

עיבוד שפה טבעית – תרגיל בית 2 - רטוב

תיאור המשימה

בתרגיל בית זה תממשו Dependency Parser (כפי שנלמד בשבוע 8) מבוסס רשתות נוירונים ללמידת דקדוק תלויות, ותנתחו את טיב הצלחתכם.

אתם מתבקשים לממש מודל המבוסס על ה-Graph-based Parser של Kiperwasser and Goldberg אשר [מאמר](#) הוצג בכיתה. אנו ממליצים לקרוא לפחות את החלקים הבאים מתוך המאמר טרם ובמהלך מימוש המודל ע"מ לוודא הבנתכם: 2.3 (Bidirectional Recurrent Neural Networks), 3 (Our Approach), 5 (Graph-based Parser Implementation Details, Hyperparameter Tuning ו-Tuning Details), 6 (Experiments and Results). כמו כן, אנא חזרו על מצגת התרגול של שבוע 8 ובפרט על אלגוריתם Chu-Liu-Edmonds. המודל שתממשו אינו זהה ב-100% לזה שמוצג במאמר אך דומה לו מאוד, ההבדלים יפורטו בהמשך ההוראות.

כפי שתראו בהמשך, בתרגיל הבית הנוכחי הופחתו מספר הרכיבים האלגוריתמים אותם אתם צריכים לממש. זאת על מנת שתתעסקו יותר בבניית הייצוג, בארכיטקטורת הרשת, מאפייניה וכו'.

סגל הקורס ממליץ להשתמש בשפת Python 3, ובספריית PyTorch למימוש ואימון רשתות נוירונים. לצורך אימון הרשתות בזמן סביר, לרשותכם עומדות מכונות GPU בסביבת ה-Azure. מצורפים לתרגיל הוראות התחברות למכונות אלו.

הסבר על מבנה הציון בתרגיל:

- **60%** - מימוש מלא של המודלים הבסיסי והמתקדם כפי שיפורטו בהמשך, אימוןם על קובץ ה-`train.labeled`, ועמידה בסף אחוז הדיוק (UAS) של 70% על קובץ ה-`test.labeled` (עבור המודל הטוב מבין השניים).
- **30%** - תחרות מבוססת אחוז דיוק (UAS) בתיג קובץ התחרות `comp.unlabeled`. הציון יתבסס על תחרות ביצועים כללית על קובץ התחרות ועל ביצועים בתוך קבוצות משפטים לפי דרגות קושי שונות. הציון על התחרות יתבסס על $\max\{accuracy(comp_m1), accuracy(comp_m2)\}$, כלומר נסתכל על התוצאה הטובה יותר מבין שני המודלים שתגישו (הבסיסי והמתקדם).
- **10%** - כתיבת דו"ח תמציתי (עד 3 עמודים) אשר יכלול את הסעיפים הנדרשים ועמידה בתנאי פורמט ההגשה (יפורטו בהמשך המסמך).

מצורף נספח בסוף הוראות התרגיל המציג דוגמא לחישוב אחוז הדיוק UAS.

נתונים:

הסבר על הקבצים המצורפים –

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

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 התו '\t'.
- העמודות היחידות הרלוונטיות למשימה שלנו הן הצבועות באדום – מיקום המילה במשפט, המילה עצמה, חלק הדיבר המתאים עבורה והראש שלה. נא להתעלם מ-Dependency Label, שכן בתרגיל זה אנו נבצע חיזוי רק לעץ התלויות של המשפט ולא לתוויות של הקשתות.
- בין כל זוג משפטים בקובץ ישנה שורה ריקה בה מופיע התו 'ח' בלבד.
- בקובץ התחרות, בעמודה Token Head (וכן בעמודה Dependency Label) יש קו תחת ('_')

אימון (Train) :

כאמור את שערך הפרמטרים תבצעו על הנתונים שבקובץ train.labeled. קשתות גרף התלויות האפשריות במשפט ימושקלו על סמך התכונות אשר יילמדו ע"י רשת הניורונים הנשנית שתממשו.

