

複数個の母集団に関する正規型多次元多重回帰分析 の手法について

伊藤, 孝一
南山大学

大崎, 紘一
岡山大学工学部

工藤, 昭夫
九州大学理学部

<https://doi.org/10.15017/1468004>

出版情報 : 九州大学大型計算機センター広報. 5 (3), pp.11-18, 1972-06-27. 九州大学大型計算機センター

バージョン :

権利関係 :

複数の母集団に関する正規型多次元多重回帰分析の手法について

* ** ***
伊藤孝一・大崎紘一・工藤昭夫

第1節 最小自乗法及びその拡張

データを直線であてはめたり、2次曲線であてはめたりする事は、最小自乗法の特殊例としてよく用いられるところである。念の為に簡単にその概略を述べてみると次のようになる。

モデルを

$$Y = a_1 X_1 + \cdots + a_k X_k = (\underline{a}' \underline{X}) \quad (1)$$

とする。

ここで、1次式 $y = \alpha + \beta x$ の場合には $a_1 = \alpha$ $a_2 = \beta$ $X_1 = 1$ $X_2 = x$ となり

2次式 $y = \alpha + \beta x + \gamma x^2$ の場合には $a_1 = \alpha$ $a_2 = \beta$ $a_3 = \gamma$ $X_1 = 1$ $X_2 = x$ $X_3 = x^2$

となる。 a_i ($i = 1, \cdots, k$) は回帰係数とよばれる。

$$\text{データとして} \begin{bmatrix} x_1^1 & \cdots & x_k^1 & y_1 \\ \cdot & & & \cdot \\ \cdot & & & \cdot \\ \cdot & & & \cdot \\ x_1^n & & x_k^n & y_n \end{bmatrix} = [\underline{X} \ \underline{y}]$$

即ち、各々が $k+1$ 個の数値よりなる大きさ n のデータが与えられた時の最小自乗法による解は、

$$\hat{\underline{a}} = (X'X)^{-1} X' \underline{y}$$

であり、 $\hat{\underline{a}}$ は、正規方程式 $X' \underline{y} = (X'X)^{-1} \underline{a}'$ の解である。

更に誤差項を考慮して、モデルを改め

$$Y = (\underline{a}' \underline{X}) + \varepsilon \quad (1)'$$

但し、 ε は平均0で、同一分散 σ^2 をもつ誤差変量と考える。 $E(\varepsilon) = 0, V(\varepsilon) = \sigma^2$ 。

$\hat{\underline{a}}$ は、不偏推定であり、推定量の分散、共分散行列は $V(\hat{\underline{a}}) = \sigma^2 (X'X)^{-1}$ の形をした $k \times k$ 次非負な対称行列で与えられ、更に σ^2 の不偏推定は、

$$\hat{\sigma}^2 = (Y - X\hat{\underline{a}})'(Y - X\hat{\underline{a}}) / (n - k) = (Y'Y - Y'X\hat{\underline{a}}) / (n - k)$$

により与えられる。

この議論を多次元のデータに拡張する事、即ち モデルを

$$\underline{Y} = \underline{a}_1 X_1 + \cdots + \underline{a}_k X_k \quad (2)$$

但し、 \underline{Y} , $\underline{a}_1, \cdots, \underline{a}_k$ は各々 p 次元のベクトル

この拡張も、初等的な統計の教科書に載っていないが、理論的には、完成している。

* 南山大学

** 岡山大学工学部

*** 九州大学理学部

多次元のモデルの理論は、とりあえず問題にしない事にして、再び1次式の場合に戻り議論する事にしよう。 実例としては、放射線の照射量と白血球の減少量、化学工業の生産工程における温度と成分量、不快指数と交通事故の数等、いろいろ考える事ができるが、私共のような、具体的な問題に対する思考が音痴に等しい者は、記号で書いてある方がかえってわかりやすいので、記号のまま説明させていただく。

