

COUNTS OF FAILURE STRINGS IN CERTAIN BERNOULLI SEQUENCES

LARS HOLST,* *Royal Institute of Technology*

Abstract

In a sequence of independent Bernoulli trials the probability for success in the k th trial is p_k , $k = 1, 2, \dots$. The number of strings with a given number of failures between two subsequent successes is studied. Explicit expressions for distributions and moments are obtained for the case in which $p_k = a/(a + b + k - 1)$, $a > 0$, $b \geq 0$. Also, the limit behaviour of the longest failure string in the first n trials is considered. For $b = 0$, the strings correspond to cycles in random permutations.

Keywords: Binomial moment; Ewens sampling formula; Hoppe's urn; Poisson distribution; Poisson–Dirichlet distribution; Pólya's urn; random permutation; record; spacing; sums of indicators

2000 Mathematics Subject Classification: Primary 60C05

Secondary 60K99

1. Introduction

In an infinite sequence of independent Bernoulli trials the probability for success in the k th trial is p_k for $k = 1, 2, \dots$. A d -string is a string $SF \cdots FS$ of $d - 1$ failures between two subsequent successes. In this paper we study the number of such strings. Explicit results are obtained for the case in which $p_k = a/(a + b + k - 1)$, $a > 0$, $b \geq 0$. To our knowledge only special cases have been studied previously.

For $a = 1$ and $b = 0$, that is $p_k = 1/k$, 1-strings correspond to double records in a record sequence. Hahlin (1995) proved that the total number of such records is $\text{Po}(1)$ (Poisson distributed with mean 1). After that, an unpublished proof by Diaconis inspired a number of studies on 1-strings; see Chern *et al.* (2000), Mori (2001), Joffe *et al.* (2004) and the references therein. For the case in which $p_k = a/(a + b + k - 1)$, explicit expressions for the binomial moments of the number of 1-strings in the first n trials have been derived in Holst (2006).

Sethuraman and Sethuraman (2004) studied d -strings for $a = 1$ and $b > 0$, and obtained the joint distribution of the number of d -strings for $d = 1, 2, \dots$. For $a > 0$ and $b = 0$, d -strings are closely connected with cycle lengths in a -biased random permutations; see, e.g. Arratia *et al.* (2003, p. 95). In Gnedin (2007) coherent sequences of random permutations are studied, where the probability of a 'lower record' in the k th permutation is $p_k = a/(a + b + k - 1)$.

In Section 2 we introduce notation and derive recursions for the binomial moments of the number of d -strings in a finite sequence for general p_k s. The special case in which $p_k = a/(a + k - 1)$, connected with random permutations, is studied in Section 3. In Section 4 we derive our main result: the joint distribution of the total number of d -strings, $d = 1, 2, \dots$, and study the limit behaviour of the longest failure string in the first n trials in an infinite Bernoulli sequence with $p_k = a/(a + b + k - 1)$.

Received 9 October 2006; revision received 11 May 2007.

* Postal address: Department of Mathematics, Royal Institute of Technology, SE-10044 Stockholm, Sweden.
Email address: lholst@math.kth.se

2. The general case: notation and moments

In the following, I_1, I_2, \dots is a sequence of independent Bernoulli random variables, i.e. I_k is $\text{Be}(p_k)$, with

$$P(I_k = 1) = 1 - P(I_k = 0) = p_k > 0.$$

The number of d -strings in the first n trials is

$$M_{dn} = \sum_{k=1}^{n-d} I_k(1 - I_{k+1}) \cdots (1 - I_{k+d-1})I_{k+d}.$$

Note that $M_{dn} = 0$ for $d \geq n$ and $\sum_{j=1}^{n-1} jM_{jn} \leq n - 1$. Implicitly, the following result gives the distribution of (M_{1n}, \dots, M_{dn}) .

