



Centrum voor Wiskunde en Informatica

REPORTRAPPORT

Lipshits distance and hierarchical clustering

M. Hazewinkel

Department of Analysis, Algebra and Geometry

AM-R9506 1995

Report AM-R9506
ISSN 0924-2953

CWI
P.O. Box 94079
1090 GB Amsterdam
The Netherlands

CWI is the National Research Institute for Mathematics and Computer Science. CWI is part of the Stichting Mathematisch Centrum (SMC), the Dutch foundation for promotion of mathematics and computer science and their applications.

SMC is sponsored by the Netherlands Organization for Scientific Research (NWO). CWI is a member of ERCIM, the European Research Consortium for Informatics and Mathematics.

Copyright © Stichting Mathematisch Centrum
P.O. Box 94079, 1090 GB Amsterdam (NL)
Kruislaan 413, 1098 SJ Amsterdam (NL)
Telephone +31 20 592 9333
Telefax +31 20 592 4199

Lipshits distance and hierarchical clustering

Michiel Hazewinkel
 CWI
 POBox 94079
 1090GB Amsterdam
 The Netherlands
 mich@cwi.nl

Abstract

This note discusses the Lipschitz distance between two different metrics (or definite) dissimilarities on a set. Given a metric space a universal lower bound is established for the Lipschitz distance between the original metric on M and the metric defined by any hierarchical classification tree on M . Finally it is shown that single link clustering attains this lower bound.

Mathematics Subject Classification 1991: 62H30, 68T10, 54E35.

Key words & Phrases: Lipschitz distance, hierarchical clustering, single link clustering.

1. Classification trees. A definite dissimilarity space is a (finite) set M together with a function $d : M \times M \rightarrow \mathbf{R}$ such that $d(x, y) = d(y, x) > 0$, $d(x, x) = 0$ for all $x, y \in M, x \neq y$. This function is called a dissimilarity. If in addition the triangle inequality is satisfied, (M, d) is a metric space.

In this note a hierarchical classification tree for M is a series of coarser and coarser partitions

$$\{\text{singletons}\} = \pi_0 \preceq \pi_1 \preceq \pi_2 \preceq \dots \preceq \pi_m = \{M\} \quad (1.1)$$

together with a (level) function on the set of partitions in the form of m strictly positive numbers

$$d_1 < d_2 < \dots < d_m. \quad (1.2)$$

These data define a new metric on M by

$$d_T(x, y) = d_j \quad (1.3)$$

if and only if j is the smallest integer such that both x, y are in a same set of the partition π_j .

Inversely, given an ultrametric which assumes the values $d_i, i = 1, \dots, m$, the balls with diameter d_i (maximal sets all of whose members have distance less than d_i to one another) define a partition π_i ; the correspondence thus set up is bijective [Hartigan, 1967; Johnson, 1967]. In this note the sequence of partitions point of view is more convenient for the proofs and constructions.

2. Lipschitz distance.

Consider a set M and two metrics (or dissimilarities), d_1, d_2 defined on it. The *distortion* of d_2 with respect to d_1 is defined by

$$\text{distor}(d_2, d_1) = \max \frac{d_2(x, y)}{d_1(x, y)} \quad (2.1)$$

where the max is taken over all $x, y \in M$, $x \neq y$. The Lipschitz distance between d_1, d_2 is now defined as

$$\delta_L(d_1, d_2) = \log(\text{distor}(d_2, d_1) \text{distor}(d_1, d_2)). \quad (2.2)$$

Note that if the two distances are proportional, their Lipschitz distance is zero. This feature is really an advantage for classification problems because a constant scalar factor should not matter.

It is easy to see that:

4.2. *Proposition* [Johnson, Lindenstrauss, and Schechtman, 1987; Matousek, 1990]. The Lipschitz distance δ_L defines a metric on isometry classes up-to a scalar factor of metrics (or definite dissimilarities) on a fixed set M .

