

線形代数で眺める確率統計

@minami106

2013年2月4日

確率統計では、データの特徴を代表して表すさまざまな指標が出現します。それらはわりとごちゃごちゃしたためんどくさい式によって定義されるのですが、ベクトルや、ベクトルの演算を用いると、びっくりするくらい簡潔かつ明快にいろんなことが記述できることにふと気づいたので、まとめてみたいと思います。

1 いろんなベクトルの定義

実際に観測されたデータを $x_1, x_2, \dots, x_n (\in \mathbb{R})$ とし、データから求めた算術平均を \bar{x} とします。

定義 1.1. データベクトル

$\mathbf{x} \in \mathbb{R}^n$ を次のように定義し、**データベクトル (data vector)** と呼ぶ。

$$\mathbf{x} \stackrel{\text{def}}{=} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

定義 1.2. 平均ベクトル

$\bar{\mathbf{x}} \in \mathbb{R}^n$ を次のように定義し、**平均ベクトル (mean vector)** と呼ぶ。 \bar{x} は、データから求めた算術平均である。

$$\bar{\mathbf{x}} \stackrel{\text{def}}{=} \begin{pmatrix} \bar{x} \\ \bar{x} \\ \vdots \\ \bar{x} \end{pmatrix}$$

定義 1.3. 偏差ベクトル

$\tilde{\mathbf{x}} \in \mathbb{R}^n$ を次のように定義し、**偏差ベクトル (deviation vector)** と呼ぶ。

$$\tilde{\mathbf{x}} \stackrel{\text{def}}{=} \mathbf{x} - \bar{\mathbf{x}} = \begin{pmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{pmatrix}$$

2 標本分散, 標本標準偏差

定義 2.1. 標本分散

次のように定義される $V_{\mathbf{x}}$ を, \mathbf{x} の標本分散 (sample variance) という.

$$V_{\mathbf{x}} = \frac{1}{n} \|\tilde{\mathbf{x}}\|^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

偏差ベクトルは, データベクトルが平均 (ベクトル) からどの向きにどのくらいずれているか? ということを表すベクトルと捉えることができます. このことから, データの標本分散というのは, データベクトルの, 平均ベクトルからのズレの大きさを, 偏差ベクトルのノルムという形で測ったものに他ならないということが分かります (図 1) *1.

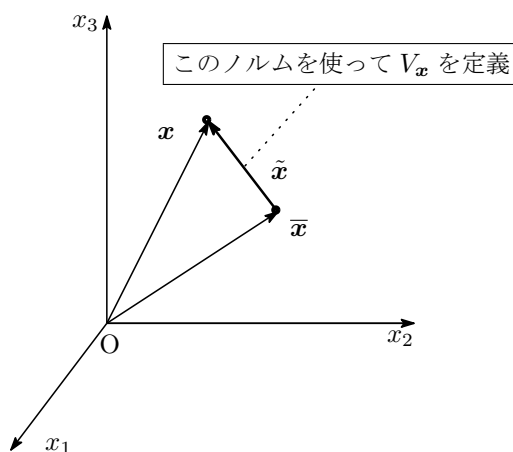


図 1 標本分散の定義のイメージ (データ数 $n = 3$ の場合)

定義 2.2. 標本標準偏差

$\sigma_{\mathbf{x}} \stackrel{\text{def}}{=} \sqrt{V_{\mathbf{x}}}$ と定義し, $\sigma_{\mathbf{x}}$ を標本標準偏差と呼ぶ.

なんか, この, ノルムによる定義 (ノルムは自分自身との内積でかけるので, 内積を使った定義でもある) に, ベクトルの内積の性質を適用して, 標本分散や標本標準偏差に関する性質をいろいろ導き出せないですかね. やって見たら面白いかもですね.

*1 データ数 n で $\|\tilde{\mathbf{x}}\|^2$ を割っているのは, データ数をたくさんにするとノルムが際限なく増えてしまうので, 正規化するためにかけている定数です. データ数で割ることにより, ひとつのデータあたり, 大体どのくらい平均からズレているかを考えていると思うとわかりやすいかもです.

