

תרגיל 4 NLP

אור מלכאי – 203569264 – ormalkai@gmail.com

איציק מור – 204240246

סטפן גולדברג – 316802222

שאלה 1

סעיף a

i. נציג מספר דוגמאות ל-ambiguity במשפט

הבא (https://www.washingtonpost.com/business/economy/mary-schapiro-to-leave-sec-next-month/2012/11/26/d61d4c74-37df-11e2-8a97-363b0f9a0ab3_story.html?noredirect=on&utm_term=.39a79ff1710c):

S.E.C. chief Mary Shapiro to leave Washington in December.

התיג המקורי הוא:

(ORG S.E.C.) chief (PER Mary Shapiro) to leave (LOC Washington) in December.

בתיג הזה S.E.C זה למעשה acronym ל-Security and Exchange Commission וזה למעשה ארגון.

בנוסף Mary Shapiro תיג כבנאדם ו-Washington תיג כמקום.

במשפט הנ"ל ניתן לראות שתי דוגמאות ל-ambiguity הראשונה, הקלה יותר, היא ש-Washington במקרה הזה זה שם של עיר, אך מאחר והעיר נקראה על שם של בנאדם, בקלות ניתן למצוא מופעים של Washington שיתויגו כבנאדם. בנוסף נשים לב שהמילה chief לא נכללה בתיג "בנאדם" של Mary Shapiro, בקלות היה ניתן לתייג את S.E.C chief Mary Shapiro בתיג אחד של בנאדם, או אולי אפילו תיג ותיגים פנימיים.

ii. הצורך בפיצ'רים נוספים מלבד המילה עצמה נועד כדי לפתור קונפליקטים שנובעים מ-ambiguity. למשל, אם היינו פותרים קונפליקטים ע"י אומד MLE של תיג למילה, על פי התיג הנפוץ ביותר ב-training היינו בטוח טועים על התיג שראינו פחות. כדי לפתור זאת עלינו להבין מה ה-context של המילה ולהסתכל על פיצ'רים נוספים. יתר על כן, הצורך בפיצ'רים נוספים נובע גם מהצורך לטפל ב"rare words" שלרוב הן named entities.

iii. ישנם מספר פיצ'רים שניתן לחשוב שיעזרו בתיג named entities למשל, המילה הקודמת/הבאה, או באופן כללי חלון של מילים, או "תקציר" של המילים בחלון מסוים. פיצ'ר נוסף יכול להיות ה-Part Of Speech של המילה. יתר על כן ניתן לחשוב גם על פיצ'רים שקשורים בסינטקס של המילה, למשל prefix או suffix, לדוגמא מילים שמסתיימות ב-ל רוב הסיכויים לא יהיו Person או Organization.

סעיף b

$$e = (1, (2w+1) * D) \quad i.$$

$$W = ((2w+1) * D, H)$$

$$U = (H, C)$$

כאשר D הוא גודל embedding, H הוא גודל השכבה החבויה ו-C הוא מספר ה-classים.

ii. לפי המימדים שתיארנו בסעיף הקודם, עבור חלון יחיד בגודל w , הסיבוכיות היא:

$$O((2w + 1)D) \text{ עבור חישוב } e^{(t)}$$

$$O((2w + 1)DH + H) \text{ עבור חישוב } h^{(t)}$$

$$O(HC + C) \text{ עבור חישוב } y^{(t)}$$

לכן, סה"כ עבור משפט באורך T נקבל $O((2w + 1)DHT + HC)$

סעיף d

i. התוצאה הטובה ביותר שקיבלנו עבור ה-development set היא:

$$F1 = 0.83, Precision = 0.82, Recall = 0.85$$

Confusion Matrix

Gold\Guess	PER	ORG	LOC	MISC	O
PER	2981	18	65	14	71
ORG	156	1592	133	75	136
LOC	47	68	1910	22	47
MISC	45	46	58	1007	112
O	47	41	16	29	42626

לפי Confusion Matrix ניתן לראות כי באופן יחסי המודל מנבא באופן לא רע, כלומר רוב האנשים מזוהים כאנשים, רוב המקומות מזוהים כמקומות, רוב ה"שונות" מזוהים כ"שונות" ורוב ה-O מזוהים כ-O. לעומת זאת, הרבה ארגונים מזוהים כאנשים, ניתן לראות זאת גם על בסיס ה-Recall הנמוך, 0.76, ברמת Token.

