

Cluster Analysis

Jan Graffelman¹

¹Department of Statistics and Operations Research
Universitat Politècnica de Catalunya
Barcelona, Spain



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

jan.graffelman@upc.edu

March 30, 2020

Contents

- 1 Introduction
- 2 Distance measures
- 3 Hierarchical agglomerative clustering
- 4 Non-hierarchical K-means clustering
- 5 Model-based clustering
- 6 Validation

Objectives

Goals:

- Discover "natural" groups of cases (or variables) in the data.
- Data reduction: from n cases to $m \ll n$ clusters

Considerations:

- The number of clusters may a priori be unknown.
- There is no categorical variable that defines the grouping.
- Cluster analysis is an exploratory tool.

Ingredients

- Distance measure
 - In order to cluster item or variables we need a measure of **similarity (proximity)** or **distance (dissimilarity)**.
- Algorithm
 - We cannot consider all possible groupings and need algorithms to produce the grouping.

Algorithms

- Hierarchical methods: cases can not change group
 - agglomerative (most common)
 - divisive
- Partitioning methods: cases can change group
 - K-means
- Model based methods
- Other

Distance measure for quantitative variables

$$\mathbf{x}' = (x_1, x_2, \dots, x_p) \quad \mathbf{y}' = (y_1, y_2, \dots, y_p)$$

Euclidian distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})'(\mathbf{x} - \mathbf{y})}$$

Weighted Euclidian distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})'\mathbf{A}(\mathbf{x} - \mathbf{y})}$$

Mahalanobis distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})'\mathbf{S}^{-1}(\mathbf{x} - \mathbf{y})}$$

Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^p |x_i - y_i|$$

Minkowski distance:

$$d(\mathbf{x}, \mathbf{y}, \lambda) = \left(\sum_{i=1}^p |x_i - y_i|^\lambda \right)^{1/\lambda}$$

Canberra distance:

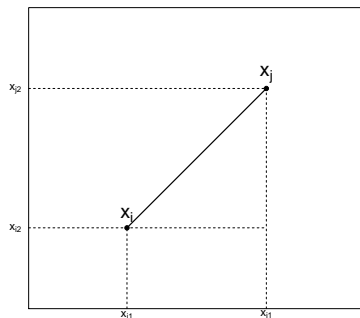
$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i}$$

Bray-Curtis distance:

$$d(\mathbf{x}, \mathbf{y}) = \frac{1}{P} \frac{\sum_{i=1}^p |x_i - y_i|}{\sum_{i=1}^p (x_i + y_i)}$$

...

Euclidean distance



In two dimensions:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{j1} - x_{i1})^2 + (x_{j2} - x_{i2})^2} = \sqrt{(\mathbf{x}_j - \mathbf{x}_i)'(\mathbf{x}_j - \mathbf{x}_i)}$$

In p dimensions:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{j1} - x_{i1})^2 + (x_{j2} - x_{i2})^2 + \cdots + (x_{jp} - x_{ip})^2} = \sqrt{(\mathbf{x}_j - \mathbf{x}_i)'(\mathbf{x}_j - \mathbf{x}_i)}$$

Similarity measures for qualitative variables

		case j		
		1	0	
case i	1	a	b	$a+b$
	0	c	d	$c+d$
		$a+c$	$b+d$	$p=a+b+c+d$

$$\frac{a+d}{p}$$

simple matching coefficient

$$\frac{a}{p}$$

only one-one matches

$$\frac{a}{a+b+c}$$

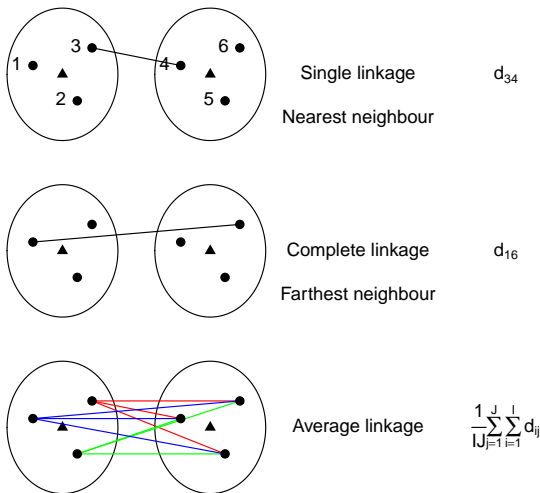
Jaccard's coefficient (0-0 irrelevant)

Example

	indicators					
case 1	1	1	0	0	1	1
case 2	0	1	1	0	0	1

- Compute the squared Euclidean distance between the cases
- What does this distance represent?

