# Class07Lab

## Joshua Mac

## 2025-02-03

### *PCA of UK food data*

First, read the .csv to open our data on food in the UK.

```
url <- "https://tinyurl.com/UK-foods"
x <- read.csv(url)
```

> 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)
```

```
[1] 17  5
```

> A1. So there are 17 rows and 5 columns.

To preview the first 6 of the 17 rows, we'll use:

```
## View all we can use View(x), which is case-sens btw
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
```

Our first column was taken as the row headers, which isn't ideal. Let's get rid of it. WE'll use `rownames()` to set that up.

```
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
```

Now, dim(x) should be:
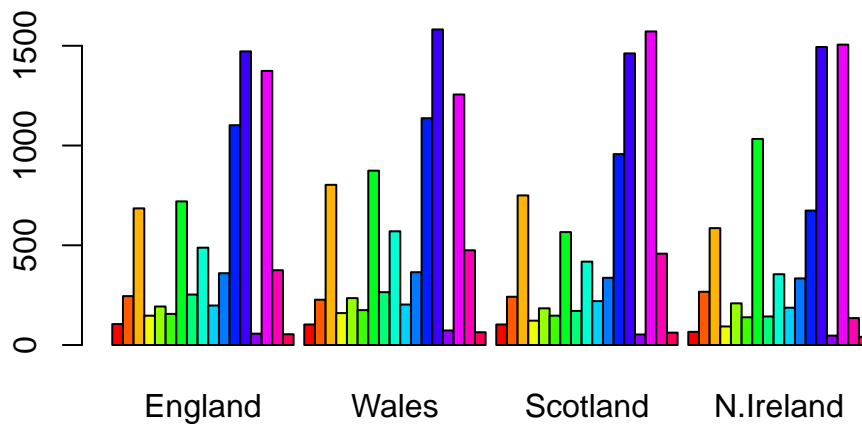
```
dim(x)
```

```
[1] 17   4
```

Alternatively, we could have added `row.names=1` to our initial read of the csv to do this.

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?

A2. I prefer using row.names=1 to make changes in the original read rather than having to write out additional code to change it. Morover, the former approach will continue to remove a column every time it is ran due to x<-x[,-1] will return the dataframe without the first column eatch time since we are not setting it to a different variable.
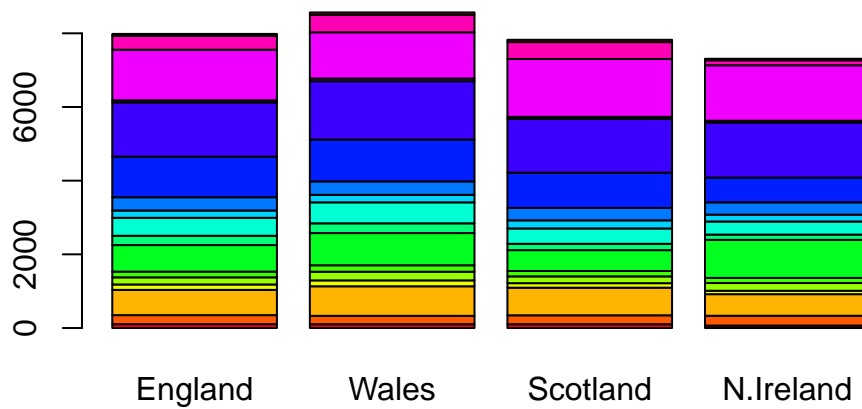
Numbers alone doesn't tell us much here. Neither do regular plots.

```
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?
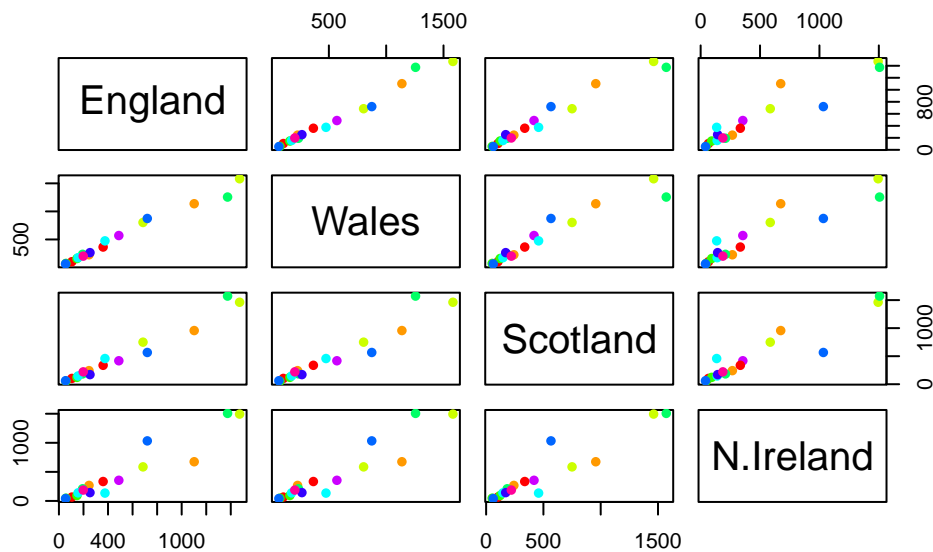
```
## Changing beside to = False will not separate each row side-by-side.
barplot(as.matrix(x), beside=F, col=rainbow(nrow(x)))
```

