

Machine Learning 1

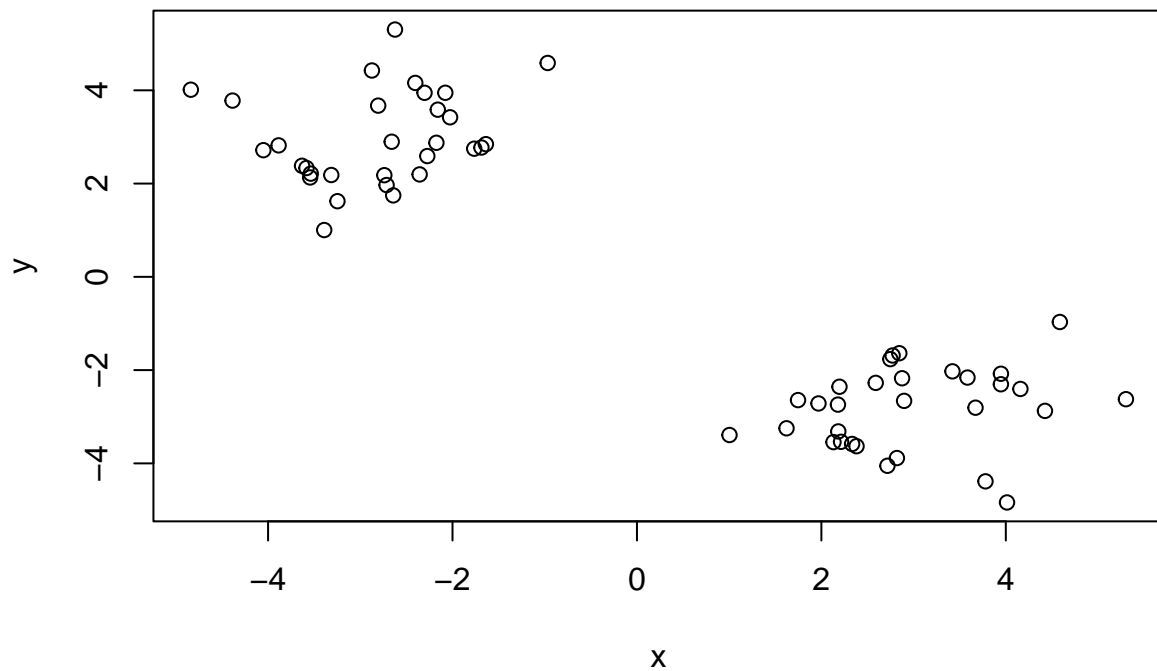
Sabrina Wu

2024-10-22

#First up kmeans()

Demo of using kmeans() function in base R. First make up some data with known structure.

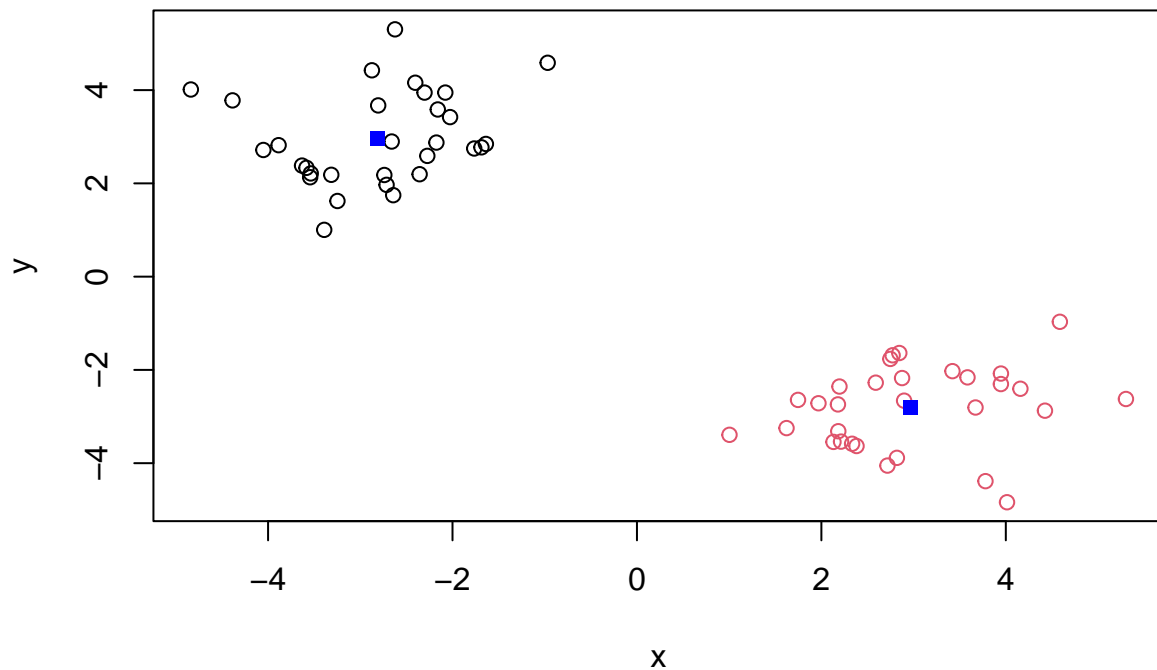
```
tmp <- c(rnorm(30,-3), rnorm(30,3))  
x <- cbind(x=tmp, y=rev(tmp))  
plot(x)
```



Now we have made up data in x let's see how kmeans works with this data

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

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



Now for hclust()

We will cluster the same data `x` with the `hclust()`. In this case `hclust()` requires a distance matrix as input.

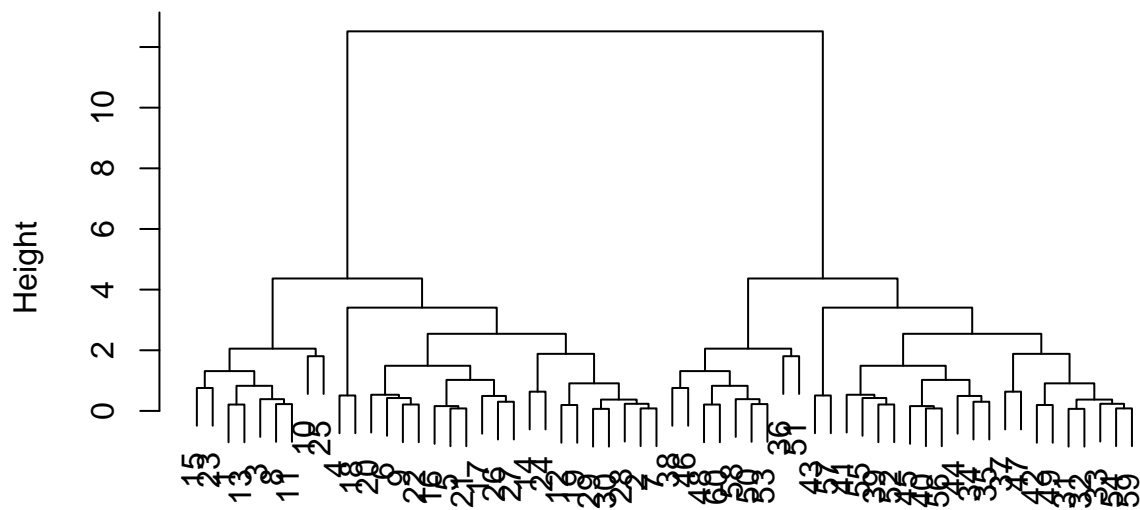
```
hc <- hclust( dist(x))
hc
```

```
##
## Call:
## hclust(d = dist(x))
##
## Cluster method   : complete
## Distance        : euclidean
## Number of objects: 60
```

