

## 制約付き主成分分析法について

高 根 芳 雄\*

ON CONSTRAINED PRINCIPAL COMPONENT ANALYSIS

Yoshio TAKANE

Constrained principal component analysis (CPCA) was proposed by Takane and Shibayama (1991) for structural analysis of multivariate data. In this method the data are first decomposed into several components according to external information. The decomposed submatrices are then subjected to principal component analysis (PCA) to explore possible structures within the submatrices. The method thus combines two major conventional multivariate analysis techniques, multiple regression analysis and PCA, in a unified framework. This paper illustrates the basic model, computational methods, various uses and extensions of CPCA. An illustrative example is given, and relative merits and demerits of CPCA are discussed in relation to the analysis of covariance structure (ACOVs) approach.

### 1. 序

データの背後にある現象を理解するために、データを構造部分と誤差成分に分けて考えるのが統計学の常套手段である。データが多変量の場合はこれに情報の縮約機能が加わる。本論文のテーマである制約付き主成分分析法（以下 Constrained Principal Component Analysis の頭文字をとって CPCA と略記する。）は与えられた多変量データを外部情報によって構造部分と誤差成分に分け（外部分析）、さらにそれぞれの部分を情報縮約の観点から構造部分と誤差成分に分解していく（内部分析）方法である。後者により元の構造部分からはその最も重要な変動部分、誤差成分からはその最も構造化し易い部分を取り出し、結果をグラフィカルに提示することができる。

本論文ではまず CPCA が適用されるデータの要件について述べる（第 2 節）。次に CPCA の基本モデルと計算法を説明する（第 3 節）。さらに基本モデルの拡張、CPCA の関連手法についても簡単に触れる。第 4 節では CPCA の適用の実際場面を応用例を通して概観する。最後に CPCA の利点、限界を共分散構造分析 (ACOVs) と比較しながら考察する。

\* マツギル大学心理学系 (Department of Psychology, McGill University)

### 2. データの要件

CPCA は簡単に言うと、多変量データが与えられた時、これを外部情報によって説明できる部分とできない部分に分け分解されたそれぞれの部分に主成分分析 (PCA) を適用する方法である。従って、CPCA を有効に用いるためには何らかの外部情報が必要である。PCA ではモデルをデータに当てはめる時、重みを付けない単純最小 2 乗基準が用いられることが多い。CPCA ではこの重みに相当する計量 (metric) 行列を指定し、一般化最小 2 乗基準を用いることもできる。そこでまず入力データ、外部データ、計量行列の説明から始めることにしよう。

#### 2.1 データ行列

データ行列を  $N \times n$  の行列  $Z$  で表す。  $Z$  の行は  $N$  人の被験者（ケースとも言う。）、列は  $n$  個の変数を表すものとする。データは原則的には多変量データであれば何でもよい。特に分布の仮定は設けない。データは数値データでも、ダミー変数化されたカテゴリー・データでもよい。もちろん両者が混在していてもかまわない。

データは粗データのままだでも、何らかの前処理を施したものでよい。ここでいう前処理とはデータの中

心化や標準化を指す。もっとも、どちらでもよいというのはどちらでも解析できるという意味であって、結果が不変だという意味ではない。むしろPCAの結果は前処理のしかたによって大きく左右されるから、どのような前処理をするかの選択は研究者の興味と目的に合わせて慎重に行われなければならない。

データが数量とカテゴリー・データの両方を含む場合、両者の間の単位の両立性が問題となる。この場合、変数を一様に標準化することが多いが、Kiers (1991) はカテゴリー変数を変数ごとに中心化した後、正規直交化することを勧めている。

## 2.2. 外部データ

外部データには主データの行に関するものと、列に関するものの2種類がある。行に関するものを行デザイン行列と言い、 $N \times p$ の行列  $G$  で表す。同様に列に関するものを列デザイン行列と呼び、 $n \times q$ の行列  $H$  で表す。どちらか一方が欠けている場合は、 $G=I$  または  $H=I$  とする。

$G$  の例としては被験者の性別、年齢、学歴などの人口統計学的変数が考えられる。この場合、主データがこれらの変数といかに関係しているかを分析できる。また  $G$  としてその要素がすべて1よりなる  $N$  次のベクトル ( $1_N$ ) をとると、被験者全体の平均的傾向を見ることができる。これをやや一般化したものが被験者の所属するグループを示すダミー変数である。このダミー変数行列を  $G$  と置くと、グループ間の平均的傾向の違いが分析できる。グループは性差のように統制できないグループもあれば、実験などで人工的に作られたグループの場合もある。

一方、 $H$  としては  $G$  と同様のものを変数について考えればよい。具体例としては、変数が何らかの刺激を表す時、刺激の特徴を示す特徴行列や刺激のデザイン行列があげられる。また変数が異なった実験条件を表す場合にはコントラスト、異なった時点での繰り返し測定を表す場合にはトレンドなどを指定することも考えられる。後述の例 (第4節) ではデータが刺激対の比較選好データであることから、対比較のデザイン行列 (どの変数がどの刺激対の比較結果を表しているかを示す行列) を  $H$  として用いている。

特定の  $G$  や  $H$  を指定することはデータの行空間、列空間を  $G$  や  $H$  の列ベクトルによって張られる空間に限定することを意味する。このことは同時にそれとは逆の部分、すなわち  $G$  や  $H$  とは無関係な部分を指

定することに一致する。この事実を積極的に用いると、データ行列の中から特定の変数の影響を取り除いた残りの部分の構造だけを取り出して分析することができる。

例えばデータの中心化は  $G=1_N$  として、被験者の平均的傾向を取り除いた残りの部分を分析することに他ならない。

## 2.3. 計量行列

計量行列もデータの行側と列側の2種類が存在する。行側の計量行列を  $N$  次の正方行列  $K$ 、列側の計量行列を  $n$  次の正方行列  $L$  で表す。どちらも対称で非負定値 (nnd) であることを仮定する。 $K$  や  $L$  が正定値 (pd) でなく、半正定値 (psd) の場合は次の条件を満たすことを仮定する

$$\text{rank}(KG) = \text{rank}(G), \quad (1a)$$

$$\text{rank}(LH) = \text{rank}(H). \quad (1b)$$

( $K, L$  が pd であれば上の条件が自動的に満足される。)

