

דו"ח תרגיל בית 4

חלק א':

בשלב הזה בהתחלה קראנו המשפטים מהקובץ jsonl. ואז אחרי זה עבור כל משפט פרקנו לטוקנים ורק מתוך הטוקנים של המשפט לקחנו הטוקנים שהם מילים בעברית ואז הכנסנו הטוקנים של המשפט הזה לרשימה שהרשימה הזאת גם הכנסנו אותה לרשימה בשם `tokenized_sentences` רשימה של רשימות הטוקנים של כל משפט. ואחרי זה עשינו סעיף ב ו-ג כמו מבוקש.

השאלות:

(1) `vector_size` הוא גודל הווקטור שמייצג המילה במרחב הווקטורים כלומר כמה מימדים יהיה הווקטור שמייצג המילה.

- חסרונות הגדלת גודל הווקטור:
- זמן האימון יהיה יותר גדול וגם זה יכול לגרום ל-`overfitting` הדבר שישפיע לרע.
- יתרונות הגדלת גודל הווקטור:
- יהיה לנו יותר מידע כלומר יותר תכונות ואז זה יוביל ליותר דיוק בפעולות.
- חסרונות הקטנת גודל הווקטור:
- יהיה לנו פחות תכונות ומידע הדבר שמקטין את הדיוק.
- יתרונות הקטנת גודל הווקטור:
- זמן אימון יהיה קצר וגם פחות סיכוי ל-`overfitting`.

(2) יש לנו למשל את הבעיה של מילים מחוברים למשל "בית משפט" או "בית חולים" כך שמשני מילים ביחד נקבל ביטוי הגיוני כי בית לבד או משפט לבד זה משהו אחר. בנוסף יש לנו בעיה של אותיות מחוברות עם מילים למשל: "ועד שהגיע" המילה "ועד" אפשר להבין אותה בשני דרכים שונים שלפי הרונטקסט אפשר לקבוע.

חלק ב':

בשביל שלב א עברנו בלולאה עבור כל מילה ב-9 המילים שנתונים לנו ואז עבור כל מילה מהם בדקנו היחד בינה לבין כל מילה בקורפוס שלנו (כלומר המילים שבמודל שבינינו) בעזרת הפקודה `similarity`, ואז אחרי הבדיקות עבור כל מילה מ-9 המילים בחרנו ה-5 מילים הכי דומות לה במודל ואז כתבנו אותם בקובץ `kneset_similar_words.txt`.

בשביל שלב ב עברנו על משפט ברשימה `tokenized_sentences` ואז בכל פעם בהתחלה מגדירים ווקטור של 50 אפסים ואז אחרי זה נעבור עבור כל מילה במשפט הזה (רשימה של טוקנים) נבדוק אם היא במודל שלנו אם כן נוסיף להווקטור את `word_vectors[word]` (כאשר `word_vectors = model.wv`) ועלה הספירה ב-1 ואחרי כל החישוב נחלק הווקטור בקאונטר.

בשביל שלב ג אחרי שבחרנו 10 המשפטים השתמשנו בפונקציה `embeddings_of_sentences` שהגדרנו (שמקבלת רשימת משפטים ומודל) כדי לחשב ה-`embedding` של 10 המשפטים שבחרנו ואחרי זה חשבנו הדמיון בעזרת `cosine_similarity` בין ה-`embedding` של 10 המשפטים וה-`embedding` של המשפטים בדאטא שלנו.

בשביל שלב ד כתבנו 4 המשפטים כאשר המילה שבאדום תהיה עם [] (כלומר [מילה_אדומה]) ואז בשימוש ב-`most_similar` נסינו עבורה הרבה נסיונות ומילים שונים בכל פעם שייתכן עבורם לקבל תוצאות מתאימות כלומר נסינו לבחור מילות שעבורם יתכן לקבל תוצאות מתאימות, ועבור כל מילה לקחנו מילה הכי מתאימה שקבלנו כך הגדרנו המילות האילו ידנית (הכי מתאימות שקבלנו עבור הפקודה `most_similar`) ושמו במקום המילות באדום (עברנו על המשפטים ועבור המילות שמתחילות ב-[ומסתיימות ב-] שהם מילות אדומות הבדלנו אותם עם המילה המתאימה שמצאנו) כתבנו זה המילות החדשות ידנית כי הפקודה תתן תוצאות שונות. אילו הם התוצאות שרואים הם טובות ומתאימות שקבלנו עם מה בחרנו:

```
[('0.9293624758720398', 'יתקיים'), ('0.9383843541145325', 'למלאה'), ('0.9472965002059937', 'שבמסגרת')]
```

```
[('0.8066898584365845', 'לבתר'), ('0.8069127798080444', 'יב'), ('0.818346381187439', 'יכולה')]
```

```
[ (0.9290666580200195 , 'המהלך') , (0.9313275218009949 , 'הדיון') , (0.9397198557853699 , 'ההחלטה') ]
```

```
[('חכם', 0.921226978302002), ('בריא', 0.9138247966766357), ('בפור', 0.9116674065589905)]
```

```
[('סיימת', 0.892031729221344), ('התייחסתי', 0.8826464414596558), ('עוצר', 0.8766990303993225)]
```

```
[('הנוכחים', 0.9081310629844666), ('רבותי', 0.9071613550186157), ('אורי', 0.9056839346885681)]
```

```
[('להודיעכם', '0.9309117197990417'), ('הנני', '0.9263291954994202'), ('מתכבדת', '0.8813149929046631')]
```

