

איחזור מידע תשע"א – 372.1.4406
סמסטר חורף מועד ב' 17.02.11
ד"ר ברכה שפירא, אורלי מורנו

משך המבחן : שעתיים וחצי
חומר עזר: מותר (לא מחשב נייד)
יש להחזיר את השאלון

מס נבחן 02 30 93

חלק א – נא לענות במחברת הבחינה

1. 26% נתון log של שאליות שמשתמשים שלחו למנוע, ולכל שאליתא סט של 30 מסמכים שחזרו לשאליתא. כמו כן, נתון לכל שאליתא על אילו מסמכים מתוך אלו שחזרו המשתמש הסתכל. כדי לנתח את האיכות של המנוע צריך לזהות מתי המשתמש חיפש באותו נושא – כלומר מבין השאליות של המשתמש- אילו מהן היו בהקשר לאותו צורך מידע בהנחה שאין זיהוי של משתמש על השאליתא (לא ידוע אילו מהשאליות התבצעו על ידי אותו משתמש ברצף).

א. 18% הצע פיתרון (אלגוריתם בפסאודו קוד) שיקבץ את השאליות של משתמש ברצף ל session – , כלומר לזהות באותו רצף את השאליות של משתמש אחד המנסה למצוא מידע בנושא מסוים ומתי מתחיל session של משתמש אחר (או אותו משתמש בנושא שונה).

ב. 8% הסבר מהן המגבלות של הפתרון – כלומר באילו מקרים הוא לא יצליח לזהות את ה Session (אם אין כאלה מה טוב..... הסבר מדוע אין מגבלות)

2. 8% נניח crawler שעובד ב batch ושומר את ה Repository באופן מבוזר. הביזור מתבצע תוך כדי ה crawling כאשר ה crawler מוריד את הדף ואז שומר אותו בצומת לפני הטיפול שמתבצע על הדף ("טיפול" למשל הוצאת הלינקים וכו'). מנה שתי סיבות שבגללן עדיף ל crawler לצורך ייעול עבודתו, לבזר את הנתונים בצמתים כך שבצומת אחת יישמרו כל הדפים של אתר מסוים – על פני ביזור של דפים שונים של אותו אתר בצמתים שונים. (הסיבות צריכות להיות קשורות לייעול עבודתו של ה crawler – לא לאינדוקס בהמשך).

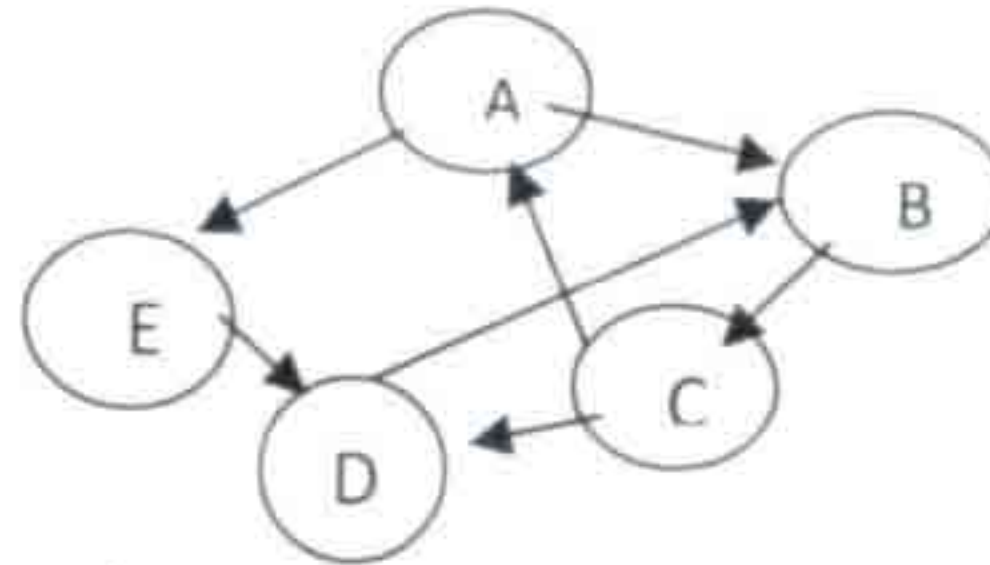
3. 16% למאגר מסוים ואוסף שאליות התבצעה הערכה של מומחים על הרלוונטיות של כל מסמך לשאליתא (לצורכי הערכה של מנועים). ההערכה כללה 3 סוגי תוצאות (במקום ה 2 הסטנדרטיות) : 0 – כאשר המסמך לא רלוונטי, 1 כאשר מסמך רלוונטי חלקית ו 2- כאשר המסמך רלוונטי לגמרי.

