

Czesław Domąński*, Andrzej S. Tomaszewicz**

NORMAL APPROXIMATION OF MULTIPLE RUNS DISTRIBUTIONS

Tests based on the number of runs can be applied to verification of many types of hypotheses, for instance the hypothesis that

- elements of the sample are independent,
- two or more samples are drawn from the same distribution,
- regression model of one or two explanatory variables is linear, etc.

Let us consider a sequence of random variables

$$X_1, X_2, \dots, X_n \quad (1)$$

which have the same discrete distribution with s values, i.e.

$$P(X_1 = x_j) = p_j, \quad j = 1, \dots, s, \quad \sum_{j=1}^s p_j = 1.$$

We also use the symbols n_j to denote the number of elements of j -th type in (1):

$$n_j = \text{card } \{i : X_i = x_j\}.$$

The object of our consideration is the number of runs in the sequence (1):

$$R_n = 1 + \text{card } \{i : 2 \leq i \leq n, X_i \neq X_{i-1}\}.$$

For instance:

- for $n = 8$, $s = 2$, using the traditional notation: $x_1 = A$, $x_2 = B$ we can have the sequence

* Associate Professor at the University of Łódź.

** Associate Professor at the University of Łódź.

A A B A A B B A

when we observe 5 runs,

- for $n = 9$, $s = 3$, $x_3 = C$ - we also have 5 runs in the sequence

A B B C C C A A C

The critical region for the above specified types of hypotheses is based on the distribution of the number of runs in (1) under the assumption that X_i 's are independent. This distribution for two kinds of elements was investigated by Stevens (1939), Mood (1940), Wald and Wolfowitz (1940), Sweden and Eisenhart (1943) and, for three and more kinds of elements - Barton and David (1957, 1960). Some new results could be found in monographies Barton (1960) and Gibbons (1987). They showed some combinatorial formulae for number-of-run distribution. Besides, the convergence to the normal distribution was proved.

Tables of number-of-run distribution for more than two kinds of elements were constructed just for small values of n (for instance, for 3 and 4 kinds of elements - up to $n = 12$ only). It was because the classical formulae are not convenient for computations. For $s = 2$ we published quite general recursive formulae (cf. Domański and Tomaszewicz (1984)) which can be used when the sequence (1) is generated by stationary Markov chain.

The number of runs is a discrete variable. Therefore (excluding exceptional cases) it is not possible to select such a critical value to make the test size exactly equal to fixed significance level α . Commonly the critical values as integer numbers

$$k^L(n, \alpha) = \max \{k: P(K_n \leq k) \leq \alpha\},$$

$$k^R(n, \alpha) = \min \{k: P(K_n \geq k) \leq \alpha\}$$

are accepted (the first of them concerns a left-hand-sided test, the second one - a right-hand-sided test. Hence, the test size is, in general, less than α :

$$P(K_n \leq k^L(n, \alpha)) < \alpha,$$

$$P(K_n \geq k^R(n, \alpha)) < \alpha.$$

For that reason the randomized tests are applied.

The randomized test we used is defined in common way:
reject H_0 when

$$K_n \leq k^L(n, \alpha) \text{ (or } K_n \geq k^R(n, \alpha)),$$

accept H_0 when

$$K_n > k^L(n, \alpha) + 1 \text{ (or } K_n < k^R(n, \alpha) - 1),$$

reject H_0 with the probability

$$p_{\text{rand}}^L(n, \alpha) = \frac{\alpha - P(K_n \leq k^L(n, \alpha))}{P(K_n = k^L(n, \alpha) + 1)}, \quad \text{when } K_n = k^L(n, \alpha) + 1$$

(or reject H_0 with the probability

$$p_{\text{rand}}^R(n, \alpha) = \frac{\alpha - P(K_n \geq k^R(n, \alpha))}{P(K_n = k^R(n, \alpha) - 1)}, \quad \text{when } K_n = k^R(n, \alpha) - 1.$$

The idea of interpolated quantiles is connected with randomized tests. We define

$$k_i^L(n, \alpha) = k^L(n, \alpha) + p_{\text{rand}}^L(n, \alpha),$$

$$k_i^R(n, \alpha) = k^R(n, \alpha) - p_{\text{rand}}^R(n, \alpha).$$

Of course, when we know the interpolated quantile we know both the integer quantile and the randomization probability.

