

# Explicit formulae for the distributions of runs

Hayato Takahashi

Random Data Lab. Inc.

Jan. 25, 2023  
IEICE IT Maebashi

## Problem: the number of the occurrences of words in finite strings

The number of the occurrences of words in finite strings plays important role in information theory, genome analysis, statistics, AI, etc.

Example: The words 10 and 00 appear in 10001110010 three times.

Run:  $0^m, m = 2, 3, \dots$

We study the enumeration (and distribution) of the number of the occurrences of runs with several types of counting in finite strings.

Remark: The distributions of the number of the occurrence of letters 1 and 0 are given by binomial distribution.

# Contents

1. Known results, generating functions.
2. Known results, runs.
3. Main theorem.
4. Outline of Proof.

## Known results, generating functions

In Regnier et.al [12], Bassino et.al [1], and Robin [13], the number of the occurrences of words given as generating functions.

$f(n, k, w)$ : the number of  $x_1 \cdots x_n$  in which  $w$  appears  $k$  times. Then

$$\sum_{n,k} f(n, k, w) z_1^n z_2^k = \frac{g(z_1, z_2)}{h(z_1, z_2)}.$$

$g, h$ : polynomial.

## Known results, generating functions 2, example

例 : Guibas and Odlyzko [5]

$$\begin{aligned}\sum_n f(n, 0, 10)z^n &= \frac{1}{(1-z)^2} \\ &= \left(\sum z^n\right)^2 \\ &= \sum (n+1)z^n.\end{aligned}$$

$$f(n, 0, 10) = n + 1 \text{ for all } n = 1, 2, \dots$$

$$000, 001, 011, 111 \text{ and } f(3, 0, 10) = 4.$$

## Known results, generating functions 3

$f(n, k, w)$ : the number of  $x_1 \cdots x_n$  in which  $w$  appears  $k$  times. Then

$$\sum_{n,k} f(n, k, w) z_1^n z_2^k = \frac{g(z_1, z_2)}{h(z_1, z_2)}, \quad g, h: \text{polynomial.}$$

注1：既存の方法では長さ  $n$  に関する再帰的な関係から生成関数を導出したために上式で  $n$  を固定した有限次元の生成関数で表すことはできない。Generating functions are derived by induction on length  $n$  and we do not have generating function with fixed  $n$ .

注2：生成関数をべき級数展開すれば  $f(n, k, w)$  が求まるが一般には有利関数のべき級数展開は簡単な場合を除いて難しい。It is difficult to expand rational function into power series except for simple cases.

注3：生成関数から  $f(n, k, w)$  の近似解や再帰的計算法を導くことができる。We have approximation of  $f(n, k, w)$  from generating function. Some recurrence formula for  $f(n, k, w)$  are derived from generating function.

## Known results, runs, various type of counting

Fu et.al (1994) [3] showed the distributions of the following statistics by Markov imbedding method.

- (i)  $E_{n,m}$ , the number of  $0^m$  of size exactly  $m$ . Mood(1940) [9].
- (ii)  $G_{n,m}$ , the number of  $0^m$  of size greater than or equal to  $m$ .
- (iii)  $N_{n,m}$ , the number of nonoverlapping consecutive  $0^m$ . Feller [2].
- (iv)  $M_{n,m}$ , the number of overlapping consecutive  $0^m$ .
- (v)  $L_n$ , the size of the longest run of 0s.

Example:  $E_{9,2} = 1$ ,  $G_{9,2} = 2$ ,  $N_{9,2} = 3$ ,  $M_{9,2} = 4$ , and  $L_9 = 4$  for

run 00 and 110000100.

## Known results, other explicit formulae

Explicit formulae for the distributions of runs are given in

$G_{n,m}$ : Makri et.al (2007) [8]

$N_{n,m}$ : Hirano (1986) [6], Phillipou et.al (1986) [11], Godbole (1990) [4],  
Muselli (1996) [10]

$M_{n,m}$ : Ling (1988) [7].

$L_n$ : Makri et.al (2007) [8]



# Main Theorem

## Theorem

Let  $P$  be an i.i.d. probability on  $\{0, 1\}^n$ .

(i)  $P(E_{n,m} = t) = (P(E'_{n+1,m} = t) - P(0)P(E'_{n,m} = t))/P(1)$  where

$P(E'_{n,m} = t) =$

$$\sum_{\substack{k_1, k_2: \\ (m+1)k_1 + (m+2)k_2 \leq n, \\ t \leq k_1 + k_2}} (-1)^{k_1 - t} \binom{n - (m+1)k_1 - (m+2)k_2 + k_1 + k_2}{k_1, k_2} \\ \times \binom{k_1 + k_2}{t} P^{k_1}(10^m) P^{k_2}(10^{m+1}).$$

## Main Theorem 2

### Theorem (Continue)

Let  $P$  be an i.i.d. probability on  $\{0, 1\}^n$ .

