

## DETECTION OF MULTIVARIATE OUTLIERS WITH DISPERSION SLIPPAGE IN ELLIPTICALLY SYMMETRIC DISTRIBUTIONS<sup>1</sup>

BY RITA DAS AND BIMAL K. SINHA

*University of Pittsburgh and University of Maryland Baltimore County*

In this paper Ferguson's univariate normal results for detection of outliers with variance slippage are extended to the multivariate elliptically symmetric case with dispersion slippage. The locally optimum test we derive possesses all three robustness properties, optimality, null and nonnull, and is based on Mardia's multivariate kurtosis statistic.

**1. Introduction.** Recently, Ferguson's (1961) pioneering work on the detection of outliers with mean slippage in samples from a univariate normal distribution has been generalized to the multivariate case. While Schwager and Margolin (1982) considered a multivariate normal distribution, Sinha (1984) treated an arbitrary elliptically symmetric multivariate distribution. The main result proved by these authors can be viewed as a robustness property of the use of Mardia's (1970) multivariate kurtosis statistic as a locally optimum test statistic to detect outliers with mean slippage.

In this paper we consider the other aspect of Ferguson's work, namely, the detection of outliers with variance slippage in a univariate normal distribution and extend this to the case of an arbitrary elliptically symmetric multivariate distribution with dispersion slippage in the same spirit as in Sinha (1984). The locally best invariant locally unbiased test we derive has some very interesting features. First, as in Ferguson (1961) in the univariate case, the same test statistic, namely, Mardia's multivariate kurtosis statistic as in the mean slippage case (Schwager and Margolin (1982); Sinha (1984)), turns out to be locally optimum also in the case of dispersion slippage. Second, this locally best invariant test is both optimality robust and null robust. Third, this test is also nonnull robust. We refer to Kariya (1981) and Kariya and Sinha (1985) for the notions and various results on optimality, null and nonnull robustness. The multivariate outlier problem with dispersion slippage is formulated in Section 2 and the main result is contained in Section 3. It may be mentioned that Schwager, in his Yale Ph.D. thesis (1979), formulated the multivariate dispersion slippage problem under normality, demonstrated the invariance discussed here, and conjectured that the Mardia kurtosis statistic was optimal.

---

Received March 1985; revised February 1986.

<sup>1</sup>This work was sponsored by the Air Force Office of Scientific Research under Contract F49620-82-K-0001. Reproduction in whole or in part is permitted for any purpose of the United States Government.

AMS 1980 *subject classifications*. Primary 62A05, 62H15; secondary 62H10, 62E15.

*Key words and phrases*. Locally best invariant, maximal invariant, variance slippage, outliers, robustness, Wijsman's representation theorem.

**2. The multivariate outlier problem with dispersion slippage.** Consider a sample of size  $n$  from a multivariate distribution. We will denote the sample by  $X: n \times p$  and assume that the following model holds:

$$(2.1) \quad X = \mathbf{1}\mu' + U\Sigma^{1/2},$$

where  $\mathbf{1}$  is the unit  $n \times 1$  vector,  $\mu$  is the common unknown  $p \times 1$  mean vector of the rows of  $X$ ,  $\Sigma > 0$  (positive definite  $p \times p$ ), and the random error component  $U$  has a distribution  $\mathcal{L}(U) \in \mathcal{F}(0, I_n \times I_p)$ , the class of  $np$ -dimensional elliptically symmetric distributions about 0 with scale matrix  $I_p$  and with density given by

$$(2.2) \quad f(u) = \phi(\text{tr } \mathbf{u}'\mathbf{u}),$$

where  $\phi: [0, \infty) \rightarrow [0, \infty)$ ,  $U \in \mathcal{U} = \{U: n \times p | \text{rank } U = p\}$ . This amounts to the specification that  $\mathcal{L}(X) \in \mathcal{F}(\mathbf{1}\mu', I_n \times \Sigma)$ . Moreover, we will assume that  $\phi$  satisfies

$$\int_{R^p} \int_{\text{Gl}(p)} \phi(\text{tr } c'c + \mathbf{u}'\mathbf{u}) |c'c|^{(n-p-1)/2} dc d\mathbf{u} < \infty.$$

