



המכללה האקדמית להנדסה סמי שמעון

מדור בחינות ומערכת שעות

המחלקה להנדסת תוכנה

25/6/2019
09:00-12:00

מבוא לכריית נתונים

מועד א'

יניב הדר

תשע"ט סמסטר ב'

חומר עזר – דף נוסחאות

השאלון מכיל 8 עמודים.

בהצלחה !

=====

חומר עזר : נא סמן במשבצת המתאימה את המתאים

X * ניתן להשתמש בכל מחשבון (ללא יכולות תכנות)

___ * לא ניתן להשתמש במחשבון Casio FX-991EX

___ * לא ניתן להשתמש במחשבון

___ * לא ניתן להשתמש בחומר עזר

X * מותר שימוש בדף נוסחאות, כמפורט: דף אחד משני הצדדים (סטודנטים מביאים עימם)

___ * הבחינה בחומר פתוח – מותר להשתמש בכל חומר עזר מודפס או כתוב

הערות

X יש לענות על כל השאלות במקומות המיועדים ע"ג טופס השאלון בלבד

X יש להחזיר את השאלון ביחד עם הכריכה/מחברת.

אחר:

1.

2.

3.

השאלון מכיל 8 עמודים (כולל עמוד זה).

בהצלחה !

=====

המכללה האקדמית להנדסה סמי שמעון

קמפוס באר שבע ביאליק פינת בזל 84100 | קמפוס אשדוד ז'בוטינסקי 77245,84 | www.sce.ac.il | חייג: *ממפסס

חלק א' (40 נקודות, 4 נקודות לכל שאלה)

יש לבחור תשובה אחת בלבד לכל שאלה ולסמן על גבי השאלון בלבד – תשובות במחברת לא ייבדקו.

1. אדם ביצע בדיקה לנטיחות מולה מסוימת. הבדיקה יצאה ודנית הוא אכן ודלה. זוהי דוגמה ל:
 - א. False negative
 - ב. False positive
 - ג. True positive
 - ד. True negative
2. המרחק בין שני אשכולות בשיטת Centroid מורשב עפ"י
 - א. המרחק בין מרכז האשכולות
 - ב. המרחק הארוך ביותר בין שני איינקטים בשני האשכולות
 - ג. המרחק הממוצע בין שני האשכולות
 - ד. המרחק הקצר ביותר בין שני איינקטים בשני האשכולות
3. חלוקת ערכים לאינטרוולים לפי עומק שווה:
 - א. מקטינה את סטיית התקן של התוצרים
 - ב. יוצרת אינטרוולים בעלי שכידות שווה
 - ג. יוצרת אינטרוולים בעלי טוחח שווה
 - ד. תשובות א' ו-ב' נכונות
4. עפ"י תורת האינפורמציה, אי ודאות האירוע שווה למינימום אם:
 - א. כל התוצאות הן שרר דטרמיניסטי (קבוע) יחיד
 - ב. התוצאות מתפלגות התפלגות משריכית
 - ג. התוצאות מתפלגות התפלגות נורמאלית
 - ד. התוצאות מתפלגות התפלגות אחידה
5. למה אלגוריתם Naïve Bayes נאריב?
 - א. התוצה שקיימת תלות זוהי בין כל משתנה למשתנה היא נאריבית
 - ב. התוצה שלא קיימת תלות בין אף אחד מהמשתנים היא נאריבית
 - ג. התוצה שניתן לסטם את המשתנים היא נאריבית
 - ד. התוצה שקבוצת על של קבוצת פריטים שכידים ודיבת להיחז שכידה היא נאריבית
6. אלגוריתם מסוג eager לעומת אלגוריתם מסוג lazy:
 - א. eager שומר את כל נתוני האימון ו-lazy שומר חלק מנתוני האימון
 - ב. eager שומר חלק מנתוני האימון ו-lazy שומר את כל נתוני האימון
 - ג. eager שומר את נתוני האימון ו-lazy שומר מודל
 - ד. eager שומר מודל ו-lazy שומר את נתוני האימון
7. בית אריזה מעניין להפריד בעזרת צילום תמונות, בין תפוחים למנצות וללא סיווג ראשוני. מדובר במשימה של:
 - א. גיחוח אשכולות (Clustering)
 - ב. סיווג (Classification)
 - ג. חיזוי (Prediction)
 - ד. למידה מונחית

