

# Machine Learning 1

Julia Napoli

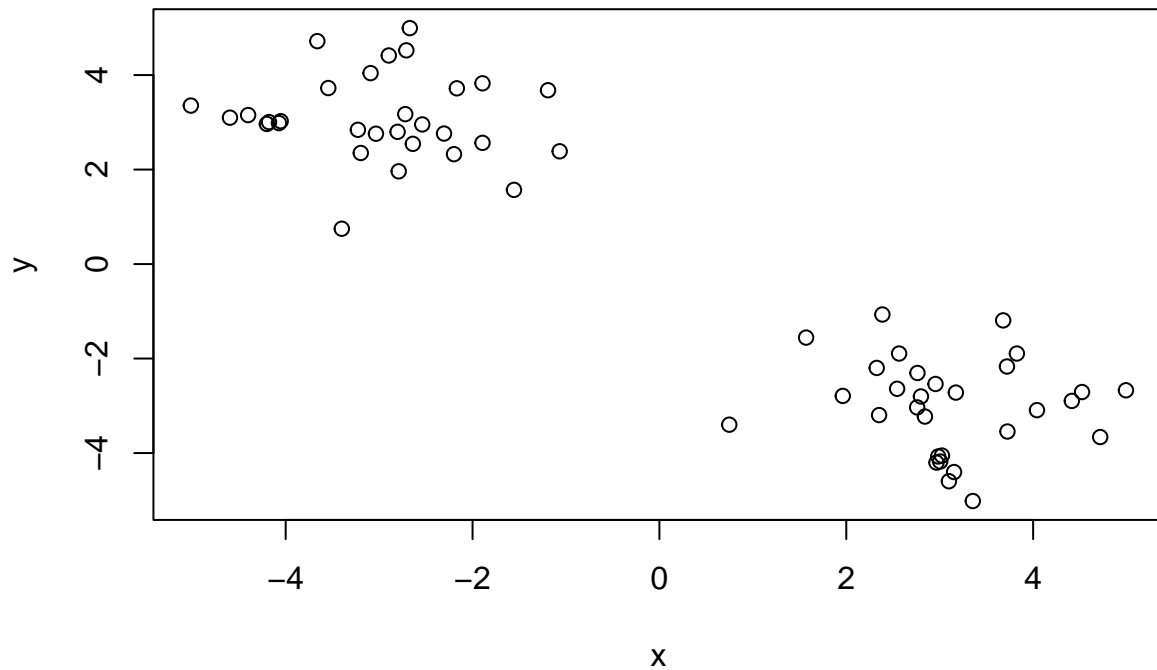
10/22/2021

## Clustering methods

Kmeans clustering in R is done with the 'kmeans()' function.

Here we make up some data to test and learn with.

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



```
hist(data)
```

A histogram showing the frequency distribution of data. The x-axis is labeled 'data' and ranges from -6 to 4. The y-axis is labeled 'Frequency' and ranges from 0 to 20. The distribution is bimodal, with peaks around -2.5 and 2.5.

data (bin center)	Frequency
-5.5	2
-4.5	12
-3.5	14
-2.5	22
-1.5	10
0.5	2
1.5	4
2.5	24
3.5	20
4.5	10

```
km <- kmeans(data, centers = 2, nstart = 20)
km
```

Q. How many points are in each cluster?

```
## [1] 30 30
```

```
km$cluster
```

km\$centers

```
plot(data, col = km$cluster)
points(km$centers, col = "blue", pch = 15, cex=2)
```



