

Robust Multivariate Regression

Peter J. ROUSSEEUW

Department of Mathematics and Computer Science
University of Antwerp
Antwerpen, Belgium
(*Peter.Rousseeuw@ua.ac.be*)

Katrien VAN DRIESSEN

Faculty of Applied Economics
University of Antwerp
Antwerpen, Belgium
(*Katrien.VanDriessen@ua.ac.be*)

Stefan VAN AELST

Department of Applied Mathematics and Computer Science
Ghent University
Ghent, Belgium
(*Stefan.VanAelst@Ugent.be*)

Jose AGULLÓ

Department of Economic Analysis
University of Alicante
Alicante, Spain
(*Agullo@merlin.fae.ua.es*)

We introduce a robust method for multivariate regression based on robust estimation of the joint location and scatter matrix of the explanatory and response variables. As a robust estimator of location and scatter, we use the minimum covariance determinant (MCD) estimator of Rousseeuw. Based on simulations, we investigate the finite-sample performance and robustness of the estimator. To increase the efficiency, we propose a reweighted estimator selected from several possible reweighting schemes. The resulting multivariate regression does not need much computation time and is applied to real datasets. We show that the multivariate regression estimator has the appropriate equivariance properties, has a bounded influence function, and inherits the breakdown value of the MCD estimator. These theoretical robustness properties confirm the good finite-sample results obtained from the simulations.

KEY WORDS: Breakdown value; Diagnostic plot; Influence function; Minimum covariance determinant; Reweighting.

1. INTRODUCTION

Suppose that we have a p -variate predictor $\mathbf{x} = (x_1, \dots, x_p)^t$ and a q -variate response $\mathbf{y} = (y_1, \dots, y_q)^t$. The multivariate regression model is given by $\mathbf{y} = \mathbf{B}'\mathbf{x} + \boldsymbol{\alpha} + \boldsymbol{\varepsilon}$, where \mathbf{B} is the $(p \times q)$ slope matrix, $\boldsymbol{\alpha}$ is the q -dimensional intercept vector, and the errors $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_q)^t$ are iid with mean $\mathbf{0}$ and with $\text{cov}(\boldsymbol{\varepsilon}) = \boldsymbol{\Sigma}_{\varepsilon}$ a positive definite matrix of size q . Denote the location of the joint (\mathbf{x}, \mathbf{y}) variables by $\boldsymbol{\mu}$ and their scatter matrix by $\boldsymbol{\Sigma}$. Partitioning $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ yields the notation

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix} \quad \text{and} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{pmatrix}.$$

Traditionally, $\boldsymbol{\mu}$ is often estimated by the empirical mean $\hat{\boldsymbol{\mu}}$ and $\boldsymbol{\Sigma}$ is often estimated by the empirical covariance matrix $\hat{\boldsymbol{\Sigma}}$. It turns out that the least squares estimators of \mathbf{B} , $\boldsymbol{\alpha}$, and $\boldsymbol{\Sigma}_{\varepsilon}$ can be written as functions of the components of $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Sigma}}$, namely

$$\hat{\mathbf{B}} = \hat{\boldsymbol{\Sigma}}_{xx}^{-1} \hat{\boldsymbol{\Sigma}}_{xy}, \quad (1)$$

$$\hat{\boldsymbol{\alpha}} = \hat{\boldsymbol{\mu}}_y - \hat{\mathbf{B}}' \hat{\boldsymbol{\mu}}_x, \quad (2)$$

and

$$\hat{\boldsymbol{\Sigma}}_{\varepsilon} = \hat{\boldsymbol{\Sigma}}_{yy} - \hat{\mathbf{B}}' \hat{\boldsymbol{\Sigma}}_{xx} \hat{\mathbf{B}} \quad (3)$$

(see, e.g., Johnson and Wichern 1998, p. 440). Multivariate regression has applications in chemometrics, engineering, econometrics, psychometrics, and other fields. Recent work on multivariate regression has been done by Barrett and Ling (1992), Breiman and Friedman (1997), Cook and Setodji (2003), Davis and McKean (1993), Gleser (1992), Koenker and Portnoy (1990), Ollila, Hettmansperger, and Oja (2002), and Ollila, Oja, and Koivunen (2003).

It is well known that classical multiple regression is extremely sensitive to outliers in the data. This problem also holds in the case of multivariate regression, as can be seen from the following example.

Example 1. We consider a dataset (Lee 1992) that contains measurements of properties of pulp fibers and the paper made from them. The aim is to investigate relations between pulp fiber properties and the resulting paper properties. The dataset contains $n = 62$ measurements of the following four pulp fiber characteristics: arithmetic fiber length, long fiber fraction, fine fiber fraction, and zero span tensile. The four paper properties that have been measured are breaking length, elastic modulus, stress at failure, and burst strength. The dataset is available at <http://allserv.ugent.be/~svaelst/data/pulpfiber.txt>.

Our goal is to predict the $q = 4$ paper properties from the $p = 4$ fiber characteristics. For this purpose, we first applied classical multivariate regression to the data.

Figure 1 represents the result of the classical analysis. It plots the Mahalanobis distances of the residuals $\mathbf{r}_i = \mathbf{y}_i - \hat{\mathbf{B}}'\mathbf{x}_i - \hat{\boldsymbol{\alpha}}$ as given by

$$d(\mathbf{r}_i) := \sqrt{\mathbf{r}_i' (\hat{\boldsymbol{\Sigma}}_{\varepsilon})^{-1} \mathbf{r}_i}$$

versus the Mahalanobis distances of the carriers,

$$d(\mathbf{x}_i) := \sqrt{(\mathbf{x}_i - \boldsymbol{\mu}_x)' (\hat{\boldsymbol{\Sigma}}_{xx})^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_x)}.$$

This diagnostic plot combines the information on regression outliers and leverage points and is much more useful than either distance separately. The horizontal and vertical lines are

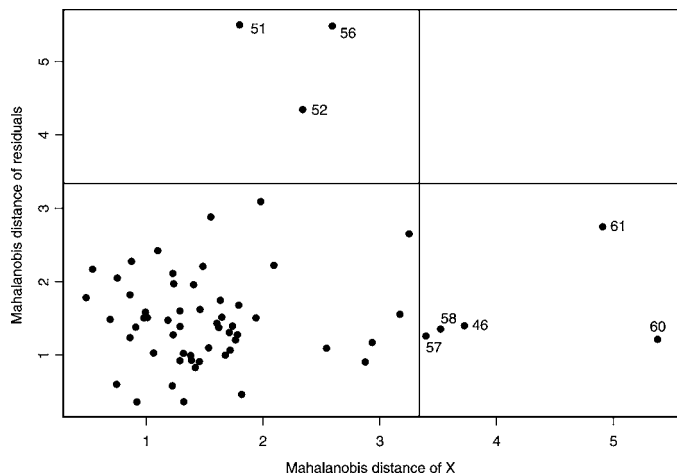


Figure 1. Plot of Mahalanobis Distances of Least Squares Residuals versus Mahalanobis Distances of the Carriers for the Pulp Fiber Data.

the usual cutoff values, $\sqrt{\chi_{p, .975}^2}$ and $\sqrt{\chi_{q, .975}^2}$, both of which equal 3.34 because $p = q = 4$ in this example. From this plot, we see that observations 51, 52, and 56 are detected as vertical outliers. On the other hand, some observations are identified as leverage points (observations 60 and 61 are the largest), but they are not considered regression outliers because they have small residual distance.

To check the result obtained by classical multivariate regression, we start by applying univariate robust regression with the same regressors but for each of the responses separately. Here we use the least trimmed squares (LTS) estimator of Rousseeuw (1984), which can be computed quickly with the FAST-LTS algorithm of Rousseeuw and Van Driessen (2000). To obtain reliable outlier identification, we use the reweighted LTS with finite-sample correction factor as proposed by Pison, Van Aelst, and Willems (2002).

Figure 2 shows the standardized residuals resulting from LTS regression with the first response (breaking length). From this plot, we immediately see that observations 51, 52, 56, and 61 are detected as outliers. Similarly, outliers can be identified from standardized LTS residuals corresponding to the other three responses. Table 1 summarizes the outliers detected by applying LTS for each of the four responses. From Table 1, we see

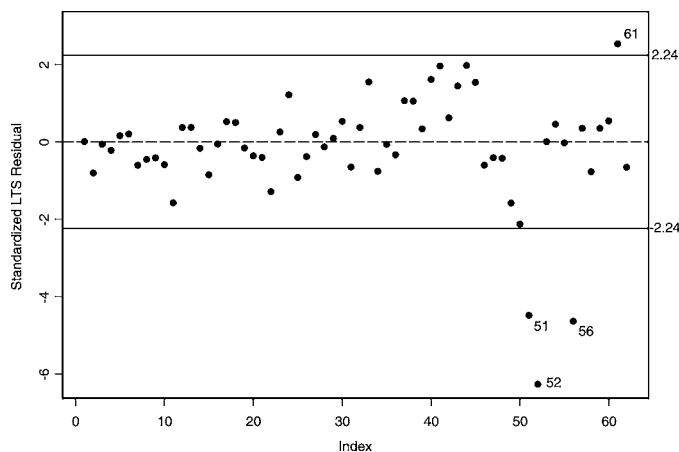


Figure 2. Plot of the Standardized LTS Residuals Corresponding to the First Response (breaking length) versus the Case Number.

Table 1. Observations in the Pulp Fiber Data That Are Detected as Outliers by Applying LTS Regression to Each of the Four Responses Separately

Response	Outliers
y_1	51, 52, 56, 61
y_2	61
y_3	52, 56, 61
y_4	51, 52, 56

that the univariate LTS regressions identify observations 51, 52, 56, and 61 as outliers. This already shows that the classical multivariate regression based on least squares in Figure 1 has been influenced by outliers, because it did not detect observation 61 as a regression outlier. Hence, clearly the least squares multivariate regression has been influenced by this leverage point. This analysis shows that we need robust estimators to investigate these data. However, applying univariate LTS regressions to each of the response variables separately does not yield a solution that is equivariant under affine transformations of the response variables. Moreover, this approach allows us to detect outliers only in the coordinate directions of the responses, not outliers that are masked in these directions. Therefore, we aim to construct a robust method for multivariate regression that allows us to detect all of the outliers and is also reasonably efficient in both the statistical and computational senses. After developing such a robust method, we further analyze these data in Section 5.

In the next section we introduce a robust method for multivariate regression based on the minimum covariance determinant (MCD) estimate of the joint (\mathbf{x}, \mathbf{y}) variables. We study the performance of the estimator by simulations. In Section 3 we investigate several reweighted versions of the estimator that improve the performance of the initial estimator and select the reweighting scheme that works best. We study the finite-sample robustness of the optimal estimator in Section 4. In Section 5 we continue the analysis of the previous example and describe an application to chemical engineering. In Section 6 we show that the robust estimator has the equivariance properties that we expect from a multivariate regression method. In Section 7 we discuss the robustness properties of the estimator and derive studentized residual distances. We summarize our conclusions in Section 8 and give all proofs in the Appendix.

2. MINIMUM COVARIANCE DETERMINANT REGRESSION

We propose using robust estimators for the center μ and scatter matrix Σ in expressions (1)–(3) to construct a robust multivariate regression method that has the equivariance properties required for a multivariate regression estimator. Many robust estimators of multivariate location and scatter have been investigated in the literature, including M estimators (Maronna 1976), the minimum volume ellipsoid and MCD estimator (Rousseeuw 1984, 1985), S estimators (Davies 1987; Rousseeuw and Leroy 1987; Lopuhaä 1989), CM estimators (Kent and Tyler 1996), and τ estimators (Lopuhaä 1991). More recently, depth-based location and scatter estimators were introduced (Zuo, Cui, and He 2001; Zuo and Cui 2002). Robust

estimators of location and scatter in high dimensions have been investigated by Woodruff and Rocke (1994), Rocke (1996), and Rocke and Woodruff (1996). In the multiple regression case, Maronna and Morgenthaler (1986) used multivariate M estimators in (1)–(3), but their method inherits the low breakdown value of M estimators. A multivariate regression method of the M type was proposed by Koenker and Portnoy (1990), who noted that their method lacks affine equivariance.

