

הוראות לנבחן בצידו השני של הדף

אין לכתוב מעבר לשוליים משני צידי הדף

חלק א' -  $\frac{8}{30}$  - 22

חלק ב' -  
שאלה 1 -  $\frac{20}{20}$   
שאלה 2 -  $\frac{20}{20}$

מס' כיתה 226 בנין 90

מס' נבחן 0085564

שאלה 15  
שאלה 15

25/03/2012 1352180 1 9/9

הנדסת מערכות מידע

אחזור מידע וספריות דיגיטליות

03721440601402



לשימוש המרצה הבודק

יחידות ועשרות ומאות

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

ציון הבחינה 78

שם המרצה ד"ר שלמה

חתימה

תאריך

המחלקה הנ"ל מא' ג' שנה 3

תאריך בחינה 25/3/12

מקצוע בחינה אקטורי גידול וסמכיות ביטליות





אוניברסיטת בן-גוריון בנגב

הוראות לנבחן

1. בהגיעך למקומך יש להניח את כרטיס הנבחן ותעודה מזהה על שולחןך.
2. אסור להביא למקום הבחינה תיקים, ספרים, מחברות, טלפון נייד או רשימות פרט למותר על פי שאלון הבחינה.
3. עזב תלמיד את האולם אחרי חלוקת השאלונים, דינו כדין "נבחן" בבחינה.
4. אסור לנדרו לשוחח בזמן הבחינה. או לעזוב את מקומו ללא נטילת רשות.



איחזור מידע תשע"א – 372.1.4406

סמסטר חורף מועד ב' 23.03.2012

פרופ' ברכה שפירא, אורלי מורנו

משך המבחן : שעתיים וחצי

חומר עזר: מותר (לא מחשב נייד), מותר מחשבון

יש להחזיר את גיליון הבחינה. המבחן כולל 4 דפים. יש לענות על כל השאלות

חלק א 30% - יש לענות על הגליון

בחר תשובה אחת נכונה

1. 4% הפרמטרים בנוסחת rocchio ל relevance feedback :

- א. קובעים האם המשתמש סימן את המסמכים כרלוונטים או לא.
- ב. נקבעים דינמית לפי תגובות המשתמש למסמכים
- ג. צריכים להסתכם ל - 1
- ד. קובעים את רמת ההתחשבות בשאילתא הנוכחית, במסמכים הרלוונטים ובמסמכים הלא רלוונטים בעדכון השאילתא
- ה. ג+ד נכונים
- ו. ב+ג נכונים

2. 6% נתונה טבלת הדירוגים הבאה של משתמשים על מוצרים בדירוג 1-5:

	User 1	User 2	User 3	User 4	User 5
Item1	1			1	
Item2		2	3		3
Item3	1	1	5	1	
Item4				4	
Item5	5		1		5

כדי לנבא דירוג של item5 עבור User4 על פי שיטת הדמיון בין המשתמשים:

- א. יש לחשב דמיון בין user4 לכל אחד מהמשתמשים האחרים ולחשב את הדירוג על פי המשתמשים הדומים ביותר
- ב. יש לחשב דמיון בין user4 ל user3 ו user1 ולהשתמש בדירוגים של שניהם ל item5 כדי לחשב את הדירוג של item5 ל user4
- ג. יש לחשב דמיון בין user4 ל user3 ו user1 ולהשתמש רק בדירוג של user1 ל item5 כדי לחשב את הדירוג של item5 ל user4
- ד. יש לחשב דמיון בין user4 ל user3 ו user1 ולהשתמש רק בדירוג של user3 ל item5 כדי לחשב את הדירוג של item5 ל user4
- ה. על סמך הנתונים הקיימים אי אפשר לנבא דירוג של item5 עבור User4
- ו. יש לחשב דמיון בין user4 ל user 5 ובין user 1 ו user 3

