

# ケースの分類

判別分析におけるケースの分類法としては2通りある。

## 1 マハラノビス汎距離による方法

この方法では、ケースと各グループの重心までのマハラノビス汎距離を計算し、もっとも近いグループに判別する。このとき使われるものが分類関数 classification function である。

$k$  個の群の母集団の平均を  $\mu_j = (\mu_{1j}, \mu_{2j}, \dots, \mu_{pj})'$  ( $j = 1, 2, \dots, k$ ), 観測値を  $X = (X_1, X_2, \dots, X_p)'$  とする。

各群の分散共分散行列を  $\Sigma_j$ , その逆行列を  $\Sigma_j^{-1}$  とするき, (1) 式による各群までのマハラノビス距離を計算し, 最も近い群に属すると判定する。

$$d_j^2 = (X - \mu_j)' \Sigma_j^{-1} (X - \mu_j) \quad (1)$$

もし, 各群の分散共分散行列が等しい, すなわち  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_k = \Sigma$  が仮定できれば (2) 式のようになる。

$$\begin{aligned} d_j^2 &= (X - \mu_j)' \Sigma^{-1} (X - \mu_j) \\ &= X' \Sigma^{-1} X - 2 X' \Sigma^{-1} \mu_j + \mu_j' \Sigma^{-1} \mu_j \end{aligned} \quad (2)$$

第1項は各群に共通, 第3項は各群ごとに異なる定数 (これを  $c_j$  とする) である。各ケースごとに異なるのは第2項のみであるから, (3) 式の計算を行えばよい。

$$2 X' \Sigma^{-1} \mu_j = a_{1j} X_1 + a_{2j} X_2 + \dots + a_{pj} X_p \quad (3)$$

係数  $a_{1j}, a_{2j}, \dots, a_{pj}$  は  $\Sigma^{-1}$  の要素を  $\sigma^{ij}$  とすれば (4) 式で求められる。

$$a_{ij} = 2(\sigma^{i1} \mu_{1j} + \sigma^{i2} \mu_{2j} + \dots + \sigma^{ip} \mu_{pj}), \quad i = 1, 2, \dots, p \quad (4)$$

(2) 式の第1項は群に関係ないため無視できるので, (5) 式の数値が最も小さい群に属すると判定すればよい。(5) 式は, **分類関数**と呼ばれる。

$$f_j = a_{1j} X_1 + a_{2j} X_2 + \dots + a_{pj} X_p + c_j \quad (5)$$

また, マハラノビス距離の大小を比較する代わりにあらゆる2群の組合わせに対して, (6) 式で表される  $kC_2$  個の**判別関数**を定義しておくこともできる。例えば, 第1群と第2群の判別関数は,

$$\begin{aligned} Z_{1,2} &= d_1^2 - d_2^2 = f_1 - f_2 \\ &= (a_{11} - a_{12}) X_1 + (a_{21} - a_{22}) X_2 + \dots + (a_{p1} - a_{p2}) X_p + (c_1 - c_2) \end{aligned} \quad (6)$$

(6) 式の判別関数は群の数が2群のときは1個の判別関数を計算すればよいので便利であるが, 群の数が3群以上のときには, 分類関数を使う方がわかりやすい。

## 2 ベイズの定理による方法

各グループの事前確率が与えられていれば、ベイズの定理を適用することにより、ケースがどのグループに属するかという事後確率によってケースの分類を行うことができる。

グループ  $G_g$  の事前確率を  $P_g$ 、 $G_g$  の判別得点の平均ベクトルを  $\bar{\mathbf{y}}_g$  とする。また、全体での平均ベクトルを  $\bar{\mathbf{y}}$ 、プールされた判別得点の分散・共分散行列を  $\hat{\Sigma}_y$  とする。このとき、任意のケース  $s$  の判別得点が  $\mathbf{y}_s$  であるならば、元のデータが多変量正規分布に従うことと、各グループの分散・共分散行列が同等であることを仮定すれば、ケース  $s$  がグループ  $G_g$  に属する確率  $\Pr(G_g|\mathbf{y}_s)$  は、

$$\Pr(G_g|\mathbf{y}_s) = \alpha P_g \exp[-(\mathbf{y}_s - \bar{\mathbf{y}}_g)' \hat{\Sigma}_y^{-1} (\mathbf{y}_s - \bar{\mathbf{y}}_g) / 2] \quad (7)$$

となり、定数  $\alpha$  を  $\sum_{j=1}^k \Pr(G_j|\mathbf{y}_s) = 1$  となるように定義すると、 $\Pr(G_g|\mathbf{y}_s)$  は、

$$\Pr(G_g|\mathbf{y}_s) = \frac{P_g \exp[-(\mathbf{y}_s - \bar{\mathbf{y}}_g)' \hat{\Sigma}_y^{-1} (\mathbf{y}_s - \bar{\mathbf{y}}_g) / 2]}{\sum_{j=1}^k P_j \exp[-(\mathbf{y}_s - \bar{\mathbf{y}}_j)' \hat{\Sigma}_y^{-1} (\mathbf{y}_s - \bar{\mathbf{y}}_j) / 2]} \quad (8)$$

となる。この確率をもっとも高いグループに判別すればよい。

なお、ケース  $s$  がグループ  $G_g$  に属するとき、判別得点  $\mathbf{y}_s$  をとる確率  $\Pr(\mathbf{y}_s|G_g)$  は、自由度  $p$  のカイ二乗分布に従い、

$$\Pr(\mathbf{y}_s|G_g) = \Pr\{\chi^2 \leq (\mathbf{y}_s - \bar{\mathbf{y}}_g)' \hat{\Sigma}_y^{-1} (\mathbf{y}_s - \bar{\mathbf{y}}_g)\} \quad (9)$$

となる。