3 標本共分散, 相関係数

今度は、ふたつのデータ列 $x_1, \dots, x_n \in \mathbb{R}, y_1, \dots, y_n \in \mathbb{R}$ を考えます。ふたつのデータ列*2に対して、それぞれのデータベクトルを \mathbf{x}, \mathbf{y} 、偏差ベクトルを $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$ と置いた時、**標本共分散**を定義することができます。

定義 3.1. 標本共分散

次のように定義される $\text{Cov}(\mathbf{x}, \mathbf{y})$ を、 \mathbf{x}, \mathbf{y} の**標本共分散 (sample covariance)** と呼ぶ

$$\text{Cov}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \frac{1}{n}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

定義 3.2. 相関係数

次のように定義される $r(\mathbf{x}, \mathbf{y})$ を、 \mathbf{x}, \mathbf{y} の**相関係数 (correlation coefficient)** と呼ぶ。

$$r(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \frac{(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})}{\|\tilde{\mathbf{x}}\| \|\tilde{\mathbf{y}}\|} = \frac{\text{Cov}(\mathbf{x}, \mathbf{y})}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \cos \theta.$$

ただし、 θ は、 $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$ が \mathbb{R}^n でなす角である。

まず、 $\cos \theta$ の性質から、相関係数の次の性質がすぐに分かります。

命題 3.1. (相関係数の値)

$$-1 \leq r(\mathbf{x}, \mathbf{y}) \leq 1.$$

さて、相関係数がどんなものかってことは、統計学の授業なんかで何となく習ったことがある方も多いのではなかろうかと思いますが、なぜそういう性質を持つのかってことは知らない人が結構いるのではないかと。実は、相関係数の意味は、**相関係数は偏差ベクトルのなす角 θ の余弦 $\cos \theta$ である**という視点に立つと、相関係数とは何たるかがものすごくスッキリわかるのです。

3.1 相関係数の意味

さて、偏差ベクトルは、**データベクトルが平均どっち向きに、どのくらいずれているかを表すベクトル**でした (図 2)。いま、偏差ベクトル $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$ と、なす角の余弦 $\cos \theta$ について、図 3 のような関係が成立します。

図 3 のように、

- 偏差ベクトルどうしが**同じ向き** のとき、 $r = \cos \theta$ は**最大 (+1)**。
- 偏差ベクトルどうしが**そっぽ向き** のとき、 $r = \cos \theta$ は**0**。
- 偏差ベクトルどうしが**逆向き** のとき、 $r = \cos \theta$ は**最小 (-1)**。

*2 たとえば、 x_i が出席番号 i 番の学生の国語の点数、 y_i が数学の点数だったりするわけです。

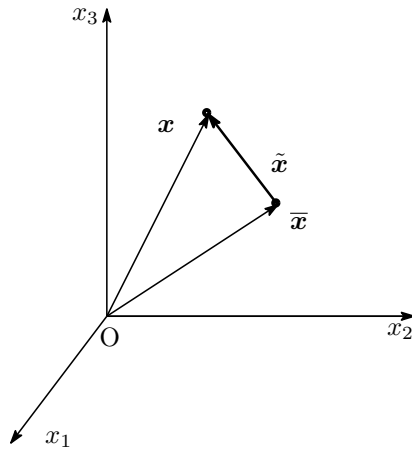


図2 偏差ベクトル (\tilde{y} についても同様のイメージ)

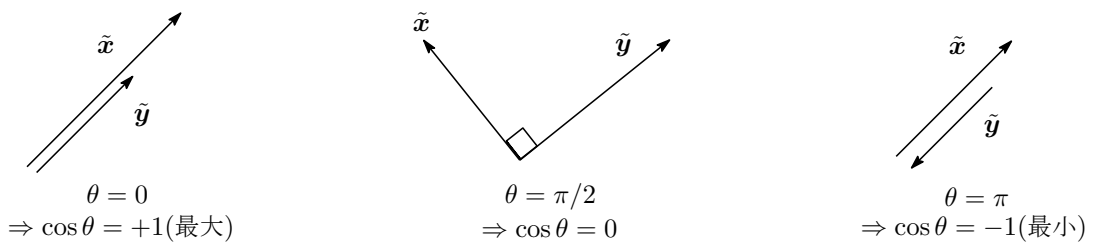


図3 \tilde{x}, \tilde{y} の向き関係と $\cos \theta$ の値

ということが言えます。つまり $\cos \theta$ は、ベクトルがどのくらい同じ向きを向いているかの指標といえるのです。さて、これから、相関係数の3つの値のパターンについて、データはどのような挙動を見せるのか、考えてみましょう。

1. $0 < r \leq +1$ のとき (相関係数が正) .

