COMBINATORICS OF GEOMETRICALLY DISTRIBUTED RANDOM VARIABLES: RUN STATISTICS

PETER J. GRABNER[†], ARNOLD KNOPFMACHER, AND HELMUT PRODINGER

ABSTRACT. For words of length n, generated by independent geometric random variables, we consider the mean and variance, and thereafter the distribution of the number of runs of equal letters in the words. In addition, we consider the mean length of a run as well as the length of the longest run over all words of length n.

1. INTRODUCTION

Let X denote a geometrically distributed random variable, i. e. $\mathbb{P}\{X = k\} = pq^{k-1}$ for $k \in \mathbb{N}$ and q = 1 - p. The combinatorics of n geometrically distributed independent random variables X_1, \ldots, X_n has attracted recent interest, especially because of applications in computer science. We mention just two areas, the skip list [1, 15, 17, 11] and probabilistic counting [5, 9, 10, 12].

In [16] the number of left-to-right maxima was investigated for words $a_1 \ldots a_n$, where the letters a_i are independently generated according to the geometric distribution. In [13] the study of left-to-right maxima was continued, but now the parameters studied were the mean value and mean position of the *r*-th maximum.

In this article we study runs of consecutive equal letters in a string of n geometrically distributed independent random letters. For example in w = 22211114431 we have 5 runs of equal letters of respective lengths 3, 4, 2, 1, 1. In the sequel we denote by $R_n(w)$ the number of runs in the word w, where w is of length n.

In section 2 we study the mean and variance of R_n . Thereafter, in section 3 we study the distribution of the number of runs, which turns out to be Gaussian. Subsequently, in section 4 we study the average length of the runs per word. Finally, in section 5 we determine the mean and variance of the length of the longest run in a word of length n. The treatment follows to some extent the classical paper [14]; compare also [7] for such longest run statistics in a slightly different context.

Date: September 30, 2003.

¹⁹⁹¹ Mathematics Subject Classification. Primary: 05A15; Secondary: 60C05.

Key words and phrases. geometric random variables, generating function, limiting distribution, runs.

[†] This author is supported by the START-project Y96-MAT of the Austrian Science Fund.

Part of this work was done during the first author's visit to the Centre for Applicable Analysis and Number Theory at the University of the Witwatersrand, Johannesburg, South Africa.

P. J. GRABNER, A. KNOPFMACHER, AND H. PRODINGER

2. Moments of number of runs

In order to determine the mean and variance of the number of runs we will make use of the following decomposition of the set of all (non-empty) words. Here $\{\geq \mathbf{k}\}$ denotes the set $\{k, k+1, \ldots\}$; for a given set A we denote

$$A^+ = \bigcup_{k=1}^{\infty} A^k, \quad A^* = \varepsilon \cup A^+,$$

where ε stands for the empty word. We decompose the set of non-empty words according to runs of 1's, separated by words consisting of larger digits only

(2.1)
$$\{\geq \mathbf{1}\}^{+} = (\varepsilon + \mathbf{1}^{+}) (\{\geq \mathbf{2}\}^{+} \mathbf{1}^{+})^{*} \{\geq \mathbf{2}\}^{+} (\varepsilon + \mathbf{1}^{+}) + \mathbf{1}^{+};$$

here we find it more convenient to write + instead of \cup .

We consider a probability generating function F(z, u), where z labels the length of the word, and u counts the number of runs. We should always have $F(z, 1) = \frac{z}{1-z}$, and a replacement of z by qz, if we increase all letters by 1.

Then (2.1) translates into the functional equation

(2.2)
$$F(z,u) = \frac{F(qz,u)}{1 - F(qz,u)\frac{pzu}{1 - pz}} \left(\frac{pzu}{1 - pz} + 1\right)^2 + \frac{pzu}{1 - pz}$$