Cluster distance



Criteria for joining clusters

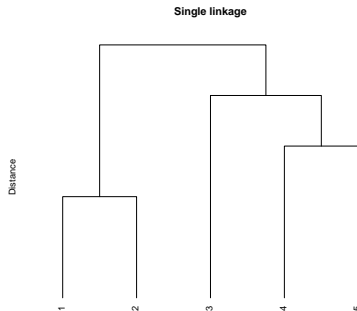
- single linkage
- complete linkage
- average linkage
- centroid distance $d_{rs}^2 = \sum_{j=1}^p (\bar{x}_{rj} - \bar{x}_{sj})^2$
(UPGMA, Unweighted Pair Group Method using Averages)
- Ward's incremental sum-of-squares

Miniature example: hierarchical agglomerative

	1	2	3	4	5
1	0	2	6	10	9
2	2	0	5	9	8
3	6	5	0	4	5
4	10	9	4	0	3
5	9	8	5	3	0

Distance	Clusters
0	1,2,3,4,5
2	(1,2),3,4,5
3	(1,2),3,(4,5)
4	(1,2),(3,4,5)
5	(1,2,3,4,5)

Dendrogram



Miniature example: continuation

D_0	1	2	3	4	5
1	0	2	6	10	9
2	2	0	5	9	8
3	6	5	0	4	5
4	10	9	4	0	3
5	9	8	5	3	0

D_1	(1,2)	3	4	5
(1,2)	0	5	9	8
3	5	0	4	5
4	9	4	0	3
5	8	5	3	0

D_2	(1,2)	3	(4,5)
(1,2)	0	5	8
3	5	0	4
(4,5)	8	4	0

D_3	(1,2)	(3,4,5)
(1,2)	0	5
(3,4,5)	5	0

D_4	(1,2,3,4,5)
(1,2,3,4,5)	0

Exercise

Given the distance matrix

	1	2	3	4	5
1	0	10	27	15	19
2	10	0	18	6	8
3	27	18	0	16	12
4	15	6	16	0	7
5	19	8	12	7	0

Write down the successive formation of cluster according to the complete linkage criterion

Ward's criterion

Given two clusters r and s we have the **within-group sums-of-squares**

$$WSS_r = \sum_{j=1}^p \sum_{i=1}^{n_r} (x_{ij} - \bar{x}_j)^2 \quad WSS_s = \sum_{j=1}^p \sum_{i=1}^{n_s} (x_{ij} - \bar{x}_j)^2$$

On joining, an new cluster t is obtained with a new WSS:

$$WSS_t = \sum_{j=1}^p \sum_{i=1}^{n_r+n_s} (x_{ij} - \bar{x}_j)^2$$

This gives an increase in WSS:

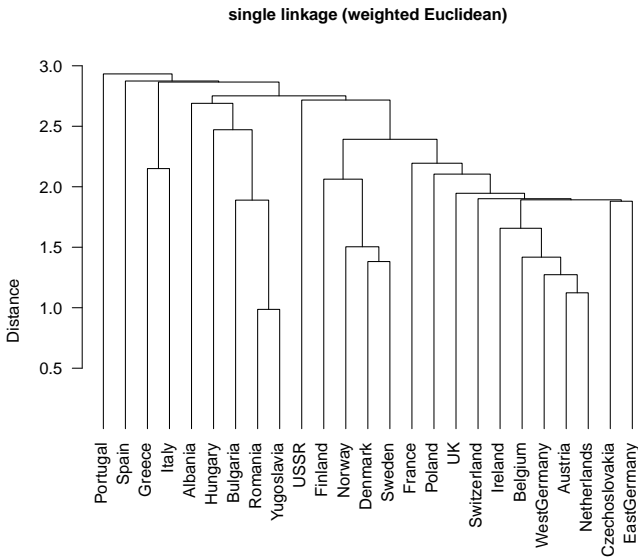
$$\Delta = WSS_t - (WSS_r + WSS_s) = \frac{n_r n_s}{n_r + n_s} d_{rs}^2$$

Join those two clusters for which Δ is minimal.

