

(要約) SAS による多変量データの解析

前川眞一
東京大学出版会 1997

本書の目的は、所謂矩形のデータ行列の分析のために SAS に用意されているプロシジャを、線形モデルと最小2乗法の観点から統一的に表現し、各方法の関連を明確にすることである。

1 多変量解析法と線形モデル

所謂 $N \times p$ のデータ行列 (N 個の個体の p 変量による測定値) を取り扱う多変量解析の方法のうち、利用度の高いものには、回帰分析 (proc reg, glm, transreg)、因子分析 (proc factor)、主成分分析 (proc princomp, prinqual)、クラスター分析 (proc fastclus, cluster, varclus)、冗長性分析 (proc transreg)、等があるが、これらの方法は線形モデルおよび最小2乗法の観点から統一的に説明する事が可能である。

Table 1: 多変量解析法の分類

		説明変数 X は観測されているか	
		観測されている (顕在変数)	観測されていない (潜在変数)
説明 変数 X は	連続	回帰分析 冗長性分析 (判別分析) (ロジスティック回帰)	因子分析 主成分分析 (IRT) (カテゴリーカルデータの 因子分析) (数量化 III 類, CA)
	連続 か 離散 か	分散分析 (対数線形モデル)	クラスター分析 (潜在クラスモデル)

表中、括弧の中は、基準変数 Y がカテゴリーカルな場合、それ以外は、基準変数が連続の場合を表す。また、説明変数が観測されている場合、その連続離散の区別はあまり重要ではない。

一般に p 変数の多変量線形モデルは

$$Y = X\beta' + E \quad (1)$$

という形で書き表されるが、ここに、 Y は $N \times p$ のデータ行列、 X は $N \times q$ の説明変数行列、 β は $p \times q$ の回帰係数行列、 E は $N \times p$ の残差行列である。説明変数行列として観測不可能な潜在変数を含ばこれらの方法は以下のように分類することが出来る。なお、表中に括弧を付した方法は基準変数 Y が離散変数の場合の方法である。また、基準変数の尺度水準に応じた最適を考えると、 Y が順序尺度の場合の取り扱いも可能となる。

1. 回帰分析： X が連続な顕在変数、
 x_i は個体 i の説明変数の値、 β_j は回帰係数
2. 分散分析： X がカテゴリカルな顕在変数、
 x_i は個体 i がどの要因の組み合わせの下で観測されたかを示すインディケイタ（デザイン行列）、 β_j は各要因の水準の効果
3. 因子分析： X が連続な潜在変数、
 x_i は個体 i の因子スコア、 β_j は変数 j の因子パターン
4. クラスタ分析： X がカテゴリカルな潜在変数、
 x_i は個体 i がどのクラスターに属するかを示すインディケイタ、 β_j は、各クラスターにおける変数 j の平均値

という具合になる。なお、通常回帰分析や分散分析の場合には、推定されるべきパラメタは $\beta_j, j=1, \dots, p$ であるが、説明変数が潜在変数の場合には、 X をも推定することになる。

2 データの最適変換を含む多変量線形モデル

2.1 回帰分析

実際にデータを収集する場合、間隔尺度水準以上の変数を用いることは必ずしも可能であるとは限らない。たとえば、ある商品が好まれる理由をその商品の持つさまざまな特性を用いて説明しようとする場合を考えてみる。この場合、基準変数 Y としては n 個の商品の好みの程度を p 人に評定させたものを用いることになるが、その評定値は果たして間隔尺水準の測定値として取り扱ってよいものであろうか。もしも、この評定値が単に好みの大小関係のみを示すものであれば、それは順序尺度水準の測定値であると呼ばれる。従って、 $y_{1j} = 2 < y_{2j} = 5 < y_{3j} = 6, \dots$ というデータは、その順序を保つように、 $y_{1j} = 6 < y_{2j} = 7 < y_{3j} = 8, \dots$ と変換してもそれが持つ情報はなんら失われたことにならない。しかし、これらの評定値を商品の特性 X に回帰する場合、どちらのデータを用いたかにより回帰係数は一般に異なったものとなる。また、商品の特性そのものも、順序尺度や名義尺度水準の評定値として得られているかもしれない。

