

How often do we reject a superior value?

— Extended abstract —

Kamilla Oliver^{1†} and Helmut Prodinger²

¹91052 Erlangen, Germany, olikamilla@gmail.com

²Department of Mathematics, University of Stellenbosch, 7602 Stellenbosch, South Africa, hprodinger@sun.ac.za

Abstract. Words $a_1a_2 \dots a_n$ with independent letters a_k taken from the set of natural numbers, and a weight (probability) attached via the geometric distribution pq^{i-1} ($p + q = 1$) are considered. A consecutive record (motivated by the analysis of a skip list structure) can only advance from k to $k + 1$, thus ignoring perhaps some larger (=superior) values. We investigate the number of these rejected superior values. Further, we study the probability that there is a single consecutive maximum and show that (apart from fluctuations) it tends to a constant.

Keywords: Combinatorics on words, records, generating functions, Rice’s method, q -series.

1 Introduction

We consider words $a_1a_2 \dots a_n$ with letters a_k taken from the set of natural numbers, and a weight (probability) attached to it by saying that the letter $i \in \mathbb{N}$ occurs with probability pq^{i-1} ($p + q = 1$) and that the letters are independent. The parameter $\mathcal{K}(a_1a_2 \dots a_n)$, which we call the number of *weak consecutive records*, has proved to be essential in the analysis of a skip list structure Louchard and Prodinger (2006).⁽ⁱ⁾ The word is scanned from left to right, and assuming that the current record (maximum) is value k , any letter different from k , $k + 1$ is ignored. If, however, the symbol scanned is one of these, we call it a weak consecutive record, and set the value of the current maximum to it. So, the current record either stays at k or advances to the next value $k + 1$. The skip list version assumes that the first letter of the word defines the first record.

For the sake of clarity, we consider the word 1311243535141234651 and underline each consecutive weak maximum: 1 3 1 1 2 4 3 5 3 5 1 4 2 1 3 4 6 5 1. The number of underlined symbols (9 in this case) is the parameter \mathcal{K} of interest.

In Louchard and Prodinger (2006), the average of the parameter $\mathcal{K}(n)$ was shown to be (with $Q = 1/q$, $L = \log Q$)

$$\mathbb{E}\mathcal{K}(n) = 1 + (Q + 1) \sum_{j=1}^n \binom{n}{j} \frac{(-1)^{j-1} (q; q)_{j-1} p^{j+1} q^j}{1 - q^{j+1}},$$

[†]K.O.’s work was partially done when she visited the Department of Mathematics of the University of Stellenbosch. She is thankful for the hospitality encountered there.

⁽ⁱ⁾ In order to learn more about this structure, the interested reader is invited to consult the earlier paper Louchard and Prodinger (2006), which also appeared in this journal; it also contains many pointers to some interesting earlier papers.

The consecutive records are printed in boldface; the ordinary records are marked by \vee , and the superior values that are neglected when scanning the word for consecutive records are marked by \bullet .

We consider this parameter in the next section, assuming random words of length n . It will turn out that both, expectation and variance, are of order $\log n$.

Some fifteen years ago, it was observed by several authors that the probability of a single winner (= a single maximum) in a word of length n does not tend to a limit, but rather oscillates around a certain value. Here are some papers about this Eisenberg et al. (1993); Baryshnikov et al. (1995); Bruss and O’Cinneide (1990); Pakes and Steutel (1997); Brands et al. (1994); Qi and Wilms (1997); Kirschenhofer and Prodinger (1996); Louchard and Prodinger (2006), and we apologize if we should have left out some relevant paper.

Now the maximum is also the left-to-right maximum (record). We investigate in the last section of this paper the analogous question related to consecutive records. Assume that k is the consecutive maximum, we consider the probability that after this consecutive record has been established, no further k ’s are read. Note, however, that earlier k ’s are possible, since they might have been ignored. **241152423121232411** has consecutive single maximum 4, but earlier some 4’s have been read (and rejected).

We encounter a similar phenomenon: there is a limiting value, but there are (tiny) oscillations around it. This constant comes out as a (quite complicated) series.