Let's plot our hclust result

```
plot(hc)
```

Cluster Dendrogram



```
dist(x)
hclust (*, "complete")
```

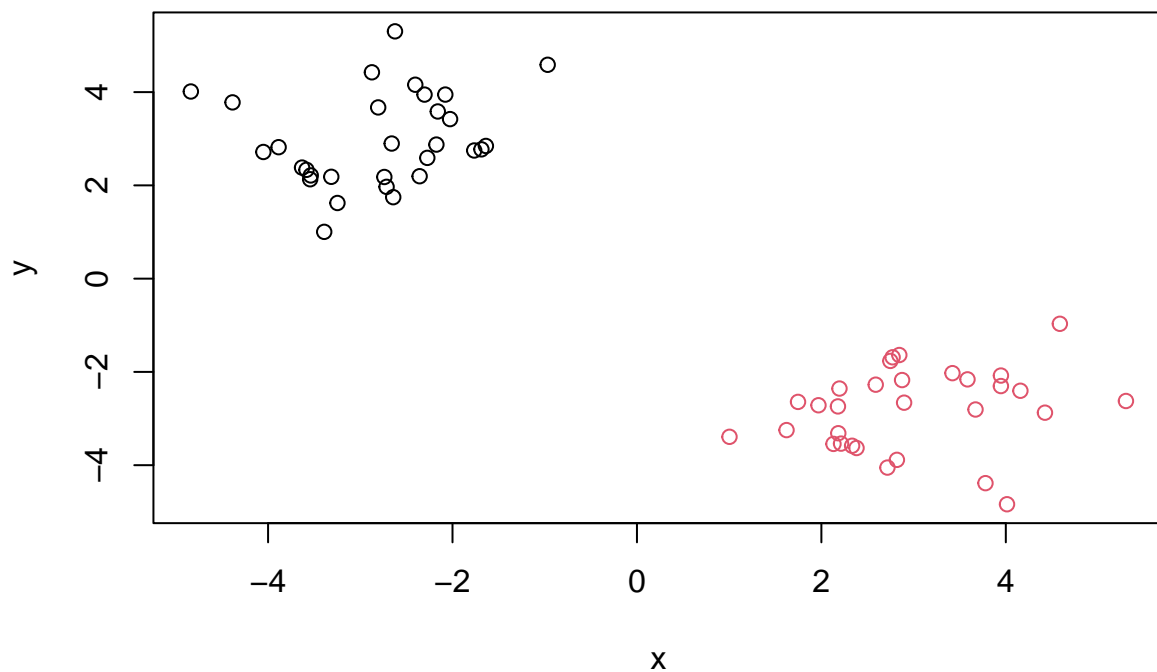
To get our cluster membership vector we need to “cut” the tree with the `cutree()`.

```
grps <- cutree(hc,h=8)
grps
```

[illegible]

Now plot our data with thte `hclust()` results

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



#Principal Component Abnalysis (PCA)

PCA of UK food data

Read data from website and try a few visualizations.

```
url <- "C:/Users/sabri/OneDrive/Desktop/BIMM 143/class07/UK_foods.csv"
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)
```

```
## [1] 17  5
```

There are 17 rows and 5 columns in the dataset.

```
## Preview the first 6 rows
```

```
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
```

```
# Moving the names from 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
```

```
#New Dimension after move
```

```
dim(x)
```

```
## [1] 17  4
```

-or- can also just set this from the beginning

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

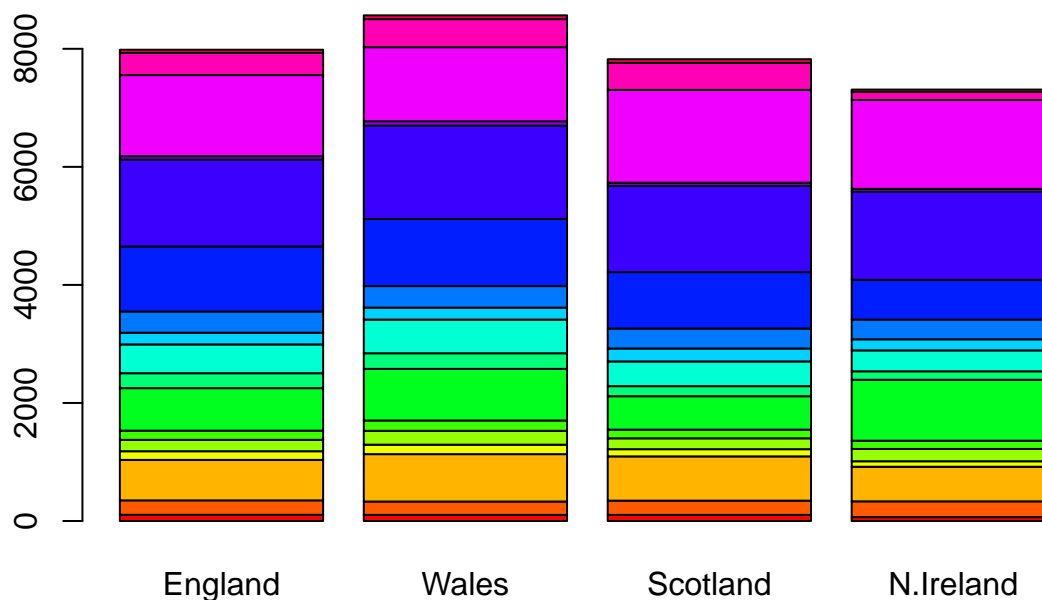
```
head(x)
```

```
##           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
```

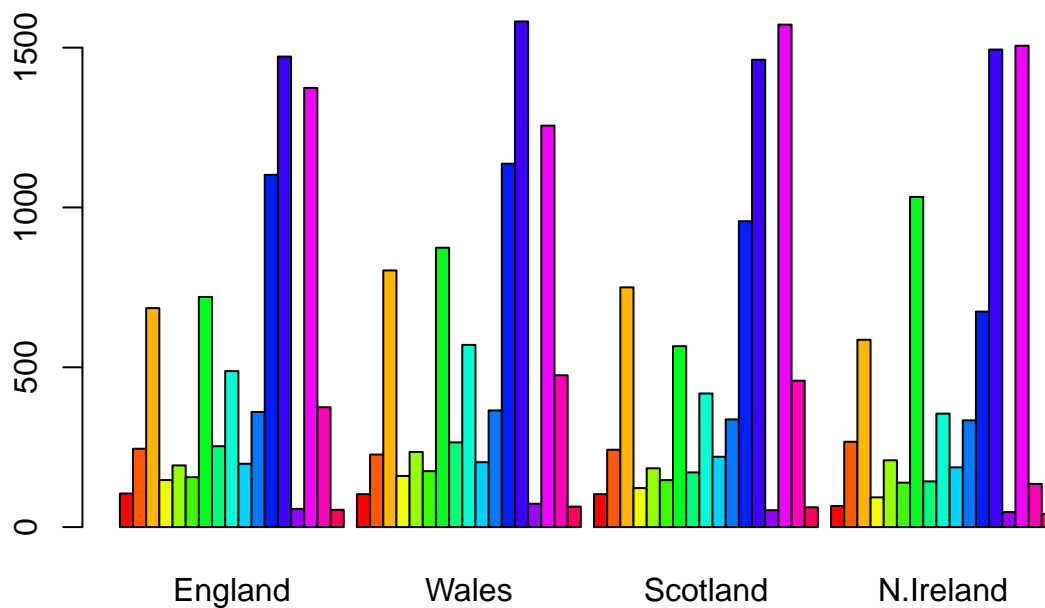
Q2. Which approach to solving the ‘row-names problem’ mentioned above do you prefer and why? Is one approach more robust than another under certain circumstances?

The second method of putting the command directly into the `read.csv()` is more robust. If you run the first approach of `x[, -1]` it will keep subtracting the first column beyond just the name column.

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



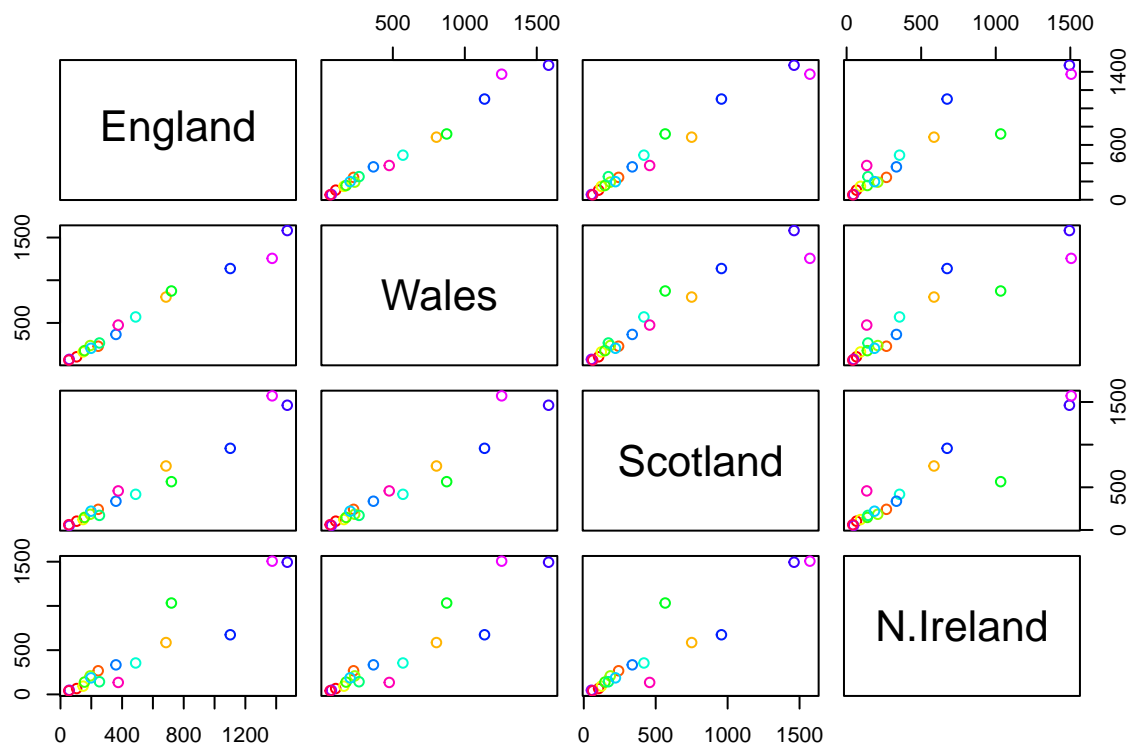
```
barplot(as.matrix(x), col=rainbow(nrow(x)), beside=TRUE)
```



