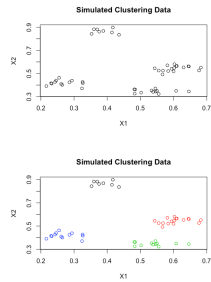


Example Data



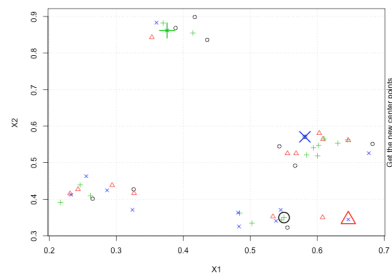
	X1	X2	True Cluster
1	0.4356460	0.839740	1
2	0.3539190	0.8451142	1
3	0.3757543	0.8616696	1
4	0.3699114	0.8822238	1
5	0.3588940	0.8830917	1
6	0.4141158	0.8501112	1
7	0.3884225	0.8661076	1
8	0.4170783	0.8983235	1
9	0.5942671	0.5407596	2
10	0.6085908	0.5644835	2
11	0.5834096	0.5707407	2
12	0.5433074	0.5447830	2
13	0.6018648	0.5471825	2
14	0.6824566	0.5313103	2
15	0.5968619	0.4933855	2
16	0.6099432	0.5650591	2
17	0.6028487	0.5806603	2
18	0.6789366	0.5259041	2
19	0.5939923	0.5257776	2
20	0.6463144	0.5608932	2
21	0.5682214	0.5251682	2
22	0.6451486	0.5623831	2
23	0.5541661	0.5253996	2
24	0.6003407	0.5188467	2
25	0.6304850	0.529242	2

Well-Known Methods

• K-means Clustering

1. Specify the number of clusters (K).
2. Define what is similar or dissimilar
→ Distance!
3. Algorithm
 1. Random assign all the points to a unique cluster (k), where $k = 1, 2, 3, \dots, K$.
 2. Calculate the mean center point for each cluster (m_1, m_2, \dots, m_K).
 3. Reassign the point to the cluster that it has closest distance to those mean center points.
 4. Repeat 2. ~ 3., until the assignment stop changing.

K-means Animation!



Definition for (Dis)similarity

- Suppose we have a measured point, x_{ij} , where

Points: $i = 1, 2, 3, \dots, N$

Attributes: $j = 1, 2, 3, \dots, p$

- Dissimilarity between point i and i'

$$D(x_i, x_{i'}) = \sum_{j=1}^p d_j(x_{ij}, x_{i'j})$$

- Euclidean Distance

$$\begin{matrix} x_i (x_{i1}, x_{i2}) \\ (x_{i1} - x_{i'1})^2 = b \\ (x_{i2} - x_{i'2})^2 = a \\ a^2 + b^2 \end{matrix}$$

$$d_{ij}(x_{ij}, x_{i'j}) = (x_{ij} - x_{i'j})^2$$

$$d(x_i, x_{i'}) = \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = \|x_i - x_{i'}\|^2$$

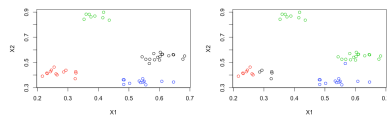
K-means
Clustering

• Things Need To Be Noticed

1. You won't always get the same results. So, do it more times!

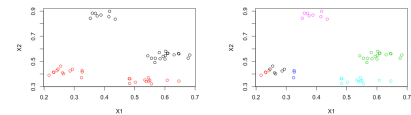
```
#### Things need to be noticed ####
set.seed(2223) # good case
k = kmeans(sim_d[,1:2], centers=4, nstart = 1)
plot(sim_d[,1:2], col = k$cluster, main = "Good Case")

set.seed(14) # bad case
k = kmeans(sim_d[,1:2], centers=4, nstart = 1)
plot(sim_d[,1:2], col = k$cluster, main = "Bad Case")
```

K-means
Clustering

• Things Need To Be Noticed

2. How to choose the correct number of clusters ?



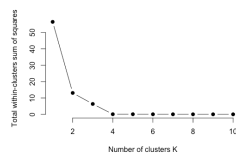
$$\sum (x_i - \bar{x})^2 = \sum (x_i - \bar{x}_k)^2 + \sum (\bar{x}_k - \bar{x})^2$$

