

Lineární regrese – Stanovení rovnice kalibrační křivky

Napište rovnici kalibrační křivky průmyslového analyzátoru se selektivní elektrodou. Měřená veličina je koncentrace rozpuštěného amoniaku na výstupu z absorpční kolony. Za předpokladu linearity napěťové odezvy elektrody na koncentraci amoniaku rozhodněte, zda

- (A) kalibrační křivka prochází počátkem;
- (B) mohlo dojít při měření k systematické chybě.
- (C) Stanovte predikční pás pro vypočtenou regresní přímku (hlinu významnosti volte 95%).
- (D) Rozhodněte, zda na této hlině významnosti lze určit koncentraci amoniaku $c = 50 \text{ mg/l}$ s přesností 2%.
- (E) Určete interval spolehlivosti pro směrnici regresní přímky. Může být $\beta_2 = 0$?

Řešení – teorie

Z teorie často víme, že naměřená data vykazují funkcionální závislost, ovšem zatíženou chybou měření. V našem příkladu je tato závislost lineární. Chtěli bychom najít takovou přímku, která „nejlépe“ (v jakém smyslu?) vystihuje naše data měřená s náhodnou chybou. Tuto úlohu lze popsat pomocí tzv. **lineární regrese**

$$Y(x) = \beta_1 + \beta_2 x + \epsilon(x),$$

kde

- pro každé x (nenáhodné!) z podmnožiny M reálných čísel je $Y(x)$ náhodná veličina,
- $\beta_1, \beta_2 \in \mathbb{R}$ jsou parametry,
- pro každé $x \in M$ je $\epsilon(x)$ náhodná veličina (chyba měření).

Ona funkcionální závislost se nazývá **regresní funkce** $\eta(x)$ a v lineárním případě je definovaná vztahem

$$\eta(x) = \mathbb{E}[Y(x)] = \beta_1 + \beta_2 x.$$

Sestavme tzv. **lineární regresní model**:

$$Y_j = \beta_1 + \beta_2 x_1 + \epsilon_j, \quad j = 1, \dots, n,$$

kde pro hodnoty nezávislé proměnné $x_j, j = 1, \dots, n$, jsme označili

$$Y_j \equiv Y(x_j), \quad \epsilon_j \equiv \epsilon(x_j), \quad j = 1, \dots, n.$$

Jeho **předpoklady** jsou

- (i) $n > 2$ a existují $i, j \in \{1, \dots, n\}$ taková, že $x_i \neq x_j$,
- (ii) $\mathbb{E} Y_j = \eta(x_j) = \beta_1 + \beta_2 x_j, \quad j = 1, \dots, n$,

- (iii) $\text{var } Y_j = \sigma^2 > 0, j = 1, \dots, n,$
- (iv) $\text{cov}(Y_i, Y_j) = 0, i, j = 1, \dots, n, i \neq j.$

Poznámka 1. Podmínky (ii)–(iv) lze přepsat v řeči chyb $\epsilon(x_j)$ následovně:

$$\mathbb{E} \epsilon_j = 0, \text{var } \epsilon_j = \sigma^2, j = 1, \dots, n, \text{cov}(\epsilon_i, \epsilon_j) = 0, i, j = 1, \dots, n, i \neq j.$$

Bodové odhady parametrů modelu

Model má tři parametry $\beta_1, \beta_2, \sigma^2$, pro něž bychom chtěli nalézt bodové odhady

$$\hat{\beta}_1 = b_1, \hat{\beta}_2 = b_2, \hat{\sigma}^2 = s^2.$$

- Pro odhady β_1, β_2 se standardně používá **metoda nejmenších čtverců** – hledáme b_1, b_2 tak, aby platilo

$$\sum_{j=1}^n (Y_j - (b_1 + b_2 x_j))^2 = \min_{\beta_1, \beta_2 \in M} \sum_{j=1}^n (Y_j - (\beta_1 + \beta_2 x_j))^2.$$

Je to aplikační úloha na lokální extrémy funkcí dvou proměnných (viz Matematiku II) a nalezená řešení b_1, b_2 jsou nejlepšími nestrannými lineárními odhady parametrů β_1, β_2 .

Poznámka 2. Bud' X_1, \dots, X_n náhodný výběr z rozdělení náhodné veličiny X . Bodový odhad $\hat{\theta}_N \equiv \hat{\theta}_N(X_1, \dots, X_n)$ parametru θ tohoto rozdělení je **nestranný**, jestliže

$$\mathbb{E} \hat{\theta}_N = \theta.$$

Nestranný bodový odhad $\hat{\theta}_{NN}$ parametru θ je **nejlepší nestranný odhad**, jestliže

$$\text{var } \hat{\theta}_{NN} = \min_{\hat{\theta}_N} \text{var } \hat{\theta}_N.$$

