

# מערכת לאגירת חדשות

במטלה זו ישנו דגש על OOP הפרדה למחלקות וירשיות.  
בחרתי להשתמש בתבנית עיצוב MVC (Model View Controller).  
הסבר על MVC:

תבנית זאת מתחלקת לשלושה חלקים: Model, View, Controller ובכך מפחיתה את התלות בין ממשק המשתמש לשאר חלקי התוכנה.

1. Model-משתמש כbackends שמכיל את כל הלוגיקה של הנתונים.
2. View- משמש כחזית או ממשק המשתמש הגרפי.
3. Controller- משמש כ"מוח" של המערכת הוא שולט על אופן הצגת הנתונים.

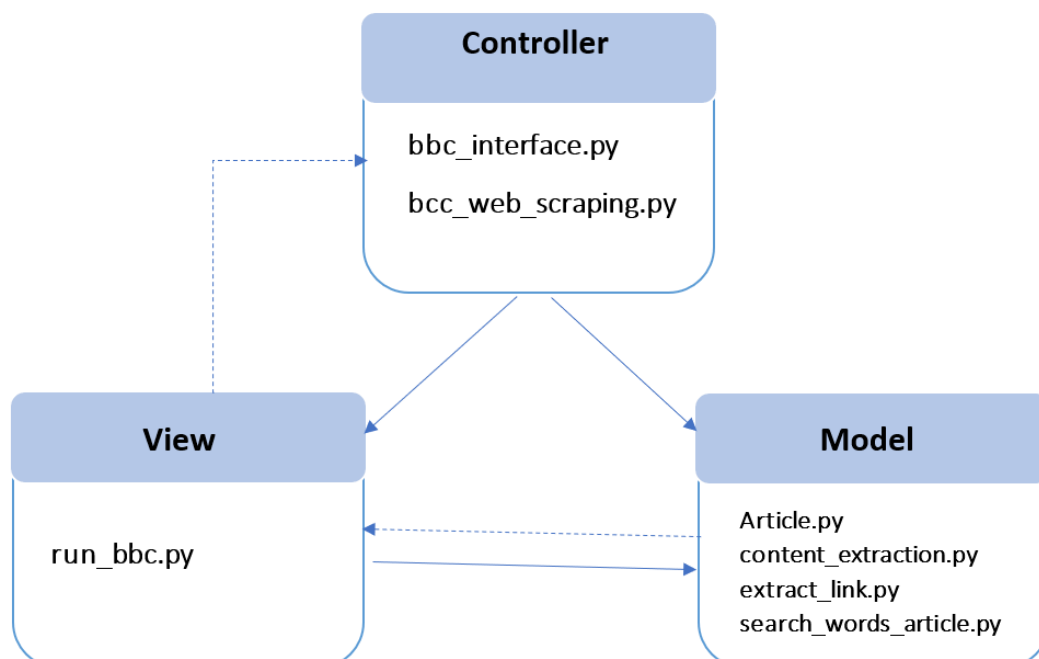
## מערכת כתבות BBC

תחילת התמקדתי בפתרון בבניית מערכת המורידה חומרי גלם מאתר:

<https://www.bbc.com/>

אתר זה מכיל כתבות בנושאים מרובים: חדשות, ספורט, תזונה ועוד..

התבנית MVC מתבטאת במערכת של האתר BBC באופן הבא:



## הסבר על המחלקות:

המחלקה `Scraping_BBC` הנמצאת בקובץ: `bbc_web_scraping.py`:

היא המחלקה הראשית במימוש המערכת על אתר [Bbc](#).

מחלקה זו מממשת את הממשק `bbc_interface.py` ויורשת `webdriver.Chrome`.

הפונקציות במחלקה עוסקות במימוש תתי משימות המפורטות בדרישות המטלה:

הורדת המאמרים מהאתר, אגירת חומרי גלם וחיפוש מילה או קומבינציה של מילים במאמרים.

### להלן הפונקציות המחלקה:

א. `land_page` – טוענת בdriver את האתר BBC.

ב. `Install_articles` –

פלט- קבצי JSON המכילים את כל חומרי הגלם מהמאמרים.

ראשית, הפונקציה שומרת ברשימה את כל הקישורים לכל הכתבות הנמצאים באתר.