אתם נדרשים לבנות שני מודלים:

1. מודל **בסיסי**, מבוסס על רשת הניורונים המוצגת במאמר של Kiperwasser and Goldberg בחלק 5 כפי שהיא ממומשת במאמר, כאשר תדרשו לממש פונקציית loss שונה מזו המתוארת במאמר. במודל שלנו נמזער פונקציית Negative log-Likelihood Loss (NLLLoss) המוגדרת באופן הבא:

$$\min_{\theta} NLLLoss(D; \theta) = \min_{\theta} \sum_{(X^i, Y^i) \in D} \sum_{(h, m) \in Y^i} -\frac{1}{|Y^i|} \cdot \log(P(S_{h,m}^i | X^i, \theta))$$

$$P(S_{h,m}^i | X^i, \theta) = \frac{\exp(S_{h,m}^i)}{\sum_{j=1}^{|Y^i|} \exp(S_{j,m}^i)} = \text{Softmax}(S_{h,m}^i)$$

Where:

- $D = \{(X^i, Y^i)\}_{i=1}^n$ is a dataset consisting of n (sentence, true tree) pairs.
- $X^i = \{x_0 = \text{ROOT}, x_1, \dots, x_{k_i}\}$ is the full sequence of words in the sentence.

- $Y^i = \{(h, m)\}$ is the set of all (head_index, modifier_index) edges in the **true** dependency tree of sentence X^i . $|Y^i| = k_i$.
- $S^i \in R^{(k_i+1)^2}$ is the score matrix for all possible (head, modifier) edges in the dependency graph of sentence X^i .
The cell $S_{h,m}^i$ refers to the score of h being the head of m in sentence X^i .
- θ are all the network's learned parameters.

2. מודל מתקדם, בו תוכלו לבצע כל שינוי שתמצאו למודל הבסיסי: בארכיטקטורת הרשת, פונקציות אקטיבציה, מימדי השכבות, Hyperparameters, פונקציית loss, ואף שימוש ב-word embeddings מאומנים מראש (כגון Word2Vec או GloVe).

במהלך האימון, תשתמשו באלגוריתם אופטימיזציה (Optimizer) מבוסס גרדיאנט (SGD, Adam וכו') של PyTorch לשערוך המשקולות האופטימליות לכל הקשתות האפשריות בגרף התלויות למשפט נתון.

כל שיפור שהכנסתם למודל המתקדם צריך להיות מוסבר היטב, כולל המוטיבציה לבצע אותו.

בנוסף, יש לפרט את ה-Hyperparameters, את זמן האימון הכולל ולצרף גרף המציג את ערך פונקציית ה-loss וערך ה-UAS (בציר y) על פני ה-epochs (בציר x - מס' המעברים על כל המשפטים שבקובץ train.labeled) לכל מודל. ניתן לייצר גרפים נפרדים ל-UAS ו-loss.

הסקה (Inference):

הסקת עצי התלויות (בהינתן המשקולות שנלמדו) תתבצע ע"י אלגוריתם Chu-Liu-Edmonds הנלמד בתרגול 8. אינכם נדרשים לממש אלגוריתם זה, מצורף לתרגיל קוד (chu_liu_edmonds.py) המממש את האלגוריתם, אנא השתמשו בו. עם זאת, אלגוריתם זה ישפיע באופן ישיר על אחוז הדיוק שתקבלו עבור המודלים, כך שאנו ממליצים להבינו היטב ע"מ להשתמש בו בצורה מיטבית.

מבחן (Test):

לכל אחד מן המודלים, יש לבצע הסקה (Inference) על הקובץ test.labeled, ולדווח את תוצאות הדיוק (UAS) ברמת מילה, כפי שנעשה בהרצאת הוידאו (שקף דוגמא מצורף בנספח בסוף הוראות התרגיל). בנוסף, יש לצרף גרף המציג את ערכי ה-loss ו-UAS (ציר y) על פני ה-epochs (ציר x). ניתן לייצר גרפים נפרדים ל-UAS ו-loss.

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.

תחרות:

לכל אחד מן המודלים, 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 מתוייגים) בקובץ המקורי זהה לסדר המשפטים בקובץ הפלט, שמספר העמודות זהה ושארף אחד מן הערכים חוץ מהעמודה ששיניתם לא נפגע.

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

קוד חיצוני המותר לשימוש:

החבילות הסטנדרטיות ב-Python 3, והחבילות הבאות:

pytorch, torchtext, ignite, numpy, scipy, matplotlib, seaborn, pandas, tabulate, jupyter, jupyterlab

מצורף לתרגיל קובץ (nlp_hw2_env.yml) המייצר סביבת Anaconda מתאימה על מכונת ה-Azure, המכילה את כל החבילות המותרות לשימוש. זו גם הסביבה בה ייבדקו התרגילים, אנו ממליצים לעבוד איתה.

