

מדד השקט הביטחוני העולמי

פרויקט סיום במסגרת הקורס "ביג דאטה ובינה עסקית למדע המדינה" - 56146
קבוצת המט"ריסטים

מגישים: קסם לופו, טל דקלו, שיר שגב, אברהם מרון, נווה מבורך



המוטיבציה והרקע לפרויקט

מדד השקט הביטחוני נבנה במטרה לשקף את העיסוק העולמי
(התקשורתי) בנושא: ביטחון.

**כך נוכל להבין עד כמה העולם היה בטוח בתאריך
מסוים.**

המדד מאפשר לנו לזהות את רגעי השקט או את רגעי הרעש
העולמיים על פני ציר זמן ולזהות אירועים חריגים בעזרת ניתוח
כתבות.



המשפך ההפוך של הפרויקט

geg_gcnlapi

source_urls

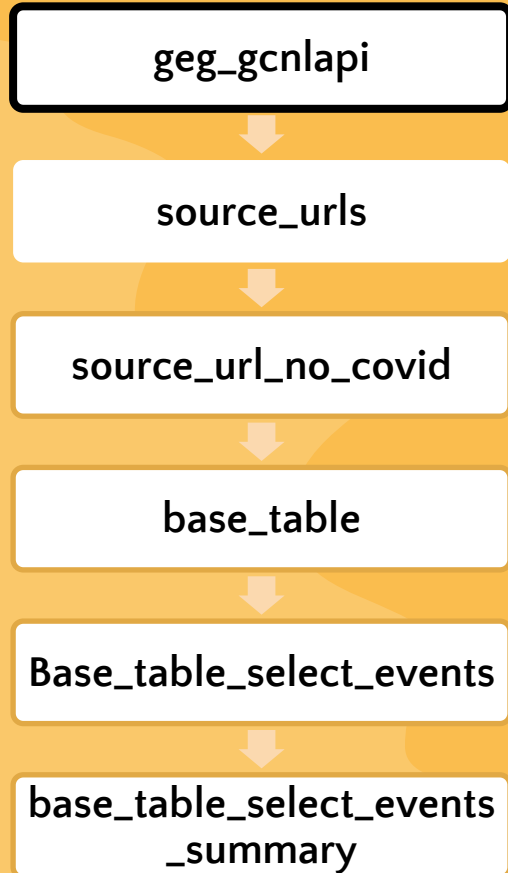
source_url no covid

base_table

Base table select
events

base table se
lect events s
ummary

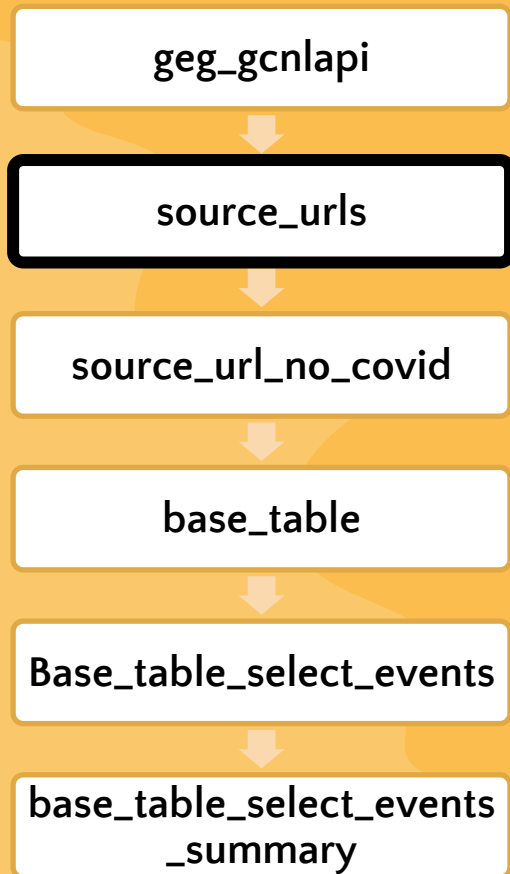
geg_gcnlapi



1. 256M שורות.
2. NLP - הטבלה הראשונה שמבוססת על עיבוד שפה טבעית
3. כל ישות שמוזכרת מקבלת סיווג- ראש מדינה, ארגון, מיקום, אירוע, תאריך וכו'.
4. 11 שפות
5. ניתן לזהות אירועים בזמן אמת: מגפות, מלחמות...
6. לכל ישות ניתן ציון בולטות בטווח 0-1



source_urls



- .1 52M שורות.
- .2 סינון לפי שפה: אנגלית בלבד.
- .3 סינון תאריכים: ינואר 21 – יוני 22.

