

THE JOINT PROBABILITY GENERATING FUNCTION FOR RUN-LENGTHS IN REGENERATIVE BINARY MARKOV CHAINS, WITH APPLICATIONS¹

BY I. J. GOOD

Virginia Polytechnic Institute and State University

Gontcharov obtained the joint probability generating function for the numbers of runs of all lengths, both of successes and failures, in a Bernoulli sequence. This is here generalized to a class of regenerative binary Markov processes. For an allied class of Markov processes, the probability generating function is obtained for a "total score" defined in terms of runs of successes only, and asymptotic formulas are derived for the expectation and variance of the score.

1. Bernoulli and other runs. Gontcharov (1943, 1944) obtained numerous properties of the distribution of the numbers of runs of various lengths in a sequence of Bernoulli trials. All his results depended on his showing first that the joint or multivariate PGF (probability generating function) for runs of 1's and 0's, of various lengths, is

$$(1) \quad \frac{(1 + \sum p^r x_r)(1 + \sum q^s y_s)}{1 - \sum p^r x_r - \sum q^s y_s}$$

where the summations are for $r, s = 1, 2, \dots$; the probability of a "success" (a 1) on one trial is $p = 1 - q$, and the coefficient of $\prod x_r^{m_r} y_s^{n_s}$ is the joint probability of m_r runs of successes of length r , and n_s runs of failures of length s ($r, s = 1, 2, \dots$) when the Bernoulli sequence is known to be of length $\sum r(m_r + n_r)$. The constant term in (1) is 1 and corresponds to a sequence of length zero. In the terminology of Good (1961), (1) is a "universal" PGF in the sense that it covers sequences of all lengths simultaneously. My purpose is to generalize (1) to some Markov chains, "regenerative" in a sense to be defined, and to draw some deductions. Any application of the results for Bernoulli trials, such as to quality control, is also a potential application for the more general models of the present work.

We consider two classes of chains. The first class is regenerative at the end of *each run* of "successes" and of "failures"; the second class is regenerative after *each failure*. The asymptotic moments of an arbitrary "score", or linear function of the numbers of runs of each length, can be deduced from the universal PGF for both classes. But for the sake of simplicity this deduction is carried out explicitly only for the mean and variance of the score for the second class of

Received February 1972; revised January 1973.

Key words and phrases. Regenerative Markov chains, binary Markov chains, runs in Markov chains.

¹ This work was supported in part by H.E.W. grant No. 1 R01 GM18770-01.

chains, where here the score is defined in terms of the frequencies of runs of successes alone.

2. The PGF for a regenerative chain. We first generalize the joint PGF (1) to a regenerative binary Markov chain defined in the following manner. The conditional probability of a run of r 1's, following the end of a run of zeros, is κ_r , and the conditional probability of a run of s 0's, following the end of a run of 1's is λ_s ($r, s = 1, 2, 3, \dots$). This defines a process having two regenerative conditions, namely the right-hand ends of runs of 1's and 0's. (Perhaps the process should be called "biregenerative". Note that the process is *not necessarily an ordinary binary Markov chain of finite order.*) We shall consider a segment of such a process, of length L , and we shall refer to this segment as *the sequence*. We shall first find a formula for the probability that the sequence, assumed to be ergodic and in the stationary part of the process, will contain m_r runs of 1's of length r , and n_s runs of 0's of length s ($r, s = 1, 2, 3, \dots$). The formula is basically a multinomial expression, but with modifications related to the beginning and end of the sequence. If the sequence begins (or ends) with precisely a 1's or 0's we shall count this as a run of length a although its probability is not equal to κ_a or λ_a .

