

תרגיל בית 1 - רטוב - דו"ח

מגישים- אליעזר שלומי 315227686, יקטרינה סופרפין 321170383

הסבר על סט המאפיינים:

השתמשנו בפיצ'רים $f_{100} - f_{105}$ של רטנפרקי והפיצ'רים הנוספים של המספרים והמילים המכילות אותיות גדולות- וזאת כנדרש בתנאי התרגיל. בנוסף הוספנו גם את הפיצ'רים $f_{106} - f_{107}$. נפרט על כל פיצ'ר:

כמות באימון מודל 2	כמות באימון מודל 1	פיצ'רים שנדרשנו לממש:
1805	15415	f_{100} - פיצ'ר של התאמה של תיוג למילה.
2556	13265	f_{102} - פיצ'ר של התאמת תיוג למילה ותחילית מסויימת של המילה.
3978	22393	f_{101} - פיצ'ר של התאמת תיוג למילה וסיומת מסויימת של המילה.
1168	8150	f_{103} - פיצ'ר טריגרם.
316	1060	f_{104} - פיצ'ר בייגרם.
32	44	f_{105} - פיצ'ר יוניגרם.
180	1248	$feature_words_with_numbers$ - פיצ'ר של הופעת מספר
420	397	$feature_words_with_capital_letters$ - פיצ'ר של מילים אותיות גדולות
		פיצ'רים שהוספנו:
3314	38164	f_{106} - פיצ'ר של התאמת תיוג לפי מילה קודמת.
3141	35613	f_{107} - פיצ'ר של התאמת תיוג לפי מילה הבאה
219	1018	feature words with capital hyphens
6	22	Feature words with capitals only
420	36386	Feature words with capital letters only
2928	36386	$\langle next\ next\ word, tag \rangle$ feature
3065	38050	$\langle pre\ pre\ word, tag \rangle$ feature
5007	82798	$\langle pre\ pre\ word, pre\ word, tag \rangle$
5007	82798	$\langle next\ next\ word, next\ word, tag \rangle$
2826	33511	$\langle pre\ pre\ word, pre\ tag, tag \rangle$ feature
1736	14719	feature 100 lower - כמו 100 רק באותיות קטנות
2467	12608	feature 101 lower - כמו 101 רק באותיות קטנות
5574	42044	feature 101 lower complete - כמו 101 רק אותיות קטנות וניקה את המשלים (הצד השני של המילה)
3593	20156	feature 102 lower - כמו 102 רק אותיות קטנות
5226	35049	feature 102 lower complete - כמו 102 רק אותיות וניקה את המשלים
16	35	Feature upper tag exist - $\langle tag \rangle$ אם יש $capital\ letter$.
7	15	Feature hyphen exist - $\langle tag \rangle$ אם יש מקף במילה.

4	6	Feature digit exist - <tag> אם יש ספרה במילה.
466	1025	מבנה מילה עבור המילה הנוכחית XXxxxdd - X אות גדולה, x אות קטנה ו - d זאת ספרה
772	1025	מבנה מילה עבור המילה הקודמת XXxxxdd - X אות גדולה, x אות קטנה ו - d זאת ספרה
466	2472	מבנה מילה עבור המילה הבאה XXxxxdd - X אות גדולה, x אות קטנה ו - d זאת ספרה

שלב האימון:

מפרט בסיסי של התוכנה עליה הרצנו מכונת Azure – Standard D4s v3 (4 vcpus, 16 GiB memory)

בשביל לאמן השתמשנו בפונקציה של scipy הנקראת fmin_l_bfgs_b שהשתמשה בפונקציה שכתבנו שמחשבת גרדיאנט ולוס.

מודל 1 – יצירת הפיצ'רים לקח 3 דקות, האימון עוד 5-7 דקות, חישוב הדיוק על הטסט 3 דקות והתיוג לקח גם 2-3 דקות.

- מדד ה-Accuracy- 0.96321

```

NN      10    26    2    19    0    8
Accuracy of train 1 test is  0.96321

Process finished with exit code 0

```

- זמן האימון- 315 איטרציות כאשר איטרציה לוקחת בערך שנייה.

- מקדם רגולריזציה- 1.

- Threshold- 1 – חשבנו שאולי יהיה לנו אוברפיט אך סט הולידציה מראה

שימוש ב threshold=1 מביא לתוצאות טובות.

- באימון מודל זה השתמשנו בכל הפיצ'רים המצוינים בטבלה. ניתן לראות בטבלה שלעיל את מספר הפיצ'רים מכל סוג ששימשו לאימון מודל 1 ו 2.