3組のデータがあり、各々直線で当てはめる事ができるとしよう。

モデルとして

$$y = \alpha_1 + \beta_1 x + \varepsilon_1$$

$$y = \alpha_2 + \beta_2 x + \varepsilon_2$$

$$y = \alpha_3 + \beta_3 x + \varepsilon_3$$

但し、 α_i 、 β_i は回帰係数 ε_i は誤差項とする。

実際に3組のデータに直面し、散布図を書いて検討する場合、次のようないろいろの疑問が起こってくる。

- (1°) 誤差項の分散は共通でなかろうか。
- (2°) 誤差項の分散が共通でなくとも、そのうちの2個、たとえば ε_1 と ε_2 の分散はほぼ同一でなかろうか。
- (3°) 回帰係数が全部相異なる、即ち3本の別々の直線で当てはめなくてはならない。
- (4°) 3本の直線は異なるにしても、互に平行である。 即ち $\beta_1 = \beta_2 = \beta_3$
- (5°) 平行でないかわりに、特定の x の値 x_0 で交わる。 即ち、モデルを

$$y = \alpha'_1 + \beta'_1 (x - x_0)$$

$$y = \alpha'_2 + \beta'_2 (x - x_0)$$

$$y = \alpha'_3 + \beta'_3 (x - x_0)$$

と書き変えると $\alpha'_1 = \alpha'_2 = \alpha'_3$ を満足する。

- (6°) 3個のデータは同一直線で書くことができる。 即ち $\alpha_1 = \alpha_2 = \alpha_3$ $\beta_1 = \beta_2 = \beta_3$
- (7°) x の値により y は変わらない。 即ち $\beta_1 = \beta_2 = \beta_3 = 0$

このような、複数個のデータの組がある場合には、簡単な直線の当てはめの問題ですら、多種多様な解析が必要となり、(1)のような一般のモデルの場合には、なおのこと非常に数多くの計算が必要になる。上述の回帰係数の推定式や分散の推定式をみると、逆行列を作る操作、行列とベクトルあるいは、行列と行列の積をとる操作のみで計算可能であるので、計算のアルゴリズムには、ある種の共通性がある。即ち、汎用アルゴリズムが可能ではなかろうかとの疑問に答えたものが〔1〕である。

このたび、この論文と同一名のライブラリープログラムを開発する事ができて、4月のプログラムライブラリー専門委員会の決定により九大センターに登録する運びとなったが同一アルゴリズムによるプログラムで、我々が経験した所では、24組のデータについて16種類、全部で約250回の計算を半年かけて計算した結果を、再度このアルゴリズムを用いて計算した所、数分で全部の計算が完了した。〔2〕統計解析の実地では多くの場合、1次元の y がデータとしてとられるのではなく、多次元のデータ

— \underline{Y} がとられる場合が多く、複数個の母集団の最小自乗法はきわめて重要である。この背後にあるアルゴリズムを統一モデルについて述べる前に、多次元の回帰理論を概観しておく。

モデル (1) を拡張したモデル (2) に、誤差項を考慮して

$$\underline{Y} = a_1 x_1 + a_2 x_2 + \dots + a_k x_k + \varepsilon = (a_1 \ a_2 \ \dots \ a_k) \underline{X} + \varepsilon = A \underline{X} + \varepsilon \quad (2)'$$

但し $\underline{X} = \begin{bmatrix} x_1 \\ \cdot \\ \cdot \\ x_k \end{bmatrix}$ k次元ベクトル

$A = (a_1, a_2, \dots, a_k)$ 大きさ $p \times k$ の行列

ε : p 次元の誤差変量で平均値ベクトル Q , 分散共分散行列 Λ をもつとする。

大きさ n のデータを考え、このモデルに従うデータ

$$\begin{bmatrix} x_1^1 & \dots & x_k^1 & y_1^1 & \dots & y_p^1 \\ x_1^2 & \dots & x_k^2 & y_1^2 & \dots & y_p^2 \\ \cdot & & \cdot & \cdot & & \cdot \\ \cdot & & \cdot & \cdot & & \cdot \\ x_1^n & \dots & x_k^n & y_1^n & \dots & y_p^n \end{bmatrix} = [XY]$$

