

איחזור מידע תשע"א – 372.1.4406

מסטר חורף מועד א' 31.01.11

ד"ר ברכה שפירא, אורלי מורנו

מספר נבחן 28/1

משך המבחן : שעתיים וחצי

חומר עזר: מותר (לא מחשב נייד) – יש להחזיר את הטופס

חלק א.  
הנחיה: 10 מסמכים (מתוכם, חלקם רלוונטיים). הוחלט  
לנסות לשפר את ביצועי המנוע על ידי הגדלת מספר המסמכים המוחזרים לכל שאילתא מ  
10 ל 1000, כלומר, יוחזרו עוד 990 מסמכים על 10 המוחזרים. (המסמך נשאר).

סמן תשובה אחת נכונה – יש לסמן תשובות על טופס המבחן.

1. 4% מנוע חיפוש מחזיר לכל שאילתא 10 מסמכים (מתוכם, חלקם רלוונטיים). הוחלט לנסות לשפר את ביצועי המנוע על ידי הגדלת מספר המסמכים המוחזרים לכל שאילתא מ 10 ל 1000, כלומר, יוחזרו עוד 990 מסמכים על 10 המוחזרים. (המסמך נשאר).

אם אף מסמך נוסף לא זוהה כרלוונטי בקבוצת המסמכים הנוספת המוחזרת אז :

- א. ה- Precision ו ה Recall של המנוע יישארו אותו דבר גם לאחר הוספת המסמכים לתוצאה  
ב. לאחר ההוספה, ה precision יקטן ו ה Recall יישאר אותו דבר  
ג. ה Precision וה Recall יקטנו לאחר ההוספה  
ד. ה r-precision יקטן  
ה. אף תשובה לא נכונה
- Recall = נכונ  
Precision = נכון

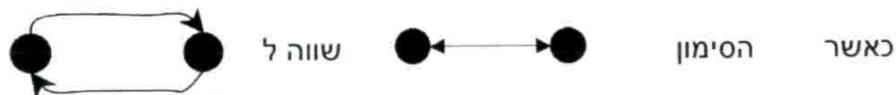
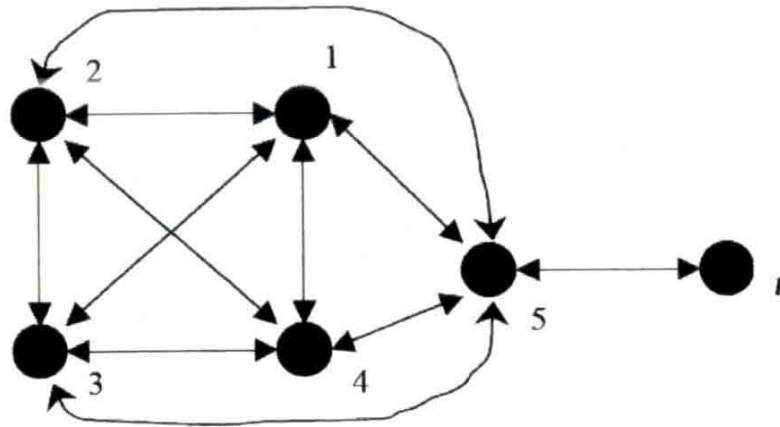
2. 4% (בהמשך ל-1). אם לפחות מסמך אחד נוסף זוהה כרלוונטי בקבוצת המסמכים הנוספת שהוחזרה, אז:

- א. ה precision לאחר ההוספה יהיה בהכרח נמוך יותר  
ב. ה precision לאחר ההוספה יהיה בהכרח גבוה יותר  
ג. ה Recall לאחר ההוספה יהיה בהכרח גבוה יותר  
ד. ב-ג נכונים  
ה. ה precision ב-10 מסמכים יגדל
- Recall = נכון  
Precision = נכון

3. 4% Interpolated average precision הוא מדד יעיל ל-  
א. השוואה בין מנועים ששונים רק בסדר המסמכים שמוחזרים

- ב. השוואה בין מנועים שמתחשב גם ב Recall  
ג. השוואה רק בין ה precision של המנועים  
ד. ציור של גרף precision-recall על פני הרבה שאילתות  
ה. ב+ד נכונים

