

# Class 7: Machine Learning (Part I)

Mirte Ciz Marieke Kuijpers

09/02/2022

## R norm function

```
# Investigate the rnorm() function using the help pages.  
?rnorm
```

```
## starting httpd help server ... done
```

```
# Test the rnorm() function  
rnorm(10)
```

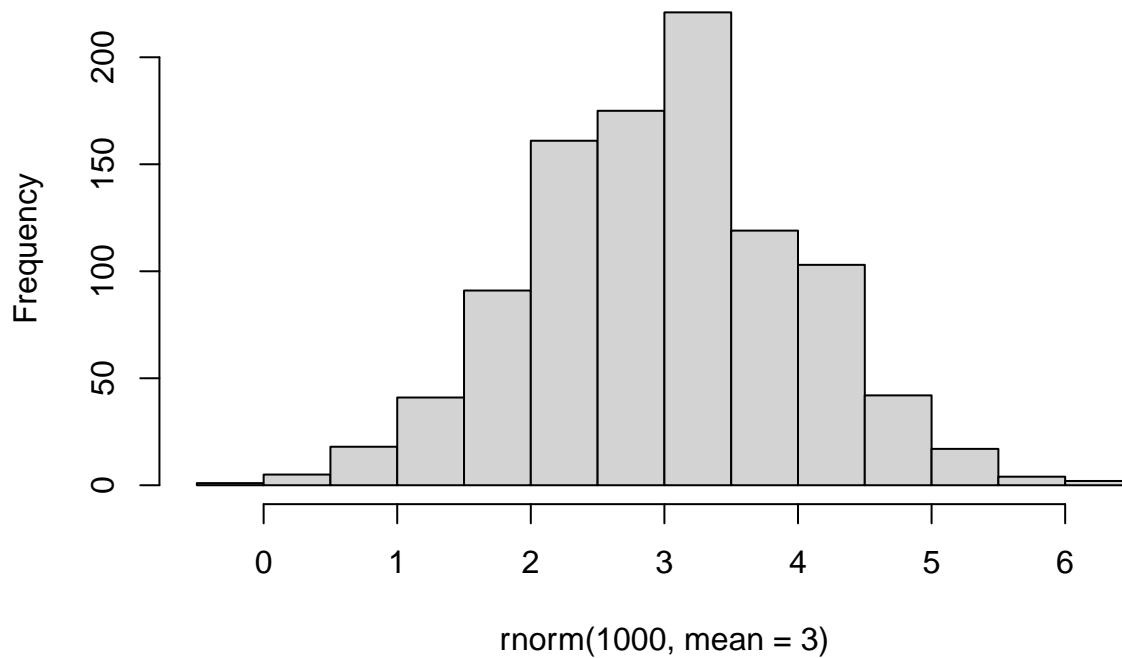
```
## [1] 0.40998733 -0.86522223 -0.30596760 -1.53294512 -0.05017689 -0.26455896  
## [7] -0.41560682 1.34609340 -1.18520498 -0.97632872
```

```
# Use some of the arguments with a default set value  
rnorm(100, mean = 3)
```

```
## [1] 3.9300633 1.3825052 3.5360795 2.3672276 2.0269070 1.8440107 2.9834218  
## [8] 2.6847882 1.8337030 4.5736426 2.6595380 2.7146259 3.3484954 3.1696855  
## [15] 3.4177116 4.2012535 1.4712388 4.0495494 2.1450748 3.4129757 3.4691338  
## [22] 3.3876914 4.3869486 2.9865578 3.1283023 1.7611415 2.8496135 1.2062842  
## [29] 1.4757827 2.1757551 3.0104139 2.1813954 3.6211891 2.2407487 3.1661245  
## [36] 1.5143296 2.9422480 2.2812253 3.4233981 3.3445857 2.5015398 3.2667491  
## [43] 2.5542825 0.8746414 4.1969366 1.5000997 2.7538626 2.1020798 3.3365488  
## [50] 2.9289114 0.7521273 2.9324975 3.1976074 4.1712349 2.0214378 2.7247771  
## [57] 1.9415595 2.6043818 2.5603758 2.7090364 2.0836651 3.0634716 4.0315803  
## [64] 3.4608436 3.5686594 3.6204355 4.6079528 4.6209572 3.4025817 2.4039090  
## [71] 3.3310361 2.8227590 2.5091150 3.1593026 3.8757895 2.0734027 4.3009078  
## [78] 2.4028432 5.6273582 4.1217195 1.2461768 2.0007498 1.7146158 0.9395673  
## [85] 3.3317002 2.5542222 2.4892682 1.3710071 4.1075185 3.6025068 3.5929752  
## [92] 1.1239289 0.6336703 2.4964514 2.7387993 3.4672177 3.3440138 2.5407716  
## [99] 2.5825720 1.2758524
```

```
# Plot some of these results  
hist(rnorm(1000, mean = 3))
```

## Histogram of rnorm(1000, mean = 3)



## Clustering methods

### Using the kmeans() function

Now that we understand `rnorm()`'s function, we can proceed to the example usage of `kmeans()`. The first step is to generate some example data to test the `kmeans()` function with.

```
# Generate some random data
```

```
tmp <- c(rnorm(30, -3), rnorm(30, 3))  
tmp
```

