

This document is an historical remnant. It belongs to the collection Skeptron Web Archive (included in Donald Broady's archive) that mirrors parts of the public Skeptron web site as it appeared on 31 December 2019, containing material from the research group Sociology of Education and Culture (SEC) and the research programme Digital Literature (DL). The contents and file names are unchanged while character and layout encoding of older pages has been updated for technical reasons. Most links are dead. A number of documents of negligible historical interest as well as the collaborators' personal pages are omitted.

The site's internet address was since Summer 1993 www.nada.kth.se/~broady/ and since 2006 www.skeptron.uu.se/broady/sec/.

IV — Introduction to Euclidean Classification

Brigitte.LeRoux@math-info.univ-paris5.fr
rouanet@math-info.univ-paris5.fr

www.math-info.univ-paris5.fr/~lerb/
www.math-info.univ-paris5.fr/~rouanet/

1 What is Euclidean Classification?

1.1 Introduction

a) Partition and hierarchy

Hierarchical classification: System of nested classes (the paradigm of natural sciences) represented by a hierarchical tree.

b) Qualities of classification: compactness and separability

c) Descending (or divisive) classification vs ascending (or agglomerative) classification

1.2 Ascending Hierarchical Classification (AHC) using Variance Criterion

Grouping property: If 2 classes are grouped together, the between–variance decreases from an amount equal to the contribution of the dipole defined by the centers of the 2 grouped classes.

Target Example (see II): Three–class partition \mathcal{A} , \mathcal{B} and \mathcal{C} with between-variance 57.43 (variance of the cloud of 3 mean points (A,B,C) of classes). If \mathcal{A} and \mathcal{B} are grouped, the between–variance of the partition in 2 classes, that is, the variance of the cloud of 2 points (barycenter of A and B, C) is equal to 38.10.

Within-contribution of the pair (A,B): $\widetilde{\frac{n_{AB}}{n}} \times (AB)^2 = 19.33$, with $AB^2 = 290$ and $\widetilde{n_{AB}} = \frac{1}{\frac{1}{2} + \frac{1}{1}} = 2/3$ (weight of dipole);

One has: $38.10 = 57.43 - 19.33$ (grouping property).

Ascending Hierarchical Classification: starting with the basic objects (one-element classes) proceed to successive aggregations, until all objects are grouped in a single class.

At each step, one groups 2 classes of the current partition.

Euclidean classification:

1. Objects = *points of Euclidean cloud*: distance between objects is Euclidean distance.
2. *Aggregation index* = variance index, that is, the contribution of the dipole associated with the 2 aggregated classes (Ward index).

At each step, the aggregated classes are those which lead to the minimal decrease of the between-variance.

1.3 Basic Algorithm

- **Step 1.** Calculate the contributions of the $9 \times 10/2 = 45$ dipoles

δ	$i1$	$i2$	$i3$	$i4$	$i5$	$i6$	$i7$	$i8$	$i9$
$i2$	2								
$i3$	11.6	4							
$i4$	6.8	3.2	4						
$i5$	14.4	6.8	2	2					
$i6$	13	17	27.4	10.6	20.2				
$i7$	13	10.6	12.2	2.6	5.8	5.2			
$i8$	14.6	9.8	8.2	1.8	2.6	10	0.8		
$i9$	29.2	20.8	13.6	8	5.2	19.4	5	2.6	
$i10$	31.4	21.8	13	9	5	23.2	6.8	3.6	0.2

Example: For dipole $\{i1, i2\}$: $\widetilde{n_{12}} = 1/(\frac{1}{1} + \frac{1}{1}) = 0.5$, squared distance = $(0 - 6)^2 + (-12 + 10)^2 = 40$, hence the absolute contribution of dipole $\frac{0.5}{10} \times 40 = 2$.

Minimum index 0.2 for the pair of points $\{i9, i10\}$ which are aggregated (fig. 1), hence the mean point ℓ_{11} and a derived *cloud of 9 points* (fig. 2).

• **Step 2.** Calculate the aggregation index between the new point ℓ_{11}

and the 8 other points

	$i1$	$i2$	$i3$	$i4$	$i5$	$i6$	$i7$	$i8$
ℓ_{11}	40.33	28.33	17.67	11.27	6.73	28.33	7.8	4.07

New minimum 0.8 for $\{i7, i8\}$ which aggregated (fig. 2), hence the new point ℓ_{12} and a derived *cloud of 8 points* (fig. 3).

