

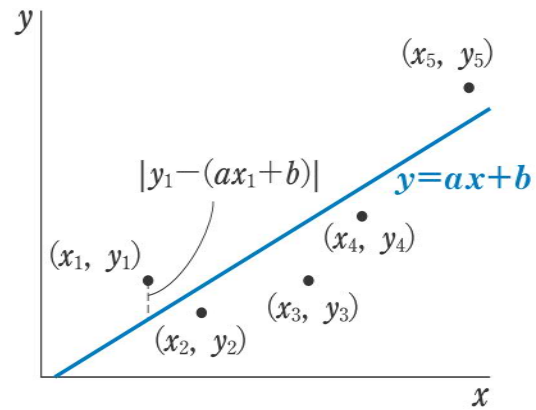
B 最小2乗法

回帰直線は、どのようにして求めることができるのだろうか。簡単のため、大きさ5のデータで考えることにする。

2つの変量 x , y のデータが、次のように与えられているとする。

$$5 \quad (x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4), (x_5, y_5)$$

x と y に直線的な相関関係があるとき、散布図の点は回帰直線の近くに分布する。



各点 (x_k, y_k) が、 $y = ax + b$ で表される直線上にあるとすると
 $y_k = ax_k + b$ であるが、実際のデータでは、ほとんどの場合 $y_k \neq ax_k + b$ である。そこで、 y_k と $ax_k + b$ の差の2乗の和

$$15 \quad \{y_1 - (ax_1 + b)\}^2 + \{y_2 - (ax_2 + b)\}^2 + \cdots + \{y_5 - (ax_5 + b)\}^2 \quad \text{①}$$

が最小となるように a , b を定め、直線 $y = ax + b$ を x , y の関係を近似する直線と考える。この直線 $y = ax + b$ が回帰直線の1つである。

このような回帰直線の求め方を、**最小2乗法** という。

20 このように、変量 x , y の間の関係をデータから統計的に推測する方法を **回帰分析** という。

2つの変量 x , y のデータの値の組が n 個与えられたとき、最小2乗法による回帰直線 $y = ax + b$ の a , b の値は、次のように求めることができる。

x, y のデータの平均値を \bar{x}, \bar{y} , 分散を s_x^2, s_y^2 , x と y の相関係数を r とするとき $a = \frac{s_y}{s_x} r, b = \bar{y} - a\bar{x}$ …… ②

練習 23

右の表は、同じ種類の5本の木の太さ x (cm) と高さ y (m) を測定した結果である。

木の番号	1	2	3	4	5
x	22	27	29	19	33
y	13	15	18	14	20

- (1) 2つの変数 x, y の散布図をかけ。
- (2) 2つの変数 x, y の回帰直線を表す1次関数を求めよ。また、その回帰直線を(1)の散布図に重ねてかけ。

参考

2つの変数 x, y のデータの値の組が n 個与えられたとき、回帰直線 $y = ax + b$ の a, b の値が上の②で求められることを確かめてみよう。

回帰直線は点 (\bar{x}, \bar{y}) を通る、すなわち $y = a'(x - \bar{x}) + \bar{y}$ と表されると推測されるが、推測が正しいとは限らないので

$$y = a'(x - \bar{x}) + \bar{y} + b'$$

とおいて、 a', b' を求めることにする。前ページの①の値は

$$\begin{aligned} \sum_{k=1}^n \{y_k - \{a'(x_k - \bar{x}) + \bar{y} + b'\}\}^2 &= \sum_{k=1}^n \{a'(x_k - \bar{x}) - (y_k - \bar{y}) + b'\}^2 \\ &= \sum_{k=1}^n \{a'^2(x_k - \bar{x})^2 - 2a'(x_k - \bar{x})(y_k - \bar{y}) + (y_k - \bar{y})^2 \\ &\quad + 2a'b'(x_k - \bar{x}) - 2b'(y_k - \bar{y}) + b'^2\} \\ &= ns_x^2 a'^2 - 2ns_{xy} a' + ns_y^2 + 0 - 0 + nb'^2 = ns_x^2 a'^2 - 2ns_x s_y r a' + ns_y^2 + nb'^2 \\ &= n(s_x a' - s_y r)^2 + ns_y^2(1 - r^2) + nb'^2 \end{aligned}$$

ゆえに、①の値は、 $a' = \frac{s_y}{s_x} r, b' = 0$ のとき最小になる。

よって、 a, b の値は $a = \frac{s_y}{s_x} r, b = \bar{y} - a\bar{x}$ である。

【補足】 記号 $\sum_{k=1}^n$ については、第1章「数列」の23～25ページを参照。

5

10

15

20