

שאלה 1

(d) בהינתן מודל שפה bigram אלמנט $S: \text{START}, w_1, w_2, \dots, w_n, \text{STOP}$
 אנו רוצים כי שם הסגנון המצויין מוצגת היכן (נסמך 1)
 אלא ההסתברות קימה הסגנון שזה מאבד להצטרף STOP
 אזי סוג הסגנון מכלל הרכבים הסוכים הוא 1.

הוכחה:

$$\forall w \quad P(\text{STOP} | w) > 0$$

$$\forall w \quad P(\text{STOP} | w) \geq p \quad \text{ע"פ} \quad \exists p > 0$$

$$A_i = \{ \text{המשטר } i \text{ במילה } i \} = \{ X_i = \text{STOP} | X_{i-1} \neq \text{STOP} \} \geq p$$

$$A = \{ \text{המשטר } i \text{ במילה } i \} = \bigcup_{i=1}^{\infty} \{ \text{המשטר } i \text{ במילה } i \} = \bigcup_{i=1}^{\infty} A_i$$

$$A^c = \{ \text{המשטר } i \text{ במילה } i \} = \bigcap_{i=1}^{\infty} \{ X_i \neq \text{STOP} | X_{i-1} \neq \text{STOP} \}$$

$$P(A^c) = P\left(\bigcap_{i=1}^{\infty} \{ X_i \neq \text{STOP} | X_{i-1} \neq \text{STOP} \}\right) =$$

$$\text{בנייה} \quad \text{בנוף הסגנון} = \prod_{i=1}^{\infty} P(\{ X_i \neq \text{STOP} | X_{i-1} \neq \text{STOP} \})$$

$$= \sum_{i=1}^{\infty} (1-p)$$

$$1+x \leq e^x \quad \Rightarrow \exp(-\sum_{i=1}^{\infty} p)$$

$$= \frac{1}{\exp(\sum_{i=1}^{\infty} p)} \longrightarrow \frac{1}{\infty} \longrightarrow 0$$

$$\Rightarrow \boxed{P(A) = 1}$$

סדר הסדר
הקבוצה מוגדרת

(b) בהק' א' הראנו כי:

שאלה 1

(d) רצינו להראות שסכום ההסתברויות של כל המשפטים הסופיים הוא 1. נשתמש ברמז ונראה זאת ע"י המאונך הימני.

יהי מילה אנטיגון: $S = \{START, w_1, \dots, w_n, \dots\}$

סוגי זוגות $\forall w \quad P(STOP | w) > 0$

כן $\forall w, \forall w' \neq STOP \quad P(w' | w) < 1$

bigram def.

כפוף:

$$P(S) = P(START, w_1, \dots, w_n, \dots) \stackrel{\text{bigram def.}}{=} \underbrace{P(w_1 | START)}_{CER} \cdot \underbrace{\prod_{i=2}^{\infty} P(w_i | w_{i-1})}_{< 1} \xrightarrow{i \rightarrow \infty} 0$$

זה נובע מכך שאם אינסופי של איברים קטנים מ-1, למדנס-0.

מאחר והזכר יהיה נכון לכל מילה אנטיגון, נסך שזמן סכום ההסתברויות של כל המשפטים האנטיגונים יהיה 0.

לפיכך מאחר וראינו כי סכום ההסתברויות המשפטים אנטיגון (STOP) הוא אפס. נסך שסכום ההסתברויות המשפטים הסופיים (יהי האנטיגון כוללם) הוא 1 בדיוק.

□

(b) נרצה לטעון שההוכחה מסתמך על תצפית במקרה למיזם שאתם מרחיבים.

נתבונן במרחב המילים $C = \{W\}$

אנטיגון

$$P(w_1 = w | w_0 = START) = e^{-1}$$

$$P(w_n = w | w_1 \dots w_{n-1} = \underbrace{w \dots w}_n) = \prod_{i=1}^n e^{-\frac{1}{i^2}} = e^{-\sum_{i=1}^n \frac{1}{i^2}} = e^{-\frac{\pi^2}{6}} = e^{-\frac{1}{\frac{6}{\pi^2}}} < 1$$

$$P(w_n = STOP | w_1 \dots w_{n-1} = w \dots w) = 1 - e^{-\sum_{i=1}^n \frac{1}{i^2}}$$

נבחין שההסתברויות בין 0 ל-1
אם חוקיות אנטיגון יהיה קן נסכומם 1

$$e^{-\sum_{i=1}^n \frac{1}{i^2}} + 1 - e^{-\sum_{i=1}^n \frac{1}{i^2}} = 1$$

בנוסף נבחין כי ההסתברות לכך שגביה תהיה אנטיגון היא

$$\lim_{n \rightarrow \infty} P(w_n = w | w_1 \dots w_{n-1} = \underbrace{w \dots w}_n) \rightarrow e^{-\frac{1}{\frac{\pi^2}{6}}}$$

קיבלנו שזמן עבור $\infty \rightarrow n$ יש הגרנטור
זוהי קיטור הסתברות חוקיות וחוקיות למאונך.
אכן תהיה אפס.
בשנה מתקרה עם המור המרחיבים.

