

6

Регресиона анализа

При мерење на повеќе обележја, често се поставува прашањето како тие зависат едно од друго. На пример, како зависи нивото на холестеролот во крвта од возраста на индивидуата, како зависи потрошувачката на електричната енергија во домаќинството од надворешната температура и месечните примања на членовите на домаќинството и слично. При тоа променливата која е предмет на истражувањето се нарекува **зависна променлива** и се означува со Y и таа е случајна променлива (на пример, нивото на холестеролот во крвта, односно потрошувачката на електричната енергија во домаќинството), додека променливата која може да се контролира и која влијае на промените на зависната променлива се нарекува **независна променлива** и во општ случај нив може да ги има повеќе од една и се означуваат со $x^{(1)}, \dots, x^{(r)}$ и тие не се случајни променливи (на пример, возраста на индивидуата, односно надворешната температура и месечните примања на членовите на домаќинството).

Одредувањето на зависноста на променливата Y од независните променливи $x^{(1)}, \dots, x^{(r)}$ е преку формирање на математички модел кој ќе ја изразува таа зависност. Бидејќи во реалните ситуации на зависноста често влијаат и некои случајни фактори (на пример, грешки при мерењата), треба и тие да бидат земени во предвид при формирањето на моделот. Општиот облик на овој модел наречен **модел на регресија** е

$$Y = f_t(x^{(1)}, \dots, x^{(r)}) + \varepsilon, \quad (6.1)$$

каде функцијата f_t се нарекува **регресиона функција** и зависи од параметарот t кој многу често е векторски параметар, а ε е **случајна компонента** со $E(\varepsilon) = 0$ и $D(\varepsilon) = \sigma^2$.

Кога имаме само една независна променлива x , станува збор за **модел на еднодимензионална регресија** т.е.

$$Y = f_t(x) + \varepsilon, \quad (6.2)$$

а кога се повеќе од една независна променлива, имаме **модел на повеќедимензионална регресија**.

Изборот на соодветниот модел на регресија зависи од конкретната ситуација. Затоа, во случај на еднодимензионална регресија, најнапред се прикажуваат графички статистичките податоци $(x_1, y_1), \dots, (x_n, y_n)$, каде x_1, \dots, x_n се вредности на независната променлива x , а y_1, \dots, y_n се соодветните вредности на случајната променлива Y добиени како резултат на набљудувања или мерења, а потоа се бара општиот облик и отстапувањата од овој облик, со цел да се најде математичкиот модел кој најдобро би го описан обликов.

Основната задача на регресионата анализа е да после изборот на моделот, врз основа на низата статистички податоци, да ги процени непознатите параметри t и σ^2 . Попрецизно кажано, во случај на еднодимензионална регресија, треба да се дефинираат оценувачи за непознатите параметри t и σ^2 , врз основа на примерокот $(x_1, Y_1), \dots, (x_n, Y_n)$, каде x_1, \dots, x_n се вредности на независната променлива x и Y_1, \dots, Y_n се соодветните случајни променливи дефинирани со

$$Y_i = f_t(x_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad (6.3)$$

при што за случајните променливи $\varepsilon_1, \dots, \varepsilon_n$ претпоставуваме дека се независни со $E(\varepsilon_i) = 0$ и $D(\varepsilon_i) = \sigma^2$, $i = 1, \dots, n$.

6.1 Линеарна регресија

Ако регресионата функција f_t во (6.1) е линеарна, станува збор за **линеарен модел на регресија**. Ако кај линеарниот модел имаме само една независна променлива x , тогаш таквият модел е **прост линеарен модел на регресија**, и тогаш $t = (a, b)$, па моделот е

$$Y = ax + b + \varepsilon, \quad (6.4)$$

каде a и b се параметри и ε е случајна променлива со $E(\varepsilon) = 0$ и $D(\varepsilon) = \sigma^2$. Да забележиме дека тогаш,

$$E(Y) = E(ax + b + \varepsilon) = ax + b + E(\varepsilon) = ax + b,$$

$$D(Y) = D(ax + b + \varepsilon) = D(\varepsilon) = \sigma^2,$$

затоа што $ax + b$ не е случајна променлива.