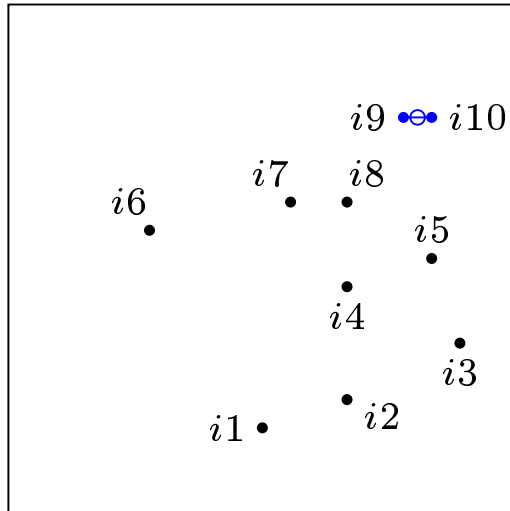


Figure 1

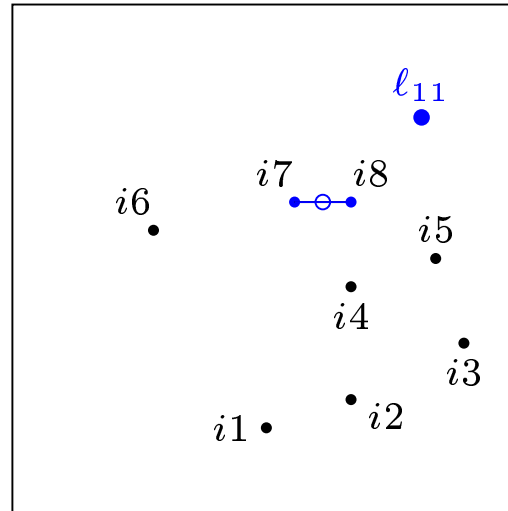


Figure 2

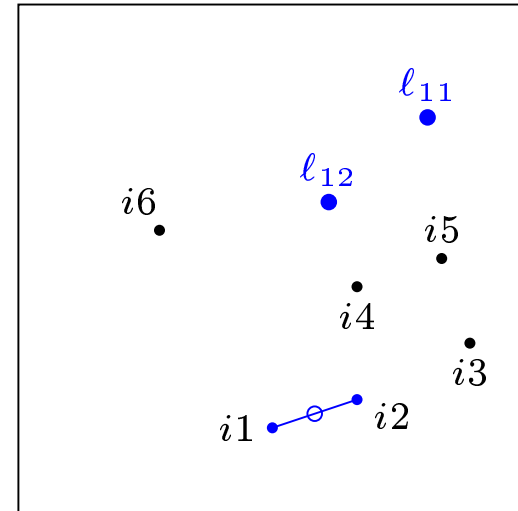


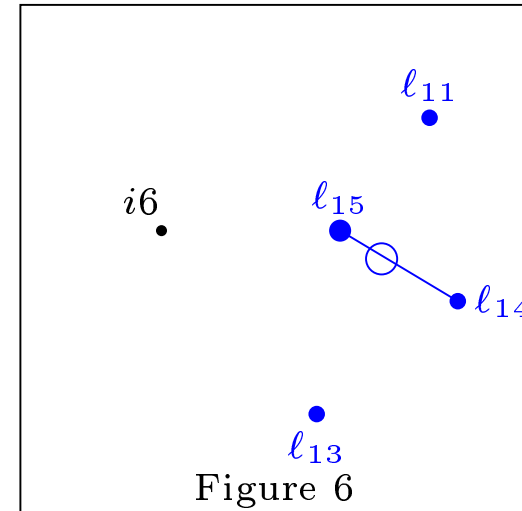
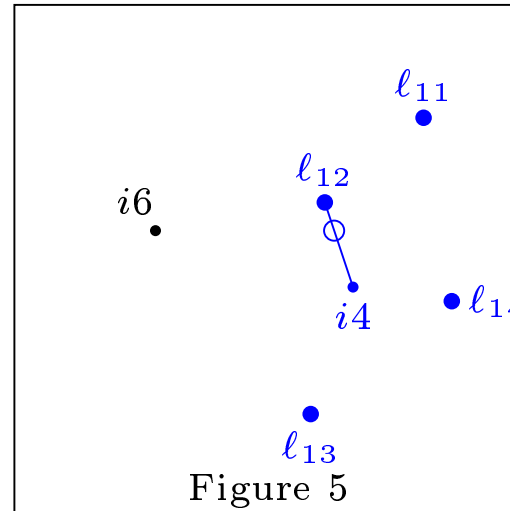
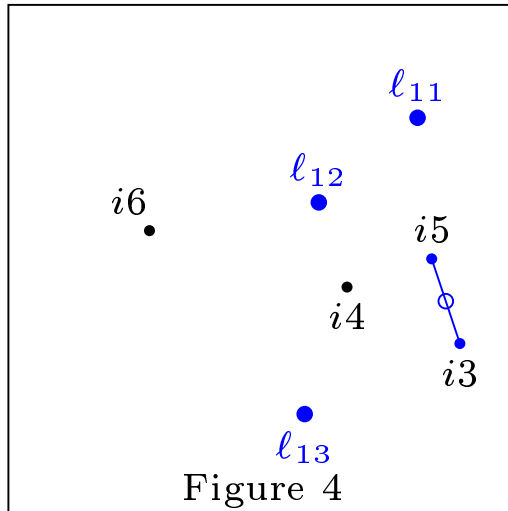
Figure 3

