

הוראות לנבחן בצידו השני של הדף

אין לכתוב מעבר לשוליים משני צידי הדף

שאלה 1 - 15
 ✓ שאלה 8
 ✓ שאלה 3

חלק 2

62/2 - 9

מס' כיתה 16 בנין 34

מס' נבחן 023093

מדבקה

1

מדבקה

המחלקה _____ שנה _____

תאריך בחינה _____

מקצוע בחינה _____

לשימוש המרצה הבודק

יחידות / עשרות / מאות

<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	0
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	1
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	2
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	3
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	4
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	5
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	6
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	7
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	8
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	9

ציון הבחינה 76

שם המרצה _____

חתימה נצה שפירא

תאריך _____

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).

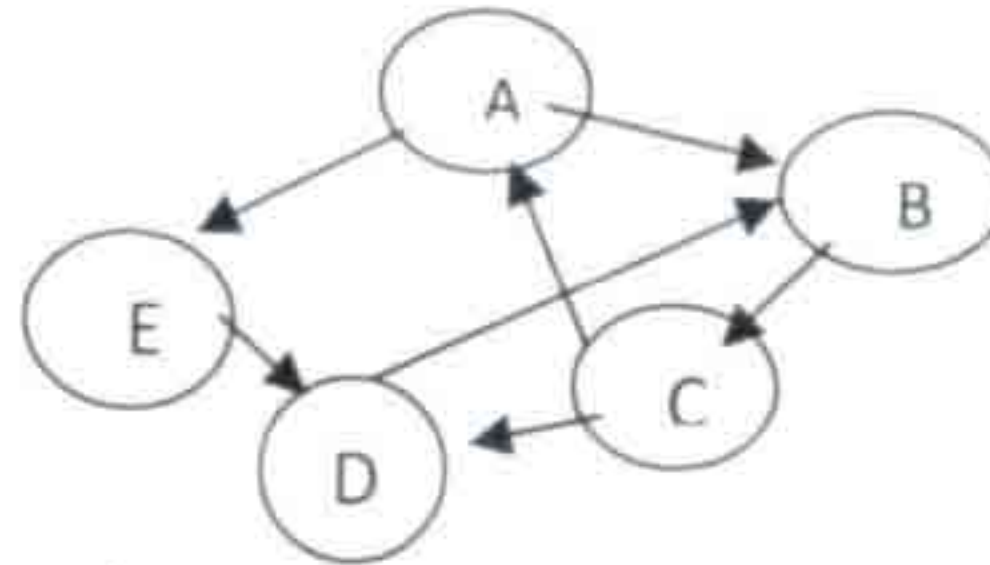
3. 16% למאגר מסוים ואוסף שאליות התבצעה הערכה של מומחים על הרלוונטיות של כל מסמך לשאליתא (לצורכי הערכה של מנועים). ההערכה כללה 3 סוגי תוצאות (במקום ה 2 הסטנדרטיות) : 0 – כאשר המסמך לא רלוונטי, 1 כאשר מסמך רלוונטי חלקית ו 2- כאשר המסמך רלוונטי לגמרי.

- 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 הקטנה)

