Dynamic Programming Method for Analyzing Biomolecular Sequences

Tao Jiang Department of Computer Science University of California - Riverside (Typeset by Kun-Mao Chao)

E-mail: jiang@cs.ucr.edu http://www.cs.ucr.edu/~jiang

Outline

- The paradigm of dynamic programming
- Sequence alignment a general framework for comparing sequences in bioinformatics
- Dynamic programming algorithms for sequence alignment
- Techniques for improving the efficiency of the algorithms
- Multiple sequence alignment

Dynamic Programming

- Dynamic programming is an algorithmic method for solving optimization problems with a compositional/recursive cost structure.
- Richard Bellman was one of the principal founders of this approach.

Two key ingredients

• Two key ingredients for an optimization problem to be suitable for a dynamic programming solution:

l. optimal substructures



substructure is optimal.

Otherwise, a divide-and-conquer approach is the choice.)

Three basic components

- The development of a dynamic programming algorithm has three basic components:
 - A recurrence relation (for defining the value/cost of an optimal solution);
 - A tabular computation (for computing the value of an optimal solution);
 - A backtracing procedure (for delivering an optimal solution).

Fibonacci numbers

The *Fibonacci numbers* are defined by the following recurrence:

$$F_0 = 0 F_1 = 1 F_i = F_{i-1} + F_{i-2} \text{ for } i > 1.$$







 Given a sequence of real numbers a₁a₂...a_n, find a consecutive subsequence with the maximum sum.

9 -3 1 7 -15 2 3 -4 2 -7 6 -2 8 4 -9

For each position, we can compute the maximum-sum interval starting at that position in O(n) time. Therefore, a naive algorithm runs in $O(n^2)$ time.

O-notation: an asymptotic upper bound

• f(n) = O(g(n)) iff there exist two positive constant *c* and n_0 such that $0 \le f(n) \le cg(n)$ for all $n \ge n_0$



For example, 5n + 108 = O(n) and $2n = O(n\log n)$.

How functions grow?

function n	30 <i>n</i>	92n log n	26 <i>n</i> ²	0.68 <i>n</i> ³	2"
100	0.003 sec.	0.003 sec.	0.26 sec.	0.68 sec.	4 x 10 ¹⁶ yr.
100,000	3.0 sec.	2.6 min.	3.0 days	22 yr.	

(Assume one million operations per second.)

For large data sets, algorithms with a complexity greater than $O(n \log n)$ are often impractical!

11

Maximum sum interval (the recurrence relation)

• Define *S*(*i*) to be the maximum sum of the intervals ending at position *i*.

$$S(i) \leftarrow a_i + \max \begin{cases} S(i-1) \\ 0 \end{cases}$$

If S(i-I) < 0, concatenating a_i with its previous interval gives less sum than a_i itself.





Defining scores for alignment columns

- *infocon* [Stojanovic *et al.*, 1999]
 - Each column is assigned a score that measures its information content, based on the frequencies of the letters both within the column and within the alignment.

CGGATCAT—GGA CTTAACATTGAA GAGAACATAGTA

1

Defining scores (cont'd)

phylogen [Stojanovic *et al.*, 1999]
 – columns are scored based on the evolutionary relationships among the sequences implied by a supplied phylogenetic tree.



Two fundamental problems we solved (joint work with Lin and Chao)

• Given a sequence of real numbers of length *n* and an upper bound *U*, find a consecutive subsequence of length at most *U* with the maximum sum --- an *O*(*n*)-time algorithm.

U = 3 9 -3 1 7 -15 2 3 -4 2 -7 6 -2 8 4 -9

17

Two fundamental problems we solved (joint work with Lin and Chao)

Given a sequence of real numbers of length *n* and a lower bound *L*, find a consecutive subsequence of length at least *L* with the maximum average --- an O(n log L)-time algorithm. This has been improved to O(n) by others later.

