# Class 7: Machine Learning 1

Benjie Miao (PID: A69026849)

## Clustering
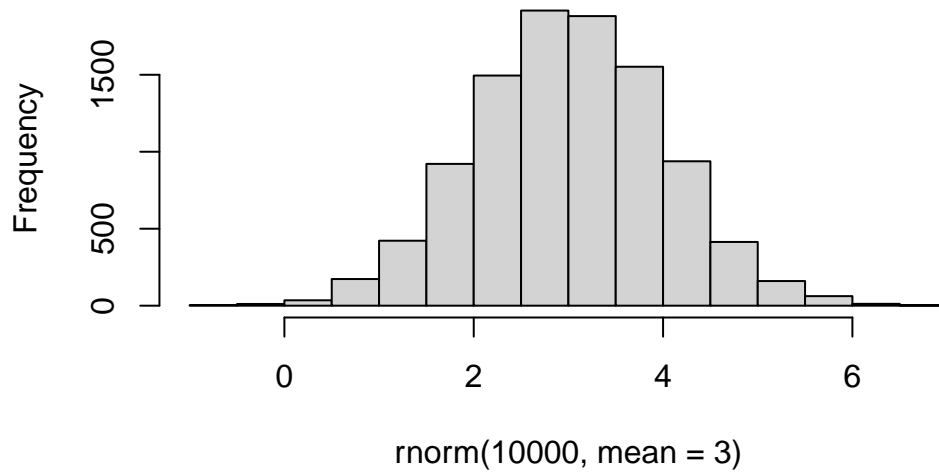
We will start with k-means clustering, one of the most prevelent of all clustering methods.

### K-means clustering

K-means clustering is the commonest algorithm for doing clustering.
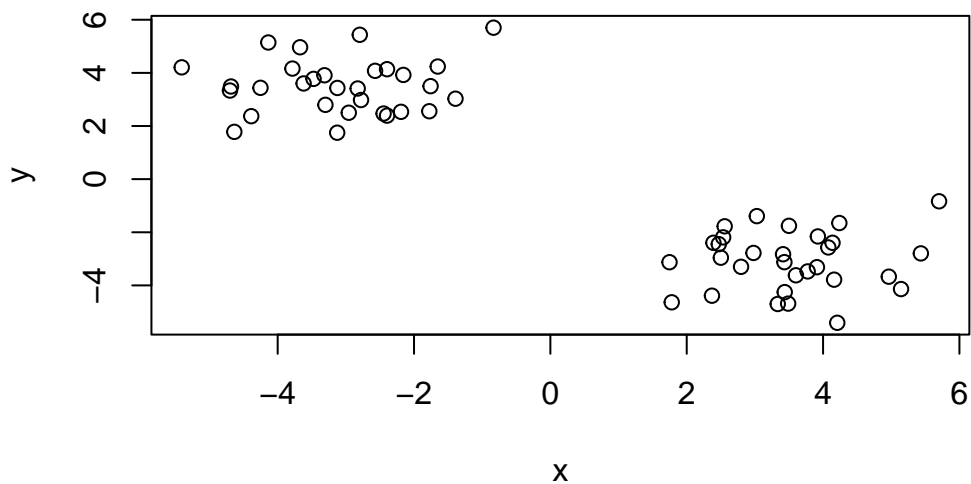
```
hist(rnorm(10000, mean = 3))
```

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

```
tmp <- c(rnorm(30, 3), rnorm (30, -3)) # we generate two bunches of points
a <- cbind(x = tmp, y = rev(tmp)) # for coordinate. two buckles of points around (3, -3) a
head(a)
```

```
          x         y
[1,] 3.437485 -4.252542
[2,] 4.077042 -2.571675
[3,] 3.412754 -2.826503
[4,] 2.369859 -4.387497
[5,] 3.499400 -1.757220
[6,] 3.334941 -4.699861
```
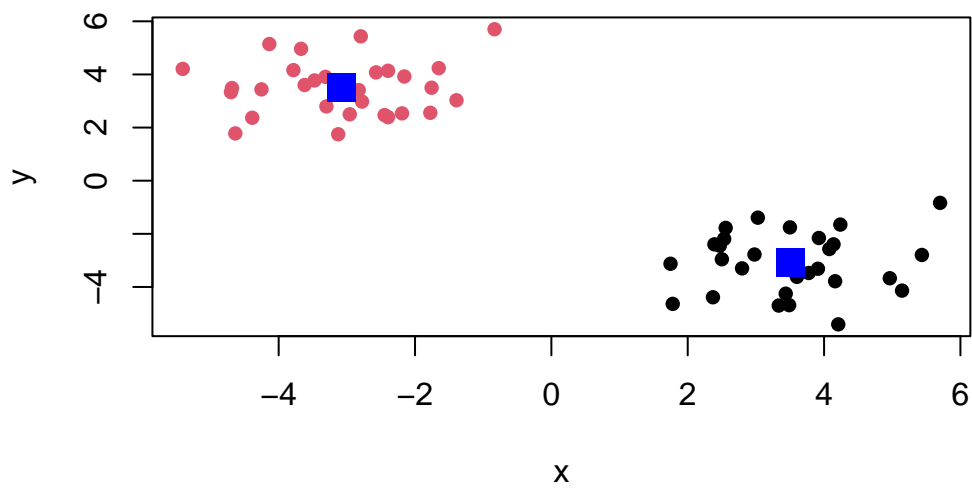
```
plot(a)
```



