

# A FAMILY OF HIERARCHICAL CLUSTERING ALGORITHMS BASED ON HIGH-ORDER DISSIMILARITIES

Helena Aidos and Ana Fred

Instituto de Telecomunicações  
Instituto Superior Técnico, Universidade de Lisboa, Portugal  
{haidos, afred}@lx.it.pt

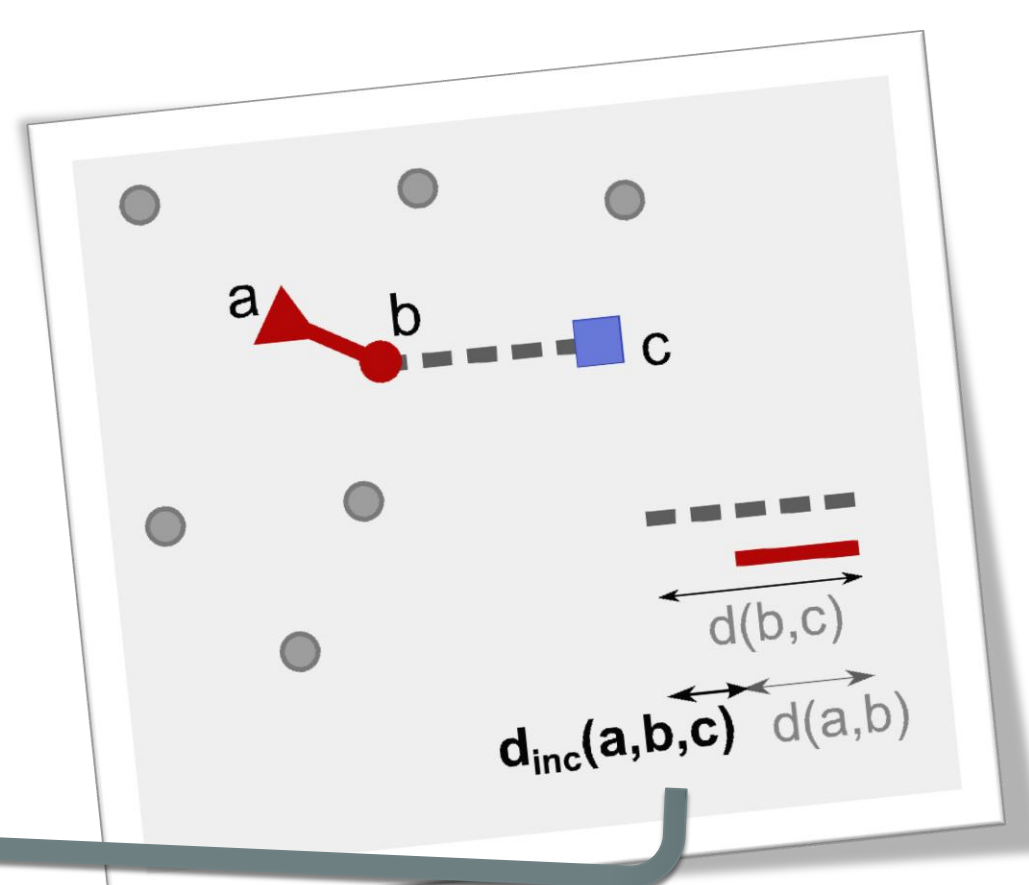
## Dissimilarity increments: definition and distribution