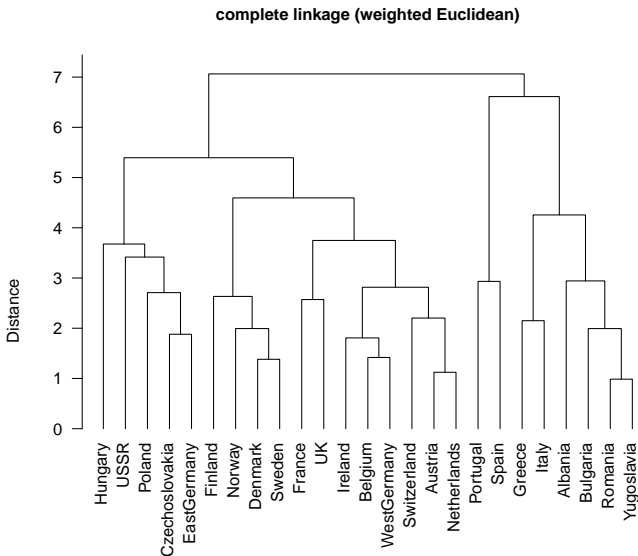
Example: protein consumption data

Country	Red meat	White meat	Eggs	Milk	Fish	Cereals	Starchy foods	Pulses, Nuts Oilseeds	Fruit Vegetables
Albania	10.1	1.4	0.5	8.9	0.2	42.3	0.6	5.5	1.7
Austria	8.9	14.0	4.3	19.9	2.1	28.0	3.6	1.3	4.3
Belgium	13.5	9.3	4.1	17.5	4.5	26.6	5.7	2.1	4.0
Bulgaria	7.8	6.0	1.6	8.3	1.2	56.7	1.1	3.7	4.2
Czechoslovakia	9.7	11.4	2.8	12.5	2.0	34.3	5.0	1.1	4.0
Denmark	10.6	10.8	3.7	25.0	9.9	21.9	4.8	0.7	2.4
EastGermany	8.4	11.6	3.7	11.1	5.4	24.6	6.5	0.8	3.6
Finland	9.5	4.9	2.7	33.7	5.8	26.3	5.1	1.0	1.4
France	18.0	9.9	3.3	19.5	5.7	28.1	4.8	2.4	6.5
Greece	10.2	3.0	2.8	17.6	5.9	41.7	2.2	7.8	6.5
Hungary	5.3	12.4	2.9	9.7	0.3	40.1	4.0	5.4	4.2
Ireland	13.9	10.0	4.7	25.8	2.2	24.0	6.2	1.6	2.9
Italy	9.0	5.1	2.9	13.7	3.4	36.8	2.1	4.3	6.7
Netherlands	9.5	13.6	3.6	23.4	2.5	22.4	4.2	1.8	3.7
Norway	9.4	4.7	2.7	23.3	9.7	23.0	4.6	1.6	2.7
Poland	6.9	10.2	2.7	19.3	3.0	36.1	5.9	2.0	6.6
Portugal	6.2	3.7	1.1	4.9	14.2	27.0	5.9	4.7	7.9
Romania	6.2	6.3	1.5	11.1	1.0	49.6	3.1	5.3	2.8
Spain	7.1	3.4	3.1	8.6	7.0	29.2	5.7	5.9	7.2
Sweden	9.9	7.8	3.5	24.7	7.5	19.5	3.7	1.4	2.0
Switzerland	13.1	10.1	3.1	23.8	2.3	25.6	2.8	2.4	4.9
UK	17.4	5.7	4.7	20.6	4.3	24.3	4.7	3.4	3.3
USSR	9.3	4.6	2.1	16.6	3.0	43.6	6.4	3.4	2.9
WestGermany	11.4	12.5	4.1	18.8	3.4	18.6	5.2	1.5	3.8
Yugoslavia	4.4	5.0	1.2	9.5	0.6	55.9	3.0	5.7	3.2

