

# Clustering in Massive Data Sets

Fionn Murtagh  
School of Computer Science  
The Queen's University of Belfast  
Belfast BT7 1NN, Northern Ireland  
f.murtagh@qub.ac.uk

August 22, 1999

## Abstract

We review the time and storage costs of search and clustering algorithms. We exemplify these, based on case-studies in astronomy, information retrieval, visual user interfaces, chemical databases, and other areas. Sections 2 to 6 relate to nearest neighbor searching, an elemental form of clustering, and a basis for clustering algorithms to follow. Sections 7 to 11 review a number of families of clustering algorithm. Sections 12 to 14 relate to visual or image representations of data sets, from which a number of interesting algorithmic developments arise.

*'Now', said Rabbit, 'this is a Search, and I've Organised it -'  
'Done what to it?' said Pooh.  
'Organised it. Which means - well, it's what you do to a Search,  
when you don't all look in the same place at once.'  
A.A. Milne, The House at Pooh Corner (1928) - M.S. Zakaria*

Topics:

1. Introduction
2. Binning or Bucketing
3. Multidimensional Binary Search or kD Tree
4. Projections and Other Bounds
5. The Special Case of Sparse Binary Data
6. Fast Approximate Nearest Neighbor Finding
7. Hierarchical Agglomerative Clustering
8. Graph Clustering
9. Nearest Neighbor Finding on Graphs
10. K-Means and Family
11. Fast Model-Based Clustering
12. Noise Modeling
13. Cluster-Based User Interfaces
14. Images from Data
15. Conclusion

# 1 Introduction

Nearest neighbor searching is considered in sections 2 to 6 for one main reason: its utility for the clustering algorithms reviewed in sections 7 to 10 and, partially, 11. They are the building blocks for the most efficient implementations of hierarchical clustering algorithms, and they can be used to speed up other families of clustering algorithms.

Sections 12 to 14 will deal with facets of visual or image representations of data sets.

The best match or nearest neighbor problem is important in many disciplines. In statistics,  $k$ -nearest neighbors, where  $k$  can be 1, 2, etc., is a method of non-parametric discriminant analysis. In pattern recognition, this is a widely-used method for unsupervised classification (see Dasarathy, 1991). Nearest neighbor algorithms are the building block of clustering algorithms based on nearest neighbor chains; or of effective heuristic solutions for combinatorial optimization algorithms such as the traveling salesman problem, which is a paradigmatic problem in many areas. In the database and more particularly data mining field, NN searching is called similarity query, or similarity join (Bennett et al., 1999).

In section 2, we begin with data structures where the objective is to break the  $O(n)$  barrier for determining the nearest neighbor (NN) of a point. A database record or tuple may be taken as a point in a space of dimensionality  $m$ , the latter being the associated number of fields or attributes. These approaches have been very successful, but they are restricted to low dimensional NN-searching. For higher dimensional data, a wide range of bounding approaches have been proposed, which remain  $O(n)$  algorithms but with a low constant of proportionality.

We assume familiarity with basic notions of similarity and distance, the triangular inequality, ultrametric spaces, Jaccard and other coefficients, normalization and standardization. For an implicit treatment of data theory and data coding, see Murtagh and Heck (1987). Useful background reading can be found in Arabie et al. (1996). In particular output representational models include discrete structures, e.g. rooted labeled trees or dendrograms, and spatial structures (Arabie and Hubert, 1996), with many hybrids.

## 2 Binning or Bucketing

In this approach to NN-searching, a preprocessing stage precedes the searching stage. All points are mapped onto indexed cellular regions of space, so that NNs are found in the same or in closely adjacent cells. Taking the plane as an example, and considering points  $(x_i, y_i)$ , the maximum and minimum values on all coordinates are obtained (e.g.  $(x_j^{\min}, y_j^{\min})$ ). Consider the mapping (Fig. 1)

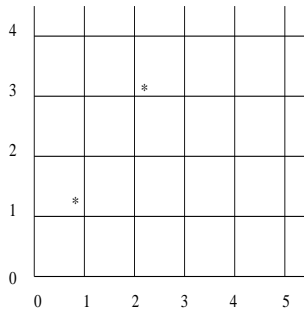
$$x_i \rightarrow \{[(x_{ij} - x_j^{\min})/r]\}$$

where constant  $r$  is chosen in terms of the number of equally spaced categories into which the interval  $[x_j^{\min}, x_j^{\max}]$  is to be divided. This gives to  $x_i$  an integer value between 0 and  $[(x_{ij}^{\max} - x_{ij}^{\min})/r]$  for each attribute  $j$ .  $O(nm)$  time is required to obtain the transformation of all  $n$  points, and the result may be stored as a linked list with a pointer from each cell identifier to the set of points mapped onto that cell.

NN-searching begins by finding the closest point among those which have been mapped onto the same grid cell as the target point. This gives a current NN point. A closer point may be mapped onto some other grid cell if the distance between target point and current NN point is greater than the distance between the target point and any of the boundaries of the cell containing it. Some further implementation details can be found in Murtagh (1993a).

A powerful theoretical result regarding this approach is as follows. For uniformly distributed points, the NN of a point is found in  $O(1)$ , or constant, expected time (see Delannoy, 1980, or Bentley et al., 1980, for proof). Therefore this approach will work well if approximate uniformity can be assumed or if the data can be broken down into regions of approximately uniformly distributed points.

$$x^{\min}, y^{\min} = 0, 0, x^{\max}, y^{\max} = 50, 40, r = 10$$



Point (22,32) is mapped onto cell (2,3); point (8,13) is mapped onto cell (0,1).

Figure 1: Example of simple binning in the plane.

Simple Fortran code for this approach is listed, and discussed, in Schreiber (1993). The search through adjacent cells requires time which increases exponentially with dimensionality (if it is assumed that the number of points assigned to each cell is approximately equal). As a result, this approach is suitable for low dimensions only. Rohlf (1978) reports on work in dimensions 2, 3, and 4; and Murtagh (1983) in the plane. Rohlf also mentions the use of the first 3 principal components to approximate a set of points in 15-dimensional space.

From the constant expected time NN search result, particular hierarchical agglomerative clustering methods can be shown to be of linear expected time,  $O(n)$  (Murtagh, 1983). The expected time complexity for Ward's minimum variance method is given as  $O(n \log n)$ . Results on the hierarchical clustering of up to 12,000 points are discussed.

The limitation on these very appealing computational complexity results is that they are only really feasible for data in the plane. Bellman's curse of dimensionality manifests itself here as always. For dimensions greater than 2 or 3 we proceed to the situation where a binary search tree can provide us with a good preprocessing of our data.

### 3 Multidimensional Binary Search or kD Tree

A binary search tree preprocesses the data to be searched through by two-way subdivision, and subdivisions continue until some prespecified number of data points is arrived at. See example in Fig. 2. We associate with each node of the decision tree the definition of a subdivision of the data only, and we associate with each terminal node a pointer to the stored coordinates of the points. Using the approximate median of projections keeps the tree balanced, and consequently  $O(\log n)$  levels, at each of which  $O(n)$  processing is required. Hence the construction of the tree takes  $O(n \log n)$  time.

The search for a NN then proceeds by a top-down traversal of the tree. The target point is transmitted through successive levels of the tree using the defined separation of the two child nodes at each node. On arrival at a terminal node, all associated points are examined and a current NN selected. The tree is then backtracked: if the points associated with any node could furnish a closer point, then subnodes must be checked out.

The approximately constant number of points associated with terminal nodes (hyper-rectangular cells in the space of points) should be greater than 1 in order that some NNs may be obtained without requiring a search of adjacent cells (other terminal nodes). Friedman et al. (1977) suggest a value of the number of points per bin between 4 and 32 based on empirical study.

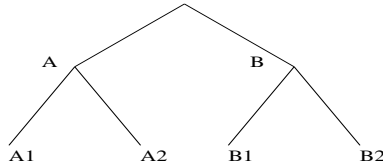
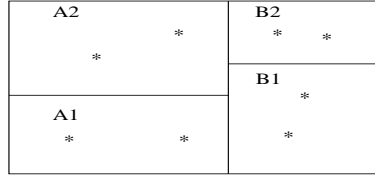


Figure 2: A MDBST using planar data.

The MDBST approach only works well with small dimensions. To see this, consider each coordinate being used once and once only for the subdivision of points, i.e. each attribute is considered equally useful. Let there be  $p$  levels in the tree, i.e.  $2^p$  terminal nodes. Each terminal node contains approximately  $c$  points by construction and so  $c2^p = n$ . Therefore  $p = \log_2 n/c$ . As sample values, if  $n = 32768$ ,  $c = 32$ , then  $p = 10$ . I.e. in 10-dimensional space, using a large number of points associated with terminal nodes, more than 30000 points will need to be considered. For high dimensional spaces, two alternative MDBST specifications are as follows.

All attributes need not be considered for splitting the data if it is known that some are of greater interest than others. Linearity present in the data may manifest itself via the variance of projections of points on the coordinates; choosing the coordinate with greatest variance as the discriminator coordinate at each node may therefore allow repeated use of certain attributes. This has the added effect that the hyper-rectangular cells into which the terminal nodes divide the space will be approximately cubical in shape. In this case, Friedman et al. (1977) show that search time is  $O(\log n)$  on average for the finding of a NN. Results obtained for dimensionalities of between 2 and 8 are reported on in Friedman et al. (1977), and in the application of this approach to minimal spanning tree construction in Bentley and Friedman (1978). Lisp code for the MDBST is discussed in Broder (1990).

The MDBST has also been proposed for very high dimensionality spaces, i.e. where the dimensionality may be greater than the number of points, as could be the case in a keyword-based system. Keywords (coordinates) are batched, and the following decision rule is used: if some *one* of a given batch of node-defining discriminating attributes is present, then take the left subtree, else take the right subtree. Large  $n$ , well in excess of 1400, was stated as necessary for good results (Weiss, 1981; Eastman and Weis, 1982). General guidelines for the attributes which define the direction of search at each level are that they be related, and the number chosen should keep the tree balanced. On intuitive grounds, our opinion is that this approach will work well if the clusters of attributes, defining the tree nodes, are mutually well separated.

An MDBST approach is used by Moore (1999) in the case of Gaussian mixture clustering. Over and above the search for nearest neighbors based on Euclidean distance, Moore allows for the Mahalanobis metric, i.e. distance to cluster centers which are “corrected” for the (Gaussian) spread or morphology of clusters. The information stored at each node of the tree includes covariances. Moore (1999) reports results on numbers of objects of around 160,000, dimensionalities of between 2 and 6, and speedups of 8-fold to 1000-fold. Pelleg and Moore (1999) discuss results on some 430,000 two-dimensional objects from the Sloan Digital Sky Survey (see section 10 below).

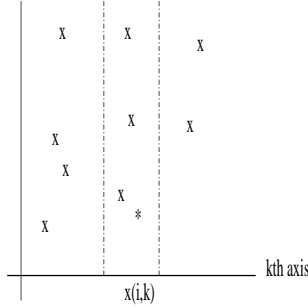


Figure 3: Two-dimensional example of projection-based bound. Points with projections within distance  $c$  of given point's (\*) projection, alone, are searched. Distance  $c$  is defined with reference to a candidate or current nearest neighbor.

## 4 Projections and Other Bounds

### 4.1 Bounding using Projection or Properties of Metrics

Making use of bounds is a versatile approach, which may be less restricted by dimensionality. Some lower bound on the dissimilarity is efficiently calculated in order to dispense with the full calculation in many instances.

Using projections on a coordinate axis allows the exclusion of points in the search for the NN of point  $x_i$ . Points  $x_k$ , only, are considered such that  $(x_{ij} - x_{kj})^2 \leq c^2$  where  $x_{ij}$  is the  $j$ th coordinate of  $x_i$ , and where  $c$  is some prespecified distance (see Fig. 3).

Alternatively more than one coordinate may be used. The prior sorting of coordinate values on the chosen axis or axes expedites the finding of points whose full distance calculation is necessitated. The preprocessing required with this approach involves the sorting of up to  $m$  sets of coordinates, i.e.  $O(mn \log n)$  time.

Using one axis, it is evident that many points may be excluded if the dimensionality is very small, but that the approach will disimprove as the latter grows. Friedman et al. (1975) give the expected NN search time, under the assumption that the points are uniformly distributed, as  $O(mn^{1-1/n})$ . This approaches the brute force  $O(nm)$  as  $n$  gets large. Reported empirical results are for dimensions 2 to 8.