Now we differentiate it w. r. t. u, plug in u = 1, set $G(z) = \frac{\partial}{\partial u}F(z, 1)$, and get

$$G(z) = G(qz)\frac{(1-qz)^2}{(1-z)^2} + \frac{pz(1-pz)}{(1-z)^2}$$

Setting $H(z) = (1-z)^2 G(z)$ yields

$$H(z) = H(qz) + pz(1 - pz) .$$

Comparing coefficients, we see that

$$[z]H(z) = 1$$
,
 $[z^2]H(z) = -\frac{p^2}{1-q^2} = -\frac{p}{1+q}$,

and that the other coefficients are zero. Consequently,

$$H(z) = z - \frac{p}{1+q}z^2 ,$$

and

$$G(z) = \frac{z - \frac{p}{1+q}z^2}{(1-z)^2}$$

This leads to

Proposition 1. The mean value of the number of runs for $n \ge 1$ is given by

$$\mu_n = \mathbb{E}R_n = [z^n]G(z) = \frac{2q}{1+q}n + \frac{p}{1+q}.$$

RUN STATISTICS

Now that we see such a simple result, we are tempted to look for a simple proof as well. Indeed, since the expectation is additive, we have n-1 times the expectation of having unequal digits at two adjacent random variables, plus 1. This expectation is given by

$$\sum_{i \neq j} pq^{i-1} pq^{j-1} = \frac{2q}{1+q}.$$

Then the expected number of runs is given by

$$(n-1)\frac{2q}{1+q} + 1 = \frac{2q}{1+q}n + \frac{p}{1+q}.$$

If we want to compute the variance, such a simple argument seems to be out of reach. Henceforth, we differentiate (2.2) twice, and use the notation $V(z) = \frac{\partial^2}{\partial u^2} F(z, 1)$. We see

$$V(z) = V(qz)\frac{(1-qz)^2}{(1-z)^2} + \frac{8q^2}{(1+q)^2(1-z)^3} + \frac{8pq^2z^2}{(1+q)^2(1-z)^2} + \frac{4p^2qz^2}{(1+q)(1-z)^2} + \frac{2p^3qz^3}{(1+q)^2(1-z)^2} - \frac{8q^2}{(1+q)^2(1-qz)(1-z)^2}$$

or, with $W(z) = (1-z)^2V(z)$,

$$W(z) = W(qz) + \frac{8q^2}{(1+q)^2(1-z)} + \frac{8pq^2z^2}{(1+q)^2} + \frac{4p^2qz^2}{(1+q)} + \frac{2p^3qz^3}{(1+q)^2} - \frac{8q^2}{(1+q)^2(1-qz)} .$$

From this we see that (with $w_n := [z^n]W(z)$) for $n \ge 4$,

$$w_n = q^n w_n + \frac{8q^2}{(1+q)^2} [1-q^n]$$

or

$$w_n = \frac{8q^2}{(1+q)^2}$$

Furthermore we see that

$$w_0 = w_1 = 0$$
, $w_2 = \frac{4q}{1+q}$, and $w_3 = \frac{2q(1+q+4q^2)}{(1+q)(1+q+q^2)}$

Hence

$$V(z) = \frac{4q}{1+q} \frac{z^2}{(1-z)^2} + \frac{2q(1+q+4q^2)}{(1+q)(1+q+q^2)} \frac{z^3}{(1-z)^2} + \frac{8q^2}{(1+q)^2} \frac{z^4}{(1-z)^3}$$

and

$$[z^{n}]V(z) = \frac{4q^{2}}{(1+q)^{2}}(n-2)(n-3) + \frac{2q(1+q+4q^{2})}{(1+q)(1+q+q^{2})}(n-2) + \frac{4q}{1+q}(n-1).$$

Adding the expectation and subtracting the square of the expectation, we obtain the variance.

Proposition 2. The variance of the number of runs is given for $n \ge 2$ by

$$\sigma_n^2 = \mathbb{V}R_n = \frac{2q(1-q)^2(2+q^2)}{(1+q)^2(1-q^3)}n - \frac{2q(1-q)^2(3-q+q^2)}{(1+q)^2(1-q^3)}$$

3. Distribution of the number of runs

In this section we prove a central limit theorem for the distribution of the number of runs. In order to do this, we have to extract further information from the functional equation (2.2). We observe that the terms on the right-hand side are all simple rational functions, except for the terms containing F(qz, u). From the definition of F(z, u) it is clear that F(z, u) can be written as

$$F(z,u) = \sum_{n \ge 1} z^n f_n(u)$$

for polynomials $f_n(u)$ with deg $f_n = n$, whose coefficients are positive and ≤ 1 . Therefore we have

