

LONGEST COMMON SUBSEQUENCES OF TWO RANDOM SEQUENCES

by

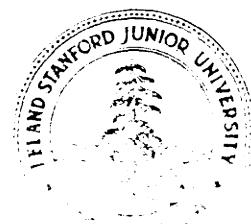
Vaclav **Chvatal**

David **Sankoff**

STAN-CS-75-477

JANUARY 1975

COMPUTER SCIENCE DEPARTMENT
School of Humanities and Sciences
STANFORD UNIVERSITY



Longest Common Subsequences of Two Random Sequences

by

Václav Chvátal
Computer Science Department
Stanford University

and

David Sankoff
Centre de recherches mathématiques
Université de Montréal

Abstract

Given two random k -ary sequences of length n , what is $f(n,k)$, the expected length of their longest common subsequence? This problem arises in the study of molecular evolution. We calculate $f(n,k)$ for all k , where $n \leq 5$, and $f(n,2)$ where $n < 10$. We study the limiting behavior of $n^{-1} f(n,k)$ and derive upper and lower bounds on these limits for all k . Finally we estimate by Monte-Carlo methods $f(100,k)$, $f(1000,2)$ and $f(5000,2)$.

American Mathematical Society 1970 subject classifications:

Primary 60C05 ; Secondary 92A10 .

Key words and phrases: common subsequences, matches, random sequences.

This research was supported in part by National Science Foundation grant GJ 36473X and by the Office of Naval Research contract NR 044-402. Reproduction in whole or in part is permitted for any purpose of the United States Government.

1. Introduction.

In the study of the evolution of long molecules such as proteins or nucleic acids, it is common practice to try to construct a large set of correspondences, or matches, between two such molecules. Mathematically, this is just the problem of finding a longest common subsequence of two given finite sequences. A quadratic algorithm for doing this is available (Sankoff (1972)). It is often difficult to judge whether this set of correspondences is significantly large, i.e. contains more correspondences than one would expect in the case of two random molecules of the same length and subunit composition. Tests of significance are unavailable, except on a Monte-Carlo basis (Sankoff and Cedergren (1973)), since nothing is known about the distribution of the length of the longest common subsequence. As a first step in the study of this distribution, this note investigates its mean value.

We introduce the following notation.

Let $\underline{a} = (a_1, a_2, \dots, a_n)$, $\underline{b} = (b_1, b_2, \dots, b_n)$ be two sequences. A common subsequence, or $(\underline{a}, \underline{b})$ -match is a set $M = \{(i_k, j_k) : 1 \leq k \leq m\}$ with $1 \leq i_1 < i_2 < \dots < i_m \leq n$, $1 \leq j_1 < j_2 < \dots < j_m \leq n$ and $a_{i_k} = b_{j_k}$ for each $(i_k, j_k) \in M$. The size of a largest $(\underline{a}, \underline{b})$ -match will be denoted by $v(\underline{a}, \underline{b})$. By a k -ary sequence we mean one whose terms come from $\{1, 2, \dots, k\}$. We shall study the function $f(n, k)$ defined as the mean value of $v(\underline{a}, \underline{b})$ over all the k^{2n} ordered pairs $(\underline{a}, \underline{b})$ of k -ary sequences of length n .

2. Exact formulae for $f(n, k)$ with small n .

Let $\underline{a} = (a_1, a_2, \dots, a_n)$ and $\underline{b} = (b_1, b_2, \dots, b_n)$ be two k -ary sequences. The pair $(\underline{a}, \underline{b})$ will be called *normal* if, setting $a_{n+j} = b_j$ for all j , we have $a_1 = 1$ and

$$a_3 \leq \max(a_1, a_2, \dots, a_{j-1}) + 1 \quad (2 \leq j \leq 2n).$$

Let $N(n, v, t)$ denote the number of normal pairs $(\underline{a}, \underline{b})$ with $v(\underline{a}, \underline{b}) = v$ and $\max\{a_1, a_2, \dots, a_{2n}\} = t$. Clearly, the number of pairs $(\underline{c}, \underline{d})$ where $\underline{c}, \underline{d}$ are k -ary sequences of length n with $v(\underline{c}, \underline{d}) = v$ is equal to

