

Z. HELLWIG (Wrocław)

WYZNACZANIE PARAMETRÓW REGRESJI LINIOWEJ METODĄ DWÓCH PUNKTÓW

Niech (ξ, η) będzie zmienną losową typu ciągłego. Zakładamy, że w rozkładzie tej zmiennej krzywe regresji typu pierwszego są liniami prostymi.

Równania linii regresji można więc napisać w postaci

$$(1) \quad y_x - \nu = \alpha_y(x - \mu),$$

$$(2) \quad x_y - \mu = \alpha_x(y - \nu),$$

gdzie μ, ν są współrzędnymi środka ciężkości populacji, α_x i α_y zaś są współczynnikami regresji. Jeżeli $\mu = 0$ i $\nu = 0$, to równania (1) i (2) przybierają prostszą postać

$$y_x = \alpha_y x,$$

$$x_y = \alpha_x y.$$

Linie (1) i (2) są liniami regresji typu pierwszego, tzn.

$$y_x = E(\eta | \xi = x),$$

$$x_y = E(\xi | \eta = y).$$

Estymując na podstawie próbki parametry występujące w (1) i (2), korzystamy na ogół z metody najmniejszej sumy kwadratów. Ujemną stroną tej metody jest to, że wymaga ona wielu żmudnych operacji rachunkowych, jak podnoszenie do kwadratu, obliczanie iloczynów, sumowanie długich kolumn liczb oraz wyciąganie pierwiastków. Poniżej proponuję inną, prostszą rachunkowo metodę wyznaczania parametrów linii regresji.

Metoda ta, zwana dalej *metodą dwóch punktów*, daje, jak zobaczymy, estymatory zgodne i nieobciążone. Efektywność ich jest mniejsza od efektywności estymatorów uzyskanych metodą klasyczną, lecz za to rachunki związane z ich wyznaczeniem są łatwiejsze. Rozkład proponowanych estymatorów jest zbieżny do rozkładu normalnego.

Przypuśćmy, że w celu oszacowania parametrów a_x i a_y , z populacji generalnej pobrano próbkę liczącą n punktów o współrzędnych (x_i, y_i) ($i = 1, 2, \dots, n$). Obliczamy średnie

$$\bar{x} = (x_1 + x_2 + \dots + x_n)/n, \quad \bar{y} = (y_1 + y_2 + \dots + y_n)/n,$$

a następnie dzielimy zbiór wszystkich punktów na dwa podzbiory, które oznaczymy liczbami rzymskimi I i II. Do podzbioru I zaliczamy punkty, których współrzędne x są większe od \bar{x} , do podzbioru II pozostałe punkty. Środki ciężkości tych podzbiorów wyznaczają prostą, na której leży środek ciężkości wszystkich punktów, tzn. punkt o współrzędnych (\bar{x}, \bar{y}) .

Otóż za estymator parametru a_y proponuję wziąć współczynnik kierunkowy tej prostej. Oznaczmy go symbolem a_y , przy czym

$$a_y = (\bar{y}_I - \bar{y}) / (\bar{x}_I - \bar{x}).$$

We wzorze tym mamy

$$\bar{y}_I = (y_{i_1} + y_{i_2} + \dots + y_{i_k})/k, \quad \bar{x}_I = (x_{i_1} + x_{i_2} + \dots + x_{i_k})/k,$$

gdzie $(x_{i_1}, y_{i_1}), \dots, (x_{i_k}, y_{i_k})$ są to te spośród punktów (x_i, y_i) , których współrzędne x_i spełniają warunek $x_i > \bar{x}$. Litera k oznacza tu liczbę punktów o odciętych większych od \bar{x} . Zwróćmy uwagę, że k jest zmienną losową, która może przybierać wartości $1, 2, \dots, n-1$.

Udowodnimy, że estymator a_y jest zgodny i nieobciążony, oraz zbadamy jego efektywność.

TWIERDZENIE 1. *Współczynnik a_y jest estymatorem zgodnym współczynnika regresji a_y .*

Dowód. Mamy udowodnić, że dla dowolnego $\varepsilon > 0$

$$P\{|a_y - a_y| > \varepsilon\} \xrightarrow{n \rightarrow \infty} 0.$$

W dowodzie wykorzystamy następujący

LEMAT 1. *Jeśli $E(\eta|\xi = x) = ax$, to dla każdego h*

$$E(\eta|x-h < \xi < x+h) = aE(\xi|x-h < \xi < x+h).$$

Istotnie, z definicji mamy

$$E(\eta|x-h < \xi < x+h) = \int_{x-h}^{x+h} \int_{-\infty}^{\infty} y f(x, y) \left[\int_{x-h}^{x+h} \left(\int_{-\infty}^{+\infty} f(x, y) dy \right) dx \right]^{-1} dy dx,$$

gdzie $f(x, y)$ jest gęstością dwuwymiarowego rozkładu prawdopodobieństwa.

