



A study on Freeman-Tukey test statistic under the symmetry model for square contingency tables

Gökçen ALTUN¹

¹ Bartın University, Department of Computer Technology and Information Systems, 74100, Bartın / TURKEY

Abstract

The symmetry model is the basic model in the analysis of square contingency tables. Multiple test statistics have been developed for the goodness of fit test. Freeman-Tukey test statistics is appropriate to be used in large samples. However, the required sample size to use the Freeman-Tukey test statistics is not clear. In this paper, the asymptotic properties of Freeman-Tukey test statistic are discussed via extensive Monte-Carlo simulation study. The Freeman-Tukey test statistic is compared with members of power-divergence family test statistic under the symmetry model. The results of simulation study are evaluated based on the Type-I error and power of a test. The results of simulation study and artificial data study show that Freeman-Tukey's T^2 test statistic does not converge to chi-squared distribution for both sparse and non-sparse square contingency tables.

Article info

History:
Received: 02.12.2020
Accepted: 31.03.2021

Keywords:
Freeman-Tukey test statistic, Power-divergence family, Goodness-of-fit test statistics, Square contingency tables, Symmetry model.

1. Introduction

Square contingency tables that arise independent samples where the row and column variables have the same level. Let n_{ij} be the observed frequency in the cell (i,j) and p_{ij} denotes the probability of the same cell. The representation of the $R \times R$ dimensional square contingency table can be found in [1]. Some specific models used in the analysis of these kinds of tables. These models are mostly in the symmetrical pattern. Then, the complete symmetry model (S) is defined by;

$$p_{ij} = p_{ji}, \quad i, j = 1, 2, \dots, R, \quad (\text{for } i \neq j) \quad (1)$$

and is based on $R(R-1)/2$ degree of freedom, where R is the dimension of the square table [2]. The likelihood estimates of expected values e_{ij} under the S model is $e_{ij} = (n_{ij} + n_{ji})/2$. This model indicates that the probability that an observation will fall in cell (i,j) is equal to probability that it falls in symmetric cell (j,i).

Goodness-of-fit tests summarize the discrepancy between the observed values and the expected values under the corresponding model. Cressie and Read (1984) introduced a class of goodness-of-fit test statistics named as power-divergence (PD) family of statistics. The PD statistic includes Pearson's χ^2 and the likelihood ratio statistic G^2 as well as other

statistics such as Freeman-Tukey's T^2 , the modified likelihood ratio statistics GM^2 , and Neyman's modified χ^2 statistics. All of these statistics are asymptotically chi-squared distributed with appropriate degrees of freedom [3,4,5].

The researchers have shown a great interest to compare goodness of test statistics for analyzing the sparse contingency tables. Cochran et. al. showed that Pearson statistic does not follow the chi-squared distribution well for small expected values [6-8]. Cochran et. al. discussed which approximation is "reasonable" for the Pearson chi-squared statistic [6-10]. "Standard rules" specify that the minimum cell expectation should be five, with a few cells possibly smaller. The true sampling distributions converge to chi-squared as $n \rightarrow \infty$, for a fixed number of cells N . The adequacy of the chi-squared approximation depends both on n and N [11]. Cochran suggested that to test independence with $df > 1$, a minimum expected value $e_{ij} \approx 1$ is permissible as long as no more than about 20% of $e_{ij} < 5$ [6]. Koehler et al. showed that it is more appropriate to use χ^2 test statistics instead of G^2 for sparse tables and small sample sizes [12-14]. When n/N (sparseness index) is less than 5 the distribution of G^2 is usually poorly approximated by chi-squared. Depending on the sparseness, p values can be either too large or too

small. As N increases, the approximation to the chi-square distribution decreases [13]. However, Haberman showed that the approximation tends to be poor for sparse tables containing both small and moderately large e_{ij} (see, Cressie and Read (1989) and Lawal (1984) for detailed information) [15-17]. Larntz compared Pearson's χ^2 , likelihood ratio statistic G^2 and Freeman-Tukey's T^2 statistics based on the five models [14]. Larntz stated that Pearson's χ^2 demonstrates the best approximation to the chi-squared distribution for small samples and the other two statistics are not approximate well [14]. Fienberg emphasized that the behavior of G^2 in large sparse multinomial structures requires serious attention [3]. Baglivo et al. stated that each statistic in the power-divergence family can be regarded as a different measure of goodness of fit and these differences make the tests useful in different situations [18]. Many authors including Bishop et al., Aitchinson et al., Simonoff and Burman concerned with smoothing such tables in order to avoid the problems associated with sparseness [19-23]. Kim et al. studied on Zelterman's D^2 statistic and compared the efficiency of Zelterman's D^2 statistic with other well-known statistics via simulation study [24,25]. Aktaş compared the power-divergence statistics based on the power values under the Quasi-Independence model in square contingency tables [26].

