

תוכנית 1 – אביב 2022/23

תרגיל מס' 8

אוסףים גנריים ו- collection framework

הנחיות כלליות:

קראו בעיון את קובץ נהלי הגשת התרגילים אשר נמצא באתר הקורס.

את התרגיל הבא צריך להגיש באופן הבא:

- הגשה במערכת Git תבוצע על פי הנחיות שראיתם בתרגול 1.
צרו את ה repository שלכם מתוך קישור הבא:

<https://classroom.github.com/a/FxCBG18J>

יש לוודא שבתיקית הגיט שלכם נמצאים הקבצים הבאים:

- קובץ פרטיים אישיים בשם details.txt המכיל את שם המשתמש שלכם ב Moodle ואת מס' תעודה זהה שלכם.
 - 2 התקיות zips ו-resources. יש להגיש את ה repository ו- את מבנה התקיות שבתוכו בדיק באוטה היררכיה שקיבלה אותם
- הגשה במערכת Moodle (<http://moodle.tau.ac.il/>): עלייכם להגיש את קובץ הטקסט assignment.txt ובו קישור לrepository git האישי שלכם.

הנחיות כלליות נוספת לתרגיל:

1. מומלץ ראשית לקרוא את כל הוראות עבר הchalk שאתם ניגשים לפטור על מנת לוודא שהמימוש שלכם יהיה מותאם לצורה הטובה ביותר לדרישות התרגיל.
2. בכל אחד מחלקי התרגיל ניתן להוסיף שירותים ומחלקות לפי הצורך, אך אין לשנות חתימות של שירותים קיימים והגדירות של מנשכים (אם אין להוסיף או להוריד throws).
3. אין הגבלות על import-importים או שימוש במבני נתונים מסוימים, אך שימו לב שאתם מוחקים import-ים מיוחדים או אלה שאינם מתקנופלים במבנה.
4. בכל חלק קיימ טסט קצר המבצע בדיקות שפויות (חפשו בגוגל sanity tests). כדאי ומומלץ להוסיף בדיקות משלכם שכן הטסטרים הם בסיסיים ביותר ולא בודקים את כל המקרים.

حلק א' (50 נק')

בתרגיל זה עליכם למשם מבנה נתונים של היסטוגרמה באמצעות אוסףים גנריים. נגיד היסטוגרמה בתור מבנה נתונים אשר סופר מופעים של עצמים מטיפוס `T` כלשהו (טיפוס גנרי). הקוד ימושך בחבילה `.il.ac.tau.cs.sw1.ex8.histogram`.

לדוגמא, עבור אוסף האיברים הבא: `1, 2, 1, 3, 2, 1`, היסטוגרמה תכיל את האיברים `1, 2, 3` ואת מספר המופעים שלהם.

יחד עם קבצי התרגיל מסופק לכם הממשק `Histogram` את השירותים הבאים (כולל החוצה של כל שירות):

```
public void addItem(T item);
public boolean removeItem(T item);
public void addAll(Collection<T> items);
public int getCountForItem(T item);
public void clear();
public Set<T> getItemsSet();
public int getCountsSum();
```

- א. השירות `addItem` מוסיף מופיע אחד של הפריט `item` להיסטוגרמה.
- ב. השירות `removeItem` מוריד מופיע אחד של הפריט `item`. השירות מחזיר `true` אם הפריט הוסר, ו-`false` אם הפריט לא היה קיים, וכן גם לא הוסר.
- ג. השירות `addAll` מוסיף אוסף של פריטים להיסטוגrama (הוסף יכול להכיל את אותו הפריט יותר מפעם אחת. בנווסף, לא ניתן להגיד דבר על הקיום של פריטים אלה בהיסטוגרמה לפני הקראיה ל `addAll`).
- ד. השירות `getCountForItem` יחזיר את מספר הפעמים שהאיבר `item` נספר. אם `item` הוא פריט שלא קיים בהיסטוגrama, יוחזר הערך `0`.
- ה. השירות `clear` ירוקן את ההיסטוריה מכל האיברים והספירות (כלומר, לאחר `clear`, השירות `getCountForItem` יחזיר ספירה `0` לכל איבר).
- ו. השירות `getItemsSet` יחזיר אוסף מטיפוס `Set` אשר מכיל את כל האיברים בהיסטוגרמה אשר מסpter המופיעים שלהם גדול מ-`0`, ללא הספירות שלהם.
- ז. השירות `getCountsSum` יחזיר את סכום המופיעים של איברי ההיסטוגרמה.