L=4

3 2 14 6 6 2 10 2 6 6 14 2 1

Another example

Given a sequence as follows: 2, 6.6, 6.6, 3, 7, 6, 7, 2

and L = 2, the highest-average interval is the squared area, which has the average value 20/3.

2, 6.6, 6.6, 3, 7, 6, 7, 2

GC-rich regions

• Our method can be used to locate a region of length at least *L* with the highest C+G ratio in *O*(*n* log *L*) time.

ATGACTCGAGCTCGTCA 00101011011011010 • Search for an interval of length at least *L* with the highest average.



A naive algorithm
A simple shift algorithm can compute the highest-average interval of a fixed length in



A pigeonhole principle

• Notice that the length of an optimal interval is bounded by 2*L*, we immediately have an O(*nL*)-time algorithm.



regments, where each of them is of length >= L.

Future Development

- Best *k* (nonintersecting) subsequences?
- Max-average with both upper and lower length bounds
- General (gapped) local alignment with length upper bound.
- Measurement of goodness?



• The longest increasing subsequence is to find a longest increasing subsequence of a given sequence of distinct integers $a_1a_2...a_n$.

e.g. 9 2 5 3 7 11 8 10 13 6

2 3 7
5 7 10 13
9 7 11
3 5 11 13
are increasing subsequences.
We want to find a longest one.









An $O(n \log n)$ method for LIS

• Define *BestEnd*[*k*] to be the smallest end number of an increasing subsequence of length *k*.



An $O(n \log n)$ method for LIS

 Define *BestEnd*[k] to be the smallest end number of an increasing subsequence of length k.

9 2	5 3	7	11	8	10	13	6	
9 2	2, 2,	2	2	2	2	2	2	← BestEnd[1]
	5 3	3,	3	3	3	3	3,	BestEnd[2]
		7	7,	7	7	7	6	← BestEnd[3]
For each position	we perfe		11	8	8,	8	8	BestEnd[4]
a binary search to t	ipdate	<i>и</i> ш			10	10,	10	← BestEnd[5]
BestEnd. Therefore	e, the					13	13	BestEnd[6]
running time is O()	$n \log n$.							

Longest Common Subsequence (LCS)

- A subsequence of a sequence S is obtained by deleting zero or more symbols from S.
 For example, the following are all subsequences of "president": pred, sdn, predent.
- The longest common subsequence problem is to find a maximum length common subsequence between two sequences.

LCS

For instance, Sequence 1: president Sequence 2: providence Its LCS is priden.



LCS

Another example: Sequence 1: algorithm Sequence 2: alignment One of its LCS is algm.



How to compute LCS?

- Let $A = a_1 a_2 \dots a_m$ and $B = b_1 b_2 \dots b_n$
- len(i, j): the length of an LCS between $a_1a_2...a_i$ and $b_1b_2...b_j$
- With proper initializations, *len(i, j)* can be computed as follows.















Sequence A: CTTAACT Sequence B: CGGATCAT

An alignment of A and B:

C---TTAACT ← Sequence A CGGATCA--T ← Sequence B







Scoring Matrices

Amino acid substitution matrices

-PAM

- BLOSUM
- DNA substitution matrices
 - DNA is less conserved than protein sequences
 - Less effective to compare coding regions at nucleotide level

PAM

- Point Accepted Mutation (Dayhoff, et al.)
- 1 PAM = PAM₁ = 1% average change of all amino acid positions
 - After 100 PAMs of evolution, not every residue will have changed
 - some residues may have mutated several times
 - some residues may have returned to their original state
 - some residues may not changed at all

					PA	M	X					
• PAI	• • •	= P /	AM									
<i>E</i> . • PAN	<i>g</i> . M ₂₅₀	PA is	AM ₂ a w	^{250 =} ide	= P. ly t	AM isec	(₁ 25) I sc) orii	1g 1	nat	rix.	
Asp D												

