

אוניברסיטה הפתוחה

20595

## **כריית מידע**

חוברת הקורס אביב 2015ב

כתבה: ד"ר מיה הרמן

מרץ 2015 - סמסטר אביב – תשע"ה

**פנימי – לא להפצה.**

© כל הזכויות שמורות לאוניברסיטה הפתוחה.

## תוכן העניינים

1	אל הסטודנט
2	1. לוח זמנים ופעילויות
4	2. תיאור המטלות
4	2.1 מבנה המטלות
4	2.2 חומר הלימוד הדרוש לפתרון המטלות
4	2.3 ניקוד המטלות
5	3. התנאים לקבלת נקודות זכות בקורס
7	ממ"ן 11
9	ממ"ן 21 (פרויקט)
13	ממ"ן 12
15	ממ"ן 22 (פרויקט)



## אל הסטודנט

אנו מקדמים את פניך בברכה עם הצטרפותך אל הלומדים בקורס "כריית מידע".

בחוברת זו תמצא את "לוח הזמנים ופעילויות", תנאים לקבלת נקודות זכות ומטלות הקורס.

לקורס קיים אתר באינטרנט בו תמצאו חומרי למידה נוספים, אותם מפרסם/מת מרכז/ת ההוראה. בנוסף, האתר מהווה עבורכם ערוץ תקשורת עם צוות ההוראה ועם סטודנטים אחרים בקורס. פרטים על למידה מתוקשבת ואתר הקורס, תמצאו באתר שה"ס בכתובת:

<http://telem.openu.ac.il>

מידע על שירותי ספרייה ומקורות מידע שהאוניברסיטה מעמידה לרשותכם, תמצאו באתר הספרייה באינטרנט [www.openu.ac.il/Library](http://www.openu.ac.il/Library).

ייעוץ ינתן ביום ד' בין השעות 11:30-9:30 בטלפון 09-7781260. פגישה יש לתאם מראש.

ניתן לפנות גם בדואר אלקטרוני [maya@openu.ac.il](mailto:maya@openu.ac.il)

אני מאחלת לך לימוד פורה ומהנה.

בברכה,

ד"ר מיה הרמן  
מרכזת הקורס

**1. לוח זמנים ופעילויות (20595 / ב'2015)**

שבוע לימוד	תאריכי שבוע הלימוד	יחידת הלימוד המומלצת	מפגשי ההנחיה*	תאריך אחרון למשלוח ממ"ן (למנחה)
1	13.3.2015-10.3.2015	יחידה 1 מבוא		
2	20.3.2015-15.3.2015	יחידה 2 תורת המידע	מפגש ראשון	
3	27.3.2015-22.3.2015	יחידה 3 הכנת נתונים		
4	3.4.2015-29.3.2015 (ו' ערב פסח)	יחידה 4 סיווג וחיזוי	מפגש שני	
5	10.4.2015-5.4.2015 (א-ו פסח)	יחידה 5 עצי החלטה- חלק א		
6	17.4.2015-12.4.2015 (ה' יום הזכרון לשואה)	יחידה 6 עצי החלטה- חלק ב	מפגש שלישי	
7	24.4.2015-19.4.2015 (ד' יום הזכרון) (ה' יום העצמאות)	יחידה 7 למידה בייסיאנית ולמידה מבוססת תצפיות		ממ"ן 11 19.4.2015
8	1.5.2015-26.4.2015	יחידה 8 חוקי הקשר – חלק א		
9	8.5.2015-3.5.2015 (ה' ל"ג בעומר)	יחידה 9 חוקי הקשר – חלק ב	מפגש רביעי	

\* התאריכים המדויקים של המפגשים הקבוצתיים מופיעים ב"לוח מפגשים ומנחים".