$$\sum_{t=1}^{2n} N(n, v, t) \cdot (k)_t$$

where $(k)_t$ is the falling factorial $k(k-1)\dots(k-t+1)$. Hence

$$\begin{aligned} f(n, k) &= \frac{1}{k^{2n}} \sum_{v=0}^n v \sum_{t=1}^{2n} N(n, v, t) \cdot (k)_t \\ &= \frac{1}{k^{2n}} \sum_{v=0}^n v \sum_{t=1}^{2n} N(n, v, t) \sum_{j=1}^t s(t, j) k^{j-2n} \\ &= \sum_{j=1}^{2n} \sum_{t=j}^{2n} s(t, j) \sum_{v=0}^n v N(n, v, t) k^{j-2n} \end{aligned}$$

where $s(t, j)$ are the Stirling numbers of the first kind (Riordan (1958)).

Note that $N(n, v, 2n) = 0$ unless $v=0$ and so

$$f(n, k) = \sum_{j=1}^{2n-1} \sum_{t=1}^{2n-1} s(t, j) \sum_{v=0}^n v N(n, v, t) k^{j-2n}.$$

Also

$$N(n, v, 2n-1) = \begin{cases} n^2 & \text{if } v=1 \\ 0 & \text{if } v>1 \end{cases}$$

and so the coefficient of $f(n, k)$ at k^{-1} is

$$s(2n-1, 2n-1) = \sum_{v=0}^n v N(n, v, 2n-1) = n^2.$$

We have evaluated $N(n, v, t)$ for $1 \leq n \leq 5$ and arrived at the following formulae.

$$f(1, k) = k^{-1},$$

$$f(2, k) = 4k^{-1} - 5k^{-2} + 3k^{-3},$$

$$f(3, k) = 9k^{-1} - 27k^{-2} + 60k^{-3} - 71k^{-4} + 32k^{-5},$$

$$f(4, k) = 16k^{-1} - 84k^{-2} + 380k^{-3} - 1146k^{-4} + 2085k^{-5} - 2018k^{-6} + 771k^{-7},$$

$$f(5, k) = 25k^{-1} - 200k^{-2} + 1500k^{-3} - 8200k^{-4} + 30640k^{-5} - 75096k^{-6} + 113748k^{-7} - 94790k^{-8} + 32378k^{-9}.$$

The values of these functions for $1 \leq k \leq 15$ are given in the table below.

	$f(1,k)$	$f(2,k)$	$f(3,k)$	$f(4,k)$	$f(5,k)$
k= 1	1 .000000	2.000000	3.000000	4 .000000	5.000000
2	.500000	1.125000	1.812500	2.523438	3.246094
3	.333333	.888889	1.477366	2.090535	2.718742
4	.250000	.734375	1.253906	1.801453	2.363899
5	.200000	.624000	1.096640	1.594317	2.108546
6	.166667	.541667	.977109	1.435968	1.912269
7	.142857	.478134	.881954	1.309838	1.754954
8	.125000	.427734	.803955	1.206201	1.625155
9	.111111	.386831	0.888889	1.119008	1.515694
10	.100000	.353000	.683220	1.044309	1.421763
11	.090909	.324568	.635470	.979404	1.340005
12	.083333	.300347	.593927	.922366	1.267999
13	.076923	.279472	.557455	.871776	1.203953
14	.071429	.261297	.525179	.826554	1.146514
15	.066667	.245333	.496417	.785862	1.094633

TABLE 1

Moreover, we have evaluated $f(n,2)$ for all $n = 1, 2, \dots, 10$; the results are given in Table 2 in proportion to n ,

n	$f(n, 2)/n$
1	0.500000
2	0.562500
3	0.604167
4	0.630859
5	0.649219
6	0.663330
7	0.674491
8	0.683640
9	0.691303
10	0.697844

TABLE 2

3. Limiting behaviour of $f(n, k)$.

Klarner and Rivest (personal communication) have observed that $f(n, k)$ is superadditive with respect to n , that is, $f(n_1 + n_2, k) \geq f(n_1, k) + f(n_2, k)$. Thus, by Fekete's theorem (Fekete (1923)),

$$\lim_{n \rightarrow \infty} n^{-1} f(n, k) = \sup_n n^{-1} f(n, k). \quad (1)$$

We shall denote the common value of (1) by c_k . Klarner and Rivest asked whether $c_2 = 1$; we shall show that this is not the case.

