

LONG COMMON SUBSEQUENCES AND THE PROXIMITY OF TWO RANDOM STRINGS*

J. MICHAEL STEELE†

Abstract. Let (x_1, x_2, \dots, x_n) and $(x'_1, x'_2, \dots, x'_n)$ be two strings from an alphabet \mathcal{A} , and let L_n denote their longest common subsequence. The probabilistic behavior of L_n is studied under various probability models for the x 's and x' 's.

1. Introduction. Long molecules such as proteins and nucleic acids can be thought of schematically as sequences from a finite alphabet \mathcal{A} . From an evolutionary point of view it is natural to compare molecules by considering their common ancestors, and in schematic terms this reduces to the problem of considering the longest common subsequence of two given sequences.

Sankoff (1972) gave an efficient algorithm for calculating the length of the longest common subsequence. Subsequently, Sankoff and Cedergren (1973), and Sankoff, Cedergren, and Lapalme (1976) considered a number of empirical cases and conducted some Monte Carlo investigations. The first formal probabilistic analysis of the problem of long common subsequences was initiated in Chvátal and Sankoff (1975). To describe their work we first introduce some notation.

By X_i and X'_i , $1 \leq i < \infty$, we denote two sequences of independent, and identically distributed random variables with values in \mathcal{A} . The random variable of main interest is

$$L_n := \max \{k : X_{i_1} = X'_{j_1}, X_{i_2} = X'_{j_2}, \dots, X_{i_k} = X'_{j_k} \text{ where} \\ 1 \leq i_1 < i_2 < \dots < i_k \leq n \text{ and } 1 \leq j_1 < j_2 < \dots < j_k \leq n\}.$$

In words, L_n is the largest cardinality of any subsequence common to the sequences (X_1, X_2, \dots, X_n) and $(X'_1, X'_2, \dots, X'_n)$.

Under the assumption that $|\mathcal{A}| = k$ and that X_i and X'_i are both uniform on \mathcal{A} , Chvátal and Sankoff proved the existence of the limit of the means,

$$(1.1) \quad \lim_{n \rightarrow \infty} E \frac{L_n}{n} = c_k.$$

Among other results, Chvátal and Sankoff obtained upper and lower bounds on c_k . These authors proved no results for $\text{Var } L_n$, but on the basis of a Monte Carlo study they were led to conjecture $\text{Var } L_n = o(n^{2/3})$.

Deken (1979) was able to sharpen the bounds on c_k , and also noted that as a consequence of Kingman's subadditive ergodic theorem (Kingman (1968)), that one actually has

$$(1.2) \quad \lim_{n \rightarrow \infty} \frac{L_n}{n} = c \quad \text{a.s.}$$

where c depends on the distributions of the processes $\{(X_i, Y_i) : 1 \leq i \leq \infty\}$.

This result naturally entails $\text{Var } L_n = o(n^2)$, but no further progress was made on the variance problem.

* Received by the editors November 28, 1980. This work was supported in part by the Office of Naval Research under contract N00014-76-C-0475.

† Department of Statistics, Stanford University, Stanford, California 94305.

The present article takes up several aspects of the study of L_n . In the second section, as an elementary application of an inequality of Efron and Stein (1981), it is proved that $\text{Var } L_n = O(n)$. This makes only modest progress on the Chvátal–Sankoff conjecture that $\text{Var } L_n = o(n^{2/3})$, but it still serves to supplement (1.2) with a rate of convergence result.

The third section takes up the question of the behavior of L_n under more general assumptions than independence. A simple complement of Kingman's subadditive ergodic theorem (Kingman (1973)) is derived and then applied to L_n . The coupling method which is used here (or the Radon–Nikodým method which is sketched) may likely be of use in many other problems where subadditivity is available but stationarity is absent.

The fourth section branches out from the explicit analysis of L_n . It addresses the question of whether there exist statistics which are more tractable than L_n , but which still reasonably measure the genetic proximity of long molecules. The principal new candidate is T_n , the total number of common subsequences. Here one can compute ET_n exactly, but we note T_n has other drawbacks to its analysis.