Let us begin at the beginning. The probability that the first element of the sequence, which we regard as a randomly selected trial, being a success is equal to $\bar{\kappa}/(\bar{\kappa} + \bar{\lambda})$, where $\bar{\kappa} = \sum r\kappa_r$ is the average length of a success run in an infinite chain, and similarly for $\bar{\lambda}$. Therefore the probability that the first element is *within* a success run of length r (in the *infinite* chain) is

$$\pi_r = \frac{r\kappa_r}{\sum r\kappa_r} \cdot \frac{\bar{\kappa}}{\bar{\kappa} + \bar{\lambda}} = \frac{r\kappa_r}{\bar{\kappa} + \bar{\lambda}}.$$

Therefore the probability that the sequence starts with exactly r consecutive 1's (which we are going to count as a run of length r) is

$$(2) \quad \sum_{\mu=r}^{\infty} \frac{\pi_{\mu}}{\mu} = \frac{K_r}{\bar{\kappa} + \bar{\lambda}},$$

and similarly for s 0's, it is

$$(3) \quad \frac{\Lambda_s}{\bar{\kappa} + \bar{\lambda}},$$

where

$$(4) \quad K_r = \sum_{\mu=r}^{\infty} \kappa_{\mu}, \quad \Lambda_s = \sum_{\nu=s}^{\infty} \lambda_{\nu}.$$

The end of the sequence is easier to deal with. If a run of 0's ends on the $(L - r)$ th element of the sequence, where $r < L$, then the probability that the sequence ends with a run of r 1's is

$$(5) \quad K_r,$$

with an obvious modification corresponding to the interchange of 0's and 1's.

By using (2) and (5) we see that the probability that the sequence begins with a run of a 1's, ends with b 1's, and contains m_r runs of 1's of length r , and n_s runs of 0's of length s ($r, s = 1, 2, \dots; \sum r(m_r + n_r) = L$), is

$$(6) \quad \frac{1}{\bar{\kappa} + \bar{\lambda}} \mathcal{C}(\mathbf{x}^m \mathbf{y}^n) (\sum \kappa_r x_r)^{M-2} K_a x_a K_b x_b (\sum \lambda_r y_r)^{M-1}$$

where $a \geq 1, b \geq 1, a + b < L, \mathbf{x}^m = x_1^{m_1} x_2^{m_2} \dots, \mathbf{y}^n = y_1^{n_1} y_2^{n_2} \dots, \mathcal{C}(\mathbf{x}^m \mathbf{y}^n)$ means "the coefficient of $\mathbf{x}^m \mathbf{y}^n$ in", $M = \sum m_r$, and the summations run from $r = 1$ to ∞ . If a and b are unspecified, the probability is therefore

$$(7) \quad \frac{1}{\bar{\kappa} + \bar{\lambda}} \mathcal{C}(\mathbf{x}^m \mathbf{y}^n) (\sum \kappa_r x_r)^{M-2} (\sum K_a x_a)^2 (\sum \lambda_r y_r)^{M-1}$$

where we can allow the middle summations to run from $a = 1$ to ∞ if $M > 1$. We shall deduce that if the initial and final elements of the sequence are unspecified the PGF is

$$(8) \quad 1 + \frac{1}{\bar{\kappa} + \bar{\lambda}} \left\{ K^* + \Lambda^* + \frac{K^2 Y + \Lambda^2 X + 2K\Lambda}{1 - XY} \right\}$$

where

$$\begin{aligned} X &= \sum \kappa_r x_r, & Y &= \sum \lambda_r y_r \\ K &= \sum K_r x_r, & \Lambda &= \sum \Lambda_r y_r \\ K^* &= \sum (K_r + K_{r+1} + \dots) x_r, & \Lambda^* &= \sum (\Lambda_r + \Lambda_{r+1} + \dots) y_r. \end{aligned}$$

In (8), the first term 1 corresponds to $L = 1$, the terms K^* and Λ^* within the braces correspond to the cases $m_L = 1$ and $n_L = 1$, and the other terms are obtained by summing expressions like (7) from $M = 2$ to ∞ . The expression for K^* is easily proved to be appropriate by writing it in the form

$$\sum (\kappa_r + 2\kappa_{r+1} + 3\kappa_{r+2} + \dots) x_r.$$

Gontcharov's PGF (1) can be deduced from (8) by writing $\kappa_r = p^{r-1}q, \lambda_r = q^{r-1}p, K_r = p^{r-1}, \Lambda_r = q^{r-1}$, and $\bar{\kappa} + \bar{\lambda} = 1/(pq)$.