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

By changing the `beside` part of the function into `beside=FALSE` or just deleting it, it makes it into a stacked column. The default is set as false therefore just deleting it will do the same thing.

```
pairs(x, 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?

The code `pair()` compares each of the countries with each other. If a give point lies on the the diagonal for a given plot, they consumed the same amount of that food item in both countries.

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

N.Ireland varies a lot more on the blue point (fresh fruits) which is much higher above the diagonal and the green point (potatoes) which is lower than the diagonal compared to the other countries of the UK.

PCA to the rescue!! The main base R PCA function is called `pcomp()` and we will need to give it the tranpose of our input data!

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

```
## Importance of components:
##
## 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
```

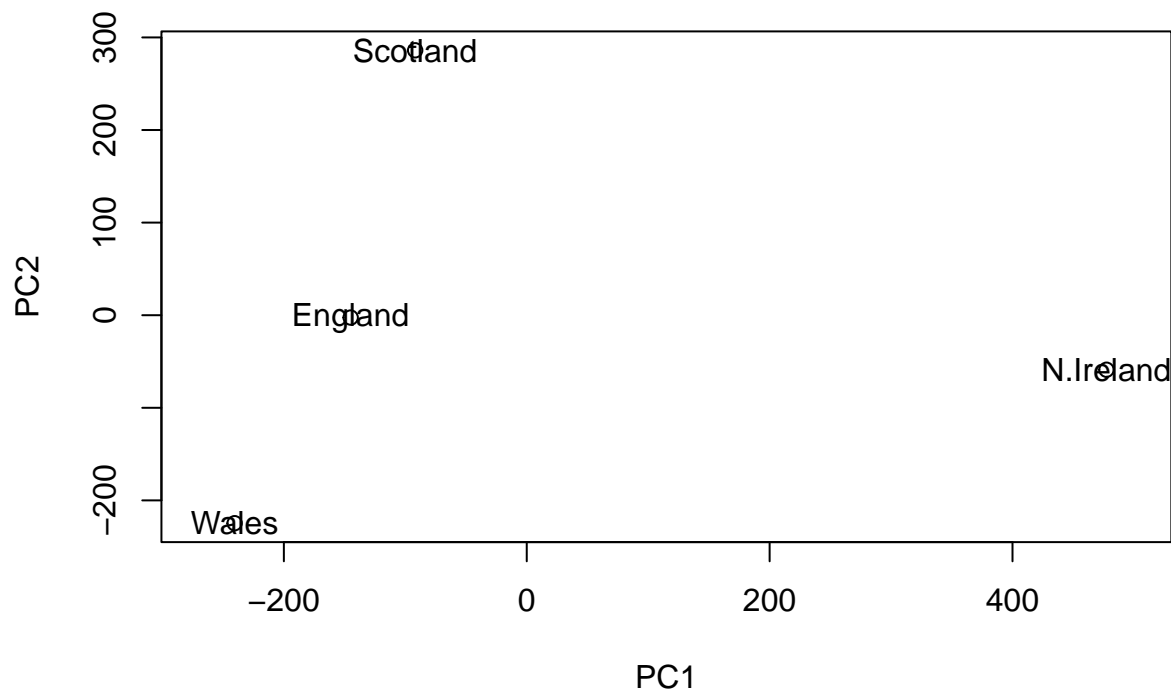
```
attributes(pca)
```

```
## $names
## [1] "sdev"      "rotation" "center"    "scale"     "x"
##
## $class
## [1] "prcomp"
```

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

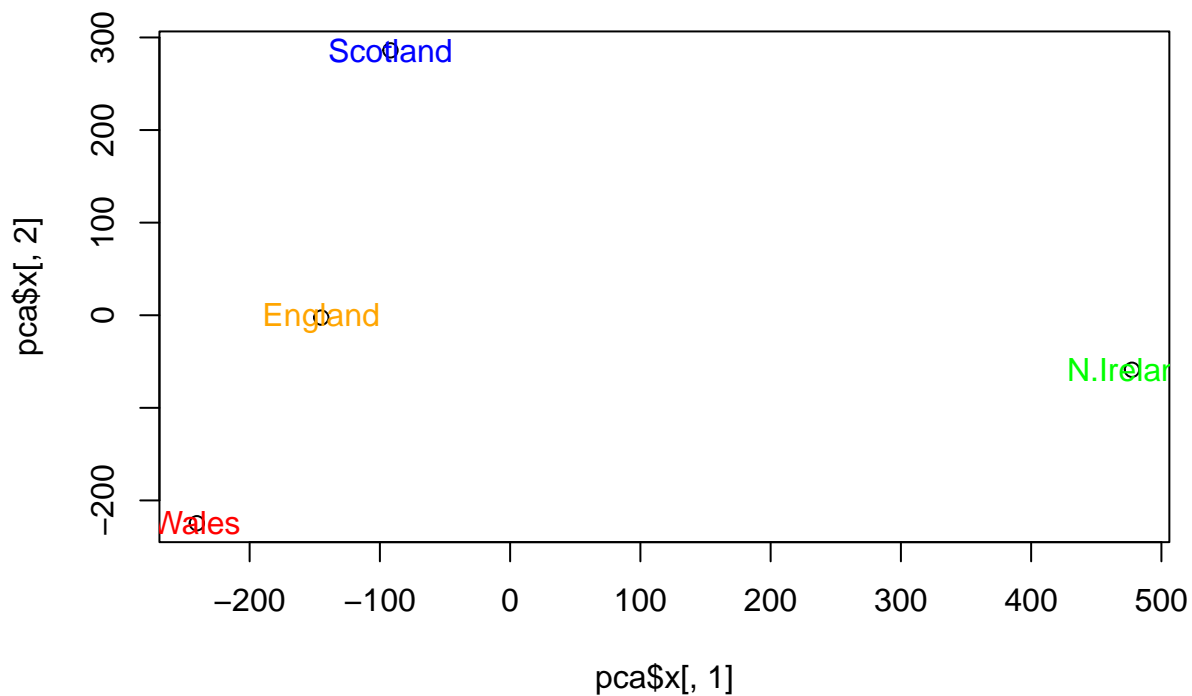
To make our new PCA plot (a.k.a. PCA score plot) we access `pca$x`

```
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.

```
country_cols <- c("orange", "red", "blue", "green")
plot(pca$x[,1], pca$x[,2])
text(pca$x[,1], pca$x[,2], colnames(x), col=country_cols)
```



Below we can use the square of `pca$sdev`, which stands for “standard deviation”, to calculate how much variation in the original data each PC accounts for.