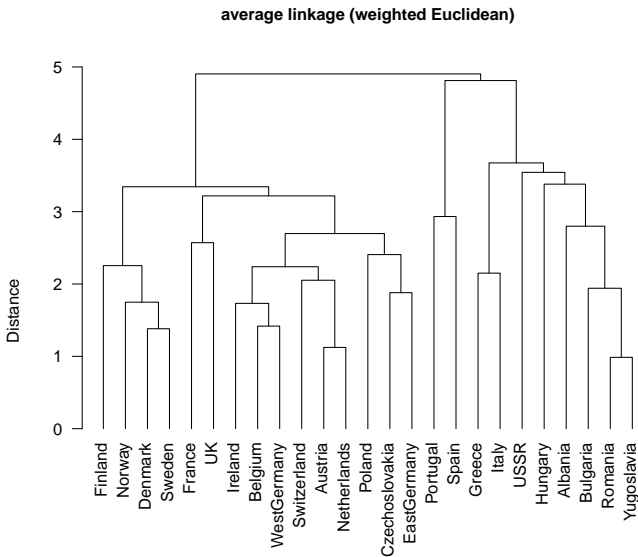
Dendrogram



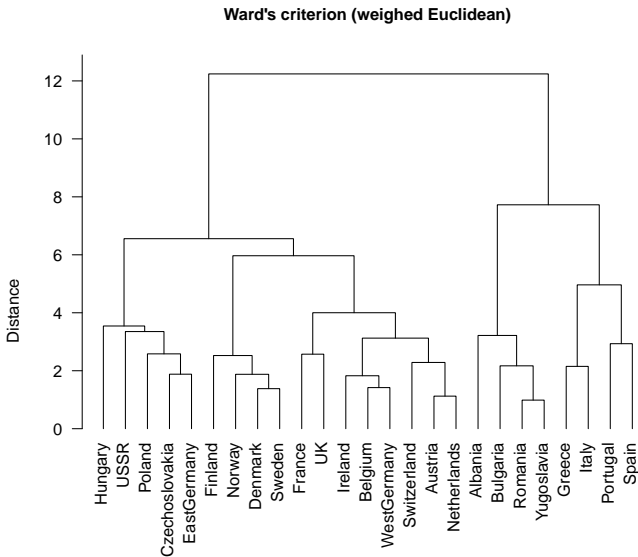
Dendrogram



Dendrogram



Dendrogram

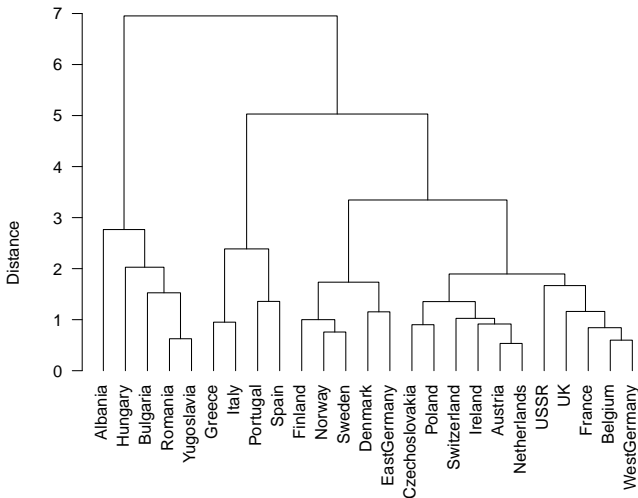


A compositional note

- The protein data set can be seen as a [compositional data set](#)
- Apply [closure](#) by dividing each row by its sum.
- Transform by taking [log-ratios](#)
- Compute the Euclidean distance of the transformed data ([Aitchison distance](#))
- Cluster with this new distance matrix.

Dendrogram

Ward's criterion (Aitchison distance)



Some considerations on cluster distance

single linkage	late inclusion of outliers can identify chain-like clusters sensitive to outliers
complete linkage	fast inclusion of outliers sensitive to outliers
average linkage	less sensitive to outliers
centroid distance	less sensitive to outliers
Ward's criterion	less sensitive to outliers tends to form equally sized clusters

Hierarchical clustering in R

```
> X <- read.table("http://www-eio.upc.es/~jan/data/MVA/protein.dat",header=TRUE)

> rownames(X) <- X[,1]
> X <- X[,-1]

> Xs <- scale(X,scale=TRUE)

> De <- dist(Xs)

> hc.ward <- hclust(De,method="ward.D2")

> plot(hc.ward,ylab="Distance",main="single linkage (weighted Euclidean)",
      xlab="",hang=-1,las=1,cex.main=1)

> clusters <- cutree(hc.ward, 5)
> table(clusters)
clusters
1 2 3 4 5
4 8 5 4 4
>
```


Non-hierarchical Clustering: K-means method

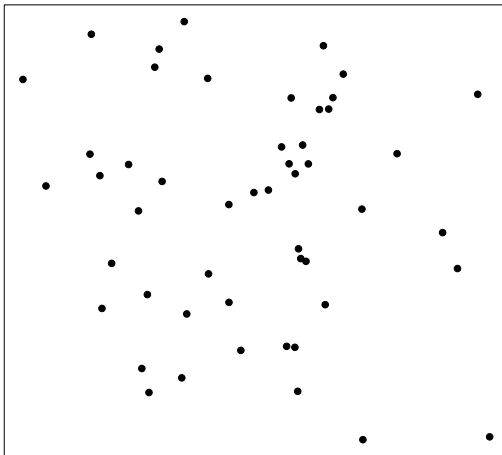
Algorithm:

- 1 Choose a value for the number of clusters K .
- 2 Partition all items into K initial clusters (at random or using seeds).
- 3 Compute the centroids of each cluster.
- 4 Assign each item to the cluster whose centroid is nearest.
- 5 Go back to 3, until there are no re-assignments.

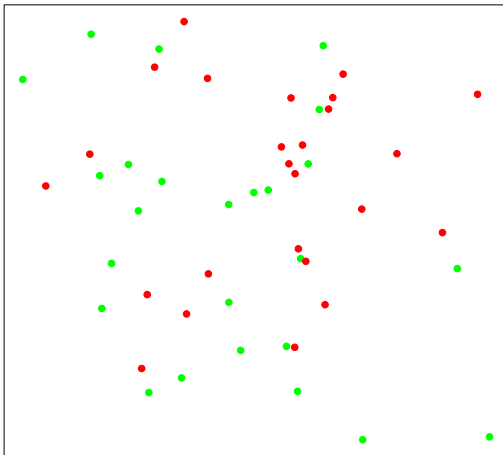
After convergence, it is recommended to

- Try different initial clusters, and compare the final clusters obtained.
- Try a different number of clusters K .
- Compare cluster means and within-cluster variances.

