

20. Линейен регресионен модел. Метод на най-малките квадрати. Теорема на Гаус-Марков.

* Линейен регресионен модел.

* Регресионен експеримент. - Данните от проведеното изследване или експеримент са наблюдения на

- отклик y и
- предиктор $x = (x_{i1}, x_{i2}, \dots, x_{in})'$

Предикторите описват условията, при които се провежда отделно наблюдение. Задачата е да предскажем резултата от наблюдение на количествена променлива, наричана отклик, при известни условия на наблюдение. Едновременно да установим зависимостта на отклика от конкретните стойности на предикторите.

* Регресионен модел. - Регресионен модел наричат y и x , който разпределят зависимостта от предикторите, обяснявателната променлива отклик и удовлетворява условията:

- $E(y|x) = f(x)$ - математическо очакване на y е ф-я на предикторите, $f(x)$ - наричат регресионна ф-я
- $\text{Var } y = \sigma^2$ - дисперсията на y е константа (не зависи от x)

Разликата $e = y - f(x)$ е сл.в. и се интерпретира като случайна грешка при наблюдаване на отклика, чийто "настинска" стойност е $f(x)$. Разпределението на грешката е с очакване $Ee = 0$ и $\text{Var } e = \sigma^2$

Регресионния модел се представя из регресионно $y-e$:

$$y = f(x) + e \quad (Ee = 0, \text{Var } e = \sigma^2)$$

* Линейен регресионен модел - Регресионен модел е линейен, когато функцията $f(x)$ е зададена с точности до неизвестен параметър $\beta' = (\beta_1, \dots, \beta_p)$ и зависи от него параметър линейно.

$$f(x) = \beta_1 f_1(x) + \dots + \beta_p f_p(x) = \sum_{j=1}^p \beta_j f_j(x)$$

За кратко писан

$$f(x) = \beta_1 x(1) + \dots + \beta_p x(p) = \sum_{j=1}^p x(j) \beta_j = x' \beta.$$

Линейния регресионен модел се представя от линейно регресионно уравнение.

$$y = \beta_1 x(1) + \dots + \beta_p x(p) + e \quad (E e = 0 \quad \text{Var } e = \sigma^2),$$

линейно пряко неизвестен параметър β_j .

* Метод на най-малки квадрати

Все регресионни експерименти е проведен n пъти при различни условия

$x_i = (x_i(1), \dots, x_i(p))'$ $i=1, \dots, n$ и наблюдаваните стойности на отклика са били y_i , $i=1, \dots, n$.

Наблюденията удовлетворяват n регресионни уравнения

$$y_i = \beta_1 x_i(1) + \dots + \beta_p x_i(p) + e_i \quad i=1, \dots, n.$$

Данните записване в матричен вид, като наблюденията на V променлива е n -мерен вектор

• $y = (y_1, \dots, y_n)'$ са наблюденията на отклика, а

• $x(j) = (x_1(j), \dots, x_n(j))'$ са наблюденията на j -тия от предикторите ($j=1, \dots, p$).

Матрицата на експерименталните данни матрицата с n редове са наблюдавания на p -те предиктора а редове са транспонирани векторите на условията на наблюдения.

$$x_i' = (x_i(1), \dots, x_i(p)) \quad i=1, \dots, n$$

$$X(n \times p) = (x_{(1)} \dots x_{(p)}) = \begin{pmatrix} x'_{11} \\ \vdots \\ x'_{1n} \end{pmatrix} \text{ Регресионните уравнения}$$

за n -те наблюдения в матричен запис са:

$$y_{(n \times 1)} = X(n \times p) \cdot b_{(p \times 1)} + e_{(n \times 1)} \text{ или накратко } Y = Xb + e$$

Със e_i - нецелесителната "грешка" при i -томо наблюдение
и $e' = (e_1, e_2, \dots, e_n)$

В системата от n уравнения нецелесителите са $p+n$ на бр-параметрите b и грешката e .

Задачата е да намерим ~~грешката~~ в каноничен подходен и разумен ограничен изглед грешките e .

Гаус разширил метод на най-малките квадрати като определя условие за грешките.

Нецелесителните параметри се определят така, че сумата от квадратите на грешките да е минимална!

Разглеждаме $y_i - \sum_{j=1}^p x_{ij} b_j = e_i, i=1, \dots, n$ наричан още остатък и сумата от квадратите им е

$$SS(b) = \sum_{i=1}^n e_i^2 = e'e.$$

Теорема 1. Ако \hat{b} е решение на системата

$$X'Xb = X'Y \quad \text{— наричана системата от нормални уравнения, то за } \hat{b} \text{ е достигната минимална}$$

сумата от квадрати $SS(b)$:

$$SS(b) \geq SS(\hat{b})$$

Забелешка 1. Нормалната система уравнения е друга форма на запис на условието за некорелираност на e и x_j $j=1 \dots p$. т.е.

$$x'e = x'(y - x\beta) = x'y - x'x\beta = 0$$

Забелешка 2. Нормалната система уравнения е еквивалентна на анулиране на първите производни на $SS(\beta)$ по параметра β .

$$\begin{aligned} \frac{\partial SS}{\partial \beta} &= \frac{\partial e'e}{\partial \beta} = \frac{\partial (y - x\beta)'(y - x\beta)}{\partial \beta} = \frac{\partial (y'y - 2\beta'x'y + \beta'x'x\beta)}{\partial \beta} \\ &= -2x'y + 2x'x\beta. \end{aligned}$$