The possibility of outliers with dispersion slippage can be detected by testing the model

$$(2.3) \quad \mathcal{L}(X) \in \mathcal{F}(\mathbf{1}\mu', I_n \times \Sigma) \quad \text{versus} \quad \mathcal{L}(X) \in \mathcal{F}(\mathbf{1}\mu', D \times \Sigma), \quad D \neq I.$$

$D$  is a function of a scalar parameter  $\Delta$ ,  $D = D_\Delta^2(\delta) = \text{diag}(\delta_1, \dots, \delta_n)$  with  $\delta_i = \exp(\Delta a_{\nu_i})$ ,  $i = 1, \dots, n$ , where  $(\nu_1, \dots, \nu_n)$  is an unknown permutation of  $(1, \dots, n)$  and  $a_1, \dots, a_n$  are arbitrary constants at least two of which are dissimilar and some of which may be zero. In this formulation, which is similar to Ferguson's (1961) in the univariate case, unless  $\Delta = 0$ , the observation  $X_i$  corresponding to the  $i$ th row of  $X$  is an outlier if  $a_i$  is nonzero.

The general multivariate outlier problem with dispersion slippage thus consists of models  $\mathcal{L}(X) \in \mathcal{F}(\mathbf{1}\mu', I_n \times \Sigma)$  and  $\mathcal{L}(X) \in \mathcal{F}(\mathbf{1}\mu', D \times \Sigma)$ ,  $D = D_\Delta^2(\delta)$ , and the null hypothesis  $H_0: \Delta = 0$  against the alternative hypothesis  $H_1: \Delta \neq 0$ .

Following Ferguson (1961), Schwager and Margolin (1982), and Sinha (1984), it is clear that the above testing problem remains invariant under the action of the group  $\mathcal{G} = \mathcal{P} \times \text{Gl}(p) \times R^p$  where  $\mathcal{P}$  denotes the (finite) group of all  $n \times n$  permutation matrices with elements  $\Gamma_\alpha$ ,  $\text{Gl}(p)$  the full linear group of  $p \times p$  nonsingular matrices with elements  $C$  and  $R^p$  the Euclidean  $p$  space. The three (sub)group operations are defined by  $g\mathbf{x} = \Gamma_\alpha \mathbf{x}C + \mathbf{1}\mu^{*\prime}$ ,  $\Gamma_\alpha \in \mathcal{P}$ ,  $C \in \text{Gl}(p)$ , and  $\mu^* \in R^p$ . In the next section we apply Wijsman's (1967) theorem to derive the distribution of a maximal invariant statistic.

**3. Main results.** By invariance of the problem we may assume without any loss of generality  $\mu = \mathbf{0}$  and  $\Sigma = I_p$ .

Let  $T = t(x)$  be a maximal invariant under the transformation  $\mathcal{G}$  and let  $P_\Delta^T$  be the distribution induced by  $T$  under  $\Delta$ . Then, using a version of Wijsman's (1967) theorem, the pdf of  $T$  under  $\Delta$  with respect to  $P_0^T$  evaluated at  $T = t(x)$  is

given by

$$(3.1) \quad (dP_{\Delta}^T/dP_0^T)(t(x)) = \frac{\int_{\mathcal{G}} f(gx|\Delta)|C'C|^{n/2} d\nu(g)}{\int_{\mathcal{G}} f(gx|\Delta = 0)|C'C|^{n/2} d\nu(g)},$$

where

$$f(x|\Delta) = \exp(-\Delta p \sum a_i/2) \phi(\text{tr}(x'D_{\Delta}^{-2}x)),$$

$$D_{\Delta}^{-2} = \text{diag}(\exp(-\Delta a_{\nu_1}), \dots, \exp(-\Delta a_{\nu_n})),$$

and  $\nu$  is a left invariant measure on  $\mathcal{G}$ . We take  $\nu = \nu_1 \times \nu_2$ , where  $\nu_1$  is the discrete uniform probability measure with mass  $1/n!$  at each of the  $n!$  elements  $\Gamma_{\alpha} \in \mathcal{P}$ , and  $d\nu_2(C, \mu^*) = dC d\mu^*/|C'C|^{(p+1)/2}$ , where  $d\mu^*$  is the Lebesgue measure on  $R^p$ . This is a left invariant measure on the affine group  $\text{Gl}(p) \times R^p$ . The following result is crucial in the derivation of an invariant test.