(3.1)
$$|f_n(u)| \le \frac{|u|^{n+1} - 1}{|u| - 1}, \text{ for } |u| > 1.$$

Using q < 1 and (3.1) we obtain that F(qz, u) is holomorphic in $|z| < \frac{1}{\sqrt{q}}$, $|u| < \frac{1}{\sqrt{q}}$. Since for u = z = 1 we have

$$1 - F(q,1)\frac{p}{1-p} = 0$$
, and $\frac{\partial}{\partial z} \left(1 - F(qz,u)\frac{pzu}{1-pz} \right) \Big|_{z=1,u=1} \neq 0$

there exists a function f(u) holomorphic in a neighbourhood of u = 1 such that $z = f(u)^{-1}$ solves

$$1 - F(qz, u)\frac{pzu}{1 - pz} = 0$$

and satisfies f(1) = 1. Furthermore, $|f(e^{it})| < 1$ for $0 < |t| < \varepsilon$ for some $\varepsilon > 0$ by an application of Rouché's theorem. Thus we can write

(3.2)
$$F(z,u) = \frac{g(z,u)}{1 - f(u)z} + R(z,u),$$

where g(z, u) and R(z, u) are holomorphic in $|z| < 1 + \delta$, $|u - 1| < \delta$ for some $\delta > 0$. Now we are in the general framework of Hwang's quasi-power theorem (cf. [8]) and can deduce the following theorem.

Theorem 1. The number of runs in words of length n produced by independent geometric random variables obeys a central limit law, more precisely

(3.3)
$$\mathbb{P}\left(R_n \le \frac{2q}{1+q}n + t\sqrt{\frac{2q(2+q^2)}{1-q^3}}\frac{1-q}{1+q}\sqrt{n}\right) = \Phi(t) + \mathcal{O}(n^{-\frac{1}{2}});$$

the error term is uniform in t.

RUN STATISTICS

4. Average length of runs

Given a string w of geometric random variables of length n with k runs we define the average length of a run to be $L_n(w) = \frac{n}{k}$. It is of interest to determine the moments and the distribution of this parameter over all strings of length n. Intuitively, one expects that the mean length of a run should be close to n divided by the mean number of runs, which is

$$\frac{n}{\frac{2q}{1+q}n + \frac{p}{1+q}} = \frac{1+q}{2q} - \frac{1-q^2}{4q^2}\frac{1}{n} + \mathcal{O}(\frac{1}{n^2}).$$

In fact we obtain

Proposition 3. For $n \ge 1$ the mean and variance of L_n are given respectively by

$$\frac{1+q}{2q} + \mathcal{O}(\frac{1}{n}), \quad \frac{(1-q^2)^2(2+q^2)}{8q^3(1-q^3)}\frac{1}{n} + \mathcal{O}(\frac{1}{n^2}).$$

Moreover, L_n obeys a central limit theorem:

$$\mathbb{P}\left(L_n - \frac{1+q}{2q} \le \frac{(1-q^2)\sqrt{2+q^2}}{\sqrt{8q^3(1-q^3)}} \frac{t}{\sqrt{n}}\right) = \Phi(t) + \mathcal{O}(n^{-\frac{1}{2}}).$$

Proof. The proof will make use of the distribution obtained for the number of runs in Theorem 1. We first write $R_n = \mu_n + \sigma_n X_n$, where X_n is a sequence of random variables with asymptotically normal distribution, and μ_n and σ_n are given by Propositions 1, 2. Then we can write

(4.1)
$$\frac{n}{R_n} = \frac{n}{\mu_n} \frac{1}{1 + \frac{\sigma_n}{\mu_n} X_n} = \frac{n}{\mu_n} - \frac{n\sigma_n}{\mu_n^2} X_n + \mathcal{O}\left(\frac{n\sigma_n^2}{\mu_n^3} |X_n|^2\right).$$

We observe that by (3.3) and $1 - \Phi(t) \sim \frac{1}{\sqrt{2\pi x}} \exp(-\frac{x^2}{2})$ (cf. [3]) we have

$$\mathbb{P}\left(|X_n| > \log n\right) = \mathcal{O}(n^{-\frac{1}{2}}).$$

This gives an error term of $\mathcal{O}(n^{-1}\log n)$ in (4.1) and yields the desired result.