その時 A 及び Λ の推定量の $\hat{A}, \hat{\Lambda}$ は

$$(1^\circ) \quad \hat{A}' = (X'X)^{-1} X'Y$$

$$(2^\circ) \quad \hat{\Lambda} = \{Y'Y - Y'X(X'X)^{-1} X'Y\} / (n-k)$$

で与えられる。

第2節 統一モデルについて

複数個の母集団からのデータが得られる場合の統一モデルについて、1次元の回帰モデルに限って述べよう。 ℓ 個の母集団があり、各々から大きさ n_1, n_2, \dots, n_ℓ の標本が得られる場合を考えよう。 この場合には ℓ 個の (1) の形をしたモデルが想定される。この場合、回帰係数、 a_1, a_2, \dots, a_k のうちいくつかは、母集団に特有なもので、母集団毎に異なり、残りのいくつかは母集団について共通であると考えることができる。

今、 ℓ 個のモデルを

$$Y = a_1^{(1)} X_1 + \dots + a_k^{(1)} X_k$$

$$Y = a_1^{(2)} X_1 + \dots + a_k^{(2)} X_k$$

$$\cdot \qquad \qquad \qquad \cdot$$

$$Y = a_1^{(\ell)} X_1 + \dots + a_k^{(\ell)} X_k$$

と書き、

ℓ 組のデータを次のように書く。

$$\begin{aligned}
 & \begin{bmatrix} x_1^1 & \cdots & x_m^1 & x_{m+1}^1 & \cdots & x_k^1 & y^1 \\ \cdot & & \cdot & \cdot & & \cdot & \cdot \\ \cdot & & \cdot & \cdot & & \cdot & \cdot \\ x_1^{n_1} & \cdots & x_m^{n_1} & x_{m+1}^{n_1} & \cdots & x_k^{n_1} & y^{n_1} \end{bmatrix} \equiv [X(1,1)X(1,2)y(1)] \\
 & \begin{bmatrix} x_1^1 & \cdots & x_m^1 & x_{m+1}^1 & \cdots & x_k^1 & y^1 \\ \cdot & & \cdot & \cdot & & \cdot & \cdot \\ \cdot & & \cdot & \cdot & & \cdot & \cdot \\ x_1^{n_2} & \cdots & x_m^{n_2} & x_{m+1}^{n_2} & \cdots & x_k^{n_2} & y^{n_2} \end{bmatrix} \equiv [X(2,1)X(2,2)y(2)] \\
 & \vdots \\
 & \begin{bmatrix} x_1^1 & \cdots & x_m^1 & x_{m+1}^1 & \cdots & x_k^1 & y^1 \\ \cdot & & \cdot & \cdot & & \cdot & \cdot \\ \cdot & & \cdot & \cdot & & \cdot & \cdot \\ x_1^{n_\ell} & \cdots & x_m^{n_\ell} & x_{m+1}^{n_\ell} & \cdots & x_k^{n_\ell} & y^{n_\ell} \end{bmatrix} \equiv [X(\ell,1)X(\ell,2)y(\ell)]
 \end{aligned}$$

はじめの m 個の回帰係数は、母集団に個有なものと考え、残りは共通なものと考え、共通な $a_j^{(i)}$ の母集団を表わす添字 i をはずして

$$\begin{aligned}
 Y &= a_1^{(1)}X_1^{(1)} + \cdots + a_m^{(1)}X_m^{(1)} + a_{m+1}X_{m+1} + \cdots + a_kX_k \\
 Y &= a_1^{(2)}X_1^{(2)} + \cdots + a_m^{(2)}X_m^{(2)} + a_{m+1}X_{m+1} + \cdots + a_kX_k \\
 &\vdots \\
 Y &= a_1^{(\ell)}X_1^{(\ell)} + \cdots + a_m^{(\ell)}X_m^{(\ell)} + a_{m+1}X_{m+1} + \cdots + a_kX_k
 \end{aligned}$$

と書き直す事ができる。

この ℓ 個のモデルをひとつのモデルで書き表わすとすれば

$$\begin{aligned}
 Y &= (a_1^{(1)}X_1^{(1)} + \cdots + a_m^{(1)}X_m^{(1)}) + (a_1^{(2)}X_1^{(2)} + \cdots + a_m^{(2)}X_m^{(2)}) + \cdots + (a_1^{(\ell)}X_1^{(\ell)} + \cdots + a_m^{(\ell)}X_m^{(\ell)}) \\
 &+ a_{m+1}X_{m+1} + a_{m+2}X_{m+2} + \cdots + a_kX_k \quad (3)
 \end{aligned}$$