## Heirarchical clustering

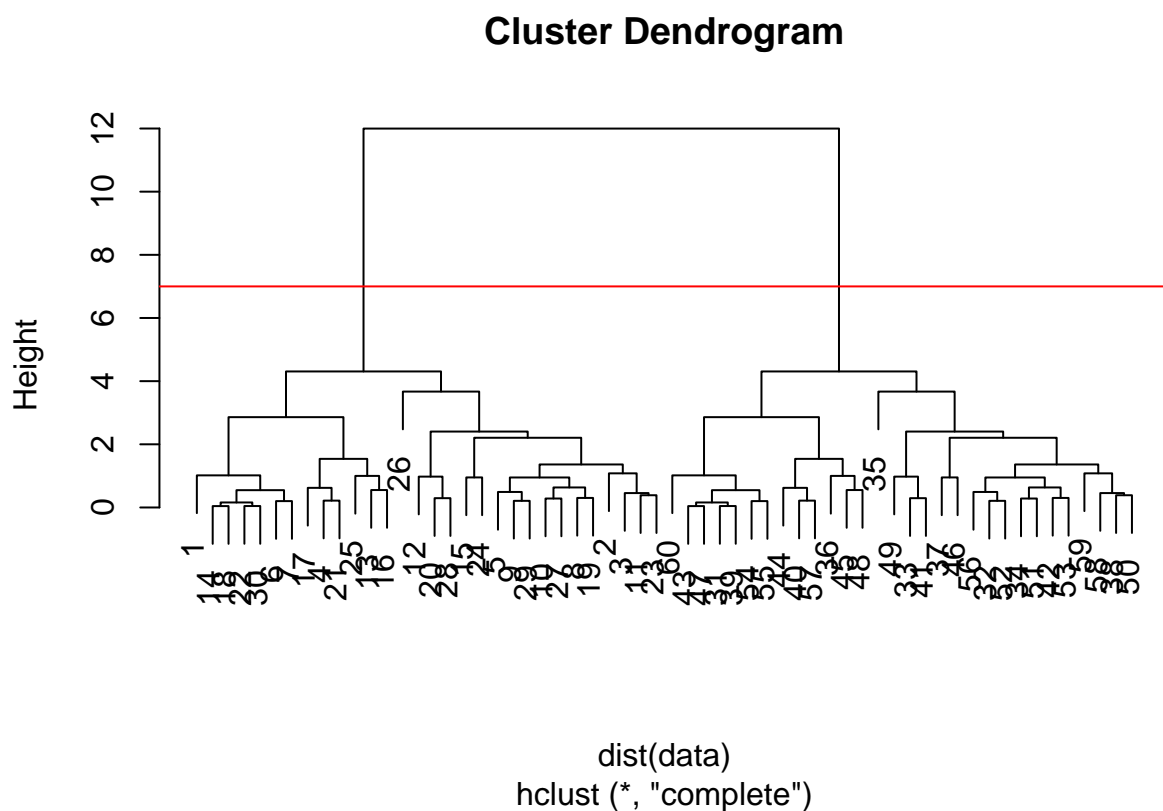
We will use the `hclust()` function on the same data as before and see how this method works.

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

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

`hclust` has a plot method

```
plot(hc)
abline(h=7, col = "red")
```



To find our membership vector we need to “cut” the tree (dendrogram) and for this we use the `cutree()` function and tell it the height to cut at.

