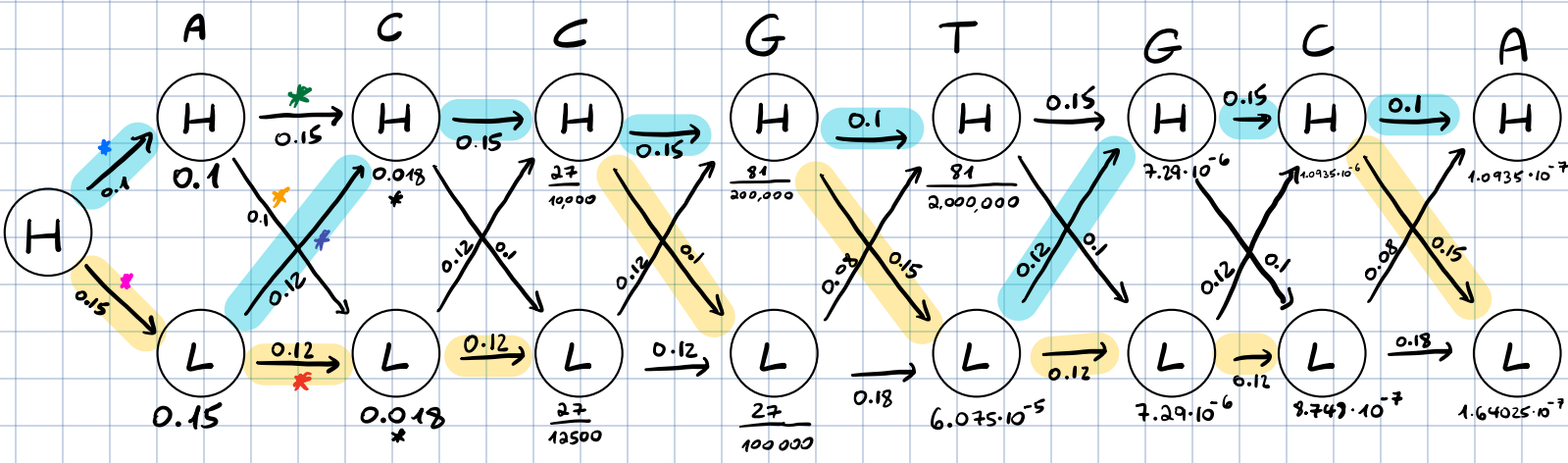


1. (10 pts) Consider this (toy) biological setup:
 A cell can be in one of two states - H , for high GC-content, and L for low GC. On each time step the cell produces one nucleotide, A,C,T or G, and might also change its state. The probability of changing from state H to L is 0.5, and from state L to H is 0.4.
 In state H the probabilities for producing nucleotides are 0.2 for A, 0.3 for C, 0.3 for G and 0.2 for T. In L the probabilities are 0.3 for A, 0.2 for C, 0.2 for G and 0.3 for T.
 Consider the nucleotide sequence $S = ACCGTGCA$. Use the Viterbi algorithm to find the best state-sequence and calculate the probability of S given this state-sequence. Assume the previous state before S was H .

1. תנאים ראשוניים: $q(H|L)=0.4 \quad q(L|L)=0.6$
 $q(H|H)=0.5 \quad q(L|H)=0.5$

• $e(A|H)=0.2 \quad e(A|L)=0.3 \quad e(G|H)=0.3 \quad e(G|L)=0.2$
 • $e(C|H)=0.3 \quad e(C|L)=0.2 \quad e(T|H)=0.2 \quad e(T|L)=0.3$



גוף מחישוב עבור המקרה הראשון:

* $q(H|H)e(A|H)=0.1$

* $q(H|L)e(C|H)=0.12$

* $q(L|H)e(A|L)=0.15$

* $q(L|L)e(C|L)=0.1$

* $q(H|H)e(C|H)=0.15$

* $q(L|L)e(C|L)=0.12$

סדרת המצבים עם ההסתברות המכה ביותר היא: $LH H H L H H L$
 ההסתברות של S בהינתן רצף המצבים זהה היא:
 $e(A|L)e(C|H)e(C|H)e(G|H)e(T|L)e(G|H)e(C|H)e(A|L)=6.561 \cdot 10^{-5}$

2. (10 pts) In class we saw the trigram HMM model and the corresponding Viterbi algorithm. We will now make two main changes. First, we will consider a four-gram tagger, where p takes the form:

$$p(x_1 \cdots x_n, y_1 \cdots y_{n+1}) = \prod_{i=1}^{n+1} q(y_i | y_{i-3}, y_{i-2}, y_{i-1}) \prod_{i=1}^n e(x_i | y_i) \quad (1)$$

We assume in this definition that $y_0 = y_{-1} = y_{-2} = *$, where $*$ is the START symbol, $y_{n+1} = STOP$, and $y_i \in \mathcal{K}$ for $i = 1 \cdots n$, where \mathcal{K} is the set of possible tags in the HMM.

