

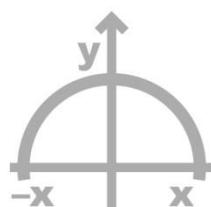
כלים כמותיים מתקדמים של תכנון סטטיסטי לאיכות



A square frame on a white background. Inside, there are four numbers: '1' at the top-left, 'sqrt(2)' at the top-right, '1' at the bottom-left, and '1' at the bottom-right. A diagonal line from the top-left corner to the bottom-right corner passes through the number 'sqrt(2)'.



A large white brace is placed over the mathematical expression $\{\sqrt{x}\}^2$ on an orange background.



תוכן העניינים

1	. מקדם המתאים (מזרק קשור) הליינארי וmobahkotno
24	. רגרסיה פשוטה
36	. רגרסיה מרובה
41	. רגרסיה - שאלות ממבחןים

כליים כמותיים מתקדמים של תכנון סטטיסטי לאיכות

פרק 1 - מקדם המתאים (מדד קשר) הלנארי וМОבהקותו

תוכן העניינים

1	1.	מדד המתאים הלנארי (פירסון)
12	2.	חישוב מקדם המתאים הלנארי (פירסון)
17	3.	בדיקה השערות על מקדם המתאים הלנארי
21	4.	בדיקה השערות על מקדם המתאים הלנארי באמצעות טבלה של ערכי קרייטיים

מקדם המתאים (מדד קשר) הלינארי וモבהקוטו

מדד הקשר הלינארי (פירסון) – מבוא

מעוניינים לבדוק עד כמה קיים קשר מסווג קשר לינארי (קו ישר) בין שני משתנים. שני המשתנים שאנו בודקים לגביים קשר צריכים להיות משתנים כמותיים. מבחינת סולמות מדידה כל משתנה נחקר צריך להיות מסולם רוחחים או מנה. בדרך כלל המשתנה המוצג כ- Y הוא המשתנה תלוי והמשנה המוצג כ- X הוא המשתנה הבלטי תלוי. תיאור גרפי לניטונים נעשה על ידי דיאגרמת פיזור. בדיאגרמת פיזור אנחנו מסמנים כל תצפית בנקודה לפי שיעור ה- X ושיעור ה- Y שלו. דיאגרמת הפיזור נותנת אינדיקציה גרפית על הקשר בין שני המשתנים.

דוגמה (פתרון בהקלטה) :

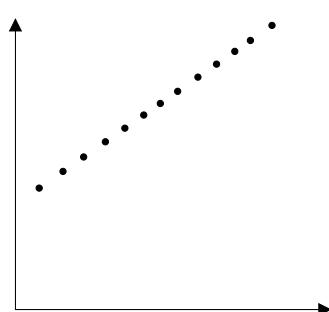
בבנייה 8 דירות בדקו לכל דירה את מספר החדרים שלה וכמו כן את מספר הנפשות הגרות בדירה. להלן התוצאות שהתקבלו :

מספר חדרים בדירה	מספר הנפשות בדירה
4	4
5	4

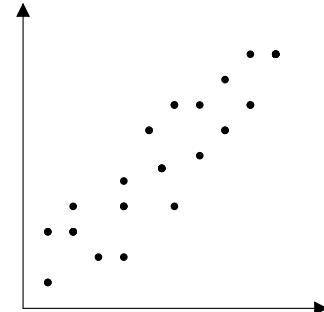
- 1) כמה תציפות ישן בדוגמה?
- 2) כמה משתנים ישנים בדוגמה, מי הם?
- 3) שרטטו לניטונים דיאגרמת פיזור.
- 4) מי המשתנה התלו依 ומיהו המשתנה הבלטי תלוי?

דיאגרמות פיזור לקשר בין משתנים וניתוחם

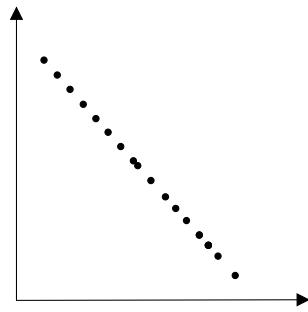
קשר לנארוי חיובי מלא



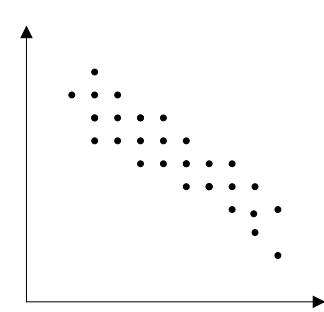
קשר לנארוי חיובי חלק



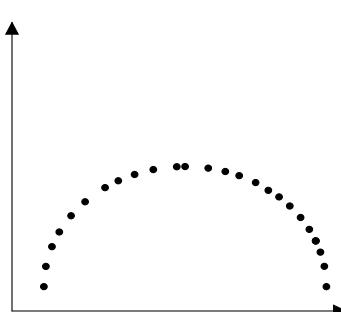
קשר לנארוי שלילי מלא



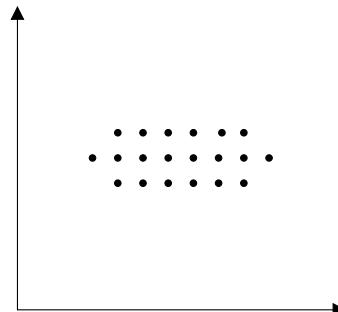
קשר לנארוי שלילי חלק



אין קשר לנארוי



אין קשר



משמעות מקדם המתאים:

כדי לבדוק עד כמה קיים קשר לנארוי בין שני המשתנים ישנו מדד קשר שנקרא גם מקדם המתאים הלינארי הידוע גם בשם מקדם המתאים של פירסון. מקדם מתאים זה מקבל ערכים בין 1 ל-1.

-1

0

1

מقدم מותאם 1-או 1 אומר שקיים קשר לינארי מלא בין המשתנים שנייתן לבטא על ידי נוסחה של קו ישר: $y = ax + b$.

מתאים חיובי מלא (מقدم מותאם 1):

קיים קשר לנארי מלא בו השיפוע a יהיה חיובי ואילו מותאם שלילי (מقدم מותאם-1) מלא אומר שקיים קשר לנארי מלא בו השיפוע a שלילי.

מתאים חיובי חלק:

כל משתנה אחד עולה לשני יש נטייה לעלות בערכו אבל לא קיימת נוסחה לינארית שמקשרת את X ל- Y באופן מוחלט ואילו מותאם שלילי חלקי אומר שככל המשתנה אחד עולה לשני יש נטייה לרדת אבל לא קיימת נוסחה לינארית שמקשרת את X ל- Y באופן מוחלט. ככל שמדובר המתאים הקרוב לאפס עצמת הקשר יותר חלשה וככל שהמדד רחוק יותר מהאפס העוצמה יותר חזקה. לsicום, מقدم המתאים בודק את עצמת הקשר הלינארי, ואת כיוון הקשר.

מقدم המתאים הלינארי אינו מושפע מייחדות המדידה. כל שינוי ביחסות המדידה של המשתנים, לא ישנה את מقدم המתאים.

מדד הקשר הלינארי באוכולוסייה, שנראה גם מقدم המתאים של פירסון או מדד הקשר של פירסון באוכולוסייה מסומן ב: r - פרמטר המאפיין את עצמת הקשר הלינארי באוכולוסייה וכיונו בין שני המשתנים הנחקרים. כאשר:

- מדד הקשר הלינארי במדגם שמהווה אומד לפרמטר r .

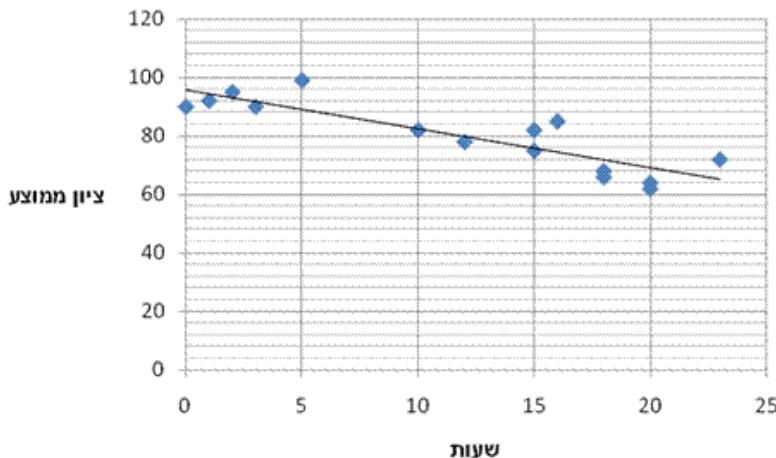
קיומו של מותאם בין שני משתנים אינו מצביע על סיבות בבחירה. למשל, אם נמצא מותאם חיובי בין כמות הסוכרזיות שאדם אוכל לבין משקל שלו אין זה אומר שהסיבה להשמנה היא הסוכרזית. מדד הקשר של פירסון הוא מדד קשר סימטרי,قولר אם נחליף את X ב- Y התוצאה תהיה זהה.

דוגמה (פתרון בהקלטה):

- מה ניתן להגיד על מועד המתאים של שני המשתנים על סמך דיאגרמת הפיזור שרטטנו?
- אם היינו משנים את הشرط כך שבציר האנכי היה המשתנה "מספר החדרים" ובציר האופקי היה "מספר הנפשות", האם הדבר היה משנה על מדד הקשר של פירסון?

שאלות

1) חוקר רצה לאפיין את הקשר בין מספר השעות בשבוע שסטודנט מקדיש לבילויים לבין הציון הממוצע שלו בסוף הסמסטר. לשם כך הוא אסף נתונים של 15 סטודנטים ויצר דיאגרמת פיזור:



- א. מיהו המשתנה הבלתי תלוי?
- ב. מה ניתן לומר על כיוון הקשר בין מספר שעות הבילוי השבועית לבין הציון הממוצע של הסמסטר? מה ניתן להגיד על עוצמת הקשר?

2) להלן טבלה המסכםת את מקדמי המתאים הליינארי בין ציוני מבחנים שונים שהתקבלו עבור תלמידים בכיתה מסוימת:

מתמטיקה	לשון	ספורט	ספורט
?	-0.7	?	ספורט
0.6	?	?	לשון
?	?	-0.1	מתמטיקה

א. השלימו את מקדמי המתאים שמשמעותם בסימן שאלה בטבלה.

ב. בין אילו שני ציוני מקצועות שונים קיים מתאם בעל העוצמה החזקה ביותר?

3) במחקר נתקשו לבדוק את הקשר בין מספר שעות התרגול של קורס לביון הציון הסופי שלו. להלן תוצאות מדגם שהתקבל:

שיעור תרגול	ציון סופי
90	20
90	25
95	30
60	15
90	30
85	20
50	10

- א. מיהו המשתנה התלו依 ומיהו המשתנה הבלתי תלוי בדוגמה זו?
- ב. שרטטו דיאגרמת פיזור לנוטונים.
- ג. מה ניתן לומר על הקשר בין המשתנים במדגם?
- ד. מסתבר שבסופו של דבר נתנו פקטור של 5 נקודות לציון הסופי. כיצד הדבר היה משנה את מקדם המתאים של המדגם?

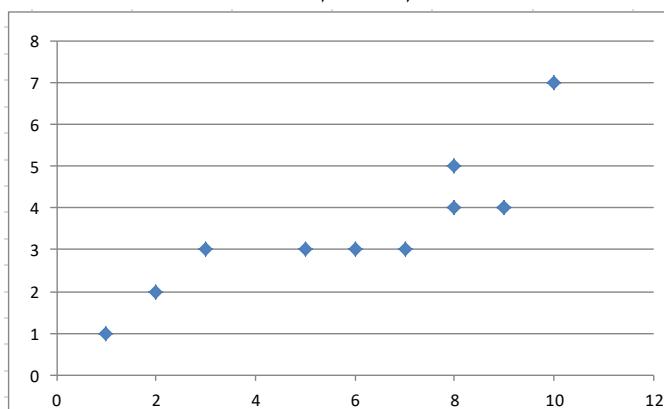
4) בتحقנה המטאורולוגית רצוי לבדוק את הקשר שבין הטמפרטורה במערכות כלזיות לכמות המשקעים במ"מ. הם אספו נתונים על 10 ימים במהלך חודש ינואר. המתאים שהתקבל היה 0.8.

