

UNBIASEDNESS OF A MULTIVARIATE OUTLIER TEST FOR ELLIPTICALLY CONTOURED DISTRIBUTIONS

BY JIAN-XIN PAN* AND XUE-REN WANG

Hong Kong Baptist College and Yunnan University

Under a class of elliptically contoured distributions, the likelihood ratio criterion (LRC) for detecting multiple outliers is established in this paper. It is shown that the LRC has the same form as in the normal case. Furthermore, the unbiasedness of the LRC is derived.

1. Introduction. Suppose that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ is a random sample from some p -dimensional distribution $F_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma} > \mathbf{0}$ are the mean vector and the covariance matrix, respectively, and both are unknown. We are often required to detect whether there are outliers in the sample and which sample points are outliers. This can be reduced to a testing hypothesis problem as follows. If we choose the null model in which there are no outliers in the sample

$$H : \mathbf{x}_i \sim F_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad i = 1, 2, \dots, n, \quad (1.1)$$

then a possible alternative model which may account for multiple outliers is the multivariate model with mean slippage

$$K : \begin{aligned} \mathbf{x}_i &\sim F_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) & (i \notin I) \\ \mathbf{x}_i &\sim F_p(\boldsymbol{\mu} + \mathbf{a}_i, \boldsymbol{\Sigma}) & (i \in I) \end{aligned} \quad (1.2)$$

where $I = \{i_1, i_2, \dots, i_k\}$ is an index subset of $\{1, 2, \dots, n\}$ with a fixed positive integer k and $\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_k}$ are unknown mean slippage parameters.

For hypotheses (1.1) and (1.2), Siotani (1959) and Wilks (1963) discussed the likelihood ratio criterion under the multivariate normal distribution $F_p = N_p$. In this paper, we extend their results to a class of elliptically contoured distributions. The null distribution of the likelihood ratio testing statistic for detecting the multiple outliers is shown to be a Wilk's distribution, which has the same distribution as in the normal case. Furthermore, the unbiasedness of the likelihood ratio test is derived. Under the assumption of an elliptically contoured distribution, Sinha (1984) provided a locally best invariant test for outliers based on multivariate sample kurtosis, which is useful for identifying

* Supported partially by the WAI TAK Investment and Loan Company LT.D. Research Scholarship of Hong Kong.

AMS 1991 Subject Classifications: Primary 62H15; Secondary 62H10.

Key words and phrases: Elliptically contoured distribution; likelihood ratio criterion; multiple outlier; unbiasedness of a test.

whether there are outliers in the sample. The likelihood ratio test for discordantly multiple outliers, however, can be used not only for determining whether there are outliers in the sample, but also for detecting which of the observations are discordant outliers.

2. Likelihood Ratio Criteria. Suppose that the characteristic function of a random matrix $\mathbf{X}_{n \times p} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^\tau$ has the form:

$$\exp\left\{i \sum_{j=1}^n \mathbf{t}_j^\tau \boldsymbol{\mu}_j\right\} \cdot \Phi\left(\sum_{j=1}^n \mathbf{t}_j^\tau \boldsymbol{\Sigma}_j \mathbf{t}_j\right), \tag{2.1}$$

then \mathbf{X} is known to have a matrix elliptically contoured distribution (See, e.g., Fang and Zhang, 1990). When there is a density function of \mathbf{X} , it must be of the form:

$$\prod_{j=1}^n |\boldsymbol{\Sigma}_j|^{-1/2} g\left(\sum_{j=1}^n \text{tr} \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_j)(\mathbf{x}_j - \boldsymbol{\mu}_j)^\tau\right), \tag{2.2}$$

where $g(\cdot)$ is a non-negative function and satisfies $\int_0^{+\infty} y^{\frac{np}{2}-1} g(y) dy < +\infty$. In this case we denote $\mathbf{X} \sim LEC_{n \times p}(\mathbf{M}; \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2, \dots, \boldsymbol{\Sigma}_n; g)$, where $\mathbf{M} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_n)^\tau$. Especially, if $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \dots = \boldsymbol{\mu}_n = \boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \dots = \boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}$ (say) we denote $\mathbf{X} \sim LEC_{n \times p}(\mathbf{1}_n \boldsymbol{\mu}^\tau; \mathbf{I}_n \otimes \boldsymbol{\Sigma}; g)$, where $\mathbf{1}_n = (1, \dots, 1)^\tau \in R^n$ and the notation \otimes represents the Kronecker product of two matrices. In this case, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, the rows of the matrix \mathbf{X} , are mutually uncorrelated and have a common mean vector $\boldsymbol{\mu}$ and a covariance matrix $\boldsymbol{\Sigma}$. Now, the tests (1.1)–(1.2) in Section 1 become

$$H : \mathbf{X} \sim LEC_{n \times p}(\mathbf{1}_n \boldsymbol{\mu}^\tau; \mathbf{I}_n \otimes \boldsymbol{\Sigma}; g) \tag{2.3}$$

