

Group HW5

CS 685-001/PPA 784-003/STA 695-001 Fall 2009

Due: Due Nov 5th, 2009

Problem 1 Suppose we have the hidden state space S' and observed state space S with the transition matrix and emission matrix:

$$\theta = \begin{matrix} & \text{fair} & \text{loaded} \\ \text{fair} & \left(\begin{matrix} 4/5 & 1/5 \\ 5/6 & 1/6 \end{matrix} \right) \\ \text{loaded} & \end{matrix}$$

$$\theta' = \begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 \\ \text{fair} & \left(\begin{matrix} 1/6 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 \\ 1/12 & 1/6 & 1/12 & 1/4 & 1/24 & 1/24 \end{matrix} \right) \\ \text{loaded} & \end{matrix}$$

respectively. The initial distribution is $2/3$ for fair die and $1/3$ for loaded die.
Compute the probability to observe the sequence 46.

Problem 2 Go back to the bar hopping example. Suppose we have the transition matrix

	A	C	G	T
Your die A	1/4	1/4	1/4	1/4
Arne's die C	1/5	1/5	2/5	1/5
Connie's die G	1/3	1/3	1/6	1/6
Dave's die T	1/6	1/3	1/3	1/6

and suppose you choose the initial bar at random (i.e., probability $1/4$ for going to each bar). If you and your friends go to the bar and record which bar you go, you observe the sequence ACGTTTCGA. What is the probability to observe the sequence.

Problem 3 Compute $P(X_{25} = A)$? Please note that you do not want to compute by hand...

Problem 4 Also compute a stationary distribution π . Is it unique and the same as the limiting distribution π ?

Problem 5 (CpG island) DiAN has three tetrahedral dies. The first die corresponds to DNA that is G + C rich, the second die corresponds to DNA that is G + C poor, and the third is a fair die.

	A	C	G	T
first die	0.15	0.33	0.36	0.16
second die	0.27	0.24	0.23	0.26
third die	0.25	0.25	0.25	0.25

(1)

The transition matrix (warning: I completely made up!!) is

	first	second	third
first die	31/100	33/100	9/25
second die	29/100	4/25	11/20
third die	1/12	1/12	5/6

The initial distribution for the hidden states is $1/2$, $1/4$, $1/4$ for the first, second and third dies, respectively. What is the probability to observe the observed sequence data ACGTTTCGA with the sequence with hidden states, third third third first first second third third third?

Problem 6 Go back to the example,

$$D = \begin{array}{c|cccccc} & 1 & 2 & 3 & 4 & 5 & 6 \\ \hline 1 & 0 & 6 & 8 & 9 & 12 & 11 \\ 2 & 6 & 0 & 6 & 7 & 10 & 9 \\ 3 & 8 & 6 & 0 & 3 & 6 & 5 \\ 4 & 9 & 7 & 3 & 0 & 5 & 4 \\ 5 & 12 & 10 & 6 & 5 & 0 & 5 \\ 6 & 11 & 9 & 5 & 4 & 5 & 0 \end{array}$$

Reconstruct the tree via the NJ method by NOT hand.

Problem 7 Download PHYLIP software from <http://www.phylip.com/> and download the data sets 5000 data sets comprising the true 40-taxon trees and the corresponding 500-bp homologous sequences from <http://www.atgc-montpellier.fr/phyml/datasets.php>

Problem 8 A software `treedist` from PHYLIP computes Robinson-Foulds symmetric distance between two trees <http://portal.litbio.org/Registered/Help/phylip/doc/treedist.html> Take the first 1,000 data from the data sets you downloaded and construct the NJ trees from these data (under Kimura 2 parameter model for computing pairwise distances). Then compare the tree topologies between each reconstructed NJ tree and the true tree via Robinson-Foulds symmetric distance (so you have to do 1,000 comparisons so you do not want to do by hand).

Problem 9 From these results, get the five number summary statistics and get the histogram of the Robinson-Foulds symmetric distances you computed in Problem 2.