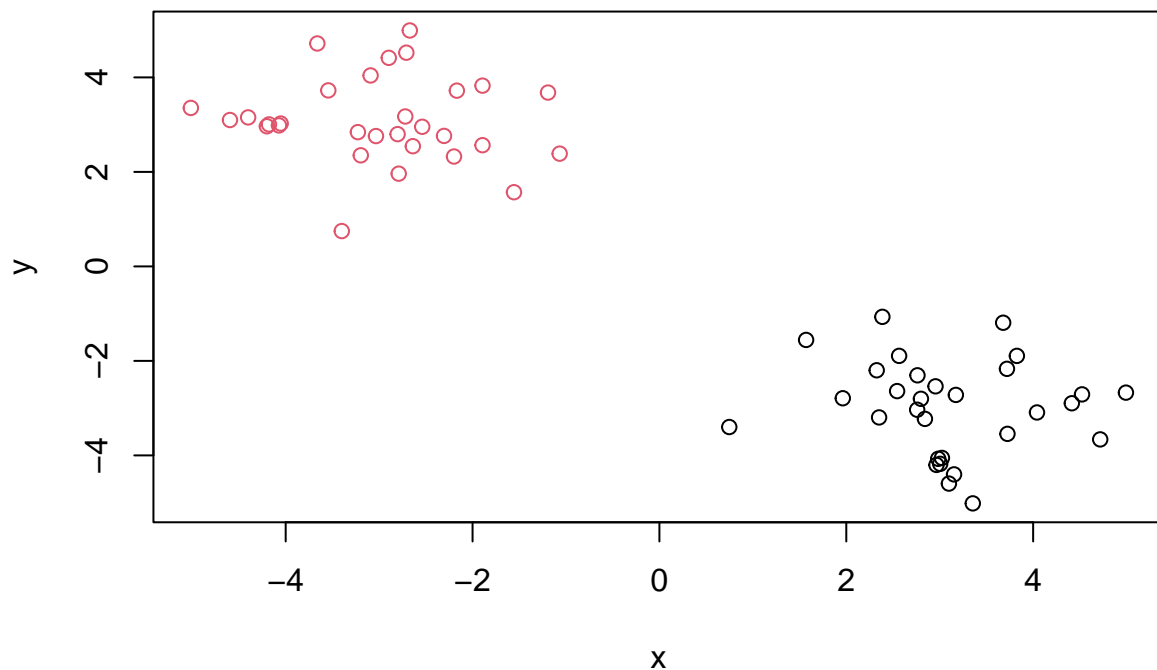
```
cutree(hc, h=7)
```

```
## [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 2 2
```

We can also use `cutree()` and state the number of clusters we want...

```
grps <- cutree(hc,k=2)
```

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



## Principal Component Analysis (PCA)

PCA is a super useful analysis method when you have lots of dimensions in your data...

### PCA of UK food data

```
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?

```
nrow(x)
```

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

```
ncol(x)
```

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

```
dim(x)
```

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

```
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
```

```
rownames(x) <- x[,1]
x <- x[,-1]
ncol(x)
```

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

```
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
## Fresh_potatoes 720  874      566      1033
## Fresh_Veg    253   265      171       143
## Other_Veg    488   570      418       355
## Processed_potatoes 198  203      220       187
## Processed_Veg 360   365      337       334
## 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
```

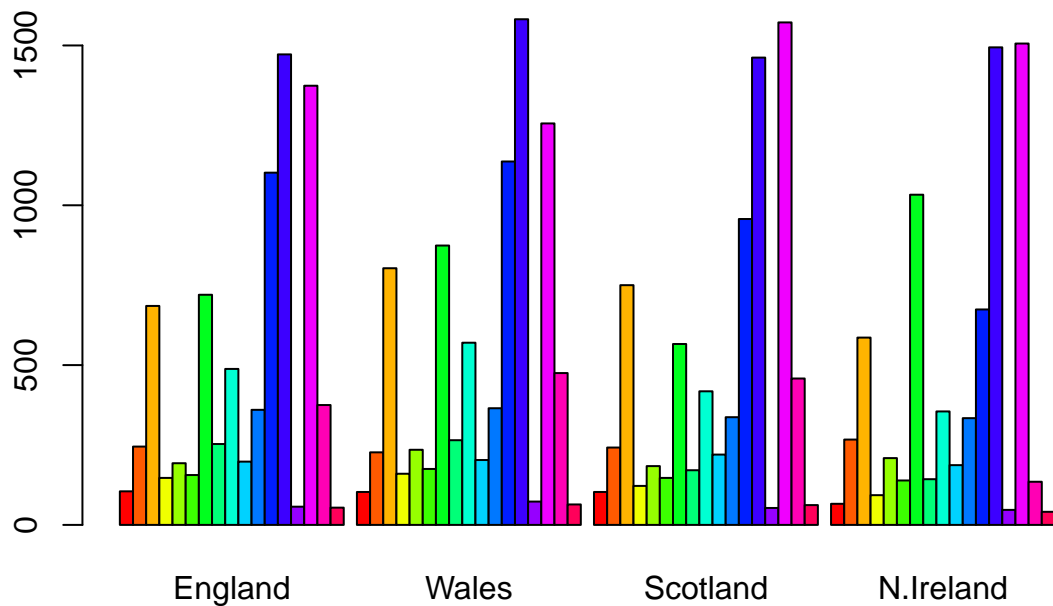
```
x <- read.csv(url, row.names=1)
head(x)
```