LEMMA 3.1. *The ratio of the pdf 's in expression (3.1) is evaluated as*

$$(3.2) \quad \tau_1 \tau_2 \sum_{\alpha} (|S_{\alpha}^*|/|S|)^{-(n-1)/2} \frac{n^{p/2}}{n!},$$

where  $S = x'x - n\bar{x}\bar{x}'$ ,  $\bar{x}$  is the sample mean vector,  $\tau_1 = \exp(-p\Delta \sum a_i/2)$ ,  $\tau_2 = \tau^{-p/2}$  with  $\tau = \mathbf{1}'D_{\Delta}^{-2}\mathbf{1} = \sum \exp(-\Delta a_i)$ , and

$$(3.3) \quad S_{\alpha}^* = x'\Gamma_{\alpha}'(D_{\Delta}^{-2} - D_{\Delta}^{-2}\mathbf{1}\mathbf{1}'D_{\Delta}^{-2}/\tau)\Gamma_{\alpha}x.$$

PROOF. The numerator  $N_{\Delta}$  (say) of (3.1) can be written as

$$(3.4) \quad N_{\Delta} = \frac{1}{n!} \tau_1 \sum_{\alpha} \int_{\text{Gl}(p)} \int_{R^p} \phi(\text{tr}(x'D_{\Delta}^{-2}x|x \rightarrow gx)) |C'C|^{(n-p-1)/2} d\mu^* dC,$$

where  $\phi(\text{tr}(\dots|x \rightarrow gx))$  stands for the value of  $\phi$  evaluated when  $x$  is replaced by  $gx$ . Using  $gx = \Gamma_{\alpha}x C + \mathbf{1}\mu^{*'}$ , the argument of  $\phi$ , after the substitution  $x = gx$ , simplifies to

$$(3.5) \quad \begin{aligned} & \text{tr}[(\Gamma_{\alpha}x C + \mathbf{1}\mu^{*'})' D_{\Delta}^{-2} (\Gamma_{\alpha}x C + \mathbf{1}\mu^{*'})] \\ &= \text{tr}[C'x'\Gamma_{\alpha}'D_{\Delta}^{-2}\Gamma_{\alpha}x C + \tau\mu^{*'}\mu^{*'} + 2C'x'\Gamma_{\alpha}'D_{\Delta}^{-2}\mathbf{1}\mu^{*'}] \\ &= \text{tr}[\tau\mathbf{c}_{\alpha}^*\mathbf{c}_{\alpha}^{*'} + C'x'\Gamma_{\alpha}'(D_{\Delta}^{-2} - D_{\Delta}^{-2}\mathbf{1}\mathbf{1}'D_{\Delta}^{-2}/\tau)\Gamma_{\alpha}x C]. \end{aligned}$$

In the last equality above,  $\mathbf{c}_{\alpha}^* = \mu^{*'} + C'x'\Gamma_{\alpha}'D_{\Delta}^{-2}/\tau$ . Since  $d\mu^* = d\mathbf{c}_{\alpha}^*$ , using a result of Dawid (1977), integration with respect to  $\mathbf{c}_{\alpha}^*$  over  $R^p$  yields

$$(3.6) \quad \begin{aligned} N_{\Delta} &= \frac{1}{n!} \tau_1 \tau_2 \sum_{\alpha} \int_{\text{Gl}(p)} \tilde{\phi}[\text{tr} C'x'\Gamma_{\alpha}'(D_{\Delta}^{-2} - D_{\Delta}^{-2}\mathbf{1}\mathbf{1}'D_{\Delta}^{-2}/\tau)\Gamma_{\alpha}x C] \\ & \quad \times |C'C|^{(n-p-1)/2} dC \end{aligned}$$

for  $\tilde{\phi}: [0, \infty) \rightarrow [0, \infty)$  given by  $\tilde{\phi}(z) \equiv \int_{R^p} \phi(z + \mathbf{u}'\mathbf{u}) d\mathbf{u}$ . Now  $x'\Gamma_{\alpha}'(D_{\Delta}^{-2} -$

