

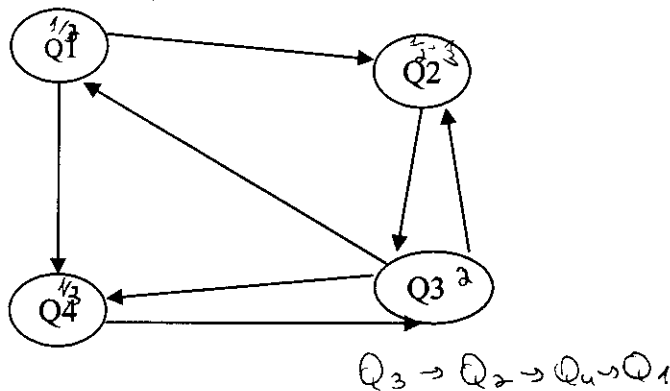
$$\begin{array}{r} 28 \\ 64 \\ \hline 48 \end{array} \quad \begin{array}{r} 62 \\ 48 \end{array}$$

המחלקה להנדסת מערכות מידע והתוכנית להנדסת תוכנה
מבחן בקורס איחזור מידע 372.1.4406 מועד א'
ד"ר ברכה שפירא, ארז שלום
13:30 – 22.07.07

משך הבחינה: שעתיים וחצי
חומר עזר מותר – לא מחשב נייד!

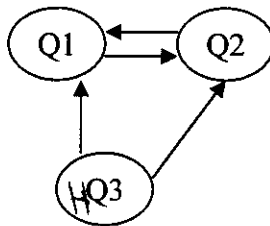
יש לענות על כל השאלות

1. 20%
I. 10% נתון הגרף הבא המתאר רשת של דפי אינטרנט- יש לדרג את הצמתים על פי ערכי ה pagerank שלהם. אין צורך לבצע חישוב של pagerank של כל צומת, אלא לנמק את הדרוג במדויק על פי ההגיון של הנוסחה.



- II. 4% הסבר כיצד לשנות את נוסחת ה pagerank כך שתהיה העדפה מלאכותית של אתרים מסויימים בדרוג על פני אחרים.
III. 6% נתון הגרף הבא המייצג רשת דפי אינטרנט: הפעל את אלגוריתם HITS ודרג את האתרים על פי ערכי ה authority וה hub- שלהם (דירוגים נפרדים ל hub – ו authority). אין צורך לחשב את הערכים המדויקים רק לנמק את הדירוג.

H $Q_3 \rightarrow Q_1 \rightarrow Q_2$
A $Q_1 \rightarrow Q_2 \rightarrow Q_3$



2. 30% נתון אינדקס הופכי. כל כניסה באינדקס עבור term מסויים כולל מופעים בפורמט הבא:
 $x:y,z$.. כאשר המשמעות היא שבמסמך x – term מופיע במיקומים y,z.

$\frac{2}{3}$ a: 1:3, 2:1,4
 everywhere: 3:1
 go: 4:6
 had: 1:2
 $\frac{1}{2}$ -lamb: 1:5,2:3,6,4:2
 $\frac{2}{3}$ little: 1:4, 2:2,5
 -mary: 1:1,3:3
 sure: 4:4
 that: 3:2
 the: 4:1
 to: 4:5
 was: 4:3
 went: 3:4

1: 1, 1,

1 1111 5
 2 1111 6
 3 111 4
 4 1111 6

$$J_3 = (0, \frac{1}{4} \cdot 4, 0, 0, 0, 0, \frac{1}{4} \cdot 2, 0, \frac{1}{4} \cdot 4, 0, 0, 0, \frac{1}{4} \cdot 4)$$

- א. 5% חשב את ה IDF של כל Term (יש להתעלם מה log בנוסחה)
 ב. 5% הראה את הווקטור של מסמך 2 (ערכי tf-idf – שוב ללא log – כאשר tf מנורמל לאורך המסמך)
 ג. 10% דרג את המסמכים במאגר לשאילתא: "mary lamb" על פי נוסחת ה cosine (ללא log בחישוב Tf-idf). בשאילתא- משקל כל term הוא 1.
 ד. 10% משתמש החזיר feedback למנוע על המסמכים שהמנוע החזיר כרלוונטים. המשתמש סימן את מסמך 2 כרלוונטי ואת מסמך 3 כלא רלוונטי. עדכן את השאילתא בהתאמה ל - feedback על פי נוסחת rocchio כאשר $\alpha=1, \beta=0.5, \gamma=0.25$

