

תרגיל בית 1 – Bash Scripting

מועד ההגשה:	יום ה', 14/11/22, בשעה 23:55	
האחראי על התרגיל:	פבל ליפשיץ	pavel@ee.technion.ac.il

בתרגיל זה ננסה לבצע אוטומציה לעבודתו של יועץ תקשורת לאיש ציבור, אשר מחפש כתבות בעיתון בהן איש הציבור מוזכר. לצורך כך, נכתוב סקריפט bash שיעשה שימוש בכלי command line כדי שנלמדו בתרגול ובסדנה במטרה למצוא את כל הכתבות היומיות באתר החדשות ynetnews.com בהן מוזכרים נתניהו ו/או לפיד ו/או גנץ ו/או בן-גביר, ולייצר סיכום בקובץ טקסט הכולל:

1. מספר הכתבות שנבדקו
2. קישור לכל כתבה + מספר הפעמים ש-Netanyahu מוזכר בכתבה, מספר הפעמים ש-Gantz מוזכר, מספר הפעמים ש-Lapid מוזכר בכתבה ומספר הפעמים ש-Ben-gvir מוזכר בכתבה באופן הבא:

60

<https://www.ynetnews.com/article/SklUnwXow>, Netanyahu, 18, Gantz, 21, Lapid, 0, Ben-gvir, 0

<https://www.ynetnews.com/article/rJBFpUicD>, -

<https://www.ynetnews.com/article/Hy3I0PGow>, Netanyahu, 6, Gantz, 2, Lapid, 0, Ben-gvir, 0

...

את התוצאה יש לשמור בקובץ results.csv.

שימו לב:

1. אין לייצר קבצי ביניים!
2. שמות המשפחה מתחילים באות גדולה.
3. במידה ואחד מאנשי הציבור לא מופיע והאחר כן, נכתוב 0 ליד שמו של זה שאינו מופיע.
4. במידה ואף אחד מאנשי הציבור לא מופיע, נכתוב את הסימן -.
5. שימו לב במיוחד לרווחים!

יש להגיש את הסקריפט בשם scrape_news.sh יחד עם קובץ התוצאה ב-git ולהגיש את הלינק במערכת ה-moodle. הקוד הקצר והיעיל יזכו את המגישים ב-0.5 נקודה בונים לציון הסופי.

בנוסף, ענו בכתב בקצרה באנגלית בקובץ answers.txt (קובץ טקסט פשוט):

1. כמה זמן להערכתכם היה לוקח לעשות זאת באופן ידני?
2. לאיזה מסקנות אתם מגיעים בעקבות התרגיל? האם יש לכם רעיונות באלו עוד משימות/יישומים ניתן ליישם רעיון מסוג זה?

3. במידה והייתי רוצה לחזור על הפעולה כל שעה, מה היה נדרש ממני? האם וכיצד הייתי יכול לבצע גם את זה באופן אוטומטי? כיצד נתמודד עם כתבות שעדיין מופיעות וכבר נסרקו על ידי הקוד שלנו?

הכוונה ורמזים

עבדו בשלבים. פצלו כל שלב ובדקו את תקינותו ורק אח"כ חברו בין כל הדברים. יתכן ותרצו לעשות את חלקי המשימה באופן "ידני" על מנת להבין טוב יותר מה מנסים לעשות.

שלב ראשון

הדבר הראשון שתבצעו יהיה "לגלוש" באמצעות פקודה wget לאתר

<https://www.ynetnews.com/category/3082>

```
wget https://www.ynetnews.com/category/3082
```

תוצר הפקודה יהיה קובץ 3082, זהו למעשה קובץ html שהדפדפן שלכם מציג את תוכנו בצורה גרפית. (נסו לגלוש לאתר באמצעות הדפדפן, לחצו על כפתור עכבר-ימני, ובחרו view-source).

שלב שני

כעת נעבור על הקובץ שהתקבל מהשלב הראשון (זהו קובץ העמוד החדשות הראשי ובו קישורים לכל כתבה/ידיעה חדשותית). אנו נחפש בו קישורים מהצורה:

<https://www.ynetnews.com/article/XXXXXXXXXX>

כאשר X יכול להיות ספרה או אות גדולה/קטנה. יתכן ותרצו לשמור רשימה זאת.

שלב שלישי

כעת ניתן לספור כמה קישורים ייחודיים יש לכם.

שלב רביעי

כעת נרצה "לגלוש" (בפקודת wget) לכל קישור כזה לידיעה חדשותית. ולבצע מעבר על התוכן מתוך מטרה לספור כמה פעמים (אם בכלל) מופיעים שמות המשפחה של אישי הציבור.

מהו סקריפט bash?

עד כה, ראיתם בתרגול ובסדנאות סדרה של פקודות command line אשר ניתן לשלבן ביחד (למשל על ידי pipe - | או הפניית קלט/פלט). אך מה קורה במידה ואנחנו רוצים להריץ סדרה של פקודות? לצורך כך ניתן לערוך קובץ טקסט. הקבוצה חייב להתחיל בשורה #!/bin/bash. לאחר מכן ניתן לתת לקובץ הרשאות הרצה:

```
$ chmod +x ./my_script.sh
```

ואם תוכן הקובץ נראה כך :

```
#!/bin/bash
echo hello everyone
# this is a comment
echo there are `ls -l | wc -l` entries in this directory
```

הרצתו על ידי :

```
$ ./my_script.sh
```

תגרום להדפסה של :

```
hello everyone
there are 5 entries in this directory
```

הוראות הגשה :

- עברו היטב על הוראות ההגשה של תרגילי הבית המופיעים באתר טרם ההגשה! ודאו כי התוכנית שלכם עומדת בדרישות הבאות :
 - קוד התוכנית קריא וברור
 - קוד התוכנית מתועד היטב לפי דרישות התיעוד אותן למדנו
 - יש להגיש לינק ל-repository המכיל את הקבצים scrape_news.sh results.csv answer.txt (שימו לב להגיש לינק ל-repository lower case). על ה-repository להיות בעל הרשאות private. בעת בדיקת התרגיל, אנו נבצע clone ל-repository שלכם, נריץ את הסקריפט scrape_news.sh ונבדוק את קובץ ה-results.csv שיתקבל בהרצה. שימו לב להגיש לפי הפורמט הבא :
- <https://github.com/your-username/repository-name>
- 0123456789 student_1_mail@campus.technion.ac.il first_name_1 last_name_1
- 0123456789 student_2_mail@campus.technion.ac.il first_name_2 last_name_2
- יש לקרוא את קובץ הנחיות שיעורי הבית באתר לפני ההגשה!
 - שאלות בנוגע לתרגיל יש ניתן לשאול בשעות הקבלה של סגל הקורס או לשלוח מייל למתרגל האחראי על התרגיל בלבד, ורק במידה והשאלה מכילה פתרון חלקי.

סיכום מפרט התרגיל :

סעיף	תיאור
נושא התרגיל	Bash scripting
תאריך ההגשה	יום ב', 14/11/2020 בשעה 23:55
המתרגל האחראי על התרגיל	פבל ליפשיץ pavel@ee.technion.ac.il

	קבצי הקוד הנתונים
	קבצי הקלט והפלט הנתונים
answers.txt results.csv scrape_news.sh	הקבצים שיש להגיש

בהצלחה!