Oznaczmy

$$C = \int_{x-h}^{x+h} \left[\int_{-\infty}^{\infty} f(x, y) dy \right] dx.$$

W takim razie jest

$$\mathbf{E}(\eta|x-h < \xi < x+h) = C^{-1} \int_{x-h}^{x+h} \int_{-\infty}^{\infty} yf(x, y) dx dy.$$

Ponieważ zachodzi równość

$$f(xy) = f(y|x)f_1(x),$$

w której $f(y|x)$ jest gęstością warunkową prawdopodobieństwa zmiennej η , a $f_1(x)$ gęstością brzegową zmiennej ξ , więc

$$\begin{aligned} \mathbf{E}(\eta|x-h < \xi < x+h) &= C^{-1} \int_{x-h}^{x+h} \int_{-\infty}^{\infty} yf(y|x)f_1(x) dy dx = \\ &= C^{-1} \int_{x-h}^{x+h} f_1(x) \left[\int_{-\infty}^{\infty} yf(y|x) dy \right] dx = \\ &= C^{-1} \int_{x-h}^{x+h} f_1(x) \mathbf{E}(\eta|\xi = x) dx = \\ &= C^{-1} \int_{x-h}^{x+h} f_1(x) ax dx = \\ &= a \left(\int_{x-h}^{x+h} xf_1(x) dx \right) / \int_{x-h}^{x+h} f_1(x) dx = \\ &= a \mathbf{E}(\xi|x-h < \xi < x+h). \end{aligned}$$

Otrzymaliśmy zatem równość

$$\mathbf{E}(\eta|x-h < \xi < x+h) = a \mathbf{E}(\xi|x-h < \xi < x+h), \quad \text{c. n. d.}$$

Z udowodnionego lematu wynika następujący
WNIOSEK.

$$\mathbf{E}(\bar{y}_I) = a \mathbf{E}(\bar{x}_I).$$

Mamy bowiem

$$(3) \quad \mathbf{E}(y_1|x > \bar{x}) = \mathbf{E}(y_2|x_2 > \bar{x}) = \dots = \mathbf{E}(y_n|x_n > \bar{x}),$$

z definicji zaś otrzymujemy

$$\mathbf{E}(\bar{y}_I) = \mathbf{E}((y_{i_1} + \dots + y_{i_k})/k).$$

Wobec tego, przy ustalonym k jest

$$\begin{aligned} &\mathbf{E}((y_{i_1} + \dots + y_{i_k})/k) = \\ &= \mathbf{E} \left(\frac{1}{k} (y_1 + \dots + y_k) \mid x_i > \bar{x}, i = 1, 2, \dots, k; \bar{x} \geq x_i, i = k+1, k+2, \dots, n \right). \end{aligned}$$

W takim razie

$$E(\bar{y}_I) = E(y_1|x_1 > \bar{x}).$$

Analogicznie, łatwo wykazać, że

$$E(\bar{x}_I) = E(x_1|x_1 > \bar{x}).$$

Dla uzasadnienia słuszności wniosku pozostaje więc do udowodnienia równość

$$E(y_1|x_1 > \bar{x}) = \alpha E(x_1|x_1 > \bar{x}).$$

Równość powyższą można również napisać inaczej, a mianowicie

$$E(y_1|x_1 > z) = \alpha E(x_1|x_1 > z),$$

gdzie x_1 i z są zmiennymi losowymi niezależnymi. Wystarczy bowiem warunek $x_1 > \bar{x}$ napisać w postaci $x_1 > (x_2 + \dots + x_n)/(n-1)$ i przyjąć $z = (x_2 + \dots + x_n)/(n-1)$.

Dla dowolnego rzeczywistego u mamy na podstawie lematu

$$E(y_1|x_1 > u) = \alpha E(x_1|x_1 > u).$$

Oznaczając przez $g(u)$ gęstość prawdopodobieństwa zmiennej losowej z widzimy, że

$$\begin{aligned} E(y_1|x_1 > \bar{x}) &= \int_{-\infty}^{\infty} E(y_1|x_1 > u)g(u)du = \\ &= \int_{-\infty}^{\infty} \alpha E(x_1|x_1 > u)g(u)du = \\ &= \alpha \int_{-\infty}^{\infty} E(x_1|x_1 > u)g(u)du = \\ &= \alpha E(x_1|x_1 > z), \quad \text{c. n. d.} \end{aligned}$$

Wiadomo, że \bar{x} zmierza według prawdopodobieństwa do $E(x)$ i że \bar{y} zmierza według prawdopodobieństwa do $E(\bar{y})$. Można również pokazać, że \bar{x}_I zmierza według prawdopodobieństwa do $E(\bar{y}_I)$ oraz \bar{y}_I zmierza według prawdopodobieństwa do $E(\bar{x}_I)$. Możemy więc korzystać z twierdzenia Słuckiego (zob. [3], str. 282), które głosi, że jeśli zmienne losowe $\xi_n, \eta_n, \dots, \varrho_n$ są zbieżne według prawdopodobieństwa odpowiednio do stałych x, y, \dots, r , to każda wymierna funkcja $R(\xi_n, \eta_n, \dots, \varrho_n)$ zdąży według prawdopodobieństwa do stałej $R(x, y, \dots, r)$, jeśli tylko wielkość $R(x, y, \dots, r)$ jest skończona.

