

מדור בחינות ומערכת שעות

הנדסת תוכנה

4/02/19 13:30-16:30

מבוא לאחזור מידע מועד ב' ליטבק מרינה

'תשע"ט סמסטר א

חומר עזר – אסור (מלבד דפי הנוסחאות המצורפים לטופס זה)

. השאלון מכיל 11 עמודים (כולל דף נוסחאות וטיוטה).

| ======================================= |
|--|
| |
| <u>חומר עזר : נא סמן במשבצת המתאימה את המתאים</u> |
| יניתן להשתמש בכל מחשבון* ניתן להשתמש בכל מחשבון Casio FX-991EX * V לא ניתן להשתמש במחשבון Casio FX-991EX |
| * לא ניתן להשתמש במחשבון * לא ניתן להשתמש במחשבון |
| · א ניתן להשתמש בחומר עזר *_V |
| * מותר שימוש בדף נוסחאות, כמפורט: |
| * הבחינה בחומר פתוח – מותר להשתמש בכל חומר עזר מודפס או כתוב |
| <u>הערות</u> _V_ יש לענות על כַל השאלות במקומות המיועדים ע"ג טופס השאלון בלבד יש להחזיר את השאלון ביחד עם הכריכה/מחברת. |
| אחר: |
| ··································· |
| .2 |
| .3 |
| |
| השאלון מכיל 11 עמודים (כולל עמוד זה). |
| בהצלחה ! |
| |



| נכי) | 25) | 1 | שאלה |
|------|-----|---|------|

יש לענות לכל השאלות •

| ן באופן ברור את התשובה <u>הנכונה ביותר</u> על גבי שאלון הבחינה | יש לסמ | • |
|---|----------|-------|
| ל יותר מתשובה אחת לאותה שאלה יקבל ציון של <u>אפס</u> | | • |
| נכלל בתחום: Vector Space Model ב-Euclidean distance נכלל בתחום: [0, 1] | | ж. |
| $[0,\infty]$ | .2 | |
| [-∞, 1] | .3 | |
| $[-\infty,\infty]$ | .4 | |
| (4/ | 3/2/1) | תשובה |
| זפקיד של-Name Entity Recognition הוא: | | ב. |
| · · | .1 | |
| | .2 | |
| | .3 | |
| אָרגונים, מיִקומים, וכד' | | |
| להוציא מילות מפתח עבור מסמך | .4 | |
| (4/ | 3/2/1) | תשובה |
| אלגוריתם למידה בשם K-Nearest Neighbors מקבל כקלט מסמכים בצורה של: | | ג. |
| ווקטור (Vector Space Model) ווקטור | | |
| (Bag of words) BOW | | |
| גרף | | |
| טקסט ללא מילות עצירה | .4 | |
| (4/ | 3/2/1) | תשובה |
| :אפשר לאפיין כ Rocchio אפשר לאפיין כ | (5נק') א | .7 |
| | .1 | |
| • | .2 | |
| (nonlinear classifier) אלגוריתם סיווג לא לינארי | | |
| (probabilistic classifier) אלגוריתם סיווג הסתברותי | .4 | |
| (4/ | 3/2/1) | תשובה |
| | | |
| תונה שאילתה L^*L . ניתן לאחזר מסמכים שעונים על השאלתה בעזרת | | .7 |
| permu ושאילתה הבאה: | | |
| \$RL\$ | | |
| \$RL | | |
| LR\$ | | |
| L\$R* | .4 | |
| ·9m3 (A) | 2/2/1 | |