既に述べたように計量行列はモデルをデータに当てはめる時の基準に深くかかわっている。データ行列の行空間、列空間の座標軸が比較可能な単位を持ち、互いに直交しているならば特別な計量行列を用いなくて済む。 $K=I, L=I$  と置いて重みの付かない単純最小2乗基準を用いればよい。しかし、例えば身長や体重のように相互に比較不能な単位で測定された変数が混在しているデータにPCAを適用する場合は単位行列以外の計量行列が必要である。第2.1節でデータの標準化について述べたが、これは  $L$  として標本分散の逆行列を対角要素とする対角行列を用いることに相当する。

データの行間に相関が存在する場合も特殊な計量行列が必要となる。通常PCAでは、データの行間に独立性が仮定される。これはデータの行がランダム・サンプルされた被験者に対応する場合は殆ど問題とならないが、多変量データを時系列的にとったような場合明らかに問題となる。これはデータの行空間を規定する座標軸が斜交していることを意味する。従って、時系列相関を何らかの方法で推定し、その逆行列を  $K$  として用いて座標軸を直交化してやらなければならない。また行間に相関がない場合でも、行によって重要性に差があったり、信頼性に差がある場合は行によって重み付けを変える働きを持つ対角行列が  $K$  として用い

られる。例えばデータが既にグループごとに集計されている場合、当然のことながらグループ内の標本数によって数値の信頼性に差が出てくるから、その信頼性を反映した重み付けが必要となる。「対応分析 (Correspondence Analysis)」と呼ばれる分割表の主成分分析 (Greenacre, 1984; Nishisato, 1980) では通常、分割表の行と列をそれぞれ行和と列和の平方根の逆数で重み付けるが、これも一種の信頼性を反映させた重み付けといえることができる。

一方、データ行列の列間に相関があるのはむしろ当然である。この場合、データの列空間を分析に先立って直交化することはまずない。PCA はそもそも列間の相関構造を少数の成分によって説明しようとするものだからである。しかし、列によってモデルを当てはめた後の誤差の大きさに差があったり、誤差行列の列間に相関がある場合は、誤差の共分散構造を何らかの形で推定し、その逆行列を  $L$  として誤差を正規直交化することが行われる。これは誤差の全体的大きさを評価する際、冗長な部分を排除したり、どの変数の誤差も大きさをそろえ平等に扱うことによってパラメータの推定精度を高める働きがある。また、付随的效果として結果の尺度不変性 (scale invariance) が得られるという利点がある (Rao, 1964, section 9)。Meredith & Millsap (1985) は誤差分散を変数の信頼性係数 (例えば再検査法による) や、Guttman (1953) のイメージ理論に基づいて推定する方法を提案している。

### 3. 方 法

本節では CPCA の概要について述べる。記述は簡潔にならざるを得ないが、詳細は Takane & Shibayama (1991), Takane (1991) を参照されたい。既に述べたように、CPCA には「外部分析」と「内部分析」の2つの柱がある。これらの分析を順を追って説明して行くことにしよう。

#### 3.1 外部分析

前節で定義されたデータ行列  $Z$ 、外部データ  $G, H$  を用いて外部分析における CPCA の基本モデルを次の様に定式化する。

$$Z = GMH' + BH' + GC + E. \quad (2)$$

ここで  $M(p \times q)$ ,  $B(N \times q)$ ,  $C(p \times n)$  はモデルのパラメータ行列、 $E(N \times n)$  は誤差行列を表す。このモデルには全部で4つの項が存在するが研究者の興味や目

的、データ解析の状況に合わせて、一部を省略してもよい。

上のモデルには冗長性がある。冗長性を取り除くために次の様な直交条件を導入する。

$$G'KB=0, \quad (3a)$$

$$CLH=0. \quad (3b)$$

ここで  $K, L$  は前節で導入した計量行列である。これらの直交条件によりモデルの各項が一義的に解釈できるようになる。第1項は  $Z$  のうち  $G$  と  $H$  の両方によって説明できる部分、第2項は  $H$  によって説明できるが  $G$  によっては説明できない部分、第3項は  $G$  によって説明できるが  $H$  によっては説明できない部分を表す。

モデルのパラメータを次の基準を最小化するように定める。

$$SS(E)_{K,L} = \text{tr}(E'KEL). \quad (4)$$

(これは一般化最小2乗基準である。) これより  $M, B, C, E$  の最小2乗 (LS) 解は

$$\hat{M} = (G'KG)^{-1}G'KZLH(H' LH)^{-1}, \quad (5a)$$

$$\hat{B} = Q_{G/K}ZLH(H' LH)^{-1} \quad (5b)$$

$$\hat{C} = (G'KG)^{-1}G'KZ Q_{H/L}', \quad (5c)$$

$$\hat{E} = Q_{G/K}ZQ_{H/L}' \quad (5d)$$

と求められる。ここで  $\hat{E}$  は一意に定まるが、 $\hat{M}, \hat{B}, \hat{C}$  は通常一意にはならない。これらが一意であるためには  $G'KG, H' LH$  が正則である必要がある。) ここで  $X$  は  $X$  の一般化逆行列、 $Q_{G/K}, Q_{H/L}$  はそれぞれ

$$P_{G/K} = G(G'KG)^{-1}G'K, \quad (6a)$$

$$P_{H/L} = H(H' LH)^{-1}H'L \quad (6b)$$

として、 $Q_{G/K} = I - P_{G/K}$ ,  $Q_{H/L} = I - P_{H/L}$  と定義される。条件 (1) のもとで  $P_{G/K}, Q_{G/K}$  は次の条件を満たす射影行列となる。(例えば、柳井・竹内 (1983) を参照のこと。)

$$(P_{G/K})^2 = P_{G/K}, (Q_{G/K})^2 = Q_{G/K}, \quad (7a)$$

$$P_{G/K} Q_{G/K} = Q_{G/K} P_{G/K} = 0, \quad (7b)$$

$$P_{G/K}' K P_{G/K} = P_{G/K}' K = K P_{G/K}. \quad (7c)$$

条件 (1b) のもとで  $P_{H/L}, Q_{H/L}$  にも同様の性質が成り立つ。 $\text{Sp}(X)$ ,  $\text{Ker}(X')$  をそれぞれ  $X$  の列ベクトルによって張られる空間、 $X'$  のゼロ空間 ( $X'u=0$  を満たす全てのベクトル  $u$  の張る空間) と定義すると、 $P_K$

は  $\text{Ker}(G'K)$  に沿った  $\text{Sp}(G)$  への射影,  $Q_{G/K}$  は  $\text{Sp}(G)$  に沿った  $\text{Ker}(G'K)$  への射影を表す.  $P_{H/L}, Q_{H/L}$  にも同様の解釈が成り立つ.  $K=I, L=I$  の時, これらの射影行列は直交射影行列に還元する. 直交射影行列は対称性を満たす.

