

2023年11月吉日

データの分析で学んだ公式からの考察
— 回帰直線と相関係数の関係 —

統計学 数理研究班

研究背景

数学Iで学んだデータの分析の単元で、相関係数 r は

$$-1 \leq r \leq 1$$

と解説してあるが、証明が与えられていない。

分布図の傾向を最もよく表す1次直線として、回帰直線を導入してるが、理論的な解説がない。

研究目的

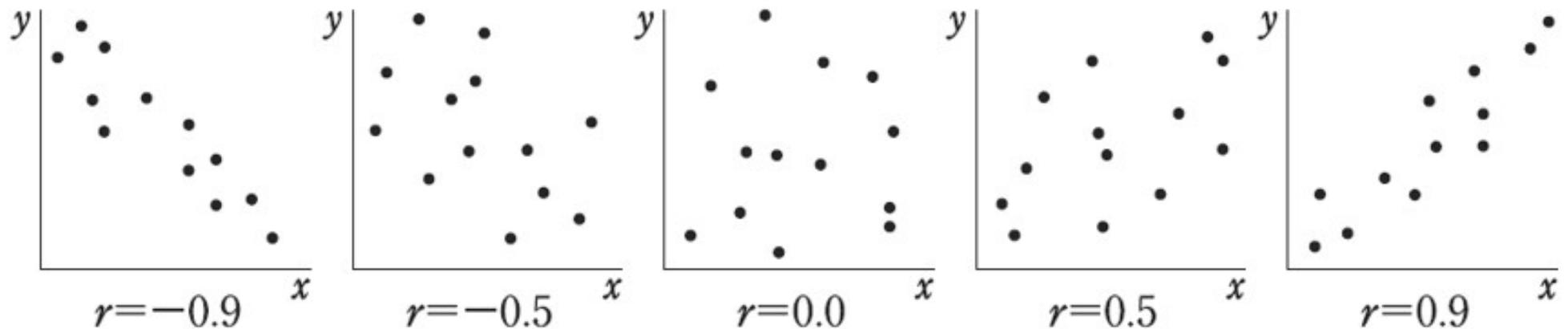
これらの理論的な根拠を学習するとともに、相関係数と回帰直線との関係性について研究を行った。

主に次の3点について研究を行った。

- 1 $-1 \leq r \leq 1$ の証明。
- 2 回帰直線の方程式を求める。
- 3 相関係数と回帰直線の傾きとの関係を研究。

正の相関，負の相関

2つの変量からなるデータにおいて，一方が増えると他方も増える傾向がみられるとき，2つの変量の間には正の相関があるという．また，一方が増えると他方が減る傾向がみられるとき，2つの変量の間には負の相関があるという．どちらの傾向もみられないときは，2つの変量の間には相関がないという．



相関係数

2つ変量 x , y からなるデータとして n 個の値の組

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

が得られているとするとき, x_1, x_2, \dots, x_n の平均値を \bar{x} , y_1, y_2, \dots, y_n の平均値を \bar{y} とおく.

$$\bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

$$\bar{y} = \frac{1}{n} (y_1 + y_2 + \dots + y_n)$$

相関係数 r は共分散 s_{xy} を， x の標準偏差 s_x と y の標準偏差 s_y の積 $s_x s_y$ で割った値として定義される。

$$\begin{aligned} r &= \frac{s_{xy}}{s_x s_y} \\ &= \frac{\frac{1}{n} \{ (x_1 - \bar{x})(y_1 - \bar{y}) + \cdots + (x_n - \bar{x})(y_n - \bar{y}) \}}{\sqrt{\frac{1}{n} \{ (x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2 \}} \sqrt{\frac{1}{n} \{ (y_1 - \bar{y})^2 + \cdots + (y_n - \bar{y})^2 \}}} \\ &= \frac{\frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{k=1}^n (y_k - \bar{y})^2}} \end{aligned}$$

$-1 \leq r \leq 1$ の証明

次の t の 2 次関数を考える ($f(t) \geq 0$).

$$\begin{aligned} f(t) &= \frac{1}{n} \sum_{k=1}^n \{t(x_k - \bar{x}) - (y_k - \bar{y})\}^2 \\ &= \frac{t^2}{n} \sum_{k=1}^n (x_k - \bar{x})^2 - \frac{2t}{n} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y}) \\ &\quad + \frac{1}{n} \sum_{k=1}^n (y_k - \bar{y})^2 \\ &= s_x^2 t^2 - 2s_{xy} t + s_y^2 \end{aligned}$$

