

Journal of Classification

Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion? --Manuscript Draft--

Manuscript Number:	CLAS276R1
Full Title:	Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion?
Article Type:	Original Research
Keywords:	Hierarchical clustering, Ward, Lance-Williams, minimum variance, statistical software.
Corresponding Author:	Fionn Murtagh UNITED KINGDOM
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	
Corresponding Author's Secondary Institution:	
First Author:	Fionn Murtagh
First Author Secondary Information:	
Order of Authors:	Fionn Murtagh Pierre Legendre
Order of Authors Secondary Information:	
Abstract:	The Ward error sum of squares hierarchical clustering method has been very widely used since its first description by Ward in a 1963 publication. It has also been generalized in various ways. Two algorithms are found in the literature and software, both announcing that they implement the Ward clustering method. When applied to the same distance matrix, they produce different results. One algorithm preserves Ward's criterion, the other does not. Our survey work and case studies will be useful for all those involved in developing software for data analysis using Ward's hierarchical clustering method.
Response to Reviewers:	Please see accompanying letter where we respond to each point raised.

Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion?

Fionn Murtagh (1) and Pierre Legendre (2)

(1) Department of Computer Science,

Royal Holloway, University of London,

Egham TW20 0EX, UK (fmurtagh@acm.org)

(2) Département de sciences biologiques, Université de Montréal,

C.P. 6128 succursale Centre-ville, Montréal, Québec,

Canada H3C 3J7 (pierre.legendre@umontreal.ca)

23 December 2012

Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion?

December 23, 2012

Abstract

The Ward error sum of squares hierarchical clustering method has been very widely used since its first description by Ward in a 1963 publication. It has also been generalized in various ways. Two algorithms are found in the literature and software, both announcing that they implement the Ward clustering method. When applied to the same distance matrix, they produce different results. One algorithm preserves Ward's criterion, the other does not. Our survey work and case studies will be useful for all those involved in developing software for data analysis using Ward's hierarchical clustering method.

Keywords: Hierarchical clustering, Ward, Lance-Williams, minimum variance, statistical software.

1 Introduction

In the literature and in software packages there is confusion in regard to what is termed the Ward hierarchical clustering method. This relates to: (i) input dissimilarities, whether squared or not; and (ii) output dendrogram heights and whether or not their square root is used. Our main objective in this work is to warn users of hierarchical clustering about this, to raise awareness about these distinctions or differences, and to urge users to check what their favourite software package is doing.

In R, the function `hclust` of `stats` with the `method="ward"` option produces results that correspond to a Ward method (Ward¹, 1963) described in terms of

¹This article is dedicated to Joe H. Ward Jr., who died on 23 June 2011, aged 84.

a Lance-Williams updating formula using a sum of dissimilarities, which produces updated dissimilarities. This is the implementation used by, for example, Wishart (1969), Murtagh (1985) on whose code the `hclust` implementation is based, Jain and Dubes (1988), Jambu (1989), in the XploRe (2007) and Clustan (www.clustan.com) software packages, and elsewhere.

An important issue though is the form of input that is necessary to give Ward’s method. For an input data matrix, \mathbf{x} , in R’s `hclust` function the following command is required: `hclust(dist(x)^2,method="ward")` although this is not mentioned in the function’s documentation file. In later sections (in particular, section 4.2) of this article we explain just why the squaring of the distances is a requirement for the Ward method. In section 5 (Experiment 4) it is discussed why we may wish to take the square roots of the agglomeration, or dendrogram node height, values.

In R again, the `agnes` function of package `cluster` with the `method="ward"` option is also presented as the Ward method in Kaufman and Rousseeuw (1990) and in Legendre and Legendre (2012), among others. A formally similar algorithm is used, based on the Lance and Williams (1967) recurrence. The results of `agnes` differ from those of `hclust` when both functions are applied to the same distance matrix.

Lance and Williams (1967) did not themselves consider the Ward method for which the updating formula was first investigated by Wishart (1969).

What is at issue for us here starts with how `hclust` and `agnes` give different outputs when applied to the same dissimilarity matrix as input. What therefore explains the formal similarity in terms of criterion and algorithms, yet at the same time yields outputs that are different?

2 Applications

Ward’s is the only one among the agglomerative clustering methods that is based on a classical sum-of-squares criterion, producing groups that minimize within-group dispersion at each binary fusion.

In addition, Ward’s method is interesting because it looks for clusters in multivariate Euclidean space. That is also the reference space in multivariate ordination methods, and in particular in principal component analysis (PCA). We will show in section 3.2 that the total sum of squares in a data table can be computed from either the original variables or the distance matrix among observations, thus establishing the relationship between distances and sum of squares (or variance). This is why one can use an update formula based on dissimilarities to minimize the within-group sum of squares during Ward hierarchical clustering.

PCA is another way of representing the variance among observations, this time in an ordination diagram, which can be seen as a “spatial” representation of the relationships among the observations. PCA is a decomposition of the total variance of the data table, followed by selection of the axes that account for the largest portion of the variance; these axes are then used for representation of the

observations in a few dimensions, usually two. From this reasoning, we can see that spatial (e.g. PCA) and clustering (e.g. Ward's) methods involve different yet complementary spatial and clustering models that are fit to the data using the same mathematical principle. This is why in practice the results of Ward's agglomerative clustering are likely to delineate clusters that visually correspond to regions of high densities of points in PCA ordination.

Similar to use in conjunction with PCA, Ward's method is complementary to the use of correspondence analysis. The latter is a decomposition of the inertia of the data table. Ward's method accommodates weights on the observations. Ward's method applied to the output of a correspondence analysis, i.e. to the factor projections, implies equiweighted observations, endowed with the Euclidean distance. See Murtagh (2005) for many application domains involving the complementary use of correspondence analysis and Ward's method.

Ward's method can also be applied to dissimilarities other than the Euclidean distance. For these dissimilarities, ordinations can be produced by principal coordinate analysis (PCoA, Gower 1966), which is also called classical multidimensional scaling. When drawn onto a PCoA ordination diagram, the Ward clustering results often delineate clusters that visually correspond to the density centres in PCoA ordination.

Ward's method shares the total error sum of squares criterion with K -means partitioning, which is widely used to directly cluster observations in Euclidean space, hence to create a partition of the observation set. This is done without any structural constraint such as cluster embeddedness, represented by a hierarchy. Since K -means partitioning is an NP-hard problem, an approximate solution is often sought by using multiple random starts of the algorithm and retaining the solution that minimizes the total error sum of squares criterion. A more direct and computer-efficient approach is to first apply Ward's minimum variance agglomerative clustering to the data, identify the partition of the objects into K groups in the dendrogram, and then use that partition as the starting approximation for K -means partitioning, since it is close to the solution that one is seeking. That solution can then be improved by iterations of the K -means algorithm.

