



A fast algorithm for computing a longest common increasing subsequence

I-Hsuan Yang^a, Chien-Pin Huang^a, Kun-Mao Chao^{a,b,*}

^a Department of Computer Science and Information Engineering, National Taiwan University, Taipei, Taiwan

^b Institute of Networking and Multimedia, National Taiwan University, Taipei, Taiwan

Received 21 May 2003; received in revised form 9 August 2004

Available online 2 December 2004

Communicated by P.M.B. Vitányi

Abstract

Let $A = \langle a_1, a_2, \dots, a_m \rangle$ and $B = \langle b_1, b_2, \dots, b_n \rangle$ be two sequences, where each pair of elements in the sequences is comparable. A common increasing subsequence of A and B is a subsequence $\langle a_{i_1} = b_{j_1}, a_{i_2} = b_{j_2}, \dots, a_{i_l} = b_{j_l} \rangle$, where $i_1 < i_2 < \dots < i_l$ and $j_1 < j_2 < \dots < j_l$, such that for all $1 \leq k < l$, we have $a_{i_k} < a_{i_{k+1}}$. A longest common increasing subsequence of A and B is a common increasing subsequence of the maximum length. This paper presents an algorithm for delivering a longest common increasing subsequence in $O(mn)$ time and $O(mn)$ space.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Algorithms; Computational biology; Longest common subsequence; Longest increasing subsequence

1. Introduction

The longest common subsequence (LCS) problem and the longest increasing subsequence (LIS) problem are both very classical problems in computer science. By using the dynamic programming technique, the LCS problem can be solved in $O(mn)$ time. Knuth [3] posed the question of whether a sub-quadratic algorithm for the LCS problem exists. Masek and Paterson [10] gave an algorithm that runs in $O(mn/\log n)$ time,

where $n \leq m$ and the subsequences are drawn from a set of bounded size. Szymanski [13] proposed an $O((n+m)\log(n+m))$ algorithm for the special case in which no element appears more than once in an input sequence. Interested readers can refer to a recent survey by Bergroth et al. [1].

On the other hand, there is a rich history for the longest increasing subsequence problem as well, e.g., see [6,11]. Schensted [12] and Knuth [8] gave an $O(n \log n)$ time algorithm for this problem where the input is an arbitrary sequence of n numbers. For a special case in which the input sequence is a permutation of $\{1, 2, \dots, n\}$, Hunt and Szymanski [7], and

* Corresponding author.

E-mail address: kmchao@csie.ntu.edu.tw (K.-M. Chao).

Bespamyatnikh and Segal [2] gave algorithms that run in $O(n \log \log n)$ time.

In this paper, we consider the longest common increasing subsequence (LCIS) problem, whose goal is to find a maximum-length common subsequence of the two sequences such that the subsequence is increasing. This problem might arise in the situation where we wish to find the longest set of MUMs whose sequences occur in ascending order in three or more genomic sequences [4,5]. Formally speaking, let A and B be two sequences $A = \langle a_1, a_2, \dots, a_m \rangle$ and $B = \langle b_1, b_2, \dots, b_n \rangle$, where $m \geq n$ and each pair of elements in the sequences is comparable. A common increasing subsequence of A and B is a subsequence $\langle a_{i_1} = b_{j_1}, a_{i_2} = b_{j_2}, \dots, a_{i_l} = b_{j_l} \rangle$, where $i_1 < i_2 < \dots < i_l$ and $j_1 < j_2 < \dots < j_l$, such that for all $1 \leq k < l$, we have $a_{i_k} < a_{i_{k+1}}$. A longest common increasing subsequence of A and B is a common increasing subsequence of the maximum length. A straightforward $O(mn^2)$ -time algorithm for this problem is to sort the shorter sequence, and then find a longest common subsequence among the two sequences and the sorted sequence. Here we give an algorithm that solves the longest common increasing subsequence problem in $O(mn)$ time and $O(mn)$ space.

2. The algorithm

This algorithm utilizes a folklore algorithm (see [9] for more details) which runs in $O(n \log n)$ time and $O(n)$ space for the LIS problem. Define an array $L[k]$ to be the smallest ending number of an increasing subsequence of length k , where $1 \leq k \leq n$. Assume that each entry in L is initially an infinite value. For each position, perform a binary search to update L and make a backtracking link to the former number. The largest k such that $L[k]$ contains non-infinite value is the length of a longest increasing subsequence. Tracing back from $L[k]$ along with the link we established before will deliver a longest increasing subsequence. Its running time and space is $O(n \log n)$ and $O(n)$, respectively. In this paper, the array L will be explored between every pair of the prefix substrings of the two sequences. Instead of using binary search, we use linear search by a given starting point to update the array L .