$M, B, C, E$  の中の独立なパラメータ数はそれぞれ  $pq, (N-p)q, p(n-q), (N-p)(n-q)$  と, 見かけのパラメータ数よりも少ない. これは (3) の直交性の条件のためである.

(5) の推定値を (2) に代入すると

$$Z = P_{G/K} Z P_{H/L}' + Q_{G/K} Z P_{H/L}' + P_{G/K} Z Q_{H/L}' + Q_{G/K} Z Q_{H/L}' \quad (8)$$

が得られる. 上式の右辺の4項は計量行列  $K, L$  のもとで直交している. (2つの行列  $X, Y$  が  $\text{tr}(X'KYL) = 0$  を満たす時,  $X$  と  $Y$  は計量行列  $K, L$  のもとで直交しているという.) これより,

$$\begin{aligned} \text{SS}(Z)_{K,L} = & \text{SS}(P_{G/K} Z P_{H/L}')_{K,L} \\ & + \text{SS}(Q_{G/K} Z P_{H/L}')_{K,L} \\ & + \text{SS}(P_{G/K} Z Q_{H/L}')_{K,L} \\ & + \text{SS}(Q_{G/K} Z Q_{H/L}')_{K,L} \end{aligned} \quad (9)$$

が成り立つ. これら  $Z$  の変動 (計量行列  $K, L$  のもとでの  $Z$  の平方和) が4つの独立な成分に分解されたことを意味する.

$K, L$  の任意の平方根分解を  $K = R_K R_K', L = R_L R_L'$  とする.

$$Z^* = R_K' Z R_L, \quad (10a)$$

$$G^* = R_K' G, \quad (10b)$$

$$H^* = R_L' H \quad (10c)$$

と定義すると,  $\text{SS}(Z)_{K,L} = \text{SS}(Z^*)_{I,I}$  が成り立つ. 従って  $\text{SS}(Z^*)_{I,I}$  を単に  $\text{SS}(Z^*)$  と表すと

$$\begin{aligned} \text{SS}(Z^*) = & \text{SS}(P_G Z^* P_H) \\ & + \text{SS}(Q_G Z^* P_H) \\ & + \text{SS}(P_G Z^* Q_H) \\ & + \text{SS}(Q_G Z^* Q_H) \end{aligned} \quad (11)$$

と表せる. ここで  $P_G, Q_G, P_H, Q_H$  はそれぞれ  $G^*, H^*$  で定義された直交射影行列 (それぞれ  $P_{G/K}, Q_{G/K}, P_{H/L}, Q_{H/L}$  に対応する.) である. これは  $K, L$  を用いた最小2乗問題が単位行列を計量行列とした場合のそれに還元できることを示している (例えば Rao, 1980).

### 3.2. 内部分析

CPCA では分解された  $Z$  の各項を個別にあるいはいくつかをまとめて主成分分析する. どの項を主成分分析するかは研究者の興味や目的による. 例えば (8) の第1項を PCA にかけて,  $Z$  のうち  $G$  と  $H$  の両方によって説明できる部分の中で, 最も重要な変動部分を取り出すことができる. また, 第4項の PCA は誤差の中に重要な情報が残っていないかどうかを調べる誤差分析として重要な意味を持つ (Gabriel, 1978; Rao, 1980; Yanai, 1970).

PCA はデータ行列を階数の低い行列で近似する問題と密接に結びついている. この問題は一意的には解けないが, 1つの解が特異値分解 (SVD と略称する.) によって与えられることは良く知られた事実である (Eckart & Young, 1936). CPCA では一般化最小2乗基準を用いているため, SVD を若干一般化しなければならない. これを一般化特異値分解 (GSVD; 例えば Greenacre, 1984 を参照) と言う.

**定義 (GSVD):**  $\text{rank}(KXL) = \text{rank}(X) = r$  を満たす行列,  $X(N \times n)$ ,  $K(N \times N)$  で  $\text{mnd}$ ,  $L(n \times n)$  で  $\text{mnd}$ , が与えられた時,  $X$  を3つの行列の積,  $UDV'$  で表す分解を計量行列  $K, L$  に基づく  $X$  の GSVD と呼び,  $\text{GSVD}(X, K, L)$  と表す. ここで  $U(N \times r)$  は  $U'KU = I$ ,  $V(n \times r)$  は  $V'LV = I$  を満たす行列,  $D$  は対角行列で  $\text{pd}$  である.

$X$  の普通の SVD, すなわち  $\text{GSVD}(X, I, I)$  を特に  $\text{SVD}(X)$  と表す.  $\text{SVD}(X)$  と  $\text{GSVD}(X, K, L)$  の違いは前者では  $U'U = I, V'V = I$  が要求されるのに対し, 後者では  $U'KU = I, V'LV = I$  が必要とされる点である.

$\text{GSVD}(X, K, L)$  は次の様にして求められる.  $K = R_K R_K', L = R_L R_L'$  として  $X^* = R_K' X R_L$  と定義する.  $X^*$  の SVD を  $X^* = U^* D^* V^{*'}$  と表す. これより  $\text{GSVD}(X, K, L)$  は  $U = (R_K')^* U^*, V = (R_L')^* V^{*'}, D = D^*$  と求められる. ここで  $A^*$  は行列  $A$  のムーア・ペンローズ逆行列を表す. このようにして求められた  $U, V, D$  が GSVD の定義の中で要求される性質を満たすことは明らかである. また, CPCA で例えば (8) の第1項,  $X = P_{G/K} Z P_{H/L}'$  が  $\text{rank}(KXL) = \text{rank}(X)$  を満たしていることは (1) の条件より明らかである.

このように内部分析では GSVD が重要な役割を果たす. GSVD を効率良く求めるための定理は重要である.

定理1:  $T(N \times t \text{ で } t \leq N)$ ,  $W(n \times w \text{ で } w \leq n)$  を  $T'T=I$ ,  $W'W=I$  を満たす準直交行列とする.  $X(t \times w)$  の SVD を  $X=U_X D_X V_X'$ ,  $TXW'$  の SVD を  $TXW'=UDV'$  とすると,  $U=TU_X (U_X=U_X T')$ ,  $V=V_X (V_X=W'V)$ ,  $D=D_X$  が成り立つ. (証明は付録を参照.)