Marimont and Shapiro (1979) extend this approach by the use of projections in subspaces of dimension greater than 1 (usually about  $m/2$  is suggested). This can be further improved if the subspace of the principal components is used. Dimensions up to 40 are examined.

The Euclidean distance is very widely used. Two other members of a family of Minkowski metric measures require less computation time to calculate, and they can be used to provide bounds on the Euclidean distance. We have:

$$d_1(x, x') \geq d_2(x, x') \geq d_\infty(x, x')$$

where  $d_1$  is the Hamming distance defined as  $\sum_j |x_j - x'_j|$ ; the Euclidean distance is given by the square root of  $\sum_j (x_j - x'_j)^2$ ; and the Chebyshev distance is defined as  $\max_j |x_j - x'_j|$ .

Kittler (1978) makes use of the following bounding strategy: *reject all points  $y$  such that  $d_1(x, y) \geq \sqrt{m}\delta$  where  $\delta$  is the current NN  $d_2$ -distance.* The more efficiently calculated  $d_1$ -distance may thus allow the rejection of many points (90% in 10-dimensional space is reported on by Kittler). Kittler's rule is obtained by noting that the greatest  $d_1$ -distance between  $x$  and  $x'$  is attained when

$$|x_j - x'_j|^2 = d_2^2(x, x')/m$$

for all coordinates,  $j$ . Hence  $d_1(x, x') = d_2(x, x')/\sqrt{m}$  is the greatest  $d_1$ -distance between  $x$  and

$x'$ . In the case of the rejection of point  $y$ , we then have:

$$d_1(x, y) \leq d_2(x, y) / \sqrt{m}$$

and since, by virtue of the rejection,

$$d_1(x, y) \geq \sqrt{m}\delta$$

it follows that  $\delta \leq d_2(x, y)$ .

Yunck (1976) presents a theoretical analysis for the similar use of the Chebyshev metric. Richetin et al. (1980) propose the use of both bounds. Using uniformly distributed points in dimensions 2 to 5, the latter reference reports the best outcome when the rule: *reject all  $y$  such that  $d_\infty(x, y) \geq \delta$*  precedes the rule based on the  $d_1$ -distance. Up to 80% reduction in CPU time is reported.

## 4.2 Bounding using the Triangular Inequality

The triangular inequality is satisfied by distances:  $d(x, y) \leq d(x, z) + d(z, y)$ , where  $x$ ,  $y$  and  $z$  are any three points. The use of a *reference point*,  $z$ , allows a full distance calculation between point  $x$ , whose NN is sought, and  $y$  to be avoided if

$$|d(y, z) - d(x, z)| \geq \delta$$

where  $\delta$  is the current NN distance. The set of all distances to the reference point are calculated and stored in a preprocessing step requiring  $O(n)$  time and  $O(n)$  space. The above cut-off rule is obtained by noting that if

$$d(x, y) \geq |d(y, z) - d(x, z)|$$

then, necessarily,  $d(x, y) \geq \delta$ . The former inequality above reduces to the triangular inequality irrespective of which of  $d(y, z)$  or  $d(x, z)$  is the greater.

The set of distances to the reference point,  $\{d(x, z) \mid x\}$ , may be sorted in the preprocessing stage. Since  $d(x, z)$  is fixed during the search for the NN of  $x$ , it follows that the cut-off rule will not then need to be applied in all cases.

The single reference point approach, due to Burkhard and Keller (1973), was generalized to multiple reference points by Shapiro (1977). The sorted list of distances to the first reference point,  $\{d(x, z_1) \mid x\}$ , is used as described above as a preliminary bound. Then the subsequent bounds are similarly employed to further reduce the points requiring a full distance calculation. The number and the choice of reference points to be used is dependent on the distributional characteristics of the data. Shapiro (1977) finds that reference points ought to be located away from groups of points. In 10-dimensional simulations, it was found that at best only 20% of full distance calculations were required (although this was very dependent on the choice of reference points).

Hodgson (1988) proposes the following bound, related to the training set of points,  $y$ , among which the NN of point  $x$  is sought. Determine in advance the NNs and their distances,  $d(y, NN(y))$  for all points in the training set. For point  $y$ , then consider  $\delta_y = \frac{1}{2}d(y, NN(y))$ . In seeking NN( $x$ ), and having at some time in the processing a candidate NN,  $y'$ , we can exclude all  $y$  from consideration if we find that  $d(x, y') \leq \delta_{y'}$ . In this case, we know that we are sufficiently close to  $y'$  that we cannot improve on it.

We return now to the choice of reference points: Vidal Ruiz (1986) proposes the storing of inter-point distances between the members of the training set. Given  $x$ , whose NN we require, some member of the training set is used as a reference point. Using the bounding approach based on the triangular inequality, described above, allows other training set members to be excluded from any possibility of being NN( $x$ ). Micó et al. (1992) and Ramasubramanian and Paliwal (1992) discuss further enhancements to this approach, focused especially on the storage requirements.

Fukunaga and Narendra (1975) make use of both a hierarchical decomposition of the data set (they employ repeatedly the k-means partitioning technique), and bounds based on the triangular

inequality. For each node in the decomposition tree, the center and maximum distance to the center of associated points (the “radius”) are determined. For 1000 points, 3 levels were used, with a division into 3 classes at each node.

All points associated with a non-terminal node can be rejected in the search for the NN of point  $x$  if the following rule (Rule 1) is not verified:

$$d(x, g) - r_g < \delta$$

where  $\delta$  is the current NN distance,  $g$  is the center of the cluster of points associated with the node, and  $r_g$  is the radius of this cluster. For a terminal node, which cannot be rejected on the basis of this rule, each associated point,  $y$ , can be tested for rejection using the following rule (Rule 2):

$$|d(x, g) - d(y, g)| \geq \delta.$$

These two rules are direct consequences of the triangular inequality.

A branch and bound algorithm can be implemented using these two rules. This involves determining some current NN (the bound) and subsequently branching out of a traversal path whenever the current NN cannot be bettered. Not being inherently limited by dimensionality, this approach appears particularly attractive for general purpose applications.

Other rejection rules are considered by Kamgar-Parsi and Kanal (1985). A simpler form of clustering is used in the variant of this algorithm proposed by Niemann and Goppert (1988). A shallow MDBST is used, followed by a variant on the branching and bounding described above.

Bennett et al. (1999) use the nearest neighbor problem as a means towards solving the Gaussian distribution mixture problem. They consider a preprocessing approach similar to Fukunaga and Narendra (1975) but with an important difference: to take better account of cluster structure in the data, the clusters are multivariate normal but not necessarily of diagonal covariance structure. Therefore very elliptical clusters are allowed. This in turn implies that a cluster radius is not of great benefit for establishing a bound on whether or not distances need to be calculated. Bennett et al. (1999) address this problem by seeking a stochastic guarantee on whether or not calculations can be excluded. Technically, however, such stochastic bounds are not easy to determine in a high dimensional space.

An interesting issue raised in Beyer et al. (1999) is discussed also by Bennett et al. (1996): if the ratio of the nearest and furthest neighbor distances converges in probability to 1 as the dimensionality increases, then is it meaningful to search for nearest neighbors? This issue is not all that different from saying that neighbors in an increasingly high dimensional space tend towards being equidistant. In section 5, we will look at approaches for handling particular classes of data of this type.

### 4.3 Fast Approximate Nearest Neighbor Finding

Kushilevitz et al. (1998), working in Euclidean and  $L_1$  spaces, propose fast approximate nearest neighbor searching, on the grounds that in systems for content-based image retrieval, approximate results are adequate. Projections are used to bound the search. Probability of successfully finding the nearest neighbor is traded off against time and space requirements.

## 5 The Special Case of Sparse Binary Data

“High-dimensional”, “sparse” and “binary” are the characteristics of keyword-based bibliographic data, with maybe values in excess of 10000 for both  $n$  and  $m$ . Such data is usually stored as list data structures, representing the mapping of documents onto index terms, or vice versa. Commercial document collections are usually searched using a Boolean search environment. Documents associated with particular terms are retrieved, and the intersection (AND), union (OR) or other

operations on such sets of documents are obtained. For efficiency, an *inverted file* which maps terms onto documents must be available for Boolean retrieval. The efficient NN algorithms, to be discussed, make use of both the document-term and the term-document files.

The usual algorithm for NN-searching considers each document in turn, calculates the distance with the given document, and updates the NN if appropriate. This algorithm is shown schematically in Fig. 4 (top). The inner loop is simply an expression of the fact that the distance or similarity will, in general, require  $O(m)$  calculation: examples of commonly used coefficients are the Jaccard similarity, and the Hamming ( $L_1$  Minkowski) distance.

If  $\bar{m}$  and  $\bar{n}$  are, respectively, the average numbers of terms associated with a document, and the average number of documents associated with a term, then an average complexity measure, over  $n$  searches, of this usual algorithm is  $O(n\bar{m})$ . It is assumed that advantage is taken of some packed form of storage in the inner loop (e.g. using linked lists).

Croft's algorithm (see Croft, 1977, and Fig. 4) is of worst case complexity  $O(nm^2)$ . However the number of terms associated with the document whose NN is required will often be quite small. The National Physical Laboratory test collection, for example, which was used in Murtagh (1982) has the following characteristics:  $n = 11429$ ,  $m = 7491$ ,  $\bar{m} = 19.9$ , and  $\bar{n} = 30.4$ . The outermost and innermost loops in Croft's algorithm use the document-term file. The center loop uses the term-document inverted file. An average complexity measure (more strictly, the time taken for best match search based on an average document with associated average terms) is seen to be  $O(\bar{n}\bar{m}^2)$ .

In the outermost loop of Croft's algorithm, there will eventually come about a situation where – if a document has not been thus far examined – the number of terms remaining for the given document do not permit the current NN document to be bettered. In this case, we can cut short the iterations of the outermost loop. The calculation of a bound, using the greatest possible number of terms which could be shared with a so-far unexamined document has been exploited by Smeaton and van Rijsbergen (1981) and by Murtagh (1982) in successive improvements on Croft's algorithm.

The complexity of all the above algorithms has been measured in terms of operations to be performed. In practice, however, the actual accessing of term or document information may be of far greater cost. The document-term and term-document files are ordinarily stored on direct access file storage because of their large sizes. The strategy used in Croft's algorithm, and in improvements on it, does not allow any viable approaches to batching together the records which are to be read successively, in order to improve accessing-related performance.

The Perry-Willett algorithm (see Perry and Willett, 1983) presents a simple but effective solution to the problem of costly I/O. It focuses on the calculation of the number of terms common to the given document  $x$  and each other document,  $y$ , in the document collection. This set of values is built up in a computationally efficient fashion.  $O(n)$  operations are subsequently required to determine the (dis)similarity, using another vector comprising the total numbers of terms associated with each document. Computation time (the same "average" measure as that used above) is  $O(\bar{n}\bar{m} + n)$ . We now turn attention to numbers of direct-access reads required.

In Croft's algorithm, all terms associated with the document whose NN is desired may be read in one read operation. Subsequently, we require  $\bar{n}\bar{m}$  reads, giving in all  $1 + \bar{n}\bar{m}$ . In the Perry-Willett algorithm, the outer loop again pertains to the one (given) document, and so all terms associated with this document can be read and stored. Subsequently,  $\bar{m}$  reads, i.e. the average number of terms, each of which demands a read of a set of documents, are required. This gives, in all,  $1 + \bar{m}$ . Since these reads are very much the costliest operation in practice, the Perry-Willett algorithm can be recommended for large values of  $n$  and  $m$ . Its general characteristics are that (i) it requires, as do all the algorithms discussed in this section, the availability of the inverted term-document file; and (ii) it requires in-memory storage of two vectors containing  $n$  integer values.



Usual algorithm:

```
Initialize current NN
For all documents in turn do:
... For all terms associated with the document do:
... .. Determine (dis)similarity
... Endfor
... Test against current NN
Endfor
```

Croft's algorithm:

```
Initialize current NN
For all terms associated with the given document do:
... For all documents associated with each term do:
... .. For all terms associated with a document do:
... .. .. Determine (dis)similarity
... .. .. Endfor
... .. Test against current NN
... Endfor
Endfor
```

Perry-Willett algorithm:

```
Initialize current NN
For all terms associated with the given document, i, do:
... For all documents, i', associated with each term, do:
... .. Increment location i' of counter vector
... Endfor
Endfor
```

Figure 4: Algorithms for NN-searching using high-dimensional sparse binary data.

