

67658) עיבוד שפה טבעית | תרגיל 2

שם: דניאל בראון 311340723, נעה וייס 205676638

16 בדצמבר 2020

שאלה 1

יהי מודל שמקיים את שתי הנחות השאלה. מהנחה מס' 1 נובע שעבור שרשרת תגיות y_1, \dots, y_n מתקיים:

$$\mathbb{P}(y_1, \dots, y_n) = \prod_{i=1}^n P(y_i | y_{i-m}, y_{i-m+1}, \dots, y_{i-1})$$

נשים לב שזוהי נוסחה שמייצגת את התגיות כשרשרת מרקובית מסדר m , כלומר ההסתברות לכל תגית אינה תלויה בתגיות האחרות בשרשרת, בהנתן n התגיות שקדמו לה. מהנחה מס' 2 נובע שעבור משפט $x = x_1, \dots, x_n$ מתקיים:

$$\mathbb{P}(x_1, \dots, x_n) = \prod_{i=1}^n e(x_i | y_i)$$

כלומר, x_i בלתי תלוי ב- x_j, y_j אחרים בהינתן y_i , ולכן ההסתברות ל- x שווה למכפלת ההסתברויות ה- $emissions$ של x_i בהינתן y_i .

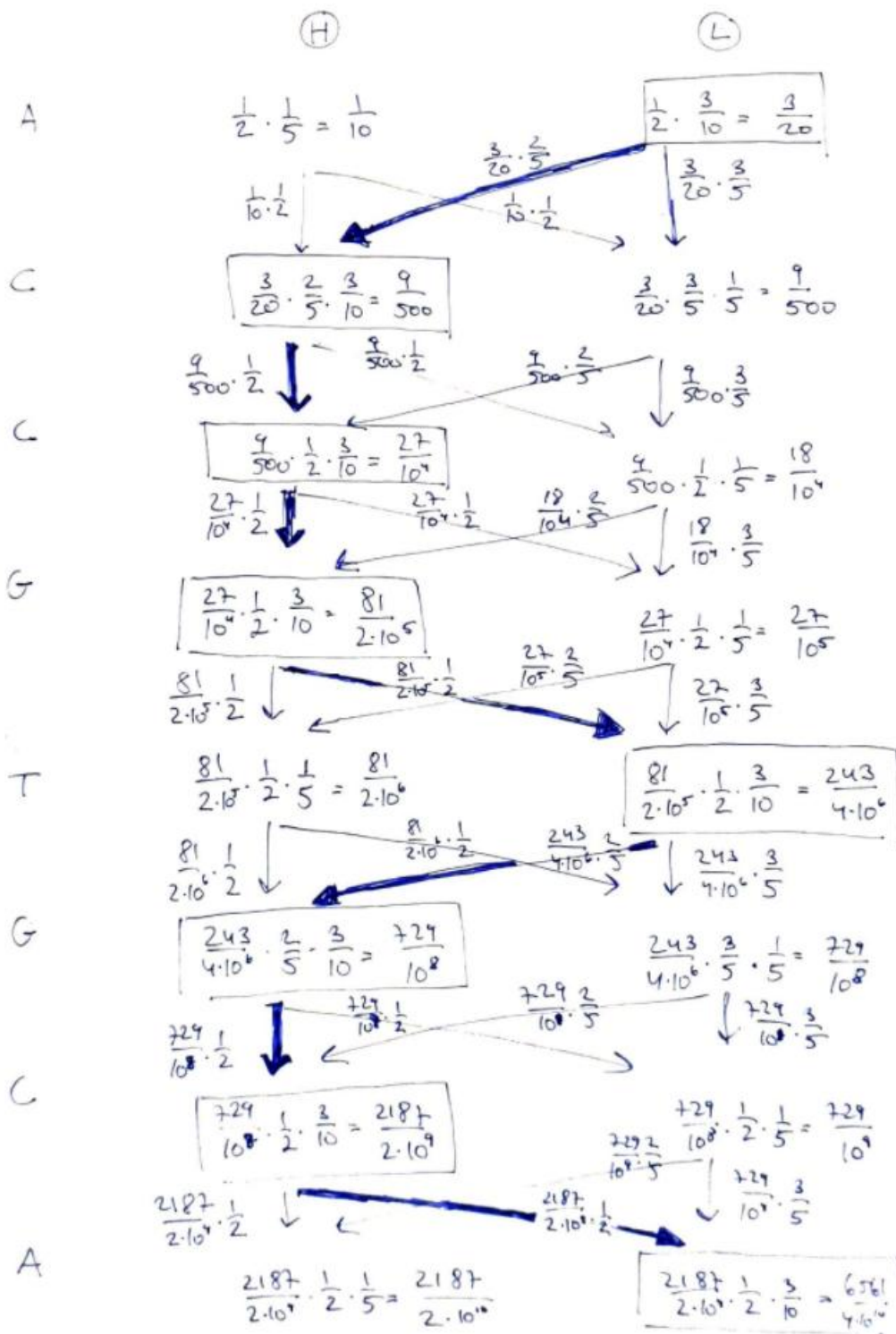
קעת, היות שההסתברות ל- y אינה תלויה כלל ב- x , וההסתברות של x_i בהינתן y_i אינה תלויה באף גורם אחר, ההסתברות המשותפת של x ו- y , כלומר ההסתברות המשותפת של משפט ושל שרשרת התגיות שתמקסם את ההסתברות למשפט זה, תהיה מכפלת ההסתברויות:

$$\begin{aligned} \mathbb{P}(x_1, \dots, x_n, y_1, \dots, y_n) &= \prod_{i=1}^n P(y_i | y_{i-m}, y_{i-m+1}, \dots, y_{i-1}) \cdot \prod_{i=1}^n e(x_i | y_i) = \\ &= \prod_{i=1}^n P(y_i | y_{i-m}, y_{i-m+1}, \dots, y_{i-1}) \cdot e(x_i | y_i) \end{aligned}$$

קיבלנו הכללה של נוסחת מודל ה- HMM עבור שרשרת תגיות מרקובית מסדר m , ולכן המודל הוא מסוג HMM , כנדרש. ■

שאלה 2

נריץ את אלגוריתם ויטרבי על הקלט הנתון:



רצף המצבים בעל ההסתברות הגבוהה ביותר לייצר את S הוא LHHHLHHL.

ההסתברות ל- S בהינתן רצף מצבים זה היא $1.64025 \cdot 10^{-7} = \frac{.6561}{4 \cdot 10^{10}}$.

שאלה 3

נשלים את ריצת האלגוריתם באופן הבא:

Initialization : Set $\pi(0, *, *, *) = 1$

Algorithm : For $k = 1, \dots, n$:

For $u \in K_{k-2}, v \in K_{k-1}, w \in K_k$:

$$\pi(k, u, v, w) = \max_{t \in K_{k-3}} (\pi(k-1, t, u, v) \times q(w|t, u, v)) \times e\left(\operatorname{argmax}_{x_k \in V} (x_k|w) | w\right)$$

Return : $\max_{u \in K_{n-2}, v \in K_{n-1}, w \in K_n} (\pi(n, u, v, w) \times q(STOP|u, v, w))$

הערכים עבור הפונקציות q, e , כלומר פונקציות ה-*transition* וה-*emission*, מחושבות ונשמרות בטבלאות לפני ריצת האלגוריתם. היות שאנו מחשבים מודל *four-gram*, הטבלה עבור הפונקציה q תהיה בגודל $|K|^3 \times |K|$. הטבלה עבור e תהיה בגודל $|K| \times |V|$, משום שאנו צריכים למצוא את המילה שבהינתן כל $w \in K$ תתן את ההסתברות הגבוהה ביותר עבור אותו w . לכן עלינו לחשב את ההסתברויות המותנות $p(x_i|w_j)$ לכל $i \in [|V|]$ ולכל $j \in [|K|]$. שמירת שתי הטבלאות הללו לפני ריצת האלגוריתם תשפר את זמן הריצה, היות שמספיק לחשבן פעם אחת. לאחר שמירת הטבלאות נוכל למצוא את q, e בסיבוכיות $O(1)$, ולכן זמן ריצת האלגוריתם יהיה $O(n \cdot |K|^3)$.

שאלה 4

טבלת תוצאות הריצות השונות:

	known words	unknown words	total
<i>MLE</i>	0.1079	0.7897	0.1857
<i>base Viterbi</i>	0.0992	0.4598	0.1772
<i>Viterbi with add - 1</i>	0.1576	0.6464	0.2634
<i>Viterbi with pseudo words</i>	0.1114	0.7219	0.2435
<i>Viterbi with pseudo words and add - 1</i>	0.1730	0.7403	0.2958

השגיאות הנפוצות ביותר לפי ה-*confusion matrix*:

NN במקום *NN* - 210 שגיאות

NP במקום *NN* - 95 שגיאות

JJ במקום *AT* - 50 שגיאות

NN במקום *AT* - 46 שגיאות

JJ במקום NN - 45 שגיאות

NN במקום NP - 43 שגיאות

המטריצה כולה מודפסת בתור מילון בקובץ ה- py . המצורף.