Proposition 1. *The binomial moments*

$$f_n(r_1, \dots, r_d) = E\left(\binom{M_{1n}}{r_1} \cdots \binom{M_{dn}}{r_d}\right)$$

disappear for $\sum_{j=1}^d jr_j \geq n$ and fulfill the following recursion:

$$\begin{aligned} f_{n+1}(r_1, \dots, r_d) &= f_n(r_1, \dots, r_d) \\ &+ p_{n+1}[f_n(r_1 - 1, r_2, \dots, r_d) - (1 - p_n)f_{n-1}(r_1 - 1, r_2, \dots, r_d)] \\ &+ p_{n+1}(1 - p_n)[f_{n-1}(r_1, r_2 - 1, r_3, \dots) - (1 - p_{n-1})f_{n-2}(r_1, r_2 - 1, r_3, \dots)] \\ &+ \cdots + p_{n+1}(1 - p_n) \cdots (1 - p_{n-d+2}) \\ &\times [f_{n-d+1}(r_1, \dots, r_{d-1}, r_d - 1) - (1 - p_{n-d+1})f_{n-d}(r_1, \dots, r_{d-1}, r_d - 1)]. \end{aligned}$$

Proof. As $\sum_{j=1}^d jM_{jn} \leq n - 1$ the first assertion follows. To obtain the recursion we use generating functions and the independence between the I_k s, i.e.

$$\begin{aligned} E[t_1^{M_{1,n+1}} \cdots t_d^{M_{d,n+1}}] &= E[t_1^{M_{1n}} \cdots t_d^{M_{dn}} (1 + (t_1 - 1)I_n I_{n+1})(1 + (t_2 - 1)I_{n-1}(1 - I_n)I_{n+1}) \cdots] \\ &= E[t_1^{M_{1n}} \cdots t_d^{M_{dn}} (1 + (t_1 - 1)I_n I_{n+1} + (t_2 - 1)I_{n-1}(1 - I_n)I_{n+1} + \cdots)] \\ &= E[t_1^{M_{1n}} \cdots t_d^{M_{dn}}] + (t_1 - 1)p_{n+1} E[t_1^{M_{1n}} \cdots t_d^{M_{dn}} (1 - (1 - I_n))] \\ &\quad + (t_2 - 1)p_{n+1} E[t_1^{M_{1n}} \cdots t_d^{M_{dn}} (1 - (1 - I_{n-1}))(1 - I_n)] + \cdots \\ &= E[t_1^{M_{1n}} \cdots t_d^{M_{dn}}] + (t_1 - 1)p_{n+1}[E(t_1^{M_{1n}} \cdots t_d^{M_{dn}}) - (1 - p_n) E(t_1^{M_{1,n-1}} \cdots t_d^{M_{d,n-1}})] \\ &\quad + (t_2 - 1)p_{n+1}(1 - p_n)[E(t_1^{M_{1,n-1}} \cdots t_d^{M_{d,n-1}}) - (1 - p_{n-1}) E(t_1^{M_{1,n-2}} \cdots t_d^{M_{d,n-2}})] \\ &\quad + \cdots \end{aligned}$$

Expansion in series around $t_1 = 1, \dots, t_d = 1$ proves the second assertion.

Including the string $SF \cdots F$ with $d - 1$ failures after the last success in the count, we obtain the random variable

$$N_{dn} = M_{dn} + I_{n-d+1}(1 - I_{n-d+2}) \cdots (1 - I_n).$$

Proposition 2. *For the binomial moments*

$$\begin{aligned} & \mathbb{E}\left(\binom{N_{1n}}{r_1} \cdots \binom{N_{dn}}{r_d}\right) \\ &= \mathbb{E}\left(\binom{M_{1n}}{r_1} \cdots \binom{M_{dn}}{r_d}\right) + \frac{1}{p_{n+1}} \left(\mathbb{E}\left(\binom{M_{1,n+1}}{r_1} \cdots \binom{M_{d,n+1}}{r_d}\right) \right. \\ & \quad \left. - \mathbb{E}\left(\binom{M_{1n}}{r_1} \cdots \binom{M_{dn}}{r_d}\right) \right). \end{aligned}$$

Proof. By the law of total probability we have

$$\mathbb{E}\left(t_1^{M_{1,n+1}} \cdots t_d^{M_{d,n+1}}\right) = p_{n+1} \mathbb{E}\left(t_1^{N_{1n}} \cdots t_d^{N_{dn}}\right) + (1 - p_{n+1}) \mathbb{E}\left(t_1^{M_{1n}} \cdots t_d^{M_{dn}}\right),$$

from which the assertion follows by an expansion in series.

