

Objective general index and its applications

東京大学・情報理工学系研究科 清 智也

Tomonari Sei

School of Information Science and Technology,
The University of Tokyo.

概要

各変量に順序（優劣）が定義されているような多変量データを，一変量の総合指数にまとめることを考える．例えば，各変量の Zスコアの和は自然な総合指数と考えられるが，データによっては結果として得られる総合指数と変量の間負の相関が生ずることがある．本稿では，Sei (2016, J. Multivariate Anal.) にならい，全ての変量と正の相関を持つような総合指数を構成する．応用として，女子七種競技の総合点を考察する．また，カテゴリカル変数の数量化，及びその極限についても議論する．

キーワード 数量化，相関行列，総合指数，ランキング，レイティング．

1 総合指数

p 変量データ $X = (x_1, \dots, x_p) \in \mathbb{R}^{n \times p}$ を考える．各 x_i は n 次元ベクトルである．

各変量には「大きいほど良い」という順序が与えられているものと仮定する．例えば n 人の学生の p 科目に関する成績データを想像すればよい．また，陸上の十種競技データでは，走り幅跳びの記録など「大きいほど良い」種目のほか，100m走の記録など「小さいほど良い」種目が混ざっているが，後者の成績にはあらかじめマイナスを掛けるものと約束する．また，各変量はあらかじめ中心化されているものとする．すなわち， $\mathbf{1}_n = (1, \dots, 1)^\top$ として， $\mathbf{1}_n^\top x_i = 0$ と仮定する．

定義 1. $g = Xw = \sum_{i=1}^p w_i x_i$ と定義される g を一般に X の総合指数 (general index) と呼ぶ．ここで $w \in \mathbb{R}^p$ は X の共分散行列 $S = n^{-1} X^\top X$ に依存してよいものとする： $w = w(S)$ ．

総合指数は，[9]の用語で言えばレイティング (rating) に相当する．一方，得られた総合指数をもとにして，個体の優劣を決定することをランキング (ranking) という．関連する研究として [1, 4, 5] などがある．

以下，ベクトル a, b に対して $a > b$ とは成分ごとに $a_i > b_i$ が成り立つことを表す．

定義 2. $w > 0$ である総合指数は**整合的**， $Sw = (n^{-1} x_i^\top g)_{i=1}^p > 0$ となる総合指数は**共分散整合的**ということにする．

例 1 (Z -スコアの和). $w_i = 1/\sqrt{S_{ii}}$ とする. これは一見自然だが $p \geq 3$ のとき共分散整合的でない. 例えば

$$S = \begin{pmatrix} 1 & a & \cdots & a \\ a & 1 & & \\ \vdots & & \ddots & \\ a & & & 1 \end{pmatrix} \quad (1)$$

という行列を考えると, S は $|a| < 1/\sqrt{p-1}$ のとき正定値となるが,

$$-\frac{1}{\sqrt{p-1}} < a < -\frac{1}{p-1}$$

のとき共分散整合的でない. 実際, Sw の第 1 成分が $1 + (p-1)a < 0$ となる.

例 2 (第一主成分). w を S の第一固有ベクトルとする. これは整合的でも共分散整合的でもない. 例えば

$$S = \begin{pmatrix} 1 & -1/2 \\ -1/2 & 1 \end{pmatrix}$$

の場合, 最大固有値 $3/2$ に対応する固有ベクトルは $w = (1, -1)$ の定数倍となり, 必ず負の成分を持つ. そして $Sw = (3/2)w$ も負の成分を持つ. 整合的でないということは, 次のような不公平が生じ得る: 「AさんはBさんより全ての科目で点数が上回っているのに, 総合成績ではBさんの方が上となる。」

これらの欠点を克服したい. そのために次の定理を利用する.

定理 1 ([10]). $S = (S_{ij})$ が正定値対称行列ならばある正のベクトル $w = (w_i)$ が存在し,

$$w_i \sum_{j=1}^p S_{ij} w_j = 1, \quad i = 1, \dots, p, \quad (2)$$

を満たす. またそのような w は一意である.

Proof. $w = (w_1, \dots, w_p)^\top \in \mathbb{R}_+^p$, とおき,

$$\psi(w) = \frac{1}{2} w^\top S w - \sum_{i=1}^p \log w_i \quad (3)$$

とおく. この ψ は狭義凸関数である. また, 各 $a \in \mathbb{R}$ に対して $M_a := \{w \mid \psi(w) \leq a\}$ はコンパクト集合となる. 実際, S の最小固有値を $\rho > 0$ とすれば $w^\top S w \geq \sum_i \rho w_i^2$ より, $M_a \subset \{w \mid \sum_{i=1}^p (\rho w_i^2 / 2 - \log w_i) \leq a\}$ となり, M_a は有界であることが分かる. よって最小点が一意に存在する. その停留条件が式 (2) となる. \square

例 3. 式 (1) の S に対して,

$$\mathbf{w} = (x, y, \dots, y)^\top, \quad x > 0, \quad y > 0,$$

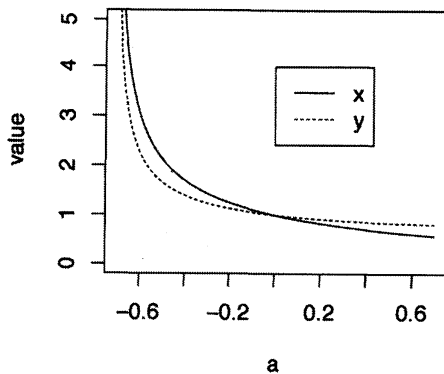
とにおいて, 式 (2) を満たす x, y を求めてみる. 式 (2) は

$$\begin{aligned} x^2 + (p-1)axy - 1 &= 0, \\ axy + y^2 - 1 &= 0 \end{aligned}$$

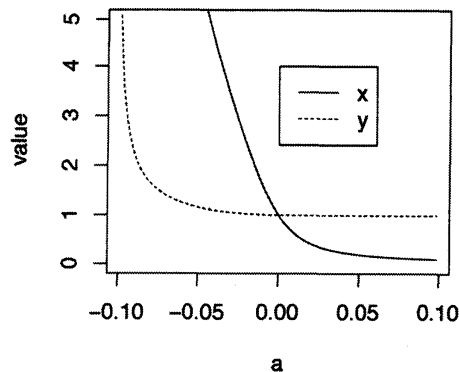
という連立代数方程式となる. 初等的な計算, あるいはグレブナー基底 ([7] など) の援用により, x と y がともに正となる解は

$$\begin{aligned} y &= \left\{ \frac{1}{2(1-(p-1)a^2)} \left(2 - (p-2)a^2 - a\sqrt{4 + (p-2)^2 a^2} \right) \right\}^{1/2}, \\ x &= ((p-1)a - 1/a)y^3 + (-(p-2)a + 1/a)y \end{aligned}$$

であることが示される. a を変化させたときの解の変化の様子を図 1 に示す.



(a) $p = 3$ のとき.



(b) $p = 100$ のとき.

図 1: パラメータ a を変化させたときの重み (x, y) の変化.

定義 3 ([13]). データ行列 X の共分散行列 S に対し, 定理 1 から定まる \mathbf{w} を用いて $\mathbf{g} = X\mathbf{w}$ とおいたものを **客観的综合指数** (Objective General Index; **OGI**) と呼ぶ. またこのときの \mathbf{w} を **客観重み** という.

OGI は整合的かつ共分散整合的である. すなわち, \mathbf{w} の各成分は正であり, 式 (2) から「各変量と総合指数のあいだの相関」も正である: $(1/n)\mathbf{x}_i^\top \mathbf{g} = 1/w_i > 0$. OGI の適用例は 2 節に示す. また客観重みの数値的な求め方は 5 節に記す.

各総合指数の幾何学的イメージを図 2 に示す。それぞれの図において、変量は 3 つあり、総合指数 $g \in \mathbb{R}^n$ が鉛直下向き（破線）に描かれている。各変量 x_i は実線で表され、スケールされた変量 $y_i = w_i x_i$ の位置は丸で記されている。この図は、原点を支点とし、位置 y_i に錘が取り付けられた「やじろべえ」をイメージすると理解しやすい。OGI は錘の位置が水平になるように調整されていることが分かる。実際、OGI に対しては