3 Ward's Agglomerative Hierarchical Clustering Method

3.1 Some Definitions

We recall that a distance is a positive, definite, symmetric mapping of a pair of observation vectors onto the positive reals which in addition satisfies the triangular inequality. For observations i, j, k we have: $d(i, j) > 0$; $d(i, j) = 0 \iff i = j$; $d(i, j) = d(j, i)$; $d(i, j) \leq d(i, k) + d(k, j)$. For an observation set, I , with $i, j, k \in I$ we can write the distance as a mapping from the Cartesian product of the observation set into the positive reals: $d : I \times I \longrightarrow \mathbb{R}^+$.

A dissimilarity is usually taken as a distance but without the triangular

inequality property $(d(i, j) \leq d(i, k) + d(k, j), \forall i, j, k)$. Lance and Williams (1967) use the term “an (i, j) -measure” for a dissimilarity.

An ultrametric, or tree distance, which defines a hierarchical clustering (and also an ultrametric topology, which goes beyond a metric geometry, or a p-adic number system) differs from a distance in that the strong triangular inequality is instead satisfied. This inequality, also commonly called the ultrametric inequality, is: $d(i, j) \leq \max\{d(i, k), d(k, j)\}$.

For observations i in a cluster q , and a distance d (which can potentially be relaxed to a dissimilarity) we have the following definitions. We may want to consider a mass or weight associated with observation i , $p(i)$. Typically we take $p(i) = 1/|q|$ when $i \in q$, i.e. 1 over cluster cardinality of the relevant cluster.

With the context being clear, let q denote the cluster (a set) and q^* the cluster’s centre. We have this centre defined as $q^* = 1/|q| \sum_{i \in q} i$. Furthermore, and again the context makes this clear, we have i used for the observation label, or index, among all observations, and the observation vector.

Some further definitions follow.

- Error sum of squares: $\sum_{i \in q} d^2(i, q^*)$.
- Variance (or centred sum of squares): $1/|q| \sum_{i \in q} d^2(i, q^*)$.
- Inertia: $\sum_{i \in q} p(i) d^2(i, q^*)$ which becomes variance if $p(i) = 1/|q|$, and becomes error sum of squares if $p(i) = 1$.
- Euclidean distance squared using norm $\|\cdot\|$: if $i, i' \in \mathbb{R}^{|J|}$, i.e. these observations have values on attributes $j \in \{1, 2, \dots, |J|\}$, J is the attribute set, $|\cdot|$ denotes cardinality, then $d^2(i, i') = \|i - i'\|^2 = \sum_j (i_j - i'_j)^2$.

Consider now a set of masses, or weights, m_i for observations i . Following Benzécri (1976, p. 185), the centred moment of order 2, $M^2(I)$ of the cloud (or set) of observations $i, i \in I$, is written: $M^2(I) = \sum_{i \in I} m_i \|i - g\|^2$ where the centre of gravity of the system is $g = \sum_i m_i i / \sum_i m_i$. The variance, $V^2(I)$, is $V^2(I) = M^2(I)/m_I$, where m_I is the total mass of the cloud. Due to Huyghen’s theorem the following can be shown (Benzécri, 1976, p. 186) for clusters q whose union make up the partition, Q :

$$M^2(Q) = \sum_{q \in Q} m_q \|q^* - g\|^2$$

$$M^2(I) = M^2(Q) + \sum_{q \in Q} M^2(q)$$

$$V(Q) = \sum_{q \in Q} \frac{m_q}{m_I} \|q^* - g\|^2$$

$$V(I) = V(Q) + \sum_{q \in Q} \frac{m_q}{m_I} V(q)$$

The $V(Q)$ and $V(I)$ definitions here are discussed in Jambu (1978, pp. 154–155). The last of the above can be seen to decompose (additively) total variance of the cloud I into (first term on the right hand side) variance of the cloud of cluster centres ($q \in Q$), and summed variances of the clusters. We can consider this last of the above relations as total variance decomposed into the sum of between and within cluster variances, or the sum of inter and intra cluster variances. This relationship will be important below.

A range of variants of the agglomerative clustering criterion and algorithm are discussed by Jambu (1978). These include: minimum of the centred order 2 moment of the union of two clusters (p. 156); minimum variance of the union of two clusters (p. 156); maximum of the centred order 2 moment of a partition (p. 157); and maximum of the centred order 2 moment of a partition (p. 158). Jambu notes that these criteria for maximization of the centred order 2 moment, or variance, of a partition, were developed and used by numerous authors, with some of these authors introducing modifications (such as the use of particular metrics). Among authors referred to are Ward (1963), Orlóci (1967), Wishart (1969), and various others.

3.2 Alternative Expressions for the Variance

As already noted in section 3.1, the cluster centre, i.e. cluster mean, q^* , is: $q^* = \frac{1}{|q|} \sum_{i \in q} i$. In the previous section, the variance was written as $1/|q| \sum_{i \in q} d^2(i, q^*)$. This is the so-called population variance. When viewed in statistical terms, where an unbiased estimator of the variance is needed, we require the sample variance: $1/(|q| - 1) \sum_{i \in q} d^2(i, q^*)$. The population quantity is used in Murtagh (1985). The sample statistic is used in Le Roux and Rouanet (2004), and by Legendre and Legendre (2012).

The sum of squares, $\sum_{i \in q} d^2(i, q^*)$, can be written in terms of all pairwise distances:

$$\begin{aligned} \sum_{i \in q} d^2(i, q^*) &= 1/|q| \sum_{i, i' \in q, i < i'} d^2(i, i'). \text{ This is proved as follows (see, e.g.,} \\ &\text{Legendre and Fortin, 2010). Write} \\ \frac{1}{|q|} \sum_{i, i' \in q, i < i'} d^2(i, i') &= \frac{1}{|q|} \sum_{i, i' \in q, i < i'} (i - i')^2 \\ &= \frac{1}{|q|} \sum_{i, i' \in q, i < i'} (i - q^* - (i' - q^*))^2 \\ &= \frac{1}{|q|} \sum_{i, i' \in q, i < i'} ((i - q^*)^2 + (i' - q^*)^2 - 2(i - q^*)(i' - q^*)) \\ &= \frac{1}{2} \frac{1}{|q|} \sum_{i \in q} \sum_{i' \in q} ((i - q^*)^2 + (i' - q^*)^2 - 2(i - q^*)(i' - q^*)) \\ &= \frac{1}{2} \frac{1}{|q|} \left(2|q| \sum_{i \in q} (i - q^*)^2 \right) - \frac{1}{2} \frac{1}{|q|} \left(\sum_{i \in q} \sum_{i' \in q} 2(i - q^*)(i' - q^*) \right) \end{aligned}$$

By writing out the right hand term, we see that it equals 0. Hence our result.