ii. הבעיה עיקרית של windowed model היא שהוא לא נעזר בפרדיקציות שבחלון בשביל לנבא פרדיקציות נוספות, בדומה למשל ל-log linear models שראינו בתרגיל הקודם עבור POS. בגלל המגבלה הזו ניתן לראות "רצפים" של entities שלא מזוהים כמו שצריך, למשל עבור המשפט:
 Starting(go:O, gu:O) on(go:O, gu:O) May(go:O, gu:O) 13(go:O, gu:O) next(go:O, gu:O) year(go:O, gu:O), (go:O, gu:O) the(go:O, gu:O) Test(go:ORG, gu:MISC) and(go:ORG, gu:O) Country(go:ORG, gu: ORG) Cricket(go: ORG, gu: ORG) Board(go: ORG, gu: ORG)

במקרה הנ"ל קל לראות ש"Test and Country Cricket Board" זה בעצם ארגון, ואם המודל היה משתמש בתיגים סביב כדי לנבא את Test ואת and הוא כנראה לא היה טועה. באופן דומה ניתן גם לומר ש-windowed model לא נעזר במידע משאר המשפט, כלומר חלקי המשפט שלא נמצאים בחלון. עבור משפטים שמכילים "לוואי" ארוך ניתן לחלץ מידע רלוונטי דווקא מחלקי משפט "רחוקים".

שאלה 2

סעיף a

i.

ישנם שני הבדלים בין מודל RNN לבין מודל window. ראשית המימד של W_e הוא (D, H) במודל RNN לעומת $((2w+1) * D, H)$ עבור W במודל window. שנית, W_h לא קיימת במודל window והמימד שלה הוא (H, H) .

ii.

עבור cell יחיד ב RNN הסיבוכיות היא:

$$O(D) \text{ עבור חישוב } e^{(t)}$$

$$O(H * H + DH + H) = O(H(H + D + 1)) = O(H(H + D)) \text{ עבור חישוב } h^{(t)}$$

$$O(HC + C) \text{ עבור חישוב } y^{(t)}$$

לכן, סה"כ עבור משפט שלם באורך T נקבל $O((D + H)HT)$ בהנחה ש $D, H \gg C$

סעיף b

i.

ראשית עלינו להבין את החישוב של F1 ובפרט את Precision ו Recall ברמת ה-entity. Precision הוא החלק היחסי של ה-named entities שהמודל חזה נכונה שמישורים עם-named entities ב-development set, מתוך כל ה-named entities שהמודל חזה. $Precision = \frac{tp}{tp+fp}$

Recall הוא החלק היחסי של ה-named entities ב-development set שהמודל חזה נכונה מתוך כל

$$Recall = \frac{tp}{tp+fn} \text{ ב-named entities ב-development set}$$

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall} \text{ הוא הממוצע ההרמוני שלהם.}$$

נתבונן במשפט הבא:

תיוגים נכונים:

"Harry(PER) Potter(PER) and(O) the(LOC) chamber(LOC) of(LOC) secretes(LOC)"

תיוגים אפשריים:

1. "Harry(PER) Potter(PER) and(O) the(LOC) chamber(LOC) of(LOC) secretes(O)"

2. "Harry(PER) Potter(PER) and(O) the(O) chamber(O) of(O) secretes(O)"

נניח ובחרנו בתיוג הראשון על פני התיוג השני, ברור כי ה-cross entropy תקטן מכיוון ש ברמת tokens אנחנו צודקים יותר. לעומת זאת ברמת ה-entity ה-F1 גם כן יקטן:

$$Precision1 = \frac{1}{1+1} = 0.5, Recall1 = \frac{1}{1+1} = 0.5, F1 = \frac{2 * 0.5 * 0.5}{0.5 + 0.5} = 0.5$$

$$Precision2 = \frac{1}{1+0} = 1, Recall2 = \frac{1}{1+1} = 0.5, F1 = \frac{2 * 1 * 0.5}{0.5 + 1} = \frac{1}{3}$$

ii.

ישנן לפחות שתי סיבות לכך שקשה לבצע אופטימיזציה ישירות על F1. הסיבה המתמטית היא שהפונקציה לא דיפרנציאבילית. הסיבה הפרקטית היא שבשביל לחשב F1 דרוש חישוב על פני כל datan וזה יקר.

סעיף d

i.

במידה ולא היינו מבצעים masking ה loss היה כולל גם את השגיאה שלנו על padding שהוספנו.

לכן הגרדיאנטים של padding שהוספו היו משפיעים על הפרמטרים של הרשת במהלך backpropagation. באמצעות השימוש בmasking למעשה loss על padding הוא 0, לכן גם הגרדיאנטים של padding יתאפסו ולא נשפיע על הפרמטרים של הרשת.