א. השלימו את המשפט:

בחודש ינואר ככל שהטמפרטורה היומית נוטה לרדת, כך כמות המשקעים נוטה _____.

ב. הוחלט להעביר את הטמפרטורה למערכות פרנהייט על מנת שיוכלו להשוות אותה לנ נתונים מארה"ב. נוסחת המעבר היא $F^0 = 32 + \frac{9}{5}C^0$. כיצד הדבר ישפיע על מקדם המתאים בין הטמפרטורה במערכות פרניאיט לכמות המשקעים במ"מ?

5) להלן דיאגרמת פיזור המראה קשר בין שני משתנים:



א. השלימו: ניתן לראות קשר הוא לינארי _____ (מלאו חלקי) כיוון שהקשר הוא (חיובי ושלילי).

ב. השלימו: אם היינו מושפעים תצפית שערך ה- X שלה הוא 4 וערך ה- Y שלה הוא 7, מקדם המתאים של פירסון היה _____ (גדלו קטו לא משתנה).

שאלות רב ברירה (יש לבחור את התשובה הנכונה):

6) חוקר אקלים דגם כמה ימים בשנה ומדד את הטמפרטורה בטורונטו שבקנדזה ואת הטמפרטורה בסידני שבאוסטרליה באותו היום. הוא חישב ומצא מקדם מתאים שלילי בין הטמפרטורה היומית בטורונטו לבין הטמפרטורה היומית בסידני. משמעות מקדם המתאים השלילי בדגם:

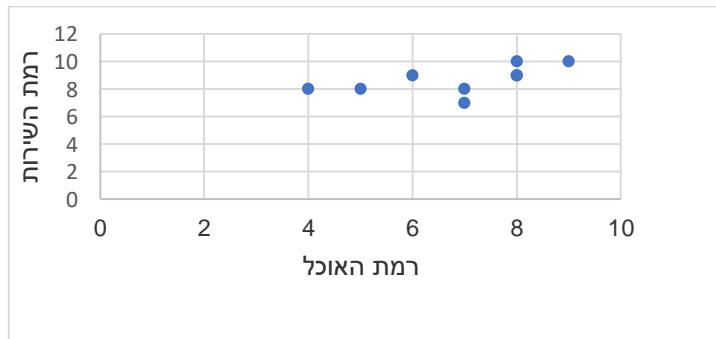
א. אין קשר בין הטמפרטורה בטורונטו לבין הטמפרטורה בסידני ביום שנדגמו.

ב. בדגם, רוב הטמפרטורות בטורונטו היו שליליות.

ג. ההפרש בין הטמפרטורה בטורונטו לבין הטמפרטורה באוסטרליה, בדגם זה, הוא שלילי.

ד. בדגם יש נטייה שהטמפרטורה יורדת בטורונטו לטמפרטורה לעלות בסידני.

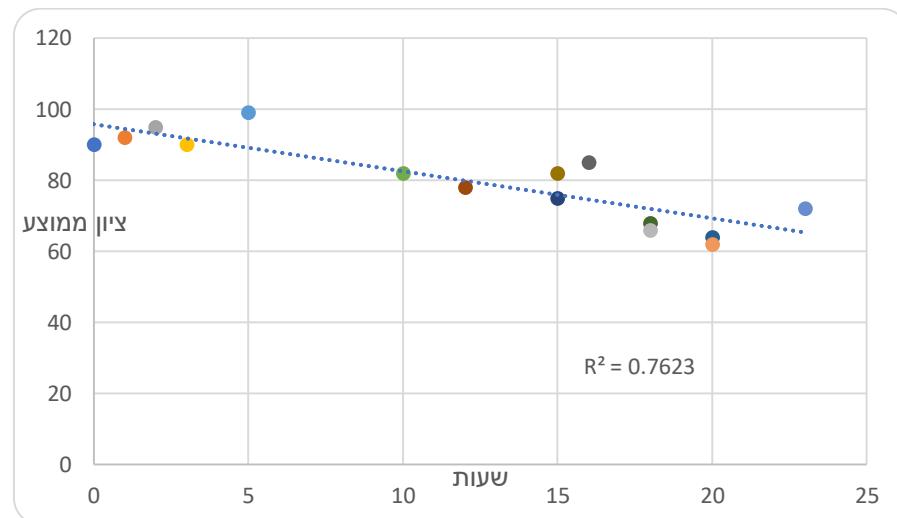
- 7) בסקר שביעות רצון שנערך בבית הקפה "fat لלחס" התבקוו הליקות לדרג את מידת שביעות הרצון שלהם (בסולם 1-10) בשני נושאים: רמת האוכל ורמת השירות.



מה יהיה ערכו של מקדם המתאים (r)?

- א. $r = -0.3$
- ב. $r = 0$
- ג. $r = 1.125$
- ד. $r = 0.593$

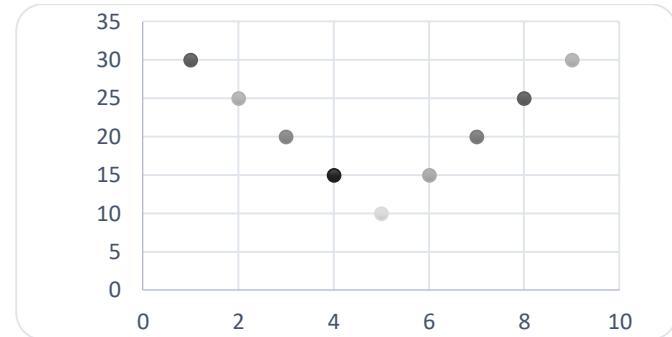
- 8) חוקר רצה לאפיין את הקשר בין מספר השעות בשבוע שסטודנט מקדיש לבילויים לבין הציון הממוצע שלו בסוף הסמסטר. לשם כך הוא אסף נתונים של 15 סטודנטים ויוצר דיאגרמת פיזור.



מה ניתן לומר על כיוון הקשר במדגם בין מספר שעות הבילוי השבועית לבין הציון הממוצע של הסמסטר?

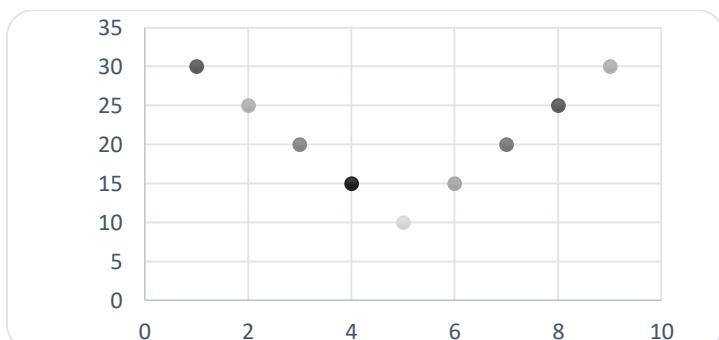
- א. ככל שմבלים יותר הציון נוטה לרדת.
- ב. אין קשר בין שעות הבילוי לציון.
- ג. ככל שմבלים פחות הציון נוטה לרדת.
- ד. ככל שהציון נוטה לרדת הסטודנט מבליה פחות.

9) התרשימים הבא מתאר קשר בין שני משתנים, איזה מהמתאים הבאים הוא המתאים ביותר לתיאור הקשר בין שני המשתנים?



- א. $1 = r$ היות ושני המשתנים יוצרים קוים ישרים.
- ב. $2 = r$ היות ויש שני קוים בעלי קשר מושלם.
- ג. $0 = r$ היות והקו יורד ולאחר מכן עולה באותו האופן.
- ד. $1 \pm 1 = r$ היות ויש קו עולה וגם קו יורד.

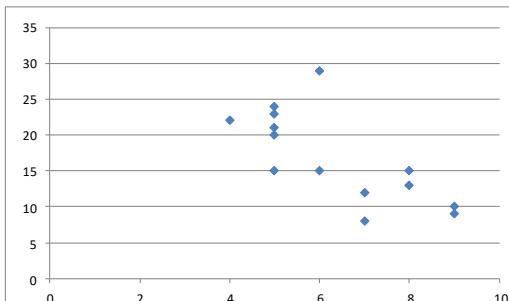
10) התרשימים הבא מתאר דיאגרמת פיזור.



איזה טענה נכונה?

- א. בתרשימים מוצג הקשר בין שני משתנים.
- ב. בתרשימים מוצג הקשר בין 9 משתנים.
- ג. בתרשימים מוצג הקשר בין 10 משתנים.
- ד. אין לדעת כמה משתנים מוצגים בתרשימים.

בגרף הבא מתוארת דיאגרמת פיזור של שני משתנים :



X - (משתנה בלתי תלוי בציר האופקי)
 ו- Y (משתנה תלוי).

במדגם התקבל $r^2 = 0.52$.

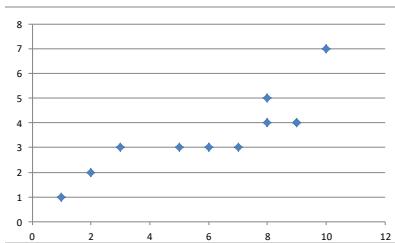
11) לאור הנתונים המופיעים בדיאגרמה, איזה מבחן הערכים הבאים מתאים להיות התוצאה של r ?

- א. -0.52
- ב. 0.72
- ג. -0.72
- ד. 0.52

12) אם מקדם המתאים בין שני משתנים הוא 1, אז :

- א. הערכים של המשתנים הם חיוביים.
- ב. עברו כל תצפית ערך של משתנה אחד שווה לערך של המשתנה השני.
- ג. הקשרlienاري הוא בעוצמה חזקה.
- ד. אף אחת מהתשובות לא בהכרח נכונה.

13) להלן דיאגרמת פיזור :
 מה יהיה מקדם המתאים בין שני המשתנים ?



- א. 1
- ב. 0.85
- ג. 0.15
- ד. 0

14) בבדיקה קשר בין שני משתנים התקבל : $-1 = r$.

א. קיימת נוסחהlienארית הקושרת בין כל התצפיות.

ב. לא קיים קשר בין שני המשתנים.

ג. ככל משתנה אחד נוטה לרדת גם לשני יש נטייה לרדת.

ד. קיים קשר בין שני המשתנים, אך לא ניתן לדעת מאיזה סוג.

15) לפי הפטגס "רחוק מהעיר, רחוק מהלב", יש קשר ____ בין קרבה פיזית לקרבה נפשית.

- א. חיובי
- ב. שלילי
- ג. אפסי
- ד. לא ניתן לדעת.

16) מבחן אמייר הינו מבחן מיוון באנגלית של המרכז הארצי לבחינות והערכתה. הציון המינימלי בבחינה הינו 150 והמаксימלי הינו 250. בקורס הכנה לבחן השתתפו 19 תלמידים. להלן הציונים שלהם על פי פلت שהתקבל:

	159
	170
	180
	185
	204
	224
	236
	212
	168
	189
	195
	163
	187
	206
	201
	223
	242
	203
	205
197.47	AVERAGE
536.25	VARPA

יש להוסיף עמודה נוספת לצד עמודות הציונים שטראה לכל תלמיד כמה נקודות חסרות לו כדי להשלים לציוון המקסימלי בבחינה.

מה יהיה מקדם המתאים בין שתי העמודות (תלמיד, מקדם המתאים בין הציון לבין הנקודות החסרות)?

- א. -1
- ב. 1
- ג. -0.5
- ד. 0.5

17) מקדם המתאים בין שטחי דירה למחיר שלהם חושב ונמצא 1.2. מה נובע לכך?

- א. ככל שהדירה גדולה יותר בשטחה כך היא יקרה יותר.
- ב. ככל שהדירה קטנה יותר בשטחה כך היא זולה יותר.
- ג. לא קיים קשר בין שטח הדירה למחיר הדירה.
- ד. מצב כזה שמתואר הנתונים לא אפשרי.