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

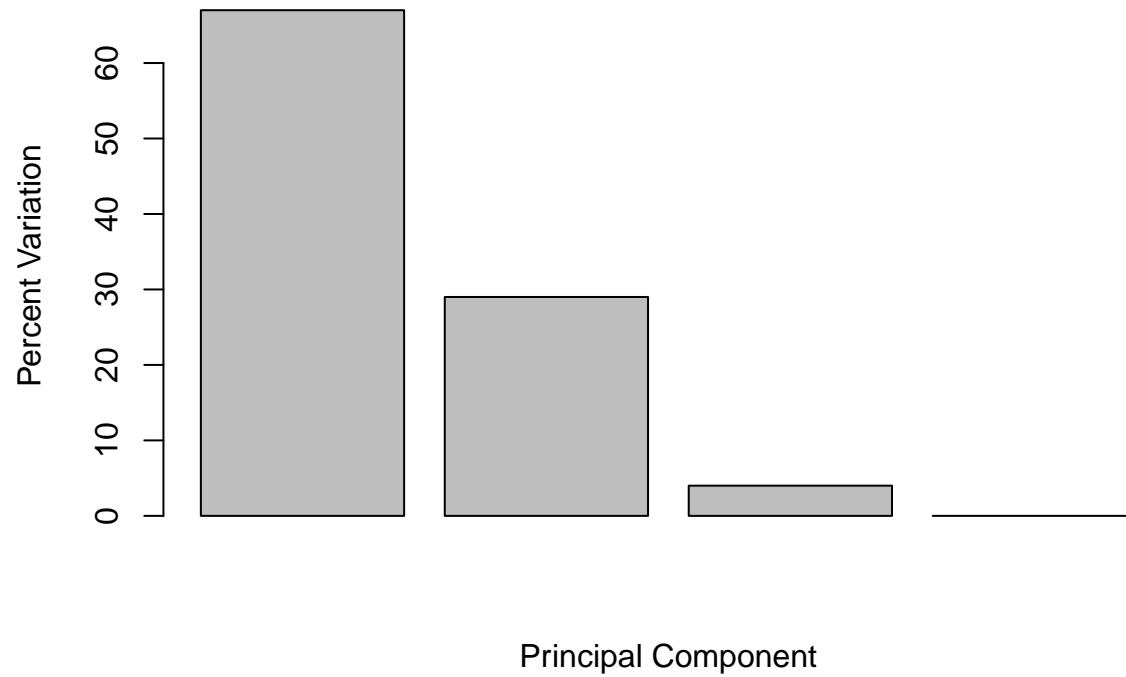
```
## [1] 67 29 4 0
```

```
## or the second row here...
```

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

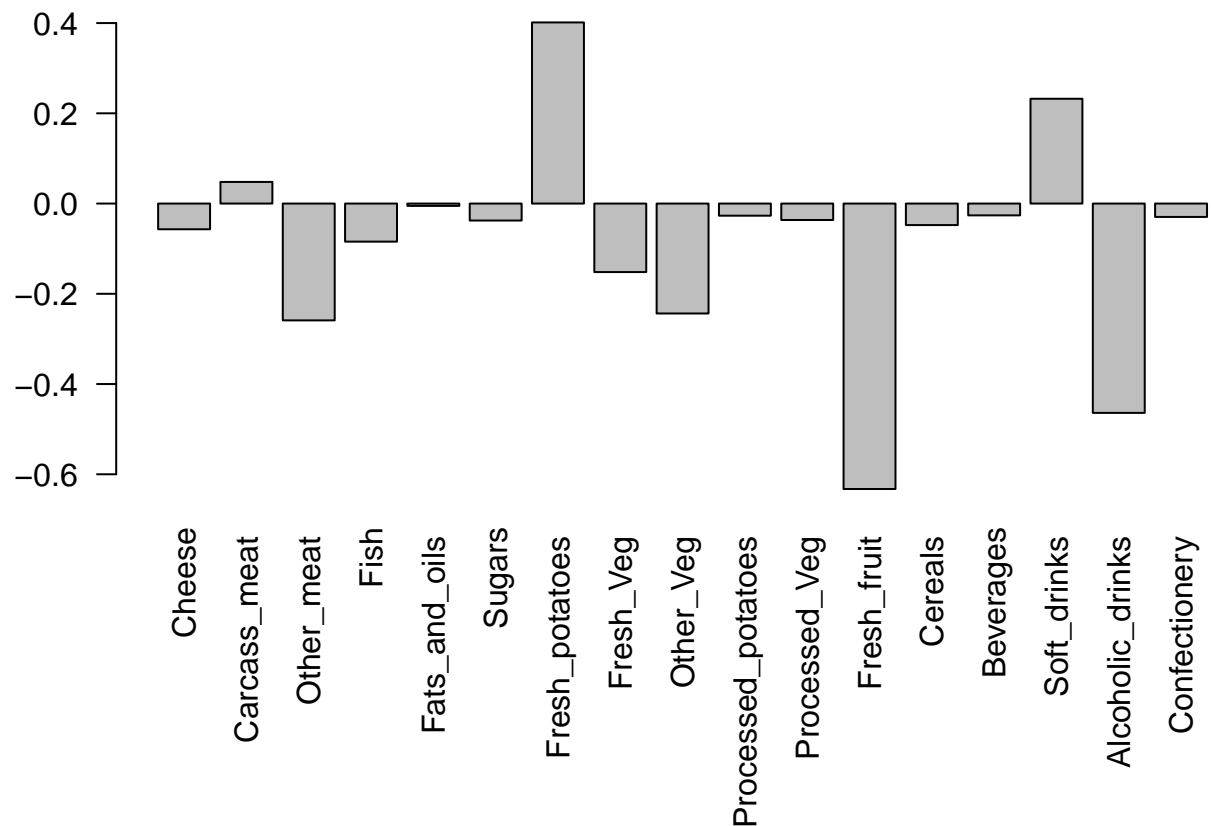
```
##              PC1      PC2      PC3      PC4
## Standard deviation 324.15019 212.74780 73.87622 3.175833e-14
## Proportion of Variance 0.67444 0.29052 0.03503 0.000000e+00
## Cumulative Proportion 0.67444 0.96497 1.00000 1.000000e+00
```

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



Digging Deeper (variable loadings)

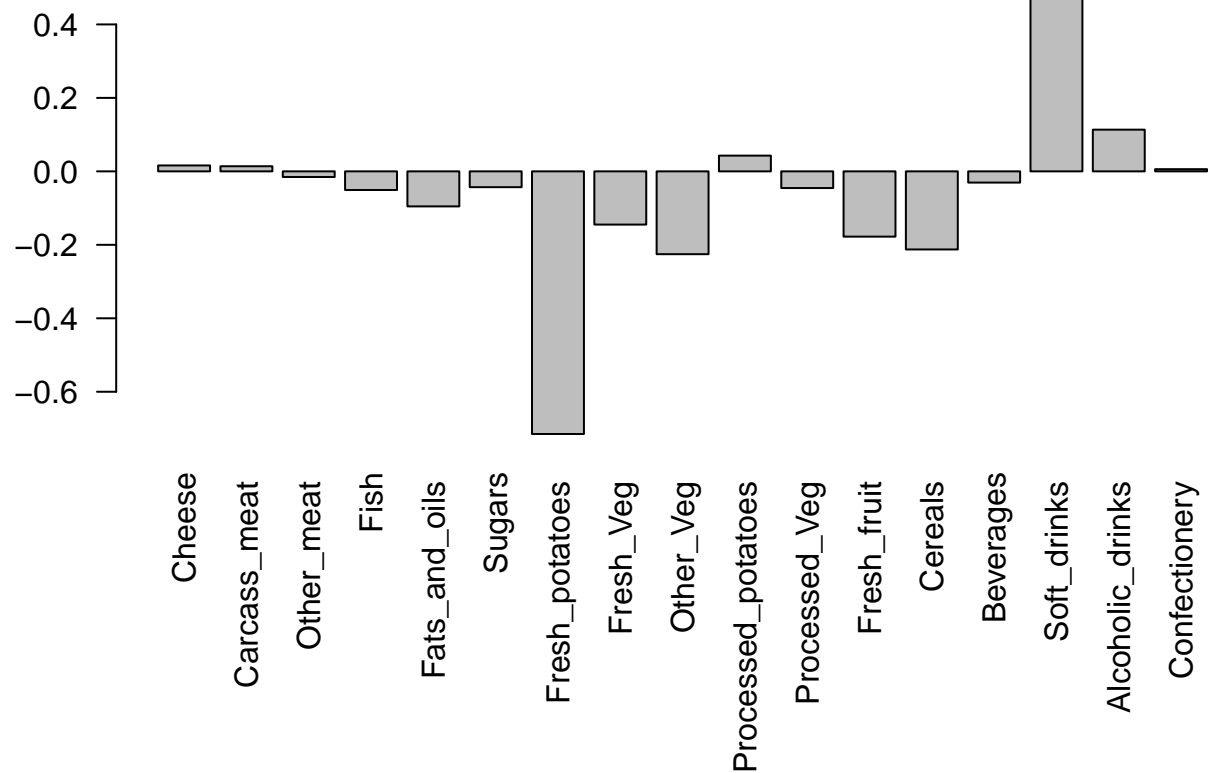
```
## Lets focus on PC1  
par(mar=c(10, 3, 0.35, 0))  
barplot( pca$rotation[,1], las=2 )
```



The largest positive loading scores that push N.Ireland to right positive side is potatoes and soft drinks. The highest negative scores that push other countries to the left side of the plot is fresh fruit and alcoholic drinks.

