

多変量解析における変数の 選択について

岡山大学 大崎 純一

1. 緒言 種々の状況のもとで、多次元のデータが測定される。このデータをもとに、次におけると予想される状況も予測するため、あるいは他の状況との比較をするために、多変量解析に含まれる各種の方法が用いられる。ところで、多次元のデータのなかにも、目的に対してきわめて有効なものもあるし、反対にほとんど有効でないものもある。

そこで、多変量解析のうちで、重回帰分析、2つの多次元正規分布の平均値の比較、そして線形判別関数を利用する際の変数の重要性について統一的に考えることとする。

2 変数の効果

1) 重回帰分析 回帰方程式における独立変数の効果は、回帰直線からの偏差平方和に含まれている。

$$S_y^2 = \left\{ \sum_{i=1}^n (y_i - \bar{y})^2 - d^T S^{-1} d \right\} / (n - k - 1) \text{ ----- (1)}$$

d : $k \times 1$ vector, S : $k \times k$ matrix.

ロ) 2つの多次元正規分布の平均値間の距離 等分散であるという仮定の下で, 変数の効果は, 次に示す距離に含まれる。

$$D = d' S^{-1} d \quad \text{----- (2)}$$

d : $k \times 1$ 平均値差ベクトル, S : $k \times k$ 分散共分散行列

ここで, D は, マハラノビスの距離である。

ハ) 線形判別関数 2群を判別するための判別関数における変数の効果は, 誤まって判別する確率に含まれる。

$$P = \Pr \{ Z > \sqrt{D/2} \} \quad \text{----- (3)}$$

$F = F(1, P) \{ Z > \alpha \}$ は, $N(0, 1)$ による確率。

(1), (2), (3) いずれにおいても, 各個の変数の効果は, 全体として, $d' S^{-1} d$ に含まれている。そこで, 各変数の効果あるいは変数の組の効果も $d' S^{-1} d$ から分離し, これらの方法を使用する際に有効な変数も選ぶ方法を次に述べる。

3 変数の選択方法

$d = (d_1, d_2, \dots, d_k)'$, $S = (s_{ij})$ とおけば,

$$(1) \text{ では, } \begin{cases} d_i = \sum_{j=1}^n (x_{ij} - \bar{x}_i)(y_j - \bar{y}) \\ s_{ij} = \sum_{k=1}^n (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j) \end{cases}$$

そして, (2), (3) では,

$$d_i = \bar{x}_{1i} - \bar{x}_{2i}$$

$$s_{ij} = \frac{\sum_{p=1}^n \sum_{l=1}^k (x_{pil} - \bar{x}_{pi})(x_{plj} - \bar{x}_{pj})}{(n_1 + n_2 - 2)}$$

である。さらに、 $r_{ij} = s_{ij} / \sqrt{s_{ii} s_{jj}}$, $\delta_i = d_i / \sqrt{s_{ii}}$

とおき、 $R = (r_{ij})$, $\delta = (\delta_1, \delta_2, \dots, \delta_k)'$ とおけば、

$$d'S^T d = \delta' R' \delta = \delta_i^2 + (\delta_{i_2} - r_{i_1 i_2} \delta_{i_1})^2 / (1 - r_{i_1 i_2}^2) + \dots$$

$$+ a_{ijij}^j (\delta_{ij} + \sum_{l=1}^{j-1} \delta_{il} a_{ijil}^j / a_{ijij}^j)^2 + \dots \quad \dots (4)$$

$$+ a_{i_k i_k}^k (\delta_{i_k} + \sum_{l=1}^{k-1} \delta_{il} a_{i_k i_l}^k / a_{i_k i_k}^k)^2$$

但し、 a_{ijil}^j は、変数 $x_{i_1}, x_{i_2}, \dots, x_{ij}$ の相関行列の逆行列の j 行 l 列の要素

(4)の性質を用いて、変数の選択方法を定式化する。重回帰分析において、N. Draper¹⁾ は、変数の選択方法の名前として、Forward Selection Procedure, Backward Elimination Procedure, をして、All Possible Selection Procedure等を使用しているので、ここでもそれらの名前を用いることとする。

解析に用いられる m 個の変数の集合を

$$C = \{x_1, x_2, \dots, x_m\}$$

とする。また、 $D(i_1, i_2, \dots, i_m)$ を m 個の変数 $x_{i_1}, x_{i_2}, \dots, x_{i_m}$ から計算した $d'S^T d$ の値とする。さらに $E(m, i_l)$ を m 個の変数 $x_{i_1}, x_{i_2}, \dots, x_{i_m}$ から変数 x_{i_l} を除き、残り $m-1$ 個の変数から計算した $d'S^T d$ の値とする。

1) Forward Selection Procedure (FSP)

$$\max_{i_m} D(I_1, I_2, \dots, I_{m-1}, i_m), x_{i_m} \in C - \{x_{I_1}, x_{I_2}, \dots, x_{I_{m-1}}\}$$

