

OS TESTES NÃO-PARAMÉTRICOS G E QUI-QUADRADO E SUAS APROXIMAÇÕES

Marcos A. Antonello Scremin

Curso de Pós-graduação em Métodos Quantitativos. Departamento de Estatística. Centro de Ciências Naturais e Exatas. UFSM. Santa Maria, RS.

Anaelena B. de Moraes Ethur e Maria Emília Camargo

Departamento de Estatística. Centro de Ciências Naturais e Exatas. UFSM. Santa Maria, RS.

RESUMO

Neste trabalho é apresentada a procedência e a aproximação de um teste não-paramétrico alternativo ao teste do Qui-quadrado (χ^2), conhecido na literatura como teste G, mas pouco divulgado. Foi verificada a origem desta estatística através de deduções e também a aproximação à estatística do Qui-quadrado, utilizada nos testes de hipóteses.

SUMMARY

SCREMIN, M. A. A.; ETHUR, A. B. de M. and CAMARGO, M. E., The no-parametries tests g and chi-square and their aproximations. Ciência e Natura, 13:55-59, 1991.

In this paper we show from where became and the aproach of a alternative to the few published no-parametric test, knowned chis-quare in the literature as G test. By deduction, we research the origin of this statistics and the aproach of the statistics of chi-square that is used at the hypothesis tests.

1 - ESTATÍSTICA G

A Estatística G é definida como sendo duas vezes o logarítmo natural da razão entre a probabilidade da amostra com todos os parâmetros estimados a partir dos dados e a probabilidade da amostra, admitindo-se que a hipótese nula é verdadeira. Esta definição origina-se do teste de razão de verossimilhança, como segue:

Considerando-se k categorias com probabilidades de ocorrência p_1, p_2, \dots, p_k com $\sum p_i = 1$ para $i = 1, 2, \dots, k$. Seja x_1, x_2, \dots, x_k , as frequências observadas naquelas k categorias em n tentativas, com $n = \sum x_i$, para $i = 1, 2, \dots, k$. Então a função de verossimilhança é:

$$L(p) = (p_1)^{x_1} \cdot (p_2)^{x_2} \cdot (p_3)^{x_3} \dots (p_k)^{x_k} \quad (1)$$

Para maximizar $L(p)$, por método de cálculo, deve-se considerar que existe k-1 parâmetros independentes, pois $\sum p = 1$.

Assim é conveniente escolher p_k como parâmetro a ser expresso em termos dos parâmetros remanescentes.

Aplicando logarítmos naturais a ambos os membros da Eq. (1), obtém-se:

$$\ln L(p) = x_1 \cdot \ln(p_1) + x_2 \cdot \ln(p_2) + \dots + x_k \cdot \ln(p_k) \quad (2)$$

Derivando-se a Eq. (2) em relação a p , obtém-se:

$$\frac{\ln L(p)}{p_i} = \frac{x_1}{p_1} + \frac{x_2}{p_2} + \frac{x_3}{p_3} + \dots + \frac{x_k}{p_k} = \frac{x_i}{p_i}$$

Como só existem $k - 1$ parâmetros independentes, deve-se eliminar o termo $\frac{x_k}{p_k}$. Então:

$$\frac{\ln L(p)}{p_i} = \frac{x_i}{p_i} - \frac{x_k}{p_k} \quad (3)$$

Para maximizar é necessário que $\frac{\ln L(p)}{p_i}$ seja igual a zero, ou seja, todas as $k - 1$ derivadas parciais desapareçam. Então:

$$\frac{x_i}{p_i} - \frac{x_k}{p_k} = 0, \text{ para } i = 1, 2, \dots, k-1$$

Logo:

$$p_i = \frac{p_k}{x_k} \cdot x_i, \text{ para } i = 1, 2, \dots, k \quad (4)$$

Aplicando-se o somatório a ambos os membros da Eq. (4) e considerando que $\sum p = 1$, tem-se:

$$1 = \sum \frac{p_k}{x_k} \cdot x_i \iff 1 = \frac{p_k}{x_k} \sum x_i$$

$$\text{Como } \sum x_i = n \iff 1 = \frac{p_k}{x_k} \cdot n, \text{ então:}$$

$$\frac{p_k}{x_k} = \frac{1}{n} \quad (5)$$

Substituindo-se a Eq. (5) na Eq. (4), obtém-se:

$$p_i = \frac{1}{n} \cdot x_i \iff p_i = \frac{x_i}{n}, \text{ } i = 1, 2, \dots, k \quad (6)$$