[('הנאמן', 0.9270415902137756), ('מהליכוד', 0.9242457151412964), ('לשאלות', 0.9224401116371155)]

ישראל - מדינת

ועדה <- החלטה

ולפעמים יש מילים שבכלל לא קשורים אבל בכל מקרה קיבלו תוצאות גדולות +0.7

אנחנו חושבים שהסיבה העיקרית היא גודל הקורפוס, אם היה לנו אחד יותר גדול היה יכולים להגיע לתוצאות יותר טוב ומילים יותר רלוונטיות.

3 כמו שרואים מתוצאות שקיבלנו שה antonyms קיבלו תוצאות 0.55 ויותר ז"א שהמרחק ביניהן קטן כמו שהסברנו בסעיף קודם

3 כמו שרואים מתוצאות שקיבלנו שה antonyms קיבלו תוצאות 0.55 ויותר ז"א שהמרחק ביניהן קטן כמו שהסברנו בסעיף קודם

שנאה, אהבה 0.9624152183532715
 קר, חם 0.925234854221344
 אחרון, ראשון 0.8305953145027161
 עצוב, שמח 0.6510326266288757
 קטן, גדול 0.5600370168685913

4) אלה התוצאות שקיבלנו:

```

1 | אכל זה אותו דבר. :most similar sentence
2 | ולכן, צריך להביא את זה בחשבון. :most similar sentence
3 | אני לא מוכן לקבל את זה. :most similar sentence
4 | אם כן, רבותי, אנחנו עוברים להצבעה. :most similar sentence
5 | זה לא דבר שהוא חדש. :most similar sentence
6 | מה התפקיד שלכם בנושא הזה? :most similar sentence
7 | בגלל שאני אומר את האמת? :most similar sentence
8 | בכל מקרה ההצבעה לא תתקיים היום. :most similar sentence
9 | אני לא כל כך מבין. :most similar sentence
10 | איך ייתכן דבר כזה? :most similar sentence

```

אני חושב שקיבלנו תוצאות טובות מאוד כך שרוב המשפטים קיבלנו באמת משהו דומה להן, למשל משפטים 1, 2, 4, 7, 3 אלה משפטים באמת דומים אחד לשני והן מהתוצאות הכי טובות שקיבלנו, אפילו יש משפטים שלא באופן ישיר דומות למשל משפט מספר 9 לא כל כך מבין ולא כל בטוח יכולות להיות אותו דבר (לא תמיד כמובן) אבל אם משהו אומר שהוא לא מבין דבר שנאמר אז הוא לא בטוח על החלטה מסוימת או לא בטוח שהוא באמת הבין מה צריכים ממנו.

קיבלנו תוצאות כי לכל משפט המילים יש להם הסתברות גדולה להופיע ביחד ז"א שאפשר למצוא אותם כמה פעמים בקורפוס.