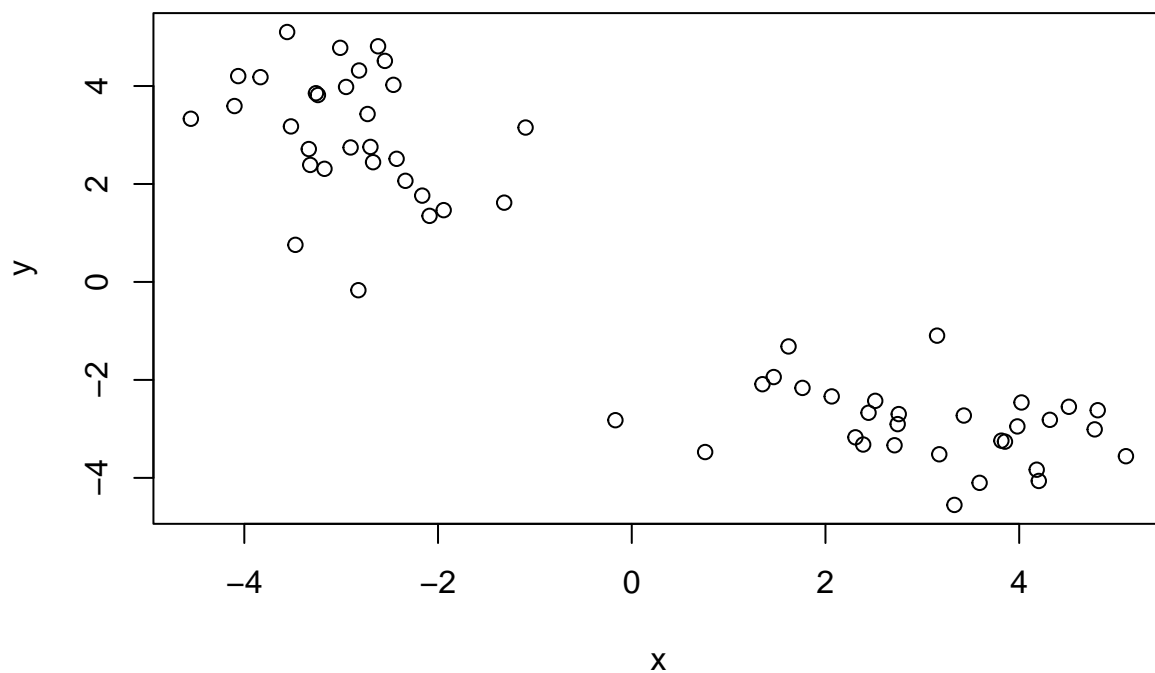
```
## [1] -2.6706498 -2.4275924 -3.5578782 -2.4606483 -1.9417082 -2.9025700  
## [7] -3.2604283 -3.8332275 -2.6979291 -2.8152201 -4.0630301 -4.1014854  
## [13] -2.7275945 -2.0875389 -1.3169864 -1.0962094 -2.5488222 -3.4726011  
## [19] -3.0099445 -2.9493427 -3.5184489 -3.3342864 -2.3369018 -3.1723572  
## [25] -2.6192688 -2.1631409 -3.3193027 -4.5520895 -3.2399989 -2.8226060  
## [31] -0.1687462 3.8161921 3.3313461 2.3892602 1.7626774 4.8110286  
## [37] 2.3093777 2.0644554 2.7130256 3.1750227 3.9811772 4.7795474  
## [43] 0.7576097 4.5129671 3.1522500 1.6186022 1.3494991 3.4284150  
## [49] 3.5912022 4.2024719 4.3166871 2.7563770 4.1809541 3.8530796  
## [55] 2.7456076 1.4656686 4.0237895 5.1022900 2.5137619 2.4446214
```

```
# Format the random data to be used by kmeans()
x <- cbind(x=tmp, y=rev(tmp))
x
```

```
##           x           y
## [1,] -2.6706498  2.4446214
## [2,] -2.4275924  2.5137619
## [3,] -3.5578782  5.1022900
## [4,] -2.4606483  4.0237895
## [5,] -1.9417082  1.4656686
## [6,] -2.9025700  2.7456076
## [7,] -3.2604283  3.8530796
## [8,] -3.8332275  4.1809541
## [9,] -2.6979291  2.7563770
## [10,] -2.8152201  4.3166871
## [11,] -4.0630301  4.2024719
## [12,] -4.1014854  3.5912022
## [13,] -2.7275945  3.4284150
## [14,] -2.0875389  1.3494991
## [15,] -1.3169864  1.6186022
## [16,] -1.0962094  3.1522500
## [17,] -2.5488222  4.5129671
## [18,] -3.4726011  0.7576097
## [19,] -3.0099445  4.7795474
## [20,] -2.9493427  3.9811772
## [21,] -3.5184489  3.1750227
## [22,] -3.3342864  2.7130256
## [23,] -2.3369018  2.0644554
## [24,] -3.1723572  2.3093777
## [25,] -2.6192688  4.8110286
## [26,] -2.1631409  1.7626774
## [27,] -3.3193027  2.3892602
## [28,] -4.5520895  3.3313461
## [29,] -3.2399989  3.8161921
## [30,] -2.8226060 -0.1687462
## [31,] -0.1687462 -2.8226060
## [32,]  3.8161921 -3.2399989
## [33,]  3.3313461 -4.5520895
## [34,]  2.3892602 -3.3193027
## [35,]  1.7626774 -2.1631409
## [36,]  4.8110286 -2.6192688
## [37,]  2.3093777 -3.1723572
## [38,]  2.0644554 -2.3369018
## [39,]  2.7130256 -3.3342864
## [40,]  3.1750227 -3.5184489
## [41,]  3.9811772 -2.9493427
## [42,]  4.7795474 -3.0099445
## [43,]  0.7576097 -3.4726011
## [44,]  4.5129671 -2.5488222
## [45,]  3.1522500 -1.0962094
## [46,]  1.6186022 -1.3169864
## [47,]  1.3494991 -2.0875389
## [48,]  3.4284150 -2.7275945
```

```
## [49,] 3.5912022 -4.1014854
## [50,] 4.2024719 -4.0630301
## [51,] 4.3166871 -2.8152201
## [52,] 2.7563770 -2.6979291
## [53,] 4.1809541 -3.8332275
## [54,] 3.8530796 -3.2604283
## [55,] 2.7456076 -2.9025700
## [56,] 1.4656686 -1.9417082
## [57,] 4.0237895 -2.4606483
## [58,] 5.1022900 -3.5578782
## [59,] 2.5137619 -2.4275924
## [60,] 2.4446214 -2.6706498
```

```
# Plot data to get a feel for the distribution
plot(x)
```