• **Step 3.** Iterate the procedure

Aggregation index between ℓ_{12} and the 7 other points

	$i1$	$i2$	$i3$	$i4$	$i5$	$i6$	ℓ_{11}
ℓ_{12}	18.13	13.33	13.33	2.67	5.33	9.87	8.2

Minimum of index = 2 for the three pairs $\{i1, i2\}$, $\{i3, i5\}$ and $\{i4, i5\}$. We choose^a to aggregate $i1$ and $i2$ (fig. 3), hence the point ℓ_{13} and a *cloud of 7 points* (fig. 4).



^aIn indeterminate cases different choices may yield different classifications.

Aggregation index between ℓ_{13} and the 6 other points

	$i3$	$i4$	$i5$	$i6$	ℓ_{11}	ℓ_{12}
ℓ_{13}	9.73	6.00	13.47	19.33	50.5	22.6

Minimum of index = 2 for the two pairs $\{i3, i5\}$ and $\{i4, i5\}$. We choose to aggregate $i3$ and $i5$ (fig. 4), hence the point ℓ_{14} and the *cloud of 6 points* (fig. 5).

Aggregation index between ℓ_{14} and the 5 other points

	$i4$	$i6$	ℓ_{11}	ℓ_{12}	ℓ_{13}
ℓ_{14}	3.33	31.07	17.33	13.00	16.4

→ aggregation of ℓ_{12} and $i4$ at level 2.67 (fig. 5), hence the point ℓ_{15} and the *cloud of 5 points* (fig. 6).

Aggregation index between ℓ_{15} and the 4 other points

	$i6$	ℓ_{11}	ℓ_{13}	ℓ_{14}
ℓ_{15}	12.03	12.49	20.61	11.33

→ aggregation of ℓ_{15} and ℓ_{14} at level 11.33 (fig. 6), hence the point ℓ_{16} and the *cloud of 4 points* (fig. 7).

Aggregation index between $\ell 16$ and the 3 other points

	$i6$	$\ell 11$	$\ell 13$
$\ell 16$	21.67	15.57	20.86

\rightarrow aggregation of $\ell 16$ and $\ell 11$ at level 15.57 (fig. 7), hence the point $\ell 17$ and the *cloud of 3 points* (fig. 8).