The primitive of calling k-means is just `kmeans()`

```
k <- kmeans(a, centers = 2, nstart = 20)
k
```
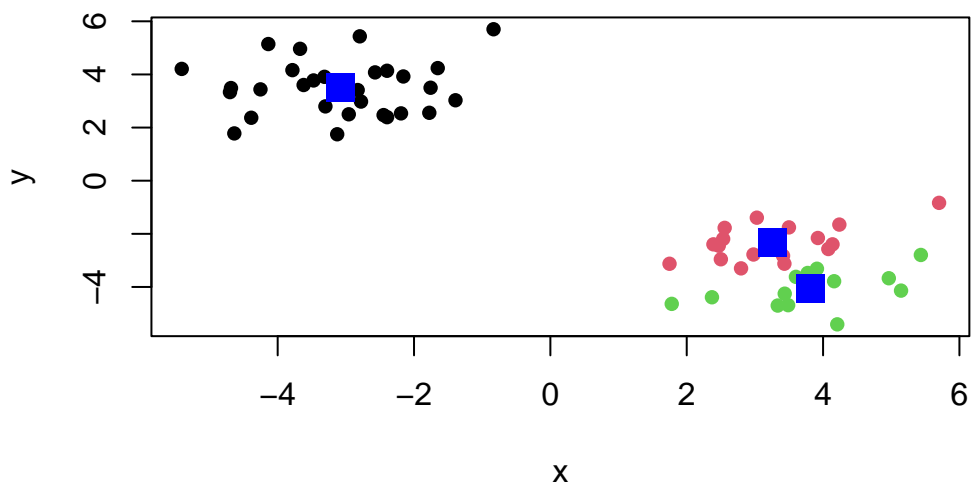
```
K-means clustering with 2 clusters of sizes 30, 30
```

2

```
Cluster means:
          x          y
1  3.501267 -3.084773
2 -3.084773  3.501267


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] 65.05629 65.05629
 (between_SS / total_SS =  90.9 %)


Available components:

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

Let's check some of the properties:

```
k$size
```

```
[1] 30 30
```

```
k$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
```

```
k$centers
```

```
          x          y
1  3.501267 -3.084773
2 -3.084773  3.501267
```

```
plot(a, col=k$cluster, pch=16)
points(k$centers, col="blue", pch=15, cex=2)
```

What if we chance the number of centers?

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

## Hierarchical clustering

HC can reveal multi-level structural information rather than just imposing the structure (k-means).

Funciton in R is `hclust()`, which requires a distance matrix as the input.
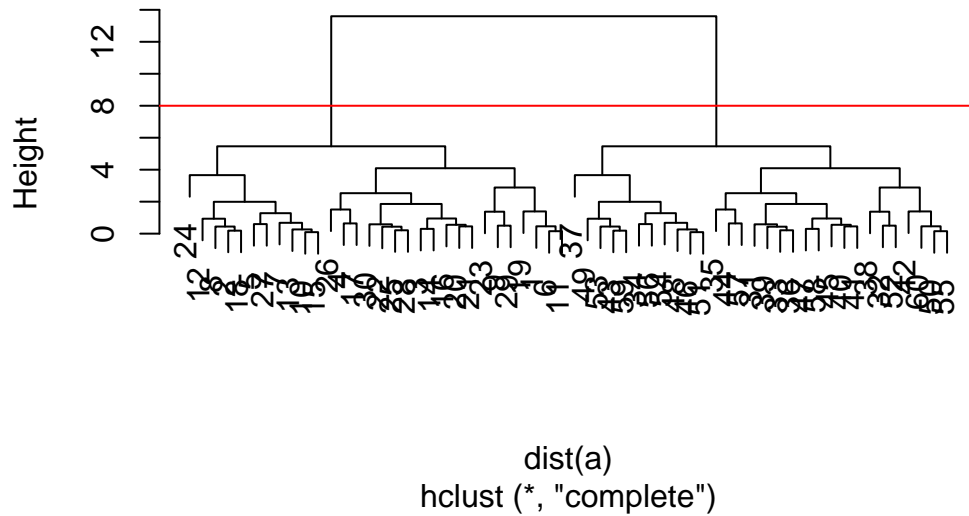
```r
hc <- hclust(dist(a))
hc
```

```
Call:
hclust(d = dist(a))

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

```r
plot(hc)
abline(h=8, col='red')
```
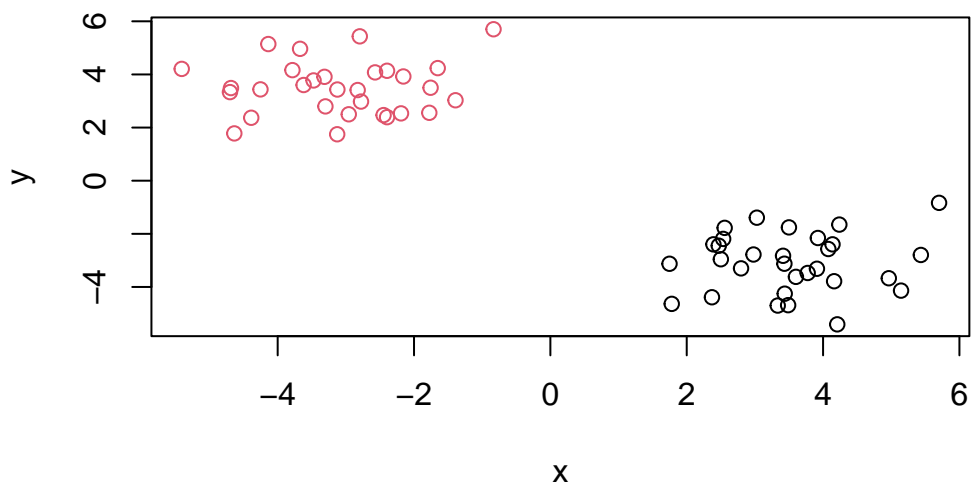
## Cluster Dendrogram



dist(a)
hclust (*, "complete")