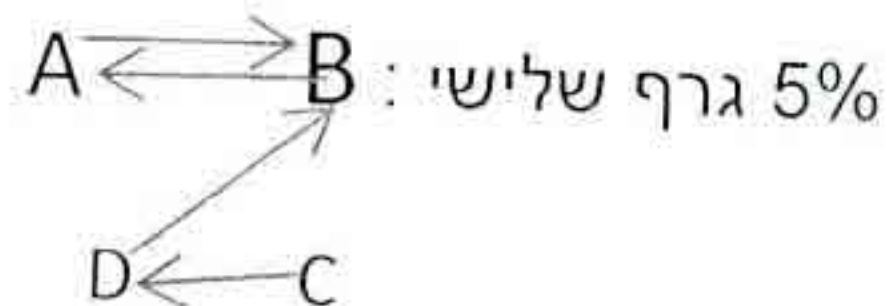
1894



3. 5% המשקל של term בווקטור של מסמך במודל הווקטורי
- מבטא את מידת הייצוגיות של ה term למסמך
  - צריך להיות מחושב באותו אופן שבו מחשבים את משקל ה Terms בשאילתא
  - ☒ מחושב תמיד על פי  $tf*idf$
  - חייב להיות בין 0-1
  - ב+ג נכונים
  - ג+ד נכונים
  - א+ד נכונים
4. Benchmark 4%
- מגדיר מדדי ביצוע
  - משמש להשוואה של ביצועי אלגוריתמים שונים מול מדדים חדשים שהוגדרו
  - משמש להשוואה של ביצועי אלגוריתמים שונים מול תוצאות נכונות ידועות
  - א+ג נכונים
  - מחזק את הקשר בין האקדמיה לתעשייה
  - ☒ א+ג+ה נכונים
5. ענה נכון או לא נכון
- 3% - ההחלטה אם להתחשב ב Robots.txt תלויה ברמת האתיקה של אלגוריתם ה Crawling
- ☒ נכון/לא נכון
6. 4% DCG הוא מדד לבדיקת יכולת של מנוע לדרג גבוה דפים רלוונטים, תוך התחשבות ברמת הרלוונטיות של הדף ☒ נכון/לא נכון
7. 4% stemming הוא תהליך התלוי תמיד בכללי שפה ולכן יש צורך להגדיר לכל שפה stemmer נפרד.
- ☒ נכון/לא נכון

## חלק ב 70%

1. link analysis 20%
- 11% נתונים 3 גרפים שונים המתארים קישוריות של צמתים ברשת:
- 4% גרף אחד:  $A \rightarrow B \rightarrow C \rightarrow D$
- 5% גרף שני: כמו הגרף הראשון, רק שבנוסף D מצביע על A (כלומר נוצר מעגל)
- 5% גרף שלישי:

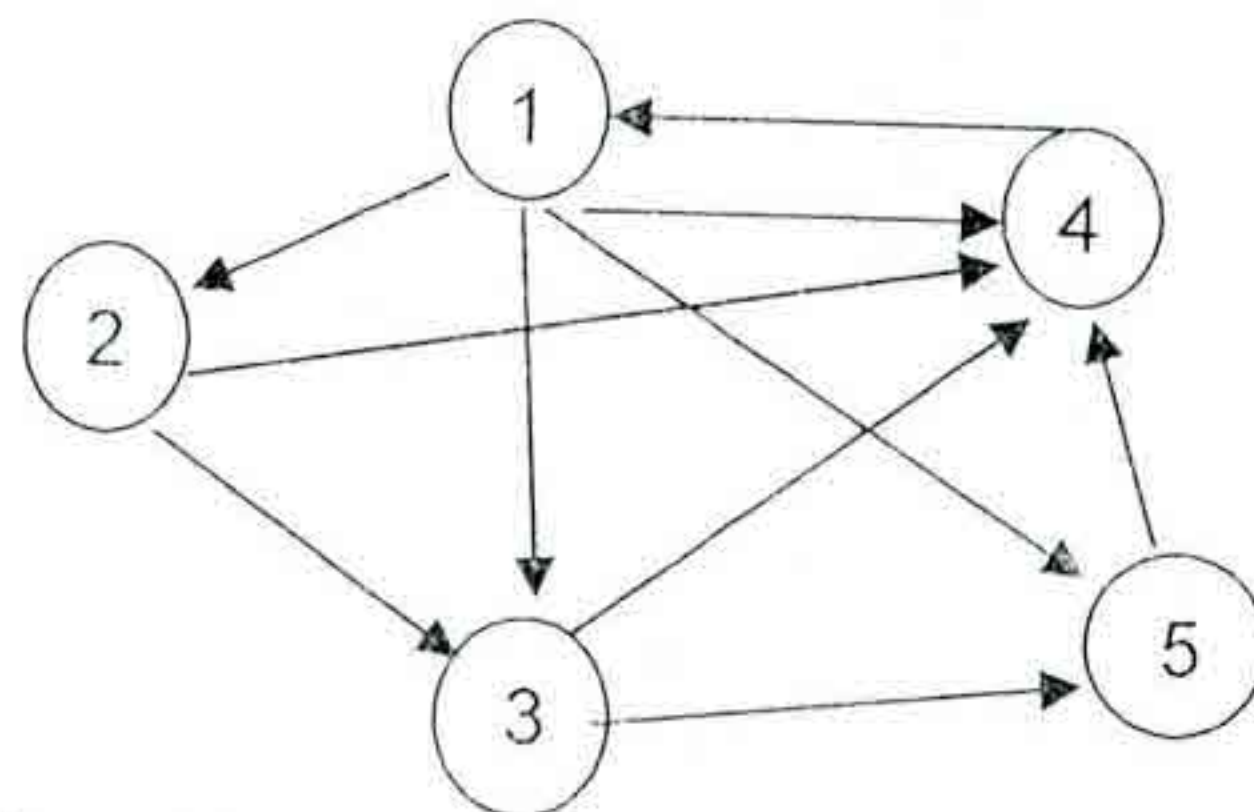


עבור כל אחד מהגרפים הסבר מה יהיו ערכי pagerank של כל הצמתים (אין לחשב ערכים מדויקים, אלא להעריך ולהסביר את היחסים בין ערכי כל הצמתים)