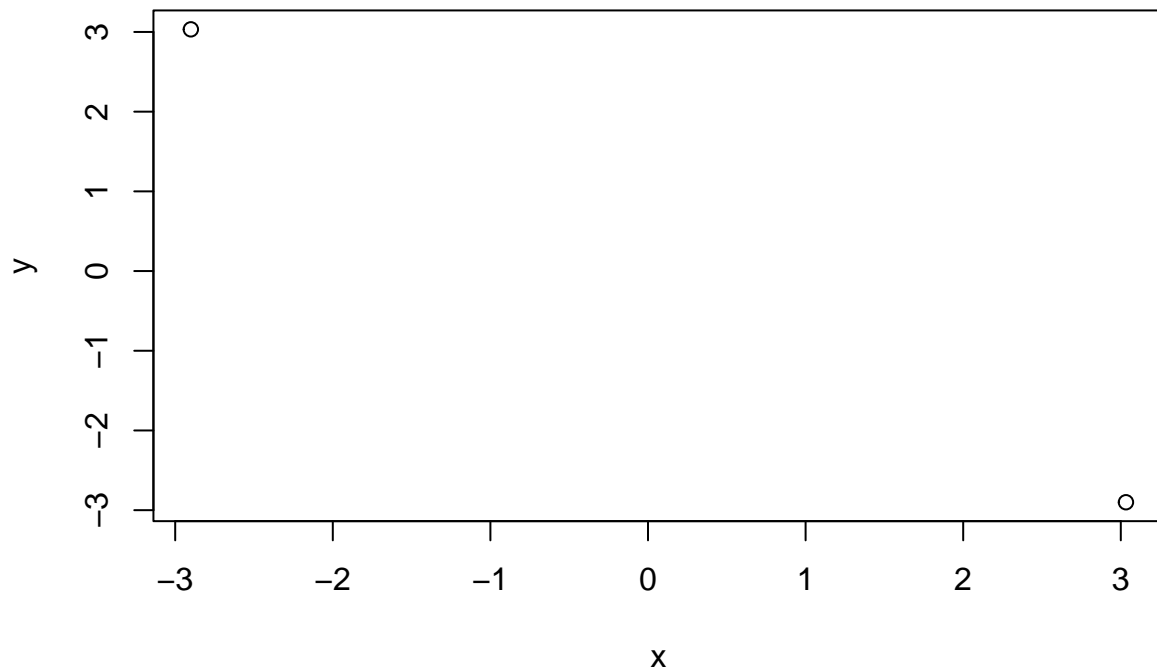


Note that the data above is randomly generated each time the code chunk is run. Therefore, the plot and all answers relying on this data will vary slightly each time the code is re-initialised. Now we can use the `kmeans()` function to cluster this random data.

```
# Cluster the data using the kmeans method
k <- kmeans(x, centers = 2, nstart = 10) # nstarts represents number of new starts (i.e. iterations) to
k
```

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





The membership vector, assigning each point to a cluster, is perhaps the most important results, as it allows you to analyse how your points are clustered.

```
# Print the cluster vector
k$cluster
```

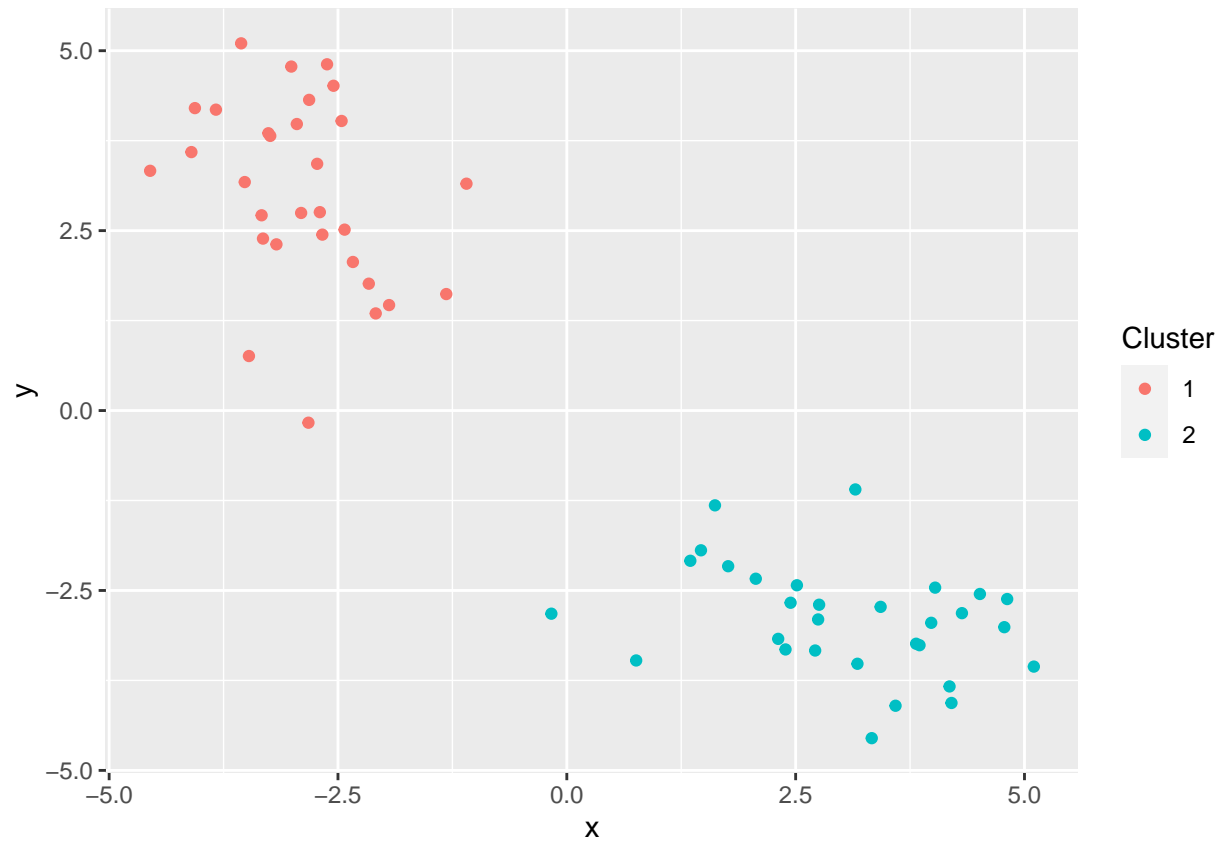
[illegible]

```
# Plot the points coloured by these clusters

# Load ggplot2 library
library(ggplot2)

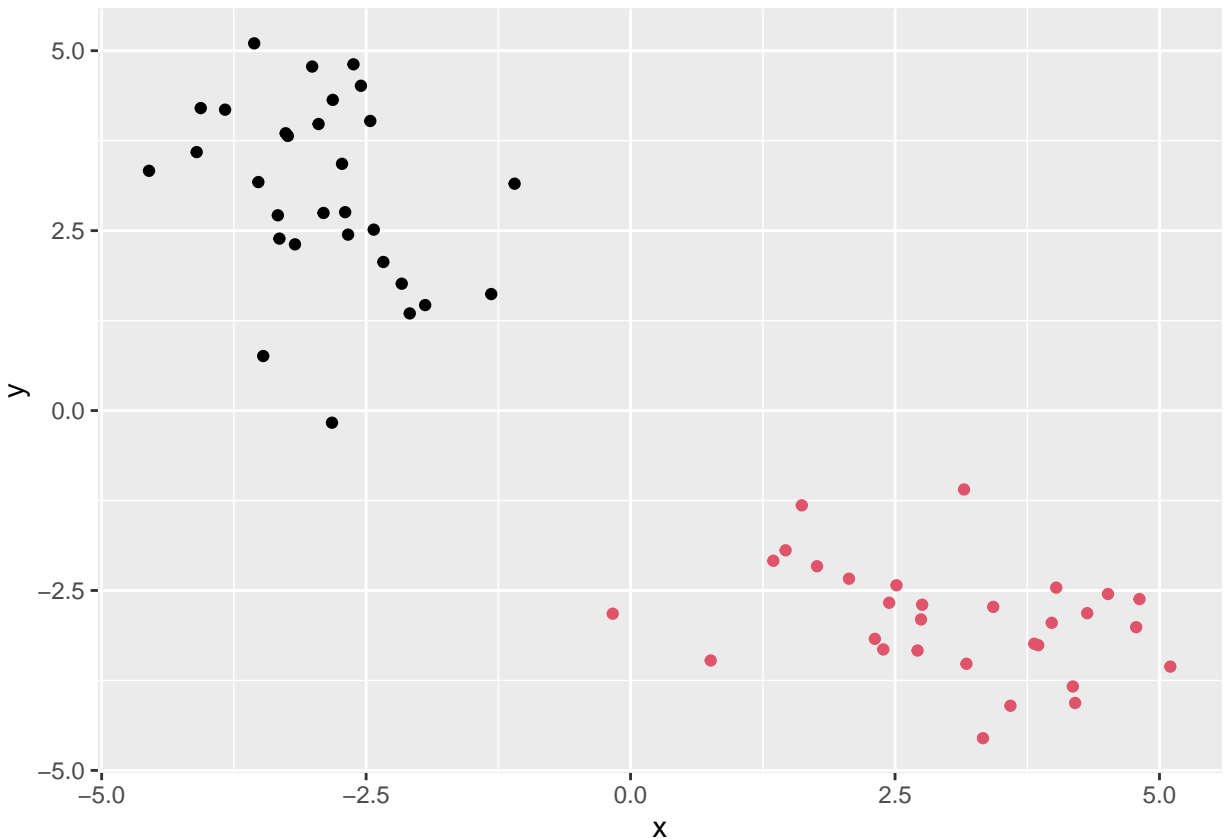
# Put data into a format ggplot can use
dat <- as.data.frame(cbind(x, k$cluster))
colnames(dat) <- c("x", "y", "cluster")

# Plot the data
ggplot(dat, aes(x, y, colour = as.factor(cluster))) +
  geom_point() +
  labs(colour = "Cluster")
```