3. Scoring. If scores of s_r and s_r' are associated with runs of length r of 1's and 0's, and the total score is $S = \sum (s_r m_r + s_r' n_r)$, then the PGF of S is obtainable from (8) by replacing x_r by $x^{s_r} z$, y_r by $x^{s_r'} z$ and then extracting the coefficient of z^l to allow for the constraint $\sum r(m_r + n_r) = L$.

It would be possible to obtain the moments of the total score S from the generating function, but the calculations would be very heavy in the general case. Instead I shall exemplify the methods by means of a slightly different problem which is almost a special case of the above one.

4. Runs of successes. In the present section we shall be concerned with runs of ones only, but we shall include runs of zero length. We shall abbreviate the expression "run of success" simply to "run". Our model will be of the "second class" (see Section 1).

Let ρ_n and ρ'_n be the probabilities that respectively *true* and *apparent* runs of length n occur at a specified place in the (infinitely long) chain, where a true run is a run of *precisely* n successes (preceded and followed by a failure) and an apparent run is one that may or may not be preceded and followed by failures. A true run is regarded as “occurring” at its first success, and a true run of length n contains $n - m + 1$ apparent runs of length $m (m \leq n)$. The distinction between true and apparent runs (and half-true ones: see Section 5) and formulas (10) to (13) were drawn to my attention in 1940 by A. M. Turing, who invented the regenerative model of the second class for the analysis of certain binary processes. We adopt the natural convention that

$$(9) \quad \rho'_0 = 1.$$

The following identities are easy to prove once they are pointed out:

$$(10) \quad \rho_n = \Delta^2 \rho'_n = \rho'_n - 2\rho'_{n+1} + \rho'_{n+2} = \nabla^2 \rho'_{n+2}$$

$$(11) \quad \rho_n + \rho_{n+1} + \rho_{n+2} + \dots = \rho'_n - \rho'_{n+1}$$

$$(12) \quad \rho_0 + \rho_1 + \rho_2 + \dots = 1 - \rho'_1$$

$$(13) \quad \rho_n + 2\rho_{n+1} + 3\rho_{n+2} + \dots = \rho'_n$$

$$(14) \quad \rho_n + 3\rho_{n+1} + 6\rho_{n+2} + 10\rho_{n+3} + \dots = \rho'_n + \rho'_{n+1} + \dots$$

$$(15) \quad \sum_{n=0}^{\infty} (n+1)^2 \rho_n = 1 + 2(\rho'_1 + \rho'_2 + \rho'_3 + \dots).$$

These formulas depend on the assumption of stationarity. They reduce to simple algebraic identities for the case of a random sequence (Bernoulli trials), for which $\rho'_n = p^n$, $\rho_n = (1-p)^2 p^n$, p being the probability of a success.

We shall assume further that the Markov chain is *regenerative after each failure*, and we shall think of the whole chain as starting with a failure. This model is almost a special case of the model of the first class considered in Section 2, and may be regarded as a compromise between that model and a Bernoulli process. But we are now further assuming that the sequence starts with a failure so as to avoid “initial end effects”. This will not affect the asymptotic results. The probability of any specified sequence, of length N , can be readily expressed in terms of the ρ_n 's.

We now define the *total score* as

$$(16) \quad S = s_0 m_0 + s_1 m_1 + s_2 m_2 + \dots,$$

where s_n is a score arising from each true run of length n . The advantage of considering total scores with the generality of S is that it includes various special cases of some interest. For example, the following formulas can be readily obtained from the subsequent formulas (20) and (21).

