

# A Fast Algorithm for the Minimum Covariance Determinant Estimator

Peter J. Rousseeuw and Katrien Van Driessen

Revised Version, 15 December 1998

## Abstract

The minimum covariance determinant (MCD) method of Rousseeuw (1984) is a highly robust estimator of multivariate location and scatter. Its objective is to find  $h$  observations (out of  $n$ ) whose covariance matrix has the lowest determinant. Until now applications of the MCD were hampered by the computation time of existing algorithms, which were limited to a few hundred objects in a few dimensions. We discuss two important applications of larger size: one about a production process at Philips with  $n = 677$  objects and  $p = 9$  variables, and a data set from astronomy with  $n = 137,256$  objects and  $p = 27$  variables. To deal with such problems we have developed a new algorithm for the MCD, called FAST-MCD. The basic ideas are an inequality involving order statistics and determinants, and techniques which we call ‘selective iteration’ and ‘nested extensions’. For small data sets FAST-MCD typically finds the exact MCD, whereas for larger data sets it gives more accurate results than existing algorithms and is faster by orders of magnitude. Moreover, FAST-MCD is able to detect an exact fit, i.e. a hyperplane containing  $h$  or more observations. The new algorithm makes the MCD method available as a routine tool for analyzing multivariate data. We also propose the distance-distance plot (or ‘D-D plot’) which displays MCD-based robust distances versus Mahalanobis distances, and illustrate it with some examples.

KEY WORDS: Breakdown value; Multivariate location and scatter; Outlier detection; Regression; Robust estimation.

---

Peter J. Rousseeuw is Professor, Department of Mathematics and Computer Science, Universitaire Instelling Antwerpen (UIA), Universiteitsplein 1, B-2610 Wilrijk, Belgium, and Katrien Van Driessen is Assistant, Faculty of Applied Economics, Universitaire Faculteiten Sint Ignatius Antwerpen (UFSIA), Prinsstraat 13, B-2000 Antwerp, Belgium. We wish to thank Doug Hawkins and José Agulló for making their programs available to us. We also want to dedicate special thanks to Gertjan Otten, Frans Van Dommelen en Herman Veraa for giving us access to the Philips data, and to S.C. Odewahn and his research group at the California Institute of Technology for allowing us to analyze their Digitized Palomar data.

# 1 Introduction

It is difficult to detect outliers in  $p$ -variate data when  $p > 2$  because one can no longer rely on visual inspection. While it is still quite easy to detect a single outlier by means of the Mahalanobis distances, this approach no longer suffices for multiple outliers because of the masking effect, by which multiple outliers do not necessarily have large Mahalanobis distances. It is better to use distances based on robust estimators of multivariate location and scatter (Rousseeuw and Leroy 1987, pages 265–269). In regression analysis, robust distances computed from the explanatory variables allow us to detect leverage points. Moreover, robust estimation of multivariate location and scatter is the key tool to robustify other multivariate techniques such as principal component analysis and discriminant analysis.

Many methods for estimating multivariate location and scatter break down in the presence of  $n/(p + 1)$  outliers, where  $n$  is the number of observations and  $p$  is the number of variables, as was pointed out by Donoho (1982). For the breakdown value of the multivariate  $M$ -estimators of Maronna (1976) see (Hampel et al. 1986, page 296). In the meantime, several positive-breakdown estimators of multivariate location and scatter have been proposed. One of these is the minimum volume ellipsoid (MVE) method of Rousseeuw (1984, page 877; 1985). This approach looks for the ellipsoid with smallest volume that covers  $h$  data points, where  $n/2 \leq h < n$ . Its breakdown value is essentially  $(n - h)/n$ .

Positive-breakdown methods such as the MVE and least trimmed squares regression (Rousseeuw 1984) are increasingly being used in practice, e.g. in finance, chemistry, electrical engineering, process control, and computer vision (Meer et al. 1991). For a survey of positive-breakdown methods and some substantive applications, see (Rousseeuw 1997).

The basic resampling algorithm for approximating the MVE, called MINVOL, was proposed in (Rousseeuw and Leroy 1987). This algorithm considers a trial subset of  $p + 1$  observations and calculates its mean and covariance matrix. The corresponding ellipsoid is then inflated or deflated to contain exactly  $h$  observations. This procedure is repeated many times, and the ellipsoid with the lowest volume is retained. For small data sets it is possible to consider all subsets of size  $p + 1$ , whereas for larger data sets the trial subsets are drawn at random.

Several other algorithms have been proposed to approximate the MVE. Woodruff and Rocke (1993) constructed algorithms combining the resampling principle with three heuristic search techniques: simulated annealing, genetic algorithms, and tabu search. Other people

developed algorithms to compute the MVE exactly. This work started with the algorithm of Cook, Hawkins and Weisberg (1992) which carries out an ingenious but still exhaustive search of all possible subsets of size  $h$ . In practice, this can be done for  $n$  up to about 30. Recently, Agulló (1996) developed an exact algorithm for the MVE which is based on a branch and bound procedure that selects the optimal subset without requiring the inspection of all subsets of size  $h$ . This is substantially faster, and can be applied up to (roughly)  $n \leq 100$  and  $p \leq 5$ . Since for most data sets the exact algorithms would take too long, the MVE is typically computed by versions of MINVOL, e.g. in S-Plus (see the function ‘cov.mve’).

Nowadays there are several reasons for replacing the MVE by the minimum covariance determinant (MCD) estimator, which was also proposed in (Rousseeuw 1984, page 877; 1985). The MCD objective is to find  $h$  observations (out of  $n$ ) whose classical covariance matrix has the lowest determinant. The MCD estimate of location is then the average of these  $h$  points, whereas the MCD estimate of scatter is their covariance matrix. The resulting breakdown value equals that of the MVE, but the MCD has several advantages over the MVE. Its statistical efficiency is better, because the MCD is asymptotically normal (Butler, Davies and Jhun 1993; Croux and Haesbroeck 1998) whereas the MVE has a lower convergence rate (Davies 1992). As an example, the asymptotic efficiency of the MCD scatter matrix with the typical coverage  $h = 0.75n$  is 44% in 10 dimensions, and the reweighted covariance matrix with weights obtained from the MCD attains 83% of efficiency (Croux and Haesbroeck 1998), whereas the MVE attains 0%. The MCD’s better accuracy makes it very useful as an initial estimate for one-step regression estimators (Simpson, Ruppert and Carroll 1992; Coakley and Hettmansperger 1993). Robust distances based on the MCD are more precise than those based on the MVE, and hence better suited to expose multivariate outliers, e.g. in the diagnostic plot of Rousseeuw and van Zomeren (1990) which displays robust residuals versus robust distances. Moreover, the MCD is a key component of the hybrid estimators of (Woodruff and Rocke 1994) and (Rocke and Woodruff 1996), and of high-breakdown linear discriminant analysis (Hawkins and McLachlan 1997).

In spite of all these advantages, until now the MCD has rarely been applied because it was harder to compute. However, in this paper we construct a new MCD algorithm which is actually much *faster* than any existing MVE algorithm. The new MCD algorithm can deal with a sample size  $n$  in the tens of thousands. As far as we know none of the existing MVE

algorithms can cope with such large sample sizes. Since the MCD now greatly outperforms the MVE in terms of both statistical efficiency and computation speed, we recommend the MCD method.

## 2 Motivating Problems

Two recent problems will be shown to illustrate the need for a fast robust method which can deal with a large number of objects ( $n$ ) and/or a large number of variables ( $p$ ), while maintaining a reasonable statistical efficiency.

**Problem 1 (Engineering).** We are grateful to Gertjan Otten for providing the following problem. Philips Mecoma (The Netherlands) is producing diaphragm parts for TV sets. These are thin metal plates, molded by a press. Recently a new production line was started, and for each of  $n = 677$  parts, nine characteristics were measured. The aim of the multivariate analysis is to gain insight in the production process and the interrelations between the nine measurements, and to find out whether deformations or abnormalities have occurred and why. Afterwards, the estimated location and scatter matrix can be used for multivariate statistical process control (MSPC).

Due to the support of Herman Veraa and Frans Van Dommelen (at Philips PMF/Mecoma, Product Engineering, P.O. Box 218, 5600 MD Eindhoven, The Netherlands) we obtained permission to analyze these data and to publish the results.

Figure 1 shows the classical Mahalanobis distance

$$\text{MD}(\mathbf{x}_i) = \sqrt{(\mathbf{x}_i - \mathbf{T}_0)' \mathbf{S}_0^{-1} (\mathbf{x}_i - \mathbf{T}_0)} \quad (2.1)$$

versus the index  $i$ , which corresponds to the production sequence. Here  $\mathbf{x}_i$  is 9-dimensional,  $\mathbf{T}_0$  is the arithmetic mean, and  $\mathbf{S}_0$  is the classical covariance matrix. The horizontal line is at the usual cutoff value  $\sqrt{\chi_{9,0.975}^2} = 4.36$ .

In Figure 1 it seems that everything is basically okay (except for a few isolated outliers). This should not surprise us, even in the first experimental run of a new production line, because the Mahalanobis distances are known to suffer from masking. That is, even if there were a group of outliers (here, deformed diaphragm parts) they would affect  $\mathbf{T}_0$  and  $\mathbf{S}_0$  in such a way as to become invisible in Figure 1. To get any further we need robust estimators

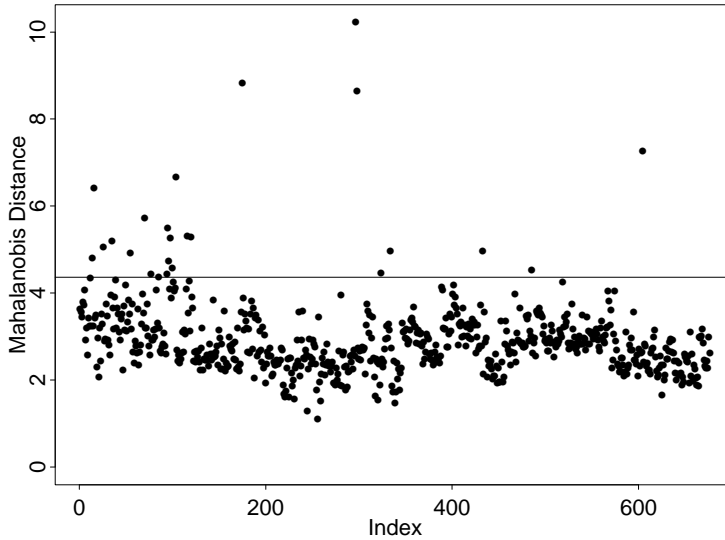


Figure 1. Plot of Mahalanobis distances for the Philips data.

