

הצגת נתונים בטבלאות ובתרשימים גרפיים

מושגים

משתנה בדיד - משתנה כמותי היכול לקבל מספר סופי של ערכים בין כל שני ערכים אפשריים. למשל: מספר הנפשות במשפחה, מספר מכשירי T.V בדירה וכדומה.

משתנה רציף - משתנה כמותי היכול לקבל אינסוף ערכים בין כל שני ערכים אפשריים. למשל: גובה, משקל, טמפרטורה וכדומה.

שכיחות - מספר הפעמים שהתקבל ערך של המשתנה. השכיחות מסומנת ב- $f(x)$

התפלגות שכיחויות - טבלה שבה מצוינים ערכי המשתנה ולצד כל ערך שכיחות הופעתו.

שכיחות יחסית של מחלקה - שכיחות המחלקה מחולקת בגודל המדגם, $f(x)/n$.

שכיחות מצטברת - סה"כ השכיחות עד לגבול העליון של המחלקה. מסומנת ב- $F(x)$.

צפיפות - צפיפות היא השכיחות ליחידה של המשתנה הנחקר. הצפיפות שווה לשכיחות המחלקה המחולקת ברוחבה. הצפיפות מסומנת ב- d_i . בהיסטוגרמה הצפיפות היא גובה המלבן.

דיאגרמת מקלות - הצגה גראפית המתאימה למשתנה איכותי או למשתנה כמותי בדיד. על הציר האופקי X יירשמו ערכי המשתנה הנחקר, ועל הציר האנכי Y נציין את השכיחות המוחלטת של הערכים או שכיחות יחסית באחוזים. מעל כל ערך שעל ציר ה- X יוצב 'מקל' שגובהו פרופורציוני לשכיחות המקרים שהערך מייצג.

היסטוגרמה - הצגה גרפית של התפלגות שכיחויות עבור משתנה המקובץ במחלקות. הצגה זו בנויה ממלבנים כאשר המחלקה מיוצגת ע"י אורך של קטע – בסיס המלבן, והשכיחות מיוצגת ע"י שטח המלבן.

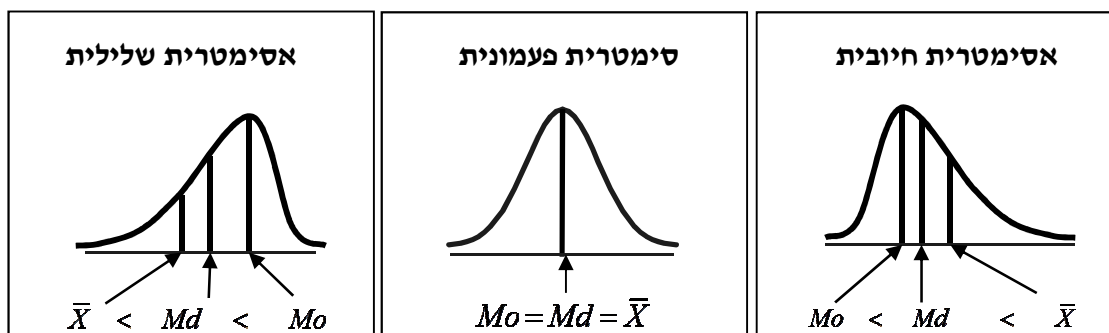
תיאור גרפי של צורת ההתפלגות:

	דיאגרמת מקלות	היסטוגרם
הנתונים	רשימת תצפיות	רשימת מחלקות
על ציר X על ציר Y	הערכים המדויקים של התצפיות <u>שכיחויות</u> הערכים	התחומים של ערכי התצפיות <u>הצפיפות</u> של המחלקות
השכיחות	מתוארת ע"י גובה המקל	מתוארת ע"י שטח המלבן