CPCA で  $G^*$ ,  $H^*$  の QR 分解 (例えば Golub & Van Loan, 1989) を  $G^*=F_G R_G'$ ,  $H^*=F_H R_H'$  (ただし  $F_G' F_G=I$ ,  $F_H' F_H=I$ ) とすると,  $P_G$ ,  $P_H$  は  $P_G=F_G F_G'$ ,  $P_H=F_H F_H'$  と書き表される. これより  $R_K' P_{G/K} Z_{P_H/L}' R_L=P_G Z_{P_H}=F_G F_G' Z_{P_H} F_H'$  が得られる.  $F_G, F_H$  は準直交行列だから SVD( $P_G Z_{P_H}$ ) を求めるにはまず SVD( $F_G' Z_{P_H}$ ) を求め, 定理1を適用すればよい.  $F_G' Z_{P_H}$  は通常  $P_G Z_{P_H}$  よりもずっとサイズが小さいので計算の効率化が期待できる. さらに GSVD( $P_{G/K} Z_{P_H/L}', K, L$ ) は GSVD の定義の下で説明した手続きに従って SVD( $P_G Z_{P_H}$ ) から簡単に求められる. 定理1を次の様に一般化すると興味深い結果が導かれる.

定理2: 定理1で  $T, W$  を必ずしも準直交行列でないものとする. ただし  $T'T$  及び  $W'W$  は正則であるものとする. GSVD( $TXW', K, L$ ) を  $UDV'$ , GSVD( $X, T'KT, W'LV$ ) を  $U_X D_X V_X'$  と表すと

$$U=TU_X \quad (U_X=(T'T)^{-1}T'U), \quad (12a)$$

$$V=V_X \quad (V_X=(W'W)^{-1}W'V), \quad (12b)$$

$$D=D_X \quad (12c)$$

の関係が成り立つ. (証明は付録を参照のこと.)

定理2で  $T=G, W=H, X=\hat{M}$  と置くと,  $TXW'GMH'=P_{G/K} Z_{P_H/L}'$ . 従って GSVD( $P_{G/K} Z_{P_H/L}', K, L$ ) と GSVD( $\hat{M}, G'KG, H'LV$ ) の間には定理2で示されているような関係が成り立つ. 後者は  $GMH'$  のうち  $\hat{M}$  にだけ PCA を適用したい場合に必要となる.  $U_X$  と  $U$  は  $T$  にかかる重み(係数)行列とそれによって求められるスコア行列の関係にある. ( $V_X$  と  $V$  も同様.)

### 3.3. CPCA の拡張

以上が CPCA の概要である. これまでの説明から明らかのように CPCA は射影と GSVD の組み合わせから成る. 本節では射影行列の分解に基づいて CPCA の拡張を試みる.

(8) 式における  $Z$  の分解はごく基本的なもので,  $G$

や  $H$  がそれぞれ複数個ある場合には  $Z$  をさらに分解できる. 例えば,  $G=[X \ Y]$  とすると条件に応じて様々な分解が可能である (Rao & Yanai, 1979; 柳井・竹内, 1983). (以下 (1) の条件が成り立つものと仮定する.) まず,  $X$  と  $Y$  が ( $K$  のもとで; 以下同様) 直交する場合は

$$P_{G/K}=P_{X/K}+P_{Y/K} \quad (13)$$

が成り立つ. これは単純に  $G$  の影響を  $X$  の影響と  $Y$  のそれに分解するものである.  $X$  と  $Y$  がそれらが交わる空間以外で直交する場合は  $P_{X/K}$  と  $P_{Y/K}$  が可換で

$$P_{G/K}=P_{X/K}+P_{Y/K}-P_{X/K}P_{Y/K} \quad (14)$$

が成り立つ. この分解で  $K=I$  と置いたものは分散分析で重要な役割を果たす.  $X$  と  $Y$  に何の制限も置かない場合は

$$P_{G/K}=P_{X/K}+P_{Q_{X/K}Y/K} \quad (15a)$$

$$(または =P_{Y/K}+P_{Q_{Y/K}X/K}) \quad (15b)$$

が成り立つ.  $P_{Q_{X/K}Y/K}$  は  $Y$  から  $X$  の影響を取り除いた行列,  $Q_{X/K}Y$  で定義される射影行列である ( $P_{Q_{Y/K}X/K}$  はその逆). この分解はまず  $X$  を当てはめ,  $X$  によって説明できる部分を取り除いてから  $Y$  を当てはめた時に得られる. また  $X$  と  $Y$  が互いに素であれば

$$P_{G/K}=X(X'Q_{Y/K}KX)^{-1}X'Q_{Y/K}K \\ +Y(Y'Q_{X/K}KY)^{-1}Y'Q_{X/K}K \quad (16)$$

が成り立つ. ((7c) より  $Q_{X/K}K$  は対称であることに注意.) これは  $X$  と  $Y$  を同時に当てはめた時に得られる分解である. 右辺の第1項は  $\text{Sp}(Q_{G/K}) \oplus \text{Sp}(Y)$  に沿った  $X$  への射影, 第2項は  $\text{Sp}(Q_{G/K}) \oplus \text{Sp}(X)$  に沿った  $Y$  への射影を表す. ここで  $\oplus$  は部分空間の直和を表す. これらの2つの項はこれまでの分解と違って直交しないので解釈には十分注意が必要である.

付加的な情報が  $G$  にかかる係数行列  $U_G$  に対する制約として与えられる場合は次の分解が有効である. いま制約  $U_G=AU_A$  と表せるものとする,

$$P_{G/K}=P_{G_A/K}+P_{G(G'KG)^{-1}B/K} \quad (17)$$

が成り立つ (Takane, Yanai & Mayekawa, 1991). ここで,  $G'KG$  が正則であることを仮定した.  $B$  は  $\text{Ker}(A')=\text{Sp}(B)$  を満たす行列である. (17) の右辺の第1項は  $G$  によって説明できる部分のうちさらに

$GA$  によって説明できる部分, 第2項は  $GA$  によって説明できない部分を表す.  $B'(G'KG)^{-1}G'KGAU_A = B'U_G = 0$  より  $U_G = AU_A$  の制約は  $B'U_G = 0$  とも表されることがわかる. また (17) で  $A' = [I \ 0]$ ,  $B' = [0 \ I]$  と置くと (15a) が導かれる.

なお (13)-(17) と同様の分解が  $H$  側にも成り立つことは明らかであろう. また以上の分解は  $G$  や  $H$  が3個以上の部分行列から成り立つ場合にも容易に拡張できる.