בנוסף לשירותים שפורטו כאן, הממשק `Histogram` מימוש את `Iterable`, כך שתצטרכו למשם את הפונקציה `iterator` כדי שיאפשר בהמשך.

סעיף 1 (25 נק')

משמעות המחלקה `HashMapHistogram` אשר מימושה את הממשק `IHistogram` עבר כל טיפוס `T` המழמם את הממשק `Comparable` (כלומר, `T` יכול לקבל ערך של כל מחלקה המழממת את הממשק `Comparable`). נזכיר כי הטיפוסים המובנים הבסיסיים כמו `Integer` ו-`String` מממשים ממשק זה). לדרישה זו יש סיבה אותה ניתן לראות בהמשך.

פרקיתית, זה אומר שנדריך את `HashMapHistogram` באופן הבא:

```
public class HashMapHistogram<T extends Comparable<T>> implements IHistogram<T>
```

משמעות הגדרה זו: הממשק `IHistogram` מחייב פרמטר גנרי. את הפרמטר הגנרי נגדיר כ-`T`, ונוסף עליו את האילוץ שהוא צריך למש את הממשק `Comparable`. ככלומר, `T` יהיה פרמטר גנרי מתאים אם הוא Comparable עם עצמים אחרים מטיפוס `T`. את התחריב הזה ניתן להסביר ה庫ורס.

השימוש יעשה באמצעות הכללה (aggregation) של `HashMap`, כלומר, כל מופע של `HashMapHistogram` יוכל שדה מטיפוס `HashMap`. שדה זה יהיה אחראי על שמירת הספירות עבור כל אובייקט מטיפוס `T`.

סעיף 2 (25 נק')

הממשק `IHistogram` ירוש מהממשק `Iterable`, מה שמחייב את `HashMapHistogram` למש את השירות `iterator()`. הגדרת הממשק היא:

```
public interface IHistogram<T> extends Iterable<Map.Entry<T, Integer>>
```

כלומר, האיטרטור שההיסטוריה אמורה להחזיר עבר על זוגות (מטיפוס [Map.Entry](#)). הזוגות מייצגים מפתחות וערכים בהיסטוריה. תזכורת - השירות [entrySet](#) של המחלקה `Map` מחזיר בנוסף של איברים מטיפוס `Map.Entry`.

האיטרטור יחזיר את הזוגות על פי הסדר הטבעי של המפתחות (האיברים מטיפוס `T`). לדוגמה – אם ההיסטוריה מכילה מחרוזות, הסידור יהיה לפי סדר לקסיקוגרפי עולה. אם ההיסטוריה מכילה מספרים שלמים, הסידור יהיה סידור המספרים בסדר עולה.

לצורך כך علينا למש:

- א. מחלקה חדשה המழממת את הממשק `Iterator`. שם המחלקה הוא `HashMapHistogramIterator`. והשלד שלו נתון לכם. אין צורך למש את השירות `remove`.
- ב. כדאי למש גם `Comparator` בשבייל סידור האיברים שהאיטרטור יחזיר. ניתן למש מחלקה זו כמחלקה פנימית במחלקה האיטרטור או כמחלקה בקובץ `Java` נפרד ממנה.

שימוש לב, ניתן, ואף כדאי, להעביר למחלקות `Comparator` וה-`Iterator` את המידע הרלוונטי מתוך המופע של `HashMapHistogram` וכן להשתמש בהכלה של אוספים לפי הצורך.