As noted in Legendre and Legendre (2012) there are many other alternative expressions for calculating $\sum_{i \in q} d^2(i, q^*)$, such as using the trace of a particular transformation of the distance matrix, and the sum of eigenvalues of a principal coordinate analysis of the distance matrix. The latter is invoking what is known as the Parseval relation, i.e. the equivalence of the norms of vectors in inner product spaces that can be orthonormally transformed, one space to the other.

3.3 Lance-Williams Dissimilarity Update Formula

Lance and Williams (1967) established a succinct form for the update of dissimilarities following an agglomeration. The parameters used in the update formula are dependent on the cluster criterion value. Consider clusters (including possibly singletons) i and j being agglomerated to form cluster $i \cup j$, and then consider redefining the dissimilarity relative to an external cluster (including again possibly a singleton), k . We have:

$$d(i \cup j, k) = a(i) \cdot d(i, k) + a(j) \cdot d(j, k) + b \cdot d(i, j) + c \cdot |d(i, k) - d(j, k)|$$

where d is the dissimilarity used – which does not have to be a Euclidean distance to start with, insofar as the Lance and Williams formula can be used as a repeatedly executed recurrence, without reference to any other or separate criterion; coefficients $a(i), a(j), b, c$ are defined with reference to the clustering criterion used (see tables of these coefficients in Murtagh, 1985, p. 68; Jambu, 1989, p. 366); and $|\cdot|$ denotes absolute value.

The Lance-Williams recurrence formula considers dissimilarities and not dissimilarities squared.

The original Lance and Williams (1967) paper did not consider the Ward criterion. It did however note that it allowed one to “generate an infinite set of new strategies” for agglomerative hierarchical clustering. Wishart (1969) brought the Ward criterion into the Lance-Williams algorithmic framework.

Even starting the agglomerative process with a Euclidean distance will not avoid the fact that the inter-cluster (non-singleton, i.e. with 2 or more members) dissimilarity does not respect the triangular inequality, and hence it does not respect this Euclidean metric property.

3.4 Generalizing Lance-Williams

The Lance and Williams recurrence formula has been generalized in various ways. See e.g. Batagelj (1988) who discusses what he terms “generalized Ward clustering” which includes agglomerative criteria based on variance, inertia and weighted increase in variance.

Jambu (1989, pp. 356 et seq.) considers the following cluster criteria and associated Lance-Williams update formula in the generalized Ward framework: centred order 2 moment of a partition; variance of a partition; centred order 2 moment of the union of two classes; and variance of the union of two classes.

When using a Euclidean distance, the Murtagh (1985) and the Jambu (1989) Lance-Williams update formulas for variance and related criteria (as discussed by Jambu, 1989) are associated with an alternative agglomerative hierarchical clustering algorithm which defines cluster centres following each agglomeration, and thus does not require use of the Lance-Williams update formula. The same is true for hierarchical agglomerative clustering based on median and centroid criteria.

As noted, the Lance-Williams update formula uses a dissimilarity, d . Székely and Rizzo (2005) consider higher order powers of this, in the Ward context: “Our proposed method extends Ward’s minimum variance method. Ward’s method minimizes the increase in total within-cluster sum of squared error. This increase is proportional to the squared Euclidean distance between cluster centres. In contrast to Ward’s method, our cluster distance is based on Euclidean distance, rather than squared Euclidean distance. More generally, we define ... an objective function and cluster distance in terms of any power α of Euclidean distance in the interval $(0,2]$... Ward’s minimum variance method is obtained as the special case when $\alpha = 2$.”

Then the authors indicate what beneficial properties the case of $\alpha = 1$ has, including: Lance-Williams form, ultrametricity and reducibility, space-dilation, and computational tractability. In Székely and Rizzo (2005, p. 164) it is stated that “We have shown that” the $\alpha = 1$ case, rather than $\alpha = 2$, gives “a method that applies to a more general class of clustering problems”, and this finding is further emphasized in their conclusion. Notwithstanding this finding of Székely and Rizzo (2005), viz. that the $\alpha = 1$ case is best, in this work our interest remains with the $\alpha = 2$ Ward method.

Our objective in this section has been to discuss some of the ways that the Ward method has been generalized. After this methodological review, we will, in the next section, come to our central theme in this article.

4 Implementations of Ward’s Method

We now come to the central part of our work, distinguishing in subsections 4.2 and 4.3 how we can arrive at subtle but important differences in relation to how the Ward method, or what is said to be the Ward method, is understood in practice, and put into software code. We consider: data inputs, the main loop structure of the agglomerative dissimilarity-based algorithms, and the output dendrogram node heights. The subtle but important differences that we uncover are further explored and exemplified in section 5.

Consider hierarchical clustering in the following form. On an observation set, I , define a dissimilarity measure. Set each of the observations, i, j, k , etc. $\in I$ to be a singleton cluster. Agglomerate the closest (i.e. least dissimilar) pair of clusters, deleting the agglomerands, or agglomerated clusters. Redefine the inter-cluster dissimilarities with respect to the newly created cluster. If n is the cardinality of observation set I then this agglomerative hierarchical clustering algorithm completes in $n - 1$ agglomerative steps.

Through use of the Lance-Williams update formula, we will focus on the updated dissimilarities relative to a newly created cluster. Unexpectedly in this work, we found a need to focus also on the form of input dissimilarities.

4.1 The Minimand or Cluster Criterion Optimized

The function to be minimized, or minimand, in the Ward2 case (see subsection 4.3), as stated by Kaufman and Rousseeuw (1990, p. 230, cf. relation (22)) is:

$$D(c_1, c_2) = \delta^2(c_1, c_2) = \frac{|c_1||c_2|}{|c_1| + |c_2|} \|c_1 - c_2\|^2 \quad (1)$$

(a term that measures the change in total sum of squares resulting from the fusion of c_1 and c_2) whereas for the Ward1 case, as discussed in subsection 4.2, we have:

$$\delta(c_1, c_2) = \frac{|c_1||c_2|}{|c_1| + |c_2|} \|c_1 - c_2\|^2 \quad (2)$$

It is clear therefore that the same criterion is being optimized. Both implementations minimize the change in variance, or the error sum of squares.

Error sum of squares, or minimum variance, or other related criteria are NP-complete optimization problems. This implies that a polynomial bound on computational complexity is not possible. Only exponential search in the solution space, because it is exhaustive, will guarantee an optimal solution. Because the error sum of squares and other related criteria are not optimized precisely in reality, we are content with good heuristics in practice, i.e. sub-optimal solutions. Such a heuristic is the sequence of two-way agglomerations carried out by a hierarchical clustering algorithm.