סעיף g

נתאר שתי מגבלות של מודל RNN:

1. מגבלה ראשונה במימוש RNN שלנו היא שהRNN עובדת משמאל לימין, כלומר מתחילת המשפט לסופו. ברור שאם היינו יודעים מידע על המשך המשפט היינו יכולים לתייג יותר טוב. פתרון פשוט, בדומה לbiLSTM ניתן לבצע מעבר גם מהסוף להתחלה. דוגמא מהdata:

May(go:O, gu:O) 15(go:O, gu:O) v(go:O, gu:O) Duke(go:ORG, gu:ORG) of(go:ORG, gu:O)
Norfolk(go:ORG, gu:ORG) 's(go:ORG, gu:O) XI(go:ORG, gu:ORG) ((go:O, gu:O) at(go:O, gu:O)
Arundel(go:LOC, gu:LOC)) (go:O, gu:O)

בדוגמא הנ"ל טעינו על s' ועל of, כנראה שלא היינו טועים אם היה לנו מידע על המשך המשפט.

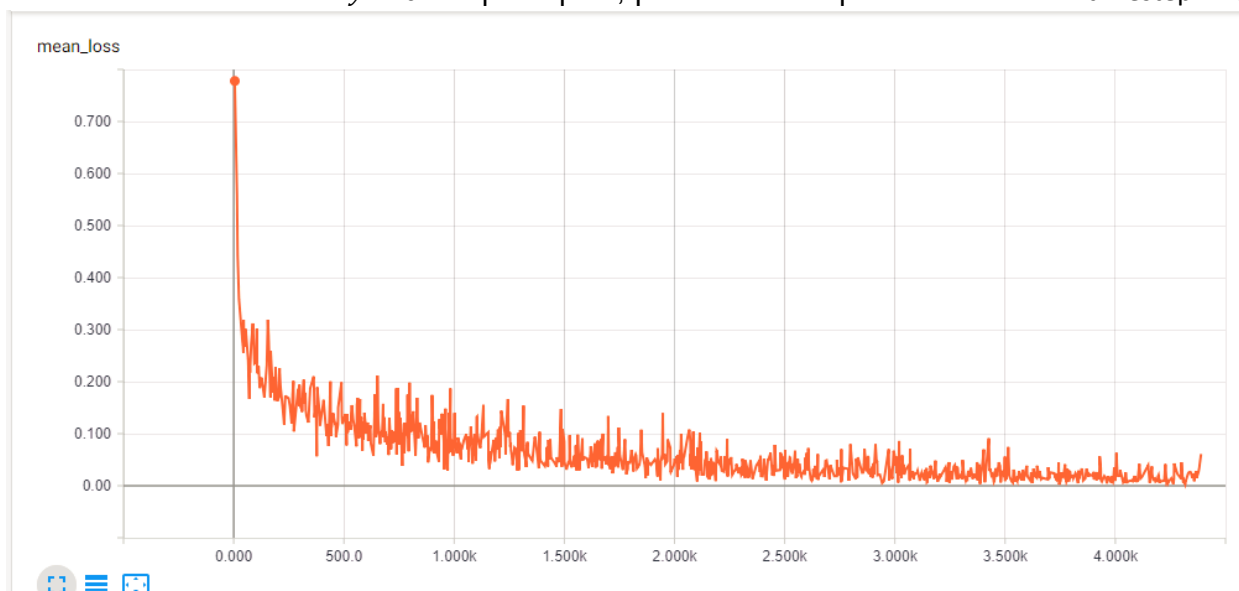
2. המגבלה השנייה במימוש שלנו היא שהיינו רוצים ל"תגמל" תיוגים רציפים של אותם entities. הדוגמא לעיל עדיין טובה, שכן היינו רוצים להאמין שבתוך היער של ORG נמשיך לתייג ORG ולא O. יכולנו לפתור זאת ע"י רגולריזציה בloss שנותנת עדיפות לתיוגים שכנים זהים. לחילופין, יכולנו להכניס קלט נוסף לכל cell שהוא התיוג הקודם.

