

תרגיל בית 2 – עיבוד שפה טבעית – חלק רטוב

תיאור המשימה

בתרגיל בית זה תממשו Dependency Parser ללמידת דקדוק תלויות (כפי שנלמד בשבוע 9), ותנתחו את טיב הצלחתכם.

כפי שתראו בהמשך, בתרגיל הבית הנוכחי הופחתו מספר הרכיבים האלגוריתמים אותם אתם צריכים לממש. זאת על מנת שתתעסקו יותר בבניית הייצוג, בהחלטה מהן תכונות (features) טובות, לפי איזה סף כדי לקצץ וכו'. מכיוון שכך, הציון על הביצועים בתחרות יהיה גבוה.

סגל הקורס ממליץ להשתמש בשפת python, אולם תוכלו להשתמש בכל שפת תכנות שתחפצו. יש לקחת בחשבון כי נפח התרגיל נקבע על סמך ההנחה שאתם עובדים עם שפה עילית, המאפשרת גמישות בעבודה עם נתונים בצורה נוחה יותר.

נתונים:

הסבר על הקבצים המצורפים –

1. train.labeled – קובץ המכיל 5000 משפטים מתווייגים. עליכם להשתמש בקובץ זה בשלב האימון (הסבר בהמשך)
2. test.labeled – קובץ המכיל 1000 משפטים מתווייגים, בפורמט זהה לפורמט של הקובץ הקודם.
3. comp.unlabeled – קובץ המכיל 1000 משפטים לא מתווייגים.

פורמט קבצי האימון (הדוג' היא המשפט השני בקובץ train):

Token Counter	Token	–	Token POS	–	–	Token Head	Dependency Label	–		–
1	Mr.	–	NNP	–	–	2	NAME	–		–
2	Vinken	–	NNP	–	–	3	VMOD	–		–
3	is	–	VBZ	–	–	0	ROOT	–		–
4	chairman	–	NN	–	–	3	VMOD	–		–
5	of	–	IN	–	–	4	NMOD	–		–
6	Elsevier	–	NNP	–	–	7	NAME	–		–
7	N.V.	–	NNP	–	–	5	PMOD	–		–
8	,	–	,	–	–	7	P	–		–
9	the	–	DT	–	–	12	NMOD	–		–
10	Dutch	–	NNP	–	–	12	NMOD	–		–
11	publishing	–	VBG	–	–	12	NMOD	–		–
12	group	–	NN	–	–	7	APPO	–		–
13	.	–	.	–	–	3	P	–		–

- כל שורה מייצגת מילה, וכוללת 10 עמודות, המופרדות ע"י 'ז'
- העמודות היחידות הרלוונטיות למשימה שלנו הן הצבועות באדום – מיקום המילה במשפט, המילה עצמה, חלק הדיבר המתאים עבורה והראש שלה.
- בין כל זוג משפטים יש שורה ריקה
- בקובץ התחרות, בעמודה token head (בכך העמודה Dependency Label) יש קו תחת ('_')

אימון (Train) :

כאמור את הערכת הפרמטרים תעשו על הקובץ train.labeled. קשתות העץ הפורש ימושקלו על סמך התכונות הבאות (כפי שראינו במאמר של McDonald) :

Unigram Features

- 1.) p-word, p-pos
- 2.) p-word
- 3.) p-pos
- 4.) c-word, c-pos
- 5.) c-word
- 6.) c-pos

Bigram Features:

- 7.) p-word, p-pos, c-word, c-pos
- 8.) p-pos, c-word, c-pos
- 9.) p-word, c-word, c-pos
- 10.) p-word, p-pos, c-pos
- 11.) p-word, p-pos, c-word
- 12.) p-word, c-word
- 13.) p-pos, c-pos

כאשר:

p-word: word of parent (head) node.

p-pos: POS of parent (head) node.

c-word: word of child (modifier) node.

c-pos: POS of child (modifier) node.