$$\frac{1}{n} \mathbf{y}_i^\top \mathbf{g} = 1, \quad i = 1, \dots, p,$$

$$\mathbf{g} = \sum_{i=1}^p \mathbf{y}_i,$$

が成り立つ。

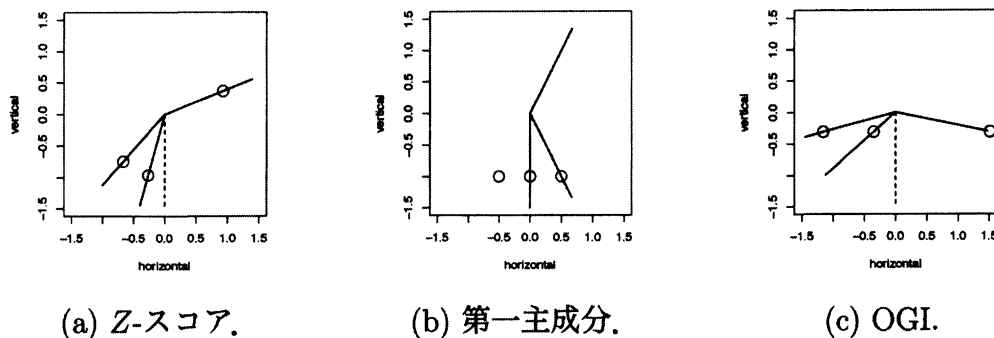


図 2: 各総合指数のイメージ。空間 \mathbb{R}^n において総合指数 g を鉛直下方向に描いた図。

2 適用例

世界陸上女子七種競技の 1991 年から 2013 年までのデータを扱う。女子七種競技は、100m ハードル (100mH), 走高跳 (HJ), 砲丸投 (SP), 200m (200m), 走幅跳 (LJ), やり投げ (Javelin), 800m (800m) の七種目からなり、その総合点を競う競技である。データは IAAF (International Association of Athletics Federations) の公式 web サイト¹ から取得可能である。1991 年から 2013 年までの競技者の人数は延べ 260 人であったが、散布図から明確に外れ値と判断される 2 人のデータを除いてから分析を行った ($n = 258$)。

IAAF のルールによると、各種目の記録はまず非線形変換され、そのあと単純和が取られる [6]。非線形変換の公式を表 1 に示す。実際のデータの一部を表 2 に、変

¹<http://www.iaaf.org/competitions/iaaf-world-championships>

換後のスコアの一部を表3にそれぞれ示す。また変換後の全データの要約値を表4、相関行列を表5に示す。

ここでは、非線形変換した後のスコアを元データと見なし、線形結合による総合指数を検討する。IAAFの総合点は単純和であるから、各種目の重みは1ということになる。Z-スコアとOGIの重み $\{w_i\}_{i=1}^7$ を計算した結果が図3である。エラーバーはブートストラップ法による標準誤差を表す。

男子十種競技に関する先行研究 [3] では、Z-スコアによる考察から、「男子十種競技ではフィールド競技がトラック競技に比べて優遇されている」と報告されている。図3(a)より、この傾向は女子七種にも見られ、特に重みをより大きくすべきなのは100mHということになる。しかし、OGIの観点からは、重みをより大きくすべきなのは800mであると考えられる(図3(b))。ただし、結果として得られるランキングについては、OGIとIAAFルールで大きな違いはない(図3(d))。

表 1: IAAF で決められた非線形変換 $\tilde{x}_i = a_i |d_i x_i - b_i|^{c_i}$ の各パラメータ。

	100mH	HJ	SP	200m	LJ	Javelin	800m
a_i	9.23076	1.84523	56.0211	4.99087	0.188807	15.9803	0.11193
b_i	26.70	75.00	1.50	42.50	210.00	3.80	254.00
c_i	1.835	1.348	1.05	1.81	1.41	1.04	1.88
d_i	1	100	1	1	100	1	1

表 2: 七種競技データの一部。

	100mH (秒)	HJ (m)	SP (m)	200m (秒)	LJ (m)	Javelin (m)	800m (秒)	総合点
1	13.32	1.91	13.62	24.49	6.67	48.66	136.09	6672
2	13.02	1.76	13.52	23.98	6.54	43.58	131.48	6493
3	13.70	1.85	13.23	24.13	6.33	40.96	125.23	6448
4	13.54	1.88	12.46	24.20	6.18	47.04	133.24	6404
5	13.34	1.79	15.40	24.32	6.15	44.62	137.17	6392
6	13.36	1.76	14.36	24.05	6.19	44.28	131.72	6391

表 3: 非線形変換後のスコアの一部. 単純和をとると総合点になる.