We use (standard) notation from q -analysis: $(x; q)_n = \prod_{i=0}^{n-1} (1 - xq^i)$ and $(x; q)_\infty = \prod_{i \geq 0} (1 - xq^i)$. Note that $(x; q)_n = (x; q)_\infty / (xq^n; q)_\infty$, and the latter form makes sense also for n a complex number. For several identities related to such quantities (Euler’s partition identities, Heine’s transformation formula, etc.), we refer to Andrews et al. (1999).

Furthermore, to say it again, we use $Q = 1/q$ and $L = \log Q$.

Our method is to set up (ordinary) generating functions, where z marks the length of the word and u the additional parameter. Moments are obtained by differentiation. It turns out that in this kind of problems the substitution $z = \frac{w}{w-1}$ makes everything much nicer. Eventually, one can translate everything back into the z -world:

$$\sum_{n \geq 1} a_n w^n = \sum_{n \geq 1} z^n \sum_{k=0}^{n-1} a_{k+1} (-1)^{k+1} \binom{n-1}{k}.$$

For the asymptotic evaluation, we use a contour integral representation of alternating sums (“Rice’s method”).

2 The difference between ordinary maxima and consecutive maxima

Since in the present setting the consecutive maximum can only advance by $+1$, we will miss (reject) some values that are larger than that. In this section, we count the number of rejected values.

Let $c_k(z, u)$ be the generating function where $[z^n u^j] c_k(z, u)$ is the probability that a random word of length n has consecutive maximum equal to k , and j better values have been rejected.

A recursion

$$c_k(z, u) = c_{k-1}(z, u) \frac{z p q^{k-1}}{1 - z[1 + (u-1)q^{k+1} - p q^k]} + \frac{z p q^{k-1}}{1 - z[1 + (u-1)q^{k+1} - p q^k]}.$$

It holds for $k \geq 1$, and we assume that $c_0 = 0$. It is good to use an abbreviation:

$$\lambda_k := 1 - z[1 + (u - 1)q^{k+1} - pq^k].$$

This follows from taking the instance with consecutive maximum $k-1$, and attaching the letter k , followed by an arbitrary sequence of letters different from $k+1$; the large ones are marked by u . The last terms reflects the situation that k is the first letter of the word.

Then

$$\frac{c_k(z, u)\lambda_1 \dots \lambda_k}{z^k p^k q^{\binom{k}{2}}} = \frac{c_{k-1}(z, u)\lambda_1 \dots \lambda_{k-1}}{z^{k-1} p^{k-1} q^{\binom{k-1}{2}}} + \frac{\lambda_1 \dots \lambda_{k-1}}{z^{k-1} p^{k-1} q^{\binom{k-1}{2}}}.$$

This first order recursion can now be solved by summation:

$$c_k(z, u) = \frac{z^k p^k q^{\binom{k}{2}}}{\lambda_1 \dots \lambda_k} \sum_{j=0}^{k-1} \frac{\lambda_1 \dots \lambda_j}{z^j p^j q^{\binom{j}{2}}} = \sum_{j=0}^{k-1} \frac{z^{k-j} p^{k-j} q^{\binom{k}{2} - \binom{j}{2}}}{\lambda_{j+1} \dots \lambda_k}.$$

Expectations

The generating function of the expectations is obtained by differentiation:

$$d_k(z) := \left. \frac{\partial}{\partial u} c_k(z, u) \right|_{u=1},$$

and is finally given by

$$\sum_{k \geq 1} d_k(z),$$

since we must take any final consecutive maximum into account. Now

$$d_k(z) = \sum_{j=0}^{k-1} \frac{z^{k-j} p^{k-j} q^{\binom{k}{2} - \binom{j}{2}}}{\lambda_{j+1} \dots \lambda_k} \Big|_{u=1} \sum_{i=j+1}^k \frac{z q^{i+1}}{1 - z(1 - pq^i)}$$

and

$$\lambda_{j+1} \dots \lambda_k \Big|_{u=1} = (1 - w)^{j-k} \prod_{i=j+1}^k (1 - wpq^i) = (1 - w)^{j-k} (wpq^{j+1}; q)_{k-j}.$$