このような順序尺度もしくは名義尺度水準のデータを用いて回帰分析を行うための方法として、SAS/STAT の PROC TRANSREG があるが、そこでは通常回帰分析のモデルを修正した、以下のようなモデルが用いられている。

$$T_j(y_{ij}) = \mathcal{E}(T_j(y_{ij}) | x_i) + e_{ij} \quad (2)$$

$$\mathcal{E}(T_j(y_{ij}) | x_i) = \mathcal{S}(x_i)' \beta_j \quad (3)$$

$$= \sum_{\ell=1}^q \mathcal{S}_\ell(x_{i\ell}) \beta_{j\ell} \quad (4)$$

ただし、 $T_j()$ は第 j 番目の予測変数を間隔尺度へ変換するための関数、また、 $S()$ は各説明変数を間隔尺度に変換するための関数 $S_\ell()$, $\ell=1,2, \dots, q$, を元のデータの形に並べたもの

$$S(\mathbf{x}_i) = [S_1(x_{i1}), S_2(x_{i2}), \dots, S_q(x_{iq})]' \quad (5)$$

である。Eq.(4) は、変換後の予測変量の \mathbf{x}_i を与えたときの条件付き期待値が、変換後の \mathbf{X} の線形関数で書き表されることを示している。たとえば、被験者 k 以外の評定値は間隔尺度水準の測定値とみなせるが、被験者 k の評定値はその対数をとることによって間隔尺度へ変換されること、また、商品の特性のうち第 m 番目以外のは間隔尺度とみなせるが、第 m 番目の特性はそれを 3 乗することにより間隔尺度へ変換可能である、ということが前もって解っていれば、上記のモデルのパラメタは、 $T_j(y_{ij}) = y_{ij}$, $j = 1, 2, \dots, k-1, k+1, \dots, p$, $T_k(y_{ik}) = \log(y_{ik})$, $S_\ell(x_{i\ell}) = x_{i\ell}$, $\ell = 1, 2, \dots, m-1, m+1, \dots, q$, $S_m(x_{im}) = x_{im}^3$, という変換を施したデータに通常の回帰分析を行えば得られる。ただ、この例のように、変換の形が前もって知られている場合は希であり、実際は、その変換の形も推定することが普通である。($T()$ や $S()$ をの形を規定する変換のパラメタをも推定する。)

上記の回帰分析モデルの場合、最小 2 乗法を用いて推定を行うとすれば、最小 2 乗基準が

$$RSS = \text{tr}((T(\mathbf{Y}) - S(\mathbf{X})\mathbf{B}')'(T(\mathbf{Y}) - S(\mathbf{X})\mathbf{B}')) \quad (6)$$

$$= \sum_{j=1}^p \sum_{i=1}^n (T_j(y_{ij}) - S(\mathbf{x}_i)' \beta_j)^2 \quad (7)$$

$$= \sum_{j=1}^p RSS_j \quad (8)$$

と書けるため、説明変数の変換、 S_ℓ , $\ell=1,2, \dots, q$, が既知の場合、RSS を最小にするような変換 $T_j()$ および回帰係数 β_j は、各予測変数ごとに、 RSS_j を最小とすように求めれば良いことが解る。また、同じ最小 2 乗基準が、第 ℓ 番目の説明変数ベクトル $\mathbf{x}_{(\ell)}$ ($n \times 1$ のベクトル) とそれを除く部分の変換 $S(\mathbf{X}_{(\ell)})$ ($n \times (q-1)$ の行列) を用いて

$$RSS = \text{tr} \left(T(\mathbf{Y}) - S(\mathbf{X}_{(\ell)})\mathbf{B}_{(\ell)}' - S_\ell(\mathbf{x}_{(\ell)})\beta_{(\ell)} \right)' [T(\mathbf{Y}) - S(\mathbf{X}_{(\ell)})\mathbf{B}_{(\ell)}' - S_\ell(\mathbf{x}_{(\ell)})\beta_{(\ell)}] \quad (9)$$

$$= [S_\ell(\mathbf{x}_{(\ell)}) - (T(\mathbf{Y}) - S(\mathbf{X}_{(\ell)})\mathbf{B}_{(\ell)}')\beta_{(\ell)} / (\beta_{(\ell)}'\beta_{(\ell)})]' \\ \times [S_\ell(\mathbf{x}_{(\ell)}) - (T(\mathbf{Y}) - S(\mathbf{X}_{(\ell)})\mathbf{B}_{(\ell)}')\beta_{(\ell)} / (\beta_{(\ell)}'\beta_{(\ell)})] \\ \times (\beta_{(\ell)}'\beta_{(\ell)}) + (S_\ell(\mathbf{x}_{(\ell)}) \text{ を含まない部分}) \quad (10)$$