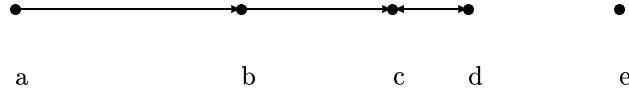


Figure 5: Five points, showing *NNs* and *RNNs*.

## 6 Hierarchical Agglomerative Clustering

The algorithms discussed in this section can be characterized as *greedy* (Horowitz and Sahni, 1979). A sequence of irreversible algorithm steps is used to construct the desired data structure.

We will not review hierarchical agglomerative clustering here. For essential background, the reader is referred to Murtagh and Heck (1987), Gordon (1999), or Jain and Dubes (1988). This section borrows on Murtagh (1992).

One could practically say that Sibson (1973) and Defays (1977) are part of the prehistory of clustering. Their  $O(n^2)$  implementations of the single link method and of a (non-unique) complete link method, respectively, have been widely cited.

In the early 1980s a range of significant improvements were made to the Lance-Williams, or related, dissimilarity update schema (de Rham, 1980; Juan, 1982), which had been in wide use since the mid-1960s. Murtagh (1985) presents a survey of these algorithmic improvements. We will briefly describe them here. The new algorithms, which have the potential for *exactly* replicating results found in the classical but more computationally expensive way, are based on the construction of *nearest neighbor chains* and *reciprocal* or mutual NNs (NN-chains and RNNs).

A NN-chain consists of an arbitrary point ( $a$  in Fig. 5); followed by its NN ( $b$  in Fig. 5); followed by the NN from among the remaining points ( $c$ ,  $d$ , and  $e$  in Fig. 5) of this second point; and so on until we necessarily have some pair of points which can be termed reciprocal or mutual NNs. (Such a pair of RNNs may be the first two points in the chain; and we have assumed that no two dissimilarities are equal.)

In constructing a NN-chain, irrespective of the starting point, we may agglomerate a pair of RNNs as soon as they are found. What guarantees that we can arrive at the same hierarchy as if we used traditional “stored dissimilarities” or “stored data” algorithms? Essentially this is the same condition as that under which no inversions or reversals are produced by the clustering method. Fig. 6 gives an example of this, where  $s$  is agglomerated at a lower criterion value (i.e. dissimilarity) than was the case at the previous agglomeration between  $q$  and  $r$ . Our ambient space has thus contracted because of the agglomeration. This is due to the algorithm used – in particular the agglomeration criterion – and it is something we would normally wish to avoid.

This is formulated as:

$$\text{Inversion impossible if: } d(i, j) < d(i, k) \text{ or } d(j, k) \Rightarrow d(i, j) < d(i \cup j, k)$$

This is essentially Bruynooghe’s *reducibility property* (Bruynooghe, 1977; see also Murtagh, 1984). Using the Lance–Williams dissimilarity update formula, it can be shown that the minimum variance method does not give rise to inversions; neither do the linkage methods; but the median and centroid methods cannot be guaranteed not to have inversions.

To return to Fig. 5, if we are dealing with a clustering criterion which precludes inversions, then  $c$  and  $d$  can justifiably be agglomerated, since no other point (for example,  $b$  or  $e$ ) could have been agglomerated to either of these.

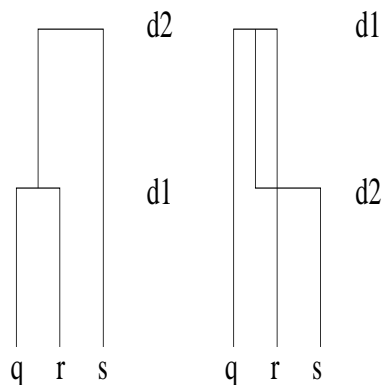


Figure 6: Alternative representations of a hierarchy with an inversion. Assuming dissimilarities, as we go vertically up, criterion values ( $d_1$ ,  $d_2$ ) decrease. But here, undesirably,  $d_2 > d_1$ .

The processing required, following an agglomeration, is to update the NNs of points such as  $b$  in Fig. 5 (and on account of such points, this algorithm was dubbed *algorithme des célibataires* in de Rham, 1980). The following is a summary of the algorithm:

#### NN-chain algorithm

**Step 1** Select a point arbitrarily.

**Step 2** Grow the NN-chain from this point until a pair of RNNs are obtained.

**Step 3** Agglomerate these points (replacing with a cluster point, or updating the dissimilarity matrix).

**Step 4** From the point which preceded the RNNs (or from any other arbitrary point if the first two points chosen in Steps 1 and 2 constituted a pair of RNNs), return to Step 2 until only one point remains.

In Murtagh (1983, 1984, 1985) and Day and Edelsbrunner (1984), one finds discussions of  $O(n^2)$  time and  $O(n)$  space implementations of Ward's minimum variance (or error sum of squares) method and of the centroid and median methods. The latter two methods are termed the UPGMC and WPGMC criteria by Sneath and Sokal (1973). Now, a problem with the cluster criteria used by these latter two methods is that the reducibility property is not satisfied by them. This means that the hierarchy constructed may not be unique as a result of inversions or reversals (non-monotonic variation) in the clustering criterion value determined in the sequence of agglomerations.

Murtagh (1984) describes  $O(n^2)$  time and  $O(n^2)$  space implementations for the single link method, the complete link method and for the weighted and unweighted group average methods (WPGMA and UPGMA). This approach is quite general vis à vis the dissimilarity used and can also be used for hierarchical clustering methods other than those mentioned.

Day and Edelsbrunner (1984) prove the exact  $O(n^2)$  time complexity of the centroid and median methods using an argument related to the combinatorial problem of optimally packing hyperspheres into an  $m$ -dimensional volume. They also address the question of metrics: results are valid in a wide class of distances including those associated with the Minkowski metrics.

The construction and maintenance of the nearest neighbor chain as well as the carrying out of agglomerations whenever reciprocal nearest neighbors meet, both offer possibilities for parallelization. Implementations on a SIMD machine were described by Willett (1989).

Evidently both coordinate data and graph (e.g., dissimilarity) data can be input to these agglomerative methods. Gillet et al. (1998) in the context of clustering chemical structure databases

refer to the common use of the Ward method, based on the reciprocal nearest neighbors algorithm, on data sets of a few hundred thousand molecules.

Applications of hierarchical clustering to bibliographic information retrieval are assessed in Griffiths et al. (1984). Ward's minimum variance criterion is favored.

From details in White and McCain (1997), the Institute of Scientific Information (ISI) clusters citations (science, and social science) by first clustering highly cited documents based on a single linkage criterion, and then four more passes are made through the data to create a subset of a single linkage hierarchical clustering.

## 7 Graph Clustering

Hierarchical clustering methods are closely related to graph-based clustering. For a start, a dendrogram is a rooted labeled tree. Secondly, and more importantly, some methods like the single and complete link methods can be displayed as graphs, and are very closely related to mainstream graph data structures.

An example of the increasing prevalence of graph clustering in the context of data mining on the web is presented in Fig. 7: Amazon.com provides information on what other books were purchased by like-minded individuals.

The single link method was referred to in the previous section, as a widely-used agglomerative, hence hierarchical, clustering method. Rohlf (1982) reviews algorithms for the single link method with complexities ranging from  $O(n \log n)$  to  $O(n^5)$ . The criterion used by the single link method for cluster formation is weak, meaning that noisy data in particular give rise to results which are not robust.

The minimal spanning tree (MST) and the single link agglomerative clustering method are closely related: the MST can be transformed irreversibly into the single link hierarchy (Rohlf, 1973). The MST is defined as of minimal total weight, it spans all nodes (vertices) and is an unrooted tree. The MST has been a method of choice for at least four decades now either in its own right for data analysis (Zahn, 1971), as a data structure to be approximated (e.g. using shortest spanning paths, see Murtagh, 1985, p. 96), or as a basis for clustering. We will look at some fast algorithms for the MST in the remainder of this section.

Perhaps the most basic MST algorithm, due to Prim and Dijkstra, grows a single fragment through  $n - 1$  steps. We find the closest vertex to an arbitrary vertex, calling these a fragment of the MST. We determine the closest vertex, not in the fragment, to any vertex in the fragment, and add this new vertex into the fragment. While there are fewer than  $n$  vertices in the fragment, we continue to grow it.

This algorithm leads to a unique solution. A default  $O(n^3)$  implementation is clear, and  $O(n^2)$  computational cost is possible (Murtagh, 1985, p. 98).

Sollin's algorithm constructs the fragments in parallel. For each fragment in turn, at any stage of the construction of the MST, determine its closest fragment. Merge these fragments, and update the list of fragments. A tree can be guaranteed in this algorithm (although care must be taken in cases of equal similarity) and our other requirements (all vertices included, minimal total edge weight) are very straightforward. Given the potential for roughly halving the data remaining to be processed at each step, not surprisingly the computational cost reduces from  $O(n^3)$  to  $O(n^2 \log n)$ .

The real interest of Sollin's algorithm arises when we are clustering on a graph and do not have all  $n(n - 1)/2$  edges present. Sollin's algorithm can be shown to have computational cost  $m \log n$  where  $m$  is the number of edges. When  $m \ll n(n - 1)/2$  then we have the potential for appreciable gains.

The MST in feature spaces can of course make use of the fast nearest neighbor finding methods studied earlier in this article. See Murtagh (1985, section 4.4) for various examples.

Other graph data structures which have been proposed for data analysis are related to the MST. We know, for example, that the following subset relationship holds:

$$\text{MST} \subseteq \text{RNG} \subseteq \text{GG} \subseteq \text{DT}$$

where RNG is the relative neighborhood graph, GG is the Gabriel graph, and DT is the Delaunay triangulation. The latter, in the form of its dual, the Voronoi diagram, has been used for analyzing the clustering of galaxy locations. References to these and related methods can be found in Murtagh (1993b).

## 8 Nearest Neighbor Finding on Graphs

Clustering on graphs may be required because we are working with (perhaps complex non-Euclidean) dissimilarities. In such cases where we must take into account an edge between each and every pair of vertices, we will generally have an  $O(m)$  computational cost where  $m$  is the number of edges. In a metric space we have seen that we can look for various possible ways to expedite the nearest neighbor search. An approach based on visualization – turning our data into an image – will be looked at below. However there is another aspect of our similarity (or other) graph which we may be able to turn to our advantage. Efficient algorithms for sparse graphs are available. Sparsity can be arranged – we can threshold our edges if the sparsity does not suggest itself more naturally. A special type of sparse graph is a planar graph, i.e. a graph capable of being represented in the plane without any crossovers of edges.

For sparse graphs, algorithms with  $O(m \log \log n)$  computational cost were described by Yao (1975) and Cheriton and Tarjan (1976). A short algorithmic description can be found in Murtagh (1985, pp. 107–108) and we refer in particular to the latter.

The basic idea is to preprocess the graph, in order to expedite the sorting of edge weights (why sorting? – simply because we must repeatedly find smallest links, and maintaining a sorted list of edges is a good basis for doing this). If we were to sort all edges, the computational requirement would be  $O(m \log m)$ . Instead of doing that, we take the edge set associated with each and every vertex. We divide each such edge set into groups of size  $k$ . (The fact that the last such group will usually be of size  $< k$  is taken into account when programming.)

Let  $n_v$  be the number of incident edges at vertex  $v$ , such that  $\sum_v n_v = 2m$ .

The sorting operation for each vertex now takes  $O(k \log k)$  operations for each group, and we have  $n_v/k$  groups. For all vertices the sorting requires a number of operations which is of the order of  $\sum_v n_v \log k = 2m \log k$ . This looks like a questionable – or small – improvement over  $O(m \log m)$ .

Determining the lightest edge incident on a vertex requires  $O(n_v/k)$  comparisons since we have to check all groups. Therefore the lightest edges incident on all vertices are found with  $O(m/k)$  operations.

When two vertices, and later fragments, are merged, their associated groups of edges are simply collected together, therefore keeping the total number of groups of edges which we started out with. We will bypass the issue of edges which, over time, are to be avoided because they connect vertices in the same fragment: given the fact that we are building an MST, the total number of such edges-to-be-avoided cannot surpass  $2m$ .