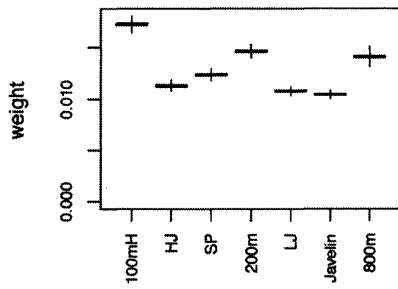
	100mH	HJ	SP	200m	LJ	Javelin	800m	総合点
1	1077	1119	769	934	1062	834	877	6672
2	1121	928	762	983	1020	736	943	6493
3	1021	1041	743	968	953	686	1036	6448
4	1044	1080	692	962	905	803	918	6404
5	1074	966	888	950	896	756	862	6392
6	1071	928	818	976	908	750	940	6391

表 4: 七種競技データの要約値. データ数は $n = 258$.

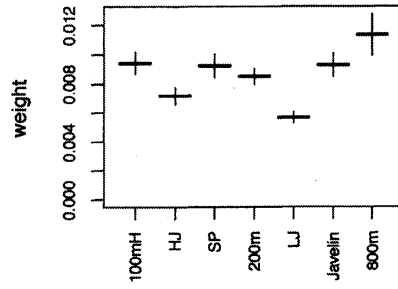
	100mH	HJ	SP	200m	LJ	Javelin	800m	総合点
最小値	783.0	678.0	553.0	751.0	680.0	482.0	607.0	4983
第1分位点	974.5	891.0	716.0	863.2	822.8	679.5	850.5	5938
中央値	1010.0	953.0	755.5	904.5	883.0	740.5	883.5	6134
平均値	1010.0	946.5	763.2	909.0	891.1	745.5	882.4	6148
第3分位点	1048.5	1003.0	811.8	959.8	955.2	809.8	928.0	6341
最大値	1147.0	1171.0	997.0	1095.0	1186.0	984.0	1036.0	7032
標準偏差	58.1	88.0	80.8	68.4	92.4	95.7	70.9	321

表 5: 七種競技データの相関行列.

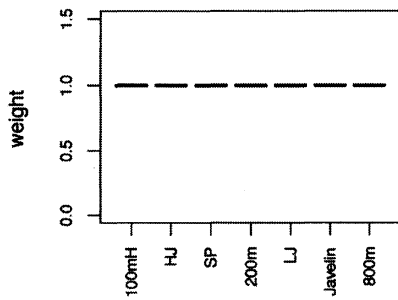
	100mH	HJ	SP	200m	LJ	Javelin	800m
100mH	1.00	0.30	0.14	0.66	0.58	0.11	0.25
HJ	0.30	1.00	0.20	0.25	0.48	0.12	0.17
SP	0.14	0.20	1.00	0.12	0.22	0.23	0.00
200m	0.66	0.25	0.12	1.00	0.57	-0.04	0.33
LJ	0.58	0.48	0.22	0.57	1.00	0.11	0.21
Javelin	0.11	0.12	0.23	-0.04	0.11	1.00	-0.13
800m	0.25	0.17	0.00	0.33	0.21	-0.13	1.00



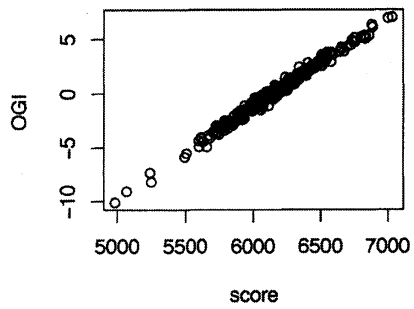
(a) Z-スコアの重み.



(b) OGI の重み.



(c) IAAF ルールの重み.



(d) OGI v.s. IAAF ルールの総合点.

図 3: 世界陸上女子七種競技の分析.

