

*Czesław Domański**

THE ESTIMATE OF POWER OF RANDOM TESTS BASED ON LENGTH OF RUNS

Abstract. Tests based on length of runs are used both in statistical inference and statistical quality assurance. The paper presents power of three tests based on:

- maximum run length on one side of the median,
- smaller from the maximum run length above and beneath the median,
- bigger from maximum run length above and beneath the median.

Key words: run tests, number and run length distributions, power of tests based on run length.

1. INTRODUCTION

The analysis of statistical tests properties from applicational point of view usually involves examining their power or resistance.

The problem of non-parametric tests power examination is relatively difficult because of the lack of general theory for this branch of science.

For the last 25 years non-parametric tests power rarely has been tested analytically but mostly by means of numerical-simulation methods (Monte Carlo) both mentioned methods were used. Generally, we arbitrarily formulate a list of some alternative hypotheses (it usually consists of the most common hypotheses in practice of statistical research). Then, we calculate the frequencies of rejection of the null hypothesis on the ground of a number generated samples, which fulfil the assumptions given by alternative hypotheses. These frequencies are the empirical power of the given test.

The paper presents results referring to three randomness tests based on:

- maximum run length on one side of the median,
- smaller from the maximum run length above and beneath the median,
- bigger from maximum run length above and beneath the median.

The aim of this paper is to formulate some conclusions about the power of randomness tests based on run length. The conclusions will be helpful

* Prof., Department of Statistical Methods, University of Łódź.

in choosing run test in practical applications. Formulated conclusions are presented for the case of Markov stationary chain of two states.

2. PROBLEM FORMULATION

Cz. Domański (1986) presented some conclusions referring to three tests to verify the hypothesis saying that the sequence of consecutive observations in sample is independent. They are based on:

- maximum run length on one side of the median (S_A),
- smaller from maximum run length above and beneath of the median (S_D),
- bigger from maximum run length above and beneath of the median (S_G).

To make the practical application of these tests possible there were given the tables of critical values for series tests S_A , S_D , S_G and an attempt to assess association between considered statistics was made.

The aim of this paper is to formulate some conclusions referring to the power of most often used run tests. We hope these conclusions will be helpful in choosing run tests in practical applications. We are confined to the case of Markov stationary chain of two states traditionally marked by A and B and the transition matrix.

$$\begin{bmatrix} p_{AA} & p_{AB} \\ p_{BA} & p_{BB} \end{bmatrix} = \begin{bmatrix} 1 - q_0 & q_0 \\ q_1 & 1 - q_1 \end{bmatrix}.$$

Let $P_{n\theta}$ be the distribution of this chain for each

$$\theta \in \Theta = \{ \{q_0, q_1\} : 0 < q_0 < 1, 0 < q_1 < 1 \}$$

and let $\Omega_n = \{A, B\}^n$ be the set of all n -element sequences made of elements A, B . We will consider the probability space

$$M_{n,\theta} = (\Omega_n, 2^{\Omega_n}, P_{n\theta}) \quad \text{for } \theta \in \Theta.$$

Conclusions formulated in the final part of this paper are based on the power of tests numerically determined for $n = 1, 2, \dots, 100$ and a few dozen pairs chosen from the set Θ . Formulating the most suitable algorithm for these calculations was the necessary stage of this research.

Until now, most often we have been considering the run elements above and beneath the median i.e. the case when the probability of all kind of elements is 0.5, if we follow the Bateman's (Bateman 1948) conclusions

who showed that in this case the power of independence test is the biggest. The results presented in article refer to Markov chain with arbitrary stationary probabilities $0 < p_A < 1$.

Combinatorial formulas on probabilities connected with run distribution are very unsuitable for numerical calculations. More effective is using recurrent formulas, particularly when calculations are made for consecutive values of n (sample size).

