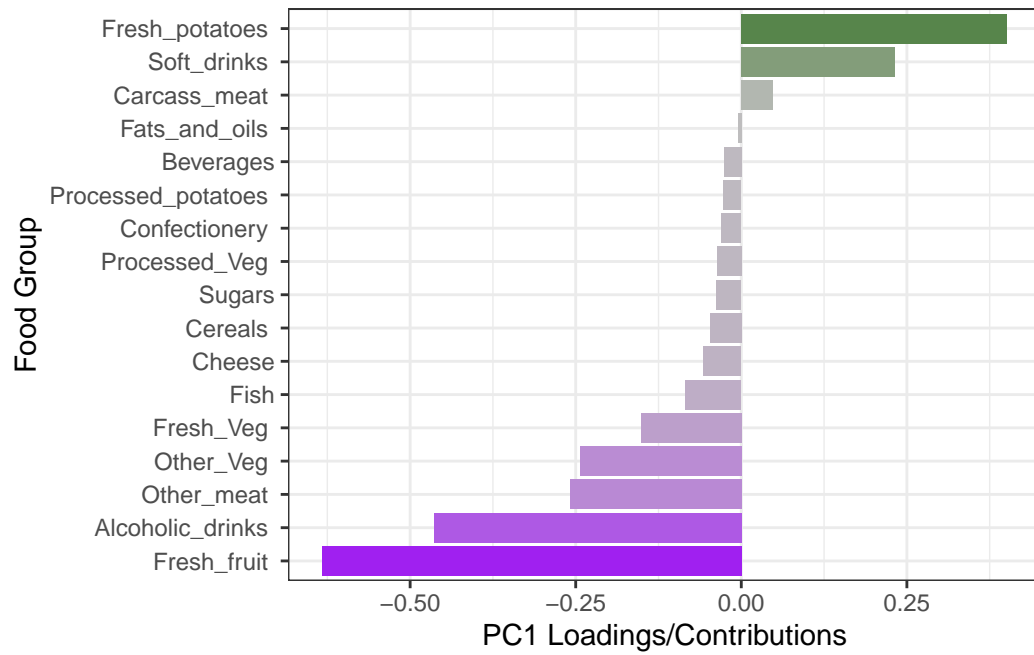


```
ld <- as.data.frame(pca$rotation)
ld_lab <- tibble::rownames_to_column(ld, "Food")

ggplot(ld_lab) +
  aes(PC1, Food) +
  geom_col()
```

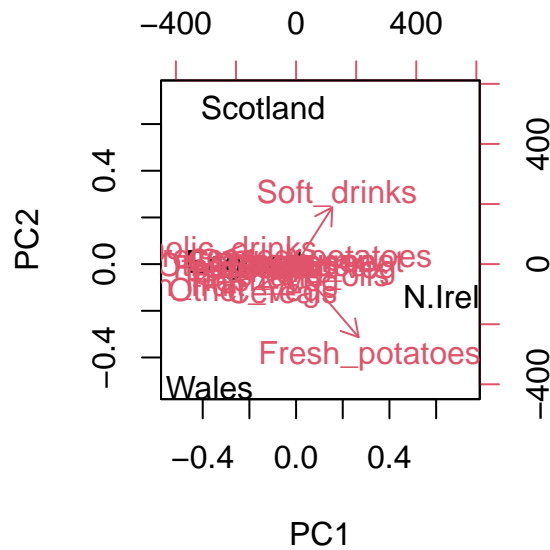


```
ggplot(ld_lab) +
  aes(PC1, reorder(Food, PC1), bg=PC1) +
  geom_col() +
  xlab("PC1 Loadings/Contributions") +
  ylab("Food Group") +
  scale_fill_gradient2(low="purple", mid="gray", high="darkgreen", guide=NULL) +
  theme_bw()
```



Biplots

```
biplot(pca)
```



2. PCA of RNA-seq data

```
url2 <- "https://tinyurl.com/expression-CSV"
rna.data <- read.csv(url2, row.names=1)
head(rna.data)
```

	wt1	wt2	wt3	wt4	wt5	ko1	ko2	ko3	ko4	ko5
gene1	439	458	408	429	420	90	88	86	90	93
gene2	219	200	204	210	187	427	423	434	433	426
gene3	1006	989	1030	1017	973	252	237	238	226	210
gene4	783	792	829	856	760	849	856	835	885	894
gene5	181	249	204	244	225	277	305	272	270	279
gene6	460	502	491	491	493	612	594	577	618	638

Q10: How many genes and samples are in this data set?

```
dim(rna.data)
```

```
[1] 100 10
```

There are 100 genes and 10 samples.