Z twierdzenia tego wynika, że $a_y = (\bar{y}_I - \bar{y})/(\bar{x}_I - \bar{x})$ zdąży według prawdopodobieństwa do $(E(\bar{y}_I) - E(\bar{y})) / (E(\bar{x}_I) - E(\bar{x}))$. Na podstawie wniosku mamy

$$\frac{E(\bar{y}_I) - E(\bar{y})}{E(\bar{x}_I) - E(\bar{x})} = \frac{\alpha E(\bar{x}_I) - \alpha E(\bar{x})}{E(\bar{x}_I) - E(\bar{x})} = \alpha_y.$$

a więc a_y zdąży według prawdopodobieństwa do a_y , czego należało dowieść.

Obecnie wykażemy, że a_y jest estymatorem nieobciążonym.

TWIERDZENIE 2.

$$E(a_y) = a_y.$$

Dowód. Niech u_1, u_2, \dots, u_n będą dowolnymi liczbami rzeczywistymi. Wobec tego, jeśli $x_1 = u_1, x_2 = u_2, \dots, x_n = u_n$, to

$$E(a_y | x_1 = u_1, \dots, x_n = u_n) = (E(\bar{y}_I) - E(\bar{y})) / (\bar{x}_I - \bar{x}).$$

Ale

$$\begin{aligned} & \frac{E(\bar{y}_I) - E(\bar{y})}{\bar{x}_I - \bar{x}} = \\ & = \frac{[E(\eta | \xi = x_1) + \dots + E(\eta | \xi = x_k)]/k - [E(\eta | \xi = x_1) + \dots + E(\eta | \xi = x_n)]/n}{\bar{x}_I - \bar{x}}, \end{aligned}$$

a ponieważ z założenia mamy $E(\eta | \xi = x) = ax$, więc

$$\frac{E(\bar{y}_I) - E(\bar{y})}{\bar{x}_I - \bar{x}} = \frac{(ax_1 + \dots + ax_k)/k - (ax_1 + \dots + ax_n)/n}{\bar{x}_I - \bar{x}} = a.$$

Oczywiście jest

$$\begin{aligned} E(a_y) &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} E(a_y | x_1 = u_1, \dots, x_n = u_n) f_1(u_1) \dots f_1(u_n) du_1 \dots du_n = \\ &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} a f_1(u_1) \dots f_1(u_n) du_1 \dots du_n = \\ &= a \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_1(u_1) \dots f_1(u_n) du_1 \dots du_n, \end{aligned}$$

przy czym $f_1(u)$ oznacza gęstość brzegową zmiennej ξ .

Ponieważ mamy

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_1(u_1) \dots f_1(u_n) du_1 \dots du_n = 1,$$

więc

$$E(a_y) = a_y, \quad \text{c. n. d.}$$

Przystąpimy obecnie do porównania wariancji estymatora a_y z wariancją estymatora A_y (symbol A_y oznacza współczynnik regresji z próbki, otrzymany metodą klasyczną).

Niech będzie $u_n \leq u_{n-1} \leq \dots \leq u_1$. Wprowadźmy oznaczenia ⁽¹⁾

$$V(A_y) = \mathcal{V}(A_y | x_1 = u_1, \dots, x_n = u_n),$$

$$V(a_y) = \mathcal{V}(a_y | x_1 = u_1, \dots, x_n = u_n)$$

oraz

$$s_u = \sqrt{\left(\sum_1^n (u - \bar{u})^2\right) / n}, \quad d_u = \left(\sum_1^n |u - \bar{u}|\right) / n.$$

TWIERDZENIE 3.

$$e(a_y | x_1 = u_1, \dots, x_n = u_n) = \frac{V(A_y)}{V(a_y)} \geq \frac{d^2 u}{s^2 u} \cdot \frac{n-1}{n}.$$

Zbadajmy, czemu jest równe $V(a_y)$. Mamy

$$(4) \quad V(a_y) = \mathcal{V}(\bar{y}_I - \bar{y}) / (\bar{x}_I - \bar{x})^2.$$

LEMAT 2.

$$(5) \quad d_x = \left(\sum_{i=1}^n |x_i - \bar{x}|\right) / n = 2k(\bar{x}_I - \bar{x}) / n,$$

gdzie k jest liczbą x -ów większych od \bar{x} .

Do wód. Oznaczając $z_i = x_i - \bar{x}$ możemy napisać równość następującą:

$$(6) \quad \sum_{i=1}^k z_i = \sum_{j=1}^l (-z_j),$$

gdzie każde $z_i > 0$, a każde $z_j \leq 0$, przy czym $k+1 = n$.