In either form the criterion (1), (2) is characterized in Le Roux and Rouanet (2004, p. 109) as either the variance index; the inertia index; the centred moment of order 2; or the Ward index (citing Ward, 1963). In the sequel, given two classes c_1 and c_2 , the variance index is defined as the contribution of the dipole of the class centre, denoted as in (2). The resulting clustering is termed a Euclidean classification by Le Roux and Rouanet (2004).

As noted by Le Roux and Rouanet (2004, p. 110), the variance index (as they term it) (2) is not guaranteed to satisfy the triangular inequality.

4.2 Implementation of Ward: Ward1

We start with (let us term it) the Ward1 algorithm as described in Murtagh (1985).

It was initially Wishart (1969) who wrote the Ward algorithm in terms of the Lance-Williams update formula. In Wishart (1969) the Lance-Williams formula is written in terms of squared dissimilarities, in a way that is formally identical to the following.

Cluster update formula:

$$\delta(i \cup i', i'') = \frac{w_i + w_i''}{w_i + w_{i'} + w_{i''}} \delta(i, i'') + \frac{w_i' + w_i''}{w_i + w_{i'} + w_{i''}} \delta(i', i'') - \frac{w_i''}{w_i + w_{i'} + w_{i''}} \delta(i, i')$$

$$\text{and } w_{i \cup i'} = w_i + w_{i'} \quad (3)$$

For the minimand (relation 1) of section 4.1, the input dissimilarities need to be as follows: $\delta(i, i') = \sum_j (x_{ij} - x_{i'j})^2$. Note the presence of the *squared Euclidean distance* in this initial dissimilarity specification. This is the Ward algorithm of Murtagh (1983, 1985 and 2000), and the way that `hclust` in R, which is based on that algorithm, ought to be used. When, however, the Ward1 algorithm is used with Euclidean distances as the initial dissimilarities, then the clustering topology can be very different, as will be seen in Section 5.

The weight w_i is the cluster cardinality, and thus for a singleton, $w_i = 1$. An immediate generalization is to consider probabilities given by $w_i = 1/n$. Generalization to arbitrary weights can also be considered. Ward implementations that take observation weights into account are available in Murtagh (2000).

$\frac{1}{2} \sum_j (x_{ij} - x_{i'j})^2$, i.e. 0.5 times Euclidean distances squared, is the population variance (cf. section 3.1) of the new cluster comprising two singletons, $i \cup i'$. To see this, note that the variance of the new cluster c formed by merging c_1 and c_2 is $(|c_1| \|c_1 - c\|^2 + |c_2| \|c_2 - c\|^2) / (|c_1| + |c_2|)$ where $|c_1|$ is both the cardinality and the mass of cluster c_1 , and $\|\cdot\|$ is the Euclidean norm. The new cluster's centre of gravity, or mean, is $c = \frac{|c_1|c_1 + |c_2|c_2}{|c_1| + |c_2|}$, where c_1 and c_2 are vectors of original data or means. By using this expression for the new cluster's centre of gravity (or mean) in the expression given for the variance, we see that we can write the variance of the new cluster c , combining c_1 and c_2 , to be $\frac{|c_1||c_2|}{|c_1| + |c_2|} \|c_1 - c_2\|^2$. So when $|c_1| = |c_2|$ we have the stated result, i.e. the variance of the new cluster equaling 0.5 times Euclidean distances squared.

The criterion that is optimized arises from the foregoing discussion (previous paragraph), i.e. the variance of the dipole formed by the agglomerands. This is the variance of new cluster c minus the variances of (now agglomerated) clusters c_1 and c_2 , which we can write as $\text{Var}(c) - \text{Var}(c_1) - \text{Var}(c_2)$. The variance of the partition containing c necessarily decreases, so we need to minimize this decrease when carrying out an agglomeration.

Murtagh (1985) also shows how this optimization criterion is viewed as achieving the (next possible) partition with maximum between-cluster variance. Maximizing between-cluster variance is the same as minimizing within-cluster variance, arising out of Huyghen's variance (or inertia) decomposition theorem. With reference to section 3.1 we are minimizing the change in B , hence maximizing B , and hence minimizing W .

Jambu (1978, p. 157) calls the Ward1 algorithm the maximum centred order 2 moment of a partition (cf. section 3.1 above). The criterion is denoted by him as δ_{mot} .

4.3 Implementation of Ward: Ward2

We now look at the Ward2 algorithm described in Kaufman and Rousseeuw (1990) and Legendre and Legendre (2012).

At each agglomerative step, the extra sum of squares caused by agglomerating clusters is minimized, exactly as we have seen for the Ward1 algorithm

above. We have the following.
Cluster update formula:

$$\begin{aligned} \delta(i \cup i', i'') = \\ \left(\frac{w_i + w_i''}{w_i + w_{i'} + w_{i''}} \delta^2(i, i'') + \frac{w_{i'} + w_i''}{w_i + w_{i'} + w_{i''}} \delta^2(i', i'') - \frac{w_i''}{w_i + w_{i'} + w_{i''}} \delta^2(i, i') \right)^{1/2} \\ \text{and } w_{i \cup i'} = w_i + w_{i'} \end{aligned} \quad (4)$$

Contrary to Ward1, the input dissimilarities are Euclidean distances (not squared). They are squared within equation (4): $\delta^2(i, i') = \sum_j (x_{ij} - x_{i'j})^2$. It is such squared Euclidean distances that interest us, since our motivation arises from the error sum of squares criterion.

A second point to note is that equation (4) relates to, on the right hand side, *the square root of a weighted sum of squared distances*. Consider how in equation (3) the cluster update formula was in terms of *a weighted sum of distances*. As a consequence, function δ is not the same in equations (3) and (4): the function produces a squared distance in the former and a distance in the latter.

A final point about equation (4) is that in the cluster update formula it is the set of δ values that we seek.

Now let us look further at the relationship between equations (4) and (3), and show their relationship. Rewriting the cluster update formula (4) after squaring both sides establishes that we have:

$$\begin{aligned} \delta^2(i \cup i', i'') = \\ \frac{w_i + w_i''}{w_i + w_{i'} + w_{i''}} \delta^2(i, i'') + \frac{w_{i'} + w_i''}{w_i + w_{i'} + w_{i''}} \delta^2(i', i'') - \frac{w_i''}{w_i + w_{i'} + w_{i''}} \delta^2(i, i') \end{aligned} \quad (5)$$

Let us use the notation $D = \delta^2$ because then, with

$$\begin{aligned} D(i \cup i', i'') = \\ \frac{w_i + w_i''}{w_i + w_{i'} + w_{i''}} D(i, i'') + \frac{w_{i'} + w_i''}{w_i + w_{i'} + w_{i''}} D(i', i'') - \frac{w_i''}{w_i + w_{i'} + w_{i''}} D(i, i') \end{aligned} \quad (6)$$

we see exactly the form of the Lance-Williams cluster update formula (section 3.3).