*# In actual fact, one does not need to consolidate the data, one can simply use `k$cluster` outside the `aes` function*

```
ggplot(as.data.frame(x), aes(x,y)) +  
  geom_point(col = k$cluster)
```



## hclust() function

kmeans() uses Euclidean distance, hclust can use various similarity/difference measures, which is useful, but requires extra input. hclust() requires data already as distances between points (a dissimilarity/similarity matrix), not the raw data. One way to do so is to use the dist() function.

```
# Use the hclust() function on the same data (x) used for investigating the kmeans() function
hc <- hclust(dist(x))
hc
```

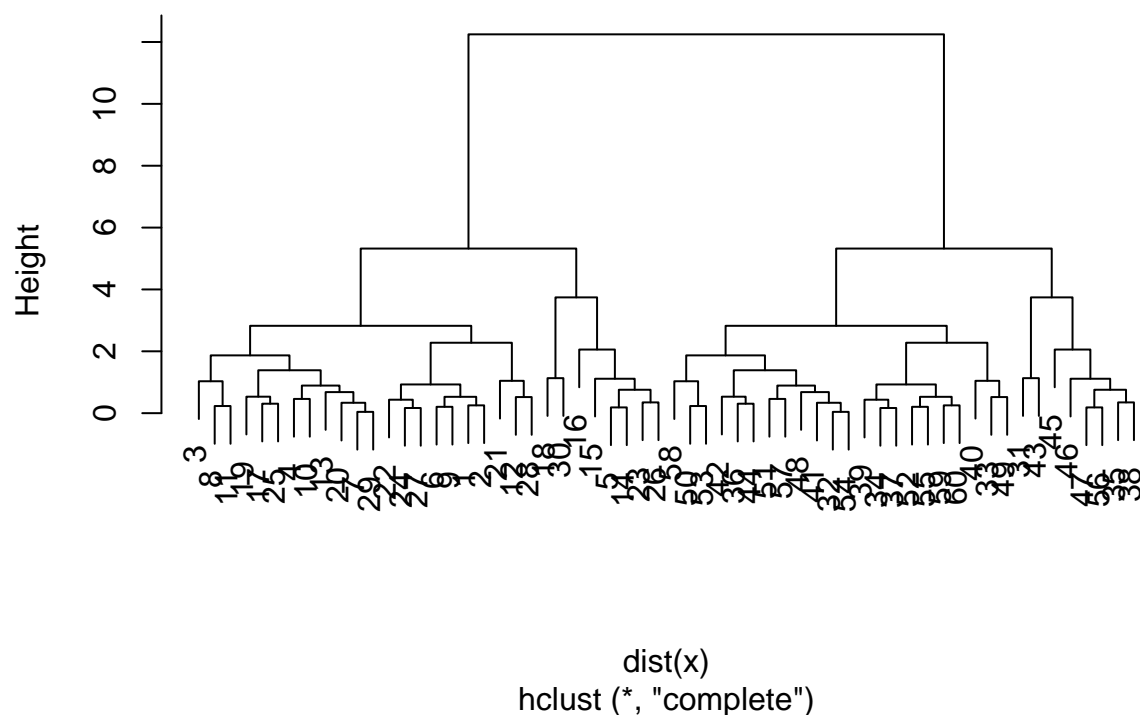
```
##
## Call:
## hclust(d = dist(x))
##
## Cluster method   : complete
## Distance         : euclidean
## Number of objects: 60
```

While the output for hclust(), when printed, is not as useful as for kmeans(), it has a custom plot method which helps interpret the data. In the dendrogram produced, the height of the crossbars represents (or is proportional to) the distance between the two groups joined by said crossbar.

```
# Plot hclust() data, as printing the output is not as useful as for kmeans()
plot(hc)
```



