

相関とは

もりしま みつる
森島 充

「相関とは要するに何だろう?」と思いながら定義式を見ていたら「cos だ」と気がつきました。これを足掛かりに、相関の図形的意味について考えます。

§1. 定義

2つのデータ,

$$x_1, x_2, x_3, \dots, x_n \quad y_1, y_2, y_3, \dots, y_n$$

に対して, n 次元空間の点,

$$X(x_1, x_2, x_3, \dots, x_n) \quad Y(y_1, y_2, y_3, \dots, y_n)$$

$$P(\bar{x}, \bar{x}, \bar{x}, \dots, \bar{x}) \quad Q(\bar{y}, \bar{y}, \bar{y}, \dots, \bar{y})$$

を考えます。

このとき, $\overrightarrow{PX} \neq \vec{0}$, $\overrightarrow{QY} \neq \vec{0}$ であれば, 相関係数 r は,

$$r = \frac{\overrightarrow{PX} \cdot \overrightarrow{QY}}{|\overrightarrow{PX}| |\overrightarrow{QY}|}$$

と書けます。

すなわち, \overrightarrow{PX} と \overrightarrow{QY} の n 次元空間における「なす角」を θ とすると,

$$r = \cos \theta$$

です。これより, $-1 \leq r \leq 1$ であり,

$$r = 1 \iff \overrightarrow{PX} \parallel \overrightarrow{QY} \text{ (同じ向き)}$$

$$r = -1 \iff \overrightarrow{PX} \parallel \overrightarrow{QY} \text{ (逆向き)}$$

$$r = 0 \iff \overrightarrow{PX} \perp \overrightarrow{QY}$$

となります。

§2. n 次元空間における図形的意味

原点を O , 点 $(1, 1, 1, \dots, 1)$ を E とすると, $\overrightarrow{PX} = (x_1 - \bar{x}, x_2 - \bar{x}, x_3 - \bar{x}, \dots, x_n - \bar{x})$ より,

$$\overrightarrow{OP} \cdot \overrightarrow{PX} = \bar{x} \overrightarrow{OE} \cdot \overrightarrow{PX}$$

$$= \bar{x} \left\{ \sum_{j=1}^n (x_j - \bar{x}) \right\}$$

$$= 0$$

です。したがって $\overrightarrow{OP} \neq \vec{0}$, $\overrightarrow{PX} \neq \vec{0}$ であれば, $\overrightarrow{OP} \perp \overrightarrow{PX}$ です。

このことは標準偏差の公式

$$\sigma = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2} = \sqrt{\frac{1}{n} \sum_{j=1}^n x_j^2 - (\bar{x})^2}$$

を変形して,

$$\sum_{j=1}^n (x_j - \bar{x})^2 = \sum_{j=1}^n x_j^2 - n(\bar{x})^2$$

$$\sum_{j=1}^n (x_j - \bar{x})^2 = \sum_{j=1}^n x_j^2 - \sum_{j=1}^n (\bar{x})^2$$

$$\sum_{j=1}^n (x_j - \bar{x})^2 + \sum_{j=1}^n (\bar{x})^2 = \sum_{j=1}^n x_j^2$$

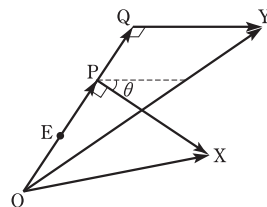
すなわち,

$$PX^2 + OP^2 = OX^2$$

であることから導かれます。

また, $\overrightarrow{OP} \perp \overrightarrow{PX}$ は, \overrightarrow{OP} が \overrightarrow{OX} の直線 OE への正射影であることも示しています。

2点 P, Q は直線 OE 上にあるので, これを図にすると右のようになります。



この図を見ると, 改めて, 相関とは点 X, Y の,

点 P, Q からの「ずれの方向」の違いを表す指標であることが分かります。

しかし, 相関が強いときはともかく, 相関が弱いときは「ずれの方向」の自由度はかなり高くなります。

上の図は3次元的な図なので, 例えば相関係数が 0 , すなわち $\overrightarrow{PX} \perp \overrightarrow{QY}$ のとき, \overrightarrow{PX} が決まれば \overrightarrow{QY} は大きさを除いて固定されるような印象がありますが, $\overrightarrow{OQ} \cdot \overrightarrow{QY} = 0$ と $\overrightarrow{PX} \cdot \overrightarrow{QY} = 0$ の2つの条件しかないので, n 次元空間であれば, \overrightarrow{QY} は $n-2$ 次元の自由度をもちます。

また, この図から, $n=2$ のときは, 相関係数は $r = \pm 1$ に限られ, $n=1$ のときは, $\overrightarrow{PX} = \vec{0}$, $\overrightarrow{QY} = \vec{0}$ となるので相関係数が存在しないことが分かります。

