

# תרגיל 1 - NLP

27 בינואר 2024

Avraham Asraf : 315774570

Schnaidman Elchanan: 316092436

## חלק 1

סעיף a

לפי ההדרכה, נוכיח כי המשלים, כלומר הסיכוי למשפטים שלא נגמרים הוא 0.  
נסמן

$$\Omega = \{s = (w_1, \dots, w_n) \mid n \in \mathbb{N} \cup \infty\}$$

מתקיים

$$\Omega = \Omega_{\mathbb{N}} \cup \Omega_{\infty}$$

נוכיח כי

$$\mathbb{P}(\Omega_{\infty}) = 0$$

מתקיים

$$\mathbb{P}(\Omega_{\infty}) = \sum_{s \in \Omega_{\infty}} \mathbb{P}(s)$$

כיון שמהאורעות זרים,

קיים  $\varepsilon > 0$  המקיים

$$\varepsilon = \min_{w \in W} \{\mathbb{P}(\text{stop}|w)\}$$

אזי לכל  $w_i, w_j \in W$  מתקיים

$$\mathbb{P}(w_i|w_1 \dots, w_j) \leq 1 - \varepsilon$$

$\mathbb{P}(\text{stop}|w_1 \dots, w_j) = \mathbb{P}(\text{stop}|w_j) \geq \varepsilon$  כיון שהסיכוי ל

א"כ לכל  $s \in \Omega_\infty$  מתקיים

$$\mathbb{P}(s) \leq \prod_{i=1}^{\infty} 1 - \varepsilon = \lim_{n \rightarrow \infty} (1 - \varepsilon)^n = 0$$

לכן

$$\mathbb{P}(\Omega_\infty) = \sum_{s \in \Omega_\infty} \mathbb{P}(s) = \sum_{s \in \Omega_\infty} 0 = 0$$

כנדרש.

**b1.**

במודל לא מרקובי המקיים את ההנחות לעיל יתכן שיהיה משפט אינסופי עם הסתברות גדולה מ-0.

דוגמה:

יהי  $P$  מודל שפה על אוצר המילים  $(\text{start}, w, \text{end})$  מוגדר בצורה הבאה:

$$P(\text{end}|s) = \begin{cases} 0 & s = (\text{start}) \\ \frac{1}{2^n} & s = (\text{start}, w, \dots, w_n) \end{cases}, P(w|s) = \begin{cases} 1 & s = (\text{start}) \\ 1 - \frac{1}{2^n} & s = (\text{start}, w, \dots, w_n) \end{cases}$$

זהו מודל מוגדר היטב:

$$\sum_{n \rightarrow \infty} P(s_n) = \frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^n} = \frac{\frac{1}{2}}{1 - \frac{1}{2}} = 1$$

כאשר  $s_n$  מסמן את מספר ה- $w$  במילה.

בכל רצף מילים קיים סיכוי חיובי להגיע ל- $\text{end}$

קיים סיכוי חיובי לטור אינסופי

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{2^n}\right)^n > \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = \frac{1}{e} > 0$$

כנדרש.

## a2.

במודל unigram קיימות שלוש אפשרויות :

$\mathbb{P}(\text{where}) = \mathbb{P}(\text{were})$  או  $\mathbb{P}(\text{where}) < \mathbb{P}(\text{were})$  או  $\mathbb{P}(\text{where}) > \mathbb{P}(\text{were})$  ואז ההחלטה בן  $\text{where}$  ל  $\text{were}$  תהיה לא דטרמיניסטית.

אם  $\mathbb{P}(\text{where}) > \mathbb{P}(\text{were})$  אזי מתקן השגיאות לא יחליף את ההופעה של  $\text{were}$  במקום הלא נכון.

אם  $\mathbb{P}(\text{where}) < \mathbb{P}(\text{were})$  מתקן השגיאות יחליף את שתי ההופעות של  $\text{where}$ . במקרה בו הסיכוי שווה לא תהיה החלטה קבועה ובכל מקרה לא נקבל פתרון נכון בוודאות.

## b2.

אם נשתמש במודל bigram

אם  $\mathbb{P}(\text{where}|\text{went}) > \mathbb{P}(\text{were}|\text{went})$  וגם  $\mathbb{P}(\text{where}|\text{there}) < \mathbb{P}(\text{were}|\text{there})$

כלומר יש יותר מקרים כאלה בקורפוס, אזי המודל יתקן נכון.

