

# מה זה בעצם JobMiner?

שימוש ב-JobMiner מאפשר חיפוש קל בפוסטים מתוך קבוצות פייסבוק המפרסמות הצעות עבודה. בזמן שפייסבוק נעשה אחת הפלטפורמות העיקריות לחיפוש עבודה, החיפוש של עבודות בפייסבוק הינו מסורבל בגלל הרבה מגבלות:

- רוב הקבוצות ממוקדות לנישה ספציפית למשל "משרות להייטק בסן פרנסיסקו" ולא ניתן לדעת מה הם הקבוצות הכי רלוונטיות לך לחיפוש עבודה.
- גם בתוך קבוצות רלוונטיות יש הרבה פוסטים שאינם קשורים, וקשה לחפש בהם תוכן ספציפי
- יש הרבה קבוצות ובשביל לחפש בכולן צריך לעבור אחת אחת.

האתר JobMiner נותן לחפש פוסטים רלוונטים, על פי 4 קטריונים חשובים - עיר, מדינה, חברה, סוג משרה, ובנוסף גם חיפוש על פי מילות מפתח.

**\*הערה: האתר מותאם כרגע לדפדפן Chrome**

## תרשים הזרימה של האתר

- האתר מרענן את הפוסטים שהוא שואב בצורה אוטומטית מדי יום.
- כאשר נכנס משתמש הוא מזין שאילתת חיפוש רלוונטית למשרה אותה הוא מחפש.
- מתוך ה-DB מוחזרים הפוסטים הרלוונטיים.
- כעת המשתמש יכול לבחור:
- על ידי לחיצה על הפוסט הרלוונטי הוא יכול להיכנס לפוסט עצמו בפייסבוק ולחפש מידע נוסף
  - על ידי לחיצה על הקבוצה ניתן לעבור לקבוצת הפייסבוק הרלוונטית בה הפוסט הופיע.
  - על ידי לחיצה על המייל נפתח חלון שליחת מייל והמשתמש יכול לשלוח מייל למפרסם ההודעה (אם מצויין בפוסט)

# תיאור מבנה בסיס הנתונים

## Groups(group\_id, group\_fb\_id, group\_name, update\_date) – טבלת

הקבוצות

מס' מזהה של הקבוצה, משמש כמפתח ראשי	group_id
מס' מזהה של הקבוצה באתר Facebook	group_fb_id
שם הקבוצה	group_name
תאריך שבו פורסם הפוסט האחרון בקבוצה	update_date

## JobPost(post\_id, post\_story\_id, publish\_date, employment\_form, working\_manner, email, full\_post\_body) – טבלת הפוסטים

מס' מזהה של הפוסט בטבלה, משמש כמפתח ראשי	post_id
מס' מזהה של הפוסט באתר Facebook	post_story_id
תאריך פרסום הפוסט בקבוצה	publish_date
סוג משרה	employment_form
אופן העבודה	working_manner
דואר אלקטרוני	email
תוכן הפוסט, מאונדקס לצורך חיפוש יעיל בטבלה	full_post_body

## Companies(company\_id, company\_name) – טבלת החברות

מס' מזהה של החברה, משמש כמפתח ראשי	company_id
שם החברה, משמש כאינדקס לצורך חיפוש יעיל בטבלה	company_name

## CitiesStates(city\_id, state\_id) – טבלה מקשרת בין ערים למדינות

מס' מזהה של העיר	city_id
מס' מזהה של המדינה	state_id

## Cities(city\_id, city\_name) – טבלת הערים

מס' מזהה של העיר, משמש כמפתח ראשי	city_id
שם העיר, משמש כאינדקס לצורך חיפוש יעיל בטבלה	city_name

## States(state\_id, state\_name) – טבלת המדינות

מס' מזהה של המדינות, משמש כמפתח ראשי	state_id
שם המדינות, משמש כאינדקס לצורך חיפוש יעיל בטבלה	state_name