Using normal approximation we obtain following estimates of interpolated quantiles

$$k_i^L(n, \alpha) \approx \hat{k}_i^L(n, \alpha) = \mu_K(n) + \phi^{-1}(\alpha)\sigma_K(n) - \frac{1}{2},$$

$$k_i^R(n, \alpha) \approx \hat{k}_i^R(n, \alpha) = \mu_K(n) - \phi^{-1}(\alpha)\sigma_K(n) + \frac{1}{2}.$$

ϕ denotes the inverse of standard normal cumulative distribution function and

$$\mu_K(n) = E(K_n), \quad \sigma_K(n) = \text{var}(K_n).$$

Randomized test based on these quantiles are defined using estimates of the integer quantiles and the randomization probabilities:

$$\hat{k}^L(n, \alpha) = \text{entier}(\hat{k}_i^L(n, \alpha)),$$

$$\hat{p}_{\text{rand}}^L(n, \alpha) = \hat{k}_i^L(n, \alpha) - \hat{k}^L(n, \alpha),$$

$$\hat{k}_i^R(n, \alpha) = -\text{entier}(-\hat{k}_i^L(n, \alpha)),$$

$$\hat{p}_{\text{rand}}^R(n, \alpha) = \hat{k}_i^R(n, \alpha) - \hat{k}_i^L(n, \alpha),$$

Thus, the size of the test based on normal approximation is

$$a^L(n, \alpha) = P(K_n \leq \hat{k}_i^L(n, \alpha)) + p_{\text{rand}}^L(n, \alpha) P(K_n^L = k(n, \alpha) + 1).$$

$$a^R(n, \alpha) = P(K_n \geq \hat{k}_i^R(n, \alpha)) + p_{\text{rand}}^R(n, \alpha) P(K_n^R = k(n, \alpha) - 1).$$

(These formulae follow immediately from the randomized test definition).

As a measure of goodness of approximation we chose the difference between the test size and the assumed significance level α :

$$\delta^L(n, \alpha) = \hat{a}^L(n, \alpha) - \alpha,$$

$$\delta^R(n, \alpha) = \hat{a}^R(n, \alpha) - \alpha.$$

It is hard to deny that the smaller the value of this difference, the better approximation.

We would like to present some evaluation of normal approximation of tests based on number-of-runs distribution in symmetric case for 3 types of elements. That means we assume that

$$s = 3,$$

and runs with different lengths are equal, i.e.

$$p_1 = p_2 = p_3 = \frac{1}{3}$$

and quantiles depend on three numbers n_1, n_2, n_3

$$n_1 = n_2 = n_3$$

are fixed. The recursive formula

$$P(K_n = k | n_1, n_2, n_3) = \frac{n_1}{n} (P(K_{n-1} = k | n_1 + 1, n_2, n_3) + \\ + P(K_{n-1} = k-1 | n_2, n_1 + 1, n_3) + P(K_{n-1} = k - 1 | n_3, \\ n_1 + 1, n_2)).$$

was used to compute the distribution. We present some results for the special cases:

$$n_1 = n_2 = n_3 = 1, 2, \dots, 40 \quad (n = 3, 6, \dots, 120).$$

Table 1 contains the interpolated quantiles k_i^L, k_i^R , their normal approximates \hat{k}_i^L, \hat{k}_i^R , and test size errors δ^L, δ^R . The dif-

Table 1

Quantiles (k^L , \hat{k}^L) and test size errors (a)

n	k_i^L , \hat{k}_i^L , α	Left-hand sided			Right-hand sided		
		$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.10$	$\alpha = 0.10$	$\alpha = 0.05$	$\alpha = 0.01$
15	$k_i^L(n, \alpha)$	6.267	7.534	9.077	12.956	14.209	14.990
	$\hat{k}_i^L(n, \alpha)$	6.568	7.720	9.077	12.923	14.280	15.432
	$\alpha^L(n, \alpha)$	0.0052	0.0095	0.0000	-0.0070	0.0037	0.0046
30	$k_i^L(n, \alpha)$	14.416	16.260	18.380	23.653	25.556	26.958
	$\hat{k}_i^L(n, \alpha)$	14.704	16.402	18.403	23.597	25.598	27.296
	$\alpha^L(n, \alpha)$	0.0029	0.0063	0.0025	-0.0068	0.0019	0.0029
60	$k_i^L(n, \alpha)$	31.940	34.460	37.467	44.567	47.338	49.615
	$\hat{k}_i^L(n, \alpha)$	32.150	34.600	37.479	44.521	47.403	49.849
	$\alpha^L(n, \alpha)$	0.0018	0.0038	0.0010	-0.0037	0.0018	0.0020
90	$k_i^L(n, \alpha)$	50.008	53.123	56.778	65.257	68.702	71.552
	$\hat{k}_i^L(n, \alpha)$	50.214	53.227	56.779	65.221	68.773	71.786
	$\alpha^L(n, \alpha)$	0.0015	0.0027	0.0000	-0.0022	0.0019	0.0015
120	$k_i^L(n, \alpha)$	68.330	71.992	76.194	85.837	89.837	93.145
	$\hat{k}_i^L(n, \alpha)$	68.588	72.078	76.191	85.809	89.922	93.419
	$\alpha^L(n, \alpha)$	0.0014	0.0019	0.0002	-0.0016	0.0020	0.0012