$f(t) = s_x^2 t^2 - 2s_{xy}t + s_y^2$ は、 $f(t) \geq 0$ であるから、
係数について

$$D/4 = s_{xy}^2 - s_x^2 s_y^2 \leq 0$$

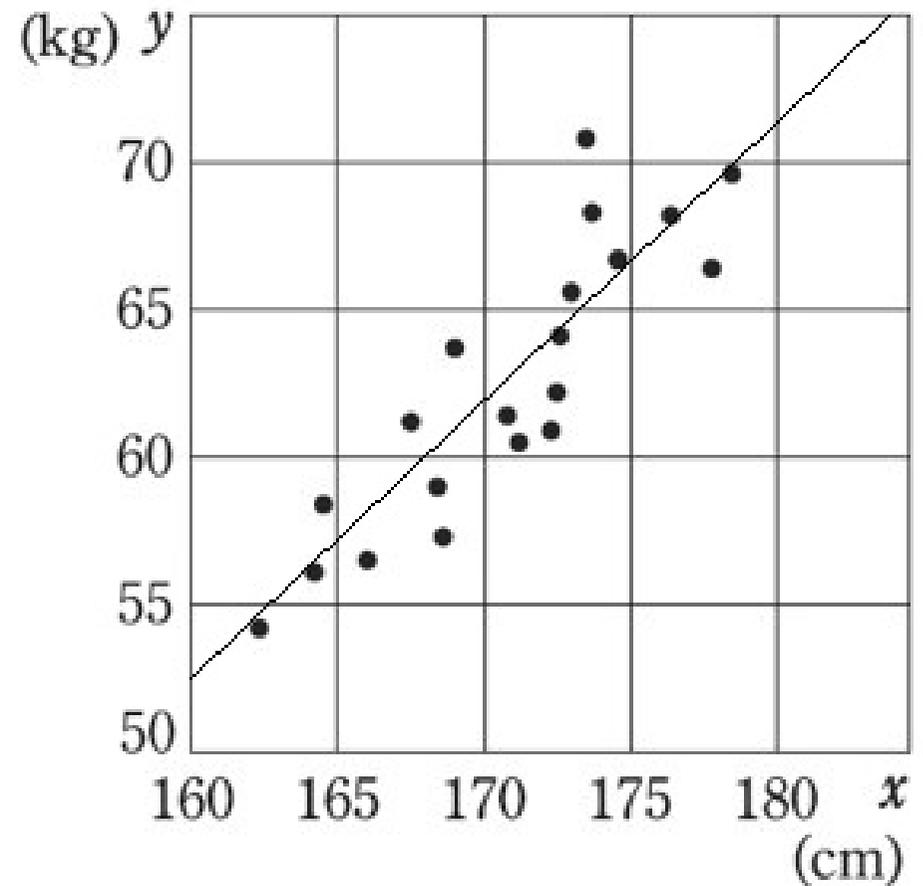
したがって

$$\frac{s_{xy}^2}{s_x^2 s_y^2} - 1 \leq 0 \quad \text{ゆえに} \quad r^2 - 1 \leq 0$$

これから $(r + 1)(r - 1) \leq 0$ よって $-1 \leq r \leq 1$

回帰直線

回帰直線とは、散布図上のデータに最もよく当てはまるよう引いた直線のことである。右の図はある高校の1年生男子の身長 x と体重 y の散布図について、引いた直線が回帰直線の例として示したものである。



2つ変量 x, y からなるデータとして n 個の値の組

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

について, $\frac{1}{n} \sum_{k=1}^n (y_k - ax_k - b)^2$ を最小にする定数 a, b を求める.

$$F(a, b) = \frac{1}{n} \sum_{k=1}^n (y_k - ax_k - b)^2$$

とおく.

F を a , b について偏微分すると

$$\begin{aligned}\frac{\partial F}{\partial a} &= -\frac{2}{n} \sum_{k=1}^n (y_k - ax_k - b)x_k \\ &= -2\{E(xy) - aE(x^2) - b\bar{x}\}\end{aligned}$$

$$\begin{aligned}\frac{\partial F}{\partial b} &= -\frac{2}{n} \sum_{k=1}^n (y_k - ax_k - b) \\ &= -2(\bar{y} - a\bar{x} - b)\end{aligned}$$

F が最小となるとき, $\frac{\partial F}{\partial a} = \frac{\partial F}{\partial b} = 0$ であるから

$$aE(x^2) + b\bar{x} = E(xy), \quad a\bar{x} + b = \bar{y}$$

したがって

$$a = \frac{E(xy) - \bar{x}\bar{y}}{E(x^2) - \bar{x}^2} = \frac{\sigma_{xy}}{\sigma_x^2}, \quad b = -a\bar{x} + \bar{y}$$

このとき，直線 $y = ax + b$ ，すなわち

$$y - \bar{y} = \frac{\sigma_{xy}}{\sigma_x^2} (x - \bar{x})$$

を回帰直線という。

結論

回帰直線

$$y - \bar{y} = \frac{\sigma_{xy}}{\sigma_x^2} (x - \bar{x})$$

によって、ある値 x_i に対する y の予測値を計算できる。また、相関係数 r と回帰直線の傾き a の符号が一致する。