להבין איזה מידע כל אובייקט צריך לקבל הוא חלק מהאתגר, חישבו על כך בעת בנית המחלקות.

שלד כל המחלקות אותן נדרש למשתמש למסמך בחרילה `hashmapTester.java` המופיעה בקובץ התרגיל.

העזרו ב `HashMapHistogramTester` בשבייל לבדוק את עצמכם, והסיפו לו בדיקות משלכם.

חלק ב' (50 נק')

בחלק זה נתרגל עבודה עם אוסףים (Collections) ע"י מימוש מחלקה אשר מנהלת קבצי טקסטים ומחזה מילים שימושיות בכל טקסט וממיינת את הטקסטים לפי הדמיון הלקסיקלי (מילולי) שלהם.

הקוד בחלק זה ימומש בחבילה `pdf8.cs1.ws1.ex8.tau.ac.il` אך ישמש גם בקוד של היסטוגרמה אחרת מימשתם בחלק א' (כלומר, ישמש בקוד שモופיע בחבילה אחרת – וודאו שני החלקים האלה מופיעים אצל אותו הפרויקט ב `eclipse`)

המחלקה שלנו תבצע את הניתוח באופן הבא: עבור תיקית קבצים כלשהי, היא תקרא את התוכן של כל קובץ ותשמר את הנתונים ההרלוונטיים לכל קובץ (אנחנו קוראים לזה – ביצוע אינדקס, indexing). פעולה זו תקרה פעם אחת בלבד, ולאחר מכן כל הפעולות תבצענה על תוכן האינדקס שתיצרו.

סעיף 1 (הבדיקה של סעיף מתבצעת דרך הסעיפים האחרים ולן לא ניתן לו ניקוד)

המתודה `indexDirectory` במחלקה `FileIndex` קוראת את הקבצים ומוסיפה אותם לאינדקס. המימוש של פונקציה זו נתון לכם חלונית ואתם רשאים לערוך אותו. קריאת המילים מן הקובץ תתבצע בעזרת `FileUtils.readFileAllTokens`, שכך נtauותה לכם. שימוש לב שהמחלקה `Zorah` שגיאה ואין לשנות זאת, ניתן במקרה זאת להיעזר ב `try catch` אם יש צורך בכך.

המטרה של שירות זה היא לקרוא את תוכן כל הקבצים, לנתח אותו ולשמור אותו כך שהשימוש של שאר השירותים במחלקה `FileIndex` יהיה יעיל ומהיר. בתרגיל זה, חלק מהשירותים צריכים להיות מאוד יעילים, כך שבשלב ה `index` נעשה את החישובים הבודדים ונשמר אותם מבני נתונים מתאימים. הגדרה נconaה של מבני הנתונים, והוצאת קוד משותף למתקודות פרטיות תהפוך את המימוש של חלק מהשירותים המוגדרים ב `fileIndex` לפעולות שליפה פשוטות מבני נתונים. אתם לא נדרדים על זמן הריצה של המתקודות, אך זהה שיטת העבודה הרצiosa בתרגיל זה.

שימוש לב, עלייכם לבחור את מבני הנתונים המתאימים לייצוג המידע הדרוש. לשם כך, עליכם לקרוא ולהבין את כל הסעיפים של חלק ב', ורק לאחר מכן לקבל את החלטה על מבני הנתונים שיישמשו אתם. עלייכם להשתמש ביעילות במבנה נתונים גנריים מתוך `collection framework Java`. בפרט, עליכם להשתמש במבנה הנתונים `HashMapHistogram` אשר מומש בחלק א' על מנת לשמור את מספר המופיעים של ה `token`-ים בכל קובץ.