$\mathbf{T}$  and  $\mathbf{S}$ , preferably with a substantial statistical efficiency so that we can be confident that any effects that may become visible are real and not due to the estimator’s inefficiency. After developing the FAST-MCD algorithm, we’ll come back to these data in Section 8.

**Problem 2 (Physical Sciences).** A group of astronomers at Cal Tech are working on the Digitized Palomar Sky Survey (DPOSS); for a full description see their report (Odehahn, Djorgovsky, Brunner, and Gal 1998). In essence, they make a survey of celestial objects (light sources) for which they record nine characteristics (such as magnitude, area, image moments) in each of three bands: blue, red, and near-infrared. They seek collaboration with statisticians to analyze their data, and gave us access to a part of their database, containing 137,256 celestial objects with all 27 variables.

We started by using Q-Q plots, Box-Cox transforms, selecting one variable out of three variables with near-perfect linear correlation, and other tools of data analysis. One of these avenues led us to study six variables (two from each band). Figure 2 plots the Mahalanobis distances (2.1) for these data (to avoid overplotting, Figure 2 shows only 10,000 randomly drawn points from the entire plot). The cutoff is at  $\sqrt{\chi_{6,0.975}^2} = 3.82$ . In Figure 2 we see two groups of outliers with  $\text{MD}(\mathbf{x}_i) \approx 9$  and  $\text{MD}(\mathbf{x}_i) \approx 12$ , plus some outliers still further away.

Returning to the data and their astronomical meaning, it turned out that these were all objects for which one or more variables fell outside the range of what is physically possible.

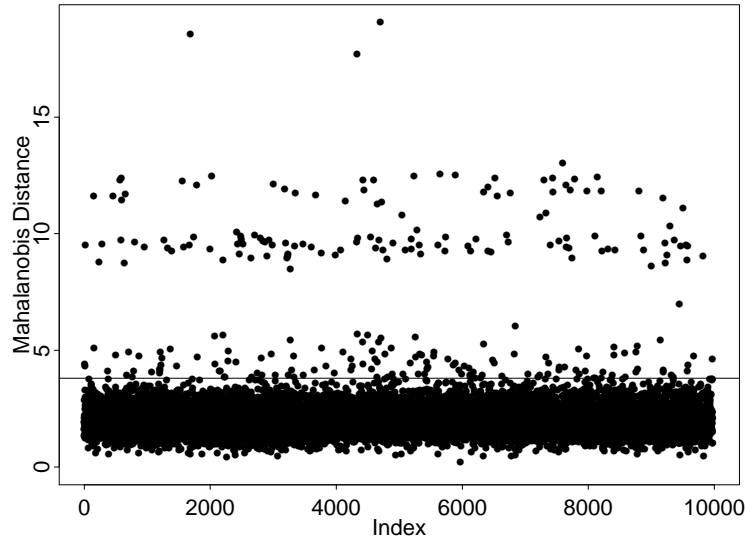


Figure 2. Digitized Palomar Data: plot of Mahalanobis distances of celestial objects, based on 6 variables concerning magnitude and image moments.

So, the  $MD(\mathbf{x}_i)$  did help us to find outliers at this stage. We then cleaned the data by removing all objects with a physically impossible measurement, which reduced our sample size to 132,402. To these data we then again applied the classical mean  $\mathbf{T}_0$  and covariance  $\mathbf{S}_0$ , yielding the plot of Mahalanobis distances in Figure 3.

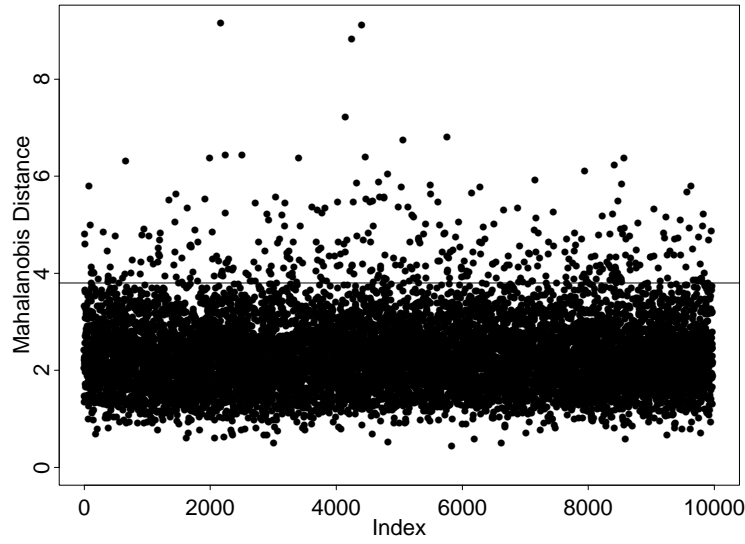


Figure 3. Digitized Palomar Data: plot of Mahalanobis distances of celestial objects as in Figure 2, after removal of physically impossible measurements.

Figure 3 looks innocent, like observations from the  $\sqrt{\chi_6^2}$  distribution, as if the data would form a homogeneous population (which is doubtful, because we know that the database contains stars as well as galaxies). To get any further, we need high-breakdown estimates  $\mathbf{T}$  and  $\mathbf{S}$ , and an algorithm that can compute them for  $n = 132,402$ . Such an algorithm will be constructed in the next sections.

### 3 Basic theorem and the C-step

A key step of the new algorithm is the fact that starting from any approximation to the MCD, it is possible to compute another approximation with an even lower determinant.

**Theorem 1.** Consider a data set  $X_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of  $p$ -variate observations. Let  $H_1 \subset \{1, \dots, n\}$  with  $|H_1| = h$ , and put  $\mathbf{T}_1 := \frac{1}{h} \sum_{i \in H_1} \mathbf{x}_i$  and  $\mathbf{S}_1 := \frac{1}{h} \sum_{i \in H_1} (\mathbf{x}_i - \mathbf{T}_1)(\mathbf{x}_i - \mathbf{T}_1)'$ . If  $\det(\mathbf{S}_1) \neq 0$  define the relative distances

$$d_1(i) := \sqrt{(\mathbf{x}_i - \mathbf{T}_1)' \mathbf{S}_1^{-1} (\mathbf{x}_i - \mathbf{T}_1)} \quad \text{for } i = 1, \dots, n.$$

Now take  $H_2$  such that  $\{d_1(i); i \in H_2\} := \{(d_1)_{1:n}, \dots, (d_1)_{h:n}\}$  where  $(d_1)_{1:n} \leq (d_1)_{2:n} \leq \dots \leq (d_1)_{n:n}$  are the ordered distances, and compute  $\mathbf{T}_2$  and  $\mathbf{S}_2$  based on  $H_2$ . Then

$$\det(\mathbf{S}_2) \leq \det(\mathbf{S}_1)$$

with equality if and only if  $\mathbf{T}_2 = \mathbf{T}_1$  and  $\mathbf{S}_2 = \mathbf{S}_1$ .

The proof is given in the Appendix. Although this theorem appears to be quite basic, we have been unable to find it in the literature.

The theorem requires that  $\det(\mathbf{S}_1) \neq 0$ , which is no real restriction because if  $\det(\mathbf{S}_1) = 0$  we already have the minimal objective value. Section 6 will explain how to interpret the MCD in such a singular situation.

If  $\det(\mathbf{S}_1) > 0$ , applying the theorem yields  $\mathbf{S}_2$  with  $\det(\mathbf{S}_2) \leq \det(\mathbf{S}_1)$ . In our algorithm we will refer to the construction in Theorem 1 as a *C-step*, where C can be taken to stand for ‘covariance’ since  $\mathbf{S}_2$  is the covariance matrix of  $H_2$ , or for ‘concentration’ since we concentrate on the  $h$  observations with smallest distances, and  $\mathbf{S}_2$  is more concentrated (has a lower determinant) than  $\mathbf{S}_1$ . In algorithmic terms, the C-step can be described as follows.

Given the  $h$ -subset  $H_{\text{old}}$  or the pair  $(\mathbf{T}_{\text{old}}, \mathbf{S}_{\text{old}})$ :

- compute the distances  $d_{\text{old}}(i)$  for  $i = 1, \dots, n$
- sort these distances, which yields a permutation  $\pi$  for which  
 $d_{\text{old}}(\pi(1)) \leq d_{\text{old}}(\pi(2)) \leq \dots \leq d_{\text{old}}(\pi(n))$
- put  $H_{\text{new}} := \{\pi(1), \pi(2), \dots, \pi(h)\}$
- compute  $\mathbf{T}_{\text{new}} := \text{ave}(H_{\text{new}})$  and  $\mathbf{S}_{\text{new}} := \text{cov}(H_{\text{new}})$ .

For a fixed number of dimensions  $p$ , the C-step takes only  $O(n)$  time (because  $H_{\text{new}}$  can be determined in  $O(n)$  operations without sorting all the  $d_{\text{old}}(i)$  distances).

Repeating C-steps yields an iteration process. If  $\det(\mathbf{S}_2) = 0$  or  $\det(\mathbf{S}_2) = \det(\mathbf{S}_1)$  we stop; otherwise we run another C-step yielding  $\det(\mathbf{S}_3)$ , and so on. The sequence  $\det(\mathbf{S}_1) \geq \det(\mathbf{S}_2) \geq \det(\mathbf{S}_3) \geq \dots$  is nonnegative and hence must converge. In fact, since there are only finitely many  $h$ -subsets there must be an index  $m$  such that  $\det(\mathbf{S}_m) = 0$  or  $\det(\mathbf{S}_m) = \det(\mathbf{S}_{m-1})$ , hence convergence is reached. (In practice,  $m$  is often below 10.) Afterwards, running the C-step on  $(\mathbf{T}_m, \mathbf{S}_m)$  no longer reduces the determinant. This is not sufficient for  $\det(\mathbf{S}_m)$  to be the global minimum of the MCD objective function, but it is a necessary condition.

Theorem 1 thus provides a partial idea for an algorithm:

$$\begin{aligned} & \textit{Take many initial choices of } H_1 \textit{ and apply C-steps to each until} \\ & \textit{convergence, and keep the solution with lowest determinant.} \end{aligned} \tag{3.1}$$

Of course, several things must be decided to make (3.1) operational: how to generate sets  $H_1$  to begin with, how many  $H_1$  are needed, how to avoid duplication of work since several  $H_1$  may yield the same solution, can't we do with fewer C-steps, what about large sample sizes, and so on. These matters will be discussed in the next sections.