שם הממד	סימון	חישוב הממד	
		ברשימת תצפיות/טבלת שכיחויות משתנה בדיד	בטבלת שכיחויות - משתנה מקובץ במחלקות
שכיח	Mo	ערך המשתנה בעל השכיחות הגבוהה ביותר	במחלקות שוות רוחב: השכיח הוא אמצע המחלקה בעלת השכיחות הגדולה ביותר. במחלקות שונות רוחב: השכיח הוא אמצע המחלקה בעלת הצפיפות הגדולה ביותר.
חציון	Md	$Md(X) = \begin{cases} X_{\left(\frac{n+1}{2}\right)} & n \text{ אי זוגי} \\ \frac{X_{\left(\frac{n}{2}\right)} + X_{\left(\frac{n}{2}+1\right)}}{2} & n \text{ זוגי} \end{cases}$	$Md = L_0 + \frac{\frac{n}{2} - F(x_{m-1})}{f(x_m)} \cdot (L_1 - L_0)$ <p>ראה פירוט למטה</p>
ממוצע חשבוני	\bar{X}	<p><u>טבלת שכיחויות</u></p> $\bar{x} = \frac{\sum x f(x)}{n}$ <p>x_{mid} - נקודות האמצע של המחלקות</p>	<p><u>רשימת תצפיות</u></p> $\bar{x} = \frac{\sum x_i}{n}$
אמצע הטווח	MR	$MR = \frac{x_{\min} + x_{\max}}{2}$	$MR = \frac{L_{\min} + L_{\max}}{2}$

* מרכיבי הנוסחה לחישוב חציון כשהמשתנה מקובץ במחלקות	
n	מספר התצפיות (הנתונים)
x_m	המחלקה בה נמצא החציון (זו המחלקה שבה לראשונה $(F(x_m) \geq \frac{n}{2})$)
L_0	גבול אמיתי תחתון של המחלקה x_m
$f(x_m)$	שכיחות המחלקה x_m
$F(x_{m-1})$	השכיחות המצטברת עד למחלקה x_{m-1} (המחלקה הקודמת ל- x_m)
בנוסחה הנ"ל אין הבחנה בין n זוגי לאי זוגי. כדי להשתמש בנוסחה יש לעבוד בגבולות אמיתיים.	

מיקום מדדי מיקום מרכזי (שכיח, חציון וממוצע) בצורות התפלגות שונות



תכונות מדדי מיקום מרכזי

המדד	תכונות
שכיח	<ul style="list-style-type: none"> השכיח קל לחישוב. תיתכן התפלגות עם יותר משכיח אחד (1,1,2,2,3), ותיתכן אף התפלגות שאין לה ערך שכיח (1,2,3,7,9) השכיח אינו מושפע מערכים קיצוניים. מוסר מעט אינפורמציה על הנתונים.
אמצע הטווח	<ul style="list-style-type: none"> אמצע הטווח רגיש מאוד לערכים הקיצוניים בהתפלגות ורק להם.
חציון	<ul style="list-style-type: none"> החציון אינו מושפע מערכים קיצוניים. החציון מושפע מסדר הערכים ולא מהערכים עצמם. לכן החציון לא ישתנה כל עוד ה"מאזן" 50:50 מתחתיו ומעליו לא יופר. כאשר התפלגות הנתונים אסימטרית עם נטיה חזקה לכיוון מסוים נעדיף את השימוש בחציון כמדד מרכזי על פני השימוש בממוצע.
ממוצע חשבוני	<ul style="list-style-type: none"> הממוצע אינו חייב להיות אחד מערכי הסדרה. למשל: מס' הילדים הממוצע במשפחה ישראלית הוא 3.5, אך לא קיימת משפחה שיש לה 3.5 ילדים. הממוצע אינו יכול להימצא מחוץ לטווח הנתונים. למשל ממוצע המספרים 20,42,38,50,29 חייב להיות בין 20 ל- 50. הממוצע מושפע מכל הערכים ובפרט מערכים קיצוניים. אם מוסיפים לסדרת ערכים ערך קטן מהממוצע, הממוצע יקטן, אם מוסיפים ערך גדול מהממוצע, הממוצע יגדל, אם מוסיפים ערך השווה לממוצע, הממוצע לא ישתנה. סכום ההפרשים של כלל הערכים בסדרת נתונים מהממוצע שלהם שווה תמיד לאפס כלומר $\sum_{i=1}^n (x_i - \bar{x}) = 0$. במילים אחרות: סכום ההפרשים מן הממוצע בסימן חיובי שווה לסכום ההפרשים מן הממוצע בסימן שלילי (בערך מוחלט).