לוח זמנים ופעילויות - המשך

שבוע הלימוד	תאריכי שבוע הלימוד	יחידת הלימוד המומלצת	מפגשי ההנחיה*	תאריך אחרון למשלוח הממ"ן (למנחה)
10	15.5.2015-10.5.2015	יחידה 10 ניתוח אשכולות- חלק א		ממ"ן 21 10.5.2015
11	22.5.2015-17.5.2015 (א יום ירושלים)	יחידה 11 ניתוח אשכולות- חלק ב	מפגש חמישי	
12	29.5.2015-24.5.2015 (א שבועות)	יחידה 12- <i>ת/ע</i> רשתות אינפו-עמומות	מפגש שישי	ממ"ן 12 24.5.2015
13	5.6.2015-31.5.2015	יחידה 13 <i>ת/ע</i> בחירת מאפיינים		
14	12.6.2015-7.6.2015	יחידה 14 <i>ת/ע</i> נושאים מתקדמים בכריית מידע		ממ"ן 22 7.6.2015
15	23.6.2015-14.6.2015	חזרה	מפגש שביעי	

מועדי בחינות הגמר יפורסמו בנפרד

\* התאריכים המדויקים של המפגשים הקבוצתיים מופיעים ב"לוח מפגשים ומנחים".

## 2. תיאור המטלות בקורס

<b>קרא היטב עמודים אלו לפני שתתחיל לענות על השאלות</b>
--

פתרון המטלות הוא חלק בלתי נפרד מלימוד הקורס – הבנה מעמיקה של חומר הלימוד דורשת תרגול רב. המטלות תבדקנה ותוחזרנה לך בצירוף הערות המתייחסות לתשובות.

### 2.1 מבנה המטלות

כל מטלה מורכבת מכמה שאלות.

את הפתרונות למטלה עליך להדפיס. רצוי להשאיר שוליים רחבים להערות המנחה.

אם השאלה בממ"ן אינה ברורה לך, אל תהסס להתקשר אל המנחה בשעות הייעוץ הטלפוני בלבד לצורך קבלת הסבר.

המטלות מלוות את יחידות הלימוד בקורס. להלן פירוט המטלות והיחידות שאליהן מתייחסת כל מטלה.

### 2.2 חומר הלימוד הדרוש לפתרון המטלות

ממ"ן 11 – יחידות לימוד 1-6 – רשות 2 נקודות.

ממ"ן 21 – יחידות לימוד 1-6 – **חובה** – 13 נקודות (פרויקט – שלב א).

ממ"ן 12 – יחידות לימוד 7-8 – רשות 2 נקודות.

ממ"ן 22 – יחידות לימוד 7-11 – **חובה** – 13 נקודות (פרויקט – שלב ב).

**ממ"נים 21 ו-22 (פרויקט):**

מכיוון ואלו **מטלות חובה** ומהווה משקל רב בציון הסופי, אין להגישן באיחור ללא קבלת אישור מראש. על – כן הקפד לשלוח את המטלות במועד.

### 2.3 ניקוד המטלות

סה"כ ניתן לצבור 26 - 30 נקודות במטלות.

מטלות החובה בקורס כוללות פרויקט המוגש בשני שלבים ומהוות יחד 26 נקודות. מומלץ להגיש את כל המטלות



### 3. התנאים לקבלת נקודות זכות בקורס

- א) הגשת מטלות מנחה 21 ו- 22 (פרויקט חובה).
- ב) ציון של לפחות 60 נקודות בפרויקט.
- ג) ציון של לפחות 60 נקודות בבחינת הגמר.
- ד) ציון סופי בקורס של 60 נקודות לפחות.

**לבחינת הגמר רשאי לגשת רק סטודנט שצבר 26 נקודות לפחות.**

#### **לתשומת לבכם!**

כדי לעודדכם להגיש לבדיקה מספר רב של מטלות הנהגנו את ההקלה שלהלן:

אם הגשתם מטלות מעל למשקל המינימלי הנדרש בקורס, **המטלות** בציון הנמוך ביותר, שציוניהן נמוכים מציון הבחינה (**עד שתי מטלות**), לא יילקחו בחשבון בעת שקלול הציון הסופי.

זאת בתנאי שמטלות אלה **אינן חלק מדרישות החובה בקורס** ושהמשקל הצבור של המטלות האחרות שהוגשו, מגיע למינימום הנדרש.