を満す変数を x_{I_m} とする。

$$m = 1, 2, \dots, k$$

2) Backward Elimination Procedure (BEP)

$$\max_{i_k} E(m, i_k), x_{i_k} \in C - \{x_{J_k}, x_{J_{k-1}}, \dots, x_{J_{k-m+1}}\}$$

を満す変数を x_{J_m} とする。

$$m = k, k-1, \dots, 2, 1$$

3) All Possible Selection Procedure (APSP)

$$\max_{(i_1, i_2, \dots, i_m)} D(i_1, i_2, \dots, i_m)$$

$$\{x_{i_1}, x_{i_2}, \dots, x_{i_m}\} \subset C$$

$$i_j \neq i_l$$

を満す変数の組を $\{x_{K_1^m}, x_{K_2^m}, \dots, x_{K_m^m}\}$ とする。

$$m = 1, 2, \dots, k$$

FSP と APSP では、効率の高い変数から順に選ばれるが、BEPでは、効率の低い変数から順に選ばれることになる。それゆえ、BEPの結果を逆にして $x_{J_1}, x_{J_2}, \dots, x_{J_m}, \dots, x_{J_k}$ として用いれば、FSPあるいはAPSPの結果と比較するのに便利である。これらの選択方法の間には、次の関係が成立す

3.

定理 1

All Possible Selection Procedure をおいて

$$\{x_{k_1^i}, x_{k_2^i}, \dots, x_{k_i^i}\} \subset \{x_{k_1^{i+1}}, x_{k_2^{i+1}}, \dots, x_{k_{i+1}^{i+1}}\} \dots (5)$$

$$i = 1, 2, \dots, k-1$$

が成立するならば, Forward Selection Procedure と

Backward Elimination Procedure の結果は一致する。

すなわち, $x_{I_i} = x_{J_i}$, $i = 1, 2, \dots, k$

証明

1) APSP と FSP との一致

 $i=1$ では, APSP と FSP の性質から $x_{I_1} = x_{k_1}$ 。 $i=2$ では, (5) と $i=1$ の結果より,

$$x_{I_1} = x_{k_1^1} = x_{k_1^2}$$

とおいても一般性を失わない。そして,

APSP には,

$$\max_{(i_1, i_2)} D(i_1, i_2) = \max_{(I_1, i_2)} D(I_1, i_2)$$

となる。よって,

$$x_{I_2} = x_{k_2^2}$$

帰納法によつて, $\{x_{I_1}, x_{I_2}, \dots, x_{I_m}\} = \{x_{k_1^m}, x_{k_2^m}, \dots, x_{k_m^m}\}$ となる。そして, $x_{k_j^{m+1}} = x_{k_j^m} = x_{j_j}$, $j = 1, 2, \dots, m-1$ とおけば, $x_{I_m} = x_{k_m^m}$ となる。

↳

2) APSP と BEP の一致と FSP との一致

$i = k-1$ では, APSP と BEP の性質から

$$\max_{(i_1, \dots, i_{k-1})} D(i_1, i_2, \dots, i_{k-1}) = \max_{i_k} E(k, i_k)$$

$$\{x_{i_1}, x_{i_2}, \dots, x_{i_{k-1}}\} \subset C, \{x_{i_k}\} \subset C$$

であり, $\{x_{k_1^{k-1}}, x_{k_2^{k-1}}, \dots, x_{k_{k-1}^{k-1}}\} = C - \{x_{J_k}\}$
となる。

APSP と FSP との一致から

$$x_{I_k} = x_{J_k}.$$

$$\text{以下帰納法によつて, } \{x_{k_1^{i+1}}, x_{k_2^{i+1}}, \dots, x_{k_{i+1}^{i+1}}\}$$

$$= C - \{x_{J_k}, \dots, x_{J_{k-i-1}}\}$$

の関係から

$$x_{I_i} = x_{J_i}$$

が成立する。

定理 2

Forward Selection Procedure を用いて,

$$D(I_1, I_2, \dots, I_m) \geq D(i_1, i_2, \dots, i_{k-1}) \quad \dots (6)$$

$$\{x_{i_1}, x_{i_2}, \dots, x_{i_{k-1}}\} = C - \{x_{I_l}\}$$

$$l = 1, 2, \dots, m-1$$

$$\text{ならば, } \{x_{I_1}, x_{I_2}, \dots, x_{I_m}\} = \{x_{k_1^m}, x_{k_2^m}, \dots, x_{k_m^m}\}$$

$$i = 2, 3, \dots, k-1$$

が成立する。

証明

$$C - \{x_{I_l}\} = G_l, \quad l=1, 2, \dots, m-1 \text{ とおく。}$$

$$\text{また, } G_1 \wedge G_2 \wedge \dots \wedge G_{m-1} = G_0 \text{ とおく。}$$

G_0 に含まれる変数の数は, $n-i+1$ 個となる。

$$C - G_0 = \{x_{I_1}, x_{I_2}, \dots, x_{I_{m-1}}\} \text{ ----- (7)}$$

となる。

$k C_m$ は, G_0 と $C - G_0$ とに含まれる変数を区別する

と,

$$k C_m = \sum_{j=0}^{m-1} k-m+1 C_{m-j} x^{m-1} C_j \text{ ----- (8)}$$

となる。

$k-m+1 C_{m-j} x^{m-1} C_j$ は, G_0 から $m-j$ 個, $C - G_0$ から j 個の変数を選び合わせて m 個の変数の組を作る場合の組み合わせである。