לאחר מכן היא עוברת על כל קישור נכנסת לכתבה ולוקחת את חומרי הגלם.

(חומרי הגלם שבחרתי הם: כותרת, קישור, ותוכן הכתבה)

ג. `-Search_from_articles`

פלט- מדפיסה טבלה המכילה את שם הכתבה והלינק שבו נמצאת המילה או קומבינציה של מילים.

נעבור על כל קבצי JSON (כל קובץ = כתבה) ונחפש אם קיים בו המילה/צירוף אם כן נוסיף לרשימה.

לבסוף נדפיס את הרשימה.

## הקובץ constants.py:

קושי- יתכן אי סינכרון בין הניתוב על המחשב שלכם לשלי.  
ניסיתי לפתור אך ללא הצלחה.  
הוא קובץ המכיל את כל הניתובים לתיקיות, דרייברים ואתרים.

## המחלקה ExtractLink הנמצאת בקובץ: bbc\_web\_scraping.py:

מחלקה זו אחראית על חילוץ כל קישורי הכתבות הנמצאות בדף הראשי של האתר והשימוש בה נעשה רק במחלקה הראשית Scraping\_BBC.

קושי: היה עלי לחפש תג html אשר זהה בין כל קישורי הכתבות.  
תחילה חשבתי שמצאתי ואז כשספרתי כמה כתבות יש בפועל באתר וכמה קישורים מצאתי ראיתי שכתבות REEL עם תג שונה ולכן בקוד ניתן לראות שהתייחסתי לשניהם.

הבנאי שלה מקבל את הDriver של אתר הבית אשר בעזרתו נחלץ את התג לקישורי הכתבות.

פונקציה אחת שמומשה: pull\_all\_links\_articles – אשר עוברת על "הילדים" של התג ששמרנו כאובייקט המחלקה ומחפשת את התכונה 'href' אשר ידוע בhtml שמסמל קישור.

## המחלקה ContentExtraction הנמצאת בקובץ:

content\_extraction.py:

מחלקה זו אחראית על הורדת חומרי הגלם והשימוש בה נעשה רק במחלקה הראשית Scraping\_BBC.

קושי: נתקלתי בכך שלא כל קבצי הJSON נשמרו טוב.  
לאחר הבנה מדוע דווקא כתבות ספציפיות לא נשמרו טוב מצאתי שיש תווים כמו: ? גורמים לבעיה.

להלן פונקציות המחלקה:

א. -pull\_all\_information\_articles

פונקציה זו מושכת את המידע מדף הHTML בעזרת תגיות

המשותפות לקטגוריות הכתבה.

לדוגמא, מצאתי שכל הכתבות מסוג News יש להם תג המשותף להן וכך שלפתי את תוכן המאמר.

ב. `-extract_article_content`

פונקציה זו נקראת ע"י הפונקציה מסעיף א.

אנחנו עוברים שורה אחר שורה ומחפשים את התג עם התכונה של טקסט משרשרים את כל תוכן הכתבה לאובייקט `string`.

ג. `-is_exit`

לאחר הפעלה חוזרת של המערכת איננו רוצים שישמרו כפילויות של קבצי JSON לכן לא נוודא שלא נוריד כתבה שכבר קיימת לנו במאגר.

ד. `-write_to_json`

פונקציה זו פותחת את קובץ ה JSON ומגדירה בו את שם הכתבה, הלינק ותוכן הכתבה.

ה. `-remove_unnecessary_letters_from_name`

פונקציה זו פותרת את הקושי שציינתי למעלה.

## המחלקה SearchWords הנמצאת בקובץ: `search_words_article.py`:

מחלקה זו אחראית על חיפוש מילה/צירוף בכל הכתבות שהורדנו והשימוש בה נעשה רק במחלקה הראשית `Scraping_BBC`.

### פונקציות המחלקה:

א. `-print_data`

פונקציה זו מדפיסה טבלה מסודרת עם שמות העמודות שהגדרתי מראש.

ב. `-is_exit_in_art`

פונקציה שבודקת האם מילה/צירוף של מילים נמצא בטקסט מסויים.