Директно доказателство следва от равенствата

$$\begin{aligned} SS(\beta) &= (y - x\beta)'(y - x\beta) = (y - x\hat{\beta} + x\hat{\beta} - x\beta)'(y - x\hat{\beta} + x\hat{\beta} - x\beta) \\ &= (y - x\hat{\beta})'(y - x\hat{\beta}) + (x\hat{\beta} - x\beta)'(x\hat{\beta} - x\beta) + 2(x\hat{\beta} - x\beta)'(y - x\hat{\beta}) \\ &= SS(\hat{\beta}) + (x\hat{\beta} - x\beta)'(x\hat{\beta} - x\beta) \quad \text{тъй като} \\ &\quad (x\hat{\beta} - x\beta)'(y - x\hat{\beta}) = 0, \text{ а} \end{aligned}$$

$$(x\hat{\beta} - x\beta)'(y - x\hat{\beta}) = (\hat{\beta} - \beta)'x'(y - x\hat{\beta}) = (\hat{\beta} - \beta)'(x'y - x'x\hat{\beta}) = 0$$

Нормалната система уравнения има единствено решение, ако ранг на X е равен на p , т.е.

стълбовете x_j са ЛНЗ (линейно независими). Решението на системата е търсената оценка: $\hat{\beta} = (X'X)^{-1}X'y$.

Нормално разпределение на грешките.

Жела за грешките предположим, че са независими и нормално разпределени.

$$e_i \sim N(0, \sigma^2), i=1, \dots, n \quad \text{или} \quad e \sim N(0, \sigma^2 I)$$

Тогава нормално е разпределението и на наблюденията $y \sim N(X\beta, \sigma^2 I)$ и оценките получени по МНК съвпадат

с оценките по метода на максималното правдоподобие и прилагаме всички техники оптимални за тях.

Теорема на Гаус-Марков.

Нека неизвестните параметри се подчиняват на условията:

$$E \epsilon_i = 0, i = 1, \dots, n;$$

$$\text{Var } \epsilon_i = \sigma^2 > 0, i = 1, \dots, n$$

$\text{cov}(\epsilon_i, \epsilon_j) = 0, i \neq j$ и нека $\hat{\beta}$ е решение на системата нормални уравнения

$$X'X\hat{\beta} = X'y$$

Тогава $\hat{\beta}$ е BLUE:

1. линейна по случайния вектор y : $\hat{\beta} = Cy$ и

2. неизместен $E\hat{\beta} = \beta$

3. с минимална дисперсия - за всяка друга линейна по y неизместенa оценка $\tilde{\beta}$ е изпълнено неравенството:

$$\text{Var } \hat{\beta} \leq \text{Var } \tilde{\beta}$$

Неравенството е в смисъл, че разликата на две ковариационни матрици

$D = \text{Var } \tilde{\beta} - \text{Var } \hat{\beta}$ е неотрицателно определена матрица за всяко x $x'Dx \geq 0$

Доказателство: Последователно доказваме трите свойства на $\hat{\beta}$:

1. $\hat{\beta}$ е линейна по вектора y

$$\hat{\beta} = (X'X)^{-1} X'y = Cy \text{ като } C = (X'X)^{-1} X';$$

2. $\hat{\beta}$ е неизместенa оценка на β

$$E\hat{\beta} = E(X'X)^{-1} X'y = (X'X)^{-1} X'Ey = (X'X)^{-1} X'XB = \beta$$

3. с минимална дисперсия

Дисперсията на $\hat{\beta}$ е:

$$\begin{aligned} \text{Var } \hat{\beta} &= \text{Var}(C \cdot y) = C \cdot (\text{Var } y) \cdot C' = C \sigma^2 I_n \cdot C' = \\ &= \sigma^2 (C' C)^{-1} (C' C)^{-1} = \sigma^2 (X' X)^{-1} X' X (X' X)^{-1} = \\ &= \sigma^2 (X' X)^{-1} \end{aligned}$$

Преди да изследваме дисперсията на $\tilde{\beta}$ ще означим $V = B - C$ и ще покажем, че $VX = 0$.

$$E(Vy) = E((B - C)y) = E(By) - E(Cy) = b - b = 0$$

Овечт мова $E(Vy) = V E y = V X b \Rightarrow VXb = 0$ за произволна стойност на параметрите b , които е възможно само при $VX = 0$. В следващото равенство използваме, че

$$\begin{aligned} \text{Var } \tilde{\beta} &= \text{Var}(By) = B(\text{Var } y) \cdot B' = B \cdot (\sigma^2 I_n) B' = \\ \sigma^2 (C + V)(C + V)' &= \sigma^2 (CC' + VC' + CV' + VV') = \\ &= \text{Var } \hat{\beta} + \sigma^2 V \cdot V'. \end{aligned}$$

Защо $VV' = 0$ е неотрицателно определена матрица.

* Коментари -

Заб 1. Теорема е в сила дори когато матрицата $X'X$ е изродена и не съществува единствена обрешка матрица $(X'X)^{-1}$.

Заб 2. При предположение за нормално разпределение на отклоненията и съответно на грешките, оценките решени по нормал. с-на $\hat{\beta}$ съвпадат с оценките получени по метода на

максимално правдоподобие, достигат равенството в неравенството на Рао-Крамър и \Rightarrow са ефективни, дисперсията им е долна граница не само за линейните, а за всички остани