In this paper, we compare the PD statistics for various λ values concerning their power values and Type-I

$$I(\lambda) = \frac{2}{\lambda(\lambda + 1)} \sum_{i=1}^R \sum_{j=1}^R n_{ij} \left[\left(\frac{n_{ij}}{e_{ij}} \right)^\lambda - 1 \right], \quad i, j = 1, 2, \dots, R, \quad \lambda \in \mathfrak{R} \tag{2}$$

The PD family of statistics, given in (2), is not valid of $\lambda = 0$ and $\lambda = -1$. Therefore, the following equations are obtained by using limit for the cases $\lambda = 0$ and $\lambda = -1$.

$$\lim_{\lambda \rightarrow 0} \frac{2}{\lambda(\lambda + 1)} \sum_{i=1}^R \sum_{j=1}^R n_{ij} \left[\left(\frac{n_{ij}}{e_{ij}} \right)^\lambda - 1 \right] = 2 \sum_{i=1}^R \sum_{j=1}^R n_{ij} \left(\log \frac{n_{ij}}{e_{ij}} \right) \tag{3}$$

$$\lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda + 1)} \sum_{i=1}^R \sum_{j=1}^R n_{ij} \left[\left(\frac{n_{ij}}{e_{ij}} \right)^\lambda - 1 \right] = 2 \sum_{i=1}^R \sum_{j=1}^R e_{ij} \left(\log \frac{e_{ij}}{n_{ij}} \right) \tag{4}$$

It is straightforward to verify that the statistic $I(\lambda)$ reduces to Pearson's χ^2 when $\lambda = 1$, likelihood ratio G^2 when $\lambda = 0$, Freeman Tukey's T^2 when $\lambda = -1/2$

2.2. Simulation study

In the simulation study, the Type-I errors of the power-divergence statistics are obtained for $\lambda = 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1(\chi^2), 0(G^2), -1/2(T^2), 2/3(C^2)$.

The powers of the tests are calculated under the S model. $N = 50,000$ sample of sizes $n =$

errors under the S model in the square contingency table where the observations are cross-classified by two variables with the same categories. The goal of this paper is to show that Freeman-Tukey's T^2 test statistic does not converge to the chi-squared distribution not only in sparse square contingency tables but also in non-sparse square contingency tables. To achieve this goal, an extensive simulation study is conducted to show the relative efficiencies of PD test statistics under the symmetry model with various dimensions of tables and different sample sizes.

The remaining part of the paper is organized as follows: In Section 2, PD family of statistics are presented with its theoretical background. Section 2 also contains the simulation study and application to the artificial data set. Some concluding remarks are given in Section 3.

2. Materials and Methods

2.1. Power-Divergence family

Cressie and Read introduced a class of goodness-of-fit test statistics that can be expressed based on a family of power-divergence statistics. All members of the statistics are members of the power-divergence family. Let n_{ij} and e_{ij} represent the observed and expected frequencies. The PD family of statistics, $I(\lambda)$, is given by [27]

and Cressie Read test statistic C^2 when $\lambda = 2/3$. All of these statistics are asymptotically chi-squared distributed with appropriate degrees of freedom.

10,15, ...,200 for $R = 3$, $n = 20,25, \dots,300$ for $R = 4$, $n = 30,35, \dots,600$ for $R = 5$ and $n =$

40,45, ...,900 for $R = 6$ are generated by means of multinomial distribution under the S model using the probability matrices given below. The results of the

study are obtained by using the R programming language.

Table 1. Probability matrices for generating random frequencies from the S model.