חלק ג':

אחרי שקראנו המילים מקובץ jsonl חלקנו המשפטים בשתי רשימות רשימה של committee ושנייה של plenary ואז בלולאה שבכל פעם לוקחת גודל chunks מבין שלוש הגדלים 1, 3, 5 כך בתחילת הלולאה נטעון את המודל דרך שימוש בפקודה Word2Vec.load שמקבלת נתיב למודל ואז נעשה כמו מה עשינו בתרגיל בית 3 שהיא חלוקה ל-chunks לפי הגודל הנוכחי ואז אחרי זה נעשה down_sample ואחרי נעשה טוקניזציה כמו מה עשינו בחלק א' שלב א ואז על רשימת הטוקניזציה נייצר sentence_embeddings בדיוק כמו שעשינו בחלק ב' סעיף ב כך sentence_embeddings זה כאילו ה-features שלנו, אחרי זה נאמן את המודל KNN ונעשה הפקודה train_test_split ואחרי זה KNN.fit ואז KNN.predict או עכשיו נוכל לשתמש ב-classification_report. התוצאות שקבלנו:

- chunk_size = 1

	precision	recall	f1-score	support
committee	0.63	0.67	0.65	4298
plenary	0.65	0.61	0.63	4298
accuracy			0.64	8596
macro avg	0.64	0.64	0.64	8596
weighted avg	0.64	0.64	0.64	8596

- chunk_size = 3 :

	precision	recall	f1-score	support
committee	0.68	0.69	0.69	5731
plenary	0.69	0.68	0.68	5731
accuracy			0.68	11462
macro avg	0.68	0.68	0.68	11462
weighted avg	0.68	0.68	0.68	11462

- chunk_size = 5 :

	precision	recall	f1-score	support
committee	0.71	0.72	0.72	6591
plenary	0.72	0.71	0.71	6590
accuracy			0.71	13181
macro avg	0.71	0.71	0.71	13181
weighted avg	0.71	0.71	0.71	13181

השאלות :

(1) התוצאות היו יותר גרועות, כך אילו התוצאות שקבלנו עכשיו (האחרון כאשר chunk size = 5):

	precision	recall	f1-score	support
committee	0.63	0.67	0.65	4298
plenary	0.65	0.61	0.63	4298
accuracy			0.64	8596
macro avg	0.64	0.64	0.64	8596
weighted avg	0.64	0.64	0.64	8596
	precision	recall	f1-score	support
committee	0.68	0.70	0.69	5731
plenary	0.69	0.67	0.68	5731
accuracy			0.69	11462
macro avg	0.69	0.69	0.69	11462
weighted avg	0.69	0.69	0.69	11462
	precision	recall	f1-score	support
committee	0.72	0.73	0.72	6591
plenary	0.72	0.71	0.72	6590
accuracy			0.72	13181
macro avg	0.72	0.72	0.72	13181
weighted avg	0.72	0.72	0.72	13181

(2) ההבדל הזה נובע לרוב בגלל שה feature vector שיש לנו הוא פחות מייצג מאשר ה feature vector שהיה לנו במשימה הקודמת.

(3) אנחנו מקבלים שגודל צ'אנק 1 הוא גרוע יותר מגודל 3 ו 3 פחות טוב מ 5. (אז 5 היה הכי טוב מבין שלושה גדלים אלו) הסיבה היא שאם גודל הצ'אנק קטן, מטריצת התכונות תהיה גדולה יותר, מה שגורם לזמן הריצה לגדול עוד יותר. בנוסף, כל צ'אנק יכול תכונות שלא מייצגות את הסוגים בצורה טובה מספיק, מכיוון שכוח features יתפזר בין הצ'אנקים הקטנים ולא יהיה מרוכז בצ'אנק אחד, מה שמקשה על משימת הסיווג.

חלק ד':

אחרי שקראנו את המשפטים מהקובץ masked_sentences.txt הגדרנו

```
tokenizer = AutoTokenizer.from_pretrained('dicta-il/dictabert')
```

```
model = AutoModelForMaskedLM.from_pretrained('dicta-il/dictabert')
```

כפי ראינו באתר ואז עבור כל משפט בקובץ הנתון אחרי שהחלפנו כל [*] ב-[MASK] עברנו על כל [MASK] במשפט הנוכחי השתמשנו ב-model(tokenizer.encode-ב) שמשמשת במודל כדי לבצע החייו עבור המשפט ואז נשמור התוצאה מזה ונחשב האינדקס של ה-[MASK] ואז ב-torch.topk שנקבל ממנה הטוקן הסביר ביותר שיחליף [MASK] ואז נשאר להחליף המזהה לטוקן מילולי המתאים ואז בזה נקבל הטוקן שיחליף [MASK]. אז נשמור אותו ברשומת הטוקנים וגם נחליף אותו עם ה-[MASK].