## JobPostGroup(post\_id, group\_id) – טבלה מקשרת בין פוסטים לקבוצות

מס' מזהה של הפוסט	post_id
מס' מזהה של הקבוצה	group_id

## JobPostCity(post\_id, city\_id) – טבלה מקשרת בין פוסטים לערים

מס' מזהה של הפוסט	post_id
מס' מזהה של העיר	city_id

## JobPostCompany(post\_id, company\_id) – טבלה מקשרת בין פוסטים לחברות

מס' מזהה של הפוסט	post_id
מס' מזהה של החברה	company_id

## JobPostState(post\_id, state\_id) – טבלה מקשרת בין פוסטים למדינות

מס' מזהה של הפוסט	post_id
מס' מזהה של המדינה	state_id

# תכנון בסיס הנתונים:

---

בפרויקט שלנו, היחס בין פוסט לבין עיר/מדינה/חברה הוא many-to-one. כלומר, לכל פוסט יכולה להתאים רק עיר אחת, רק מדינה אחת, ורק חברה אחת (כך הוחלט שה-Parser יעבוד). מהסיבה הזו חילקנו את המידע למספר טבלאות (במקום לאכלס את כל המידע בטבלה אחת).

## אפשרויות לשיפור:

אם היתה לנו האפשרות להרחיב ולשפר את הפרויקט, היינו מעוניינים לייצר יחס של many-to-many בין הפוסטים לבין ערים/מדינות/חברות. כמו כן היינו רוצים להשתמש בבינה מלאכותית וב-crowd-sourcing ע"מ לשפר את ביצועי ה-Parser. ויתרנו על כך לעת עתה, בין השאר בשל הדרישה ליצור בסיס נתונים גדול – השיפורים הנ"ל גורמים לריצת ה-Parser להיות איטית בפקטור משמעותי, לכן תהליכי אכלוס בסיס נתונים בגודל המבוקש ידרוש הרבה יותר זמן.

# האופטימיזציות שבוצעו

---

הוגדרו אינדקסים מסוג B-Tree על השדות הטקסטואליים בטבלאות הסטטיות ע"מ לאפשר חיפוש יעיל לפי השמות של ערים, מדינות וחברות – שדות city\_name, state\_name ו-company\_name בהתאמה.

בנוסף הוגדר אינדקס מסוג FULLTEXT על השדה full\_post\_body בטבלה JobPost ע"מ לאפשר חיפוש יעיל בעזרת AGAINST() MATCH() (נציין שהגרסא שמותקנת בשרת בפקולטה אינה מאפשרת הגדרת אינדקס כ-FULLTEXT -ולכן שינינו את ה-ENGINE ל-MyISAM)

# תחזוקת בסיס הנתונים

הקובץ updateTimer.py מריץ עדכון כל 24 שעות. כל עדכון עובר על פוסטים שנוצרו ב 48 השעות האחרונות (עבור בניית בסיס הנתונים ההתחלתי שאבנו פוסטים שנוצרו ב 1500 הימים האחרונים).

החלק המכיל את קוד עדכון ה DB מחולק לשלושה חלקים:

- getPostsFromGroups** - האחראי על איטרציה על כל קבוצות הפייסבוק ששמורות בבסיס הנתונים, ובידוד פוסט יחיד.
- Updater** - האחראי על עדכון פוסט יחיד בטבלאות בבסיס הנתונים.
- Parser** - האחראי על סריקת ההודעה עצמה שבפוסט, וזיהוי אלמנטים מעניינים.

## getPostsFromGroups

מודול זה מעדכן את בסיס הנתונים. המודול משתמש ב-long-lived-access-token, אותו השגנו באמצעות Facebook Graph API. `startUpdate()` (הפונקציה הראשית) עוברת על הקבוצות שמופיעות בטבלה Groups ושולחת את המזהה של כל קבוצה לפונקציית `getGroupFeed()`. `getGroupFeed(group_id)` - יוצרת url לפי ה Facebook Graph API אשר שולף פוסטים של הקבוצה `group_id` שנוצרו ביומיום האחרונים ומכיל את השדות `message`, `created_time`, `id` של כל פוסט. מקבלים json עם חלק מהתוצאות, כאשר בסוף ה json מופיע `next_url` שמתאים לתוצאות הבאות. ב-`next_url` ע"מ לשלוף את התוצאות הבאות וממשיכים כך עד מתקבלות כל התוצאות. כל 'chunk' של json נשלח לפונקציה `getPosts(json_chunk, group_id)`. `getPosts(json_chunk, group_id)` - הופכת את `json_chunk` למילון של json, ושולחת כל json לפונקציה `UpdatePost(post, group_id)` של המודול **Updater**. לאחר מכן מגדיר המודול את זמן עדכון כל קבוצה להיות הזמן הכי מאוחר שפוסט בה נכתב (בעזרת שאילתת UPDATE).