と書き表せるため、基準変数の変換 $T(\mathbf{Y})$ 、回帰係数行列 \mathbf{B} 、および $S(\mathbf{X}_{(\ell)})$ が既知の場合には、第 ℓ 番目の説明変数の変換 $S_\ell(\mathbf{x}_{(\ell)})$ は上記の 2 次形式を最小とすように求めればよいことが解る。

このように、モデルに含まれるパラメタをいくつかのサブセットに分轄し、パラメタの各サブセットごとに同じ最小 2 乗基準を最小化する様な推定を行うことを繰り返す方法は、交互最小 2 乗法 (ALS) と呼ばれ、複雑なモデルのパラメタを推定する際の常套手段となっている。回帰分析モデルの場合は、データの変換が与えられた場合の回帰係数の推定は、

$$\hat{\mathbf{B}}' = (S(\mathbf{X})'S(\mathbf{X}))^{-1}S(\mathbf{X})'T(\mathbf{Y}) \quad (11)$$

となることは明かである。また、Eq.(7) や Eq.(10) は、 (u_i, v_i) を既知のデータ、 \mathcal{F} を推定すべき変換とすれば、

$$RSS = \sum_{i=1}^n (\mathcal{F}(v_i) - u_i)^2 + (\mathcal{F} \text{ のパラメタを含まない部分}) \quad (12)$$

という形をしているため、 $\mathcal{F}()$ に含まれるパラメタの推定は、 ou の v へ回帰する 1 変量の回帰分析の問題となる。特に、PROC TRANSREG においては、データの変換として、スプライン関数が用いられる

ことが多いが、これは変換のパラメタに関して線形に書き表すことが出来るため、その推定は簡単である。また、 $\mathcal{F}()$ として任意の単調変換を仮定すれば、Eq.(12) を最小にするようなものは、Kruskal の最小2乗単調変換として計算することが出来るし、名義尺度のデータは Fisher の方法を用いてカテゴリの平均値を計算すればよい。このように、変数の尺度水準を考慮しながら、その変数の持つ情報を損なわない形で最小2乗基準を最小にするようなデータの変換を求めることを、データの最適変換と呼ぶ。なお、スプライン関数を用いてデータを変換する場合は、単にその変換が連続で滑らかなものであることのみを仮定していることになる。

PROC TRANSREG は、きわめて一般的なプロシジャであるが、主に、説明変数として名義尺度の変数を用いたコンジョイント分析としてマーケティングの分野で使用されることが多い。また、順序データの PREFMAP モデルによる分析にも利用することが出来る。

2.2 因子分析

上記の方法は、そのまま因子分析モデルに応用することが出来る。すなわち、データを変換したものが、潜在変数である因子スコアと線形に関係していることを仮定するモデル、

$$\mathcal{E}(\mathcal{T}_j(y_{ij})|\mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}_j \quad (13)$$

である。これは、まとめて書き表せば

$$\mathcal{E}(\mathcal{T}(\mathbf{Y})|\mathbf{X}) = \mathbf{X}\mathbf{B}' \quad (14)$$

となるが、 \mathbf{X} は潜在変数であるため、その変換 $\mathcal{S}()$ をモデルにいれる必要はない。回帰分析と同様に ALS による最小2乗推定を行うことにすれば、最小2乗基準は以下のように書き表される。

$$\text{RSS} = \text{tr}((\mathcal{T}(\mathbf{Y}) - \mathbf{X}\mathbf{B}')'(\mathcal{T}(\mathbf{Y}) - \mathbf{X}\mathbf{B}')) \quad (15)$$

$$= \sum_{j=1}^p \sum_{i=1}^n (\mathcal{T}_j(y_{ij}) - \mathbf{x}_i' \boldsymbol{\beta}_j)^2 \quad (16)$$