שאלה 2

ע"ט חוגר Unigram.

$$\forall w \in C \quad P(w) = \frac{\#w}{N}$$

(d) המשפט - Unigram:

בהינתן קורפוס בגודל N וזאת w $S = (w_1, \dots, w_n)$ (שני מילים הסתברות המילה)
באופן הבא:

$$P(S) = \prod_{i=1}^n P(w_i) = \prod_{i=1}^n \frac{\text{count}(w_i)}{N}$$

כלומר, כל מילה מקבלת הסתברות בלתי תלויה. היחסות הנתון $\#$ כל המילים בקורפוס.
• נגד המילה w Where

$$P(\text{Where}) = \frac{\text{count}(\text{Where})}{N} > \frac{\text{count}(\text{Were})}{N} = P(\text{Were})$$

• נגד המילה w Were

$$P(\text{Where}) = \frac{\text{count}(\text{Where})}{N} < \frac{\text{count}(\text{Were})}{N} = P(\text{Were})$$

כלומר בשני המקרים, נבחרת המילה שיש לה שכיחות גבוהה יותר, חלילה
לפי count גבוה יותר בקורפוס.

סבור המשפט שלנו לא יתכן שנבחר נכונה באי המופיע בו זמנית מחזור ולא
יתכן כי $P(w) < P(w)$ ו $P(w) > P(w)$.

סעיף 4

מוקד ה bigram יוצר באופן הבא:

עבור קורסוס כל מילה וקדם את הסכימה המילה Where, were בהקשר
למילה א המופיעה אחריו

$$P(\text{Where} | W) = \frac{\text{Count}(\text{Where}, W)}{\text{Count}(W)}$$

כלומר את סכימת המילה Where למחרת W בהיחס לסך המילה W.
נחשב אותו קודם עבור Were וכל מילה אחרת במסד הנתונים.

- מוקד ה bigram יכל להיות אובייקט וזו מכיוון שעבור הקצאה שלנו יתכן ונבחר את
המילה נכונה בשני האפשרויות. זאת כאשר יתקיים:

$$P(\text{Where} | \text{went}) = \frac{\text{Count}(\text{Where}, \text{went})}{\text{Count}(\text{went})} > \frac{\text{Count}(\text{were}, \text{went})}{\text{Count}(\text{went})} = P(\text{were} | \text{went})$$

ע"פ:

$$P(\text{Where} | \text{there}) = \frac{\text{Count}(\text{Where}, \text{there})}{\text{Count}(\text{there})} > \frac{\text{Count}(\text{were}, \text{there})}{\text{Count}(\text{there})} = P(\text{were} | \text{there})$$

- "גרנט המשל" יקבל הסכימות שבה מילים אחרת המשל לא יופיעו כלל
בקורסוס ע"ה לאפס את הסכימות המשל כלל.

$$P(S) = \prod_{i=1}^n P(w_i | w_{i-1}) \quad \text{הוא: ע"פ המודל}$$

לכן אם יש צמד מילים w_{i-1}, w_i אשר לא הופיעה (בסימולציה) בקורסוס
המיומן שלו $P(w_i | w_{i-1}) = 0$.

כאשר נחשב את המשל הכולל ונקבל שכן קיבל הסכימות 0 למחרת ולא נחשב
קודם כן מילה אחרת אחרת.