18) אם ניקח 10 אנשים וונרשום לכל אדם את הגובה במטר וכמה כו' את הגובה בס"מ. מה יהיה מקדם המתאים בין גובה האדם במטר לגובה האדם בס"מ?

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

- 19)** נמצא מתאים חיובי בעוצמה גבוהה בין X – ציון בගראות בלשון ל Y – ציון בගראות במתמטיקה. אילו מהמשפטים הבאים נכון?
- ניתן לומר שאחת מהסיבות להבדלים שיש לסטודנטים במתמטיקה נובעים מההבדלים שיש להם בלשון.
 - קיימת נוסחה של קו ישר שקשורה בין ציון בගראות במתמטיקה לציון בගראות בלשון.
 - לא יוצא מן הכלל, ניתן להגיד שככל תלמיד שמציל יותר מטלמיד אחר בלשונו גם יצליח יותר מאותו תלמיד במתמטיקה.
 - אף אחד מהטענות שהוצעו אינה בהכרח נכונה.
- 20)** עברו סדרה של תצפיות מדדו את X ואת Y . נמצא שעבור כל התצפיות שהערך של Y ירד הערך של X בהכרח ירד ללא יוצא מן הכלל. מקדם המתאים של פירסון יהיה בהכרח :
- 1
 - 1
 - 0
 - אף אחת מהתשובות.

תשובות סופיות

- ב. הקשר חלקי, כיון הקשר שלילי.
ב. ספורט ולשון.

- (1) א. שעות בילוי.
(2) א. להלן טבלה:

מתמטיקה	לשון	ספורט	ספורט
ספורט	0.1	-0.7	1
לשון	0.6	1	-0.7
מתמטיקה	1	0.6	-0.1

- ב. ראה גרפ' בפתרון וידאו.
ד. מקדם המתאים לא היה משתנה.
ב. לא ישפיע על מקדם המתאים.
ב. קטן.

- (3) א. ב'ית- מס' שעות התרגול, תלוי- ציון.
ג. קשר לינארי חיובי חלקי.

- (4) א. עלות.
(5) א. חלקי, חיובי.
(6) ד'. ד'. א'. א'. ג'. (10)
(11) ג'. א'. א'. א'. (14) (15) (12) (13) ב'. ד'. ד'. (16) א'. א'. א'. (17) (18) א'. ד'. ד'. (20) (19)

מדדי קשר – מדד הקשרlienاري (פירסון) – רקע

המטרה היא לבדוק האם קיים קשר (קורלציה, מתאים) של קו ישר בין שני משתנים כמותיים. מבחינת סולמות המודיעה קשר בין סולמות רוחניים ומנה. בדרך כלל, X הוא המשתנה המסביר (הבלתי תלוי) ו- Y הוא המשתנה המוסבר (התלויה).

דוגמה:

נרצה להסביר כיצד השכלה של אדם הנ마다 בשנות לימוד – X מסביר את ההכנסה שלו Y . במקרה זה שנות ההשכלה זהו המשתנה המסביר (או הבלתי תלוי) ואנחנו מעוניינים לבדוק כיצד שינויים בשנות ההשכלה של אדם יכולים להשיבר את השינויים שלו בהכנסה, וכך רמת ההכנסה זהו המסביר התלויה במשתנה המסביר אותו.

שלב ראשון: נהוג לשרטט דיאגרמת פיזור. זו דיאגרמה שנוננת אינדיקטיבית ויזואלית על טיב הקשר בין שני המשתנים.

דוגמה:

מספר דירה	X	Y
1	3	2
2	2	2
3	4	3
4	3	3
5	5	4

בבנייה של 5 דירות בדקנו את הנתונים הבאים :

X - מספר חדרים בדירה. Y - מס' נפשות הגרות בדירה.

להלן התוצאות שהתקבלו :

נשרטט מנתונים אלה דיאגרמת פיזור (הDİAGRAM המלאה בסרטון). נתבונן בכמה מקרים של דיאגרמות פיזור וננתח אותן (הDİAGRAMS המלאות בסרטון).

שלב שני: מחשבים את מקדם המתאים (מדד הקשר) שבזוק עד כמה קיים קשרlienاري בין שני המשתנים. המדד (נקרא גם מדד הקשר של פירסון) מכמת את מה שנראה בשלב הראשון רק בעין.

המדד בודק את כיוון הקשר (חיובי או שלילי) ואת עוצמת הקשר (חלש עד חזק).

מקדם מתאים זה מקבל ערכאים בין 1- -1.

מקדם מתאים 1- או 1 אומר שקיים קשרlienاري מוחלט ומלא בין המשתנים שניינו לבטא על ידי הנוסחה : $y = bx + a$.

מתאים חיובי מלא (מקדם מתאים 1):

קיים קשר לנארוי מלא בו השיפוע b יהיה חיובי ואילו מתאים שלילי מלא אומר שקיים קשר לנארוי מלא בו השיפוע b שלילי (מקדם מתאים 1-).

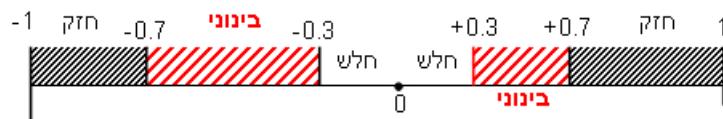
מתאים חיובי חלקי:

ככל שמשתנה אחד עולה לשני יש נטייה לעלות בערכו אבל לא קיימת נוסחה לינארית שמקשרת את X ל- Y באופן מוחלט.

מתאים שלילי חלקי:

ככל שמשתנה אחד עולה לשני יש נטייה לרדת אבל לא קיימת נוסחה לינארית שמקשרת את X ל- Y באופן מוחלט.

ככל שערך מקדם המתאים קרוב לאפס נאמר שעוצמת הקשר חלה יותר וככל שמקדם המתאים רחוק מהאפס נאמר שעוצמת הקשר חזקה יותר :



מקדם המתאים יסומן באות r .

כדי לחשב את מקדם המתאים, יש לחשב את סטיות התקן של כל משתנה ואת השונות המשותפת.

$$COV(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n} = \frac{\sum xy}{n} - \bar{x} \cdot \bar{y}$$

$$s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} = \frac{\sum_{i=1}^n x_i^2}{n} - \bar{x}^2 : \text{שונות של המשתנה } X$$

$$S_y^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n} = \frac{\sum_{i=1}^n y_i^2}{n} - \bar{y}^2 : \text{שונות המשתנה } Y$$

$$\text{מקדם המתאים הלינארי : } r_{xy} = \frac{COV(x, y)}{S_x \cdot S_y}$$

שאלות

1) להלן נתונים לגבי שישה תלמידים שנגשו ל מבחון. בדקו לגבי כל תלמיד את הציון שלו בסוף הקורס וכמו כן את מספר החיסורים שלו מהקורס.

מספר חיסורים	ציון
4	70
3	70
2	90
0	90
1	90
2	80

א. שרטטו דיאגרמת פיזור לנ נתונים. מה ניתן להסיק מהדיאגרמה על טיב הקשר בין מספר החיסורים של תלמיד לציונו? מיהו המשתנה הבלטי תלוי ומיהו המשתנה התלווי?

ב. חשבו את ממד הקשר של פירסון. האם התוצאה מתוישבת עם תשובה לסעיף א'?

ג. הסבירו, ללא חישוב, כיצד מקדם המתאים היה משתנה אם היה מתווסף תלמיד שהיחסיר 4 פעמים וקיבל ציון 80?

X	Y
10	12
14	15
15	15
18	17
20	21

2) במחקר רפואי רצוי לבדוק האם קיים קשר בין רמת ההורמון X בدم החולים לרמת ההורמון Y שלו. לצורך כך מדדו את רמת ההורמוניים ההלו עבור חמישה חולים. להלן התוצאות שהתקבלו:

א. מה הממוצע של כל רמת ההורמו?

ב. מהו מקדם המתאים בין ההורמוניים? ומה המשמעות ההתואמת?

3) נסמן ב- X את ההכנסה של משפחה באלפי ש. נסמן ב- Y את ההוצאות של משפחה באלפי ש. נלקחו 20 משפחות והתקבלו התוצאות הבאות:

$$\sum_{i=1}^{20} Y_i = 200 \quad \sum_{i=1}^{20} X_i = 240$$

$$\sum_{i=1}^{20} (Y_i - \bar{Y})^2 = 76 \quad \sum_{i=1}^{20} (X_i - \bar{X})^2 = 76$$

$$\sum_{i=1}^{20} (X_i - \bar{X})(Y - \bar{Y}) = 60.8$$

א. חשב את ממד הקשר הליינארי בין X ל- Y. מיהו המשתנה התלווי?

ב. מה המשמעות של התוצאה שקיבלת בסעיף א'?

4) נסמן ב- X את ההכנסה של משפחה באלפי נק. נסמן ב- Y את ההוצאות של משפחה באלפי נק. נלקחו 20 משפחות והתקבלו התוצאות הבאות:

$$\sum_{i=1}^{20} Y_i = 200 \quad \sum_{i=1}^{20} X_i = 240$$

$$\sum_{i=1}^{20} Y_i^2 = 2080 \quad \sum_{i=1}^{20} X_i^2 = 2960$$

$$\sum_{i=1}^{20} X_i Y_i = 2464$$

חשבו את ממד הקשרlienاري בין X ל- Y .

5) במוסד אקדמי ציון ההתאמה מחושב כך: מכפילים את הציון הממוצע בוגרות ב-3 ומחיתנים 2 נקודות. ידוע שעבור 40 מועמדים סטיטית התקן של ממוצע הציון בוגרות הייתה 2. מה מגדם המתאים בין ציון ההתאמה לציון הממוצע בוגרות שלהם?

- 6)
- הלו רשימה טענות, לגבי כל טענה קבעו נכון/לא נכון ונמקו.
 - א. מתוויך דירות המיר מחירי דירות מדולר לשקל. נניח שдолר אחד הוא 3.5 נק. אם מתוויך הדירות יחשב את ממד הקשר של פירסון בין מחיר הדירה בשקלים למחיר הדירה בדולרים הוא יקבל 1.
 - ב. לסדרה של נתונים התקבל $S_x = S_y = 1$, $\bar{X} = \bar{Y}$. לכן, ממד הקשר של פירסון יהיה 1.
 - ג. אם השונות המשותפת של X ושל Y הינה 0 אז בהכרח גם מגדם המתאים של פירסון יהיה 0.

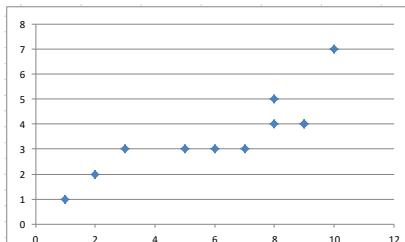
שאלות רב-ברירה:

- 7) נמצא שקיים מגדם מתאים שלילי בין הציון בעברית לחשבון בבחינה לכן:
- א. הדבר מעיד שהציונים בכתה היו שליליים.
 - ב. ככל שהציון של תלמיד יורך בחשבון יש לו נטייה לרדת בעברית.
 - ג. ככל שהציון של תלמיד עולה בחשבון יש לו נטייה לרדת בעברית.
 - ד. אף אחת מהתשובות לא נכונה.

8) נלקחו 20 מוצרים ונבדק ביום מסוים המחיר שלהם בדולרים והמחיר שלהם בש"ח (באותו היום ערך הדולר היה-2.4₪). מהו מקדם המתאים בין המחיר בדולר למחיר בש"ח?

- א. 1
- ב. 0
- ג. 4.2
- ד. לא ניתן לדעת.

9) להלן דיאגרמת פיזור:
מה יהיה מקדם המתאים בין שני המשתנים?



- א. 1
- ב. 0.85
- ג. 0.15
- ד. 0

תשובות סופיות

- 1) א. משתנה תלוי : ציון, משתנה ב"ת : מס' חיסורים. ראה דיאגרמה בוידאו. ניתן להסיק שקיים קשר לינארי שלילי וחליqi בין מספר החיסורים לציון התלמיד.
 ב. $r_{xy} = -0.9325$.
 ג. הקשר יישאר לינארי שלילי חליqi אך עוצמתו תחלש.
- 2) א. $r_{xy} = 0.96$ ב. $\bar{x} = 15.4$, $\bar{y} = 16$.
 3) א. 0.8
 4) ב. 0.8
 5) ג. 1.
 6) א. נכון.
 ב. לא נכון.
 ג. נכון.
 7) א'. ג'.
 8) א'.
 9) ב'.