CPCA に入れ子型 (nested) の高次構造を導入することも可能である. 既に(17)の分解がその一例になっているが, その他に次の様な例があげられる. (2) の  $M$  に  $Z$  と類似した構造

$$M = AM_{AB}B' + M_B B' + AM_A + E^* \quad (18)$$

を仮定し, 説明を簡略化するために  $Z$  には (2) の第1項のみを当てはめるモデル (これを成長モデルともいう—Potthoff & Roy, 1964) を考えると

$$Z = G(AM_{AB}B' + M_B B' + AM_A + E^*)H' + E \quad (19)$$

が得られる. このモデルは  $Z$  を全部で5つの部分に分解する. このモデルは少数の基本変数によって構成された刺激の対比較データなどを分析するのに有効である. (例えば,  $G = A = I$ , 刺激のデザイン行列を  $B$ , 対比較のデザイン行列を  $H$  と置く.)

高次の構造で面白いのは  $Z$ , 或いは分解された  $Z$  の一部に回帰分析によって外部情報を「埋め込む」方法である. この手法は多次元尺度法(MDS), PCA, 対応分析などの結果の解釈のためによく用いられる (Carroll, 1972; Lebart, Morineau & Warwick, 1984). この手法は  $M$  を GSVD で分解した後,  $A$  によって説明できる部分とできな部分に分解するものと解釈できる. これも高次の構造の一例である.

高次の構造は一般に

$$Z = (\prod_i G_i^*) M^* (\prod_j H_j^*) \quad (20)$$

と書き表せる (Takane, 1990). ここで  $G_i^*$ ,  $H_j^*$  は外部データを部分行列とする超行列,  $M^*$  はパラメータ行列を対角要素とするブロック対角行列である. (20) は McDonald (1978) によって提案された共分散構造分析のためのモデル, COSAN に類似しているが, COSAN では共分散を分析する (従って  $Z$  は nnd,  $M^*$  も nnd,  $G_i^*$  と  $H_j^*$  は一致しなければならない.) のに

対し, CPCA ではモデルをデータに直接当てはめるので,  $Z, M$  は一般に矩形行列,  $G_i^*$  と  $H_j^*$  も一致しないなどの違いがある.

### 3.4. 関連手法

CPCA は様々な方法を特別の場合として含む非常に一般的な方法である. 本節では CPCA の関連手法について簡単に触れる.

(1) 制約付き (一般化) 固有値問題を最初に論じたのは Rao (1964; section 2) である. Golub (1973) 及び Böckenholt & Böckenholt (1990) はこれを SVD に拡張した. また Besse & Ramsay (1986) も「関数データの解析」で同様の方法を提案している.

線形の制約はパラメータのゼロ空間を指定する方法と元のパラメータを少数のパラメータで表現し直す方法がある. 上記の方法は前者に属するが, これらの方法と (4) で述べる方法との関係については Takane, Yanai & Mayekawa (1991) で詳しく論じられている.

(2) 従来の PCA は  $Z$  の列ベクトルによって張られる空間の中に「成分」を構成する. Guttman (1944) の方法, Rao (1964, section 11) の方法もこれに属する. これらの方法は CPCA では  $H$  のみを用いる場合 (特に Rao の場合は  $H = Z'G$  と置く.) に相当する.

(3)  $H = I$  の時, CPCA は冗長性分析 (redundancy analysis; Van den Wollenberg, 1977) と呼ばれる方法に一致する. この方法は階数制約付き回帰分析 (reduced rank regression; Anderson, 1951), 補助変数の PCA (PCA of instrumental variables; Rao, 1964, section 8) などとも呼ばれる. この方法に関連した文献は数多いが (Van der Leeden, 1990), この方法の効用を正準相関分析との関係で論じたものに Lambert, Wildt & Durand (1988) がある.

(4) CPCA は特殊な場合正準相関分析に一致する (Takane, 1991). これには次の3つの場合が考えられる. 1)  $Z = I$ , 2)  $H = I$ ,  $Z'Z = I$ , 3)  $H = I$ ,  $L = (Z'Z)^{-1}$ . この場合, CPCA は当然正準相関分析の特別の場合に当たる重判別分析, MANOVA, 対応分析 (双対尺度法, 数量化3類) などとも一致する. また, 正準相関分析で付加的な制約が付く場合の方法が Yanai & Takane (1991) によって, 対応分析でも同様の方法が ter Braak (1986) によって提案されている. 数量化2類 (Hayashi, 1952) もこの部類に属する (Takane, Yanai, & Mayekawa, 1991).

(5) (4) とも関連するが, 西里は一連の研究 (Ni-

shisato, 1980) で多肢選択データ, 対比較データ, 評定データなど様々なタイプのデータによって要求される特殊な制約を組み込んだ双対尺度法を提案している。

(6) CPCA の基本モデル(第2式)で第1項のみを当てはめるモデルに CANDELINC (Carroll, Pruzansky, Kruskal, 1980) と成長モデル (Potthof & Roy, 1964) がある。Takane & Shibayama (1991) は前者と CPCA の関連について, また Takane (1991) は後者と CPCA の関連について論じている。

(7) CPCA は BTC モデル (Bechtel, Tucker & Chang, 1971), THL モデル (Heiser & de Leeuw, 1981; Takane, 1980, 1987), WVM モデル (De Soete & Carroll, 1983) などのベクトル選好モデルを特別の場合として含む (Takane & Shibayama, 1991)。次節の例はベクトル選好モデルを用いた例である。

(8) 本論文で定義された GSVD の用法は主としてフランスのデータ解析学派によるものである (例えば Escoufier, 1987; Greenacre, 1984; Ramsay, ten Berge & Styau, 1984)。Bojanczyk, Ewerbring, Luk & Van Dooren (1991) は同じ分解を HK-SVD と呼んでいる。

一方, 同じ GSVD という名称が数値解析の分野 (例えば Golub & Van Loan, 1989) では別の意味で用いられてきた (Van Loan, 1976)。数値解析でいう GSVD (最近これを QSVD (Quotient SVD) と呼ぼうという提案が数値解析の専門家 (De Moor & Golub, 1991) の間でなされている。) は一般化固有値問題 ( $A'A - sB'B)x=0$  を行列の積を取らず  $A$  と  $B$  のままで計算する方法である。この方法は最近 RSVD (Restricted SVD) と呼ばれる3つの行列を同時に考慮する分解 (CPCA では  $Z$  と  $G$  と  $H$  に相当。) に拡張された (Zha, 1991)。RSVD は計算法に重点を置いているが概念的には CPCA と非常に良く似ている。

#### 4. 適用例

本節では CPCA の適用例を紹介する。紙面の制約からここでは1例しか紹介できないが, 他の例については Takane (1990) を参照されたい。

ここで用いるデータは9つの刺激に対する対比較選好データである。刺激は3つの異なった領域から選ばれた9人の有名人 (1. Brian Mulroney (カナダの首相) 2. Ronald Reagan (アメリカ合衆国の大統領; 当時) 3. Margaret Thatcher (イギリスの首相; 当時)