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 )
```



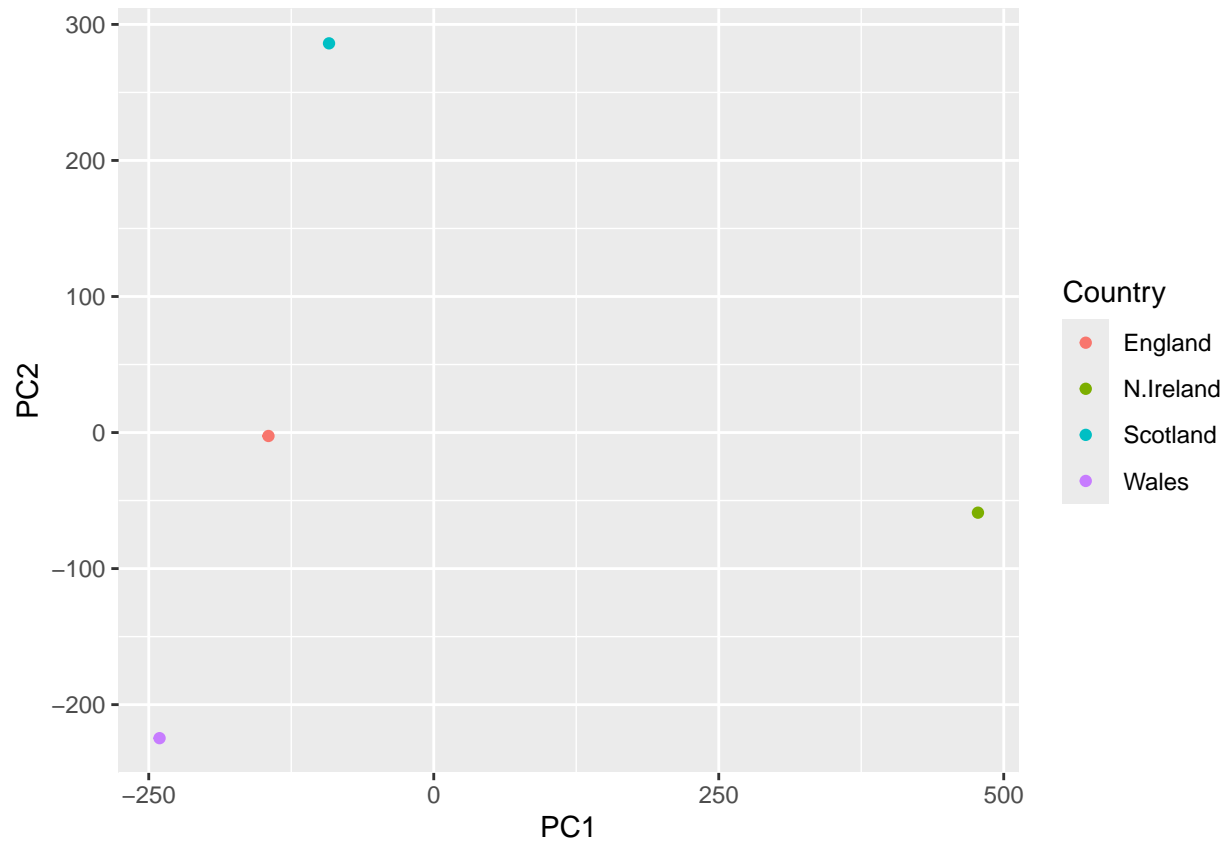
The largest positive loading scores in PC2 that push N.Ireland to right positive side is soft drinks. The highest negative scores that push other countries to the left side of the plot is fresh potatoes.

Using ggplot

```
library(ggplot2)

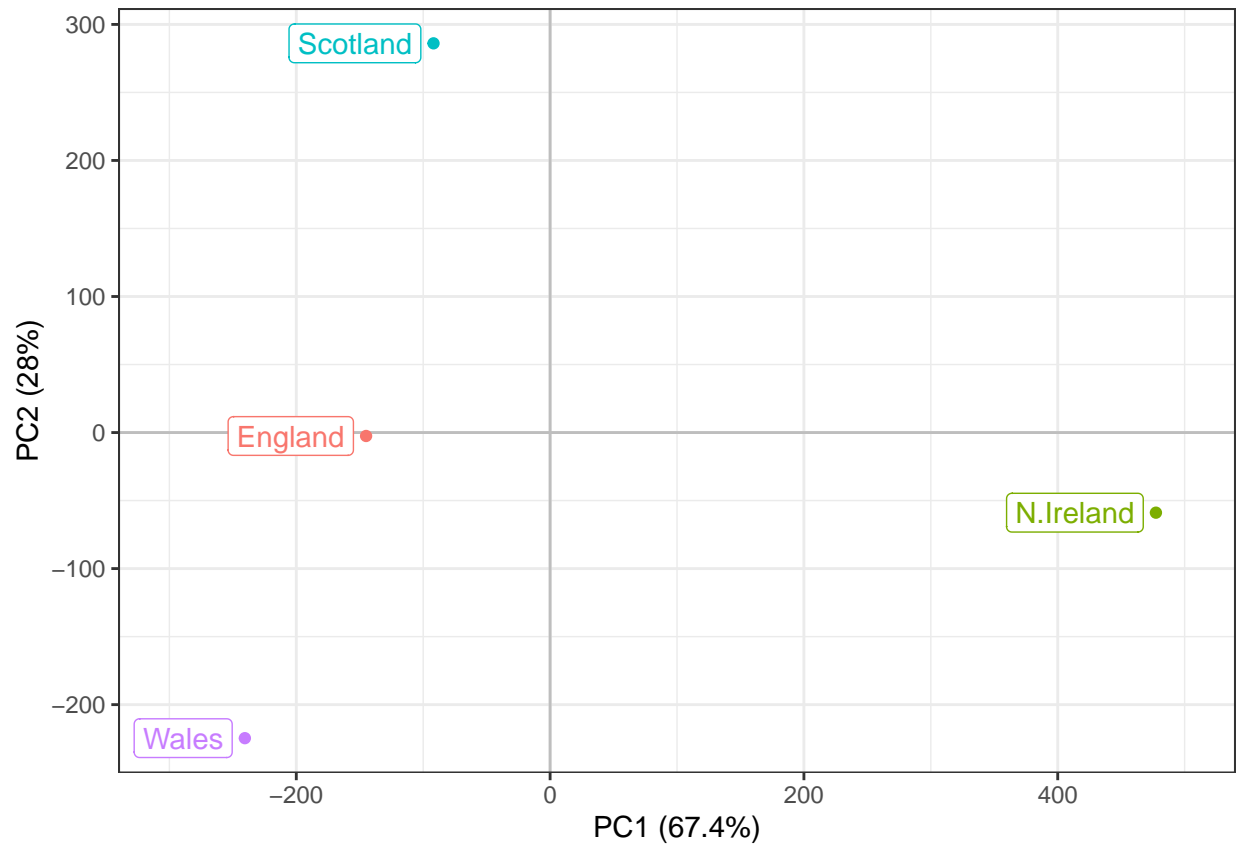
df <- as.data.frame(pca$x)
df_lab <- tibble::rownames_to_column(df, "Country")