$(a, b, c)$  – triplet of nearest neighbors

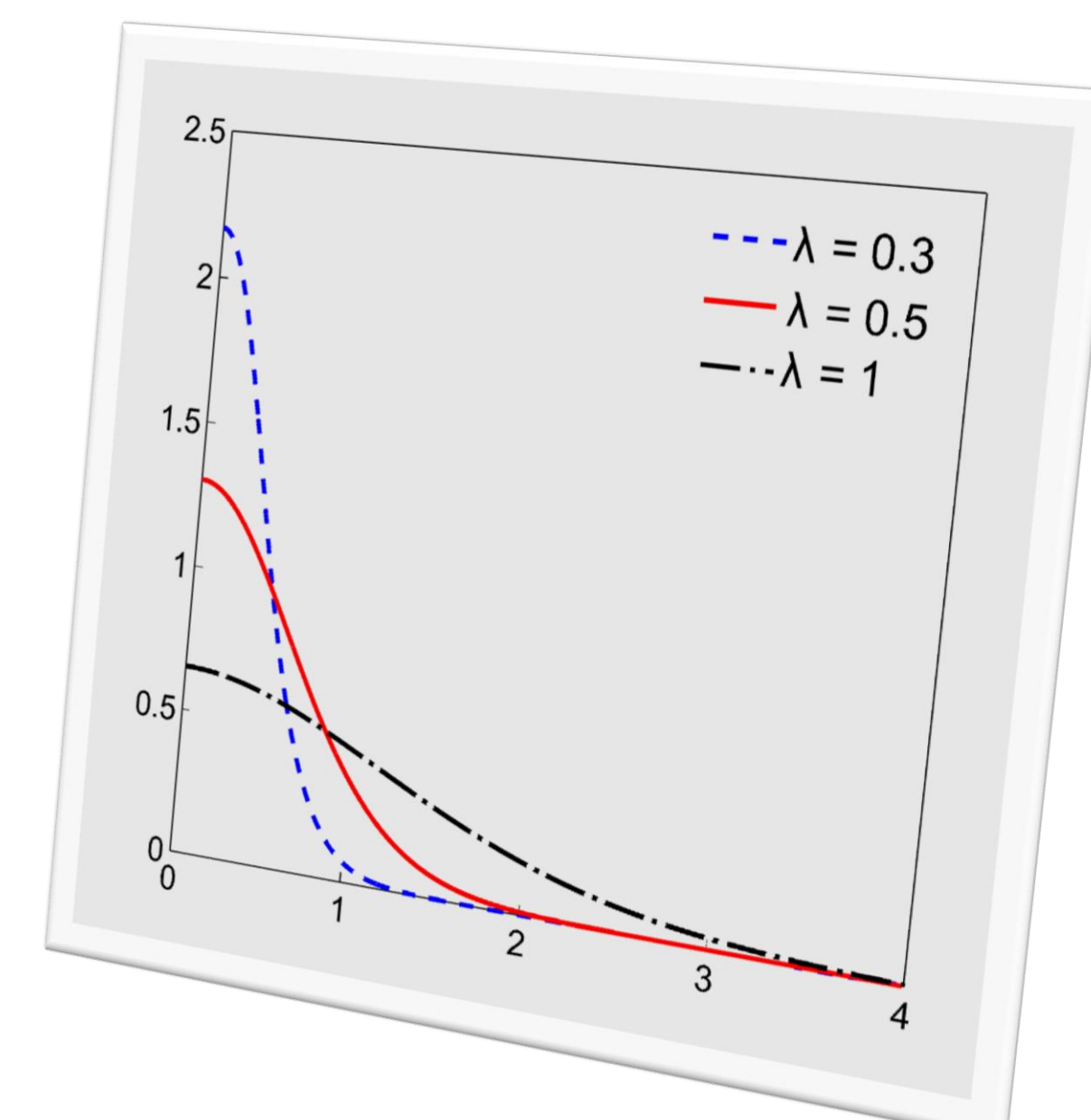
- $a$  is the nearest neighbor of  $b$
- $c$  is the nearest neighbor of  $b$  (different from  $a$ )