שאלה 3:

א) נראה שמגדל ק"מ הסדר:

$$\sum_{c=1}^{C_{max}} \left(\frac{(C+1)N_{C+1}}{N_C \cdot N} \cdot N_C \right) = 1 - P_{unseen}$$

$\xrightarrow{\text{התקף של } N_C \text{ באתר שנבדק } N_C \text{ פעמים}}$
 $\xrightarrow{N_C \text{ הניסויים}}$

$$\Leftrightarrow \frac{1}{N} \sum_{c=1}^{C_{max}} (C+1)N_{C+1} = \frac{N - N_1}{N}$$

$$\Leftrightarrow \sum_{c=1}^{C_{max}} (C+1)N_{C+1} = N - N_1$$

$$N_1 + 2N_2 + 3N_3 + \dots + C_{max}N_{C_{max}} + \cancel{C_{max+1} \cdot N_{C_{max+1}}} = N$$

$$N_1 + 2N_2 + 3N_3 + \dots + C_{max}N_{C_{max}} = N$$

ואכן בשיון נכון למציג N

ב) חיינו לרשום את המשוואה של Add-On של מילה הנוספה C פעמים.

$$q_{add-1}(w) = \frac{C(w)+1}{\sum_{w' \in Y} (C(w')+1)} = \frac{C(w)+1}{C() + |Y|} = \frac{C+1}{N + \sum_{i=0}^{C_{max}} N_i}$$

$\xrightarrow{\text{מספר } C \text{ פעמים מילוי התקף } W}$
 $\xrightarrow{\text{כמות המילים בשפה}}$
 $\xrightarrow{\text{מילוי הקופסה מאשר } N}$

$$\xrightarrow{MLE} \frac{C+1}{N+|Y|} > \frac{C}{N} \Leftrightarrow \frac{C+1}{C} > \frac{N+|Y|}{N} \Leftrightarrow 1 + \frac{1}{C} > 1 + \frac{|Y|}{N} \Leftrightarrow \boxed{\frac{N}{|Y|} > C}$$

כלומר מספר ההופעות של מילה קטן $\frac{N}{|Y|}$ של נקודת התקפה של ה MLE.
 ואם $\mu = \frac{N}{|Y|}$ מילה של הניסויים.

(C) חיינו להראות שהמילה מסווגת כי נכונה להחליף מסווגת.

נראה שיש ע"י דוגמה נשגת:

נשגת $C_{max} = 5$ ולכן C בין 1-5 נשגת $N_C = C$.
 אזי עבדי מילה W הנוספה C פעמים בקופסה יתקיימ:

$$\frac{(C+1) \cdot N_{C+1}}{N \cdot N_C} = \frac{(C+1)^2}{C \cdot N} = \frac{C^2 + 2C + 1}{C \cdot N} > \frac{C(C+2)+1}{C \cdot N} > \frac{C(C+2)}{C \cdot N} = \frac{C+2}{N} > \frac{C}{N} = MLE$$

כלומר אם מספר התקפה C , והנשגת הנ"ל, גרמי אומדן ההתקפה יהיה גדול יותר מה MLE
 ולכן סבירה לרשום מסווגת 2.

שאלה 4:

(a) המשמאה למודל triagram היא:

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

ההנחה שהמודל הוא הן שלמה ולכן הן בסגור המילים שהיא רפניה וזה כאלו אחר אחר.

(b) עברית: הכלבים נובכים על הכדור. The dogs barks on the ball.

המילה "נובכים" צמחה לכלבים לכן המודל ולמה אחר.

(c) אנגלית: הכלבים רצו במהירות של ראשון היפה ונכך חזר.

המילה "נכח" חלוצה בצורה הראשונה היפה. ואינה חלוצה במילים כלבים או רצו לכן המודל לא יודע את נכחו למידה שבו שזה יקרה יותר.

או המודל היה צ-ערס כמראה שהמודל היה מבין שזה "נכחו".

The dogs ran to the beautiful lake, and bark on it.

המילה bark כחומר וטובה - barked, אך אין מהמיל כוחן כך 3 מילים אחרות הוא לא ישים שפירגניה סבבו. חל המודל היה ט-ערס הוא היה מבין שאתנו כלטון עבר ומקן נחמס.

שאלה 5

הכלבים רצו במהירות של ראשון היפה ונכך חזר.

המילה חזר במהירות של ראשון היפה ונכך חזר.

נמין כי כל שהמודל המוקדני קלל ח לאלו זמרי הסתנו אלו הרכיבי משל לא מקן קאנה.

פרויקט של מלח תכנות:

*** Question 2 ***

The predicted sentence using bigram model is:

I have a house in the

*** Question 3 *** , using the bigram model:

Question 3 (a) - The probability of the following two sentences:

Probability for question 3a: -inf

Probability for question 3b: -29.667

Question 3 (b) - The perplexity of both the following two sentences:

Perplexity for question 3b: inf

*** Question 4 *** - Estimate a new model using linear interpolation smooth between the bigram model and unigram model:

The probability of 1th sentence is: -36.192

The probability of 2nd sentence is: -30.993

The perplexity of sentences 1 and 2 is: 270.076

Process finished with exit code 0