בדיקות השערות על מקדם המתאיםlienاري – רקע

מדד הקשרlienاري באוכולוסייה, שנראה גם מקדם המתאים של פירסון או מדד הקשר של פירסון באוכולוסייה מסומן ב: r - פרמטר המאפיין את עצמת הקשרlienاري וכיונו בין שני המשתנים הנחקרים באוכולוסייה. כאשר:

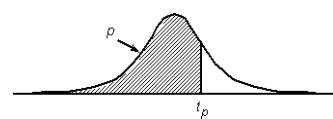
- מדד הקשרlienاري במדגם שמהווה אומד לפרמטר r .

השערת האפס: תהיה שבאוכולוסייה לא קיים כלל קשרlienاري בין שני המשתנים $H_0: \rho = 0$.
ההנחה שעלייה אנו מtabסים בתחילת היא שני המשתנים הנחקרים מתפלגים דו נורמלית.

$$\text{סטטיסטי המבחן: } t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t(n-2)$$

סטטיסטי זה מתפלג t עם $n-2$ דרגות חופש.

$H_0: \rho = 0$	$H_0: \rho = 0$	$H_0: \rho = 0$	השערת האפס :
$H_1: \rho > 0$	$H_1: \rho < 0$	$H_1: \rho \neq 0$	השערת המבחן :
$t \geq t_{1-\alpha}$	$t \leq -t_{1-\alpha}$	$t \geq t_{1-\alpha}$ או $t \leq -t_{1-\alpha}$	כל ההכרעה: אזור דחייה של השערת האפס

טבלת ערכים קרייטיים של ζ - נספח: טבלת התפלגות T
P

דרגות חופש	0.75	0.90	0.95	0.975	0.99	0.995	0.9995
1	1.000	3.078	6.314	12.709	31.821	63.657	636.619
2	0.816	1.886	2.920	4.303	6.965	9.925	31.598
3	0.765	1.638	2.353	3.182	4.541	5.841	12.941
4	0.741	1.533	2.132	2.776	3.747	4.604	8.610
5	0.727	1.476	2.015	2.571	3.365	4.032	6.859
6	0.718	1.440	1.943	2.447	3.143	3.707	5.959
7	0.711	1.415	1.895	2.365	2.998	3.499	5.405
8	0.706	1.397	1.860	2.306	2.896	3.355	5.041
9	0.703	1.383	1.833	2.262	2.821	3.250	4.781
10	0.700	1.372	1.812	2.228	2.764	3.169	4.587
11	0.697	1.363	1.796	2.201	2.718	3.106	4.437
12	0.695	1.356	1.782	2.179	2.681	3.055	4.318
13	0.694	1.350	1.771	2.160	2.650	3.012	4.221
14	0.692	1.345	1.761	2.145	2.624	2.977	4.140
15	0.691	1.341	1.753	2.131	2.602	2.947	4.073
16	0.690	1.337	1.746	2.120	2.583	2.921	4.015
17	0.689	1.333	1.740	2.110	2.567	2.898	3.965
18	0.688	1.330	1.734	2.101	2.552	2.878	3.922
19	0.688	1.328	1.729	2.093	2.539	2.861	3.883
20	0.687	1.325	1.725	2.086	2.528	2.845	3.850
21	0.686	1.323	1.721	2.080	2.518	2.831	3.819
22	0.686	1.321	1.717	2.074	2.508	2.819	3.792
23	0.685	1.319	1.714	2.069	2.500	2.807	3.767
24	0.685	1.318	1.711	2.064	2.492	2.797	3.745
25	0.684	1.316	1.708	2.060	2.485	2.787	3.725
26	0.684	1.315	1.706	2.056	2.479	2.779	3.707
27	0.684	1.314	1.703	2.052	2.473	2.771	3.690
28	0.683	1.313	1.701	2.048	2.467	2.763	3.674
29	0.683	1.311	1.699	2.045	2.462	2.756	3.659
30	0.683	1.310	1.697	2.042	2.457	2.750	3.646
40	0.681	1.303	1.684	2.021	2.423	2.704	3.551
60	0.679	1.296	1.671	2.000	2.390	2.660	3.460
120	0.677	1.289	1.658	1.980	2.358	2.617	3.373
∞	0.674	1.282	1.645	1.960	2.326	2.576	3.291

שאלות

1) להלן נתונים על הוווטק בעבודה (בשנים) ועל השכלה (בשנים) במדגם של 10 עובדים :

10	9	8	7	6	5	4	3	2	1	נבדק
24	17	28	5	9	16	8	2	18	13	X - הווטק
15	12	8	13	12	11	8	17	14	12	Y - השכלה

מقدم המתאים חושב והתקבל : 0.31 --.

א. האם קיימים מתאים בין וווטק העובד להשכלה? בדקו ברמת מובהקות של 5%?

ב. אם הווטק של העובד היה נמדד בחודשים האם התשובה לסעיף א' הייתה משתנה?

2) מחקר התעניין לבדוק את הקשר בין גיל נשים בהריאן לרמת ההמוגולובי שלחן בדם בזמן הריאן. נדרגו 7 נשים והתקבלו התוצאות הבאות :

גיל	1	2	3	4	5	6	7	נבדק
המוגולובי	14.7	13.5	9.7	12	10.8	13	10.3	
גיל	39	34	30	29	28	26	23	

במדגם חושב מדד הקשר של פירסון להיות 0.7.

א. האם ניתן לומר שבמדגם אם איש היא יותר מבוגרת אזי בהכרח יש לה יותר המוגולובי בדם?

ב. האם ניתן לומר, ברמת מובהקות של 5%, שקיים מתאם בין גיל האישה שהריאן לבין רמת ההמוגולובי שלה בדם?

3) בתחנה המטאורולוגית רצוי לבדוק את הקשר שבין הטמפרטורה במעלות צלזיות לכמות המשקעים במ"מ. הם אספו נתונים על 10 ימים במהלך חודש ינואר. המתאים שהתקבל היה 0.8.-.

א. בדקו ברמת מובהקות של 2.5% האם קיים קשר לינארי שלילי בחודש ינואר בין הטמפרטורה במעלות צלזיות לבין המשקעים במעלות צלזיות.

ב. כיצד הייתה המשתנה התשובה לסעיף א' אם היו מוסיפים עוד תצפיות למדגם?

ג. על סמך טבלת D המצורפת עבור אילו רמות מובהקות ניתן להחליט שקיים קשר לינארי שלילי מובהק?

4) מtower דירות חישב את מועד המתאים בין שטח דירה במרכז תל אביב לבין המחיר של הדירה עבור 17 דירות. מועד המתאים שקיבל היה 0.6.

א. בדקו ברמת מובהקות של 5% האם ניתן להגיד שקיים קשר ישר עולה בין שטח הדירה לבין מחיר הדירה במרכז תל אביב?

ב. מהי מובהקות התוצאה לבדיקת השערת שקיים קשר ישר עולה בין שטח הדירה לבין מחיר הדירה בתל אביב.

תשובות סופיות

- ב. לא תשתנה. 1) א. לא נדחה את H_0 .
- ב. לא נדחה את H_0 . 2) א. לא
- ב. לא ניתן לדעת. 3) א. נדחה את H_0 .
- ב. $0.005 < P_v < 0.01$. 4) א. נדחה את H_0 .

בדיקת השערות על מقدم המתאםlienاري (באמצעות טבלה של ערכים קרייטיים) – רקע

מדד הקשרlienاري באוכולוסייה, שנראה גם מقدم המתאם של פירסון או מדד הקשר של פירסון באוכולוסייה מסומן ב: r - פרמטר המאפיין את עצמת הקשרlienاري וכיונו בין שני המשתנים הנחקרים באוכולוסייה. כאשר :

- מדד הקשרlienاري במדגם שמהווה אומד לפרמטר r .

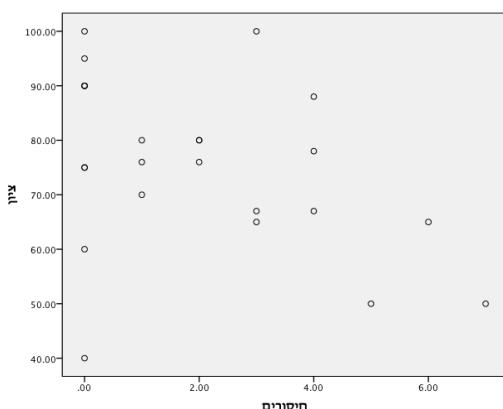
השערת האפס : תהיה שבאוכולוסייה לא קיים כלל קשרlienاري בין שני המשתנים: $H_0: \rho = 0$.
הנחה שעלייה אלו מתבססים בתהlik היא שני המשתנים הנחקרים מתפלגים דו-נורמלית.

את מقدم המתאם הקרייטי, שנסמנו ב- r_c , נוציא מתוך טבלה של ערכים קרייטיים שמצוירת בהמשך.

$H_0: \rho = 0$	$H_0: \rho = 0$	$H_0: \rho = 0$	השערת האפס:
$H_1: \rho > 0$	$H_1: \rho < 0$	$H_1: \rho \neq 0$	השערת המחקר:
$r \geq r_c$	$r \leq -r_c$	$r \geq r_c$ או $r \leq -r_c$	כל הנסיבות: אזרז דחיה של השערת האפס

דוגמה (פתרון בהקלטה) :

הזיכון ביקש לדגום סטודנטים כדי לבדוק את הקשר בין ציון הסטודנט בקורס למספר הפעמים שהוא החסיר שיעור בקורס.
דיאגרמת הפיזור שהתקבל במדגם שבוצע :



מייהו המשתנה תלוי ומיהו המשתנה הבלתי תלוי במחקר?
מה ניתן לראות לגבי הקשרlienاري בין המשתנים שהתקבל במדגם?

חוושב האומד למقدم המתאםlienاري על סמך 24 הסטודנטים שנדגמו והתקבל: -0.389 .

מה משמעות של מقدم המתאם שהתקבל במדגם?

אם ניתן להגיד ברמת מובהקות של 5% שקיים מתאםlienاري שלילי בין מספר החיסורים של הסטודנטים מהקורס לבין הציון של הסטודנטים בקורס?

טבלת ערכים קרייטיים של מקדם המתאם הלינארי



0.0005	0.005	0.025	0.05	α
n				
0.999	0.990	0.950	0.900	4
0.991	0.959	0.878	0.805	5
0.974	0.917	0.811	0.729	6
0.951	0.875	0.754	0.669	7
0.925	0.834	0.707	0.621	8
0.898	0.798	0.666	0.582	9
0.872	0.765	0.632	0.549	10
0.847	0.735	0.602	0.521	11
0.823	0.708	0.576	0.497	12
0.801	0.684	0.553	0.476	13
0.780	0.661	0.532	0.458	14
0.760	0.641	0.514	0.441	15
0.742	0.623	0.497	0.426	16
0.725	0.606	0.482	0.412	17
0.708	0.590	0.468	0.400	18
0.693	0.575	0.456	0.389	19
0.679	0.561	0.444	0.378	20
0.665	0.549	0.433	0.369	21
0.652	0.537	0.423	0.360	22
0.640	0.526	0.413	0.352	23
0.629	0.515	0.404	0.344	24
0.618	0.505	0.396	0.337	25
0.607	0.496	0.388	0.330	26
0.597	0.487	0.381	0.323	27
0.588	0.479	0.374	0.317	28
0.579	0.471	0.367	0.311	29
0.570	0.463	0.361	0.306	30
0.532	0.430	0.334	0.283	35

שאלות

1) להלן נתונים על הוווטק בעבודה (בשנים) ועל השכלה (בשנים) במדגם של 10 עובדים :

10	9	8	7	6	5	4	3	2	1	נחקר
24	17	28	5	9	16	8	2	18	13	X - וווטק
15	12	8	13	12	11	8	17	14	12	Y - השכלה

מדם המתאים חושב והתקבל : -0.31 .

