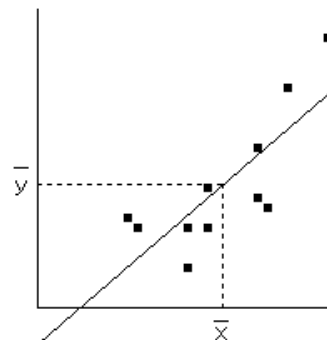


回帰直線

2つのデータ $X = \{x_1, x_2, \dots, x_n\}$, $Y = \{y_1, y_2, \dots, y_n\}$ を考えます. xy 平面に点 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ を描いた図を散布図といいます. 2つのデータの相関が強いとき, 散布図を見ると, ある直線の近くに点が集まっています. このような直線を回帰直線と言います.



それぞれのデータの, 平均を \bar{x}, \bar{y} , 標準偏差を σ_x, σ_y , 2つのデータの, 共分散を σ_{xy} , 相関係数を r とします. これらは左の式で定義され, 右の性質を満たします.

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{k=1}^n x_k & \sum_{k=1}^n (x_k - \bar{x}) &= 0 \\ \sigma_x^2 &= \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2 & \sigma_x^2 &= \frac{1}{n} \sum_{k=1}^n x_k^2 - \bar{x}^2 \\ \sigma_{xy} &= \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y}) & \sigma_{xy} &= \frac{1}{n} \sum_{k=1}^n x_k y_k - \bar{x} \bar{y} \\ r &= \frac{\sigma_{xy}}{\sigma_x \sigma_y} & \sigma_{xy} &= \sigma_x \sigma_y r \end{aligned}$$

上の性質の証明

$$\sum_{k=1}^n (x_k - \bar{x}) = \sum_{k=1}^n x_k - n\bar{x} = n\bar{x} - n\bar{x} = 0$$

$$\begin{aligned} \sigma_x^2 &= \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{1}{n} \sum_{k=1}^n (x_k^2 - 2\bar{x}x_k + \bar{x}^2) = \frac{1}{n} \sum_{k=1}^n x_k^2 - \frac{2}{n} \bar{x} \sum_{k=1}^n x_k + \bar{x}^2 \\ &= \frac{1}{n} \sum_{k=1}^n x_k^2 - 2\bar{x}^2 + \bar{x}^2 = \frac{1}{n} \sum_{k=1}^n x_k^2 - \bar{x}^2 \end{aligned}$$

$$\begin{aligned} \sigma_{xy} &= \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y}) = \frac{1}{n} \sum_{k=1}^n (x_k y_k - x_k \bar{y} - \bar{x} y_k + \bar{x} \bar{y}) \\ &= \frac{1}{n} \sum_{k=1}^n x_k y_k - \frac{1}{n} \sum_{k=1}^n x_k \bar{y} - \frac{1}{n} \bar{x} \sum_{k=1}^n y_k + \bar{x} \bar{y} = \frac{1}{n} \sum_{k=1}^n x_k y_k - \bar{x} \bar{y} - \bar{x} \bar{y} + \bar{x} \bar{y} \\ &= \frac{1}{n} \sum_{k=1}^n x_k y_k - \bar{x} \bar{y} \end{aligned}$$

最小2乗直線

直線 $y = f(x) = mx + n$ について、各 x_k に対する誤差 ($f(x_k)$ と実際の値 y_k との差) $f(x_k) - y_k$ の自乗の総和 $E = \sum_{k=1}^n (f(x_k) - y_k)^2$ を考えます。これが最小になるような直線 $y = f(x)$ が回帰直線の1つで、**最小2乗直線**といいます。

回帰直線は、中心 (\bar{x}, \bar{y}) を通る、すなわち、

$$y = f(x) = A(x - \bar{x}) + \bar{y}$$

と表されると推測されますが、推測が正しいとは限らないので

$$y = f(x) = A(x - \bar{x}) + \bar{y} + B$$

とおいて、 A, B を求めることにします。

$$\begin{aligned} E &= \sum_{k=1}^n (f(x_k) - y_k)^2 \\ &= \sum_{k=1}^n (A(x_k - \bar{x}) - (y_k - \bar{y}) + B)^2 \\ &= \sum_{k=1}^n (A^2(x_k - \bar{x})^2 - 2A(x_k - \bar{x})(y_k - \bar{y}) + (y_k - \bar{y})^2 + 2AB(x_k - \bar{x}) - 2B(y_k - \bar{y}) + B^2) \\ &= n\sigma_x^2 A^2 - 2n\sigma_{xy}A + n\sigma_y^2 + 0 - 0 + B^2 \\ &= n\sigma_x^2 A^2 - 2n\sigma_x\sigma_y rA + n\sigma_y^2 + B^2 \\ &= n(\sigma_x A - \sigma_y r)^2 + n\sigma_y^2(1 - r^2) + B^2 \end{aligned}$$

ゆえに、 E は、 $A = \frac{\sigma_y}{\sigma_x} r$ 、 $B = 0$ のときに最小になります。したがって、最小2乗直線は

$$y = \frac{\sigma_y}{\sigma_x} r(x - \bar{x}) + \bar{y} \tag{1}$$

です。相関係数 r が ± 1 のときは、 $E = 0$ すなわち、すべての点 (x_k, y_k) がこの直線上にあります。

回帰直線は、これだけではありません。今求めた直線は、 x を独立変数、 y を従属変数としていますが、逆に y を独立変数とすると、

$$\begin{aligned} x &= \frac{\sigma_x}{\sigma_y} r(y - \bar{y}) + \bar{x} \\ y &= \frac{\sigma_y}{\sigma_x} \frac{1}{r}(x - \bar{x}) + \bar{y} \end{aligned} \tag{2}$$

となります。(1) と (2) の傾きの相乗平均をとると

$$\begin{aligned} y &= \frac{\sigma_y}{\sigma_x}(x - \bar{x}) + \bar{y} \quad (r > 0 \text{ すなわち、正の相関があるとき}) \\ y &= -\frac{\sigma_y}{\sigma_x}(x - \bar{x}) + \bar{y} \quad (r < 0 \text{ すなわち、負の相関があるとき}) \end{aligned}$$

となります。