BLOSUM

- Blocks Substitution Matrix
- Scores derived from *observations* of the frequencies of substitutions in blocks of local alignments in related proteins
- Matrix name indicates evolutionary distance
 - BLOSUM62 was created using sequences sharing no more than 62% identity



An optimal alignment -- an alignment of maximum score

- $S_{i,i}$ the score of an optimal alignment between

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + w(a_i, -) \\ s_{i,j-1} + w(-, b_j) \\ s_{i-1,j-1} + w(a_i, b_j) \end{cases}$$





			2	5 _{3,5}	=	?			
		С	G	G	А	Т	С	А	Т
	0	-3	-6	-9	-12	-15	-18	-21	-24
С	-3		5	2	-1	-4	-7	-10	-13
Т	-6		3	0	-3	7	4	1	-2
Т	-9		0	-2	-5	?			
A	-12								
A	-15								
С	-18								
Т	-21								

			5	5 _{3,5}	=	?				
		С	G	G	Α	Т	С	А	Т	
	0	-3	-6	-9	-12	-15	-18	-21	-24	
С	-3	8	5	2	-1	-4	-7	-10	-13	
т	-6	5	3	0	-3	7	4	1	-2	
т	-9	2	0	-2	-5	5	-1	-4	9	
А	-12	-1	-3	-5	6	3	0	7	6	
Α	-15	-4	-6	-8	3	1	-2	8	5	
С	-18	-7	-9	-11	0	-2	9	б	3 ફ	optimal score
Т	-21	-10	-12	-14	-3	8	6	4	14 ⁴	54

ТТ	A i	A C	-	Т					
GG	Α .	ГС	A	Т					
- 5 - 5	+8 -	$\frac{5+8}{7}$	-3 -	+8 =	14				
		C	G	G	A	T	<u> </u>	A	T
		-3	-6	-9	-12	-15	-18	-21	-24
С	-3	8	5	2	-1	-4	-7	-10	-13
т	-6	5	3	0	-3	7	4	1	-2
Т	-9	2	0		-5	5	-1	-4	9
Α	-12	-1	-3	-5			0	7	6
Α	-15	-4	-6	-8	3		-2	8	5
С	-18	-7	-9	-11	0	-2		6	3
Т	-21	-10	-12	-14	-3	8	6	4	





Match: 8 Mismatch: -5		10	oca	l a	lig	nn	ien	t	
Gap symbol: -3		С	G	G	A	Т	С	A	Т
	0	0	0	0	0	0	0	0	0
С	0	8	5	2	0	0	8	5	2
т	0	5	3	0	0	8	5	3	13
т	0	2	0	0	0	8	5	2	11
А	0	0	0	0	8	5	3	?	
А	0								
С	0								
т	0								

Match: 8 Mismatch: -5		10	oca	.l a	lig	nm	ien	t		
Gap symbol: -3		С	G	G	Α	Т	С	A	Т	
	0	0	0	0	0	0	0	0	0	
С	0	8	5	2	0	0	8	5	2	
т	0	5	3	0	0		5	3	13	
т	0	2	0	0	0		5	2	11	
А	0	0	0	0	8		3	13	10	
А	0	0	0	0	8		2	11		the
С	0	8	5	2	5		13	10		best score
т	0	5	3	0	2	13	10	8	18°	59

A - C - A T C A 8-3+8-3-	т т +8 =	- 18								
		С	G	G	А	Т	С	А	Т	
	0	0	0	0	0	0	0	0	0	
С	0	8	5		0	0	8	5		
т	0	5	3	0	0	8	5	3	13	
т	0	2	0	0	0	8	5	2	11	
А	0	0	0		8	5	3	13	10	(not always at the corner)
А	0	0	0	0	8	5	2	11		the
С	0	8	5	2	5	3	13	10		best score
Т	0	5	3	0	2	13	10	8	18°	60
										80





Affine gap penalties

- Let D(i, j) denote the maximum score of any alignment between a_ja₂...a_i and b_jb₂...b_j ending with a deletion.
- Let *I*(*i*, *j*) denote the maximum score of any alignment between *a₁a₂...a_i* and *b₁b₂...b_j* ending with an insertion.
- Let S(i, j) denote the maximum score of any alignment between a₁a₂...a_i and b₁b₂...b_i.