3 OGI数量化

順序付きカテゴリカル変数の場合もOGIは定義される。ここでは簡単のため、1変数の場合を考える。多変数の場合は[13]を参照されたい。また関連研究として[2, 12]がある。

いま確率変数 X は $\{0, 1, \dots, K\}$ に値をとると仮定し、次のダミー変数を用意する：

$$h_k(x) = 1_{\{x \geq k\}} - P(X \geq k), \quad k = 1, \dots, K.$$

h_k の凸結合は階段型の単調増加関数となる。

$h_1(X), \dots, h_K(X)$ に対して、その共分散行列

$$\begin{aligned} S_{kl} &= E[h_k(X)h_l(X)] \\ &= P(X \geq \max(k, l)) - P(X \geq k)P(X \geq l) \\ &= P(X < \min(k, l)) - P(X < k)P(X < l) \end{aligned}$$

から客観重み w を求めれば、

$$y(x) = \sum_{k=1}^K w_k h_k(x)$$

が数量化を与える。ただし、式(2)を

$$w_k \sum_{l=1}^K S_{kl} w_l = \frac{1}{K} \quad (4)$$

に変更する²。このようにして得られる $y(x)$ を x のスコアと呼ぶことにしよう。

例として、次の人工的な成績データを考える：

評価 (x)	優 (3)	良 (2)	可 (1)	不可 (0)	合計
人数	9	13	8	3	33

このデータの場合、 $K = 3$ であり、

$$P(X = 0) = \frac{3}{33}, \quad P(X = 1) = \frac{8}{33}, \quad P(X = 2) = \frac{13}{33}, \quad P(X = 3) = \frac{9}{33}$$

となる。 $h_1(X), h_2(X), h_3(X)$ の共分散行列は

$$\begin{aligned} S &= \begin{pmatrix} \frac{3}{33} - \left(\frac{3}{33}\right)^2 & \frac{3}{33} - \left(\frac{3}{33}\right)\left(\frac{11}{33}\right) & \frac{3}{33} - \left(\frac{3}{33}\right)\left(\frac{24}{33}\right) \\ \frac{3}{33} - \left(\frac{11}{33}\right)\left(\frac{3}{33}\right) & \frac{11}{33} - \left(\frac{11}{33}\right)^2 & \frac{11}{33} - \left(\frac{11}{33}\right)\left(\frac{24}{33}\right) \\ \frac{3}{33} - \left(\frac{24}{33}\right)\left(\frac{3}{33}\right) & \frac{11}{33} - \left(\frac{24}{33}\right)\left(\frac{11}{33}\right) & \frac{24}{33} - \left(\frac{24}{33}\right)^2 \end{pmatrix} \\ &= \begin{pmatrix} 0.0826 & 0.0606 & 0.0248 \\ 0.0606 & 0.2222 & 0.0909 \\ 0.0248 & 0.0909 & 0.1983 \end{pmatrix} \end{aligned}$$

²1 変量カテゴリカルデータに限れば、式(2)と式(4)は定数倍の違いしかないが、多変数の場合は本質的な違いが現れる。前述の通り、本稿では多変量カテゴリカルデータは扱わない。

となり，式(4)を満たす客観重み w は数値的に

$$(w_1, w_2, w_3) = (1.590, 0.869, 1.032)$$

と求まる．結果として得られるスコアは

評価	優	良	可	不可
スコア	1.185	0.153	-0.717	-2.306

となる．

4 OGI数量化の極限

1変量の順序付きカテゴリカルデータにおいて，カテゴリ数を n とし，それぞれが観測される確率は $1/n$ として， $n \rightarrow \infty$ の極限を考える．図4は結果として得られるスコア $y(x)$ をプロットしたものである．ただし，カテゴリカル変数 x の値は $\{0, \dots, n-1\}$ でなく n 個の区間 $(\frac{0}{n}, \frac{1}{n}), \dots, (\frac{n-1}{n}, \frac{n}{n})$ で表した．図より， $y(x)$ は標準正規分布の分位点関数 $\Phi^{-1}(x)$ に近づいていく．ここで Φ は標準正規分布の累積分布関数である．