4. 8% נתונה הרשת הבאה :



לרשת שכללה רק את הצמתים 1-5, נוספה הצומת  $t$ , כיצד ישתנו ערכי ה  $\text{pagerank}$  של הצמתים ברשת בעקבות ההוספה:

- כל ערכי ה  $\text{pagerank}$  של כל הצמתים ישארו שווים לערכם לפני ההוספה, גם ערך צומת  $t$  יהיה שווה לערכי צמתים. ☒ א.
- הערך של הצמתים 1-4 יקטנו במידה שווה. ☒ ב.
- הערך של צומת 5 יגדל ויהיה הגבוה ביותר ברשת. ☒ ג.
- הערך של צומת  $t$  יהיה הקטן ביותר ברשת. ☒ ד.
- הערך של צומת 5 יקטן. ☒ ה.
- ה+ב נכונים. ☒ ו.
- ב+ג+ד נכונים. ☒ ז.
- רק ב+ג נכונים. ☐ ח.

5. 16% ציין נכון או לא נכון ליד כל אחד מהמשפטים הבאים:
- 4% במנוע חיפוש באינטרנט F-measure הוא מדד חשוב. ☒ א. נכון
- 4% נניח מאגר מסמכים המכיל את המסמך השגוי הבא: "שלום טמבל מה העניינים ואיך אתה מרגיש" - נניח מנוע הפועל לפי המודל הווקטורי. נניח שהטעות (מופע של המילה "טמבל" לא רלוונטית למסמך) התגלתה. תיקון הטעות ישפיע רק על המסמך השגוי. ☒ ב. נכון
- 4% שאילתא הכוללת מילה שיש לה מילים נרדפות רבות בהכרח תסבול מ  $\text{precision}$  נמוך. ☒ ג. נכון
- 4% תור של Crawler של מנוע המיועד לבניית מאגר בנושא כלשהו יכלול רק דפים הרלוונטיים לנושא המאגר. ☒ ד. נכון
- המאגר הסיכומים
- topic driven

חלק ב.  
ענה על השאלות הבאות:

6. 7% נתונה טבלת הדירוגים הבאה של 5 משתמשים (user1..user5) לחמישה items (A..E). כמו כן נתונה הקרבה  $(W_{5,i})$  בין כל אחד מה-users ל user5 (המוגדר כ Active user). יש לבנא את הרלוונטיות של D item ל user 5 על פי שני השכנים הקרובים אליו ביותר (יש לחשב את  $(P_{5,D})$ ).

Item	User1	User2	User3	User4	User5 (active user)
A	10	5	9		9
B	6	9		5	5
C	2	7	3		1
D	4	8	3	3	2.526
E	8	1	9	2	
דמיון ל user5 $W_{5,i}$	1	-0.5	0.9	0.7	

7. 15% הצע מדד Discount Cumulative Gain (DCG) קבוצתי, שמפחית ציון של מסמכים בקבוצות של 5. כלומר הפחתה מסוימת ל- 5 המסמכים הראשונים, הפחתה מוגדלת ל-5 מסמכים השניים וכו'. במסמכים שבתוך כל קבוצה ההפחתה שווה. כתוב נוסחה מדויקת ל DCG קבוצתי לדרגה K כלשהי, הנוסחה צריכה לכלול תמיכה ב- K שאינו בכפולות של 5. יש להסביר בקצרה את ההיגיון של הנוסחה.

$$DCG = \sum_{j=0}^{\lfloor \frac{K}{5} \rfloor} \sum_{i=1}^5 \frac{rel(r_{5+j \cdot 5 + i})}{\log_2(j+2)} + \sum_{i=1}^{K-y} \frac{rel(r_{5+j \cdot 5 + i})}{\log_2(y+3)}$$

