

## תרגיל 2

יוסי גואטה 032528267  
נטע זינגר 201111648

### חלק א

סעיף	gold	tain
1	11282	127884
2	3171	15986
3	11282	127884
4	36	37
5	1.07978555661	1.1349305642

6

- ככל שמדד העמימות עולה כך סיכויי המתייג פוחתים לתת חיזוי טוב
- מסעיף 4 ניתן לראות שכמות המזהים שונה מה שמצביע על מזהים שונים בין קובץ הgold והtrain , דבר אשר עלול לפגוע בביצועי המתייג שכן אנו עלולים לתת משקל יתר לתיוגים אשר אינם "קיימים"
- 

### חלק ב

- הגדירו במדויק את סכמת הפרמטרים במודל ? קובץ טקסטואלית בסיומת "baseline" אשר מכיל בכל שורה סגמנט ואת התיוג הנפוץ עבורו בקורפוס. `<tag> " " <segment>`

```
1  ALMLA IN
2  HWFG VB
3  HXLJWJ NN
4  HWLLIM JJ
5  1,776 CD
6  ATRIM NN
7  ABZTRIM NN
8  FQWPIM JJ
```

- הגדירו במפורש את נוסחאות המשערכים של הפרמטים במודל ?

$$T(seg) = \underset{tag}{\operatorname{argmax}} count(tag_{seg})$$

הגדירו במדויק את סיבוכיות זמן הריצה של המודל?

עבור מספר המילים num\_of\_words ומספר התגים num\_of\_tag, מספר המשפטים num\_of\_sentences

• אימון:

$O(\text{num\_of\_words} + \text{num\_of\_tag})$

• תיוג:

$O(\text{num\_of\_sentences})$

דווחו מהו הדיוק הכולל עבור תיוג קובץ הבדיקה

macro avg word - 0.830969686226

macro avg sentence - 0.106

## חלק ג

מהי פונקציית המטרה של המודל ?

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(t_1^n | w_1^n) \approx \underset{t_1^n}{\operatorname{argmax}} \prod_{i=1}^n \overbrace{P(w_i | t_i)}^{\text{emission}} \overbrace{P(t_i | t_{i-1})}^{\text{transition}}$$

מהן נוסחאות הפרמטרים של המודל ?

$|V| = \text{size of corpus}$

$$P(t_i) = \frac{\text{count}(t_i)}{|V|}$$

$$P(t_i, t_{i-1}) = \frac{\text{count}((t_i, t_{i-1})) + \delta}{\text{count}(t_i) + \delta|V|}$$

$$P(w_i) = \left[ \frac{\text{count}((w_1, t_1)) + \delta}{\text{count}(t_1) + \delta|V|}, \dots, \frac{\text{count}((w_n, t_n)) + \delta}{\text{count}(t_n) + \delta|V|} \right]$$

מהן נוסחאות המשערכים של הפרמטרים במודל?

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

$$P(w_i | t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

מהי סיבוכיות זמן הריצה של האלגוריתם האימון?  
עבור T מספר המשפטים Si מספר הסגמנטים

$$O\left(\binom{|S|}{2} \times T\right) = |S|^2 \times T$$

מהי סיבוכיות זמן הריצה של אלגוריתם התיג?  
עבור T מספר המשפטים Si מספר הסגמנטים

$$O(T \times |S|^2)$$

מהו הדיוק macro avg עבור קובץ הבדיקה?  
macro avg word - 0.854192519057 •  
macro avg sentence - 0.112 •

במקום לציין כחח ממשו טכניקת החלקה למילים לא ידועות האם תוצאותיכם השתפרו?  
macro avg word - 0.858890267683 •  
macro avg sentence - 0.12 •  
כפי שניתן לראות התוצאות השתפרו במעט

ממשו בנוסף טכניקת החלקה למעברים לא ידועים האם תוצאותיכם השתפרו ביחס לסעיף הקודם?

$$P(t_i, t_{i-1}) = \frac{\text{count}(t_i, t_{i-1}) + \delta}{\text{count}(t_i) + \delta|V|}$$