## Cluster Dendrogram

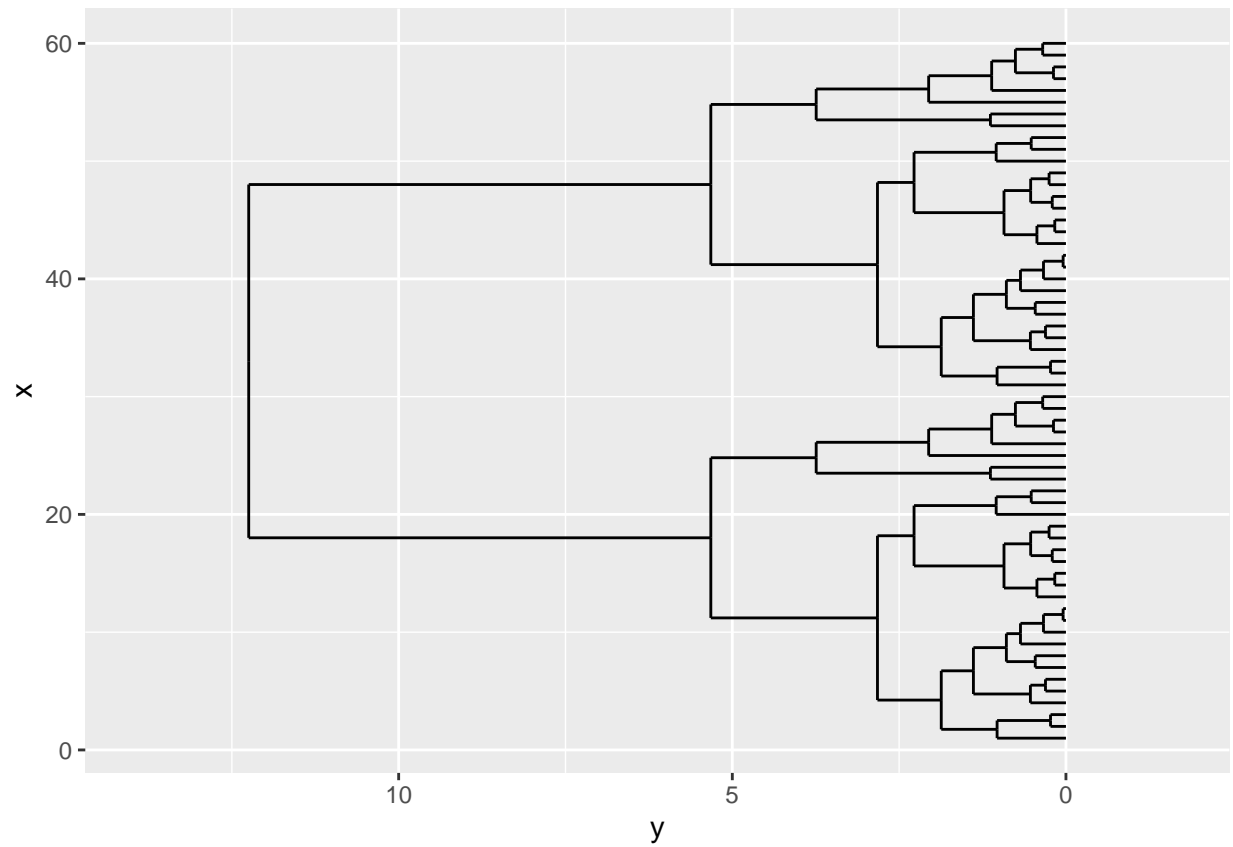


*# One can also plot with ggplot, if a wrapper for dendrograms is also loaded:*  
library(ggdendro)

*# Convert the data into a format ggplot can use - code borrowed from the following tutorial <https://cran.r-project.org/web/packages/ggdendro/vignettes/ggdendro.html>*  
ghc <- as.dendrogram(hc)

ghc.dat <- dendro\_data(ghc, type = "rectangle")

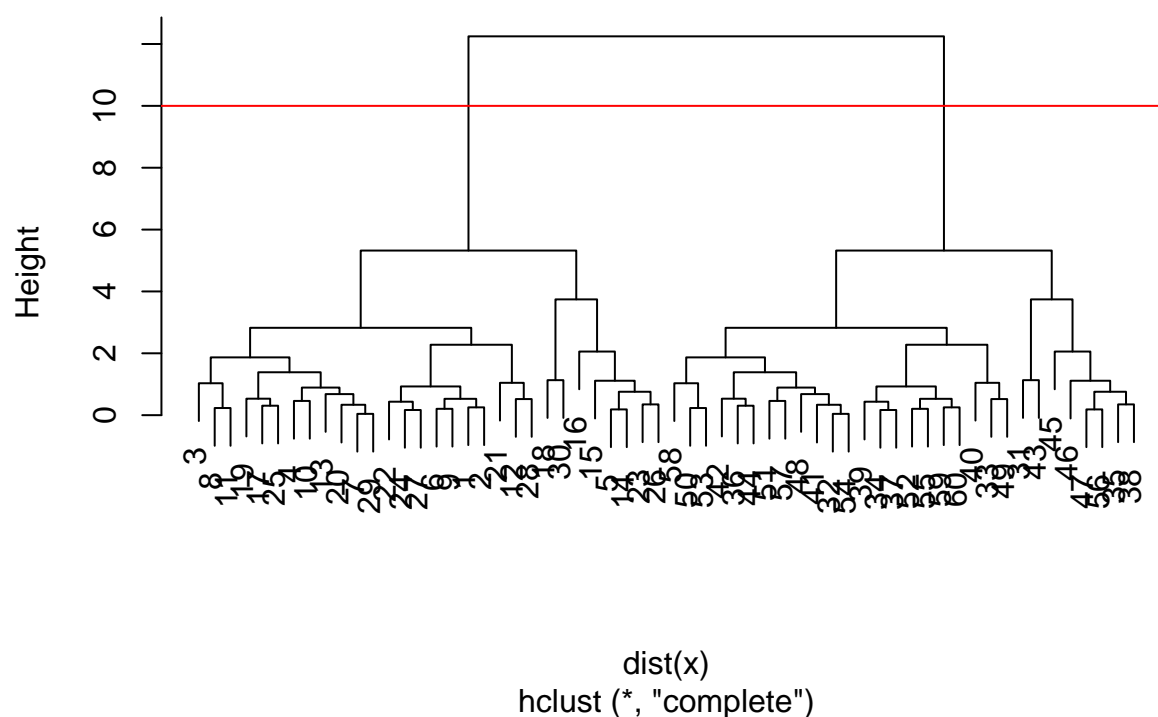
```
ggplot(segment(ghc.dat)) +  
  geom_segment(aes(x = x, y = y, xend = xend, yend = yend)) +  
  coord_flip() +  
  scale_y_reverse(expand = c(0.2, 0))
```



To determine groups, one can set a horizontal cut-off line, which separates points connected below that cut-off into clusters.

```
plot(hc)
abline(h = 10, col = "Red")
```