We can 'cut the tree' using `cutree`:

```r
groups <- cutree(hc, h = 8)
groups
```

```
 [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 1 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
```

```r
plot(a, col=groups)
```

Seems cool!

## PCA

Now we use the UK food dataset for the PCA exploration.

```
#load the dataset
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
```
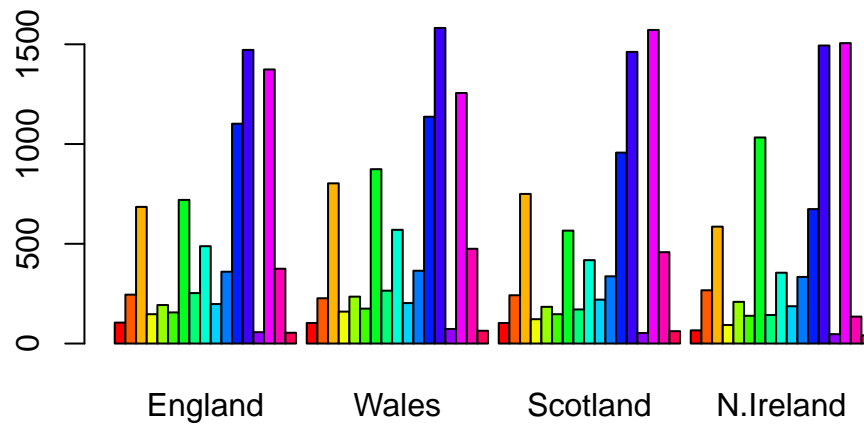
```
dim(x)
```

```
[1] 17   5
```

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

17 and 4. `dim`

Let's set the index to be the first column:

```
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
```

Alternatively, we can use the parameter in `read.csv`

```
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 prefer the latter method, since it preserve when running multiple times.
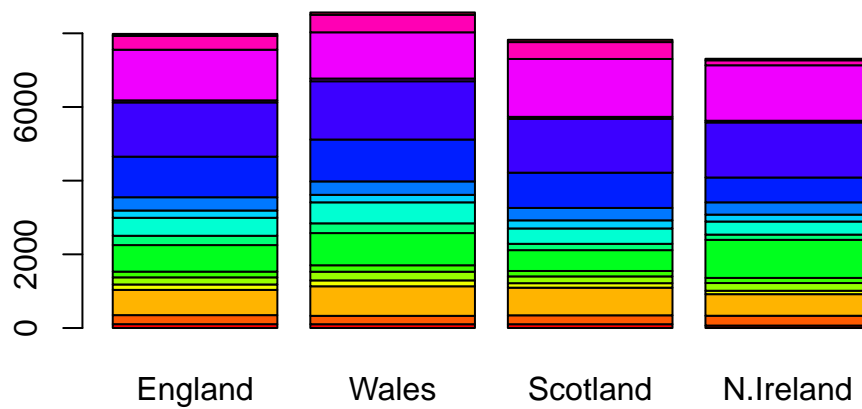
**Spotting major differences and trends**

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



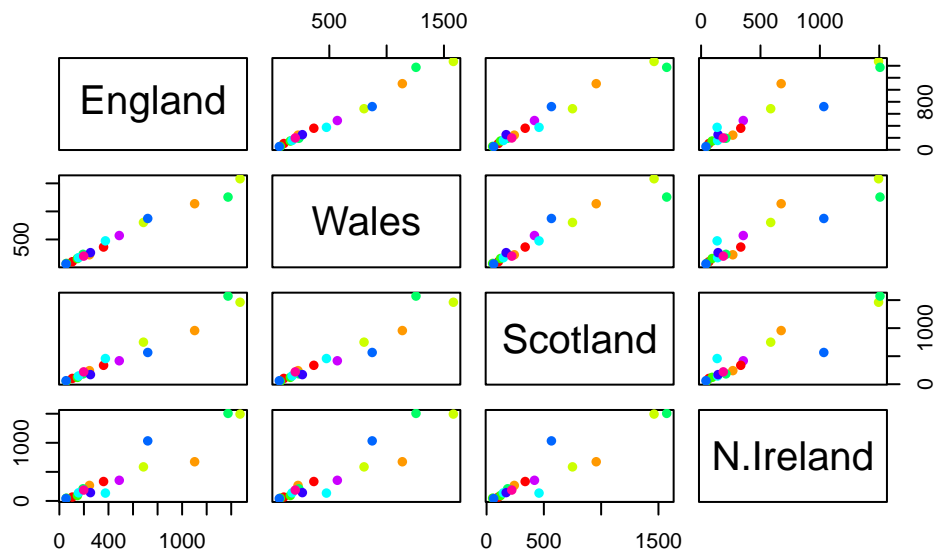**Q3:** **Changing what optional argument in the above barplot() function results in the following plot?**

We can use `beside=FALSE`