A3. Changing `beside=T` to `beside=F` will produce the graph above.

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?

```
## The code in question
pairs(x, col=rainbow(10), pch=16)
```

A5. It seems there are 4 figures that appear after each country (bar N. Ireland), but otherwise, it is quite confusing because what do those 4 figures represent. pairs(x) creates the scatterplot where each variable in x is plotted with one of the other variables for every variable. col=rainbow(10) sets the color to 10 different colors of the rainbow for visualizaiton. pch=16 uses solid circles for points. If a given point lies on the diagonal of a given plot, the correlation/relationship between the two variables of that given plot may be very strong.

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

A6. The main differences are that N. Ireland eats way more fresh_potatoes and consumes less alcoholic_drinks, cheese, fresh_fruit, and other_meat.

**PCA to visualize** PCA will help us better visualize. Using `prcomp()`, we will need to transpose using `t()` because its default assumes observations is rows and variables is columns.

```
# Use the prcomp() PCA function
pca <- prcomp( t(x) )
summary(pca)
```
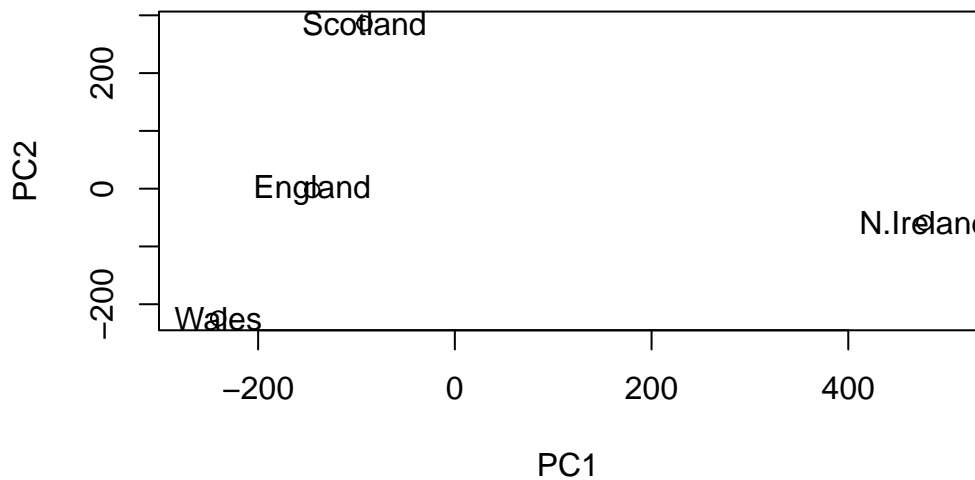
```
Importance of components:
                          PC1      PC2      PC3      PC4
```

```
Standard deviation      324.1502 212.7478 73.87622 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
```
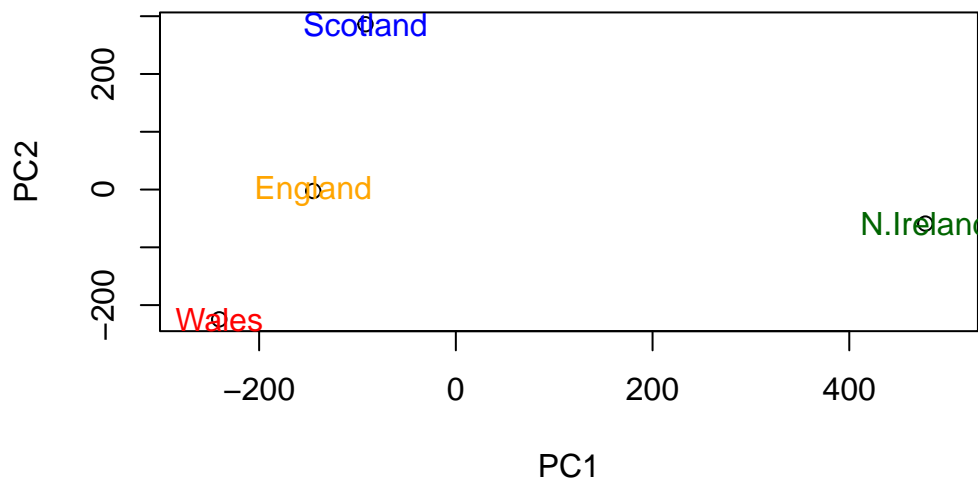
Q7. Complete the code below to generate a plot of PC1 vs PC2. The second line adds text labels over the data points.