Let r_0 and r_1 be the numbers of failures and successes in the segment. Then, for large N ,

$$(17) \quad E(r_1) = NA + (A + B)(1 - \frac{1}{2}A) + o(1)$$

$$(18) \quad \text{Var } r_0 = \text{Var } r_1 \sim N(B - A^2 - 2AB + A^3 + A^2B),$$

where

$$A = \sum n\rho_n, \quad B = \sum n^2\rho_n.$$

The variance of the number of runs of length ν is asymptotically

$$(19) \quad N\{\rho'_\nu + 2(\rho'_{\nu+1} + \rho'_{\nu+2} + \dots) - (2\nu + 1)\rho'_\nu{}^2 - 4\rho'_\nu(\rho'_{\nu+1} + \rho'_{\nu+2} + \dots) + 2(\rho'_1 + \rho'_2 + \dots)\rho'_\nu{}^2\}.$$

More generally, under assumptions (38),

$$(20) \quad E(S) = N \sum_{n=0}^{\infty} t_n + \sum (n + 1)t_n - \sum t_n \sum_{n=1}^{\infty} \rho'_n + o(1)$$

$$(21) \quad \text{Var}(S) = N\{\sum s_n t_n - \sum (2n + 1)t_n \sum t_n + 2(\rho'_1 + \rho'_2 + \dots)(\sum t_n)^2\} + O(1)$$

where

$$(22) \quad t_n = \rho_n s_n.$$

5. A generating function for the score. We can regard “the sequence” as composed entirely of subsequences (“half-true runs”) all of the form $0\ 1\ 1\ 1\ \dots\ 1$, juxtaposed, where each subsequence ends with a run of 1’s but this run can be of zero length. By virtue of the regenerative assumption, this way of regarding the sequence reduces it to a random sequence, not of course to be confused with the Bernoulli sequence that we mentioned before as a special case. (Our “alphabet” is countably infinite.)

Let κ_n be the probability of the selection of a subsequence with n 1’s on it. Then

$$(23) \quad \kappa_n = \rho_n / (1 - \rho_1')$$

$$(24) \quad \kappa_0 + \kappa_1 + \kappa_2 + \dots = 1$$

$$(25) \quad \kappa_1 + 2\kappa_2 + 3\kappa_3 + \dots = \rho_1' / (1 - \rho_1')$$

$$(26) \quad m_0 + 2m_1 + 3m_2 + \dots = N,$$

provided we count the right-most subsequence as a “true” run (but we shall allow for its correct probability). Let

$$(27) \quad m_0 + m_1 + m_2 + \dots = M,$$

but we regard N as given, rather than M .

The probability of our sequence, if it ends with just c 1’s, is

$$(28) \quad \kappa_0{}^{m_0} \kappa_1{}^{m_1} \dots \kappa_{c-1}{}^{m_{c-1}} \kappa_c{}^{m_c-1} \kappa_{c+1}{}^{m_{c+1}} \dots (\kappa_c + \kappa_{c+1} + \dots).$$

The PGF (probability generating function) of (m_0, m_1, m_2, \dots) is therefore

$$(29) \quad \sum_{\substack{m_0+2m_1+\dots=N \\ m_0, m_1, \dots}} \frac{(M-1)!}{m_0! m_1! \dots m_{c-1}! (m_c-1)! m_{c+1}! \dots} (\kappa_0 x_0)^{m_0} \dots (\kappa_{c-1} x_{c-1})^{m_{c-1}} \{(\kappa_c x_c)^{m_c-1} (\kappa_{c+1} x_{c+1})^{m_{c+1}} \dots (\kappa_c + \kappa_{c+1} + \dots) x_c\}.$$