A sequence (s_1, s_2, \dots, s_m) is said to be a subsequence of a sequence (a_1, a_2, \dots, a_n) if there is a mapping $\varphi: \{1, 2, \dots, m\} \rightarrow \{1, 2, \dots, n\}$ such that

$$i < j \Rightarrow \varphi(i) < \varphi(j)$$

and such that

$$a_{\varphi(i)} = b_i \quad \text{for all } i = 1, 2, \dots, m.$$

LEMMA 1. Let s be a k -ary sequence of length m , let n be an integer with $n \geq m$ and let $F(n, \underline{s}, k)$ denote the number of k -ary sequences of length n containing \underline{s} as a subsequence. Then

$$F(n, \underline{s}, k) = \sum_{j=m}^n \binom{n}{j} (k-1)^{n-j}.$$

Proof. The formula holds trivially if $m = 1$ or $m = n$. To prove that it holds for all choices of $\underline{s} = (s_1, s_2, \dots, s_m)$, k and n , we shall proceed by induction on $m+n$. Let \hat{s} denote the sequence $(s_1, s_2, \dots, s_{m-1})$; for every sequence $\underline{a} = (a_1, a_2, \dots, a_n)$, let $\hat{\underline{a}}$ denote the sequence $(a_1, a_2, \dots, a_{n-1})$. Let A^+ , resp. A' , denote the set of all the k -ary sequences (a_1, a_2, \dots, a_n) containing \underline{s} as a subsequence and such that $a_n = s_m$, resp. $a_n \neq s_m$. Clearly, $a \in A^+$ if and only if $\hat{\underline{a}}$ contains $\hat{\underline{s}}$ and $a_n = s_m$; similarly, $a \in A'$ if and only if $\hat{\underline{a}}$ contains $\hat{\underline{s}}$ and $a_n \neq s_m$. Hence

$$F(n, \underline{s}, k) = |A^+| + |A'| = F(n-1, \hat{\underline{s}}, k) + (k-1)F(n-1, \hat{\underline{s}}, k)$$

The rest follows by the induction hypothesis.

Note that

$$\binom{n}{j} (k-1)^{n-j} \geq \binom{n}{j+1} (k-1)^{n-j-1}$$

whenever $j \geq n/k$. Hence

$$F(n, s, k) \leq n \sum_0^n (k-1)^{n-m} \quad \text{for } m \geq n/k. \quad (2)$$

For every real x with $1/k < x < 1$, we shall set

$$h_k(x) = \frac{k^{x/2-1} k^{\frac{1}{x}}}{x^{\frac{1}{x}} (1-x)^{1-x}}.$$

LEMMA 2. Let $g(n, m, k)$ denote the number of ordered pairs (\tilde{a}, \tilde{b}) of k -ary sequences of length n with $v(\tilde{a}, \tilde{b}) \geq m$. If x is a real number with

$$1/k < x < 1, \quad h_k(x) < 1$$

then

$$g(n, [xn], k) = o(k^{2n}) \quad (n \rightarrow \infty).$$

Proof. Let $G(n, m, k)$ denote the number of ordered triples $(\tilde{a}, \tilde{b}, \tilde{s})$ such that \tilde{a}, \tilde{b} are k -ary sequences of length n , \tilde{s} is a k -ary sequence of length m and \tilde{s} is a subsequence of both \tilde{a} and \tilde{b} .

Clearly,

$$g(n, m, k) \leq G(n, m, k) \quad (3)$$

and

$$G(n, m, k) = \sum (F(n, s, k))^2$$

with the summation extending over all the k -ary sequences \tilde{s} of length m .

