

הראינו בט"ס"ר שלמה

$$a_{LS} = \bar{y} - b_{LS} \bar{x}$$

$$b_{LS} = \frac{Cov_{mp}(X, Y)}{Var_{mp}(X)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

מסקנה מעניינת: הנקודה (\bar{x}, \bar{y}) תמיד נמצאת על הישר הטוב ביותר.

נראה כי הנקודות המקסימליות (ט"ס"ר) הן:

$$a_{LS} = 56 \text{ cm} \quad \bar{x} = 170$$

$$b_{LS} = 0.7 \quad \bar{y} = 175$$

→ (ממשל) זקנה סוף שיהיה קולר בין
ט"ס"ר הט"ס זקנה הקולר

(*) נזכר בעבר שהייתה גם בט"ס"ר הקצרים. שני רגלים שטענוי היתרונות הקצרים "ים" בעצם משלים מ"מ"ר כי קו הרזותם הפחית את משקלם וזוהי הסיבה שהקו החסון.

מדידת איכות ההתאמה

(*) המדידה בהכרח לא היא המדידה של התאמה.

$$RSS_{LS} = \sum_{i=1}^n (y_i - a_{LS} - b_{LS} x_i)^2$$

ההפרדה המינימלית

(*) גם ט"ס"ר המדידה משלם מ"מ"ר המדידה.

למשל, אחרת: ניקח את a_{LS} כמדידה של Y ונראה שהיא "הכי טובה" (כלומר b_{LS} גדול) ונראה שהיא הטובה ביותר של X על Y .

אם כן, אפקטיות אפילו בקנה שאלה: "מהן תכונות המדידה טובה אפילו בקנה X ו- Y ?"

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

מדידת סטייה: מדד המאפיין את פארסון (Pearson)

$$r(x, y) = \frac{Cov(X, Y)}{\sqrt{Var(X) \cdot Var(Y)}}$$

$$b_{LS} = \frac{Cov(X, Y)}{Var(X)}$$

מחזורי

$$b_{LS} = r_{xy} \cdot \frac{sd(Y)}{sd(X)}$$

גבולות: r

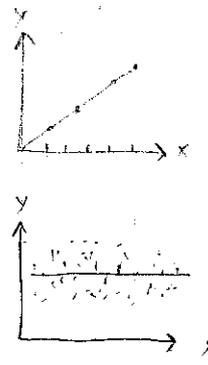
$$r(x, y) = r(y, x) \quad 1$$

$$r(x, a+bx) = r(x, y) \quad \text{אם } b > 0 \quad 2$$

$$r(a+bx, y) = -r(x, y) \quad \text{אם } b < 0 \quad 2$$

$$r(a+bx, y) = \begin{cases} 1 & b > 0 \\ 0 & b = 0 \\ -1 & b < 0 \end{cases} \quad 3$$

$$-1 \leq r(x, y) \leq 1 \quad 4$$



$Y_i = a + bX_i$, X_1, \dots, X_n נתון

⊕ כמות $r=1$ אם יש תלות מושלמת בין הנקודות.

אם $r=0$, $b_{LS} = 0$, הריבוע הריבועי של X הוא Y , כלומר $Y = \text{const}$ (אנטי-קורלציה)

יש להיזהר: אם יש קו שטוח בלבד, זה לא אומר תלות. רק אם $Y = \text{const}$, כלומר התקרה כוללת.

רצף

$\hat{Y}_i = a_{LS} + b_{LS} X_i$ (רגריסיה)

$RSS(a_{LS}, b_{LS}) = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ (ריבועי)

$r^2(x, y) = 1 - \frac{RSS(a_{LS}, b_{LS})}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{\text{Var}(Y) - RSS(a_{LS}, b_{LS})}{\text{Var}(Y)}$

$\text{Var}(Y)$

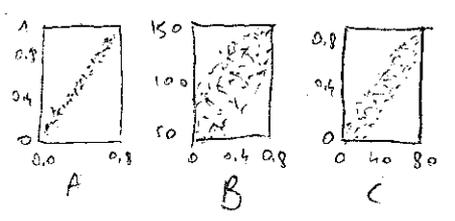
רצף הריבועי המוסכמת

אם $b_{LS} = 0$, כל (אנטי-קורלציה) : $r = r^2 = 0$, כלומר אין תלות בין X ל- Y .