כפי שניתן לראות התוצאות השתפרו במעט

## חלק ד

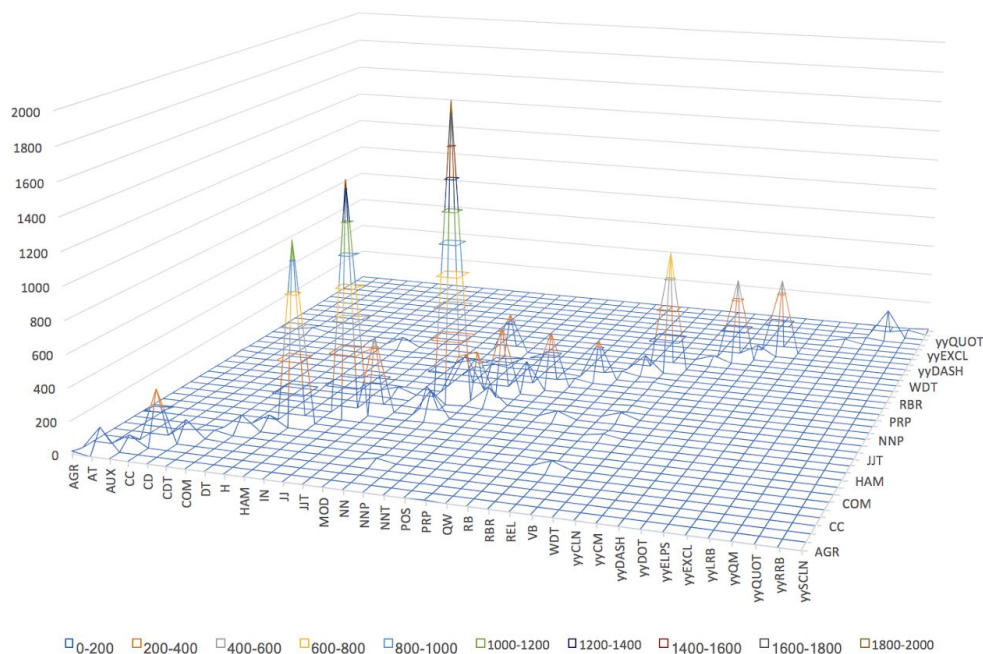
עבור המודל הטוב ביותר שלך חשב את מטריצה הבלבול מהן 3 השגיאות הנפוצות ביותר לפי המטריצה?  
עבור התג NN, התגים שקיבלו את ערך השגיאה המקסימלי הם

תג שצפינו	תג שקיבלנו	כמות
NN	VB	254
NN	NNT	221
NN	NNP	245

## מטריצת הבילבול:

	AGR	AT	AUX	CC	CD	CDT	COM	DT	H	HAM	IN	J	JT	MOD	NN	NNP	NNT	POS	PRP	QW	RB	RBR	REL	VB	WDT	WYCLN	WYCM	WYDASH	WYELPS	WYEXCL	WYLRB	WYQCM	WYQUOT	WYRRB	WYSCLN
AGR	19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
AT	0	19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
AUX	0	0	67	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
CC	0	0	0	325	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
CD	0	0	0	0	1127	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
CDT	0	0	0	0	0	180	21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
COM	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
DT	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
HAM	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
IN	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
J	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
JT	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
MOD	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
NN	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
NNP	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
NNT	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
POS	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
PRP	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
QW	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
RB	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
RBR	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
REL	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
VB	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WDT	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYCLN	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYCM	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYDASH	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYELPS	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYEXCL	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYLRB	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WYQCM	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				

## מטריצת הבילבול



תיגו את המשפט הבא ידנית ולאחר מכן תייגו אותו באמצעות המתייג שכתבתם ?