3. A universal lower bound.

A path P from x to y in M is simply a sequence of points $x = x_0, x_1, \dots, x_n = y$. The step length (sl) of P is equal to

$$sl(P) = \max_{i=0, \dots, n} \{d(x_i, x_{i+1})\}. \quad (3.1)$$

The step-across-separation (sas) of two unequal points $x, y \in M$ is defined as

$$sas(x, y) = \min_P \{sl(P)\} \quad (3.2)$$

where the minimum is taken over all paths P from x to y . Finally the ‘‘distance step across separation quotient’’ (dsq) of M is defined as

$$dsq(M) = \max_{x \neq y} \frac{d(x, y)}{sas(x, y)}. \quad (3.3)$$

3.4. *Theorem.* Let $\pi = (\pi_0, \pi_1, \dots, \pi_n)$ be a classification tree on M with associated distance d_T . Then $\delta_L(d, d_T) \geq \log(dsq(M))$.

Proof. Let $x, y \in M$ be such that the maximum in (3.3) is assumed for this pair. Moreover, let $x = x_0, x_1, \dots, x_n = y$ be a path for which the step length is equal to $sas(x, y)$. Let $d_T(x, y) = d_j$. Then

$$\text{distor}(d, d_T) \geq \frac{d(x, y)}{d_j}. \quad (3.5)$$

Now consider the chain $x = x_0, x_1, \dots, x_n = y$. Let $\pi_j = \{C_{1,j}, \dots, C_{m,j}\}$. We can assume that $x, y \in C_1$. Suppose that in π_{j-1} the set C_1 splits up into the disjoint sets $C_{11,j-1}, \dots, C_{1k,j-1}$. We can assume that $x \in C_{11}$, and then $y \notin C_{11}$. Therefore there is a first x_i that is not in C_{11} . There are two possibilities. If $x_i \in C_1$, then $d_T(x_{i-1}, x_i) = d_j$; if $x_i \notin C_1$, then the finest partition with a set in it that contains both x_{i-1}, x_i has index greater than j , so that

$$d_T(x_{i-1}, x_i) \geq d_{j+1} > d_j.$$

Also $d(x_{i-1}, x_i) \leq sas(x, y)$. Combining these two we see that

$$(3.6) \quad \text{distor}(d_T, d) \geq \frac{d_j}{\text{sas}(x, y)}. \quad (3.6)$$

Inequalities (3.5) and 3.6) together prove the theorem.

4. Single link clustering.

Single link clustering yields the following hierarchy. Let $d_1 < d_2 < \dots < d_s$ be the distances that actually occur in the original metric. Let G_j be the graph on M which has two vertices linked if and only if their distance in M is $\leq d_j$. Then the partition π_j has level d_j and it consists of the connected components of the graph G_j . Let d_{sl} , the single link clustering hierarchy distance, be the distance defined by this classification tree.

4.1. *Theorem.* Single link clustering is optimal with respect to Lipschitz distance. More precisely:

$$(4.2) \quad \delta_L(d_{sl}, d) = \log(dsq(M)).$$

Proof. In view of theorem 3.4, it suffices to show that

$$(4.3) \quad \delta_L(d_{sl}, d) \leq \log(dsq(M)).$$

Let $x, y \in M, x \neq y$. Then $\text{sas}(x, y) = d_{sl}(x, y)$. It follows that

$$(4.4) \quad \text{distor}(d, d_{sl}) = \max \frac{d(x, y)}{d_{sl}(x, y)} \leq dsq(M)$$

On the other hand, as is rather obvious and also well known, [Jardine and Sibson, 1971, p. 63], $d_{sl}(x, y) \leq d(x, y)$ for all $x, y \in M$. Hence

$$\text{distor}(d_{sl}, d) \leq 1. \quad (4.5)$$

The combination of (4.4) and 4.5) establishes (4.3) and hence (4.2) and the theorem.

4.6. *Remark.* In [Jardine, Jardine, and Sibson, 1967] it is shown that d_u is subdominant; i.e. that it is maximal with respect to the ordering $d \preceq d' \Leftrightarrow \forall_{x,y} d(x, y) \leq d'(x, y)$ among the ultrametrics that are no greater in this ordering than the original metric. This, however, carries no implications with respect to its Lipschitz distance to the original metric.

