

A c -SAMPLE NON-PARAMETRIC TEST FOR LOCATION IN A MIXED MODEL OF CONTINUOUS AND DISCRETE VARIABLES*

BY SHASHIKALA SUKHATME

New Delhi

1. INTRODUCTION

LET $Z_i = (Y_i, X_{1i}, X_{2i}, \dots, X_{ci})$ $i = 1, 2, \dots, N$ be N independent observations from a $(c + 1)$ variate distribution where for each i , $X_{ji} = 0$ or 1 , $\sum_{j=1}^c X_{ji} = 1$, $P\{X_{ji} = 1\} = p_j$, $P\{X_{ji} = 0\} = q_j = 1 - p_j$,

$\sum_{j=1}^c p_j = 1$, and $P\{Y \leq y | X_j = 1\} = F_j(y)$, $j = 1, 2, \dots, c$. The distribution functions F_1, \dots, F_c are assumed to be absolutely continuous. In this paper we propose a median test for testing the hypothesis $H_0: F_1 = \dots = F_c$. For this purpose, divide the observations Y_1, \dots, Y_N into c sets according as $X_{ji} = 1$, $j = 1, 2, \dots, c$.

Let U_{j1}, \dots, U_{jn_j} ($n_j > 0$ for each j , $\sum_{j=1}^c n_j = N$) denote those Y_j 's for which the corresponding $X_{ji} = 1$. For given n_1, \dots, n_c the problem of testing the hypothesis H_0 reduces to testing the hypothesis that the c independent samples of U_{ji} 's ($i = 1, 2, \dots, n_j$; $j = 1, 2, \dots, c$) come from the same distribution. However, the problem under consideration differs from the usual c -sample problem in that the sample sizes n_1, \dots, n_c are random variables having a multinomial distribution with parameters p_1, \dots, p_c .

We assume that F_j 's differ only in location. Let $F_j(y) = F(y + \theta_j)$, $j = 1, 2, \dots, c$ for some arbitrary choice of real numbers $\theta_1, \dots, \theta_c$. Further we denote by H_N the hypothesis which specifies that $F_j(y) = F(y + \theta_j/\sqrt{N})$, $j = 1, 2, \dots, c$ and for some pair (i, j) $\theta_i \neq \theta_j$.

Let \tilde{W} denote the sample median of Y observations and m_j the number of U_{ji} 's ($i = 1, 2, \dots, n_j$) that are less than \tilde{W} . Assume

* This research was performed while the author was a graduate student at Michigan State University, East Lansing, Michigan, and was sponsored in part by the Office of Ordnance Research,

$N = 2k + 1$. Clearly $\sum_{j=1}^c m_j = k$. The test statistic proposed for testing the hypothesis $H: F_1 = \dots = F_c$ is then defined as

$$M = \sum_{j=1}^c \left(\frac{m_j - kp_j}{\sqrt{Np_j}} \right)^2, \tag{1.1}$$

when p_1, \dots, p_c are known, and as

$$\hat{M} = \left(\frac{m_j - k\hat{p}_j}{\sqrt{N\hat{p}_j}} \right)^2, \tag{1.2}$$

where $\hat{p}_j = n_j/N$, when p_1, \dots, p_c are unknown. The test consists in rejecting the hypothesis if $M(\hat{M})$ is large.

In Section 2 we find the joint distribution of m_1, \dots, m_c and \tilde{W} and in Section 3 the limiting distribution of M . In Section 4 the relative asymptotic efficiency of the median test based on M with respect to a corresponding parametric test based on multiple correlation coefficient is evaluated. Section 5 deals with the case when p_1, \dots, p_c are unknown and gives the asymptotic distribution of \hat{M} under the hypothesis H_c , from which we conclude that the test based on \hat{M} is asymptotically distribution-free.

2. JOINT DISTRIBUTION OF m_1, \dots, m_c AND \tilde{W}

Henceforth $f(\cdot)$ denotes the probability density function of the random variables written in the parentheses.

LEMMA 2.1. The joint distribution of m_1, \dots, m_c and \tilde{W} is given by

$$\begin{aligned} f(m_1, \dots, m_c, \tilde{w}) &= \frac{N!}{k! \prod_{j=1}^c m_j!} \prod_{j=1}^c [p_j F_j(\tilde{w})]^{m_j} \\ &\quad \times \left[1 - \sum_{j=1}^c p_j F_j(\tilde{w}) \right]^k \left[\sum_{j=1}^c p_j F_j'(\tilde{w}) \right], \end{aligned} \tag{2.1}$$

where m_1, \dots, m_c is a partition of k , $\sum_{j=1}^c m_j = k$.

Proof.—Noting that the conditional probability density of m_1, \dots, m_c and \tilde{W} for fixed values of n_1, \dots, n_c is given by

$$\begin{aligned}
 & f(m_1, \dots, m_c, \tilde{w} | n_1, \dots, n_c) \\
 &= \left[\sum_{j=1}^c \frac{(n_j - m_j)}{[1 - F_j(\tilde{w})]} F_j'(\tilde{w}) \right] \\
 & \quad \times \left[\prod_{j=1}^c \binom{n_j}{m_j} (F_j(\tilde{w}))^{m_j} (1 - F_j(\tilde{w}))^{n_j - m_j} \right] \quad (2.2)
 \end{aligned}$$

and that n_1, \dots, n_c have the multinomial distribution $m(N; p_1, \dots, p_c)$ given by

$$f(n_1, \dots, n_c) = \frac{N!}{\prod_{j=1}^c n_j!} \prod_{j=1}^c p_j^{n_j}, \quad (2.3)$$

we obtain by using

$$\begin{aligned}
 & f(m_1, \dots, m_c, \tilde{w}) \\
 &= \sum_{n_1, \dots, n_c} f(m_1, \dots, m_c, \tilde{w} | n_1, \dots, n_c) f(n_1, \dots, n_c)
 \end{aligned}$$

the required joint probability density given by (2.1).

Summing (2.1) over m_1, \dots, m_c we obtain the marginal distribution of \tilde{W} ,

$$\begin{aligned}
 f(\tilde{w}) &= \frac{N!}{k! k!} \left[\sum_{j=1}^c p_j F_j(\tilde{w}) \right]^k \left[1 - \sum_{j=1}^c p_j F_j(\tilde{w}) \right]^k \\
 & \quad \times \left[\sum_{j=1}^c p_j F_j'(\tilde{w}) \right].
 \end{aligned}$$

Under $H_c: F_1 = F_2 = \dots = F_c$, integration over the domain $0 \leq F_j(\tilde{w}) \leq 1$ in (2.1) yields the distribution of m_1, \dots, m_c as

$$f(m_1, \dots, m_c) = \frac{k!}{\prod_{j=1}^c m_j!} \prod_{j=1}^c p_j^{m_j}$$

which is multinomial distribution $m(k; p_1, \dots, p_c)$.

Also note that (m_1, \dots, m_c) and \tilde{W} are independent under the hypothesis H_c .

3. ASYMPTOTIC DISTRIBUTION OF M

We first prove the following lemma which gives the joint limiting distribution of m_1, \dots, m_c and \tilde{W} .

LEMMA 3.1. Let

$$v_j = \frac{m_j - Np_j F_j(\xi)}{\sqrt{Np_j F_j(\xi)}}, \quad j = 1, 2, \dots, c; \quad \eta = \sqrt{N}(\tilde{w} - \xi),$$

where ξ is such that

$$\sum_{j=1}^c p_j F_j(\xi) = \frac{1}{2}. \tag{3.1}$$

Assume that in some neighbourhood of ξ the density function $F_j'(y) = f_j(y)$ ($j = 1, 2, \dots, c$) has a continuous derivative. Then the asymptotic joint distribution of v_1, \dots, v_{c-1} and η is c -variate normal distribution with zero mean vector and covariance matrix Σ given by $\Sigma^{-1} = \Lambda = (\lambda_{ij})$ where

$$\lambda_{ii} = 1 + \frac{p_i F_i(\xi)}{p_c F_c(\xi)}, \quad i = 1, 2, \dots, (c-1);$$