Let $L_j^i[k]$ be the smallest ending number of a longest common increasing subsequence of length k

between the sequences $\langle a_1, a_2, \dots, a_i \rangle$ and $\langle b_1, b_2, \dots, b_j \rangle$, where $1 \leq i \leq m$ and $1 \leq j, k \leq n$. In order to keep the backtracking information, we need two more variables. Specifically, let $L_index_j^i[k]$ ($1 \leq i \leq m, 1 \leq j, k \leq n$) record the index pair (x, y) such that $L_j^i[k]$ is the element a_x and b_y , and let $Prev[i, j]$ ($1 \leq i \leq m, 1 \leq j \leq n$) record the index pair (x, y) such that there is a link from (i, j) to (x, y) . A link will be made in the function LIS_insert and used in the backtracking process.

It can be shown that the difference between arrays L_j^i and L_j^{i-1} as well as between L_j^i and L_{j-1}^i is at most one entry, which will be proved later in Lemma 1. We therefore use a variable χ_j^i for keeping the index of the difference between L_j^i and L_j^{i-1} as follows.

$$\chi_j^i = \begin{cases} l & \text{if } L_j^i[l] \neq L_j^{i-1}[l] \text{ for some } l, \\ -1 & \text{if } L_j^i = L_j^{i-1}. \end{cases}$$

Variable χ_j^i is used while L_j^{i-1} and L_{j-1}^i are merged in the *mismatch* case.

For $i \geq 2$ and $j \geq 2$, L_j^i can be derived from L_j^{i-1} and L_{j-1}^i . There are two cases for computing L_j^i : the *match* case if $a_i = b_j$, and the *mismatch* case if $a_i \neq b_j$. In the *match* case, we do the following:

$$L_j^i = \text{Insert } a_i \text{ into } L_j^{i-1},$$

and in the *mismatch* case, we do the following:

$$L_j^i = \text{merge}(L_{j-1}^i, L_j^{i-1}).$$

Here, *inserting* a number into L refers to the process of calling the function $LIS_insert(L, L_index, Prev, a, p, i, j)$, which inserts an element a into L , makes a link $Prev[i, j]$ to the former number in L , and returns the insertion index. Variable p stores the index of the last

$LIS_insert(\text{array } L, \text{array } L_index, \text{array } Prev, \text{element } a, \text{integer } p, \text{integer } i, \text{integer } j)$

1. $x := p$
 2. **while** ($L[x] < a$) **do**
 3. $x := x + 1$
 4. $L[x] := a$
 5. $L_index[x] := (i, j)$
 6. **if** ($x \neq 1$) **then**
 7. $Prev[i, j] := L_index[x - 1]$
 8. **return** x
-

Fig. 1. The function LIS_insert .

inserted point, and is used as the starting index for the next linear search (see Fig. 1).

In the *merging* process, $L_j^i = \text{merge}(L_{j-1}^i, L_j^{i-1})$ is defined to be for any $k \geq 1$, $L_j^i[k] = \min(L_{j-1}^i[k], L_j^{i-1}[k])$. Our algorithm first assigns L_j^{i-1} to L_j^i , and then compares $L_j^i[\chi_{j-1}^i]$ with $L_{j-1}^i[\chi_{j-1}^i]$:

$$L_j^i := L_j^{i-1},$$

$$L_j^i[\chi_{j-1}^i] := \min\{L_j^i[\chi_{j-1}^i], L_{j-1}^i[\chi_{j-1}^i]\}.$$

The resulting L_j^i is equal to $\text{merge}(L_{j-1}^i, L_j^{i-1})$, which will be proved later.

The algorithm for computing a longest common increasing subsequence is given in Fig. 2. One observation is that after finishing computing L_j^i , the array L_j^{i-1} can be discarded. Therefore, all the L_j in different rows can use the same memory space and thus avoid the time for copying L_j^{i-1} to L_j^i . Besides, χ_j^i can also be recycled to save space. Thus, L_j^i , χ_j^i , and

$L_index_j^i$ are denoted by L_j , χ , and L_index_j , respectively. Once this algorithm reaches line 17 of Fig. 2, the last non-infinite element of array L_n stores the LCIS's ending number. A longest common increasing subsequence can be delivered by tracing back along the links.