The final section makes brief comment on some open problems and related literature.

2. A variance bound. Let $S(v_1, v_2, \dots, v_{n-1})$ denote any real-valued function of $n-1$ vectors $v_i \in \mathbb{R}^d$; and suppose V_i , $1 \leq i < \infty$, is any sequence of independent, identically distributed, random vectors in \mathbb{R}^d . We then define new random variables $S_i = S(V_1, V_2, \dots, V_{i-1}, V_{i+1}, \dots, V_n)$ for $1 \leq i \leq n$, and we further set $S = (1/n) \sum_{i=1}^n S_i$. Tukey's jackknife estimate for the variance of S is $\sum_{i=1}^n (S_i - S)^2$, and Efron and Stein (1981) have proved the very useful inequality,

$$(2.1) \quad \text{Var}(S) \leq E \sum_{i=1}^n (S_i - S)^2.$$

The main point of this section is to show that (2.1) leads to the bound

$$(2.2) \quad \text{Var } L_n = O(n),$$

under the general assumption that $V_i = (X_i, X_i')$ are independent and identically distributed. In fact, one can prove the following result.

THEOREM 1. *For each n , suppose there is defined a function $S(x_1, x_2, \dots, x_n)$ from $(\mathbb{R}^d)^n$ to \mathbb{R} . Suppose also that V_i , $1 \leq i < \infty$, is any sequence of independent random vectors in \mathbb{R}^d , and for $1 \leq i \leq n$, $1 \leq n < \infty$ set*

$$(2.3) \quad S_{i,n} = S(V_1, V_2, \dots, V_{i-1}, V_{i+1}, \dots, V_n).$$

If $E(S_{i,n} - S_{j,n})^2$ is bounded for all $1 \leq i < j \leq n$ and $1 \leq n < \infty$, then

$$(2.4) \quad \text{Var } S(V_1, V_2, \dots, V_n) = O(n).$$

Proof. Let the bound on $E(S_{i,n} - S_{j,n})^2$ be B . Fix n , define $S = (1/n) \sum_{i=1}^n S_{i,n}$, and let

$$(2.5) \quad \begin{aligned} D_n &= S(V_1, V_2, \dots, V_{n-1}) - S \\ &= \frac{1}{n} \sum_{i=1}^n (S(V_1, V_2, \dots, V_{n-1}) - S(V_1, V_2, \dots, V_{i-1}, V_{i+1}, \dots, V_n)). \end{aligned}$$

By Schwarz' inequality,

$$(2.6) \quad \text{Var } S(V_1, V_2, \dots, V_n) \leq \text{Var}(S) + \text{Var } D_n + 2(\text{Var } S)^{1/2}(\text{Var } D_n)^{1/2},$$

and

$$(2.7) \quad \text{Var } D_n \leq ED_n^2 \leq B.$$

Since $E(S_{i,n} - S_{i,n})^2 \leq B$ one also has $E((S_{i,n} - S.)^2) \leq B$. So inequalities (2.1), (2.6), and (2.7) entail

$$(2.8) \quad \text{Var } S(V_1, V_2, \dots, V_{n-1}) \leq nB + B + 2(nB)^{1/2}B^{1/2} = B(n^{1/2} + 1)^2.$$

This completes the proof of the theorem with a very specific form of the $O(n)$ term. \square

Returning to L_n we note that, for $V_i = (X_i, X'_i)$ and $\mathcal{A} = \{1, 2, \dots\}$, Theorem 1 is applicable to $S(V_1, V_2, \dots, V_n) = L_n(V_1, V_2, \dots, V_n) \equiv L_n$. Since

$$(2.9) \quad 0 \leq L(V_1, V_2, \dots, V_n) - L(V_1, V_2, \dots, V_{i-1}, V_{i+1}, \dots, V_n) \leq 1,$$

it is trivial that (2.8) can be taken with $B = 1$. In summary we have the following bound.