$$\lambda_{cc} = \sum_{i=1}^c \frac{p_i f_i^2(\xi)}{F_i(\xi)} + 2 \left[\sum_{i=1}^c p_i f_i(\xi) \right]^2;$$

$$\lambda_{ij} = \frac{[p_i p_j F_i(\xi) F_j(\xi)]^{\frac{1}{2}}}{p_c F_c(\xi)}, \quad i \neq j = 1, 2, \dots, (c-1).$$

$$\lambda_{ic} = f_i(\xi) \sqrt{\frac{p_i}{F_i(\xi)}} - \frac{f_c(\xi)}{F_c(\xi)} \sqrt{p_i F_i(\xi)}, \quad i = 1, 2, \dots, (c-1).$$

Proof.—Throughout this proof for convenience set $F_i = F_i(\xi)$ and $f_i = f_i(\xi)$. Using Taylor's expansion

$$F_j(\tilde{w}) = F_j\left(\xi + \frac{\eta}{\sqrt{N}}\right) = F_j + \frac{\eta}{\sqrt{N}} f_j + o\left(\frac{\eta^2}{N}\right),$$

$$j = 1, 2, \dots, c;$$

$$1 - \sum_{j=1}^c p_j F_j(\tilde{w}) = \frac{1}{2} - \frac{\eta}{\sqrt{N}} \sum_{j=1}^c p_j f_j + o\left(\frac{\eta^2}{N}\right).$$

and substituting these in (2.1) we get

$$\begin{aligned}
 & f(m_1, \dots, m_c, \bar{w}) \\
 &= \left\{ \frac{N(2k)!}{k! k! 2^{2k}} \right\} \left\{ \frac{k!}{\prod_{i=1}^c m_i!} \prod_{j=1}^c (2p_j F_j)^{m_j} \right\} \\
 & \quad \times \left\{ \left(\prod_{i=1}^c \left[1 + \frac{\eta}{\sqrt{N}} \frac{f_i}{F_i} + o\left(\frac{\eta^2}{N}\right) \right]^{m_i} \right) \right. \\
 & \quad \times \left. \left(1 - \frac{-2\eta}{\sqrt{N}} \sum_{i=1}^c p_i f_i - o\left(\frac{\eta^2}{N}\right) \right)^k \right\} \\
 & \quad \times \left\{ \sum_{j=1}^c p_j f_j \left(\xi + \frac{\eta}{\sqrt{N}} \right) \right\} \\
 &= \{A_1\} \{A_2\} \{A_3\} \{A_4\}. \tag{3.2}
 \end{aligned}$$

Note that v_j 's satisfy the relation

$$\sum_{j=1}^c v_j (p_j F_j)^{\frac{1}{2}} = 0 \tag{3.3}$$

Now consider the region S defined by

$$\begin{aligned}
 S = \{ & (v_1, \dots, v_{c-1}, \eta) : a_1 \leq v_1 \leq b_1, a_2 \leq v_2 \leq b_2, \dots, \\
 & a_c \leq \eta \leq b_c \}.
 \end{aligned}$$

Using Stirling's approximation for $n!$

$$A_1 \sim \frac{N}{\sqrt{k\pi}}.$$

A_2 is independent of η and because of convergence of multinomial distribution to normal distribution, uniformly in S

$$\begin{aligned}
 A_2 \sim & [(2\pi)^{(c-1)} (2p_c F_c) \prod_{j=1}^c (2kp_j F_j)]^{-\frac{1}{2}} \\
 & \times \exp \left\{ -\frac{1}{2} \left[\sum_{i=1}^{c-1} v_i^2 \left(1 + \frac{p_i F_i}{p_c F_c} \right) + \sum_{i \neq j=1}^{c-1} v_i v_j \frac{\sqrt{p_i p_j F_i F_j}}{p_c F_c} \right] \right\}
 \end{aligned}$$

Using series expansion for $\log(1+x)$, uniformly in S

$$\begin{aligned} \log A_3 = & -\frac{\eta^2}{2} \left[\sum_{i=1}^c \frac{p_i f_i^2}{F_i} + 2 \left(\sum_{i=1}^c p_i f_i \right)^2 \right] \\ & + \sum_{i=1}^c \eta v_i f_i \left(\frac{p_i}{F_i} \right)^{\frac{1}{2}} + o(1). \end{aligned}$$