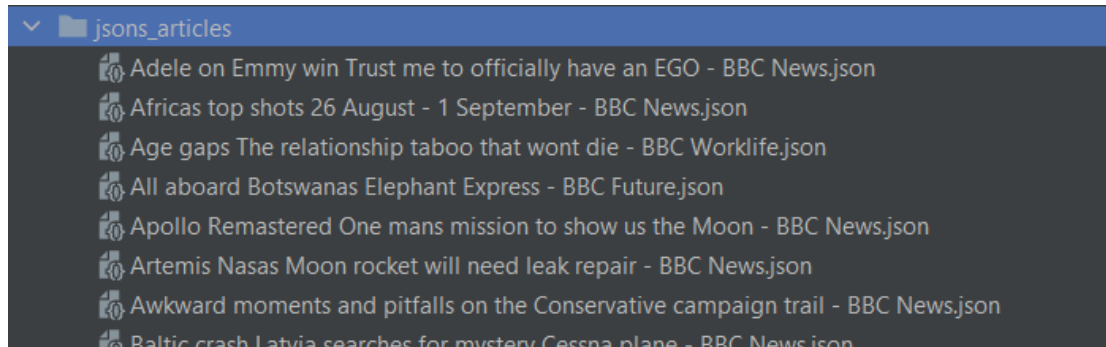
ג. `-search_in_all_articles`

פונקציה זו עוברת על כל הקבצים המסתיים בסיומת של `.json`. בתיקייה ספציפית ובודקת האם קיים טקסט המכיל מילה מסוימת/ צירוף.

אם כן, נשמור את כל הקבצים המכילים זאת באובייקט מסוג רשימה.

## צילומי מסך ודוגמאות הרצה:

לאחר הרצה נשמרו כל הכתבות בתיקיה: `jsons_articles` שם הקובץ הוא שם הכתבה.



התוכן שנשמר הוא: שם הכתבה, קישור לכתבה ותוכן הכתבה.

לאחר הרצה של פונקציית החיפוש על המילים הבאות:

1. USA

Article Name	Article Link
Age gaps The relationship taboo that wont die - BBC Worklife	<a href="https://www.bbc.com/worklife/article/20220317-age-gaps-the-relationship-taboo-that-wont-die">https://www.bbc.com/worklife/article/20220317-age-gaps-the-relationship-taboo-that-wont-die</a>
All aboard Botswanas Elephant Express - BBC Future	<a href="https://www.bbc.com/future/article/20220816-the-botswana-buses-tackling-human-elephant-conflict">https://www.bbc.com/future/article/20220816-the-botswana-buses-tackling-human-elephant-conflict</a>
Australian artist removes Ukraine and Russia mural after backlash - BBC News	<a href="https://www.bbc.com/news/world-australia-62751089">https://www.bbc.com/news/world-australia-62751089</a>
Chinese motorcyclists join Chongqing wildfire battle - BBC News	<a href="https://www.bbc.com/news/world-asia-china-62770119">https://www.bbc.com/news/world-asia-china-62770119</a>
In pictures The life of ex-Soviet leader Mikhail Gorbachev - BBC News	<a href="https://www.bbc.com/news/world-europe-62732829">https://www.bbc.com/news/world-europe-62732829</a>
New prime minister Seven big questions for the new leader - BBC News	<a href="https://www.bbc.com/news/uk-politics-62769937">https://www.bbc.com/news/uk-politics-62769937</a>
Tamil Nadu India womans food paintings you will want to eat - BBC News	<a href="https://www.bbc.com/news/world-asia-india-62763005">https://www.bbc.com/news/world-asia-india-62763005</a>
The spongy cities of the future - BBC Future	<a href="https://www.bbc.com/future/article/20220823-how-auckland-worlds-most-spongy-city-tackles-floods">https://www.bbc.com/future/article/20220823-how-auckland-worlds-most-spongy-city-tackles-floods</a>
UK looks to Sweden for a solution to nuclear waste - BBC News	<a href="https://www.bbc.com/news/business-62677934">https://www.bbc.com/news/business-62677934</a>
Why its time to talk about poo - BBC Future	<a href="https://www.bbc.com/future/article/20220830-the-new-science-of-recycling-human-poo">https://www.bbc.com/future/article/20220830-the-new-science-of-recycling-human-poo</a>