R	Type I error	Power of a test
3	$\begin{bmatrix} 0,05 & 0,13 & 0,17 \\ 0,13 & 0,05 & 0,125 \\ 0,17 & 0,125 & 0,05 \end{bmatrix}$	$\begin{bmatrix} 0,05 & 0,10 & 0,15 \\ 0,30 & 0,05 & 0,19 \\ 0,05 & 0,06 & 0,05 \end{bmatrix}$
4	$\begin{bmatrix} 0,02 & 0,08 & 0,10 & 0,06 \\ 0,08 & 0,02 & 0,05 & 0,11 \\ 0,10 & 0,05 & 0,02 & 0,06 \\ 0,06 & 0,11 & 0,06 & 0,02 \end{bmatrix}$	$\begin{bmatrix} 0,02 & 0,04 & 0,15 & 0,03 \\ 0,12 & 0,02 & 0,085 & 0,055 \\ 0,05 & 0,035 & 0,02 & 0,09 \\ 0,09 & 0,145 & 0,03 & 0,02 \end{bmatrix}$
5	$\begin{bmatrix} 0,01 & 0,05 & 0,08 & 0,04 & 0,045 \\ 0,05 & 0,01 & 0,03 & 0,09 & 0,01 \\ 0,08 & 0,03 & 0,01 & 0,02 & 0,06 \\ 0,04 & 0,09 & 0,02 & 0,01 & 0,05 \\ 0,045 & 0,01 & 0,06 & 0,05 & 0,01 \end{bmatrix}$	$\begin{bmatrix} 0,01 & 0,075 & 0,04 & 0,04 & 0,0225 \\ 0,025 & 0,01 & 0,045 & 0,045 & 0,015 \\ 0,12 & 0,015 & 0,01 & 0,01 & 0,09 \\ 0,02 & 0,135 & 0,03 & 0,01 & 0,025 \\ 0,0675 & 0,005 & 0,03 & 0,075 & 0,01 \end{bmatrix}$
6	$\begin{bmatrix} 0,001 & 0,017 & 0,02 & 0,03 & 0,012 & 0,014 \\ 0,017 & 0,001 & 0,07 & 0,05 & 0,02 & 0,013 \\ 0,02 & 0,07 & 0,001 & 0,09 & 0,009 & 0,06 \\ 0,03 & 0,05 & 0,09 & 0,001 & 0,044 & 0,008 \\ 0,012 & 0,02 & 0,009 & 0,044 & 0,001 & 0,04 \\ 0,014 & 0,013 & 0,06 & 0,008 & 0,04 & 0,001 \end{bmatrix}$	$\begin{bmatrix} 0,001 & 0,0255 & 0,01 & 0,045 & 0,018 & 0,007 \\ 0,0085 & 0,001 & 0,105 & 0,025 & 0,03 & 0,0065 \\ 0,03 & 0,035 & 0,001 & 0,045 & 0,0135 & 0,09 \\ 0,015 & 0,075 & 0,135 & 0,001 & 0,066 & 0,012 \\ 0,006 & 0,01 & 0,0045 & 0,022 & 0,001 & 0,02 \\ 0,021 & 0,0195 & 0,03 & 0,004 & 0,06 & 0,001 \end{bmatrix}$

Since the type 1 error is calculated under the accuracy of the H_0 hypothesis, the matrices generated for the type 1 error are symmetrical. For the strength of the test, the deterioration in the symmetrical structure of the matrix is made on the condition that the symmetrical cells with respect to the diagonal are approximately 1 to 3. The convergences of test statistics to chi-squared distribution are evaluated based on the critical point, 0.06, for the Type-I error. When the obtained Type-I error is lower than 0.06 value, the corresponding test statistic is asymptotically distributed as chi-squared distribution. Figures 1.a and 1.b display the simulation results for $R=3$.

As seen from Figure 1.a, when the sample size is lower than 50, all test statistics, except G^2 , T^2 , $I(0.1)$ and $I(0.2)$, converge to the chi-squared distribution. When the sample size is higher than 50, all test statistics, except for T^2 converge to the chi-squared distribution. It is clear that the test statistic with the highest power of a test is T^2 (Figure 1.b). However, the T^2 test statistic converges to the chi-squared distribution after the sample size is approximately 80. This issue can be expressed based on the sparseness index. It is easy to see that T^2 test statistic converges to chi-squared distribution when the sparseness index is higher than 9 for $R=3$.

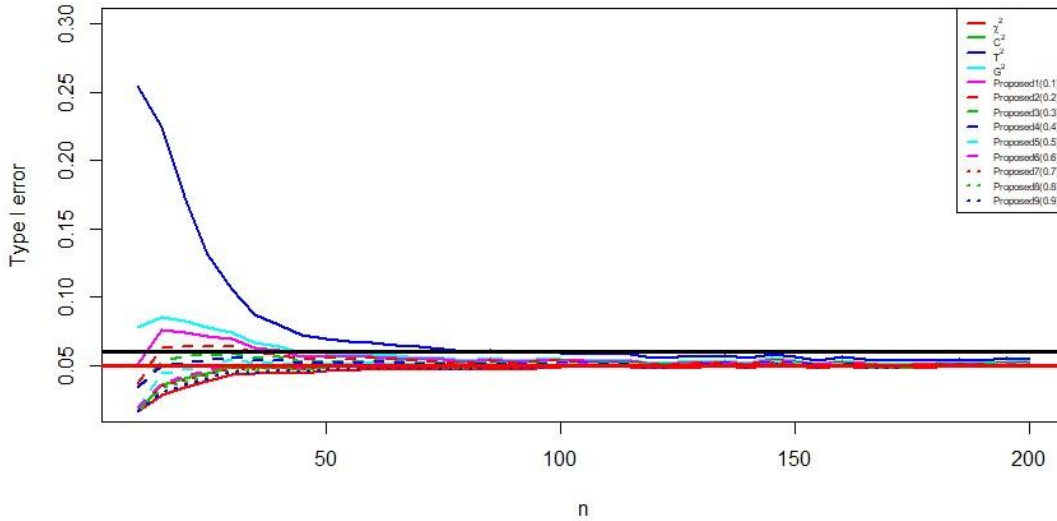


Figure 1a. The Type-I errors of the corresponding test statistics for 3×3 square contingency tables.