By (2), we now have

$$G(n, m, k) \leq k^m \left(n \binom{n}{m} (k-1)^{n-m} \right)^2 \quad (4)$$

whenever $m \geq n/k$. Let $m = [xn]$. By Stirling's formula, we have

$$\begin{aligned} & \lim_{n \rightarrow \infty} \left(k^m \left(n \binom{n}{m} (k-1)^{n-m} \right)^2 k^{-2n} \right)^{1/n} \\ &= \lim_{n \rightarrow \infty} \left(k^{xn} \left(\binom{n}{xn} (k-1)^{n-xn} \right)^2 k^{-2n} \right)^{1/n} \\ &= (h_k(x))^2 < 1 \end{aligned}$$

and so

$$k^m \left(n \binom{n}{m} (k-1)^{n-m} \right)^2 = o(k^{2n}) \quad (n \rightarrow \infty)$$

The rest follows by (3) and (4).

Note that

$$h_k(1/k) = k^{1/2k} > 1, \quad \lim_{x \rightarrow 1} h_k(x) = k^{-1/2} < 1$$

and

$$\frac{d}{dx} h_k(x) = h_k(x) \log \left(\frac{(1-x)k^{1/2}}{x(k-1)} \right),$$

so that h_k first increases and then decreases in the interval $[1/k, 1]$.

Hence there is a unique solution of

$$h_k(x) = 1, \quad 1/k < x < 1;$$

we shall denote this solution by y_k . Values of y_k with $2 \leq k \leq 15$ are shown in the following table, to six-decimal accuracy.

k	y_k
2	0.866595
3	0.786473
4	0.729705
5	0.686117
6	0.650984
7	0.621719
8	0.596756
9	0.575075
10	0.555971
11	0.538945
12	0.523625
13	0.511667
14	0.497038
15	0.485378

Table 3

THEOREM 1. If $k \geq 2$ then $c_k \leq y_k$.

Proof. For every positive ϵ with $y_k + \epsilon < 1$, we have $h_k(y_k + \epsilon) < 1$. Lemma 2 implies that

$$g(n, \lceil (y_k + \epsilon)n \rceil, k) = o(k^{2n})$$

and so

$$f(n, k) = k^{-2n} \sum \nu(a, b) \leq (1 - o(1)) \lceil (y_k + \epsilon)n \rceil + o(1)n$$

Hence

$$c_k = \lim_{n \rightarrow \infty} f(n, k)/n \leq y_k + \epsilon$$

and the desired conclusion follows.

4. Lower bounds on c_k .

For each pair $(\underline{a}, \underline{b})$ of k -ary sequences of length n , we shall construct a certain $(\underline{a}, \underline{b})$ -match M of size $v'(\underline{a}, \underline{b})$ and show that $f'(n, k)$, the average of $v'(\underline{a}, \underline{b})$ over all k^{2n} ordered pairs $(\underline{a}, \underline{b})$, satisfies

$$\lim_{n \rightarrow \infty} n^{-1} f'(n, k) = \frac{2k^2}{k^3 + 2k - 1}. \quad (2)$$

The construction of M is described below. The main idea is to begin by looking for the "first" matching pair (a_3, b_j) where $i = 1$ or $j = 1$. For example, suppose we examine the pairs $(a_1, b_1), (a_1, b_2), (a_2, b_1), (a_1, b_3)$ and finally find the first matching pair, namely (a_3, b_1) . Then we include (a_3, b_1) in M and proceed to look for the "first" matching pair in the sequences a_4, a_5, \dots, a_n and b_3, b_4, \dots, b_n . We continue until one or both sequences are exhausted.

STEP 0. Let $a_i = \underline{a}_i$, $b_i = \underline{b}_i$ and $S(i) = T(i) = i$ for all $i = 1, 2, \dots, n$. Let FLAG = 1 and $M = \emptyset$.

STEP 1. If FLAG = 1, check successively

$$(\alpha_1, \beta_1), (\alpha_1, \beta_2), (\alpha_2, \beta_1), \dots, (\alpha_1, \beta_d), (\alpha_d, \beta_1), \dots$$

until α or β is exhausted or until we find a pair with $\alpha_i = \beta_j$. If FLAG = -1, check the pairs in the order

$$(\alpha_1, \beta_1), (\alpha_2, \beta_1), (\alpha_1, \beta_2), \dots, (\alpha_d, \beta_1), (\alpha_1, \beta_d), \dots.$$

In the case of exhaustion, stop; otherwise add the pair $(S(i), T(j))$ to M .

