









Du

407 183

Reproduction in whole or part is permitted for any purposes of the U. S. Government. Distribution of this document is unlimited.

> DISTRIBUTION STATEMENT A Approved for public release; Distribution Unlimited

## CHI-SQUARE TESTS FOR THE MULTINOMIAL DISTRIBUTION

Khursheed Alam

Clemson University

A simple proof of the asymptotic property of Chi-square tests, commonly used in the analysis of categorical data, is given for use as a note for instruction to first-year graduate students.

> Key Words: Multinomial Distribution; Chi-square Tests; Contigency Tables.

## 1. Introduction

The Chi-square tests associated with the multinomial distribution are commonly used in the analysis of categorical data with reference to problems of specification, homogeneity of parallel samples, independence of attributes, etc. The asymptotic property of the tests, that is, the Chi-square distribution of the test statistics in large samples is generally known. However, it has been our observation that many applied statisticians tacitly accept the asymptotic result without satisfying themselves with its proof. It is also true that nearly all the text books in use on elementary and higher statistics either omit the proof or barely sketch it. In this paper we outline a fairly simple proof of the fundamental result, for use as a note for the instruction to first-year graduate students and as a needed theory for the frequent application of Chi-square tests by applied statisticians. The proof is essentially based on the contents of Chapters 5 and 6 of Rao (1966).

Consider a multinomial distribution M(n,p) with K cells, where  $p = (p_1, \ldots, p_K)'$  denotes the vector of cell probabilities,  $n = (n_1, \ldots, n_K)'$  denotes the vector of cell frequencies resulting from n independent trials,  $\sum_{i=1}^{K} p_i = 1$  and  $\sum_{i=1}^{n} n_i = n$ . A general problem in judging goodness of fit is to test whether the cell probabilities are specified functions of a fewer number of parameters whose values may be unknown. Let the cell probabilities be given functions  $p_1(\theta), \ldots, p_K(\theta)$ of an unknown vector  $\theta' = (\theta_1, \ldots, \theta_r)$ , where r < K. To test the specification, it is a standard method to use the statistic

(1.1) 
$$T = \sum_{i=1}^{K} (n_i - np_i(\hat{\theta}))^2 / np_i(\hat{\theta})$$

where  $\theta$  is a consistent estimator of  $\theta$ , usually the maximum likelihood estimator. Next, assuming that the specification is true, consider the hypothesis that  $\theta$  is given by  $\theta_i = g_i(\alpha)$ , where  $g_1, \ldots, g_r$  are given functions,  $\alpha' = (\alpha_1, \ldots, \alpha_s)$  and s < r. This hypothesis arises in a test of homogeneity of parallel samples and of independence in a contingency table. The statistic

(1.2) 
$$T^* = n \sum_{i=1}^{K} (p_i(\hat{\theta}) - p_i(\hat{\alpha}))^2 / p_i(\hat{\alpha})$$

is used to test the hypothesis, where a denotes an estimate of a and  $p_i(a)$  denotes the value of  $p_i(\theta)$  as a function of a, under the given hypothesis. It is shown below that T and T\* are asymptotically distributed for large n according to the Chi-square distribution under certain conditions.

DISTRIBUTION ADDITY CODES

or SPECIAL

2. Asymptotic Distribution of T and T\*

First, consider the specification that the multinomial cell probabilities are given functions  $p_1(\theta), \ldots, p_K(\theta)$  of an unknown vector  $\theta = (\theta_1, \ldots, \theta_r)'$ , where r < K. Let  $\theta^0$  denote the true value of  $\theta$ . We make the following assumptions:

(i) The functions  $p_i(\theta)$  are continuous in  $\theta$ , admitting first order partial derivatives which are continuous at  $\theta^0$ .

