

# Numerical Studies on Convergence of Multinomial Goodness-of-fit Statistics to Chisquare Distribution

Zh. Assylbekov

*Graduate School of Science, Hiroshima University*

October 1, 2008

## Abstract

Let  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_k)'$  be a random vector with multinomial distribution. In this report we investigate numerically the convergence rate of so-called power divergence family of statistics  $\{I^\lambda(\mathbf{Y}), \lambda \in \mathbb{R}\}$  introduced by Cressie and Read (1984) to chi-square distribution for  $k = 4, 5, 6$ .

## 1 Introduction

Let  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_k)'$  be a random vector with the multinomial distribution  $M_k(n, \boldsymbol{\pi})$ , i.e.,

$$\Pr(Y_1 = n_1, Y_2 = n_2, \dots, Y_k = n_k) = \begin{cases} n! \prod_{j=1}^k \frac{\pi_j^{n_j}}{n_j!} & \sum_{j=1}^k n_j = n \\ 0 & \text{otherwise,} \end{cases}$$

where  $n_j = 0, 1, \dots, n$ ,  $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_k)'$ ,  $\pi_j > 0$ ,  $\sum_{j=1}^k \pi_j = 1$ . For testing the simple hypothesis  $H : \boldsymbol{\pi} = \mathbf{p}$  ( $\mathbf{p}$  is a fixed vector) against  $K : \boldsymbol{\pi} \neq \mathbf{p}$  the power divergence statistics (introduced by Cressie and Read in [2]) can be used:

$$2nI^\lambda = \frac{2}{\lambda(\lambda + 1)} \sum_{j=1}^k Y_j \left( \left( \frac{Y_j}{np_j} \right)^\lambda - 1 \right), \quad \lambda \in \mathbb{R},$$

where  $\mathbf{p} = (p_1, p_2, \dots, p_k)'$ ,  $p_j > 0$  ( $j = 1, 2, \dots, k$ ) and  $\sum_{j=1}^k p_j = 1$ .

Throughout this paper we will use the following notation:

$$\begin{aligned}\mathbf{x} &= (x_1, \dots, x_r)', \\ \mathbf{x}^* &= (x_1, \dots, x_{l-1}, x_{l+1}, \dots, x_r)',\end{aligned}$$

For any  $B \subset \mathbb{R}^r$  and for any  $l \in \{1, \dots, r\}$  denote

$$B_l = \{\mathbf{x}^* : \mathbf{x} \in B\}.$$

**Definition 1.** A set  $B \subset \mathbb{R}^r$  is called an *extended convex set*, if  $B$  has the following representation for every  $l \in \{1, 2, \dots, r\}$ :

$$B = \{\mathbf{x} : \lambda_l(\mathbf{x}^*) < x_l < \theta_l(\mathbf{x}^*), \mathbf{x}^* \in B_l\},$$

where  $\lambda_l, \theta_l$  are continuous functions on  $B_l$ .

It is known (see Cressie, Read [2]), that under the null hypothesis  $2nI^\lambda$  has the chisquare distribution with  $r = k - 1$  degrees of freedom in the limit. Moreover the distribution function of  $2nI^\lambda$  has the following expansion:

$$\Pr(2nI^\lambda < c) = \Pr(\chi_r^2 < c) + J_2 + O(n^{-1}), \quad (1)$$

where

$$\begin{aligned}J_2 &= -\frac{1}{\sqrt{n}} \sum_{l=1}^r n^{-\frac{r-l}{2}} \sum_{x_{l+1} \in L_{l+1}} \dots \sum_{x_r \in L_r} \\ &\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \chi_{B_l^\lambda(\mathbf{x}^*)} [S_1(\sqrt{n}x_l + pln) \phi(\mathbf{x})]_{\lambda_l(\mathbf{x}^*)}^{\theta_l(\mathbf{x}^*)} dx_1 \dots dx_{l-1}, \quad (2)\end{aligned}$$

$$L_j = \{x_j : x_j = \frac{1}{\sqrt{n}}(n_j - np_j), n_j \in \mathbb{Z}\}, \quad (3)$$

$$S_1(x) = x - [x] - \frac{1}{2},$$

$$\begin{aligned}[h(\mathbf{x})]_{\lambda(\mathbf{x}^*)}^{\theta(\mathbf{x}^*)} &= h(x_1, \dots, x_{l-1}, \theta_l(\mathbf{x}^*), x_{l+1}, \dots, x_r) \\ &- h(x_1, \dots, x_{l-1}, \lambda_l(\mathbf{x}^*), x_{l+1}, \dots, x_r),\end{aligned}$$

$$\phi(\mathbf{x}) = (2\pi)^{-\frac{r}{2}} |\Omega|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \mathbf{x}' \Omega^{-1} \mathbf{x}\right),$$