Although the agglomerative clustering algorithm is not fully specified as such in Cailliez and Pagès (1976), it appears that the Ward2 algorithm is the one attributed to Ward (1963). See their criterion d_9 (Cailliez and Pagès, 1976, pp. 531, 571).

With the appropriate choice of δ , different for Ward1 and for Ward2, what we have here is the identity of the algorithms Ward1 and Ward2, although they are implemented to a small extent differently. We show this as follows.

Take the Ward2 algorithm one step further than above, and write the input dissimilarities and cluster update formula using $D(i, i') = \delta^2(i, i') = \sum_j (x_{ij} - x_{i'j})^2$. Square both sides of equation (4). The cluster update formula is now:

$$\begin{aligned}
D(i \cup i', i'') = & \\
& \frac{w_i + w_i''}{w_i + w_{i'} + w_{i''}} D(i, i'') + \frac{w_i' + w_i''}{w_i + w_{i'} + w_{i''}} D(i', i'') - \frac{w_i''}{w_i + w_{i'} + w_{i''}} D(i, i') \\
& \text{and } w_{i \cup i'} = w_i + w_{i'}
\end{aligned} \tag{7}$$

In this form, equation (7), implementation Ward2 (equation (4)) is *identical* to implementation Ward1 (equation (3)). We conclude that we *can* have Ward1 and Ward2 implementations such that the outputs are identical.

5 Case Studies: Ward Implementations and Their Relationships

The hierarchical clustering programs used in this set of case studies are:

- `hclust` in package `stats`, “The R Stats Package”, in R version 2.15. Based on code by F. Murtagh (Murtagh, 1985), included in R by Ross Ihaka and Fritz Leisch.
- `agnes` in package `cluster`, “Cluster Analysis Extended Rousseeuw et al.”, in R, by L. Kaufman and P.J. Rousseeuw.
- `hclust.PL`, an extended version of `hclust` in R, by P. Legendre. In this function, the Ward1 algorithm is implemented by `method="ward.D"` and the Ward2 algorithm by `method="ward.D2"`. That function is available on the web page numerationecology.com, in the section on R-language functions.

We ensure reproducible results by providing all code used and, to begin with, by generating an input data set as follows.

```
# Fix the seed of the random number generator in order
# to have reproducible results.
set.seed(19037561)
# Create the input matrix to be used.
y <- matrix(runif(20*4),nrow=20,ncol=4)
# Look at overall mean and column standard deviations.
mean(y); sd(y)
0.4920503      # mean
0.2778538 0.3091678 0.2452009 0.2918480 # standard deviations
```

5.1 Experiment 1: agnes and Ward2 Implementation, hclust.PL

In experiment 1, non-squared distances are used as input. The R code used is shown in the following, with output produced. In all of these experiments, we used the dendrogram node heights, associated with the agglomeration criterion values, in order to quickly show numerical equivalences. This is then followed up with displays of the dendrograms.

```
# EXPERIMENT 1 -----
X.hclustPL.wardD2 = hclust.PL(dist(y),method="ward.D2")
X.agnes.wardD2 = agnes(dist(y),method="ward")

sort(X.hclustPL.wardD2$height)
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163

sort(X.agnes.wardD2$height)
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

This points to: `hclust.PL` with the `method="ward.D2"` option being identical to: `agnes` with the `method="ward"` option. This shows that `agnes` implements the Ward2 algorithm.

Figure 1 displays the outcome, and we see the same visual result in both cases. That is, the two dendrograms are identical except for inconsequential swiveling of nodes. In group theory terminology we say that the trees are wreath product invariant.

To fully complete our reproducibility of research agenda, this is the code used to produce Figure 1:

```
par(mfrow=c(1,2))
plot(X.hclustPL.wardD2,main="X.hclustPL.wardD2",sub="",xlab="")
plot(X.agnes.wardD2,which.plots=2,main="X.agnes.wardD2",sub="",xlab="")
```

5.2 Experiment 2: hclust and Ward1 Implementation, hclust.PL

In this experiment, squared distances are used as input. The code used is as follows, with output shown.

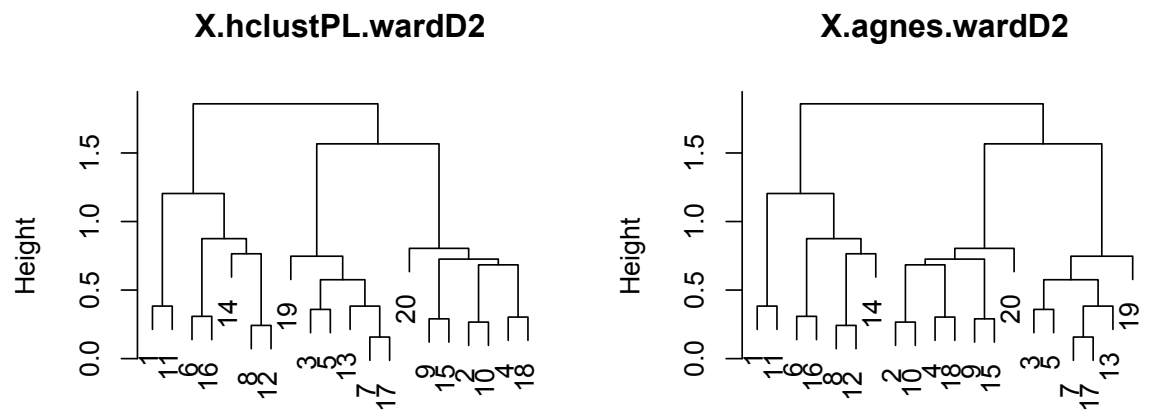


Figure 1: Experiment 1 outcome. The two dendrograms are morphologically identical.

```
# EXPERIMENT 2 -----
X.hclust = hclust(dist(y)^2, method="ward")
X.hclustPL.sq.wardD = hclust.PL(dist(y)^2, method="ward.D")

sort(X.hclust$height)
0.02477046 0.05866380 0.07097546 0.08420102 0.09184743
0.09510249 0.12883390 0.14671052 0.14684403 0.33106478
0.46791879 0.52680768 0.55799612 0.58483318 0.64677705
0.76584542 1.45043423 2.45393902 3.45371103

sort(X.hclustPL.sq.wardD$height)
0.02477046 0.05866380 0.07097546 0.08420102 0.09184743
0.09510249 0.12883390 0.14671052 0.14684403 0.33106478
0.46791879 0.52680768 0.55799612 0.58483318 0.64677705
0.76584542 1.45043423 2.45393902 3.45371103
```

This points to: `hclust`, with "ward" option, on squared input being identical to: `hclust.PL` with `method="ward.D"` option, on squared input.

The clustering levels shown here in Experiment 2 are the squares of the clustering levels produced by Experiment 1.

Figure 2 displays the outcome, and we see the same visual result in both cases. This is the code used to produce Figure 2:

```
par(mfrow=c(1,2))
plot(X.hclust, main="X.hclust",sub="",xlab="")
plot(X.hclustPL.sq.wardD, main="X.hclustPL.sq.wardD",sub="",xlab="")
```

5.3 Experiment 3: Non-Ward Result Produced by `hclust` and `hclust.PL`

In this experiment, with non-squared distances used as input, we achieve a well-defined hierarchical clustering, but one that differs from Ward. Code used is as follows, with output shown.

