

On the Robustness of Chi-square Test Procedure for a Compounded Multivariate Normal Mean

Hea-Jung Kim¹⁾

Abstract

The robustness of one sample Chi-square test for multivariate normal mean vector is investigated when the multivariate normal population is mixed with another multivariate normal population with differing in the mean vector. Explicit expressions for the level of significance and power of the test are derived. Some numerical results indicate that the Chi-square test procedure is quite robust against slight mixtures of multivariate normal populations differing in location parameters.

1. Introduction

Consider a random sample of size N from the mixture of multivariate normal population with probability density function(pdf)

$$f(\mathbf{x}) = w\phi(\mathbf{x}; \mu_1, \Sigma) + (1-w)\phi(\mathbf{x}; \mu_2, \Sigma), \quad (1)$$

where $0 \leq w \leq 1$ and $\phi(\mathbf{x}; \mu, \Sigma)$ is the multivariate normal pdf with mean vector μ and known variance covariance matrix Σ . Day(1969), Wolfe(1970) and Johnson and Kotz(1972) have studied the distributions of some statistics derived from (1) and its general form.

The purpose of the present study is to consider the effects of nonnormality of the type (1) on familiar one sample mean test(Chi-square test) of significance when Σ is known, i.e. the robustness of the one sample Chi-square test procedure against slight contamination of the population with another multivariate normal population having a different mean vector. The reader is referred to Subrahmaniam, et al.(1975) and Blumenthal and Govindarajulu(1977) for univariate robustness studies of this kind.

2. Distribution of the Chi-square test statistic

Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ be a sequence of iid p -dimensional random vectors with common density function

1) Department of Statistics, Dongguk University, Seoul 100-715, KOREA.

$$f(\mathbf{x}; w, \mu_1, \mu_2, \Sigma) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \left\{ w \exp\left[-\frac{1}{2} (\mathbf{x} - \mu_1)' \Sigma^{-1} (\mathbf{x} - \mu_1)\right] + (1-w) \exp\left[-\frac{1}{2} (\mathbf{x} - \mu_2)' \Sigma^{-1} (\mathbf{x} - \mu_2)\right] \right\}, \tag{2}$$

where $0 \leq w \leq 1$ and $\Sigma > 0$ is known. The common expected value is

$$\theta = E(\mathbf{X}) = w\mu_1 + (1-w)\mu_2 \tag{3}$$

and the variance covariance matrix is

$$\Omega = V(\mathbf{X}) = \Sigma + w(1-w)(\mu_1 - \mu_2)(\mu_1 - \mu_2)' . \tag{4}$$

To test the hypothesis of $H_0 : \theta = \theta_0$ vs $H_1 : \theta \neq \theta_0$, assuming the normality of \mathbf{X}_i 's, the usual Chi-square test procedure with the critical region $T^2 > \chi_{\alpha, p}^2$ is used where α is the given type I error probability, $\chi_{\alpha, p}^2$ is the $100(1-\alpha)$ th percentile of the χ^2 distribution with p degrees of freedom and T^2 is defined by

$$T^2 = N(\bar{\mathbf{X}} - \theta_0)' \Sigma^{-1} (\bar{\mathbf{X}} - \theta_0), \tag{5}$$

where $\bar{\mathbf{X}} = \frac{1}{N} \sum_{i=1}^N \mathbf{X}_i$. For obtaining the distribution of T^2 under the distribution (2), it will be convenient to think of the sample as coming from two multivariate normal populations; $\Pi_1 \sim N_p(\mu_1, \Sigma)$ and $\Pi_2 \sim N_p(\mu_2, \Sigma)$. With probability w an observation comes from Π_1 and with probability $(1-w)$ it comes from Π_2 . We let R denote the random unobservable number of observations among $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ which come from Π_1 . Then we have the following result.

Lemma 1. Given $R = r$,

$$\Pr(T^2 \leq t | R) = e^{-\lambda/2} \sum_{j=0}^{\infty} \frac{\lambda/2^j}{j!} \frac{1}{2^{p/2+j} \Gamma(p/2+j)} \int_0^t y^{p/2+j-1} e^{-y/2} dy, \tag{6}$$

where

$$\lambda = \frac{1}{N} (N\mu_2 + r(\mu_1 - \mu_2) - N\theta_0)' \Sigma^{-1} (N\mu_2 + r(\mu_1 - \mu_2) - N\theta_0).$$

Proof. Given $R = r$, $\bar{\mathbf{X}}$ is distributed as

$$\bar{\mathbf{X}} \sim N_p\left(\frac{1}{N} (r\mu_1 + (N-r)\mu_2), \frac{1}{N} \Sigma\right)$$

so that $N^{1/2} \Sigma^{-1/2} (\bar{\mathbf{X}} - \theta_0) \sim N_p(N^{-1/2} \Sigma^{-1/2} \{r\mu_1 + (N-r)\mu_2 - N\theta_0\}, I_p)$. Therefore, the conditional distribution of T^2 is a noncentral Chi-square with p degrees of freedom and noncentrality parameter λ . Using the cumulative distribution function of the noncentral Chi-square(cf. Johnson and Kotz, 1972), we have the result.