$$\Omega = \text{diag}(p_1, \dots, p_r) - (p_1, \dots, p_r)'(p_1, \dots, p_r).$$

Here  $\chi_A(x)$  is an indicator function,  $\theta_l(\mathbf{x}^*)$  and  $\lambda_l(\mathbf{x}^*)$  are continuous functions from Definition 1 for the set

$$B^\lambda = \{\mathbf{x} : 2nI^\lambda(\mathbf{x}) < c\} \quad (4)$$

with

$$2nI^\lambda(\mathbf{x}) = \frac{2}{\lambda(\lambda+1)} \sum_{j=1}^k (np_j + \sqrt{n}x_j) \left( \left(1 + \frac{x_j}{\sqrt{np_j}}\right)^\lambda - 1 \right), \quad (5)$$

$$x_k = -(x_1 + \cdots + x_r).$$

It follows from Yarnold's result [5] that

$$J_2 = O(n^{-1/2}).$$

Zubov and Ulyanov in [6] showed that

$$J_2 = O(n^{-1+\frac{1}{r+1}}).$$

This was improved by the author in [1], where it was shown that

$$J_2 = O(n^{-1+\mu(r)}),$$

with

$$\mu(r) = \begin{cases} 6/(7r+4) & \text{for } 3 \leq r \leq 7, \\ 5/(6r+2) & \text{for } r \geq 8. \end{cases}$$

In the present report we investigate numerically whether the upper bound for  $J_2$  can be improved.

## 2 Preliminaries

By definition put

$$L = \left\{ \mathbf{x} : x_j = \frac{1}{\sqrt{n}}(m_j - np_j), m_j \in \mathbb{Z}, j = \overline{1, r} \right\},$$

This means that  $L$  is an  $r$ -dimensional lattice in  $\mathbb{R}^r$  and lattice spacing of  $L$  is  $\frac{1}{\sqrt{n}}$ . Let  $N^\lambda$  be the number of lattice points in the ellipsoid  $B^\lambda$ , i.e,  $N^\lambda = \#(L \cap B^\lambda)$ . Let  $V^\lambda$  be the volume of  $B^\lambda$ .

**Lemma 1.** *Let  $J_2$  be the term defined by (2); then*

$$J_2 = dn^{-\frac{r}{2}}(N^\lambda - n^{\frac{r}{2}}V^\lambda) + O(n^{-1}), \quad (6)$$

where  $d = \text{const} > 0$ .

*Proof.* The proof is given in [1], Proposition 1. □

### 3 Numerical studies

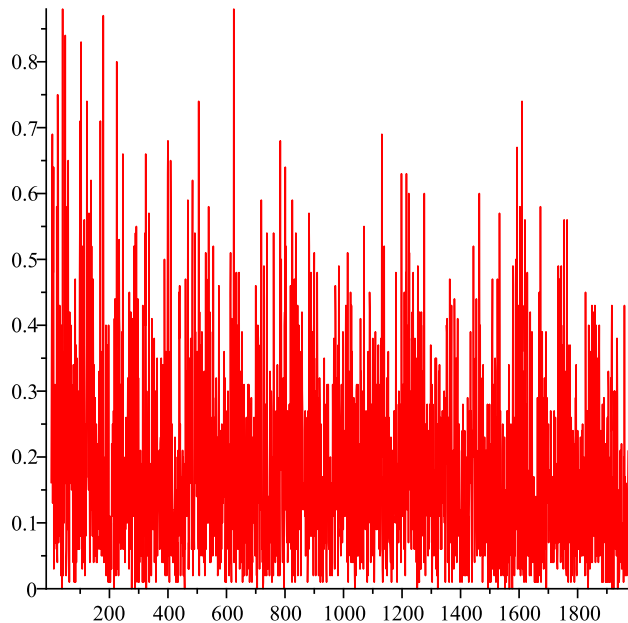
In [1] it was shown that  $B^\lambda$  is a convex body which has smooth boundary with nonvanishing and bounded Gaussian curvature throughout. Hence Lemma 1 reduces the original problem of estimating  $J_2$  (see (2)) to the lattice point problem inside convex body  $B^\lambda$  (see (6)). In [4], W. Müller made a conjecture, which in terms of our problem reads

$$n^{-\frac{r}{2}}(N^\lambda - n^{\frac{r}{2}}V^\lambda) = \begin{cases} O(n^{-1+\varepsilon}), & r = 3, 4, \\ O(n^{-1}), & r \geq 5. \end{cases}$$