אתם נדרשים לבנות שני מודלים:

1. מודל בסיס, המורכב מ 1 עד 13, מלבד 7,9,11,12 (אלו תכונות עם דלילות גבוהה)
2. מודל מורכב, בו תוכלו לבחור כל תכונות שתמצאו. (תתפרעו)

במהלך האימון, תשתמשו בStructured Perceptron כפי שהוצג בהרצאות הוידאו:

Perceptron($\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|\mathcal{T}|}$)

1. $\mathbf{w}^{(0)} = \mathbf{0}; k = 0$
2. for $n : 1..N$
3. for $t : 1..T$
4. Let $\mathbf{y}' = \arg \max_{\mathbf{y}'} \mathbf{w}^{(k)} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$
5. if $\mathbf{y}' \neq \mathbf{y}_t$
6. $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$
7. $k = k + 1$
8. return $\mathbf{w}^{(k)}$

כאשר t מייצג את האינדקס של משפט מסוים, ו n מייצג את מספר האיטרציה.

בסיכום שתכינו (עוד על כך בהמשך) יש לציין את מספר המאפיינים שמכל סוג (לדוג' 217 מאפיינים מסוג' p-word, p-pos). את המאפיינים שבחרתם למודל המורכב יש להגדיר במפורש.

כל שיפור שהכנסתם למודל (לדוג' קיצוץ של מאפיינים ש"לא הופיעו מספיק") צריך להיות מוסבר היטב, כולל המוטיבציה לבצע אותו. אם השתמשתם בחבילת קוד חיצונית, יש לפרט במפורש היכן ולאילו מטרות.

בנוסף, יש לפרט כמה זמן לקח לאמן כל מודל (וכן מפרט בסיסי של החומרה עליה הרצתם)

הסקה (Inference):

ההסקה תתבצע ע"י אלגוריתם Chu-Liu-Edmonds הנלמד בתרגול 9.

מבחן (test):

לכל אחד מן המודלים, יש לבצע הסקה (Inference) על הקובץ test.labeled, ולדווח את תוצאות הדיוק (accuracy) ברמת מילה, כפי שנעשה בהרצאה (ומופיע בתרגול 9). יש לדווח תוצאות על קובץ test עבור $N=20,50,80,100$ (מספר האיטרציות באלגוריתם ה Perceptron).

התייחסו להבדל בביצועים בין המודלים השונים, והעלו מספר סיבות שעשויות לגרום להבדלים אלו.

בנוסף, אנא ציינו כמה זמן לקח לתייג את הקובץ לפי כל אחד מהמודלים (וכן מפרט בסיסי של החומרה עליה הרצתם את הקוד)

תחרות:

לכל אחד מן המודלים (כאשר מספר האיטרציות הוא לשיקולכם), יש לבצע הסקה (Inference) על הקובץ comp.unlabeled (אשר אינו כולל תיוגים), ולכתוב את תוצאות התיוג לתוך קובץ חדש בפורמט labeled (כמו קבצי האימון) (שמות הקבצים הרצויים מופיעים בהמשך). לדוג', עבור המשפט:

Token Counter	Token	—	Token POS	—	—	Token Head	Dependency Label	—		—
1	The	—	DT	—	—	—	—	—		—
2	Boy	—	NNP	—	—	—	—	—		—

יש לבצע הסקה, שתיתן לכם את התלויות. בהנחה שהתלויות שמצאתם הן $1 \rightarrow 0$, $2 \rightarrow 1$, תכתבו אותן לקובץ ההגשה באופן הבא -

Token Counter	Token	—	Token POS	—	—	Token Head	Dependency Label	—		—
1	The	—	DT	—	—	0	—	—		—
2	Boy	—	NNP	—	—	1	—	—		—

שימו לב שסדר המשפטים (הלא מתויגים) בקובץ המקורי זהה לסדר המשפטים בקובץ הפלט, שמספר העמודות זהה ושארף אחד מן הערכים מחוץ לשתי העמודות ששיניתם לא נפגע.

יש לתאר במפורש מה עשיתם כדי לקבל את התוצאות שקיבלתם (שינויים שביצעתם בלמידה, בהסקה וכו').

בנוסף, אתם מתבקשים לכתוב: תחזית של אחוז הדיוק שאתם צופים לקבל, וכן להסביר מדוע עשוי להיות הבדל בין הדיוק על קובץ ה test. הסברים חכמים אף עשויים לקבל בונוס.