と書ける。ここで $x_j^{(i)} (i=1, \dots, \ell)$ は観測値が第 i 母集団から来た時には $X_j^{(i)} = X_j$ 他は全部 0 とおく事にする。モデル (3) で $(n_1 + n_2 + \cdots + n_\ell)$ 個のデータから成る 1 組のデータを考えると

$$\begin{bmatrix} x_1^1 & \vdots & x_m^1 & \vdots & x_{m+1}^1 & \cdots & x_k^1 & y^1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_1^{n_1} & \vdots & x_m^{n_1} & \vdots & x_{m+1}^{n_1} & \cdots & x_k^{n_1} & y^{n_1} \\ & & x_1^1 & \cdots & x_m^1 & & x_{m+1}^1 & \cdots & x_k^1 & y^1 \\ & & \vdots & & \vdots & & \vdots & & \vdots & \vdots \\ & & x_1^{n_2} & \cdots & x_m^{n_2} & & x_{m+1}^{n_2} & \cdots & x_k^{n_2} & y^{n_2} \\ & & & & & & \vdots & & \vdots & \vdots \\ & & & & & & \vdots & & \vdots & \vdots \\ & & & & & & x_1^1 & \cdots & x_m^1 & x_{m+1}^1 & \cdots & x_k^1 & y^1 \\ & & & & & & \vdots & & \vdots & \vdots & & \vdots & \vdots \\ & & & & & & x_1^{n_\ell} & \cdots & x_m^{n_\ell} & x_{m+1}^{n_\ell} & \cdots & x_k^{n_\ell} & y^{n_\ell} \end{bmatrix}$$

と書ける。上の表から明らかなように、統一モデル(3)に必要な行列やベクトルは、 ℓ 組のデータの各々から作られた行列やベクトルで書く事ができる。

$$X'X = \begin{bmatrix} X'(1,1)X(1,1) & & & X'(1,1)X(1,2) \\ & X'(2,1)X(2,1) & & \\ & & & \\ & & X'(\ell,1)X(\ell,1) & X'(\ell,1)X(\ell,2) \\ \hline X'(1,2)X(1,1) & & X'(\ell,2)X(\ell,1) & \sum_{i=1}^{\ell} X'(i,2)X(i,2) \end{bmatrix}$$

: 大きさ $m \times \ell + (n-k)$ の正方行列

$$X'Y = \begin{bmatrix} X'(1,1)Y^1 \\ X'(2,1)Y^2 \\ \vdots \\ X'(\ell,1)Y^\ell \\ \sum_{i=1}^{\ell} X'(i,2)Y^i \end{bmatrix}$$

: 大きさ $m \times \ell + (m-k)$ のベクトル

$$Y'Y = \sum_{i=1}^{\ell} Y^i Y^i \quad \text{: スカラー}$$

となり最初に述べた公式により未知母数の推定が可能である。多次元のモデルの場合には、 $X'Y$ はベクトルではなく、 $\{m \times \ell + (m-k)\} \times P$ 次の行列、 $Y'Y$ は $P \times P$ の行列となる。

モデル(2)に戻って考えると、推定量の計算のアルゴリズムは、 $(m+p)$ 次の正方行列

$$\begin{bmatrix} X'X & X'Y \\ Y'X & Y'Y \end{bmatrix}$$

の左上の m 次の正方行列を“sweep out”する事によって求められる。

ℓ 組のデータがある場合には、このような行列が ℓ 枚できる。これを我々は、データ行列“data matrix”とよんでいる。更に統一モデルを媒介として、 ℓ 組のデータを1組のデータと考え直した場合には、sweep out すべき行列は、 ℓ 枚の行列をいわば適当に重ね合わせる事によって作ることができる。これは我々は“working matrix.”とよんでいる。

我々が書いたプログラムは、出来るだけ多様性がある解析を一回の計算で遂行出来るように考え作成したものである。たとえば data matrix を記憶したあとで、working matrix を作る時に、data matrix の全体を用いて作る必要はない。即ちおのおの data matrix のなかで、利用しない部分があってもかまはない。このことは、第1節の(5')で記したような疑問があらかじめわかっているときには、 x の外に $(x-x_0)$ もデータの一つのコンポーネントとして取扱い、data matrix を作っておけば、同じ直線をつぎの二つの式

$$y = \alpha + \beta x$$

$$y = \alpha' + \beta(x - x_0)$$

を用いてあてはめた場合に (α, β) を計算することと全く同様な手順で (α', β) を計算することが出来る。このような時には、data matrix 自身は singular になり、逆行列を計算するときには overflow