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

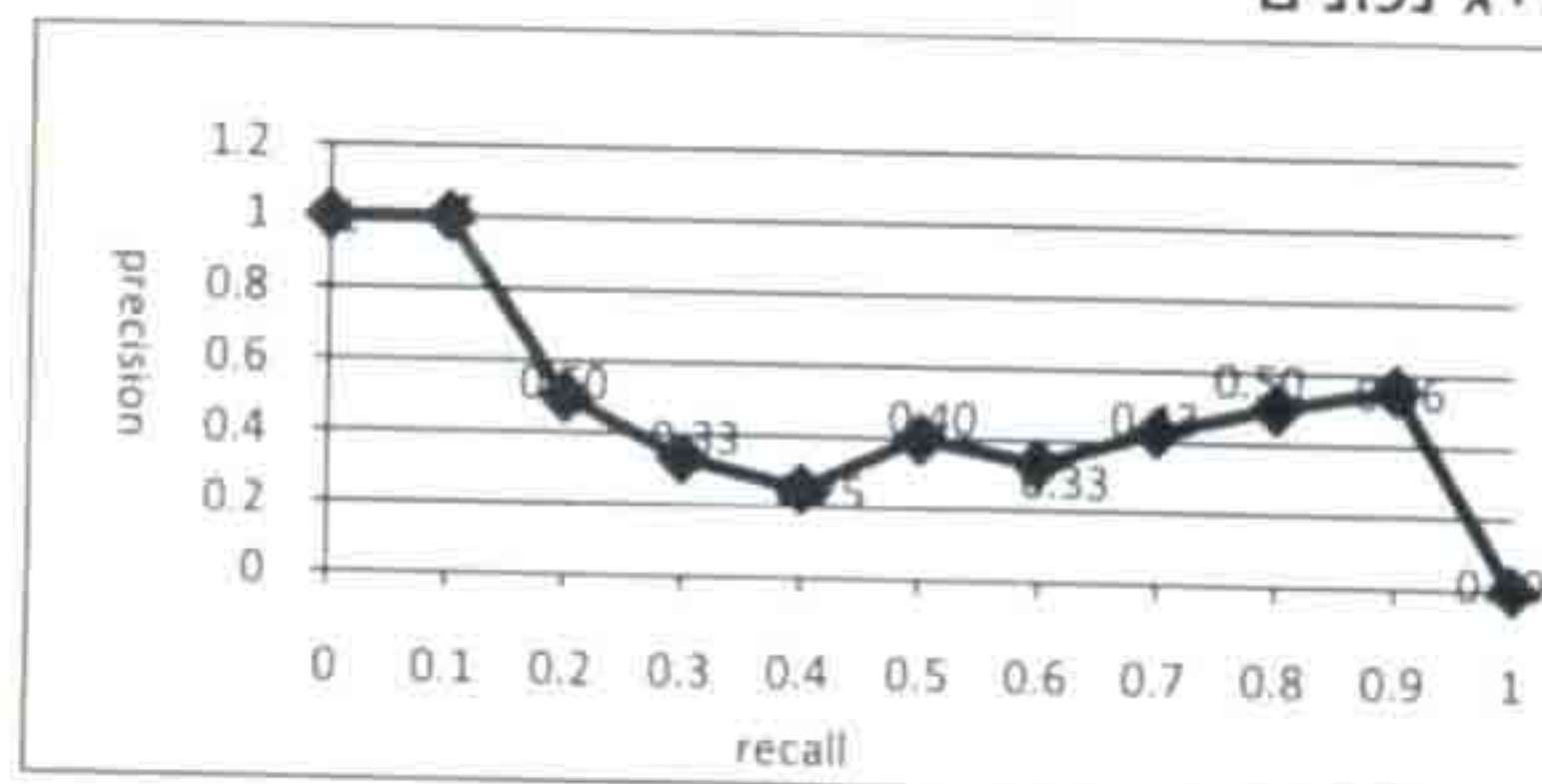
12% סמן תשובה אחת נכונה לשאלות הבאות

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

(4) תחילה נחשב את הצומת בעל ערך ה-hub העליון ביותר לפי HITS:

$$H(A) = A(B) + A(E)$$

$$H(B) = A(C)$$

$$H(C) = A(A) + A(C)$$

$$H(D) = A(B)$$

$$H(E) = A(D)$$

צמתים A ו-C מצביעים על יותר צמתים מאשר שאר הצמתים האחרים ולכן הם המועמדים להיות ה-HUB. נבדוק רק בשניהם ונראה שהבין מי יותר גדול:

$$(*) H(A) = A(B) + A(E) = 2H(A) + H(D)$$

$$A(B) = H(A) + H(D) \quad A(E) = H(A)$$

$$H(C) = A(A) + A(C) = H(C) + H(B)$$

$$A(A) = H(C) \quad A(C) = H(B)$$

ניתן לראות כי A מצביעה על יותר צמתים והנוסף יותר צמתים מצביעים אליה ולכן זוהי הצומת בעל ערך ה-hub העליון ביותר.