מעבר לחבילות המפורטות בקובץ זה, **אין להשתמש** באף חבילה או ספריה חיצונית אחרת, ובפרט שום ספריה או חבילה שעושה עיבוד על טקסט או מממשת Dependency Parser.

עבור ההסקה מומלץ להשתמש בקוד המצורף לתרגיל (chu_liu_edmonds.py) המממש את אלגוריתם Chu-Liu-Edmonds.

במודל המתקדם, אם תרצו להשתמש ב-word embeddings מאומנים מראש, תוכלו להיעזר בחבילת torchtext.

הגשה:

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

1. **דו"ח קצר** (עד 3 עמודים בפורמט PDF) המכיל הסברים תמציתיים, דיווח וניתוח תוצאות, הכולל:
 - a. שמות המחברים ות"ז
 - b. **אימון** - דיווח אחוז דיוק (UAS) על קובץ האימון, גרפים והערות על תהליך אימון המודלים (לפי הדגשים בסעיף "אימון")
 - c. **הסקה** - הערות על שימוש באלגוריתם ההסקה (לפי הדגשים בסעיף "הסקה")
 - d. **מבחן** - דיווח אחוז דיוק (UAS) על קובץ המבחן, גרפים והערות עבור כל מודל (לפי הדגשים בסעיף "מבחן").
 - e. **תחרות** - הסבר קצר על שיפורים שעשיתם למודלים עבור תיוג קבצי התחרות (לפי הדגשים בסעיף "תחרות")
 - f. הסבר קצר 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 להיות מסוגל לרוץ על מכונת Azure עם סביבת העבודה המתאימה. אנא כתבו ממשקי הרצה פשוטים לאימון, מבחן וייצור קבצי התחרות המתויגים.
3. **קבצי התחרות מתויגים** – על קבצי התוצאות להיות בפורמט labeled (כפי שמפורט בחלק "אימון"). על מנת להימנע מאי נעימויות, אנא ודאו כי אם שמים '_' בעמודות ששיניתם מקובץ התחרות 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.labeled – קובץ labeled שאומן על ידי המודל הבסיסי.
 - b. comp_m2_123456789.labeled – קובץ labeled שאומן על ידי המודל המתקדם.
4. **ממשק לתיג קבצי התחרות** - על קבצי התחרות להיות ניתנים לשחזור (Reproducible). הדרישה היא שניתן יהיה לקחת את הקוד והמודלים המאומנים שהגשתם ולייצר באמצעותם קבצי תחרות מתויגים זהים לחלוטין לקבצים שהגשתם. לטובת שחזור הקבצים, יש לכתוב ממשק הרצה פשוט, בקובץ נפרד בעל השם – generate_comp_tagged.py להרצת Inference בלבד על המודלים **המאומנים** ויצירת קבצי התחרות המתויגים ע"י כל מודל.

בסה"כ מבנה קובץ ההגשה צריך להיראות כך (דומה לתרגיל 1):

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

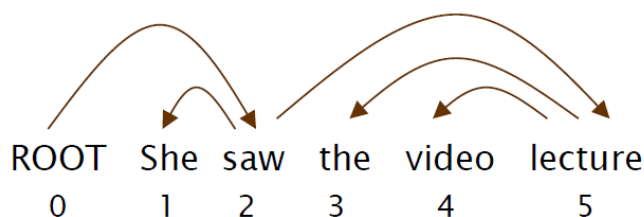
העתקות:

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 להעתיק קטעי קוד מוכנים מהאינטרנט, ובכלל אין להסתמך על שום מקור אחר לקוד מלבד פרי יצירכם והחבילות החיצוניות אשר צוינו בסעיף הרלוונטי.

Christopher Manning



Evaluation of Dependency Parsing: (labeled) dependency accuracy



$$\text{Acc} = \frac{\# \text{ correct deps}}{\# \text{ of deps}}$$

$$\text{UAS} = 4 / 5 = 80\%$$

$$\text{LAS} = 2 / 5 = 40\%$$

Gold

1	2	She	nsubj
2	0	saw	root
3	5	the	det
4	5	video	nn
5	2	lecture	obj

Parsed

1	2	She	nsubj
2	0	saw	root
3	4	the	det
4	5	video	nsubj
5	2	lecture	ccomp