Using continuity of f_i , we have, uniformly in S ,

$$\begin{aligned} & f(m_2, \dots, m_c, \bar{W}) \\ & \sim N \left(\sum_{j=1}^c p_j F_j \right) \left[k\pi (2\pi)^{(c-1)} (2p_c F_c) \prod_{j=1}^c (2kp_j F_j) \right]^{-\frac{1}{2}} \\ & \times \exp. -\frac{1}{2} \left[\sum_{i=1}^{c-1} v_i^2 \left(1 + \frac{p_i F_i}{p_c F_c} \right) \right. \\ & + \eta^2 \left\{ \sum_{i=1}^{c-1} \frac{p_i f_i^2}{F_i} + 2 \left(\sum_{i=1}^c p_i f_i \right)^2 \right\} \\ & + \sum_{i \neq j=1}^{c-1} v_i v_j \frac{\sqrt{p_i p_j F_i F_j}}{p_c F_c} \\ & \left. - 2\eta \sum_{i=1}^{c-1} v_i \left\{ f_i \left(\frac{p_i}{F_i} \right)^{\frac{1}{2}} - \frac{f_c}{F_c} (p_i F_i)^{\frac{1}{2}} \right\} \right]. \end{aligned}$$

Now making the transformation $(m_j, \dots, m_c, \bar{W}) \rightarrow (v_1, \dots, v_{c-1}, \eta)$ it is seen that

$$\begin{aligned} & \lim_{N \rightarrow \infty} P \{ a_1 \leq v_1 \leq b_1, \dots, a_c \leq \eta \leq b_c \} \\ & = \int_{a_1}^{b_1} \int_{a_2}^{b_2} \dots \int_{a_c}^{b_c} f(v_1, \dots, v_{c-1}, \eta) dv_1 dv_2 \dots d\eta \end{aligned}$$

where $f(v_1, \dots, v_{c-1}, \eta)$ is the probability density function of normal distribution described in the present theorem.

The following lemma gives the asymptotic joint distribution of v_1, \dots, v_{c-1} and η under the hypothesis H_N which specifies that $F_j(y) = F(y + \theta_j/\sqrt{N})$, $j = 1, 2, \dots, c$.

LEMMA 3.2. Under the hypothesis H_N the asymptotic joint distribution of $(v_1, \dots, v_{c-1}, \eta)$ is the c -variate normal distribution given by

$$\begin{aligned}
 & f(v_1, \dots, v_{c-1}, \eta) \\
 &= \frac{f(\xi)}{\pi^{c/2} 2^{(c-2)/2} p_c^{1/2}} \times \exp. - \frac{1}{2} \left[4\eta^2 f^2(\xi) \right. \\
 & \quad \left. + \sum_{i=1}^{c-1} v_i^2 \left(1 + \frac{p_i}{p_c} \right) + \sum_{i \neq j=1}^{c-1} v_i v_j \frac{\sqrt{p_i p_j}}{p_c} \right].
 \end{aligned}$$

Proof.—It is similar to that of Lemma 3.1.

Now we are in a position to obtain the limiting distribution of M defined by (1.1) under the hypothesis H_N .

THEOREM 3.1. Under the hypothesis H_N the asymptotic distribution of $2M$ is non-central χ^2 with $(c - 1)$ degrees of freedom and non-centrality parameter

$$\lambda = 2 [F'(\xi)]^2 \sum_{j=1}^c p_j (\theta_j - \bar{\theta})^2, \tag{3.4}$$

where

$$\bar{\theta} = \sum_{j=1}^c p_j \theta_j.$$

Proof.—Write

$$\begin{aligned}
 u_i &= \frac{m_i - kp_i}{\sqrt{Np_i}} = \frac{[m_i - Np_i F_i(\xi)] \sqrt{F_i(\xi)}}{\sqrt{Np_i F_i(\xi)}} \\
 & \quad - \frac{kp_i - Np_i F_i(\xi)}{\sqrt{Np_i}},
 \end{aligned}$$