Second, we consider a version of the Viterbi algorithm that takes as input **an integer** n (and not a sentence $x_1 \cdots x_n$ as we saw in class) and finds

$$\max_{y_1 \cdots y_{n+1}, x_1 \cdots x_n} p(x_1 \cdots x_n, y_1 \cdots y_{n+1})$$

for a four-gram tagger, as defined in Equation 1. $x_1 \cdots x_n$ may range over the values of some fixed vocabulary \mathcal{V} . Complete the following pseudo-code of this version of the Viterbi algorithm for this model. The pseudo-code must be efficient.

Input: An integer n , parameters $q(w|t, u, v)$ and $e(x|s)$.

Definitions: Define \mathcal{K} to be the set of possible tags. Define $\mathcal{K}_{-2} = \mathcal{K}_{-1} = \mathcal{K}_0 = \{*\}$, and $\mathcal{K}_k = \mathcal{K}$ for $k = 1 \cdots n$. Define \mathcal{V} to be the set of possible words.

Initialization: ...

Algorithm: ...

Return: ...

$$(n+1) \times (|K|+1) \times (|K|+1) \times (|K|+1) \quad \text{ריבוע } \pi \quad \text{הבט } \int_{\partial K} \omega_K$$

$$\underset{x \in V}{\operatorname{argmax}} e(x|v) = x_v \quad \text{K3NJ} \quad v \in \bigcap \delta\delta \quad .\pi(0, *, *, *) = 1 \quad \text{InjKJ}$$

for $k=1, \dots, n$

for $w \in K_{k-2}$ $u \in K_{k-1}$ $v \in K_k$

$$\pi(k, w, u, v) = \max_{z \in K_{k-3}} (\pi(k-1, z, w, u) \times q(v|z, w, u) \times e(x_u|v))$$

Return $\max_{\substack{w \in K_{k-2} \\ u \in K_{k-1} \\ v \in K_k}} (\pi(u, w, u, v) \times q(\text{stop} | w, u, v))$

בגלגולים עבדו כל ימי אנו עובדים על $\text{ker } \nu$ הק"מ לפי,
עם המילה שצאנו מהקומה $e(xiv)$.

באגרוף אין דזעם טעגלעכע $(\text{זאגט}) - (\text{זאגט})$
וועגן א גענעם מוזאיקע געווען מיליאנען
שטחן $(\text{זאגט}) - (\text{זאגט})$.

נשים געבן איר העלפער וואו (זאגט) וואו איר זענט
"קא" (וואו) משה שטען מעגלעך אז אלע ווערען א אלע הענט
- (זאגט) .

בהחזרה אין עזרים אז (זאגט) תשים.
סה"כ: זען הייבט א האלדזשען היינט (זאגט) .

b ii)

Error rate for known words is 0.07044
 Error rate for unknown words is 0.74346
 Total error rate is 0.14731

c iii)

Error rate for known words viterbi is 0.19838
 Error rate for unknown words viterbi is 0.74346
 Total error rate viterbi is 0.26064

d ii)

Error rate for known words viterbi add one is 0.16699
 Error rate for unknown words viterbi add one is 0.73386
 Total error rate viterbi add one is 0.23174

e ii)

Error rate for known words viterbi pseudo is 0.21738
 Error rate for unknown words viterbi pseudo is 0
 Total error rate viterbi pseudo is 0.21738

e iii)

Error rate for known words viterbi add one pseudo is 0.20941
 Error rate for unknown words viterbi add one pseudo is 0
 Total error rate viterbi add one pseudo is 0.20941

| | START | AT | NP | NN | JJ | VBD | ... | DTX | BEM | WP | UH | QLP | |
|-------|-------|-------|-------|--------|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| START | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| AT | 0.0 | 856.0 | 0.0 | 1.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| NP | 0.0 | 8.0 | 257.0 | 248.0 | 9.0 | 1.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| NN | 0.0 | 18.0 | 27.0 | 1377.0 | 27.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| JJ | 0.0 | 12.0 | 8.0 | 180.0 | 321.0 | 2.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| BEM | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| WP | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| UH | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| QLP | 0.0 | 0.0 | 0.0 | 4.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

מה תבינות עמוקה במטריצה שמן לב כי הפסיאית הנפוצות ביותר שן
קורות עבור מילים שחזינו להן תום מממ א ש שיה שן תום
המקורי.