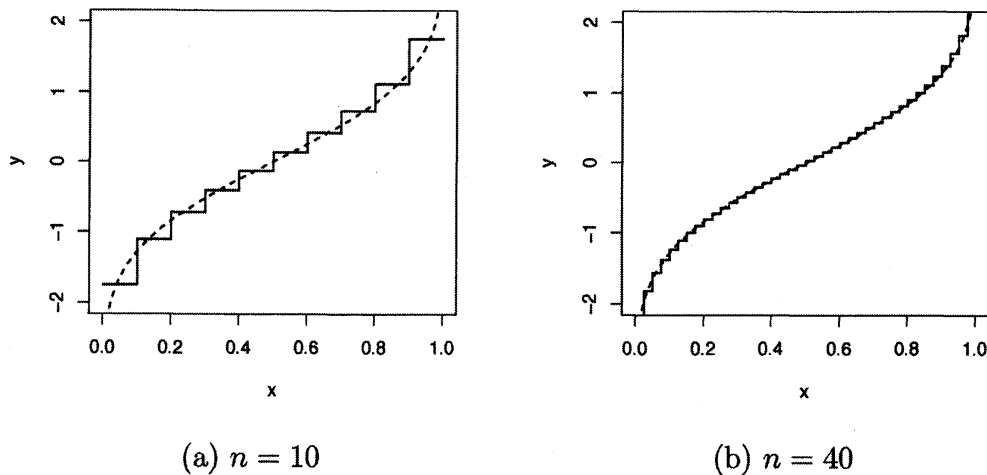


図4: OGI数量化の収束．実線はスコア $y(x)$ ，破線は $\Phi^{-1}(x)$ を表す．

この理由を説明する（厳密な証明ではない）．まず，この場合の共分散行列は

$$S_{ij} = \frac{i}{n} \left(1 - \frac{j}{n}\right), \quad 1 \leq i \leq j \leq n-1,$$

で与えられる。ここで、標準正規分布の密度関数を $\phi(z)$ として、重み w_i を

$$w_i = \frac{1}{n\phi(\Phi^{-1}(\frac{i}{n}))}, \quad 1 \leq i \leq n-1,$$

とおいてみる。すると式(4)の左辺は

$$\begin{aligned} \sum_{j=1}^{n-1} w_i S_{ij} w_j &= \frac{1}{n^2 \phi(\Phi^{-1}(\frac{i}{n}))} \left(\sum_{j=1}^i \frac{j}{n} \left(1 - \frac{i}{n}\right) \frac{1}{\phi(\Phi^{-1}(\frac{j}{n}))} + \sum_{j=i+1}^{n-1} \frac{i}{n} \left(1 - \frac{j}{n}\right) \frac{1}{\phi(\Phi^{-1}(\frac{j}{n}))} \right) \\ &\simeq \frac{1}{n\phi(\Phi^{-1}(x))} \left(\int_0^x \frac{y(1-x)}{\phi(\Phi^{-1}(y))} dy + \int_x^1 \frac{x(1-y)}{\phi(\Phi^{-1}(y))} dy \right) \quad (x = \frac{i}{n}, y = \frac{j}{n}) \\ &= \frac{1}{n} \end{aligned}$$

と近似できる。最後の等号は部分積分による。よって近似的に式(4)が成り立つことが分かる。対応するスコア $y(i/n)$ は

$$\begin{aligned} y(i/n) &= \sum_{j=1}^{n-1} w_j (1_{\{i \geq j\}} - (n-j)/n) \\ &\simeq \int_0^1 \frac{1}{\phi(\Phi^{-1}(s))} (1_{\{s \leq x\}} - 1 + s) ds \quad (x = \frac{i}{n}) \\ &= \Phi^{-1}(x) \end{aligned}$$

となる。よってスコアが近似的に標準正規分布の分位点関数となることが示された。

[13]では、別のアプローチとして、式(4)の極限が微分積分方程式となることを示し、その解が $\Phi^{-1}(x)$ であることを示している。

5 客観重みの数値計算

与えられた共分散行列 S に対し、式(2)を満たすベクトル w を求めるには、式(3)で定義される関数 $\psi(w)$ の最小化問題を汎用アルゴリズムで解くのが簡単である。表6にRの関数 `optim` を用いた場合のコードの例を示す。

他には、ホモトピー法によって解く方法 [11] や、逐次的に二次方程式を解く方法 [13] がある。後者は、式(3)の $\psi(w)$ を w_i の1変数関数と見なした時、その最小点が二次方程式

$$S_{ii} w_i^2 + \left(\sum_{j \neq i} S_{ij} w_j \right) w_i - 1 = 0$$

の解であることを利用して、これを $i = 1, \dots, p$ について繰り返し解いていく、という方法である。汎用アルゴリズムを使えない場合や、 S がスパースな場合には有用と考えられる。