א. האם קיימים מתאים בין וווטק העובד להשכלה? בדקו ברמת מובהקות של 5%.

ב. אם הווטק של העובד היה נמדד בחודשים האם התשובה לסעיף א' הייתה משתנה?

2) מחקר התענין לבדוק את הקשר בין גיל נשים בהריון לרמת המוגולובין שלהן בדם בזמן הריאון. נדגמו 7 נשים והתקבלו התוצאות הבאות :

נחקרת	1	2	3	4	5	6	7
המוגולובי	14.7	13.5	9.7	12	10.8	13	10.3
גיל	39	34	30	29	28	26	23

במדגם חושב מדד הקשר של פירסון להיות 0.7 .

א. האם ניתן לומר שבמדגם אם אישת היא יותר מבוגרת אזי היא בהכרח יש לה יותר המוגולובי בדם?

ב. האם ניתן לומר, ברמת מובהקות של 5%, שהמתאים בין גיל האישה שהריון לבין רמת המוגולובי שלה בדם הוא חיובי?

3) בתחנה המטאורולוגית רצוי לבדוק את הקשר שבין הטמפרטורה במעלות צלזיאוס לכמות המשקעים במ"מ. הם אספו נתונים על 10 ימים במהלך חודש ינואר. המתאים שהתקבל היה -0.8 .

א. בדוק ברמת מובהקות של 2.5% האם קיים קשרلينי שילילי בחודש ינואר בין הטמפרטורה במעלות צלזיאוס לבין המשקעים במ"מ?

ב. כיצד הייתה משתנה התשובה לסעיף א' אם היו מוסףים עוד תצפיות למדגם?

תשובות סופיות

- (1) א. לא נדחה את H_0 .
 ב. לא משתנה.
- (2) א. לא.
 ב. נדחה את H_0 .
- (3) א. נדחה את H_0 .
 ב. לא ניתן לדעת.

כליים כמותיים מתקדמים של תכנון סטטיסטי לאיכות

פרק 2 - רגרסיה פשוטה

תוכן העניינים

- 24 1. כללי

רגرسיה פשוטה:

רקע:

רגרסיה ליניארית פשוטה מסתמכת על המתאם הליניארי בין המשתנה הבלתי (המנובא) לב'ית (המנבא).

$$r = \frac{\text{cov}(x, y)}{S_x \cdot S_y} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right)} \cdot \sqrt{\left(\sum_{i=1}^n y_i^2 - n \bar{y}^2 \right)}} = \frac{SXY}{\sqrt{SXX} \cdot \sqrt{SYY}}$$

מקדם המתאים :

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

כאשר :

β_0 הוא החותך.

β_1 הוא שיפוע.

ε_i הינו גורם הטעות מסביב לקו הליניארי.

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

לסיכום :

1. במודל $y_i = \alpha + \beta x_i + \varepsilon_i$, α ו- β הם מספרים קבועים אך לא ידועים.

אנו יכולים להעריך אותם ולקבל אומדיים (תהליך קבלת האומדיים נקרא אמידה).

2. $\hat{\alpha}$ הוא האומד ל- α . $\hat{\beta}$ הוא האומד ל- β .

3. אומדי ריבועים פחותים (אר"פ) הם אומדיים שchosבו בשיטת הריבועים הפחותים. אומדי הריבועים הפחותים מסוימים בד"כ ע"י 'קובע' - $\hat{\alpha}$, $\hat{\beta}$.

4. בעוד α ו- β הם קבועים, $\hat{\alpha}$ ו- $\hat{\beta}$ הם משתנים מקריים. מדוע? מפני שבכל מודגש מתקבלים $\hat{\alpha}$ ו- $\hat{\beta}$ אחרים.

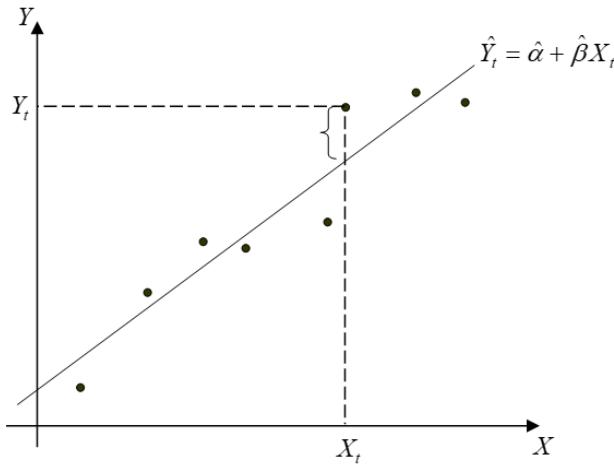
5. את α ו- β אי אפשר לדעת, ולכן אין אפשרות לדעת מהו הקו האמתי, וכן אין אפשרות לדעת את ε .

6. אפשר לדעת את e , שהיא הסטייה מקו הרגרסיה. נגיד זאת באופן הבא :

- עבור X_i , הערך הצפוי של המשתנה המוסף (\hat{Y}_i) המתקבל לפי הרגרסיה

$$\text{הוא : } \hat{Y}_t = \hat{\alpha} + \hat{\beta} X_t$$

- הסטייה של התצפית (ε_i) מהערך הצפוי לפי הרגרסיה (\hat{Y}_i) היא :



האומדיים של הרגרסיה $(\hat{\alpha}, \hat{\beta})$:

שיטת האמידה של α ושל β נקראת שיטת הריבועים הפחותים Ordinary Least Squares (OLS).

השאלה הנשאלת בשיטת אמידה זו היא:

אייזה $\hat{\alpha}$ ו- $\hat{\beta}$ יביאו למינימום את סכום ריבועי טעויות האמידה.

$\min_{\hat{\alpha}\hat{\beta}} \sum e_t^2 = \min_{\hat{\alpha}\hat{\beta}} \sum (y_t - \hat{y}_t)^2 = \min_{\hat{\alpha}\hat{\beta}} \sum [y_t - (\hat{\alpha} + \hat{\beta}x_t)]^2 = ?$

ובתרגם מתמטי:

מתוך גזירת הפונקציה הזו מתקבלים האומדיים $\hat{\alpha}$ ו- $\hat{\beta}$:

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{SXY}{SXX} = \frac{COV(X, Y)}{V(X)} = r \frac{S_y}{S_x}$$

מבחני המובהקות :

$$H_0 : \beta = 0$$

השערות :

$$H_1 : \beta \neq 0$$

ברגרסיה פשוטה בה יש לנו רק מבוא אחד: ניתן לבצע מבחן F ל모בהקות משווהות

הרגרסיה או מבחן T לモבהקות מקדים הרגרסיה (הביתא).

משמעות דחיתת השערת האפס: משווהות הרגרסיה מובהקת, מקדים הרגרסיה

מוגהבק, הקשר בין X ל- Y מוגהבק.

ולහיפך – אם השערת האפס לא נדחתת: אין הוכחה לכך בין המשתנים X ו- Y ,

משווהות הרגרסיה אינה מוגהבקת וכך גם מקדים הרגרסיה.

$$\hat{\sigma}^2 = MSE = \frac{SSE}{n-2} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2} = \frac{(1-r^2)SST}{n-2}$$

אמידת שונות הטעויות:

מבחן F:

מבחן זה נעשה על מנת לבדוק האם משווהת הרגרסיה מובהקת.

המבחן מתבסס על פירוק סכום הריבועים :

$$SST = SSR + SSE$$

$$S_Y^2 = r^2 S_Y^2 + (1-r^2) S_Y^2$$

טבלת ניתוח שונות (טבלת ANOVA) :

מקור	סכום ריבועים SS	דרגות חופש $d.f.$	ממוצע סכום ריבועים $MS = \frac{SS}{d.f.}$	F
מודל הרגרסיה	SSR	1	$MSR = \frac{SSR}{1}$	$\frac{MSR}{MSE}$
שאריות	SSE	$n-2$	$MSE = \frac{SSE}{n-2}$	
סה"כ	SST	$n-1$		

כלל הכרעה :

אם : $F_{st} > F_c \alpha(1, n-2)$ נדחה את השערת האפס.

מבחן t:

מבחן זה נעשה על מנת לבדוק האם מקדם הרגרסיה מובהק.

$$\text{סטטיסטי המבחן : } t_{st} = \frac{\hat{\beta}_1 - \beta_{1,0}}{s.e.(\hat{\beta}_1)} \sim t_{c(n-2)}$$

$$s.e.(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{SXX}}$$

$$\text{אם השערת האפס מתiyחשת ל- } \beta_0 = \hat{\beta}_1 - \frac{r^2 \sqrt{n-2}}{\sqrt{1-r^2}} \text{ : (בדר"כ)}$$

כלל הכרעה :

השערה דו צדדית $H_1: \beta_1 \neq \beta_{1,0}$	השערה חד צדדית שמאלית $H_1: \beta_1 < \beta_{1,0}$	השערה חד צדדית ימנית $H_1: \beta_1 > \beta_{1,0}$	
$t_{statistic} = \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\frac{s.e.(\hat{\beta}_1)^2}{SXX}}} = \frac{r^2 \sqrt{n-2}}{\sqrt{1-r^2}}$			סטטיסטי המבחן
$ t_{statistic} \geq t_{n-2,1-\alpha/2}$	$t_{statistic} \leq -t_{n-2,1-\alpha}$	$t_{statistic} \geq t_{n-2,1-\alpha}$	אזור דחיפה
$2 * P(t_{n-2} > t_{statistic})$	$P(t_{n-2} > t_{statistic})$	$P(t_{n-2} > t_{statistic})$	P-VALUE

- שימושו לב Ci במודל של רgression ליניארית פשוטה ערך ה- t סטטיסטי

שתחזק שווה בדיקת שורש של ערך F המוחושב:
 $Pvalue = Pvalue$

רוח סמך לאמידת β : $\alpha = 1 - \text{גבול תחתון} \leq \beta \leq \text{גבול עליון} p$.

$$\cdot \hat{\beta}_1 \pm t_{n-2,1-\frac{\alpha}{2}} \cdot s.e.(\hat{\beta}_1)$$

מדד טיב ההתאמה : R^2

מדד שנותן את פרופורציות השונות המוסברת. כמה מהשונות של Y מוסברת על ידי השונות של X :

(X מסביר את כל השונות של Y) : $0 \leq R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} \leq 1$ (X לא מסביר כלל מהשונות של Y).

נרצה פרופורציות שונות מוסברת קרובה ככל האפשר ל-1.

אחווש השונות המוסברת : $R^2 \cdot 100$.

רב"ס שמטרתו לאמוד את תוחלת ערכי המשתנה תלוי (μ_0) עבור ערך מסוים של המשתנה הב"ת (x_0). במקרה אחריות אנו מתבקשים לאמוד את הניבוי באוכלויסיה עבור ערך מסוים של X.
 האומד הנקודתי (הסטטיסטי) סביבו בניוי הרב"ס הוא הניבוי במדגם עבור אותו

$$\text{ח- } X : \hat{\mu}_0 = \hat{y}_0 = \ddot{\alpha} + \hat{\beta}x_0$$

$$\cdot \hat{\mu}_0 \pm t_{n-2} \left(\frac{\alpha}{2} \right) \cdot \sqrt{MSRES \cdot \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SSX} \right)}$$

נוסחת הרב"ס :

טעות התקן/גודל הרב"ס מושפעים מ-4 גורמים :

1. $MSRES$ - האומד לשונות הטיעויות. ככל שגדל, טעות התקן/הרבי"ס גדלים ולהפוך.
2. n - גודל המדגם. ככל שגדל, טעות התקן/הרבי"ס קטנים ולהפוך.
3. SSX - מונה השונות של X (קשרו לתופעת קיצוץ תחום). ככל שגדל, טעות התקן/הרבי"ס קטנים ולהפוך.
4. $(x_0 - \bar{x})$ - הסטייה של ערך X המסוים מהממוצע של X . ככל שגדלה טעות התקן/הרבי"ס גדלים ולהפוך.