We use the MCD to estimate the center and scatter matrix of the joint (\mathbf{x}, \mathbf{y}) variables because the MCD is a robust estimator with high breakdown value and bounded influence function (Croux and Haesbroeck 1999). Moreover, the MCD estimator is asymptotically normal (Butler, Davies, and Jhun 1993). We call the resulting robust multivariate regression method *MCD regression*.

Consider a dataset $\mathbf{Z}_n = \{\mathbf{z}_i; i = 1, \dots, n\} \in \mathbb{R}^{p+q}$. The MCD looks for the subset $\{\mathbf{z}_{i_1}, \dots, \mathbf{z}_{i_h}\}$ of size h whose covariance matrix has the smallest determinant, where $[n/2] \leq h \leq n$. We denote $\gamma = (n - h)/n$, so $0 \leq \gamma \leq .5$. The estimate for the center is then defined as the mean $\mathbf{t}_n = \frac{1}{h} \sum_{j=1}^h \mathbf{z}_{i_j}$ and the covariance estimate is given by $\mathbf{C}_n = c_n c_\gamma \frac{1}{h} \sum_{j=1}^h (\mathbf{z}_{i_j} - \mathbf{t}_n) \times (\mathbf{z}_{i_j} - \mathbf{t}_n)^t$, where c_γ is a consistency factor and c_n is a small-sample correction factor (see Pison et al. 2002). The MCD estimator has breakdown value approximately equal to γ . Two common choices for h are $h = [(n + p + q + 1)/2] \approx n/2$ so $\gamma \approx .5$, which yields the highest possible breakdown value, and $h \approx 3n/4$ (i.e., $\gamma \approx .25$), which gives a better compromise between efficiency and breakdown. Recently, Rousseeuw and Van Driessen (1999) constructed a fast algorithm for computing the MCD. This FAST-MCD algorithm made the MCD very useful for analyzing large datasets, for example, with n in the hundred thousands. Other robust methods for analyzing large datasets have been developed and used by Knorr, Ng, and Zamar (2001), Alqallaf, Konis, Martin, and Zamar (2002), and Maronna and Zamar (2002).

Because computation of the MCD regression estimates involves computation of the MCD of the joint (\mathbf{x}, \mathbf{y}) variables followed by standard matrix operations, we obtain a computationally efficient method. Moreover, from (1)–(3), we immediately see that regressions of all possible splits in \mathbf{x} and \mathbf{y} variables can be carried out once the MCD of the joint (\mathbf{x}, \mathbf{y}) variables has been computed. It has been shown that observations that lie far from the center can have only a small effect on the MCD estimates. Therefore, both leverage points (which have a large \mathbf{x} distance) and regression outliers (which are deviating in \mathbf{y} space) can have only a small effect on the MCD regression estimates. However, it has been noted that the MCD can have a low efficiency (Croux and Haesbroeck 1999).

To investigate the efficiency of the MCD regression, we performed the following simulation study. For various sample sizes n and for different choices of p and q , we generated m datasets of size n from the multivariate standard Gaussian distribution $N(\mathbf{0}, \mathbf{I}_{p+q})$, which corresponds to putting $\mathbf{B} = \mathbf{0}$ and $\boldsymbol{\alpha} = \mathbf{0}$. For each dataset $\mathbf{Z}^{(l)}$, $l = 1, \dots, m$, we carried out MCD regression, yielding the $(p \times q)$ slope matrix estimate $\hat{\mathbf{B}}^{(l)}$, the intercept vector $\hat{\boldsymbol{\alpha}}^{(l)} \in \mathbb{R}^q$, and the $(q \times q)$ covariance matrix estimate $\hat{\boldsymbol{\Sigma}}_\epsilon^{(l)}$ of the errors.

The Monte Carlo variance of a slope coefficient $\hat{\mathbf{B}}_{jk}$ is measured as

$$\text{var}(\hat{\mathbf{B}}_{jk}) = n \text{var}_l(\hat{\mathbf{B}}_{jk}^{(l)}) \quad \text{for } j = 1, \dots, p \text{ and } k = 1, \dots, q. \quad (4)$$

The overall variance of the estimated matrix $\hat{\mathbf{B}}$ is defined as $\text{var}(\hat{\mathbf{B}}) = \text{ave}_{j,k}(\text{var} \hat{\mathbf{B}}_{jk})$. The corresponding finite-sample efficiency of the slope is then given by $1/\text{var}(\hat{\mathbf{B}})$. Analogously, we compute the finite-sample efficiency of the intercept vector. To measure the accuracy of the error scatter matrix, we use the standardized variance (Bickel and Lehmann 1976) of the elements of the error covariance matrix, defined as

$$\text{Stvar}((\hat{\boldsymbol{\Sigma}}_\epsilon)_{jk}) = \frac{n \text{var}_l((\hat{\boldsymbol{\Sigma}}_\epsilon^{(l)})_{jk})}{[\text{ave}_l \text{ave}_j((\hat{\boldsymbol{\Sigma}}_\epsilon^{(l)})_{jj})]^2} \quad \text{for } j = 1, \dots, q \text{ and } k = 1, \dots, q. \quad (5)$$

The overall finite-sample efficiency of the off-diagonal elements is then given by $1/\text{ave}_{j \neq k}(\text{Stvar}((\hat{\boldsymbol{\Sigma}}_\epsilon)_{jk}))$. For the diagonal elements, the finite-sample efficiency is given by $2/\text{ave}_j(\text{Stvar}((\hat{\boldsymbol{\Sigma}}_\epsilon)_{jj}))$, because the Fisher information equals 2 in this case.

The top part of Table 2 shows the simulation results for $p = 4$ and $q = 4$, but the results were similar for many other choices of p and q . The table contains sample sizes between 50 and 500. All simulations were done with $m = 1,000$ replications. Given are the finite-sample efficiencies of $\hat{\mathbf{B}}$, $\hat{\boldsymbol{\alpha}}$, the diagonal elements of $\hat{\boldsymbol{\Sigma}}_\epsilon$, and the off-diagonal elements of $\hat{\boldsymbol{\Sigma}}_\epsilon$. We see that the finite-sample efficiencies are very low for $\gamma = .5$ and are somewhat better for $\gamma = .25$. In the next section we propose using reweighted estimators to improve these efficiencies.

3. REWEIGHTED MULTIVARIATE REGRESSION

To increase the efficiencies obtained in the previous section, we now consider reweighted versions of our estimator. These reweighted estimators inherit the robustness properties of the initial estimator while attaining a higher efficiency. We consider three versions, one based on reweighting the location estimator, one based on reweighting the regression estimator, and one based on reweighting both.

3.1 Reweighting the Location Estimator

To increase the efficiency of the location and scatter estimators, it is customary to compute one-step reweighted versions (Rousseeuw and Leroy 1987; Lopuhaä 1999; Zuo et al. 2001; Zuo and Cui 2002). The one-step reweighted MCD estimates with nominal trimming portion δ_l are defined as

$$\begin{aligned} \mathbf{t}_n^1 &= \frac{\sum_{i=1}^n w(d^2(\mathbf{z}_i)) \mathbf{z}_i}{\sum_{i=1}^n w(d^2(\mathbf{z}_i))} \quad \text{and} \\ \mathbf{C}_n^1 &= d_{\delta_l} \frac{\sum_{i=1}^n w(d^2(\mathbf{z}_i)) (\mathbf{z}_i - \mathbf{t}_n^1) (\mathbf{z}_i - \mathbf{t}_n^1)^t}{\sum_{i=1}^n w(d^2(\mathbf{z}_i))}, \end{aligned} \quad (6)$$

where d_{δ_l} is a consistency factor. The weights are computed as $w(d^2(\mathbf{z}_i)) = I(d^2(\mathbf{z}_i) \leq q_{\delta_l})$, where $q_{\delta_l} = \chi_{p+q, 1-\delta_l}^2$ and $d(\mathbf{z}_i) = ((\mathbf{z}_i - \mathbf{t}_n)^t \mathbf{C}_n^{-1} (\mathbf{z}_i - \mathbf{t}_n))^{1/2}$ is the robust distance of observation \mathbf{z}_i based on the initial MCD estimates $(\mathbf{t}_n, \mathbf{C}_n)$. It is

Table 2. Finite-Sample Efficiencies of the Slope Matrix, Intercept Vector, and Error Covariance Matrix of the Four Types of MCD Regression, for $p = 4$ and $q = 4$

		n				
γ		50	100	300	500	∞
MCD regression: $\hat{\mathbf{B}}, \hat{\alpha}, \hat{\Sigma}_e$						
.50	Slope	.176	.167	.166	.169	.166
	Intercept	.268	.290	.300	.298	.307
	Σ_{diag}	.211	.205	.196	.190	.182
	Σ_{offdiag}	.194	.183	.166	.172	.166
.25	Slope	.371	.387	.401	.410	.403
	Intercept	.506	.545	.568	.543	.578
	Σ_{diag}	.401	.431	.439	.432	.430
	Σ_{offdiag}	.387	.415	.393	.401	.403
MCD regression with reweighted location: $\hat{\mathbf{B}}^L, \hat{\alpha}^L, \hat{\Sigma}_e^L$						
.50	Slope	.200	.354	.662	.762	.851
	Intercept	.303	.525	.811	.838	.934
	Σ_{diag}	.245	.391	.677	.727	.794
	Σ_{offdiag}	.222	.384	.664	.735	.851
.25	Slope	.403	.598	.772	.830	.864
	Intercept	.545	.747	.883	.877	.936
	Σ_{diag}	.434	.613	.793	.798	.812
	Σ_{offdiag}	.427	.629	.782	.813	.864
MCD regression with reweighted regression: $\hat{\mathbf{B}}^R, \hat{\alpha}^R, \hat{\Sigma}_e^R$						
.50	Slope	.245	.465	.812	.902	.957
	Intercept	.338	.582	.862	.875	.959
	Σ_{diag}	.251	.387	.685	.735	.858
	Σ_{offdiag}	.232	.399	.684	.763	.880
.25	Slope	.538	.758	.895	.948	.960
	Intercept	.622	.804	.927	.906	.961
	Σ_{diag}	.463	.627	.820	.815	.874
	Σ_{offdiag}	.462	.665	.812	.841	.892
MCD regression with reweighted location and regression: $\hat{\mathbf{B}}^{LR}, \hat{\alpha}^{LR}, \hat{\Sigma}_e^{LR}$						
.50	Slope	.233	.628	.906	.955	.961
	Intercept	.332	.721	.928	.920	.962
	Σ_{diag}	.252	.501	.826	.829	.881
	Σ_{offdiag}	.233	.542	.824	.860	.900
.25	Slope	.508	.801	.913	.959	.961
	Intercept	.614	.849	.942	.924	.962
	Σ_{diag}	.458	.680	.864	.839	.881
	Σ_{offdiag}	.459	.728	.854	.872	.900

NOTE: The number of replications was $m = 1,000$.

customary to take $\delta_l = .025$ (Rousseeuw and Van Driessen 1999). The robustness properties of the one-step reweighted MCD estimators are similar to those of the initial MCD (Lopuhaä and Rousseeuw 1991; Lopuhaä 1999). Other methods for increasing the efficiency of the MCD location and scatter include one-step M estimators and cross-checking (He and Wang 1996).

We can now compute the multivariate regression estimates (1), (2), and (3) based on the reweighted location and scatter $(\mathbf{t}_n^L, \mathbf{C}_n^L)$. We denote the resulting regression by $\hat{\mathbf{B}}^L, \hat{\alpha}^L$, and $\hat{\Sigma}_e^L$, where “L” indicates that the reweighting was done in the location stage. The simulation results for the reweighted location estimators are given in the second part of Table 2. We see that multivariate regression estimates based on the reweighted MCD have a much higher efficiency than those based on the initial unweighted MCD.

