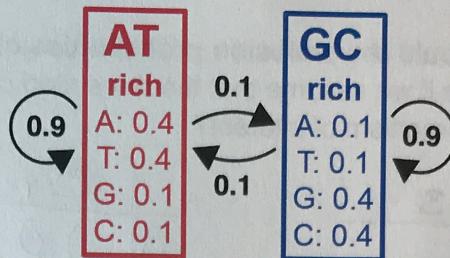


GS 373 Homework 7 - KEY

Bioinformatics Questions (80 points)

1. (30 points) Consider the hidden Markov model that describes DNA sequences in two different states, AT rich and GC rich. The numbers above the arrows denote transition probabilities, and the numbers within the boxes denote emission probabilities. Assume the initial probabilities of the two states are the same.



- a. (15 points) Based on this model, calculate the joint probability of observing the sequence **ACT** with each of the following state paths:

i. The state path AT-rich, GC-rich, AT-rich

$$\begin{aligned}
 &= P(\text{init in AT-rich}) * P(\text{emit A | AT-rich}) * P(\text{AT-rich} \rightarrow \text{GC-rich}) * P(\text{emit C | GC-rich}) * \\
 &\quad P(\text{GC-rich} \rightarrow \text{AT-rich}) * P(\text{emit T | AT-rich}) \\
 &= 0.5(0.4)(0.1)(0.1)(0.4) = \boxed{0.00032}
 \end{aligned}$$

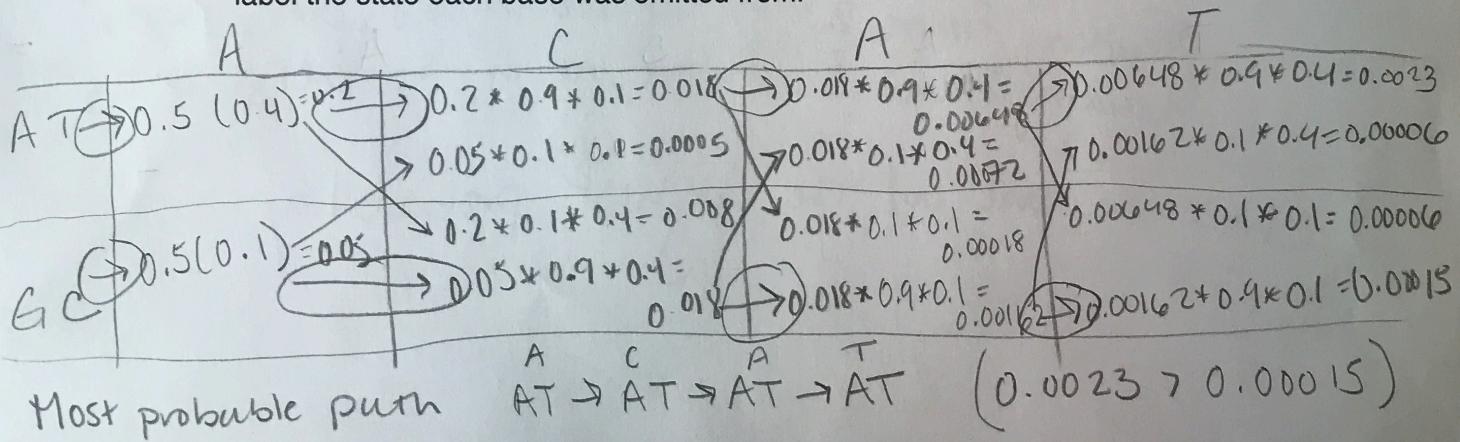
ii. The state path AT-rich, AT-rich, AT-rich

$$\begin{aligned}
 &= P(\text{init AT}) * P(\text{emit A | AT}) * P(\text{AT} \rightarrow \text{AT}) * P(\text{emit C | AT}) * \\
 &\quad P(\text{AT} \rightarrow \text{AT}) * P(\text{emit T | AT}) \\
 &= 0.5(0.4)(0.9)(0.1)(0.9)(0.4) = \boxed{0.00648}
 \end{aligned}$$

iii. The state path GC-rich, GC-rich, GC-rich

$$\begin{aligned}
 &= P(\text{init GC}) * P(\text{emit A | GC}) * P(\text{GC} \rightarrow \text{GC}) * P(\text{emit C | GC}) * \\
 &\quad P(\text{GC} \rightarrow \text{GC}) * P(\text{emit T | GC}) \\
 &= 0.5(0.1)(0.9)(0.4)(0.9)(0.1) = \boxed{0.00162}
 \end{aligned}$$

- b. (15 points) Use the Viterbi algorithm to find the most likely state path for the observed sequence **ACAT** according to this model. Show your work and be sure to label the state each base was emitted from.