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

Below, we will calculate how much variation in the data each PC acounts for using:

```
v <- round( pca$sdev^2/sum(pca$sdev^2) * 100 )
v
```
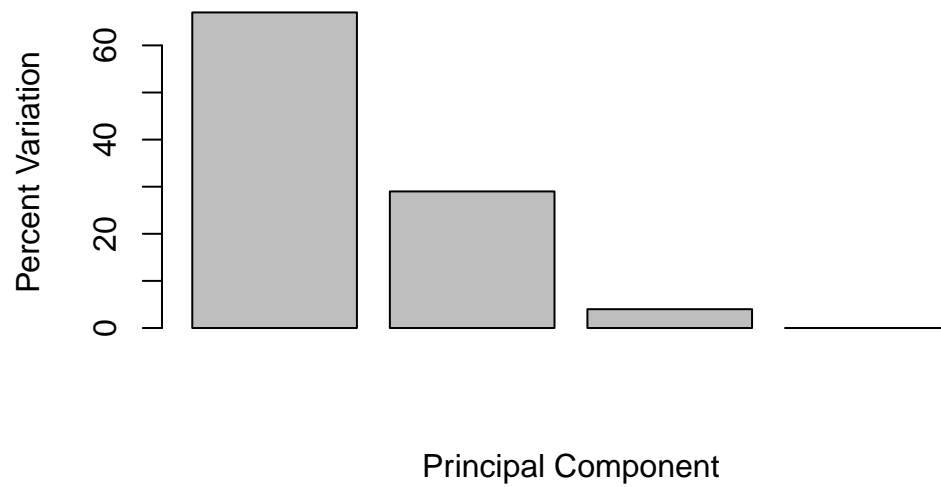
```
[1] 67 29  4  0
```

Which returned the proportion of variance from earlier.

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

|  | PC1 | PC2 | PC3 | PC4 |
|---|---|---|---|---|
| Standard deviation | 324.15019 | 212.74780 | 73.87622 | 2.921348e-14 |
| Proportion of Variance | 0.67444 | 0.29052 | 0.03503 | 0.000000e+00 |
| Cumulative Proportion | 0.67444 | 0.96497 | 1.00000 | 1.000000e+00 |

The information can be summarized in a plot of eigenvalues (variances).

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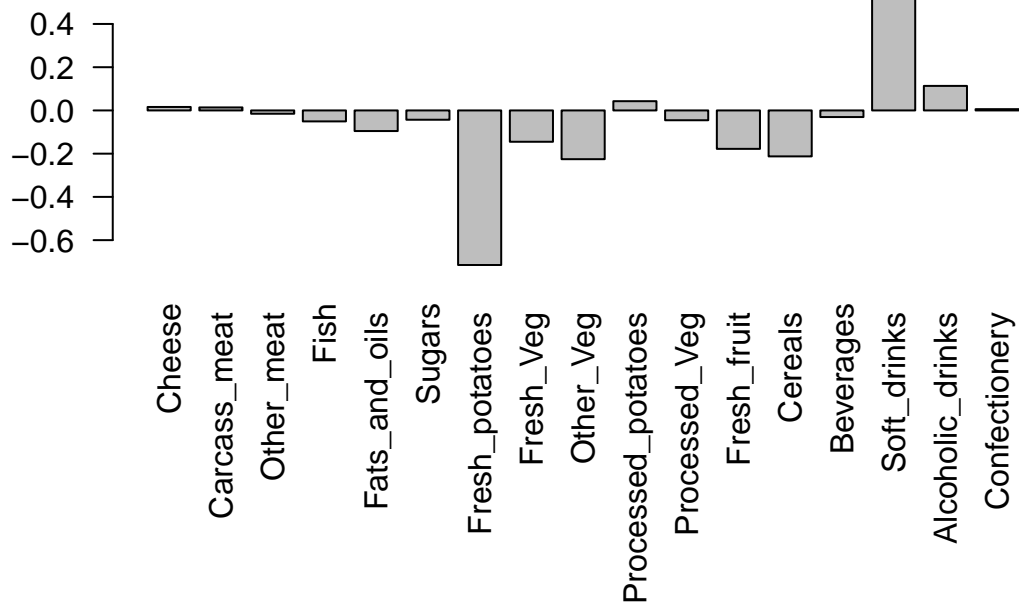
Principal Component

Moving on to variable loadings, which considers the influence of each of the origianl variables on the PC.

```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?