3.2 Reweighting the Regression

In a regression analysis, it is natural to use weights based on the residuals corresponding to the initial fit (Rousseeuw

and Leroy 1987). Denote the residual of the observation \mathbf{z}_i by $\mathbf{r}_i = \mathbf{y}_i - \hat{\mathbf{B}}^L \mathbf{x}_i - \hat{\alpha}^L$. We now define the reweighted regression estimators

$$\mathbf{T}_n^R = \left(\sum_{i=1}^n w(d^2(\mathbf{r}_i)) \mathbf{u}_i \mathbf{u}_i^t \right)^{-1} \sum_{i=1}^n w(d^2(\mathbf{r}_i)) \mathbf{y}_i \mathbf{u}_i \quad (7)$$

and

$$\hat{\Sigma}_e^R = d_{\delta_r} \frac{\sum_{i=1}^n w(d^2(\mathbf{r}_i)) (\mathbf{r}_i^R) (\mathbf{r}_i^R)^t}{\sum_{i=1}^n w(d^2(\mathbf{r}_i))}, \quad (8)$$

where $\mathbf{T}_n^R = ((\hat{\mathbf{B}}^R)^t, \hat{\alpha}^R)^t$, $\mathbf{u}_i = (\mathbf{x}_i^t, 1)^t$, $(\mathbf{r}_i^R)_i = \mathbf{y}_i - (\hat{\mathbf{B}}^R)^t \mathbf{x}_i - \hat{\alpha}^R$, δ_r is the trimming portion, and d_{δ_r} is a consistency factor. Following Rousseeuw and Leroy (1987), we take $\delta_r = .01$ as our default. The superscript “R” denotes that the weights were based on the initial regression. In particular, the weights are computed as $w(d^2(\mathbf{r}_i)) = I(d^2(\mathbf{r}_i) \leq q_{\delta_r})$, where $q_{\delta_r} = \chi_{q, 1-\delta_r}^2$ and $d(\mathbf{r}_i) = (\mathbf{r}_i^t (\hat{\Sigma}_e^L)^{-1} \mathbf{r}_i)^{1/2}$ is the robust distance of the i th residual. The robustness properties of these reweighted regression estimators follow from the properties of

the initial regression estimators. Note that the weights now depend only on the size of the residual distance $w(d^2(\mathbf{r}_i))$ so in contrast with the initial estimates, good leverage points (which have large distance in x -space but small residual distance and thus are not outliers for the regression model) are no longer downweighted.

The third part of Table 2 gives the simulation results for the reweighted regression estimators. We see that the reweighted multivariate regression estimates have a much higher efficiency than the initial estimates based on MCD. Moreover, the efficiency of the reweighted regression estimates is also higher than the efficiency of the estimates based on the reweighted MCD.

3.3 Reweighting Both Location and Regression

A further possibility is to use the robust distances $d(\mathbf{r}_i^L) = ((\mathbf{r}_i^L)^t (\hat{\Sigma}_\epsilon^L)^{-1} \mathbf{r}_i^L)^{1/2}$ in (7) and (8), where $\mathbf{r}_i^L = \mathbf{y}_i - (\hat{\mathbf{B}}^L)^t \times \mathbf{x}_i - \hat{\alpha}^L$. This yields a weighted regression estimator with weights based on the residuals of the method in Section 3.1. We denote the resulting estimators by $\mathbf{T}_n^{LR} = ((\hat{\mathbf{B}}^{LR})^t, \hat{\alpha}^{LR})^t$ and $\hat{\Sigma}_\epsilon^{LR}$. Also here good leverage points are no longer downweighted. The simulation results for the reweighted location estimators are given in the last part of Table 2. From this table, we see that the efficiency of location regression (LR) weighting is comparable for small samples ($n = 50$) and clearly better for larger samples than the efficiency of the other reweighting schemes. Overall, we also see that $\gamma = .25$ consistently outperformed $\gamma = .50$, with the difference being larger for small samples. Hence, from an efficiency standpoint, LR-weighted MCD regression with $\gamma = .25$ comes out best. We demonstrate in Section 6 that the breakdown value of MCD regression is approximately equal to γ , so Table 2 shows that there is a trade-off between efficiency and breakdown. In practice, data with more than 20% of outliers rarely occur, so we recommend using the LR-weighted MCD regression with $\gamma = .25$ as the default to obtain a better efficiency. If the data are of very low quality such that a higher level of outliers can be expected, using LR-weighted MCD regression with $\gamma = .50$ is more appropriate.

4. FINITE-SAMPLE ROBUSTNESS

To study the finite-sample robustness, we carried out simulations with datasets contaminated by different types of outliers. A point $(\mathbf{x}_i, \mathbf{y}_i)$ that does not follow the linear pattern of the majority of the data but whose \mathbf{x}_i is not outlying is called a *vertical outlier*. A point $(\mathbf{x}_i, \mathbf{y}_i)$ whose \mathbf{x}_i is outlying is called a *leverage point*. We say that such an $(\mathbf{x}_i, \mathbf{y}_i)$ is a bad leverage point when it does not follow the pattern of the majority; otherwise, it is a good leverage point (which does not harm the fit).

Because regression estimators often break down in the presence of vertical outliers or bad leverage points, we generated datasets with both types of outliers. For sample sizes between $n = 50$ and $n = 500$ and with $p = 4$ and $q = 4$, we generated $m = 1,000$ datasets from the multivariate standard Gaussian distribution $N(\mathbf{0}, \mathbf{I}_{p+q})$. (This is the same situation as described in Sec. 2.) We then replaced 10% of the data as follows. The \mathbf{x}_i are kept, but the q response variables are distributed as $N(2\sqrt{\chi_{p+q,.99}^2}, .1)$. This yields vertical outliers, because only

their responses are outlying. We also replaced 10% of the data with bad leverage points for which the p independent variables are generated according to $N(2\sqrt{\chi_{p,.99}^2}, .1)$ and the q dependent variables are generated from $N(2\sqrt{\chi_{q,.99}^2}, .1)$.

As in the previous simulations, for each dataset $\mathbf{Z}^{(l)}$, $l = 1, \dots, m$, we computed the $(p \times q)$ slope matrix $\hat{\mathbf{B}}^{(l)}$, the intercept vector $\hat{\alpha}^{(l)} \in \mathbb{R}^q$, and the $(q \times q)$ covariance matrix $\hat{\Sigma}_\epsilon^{(l)}$ of the errors. To measure robustness, we used the bias and the mean squared error (MSE). As commonly defined, the bias and MSE of a univariate component T are given by

$$\text{bias}(T) = \text{ave}_l (T^{(l)} - \theta)$$

and

$$\text{MSE}(T) = n \text{ave}_l (T^{(l)} - \theta)^2,$$

with θ the true value of the parameter. The bias and MSE of the slope are defined as

$$\text{bias}(\hat{\mathbf{B}}) = \sqrt{\text{ave}_{j,k} (\text{bias}(\hat{\mathbf{B}}_{jk})^2)}$$

and

$$\text{MSE}(\hat{\mathbf{B}}) = \text{ave}_{j,k} (\text{MSE}(\hat{\mathbf{B}}_{jk}))$$

and similarly for the intercept $\hat{\alpha}$ and for the diagonal and off-diagonal elements of $\hat{\Sigma}_\epsilon$.

Table 3 gives the simulation results when the estimates of the slope matrix, intercept vector, and error covariance matrix were obtained from the LR-weighted method with $\gamma = .25$ and from the classical multivariate least squares regression. Simulations for other sample sizes n and different dimensions p and q gave similar results. In Table 3 we see that in the presence of vertical outliers and bad leverage points, both the bias and MSE obtained from the LR-weighted MCD regression are much lower than those obtained from least squares regression. The low bias and MSE values of the LR-weighted method are in line with the asymptotic robustness properties in Section 6.

To compare the MCD regression with the univariate robust regressions approach used in Section 1, we used the foregoing simulation setup but we now generated correlated multivariate Gaussian responses with correlation $r_{jk} = .5$, $j \neq k$. Thus we obtained a regression model with correlated errors. We generated 10% of vertical outliers and 10% of bad leverage points as before.

The results for the LR-weighted MCD regression ($\gamma = .25$) in Table 4 are comparable with the results in Table 3, as expected from the equivariance of the estimator. Table 4 shows that the LR-weighted MCD regression in general outperforms the coordinatewise LTS regressions both in bias and MSE. The differences are largest for the slope estimates. Note that Table 4 does not contain results for the off-diagonal elements of the error covariance matrix, because these elements are not estimated in the univariate LTS approach. Hence, another advantage of the multivariate MCD regression method is that it gives a robust estimate of the full error covariance matrix.

Based on the performance results in the previous section and the robustness results here, we recommend using the LR-weighted method with $\gamma = .25$ in practice to identify all outliers and robustly estimate the full error covariance matrix.

Table 3. Bias and MSE of the Slope Matrix, Intercept Vector, and Error Covariance Matrix Obtained by the LR-Weighted MCD Regression With $\gamma = .25$ and Multivariate Least Squares Regression

	n					
	50		100		500	
	Bias	MSE	Bias	MSE	Bias	MSE
LR-weighted MCD regression ($\gamma = .25$)						
Slope	.0066	1.637	.0038	1.462	.0013	1.307
Intercept	.0104	1.501	.0036	1.415	.0021	1.336
Σ_{diag}	.1349	3.326	.0704	3.240	.0104	2.845
Σ_{offdiag}	.0050	1.245	.0049	1.319	.0021	1.369
Least squares regression						
Slope	.2068	10.883	.2071	12.233	.2064	28.806
Intercept	1.0225	54.275	1.0214	105.876	1.0243	525.924
Σ_{diag}	6.5387	2,156.323	6.8122	4,655.881	7.0394	24,788.392
Σ_{offdiag}	6.8076	2,332.350	7.0464	4,977.298	7.2440	26,246.846

NOTE: The data contained 20% of outliers. The number of replications was $m = 1,000$.

5. EXAMPLES

Example 1 (Continued). We now continue the analysis of Example 1 in Section 1 by applying the LR-weighted robust multivariate regression method with $\gamma = .25$ to these data. Figure 3 shows the diagnostic plot corresponding to the robust analysis. This plot is a generalization of the diagnostic plot for multiple regression due to Rousseeuw and van Zomeren (1990). In this display the robust distances of the q -dimensional residuals $d(\mathbf{r}_i^{LR}) = ((\mathbf{r}_i^{LR})^t (\hat{\Sigma}_e^{LR})^{-1} \mathbf{r}_i^{LR})^{1/2}$ are plotted versus the robust distances of the p -dimensional \mathbf{x}_i given by $d(\mathbf{x}_i) = ((\mathbf{x}_i - (\mathbf{t}_n^1)_x)^t ((\mathbf{C}_n^1)_{xx})^{-1} (\mathbf{x}_i - (\mathbf{t}_n^1)_x))^{1/2}$. The plot enables us to classify the data points into regular observations, vertical outliers, good leverage points, and bad leverage points. Moreover, it allows us to see whether a point is an extreme outlier or merely a borderline case. Being a graphical tool, this plot also allows us to discover unexpected structure in the data. Note that all of the estimates on which the plot is based are byproducts of the robust multivariate regression algorithm, so the plot requires very little computation time.

From Figure 3, we see that 13 observations have residuals with robust distance above the horizontal cutoff line at $\sqrt{\chi_{4, .975}^2} = 3.34$ and thus are detected as regression outliers. Eight of these points also have a large x -distance and thus are bad leverage points. Note that classical multivariate regression detected only three of these outliers (51, 52, and 56) and considered four of the outliers (46, 58, 60, and 61) to be good

leverage points. Moreover, by applying LTS for each of the responses separately, we detected only one additional outlier (61) but nine other outliers remained hidden, among which the bad leverage points 59, 60, and 62 are the most severe.

By exploring the origin of the collected data, we found out that all but the last four pulp samples (observations 59–62) were produced from fir wood. Moreover, most of the outlying samples were obtained using different pulping processes. For example, observation 62 is the only sample from a chemi-thermomechanical pulping process, observations 60 and 61 are the only samples from a solvent pulping process, and observations 51, 52, and 56 are obtained from a kraft pulping process. Finally, the smaller outliers (22, 46–48, and 58) all were Douglas fir samples.

