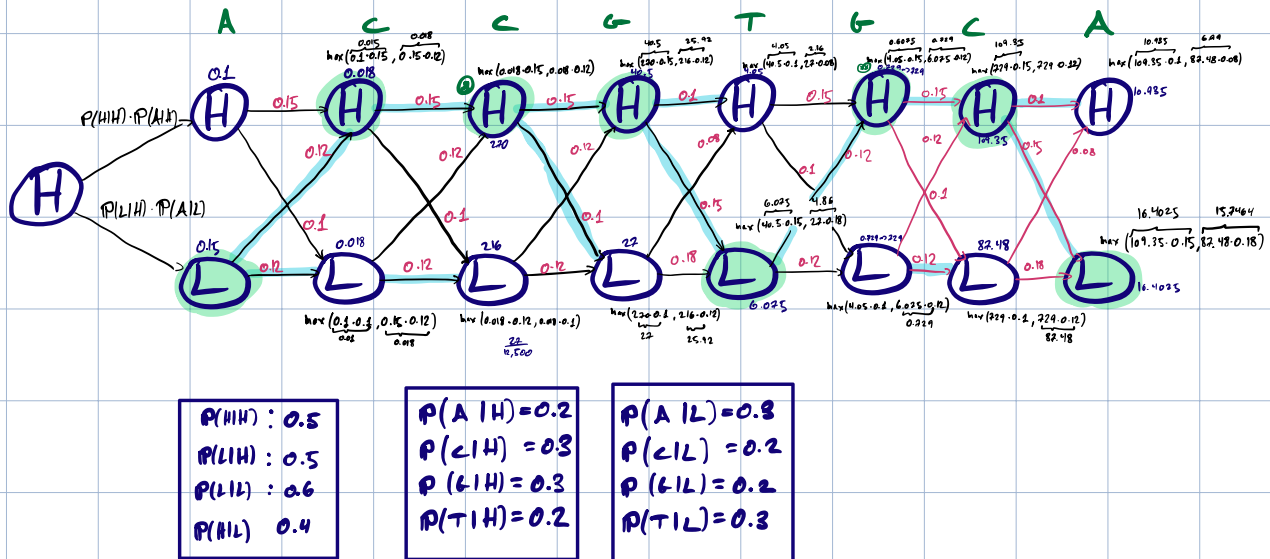1. (10 pts) Consider this (toy) biological setup:
A cell can be in one of two states - $H$, for high GC-content, and $L$ for low GC. On each time step the cell produces one nucleotide, A,C,T or G, and might also change its state. The probability of changing from state $H$ to $L$ is 0.5, and from state $L$ 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 נבנה את כל המסלול של מעבר בין המצבים:



| | |
|---|---|
| **P(H\|H) : 0.5** | **P(A \|H)=0.2** | **P(A \|L)=0.3** |
| **P(L\|H) : 0.5** | **P(C\|H) =0.3** | **P(C\|L) =0.2** |
| **P(L\|L) : 0.6** | **P(G\|H)= 0.3** | **P(G\|L)=0.2** |
| **P(H\|L)  0.4** | **P(T\|H)=0.2** | **P(T\|L)=0.3** |

— המסלול, נבחר ③ המסלול את כל בכפולה של 100,000 כדי לקבל את מספר גדול יותר.

‎ומהו‎ ‎מען‎ ‎אפל‎ ‎שלל‎ ‎כה:‎

$$L-H-H-H-L-H-H-L$$

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:

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

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:** $\cdots$
**Algorithm:** $\cdots$
**Return:** $\cdots$

② ‎כלנו‎ ‎ישאל‎ ‎תאור‎ ‎היות‎ High-level ‎3‎ ‎בשלבים‎.

‎בשלבים‎ ‎נוס‎ ‎חלק‎ ‎אף‎ ‎vtebi;‎ ‎אך‎ ‎בלל‎ ‎מק‎ ‎שישת‎ ‎ל‎ ‎קומת‎ ‎קדום‎.

‎בכל‎ ‎אלא‎ ‎שלב‎ ‎labels‎ ‎שומל‎ ‎אותוש‎ ‎זת‎ ‎אלא‎ ‎שאית‎ ‎ל‎ ‎label-‎ם.

‎מקורים‎ ‎דכסי‎ ‎של‎ ‎4-gram.

‎נוש‎ ‎זל‎ ‎ל‎ ‎קולית‎ ‎דס‎ ‎מען‎ ‎אותו‎ $x$ ‎בשיל‎, ‎שלא‎ $x \in \mathcal{V}$

!‎ושב‎ ‎ת‎ ‎מלא‎ ‎ה‎ ‎במשינ‎.