```r
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)
```

Points on the diagonal means two countries have the same value on a certain category.

**Q6. What is the main differences between N. Ireland and the other countries of the UK in terms of this data-set?**

They are further away from the diagonal.

**PCA to the rescue.**

```
# 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
```

**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.**

```
color <- c("orange", "blue", "pink", "red")

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

We can further analyze the standard deviation

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

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

```
## or the second row here...
z <- summary(pca)
z$importance
```
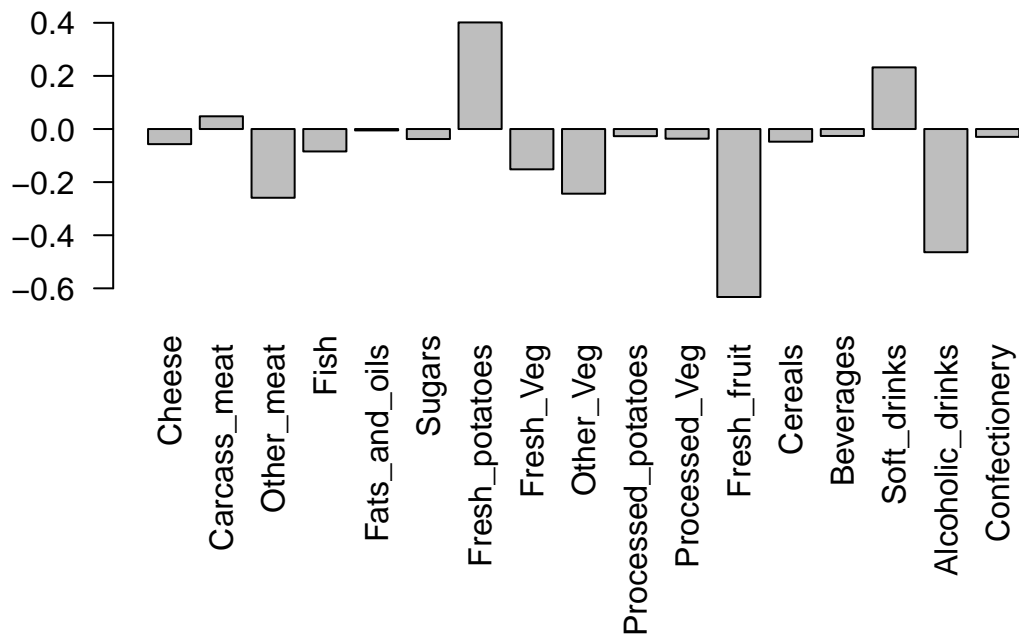
|                        | PC1      | PC2      | PC3      | PC4         |
|------------------------|----------|----------|----------|-------------|
| Standard deviation     | 324.15019| 212.74780| 73.87622 | 3.175833e-14|
| Proportion of Variance | 0.67444  | 0.29052  | 0.03503  | 0.000000e+00|
| Cumulative Proportion  | 0.67444  | 0.96497  | 1.00000  | 1.000000e+00|

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

## Digging deeper (variable loadings)

```r
## Lets focus on PC1 as it accounts for > 90% of variance
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 prominantely and what does PC2 maninly tell us about?**