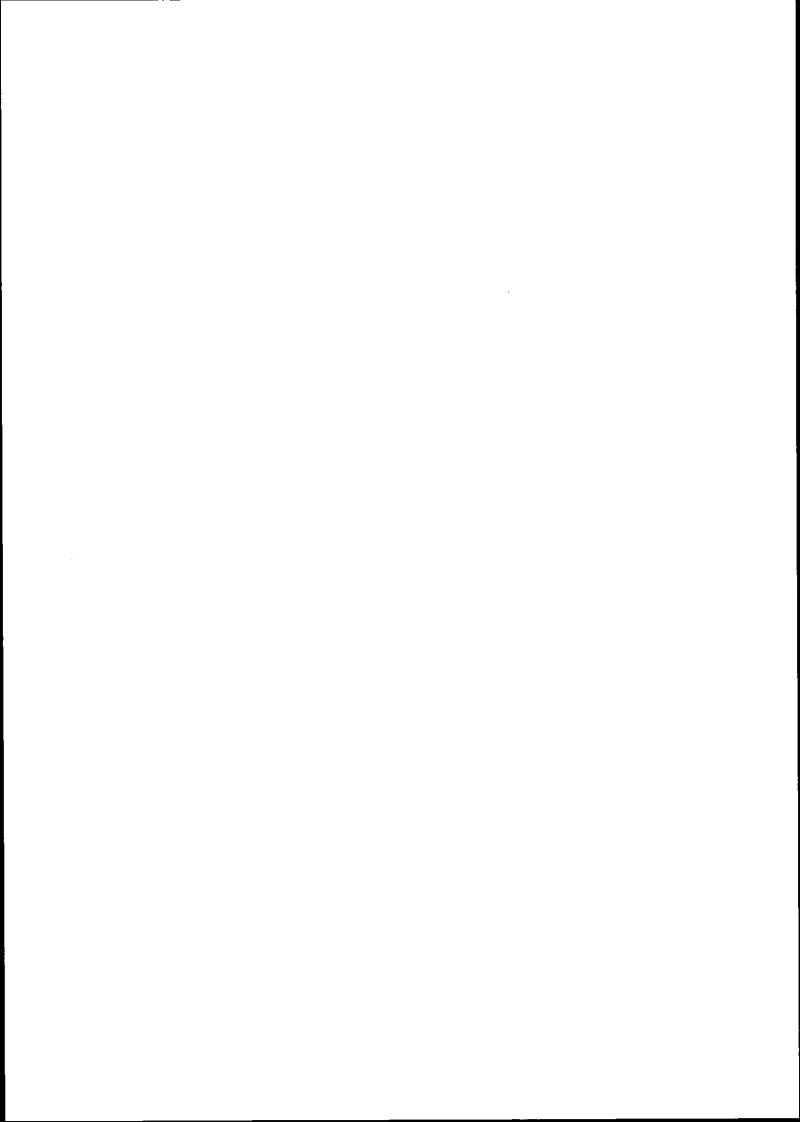


| d h d :D5 | | |
|---|----------------|--------------|
| h c a h c a :D6 | | |
| מסווגים לשתי קטגוריות: P and N באופן הבא: | | |
| | קטגוריה | מסמך |
| | P | D1 |
| | N | D2 |
| | N | D3 |
| - | P | D4 |
| - | N | D5 |
| | P | D6 |
| 1 נק') לא נתונה רשימת מילות מפתח (keywords), אך ישנה דרישו ן לעשות זאת? תציע שיטה ותראה את היישום שלה על המאגר הנתון <u>ל מסמד.</u> | | |
| | | |
| :D1 | | |
| :D2 | | |
| :D3 | | |
| :D4 | | |
| :D5 | | |
| ;D6 | | |
| נק') יש לבצע בחירת מאפיינים (feature selection) בעזרת שיטו | Information הד | Mutual l |
| · | | |
| | | |
| | <u>—</u> | |
| | | - |
| _ | | |
| | | |
| - | | |
| | | |
| | | |
| | | |
| | | |
| | | |

('נק') שאלה 2 (25 נק'

נתונים ששה מסמכים (a-h הם המילים):

a h a b a h c b :D1 h b e b h :D2 d h b h d :D3 h d d a h a e :D4



שאלה 3 (30 נק')

| ון מאגר של ששה מסמכים (ראה שאלה 2). (5 נק') יש לבנות אינדקס הפוך (inverted index) עבור כל המאגר, על בסיס מילים בודדות | נת א. |
|--|-------------------|
| | _ _ |
| | <u> </u> |
| bi-) עבור כל המאגר, על בסיס זוגות מילים (inverted index) עבור כל המאגר, על בסיס זוגות מילים (grams | |
| | <u> </u> |
| | _ _ _ |
| תאר תאליך ? a AND d AND c תאר שאילתה? מסיף א' או ב') תשתמש עבור שאילתה אינדקס (מסיף א' או ב') תשתמש עבור שאילתה אחזור ופלט שלו. האם תשובה תמיד תהיה מדוייקת? | |
| 2 במה מוגבל זמן ריצה של תהליך אחזור עבור שאילתות מסוג a AND b AND c? האם ניתן לעייל את התהליך הזה? אם כן, כיצד? מה אינדקס אמור להכיל עבור זה? | |
| תאר תהליך אחזור ופלט (מסיף א' או ב') תשתמש עבור שאילתה a h c? מאר תהליך אחזור ופלט שלו. האם תשובה תמיד תהיה מדוייקת? נמק. | _ _ _ ה. |
| | _ _ _ |
| | _ _ _ |

| . (4 נק') אם תשובה בסעיף הקודם "אינה תמיד מדוייקת", תציע עדכונים לאינקס לתיקון הבעייה. אם |
|--|
| התשובה בסעיף הקודם "תמיד מדוייקת" תראה מה באינקס מבטיח את זה. |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| שאלה 4 (20 נק') |
| נתון מאגר של ששה מסמכים (ראה שאלה 2). יש לדרג את המסמכים לפי שיטת TextRank באופן הבא: |
| ו. (10 נק') יש לבנות גרף של מסמכים כך שקדקודים מייצגים מסמכים וקשתות מייצגות דמיון ביניים. |
| יש להשתמש ב-Jaccard similarity |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| · |
| |
| |
| |
| (0.8=0.8)עבור הדפים (חשבו שלוש איטרציות, מקדם השיכוך PageRank ב: |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |

 $Jaccard(A,B) = |A \cap B| / |A \cup B|$

Cosine similarity:
$$\cos(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{\left|\vec{q}\right| \left|\vec{d}\right|} = \frac{\vec{q}}{\left|\vec{q}\right|} \cdot \frac{\vec{d}}{\left|\vec{d}\right|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

Term frequency:

$$tf(t,d) = \frac{f(t,d)}{\max\{f(w,d) : w \in d\}}$$

Inverse document frequency (idf):

$$\operatorname{idf}(t,D) = \log rac{N}{|\{d \in D : t \in d\}|}$$

Page Rank:

$$PR(A) = (1-d) + d \sum_{i=1}^{n} \frac{PR(T_i)}{C(T_i)}$$

Euclidean distance (in n-space):

$$egin{align} d(\mathbf{p},\mathbf{q}) &= d(\mathbf{q},\mathbf{p}) = \sqrt{(q_1-p_1)^2 + (q_2-p_2)^2 + \dots + (q_n-p_n)^2} \ &= \sqrt{\sum_{i=1}^n (q_i-p_i)^2}. \end{array}$$

Mutual Information:

$$I(w,c) = \sum_{e_w \in \{0,1\}} \sum_{e_c \in \{0,1\}} p(e_w, e_c) \log \frac{p(e_w, e_c)}{p(e_w)p(e_c)}$$



טיוטה