$i = 1, 2, \dots, c.$

Under the hypothesis H_N using Lemma 3.2 it follows that the asymptotic joint distribution of (u_1, \dots, u_c) is c -variate normal with means $\mu_i = \theta_i F'(\xi) \sqrt{p_i}$, and covariance matrix $\Sigma = (\sigma_{ij})$ of rank $(c - 1)$ where $\sigma_{ii} = (1 - p_i)/2$, $i = 1, 2, \dots, c$, $\sigma_{ij} = -\sqrt{p_i p_j}/2$, $i \neq j = 1, 2, \dots, c.$

Hence noting that $\sum_{j=1}^c \sqrt{p_j} u_j = 0$ it follows that the limiting distribution of

$$2M = 2 \left[\sum_{i=1}^{c-1} u_i^2 \left(1 + \frac{p_i}{p_c} \right) + \sum_{i \neq j=1}^c u_i u_j \frac{\sqrt{p_i p_j}}{p_c} \right]$$

is $\chi^2_{c-1}(\lambda)$, where λ is given by (3.4).

4. ASYMPTOTIC EFFICIENCY

Let $F_j(y) = F(y + \theta_j)$, then H_c is true when $\theta_j = 0$. We now find the relative asymptotic efficiency of the c -sample median test with respect to the corresponding parametric test, when F_j ($j = 1, 2, \dots, c$) is a normal distribution with mean μ_j and variance σ^2 . The hypothesis H_c is true if and only if $\rho^2_{Y(X_1, \dots, X_c)} = 0$, (Olkin and Tate¹) where $\rho_{Y(X_1, \dots, X_c)}$ is the multiple correlation coefficient between Y and X . Let R denote the sample multiple correlation coefficient between Y and X . If

$$\bar{U}_{..} = \frac{(\sum_{i,j} U_{ij})}{N}, \quad \bar{U}_j = \frac{(\sum_{i=1}^{n_j} U_{ji})}{n_j},$$

then

$$T^2 = \frac{R^2}{1 - R^2} = \frac{\sum_{j=1}^c n_j (\bar{U}_j - \bar{U}_{..})^2}{\sum_{j,i} (U_{ji} - \bar{U}_{..})^2 - \sum_{j=1}^c (\bar{U}_j - \bar{U}_{..})^2}$$

Also

$$\rho^2_{Y(X_1, \dots, X_c)} = \frac{\sum_{j=1}^c \frac{[(\mu_j - \bar{\mu})^2 p_j]}{\sigma^2}}{1 + \sum_{j=1}^c \frac{[(\mu_j - \bar{\mu})^2 p_j]}{\sigma^2}}$$

where

$$\bar{\mu} = \sum_{j=1}^c p_j \mu_j$$

Following Fisher² it is seen that under the hypothesis H_N the asymptotic distribution of $(N - c) T^2/(c - 1)$ is $\chi^2_{(c-1)}(\lambda')$ where the non-centrality parameter is given by

$$\lambda' = \sum_{j=1}^c p_j \frac{(\mu_j - \bar{\mu})^2}{\sigma^2}.$$

Also it is proved that the limiting distribution of $2M$ is $\chi^2_{c-1}(\lambda)$, where λ is given by

$$\lambda = 2 [F'(\xi)]^2 \sum_{j=1}^c p_j \frac{(\mu_j - \bar{\mu})^2}{\sigma^2}.$$

Since the two test statistics are asymptotically distributed as a non-central χ^2 with the same number of degrees of freedom, following Andrews³ and Hannan⁴ it is seen that the asymptotic efficiency is given by the ratio of the two non-centrality parameters. Hence the asymptotic relative efficiency is found to be

$$e(M, R) = 2\sigma^2 [F'(\xi)]^2 = \frac{1}{\pi}.$$

5. CASE WHEN p_1, \dots, p_c ARE UNKNOWN

In this case we estimate p_j by $\hat{p}_j = n_j/N, j = 1, 2, \dots, c$ and consider the test based on \hat{M} defined by (1.2). It is interesting to note that the test of H_c based on \hat{M} is asymptotically distribution-free, which is seen from Theorem 5.1.