מדדי פיזור (עמודים 35-42)

מדדי הפיזור מציינים את מידת פיזור הנתונים בהתפלגות בשני אופנים: טווח הפיזור או פיזור ביחס לממוצע.

- **טווח (Range)** - ההפרש בין הערך הגבוה ביותר בהתפלגות לבין הערך הנמוך ביותר בהתפלגות.
- **טווח בינרבעוני (Interquartile Range)** - טווח הכולל את 50% הערכים הנמצאים במרכז ההתפלגות בין הרבעון התחתון (Q_1) לרבעון העליון (Q_3).
- **שונות (Variance)** - ממוצע של ריבועי הסטיות מן הממוצע.
- **סטיית תקן (Standard Deviation)** - השורש הריבועי (החיובי) של השונות.

שם הממד	סימון	חישוב הממד	
		טבלת שכיחויות - משתנה מקובץ במחלקות	רשימת תצפיות/טבלת שכיחויות משתנה בדיד
טווח	R	$R = L_{\max} - L_{\min}$	$R = x_{\max} - x_{\min}$
טווח בין רבעוני	IQR	$Q_1 = L_0 + \frac{\frac{n}{4} - F(x_{m-1})}{f(x_m)} \cdot (L_1 - L_0)$ $Q_3 = L_0 + \frac{\frac{3n}{4} - F(x_{m-1})}{f(x_m)} \cdot (L_1 - L_0)$ \Downarrow $IQR = Q_3 - Q_1$	
שונות	S^2	<p><u>טבלת שכיחויות</u></p> $S_x^2 = \frac{\sum (x_{mid} - \bar{x})^2 f(x)}{n}$ <p>או</p> $S_x^2 = \frac{\sum x_{mid}^2 f(x)}{n} - \bar{x}^2$ <p>x_{mid} - נק' האמצע של המחלקות</p>	<p><u>רשימת תצפיות</u></p> $S_x^2 = \frac{\sum (x_i - \bar{x})^2}{n}$ <p>או</p> $S_x^2 = \frac{\sum x_i^2}{n} - \bar{x}^2$
סטיית תקן	S	$S_x = +\sqrt{S_x^2}$	$S_x = +\sqrt{S_x^2}$

תכונות של מדדי הפיזור

המדד	תכונות
טווח	<ul style="list-style-type: none"> הטווח מושפע רק מהערכים הקיצוניים בהתפלגות (min/max), ולא מתחשב בפיזור בין יתר הערכים.
טווח בינרבעוני	<ul style="list-style-type: none"> הטווח הבינרבעוני מתחשב רק ב- 50% מהנתונים הנמצאים במרכז ההתפלגות (ומתעלם ממחצית מנתוני ההתפלגות הנמצאים בקצוות). כאשר ההתפלגות אסימטרית עם נטיה חזקה לכיוון מסוים נעדיף את השימוש בטווח הבינרבעוני על פני סטיית התקן.
שוונות/סטיית תקן	<ul style="list-style-type: none"> השוונות נמדדת ביחידות של המשתנה בריבוע. למשל במדידת גובה בס"מ יחידות המדידה של השוונות הן ס"מ². מאידך סטיית התקן נמדדת ביחידות של המשתנה. סטיית התקן מבטאת בקירוב את המרחק הממוצע של הנתונים מממוצע ההתפלגות. סטיית התקן מתחשבת בסטיות כל הנתונים מן הממוצע ונותנת משקל יחסי מתאים לכל ערך בהתאם למרחקו מן הממוצע ובהתאם לשכיחותו. כל סטייה תורמת כריבוע גודלה, כך שלסטיות של ערכים קיצוניים יש השפעה גדולה יותר על המדד.