このとき、 $\cos \theta > 0$ なので、偏差ベクトル \tilde{x}, \tilde{y} は、比較的同じ方向を向いているということが言えます。つまり、 x が平均からズレたのと大体同じ方向に、 y もズレる傾向があるということです (図4)。こんな傾向にのっとって、データが生起しているということを頭に入れて、データをプロットしまくっ

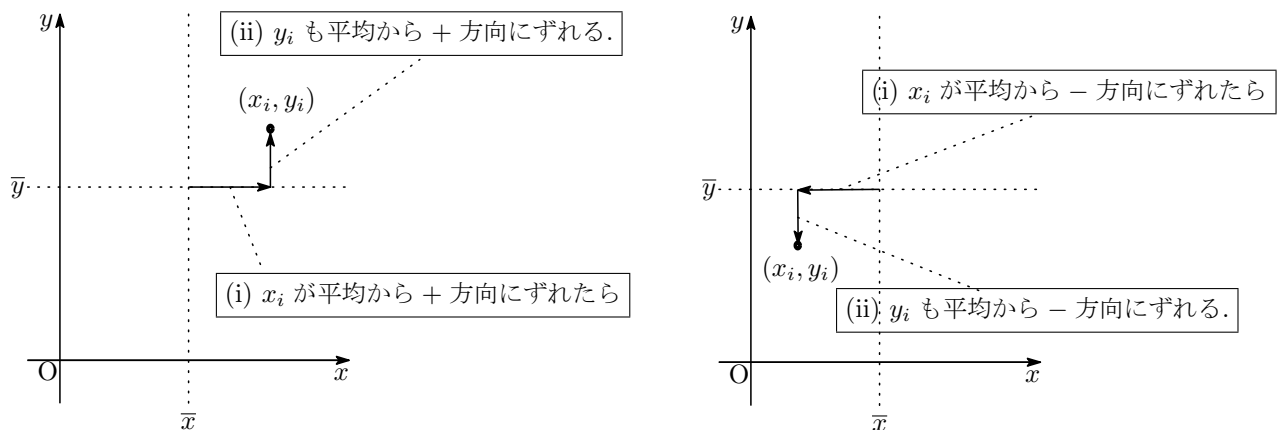


図4 $r = \cos \theta > 0$ のときのデータの傾向.

てみると、データはおおよそ図5のように分布するだろうということが分かります。データが x_i, y_i が

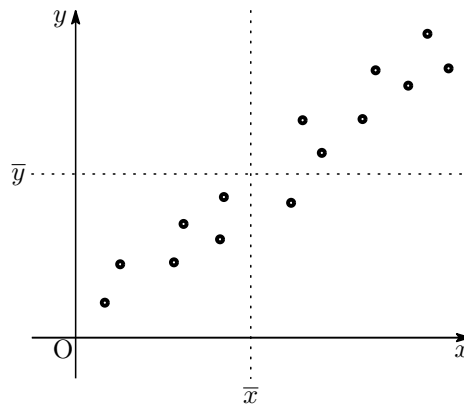


図5 $r = \cos \theta > 0$ のときのデータの分布の感じ

こんな感じに分布することを、**正の相関**があるといいます。

定義 3.3. 正の相関

$r(x, y) > 0$ のとき、データ x, y は**正の相関 (plus correlation)** を持つといいます。

2. $-1 \leq r < 0$ のとき (相関係数が負) .

このとき、 $\cos \theta < 0$ なので、偏差ベクトル \tilde{x}, \tilde{y} は、比較的逆の方向を向いているということが言えます。つまり、 x が平均からズレたのと大体逆の方向に、 y もズレる傾向があるということです (図6)。こんな傾向にのっかって、データが生起しているということを頭に入れて、データをプロットしまくっ