5. Longest runs

In this section we study the mean of the longest run M_n of equal digits in a string of length n. For this purpose we introduce the probability generating function G_h of all strings that have runs only of length less than h. Similar arguments as in the proof of (2.2) show that G_h satisfies

(5.1)
$$G_h(z) = \left(\frac{1 - (pz)^h}{1 - pz}\right)^2 \frac{G_h(qz)}{1 - G_h(qz)\frac{pz}{1 - pz}\left(1 - (pz)^{h-1}\right)} + pz\frac{1 - (pz)^{h-1}}{1 - pz}$$

In order to extract the asymptotic behaviour of the probability that a string of length n has runs of length at most h, we have to find the singularities of $G_h(z)$. We start with simple a priori estimates for the power series coefficients of G_h :

(5.2)
$$\frac{z}{1-z} - \frac{1}{(1-z)^2} \sum_{k=0}^{\infty} \frac{(pq^k z)^h}{1-pq^k z} \preceq G_h(z) \preceq \frac{z}{1-z},$$

where

$$A(z) \preceq B(z) \Leftrightarrow \forall n \in \mathbb{N} : [z^n] A(z) \le [z^n] B(z)$$

The left hand side of (5.2) is obtained by subtracting all strings which contain at least one run of 1's of length $\geq h$, or one run of 2's of length $\geq h$, and so on; clearly this is not a disjoint union, therefore it only gives an estimate. The upper estimate is trivial.

Now we investigate the solution of the equation

(5.3)
$$1 - G_h(qz) \frac{pz}{1 - pz} \left(1 - (pz)^{h-1}\right) = 0.$$

Since all the occurring power series have positive coefficients, we could apply Rouché's theorem to conclude that the positive real root ρ_h of (5.3) is the root of smallest modulus. Then inserting the two estimates (5.2) and multiplying out, we obtain

$$1 + q\rho_h(p\rho_h)^h \le \rho_h \le 1 + q\rho_h(p\rho_h)^h + \frac{p\rho_h - (p\rho_h)^h}{1 - q\rho_h} \sum_{k=1}^{\infty} \frac{(pq^k\rho_h)^h}{1 - pq^k\rho_h}.$$

By bootstrapping we easily see that this implies $\rho_h = 1 + pq^h + \mathcal{O}((pq)^h)$. Since the poles of the other terms in (5.1) are $\frac{1}{p}$ and $\frac{1}{q}$, ρ_h is the dominant singularity of the function G_h . Furthermore, we have

$$G_{h}(q\rho_{h}) = \frac{q}{p} + \mathcal{O}(q^{h}), \quad G_{h}'(q\rho_{h}) = \frac{1}{p^{2}} + \mathcal{O}(hq^{h})$$
$$\frac{d}{dz} \left(1 - G_{h}(qz) \frac{pz}{1 - pz} \left(1 - (pz)^{h-1} \right) \right) \Big|_{z=\rho_{h}} = -\frac{1}{pq} + \mathcal{O}(hp^{h})$$

Putting everything together we obtain

(5.4)
$$\mathbb{P}(M_n < h) = (1 - pq^h)^n + \mathcal{O}(hq^h) + \mathcal{O}(hq^h)$$

Using (5.4) and Abel summation we derive that the first and second moment of the longest run are given by

(5.5)
$$\mathbb{E}M_n = \sum_{h \ge 1} (1 - \mathbb{P}(M_n < h)) = \sum_{h \ge 1} (1 - (1 - pq^h)^n) + \mathcal{O}(1)$$
$$\mathbb{E}M_n^2 = 2\sum_{h \ge 1} h (1 - \mathbb{P}(M_n < h)) - \mathbb{E}M_n = \sum_{h \ge 1} (2h - 1) (1 - (1 - pq^h)^n) + \mathcal{O}(1).$$

In order to compute the asymptotic behaviour of these two moments, we use the now classical exponential approximation technique (cf. [6]). We replace the terms $(1-pq^h)^n$ in the two sums by $\exp(-npq^h)$ and estimate the error. For this purpose we split the range of summation into three parts: $h < \frac{3}{4} \log_{\frac{1}{q}} n$, $|h - \log_{\frac{1}{q}} n| \le \frac{1}{4} \log_{\frac{1}{q}} n$, and $h > \frac{5}{4} \log_{\frac{1}{q}} n$.