$r^2 = \frac{(\sum (X_i - \bar{X})(Y_i - \bar{Y}))^2}{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2} = \frac{(\sum X_i Y_i - n \bar{X} \bar{Y})^2}{(\sum X_i^2 - n \bar{X}^2)(\sum Y_i^2 - n \bar{Y}^2)}$

הצורה של r^2 היא תמיד בין 0 ל-1. יש 3 מקרים:

- 1. $r^2 = 1$: תלות מושלמת (A, B, C)
- 2. $r^2 = 0$: אין תלות (A, B, C)
- 3. $0 < r^2 < 1$: תלות חלקית (A, B, C)



התבונה
התבונה
התבונה

תבונה ומודעות

r^2 הוא מדד לכמות התלות. אם $r^2 = 1$, זה אומר תלות מושלמת. אם $r^2 = 0$, זה אומר אין תלות.

אם $r^2 = 1$, זה אומר תלות מושלמת. אם $r^2 = 0$, זה אומר אין תלות.

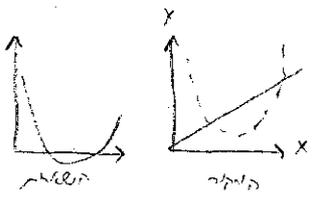
המסקנה היא שהקשר בין המשתנים הוא ליניארי

ניתוח משתנה וטרינספונדנציה ליניארית

$$(\hat{y}_i = a_{LS} + b_{LS} x_i \quad \text{דגם 1}) \quad v_i = y_i - \hat{y}_i$$

אם המשתנה הוא "זמן" או "מרחק" אזי הקשר בין x ו- y קרוב
 ליניארי

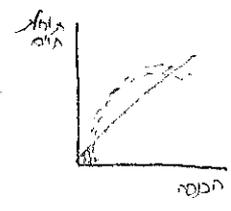
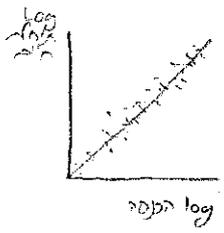
דוגמה $y_i = x_i^2$ (נניח)



אם ננסה להשתמש ב- $\hat{y}_i = a_{LS} + b_{LS} x_i$ (דגם 1) אזי
 המודל לא יתאים טוב

אם ננסה להשתמש ב- $\hat{y}_i = a_{LS} + b_{LS} z_i$ (דגם 2) אזי
 $z_i = x_i^2$ יתאים טוב יותר
 $a_{LS} = 0$ ו- $b_{LS} = 1$

2. Graphinder



←

דגמי תחום (ב) שיעור

אם $p > 0$ אזי y^p הוא פונקציה קמורה (convex) ו- x^p הוא פונקציה קמורה (convex) גם כן.
 אם $p < 0$ אזי y^p הוא פונקציה קעורה (concave) ו- x^p הוא פונקציה קעורה (concave) גם כן.

אם $p > 1$ אזי y^p הוא פונקציה קמורה (convex) ו- x^p הוא פונקציה קמורה (convex) גם כן.

אם $p < 1$ אזי y^p הוא פונקציה קעורה (concave) ו- x^p הוא פונקציה קעורה (concave) גם כן.

דגמי תחום (ב) שיעור

נסתכל על המשתנים $(x_1, y_1), \dots, (x_n, y_n)$ כאשר $x_i \in \mathbb{R}^p$ ו- $p \geq 1$

אנחנו נכתוב: $X_i = (x_{i1}, \dots, x_{ip})$ (וקטור עמודה)

המודל שנתבונן בו הוא: $\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$ (כאשר β_0, \dots, β_p הם הפרמטרים שאנחנו רוצים למצוא)

הפונקציה $RSS(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$: RSS

המטרה היא למצוא את הערכים של β_0, \dots, β_p המינימיים את הפונקציה

(*) $RSS(\beta) = \|y - X\beta\|_2^2$
 (כאשר y הוא וקטור עמודה של y_i , X הוא מטריצה של x_{ij} ו- β הוא וקטור עמודה של β_j)

אם $p > 0$ אזי y^p הוא פונקציה קמורה (convex) ו- x^p הוא פונקציה קמורה (convex) גם כן.
 אם $p < 0$ אזי y^p הוא פונקציה קעורה (concave) ו- x^p הוא פונקציה קעורה (concave) גם כן.

$$y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$$

$$X_{n \times (p+1)} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{p1} & x_{p2} & \dots & x_{pp} \end{pmatrix}$$

הבעיה של ריבועי אטיות (RSS) של $\hat{\beta}$ מוגדרת כ-

$$\frac{\partial}{\partial \beta} \ell(\beta) = 2(X^T y - X^T X \beta) = 0$$

$$\Rightarrow X^T X \hat{\beta} = X^T y \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

← פתרון המערכת
המשוואות
ע"י שימוש
במטריצה
ההפוכה

$$b_{LS} = \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad a_{LS} = \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