רב"ס לערכי Y עבור ערך מסוים של X :

רב"ס שמטרתו לאמוד את כל טווח ערכי Y (y_0) עבור ערך X מסוים (x_0).

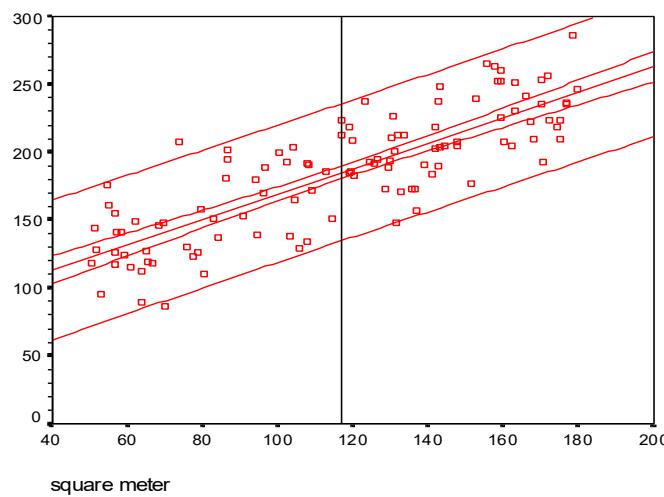
$$\hat{\mu}_0 \pm t_{n-2} \left(\frac{\alpha}{2} \right) \cdot \sqrt{MSRES \cdot \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SSX} \right)}$$

נוסחת הרב"ס :

ניתן לראות כי גם רב"ס זה בניו סביב האומדן הנקודתי לתוכלת ערכי Y עבור ערך X מסוים ($\hat{\mu}_0$).

ההבדל בין רב"ס לערכי Y לבין הרבי"ס לתוכלת ערכי Y בא לידי ביטוי בטיעות התקן. ניתן לראות כי טיעות התקן של הרבי"ס לערכי Y גדולה יותר מטיעות התקן של הרבי"ס לתוכלת ערכי Y . כאשר כל יתר הפרמטרים נשארים קבועים רב"ס זה יהיה רחב יותר מן הרבי"ס לתוכלת.

התרשימים הבא מתאר רב"ס לתוכלת ולערך המשתנה תלוי וממחיש זאת בבירור :



שאלות:**קו הרגרסיה:**

- 1) מתוויך דירות בתל אביב רצה לבדוק איך משפיע גודלה של דירה על המחיר שבו היא נמכרת. הוא הניח 2 הנחות מקדיומות:
1. רק גודל הדירה משפיע על מחיר הדירה באופן שיטתי. כל שאר הדברים המשפיעים על מחיר הדירה הם אקראיים ולא ניתנים לחיזוי.
 2. ההשפעה של גודל הדירה על מחיר הדירה היא ליניארית.
- גודל הדירה הינו X ומחיר הדירה הינו Y . מודל המתווך: $y_i = \alpha + \beta x_i + \varepsilon$.
- המתווך אוסף נתונים על 6 דירות, שנמכרו בחודש האחרון באותואזור:

מספר הדירה	גודל הדירה ב' m^2 '	מחיר הדירה באלפי'\$
1	$X_1 = 70$	$Y_1 = 190$
2	$X_2 = 70$	$Y_2 = 210$
3	$X_3 = 80$	$Y_3 = 250$
4	$X_4 = 100$	$Y_4 = 290$
5	$X_5 = 120$	$Y_5 = 360$
6	$X_6 = 120$	$Y_6 = 380$

- א. מקדם המתאים בין גודל הדירה למחיר הדירה. מהמשמעותו?
- ב. קו הרגרסיה לניבוי מחיר הדירה באמצעות גודל הדירה ופרשו את משמעות המקדים.
- ג. המחיר החזווי על פי קו הרגרסיה של דירה בגודל 100 מ'ר.

מבחן F:

- 2) בצעו מבחן F לבדיקת הקשר בין גודל הדירה למחיר ברמת מובהקות של 1%.

מבחן t:

- 3) בהמשך לדוגמא הניל:
- א. בצעו מבחן t למובהקות מקדם הרגרסיה ברמת מובהקות של 1%.
 - ב. בדקו את הטענה כי עליה במי'ר אחד עולה את מחיר הדירה ביותר מ-2000\$. מהו ה-*pvalue* של מובהקות הקשר בין גודל הדירה למחיר. מהמשמעותו?

קשר בין מבחן F ל מבחן t :

- 4) חשבו את סטטיסטי המבחן F על סמך סטטיסטי המבחן t שקיבלתם.
מה ה-*pvalue* של מבחן F?
- 5) חשבו רבי"ס לאמידת מקדם הרגסיה ברמת סמך של 0.99.
השו עם תוצאות מבחן t.
- 6) חשבו את אחוז השונות המוסברת של מחיר הדירה על ידי גודלה.

רוח בר סמך לתוצאות:

- 7) השאלה מבוססת על נתונים דוגמא מס' 2 (ראו סרטון) והפליטים הבאים:

Case Summaries

	N	Mean	Std. Deviation	Minimum	Maximum
SIZE square meter	112		39.13942	50.46	179.76
PRICE thousands \$	112	185.0664	44.45345	86.20	286.56

Coefficients^a

Model	Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95% Confidence Interval for B	
	B	Std. Error				Lower Bound	Upper Bound
1 (Constant)				15.173	.000	60.979	91.015
SIZE square meter	.062	.823					

a. Dependent Variable: PRICE thousands \$

Descriptive Statistics

	Mean	Std. Deviation	N
SIZE square meter	116.740	39.139	112
PRICE thousands \$	185.066	44.453	112
PRE_1 Unstandardized Predicted Value	185.066	36.568	112
RES_1 Unstandardized Residual	.000	25.277	112

חשב רבי"ס ברמת סמך של 95% לתוצאה מחיר הדירה כאשר שטח הדירה הוא 100 מ"ר.

רוח בר סמן לערכי נעלם:

- 8) חשב רב"ס ברמת ביטחון של 95% למחיר הדירה עבור שטח דירה של 100 מ"ר. מה ההבדל בין רב"ס זה לרב"ס הקודם?

תרגול מסכם:

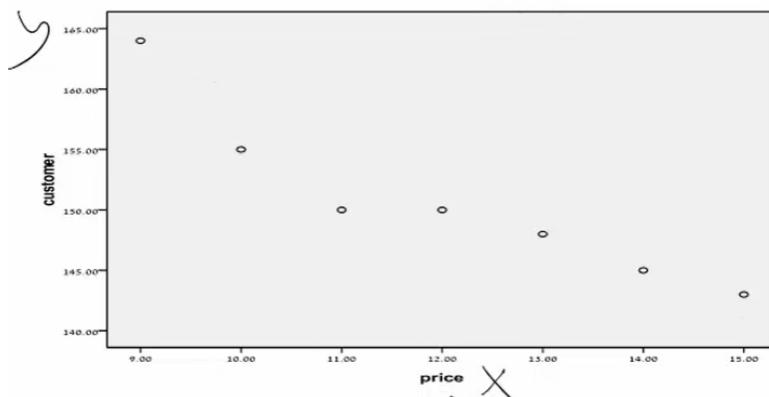
- 9) בפיוצציות "שלמה המלך" חושדים כי מספר הלקוחות המבקרים בפיוצציה תלוי במחיר המכירה של הבירה במקום. לשם בדיקת הנושא ערכו ניסוי בו בכל שבוע שינו את מחיר הבירה במקום ומנו את מספר הלקוחות שהגיעו במשך Wochen. משך הניסוי 7 שבועות עוקבים. להלן נתוני הניסוי:

שבוע 7	שבוע 6	שבוע 5	שבוע 4	שבוע 3	שבוע 2	שבוע 1	מחיר הבירה	כמות הלקוחות
9	10	11	12	13	14	15		
164	155	150	150	148	145	143		

- א. אמדאו את מודל הרגרסיה ע"י חישוב מקדמי הרגרסיה.
- ב. חשבו את מקדם המתאים r_{xy} .
- ג. אמדאו את השונות של שאריות המודל.
- ד. בצעו בדיקה גראפית של אקרראיות השאריות.
- ה. חשבו את אחוז השונות המוסברת. מה ממשמעותה?
- ו. בצעו חיזוי לכמות הלקוחות אם מחיר הבירה יהיה 16 ש"ח. האם להערכתכם ניתן להסתמך על חיזוי זה?
- ז. בצעו מבחון F לבדיקה האם קיים קשר בין מחיר הבירה לבין כמות הלקוחות.
- ח. בצעו מבחון t לבדיקה האם קיים קשר בין מחיר הבירה לבין כמות הלקוחות המבקרים בפיוצציה ברמת מובהקות 5%. השוו את התוצאות.
- ט. אמדאו את מקדם הרגרסיה ברמת סמן של 0.95. השוו את התוצאה עם הסעיף הקודם.
- י. כתבו דו"ח קצר על הממצאים.

קריאת פלטים של SPSS

10) על סמך הנתונים של השאלה הקודמת התקבלו הפלטים הבאים :

דיאגרמת הפיזור (scatter plot)**סטטיטטיקה תיאורית (descriptive statistics)**

Descriptive Statistics

	Mean	Std. Deviation	N
customer	150.7143	7.01699	7
Price	12.0000	2.16025	7

פלט מקדם המתאים (correlations)

Correlations

		customer	Price
Pearson Correlation	customer	1.000	-.935
	price	-.935	1.000
Sig. (1-tailed)	customer	.	.001
	price	.001	.
N	customer	7	7
	price	7	7

פלט model summary**Model Summary^b**

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.935 ^a	.873	.848	2.73470

a. Predictors: (Constant), price

b. Dependent Variable: customer

פלט ניתוח שונות (ANOVA)**ANOVA^b**

Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression	258.036	1	258.036	34.503	.002 ^a
Residual	37.393	5	7.479		
Total	295.429	6			

a. Predictors: (Constant), price

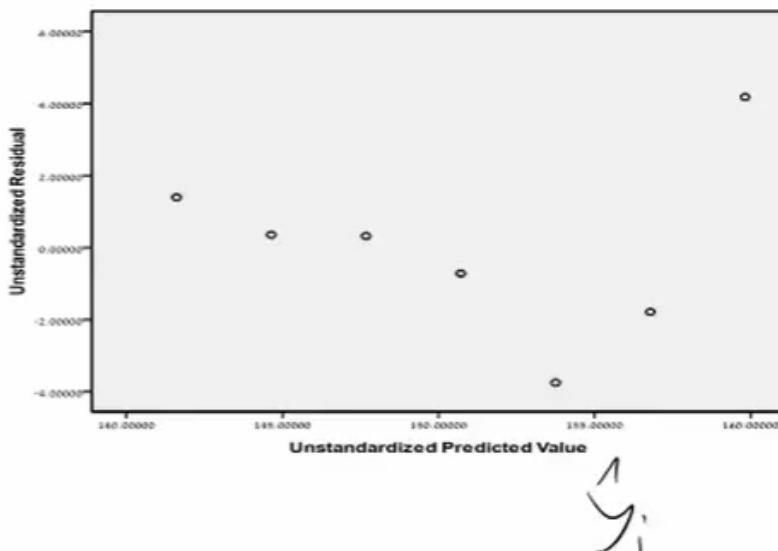
b. Dependent Variable: customer

פלט מקדמי הרוגרסייה (coeffitients)**Coefficients^a**

Model	Unstandardized Coefficients			t	Sig.
	B	Std. Error	Beta		
1 (Constant)	187.143	6.287		29.765	.000
	-3.036	.517	-.935	-5.874	.002

a. Dependent Variable: customer

גרף ניתוח שאריות:



על סמך הפלטים הנתוניים :

- א. מהו מודל הרגסיה שנאמד?
- ב. מהו מקדם המתאים r_{xy} ?
- ג. מהי השונות של שארית המודל?
- ד. האם נמצא דפוס מיוחד בשאריות?
- ה. מהו אחוז השונות הקשורות?
- ו. על פי מבחן F: האם קיים קשר בין מחיר הבירה לבין כמות הלקוחות המבקרים בפייצוצייה ברמת מובהקות 5%?
- ז. על פי מבחן t: האם קיים קשר בין מחיר הבירה לבין כמות הלקוחות המבקרים בפייצוצייה ברמת מובהקות 5%? השוו את התוצאות.
- ח. מה ה-pvalue של המבחן הסטטיסטיים? מה משמעותו?
- ט. בדקו האם קיים קשר חיובי מובהק בין המשתנים ברמת מובהקות 5%?