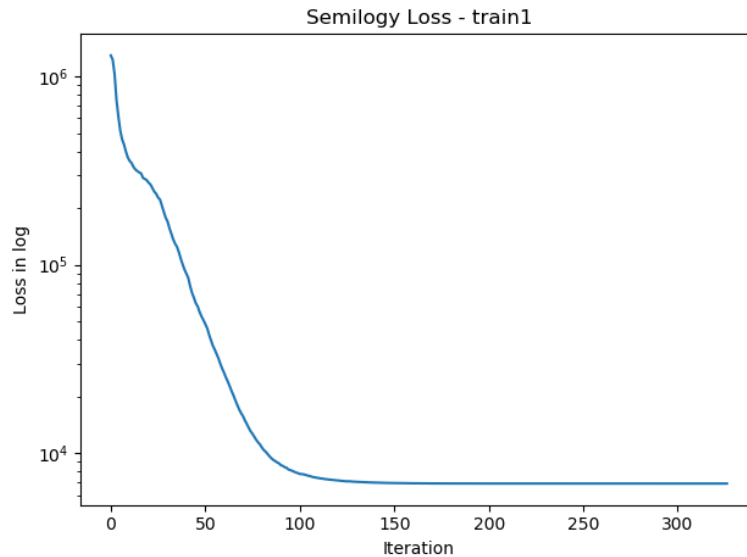
- בשביל מימוש יעיל השתמשנו בפעולות וקטוריות ומטריצה sparse של scipy. משתנים שיכולנו לחשב פעם אחת

לפני האימון חושבו פעם אחת והוכנסו כפרמטר (empirical count).

בהניתן ה Confusion Matrix ניתן לראות שיש בלבול רב בין ה JJ ו - NN. הייתי מנסה לתת למודל כוח להבדיל ביניהם ע"י זה שהייתי מוסיף פיצ'ר של תחילת המשפט. לדעתי פעמים רבות אנו מזכירים שם של מילה בתחילת המשפט ובהמשך מוסיפים עליו תארים, אולי בדרך זו ניתן יהיה להבין מתי אנחנו מתכוונים לשם עצם ומתי לשם תואר. עוד רעיון שחשבתי עליו (שהוא מתקשר בעיקר למילים לא מוכרות ויש לו קשר גם לטבלה) יהיה לנסות לעשות stemming או משהו דומה כדי שאם המודל יראה מילה לא מוכרת הוא ימצא את הגרעין שלה (למקרה שהמודל לא הכיר את המילה רק בגלל שהיא וראיציה של מילה אחרת) והוא יוכל לנחש מה התפקיד שלה בהינתן סיומת/גרעין וכו'.

Confusion Matrix

	VB	VBG	VBD	NNP	NNPS	RB	VRN	IN	JJ	NN
VB	554	0	1	1	0	1	0	1	2	15
VBG	0	368	0	2	0	0	0	0	1	29
VBD	1	0	795	0	0	0	34	0	0	2
NNP	0	0	0	1940	18	0	0	0	10	12
NNPS	0	0	0	29	32	0	0	0	4	0
RB	1	0	0	0	0	732	0	17	16	6
VRN	0	0	31	0	0	0	434	0	30	3
IN	0	1	0	2	0	48	0	2492	1	1
JJ	3	6	6	8	0	11	27	2	1367	65
NN	10	26	2	19	0	8	1	3	79	3173



מודל 2 – יצירת הפיצ'רים לקח דקה, 2-3 דקות לאימון ועוד 1-2 דקות לתיג התחרות.

- זמן האימון- 97 איטרציות
- מקדם רגולריזציה- 1.0
- Threshold=1
- באימון מודל זה השתמשנו באותם פיצ'רים כמו במודל 1
- ניתן לראות בטבלה שלעיל את מספר הפיצ'רים מכל משפחת פיצ'רים ששימשו לאימון מודל 2.

שלב ההסקה:

את שלב ההסקה ביצענו ע"י אלגוריתם Viterbin, תוך שימוש ב-

beam search. ביצענו את האלגוריתם ללא חריגות ושינויים דרסטיים מעבר לשימוש ב-*beam search*. לצורך שיפור זמן הריצה נעשה שימוש בחישובים וקטוריים ומטריצות ספרסיות.

תחרות:

- הוספנו הרבה מאוד פיצ'רים, יש לנו בסה"כ 543,060 פיצ'רים ו 28 משפחות פיצ'רים בסה"כ.
- ביצענו אופטימיזציה על ה hyperparameters כמו ה threshold ומקדם הרגולריזציה.
- קראנו מאמרים כדי לחשוב אילו פיצ'רים כדאי להוסיף (למשל 97% to 100%).
- בדקנו את הפיצ'רים שהוספנו כל פעם והאם הם באמת עוזרים למודל.

חלוקת עבודה:

החשיבה התבצעה ביחד אך בכתובת הקוד והמימוש אליעזר עשה את שלב האימון, את האופטימיזציה והפיצ'רים הבסיסיים כדי לעבור את אחוזי הדיוק הבסיסיים ואילו קטיה מימשה את ה Viterbi והוסיפה משפחות פיצ'רים למודל כדי לשפר את המודל לתחרות ולהגיע למקומות גבוהים. החשיבה, הרעיונות וחלוקת העבודה התבצעו יחד.