Recurrent formulas referring to length run distribution in the case when the consecutive observations in sample are generated by the Markov stationary chain of two states were presented in Domański's paper (Domański 1986). Let us observe that it is fundamental to compare the power of tests based on run length with the power of tests including run number. Therefore, we give recurrent formulas to calculate the function of total number and series length probability function. In studies that have been made so far, only formulas for one-dimensional distributions were used.

3. RECURRENT FORMULA FOR THREE-DIMENSIONAL DISTRIBUTION OF RUNS

Let us assign to every sequence

$$\omega = (x_1, x_2, \dots, x_n) \in \Omega_n$$

the following numbers:

$N_A(\omega)$ – number of A elements in sequence ω ,

$L_A(\omega)$ – number of run made of A elements,

$L(\omega)$ – total number of runs,

$S_A(\omega), S_B(\omega)$ – maximum run length composed of A elements (respectively B),

$K_A(\omega), K_B(\omega)$ – number of elements A (respectively B) located at the end of sequence ω ,

$Z_A(\omega), Z_B(\omega)$ – maximum run length composed of elements A (respectively B) without taking into consideration the last run.

Let us assume that sequences $\omega \in \Omega_n$ are the realizations of Markov stationary chain of transition matrix

$$\begin{bmatrix} p_{AA} & p_{BA} \\ p_{AB} & p_{BB} \end{bmatrix} \quad (1)$$

where $0 < p_{AB}, p_{BA} < 1$.

Therefore stationary probabilities are given by the formula:

$$p_A = P(X_j = A) = \frac{p_{AB}}{p_{AB} + p_{BA}} \quad \text{for } j = 1, 2, \dots, n$$

$$p_B = P(X_j = B) = \frac{p_{BA}}{p_{AB} + p_{BA}} \quad \text{for } j = 1, 2, \dots, n$$
(2)

Taking into account these assumptions, the probability distribution on the set Ω_n can be given by the formula:

$$P(\omega) = \frac{1}{p_{AB} + p_{BA}} p_{AA}^{n_A - l_A} p_{AB}^{l_A} p_{BA}^{l_B - l_A} p_{BB}^{n_B - l_B - l_A} \quad (3)$$

where to simplify things, it was assumed $n_A = N_A(\omega)$, $l = L(\omega)$, $l_A = L_A(\omega)$. In fact,

$$P(\omega) = P(X_1 = x_1)P(X_2 = x_2|X_1 = x_1) \dots P(X_n = x_n|X_{n-1} = x_{n-1}) \quad (4)$$

and at the same time

$$P(X_1 = x_1) = \begin{cases} p_A & \text{if } x_1 = A, \\ p_B & \text{if } x_1 = B. \end{cases}$$

Let us observe that l_A is the number of these A elements, which create new run i.e. they follow B (perhaps without the first element). That is why on the right side (4) there is l_A factor equal to p_{AB} (including also factor (5) in a shape (2) when $x_1 = A$). The number of A elements not creating new run and at the same time following A , amounts $(n_A - l_A)$, so there is the same number of factors p_{AA} on the right side (4). Similarly, we show that numbers of factors p_{BB} and p_{BA} equal respectively $(n_B - l_B)$ and l_B (including also factor (5) in formula (2), if $x_1 = B$). Both if $x_1 = A$, and if $x_1 = B$, on the right side occurs one factor $\frac{1}{p_{AB} + p_{BA}}$.

Let us consider, for fixed n , total three-dimensional distribution (L, S_A, S_B) of number of run L , maximum length of run consisting of A elements and maximum length of run consisting of B elements.

Let us assign

$$M(n, l, s, t, u) = \text{card}\{\omega \in \Omega_n : l = l(\omega), \quad s = Z_A(\omega), \quad t = S_B(\omega), \quad u = K_A(\omega)\}$$
(6)

the following formulas are true:

$$M(n, l, s, t, u) = \begin{cases} M(n-1, l, s, t, u-1) & \text{for } u > 1 \\ M(n-1, l-1, s, t, 0) & \text{for } u = 1 \\ \sum_{v=0}^{t-1} M(n, l, v, s, t) + \sum_{w=1}^t M(n, l, t, s, w) & \text{for } u = 1 \end{cases} \quad (7)$$

The first two formulas are obvious. We get them by adding n -th element A to $(n-1)$ element sequence. In case $u = 0$, transforming element A to B and inversely, we get $M(n, l, s, t, u) = \sum_{v, w} M(n, l, v, s, w)$, where the summation is over these pairs (v, w) for which $\max\{v, w\} = t$.