(ii)  $P(G_{n,m} = t) = (P(G'_{n+1,m} = t) - P(0)P(G'_{n,m} = t))/P(1)$  where

$$P(G'_{n,m} = t) = \sum_{k: t \leq k, (m+1)k \leq n} (-1)^{k-t} \binom{n - (m+1)k + k}{t, k-t} P^k (10^m).$$

## Main Theorem 3

### Theorem (Continue)

Let  $P$  be an i.i.d. probability on  $\{0, 1\}^n$ .

(iii)  $P(N_{n,m} = t) = (P(N'_{n+1,m} = t) - P(0)P(N'_{n,m} = t))P^{-1}(1)$  where

$P(N'_{n,m} = t) =$

$$\sum_{\substack{r, k_1, \dots, k_T: \\ \sum_i (mi+1)k_i \leq n, 0 \leq r \leq \sum_i k_i \\ t = \sum_i ik_i - r}} (-1)^r \binom{n - \sum_i (mi+1)k_i + \sum_i k_i}{k_1, \dots, k_{n-m}} \binom{\sum_i k_i}{r} \\ \times \prod_{i=1}^T P^{k_i}(10^{im}).$$

and  $T$  is a maximum integer such that  $Tm + 1 \leq n$ .

## Main Theorem 4

### Theorem (Continue)

Let  $P$  be an i.i.d. probability on  $\{0, 1\}^n$ .

(iv)  $P(M_{n,m} = t) = (P(M'_{n+1,m} = t) - P(0)P(M'_{n,m} = t))P^{-1}(1)$  where

$$P(M'_{n,m} = t) = \sum_{\substack{r, k_1, \dots, k_{n-m}: \\ \sum_i (m+i)k_i \leq n, 0 \leq r \leq \sum_i k_i \\ t = \sum_i ik_i - r}} (-1)^r \binom{n - \sum_i (m+i)k_i + \sum_i k_i}{k_1, \dots, k_{n-m}} \\ \times \binom{\sum_i k_i}{r} \prod_{i=1}^{n-m} P^{k_i} (10^{m+i-1}).$$

(v)  $P(L_n = t) = P(N_{n,t+1} = 0) - P(N_{n,t} = 0)$ .

## Outline of Proof of main theorem (iv).

$\mathbf{N}(w_1, \dots, w_l; X_1^n)$  : the number of the overlapping appearances of  $w_1, w_2, \dots, w_l$  in  $X_1^n$ .

Suppose that

$w_1 \sqsubset w_2 \sqsubset \dots \sqsubset w_l$  and  $\mathbf{N}(w_1, \dots, w_l; X_1^n) = (s_1, \dots, s_l)$ .

Then

$s_1$  is the number of the appearances of  $w_1$  and  $w_2$ ,

$s_2$  is the number of the appearances of  $w_2$  and  $w_3$ ,

....

$$\mathbf{N}'(w_1, \dots, w_l; X_1^n) := (s_1 - s_2, s_2 - s_3, \dots, s_l)$$

$$\text{if } \mathbf{N}(w_1, \dots, w_l; X_1^n) = (s_1, s_2, \dots, s_l).$$

Example:  $\mathbf{N}(00, 000; 1010001) = (2, 1)$  and  $\mathbf{N}'(00, 000; 1010001) = (1, 1)$ .

# Outline of Proof of main theorem (iv).

## Lemma

Let  $w_1 \sqsubset w_2 \cdots \sqsubset w_l$  be an increasing sequence of words,

$$A(k_1, \dots, k_l) := \binom{n - \sum_i m_i k_i + \sum_i k_i}{k_1, \dots, k_l} \prod_{i=1}^l P^{k_i}(w_i),$$

$$B(k_1, \dots, k_l) := P(\mathbf{N}'(w_1, \dots, w_l; X^n) = (k_1, k_2, \dots, k_l)),$$

$$F_A(z_1, \dots, z_l) := \sum_{\substack{k_1, \dots, k_l: \\ \sum_i m_i k_i \leq n}} A(k_1, \dots, k_l) z^{k_1} \cdots z^{k_l}.$$

Then

$$F_A(Y-1, (Y-1)Y, \dots, (Y-1)Y^{l-1}) = \sum_{\substack{k_1, \dots, k_l: \\ \sum_i m_i k_i \leq n}} B(k_1, \dots, k_l) Y^{\sum ik_i}.$$

## Outline of Lemma

Takahashi [15, 14].

$$F_A(z_1, \dots, z_l) = F_B(z_1 + 1, z_1 + z_2 + 1, \dots, \sum_i z_i + 1). \quad (1)$$