הערות נוספות:

- שם תקין הקבצים יהיה שם חוקי של תיקיה המכילה לפחות קובץ אחד.
- השירות `readAllTokens` של `FileUtil` מבטל סימני פיסוק ומוחזיר מילים שאין ריקות, אין לבצע עיבוד או סינון נוסף בגין המימוש שלהם: כל המילים שחוזרות ע"י `readAllTokens` הן חוקיות מבחינתכם.
- הנacho כי כל קובץ מכיל לפחות מילה חוקית אחת.
- החוזה של כל שירות ב `fileIndex` נתון בלבד הקוד – קריית החוזה היא חלק מהתרגום ומהוा השלמה לפירוט הקצר שמווני עבור כל שירות.

סעיף 2 (10 נק')

משוו את שלושת השירותים הבאים:

1. השירות `getCountInFile` אשר מקבל מחרוזת `fileName` ומחרוזת `word` ומוחזיר את מספר המופיעים של המילה `word` בקובץ `fileName`. עבור מילה שאינה מופיעה בקובץ יוחזר ערך 0.
2. השירות `getNumOfUniqueWordsInFile` אשר מקבל מחרוזת `filename` ומוחזיר את מספר המילים השונות שמופיעות בקובץ `filename`.
3. השירות `getNumOfFilesInIndex` אשר מוחזיר את מספר הקבצים שמופיעים באינדקס.

חתימת השירותים:

```
public int getCountInFile(String word, String fileName) throws FileIndexException
public int getNumOfUniqueWordsInFile(String fileName) throws FileIndexException
public int getNumOfFilesInIndex()
```

הנחיות כלויות לסעיף זה והסעיפים הבאים:

- בכל שירות המתקבל שם של קובץ, המחרוזת `fileName` מכילה שם קובץ בלבד (לא נתיב), יש לחפש אותו בתיקיה עליה בוצע שלב ה `index` (ראו דוגמת שימוש בחלוקת הטסטר). עליים לקרוא את התוכן של כל קובץ פעם אחת בלבד, בשלב האינדקס. לאחר מכן, עליים לעשות שימוש באינדקס בשביל למש את השירותים של המחלוקת.
- בכל שירות המתקבל שם של קובץ, במידה ושם הקובץ אינם קיימים בתיקיה זו, יש לזרוק חריג מטיפוס `FileIndexException` (מומש עבורכם) עם הודעה אינפורמטיבית לבחירתכם.
- בכל שירות שמקבל מילה `word` יש להמירה ל `lowercase` לצורך ביצוע החיפוש באינדקס.

סעיף 3 (5 נק')

בסעיף זה נתחל למשמש את הממד **tf-idf** שבו נשתמש בהמשך התרגילים.
עבור אוסף מסמכים (כל מסמך מופיע בקובץ נפרד), נגידר את הממדים הבאים:
המדד (term frequency) tf מחושב עבור מילה (word) ומסמך (doc) באופן הבא:

$$tf(word, doc) = \frac{\text{number of repetitions of } word \text{ in } doc}{\text{number of words (with repetitions) in } doc}$$

המדד (inverse document frequency) idf מחושב עבור מילה (word) באופן הבא:

$$idf(word) = \log\left(\frac{\text{number of documents}}{\text{number of document containing } word}\right)$$

המדד **tf-idf** הוא מודד אשר נותן ציון רלוונטיות לכל מילה word במסמך doc, והוא מחושב ע"י מכפלת **tf** ב **idf**. ככל שמילה מסוימת מופיעה יותר פעמים במסמך מסוים, היא יותר משמעותית במסמך זה - לזה אחראי החלק של **tf**. מצד שני, אם המילה מופיעה בכל המסמכים שיש לנו במאגר, זה אומר שהיא פחותה "מיוחדת" עבור המסמך עליו אנחנו מסתכלים – ועל החלק זהה אחראי ה **idf**. אם המילה מופיעה בכל מהמסמכים במאגר, ה **tf** שלה הוא (1) log 0. לכן, המודד **tf-idf** יגדל ככל שהמילה מופיעה בפחות במסמכים.