Initial conditions for formula (7) are as follows:

$$M(1, l, s, t, u) = \begin{cases} 1 & \text{for } l = u = 1, \quad s = t = 0 \\ 0 & \text{for the other} \end{cases} \quad (8)$$

Cz. Domański (1986) proved:

Theorem 1. The total distribution of variables (L, S_A, S_B) specified on probability space $M_{n, \theta}$ takes the following form:

$$P(L = l, S_A = s, S_B = t) = Q_0(n, l, s, t) + Q_1(n, l, t, s) \quad (9)$$

where for $n = 0.1$

$$Q_n(n, l, s, t) = \sum_{v=0}^{t-1} Q_{1-n}(n-t, l-1, v, s) q_n (1 - q_{1-n})^{t-1} + \sum_{w=1}^t Q(n-w, l-1, s, w) q_n (1 - q_{1-n})^{n-1},$$

$$Q_n(0, l, s, t) = \begin{cases} \frac{1 - q_n}{(q_0 + q_1) q_{1-n}} & \text{for } l = s = t = 0 \\ 0 & \text{for the other} \end{cases}$$

with initial conditions for $h = 0.1$.

From Theorem 1 follow the current formulas for one-dimensional distribution probabilities L_A, L_B, L_G .

Theorem 2. Variable distribution L_A specified on space $M_{n,\theta}$ is given by the formula

$$P(S_1 = s) = Q_0^A(n, s) + Q_1^A(n, s)$$

$$Q_0^A(n, s) = \sum_{v=1}^{n-1} Q_1^A(n-v, s) q_0 (1-q_1)^{v-1} \quad (10)$$

$$Q_1^A(n, s) = \sum_{v=0}^{s-1} Q_0^A(n-s, v) q_1 (1-q_0)^{l-1} + \sum_{w=1}^l Q_0^A(n-w, s) q_1 (1-q_0)^{w-1}$$

with initial conditions

$$Q_0^A(0, 0) = Q_1^A(0, 0) = \frac{1}{q_0 + q_1}.$$

Transforming A into B and 0 into 1, we get formula for distribution S_B .

Theorem 3. Distribution of variable S_G specified on $M_{n,\theta}$ is given by the recurrent formula

$$P(S_G = s) = Q_0^G(n, s) + Q_1^G(n, s) \quad (11)$$

where for $h = 0.1$

$$Q_h^G(n, s) = \sum_{v=0}^{s-1} Q_{1-h}^G(n-s, v) q_h (1-q_{1-h})^{s-1} + \sum_{w=1}^s Q_{1-h}^G(n-w, s) q_h (1-q_{1-h})^{w-1}$$

with initial conditions

$$Q_0^G(0, 0) = Q_1^G(0, 0) = \frac{1}{q_0 + q_1}.$$

Let us observe that distribution S_D can be assigned as follows

$$P(S_D \leq s) = P(S_A \leq s) + P(S_B \leq s) - P(S_G \leq s).$$

Theorem 4. Distribution of variable L specified on $M_{n,\theta}$ is given by recurrent formula

$$P(L = 1) = Q_0^L(n, 1) + Q_1^L(n, 1) \quad (12)$$

where for $h = 0, 1$

$$Q_h^L(n, 1) = Q_h^G(n-1, 1)(1 - q_{1-h}) + Q_{1-h}^L(n-1, s-1)q_h$$

with initial conditions

$$Q_0^L(0, 0) = Q_1^L(0, 0) = \frac{1}{q_0 + q_1}.$$

4. ESTIMATION OF THE POWER OF RUN TESTS

Taking into account recurrent formulas the power of test for $n = 1, 2, \dots, 1000$ and some pairs $\theta = (q_0, q_1) \in \Theta$ was assigned.