## Cluster Dendrogram

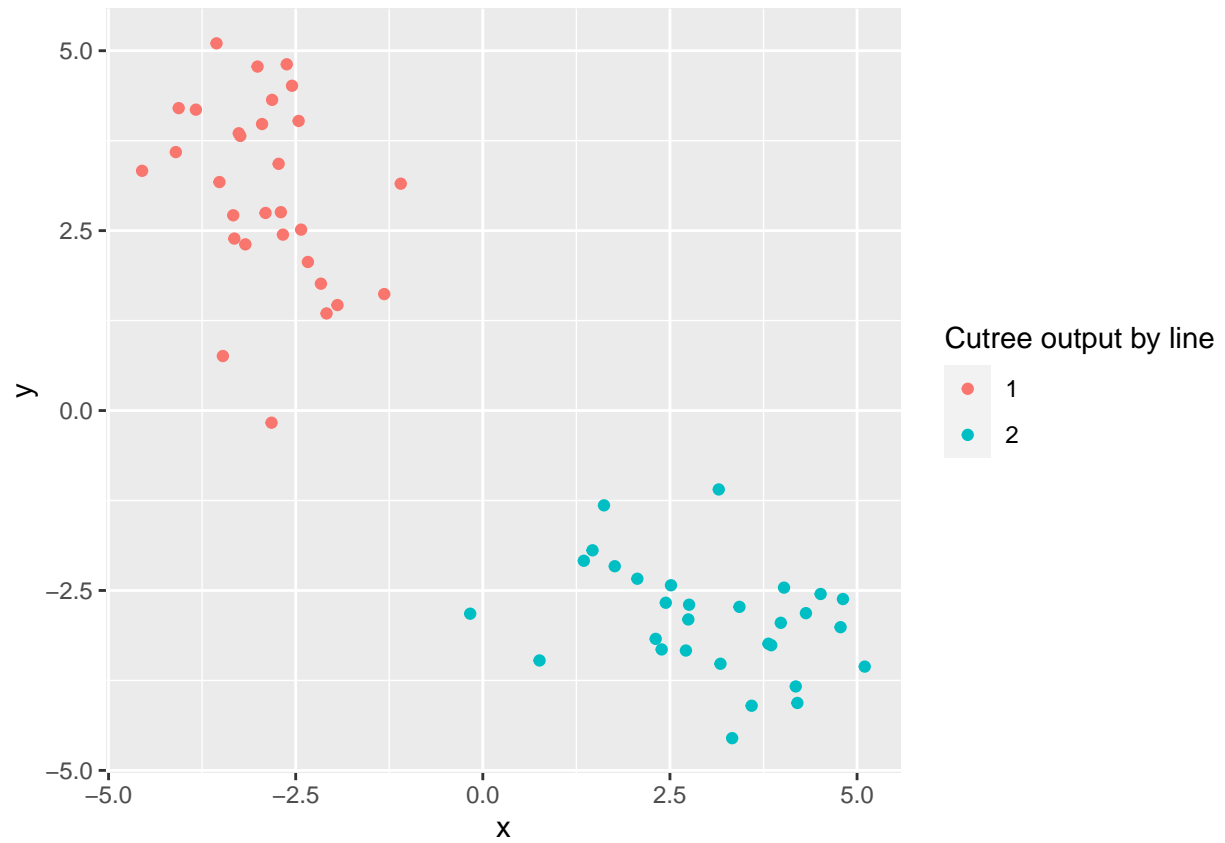


Instead of adding this cut-off line, we can use `cutree()`, a function specifically for this.

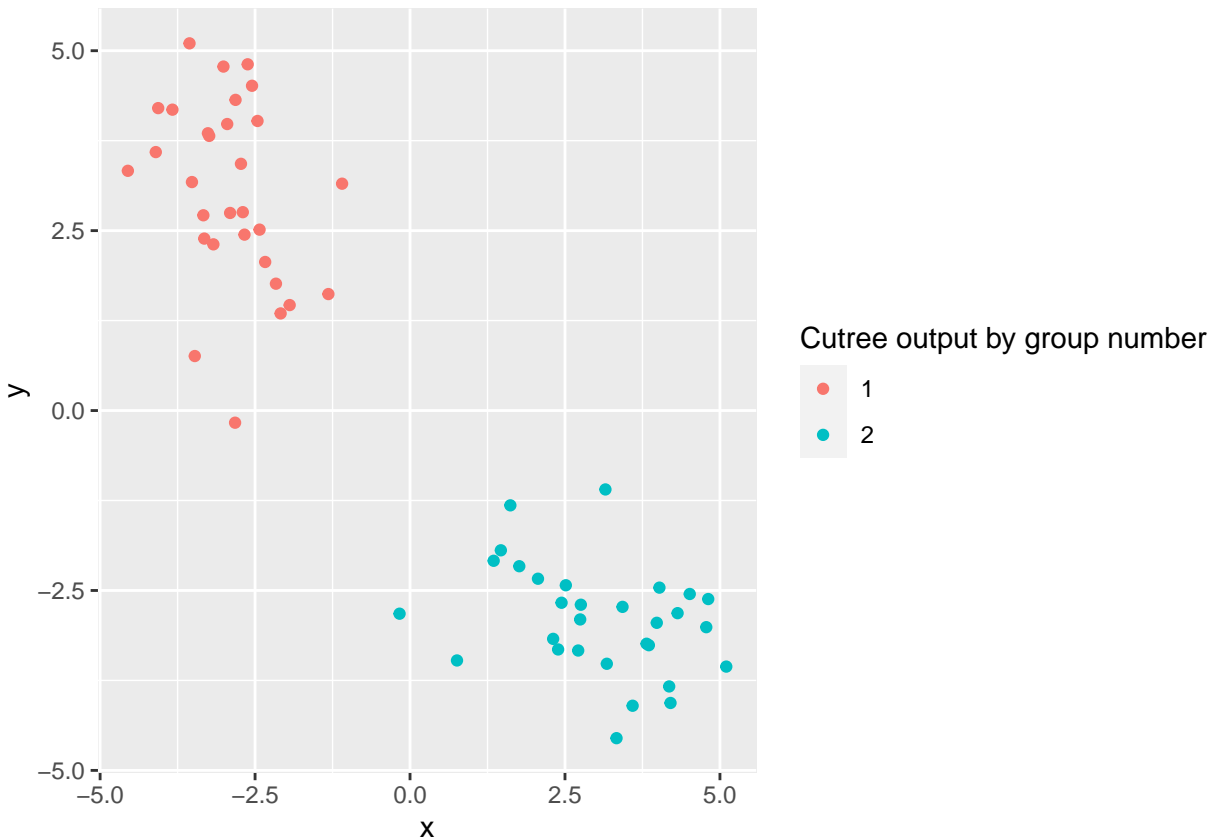
```
cth <- cutree(hc, h=10)
#or specify the number of groups
ctk <- cutree(hc, k=2)

# Put data into a format ggplot can use
dat <- as.data.frame(cbind(x, cth, ctk))
colnames(dat) <- c("x", "y", "cth", "ctk")

# Plot the data
ggplot(dat, aes(x, y, colour = as.factor(cth))) +
  geom_point() +
  labs(colour = "Cutree output by line")
```



```
ggplot(dat, aes(x, y, colour = as.factor(ctk))) +  
  geom_point() +  
  labs(colour = "Cutree output by group number")
```



## Principle Component Analysis

For this we are going to use data on the food consumption of citizens of different countries within the United Kingdom. First download the data:

```
# Download and assign data to an r object
url <- "https://tinyurl.com/UK-foods"
f.dat <- read.csv(url)
```

```
# Inspect the data
str(f.dat)
```

```
## 'data.frame': 17 obs. of 5 variables:
## $ X : chr "Cheese" "Carcass_meat " "Other_meat " "Fish" ...
## $ England : int 105 245 685 147 193 156 720 253 488 198 ...
## $ Wales : int 103 227 803 160 235 175 874 265 570 203 ...
## $ Scotland : int 103 242 750 122 184 147 566 171 418 220 ...
## $ N.Ireland: int 66 267 586 93 209 139 1033 143 355 187 ...
```

```
# The first column would be better set as rownames than its own column
rownames(f.dat) <- f.dat[,1]
f.dat <- f.dat[,-1]
```

```
# Check this has worked
head(f.dat)
```

```
##              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(f.dat)
```