$$= \sum_{j=1}^p \text{RSS}_j \quad (17)$$

したがって、変換を既知とした場合の因子スコア \mathbf{X} および因子パターン \mathbf{B} の推定値は、 $\mathcal{T}(\mathbf{Y})$ をあたかもデータ行列のように取り扱った因子分析の解（いわゆる主成分解および主成分スコア）

$$\hat{\mathbf{B}} = \mathbf{P}\boldsymbol{\Lambda}^{\frac{1}{2}} \quad (18)$$

ただし、 $\boldsymbol{\Lambda}$ と \mathbf{P} は $\frac{1}{n}\mathcal{T}(\mathbf{Y})'\mathcal{T}(\mathbf{Y})$ の最大 q 個の固有値と固有ベクトル行列、

$$\mathbf{X} = \mathcal{T}(\mathbf{Y})\mathbf{B}(\mathbf{B}'\mathbf{B})^{-1} = \mathcal{T}(\mathbf{Y})\mathbf{P}\boldsymbol{\Lambda}^{-\frac{1}{2}} \quad (19)$$

として与えられる。また、データの最適変換は、各変数毎に、Eq.(16) を最小とするように回帰分析の場合と全く同様にして求めればよい。この、データの最適変換を含む因子分析モデルによる分析は、SAS の PROC PRINQUAL を用いて行うことが出来る。データの最適変換という過程が含まれるため、通常の因子分析を行う場合よりも少ない因子の数でデータを説明することが可能となる。

3 類似度データのモデル

他方、多変量解析には、上記のようなプロフィールデータの分析ではなく、個体相互間の関係を取り扱う方法も存在する。たとえば、個体 i と個体 j の類似度の測定値を、 o_{ij} , $i, j=1, \dots, n$ とすれば、その期

待値の構造として個体の q 次元の尺度値（潜在変数による個体の測定値） $\mathbf{x}_i, i=1, 2, \dots, n$, から計算された距離を当てはめるモデル

$$\mathcal{E}(\mathcal{T}(o_{ij})|\mathbf{x}_i, \mathbf{x}_{(j)}) = \sqrt{(\mathbf{x}_i - \mathbf{x}_{(j)})'(\mathbf{x}_i - \mathbf{x}_{(j)})} \quad (20)$$

を考えることができる。この場合、もはやモデルはパラメタに関して線形とはならないが、ALS を用いれば、

$$\text{RSS} = \sum_{i>j} \left(\mathcal{T}(o_{ij}) - \sqrt{(\mathbf{x}_i - \mathbf{x}_{(j)})'(\mathbf{x}_i - \mathbf{x}_{(j)})} \right)^2 \quad (21)$$

を最小とするデータの変換は前節と同様にして求められることが解る。また、変換が既知の場合の尺度値の推定は、例えば、Gauss-Newton 法を用いて数値的に行うことができる。SAS の PROC MDS はこのモデルの下での解析を行うためのもので、類似度が多数の個人から得られている場合には、個人差を各尺度値にかかる重みの大小という形でモデルに組み込んだ、いわゆる重み付きユークリッドモデル（個人差 MDS モデル）

$$\mathcal{E}(\mathcal{T}_k(o_{ijk})|\mathbf{x}_i, \mathbf{x}_{(j)}, \mathbf{O}W_k) = \sqrt{(\mathbf{x}_i - \mathbf{x}_{(j)})' \mathbf{O}W_k (\mathbf{x}_i - \mathbf{x}_{(j)})} \quad (22)$$

による分析も行うことができる。

4 最小2乗法によるパラメタの推定

データとモデルの値との残差の2乗和、

$$\text{RSS} = \sum_{j=1}^p \sum_{i=1}^N e_{ij}^2 = \sum_{j=1}^p \sum_{i=1}^N (y_{ij} - \mathbf{x}_i \boldsymbol{\beta}_j)^2 \quad (23)$$

で定義される最小2乗基準 RSS を最小とするようにモデルのパラメタ、 $\boldsymbol{\beta}$ を推定する方法を最小2乗法と呼び、推定されたパラメタの値を $\hat{\boldsymbol{\beta}}$ と記す。