3. Correctness

Assume that L_j^{i-1} and L_{j-1}^i keep the correct ending numbers. For convenience, L_j^i is first assigned as L_j^{i-1} , and then combined with L_{j-1}^i . There are two cases: the mismatch case and the match case.

3.1. The mismatch case

In the mismatch case ($a_i \neq b_j$), we do the following:

$$L_j^i = \text{merge}(L_{j-1}^i, L_j^{i-1}).$$

INPUT: Two sequences $A = \langle a_1, a_2, \dots, a_m \rangle$ and $B = \langle b_1, b_2, \dots, b_n \rangle$, where $m \geq n$.

```

//initialization
1. for j = 1 to n do
2.   for k = 1 to n do
3.     L_j[k] := ∞
4. for j = 1 to m do
5.   for k = 1 to n do
6.     Prev[j, k] := (-1, -1)
//main program
7. for i = 1 to m do
8.   χ := -1, p := 1 //default value
9.   for j = 1 to n do
10.    if (a_i = b_j) then //the match case
11.      χ, p := LIS_insert(L_j, L_index_j, Prev, a_i, p, i, j) //return the insertion position
12.    else //the mismatch case
13.      if ((χ ≠ -1) and (L_{j-1}[χ] < L_j[χ])) then
14.        L_j[χ] := L_{j-1}[χ]
15.      else
16.        χ := -1
//recover a longest common increasing subsequence in reverse order
17. x := the largest x such that L_n[x] ≠ ∞ (if x does not exist, print "NULL" and exit)
18. (y_1, y_2) := L_index_n[x]
19. print a_{y_1}
20. while (Prev[y_1, y_2] ≠ (-1, -1)) do
21.   (y_1, y_2) := Prev[y_1, y_2]
22. print a_{y_1}

```

Fig. 2. The algorithm for computing a longest common increasing subsequence.

Any common increasing subsequence of length k in the sequences $\langle a_1, a_2, \dots, a_i \rangle$ and $\langle b_1, b_2, \dots, b_{j-1} \rangle$ or in the sequences $\langle a_1, a_2, \dots, a_{i-1} \rangle$ and $\langle b_1, b_2, \dots, b_j \rangle$ will still exist in the sequences $\langle a_1, a_2, \dots, a_i \rangle$ and $\langle b_1, b_2, \dots, b_j \rangle$. By the definition of $L_j^i[k]$, we have to choose the “smallest” one. Moreover, because $a_i \neq b_j$, there will be no more new elements added into the list, which means that a common increasing subsequence of length k will still be a common subsequence of length exactly k in the sequences $\langle a_1, a_2, \dots, a_i \rangle$ and $\langle b_1, b_2, \dots, b_j \rangle$.

Now, we prove that L_j^i computed by our algorithm is the same as $\text{merge}(L_{j-1}^i, L_j^{i-1})$.

Lemma 1. *For any $i, j > 1$ in the dynamic-programming table, there is at most one different entry between L_j^i and L_{j-1}^{i-1} (or between L_j^i and L_{j-1}^i). If the different entry is $L_j^i[k]$, then $L_j^i[k] < L_{j-1}^{i-1}[k]$ (or $L_j^i[k] < L_{j-1}^i[k]$).*

Proof. Between L_j^i and L_{j-1}^{i-1} , if a_i matches one of the letter in sequence $\langle b_1, b_2, \dots, b_j \rangle$, we will insert at most one a_i into L_{j-1}^{i-1} and cause at most one entry smaller. Similar arguments hold for L_j^i and L_{j-1}^i . \square

Lemma 2. *For any $i, j > 1$ in the dynamic-programming table, there are at most two different entries between L_{j-1}^i and L_j^{i-1} . If there are two differences, then one of them is at index χ_{j-1}^i such that $L_{j-1}^i[\chi_{j-1}^i] < L_j^{i-1}[\chi_{j-1}^i]$, and the other difference is at index λ , where $L_{j-1}^i[\lambda] > L_j^{i-1}[\lambda]$. If there is only one difference, it is at index χ_{j-1}^i .*

