

# 主成分分析と合成変数

多変量解析では、複数個の変数の重み付き合計値を用いる。解析の目的により各種の重み付けを行う。 $p$  個の変数  $x_1, x_2, \dots, x_p$  が、重み  $w_1, w_2, \dots, w_p$  で重み付けされた  $f$  を、合成変数と呼ぶ。

$$f = w_1 x_1 + w_2 x_2 + \dots + w_p x_p \quad (1)$$

ここで、 $x_i$  はそれぞれ標準化されており、それぞれの変数が  $n$  ケースについて測定されているとき、 $n \times p$  の大きさのデータ行列を  $\mathbf{X}$  で表すことにする。

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & x_{ij} & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \quad (2)$$

変数  $x_j$  の平均値  $E(x_j)$  と分散  $V(x_j)$  はそれぞれ定義により、

$$E(x_j) = \frac{1}{n} \sum_{i=1}^n x_{ij} = 0 \quad (3)$$

$$V(x_j) = \frac{1}{n} \sum_{i=1}^n x_{ij}^2 = 1 \quad (4)$$

となる。

また、変数  $x_j$  と  $x_k$  の共分散  $Cov_{jk}$  と相関係数  $r_{jk}$  は、

$$Cov_{jk} = \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} \quad (5)$$

$$r_{jk} = \frac{Cov(x_j, x_k)}{\sqrt{V(x_j) V(x_k)}} = Cov(x_j, x_k) \quad (6)$$

となる。相関係数行列  $\mathbf{R}$  は、データ行列  $\mathbf{X}$  を用いて、

$$\mathbf{R} = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1p} \\ r_{21} & r_{22} & \dots & r_{2p} \\ \vdots & \vdots & r_{ij} & \vdots \\ r_{p1} & r_{p2} & \dots & r_{pp} \end{pmatrix} = \frac{1}{n} \mathbf{X}' \mathbf{X} \quad (7)$$

のように表される ( $\mathbf{X}'$  は転置行列を表す)。

**重みベクトル**を  $\mathbf{w}' = (w_1, w_2, \dots, w_p)$  として、

$$\mathbf{f} = \mathbf{X} \mathbf{w} \quad (8)$$

で表すことにする。

合成変数  $f$  の平均値と分散は,

$$\begin{aligned}
 E(f) &= \frac{1}{n} \sum_{i=1}^n f_i \\
 &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p w_j x_{ij} \right) \\
 &= \sum_{j=1}^p w_j \left( \frac{1}{n} \sum_{i=1}^n x_{ij} \right) \\
 &= 0
 \end{aligned} \tag{9}$$

$$\begin{aligned}
 V(f) &= \frac{1}{n} \sum_{i=1}^n f_i^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p w_j x_{ij} \right)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p w_j x_{ij} \right) \left( \sum_{k=1}^p w_k x_{ik} \right) \\
 &= \sum_{j=1}^p \sum_{k=1}^p w_j w_k \left( \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} \right) \\
 &= \sum_{j=1}^p \sum_{k=1}^p w_j w_k r_{jk} \\
 &= \mathbf{w}' \mathbf{R} \mathbf{w}
 \end{aligned} \tag{10}$$

となる。

標準化された合成変数ベクトル  $\mathbf{g}$  は,

$$\mathbf{g} = \frac{\mathbf{f}}{\sqrt{V(f)}} = \frac{\mathbf{X}\mathbf{w}}{\sqrt{\mathbf{w}'\mathbf{R}\mathbf{w}}} \tag{11}$$

となり, 合成変数  $g$  と, もとの変数  $x_j$  との相関係数は,

$$a_j = r(g, x_j) = r(f, x_j) = \frac{1}{n} \sum_{i=1}^n g_i x_{ij} \tag{12}$$

となるので, 全ての変数との相関係数ベクトルは,

$$\begin{aligned}
 \mathbf{a} &= \frac{1}{n} \mathbf{X}' \mathbf{g} \\
 &= \frac{1}{n} \mathbf{X}' \frac{\mathbf{X}\mathbf{w}}{\sqrt{\mathbf{w}'\mathbf{R}\mathbf{w}}} \\
 &= \frac{\mathbf{R}\mathbf{w}}{\sqrt{\mathbf{w}'\mathbf{R}\mathbf{w}}}
 \end{aligned} \tag{13}$$

ここで, データ行列から標準化された合成変数を直接求めることのできる重みベクトル (標準重みベクトル) を  $\mathbf{w}_s$  とする。すなわち,

$$\mathbf{g} = \mathbf{X}\mathbf{w}_s \tag{14}$$

とすると, (11) 式から,  $\mathbf{w}$  が既知であるときは,

$$\mathbf{w}_s = \frac{\mathbf{w}}{\sqrt{\mathbf{w}'\mathbf{R}\mathbf{w}}} \tag{15}$$

となる。

また、 $\mathbf{a}$  が既知のときは、(14) 式の両辺に左から  $\frac{1}{n}\mathbf{X}'$  をかけて、

$$\begin{aligned} \mathbf{g} &= \mathbf{X}\mathbf{w}_s \\ \frac{1}{n}\mathbf{X}'\mathbf{g} &= \frac{1}{n}\mathbf{X}'\mathbf{X}\mathbf{w}_s \\ \mathbf{a} &= \mathbf{R}\mathbf{w}_s \end{aligned} \tag{16}$$

さらに、両辺に  $\mathbf{R}$  の逆行列  $\mathbf{R}^{-1}$  を左からかけることにより、

$$\mathbf{R}^{-1}\mathbf{a} = \mathbf{R}^{-1}\mathbf{R}\mathbf{w}_s = \mathbf{w}_s \tag{17}$$

が得られる。

相関係数行列  $\mathbf{R}$  の固有値を  $\lambda$ 、固有ベクトルを  $\mathbf{u}$  とすると、(18) 式に示すような性質がある。

$$\mathbf{R}\mathbf{u} = \lambda\mathbf{u}, \quad \lambda = \mathbf{u}'\mathbf{R}\mathbf{u}, \quad \mathbf{u}'\mathbf{u} = 1 \tag{18}$$

これらと前述の関連式を組み合わすと、主成分分析の場合には以下のような関連があることを示せる。

$$V(f) \equiv \lambda \quad \text{固有値} \tag{19}$$

$$\mathbf{w} \equiv \mathbf{u} \quad \text{固有ベクトル} \tag{20}$$

$$\mathbf{a} = \frac{\mathbf{R}\mathbf{w}}{\sqrt{\mathbf{w}'\mathbf{R}\mathbf{w}}} = \frac{\lambda\mathbf{w}}{\sqrt{\lambda}} = \sqrt{\lambda}\mathbf{w} \quad \text{主成分負荷量} \tag{21}$$

$$\mathbf{w}_s = \frac{\mathbf{w}}{\sqrt{\lambda}} \quad \text{主成分得点係数} \tag{22}$$