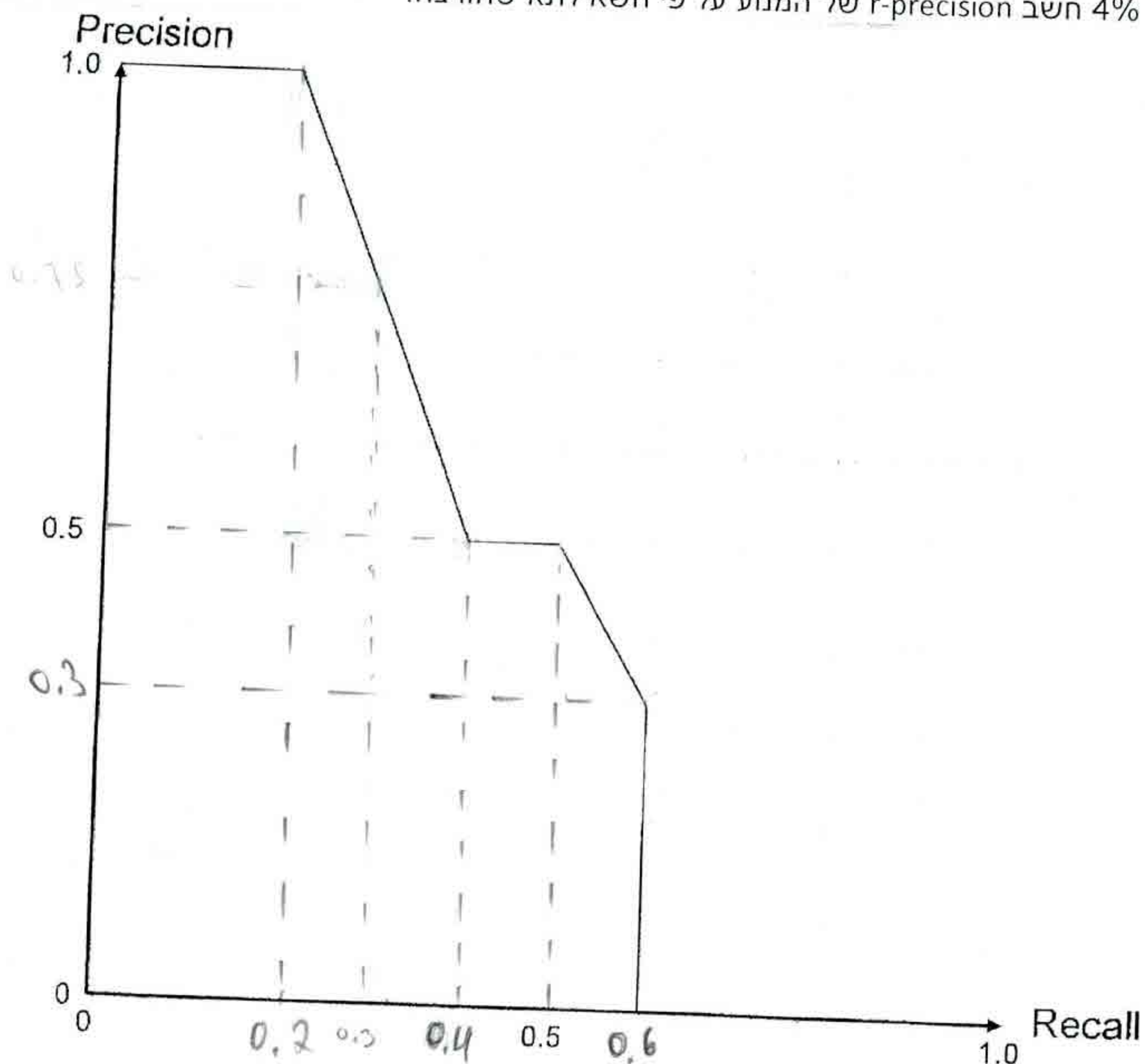


ב. 6% נתון הגרף הבא:



הראה את מטריצת המעברים (transition matrix) שעל פיה יחושב Pagerank לדפים (ללא התחשבות בפרמטר d).

2. 20% נתון גרף Precision-recall על תוצאות של הרצת מנוע מסויים על שאילתא מסויימת. המנוע החזיר 20 תוצאות. ידוע שבמאגר 10 מסמכים רלוונטים לשאילתא.
  - א. 4% מהו ה precision לאחר שהמנוע החזיר 4 מסמכים רלוונטים.
  - ב. 8% הראה את הרשימה שהמנוע החזיר ( הצג רשימה של 20 מסמכים מדורגים וסמן מי מהם רלוונטי ומי לא)
  - ג. 4% חשב את ה map של המנוע על פי השאילתא שהורצה.
  - ד. 4% חשב r-precision של המנוע על פי השאילתא שהורצה.







3. 15% א. 10% הועלה רעיון להשתמש ב jaccard similarity לצורך דירוג מסמכים לשאילתא (כלומר, לדרג את המסמך על פי מספר המילים המשותפות בין שאילתא ומסמך מחולק בסכום של מספר המילים של השאילתא והמסמך). אפשר להניח שכל המילים בשאילתא ייחודיות (אין חזרה על מילים) ושבמסמך סופרים כל מילה רק פעם אחת. באופן פורמאלי הדרוג למסמך D נקבע על פי:
- $$jaccard(Q,D) = \frac{|Q \cap D|}{|Q \cup D|}$$
- כאשר Q הוא מספר המילים הייחודיות בשאילתא, ו D הוא מספר המילים הייחודיות במסמך.
- תן 2 סיבות לכך שהרעיון אינו מוצלח, כלומר שתי סיבות לכך שהשיטה אינה מדרגת טוב מסמכים.
- ב. 5% הדגם את חסרונה של השיטה מסעיף א. כלומר, תן דוגמא ספציפית לשאילתא Q כלשהי ושני מסמכים D1 ו D2. כאשר ברור ש D1 רלוונטי לשאילתא ו D2 לא רלוונטי. בחישוב על פי jaccard יוצא ש D2 יותר רלוונטי מ D1.
4. 15% בניסיון לשפר תוצאות של מנוע שעובד לפי השיטה הווקטורית, הוצע להוסיף לשאילתא מילים נרדפות למילות השאילתא מתוך מילון. לכל מילה בשאילתא הוצע להוסיף עד 10 מילים נרדפות.
- למשל: אם השאילתא היא: אחזור מידע – המילים הנרדפות לאחזור במילון הן: דלייה, אשוב, אשנה. מילים נרדפות למידע הן: ידיעה, אינפורמציה. השאילתא שתישלח למנוע תהיה:
- אחזור, דלייה, אשוב, אשנה, מידע, ידיעה, אינפורמציה
- א. 8% תן סיבה אחת לסיכוי שהשיטה תשפר תוצאות וסיבה אחת לסיכוי שהשיטה תפגע באיכות התוצאות
- ב. 7% הצע רעיון לשיפור השיטה שתגדיל את סיכוייה לשפר את התוצאות, הסבר את ההיגיון ברעיון שאתה מציע.

