

2 527 1

207766 u78 - 7'23 711C

2091KK013 - גן גיבון ז'ק

1. (10 pts) Consider this (toy) biological setup:

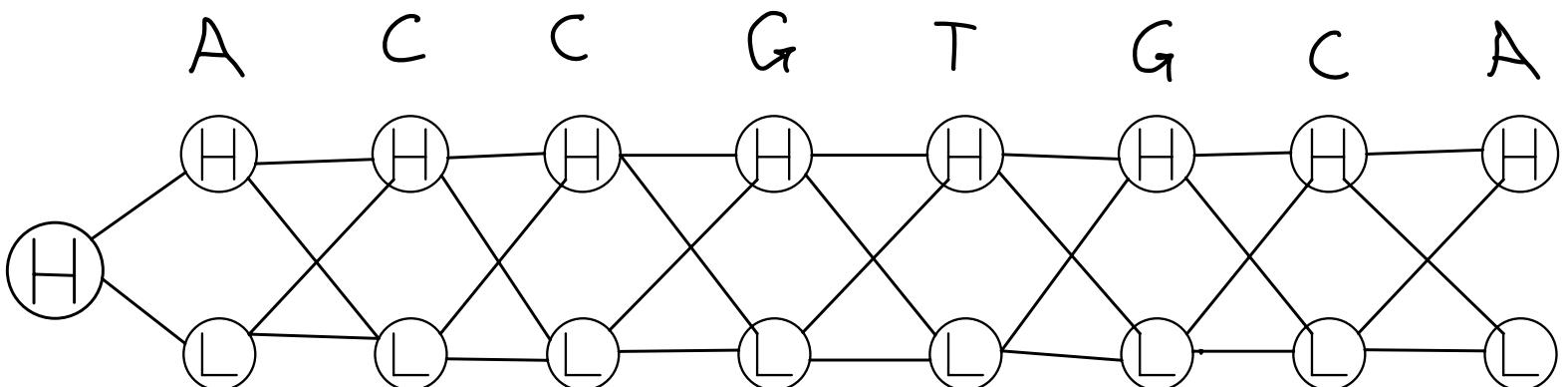
A cell can be in one of two states - H , for high GC-content, and L for low GC. On each time step the cell produces one nucleotide, A,C,T or G, and might also change its state. The probability of changing from state H to L is 0.5, and from state L to H is 0.4.

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 .

Sequence-state \rightarrow No 73 (Sequence-state).

: H ↗ State ↗



$$\therefore N = \{A, G, C, T\}, S = S_1, \dots, S_8 \quad (20)$$

$\forall \alpha, \beta \in \{H, L\} P(\alpha | \beta) : \beta \text{ state-N } \propto \text{state-L}$

$\forall \alpha \in \{H, L\}, n \in N \quad P(S_i = n | \alpha) \propto \text{exp}(-\alpha S_i)$

Given C , all α_i for $1 \leq i \leq g$ are known from the states s_1, \dots, s_g .

$\prod_{j=1}^i P(s_j | z_j) \cdot P(x_i | x_{i-1})$; States \rightarrow for Bigram HMM x_i often occurs in s_i right after s_{i-1}

new state \rightarrow if $q_i \in S_0 = n$: e unknown in S_0

States \rightarrow 8 states S_1, \dots, S_8 - $\pi(i, H) \quad 1 \leq i \leq 8$ over $S_i = S_1, \dots, S_8$: S is the set of states at time i .

$\pi(i, L) \rightarrow \alpha_i^* = L$: $\alpha_i^* \in S_i$ \rightarrow $\pi(i, L)$ is the probability of being in state $\alpha_i^*, \dots, \alpha_8^*$ at time i .

new S_i from S_{i-1} . $\alpha_i^* = L$: $\alpha_i^* \in S_i$ \rightarrow $\pi(i, L)$ is the probability of being in state $\alpha_i^*, \dots, \alpha_8^*$ at time i .