Source: The author's calculations.

ferences between the quantiles seem to be not large. But, in our opinion, a much better measure of goodness of approximation is the test size error δ^L and δ^R . In Table 2 values of δ^L , δ^R are shown. They are absolute values of δ in some intervals of sample size n:

$$\delta_*^L(n_1, n_2, \alpha) = \max_{n_1 \leq n \leq n_2} |\delta^L(n, \alpha)|,$$

$$\delta_*^R(n_1, n_2, \alpha) = \max_{n_1 \leq n \leq n_2} |\delta^R(n, \alpha)|$$

Table 2

Test size errors

$n_1 - n_2$	Left-hand sided			Right-hand sided		
	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.10$	$\alpha = 0.10$	$\alpha = 0.05$	$\alpha = 0.01$
15-30	0.0052	0.0095	0.0033	0.0084	0.0048	0.0046
33-45	0.0030	0.0049	0.0014	0.0056	0.0048	0.0029
48-60	0.0022	0.0040	0.0010	0.0039	0.0033	0.0025
63-75	0.0019	0.0037	0.0008	0.0035	0.0030	0.0022
78-90	0.0017	0.0031	0.0009	0.0022	0.0020	0.0019
93-105	0.0015	0.0026	0.0003	0.0024	0.0025	0.0017
108-120	0.0014	0.0026	0.0003	0.0022	0.0020	0.0015

Source: The author's calculations.

Some authors suggest that the normal approximation is good enough even for $n = 15$. For very rough statistical investigation maybe it is not a big difference whether the test size is 5% or 6%. But in many kinds of statistical analysis (for instance test power investigation) more accuracy is needed. Thus, the error values 1 or 2 pro mille which we observe for sample size 100-120 is not satisfactory.

Therefore we cannot base only on normal approximation of number-of-runs distribution. Even for moderate sample sizes the exact distribution should be applied.

REFERENCES

- Barton D. E., David F. N. (1957), *Runs Multiple*, "Biometrika", No. 44.
- Barton D. E., David F. N. (1960), *Runs in a Ring*, "Biometrika", No. 45.
- Barton D. E. (1966), *Combinatorial Chance*, Hofner Publishing Company, New York.
- Domański C., Tomaszewicz A. S. (1984), *Recursive Formulae for Runs Distributions*, "Acta Universitatis Lodzienensis", No. 34.
- Gibbons J. (1987), *Nonparametric Statistics Inference*, Mc Graw-Hill Book Company, New York.

- Mood A. H. (1940), *The Distribution Theory of Runs*, "Annals of Mathematical Statistics", No. 11.
- Stevens W. L. (1939), *Distribution of Groups in a Sequence of Alternatives*, "Annals of Eugenics", No. 9.
- Swed F. S., Eisenhart C. (1943), *Tables for Testing Randomness*, "Annals of Mathematical Statistics", No. 14.
- Wald A., Wolfowitz J. (1940), *On a Test Whether Two Samples are From the Same Population*, "Annals of Mathematical Statistics", No. 11.
- Walsh J. E. (1962), *Handbook of Nonparametric Statistics*, D. von Nostrand Co. Inc., Princeton.

Czesław Domaniński, Andrzej S. Tomaszewicz

ROZKŁADY DŁUGOŚCI I LICZBY SERII WIELOKROTNYCH

Rozkłady liczby i długości serii dla dwóch rodzajów elementów zostały stosunkowo dobrze poznane. Znacznie mniej natomiast wiadomo o własnościach rozkładów liczby bądź długości serii dla trzech lub więcej elementów.

W artykule prezentujemy niektóre wyniki dotyczące własności testów opartych na seriach złożonych z trzech lub więcej rodzajów elementów, weryfikujących hipotezę o niezależności obserwacji w próbie. Ze względu na to, że rozkłady badanych statystyk są dyskretne, analizowano testy zrandomizowane i kwantyle interpolowane.