צמצום של 80% מהטבלה הקודמת.

source_url_no_covid

geg_gcnlapi



source_urls



source_url_no_covid



base_table



Base_table_select_events



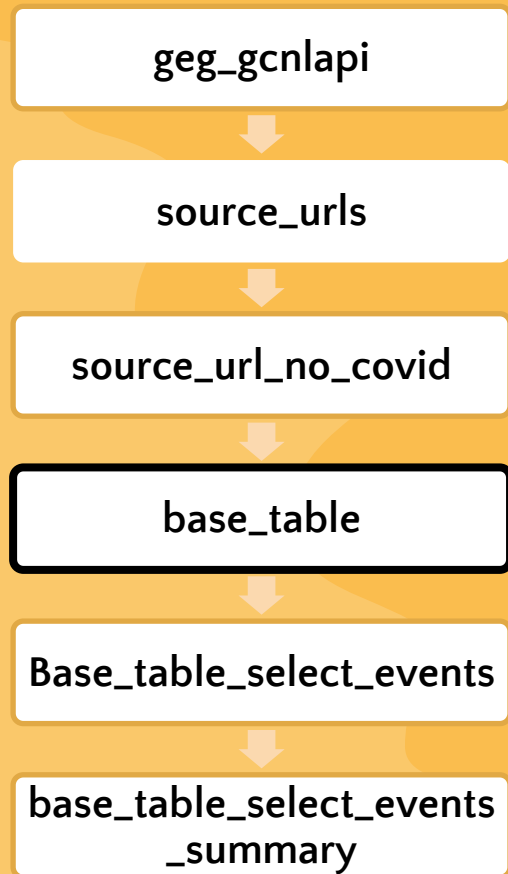
base_table_select_events_summary

.1 44M שורות.

.2 הסרת כתבות בנושא הקורונה- ממאגר הכתבות COVID 19

צמצום של 16% מהטבלה הקודמת.

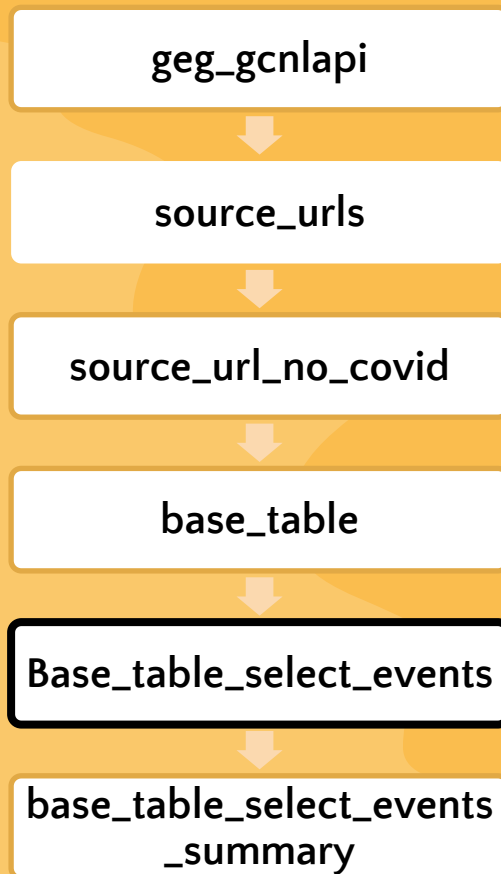
Base_table



White list: רשימת מילים שנרצה שיופיעו בכתבה - בנושא ביטחון
Black list: רשימת מילים שאסור שיופיעו בכתבה

1. 53M שורות.
2. ישות מסוג אירוע.
3. אין מילה שמופיעה בblack list.
4. לפחות מילה 1 מופיעה בwhite list.
5. רמת בולטות 0.001 של האירוע.
6. כל שורה – מכילה את כל האירועים של הכתבה ברמת בולטות מעל 0.001 (לכן גדלנו במספר)