To find what to merge next, again  $O(m/k)$  processing is required. Using Sollin's algorithm, the total processing required in finding what to merge next is  $O(m/k \log n)$ . The total processing required for grouping the edges, and sorting within the edge-groups, is  $O(m \log k)$ , i.e. it is one-off and accomplished at the start of the MST-building process.

The total time is  $O(m/k \log n) + O(m \log k)$ . Let's fix  $k = \log n$  (aha!). Then the second term dominates and gives overall computational complexity as  $O(m \log \log n)$ .

This result has been further improved to near linearity in  $m$  by Gabow et al. (1985), who develop an algorithm with complexity  $O(m \log \log \log \dots n)$  where the number of iterated log terms is bounded by  $m/n$ .

Motwani and Raghavan (1995, chapter 10) base a stochastic  $O(m)$  algorithm for the MST on random sampling to identify and eliminate edges that are guaranteed not to belong to the MST.

Let's turn our attention now to the case of a planar graph. For a planar graph we know that  $m \leq 3n - 6$  for  $m > 1$ . (For proof, see for example Tucker, 1980, or any book on graph theory).

Referring to Sollin's algorithm, described above,  $O(n)$  operations are needed to establish a least cost edge from each vertex, since there are only  $O(n)$  edges present. On the next round, following fragment-creation, there will be at most  $\text{ceil}(n/2)$  new vertices, implying of the order of  $n/2$  processing to find the least cost edge. The total computational cost is seen to be proportional to:  $n + n/2 + n/4 + \dots = O(n)$ .

So determining the MST of a planar graph is linear in numbers of either vertices or edges.

Before ending this review of very efficient clustering algorithms for graphs, we note that algorithms discussed so far have assumed that the similarity graph was undirected. For modeling transport flows, or economic transfers, the graph could well be directed. Components can be defined, generalizing the clusters of the single link method, or the complete link method. Tarjan (1983) provides an algorithm for the latter agglomerative criterion which is of computational cost  $O(m \log n)$ .

## 9 K-Means and Family

The non-technical person more often than not understands clustering as a partition. K-means looked at in this section, or the distribution mixture approach looked at section 10, provide solutions.

A mathematical definition of a partition implies no multiple assignments of observations to clusters, i.e. no overlapping clusters. Overlapping clusters may be faster to determine in practice, and a case in point is the one-pass algorithm described in Salton and McGill (1983). The general principle followed is: make one pass through the data, assigning each object to the first cluster which is close enough, and making a new cluster for objects that are not close enough to any existing cluster.

Broder et al. (1997) use this algorithm for clustering the web. A feature vector is determined for each HTML document considered, based on sequences of words. Similarity between documents is based on an inverted list, using an approach like those described in section 5. The similarity graph is thresholded, and components sought.

Broder (1998) solves the same clustering objective using a thresholding and overlapping clustering method similar to the Salton and McGill one. The application described is that of clustering the Altavista repository in April 1996, consisting of 30 million HTML and text documents, comprising 150 GBytes of data. The number of serviceable clusters found was 1.5 million, containing 7 million documents. Processing time was about 10.5 days. An analysis of the clustering algorithm used by Broder can be found in Borodin et al. (1999), who also considers the use of approximate minimal spanning trees.

The threshold-based pass of the data, in its basic state, is susceptible to lack of robustness. A bad choice of threshold leads to too many clusters or two few. To remedy this, we can work on a well-defined data structure such as the minimal spanning tree. Or, alternatively, we can iteratively refine the clustering. Partitioning methods, such as k-means, use iterative improvement of an initial estimation of a targeted clustering.

A very widely used family of methods for inducing a partition on a data set is called k-means, c-means (in the fuzzy case), Isodata, competitive learning, vector quantization and other more general names (non-overlapping non-hierarchical clustering) or more specific names (minimal distance or exchange algorithms).

The usual criterion to be optimized is:

$$\frac{1}{|I|} \sum_{q \in Q} \sum_{i \in q} \|\vec{i} - \vec{q}\|^2$$

where  $I$  is the object set,  $|\cdot|$  denotes cardinality,  $q$  is some cluster,  $Q$  is the partition, and  $q$  denotes a set in the summation, whereas  $\vec{q}$  denotes some associated vector in the error term, or metric norm. This criterion ensures that clusters found are compact, and therefore assumed homogeneous. The optimization criterion, by a small abuse of terminology, is often referred to as a minimum variance one.

A necessary condition that this criterion be optimized is that vector  $\vec{q}$  be a cluster mean, which for the Euclidean metric case is:

$$\vec{q} = \frac{1}{|q|} \sum_{i \in q} \vec{i}$$

A batch update algorithm, due to Lloyd (1957), Forgy (1965), and others, makes assignments to a set of initially randomly-chosen vectors,  $\vec{q}$ , as step 1. Step 2 updates the cluster vectors,  $\vec{q}$ . This is iterated. The distortion error, equation 1, is non-increasing, and a local minimum is achieved in a finite number of iterations.

An online update algorithm is due to MacQueen (1967). After each presentation of an observation vector,  $\vec{i}$ , the closest cluster vector,  $\vec{q}$ , is updated to take account of it. Such an approach is well-suited for a continuous input data stream (implying “online” learning of cluster vectors).

Both algorithms are gradient descent ones. In the online case, much attention has been devoted to best learning rate schedules in the neural network (competitive learning) literature: Darken and Moody (1991, 1992), Darken et al. (1992), Fritzke (1997).

A difficulty, less controllable in the case of the batch algorithm, is that clusters may become (and stay) empty. This may be acceptable, but also may be in breach of our original problem formulation. An alternative to the batch update algorithm is Späth’s (1985) exchange algorithm. Each observation is considered for possible assignment into any of the *other* clusters. Updating and “downdating” formulas are given by Späth. This exchange algorithm is stated to be faster to converge and to produce better (smaller) values of the objective function. Over decades of use, we have also verified that it is a superior algorithm to the minimal distance one.

K-means is very closely related to Voronoi (Dirichlet) tessellations, to Kohonen self-organizing feature maps, and various other methods.

The batch learning algorithm above may be viewed as

1. An assignment step which we will term the E (estimation) step: estimate the posteriors,

$$P(\text{observations} \mid \text{cluster centres})$$

2. A cluster update step, the M (maximization) step, which maximizes a cluster center likelihood.

Neal and Hinton (1998) cast the k-means optimization problem in such a way that the both E- and M-steps monotonically increase the maximand’s values. The EM algorithm may, too, be enhanced to allow for online as well as batch learning (Sato and Ishii, 1999).

In Thiesson et al. (1999), k-means is implemented (i) by traversing blocks of data, cyclically, and incrementally updating the sufficient statistics and parameters, and (ii) instead of cyclic traversal, sampling from subsets of the data is used. Such an approach is admirably suited for very large data sets, where in-memory storage is not feasible. Examples used by Thiesson et al. (1999) include the clustering of a half million 300-dimensional records.

## 10 Fast Model-Based Clustering

It is traditional to note that models and (computational) speed don't mix. We review recent progress in this section.

### 10.1 Modeling of Signal and Noise

A simple and applicable model is a distribution mixture, with the signal modeled by Gaussians, in the presence of Poisson background noise.

Consider data which are generated by a mixture of  $(G - 1)$  bivariate Gaussian densities,  $f_k(x; \theta) \sim \mathcal{N}(\mu_k, \Sigma_k)$ , for clusters  $k = 2, \dots, G$ , and with Poisson background noise corresponding to  $k = 1$ . The overall population thus has the mixture density

$$f(x; \theta) = \sum_{k=1}^G \pi_k f_k(x; \theta)$$

where the mixing or prior probabilities,  $\pi_k$ , sum to 1, and  $f_1(x; \theta) = \mathcal{A}^{-1}$ , where  $\mathcal{A}$  is the area of the data region. This is the basis for *model-based clustering* (Banfield and Raftery, 1993, Dasgupta and Raftery, 1998, Murtagh and Raftery, 1984, Banerjee and Rosenfeld, 1993).

The parameters,  $\theta$  and  $\pi$ , can be estimated efficiently by maximizing the mixture likelihood

$$L(\theta, \pi) = \prod_{i=1}^n f(x_i; \theta),$$

with respect to  $\theta$  and  $\pi$ , where  $x_i$  is the  $i$ -th observation.

Now let us assume the presence of two clusters, one of which is Poisson noise, the other Gaussian. This yields the mixture likelihood

$$L(\theta, \pi) = \prod_{i=1}^n \left[ \pi_1 \mathcal{A}^{-1} + \pi_2 \frac{1}{2\pi\sqrt{|\Sigma|}} \exp \left\{ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\} \right],$$

where  $\pi_1 + \pi_2 = 1$ .

An iterative solution is provided by the expectation-maximization (EM) algorithm of Dempster et al. (1977). We have already noted this algorithm in informal terms in the last section, dealing with k-means. Let the "complete" (or "clean" or "output") data be  $y_i = (x_i, z_i)$  with indicator set  $z_i = (z_{i1}, z_{i2})$  given by  $(1, 0)$  or  $(0, 1)$ . Vector  $z_i$  has a multinomial distribution with parameters  $(1; \pi_1, \pi_2)$ . This leads to the *complete data log-likelihood*:

$$l(y, z; \theta, \pi) = \sum_{i=1}^n \sum_{k=1}^2 z_{ik} [\log \pi_k + \log f_k(x_k; \theta)]$$

The E-step then computes  $\hat{z}_{ik} = E(z_{ik} | x_1, \dots, x_n, \theta)$ , i.e. the posterior probability that  $x_i$  is in cluster  $k$ . The M-step involves maximization of the *expected complete data log-likelihood*:

$$l^*(y; \theta, \pi) = \sum_{i=1}^n \sum_{k=1}^2 \hat{z}_{ik} [\log \pi_k + \log f_k(x_i; \theta)].$$

The E- and M-steps are iterated until convergence.

For the 2-class case (Poisson noise and a Gaussian cluster), the complete-data likelihood is

$$L(y, z; \theta, \pi) = \prod_{i=1}^n \left[ \frac{\pi_1}{\mathcal{A}} \right]^{z_{i1}} \left[ \frac{\pi_2}{2\pi\sqrt{|\Sigma|}} \exp \left\{ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\} \right]^{z_{i2}}$$

The corresponding expected log-likelihood is then used in the EM algorithm. This formulation of the problem generalizes to the case of  $G$  clusters, of arbitrary distributions and dimensions.



Fraley (1999) discusses implementation of model-based clustering, including publicly available software.

In order to assess the evidence for the presence of a signal-cluster, we use the *Bayes factor* for the mixture model,  $M_2$ , that includes a Gaussian density as well as background noise, against the “null” model,  $M_1$ , that contains only background noise. The Bayes factor is the posterior odds for the mixture model against the pure noise model, when neither is favored a priori. It is defined as  $B = p(x|M_2)/p(x|M_1)$ , where  $p(x|M_2)$  is the *integrated likelihood* of the mixture model  $M_2$ , obtained by integrating over the parameter space. For a general review of Bayes factors, their use in applied statistics, and how to approximate and compute them, see Kass and Raftery (1995).

We approximate the Bayes factor using the *Bayesian Information Criterion* (BIC) (Schwartz, 1978). For a Gaussian cluster and Poisson noise, this takes the form:

$$2 \log B \approx BIC = 2 \log L(\hat{\theta}, \hat{\pi}) + 2n \log \mathcal{A} - 6 \log n,$$

where  $\hat{\theta}$  and  $\hat{\pi}$  are the maximum likelihood estimators of  $\theta$  and  $\pi$ , and  $L(\hat{\theta}, \hat{\pi})$  is the maximized mixture likelihood.

A review of the use of the BIC criterion for model selection – and more specifically for choosing the number of clusters in a data set – can be found in Fraley and Raftery (1998).

An application of mixture modeling and the BIC criterion to gamma-ray burst data can be found in Mukherjee et al. (1998). So far around 800 observations have been assessed, but as greater numbers become available we will find the inherent number of clusters in a similar way, in order to try to understand more about the complex phenomenon of gamma-ray bursts.

## 10.2 Application to Thresholding

Consider an image or a planar or 3-dimensional set of object positions. For simplicity we consider the case of setting a single threshold in the image intensities, or the point set’s spatial density.

We deal with a combined mixture density of two *univariate* Gaussian distributions  $f_k(x; \theta) \sim \mathcal{N}(\mu_k, \sigma_k)$ . The overall population thus has the mixture density

$$f(x; \theta) = \sum_{k=1}^2 \pi_k f_k(x; \theta)$$

where the mixing or prior probabilities,  $\pi_k$ , sum to 1.