STEP 2. Note that $i = 1$ or $j = 1$ or both

If $i \leq 2$ and $j \leq 2$, set

$$i' = i+1, j' = j+1.$$

If $i = 1$ and $j \geq 3$, set

$$i' = \begin{cases} j-1 & (\text{FLAG} = 1) \\ j & (\text{FLAG} = -1) \end{cases}, \quad j' = j+1.$$

If $i \geq 3$ and $j = 1$, set

$$i' = i+1, \quad j' = \begin{cases} i & (\text{FLAG} = 1) \\ i-1 & (\text{FLAG} = -1) \end{cases}.$$

STEP 3. Let $p = S(i')-1, q = T(j')-1$ and redefine

$$S(i) = p+i, \quad a_i = a_{S(i)}$$

$$T(j) = q+j, \quad b_j = b_{T(j)}$$

for all i, j with $1 \leq i \leq n-p, 1 \leq j \leq n-q$.

Reverse the sign of FLAG and go to Step 1.

LEMMA 2. For infinite sequences \mathbf{a}^* and \mathbf{b}^* , we have ,

$$E(i'+j'-2) = \frac{k^3+2k-1}{k^2}$$

where i', j' are defined as in the preceding algorithm and $E(\cdot)$ denotes mathematical expectation.

Proof. Consider the sequence of pairs in case $\text{FLAG} = 1$, that is,

$$(\alpha_1, \beta_1), (\alpha_1, \beta_2), (\alpha_2, \beta_1), \dots, (\alpha_1, \beta_d), (\alpha_d, \beta_1), \dots .$$

The event that any of these pairs contains equal terms has probability $1/k$ and this is also the conditional probability given any or all the preceding pairs. Hence the probability that the r -th pair will be the first equal one is $(k-1)^{r-1}/k^r$. Now,

$$i' + j' - 2 = \begin{cases} 2 & \text{if } r = 1, \\ 3 & \text{if } r = 2, \\ r & \text{if } r \geq 3. \end{cases}$$

Therefore

$$E(i' + j' - 2) = 2 \cdot \frac{1}{k} + 3 \cdot \frac{1}{k^2} + \sum_{r=3}^{\infty} r \cdot \frac{(k-1)^r}{k^r} = \frac{k^3 + 2k - 1}{k^2} .$$

The same can be shown for case $\text{FLAG} = -1$.

THEOREM 2. For all k , we have $c_k \geq \frac{2k^2}{k^3 + 2k - 1}$.

Proof. Obviously, it will suffice to prove (2). Let x_1, x_2, \dots be successive values of $i' + j' - 2$ found by the algorithm when applied to the infinite sequences \mathbf{a}^* and \mathbf{b}^* . It is clear that the x_i 's are independent, identically distributed random variables (indeed, in each cycle, equality or inequality of pairs is independent of all previous cycles). Let

$$x_k = \frac{2k^2}{k^3 + 2k - 1} .$$

The symmetry ensured by the alternation of sign of FLAG ensures that after $w = 2u$ cycles of the algorithm, the total number p (resp. q) of

the a^*_i 's (resp. b^*_j 's) that have been used up satisfies

$$E(p) = E(q) = \frac{w}{2} E(i'+j'-2) = w/x_k.$$

Furthermore,

$$\Pr \left(\left| \frac{p}{w} - \frac{1}{x_k} \right| > \epsilon \right) = \Pr \left(\left| \frac{q}{w} - \frac{1}{x_k} \right| > \epsilon \right) = o(w)$$

by the law of large numbers. Now a pair $(\underline{a}, \underline{b})$ of random sequences of length n can be considered as being the first n terms of \underline{a}^* and \underline{b}^* . If the algorithm (applied to $\underline{a}, \underline{b}$) halts during the $(w+1)$ -st cycle then the first w cycles are the same as the first w cycles of the algorithm applied to \underline{a}^* and \underline{b}^* . Now, after $\lfloor nx_k \rfloor$ cycles of the algorithm applied to $\underline{a}^*, \underline{b}^*$, we have

$$\begin{aligned} \Pr(p > n(1+\epsilon) \text{ or } p < n(1-\epsilon)) &= \Pr \left(\left| \frac{p}{\lfloor nx_k \rfloor} - \frac{1}{x_k} \right| > \frac{\epsilon}{x_k} \right) \\ &= o(n) \end{aligned}$$

and so

$$\Pr(n(1-\epsilon) \leq p \leq n(1+\epsilon) \text{ and } n(1-\epsilon) \leq q \leq n(1+\epsilon)) = 1-o(n).$$