יתכן שמודל כזה יתן סיכוי 0 לצמד מילים שלא נמצאים בקורפוס, גם אם הוא הגיוני.

כמו כן מילה אחת לפני לא מספיקה בשביל לדעת האם המשפט נכון או לא. ויתכן שצימדי מילים שנמצאות הרבה יקבלו סיכוי גבוה

ויחזו כנכונים גם במקרים לא הגיוניים.

## 3. נענה לפי הסעיפים

א. עבור מילה שהופיע  $c$  פעמים, יש  $N_c$  מילים כאלה בקורפוס, לכן סכום התדירויות של סוג

$$\text{מילה כזה (כל המילים) הוא } \frac{(c+1) \cdot N_{c+1}}{N} \cdot N_c = \frac{(c+1) \cdot N_{c+1}}{N_c \cdot N} \cdot N_c$$

כדי לקבל את סכום התדירויות של כלל המילים בקורפוס, נסכום על כל האפשרויות של  $c$ .

$$\begin{aligned} \sum_{c \in [c_{max}]} \frac{(c+1) \cdot N_{c+1}}{N} &= \frac{1}{N} * \sum_{c \in [c_{max}]} \frac{(c+1) \cdot N_{c+1}}{N} \stackrel{*}{=} \\ &= \frac{1}{N} * \sum_{c \in [c_{max}-1]} \frac{(c+1) \cdot N_{c+1}}{N} \stackrel{**}{=} \frac{1}{N} * \sum_{c \in [2, c_{max}]} \frac{c \cdot N_c}{N} \stackrel{***}{=} \\ &= \frac{1}{N} * (N - 1 * N_1) = \frac{N}{N} - \frac{N_1}{N} = 1 - p_{unseen} \end{aligned}$$

מעבר \*: מכיוון ש  $N_{c+1} = 0$  ניתן להשמיט את האיבר האחרון

מעבר \*\*: במקום להתחיל לספור מ-1 ולסיים ב-  $c_{max} - 1$ , מתחילים ומסיים במספר העוקב ומפחיתים איד בכל מופע של  $c$ .

מעבר \*\*\*: מתבצעת שם סכימה של כל המילים שמופיעים לפחות פעמיים, חוץ מהמילים שמופיעות פעם אחת.

ב. הנוסחה ל-Add-One הינה  $q_{Add-One} = \frac{c+1}{N+|V|}$ . הנוסחה ל-MLE הינה  $q_{MLE} = \frac{c}{N}$ .

נראה מתי מתקיים  $q_{MLE} > q_{Add-One}$ .

$$\begin{aligned} q_{MLE} > q_{Add-One} &\leftrightarrow \frac{c}{N} > \frac{c+1}{N+|V|} \leftrightarrow cN + c|V| > cN + N \leftrightarrow c|V| > N \leftrightarrow \frac{c}{N} > \frac{1}{|V|} \\ &\leftrightarrow q_{MLE} > \frac{1}{|V|} \end{aligned}$$

כלומר, קיבלנו שכאשר  $q_{MLE} > \frac{1}{|V|}$  אזי מתקיים  $q_{MLE} > q_{Add-One}$ .

נראה מתי מתקיים  $q_{MLE} < q_{Add-One}$ .

$$q_{MLE} < q_{Add-One} \leftrightarrow \frac{c}{N} < \frac{c+1}{N+|V|} \leftrightarrow cN + c|V| < cN + N \leftrightarrow c|V| < N \leftrightarrow \frac{c}{N} > \frac{1}{|V|}$$

$$\leftrightarrow q_{MLE} < \frac{1}{|V|}$$

כלומר, קיבלנו שכאשר  $q_{MLE} < \frac{1}{|V|}$  אזי מתקיים  $q_{MLE} < q_{Add-One}$ .

מחיבור שני החלקים, מצאנו שיש  $threshold$   $\mu = \frac{1}{|V|}$  שמקיים את מה שהיה צריך להוכיח.

ג. הנוסחה ל-GTE הינה  $q_{GTE} = \frac{(c+1)*N_{c+1}}{N_c*N}$ . הנוסחה ל-MLE הינה  $q_{MLE} = \frac{c}{N}$ .

