CS 466: Final Review

December 9, 2020

Review Session Topics

- 1. Solutions to HW5 (Genome Assembly)
- 2. RNA Secondary Structure Prediction
- 3. Phylogeny Inference
 - Distance-based
 - Character-based
- 4. Hidden Markov Models

Genome Assembly

Homework 5 Solution Review

Clarification

- Theorem from class slides:
 - A note is semi-balanced if indegree differs from outdegree by 1.
 - A directed, connected graph is Eulerian if and only if it has at most 2 semi-balanced nodes and all other nodes are balanced.



• Here is the full theorem:

A directed graph has an Eulerian trail if and only if at most one vertex has (out-degree) – (in-degree) = 1, at most one vertex has (in-degree) – (out-degree) = 1, every other vertex has equal in-degree and out-degree, and all of its vertices with nonzero degree belong to a single connected component of the underlying undirected graph.

RNA Secondary Structure Prediction

Things to Remember

- RNA can fold into structures due to nucleotide complementarity
 - A <--> U and G <--> C
- We can use dot-parenthesis format to represent secondary structures
- Find a pseudoknot-free secondary structure with the maximum number of complementary base pairings
 - Use dynamic programming with the Nussinov Algorithm
 - Recall that in homework we also looked at how to fill out a backtrace given an optimal solution.
 - Watch out for bifurcation when calling $\max_{i < k < j} \{s[i,k] + s[k+1,j]\}$
 - Watch out for more than one optimal solution



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

	0,	if $i \geq j$,	
	$\overline{s[i+1,j-1]+1},$	if $i < j$ and $(v_i, v_j) \in \Gamma$,	(1)
mov	s[i+1,j-1],	if $i < j$ and $(v_i, v_j) \notin \Gamma$,	(1*)
max	s[i+1,j],	if i < j,	(2)
	s[i,j-1],	if i < j,	(3)
	$\max_{i < k < j} \{ s[i,k] + s[k+1,j] \},\$	if i < j,	(4)





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

max 〈	$ \begin{cases} 0, \\ \hline s[i+1,j-1]+1, \\ s[i+1,j-1], \\ s[i+1,j], \\ s[i,j-1], \\ \max_{i \le k \le j} \{s[i,k]+s[k+1,j]\}, \end{cases} $	$\begin{array}{l} \text{if } i \geq j, \\ \\ \text{if } i < j \text{ and } (v_i, v_j) \in \Gamma, \\ \\ \text{if } i < j \text{ and } (v_i, v_j) \not\in \Gamma, \\ \\ \\ \text{if } i < j, \\ \\ \\ \text{if } i < j, \\ \\ \\ \\ \text{if } i < j, \end{array}$	(1) (1*) (2) (3) (4)
	$\max_{i < k < j} \{s[i,k] + s[k+1,j]\},\$	if i < j,	(4)



9

10

G	U	С	С	U	Α
З	4	5	6	7	8
(-)	-	()

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

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

	G	G	G	U	C	C	U	A	U	C		
G	0	0	0	0	1	2					G	1
G	0	0	0	0	1	2	2				G	2
G	0	0	0	0	1	1	1	2			G	3
U	0	0	0	0	0	0	0	1	1		υ	4
С	0	0	0	0	0	0	0	1	1	1	с	5
С	0	0	0	0	0	0	0	1	1	1	с	6
U	0	0	0	0	0	0	0	1	1	1	υ	7
Α	0	0	0	0	0	0	0	0	1	1	A	8
U	0	0	0	0	0	0	0	0	0	0	υ	9
С	0	0	0	0	0	0	0	0	0	0	с	10