Because of better interpretation of parameters

$$p = p_A = \frac{q_1}{q_0 + q_1} \quad \text{and} \quad p = 1 - q_0 - q_1$$

adequately expressing Markov chain stationary probability and its autocorrelation coefficient, these pairs (q_0, q_1) , were chosen for which hypothesis revision $H_0: \rho = 0$ is as follows:

- if $S^A < s_\alpha^A - 1$, hypothesis H_0 is accepted,
- if $S^A \geq s_\alpha^A$, hypothesis H_0 is rejected (for the alternative $H_1': \rho > 0$),
- if $S^A = s_\alpha^A - 1$, hypothesis H_0 is accepted with probability r_α^A .

Randomized tests based on statistics S_B, S_D, S_G were analogously assigned.

The critical value for test based on run number is defined by the formula:

$$S_\alpha^L = \max\{a: F_L(s) \leq \alpha\},$$

where:

$$F_L(s) = P(S_L \leq s);$$

adequately randomized probability equals:

$$r_\alpha^L = \frac{F_L(s_\alpha)}{F_L(s_\alpha + 1) - F_L(s_\alpha)}.$$

For statistics S_L we took left-sided, and for other statistics right-sided critical regions.

Let us take for fixed n , p , significance level α and $\rho = 0$,

$$F_A(s) = P(S_A \leq s) \quad \text{for } s = 0, 1, \dots$$

The critical value of the test based on statistics S_A will be

$$s_\alpha^A = \min\{s : F_A(s) \geq 1 - \alpha\}.$$

This value corresponds to randomized probability

$$r_\alpha^A = \frac{F_A(s_\alpha) - (1 - \alpha)}{F_A(s_\alpha) - F_A(s_\alpha - 1)}.$$

For fixed n , p , and significance level α and simple alternative hypothesis of type $H_1: \rho = \rho_1$ given results allow (at least approximately) determine the power of independence tests (which verify hypothesis $H_0: \rho = 0$) based on statistics S_A , S_B , S_D , S_G , L , and, at the same time, to choose the strongest of them.

None of considered test is stronger than other tests, taking into account general alternative hypothesis $H_1: \rho > 0$. It is obvious that we must be careful with that generalizations based on numerical results but even on this stage we can formulate the following conclusions which are useful in choosing run test in applications (see Tab. 1, 2).

1. The test based on statistics S_A is stronger than the test based on statistics S_B for $\rho > 0.5$ except cases of strong asymmetry ($\rho > 0.6$) and very strong autocorrelation ($\rho > 0.7$).

2. The test S_G proved to be stronger than the tests S_A and S_B except cases of big asymmetry ($p > 0.6$) and very strong autocorrelation ($\rho > 0.7$). This test is also stronger than tests S_D , but only for p close to 0.5 and not very strong autocorrelation in relatively small sample size.

3. With the increase of p differences between the power of tests S_A and S_G and S_B and S_D decrease very fast (except cases of very small n ($n < 15$)). These differences are significant only in case of strong autocorrelation ($\rho > 0.5$).

The test S_L proved to be stronger than tests S_D and S_A in every considered case. It is also stronger than tests S_D and S_B except cases of strong asymmetry ($p > 0.7$) and not very strong autocorrelation ($\rho < 0.5$) for not very numerous samples ($n < 80$).

The power of the tests examined is the biggest for $p = 0.5$.