COROLLARY 1. *If (X_i, X'_i) are i.i.d. with values in $\mathcal{A} \times \mathcal{A}$ then*

$$(2.10) \quad \text{Var } L_{n-1} \leq (n^{1/2} + 1)^2.$$

By the usual Borel–Cantelli and subsequence arguments together with (2.9) and (2.10) one can prove a rate result.

COROLLARY 2. *We have for all $\epsilon > 0$ that*

$$(2.11) \quad L_n - EL_n = o(n^{3/4+\epsilon}) \quad \text{with probability one.}$$

Since the techniques for proving (2.11) are well known and since the result is not the best possible, there is no reason to include the proof. This is nevertheless the first rate result available on L_n , since such rates cannot be obtained in general from the subadditive ergodic theorem (cf. Hammersley (1978, p. 670)).

3. Nonstationary sequences. By Deken’s observation we know Kingman’s theorem implies that L_n/n converges almost surely under the assumption that the V_i , $1 \leq i < \infty$ form a stationary sequence. The point of this section is to give a very simple illustration of how Kingman’s theorem can also be used for nonstationary processes. Naturally, one must appeal to some underlying asymptotic stationarity, but the resulting class of results seem useful enough to merit recording. In particular, one should compare the present result to the “substationary” subadditive ergodic theorem of Abid (1979). That result apparently does not suffice for the application to L_n given here, it is considerably more complicated.

By a *subadditive sequence of functions* on E we denote a sequence $h_n : E^n \rightarrow \mathbb{R}$ which satisfies

$$(3.1) \quad h_{m+n}(e_1, e_2, \dots, e_{n+m}) \leq h_m(e_1, e_2, \dots, e_m) + h_n(e_{m+1}, e_{n+2}, \dots, e_{n+m}).$$

As an example, we note that if $E = \mathcal{A} \times \mathcal{A}$ and $e_i = (a_i, a'_i)$, then letting $h_n(e_1, e_2, \dots, e_n)$ denote the length of the longest common subsequence of (a_1, a_2, \dots, a_n) and $(a'_1, a'_2, \dots, a'_n)$ one has (3.1). Because of the applications we have in view, we will also focus on *monotone subadditive functions*, i.e., those functions which satisfy (3.1) as well as

$$(3.2) \quad h_{n-m}(x_{m+1}, x_{m+2}, \dots, x_n) \leq h_n(x_1, x_2, \dots, x_n) \quad \text{for all } m \leq n \text{ and } \{x_1, x_2, \dots, x_n\}.$$

We will say that a stochastic process $\{X_i\}_{i=1}^\infty$ on the discrete state space E has a *stationary ergodic coupling* if there is a stationary ergodic process $\{X'_i\}_{i=1}^\infty$ on the same

probability space such that $Z_i = (X_i, X'_i)$ is a coupling, i.e., such that the stopping time $\tau = \min \{i: X_i = X'_i\}$ is finite with probability one, and $X_i = X'_i$ for all $i \geq \tau$.

It is well known that couplings are a convenient and powerful way of expressing the asymptotic properties of stochastic processes (see, e.g., Griffeath (1978)). The next result illustrates this ease of application.

THEOREM 2. *Suppose that h_n is a positive and monotone sequence of subadditive functions on E . If $\{X_i\}_{i=1}^\infty$ is a stochastic process with state space E for which there is stationary ergodic coupling, then*

$$(3.3) \quad \lim_{n \rightarrow \infty} \frac{h_n(X_1, X_2, \dots, X_n)}{n} = c \quad \text{a.s.}$$

for some constant c .

Proof. Let $\{X'_i\}_{i=1}^\infty$ denote the stationary ergodic process to which $\{X_i\}_{i=1}^\infty$ may be coupled, and let τ be the coupling time, i.e., $\tau = \min \{i: X_i = X'_i\}$. The doubly indexed process $Y_{st} = h_{t-s}(X'_{s+1}, X'_{s+2}, \dots, X'_t)$ is easily checked to have the properties:

$$(3.4a) \quad Y_{su} \leq Y_{st} + Y_{tu} \text{ whenever } s < t < u.$$

$$(3.4b) \quad \text{The joint distributions of the shifted process } \{Y_{s+1, t+1}\} \text{ are the same as those of the unshifted process,}$$

$$(3.4c) \quad \text{The expectations } g_t = ET_{0t} \text{ satisfy } g_t \geq -At \text{ for some } A \text{ and for all } t.$$