2. (30 points) The hidden Markov model below is a variation of the one proposed in class to identify open reading frames in DNA sequence, where the subscripts denote the position of a nucleotide within a codon. "s" refers to the start codon, "e" refers to the stop nucleotide in a start codon, ORF₂ refers to the 2nd nucleotide in a codon between the start and stop codons, etc.

- a. (20 points) What should the emission probabilities of each of the hidden states e₁, e₂, e₃ and s₁, s₂, s₃ be if we assume that the three stop codons UAA, UAG, and UGA are equally likely in the organism of interest?

	s ₁	s ₂	s ₃
A	1.0	0.0	0.0
C	0.0	0.0	0.0
G	0.0	1.0	0.0
T	0.0	0.0	1.0

	e ₁	e ₂	e ₃
A	0.0	1/3	2/3
C	0.0	0.0	0.0
G	0.0	2/3	1/3
T	1.0	0.0	0.0

- b. (10 points) Suppose you know that a specific nucleotide in an mRNA sequence is part of a coding region because it was shown experimentally to be directly bound by ribosomes. Describe how you could use the HMM to infer the most likely reading frame for the mRNA sequence.

Use forward-back algorithm to calculate the probability the nucleotide is in reading frame 1, 2, or 3.

Reading Frame 1: is the nt in s₁, ORF₁, or e₁?

Reading Frame 2: is the nt in s₂, ORF₂, or e₂?

Reading Frame 3: is the nt in s₃, ORF₃, or e₃?

Choose the one w/ the highest probability

3. (20 points) Train the parameters of an HMM with the two states A-rich and T-rich that describe sequences of A's and T's using the training sequence below, whose state sequence is also known.

Sequence	A	T	T	T	T	A	A	A	A	A	A	A	A	A	A	A	A	A	A	T	T	T	T
State	t	t	t	t	a	a	a	a	a	a	a	a	a	a	a	a	a	a	t	t	t	t	

- a. (10 points) Compute the transition probabilities for the states learned from the above sequence.

- There are 4 transition probabilities: $t \rightarrow t, t \rightarrow a$
 $a \rightarrow a, a \rightarrow t$

- The probabilities must sum to 1.0 coming out of a given state

transition	formula	answer
$t \rightarrow t$	$(\#t \rightarrow t) / (\#t \rightarrow t + \#t \rightarrow a)$	$6 / (6+2) = 6/8$
$t \rightarrow a$	$(\#t \rightarrow a) / (\#t \rightarrow a + \#t \rightarrow t)$	$2 / (6+2) = 2/8$
$a \rightarrow a$	$(\#a \rightarrow a) / (\#a \rightarrow a + \#a \rightarrow t)$	$8 / (8+2) = 8/10$
$a \rightarrow t$	$(\#a \rightarrow t) / (\#a \rightarrow t + \#a \rightarrow a)$	$2 / (8+2) = 2/10$

- b. (10 points) Compute the emission probabilities for the states learned from the above sequence.

- There are 4 emission probs
- $P(\text{emit } T | t)$
 $P(\text{emit } A | t)$
 $P(\text{emit } T | a)$
 $P(\text{emit } A | a)$

- The prob must sum to 1.0 for a given state

emission	formula	answer
$\text{emit } T t$	$\#\text{emit } T t / (\#\text{emit } T t + \#\text{emit } A t)$	$6 / (6+3) = 6/9$
$\text{emit } A t$	$\#\text{emit } A t / (\#\text{emit } T t + \#\text{emit } A t)$	$3 / (6+3) = 3/9$
$\text{emit } T a$	$\#\text{emit } T a / (\#\text{emit } T a + \#\text{emit } A a)$	$2 / (2+8) = 2/10$
$\text{emit } A a$	$\#\text{emit } A a / (\#\text{emit } A a + \#\text{emit } T a)$	$8 / (2+8) = 8/10$