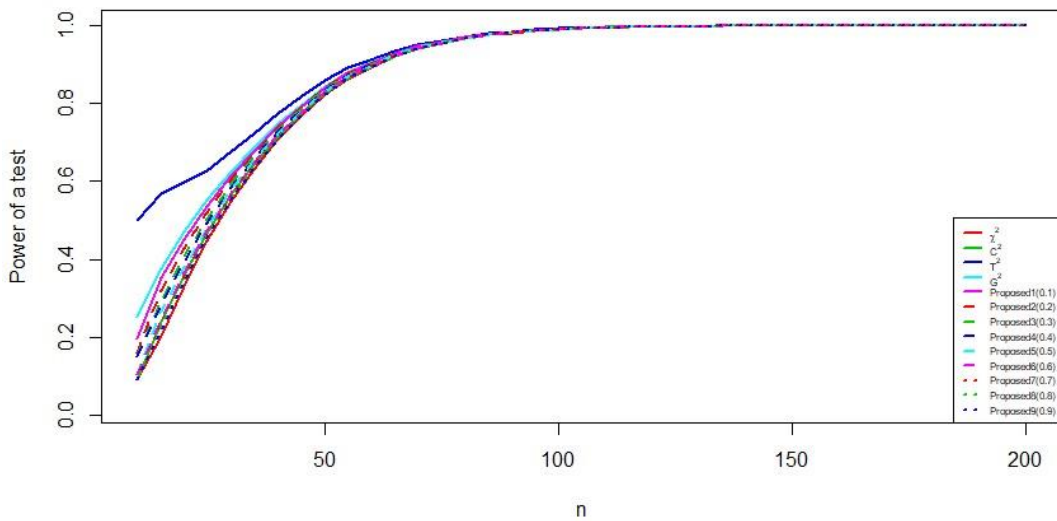


Figure 1b. The power of a test results of the corresponding test statistics for 3×3 square contingency tables.

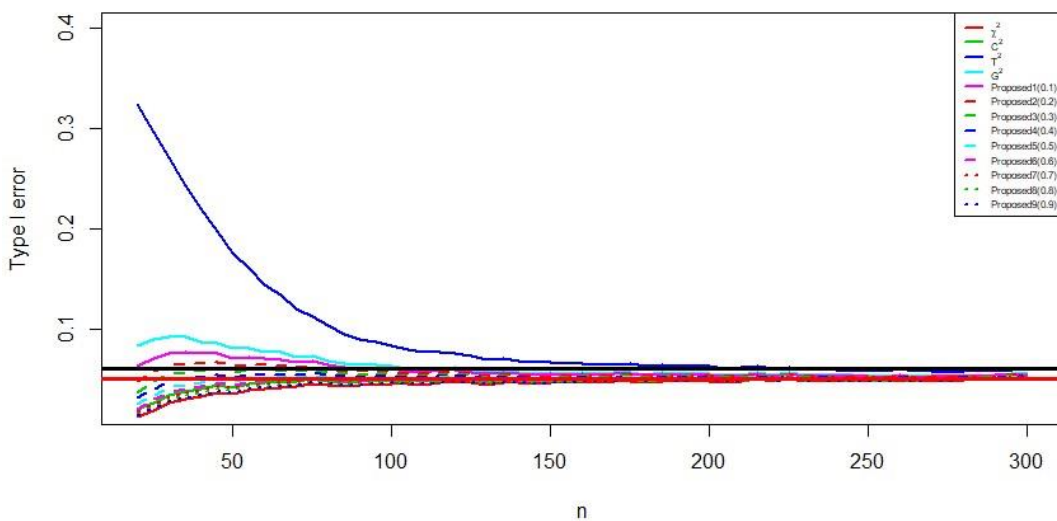


Figure 2a. The Type-I errors of the corresponding test statistics for 4×4 square contingency tables.

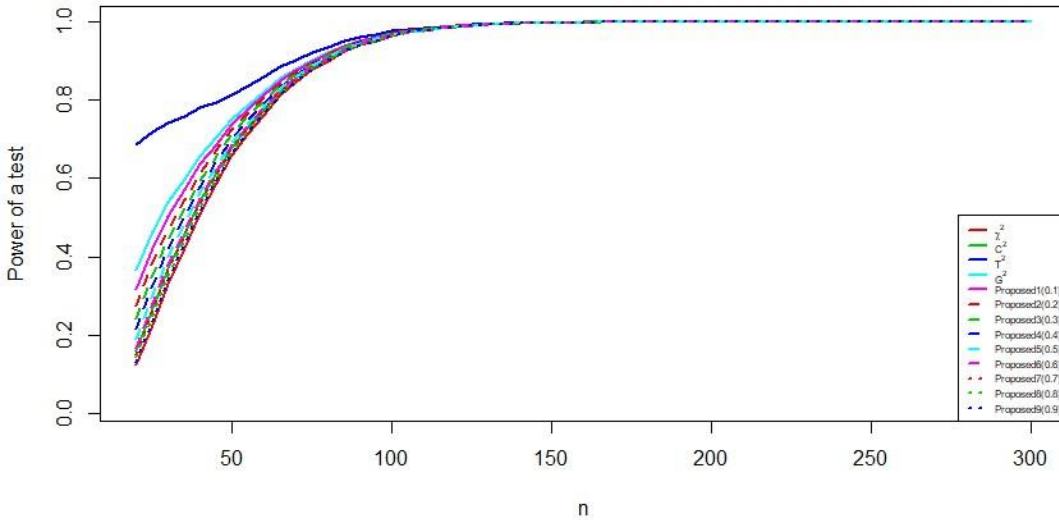


Figure 2b. The power of a test results of the corresponding test statistics for 4×4 square contingency tables.