קבוצת מסמכים  $y = \frac{K}{5}$

דרגה  $y = \lfloor \frac{K}{5} \rfloor$

8. 42% במאגר של 25000 מסמכים צירפו לכל מסמך רשימה מדורגת (לא ריקה) הכוללת 1-10 מילים לכל מסמך של מילות מפתח מתוך מאגר של 100 מילות מפתח הרלוונטיות למאגר (המילים לא בהכרח מופיעות במסמך אליו הן מצורפות). סדר המילים המצורפות למסמך מעידה על חשיבותן למסמך (המילה הראשונה חשובה ביותר, השנייה פחות מהראשונה, השלישית פחות מהשנייה וכו'). שאילתא במנוע שיועד למאגר היא רשימה של מילות מפתח מתוך 100 המילים (ללא חשיבות לסדר המילים בשאילתא), אין משקל למילים בשאילתא.
- פונקציית הדירוג של מסמך לשאילתא צריכה :
- עבור שאילתא של מילה אחת להחזיר מסמכים שברשימה של מילות המפתח שלהם מופיעה מילת המפתח של השאילתא ולדרג את המסמכים המוחזרים לפי המיקום של המילה ברשימת מילות המפתח של המסמך.
- עבור שאילתא של יותר ממילה אחת, פונקציית הדירוג צריכה להתחשב במספר המילים המשותפות בין מילות המפתח של המסמך ומילות השאילתא וכן בסדר של מילות המפתח במסמך.
- (לדוגמא, מילות המפתח של מסמך A כוללות שתי מילים של השאילתא- מילה אחת של השאילתא היא ראשונה במילות המפתח של המסמך, והמילה השנייה של השאילתא היא השנייה במילות המפתח של מסמך A,
- מילות המפתח של מסמך B גם כוללות שתי מילים מהשאילתא- מילה אחת של השאילתא נמצאת במקום השלישי במילות המפתח של המסמך, והשנייה במקום ה-5 במסמך. A ידורג יותר גבוה מ B.)
- א. 10% הצע מודל ווקטורי (עם 100 ממדים) שיתמוך בפונקציית הדירוג הנ"ל: תאר מה יכלול ווקטור של מסמך. תאר בקיצור כיצד נקבע המשקל של מילת מפתח בווקטור של מסמך (כלומר, מהי פונקציית המשקול של מילה במסמך). אין צורך לחשוב על מודל חסכני בזכרון.
- ב. 12% כתוב פונקציות דרוג מדוייקות שיבטאו את דרישת הדרוג לעיל. פונקציה אחת לשאילתאות של מילה אחת, ופונקציה שניה לשאילתאות עם יותר ממילה 1.
- ג. 14% כיצד ניתן לשלב את המודל שהצעת ב-א עם מודל ווקטורי סטנדרטי – כלומר מסמך ייוצג גם על ידי מילות מפתח מהרשימה וגם על ידי המילים המופיעות במסמך ושאילתא תוכל לכלול גם מילים מהמסמך וגם מילות מפתח.
- 6% כיצד ייראה ווקטור של מסמך במודל המשולב
- 8% כיצד יחושב הדמיון בין שאילתא ומסמך (נוסחא מדוייקת)
- ד. 6% ציין שיפור אפשרי של ביצועי מנוע המבוסס על המודל המוצע (כלומר, מדוע כדאי להשתמש במודל כזה, מה הוא יכול לשפר).

בהצלחה - ברכה ואורלי

$$P_{5,0} = 5 + \frac{(4-6) \cdot 1 + (3-6) \cdot 0.9}{1+0.9} = 5 - 2.473 = 2.526$$

\* הנסחה גם בתורה ברוך  
\* יתר הטלס

$$y = \frac{k}{L5}$$

2. 3. 2020

$$DCG = \sum_{j=0}^{K-1} \left[ \sum_{i=1}^j \frac{rel(s_{i,j}+i)}{\log_2(j+1)} + \sum_{i=j+1}^{K-1} \frac{rel(g_{i,j}+i)}{\log_2(j+2)} \right]$$

בשם ה' אלהינו  $\int_0^1 \frac{x^k}{1-x} dx = 1$  כי אם  $\int_0^1 \frac{x^k}{1-x} dx = 0$  אזי  $\int_0^1 \frac{x^k}{1-x} dx = 1$  כי אם  $\int_0^1 \frac{x^k}{1-x} dx = 0$  אזי  $\int_0^1 \frac{x^k}{1-x} dx = 1$

