מספר התלמיד הנבחן

רשום את כל תשע הספרות

# האוניברסיטה

הדבק כאן את מדבקת הנבחן

י"ב בתמוז תשע"ו

18

ביולי 2016

סמסטר 2016ב

86 מס' מועד

מס' שאלון - 535

20595 / 4

שאלון בחינת גמר

20595 - כריית מידע

משך בחינה: שעות

> בשאלון זה 5 עמודים

> > מבנה הבחינה:

בבחינה שני חלקים.

יש לענות על כל השאלות.

:חלק א

ענו על ארבע השאלות שבחלק זה.

90 נקודות לחלק זה.

חלק ב:

ענו על השאלה שבחלק זה.

10 נקודות לחלק זה.

משקל כל שאלה מפורט בגוף השאלון.

אפשר להשתמש בהנחות במהלך הפתרון, אם תזדקקו לכך.

פרטו את הנחותיכם!

חומר עזר:

כל חומר עזר מותר בשימוש.

אסור בשימוש כל מכשיר אלקטרוני שבאמצעותו ניתן לאצור מידע

לרבות מכשיר טלפון נייד, מחשב נישא, שעון חכם וכד'.

החזירו

למשגיח את השאלון

וכל עזר אחר שקיבלתם בתוך מחברת התשובות



בהצלחה !!!

## חלק א (90 נקודות)

בחלק זה ארבע שאלות. ענו במפורט על ארבע השאלות.

שאלה 1 (30 נקודות) – עצי החלטה (decision tree)

נתונה טבלת נתוני אימון:

| סוג טיפול | פרופיל שומנים | סכרת | לחץ דם    | מדד מסת גוף | מס' נבדק |
|-----------|---------------|------|-----------|-------------|----------|
| דיאטה     | תקין          | לא   | תקין      | 27          | 1        |
| ללא טיפול | לא תקין       | לא   | גבוה      | 22          | 2        |
| תרופתי    | תקין          | לא   | תקין      | 34          | 3        |
| ללא טיפול | תקין          | כן   | גבוה מאוד | 25          | 4        |
| תרופתי    | לא תקין       | כן   | תקין      | 29          | 5        |
| תרופתי    | תקין          | לא   | תקין      | 32          | 6        |
| תרופתי    | תקין          | לא   | גבוה      | 31          | 7        |
| דיאטה     | תקין          | לא   | גבוה מאוד | 30          | 8        |
| דיאטה     | תקין          | כן   | תקין      | 26          | 9        |
| ניתוחי    | לא תקין       | כן   | גבוה      | 33          | 10       |
| תרופתי    | תקין          | לא   | גבוה מאוד | 32          | 11       |
| ללא טיפול | תקין          | כן   | גבוה      | 18          | 12       |
| דיאטה     | תקין          | לא   | תקין      | 28          | 13       |
| דיאטה     | לא תקין       | כן   | גבוה מאוד | 30          | 14       |

א. בנו עץ החלטה חלקי, הכולל את רמת השורש ורמה אחת נוספת בלבד, עבור נתוני האימון שבטבלה לחיזוי החלטה מהו הטיפול שעל הנבדק לקבל לאור מדד מסת הגוף שלו. בתשובתכם הדגימו את שלבי בחירת התכונה המפצלת בעץ.

הערה: יש לכלול חישוב של אחד המדדים כדוגמת אנטרופיה, Gain ratio, מדד גיני.

ב. האם ניתן לגזום את עץ ההחלטה שבניתם בסעיף א! נמקו תשובתכם.

#### המשך הבחינה בעמוד הבא

## (association rules) – חוקי הקשר – חוקי נקודות) – אלה 2

: נתון בסיס הנתונים

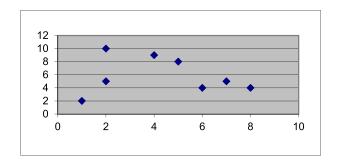
| דרגת חומרת<br>התאונה | חגירת חגורה | עבירות תנועה         | נהיגה בשכרות | תנאי מזג אויר |
|----------------------|-------------|----------------------|--------------|---------------|
| חמורה                | לא          | נהיגה במהירות מופרזת | כן           | טוב           |
| קלה                  | כן          | אין                  | לא           | גרוע          |
| קלה                  | כן          | אי ציות לתמרור עצור  | לא           | טוב           |
| חמורה                | כן          | נהיגה במהירות מופרזת | לא           | טוב           |
| חמורה                | לא          | נהיגה ברמזור אדום    | לא           | גרוע          |
| קלה                  | כן          | אי ציות לתמרור עצור  | כן           | טוב           |
| חמורה                | כן          | אין                  | כן           | גרוע          |
| חמורה                | כן          | נהיגה ברמזור אדום    | לא           | טוב           |
| חמורה                | לא          | אין                  | כן           | טוב           |
| חמורה                | לא          | נהיגה ברמזור אדום    | לא           | גרוע          |
| חמורה                | כן          | נהיגה במהירות מופרזת | כן           | טוב           |
| קלה                  | כן          | אי ציות לתמרור עצור  | לא           | גרוע          |

- א. בצעו בינאריזציה לבסיס הנתונים והציגו אותו.
- ב. הגדירו מהו הרוחב המקסימלי עבור כל אחת מהטרנזקציות בבסיס הנתונים הבינארי!
  - ג. בהנחה %min\_support= 30 מצאו את כל הקבוצות התדירות.

המשך הבחינה בעמוד הבא

## שאלה 3 (25 נקודות) – ניתוח אשכולות (clustering)

.k-means בצעו אשכול לשמונה הנקודות המתוארות באיור הבא תוך שימוש באלגוריתם



#### בתשובתכם הניחו:

- k = 3 .
- (2,10) (5,8) (1,2) ב. מרכזי האשכולות ההתחלתיים הם:
  - ג. השתמשו במרחק אוקלידי
  - ד. תנאי סיום: לא חלים שינויים באשכול

### שימו לב,

בתשובתכם הדגימו את שלושת האשכולות ההתחלתיים וכן את כל השלבים והגדירו את שלושת האשכולות הסופיים.

### שאלה 4 (20 נקודות) – למידה בייאסינית

במחשב הביתי התקבלו הודעות המייל הבאות:

| D1:"Send us your password" | spam |
|----------------------------|------|
| D2:"Send us your review"   | ham  |
| D3:"Review your password"  | ham  |
| D4:"Review us"             | spam |
| D5:"Send your password"    | spam |
| D6:"Send us your account"  | spam |

. naïve Bayes בעזרת המסווג (spam) יש לזהות את התולעת

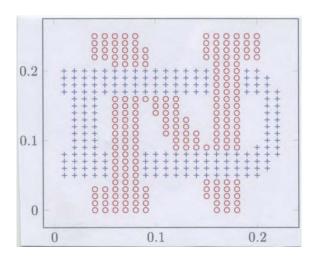
#### המשך הבחינה בעמוד הבא

## חלק ב (10 נקודות)

בחלק זה שאלה אחת. ענו עליה במפורט.

### שאלה 5 (10 נקודות)

את מידע מידע בקורס כריית הסטודנטים את נתוניו שיש לסווג את מדעם בסיס נתונים שיש לסווג את מחונים בסיס נתונים שיש לסווג את מחונים הבאים: עץ החלטה, המסווגים הבאים: עץ החלטה,



בחנו את ביצועי שלושת המסווגים וציינו מיהו המסווג בעל הביצועים הגבוהים ביותר.

## בהצלחה!