2. UK businesses

Article Name	Article Link
Gas prices soar 26% after Russia keeps key pipeline closed - BBC News	<a href="https://www.bbc.com/news/business-62789675">https://www.bbc.com/news/business-62789675</a>

3. a

Article Name	Article Link
Adele on Emmy win Trust me to officially have an EGO - BBC News	<a href="https://www.bbc.com/news/entertainment-arts-62793207">https://www.bbc.com/news/entertainment-arts-62793207</a>
Africas top shots 26 August - 1 September - BBC News	<a href="https://www.bbc.com/news/world-africa-62751074">https://www.bbc.com/news/world-africa-62751074</a>
Age gaps The relationship taboo that wont die - BBC Worklife	<a href="https://www.bbc.com/worklife/article/20220317-age-gaps-the-relationship-taboo-that-wont-die">https://www.bbc.com/worklife/article/20220317-age-gaps-the-relationship-taboo-that-wont-die</a>
All aboard Botswanas Elephant Express - BBC Future	<a href="https://www.bbc.com/future/article/20220816-the-botswana-buses-tackling-human-elephant-conflict">https://www.bbc.com/future/article/20220816-the-botswana-buses-tackling-human-elephant-conflict</a>
Apollo Remastered One mans mission to show us the Moon - BBC News	<a href="https://www.bbc.com/news/science-environment-62662685">https://www.bbc.com/news/science-environment-62662685</a>
Artemis Nasas Moon rocket will need leak repair - BBC News	<a href="https://www.bbc.com/news/science-environment-62758462">https://www.bbc.com/news/science-environment-62758462</a>
Australian artist removes Ukraine and Russia mural after backlash - BBC News	<a href="https://www.bbc.com/news/world-australia-62751089">https://www.bbc.com/news/world-australia-62751089</a>
Awkward moments and pitfalls on the Conservative campaign trail - BBC News	<a href="https://www.bbc.com/news/uk-politics-62753199">https://www.bbc.com/news/uk-politics-62753199</a>
Baltic crash Latvia searches for mystery Cessna plane - BBC News	<a href="https://www.bbc.com/news/world-europe-62789827">https://www.bbc.com/news/world-europe-62789827</a>
Boris Johnsons next move Making millions or a comeback - BBC News	<a href="https://www.bbc.com/news/uk-politics-62547853">https://www.bbc.com/news/uk-politics-62547853</a>
Canada stabbings Police hunt suspects after killing spree in Saskatchewan - BBC News	<a href="https://www.bbc.com/news/world-us-canada-62790909">https://www.bbc.com/news/world-us-canada-62790909</a>
Chile constitution Voters overwhelmingly reject radical change - BBC News	<a href="https://www.bbc.com/news/world-latin-america-62752025">https://www.bbc.com/news/world-latin-america-62752025</a>
China quake Deadly tremor rocks Sichuan city in lockdown - BBC News	<a href="https://www.bbc.com/news/world-asia-china-62764223">https://www.bbc.com/news/world-asia-china-62764223</a>
Chinese motorcyclists join Chongqing wildfire battle - BBC News	<a href="https://www.bbc.com/news/world-asia-china-62770119">https://www.bbc.com/news/world-asia-china-62770119</a>
Crypto.com pulls out of Uefa Champions League deal - BBC News	<a href="https://www.bbc.com/news/technology-62764654">https://www.bbc.com/news/technology-62764654</a>
Flexible firms take break from fixed bank holidays - BBC News	<a href="https://www.bbc.com/news/business-62491321">https://www.bbc.com/news/business-62491321</a>
Foo Fighters Girl, 12, drums at Taylor Hawkins memorial gig - BBC News	<a href="https://www.bbc.com/news/uk-england-wuffile-62793379">https://www.bbc.com/news/uk-england-wuffile-62793379</a>
Gas prices soar 26% after Russia keeps key pipeline closed - BBC News	<a href="https://www.bbc.com/news/business-62789675">https://www.bbc.com/news/business-62789675</a>
How to stay cool The Japanese way - BBC News	<a href="https://www.bbc.com/news/business-62691192">https://www.bbc.com/news/business-62691192</a>
In pictures The life of ex-Soviet leader Mikhail Gorbachev - BBC News	<a href="https://www.bbc.com/news/world-europe-62732829">https://www.bbc.com/news/world-europe-62732829</a>
Japans female bosses mapping a course for other women - BBC News	<a href="https://www.bbc.co.uk/news/business-62719779">https://www.bbc.co.uk/news/business-62719779</a>
Legionnaires suspected cause of Argentina pneumonia deaths - BBC News	<a href="https://www.bbc.com/news/world-latin-america-62785348">https://www.bbc.com/news/world-latin-america-62785348</a>
Twentynine 2022 Meet the esteemed film festival jury - BBC Culture	<a href="https://www.bbc.com/culture/articles/C0T90836-2022-meet-the-esteemed-film-festival-jury">https://www.bbc.com/culture/articles/C0T90836-2022-meet-the-esteemed-film-festival-jury</a>