שאלה 3

סעיף c

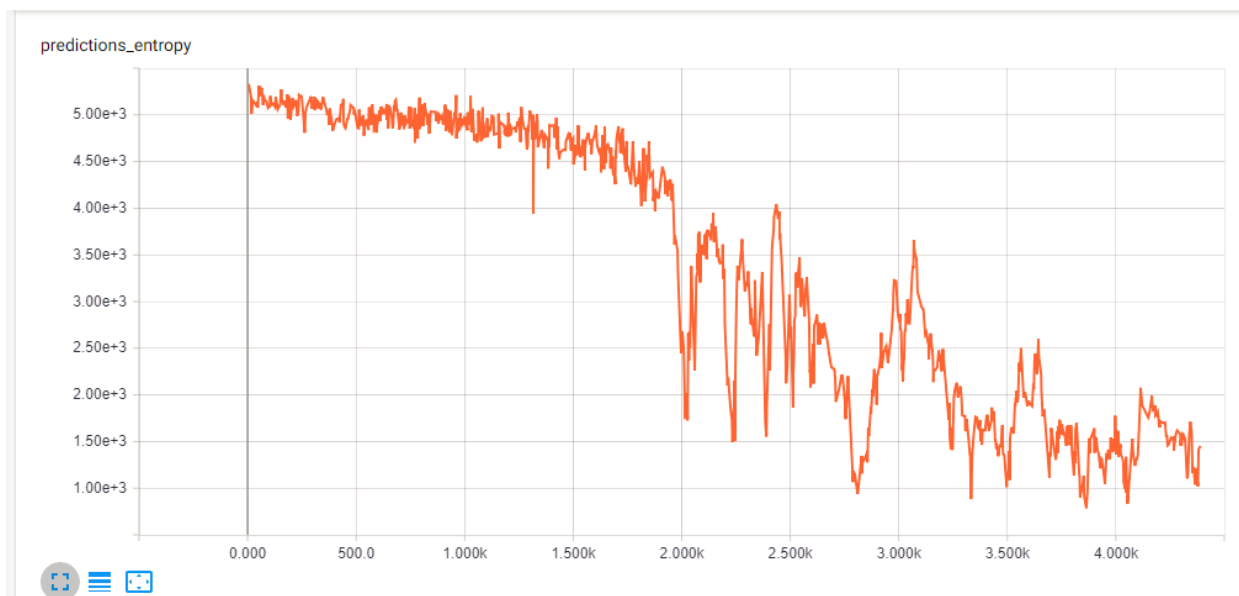
i.