Смислата на моделот (6.4) е следната. Случајната променлива Y зависи од променливата x , на тој начин што за секој $x = x_i$, $Y_i = ax_i + b + \varepsilon_i$ има еден детерминистички собирок кој е линеарна функција од x_i т.е. $ax_i + b$ и

Ирена Стојковска

стохастички собирок ε_i кој претставува случајна осцилација околу детерминистичкиот собирок затоа што $E(\varepsilon_i) = 0$. При едно реализирано мерење, добиваме една реализација на примерокот $(x_1, Y_1), \dots, (x_n, Y_n)$, која ја означуваме со $(x_1, y_1), \dots, (x_n, y_n)$. Нејзиниот графички приказ има тенденција на групирањето околу правата $y = ax + b$. Групирањето зависи од распределбата на $\varepsilon_1, \dots, \varepsilon_n$, имено колку нивната дисперзија $D(\varepsilon_i) = \sigma^2$, $i = 1, \dots, n$ е поголема, толку растурањето на тие точки околу правата $y = ax + b$ е поголемо.

6.1.1 Оценување на параметрите на регресија

Врз основа на примерокот $(x_1, Y_1), \dots, (x_n, Y_n)$, треба да најдеме оценувачи \hat{a} , \hat{b} , $\hat{\sigma}^2$ за непознатите параметри a , b , σ^2 на линеарниот модел на регресија (6.4). Без воведување на дополнителни претпоставки за распределбата на случајните променливи ε_i , со **методот на најмали квадрати** може да се најдат оценувачите \hat{a} и \hat{b} .

Нека $(x_1, y_1), \dots, (x_n, y_n)$ е реализација на примерокот $(x_1, Y_1), \dots, (x_n, Y_n)$, тогаш за реализациите \hat{a}' и \hat{b}' на оценувачите \hat{a} и \hat{b} , според методот на најмали квадрати треба да важи

$$\min_{(a,b) \in \mathbb{R}^2} \sum_{i=1}^n (y_i - ax_i - b)^2 = \sum_{i=1}^n (y_i - \hat{a}'x_i - \hat{b}')^2.$$

Па, слично како претходно (види Дескриптивна статистика, Дводимензионални обележја, определување на коефициентите на правата на регресија), со решавање на системот

$$\frac{\partial S(a, b)}{\partial a} = 0, \quad \frac{\partial S(a, b)}{\partial b} = 0,$$

каде

$$S(a, b) = \sum_{i=1}^n (y_i - ax_i - b)^2,$$

се добиваат решенијата (оценките на a и b добиени со методот на најмали квадрати)

$$\hat{a}' = \frac{s_{xy}}{s_x^2}, \quad \hat{b}' = \bar{y} - \hat{a}'\bar{x},$$

каде

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2, \\ s_{xy} &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \end{aligned}$$

(покажи!). Ако означиме со

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i,$$

тогаш оценувачите добиени со метод на најмали квадрати, статистики кои зависат од примерокот $(x_1, Y_1), \dots, (x_n, Y_n)$ се

$$\hat{a} = \frac{1}{ns_x^2} \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}) = \frac{1}{ns_x^2} \sum_{i=1}^n (x_i - \bar{x})Y_i, \quad (6.5)$$

$$\hat{b} = \bar{Y} - \hat{a}\bar{x} = \frac{1}{n} \sum_{i=1}^n \left(1 - \frac{\bar{x}}{s_x^2} (x_i - \bar{x}) \right) Y_i, \quad (6.6)$$

(покажи!). Исто така, според методот на најмали квадрати, како оценувач за σ^2 може да се земе

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{a}x_i + \hat{b}))^2. \quad (6.7)$$

Бидејќи $E(Y_i) = ax_i + b$ и $D(Y_i) = \sigma^2$, за оценувачите \hat{a} и \hat{b} имаме

$$\begin{aligned} E(\hat{a}) &= E\left(\frac{1}{ns_x^2} \sum_{i=1}^n (x_i - \bar{x})Y_i\right) = \frac{1}{ns_x^2} \sum_{i=1}^n (x_i - \bar{x})E(Y_i) = \frac{1}{ns_x^2} \sum_{i=1}^n (x_i - \bar{x})(ax_i + b) = \\ &= \frac{1}{s_x^2} \left(\frac{a}{n} \sum_{i=1}^n (x_i^2 - x_i\bar{x}) + \frac{b}{n} \sum_{i=1}^n (x_i - \bar{x}) \right) = \frac{1}{s_x^2} (as_x^2 + b \cdot 0) = a, \\ E(\hat{b}) &= E(\bar{Y} - \hat{a}\bar{x}) = E(\bar{Y}) - E(\hat{a})\bar{x} = \frac{1}{n} \sum_{i=1}^n E(Y_i) - a\bar{x} = \\ &= \frac{1}{n} \sum_{i=1}^n (ax_i + b) - a\bar{x} = a\bar{x} + b - a\bar{x} = b, \end{aligned}$$