If $N \geq 1$, the PGF is the coefficient of y^N in the expression obtained from (29) by replacing x_n by $x_n y^{n+1}$ ($n = 0, 1, 2, \dots$), and where the condition

$$m_0 + 2m_1 + \dots = N$$

is no longer required. Hence the PGF equals

$$\mathcal{E}(y^N)G(x_0, x_1, x_2, \dots; y),$$

where $\mathcal{E}(y^N)$ means “the coefficient of y^N in” and

$$(30) \quad G(x_0, x_1, \dots; y) = \frac{\sum_{n=0}^{\infty} K_n x_n y^{n+1}}{1 - \sum_{n=0}^{\infty} \kappa_n x_n y^{n+1}}$$

where

$$(31) \quad K_n = \kappa_n + \kappa_{n+1} + \dots = \frac{\rho_n' - \rho_{n+1}'}{1 - \rho_1'}.$$

From (24) we have

$$(32) \quad (y + y^2 + y^3 + \dots)(1 - \sum \kappa_n y^{n+1}) = \sum K_n y^{n+1},$$

so that

$$(33) \quad G(1, 1, 1, \dots; y) = y + y^2 + y^3 + \dots,$$

which is a check of (30).

If we replace x_n by x^{s_n} we deduce that the PGF of S is $\mathcal{E}(y^N)F(x, y)$ where

$$(34) \quad F(x, y) = \frac{\sum_{n=0}^{\infty} K_n x^{s_n} y^{n+1}}{1 - \sum_{n=0}^{\infty} \kappa_n x^{s_n} y^{n+1}}.$$

6. The moments of S : outline of derivations. The moment generating function of S is $\mathcal{E}(y^N)F(e^u, y)$ where u is the dummy variable. The derivation of the moments involves some points of mathematical rigor, but we give only an outline because the kinds of mathematics required are already in the literature (see, for example, Smith (1957)) and because the full details are decidedly intricate. We first mention the easily proved lemma that if

$$f(y) = \sum_{n=0}^{\infty} (\alpha n + \beta + \varepsilon_n)y^n,$$

where $\varepsilon_n \rightarrow 0$ as $n \rightarrow \infty$, then $\alpha = \lim_{y \rightarrow 1} (1 - y)^2 f(y)$. Now we can show, by means of (32), that

$$(35) \quad \left. \frac{\partial}{\partial u} F(e^u, y) \right|_{u=0} = \frac{\sum K_n s_n y^{n+1} + (y + y^2 + \dots) \sum \kappa_n s_n y^{n+1}}{1 - \sum \kappa_n y^{n+1}}.$$

Hence, assuming $E(S)$ is of the form

$$(36) \quad E(S) = \alpha N + \beta + o(1),$$

we can deduce $\alpha = \sum t_n$ from the lemma. By a more lengthy but similar argument we can show that

$$(37) \quad \beta = \sum (n + 1)t_n - \sum t_n \sum_{n=1}^{\infty} \rho_n'.$$

By means of an even more lengthy argument (which makes use of Abel's theorem for power series), we can derive formula (21) for $\text{Var}(S)$ under the assumptions

$$(38) \quad \begin{array}{ll} \kappa_n = O(\lambda^n) & \text{for some } \lambda < 1 \\ |s_n| < n^c & \text{for some } c. \end{array}$$

REFERENCES

- [1] GONTCHAROV, W. (1943). On runs of events in a series of independent trials of the Bernoullian type. *C. R. (Doklady) Acad. Sci. URSS* **38** 283-285. *MR* **5** 124.
- [2] GONTCHAROV, W. (1944). On the field of combinatory analysis. *Bull. Acad. Sci. URSS Ser. Math.* **8**, no. 4 1-48. *MR* **6** 88.
- [3] GOOD, I. J. (1961). The frequency count of Markov chain and the transition to continuous time. *Ann. Math. Statist.* **32** 41-48. *MR* 23A no. 4241.
- [4] SMITH, W. L. (1957). Contribution to the discussion of a paper by J. G. Skellam and L. R. Shenton. *J. Roy. Statist. Soc. Ser. B* **19** 114-115. *MR* 19 990.

DEPARTMENT OF STATISTICS
VIRGINIA POLYTECHNIC INSTITUTE
& STATE UNIVERSITY
BLACKSBURG, VIRGINIA 24061