In an infinite sequence the total number of d -strings

$$M_{d\infty} = \sum_{k=1}^{\infty} I_k(1 - I_{k+1}) \cdots (1 - I_{k+d-1}) I_{k+d} < +\infty$$

with probability 1 if and only if

$$\mathbb{E}(M_{d\infty}) = \sum_{k=1}^{\infty} p_k(1 - p_{k+1}) \cdots (1 - p_{k+d-1}) p_{k+d} < +\infty.$$

Indeed, by splitting the series for $M_{d\infty}$ into $d + 1$ (independent) series this follows from the Borel–Cantelli lemmas; see Mori (2001, p. 834).

3. The case in which $p_k = a/(a + k - 1)$

Following Knuth (1992) we denote descending and ascending factorials by

$$x^n = x(x - 1) \cdots (x - n + 1), \quad x^{\bar{n}} = x(x + 1) \cdots (x + n - 1) = \sum_{k=1}^n \left[\begin{matrix} n \\ k \end{matrix} \right] x^k,$$

where $\left[\begin{matrix} n \\ k \end{matrix} \right]$ is a cycle number or signless Stirling number of the first kind. In the rest of this section we assume that $p_k = a/(a + k - 1)$ with $a > 0$. Closed simple formulae can be obtained for the binomial moments. Note that $\sum_{j=1}^{n-1} j M_{jn} \leq n - 1$ and $\sum_{j=1}^n j N_{jn} = n$.

Proposition 3. *With $m = \sum_{j=1}^d j r_j$,*

$$\mathbb{E}\left(\binom{M_{1n}}{r_1} \cdots \binom{M_{dn}}{r_d}\right) = f_n(r_1, \dots, r_d) = I(m \leq n - 1) \frac{(n - 1)^m}{(a + n - 1)^m} \prod_{j=1}^d \frac{(a/j)^{r_j}}{r_j!}.$$

Proof. By telescoping sums we obtain, for $d \leq n - 1$,

$$\begin{aligned} E(M_{dn}) &= \sum_{k=1}^{n-d} \frac{a}{a+k-1} \left(1 - \frac{a}{a+k}\right) \cdots \left(1 - \frac{a}{a+k+d-2}\right) \frac{a}{a+k+d-1} \\ &= \frac{a}{d} \sum_{k=1}^{n-d} \left[\left(1 - \frac{a}{a+k+d-1}\right) - \left(1 - \frac{a}{a+k-1}\right) \right] \\ &\quad \times \left(1 - \frac{a}{a+k}\right) \cdots \left(1 - \frac{a}{a+k+d-2}\right) \\ &= \frac{a}{d} \left(1 - \frac{a}{a+n-d}\right) \cdots \left(1 - \frac{a}{a+n-1}\right) \\ &= \frac{a}{d} \frac{(n-1)^d}{(a+n-1)^d}. \end{aligned}$$

Hence, the proposition holds for $E(M_{dn})$.

With $p_k = a/(a+k-1)$ and f_n as in the assertion, we can verify that the right-hand side of the recursion in Proposition 1 can be written

$$\begin{aligned} f_n(r_1, \dots, r_d) + \prod_{j=1}^d \frac{(a/j)^{r_j}}{r_j!} \left(\frac{n}{a+n} - \frac{n-m}{a+n-m} \right) \frac{(n-1)^{m-1}}{(a+n-1)^{m-1}} \\ = f_n(r_1, \dots, r_d) + (f_{n+1}(r_1, \dots, r_d) - f_n(r_1, \dots, r_d)) \\ = f_{n+1}(r_1, \dots, r_d). \end{aligned}$$

Hence, the recursion is satisfied by f_n with $E(M_{dn})$ as starting values. The assertion follows from this.