**זכרו!** ציון סופי מחושב רק לסטודנטים שעברו את בחינת הגמר בציון 60 ומעלה והגישו מטלות כנדרש באותו קורס.

**הכנת המטלות חייבת להעשות על-ידי כל סטודנט בנפרד.**

**מטלות שלא יבוצעו באופן עצמאי – יפסלו!!!**



# מטלת מנחה (ממ"ן) 11

הקורס: 20595 - כריית מידע

חומר הלימוד למטלה: יחידות 1-6

משקל המטלה: 2 נקודות

מספר השאלות: 2

מועד אחרון להגשה: 19.4.2015

סמסטר: 2015

אנא שים לב:

מלא בדיוקנות את הטופס המלווה לממ"ן בהתאם לדוגמה שלפני המטלות.  
העתק את מספר הקורס ומספר המטלה הרשומים לעיל.

בממ"ן זה שתי שאלות המתייחסות לטבלת נתוני האימון הרצ"ב.

ניתן להשתמש בתוכנת WEKA וכן בגיליון אלקטרוני EXCEL. ענו במפורט על שתי השאלות.

בתשובתכם, ציינו כל הנחה שבצעתם וצרפו את התוכנות והקבצים שערכתם בהם שימוש בחישובים.

נתונה טבלת נתוני אימון:

מס' נבדק	מדידת חלופית	המתנה נוחה	שישי/שבת	כאבים	תפוסת חולים במרפאה	עלות	נקבע תור מראש	מומחיות רופא	זמן המתנה משוער	גשום	ימתין
1	כ	ל	ל	כ	בינונית	גבוהה	כ	פנימאי	0-10	?	כ
2	כ	ל	ל	כ	מלאה	נמוכה	ל	אורטופד	30-60	ל	ל
3	ל	כ	?	ל	חלקית	נמוכה	ל	נשים	0-10	ל	כ
4	כ	ל	כ	כ	מלאה	נמוכה	ל	אורטופד	10-30	ל	כ
5	כ	ל	כ	ל	מלאה	גבוהה	כ	פנימאי	60>	ל	ל
6	ל	כ	ל	כ	חלקית	בינונית	כ	ילדים	ככככ	כ	כ
7	ל	כ	ל	ל	אפסית	נמוכה	ל	נשים	0-10	כ	ל
8	ל	ל	ל	כ	99999	בינונית	כ	אורטופד	0-10	כ	כ
9	ל	כ	כ	ל	מלאה	נמוכה	ל	נשים	60>	כ	ל
10	כ	כ	כ	כ	מלאה	גבוהה	כ	ילדים	10-30	ל	ל
11	ל	ל	ל	ל	אפסית	נמוכה	ל	אורטופד	0-10	ל	ל
12	כ	כ	כ	כ	מלאה	נמוכה	ל	?	30-60	ל	כ

### שאלה 1 (25 נקודות)

ציינו והדגימו את שלבי הכנת הנתונים לביצוע כריית מידע כדוגמת טיפול בערכים חסרים, ערכים שגויים ועוד. בסיום, בנו בסיס נתונים מטוייב.

### שאלה 2 (75 נקודות)

א. בנו עץ החלטה חלקי, הכולל את רמת השורש ורמה אחת נוספת בלבד, עבור נתוני האימון שבטבלה לחיזוי החלטה האם הנבדק ימתין לרופא. בתשובתכם הדגימו את שלבי בחירת התכונה המפצלת בעץ.

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

ב. איזה מבין התכונה/תכונות ניתן להסיר ומדוע? אם אין תכונה הניתנת להסרה, יש לציין זאת מפורשות.

# מטלת מנחה (ממ"ן) 21 - פרויקט גמר

הקורס: 20595 - כריית מידע

חומר הלימוד למטלה: יחידות 1-6

מספר השאלות: 2

משקל המטלה: 13 נקודות

סמסטר: 2015

מועד אחרון להגשה: 10.5.2015

אנא שים לב:

מלא בדייקנות את הטופס המלווה לממ"ן בהתאם לדוגמה שלפני המטלות.  
העתק את מספר הקורס ומספר המטלה הרשומים לעיל.

## הנחיות

נתון קובץ נתוני שיחות שיווק בנקאי טלפוני המצוי בכתובת:

<http://archive.ics.uci.edu/ml/datasets/Bank+Marketing>

נתוני הקובץ תוכלו למצוא בכתובת:

<http://archive.ics.uci.edu/ml/machine-learning-databases/00222>

מטרות הפרויקט:

- חיזוי לקוח שיחתום על תנאי הפקדון שהוצעו לו במהלך שיחות השיווק הטלפוני – ממ"ן 21.
- איפיון פרופיל לקוח שיחתום על תנאי הפקדון שהוצעו לו במהלך שיחות השיווק הטלפוני – ממ"ן 22.