```
par(mar=c(10, 3, 0.35, 0))
barplot( pca$rotation[,2], las=2 )
```
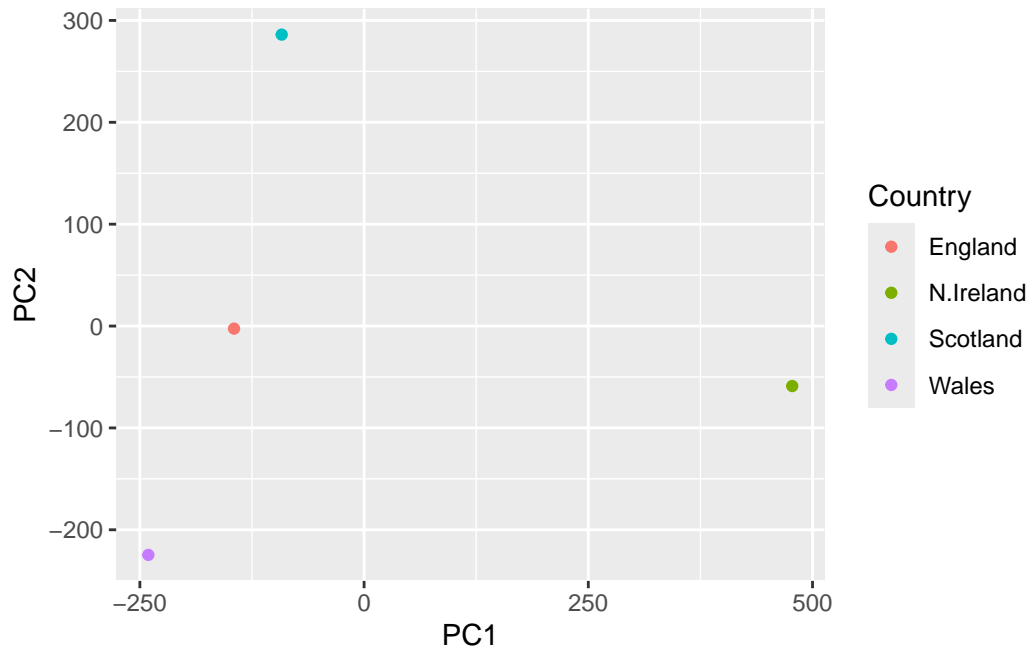
A9. The two groups that feature prominently are fresh_potatoes and soft_drinks negatively and positively pushing, respectively.

**Using ggplot for the figures!**
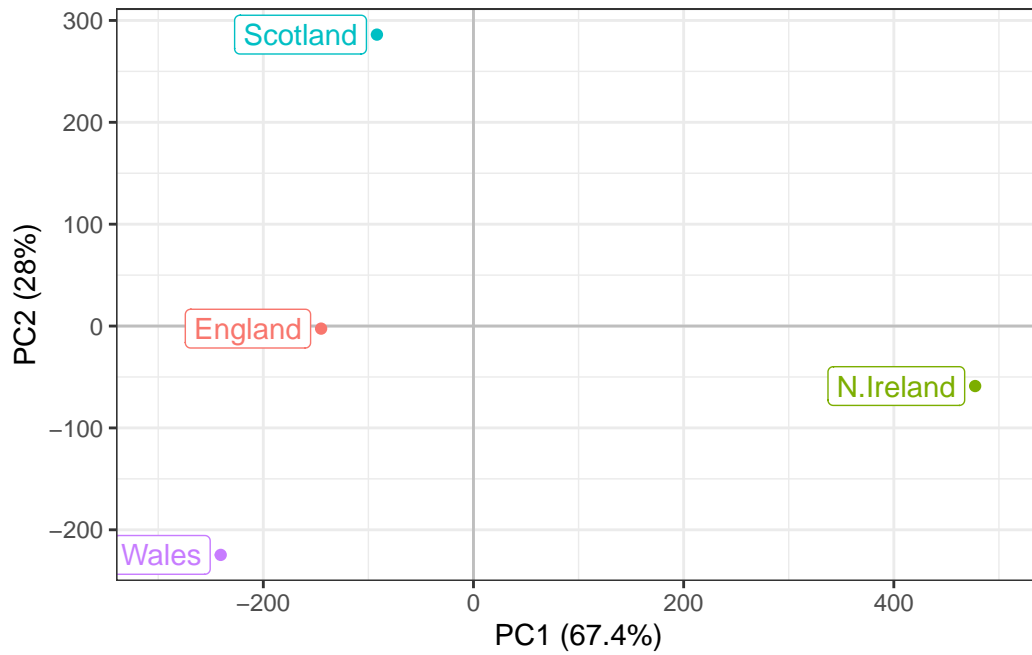
```r
library(ggplot2)