Hence with probability $1-o(n)$, at least $\lfloor nx_k \rfloor - n\epsilon$ and at most $\lfloor nx_k \rfloor + n\epsilon$ cycles of the algorithm (applied to $\underline{a}^*, \underline{b}^*$) operate within \underline{a} and \underline{b} since $n\epsilon$ successive terms in a sequence can give rise to at most $n\epsilon$ cycles of the algorithm. Equivalently,

$$\Pr(|v'(\underline{a}, \underline{b}) - \lfloor nx_k \rfloor| \leq n\epsilon) = 1-o(n)$$

and so $\lim_{n \rightarrow \infty} n^{-1} f'(n, k) = x_k$.

Values of x_k with $2 \leq k \leq 15$ are given in the table below.

k	x_k
2	0.727273
3	0.562500
4	0.450704
5	0.373134
6	0.317181
7	0.275281
8	0.242884
9	0.217158
10	0.196271
11	0.178994
12	0.164477
13	0.152115
14	0.141465
15	0.132197

TABLE 4

5. Monte-Carlo estimates for $f(100, k)$ and σ_2 .

To obtain further information about c_k , we carried out two series of Monte-Carlo simulations. First, for $n = 100$ and for each $k = 2, \dots, 15$, we generated 100 pairs (a, b) of random k -ary sequences and calculated $v(a, b)$ in each case. We denote by $m_{k, n}$ the average value of $n v(a, b)$ in a given sample. For large n , this quantity may be considered an estimate of c_k . Values of $m_{k, 100}$ are tabulated in Table 5, and may be compared with the upper and lower bounds in Tables 2 and 4. Table 5 also contains $s_{k, 100}$, where

$$s_{k,n}^2 = \sum_{\substack{(a,b) \\ \text{in} \\ \text{sample}}} (n^{-1}v(a,b) - m_{k,n})^2 / (\text{sample size} - 1)$$

is an unbiased estimator of the variance of $n^{-1}v(a,b)$.

k	$m_{k,100}$	$s_{k,100}$
2	0.7814	0.0243
3	0.6855	0.0210
4	0.6242	0.0176
5	0.5778	0.0211
6	0.5332	0.0208
7	0.5065	0.0214
8	0.4812	0.0219
9	0.4593	0.0211
10	0.4423	0.0208
11	0.4268	0.0200
12	0.4126	0.0193
13	0.4003	0.0212
14	0.3827	0.0212
15	0.3712	0.0198

TABLE 5

To estimate c_2 more closely, a second series of simulations were carried out for $k=2$ and $n = 10, 100, 1000$, and 5000 . Table 6 lists $m_{2,n}$ and $s_{2,n}$, as well as the size of the sample used to make these estimates.

n	$m_{k,n}$	$s_{k,n}$	sample size
10	0.6991	0.1079	1000
100	0.7806	0.0238	100
1000	0.80529	0.00468	100
5000	0.8082	0.0015	6

TABLE 6

On the basis of these simulations, it seems fair to conjecture that $c_2 > 4/5$ and that the variance of $v(\mathbf{a}, \mathbf{b})$ is $o(n^{2/3})$.

REFERENCES

- [1] Fekete, M. (1923). Über die Verteilung der Wurzeln bei gewissen algebraischen Gleichungen mit ganzzahligen Koeffizienten. *Math. Zeit.* 17 228-249.
- [2] Riordan, J. (1958). *An Introduction to Combinatorial Analysis*. Wiley, New York.
- [3] Sankoff, D. (1972). Matching sequences under deletion/insertion constraints. *Proc. Nat. Acad. Sci. U.S.A.* 69 4-6.
- [4] Sankoff, D. and Cedergren, R. J. (1973). A test for nucleotide sequence homology. *J. Mol. Biol.* 77 159-164.