残差の2乗和 RSS をモデルのパラメタの関数と見た場合、その下限はゼロであるから必ず最小値が存在し、その点では RSS のすべてのパラメタに関する偏微分がゼロになっているはずである。従って、最小2乗推定値を求めるためには、正規方程式

$$\frac{\partial \text{RSS}}{\partial \beta_{j\ell}} = 0, \quad j = 1, 2, \dots, p, \ell = 1, 2, \dots, q \quad (24)$$

を解けばよい。また、モデルが潜在変数を含む場合には \mathbf{X} をもモデルのパラメタとして取り扱い

$$\frac{\partial \text{RSS}}{\partial x_{i\ell}} = 0, \quad i = 1, 2, \dots, N, \ell = 1, 2, \dots, q \quad (25)$$

を満たす \mathbf{X} の推定値 $\hat{\mathbf{X}}$ をも同時に求めることになる。ただ、因子分析モデルの場合、 \mathbf{X} と $\boldsymbol{\beta}$ の間に乗算的不定性が存在するため、一意な解を求めるためには何らかの制約を課す必要が生じる。また、クラスター分析モデルの場合には、 \mathbf{X} の各行は一つの1と $q - 1$ 個のゼロからなるという制約を付ける。

最後に、データの最適変換を含むモデルの場合には、最適変換のパラメタを $\boldsymbol{\Gamma}, \boldsymbol{\Lambda}$ とすれば、RSS は

$$\text{RSS} = \sum_{j=1}^p \sum_{i=1}^N (\mathcal{T}_j(y_{ij}|\boldsymbol{\gamma}_j) - \mathcal{S}(\mathbf{x}_i|\boldsymbol{\Lambda})\boldsymbol{\beta}_j)^2 \quad (26)$$

と書き表されるが、この RSS を最小とするように最適変換のパラメタとモデルのパラメタを同時に推定することになる。ただ、一般に、基準変数の最適変換を含む場合には、 $\|\mathcal{T}_j(\mathbf{Y}|\boldsymbol{\gamma}_j)\|^2$ を小さくすることにより RSS は限りなくゼロに近づくため、RSS を $\|\mathcal{T}(\mathbf{Y}|\boldsymbol{\Gamma})\|^2$ で基準化した、規準化最小2乗基準

$$\text{RSS}_1 = \sum_{j=1}^p \frac{\sum_{i=1}^N (\mathcal{T}_j(y_{ij}|\boldsymbol{\gamma}_j) - \mathcal{S}(\mathbf{x}_i|\boldsymbol{\Lambda})\boldsymbol{\beta}_j)^2}{\text{var}(\mathcal{T}_j(y_{ij}|\boldsymbol{\gamma}_j))} \quad (27)$$

を最小化することが多い。すなわち、正規方程式

$$\frac{\partial \text{RSS}_1}{\partial \beta_{j\ell}} = 0, \quad j = 1, 2, \dots, p, \ell = 1, 2, \dots, q \quad (28)$$

$$\left(\frac{\partial \text{RSS}_1}{\partial x_{i\ell}} = 0, \quad i = 1, 2, \dots, N, \ell = 1, 2, \dots, q \right) \quad (29)$$

$$\frac{\partial \text{RSS}_1}{\partial \gamma_j} = 0, \quad j = 1, 2, \dots, p \quad \text{and} \quad \frac{\partial \text{RSS}_1}{\partial \lambda_\ell} = 0, \quad \ell = 1, 2, \dots, q \quad (30)$$

$$(31)$$

が解が求める最小2乗推定値となる。この連立方程式が解析的に解けない場合には交互最小2乗法（ALS）を用いて数値的に解を求めることになる。

5 モデルの適合度

モデルの最小2乗推定値を最小2乗法でパラメタを推定する場合、データへのモデルの適合度を表す指標として、分散説明率、 R^2 を考えると都合がよい。いま、モデルの最小2乗推定値を

$$\hat{\mathbf{y}}_{(j)} = \mathbf{X}\hat{\boldsymbol{\beta}}_j \quad \text{または} \quad \widehat{\mathbf{X}}\hat{\boldsymbol{\beta}}_j \quad (32)$$