**Corollary 1.** The MCD subset  $H$  of  $X_n$  is separated from  $X_n \setminus H$  by an ellipsoid.

*Proof.* For the MCD subset  $H$ , and in fact any limit of a C-step sequence, applying the C-step to  $H$  yields  $H$  itself. This means that all  $\mathbf{x}_i \in H$  satisfy  $(\mathbf{x}_i - \mathbf{T})' \mathbf{S}^{-1} (\mathbf{x}_i - \mathbf{T}) \leq m = \{(\mathbf{x} - \mathbf{T})' \mathbf{S}^{-1} (\mathbf{x} - \mathbf{T})\}_{h:n}$  whereas all  $\mathbf{x}_j \notin H$  satisfy  $(\mathbf{x}_j - \mathbf{T})' \mathbf{S}^{-1} (\mathbf{x}_j - \mathbf{T}) \geq m$ . Take the ellipsoid  $E = \{\mathbf{x} \in \mathbb{R}^p; (\mathbf{x} - \mathbf{T})' \mathbf{S}^{-1} (\mathbf{x} - \mathbf{T}) \leq m\}$ . Then  $H \subset E$  and  $X_n \setminus H \subset \text{closure}(E^c)$ . Note that there is at least one point  $\mathbf{x}_i \in H$  on the boundary of  $E$ , whereas there may or may not be a point  $\mathbf{x}_j \notin H$  on the boundary of  $E$ .  $\square$



The same result was proved by Butler, Davies and Jhun (1993) under the extra condition that a density exists. Note that the ellipsoid in Corollary 1 contains  $h$  observations but is not necessarily the smallest ellipsoid to do so, which would yield the MVE. We know of no technique like the C-step for the MVE estimator, hence the latter estimator cannot be computed faster in this way.

Independently of our work, Hawkins (1997) discovered a version of Corollary 1 in the following form: ‘A necessary condition for the MCD optimum is that if we calculate the distance of each case from the location vector using the scatter matrix, each covered case must have smaller distance than any uncovered case.’ This necessary condition could be called the ‘C-condition’, as opposed to the C-step of Theorem 1 where we proved that a C-step always decreases  $\det(\mathbf{S})$ . In the absence of Theorem 1, Hawkins (1997) used the C-condition as a preliminary screen, followed by case swapping as a technique for decreasing  $\det(\mathbf{S})$ , as in his feasible solution approach (Hawkins 1994) which will be described in Section 7 below. The C-condition did not reduce the time complexity of his approach, but it did reduce the actual computation time in experiments with fixed  $n$ .

## 4 Construction of the new algorithm

### 4.1 Creating initial subsets $H_1$

In order to apply the algorithmic idea (3.1) we first have to decide how to construct the initial subsets  $H_1$ . Let us consider the following two possibilities:

- (a) Draw a random  $h$ -subset  $H_1$ .
- (b) Draw a random  $(p + 1)$ -subset  $J$ , and then compute  $\mathbf{T}_0 := \text{ave}(J)$  and  $\mathbf{S}_0 := \text{cov}(J)$ . If  $\det(\mathbf{S}_0) = 0$  then extend  $J$  by adding another random observation, and continue adding observations until  $\det(\mathbf{S}_0) > 0$ . Then compute the distances  $d_0^2(i) := (\mathbf{x}_i - \mathbf{T}_0)' \mathbf{S}_0^{-1} (\mathbf{x}_i - \mathbf{T}_0)$  for  $i = 1, \dots, n$ . Sort them into  $d_0(\pi(1)) \leq \dots \leq d_0(\pi(n))$  and put  $H_1 := \{\pi(1), \dots, \pi(h)\}$ .

Option (a) is the simplest, whereas (b) starts like the MINVOL algorithm (Rousseeuw and Leroy 1987, pages 259–260). It would be useless to draw fewer than  $p + 1$  points, for then  $\mathbf{S}_0$  is always singular.

When the data set doesn't contain outliers or deviating groups of points, it makes little difference whether (3.1) is applied with (a) or (b). But since the MCD is a very robust estimator, we have to consider contaminated data sets in particular. For instance, we generated a data set with  $n = 400$  observations and  $p = 2$  variables, in which 205 observations were drawn from the bivariate normal distribution

$$N_2 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \right)$$

and the other 195 observations were drawn from

$$N_2 \left( \begin{bmatrix} 10 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \right).$$

The MCD has its highest possible breakdown value when  $h = \lceil (n + p + 1)/2 \rceil$  (see Lopuhaä and Rousseeuw 1991) which becomes  $h = 201$  here. We now apply (3.1) with 500 starting sets  $H_1$ . Using option (a) yields a resulting  $(\mathbf{T}, \mathbf{S})$  whose 97.5% tolerance ellipse is shown in Figure 4a. Clearly, this result has broken down due to the contaminated data. On the other hand, option (b) yields the correct (robust) result in Figure 4b.

The situation in Figure 4 is extreme, but it is useful for illustrative purposes. (The same effect also occurs for smaller amounts of contamination, especially in higher dimensions.) Approach (a) has failed because each random subset  $H_1$  contains a sizeable number of points from the majority group as well as from the minority group, which follows from the law of large numbers. When starting from a bad subset  $H_1$  the iterations will not converge to the good solution. On the other hand, the probability of a  $(p + 1)$ -subset without outliers is much higher, which explains why (b) yields many good initial subsets  $H_1$  and hence a robust result. From now on, we will always use (b).

*Remark.* For increasing  $n$ , the probability of having at least one 'clean'  $(p + 1)$ -subset among  $m$  random  $(p + 1)$ -subsets tends to

$$1 - (1 - (1 - \epsilon)^{p+1})^m > 0 \tag{4.1}$$

where  $\epsilon$  is the percentage of outliers. In contrast, the probability of having at least one clean  $h$ -subset among  $m$  random  $h$ -subsets tends to zero.

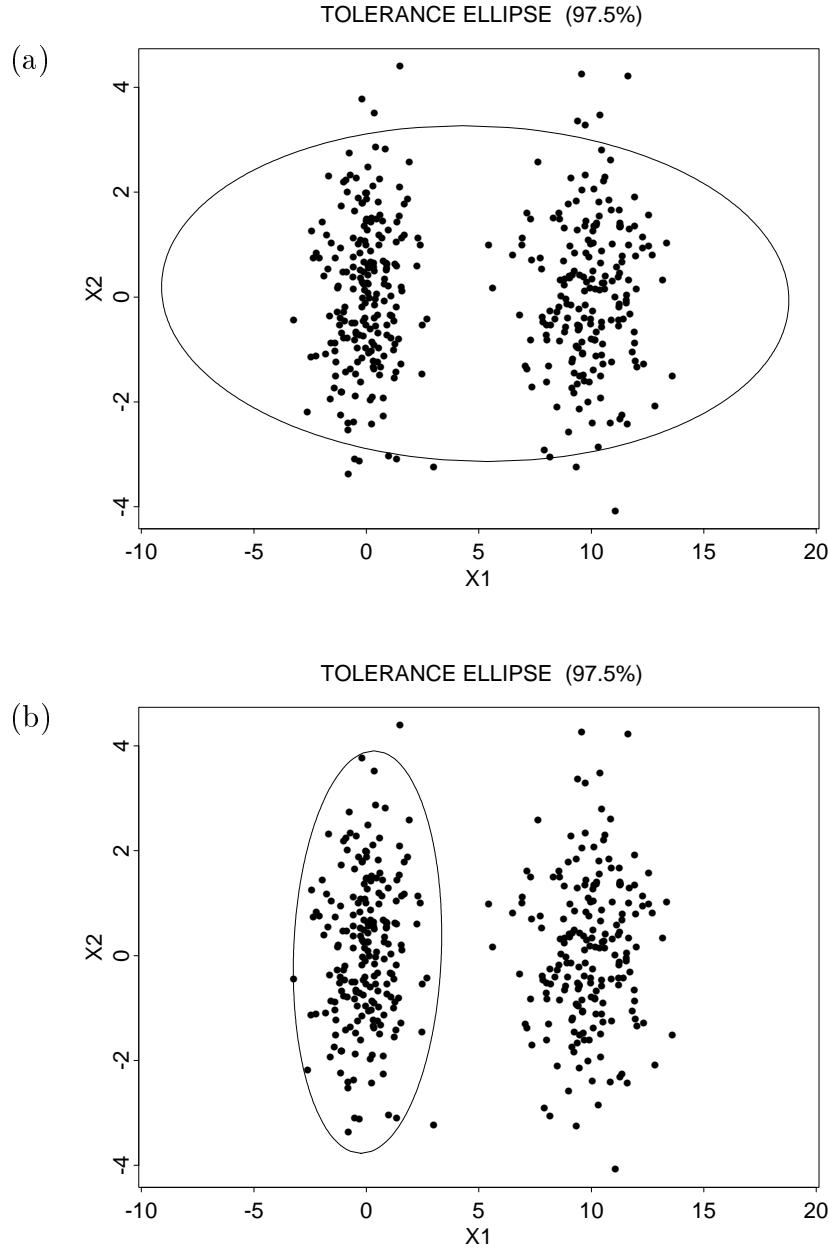


Figure 4. Results of iterating C-steps starting from 500 random subsets  $H_1$  of (a) size  $h=201$ ; and (b) size  $p+1=3$ .