```
# EXPERIMENT 3 -----
X.hclustPL.wardD = hclust.PL(dist(y),method="ward.D")
X.hclust.nosq = hclust(dist(y),method="ward")

sort(X.hclustPL.wardD$height)
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3832023 0.4018957 0.5988721
0.7443850 0.7915592 0.7985444 0.8016877 0.8414950
0.9273739 1.4676446 2.2073106 2.5687307
```

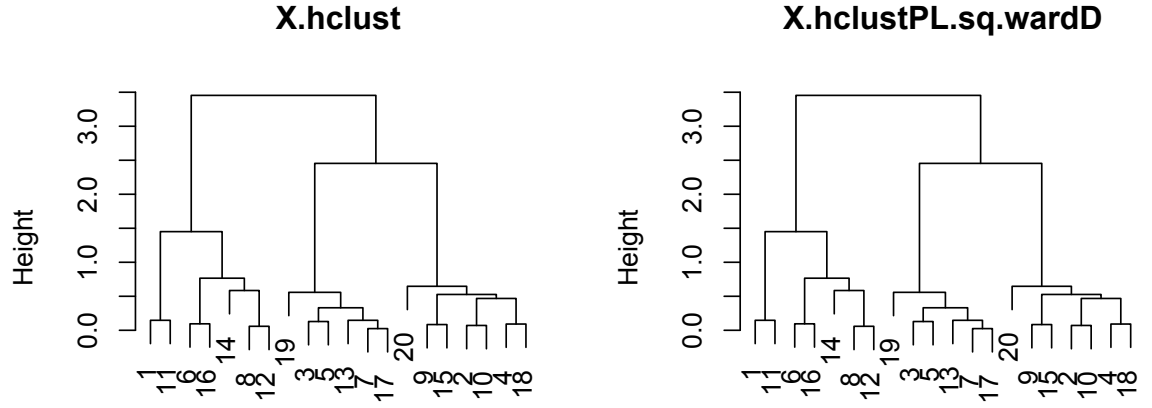



Figure 2: Experiment 2 outcome.

```
sort(X.hclust.nosq$height)
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3832023 0.4018957 0.5988721
0.7443850 0.7915592 0.7985444 0.8016877 0.8414950
0.9273739 1.4676446 2.2073106 2.5687307
```

This points to: `hclustPL` with `method="wardD"` option being the same as: `hclust` with `method="ward"` option. It shows that `hclust` with the “ward” option, used with non-squared distances as recommended in the R documentation file, implements the Ward1 algorithm.

Note: there is no squaring of inputs in the latter, nor in the former either. The clustering levels produced in this experiment using non-squared distances as input differ from, and are not monotonic relative to, those produced in Experiments 1 and 2.

Figure 3 displays the outcome, and we see the same visual result in both cases. This is the code used to produce Figure 3:

```
par(mfrow=c(1,2))
plot(X.hclustPL.wardD, main="X.hclustPL.wardD",sub="",xlab="")
plot(X.hclust.nosq, main="X.hclust.nosq",sub="",xlab="")
```

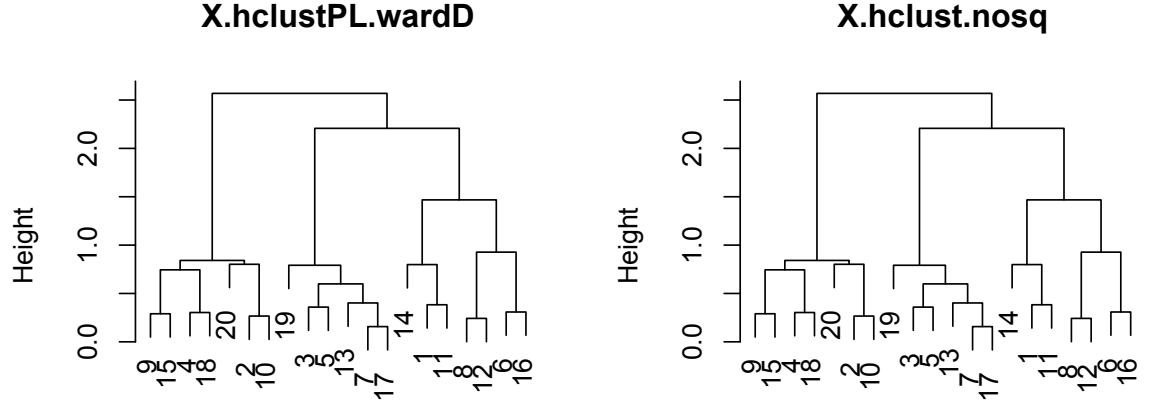


Figure 3: Experiment 3 outcome.

5.4 Experiment 4: Modifying Inputs and Options so that Ward1 Output is Identical to Ward2 Output

In this experiment, given the formal equivalences of the Ward1 and Ward2 implementations in sections 4.2 and 4.3, we show how to bring about identical output. We do this by squaring or not squaring input dissimilarities, and by playing on the options used.

```
> # EXPERIMENT 4 -----
X.hclust = hclust(dist(y)^2, method="ward")
X.hierclustPL.sq.wardD = hclust.PL(dist(y)^2, method="ward.D")
X.hclustPL.wardD2 = hclust.PL(dist(y), method="ward.D2")
X.agnes.wardD2 = agnes(dist(y),method="ward")
```

We will ensure that the node heights in the tree are in “distance” terms, i.e. in terms of the initial, unsquared Euclidean distances as used in this article. Of course, the agglomerations redefine such distances to be dissimilarities. Thus it is with unsquared dissimilarities that we are concerned.

While these dissimilarities are inter-cluster measures, defined in any given partition, the inter-node measures that are defined on the tree are ultrametric.

```
sort(sqrt(X.hclust$height))
```

```
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

```
sort(sqrt(X.hierclustPL.sq.wardD$height))
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

```
sort(X.hclustPL.wardD2$height)
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

```
sort(X.agnes.wardD2$height)
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

There is no difference of course between sorting the square roots of the agglomeration or height levels, versus sorting them and then taking their square roots. Consider the following examples, the first repeated from the foregoing (Experiment 4) batch of results.