3. 20%

- I. 6% סמן על כל אחד מהמשפטים הבאים "נכון" או "לא נכון" ונמק במשפט אחד את תשובתך:

- i. הפעלת stemming תמיד מעלה את ה recall
 ii. הפעלת stemming תמיד מורידה את ה precision
 iii. כדי לשפר את תהליך האיחזור ניתן לבצע stemming רק למסמכי המאגר ללא stemming לשאילתא.
 II. 14% להלן סט של 3 חוקים מתוך אלגוריתם ה Stemming של פורטר. החוקים מיוצגים על ידי תבניות של $s1 \rightarrow s2$ (s1 מוחלף ב s2) בביצוע ה stemming מתבצע תמיד רק אחד מהחוקים בסט, על פי ההתאמה ה"ארוכה ביותר" לתבנית (longest match).

IES->I

SS->SS

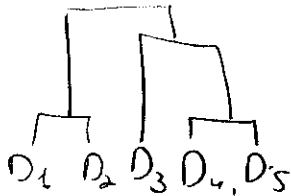
S-> null

- i. 3% מדוע יש צורך בחוק: SS->SS
 ii. 3% הראה את תוצאות ה Stem על המילים הבאות: PONIES CARESS, TIES
 iii. 4% אילו חוקים יש להוסיף לסט כדי לבצע stem נכון על המילים הבאות: CARESSES, PONY
 iv. 4% תוצאת ה Stem של המילים PONIES ו-PONY יוצרות מילים לא "תקניות" באנגלית. האם לעובדה זו יש השפעה על תוצאות האיחזור? – נמק (במשפט)

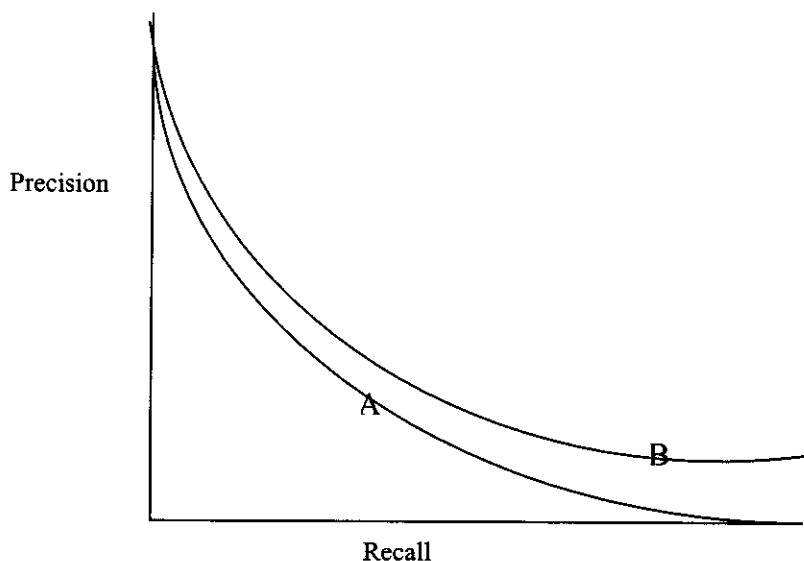
4. clustering: 14%

I. נתונה מטריצת הדמיון (לא מטריצת מרחקים!) הבאה בין מסמכים:
הראה את ה Clusters ההיררכיים הצוברים (דנדוגרם) שיווצרו מהמטריצה הזו כאשר הדמיון
בין ה clusters מחושב על פי complete link

D1	D2	D3	D4	D5	
1.00	0.42	0.00	0.1	0.1	D1
-	1.00	0.00	0.07	0.2	D2
-	-	1.00	0.33	0.33	D3
-	-	-	1.00	0.36	D4
-	-	-	-	1.00	D5



5. 16% נתונים גרפים Precision-Recall של שתי מערכות איחזור.
סמן על כל אחד מהמשפטים הבאים המתייחסים לתמונה: "נכון", "לא נכון" או "אי אפשר
לדעת על פי הגרף" (ללא הסבר)



- I. מערכת B יעילה יותר ממערכת A
- II. מערכת B איכותית (אפקטיבית) יותר ממערכת A
- III. מערכת A איכותית יותר ממערכת B למאגרי מידע סטטיים (כגון מאגרים של רפואה)
- IV. מערכת A לא הצליחה למצוא את כל המסמכים הרלוונטים לשאלות שאותן הגרף מתאר

בהצלחה
ברכה וארז

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

4

(3)

(1)

I - הפעלה - stemming לכן מזהה את ה recall משהו שהוא
 גורם לעצם מסמכים למחצה. בלב סל של מסמכים יש אתר מסמכים
 של בלונטלים שזה נמצא, לכן בהכרח, במידה הרלוונטיות נאמרה סה"כ (ללא)
 יש מקרים שגורמים את התבונה: כגון, אם ה stem לא נאמרה מסמכים או
 שפיר ל המסמכים הרלוונטיות נמצא לפני ה stem, אז ה recall
 לא מזהה, מי שזמן מקרים עלו (או מקרי קרה אחרים לפני) קהל את
 התקלה.