and

$$K : \mathbf{X} \sim LEC_{n \times p}(\mathbf{1}_n \boldsymbol{\mu}^\tau + \mathbf{D}\mathbf{A}; \mathbf{I}_n \otimes \boldsymbol{\Sigma}; g) \tag{2.4}$$

where $\mathbf{D} = (\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_k})$ and $\mathbf{A} = (\mathbf{a}_{i_1}, \mathbf{a}_{i_2}, \dots, \mathbf{a}_{i_k})^\tau$ are $n \times k$ and $k \times p$ matrices, respectively, \mathbf{e}_i ($i \in I$) is a n -dimensional vector whose i -th element is one and the remaining elements are zeros, and \mathbf{a}_i ($i \in I$) is an unknown mean slippage p -variate vector parameter.

THEOREM 2.1. Suppose the function $g(y)$ in (2.2) is continuous and decreasing for $y \in [0, +\infty)$, and $n > p + k + 1$. Then the likelihood ratio criterion of the tests (2.3)–(2.4) is equivalent to rejecting H if

$$\Lambda_I = |\mathbf{S}_{(I)}|/|\mathbf{S}| \leq C_\alpha \tag{2.5}$$

where C_α is the lower $100\alpha\%$ critical point of the Wilks' distribution $\Lambda(p, n - k - 1, k)$, $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$, $\mathbf{S} = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\tau$, $\bar{\mathbf{x}}_{(I)} = \frac{1}{n-k} \sum_{i \notin I} \mathbf{x}_i$ and $\mathbf{S}_{(I)} = \sum_{i \notin I} (\mathbf{x}_i - \bar{\mathbf{x}}_{(I)})(\mathbf{x}_i - \bar{\mathbf{x}}_{(I)})^\tau$.

PROOF. It can be shown that the likelihood ratio testing statistic is equivalent to $\Lambda_I = |S_{(I)}|/|S|$. In the following only the null distribution of Λ_I is derived. Let $\Lambda_I(\mathbf{X})$ denote the statistic Λ_I for the data matrix \mathbf{X} . Under the null hypothesis H , it is obvious that $\Lambda_I(\mathbf{X} - \mathbf{1}_n\boldsymbol{\mu}^\tau) \equiv \Lambda_I(\mathbf{X})$ and $\Lambda_I(\alpha\mathbf{X}) \equiv \Lambda_I(\mathbf{X})$ are true for all $\boldsymbol{\mu} \in R^p$ and $\alpha > 0$, respectively. Therefore, the null density distribution of Λ_I does not depend on the choice of the function $g(y)$ (See, e.g., Fang and Zhang, 1990). In other words, it is distribution-free or distribution-robust on the class of elliptically contoured distributions. Especially, the null density distribution of Λ_I is the same that of Λ_I as $g(y) = (2\pi)^{-pn/2}e^{-y/2}$ ($y \geq 0$), which corresponds to the matrix normal distribution $\mathbf{X} \sim N_{n,p}(\mathbf{1}_n\boldsymbol{\mu}^\tau, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$. Noting that $\mathbf{S} = \mathbf{X}^\tau \mathbf{D}_n \mathbf{X}$ and $\mathbf{S}_{(I)} = \mathbf{X}^\tau \mathbf{C}_1 \mathbf{X}$ where $\mathbf{D}_n = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\tau$, $\mathbf{C}_1 = \text{diag}(\mathbf{D}_{n-k}, \mathbf{0})$ and $\mathbf{X} = (\mathbf{X}_{(I)}^\tau, \mathbf{X}_I^\tau)^\tau$, we know that $\mathbf{C}_1^2 = \mathbf{C}_1$, $\mathbf{C}_1^\tau = \mathbf{C}_1$, $rk(\mathbf{C}_1) = n - k - 1$ and $\mathbf{C}_1(\mathbf{1}_n\boldsymbol{\mu}^\tau) = \mathbf{0}$ so that $\mathbf{S}_{(I)} \sim W_p(n - k - 1, \boldsymbol{\Sigma})$ under the multivariate normal assumption. Let $\mathbf{C}_2 = \mathbf{D}_n - \mathbf{C}_1$ and $\mathbf{E}_{(I)} = \mathbf{S} - \mathbf{S}_{(I)}$ then $\mathbf{E}_{(I)} = \mathbf{X}^\tau \mathbf{C}_2 \mathbf{X}$. It can be verified that $\mathbf{C}_2^2 = \mathbf{C}_2$, $\mathbf{C}_2^\tau = \mathbf{C}_2$, $\mathbf{C}_1\mathbf{C}_2 = \mathbf{0}$, $rk(\mathbf{C}_2) = k$ and $\mathbf{C}_2(\mathbf{1}_n\boldsymbol{\mu}^\tau) = \mathbf{0}$. Therefore when $\mathbf{X} \sim N_{n,p}(\mathbf{1}_n\boldsymbol{\mu}^\tau, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$, $\mathbf{E}_{(I)} = \mathbf{S} - \mathbf{S}_{(I)}$ is independent of $\mathbf{S}_{(I)}$ and $\mathbf{E}_{(I)} \sim W_p(k, \boldsymbol{\Sigma})$. Furthermore,

$$\Lambda_I = |\mathbf{S}_{(I)}|/|\mathbf{S}_{(I)} + \mathbf{E}_{(I)}| \sim \Lambda(p, n - k - 1, k) \tag{2.6}$$

under the null hypothesis H . The proof is complete.