When the mixing proportions are assumed equal, the log-likelihood takes the form

$$l(\theta) = \sum_{i=1}^n \ln \left[ \sum_{k=1}^2 \frac{1}{2\pi\sqrt{|\sigma_k|}} \exp \left\{ -\frac{1}{2\sigma_k} (x_i - \mu_k)^2 \right\} \right]$$

The EM algorithm is then used to iteratively solve this (see Celeux and Govaert, 1995). This method is used for appraisals of textile (jeans and other fabrics) fault detection in Campbell et al. (1999). Industrial vision inspection systems potentially produce large data streams, and fault detection can be a good application for fast clustering methods. We are currently using a mixture model of this sort on SEM (scanning electron microscope) images of cross-sections of concrete to allow for subsequent characterization of physical properties.

Image segmentation, per se, is a relatively straightforward application, but there are novel and interesting aspects to the two studies mentioned. In the textile case, the faults are very often perceptual and relative, rather than “absolute” or capable of being analyzed in isolation. In the SEM imaging case, a first phase of processing is applied to de-speckle the images, using multiple resolution noise filtering.

Turning from concrete to cosmology, the Sloan Digital Sky Survey (SDSS, 1999) is producing a sky map of more than 100 million objects, together with 3-dimensional information (redshifts) for a million galaxies. Pelleg and Moore (1999) describe mixture modeling, using a k-D tree preprocessing to expedite the finding of the class (mixture) parameters, e.g. means, covariances.

## 11 Noise Modeling

In Starck et al. (1998) and in a wide range of papers, we have pursued an approach for the noise modeling of observed data. A multiple resolution scale vision model or data generation process is used, to allow for the phenomenon being observed on different scales. In addition, a wide range of options are permitted for the data generation transfer path, including additive and multiplicative, stationary and non-stationary, Gaussian (“read out” noise), Poisson (random shot noise), and so on.

Given point pattern clustering in two- or three-dimensional spaces, we will limit our overview here to the Poisson noise case.

### 11.1 Poisson noise with few events using the $\grave{a}$ trous transform

If a wavelet coefficient  $w_j(x, y)$  is due to noise, it can be considered as a realization of the sum  $\sum_{k \in K} n_k$  of independent random variables with the same distribution as that of the wavelet function ( $n_k$  being the number of events used for the calculation of  $w_j(x, y)$ ). This allows comparison of the wavelet coefficients of the data with the values which can be taken by the sum of  $n$  independent variables.

The distribution of one event in wavelet space is then directly given by the histogram  $H_1$  of the wavelet  $\psi$ . As we consider independent events, the distribution of a coefficient  $w_n$  (note the changed subscripting for  $w$ , for convenience) related to  $n$  events is given by  $n$  autoconvolutions of  $H_1$ :

$$H_n = H_1 \otimes H_1 \otimes \dots \otimes H_1$$

For a large number of events,  $H_n$  converges to a Gaussian.

Fig. 8 shows an example of where point pattern clusters – density bumps in this case – are sought, with a great amount of background clutter. Murtagh and Starck (1998) refer to the fact that there is no computational dependence on the number of points (signal or noise) in such a problem, when using a wavelet transform with noise modeling.

Some other alternative approaches will be briefly noted. The Haar transform presents the advantage of its simplicity for modeling Poisson noise. Analytic formulas for wavelet coefficient distributions have been derived by Kolaczyk (1997), and Jammal and Bijaoui (1999). Using a new wavelet transform, the Haar  $\grave{a}$  trous transform, Zheng et al. (1999) appraise a denoising approach for financial data streams, – an important preliminary step for subsequent clustering, forecasting, or other processing.

### 11.2 Poisson noise with nearest neighbor clutter removal

The wavelet approach is certainly appropriate when the wavelet function reflects the type of object sought (e.g. isotropic), and when superimposed point patterns are to be analyzed. However, non-superimposed point patterns of complex shape are very well treated by the approach described in Byers and Raftery (1998). Using a homogeneous Poisson noise model, they derive the distribution of the distance of a point to its  $k$ th nearest neighbor.

Next, Byers and Raftery (1998) consider the case of a Poisson process which is signal, superimposed on a Poisson process which is clutter. The  $k$ th nearest neighbor distances are modeled as a mixture distribution: a histogram of these, for given  $k$ , will yield a bimodal distribution if our assumption is correct. This mixture distribution problem is solved using the EM algorithm. Generalization to higher dimensions, e.g. 10, is also discussed.

Similar data was analyzed by noise modeling and a Voronoi tessellation preprocessing of the data in Allard and Fraley (1997). It is pointed out there how this can be a very useful approach with the Voronoi tiles have meaning in relation the morphology of the point patterns. However, it does not scale well to higher dimensions, and the statistical noise modeling is approximate.

Ebeling and Wiedenmann (1993), reproduced in Dobrzycki et al. (1999), propose the use of a Voronoi tessellation for astronomical X-ray object detection and characterization.

## 12 Cluster-Based User Interfaces

Information retrieval by means of “semantic road maps” was first detailed by Doyle (1961). The spatial metaphor is a powerful one in human information processing. The spatial metaphor also lends itself well to modern distributed computing environments such as the web. The Kohonen self-organizing feature map (SOM) method is an effective means towards this end of a visual information retrieval user interface. We will also provide an illustration of web-based semantic maps based on hyperlink clustering.

The Kohonen map is, at heart, k-means clustering with the additional constraint that cluster centers be located on a regular grid (or some other topographic structure) and furthermore their location on the grid be monotonically related to pairwise proximity (Murtagh and Hernández-Pajares, 1995). The nice thing about a regular grid output representation space is that it lends itself well as a visual user interface.

Fig. 9 shows a visual and interactive user interface map, using a Kohonen self-organizing feature map (SOM). Color is related to density of document clusters located at regularly-spaced nodes of the map, and some of these nodes/clusters are annotated. The map is installed as a clickable imagemap, with CGI programs accessing lists of documents and – through further links – in many cases, the full documents. In the example shown, the user has queried a node and results are seen in the right-hand panel. Such maps are maintained for (currently) 12000 articles from the *Astrophysical Journal*, 7000 from *Astronomy and Astrophysics*, over 2000 astronomical catalogs, and other data holdings. More information on the design of this visual interface and user assessment can be found in Poinçot et al. (1998, 1999).

Guillaume (Guillaume and Murtagh, 1999) developed a Java-based visualization tool for hyperlink-based data, consisting of astronomers, astronomical object names, article titles, and with the possibility of other objects (images, tables, etc.). Through weighting, the various types of links could be prioritized. An iterative refinement algorithm was developed to map the nodes (objects) to a regular grid of cells, which as for the Kohonen SOM map, are clickable and provide access to the data represented by the cluster. Fig. 10 shows an example for an astronomer (Prof. Jean Heyvaerts, Strasbourg Astronomical Observatory).

These new cluster-based visual user interfaces are not computationally demanding. They are not however scalable in their current implementation. Document management (see e.g. Cartia, 1999) is less the motivation as is instead the interactive user interface.

## 13 Images from Data

It is quite impressive how 2D (or 3D) image signals can handle with ease the scalability limitations of clustering and many other data processing operations. The contiguity imposed on adjacent pixels bypasses the need for nearest neighbor finding. It is very interesting therefore to consider the feasibility of taking problems of clustering massive data sets into the 2D image domain. We will look at a few recent examples of work in this direction.

Church and Helfman (1993) address the problem of visualizing possibly millions of lines of computer program code, or text. They consider an approach borrowed from DNA sequence analysis. The data sequence is tokenized by splitting it into its atoms (line, word, character, etc.) and then placing a dot at position  $i, j$  if the  $i$ th input token is the same as the  $j$ th. The resulting dotplot, it is argued, is not limited by the available display screen space, and can lead to discovery of large-scale structure in the data.

When data do not have a sequence we have an invariance problem which can be resolved by finding some row and column permutation which pulls large array values together, and perhaps

furthermore into proximity to an array diagonal. Berry et al. (1996) have studied the case of large sparse arrays. Gathering larger (or nonzero) array elements to the diagonal can be viewed in terms of minimizing the envelope of nonzero values relative to the diagonal. This can be formulated and solved in purely symbolic terms by reordering vertices in a suitable graph representation of the matrix. A widely-used method for symmetric sparse matrices is the Reverse Cuthill-McKee (RCM) method.

The complexity of the RCM method for ordering rows or columns is proportional to the product of the maximum degree of any vertex in the graph representing the array values and the total number of edges (nonzeroes in the matrix). For hypertext matrices with small maximum degree, the method would be extremely fast. The strength of the method is its low time complexity but it does suffer from certain drawbacks. The heuristic for finding the starting vertex is influenced by the initial numbering of vertices and so the quality of the reordering can vary slightly for the same problem for different initial numberings. Next, the overall method does not accommodate dense rows (e.g., a common link used in every document), and if a row has a significantly large number of nonzeroes it might be best to process it separately; i.e., extract the dense rows, reorder the remaining matrix and augment it by the dense rows (or common links) numbered last. Elapsed CPU times for a range of arrays and permuting methods are given in Berry et al. (1996), and as an indication show performances between 0.025 to 3.18 seconds for permuting a  $4000 \times 400$  array. A review of public domain software for carrying out SVD and other linear algebra operations on large sparse data sets can be found in Berry et al. (1999, section 8.3).

Once we have a sequence-respecting array, we can immediately apply efficient visualization techniques from image analysis. Murtagh et al. (1999) investigate the use of noise filtering (i.e. to remove less useful array entries) using a multiscale wavelet transform approach.

An example follows. From the Concise Columbia Encyclopedia (1989 2nd ed., online version) a set of data relating to 12025 encyclopedia entries and to 9778 cross-references or links was used. Fig. 11 shows a  $500 \times 450$  subarray, based on a correspondence analysis (i.e. ordering of projections on the first factor).

This part of the encyclopedia data was filtered using the wavelet and noise-modeling methodology described in Murtagh et al. (1999) and the outcome is shown in Fig. 12. Overall the recovery of the more apparent alignments, and hence visually stronger clusters, is excellent. The first relatively long “horizontal bar” was selected – it corresponds to column index (link) 1733 = `geological era`. The corresponding row indices (articles) are, in sequence:

SILURIAN PERIOD  
 PLEISTOCENE EPOCH  
 HOLOCENE EPOCH  
 PRECAMBRIAN TIME  
 CARBONIFEROUS PERIOD  
 OLIGOCENE EPOCH  
 ORDOVICIAN PERIOD  
 TRIASSIC PERIOD  
 CENOZOIC ERA  
 PALEOCENE EPOCH  
 MIOCENE EPOCH  
 DEVONIAN PERIOD  
 PALEOZOIC ERA  
 JURASSIC PERIOD  
 MESOZOIC ERA  
 CAMBRIAN PERIOD  
 PLIOCENE EPOCH  
 CRETACEOUS PERIOD

The work described here is based on a number of technologies: (i) data visualization techniques;

(ii) the wavelet transform for data analysis; and (iii) data matrix permuting techniques. The wavelet transform has linear computational cost in terms of image row and column dimensions, and is independent of the pixel values.

## 14 Conclusion

Viewed from a commercial or managerial perspective, one could justifiably ask where we are now in our understanding of problems in this area relative to where we were back in the 1960s? Depending on our answer to this, we may well proceed to a second question: Why have all important problems not been solved by now in this area – are there major outstanding problems to be solved?

As described in this chapter, a solid body of experimental and theoretical results have been built up over the last few decades. Clustering remains a requirement which is a central infrastructural element of very many application fields.

There is continual renewal of the essential questions and problems of clustering, relating to new data, new information, and new environments. There is no logjam in clustering research and development simply because the rivers of problems continue to broaden and deepen. Clustering and classification remain quintessential issues in our computing and information technology environments (Murtagh, 1998).

## Acknowledgements

Some of this work, in particular in sections 9 to 12, represents various collaborations with the following: J.L. Starck, CEA; A. Raftery, University of Washington; C. Fraley, University of Washington and MathSoft Inc.; D. Washington, MathSoft, Inc.; Ph. Poinçot and S. Lesteven, Strasbourg Observatory; D. Guillaume, Strasbourg Observatory, University of Illinois and NCSA.