4 5 6 7 8

1

2

3



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

max ‹	$\begin{cases} 0, \\ s[i+1, j-1] + 1, \\ s[i+1, j-1], \\ s[i+1, j], \\ s[i, j-1], \\ \max_{i \le k \le j} \{s[i, k] + s[k+1, j]\}, \end{cases}$	$ \begin{array}{l} \text{if } i \geq j, \\ \text{if } i < j \text{ and } (v_i, v_j) \in \Gamma, \\ \text{if } i < j \text{ and } (v_i, v_j) \not\in \Gamma, \\ \text{if } i < j, \end{array} $	(1) (1*) (2) (3) (4)
	$\left(\max_{i < k < j} \{s[i,k] + s[k+1,j]\}\right),\$	If $i < j$,	(4)

Phylogeny Inference

Distance-Based Methods

Things to Remember

- Tree topology represents similarity/distance between sequences
- Hierarchical clustering is one way to get tree from distance matrix
 - Greedy algorithm that merges closest clusters until just one cluster remains
 - The definition of distance between clusters, known as the linkage criterion, affects clustering result
 - Complete linkage take the pairwise max
 - Single linkage take the pairwise min
 - Mean linkage take the pairwise average
 - Neighbor joining also hierarchical clustering. Intuitively, clusters that are more similar relative to their distance to other clusters are merged
 - Produces correct tree when matrix is additive













(a,b)

d

b

(c,d)

a

С



 $\min\{D(a,(c,d)),D(b,(c,d))\}$



Things to Remember

- Tree topology represents similarity/distance between sequences
- We also saw how construct a tree from an additive matrix
 - Recall how to check if a matrix is additive
 - Check if distance matrix (non-negativity, zero diagonal, symmetry, triangle inequality)
 - Check Four Point Condition (for every set of four leaves, compute 3 sums, 2 should be equal, one should be less than or equal)
 - Small Additive Distance Problem: leaf-labeled tree + additive matrix



- Large Additive Distance Problem: additive matrix only
 - Find trimming parameter and remove degenerate triplet
 - Remember to decrease entries in matrix by two times the trimming parameter

Check if Matrix is Additive Example

	а	b	С	d	
а	0	6	7	8	
b	6	0	11	12	
С	7	11	0	5	
d	8	12	5	0	

- Is it a distance matrix?
 - Non-negative 🗹
 - Zero diagonal 🗹
 - Symmetric
 - Triangle inequality ?

 \checkmark

Check if Matrix is Additive Example

	а	b	С	d
а	0	6	7	8
b	6	0	11	12
С	7	11	0	5
d	8	12	5	0

- Is it a distance matrix?
 - Non-negative 🗹
 - Zero diagonal 🗹
 - Symmetric
 - Triangle inequality ?
- Check Four Point Condition

- $d_{ab} + d_{cd} = 6 + 5 = 11$
- $d_{ac} + d_{bd} = 7 + 12 = 19$
- $d_{ad} + d_{bc} = 8 + 11 = 19$

Check if Matrix is Additive Example

	а	b	С	d
а	0	6	7	8
b	6	0	11	12
С	7	11	0	5
d	8	12	5	0

- Is it a distance matrix?
 - Non-negative 🗹
 - Zero diagonal 🗹
 - Symmetric
 - Triangle inequality ?
- Check Four Point Condition

 \checkmark

- $d_{ab} + d_{cd} = 6 + 5 = 11$
- $d_{ac} + d_{bd} = 7 + 12 = 19$
- $d_{ad} + d_{bc} = 8 + 11 = 19$



• $11 = d_{ab} + d_{cd} \le d_{ac} + d_{bd} = d_{ad} + d_{bc} = 19$

Phylogeny Inference

Character-Based Methods

Things to Remember

- Tree topology represents fewest state changes along edges
- Small Parsimony Problem: leaf-labeled tree + character matrix
 - Find labels of internal nodes in given tree maximizing parsimony
 - Recall that characters can be solved independently
 - Apply Sankoff dynamic programming algorithm to each character
 - Subproblem: minimum parsimony score of the subtree rooted at vertex v if v has character state t
- Large Parsimony Problem: character matrix only
 - In general, the problem is NP-hard (including multi-state perfect phylogeny)
 - Special case: two-state (i.e., binary) perfect (i.e., infinite sites) phylogeny
 - Check if conflict-free: no pair of columns contain the three pairs (0, 1), (1, 0) and (1, 1)
 - To reconstruct tree, sort columns and apply graph algorithm