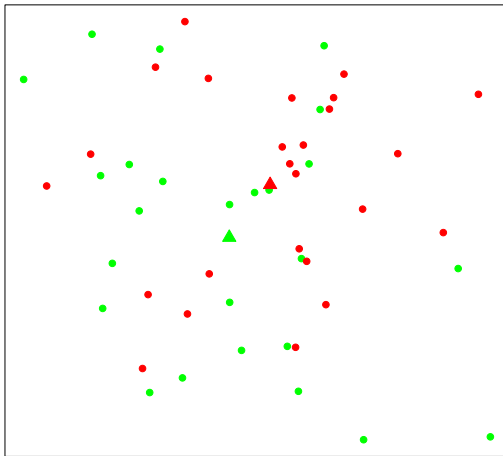
K means graphical illustration



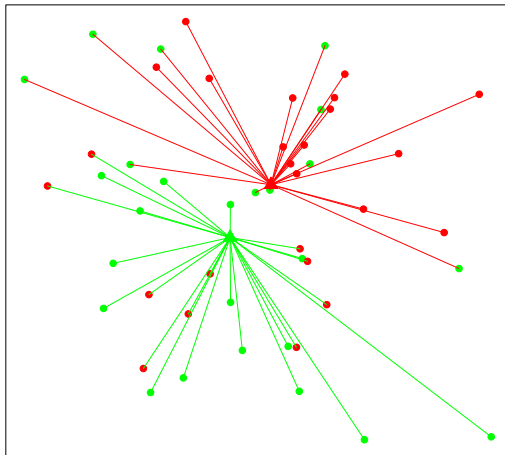
K means graphical illustration



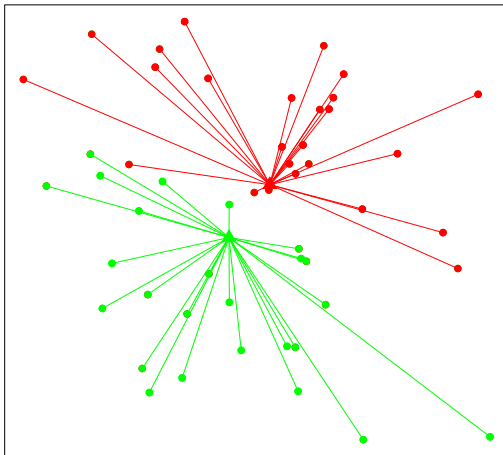
K means graphical illustration



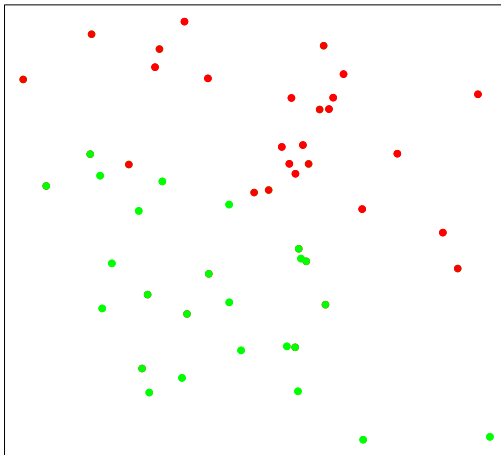
K means graphical illustration



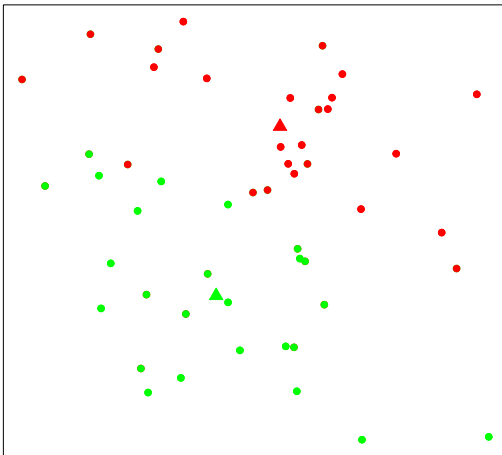
K means graphical illustration



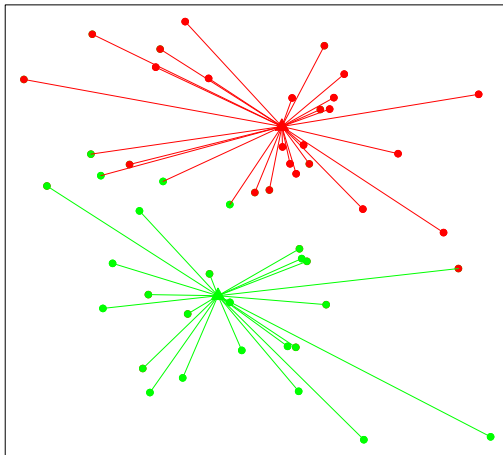
K means graphical illustration



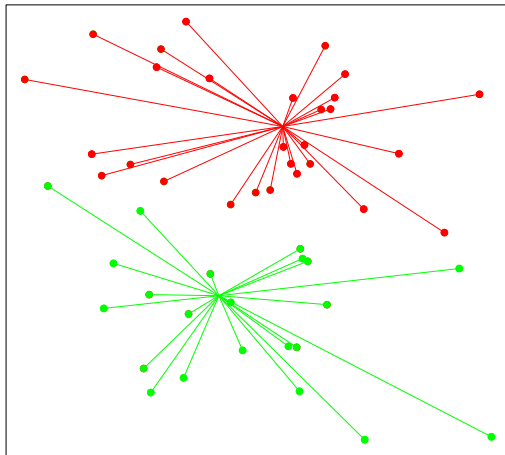
K means graphical illustration



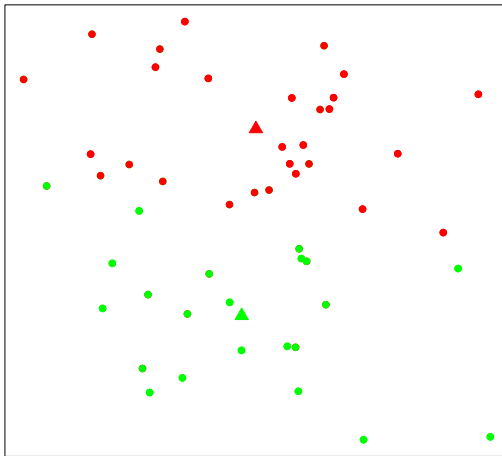
K means graphical illustration



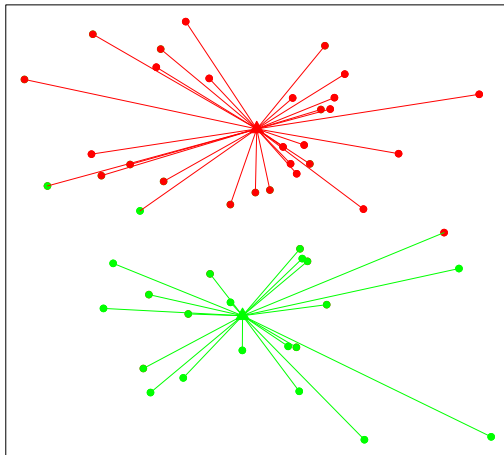
K means graphical illustration



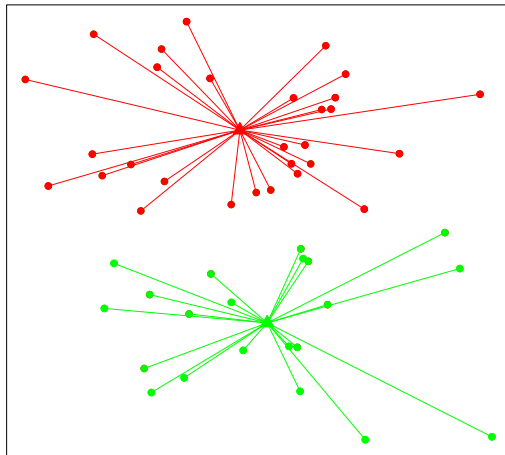
K means graphical illustration



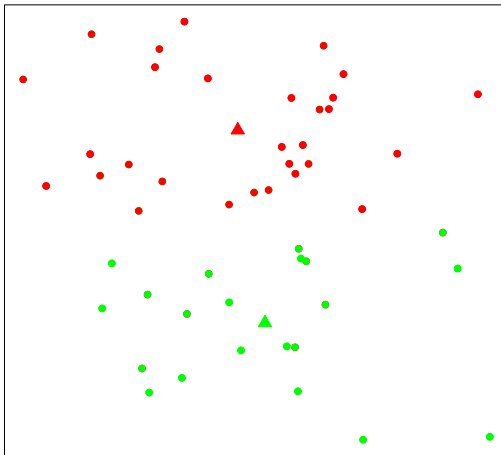
K means graphical illustration



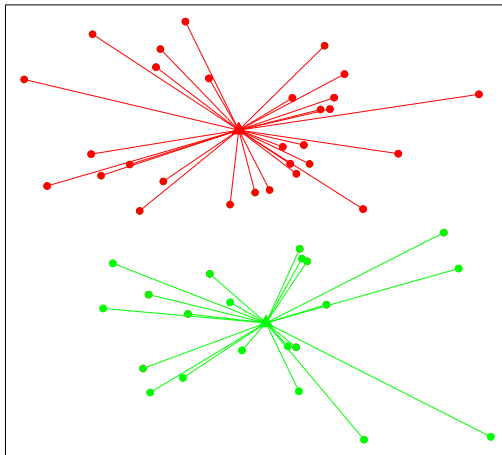
K means graphical illustration



K means graphical illustration



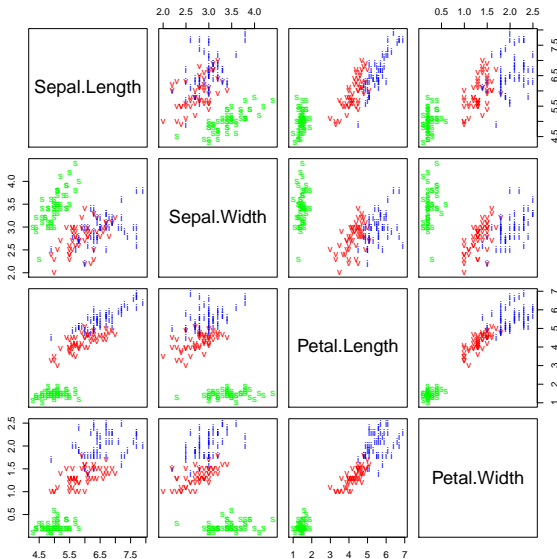
K means graphical illustration



Fisher's Iris data

	Species	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
1	setosa	5.1	3.5	1.4	0.2
2	setosa	4.9	3.0	1.4	0.2
3	setosa	4.7	3.2	1.3	0.2
⋮	⋮	⋮	⋮	⋮	⋮
51	versicolor	7.0	3.2	4.7	1.4
52	versicolor	6.4	3.2	4.5	1.5
53	versicolor	6.9	3.1	4.9	1.5
⋮	⋮	⋮	⋮	⋮	⋮
101	virginica	6.3	3.3	6.0	2.5
102	virginica	5.8	2.7	5.1	1.9
103	virginica	7.1	3.0	5.9	2.1
⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮

Scatterplot matrix



29 / 42

Model-based clustering

- Previous approaches do not make any **distributional assumptions**