כאשר $p=1$ קיטוני קיצוני

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad X^T X = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}$$

$$(X^T X)^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i^2 - n\bar{x} & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix} \quad X^T y = \begin{pmatrix} n\bar{y} \\ \sum x_i y_i \end{pmatrix}$$

$$(X^T X)^{-1} X^T y = \frac{1}{n \sum (x_i - \bar{x})^2} \begin{pmatrix} n\bar{y} \sum x_i^2 - n\bar{x} \sum x_i y_i \\ -n^2 \bar{y} \bar{x} + n \sum x_i y_i \end{pmatrix}$$

$$\hat{\beta}_0 = ?$$

$$\hat{\beta}_1 = \frac{n \sum x_i y_i - n^2 \bar{y} \bar{x}}{n \sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \checkmark$$

107) הרחבה בהצגה

נתון הרגסיה עם משתנים קטגוריים
 X המשתנה - קיצוני משתנים
 y קטגוריים - תמיד שקוף (קואסי-פיקטיב), בעזרתם רואים
 מה המשמעות של הרגסיה כמאפיין

מכאן הרגסיה איננה
 הרגסיה של אינטיות
 שניתן לקבלתם שחיים - RSS
 תמיד קטנים שרובם בעיטת התקנות

הסקה סטטיסטית על תוצאות הרגסיה

תכונות פונקציית התפלגות

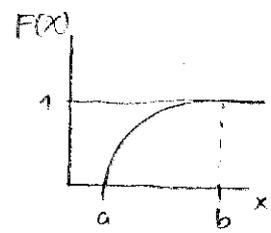
לכל התפלגות נתונה

$\forall \omega \in \Omega \quad F(\omega) \geq 0 \quad \sum_{\omega \in \Omega} F(\omega) = 1$ פונקציית התפלגות F בתחום $\Omega \subset \mathbb{R}$

התפלגות הדיסקרטית, נקראת $\Omega \subset \mathbb{R}$ אם היא מתחילה בהתפלגות הדיסקרטית (אם לא):

$$\left. \begin{array}{l} \Omega \text{ הוא סדר} \\ \text{הוא } \mathbb{R} \\ \text{בין } a \text{ ו- } b \end{array} \right\} \begin{array}{l} \Omega = \mathbb{R} \quad \text{א} \\ \Omega = (-\infty, 0] \quad \text{ב} \quad \Omega = [0, \infty) \quad \text{ג} \\ \mathbb{R} \ni a < b \in \mathbb{R} \quad \Omega = [a, b] \quad \text{ד} \end{array}$$

הפונקציית התפלגות F מתאפיינת על ידי:



$F(a) = 0$ א

$F(b) = 1$ ב

הפונקציית התפלגות F ד

אם $\forall x \in \Omega$ אז $(X \sim F)$ פונקציית התפלגות F של X מתאפיינת על ידי $P_X(X \leq x) = F(x)$

$Pr(a \leq X \leq b) = F(b) - F(a)$ עבור $a, b \in \Omega$

הפונקציית התפלגות F של X מתאפיינת על ידי $F(x) = Pr(X \leq x)$

$f_X(x) = \frac{dF_X(x)}{dx}$ פונקציית התפלגות

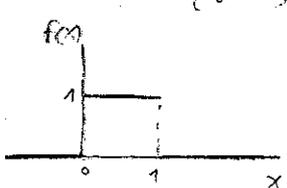
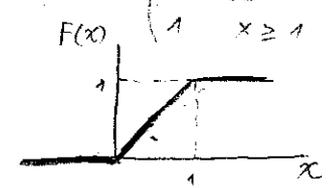
תכונות

$\forall x \in \Omega \quad f_X(x) \geq 0$ א

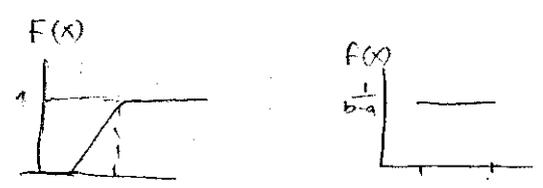
$\int_a^b f(x) dx = F(x) \Big|_a^b = 1$ ב

דוגמה

$X \sim U(0,1)$ פונקציית התפלגות $F_X(x) = \begin{cases} 0 & x \leq 0 \\ x & 0 \leq x \leq 1 \\ 1 & x \geq 1 \end{cases}$ פונקציית התפלגות $f_X(x) = \begin{cases} 1 & 0 \leq x \leq 1 \\ 0 & \text{אחרת} \end{cases}$ $\Omega = [0,1]$



התפלגות $[a,b]$ פונקציית התפלגות $F_X(x)$



$F_X(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & a \leq x \leq b \\ 1 & x > b \end{cases}$ $f_X(x) = \begin{cases} 0 & x < a \\ \frac{1}{b-a} & a \leq x \leq b \\ 0 & x > b \end{cases}$