```
## Lets focus on PC1 as it accounts for > 90% of variance
par(mar=c(10, 3, 0.35, 0))
barplot( pca$rotation[,2], las=2 )
```

The two dominant factors are `Fresh_potatoes` and `Soft_drinks`. From the previous plot, we can see that Wales and Scotland are at two ends of the PC2, therefore, that means those two factors can explain their difference.

## 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()
```
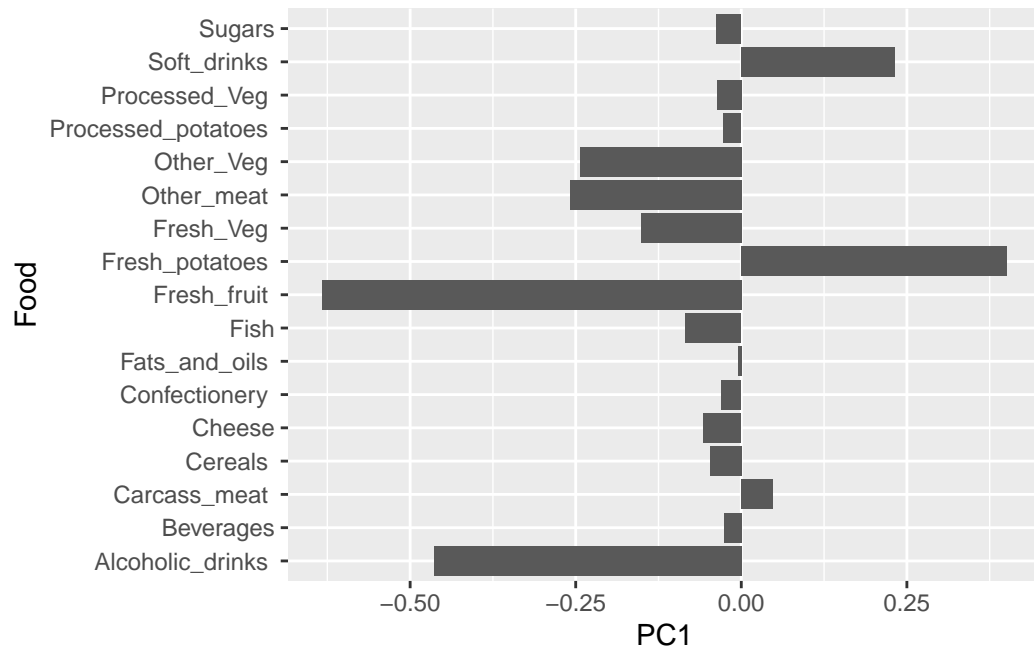
```r
ggplot(df_lab) +
  aes(PC1, PC2, col=Country, label=Country) +
  geom_hline(yintercept = 0, col="gray") +
  geom_vline(xintercept = 0, col="gray") +
  geom_point(show.legend = FALSE) +
  geom_label(hjust=1, nudge_x = -10, show.legend = FALSE) +
  expand_limits(x = c(-300,500)) +
  xlab("PC1 (67.4%)") +
  ylab("PC2 (28%)") +
  theme_bw()
```
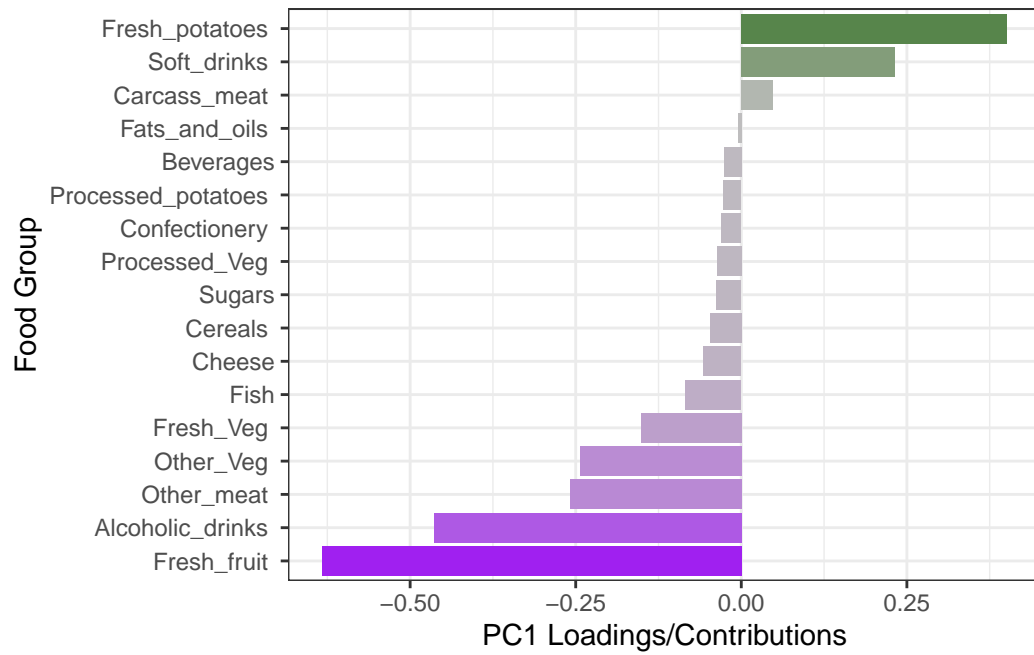
```
ld <- as.data.frame(pca$rotation)
ld_lab <- tibble::rownames_to_column(ld, "Food")

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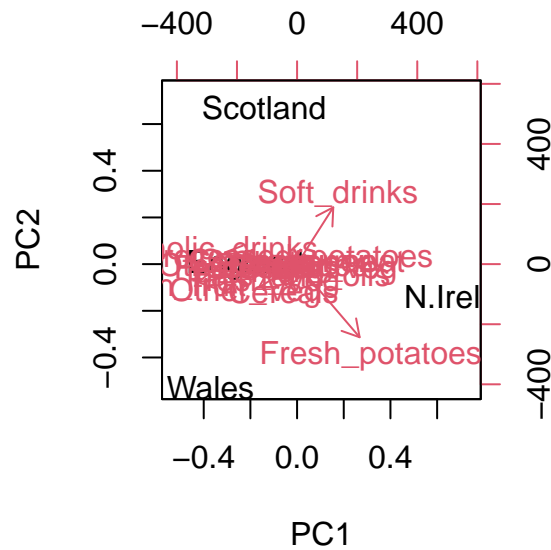
```
ggplot(ld_lab) +
  aes(PC1, reorder(Food, PC1), bg=PC1) +
  geom_col() +
  xlab("PC1 Loadings/Contributions") +
  ylab("Food Group") +
  scale_fill_gradient2(low="purple", mid="gray", high="darkgreen", guide=NULL) +
  theme_bw()
```

## Biplot

```r
## The inbuilt biplot() can be useful for small datasets
biplot(pca)
```

## 2. PCA of RNA-seq data

Again load the dataset:

```
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
```

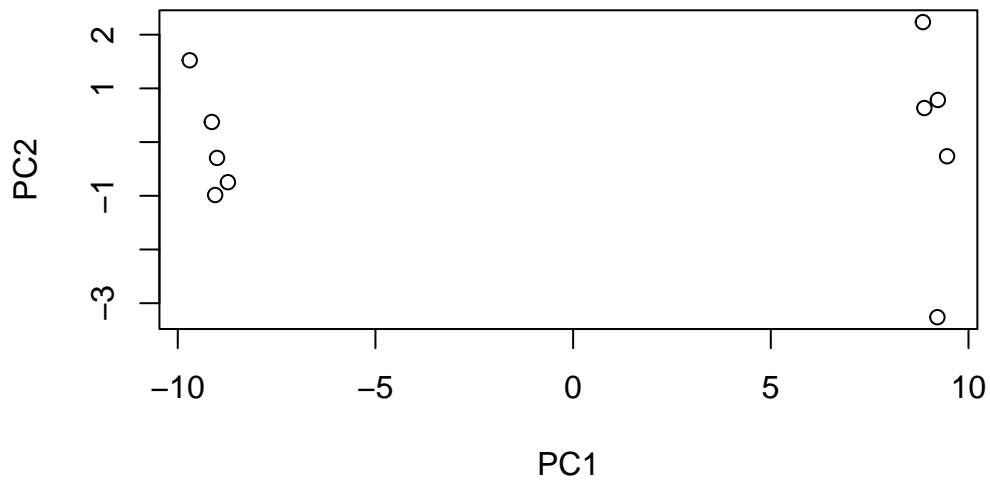**Q10: How many genes and samples are in this data set?**

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

```
[1] 100   10
```

10 samples, 100 genes.

```
## Again we have to take the transpose of our data
pca <- prcomp(t(rna.data), scale=TRUE)

## Simple un polished plot of pc1 and pc2
plot(pca$x[,1], pca$x[,2], xlab="PC1", ylab="PC2")
```



```
summary(pca)
```
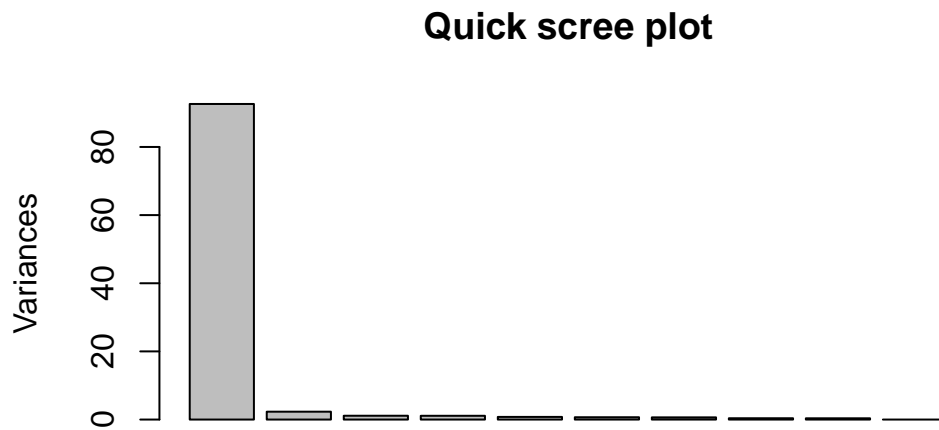
```
Importance of components:
                          PC1    PC2     PC3     PC4     PC5     PC6     PC7
Standard deviation     9.6237 1.5198 1.05787 1.05203 0.88062 0.82545 0.80111
Proportion of Variance 0.9262 0.0231 0.01119 0.01107 0.00775 0.00681 0.00642
Cumulative Proportion  0.9262 0.9493 0.96045 0.97152 0.97928 0.98609 0.99251
                          PC8     PC9       PC10
Standard deviation     0.62065 0.60342 3.457e-15
Proportion of Variance 0.00385 0.00364 0.000e+00
Cumulative Proportion  0.99636 1.00000 1.000e+00
```

```
plot(pca, main="Quick scree plot")
```
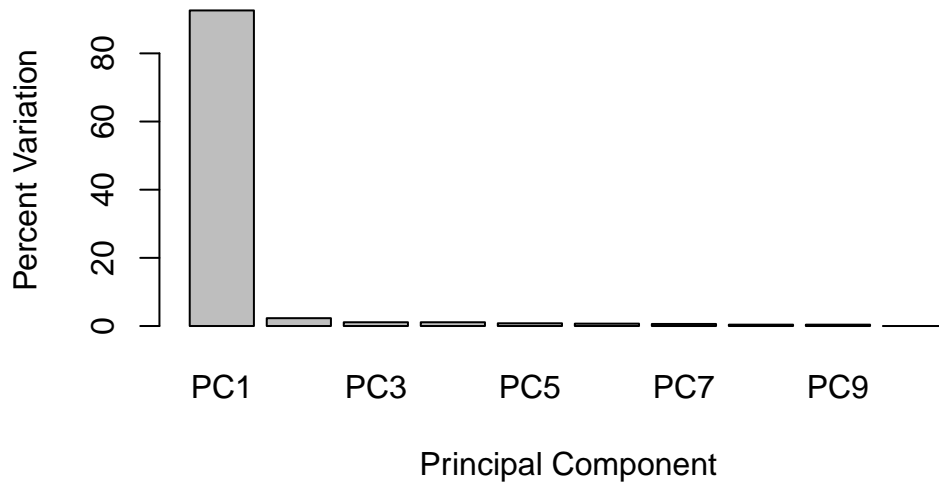
## Quick scree plot