4. Jacqueline Gareau (女子マラソン選手) 5. Wayne Gretzky (プロのアイスホッケー選手) 6. Steve Podborski (プロのスキー選手) 7. Paul Anka (歌手) 8. Tommy Hunter (カントリー・ウェスタン歌手) 9. Ann Murray (歌手)) である。被験者には刺激を対にして提示し, どちらがどの位好ましいかを25点尺度で評定してもらった。被験者は100人のカナダ人学生でそのうち約半数の学生が英語系, 残りが仏語系, その他であった。

このデータに  $G=I, H=A, K=I, L=I$  として CPCA を適用した。ここで  $A$  は対比較のデザイン行列である (これは前節3.4, (7) で述べた BTC モデルに一致する。) この分析ではまず  $ZA(A'A)^+ = ZA/n$  を求め, この行列に PCA を適用する。図1はこの結果得られた二次元の刺激布置を表している。

図中1から9までの番号が9つの刺激の位置を示している。これより9つの刺激はおおよそ刺激が選ばれた3つの異なった分野に分かれて空間内に配置されていることがわかる。(例外は4の JG でスポーツ選手でありながら歌手のグループに近い位置を占めている。) ラベルのない10個の矢印は100人のうち最初の10人の被験者の選好ベクトルを表す。個人の選好ベクトルに空間内に配置された刺激を射影するとその個人の刺激に対する選好価が得られる。

9つの刺激点を取り巻く楕円はブートストラップ法 (Efron, 1979) によって求められた95% 信頼領域を表している。これは元のデータから繰り返し100のサンプルを抽出し, それを個々に CPCA で分析して推定値の分散・共分散行列を求め, 推定量の漸近正規性を仮定して描いたものである。

図中 V の印が付いたベクトルは被験者全体の平均的傾向を表す。これより平均的に WG, SP (共にスポーツ選手) が一番好まれ, 次に政治家, 最後に歌手 (JG はここでも例外) の順で好まれていることがわかる。また, E と N の印が付いたベクトルはそれぞれ英語系学生, 仏語系学生その他の平均選好ベクトルを表す。E の方が若干政治家好みの傾向を示している。(政治家は3人とも英語系である。) しかし, これら2つのベクトルの信頼領域が重なり合っていることから2つのグループの間には統計的有意差はないと考えられる。

#### 5. 考 察

以上, CPCA の概要をデータの要件, 解析法, 応用例を通して概観した。本節ではまとめの意味で CPCA

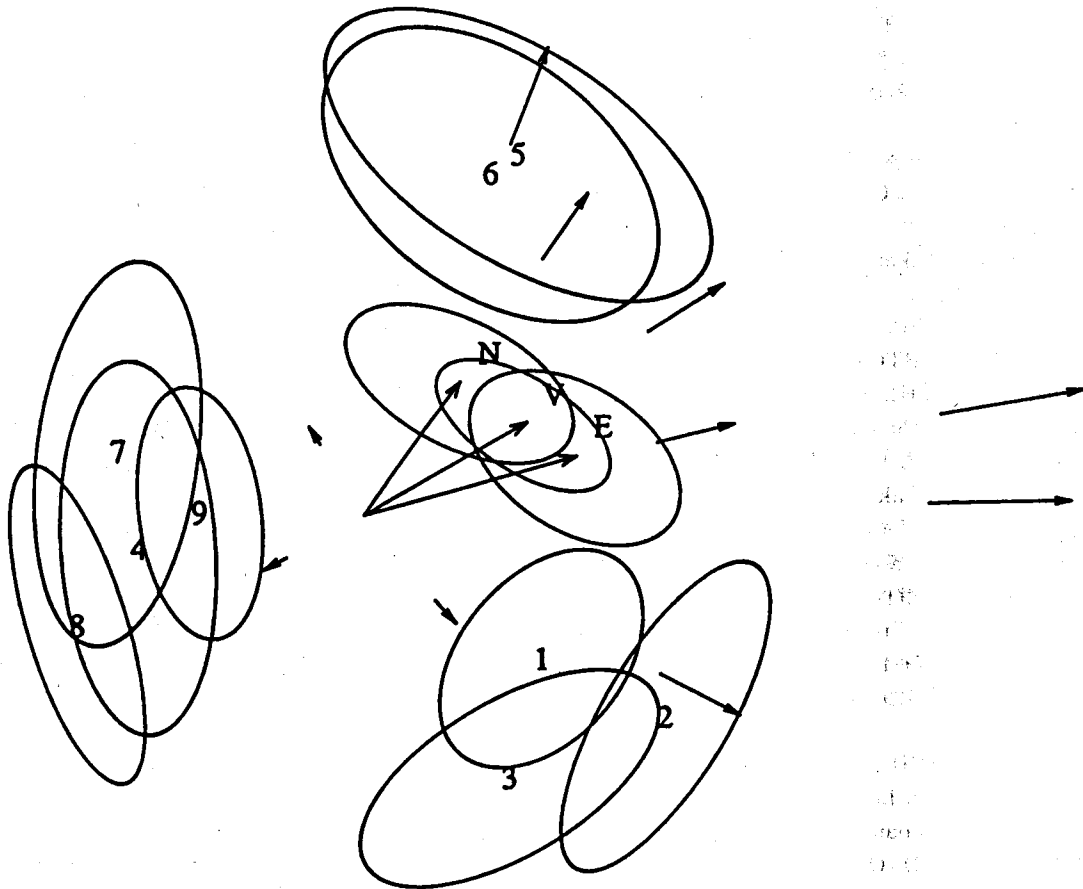


図1. BTCモデルから導かれた選好データの刺激布置と95%信頼領域

の特徴を共分散構造分析法 (ACOVs; Bock & Bargman, 1966; Jöreskog, 1970) との関連で述べてみたい。CPCA の利点としては次の様な点があげられる。

1. 不自然な分布の仮定を置かない。既成のACOVs のほとんどが多変量正規分布の仮定に基づいている。この仮定を楕円分布に一般化する試みがあるが、この分布は正規分布と大差ない。またADF (asymptotically distribution free) 法は膨大なサンプルサイズを必要とするため余り実用的とは言えない。