תשובות סופיות:

- ג. 301.68 אלף דולר ב. $\hat{Y}_t = -27.32 + 3.29 X_t$ א. $r = 0.987$ **(1)**
- . יש עדות לקשר מובהק. ב. $F_{st} = 21.198$ **(2)**
- ג. $pvalue < 0.001$ ב. יש עדות לכך. א. $t_{st} = 4.604$ **(3)**
- . $pvalue < 0.001$, $t_{st}^2 = 147$ **(4)**
- . $p(2.061 \leq \beta \leq 4.519) = 0.99$ **(5)**
- . 97.4% **(6)**
- . $p(163.889 \leq \mu_{100} \leq 174.24) = 0.95$ **(7)**
- . $p(119.036 \leq Y_{100} \leq 219.09) = 0.95$ **(8)**
- ג. $\hat{\sigma}^2 = 7.4785$ ב. $r = -0.93457$ נ. $\hat{y}_i = 187.143 - 3.0357x_i$ **(9)**
- ו. $\hat{y} = 138.5714$, כן. ח. $R^2 = 0.873$ ד. ראו סרטון.
- ח. $t_{st} = -5.87395$ י. יש עדות לכך. ג. $F_{st} = 34.5 > F_c(0.05(1,5)) = 6.6$ ט. ראו סרטון.
- ג. $MSE = 7.479$ ב. $r_{xy} = 0.935$ א. $\hat{y}_i = 187.143 - 3.036x_i$ **(10)**
- ו. $F = 34.503$ ח. $R^2 = 0.874$ ד. לא.
- ט. ראו סרטון. נ. $pvalue = 0.002$ ז. $t = -5.874$

כליים כמותיים מתקדמים של תכנון סטטיסטי לאיכות

פרק 3 - רגרסיה מרובה

תוכן העניינים

36 1. כללי

גרסיה מרובה:

רקע:

ניבוי המשנה תלוי באמצעות יותר ממשנה ב"ית אחד.

המודל אוכלוסייה: $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$.

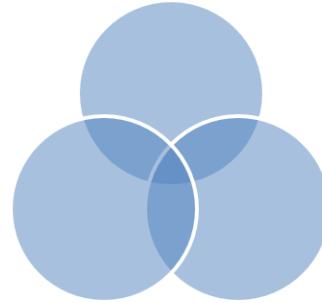
מקדמי מודל הרוגסיה המרובה:

α = חותך אחד שמשמעותו: הציון המנווא כאשר כל המשתנים הב"ית = 0.

β_1, β_2, \dots = מקדמי השיפוע. מס' הבנות = מספר המשתנים הב"ית במודל.

משמעות מקדם השיפוע β_j : ההשפעה הייחודית של המשנה הב"ית מסוים לניבוי

המשנה תלוי, בינוי השפעתם של כל יתר המשתנים הב"ית האחרים המוצאים
במשוואת הרוגסיה.



אמידת מודל הרוגסיה המרובה:

ברוגסיה מרובה, כמו ברוגסיה פשוטה, שיטת האמידה הטובה ביותר היא שיטת הריבועים הפחותים. כמובן, נרצה להביא את סכום הטיעיות בניובי למינימום.

מפתרונו פונקציית הריבועים הפחותים קיבל את אומדי הרוגסיה: $\hat{\alpha}, \hat{\beta}_1, \dots, \hat{\beta}_k$.

מבחני מובהקות:

1. מבחן F ל證明keit הרוגסיה:

בדיקה האם קיים קשר ליניארי בין המשנה תלוי לבין לפחות אחד מהמשתנים המסבירים.

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

ההשערות הן: $H_1 : \text{Not } H_0 = \text{at least one of the } \beta's \text{ is not } 0$

טבלת ניתוח שונות (ANOVA) :

מקור	סכום ריבועים SS	דרגות חופש $d.f.$	ממוצע סכום ריבועים $MS = \frac{SS}{d.f.}$	$F_{st} \sim F_{k,n-k-1}$
מודל הרגסיה	SSR	K	$MSR = \frac{SSR}{K}$	$F_{st} = \frac{MSR}{MSE}$
שאריות	SSE	$n-k-1$	$MSE = \frac{SSE}{(n-k-1)}$	
סה"כ	TSS	$n-1$		

סטטיסטי המבחן : $F_{st} = \frac{MSR}{MSE}$

כלל הכרעה : נדחה את H_0 אם : $F_{st} \geq F_{k,n-k-1}^{1-\alpha}$

חישוב סכומי הריבועים :

$$\begin{aligned} TSS &= \sum_{i=1}^n y_i^2 - n\bar{y}^2 \\ SSR &= R^2 \cdot TSS \\ SSE &= (1-R^2)TSS \end{aligned}$$

פרופורציות השונות המוסברת R^2 :

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

ברגרסיה מרובה אומד זה לפרופורציות השונות המוסברת הוא בעייתי שכן הוא מושפע ממספר המשתנים הב"ת במודל. אומד זה יכול רק לגדול בהוספה משתנים ב"ת למודל ולכן לא ניתן לנו אינדיקציה האם כדאי להוסיף אותם למודל או לא.

האומד המתוקן לפרופורציות השונות המוסברת $AdjR^2$:

$$\bar{R}^2 = 1 - \left[\frac{(1-R^2)(n-1)}{n-k-1} \right]$$

בניגוד ל- R^2 לוקח בחשבון את מספר המשתנים הב"ת במודל. יכול שלא לגדול אף לקטונו בהוספה משתנה ב"ת שלא תורם תרומה משמעותית לניבוי.

2. מבוחן t לMOVבקות משתנה ב'ית יחיד :

$$\begin{array}{l} H_0 : \beta_j = 0 \\ \text{השענות :} \\ H_1 : \text{else} \end{array}$$

סטטיסטי המבחן וכלל הכרעת השערת האפס :

$$\cdot \left| t_{\hat{\beta}_j} \right| = \left| \frac{\hat{\beta}_j}{S_{\hat{\beta}_j}} \right| > t_{(T-k-1, 1-\frac{\alpha}{2})}$$

$$\cdot \hat{\beta}_j \pm t_{(n-k-1, 1-\frac{\alpha}{2})} s.e.(\hat{\beta}_j) : \beta_j$$

שאלות:

1) לצורך בדיקת ההשערה שקיים קשר בין מספר המוניות בעירobar שבע (y) לבין מספר התושבים בעיר באלפים (x_1) ומספר הרכיבים הפרטיים באלפים (x_2).

הוחלט לבנות מודל רגרסיה מהצורה : $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$, על סמך

$$\text{הנתונים הבאים : } MSE = 119.789, \sum_i y_i^2 = 338,657, \sum_i y_i = 1673$$

א. ע"י הנתונים הנ"ל, השלימו את טבלת ניתוח השונות הבאה.
אייזו השערה ניתן לבדוק באמצעותה? כתוב את ההשערה ובוחן אותה.

SOURCE	SS	DF	MS	F
Regression				
Error				
Total		8		

ב. חשבו את מדד טיב ההתאמה. הסבר את משמעותו.

ג. נתונה טבלת המקדים (החלקית) הבאה :

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	-511.727	114.9476				
X 1	9.208785		3.732167			
X 2	-8.79921	4.420456				

i. רשמו את האומדן לשווות הרגרסיה ופרשו את מקדמיה.

ii. בוחנו את ההשערה כי קיים קשר בין מספר הרכיבים הפרטיים לבין
מספר המוניות ברמת MOVבקות של 5%.

iii. בנו רוח סמך למקדם של מספר התושבים בעיר

- iv. ענה ללא חישוב (על סמך הסעיפים הקודמים) – האם קיים קשר בין מספר התושבים לבין מספר המוניות ברמת מובהקות 5%?
- v. מהי תחזית מס' המוניות בbara שבע עבור 100,000 תושבים ו-52,000 מוניות פרטיות?
- vi. האם ניתן לסמן על תחזית זאת?

תרגול מסכם:

2) מעוניינים למצוא קשר בין מחיר הדירה (\$-ב-) לבין ארבעה משתנים מסבירים: (1) שטח הדירה ו-(2) גודל שטח האמבטיה (ב-Sqft) וכן (3) מרחק הדירה מהים ו-(4) מהאוניברסיטה (במיילים).
 לשם כך נדגו מס' דירות והריצו רגרסיה אשר בה המשתנה המוסף הוא מחיר הדירה.
 להלן פلت הרגרסיה שהתקבל:

Model Summary

Mode	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.952 ^a			

a. Predictors: (Constant), Sea_Dist, Apartment, Bath

ANOVA^b

Model	Sum of Squares	df	Mean Square	F	Sig.
1 Regression					
Residual					
Total	1940484.615	25			.000 ^a

a. Predictors: (Constant), Univ_Dist, Bath, Sea_Dist, Apartment

b. Dependent Variable: Price

Coefficients^a

Model	Unstandardized Coefficients			t	Sig.
	B	Std. Error	Beta		
1 (Constant)	-265.514	146.673		-1.810	.085
Apartment					
Bath	4.256	.449	.722	6.572	
Sea_Dist	-32.114	11.090	.297	2.687	.014
Univ_Dist	11.746	9.439	-.223		.009
			.095	1.244	.227

a. Dependent Variable: Price

ענו על הסעיפים הבאים:

- א. מלאו את התאים החסרים בטבלה (אם לא ניתן למלא את כל התאים החסרים באופן מלא נמקו באופן מפורש מדובר לא ניתן).

ב. כתבו את האומדן למשוואת מחיר הדירה בצורה מפורשת על סמך הפלט הניל. פרשו את מקדמי הרגרסיה.

ג. בדקו האם ארבעת הגורמים ביחד אכן מסבירים את מחיר הדירה. הסבירו את המסקנה שהגעתם אליה. השתמשו ברמת מובהקות 5%.

ד. הסבירו מהו ערך $Pvalue$ ומה ניתן להסיק ממנו לגבי המשתנים המסבירים?

ה. בנו רוחסמן למקדם גודל שטח האמבטיה. השתמשו ברמת מובהקות של 2%.

ו. ברמת מובהקות של 5% יש לבדוק האם המרחק מהאוניברסיטה אכן משפיע על מחיר הדירה.

ז. האם במודל הרגרסיה הנוכחי ניתן לוותר על גורם המרחק מהים? השתמשו ברמת מובהקות 1%

ח. בדקו את ההשערה כי קיים קשר חיובי בין גודל הדירה למחירה ברמת מובהקות של 5%.

תשובות סופיות:

. א. השערת $H_0 : \beta_1 = \beta_2 = 0$ (1) .
 $H_1 : \text{at least one of the } \beta's \text{ is not 0}$

SOURCE	SS	DF	MS	F
Regression	26945.784	2	13472.892	112.414
Error	718.733	6	119.789	
Total	27664.517	8		

- ב. $\hat{y}_i = -511.727 + 9.208x_{1i} - 8.799x_{2i}$ ג. i.e. 97.4%.

ii. אין עדות לכך. iii. $p(3.17 \leq \beta_1 \leq 15.24) = 0.95$.

iv. כן. v. 321. vi. כן.

(2) א. ראו סרטון. ב. $\hat{y}_i = -256.514 + 2.95x_{1i} + 4.256x_{2i} - 32.114x_{3i} + 11.746x_{4i}$

ג. לפחות אחד מהמשתנים הב'ית שונה מאפס באוכלוסייה.

ד. ראו סרטון. ה. $p(1.016 \leq \beta_2 \leq 7.496) = 0.98$.

ו. לא. ז. לא. ח. יש עדות לכך.

