

COUNTS OF FAILURE STRINGS IN CERTAIN BERNOULLI SEQUENCES

LARS HOLST

Department of Mathematics, Royal Institute of Technology
SE-100 44 Stockholm, Sweden

E-mail: lholst@math.kth.se

October 6, 2006

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 given number of failures between two subsequent successes is studied. Explicit expressions for distributions and moments are obtained for $p_k = a/(a + b + k - 1)$, $a > 0$, $b \geq 0$. 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 moments; Ewens Sampling Formula; Hoppe's urn; Poisson distribution; Poisson–Dirichlet distribution; Pólya's urn; random permutations; records; spacings; sums of indicators

AMS 2000 SUBJECT CLASSIFICATION: PRIMARY 60C05
SECONDARY 60K99

Running title: Counts of failure strings

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 $SF \dots FS$ of $d - 1$ failures between two subsequent successes. We will study the number of such strings. Explicit results are obtained for $p_k = a/(a + b + k - 1)$, $a > 0$, $b \geq 0$. To our knowledge only special cases have been studied previously.

For $a = 1, b = 0$, 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 After that, an unpublished proof by Diaconis inspired a number of studies on 1-strings, see Chern *et al* (2000), Mori (2001), Joffe *et al* (2004) and the references therein.

Sethuraman and Sethuraman (2004) studied d -strings for $a = 1, b > 0$, 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 random permutations, see e.g. Arratia, Barbour and Tavaré (2003) page 95.

In Section 2 we introduce notations 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 $p_k = a/(a + k - 1)$, connected with random permutations, is studied in Section 3. In Section 4 we derive 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 $p_k = a/(a + b + k - 1)$.

2 General case: notations and moments

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

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

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

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

with mean

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

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

Proposition 2.1 *For the binomial moments*

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

with $d \leq n - 1$ and $\sum_{j=1}^d j r_j \leq n - 1$, the recursion holds:

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

Proof. Using generating functions and the independence between the I_k 's we get

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

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

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

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

with mean

$$E(N_{dn}) = E(M_{dn}) + p_{n-d+1} (1 - p_{n-d+2}) \cdots (1 - p_n).$$

Proposition 2.2 For $p_{n+1} > 0$ it holds for the binomial moments:

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

Proof. By the law of total probability we have

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

from which the assertion follows by expansion in series. \square

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 one if and only if

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

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

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

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

$$x^n = x(x - 1) \cdots (x - n + 1), \quad x^{\bar{n}} = x(x + 1) \cdots (x + n - 1) = \sum_{k=1}^n \begin{bmatrix} n \\ k \end{bmatrix} x^k,$$

where $\begin{bmatrix} n \\ k \end{bmatrix}$ is a cycle number or signless Stirling number of the first kind. We assume in the rest of this section that $p_k = a/(a + k - 1)$ with $a > 0$. Closed simple formulas can be obtained for the binomial moments. Note that $\sum_{j=1}^{n-1} j M_{jn} \leq n - 1$ and $\sum_{j=1}^n j N_{jn} = n$ (with probability one).

Proposition 3.1 For $m = \sum_{j=1}^d j r_j \leq n - 1$:

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

Proof. By telescoping sums we get

$$E(M_{dn}) = \sum_{k=1}^{n-d} \frac{a}{a+k-1} \left(1 - \frac{a}{a+k}\right) \cdots \left(1 - \frac{a}{a+k+d-2}\right) \frac{a}{a+k+d-1}$$

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

Hence the assertion holds for $E(M_{dn})$. By elementary calculations the recursion in Proposition 2.1 is verified. The proof is finished by induction. \square

Proposition 3.2 For $d \leq n$ and $m = \sum_{j=1}^d jr_j \leq n$:

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

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 2.2 and 3.1 the first assertion follows from an elementary calculation.

We have using generating functions that

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

As $\sum_1^n jN_{jn} = n$ the binomial moments disappears 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. \square

The distribution of (N_{1n}, \dots, N_{nn}) is the famed Ewens Sampling Formula. Furthermore, N_{dn} is the number of d -strings in $1I_2I_3 \dots I_n 1$. Using this, the last proposition can be derived by combinatorial arguments, cf. Arratia *et al* (2003) page 95. In that 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 the distribution

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

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, well known for a -biased random permutations, see Aratia *et al* (2003) page 96.