Theorem 2. The unconditional cumulative distribution function of T^2 is given by

$$\Pr(T^2 \leq t) = \sum_{j=0}^{\infty} \sum_{r=0}^N \binom{N}{r} w^r (1-w)^{N-r} \left\{ \frac{(\lambda/2)^j}{j!} e^{-\lambda/2} \right\} \Pr(\chi^2_{p+2j} \leq t), \quad (7)$$

where χ^2_{p+2j} denotes the central Chi-square variate with $p+2j$ degrees of freedom.

Proof. $\Pr(T^2 \leq t) = E \Pr(T^2 \leq t | R)$ with the binomial distribution R used for the expectation. Using the binomial distribution of R , we have

$$E \Pr(T^2 \leq t | R) = \sum_{r=0}^N \binom{N}{r} w^r (1-w)^{N-r} \Pr(T^2 \leq t | R).$$

Expressing $\Pr(T^2 \leq t | R)$ in Lemma 1 as a weighted sum of central Chi-square distribution probabilities with weights equal to the probabilities of a Poisson distribution with expected value $\lambda/2$, we have the result.

As a result, when we test the common mean vector of (2) via the Chi-square test procedure, Theorem 2 yields the true level of significance of the test. The true level of significance α^* , for the nominal value α , is given by

$$\alpha^* = 1 - \Pr(T^2 \leq \chi^2_{\alpha, p}). \quad (8)$$

Here $\chi^2_{\alpha, p}$ is the critical value of the usual Chi-square size α test with p degrees of freedom. Thus one can use (8) to evaluate either size or power of the test based on T^2 by appropriate choice of the value of λ in (7).

3. Numerical results and Conclusion

Robustness of the Chi-square test for the compounded multivariate normal population mean is investigated. This is done by calculating the true size α^* of the test $H_0: \theta = \theta_0$ and then by comparing it with the nominal value α in (8). In order to obtain the true size (level of significance) α^* of the Chi-square test procedure for testing $H_0: \theta = \theta_0$, without loss of generality, we set $\theta_0 = 0$ in (5) so that $T^2 = N \bar{\mathbf{X}}' \Sigma^{-1} \bar{\mathbf{X}}$. In this case, the noncentrality parameter for the distribution of T^2 given $R = r$ is

$$\lambda = \frac{1}{N} (N \mu_2 + r(\mu_1 - \mu_2))' \Sigma^{-1} (N \mu_2 + r(\mu_1 - \mu_2)).$$

The results of our study are given in Table 1. An examination of this table shows that the Chi-square test is robust against small contamination depicted by (1). In particular we have the following comments : (i) In either of the cases where the proportion of contamination is large ($w < .90$) or the distance(Mahalanobis distance) between the populations is large, the test would be inappropriate. In other words, if the departure from multivariate normality is greatly accentuated, the Chi-square test does not give the desired protection. (ii) The increase in the nominal value α and the dimension p , however, do not affect the robustness of the test procedure (iii) An increase in N , the sample size, can lead to a slight worsening of the situation. This is due to the fact that the distribution of T^2 is noncentral Chi-square with the noncentrality parameter λ being increasingly affected by N .

In this paper we have assumed that the variance covariance matrix Σ in (1) is known. In case where Σ is unknown, the effects of nonnormality of the type (1) on the usual mean test (Hotelling's T^2 test) may be another problem to be studied. It is left as a future research topic of interest.

References

- [1] Day, N. E.(1969). Estimating the Components of a Mixture of Normal Distributions, *Biometrika*, Vol. 56, 463-474.
- [2] Blumenthal, S. and Govindarajulu, Z.(1977). Robustness of Stein's Two-Stage Procedure for Mixtures of Normal Populations, *Journal of the American Statistical Association*, Vol. 72, 192-196.
- [3] Johnson, N. L. and Kotz, S.(1972). *Distributions in Statistics : Continuous Multivariate Distributions*, John Wiley & Sons, Inc., New York.
- [4] Subrahmaniam, K, Subrahmaniam, K and Messeri, J. Y.(1975). On the Robustness of Some Tests of Significance in Sampling from a Compound Normal Population, *Journal of the American Statistical Association*, Vol. 70, 435-438.
- [5] Wolfe, J. H.(1970). Pattern Clustering by Multivariate Mixture Analysis, *Multivariate Behavioral Research*, Vol. 5, 329-350.