Considerando o teste da razão de verossimilhança para testar a hipótese:

$$H_0 : p_i = p_{i0}, \text{ } i = 1, 2, \dots, k$$

Como não existe nenhum parâmetro não específico remanescente quando H_0 é verdadeiro, segue-se que a razão de verossimilhança é dada por:

$$\lambda = \frac{L(\hat{p})}{L(p)} = \frac{(\hat{p}_1)^{x_1} \cdot (\hat{p}_2)^{x_2} \cdot \dots \cdot (\hat{p}_k)^{x_k}}{(p_1)^{x_1} \cdot (p_2)^{x_2} \cdot \dots \cdot (p_k)^{x_k}} \quad (7)$$

pode-se observar que a Eq. (7) pode ser expressa como sendo a razão de duas distribuições multinomiais, então:

$$\lambda = \frac{\frac{n!}{x_1! x_2! \dots x_k!} (\hat{p}_1)^{x_1} \cdot (\hat{p}_2)^{x_2} \cdot \dots \cdot (\hat{p}_k)^{x_k}}{\frac{n!}{x_1! x_2! \dots x_k!} (p_1)^{x_1} \cdot (p_2)^{x_2} \cdot \dots \cdot (p_k)^{x_k}}$$

De (8), tem-se:

$$\lambda = \prod_{i=1}^k \left[\frac{\hat{p}_i}{p_i} \right]^{x_i} \quad (9)$$

Se as proporções p e \hat{p} forem iguais, então a razão será igual a 1. Quanto maior a diferença entre p e \hat{p} , com maior intensidade λ será diferente de 1. Isto indica que a razão dessas probabilidades pode ser usada como estatística para medir a conformidade entre as frequências observadas e esperadas. Um teste baseado em tal razão é denominado de teste de razão de probabilidade.

Sendo a frequência esperada $e_i = n \cdot \hat{p}_i$ e a frequência observada $x_i = n \cdot p_i$, então:

$$\hat{p}_i = \frac{e_i}{n} \quad \text{e} \quad p_i = \frac{x_i}{n}$$

substituindo estes resultados na eq. (9), tem-se:

$$\lambda = \prod_{i=1}^k \left[\frac{e_i}{x_i} \right]^{x_i} \quad (10)$$

Como a distribuição teórica dessa razão é complexa e pouco conhecida, considera-se o seguinte teorema "Sob certas condições de regularidade, a variável aleatória $-2 \ln \lambda$, onde λ é dada por (7), tem uma distribuição que se aproxima da variável X^2 , quando n é finito, com graus de liberdade iguais ao número de parâmetros que são determinados pela hipótese H_0 " [8], ou seja, os graus de liberdade para o teste são iguais ao do teste de qui-quadrado.

$$-2 \ln \lambda = -2 \cdot \sum_{i=1}^k x_i \ln \left[\frac{e_i}{x_i} \right] = 2 \cdot \sum_{i=1}^k x_i \ln \left[\frac{e_i}{x_i} \right]^{-1}$$

Logo:

$$G = 2 \cdot \sum_{i=1}^k x_i \ln \left[\frac{x_i}{e_i} \right] \quad (11)$$

Esta estatística é pouco conhecida e utilizada devido a ser mais trabalhoso aplicar um teste utilizando logaritmos, mas com o uso do computador torna-se mais fácil sua aplicação.

Em geral, G será numericamente similar a X^2 , sendo que o emprego do teste G requer nenhuma sofisticação matemática especial e tem algumas vantagens consideráveis em projetos mais complexos sobre os testes convencionais que utilizam a estatística X^2 , pois permite certos tipos de análises mais detalhadas que são possíveis usando-se a estatística de qui-quadrado [13].

Mesmo com o uso de logaritmos para o teste G , este muitas vezes requer muito menos esforço computacional do que é requerido para o teste de X^2 , isto é verdadeiro para alguns testes [13].

Portanto, pode-se empregar o teste G , sempre que puder ser utilizado um computador, pois os resultados numéricos entre G e X^2 podem diferir em alguns casos, sugerindo uma análise mais cuidadosa dos dados coletados [13].

2- ORIGEM DA ESTATÍSTICA QUI-QUADRADO A PARTIR DA ESTATÍSTICA G

Em 1960, Karl Pearson propôs o seguinte teste estatístico, o qual é função dos quadrados dos desvios das frequências observadas, com relação aos valores esperados respectivos.

$$X^2 = \sum_{i=1}^k \frac{(x_i - e_i)^2}{e_i} \quad (12)$$

Pode-se demonstrar a origem da estatística X^2 a partir da estatística G , da seguinte maneira:

Considerando-se a Estatística G , a Eq. (11), e fazendo-se $y_i = x_i - e_i$, a diferença entre as frequências observadas e esperadas na i -ésima cela, tem-se:

$$G = 2 \sum (e_i + Y_i) \ln \frac{e_i + Y_i}{e_i} = 2 \sum (e_i + Y_i) \ln \left[1 + \frac{Y_i}{e_i} \right]$$