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



Make it look nicer

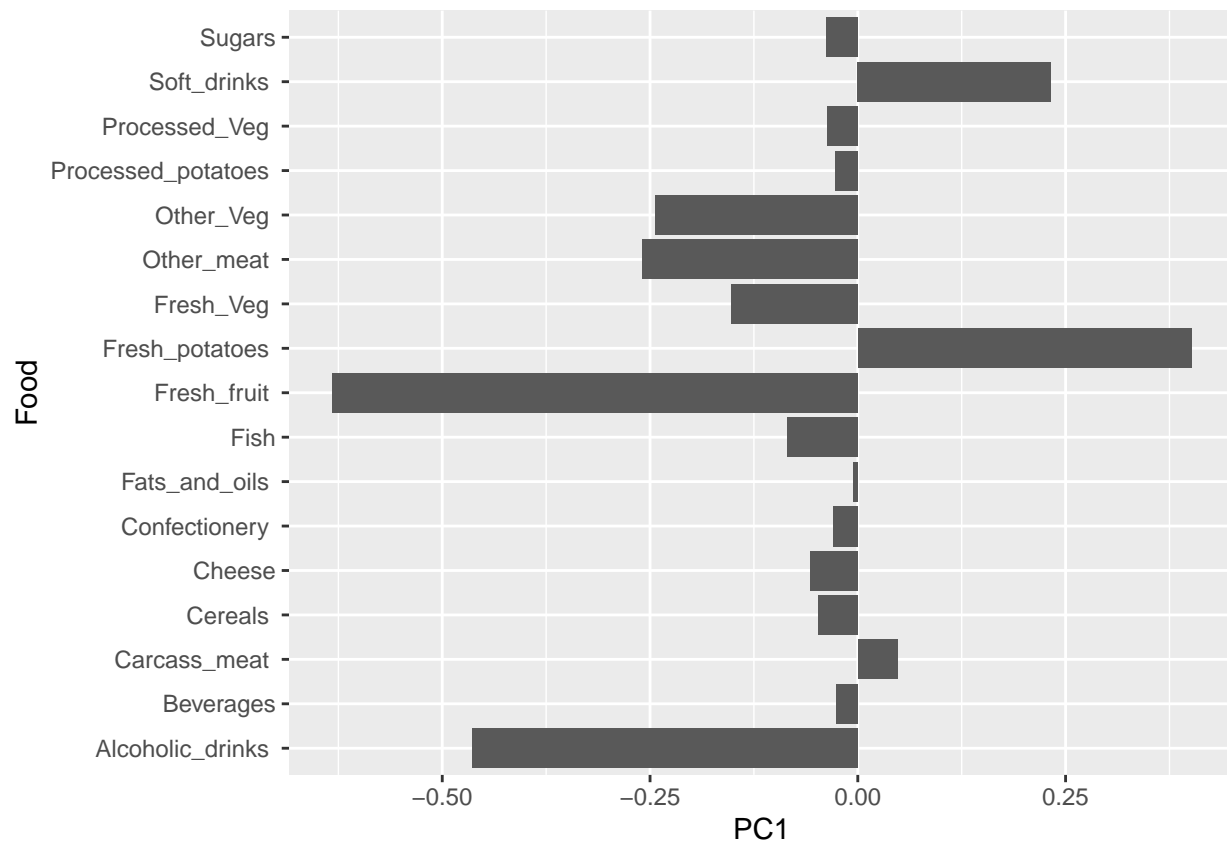
```
ggplot(df_lab) +  
  aes(PC1, PC2, col=Country, label=Country) +  
  geom_hline(yintercept = 0, col="gray") +  
  geom_vline(xintercept = 0, col="gray") +  
  geom_point(show.legend = FALSE) +  
  geom_label(hjust=1, nudge_x = -10, show.legend = FALSE) +  
  expand_limits(x = c(-300,500)) +  
  xlab("PC1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



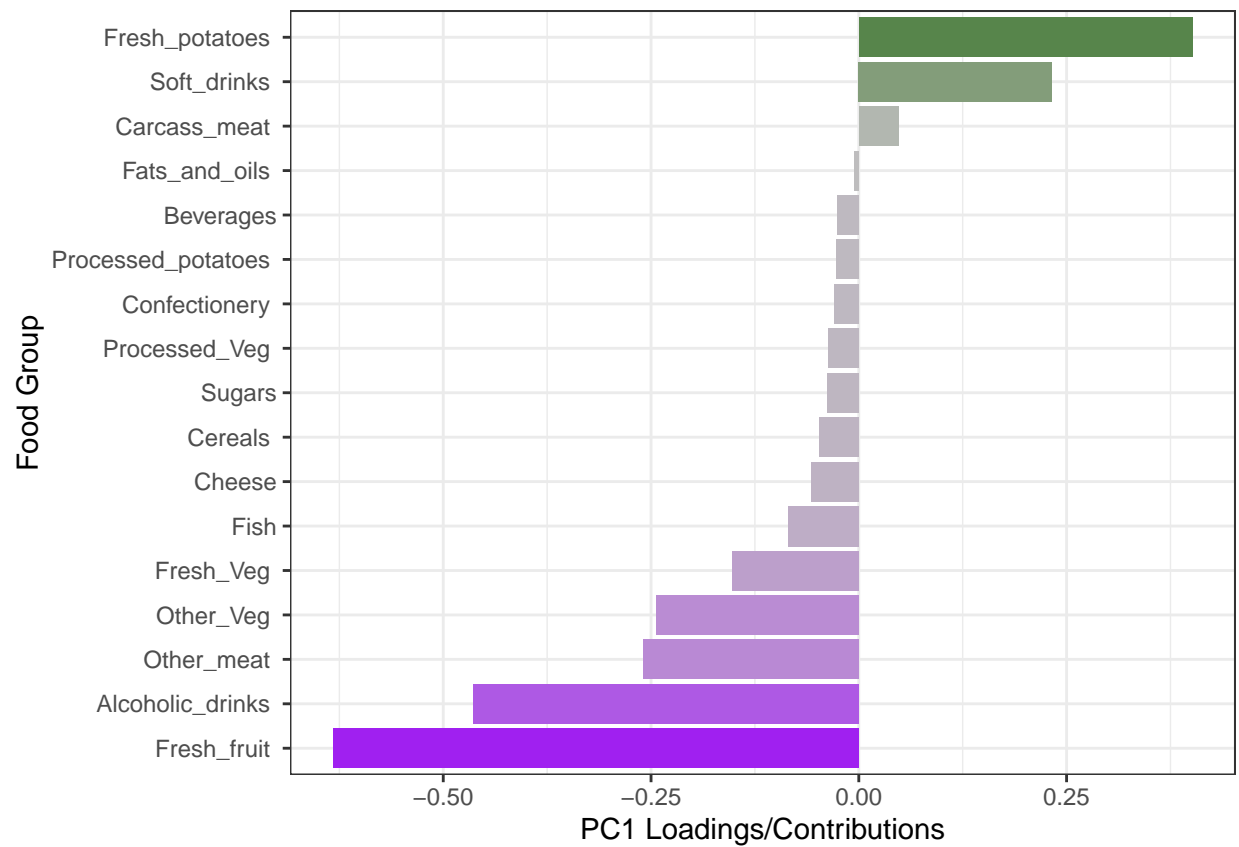
Graphing the loadings

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

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

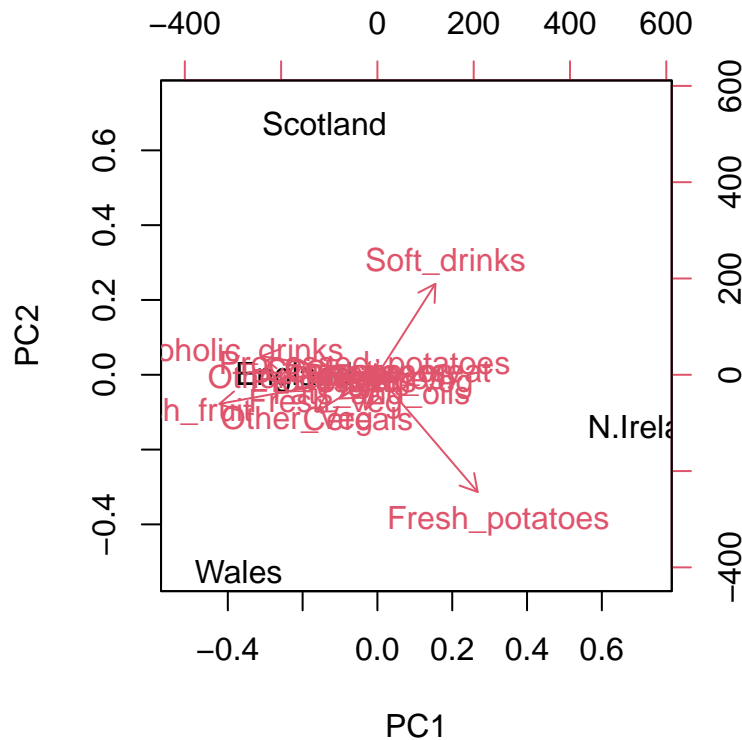



```
ggplot(ld_lab) +
  aes(PC1, reorder(Food, PC1), bg=PC1) +
  geom_col() +
  xlab("PC1 Loadings/Contributions") +
  ylab("Food Group") +
  scale_fill_gradient2(low="purple", mid="gray", high="darkgreen", guide=NULL) +
  theme_bw()
```



###Biplots

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



PCA of RNA-seq data

```
url2 <- "C:/Users/sabri/OneDrive/Desktop/BIMM 143/class07/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?