	c1	c2	c3	c4
а	0	0	1	0
b	1	0	0	0
С	0	1	1	0
d	0	1	1	1

	c1	c2	c3	c4		c3	c2	c1	c4
а	0	0	1	0	а	1	0	0	0
b	1	0	0	0	Sort b	0	0	1	0
С	0	1	1	0	c1 = 1 c2 = 2 C	1	1	0	0
d	0	1	1	1	$\begin{vmatrix} c3 \\ c3 \end{vmatrix} = 3 \\ c4 \end{vmatrix} = 1$ d	1	1	0	1

	c1	c2	c3	c4	1	c3	c2	c1	c4	0, 0, 0, 0
а	0	0	1	0	а	1	0	0	0	
b	1	0	0	0	Sort b	0	0	1	0	0, 0, 1, 0
С	0	1	1	0	c1 = 1 c2 = 2 C	1	1	0	0	
d	0	1	1	1	c3 = 3 c4 = 1 d	1	1	0	1	а

	c1	c2	c3	c4	1	c3	c2	c1	c4	0, 0, 0, 0	
а	0	0	1	0		а	1	0	0	0	
b	1	0	0	0	Sort	b	0	0	1	0	1, 0, 0, 0 /
С	0	1	1	0	c1 = 1 c2 = 2	С	1	1	0	0	
d	0	1	1	1	c3 = 3 c4 = 1	d	1	1	0	1	b a



С



+ d

	c1	c2	c3	c4		c3	c2	c1	c4
а	0	0	1	0	а	1	0	0	0
b	1	0	0	0	Sort b	0	0	1	0
С	0	1	1	0	c1 = 1 c2 = 2 C	1	1	0	0
d	0	1	1	1	c3 = 3 c4 = 1 d	1	1	0	1

Reminder – you can also check if the matrix is conflict free (i.e., no forbidden submatrix) to determine if it is a perfect phylogeny.





Hidden Markov Models

Things to Remember

- The HMM setup requires the following
 - Set of hidden states Q
 - Transition probability matrix A
 - Future state only depends on current state
 - Set of emitted symbols $\boldsymbol{\Sigma}$
 - Emission probability matrix *E*



• We observe the emitted symbols but not the hidden states

Things to Remember

- We considered three questions about HMM systems
 - What is the most probable path π^* that generated observations x?
 - Viterbi algorithm
 - What is probability of observations x generated by any path π ?
 - Forward algorithm
 - What is the probability of observation x_i generated by state s?
 - Forward and backward algorithm
- Recall that we often do computation in logspace to avoid underflow

What is the probability of observing $\mathbf{x} = [H, T]$?





$$f[s,i] = \begin{cases} a_{0,s}e_{s,x_1}, & \text{if } i = 1, \\ e_{s,x_i} \sum_{t \in Q} \{f[t,i-1] \cdot a_{t,s}\}, & \text{if } i > 1. \end{cases}$$

What is the probability of observing $\mathbf{x} = [H, T]$?



What is the probability of observing $\mathbf{x} = [H, T]$?

Н Т .25 F .13 В .375 .09 $f[F,2] = e_{F,T} \sum_{t \in Q} f[t,1] \cdot a_{t,F}$ $= \frac{1}{2} \cdot \left((.25 \cdot .9) + (.375 \cdot .1) \right) \approx 0.13$ $f[B,2] = e_{B,T} \sum_{t \in Q} f[t,1] \cdot a_{t,B}$ $=\frac{1}{4} \cdot ((.25 \cdot .1) + (.375 \cdot .9)) \approx 0.09$