Base_table_select_events

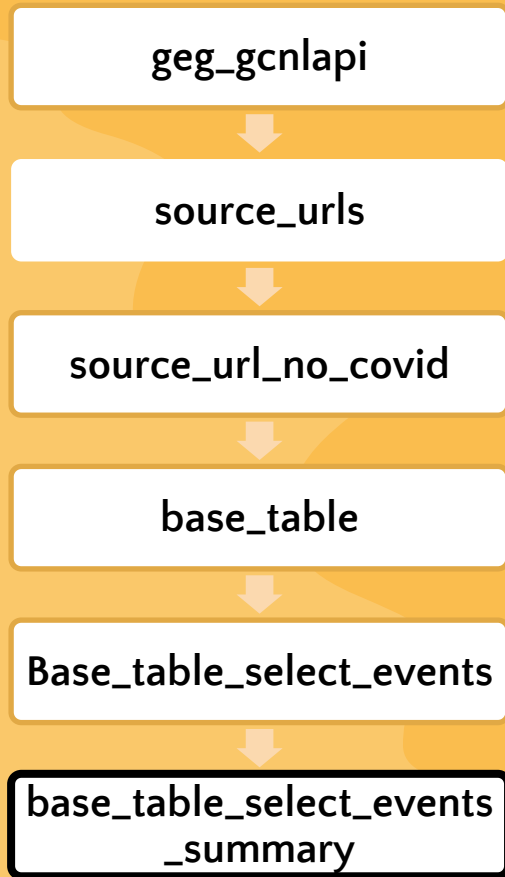


.1 9M שורות.

.2 כל שורה מכילה URL, אירוע מwhite list, רמת בולטות.

.3 לוקחים רק את האירועים הרלוונטיים.

base_table_select_events_summary



- .1 534 שורות.
- .2 מקבצים את הכתבות לפי תאריכים.
- .3 כל שורה = יום.

base table select events summary

Row	yyyyymmdd	count_distinct_urls	sum_numMentions	avg_avgSalienc
1	2021-01-01	8384	14714	0.011296849944236036
2	2021-01-02	7546	13667	0.009286360058309037
3	2021-01-03	8361	14850	0.010949649659182038
4	2021-01-04	11182	18935	0.009956830177245874
5	2021-01-05	11768	20055	0.01020762288013797
6	2021-01-06	11480	19339	0.008743130456963778

normalized

count_distinct_urls	count_distinct_urls_normalized	sum_numMentions	sum_numMentions_normalized
8384	1	14714	1
7546	1	13667	1
8361	1	14850	1
11182	4	18935	4
11768	6	20055	4
11480	5	19339	4
11477	6	25617	6

1. תאריך

2. מספר הכתבות באותו היום בנושא ביטחון

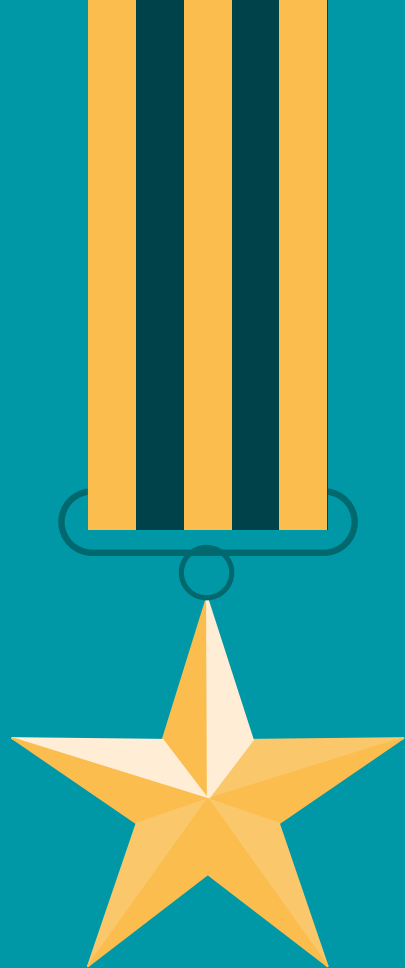
3. מספר האזכורים של white words בכתבות

4. רמת הבולטות הממוצעת- ממוצע השקט

5. ככל שמדד הרעש גבוה יותר- יש פחות בטחון.

6. מקרה חריג: קפיצה בכמות הכתבות, קפיצה

במספר האזכורים, וקפיצה במדד.