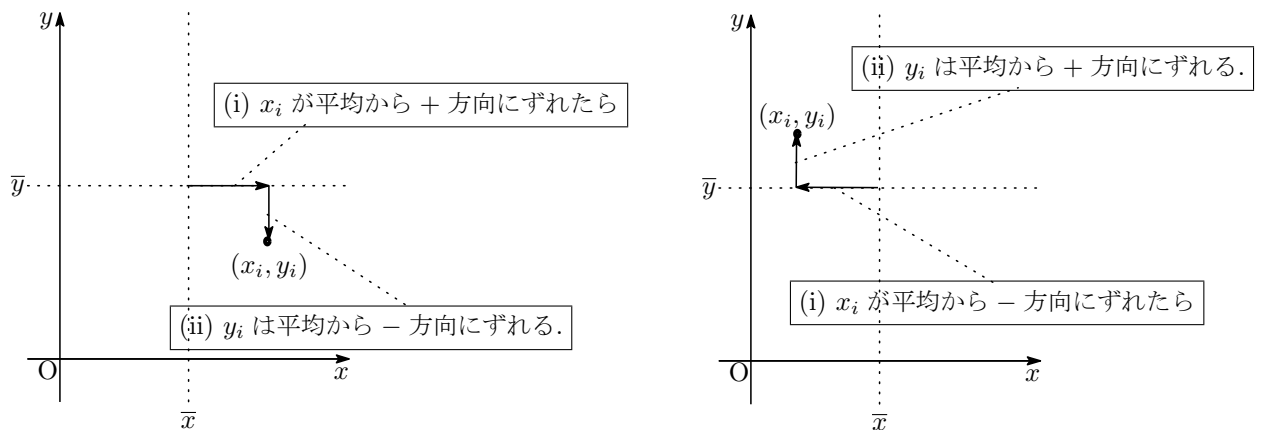


図6 $r = \cos \theta < 0$ のときのデータの傾向.

てみると、データはおおよそ図7のように分布するだろうということが分かります。データが x_i, y_i がこんな感じに分布することを、**負の相関**があるといいます。

定義 3.4. 負の相関

$r(x, y) < 0$ のとき、データ x, y は**負の相関 (minus correlation)** を持つといいます。

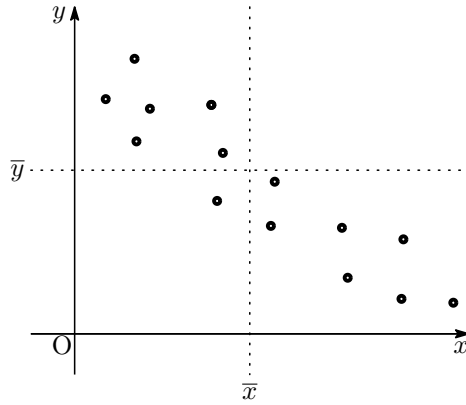


図7 $r = \cos \theta < 0$ のときのデータの分布の感じ

3. $r = 0$ のとき (相関係数が 0)

このとき, $\cos \theta = 0$ ということなので, x がどっち向きに平均からズレようが, y がどっち向きにずれるのかは分かりません. このとき, データの分布は図8のようになります.

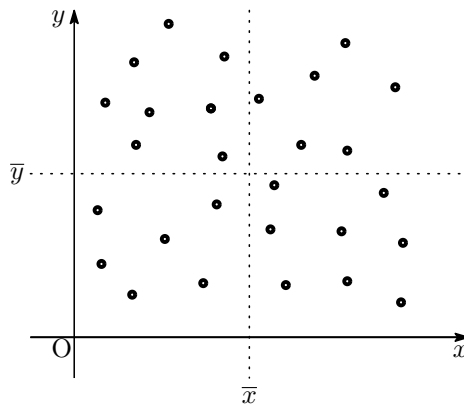


図8 $r = 0$ のときのデータの分布の感じ.

データが x_i, y_i がこんな感じに分布することを, **相関を持たない**といいます.

定義 3.5. 相関を持たない

$r(x, y) = 0$ のとき, データ x, y は相関を持たないといいます.

こんなふうに, 相関係数を偏差ベクトルのなす角の余弦と捉えることによって, 相関係数が持つ意味がこんなにもスッキリと見えてしまいました. 以下に, 相関係数の性質をまとめましょう.

命題 3.2. 相関係数の意味

相関係数 $r(x, y)$ について,

$$\begin{cases} 0 < r \leq 1 & \Rightarrow \text{データ } x, y \text{ は正の相関を持つ.} \\ r = 0 & \Rightarrow \text{データ } x, y \text{ は相関を持たない.} \\ -1 \leq r < 0 & \Rightarrow \text{データ } x, y \text{ は負の相関を持つ.} \end{cases}$$

こんなふうに、線形代数的な視点で確率統計で出てくる重要な指標を眺めてみると、随分明快な理解が得られるのだなあ実感してしまいます。この資料は相関係数のところについて書きたかったという理由で書いてみたものですが、他にも確率統計を線形代数のレンズで調べてみたときに得られる面白い結果というのは色々あるのでしょうか。気になるところですね～。面白い話があったら是非是非教えてほしいです。

参考

- Wikipedia - 相関係数
<http://ja.wikipedia.org/wiki/%E7%9B%B8%E9%96%A2%E4%BF%82%E6%95%B0>