Figures 2.a and 2.b displays the simulation results for $R=4$. As seen from the Figure 2.a, all test statistics, except $G^2, T^2, I(0.1)$ and $I(0.2)$, converge to the chi-squared distribution when the sample size is lower than 100. When the sample size is higher than 100, all test statistics except for T^2 converge to the chi-squared distribution. From Figure 2.b, it is clear that the test statistic with the highest power of a test is Freeman Tukey's T^2 . However, the T^2 test statistic converges to the chi-squared distribution when the sample size is approximately 200. It can be expressed based on the sparseness index. Freeman Tukey's T^2 converges to chi-squared distribution when the sparseness index is higher than 12 for $R=4$. Figures 3.a and 3.b display the simulation results for $R=5$. From the Figure 3.a, all test

statistics, except $G^2, T^2, I(0.1)$ and $I(0.2)$ converge to the chi-squared distribution when the sample size is approximately 200. When the sample size is higher than 200, all test statistics except for G^2 and T^2 approximate to the chi-squared distribution. Likelihood ratio test statistic G^2 converges to the chi-squared distribution after the sample size is approximately 300. From Figure 3.b, it is clear that the test statistic with the highest power of a test is Freeman Tukey's T^2 .

However, the T^2 test statistic converges to the chi-squared distribution after the sample size is approximately 600. In other words, Freeman Tukey's T^2 converges to chi-squared distribution when the sparseness index is higher than 24 for $R=5$.

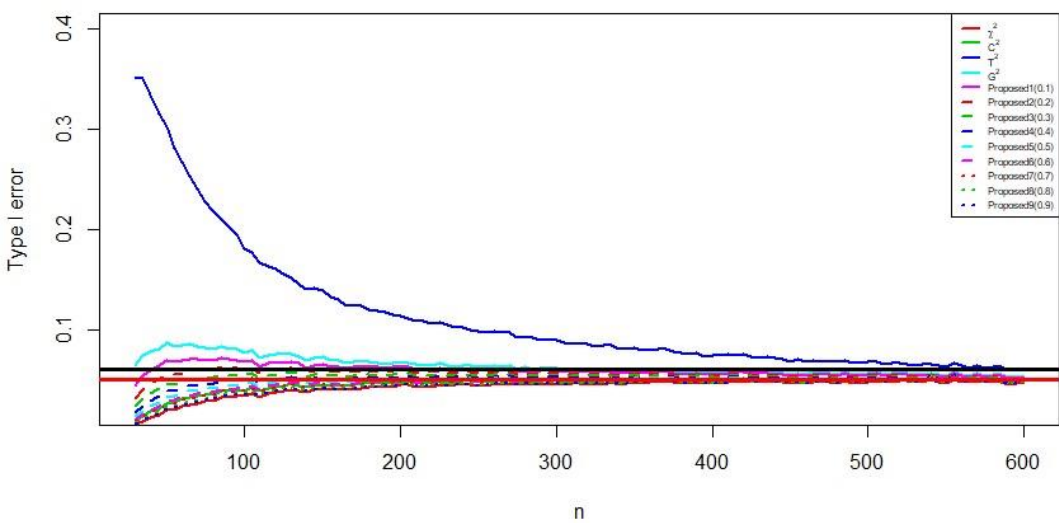


Figure 3a. The Type-I errors of the corresponding test statistics for 5×5 square contingency tables.