- Probabilistic models can be used in clustering and this is called **model-based clustering**
- Finite mixture model

$$g(x|\pi, \theta) = \pi_1 f_1(x|\theta_1) + \pi_2 f_2(x|\theta_2) + \cdots \pi_k f_k(x|\theta_k)$$

- With $\pi_i > 0$ and $\sum_{i=1}^k \pi_i = 1$
- Each f_i is a probability distribution for the i th cluster.
- Usually $f_i \sim N(\mu, \Sigma)$, but not necessarily so.
- The **posterior probabilities** that observation x_j pertains to the i th cluster can be calculated

$$\frac{\pi_i f_i(x_j|\theta_i)}{\sum_{i=1}^k \pi_i f_i(x_j|\theta_i)}$$

Procedure

- A value or estimate of the number of clusters k is needed
- The finite mixture model is estimated by maximum likelihood
- For each observation, the posterior probabilities of pertaining to j cluster are calculated
- Each observation is assigned to the cluster for which it has the largest posterior probability

Model-based clustering in R

```
> library(mclust)
> head(iris)
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1         5.1         3.5         1.4         0.2  setosa
2         4.9         3.0         1.4         0.2  setosa
3         4.7         3.2         1.3         0.2  setosa
4         4.6         3.1         1.5         0.2  setosa
5         5.0         3.6         1.4         0.2  setosa
6         5.4         3.9         1.7         0.4  setosa
> species <- iris$Species
> X <- iris[, -5]
> model1 <- Mclust(X, G=3)
fitting ...
|=====| 100%
> summary(model1, parameters = TRUE)
-----
Gaussian finite mixture model fitted by EM algorithm
-----