- א. 13% כתוב נוסחה מדויקת שמגדירה את מדד precision@k בהתאמה לתוצאות כאלו.
- ב. 3% הדגם את המדד precision@5 שהגדרת על התוצאות הבאות (משמאל לימין):

1,0,2,0,0,2,1,1

חלק ב – יש לענות במחברת הבחינה במקומות המסומנים

4. 15% נתונה רשת :



השלם: (תשובה לא נכונה קונסת בחצי נקודה – כלומר על תשובה לא נכונה יופחתו 4.5 נקודות).

5% הצומת בעלת ערך ה hub הגבוה ביותר על פי אלגוריתם HITS הוא: A
 5% בין הצמתים E, B, D הצומת בעלת ערך authority נמוך ביותר על פי אלגוריתם HITS היא E
 5% 2 הצמתים בעלי ה pagerank הגבוה ביותר הם B, D

סמן נכון או לא נכון (הסמנים במחברת הבחינה)

5. 17%

- א. 3% ערך $tf*idf$ של term במסמך יכול להיות גבוה יותר מ-1. נכון/לא נכון
 ב. 3% שימוש ב Stemmer יכול להשפיע על שביעות רצון המשתמש מהמנוע. נכון/לא נכון
 ג. צומת בגרף מסוים שהיא בעלת ה pagerank הגבוה ביותר בגרף, בהכרח תהיה גם בעלת ה Authority הגבוה ביותר. נכון/לא נכון
 ד. 4% הנח שתי מילים "teach" ו "taught" ששקללו אותן כשתי מילים נפרדות במסמכים בהן הן הופיעו, על פי $tf*idf$ סטנדרטי. לאחר מכן הסתבר ששתי המילים הן בעצם stem של אותה מילה ושאפשר לייצג אותן במשותף כמילה אחת (בהנחה שהמנוע תומך ב stem). כדי לתקן את הטעות, צריך לסכום את ה $tf*idf$ שלהן לייצוג המשותף. נכון/לא נכון

ה. 4% נתון פרופיל משתמש לאחר עדכון על פי תגובת משתמש על פי אלגוריתם rocchio לאחר ניתוח תגובתו למסמכים d1 d2 d3. בהנחה שאין הפחתה שלילית (כלומר $\gamma=0$). אפשר לראות שהמשתמש העדיף פוליטיקה ומדע על פני מוזיקה. נכון/לא נכון

ספורט	מוזיקה	פוליטיקה	מדע	רמטכ"ל	לימודים
0.3	0.2	0.8	0.8	0	פרופיל לאחר עדכון
0.2	0	0.8	0.9	0	מסמך d1
0	0.9	0.2	0	0.1	מסמך d2
0.4	0	0.8	1	0	מסמך d3

6. 6% נתון מאגר ובו 4 מסמכים :

D1- Tibet Tibet Malaga

D2- Malaga Rimini Salvador Tibet Tibet Tibet

D3- mexico Sun

D4- Mexico Malaga Tibet Sun

נתונים שלושה מנועים – E1 עובד לפי המודל הבוליאני הטהור, E2 לפי המודל הווקטורי, E3 לפי מודל בוליאני מורחב שמדרג את המסמכים לאחר הפעלת האופרטורים הבוליאנים בשאלתא.

השאלתא Malaga and Tibet נשלחה ל E1 ול E3

השאלתא Malaga Tibet נשלחה ל E2

- א. 3% E1 ו E2 יחזירו תשובה זהה לשאלתות שנשלחו אליהם. נכון/לא נכון
 ב. 3% E1 ו E3 יחזירו את אותם מסמכים לשאלתא שנשלחה אליהם. נכון/לא נכון