Example 2. This example describes an actual dataset obtained from Shell's polymer laboratory in Ottignies, Belgium (courtesy of Dr. Christian Ritter). For reasons of confidentiality, all variables have been standardized, and their exact meanings are not given. The dataset comprises of $n = 217$ observations with $p = 4$ predictor variables and $q = 3$ response variables. The predictor variables describe the chemical characteristics of a piece of foam, and the response variables measure its physical

Table 4. Bias and MSE of the Slope Matrix, Intercept Vector and Error Variances Obtained by the LR-Weighted MCD Regression With $\gamma = .25$ and Univariate LTS Regressions

	n					
	50		100		500	
	Bias	MSE	Bias	MSE	Bias	MSE
LR-weighted MCD regression ($\gamma = .25$)						
Slope	.0044	1.644	.0043	1.483	.0017	1.319
Intercept	.0069	1.512	.0017	1.364	.0021	1.313
Σ_{diag}	.1316	3.357	.0676	3.200	.0097	2.900
Combination of univariate LTS regressions						
Slope	.2360	4.630	.2403	7.270	.2406	30.214
Intercept	.0264	2.144	.0342	1.880	.0331	2.272
Σ_{diag}	.0976	4.166	.0868	6.158	.1004	9.044

NOTE: The data contained 20% of outliers. The number of replications was $m = 1,000$.

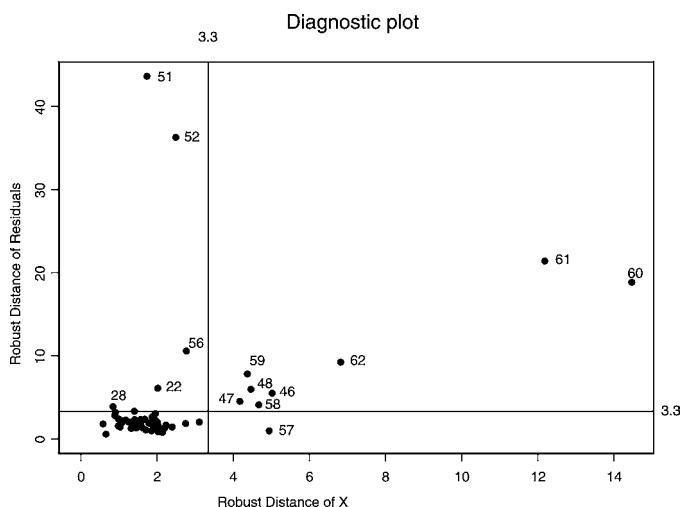


Figure 3. Plot of Robust Distances of Residuals versus Robust Distances of the Carriers for the Pulp Fiber Data.

properties, such as tensile strength. Foam product specifications are expressed in terms of physical properties. Production units around the world have to meet the prescribed physical requirements. The physical properties are determined by the chemical composition used in the production process. However, different chemical compositions will lead to foams that meet all required specifications. Moreover, depending on the location of the production unit, there is a strong variation in the price of the necessary chemicals. Therefore, the goal is to establish a relationship between the chemical inputs and the resulting physical properties, which then can be used to determine the cheapest chemical composition resulting in foams that meet all physical requirements. We used multivariate regression to determine the relationship between the chemical inputs and the physical properties. A few cases with missing values had been omitted in advance. After an initial exploratory study of the seven variables, including their Q-Q plots, we applied Box-Cox transformations to them. We then ran a robust multivariate regression using the LR-weighted method with $\gamma = .25$. This computation took only 43 seconds on a Sun SparcStation 20/514.

Figure 4 shows the diagnostic plot of the Shell foam data (robust distances of the residuals r_i^{LR} versus the robust distances of the \mathbf{x}_i). Observations 215 and 110 lie far from both the horizontal cutoff line at $\sqrt{\chi_{3,.975}^2} = 3.06$ and the vertical cutoff line at $\sqrt{\chi_{4,.975}^2} = 3.34$. These two observations can thus be classified as bad leverage points. Several observations lie substantially above the horizontal cutoff but not to the right of the vertical cutoff, which means that they are vertical outliers (i.e., their residuals are outlying but their x -values are not).

When this list of special points was presented to the scientists who had made the measurements, we learned that eight observations in Figure 4 were produced with a different production technique and hence belong to a different population with other characteristics. These include the observations 210, 212, and 215. We therefore remove these eight observations from the data and retain only observations from the intended population.

Running the method again yielded the diagnostic plot shown in Figure 5. Observation 110 is still a bad leverage point, and several of the vertical outliers also remain. No chemical/physical mechanism was found to explain why these points

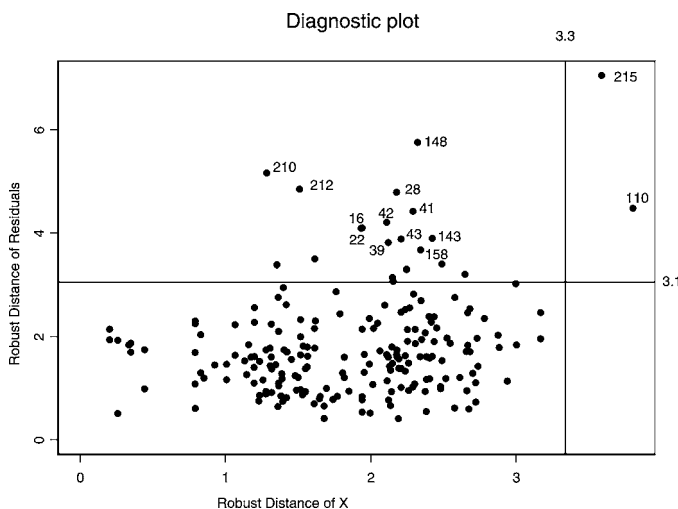


Figure 4. Diagnostic Plot of Robust Residuals versus Robust Distances of the Carriers for the Foam Data.

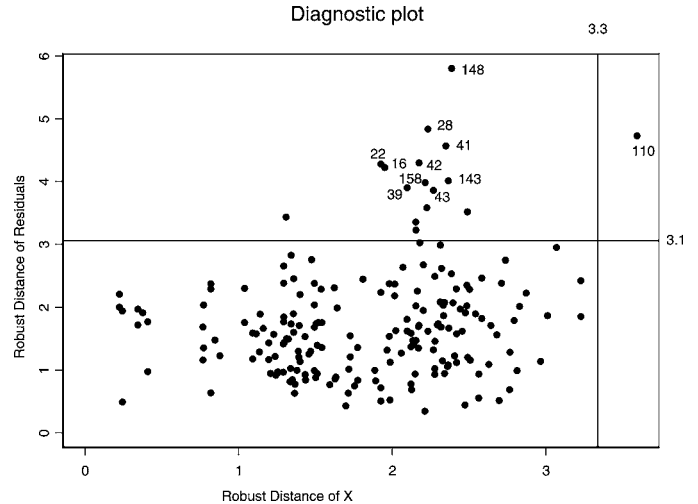


Figure 5. Diagnostic Plot of Robust Residuals versus Robust Distances for the Corrected Foam Data.

are outliers, leaving open the possibility of some large measurement errors. But the detection of these substantial outliers at least gives us the option to choose whether or not to allow them to affect the final result.

6. EQUIVARIANCE AND ROBUSTNESS PROPERTIES

The theorems in this section demonstrate that the proposed LR-weighted method based on MCD has the natural equivariance properties of multivariate regression estimators and is robust. They generalize the regression, scale, and affine equivariance (see Rousseeuw and Leroy 1987, p. 116) and robustness of multiple regression estimators. All proofs are given in the Appendix.

Denote $\mathbf{T}_n(\mathbf{X}, \mathbf{Y}) = (\hat{\mathbf{B}}^t, \hat{\alpha})^t$, where the matrix \mathbf{X} is $(n \times p)$ and \mathbf{Y} is $(n \times q)$. The estimator \mathbf{T}_n is called *regression equivariant* if

$$\mathbf{T}_n(\mathbf{X}, \mathbf{Y} + \mathbf{X}\mathbf{D} + \mathbf{1}_n\mathbf{w}^t) = \mathbf{T}_n(\mathbf{X}, \mathbf{Y}) + (\mathbf{D}^t, \mathbf{w})^t, \quad (9)$$

where \mathbf{D} is any $(p \times q)$ matrix, \mathbf{w} is any $(q \times 1)$ vector, and $\mathbf{1}_n = (1, 1, \dots, 1)^t \in \mathbb{R}^n$. Regression equivariance means that if we add a linear function of the explanatory variables to the responses, then the coefficients of this linear function are also added to the estimator.

The estimator \mathbf{T}_n is said to be *y-affine equivariant* if

$$\mathbf{T}_n(\mathbf{X}, \mathbf{Y}\mathbf{C} + \mathbf{1}_n\mathbf{d}^t) = \mathbf{T}_n(\mathbf{X}, \mathbf{Y})\mathbf{C} + (\mathbf{O}_{pq}^t, \mathbf{d})^t, \quad (10)$$

where \mathbf{C} is any nonsingular $(q \times q)$ matrix, \mathbf{d} is any $(q \times 1)$ vector, and \mathbf{O}_{pq} is the $(p \times q)$ matrix consisting of 0's. If the response variables are transformed linearly, then y-affine equivariance implies that the estimator \mathbf{T} transforms accordingly.

We say that the estimator \mathbf{T}_n is *x-affine equivariant* if

$$\mathbf{T}_n(\mathbf{X}\mathbf{A}^t + \mathbf{1}_n\mathbf{v}^t, \mathbf{Y}) = (\hat{\mathbf{B}}^t\mathbf{A}^{-1}, \hat{\alpha} - \hat{\mathbf{B}}^t\mathbf{A}^{-1}\mathbf{v})^t \quad (11)$$

for any nonsingular $(p \times p)$ matrix \mathbf{A} and any column vector $\mathbf{v} \in \mathbb{R}^{p \times 1}$. If the explanatory variables are transformed linearly, then x-affine equivariance says that the estimator \mathbf{T}_n transforms correctly.

Theorem 1. The LR-weighted multivariate MCD regression estimator $\mathbf{T}_n = ((\hat{\mathbf{B}}^{LR})^t, \hat{\boldsymbol{\alpha}}^{LR})^t$ is regression, y -affine, and x -affine equivariant.

We also study the theoretical robustness properties of the estimator in terms of its breakdown value and its influence function, which also yields its asymptotic variance. These theoretical properties will confirm the finite-sample results obtained in Sections 3.3 and 4.

The finite-sample breakdown value (Donoho and Huber 1983) of a regression estimator \mathbf{T}_n at a dataset $\mathbf{Z}_n = (\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{n \times (p+q)}$ is defined as the smallest fraction of observations of \mathbf{Z}_n that need to be replaced to carry \mathbf{T}_n beyond all bounds. Formally, this is expressed as

$$\varepsilon_n^*(\mathbf{T}_n, \mathbf{Z}_n) = \min_{1 \leq m \leq n} \left\{ \frac{m}{n} : \sup_{\mathbf{Z}'_n} \|\mathbf{T}_n(\mathbf{Z}_n) - \mathbf{T}_n(\mathbf{Z}'_n)\| = \infty \right\}, \quad (12)$$

where the supremum is over all possible collections \mathbf{Z}'_n that differ from \mathbf{Z}_n in at most m points. The following theorem shows that the LR-weighted MCD regression estimator $\mathbf{T}_n = ((\hat{\mathbf{B}}^{LR})^t, \hat{\boldsymbol{\alpha}}^{LR})^t$ inherits the breakdown value of the initial MCD location and scatter estimators applied to the $(p+q)$ -dimensional dataset \mathbf{Z}_n . Note that the breakdown value of a covariance estimator is the smallest fraction of outliers that can make the largest eigenvalue arbitrarily large or the smallest eigenvalue arbitrarily small.

Theorem 2. Let \mathbf{Z}_n be a set of $n \geq p+q+1$ observations and let \mathbf{t}_n^1 and \mathbf{C}_n^1 be the reweighted MCD estimators of location and scatter with $\min\{\varepsilon_n^*(\mathbf{t}_n^1, \mathbf{Z}_n), \varepsilon_n^*(\mathbf{C}_n^1, \mathbf{Z}_n)\} = \lceil n\gamma \rceil/n$, where $\gamma = (n-h)/n \leq (n-(p+q))/(2n)$. Then the multivariate regression estimator $\mathbf{T}_n = ((\hat{\mathbf{B}}^{LR})^t, \hat{\boldsymbol{\alpha}}^{LR})^t$ also satisfies $\varepsilon_n^*(\mathbf{T}_n, \mathbf{Z}_n) = \lceil n\gamma \rceil/n$.