```
##           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
```

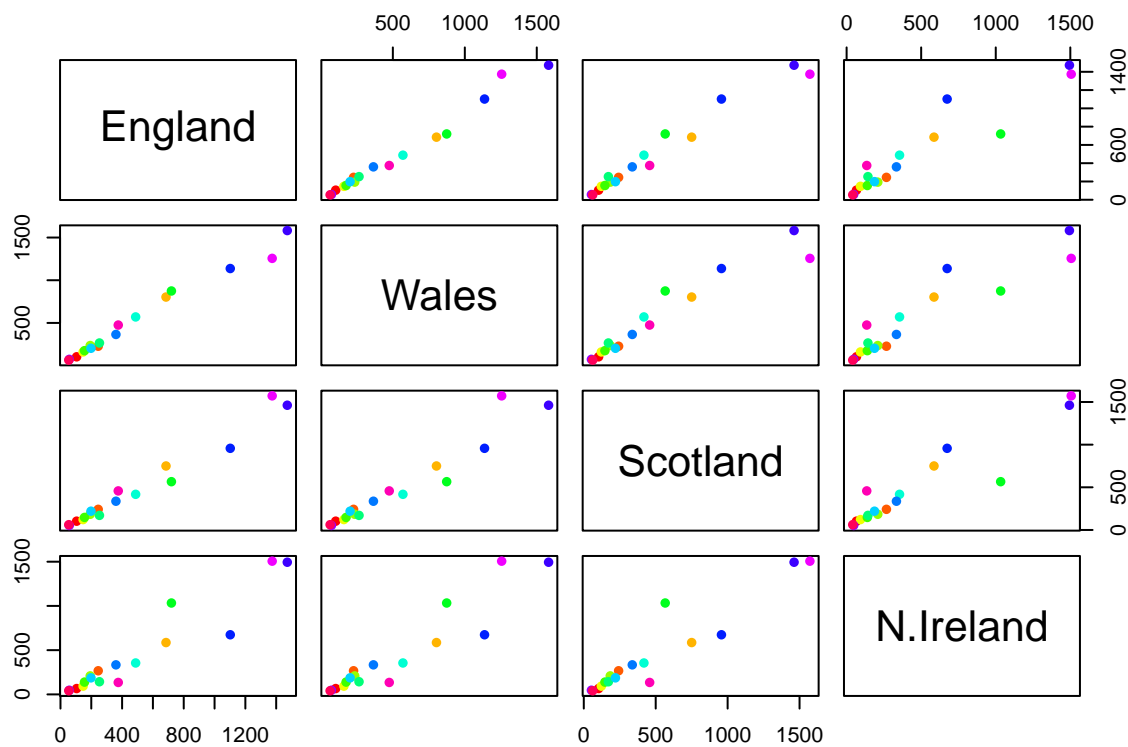
```
dim(x)
```

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

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