Odchylenie przeciętne d_x jest równe wyrażeniu

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}| = \frac{1}{n} \sum_{i=1}^n |z_i| = \frac{1}{n} \left(\sum_{i=1}^k |z_i| + \sum_{i=1}^l |z_i| \right).$$

Ze względu na (6) mamy

$$(7) \quad d_x = \frac{2}{n} \sum_{i=1}^k z_i = \frac{2}{n} \sum_{j=1}^l (-z_j).$$

Zbadajmy obecnie, czemu jest równa różnica $\bar{x}_I - \bar{x}$. Mamy

$$\begin{aligned} \bar{x}_I - \bar{x} &= \frac{1}{k} \sum_{i=1}^k (\bar{x} + z_i) - \bar{x} = \frac{1}{k} \left(k\bar{x} + \sum_{i=1}^k z_i \right) - \bar{x} = \\ &= \frac{1}{k} \left(k\bar{x} + \sum_{i=1}^k z_i - k\bar{x} \right) = \frac{1}{k} \sum_{i=1}^k z_i. \end{aligned}$$

⁽¹⁾ Korzystamy tutaj z wariacji warunkowych, gdyż jak wiadomo ([2], str. 156), nie zważa to ogólności rozważań.

Lemat zostanie udowodniony, gdy otrzymany wynik podstawimy do (7).
Mamy bowiem

$$\sum_{i=1}^k z_i = k(\bar{x}_I - \bar{x}),$$

skąd

$$d_x = 2k(\bar{x}_I - \bar{x})/n, \quad \text{c. n. d.}$$

Na mocy lematu 2 wzór (4) można również napisać w innej postaci,
mianowicie

$$(8) \quad V(a_y) = \frac{4k^2}{n^2 d_x^2} \mathcal{V}(\bar{y}_I - \bar{y}).$$

We wzorze tym

$$\begin{aligned} \mathcal{V}(\bar{y}_I - \bar{y}) &= \mathcal{V}\left(\frac{y_1 + y_2 + \dots + y_k}{k} - \frac{y_1 + y_2 + \dots + y_k + y_{k+1} + \dots + y_n}{n}\right) = \\ &= \mathcal{V}\left(\frac{n \sum_{r=1}^k y_r}{kn} - \frac{k \sum_{r=1}^k y_r}{kn} - \frac{k \sum_{r=k+1}^n y_r}{kn}\right) = \\ &= \mathcal{V}\left((n-k) \frac{\sum_{r=1}^k y_r}{kn} - \frac{k \sum_{r=k+1}^n y_r}{kn}\right) = \\ &= \frac{(n-k)^2}{k^2 n^2} \sum_{r=1}^k \mathcal{V}(y_r) + \frac{k^2}{k^2 n^2} \sum_{r=k+1}^n \mathcal{V}(y_r). \end{aligned}$$

Ponieważ jednak mamy

$$\mathcal{V}(y_r) = \mathcal{V}(\eta|\xi = x_r) = \mathcal{V}(ax + \beta + \Delta), \quad \text{gdzie } \Delta = y - ax - \beta, \quad \beta = v - a\mu,$$

więc

$$\mathcal{V}(y_r) = \mathcal{V}(\Delta).$$

W takim razie

$$\begin{aligned} \mathcal{V}(y_I - \bar{y}) &= \frac{(n-k)^2}{k^2 n^2} k \mathcal{V}(\Delta) + \frac{k^2}{k^2 n^2} (n-k) \mathcal{V}(\Delta) = \\ &= \mathcal{V}(\Delta) \left(\frac{n^2 k - 2nk^2 + k^3 + nk^2 - k^3}{k^2 n^2} \right) = \\ &= \mathcal{V}(\Delta) \left(\frac{n^2 k - nk^2}{k^2 n^2} \right) = \mathcal{V}(\Delta) \left(\frac{1}{k} - \frac{1}{n} \right). \end{aligned}$$

Wobec tego wzór (8) przybierze postać

$$V(a_y) = \frac{4k^2}{n^2 d_x^2} \mathcal{V}(\Delta) \left(\frac{1}{k} - \frac{1}{n} \right),$$

czyli

$$(9) \quad V(a_y) = \frac{\mathcal{V}(\Delta)}{n d_x^2} \left(\frac{4k}{n} - \frac{4k^2}{n^2} \right).$$

Zajmijmy się wariancją $V(A_y)$; otrzymujemy

$$\begin{aligned} V(A_y) &= \mathcal{V} \left(\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right) / \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^2 = \\ &= \mathcal{V}[(x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})] / \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^2 = \\ &= \frac{n-1}{n} \left[\sum_{i=1}^n (x_i - \bar{x}) \mathcal{V}(\Delta) \right] / \left[\sum_{i=1}^n (x_i - \bar{x})^2 \right]^2 = \\ &= \frac{n-1}{n} \mathcal{V}(\Delta) / \sum_{i=1}^n (x_i - \bar{x})^2, \end{aligned}$$

czyli

$$V(A_y) = \frac{\mathcal{V}(\Delta)}{n s_x^2} \cdot \frac{n-1}{n}.$$

Stąd mamy

$$\begin{aligned} e(a_y | x_1 = u_1, \dots, x_n = u_n) &= \frac{(\mathcal{V}(\Delta) / n s_x^2) ((n-1)/n)}{(\mathcal{V}(\Delta) / n d_x^2) (4k/n - 4k^2/n^2)} = \\ &= \frac{d_x^2}{s_x^2} \cdot \frac{n^2}{4k(n-k)} \cdot \frac{n-1}{n}. \end{aligned}$$

Wyrażenie

$$(10) \quad n^2 / 4k(n-k)$$

osiąga minimum ze względu na k , gdy $k = \frac{1}{2}n$. Podstawiając $k = \frac{1}{2}n$ do (10) otrzymujemy

$$n^2 / (4 \cdot \frac{1}{2}n(n - \frac{1}{2}n)) = 1,$$

a wobec tego

$$(11) \quad e(a_y | x_1 = u_1, \dots, x_n = u_n) \geq (d_x^2 / s_x^2) ((n-1)/n), \quad \text{c. n. d.}$$

Stosunek $e = e(a_y | x_1 = u_1, \dots, x_n = u_n)$ wariancji $V(A_y)$ i $V(a_y)$ możemy traktować jako zmienną losową ze względu na x_1, \dots, x_n .