## 4.2 Selective iteration

Each C-step calculates a covariance matrix, its determinant, and all relative distances. Therefore, reducing the number of C-steps would improve the speed. But is this possible without losing the effectiveness of the algorithm? It turns out that often the distinction between good (robust) solutions and bad solutions already becomes visible after two or

three C-steps. For instance, consider the data of Figure 4 again. The inner workings of the algorithm (3.1) are traced in Figure 5. For each starting subsample  $H_1$ , the determinant of the covariance matrix  $S_j$  based on  $h = 201$  observations is plotted versus the step number  $j$ . The runs yielding a robust solution are shown as solid lines, whereas the dashed lines correspond to non-robust results. To get a clear picture, Figure 5 only shows the first 100

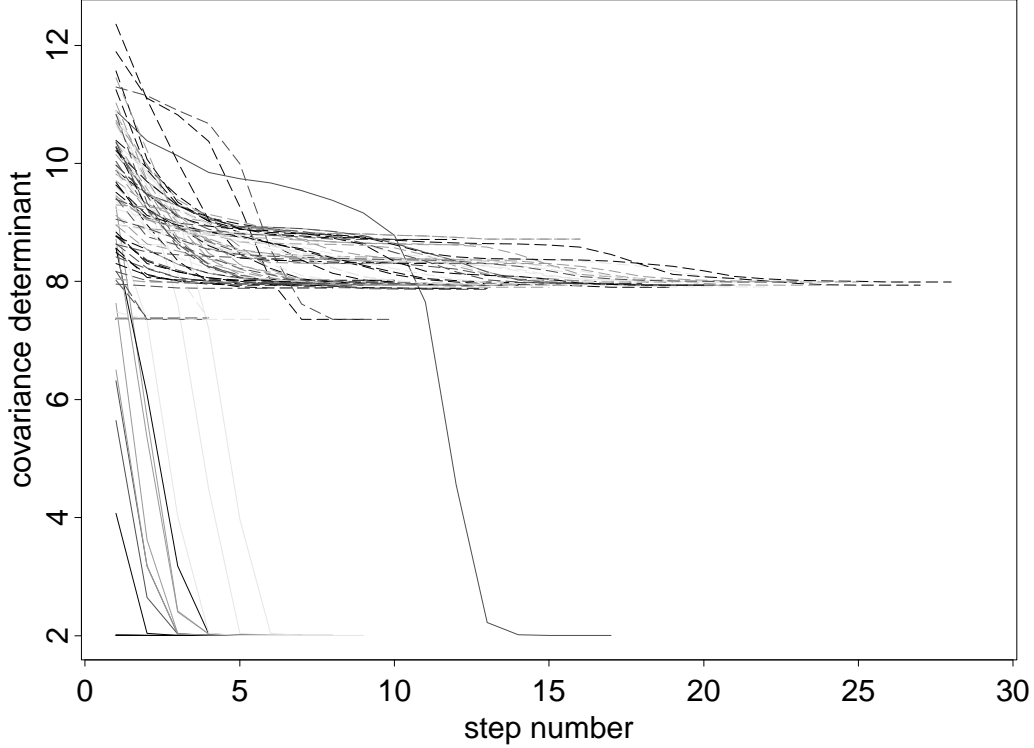


Figure 5. Covariance determinant of subsequent C-steps in the data set of Figure 4. Each sequence stops when no further reduction is obtained.

starts. After two C-steps (i.e. for  $j = 3$ ), many subsamples  $H_3$  that will lead to the global optimum already have a rather small determinant. The global optimum is a solution which contains none of the 195 ‘bad’ points. By contrast, the determinants of the subsets  $H_3$  leading to a false classification are considerably larger. For that reason, we can save much computation time and still obtain the same result by taking just two C-steps and retaining only the (say, 10) best  $H_3$  subsets to iterate further. Other data sets, also in more dimensions, confirm these conclusions. Therefore, from now on we will take only two C-steps from each initial subsample  $H_1$ , select the 10 different subsets  $H_3$  with the lowest determinants, and only for these 10 we continue taking C-steps until convergence.

### 4.3 Nested extensions

For a small sample size  $n$  the above algorithm does not take much time. But when  $n$  grows the computation time increases, mainly due to the  $n$  distances that need to be calculated each time. To avoid doing all the computations in the entire data set, we will consider a special structure. When  $n > 1500$ , the algorithm generates a nested system of subsets which looks like Figure 6, where the arrows mean ‘is a subset of’. The five subsets

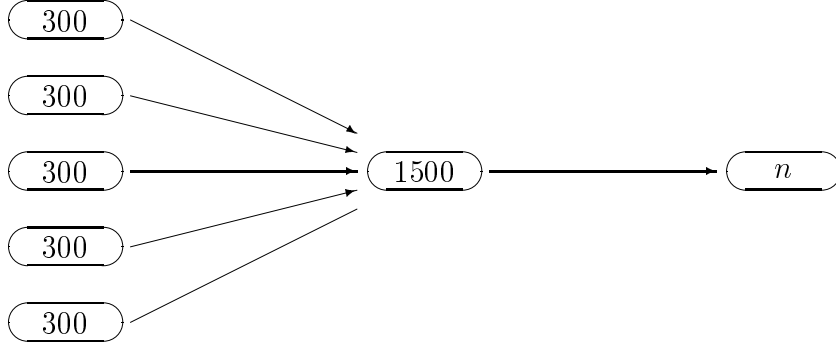


Figure 6. Nested system of subsets generated by the FAST-MCD algorithm.

of size 300 do not overlap, and together they form the merged set of size 1500, which in turn is a proper subset of the data set of size  $n$ . (Already the algorithm of Woodruff and Rocke (1994) made use of partitioning for this purpose. The only difference with the nested extensions in Figure 6 is that we work with two stages, hence our use of the word ‘nested’, whereas Woodruff and Rocke partition the entire data set which yields more and/or larger subsets.) To construct Figure 6 the algorithm draws 1500 observations, one by one, without replacement. The first 300 observations it encounters are put in the first subset, and so on. Because of this mechanism each subset of size 300 is roughly representative for the data set, and the merged set with 1500 cases is even more representative.

When  $n \leq 600$  we will keep the algorithm as in the previous section, while for  $n \geq 1500$  we will use Figure 6. When  $600 < n < 1500$  we will partition the data into at most 4 subsets of 300 or more observations, so that each observation belongs to a subset and such that the subsets have roughly the same size. For instance, 601 will be split as 300+301 and 900 as 450+450. For  $n = 901$  we use 300+300+301, and we continue until  $1499 = 375 + 375 + 375 + 374$ . By splitting 601 as 300+301 we don’t mean that the first subset

contains the observations with case numbers 1, ..., 300 but that its 300 case numbers were drawn randomly from 1, ..., 601.

Whenever  $n > 600$  (and whether  $n < 1500$  or not), our new algorithm for the MCD will take two C-steps from several starting subsamples  $H_1$  within each subset, with a total of 500 starts for all subsets together. For every subset the best 10 solutions are stored. Then the subsets are pooled, yielding a merged set with at most 1500 observations. Each of these (at most 50) available solutions  $(\mathbf{T}_{\text{sub}}, \mathbf{S}_{\text{sub}})$  is then extended to the merged set. That is, starting from each  $(\mathbf{T}_{\text{sub}}, \mathbf{S}_{\text{sub}})$  we continue taking C-steps which now use all 1500 observations in the merged set. Only the best 10 solutions  $(\mathbf{T}_{\text{merged}}, \mathbf{S}_{\text{merged}})$  will be considered further. Finally, each of these 10 solutions is extended to the full data set in the same way, and the best solution  $(\mathbf{T}_{\text{full}}, \mathbf{S}_{\text{full}})$  is reported.

Since the final computations are carried out in the entire data set, they take more time when  $n$  increases. In the interest of speed we can limit the number of initial solutions  $(\mathbf{T}_{\text{merged}}, \mathbf{S}_{\text{merged}})$  and/or the number of C-steps in the full data set as  $n$  becomes large.

The main idea of this subsection was to carry out C-steps in a number of nested random subsets, starting with small subsets of around 300 observations and ending with the entire data set of  $n$  observations. Throughout this subsection we have chosen several numbers such as 5 subsets of 300 observations, 500 starts, 10 best solutions, and so on. These choices were based on various empirical trials (not reported here). We implemented our choices as defaults so the user doesn't have to choose anything, but of course the user may change the defaults.

## 5 The resulting algorithm FAST-MCD

Combining all the components of the preceding sections yields the new algorithm, which we will call FAST-MCD. Its pseudocode looks as follows.

1. The default  $h$  is  $\lceil (n + p + 1)/2 \rceil$ , but the user may choose any integer  $h$  with  $\lceil (n + p + 1)/2 \rceil \leq h \leq n$ . The program then reports the MCD's breakdown value  $(n - h + 1)/n$ . If you are sure that the data contains less than 25% of contamination, which is usually the case, a good compromise between breakdown value and statistical efficiency is obtained by putting  $h = \lceil 0.75n \rceil$ .

2. If  $h = n$  then the MCD location estimate  $\mathbf{T}$  is the average of the whole data set, and the MCD scatter estimate  $\mathbf{S}$  is its covariance matrix. Report these and stop.
3. If  $p = 1$  (univariate data) compute the MCD estimate  $(\mathbf{T}, \mathbf{S})$  by the exact algorithm of Rousseeuw and Leroy (1987, pages 171–172) in  $O(n \log n)$  time; then stop.
4. From here on,  $h < n$  and  $p \geq 2$ . If  $n$  is small (say,  $n \leq 600$ ) then
  - repeat (say) 500 times:
    - \* construct an initial  $h$ -subset  $H_1$  using method (b) in Subsection 4.1, i.e. starting from a random  $(p + 1)$ -subset
    - \* carry out two C-steps (described in Section 3)
  - for the 10 results with lowest  $\det(\mathbf{S}_3)$ :
    - \* carry out C-steps until convergence
  - report the solution  $(\mathbf{T}, \mathbf{S})$  with lowest  $\det(\mathbf{S})$
5. If  $n$  is larger (say,  $n > 600$ ) then
  - construct up to five disjoint random subsets of size  $n_{\text{sub}}$  according to Section 4.3 (say, 5 subsets of size  $n_{\text{sub}} = 300$ )
  - inside each subset, repeat  $500/5 = 100$  times:
    - \* construct an initial subset  $H_1$  of size  $h_{\text{sub}} = \lceil n_{\text{sub}}(h/n) \rceil$
    - \* carry out two C-steps, using  $n_{\text{sub}}$  and  $h_{\text{sub}}$
    - \* keep the 10 best results  $(\mathbf{T}_{\text{sub}}, \mathbf{S}_{\text{sub}})$
  - pool the subsets, yielding the merged set (say, of size  $n_{\text{merged}} = 1500$ )
  - in the merged set, repeat for each of the 50 solutions  $(\mathbf{T}_{\text{sub}}, \mathbf{S}_{\text{sub}})$ :
    - \* carry out two C-steps, using  $n_{\text{merged}}$  and  $h_{\text{merged}} = \lceil n_{\text{merged}}(h/n) \rceil$
    - \* keep the 10 best results  $(\mathbf{T}_{\text{merged}}, \mathbf{S}_{\text{merged}})$
  - in the full data set, repeat for the  $m_{\text{full}}$  best results:
    - \* take several C-steps, using  $n$  and  $h$
    - \* keep the best final result  $(\mathbf{T}_{\text{full}}, \mathbf{S}_{\text{full}})$

Here,  $m_{\text{full}}$  and the number of C-steps (preferably, until convergence) depend on how large the data set is.