で表す。ただし、 $\hat{\mathbf{y}}_{(j)}$ は第 j 番目の変数のモデルの最小2乗推定値を要素として持つ $N \times 1$ のベクトルである。この記法および列中心化行列 $\mathbf{J} = \mathbf{I} - \frac{1}{N}\mathbf{1}\mathbf{1}'$ を用いれば、第 j 番目の変数の分散説明率 R_j^2 は

$$R_j^2 = \frac{\text{var}(\hat{\mathbf{y}}_{(j)})}{\text{var}(\mathbf{y}_{(j)})} = \frac{\hat{\mathbf{y}}_{(j)}\mathbf{J}\hat{\mathbf{y}}_{(j)}}{\mathbf{y}_{(j)}\mathbf{J}\mathbf{y}_{(j)}} \quad (33)$$

と定義されるが、これは、モデルの推定値の分散のデータとして得られる分散に対する比に他ならない。すなわち、 R_j^2 が 1 の場合には変数 j の分散の全てがモデルによって説明されたことになる。また、この値は、 $\mathbf{y}_{(j)}$ と $\hat{\mathbf{y}}_{(j)}$ の間の相関係数の2乗

$$R_j^2 = \left(\text{corr}(\mathbf{y}_{(j)}, \hat{\mathbf{y}}_{(j)}) \right)^2 \quad (34)$$

となっている。

基準変数が p 個ある場合は全体の分散説明率をいた通りに定義することが出来る。一つは

$$R_T^2 = \frac{\sum_{j=1}^p \text{var}(\hat{\mathbf{y}}_{(j)})}{\sum_{j=1}^p \text{var}(\mathbf{y}_{(j)})} = \frac{\sum_{j=1}^p \text{var}(\mathbf{y}_{(j)})R_j^2}{\sum_{j=1}^p \text{var}(\mathbf{y}_{(j)})} = \sum_{j=1}^p \frac{\text{var}(\mathbf{y}_{(j)})}{\sum_{k=1}^p \text{var}(\mathbf{y}_{(k)})} R_j^2 \quad (35)$$

で定義される、全分散説明率 R_T^2 であり、もう一つは

$$R_A^2 = \sum_{j=1}^p R_j^2 / p = \sum_{j=1}^p \frac{\text{var}(\hat{\mathbf{y}}_{(j)})}{\text{var}(\mathbf{y}_{(j)})} / p \quad (36)$$

で定義される平均分散説明率 (Average r-squared R_A^2 である。

一般に、パラメタを最小2乗法を用いて推定した場合、その値は分散説明率を最大とするようなものになっていることがいえるが、データの最適変換がなくモデルが潜在変数を含まない場合には、RSS を最小とするパラメタを求めることにより R_T^2 と R_A^2 の両方が最大化され、どちらかを含む場合には、 RSS_1 の最小化により R_A^2 が最大化されることになる。

6 説明変数の寄与率

ひとつの基準変数を q 個の説明変数で説明する際に、各説明変数がどの程度貢献しているかの測度を説明変数の寄与率と呼ぶ。寄与率の定義には様々なものが考えられるが、モデルから $\mathbf{x}_{(k)}$ を除いたときの R^2 の減少量を持って寄与率の測度とすることが考えられる。この場合、第 k 番目の説明変数 ($\mathbf{x}_{(k)}$) の寄与率 C_k を求めるためには、説明変数及びパラメタを

$$\mathbf{X} = [\mathbf{X}_{(K)} \ \mathbf{x}_{(k)}] \quad \text{and} \quad \boldsymbol{\beta} = [\boldsymbol{\beta}_K, \beta_k] \quad (37)$$

ただし、 $\mathbf{X}_{(K)}$ は全説明変数から $\mathbf{x}_{(k)}$ を除いた $N \times q - 1$ の行列、と分割し、モデル

$$\mathbf{y} = \mathbf{X}_{(K)}\boldsymbol{\beta}'_K \quad (38)$$

で計算される R^2 を R_1^2 、モデル

$$\mathbf{y} = \mathbf{X}_{(K)}\boldsymbol{\beta}_K + \mathbf{x}_{(k)}\beta_k \quad (39)$$

で計算される R^2 を R_2^2 とすれば、

$$C_k = R_2^2 - R_1^2 \quad (40)$$

$$= \frac{\mathbf{y}'(\mathbf{P}_X \mathbf{J} \mathbf{P}_X - \mathbf{P}_K \mathbf{J} \mathbf{P}_K)\mathbf{y}}{N \text{var}(\mathbf{y})} \quad (41)$$

$$= (\text{corr}(\mathbf{y}, \mathbf{Q}_K \mathbf{x}_{(k)} \mathbf{y}))^2 \quad (42)$$