Stosując cytowane wyżej twierdzenie Slutskiego do prawej strony wzoru (11) otrzymujemy, że e zmierza według prawdopodobieństwa do D_x^2/σ_x^2 , gdzie

$$D_x = \int_{-\infty}^{\infty} |x| f_1(x) dx, \quad \sigma_x^2 = \int_{-\infty}^{\infty} x^2 f_1(x) dx.$$

Jeżeli populacja ma rozkład normalny, to e zmierza według prawdopodobieństwa do $2/\pi$.

Zajmiemy się obecnie zbadaniem rozkładu estymatora a_y . Udowodnimy następujące

TWIERDZENIE 4. *Rozkład zmiennej losowej*

$$t = \frac{(a_y - a_y)}{\sqrt{V(a_y)}}$$

jest dla $n \rightarrow \infty$ zbieżny do rozkładu normalnego $N(0, 1)$.

Dowód. Zauważmy najpierw, że na podstawie twierdzenia 2 jest $E(t) = 0$, a na podstawie wzoru (9), $V(t) = 1$. Po prostych przekształceniach, w których wykorzystuje się wzór (5), otrzymujemy

$$(12) \quad t = \frac{\sqrt{n}[(\bar{y}_1 - \bar{y})k/n - \frac{1}{2}a_y d_x]}{\sqrt{\mathcal{V}(\Delta)} \sqrt{(1 - k/n)k/n}}$$

LEMAT 3. *Jeżeli $n \rightarrow \infty$, to*

$$P\{|k/n - [1 - F(\mu + 0)]| \rightarrow 0\} = 1,$$

gdzie $F(x)$ oznacza dystrybuantę brzegową zmiennej ξ .

Dowód. Na podstawie twierdzenia Gliwienki ([1], str. 282) mamy

$$P\left\{\sup_{-\infty < x < \infty} |F_n(x) - F(x)| \rightarrow 0\right\} = 1,$$

gdy $n \rightarrow \infty$; we wzorze tym $F_n(x)$ oznacza empiryczną dystrybuantę próbki, a $F(x)$ oznacza dystrybuantę populacji.

Wobec tego jest

$$P\left\{\sup_{-\infty < x < \infty} |F_n(\mu) - F(\mu)| \xrightarrow{n \rightarrow \infty} 0\right\} = 1.$$

Ponieważ z drugiej strony

$$P\left\{|\bar{x}_n - \mu| \xrightarrow{n \rightarrow \infty} 0\right\} = 1,$$

więc ze względu na ciągłość dystrybuanty $F(x)$, otrzymujemy, że

$$P\left\{\sup_{-\infty < x < \infty} |F_n(\bar{x}_n) - F(\mu)| \xrightarrow{n \rightarrow \infty} 0\right\} = 1.$$

Ale

$$k/n = 1 - F_n(\bar{x} + 0),$$

skąd

$$P\left\{ \sup_{-\infty < x < \infty} |k/n - [1 - F(\mu + 0)]| \rightarrow 0 \right\} = 1.$$

Lemat został więc udowodniony.

We wzorze (12) wartości $\mathcal{V}(\Delta)$ i α_y są stałe, a k/n i d_x zbiegają według prawdopodobieństwa do stałych.

Niech $H(t)$ będzie dystrybuantą zmiennej t ; wtedy

$$H(t) = P(k=1)H(t|k=1) + \dots + P(k=n-1)H(t|k=n-1),$$

gdzie $H(t|k=k_0)$ oznacza warunkową dystrybuantę zmiennej t pod warunkiem, że $k=k_0$. Niech m będzie taką liczbą naturalną, że dla każdego k , spełniającego nierówność $m < k < n-m$, jest

$$(13) \quad \sup |\Phi(t) - F(t|k)| < \varepsilon.$$

We wzorze tym mamy $\Phi(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx$, ε zaś jest dowolną liczbą rzeczywistą, większą od zera.

Nierówność (13) możemy napisać dlatego, że z definicji zmiennej t wynika, iż rozkład warunkowy $H(t|k)$, gdy n , k i $(n-k)$ rosną nieograniczenie, zmierza do $\Phi(t)$.

