

תרגיל בית 1 – קורפוסים

חיים שאללשוילי 200832780

צחי כפיר 200681476

חלוקה למשפטים

האלגוריתם הנבחר לחלוקת משפטים:

כל משפט יסתיים במחרוזת המתאימה לביטוי הרגולרי הבא: $[!][?][*]["]$. בתנאי והמחרוזת אינה נקודה המוקפת בספרות משני צדיה או נקודה המוקפת באותיות משני צדיה. ביטוי זה יפריד כל משפט המסתיים בנקודה, סימן קריאה וסימן שאלה, משפט המסתיים ברצף הבנוי מאותם תווים, ומשפט המסתיים ברצף המסתיים באחד מהתווים הללו ואחריהם מירכאות כפולות.

הסיבה לכך היא שבהבאת דיבור ישיר בין מירכאות מצטרפות שתי מערכות פיסוק: הפיסוק של הטקסט המצוטט והפיסוק של הכותב. ראוי למנוע גיבוב של סימני פיסוק רצופים, ועל כן יש לוותר על סימן פיסוק השייך לאחת המערכות. כדי לבצע זאת סימן הפיסוק יבוא לפני המירכאות למשל:

הוא פרש זרועותיו וקרא: "ריבוננו של עולם!"

קשיים:

- טיפול בתו "; – סימן זה יכול להגיע במקום נקודה כאשר ישנו צורך להפריד בין שני משפטים שיש ביניהם קשר ענייני אמיץ ובמרום פסיק כאשר יש לפרט עניינים ין האיברים השונים בחלקי המחרוזת. בשל כפל המשמעות הוחלט לא להפריד משפטים על פי תו זה בשלב זה כאשר אנו לא משתמשים בידע מורפולוגי וסמנטי על הטקסט.
- טיפול בתו " – במקרים מסוימים סימן זה מופיע לפני ציטוט או דיבור ישיר ובמקרה זה אין צורך לרדת שורה. במקרים אחרים שבהם התו נקודתיים מגיע לפני פירוט של נקודות נרצה להפריד למשפט חדש. בשל כפל המשמעות הוחלט לא להפריד משפטים על פי תו זה בשלב זה כאשר אנו לא משתמשים בידע מורפולוגי וסמנטי על הטקסט.
- החלטה זו לא פוגמת ברוב הטקסטים שכן טקסטים שבהם קיימים רשימות של נקודות קיימים כבר סימני ירידת שורה שאנו לא מוחקים באלגוריתם שלנו. לכן בכל מקרה תתבצע חלוקה למשפט חדש.
- כן האלגוריתם לקח בחשבון כי בטקסט המקורי קיימת הפרדה בסיסית של רווח בין משפט למשפט. שכן במקרה ולא, קיימת האפשרות כי נגנוב תווים מהמשפט הבא לדוגמא מירכאות כפולות.
- סימנים שונים של ירידת שורה עבור כתבנים שונים

טוקניזציה

תהליך הטוקניזציה ממוש ע"י הוספה של רווח (') לפני ואחרי כל תו שאיננו אות, פרט למיקרים הבאים:

- כאשר התו '-' מופיע בין שתי אותיות. לדוגמא: בצירופים תל-אביב, בית-ספר אין צורך להוסיף רווחים, אך במידה ומופיע הטקסט הבא: "מקרה ראשון-" נרצה להפריד את "ראשון" ואת "-".
 - כאשר התו "י" מופיע בין שתי אותיות. לדוגמא: בראשי התיבות צה"ל, חז"ל אין צורך להוסיף רווחים, אך במידה ומופיע ציטוט, דני אמר לי: "מסור לי את הכדור", נרצה להפריד בין ה-"י" "מסור" "כדור" "י".
 - כאשר התו "י" מופיע בין שתי אותיות. לדוגמא: במילים ג'יפ, ג'ירפה, אין צורך להוסיף רווחים, אך במידה ומופיע הטקסט הבא: קראתי את הספר 'מלחמה ושלום', נרצה להפריד בין "י" "ושלום".
- ישנו קושי להבדיל בין שלום' מהמשפט הנ"ל לבין במילה סנדויץ'. בשני המיקרים התו "י" לא מופיע בין 2 אותיות אך פעם אחת צריך להוסיף רווח ופעם אחת לא. לבסוף החלטנו כי רק כאשר התו "י" מופיע בין 2 אותיות לא להוסיף רווחים, כיוון שבעברית מילה לא יכולה להתחיל בתו זה, ובמיקרים נדירים מאוד מילה מסתיימת בתו, וגם לאחר הוספת הרווח לא נוצר בלבול עם מילה אחרת.
- כאשר התו '.' מופיע בין שתי ספרות. במקרה של נקודה עשרונית $\pi = 3.141592653589$, אין צורך להוסיף רווחים, אך במקרה של נקודה בסוף משפט נרצה להוסיף רווחים. כמו כן במקרה והתו מופיע בהקשר של תאריך 14.3.1989
 - כאשר התו ',' מופיע בין שתי ספרות. במקרה של ייצוג מספרי של מספר גדול 1,000,000, אין צורך להוסיף רווחים, אך במקרה של פיסוק במשפט נרצה להוסיף רווח.
 - כאשר התו '/' מופיע בין שתי ספרות. במקרה תאריך 14/3/89, אין צורך להוסיף רווחים, במקרה ומופיע בין שתי אותיות נרצה להוסיף רווח. לדוגמא כאשר מופיע הטקסט: זכר/נקבה נרצה להפריד בין "זכר" "/" "נקבה".

המילים נפוצות ביותר בקורפוס הבדיקה

52587	,	1	○
43495	.	2	○
18151	של	3	○
16331	"	4	○
13148	-	5	○
12395	*	6	○
11159	את	7	○
10513	(8	○
10493)	9	○
9709	על	10	○
5089	:	11	○
5036	הוא	12	○
3553	עם	13	○
3459	ב	14	○
2837	גם	15	○
2627	לא	16	○
2593	בשנת	17	○
2446	או	18	○
2422	'	19	○
2283	היא	20	○
2229	היה	21	○
2100	ידי	22	○
2068	לאחר	23	○
2000	בין	24	○
1955	זה	25	○
1650	כי	26	○
1502	אך	27	○
1450	יותר	28	○
1429	כל	29	○
1386	אשר	30	○
1344	עד	31	○
1305	ה	32	○
1239	זו	33	○
1218	הם	34	○
1199	ישראל	35	○
1139	הייתה	36	○
1136	בית	37	○
1100	כמו	38	○
1063	חיצוניים	39	○
1047	קישורים	40	○
1041	ביותר	41	○
997	;	42	○
992	אחד	43	○
982	שם	44	○
967	הראשון	45	○
964	בו	46	○
940	מספר	47	○
924	אותו	48	○
874	היו	49	○
867	יש	50	○

אחוז המילים ששכיחותן מעל 50 הוא 1.8%
אחוז המילים ששכיחותן מתחת ל-10 הוא 89.7%
אחוז המילים ששכיחותן בדיוק 1 הוא 54%

ישנם 197981 מופעים של 10 המילים הראשונות, 21.4% מכלל הטקסט
ישנם 269433 מופעים של 50 המילים הראשונות, 29.1% מכלל הטקסט
ישנם 301497 מופעים של 100 המילים הראשונות, 32.6% מכלל הטקסט

הנתונים מתיישבים עם חוק Zipf כפי שניתן לראות בגרף:

