

# On a non-parametric discriminant procedure

阪大 葦蕨工 山口 光 代

2変数の2標本 location問題に対して, S. K. Chatterjee  
と P. K. Sen (1964) が提案したある rank-sum statistics  
を判別問題に応用して, 判別関数を作ることを試み, その場  
合の判別方程式, consistency が成立していることを示し  
た。又, 標本の大きさについての考察を加えた。

## §1. 序

次のような2変量判別問題を考える。  $\pi_i (i=1, \dots, k)$  は,  
 $k$ 個の異った母集団とする。  $(X_i, Y_i)$  は, それぞれ, 分布関  
数  $F_i(x, y)$  とする確率変数の  $k$ 個の組とする。  $F_i(x, y)$  は,  
 $\pi_i$  における分布関数を表わし, 連続関数であり得る。 各  $i$   
( $i=1, \dots, k$ ) について, この  $F_i(x, y)$  については, 何等の知  
識がなく, 母集団  $\pi_i$  から取り出されたことが, 事前に判って  
いる大きさ  $m_i$  の標本値を持つだけとする。 この標本値を  
 $(x_{i1}, y_{i1}), \dots, (x_{im_i}, y_{im_i})$  とする。 之以外, 他の母集団

$\pi_0$  における連続分布関数  $F_0(x, y)$  にとりておける確率変数  $(X_0, Y_0)$  がある。この母集団  $\pi_0$  は、 $\pi_1, \dots, \pi_k$  の  $k$  個の母集団のうちの一つと一致していることが解らる。  $\pi_0$  から、又  $2n_0$  の標本組を得るとして、その標本組を、 $(x_{01}, y_{01}), \dots, (x_{0n_0}, y_{0n_0})$  とする。これらの得られた観測値の組を比較して、 $\pi_0$  は、どの母集団と一致しているかを判定する。

このような non-parametric な判別問題は、最初、Fix & Hodges (1951, 1952) によりて試みられた。

Stoller (1954) は、1変数の場合の2母集団への判別を行っている。 Das Gupta (1964) は、 $p$ 変数の場合の  $k$  個の母集団への判別を、経験分布関数の、Kolmogorov の距離を用いて行なう。その判別方法の consistency を示している。多変数の場合、tolerance region を用いてこれを試みている Quessenberry & Gessaman (1968) の論文、そして、1変数2母集団への判別で、2種の判別誤差の最大のものである、与えた他でおおむねのような sequential な方法を取り入れた Wainsky & Kurz (1969) の論文は興味深いものである。その他にも、Ryzin (1966), Hudimeto (1963, 1964, 1968), Peltó (1969) 等、いくつかの論文がある。

## §2. 判別方式.

母集団  $\pi_i$  ( $i=1, \dots, k$ ) における分布関数  $F_i(x, y)$  は、次の条件をみたすものと可なり。

(1)  $F_1(x, y), \dots, F_k(x, y)$  は、連続関数である。

(2) ある単調関数  $g(x)$  が存在して、 $Y_i = g(X_i), \dots, Y_k = g(X_k)$  とする確率は、それぞれ  $0 < \dots < 1$  である。

(3)  $i \neq j = 1, \dots, k$  に対して、

$$P(X_i > X_j) \neq P(X_i \leq X_j)$$

$$P(Y_i > Y_j) \neq P(Y_i \leq Y_j)$$

観測値の大きさ  $n_0, n_1, \dots, n_k$  は、 $N_i = n_0 + n_i$  ( $i=1, \dots, k$ ) と可なり、 $n_0/N_i \rightarrow x_i$  ( $0 < x_i < 1$ ) ( $i=1, \dots, k$ ) なる条件をみたすように大きくなり得る。

各  $i$  ( $= 1, \dots, k$ ) に対して、観測値の組、

$$\pi_0: (x_{01}, y_{01}), \dots, (x_{0n_0}, y_{0n_0})$$

$$\pi_i: (x_{i1}, y_{i1}), \dots, (x_{in_i}, y_{in_i})$$

を比較可なり。  $(x_{01}, \dots, x_{0n_0}, x_{i1}, \dots, x_{in_i})$  を大きさ  $n$  のおりに並び、その順序は  $(Y_{n_1}, \dots, Y_{n_0+n_1}, \dots, Y_{n_i})$  と可なり。同様に  $(y_{01}, \dots, y_{0n_0}, y_{i1}, \dots, y_{in_i})$  を大きさ  $n$  のおりに並び、その順序は  $(A_{1i}, \dots, A_{n_0+n_1}, A_{n_0+n_1+1}, \dots, A_{n_i})$

$$\text{と可なり。 } \sum_{\alpha=1}^{n_0} Y_{\alpha i}, \quad \sum_{\alpha=1}^{n_0} A_{\alpha i}$$

4

$$C_i = \frac{12}{N_i(N_i^2-1)} \sum_{\alpha=1}^{N_i} \left( x_{\alpha i} - \frac{N_i+1}{2} \right) \left( \lambda_{\alpha i} - \frac{N_i+1}{2} \right)$$

とおき、

$$R_i = \frac{12}{(1-C_i^2)n_0n_i(N_i+1)} \left[ \left( S_{1i} - \frac{n_0(N_i+1)}{2} \right)^2 + \left( S_{2i} - \frac{n_0(N_i+1)}{2} \right)^2 - 2C_i \left( S_{1i} - \frac{n_0(N_i+1)}{2} \right) \left( S_{2i} - \frac{n_0(N_i+1)}{2} \right) \right]$$

とおく。この統計量に基づいて、

$$R_i = \min_{1 \leq j \leq k} R_j \quad \text{ならば} \quad \pi_0 = \pi_i$$

と判別する判別方式を考える。

§3. 判別方式の consistency

$$A(x) = \begin{cases} 1 & : x > 0 \\ 0 & : x = 0 \\ -1 & : x < 0 \end{cases} \quad \text{とする。}$$

$$M_{j1}(1,1) = \sum_{\alpha}^{n_0} \sum_{\beta}^{n_j} A(x_{\alpha\alpha} - x_{j\beta})^2 / n_0 n_j$$

$$M_{j2}(1,1) = \sum_{\alpha \neq \alpha'}^{n_0} \sum_{\beta}^{n_j} A(x_{\alpha\alpha} - x_{j\beta}) A(x_{\alpha'\alpha'} - x_{j\beta}) / n_0(n_0-1)n_j$$

$$M_{j3}(1,1) = \sum_{\alpha}^{n_0} \sum_{\beta \neq \beta'}^{n_j} A(x_{\alpha\alpha} - x_{j\beta}) A(x_{\alpha\alpha} - x_{j\beta'}) / n_0 n_j (n_j-1)$$

$$M_{j4}(1,1) = \sum_{\alpha \neq \alpha'} \sum_{\beta \neq \beta'}^{n_0} \sum_{\beta'}^{n_j} \Delta(x_{\alpha\alpha} - x_{j\beta}) \Delta(x_{\alpha\alpha'} - x_{j\beta'}) / n_0 n_j (n_0 - 1) (n_j - 1)$$

$$M_{j1}(2,2) = \sum_{\alpha} \sum_{\beta}^{n_0} \sum_{\beta'}^{n_j} \Delta(y_{\alpha\alpha} - y_{j\beta})^2 / n_0 n_j$$

$$M_{j2}(2,2) = \sum_{\alpha \neq \alpha'} \sum_{\beta}^{n_0} \sum_{\beta'}^{n_j} \Delta(y_{\alpha\alpha} - y_{j\beta}) \Delta(y_{\alpha\alpha'} - y_{j\beta'}) / n_0 (n_0 - 1) n_j$$

$$M_{j3}(2,2) = \sum_{\alpha} \sum_{\beta \neq \beta'}^{n_0} \sum_{\beta'}^{n_j} \Delta(y_{\alpha\alpha} - y_{j\beta}) \Delta(y_{\alpha\alpha} - y_{j\beta'}) / n_0 n_j (n_j - 1)$$

$$M_{j4}(2,2) = \sum_{\alpha \neq \alpha'} \sum_{\beta \neq \beta'}^{n_0} \sum_{\beta'}^{n_j} \Delta(y_{\alpha\alpha} - y_{j\beta}) \Delta(y_{\alpha\alpha'} - y_{j\beta'}) / n_0 n_j (n_0 - 1) (n_j - 1)$$

$$M_{j1}(1,2) = \sum_{\alpha} \sum_{\beta}^{n_0} \sum_{\beta'}^{n_j} \Delta(x_{\alpha\alpha} - x_{j\beta}) \Delta(y_{\alpha\alpha} - y_{j\beta'}) / n_0 n_j$$

$$M_{j2}(1,2) = \sum_{\alpha \neq \alpha'} \sum_{\beta}^{n_0} \sum_{\beta'}^{n_j} \Delta(x_{\alpha\alpha} - x_{j\beta}) \Delta(y_{\alpha\alpha'} - y_{j\beta'}) / n_0 (n_0 - 1) n_j$$

$$M_{j3}(1,2) = \sum_{\alpha} \sum_{\beta \neq \beta'}^{n_0} \sum_{\beta'}^{n_j} \Delta(x_{\alpha\alpha} - x_{j\beta}) \Delta(y_{\alpha\alpha} - y_{j\beta'}) / n_0 n_j (n_j - 1)$$

$$M_{j4}(1,2) = \sum_{\alpha \neq \alpha'} \sum_{\beta \neq \beta'}^{n_0} \sum_{\beta'}^{n_j} \Delta(x_{\alpha\alpha} - x_{j\beta}) \Delta(y_{\alpha\alpha'} - y_{j\beta'}) / n_0 n_j (n_0 - 1) (n_j - 1)$$

$\varepsilon \text{ お } < \varepsilon \quad R_j \text{ は}$

$$R_j = \frac{3}{(1 - \varepsilon_j^2) n_0 n_j (N_j + 1)} \left[ n_0 n_j (M_{j1}(1,1) + M_{j1}(2,2) - 2G_j M_{j1}(1,2)) \right]$$

$+ n_0(n_0-1)n_j (M_{j2}(1,1) + M_{j2}(2,2) - 2C_j M_{j2}(1,2))$   
 $+ n_0 n_j (n_j-1) (M_{j3}(1,1) + M_{j3}(2,2) - 2C_j M_{j3}(1,2))$   
 $+ n_0 n_j (n_0-1)(n_j-1) (M_{j4}(1,1) + M_{j4}(2,2) - 2C_j M_{j4}(1,2))$

$\varepsilon$ . U-統計量を用いて書き直せる。U-統計量の性質より、

$$A_{j1} = M_{j1}(1,1) + M_{j1}(2,2) - 2C_j M_{j1}(1,2), A_{j2} = M_{j2}(1,1) + M_{j2}(2,2) - 2C_j M_{j2}(1,2),$$

$$A_{j3} = M_{j3}(1,1) + M_{j3}(2,2) - 2C_j M_{j3}(1,2),$$

$$A_{j4} = M_{j4}(1,1) + M_{j4}(2,2) - 2C_j M_{j4}(1,2) \text{ 等. } \text{それらより}$$

$$\bar{A}_{j1} = 2(1 - \bar{F}_j) \int_{E^2} (2F_j(x, \infty) - 1)(2F_j(\infty, y) - 1) dF_0(x, y),$$

$$\bar{A}_{j2} = \int_{E^1} (2F_0(x, \infty) - 1)^2 dF_j(x, \infty) - 2\bar{F}_j \int_{E^2} (2F_0(x, \infty) - 1)(2F_0(\infty, y) - 1)$$

$$dF_j(x, y) + \int_{E^1} (2F_0(\infty, y) - 1)^2 dF_j(\infty, y),$$

$$\bar{A}_{j3} = \int_{E^1} (2F_j(x, \infty) - 1)^2 dF_0(x, \infty) + \int_{E^1} (2F_j(\infty, y) - 1)^2 dF_0(\infty, y)$$

$$- 2\bar{F}_j \int_{E^2} (2F_j(x, \infty) - 1)(2F_j(\infty, y) - 1) dF_0(x, y),$$

$$\bar{A}_{j4} = \left[ \int_{E^1} (2F_j(x, \infty) - 1) dF_0(x, \infty) - \bar{F}_j \int_{E^1} (2F_j(\infty, y) - 1) dF_0(\infty, y) \right]^2$$

$$+ \left[ \int_{E^1} (2F_j(\infty, y) - 1) dF_0(\infty, y) \right]^2 (1 - \bar{F}_j^2)$$

$\varepsilon$ .  $\Pi_0 = \Pi_1$  の下で、確率収束することがいえる。こゝで、

$C_j$  については、Chatterjee と Sen の論文 [9] で、

$$\bar{F}_j = 3 \int_{E^2} (2\bar{F}_j(x, \infty) - 1)(2\bar{F}_j(\infty, y) - 1) d\bar{F}_j(x, y)$$

ただし,  $\bar{F}_j(x, y) = X_j F_0(x, y) + (1 - X_j) F_j(x, y)$  であり.  
 とすると,  $C_j$  は  $\bar{F}_j$  に確率収束することの証明に用いられる。  
 ここで,  $\bar{F}_j = \pm 1$  と仮定するのは, 単調関数  $g(x)$  が存在して,  
 $Y_j = g(X_j)$ ,  $Y_0 = g(X_0)$  と書ける場合, かつ, この場合に限定する。  
 以上により, 条件 (2), (3) の下では,

$$R_j = \frac{3A_{j1}}{(1-C_j^2)(N_j+1)} + \frac{3(n_0-1)A_{j2}}{(1-C_j^2)(N_j+1)} + \frac{3(n_j-1)A_{j3}}{(1-C_j^2)(N_j+1)} \\ + \frac{3(n_0-1)(n_j-1)}{(1-C_j^2)(N_j+1)} A_{j4}$$

の第1項が0, 第2項が  $\frac{3\bar{A}_{j2} X_j}{(1-\bar{F}_j^2)}$ , 第3項が  $\frac{3\bar{A}_{j3} (1-X_j)}{(1-\bar{F}_j^2)}$   
 第4項は  $\frac{3\bar{A}_{j4}}{(1-\bar{F}_j^2)}$  に確率収束することから  $\frac{(n_0-1)(n_j-1)}{N_j+1} \rightarrow \infty$  になり,  
 任意の定数  $C$  に対して,  $\pi_0 = \pi_i$  の下では,

$$\text{Prob}(R_j > C) \rightarrow 1 \quad (j \rightarrow \infty)$$

と存する。

一方,  $R_i$  は,  $F_i(x, y)$  の下では, その極限分布として, 自由度2の  $\chi^2$ -分布を持つことから [9] で証明されている。

よって,  $\pi_0 = \pi_i$  の下では,

$$\text{Prob}(R_i = \min_{1 \leq j \leq k} R_j) \rightarrow 1$$

となり, この判別方式は, consistent である。

§4. 標本の大きさについて.

先に述べた判別式は、1変数の場合には、

$$R_i = \frac{12}{n_i n_j (N_i + 1)} \left( S_i - \frac{n_i (N_i + 1)}{2} \right)^2 \quad i=1, \dots, k$$

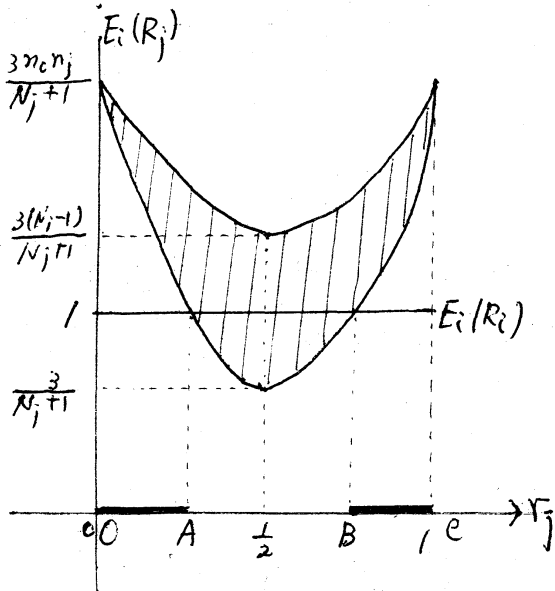
となる。\$S\_i\$ は、\$(x\_{01}, \dots, x\_{0n\_0}, x\_{i1}, \dots, x\_{in\_i})\$ の順位 \$(r\_{01}, \dots, r\_{0n\_0}, r\_{i1}, \dots, r\_{in\_i})\$ の最初の \$n\_0\$ 個の和である。

\$i \neq j\$ のとき、\$R\_j\$ は合序関数 \$F\_i(x)\$ の下で平均をとると、

$$\frac{3}{N_j + 1} [4(n_0 n_j - 1)(r_j^2 - r_j) + n_0 n_j] \leq E_i(R_j) \leq$$

$$\frac{3}{N_j + 1} [4(n_0 - 1)(n_j - 1)(r_j^2 - r_j) + n_0 n_j]$$

となる。\$\therefore R\_j = \int\_{E\_i} F\_j(x) dF\_i(x)\$。よって、\$R\_j\$ の合序関数 \$F\_i(x)\$ の下での平均区は、下図の斜線部分にある。



A, B の座標は、

$$\frac{1}{2} - \sqrt{\frac{N_j - 2}{12(n_0 n_j - 1)}}, \quad \frac{1}{2} + \sqrt{\frac{N_j - 2}{12(n_0 n_j - 1)}}$$

である。

\$E\_i(R\_i) = 1\$ であるから、

平均的にみれば、図の OA

BC の部分には \$r\_j\$ が落ち

るように、標本の大きさ

\$(n\_0, n\_j)\$ を決めなければならない



思われる。これをすべて  $j=1, \dots, k$  に行なって、標本の大きさ  $(n_0, n_1, \dots, n_k)$  を決めるのが適当であろう。non-parametric な場合、 $V_j$  は未知であるが、標本の大きさが大きくなるにつれ、OA, BC の長さが大きくなるので、十分な大きさの標本の大きさとすれば、未知の  $V_j$  が、この区間に落ちる確率も大きくなる。これは、又、§3 で証明された consistency のことも当然である。

## 参考文献

- [1]. E. Fix and J. L. Hodges (1951)  
Discriminatory Analysis, Non-Parametric Discrimination: Consistency Properties. USAF. school of aviation Medicine, Report No.4
- [2]. E. Fix and J. L. Hodges (1952)  
Discriminatory Analysis, Non-Parametric Discrimination: Small Sample Performance. USAF. school of aviation Medicine, Report No.11
- [3]. S. Das-Gupta (1964)  
Non-Parametric Classification Rules. Sankhyā A 24. p.p. 25-30
- [4]. H. Hudimoto (1963)  
分類について - I. 二群への分類 AISM. p.p. 31-38
- [5]. H. Hudimoto (1964)  
On a Distribution-free Two-Way Classification. AISM. 16 p.p. 247-253
- [6]. H. Hudimoto (1968)  
On the Empirical Bayes Procedure (1). AISM. 20 p.p. 169-185
- [7]. C. R. Pelto (1969)  
Adaptive Non-Parametric Classification. Technometrics 11.
- [8]. C. P. Queenberry and M. P. Gassaman (1968) Non-Parametric  
Discrimination Using Tolerance regions. AMS. 39 p.p. 664-673
- [9]. J. Van Ryzin (1966)  
Bayes Risk Consistency of Classification Procedures Using Density

Estimation. Sankhyā A 26. p.p. 261-270

[10]. S. K. Chatterjee and P. K. Sen (1964)

Non-Parametric Tests for the Bivariate Two-Sample Location Problems.

CSAB 14 p.p. 18-58

[11]. D. S. Stoller (1954)

Univariate Two-Population Distribution-Free Discrimination.

TASA 49 p.p. 770-777

[12]. M. N. Woinsky and L. Kurz (1969)

Sequential Non-Parametric Two-Way Classification with a Prescribed

Maximum Asymptotic Error Probability. AMS. 40 445-455.