The influence function of an estimator \mathbf{T} at a distribution H measures the effect on \mathbf{T} of an infinitesimal contamination at a single point (Hampel, Ronchetti, Rousseeuw, and Stahel 1986). If we denote the point mass at $\mathbf{z} = (\mathbf{x}^t, \mathbf{y}^t)^t$ by $\Delta_{\mathbf{z}}$ and write $H_\varepsilon = (1-\varepsilon)H + \varepsilon\Delta_{\mathbf{z}}$, then the influence function is given by

$$IF(\mathbf{z}, \mathbf{T}, H) = \lim_{\varepsilon \downarrow 0} \frac{\mathbf{T}(H_\varepsilon) - \mathbf{T}(H)}{\varepsilon} = \frac{\partial}{\partial \varepsilon} \mathbf{T}(H_\varepsilon) \Big|_{\varepsilon=0}. \quad (13)$$

The following theorem gives the influence functions of the LR-weighted MCD regression estimators at the standard Gaussian distribution. The influence function at general Gaussian distributions then follows from the equivariance properties in Section 6.

Theorem 3. The influence functions of $\hat{\mathbf{B}}^{LR}$, $\hat{\boldsymbol{\alpha}}^{LR}$, and $\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR}$ at the standard Gaussian distribution $H = N(\mathbf{0}, \mathbf{I}_{p+q})$ are given by

$$\begin{aligned} IF(\mathbf{z}, \hat{\mathbf{B}}_{jk}^{LR}, H) &= [c_1 I(\|\mathbf{z}\|^2 \leq q_\gamma) + c_2 I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \\ &\quad + c_3 I(\|\mathbf{y}\|^2 \leq q_{\delta_r})] y_j y_k, \\ IF(\mathbf{z}, \hat{\boldsymbol{\alpha}}_j^{LR}, H) &= [c_4 I(\|\mathbf{z}\|^2 \leq q_\gamma) + c_5 I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) + c_3 I(\|\mathbf{y}\|^2 \leq q_{\delta_r})] y_j, \end{aligned}$$

$$\begin{aligned} IF(\mathbf{z}, (\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR})_{jk}, H) &= [c_6 I(\|\mathbf{z}\|^2 \leq q_\gamma) + c_7 I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \\ &\quad + c_3 I(\|\mathbf{y}\|^2 \leq q_{\delta_r})] y_j y_k, \end{aligned}$$

and

$$\begin{aligned} IF(\mathbf{z}, (\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR})_{jj}, H) &= [c_8 I(\|\mathbf{z}\|^2 \leq q_\gamma) + c_9 I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) + c_3 I(\|\mathbf{y}\|^2 \leq q_{\delta_r})] y_j^2 \\ &\quad + g_1(\|\mathbf{z}\|, \|\mathbf{y}\|) I(\|\mathbf{z}\|^2 \leq q_\gamma) + g_2(\|\mathbf{y}\|) I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \\ &\quad + c_{10} I(\|\mathbf{y}\|^2 \leq q_{\delta_r}) + c_{11}. \end{aligned}$$

Here $q_\gamma = \chi_{p+q, 1-\gamma}^2$. The constants c_1, \dots, c_{11} and the functions g_1 and g_2 are given in the Appendix. Note that the influence functions of the slope and intercept become 0 as soon as $\|\mathbf{y}\|$ becomes large, so vertical outliers as well as bad leverage points have no effect on the regression estimates. For the covariance of the errors, the influence on the off-diagonal elements becomes 0, and the influence on the diagonal elements becomes constant for observations with large $\|\mathbf{y}\|$ so the effect of outliers and leverage points is bounded. On the other hand, good leverage points (which have large $\|\mathbf{x}\|$ but small $\|\mathbf{y}\|$ and thus are not outliers for the regression model) are not downweighted.

Figure 6 shows the influence functions of the LR-weighted MCD regression estimators with $\gamma = .25$ at the bivariate Gaussian distribution $H = N_2(\mathbf{0}, \mathbf{I})$ ($p = q = 1$). The influence functions of the slope $\hat{\boldsymbol{\beta}}^{LR} = \hat{\mathbf{B}}^{LR}$ and the intercept $\hat{\boldsymbol{\alpha}}^{LR}$ are shown in Figures 6(a) and 6(b). The influence function of the error scale $(\hat{\sigma}^{LR})^2 = \hat{\boldsymbol{\Sigma}}_\varepsilon^{LR}$ is shown in Figure 6(c).

From the influence function, we can compute the asymptotic variance of the elements of the slope matrix $\hat{\mathbf{B}}^{LR}$ at the standard Gaussian distribution as

$$ASV(\hat{\mathbf{B}}_{ij}^{LR}, H) = E_H[IF(\mathbf{z}, \hat{\mathbf{B}}_{jk}^{LR}, H)^2] \quad (14)$$

(see Hampel et al. 1986), and similarly for $\hat{\boldsymbol{\alpha}}$ and $\hat{\boldsymbol{\Sigma}}_\varepsilon$. It can be easily shown that the asymptotic variances of the slope and intercept elements of the least squares estimator equal 1. For the least squares estimator of the error covariance, it holds that the asymptotic variance equals 1 for the off-diagonal elements and 2 for the diagonal elements. Therefore, the asymptotic relative efficiency (ARE) of the slope $\hat{\mathbf{B}}^{LR}$ relative to the least squares slope $\hat{\mathbf{B}}_{LS}$ is given by

$$ARE(\hat{\mathbf{B}}^{LR}, H) = \frac{1}{ASV(\hat{\mathbf{B}}_{ij}^{LR}, H)} \quad (15)$$

and similarly for the intercept $\hat{\boldsymbol{\alpha}}^{LR}$ and off-diagonal elements of $\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR}$. The ARE of the diagonal elements of $\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR}$ equals

$$ARE((\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR})_{jj}, H) = \frac{2}{ASV((\hat{\boldsymbol{\Sigma}}_\varepsilon^{LR})_{jj})}. \quad (16)$$

For $p = 4$ and $q = 4$, the ARE of slope, intercept, and diagonal and off-diagonal elements of the error covariance are given in Table 2 under $n = \infty$. For the initial MCD regression and the L- and R-weighted methods, the efficiencies can be obtained from additional results in the Appendix. It is reassuring to note that the finite-sample efficiencies correspond quite well to the

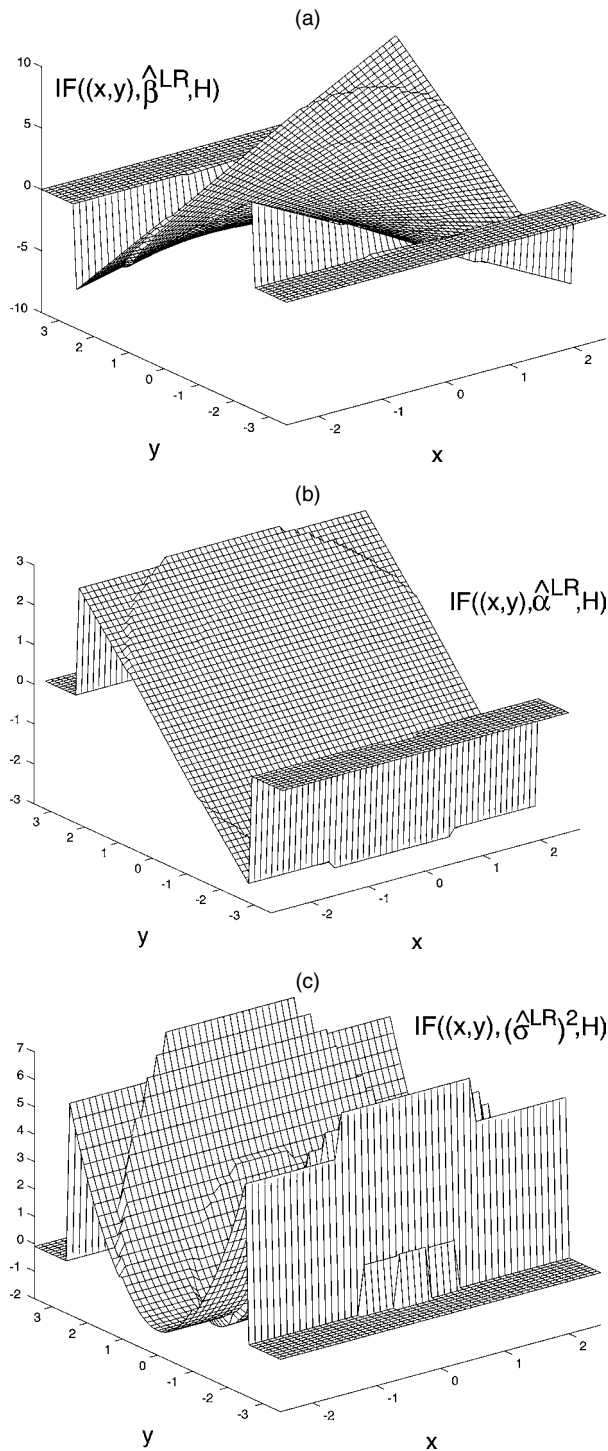


Figure 6. Influence Functions at the Bivariate Gaussian Distribution of (a) Slope, (b) Intercept, and (c) Error Scale of LR-Weighted MCD Regression.

asymptotic efficiencies. The difference is often negligible already for $n = 500$.

To obtain outlier diagnostics that take into account the residual error and the location of the observation in x -space, we now introduce studentized robust residual distances. These studentized residual distances generalize the studentized residuals for univariate robust regression (McKean, Sheather, and Hettmansperger 1990, 1993) to multivariate regression. They also extend the studentized residual distances for multivariate

least squares regression (Caroni 1987) to robust multivariate regression.

Consider the asymptotic representation of the estimator given by the influence function,

$$\mathbf{T}_n = \boldsymbol{\theta} + \frac{1}{n} \sum_{j=1}^n IF(\mathbf{z}_j, \mathbf{T}_n, G) + o(n^{-1/2}),$$

where $\boldsymbol{\theta} = (\mathbf{B}^t, \boldsymbol{\alpha})^t$ and G is the joint distribution of $\mathbf{z} = (\mathbf{x}^t, \mathbf{y}^t)^t$. We then obtain the first-order approximation for the residuals

$$\mathbf{r}_i \doteq \boldsymbol{\varepsilon}_i - \frac{1}{n} \sum_{j=1}^n [IF(\mathbf{z}_j, \hat{\mathbf{B}}^{LR}, G)^t \mathbf{x}_i + IF(\mathbf{z}_j, \hat{\boldsymbol{\alpha}}^{LR}, G)], \quad (17)$$

from which the covariance matrix $\text{cov}(\mathbf{r}_i)$ can be derived as outlined in the Appendix. Studentized residual distances are now defined as

$$sd_i = \sqrt{\mathbf{r}_i^t (\widehat{\text{cov}(\mathbf{r}_i)})^{-1} \mathbf{r}_i}.$$

Here $\widehat{\text{cov}(\mathbf{r}_i)}$ is the estimated covariance matrix for residual \mathbf{r}_i obtained by replacing the unknown error covariance matrix $\boldsymbol{\Sigma}_{\boldsymbol{\varepsilon}}$ with an estimate $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}$. If the estimate $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}$ is derived from the fitted model based on all data points, then we obtain *internally* studentized residual distances. If $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}$ comes from the model using all data points except \mathbf{z}_i when computing sd_i , then we obtain *externally* studentized residual distances. For large outlying points, there will be little difference between internally and externally studentized residual distances, because large outliers have only small influence on the LR-weighted MCD regression estimates, but for intermediate points externally studentized residuals will be larger than internally studentized residuals. To identify outliers, we compare the squared studentized residuals with quantiles of the χ_q^2 distribution. Figure 7 shows the externally studentized residuals for the pulp fiber and foam datasets analyzed earlier. The horizontal line in both plots is the square root of the 97.5% quantile of the corresponding chi-squared distribution. The labeled points in Figure 7 even lie above the 99.5% quantile of the chi-squared distribution. These outliers have also been labeled in the diagnostic plots (Figs. 3 and 4) in Section 5.

7. CONCLUSIONS

Least squares multivariate regression is sensitive to outliers in the dataset. Therefore, alternative methods that can detect and resist outliers are needed so that reliable results can be obtained also in the presence of outliers. Substantial work has been done to develop influence measures for multivariate regression (Hossain and Naik 1989; Barrett and Ling 1992; Hadi, Jones, and Ling 1995; Kim 1995; Seaver, Blankenship, and Triantis 1998). Much less has been done so far to develop robust estimators with bounded influence and/or high breakdown value. Singer and Sen (1985) and Koenker and Portnoy (1990) proposed robust methods based on M estimators. Methods based on affine equivariant sign and ranks were recently proposed by Ollila et al. (2002, 2003); however, these methods still have zero breakdown value.