$D_{\Delta}^{-2}\mathbf{11}'D_{\Delta}^{-2}/\tau)\Gamma_{\alpha}x = S_{\alpha}^*$  (say) is p.d. by our assumption  $n \geq p + 1$ . Writing  $S_{\alpha}^* = S_{\alpha}^{*1/2}S_{\alpha}^{*1/2}$  where  $S_{\alpha}^{*1/2}$  is the positive square root of  $S_{\alpha}^*$  and making the transformation  $C \rightarrow S_{\alpha}^{*1/2}C$ ,  $N_{\Delta}$  reduces to

$$(3.7) \quad N_{\Delta} = \frac{1}{n!} \tau_1 \tau_2 \left\{ \sum_{\alpha} |S_{\alpha}^*|^{-(n-1)/2} \right\} \int_{\text{Gl}(p)} \tilde{\phi}[\text{tr } C'C] |C'C|^{(n-p-1)/2} dC.$$

Since the denominator of (3.1) corresponds to  $N_{\Delta}$  with  $\Delta = 0$  and since  $S_{\alpha}^*$  at  $\Delta = 0$  is  $x'\Gamma_{\alpha}'(I_n - \mathbf{11}'/n)\Gamma_{\alpha}x = x'x - n\bar{x}\bar{x}' = S$ , the sample sum of squares and products matrix, the lemma follows.  $\square$

**REMARK 3.1.** Since  $dP_{\Delta}^T/dP_0^T(t(x))$  is independent of  $\phi$ , it follows that any null robust invariant test is also nonnull robust. In particular, the locally optimum invariant test derived below is null and hence nonnull robust. This is yet another example of a test for covariance structure that is nonnull robust (see Kariya and Sinha (1985)).

We now proceed to evaluate the expression in (3.2). It turns out that there is no uniformly most powerful invariant test for this problem. To derive a locally best invariant test, we expand the expression in (3.2) locally in  $\Delta$  around  $\Delta = 0$ . Toward this end, note that

$$(3.8) \quad D_{\Delta}^{-2} - I_n = \Delta D(\mathbf{a}) + \Delta^2 D(\mathbf{a}^2)/2 + o(\Delta^2),$$

where  $\mathbf{a}^j = (a_1^j, \dots, a_n^j)'$ ,  $D(\mathbf{a}^j) = \text{diag}(a_1^j, \dots, a_n^j)$ ,  $j = 1, 2$ , and

$$(3.9) \quad D_{\Delta}^{-2}\mathbf{11}'D_{\Delta}^{-2}/\tau = [\mathbf{11}' - \Delta(\mathbf{a1}' + \mathbf{1a}' - \bar{a}\mathbf{11}') + \Delta^2(\mathbf{a} - \bar{a}\mathbf{1})(\mathbf{a} - \bar{a}\mathbf{1})' + \Delta^2(\mathbf{a}^2\mathbf{1}' + \mathbf{1a}^2' - \bar{a}^2\mathbf{11}')/2 + o(\Delta^2)]/n,$$

where  $\bar{a} = \sum a_i/n$ ,  $\bar{a}^2 = \sum a_i^2/n$ .

Using (3.3), (3.8) and (3.9), and writing  $\tilde{x}_{\alpha} = \Gamma_{\alpha}(x - \mathbf{1}\bar{x}')$ ,

$$(3.10) \quad S_{\alpha}^* = S - \Delta \tilde{x}'_{\alpha} D(\mathbf{a}) \tilde{x}_{\alpha} + (\Delta^2/2) \tilde{x}'_{\alpha} D(\mathbf{a}^2) \tilde{x}_{\alpha} - \Delta^2 \tilde{x}'_{\alpha} \mathbf{a} \mathbf{a}' \tilde{x}_{\alpha} + \text{a remainder term } R, \text{ which for every fixed } x \text{ is } o(\Delta^2),$$

implying

$$(3.11) \quad |S_{\alpha}^*|/|S| = |I_p - \Delta S^{-1/2} \tilde{x}'_{\alpha} D(\mathbf{a}) \tilde{x}_{\alpha} S^{-1/2} + (\Delta^2/2) S^{-1/2} \tilde{x}'_{\alpha} D(\mathbf{a}^2) \tilde{x}_{\alpha} S^{-1/2} - \Delta^2 S^{-1/2} \tilde{x}'_{\alpha} \mathbf{a} \mathbf{a}' \tilde{x}_{\alpha} S^{-1/2} + \text{a remainder term, which uniformly in } x \text{ is } o(\Delta^2)|.$$

To evaluate this determinant (locally in  $\Delta$ ), we use the following result, whose proof is easy and hence omitted.

