

תרגיל 1 קורס: עיבוד שפה טבעית משימה: למידה כהן, נריה בן דוד חלק תיאורטי

1. א. נגדיר  $W = \{w_1, \dots, w_m\}$  אוצר המילים  $start, stop$ .

$$\text{לפני: } \forall 1 \leq i \leq m \sum_{j=1}^m P(w_j | w_i) + P(stop | w_i) = 1$$

$$\forall 1 \leq i \leq m \exists \epsilon_i > 0 \quad P(stop | w_i) > \epsilon_i$$

נגדיר  $X_k$  משתנה מקרי המייצג את המילה ה- $k$  במשפט.

$$A = \{ \exists k > 0 \quad X_k = stop \} \quad \text{ד"ר. } P(A) = 1$$

$$A^c = \{ \forall k > 0 \quad X_k \neq stop \} \quad P(A^c) = 0$$

נחק שנומסר סופי על מילים  $\epsilon = \min \{ \epsilon_i \mid 1 \leq i \leq m \}$  כך שכל  $1 \leq i \leq m$

$$P(stop | w_i) > \epsilon > 0$$

$$P(A^c) = \prod_{k=1}^{\infty} P(X_k \neq stop \mid X_{k-1} \neq stop) \leq \prod_{i=1}^{\infty} (1 - \epsilon) \leq \prod_{i=1}^{\infty} e^{-\epsilon} = e^{-\sum_{i=1}^{\infty} \epsilon} = 0$$

נוסח המטרה  
ה- bigram

הנחה על הפסגה המקב

ב. דא, נראה דוגל נגדית.

נבנה מודל שבו ההסתברות לעבור ל- $stop$  תלויה באורך המשפט עד

רבע המעבר.

$$P(X_k = stop \mid X_1, \dots, X_{k-1} \in W) = e^{-\frac{1}{2^k}} \quad \text{כלומר:}$$

$$P(A^c) = \prod_{k=1}^{\infty} P(X_k \neq stop \mid X_1, \dots, X_{k-1} \in W) = \prod_{k=1}^{\infty} e^{-\frac{1}{2^k}} = e^{-\sum_{k=1}^{\infty} \frac{1}{2^k}} = e^{-1} > 0$$

דפ  $P(A) < 1$

כלומר קיים מילה אינסופית שאינה מסיימת ל- $stop$ .

ג. א. נגזיר מודל unigram כאופן הבא:

בהינתן מפתח המודל מחזיר את ההסתברות של המפתח

על פי כמות ההופעות של המילים במט אימון.

לוח: 
$$P(w_i) = \frac{C(w_i)}{\sum_{w \in W} C(w)}$$
 עבור כל מילה.

ועבור כל מפתח: 
$$P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i)$$

בהינתן המפתח המבוקש,

כמאורע הכאן שני נהיה צודקם במקרה בו  $C(w_{\text{where}}) > C(w_{\text{were}})$

כמאורע השני שני נהיה צודקם במקרה בו  $C(w_{\text{where}}) < C(w_{\text{were}})$

לעומת זאת לא נצדוק בשני המאורעות מכיוון שהתנאים

$C(w_{\text{where}}) > C(w_{\text{were}})$ ,  $C(w_{\text{where}}) < C(w_{\text{were}})$  אינם יכולים להתקיים במקביל.

ב. נגזיר מודל bigram כאופן הבא:

בהינתן מפתח המודל מחזיר את ההסתברות של המפתח

על פי כמות ההופעות של המילים במט אימון תוך תלות במילה

הקודמת להן. לוח: 
$$P(w_i | w_{i-1}) = \frac{C(w_i, w_{i-1})}{\sum_w C(w_i, w_{i-1})}$$
 עבור כל מילה.

ועבור כל מפתח: 
$$P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{i-1})$$

המודל יכול להציג תוצאה יותר טובת מהמודל בסעיף א'

משום שבניגוד למודל הקודם, מודל זה מתחשב גם בהקשר.

ואכן יכול להבדיל בין מקרים בהם דקדוקיות נכונה יותר על

where לעומר were.

. במורז זה משפטים יכולים לקבל הסתברות 0, כיוון

שלמחר שכל המילים מופיעות לפחות פעם אחת.

"יתכן כי צמחים לא מילים בשרת לא יופיעו כלל בסט אימון

ומכאן כאשר נחשב את הסתברות המשפט הצמח לא יופיעו

ולכן נקבל הסתברות אפס.

למשל: סט אימון: the dog ate the ומכאן כי  $P(\text{the ate}) = 0$ .

. זו יכולה להיות בעיה כיוון שמשפט שאחרי לקבל הסתברות

גבוהה מהמורז יכול להתאים כיוון שהיו בו שני מילים סמוכות שאינן

הופיעו, ומכאן הישג ה-bigram איבד את כל הסתברות המשפט.