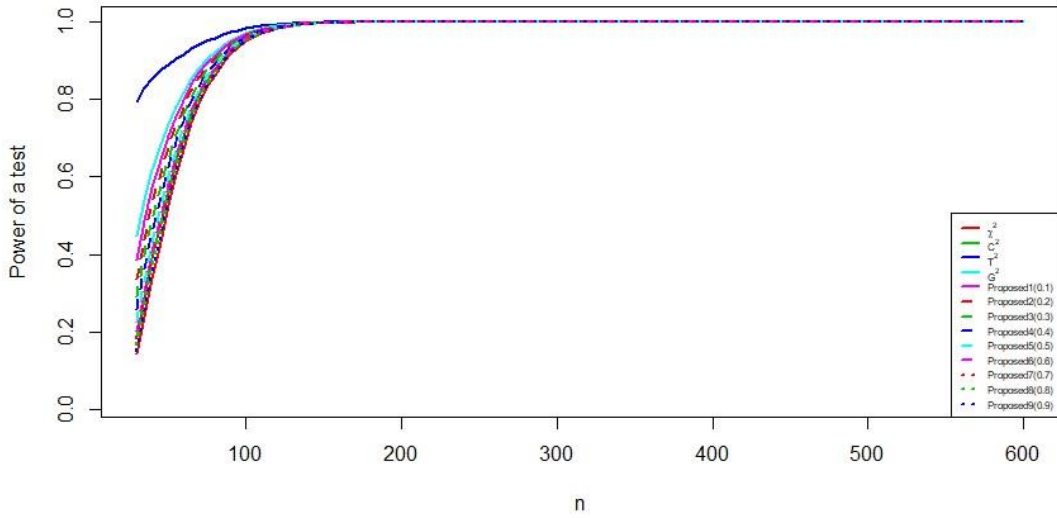


Figure 3b. The power of a test results of the corresponding test statistics for 5×5 square contingency tables

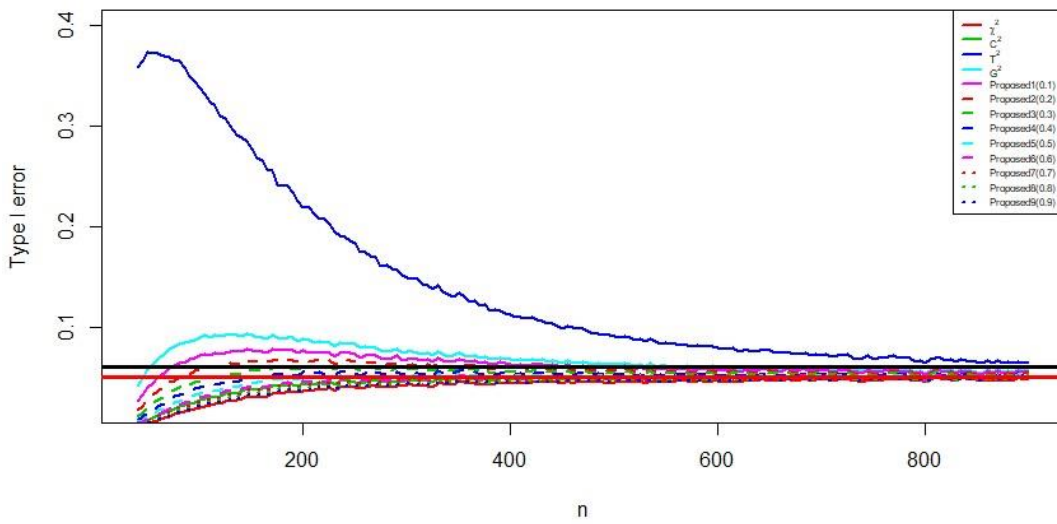


Figure 4a. The Type-I errors of the corresponding test statistics for 6×6 square contingency tables.

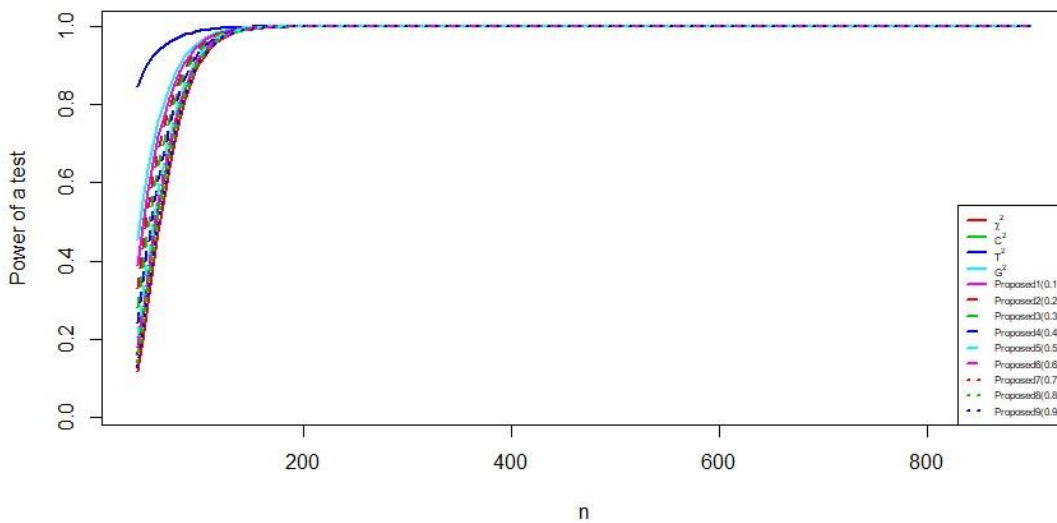


Figure 4b. The power of a test results of the corresponding test statistics for 6×6 square contingency tables.