בהצלחה

ברכה ואורלי





שם / מס' 9/9

אברהם יהונתן

WINTER 2018

12/1/18

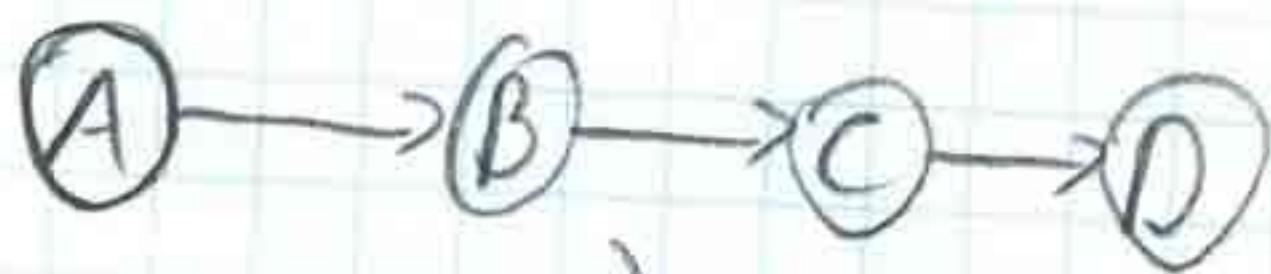
2/2



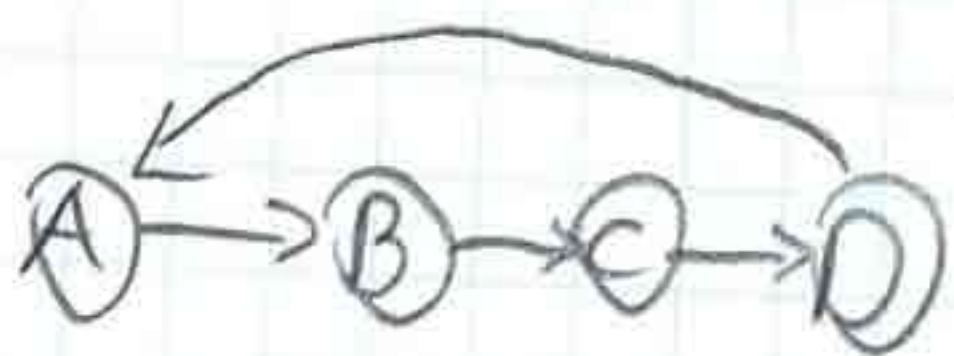
# חלק ב

שאלה 4

א. 0



ה - PR של A שונה מ-PR של D (ההסתברות שנגיע לבית A באופן ישר) דבר שבו מסווג את כל ה-PR של הציורים באופן יחסי.



