

## מדור בחינות ומערכת שעות

הנדסת תוכנה

14/01/19 09:00-12:00

## **מבוא לאחזור מידע** מועד א' ליטבק מרינה

'תשע"ט סמסטר א

חומר עזר – אסור (מלבד דפי הנוסחאות המצורפים לטופס זה)
השאלון מכיל 13 עמודים (כולל נספח וטיוטה).
•.
<u>חומר עזר : נא סמן במשבצת המתאימה את המתאים</u>
ניתן להשתמש בכל מחשבון
Casio FX-991EX לא ניתן להשתמש במחשבון * _V_
* לא ניתן להשתמש במחשבון
_V_ <mark>* לא</mark> ניתן להשתמש בחומר עזר
* מותר שימוש בדף נוסחאות, כמפורט:
* הבחינה בחומר פתוח – מותר להשתמש בכל חומר עזר מודפס או כתוב
הערות
<u>הערות</u> _V_ יש לענות על <u>כל</u> השאלות במקומות המיועדים ע"ג טופס השאלון בלבד
יש להחזיר את השאלון ביחד עם הכריכה/מחברת.
אחר:

השאלון מכיל 🧕 עמודים (כולל עמוד זה).

בהצלחה!



נק")	25)	1	שאלה

y/w' •	מות לכל וימאלוע
יש לס •	זמן באופן ברור את התשובה <u>הנכונה ביותר</u> על גבי <b>שאלון הבחינה</b>
• סימון	של יותר מתשובה אחת לאותה שאלה יקבל ציון של <u>אפס</u>
א. (5נק')	) ערך אפשרי של Jaccard similarity נכלל בתחום הבא:
• • •	[0, 1]
<del>-</del> 3	$ \begin{bmatrix} 0, \infty \end{bmatrix} $ $ \begin{bmatrix} -\infty, 1 \end{bmatrix} $
<u>1</u>	
7	[-w,w <sub>]</sub> .
תשובה (1	(4/ 3/ 2/
/ 5\ ·	t
	) תפקיד של-clustering הוא:
	. לחשב שכיחות של המילים במסמך
2	2. לצמצם כל מילה לשורש-מילה
3	לזהות קבוצות של מסמכים דומים
4	. להוציא מילות מפתח עבור מסמך
תשובה (1	
ג. (5נק')	) אלגוריתם למידה בשם NaiveBayes מקבל כקלט מסמכים בצורה של:
	(Vector Space Model) ווקטור
	(Bag of words) BOW
	גרף גרף
	 2. טקסט ללא מילות עצירה
•	
תשובה (1	(2 / 3 / 4 / 3 ) נימק:
ד. (5נק')	:מילות עצירה (stop words) אלא מילים שניתן לאפיין על-ידי
(1/13) ii	idf (inverse document frequency) . גבוה ו- (term frequency) .
2	idf (inverse document frequency) גבוה (term frequency) גבוה idf (inverse document frequency) גבוה
2	77123 idf (inverse document frequency) -17123 if (term frequency)
	idf (inverse document frequency) - נמוך ו- (term frequency) נמוך ו- (לוישים לא לוישים
4	idf (inverse document frequency) - נמוך ו- tf (term frequency) .
תשובה (1	(2 / 3 / 4 / 3 בימק:
	and the same of th
	) נתון ביטוי "Natural Language Processing". ניתן לאחזר מסמכים שמכילינ
אותו נ	בעזרת:
1	בנוי על מילים בודדות ושאילטה (inverted index). אינדקס אפוך
	"Natural AND Language AND Processing"
2	ושאילטה (bi-grams) אינדקס אפוך בנוי על זוגות המילים
	"Natural Language AND Language Processing"
3	מוניסטסיר במודעות ביווי על מילים בודדות עם מיקומם בתוך מסמך ושאילטה
-	"Natural AND Language AND Processing"
4	ר באוק ב-cosine similarity בין ווקטורים של מסמכים ושאילטה cosine similarity
•	
תשובה (1	(4 / 3 / 2 / 2 ) נימק:
•	,



		h c a h c a :D6
		מסווגים לשלוש קטגוריות: P, B and S באופן הבא:
מסמך	קטגוריה	
D1	P	
D2	N	
D3	N	
D4	P	
D5	N	
D6	P	
לסנו אותם	אד ישנה דרישה .	א. (10 נק') לא נתונה רשימת מילות עצירה (stop words).
		מהסמכים. כיצד ניתן לעשות זאת? תציע שיטה ותראה או
	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	יש להראות את התוכן המסמכים לאח <u>ר סינון.</u>
		The original property for the first that the
		<del></del> -
•	<u> </u>	
<del></del>		
		:D1
		:D2
		:D3
		:D4
		:D5
		:D6
2	רעזרת שייוה <sup>2</sup> .	Do
^	(	ב. (כו נק) ש קבבע בוווון נואב ב ב (נוסוסססססס
		<del></del>
	<del></del>	
	-	
	<del></del>	
·		<del></del>
	·	