Proposition 4. With $m = \sum_{j=1}^d jr_j$,

$$E\left(\binom{N_{1n}}{r_1} \cdots \binom{N_{dn}}{r_d}\right) = I(m \leq n) \frac{n^m}{(a+n-1)^m} \prod_{j=1}^d \frac{(a/j)^{r_j}}{r_j!},$$

and, for $\sum_{j=1}^n jx_j = n$,

$$P(N_{1n} = x_1, \dots, N_{nn} = x_n) = \frac{n!}{a^n} \prod_{j=1}^n \frac{(a/j)^{x_j}}{x_j!}.$$

Proof. Using Propositions 2 and 3, the first assertion follows from an elementary calculation. Using generating functions, we have

$$\begin{aligned} E(t_1^{N_{1n}} \cdots t_n^{N_{nn}}) &= E((1 + (t_1 - 1))^{N_{1n}} \cdots (1 + (t_n - 1))^{N_{nn}}) \\ &= \sum E\left(\binom{N_{1n}}{r_1} \cdots \binom{N_{nn}}{r_n}\right) (t_1 - 1)^{r_1} \cdots (t_n - 1)^{r_n} \\ &= \sum \sum E\left(\binom{N_{1n}}{r_1} \cdots \binom{N_{nn}}{r_n}\right) (-1)^{r_1-x_1} \cdots (-1)^{r_n-x_n} \\ &\quad \times \binom{r_1}{x_1} \cdots \binom{r_n}{x_n} t_1^{x_1} \cdots t_n^{x_n}. \end{aligned}$$

As $\sum_1^n jN_{jn} = n$, the binomial moments disappear for $\sum_1^n jr_j > n$. Therefore, for $\sum_1^n jx_j = n$, we have $r_j = x_j$ in the summation. Thus,

$$P(N_{1n} = x_1, \dots, N_{nn} = x_n) = E\left(\binom{N_{1n}}{x_1} \cdots \binom{N_{nn}}{x_n}\right),$$

proving the second assertion.

The distribution of (N_{1n}, \dots, N_{nn}) is the famous *Ewens Sampling Formula*. Furthermore, N_{dn} is the number of d -strings in $1I_2I_3 \cdots I_n1$. Using this, Propositions 3 and 4 can be derived by combinatorial arguments; see Arratia *et al.* (2003, p. 95). In this context, N_{dn} is interpreted as the number of cycles of length d in a random permutation of $1, 2, \dots, n$ biased by a^{K_n} , where $K_n = \sum_{k=1}^n I_k$ is the number of cycles with

$$P(K_n = j) = \binom{n}{j} \frac{a^j}{a^n}, \quad j = 1, 2, \dots, n.$$

The moment convergence

$$E\left(\binom{M_{1n}}{r_1} \cdots \binom{M_{dn}}{r_d}\right) \rightarrow \prod_{j=1}^d \frac{(a/j)^{r_j}}{r_j!}, \quad n \rightarrow \infty,$$

implies the following result, which is well known for a -biased random permutations; see Arratia *et al.* (2003, p. 96).

Proposition 5. *The number of strings $M_{1\infty}, M_{2\infty}, \dots$ are independent Poisson random variables with $E(M_{d\infty}) = a/d$.*

4. The case in which $p_k = a/(a + b + k - 1)$

In this section we assume that $p_k = a/(a + b + k - 1)$ with $a > 0$ and $b > 0$. Clearly

$$M_{d\infty} = \sum_{k=1}^{\infty} I_k(1 - I_{k+1}) \cdots (1 - I_{k+d-1})I_{k+d} < +\infty$$

with probability 1. Mori (2001) derived the distribution of $M_{1\infty}$. For the special case in which $a = 1$, Sethuraman and Sethuraman (2004) obtained the joint distribution of $M_{1\infty}, M_{2\infty}, \dots$. Using different methods, we generalise their result to any $a > 0$. Let U be Beta(a, b), that is a random variable with density

$$f_U(u) = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} u^{a-1}(1 - u)^{b-1}, \quad 0 < u < 1.$$