Figures 4.a and 4.b display the simulation results for R=6. From the Figure 4.a, all test statistics, except $G^2, T^2, I(0.1)$ and $I(0.2)$ converge to the chi-squared distribution when the sample size is lower than 400. When the sample size is higher than 400, all test statistics except for $G^2, I(0.1)$ and T^2 converge to the chi-squared distribution. Likelihood ratio G^2 and $I(0.1)$ test statistic converges to the chi-squared distribution when the sample size is approximately 500. From Figure 3.b, it is clear that the test statistic with the highest power of a test is Freeman Tukey's T^2 . However, the T^2 test statistic converges to the chi-squared distribution when the sample size is approximately 900. Based on the sparseness index, Freeman Tukey's T^2 converges to chi-squared distribution when the sparseness index is higher than 25 for R=6.

As seen from simulation results, the goodness-of-fit test statistics and test statistics corresponding to different λ values introduced from the PD family of statistics do not converge to the chi-squared distribution in very small or sparse samples. When the sample size increases, it is clear that all of the test statistics converge to the chi-squared distribution. However, the key point here is that the Freeman

Tukey's T^2 test statistic does not converge to the chi-squared distribution in large samples. It is also possible to say the same comment for Likelihood ratio G^2 test statistic up to a certain sample size. The sample size should be very large to use the Freeman Tukey's T^2 test statistic in analysis of square contingency tables. In other words, if the sparseness index is lower than 5, all test statistics do not show well-convergence to the chi-squared distribution. When the sparseness index is higher than 5, some of the test statistics show well-convergence to the chi-squared distribution. However, the sparseness index should be very large for the convergence of Freeman Tukey's T^2 test statistic to the chi-squared distribution. Therefore, the use of Freeman Tukey's T^2 in the analysis of square contingency tables requires a very large sparseness index or frequencies.

2.3. Artificial data study

Here, an artificial data set is generated under the S model for n=200 and R=5 to demonstrate the efficiencies of test statistics. Since the data set is generated under the S model, the expectation is that the S model holds for the generated data. Table 2 shows the simulated data set.

Table 2. Artificial data for 5 X 5 square contingency table

	[1]	[2]	[3]	[4]	[5]	Total
[1]	9 (9)	1 (7.5)	10 (10)	7 (7.5)	8 (8)	35
[2]	14 (7.5)	8 (8)	10 (9.5)	12 (9)	10 (8.5)	54
[3]	10 (10)	9 (9.5)	5 (5)	5 (6.5)	10 (9)	39
[4]	8 (7.5)	6 (9)	8 (6.5)	9 (9)	7 (6)	38
[5]	8 (8)	7 (8.5)	8 (9)	5 (6)	6 (6)	34
Total	49	31	41	38	41	200

The artificial data set is fitted to the S model by using the test statistics given in Section 2. Table 3 lists the estimated test values and corresponding p-values. As seen from Table 3, all test statistics, expect T^2 , reveal that the corresponding data is well-fitted with the S model. The T^2 test statistic rejects the H_0 hypothesis with p-value=0.0282. Here, H_0 represents that the S model holds for the corresponding data set. Since the

T^2 test statistic does not converge to chi-squared distribution not only small samples but also large samples, the T^2 test statistic causes the false decision in hypothesis testing which is called a Type-I error.

Table 3. Results for artificial data set under the S model

Test statistics and λ values	values	p-value
$\chi^2(1)$	15.1632	0.1262
$G^2(0)$	17.3930	0.0660
$T^2(-1/2)$	20.1069	0.0282
$C^2(2/3)$	15.5816	0.1123
0.1	17.0151	0.0740
0.2	16.6800	0.0818
0.3	16.3841	0.0892
0.4	16.1237	0.0961
0.5	15.8959	0.1027
0.6	15.6980	0.1086
0.7	15.5277	0.1140
0.8	15.3830	0.1187
0.9	15.2620	0.1228

3. Results and Discussion

In this paper, the goodness-of-fit test statistics which are commonly used in the literature such as χ^2 , G^2 , T^2 and power-divergence statistics for various λ values are compared in terms of their Type-I error and powers under the S model in square contingency table. The square contingency table analyses in the literature are mostly for large sample sizes. In small samples, there is not enough square contingency table study. It is stated in the studies conducted that T^2 test statistics should be used in large samples. However, how large the sample should be is not included. As a result of the simulation study, it is concluded that all test statistics in large samples asymptotically converge to the chi-square distribution, but the sample size should be very large for T^2 test statistics converge to the chi-square distribution. As a result of the study, it is stated that how much the sample size should be. A simulation study is conducted to demonstrate the converge of these statistics to chi-squared distribution for small and large samples. It is concluded that Freeman-Tukey's T^2 test statistics are not well in convergency to chi-squared distribution not only in sparse square contingency tables but also in non-sparse square contingency tables. We hope that the results given in the paper will be very useful to practitioners and academicians studying in this field.