2. 計算が簡単である。CPCA は射影とGSVDの組み合わせから成る。いずれも効率の良い計算法が存在する (Golub & Van Loan, 1989)。

3. 解析的に解が求まる。ACOVs が一部の例外を除いては殆ど逐次解を採用しているのに対し、CPCA では解析的に解が求められる。逐次解では非最適解へ

の収束の危険性が常にあるが、CPCA ではそのような心配がない。

4. 不適解が生じない。ACOVs では分散・共分散行列の推定値が  $n \geq m$  でないことがあるが、CPCA ではそのような問題は生じない。最尤解ではさらに分散・共分散行列が  $pd$  であることが要求されるが、 $pd$  の制約をかけるのは  $n \geq m$  の制約をかけるよりもずっと困難である。

5. 成分得点が一義的に得られる。CPCA では因子得点の推定にまつわる因子不定性の問題が生じない。

6. 次元数や増やしても低次元解が保持される。従って次元数の選択にそれ程神経を使わないで済む。

この他、Velicer & Jackson (1990) は因子分析と比較したPCAの効用を実証的研究に基づいて議論している。



一方、CPCAの欠点とでも言うべき点は

1. CPCAは記述的に過ぎない。無理な分布の仮定を置かない反面、モデル評価の基準がないため統計的推測ができない。この議論は半分しか正しくない。既に応用例で見たように、CPCAによって得られた結果の信頼性はブートストラップ法 (Efron, 1979) などによって調べることができる。またブートストラップ法によって得られたパラメータの推定値の分散・共分散行列の推定値を用いて、一部の仮説検定を行うことも可能である。モデル評価の方法としてはストーン・ガイサーの方法 (Geisser, 1975; Stone, 1974) なども有望である。次元数を統計的に決めたい場合は交差妥当性に基づいた Eastment & Krzanowski (1982) の方法なども使える。

2. CPCAは測定誤差を考慮しない。この批判も半分しか正しくない。内部分析で小さい特異値に対応する成分を取り除くことは測定誤差を選り分ける働きを持つ (Gleason & Staelin, 1973)。変数の共分散のみを解析目的とするならば問題はないが、因子分析などで測定誤差と同時に特殊因子の効果も取り除いてしまうことの方が問題である。

3. CPCAは尺度の不変性がない。最尤法、一般化最小2乗法に基づいた ACOVS は入力変数の単位を変えても結果は不変であるのに対し、CPCAには単位の不変性がない。これも半分しか正しくない。既に見たように、計量行列  $L$  を導入すれば CPCA でも単位の不変性が得られる。

4. CPCAは柔軟性に欠ける。CPCAは確かに ACOVS に比べると扱い得る制約の種類などで劣っている。しかし、CPCAの他の利点はこの欠点を補って余りあると考えられる。

本論文では紙面の制約上すべてを語り尽くすことができなかった (例えば欠測値やはずれ値の問題など)。にもかかわらず CPCA の本質的な部分は紹介できたものと信じる。語り残した部分に関しては成書を準備中である。

## 付 録

### 定理1の証明.

$X=U_X D_X V_X'$  の両辺に前と後ろから  $T, W'$  をかけると  $TXW'=TU_X D_X V_X' W'$  が得られる。  $U=U_X D_X V_X'$  と置くと  $TXW'=UDV'$  が得られる。あとは  $U, D, V$  が SVD の構成要素が満

たすべき条件 ( $U'U=I, V'V=I, D$  は対角で pd) を満たすことを示せばよい。  $T$  は準備直交行列、  $U_X$  は  $X$  の左特異ベクトル行列だから  $U'U=U_X' T' T U_X=I$  が成り立つ。同様に  $V'V=V_X' W' W V_X=I$  が成り立つ。  $D$  が対角で pd であることは明らか。

逆に  $TXW'=UDV'$  の両辺に前と後ろから  $T'$  と  $W$  をかけると  $T' TXW' W=X=T' U D V' W$ 。ここで  $U_X=T' U, V_X=W' V, D=D_X$  と置くと  $X=U_X D_X V_X'$  が得られる。あとは  $U_X' U_X=I, V_X' V_X=I$  であることを示せばよい。 ( $D_X$  が対角で pd であることは自明。) 前者は  $P_T U=U$  (ここで  $P_T=TT'$  は  $T$  によって定義される直交射影行列)、後者は  $P_W V=V$  より簡単に示される。

### 定理2の証明.

定理2も定理1と同様にして示される。

$X=U_X D_X V_X'$  の両辺に前後から  $T, W'$  をかけると  $TXW'=TU_X D_X V_X' W'=UDV'$  が得られる。このとき、  $U=U_X D_X V_X'$  が  $U'KU=I, V'LV=I$  を満たし、  $D=D_X$  が対角で pd であることを示せばよい。実際

$$\begin{aligned} U'KU &= U_X' T' K T U_X = I, \\ V'LV &= V_X' W' L W V_X = I \end{aligned}$$

が成り立つ。 ( $D$  が対角で pd であることは自明。)

逆を示すには  $TXW'=UDV'$  の両辺に前後から  $(T'T)^{-1}T', W(W'W)^{-1}$  をかけ、  $U_X=(T'T)^{-1}T' U, V_X=(W'W)^{-1}W' V$  が  $U_X' T' K T U_X=I, V_X' W' L W V_X=I$  を満たすことを示せばよい。 ( $D_X$  が対角で pd であることは自明。) これらは  $P_T U=T(T'T)^{-1}T' U=U, P_W V=W(W'W)^{-1}W' V=V$  より簡単に示される。

## 参 考 文 献

- Anderson, T.W. (1951) Estimating linear restrictions on regression coefficients for multivariate normal distribution. *Annals of Mathematical Statistics*, 22, 327-351.
- Bechtel, G.G., Tucker, L.R., & Chang, W. (1971) A scalar product model for the multidimensional scaling of choice. *Psychometrika*, 36, 369-387.
- Besse, P., & Ramasy, J.O. (1986) Principal component analysis of sampled functions. *Psychometrika*, 51, 285-311.
- Bock, R.D., & Bargman, R.E. (1966) Analysis of covariance structures. *Psychometrika*, 31, 507-

- 534.
- Böckenholt, U., & Böckenholt, I. (1990) Canonical analysis of contingency tables with linear constraints. *Psychometrika*, 55, 633-639.
- Bojanczyk, A.W., Ewerbring, M., Luk, F.T., & Van Dooren, P. (1991) An accurate product SVD algorithm. In R.J. Vaccaro (Ed.), *SVD and signal processing II*, (pp.113-131). Amsterdam: Elsevier.
- Carroll, J.D. (1972) Individual differences and multidimensional scaling. In R.N. Shepard, A. K. Romney & S.B. Nerlove (Eds.), *Multidimensional scaling, Vol. I* (pp.105-155). New York: Seminar Press.
- Carroll, J.D., Pruzansky, S., & Kruskal, J.B. (1980) CANDELINC: a general multidimensional analysis of many-way arrays with linear constraints on parameters. *Psychometrika*, 45, 3-24.
- De Moor, B.L.R., & Golub, G.H. (1991) The restricted singular value decomposition: properties and applications. *SIAM Journal: Matrix Analysis and Applications*, 12 401-425.
- De Soete, G., & Garroll, J.D. (1983) A maximum likelihood method for fitting the wandering vector model. *Psychometrika*, 48, 553-566.
- Eastment, H.T., & Krzanowski, W.J. (1982) Cross-validatory choice of the number of components from a principal component analysis. *Technometrics*, 24, 73-77.
- Eckart, C., & Young, G. (1936) The approximation of one matrix by another of lower rank. *Psychometrika*, 1, 211-218.
- Efron, B. (1979) Bootstrap methods: another look at the Jackknife. *Annals of Statistics*, 1, 1-26.
- Escoufier, Y. (1987) The duality diagram: a means for better practical applications. In P. Legendre & L. Legendre (Eds.), *Development in numerical ecology* (pp.139-156). Berlin: Springer.
- Gabriel, K.R. (1978) Least squares approximation of matrices by additive and multiplicative models. *Journal of Royal Statistical Society, Series B*, 40, 186-196.
- Geisser, S. (1975) The predictive sample reuse method with applications. *Journal of the American Statistical Association*, 70, 320-328.
- Gleason, T.C., & Staelin, R. (1973) Improving the metric quality of questionnaire data. *Psychometrika*, 38, 393-410.
- Golub, G.H. (1973) Some modified eigenvalue problems. *SIAM Journal: Review*, 15, 318-335.
- Golub, G.H., & Van Loan, C.F. (1989) *Matrix computations*. (2nd ed.) Baltimore: Johns Hopkins University Press.
- Greenacre, M.J. (1984) *Theory and applications of correspondence analysis*. London: Academic Press.
- Guttman, L. (1944) Generalized theory and methods for matrix factoring. *Psychometrika*, 9, 1-16.
- Guttman, L. (1953) Image theory for the structure of quantitative variables. *Psychometrika*, 18, 277-296.
- Hayashi, C. (1952) On the prediction of phenomena from qualitative data and the quantification of qualitative data from mathematico-statistical point of view. *Annals of the Institute of Statistical Mathematics*, 3, 69-98.
- Heiser, W., & de Leeuw, J. (1981) Multidimensional mapping of preference data. *Mathématique et sciences humaines*, 19, 39-96.
- Jöreskog, K.G. (1970) A general method for analysis of covariance structures. *Biometrika*, 57, 239-251.
- Kiers, H.A.L. (1991) Simple structure in component analysis techniques for mixtures of qualitative and quantitative variables. *Psychometrika*, 56, 197-212.
- Lambert, Z.V., Wildt, A.R., & Durand, R.M. (1988) Redundancy analysis: an alternative to canonical correlation and multivariate multiple regression in exploring interest association. *Psychological Bulletin*, 104, 282-289.
- McDonald, R.P. (1978) A simple comprehensive model for the analysis of covariance structures. *British Journal of Mathematical and Statistical Psychology*, 31, 59-72.
- Meredith, W., & Millsap, R.E. (1985) On component analysis. *Psychometrika*, 50, 495-507.
- Nishisato, S. (1980) *Analysis of categorical data: dual scaling and its applications*. Toronto: University of Toronto Press.
- Potthof, R.F., & Roy, S.N. (1964) A generalized multivariate analysis of variance model useful especially for growth curve problems. *Biometrika*, 51, 313-326.
- Ramsay, J.O., ten Berge, J., & Stylian, G.P.H. (1984) Matrix correlation. *Psychometrika*, 49, 403-423.
- Rao, C.R. (1964) The use and interpretation of principal component analysis in applied research. *Sankhya A*, 26, 329-358.
- Rao, C.R. (1980) Matrix approximation and reduction of dimensionality in multivariate statistical analysis. In P.R. Krishnaiah (Ed.), *Multivariate analysis V*. Amsterdam: North Holland.
- Rao, C.R., & Yanai, H. (1979) General definition and decomposition of projectors and some

- applications to statistical problems. *Journal of Statistical Inference and Planning*, **3**, 1-17.
- Stone, M. (1974) Cross-validators choice and assessment of statistical predictions (with discussion). *Journal of the Royal Statistical Society, Series B*, **36**, 111-147.
- Takane, Y. (1980) Maximum likelihood estimation in the generalized case of Thurstone's model of comparative judgment, *Japanese Psychological Research*, **22**, 188-196.
- Takane, Y. (1987) Analysis of covariance structures and binary choice data, *Communication and cognition*, 45-62.
- Takane, Y. (1990) Constrained principal component analysis and its applications. Paper submitted for publication.
- Takane, Y. (1991) Principal component analysis with linear constraints on both rows and columns of a data matrix. Paper submitted for publication.
- Takane, Y., & Shibayama, T. (1991) Principal component analysis with external information on both subjects and variables. *Psychometrika*, **56**, 97-120.
- Takane, Y., Yanai, H., & Mayekawa, S. (1991) Relationships among several methods of linearly constrained correspondence analysis. *Psychometrika*, **56**, 667-684.
- Ter Braak, C.J.F. (1986) Canonical correspondence analysis: a new eigenvector technique for multivariate direct gradient analysis. *Ecology*, **67**, 1167-1179.
- Van den Wollenberg, A.L. (1977) Redundancy analysis: an alternative for canonical analysis. *Psychometrika*, **42**, 207-219.
- Van der Leeden, R. (1990) Reduced rank regression with structured residuals. Leiden, The Netherlands: DSWO Press.
- Van Loan, C.F. (1976) Generalizing the singular value decomposition. *SIAM Journal: Numerical Analysis*, **13**, 76-83.
- Velicer, W.F., & Jackson, D.N. (1990) Component analysis versus common factor analysis: some issues in selecting an appropriate procedure. *Multivariate Behavioral Research*, **25**, 1-28.
- Yanai, H. (1970) Factor analysis with external criteria. *Japanese Psychological Research*, **12**, 143-153.
- Yanai, H., & Takane, Y. (1991) Canonical correlation analysis with linear constraints. To appear in *Linear Algebra and Applications*.
- 柳井晴夫・竹内 啓 (1983) 「射影行列・一般逆行列・特異値分解」東京：東京大学出版会。
- Zha, H. (1991) The restricted singular value decomposition of matrix triplets. *SIAM Journal: Matrix Analysis and Applications*, **12**, 172-194.