Theorem 1. *Conditional on a Beta(a, b) random variable U , the number of strings $M_{1\infty}, M_{2\infty}, \dots$ are independent Poisson random variables with*

$$E(M_{d\infty} | U) = \frac{a}{d}(1 - (1 - U)^d), \quad d = 1, 2, \dots$$

Proof. We introduce the following mixture of Pólya’s and Hoppe’s urn models. An urn contains initially one white ball and one black ball of weights a and b , respectively. Balls are

drawn at random proportional to weights. The white ball and the black ball are replaced together with a new ball of a colour not present in the urn, other balls are replaced together with one new ball of the same colour. All new balls have weight 1. Let $I_k = 1$ if we obtain the white ball at drawing k , else $I_k = 0$. Obviously the I_k s are independent and $P(I_k = 1) = a/(a + b + k - 1)$.

Next, generate a sequence of W s and B s. We obtain a W if drawing the white ball or a ball of a colour emanating from a draw of the white, else we get a B . This sequence is as drawing from an ordinary Pólya urn. Note that the sequence is exchangeable. Therefore, by de Finetti's theorem the sequence can be thought of as having been generated by first observing a Beta(a, b) random variable U and then, conditional on the outcome $U = u$, generating a sequence of independent Be(u) random variables, with W corresponding to 1 and B to 0.

In the *subsequence* of W s in the original sequence I_1, I_2, \dots , the probability of getting the white ball at the j th trial is $p_j^* = a/(a + j - 1)$. According to Proposition 5 the number of d -strings in the subsequence, $M_{d\infty}^*$, is Po(a/d) and $M_{1\infty}^*, M_{2\infty}^*, \dots$ are independent.

Now recall the following well-known fact about disintegration of a Poisson process. If the random variable ξ is Po(μ) and independent of the independent Be(p) random variables $\varepsilon_1, \varepsilon_2, \dots$, then $\sum_{j=1}^{\xi} \varepsilon_j$ and $\sum_{j=1}^{\xi} (1 - \varepsilon_j)$ are independent Po(μp) and Po($\mu(1 - p)$), respectively.

Consider the $M_{1\infty}^*$ 1-strings in the *subsequence* of W s. Each such 1-string is also a 1-string in the *original sequence* I_1, I_2, \dots , provided it is *not interrupted* by a B . Conditional on $U = u$ the sequence of W s and B s can be considered as independent Be(u)-trials. Hence, the probability for interruption is $1 - u$. As $M_{1\infty}^*$ is Po(a), it follows from the above disintegration with $\xi = M_{1\infty}^*$ and the ε s independent Be(u) random variables, that the number of 1-strings in the *original sequence*, $M_{1\infty}$, is Po(au) and independent of the number of interrupted 1-strings, $M_{1\infty}^* - M_{1\infty}$, which is Po($a(1 - u)$).

For the number of 2-strings, $M_{2\infty}$, we can argue in a similar way. First we have the $M_{1\infty}^* - M_{1\infty}$ interrupted 1-strings in the *subsequence* of W s. Such a 1-string is a 2-string in the *original sequence* provided it is *only* interrupted by *one* B . The probability for only one interruption is $(1 - u)u$, so by the disintegration the number of such strings is Po($au(1 - u)$). The number of 2-strings in the *subsequence* of W s, $M_{2\infty}^*$, is Po($a/2$). Such a 2-string is also a 2-string in the *original sequence* if it is *not interrupted*. The probability for this is u^2 . Hence, by disintegration the number of such 2-strings is Po($au^2/2$). As $M_{1\infty}^*$ and $M_{2\infty}^*$ are independent, it follows that the random variable $M_{2\infty}$ is Poisson with mean

$$a(1 - u)u + \frac{a}{2}u^2 = \frac{a}{2}(1 - (1 - u)^2),$$

and independent of $M_{1\infty}$.

The arguments extend to show that the random variable $M_{d\infty}$ conditional on $U = u$ is Poisson with mean

$$\begin{aligned} \frac{a}{d}u^d + \frac{a}{d-1} \binom{d-1}{1} u^{d-1}(1-u) + \frac{a}{d-2} \binom{d-1}{2} u^{d-2}(1-u)^2 + \dots + au(1-u)^{d-1} \\ = \frac{a}{d}(1 - (1-u)^d), \end{aligned}$$

and independent of $M_{1\infty}, M_{2\infty}, \dots, M_{d-1,\infty}$.