Istotnie, występująca w liczniku prawej strony wzoru (12) różnica średnich

$$\begin{aligned} \bar{y}_I - \bar{y} &= \frac{(n-k) \sum_1^k (y|x > \bar{x}) - k \sum_{k+1}^n (y|x \leq \bar{x})}{nk} = \\ &= \frac{n-k}{n} \left[\frac{1}{k} \sum_1^k (y|x > \bar{x}) - \frac{1}{n-k} \sum_{k+1}^n (y|x \leq \bar{x}) \right] = \\ &= \frac{n-k}{n} [\bar{Y}_I - \bar{Y}_{II}], \end{aligned}$$

gdzie \bar{Y}_I i \bar{Y}_{II} są to niezależne od siebie średnie arytmetyczne niezależnych zmiennych losowych. Jeżeli $n \rightarrow \infty$, $k \rightarrow \infty$, $n-k \rightarrow \infty$, to na podstawie centralnego twierdzenia granicznego rozkład różnicy $\bar{Y}_I - \bar{Y}_{II}$ zmierza do rozkładu normalnego, a na mocy lematu 3 wyrażenie $(n-k)/n$ zdąży według prawdopodobieństwa do stałej. Wobec tego na podstawie twierdzenia o zbieżności ([3], str. 281) mamy $H(t|k) \rightarrow \Phi(t)$.

LEMAT 4. Dla każdego m jest

$[P(k=1)+\dots+P(k=m)+P(k=n-m)+\dots+P(k=n-1)] \rightarrow 0$,
gdy $n \rightarrow \infty$.

Dowód. Na podstawie lematu 3 istnieje takie n_0 , że dla $n > n_0$ mamy

$$P\{|k/n - [1 - F(\mu+0)]| < \varepsilon\} > 1 - \varepsilon.$$

Dla wygody oznaczmy

$$1 - F(\mu+0) = q, \quad F(\mu+0) = p;$$

będzie wtedy

$$P\{k < n(q-\varepsilon) \text{ lub } k > n(q+\varepsilon)\} < \varepsilon.$$

Widać stąd, że dla spełnienia nierówności

$$P\{k < m \text{ lub } k > n-m\} < \varepsilon,$$

trzeba, żeby było $n(q-\varepsilon) > m$ oraz $n(q+\varepsilon) < n-m$, czyli żeby n było większe od większej z dwóch liczb $m/(q-\varepsilon)$, $m/(p-\varepsilon)$. Wynika stąd, że zawsze można dobrać takie n , że będzie spełniona nierówność

$$P\{k < m \text{ lub } k > n-m\} < \varepsilon.$$

Lemat został więc udowodniony. Niech

$$S(t) = \sum_{k=m+1}^{n-m-1} p(k)H(t|k).$$

W takim razie na mocy lematu 4 istnieje zawsze takie n_0 , zależne od m i ε , że dla $n > n_0$ spełniona jest nierówność

$$(14) \quad |H(t) - S(t)| < \varepsilon.$$

Z wzorów (13) i (14) wynika nierówność

$$|\Phi(t) - H(t)| < 2\varepsilon,$$

co dowodzi twierdzenia.

Na zakończenie przytoczymy przykład liczbowy. Pozwoli on zaznaczyć się z techniką znajdowania parametrów regresji metodą dwóch punktów.

Kolumna x -ów w podanej tabelicy (str. 78 i 79) przedstawia produkcję prądu elektrycznego mierzoną na stykach generatorów, natomiast kolumna y -ów zawiera dane liczbowe dotyczące mialu węglowego zużytego w związku z tą produkcją.

Materiał pochodzi z elektrowni hutniczej, mającej dwa turbogeneratory systemu Škoda o mocy 3100 kW każdy oraz dwa kotły wodnorurkowe systemu Duquesne, opalane pyłem węglowym.

Zużycie węgla (miału) podane jest brutto, tak jak wykazała waga kontrolna. Produkcja prądu jest również podana brutto, gdyż pomiarów prądu dokonano na zaciskach generatorów.

Używając jakichkolwiek znaków wyróżniających (w niniejszym tekście użyto znaku +) łatwo podzielić zbiory liczb zawartych w kolumnie x -ów i y -ów na podzbiory. Jeżeli po podziale okaże się, że liczność podzbiorów jest niejednakowa — należy w obliczeniach oprzeć się na podzbiorach mniej licznych, gdyż skraca to rachunki.

Tablica i związane z nią obliczenia są proste i nie wymagają dalszych wyjaśnień.

Współczynniki regresji, znalezione metodą dwóch punktów, wyrażają się liczbami

$$a_y = 0,694, \quad a_x = 1,178.$$

Przy użyciu metody klasycznej otrzymalibyśmy następujące rezultaty:

$$A_y = 0,716, \quad A_x = 1,288.$$

Wprowadzimy następującą miarę rozbieżności między współczynnikami regresji, otrzymanymi obu metodami:

$$\Delta a_y / A_y = (a_y - A_y) / A_y;$$

nazywać ją będziemy *rozbieżnością względną współczynników regresji*.

W przytoczonym przykładzie jest

$$\frac{\Delta a_y}{A_y} = \frac{0,022}{0,716} = 0,031; \quad \frac{\Delta a_x}{A_x} = \frac{0,110}{1,288} = 0,085.$$

Wyniki uzyskane za pomocą obu metod mało się różnią między sobą⁽²⁾, natomiast obliczenia wymagane przez metodę dwóch punktów są znacznie mniej uciążliwe od obliczeń, związanych ze stosowaniem metody klasycznej.