8. כדי לדעת מה הסיכויים שקבוצת פריטים אכן תיבדק ביחד, המדד שצריך לחשב הוא:

- א. confidence
- ב. mutual information
- ג. conditional probability
- ד. support

9. בבעיית סיווג בינארי, דיוק האימון של מודל בעל אנטרופיה השווה ל-1 הוא:

- א. 50%
- ב. 0%
- ג. 100%
- ד. לא ניתן לדעת

10. נתונה הפונקציה $Y=2X+3$. האנטרופיה של Y בהינתן X שווה ל:

- א. 1
- ב. 0.5
- ג. 0
- ד. לא ניתן לדעת

חלק ב' (60 נקודות)

בכל הסעיפים בחלק זה יש להראות את כל החישובים והלוגיקים במתכונת יש לרשום תשובה סופית על גבי השאלון בלבד. הניקוד הוא על תשובה סופית מדויקת בלבד. אולם תשובה ללא חישובים ולוגיקים במתכונת תיפסל.

זואולוגים מעניינים לחזות האם נחש הוא ארסי. לשם כך, נאסף מידע מ 10 נחשים והכיל 5 משתנים מתעמדים ומשתנה מסרה אחד – סיווג הנחשים בסכנה:

| סיווג | ניבים | צבע אחיד | צבע קשקשים עיקרי | אורך | מין | רשומה |
|-------|-------|----------|------------------|------|--------|-------|
| yes | TRUE | TRUE | yellow | 24 | male | 1 |
| no | FALSE | FALSE | black | 32 | male | 2 |
| no | TRUE | FALSE | black | 15 | male | 3 |
| yes | TRUE | TRUE | brown | 18 | male | 4 |
| yes | FALSE | FALSE | yellow | 27 | female | 5 |
| yes | FALSE | TRUE | yellow | 9 | male | 6 |
| no | TRUE | FALSE | black | 3 | female | 7 |
| yes | FALSE | FALSE | black | 18 | female | 8 |
| no | FALSE | TRUE | black | 22 | male | 9 |
| yes | FALSE | FALSE | brown | 21 | male | 10 |