We will refer to the above as the FAST-MCD algorithm. Note that it is affine equivariant: when the data are translated or subjected to a linear transformation, the resulting  $(\mathbf{T}_{\text{full}}, \mathbf{S}_{\text{full}})$  will transform accordingly. For convenience, the computer program contains two more steps:

6. In order to obtain consistency when the data come from a multivariate normal distribution, we put

$$\mathbf{T}_{\text{MCD}} = \mathbf{T}_{\text{full}} \quad \text{and} \quad \mathbf{S}_{\text{MCD}} = \frac{\text{med}_i d_{(\mathbf{T}_{\text{full}}, \mathbf{S}_{\text{full}})}^2(i)}{\chi_{p,0.5}^2} \mathbf{S}_{\text{full}}$$

7. A one-step reweighted estimate is obtained by

$$\begin{aligned} \mathbf{T}_1 &= \left( \sum_{i=1}^n w_i \mathbf{x}_i \right) / \left( \sum_{i=1}^n w_i \right) \\ \mathbf{S}_1 &= \left( \sum_{i=1}^n w_i (\mathbf{x}_i - \mathbf{T}_1)(\mathbf{x}_i - \mathbf{T}_1)' \right) / \left( \sum_{i=1}^n w_i - 1 \right) \end{aligned}$$

where

$$\begin{aligned} w_i &= 1 && \text{if } d_{(\mathbf{T}_{\text{MCD}}, \mathbf{S}_{\text{MCD}})}(i) \leq \sqrt{\chi_{p,0.975}^2} \\ &= 0 && \text{otherwise.} \end{aligned}$$

The program FAST-MCD has been thoroughly tested and can be obtained from our website <http://win-www.uia.ac.be/u/statis/index.html>. It has been incorporated into S-Plus 4.5 (as the function “cov.mcd”) and it will also be in SAS/IML 7.01 (as the function “MCD”).

## 6 Exact fit situations

An important advantage of the FAST-MCD algorithm is that it allows for exact fit situations, i.e. when  $h$  or more observations lie on a hyperplane. Then the algorithm still yields the MCD location  $\mathbf{T}$  and scatter matrix  $\mathbf{S}$ , the latter being singular as it should be. From  $(\mathbf{T}, \mathbf{S})$  the program then computes the equation of the hyperplane.



When  $n$  is larger than (say) 600, the algorithm performs many calculations on subsets of the data. In order to deal with the combination of large  $n$  and exact fits, we added a few things to the algorithm. Suppose that, during the calculations in a subset, we encounter some  $(\mathbf{T}_{\text{sub}}, \mathbf{S}_{\text{sub}})$  with  $\det(\mathbf{S}_{\text{sub}}) = 0$ . Then we know that there are  $h_{\text{sub}}$  or more observations on the corresponding hyperplane. First we check whether  $h$  or more points of the full data set lie on this hyperplane. If so we compute  $(\mathbf{T}_{\text{full}}, \mathbf{S}_{\text{full}})$  as the mean and covariance matrix of all points on the hyperplane, report this final result and stop. If not, we have to continue. Since  $\det(\mathbf{S}_{\text{sub}}) = 0$  is the best solution for that subset, we know that  $(\mathbf{T}_{\text{sub}}, \mathbf{S}_{\text{sub}})$  will be among the 10 best solutions that are passed on. In the merged set we take the set  $H'_1$  of the  $h_{\text{merged}}$  observations with smallest orthogonal distances to the hyperplane, and start the next C-step from  $H'_1$ . Again, it is possible that during the calculations in the merged set we encounter some  $(\mathbf{T}_{\text{merged}}, \mathbf{S}_{\text{merged}})$  with  $\det(\mathbf{S}_{\text{merged}}) = 0$ , in which case we repeat the above procedure.

As an illustration, the data set in Figure 7 consists of 45 observations generated from a bivariate normal distribution, plus 55 observations that were generated on a straight line (using a univariate normal distribution). The FAST-MCD program (with default value  $h = 51$ ) finds this line within 0.3 seconds. A part of the output is shown below.

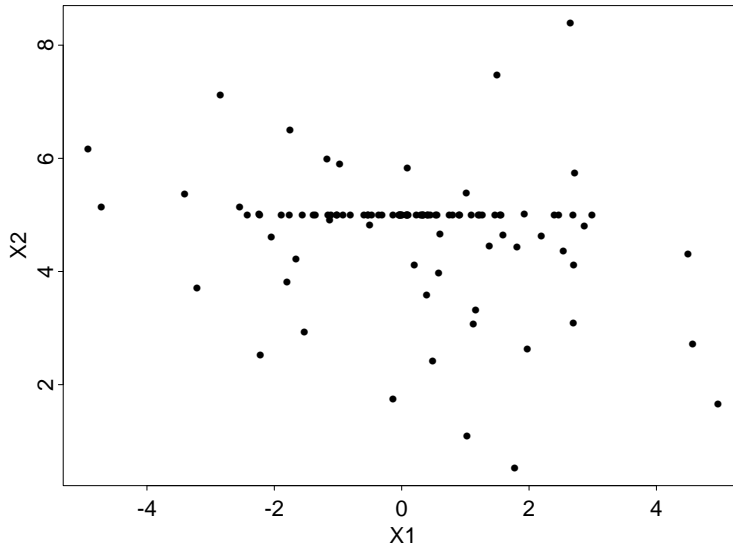


Figure 7. Exact fit situation ( $n = 100$ ,  $p = 2$ ).

There are 55 observations in the entire dataset of 100 observations

that lie on the line with equation

$$0.000000(\mathbf{x}_{i1} - \mathbf{m}_1) + 1.000000(\mathbf{x}_{i2} - \mathbf{m}_2) = 0$$

where the mean  $(\mathbf{m}_1, \mathbf{m}_2)$  of these observations is the MCD location :

0.10817

5.000000

and their covariance matrix is the MCD scatter matrix :

1.40297 0.00000

0.00000 0.00000

Therefore the data are in an "exact fit" position.

In such a situation the MCD scatter matrix has determinant zero,

and its tolerance ellipse becomes the line of exact fit.

If the original data were in  $p$  dimensions and it turns out that most of the data lie on a hyperplane, it is possible to apply FAST-MCD again to the data in this  $(p - 1)$ -dimensional space.

## 7 Performance of FAST-MCD

To get an idea of the performance of the overall algorithm, we start by applying FAST-MCD to some small data sets taken from (Rousseeuw and Leroy 1987). To be precise, these were all regression data sets but we ran FAST-MCD only on the explanatory variables, i.e. not using the response variable. The first column of Table 1 lists the name of each data set, followed by  $n$  and  $p$ . We stayed with the default value of  $h = [(n + p + 1)/2]$ . The next column shows the number of starting  $(p + 1)$ -subsets used in FAST-MCD, which is usually 500 except for two data sets where the number of possible  $(p + 1)$ -subsets out of  $n$  was fairly small, namely  $\binom{12}{3} = 220$  and  $\binom{18}{3} = 816$ , so we used all of them.

The next entry in Table 1 is the result of FAST-MCD, given here as the final  $h$ -subset. By comparing these with the exact MCD algorithm of Agulló (personal communication) it turns out that these  $h$ -subsets do yield the exact global minimum of the objective function. The next column shows the running time of FAST-MCD, in seconds on a Sun Ultra 2170. These times are much shorter than those of our MINVOL program for computing the MVE

Table 1. Performance of the FAST-MCD and FSA algorithms on some small data sets.

data set	$n$	$p$	starts	best $h$ -subset found	time (seconds)	
					FAST-MCD	FSA
Heart	12	2	220	1 3 4 5 7 9 11	0.6	0.6
Phosphor	18	2	816	3 5 8 9 11 12 13 14 15 17	1.8	3.7
Stackloss	21	3	500	4 5 6 7 8 9 10 11 12 13 14 20	2.1	4.6
Coleman	20	5	500	2 3 4 5 7 8 12 13 14 16 17 19 20	4.2	8.9
Wood	20	5	500	1 2 3 5 9 10 12 13 14 15 17 18 20	4.3	8.2
Salinity	28	3	500	1 2 6 7 8 12 13 14 18 20 21 22 25 26 27 28	2.4	8.6
HBK	75	3	500	15 16 17 18 19 20 21 22 23 24 26 27 31 32 33 35 36 37 38 40 43 49 50 51 54 55 56 58 59 61 63 64 66 67 70 71 72 73 74	5.0	71.5

estimator. We may conclude that for these small data sets FAST-MCD gives very accurate results in little time.

Let us now try the algorithm on larger data sets, with  $n \geq 100$ . In each data set over 50% of the points was generated from the standard multivariate normal distribution  $N_p(\mathbf{0}, \mathbf{I}_p)$ , and the remaining points were generated from  $N_p(\boldsymbol{\mu}, \mathbf{I}_p)$  where  $\boldsymbol{\mu} = (b, b, \dots, b)'$  with  $b = 10$ . This is the model of ‘shift outliers’. For each data set Table 2 lists  $n$ ,  $p$ , the percentage of good points, and the percentage of contamination. The algorithm always used 500 starts and the default value of  $h = [(n + p + 1)/2]$ .

The results of FAST-MCD are given in the next column, under ‘robust’. Here ‘yes’ means that the correct result is obtained, i.e. corresponding to the first distribution (as in Figure 4b), whereas ‘no’ stands for the nonrobust result, where the estimates describe the entire data set (as in Figure 4a). Table 2 lists data situations with the *highest* percentage of outlying observations still yielding a robust result with FAST-MCD, as was suggested by a referee. That is, the table says which percentage of outliers the algorithm can handle for given  $n$  and  $p$ . Increasing the value of  $b$  or the number of starts only slightly improves this

Table 2. Performance of the FAST-MCD and FSA algorithms on larger data sets, with time in seconds.

$n$	$p$	% good	% bad	FAST-MCD		FSA	
				robust	time	robust	time
100	2	51	49	yes	2	yes	50
	5	53	47	yes	5	no	80
	10	63	37	yes	40	no	110
	20	77	23	yes	70	no	350
500	2	51	49	yes	7	no	2,800
	5	51	49	yes	25	no	3,800
	10	64	36	yes	84	no	4,100
	30	77	23	yes	695	no	8,300
1,000	2	51	49	yes	8	no	20,000
	5	51	49	yes	20	—	—
	10	60	40	yes	75	—	—
	30	76	24	yes	600	—	—
10,000	2	51	49	yes	9	—	—
	5	51	49	yes	25	—	—
	10	63	37	yes	85	—	—
	30	76	24	yes	700	—	—
50,000	2	51	49	yes	15	—	—
	5	51	49	yes	45	—	—
	10	58	42	yes	140	—	—
	30	75	25	yes	890	—	—

percentage. The computation times were quite low for the given values of  $n$  and  $p$ . Even for a sample size as high as 50,000 a few minutes suffice, whereas no previous algorithm we know of could handle such large data sets.