Rewriting in w -notation:

$$d_k(z) = \sum_{0 \leq j < i \leq k} \frac{w^{k-j} (-1)^{k-j-1} p^{k-j} q^{\binom{k}{2} - \binom{j}{2}}}{(wpq^{j+1}; q)_{k-j}} \frac{w q^{i+1}}{1 - wpq^i}.$$

Summing up leads after a lengthy computation to:

$$\sum_{k \geq 1} d_k(z) = \frac{1}{p} \sum_{j \geq 0} \sum_{i \geq 1} q^{\binom{i}{2}} \frac{(-pq^{j+1}w)^{i+1}}{(wpq^{j+1}; q)_i}.$$

Note (this will be discussed later):

$$\sum_{h \geq 1} \frac{(-q^d w)^h q^{\binom{h}{2}}}{(w; q)_h} = - \sum_{n \geq 1} q^d \frac{(q; q)_{n+d-1}}{(q; q)_d} w^n. \quad (1)$$

Therefore

$$\sum_{k \geq 1} d_k(z) = -\frac{1}{p} \sum_{n \geq 1} (q; q)_{n-1} (wpq)^{n+1} \frac{1}{1 - q^{n+1}}.$$

Hence the coefficient of w^n in this expression is

$$\frac{p^{n-1} q^n}{1 - q^n} (q; q)_{n-2}$$

for $n \geq 2$ and 0 otherwise. And now we rewrite this in terms of coefficients z^n :

$$[z^n] \sum_{j \geq 1} d_j = (-1)^n [w^n] (1 - w)^{n-1} \sum_{j \geq 1} d_j = \sum_{k=1}^{n-1} \binom{n-1}{k} (-1)^{k+1} \frac{p^k q^{k+1}}{1 - q^{k+1}} (q; q)_{k-1}.$$

Second factorial moments

Let

$$f_k(z) := \left. \frac{\partial^2}{\partial u^2} c_k(z, u) \right|_{u=1}.$$

We have

$$\left. \frac{\partial^2}{\partial u^2} \lambda_{j+1} \cdots \lambda_k \right|_{u=1} = \frac{2}{(1-w)^{j-k} (wpq^{j+1}; q)_{k-j}} \sum_{j+1 \leq h \leq i \leq k} \frac{w^2 q^{h+i+2}}{(1-wpq^h)(1-wpq^i)}.$$

Combined with the rest,

$$f_k(z) = \sum_{0 \leq j < h \leq i \leq k} (-pw)^{k-j} q^{\binom{k}{2} - \binom{j}{2}} \frac{2}{(wpq^{j+1}; q)_{k-j}} \frac{w^2 q^{h+i+2}}{(1-wpq^h)(1-wpq^i)}.$$

And this must be summed (the long computation is not shown here):

$$\sum_{k \geq 1} f_k(z) = -\frac{2q}{p^2} \sum_{N \geq 3} \frac{(q; q)_{N-2} (wpq)^N}{1 - q^N} \sum_{n=1}^{N-2} \frac{1}{1 - q^n}.$$

The coefficient of w^N is

$$-\frac{2(q; q)_{N-2} p^{N-2} q^{N+1}}{1 - q^N} \sum_{n=1}^{N-2} \frac{1}{1 - q^n}$$

for $N \geq 1$, otherwise 0. Now we rewrite this:

$$[z^n] \sum_{k \geq 1} f_k(z) = \sum_{j=1}^{n-1} \binom{n-1}{j} (-1)^j \frac{2(q; q)_{j-1} p^{j-1} q^{j+2}}{1 - q^{j+1}} \sum_{m=1}^{j-1} \frac{1}{1 - q^m}.$$

Let us summarize what we found so far:

Theorem 3 *The expectation and the second factorial moment of the number of superior elements that we reject when we consider the consecutive maxima instead of the true maxima, assuming random words of length n , are given by*

$$\sum_{k=1}^{n-1} \binom{n-1}{k} (-1)^{k+1} \frac{p^k q^{k+1}}{1 - q^{k+1}} (q; q)_{k-1}$$

and

$$\sum_{j=1}^{n-1} \binom{n-1}{j} (-1)^j \frac{2(q; q)_{j-1} p^{j-1} q^{j+2}}{1 - q^{j+1}} \sum_{m=1}^{j-1} \frac{1}{1 - q^m}.$$

□

In various places of these computations, the formula (1) was used. It can be shown using Heine's transform Andrews (1976); a proof is included in the full paper.