כ' חזק (היטל) ה' רע' י' ז' ה' ס' א' ; ז' ב' ה' א' ח' כ' א' ח' א'

$\log_2(2) = 1$  - א' (אחד)

[illegible]

נסים מאת הרב ורקה DCG אבוי K נקוניץ עקבא

.5 fe



שורה 8

	Term 1	Term 2	...	Term 100
Doc 1	$\frac{1}{doc(1,1)}$			
Doc 2				
...				
Doc 25,000				

(א)

אנחנו נחשב את המטריצה  $W$  בגודל 100 מילים בלבד. כלומר  $W_{ij} = \frac{1}{doc(i,j)}$

$$W_{ij} = \frac{1}{doc(i,j)}$$

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

✓ במילים אחרות,  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ . כלומר  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

(ב) נניח שיש לנו מטריצה  $W$  בגודל 100 מילים בלבד. כלומר  $W_{ij} = \frac{1}{doc(i,j)}$

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

$$Q = (x, y, \dots, z)$$

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

(2)

הערות:  $Term$  שבו  $i$  הוא המילה,  $j$  הוא המספר.  $doc(i,j)$  הוא מספר המילים  $i$  ב- $j$ .

② נדמה שיש לנו מודל וקטורי כזה  $\vec{v}$  מתוך מרחב נוסף  
 המילים מפתח בנוסף למילים המילות (אם למחנך ואם לשאלות).

$$\text{doc}_i = (x, y, \dots, z)$$

בגורם האחד הוא מספר המילים במאמר הכתוב אחריו את המילים  
 מפתח (לא נרצה כמילות מילים). בדרך אחרת יחושך בקי.

$$w_{ij} = \begin{cases} d(t_i \times idf) + (1-d) \cdot k_{ij} & \text{אם מילים במאמר מילות מפתח} \\ (t_i \times idf) & \text{אם מילה לא הופיעה במילות מפתח} \end{cases}$$

$k_{ij}$  - צורה המילה  
 מילות הקולט  
 $d$  - שדה המילה  
 שונה מפתח למילה  
 מפתח ומילים  
 במאמר

אם בהמשך נרצה למצוא מילה שבו נשנה רק מילות מפתח נאמר  $d=0$   
 ואז  $w_{ij}$  וקטורי כזה נאמר  $d=1$ . וזה לא אחרי נותן לי אפשרות  
 מספר את שילוב 2 המונחים. בדרך נרצה להמשיך יהיה מומלץ.

מבין שיהיה התחשבו במילות מפתח במסמכים שני מוצגים את עת  
 היה זה מפתח בשאלות וכן המילתה יראה כמו במיל וקטורי

כזה  $\vec{q} = (x, y, \dots, z)$   
 כאשר אורכו הוא במספר ה-terms במאמר כולו.  
 ואין חשש ה similarity יתגש שוב בעצרת  $\cos \theta$   
 ומהמילים?  
 המילים המילים?  
 אחריות?

$$\cos \theta = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| \cdot |\vec{q}|}$$

יפה!  
 מאלו!

③ Web מודל מפתח - keywords (בז' 15-10) מילים

בדרך אחרת למעט חיפוש דומה מה הנשא של הדף. זה  
 מילוי המילוי של הדף מילוי זה הנשא של הדף. זה  
 מילוי זה מילוי מילים שמחזקות את מילוי זה וכן נרצה עת  
 key-words שילובי זה את הדף שבו למעט.

אחת כזו המילה כמילה זה בעל המילים אג המילה  
 המילה. בלוי במילים המילים כמו Web במילי מילוי

כאן מופיעים כללים אחרים. ייתכן שיש להם חשיבות רבה יותר מאשר  
האחרים. (אם כן).

(2-)  
הייתה יפה  
אך קצת חזקה  
הרעיון.  
לה כן שם אם  
תוכלו לזה סבסטיאנוס  
על סבסטיאנוס  
על סבסטיאנוס  
היה סבסטיאנוס  
לדא וסבסטיאנוס  
גם הסבסטיאנוס  
תקני הוא התיאור concept

אברהם יצחק