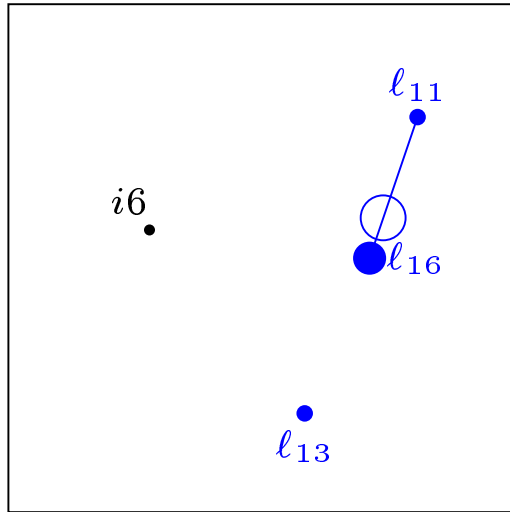


Figure 7

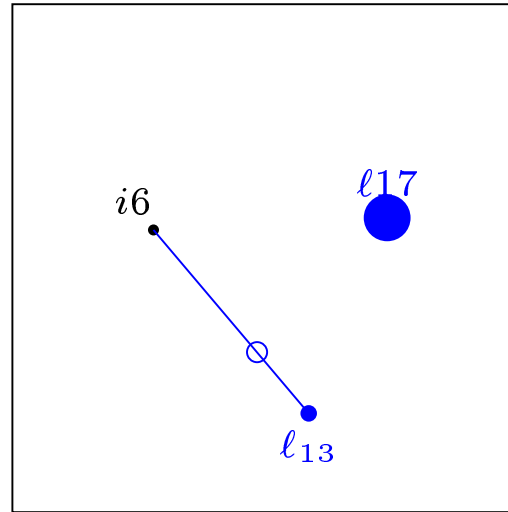


Figure 8

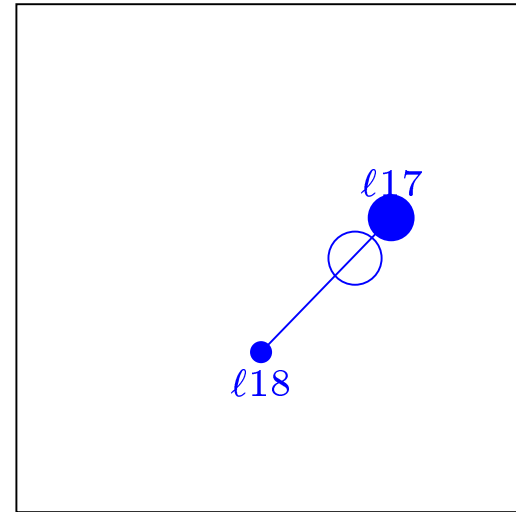


Figure 9

The three-class partition $\mathcal{A}, \mathcal{B}, \mathcal{C}$ (already studied in II) with mean points A ($\ell 13$), B ($i6$), C ($\ell 17$) (fig. 8).

Table of the within-contributions of the 3 pairs of points

(distance) ²	weight	Contribution
$AB^2 = 290$	$\widetilde{n_{AB}} = \frac{1}{\frac{1}{2} + \frac{1}{1}} = 2/3$	$Cta_{(A,B)} = \frac{2/3}{10} \times 290 = 19.33$
$AC^2 = 226.33$	$\widetilde{n_{AC}} = \frac{1}{\frac{1}{2} + \frac{1}{7}} = 14/9$	$Cta_{(A,C)} = \frac{14/9}{10} \times 226.33 = 35.21$
$BC^2 = 284.90$	$\widetilde{n_{BC}} = \frac{1}{\frac{1}{1} + \frac{1}{7}} = 7/8$	$Cta_{(B,C)} = \frac{7/8}{10} \times 284.90 = 24.93$

At this step, we group A and B at level 19.33 (fig. 9).

Successive steps of the AHC

ℓ	δ_ℓ	classes		n	class description
ℓ_{19}	38.095	ℓ_{18}	ℓ_{17}	10	$i_9 i_{10} i_3 i_5 i_4 i_7 i_8 i_6 i_1 i_2$
ℓ_{18}	19.333	ℓ_{13}	ℓ_6	3	$i_6 i_1 i_2$
ℓ_{17}	15.571	ℓ_{16}	ℓ_{11}	7	$i_9 i_{10} i_3 i_5 i_4 i_7 i_8$
ℓ_{16}	11.333	ℓ_{15}	ℓ_{14}	5	$i_3 i_5 i_4 i_7 i_8$
ℓ_{15}	2.667	ℓ_{12}	ℓ_4	3	$i_4 i_7 i_8$
ℓ_{14}	2.	ℓ_5	ℓ_3	2	$i_3 i_5$
ℓ_{13}	2.	ℓ_2	ℓ_1	2	$i_1 i_2$
ℓ_{12}	0.8	ℓ_8	ℓ_7	2	$i_7 i_8$
ℓ_{11}	0.2	ℓ_{10}	ℓ_9	2	$i_9 i_{10}$