Asymptotics

The main ingredient is *Rice's formula* Flajolet and Sedgewick (1995) which allows to write an alternating sum as a contour integral:

$$\sum_{k=1}^n \binom{n}{k} (-1)^k f(k) = -\frac{1}{2\pi i} \int_{\mathcal{C}} \frac{n! \Gamma(-z)}{\Gamma(n+1-z)} f(z) dz.$$

The positively oriented curve \mathcal{C} encircles the poles $1, 2, \dots, n$ and no others. (The lower summation index 1 could be replaced by another constant, say a .) Changing the contour, there are extra residues, which have to be taken into account with the opposite sign, and they constitute the asymptotic expansion of the sum. In our instance, this residue is at 0. There are also contributions from poles at $\frac{2\pi i k}{L}$, with $L = \log Q$, $Q = \frac{1}{q}$, which we do not compute explicitly. When one collects them, they constitute a (tiny) periodic oscillation.

The function $f(z)$ extrapolates the sequence $f(k)$. This is often obvious, but sometimes requires some work to provide such a function. For our expected value, this discussion leads to

$$\sum_{k=1}^{n-1} \binom{n-1}{k} (-1)^{k+1} \frac{p^k q^{k+1}}{1 - q^{k+1}} (q; q)_{k-1} = \frac{1}{2\pi i} \int_{\mathcal{C}} \frac{(n-1)! \Gamma(-z)}{\Gamma(n-z)} \frac{p^z q^{z+1}}{1 - q^{z+1}} (q; q)_{z-1} dz.$$

It is not *a priori* clear what $(q; q)_{z-1} = \frac{(q; q)_z}{1 - q^z}$ should be. However, note that

$$(q; q)_z = \frac{(q; q)_\infty}{(q; q^{z+1})_\infty}.$$

The expansion

$$(q; q)_z \sim 1 - \alpha Lz + \frac{\alpha^2 + \beta}{2} L^2 z^2$$

with

$$\alpha = \sum_{j \geq 1} \frac{q^j}{1 - q^j} \quad \text{and} \quad \beta = \sum_{j \geq 1} \frac{q^j}{(1 - q^j)^2}$$

is not hard to obtain. Providing this expansion, the rest can be done by a computer, and the expectation comes out as

$$\frac{q}{p} \log_Q n + \frac{q}{p} \log(p) - \frac{q}{p} \alpha + \frac{q\gamma}{pL} - \frac{q(1+q)}{2p^2}.$$

For the second factorial moment, things are more involved, and we must find the continuation of

$$\sum_{m=1}^{j-1} \frac{1}{1 - q^m} = j - 1 + \sum_{m=1}^{j-1} \frac{q^m}{1 - q^m}.$$

But

$$\sum_{m=1}^{j-1} \frac{q^m}{1 - q^m} = \sum_{m=1}^j \frac{q^m}{1 - q^m} - \frac{q^j}{1 - q^j} = \alpha - \sum_{m \geq 1} \frac{q^{m+j}}{1 - q^{m+j}} - \frac{q^j}{1 - q^j},$$

and thus we may take

$$z - 1 + \alpha - \sum_{m \geq 1} \frac{q^{m+z}}{1 - q^{m+z}} - \frac{q^z}{1 - q^z}.$$

Near $z = 0$, this function can be expanded as

$$z - 1 + \beta Lz - \frac{1}{Lz} + \frac{1}{2} - \frac{Lz}{12}.$$

Again, the rest can be done by a computer, and the second factorial moment can be obtained. It is not displayed here. From this, the variance is computed by adding the expectation and subtracting the square of the expectation. We summarize the results.