(ii) Given  $\delta > 0$ , there exists  $\varepsilon > 0$  such that  $\inf_{\substack{\theta = \theta^{\circ} \\ 0 < \delta}} N(\theta) > \varepsilon$ ,

where

 $N(\theta) = \sum_{i=1}^{K} p_i(\theta^\circ) \log (p_i(\theta^\circ)/p_i(\theta)).$ 

Let

$$M(\theta) = \sum_{i=1}^{K} \frac{n_i}{n} \quad \log (p_i(\theta^\circ)/p_i(\theta)).$$

Consider the function N( $\theta$ ) on the sphere  $|\theta - \theta^{\circ}| = \delta$ . Since N( $\theta$ ) is continuous in  $\theta$ , the infimum of N( $\theta$ ) is attained on the sphere. Therefore, in view of (ii), N( $\theta$ )  $\geq \varepsilon$  for every point on the sphere. Since  $\frac{n_i}{n}$  converges in probability to  $p_i(\theta^{\circ})$  as  $n \neq \infty$ , it follows that M( $\theta$ ) > 0 for all points on the sphere with probability approaching 1 as  $n \neq \infty$ .

The log likelihood function is proportional to  $\sum_{i=1}^{K} n_i \log p_i(\theta)$ . In view of (i) and the result given above, we have that for sufficiently large n, the likelihood function has a local maximum inside the open sphere  $|\theta - \theta^0| < \delta$  at a point  $\hat{\theta}$ , say, which is a solution of the likelihood equation

(2.1) 
$$\sum_{i=1}^{K} \frac{n_i}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, j = 1, \dots, r.$$

Since  $\delta$  can be made arbitrarily small, the maximum likelihood estimator  $\hat{\theta}$  is a consistent solution of the likelihood equation.

3.

Let  $I(\theta) = (M'(\theta)) (M(\theta))$  denote the information matrix of the multinomial distribution, where

$$M(\theta) = \left(\frac{1}{\sqrt{p_{i}(\theta)}} - \frac{\partial p_{i}(\theta)}{\partial \theta_{j}}\right)$$

is a K x r matrix, and let Z = M V, where M = M( $\theta^{\circ}$ ) and  $V' = \left(\frac{n_1 - np_1(\theta^{\circ})}{\sqrt{np_1(\theta^{\circ})}}\right), \dots, \frac{n_K - np_K(\theta^{\circ})}{\sqrt{np_K(\theta^{\circ})}}\right).$ 

By the central limit theorem, the asymptotic distribution of V is multivariate normal N(0,I- $\phi\phi'$ ), where 0 denotes the null vector, I denotes the identity matrix and  $\phi' = (\sqrt{p_1(\theta^\circ)}, \dots, \sqrt{p_K(\theta^\circ)})$ . The asymptotic distribution of Z is N(0,I), where I = I( $\theta^\circ$ ). Note that M' $\phi$  = 0.

Substituting  $\theta$  for  $\theta$  in (2.1), the jth equation can be written as

(2.2) 
$$\sum_{i=1}^{K} \frac{n_{i} - np_{i}(\theta^{\circ})}{\sqrt{n}p_{i}(\theta)} \quad \frac{\partial p_{i}(\theta)}{\partial \theta_{j}} = \sum_{i=1}^{K} \frac{\sqrt{n}(p_{i}(\theta) - p_{i}(\theta^{\circ}))}{p_{i}(\theta)} \frac{\partial p_{i}(\theta)}{\partial \theta_{j}}.$$

In view of (i) we have

(2.3) 
$$\mathbf{p}_{\mathbf{i}}(\hat{\theta}) - \mathbf{p}_{\mathbf{i}}(\theta^{\circ}) = \sum_{\ell=1}^{r} (\hat{\theta}_{\ell} - \theta_{\ell}^{\circ}) \frac{\partial \mathbf{p}_{\mathbf{i}}(\theta^{\circ})}{\partial \theta_{\ell}^{\circ}} + \eta |\hat{\theta} - \theta^{\circ}|$$

where  $n \neq 0$  as  $\theta \neq \theta^{\circ}$ . Since  $\theta$  is a consistent estimator of  $\theta$ , as shown above, the left side of (2.2) is asymptotically equivalent (converging in probability) to  $Z_j$ . Therefore, by the substitution of (2.3) on the right side of (2.2) we have that

$$z_{j} = \sum_{\ell=1}^{r} \sqrt{n} \left( \hat{\theta} - \theta_{\ell}^{\circ} \right) I_{\ell j}$$

where  $\frac{a}{2}$  means "asymptotically equivalent to", and  $I_{2j}$  denotes the ljth element of I. Hence

$$z \stackrel{a}{=} \sqrt{n} I (\hat{\theta} - \theta^{\circ})$$

or

(2.4)  $I^{-} Z \stackrel{a}{\rightarrow} \sqrt{n} (\hat{\theta} - \theta^{\circ})$ 

where I denotes a generalized inverse of I, given by I I I = I, and is equal to  $I^{-1}$  if I is non-singular.