Actions $L \& H$ via α_i^* new current states \rightarrow first α_i^* state $\rightarrow \mu_i^L \& \mu_i^H$ current α_{i-1}^* , α_i^*

$\therefore \mu_i^H = \mu_i^L = H$ $\text{and } \alpha_0 = H$ initial State \rightarrow first α_0 , α_1 , α_2 , α_3 , α_4 , α_5 , α_6 , α_7 , α_8

$$\pi(1, H) = P(S_1 = A | H) \cdot P(H | H) = 0.2 \cdot (1 - 0.5) = 0.1$$

$$\pi(1, L) = P(S_1 = A | L) \cdot P(L | H) = 0.3 \cdot 0.5 = 0.15$$

for $i < 8$ and $\alpha_i^* = C$

$$\pi(i, B) = \max_{\alpha \in \{H, L\}} \{\pi(i-1, \alpha) \cdot P(S_i | B) \cdot P(B | \alpha)\} \quad \therefore \mu_i^B = \arg \max_{\alpha \in \{H, L\}} \{\pi(i-1, \alpha) \cdot P(S_i | B) \cdot P(B | \alpha)\}$$

$$\underline{\pi(2, H)} = \max \{\pi(1, H) \cdot P(S_2 = C | H) \cdot P(H | H), \pi(1, L) \cdot P(S_2 = C | H) \cdot P(H | L)\} =$$

$$= \max \{0.1 \cdot 0.3 \cdot 0.5, 0.15 \cdot 0.3 \cdot 0.4\} = \max \{0.015, 0.018\} = 0.018 \Rightarrow \underline{\mu_2^H = L}$$

$$\underline{\pi(2, L)} = \max \{\pi(1, H) \cdot P(S_2 = C | L) \cdot P(L | H), \pi(1, L) \cdot P(S_2 = C | L) \cdot P(L | L)\} =$$

$$= \max \{0.1 \cdot 0.2 \cdot 0.5, 0.15 \cdot 0.2 \cdot 0.6\} = \max \{0.01, 0.018\} = 0.018 \Rightarrow \underline{\mu_2^L = L}$$

$$\underline{\pi(3, H)} = \max \{\pi(2, H) \cdot P(S_3 = C | H) \cdot P(H | H), \pi(2, L) \cdot P(S_3 = C | L) \cdot P(H | L)\} =$$

$$= \max \{0.018 \cdot 0.3 \cdot 0.5, 0.018 \cdot 0.3 \cdot 0.4\} = \max \{0.0027, 0.00216\} = 0.0027 \Rightarrow \underline{\mu_3^H = H}$$

$$\underline{\pi(3, L)} = \max \{\pi(2, H) \cdot P(S_3 = C | L) \cdot P(L | H), \pi(2, L) \cdot P(S_3 = C | L) \cdot P(L | L)\} =$$

$$= \max \{0.018 \cdot 0.2 \cdot 0.5, 0.018 \cdot 0.2 \cdot 0.6\} = \max \{0.0018, 0.00216\} = 0.00216 \Rightarrow \underline{\mu_3^L = L}$$

$$\underline{\pi(4, H)} = \max \{\pi(3, H) \cdot P(S_4 = G | H) \cdot P(H | H), \pi(3, L) \cdot P(S_4 = G | H) \cdot P(H | L)\} =$$

$$= \max \{0.0027 \cdot 0.3 \cdot 0.5, 0.00216 \cdot 0.3 \cdot 0.4\} = \max \{0.000405, 0.0002592\} = 0.000405 \Rightarrow \underline{\mu_4^H = H}$$

$$\underline{\pi(4, L)} = \max \{\pi(3, H) \cdot P(S_4 = G | L) \cdot P(L | H), \pi(3, L) \cdot P(S_4 = G | L) \cdot P(L | L)\} =$$

$$= \max \{0.0027 \cdot 0.2 \cdot 0.5, 0.00216 \cdot 0.2 \cdot 0.6\} = \max \{0.00027, 0.0002592\} = 0.00027 \Rightarrow \underline{\mu_4^L = H}$$

$$\pi(5, H) = \max \{ \pi(4, H) \cdot P(S_5=H|H), \pi(4, L) \cdot P(S_5=H|L) \} =$$

$$= \max \{ 0.0000105, 0.0000216 \} = 0.0000105 \Rightarrow \mu_s^H = H$$

$$\pi(5, L) = \max \{ \pi(4, H) \cdot P(S_5=L|H), \pi(4, L) \cdot P(S_5=L|L) \} =$$