## מערכת לוח טיסות נתב"ג

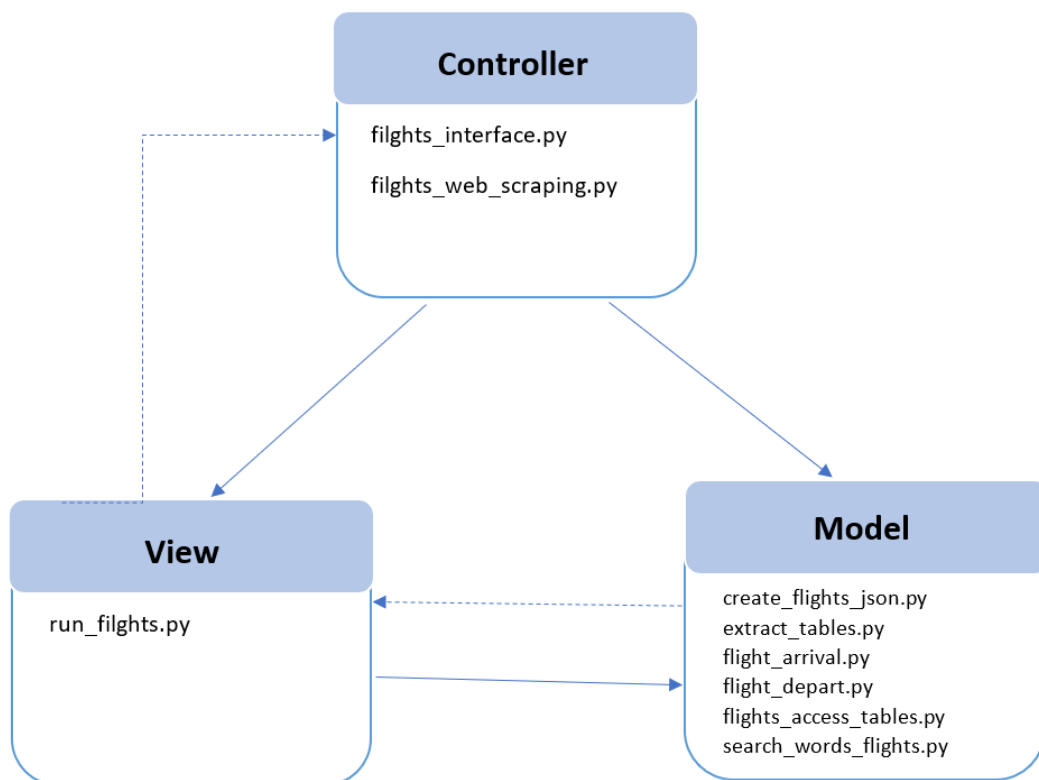
שנית עסקתי בבניית מערכת המורידה חומרי גלם מאתר:

<http://www.iaa.gov.il/he-IL/airports/BenGurion/Pages/OnlineFlights.aspx>

יש לציין שאתר זה נותן מידע על טיסות אשר יוצאות/נכנסות בשדה תעופה בן גוריון ישראל.

בחלק זה התבקשנו לאסוף את נתוני הטיסות בReal-Time, נעזרתי באתר הבא: [לינק](#)

התבנית MVC מתבטאת במערכת של האתר הנ"ל באופן הבא:



הסבר על המחלקות:

## המחלקה Scraping\_Flights הנמצאת בקובץ: flights\_web\_scraping.py

היא המחלקה הראשית במימוש המערכת על אתר הנ"ל.