מדד זה שימושי לכל מני משימות בעיבוד טקסטים, כולל באחזור (retrieval) של מסמכים ומיציאת דמיון בין מסמכים. הסבר מפורט יותר, כולל הסבר מלא לדוגמא שבה השתמשנו בחלק הראשון של תוכנית הבדיקות **FileIndexTester** ניתן למצוא כאן:

<https://medium.com/analytics-vidhya/tf-idf-term-frequency-technique-easiest-explanation-for-text-classification-in-nlp-with-code-8ca3912e58c3>

ממשו את השירות **getTF** אשר מקבל מילה **word** ושם של קובץ **fileName** ומחזיר את ערך ה **tf**.

חתימת השירות:

```
public double getTF(String fileName, String word) throws FileIndexException
```

העזר בשירות **calcTF** אשר מבצע את חישוב ה **TF** עבור שני הפרמטרים המציגים את המונה והמכנה בנוסחה. השירות מקבל שני שלמים ומחזיר **double**, והוא ניתן לכם על מנת שכל החישובים העשויים שלכם שלכם יבוצעו באותה הצורה.

סעיף 4 (5 נק')

משו את השירות `getIDF` אשר מקבל מילה `word` ומחזירה עבורה את ערך ה `idf`.

חתימת השירות:

```
public double getIDF(String word)
```

כמו בסעיף הקודם, נתונה לכם פונקציית עזר בשם `calcIDF` אשר מבצעת את החישובים העשויים ומקבלת שני שלמים המציגים את המונה והמכנה בשבר המופיע בנוסחתה `idf`.

בשלב זהה אתם יכולים להריץ את החלקה הראשונית בתוכנית הבדיקה `fileIndexTester`. החלק הראשוני מופיע בפונקציה `testFileIndexFirstPart`.

סעיף 5 (10 נק')

משו את השירות `getTopKMostSignificantWords` אשר מקבל שם של קובץ `filename` ומחזיר רשימה של k זוגות של מילה-ציוון `tf-idf`. הזוגות שיחזו יהיו זוגות של המילים היכי "משמעויות" (בעלות ציוון `tf-idf` היכי גבוהה) ב `filename`, מוסדרים בסדר יורד של ציוון `tf-idf`. אם יש שתי מילים בעלות ציוון `tf-idf` זהה, המילה שתופיע קודם היא המילה שמופיעה קודם בסידור לקובץ קוגרפי.

邏輯 של שימושה של שירות זה מצריך מיון, ומיין הוא פעולה יקרה, נרצה לבצע מיון יחיד עבור כל קובץ. אפשרות אחת היא לבצע את המיון בשלב האינדקס, ואפשרות נוספת היא לבצע את המיון בקריאה הראשונה לשירות זה עבור קובץ מסוים. בכל מקרה, עליכם להימנע מביצוע מיונים במשום שירות זה.

```
public List<Map.Entry<String, Double>>
    getTopKMostSignificantWords(String fileName, int k) throws FileIndexException
```

הערה: השירות `getTFIDF` כבר מומש עבורכם (שירות זה עושה שימוש בשירותים שמיימשתם בסעיפים 3+4).

סעיף 6 (10 נק')

נרצה לבדוק דמיון מילולי בין הקבצים השונים המופיעים באינדקס.

בתרג'il בית 5 יציגת מילים באמצעות מערכים של מספרים, וחישבותם cosine similarity בין ייצוגים של מילים שונות בשילוב למצוא מילים דומות מבחינת הקשרים שבהם הן מופיעות (שזהו בעצם דמיון סמנטי – דמיון במשמעות). כעת, נשתמש בשיטה דומה בשילוב לחשב דמיון בין מסמכים. על מנת להשתמש ב cosine similarity, علينا ליצג כל מסמך ע"י וקטור מספרי, ואת הייצוג הזה לבנה ע"י שימוש ב tf-idf.

ביצוג ע"י וקטור אנחנו צריכים לאוצר המילים (שקובע את גודל ה-וקטור) ולתוקן של כל תא. לשם פשוטות, נניח כי אוצר המילים שלנו כולל את כל המילים שהופיעו בכל הקבצים שופיעים באינדקס, וכן אורך וקטור יהיה כמספר המילים השונות שופיעו בכל הקבצים, נסמן מספר זה ב N.

כל מסמך יוצג, אם כך, ע"י וקטור באורך N. עבור מסמך doc כלשהו, התא ה-n בוקטור יכול את ערך ה tf-idf עבור המילה ה n באוצר המילים.

שימוש לב – אין באמת צורך לבנות וקטורים בגודל כל אוצר המילים. הנוסחה של cosine similarity היא:

$$\text{cosSim}(A, B) = \frac{\sum_{i=0}^n A_i * B_i}{\sqrt{\sum_{i=0}^n A_i^2 * \sum_{i=0}^n B_i^2}}$$

התאים שמעניינים אותנו בוקטורים שמייצגים את A ו B הם שבהם יש ערך שאינו 0 או ב A או ב B. אחרת, אם בתא ה n יש 0 גם ב A וגם ב B, הם לא משפיעים על תוצאה החישוב. בנוסף, אין חשיבות לסדר האיברים ב- וקטור. מה שצורך הוא לוודא שביחסוב המונה, הערך המתאים למילה a כלשהו במסמך A יוכפל בערך המתאים למילה a כלשהו במסמך B.

מסיבה זו, אתם לא נדרש למשר שירות אשר מייצג מסמך כווקטור. בחישובים שלכם אתם יכולים לעשות שימוש ישיר בהיסטוגרמות או במקרים אחרים שאתם מחזיקים בשילוב לחשב את ה cosine similarity.

ממשו את השירות getCosineSimilarity המקבל שמות של שני קבצים, ומחזיר את תוצאה חישוב ה cosine similarity בין ייצוגי df-idf שלהם.

```
public double getCosineSimilarity(String fileName1, String fileName2)
    throws FileNotFoundException
```

סעיף 7 (10 נק')

כעת, נרצה לקבל את המסמכים היכי דומים למסמך כלשהו על פי דמיון בקטורי fdf-tf . בדומה לסעיף 5, נממש שירות אשר מחזיר את \mathbf{A} המסמכים היכי דומים למסמך `filename` כלשהו, יחד עם הציון של מידת הדמיון (ערך ה *cosine similarity*).

חתימת השירות:

```
public List<Map.Entry<String, Double>>
getTopKClosestDocuments(String fileName, int k) throws FileIndexException
```

הנchina: ניתן, אך לא חובה (בניגוד לסעיף 5) לשמר את תוצאת המין עבור כל קובץ, כך שיבצעו פעם אחת בלבד. לחילופין, ניתן לבצע את החישובים בכל פעם מחדש לפונקציה (אנחנו מניחים שמספר הקבצים הוא קטן ולכן שמיירת תוצאות המינויים מראש לא תשפר משמעותית את זמן הריצה של השירות).

במציאות, את ההחלטה לגבי ביצוע פעולות בזמן ריצה מול שמיירת תוצאות החישובים נבצע לפי דרישות המערכת וגודל בסיס הנתונים. יש לנו 2 "ישויות" עליון אנחנו צריכים לבצע חישובים – חישוב ברמת המילה (המילים היכי "משמעותיות") וחישוב ברמת הקובץ (הקבצים היכי "דומים"), וכן גודל אוצר המילים ומספר הקבצים הם ממשמעותיים בקבלת החלטה. כמו כן, צריך להבין מה הסיכוי לעשוות שימוש חוזר בערכים שנשמרים מראש.

טיטרים:

כמו בכל תרגיל אחר, הטיטרים הם טיטרים בסיסיים שאינם בודקים את כל המקדים, ורק מוצלחת שלהם מהוות תנאי הכרחי אך לא מספיק בשבייל לוודא שהתרגיל שלכם עובד כנדרש. הוסיפו בדיקות משלכם!

(אין בעיה לשנות את קובץ הבדיקות הקיימים, רצוי להוסיף בדיקות חדשות ולא לשנות את הקייםות)

בהצלחה!