We have shown that substituting robust estimates of location and scatter in the classical expressions for the slope, intercept and error scale yields a robust multivariate regression

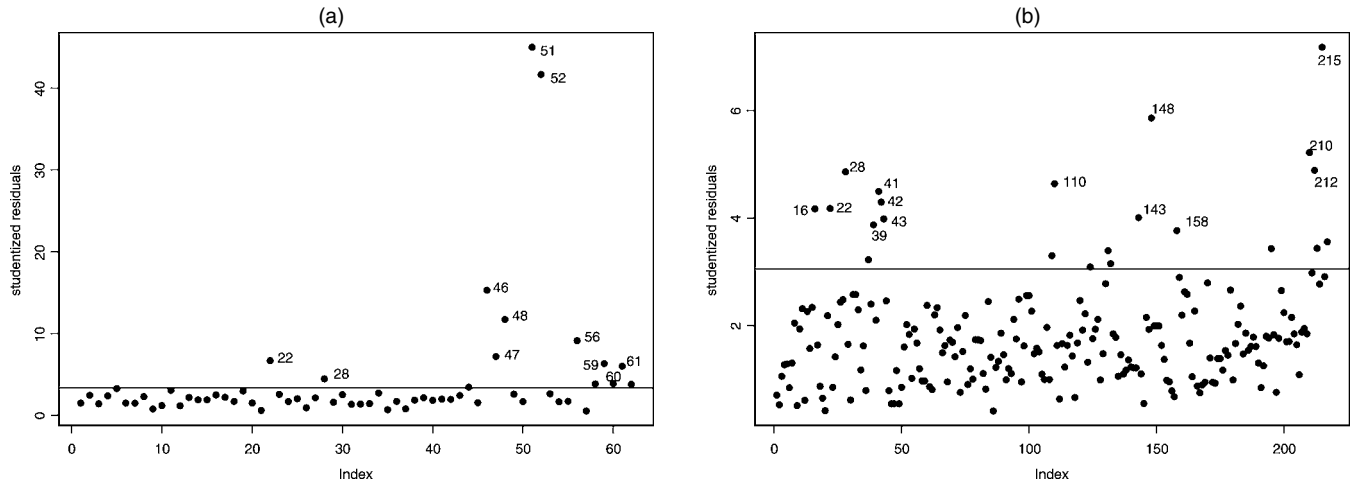


Figure 7. Externally Studentized Robust Residuals for (a) the Pulp Fiber Data and (b) the Foam Data. The horizontal line is the 97.5% quantile of the χ^2_q distribution.

method. By inserting the MCD estimator of location and scatter, we obtain a positive-breakdown and bounded-influence method, albeit with a rather low efficiency. To improve the efficiency, we have studied several types of reweighting schemes. We found that the best result is obtained by using the MCD-based robust distances to form a reweighted estimator of location and scatter, which then yields the initial regression. The robust residuals from this initial regression then give us the weights for the final regression. We call this the LR-weighted MCD regression. This approach gave the best finite-sample performance in our simulations and also yielded the highest asymptotic efficiency. Moreover, simulations with contaminated datasets indicated that its robustness properties also hold at finite samples. These simulations also showed that the LR-weighted MCD regression clearly outperforms classical least squares regression as well as univariate LTS regression applied to each of the responses separately. We illustrated the proposed method on two real data applications, where a new diagnostic plot turned out to be a very useful graphical tool to detect special points. Formal outlier diagnostics have been constructed based on studentized robust residual distances. MCD regression also is an essential part of robust principal component regression (Hubert and Verboven 2003) and robust partial least squares regression (Hubert and Vanden Branden 2003) procedures that are used to analyze high-dimensional data from spectra with several responses.

APPENDIX: PROOFS

To prove Theorem 1, we need the following lemma.

Lemma A.1. From the affine equivariance of the reweighted MCD location and scatter estimators $(\mathbf{t}_n^1, \mathbf{C}_n^1)$, it follows that the L-weighted MCD regression estimator $\mathbf{T}_n^L = ((\hat{\mathbf{B}}^L)^t, \hat{\boldsymbol{\alpha}}^L)^t$ is regression, y-affine, and x-affine equivariant.

Proof of Lemma A.1. Affine equivariance of $(\mathbf{t}_n^1, \mathbf{C}_n^1)$ means that for any nonsingular $(p+q) \times (p+q)$ matrix \mathbf{M} and any vector $\mathbf{a} \in \mathbb{R}^{p+q}$, it holds that $\mathbf{t}_n^1(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) = \mathbf{M} \mathbf{t}_n^1(\mathbf{Z}) + \mathbf{a}$ and

$\mathbf{C}_n^1(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) = \mathbf{M} \mathbf{C}_n^1 \mathbf{M}^t$. To prove regression equivariance, we take

$$\mathbf{M} = \begin{pmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{D}^t & \mathbf{I}_q \end{pmatrix} \quad \text{and} \quad \mathbf{a} = \begin{pmatrix} \mathbf{0} \\ \mathbf{w} \end{pmatrix}.$$

Then $\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t = (\mathbf{X}, \mathbf{Y})\mathbf{M}^t + \mathbf{1}_n \mathbf{a}^t = (\mathbf{X}, \mathbf{Y} + \mathbf{XD} + \mathbf{1}_n \mathbf{w}^t)$ and

$$(\mathbf{t}_n^1)_x(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) = (\mathbf{t}_n^1)_x(\mathbf{Z}),$$

$$(\mathbf{t}_n^1)_y(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) = (\mathbf{t}_n^1)_y(\mathbf{Z}) + \mathbf{D}^t(\mathbf{t}_n^1)_x(\mathbf{Z}) + \mathbf{w},$$

$$(\mathbf{C}_n^1)_{xx}(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) = (\mathbf{C}_n^1)_{xx}(\mathbf{Z}),$$

and

$$(\mathbf{C}_n^1)_{xy}(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) = (\mathbf{C}_n^1)_{xx}(\mathbf{Z})\mathbf{D} + (\mathbf{C}_n^1)_{xy}(\mathbf{Z}).$$

Therefore, we obtain

$$\begin{aligned} \hat{\mathbf{B}}^L(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) &= (\mathbf{C}_n^1)_{xx}^{-1}(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t)(\mathbf{C}_n^1)_{xy}(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) \\ &= \mathbf{D} + (\mathbf{C}_n^1)_{xx}^{-1}(\mathbf{Z})(\mathbf{C}_n^1)_{xy}(\mathbf{Z}) \\ &= \mathbf{D} + \hat{\mathbf{B}}^L(\mathbf{Z}) \end{aligned}$$

and

$$\begin{aligned} \hat{\boldsymbol{\alpha}}^L(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) &= (\mathbf{t}_n^1)_y(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) \\ &\quad - \hat{\mathbf{B}}^L(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t)(\mathbf{t}_n^1)_x(\mathbf{ZM}^t + \mathbf{1}_n \mathbf{a}^t) \\ &= (\mathbf{t}_n^1)_y(\mathbf{Z}) + \mathbf{D}^t(\mathbf{t}_n^1)_x(\mathbf{Z}) + \mathbf{w} \\ &\quad - (\mathbf{D} + \hat{\mathbf{B}}^L(\mathbf{Z}))^t(\mathbf{t}_n^1)_x(\mathbf{Z}) \\ &= (\mathbf{t}_n^1)_y(\mathbf{Z}) - \hat{\mathbf{B}}^L(\mathbf{Z})^t(\mathbf{t}_n^1)_x(\mathbf{Z}) + \mathbf{w} \\ &= \hat{\boldsymbol{\alpha}}^L(\mathbf{Z}) + \mathbf{w}, \end{aligned}$$

which is the desired result. To prove y-affine equivariance, we put

$$\mathbf{M} = \begin{pmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{C}^t \end{pmatrix} \quad \text{and} \quad \mathbf{a} = \begin{pmatrix} \mathbf{0} \\ \mathbf{d} \end{pmatrix}.$$

Finally, to prove x -affine equivariance, we put

$$\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_q \end{pmatrix} \quad \text{and} \quad \mathbf{a} = \begin{pmatrix} \mathbf{v} \\ 0 \end{pmatrix}.$$

Proof of Theorem 1

From Lemma A.1, we already have that the equivariance properties hold for the regression estimator based on reweighted MCD. It immediately follows that the reweighted regression estimator is also regression, x -affine, and y -affine equivariant, because the weights $d(\mathbf{r}_i^L)$ are invariant under these transformations, which can be proved similarly as in Lemma A.1.

Lemma A.2. Let \mathbf{Z}_n be a set of $n \geq p + q + 1$ observations and let \mathbf{t}_n^1 and \mathbf{C}_n^1 be the reweighted MCD estimators of location and scatter with $\min\{\varepsilon_n^*(\mathbf{t}_n^1, \mathbf{Z}_n), \varepsilon_n^*(\mathbf{C}_n^1, \mathbf{Z}_n)\} = \lceil n\gamma \rceil / n$, where $\gamma = (n - h)/n \leq (n - (p + q))/(2n)$. Then the estimator $\mathbf{T}_n^L = ((\hat{\mathbf{B}}^L)^t, \hat{\alpha}^L)^t$ also satisfies $\varepsilon_n^*(\mathbf{T}_n^L, \mathbf{Z}_n) = \lceil nr \rceil / n$.

Proof of Lemma A.2. Because the estimator \mathbf{T}_n^L is regression, y -affine, and x -affine equivariant (Lemma A.1), we may assume without loss of generality that $\mathbf{t}_n^1(\mathbf{Z}_n) = \mathbf{0}$. Let $\tilde{\mathbf{Z}}_n$ be a dataset obtained by replacing $m < \lceil n\gamma \rceil / n$ points from the original dataset \mathbf{Z}_n by arbitrary values. We first show that the slope $\hat{\mathbf{B}}^L(\tilde{\mathbf{Z}}_n)$ remains bounded. Denote the eigenvalues of $(\mathbf{C}_n^1)_{\mathbf{xx}}(\tilde{\mathbf{Z}}_n)$ by $\lambda_1((\mathbf{C}_n^1)_{\mathbf{xx}}(\tilde{\mathbf{Z}}_n)) \leq \dots \leq \lambda_p((\mathbf{C}_n^1)_{\mathbf{xx}}(\tilde{\mathbf{Z}}_n))$. Note that $\|\hat{\mathbf{B}}^L(\tilde{\mathbf{Z}}_n)\| = \|(\mathbf{C}_n^1)_{\mathbf{xx}}^{-1}(\tilde{\mathbf{Z}}_n)(\mathbf{C}_n^1)_{\mathbf{xy}}(\tilde{\mathbf{Z}}_n)\| \leq \|(\mathbf{C}_n^1)_{\mathbf{xx}}^{-1}(\tilde{\mathbf{Z}}_n)\| \times \|(\mathbf{C}_n^1)_{\mathbf{xy}}(\tilde{\mathbf{Z}}_n)\|$. Now we have that

$$\begin{aligned} \|(\mathbf{C}_n^1)_{\mathbf{xx}}^{-1}(\tilde{\mathbf{Z}}_n)\| &= \sup_{\|\mathbf{x}\| \neq 0} \frac{\|(\mathbf{C}_n^1)_{\mathbf{xx}}^{-1}(\tilde{\mathbf{Z}}_n)\mathbf{x}\|}{\|\mathbf{x}\|} \\ &= \left(\inf_{\|\mathbf{x}\| \neq 0} \frac{\|(\mathbf{C}_n^1)_{\mathbf{xx}}(\tilde{\mathbf{Z}}_n)\mathbf{x}\|}{\|\mathbf{x}\|} \right)^{-1} \\ &= \frac{1}{\lambda_1((\mathbf{C}_n^1)_{\mathbf{xx}}(\tilde{\mathbf{Z}}_n))}, \end{aligned}$$

which is bounded because the covariance matrix $\mathbf{C}_n^1(\tilde{\mathbf{Z}}_n)$ does not break down for $m < \lceil n\gamma \rceil / n$. Denote $\lambda_1(\tilde{\mathbf{Z}}_n) \leq \dots \leq \lambda_{p+q}(\tilde{\mathbf{Z}}_n)$ as the eigenvalues of $\mathbf{C}_n^1(\tilde{\mathbf{Z}}_n)$; then we have that $\|(\mathbf{C}_n^1)_{\mathbf{xy}}(\tilde{\mathbf{Z}}_n)\| \leq \|(\mathbf{C}_n^1)_{\mathbf{xx}}(\tilde{\mathbf{Z}}_n)\| \leq \lambda_{p+q}(\tilde{\mathbf{Z}}_n)$, which is also bounded for $m < \lceil n\gamma \rceil / n$. For the intercept, it clearly holds that $\|\hat{\alpha}^L(\tilde{\mathbf{Z}}_n)\| = \|(\mathbf{t}_n)_y(\tilde{\mathbf{Z}}_n) - (\hat{\mathbf{B}}^L)^t(\tilde{\mathbf{Z}}_n)(\mathbf{t}_n)_x(\tilde{\mathbf{Z}}_n)\| \leq \|(\mathbf{t}_n)_y(\tilde{\mathbf{Z}}_n)\| + \|(\hat{\mathbf{B}}^L)^t(\tilde{\mathbf{Z}}_n)\| \|(\mathbf{t}_n)_x(\tilde{\mathbf{Z}}_n)\|$ is bounded for $m < \lceil n\gamma \rceil / n$ because $\|(\hat{\mathbf{B}}^L)^t(\tilde{\mathbf{Z}}_n)\|$ and $\|(\mathbf{t}_n)_x(\tilde{\mathbf{Z}}_n)\|$ are bounded.