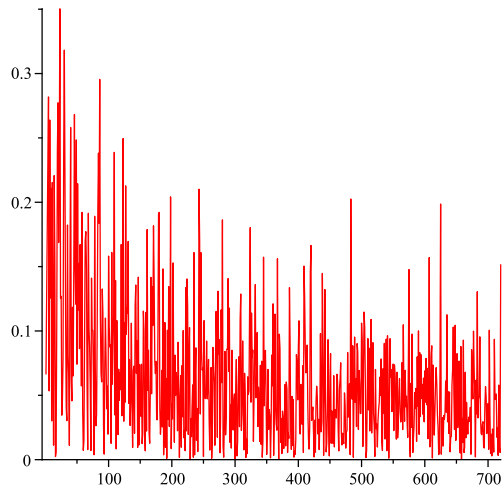
In order to test this conjecture we calculate the expressions

$$\begin{cases} n^{-\frac{r}{2}} |N^\lambda - n^{\frac{r}{2}}V^\lambda| n^{0.9}, & r = 3, 4, \\ n^{-\frac{r}{2}} |N^\lambda - n^{\frac{r}{2}}V^\lambda| n, & r = 5. \end{cases} \quad (7)$$

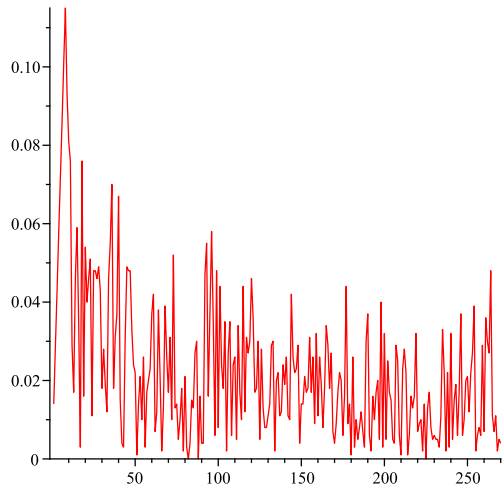
If Müller's conjecture is true then the expressions in (7) should be bounded by above. The results of our computations for  $r = 3$  and  $n = \overline{1, 2000}$  are given below



For  $r = 4$  and  $n = \overline{1,720}$  we have



And, finally, for  $r = 5$  and  $n = \overline{1,270}$  we have



As it is seen on these pictures, our results fit Müller's conjecture. The algorithm of computations is rather straightforward and is given for  $r = 3$  in Appendix A.

## A Algorithm of computations in C

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#define C 1.0
#define PI 3.14159265358979323846
#define K 4

double T(int n, double *x, double *p)
{
    double result = 0;
    int j;
    x[K-1] = 0;
    for (j = 0; j < K - 1; j++)
    {
        result += 2 * (n*p[j] + sqrt(n)*x[j])*log(1 + x[j]/sqrt(n)/p[j]);
        x[K-1] -= x[j];
    }
    result += 2 * (n*p[K-1] + sqrt(n)*x[K-1])*log(1 + x[j]/sqrt(n)/p[K-1]);
    return result;
}

int N(int n, double *p)
{
    double x[K];
    int i, j, k, count = 0;
    for (i = 1; i < n; i++)
        for (j = 1; j < n - i; j++)
            for (k = 1; k < n - i - j; k++)
            {
                x[0] = 1/sqrt(n)*(i - n*p[0]);
                x[1] = 1/sqrt(n)*(j - n*p[1]);
                x[2] = 1/sqrt(n)*(k - n*p[2]);
                if (T(n,x,p) < C)
                {
                    count++;
                }
            }
    return count;
}
```

```

double V(int n, double *p)
{
    return 4*pow(PI*C, 3.0/2)*sqrt(p[0]*p[1]*p[2]*p[3])/3/sqrt(PI);
}

int main(int argc, char *argv[])
{
    double p[K], volume;
    int n;

    p[0] = 0.1;
    p[1] = 0.1;
    p[2] = 0.3;
    p[K-1] = 0.5;

    volume = V(n,p);

    for(n = 1; n < 1000; n++)
    {
        printf("%f\n", (pow(n, -3.0/2)*N(n, p)-volume));
        fprintf(stderr, "%d ", n);
    }

    return 0;
}

```

## References

- [1] ZH. ASSYLBEKOV, Convergence of multinomial goodness-of-fit statistics to chisquare distribution, submitted.
- [2] N. C. CRESSIE AND T. R. C. READ, Multinomial goodness-of-fit tests, *J. R. Statist. Soc. B* (1984) 46, No. 3, 440-464.
- [3] E. HLAWKA, Über Integrale auf konvexen Körpern I, II, *Monatsh. Math.* 54, 1-36, 81-99 (1950).
- [4] W. MÜLLER, Lattice points in large convex bodies, *Mh. Math.* 128, 315–330 (1999).
- [5] J. K. YARNOLD, Asymptotic approximations for the probability that a sum of lattice random vectors lies in a convex set, *The Annals of Mathematical Statistics* 1972, Vol. 43, No. 5, 1566–1580.
- [6] V. N. ZUBOV, V. V. ULYANOV, Refinement on the convergence of one family of goodness-of-fit statistics to chi-squared distribution, submitted.