II stem - מזהה מסמכים נכונים, ושכחתי חלקם רלוונטיות, אז
 נמצא ה precision ירד, תלוי בקיום שבין הרלוונטיות ללא רלוונטיות לפני
 ה stem לרוב אחריו. למשל אם לפני ה stem הן
 נמצא אתר מתוך 2 (בלוטי) $PR = \frac{1}{2}$ בסתירה ה stem חזרה אל
 5 מסמכים שמסמך 4 רלוונטי, היחס יהיה $\frac{5}{7}$ שזה יותר מ $\frac{1}{2}$.

III stemming ללא סמנטי (הן למסמכים) אין חלק, כי לפני המילה
 בין השלל למסמך לכן במקום שזה לא מספר.
 (2)
 I החלק $ss \rightarrow ss$ מזהה בין מילים קרובות (כגון students)
 שפיר תהיה קרובים את S היקפי לפני מילים שיש להם שני S
 כחלק מהמילה (בלוטי) ה-S היא לא של היקפיו ולכן לא נכרה
 תהיה S.
 II $caress \rightarrow caress$, $ponies \rightarrow pony$, $tigs \rightarrow ti$;
 III בין המילים: $I \rightarrow Y$, $sses \rightarrow ss$ (למשל $es \rightarrow es$ בלתי קושר כמו dates)
 IV - אין השפעה של התוצאה, לא משהו מה מילים בשנייה, השפעה שהשפעה למסמך

שאלה 8

(b) מספר המסמכים $N = 4$ (הנוסחא הזו) $idf_j = \log_2(N/df_j) : T_j$
 N היא מספר המסמכים במילר, df_j הינו מספר המסמכים 12 term j מופיע.
 $idf_j = \frac{N}{df_j}$ כל המילים מה-log נקרא

הביטויים המופיעים רק במסמך אחד הינם:
 everywhere, go, had, sure, that, the, to, was, want.

הביטויים המופיעים בשני מסמכים הם:
 $\frac{4}{1} = 4$ זהו קבוצה 15 ה-idf יהיה ככה וזהו: $\frac{4}{1} = 4$

a, little, mary

lamb

המילים ה-idf הינו: $\frac{4}{2} = 2$
 הביטויים המופיעים בשלושה מסמכים הינם:
 זהו ה-idf הינו $\frac{4}{3} = 1 \frac{1}{3}$

(2) אורך מספר 6 הינו 6. ישנם שלושה ביטויים המופיעים 12. נחשב tf זהו:

$$a: tf = \frac{2}{6} = \frac{1}{3} \quad tf \cdot idf = \frac{2}{3}$$

$$lamb: tf = \frac{2}{6} = \frac{1}{3} \quad tf \cdot idf = \frac{1}{3} \cdot \frac{4}{3} = \frac{4}{9}$$

$$little: tf = \frac{2}{6} = \frac{1}{3} \quad tf \cdot idf = \frac{1}{3} \cdot 2 = \frac{2}{3}$$

לפי שאת הביטויים שלא מופיעים במסמך $tf=0$ וכן $tf \cdot idf=0$.

אין וקטור המסמך "ראש כן": $Vd_2 = (\frac{2}{3}, 0, 0, 0, \frac{4}{9}, \frac{2}{3}, 0, 0, 0, 0, 0, 0, 0, 0)$

$$CosineSim(d_j, q) = \frac{d_j \cdot q}{|d_j| \cdot |q|}$$

(3) נוסחה ה-cosine:

$$d_1 \cdot q = \frac{1}{3} \cdot 2 + \frac{1}{3} \cdot \frac{4}{3} = \frac{2}{3}$$

