

2変数統計データの線形近似

Theorem. n 個の2次元データ $(x_1, y_1), \dots, (x_n, y_n)$ について, 最小二乗法を用いて $y = ax + b$ ($a, b \in \mathbb{R}$) で近似したとき

$$a = \frac{s_{xy}}{s_x^2} = r_{xy} \cdot \frac{s_y}{s_x}, b = \bar{y} - a\bar{x}$$

をみtas. ただし, \bar{x} は x の平均, s_x は x の標準偏差, s_{xy} は x, y の共分散, r_{xy} は x, y の相関係数を表すとする.

Proof. 2乗誤差関数 E を

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

とし, E を最小にする a, b を x, y を用いて表せばよい. E を a, b で偏微分すると

$$\frac{\partial E}{\partial a} = -2 \sum_{i=1}^n x_i (y_i - ax_i - b) = 0 \quad (1)$$

$$\frac{\partial E}{\partial b} = -2 \sum_{i=1}^n (y_i - ax_i - b) = 0 \quad (2)$$

となる. (2) 式の両辺を $-2n$ で割ることにより $\bar{y} = a\bar{x} + b$ が得られる. (1) 式も両辺 $-2n$ で割ることにより

$$\frac{1}{n} \sum_{i=1}^n x_i y_i - a \cdot \frac{1}{n} \sum_{i=1}^n x_i^2 - b\bar{x} = 0$$

が得られる. ここで, $s_x^2 = \overline{x^2} - \bar{x}^2, \bar{y} = a\bar{x} + b$ を代入すると

$$\frac{1}{n} \sum_{i=1}^n x_i y_i - a(s_x^2 + \bar{x}^2) - b\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i y_i - as_x^2 - \bar{x}(a\bar{x} + b) = \frac{1}{n} \sum_{i=1}^n x_i y_i - as_x^2 - \bar{x} \cdot \bar{y} = 0$$

となるから, $a = \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \cdot \bar{y} \right)$ となる. ここで, もう少し式変形を進めると

$$\begin{aligned} a &= \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \cdot \bar{y} \right) \\ &= \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n \bar{x} y_i \right) \\ &= \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) y_i \right) \\ &= \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right) + \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) \bar{y} \right) \\ &= \frac{1}{s_x^2} \left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right) \\ &= \frac{s_{xy}}{s_x^2} \end{aligned}$$

となる. ここで, $r_{xy} = \frac{s_{xy}}{s_x s_y}$ より, $s_{xy} = r_{xy} s_x s_y$ を代入すると, $a = r_{xy} \cdot \frac{s_y}{s_x}$ となる. ■