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()
```
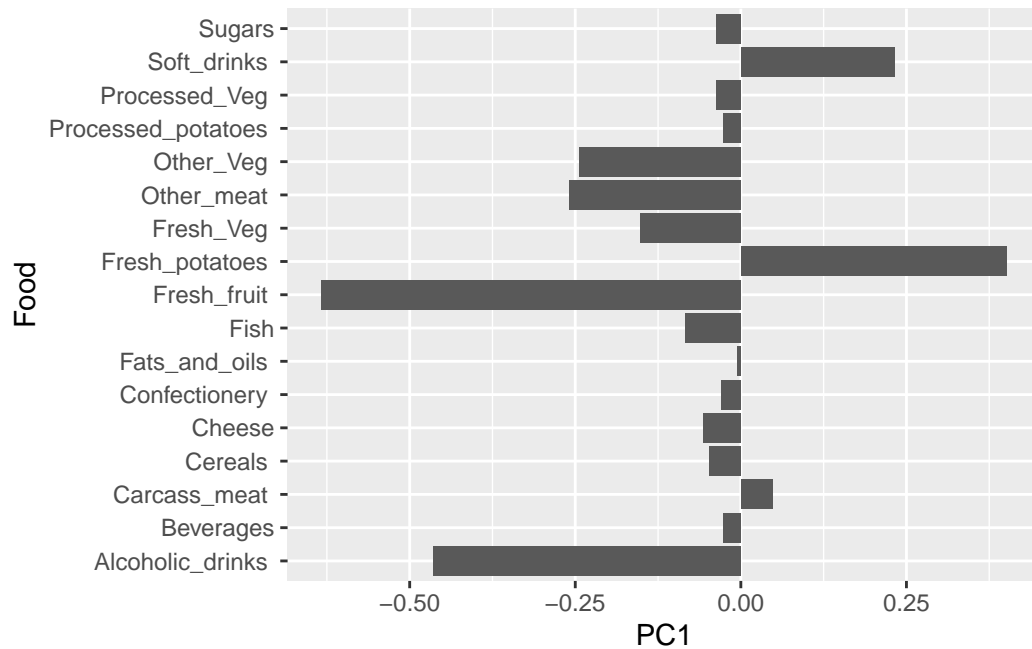
Adding more:

```
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()
```

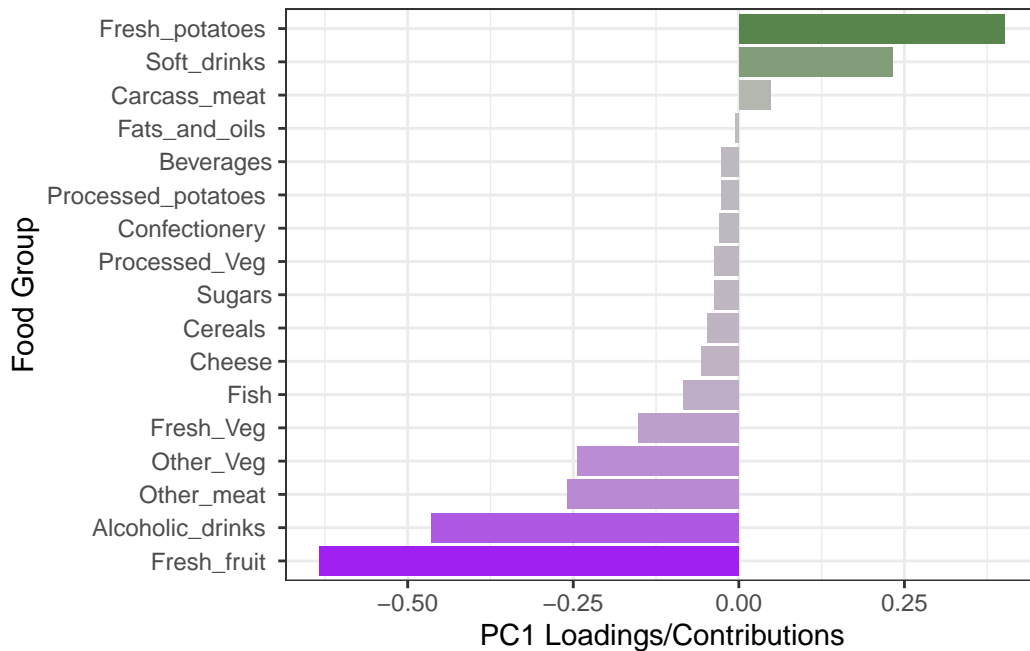Let's do the same with the loadings figure!

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

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

Adding more:

```
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()
```

13

## PCA of RNA-seq data

```r
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?
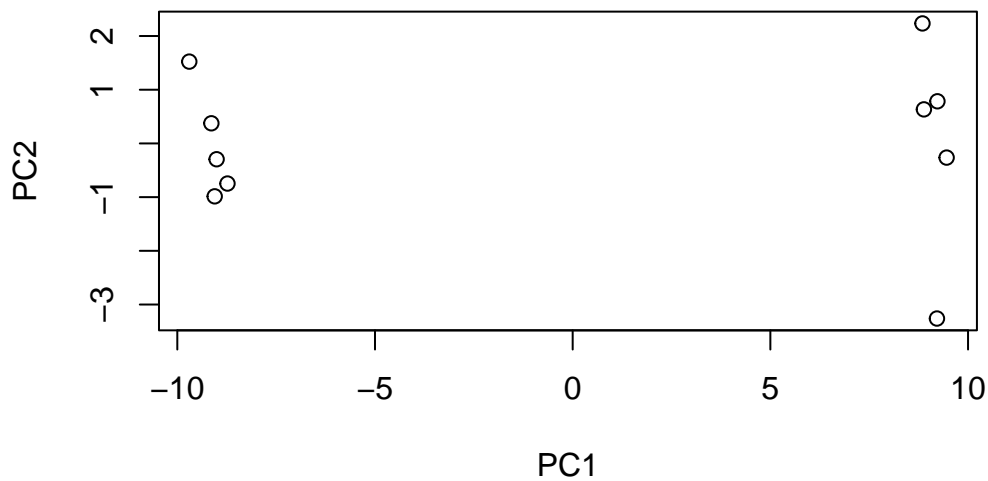
```r
dim(rna.data)
```

```
[1] 100  10
```

A10. If genes are rows and samples are columns, there are 100 genes and 10 samples.

To make the plot:

```
## 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")
```



How much variation?
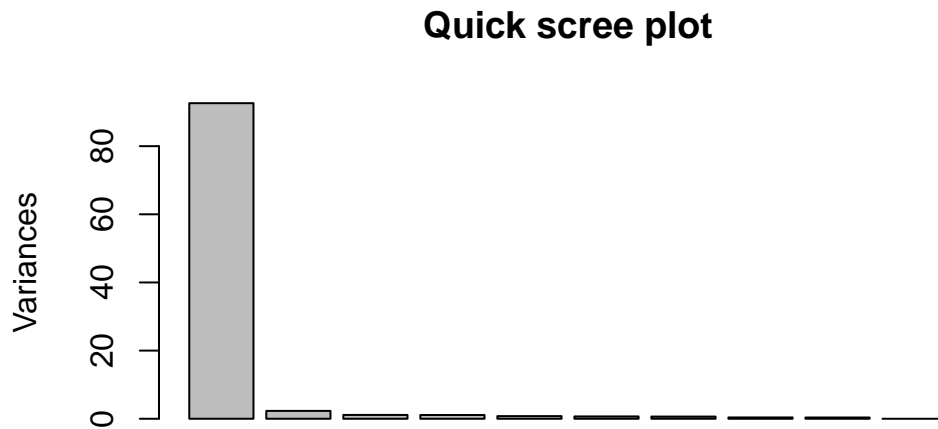
```
summary(pca)
```

```
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 3.345e-15
Proportion of Variance 0.00385 0.00364 0.000e+00
Cumulative Proportion  0.99636 1.00000 1.000e+00
```

PC1 is the dimension of interest (92.6%), accounting for the most variation.

To quickly plot:

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

## Quick scree plot



Variation!