Conflicts of interest

The authors state that did not have a conflict of interests.

References

- [1] Altun G., Karesel Olumsuzluk Tablolarında Model Uyumunun Sapma Ölçüsü ile Belirlenmesi, PD Thesis, Hacettepe Üniversitesi, Fen Bilimleri Enstitüsü, (2018).
- [2] Goodman L. A., Multiplicative models for square contingency tables with ordered categories, *Biometrika.*, 66(3) (1979) 413-418.
- [3] Fienberg S. E., The use of chi-squared statistics for categorical data problems, *Journal of the Royal Statistical Society. Series B (Methodological.)*, 41(1) (1979) 54-64.
- [4] Horn S. D., Goodness-of-fit tests for discrete data: a review and an application to a health impairment scale, *Biometrics.*, 33(1) (1977) 237-247.
- [5] Watson G. S., Some recent results in chi-square goodness-of-fit tests, *Biometrics.*, 15(3) (1959) 440-468.
- [6] Cochran W. G., The χ^2 test of goodness of fit, *The Annals of Mathematical Statistics.*, 23(3) (1952) 315-345.
- [7] Tate M. W., Hyer L. A., Inaccuracy of the X^2 test of goodness of fit when expected frequencies are small, *Journal of the American Statistical Association.*, 68(344) (1973) 836-841.
- [8] Yarnold J. K., The minimum expectation in X^2 goodness of fit tests and the accuracy of approximations for the null distribution, *Journal of the American Statistical Association.*, 65 (330) (1970) 864-886.
- [9] Fisher R. A., Statistical methods for research workers, 13th ed. New York: Hafner Publishing Co, (1958) 356.
- [10] Roscoe J. T., Byars J. A., An investigation of the restraints with respect to sample size commonly imposed on the use of the chi-square statistic, *Journal of the American Statistical Association.*, 66(336) (1971) 755-759.
- [11] Agresti A, Categorical data analysis, 2nd ed. New Jersey: John Wiley & Sons, (2003) 482.
- [12] Koehler K. J., Goodness-of-fit tests for log-linear models in sparse contingency tables, *Journal of the American Statistical Association.*, 81 (394) (1986) 483-493.

- [13] Koehler, K. J., Larntz, K., An empirical investigation of goodness-of-fit statistics for sparse multinomials, *Journal of the American Statistical Association.*, 75(370) (1980) 336-344.
- [14] Larntz K., Small-sample comparisons of exact levels for chi-squared goodness-of-fit statistics, *Journal of the American Statistical Association.*, 73(362) (1978) 253-263.
- [15] Haberman S. J., A warning on the use of chi-squared statistics with frequency tables with small expected cell counts, *Journal of the American Statistical Association.*, 83(402) (1988) 555-560.
- [16] Cressie N., Read T. R., Pearson's X^2 and the loglikelihood ratio statistic G^2 : a comparative review, *International Statistical Review/Revue Internationale de Statistique.*, 57(1) (1989) 19-43.
- [17] Lawal H. B., Comparisons of the X^2 , Y^2 , Freeman-Tukey and Williams's improved G^2 test statistics in small samples of one-way multinomials, *Biometrika*, 71(2) (1984) 415-418.
- [18] Baglivo J., Olivier D., Pagano M., Methods for exact goodness-of-fit tests, *Journal of the American Statistical Association.*, 87(418) (1992) 464-469.
- [19] Bishop Y. M. M., Fienberg S. E., Holland P. W., Discrete multivariate analysis: Theory and practice, Cambridge: The Massachusetts Institute of Technology Press Google Scholar., (1975).
- [20] Aitchison J., Aitken C. G., Multivariate binary discrimination by the kernel method, *Biometrika.*, 63(3) (1976) 413-420.
- [21] Simonoff J. S., A penalty function approach to smoothing large sparse contingency tables, *The Annals of Statistics.*, 11(1) (1983) 208-218
- [22] Simonoff J. S., Probability estimation via smoothing in sparse contingency tables with ordered categories, *Statistics & Probability Letters.*, 5(1) (1987) 55-63.
- [23] Burman P., Central limit theorem for quadratic forms for sparse tables, *Journal of Multivariate Analysis.*, 22(2) (1987) 258-277.
- [24] Kim S. H., Choi H., Lee S., Estimate-based goodness-of-fit test for large sparse multinomial distributions, *Computational Statistics & Data Analysis.*, 53(4) (2009) 1122-1131.
- [25] Zelterman D., Goodness-of-fit tests for large sparse multinomial distributions, *Journal of the American Statistical Association.*, 82(398) (1987) 624-629.
- [26] AKTAŞ S., Power Divergence Statistics under Quasi Independence Model for Square Contingency Tables, *Sains Malaysiana.*, 45(10) (2016) 1573-1578.
- [27] Cressie N., Read T. R., Multinomial goodness-of-fit tests, *Journal of the Royal Statistical Society, Series B (Methodological).*, 46(3) (1984) 440-464.