Bodový odhad $\hat{\theta}$ je **lineární**, pokud existují konstanty $a_0, a_1, \dots, a_n \in \mathbb{R}$ tak, že

$$\hat{\theta} = a_0 + \sum_{i=1}^n a_i X_i.$$

Ona řešení β_1, β_2 lze vyjádřit ve tvaru

$$b_1 = \bar{Y} - b_2 \bar{x}, \tag{1}$$

$$b_2 = \frac{n \sum_{i=1}^n x_i Y_i - \sum_{i=1}^n x_i \sum_{i=1}^n Y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}, \tag{2}$$

kde

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

- Bodové odhady chyb $\hat{\epsilon}_j$, $j = 1, \dots, n$, získané na základě metody nejmenších čtverců

$$\hat{\epsilon}_j = Y_j - b_1 - b_2 x_j$$

se nazývají **rezidua** a součet jejich kvadrátů (značí se S_e)

$$S_e = \sum_{j=1}^n \hat{\epsilon}_j^2 = \sum_{j=1}^n (Y_j - b_1 - b_2 x_j)^2$$

pak **reziduální součet čtverců**. Snadnou úpravou lze získat tvar vhodný pro výpočty

$$S_e = \sum_{j=1}^n Y_j^2 - b_1 \sum_{j=1}^n Y_j - b_2 \sum_{j=1}^n x_j Y_j.$$

Využitím reziduálního součtu čtverců obdržíme nestranný odhad rozptylu

$$s^2 = \hat{\sigma}^2 = \frac{1}{n-2} S_e,$$

který se nazývá **reziduální rozptyl**.

Intervaly spolehlivosti – pás spolehlivosti a predikční pás

Pokud bychom chtěli konstruovat intervaly spolehlivosti pro parametry regresního modelu, potřebujeme dodatečné předpoklady na chyby:

- (i) ϵ_j , $j = 1, \dots, n$, jsou vzájemně nezávislé náhodné veličiny,
- (ii) ϵ_j , $j = 1, \dots, n$ mají normální rozdělení $\mathcal{N}(0, \sigma^2)$.

Pak $(1 - \alpha)100\%$ oboustranný interval spolehlivosti pro

- β_1 je

$$[b_1 - t_{1-\frac{\alpha}{2}}(n-2)s_{b_1}, b_1 + t_{1-\frac{\alpha}{2}}(n-2)s_{b_1}],$$

- β_2 je

$$[b_2 - t_{1-\frac{\alpha}{2}}(n-2)s_{b_2}, b_2 + t_{1-\frac{\alpha}{2}}(n-2)s_{b_2}], \quad (3)$$

- regresní funkci $\eta(x)$ je

$$[b_1 + b_2 x - t_{1-\frac{\alpha}{2}}(n-2)s_{\hat{\eta}(x)}, b_1 + b_2 x + t_{1-\frac{\alpha}{2}}(n-2)s_{\hat{\eta}(x)}],$$

kde

$$\begin{aligned} s_{b_1}^2 &= s^2 \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}, \\ s_{b_2}^2 &= s^2 \frac{n}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}, \\ s_{\hat{\eta}(x)}^2 &= s^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \right), \end{aligned}$$

a $t_{1-\frac{\alpha}{2}}(n-2)$ je $(1 - \frac{\alpha}{2})$ 100% kvantil t-rozdělení s $n-2$ stupni volnosti.

Jestliže se díváme na horní a dolní meze intervalu spolehlivosti pro regresní funkci jako na funkce proměnné x , pak plocha ohraničená grafy těchto funkcí se nazývá **pás spolehlivosti kolem regresní přímky**. Obdobně lze zkonstruovat tzv. **predikční pás kolem regresní přímky**, což je plocha ohraničená mezemi $(1 - \alpha)100\%$ intervalů zkonstruovaných přímo pro náhodnou veličinu $Y(x)$

$$[b_1 + b_2 x - t_{1-\frac{\alpha}{2}}(n-2)s_{\hat{Y}(x)}, b_1 + b_2 x + t_{1-\frac{\alpha}{2}}(n-2)s_{\hat{Y}(x)}],$$

kde

$$s_{\hat{Y}(x)}^2 = s^2 \left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \right) = s^2 + s_{\eta(x)}^2.$$

Vlastní řešení příkladu

Z vzorců (1), (2) získáme rovnici kalibrační křivky $y = b_1 + b_2 x$.

Pro zodpovězení úlohy (A) chceme otestovat, zdali regresní přímka prochází počátkem, tj. testujeme

$$H_0 : \beta_1 = 0, \quad H_1 : \beta_1 \neq 0.$$

Predikční pás kolem regresní přímky (úloha (C)) lze zkonstruovat pomocí softwaru. Pak můžeme pohledem na obrázek zodpovědět otázku (B), zdali mohlo dojít k systematické chybě měření. Pokud 95% dat leží v predikčním pásu, lze usoudit, že data nevykazují systematickou chybu v měření.

V úloze (D) jde o to posoudit, zdali je graf regresní funkce $y = b_1 + b_2 x$ pro $x \in [50 - 0,02 \cdot 50; 50 + 0,02 \cdot 50]$ obsažený v pásu spolehlivosti. Pokud tomu tak je, lze určit koncentraci 50 mg/l s přesností 2%.

Interval spolehlivosti z úlohy (E) sestavíme podle vzorce (3). To, co chceme otestovat, je

$$H_0 : \beta_2 = 0, \quad H_1 : \beta_2 \neq 0,$$

a to na základě právě zkonstruovaného intervalu spolehlivosti. Pokud je 0 mimo interval, nulovou hypotézu zamítneme, a tak je nenulovost směrnice statisticky významná. V opačném případě bychom pouze mohli konstatovat, že nenulovost směrnice není významná, což ovšem neznamená, že by její nulovost významná byla!

Poznámka 3. Je vhodné se alespoň některou z grafických metod (např. Q-Q plot; pro histogram je v našem příkladě málo dat) přesvědčit, že by naměřené chyby ϵ_j , $j = 1, \dots, n$, mohly pocházet z normálního rozdělení. Též je užitečné si hned na začátku udělat první představu o povaze dat pomocí box-plotu a zkusit vykreslit graf reziduí, zdali se v datech nevyskytuje systematická chyba.