מבין הצמתים E, B, D - הצומת בעל Authority הנמוך ביותר היא:

$$A(B) = H(A) + H(D)$$

$$A(D) = H(E) + H(C)$$

$$A(E) = H(A)$$

ניתן לראות באופן מיידי כי $A(E) < A(B)$ ולכן רק נשאר לבדוק מה היחס בין $A(E)$ ל- $A(D)$.

$$(*) A(E) = H(A) = A(B) + A(E)$$

$$(*) A(D) = H(E) + H(C) = A(D) + A(A) + A(C)$$

לכן ניתן לראות שצומת E בעל Authority הנמוך יותר כי היא מצביעה על פחות צמתים וזו פחות צמתים מצביעים אליה.

2 השמות בלבד - PageRank ה' עמוד ביותר הם:

$$PR(A) = \frac{PR(C)}{2}$$

$$PR(B) = \frac{PR(A)}{2} + PR(D)$$

$$PR(C) = PR(B)$$

$$PR(D) = \frac{PR(C)}{2} + PR(E)$$

$$PR(E) = PR(A)$$

ניתן לראות כי צומח B מקבל את כל ה-PR מצומח D וכל צומח

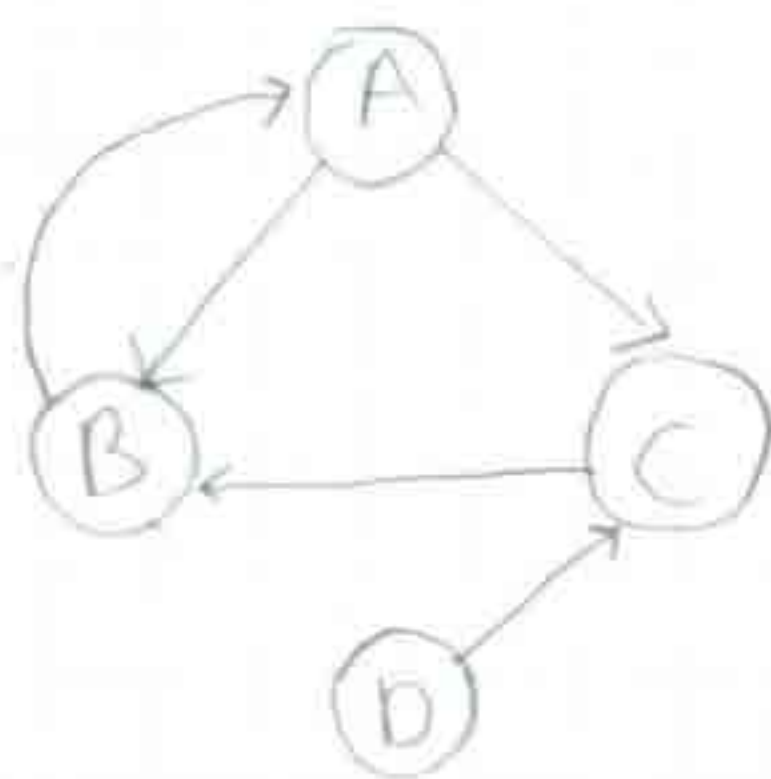
מה-PR של A וכל צומח D מקבל את כל ה-PR(E) וכל צומח

מה-PR(C), כל צומחים האלו מקבלים את ה-PR ה' עמוד בלבד

שאלה 5

א. נכון, דרך זה יכול להיות שצדד $N-1$ וזה תלוי אם מחשבים אנחנו את tf/idf או את שניהם יחד. אם בלבדית הבית התבקשנו לחשב idf מבלי לחשב את tf - \log ולכן קיבלנו דרך idf לצדד $N-1$ מה שהביא אותנו לתוצאה של w_{ij} שצדד $N-1$. אם זה תלוי מימוש.

ב. נכון, שימוש ב-Stemmer בדירג מבוא יותר תוצאות ראיונותיות כי עם מילה בשפה יש הסתם דבר ואם נדע להתייחס רק לשורש המילה אז נוכל לקבל הרבה יותר תוצאות מתאימות. כחוקן שכאשר ישנן יותר תוצאות ראיונותיות אז גם שכיחות הדיבר של המשתמש מהמנוע אלה. יש לבין שבמקרים מסוימים Stemmer יכול גם לעזור מילה מסוימת ובתוצאה מכך תתקבל תוצאה לא ראיונותית אבל ברוב המקרים Stemmer דווקא מזהר על שפה את הראיונותיות.



$$P(A) = P(B)$$

$$P(B) = \frac{P(A)}{2} + P(C)$$

$$P(C) = P(A) + P(D)$$

$$P(D) =$$

ג. לא נכון.