The currently most well-known algorithm for approximating the MCD estimator is the Feasible Subset Algorithm (FSA) of Hawkins (1994). Instead of C-steps it uses a different kind of steps, which for convenience we will baptise ‘I-steps’ where the ‘I’ stands for

‘interchanging points’. An I-step proceeds as follows:

Given the  $h$ -subset  $H_{\text{old}}$  with its average  $\mathbf{T}_{\text{old}}$  and its covariance matrix  $\mathbf{S}_{\text{old}}$ :

- repeat for each  $i \in H_{\text{old}}$  and each  $j \notin H_{\text{old}}$ :
  - \* put  $H_{i,j} = (H_{\text{old}} \setminus \{i\}) \cup \{j\}$   
(i.e., remove point  $i$  and add point  $j$ )
  - \* compute  $\Delta_{i,j} = \det(\mathbf{S}_{\text{old}}) - \det(\mathbf{S}(H_{i,j}))$
- keep the  $i'$  and  $j'$  with largest  $\Delta_{i',j'}$
- if  $\Delta_{i',j'} \leq 0$  put  $H_{\text{new}} = H_{\text{old}}$  and stop;
- if  $\Delta_{i',j'} > 0$  put  $H_{\text{new}} = H_{i',j'}$ .

An I-step takes  $O(h(n-h)) = O(n^2)$  time because all pairs  $(i, j)$  are considered. If we would compute each  $\mathbf{S}(H_{i,j})$  from scratch the complexity would even become  $O(n^3)$ , but Hawkins (1994, page 203) uses an update formula for  $\det(\mathbf{S}(H_{i,j}))$ .

The I-step can be iterated: if  $\det(\mathbf{S}_{\text{new}}) < \det(\mathbf{S}_{\text{old}})$  we can take another I-step with  $H_{\text{new}}$ , otherwise we stop. The resulting sequence  $\det(\mathbf{S}_1) \geq \det(\mathbf{S}_2) \geq \dots$  must converge after a finite number of steps, i.e.  $\det(\mathbf{S}_m) = 0$  or  $\det(\mathbf{S}_m) = \det(\mathbf{S}_{m-1})$ , so  $\det(\mathbf{S}_m)$  can no longer be reduced by an I-step. This is again a necessary (but not sufficient) condition for  $(\mathbf{T}_m, \mathbf{S}_m)$  to be the global minimum of the MCD objective function. In our terminology, Hawkins’ FSA algorithm can be written as:

- repeat many times:
  - \* draw an initial  $h$ -subset  $H_1$  at random
  - \* carry out I-steps until convergence, yielding  $H_m$
- keep the  $H_m$  with lowest  $\det(\mathbf{S}_m)$
- report this set  $H_m$  as well as  $(\mathbf{T}_m, \mathbf{S}_m)$ .

In Tables 1 and 2 we have applied the FSA algorithm to the same data sets as FAST-MCD, using the same number of starts. For the small data sets in Table 1 the FSA and FAST-MCD yielded identical results. This is no longer true in Table 2, where the FSA begins to find nonrobust solutions. This is because of

- (1) The FSA starts from randomly drawn  $h$ -subsets  $H_1$ . Hence for sufficiently large  $n$  all the FSA starts are nonrobust, and subsequent iterations do not get away from the corresponding local minimum.

We saw the same effect in Section 4.1, which also explained why it is safer to start from random  $(p + 1)$ -subsets as in MINVOL and in FAST-MCD.

The tables also indicate that the FSA needs more time than FAST-MCD. In fact,  $\text{time(FSA)}/\text{time(FAST-MCD)}$  increases from 1 to 14 for  $n$  going from 12 to 75. In Table 2 the timing ratio goes from 25 (for  $n = 100$ ) to 2500 (for  $n = 1000$ ), after which we could no longer time the FSA algorithm. The FSA algorithm is more time-consuming than FAST-MCD because:

- (2) An I-step takes  $O(n^2)$  time, compared to  $O(n)$  for the C-step of FAST-MCD.
- (3) Each I-step swaps only 1 point of  $H_{\text{old}}$  with 1 point outside  $H_{\text{old}}$ . In contrast, each C-step swaps  $h - |H_{\text{old}} \cap H_{\text{new}}|$  points inside  $H_{\text{old}}$  with the same number outside of  $H_{\text{old}}$ . Therefore more I-steps are needed, especially for increasing  $n$ .
- (4) The FSA iterates I-steps until convergence, starting from each  $H_1$ . On the other hand, FAST-MCD reduces the number of C-steps by the selective iteration technique of Section 4.2. The latter would not work for I-steps because of (3).
- (5) The FSA carries out all its I-steps in the full data set of size  $n$ , even for large  $n$ . In the same situation FAST-MCD applies the nested extensions method of Section 4.3, so most C-steps are carried out for  $n_{\text{sub}} = 300$ , some for  $n_{\text{merged}} = 1500$ , and only a few for the actual  $n$ .

In conclusion, we personally prefer the FAST-MCD algorithm because it is both robust and fast, even for large  $n$ .

## 8 Applications

Let us now look at some applications to compare the FAST-MCD results with the classical mean and covariance matrix. At the same time we will illustrate a new tool, the *distance-distance plot*.

**Example 1.** We start with the coho salmon data set (see Nickelson 1986) with  $n = 22$  and  $p = 2$ , as shown in Figure 8a. Each data point corresponds to one year. For 22 years the production of coho salmon in the wild was measured, in the Oregon Production Area. The  $x$ -coordinate is the logarithm of millions of smolts, and the  $y$ -coordinate is the logarithm of millions of adult coho salmons. We see that in most years the production of smolts lies between 2.2 and 2.4 on a logarithmic scale, while the production of adults lies between -1.0 and 0.0. The MCD tolerance ellipse excludes the years with a lower smolts production, thereby marking them as outliers. In contrast, the classical tolerance ellipse contains nearly the whole data set and thus does not detect the existence of far outliers.

Let us now introduce the distance-distance plot (D-D plot), which plots the robust distances (based on the MCD) versus the classical Mahalanobis distances. On both axes in Figure 8b we have indicated the cutoff value  $\sqrt{\chi_{p,0.975}^2}$  (here  $p = 2$ , yielding  $\sqrt{\chi_{2,0.975}^2} = 2.72$ ). If the data were not contaminated (say, if all the data would come from a single bivariate normal distribution) then all points in the distance-distance plot would lie near the dotted line. In this example many points lie in the rectangle where both distances are regular, whereas the outlying points lie higher. This happened because the MCD ellipse and the classical ellipse have a different orientation.

Naturally, the D-D plot becomes more useful in higher dimensions, where it is not so easy to visualize the data set and the ellipsoids.

**Problem 1 (continued).** Next we consider Problem 1 in Section 2. The Philips data represent 677 measurements of metal sheets with 9 components each, and the Mahalanobis distances in Figure 1 indicated no groups of outliers. However, the MCD-based robust distances  $\text{RD}(\mathbf{x}_i)$  in Figure 9a tell a different story. We now see a strongly deviating group of outliers, ranging from index 491 to index 565.

Something happened in the production process, which was not visible from the classical distances shown in Figure 1 in Section 2. Figure 9a also shows a remarkable change after the first 100 measurements. These phenomena were investigated and interpreted by the engineers at Philips. Note that the D-D plot in Figure 9b again contrasts the classical and robust analysis.

**Problem 2 (continued).** We now apply FAST-MCD to the same  $n = 132,402$  celestial objects with  $p = 6$  variables as in Figure 3, which took only 2.5 minutes. (In fact, running

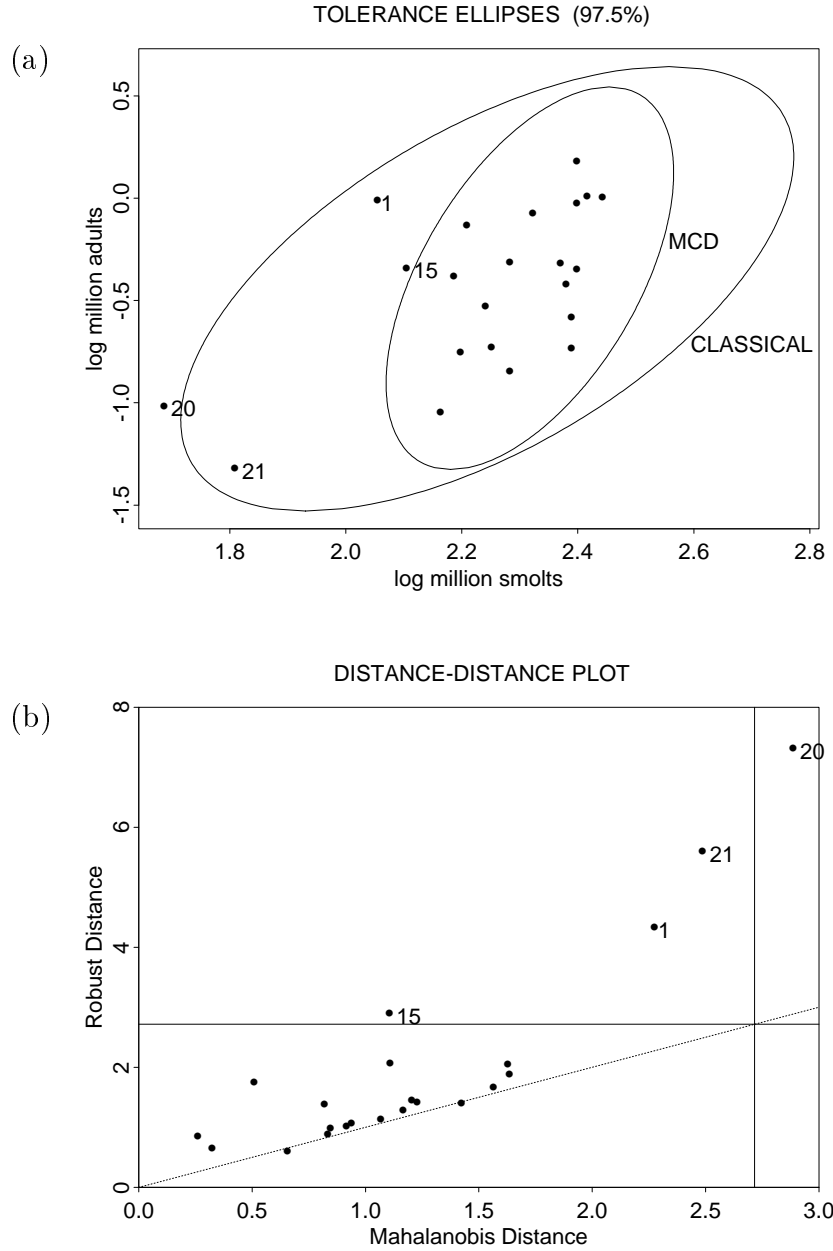


Figure 8. Coho Salmon data: (a) scatterplot with 97.5% tolerance ellipses describing the MCD and the classical method; (b) distance-distance plot.

the program on the same objects in all 27 dimensions took only 18 minutes!) Figure 10a plots the resulting MCD-based robust distances. In contrast to the homogeneous-looking Mahalanobis distances in Figure 3, the robust distances in Figure 10a clearly show that there is a majority with  $RD(\mathbf{x}_i) \leq \sqrt{\chi_{6,0.975}^2}$  as well as a second group with  $RD(\mathbf{x}_i)$  between 8 and 16. By exchanging our findings with the astronomers at Cal Tech it turned out that the lower group consists mainly of stars, and the upper group mainly of galaxies.