2

$$PR(A) = PR(D)$$

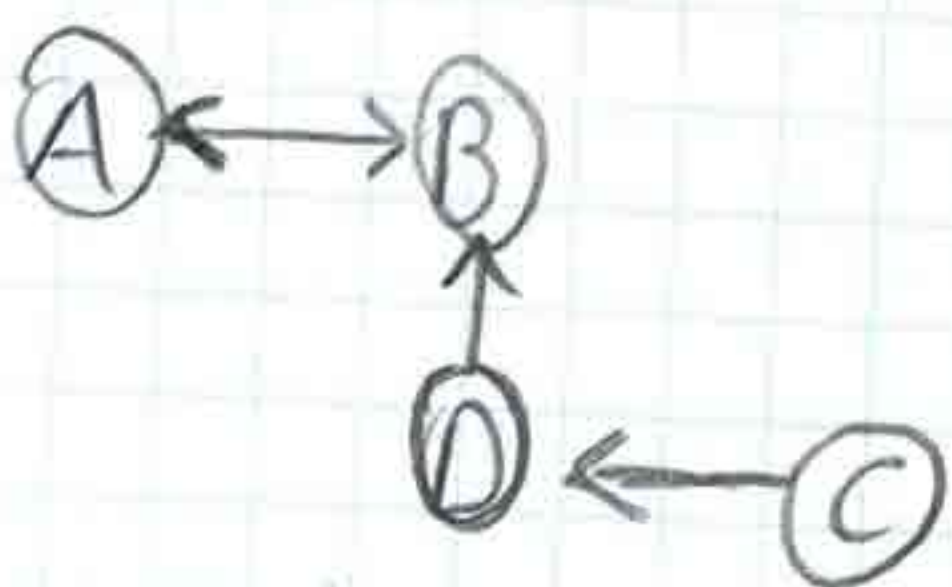
$$PR(B) = PR(A)$$

$$\Rightarrow PR(A) = PR(B) = PR(C) = PR(D)$$

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

$$PR(D) = PR(C)$$

נניח שכל ציור נצפה בצורה אחידה ובסוף ציור D נצפה בצורה אחידה. דבר אחר אינו מיוען בין כל ה-PR של הציורים.



3

PR של C שונה מ-PR של D (כמו במסלול A), ה-PR של D שונה מ-PR של C דבר אשר יאפשר לנו לדבר על PR של A ו-B.

$$PR(A) = PR(B)$$

$$PR(B) = PR(A) - PR(D)$$

אם נהיה בטוחים ש-PR של D שונה מ-PR של A ו-B, נאמר ש-PR של A ו-B שונה מ-PR של D.

בטוחים ש-PR של A ו-B שונה מ-PR של D, ואכן  $PR(A) = PR(B)$ .



பி. ௨

பி. ௧

பி. ௧

பி. ௧

பி. ௧



# המשן של 1

מט כוונת המשברים: (כ)

משנים	1	2	3	4	5
1	0	0	0	1	0
2	$\frac{1}{4}$	0	0	0	0
3	$\frac{1}{4}$	$\frac{1}{2}$	0	0	0
4	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{2}$	0	1
5	$\frac{1}{4}$	0	$\frac{1}{2}$	0	0
	"	"	"	"	"
	1	1	1	1	1



בק'קה



0.001 0.001

③

0  
0.1  
0.2  
0.3  
0.4  
0.5  
0.6



תאריך

## שאלה 2

(תנאים: התורם החזיר 20 תוצאות  
במאגר יש 10 מסמכים רלוונטיים לשאלה)

א.  $Precision = \frac{4}{8} = 0.5$  אחי 4 מסמכים רלוונטיים



א.

1	✓	✓
2	✓	✓
3	x	
4	✓	✓
5	x	
6	x	
7	x	
8	✓	✓
9	x	
10	✓	✓
11	x	
12	x	
13	x	
14	x	
15	x	
16	x	
17	x	
18	✓	✓
19	x	
20	x	

✓ - רלוונטי

ב.

x - לא רלוונטי

ג.  $MAP = \frac{1 + 1 + 0.75 + 0.5 + 0.5 + \frac{1}{3}}{6} = 0.6805$



ד.  $R-precision = \frac{5}{10} = 0.5$



ה.