The properties (3.4a-c) are exactly the hypotheses of Kingman's theorem (Kingman (1973)), so by its conclusion we have

$$(3.5) \quad \lim_{n \rightarrow \infty} \frac{Y_{0,m}}{n} = \lim_{n \rightarrow \infty} \frac{h_n(X'_1, X'_2, \dots, X'_n)}{n} = c \quad \text{a.s.}$$

Here, to conclude that the limit is indeed a constant we have made use of the fact that Kingman's theorem assures that the limit is shift invariant and we have assumed that $\{X'_i\}_{i=1}^\infty$ is ergodic.

Now we have by (3.1), (3.2), and the definition of τ that

$$(3.6) \quad \begin{aligned} h_n(X_1, X_2, \dots, X_n) &\leq h_\tau(X_1, X_2, \dots, X_\tau) + h_{n-\tau}(X'_{\tau+1}, X'_{\tau+2}, \dots, X'_n) \\ &\leq h_\tau(X_1, X_2, \dots, X_\tau) + h_n(X'_1, X'_2, \dots, X'_n). \end{aligned}$$

Since $\tau < \infty$ with probability one, (3.4) and (3.6) yield

$$(3.7) \quad \overline{\lim}_{n \rightarrow \infty} \frac{h_n(X_1, X_2, \dots, X_n)}{n} \leq c.$$

To handle the limit infimum we need only consider the analogous inequality with the variables reversed, i.e.,

$$h_n(X'_1, X'_2, \dots, X'_n) \leq h_\tau(X'_1, X'_2, \dots, X'_\tau) + h(X_1, X_2, \dots, X_n),$$

and we obtain

$$c \leq \underline{\lim}_{n \rightarrow \infty} \frac{h_n(X_1, X_2, \dots, X_n)}{n}$$

to complete the proof. \square

COROLLARY 1. *If $V_i, 1 \leq i < \infty$, is an irreducible, aperiodic, positive recurrent Markov chain with state space $\mathcal{A} \times \mathcal{A}$ then no matter what the initial distribution*

$\pi(v) = P(V_1 = v)$, one has with probability one

$$\lim_{n \rightarrow \infty} L_n(V_1, V_2, \dots, V_n)/n = c$$

for some constant c .

Proof. To prove the corollary one only has to exhibit an appropriate coupling; and, in this case, the existence of such a coupling is well known (see, e.g., Hoel, Port, and Stone (1972)).

One can also prove the above corollary without recourse to coupling; one can use an absolute continuity argument. Under the hypotheses of the corollary there is a stationary measure π' . Moreover, the initial measure π is absolutely continuous with respect to π' (since by the irreducibility and positive recurrence $\pi'(a_1, a_2) > 0$, for all $(a_1, a_2) \in \mathcal{A} \times \mathcal{A}$). If $\{V'_i: 1 \leq i < \infty\}$ is the process with initial distribution π' and the same transition function as $\{V_i: 1 \leq i < \infty\}$, it is further true that the measure \mathcal{P} for the infinite process $\{V_i: 1 \leq i < \infty\}$ is absolutely continuous with respect to that \mathcal{P}' for $\{V'_i: 1 \leq i < \infty\}$. Since $L(V'_1, V'_2, \dots, V'_n)$ satisfies the hypotheses of Kingman's subadditive ergodic theorem, $\{\omega: \lim_{n \rightarrow \infty} L(V'_1, V'_2, \dots, V'_n)/n = c\}$ is a set of \mathcal{P}' measure one. By absolute continuity of $\mathcal{P} \ll \mathcal{P}'$ the set $\{\omega: \lim_{n \rightarrow \infty} L(V_1, V_2, \dots, V_n)/n = c\}$ has \mathcal{P} measure one. This is precisely the conclusion of the corollary. \square

4. Alternative statistics. The random variable $L(V_1, V_2, \dots, V_n)$ certainly is an interesting measure of genetic proximity, but it appears to be hard to handle. In such a situation it is natural to look for suitable alternatives.