3. נניח שהמסמך d_1 שניהם הופיעו - teach פה d_1 - taught פה d_1

$$tf(teach) = \frac{2}{4}$$

$$tf(taught) = \frac{1}{4}$$

$$idf(teach) = \log_2 \frac{5}{2}$$

$$idf(taught) = \log_2 \frac{5}{3}$$

$$\left. \begin{aligned} d_1 \rightarrow teach \quad tf \cdot idf &= 0.66 \\ taught \quad tf \cdot idf &= 0.18 \end{aligned} \right\} \begin{aligned} &0.84 \\ &idf \end{aligned}$$

נניח שסך אותה המסמך הוא 4

הנרמול הוא לפי אותה המסמך

אם המילה taught ו-teach

היו באותו מסמך אז צריך

לספור אותה לאחד השני

הם לא אותה ספור ה-idf

ולעומת זאת לפני השני

צריך לספור פה

$$tf = \frac{3}{4}, idf = \log_2 \frac{5}{3} \Rightarrow tf \cdot idf = 0.55$$

אחרי
השני

סל"כ הבלתי עתידית מקרה של חיבור התוצאות שנה ועם התשובה היא לא נטן

$$\vec{q}_m = (0, 0, 0.8, 0.8, 0.2, 0.3)$$

(ה)

שאלה 6

E_1 ו- E_3 מקדים לפי מודל בוליאני טייר ומודל בוליאני מורחק בהתאמה.

E_2 ו- E_1 Malaga and Tibet נשלח לפי E_2

Malaga Tibet נשלח לפי E_2

(א) מנוע E_1 יחסי אל D_1, D_2, D_4

נבדוק מה יחסי מנוע E_2 לפי "Malaga Tibet":

נניח שהמודל הקטור מניח אל E_2 לפי אורך המסמך:

$$\text{Sim}(D_1, q) = \frac{2}{3} \cdot 1 + \frac{1}{3} \cdot 1 = 1$$

$$\text{Sim}(D_2, q) = \frac{1}{7} \cdot 1 + \frac{4}{7} \cdot 1 = \frac{5}{7} = 0.71$$

$$\text{Sim}(D_4, q) = \frac{1}{4} \cdot 1 + \frac{1}{4} \cdot 1 = 0.5$$

מנוע E_2 יחסי D_1, D_2, D_4

(ב) E_1 ו- E_3 הם מנועים שמקדים לפי מודל בוליאני ולפי קטלג החזרה

מסמכים מתאימים הם יחשבו בהתאמה מסאה אל המסמכים שיש בהם גם

Malaga ו-Tibet ← הם יחשבו אל אותם מסמכים בדיוק

אבל לא בטוח שלה יהיה קאולי סדרי כי לה תלוי בדיוק שיקדם המודל

הבוליאני המורחק.

שאלה 7

לא תמיד ה- precision יתפר בעקבות הפיסקה שיקדם מהמסמך כי

לה תלוי גם במנוע עצמו ובאופן שבו הוא אינטרנט אל המסמכים. בבוטנציאל הרצון

של הפיסקה הוא אלפר אל ה- recall וה- precision אך זה לא תמיד

קורה

שאלה 8

$$\text{MAP}_1 = \frac{\frac{1}{1} + \frac{2}{3} + \frac{3}{4} + \frac{4}{10}}{4} = 0.6$$

(מסמך אל המודלים:

$$\text{MAP}_2 = \frac{\frac{1}{2} + \frac{2}{5} + \frac{3}{6} + \frac{4}{7}}{4} = 0.49$$

$$R\text{-precision}(E_1) = \frac{2}{4} = 0.5$$

$$R\text{-Precision}(E_2) = \frac{1}{4} = 0.25$$

תשובה א+ג (מונה).