REMARK 2.1. When $I = \{i\}$ ($i = 1, 2, \dots, n$), the i -th observation is declared as a discordantly single outlier of size α if

$$T_i^2 \geq \frac{p(n - 1)C_\alpha^*}{n(n - p - 1) + npC_\alpha^*} \tag{2.7}$$

where $T_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}})^\tau \mathbf{S}^{-1}(\mathbf{x}_i - \bar{\mathbf{x}})$ and C_α^* is the upper 100 α % critical point of $F_{(p, n-p-1)}$. When $I = \{i, j\}$ ($i, j = 1, 2, \dots, n$), the (i, j) -th observation is declared as a discordant outlier pair of size α if

$$\sqrt{\Lambda_{ij}} \leq \frac{n - p - 3}{pC_{\alpha}^{**} + (n - p - 3)} \tag{2.8}$$

where C_{α}^{**} is the upper 100 α % critical point of $F_{(2p, 2(n-p-3))}$.

REMARK 2.2. In most practical problems, the index subset I can not be given in advance. Suppose k is fixed, to detect a set of k outliers, a reasonable testing statistic is $\Lambda_{\min}^k = \min_I\{\Lambda_I\}$ where I runs over all subsets including k indexes. Unfortunately, the exact null distribution of Λ_{\min}^k is unknown. In this case, Bonferroni's principle in multiple comparisons can be recommended.

THEOREM 2.2. Suppose the function $g(y)$ in (2.2) is continuous and decreasing for $y \in [0, +\infty)$, and $n > p + k + 1$. Then the likelihood ratio test (2.5) for detecting multiple outliers of a sample in elliptically contoured distribution is unbiased.

PROOF. Obviously, the power function of the test (2.5) $\Pr\{\Lambda_I \leq C_\alpha\}$ is a function of the mean shift parameter \mathbf{A} . It can be shown, however, $\Pr\{\Lambda_I \leq C_\alpha\}$ depends upon \mathbf{A} only through $\|\delta_i^*\| = \sqrt{\delta_i^{*\tau} \delta_i^*}$ ($i = 1, 2, \dots, k$) so that we denote it as $\beta(\|\delta_1^*\|, \|\delta_2^*\|, \dots, \|\delta_k^*\|)$, where $\delta_i^{*\tau}$ is the i -th row of $\Delta^* = \Gamma \mathbf{A} \Sigma^{-1/2}$ and Γ is a $k \times k$ orthogonal matrix with the first row $\mathbf{1}_k^T / \sqrt{k}$. Furthermore, $\beta(\|\delta_1^*\|, \|\delta_2^*\|, \dots, \|\delta_k^*\|)$ can be shown to be a monotonically increasing function of $\|\delta_i^*\|$ ($i = 1, 2, \dots, k$), thus $\beta(\|\delta_1^*\|, \|\delta_2^*\|, \dots, \|\delta_k^*\|) \geq \beta(0, 0, \dots, 0)$, which implies that the power function of the test (2.5) achieves its minimum at the null hypothesis $H : \mathbf{A} = \mathbf{0}$. In other words, the likelihood ratio test (2.5) is unbiased and the proof is complete.

Acknowledgement. The authors would like to thank Professor K.T. Fang and Dr. F.J. Hickernell for their valuable suggestions and discussions. Also, the referees' helpful comments are appreciated.

REFERENCES

- FANG, K. T. and ZHANG, Y. T. (1990). *Generalized Multivariate Analysis*. Springer-Verlag and Science Press, Berlin and Beijing.
- SINHA, B. K. (1984). Detection of multivariate outliers in elliptically symmetric distributions. *Ann. of Stat.* **12**, 4, 1558–1565.
- SIOTANI, M. (1959). The extreme value of the generalized distances of the individual points in the multivariate normal sample. *Ann. Inst. Stat. Math.*, Tokyo **10**, 10, 183–208.
- WILKS, S. S. (1963). Multivariate statistical outliers. *Sankhya A* **25**, 406–427.

DEPARTMENT OF MATHEMATICS
HONG KONG BAPTIST COLLEGE
224 WATERLOO ROAD, KOWLOON
HONG KONG

DEPARTMENT OF STATISTICS
YUNNAN UNIVERSITY
KUNMING 650091, CHINA