$$= \max \{ 0.00006075, 0.0000486 \} = 0.00006075 \Rightarrow \mu_s^L = H$$

$$\pi(6, H) = \max \{ \pi(5, H) \cdot P(S_6=G|H), \pi(5, L) \cdot P(S_6=G|L) \} =$$

$$= \max \{ 0.00006075, 0.00000729 \} = 0.00000729 \Rightarrow \mu_s^H = L$$

$$\pi(6, L) = \max \{ \pi(5, H) \cdot P(S_6=G|L), \pi(5, L) \cdot P(S_6=G|H) \} =$$

$$= \max \{ 0.00000105, 0.0000729 \} = 0.0000729 \Rightarrow \mu_s^L = L$$

$$\pi(7, H) = \max \{ \pi(6, H) \cdot P(S_7=C|H), \pi(6, L) \cdot P(S_7=C|L) \} =$$

$$= \max \{ 0.000010935, 0.000009748 \} = 0.000010935 \Rightarrow \mu_s^H = H$$

$$\pi(7, L) = \max \{ \pi(6, H) \cdot P(S_7=C|L), \pi(6, L) \cdot P(S_7=C|H) \} =$$

$$= \max \{ 0.0000010935, 0.000009748 \} = 0.000009748 \Rightarrow \mu_s^L = L$$

$$\pi(8, H) = \max \{ \pi(7, H) \cdot P(S_8=A|H), \pi(7, L) \cdot P(S_8=A|L) \} =$$

$$= \max \{ 0.0000010935, 0.0000009984 \} = 0.0000010935 \Rightarrow \mu_s^H = H$$

$$\pi(8, L) = \max \{ \pi(7, H) \cdot P(S_8=A|L), \pi(7, L) \cdot P(S_8=A|H) \} =$$

$$= \max \{ 0.0000010935, 0.0000009984 \} = 0.0000010935 \Rightarrow \mu_s^L = H$$

לפנינו $\pi(8, L) > \pi(8, H)$ ומכאן S ממליצה על $\alpha^*, \dots, \alpha_8^*$ כטורים סטטוס -> בזין וטבש נזון. מתקיים

$$\alpha_8^* = \arg \max_{\alpha \in \{H, L\}} \{ \pi(8, \alpha) \} = L$$

לפנינו $\pi(8, L) > \pi(8, H)$ ומכאן S ממליצה על $\alpha^*, \dots, \alpha_8^*$ כטורים סטטוס -> בזין וטבש נזון. מתקיים $\mu_i^{**} = \mu_i^{**}$ ו-

L, H, H, H, L, H, H, L

לפנינו S ממליצה על $\alpha^*, \dots, \alpha_8^*$ כטורים סטטוס -> בזין וטבש נזון.

$$P(AIL) \cdot P(C|H) \cdot P(C|H) \cdot P(G|H) \cdot P(T|L) \cdot P(G|H) \cdot P(C|H) \cdot P(AIL) = 0.3^8 = 0.0006561$$

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 $i = 1 \dots n$, 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, $x_1 \dots x_n$ 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 $\mathcal{K}_{-2} = \mathcal{K}_{-1} = \mathcal{K}_0 = \{\ast\}$, and $\mathcal{K}_k = \mathcal{K}$ for $k = 1 \dots n$. Define \mathcal{V} to be the set of possible words.

Initialization: ...

Algorithm: . . .

Return: ...

→ 23rd rev. $\pi(i, c, b, a) \rightarrow c \in k_{i-2} - 1$ $b \in k_{i-1}$, $a \in k_i$ for $1 \leq i \leq n$ (fix rev.)

בנוסף לכך, אם $y_{i-2} = c$, $y_{i-1} = b$ ו- $y_i = a$, אז $y_{i+1} = d$.

$$i=1 \quad \pi(1, c, b, a) = \max_{x \in V} \{ q(a | *, *, *) \cdot e(x | a) \}$$

$$i \geq 2 \quad \pi(i, c, b, a) = \max_{\substack{d \in k_{i-3} \\ x \in V}} \{ \pi(i-1, d, c, b) \cdot q(a|d, c, b) \cdot e(x|a) \}$$