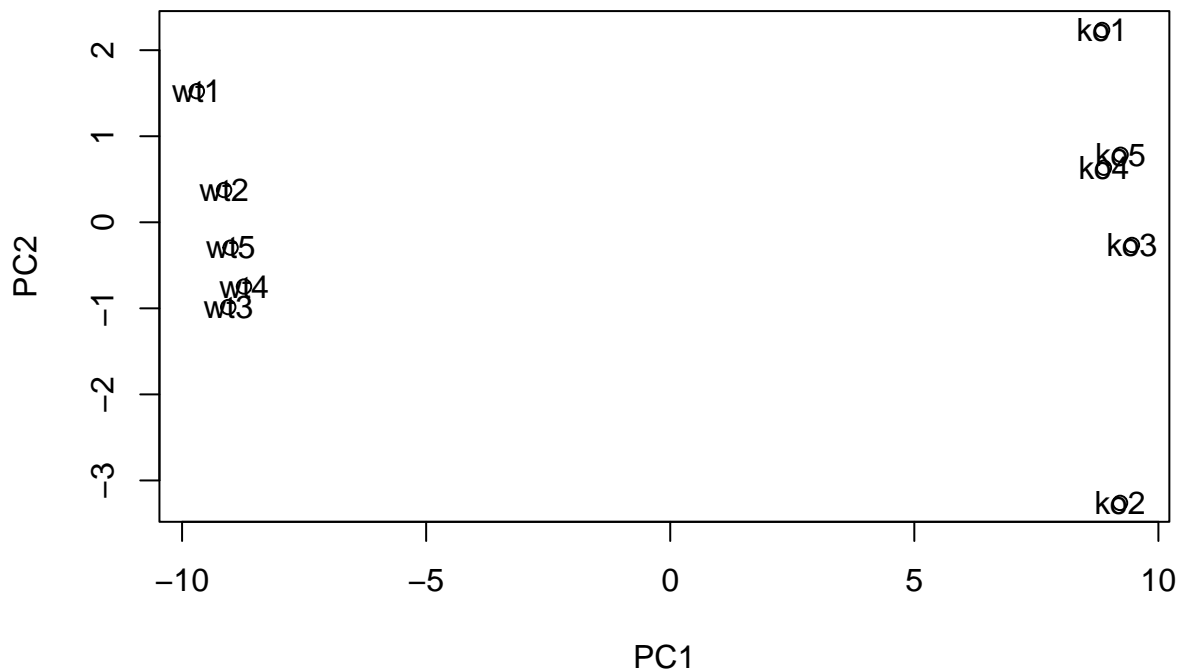
```
dim(rna.data)
```

```
## [1] 100  10
```

There are 100 genes and 10 samples in the dataset.

```
## 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")
text(pca$x[,1], pca$x[,2], colnames(rna.data))
```



```
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
```

92.6% of the variance is captured by PC1, and the first two PCs captures 94.9% of variance.

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

Quick scree plot



PC1 accounts for the majority.

Plotting the variance accounted by the difference PC

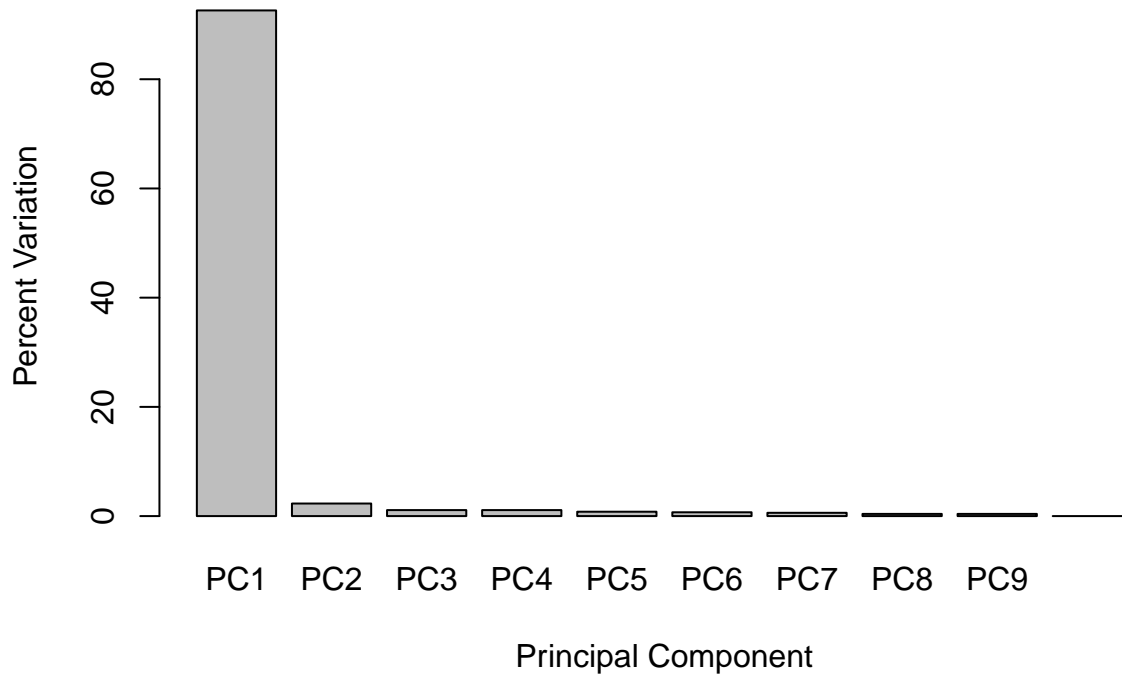
```
## Variance captured per PC  
pca.var <- pca$sdev^2
```

```
## 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")
```

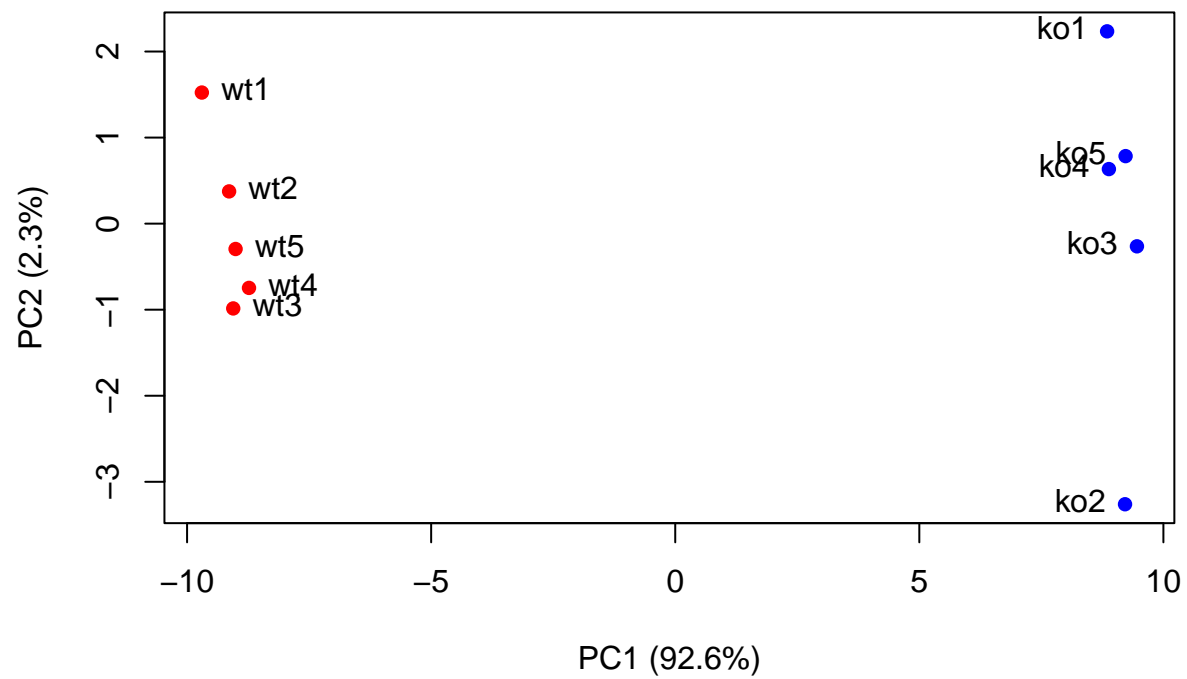
Scree Plot



```
## 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)))
```