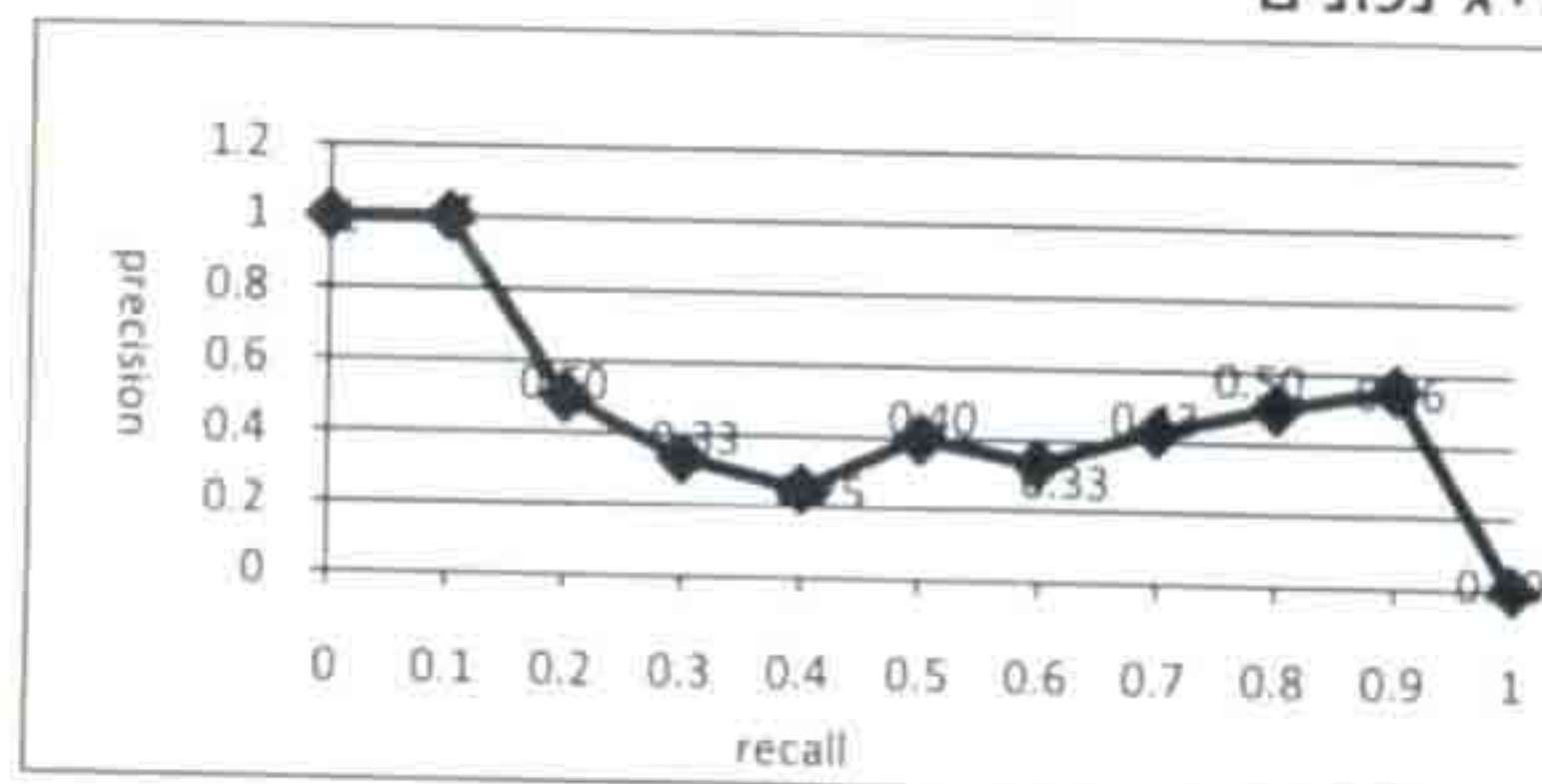
12% סמן תשובה אחת נכונה לשאלות הבאות

7. 4% מנועים מסוימים משתמשים בתגובות משתמשים לתשובות על שאלת q1 ומריצים בעקבותיה שאלת q2 נוספת q2. (relevance feedback). שיטה זו:
- משפרת תמיד את ה Recall של שאלת q1
 - עשויה לשפר את ה precision של שאלת q1
 - אינה משפיעה על ביצועי המנוע אלא רק משפרת את שביעות רצון המשתמשים
 - משפרת precision ו recall של שאלת q2

8. 4% לשאלת Q 4 מסמכים רלוונטים במאגר. להלן תוצאות שחזרו משני מנועים: E1 ו-E2. התוצאות משמאל לימין כשאר N מסמל מסמך לא רלוונטי, R מסמן מסמך רלוונטי:
- $R = 4$
- E1: R N R N N N N R R
 - E2: N R N N R R R N N

- ל E1 MAP גבוה יותר מאשר ל E2
- R-precision של שני המנועים שווה
- R-precision של E1 גבוה מ R-precision של E2
- MAP של E2 גבוה מ MAP של E1
- א+ג נכונים
- א+ב נכונים
- ג+ד נכונים

9. 4% נתון גרף precision-recall. הגרף משקף תוצאות של שאלת אחת.
- הגרף התקבל בהכרח מנתונים אמיתיים ללא אינטרפולציה
 - הגרף התקבל בהכרח לאחר אינטרפולציה של ערכים
 - כל המסמכים הרלוונטים לשאלת חזרו.
 - 90% לפחות מהמסמכים הרלוונטים לשאלת חזרו
 - ב+ד נכונים
 - ב+ג נכונים



בהצלחה

ברכה ואורלי