לעתות בסיסי לא ניתן לחלק את המספרים α ו- β למספרים נורמליים $x \in V - 1$ ו- $y \in k_{\mathbb{F}_3}$ מאחר שקיימים מקרים שונים:

a, b, c ፩ ንዑስ ስት ወጪ ነው እና $1 \leq i \leq n$ ሲሆን $\pi(i, a, b, c)$ ጥሩ ይህን

∴ if y_{n-1}, y_n & y_{n+1} be even no then

$$y_n, y_{n-1}, y_{n-2} = \arg \max_{\substack{a \in K_n \\ b \in K_{n-1} \\ c \in K_{n-2}}} \{ \pi(n, c, b, a) \cdot q(\text{STOP} | c, b, a) \}$$

new y_1, \dots, y_{n-3} to lengths μ^{α} instead of 8.2310 and 7.81.

$$x_i = \sigma^{(i, y_{i-2}, y_{i-1}, y_i)} : e \supset x_1, \dots, x_n \rightarrow \text{new } y_{221}$$

• $\max_{y_1, \dots, y_n} P(y_1, \dots, y_n, x_1, \dots, x_n)$ 를 보면 모든 가능한 $y_1, \dots, y_n, x_1, \dots, x_n$ 에 대해서 계산, 최대값

$$q(y_1, \dots, y_n, x_1, \dots, x_n) = \prod_{i=1}^{n-1} q_i(y_i | y_{i-3}, y_{i-2}, y_{i-1}) \cdot \prod_{i=1}^n e(x_i | y_i) \quad (y_{n+1} = \text{STOP}, \quad y_0 = y_{-1} = y_{-2} = *)$$

(b) Implementation of the most likely tag baseline

- i. Using the training set, compute for each word the tag that maximizes $p(\text{tag}|\text{word})$, 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., $1 - \text{accuracy}$) for known words and for unknown words, as well as the total error rate.

לעתה כ. ווג' מקסימום ליקיילhood'estimation: ב- $\hat{\theta}$ מינימיזה פונקציית האנפיה.

ההנפשה מילויים נסבטיים. מילויים אלו מוגדרים כטבילים או מילויים נסבטיים. מילויים אלו מוגדרים כטבילים או מילויים נסבטיים.

Total Error Rate: 0.14731 סך כל שגיאות בCFG נולדה, רושם תבניות או מילוי תווים, ונאלה מיפוי של חלקים של CFG כחלק מ טרמינל או לא-טרמינל, הטעות הגדולה ביותר היא שגיאת מיפוי טרמינל לא-טרמינל.

(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).

: Bigram Hidden Markov Model ស្ថិតិមាលានៅក្នុងបញ្ជីរដូចនេះ

Known Words Error Rate: 0.13795

Unknown Words Error Rate: 0.743%

Total Error Rate: 0.20712

אך מילוי פיקטורי (מילויים נטולים) יתבצע על ידי מילויים נטולים (ללא פיקטוריים).

Sparsity -> גודל גרעין ה- ϵ מוגבל ב- $\frac{1}{\sqrt{n}}$. נסמן $\delta = \frac{1}{\sqrt{n}}$.

רְמִזְרָקָה בְּגִרְאָם, בְּגִרְאָם בְּגִרְאָם (בְּגִרְאָם, בְּגִרְאָם, בְּגִרְאָם) הַמְּלֵאָה הַגְּדוֹלָה וְהַמְּלֵאָה הַמְּלֵאָה.

ՀԱՅԱՍՏԱՆԻ ՀԱՆՐԱՊԵՏՈՒԹՅԱՆ ԿԱռավարության կողմէ ՀԱՅԱՍՏԱՆԻ ՀԱՆՐԱՊԵՏՈՒԹՅԱՆ ՀԱՆՐԱՊԵՏԱԿԱՆ ՎԵՐԱԲԵՐՅԱԼ ՀԱՆՐԱՊԵՏՈՒԹՅԱՆ ՀԱՆՐԱՊԵՏԱԿԱՆ ՎԵՐԱԲԵՐՅԱԼ

(d) Using Add-one smoothing

- Training phase: Compute the emission probabilities of a bigram HMM tagger directly on the training set using (Laplace) Add-one smoothing.
- Using the new probabilities, run the algorithm from c)ii) on the test set. Compute the error rates and compare to the results from b)ii) and c)iii).

