

Kétváltozós lineáris és lineárisra visszavezethető regresszió

Kétváltozós lineáris regressziós modell

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

β_0, β_1 becslése legkisebb négyzetek módszerével:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n d_{x_i} d_{y_i}}{\sum_{i=1}^n d_{x_i}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i y_i - n \cdot \bar{x} \cdot \bar{y}}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Varianciafelbontás a kétváltozós lineáris modellre:

$$SST = SSR + SSE$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum d_y^2 \quad \text{teljes négyzetösszeg}$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2 \quad \text{belső négyzetösszeg, a hiba okozta (reziduális) négyzetösszeg}$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad \text{külső négyzetösszeg, regressziós vagy magyarázott négyzetösszeg}$$

A *determinációs együttható* azt mutatja, hogy a regressziós modellel az y_i adatokban meglévő variancia hány százaléka szüntethető/magyarázható meg:

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

$$R_{\text{adjusted}}^2 = 1 - \frac{n-1}{n-k-1}(1-R^2)$$

$R^2 \approx 1$ - jó illeszkedés, nagy magyarázó erő

$R^2 \approx 0$ - gyenge modelteljesítmény.

A mintából számolt becslt *lineáris korrelációs együttható* a magyarázó- és eredményváltozó között:

$$r = \frac{\sum d_x d_y}{\sqrt{\sum d_x^2 \cdot \sum d_y^2}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - n \bar{x}^2\right) \cdot \left(\sum_{i=1}^n y_i^2 - n \bar{y}^2\right)}} = \hat{\beta}_1 \frac{s_x}{s_y}$$

CSAK kétváltozós lineáris esetben: $R^2 = r^2$.

Az *elaszticitás* (rugalmasság) azt méri, hogy az X változó 1%-os növekedése hány százalékos növekedést/csökkenést eredményez az Y változónál. Az elaszticitás kiszámítása a becslt eredményváltozóra:

$$El(\hat{y}, x) = \frac{\partial \hat{y}}{\partial x} \cdot \frac{x}{\hat{y}}$$

Kétváltozós lineáris esetben:

$$El(\hat{y}, x) = \hat{\beta}_1 \cdot \frac{x}{\hat{y}} = \frac{\hat{\beta}_1 x}{\hat{\beta}_0 + \hat{\beta}_1 x}$$

Intervallumbecslés a függvényértékekre:

- az átlagos értékre

$$Int_{1-\alpha}(E(Y_*)) = \hat{y}_* \pm t_{1-\frac{\alpha}{2}}(n-2)s_e \sqrt{\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum d_x^2}},$$

- az egyedi értékre

$$Int_{1-\alpha}(Y_*) = \hat{y}_* \pm t_{1-\frac{\alpha}{2}}(n-2)s_e \sqrt{1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum d_x^2}},$$

ahol $s_e = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-2}}$ a *korrigált reziduális szórás*.

1. A hektáronkénti szőlőtermést befolyásolja az évenkénti permetezések száma. Az alábbi táblázat öt év adatait mutatja:

Permetezések száma (db)	Szőlőtermés (q)
5	10
6	10
7	12
7	10
5	8
Σ	Σ

- (a) Határozza meg és értelmezze a lineáris regresszió paramétereit!

	$d_x = x_i - \bar{x}$	$d_y = y_i - \bar{y}$	d_x^2	d_y^2	$d_x d_y$	\hat{y}	$e_i = y_i - \hat{y}$
1							
2							
3							
4							
5							
Σ							

- (b) Számítsa ki a két változó lineáris korrelációs együtthatóját!
 - (c) Számítsa ki és értelmezze a lineáris regresszió determinációs együtthatóját!
 - (d) Számítsa ki és értelmezze a korrigált reziduális szórást!
 - (e) Számolja ki és értelmezze az elaszticitást az átlagos permetezésszáma!
 - (f) Adjon becslést 6 permetezés esetén a szőlőtermés átlagos mennyiségére, majd szerkesszen konfidenciaintervallumot ugyanerre 95%-os megbízhatósági szinten!
 - (g) Adjon 95%-os konfidenciaintervallumot egy 6-szor permetezett szőlőültetvény szőlőtermésére!
2. Egy bank 10 ügyfelét vizsgálva az életkor (X , év) és a havi jövedelem (Y , eFt-ban) kapcsolatát elemzi. Az alábbi részeredményeket kapták:

$$\sum_{i=1}^n x_i = 260, \quad \sum_{i=1}^n y_i = 2040, \quad \sum_{i=1}^n x_i y_i = 53\,754,$$

$$\sum_{i=1}^n x_i^2 = 6\,862, \quad \sum_{i=1}^n y_i^2 = 421\,866, \quad \sum_{i=1}^n e_i^2 = 708,$$

$$\sum_{i=1}^n \log x_i = 50, \quad \sum_{i=1}^n \log y_i = 90, \quad \sum_{i=1}^n x_i \log y_i = 675.$$

- Határozza meg és értelmezze a lineáris regresszió paramétereit!
 - Számítsa ki a két változó lineáris korrelációs együtthatóját!
 - Számítsa ki és értelmezze a lineáris regresszió determinációs együtthatóját!
 - Adjon becslést a 30 évesek átlagos jövedelmére, majd szerkesszen konfidenciaintervallumot ugyanerre 98%-os megbízhatósági szinten!
 - Adjon 98%-os konfidenciaintervallumot egy 30 éves ügyfél egyedi jövedelmére!
 - Számolja ki és értelmezze az elaszticitást az átlagos életkorra!
3. Egy taxivállalat 15 véletlenszerűen kiválasztott fuvar alapján vizsgálja, hogy hogyan függ a menetidő a távolságtól (megtett km-től). A 15 fuvar esetén a távolság és a menetidő:

távolság (km)	menetidő (perc)	távolság (km)	menetidő (perc)
3	8	9	20
4	19	12	23
4	13	15	44
6	21	16	47
6	11	16	41
7	19	20	46
8	14	26	48
8	19		

$$\sum y = 393, \quad \sum xy = 5\,433, \quad \sum x \ln y = 545.8033,$$

$$\sum \ln y = 46.7381, \quad \sum x = 160, \quad \sum x^2 = 2\,328,$$

$$\sum \ln x \ln y = 106.6887, \quad \sum (\ln x)^2 = 77.2063, \quad \sum \ln x = 32.7487.$$

- Jellemezze a távolság és a menetidő közötti lineáris, exponenciális illetve hatvány kapcsolatot és értelmezze a paramétereiket!
- Becsülje meg mindegyik modell alapján, hogy egy 15 km távolságú út hány percet vesz igénybe!
- Mekkora az elaszticitás az átlagos és a 15 km távolságú fuvar környezetében?

SPSS: Graphs \Rightarrow Scatter/Dot

Analyze \Rightarrow Regression \Rightarrow Curve Estimation: Linear, Compound, Power

Transform \Rightarrow Compute variable: becslés, hiba

Analyze \Rightarrow Regression \Rightarrow Linear