Let A be a symmetric matrix with real elements, and let  $X \stackrel{d}{\circ} N(\mu, \Sigma)$ , where  $\stackrel{d}{\circ}$  means "distributed as". If the covariance matrix  $\Sigma$  is non-singular then it is known that the quadratic form X' A X  $\stackrel{d}{\circ} \chi^2_{\nu,\delta}$  (non-central Chisquare with  $\nu$  degrees of freedom and non-centrality parameter  $\delta$ ) if and only if A  $\Sigma$  is idempotent, where  $\nu = \text{Rank}$  A and  $\delta = \frac{1}{2}\mu'A\mu$ . If  $\Sigma$  is singular, then the given condition is only sufficient and  $\nu = \text{Rank}$  A $\Sigma$ (see e.g. Graybill (1976), Theorem 4.7.1). If A =  $\Sigma^-$  is a generalized inverse of  $\Sigma$ , given by  $\Sigma \Sigma^- \Sigma = \Sigma$ , then A $\Sigma$  is idempotent and Rank A  $\Sigma = \text{Rank} \Sigma$ . Therefore, X'  $\Sigma^- X \stackrel{d}{\circ} \chi^2_{\nu,\delta}$ , where  $\nu = \text{Rank} \Sigma$ and  $\delta = \frac{1}{2}\mu' \Sigma^- \mu$ .

Let  $W = (W_1, \ldots, W_K)'$  where

$$W_i = \sqrt{n} (p_i(\hat{\theta}) - p_i(\theta^\circ)) / \sqrt{p_i(\theta^\circ)}, i = 1, \dots, K.$$

From (2.3) we have that

(2.5)

$$W \stackrel{a}{\sim} \sqrt{n} M (\hat{\theta} - \theta^{\circ})$$
$$\stackrel{a}{\sim} M I \stackrel{T}{\sim} Z \qquad by (2.4)$$

Note that  $M \stackrel{f}{=} M'$  is idempotent. From (1.1) and (2.5) we have

$$\mathbf{T}_{z}^{\mathbf{a}} (\mathbf{V} - \mathbf{W}) \quad (\mathbf{V} - \mathbf{W})$$

$$\frac{a}{2}V'(I - MIM')V.$$

Now, V is asymptotically distributed as N(0, I -  $\phi \phi'$ ), (I -M<sub>I</sub><sup>-</sup>M') (I- $\phi \phi'$ ) = I - MI<sup>-</sup>M'- $\phi \phi'$  is idemptotent and

Rank (I-MI<sup>-</sup>M<sup>-</sup>
$$\phi\phi$$
) = Trace (I-MI<sup>-</sup>M<sup>-</sup> $\phi\phi$ )  
= K - 1 - Rank I =  $\beta$ , say.

Therefore, T is asymptotically distributed as  $\chi^2_\beta$ . If I is of full rank then  $\beta = K - r - 1$ .

Next, consider the hypothesis that  $\theta$  is given by  $\theta_i = g_i(\alpha)$ , i = 1,...,r, as described in the previous section. Let  $\nabla(\alpha) = (\partial \theta_i / \partial \alpha_j)$ denote the r x s matrix of the derivatives, and let  $I^*(\alpha)$  denote the information matrix under the given hypothesis. Then

(2.6) 
$$I^{\star}(\alpha) = (\nabla(\alpha)) I(\theta) \nabla(\alpha)$$

with I ( $\theta$ ) being expressed as a function of  $\alpha$ . Let  $\alpha^{\circ}$  denote the true value of  $\alpha$ . Similarly, as in the preceding we have that

 $\mathbf{T}^{\star} \stackrel{\mathbf{a}}{z} \quad \mathbf{V}^{\prime} (\mathbf{M}\mathbf{I}^{\mathsf{T}}\mathbf{M}^{\prime} - \mathbf{M}\nabla (\nabla^{\prime}\mathbf{I}\nabla)^{\mathsf{T}}\nabla^{\prime}\mathbf{M}^{\prime})\mathbf{V}$ 

where  $\nabla = \nabla(\alpha^{\circ})$ . The matrix  $(MI^{-}M\nabla(\nabla'I\nabla)^{-}\nabla'M') ((I - \phi\phi') = (MI^{-}M' - M\nabla(\nabla'I\nabla)^{-}\nabla'M')$  is idempotent and

Rank  $(MIM'-M\nabla(\nabla'I\nabla)\nabla'M') = Rank I-Rank (\nabla'I\nabla)$ 

= Y, say.

