

### V-3-3. マハラノビスの距離

多変量解析では、しばしば、データの類似性が議論されます。一般的には、類似性は多次元空間における距離によって評価されます。したがって、2つのデータの類似性が高い場合、それらのデータは多次元空間で近い場所にプロットされます。しかし、測定項目間に相関があったり分散が大きく違っている場合には、それらを考慮しなければ、距離をそのまま類似性だと考えることはできません。

二つのデータ  $\mathbf{a}$  と  $\mathbf{b}$  を次のようにベクトルで表します。

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

このように表現すると、このデータ間の距離はベクトルの矢印の先端間の距離  $|\mathbf{a} - \mathbf{b}|$  となります。

$$\mathbf{a} - \mathbf{b} = \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix}$$

$$\begin{aligned} |\mathbf{a} - \mathbf{b}| &= \sqrt{(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})} \\ &= \sqrt{\begin{pmatrix} a_1 - b_1 & a_2 - b_2 & a_3 - b_3 \end{pmatrix} \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix}} \\ &= \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + (a_3 - b_3)^2} \end{aligned}$$

これらのベクトルが行列  $\mathbf{A}$  によって変換され  $\mathbf{Aa}$ 、 $\mathbf{Ab}$  になったとすると、変換されデータ間の距離は次のようになります。

$$\begin{aligned} |\mathbf{Aa} - \mathbf{Ab}| &= |\mathbf{A}(\mathbf{a} - \mathbf{b})| \\ &= \sqrt{(\mathbf{A}(\mathbf{a} - \mathbf{b}))^T \mathbf{A}(\mathbf{a} - \mathbf{b})} \\ &= \sqrt{\begin{pmatrix} a_1 - b_1 & a_2 - b_2 & a_3 - b_3 \end{pmatrix} \mathbf{A}^T \mathbf{A} \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix}} \\ |\mathbf{A}(\mathbf{a} - \mathbf{b})|^2 &= \begin{pmatrix} a_1 - b_1 & a_2 - b_2 & a_3 - b_3 \end{pmatrix} \mathbf{A}^T \mathbf{A} \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ a_3 - b_3 \end{pmatrix} \end{aligned}$$

$$|(\mathbf{a} - \mathbf{b})\mathbf{A}|^2 = (\mathbf{a} - \mathbf{b})^T \mathbf{A}^T \mathbf{A} (\mathbf{a} - \mathbf{b})$$

ここで  $\mathbf{a} - \mathbf{b} = \mathbf{x}$  と表すことにします。

$$\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}$$

ここで、 $\mathbf{A}$  が分散共分散行列の平方根であったとします。

$$\mathbf{A} = \boldsymbol{\Sigma}^{-\frac{1}{2}}$$

$\Sigma$ : 分散共分散行列

$$\begin{aligned} |Ax| &= \sqrt{\left(\Sigma^{-\frac{1}{2}}x\right)^T \Sigma^{-\frac{1}{2}}x} \\ &= \sqrt{x^T \left(\Sigma^{-\frac{1}{2}}\right)^T \Sigma^{-\frac{1}{2}}x} = \sqrt{x^T \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} x} = \sqrt{x^T \Sigma^{-1} x} \end{aligned}$$

$$(\because \Sigma^{-\frac{1}{2}} \text{ は対称行列ですから, } \left(\Sigma^{-\frac{1}{2}}\right)^T = \Sigma^{-\frac{1}{2}})$$

これをマハラノビスの距離と言います。

$$D_{a-b} = \sqrt{(a-b)^T \Sigma^{-1} (a-b)}$$

( $\Sigma$ : 分散共分散行列)

式 70

マハラノビスの距離は分散共分散による歪みを補正した標準化された距離です。項目間に相関があったり、項目間の分散違いが大きい場合に、そのまま距離を計算する（ユークリッド距離）ことが適切でないと判断して、マハラノビス距離を計算することは、クラスター分析等でしばしば行われます。

Prasanta Chandra Mahalanobis はインドの統計学者で、インド統計研究所を創設しました。Fisher や Pearson とも交友があり、インドの社会主義的な計画経済の政策立案にも関与し、理論とともに実践面でも活躍しました。