を計算すればよい。ただし、 \mathbf{P}_X および \mathbf{P}_K は \mathbf{X} と $\mathbf{X}_{(K)}$ の列空間への射影行列また、 $\mathbf{Q}_K = \mathbf{I} - \mathbf{P}_K$ である。なお、この $\mathbf{Q}_K \mathbf{x}_{(k)}$ は、 $\mathbf{x}_{(k)}$ を $\mathbf{X}_{(K)}$ の線形モデルで説明した場合の最小 2 乗残差となっているが、上記の寄与率の定義は、 $\mathbf{x}_{(k)}$ から他の説明変数で説明する部分を除いたもの（残差）と基準変数との相関係数の 2 乗となっている。その意味で、寄与率は片側偏相関係数の 2 乗 (squared semipartial correlation coefficient) もしくは、部分相関係数の 2 乗 (squared part correlation coefficient) と呼ばれることがある。なお、Eq.(42) の寄与率の分子は、いわゆる調整済み平方和、SAS の type II SS (または type III SS) であるが、このように定義された寄与率は加算性を持たない。すなわち、全ての寄与率の和に $\text{var}(\hat{\mathbf{e}})/\text{var}(\mathbf{y})$ を加えたものが 1 とはならない。加算性を持つ寄与率は SAS の type I SS をもとにして定義する事が出来るが、この場合には順序依存性という欠点が生じる。すなわち、同じ変数が、説明変数行列の何番目に位置するかによって、その寄与率が異なってくる。

また、 $\mathbf{x}_{(k)}$ が $\mathbf{X}_{(K)}$ の全ての列と直交する場合は $\mathbf{Q}_K \mathbf{x}_{(k)} = \mathbf{x}_{(k)}$ となるため、

$$C_k = (\text{corr}(\mathbf{y}, \mathbf{x}_{(k)}))^2 \quad (43)$$

となり、単純な寄与率と一致する。

基準変数が複数個ある場合の説明変数の寄与率は、それを説明変数を除いたときの、 R_T^2 の減少量と定義するか R_A^2 の減少量と定義するかで異なってくるが、 R_T^2 を用いれば、

$$C_{Tk} = \frac{\sum_{j=1}^p (\text{var}(\mathbf{y}_{(j)}) \times (\text{corr}(\mathbf{y}_{(j)}, \mathbf{Q}_K \mathbf{x}_{(k)}))^2)}{\sum_{j=1}^p \text{var}(\mathbf{y}_{(j)})} \quad (44)$$

となり、 R_A^2 を用いれば、

$$C_{Ak} = \frac{1}{p} \sum_{j=1}^p (\text{corr}(\mathbf{y}_{(j)}, \mathbf{Q}_K \mathbf{x}_{(k)}))^2 \quad (45)$$

となる。従って、基準変数の分散が同一の場合には $C_{Tk} = C_{Ak}$ であり、説明変数の各列が直交している場合には、

$$C_{Tk} = \frac{\sum_{j=1}^p \left(\text{var}(\mathbf{y}_{(j)}) \times (\text{corr}(\mathbf{y}_{(j)}, \mathbf{x}_{(k)}))^2 \right)}{\sum_{j=1}^p \text{var}(\mathbf{y}_{(j)})} \quad \text{and} \quad C_{Ak} = \frac{1}{p} \sum_{j=1}^p (\text{corr}(\mathbf{y}_{(j)}, \mathbf{x}_{(k)}))^2 \quad (46)$$

となる。

7 統計的推測について

デザイン行列 X をパラメタとして取り扱うモデル（因子分析やクラスター分析）には、観測値数を増やすとそれだけパラメタ数も増えるという欠点がある。この問題を解決するためには、デザイン行列の各行が同一分布からのランダムサンプルであることを仮定した後にそれらを周辺化した β のみの周辺尤度を最大とするような解を求めればよい。因子分析の場合は通常用いられている最尤解がこれにあたり、クラスター分析の場合には、混合正規分布のパラメタの推定問題となる。