נראה מתי מתקיים  $q_{MLE} > q_{GTE}$ .

$$q_{MLE} > q_{GTE} \leftrightarrow \frac{c}{N} > \frac{(c+1)*N_{c+1}}{N_c*N} \leftrightarrow \frac{c}{c+1} > \frac{N_{c+1}}{N_c}$$

נניח שקיים  $threshold$   $\mu$  שמקיים שאם  $q_{MLE} > \mu$  אזי המשפט מתקיים.

כלומר,  $q_{MLE} > \mu \leftrightarrow \frac{c}{N} > \mu \leftrightarrow c > N * \mu$ .

נראה כי ניתן לבחור  $c$  כרצוננו בתחום הנ"ל שפעם האי שוויון  $\frac{c}{c+1} > \frac{N_{c+1}}{N_c}$  יתקיים ופעם לא.  
(בעקבותיו גם  $q_{MLE} > q_{GTE}$ )

$c$	$c+1$	$\frac{c}{c+1}$	$N_{c+1}$	$N_c$	$\frac{N_{c+1}}{N_c}$	האי שוויון מתקיים?
$N * \mu + 1$	$N * \mu + 2$	$\frac{N * \mu + 1}{N * \mu + 2}$	$N * \mu + 1$	$N * \mu + 10$	$\frac{N * \mu + 1}{N * \mu + 10}$	כן
$N * \mu + 7$	$N * \mu + 8$	$\frac{N * \mu + 7}{N * \mu + 8}$	$N * \mu$	$N * \mu$	$\frac{N * \mu}{N * \mu} = 1$	לא

כלומר, הראינו כי לא קיים  $threshold$  שבו בהכרח יתקיים תמיד  $q_{MLE} > q_{GTE}$ .

4. לפי סעיפים

א. הנוסחה למודל  $trigram$  הינה:  $p(w_1, w_2, \dots, w_n) = \prod_{i=1}^n p(w_i | w_{i-2}, w_{i-1})$ .

ההנחה שההסתברות להופעת כל מילה מושפעת רק משתי המילים הקודמות ולא יותר.

ב. משפט בעברית: הכלב אוהב לרוץ.

משפט באנגלית: *The dog likes to run*.

הפועל (*likes*) הוא צמוד למילת הנושא (*dog*), בגלל שהמודל  $trigram$  מניח שמילה מושפעת משתי המילים הקודמות, וכאן הנושא הוא פחות משתי מילים לפני הפועל, לכן כנראה המודל ינבא היטב את ההטיה הנכונה.

ג. משפט בעברית: הכלב עם השיניים הגדולות אוהב לרוץ.  
משפט באנגלית: *The dog with the big teeth likes to run*.

הפועל (*likes*) הוא לא צמוד מספיק למילת הנושא (*dog*), בגלל שהמודל trigram מניח שמילה מושפעת **משתי** המילים הקודמות, וכאן הנושא הוא יותר משתי מילים לפני הפועל, לכן כנראה המודל לא ינבא היטב את ההטיה הנכונה.

בגלל שבמשפט שבחרנו הנושא הוא 5 מילים לפני הפועל, נדרש מודל 6-gram כדי לחזות כהלכה את ההטיה. (במשפט בעברית זה יהיה מודל 5-gram).

5. עבור צמדים: **הילד** אכל ואז הלכה לבית הספר  
עבור שלישיות: **הילד** אכל ולאחר מכן הלכה לבית הספר  
עבור רביעיות: **הילד** עם העיניים הכחולות הלכה לבית הספר

בצמדים – כל צמד הוא חלק ממשפט תקף ורק המשפט המלא לא נכון תחבירית כי יש הטיה של נקבה עבור הילד (זכר).  
אותו הדבר עבור שלישיות ורביעיות.

ככל שהמודל המרקובי מניח תלות ביותר מילים קודמות, ככה נלקח בחשבון גם יותר מידע הקשרי. לכן, המודל יהיה יותר מדויק ויתפוס יותר הקשרים בתוך המשפט, כאשר הוא ייקח בחשבון יותר מילים לפני המילה הנוכחית. לעשות משפט שהוא לא תקף בכללותו אבל במקטעיו הוא כן, ככל המודל הוא בסדר יותר גבוה, ככה יותר קשה כי ההקשרים במשפט נוטים להיות יחסית סמוכים.