השפעת טרנספורמציה ליניארית על מדדי המיקום המרכזי ומדדי הפיזור

<p>נתונה סדרת ערכים: $x_1, x_2, x_3, \dots, x_n$. אם נבצע על כל ערך בסדרה את הפעולה הבאה $bx_i + a$ (הכפלה בקבוע b והוספת קבוע a) נקבל את הסדרה: $bx_1 + a, bx_2 + a, bx_3 + a, \dots, bx_n + a$. הפעולה שביצענו היא: $x' = bx + a$ נקראת טרנספורמציה ליניארית (המרה קווית).</p> <ul style="list-style-type: none"> כאשר $b=1$ השינוי הליניארי כולל הוספת¹ קבוע a בלבד לכל ערך בסדרה $(x' = x + a)$. כאשר $a=0$ השינוי הליניארי כולל הכפלת כל ערך בסדרה פי b בלבד² $(x' = bx)$. <p>בביצוע שינוי ליניארי מהצורה $x' = bx + a$ נוכל לקבל את המדדים של x' מתוך המדדים של x מבלי לבנות את ההתפלגות של x' באופן הבא:</p>											
<table border="1"> <thead> <tr> <th>מדדי מיקום מרכזי</th><th>מדדי פיזור</th></tr> </thead> <tbody> <tr> <td>$Mo_{x'} = b \cdot Mo_x + a$</td><td>$R_{x'} = b \cdot R_x$</td></tr> <tr> <td>$Md_{x'} = b \cdot Md_x + a$</td><td>$S_{x'} = b \cdot S_x$</td></tr> <tr> <td>$\bar{x}' = b \cdot \bar{x} + a$</td><td>$S_{x'}^2 = b^2 \cdot S_x^2$</td></tr> <tr> <td>$MR_{x'} = b \cdot MR_x + a$</td><td></td></tr> </tbody> </table>	מדדי מיקום מרכזי	מדדי פיזור	$Mo_{x'} = b \cdot Mo_x + a$	$R_{x'} = b \cdot R_x$	$Md_{x'} = b \cdot Md_x + a$	$S_{x'} = b \cdot S_x$	$\bar{x}' = b \cdot \bar{x} + a$	$S_{x'}^2 = b^2 \cdot S_x^2$	$MR_{x'} = b \cdot MR_x + a$		
מדדי מיקום מרכזי	מדדי פיזור										
$Mo_{x'} = b \cdot Mo_x + a$	$R_{x'} = b \cdot R_x$										
$Md_{x'} = b \cdot Md_x + a$	$S_{x'} = b \cdot S_x$										
$\bar{x}' = b \cdot \bar{x} + a$	$S_{x'}^2 = b^2 \cdot S_x^2$										
$MR_{x'} = b \cdot MR_x + a$											
<ul style="list-style-type: none"> הוספת קבוע a לכל הערכים בהתפלגות תגרום להגדלת כל מדדי המיקום המרכזי ב- a. הפחתת קבוע a מכל הערכים בהתפלגות תגרום להקטנת כל מדדי המיקום המרכזי ב- a. מדדי הפיזור אינם מושפעים מהוספת או הפחתת קבוע. הכפלת כל הערכים בהתפלגות בקבוע חיובי b תגרום להכפלת כל מדדי המיקום המרכזי פי b. גם מדדי הפיזור, למעט השוונות, יוכפלו פי b. (השוונות תוכפלו פי b^2) 											

¹ או הפחתת קבוע (a שלילי)

² למשל, אם x מייצג משקל בק"ג ו- $x' = 1000x$ משקל בגרמים אזי

חישוב ממוצע וסטיית התקן כאשר מאחדים שתי קבוצות או יותר

נתונות שתי קבוצות וידועים הגדלים הבאים:

II	I	
n_2	n_1	מספר נתונים
\bar{x}_2	\bar{x}_1	ממוצע
s_2^2	s_1^2	שונות

מאחדים את שתי הקבוצות לקבוצה אחת. מטרתנו היא לחשב את הממוצע המשוקלל והשונות המצורפת של הקבוצה המאוחדת כשידועים לנו הנתונים לעיל.