טבלה 9

	recall	Precision	הגדלה היה חייב להתקדם אחרי תהליך אינטרפולציה ואבסור לחיות לפי הטקסט ההתחלתי שנתתי.
1	$\frac{1}{10} = 0.1$	$\frac{1}{1} = 1$	
2			
3			
4	$\frac{2}{10} = 0.2$	$\frac{2}{4} = 0.5$	אם תוצאות הגדלה רואים לפחות 90% מהתוצאות חדש.
5			
6			
7			
8			
9	$\frac{3}{10} = 0.3$	$\frac{3}{9} = 0.33$	
10			
11			
12			
13			
14			
15			
16	$\frac{4}{10}$	$\frac{4}{16} = 0.25$	
17			
18			
19			
20			

הסיבות שבגללן צריך - Crawler שמוציא את כל הדפים

מאתר מסוים:

(1) גלישת URLs - את ה-Crawler ישמוציא את כל הדפים של אתר מסוים בצומת אחד אלא יהיה לו קוד יותר מתחילת של הוצאת הדפים מהדפים כי במקרה ויש לך מיליון אתרים באתר עצמו אתר האתר (דבר שהוא מאוחר נפוץ באתרי web) אלא אם הוציא האתר יהיה מחזיק בצומת אחד ייתכן יותר זמן - Crawler עצמו צריך הוצאה מה מנין סיבות: תלבושת רשת, ממשל על הצומת וכו'. ברוב המוחלט של המקרים, ישנו קשר חזק בין דפים ששייכים לאותו אתר אינטרנט וכמות הדפים שיהיו בתורם לדפים אחרים באותו אתר כנראה שיהיה גבוה.

(2) ציור מהוראות הדפים (אתיקה) - בעלי אתרים נוהגים לשים קובץ robots.txt או עדיין ב-Metadata של דף מסוים האם מותר לאנדרקס

דף מסוים ואיזה חלקים מתוכו בעצם מותר לאנדרקס, אם נכיר את האתר על פני כמה צמתים אלא יהיה לנו קשה לנהל את הציור מהוראות הדפים כי גם צומת נכסירק לעדיין את התיקיות שאסור עלינו לאנדרקס, התחזוק של הנושא הזה יהיה הרבה יותר קשה מושיא אם הם יהיה על צומת אחד ונדע בדיוק מה מותר ומה אסור לאנדרקס גם האתר בעצמו.

באופן כללי, התחזוקה של ביטוי אתר אחד על פני כמה אתרים הוא קשה במיוחד, צריך לעצור להתנהל במקרה של הזמן ואפילו יש בעיה שלא תלויים ב-Crawler שלא יאפשרו לו לעשות את זה.

לפעמים יטעו להיות מצב, שבו צומת אחד נפגע ואיחסן של חלק מהאתר וכלל לא נואל להשיק בפעולה ה-Crawler כי הוא תלוי בדף הזה שמתחיל באתר שנופל.

שאלה 3

(א) $Precision = \frac{\text{מס' ראיונות שחזרו}}{\text{מס' המסמכים שחזרו}}$ עדי היום חישובנו

מכיוון שכבר צייק לקחת בחשבון גם את הסקאלה של הדרכת מסמכים
אז נתקסם על אותה נוסחה שאנחנו מכירים רק שכבר במקרה נתחכם למס' הראיונות
שחזרו באופן קצת שונה:

במצב (רמזתי את ההערכה

$$\left\{ \begin{array}{l} (1) \text{ אם ההערכה היא 0 אז מס' ראיונות שחזרו} = 0 \\ (2) \text{ אם ההערכה היא 1 אז מס' הראיונות} = 0.5 \\ (3) \text{ אם ההערכה היא 2 אז מס' הראיונות} = 1 \end{array} \right.$$

עפי הציון העקוב ביותר.

$$Precision @ K = \frac{\sum (\text{ראיונות עד ענק } K)}{K}$$

ראיונות עד ענק K ייספו באופן שלביותי למעלה.

הרציון הוא שאם יש מסמך שמתאים רק באופן חלקי אז בלצמ הוא יוצי ראיונותי
ועם ניתנים לו משקל נמוך יותר בספירה.

(ב) נציג את האלגוריתם של $Precision @ 5$.

תוצאות השאלות (עד $K=5$)

1, 0, 2, 0, 0

אחרי נרמול המשקל

0.5, 0, 1, 0, 0

$$Precision @ 5 = \frac{0.5 + 1}{5} = 0.3$$