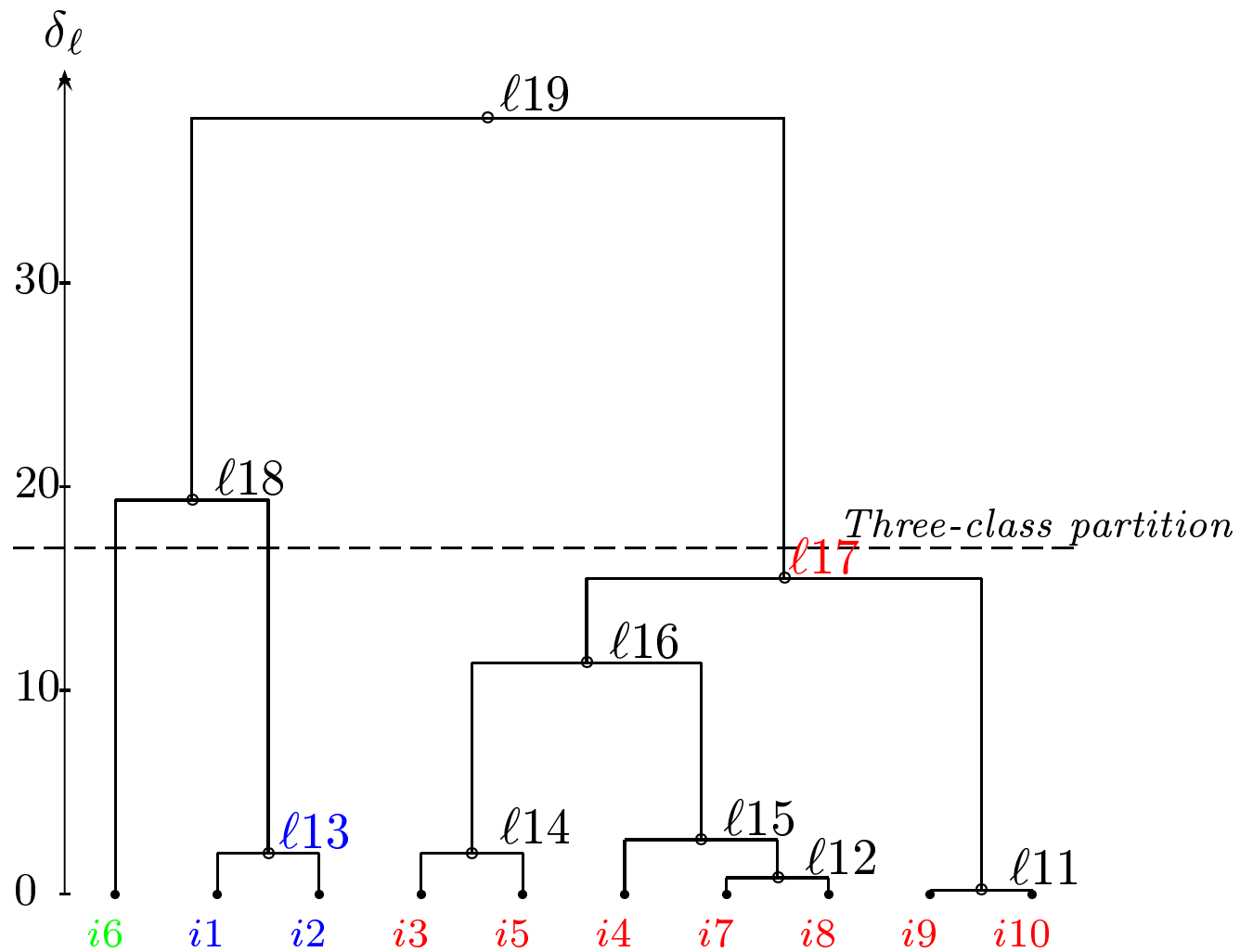
<i>Between</i> Var		η_ℓ^2
ℓ_{19}	38.10	.414
ℓ_{18}	57.43	.624
ℓ_{17}	73.00	.793
ℓ_{16}	84.33	.917
ℓ_{15}	87.00	.957
ℓ_{14}	89.00	.967
ℓ_{13}	91.90	.989
ℓ_{12}	91.80	.998
ℓ_{11}	92.00	1

The sum of the nine level indices δ_ℓ is 92 (total variance of the cloud).

Between-variance of the 2-class partition 38.095.

Between-variance of the 3-class partition $38.095 + 19.333 = 57.43$, etc.

Target example: hierarchical tree



References

- BENZÉCRI J-P. (1992) *Correspondence Analysis Handbook*, (Part V), New York: Dekker (p. 561-635).
- BOURDIEU P. (1999). Une révolution conservatrice dans l'édition [A conservative revolution in publishing], *Actes de la Recherche en Sciences Sociales*, Vol. 126-127, 3-28.
- LE ROUX B. & ROUANET H. (2003). Geometric Analysis of Individual Differences in Mathematical Performance for EPGY Students in the Third Grade. www-epgy.stanford.edu/research/.
- LE ROUX B. & ROUANET H. (2004), *Geometric Data Analysis: from Correspondence Analysis to Structured Data Analysis* (chapter 3, p.106-116), Dordrecht: Kluwer.