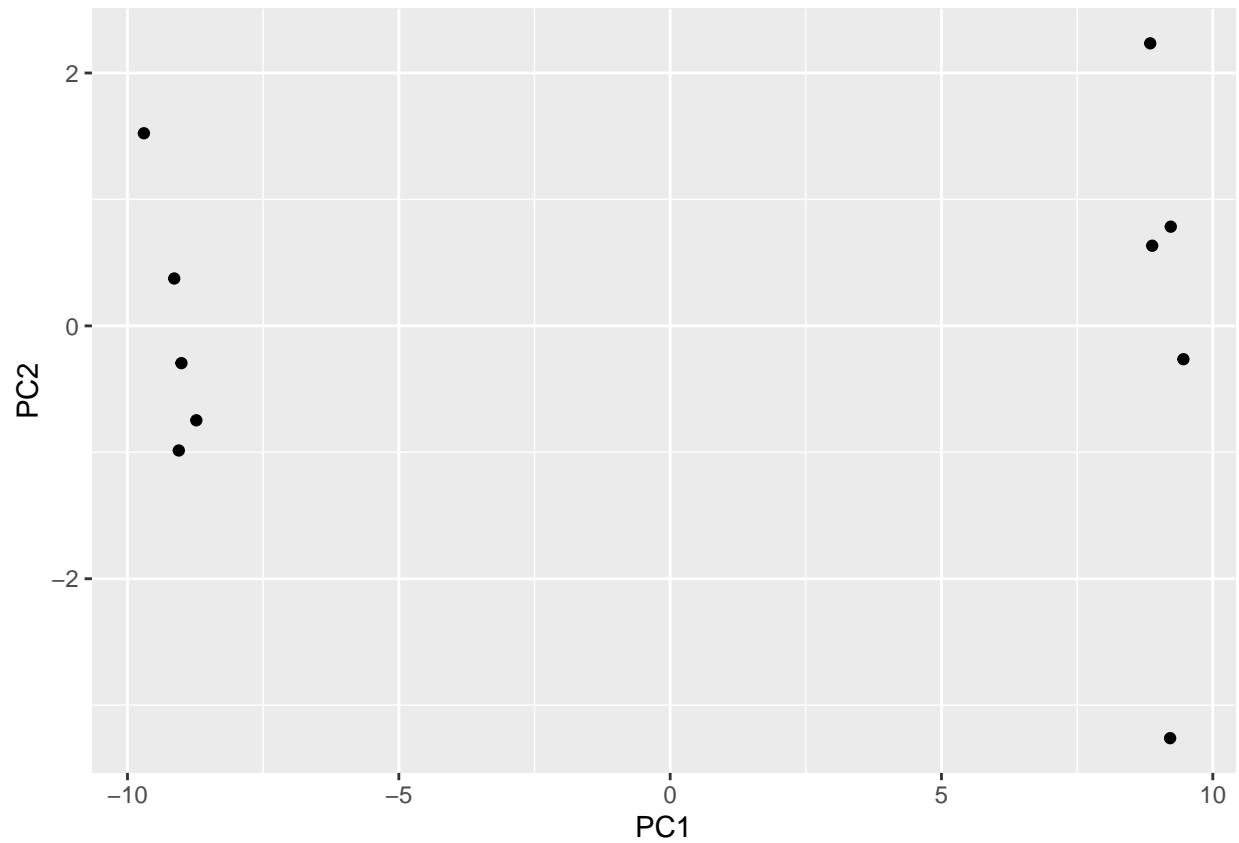


Using ggplot

```
library(ggplot2)

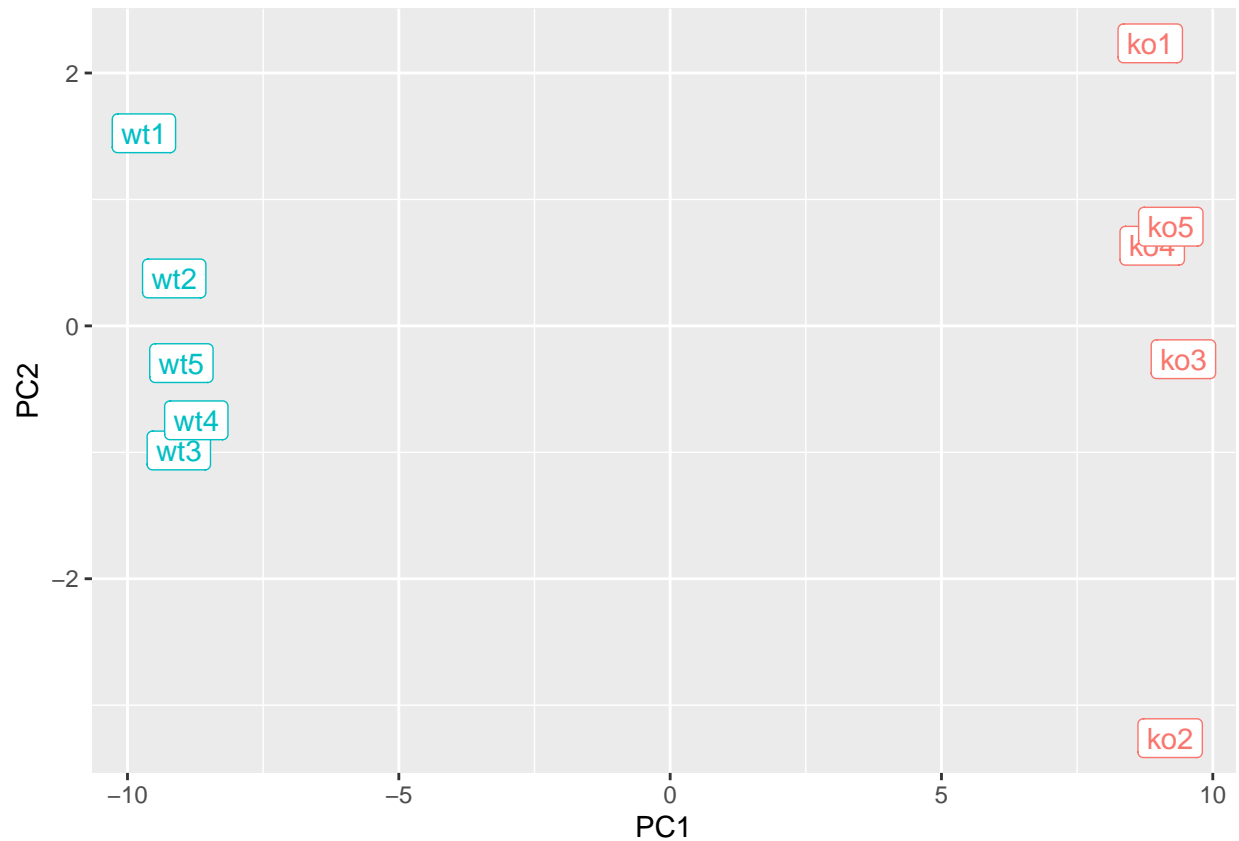
df <- as.data.frame(pca$x)

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



```
# 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
```

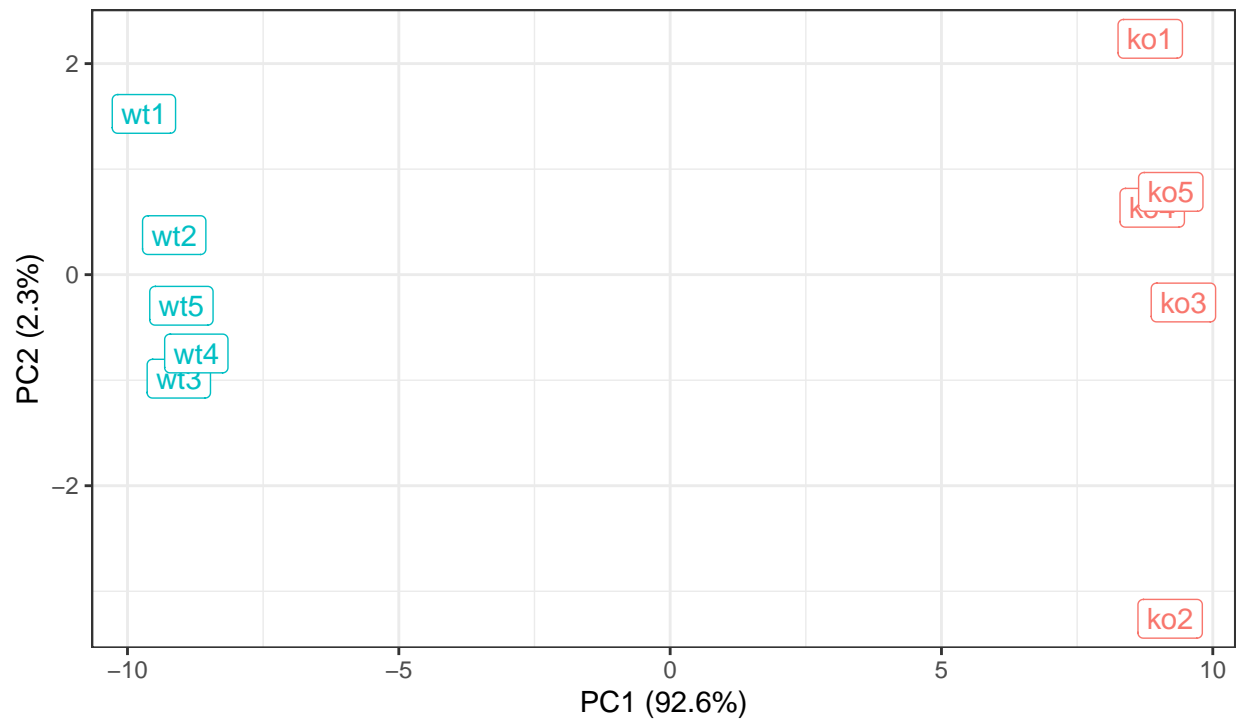



Adding titles

```
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 clearly separates wild-type from knock-out samples



Class example data

Finding the loading that contributes most to pc1

```
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"
```