Mclust VEV (ellipsoidal, equal shape) model with 3 components:

log-likelihood   n df      BIC      ICL
-186.074 150 38 -562.5522 -566.4673

Clustering table:
  1  2  3
50 45 55

Mixing probabilities:
      1      2      3
0.3333333 0.3005423 0.3661243
```

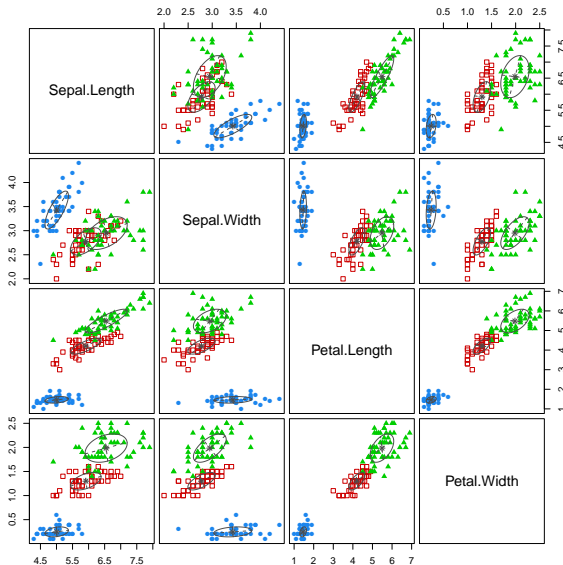
Model-based clustering in R

```
Means:
      [,1]      [,2]      [,3]
Sepal.Length 5.006 5.915044 6.546807
Sepal.Width  3.428 2.777451 2.949613
Petal.Length 1.462 4.204002 5.482252
Petal.Width  0.246 1.298935 1.985523

Variances:
[,1]
      Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length 0.13320850 0.10938369 0.019191764 0.011585649
Sepal.Width  0.10938369 0.15495369 0.012096999 0.010010130
Petal.Length 0.01919176 0.01209700 0.028275400 0.005818274
Petal.Width  0.01158565 0.01001013 0.005818274 0.010695632
[,2]
      Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length 0.22572159 0.07613348 0.14689934 0.04335826
Sepal.Width  0.07613348 0.08024338 0.07372331 0.03435893
Petal.Length 0.14689934 0.07372331 0.16613979 0.04953078
Petal.Width  0.04335826 0.03435893 0.04953078 0.03338619
[,3]
      Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length 0.42943106 0.10784274 0.33452389 0.06538369
Sepal.Width  0.10784274 0.11596343 0.08905176 0.06134034
Petal.Length 0.33452389 0.08905176 0.36422115 0.08706895
Petal.Width  0.06538369 0.06134034 0.08706895 0.08663823
> plot(model1, what="classification")
> table(species, model1$classification)

species      1  2  3
setosa      50  0  0
versicolor  0 45  5
virginica   0  0 50
```