קוד חיצוני המותר לשימוש:

החבילות הסטנדרטיות בשפה בה בחרתם.

לדוג', בשפת python החבילות בהן מותר להשתמש הן numpy ו scipy בלבד.

בשלב האימון ניתן להשתמש בקוד חיצוני המממש את Chu-Liu-Edmonds. אם בחרתם להשתמש בpython, תוכלו להיעזר בקוד המצורף.

למען הסר ספק - אִיֹּר להשתמש ב:

1. חבילות הממשות בנייה של תכונות (features)
2. חבילות הממשות וריאציה של Perceptron
3. חבילות העושות עיבוד על טקסט – ספירת חזרות, uni\tri\bi gram וכו'.

הגשה:

קובץ zip בלבד, בשם HW2-Wet_123456789_987654321.zip (עבור שני סטודנטים שמספרי הזהות שלהם 123456789 ו 987654321). הקובץ הנ"ל יכלול:

1. קובץ הסברים וניתוח תוצאות, הכולל בין היתר:
 - a. שמות המחברים ות"ר
 - b. הערות על אימון המודל, לכל אחד מן המודלים (לפי הדגשים בסעיף "אימון")
 - c. הערות על אלגוריתם ההסקה (לפי הדגשים בסעיף "הסקה")
 - d. ניתוח תוצאות על קובץ המבחן (לפי הדגשים בסעיף "מבחן")
 - e. סיכום שלכם על המשימה
 - f. ניתוח של מה שעשיתם בקובץ התחרות (לפי ההסברים כפי שמפורטים בסעיף "תחרות")

- g. הסבר על חלוקת העבודה בין שני חברי הקבוצה – איזה חלק עשה\ביצע\מימש כל אחד
2. קבצי הקוד של התרגיל. על הקוד להיות מתועד וקריא. בנוסף, הקוד צריך מסוגל לרוץ על כל מכונה שהיא. אנא צרו ממשק פשוט להרצת התוכנית המייצרת את קבצי התחרות.
3. קבצי התחרות – על קבצי התוצאות להיות בפורמט labeled (כפי שמפורט בחלק "אימון"). על מנת להימנע מאי נעימות בנוגע לציון, אנא ודאו כי אם שמים '_' בעמודות ששניתם מקבלים בדיוק את הקובץ comp.unlabeled (אותן שורות לפי אותו סדר). חוסר התאמה פירושו ציון 0 בחלק הזה.
- על שמות הקבצים להיות – (123456789 הוא ת"ז של אחד הסטודנטים)
- a. comp_m1_123456789.wtag – קובץ labeled שאומן על ידי המודל הבסיסי.
- b. comp_m2_123456789.wtag – קובץ labeled שאומן על ידי המודל המורכב.

הציון על התחרות יתבסס על $\max\{accuracy(comp_m1), accuracy(comp_m2)\}$, כלומר נסתכל על התוצאה הטובה יותר מבין שתי ההגשות.

על קבצי התחרות להיות ניתנים לשחזור (Reproducible). הדרישה היא שניתן יהיה לקחת את הקוד שהגשתם ולבנות באמצעותו קובץ זהה לחלוטין לקובץ שהגשתם. במקרים חריגים בודק התרגילים יבדוק את ההתאמה הנ"ל.

```
HW1-Wet_123456789_987654321.zip/
report (.pdf, .docx, etc.)
Code_Directory/
...
comp_m1_123456789.wtag/
comp_m2_123456789.wtag/
```

בסה"כ קובץ ההגשה צריך להיראות כך:

העתקות:

בשל אופי המשימה והמורכבות שלה, קל לבדוק העתקות של קטעי קוד \ קבצים מלאים. למען הסר הספק אנו מדגישים כי אין להעביר קוד בין סטודנטים, בין אם להגשה ובין אם לא. אין לקחת קטעי קוד מוכנים מהאינטרנט, ובכלל אין להסתמך על שום מקור אחר לקוד מלבד פרי יצירכם והחבילות החיצוניות אשר צוינו בסעיף הרלוונטי.