```
mycols <- rainbow(nrow(x))
pairs(x, col=mycols, pch = 16)
```



## PCA to the rescue!

Here we will use the base R function for PCA, which is called `prcomp()`. This function wants the transpose of our data.

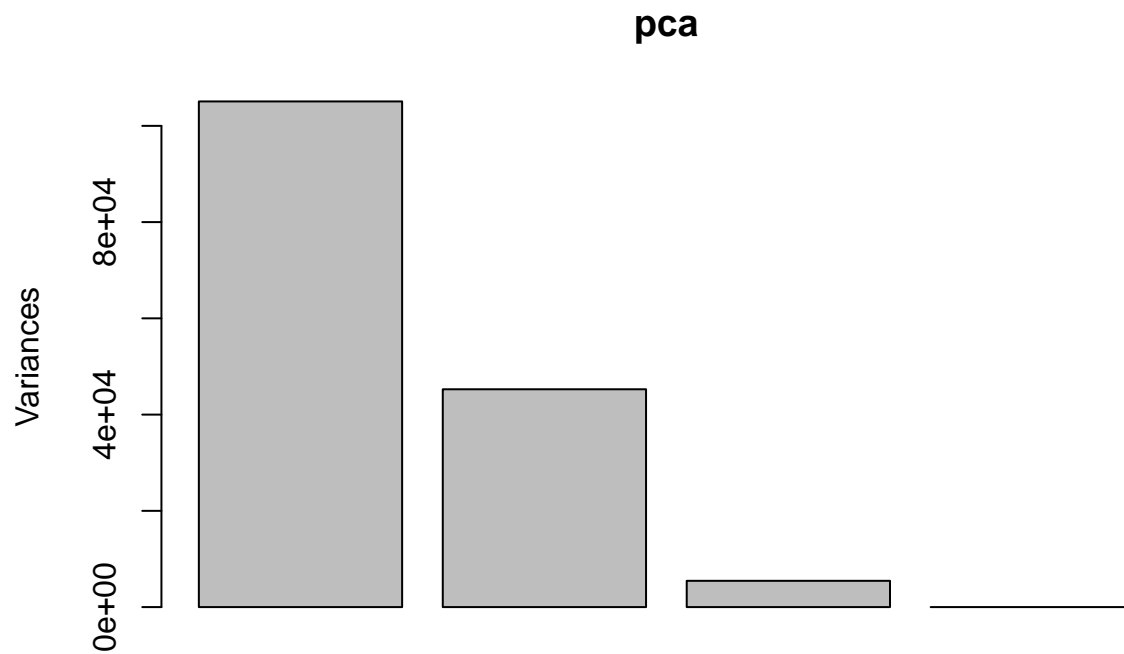
```
# t(x)
pca <- prcomp(t(x))
summary(pca)
```

## Importance of components:

##	PC1	PC2	PC3	PC4
## Standard deviation	324.1502	212.7478	73.87622	4.189e-14
## Proportion of Variance	0.6744	0.2905	0.03503	0.000e+00
## Cumulative Proportion	0.6744	0.9650	1.00000	1.000e+00

