

No.4 2次元のデータ

散布図と相関

例 10人の学生の数学と物理の得点

	1	2	3	4	5	6	7	8	9	10
数学	20	70	30	45	30	55	85	50	70	60
物理	15	50	25	50	20	50	80	40	85	75

- 平面上に縦軸と横軸をとり、横軸と縦軸にそれぞれ別の量をとり、データが当てはまる場所に点を打って（プロットするという）示したもの。

- ひとつのデータが高いと、もう一つのデータが高い傾向があるとき**正の相関**があるという。
- ひとつのデータが増える、もう一つのデータが小さくなる傾向があるとき**負の相関**があるという。

相関を表す量

- 相関を数値的に評価するとき、**相関係数**が用いられる。
- データを $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ とする。

- 次の量を**共分散**という $s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

$s_{xy} > 0$ のとき, x が大きければ y も大きい

$s_{xy} < 0$ のとき, x が大きければ y は小さい

と考えられる。

- 共分散の計算方法

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

- $\sigma_x = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$: x_1, \dots, x_n の標準偏差
- $\sigma_y = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$: y_1, \dots, y_n の標準偏差
- $r = \frac{s_{xy}}{\sigma_x \sigma_y}$ を**相関係数**という.

相関係数の性質

- 2つのデータ $\{x_i\}_i, \{y_i\}_i$ に対して, $z_i = tx_i - y_i$ の分散について考える:

$$\begin{aligned}
 0 \leq v_z &= \frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \{(tx_i - y_i) - (t\bar{x} - \bar{y})\}^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \{(x_i - \bar{x})t - (y_i - \bar{y})\}^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \{(x_i - \bar{x})^2 t^2 - 2(x_i - \bar{x})(y_i - \bar{y})t + (y_i - \bar{y})^2\} \\
 &= t^2 \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 - 2t \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\
 &= t^2 v_x - 2t s_{xy} + v_y
 \end{aligned}$$

- すべての t に対して, 0 以上より, D を判別式とすると,

$$\begin{aligned}
 D/4 &= (s_{xy})^2 - v_x v_y \leq 0 \\
 (s_{xy})^2 &\leq v_x v_y \\
 \frac{(s_{xy})^2}{v_x v_y} &\leq 1
 \end{aligned}$$

- よって, 相関係数 $r = \frac{s_{xy}}{\sigma_x \sigma_y}$ は $-1 \leq r \leq 1$ を満たす.
- 特に, 判別式 D が 0 ($|\frac{s_{xy}}{\sigma_x \sigma_y}| = 1$) のとき,

$$0 = v_z = t_0^2 v_x - 2t_0 s_{xy} + v_y$$

を満たす t_0 が存在する. 分散の定義から, すべての i について,

$$\begin{aligned}
 0 &= z_i - \bar{z} \\
 &= (t_0 x_i - y_i) - (t_0 \bar{x} - \bar{y}) \\
 y_i - \bar{y} &= t_0 (x_i - \bar{x})
 \end{aligned}$$

2つのデータ $\{x_i\}_i, \{y_i\}_i$ は 1 次式の関係にある.

例えば, $\frac{s_{xy}}{\sigma_x\sigma_y} = 1$, 即ち, $s_{xy} = \sigma_x\sigma_y$ のとき,

$$\begin{aligned}t_0^2 v_x - 2t_0 s_{xy} + v_y &= 0 \\t_0^2 v_x - 2t_0 \sigma_x \sigma_y + v_y &= 0 \\(t_0 \sigma_x - \sigma_y)^2 &= 0 \\t_0 &= \frac{\sigma_y}{\sigma_x}\end{aligned}$$

よって,

$$\begin{aligned}y_i - \bar{y} &= t_0(x_i - \bar{x}) = \frac{\sigma_y}{\sigma_x}(x_i - \bar{x}) = \frac{\sigma_x \sigma_y}{\sigma_x^2}(x_i - \bar{x}) = \frac{s_{xy}}{v_x}(x_i - \bar{x}) \\y_i &= \frac{s_{xy}}{v_x}x_i + \bar{y} - \frac{s_{xy}}{v_x}\bar{x}\end{aligned}$$

これを, y の x による **回帰直線** という.