## References

1. Allard, D. and Fraley, C., “Non-parametric maximum likelihood estimation of features in spatial point processes using Voronoi tessellation”, *Journal of the American Statistical Association*, 92, 1485–1493, 1997.
2. Arabie, P. and Hubert, L.J., “An overview of combinatorial data analysis”, in Arabie, P., Hubert, L.J. and De Soete, G., Eds., *Clustering and Classification*, World Scientific, Singapore, 5–63, 1996.
3. Arabie, P., Hubert, L.J. and De Soete, G., Eds., *Clustering and Classification*, World Scientific, Singapore, 1996.
4. Banerjee, S. and Rosenfeld, A., “Model-based cluster analysis”, *Pattern Recognition*, 26, 963–974, 1993.
5. Banfield, J.D. and Raftery, A.E., “Model-based Gaussian and non-Gaussian clustering”, *Biometrics*, 49, 803–821, 1993.
6. Bennett, K.P., Fayyad, U. and Geiger, D., “Density-based indexing for approximate nearest neighbor queries”, Microsoft Research Technical Report MSR-TR-98-58, 1999.
7. Bentley, J.L. and Friedman, J.H., “Fast algorithms for constructing minimal spanning trees in coordinate spaces”, *IEEE Transactions on Computers*, C-27, 97–105, 1978.
8. Bentley, J.L., Weide, B.W. and Yao, A.C., “Optimal expected time algorithms for closest point problems”, *ACM Transactions on Mathematical Software*, 6, 563–580, 1980.

9. Berry, M.W., Hendrickson, B. and Raghavan, P., "Sparse matrix reordering schemes for browsing hypertext", in Renegar, J., Shub, M. and Smale, S., Eds., *Lectures in Applied Mathematics (LAM) Vol. 32: The Mathematics of Numerical Analysis*, American Mathematical Society, 99–123, 1996.
10. Berry, M.W., Drmač, Z. and Jessup, E.R., "Matrices, vector spaces, and information retrieval", *SIAM Review*, 41, 335–362, 1999.
11. Beyer, K., Goldstein, J., Ramakrishnan, R. and Shaft, U., "When is nearest neighbor meaningful?", in *Proceedings of the 7th International Conference on Database Theory (ICDT)*, Jerusalem, Israel, 1999, in press.
12. Borodin, A., Ostrovsky, R. and Rabani, Y., "Subquadratic approximation algorithms for clustering problems in high dimensional spaces", *Proc. 31st ACM Symposium on Theory of Computing (STOC-99)*, 1999.
13. Broder, A.J., "Strategies for efficient incremental nearest neighbor search", *Pattern Recognition*, 23, 171–178, 1990.
14. Broder, A.Z., Glassman, S.C., Manasse, M.S. and Zweig, G., "Syntactic clustering of the web", *Proc. Sixth International World Wide Web Conference*, 391–404, 1997.
15. Broder, A.Z., "On the resemblance and containment of documents", In *Compression and Complexity of Sequences (SEQUENCES'97)*, 21–29, IEEE Computer Society, 1998.
16. Bruynooghe, M., "Méthodes nouvelles en classification automatique des données taxinomiques nombreuses", *Statistique et Analyse des Données*, no. 3, 24–42, 1977.
17. Burkhard, W.A. and Keller, R.M., "Some approaches to best-match file searching", *Communications of the ACM*, 16, 230–236, 1973.
18. Byers, S.D. and Raftery, A.E., "Nearest neighbor clutter removal for estimating features in spatial point processes", *Journal of the American Statistical Association*, 93, 577–584, 1998.
19. Campbell, J.G., Fraley, C., Stanford, D., Murtagh, F. and Raftery, A.E., "Model-based methods for textile fault detection", *International Journal of Imaging Science and Technology*, 10, 339–346, 1999.
20. Cartia, Inc., *Mapping the Information Landscape*, client-server software system, <http://www.cartia.com>, 1999.
21. Celeux, G. and Govaert, G., "Gaussian parsimonious clustering models", *Pattern Recognition*, 28, 781–793, 1995.
22. Cheriton, D. and Tarjan, D.E., "Finding minimum spanning trees", *SIAM Journal on Computing*, 5, 724–742, 1976.
23. Church, K.W. and Helfman, J.I., "Dotplot: a program for exploring self-similarity in millions of lines of text and code", *Journal of Computational and Graphical Statistics*, 2, 153–174, 1993.
24. Croft, W.B., "Clustering large files of documents using the single-link method", *Journal of the American Society for Information Science*, 28, 341–344, 1977.
25. Darken, C. and Moody, J., "Note on learning rate schedules for stochastic optimization", *Advances in Neural Information Processing Systems 3*, Lippmann, Moody and Touretzky, Eds., Morgan Kaufmann, Palo Alto, 1991.

26. Darken, C., Chang, J. and Moody, J., "Learning rate schedules for faster stochastic gradient search", *Neural Networks for Signal Processing 2*, Proceedings of the 1992 IEEE Workshop, IEEE Press, Piscataway, 1992.
27. Darken, C. and Moody, J., "Towards faster stochastic gradient search", *Advances in Neural Information Processing Systems 4*, Moody, Hanson and Lippmann, Eds., Morgan Kaufman, San Mateo, 1992.
28. Dasarathy, B.V., *Nearest Neighbor (NN) Norms: NN Pattern Classification Techniques*, IEEE Computer Society Press, New York, 1991.
29. Dasgupta, A. and Raftery, A.E., "Detecting features in spatial point processes with clutter via model-based clustering", *Journal of the American Statistical Association*, 93, 294–302, 1998.
30. Day, W.H.E. and Edelsbrunner, H., "Efficient algorithms for agglomerative hierarchical clustering methods", *Journal of Classification*, 1, 7–24, 1984.
31. Defays, D., "An efficient algorithm for a complete link method", *Computer Journal*, 20, 364–366, 1977.
32. Delannoy, C., "Un algorithme rapide de recherche de plus proches voisins", *RAIRO Informatique/Computer Science*, 14, 275–286, 1980.
33. Dempster, A.P., Laird, N.M. and Rubin, D.B., "Maximum likelihood from incomplete data via the EM algorithm", *Journal of the Royal Statistical Society, Series, B* 39, 1–22, 1977).
34. Dobrzycki, A., Ebeling, H., Glotfelty, K., Freeman, P., Damiani, F, Elvis, M. and Calderwood, T., *Chandra Detect 1.0 User Guide*, Chandra X-Ray Center, Smithsonian Astrophysical Observatory, Version 0.9, 1999.
35. Doyle, L.B., "Semantic road maps for literature searchers", *Journal of the ACM*, 8, 553–578, 1961.
36. Eastman, C.M. and Weiss, S.F., "Tree structures for high dimensionality nearest neighbor searching", *Information Systems*, 7, 115–122, 1982.
37. Ebeling, H. and Wiedenmann, G., "Detecting structure in two dimensions combining Voronoi tessellation and percolation", *Physical Review E*, 47, 704–714, 1993.
38. Forgy, E., "Cluster analysis of multivariate data: efficiency vs. interpretability of classifications", *Biometrics*, 21, 768, 1965.
39. Fraley, C., "Algorithms for model-based Gaussian hierarchical clustering", *SIAM Journal of Scientific Computing*, 20, 270–281, 1999.
40. Fraley, C. and Raftery, A.E., "How many clusters? Which clustering method? Answers via model-based cluster analysis", *The Computer Journal*, 41, 578–588, 1999.
41. Friedman, J.H., Baskett, F. and Shustek, L.J., "An algorithm for finding nearest neighbors", *IEEE Transactions on Computers*, C-24, 1000–1006, 1975.
42. Friedman, J.H., Bentley, J.L. and Finkel, R.A., "An algorithm for finding best matches in logarithmic expected time", *ACM Transactions on Mathematical Software*, 3, 209–226, 1977.
43. Fritzke, B., "Some competitive learning methods", <http://www.neuroinformatik.ruhr-uni-bochum.de/ini/VDM/research/gsn/JavaPaper>

44. Fukunaga, K. and Narendra, P.M., "A branch and bound algorithm for computing k-nearest neighbors", *IEEE Transactions on Computers*, C-24, 750–753, 1975.
45. Gabow, H.N., Galil, Z., Spencer, T., and Tarjan, R.E., "Efficient algorithms for finding minimum spanning trees in undirected and directed graphs", *Combinatorica*, 6, 109–122, 1986.
46. Gillet, V.J., Wild, D.J., Willett, P. and Bradshaw, J., "Similarity and dissimilarity methods for processing chemical structure databases", *The Computer Journal*, 41, 547–558, 1998.
47. Gordon, A.D., *Classification*, 2nd ed., Chapman and Hall, 1999.
48. Griffiths, A., Robinson, L.A. and Willett, P., "Hierarchic agglomerative clustering methods for automatic document classification", *Journal of Documentation*, 40, 175–205, 1984.
49. Guillaume, D. and Murtagh, F., "Clustering of XML documents", *Computer Physics Communications*, submitted, 1999.
50. Hodgson, M.E., "Reducing the computational requirements of the minimum-distance classifier", *Remote Sensing of Environment*, 25, 117–128, 1988.
51. Horowitz, E. and Sahni, S., *Fundamentals of Computer Algorithms*, Chapter 4 The Greedy Method, Pitman, London, 1979.
52. Jain, A.K. and Dubes, R.C, *Algorithms for Clustering Data*, Prentice-Hall, Englewood Cliffs, 1988.
53. Jammal, G. and Bijaoui, A., "Multiscale image restoration for photon imaging systems", *SPIE Conference on Signal and Image Processing: Wavelet Applications in Signal and Image Processing VII*, July 1999.
54. Juan, J., "Programme de classification hiérarchique par l'algorithme de la recherche en chaîne des voisins réciproques", *Les Cahiers de l'Analyse des Données*, VII, 219–225, 1982.
55. Kamgar-Parsi, B. and Kanal, L.N., "An improved branch and bound algorithm for computing k-nearest neighbors", *Pattern Recognition Letters*, 3, 7–12, 1985.
56. Kass, R.E. and Raftery, A.E., "Bayes factors", *Journal of the American Statistical Association*, 90, 773–795, 1995.
57. Kittler, J., "A method for determining k-nearest neighbors", *Kybernetes*, 7, 313–315, 1978.
58. Kolaczyk, E.D., "Nonparametric estimation of gamma-ray burst intensities using Haar wavelets", *Astrophysical Journal*, 483, 340–349, 1997.
59. Kushilevitz, E., Ostrovsky, R. and Rabani, Y., "Efficient search for approximate nearest neighbors in high-dimensional spaces", *Proc. of 30th ACM Symposium on Theory of Computing (STOC-30)*, 1998.
60. Lloyd, P., "Least squares quantization in PCM." Technical note, Bell Laboratories, 1957. Published in *IEEE Transactions on Information Theory*, 1982.
61. MacQueen, J., "Some methods for classification and analysis of multivariate observations". *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1, pp. 281–297, Berkeley, University of California Press, 1976.
62. Marimont, R.B. and Shapiro, M.B., "Nearest neighbor searches and the curse of dimensionality", *Journal of the Institute of Mathematics and Its Applications*, 24, 59–70, 1979.