ХИТЛОУР-ХИТЛОУР

1998

ХИТЛОУР

ХИТЛОУР

ХИТЛОУР



## שאלה 4

א. ס'בה לשיפור - במקרה זה כאשר עשאלתא מוספים מילים  
נכדיות לכל מילה בשאלתא התנוע ותפיר מסמכים  
אשר קשורים לשאלתא אך לא בוקשו בדוכה נפותר ע"י  
המפתח, הכי כאשר מפתח שאלתא עשאלתא אדואמא:  
אחזור מידע. אן זה משנה אם הוא יקבל מסמכים אשר  
כוללים את המילים אחזור אקדוממיה, מפני שאקדוממיה  
ומידע ישנה אותה משמעות, דבר זה מקבל את כמות  
כמות המסמכים הדיוקנים אשר התנוע ינוע להפציר.

ס'בה שהמטרה תפסל באמת - אם המפתח שאלתא עשאלתא  
מסוימת ותו חרזה לקבל מסמכים אשר המילים הספליפית  
עשאלתא מופיעות המסמכים, הנספח המילים תנדפית ציבור  
לפניו במסלף המסמכים וכך המסמכים הדיוקנים עשאלתא  
הספליפית יכל להיות לא יתפרו בכל ע"י התנוע או שיפצרו  
אחרי הדוכה מסמכים לא דיוקנים קבר אשר יפגע באיכות  
התוצאות.

ב. כאשר המפתח שאלתא שאלתא אפס "לבקש" המפתח  
לבתור חס הוא חרזה את המילים הספליפית הוא או הוא  
לחפס הם מילים נכדיות, כאשר הם פתח יבחר את  
התשובה הרצויה לו, לפני שהשאלתא תפסל לתנוע נסע  
הוא לחפס מילים נכדיות במלון ולחפסן עשאלתא או  
לשאלתא את השאלתא כמו שהמפתח שאלתא בהתחלה.  
הדבר יוכל לתפסל על הדיוקנים שבתפסל בסוף עשאלתא  
את התוצאות.



a. prob H  
Dog

3.  
D<sub>1</sub> Dog Dog Dog Cat Eat  
D<sub>2</sub> Dog Cat

$$D_1 q = \frac{1}{2}$$

$$D_2 q = \frac{1}{2}$$





## שאלה 3

א. סיבה 1- כאשר ישנה שאילתה בעלת מילה מסוימת  
אשר נמצאת במסמך אחד מהם של פעמים ובמסמך  
שני מהם מוצא של פעמים אזי מסמך אחד יותר רלוונטי  
✓ אולם אילתה ינהיגה העניין אך יכול להיות שהצב שבו  
המסמך העניין מקבל את אותו היכרות כמו המסמך הראשון  
נולד צירוף צורה יותר

סיבה 2- אורך המסמך מאוד משפיע על צירוף  
המסמך השטח Jaccard, כאשר מסמך הוא ארוך וקצר  
הכבה מיליון שעות ה- QUD יתה צבול וכן יקטן  
צירוף המסמך, זה אם מילה גם מוכה שאילתה חזית ע  
אצנה גם זה של כסמים עדיין ה- QUD נשאר ע צול  
ומקטן את צירוף המסמך. ✓

ב.

Dog : q<sub>1</sub> - שאילתה

D<sub>1</sub> : Dog Dog Dog Dog Cat dog boy dog

D<sub>2</sub> : Dog Cat Cat Cat Cat

$$\text{jaccard}(Q, D_1) = \frac{1}{3} = \frac{1}{3}$$

$$\text{jaccard}(Q, D_2) = \frac{1}{2} = \frac{1}{2}$$

ולכן מסמך D<sub>1</sub> יותר רלוונטי לשאלתה מ-D<sub>2</sub> אך D<sub>2</sub> מקבל  
משקל יותר גבוהה מ-D<sub>1</sub>.



ХИТТОУОРТИ-ХИИЛТИРТ

2.3.2.2

1.

2.

3.





ХИТРОУЌИ ОУЌИ РЕШЕЊИ ЗА РЕШЕЊЕ







Mathematics for the 21st Century







WILSON'S







Mathematics for the 21st Century