Proposition 3.3 *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 $p_k = a/(a + b + k - 1)$

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 one. Mori (2001) derived the distribution of $M_{1\infty}$. For the special case $a = 1$ Sethuraman and Sethuraman (2004) obtained the joint distribution of $M_{1\infty}, M_{2\infty}, \dots$. Using different methods we generalize 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 4.1 *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 and one black ball of weights a and b respectively. Balls are drawn at random proportional to weights. The white 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 one. Obviously, the probability of drawing the white ball at drawing k is $p_k = a/(a + b + k - 1)$.

Generate a sequence of W 's and B 's. We get 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 the $\text{Beta}(a, b)$ random variable U and then, conditional on the outcome $U = u$, generating a sequence of independent $\text{Be}(u)$ random variables, with 1 corresponding to W and 0 to B .

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 3.3 the number of d -strings in the subsequence, $M_{d\infty}^*$, is $\text{Po}(a/d)$ and $M_{1\infty}^*, M_{2\infty}^*, \dots$ are independent.

Recall the following well known fact. If the random variable ξ is $\text{Po}(\mu)$ and independent of the independent $\text{Be}(p)$ random variables $\varepsilon_1, \varepsilon_2, \dots$, then $\sum_{j=1}^{\xi} \varepsilon_j$ and $\sum_{j=1}^{\xi} (1 - \varepsilon_j)$ are independent $\text{Po}(\mu p)$ and $\text{Po}(\mu(1 - p))$ respectively.

Consider the 1-strings in the subsequence of W 's. Each such 1-string is also a 1-string in the original sequence, provided it was not interrupted by a B . The probability for interruption is $1 - u$. As $M_{1\infty}^*$ is $\text{Po}(a)$ it follows from the fact above and the independence, that conditional on $U = u$, the total number of 1-strings in the original sequence, $M_{1\infty}$, is $\text{Po}(au)$ and independent of $M_{1\infty}^* - M_{1\infty}$ which is $\text{Po}(a(1 - u))$.

For the 2-strings we can argue in a similar way as above. Conditional on $U = u$, the random variable $M_{2\infty}$ is Poisson with mean

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

and independent of $M_{1\infty}$. The argument extends, $M_{d\infty}$ is Poisson with (conditional) mean

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

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 set $A_{1n}=0$. 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

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

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 last success, A_{2n} , that

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

in distribution, where U_1, U_2 are independent $\text{Beta}(1, a)$ random variables. The procedure can be repeated in like manner.

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 size ordered A 's is the Poisson–Dirichlet distribution with parameter a . In particular we have:

Theorem 4.2 *For the longest string of failures in the first n trials:*

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

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

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

References

- [1] ARRATIA, R., BARBOUR, A.D. AND TAVARÉ, S. (2003). *Logarithmic Combinatorial Structures: a Probabilistic Approach*. European Mathematical Society Publishing House, ETH-Zentrum, Zürich.
- [2] CHERN, H.-H., HWANG, H.-K. AND YEH, Y.-N. (2000). Distribution of the number of consecutive records. *Random Structures & Algorithms*, **17** 169–196.
- [3] HAHLIN, L.O. (1995). Double Records. *Uppsala University Department of Mathematics Report*, **1995:12**. Licentiat thesis.
- [4] JOFFE, A., MARCHAND, E., PERRON, F. AND POPADIUK, P. (2004). On sums of products of Bernoulli variables and random permutations. *Journal of Theoretical Probability*, **17** 285–292.
- [5] KNUTH, D. (1992). Two notes on notations. *The American Mathematical Monthly*, **99** 403–422.
- [6] MORI, T.F. (2001). On the distribution of sums of overlapping products. *Acta Scientiarum Mathematica (Szeged)*, **67** 833–841.
- [7] 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 Monograph Series*, **45** 140–152, Institute of Mathematical Statistics, Beachwood, Ohio.