Table 1

Power of run tests for $\rho = 0.5, \rho = 0.3, 0.5, 0.7, 0.9, \alpha = 0.05$ (in %)o)

n	$\rho = 0.3$				$\rho = 0.5$				$\rho = 0.7$				$\rho = 0.9$			
	tests				tests				tests				tests			
	$s_A = s_B$	s_G	s_D	s_L	$s_A = s_B$	s_G	s_D	s_L	$s_A = s_B$	s_G	s_D	s_L	$s_A = s_B$	s_G	s_D	s_L
5	118	143	59	143	190	253	56	253	289	418	43	418	420	625	18	652
10	145	179	120	215	243	344	175	430	360	579	202	712	467	862	124	956
15	165	209	159	291	292	419	261	598	436	701	335	888	524	954	234	997
20	184	233	197	362	336	480	350	723	505	780	468	959	583	980	350	1 000
25	197	249	217	427	371	519	397	812	560	827	543	986	634	991	436	1 000
30	212	267	235	487	406	562	443	875	613	869	614	995	682	996	517	1 000
40	232	292	289	599	455	618	550	950	686	914	741	1 000	756	999	654	1 000
50	250	313	304	683	500	665	588	980	746	944	796	1 000	814	1 000	741	1 000
60	268	337	333	753	542	711	646	992	795	966	853	1 000	859	1 000	812	1 000
80	290	360	385	855	595	756	730	999	853	981	917	1 000	916	1 000	899	1 000
100	313	387	409	918	644	802	773	1000	898	991	942	1 000	950	1 000	944	1 000

Source: own calculations.

Table 2

Power of run tests for $p = 0.7$, $\rho = 0.6, 0.7, 0.8, 0.9$, $\alpha = 0.05$ (in %)o

n	$\rho = 0.6$					$\rho = 0.7$					$\rho = 0.8$					$\rho = 0.9$				
	S_A	S_B	S_G	S_D	S_L	S_A	S_B	S_G	S_D	S_L	S_A	S_B	S_G	S_D	S_L	S_A	S_B	S_G	S_D	S_L
5	231	270	307	45	307	143	241	175	51	175	95	194	105	67	105	67	134	70	81	70
10	351	356	507	203	669	353	338	406	211	502	213	296	221	237	221	98	193	99	176	99
15	414	436	603	343	870	386	415	443	366	688	373	359	379	331	471	143	246	143	241	143
20	482	507	679	458	947	432	482	497	448	825	390	423	397	411	550	207	299	208	298	208
25	535	574	732	553	978	485	545	551	523	920	416	483	422	478	728	302	349	302	348	302
30	584	615	775	606	992	532	602	594	590	961	447	537	453	535	767	391	395	391	395	439
40	659	692	835	718	999	605	678	663	685	990	517	631	522	631	897	404	479	404	479	495
50	716	756	873	785	1 000	664	734	717	737	998	576	689	581	689	960	428	552	428	552	657
60	763	804	903	826	1 000	711	782	759	784	1 000	623	733	626	733	982	458	612	458	612	707
80	828	861	938	894	1 000	781	856	820	856	1 000	696	806	696	806	997	526	682	526	682	861
100	875	984	961	929	1 000	830	906	861	906	1 000	751	861	753	861	1 000	588	741	588	741	920

Source: own calculations.

REFERENCES

- Bateman G. (1948), *On the Power Function of the Longest Run as a Test for Randomness in a Sequence of Alternatives*, "Biometrika", **35**, 97–112.
- Domański Cz. (1986), *Teoretyczne podstawy testów nieparametrycznych i ich zastosowanie w naukach ekonomiczno-społecznych*, Uniwersytet Łódzki, Łódź.
- Olmstead P. S. (1958), *Runs Determined in a Sample by an Arbitrary Cut*, "Bell System Technical Journal", **37**, 55–82.

Czesław Domański

OCENA MOCY TESTÓW LOSOWOŚCI OPARTYCH NA DŁUGOŚCI SERII

Testy oparte na długości serii stosowane są zarówno we wnioskowaniu statystycznym, jak również w statystycznej kontroli jakości.

W pracy przedstawiono moc trzech testów opartych na:

- maksymalnej długości serii z jednej strony mediany,
- mniejszej z maksymalnych długości serii powyżej i poniżej mediany,
- większej z maksymalnych długości serii powyżej i poniżej mediany.