ממצאים

Salience Score

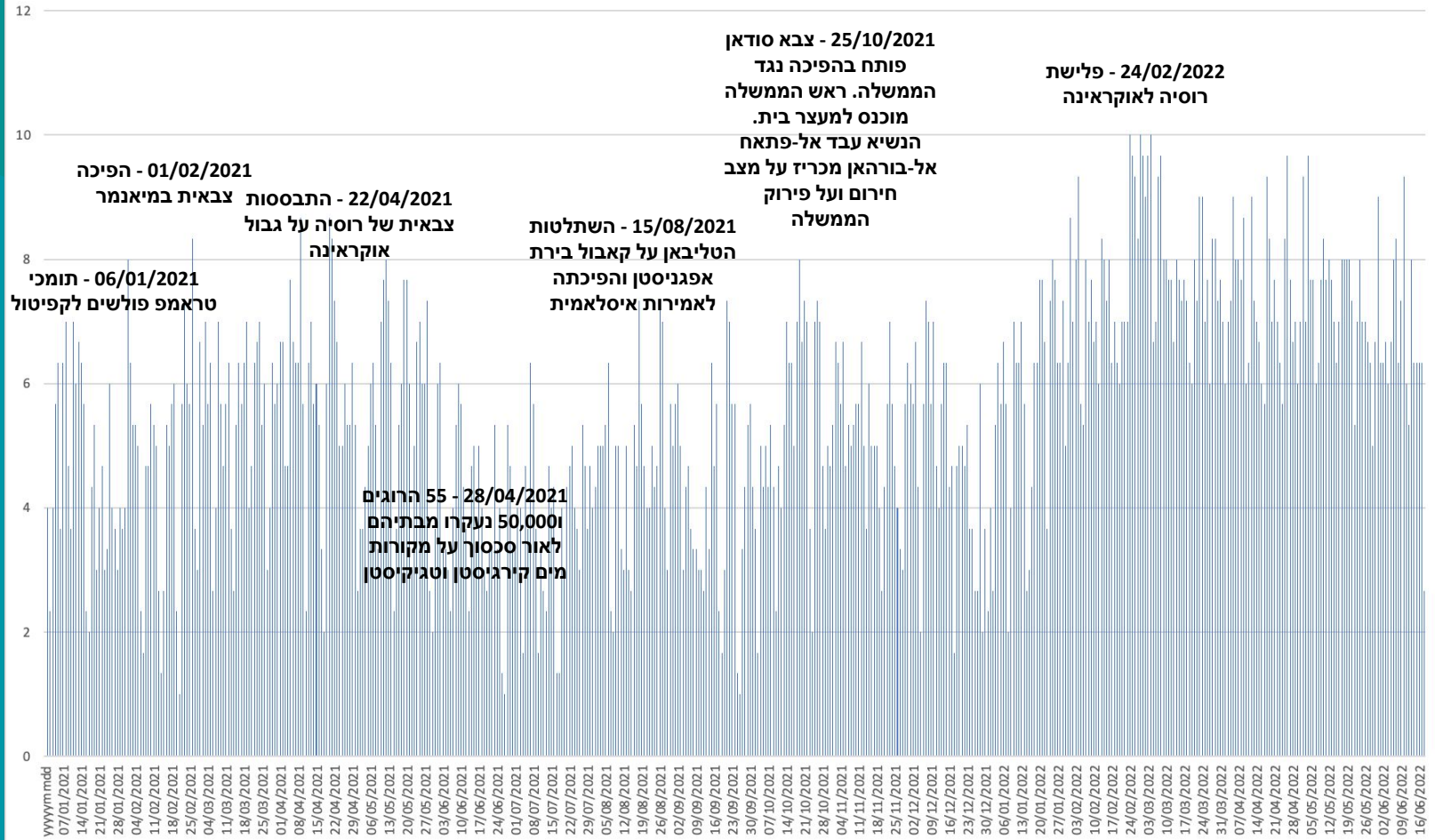


Chart Title



■ Friday ■ Monday ■ Saturday ■ Sunday ■ Thursday ■ Tuesday ■ Wednesday

צעדים להמשך



01

בחינת רשימת
המילים מחדש
white words.
דרישה של יותר
ממילה אחת בכתבה

02

סיווג הדאטה לפי
מדינות והצגת
ממצאים על גבי מפה

03

סיווג המילים עצמן
לנושאים – פיגועים,
תאונות, אלימות
ומלחמות..