```
plot(pca)
```





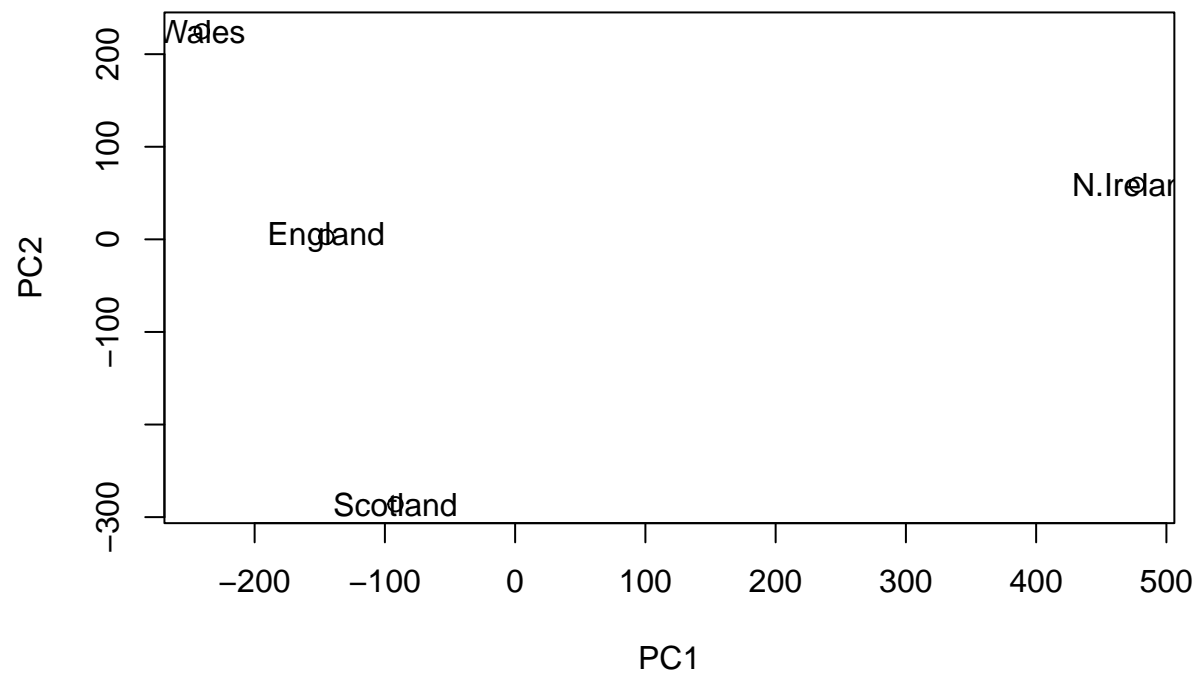
We want to score plot (aka PCA plot). Basically of PC1 vs PC2.

```
attributes(pca)
```

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

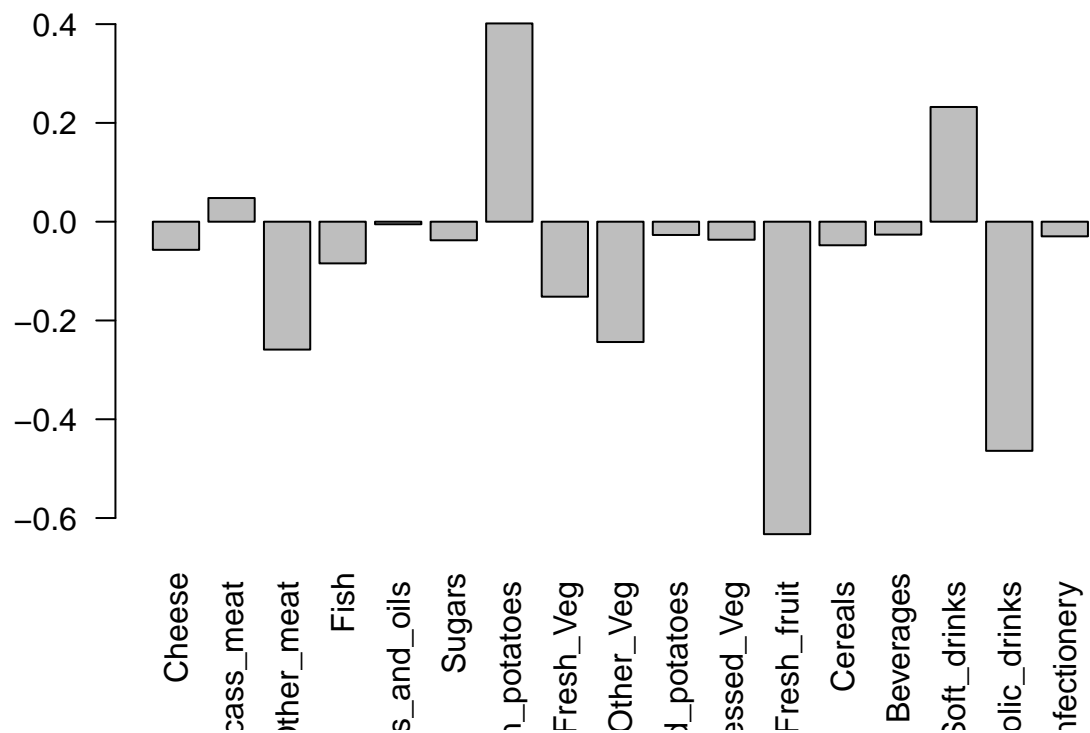
We are after the `pca$x` component for this plot...