63. Micó, L., Oncina, J. and Vidal, E., “An algorithm for finding nearest neighbors in constant average time with a linear space complexity”, in 11th International Conference on Pattern Recognition, Volume II, IEEE Computer Science Press, New York, 557–560, 1992.
64. Moore, A., “Very fast EM-based mixture model clustering using multiresolution kd-trees”, Neural Information Processing Systems, December 1998, forthcoming.
65. Motwani, R. and Raghavan, P., Randomized Algorithms, Cambridge University Press, 1995.
66. Mukherjee, S., Feigelson, E.D., Babu, G.J., Murtagh, F., Fraley, C. and Raftery, A., “Three types of gamma-ray bursts”, The Astrophysical Journal, 508, 314–327, 1998.
67. Murtagh, F., “A very fast, exact nearest neighbor algorithm for use in information retrieval”, Information Technology, 1, 275–283, 1982.
68. Murtagh, F., “Expected time complexity results for hierarchic clustering algorithms which use cluster centers”, Information Processing Letters, 16, 237–241, 1983.
69. Murtagh, F., “Complexities of hierarchic clustering algorithms: state of the art”, Computational Statistics Quarterly, 1, 101–113, 1984.
70. Murtagh F. and A.E. Raftery, A.E., “Fitting straight lines to point patterns”, Pattern Recognition, 17, 479–483, 1984.
71. Murtagh, F., Multidimensional Clustering Algorithms, Physica-Verlag, Würzburg, 1985.
72. Murtagh, F. and Heck, A., Multivariate Data Analysis, Kluwer Academic, Dordrecht, 1987.
73. Murtagh, F., “Comments on ‘Parallel algorithms for hierarchical clustering and cluster validity’”, IEEE Transactions on Pattern Analysis and Machine Intelligence, 14, 1056–1057, 1992.
74. Murtagh, F., “Search algorithms for numeric and quantitative data”, in A. Heck and F. Murtagh, Eds., Intelligent Information Retrieval: The Case of Astronomy and Related Space Sciences, Kluwer Academic, Dordrecht, 49–80, 1993a.
75. Murtagh, F., “Multivariate methods for data analysis”, in Aa. Sandqvist and T.P. Ray, Eds., Central Activity in Galaxies: From Observational Data to Astrophysical Diagnostics, Springer-Verlag, Berlin, 1993b, pp. 209–235.
76. Murtagh, F. and Hernández-Pajares, M., “The Kohonen self-organizing map method: an assessment”, Journal of Classification, 12, 165–190, 1995.
77. Murtagh, F., “Foreword to the Special Issue on Clustering and Classification”, The Computer Journal, 41, 517, 1998.
78. Murtagh, F. and Starck, J.L., “Pattern clustering based on noise modeling in wavelet space”, Pattern Recognition, 31, 847–855, 1998.
79. Murtagh, F., Starck, J.L. and Berry M., “Overcoming the curse of dimensionality in clustering by means of the wavelet transform”, The Computer Journal, 1999, submitted.
80. Neal, R. and Hinton, G., “A view of the EM algorithm that justifies incremental, sparse, and other variants”, in M. Jordan, Ed., Learning in Graphical Models, Kluwer, Dordrecht, 1998, pp. 355–371.
81. Niemann, H. and Goppert, R., “An efficient branch-and-bound nearest neighbor classifier”, Pattern Recognition Letters, 7, 67–72, 1988.

82. Pelleg, D. and Moore, A., "Accelerating exact k-means algorithms with geometric reasoning", Proceedings KDD-99, Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, August, San Diego, 1999.
83. Perry, S.A. and Willett, P., "A review of the use of inverted files for best match searching in information retrieval systems", *Journal of Information Science*, 6, 59–66, 1983.
84. Poinçot, Ph., Lesteven, S. and Murtagh, F., "A spatial user interface to the astronomical literature", *Astronomy and Astrophysics Supplement*, 130, 183–191, 1998.
85. Poinçot, Ph., Lesteven, S. and Murtagh, F. "Maps of information spaces: assessments from astronomy", *Journal of the American Society for Information Science*, submitted, 1999.
86. Ramasubramanian, V. and Paliwal, K.K., "An efficient approximation-algorithm for fast nearest-neighbor search based on a spherical distance coordinate formulation", *Pattern Recognition Letters*, 13, 471–480, 1992.
87. de Rham, C., "La classification hiérarchique ascendante selon la méthode des voisins réciproques", *Les Cahiers de l'Analyse des Données*, V, 135–144, 1980.
88. Richetin, M., Rives, G. and Naranjo, M., "Algorithme rapide pour la détermination des k plus proches voisins", *RAIRO Informatique/Computer Science*, 14, 369–378, 1980.
89. Rohlf, F.J., "Algorithm 76: Hierarchical clustering using the minimum spanning tree", *The Computer Journal*, 16, 93–95, 1973.
90. Rohlf, F.J., "A probabilistic minimum spanning tree algorithm", *Information Processing Letters*, 7, 44–48, 1978.
91. Rohlf, F.J., "Single link clustering algorithms", in P.R. Krishnaiah and L.N. Kanal, Eds., *Handbook of Statistics*, Vol. 2, North-Holland, Amsterdam, 267–284, 1982.
92. Salton, G. and McGill, M.J., *Introduction to Modern Information Retrieval*, McGraw-Hill, New York, 1983.
93. M. Sato and S. Ishii, "Reinforcement learning based on on-line EM algorithm", *Advances in Neural Information Processing Systems 11*, M.S. Kearns, S.A. Solla and D.A. Cohn (Eds.), 1052-1058, MIT Press, Cambridge, 1999.
94. Schreiber, T., "Efficient search for nearest neighbors", in A.S. Weigend and N.A. Gershenfeld, eds., *Predicting the Future and Understanding the Past: A Comparison of Approaches*, Addison-Wesley, New York, 1993.
95. G. Schwarz, "Estimating the dimension of a model", *The Annals of Statistics*, 6, 461–464, 1978.
96. SDSS, Sloan Digital Sky Survey, <http://www.sdss.org>
97. Shapiro, M., "The choice of reference points in best-match file searching", *Communications of the ACM*, 20, 339–343, 1977.
98. Sibson, R., "SLINK: an optimally efficient algorithm for the single link cluster method", *The Computer Journal*, 16, 30–34, 1973.
99. Smeaton, A.F. and van Rijsbergen, C.J., "The nearest neighbor problem in information retrieval: an algorithm using upperbounds", *ACM SIGIR Forum*, 16, 83–87, 1981.
100. Sneath, P.H.A. and Sokal, R.R., *Numerical Taxonomy*, W.H. Freeman, San Francisco, 1973.

101. Späth, H., Cluster Dissection and Analysis: Theory, Fortran Programs, Examples, Ellis Horwood, Chichester, 1985.
102. Starck, J.L., Murtagh, F. and Bijaoui, A., Image and Data Analysis: The Multiscale Approach, Cambridge University Press, New York, 1998.
103. Tarjan, R.E., “An improved algorithm for hierarchical clustering using strong components”, Information Processing Letters, 17, 37–41, 1983.
104. Thiesson, B., Meek, C. and Heckerman, D., “Accelerating EM for large databases”, Microsoft Research Technical Report MST-TR-99-31, 1999.
105. Tucker, A., Applied Combinatorics, Wiley, New York, 1980.
106. Vidal Ruiz, E., “An algorithm for finding nearest neighbors in (approximately) constant average time”, Pattern Recognition Letters, 4, 145–157, 1986.
107. Weiss, S.F., “A probabilistic algorithm for nearest neighbor searching”, in R.N. Oddy et al., Eds., Information Retrieval Research, Butterworths, London, 325–333, 1981.
108. White, H.D. and McCain, K.W., “Visualization of literatures”, in M.E. Williams, Ed., Annual Review of Information Science and Technology (ARIST), Vol. 32, 99–168, 1997.
109. Willett, P., “Efficiency of hierarchic agglomerative clustering using the ICL distributed array processor”, Journal of Documentation, 45, 1–45, 1989.
110. Yao, A.C., “An  $O(|E| \log \log |V|)$  algorithm for finding minimum spanning trees”, Information Processing Letters, 4, 21–23, 1975.
111. Yunck, T.P., “A technique to identify nearest neighbors”, IEEE Transactions on Systems, Man, and Cybernetics, SMC-6, 678–683, 1976.
112. Zahn, C.T., “Graph-theoretical methods for detecting and describing Gestalt clusters”, IEEE Transactions on Computers, C-20, 68–86, 1971.
113. Zheng, G., Starck, J.L., J.G. Campbell and Murtagh, F., “Multiscale transforms for filtering financial data streams”, Journal of Computational Intelligence in Finance, 7, 1999.

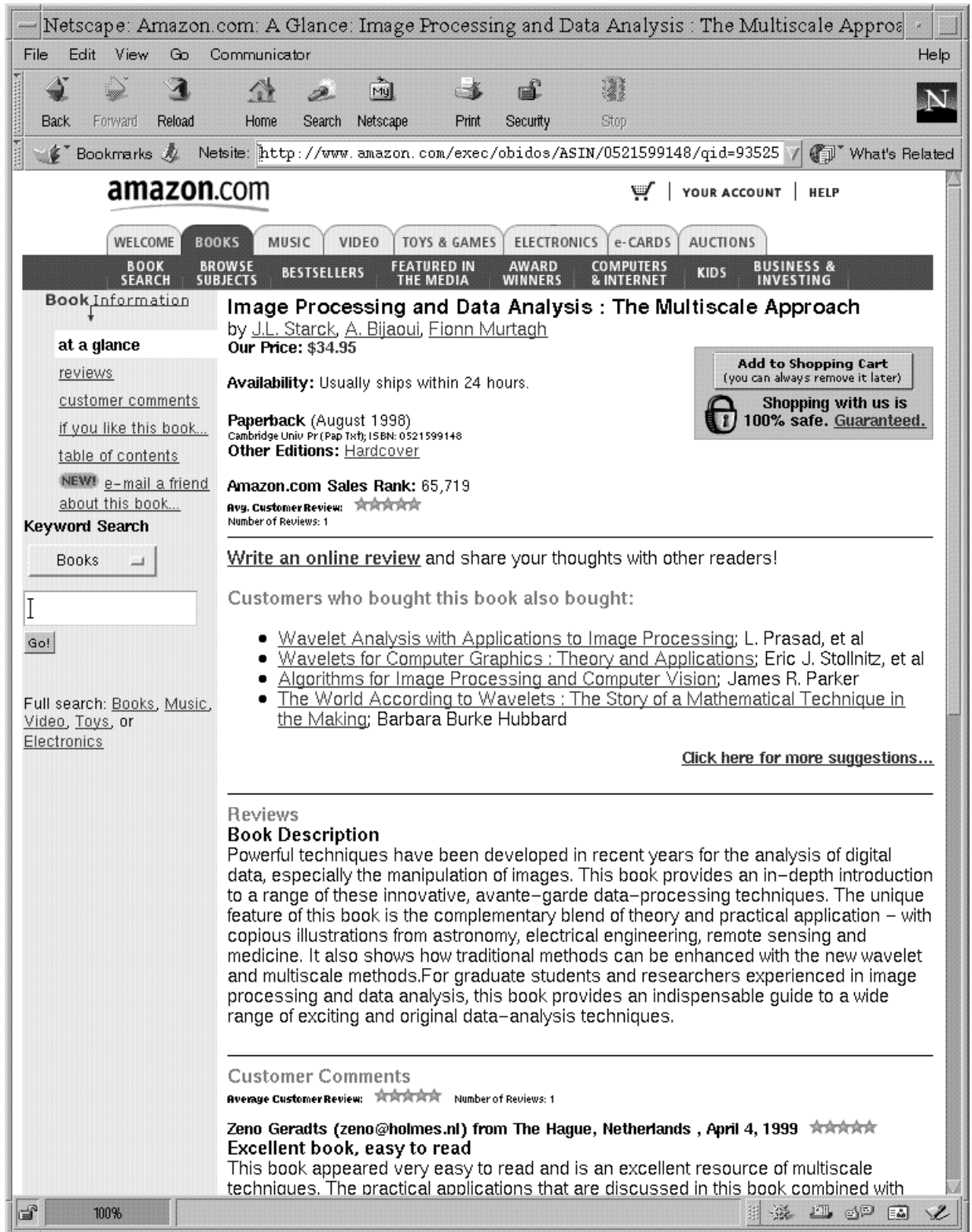


Figure 7: Example of graph clustering in a data mining perspective at Amazon.com: "Customers who bought this book also bought..."

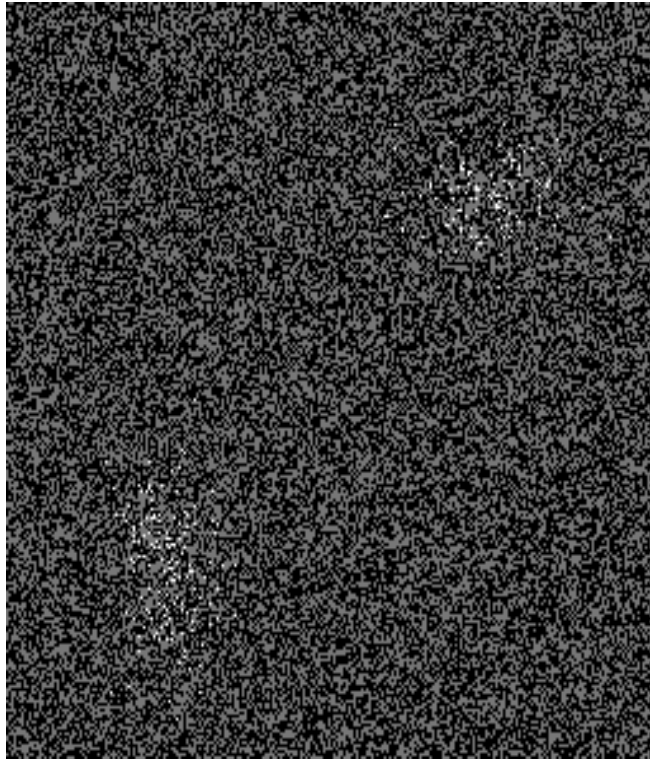


Figure 8: Data in the plane. The  $256 \times 256$  image shows 550 “signal” points – two Gaussian-shaped clusters in the lower left and in the upper right – with in addition 40,000 Poisson noise points added. Details of recovery of the clusters is discussed in Murtagh and Starck (1998).

