CS 466 Introduction to Bioinformatics Lecture 6

Mohammed El-Kebir

September 19, 2018



Course Announcements

Instructor:

- Mohammed El-Kebir (melkebir)
- Office hours: Mondays, 3:15-4:15pm

TA:

- Anusri Pampari (pampari2)
- Office hours: Thursdays, 11:00-11:59am in SC 4105

TA office hour canceled on Sept. 20

Solution to HW 1 released on Sept. 21

Midterm on Oct. 10, 7-9pm, 1310 DCL

Outline

- Two open questions on Hirschberg algorithm
- RNA secondary structure

Reading:

• Topics are not in Jones and Pevzner book but in lecture notes and slides [Based on Chapter 10 in "Biological sequence analysis" by Durbin et al.]

Hirschberg Algorithm: Reversing Edges Necessary?

m

Max weight path from (0,0) to (m,n) through $(i^*, n/2)$ n/2 $i^* = \arg \max\{ \operatorname{prefix}(i) + \operatorname{suffix}(i) \}$ $i^* = \arg \max\{ \operatorname{prefix}(i) + \operatorname{suffix}(i) \}$ Compute $\{\operatorname{prefix}(i) \mid 0 \le i \le m\}$ in O(mn) time and O(m) i^* space, by starting from (0,0) to $\left(m, \frac{n}{2}\right)$ keeping only two
columns in memory. [single-source multiple destinations]

Hirschberg Algorithm: Reversing Edges Necessary?



Doing a shortest path from each $(i, \frac{n}{2})$ to (m, n) (for all $0 \le i \le m$) will not achieve desired running time!

Reversing edges enables single-source multiple destination computation in desired time and space bound!

Hirschberg Algorithm: Reconstructing Alignment



 A
 T
 G
 T
 C

 A
 T
 C
 G
 C

Hirschberg(i, j, i', j')

- **1.** if j' j > 1
- 2. $i^* \leftarrow \underset{0 \le i \le m}{\operatorname{arg max wt}(i)}$

3. Report
$$(i^*, j + \frac{j'-j}{2})$$

4. Hirschberg
$$(i, j, i^*, j + \frac{j'-j}{2})$$

5. Hirschberg
$$(i^*, j + \frac{j'-j}{2}, i', j')$$

Problem: Given reported vertices and scores $\{(i_0, 0, s_0), \dots, (i_n, n, s_n)\}$, find intermediary vertices.

Transposing matrix does not help, because gaps could occur in both input sequences

Outline

• Two open questions on Hirschberg algorithm

• RNA secondary structure

Reading:

• Topics are not in Jones and Pevzner book but in lecture notes and slides [Based on Chapter 10 in "Biological sequence analysis" by Durbin et al.]

Central Dogma of Molecular Biology

Three fundamental molecules:

1. DNA

Information storage.

2. RNA

Old view: Mostly a "messenger". New view: Performs many important functions, through **3-D structure**!

3. Protein

Perform most cellular functions (biochemistry, signaling, control, etc.)

$\mathsf{DNA} \xrightarrow{} \mathsf{RNA} \xrightarrow{} \mathsf{Protein}$



Copyright © 1997, by John Wiley & Sons, Inc. All rights reserved.

RNA



• Single-stranded

- A (adenine)
- C (cytosine)
- U (uracil)
- G (guanine)
- Can fold into structures due to nucleotide complementarity.
 A <--> U, C <--> G
- Comes in many flavors:

mRNA, rRNA, tRNA, tmRNA, snRNA, snoRNA, scaRNA, aRNA, asRNA, piwiRNA, etc.

RNA – Nucleotide Complementarity

RNA can fold into structures due to nucleotide complementarity: A <--> U and G <--> C

A <--> U (2 hydrogen bonds) is slightly weaker than G <--> C (3 hydrogen bonds)

G <--> U also observed but not as stable



transfer RNA (tRNA) Secondary Structure



http://hyperphysics.phy-astr.gsu.edu/hbase/Organic/trna.html#c3



http://bioinfo.bisr.res.in/project/crat/pictures/codon.jpg

RNA Secondary Structure Elements



Nesting and Pseudoknot



Most RNA molecules consist of nested base pairs

Nesting and Pseudoknot – Examples

Nesting

5'-GCGGAUUCUGCCCCAAUUCGCACCA-3' (((((((----)))))))----



Nesting and Pseudoknot – Examples



RNA can fold into structures due to nucleotide complementarity:

A <--> U and G <--> C

Secondary structure is determined by a set of non-overlapping complimentary base pairs

RNA can fold into structures due to nucleotide complementarity:

A <--> U and G <--> C

Secondary structure is determined by a set of non-overlapping complimentary base pairs

Question: How to find maximum number of such pairs?

RNA can fold into structures due to nucleotide complementarity:

A <--> U and G <--> C

Secondary structure is determined by a set of non-overlapping complimentary base pairs

Question: How to find maximum number of such pairs?

Need to constrain space of feasible solutions!

RNA can fold into structures due to nucleotide complementarity:

A <--> U and G <--> C

Secondary structure is determined by a set of non-overlapping complimentary base pairs

Question: How to find maximum number of such pairs?

SIAM J. APPL, MATH. Vol. 35, No. 1, July 1978 Need to constrain space of feasible solutions!