Therefore, T\* is asymptotically distributed as  $\chi^2_{\gamma}$ . If I and V are of full rank then  $\gamma = r-s$ .

We have shown that T and T\* are asymptotically distributed according to the Chi-square distribution under the identifiability condition (i) and the continuity assumption (ii).

Remark: For the goodness of fit test where the cell probabilities are completely specified we have  $\beta = K-1$ . In this case T  $\frac{a}{z} \vee' \vee \frac{d}{z} \chi^2_{K-1}$ , asymptotically. For testing homogeneity of r parallel samples or independence of attributes in r x K contigency tables, we have  $\gamma = (r-1)$  (K-1).

## References.

- Graybill, R. A. (1976). The Theory and Application of the Linear Model.
   Wadsworth Publishing Co., Belmont, California.
- [2] Rao, C. R. (1966). Linear Statistical Inference and its Applications. Wiley Publications in Statistics.

| REPORT DOCUMENTATION PAGE  | READ INSTRUCTIONS  |
|--|--|
| REPORT NUMBER  | NO. 3. RECIPIENT'S CATALOG NUMBER  |
| N104   |  |
| TITLE (and Subtitio)   | S. TYPE OF REPORT & PERIOD COVERE  |
| hi-square tests for the multinomial  | Technical Report   |
| listuibution   |  |
| listribution   | 6. PERFORMING ORG. REPORT NUMBER   |
| AUT=204/0  | 296  |
|  |  |
|  | N00014-75-C-0451   |
|  |  |
| PERFORMING ORGANIZATION NAME AND ADDRESS   | 10. PROGRAM ELEMENT, PROJECT, TASK   |
| Clemson University   |  |
| Dept. of Mathematical Sciences   | NR 042-271   |
| Clemson, South Carolina 29631  |  |
| CONTROLLING OFFICE NAME AND ADDRESS  | 12-5-1978  |
| Code 436   | D. NUMBER OF PAGES   |
| Arlington, Va. 22217   | 7  |
| MONITORING AGENCY NAME & ADDRESS(It different from Controlling Office  | ) 15. SECURITY CLASS. (of this report)   |
|  |  |
|  | Unclassified   |
|  | 154. DECLASSIFICATION DOWNGRADING<br>SCHEDULE  |
|  |  |
| DISTRIBUTION STATEMENT (of this Report)<br>Approved for public release; distribution unlim   | dited.   |
| DISTRIBUTION STATEMENT (of the eservect entered in Block 20, if different  | trom Report)   |
| DISTRIBUTION STATEMENT (of the ebetract entered in Block 20, if different  | trom Report)   |
| DISTRIBUTION STATEMENT (of this Report)<br>Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the ebetract entered in Block 20, if different<br>SUPPLEMENTARY NOTES   | ited.  |
| DISTRIBUTION STATEMENT (of this Report)<br>Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the observact entered in Block 20, if different<br>SUPPLEMENTARY NOTES  | tron Report)   |
| Approved for public release; distribution unlim<br>OISTRIBUTION STATEMENT (of the ebetract entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS (Continue on reverse side if necessary and identify by block number<br>fultinomial distributions; Chi-square te   | <pre>ited. /rom Report) /or/ St;</pre>   |
| Approved for public release; distribution unlim<br>OISTRIBUTION STATEMENT (of the observace entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY NORDS (Continue on reverse side if necessary and identify by block number<br>fultinomial distributions; Chi-square te<br>Contingency table.  | rom Report)  |
| Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the exertact entered in Block 20, if different<br>DISTRIBUTION STATEMENT (of the exertact entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS Continue on reverse side if necessary and identify by block number<br>fultinomial distributions; Chi-square te<br>Contingency table.   | <pre>ited. /rom Report) /or) St;</pre>   |
| DISTRIBUTION STATEMENT (of this Report)<br>Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the obserract entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS Continue on reverse side if necessary and identify by block number<br>Aultinomial distributions; Chi-square te<br>Contingency table.  | <pre>ited. // // // // // // // // // // // // //</pre>  |