ところで, (7) より, $C - G_0$ の任意の部分集合は, G_l の少なくとも一つに含まれる。以上は,

$j=1, 2, \dots, m-2$ において成立する。

$k-m+1 C_1 C_{m-1}$ は, FSP における組み合わせである。

$$\text{また, (4) より, } D(i_1, i_2, \dots, i_p) \geq D(j_1, j_2, \dots, j_p) \\ \{x_{j_1}, x_{j_2}, \dots, x_{j_p}\} \subset \{x_{i_1}, x_{i_2}, \dots, x_{i_p}\} \\ j \leq p$$

が成立する。

よって, (6), (8), (9)より, FSP と APSP とは一致する。

定理 3

Backward Elimination Procedure について,

$$E(k-m, J_{k-m}) \geq D(J_k, J_{k-1}, \dots, J_{k-m}, j_1, j_2, \dots, j_{k-m-2})$$

$$\{x_{j_1}, x_{j_2}, \dots, x_{j_{k-m-2}}\} \subset C = \{x_{J_k}, x_{J_{k-1}}, \dots, x_{J_{k-m}}\}$$

$$\text{ならば, } C = \{x_{J_k}, \dots, x_{J_{k-m}}\} = \{x_{K_1^{k-m}}, x_{K_2^{k-m}}, \dots, x_{K_{k-m}^{k-m}}\}$$

$$m = 0, 1, 2, \dots, k-2$$

証明は, 定理 2 とほぼ同様に行なうことが出来る。

4. APSP の計算アルゴリズム

APSP において, すべての組み合わせについて計算を行なうと, $2^k - 1$ 回行なわなければならない。それゆえ, 工藤, 垂水^{21,22)}等は, 計算時間も短くするのための APSP の計算アルゴリズムについて述べている。

ここでは, 中間結果を蓄える領域を 2箇所にし, 組み合わせを作る回数を少なくした計算アルゴリズムについて述べる。

k 個の変数のうちで, m 個の変数を用いる場合に, 残りの $k-m$ 個の変数についても同時に考えることにする。 m 個の

変数を $\{x_{i_1}, x_{i_2}, \dots, x_{i_m}\}$ とおくと, 残りの $k-m$ 個の変数は

$$\{x_{j_1}, x_{j_2}, \dots, x_{j_{k-m}}\} = C - \{x_{i_1}, x_{i_2}, \dots, x_{i_m}\} \quad (6)$$

となる。(6)において, 右辺の m 個の変数を k 個のうちから任意に選ぶ選ぶ方は, ${}^k C_m$ 通りある。右辺の m 個の変数に対応する δ, R の要素を一箇所に, 残りの $k-m$ 個のものも他の箇所に集める。そして

$$\max_{(i_1, i_2, \dots, i_m)} D(i_1, i_2, \dots, i_m), \quad \{x_{i_1}, x_{i_2}, \dots, x_{i_m}\} \subset C$$

$$\max_{(j_1, j_2, \dots, j_{k-m})} D(j_1, j_2, \dots, j_{k-m}), \quad \{x_{j_1}, x_{j_2}, \dots, x_{j_{k-m}}\} = C - \{x_{i_1}, x_{i_2}, \dots, x_{i_m}\}$$

を考へる。二項定理より, ${}^k C_m = {}^k C_{k-m}$ となることから,

${}^k C_m$ 個の変数の組み合わせから

$$\{x_{k_1^m}, x_{k_2^m}, \dots, x_{k_m^m}\}, \{x_{k_1^{k-m}}, x_{k_2^{k-m}}, \dots, x_{k_{k-m}^{k-m}}\}$$

$$m = 1, 2, \dots, k'$$

を同時に求めることができる。但し, k が奇数の場合は,

$k' = (k-1)/2$, 偶数の場合は, $k' = k/2$ である。

参考文献

- 1). Draper, N. and Smith, H. (1966), Applied Regression Analysis, John Wiley Sons, Inc.
- 2) Kudō, A. (1963). Mahalanobis measure as a criterion for the selection of variables. Memoirs of the Faculty of Science, Kyushu University, Series A.

17, 63-75

- 3) Tarumi, T and Kudô, A. (1974). An algorithm related to all possible regression and discriminant analysis, Journ. Japan Statist. Soc. 4.2, 47-56