8 おわりに

本稿で紹介したデータ解析の方法は、特に、因子分析モデルや MDS モデルは、分析の結果として得られるパラメタの値の大小を議論することよりも、それを基にして個体と変数の関係に視覚的に表現を与えることを主眼とした方法である。その意味でこれらの方法は探索的な方法であると言える。また、推定されたデータ変換の意味を考える為にはそのプロットは必須である。

訂正表

前川眞一

970819,1016,010329,031114

第2刷の修正は頁番号に * を付けて示した。

第 2 章

1. p22, 表 2.6: 変数 X の相関行列も掲示されているがそれは無視。
2. p23, 表 2.7: 変数 X を含む相関係数も掲示されているがそれは無視。
3. p25, lb: データセットの転置の説明は p29 を参照のこと。
4. p32, lb8: $\text{var } x_{1-98}$ とすると全ての相関係数を計算する。

第 3 章

1. p63, lb1: $\mathbf{x}_i \leftarrow \mathbf{b}$
2. p63, lb4: $\mathbf{x}_i \leftarrow \mathbf{b}$
3. p70, lb2: M 次Pスプラインで近似する
4. p76, lb9: $\mathbf{H}^{(p)}(\mathbf{X}) = [\mathbf{g}_2^+(\mathbf{X}) \quad -2\mathbf{X}]$
5. p110, lb9: &NFACT が未定義となっている。例えば、`%let nfact=2` とする。
6. p110, lb9: &NFACT が未定義となっている。例えば、`%let nfact=2` とする。
7. p119, 図 3.13(a): 非階層的方法の個体 4 と 5 を横線でつなぐ。
8. p125, lb2: その解の R^2 の値は、その直線より下方に位置するクラスターを定義する横線の中で最も高いものの縦軸の座標として示される。

第 4 章

1. p150, lb6:

$$\text{RSS} = \text{tr} \left[(\mathbf{V} - \mathbf{U}\hat{\mathbf{A}}\hat{\mathbf{B}}')' (\mathbf{V} - \mathbf{U}\hat{\mathbf{A}}\hat{\mathbf{B}}') \right]$$

2. p150, lb13:

$$\text{Rss}_w = \|(\mathbf{V} - \mathbf{U}\mathbf{A}\mathbf{B}')\mathbf{W}^*\|^2$$

3. p156, 図 4.1: 縦軸は \mathbf{v} と $\hat{\mathbf{v}}$ を、横軸は \mathbf{u} を表す。
4. p156, 図 4.2: 縦軸は \mathbf{v} と $\hat{\mathbf{v}}$ を、横軸は \mathbf{u} を表す。
5. p175, lb9

$$= -2\mathbf{U}'\mathbf{v} + 2\mathbf{U}'\mathbf{U}\mathbf{x}_i - 2\lambda\mathbf{U}'\mathbf{U}\mathbf{x}_i$$

第 5 章

1. p184* 式 (5.2) 三角不等式の訂正

$$d_{ij} + d_{jk} \geq d_{ik} \tag{47}$$

2. p186, lb10: INDSCAL Model)
3. p204, 表 5.2 (中) : 非類似度データ行列

第 6 章

1. p205, lb10: 回帰係数と定数項を $\hat{\beta}$, $\hat{\mu}$ として、

付録 A

1. p213, lb11-10: 「この場合、 $AA' = I$ となる。」を削除。

付録 C

1. p231, lb21: print ステートメントに obs= データセットオプションを付けると ... (例は p33 に示した。)
2. p245, lb2: abc="MACV"; bcd='hello'; xmean=12.6;
3. p246, lb12: out* で始まる行は削除。
4. p248, print1 を用いて n 個のオブザベーションのみをプリントしたい場合には
print1 datasetname(OBS=n)
とする。
5. p254, lb1: ぜならば、マクロ変数を利用した場合は、
6. p254, lb2: ル積は次節の例で示すように
7. p256, lb6: プログラムの最後に ; を付ける。
8. p262, lb16: 列 T, および RSS を変数の数と因子の数の積で除したものの平方根 RMSE が返される。

参考文献

1. p273: Crawford, C. B. (1975)

索引

1. p286: proc iml 7, 32, 92, 160, 162, 167, 180, 190, 196, 250