Proof of Theorem 2

Lemma A.2 shows that the L -weighted MCD regression estimator \mathbf{T}_n^L inherits the breakdown value of the reweighted MCD estimators. It now easily follows (under certain regularity conditions of the design matrix) that the reweighted regression estimator \mathbf{T}_n^{LR} inherits the breakdown value of the initial regression estimator \mathbf{T}_n^L .

Lemma A.3. Denote by \mathbf{t} and \mathbf{C} the functionals corresponding to the reweighted MCD location and scatter estimators; then the influence functions of $\hat{\mathbf{B}}^L$, $\hat{\alpha}^L$, and $\hat{\Sigma}_\varepsilon^L$ satisfy

$$IF(\mathbf{z}, \hat{\mathbf{B}}^L, H) = IF(\mathbf{z}, \mathbf{C}_{\mathbf{xy}}^1, H), \quad (\text{A.1})$$

$$IF(\mathbf{z}, \hat{\alpha}^L, H) = IF(\mathbf{z}, \mathbf{t}_y^1, H), \quad (\text{A.2})$$

and

$$IF(\mathbf{z}, \hat{\Sigma}_\varepsilon^L, H) = IF(\mathbf{z}, \mathbf{C}_{\mathbf{yy}}^1, H). \quad (\text{A.3})$$

Proof of Lemma A.3. First, we derive the influence function of the slope $\hat{\mathbf{B}}^L$. Because $\hat{\mathbf{B}}^L(H_\varepsilon) = (\mathbf{C}_{\mathbf{xx}}^1)^{-1}(H_\varepsilon)\mathbf{C}_{\mathbf{xy}}^1(H_\varepsilon)$, we obtain that

$$\begin{aligned} IF(\mathbf{z}, \hat{\mathbf{B}}^L, H) &= \frac{\partial}{\partial \varepsilon} \left((\mathbf{C}_{\mathbf{xx}}^1)^{-1}(H_\varepsilon)\mathbf{C}_{\mathbf{xy}}^1(H_\varepsilon) \right) \Big|_{\varepsilon=0} \\ &= IF(\mathbf{z}, (\mathbf{C}_{\mathbf{xx}}^1)^{-1}, H)\mathbf{C}_{\mathbf{xy}}^1(H) \\ &\quad + (\mathbf{C}_{\mathbf{xx}}^1)^{-1}(H)IF(\mathbf{z}, \mathbf{C}_{\mathbf{xy}}^1, H) \\ &= IF(\mathbf{z}, \mathbf{C}_{\mathbf{xy}}^1, H), \end{aligned}$$

because consistency of \mathbf{C}^1 yields $\mathbf{C}^1(H) = \mathbf{I}_{p+q}$. Similarly, with $\hat{\alpha}^L(H_\varepsilon) = \mathbf{t}_y^1(H_\varepsilon) - (\hat{\mathbf{B}}^L)^t(H_\varepsilon)\mathbf{t}_x^1(H_\varepsilon)$, we have that

$$\begin{aligned} IF(\mathbf{z}, \hat{\alpha}^L, H) &= \frac{\partial}{\partial \varepsilon} \left(\mathbf{t}_y^1(H_\varepsilon) - (\hat{\mathbf{B}}^L)^t(H_\varepsilon)\mathbf{t}_x^1(H_\varepsilon) \right) \Big|_{\varepsilon=0} \\ &= IF(\mathbf{z}, \mathbf{t}_y^1, H) - IF(\mathbf{z}, \hat{\mathbf{B}}^L, H)^t \mathbf{t}_x^1(H) \\ &\quad - (\hat{\mathbf{B}}^L)^t(H)IF(\mathbf{z}, \mathbf{t}_x^1, H) \\ &= IF(\mathbf{z}, \mathbf{t}_y^1, H), \end{aligned}$$

because $\mathbf{t}^1(H) = \mathbf{0}$ and $\hat{\mathbf{B}}^L(H) = \mathbf{0}$. Finally, $\hat{\Sigma}_\varepsilon^L(H_\varepsilon) = \mathbf{C}_{\mathbf{yy}}^1(H_\varepsilon) - (\hat{\mathbf{B}}^L)^t(H_\varepsilon)\mathbf{C}_{\mathbf{xx}}^1(H_\varepsilon)(\hat{\mathbf{B}}^L)(H_\varepsilon)$ yields

$$\begin{aligned} IF(\mathbf{z}, \hat{\Sigma}_\varepsilon^L, H) &= \frac{\partial}{\partial \varepsilon} \left(\mathbf{C}_{\mathbf{yy}}^1(H_\varepsilon) - (\hat{\mathbf{B}}^L)^t(H_\varepsilon)\mathbf{C}_{\mathbf{xx}}^1(H_\varepsilon)\hat{\mathbf{B}}^L(H_\varepsilon) \right) \Big|_{\varepsilon=0} \\ &= IF(\mathbf{z}, \mathbf{C}_{\mathbf{yy}}^1, H) - IF(\mathbf{z}, \hat{\mathbf{B}}^L, H)^t \mathbf{C}_{\mathbf{xx}}^1(H)\hat{\mathbf{B}}^L(H) \\ &\quad - (\hat{\mathbf{B}}^L)^t(H)IF(\mathbf{z}, \mathbf{C}_{\mathbf{xx}}^1, H)\hat{\mathbf{B}}^L(H) \\ &\quad - (\hat{\mathbf{B}}^L)^t(H)\mathbf{C}_{\mathbf{xx}}^1(H)IF(\mathbf{z}, \hat{\mathbf{B}}^L, H) \\ &= IF(\mathbf{z}, \mathbf{C}_{\mathbf{yy}}^1, H), \end{aligned}$$

because $\hat{\mathbf{B}}^L(H) = \mathbf{0}$.

Proof of Theorem 3

Combining Lemma A.3 with the results of Croux and Haesbroeck (1999), we obtain that the influence functions of $\hat{\mathbf{B}}^L$, $\hat{\alpha}^L$, and $\hat{\Sigma}_\varepsilon^L$ equal

$$\begin{aligned} IF(\mathbf{z}, \hat{\mathbf{B}}_{jk}^L, H) &= \left[\frac{a_2}{c_2} I(\|\mathbf{z}\|^2 \leq q_\gamma) + \frac{1}{d_1} I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \right] x_j y_k, \\ IF(\mathbf{z}, (\hat{\alpha}^L)_j, H) &= \left[\left(1 - \frac{d_1}{1 - \delta_l} \right) \frac{1}{c_1} I(\|\mathbf{z}\|^2 \leq q_\gamma) \right. \\ &\quad \left. + \frac{1}{1 - \delta_l} I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \right] y_j, \\ IF(\mathbf{z}, (\hat{\Sigma}_\varepsilon^L)_{jk}, H) &= \left[\frac{a_2}{c_2} I(\|\mathbf{z}\|^2 \leq q_\gamma) + \frac{1}{d_1} I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \right] y_j y_k, \end{aligned}$$

and

$$\begin{aligned}
 IF(\mathbf{z}, (\hat{\Sigma}_\varepsilon^L)_{jj}, H) &= \left[\frac{a_2}{c_2} I(\|\mathbf{z}\|^2 \leq q_\gamma) + \frac{1}{d_1} I(\|\mathbf{z}\|^2 \leq q_{\delta_l}) \right] y_j^2 \\
 &+ \frac{a_2}{2c_2} \|\mathbf{y}\|^2 I(\|\mathbf{z}\|^2 \leq q_\gamma) - 1 \\
 &+ \frac{p+q+2}{2} \frac{a_2}{c_2[a_3 - (p+q)a_4]} \\
 &\times \left[a_4 \|\mathbf{z}\|^2 I(\|\mathbf{z}\|^2 \leq q_\gamma) \right. \\
 &\left. + \frac{a_3}{p+q} q_\gamma (1 - \gamma - I(\|\mathbf{z}\|^2 \leq q_\gamma)) - 1 \right].
 \end{aligned}$$

Denote the incomplete gamma function by $\Gamma(u; v) = \Gamma(u)^{-1} \times \int_0^v t^{u-1} e^{-t} dt$. Then the constants a_1 , a_2 , a_3 , and a_4 are given by $a_1 = 1/d_1$, $a_2 = (d_1 - d_2)/d_1$, $a_3 = c_2/c_1$, and $a_4 = \frac{1}{2} - \frac{1}{2c_1} [c_2 - \frac{q_\gamma}{p+q} (c_1 + \gamma - 1)]$, where $c_1 = \Gamma(\frac{p+q}{2} + 1; \frac{q_\gamma}{2})$, $c_2 = \Gamma(\frac{p+q}{2} + 2; \frac{q_\gamma}{2})$, $d_1 = \Gamma(\frac{p+q}{2} + 1; \frac{q_{\delta_l}}{2})$, and $d_2 = \Gamma(\frac{p+q}{2} + 2; \frac{q_{\delta_l}}{2})$.

It can be easily shown that the influence functions of reweighted regression estimators defined by (7) and (8) are connected to the influence functions of the initial regression estimators $\hat{\mathbf{B}}^L$, $\hat{\alpha}^L$, and $\hat{\Sigma}_\varepsilon^L$ through

$$\begin{aligned}
 IF(\mathbf{z}, \hat{\mathbf{B}}^{LR}, H) &= \left(1 - \frac{d_1^R}{1 - \delta_r} \right) IF(\mathbf{z}, \hat{\mathbf{B}}^L, H) + \frac{I(\|\mathbf{y}\|^2 \leq q_{\delta_r})}{1 - \delta_r} \mathbf{xy}^t, \\
 IF(\mathbf{z}, \hat{\alpha}^{LR}, H) &= \left(1 - \frac{d_1^R}{1 - \delta_r} \right) IF(\mathbf{z}, \hat{\alpha}^L, H) + \frac{I(\|\mathbf{y}\|^2 \leq q_{\delta_r})}{1 - \delta_r} \mathbf{y},
 \end{aligned}$$

and

$$\begin{aligned}
 IF(\mathbf{z}, \hat{\Sigma}_\varepsilon^{LR}, H) &= \frac{d_1^R - d_2^R}{d_1^R} \left(IF(\mathbf{z}, (\hat{\Sigma}_\varepsilon^L)_{jj}, H) + \frac{1}{2} \text{tr}(IF(\mathbf{z}, \hat{\Sigma}_\varepsilon^L, H)) \mathbf{I}_q \right) \\
 &+ \frac{I(\|\mathbf{y}\|^2 \leq q_{\delta_r})}{d_1^R} \mathbf{yy}^t - \mathbf{I}_q.
 \end{aligned}$$

The constants d_1^R and d_2^R are given by $d_1^R = \Gamma(\frac{q}{2} + 1; \frac{q_{\delta_r}}{2})$ and $d_2^R = \Gamma(\frac{q}{2} + 2; \frac{q_{\delta_r}}{2})$. Note that the foregoing results extend the expressions for the influence functions of reweighted multivariate location and scatter functionals given by Lopuhaä (1999).