| DISTRIBUTION STATEMENT (of this Report)<br>Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the ebetract entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS (Continue on reverse side if necessary and identify by block number<br>fultinomial distributions; Chi-square te<br>Contingency table.<br>ABSTRACT (Continue on reverse side if necessary and identify by block number<br>this paper gives a proof of the asymptot  | <pre>ited. //om Report) /** /** /** /** /** /** /** /** /** /*</pre>   |
| Approved for public release; distribution unlim<br>- DISTRIBUTION STATEMENT (of the observer entered in Block 20, if different<br>- DISTRIBUTION STATEMENT (of the observer entered in Block 20, if different<br>- SUPPLEMENTARY NOTES<br>- SUPPLEMENTARY NOTES<br>- SUPPLEMENTARY NOTES<br>- SUPPLEMENTARY NOTES<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary table.<br>- ABSTRACT (Continue on reverse side if necessary and identify by block number<br>- Supplementary and the supplementary and identify by block number<br>- Supplementary and the supplementary and identify by block number<br>- Supplementary and the supplementary and the supplementary and identify by block number<br>- Supplementary and the supplementary and identify by block number<br>- Supplementary and the supplementary and identify by block number<br>- Supplementary and the s | <pre>ited. /rom Report) ** ** ** ** ** ** ** ** ** ** ** ** **</pre>   |
| Approved for public release; distribution unlim<br>OUSTRIBUTION STATEMENT (of the observace entered in Block 20, if different<br>OUSTRIBUTION STATEMENT (of the observace entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS (Continue on reverse side if necessary and identify by block number<br>fultinomial distributions; Chi-square te<br>Contingency table.<br>Asstract (Continue on reverse side if necessary and identify by block number<br>This paper gives a proof of the asymptot<br>square tests associated with the multino<br>generally used in the analysis of catego  | <pre>ited. // // // // // // // // // // // // //</pre>  |
| Approved for public release; distribution unlim<br>OISTRIBUTION STATEMENT (of the ebetract entered in Block 20, if different<br>OISTRIBUTION STATEMENT (of the ebetract entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS (Continue on reverse side if necessary and (dentify by block number<br>fultinomial distributions; Chi-square te<br>Contingency table.<br>ABSTRACT (Continue on reverse side if necessary and identify by block number<br>This paper gives a proof of the asymptot<br>square tests associated with the multino<br>generally used in the analysis of catego<br>tingency tables.  | <pre>ited. //on Report) *** st; *** ic property of the Chi- mial distribution, prical data, such as con-</pre> |
| Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the observate entered in Block 20, if different<br>DISTRIBUTION STATEMENT (of the observate entered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY HORDS Continue on reverse side if necessary and identify by block number<br>Aultinomial distributions; Chi-square te<br>Contingency table.<br>ABSTRACT (Continue on reverse side if necessary and identify by block number<br>This paper gives a proof of the asymptot<br>square tests associated with the multino<br>generally used in the analysis of catego<br>tingency tables.   | <pre>ited. /rom Report st; ** ic property of the Chi- mial distribution, rical data, such as con-</pre>        |
| DISTRIBUTION STATEMENT (of this Report)<br>Approved for public release; distribution unlim<br>DISTRIBUTION STATEMENT (of the exercise intered in Block 20, if different<br>SUPPLEMENTARY NOTES<br>KEY CORDS (Continue on reverse side if necessary and identify by block number<br>fultinomial distributions; Chi-square te<br>Contingency table.<br>ABSTRACT (Continue on reverse side if necessary and identify by block number<br>his paper gives a proof of the asymptot<br>iquare tests associated with the multino<br>penerally used in the analysis of catego<br>ingency tables.  | <pre>ited. //om Report // st; // ic property of the Chi- mial distribution, // orical data, such as con-</pre> |