$$d_j \cdot q = \sum w_{ij} \cdot w_{iq}$$

$$d_2 \cdot q = 0 + \frac{4}{9} = \frac{4}{9}$$

$$w_{ij} = tf \cdot idf: זהו$$

$$d_3 \cdot q = \frac{1}{4} \cdot 2 + 0 = \frac{1}{2}$$

$$d_4 \cdot q = 0 + \frac{1}{6} \cdot \frac{4}{3} = \frac{2}{9}$$

זהו המכנה נחשב $|q|$

$$|q| = \sqrt{\sum w_{iq}^2} = \sqrt{1+1} = \sqrt{2}$$

2021

W-8

$$|d_1| = \sqrt{\left(\frac{1}{5} \cdot 2\right)^2 + \left(\frac{1}{5} \cdot 4\right)^2 + \left(\frac{1}{5} \cdot \frac{4}{3}\right)^2 + \left(\frac{1}{5} \cdot 2\right)^2 + \left(\frac{1}{5} \cdot 2\right)^2}$$

$$= \sqrt{\frac{4}{25} + \frac{16}{25} + \frac{16}{225} + \frac{4}{25} + \frac{4}{25}} = \sqrt{\frac{28}{25} + \frac{16}{225}} = \sqrt{\frac{268}{225}}$$

$$|d_2| = \sqrt{\left(\frac{1}{3} \cdot 2\right)^2 + \left(\frac{1}{3} \cdot \frac{4}{3}\right)^2 + \left(\frac{1}{3} \cdot 2\right)^2} = \sqrt{\frac{4}{9} + \frac{16}{81} + \frac{4}{9}} = \sqrt{\frac{8}{9} + \frac{16}{81}} = \sqrt{\frac{88}{81}}$$

$$|d_3| = \sqrt{\left(\frac{1}{4} \cdot 4\right)^2 + \left(\frac{1}{4} \cdot 2\right)^2 + \left(\frac{1}{4} \cdot 2\right)^2 + \left(\frac{1}{4} \cdot 4\right)^2} = \sqrt{1 + \frac{1}{4} + 1 + 1} = \sqrt{\frac{13}{4}}$$

$$|d_4| = \sqrt{\left(\frac{1}{6} \cdot 4\right)^2 + \left(\frac{1}{6} \cdot \frac{4}{3}\right)^2 + \left(\frac{1}{6} \cdot 4\right)^2 + \left(\frac{1}{6} \cdot 4\right)^2 + \left(\frac{1}{6} \cdot 4\right)^2 + \left(\frac{1}{6} \cdot 4\right)^2}$$

$$= \sqrt{\frac{4}{9} + \frac{4}{81} + \frac{4}{9} + \frac{4}{9} + \frac{4}{9} + \frac{4}{9}} = \sqrt{\frac{20}{9} + \frac{4}{81}} = \sqrt{\frac{184}{81}}$$

$$\text{CosSim}(d_1, q) = \frac{\frac{2}{3}}{\sqrt{\frac{268}{225} \cdot 2}} = 0.431$$

$$\text{CosSim}(d_2, q) = \frac{\frac{4}{9}}{\sqrt{\frac{88}{81} \cdot 2}} = 0.301$$

$$\text{CosSim}(d_3, q) = \frac{\frac{1}{2}}{\sqrt{\frac{13}{4} \cdot 2}} = 0.196$$

$$\text{CosSim}(d_4, q) = \frac{\frac{2}{9}}{\sqrt{\frac{184}{81} \cdot 2}} = 0.104$$

d_1
 \downarrow
 d_2
 \downarrow
 d_3
 \downarrow
 d_4

רשימת המילים Cosine - המילה יחידה:

$$q_m = q_0 + 0.5 \cdot \frac{1}{|D_r|} \cdot \sum_{d_j \in D_r} d_j - 0.25 \cdot \frac{1}{|D_{nr}|} \cdot \sum_{d_j \in D_{nr}} d_j$$

③

$|D_r| = 1$ זכרון ארוך, זכרון קצר
 $|D_{nr}| = 1$ זכרון ארוך, זכרון קצר

$$q_m = q_0 + 0.5 \cdot d_2 - 0.25 \cdot d_3$$

זכרון ארוך, זכרון קצר, זכרון ארוך, זכרון קצר

$$d_3 = (0, 1, 0, 0, 0, 0, \frac{1}{8}, 0, 1, 0, 0, 0, 1)$$

$$d_2 = (\frac{2}{3}, 0, 0, 0, \frac{4}{9}, \frac{2}{3}, 0, 0, 0, 0, 0, 0, 0)$$

$$d_0 = (0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0)$$

↓

$$q_m = (\frac{1}{2} \cdot \frac{2}{3}, -\frac{1}{4}, 0, 0, 1 + \frac{1}{2} \cdot \frac{4}{9}, \frac{1}{2} \cdot \frac{2}{3}, 1 - \frac{1}{4} \cdot \frac{1}{2}, 0, -\frac{1}{4}, 0, 0, 0, -\frac{1}{4})$$