4.7. *Corollary.* As noted during the proof

$$\text{sas}(x, y) = d_{sl}(x, y)$$

5. Example.

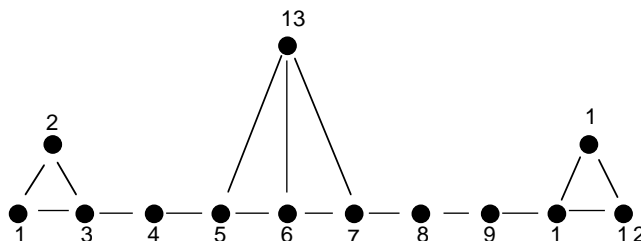


Figure 1

Though the single link classification hierarchy is optimal with respect to Lipschitz distance, it need not be the only Lipschitz nearest hierarchy. The following example shows this and also

illustrates why some intuitively appealing clusterings still have no smaller Lipschitz distance to the original metric.

Let M be the metric space defined by the network shown above in Figure 1. Here the short edges are of length 1 and the longer ones, i.e. the ones to point 13, have length 2. The distance between two points is the length of the shortest path connecting them. For this space $dsq(M) = 9$. Figure 2

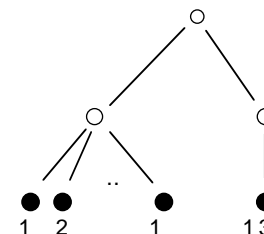


Figure 2

gives the single link hierarchy clustering for this space. Two other clusterings are in Figures 3 and 4. All three of these clusterings have Lipschitz distance $\log(9)$ to the original metric. The reason that this is the case in the case of the clusterings of Figures 4 and 5 is that, though the clusters themselves are perhaps better, one pays a high price for breaking up nearest pairs like 4, 5 and 8, 9 in Figure 3 and 4, 5 and 7, 8 in Figure 4.

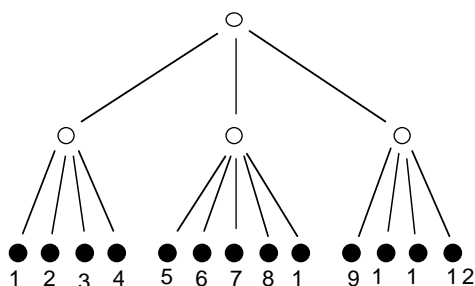


Figure 3

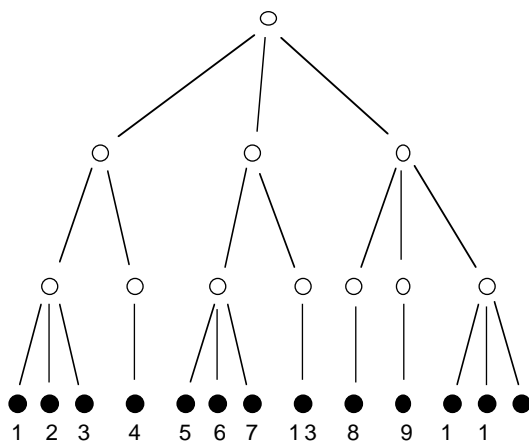


Figure 4

References.

[Hartigan, 1967 #57; Jardine, 1967 #123; Johnson, 1967 #73; Johnson, 1987 #114; Matousek, 1990 #13; Jardine, 1971 #70]

1. N. Jardine, R.Sibson, *Mathematical taxonomy*, Wiley, 1971.
2. W.B. Johnson, J. Lindenstrauss, G.Schechtman, *On Lipschitz embeddings of finite metric spaces in low dimensional normed spaces*, In: J Lindenstrauss, V.D. Millman (eds), *Geometrical aspects of functional analysis*, Springer, 1987.
3. J. Matousek, *Bi-Lipschitz embeddings into low dimensional Euclidean spaces*, *Comm. Math. Univ. Carolinae* 31:3(1990), 589-600.