השאלות:

(1) התוצאות הגיניות מאוד גם מבחינה תחבירית גם מבחינת משמעות המשפט. (התוצאות בתחילת העמוד הבא)

1	Original sentence: יש צורך [*] ביצירת מקומות עבודה רבים יותר, בתשלום שכר [*] בעד העבודה.
2	DictaBERT sentence: יש צורך גם ביצירת מקומות עבודה רבים יותר, בתשלום שכר גבוה בעד העבודה.
3	DictaBERT tokens: גם,גבוה
4	
5	Original sentence: אבל הנושא הוא [*] אם אתה בעד [*] ההתנתקות או נגד.
6	DictaBERT sentence: אבל הנושא הוא לא אם אתה בעד תוכנית ההתנתקות או נגד.
7	DictaBERT tokens: לא,תוכנית
8	
9	Original sentence: [*] לכן, הממשלה מתנגדת [*] החוק הזאת, [*] להביא לפתרון הבעיה בדרך.
10	DictaBERT sentence: לכן, הממשלה מתנגדת להצעת החוק הזאת, כדי להביא לפתרון הבעיה בדרך אחרת.
11	DictaBERT tokens: להצעת,כדי,אחרת
12	
13	Original sentence: ההצמדה [*] [*] לא הוראת שעה.
14	DictaBERT sentence: ההצמדה למדד היא לא הוראת שעה.
15	DictaBERT tokens: למדד,היא
16	
17	Original sentence: עזבי את [*] 84, אני לא מכיר את חוק [*] והבנייה.
18	DictaBERT sentence: עזבי את סעיף 84, אני לא מכיר את חוק התכנון והבנייה.
19	DictaBERT tokens: סעיף,התכנון
20	
21	Original sentence: [*] עלות החוק [*] מיליארד.
22	DictaBERT sentence: עלות החוק היא מיליארד שקל.
23	DictaBERT tokens: היא,שקל
24	
25	Original sentence: מסיבות שונות, [*] 1,275 התיקים שנפתחו בשנת [*] - טרם הוגשו כתבי אישום.
26	DictaBERT sentence: מסיבות שונות, מתוך 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום.
27	DictaBERT tokens: מתוך,שנפתחו
28	

אבל היה לנו גם תוצאות שהן פחות טובות ופחות הגיוניות:

Original sentence: מסיבות שונות, [*] 1,275 התיקים שנפתחו בשנת [*] - טרם הוגשו כתבי אישום.
DictaBERT sentence: מסיבות שונות, מתוך 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום.
DictaBERT tokens: מתוך,שנפתחו

2) כן כמו רואים היה שיפור בתוצאות כך עכשיו לרוב קבלנו תוצאות שיותר מתאימות ממה קבלנו בתרגיל בית 2 אבל גם יש משפטים שבהם בתרגיל בית 2 קבלנו תוצאה יותר טובה אבל זה במעט מהמקרים, אז בכללי לרוב כן קבלנו שיפור בתוצאות.

3) כן נקבל אותו הדבר:

dictabert_results2.txt	dictabert_results.txt
1 Original sentence: יש צורך [*] ביצירת מקומות עבודה רבים יותר, בתשלום שכר [*] בעד העבודה.	1 Original sentence: יש צורך גם ביצירת מקומות עבודה רבים יותר, בתשלום שכר גבוה בעד העבודה.
2 DictaBERT sentence: יש צורך גם ביצירת מקומות עבודה רבים יותר, בתשלום שכר גבוה בעד העבודה.	2 DictaBERT sentence: יש צורך גם ביצירת מקומות עבודה רבים יותר, בתשלום שכר גבוה בעד העבודה.
3 DictaBERT tokens: גם,גבוה	3 DictaBERT tokens: גם,גבוה
4	4
5 Original sentence: אבל הנושא הוא [*] אם אתה בעד [*] ההתנתקות או נגד.	5 Original sentence: אבל הנושא הוא [*] אם אתה בעד [*] ההתנתקות או נגד.
6 DictaBERT sentence: אבל הנושא הוא לא אם אתה בעד תוכנית ההתנתקות או נגד.	6 DictaBERT sentence: אבל הנושא הוא לא אם אתה בעד תוכנית ההתנתקות או נגד.
7 DictaBERT tokens: לא,תוכנית	7 DictaBERT tokens: לא,תוכנית
8	8
9 Original sentence: [*] לכן, הממשלה מתנגדת [*] החוק הזאת, [*] להביא לפתרון הבעיה בדרך.	9 Original sentence: [*] לכן, הממשלה מתנגדת [*] החוק הזאת, [*] להביא לפתרון הבעיה בדרך.
10 DictaBERT sentence: לכן, הממשלה מתנגדת להצעת החוק הזאת, כדי להביא לפתרון הבעיה בדרך אחרת.	10 DictaBERT sentence: לכן, הממשלה מתנגדת להצעת החוק הזאת, כדי להביא לפתרון הבעיה בדרך אחרת.
11 DictaBERT tokens: להצעת,כדי,אחרת	11 DictaBERT tokens: להצעת,כדי,אחרת
12	12
13 Original sentence: ההצמדה [*] [*] לא הוראת שעה.	13 Original sentence: ההצמדה [*] [*] לא הוראת שעה.
14 DictaBERT sentence: ההצמדה למדד היא לא הוראת שעה.	14 DictaBERT sentence: ההצמדה למדד היא לא הוראת שעה.
15 DictaBERT tokens: למדד,היא	15 DictaBERT tokens: למדד,היא
16	16
17 Original sentence: עזבי את [*] 84, אני לא מכיר את חוק [*] והבנייה.	17 Original sentence: עזבי את [*] 84, אני לא מכיר את חוק [*] והבנייה.
18 DictaBERT sentence: עזבי את סעיף 84, אני לא מכיר את חוק התכנון והבנייה.	18 DictaBERT sentence: עזבי את סעיף 84, אני לא מכיר את חוק התכנון והבנייה.
19 DictaBERT tokens: סעיף,התכנון	19 DictaBERT tokens: סעיף,התכנון
20	20
21 Original sentence: [*] עלות החוק [*] מיליארד.	21 Original sentence: [*] עלות החוק [*] מיליארד.
22 DictaBERT sentence: עלות החוק היא מיליארד שקל.	22 DictaBERT sentence: עלות החוק היא מיליארד שקל.
23 DictaBERT tokens: היא,שקל	23 DictaBERT tokens: היא,שקל
24	24
25 Original sentence: מסיבות שונות, [*] 1,275 התיקים שנפתחו בשנת [*] - טרם הוגשו כתבי אישום.	25 Original sentence: מסיבות שונות, מתוך 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום.
26 DictaBERT sentence: מסיבות שונות, מתוך 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום.	26 DictaBERT sentence: מסיבות שונות, מתוך 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום.
27 DictaBERT tokens: מתוך,שנפתחו	27 DictaBERT tokens: מתוך,שנפתחו
28	28
29 Original sentence: לך, [*], על השוברת המקיפה [*].	29 Original sentence: לך, [*], על השוברת המקיפה [*].
30 DictaBERT sentence: לך, [*], על השוברת המקיפה [*].	30 DictaBERT sentence: לך, [*], על השוברת המקיפה [*].
31 DictaBERT tokens: לך,מקיפה	31 DictaBERT tokens: לך,מקיפה
32	32

זה משום ש מודל BERT אינו מודל רנדומלי ז"א אותו קלט עם אותם פרמטריים תמיד יניבו לנו תוצאות דומות וזה באמת מה שמקבלים.

4) כן כמו שהראינו בסעיף 1:

```
Original sentence: מסיבות שונות , [*] 1,275 התיקים שנפתחו בשנת [*] - טרם הוגשו כתבי אישום .  
DictaBERT sentence: מסיבות שונות , מתוך 1,275 התיקים שנפתחו בשנת שנפתחו - טרם הוגשו כתבי אישום .  
DictaBERT tokens: מתוך,שנפתחו
```

למשל המשפט הזה, אז כנראה המודל הוא פחות טוב בחיזוי מספרים כמו שנה חודש או יום זה יכול גם להיות כי זה יכול איזה שנה והוא לא יכול איזה בדיוק לכן גם אם נתן מספר אולי יהיה לא נכון כי יש אינסוף אפשרויות.(זה למרות שחיזה טוקן תקין אבל הוא לא נכון).

```
Original sentence: [*] בסופו של דבר , גם אני הייתי בשתי אליפויות העולם בצום שבע שנים , וכן יצאנו עם דגל ישראל ואני יודעת שחזר מאימא שלי , אי אוד לא יודע שכאמך אנתוני מייצגים את ישראל .  
DictaBERT sentence: " בסופו של דבר , גם אני הייתי בשתי אליפויות העולם בצום שבע שנים , וכן יצאנו עם דגל ישראל ואני יודעת שחזר מאימא שלי , אי אוד לא יודע שכאמך אנתוני מייצגים את ישראל .  
DictaBERT tokens: "בשתי,אודי"
```

בדוגמה זו היה לני גם את ' ' ' שהיה בלי הסוגר שלו. שזה לא נכון מבחינת כתיבה.