$$\bar{x} = \frac{n_1 \cdot \bar{x}_1 + n_2 \cdot \bar{x}_2}{n_1 + n_2} \quad \text{ממוצע משוקלל של שתי הקבוצות:}$$

$$S_c^2 = \frac{n_1 \cdot (\bar{x}_1^2 + s_1^2) + n_2 \cdot (\bar{x}_2^2 + s_2^2)}{n_1 + n_2} - \bar{x}^2 \quad \text{שונות מצורפת של שתי הקבוצות:}$$

$$S_c = +\sqrt{S_c^2}$$

שים לב!

- כאשר שתי הקבוצות שוות גודל ($n_1 = n_2$), הממוצע המשוקלל הוא הממוצע הפשוט של שתי הקבוצות - $\bar{x} = \frac{\bar{x}_1 + \bar{x}_2}{2}$
- כאשר ממוצעי שתי הקבוצות שווים, השונות המצורפת היא ממוצע משוקלל של שונויות שתי הקבוצות - $S_c^2 = \frac{n_1 \cdot s_1^2 + n_2 \cdot s_2^2}{n_1 + n_2}$
- ניתן להכליל את הנוסחאות הנ"ל גם ל- k קבוצות:

$$\bar{X} = \frac{\sum_{j=1}^k n_j \cdot \bar{x}_j}{n} \quad \text{ממוצע משוקלל:}$$

$$S_c^2 = \frac{\sum_{j=1}^k n_j \cdot (\bar{x}_j^2 + s_j^2)}{n} - \bar{X}^2 \quad \text{שונות מצורפת:}$$

מקרא: k - מספר הקבוצות, j - אינדקס הקבוצות, n_j - מספר הנתונים בקבוצה ה- j , \bar{x}_j - ממוצע

הקבוצה ה- j . n - גודל הקבוצה המצורפת: $n = \sum_{j=1}^k n_j$

המאון ה- k הוא ערך של המשתנה (C_K) ש $k\%$ מהנתונים בהתפלגות קטנים ממנו (או שווים לו).

נוסחאות לחישוב C_K ו- K_C כשנתונה טבלת שכיחויות של משתנה המקובץ במחלקות:

$$C_K = \frac{\frac{n \cdot k}{100} - F(x_{m-1})}{f(x_m)} \cdot (L_1 - L_0) + L_0 \quad - \quad \text{נוסחה לחישוב } C_K \text{ (קידוע)}$$

$$K_C = \left[\frac{C_K - L_0}{L_1 - L_0} \cdot f(x_m) + F(x_{m-1}) \right] \cdot \frac{100}{n} \quad - \quad \text{נוסחה לחישוב } K_C \text{ (קידוע)}$$

הערות:

- שיטת החישוב בעזרת הנוסחה זהה לזו שבנוסחת החציון בטבלת שכיחויות של משתנה מקובץ.
- כדי להשתמש בנוסחה יש לעבוד בגבולות אמיתיים.

ציוני תקן

ציון התקן של X מסומן ב- Z_x ומבטא בכמה סטיות תקן רחוק X ממוצע ההתפלגות.

$$Z_x = \frac{x - \bar{x}}{s_x} \quad - \quad \text{נוסחה לחישוב}$$

תכונות של ציוני תקן:

- ציוני התקן הם מספרים טהורים, בלתי תלויים ביחידות המדידה, ולכן ניתן להשוות בעזרתם מיקום יחסי של תצפיות מהתפלגויות של משתנים הנמדדים ביחידות שונות (למשל: משקל בק"ג עם גובה בס"מ או שכר בש"ח עם וותק בשנים).
- $Z_x > 0$ כאשר התצפית נמצאת מעל הממוצע $x > \bar{x}$
- $Z_x < 0$ כאשר התצפית נמצאת מתחת לממוצע $x < \bar{x}$
- $Z_x = 0$ כאשר התצפית שווה לממוצע $x = \bar{x}$
- ציוני תקן הם שינוי ליניארי על ערכי הסדרה: $Z_x = \frac{1}{s_x} \cdot x - \frac{\bar{x}}{s_x}$ $\left(b = \frac{1}{s_x}, a = -\frac{\bar{x}}{s_x} \right)$
- הממוצע של ציוני התקן שווה תמיד לאפס ($\bar{z} = 0$), והשונות/סטיות התקן שלהם שווה תמיד ל-1 ($S_z^2 = S_z = 1$).