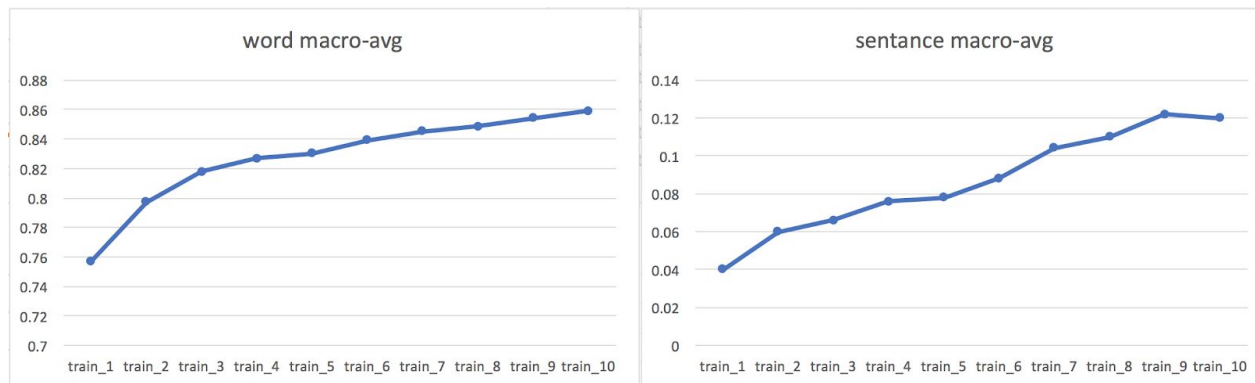
סוג המתייג	אישה	נעלה	נעלה	נעלה שלה	נעלה	את	ה דלת	בפני	בעל שלה
אנושי	N	adj	VB	N	VB	conj	N	conj	adj
מכונה	NN	NN	NN	NN	NN	AT	NN	IN	NN

1	AFH	NN
2	NELH	NN
3	NELH	NN
4	NELH	NN
5	yyCM	yyCM
6	NELH	NN
7	AT	AT
8	HDLT	NN
9	BPNI	IN
10	BELH	NN
11		

ניתן לראות כי המתייג דיי טעה בתיוג השפה ניתן להבין זאת , מדובר במשפט עם עמימות לשונים רבה.

ביצועי המתייג כתלות בגודל קובץ האימון

tain batch	word macro-avg	sentance macro-avg
train_1	0.756957986	0.04
train_2	0.797110441	0.06
train_3	0.817851445	0.066
train_4	0.826803758	0.076
train_5	0.830083319	0.078
train_6	0.839035632	0.088
train_7	0.844974295	0.104
train_8	0.848519766	0.11
train_9	0.854015246	0.122
train_10	0.858890268	0.12



**מה יקרה להערכתכם אם יתווספו עוד טקסטים מתוייגים לאימון?**  
 ניתן לראות שהנגזרת שואפת ל-0 ולכן אנו נראה שיש שיפור ממזערי והתכנסות לאזור ה-86 אחוזי הצלחה

## חלק ו

### תנו תיאור סכמטי כולל של הקוד שמימשתם ?

src/Baseline.py - מכיל את כל פונקציות העזר למימוש מתייג baseline (חלק ב) המצמיד לסגמנט את התג השכיח ביותר עבור אותו סגמנט בקורפוס  
 src/parse\_data.py - קובץ המכיל פונקציות עזר לקריאה וכתובת קבצים  
 src/const.py - קובץ המכיל את הקבועים במערכת  
 src/utils.py - קובץ המכיל פונקציית עזר עבור חלק ד הכולל חישוב מטריצות בלבול וחלוקה של הקורפוס  
 src/viterbi.py - קובץ המכיל את כל פונקציות העזר למימוש מתייג BI GRAM באמצעות אלגוריתם viterbi  
 src/ex1/\*.py - קבצים להרצת חלק א בתרגיל

### תארו ממעוף הציפור את התוצאות והמסקנות מהתרגיל ?

ניתן לראות לפי תוצאות התיוג שאף על פי שנעשה שימוש באלגוריתם מתקדם (viterbi המבוסס תכנון דינאמי) ישנו גבול יכולת חיזוי לתיוג חלקי המשפט, על פי המסקנות ככל הנראה נצטרך להוסיף פיצרים נוספים כדי לחדד את יכולת החיזוי.  
 בנוסף ראינו כי שימוש בטכניקת ההחלקה משפרת את תוצאות המתייג, לו היינו ממשים החלקה מתקדמת יותר היינו מקבלים תוצאות טובות יותר.

**פרט נסו לענות על השאלות המחקריות לגבי המתייג בעברית שהוצגו בהקדמה לתרגיל והציעו בכתב דרכים אפשריות לטיוב ושיפור עתידי של המתייג שלכם**

- שימוש בהחלקה טובה יותר כגון: good turing
- שימוש בטריאגם כשמבצעים את שהשערוך באלגוריתם viterbi
- בחינת משקלים שונים עבור emission ו-transition

- שימוש באלגוריתמים שונים כגון רשת ניוונים, לוג לינארי מודל
- מציאת גודל קורפוס אופטומלי אשר יתר תוצאות טובות מצד אחד ומצד שני לא יכביד על משך זמן האימון