In the range $h < \frac{3}{4} \log_{\frac{1}{q}} n$ we estimate

$$(1 - pq^h)^n \le (1 - pn^{-\frac{1}{4}})^n \le \exp\left(-pn^{\frac{3}{4}}\right)$$
 and $\exp\left(-npq^h\right) \le \exp\left(-pn^{\frac{3}{4}}\right)$,

which yields

(5.6)
$$\left| \left(1 - pq^h \right)^n - \exp\left(-npq^h \right) \right| \le \exp\left(-pn^{\frac{3}{4}} \right).$$

In the range $|h - \log_{\frac{1}{q}} n| \le \frac{1}{4} \log_{\frac{1}{q}} n$ we write $h = \log_{\frac{1}{q}} n + t$ and approximate by the Taylor expansion to obtain

$$(1 - pq^h)^n = \exp\left(n\log\left(1 - \frac{p}{n}q^t\right)\right) = \exp\left(-pq^t + \mathcal{O}(\frac{q^{2t}}{n})\right).$$

Observing that the error term in this equation tends to 0 we obtain

(5.7)
$$\left| \left(1 - pq^h \right)^n - \exp\left(-npq^h \right) \right| = \mathcal{O}\left(\frac{q^{2t}}{n} \right) = \mathcal{O}(n^{-\frac{1}{2}}).$$

For the range $h > \frac{5}{4} \log_{\frac{1}{q}} n$ we use the Taylor expansion again to obtain

(5.8)
$$\left| \left(1 - pq^h \right)^n - \exp\left(-npq^h \right) \right| = \mathcal{O}(nq^{2h}).$$

Inserting (5.6), (5.7), and (5.8) into (5.5) we obtain

(5.9)
$$\mathbb{E}M_n = \sum_{h \ge 1} \left(1 - \exp(-npq^h) \right) + \mathcal{O}(1)\mathcal{O}\left(e^{-n^{\frac{3}{4}}}\log n\right) + \mathcal{O}(n^{-\frac{1}{2}}\log n) + \mathcal{O}(n^{-\frac{3}{2}}).$$

We now apply the Mellin transform (cf. [2, 4]) to the function $f(t) = \sum_{h \ge 1} (1 - \exp(-tpq^h))$. This yields the transformed function

(5.10)
$$f^*(s) = -\Gamma(s)p^{-s}\frac{1}{q^s - 1}, \text{ for } -1 < \Re s < 0.$$

Application of the Mellin inversion formula, shifting the line of integration to the right and collecting residues yields

(5.11)
$$f(t) = -\frac{1}{2\pi i} \int_{-\frac{1}{2}-i\infty}^{-\frac{1}{2}+i\infty} \Gamma(s) p^{-s} \frac{1}{q^s - 1} t^{-s} ds = -\sum_{k \in \mathbb{Z}} \operatorname{Res} f_1^*(s) t^{-s} \Big|_{s = \frac{2k\pi i}{\log \frac{1}{q}}} \\ - \frac{1}{2\pi i} \int_{\frac{1}{2}-i\infty}^{\frac{1}{2}+i\infty} \Gamma(s) p^{-s} \frac{1}{q^s - 1} t^{-s} ds.$$

The residues are easily computed to be

(5.12)
$$\operatorname{Res} \Gamma(s) p^{-s} \frac{1}{q^s - 1} t^{-s} \bigg|_{s=0} = \log_{\frac{1}{q}} t - \frac{\log p}{\log q} - \frac{1}{2} - \frac{\gamma}{\log q}$$
$$\operatorname{Res} \Gamma(s) p^{-s} \frac{1}{q^s - 1} t^{-s} \bigg|_{s=\frac{2k\pi i}{\log \frac{1}{q}}} = \Gamma\left(\frac{2k\pi i}{\log \frac{1}{q}}\right) (pt)^{-\frac{2k\pi i}{\log \frac{1}{q}}} \frac{1}{\log q}.$$

This yields

(5.13)
$$f(n) = \mathbb{E}M_n + \mathcal{O}(1) = \log_{\frac{1}{q}} n + \mathcal{O}(1)$$

Similarly, we can obtain an expression for the second moment

(5.14)
$$\mathbb{E}M_n^2 + \mathcal{O}(1) = \log_{\frac{1}{a}}^2 n + \mathcal{O}(\log n)$$

by applying the same analysis to the function $\sum_{h\geq 1}(2h-1)(1-(1-pq^h)^n)$. Unfortunately, the error term in (5.5) is too weak to obtain the main term in the asymptotic of the variance of M_n .