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

```
#N.B. a better way to have dealt with this problem is to set rownames = 1 in the original read.csv
food <- read.csv(url, row.names = 1)
str(food)
```

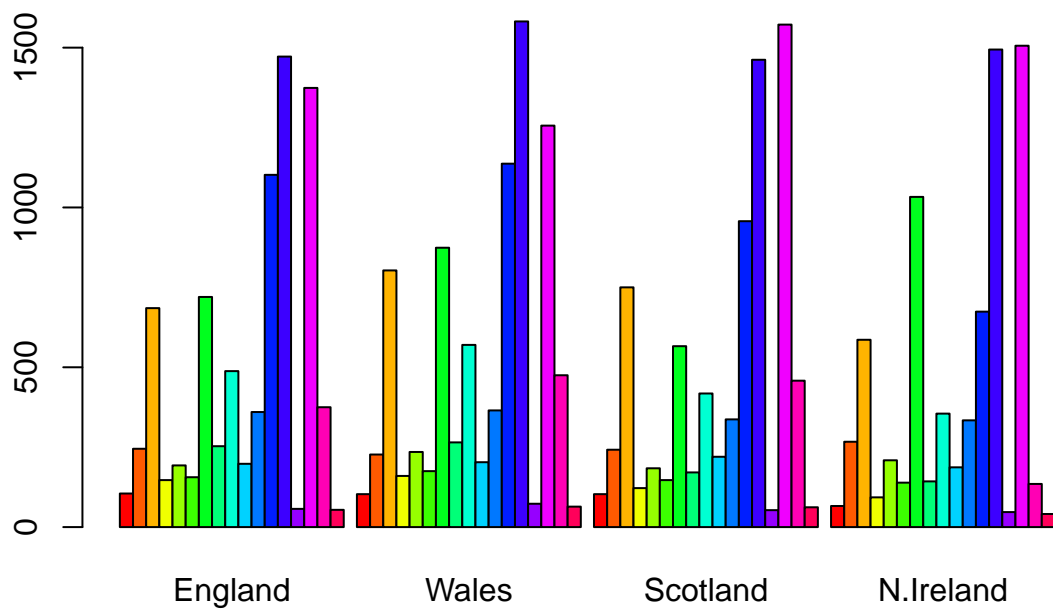
```
## 'data.frame':  17 obs. of  4 variables:
## $ England : int  105 245 685 147 193 156 720 253 488 198 ...
## $ Wales   : int  103 227 803 160 235 175 874 265 570 203 ...
## $ Scotland: int  103 242 750 122 184 147 566 171 418 220 ...
## $ N.Ireland: int  66 267 586 93 209 139 1033 143 355 187 ...
```

```
head(food)
```

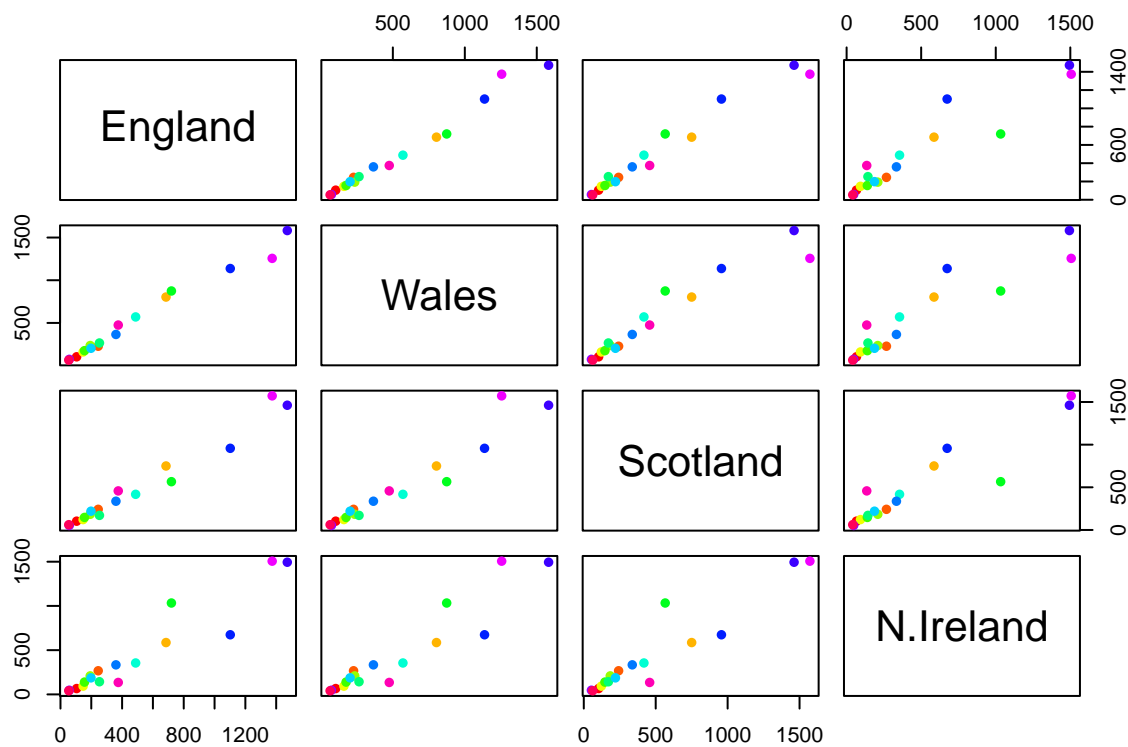
```
##              England Wales Scotland N.Ireland
## Cheese      105    103      103        66
## Carcass_meat 245    227      242       267
## Other_meat   685    803      750       586
## Fish         147    160      122        93
## Fats_and_oils 193    235      184       209
## Sugars       156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

The data (food), is now in an appropriate form for analysis.

```
# Plot the data
barplot(as.matrix(food), beside = TRUE, col = rainbow(nrow(food)))
```



```
pairs(food, col=rainbow(nrow(food)), pch=16)
```



Analysis through plotting is not very helpful, thus we move onto PCA. We will begin with the basic PCA from base r, packages with other PCA functions exist, but they are often specialized PCA functions for specific data types or circumstances. The base r PCA function is `prcomp()`.