```
sort(sqrt(X.hierclustPL.sq.wardD$height))
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

```
sqrt(sort(X.hierclustPL.sq.wardD$height))
0.1573864 0.2422061 0.2664122 0.2901741 0.3030634
0.3083869 0.3589344 0.3830281 0.3832023 0.5753823
0.6840459 0.7258152 0.7469914 0.7647439 0.8042245
0.8751259 1.2043397 1.5665054 1.8584163
```

Our Experiment 4 points to:

output of `hclustPL`, with the `method="ward.D2"` option
and

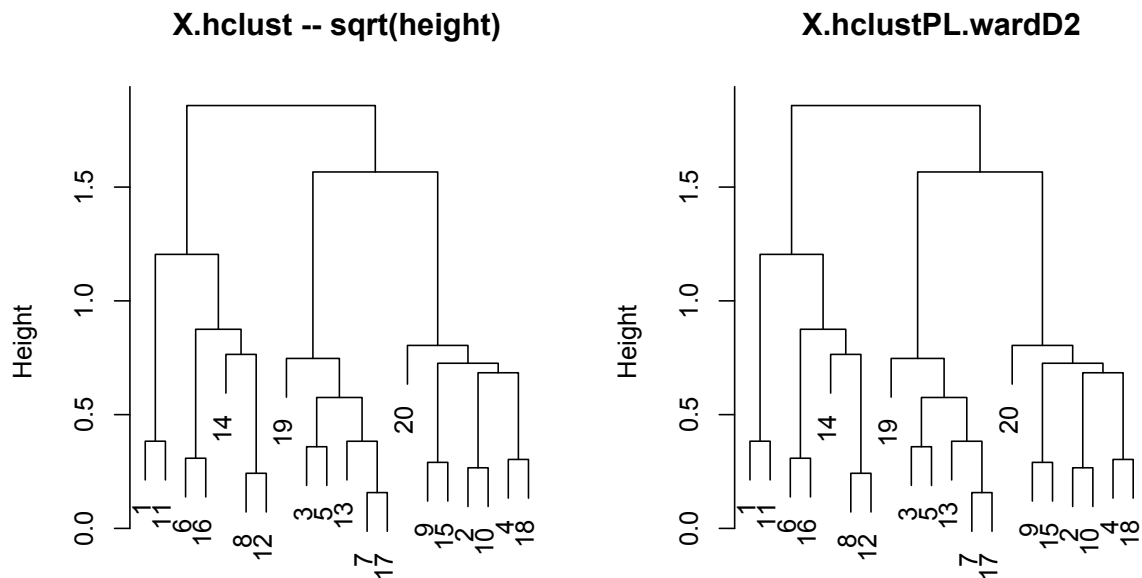


Figure 4: Experiment 4 outcome.

output of `agnes`, with the `method="ward"` option

being the same as both of the following with node heights square rooted:

`hclust`, with the `"ward"` option on squared input,

`hclust.PL`, with the `method="ward.D"` option on squared input.

Figure 4 displays two of these outcomes, and we see the same visual result in both cases, in line with the numerical node (or agglomeration) “height” values. This is the code used to produce Figure 4:

```
par(mfrow=c(1,2))
temp <- X.hclust
temp$height <- sqrt(X.hclust$height)
plot(temp, main="X.hclust -- sqrt(height)", sub="", xlab="")
plot(X.hclustPL.wardD2, main="X.hclustPL.wardD2", sub="", xlab="")
```

6 Discussion

6.1 Where Ward1 and Ward2 Implementations Lead to an Identical Result

A short discussion follows on the implications of this work. In Experiments 1 and 2, we see the crucial importance of inputs (squared or not) and options used. We set out, with Experiment 1, to implement Ward2. With Experiment 2, we set out to implement Ward1. Experiment 3 shows that Ward1 algorithms produce well-defined but non-Ward hierarchical clustering results. Finally Experiment 4 shows the underlying equivalence of the Experiment 1 and Experiment 2 results, i.e. respectively the Ward2 and Ward1 implementations.

On looking closely at the Experiment 1 and Experiment 2 dendrograms, Figures 1 and 2, we can see that the morphology of the dendrograms is the same. However the cluster criterion values – the node heights – are not the same.

From section 4.3, Ward2 implementation, the cluster criterion value is most naturally the square root of the same cluster criterion value as used in section 4.2, Ward1 implementation. From a dendrogram morphology viewpoint, this is not important because one morphology is the same as the other (as we have observed above). From an optimization viewpoint (section 4.1), it plays no role either since one optimand is monotonically related to the other.

6.2 How and Why the Ward1 and Ward2 Implementations Can Differ

Those were reasons as to why it makes no difference to choose the Ward1 implementation versus the Ward2 implementation, as long as squared distances are used as input to Ward1. Next, we will look at some practical differences.

Looking closer at forms of the criterion in (1) and (2) in section 4.1 – and contrasting these forms of the criterion with the input dissimilarities in sections 4.2 (Ward1) and 4.3 (Ward2) leads us to the following observation. The Ward2 criterion values are “on a scale of distances” whereas the Ward1 criterion values are “on a scale of distances squared”. Hence to make direct comparisons between the ultrametric distances read off a dendrogram, and compare them to the input distances, it is preferable to use the Ward2 form of the criterion. Thus, the use of cophenetic correlations can be more directly related to the dendrogram produced. Alternatively, with the Ward1 form of the criterion applied to squared distances, we can just take the square root of the dendrogram node heights. This we have seen in the generation of Figure 4.

6.3 Other Software: A Look at Six Other Packages

Readers may legitimately wonder which of the two Ward algorithms is implemented in various commercial statistical software packages. We ran examples in several packages; the selected data sets produced dendrogram topologies that

differed between the two algorithms. We found the following results at the time the final version of this paper was written (December 2012):

- Packages Statistica and Systat implement the Ward1 algorithm.
- Packages Matlab, SAS and JMP implement the Ward2 algorithm.
- The SPSS package implements the Ward1 algorithm but warns users that, for a correct Ward clustering, squared Euclidean distances should be used instead of Euclidean distances. Indeed, this produces the correct dendrogram topology, as shown in this paper. The fusion levels are, however, squared distances.
- As we have seen throughout this article, in the R language, function `hclust()` of package `stats` implements the Ward1 algorithm whereas function `agnes()` of package `cluster` implements Ward2.

6.4 Other Implementations Based on the Lance-Williams Update Formula

The algorithms Ward1 and Ward2 can be used for “stored dissimilarity” and “stored data” implementations, a distinction first made in Anderberg (1973). The latter is where the dissimilarity matrix is not used, but instead the dissimilarities are computed on the fly.