表 6: R の関数 optim を使った客観重みの計算例.

```

ogi.optim = function(S){
  f = function(w, S){ sum(w * (S%*%w))/2 - sum(log(w)) }
  gr = function(w, S){ S%*%w - 1/w }
  p = nrow(S)
  optim(rep(1,p), f, gr, S, method="L-BFGS-B",
        lower=rep(0,p), upper=rep(Inf,p))
}

# 例
S = matrix(c(1,-1/2,-1/2, -1/2,1,0, -1/2,0,1), 3, 3, byrow=TRUE)
w = ogi.optim(S)$par

show(w) # 結果
show(w*(S%*%w)) # 確認

```

また、もし逆行列 S^{-1} に対する客観重み v を求めることができれば、元の S の重みは $w = 1/v$ (要素ごとの逆数) で与えられる。なぜならば、 $S^{-1}v = 1/v$ の両辺に S を掛ければ $v = S(1/v)$ 、よって $S(1/v) = 1/(1/v)$ が成り立つからである。この事実は、 S^{-1} がスパースな場合 (つまり多くの偏相関が 0 の場合) に有用と考えられる。

6 今後の課題

4 節ではカテゴリカル変数に対する数量化の極限として正規分布が現れることを指摘した。しかし、実際に (何らかの位相で) 収束することは厳密には証明できていない。これを証明するという課題が残っている。

また、推測統計的な考察、すなわち真の重みパラメータを標本から推定する方法やその性質を調べる必要がある。特に、高次元小標本 ($p \rightarrow \infty$ かつ n 固定) の枠組みで、OGI の良い推定量が得られるかどうかは興味深い話題である。

別の観点としては、アローの不可能性定理に代表される社会的選択理論との関連も明らかにする必要がある (例えば [9])。またデータ可視化の立場からは、客観重みを用いた平行座標プロットと、textile plot [8] との比較あるいは融合も課題である。

謝辞

本研究は JSPS 科研費 26540013, 26108003 の助成を受けたものである。

参考文献

- [1] Baker, R. J., (1974). Selection indexes without economic weights for animal breeding. *Canad. J. Animal Sci.*, 54 (1), 1–8.
- [2] Bradley, R. A., Katti, S. K., Coons, I. J. (1962). Optimal scaling for ordered categories. *Psychometrika*, 27 (4), 355–374.
- [3] Cox, T. F., Dunn, R. T., (2002). An analysis of decathlon data. *J. Roy. Statist. Soc.*, Ser. D 51 (2), 179–187.
- [4] Dill, D. D., Soo, M. (2005). Academic quality, league tables, and public policy: A cross-national analysis of university ranking systems. *Higher Education*, 49 (4), 495–533.
- [5] Elston, R. C. (1963). A weight-free index for the purpose of ranking or selection with respect to several traits at a time. *Biometrics*, 19 (1), 85–97.
- [6] IAAF Council (Ed.), (2001). *IAAF Scoring Tables for Combined Events*. IAAF.
- [7] JST CREST 日比チーム (2011). 「グレブナー道場」, 共立出版.
- [8] Kumasaka, N., Shibata, R. (2008). High-dimensional data visualisation: The textile plot. *Comput. Statist. Data Anal.*, 52, 3616–3644.
- [9] Langville, A. N. and Meyer, C. D. (岩野 和生, 中村 英史, 清水 咲里 訳) (2015). レイティング・ランキングの数理 — No.1 は誰か? —, 共立出版.
- [10] Marshall, A. W., Olkin, I. (1968). Scaling of matrices to achieve specified row and column sums. *Numer. Math.*, 12, 83–90.
- [11] O'Leary, D. P. (2003). Scaling symmetric positive definite matrices to prescribed row sums. *Linear Algebra Appl.*, 370, 185–191.
- [12] Saito, T., Otsu, T. (1988). A method of optimal scaling for multivariate ordinal data and its extensions. *Psychometrika*, 53 (1), 5–25.
- [13] Sei, T. (2016). An objective general index for multivariate ordered data, *J. Multivariate Anal.*, 147, 247–264.