מקדם המתאם של פירסון – r

מדד לעוצמת הקשר הלינארי בין שני משתנים כמותיים

טווח הערכים של המדד: $-1 \leq r \leq +1$
חישוב r :

$$r = \frac{\text{cov}(x, y)}{s_x \cdot s_y} \quad (I)$$

נוסחאות נוספות שקולות:

$$r = \frac{\sum_{i=1}^n x_i y_i - n \cdot \bar{x} \cdot \bar{y}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) \left(\sum_{i=1}^n y_i^2 - n \bar{y}^2 \right)}} \quad (III)$$

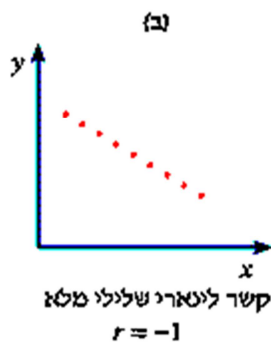
(III)

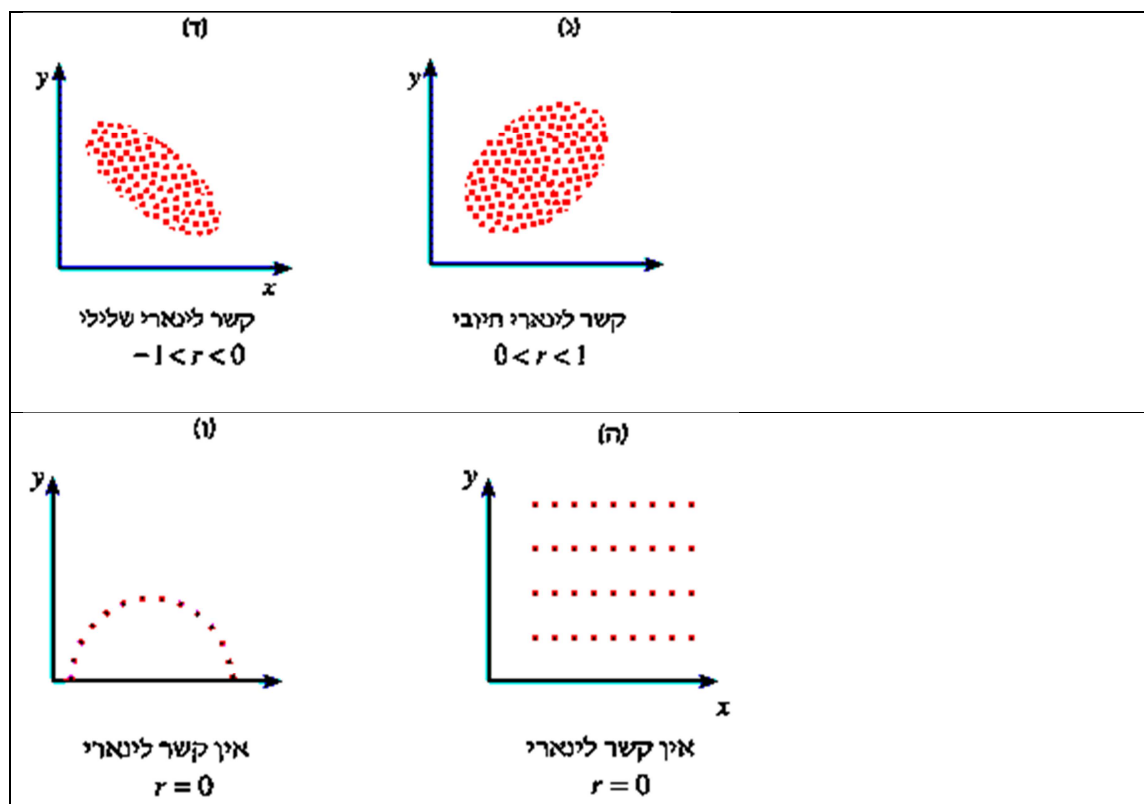
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (II)$$