што значи дека \hat{a} и \hat{b} се непристрасни оценувачи за a и b соодветно. Понатаму, од независноста на ε_i , $i = 1, \dots, n$ следи и независноста на $Y_i = ax_i + b + \varepsilon_i$, $i = 1, \dots, n$ (покажи!). Затоа, за дисперзиите на \hat{a} и \hat{b} имаме

$$D(\hat{a}) = \frac{\sigma^2}{ns_x^2} \rightarrow 0, \quad D(\hat{b}) = \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{s_x^2} \right) \rightarrow 0,$$

кога $n \rightarrow \infty$ (покажи!), од каде заклучуваме дека \hat{a} и \hat{b} се конзистентни оценувачи за a и b соодветно.

Ирена Стојковска

Ако сега ја разгледаме статистиката $\hat{a}x + \hat{b}$ како оценувач за $ax + b$, бидејќи

$$E(\hat{a}x + \hat{b}) = ax + b, \quad D(\hat{a}x + \hat{b}) = \frac{\sigma^2}{n} \left(1 + \frac{(x - \bar{x})^2}{s_x^2} \right) \rightarrow 0,$$

кога $n \rightarrow \infty$ (покажи!), заклучуваме дека $\hat{a}x + \hat{b}$ е непристрасен и конзистентен оценувач за $ax + b$.

Но, за оценувачот $\hat{\sigma}^2$ дефиниран со (6.7) имаме

$$E(\hat{\sigma}^2) = E\left(\frac{1}{n} \sum_{i=1}^n (Y_i - (\hat{a}x_i + \hat{b}))^2\right) = \frac{n-2}{n} \sigma^2$$

(покажи!), од каде заклучуваме дека тој не е непристрасен оценувач за σ^2 . Затоа, неговата корекција, оценувачот

$$\check{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - (\hat{a}x_i + \hat{b}))^2 \quad (6.8)$$

е непристрасен оценувач за σ^2 .

За да ги користиме најдените оценувачи при барање на **интервали на доверба и тестирање на хипотези** за параметрите на регресија, потребно е да ги знаеме и нивните распределби, точни или асимптотски. Така, **во случај на голем примерок**, од централната гранична теорема, следи дека оценувачите \hat{a} и \hat{b} се асимптотски нормални оценувачи за a и b соодветно, односно

$$\hat{a} \xrightarrow{\text{dist}} \mathcal{N}\left(a, \frac{\sigma^2}{ns_x^2}\right), \quad \hat{b} \xrightarrow{\text{dist}} \mathcal{N}\left(b, \frac{\sigma^2}{n} \left(1 + \frac{\bar{x}^2}{s_x^2}\right)\right).$$

Истото важи и за оценувачот $\hat{a}x + \hat{b}$, за кој се покажува дека е асимптоцки нормален оценувач за $ax + b$, односно

$$\hat{a}x + \hat{b} \xrightarrow{\text{dist}} \mathcal{N}\left(ax + b, \frac{\sigma^2}{n} \left(1 + \frac{(x - \bar{x})^2}{s_x^2}\right)\right).$$

Па, бри барање на интервали на доверба и тестирање на хипотези за a , b и $ax + b$ врз основа на голем примерок, се користиме со гореспоменатите распределби, при што непознатиот параметар σ^2 го заменуваме со реализација на оценувачот дефиниран со (6.7) или (6.8), што во пракса нема големо значење која оценка ќе се одбере затоа што за големи вредности на n разликата меѓу двете оценки е незначителна.

За да ги најдеме точните распределби на оценувачите \hat{a} и \hat{b} , **во случај на мал примерок** потребни се дополнителни претпоставки за случајните променливи ε_i , $i = 1, \dots, n$ во линеарниот модел на регресија $Y_i = ax_i + b + \varepsilon_i$, $i = 1, \dots, n$.

