

# Natural Language Processing - Exercise 3

Due: 22.12.2024 23:59

Please submit **two** files – a pdf file for the answers to the questions and a zip archive. The zip file should contain your code files and a README txt. Please do not put the pdf in the zip.

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.

# 3. Hidden Markov Model - NLP

## קורס פונקציונלי ואינטראקציה

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. נשאלת מהי הסדרה הטובה ביותר של מצבים  $H$  ו- $L$  עבור הסדרה  $S = ACCGTGCA$ .

הנני מניח  $n = 8$  ו- $k = 0, 1, \dots, 7$

הסדרה  $S_k$  היא  $S_k = (s_1, s_2, \dots, s_k)$  כאשר  $s_i \in \{A, C, G, T\}$

המרחב  $\mathcal{S} = \{H, L\}^k$  הוא מרחב המצבים

הפונקציה  $r$  היא:

$$r(y_1, \dots, y_k) = \prod_{i=1}^k q(y_i | y_{i-1}) \cdot \prod_{i=1}^k e(x_i | y_i)$$

הפונקציה  $e$  היא פונקציית הפליטה (emission) ו- $q$  היא פונקציית המעבר (transition).

הפונקציה  $r$  היא פונקציית הסתברות של סדרת המצבים  $y_1, \dots, y_k$  בהינתן הסדרה  $x_1, \dots, x_k$ .

$$\text{angle} \rightarrow \pi, \quad \pi(k, t) \rightarrow \int_{-10}^{10} \rightarrow \text{say } L \text{ in}$$
$$\pi(k, t) = \max_{\langle y_0, y_1, \dots, y_k \rangle: y_0 = H, y_k = t} (r(y_0, \dots, y_k))$$

A . C . C . G . T . G . C . A

$\Gamma \cdot L \cdot$  קושיק  $\cdot$  נשמות  $\cdot$  פאליגור  $\cdot$  מן  $\cdot$  וואלן  $\cdot$  כן  $\cdot$  גיל  $\cdot$   $S_8 \rightarrow$

$\frac{1064025}{10^7} \rightarrow 10^7, 10^7 \cdot S_7 = (H, L, H, H, H, L, H, H, L) \rightarrow 10^7, 10^7 \cdot P_3 \rightarrow 10^7$

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:** ...

(2)  $n \times |\mathcal{V}| \times |\mathcal{K}|$   $\int_{\mathcal{V}} \int_{\mathcal{K}} \dots$

$\pi(i, y, x)$   $\rightarrow$   $\pi(i, y, x)$

$\pi(i, y, x)$   $\rightarrow$   $\pi(i, y, x)$

$\pi(0, y, x) = 1$   $\rightarrow$   $\pi(0, y, x) = 1$

$y \in \mathcal{K}, x \in \mathcal{V}$

$\pi(i, y, x)$

$V_{-3} = V_{-2} = V_{-1} = V_0 = \{*\}, V_i = \mathcal{V}$  for  $i \in \{1, \dots, n\}$

$K_{-3} = K_{-2} = K_{-1} = K_0 = \{*\}$ ,  $K_i = \mathcal{K}$  for  $i \in \{1, \dots, n\}$   
 empty tag for \*

אלגוריתם:

for  $i = 1, \dots, n$

for  $y \in K$

for  $x \in V$

$$\pi(i, y, x) = \max_{\substack{y' \in K_i, y'' \in K_{i-1}, \\ y' \in K_{i-2}, y'' \in K_{i-1}, \\ x \in V_i}} \pi(i-1, y, x) \cdot q(y | y'', y', y) \cdot e(x_k | y)$$

Return  $\max_{x \in V_n, y \in K_n} (\pi(n, x, y))$

הסבר קצר: בשל האלימות יתנו מילויים של האלמנטים

נקרא אלמנטים של מילויים של "מילוי" זה ייתכן שיש מילויים

"מילויים" היתכונים ויש להם אלמנטים של מילויים של מילויים של

ערכי x והמילויים של (צדק) איזוהי של האלמנטים של מילויים של

הם והיתכונים הם של מילויים של מילויים של

אלגוריתם רצוני:

