

איחזור מידע תשע"ג – 372.1.4406
סמסטר חורף מועד א' 19.01.2014
פרופ' ברכה שפירא, אביגיל פרדיס, יוסי בן-שלמה

משך המבחן: שעותיים וחצי
חומר עזר: מותר (לא מחשב נייד), מותר מחשבון

יש להחזיר את גיליון הבחינה. יש לענות על כל השאלות, למבחן 3 דפים

1. 16% יש לענות לשאלה זו על גבי גיליון הבחינה
סמן תשובה אחת- נכונה (נא לא לצרף הסברים - הם לא ייקראו)

- 1.1 4% החיסכון המשמעותי מביצוע stem נובע מ:
- חיסכון בגודל המילון ובזמן חיפוש עליו (dictionary)
 - חיסכון בגודל ובזמן חיפוש ב posting
 - חיסכון בזמן של ביצוע פעולות Text על המסמכים בשלב יצירת האינדקס
 - א+ב
 - כל התשובות נכונות

- 1.2 4% אתיקה של crawler של מנוע חיפוש מתבטאת ב:
- הפעלה של איסוף נתונים מאתרים על פי זמנים שנוחים לאתר על פי הגדרתו (של האתר)
 - ציות למה שמוגדר ב Robots.txt באתר
 - הגבלה של מספר הדפים הנסרקים בו זמנית מאותו אתר
 - השמטה של נתונים פרטיים של משתמשים מהדפים שנאספים
 - בקרה על תכיפות הגישה לאותו דף
 - כל התשובות נכונות
 - א+ב+ג נכונים
 - ב+ג+ה נכונים

- 1.3 4% בעת ביצוע ההחלקה לצורך חישוב דמיון של שאילתא למסמך באמצעות מודל שפה, משתמשים במודל השפה של כלל המסמכים במאגר של המנוע ולא בקבוע כלשהו או במאגר מסמכים חיצוני כלשהו. משום ש:
- לא כל מונחי השאילתה בהכרח מופיעים במסמך
 - רוצים לייצג באופן טוב את סביבת הנתונים בה אנחנו פועלים – כלומר את המאגר של המנוע
 - אמנם כל מאגר מסמכים גדול יתאים, אבל המאגר של המנוע זמין ללא צורך בעיבוד נוסף
 - השימוש במודל השפה של המאגר של המנוע חוסך את הצורך בהגדרת פרמטר חופשי נוסף
 - השימוש במודל השפה של המאגר של המנוע מונע הסתברות 0 למונח כלשהו
 - כל התשובות נכונות

1.4 4% כאלטרנטיבה ל IDF (inverse Document frequency) כמרכיב למשקול term הוצע להשתמש ב ICF (inverse collection frequency) – שהוא המספר ההופכי של כמות מופעים במאגר.

- א. הסיבה להעדיף IDF על פני ICF היא כי אפשר לחשב IDF מראש
 ב. אין שום עדיפות ל IDF על פני ICF
 ג. IDF עדיף משום ש ICF לא מדגיש ומבדיל את יחודיות ה term למסמך ויכול להיות מוטה מ terms שמופיעים הרבה פעמים במעט מסמכים.
 ד. ICF עדיף משום שהוא מודד כמה ה term ייחודי במאגר על פי מספר הופעותיו הכולל, ככל שהוא מופיע פחות, הוא יותר ייחודי למסמכים שבהם הוא מופיע.
 ה. א+ג נכונים

2. 16% על שני מנועי חיפוש: e1, e2 הורצו שתי שאילתות: Q1, Q2. ידוע שלשאלת Q1 20 מסמכים רלוונטיים במאגר, לא ידוע מספר המסמכים הרלוונטיים במאגר לשאלת Q2. להלן תוצאות ההרצה של השאלת לשני המנועים (+ מסמן מסמך רלוונטי שהוחזר, - מסמן מסמך לא רלוונטי שהוחזר):

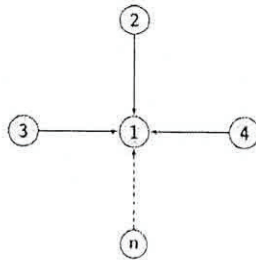
מנוע E2		מנוע E1		
Q2	Q1	Q2	Q1	
+	-	+	-	1
-	+	+	-	2
-	-	+	+	3
+	-	+	+	4
-	+	-	+	5
+	+	+	-	6
+	-	-	-	7
+	+	-	-	8
+	+	-	-	9
-	+	-	-	10

- א. 5% הסבר – האם (ואם כן, כיצד) ניתן להעריך את מספר המסמכים הרלוונטיים הנמצאים במאגר לשאלת Q2 (בהנחה שאפשר להריץ עוד שאילתות על כל אחד מהמנועים)?
 ב. 11% חשב, או הסבר איזה נתונים חסרים לחישוב, לערכים הבאים:
 1. 2% mean average precision על פני שתי השאלות למנוע E1.
 2. 2% DCG לשאלת Q1 למנוע E2.
 3. 2% E-measure לשאלת Q1 למנוע E1 כאשר ניתן משקל כפול ל precision לעומת ה recall.
 4. 2% R-precision למנוע E1.
 5. 3% Interpolated average precision למנוע E2 ושאלת Q1.

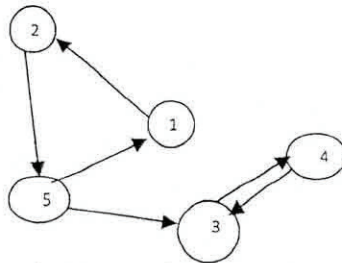
3. 15% חברה רוצה לפתח מנוע חיפוש שיחזיר תוצאות רלוונטיות אבל גם מגוונות. הכוונה בגיוון היא שהמסמכים החוזרים לא יהיו דומים יותר מדי אחד לשני, שהמשתמש ירגיש שהוא מקבל את כלנקודות המבט על השאלת ששאל. הגדר מדד חדש שיוכל לשמש להערכת המנוע הנ"ל (יש להסביר כיצד המדד שהוגדר משלב בין רלוונטיות וגיוון, וכיצד הוא מודד גיוון). יש להגדיר נוסחא מדויקת למדד.

4. 20% ניתוח הרשת :

4.1 5% הנח רשת במבנה של כוכב (לפי הצורך, יכולים להיות עד n צמתים לרשת). חשב את ה Hub וה $authority$ של הצמתים כפונקציה של n .



4.2 5% הסבר מה יהיו ערכי ה $pagerank$ של הרשת הבאה- הנח $d=0.85$. אין צורך לחשב, אלא להעריך מה יהיו הערכים.



4.3 10% הסבר האם אלגוריתם HITS מצריך טיפול מיוחד ב $dead\ ends$ וב $leaks$ (spider trap) בדומה לאלגוריתם $pagerank$ (תן הסבר גם אם התשובה חיובית או שלילית).

5. Spelling 13%

5.1 6% מה המרחק (levenstein distance) בין זוגות המילים (חשבו את המעבר של המילה הכתובה משמאל למילה הכתובה מימין):

apfel, apple
apfel, apples

פרט את הפעולות שבצעת על כל זוג מילים כדי לחשב את המרחק.

5.2 7% מצא מילה נוספת עם מרחק זהה למרחק שבין $apple$ ו $apfel$. הנח שמשמש טעה (שגיאות spell) שלח את $apfel$ כשאיילת למנוע. בהנחה שרק שתי המילים ($apple$ והמילה השנייה שמצאת) נמצאות במילון של מנוע. הסבר בדיוק איך תדע לאיזו מילה מבין שתי המילים האלה התכוון המשתמש לשלוח כשאיילת.

6. 8% נניח שמאגר מסמכים כולל רק 4 מילים שכיחות: a, b, c, d (ואף לא מילה אחרת). סדר השכיחויות: $a > b > c > d$. מספר המילים הכולל במאגר הוא 5000. הנח שהמאגר מתנהג בדיוק לפי חוק Zipf. מה השכיחות של כל אחת מהמילים במאגר?

7. 12% כתוב אלגוריתם לחישוב הדמיון בין $terms$ במאגר (בין ה $terms$ לעצמם) (יש לכתוב פסאודו-קוד מדויק), תן דוגמא לשימוש באלגוריתם כזה.

בהצלחה – ברכה אביגיל ויוסי