Proposition 4. The mean value of the length of the longest run M_n in a string of n geometric random variables satisfies

$$\mathbb{E}M_n = \log_{\frac{1}{q}} n + \mathcal{O}(1).$$

Acknowledgement. This work was done during the first author's visit to the Centre for Applicable Analysis and Number Theory at the University of the Witwatersrand, Johannesburg, South Africa. He is indebted to Arnold Knopfmacher, Doron S. Lubinsky and Helmut Prodinger for their great hospitality.

References

- 1. L. Devroye, A limit theory for random skip lists, Annals of Applied Probability 2 (1992), 597–609.
- 2. G. Doetsch, Handbuch der laplace transformation, Birkhäuser, Basel, 1955.
- 3. W. Feller, An introduction to probability theory and its applications, vol. 1, J. Wiley, 1950.
- P. Flajolet, X. Gourdon, and P. Dumas, *Mellin transforms and asymptotics: Harmonic sums*, Theoretical Computer Science 144 (1995), 3–58.
- 5. P. Flajolet and G. N. Martin, *Probabilistic counting algorithms for data base applications*, Journal of Computer and System Sciences **31** (1985), 182–209.
- 6. J. Galambos, The asymptotics of extreme order statistics, J. Wiley, 1987.
- L. J. Guibas and A. M. Odlyzko, Long repetitive patterns in random sequences, Zeitschrift f
 ür Wahrscheinlichkeitstheorie 53 (1980), 241–262.
- H.-K. Hwang, On convergence rates in the central limit theorems for combinatorial structures, European Journal of Combinatorics 19 (1998), 329–343.
- P. Kirschenhofer and H. Prodinger, On the analysis of probabilistic counting, Number-theoretic Analysis (E. Hlawka and R.F. Tichy, eds.), Lecture Notes in Mathematics, vol. 1452, 1990, pp. 117– 120.
- 10. _____, A result in order statistics related to probabilistic counting, Computing 51 (1993), 15–27.
- 11. _____, The path length of random skip lists, Acta Informatica **31** (1994), 775–792.
- P. Kirschenhofer, H. Prodinger, and W. Szpankowski, Analysis of a splitting process arising in probabilistic counting and other related algorithms, Random Structures and Algorithms 9 (1996), 379–401.
- 13. A. Knopfmacher and H. Prodinger, Combinatorics of geometrically distributed random variables: Value and position of the rth left-to-right maximum, submitted.
- 14. D. E. Knuth, *The average time for carry propagation*, Indagationes Mathematicae **40** (1978), 238–242.
- T. Papadakis, I. Munro, and P. Poblete, Average search and update costs in skip lists, BIT 32 (1992), 316–332.
- H. Prodinger, Combinatorics of geometrically distributed random variables: Left-to-right maxima, Discrete Mathematics 153 (1996), 253–270.
- 17. W. Pugh, *Skip lists: a probabilistic alternative to balanced trees*, Communications of the ACM **33** (1990), 668–676.

RUN STATISTICS

Peter Grabner, Institut für Mathematik A, Technische Universität Graz, Steyrergasse 30, 8010 Graz, Austria

E-mail address: grabner@weyl.math.tu-graz.ac.at

ARNOLD KNOPFMACHER, CENTRE FOR APPLICABLE ANALYSIS AND NUMBER THEORY, DEPART-MENT OF COMPUTATIONAL AND APPLIED MATHEMATICS, UNIVERSITY OF THE WITWATERSRAND, P. O. WITS, 2050 JOHANNESBURG, SOUTH AFRICA

E-mail address: arnoldk@gauss.cam.wits.ac.za

Helmut Prodinger, Centre for Applicable Analysis and Number Theory, Department of Mathematics, University of the Witwatersrand, P. O. Wits, 2050 Johannesburg, South Africa

E-mail address: helmut@gauss.cam.wits.ac.za