するので、その点だけ注意すればよい。また、記憶した data matrix のうち、たとえば、1 番目と 3 番目の data matrix だけをえらび、working matrix をつくることもできる。たとえば、第 1 節で述べた 3 本の直線のあてはめの例では、第 1 の母集団と第 3 の母集団とのあいだには、平行線でのあてはめができるとの関係を仮定して解析を行なう事ができる。

このように、統計解析の問題に応じて、data matrix から working matrix をつくる事が、出来るだけ自由に行なえるよう工夫したつもりであるが、吾々 3 人にはよくわかって、利用者の側からは、理解しにくい面があり、充分の能力を発揮しない事になりはしないかとも心配しており、利用者側からの素直な批判を、とくをお願いする次第である。

第 3 節 検定について

このプログラムでは、次の 3 種類の検定ができるように工夫されている。

- (1°) 誤差項の分散 σ^2 、あるいは分散、共分散行列 Λ が、母集団に関して一様であるか否か
- (2°) 特定の回帰係数 a_i (あるいは回帰ベクトル a_j) が 0 であるか否か
- (3°) 特定の回帰係数 (あるいは回帰ベクトル)、たとえば $a_{m+1}, a_{m+2}, \dots, a_k$ が母集団に関して一様であるか否か、即ち第 2 節で説明した統一モデルが妥当なものであるか否か

検定手法としては、問題 2 については t 検定 (多次元の場合、Hotelling の T^2 統計量) を計算するようになっている。(1), (3) の場合は、仮説を仮定した場合の誤差分散の推定 $\hat{\sigma}^2$ (多次元の場合は $\hat{\Lambda}$) と、仮定しない場合の誤差分散の推定値を比較する事により、検定が可能であり、T. W. Anderson の著書 [3] などを参照しながら検定の為の計算に必要な値が出力されるように工夫してある。多次元の場合に、このプログラム特有な事として次のような工夫がなされている。

たとえば、(2) の問題で $i=1$ としたとき

$$\text{仮説 } H : \underline{a}_1 = 0 \quad (\underline{a}_1' = (a_{11} \ a_{12} \ \dots \ a_{1p}))$$

の検定統計量を計算する際の sweep out の各段階に於て副産物としてでてくる p 個の仮説

$$H_1 : a_{11} = 0$$

$$H_2 : a_{11} = a_{12} = 0$$

⋮

$$H_j : a_{11} = a_{12} = \dots = a_{1j} = 0$$

⋮

$$H = H_p : a_{11} = a_{12} = \dots = a_{1i} = \dots = a_{1p} = 0$$

を検定する統計量を順次に出力する事が工夫されており、多変量解析を実地に使う場合には、便利であらうと期待している。

この計算のアルゴリズムについて説明すると次の通りである。Hotelling の T^2 統計量は、通常次のような形をする統計量として書かれている。 θ を p 次元のベクトル、 $\hat{\Sigma}$ を p 次の正定値対称行列とすると T^2 は $\theta' \hat{\Sigma}^{-1} \theta$ に比例する量であり、これは $p+1$ 次の正方行列 $\begin{bmatrix} \hat{\Sigma} & \theta \\ \theta' & 0 \end{bmatrix}$ の左上の $\hat{\Sigma}$

と書かれている部分を sweep out した後 (p+1, p+1) 要素の符号を変えたものである。

この事の解説は、専門書にゆずる。(〔4〕Chapter を参照)