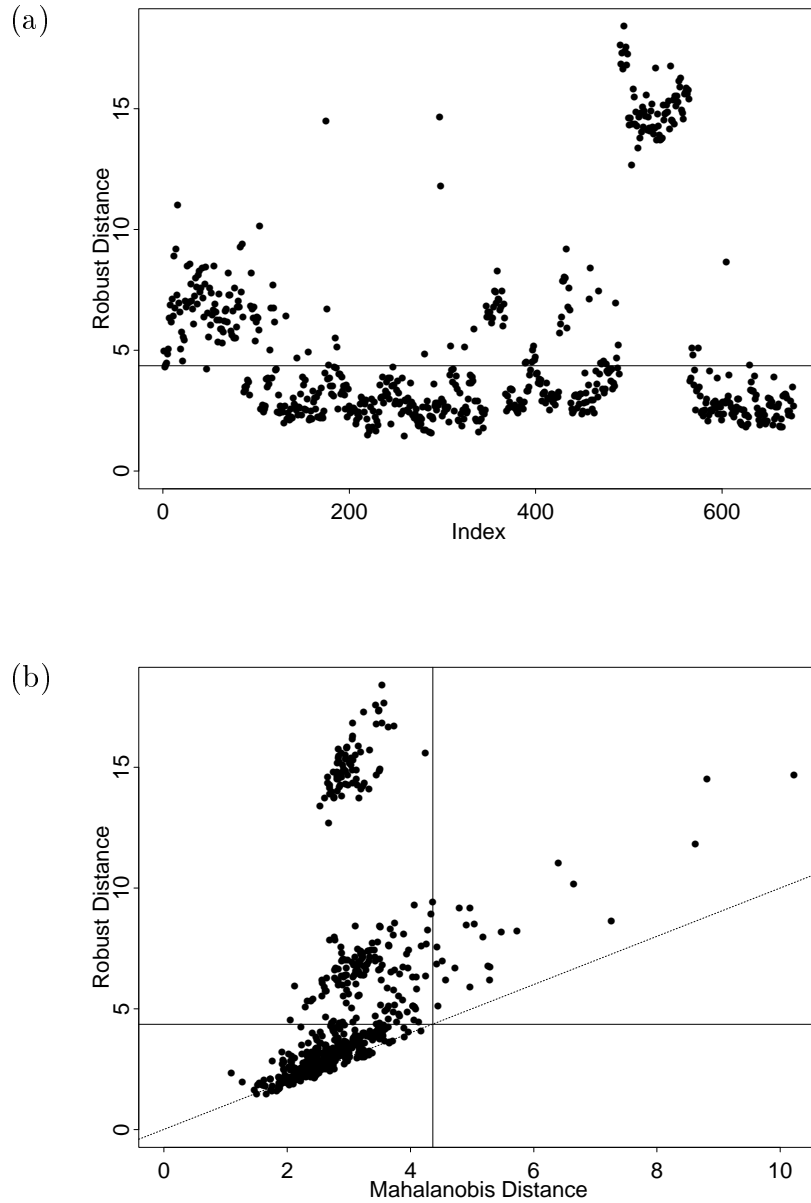


Figure 9. Philips data: (a) plot of robust distances; (b) distance-distance plot.

Our main point is that the robust distances separate the data in two parts, and thus provide more information than the Mahalanobis distances. That these distances behave differently is illustrated in Figure 10b, where we see the stars near the diagonal line and the galaxies above it.

Of course our analysis of these data was much more extensive, and also used other data-analytic techniques not described here, but the ability to compute robust estimates of location and scatter for such large data sets was a key tool. Based on our work, these as-

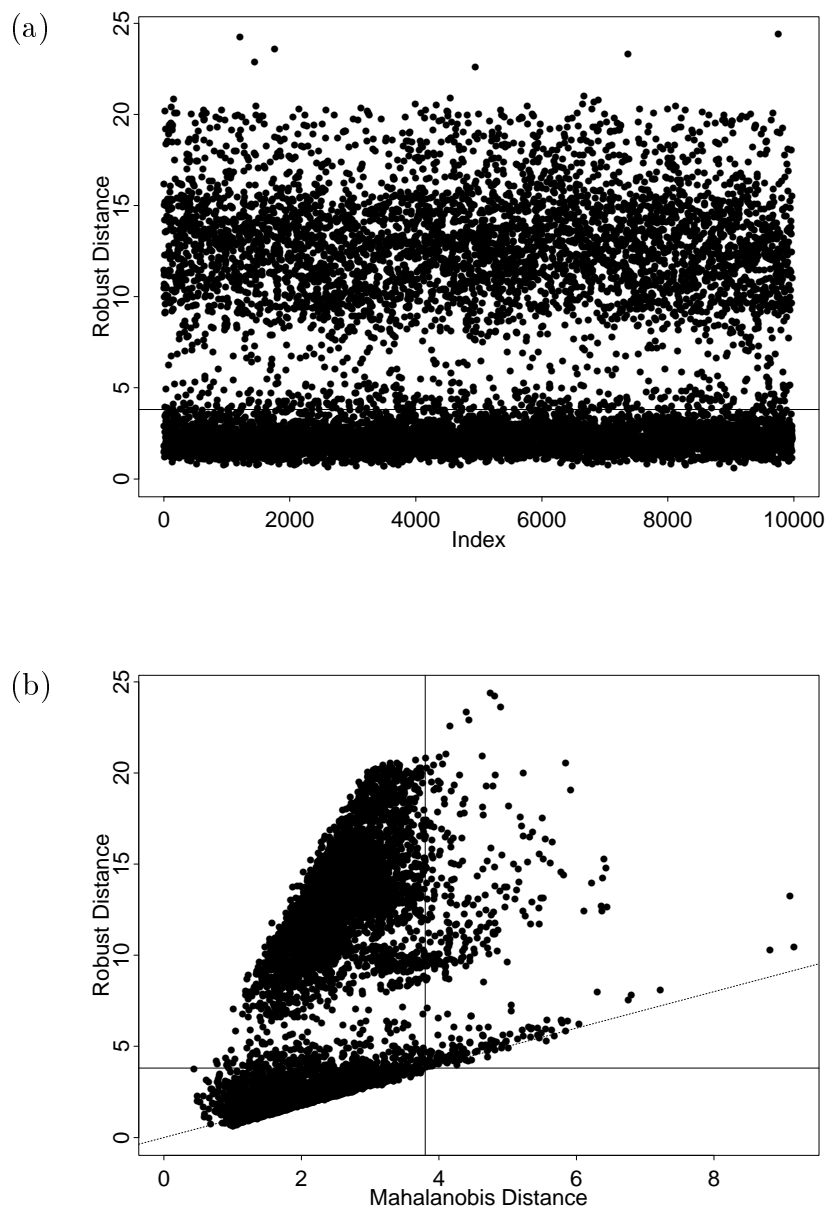


Figure 10. Digitized Palomar data: (a) plot of robust distances of celestial objects; (b) their distance-distance plot.

tronomers are thinking about modifying their classification of objects into stars and galaxies, especially for the faint light sources which are difficult to classify.

**Example 2.** We end this section by combining robust location/scatter with robust regression. The fire data (Andrews and Herzberg 1985) report the incidences of fires in 47 residential areas of Chicago. One wants to explain the incidence of fire by the age of the houses, the

income of the families living in them, and the incidence of theft. For this we apply the least trimmed squares (LTS) method of robust regression, with the usual value of  $h = \lceil \frac{3}{4}n \rceil = 35$ . In S-Plus 4.5 the function ‘ltsreg’ now automatically calls the function ‘cov.mcd’ which runs the FAST-MCD algorithm, in order to obtain robust distances in  $\mathbf{x}$ -space based on the MCD with the same  $h$ . Moreover, S-Plus automatically provides the diagnostic plot of Rousseeuw and van Zomeren (1990), which plots the robust residuals versus the robust distances. For the fire data this yields Figure 11 which shows the presence of one vertical outlier, i.e. an observation with a small robust distance and a large LTS residual. We also see two bad leverage points, which are outlying observations in  $\mathbf{x}$ -space that do not follow the linear trend of the majority of the data. The other observations with robust distances to the right of the vertical cutoff line are good leverage points, since they have regular LTS residuals and hence follow the same linear pattern as the main group. In Figure 11 we see that most of these points are merely boundary cases, except for the two leverage points which are really far out in  $x$ -space.

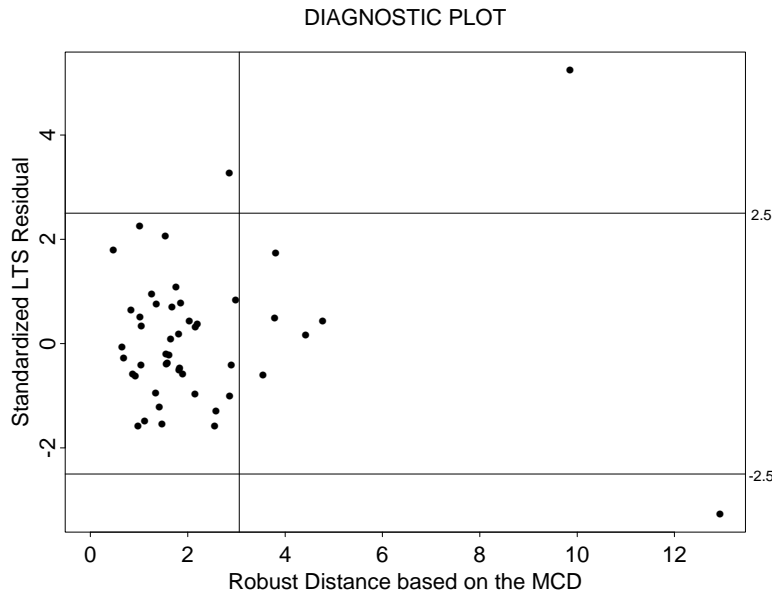


Figure 11. Diagnostic plot of the fire data set.