Proof. By Lemma 1, there is at most one different entry between L_{j-1}^i and L_j^{i-1} , and if this happens, let us assume $L_{j-1}^i[\lambda] > L_j^{i-1}[\lambda]$. There is also at most one different entry between L_{j-1}^{i-1} and L_j^i , and if this happens, let us assume $L_{j-1}^{i-1}[\chi_{j-1}^i] > L_j^i[\chi_{j-1}^i]$. If $\chi_{j-1}^i \neq \lambda$, then there are two differences and the lemma holds. If $\chi_{j-1}^i = \lambda$, then there is at most one different entry, which is recorded by χ_{j-1}^i . \square

By Lemma 2, while merging L_{j-1}^{i-1} with L_j^i , we need to compare at most two different entries.

There are four cases. First, if $L_{j-1}^{i-1} = L_{j-1}^i$ (i.e., $\chi_{j-1}^i = -1$) and $L_{j-1}^{i-1} = L_j^{i-1}$, then $L_j^i = L_{j-1}^i$. So we have $L_j^i = L_{j-1}^i = \text{merge}(L_{j-1}^i, L_j^{i-1})$. Second, if $L_{j-1}^{i-1} = L_{j-1}^i$ (i.e., $\chi_{j-1}^i = -1$) and $L_{j-1}^{i-1} \neq L_j^{i-1}$ where $L_{j-1}^{i-1}[k] > L_j^{i-1}[k]$ (Lemma 1), then for all k , $L_{j-1}^{i-1}[k] \leq L_j^{i-1}[k] = L_{j-1}^i[k]$. It follows $L_j^i = L_{j-1}^i = \text{merge}(L_{j-1}^i, L_j^{i-1})$. Third, if $L_{j-1}^{i-1} \neq L_{j-1}^i$ (i.e., $\chi_{j-1}^i \neq -1$) where $L_{j-1}^{i-1}[\chi_{j-1}^i] < L_{j-1}^i[\chi_{j-1}^i]$ and $L_{j-1}^{i-1} = L_j^{i-1}$, then for all $k \neq \chi_{j-1}^i$, $L_j^i[k] = L_{j-1}^{i-1}[k] = L_{j-1}^i$ and $L_{j-1}^{i-1}[\chi_{j-1}^i] < L_{j-1}^i[\chi_{j-1}^i] = L_{j-1}^i[\chi_{j-1}^i]$. We set $L_j^i = L_{j-1}^i$ and $L_j^i[\chi_{j-1}^i] = L_{j-1}^{i-1}[\chi_{j-1}^i]$. Thus L_j^i will be equal to $\text{merge}(L_{j-1}^i, L_{j-1}^{i-1})$. Finally, consider the case where $L_{j-1}^{i-1} \neq L_{j-1}^i$ (i.e., $\chi_{j-1}^i \neq -1$) and $L_{j-1}^{i-1} \neq L_j^{i-1}$. If there are one or two different entries between L_{j-1}^{i-1} and L_{j-1}^i , we let $L_j^i = L_{j-1}^{i-1}$ and $L_j^i[\chi_{j-1}^i] = \min(L_{j-1}^{i-1}[\chi_{j-1}^i], L_{j-1}^i[\chi_{j-1}^i])$. By Lemma 2, L_j^i is equal to $\text{merge}(L_{j-1}^i, L_{j-1}^{i-1})$. If there is no different entry between L_{j-1}^{i-1} and L_{j-1}^i , then $L_{j-1}^{i-1} = L_{j-1}^i$. Consequently, $L_j^i = L_{j-1}^i = \text{merge}(L_{j-1}^i, L_j^{i-1})$.

3.2. The match case

In the match case ($a_i = b_j$), we do the following:

$$L_j^i = \text{Insert } a_i \text{ into } L_{j-1}^{i-1}.$$

For the two sequences $\langle a_1, a_2, \dots, a_i \rangle$ and $\langle b_1, b_2, \dots, b_j \rangle$, $a_i = b_j$ is the last “common” element. That means that it can be added to any common increasing subsequence with ending number smaller than a_i . Thus, inserting a_i into L_{j-1}^{i-1} will maintain the invariance of L_j^i . Whether a_i is already in L_{j-1}^{i-1} or not, the resulting L_j^i is correct because the insertion process will do nothing if the element is already in it.