Właśnie dlatego metoda dwóch punktów może oddać znaczne usługi praktyczne, zwłaszcza w tych przypadkach, gdy badanie związku korelacyjnego oparte jest na próbkach o dużej liczności.

(2) Warto podkreślić, że w praktyce na ogół większe znaczenie od rozbieżności między współczynnikami regresji ma rozbieżność między rzędnymi linii regresji otrzymanymi metodą klasyczną i metodą dwóch punktów.

Oto względna miara tej rozbieżności, dotycząca linii regresji y względem x :

$$\Delta y_x / Y_x = |y_x - Y_x| / Y_x,$$

gdzie

$$y_x = a_y x + y - a_y \bar{x}; \quad Y_x = A_y x + \bar{y} - A_y \bar{x}.$$

Łatwo wykazać, że jest $\Delta y_x / Y_x \leq \Delta a_y / A_y$.



TABLICA 1

| L. p. | x dziesiątków ton miesięcznie | y dziesiątków tys. kWh miesięcznie | $x > \bar{x}$ | $y x > \bar{x}$ | $y > \bar{y}$ | $x y > \bar{y}$ |
|-------|--|---|---------------|-----------------|---------------|-----------------|
| 1 | 183 | 175 | | | | |
| 2 | 184 | 172 | | | | |
| 3 | 180 | 168 | | | | |
| 4 | 164 | 156 | | | | |
| 5 | 177 | 190+ | | | 190 | 177 |
| 6 | 159 | 160 | | | | |
| 7 | 147 | 142 | | | | |
| 8 | 151 | 153 | | | | |
| 9 | 164 | 149 | | | | |
| 10 | 122 | 128 | | | | |
| 11 | 167 | 167 | | | | |
| 12 | 188 | 172 | | | | |
| 13 | 180 | 162 | | | | |
| 14 | 156 | 160 | | | | |
| 15 | 163 | 154 | | | | |
| 16 | 175 | 160 | | | | |
| 17 | 173 | 179 | | | | |
| 18 | 158 | 144 | | | | |
| 19 | 190 | 193+ | | | 193 | 190 |
| 20 | 180 | 169 | | | | |
| 21 | 196+ | 177 | 196 | 177 | | |
| 22 | 206+ | 190+ | 206 | 190 | 190 | 206 |
| 23 | 199+ | 186+ | 199 | 186 | 186 | 199 |
| 24 | 201+ | 180+ | 201 | 180 | 180 | 201 |
| 25 | 207+ | 182+ | 207 | 182 | 182 | 207 |
| 26 | 209+ | 190+ | 209 | 190 | 190 | 209 |
| 27 | 184 | 164 | | | | |
| 28 | 165 | 164 | | | | |
| 29 | 142 | 149 | | | | |
| 30 | 116 | 133 | | | | |
| 31 | 147 | 164 | | | | |
| 32 | 175 | 168 | | | | |
| 33 | 197+ | 176 | 197 | 176 | | |
| 34 | 202+ | 186+ | 202 | 186 | 186 | 202 |
| 35 | 189 | 183+ | | | 183 | 189 |
| 36 | 190 | 176 | | | | |
| 37 | 180 | 191+ | | | 191 | 180 |
| 38 | 170 | 167 | | | | |
| 39 | 182 | 161 | | | | |
| 40 | 189 | 180+ | | | 180 | 189 |
| 41 | 213+ | 191+ | 213 | 191 | 191 | 213 |
| 42 | 301+ | 264+ | 301 | 264 | 264 | 301 |
| 43 | 225+ | 202+ | 225 | 202 | 202 | 225 |
| 44 | 234+ | 214+ | 234 | 214 | 214 | 234 |

TABLICA 1 (cd.)

| L. p. | \bar{x} dziesiątków ton miesięcznie | \bar{y} dziesiątków tys. kWh miesięcznie | $x > \bar{x}$ | $y/x > \bar{x}$ | $y > \bar{y}$ | $x/y > \bar{y}$ |
|-------|--|---|-----------------------------|-----------------------------|---------------|-----------------|
| 45 | 203+ | 184+ | 203 | 184 | 184 | 203 |
| 46 | 192+ | 189+ | 192 | 189 | 189 | 192 |
| 47 | 191+ | 179 | 191 | 179 | | |
| 48 | 146 | 155 | | | | |
| 49 | 193+ | 173 | 193 | 173 | | |
| 50 | 187 | 166 | | | | |
| 51 | 187 | 182+ | | | 182 | 187 |
| 52 | 212+ | 205+ | 212 | 205 | 205 | 212 |
| 53 | 251+ | 216+ | 251 | 216 | 216 | 251 |
| 54 | 220+ | 201+ | 220 | 201 | 201 | 220 |
| 55 | 180 | 164 | | | | |
| 56 | 207+ | 191+ | 207 | 191 | 191 | 207 |
| 57 | 190+ | 173 | 190 | 173 | | |
| 58 | 185 | 174 | | | | |
| 59 | 186 | 170 | | | | |
| 60 | 181 | 166 | | | | |
| 61 | 192+ | 179 | 192 | 179 | | |
| 62 | 203+ | 191+ | 203 | 191 | 191 | 203 |
| 63 | 277+ | 247+ | 277 | 247 | 247 | 277 |
| 64 | 299+ | 257+ | 299 | 257 | 257 | 299 |
| 65 | 215+ | 206+ | 215 | 206 | 206 | 215 |
| 66 | 200+ | 188+ | 200 | 188 | 188 | 200 |
| 67 | 192+ | 167 | 192 | 167 | | |
| 68 | 187 | 183+ | | | 183 | 187 |
| 69 | 194+ | 190+ | 194 | 190 | 190 | 194 |
| 70 | 194+ | 182+ | 194 | 182 | 182 | 194 |
| 71 | 190 | 178 | | | | |
| 72 | 278+ | 241+ | 278 | 241 | 241 | 278 |
| Razem | 13712 $\bar{x} = 190,4$ | 12888 $\bar{y} = 179$ | 6693 $\bar{x}_I = 215,9$ | 6097 $\bar{y}_I = 196,7$ | 6175 | 6641 |