The **dissimilarity increments** between neighboring patterns is defined as

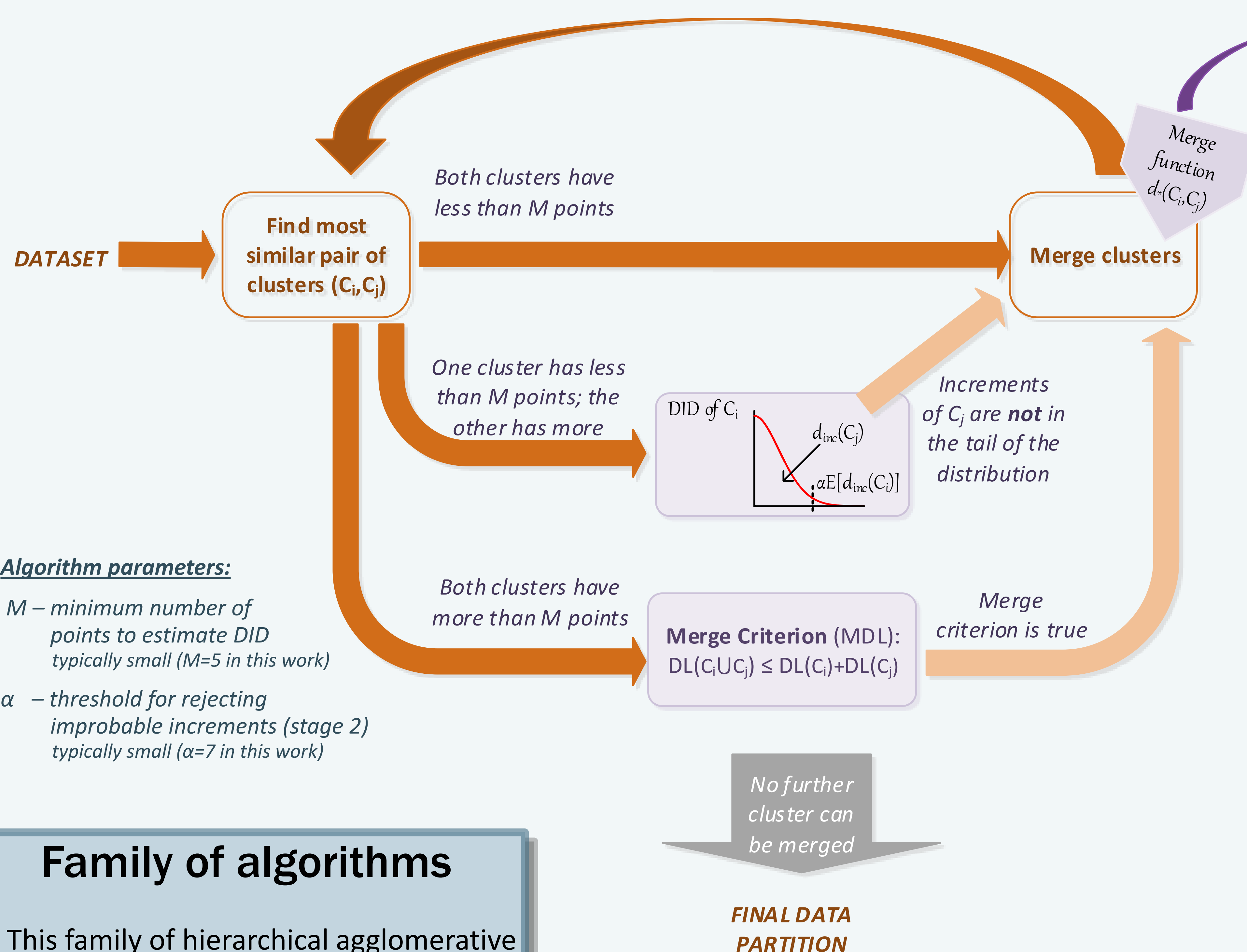
$$d_{inc}(a, b, c) = |d(a, b) - d(b, c)|$$



The **dissimilarity increments distribution (DID)** is a function of the mean value of the dissimilarity increments



**PROPOSAL:** A family of agglomerative hierarchical methods, integrating dissimilarity increments in traditional linkage algorithms



### Algorithm parameters:

- $M$  – minimum number of points to estimate DID typically small ( $M=5$  in this work)
- $\alpha$  – threshold for rejecting improbable increments (stage 2) typically small ( $\alpha=7$  in this work)

### MERGE FUNCTION

Consider the new formed cluster  $C_k = C_i \cup C_j$ , obtained by merging  $C_i$  and  $C_j$ , and  $C_k$  is one of the remaining clusters formed in previous steps. Lets consider  $|C_k|$  the number of points in cluster

#### • SLDID

$$d_{inc}(C_i, C_j, C_k) = |d(C_i, C_j) - d(C_i, C_k)|$$

#### • ALDID

$$d_{inc}(C_i, C_j, C_k) = |d(C_i, C_j) - d(C_j, C_k)|$$

#### • CLDID

$$d_{inc}(C_i, C_j, C_k) = |d(C_i, C_j) - d(C_i, C_k)|$$

#### • WLDID

$$d_{inc}(C_i, C_j, C_k) = |d(C_i, C_j) - d(C_j, C_k)|$$

$$d_{inc}(C_i, C_j, C_k) = |d(C_i, C_j) - d(C_i, C_k)|$$

$$d_{inc}(C_i, C_j, C_k) = |d(C_i, C_j) - d(C_j, C_k)|$$

## Family of algorithms

➤ This family of hierarchical agglomerative algorithms is able to **automatically find the number of clusters** using a minimum description length criterion based on the dissimilarity increments distribution (DID)

➤ Each algorithm of the proposed family is able to find classes as unions of clusters, leading to the **identification of internal structures of classes**

**DATASETS:**  
36 real-world datasets from the UCI Machine Learning Repository.

**EVALUATION:**  
Percentage of correctly clustered points assuming that one class can be represented as the union of several clusters.

## Experimental results