```
## Variance captured per PC
pca.var <- pca$sdev^2

## Percent variance is often more informative to look at
pca.var.per <- round(pca.var/sum(pca.var)*100, 1)
pca.var.per
```

```
[1] 92.6  2.3  1.1  1.1  0.8  0.7  0.6  0.4  0.4  0.0
```

```
barplot(pca.var.per, main="Scree Plot",
        names.arg = paste0("PC", 1:10),
        xlab="Principal Component", ylab="Percent Variation")
```
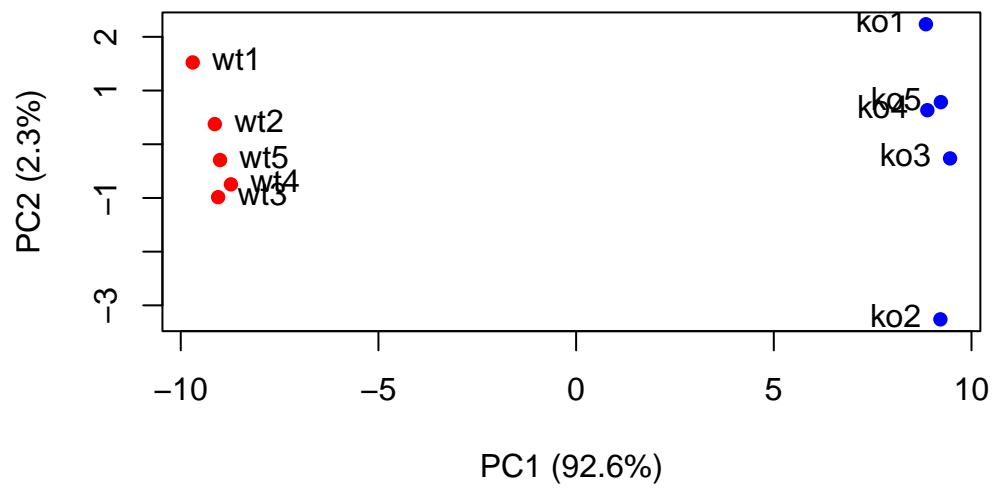
## Scree Plot



```r
## A vector of colors for wt and ko samples
colvec <- colnames(rna.data)
colvec[grep("wt", colvec)] <- "red"
colvec[grep("ko", colvec)] <- "blue"

plot(pca$x[,1], pca$x[,2], col=colvec, pch=16,
     xlab=paste0("PC1 (", pca.var.per[1], "%)"),
     ylab=paste0("PC2 (", pca.var.per[2], "%)"))

text(pca$x[,1], pca$x[,2], labels = colnames(rna.data), pos=c(rep(4,5), rep(2,5)))
```
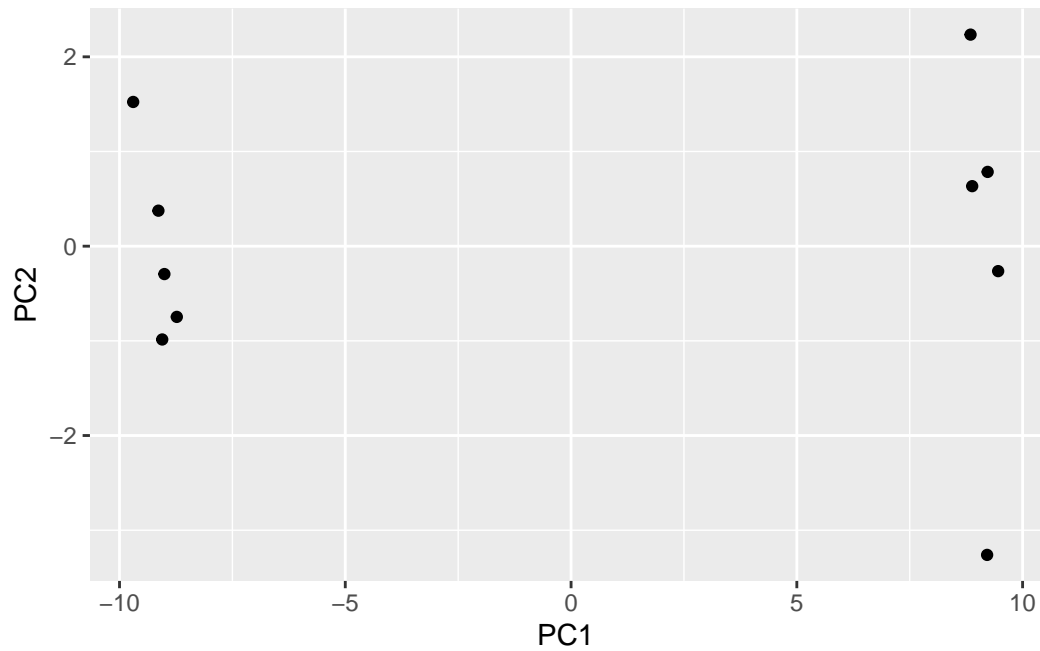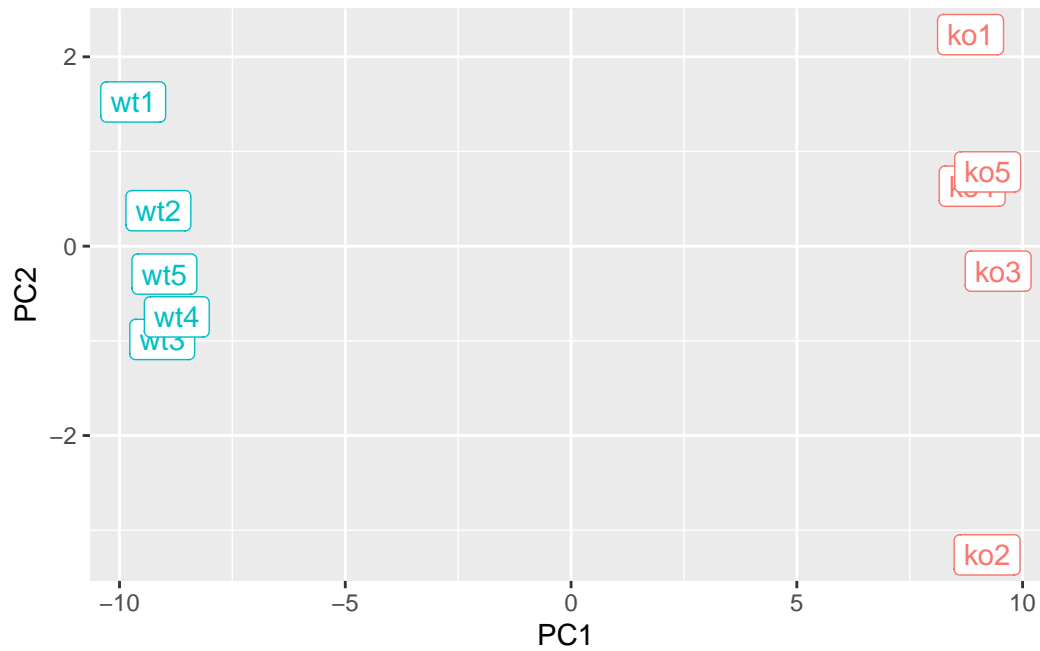
```
library(ggplot2)