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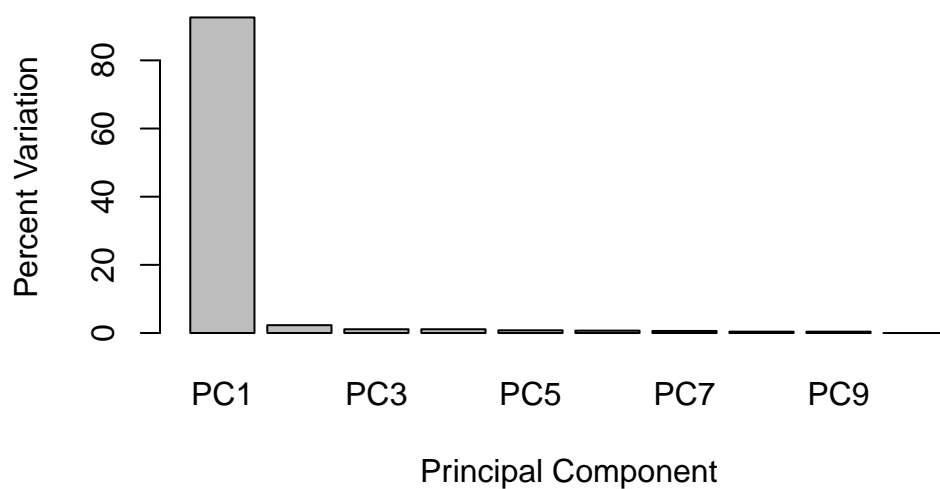
```
 [1] 92.6  2.3  1.1  1.1  0.8  0.7  0.6  0.4  0.4  0.0
```

Again, returning what we saw in the summary.

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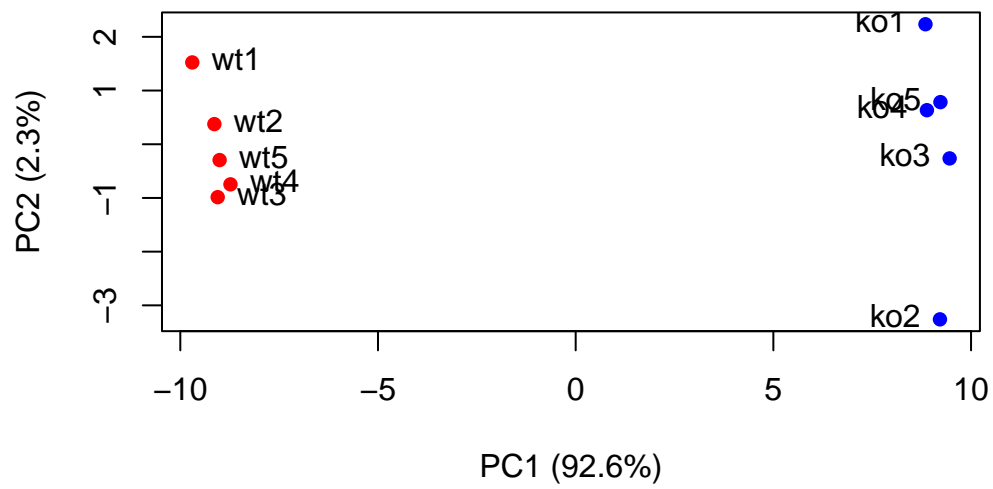
## Scree Plot