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

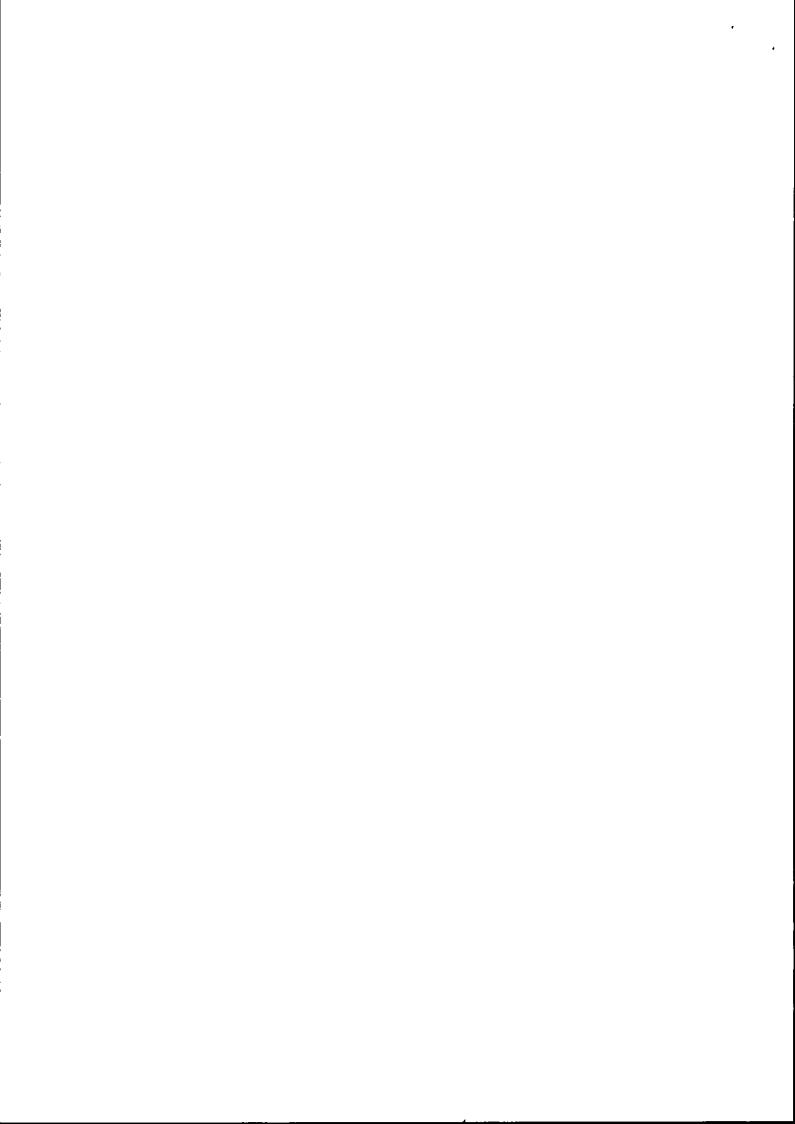
d h d :D5



	ולה 3 (00 נית ייייור
מסמכים (ב. ב.) ניון לבנות מובל במוסף אונים אונים תנונו מתחבות	
עבור סיווג מסמכים NaiveBayes עבור מודל 15)	Х.
. 33 . 3333333	
	-
(5 נק') נמק למה בחרת בסוג האלגוריתם (סוג של-NB) שהשתמשת בו? אילו עוד סוגייש?	ב.
(3נק') סווגו את שלושת מסמכי המבחן לפי מודל הנלמד בסעיף ב'	ړ.
bhcde:D7 .i bhd:D8 .ii	
a c d a e :D9 .iii	
acdae.D9 .III	
	<del></del>
על שלושת מסמכי המבחן (test accuracy) על שלושת מסמכי המבחן אם ידוע סיווג (5 נק')	



מאלה 4 (20 נקי) נתון מאגר של ששה מסמכים המתוארים בשאלה מסי 2 ושאילתה a h d h a :Q. יש לדרג שלושה מסמכים הראשונים (D1, D2, D3) ביחס לשאילתה תוך שימוש ב- tf-idf ו- cosine similarity .idf יש להתחשב בכל המאגר לחישוב.
<del></del>
<del></del>
<del></del>



NaiveBayes Classifier:

$$1 = \arg\max_{c} \hat{P}(c) \prod_{i} \hat{P}(x_{i} \mid c)$$

$$\hat{P}(c_{j}) = \frac{N(C = c_{j})}{N}$$

$$\hat{P}(x_{i} \mid c_{j}) = \frac{N(X_{i} = x_{i}, C = c_{j}) + 1}{N(C = c_{j}) + k}$$

Multinomial NB = One feature  $X_i$  for each word position in document

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

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

Term frequency:

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

Inverse document frequency (idf):

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

Chi-square ( ): 
$$X^{2}(\mathbb{D}, t, c) = \sum_{e_{t} \in \{0,1\}} \sum_{e_{c} \in \{0,1\}} \frac{(N_{e_{t}e_{c}} - E_{e_{t}e_{c}})^{2}}{E_{e_{t}e_{c}}}$$



:טיוטה