Имено, покрај независноста се претпоставува и дека $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$, $i = 1, \dots, n$. Од тука следи дека и случајните променливи Y_i , $i = 1, \dots, n$ се независни и за нивните распределби важи

$$Y_i \sim \mathcal{N}(ax_i + b, \sigma^2), \quad i = 1, \dots, n.$$

Па тогаш, за точните распределби на оценувачите \hat{a} и \hat{b} дефинирани со (6.5) и (6.6) соодветно, имаме

$$\hat{a} \sim \mathcal{N}\left(a, \frac{\sigma^2}{ns_x^2}\right), \quad \hat{b} \sim \mathcal{N}\left(b, \frac{\sigma^2}{n}\left(1 + \frac{\bar{x}^2}{s_x^2}\right)\right),$$

бидејќи претставуваат линеарна комбинација од случајните променливи Y_i , $i = 1, \dots, n$. Во овој случај, случајната променлива $\hat{\sigma}^2$ дефинирана со (6.7), а исто така и корегираниот оценувач $\check{\sigma}^2$ дефиниран со (6.8), претставуваат збир од квадрати на одредени линеарни комбинации од нормално распределени случајни променливи Y_i , \hat{a} и \hat{b} , кои не се независни. Но, се покажува дека, и во тој случај, распределбата на така претставената случајна променлива е хи-квадрат распределба. Имено, важи

$$\frac{n}{\sigma^2} \hat{\sigma}^2 \sim \chi_{n-2}^2, \text{ односно } \frac{n-2}{\sigma^2} \check{\sigma}^2 \sim \chi_{n-2}^2.$$

Исто така, се покажува дека \hat{a} и $\hat{\sigma}^2$, но и \hat{b} и $\hat{\sigma}^2$ се независни, што помага при одредување на точните распределби на статистиките кои се користат при наоѓање на интервалите на доверба и тестирање на хипотезите за параметрите a , b и $ax + b$, имено важи

$$T_1 = \frac{(\hat{a} - a) \sqrt{(n-2)s_x^2}}{\hat{\sigma}} = \frac{(\hat{a} - a) \sqrt{ns_x^2}}{\check{\sigma}} \sim t_{n-2},$$

$$T_2 = \frac{(\hat{b} - b) \sqrt{(n-2)s_x^2}}{\hat{\sigma} \sqrt{s_x^2 + \bar{x}^2}} = \frac{(\hat{b} - b) \sqrt{ns_x^2}}{\check{\sigma} \sqrt{s_x^2 + \bar{x}^2}} \sim t_{n-2},$$

$$T_3 = \frac{(\hat{a}x + \hat{b} - ax - b) \sqrt{(n-2)s_x^2}}{\hat{\sigma} \sqrt{s_x^2 + (x - \bar{x})^2}} = \frac{(\hat{a}x + \hat{b} - ax - b) \sqrt{ns_x^2}}{\check{\sigma} \sqrt{s_x^2 + (x - \bar{x})^2}} \sim t_{n-2}$$

(покажи!).

Оценувачи за параметрите на регресија a , b и σ^2 може да се најдат и со **метод на максимална подобност**. Но, и во тој случај потребни се дополнителните претпоставки за нормална распределеност на случајните променливи ε_i , $i = 1, \dots, n$ во линеарниот модел на регресија $Y_i = ax_i + b + \varepsilon_i$, $i = 1, \dots, n$, имено покрај независноста се претпоставува и дека $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$, $i = 1, \dots, n$.

Ирена Стојковска

Го означуваме со $t = (a, b, \sigma^2)$ векторот од непознати параметри кој припаѓа на просторот од параметри

$$\Theta = \{(a, b, \sigma^2) \mid a \in \mathbb{R}, b \in \mathbb{R}, \sigma^2 > 0\}.$$

Тогаш, функцијата на подобност е

$$L(x, y, t) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - ax_i - b)^2\right\}.$$

Со решавање на системот равенки на подобност

$$\frac{\partial \ln L(x, y, t)}{\partial a} = 0, \quad \frac{\partial \ln L(x, y, t)}{\partial b} = 0, \quad \frac{\partial \ln L(x, y, t)}{\partial \sigma^2} = 0,$$

по a , b и σ^2 се добиваат решенијата

$$a = \frac{s_{xy}}{s_x^2}, \quad b = \bar{y} - a\bar{x}, \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n (ax_i + b - y_i)^2$$

(покажи!), од каде заклучуваме дека максимално подобни (ML) оценувачи за a , b и σ^2 , се \hat{a} , \hat{b} и $\hat{\sigma}^2$, дефинирани со (6.5), (6.6) и (6.7) соодветно.