今一般に、 $q_1 + q_2 + q_3$ 次元の対称正方行列があり、次元の分割に対応させ、行列を分割して

$$\begin{bmatrix} A & B & C \\ B' & D & E \\ C' & E' & F \end{bmatrix}$$

と書くとする。今はじめて q_1 列を sweep out したとする。これは次の形の行列演算を行なった事に他ならない。

$$\begin{bmatrix} A^{-1} & 0 & 0 \\ -B'A^{-1} & I & 0 \\ -D'A^{-1} & 0 & I \end{bmatrix} \begin{bmatrix} A & B & C \\ B' & D & E \\ C' & E' & F \end{bmatrix} = \begin{bmatrix} I & A^{-1}B & A^{-1}C \\ 0 & C-B'A^{-1}B & E'-B'A^{-1}C \\ 0 & E'-D'A^{-1}B & F-C'A^{-1}C \end{bmatrix}$$

はじめの $q_1 + q_2$ 列を sweep out する時には、計算機の内部では、上式の右辺の形が表われ、直ちに、消滅する事になる。

今 $q_1 + q_2 = p$, $q_3 = 1$, $\Sigma = \begin{bmatrix} A & B \\ B' & D \end{bmatrix}$, $\theta = \begin{bmatrix} C' \\ E \end{bmatrix}$, $F = 0$ とおけば、上式から直ちにわかるよ

うに、全ての q_1 ($1 \leq q_1 \leq p$) に対して $-\hat{\theta}(q_1) \hat{\Sigma}^{-1}(q_1) \hat{\theta}(q_1) = -C'A^{-1}C$

即ち、 p 個の仮説のうちの、第 q_1 番目の仮説

$$H_{q_1} : a_{11} = a_{12} = \dots = a_{1q_1} = 0$$

を検定する為の Hotelling の T^2 統計量が $H = H_{q_1}$ の検定の為の統計量の計算のときの副産物として、同時に計算ができることになる。

(1), (3) の問題についても、行列式の計算過程を段階別に考えればわかるように、 P 次の正定値対称行列の行列式の計算をする時には、この principal minors の数値も副産物として出てくるが、それを利用したものである。

第4節 結論

計算機の発達に伴い、大量のデータの解析が可能になったばかりでなく、高度に複雑な計算が可能になった。その為には、ターンアンドタイムや、入出力に要する時間や、計算実行の所要時間などを考え、なるべく一度に多種類の解析を行なう事が、効果的である。それには、統計推論の論理的順序ではなく、計算過程を分解して考察し、計算機の内部で発生した情報が、出力されないうまま、消滅してしまう事がないよう特に注意する必要がある。

統計解析の理論は、一つ一つの手法については、精密な議論がなされている。しかしながら、幾つかの手法を逐次に適用し、ある種の判断に到達するのが実際のすがたである。とくに複数母集団が存在する場合の、回帰分析には、実質科学の知見から予想あるいは決定されるいくつかの可能性があり、その知見に導かれ、統計解析の手法を用いてデータを解析し、データが何を語っているかを、調べるのが統計解析の手法の目標とする所であり、白紙の状態で、データに接する事は有り得ない。

他方、かりに白紙の状態でデータに接した時、今迄とは異なる全く新しい知見が得られる事の可能性も否定できない。それ故に、データが、得られた時は、出来るだけひろい可能性を求めて、可能な限りの解析を行ない、その結果を、白紙に帰って読み直す事が望ましい。ここで紹介した、プログラムが、開発されたのはその要求をみたすためである。

この論文は、伊藤・大崎・工藤がプログラム作製作業中に行なった討論に基いて、工藤が執筆したものであるが、時間的制約などの関係上3名が会合して執筆したものではない為、文章には色々の不十分な点が多いのではないかと考えられる。そのような不充分さは、あくまで工藤の責任であり、将来このプログラム及びアルゴリズムを英文で発表する際には、訂正改善したいと考えている。

この研究は昭和46年度文部省科学研究補助金 試験研究(1), 課題番号30001によるものである。

REFERENCES

- [1] Schull, W.J. and Kudô, A. (1962) Certain multivariate problems arising in human genetics, *Bulletin of Math. Statist.*, Vol. 10, No. 3/4 pp. 77-88,
- [2] Koichi, I and Kudô, A. (1964) *The Jap. Jour. Human Genetics* Vol. 9, No. 4. pp. 216-220,
- [3] Anderson, T. W. (1958) *An introduction to Multivariate Statistical Analysis*, New York; John Wiley & Sons, Inc.,
- [4] Rao, C. R. (1965) *Linear Statistical Inference and Its Applications*, New York; John Wiley & Sons, Inc.,