Murtagh (2000) has implementations of the “stored dissimilarity” (programs `hc.f`, `HCL.java`, see Murtagh, 2000) as well as the “stored data” (`hcon2.f`) algorithms. For both, Murtagh (1985) lists the formulas. The nearest neighbour and reciprocal nearest neighbour algorithms can be applied to bypass the need for a strict sequencing of the agglomerations. See Murtagh (1983, 1985). These algorithms provide for provably worst case $O(n^2)$ implementations, as first introduced in de Rham and Juan, and published in, respectively, 1980 and 1982. Cluster criteria such as Ward’s method must respect Bruynooghe’s reducibility property (Bruynooghe, 1977) if they are to be reversal-free or inversion-free (or with monotonic variation in cluster criterion value through the sequence of agglomerations). Apart from computational reasons, the other major advantage of such algorithms (nearest neighbour chain, reciprocal nearest neighbour) is use in distributed computing (including virtual memory) environments (Murtagh, 1992).

7 Conclusions

Having different yet very close implementations that differ by just a few lines of code (in any high level language), yet claiming to implement a given method, is confusing for the learner, for the practitioner and even for the specialist. In this work, we have first of all reviewed all relevant background. Then we have laid out in very clear terms the two, differing implementations. Additionally,

with differing inputs, and with somewhat different processing driven by options set by the user, in fact our two different implementations had the appearance of being quite different methods.

The two algorithms at issue here only differ in terms of values in one algorithm being squared relative to the other algorithm (or, clearly equivalently, a square root of terms in one case relative to the other case). An upshot of this is that there is no difference from the point of view of computational scaling, i.e. order of magnitude computational complexity. For the nearest neighbour chain algorithm using either “stored data” or “stored dissimilarities” implementations, the algorithms discussed by us here are of $O(n^2)$ computational complexity for the clustering of n observations. Thus the current implementation of `hclust()` in package `stats` has complexity equal to $O(n^2)$, as can be shown by numerical simulations.

Two algorithms, Ward1 and Ward2, are found in the literature and software, both announcing that they implement the Ward (1963) clustering method. When applied to the same distance matrix, they produce different results. This article has shown that when they are applied to the same dissimilarity matrix, only Ward2 minimizes the Ward clustering criterion and produces the Ward method. The Ward1 and Ward2 algorithms can be made to optimize the same criterion and produce the same clustering topology by using Ward1 with squared distances, and Ward2 with the distances themselves. Furthermore, taking the square root of the clustering levels produced by Ward1 used with squared distances produces the same clustering levels as Ward2 used with the distances themselves. The constrained clustering package of Legendre (2011), `const.clust`, derived from `hclust` in R, offers both the Ward1 and Ward2 options.

We have shown in this article how close these two implementations are, in fact. Furthermore we discussed in detail what the implications are for the few, differing lines of code. Software developers who only offer the Ward1 algorithm are encouraged to explain clearly how the Ward2 output is to be obtained, as described in the previous paragraph.

Acknowledgements

We are grateful to the following colleagues who ran example data sets in statistical packages and sent us the results: Guy Cucumel, Pedro Peres-Neto and Yves Prairie. Our thanks also to representatives of Statistica, Systat and SAS who provided information on the Ward algorithm implemented in their package.

References

1. ANDERBERG, M.R. (1973), *Cluster Analysis for Applications*, Academic.
2. BATAGELJ, V. (1988), Generalized Ward and related clustering problems, in H.H. Bock, ed., *Classification and Related Methods of Data Analysis*, North-Holland, pp. 67–74.

3. BENZÉCRI, J.P. (1976), *L'Analyse des Données, Tome 1, La Taxinomie*, Dunod, (2nd edn.; 1973, 1st edn.).
4. BRUYNOOGHE, M. (1977), "Méthodes nouvelles en classification automatique des données taxinomiques nombreuses", *Statistique et Analyse des Données*, no. 3, 24–42.
5. CAILLIEZ, F. and PAGÈS, J.-P. (1976), *Introduction à l'Analyse des Données*, SMASH (Société de Mathématiques Appliquées et Sciences Humaines).
6. FISHER, R.A. (1936), The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7, 179–188.
7. GOWER, J.C. (1966), Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika*, 53, 325–338.
8. JAIN, A.K. and DUBES, R.C. (1988), *Algorithms for Clustering Data*, Prentice-Hall.
9. JAMBU, M. (1978), *Classification Automatique pour l'Analyse des Données. I. Méthodes et Algorithmes*, Dunod.
10. JAMBU, M. (1989), *Exploration Informatique et Statistique des Données*, Dunod.
11. KAUFMAN, L. and ROUSSEEUW, P.J. (1990), *Finding Groups in Data: An Introduction to Cluster Analysis*, Wiley.
12. LANCE, G.N. and WILLIAMS, W.T. (1967). A general theory of classificatory sorting strategies. 1. Hierarchical systems, *Computer Journal*, 9, 4, 373–380.
13. LEGENDRE, P. and FORTIN, M.-J. (2010), Comparison of the Mantel test and alternative approaches for detecting complex relationships in the spatial analysis of genetic data. *Molecular Ecology Resources*, 10, 831–844.
14. LEGENDRE, P. (2011), `const.clust`, R package to compute space-constrained or time-constrained agglomerative clustering.
<http://www.bio.umontreal.ca/legendre/indexEn.html>
15. LEGENDRE, P. and LEGENDRE, L. (2012), *Numerical Ecology*, 3rd English ed., Elsevier.
16. LE ROUX, B. and ROUANET, H. (2004), *Geometric Data Analysis: From Correspondence Analysis to Structured Data Analysis*, Kluwer.
17. MURTAGH, F. (1983), A survey of recent advances in hierarchical clustering algorithms, *The Computer Journal*, 26, 354–359.

18. MURTAGH, F. (1985), *Multidimensional Clustering Algorithms*, Physica-Verlag.
19. MURTAGH, F. (1992), Comments on: Parallel algorithms for hierarchical clustering and cluster validity, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14, 1056–1057.
20. MURTAGH, F. (2000), Multivariate data analysis software and resources, <http://www.classification-society.org/csna/mda-sw>
21. MURTAGH, F. (2005), *Correspondence Analysis and Data Coding with R and Java*, Chapman & Hall/CRC.
22. ORLÓCI, L. (1967), An agglomerative method for classification of plant communities, *Journal of Ecology*, 55, 193–206.
23. SZÉKELY, G.J. and RIZZO, M.L. (2005), Hierarchical clustering via joint between-within distances: extending Ward’s minimum variance method, *Journal of Classification*, 22 (2), 151–183.
24. WARD, J.H. (1963), Hierarchical grouping to optimize an objective function, *Journal of the American Statistical Association*, 58, 236–244.
25. WISHART, D. (1969), An algorithm for hierarchical classifications, *Biometrics* 25, 165–170.
26. XploRe (2007), version 4.8, Collaborative Research Center 649, Humboldt-Universität zu Berlin, Germany. http://sfb649.wiwi.hu-berlin.de/fedc_homepage/xplore.php