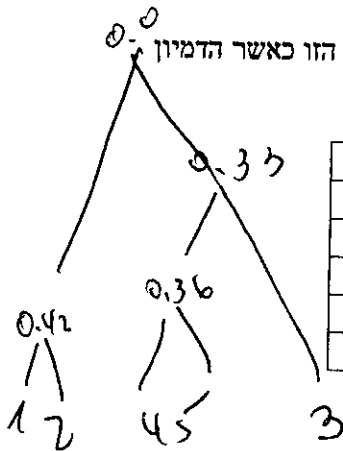
$$q_m = (\frac{1}{3}, -\frac{1}{4}, 0, 0, \frac{11}{9}, \frac{1}{3}, \frac{7}{8}, 0, -\frac{1}{4}, 0, 0, 0, -\frac{1}{4})$$

4. clustering: 14%

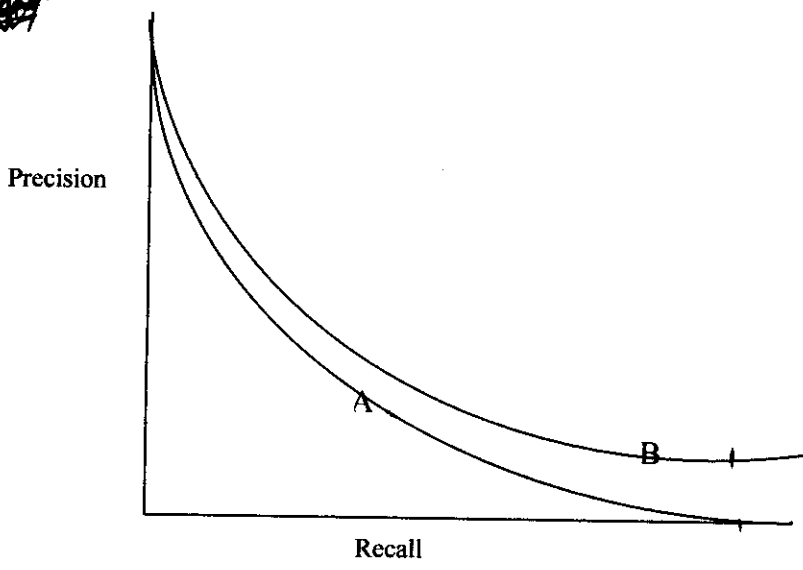
I. נתונה מטריצת הדמיון (לא מטריצת מרחקים!) הבאה בין מסמכים:

הראה את ה Clusters ההיררכיים הצוברים (דנדוגרם) שיווצרו מהמטריצה הזו כאשר הדמיון בין ה clusters מחושב על פי complete link

	D1	D2	D3	D4	D5	
D1	1.00	0.42	0.00	0.1	0.1	D1
D2	-	1.00	0.00	0.07	0.2	D2
D3	-	-	1.00	0.33	0.33	D3
D4	-	-	-	1.00	0.36	D4
D5	-	-	-	-	1.00	D5



5. 16% נתונים גרפים Precision-Recall של שתי מערכות איחזור. סמן על כל אחד מהמשפטים הבאים המתייחסים לתמונה: "נכון", "לא נכון" או "אי אפשר לדעת על פי הגרף" (ללא הסבר)



I. מערכת B יעילה יותר ממערכת A - אי אפשר לדעת

II. מערכת B איכותית (אפקטיבית) יותר ממערכת A
III. מערכת A איכותית יותר ממערכת B למאגרי מידע סטטיים (כגון מאגרים של רפואה) לא נכון
IV. מערכת A לא הצליחה למצוא את כל המסמכים הרלוונטים לשאלות שאותן הגרף מתאר נכון

בהצלחה ברכה וארו
= ואזיל מתחמם לקיומם של מהירות, שטח לחסון וכו'.
אילו אי אפקטיואל מתחמם לאילם וקלאם האחר, בלוח
ph ו Recall.

= מניין לאחרי A ופ $PR=0$ במקרה ה Recall הזלזל
(קריב 1-1 וק-1) היא אל הרלויה לקרא לא כל האסמכה