To introduce one such alternative let (X_1, X_2, \dots, X_n) and $(X'_1, X'_2, \dots, X'_n)$ denote two sequences of values from \mathcal{A} . By A, B we denote subsets of $\{1, 2, \dots, n\}$, say $A = \{i_1, i_2, \dots, i_h\}$ and $B = \{j_1, j_2, \dots, j_k\}$ if $|A| = |B| = k$. Next we set

$$(4.1) \quad \rho(A, B) = \begin{cases} 1 & \text{if } X_{i_1} = X'_{j_1}, \quad X_{i_2} = X'_{j_2}, \dots, \quad X_{i_k} = X'_{j_k}, \\ 0 & \text{otherwise.} \end{cases}$$

The statistic of interest in this section is

$$(4.2) \quad T_n = \sum_{A, B} \rho(A, B),$$

where the sum is over two pairs of subsets of $\{1, 2, \dots, n\}$ and it is understood that $\rho(A, B)$ is taken to be zero if the cardinalities of A and B differ, i.e., $|A| \neq |B|$.

If the $X_i, 1 \leq i < \infty$ and the $X'_i, 1 \leq i < \infty$ are all independent, and $P(X_i = a_j) = p_j, P(X'_i = a_j) = p'_j$ for all i, j , it is easy to see that

$$(4.3) \quad ET_n = \sum_{k=0}^n \binom{n}{k}^2 \left(\sum_{j=1}^{\infty} p_j p'_j \right)^k.$$

This explicit formula is quite a contrast to the mystery surrounding EL_n under similar hypotheses. A number of qualitative properties of ET_n are also evident from (4.3). In particular, if we set $p_i = p'_i$ for $1 \leq i \leq |\mathcal{A}| = a$ and take \mathcal{A} finite, then

$$(4.4) \quad \phi(\vec{p}) = E_p T_n = \sum_{k=0}^n \binom{n}{k} \left(\sum_{j=1}^a p_j^2 \right)^k$$

is easily checked to be a Schur-convex function, i.e., $\phi(\vec{p}) \leq \phi(\vec{p}')$ whenever \vec{p} is majorized by \vec{p}' . (For an elaboration of this terminology see Hardy, Littlewood and Polya (1951), and for an elaboration of the many consequences of Schur convexity see the treatise by Olkin and Marshall (1979)).

Despite the mathematical simplicity of T_n as evidenced by (4.3) and (4.4), it provides only a partial surrogate for L_n . In the first place T_n tends to be very large, and there is no efficient algorithm for finding T_n . Thus, from a computational view point, L_n is a superior statistic. Also, as yet, there is no information at all about the variance of T_n or of its limit properties.

5. Open problems. The main open problems concern the expectations

$$(5.1) \quad \psi_n(p) = EL_n$$

under the hypotheses of independence and identical distribution as applied in (4.4).

For one explicit conjecture, it seems inevitable that $\psi_n(p)$ is Schur convex (just as $\phi(p)$ was proved to be). Perhaps it would be easier to consider the limit,

$$(5.2) \quad \psi(p) = \lim_{n \rightarrow \infty} \frac{\psi_n(p)}{n}.$$

Again, it must be true that $\psi(p)$ is Schur convex, but so far even this has not been proved.

The older problems concern the numerical values of $\psi(p)$. Perhaps progress can be made on this problem by taking a more algorithmic point of view. Is there an efficient algorithm for computing the approximate value of $\psi(p)$ or $\psi_n(p)$ with a guaranteed error bound?

Given the results of § 2, it is very interesting to see if one can improve (2.10) to show $\text{Var } L_n = o(n)$. This would be the first really nontrivial step toward the Chvátal-Sankoff conjecture, and it would seem to require some genuinely new combinatorial insight to settle the point one way or the other.