**LEMMA 3.2.**  $|I_p - \Delta B| = 1 - \Delta \text{tr } B + (\Delta^2/2)((\text{tr } B)^2 - \text{tr } B^2) + o(\Delta^2).$

Identifying  $B$  and applying this result, (3.11) is evaluated as

$$(3.12) \quad |S_{\alpha}^*|/|S| = 1 - \Delta \text{tr } D(\mathbf{a}) \tilde{x}_{\alpha} + (\Delta^2/2) \text{tr } D(\mathbf{a}^2) \tilde{x}_{\alpha} - \Delta^2 \tilde{x}'_{\alpha} \mathbf{a} \mathbf{a}' \tilde{x}_{\alpha} + (\Delta^2/2) \{ (\text{tr } D(\mathbf{a}) \tilde{x}_{\alpha})^2 - \text{tr } D(\mathbf{a}) \tilde{x}_{\alpha} D(\mathbf{a}) \tilde{x}_{\alpha} \} + R_{\alpha}(x, \Delta),$$

where  $\tilde{x}_\alpha = \tilde{x}_\alpha S^{-1} \tilde{x}'_\alpha$ , and  $\sup_x R_\alpha(x, \Delta) = o(\Delta^2)$  for all  $\Gamma_\alpha \in P$ . Finally, we need the following results for a complete (local) evaluation of (3.2). We write  $(\mathbf{x}_i - \bar{\mathbf{x}})' S^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) = b_{ii}$ ,  $i = 1, \dots, n$ .

$$(i) \quad \sum_\alpha \text{tr} D(\mathbf{a}) \tilde{x}_\alpha = \sum_\alpha \sum_i a_{v_i} b_{ii} = pn! \bar{a},$$

using Ferguson ((1961), page 258) and  $\sum_i b_{ii} = p$ .

$$(ii) \quad \sum_\alpha \text{tr} D(\mathbf{a}^2) \tilde{x}_\alpha = pn! \bar{a}^2, \quad \text{by (i).}$$

$$(iii) \quad \sum_\alpha \mathbf{a}' \tilde{x}_\alpha \mathbf{a} = \text{tr} \left[ S^{-1} \sum_{i,j} a_i a_j \Gamma_\alpha(\mathbf{x}_{v_i} - \bar{\mathbf{x}}) (\mathbf{x}_{v_j} - \bar{\mathbf{x}})' \right]$$

$$= \text{tr} \left[ S^{-1} \left\{ n(n-2)! \sum (a_i - \bar{a})^2 \sum (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})' \right. \right.$$

$$\left. \left. + (n-2)! \left( \sum (\mathbf{x}_i - \bar{\mathbf{x}}) \right) \left( \sum (\mathbf{x}_i - \bar{\mathbf{x}}) \right)' \left( \left( \sum a_i \right)^2 - \sum a_i^2 \right) \right\} \right],$$

(3.13) (using a straightforward generalization of Ferguson's result ((1961), page 258))

$$= pn(n-2)! \sum (a_i - \bar{a})^2.$$

$$(iv) \quad \sum_\alpha (\text{tr} D(\mathbf{a}) \tilde{x}_\alpha)^2 = \sum_\alpha \sum_i a_{v_i} b_{ii}^2 = n(n-2)! \sum_i b_{ii}^2 \sum_i (a_i - \bar{a})^2$$

$$+ (n-2)! p^2 \left( \left( \sum a_i \right)^2 - \sum a_i^2 \right),$$

using Ferguson ((1961), page 258).