Model-based clustering in R



Cluster validity indices

- Choose the optimal number of clusters according to some (numerical) criterion
- Several criteria have been developed
- Some popular criteria:
 - Pseudo F -statistics (Calinski-Harabasz, 1974)
 - Silhouette coefficient
 -

Pseudo F statistics

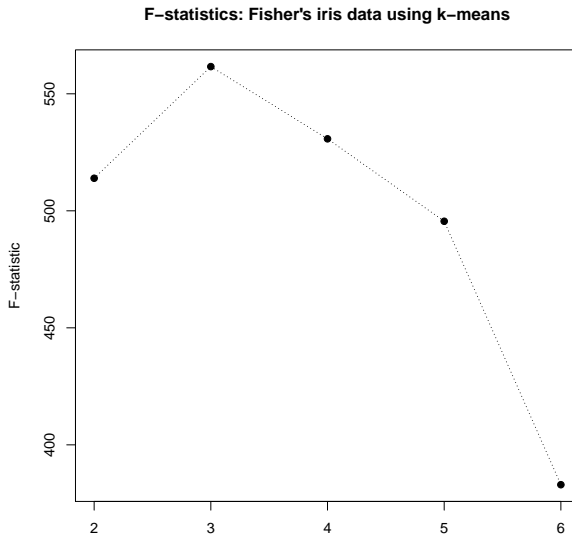
$$F = \frac{GSS/(K - 1)}{WSS/(N - 1)}$$

with:

- K = number of groups
- N = sample size
- GSS = between-group sum-of-squares
- WSS = within-group sum-of-squares

Choose the number of clusters that maximizes F

Example F statistics



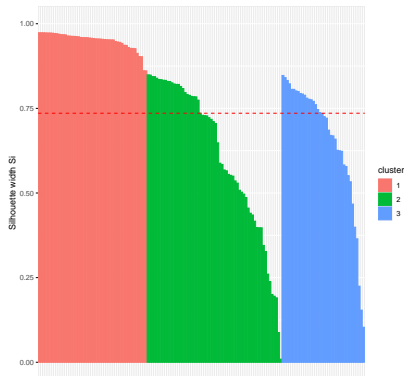
Silhouette scores and coefficient

- Let C_i represent cluster i with $i = 1, \dots, k$
- Let a_i be the average distance between i and all other points in the same cluster
- Let b_i be the minimal average distance between i and all other points in another cluster.
- The silhouette score is defined as

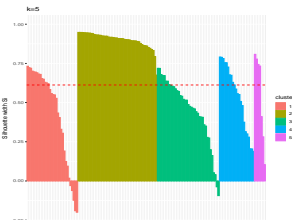
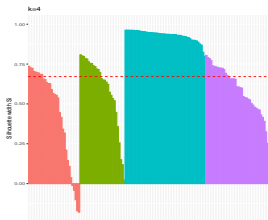
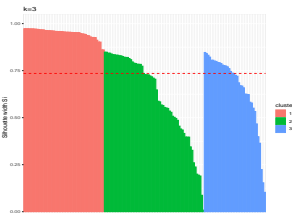
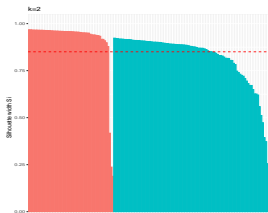
$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)} \text{ and satisfies } -1 \leq s_i \leq 1$$

- s_i measures how well a case matches its cluster.
- s_i can be averaged over all observations to give the average silhouette score.
- Choose the number of clusters that maximizes this average.

Silhouette scores iris data ($k = 3$ with k-means)



For varying k



Final remarks

After obtaining the clusters, some questions remain

- Are the clusters really different? (manova/anova)
- How homogeneous is each cluster?
- Which variables discriminate the clusters? (descriptive statistics per group, LDA/QDA)

In cluster analysis, the user has to make several choices:

- 1 the variables to include
- 2 possible transformations
- 3 the algorithm to use (hierarchical, divisive, model-based, ...)
- 4 the distance measure to use (Euclidean, City-Block, Mahalanobis, ...)
- 5 a measure of distance between clusters.
- 6

References

- Manly, B.F.J. (1989) Multivariate statistical methods: a primer. 3rd edition. Chapman and Hall, London.
- Johnson & Wichern (2002) Applied Multivariate Statistical Analysis. 5th edition. Prentice Hall, Chapter 12.
- Everitt, B.S., Landau, S., Lees, M. & Stahl, D. (2011) Cluster Analysis. 5th edition. Wiley.