TOTAL Sum of Squares = WITHIN Groups SS + BETWEEN Groups SS
(TOTAL variance = WITHIN variance + BETWEEN variance)

K-means
Clustering

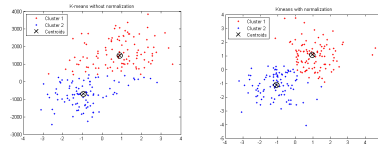
• An Intuitive Method for finding the best k

```
# Find the best number of cluster
k.max <- 10
wss <- sapply(1:k.max, function(k){
  kmeans(sim_d, k, nstart=30, iter.max = 15)$tot.withinss
})
plot(1:k.max, wss, type="b", pch = 19, frame = FALSE, xlab="Number of clusters K",
     ylab="Total within-clusters sum of squares")
```

K-means
Clustering

• Things Need To Be Noticed

3. Standardize the data or not?



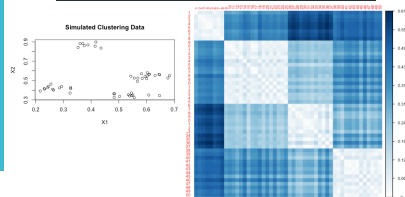
The source of pics : [here](#)

Hierarchical Clustering

• Dissimilarity Matrix

```
#### Dissimilarity Matrix ####
library("corplot")

dst = dist(sim_d[,1:2], method = "euclidean")
corplot(as.matrix(dst), is.corr = FALSE, method = "color",
        tl.cex = 0.5, cl.cex = 0.5, mar = c(0.5, 0.5, 0.5, 0.5))
```

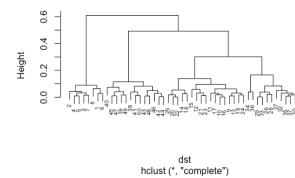


Hierarchical Clustering

• Algorithm

```
## Hierarchical clustering
dst = dist(sim_d[,1:2], method = "euclidean")
hc = hclust(dst)
plot(hc, cex = 0.6)
rect.hclust(hc, k = 4, border = 2:5)
```

Cluster Dendrogram



Hierarchical Clustering

• Find the Best k

```
#### Find the best cluster ####
library(dtw)
dst = dist(sim_d[,1:2], method = "euclidean")
hc = hclust(dst)

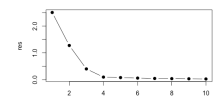
wss <- function(d) {
  sum(scale(d, scale = FALSE)^2)
}
wrap = function(l, hc, k) {
  cl = cutree(hc, l)
  spl = split(x, cl)
  wss = sum(sapply(spl, wss))
  wss
}

res = sapply(seq.int(1, 10), wrap, h = hc, x = sim_d[,1:2])
plot(seq_along(res), res, type = "b", pch = 19)

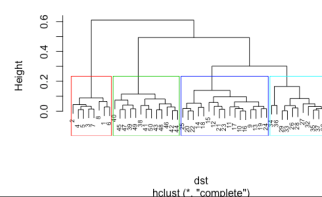
k = 4
plot(hc, cex = 0.6)
rect.hclust(hc, k, border = 2:5)
cutree(hc, k)
```

Hierarchical Clustering

• Find the Best k



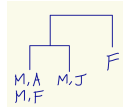
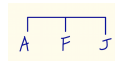
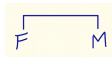
Cluster Dendrogram



Hierarchical Clustering

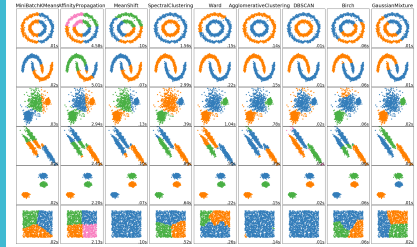
Things Need To Know

1. No need to specify the number of clusters. One Dendrogram can show many clusters.
2. It assumes that all the clusters are nested with a hierarchical structure, but it may not be the case.
Ex:
x1 = { Male, Female }, x2 = { Americans, Japanese, French }
Best 2 clusters : by Gender
Best 3 clusters : by Nationality