```
# prcomp() expects the transpose of the way our data currently is, with the columns as the categories w
pca <- prcomp(t(food))

# Printing pca gives the PC values for each of our food categories, it is thus easier to look at a summ
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
```

```
attributes(pca)
```

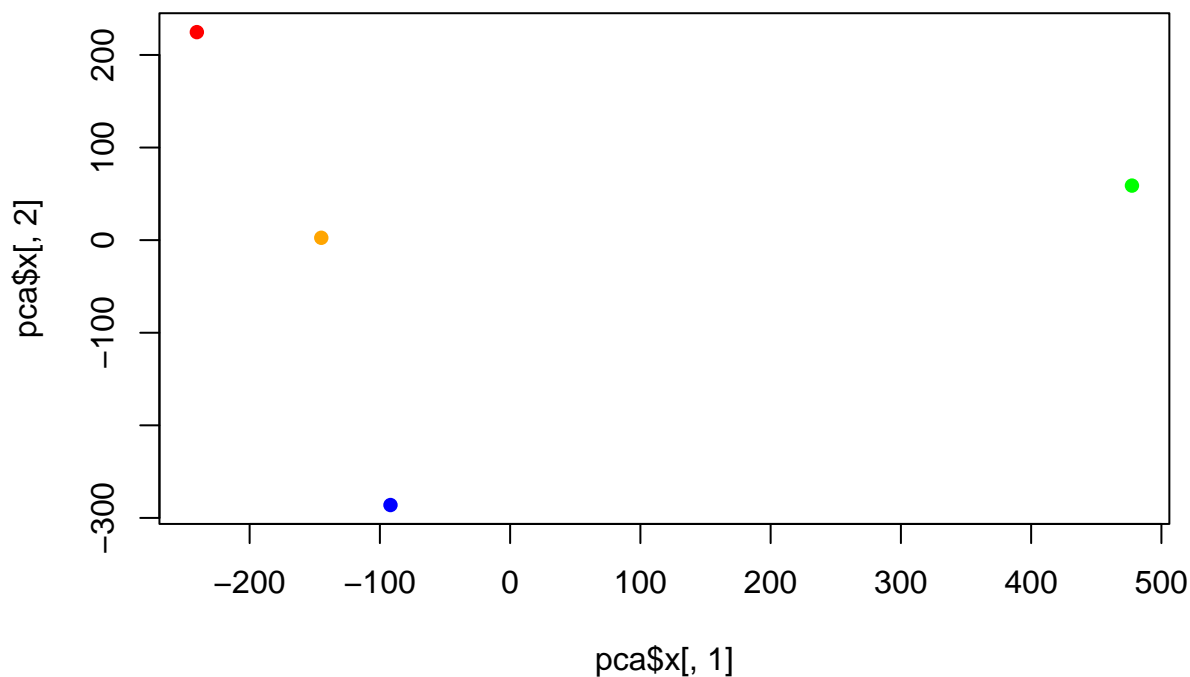
```
## $names
## [1] "sdev"      "rotation" "center"    "scale"     "x"
##
## $class
## [1] "prcomp"
```

We can now view a plot of the data along the two most important axes as found through PCA. A plot of

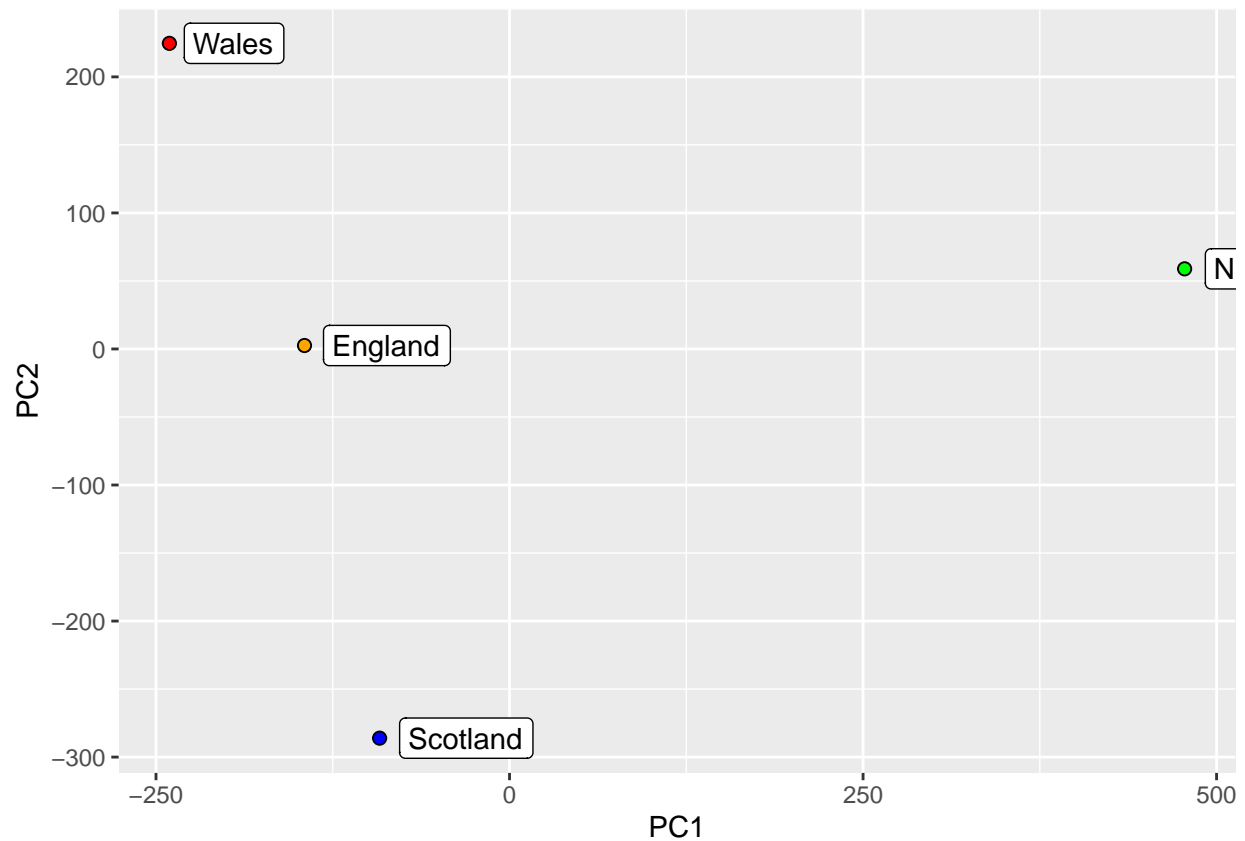


PC1 vs PC2 is often called a PCA plot or “score plot”. The x component of the `prcomp()` output gives the data for such a score plot.

```
plot(pca$x[,1], pca$x[,2], col = c("Orange", "Red", "Blue", "Green"), pch = 16) #Note order is England,
```



```
##Need to try plot with ggplot...  
p <- as.data.frame(pca$x)  
  
ggplot(p, aes(PC1, PC2)) +  
  geom_point(fill = c("Orange", "Red", "Blue", "Green"), pch = 21, size = 2) +  
  geom_label(label = c("England", "Wales", "Scotland", "N.Ireland"), hjust = -0.15)
```



The four points are the four countries. Note the the large distance between one point (N.Ireland) versus the other three (England, Wales and Scotland) on the PC1 axis, which explains the most variation. Note that another important point is to show how much these axes actually contribute to the variation (the loadings), these are found in the component rotation. As PC1 explains most of the variation, we will focus on this.

```
par(mar=c(10, 3, 0.35, 0)) #change positioning of the plot
barplot(pca$rotation[,1], las=2)
```