Set  $z_1 = X, z_2 = X(X + 1), \dots, z_l = X(X + 1)^{l-1}$  in (1). Then

$$F_A(X, X(X + 1), \dots, X(X + 1)^{l-1}) = F_B(X + 1, (X + 1)^2, \dots, (X + 1)^l)$$

$$\begin{aligned} F_A(Y - 1, (Y - 1)Y, \dots, (Y - 1)Y^{l-1}) &= F_B(Y, Y^2, \dots, Y^l) \\ &= \sum_{\substack{k_1, \dots, k_l: \\ \sum_i m_i k_i \leq n}} B(k_1, \dots, k_l) Y^{\sum i k_i}. \end{aligned}$$

## Lemma

Let

$$C_{n,t} := \{x_1^n \mid t = \sum ik_i, (k_1, \dots, k_{n-m}) = \mathbf{N}'(w_1, \dots, w_{n-m}; x_1^n)\},$$

$$C_{n,t}^0 := \{0x_1^n \mid t = \sum ik_i, (k_1, \dots, k_{n+1-m}) = \mathbf{N}'(w_1, \dots, w_{n+1-m}; 0x_1^n)\},$$

$$C_{n,t}^1 := \{1x_1^n \mid t = \sum ik_i, (k_1, \dots, k_{n+1-m}) = \mathbf{N}'(w_1, \dots, w_{n+1-m}; 1x_1^n)\},$$

$$D_{n,t} := \{x_1^n \mid \mathbf{N}(0^m; x_1^n) = t\}.$$

Then

$$C_{n+1,t} = C_{n,t}^0 \cup C_{n,t}^1, \quad C_{n,t}^0 \cap C_{n,t}^1 = \emptyset,$$

$$C_{n,t}^0 = \{0x_1^n \mid x_1^n \in C_{n,t}\}, \quad \text{and} \quad C_{n,t}^1 = \{1x_1^n \mid x_1^n \in D_{n,t}\}.$$



# Reference I

- [1] F. Bassino, J. Clément, and P. Mico-dème.  
Counting occurrences for a finite set of words: combinatorial methods.  
*ACM Trans. Algor.*, 9(4):Article No. 31, 2010.
- [2] W. Feller.  
*An Introduction to probability theory and its applications Vol. 1.*  
Wiley, 3rd edition, 1970.
- [3] J. C. Fu and M. V. Koutras.  
Distribution theory of runs: a Markov chain approach.  
*J. Amer. Statist. Assoc.*, 89(427):1050–1058, 1994.
- [4] A. P. Godbole.  
Specific formulae for some success run distributions.  
*Statist. Probab. Lett.*, 10:119–124, 1990.
- [5] L. Guibas and A. Odlyzko.  
String overlaps, pattern matching, and nontransitive games.  
*J. Combin. Theory Ser. A*, 30:183–208, 1981.

## Reference II

- [6] K. Hirano.  
Some properties of the distributions of order  $k$ .  
pages 43–53, 1986.  
Fibonacci Numbers and their Applications, A. N. Phillipou, A. F. Horadam and G. E. Bergum eds, Reidel.
- [7] K. D. Ling.  
On binomial distributions of order  $k$ .  
*Statist. Probab. Letters*, 6:247–250, 1988.
- [8] F. S. Makri, A. N. Philippou, and Z. M. Psillakis.  
Shortest and longest length of success runs in binary sequences.  
*J. Statist. Plan. Inference*, 137:2226–2239, 2007.
- [9] A. M. Mood.  
The distribution theory of runs.  
*Ann. Math. Statist*, 11(4):367–392, 1940.
- [10] M. Muselli.  
Simple expressions for success run distributions in Bernoulli trials.  
*Statist. Probab. Lett.*, 31:121–128, 1996.

## Reference III

- [11] A. N. Phillipou and F. S. Makri.  
Success, runs and longest runs.  
*Statist. Probab. Lett.*, 4:211–215, 1986.
- [12] M. Régnier and W. Szpankowski.  
On pattern frequency occurrences in a markovian sequence.  
*Algorithmica*, 22(4):631–649, 1998.
- [13] S. Robin and J. J. Daudin.  
Exact distribution of word occurrences in a random sequence of letters.  
*J. Appl. Prob.*, 36(1):179–193, 1999.
- [14] H. Takahashi.  
The distributions of sliding block patterns in finite samples and the inclusion-exclusion principles for partially ordered sets, Dec. 2018.  
Probability Symposium, Kyoto Univ. Japan [arxiv:1811.12037v1](https://arxiv.org/abs/1811.12037v1).
- [15] H. Takahashi.  
Inclusion-exclusion principles on partially ordered sets and the distributions of the number of pattern occurrences in finite samples, Sep. 2018.  
Mathematical Society of Japan, Statistical Mathematics Session, Okayama Univ. Japan.