Finally, consider long strings of failures. Let the last success in the first n trials occur at trial $n + 1 - A_{1n}$; if there is no success let $A_{1n} = 0$. We have, for $j = 1, 2, \dots, n$,

$$\begin{aligned} P(A_{1n} > j) &= \left(1 - \frac{a}{a+b+n-j}\right) \cdots \left(1 - \frac{a}{a+b+n-1}\right) \\ &= \frac{(b+n-j)^{\bar{j}}}{(a+b+n-j)^{\bar{j}}} = \frac{\Gamma(b+n)}{\Gamma(b+n-j)} \frac{\Gamma(a+b+n-j)}{\Gamma(a+b+n)}. \end{aligned}$$

For $j, n \rightarrow \infty$ such that $j/n \rightarrow x$, $0 < x < 1$, Stirling's formula gives

$$P\left(\frac{A_{1n}}{n} > \frac{j}{n}\right) \rightarrow (1-x)^a, \quad n \rightarrow \infty,$$

that is A_{1n}/n converges in distribution to $\text{Beta}(1, a)$.

In a similar way, we find, for the number of trials between the last and the second to last success, A_{2n} , that

$$\frac{(A_{1n}, A_{2n})}{n} \rightarrow (U_1, (1-U_1)U_2), \quad n \rightarrow \infty,$$

in distribution, where U_1, U_2 are independent $\text{Beta}(1, a)$ random variables. The procedure can be repeated in a like manner for the second to last and the third to last success, etc.

The limit behaviour of the long strings is as if A_{1n}, A_{2n}, \dots had been cycle lengths in an a -biased random permutation; see Arratia *et al.* (2003, Section 5.4). The limit distribution of the normalized size ordered A s is the Poisson–Dirichlet distribution with parameter a . In particular, we have the following result.

Theorem 2. *For the longest string of failures in the first n trials, we find that*

$$\frac{\max(A_{1n}, A_{2n}, \dots)}{n} \rightarrow L_1 = \max(U_1, (1-U_1)U_2, (1-U_1)(1-U_2)U_3, \dots), \quad n \rightarrow \infty,$$

in distribution, where U_1, U_2, \dots are independent $\text{Beta}(1, a)$ random variables.

Various formulae connected with the random variable L_1 can be found in Arratia *et al.* (2003, Section 5.5).

References

- ARRATIA, R., BARBOUR, A. D. AND TAVARÉ, S. (2003). *Logarithmic Combinatorial Structures: a Probabilistic Approach*. European Mathematical Society Publishing House, ETH-Zentrum, Zürich.
- CHERN, H.-H., HWANG, H.-K. AND YEH, Y.-N. (2000). Distribution of the number of consecutive records. *Random Structures Algorithms* **17**, 169–196.
- GNEDIN, A. (2007). Coherent random permutations with record statistics. To appear in *Discrete Math. Theoret. Comput. Sci.*
- HAHLIN, L. O. (1995). Double Records. Res. Rep. 1995:12, Department of Mathematics, Uppsala University.
- HOLST, L. (2006). On the number of consecutive successes in Bernoulli trials. Preprint.
- JOFFE, A., MARCHAND, E., PERRON, F. AND POPADIUK, P. (2004). On sums of products of Bernoulli variables and random permutations. *J. Theoret. Prob.* **17**, 285–292.
- KNUTH, D. (1992). Two notes on notations. *Amer. Math. Monthly* **99**, 403–422.
- MORI, T. F. (2001). On the distribution of sums of overlapping products. *Acta Scientiarum Mathematica (Szeged)* **67**, 833–841.
- SETHURAMAN, J. AND SETHURAMAN, S. (2004). On counts of Bernoulli strings and connections to rank orders and random permutations. In *A festschrift for Herman Rubin* (IMS Lecture Notes Monogr. Ser. **45**), Institute of Mathematical Statistics, Beachwood, OH, pp. 140–152.