Theorem 4 *The expectation and the variance of the number of superior elements that we reject when we consider the consecutive maxima instead of the true maxima, assuming random words of length n , are given by*

$$\text{expectation} \sim \frac{q}{p} \log_Q n + \frac{q}{p} \log(p) - \frac{q}{p} \alpha + \frac{q\gamma}{pL} - \frac{q(1+q)}{2p^2} + \delta.(\log_Q n).$$

and

$$\begin{aligned} \text{variance} &\sim \frac{q \log_Q n}{p^2} \\ &+ \frac{q \log(p)}{p^2 L} - \frac{q\alpha}{p^2} + \frac{q\gamma}{p^2 L} - \frac{q^2 \beta}{p^2} + \frac{q^2 \pi^2}{6p^2 L^2} - \frac{2q^2}{p^2 L} \\ &+ \frac{q(q^3 - 6 + 16q^2 + q)}{12p^4} + \delta.(\log_Q n). \end{aligned}$$

Here, $\delta.(x)$ denotes an unspecified tiny periodic function of period 1. It might differ in different places. The leading term in the variance of order $\log n$ has no fluctuating component, since they cancel out. \square

3 Probability of a single (consecutive) winner

Exact enumeration

Recall that the generating function of words with consecutive maximum = k is given by

$$b_k(z) = c_k(z, 1) = \sum_{j=0}^{k-1} \frac{q^{\binom{k}{2} - \binom{j}{2}} (-pw)^{k-j}}{(wpq^{j+1}; q)_{k-j}}.$$

Then

$$e_k(z) = (1 + b_k(z)) \frac{zp^k}{1 - z(1 - pq^k - pq^{k+1})} = - \left(1 + \sum_{j=0}^{k-1} \frac{q^{\binom{k}{2} - \binom{j}{2}} (-pw)^{k-j}}{(wpq^{j+1}; q)_{k-j}} \right) \frac{wp^k}{1 - wp(1 + q)q^k}$$

is the generating function where the maximum is = $k + 1$, and it is a single consecutive maximum. This must be summed over $k + 1 \geq 1$, to get the generating function of words with a single consecutive maximum (details are in the full paper):

$$\begin{aligned} \sum_{k \geq 0} e_k(z) &= - \sum_{m \geq 0} (1 + q)^m (wp)^{m+1} \sum_{j \geq 0} q^{j(m+1)} \left(1 - \sum_{n \geq 1} q^m \frac{(q; q)_{n+m-1}}{(q; q)_m} (wpq^{j+1})^n \right) \\ &=: -A + B. \end{aligned}$$

Then

$$A = \sum_{m \geq 0} (1 + q)^m (wp)^{m+1} \frac{1}{1 - q^{m+1}}$$

and

$$B = \frac{1}{q} \sum_{m \geq 0} (1 + q)^m \sum_{n \geq 1} \frac{(q; q)_{n+m-1}}{(q; q)_m} (wpq)^{m+1+n} \frac{1}{1 - q^{m+1+n}}.$$

The coefficient of w^N in A is $\frac{p^N (1+q)^{N-1}}{1 - q^N}$ for $N \geq 1$ and 0 otherwise.

The coefficient of w^N in B is

$$p^N q^{N-1} \frac{(q; q)_{N-2}}{1 - q^N} \sum_{m=0}^{N-2} \frac{(1 + q)^m}{(q; q)_m}$$

for $N \geq 2$ and 0 otherwise.

Therefore we get an explicit expression.

Theorem 5 *The probability that a random word of length n has a single consecutive maximum is given (exactly) by*

$$\sum_{k=0}^n \binom{n}{k} (-1)^k \frac{p^{k+1} (1 + q)^k}{1 - q^{k+1}} + \sum_{k=1}^n \binom{n}{k} (-1)^{k+1} p^{k+1} q^k \frac{(q; q)_{k-1}}{1 - q^{k+1}} \sum_{m=0}^{k-1} \frac{(1 + q)^m}{(q; q)_m}.$$

□

Asymptotics

The first sum is $O(1/n)$, so we concentrate on the second one. The continuation of

$$\sum_{m=0}^{k-1} \frac{(1+q)^m}{(q; q)_m}$$

isn't as easy as before, since we cannot push it to infinity, and pulling out dominant terms isn't too obvious either. So we proceed in a different way (computation not shown in this abstract):

$$\sum_{m=0}^{n-1} \frac{(1+q)^m}{(q; q)_m} = \frac{1}{(q; q)_\infty} \sum_{j \geq 0} \frac{(-1)^j q^{\binom{j+1}{2}}}{(q; q)_j} \frac{1 - (1+q)^n q^{jn}}{1 - (1+q)q^j}.$$