Spice it up a bit:

```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)))
```
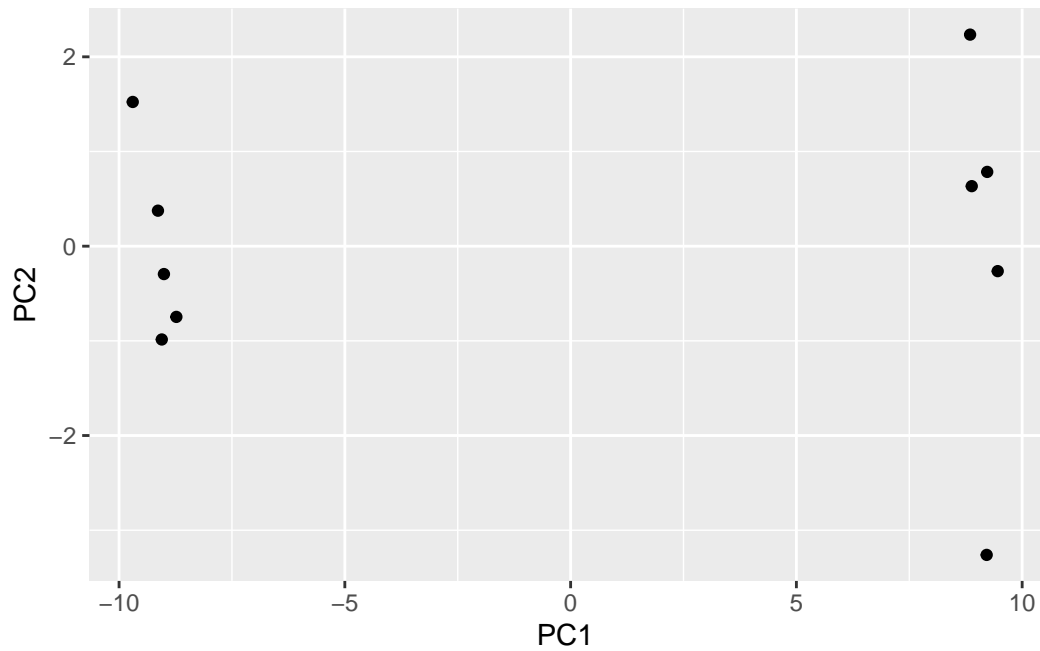
Using ggplot:

```
library(ggplot2)

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

# Our first basic plot
ggplot(df) +
  aes(PC1, PC2) +
  geom_point()
```
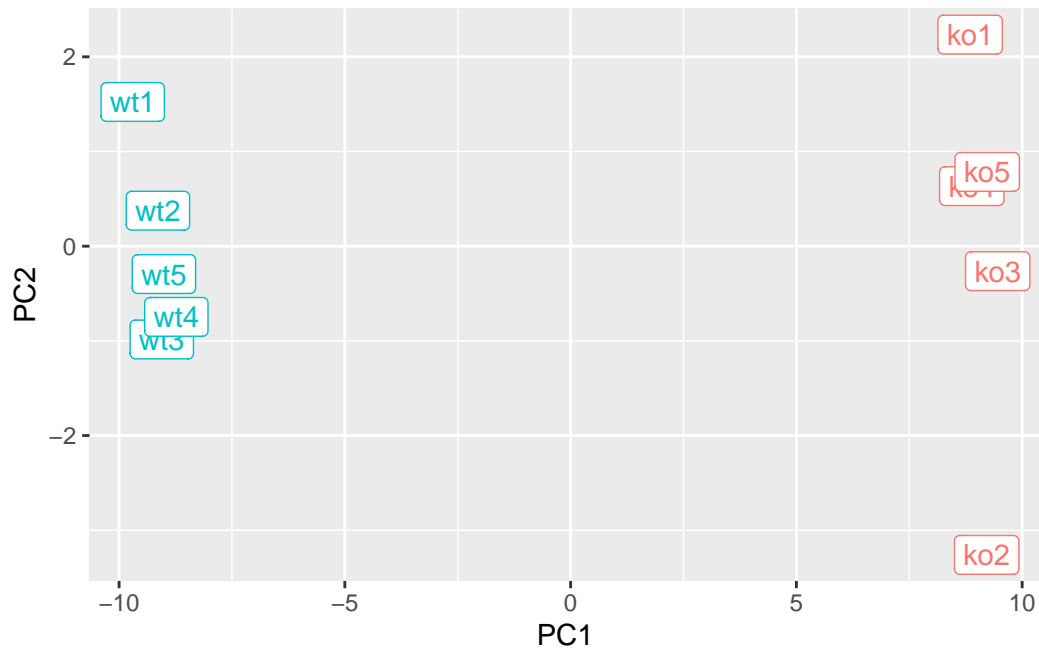
To add the information from our initial dataframe

```r
# 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
```
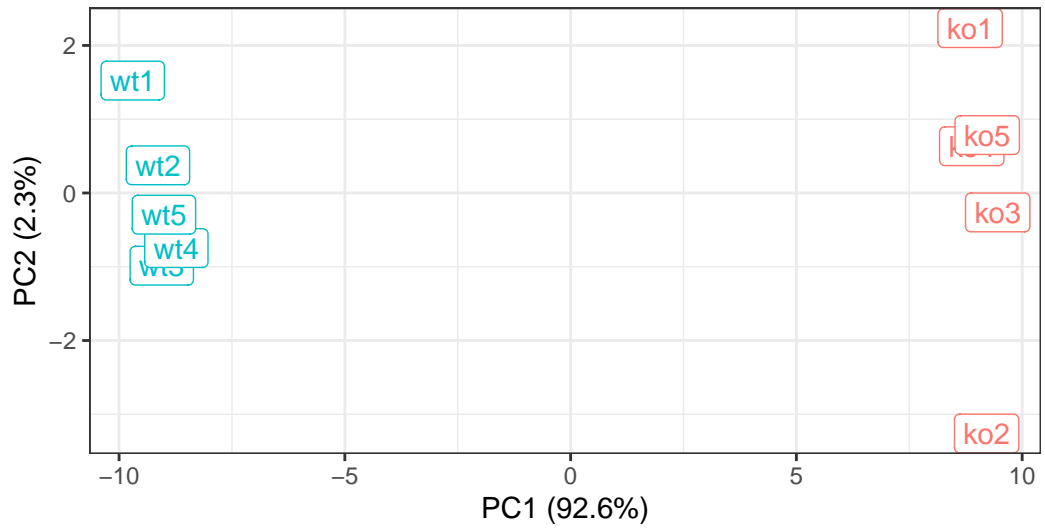
Finally

```r
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

What a beaut.