בהמשך, את $V=\{x_1,x_2\}$, $S_L$, נבנה גרף כזה בכל פעם בכל:

$k-1$

$s_i,s_j,s_l \in S$

בנה גרף כזה את מסלול קצר:

**:Init**

$$\text{Set } \pi(0,*,*,*)=1$$

**:Algorithm**

For $k=1,\ldots\ldots h:$

For $u_2 \in S_{k-2}$, $u_1 \in S_{k-1}$, $u \in S_k$, $x \in V:$

$$\pi(k, u_2, u_1, u, x) = \max_{\substack{w \in S_{k-3} \\ x_1 \in V}} \left( \pi(k-1, w, u_2, u_1, x_1) \cdot q(u \mid w, u_2, u_1) \cdot e(x \mid u) \right)$$

· Return

$$\text{Return} \quad \max_{\substack{u_2 \in S_{n-2}, u_1 \in S_{n-1} \\ u \in S_n, x \in V}} \left( \pi(n, u_2, u_1, u, x) \cdot q(\text{STOP} \mid u_2, u_1, u) \right)$$

(b) **Implementation of the most likely tag baseline**

    i. Using the training set, compute for each word the tag that maximizes $p(tag|word)$, based on the maximum likelihood estimation. Assume that the most likely tag of all the unknown words is "NN". (Unknown words are words that appear in the test set but not in the training set.)

    ii. Using the test set, compute the error rate (i.e., $1-accuracy$) for known words and for unknown words, as well as the total error rate.

```
Error rate for known words is: 0.0704399684933048
Error rate for unknown words is: 0.743455497382199
total Eroor rate is: 0.14731386424798165
```

(c) **Implementation of a bigram HMM tagger**

    i. Training phase: Compute the transition and emission probabilities of a bigram HMM tagger directly on the training set using maximum likelihood estimation.

    ii. Implement the Viterbi algorithm corresponding to the bigram HMM model. (Choose an arbitrary tag for unknown words.)

    iii. Run the algorithm from c)ii) on the test set. Compute the error rates and compare to the results from b)ii).

```
Error rate for known words is: 0.16057162146956228
Error rate for unknown words is: 0.7382198952879582
total Eroor rate is: 0.22655237715538723
```

כמו כן יודע השתמשם 3 האחד בכתב שאנו נשתמש בהם כולם כמו כאן.

אם התוך, חים, נריח ל Bigram-Hmm-model ו באותת ו'ל, ויא ראש emission יש
שתשתם ף ל P(x/y) כמו חות שומת ראש כתוב חתמה.

אם כך, גם ל Shom משתמש, ל שאו כתבתם, Unknown words.

ⓓ

```
With Smoothing:
Error rate for known words is: 0.16293462360751665
Error rate for unknown words is: 0.7347294938917975
Total Eroor rate is: 0.22824678560749523
```

לפי ן על גם אבר את ל ל שתם שרם פיע בכרות, Unknown words, בשות.
שאות כמו Bigram-Hmm בכלל.

כאו ל כ ישלש כל Bigram-HMM כאן באות וודה emission כאן כאן ;על
שאים כלך ין Total error כמו שאו ישתמש כ ל.

(שאות שאו ל את Unknown words ל אבת שאות).

(e) **Using pseudo-words**

    i. Design a set of pseudo-words for unknown words in the test set and low-frequency words in the training set.

    ii. Using the pseudo-words as well as maximum likelihood estimation (as in c)i)), run the Viterbi algorithm on the test set. Compute the error rates and compare to the results from b)ii), c)iii) and d)ii).

    iii. Using the pseudo-words as well as Add-One smoothing (as in d)i)), run the Viterbi algorithm on the test set. Compute the error rates and compare to the results from b)ii), c)iii), d)ii) and e)ii). For the results obtained using both pseudo-words and Add-One smoothing, build a confusion matrix and investigate the most frequent errors. A confusion matrix is an $|\mathcal{K}|$ over $|\mathcal{K}|$ matrix, where the $(i, j)$ entry corresponds to the number of tokens which have a true tag $i$ and a predicted tag $j$.

```
With Pseudo Words:
Error rate for known words is: 0.17321986426801683
Error rate for unknown words is: 0.6573333333333333
Total Eroor rate is: 0.20940895046347052
```

שאלה e

נשים לב שכל מה שבב אבר עבדנו שיפור יותר נאשות.

ב error לגב e על unknown words נמאלה הקטנה לעמצעו נקודה.

כמו כן, בפראות עולה ב c וכל, העים הרכבי נמצמ לעפצ עם הנאות

יותר הקטנה.

נאמיית e error העים ב b ומ נאל כי אכ, אל נאמיית

e error unknown words העים, נאלות הנאל. וכן הכי אל

```
With Pseudo Words and Add One smoothing:
Error rate for known words is: 0.1361628783798341
Error rate for unknown words is: 0.6253333333333333
Total Eroor rate is: 0.1727299910296023
```

ניתן לראות אוולי נאלם הרכבי נמצם עבל אכ ב error אל

ה ב unknown words.

נמדית ה error,Total ונים 6 ל) לוי אחן גה אל,

דלים מן של נ'או, הנטל ,דלעני ברים נוטף את נשל

גה אחין של ה Shown unknown.

נתחיון ה :cofusion matix

\* קחה נקה נפות ,ביחה את השואת גו בדף בטה גתוי.zip