כלים כמותיים מתקדמים של תכנון סטטיסטי לאיכות

פרק 4 - רגרסיה - שאלות מבחנים

תוכן העניינים

41	1. מבחן 1
45	2. מבחן 2

מבחן 1: **שאלות:**

1) להלן תוצאות הרצת רגסיה של Y בתלות ב- X עבור 10 תצפיות (חלק מהנתונים הושמו בכונה מהפלט, אך ניתנים לחישוב על ידך).

$$\text{נתון כי: } \sum (X_i - \bar{X})^2 = 1475.6.$$

מקור SS	סכום ריבועים
SSR = 2148.6	רגסיה
SSE = ?	שאריות
SST = ?	סה"כ

משתנה	מקדם b_i	טעות תיקן s_{bi}	ערך סטטיסטי(מתוקן) t	מבחן מובהקות	p-value
קבוע (חותם)	-24.7	11.3	?	?	?
X	1.20	?	10.5		

א. מהו SST?

.i. לא ניתן לקבוע.

.ii. 1994.42

.iii. 2304.1

.iv. 1629.78

ב. האם הרגסיה מובהקת? בדוק לפי p value.

.i. הרגסיה מובהקת.

.ii. הרגסיה אינה מובהקת.

(2) לפניך פلت רגרסיה פשוטה (ממנו הושמו נתונים נתוניים שבאפשרותך להשלים), המתאר את ציון המבחן כפונקציה של מספר התרגילים שהגיש הסטודנט במהלך הסמסטר, ידוע כי כל הנחות המודל תקפות.

SUMMARY OUTPUT

Regression Statistics

Multiple R 0.842105598

R Square 0.709141839

Adjusted R Square 0.688366256

Standard Error 5.315523758

Observations 16

ANOVA

	Df	SS	MS	F	Significance F
Regression	1	964.4329004	964.4329004	34.13343	4.28E-05
Residual	14	395.5670996	28.25479283		
Total	15	1360			

Coefficients Standard Error t Stat P-value

Intercept	50.99134199	4.318918368	11.80650747	1.15E-08
	4.086580087	0.699471573	5.842381939	4.28E-05

א. ע"פ הנתונים, אחוז השונות של ציוני המבחן הקשורות ע"י מספר התרגילים שהגיש הסטודנט, היא _____. אם נוסיף משתנים נוספים,

אחוז השונות הקשורות _____ Radjusted, ו- _____.

.i. 84%, יגדל, לא ניתן לקבוע ללא נתונים נוספים.

.ii. 84%, יקטן, יגדל.

.iii. 70.9%, יגדל, לא ניתן לקבוע ללא נתונים נוספים.

.iv. 70.9%, יקטן, יקטן.

.v. 68.8%, יגדל, יגדל.

.vi. 68.8%, יקטן, יקטן.

ב. מהו הרבי"ס של שיפוע הרגסיה β ? (בר"מ של 1%).

.i. $2.58 < \beta < 5.58$

.ii. $.2 < \beta < 6.17$

.iii. $1.74 < \beta < 6.88$

.iv. $2.86 < \beta < 5.3$

במטרה לנבד בצורה טובה יותר את הצלחת הסטודנטים בבחינה, החלטיט החוקר להוסיף 2 משתנים נוספים לניתוח הרגסיה.
מספר השיעורים בהם נכח הסטודנט, ומספר השעות שלמד בבחינה.
לפניכם הפלט החסר:

SUMMARY OUTPUT

Regression Statistics

Multiple R 0.880163577

R Square 0.774687922

Adjusted R Square 0.718359902

Standard Error 5.053253294

Observations 16

ANOVA

	df	SS	MS	F
Regression	3	?	351.1918579	13.75315
Residual	12	306.4244263	25.53536886	
Total	15	?		

	Coefficients	Standard Error	t Stat	P-value
Intercept	37.08959571	8.726136945	4.250402663	0.001127
	2.610000755	1.429904588	1.82529714	0.092932
	3.198068014	1.239162836	2.58082951	0.024061
	-0.108373802	1.021202927	-0.106123669	0.917238

- ג. בדוק את ההשערה כי הרגרסיה מובהקת לכל אחד מהמשתנים המסבירים בר"מ של 1%.
- i. הרגרסיה אינה מובהקת לכל המשתנים שנבדקו.
 - ii. לא ניתן לקבוע מהנתונים האם הרגרסיה מובהקת.
 - iii. הרגרסיה מובהקת למשתנה מספר התרגילים, אך אינה מובהקת למשתנים מספר השיעורים ומספר שעות הלימוד בבחינה.
 - iv. הרגרסיה מובהקת למשתנה מספר התרגילים ומספר השיעורים, אך אינה מובהקת למשתנה שעות הלימוד בבחינה.
- ד. מהו SSR של הרגרסיה המרובה?
- .i. $SSR = 964.4$
 - .ii. $SSR = 1053.57$
 - .iii. $SSR = 694.57$
 - .iv. $SSR = 853.57$
 - v. לא ניתן לחשב את SSR מהנתונים שהתקבלו.

תשובות סופיות:

- | | | | | | |
|--------|-------|--------|--------|---------|-----|
| ד. ii. | ג. i. | ב. ii. | ב. .i. | א. iii. | (1) |
| ד. ii. | ג. i. | ב. ii. | ב. .i. | א. iii. | (2) |

מבחן 2:

שאלות:

- 1) הoulתת השערה שהוצאות האחזקה של מערכת לעיבוד נתונים קשורות למספר שעות השימוש השבועיות במערכת. להלן תוצאות חלקיות של פلت EXCEL של ניתוח וגרסיה בין Y הוצאות אחזקה שנתיות (במאות \$) ו- X מספר שעות השימוש בשבועיות.

SUMMARY OUTPUT

Regression Statistics					
	Multiple R	R Square	Adjusted R Square	Standard Error	Observations 10

ANOVA					
	df	SS	MS	F	Significance F
Regression	1	860.051	860.051	860.051	0.00012
Residual	8	144.525	18.065		
Total	9	1004.525			

	Coefficients	Standard Error	t Stat	P-value
Intercept	10.528	3.745	2.811	0.023
שעות שימוש	0.953	0.227	4.201	6.901

A. לאור התוצאות ה-p-value בבדיקה ההשערה:
 $H_0 : \beta_1 \leq 0.5$
 $H_1 : \beta_1 > 0.5$

- . Pv < 0.005 .i
- . 0.005 < Pv < 0.01 .ii
- . 0.01 < Pv < 0.025 .iii
- . 0.025 < Pv < 0.05 .iv
- . Pv > 0.05 .v

ב. אם היינו מרייצים רgresיה בה Σ מספר שעות השימוש השבועיות ואילו המשתנה המסביר X, הוצאות אחזקה שנתיות (במאות \$), אז השיפוע של קו הרgresיה יהיה:

- .i. 0.953
- .ii. 1.049
- .iii. 10.528
- .iv. 0.095
- .v. 0.898

(2) הoulתת השערה שמספר התקנות ברכב Y קשורה לגיל הנהג X. לשם כך נלקח מדגם של 10 נהגים.

כמו כן חושבו הסכומים הבאים:

$$\sum X_i^2 = 14,227, \sum X_i = 363, \sum Y_i = 13, \sum Y_i^2 = 29, \sum X_i Y_i = 366$$

$$\text{משוואת קו הרgresיה נתונה על ידי: } \hat{Y} = 4.96 - 0.1X_i$$

א. ערכו של מקדם המתאים הליניארי בין מספר התקנות לבין גיל הנהג הוא:

- .i. -10.59
- .ii. -0.9395
- .iii. 0.8826
- .iv. -0.1

החוקר לא היה מרוצה מעצמת הקשר ולכן הוסיף לרgresיה את המשתנים הבאים: מספר הק"מ שהמכונית נסעה (באלפי ק"מ) וסוג הרכב. במדגם נכללו 2 סוגי רכבים: A ו-B כאשר סוג B קודד כערך 0 וסוג רכב A קודד כערך 1. להלן פלט הרgresיה המרובה. שימו לב כי חלק מהערכים בפלט חסרים:

SUMMARY OUTPUT

Regression Statistics	
Multiple R	
R Square	
Adjusted R Square	
Standard Error 0.204047	
Observations 10	

ANOVA

	df	SS	MS	F
Regression	3			0.00002
Residual	6	0.249811		
Total	9			

	Coefficients	Standard Error	t Stat	P-value
Intercept	1.407943	0.71094		
X1 גיל הנ抬起头	-0.03094	0.014611		
X2 סוג הרכב	0.239574	0.152994		
X3 ק"מ באלפים	0.027666	0.005329		

- ב. להלן מספר טענות לגבי מובהקות הרgresיה ומובהקות המשתנה סוג הרכיב ברמת מובהקות 0.1 :
- .i. הרgresיה מובהקת והמשתנה סוג הרכיב מובהק ברמת מובהקות 0.1.
 - .ii. הרgresיה מובהקת, אך המשתנה סוג הרכיב אינו מובהק ברמת מובהקות 0.1.
 - .iii. הרgresיה אינה מובהקת, אך המשתנה סוג הרכיב מובהק ברמת מובהקות 0.1.
 - .iv. הרgresיה והמשתנה סוג הרכיב אינם מובהקים ברמת מובהקות 0.1.
- ג. SST בפלט הרgresיה המרובה שווה ל :
- .12.1 .i
 - .16.0 .ii
 - .20.0 .iii
 - .iv. אין מספיק נתונים לחשבו.
- ד. לאור התוצאות, רב"ס למקד המשתנה מספר הקילומטרים, בר"מ 5% :
- .i. (0.015,0.04)
 - .ii. (0.017,0.038)
 - .iii. (0.02,0.035)
 - .iv. אין מספיק נתונים לחשבו.

ה. אם היינו מקודדים את סוג רכב B כערך 1 וסוג רכב A כערך 0 משווהת הרגרסיה הייתה:

i. נשארת ללא שינוי.

$$\text{ii. } \hat{Y} = 1.4079 - 0.0309X_{1i} - 0.2395X_{2i} + 0.0276X_{3i}$$

$$\text{iii. } \hat{Y} = 1.6475 - 0.0309X_{1i} - 0.2395X_{2i} + 0.0276X_{3i}$$

iv. לא ניתן לדעת ללא הרצה מחדש.

ו. החוקר רצה להוציא משטנה מסביר נוסף X_4 מספר השנים שהלפו מאז קבלת רישיון הנהיגה. להלן פلت הרגרסיה (חלק מהנתונים חסרים) עם 4 המשתנים המסבירים:

	Coefficients	Standard Error	t Stat	P-value
Intercept	1.5	0.71		
X1 גיל הנגיג	5.1	14.1		
X2 סוג הרכב	0.25	0.2		
X3 ק"מ באלים	0.02	0.008		
X4 מס' השנים שהלפו מאז קבלת רישיון	-5.13094	0.8		

על סמך נתונים אלו:

i. חשש סביר למולטיקוליניאריות.

ii. התחזית הנקודתית של מספר התקלות תהיה דומה לו של הרגרסיה הקודמת (עם 3 המשתנים המסבירים).

iii. קיימת קורלציה גבוהה בין חלק מהמשתנים המסבירים.

iv. כל התשובות נכונות.

(3) בפועל מסוים הורץ מודל של רגרסיה ליניארית פשוטה על 8 עובדים כאשר Y תפוקת העובד ו- X גיל העובד. נמצא שהחלק המוסבר על ידי הרגרסיה הוא 74%. הטעויות שהתקבלו מופיעות בחלקן בטבלה שלהלן:

e_8	e_7	e_6	e_5	e_4	e_3	e_2	e_1
0	2	-2	2	?	2	-3	0

אחת מהטענות שלහן נכונה:

א. מקדם ההסבר בין X ל- Y הוא 0.74.

ב. לא ניתן לחשב את מקדם המתאים המרובה.

ג. $\text{SSR} = 18$.

ד. $\text{SSE} = 74$.

ה. אף אחת מהטענות אינה נכונה.

תשובות סופיות:

- (1) א. ii. ב. v.
(2) א. ii. ב. v. ג. v. ד. i. ה. iii. ו. iv.
(3) א.ii.