Otrzymujemy

$$a_y = \frac{196,7 - 179}{215,8 - 190,4} = \frac{17,7}{25,4} = 0,694;$$

podobnie

$$a_x = \frac{214,2 - 190,4}{199,2 - 179} = \frac{23,8}{20,2} = 1,178,$$

gdzie $214,2 = 6641/31$ oraz $199,2 = 6175/31$.

Prace cytowane

- [1] Б. В. Гнеденко, *Курс теории вероятностей*, Москва-Ленинград 1950.
 [2] M. G. Kendall, *The advanced theory of statistics*, II, London 1946.
 [3] К. Крамер, *Математические методы статистики*, Москва 1948.

Praca wpłynęła 26. 3. 1955

В. ГЕЛЛВИГ (Вроцлав)

ОПРЕДЕЛЕНИЕ ЛИНЕЙНОЙ РЕГРЕССИИ МЕТОДОМ ДВУХ ТОЧЕК

РЕЗЮМЕ

Чтобы определить расположение линии регрессии η относительно ξ методом двух точек, требуется найти центр тяжести точек, абсциссы которых больше (или меньше) \bar{x} .

Расположение линии регрессии η относительно ξ определяется при помощи центра тяжести всех точек и центра тяжести точек, абсциссы которых больше (или меньше) \bar{x} . Подобным образом определяется расположение линии регрессии ξ относительно η .

Этим методом получены следующие коэффициенты регрессии:

$$a_y = (\bar{y}_I - \bar{y}) / (\bar{x}_I - \bar{x}), \quad a_x = (\bar{x}_I - \bar{x}) / (\bar{y}_I - \bar{y}).$$

В статье доказано, что эти эstimаторы состоятельны и несмещены. Их относительная эффективность, при предположении, что распределение (ξ, η) нормальное, стремится по вероятности к $2/\pi$.

Распределение эstimаторов при возрастании выборки стремится к нормальному распределению.

Доказано автором, что

$$d_x = \left(\sum_{i=1}^n |x_i - \bar{x}| \right) / n = 2k|\bar{x}_I - \bar{x}| / n.$$

Полученная формула значительно упрощает исчисления связанные с вычислением среднего отклонения, особенно тогда, когда при вычислении отклонения используется упорядоченный ряд.

Определение параметров линии регрессии методом двух точек, благодаря своей простоте, может иметь большое практическое значение, главным образом в тех случаях, когда оценка коэффициентов регрессии основывается на большой выборке.

Z. HELLWIG (Wrocław)

DETERMINING LINEAR REGRESSION PARAMETERS BY MEANS OF THE TWO POINT METHOD

SUMMARY

In order to determine the position of the line of regression of η with respect to ξ by the two point method we determine the centre of gravity of those points whose abscissae are greater (or less) than \bar{x} .

The position of the line of regression of η with respect to ξ is determined by means of the centre of gravity of all points and the centre of gravity of those points whose abscissae are greater (less) than \bar{x} . In the same way we determine the position of the line of regression of ξ with respect to η .

The regression coefficients obtained by this method are expressed by the formulae

$$a_y = (\bar{y}_I - \bar{y}) / (\bar{x}_I - \bar{x}), \quad a_x = (\bar{x}_I - \bar{x}) / (\bar{y}_I - \bar{y}).$$

The author shows that these estimators are consistent and unbiased. Under the assumption that the distribution (ξ, η) is normal, their relative efficiency is convergent in probability, to $2/\pi$. With the growth of the sample the distribution of the estimators tends to a normal distribution.

In examining the efficiency of the estimators the author proves the following lemma:

$$d_x = \left(\sum_{i=1}^n |x_i - \bar{x}| \right) / n = 2k |\bar{x}_I - \bar{x}| / n.$$

The formula deduced greatly simplifies calculations connected with computing the mean deviation, particularly if we make use of an ordered series in calculating the deviation.

Owing to its simplicity, the method of determining regression line parameters may be of considerable practical value particular in those cases where the estimation of regression coefficients is made of the basis of a large sample.