מחלקה זו מממשת את הממשק flights\_interface.py ויורשת  
webdriver.Chrome.

הפונקציות במחלקה עוסקות במימוש תתי משימות המפורטות בדרישות  
המטלה:

הורדה של נתוני לוח טיסות בזמן אמת, אגירת חומרי גלם וחיפוש מילה או  
קומבינציה של מילים במידע של הטיסות שאגרנו.

### להלן הפונקציות המחלקה:

א. land\_page – טוענת בdriver את האתר טיסות בבן גוריון.

ב. -install\_flights

פלט- קבצי JSON המכילים את כל חומרי הגלם מלוח הטיסות.  
ניתן לראות כי ב-GUI יש שני קטגוריות לטיסות (נחיתות והמראות)  
ובנוסף הטבלה לא מוצגת במלואה וישנו כפתור אז פותח עוד ועוד  
פרטים על טיסות.

אני רוצה לאגור את כל נתוני הטיסות בלוח הטיסות לכן עלי לפתוח  
את לוח הטיסות בשני הקטגוריות ולהציגו במלואו.  
לכן, ראשית התחלתי עם לוח הנחיתות פתחתי עד סוף את  
הרשימה.

לאחר מכן, עצרתי את ה"עדכון אוטומטי" מאחר שבזמן אגירת פרטי  
הטיסה אינני רוצה לפספס שום פרטים שהתעדכנו.  
וכך עשיתי גם בלוח ההמראות.

לבסוף, אספתי את כל המידע על כל טיסה הנמצאת בלוח הטיסות.

ג. -search\_from\_flights

פלט- מדפיסה טבלה המכילה את שם חברת התעופה, מס טיסה,  
יעד, מס טרמינל, זמן הגעה, זמן עדכני וסטטוס שבו נמצאת המילה  
או קומבינציה של מילים.

נעבור על כל קבצי ה-JSON (כל קובץ = פרטי טיסה) ונחפש אם קיים  
בו המילה/צירוף אם כן נוסיף לרשימה.

לבסוף נדפיס את הרשימה.

## המחלקה AccessExtractTables הנמצאת בקובץ: flights\_access\_tables.py:

מחלקה זו אחראית על חילוץ כל פרטי הטיסה בלוח הטיסות של נתב"ג והשימוש בה נעשה רק במחלקה הראשית Scraping\_Flights. הבנאי שלה מקבל את הDriver של אתר הבית אשר בעזרתו נחלץ את התג לכפתורים שנרצה ללחוץ עליהם ("הצג תוצאות נוספות", "עצור עדכון אוטומטי")

### להלן פונקציות המחלקה:

- א. get tables - פונקציה זו היא הפונקציה הראשית אשר אחראית על חילוץ חומרי הגלם.
- ב. get arrival table - ראינו כי לוח הטיסות בעמוד הראשי קצר. על מנת לראות את כל הטיסות עלינו ללחוץ על כפתור שפותח עוד ועוד פרטים על טיסות. לכן פונקציה זו נועדה לפתוח טבלת הטיסות הנכנסות. בנוסף, עצרתי את הלוח מלהתעדכן כדי שבזמן אגירת המידע לוח הטיסות ישתנה.
- ג. get depart table - באופן זהה לפונקציה מסעיף ב רק שפה אנחנו צריכים ללחוץ על כפתור אשר מציג את לוח הטיסות היוצאות ולהמתאים שיפתח.

## המחלקה ExtractTables הנמצאת בקובץ: extract\_tables.py:

מחלקה זו אחראית על הורדת חומרי הגלם והשימוש בה נעשה רק



## במחלקה `AccessExtractTables`.

הבנאי שלה מקבל את ה `Driver` של אתר הבית אשר בעזרתו נחלץ את כל המידע על הטיסה בעזרת התגים המתאימים.

### להלן פונקציות המחלקה:

א. `-extract_flights`

הפונקציה הראשית של המחלקה אשר משתמשת בפונקציות עזר שאפרט מטה.

אחראית על הוצאת המידע מהטיסות ושולחת למחלקה שתפקידה לשמור את המידע שהוצאנו בקבצי `JSON`.

ב. `- pull_flight_arrival_information/ depart`

פונקציה זאת עוברת שורה אחר שורה בטבלת הטיסות הנכנסות והיוצאות ושולחת לפונקציית עזר שתשלוף את הטקסט מהשורה.