$$Q = \{F, B\}$$

$$F = B$$

$$A = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix} \stackrel{F}{B}$$

$$\Sigma = \{H, T\}$$

$$H = \begin{pmatrix} 0.5 & 0.5 \\ 0.75 & 0.25 \end{pmatrix} \stackrel{F}{B}$$

$$O.9 \qquad 0.9$$

$$O.1 \qquad B$$

$$F = \begin{pmatrix} 0.1 & B \\ 0.1 & B \\ 0.1 & B \\ 0.5 & 0.5 & 0.75 & 0.25 \end{pmatrix}$$

$$f[s,i] = \begin{cases} a_{0,s}e_{s,x_1}, & \text{if } i = 1, \\ e_{s,x_i} \sum_{t \in Q} \{f[t,i-1] \cdot a_{t,s}\}, & \text{if } i > 1. \end{cases}$$

What is the probability of observing $\mathbf{x} = [H, T]$?



$$Q = \{F, B\}$$

$$F = B$$

$$A = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix} F$$

$$E = \begin{pmatrix} 0.5 & 0.5 \\ 0.75 & 0.25 \end{pmatrix} F$$

$$O.9 \qquad 0.9$$

$$O.1 \qquad 0.1$$

$$F = \begin{pmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.5 \\ 0.$$

$$f[s,i] = \begin{cases} a_{0,s}e_{s,x_1}, & \text{if } i = 1, \\ e_{s,x_i} \sum_{t \in Q} \{f[t,i-1] \cdot a_{t,s}\}, & \text{if } i > 1. \end{cases}$$

What is the probability of observing $\mathbf{x} = [H, T]$?

$$\begin{array}{c|c} \mathbf{H} & \mathbf{T} \\ \mathbf{F} & \boxed{.25 & .13} \\ \mathbf{B} & \boxed{.375 & .09} \end{array} \end{array} \begin{array}{c} \mathbb{W} \text{e sum the final column} \\ \text{to find the total probability} \\ .13 + .09 & = .22 \end{array} \end{array} \qquad \begin{array}{c} Q = \{F, B\} \\ A = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix} \stackrel{F}{B} \\ A = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix} \stackrel{F}{B} \\ \Sigma = \{H, T\} \\ H & T \\ E = \begin{pmatrix} 0.5 & 0.5 \\ 0.75 & 0.25 \end{pmatrix} \stackrel{F}{B} \\ 0.5 & 0.5 & 0.75 & 0.25 \end{array}$$

$$\begin{array}{c} = \frac{1}{2} \cdot \left((.25 \cdot .9) + (.375 \cdot .1)\right) \approx 0.13 \end{array} \qquad \begin{array}{c} f[B, 2] = e_{B,T} \sum_{t \in Q} f[t, 1] \cdot a_{t,B} \\ = \frac{1}{4} \cdot \left((.25 \cdot .1) + (.375 \cdot .9)\right) \approx 0.09 \end{array} \qquad \begin{array}{c} \log_2 f[s, i] = \begin{cases} \log_2 a_{0,s} + \log_2 e_{s,x_1} \\ \log_2 f[s, i] = \begin{cases} \log_2 a_{0,s} + \log_2 e_{s,x_1} \\ \log_2 e_{s,x_1} + \log_2 \sum_{t \in Q} 2^{\log_2 f[t, t-1]} \cdot 2^{\log_2 a_{t,s}} \\ \log_2 f[s, i] = \begin{cases} \log_2 e_{s,x_1} + \log_2 \sum_{t \in Q} 2^{\log_2 f[t, t-1]} \cdot 2^{\log_2 a_{t,s}} \\ \log_2 a_{0,s} + \log_2 \sum_{t \in Q} 2^{\log_2 f[t, t-1]} \cdot 2^{\log_2 a_{t,s}} \\ \log_2 a_{0,s} + \log_2 \sum_{t \in Q} 2^{\log_2 f[t, t-1]} \cdot 2^{\log_2 a_{t,s}} \\ \log_2 a_{0,s} + \log_2 \sum_{t \in Q} 2^{\log_2 f[t, t-1]} \cdot 2^{\log_2 a_{t,s}} \end{array}$$

Exam Released Tonight at 7pm

Good luck!!