df <- as.data.frame(pca$x)

# Our first basic plot
ggplot(df) +
  aes(PC1, PC2) +
  geom_point()
```

```
# Add a 'wt' and 'ko' "condition" column
df$samples <- colnames(rna.data)
df$condition <- substr(colnames(rna.data),1,2)

p <- ggplot(df) +
        aes(PC1, PC2, label=samples, col=condition) +
        geom_label(show.legend = FALSE)
p
```
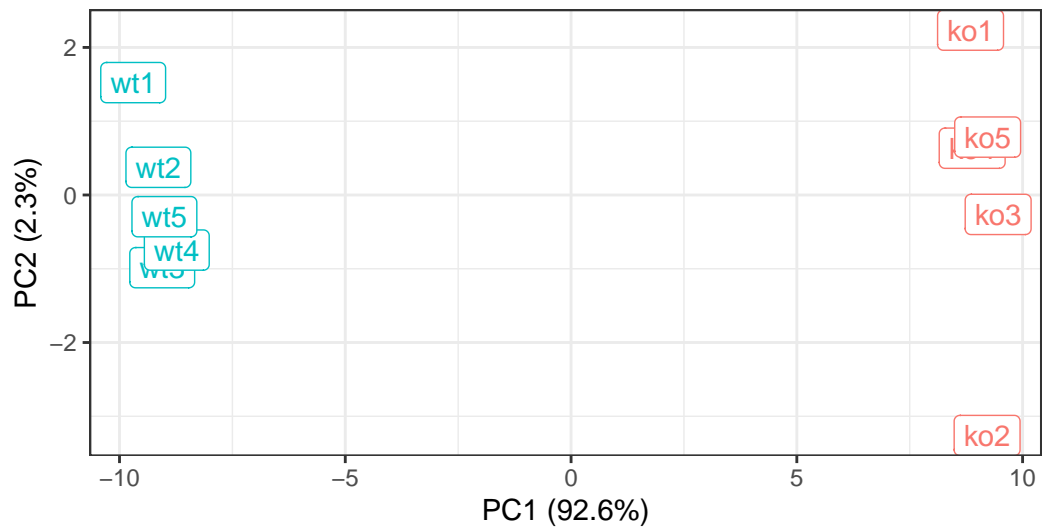
```
p + labs(title="PCA of RNASeq Data",
        subtitle = "PC1 clealy seperates wild-type from knock-out samples",
        x=paste0("PC1 (", pca.var.per[1], "%)"),
        y=paste0("PC2 (", pca.var.per[2], "%)"),
        caption="Class example data") +
    theme_bw()
```

## PCA of RNASeq Data

PC1 clealy seperates wild–type from knock–out samples



Class example data

## (Optional) Gene loading:

```
loading_scores <- pca$rotation[,1]

## Find the top 10 measurements (genes) that contribute
## most to PC1 in either direction (+ or -)
gene_scores <- abs(loading_scores)
gene_score_ranked <- sort(gene_scores, decreasing=TRUE)

## show the names of the top 10 genes
top_10_genes <- names(gene_score_ranked[1:10])
top_10_genes
```

```
[1] "gene100" "gene66"  "gene45"  "gene68"  "gene98"  "gene60"  "gene21"
[8] "gene56"  "gene10"  "gene90"
```