## Updater

מודול זה מעדכן פרטי פוסט יחיד בבסיס הנתונים. לאחר שנוצר עצם מסוג זה, ניתן לקרוא למתודת **UpdatePost** שלו המקבלת את הפוסט ואת ה- `group_fb_id` של הקבוצה ממנה הוא נשאב (ובה פורסם). בהינתן המילון של הפוסט אותו היא מקבלת כפרמטר, היא שולפת את ההודעה (אם אין - פוסט זה אינו חוקי). כעת היא קוראת ל- **Parser** כדי לקבל מבנה מסוג **JobPostDetails**. לאחר מכן היא מעדכנת בטבלאות את הפרטים של הפוסט אם הוא לא מעודכן שם.

## Parser

תפקידו של ה-**Parser** הוא ניתוח טקסטואלי של הודעות מתוך פוסטים. ה-**Parser** מקבל תוכן של הודעה מתוך פוסט (ההודעה נמצאת בחלק ה-"message" של ה-**Json**, דוגמא בנספחים) ומעבד את המידע הטקסטואלי הנמצא בה. הפונקציה המרכזית של ה-**Parser** היא הפונקציה **parse**, שמקבלת תוכן של הודעה מתוך פוסט (ההודעה נמצאת בחלק ה-"message" של ה-**Json**) ומחזירה אובייקט מסוג **JobPostDetails**, אשר מכיל את השדות הבאים:

- ❖ **city\_id** - מס' מזהה של העיר מהפוסט.
- ❖ **state\_id** - מס' מזהה של המדינה מהפוסט.
- ❖ **company\_id** - מס' מזהה של החברה מהפוסט.
- ❖ **emp\_form\_id** - סוג משרה שפורסמה. 0 - משרה מלאה, 1 - משרה חלקית, 2 - משרה זמנית.
- ❖ **work\_manner** - אופן העבודה. 0 - עבודה במשרד, 1 - עבודה מהבית.
- ❖ **email** - כתובת הדואר האלקטרוני בפוסט.

עוברת על המילים בהודעה ומנסה להתאים כל מילה (או צירוף של מילים) לביטוי שמופיע באחת הטבלאות הסטטיות הבאות: Companies, Cities, States. בנוסף, הפונקציה מנסה למצוא ביטויים מתאימים ע"מ לעדכן את השדות **emp\_form\_id**, **work\_manner**, **email**. כאשר נמצאת התאמה, האובייקט **JobPostDetails** מעודכן בהתאם. יתכן שחלק מהשדות יקבלו ערך **None** אם ה-**Parser** לא מצא ערכים עבורם בהודעה.

# שרת ה-WEB

---

## :Server

השרת בנוי בעזרת FLASK והוא מטפל ב-2 סוגים של בקשות:  
(1) 3 בקשות "סטטיות" - עיר, מדינה וחברה:  
בהינתן מחרוזת השרת מחזיר רשימה של פריטים מהטבלה הרלוונטית שמתחילים באותה מחרוזת – מאפשר השלמה אוטומטית בצד הלקוח  
(2) בקשת חיפוש – בקשת חיפוש יכולה לכלול כל אחד מהפריטים הבאים ושילובים ביניהם: מילות חיפוש, סוג משרה – מלאה, חלקית או זמנית. מדינות וערים. נכון לעכשיו החיפוש מתמקד רק בתוך ארה"ב. החיפוש בצד השרת נעשה על ידי ביצוע Join בצורה חכמה – קודם כל מאחדים טבלאות שמסננים על פיהם, על ידי Join רגיל – כדי לסנן רשומות, ורק לאחר מכן מבצעים LEFT JOIN עם שאר הטבלאות.  
כמו כן השרת מחזיר תמיד את המשרות הרלוונטיות כשהן ממוינות מהחדש לישן.

## :Client

הלקוח בנוי ב-Javascript בעזרת ספריית React של Facebook המאפשר רענון חכם של הדפים בהתאם למידע שהתקבל וחלוקה היררכית של האלמנטים בדף.

# חבילות python בהן השתמשנו

---

pymysql – added this to our directory on nova

datetime

urllib

urllib2

json

sched

time

threading

unicodedata

re

sys

os

itertools

flask