$$(v) \quad \sum_\alpha \text{tr} D(\mathbf{a}) \tilde{x}_\alpha D(\mathbf{a}) \tilde{x}_\alpha = \sum_\alpha \text{tr} \left\{ \sum_i a_{v_i} S^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})' \right\}^2$$

$$= n(n-2)! \sum (a_i - \bar{a})^2 \sum b_{ii}^2$$

$$+ (n-2)! p \left( \left( \sum a_i \right)^2 - \sum a_i^2 \right).$$

Using (3.12) and (3.13) (i)-(v), it follows that

$$(3.14) \quad \sum_\alpha \{ |S_\alpha^*| / |S| \}^{-(n-1)/2}$$

$$= n! \left[ 1 + (n-1) p \bar{a} \Delta / 2 \right.$$

$$\left. + (\Delta^2 / 8) \left\{ (n+1) \sum (a_i - \bar{a})^2 \sum b_{ii}^2 + K_1 \right\} + R(x, \Delta) \right],$$

where  $K_1$  is a constant and  $\sup_x R(x, \Delta) = o(\Delta^2)$ . This leads to (see (3.2))

$$(3.15) \quad dP_\Delta^T / dP_0^T(t(x)) = 1 + (n+1) \frac{1}{8} \Delta^2 \sum (a_i - \bar{a})^2 \left( \sum b_{ii}^2 \right)$$

$$+ K_2 \Delta^2 + \bar{R}(x, \Delta),$$

where  $\sup_x \bar{R}(x, \Delta) = o(\Delta^2)$  and  $K_2$  is a constant.

Consider now an invariant test  $\psi(t)$  of size  $\alpha$ . Then its local power is evaluated as

$$(3.16) \quad \int \psi(t) dP_{\Delta}^T(t(x)) = \alpha + (n+1)\frac{1}{8}\Delta^2 \sum (a_i - \bar{a})^2 \\ \times \int \psi(t) \left( \sum b_{ii}^2 \right) dP_0^T(t(x)) + K_2 \Delta^2 \alpha + o(\Delta^2).$$

Our main result is the following.

**THEOREM 3.1.** *The locally best invariant test for  $H_0: \Delta = 0$  versus  $H_1: \Delta \neq 0$  under the model (2.1)–(2.3) rejects  $H_0$  for large values of  $\sum_1^n \{(\mathbf{x}_i - \bar{\mathbf{x}})' S^{-1}(\mathbf{x}_i - \bar{\mathbf{x}})\}^2$ , whatever  $a_1, \dots, a_n, \sum_1^n (a_i - \bar{a})^2 > 0$ .*

**PROOF.** An application of (3.16) and the Neyman–Pearson lemma completes the proof of the theorem.  $\square$

**REMARK 3.2.** That the test based on  $\sum_1^n b_{ii}^2$  is null robust follows from Kariya (1981) and Sinha (1984).

## REFERENCES

- DAWID, A. P. (1977). Spherical matrix distributions and a multivariate model. *J. Roy. Statist. Soc. Ser. B* **39** 254–261.
- FERGUSON, T. S. (1961). On the rejection of outliers. *Proc. Fourth Berkeley Symp. Math. Statist. Probab.* **1** 253–287. Univ. California Press.
- KARIYA, T. (1981). Robustness of multivariate tests. *Ann. Statist.* **9** 1267–1275.
- KARIYA, T. and SINHA, B. K. (1985). Nonnull and optimality robustness of some tests. *Ann. Statist.* **13** 1182–1197.
- SCHWAGER, S. J. (1979). Detection of multivariate outliers. Ph.D. thesis, Yale University.
- SCHWAGER, S. J. and MARGOLIN, B. H. (1982). Detection of multivariate normal outliers. *Ann. Statist.* **10** 943–954.
- SINHA, B. K. (1984). Detection of multivariate outliers in elliptically symmetric distributions. *Ann. Statist.* **12** 1558–1565.
- WIJSMAN, R. A. (1967). Cross-sections of orbits and their application to densities of maximal invariants. *Proc. Fifth Berkeley Symp. Math Statist. Probab.* **1** 389–400. Univ. California Press.

DEPARTMENT OF MATHEMATICS  
AND STATISTICS  
UNIVERSITY OF PITTSBURGH  
PITTSBURGH, PENNSYLVANIA 15260

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF MARYLAND BALTIMORE COUNTY  
CATONSVILLE, MARYLAND 21228