## 9 Conclusions

The algorithm FAST-MCD proposed in this paper is specifically tailored to the properties of the MCD estimator. The basic ideas are the C-step (Theorem 1 in Section 3),

the procedure for generating initial estimates (Section 4.1), selective iteration (Section 4.2), and nested extensions (Section 4.3). By exploiting the special structure of the problem, the new algorithm is faster and more effective than general-purpose techniques such as reducing the objective function by successively interchanging points. Simulations have shown that FAST-MCD is able to deal with large data sets, while outrunning existing algorithms for MVE and MCD by orders of magnitude. Another advantage of FAST-MCD is its ability to detect exact fit situations.

Due to the FAST-MCD algorithm, the MCD becomes accessible as a routine tool for analyzing multivariate data. Without extra cost we also obtain the distance-distance plot (D-D plot), a new data display which plots the MCD-based robust distances versus the classical Mahalanobis distances. This is a useful tool to explore structure(s) in the data. Other possibilities include an MCD-based PCA, and robustified versions of other multivariate analysis methods.

## Appendix: Proof of Theorem 1

*Proof.* Assume that  $\det(\mathbf{S}_2) > 0$ , otherwise the result is already satisfied. We can thus compute  $d_2(i) = d_{(\mathbf{T}_2, \mathbf{S}_2)}(i)$  for all  $i = 1, \dots, n$ . Using  $|H_2| = h$  and the definition of  $(\mathbf{T}_2, \mathbf{S}_2)$  we find

$$\begin{aligned} \frac{1}{hp} \sum_{i \in H_2} d_2^2(i) &= \frac{1}{hp} \text{tr} \sum_{i \in H_2} (\mathbf{x}_i - \mathbf{T}_2) \mathbf{S}_2^{-1} (\mathbf{x}_i - \mathbf{T}_2)' \\ &= \frac{1}{hp} \text{tr} \sum_{i \in H_2} \mathbf{S}_2^{-1} (\mathbf{x}_i - \mathbf{T}_2) (\mathbf{x}_i - \mathbf{T}_2)' = \frac{1}{p} \text{tr} \mathbf{S}_2^{-1} \mathbf{S}_2 = \frac{1}{p} \text{tr}(\mathbf{I}) = 1. \end{aligned} \quad (\text{A.1})$$

Moreover, put

$$\lambda := \frac{1}{hp} \sum_{i \in H_2} d_1^2(i) = \frac{1}{hp} \sum_{i=1}^h (d_1^2)_{i:n} \leq \frac{1}{hp} \sum_{j \in H_1} d_1^2(j) = 1, \quad (\text{A.2})$$

where  $\lambda > 0$  because otherwise  $\det(\mathbf{S}_2) = 0$ . Combining (A.1) and (A.2) yields

$$\frac{1}{hp} \sum_{i \in H_2} d_{(\mathbf{T}_1, \lambda \mathbf{S}_1)}^2(i) = \frac{1}{hp} \sum_{i \in H_2} (\mathbf{x}_i - \mathbf{T}_1)' \frac{1}{\lambda} \mathbf{S}_1^{-1} (\mathbf{x}_i - \mathbf{T}_1) = \frac{1}{\lambda hp} \sum_{i \in H_2} d_1^2(i) = \frac{\lambda}{\lambda} = 1.$$

Grübel (1988) proved that  $(\mathbf{T}_2, \mathbf{S}_2)$  is the unique minimizer of  $\det(\mathbf{S})$  among all  $(\mathbf{T}, \mathbf{S})$  for which  $\frac{1}{hp} \sum_{i \in H_2} d_{(\mathbf{T}, \mathbf{S})}^2(i) = 1$ . This implies that  $\det(\mathbf{S}_2) \leq \det(\lambda \mathbf{S}_1)$ . On the other hand it

follows from the inequality (A.2) that  $\det(\lambda \mathbf{S}_1) \leq \det(\mathbf{S}_1)$ , hence

$$\det(\mathbf{S}_2) \leq \det(\lambda \mathbf{S}_1) \leq \det(\mathbf{S}_1). \quad (\text{A.3})$$

Moreover, note that  $\det(\mathbf{S}_2) = \det(\mathbf{S}_1)$  if and only if both inequalities in (A.3) are equalities. For the first, we know from Grubel’s result that  $\det(\mathbf{S}_2) = \det(\lambda \mathbf{S}_1)$  if and only if  $(\mathbf{T}_2, \mathbf{S}_1) = (\mathbf{T}_1, \lambda \mathbf{S}_1)$ . For the second,  $\det(\lambda \mathbf{S}_1) = \det(\mathbf{S}_1)$  if and only if  $\lambda = 1$ , i.e.  $\mathbf{S}_1 = \lambda \mathbf{S}_1$ . Combining both yields  $(\mathbf{T}_2, \mathbf{S}_2) = (\mathbf{T}_1, \mathbf{S}_1)$ .  $\square$

## References

- Agullo, J. (1996), “Exact Iterative Computation of the Multivariate Minimum Volume Ellipsoid Estimator with a Branch and Bound Algorithm,” in *Proceedings in Computational Statistics*, edited by Albert Prat. Heidelberg: Physica-Verlag, 175–180.
- Andrews, D.F. and Herzberg, A.M. (1985), *Data*, Springer-Verlag, New York, 407–412.
- Butler, R.W., Davies, P.L., and Jhun, M. (1993), “Asymptotics for the Minimum Covariance Determinant Estimator,” *The Annals of Statistics*, 21, 1385–1400.
- Coakley, C.W. and Hettmansperger, T.P. (1993), “A Bounded Influence, High Breakdown, Efficient Regression Estimator,” *Journal of the American Statistical Association*, 88, 872–880.
- Cook, R.D., Hawkins, D.M., and Weisberg, S. (1992), “Exact Iterative Computation of the Robust Multivariate Minimum Volume Ellipsoid Estimator,” *Statistics and Probability Letters*, 16, 213–218.
- Croux, C., and Haesbroeck, G. (1998), “Influence Function and Efficiency of the Minimum Covariance Determinant Scatter Matrix Estimator,” Preprint, University of Brussels.
- Davies, L. (1992), “The Asymptotics of Rousseeuw’s Minimum Volume Ellipsoid Estimator,” *The Annals of Statistics*, 20, 1828–1843.
- Donoho, D.L. (1982), “Breakdown properties of multivariate location estimators,” Ph.D. Qualifying Paper, Harvard University, Boston.

- Grübel, R. (1988), “A Minimal Characterization of the Covariance Matrix,” *Metrika*, 35, 49–52.
- Hampel, F.R., Ronchetti, E.M., Rousseeuw, P.J., and Stahel, W.A. (1986), *Robust Statistics, The Approach based on Influence Functions*, New York: John Wiley.
- Hawkins, D.M. (1994), “The Feasible Solution Algorithm for the Minimum Covariance Determinant Estimator in Multivariate Data,” *Computational Statistics and Data Analysis*, 17, 197–210.
- Hawkins, D.M. (1997), “Improved Feasible Solution Algorithms for High Breakdown Estimation,” Technical Report, University of Minnesota, September 1997.
- Hawkins, D.M. and McLachlan, G.J. (1997), “High-Breakdown Linear Discriminant Analysis,” *Journal of the American Statistical Association*, 92, 136–143.
- Lopuhaä, H.P. and Rousseeuw, P.J. (1991), “Breakdown Points of Affine Equivariant Estimators of Multivariate Location and Covariance Matrices,” *The Annals of Statistics*, 19, 229–248.
- Maronna, R.A. (1976), “Robust  $M$ -estimators of Multivariate Location and Scatter,” *The Annals of Statistics*, 4, 51–56.
- Meer, P., Mintz, D., Rosenfeld, A. and Kim, D. (1991), “Robust Regression Methods in Computer Vision: A Review,” *International Journal of Computer Vision*, 6, 59–70.
- Nickelson, T.E. (1986), “Influence of Upwelling, Ocean Temperature, and Smolt Abundance on Marine Survival of Coho Salmon (*Oncorhynchus Kisutch*) in the Oregon Production Area,” *Canadian Journal of Fisheries and Aquatic Sciences*, 43, 527–535.
- Odewahn, S.C., Djorgovski, S.G., Brunner, R.J., and Gal, R. (1998), “Data From the Digitized Palomar Sky Survey,” Technical Report, California Institute of Technology.
- Rocke, D.M. and Woodruff, D.L. (1996), “Identification of Outliers in Multivariate Data,” *Journal of the American Statistical Association*, 91, 1047–1061.
- Rousseeuw, P.J. (1984), “Least Median of Squares Regression,” *Journal of the American Statistical Association*, 79, 871–880.

- Rousseeuw, P.J. (1985), “Multivariate Estimation with High Breakdown Point,” in *Mathematical Statistics and Applications, Vol B*, eds. W. Grossmann, G. Pflug, I. Vincze, and W. Wertz. Dordrecht: Reidel, 283–297.
- Rousseeuw, P.J. (1997), “Introduction to Positive-Breakdown Methods,” in *Handbook of Statistics, Vol. 15: Robust Inference*, eds. G.S. Maddala and C.R. Rao. Amsterdam: Elsevier, 101–121.
- Rousseeuw, P.J. and Leroy, A.M. (1987), *Robust Regression and Outlier Detection*, New York: John Wiley.
- Rousseeuw, P.J. and van Zomeren, B.C. (1990), “Unmasking Multivariate Outliers and Leverage Points,” *Journal of the American Statistical Association*, 85, 633–639.
- Simpson, D.G., Ruppert, D., and Carroll, R.J. (1992), “On One-Step *GM*-estimates and Stability of Inferences in Linear Regression,” *Journal of the American Statistical Association*, 87, 439–450.
- Woodruff, D.L. and Rocke, D.M. (1993), “Heuristic Search Algorithms for the Minimum Volume Ellipsoid,” *Journal of Computational and Graphical Statistics*, 2, 69–95.
- Woodruff, D.L. and Rocke, D.M. (1994), “Computable Robust Estimation of Multivariate Location and Shape in High Dimension using Compound Estimators,” *Journal of the American Statistical Association*, 89, 888–896.