Expandindo $\ln \left[1 + \frac{Y_i}{e_i} \right]$ em Série de Taylor em torno de $x = 0$ obtém-se:

$$G = 2 \sum (Y_i + e_i) \cdot \left[\frac{Y_i}{e_i} - \frac{1}{2} \left[\frac{Y_i}{e_i} \right]^2 + \frac{1}{3} \left[\frac{Y_i}{e_i} \right]^3 - \dots \right]$$

Efetuando, resulta:

$$G = \sum \frac{Y_i^2}{e_i} - \frac{1}{3} \sum \frac{Y_i^3}{e_i} + \dots \quad (13)$$

A variável y é uma variável binomial com média $\mu_i = n \cdot p_i = e_i$ e variância $\sigma_i^2 = np_i(1-p)$; conseqüentemente, a variável $\frac{Y_i}{e_i}$ pode ser expressa na forma:

$$\frac{Y_i}{e_i} = \frac{x_i - e_i}{e_i} = \frac{x_i - \mu_i}{\sigma_i} \sqrt{(1 - p_i)/np_i} \quad (14)$$

Do Teorema "Se X é normalmente distribuída com média e variância σ^2 e uma amostra aleatória de tamanho n é retirada, então a média amostral \bar{X} será normalmente distribuída com média μ e variância $\frac{\sigma^2}{n}$ " [8], a variável $(x_i - \mu_i)/\sigma_i$ tem distribuição aproxima-

da à de uma variável normal padrão quando $n \rightarrow \infty$. Sendo assim o fator da raiz quadrada da Eq. (14) se aproxima de zero na mesma taxa que $1/\sqrt{n}$. Então, para n muito grande y_i/e_i será muito pequeno e da ordem de $1/\sqrt{n}$. Como conseqüência os termos sucessivos na expansão anterior serão da ordem de $1/\sqrt{n}$ vezes o termo precedente.

O resultado aproximado de $-2 \ln \lambda$ para amostras grandes é dado pelo 1º termos da direita da Eq. (13). Então:

$$-2 \ln \lambda \approx \sum \frac{Y_i^2}{e_i} = \sum_{i=1}^k \frac{(x_i - e_i)^2}{e_i} \quad (15)$$

Como $-2 \ln \lambda$ possui distribuição X^2 aproximada, então o termo da direita da Eq. (15) também possui distribuição X^2 .

3- CONCLUSÕES:

Através da literatura citada chegou-se as seguintes conclusões:

- Em projetos mais complexos a aplicação da estatística G tem algumas vantagens em relação a aplicação da estatística X^2 , pois possibilita uma análise mais detalhada dos dados;

- A estatística G é pouco utilizada devido, principalmente, a falta de informações sobre o assunto causada pela escassa li-

teratura e também pela maior dificuldade de cálculos pelo emprego de logarítmos, o que no momento não é mais problema com o uso do computador.

4- REFERÊNCIAS BIBLIOGRÁFICAS

01. AYRES, M. & AYRES JUNIOR, M Aplicações Estatísticas em Basic. São Paulo, Mc Graw-Hill, 1987.
02. BISHOP, Y. M. et alii. Discrete Multivariate Analysis: Theory and Practice. The mit Press, 1975.
03. CHRISTMANN, Raul U. Estatística Aplicada. 2ª ed. São Paulo, Edgard Blücher Ltda, 1978.
04. CONOVER, W. J. Practical Nonparametric Statistic. New York, John Wiley & Sons, 1971.
05. COSTA NETO, Pedro L. de O. Estatística. São Paulo, Edgard Blücher Ltda, 1977.
06. FONSECA, Jairo S. da, MARTINS, Gilberto de A. Curso de Estatística. 3ª ed. São Paulo, Atlas, 1982.
07. GOODMAN, Richard. Estatística. São Paulo, Livraria Editora Pioneira, 1965.
08. HOEL, Paul G. Estatística Matemática. 4ª ed. Rio de Janeiro, Guanabara Dois S.A., 1980.
09. MENDENHALL, William. Probabilidade e Estatística. Rio de Janeiro, Campus, 1985, v. 2.
10. RODRIGUES, Aroldo. A Pesquisa Experimental em Psicologia e Educação. 2ª ed. Rio de Janeiro, Vozes, 1976.
11. SCREMIN, Marcos A. A. Comparação Entre os Testes G e Qui-Quadrado Através de um Programa Computacional. Monografia do Curso de Pós-Graduação em Métodos Quantitativos, Universidade Federal de Santa Maria, Santa Maria-RS, 1989.
12. SIEGEL, Sidney. Estatística Não-paramétrica. São Paulo, McGraw-Hill do Brasil, Ltda, 1975.
13. SOKAL, R. and ROHLF, F. J. BIOMETRY. San Francisco, Freeman and Company, 1969.
14. STEVENSON, Willian J. Estatística Aplicada à Administração. São Paulo, Harper & Row do Brasil Ltda, 1981.
15. YAMANE, Taro. Estatística. 3ª ed. México, Harla, 1974.

Recebido em outubro, 1990; aceito em abril, 1991.