3. כ.

$$A_c = \{w \mid \text{בקורסי } c \text{ מופיעה } w \text{ } \} \text{ עדיף } \text{מ'נס}$$

$$* \sum_{c=1}^{C_{\max}} \sum_{w \in A_c} \frac{(c+1)N_{c+1}}{N_c \cdot N} = \frac{1}{N} \sum_{c=1}^{C_{\max}} \sum_{w \in A_c} \frac{(c+1)N_{c+1}}{N_c} = \frac{1}{N} \sum_{c=1}^{C_{\max}} N_c \cdot \frac{(c+1) \cdot N_{c+1}}{N_c} = \frac{1}{N} \sum_{c=1}^{C_{\max}} (c+1) \cdot N_{c+1} =$$

$$\frac{1}{N} (2N_2 + 3N_3 + \dots + C_{\max} N_{C_{\max}} + 0) = \frac{1}{N} (N - N_1) = 1 - \frac{N_1}{N} = 1 - P_{\text{unseen}}$$

\*  $\int_m$   $N - C_{\max} + 1$  מופיע נאליס באיכנס עליו.

$$. q_{\text{add}-1}(w) = \frac{c+1}{\sum_{w'} (c(w')+1)} = \frac{c+1}{N+|V|} \quad , \quad q_{ML}(w) = \frac{c}{N} \quad .2$$

$\downarrow$   
-V קורסי מופיעים

$$\mu = \frac{N}{|V|} \quad \text{גדל}$$

$$* q_{\text{add}-1}(w) = \frac{c+1}{N+|V|} \leq \frac{c}{N} = q_{ML}(w) \Leftrightarrow (c+1)N \leq c(N+|V|) \Leftrightarrow$$

$$\Leftrightarrow cN + N \leq cN + |V|c \Leftrightarrow N \leq |V| \cdot c \Leftrightarrow \mu := \frac{N}{|V|} \leq c$$

$$q_{\text{add}-1}(w) \leq q_{ML} : c \geq \mu \quad \text{גדל}$$

$$* q_{\text{add}-1}(w) = \frac{c+1}{N+|V|} \geq \frac{c}{N} = q_{ML}(w) \Leftrightarrow (c+1)N \geq c(N+|V|) \Leftrightarrow$$

$$\Leftrightarrow cN + N \geq cN + |V|c \Leftrightarrow N \geq |V| \cdot c \Leftrightarrow \mu := \frac{N}{|V|} \geq c$$

$$q_{\text{add}-1}(w) \geq q_{ML} : c \leq \mu \quad \text{גדל}$$

ד.

$$q_{GT}(w) \leq q_{ML}(w) \Leftrightarrow \frac{(C+1)N_{C+1}}{N_C \cdot N} \leq \frac{C}{N} \Leftrightarrow (C+1)N \cdot N_{C+1} \leq C \cdot N_C \cdot N \Leftrightarrow$$

$$(C+1)N_{C+1} \leq C \cdot N_C \Leftrightarrow CN_{C+1} + N_{C+1} \leq C \cdot N_C \Leftrightarrow N_{C+1} \leq CN_C - CN_{C+1} \Leftrightarrow$$

$$N_C - N_{C+1} = 0 \quad \text{אם} \cdot \quad \frac{N_{C+1}}{N_C - N_{C+1}} \leq C \quad N_C - N_{C+1} > 0 \quad \text{אם} \cdot$$

$$N_{C+1} \leq 0$$

$$\frac{N_{C+1}}{N_C - N_{C+1}} \geq C \quad N_C - N_{C+1} < 0 \quad \text{אם} \cdot$$

מכאן כי לא קיים threshold מ קבוע עבור מתקנים אי השינויים  
הרצויים. (כ)  $N_C$  יכול להגיע שונה עבור  $C$  ולכן לא יהיה מסת קבוע

4. ג. תוססה עבור מודל trigram:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | x_{i-1}, x_{i-2})$$

ההנחה במודל זה על התלג בין מילים היא של תלויה אך ורק בשני

המילים המופיעות לפנייה.

ג. הם אוכלים ארוחה.

במקרה זה המודל יתפשט בצורת הפוסל אוכלים ברכה

כיוון שלפניו באה המילה הם, ומודל trigram מתפשט  
השתי המילים האחרונות ובמקרה זה כפרט במילה האחרונה.

She always sings.

במקרה זה המודל יתפשט בהטיה התכונה של sings

משל לפניו באות המילים she, always

ומודל trigram מתפשט ב-2 המילים האחרונות.

ד. The book that I bought yesterday is very interesting.

מודל שמבחן בין שמוש בין is - are בצורה

נכונה במשפט צריך להתפשט במילה אסבט שזו המילה המופיעה  
5 מילים לפני המילה is.

מודל n-gram שיבחן בכך יהיה על הפח 6-gram.

קנינו ספרים כדי שיוכל לקרוא אותם.

מודל שמבחן בצורה אחרת צריך להתפשט במילה ספרים.

המילה ספרים מופיעה 4 מילים לפני מילה אדם אכן מודל n-gram

מתפשט בכך יהיה על הפח 5-gram.

5. היא פתחה את הדלת שדולה

א שתי מילים עוקבת נטות תחבירית של המשפט סיני יתכן  
תחבירית.

משכן שדור שדול נא

א ששה מילים עוקבת נטות תחבירית של המשפט סיני יתכן  
תחבירית.

שתי ומילה אחת כבוקר הולכת לבית-ספר.

א 4 מילים עוקבת נטות תחבירית של המשפט סיני יתכן  
תחבירית.

מכאן שנת לאט א לא מת שמוציא מרקוב יהיו  
אוכה כחוצי שבה בהינתן משפט באורך  $n$  נרצה להשיג  
כחוצי  $n$ -gram, אומר בשל שמוציא מרקוב יהיו מבויקים  
יה צריכים להת דעל סיכוט שבוהה.

the 2

.3 2 -29.666 :2 62N .- ∞ :1 62N .b

.2 .perplexity=∞

.4 2 -36.19 :1 62N -30.99 :2 62N .b

.2 270.076