```
plot(pca$x[,1:2])
text(pca$x[,1:2], labels = colnames(x))
```



We can also examine the PCA “loadings”, which tell us how much the original variables contribute to each new PC...

```
barplot(pca$rotation[,1], las = 2)
```



## One more PCA for today

```
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
```

```
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
```

```
colnames(rna.data)
```

```
## [1] "wt1" "wt2" "wt3" "wt4" "wt5" "ko1" "ko2" "ko3" "ko4" "ko5"
```

```
ncol(rna.data)
```

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

Let's do some RNA seq...

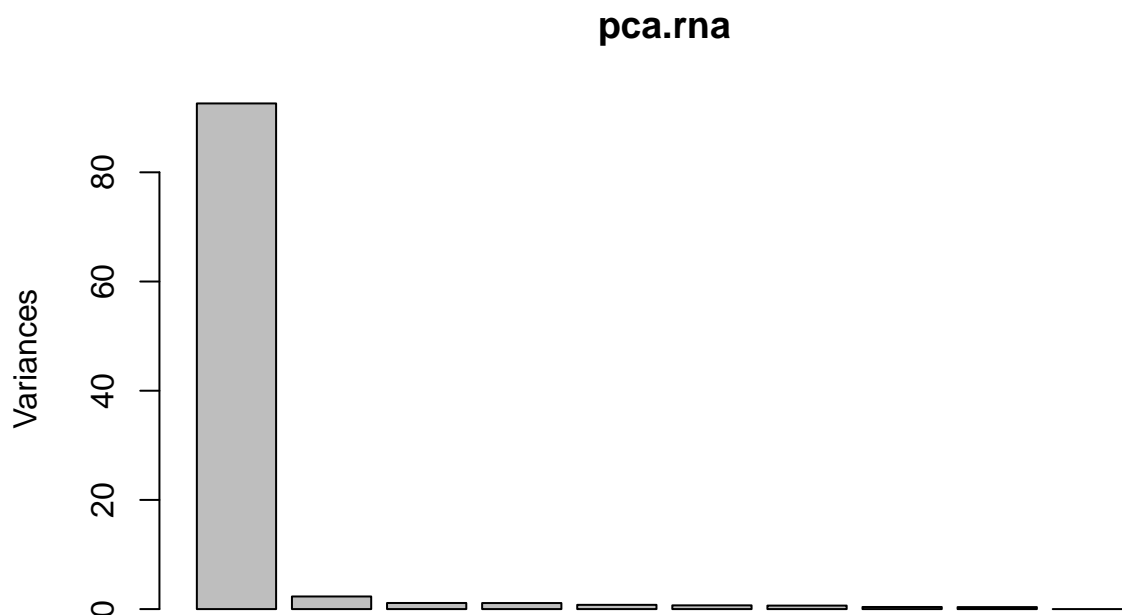
```
pca.rna <-prcomp(t(rna.data), scale = TRUE)
```

```
summary(pca.rna)
```

```
## 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 3.348e-15
## Proportion of Variance 0.00385 0.00364 0.000e+00
## Cumulative Proportion 0.99636 1.00000 1.000e+00
```

```
plot(pca.rna)
```



```
plot(pca.rna$x[,1], pca.rna$x[,2], xlab="PC1", ylab="PC2")
text(pca.rna$x[,1:2], labels = colnames(rna.data))
```

