

Linearna regresia

V texte su vynechane vektorove znacenia. Veliciny θ, y, ϵ su vektorove veliciny.

Linearna regresia je odhad vektoru parametrov θ pomocou linearného modelu $m = \mathbf{A}\theta$, o ktorom dúfame, že pre odhadnute $\hat{\theta}$ bude $m \approx y$, kde y su data. \mathbf{A} je modelova matica a pre tento problem ma tvar:

$$\mathbf{A} = \begin{pmatrix} 1 & x_0^2 \\ \vdots & \vdots \\ 1 & x_N^2 \end{pmatrix} \quad (1)$$

V linearnom prípade su odhady parametrov vypočítané ako:

$$\hat{\theta} = (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} y, \quad (2)$$

kde \mathbf{W} je vahova matica. V tomto probleme sa neuplatní, bude identita, keďže body majú rovnakú váhu. Vzorec je možné odvodiť z minimalizácie $g = \sum_{i=1}^N (y_i - m_i) w_{ij} (y_j - m_j)$.

Chyby a korelácie parametrov su dane kovariančnou maticou:

$$c_{ij} \equiv \text{Cov}(x_i, x_j) = E((x_i - E(x_i))(x_j - E(x_j))) = E(x_i x_j) - E(x_i)E(x_j) \quad (3)$$

Oznacím:

$$\xi_{ij} \equiv \left((\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \right)_{ij} \quad (4)$$

potom dosadením odhadov $\hat{\theta}$ (uz s $W = I$):

$$\text{Cov}(\hat{\theta}_i, \hat{\theta}_j) = E(\xi_{ki} y_i \xi_{lj} y_j) - E(\xi_{ki} y_i) E(\xi_{lj} y_j) = \quad (5)$$

$$\xi_{ki} \text{Cov}(y_i, y_j) \xi_{jl} \quad (6)$$

Predpokladáme, že $\mathbf{C}(y, y)$ je diagonálna, ale s neznámym rozptylom (rovnaký pre všetky body). Nevychýlený odhad rozptylu je:

$$\sigma_y^2 \approx s_y^2 = \frac{\epsilon_i \epsilon_i}{N - p}, \quad (7)$$

kde N je počet bodov y a p je počet fitovaných parametrov.

Intervalovy odhad

Odhadované parametre su tiež náhodne premenne. Rozdelenie ϵ je náhodne s nulovým prvým momentom a rozptylom σ_y^2 . Z toho vyplýva, že rozdelenie $\hat{\theta} - \theta$ je tiež normálne, keďže je dane lineárnym vzťahom z ϵ . Rozdelenie $\frac{\epsilon_i \epsilon_i}{\sigma_y^2} \equiv R$ je χ^2 z definície χ^2 . Zaujima nás, z akeho rozdelenia bude θ . Odhadovaný rozptyl parametrov $s_{\hat{\theta}}^2$ je daný diagonálou $\mathbf{C}(\hat{\theta}, \hat{\theta})$ a súvisí s odhadnutým rozptylom s_y^2 :

$$s_{\hat{\theta}i}^2 = \text{Cov}(\hat{\theta}_i, \hat{\theta}_i) = \xi_{ik} \text{Cov}(y_k, y_l) \xi_{li} = s_y^2 \xi_{ik} \xi_{ki} \quad (8)$$

Potom vyraz:

$$t_i = \frac{\theta_i - \hat{\theta}_i}{s_{\theta i}} = \frac{\theta_i - \hat{\theta}_i}{\xi_{ik} \xi_{ki} s_y} = \frac{\theta_i - \hat{\theta}_i}{\xi_{ik} \xi_{ki}} \sqrt{\frac{(N-p)\sigma_y^2}{R}} \quad (9)$$

je zo Studentovho t -rozdelenia a je možné ho použiť pre konštrukciu intervalu spoľahlivosti. Konkrétne:

$$\Pr(-\tau < t < \tau) = 1 - \alpha, \quad (10)$$

kde τ je hodnota (kvantil), pre ktorú $T(\tau) = 1 - \frac{\alpha}{2}$, kde T je (kumulatívna) Studentova distribučná funkcia. Vyraz je možné pretvoriť na tvrdenie o θ :

$$\Pr(\hat{\theta} - s_{\theta}\tau < \theta < \hat{\theta} + s_{\theta}\tau) = 1 - \alpha, \quad (11)$$

čo by malo byť ekvivalentné tvrdeniu, že θ ležia s pravdepodobnosťou $1 - \alpha$ v intervale:

$$\theta \in [\hat{\theta} - s_{\theta}\tau, \hat{\theta} + s_{\theta}\tau] \quad (12)$$

Pas spoľahlivosti

Pre pas spoľahlivosti okolo celej krivky existuje podobný postup ako vyššie (podrobne v (Casella 2002) a (Michael H Kutner 2005), alebo tiež).

$$\text{Var}(a_{ij}\theta_j) = a_{ij}a_{ik}\text{Cov}(\theta_j, \theta_k) \quad (13)$$

alebo v maticovom zápise:

$$s_m^2 = \text{diag}(\mathbf{AC}(\theta, \theta)\mathbf{A}^T) \quad (14)$$

Znova potrebujeme najst rozdelenie vyrazu:

$$f_i = \frac{A\theta - A\hat{\theta}}{s_m} \quad (15)$$

pre vycislenie tvrdenia:

$$\Pr(\mathbf{A}\hat{\theta} - s_m\phi < \mathbf{A}\theta < \mathbf{A}\hat{\theta} + s_m\phi \text{ for all } x) = 1 - \alpha \quad (16)$$

Ukazuje sa, ze ϕ^2 je z F -rozdelenia:

$$\phi = \sqrt{2F_{\alpha;p,N-p}} \quad (17)$$

Intervalovy odhad chyby

Ako bolo pouzite vyssie rozptyl data je odhadnuty z:

$$\sigma_y^2 \approx s_y^2 = \frac{\epsilon_i \epsilon_i}{N - p} \quad (18)$$

a vyraz:

$$c = \frac{\epsilon_i \epsilon_i}{\sigma_y^2} = \frac{s_y^2(N - p)}{\sigma_y^2} \quad (19)$$

je z χ^2 -rozdelenia ($N - p$ stupnov volnosti). To znamena, ze:

$$\Pr(\sigma_l \leq \sigma_y \leq \sigma_u) = 1 - \alpha \quad (20)$$

$$\Pr(\sigma_l^2 \leq \sigma_y^2 \leq \sigma_u^2) = 1 - \alpha \quad (21)$$

$$\Pr\left(\frac{(N - p)s_y^2}{\sigma_l^2} \geq \frac{(N - p)s_y^2}{\sigma_y^2} \geq \frac{(N - p)s_y^2}{\sigma_u^2}\right) = 1 - \alpha \quad (22)$$

Krajne hodnoty teda budu:

$$\sigma_u = \sqrt{\frac{(N - p)s_y^2}{c_{\alpha/2}}} \quad (23)$$

$$\sigma_l = \sqrt{\frac{(N - p)s_y^2}{c_{1-\alpha/2}}}, \quad (24)$$

kde c_β je definovana ako:

$$\beta = \int_0^{c_\beta} \chi^2(x) dx \quad (25)$$

Casella, George. 2002. *Statistical Inference*. Australia Pacific Grove, CA: Thomson Learning.

Michael H Kutner, John Neter, Christopher J. Nachtsheim. 2005. *Applied Linear Statistical Models*. 5th ed. The McGraw-Hill/Irwin Series Operations and Decision Sciences. McGraw-Hill Irwin.