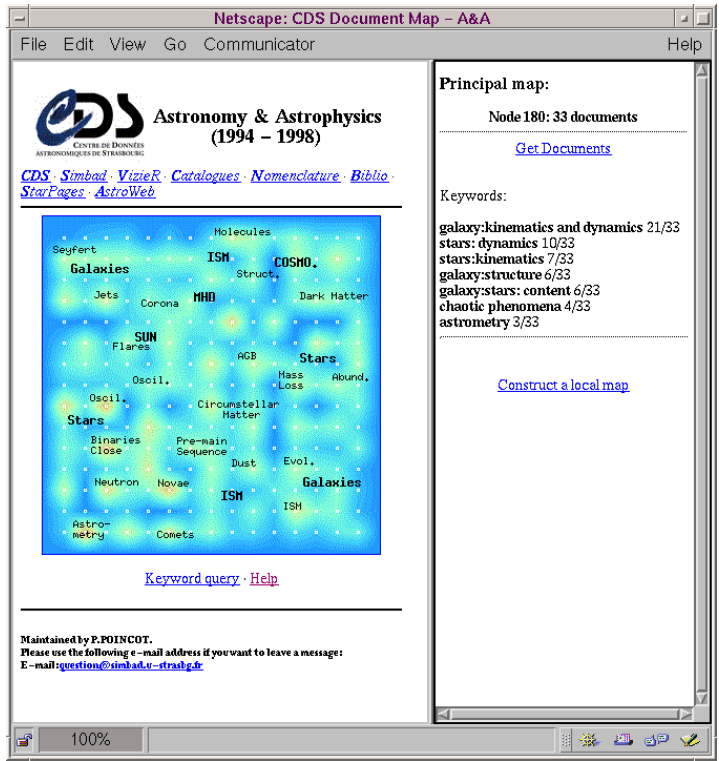


Figure 9: Visual interactive user interface to the journal *Astronomy and Astrophysics*. Original in color.

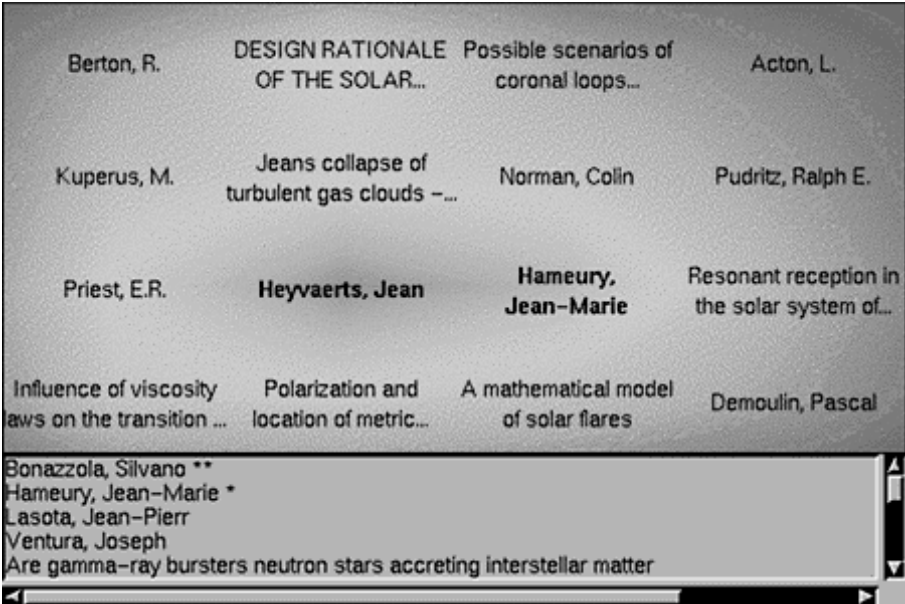


Figure 10: Visual interactive user interfaces, based on graph edges. Map for astronomer Jean Heyvaerts. Original in color.

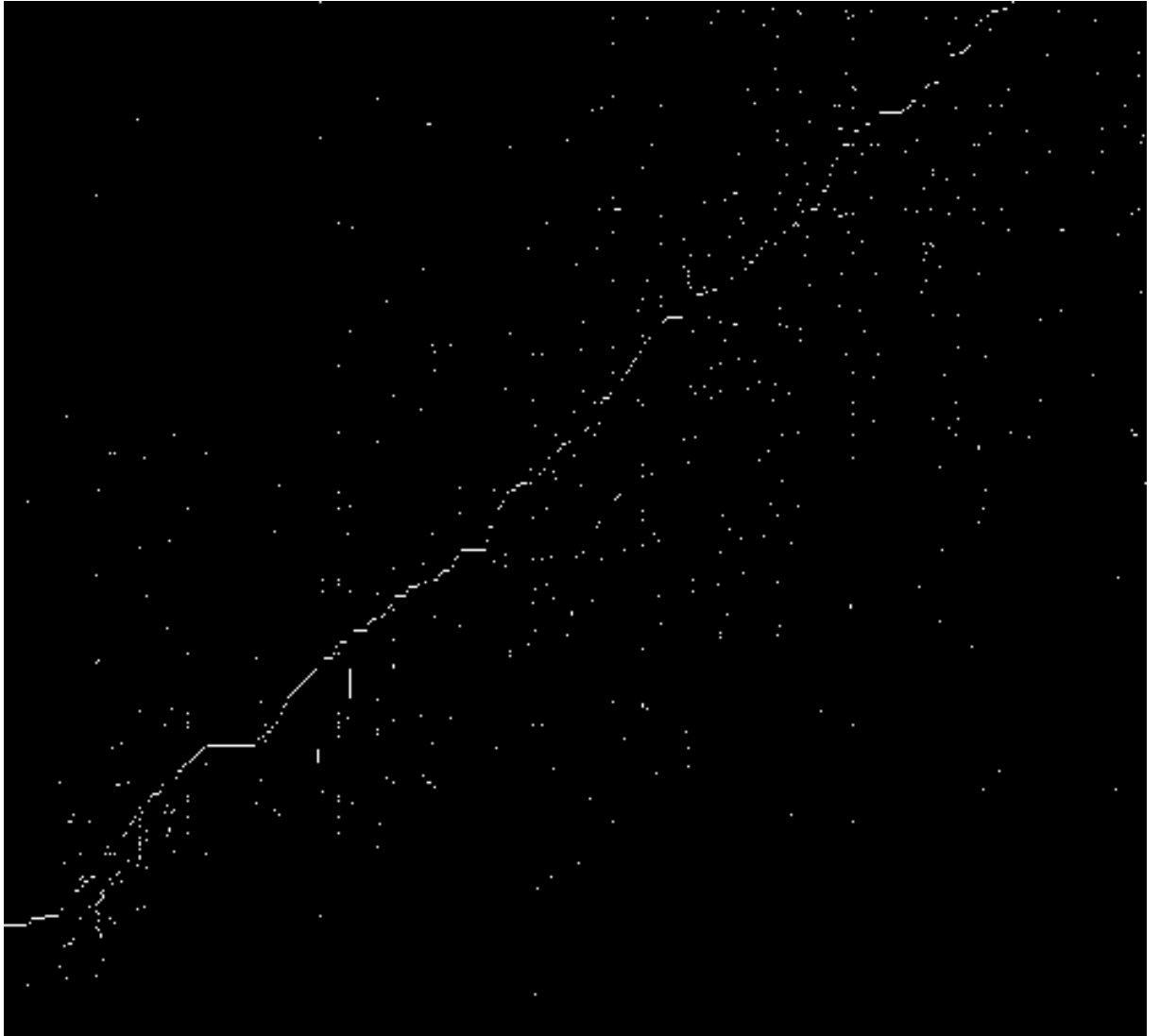


Figure 11: Part ( $500 \times 450$ ) of original encyclopedia incidence data array.

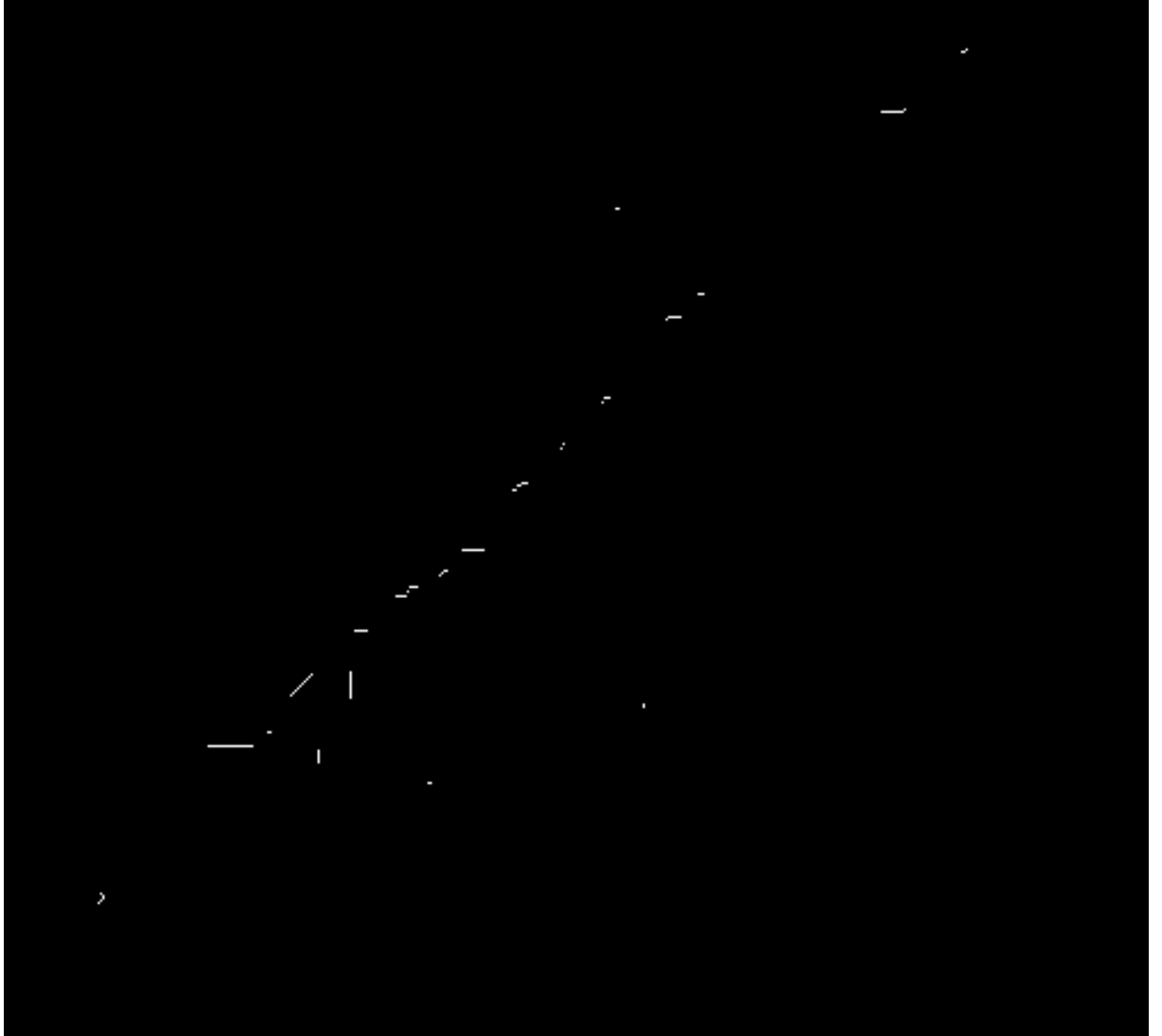


Figure 12: End-product of the filtering of the array shown in the previous Figure.