Therefore we continue

$$p^{k+1} q^k \frac{(q; q)_{k-1}}{1 - q^{k+1}} \sum_{m=0}^{k-1} \frac{(1+q)^m}{(q; q)_m}$$

via

$$p^{z+1} q^z \frac{(q; q)_{z-1}}{1 - q^{z+1}} \frac{1}{(q; q)_\infty} \sum_{j \geq 0} \frac{(-1)^j q^{\binom{j+1}{2}}}{(q; q)_j} \frac{1 - (1+q)^z q^{jz}}{1 - (1+q)q^j}.$$

At $z = 0$, this function is regular, with the value

$$\frac{1}{(q; q)_\infty} \sum_{j \geq 0} \frac{(-1)^j q^{\binom{j+1}{2}}}{(q; q)_j} \frac{j - \log_Q(1+q)}{1 - (1+q)q^j}.$$

This follows from

$$\frac{1 - (1+q)^z q^{jz}}{1 - q^z} \sim \frac{1 - (1+z \log(1+q))(1 - Ljz)}{Lz} \sim j - \log_Q(1+q), \quad (z \rightarrow 0).$$

So, at $z = 0$, there is ultimately a simple pole, and, apart from the complicated value, the negative residue is one.

Theorem 6 *The probability that the consecutive record is single, is asymptotically, as $n \rightarrow \infty$, given by*

$$\frac{1}{(q; q)_\infty} \sum_{j \geq 0} \frac{(-1)^j q^{\binom{j+1}{2}}}{(q; q)_j} \frac{j - \log_Q(1+q)}{1 - (1+q)q^j} + \delta.(\log_Q n),$$

with a periodic function of period 1 and tiny oscillation $\delta.(x)$. □

One could also consider a random variable, namely the number of times the consecutive maximum occurs, and compute moments. However, we decided not to go that route here.

Acknowledgements

We thank Guy Louchard for valuable remarks.

References

- G. Andrews. *The Theory of Partitions*, volume 2 of *Encyclopedia of Mathematics and its Applications*. Addison–Wesley, 1976.
- G. Andrews, R. Askey, and R. Roy. *Special Functions*, volume 71 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, 1999.
- Y. Baryshnikov, B. Eisenberg, and G. Stengle. A necessary and sufficient condition for the existence of the limiting probability of a tie for first place. *Statist. Probab. Lett.*, 23:203–209, 1995.
- J. J. A. M. Brands, F. W. Steutel, and R. J. G. Wilms. On the number of maxima in a discrete sample. *Statist. Probab. Lett.*, 20:209–217, 1994.
- T. Bruss and C. A. O’Cinneide. On the maximum and its uniqueness for geometric random samples. *J. Appl. Probab.*, 27:598–610, 1990.
- B. Eisenberg, G. Stengle, and G. Strang. The asymptotic probability of a tie for first place. *Ann. Appl. Probab.*, 3:731–745, 1993.
- P. Flajolet and R. Sedgewick. Mellin transforms and asymptotics: Finite differences and Rice’s integrals. *Theoretical Computer Science*, 144:101–124, 1995.
- P. Kirschenhofer and H. Prodinger. The number of winners a in discrete geometrically distributed sample. *Annals in Applied Probability*, 6:687–694, 1996.
- G. Louchard and H. Prodinger. Analysis of a new skip list variant. *Discrete Mathematics and Theoretical Computer Science*, proc. AG:365–374, 2006.
- G. Louchard and H. Prodinger. The number of elements close to near-records in geometric samples. *Quaestiones Mathematicae*, 29:447–470, 2006.
- K. Oliver and H. Prodinger. Consecutive records in geometrically distributed words. *submitted*, 2009.
- A. Pakes and F. W. Steutel. On the number of records near the maximum. *Austral. J. Statist.*, 39:179–192, 1997.
- Y. Qi and R. J. G. Wilms. The limit behavior of maxima modulo one and the number of maxima. *Statist. Probab. Lett.*, 34:75–84, 1997.