עבור timestep יחיד האנטרופיה המקסימלית היא אינסוף, במקרה שקיבלנו $\hat{y} = 0$

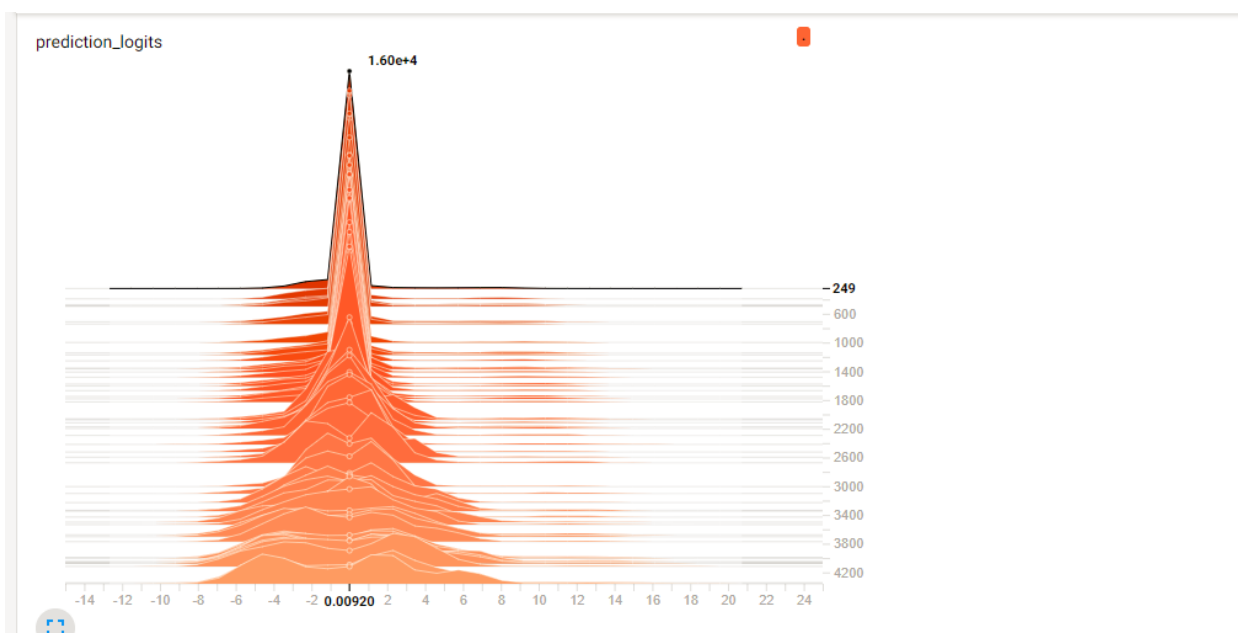
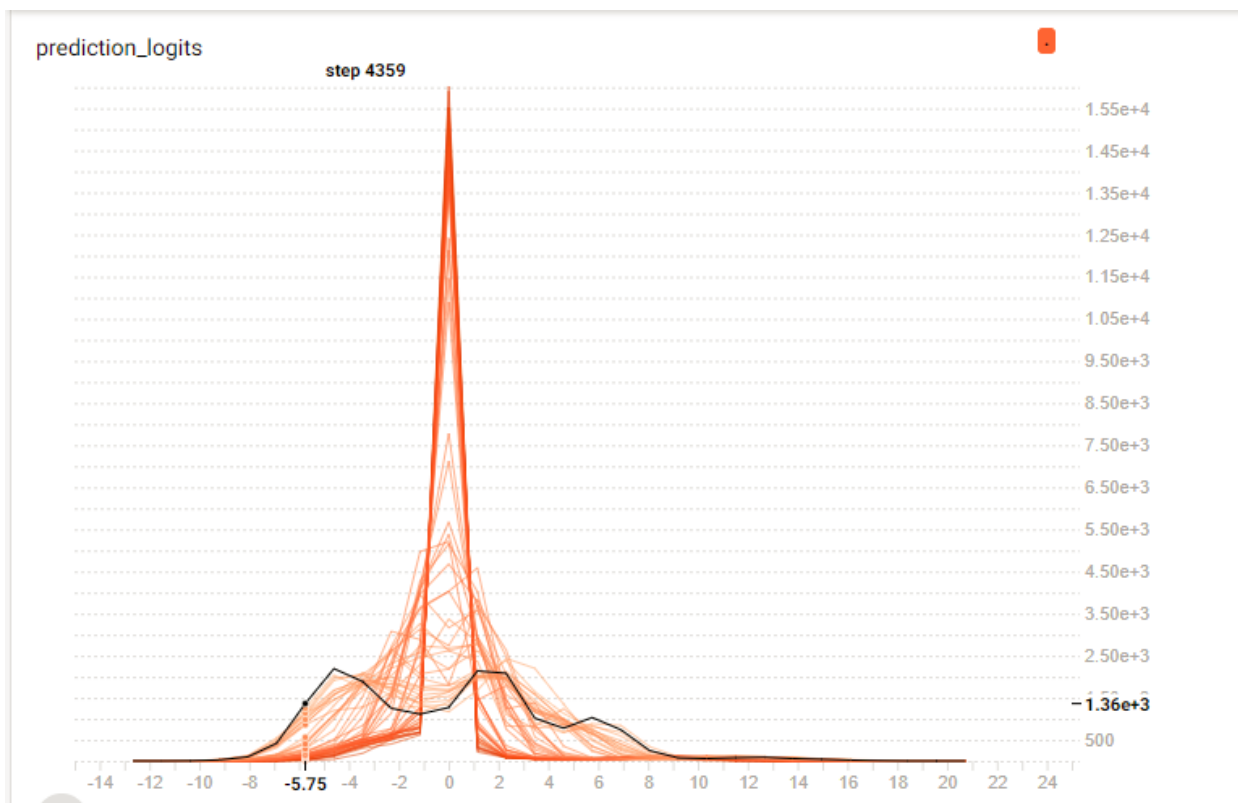


ii.

בגרף לעיל קל לראות שהloss יורד ככל שהtimesteps עולים, בדיוק כפי שהיינו מצפים.



בגרף לעיל ניתן לראות שהאנטרופיה יורדת ככל שהtimesteps עולים.



בשני הגרפים לעיל ניתן לראות שההיסטוגרמה של הפרדיקציות הופכת לשטוחה ורחבה ככל שהstep עולים.

ניתוח הגרפים:

האנטרופיה למעשה מכמתת את אי הוודאות שיש לנו בdata. מתורת האינפורמציה האנטרופיה למעשה מייצגת את מספר הביטים הדרוש כדי לתאר את y . לעומת זאת cross entropy למעשה מכמתת את מספר הביטים הדרושים כדי לעשות encoding ל y על ידי אומד \hat{y} (מילים אחרות ללומר שזו דרך למדידת מרחק בין y ל \hat{y}).