ג. `- get_arrival_row_data/ depart`

מכניסים לאובייקט טיסה נכנסת/יוצאת את המידע שנשלוף בעזרת פונקציות ה `GET`.

ד. פונקציות `Get` למיניהם-

פונקציות אלו עוברת על לוח הטיסות בדף `HTML` ולפי תגים שבהם נמצאים הטקסט מחזירות את המידע הרצוי.

## המחלקה `CreateJson` הנמצאת בקובץ: `create_flights_json.py`:

מחלקה זו אחראית על שמירת נתוני הטיסות בקבצי ג'ייסון. והשימוש בה נעשה רק במחלקה `ExtractTables`.

הבנאי שלה מקבל רשימה של טיסות וסוג טיסה (המראה או נחיתה).

להלן פונקציות המחלקה:

- א. – write arrival to json / write depart to json  
מגדירה את המבנה שבו ארצה שהקובץ JSON יבנה (key:value).
- ב. -create file  
בפונקציה זו נפתח קובץ JSON שאת שמו בחרתי להיות מספר הטיסה.

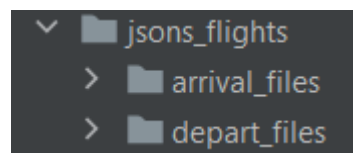
המחלקה SearchWordsFlights הנמצאת בקובץ:

:search\_words\_article.py

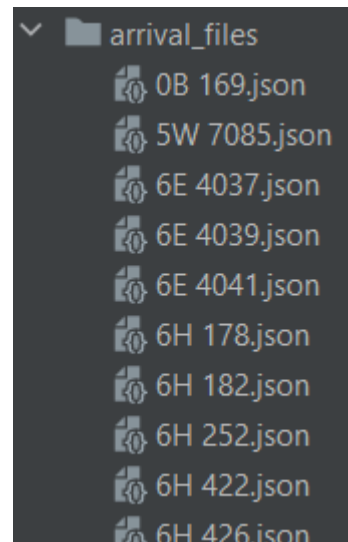
מחלקה זו אחראית על חיפוש מילה/צירוף בכל פרטי הטיסות שהורדנו והשימוש בה נעשה רק במחלקה הראשית **Scraping\_Flights**.  
מחלקה זאת דומה למחלקת חיפוש הכתבות BBC השוני היחידי הוא שפה נעבור על כל העמודות בקובץ הJSON של טיסה.

צילומי מסך ודוגמאות הרצה:

לאחר הרצה נשמר תוכן לוח הטיסות בתיקיה: `jsons_flights` שם הקובץ הוא מספר טיסה.  
התוכן שנשמר הוא: שם חברת תעופה, מספר טיסה, מאיזה/לאיזה מדינה, מספר טרמינל, זמן מתוכנן, זמן עדכני וסטטוס טיסה .



קבצי הjson מסודרים באופן הבא:



לאחר הרצה של פונקציית החיפוש על המילה: "Blue" קיבלנו-

1. נחיתות

Airline Name	Flight Num	From -> Israel	Terminal	schedule_time	updated_time	Status
blue air aviation	0B 170	בוקרסט	3	05/09 22:55	23:25	בזמן
blue bird airways	BZ 112	לרנקה	3	06/09 09:20	09:20	בזמן
blue bird airways	BZ 702	אתונה	3	06/09 09:20	09:20	בזמן
blue bird airways	BZ 752	הרקליון	3	06/09 08:40	08:40	בזמן
blue bird airways	BZ 754	הרקליון	3	06/09 14:00	14:00	בזמן

2. המראות

Airline Name	Flight Num	Israel -> To	Terminal	schedule_time	updated_time	Status
blue air aviation	0B 169	בוקרסט	3	05/09 21:55	22:45	סופי
blue bird airways	BZ 111	לרנקה	3	06/09 13:30	13:30	לא סופי
blue bird airways	BZ 701	אתונה	3	06/09 15:30	15:30	לא סופי
blue bird airways	BZ 753	הרקליון	3	06/09 13:00	13:00	לא סופי
blue bird airways	BZ 777	סנטורני	3	05/09 19:30	20:53	נחתה