Lemma 3. *In the same row L^i , the indices p of all insertion points are non-decreasing from the index $j = 1$ to $j = n$.*

Proof. Observe that, in the same row, every element to be inserted is equal to a_i . For all $k \geq 1$, from L_1^i to L_n^i , by Lemma 1, each change between $L_j^i[k]$ and

$L_{j+1}^i[k]$ is always decreasing. Assume that the index of the last insertion point of L_j^i is p_1 , it follows $L_j^i[1], L_j^i[2], \dots, L_j^i[p_1 - 1]$ are smaller than a_i . In the next match case, this property still holds. Thus, the insertion point will not be in the range of 1 to $p_1 - 1$. \square

Lemma 3 implies that we can perform linear search from the previously recorded index p .

4. Time and space complexity

There are two loops in the algorithm of Fig. 2: an inner loop **for** $j = 1$ **to** n (line 9) and an outer loop **for** $i = 1$ **to** m (line 7). During each execution of the outer loop body (the same index i), there are n steps with two cases: *mismatch* and *match*. In the *mismatch* case, only one comparison is required. In the *match* case, the procedure *LIS_insert* will do linear search from p to the insertion point and update p . By Lemma 3, p is non-decreasing for each outer loop, and $1 \leq p \leq n$. Therefore, for each outer loop, it takes $O(n)$ time. The overall time complexity is $O(mn)$.

As Fig. 2 shows, for all $1 \leq i \leq m$, L_j^i and $L_index_j^i$ share the memory space L_j and $L_index_j^i$, respectively, and thus requires $O(n^2)$ space. Variables $Prev[i, j]$ are used for the range $1 \leq i \leq m$, $1 \leq j \leq n$, which would require $O(mn)$ space. It follows that this algorithm runs in $O(mn)$ space.

5. Discussion

We presented an algorithm for computing a longest common increasing subsequence in quadratic time. Can one do it in sub-quadratic time?

Acknowledgements

We thank the reviewers for their helpful comments that improve the presentation of the paper. Kun-Mao Chao was supported in part by an NSC grant 92-2213-E-002-073 from the National Science Council, Taiwan.

References

- [1] L. Bergroth, H. Hakonen, T. Raita, A survey of longest common subsequence algorithms, in: Proc. 7th Internat. Symp. on String Processing Information Retrieval (SPIRE'00), Spain, 2000, pp. 39–48.
- [2] S. Bspamyatnikh, M. Segal, Enumerating longest increasing subsequences and patience sorting, Inform. Process. Lett. 76 (2000) 7–11.
- [3] V. Chvatal, D.A. Klarner, D.E. Knuth, Selected combinatorial research problems, Technical Report, STAN-CS-72-292, Computer Science Department, Stanford University, 1972.
- [4] A.L. Delcher, S. Kasif, R.D. Fleischmann, J. Peterson, O. White, S.L. Salzberg, Alignment of whole genomes, Nucleic Acids Res. 27 (1999) 2369–2376.
- [5] A.L. Delcher, A. Phillippy, J. Carlton, S.L. Salzberg, Fast algorithms for large-scale genome alignment and comparison, Nucleic Acids Res. 30 (2002) 2478–2483.
- [6] M.L. Fredman, On computing the length of longest increasing subsequences, Discrete Math. 11 (1975) 29–35.
- [7] J. Hunt, T. Szymanski, A fast algorithm for computing longest common subsequences, Comm. ACM 20 (1977) 350–353.
- [8] D.E. Knuth, Sorting and Searching, The Art of Computer Programming, vol. 3, Addison-Wesley, Reading, MA, 1973.
- [9] U. Manber, Introduction to Algorithms—A Creative Approach, Addison-Wesley, Reading, MA, 1989.
- [10] W.J. Masek, M.S. Paterson, A faster algorithm computing string edit distances, J. Comput. System Sci. 20 (1) (1980) 18–31.
- [11] J. Matousek, E. Welzl, Good splitters for counting points in triangles, J. Algorithms 13 (1992) 307–319.
- [12] C. Schensted, Longest increasing and decreasing subsequences, Canad. J. Math. 13 (1961) 179–191.
- [13] T. Szymanski, A special case of the maximal common subsequence problem, Technical Report, TR-170, Computer Science Laboratory, Princeton University, 1975.