הפרויקט כולל שימוש בחבילת תוכנה WEKA לכריית מידע.

הפרויקט יבוצע בשני שלבים:

- א. ממ"ן 21 – בשלב הראשון תידרשו להגדיר את הבעיה, להכין את הנתונים ולפתור את הבעיה בעזרת שיטות סיווג וחיזוי.
- ב. ממ"ן 22 – בשלב השני תידרשו לפתור את הבעיה בעזרת חוקי הקשר וניתוח אשכולות. כמו כן, יהיה עליכם לסכם את עיקרי התוצאות שקיבלתם בממ"ן 21 וכן בממ"ן 22 והמסקנות.

על מנת לסייע לכם בפתרון הפרויקט מורכב הממ"ן ממספר שאלות הנחיה. **הינכם נדרשים לענות על כל השאלות לפי סדר הופעתן.** לצורך פתרון השאלות הינכם רשאים להניח כל הנחה ו/או

הפשטה סבירה שתידרשו לה. יש לציין במפורש בתחילת הפרויקט את ההנחות ו/או ההפשטות בהן הנכם משתמשים. כמו כן, יש לציין כל הנחה ו/או הפשטה עליה אתם מתבססים במהלך הפרויקט.

במסגרת הפרויקט יש להשתמש בחבילת תוכנה חופשית WEKA המצויה בכתובת :

WEKA: <http://www.cs.waikato.ac.nz/ml/weka/index.html>

### **צורת הגשה:**

יש להגיש את הפרויקט מודפס. הקפידו על כתיבה בהירה ומאורגנת וכן על תרשימים ברורים וקריאים. יש להקפיד על תיעוד מפורט של כל שלבי הפרויקט. **אין צורך** לצרף לפרויקט נתונים טכניים של חבילת התוכנה WEKA.

## 1. הגדרת הבעיה והכנת הנתונים (50%)

- א. (6%) הגדירו את הנתונים בהם השתמשתם בפרויקט כדוגמת: תכונות, סוג הנתונים, נתונים חסרים, תחומי ערכים ועוד.
- ב. (7%) הגדירו את מטרות כריית המידע. ציינו את ההנחות וההפשטות בהן השתמשתם.
- ג. (7%) בהמשך לסעיפים א ו-ב, הגדירו ותארו את שלבי ה-KDD עבור הבעיה הנתונה.
- ד. (15%) בהמשך לסעיפים א ו-ב ערכו סקירה השוואתית לכלל החלופות האפשריות (לפחות 4 חלופות) לביצוע כריית מידע. בתשובתכם יש להתייחס ליתרונות/חסרונות כל אחת מהחלופות בהקשר לבעיה הנתונה.
- ה. (15%) תארו את שלבי הכנת הנתונים. בתשובתכם יש להתייחס לבעיות באיכות הנתונים כולל טיפול בערכים חסרים, תצוגה גרפית של הנתונים, ניקוי הנתונים, שילוב והמרה של נתונים ועוד.

### הערה:

בשלב זה ניתן להשתמש בתוכנת Excel.

## 2. סיווג וחיזוי (50%)

- א. (5%) בחרו שתי שיטות לחיזוי הנתונים. הסבירו את השיטות ונמקו את בחירתכם.
- ב. (15%) תארו את שלבי השיטות שבחרתם בסעיף ב.
- ג. (8%) עבור כל שיטה דווחו את תוצאות הניתוחים.
- ד. (7%) העריכו את מידת הדיוק של כל שיטה.
- ה. (15%) נתחו השוואתית את התוצאות והסיקו מסקנות כולל הצעות לשיפורים.





# מטלת מנחה (ממ"ן) 12

הקורס: 20595 - כריית מידע

חומר הלימוד למטלה: יחידות 7-11

מספר השאלות: 3

משקל המטלה: 2 נקודות

סמסטר: 2015ב

מועד אחרון להגשה: 24.5.2015

אנא שים לב:

מלא בדייקנות את הטופס המלווה לממ"ן בהתאם לדוגמה שלפני המטלות.  
העתק את מספר הקורס ומספר המטלה הרשומים לעיל.