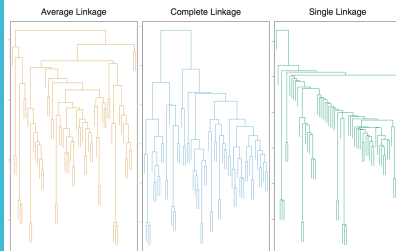
Further Topics

The source of the graph : [here](#)



Further Topics

The source of the graph : [here](#)



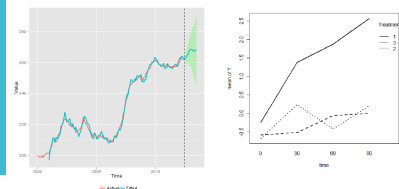
Time Series Data

General Definition

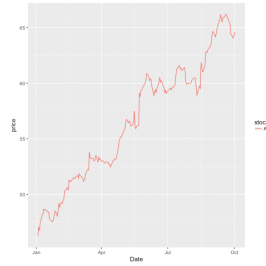
The data is taken by time and believe that time will be an influential factor to the data.

Ex: [Sales data](#)

Difference in temperature under different treatments



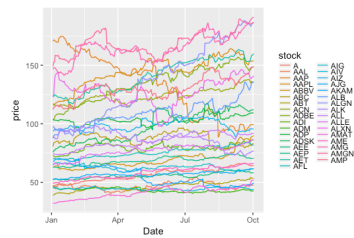
Example Data



Example Data



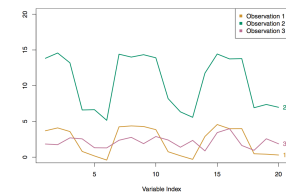
Example Data



Distance
between Time
Series

- Euclidean Distance may not be what we want.

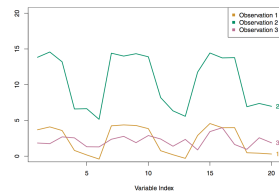
(The source of this pic : [here](#))



Distance between Time Series

- Think about Correlation

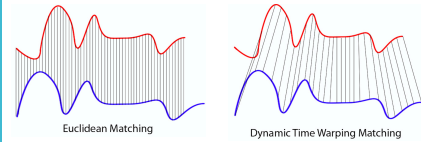
$$\rho(x_i, x_{i'}) = \frac{\sum_j (x_{ij} - \bar{x}_i)(x_{i'j} - \bar{x}_{i'})}{\sqrt{\sum_j (x_{ij} - \bar{x}_i)^2 \sum_j (x_{i'j} - \bar{x}_{i'})^2}}$$



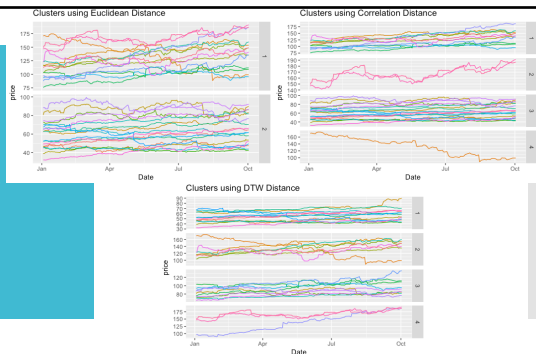
Distance between Time Series

- Think about Dynamic Time Warping (DTW) Distance

(The source of this pic : [here](#))



Results



Applications

- Try to find out the similarity of quantity sold between different products, different Site and Fulfillment Type.
- Difficulties
 - Different length of time series. New vs. Old
 - As time goes by, the amount of data will rapidly increase.
 - Subset the data
 - Extract the features of time series

References and Resources

• Books

- [An Introduction to Statistical Learning with Applications in R](#)
- [The Elements of Statistical Learning: Data Mining, Inference and Prediction](#)



• R Packages

- [MixSim](#) – Generate clustering data
- [corrplot](#) – Plot the dissimilarity matrix
- [tidyverse](#) – Data Manipulation
- [ggplot2](#) – Data Visualization
- [amap](#) – K-means with Correlation Distance

Q and A

• My personal website and Github

- <http://yintingchou.com/>
- https://github.com/choux130/EU_TSclustering_101117

Demo

• Rstudio!