א. (4 נק') יש לבצע נרמול בשיטת decimal scaling לאורך הוודאים

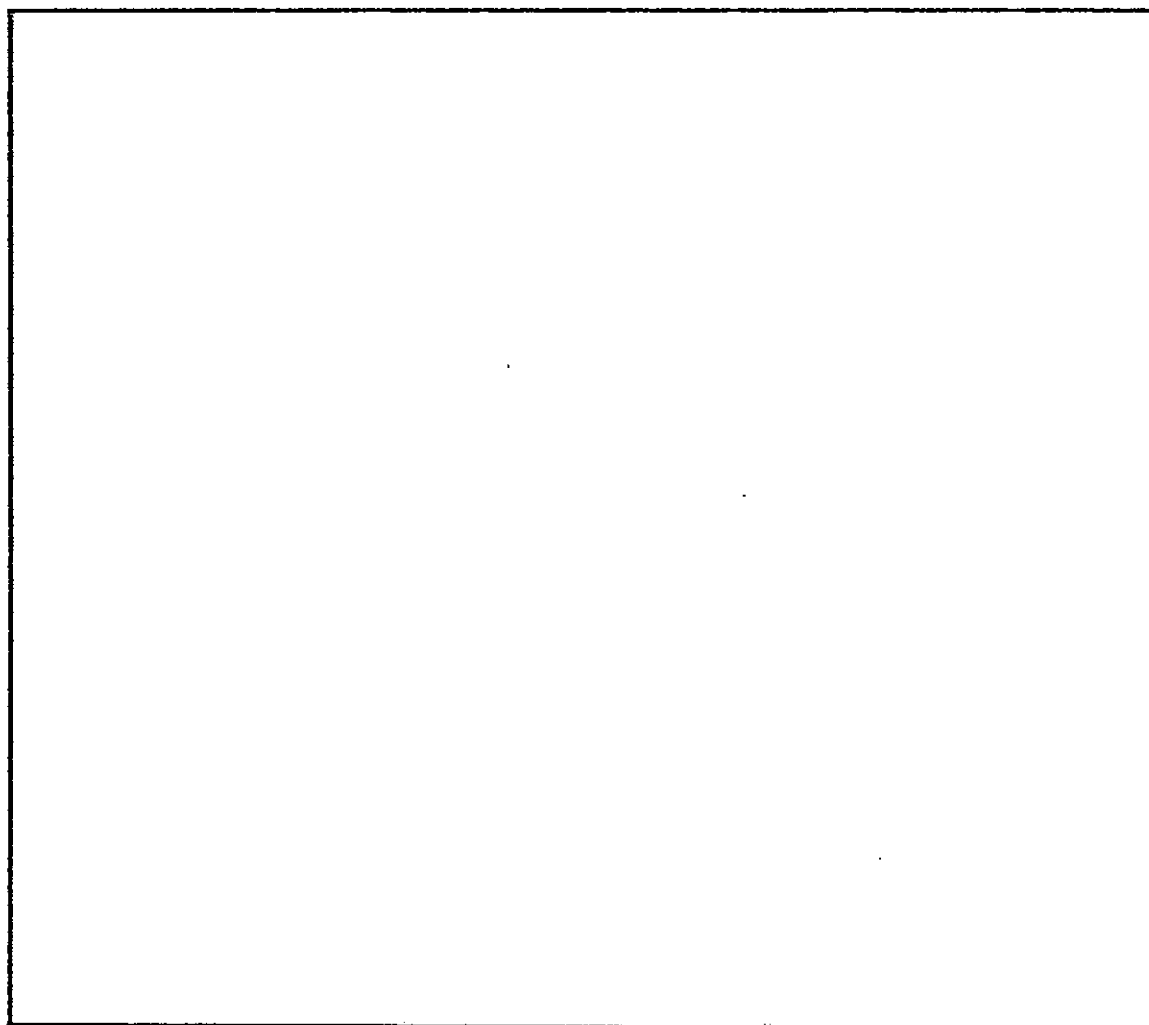
ב. (4 נק') יש לבצע דיסקרטיוזיה לשני טוחמים (BINS) בשיטת צימק שחה לערכי אחד הוחלים המנורמל.

ג. (4 נק') מהו דיוק חזק הרוב בנתונים שלה?

ד. (6 נק') מהו האנטרופיה של משתנה המטרה סידוג?

ה. (6 נק') האם מין בשך male ו- צבע אחיד בשך FALSE הם חבית שכחה?
(Minimum support = 30%)

1. (12 נק) יש לחשב את מדד ה-Gini (Gini index) של משתנה המסווג – "סיווג" עבור כל אחת מהתכונות "גברים" ו- "צבע אחיד" ולהחליט לפי איזו תכונה עדיף לפצל את קדקוד השורש.



ז. (12 נק) יש לחשב את מדד ה- $gain\ ratio$ של משתנה המטרה – "סייג" עבור כל אחת מהתכונות "צבע קשקשים" ו- "צבע אחיד" ולהחליט לפי איזו תכונה עדיף לפצל את קודקוד השורש.

ח. (12 נק') יש להשתמש באלגוריתם KNN (K nearest neighbors) על מנת לסווג את התצפית הבאה (K=1):

| סינוג | גיבים | צבע אחיד | צבע קשקשים עיקרי | אורך | מין | רשומה |
|-------|-------|----------|------------------|------|--------|-------|
| ??? | TRUE | FALSE | yellow | 50 | female | 1 |

על מנת לחשב מרחק בין תצפיות יש להשתמש במרחק מנהטן כאשר המרחקים מגודרים באופן הבא:

- עבור המשתנים מין, גיבים וצבע אחיד המרחק בין ערכים שונים הוא 1.
- עבור המשתנה אורך יש להשתמש במשתנה לאורך הדיסקרטיוזיה (מסעיף ב') – המרחק בין אינטרוולים שונים הוא 1.
- עבור המשתנה "צבע קשקשים עיקרי" יש להשתמש בטבלת המרחקים הבאה:

| צבע קשקשים עיקרי | brown | yellow | black |
|------------------|-------|--------|-------|
| brown | 0 | | |
| yellow | 1 | 0 | |
| black | 2 | 1 | 0 |