Affine gap penalties	
$D(i, j) = \max\begin{cases} D(i-1, j) - y\\ S(i-1, j) - x - y\\ \left[I(i, j-1) - y \right] \end{cases}$	
$I(i, j) = \max \begin{cases} S(i, j-1) - x - y \\ S(i, j-1) - x - y \end{cases}$	
$\begin{bmatrix} D(i,j) \\ I(i,j) \end{bmatrix}$	



k best local alignments

- Smith-Waterman (Smith and Waterman, 1981; Waterman and Eggert, 1987)
- FASTA (Wilbur and Lipman, 1983; Lipman and Pearson, 1985)
- BLAST (Altschul et al., 1990; Altschul et al., 1997)
- BLAST and FASTA are key genomic database search tools.

k best local alignments

- Smith-Waterman
 - Smith and Waterman, 1981; Waterman and Eggert, 1987)
 - linear-space version : sim (Huang and Miller, 1991)
 - linear-space variants : sim2 (Chao et al., 1995); sim3 (Chao et al., 1997)
- FASTA
 - (Wilbur and Lipman, 1983; Lipman and Pearson, 1985)
 linear-space band alignment (Chao et al., 1992)
- BLAST
 (Altschul et al., 1990; Altschul et al., 19

67

FASTA

- Find runs of identities, and identify regions with the highest density of identities.
- 2) Re-score using PAM matrix, and keep top scoring segments.
- 3) Eliminate segments that are unlikely to be part of the alignment.
- 4) Optimize the alignment in a band.

Its running time is O(n).









BLAST

- 1) Build the hash table for sequence A (the database sequence).
- 2) Scan sequence B for hits.
- 3) Extend hits.

Also O(n) time.

BLAST

Step 1: Build the hash table for se	equence A. (3-tuple example)
For DNA sequences:	For protein sequences:
Seg. A = AGATCGAT 12345678 AAA AAC AGA =1 MTC =3 GGA =5 GGA =5 TCG =4 	Seq. A = ELVIS Add xyz to the hash table if $Score(xyz, ELV) \ge T$; Add xyz to the hash table if $Score(xyz, LVI) \ge T$; Add xyz to the hash table if $Score(xyz, VIS) \ge T$;



Remarks

- Filtering is based on the observation that a good alignment usually includes short identical or very similar fragments.
- The idea of filtration was used in both FASTA and BLAST to achieve high speed









- Space: *O*(*m*+*n*)
- Time: $O(mn)^*(1+\frac{1}{2}+\frac{1}{4}+...) = O(mn)$ 2







Multiple sequence alignment (MSA)

• The multiple sequence alignment problem is to simultaneously align more than two sequences.

GCTC	GC-1
AC	
GATC	G-A
	GCTC AC GATC

85

How to score an MSA?

• *Sum-of-Pairs* (SP-score)





• an $O(n^3)$ algorithm



General MSA

- For *k* sequences of length *n*: $O(n^k)$
- NP-Complete (Wang and Jiang)
- The exact multiple alignment algorithms for many sequences are not feasible.
- Some approximation algorithms are given. (*e.g.*, 2- *l/k* for any fixed *l* by Bafna *et al.*)

Progressive alignment

- A heuristic approach proposed by Feng and Doolittle.
- It iteratively merges the most similar pairs.
- "Once a gap, always a gap"



The time for progressive alignment in most cases is roughly the order of the time for computing all pairwise alignment, i.e., $O(k^2n^2)$.

Concluding remarks

- Three essential components of the dynamic programming approach:
 - the recurrence relation
 - the tabular computation
 - the backtracing
- The dynamic-programming approach has been used in a vast number of computational problems in bioinformatics.