

# תרגיל בית 1

שיטות בעיבוד שפה טבעית

עומר כהן  
נמרוד אדמוני

308428127

203860721

## אימון:

סוגי המאפיינים בהם השתמשנו הינם :

מודל	מודל 1 (גדול)	מודל 2 (קטן)
מאפיינים	$f_{100} - f_{107}$ $f_{108} - \text{capital letter feature}$ $f_{109} - \text{only number feature}$	$f_{100} - f_{107}$

כאשר המאפיינים  $f_{108}, f_{109}$  ממומשים כך :

$f_{108}$  הוא מאפיין שמזהה שמילה מסוימת היא מספר :

$$f_{108} = \begin{cases} 1 & \text{if current word } w_i \text{ is a Number and tag is } t = CD \\ 0 & \text{otherwise} \end{cases}$$

$f_{109}$  הוא מאפיין שמזהה שהמילה הנוכחית מתחילה באות גדולה :

$$f_{109} = \begin{cases} 1 & \text{if current word } w_i \text{ begins with } T \text{ and } t = NN \\ 0 & \text{otherwise} \end{cases}$$

$$f_{110} = \begin{cases} 1 & \langle pp_{word}, p_{word}, c_{word}, c_{tag} \rangle = \langle w_1, w_2, w_3, tag \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$f_{111} = \begin{cases} 1 & \langle suffix, c_{word}, c_{tag} \rangle = \langle s, w_3, tag \rangle \\ 0 & \text{otherwise} \end{cases}$$

$$f_{112} = \begin{cases} 1 & \text{if word } w_i \text{ length is } > 4 \text{ and tag is } t \\ 0 & \text{otherwise} \end{cases}$$

$$f_{113} = \begin{cases} 1 & \text{if current word } w_i \text{ contains punctuation and tag is } t \\ 0 & \text{otherwise} \end{cases}$$

מספר המאפיינים בכל מודל הם :

מודל 1 (גדול)	מודל 2 (קטן)
f100 15415	75 f100
f101 13265	333 f101
f102 22393	349 f102
f103 8150	118 f103
f104 1060	107 f104
f105 44	21 f105
f106 32132	70 f106
f107 30793	66 f107
f108 46	22 f108
f109 76	29 f109
	4 f110
	57 f111
	34 f112
	24 f113

השיפורים שהוספנו למודלים הם :

מודל 1 – אותיות גדולות וזיהוי מספרים

מודל 2 – בנוסף לזיהוי אותיות גדולות ומספרים גם את משפחות המאפיינים 110-113.

זמן האימון עבור מודל 1 הוא 5 דקות עם אחוז דיוק של על קובץ המבחן . עבור מודל 2 זמן האימון הוא שניות בודדות ואחוז הדיוק המשוערך הוא 83.8% , בפועל המודל מבצע אופטימיזציה על קובץ האימון ל 88% שני המודלים אומנו על 2021 Macbook pro 13'' עם ערכת שבבים M1.

#### הסקה:

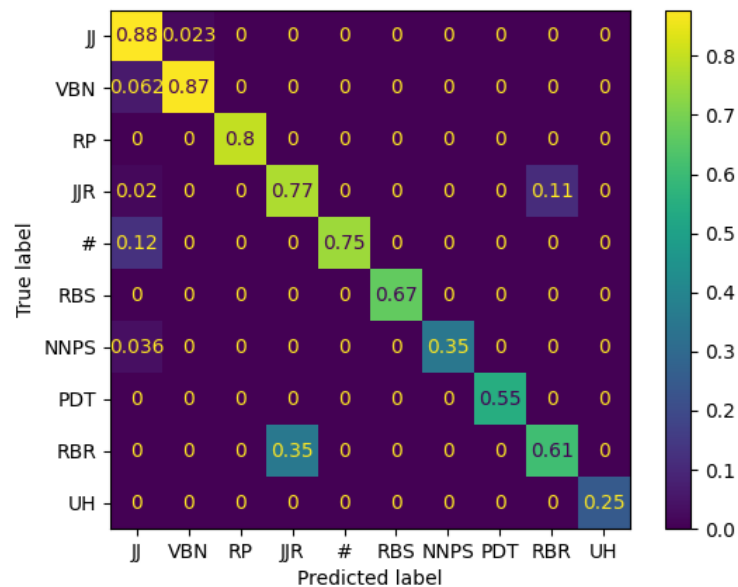
ההסקה התבצעה על אלגוריתם *Viterbi* עם  $beam = 2$  עבור מודל 1 ו  $beam =$  עבור מודל 2 ללא שינויים.

#### מבחן:

עבור מודל 1 אחוז הדיוק עבור קובץ המבחן *test1.wtag* הוא 95.8% :

```
you have 123374 features!
123374
f100 15415
f101 13265
f102 22393
f103 8150
f104 1060
f105 44
f106 32132
f107 30793
f108 46
f109 76
the following tags the model confused the most
JJ confused 12.224448897795593% of the time.
VBN confused 12.650602409638555% of the time.
JJR confused 23.0% of the time.
RBS confused 33.333333333333336% of the time.
# confused 25.0% of the time.
RP confused 20.408163265306122% of the time.
RBR confused 39.21568627450981% of the time.
PDT confused 45.45454545454546% of the time.
NNPS confused 65.06024096385542% of the time.
UH confused 75.0% of the time.
Accuracy is: 0.9589000591366056
```

קיבלנו עבור המבחן את ה *Confusion – Matrix* הבאה :



ניתן לראות כי המודל מתבלבל בהסתברות גבוהה בין התיוגים JJR ו-RBR, נוסף מאפיין שמחזיק את המילה, סיומת בת 2/3 אותיות והתיוג וכך נוכל להבדיל בין תיוג JJR שנגמרים ב'-er' לבין תיוג RBR שנגמרים ב'-ly'.

עבור מודל 2 נרצה להגדיל את ה-lambda כדי להגדיל רגולריזציה ולתת עדיפות ל - varians נמוך על פני bias נמוך. בנוסף, נרצה גם להגדיל את ערך הסף עבור המאפיינים כדי להשאיר את המאפיינים הדומיננטיים ולהימנע מ-overfit על ה-dataset. כדי להתמודד עם המחסור ב-data להערכת ביצועי המודל, פתרון אפשרי הוא להקצות חלק מה-data שברשותנו ל-validation. לדוגמא, לחלק את ה-dataset ל-80% אימון ו-20% validation. נרצה לאמן מספר פעמים את המודל עם הדאטה המחולק כך שכל פעם התוכן בכל חלק יהיה שונה, ולבצע ממוצע של הערכת הביצועים כדי לקבל שיערוך של ביצועי המודל. לבסוף המודל האמיתי יתאמן על כל ה-dataset.

#### תחרות:

עבור שני המודלים אנחנו מצפים לקבל אחוזי דיוק כמו על קבצי testn מכיון שאת שניהם לא ראינו באימון. המודלים שהשתמשנו להם זהים למודלים שדווחו באימון. אם זאת במהלך האימון של המודל נקודת ההשוואה שלנו היא קובץ ה-testn והשינויים שעשינו תאמו לתוצאות שקיבלנו עליו, לכן הוא למעשה evaluation.

#### חלוקת עבודה:

כל העבודה על תרגיל הבית נעשתה במשותף.