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

Known Words Error Rate: 0.14869

:Add-One Smoothing -ב

Unknown Words Error Rate: 0.71291

Total Error Rate: 0.20871

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

$\pi_{add-1} = \frac{C(x|y) + 1}{C(x) + N}$.
הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

הנתק ב' שמיינד אוניברסיטת תל אביב מילויים של מילים ב'

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

Word Class:	detection	
Capitalized	first word is uppercase (and not all word)	Adverb ends with "able" or "ible"
Acronym	all uppercase	Adjective starts with ly
Hyphenated	contains "-"	
Gerund	ends with "ing"	
Plural	ends with s (and not ss)	kill ↗ <class word> ↗ i, y, u, n, e, n, s
Past	ends with "ed"	now ↗ a, n, d ↗ P, L, Y, N, E, R ↗ i, u, d, k, e
Numerical	all chars are digits	
Currency	contains "\$"	
Percentage	contains "%"	
Alphanumeric	contains numbers and not belong to previous options	
Unknown	others	

- 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).

Known Words Error Rate: 0.12580

: Pseudo-words -?

Unknown Words Error Rate: 0.34904

Total Error Rate: 0.1513

Psuedo-words : ? \in $\{ \text{a}, \text{b}, \text{c} \}$ \cup $\{ \text{aa}, \text{bb}, \text{cc} \}$ \cup $\{ \text{aaa}, \text{bbb}, \text{ccc} \}$ \cup $\{ \text{aaaa}, \text{bbbb}, \text{cccc} \}$ \cup $\{ \text{aaaaa}, \text{bbbbb}, \text{ccccc} \}$

בנאות, ברִמְנוֹת נעֲמָקָם מעַמְקָה ועַמְקָה נעֲמָקָם

תְּנַשֵּׁאָה מִלְּבָדָה תְּנַשֵּׁאָה מִלְּבָדָה תְּנַשֵּׁאָה מִלְּבָדָה

לעומת הטענה הנוכחית, לא מדובר במקרה אחד, אלא ב证实 (证实) של מילוי אחד במקומו של אחר.

- 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 .

:Add-One Smoothing -> Pseudo-words ->

Known Words Error Rate: 0.11084

Unknown Words Error Rate: 0.33595

Total Error Rate: 0.133595

סִגְנָה אֲלֵיכֶךָ גַּוְיִם הַנִּזְבְּחָה בְּנֵי יִשְׂרָאֵל כִּי תְּמִימָה נְקַדְּמָה בְּנֵי יִשְׂרָאֵל

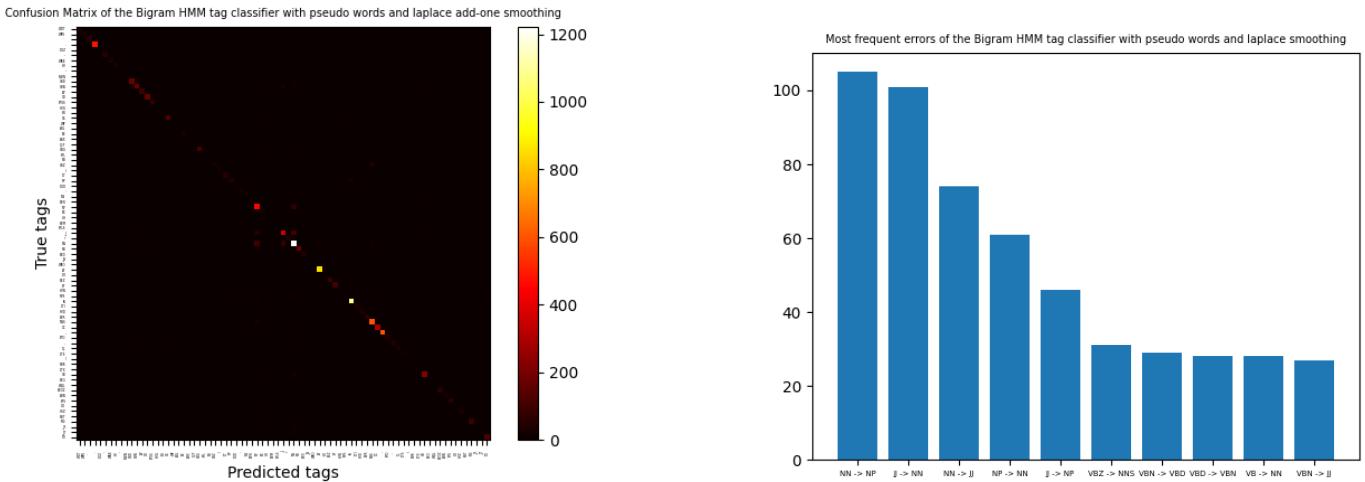
የፍቅር ወጪዎች በተመለከተ ከፍተኛ የደንብ ደንብ የሚያስፈልግ ይገልጻል

הסגרה הינה מושג של מילוי כל אחד מהלבים.