関連して、

$$\sigma = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2} = \frac{|\overrightarrow{PX}|}{\sqrt{n}}$$

ですから、標準偏差 σ は、点 X の点 P からの「ずれの大きさ」を表していることが分かります。 \sqrt{n} で割っているのは、データの値の個数が変化しても指標として比較できるように調整するためだと考えることができます。

次に散布図における相関の図形的意味を考えます。

§3. 散布図における図形的意味

散布図とは、2つのデータ、

$$x_1, x_2, x_3, \dots, x_n \quad y_1, y_2, y_3, \dots, y_n$$

に対して、 xy 平面の点、

$$A_j(x_j, y_j) \quad (j=1, 2, 3, \dots, n)$$

を図示したものでした。

点 (\bar{x}, \bar{y}) を G とすると、点 G は質量が等しい n 個の質点 A_j の重心になります。

ここで、簡単のためにデータを「平行移動」して、以下は点 G が原点 O に重なるようにします。

すなわち、改めて $\bar{x}=0, \bar{y}=0$ であるデータ

$$x_1, x_2, x_3, \dots, x_n \quad y_1, y_2, y_3, \dots, y_n$$

について、

$$A_j(x_j, y_j) \quad (j=1, 2, 3, \dots, n)$$

を定めることにします。

このとき、§1 の n 次元空間の点 P, Q も原点 O に重なるので、

$$r = \pm 1 \text{ のとき、} \overrightarrow{OX} \parallel \overrightarrow{OY}$$

となります。したがって、 k を実数として、

$$\overrightarrow{OY} = k \overrightarrow{OX}$$

と書けますから、 xy 平面上の点 A_j の座標は、

$$(x_j, y_j) = (x_j, kx_j)$$

となり、点 A_j は直線 $y=kx$ 上に並びます。相関係数を考える際は、 $\overrightarrow{OX} \neq \vec{0}, \overrightarrow{OY} \neq \vec{0}$ ですから、 $k \neq 0$ です。

データのスケールを適当に変えれば、 k の値も変えられるので、 $k=1$ か $k=-1$ に限定することもできます。ですから相関図の「傾き」にはあまり意味がなさそうですが、経年変化などを考える場合には意味があるといえます。

次に、 $r=0$ のときを考えます。

このとき、 n 次元空間において、

$$\overrightarrow{OX} \cdot \overrightarrow{OY} = \sum_{j=1}^n x_j y_j = 0$$

となりますが、 n 次元空間で $\overrightarrow{OX} \perp \overrightarrow{OY}$ であることを、2次元平面の散布図で考えようと思っても簡単ではなさそうです。しかし、内積のもう1つの意味が「面積の和」であることを利用すると簡単にイメージがもてます。

点 A_j を頂点とする右の図のような長方形の面積を

S_j とすると、 $x_j y_j$ の値は

第1・3象限で

$$x_j y_j = S_j$$

第2・4象限で

$$x_j y_j = -S_j$$

です。

したがって、 $r=0$ とは、

「点 $A_j (j=1, 2, 3, \dots, n)$ の重心を原点として、第1・3象限にできる長方形の面積 S_j の和と第2・4象限にできる長方形の面積 S_j の和が等しいことである」…(*)

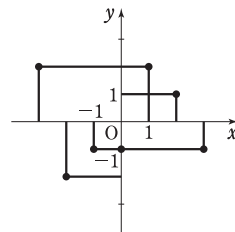
といえます。

これによって、 $r=0$ のとき、点 A_j は原点の周りにバランスよく散らばることになります。

しかし、(*)はそれ以上の意味をもたないので、同じ $r=0$ でも様々な例を考えることができます。

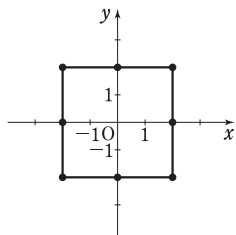
例1

x	-3	-2	-1	0	1	2	3
y	2	-2	-1	-1	2	1	-1



例 2

x	-2	-2	0	2	2	2	0	-2
y	0	2	2	2	0	-2	-2	-2



例 3

$$(x_j, y_j) = \left(\cos \frac{2j}{n} \pi, \sin \frac{2j}{n} \pi \right)$$

$$(n \geq 3, j = 0, 1, 2, \dots, n-1)$$

例 2, 例 3 は, とても「相関係数 0」とは思えません, $\overrightarrow{OX} \perp \overrightarrow{OY}$ であっても n 次元空間では $n-2$ 次元の自由度をもつことや, 散布図では (*) の意味しかないことを考えると理解できます。

逆に, 相関が弱いからといって「法則性がない」とは判断できない, という点に注意する必要があります。

(東京都立調布南高等学校)