Finally, the main scientific problem is to find a replacement for L_n which still has a genetic justification. The null distributions of L_n seem likely to be always out of reach, and major progress will be made when L_n finds a suitable substitute. The statistic T_n is a reasonable first choice, but it leads to its own problems. For example, what is the order of the growth of $\text{Var } T_n$?

In the search for surrogates for L_n , it may be critical to consider the variety of problems to which it has been applied. In addition to the application to molecule comparisons noted previously, there is a natural application in communications. In particular, Bradley and Bradley (1978) have applied L_n in the study of bird songs.

There are also a variety of potential uses in computer science, and for an introduction there it seems useful to refer to the papers of Aho, Hirschberg and Ullman (1976), Okuda, Tanaka and Kasai (1975), Selkow (1977) and Wagner and Fischer (1974). In at least some of these papers in which L_n has been used, it seems there must exist a more tractable substitute.

Acknowledgment. The observation that absolute continuity provides a second proof of Corollary 1 in § 3 is due to Steve Lalley who kindly commented on an earlier draft of this article.

REFERENCES

- A. V. AHO, D. S. HIRSCHBERG AND J. D. ULLMAN (1976), *Bounds on the complexity of the longest common subsequences problem*, J. Assoc. Comput. Mach., 23, pp. 1-12.
 M. ABID (1978), *Un théorème ergodique pour des processus sous-additifs et sur-stationnaires*. C. R. Acad. Sci. Paris Sér. A, 217, pp. 149-152.
 D. W. BRADLEY AND R. A. BRADLEY (1978), *Application of sequence comparison to the study of bird songs*, Tech. Rep. Dept. of Data Processing, California State Univ., Long Beach, CA.

- V. CHVÁTAL AND D. SANKOFF (1975), *Longest common subsequences of two random sequences*, J. Appl. Probab., 12, pp. 306–315.
- J. G. DEKEN (1979), *Some limit results for longest common subsequences*, Discrete Math., 26, pp. 17–31.
- B. EFRON AND C. STEIN (1981), *The jackknife estimate of variance*, Ann. Statist., 9, pp. 586–596.
- D. GRIFFAETH (1978), *Coupling methods for Markov processes*, in Studies in Probability and Ergodic Theory, Advances in Mathematics Supplementary Studies, Vol. 2, G.-C. Rota, ed., Academic Press, New York.
- J. M. HAMMERSLEY (1974), *Postulates for subadditive processes*, Ann. Probab., 2, pp. 652–680.
- G. H. HARDY, J. E. LITTLEWOOD AND G. POLYÁ (1951), *Inequalities*, Cambridge University Press, Cambridge.
- P. G. HOEL, S. C. PORT AND C. J. STONE (1972), *Introduction to Stochastic Processes*, Houghton-Mifflin, Palo Alto, CA.
- J. F. C. KINGMAN (1973), *Subadditive ergodic theory*, Ann. Probab., 1, pp. 883–909.
- A. W. MARSHALL AND I. OLKIN (1979), *Inequalities: Theory of Majorization and Its Applications*, Academic Press, New York.
- T. OKUDA, E. TANAKA AND T. KASAI (1976), *A method for correction of garbled words based on the Levenshtein metric*, IEEE Trans. Comput., C-25, pp. 172–177.
- D. SANKOFF (1972), *Matching sequences under deletion/insertion constraints*, Proc. Nat. Acad. Sci. U.S.A., 69, pp. 4–6.
- D. SANKOFF AND R. J. CEDERGREN (1973), *A test for nucleotide sequence homology*, J. Molecular Biol., 77, pp. 159–164.
- D. SANKOFF, R. J. CEDERGREN AND G. LAPALME (1976), *Frequency of insertion-deletion, transversion, and transition in evolution of 5S ribosomal RNA*, J. Molecular Evolution, 7, pp. 133–149.
- S. M. SELKOW (1977), *The tree to tree editing problem*, Inform. Processing letters, 6, pp. 1–7.
- R. A. WAGNER AND M. J. FISCHER (1974), *The string to string correction problem*, J. Assoc. Comput. Mach., 21, pp. 168–173.