(II)

דיאגרמות פיזור וערכו של R

מטרת דיאגרמת הפיזור היא לקבל רושם ראשוני מהי צורת הקשר בין שני המשתנים (לינארי/לא לינארי), מהי מגמתו (חיובית/שלילית) ומהי עוצמתו (חלשה/חזקה).
להלן מספר דוגמאות לדיאגרמות פיזור. לצד כל דיאגרמה מצוין ערכו המתאים של r .





במקרים א' ו-ב' כל התצפיות – הנקודות בדיאגרמת הפיזור – מונחות על קו ישר אחד. הקשר הוא קשר ליניארי מלא.

במקרה א' הקשר הוא במגמה חיובית ובמקרה ב' - במגמה שלילית.

במקרה ג' התצפיות מפוזרות סביב קו ישר במגמה חיובית. הקשר הוא ליניארי לא מלא.

במקרה ד' התצפיות מפוזרות סביב קו ישר במגמה שלילית. הקשר הוא ליניארי לא מלא. עוצמת הקשר הליניארי ב – ד' גבוהה מזו שב – ג' כיוון שהנקודות "קרובות" יותר לקו.

במקרה ה' אין קשר בין שני המשתנים

במקרה ו' אין קשר ליניארי בין שני המשתנים אך אין זה אומר שאין קשר בין שני המשתנים. מהתרשים אפשר לראות שיש קשר לא ליניארי (ממעלה שניה).

תכונות של r

1. טווח הערכים של המדד: $-1 \leq r \leq +1$
2. כאשר $r=+1$ הקשר בין x ל- y חיובי מלא. כל הנקודות בדיאגרמת הפיזור נמצאות על קו ישר אחד עולה. (ראה תרשים א). בין המשתנים מתקיים הקשר $y = bx + a$; $b > 0$.
- כאשר $r=-1$ הקשר בין x ל- y שלילי מלא. כל הנקודות בדיאגרמת הפיזור נמצאות על קו ישר אחד יורד. (ראה תרשים ב). בין המשתנים מתקיים הקשר $y = bx + a$; $b < 0$.
3. ככל ש- r גדל בערכו המוחלט ומתקרב ל- 1, הקשר הליניארי בין שני המשתנים חזק יותר.
4. אם $r=0$, פירוש הדבר שאין מתאם ליניארי, אך ייתכן שיש קשר אחר לא ליניארי.
5. השפעת טרנספורמציה לינארית על מקדם המתאם:
נתונים שני משתנים X ו- Y . ומקדם המתאם ביניהם: $r_{X,Y}$
נגדיר: $X' = b \cdot X + a$ ו- $Y' = c \cdot Y + d$ (טרנספורמציה לינארית)
אזי מקדם המתאם ביניהם (לאחר הטרנספורמציה) יהיה: $r_{X',Y'} = \frac{b \cdot c}{|b \cdot c|} \cdot r_{X,Y}$
6. מקדם המתאם הוא מספר טהור ואינו תלוי ביחידות המדידה.

שונוות משותפת

השונוות המשותפת של x ו- y - $\text{cov}(x, y)$ - מדד למידת ההשתנות המשותפת של שני משתנים כמותיים.

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{n} = \frac{\sum_{i=1}^n x_i y_i}{n} - \bar{x} \cdot \bar{y}$$

נוסחאות חישוב:

תכונות:

1. אם $\text{cov}(x, y) > 0$ אזי $r > 0$
- אם $\text{cov}(x, y) < 0$ אזי $r < 0$
- אם $\text{cov}(x, y) = 0$ אזי $r = 0$
2. $\text{cov}(x, x) = S_x^2$
3. השונוות המשותפת תלויה ביחידות המדידה של המשתנים.