THEOREM 5.1. Under the hypothesis $H_c, 4\hat{M}$ is asymptotically distributed as a χ^2 variable with $(c - 1)$ degrees of freedom.

Proof.—Write

$$v_j = \frac{m_j - k\hat{p}_j}{\sqrt{N\hat{p}_j}} = \left(\frac{p_j}{\hat{p}_j}\right)^{\frac{1}{2}} \frac{(m_j - k\hat{p}_j)}{\sqrt{Np_j}} = \left(\frac{p_j}{\hat{p}_j}\right)^{\frac{1}{2}} w_j,$$

where

$$w_j = \frac{m_j - kp_j}{\sqrt{Np_j}} - \frac{k(\hat{p}_j - p_j)}{\sqrt{Np_j}}.$$

Let $v = (v_1, \dots, v_c)$ and $w = (w_1, \dots, w_c)$, then $v = wD$, where D is a diagonal matrix with $\sqrt{p_j/\hat{p}_j}$ as its diagonal elements. Since $\text{plim}_{N \rightarrow \infty} \hat{p}_j = p_j$, it follows that $\text{plim}_{N \rightarrow \infty} (\sqrt{p_j/\hat{p}_j}) = 1$ and hence the matrix

D converges in probability (element-wise) to identity matrix. An application of lemma of [5, Lemma 1] yields that the vectors *v* and *w* have the same limiting distribution. Also it can be proved that the asymptotic distribution of *w* is *c*-variate normal with zero mean vector and covariance matrix $\Sigma = (\sigma_{ij})$ of rank $(c - 1)$ with $\sigma_{jj} = (1 - p_j)/4$, $j = 1, 2, \dots, c$ and $\sigma_{ij} = -\sqrt{p_i p_j}/4$, $i \neq j, 1, 2, \dots, c$. Noting that $\sum_{j=1}^c \sqrt{p_j} v_j$ converges in probability to zero as $N \rightarrow \infty$, the asymptotic distribution of v_1, \dots, v_{c-1} is given by

$$f(v_1, \dots, v_{c-1}) = \frac{1}{(2\pi)^{(c-1)/2} \left(\frac{1}{4}\right)^{(c-1)/2} p_c^{1/2}} \times \exp. - 2 \left[\sum_{i=1}^{c-1} v_i^2 \left(1 + \frac{p_i}{p_c}\right) + \sum_{i \neq j=1}^{c-1} v_i v_j \frac{\sqrt{p_i p_j}}{p_c} \right].$$

Hence

$$4\hat{M} = 4 \left[\sum_{i=1}^{c-1} v_i^2 \left(1 + \frac{\hat{p}_i}{\hat{p}_c}\right) + \sum_{i \neq j=1}^{c-1} v_i v_j \frac{\sqrt{\hat{p}_i \hat{p}_j}}{\hat{p}_c} \right]$$

has the asymptotic distribution stated in the theorem.

6. ACKNOWLEDGEMENT

It is a pleasure to record my sincere thanks to Professor Ingram Olkin for suggesting the problem and for his keen interest in its solution.

7. REFERENCES

1. Olkin, I. and Tate, R. F. . . . "Multivariate correlation models with mixed discrete and continuous variables," *Ann. of Math. Stat.*, 1961, 27, 446-465.
2. Fisher, R. A. . . . "The general sampling distribution of the multiple correlation coefficient," *Proc. Roy. Soc. Lond.*, 1928, 7 A (121), 654-73.
3. Fred C. Andrews . . . "Asymptotic behaviour of some rank tests for analysis of variance," *Ann. of Math. Stat.*, 1954, 25, 724-36.
4. Hannan, E. J. . . . "The asymptotic power of certain tests based on multiple correlation coefficient," *Journal Royal Stat. Soc., London*, 1956, 18 B, 227-33.
5. Chin Long Chiang . . . "On regular best asymptotically normal estimates," *Ann. of Math. Stat.*, 1956, 27, 336-51.