תְּנַשֵּׁא מִגְּדָלָה וְתַּחֲזֶק אֶת־מִצְבָּה־יְהוָה כִּי־בְּעַמְּךָ תְּהִלָּתָךְ

רְבָבָה, קֶלֶל, וְעַמְּדָה נִזְמָנָה רְבָבָה, קֶלֶל, וְעַמְּדָה נִזְמָנָה

$$q_{\text{add}-\delta}(x|y) = \frac{c(x|y)^{\gamma-\delta}}{c(y)^{\gamma} + \delta|x|^\gamma}$$



הטבָּה רְאֵתָה שְׁמַנִּית לְעֵנֶיךָ, יְהוָה כָּל הַיּוֹם תְּהִגֵּן עַל־עֲמָךְ וְעַל־עֲמָךְ בְּנָיו.

כִּיראָה שְׁעָרֶת יְהוָה אֵת הַמִּזְבֵּחַ וְאֵת הַמִּזְבֵּחַ וְאֵת הַמִּזְבֵּחַ

```
##### MLE Most Likely Tag Model POS Tagger Results #####
Error rate for known words: 0.07044
Error rate for unknown words: 0.74346
Total error rate: 0.14731
#####
Bigram HMM Tagger POS Tagger Results #####
Error rate for known words: 0.13795
Error rate for unknown words: 0.74346
Total error rate: 0.20712
#####
Bigram HMM Tagger with Add-One Smoothing POS Tagger Results #####
Error rate for known words: 0.14369
Error rate for unknown words: 0.71291
Total error rate: 0.20871
#####
Bigram HMM Tagger with Pseudo Words POS Tagger Results #####
Error rate for known words: 0.12580
Error rate for unknown words: 0.34904
Total error rate: 0.15130
#####
Bigram HMM Tagger with Pseudo Words and Add-One Smoothing POS Tag
Error rate for known words: 0.11084
Error rate for unknown words: 0.33595
Total error rate: 0.13655
```

Natural Language Processing - Exercise 2

Due: 11.2.2024 23:59

Please submit a single zip file. The zip file should contain your code files, a README txt file and a single pdf file for the answers to the questions

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 .

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: ...

3. (80 pts) In this programming exercise with Python, we will implement several versions of an HMM POS tagger.

Note: For all the below sections, unknown words are words that appear in the test set but did not appear in the training set.

- (a) Use the NLTK toolkit for importing the Brown corpus. This corpus contains text from 500 sources, and the sources have been categorized by genre. Here we will use a portion of the corpus: the “news” category. Load the tagged sentences for this portion of the corpus. Then, divide the obtained corpus into training set and test set such that the test set is formed by the last 10% of the sentences.
- (b) **Implementation of the most likely tag baseline**
 - i. Using the training set, compute for each word the tag that maximizes $p(\text{tag}|\text{word})$, 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., $1 - \text{accuracy}$) for known words and for unknown words, as well as the total error rate.
- (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).
- (d) **Using Add-one smoothing**
 - i. Training phase: Compute the emission probabilities of a bigram HMM tagger directly on the training set using (Laplace) Add-one smoothing.
 - ii. Using the new probabilities, run the algorithm from c)ii) on the test set. Compute the error rates and compare to the results from b)ii) and c)iii).
- (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 .

Note: NLTK is a Python package for NLP. Please see here <http://www.nltk.org/> for downloading and documentation. You can download the Brown corpus using `nltk.download('brown')` and traverse its sentences using the `tagged_sents()` method. Please refer to the NLTK documentation for further information. Please use NLTK just for loading and reading the corpus.

A given word w may be annotated with a complex tag t , containing the symbols '+' and/or '-'. When encountering such a complex tag t , consider only the prefix of t that comes before the first occurrence of '+' or '-' in t as the POS tag of w .