© Society for Industrial and Applied Mathematics 0036-1399/78/3501-0006 \$01.00/0

ALGORITHMS FOR LOOP MATCHINGS*

RUTH NUSSINOV,[†] GEORGE PIECZENIK,[‡] JERROLD R. GRIGGS¶ AND DANIEL J. KLEITMAN§

Problem: Given RNA sequence $\mathbf{v} \in \{A, U, C, G\}^n$, find a *pseudoknot-free secondary structure* with the maximum number of complementary base pairings



Nussinov Algorithm – Dynamic Programming

Problem: Given RNA sequence $\mathbf{v} \in \{A, U, C, G\}^n$, find a *pseudoknot-free secondary structure* with the maximum number of complementary base pairings

Let s[i, j] denote the maximum number of pseudoknot-free complementary base pairings in subsequence v_i, \dots, v_j

Nussinov Algorithm – Dynamic Programming

Problem: Given RNA sequence $\mathbf{v} \in \{A, U, C, G\}^n$, find a *pseudoknot-free secondary structure* with the maximum number of complementary base pairings

Let s[i, j] denote the maximum number of pseudoknot-free complementary base pairings in subsequence v_i, \dots, v_j



Nussinov Algorithm – Dynamic Programming

Problem: Given RNA sequence $\mathbf{v} \in \{A, U, C, G\}^n$, find a *pseudoknot-free secondary structure* with the maximum number of complementary base pairings

Let s[i, j] denote the maximum number of pseudoknot-free complementary base pairings in subsequence v_i, \dots, v_j



Γ, **(1)**

Γ, **(1*)**

(2)

(3)

(4)

$$s[i, j] = \max \begin{cases} 0, & \text{if } i \ge j, \\ s[i+1, j-1] + 1, & \text{if } i < j \text{ and } (v_i, v_j) \in \\ s[i+1, j-1], & \text{if } i < j \text{ and } (v_i, v_j) \notin \\ s[i+1, j], & \text{if } i < j, \\ s[i, j-1], & \text{if } i < j, \\ \max_{i < k < j} \{s[i, k] + s[k+1, j]\}, & \text{if } i < j, \end{cases}$$

Develop Intuition [Spreadsheet/Whiteboard]

Nussinov Algorithm – Traceback Step



Question: Will this return one alignment? Or all alignments?

Question: Can we do this recursively?

We only need to know matches.

	G	G	G	Α	Α	Α	U	С	С
G	0								
G	0	0							
G	0	0	0						
Α	0	0	0	0					
Α	0	0	0	0	0				
Α	0	0	0	0	0	0			
U	0	0	0	0	0	0	0		
С	0	0	0	0	0	0	0	0	
С	0	0	0	0	0	0	0	0	0



	G	G	G	Α	Α	Α	U	С	С
G	0	0							
G	0	0	0						
G	0	0	0	0					
Α	0	0	0	0	0				
Α	0	0	0	0	0	0			
Α	0	0	0	0	0	0	1		
U	0	0	0	0	0	0	0	0	
С	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0





if i < j and $(v_i, v_j) \in \Gamma$, (1)if i < j and $(v_i, v_j) \notin \Gamma$, (1*)if i < j,if i < j,if i < j,if i < j,(3)if i < j,(4)

(4)

k k+1

bifurcation







Nussinov Algorithm – Alternative Solutions

	G	G	G	Α	Α	Α	U	С	С
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	2
Α	0	0	0	0	0	0	1	1	1
Α	0	0	0	0	0	0	1	1	1
Α	0	0	0	0	0	0	1	1	1
U	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0



	G	G	G	Α	Α	Α	U	С	С
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	2
Α	0	0	0	0	0	0	1	1	1
Α	0	0	0	0	0	0	1	1	1
Α	0	0	0	0	0	0	1	1	1
U	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0



	G	G	G	Α	Α	Α	U	С	С
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	2
Α	0	0	0	0	0	0	1	1	1
Α	0	0	0	0	0	0	1	1	1
Α	0	0	0	0	0	0	1	1	1
U	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0
С	0	0	0	0	0	0	0	0	0



Nussinov Algorithm – Example With Bifurcation



GCACGACG

Does this make sense?







Adenosine (A)









GCACGACG ()•((•))

Extension: Hairpin Loops with Minimum Length ℓ





Extension: Hairpin Loops with Minimum Length ℓ



$$s[i,j] = \max \begin{cases} 0, & \text{if } i + \ell \ge j, \\ s[i+1,j-1] + 1, & \text{if } i + \ell < j \text{ and } (v_i, v_j) \in \Gamma, \\ s[i+1,j-1], & \text{if } i + \ell < j \text{ and } (v_i, v_j) \notin \Gamma, \\ s[i+1,j], & \text{if } i + \ell < j, \\ s[i,j-1], & \text{if } i + \ell < j, \\ \max_{i+\ell < k < j} \{s[i,k] + s[k+1,j]\}, & \text{if } i + \ell < j, \end{cases}$$
(2)
(3)
(4)

RNA Secondary Structure Prediction in Practice

Rather than maximize number of compl. base pairs, minimize free energy (FE)

Zuker's algorithm: Dynamic programming w/ three matrices similar to affine gap penalties

- V(i,j): FE of optimal structure of s[i..j] assuming i,j form a base pair
- VBI(i,j): FE of optimal structure of s[i..j] assuming i,j closes a bulge or internal loop
- VM(i,j): FE of optimal structure of s[i..j] assuming i,j closes a multibranch loop



FE minimization with pseudoknots is NP-hard [Lyngso and Pedersen, RECOMB 2000]

Summary

- RNA is a sequence of four bases/nucleotides {A, U, C, G}
- RNA folds into structures due to base/nucleotide complementarity
 - A <--> U and C <--> G
- RNA secondary structure is defined by a set of non-overlapping complementary nucleotide pairs
- Pseudoknot-free structures have no "crossing" pairs
- Nussinov Algorithm: Dynamic programming to find pseudoknot-free structure with maximum number of complementary nucleotide pairs

Reading:

• Topics are not in Jones and Pevzner book but in lecture notes and slides [Based on Chapter 10 in "Biological sequence analysis" by Durbin et al.]