שאלה 1 (30%):

בטבלה הבאה נתונים נתוני תצפיות  $X_1, X_2, X_3$  המנבאים את  $Y$ .

$Y$	$X_3$	$X_2$	$X_1$
0	0	0	0
1	0	0	0
0	1	1	0
0	0	1	1
1	0	1	1
1	1	0	1
1	1	1	1

א. מה הן ההסתברויות החיוניות למסווג נאיב בייס.

ב. השתמשו בנאיב בייס לסיווג ערכו של  $Y$  עבור  $(X_1, X_2, X_3) = (0, 0, 0)$

**שאלה 2 (30%):**

נתונה הטבלה:

Items	TransID
A, B, C, D	T100
A, B, C, E	T200
A, B, E, F, H	T300
A, C, H	T400

הניחו:

Min\_support=50%

Min\_confidence=60%

חשבו את מדד הענין (lift) עבור כל הקבוצות התדירות.

**שאלה 3 (40%):**

בצעו אשכול לשבע הנקודות הבאות תוך שימוש באלגוריתם k-means.

A1(0, 0), A2(8,0), A3(8,6)

B1(16,0), B2(0,6)

C1(16,6), C2(16,0)

בתשובתכם הניחו:

א. יש לבצע אשכול לשלושה אשכולות בלבד.

ב. מרכזי האשכולות בשלב ההתחלתי הם A1, B1, C1

ג. שימוש במרחק אוקלידי (Euclidian distance) בריבוע.

**שימו לב,**

בתשובתכם הסופית יש להדגים את כל השלבים וכן להגדיר את האשכולות הסופיים.

# מטלת מנחה (ממ"ן) 22 - פרויקט גמר

הקורס: 20595 - כריית מידע

חומר הלימוד למטלה: יחידות 7-11

משקל המטלה: 13 נקודות

מספר השאלות: 3

מועד אחרון להגשה: 7.6.2015

סמסטר: 2015ב

אנא שים לב:

מלא בדייקנות את הטופס המלווה לממ"ן בהתאם לדוגמה שלפני המטלות.  
העתק את מספר הקורס ומספר המטלה הרשומים לעיל.

הנחיות:

בהמשך לממ"ן 21 השלימו במסגרת מטלה זו את הפרויקט בקורס. מומלץ לשוב ולעיין בהנחיות שניתנו בממ"ן 21.

ענו על השאלות הבאות:

## 1. חוקי הקשר (45%)

- (10%) א. הניחו מינימום confidence & support ובהתאמה, הגדירו חוקים המקיימים זאת. הסבירו ונמקו.
- (4%) ב. הגדירו את סט הפרמטרים. הסבירו ונמקו.
- (10%) ג. בחרו שני אלגוריתמים של חוקי הקשר. תארו את האלגוריתמים ונמקו את בחירתכם.
- (4%) ד. הריצו ודווחו את התוצאות של שני האלגוריתמים.
- (7%) ה. לאור סעיף ד הגדירו אילו מהחוקים שהגדרתם בסעיף א יש בהם עניין ועשויים להיות שימושיים.
- (10%) ו. נתחו השוואתית את התוצאות של שני האלגוריתמים והסיקו מסקנות.

## **2. ניתוח אשכולות (45%)**

- (2%) א. הגדירו מהו ניתוח אשכולות.
- (2%) ב. תארו את אופן מדידת איכות האשכולות.
- (12%) ג. בחרו שתי גישות לניתוח אשכולות.  
בתשובתכם יש לכלול הסבר אודות כל גישה וכן נימוק לבחירתכם.
- (12%) ד. תארו את שלבי ניתוח האשכולות עבור 2 הגישות שצינתם בסעיף ג.  
בתשובתכם יש להתייחס בין היתר לאופן הכנת הנתונים, מה הם הפרמטרים, ערכי הפרמטרים ועוד.
- (7%) ה. עבור כל גישה דווחו את תוצאות הניתוחים.
- (10%) ו. נתחו השוואתית את התוצאות והסיקו מסקנות.

## **3. סיכום ומסקנות (10%)**

סכמו בקצרה את עיקרי התוצאות שהתקבלו בממ"ן 21 וכן בסעיפים הקודמים בממ"ן הנוכחי ואת המסקנות שניתן לקבל מתוצאות אלה.