Studentized Residual Distances. First, note that for any $(\mathbf{x}^t, \mathbf{y}^t)\mathbf{t}$ such that

$$\mathbf{x} = \Sigma_{\mathbf{xx}}^{1/2} \mathbf{u} + \mu_{\mathbf{x}}$$

and

$$\mathbf{y} = \mathbf{B}^t \mathbf{x} + \alpha + \Sigma_\varepsilon^{1/2} \varepsilon$$

with $(\mathbf{u}^t, \varepsilon^t)^t \sim H$, that is, the standard Gaussian distribution, it follows from Theorem 3 and the equivariance properties in

Theorem 1 that the influence function of $\hat{\mathbf{B}}^{LR}$ and $\hat{\alpha}^{LR}$ at the joint distribution G of $(\mathbf{x}^t, \mathbf{y}^t)^t$ can be written as

$$\begin{aligned}
 IF(\mathbf{z}, \hat{\mathbf{B}}^{LR}, G) &= \left(\left(1 - \frac{d_1^R}{1 - \delta_r} \right) \right. \\
 &\times \left[\frac{a_2}{c_2} I(d^2(\mathbf{z}) \leq q_\gamma) + \frac{1}{d_1} I(d^2(\mathbf{z}) \leq q_{\delta_l}) \right] \\
 &\left. + \frac{I(d^2(\mathbf{r}) \leq q_{\delta_r})}{1 - \delta_r} \right) \Sigma_{\mathbf{xx}}^{-1} (\mathbf{x} - \mu_{\mathbf{x}}) \mathbf{r}^t
 \end{aligned}$$

and

$$\begin{aligned}
 IF(\mathbf{z}, \hat{\alpha}^{LR}, G) &= \left(\left(1 - \frac{d_1^R}{1 - \delta_r} \right) \right. \\
 &\times \left[\left(1 - \frac{d_1}{1 - \delta_l} \right) \frac{1}{c_1} I(d^2(\mathbf{z}) \leq q_\gamma) \right. \\
 &\left. + \frac{1}{1 - \delta_l} I(d^2(\mathbf{z}) \leq q_{\delta_l}) \right] \\
 &\left. + \frac{I(d^2(\mathbf{r}) \leq q_{\delta_r})}{1 - \delta_r} \right) \mathbf{r} - IF(\mathbf{z}, \hat{\mathbf{B}}^{LR}, G)^t \mu_{\mathbf{x}},
 \end{aligned}$$

where $d^2(\mathbf{z})$ and $d^2(\mathbf{r})$ are the squared robust distances of the point \mathbf{z} and its corresponding residual. Substituting the foregoing expressions for the influence functions in the right side of (17), the following approximation for the covariance matrix of residual \mathbf{r}_i can be obtained:

$$\begin{aligned}
 \text{cov}(\mathbf{r}_i) &\doteq \left(1 - \frac{2}{n} \left[f_i(\mathbf{z}_i) + \frac{d_1^R}{1 - \delta_r} (d^2(\mathbf{x}_i) + 1) \right] \right. \\
 &+ \frac{1}{n^2} \sum_{j=1}^n \left[f_i^2(\mathbf{z}_j) \frac{d_1^R}{(1 - \delta_r)^2} (d_{ji} + 1)^2 \right. \\
 &\left. \left. + \frac{2d_1^R}{1 - \delta_r} f_i(\mathbf{z}_j) (d_{ji} + 1) \right] \right) \Sigma_\varepsilon,
 \end{aligned}$$

where $d_{ji} = (\mathbf{x}_j - \mu_{\mathbf{x}})^t \Sigma_{\mathbf{xx}}^{-1} (\mathbf{x}_i - \mu_{\mathbf{x}})$ and $f_i(\mathbf{z}_j) = (1 - \frac{d_1^R}{1 - \delta_r}) \times [\frac{a_2}{c_2} I(d^2(\mathbf{z}) \leq q_\gamma) + \frac{1}{d_1} I(d^2(\mathbf{z}) \leq q_{\delta_l})] d_{ji} + (1 - \frac{d_1^R}{1 - \delta_r}) [(1 - \frac{d_1}{1 - \delta_l}) \frac{1}{c_1} I(d^2(\mathbf{z}) \leq q_\gamma) + \frac{1}{1 - \delta_l} I(d^2(\mathbf{z}) \leq q_{\delta_l})]$.

[Received June 2000. Revised September 2003.]

REFERENCES

- Alqallaf, F. A., Konis, K. P., Martin, R. D., and Zamar, R. H. (2002), "Scalable Robust Covariance and Correlation Estimates for Data Mining," in *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Edmonton, Alberta, Canada, pp. 14–23.
- Barrett, B. E., and Ling, R. F. (1992), "General Classes of Influence Measures for Multivariate Regression," *Journal of the American Statistical Association*, 87, 184–191.
- Bickel, P. J., and Lehmann, E. L. (1976), "Descriptive Statistics for Nonparametric Models III: Dispersion," *The Annals of Statistics*, 4, 1139–1159.
- Breiman, L., and Friedman, J. H. (1997), "Predicting Multivariate Responses in Multiple Linear Regression," *Journal of the Royal Statistical Society, Ser. B*, 59, 3–54.
- Butler, R. W., Davies, P. L., and Jhun, M. (1993), "Asymptotics for the Minimum Covariance Determinant Estimator," *The Annals of Statistics*, 21, 1385–1400.
- Caroni, C. (1987), "Residuals and Influence in the Multivariate Linear Model," *The Statistician*, 36, 365–370.

- Cook, D. R., and Setodji, M. C. (2003), "A Model-Free Test for Reduced Rank in Multivariate Regression," *Journal of the American Statistical Association*, 98, 340–351.
- Croux, C., and Haesbroeck, G. (1999), "Influence Function and Efficiency of the Minimum Covariance Determinant Scatter Matrix Estimator," *Journal of Multivariate Analysis*, 71, 161–190.
- Davies, L. (1987), "Asymptotic Behavior of S-Estimators of Multivariate Location Parameters and Dispersion Matrices," *The Annals of Statistics*, 15, 1269–1292.
- Davis, J. B., and McKean, J. W. (1993), "Rank-Based Methods for Multivariate Linear Models," *Journal of the American Statistical Association*, 88, 245–251.
- Donoho, D. L., and Huber, P. J. (1983), "The Notion of Breakdown Point," in *A Festschrift for Erich Lehmann*, eds. P. J. Bickel, K. A. Doksum, and J. L. Hodges, Belmont, CA: Wadsworth, pp. 157–184.
- Gleser, L. J. (1992), "The Importance of Assessing Measurement Reliability in Multivariate Regression," *Journal of the American Statistical Association*, 87, 696–707.
- Hadi, A. S., Jones, W. D., and Ling, R. F. (1995), "A Unifying Representation of Some Case-Deletion Influence Measures in Univariate and Multivariate Linear Regression," *Journal of Statistical Planning and Inference*, 46, 123–135.
- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., and Stahel, W. A. (1986), *Robust Statistics: The Approach Based on Influence Functions*, New York: Wiley.
- He, X., and Wang, G. (1996), "Cross-Checking Using the Minimum Volume Ellipsoid Estimator," *Statistica Sinica*, 6, 367–374.
- Hossain, A., and Naik, D. N. (1989), "Detection of Influential Observations in Multivariate Regression," *Journal of Applied Statistics*, 16, 25–37.
- Hubert, M., and Vanden Branden, K. (2003), "Robust Methods for Partial Least Squares Regression," *Journal of Chemometrics*, 17, 537–549.
- Hubert, M., and Verboven, S. (2003), "A Robust PCR method for High-Dimensional Regressors," *Journal of Chemometrics*, 17, 438–452.
- Johnson, R. A., and Wichern, D. W. (1998), *Applied Multivariate Statistical Analysis* (4th ed.), Upper Saddle River, NJ: Prentice-Hall.
- Kent, J. T., and Tyler, D. E. (1996), "Constrained M-Estimation for Multivariate Location and Scatter," *The Annals of Statistics*, 24, 1346–1370.
- Kim, M. G. (1995), "Local Influence in Multivariate Regression," *Communications in Statistics, Part A—Theory and Methods*, 24, 1271–1278.
- Knorr, E. M., Ng, R. T., and Zamar, R. H. (2001), "Robust Space Transformation for Distance-Based Operations," in *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco.
- Koenker, R., and Portnoy, S. (1990), "M Estimation of Multivariate Regressions," *Journal of the American Statistical Association*, 85, 1060–1068.
- Lee, J. (1992), "Relationships Between Properties of Pulp-Fibre and Paper," unpublished doctoral thesis, University of Toronto, Faculty of Forestry.
- Lopuhaä, H. P. (1989), "On the Relation Between S-Estimators and M-Estimators of Multivariate Location and Covariance," *The Annals of Statistics*, 17, 1662–1683.
- (1991), "Multivariate τ -Estimators for Location and Scatter," *Canadian Journal of Statistics*, 19, 307–321.
- (1999), "Asymptotics of Reweighted Estimators of Multivariate Location and Scatter," *The Annals of Statistics*, 27, 1638–1665.
- 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-Estimates of Multivariate Location and Scatter," *The Annals of Statistics*, 4, 51–67.
- Maronna, R. A., and Morgenthaler, S. (1986), "Robust Regression Through Robust Covariances," *Communications in Statistics, Part A—Theory and Methods*, 15, 1347–1365.
- Maronna, R., and Zamar, R. H. (2002), "Robust Multivariate Estimates for High-Dimensional Data Sets," *Technometrics*, 44, 307–317.
- McKean, J. W., Sheather, S. J., and Hettmansperger, T. P. (1990), "Regression Diagnostics for Rank-Based Methods," *Journal of the American Statistical Association*, 85, 1018–1028.
- (1993), "The Use and Interpretation of Residuals Based on Robust Estimation," *Journal of the American Statistical Association*, 88, 1254–1263.
- Ollila, E., Hettmansperger, T. P., and Oja, H. (2002), "Estimates of Regression Coefficients Based on Sign Covariance Matrix," *Journal of the Royal Statistical Society, Ser. B*, 64, 447–466.
- Ollila, E., Oja, H., and Koivunen, V. (2003), "Estimates of Regression Coefficients Based on Rank Covariance Matrix," *Journal of the American Statistical Association*, 98, 90–98.
- Pison, G., Van Aelst, S., and Willems, G. (2002), "Small-Sample Corrections for LTS and MCD," *Metrika*, 55, 111–123.
- Rocke, D. M. (1996), "Robustness Properties of S-Estimators of Multivariate Location and Shape in High Dimension," *The Annals of Statistics*, 24, 1327–1345.
- 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.
- (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, pp. 283–297.
- Rousseeuw, P. J., and Leroy, A. M. (1987), *Robust Regression and Outlier Detection*, New York: Wiley.
- Rousseeuw, P. J., and Van Driessen, K. (1999), "A Fast Algorithm for the Minimum Covariance Determinant Estimator," *Technometrics*, 41, 212–223.
- (2000), "An Algorithm for Positive-Breakdown Regression Based on Concentration Steps," in *Data Analysis: Scientific Modeling and Practical Application*, eds. W. Gaul, O. Opitz, and M. Schader, New York: Springer-Verlag, pp. 335–346.
- Rousseeuw, P. J., and van Zomeren, B. C. (1990), "Unmasking Multivariate Outliers and Leverage Points," *Journal of the American Statistical Association*, 85, 633–651.
- Seaver, B., Blankenship, A., and Triantis, K. P. (1998), "Identifying Potentially Influential Subsets in Multivariate Regression," technical report, University of Tennessee, Dept. of Statistics.
- Singer, J. M., and Sen, P. K. (1985), "M-Methods in Multivariate Linear Models," *Journal of Multivariate Analysis*, 17, 168–184.
- 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.
- Zuo, Y., and Cui, H. (2002), "Depth-Weighted Scatter Estimators," *The Annals of Statistics*, to appear.
- Zuo, Y., Cui, H., and He, X. (2001), "On the Stahel-Donoho Estimator and Depth-Weighted Means of Multivariate Data," *The Annals of Statistics*, 32, 167–188.