

R. ニャナデシカンの一般化主成分分析について

新井 宏 尚

1. 初めに

本ノートでは K. Pearson (1902) により最初提案され、後に H. Hotelling (1933) が現在ある形に発展させた主成分分析をさらに拡張したと言う R. ニャナデシカンの一般化主成分分析法について考察してみる。

2. 通常の主成分分析と一般化主成分分析

通常の主成分分析法の概要をここで説明すれば次のようになり。いま p 個の特性を持つ変量 x_1, x_2, \dots, x_p に関して n 個のサンプルが与えられるならばベクトル \mathbf{x} 、マトリックス \mathbf{X} は次のようになり

$$\mathbf{x}' = (x_1, x_2, x_3, \dots, x_p)$$

\mathbf{X} , \mathbf{X}' の行で n 個の p 次元データを示す

したが、標本平均ベクトル ($\bar{\mathbf{x}}$)、分散、共分散行列 (\mathbf{S})、相関行列 (\mathbf{R}) は(1), (2), (3)のごとく定義できよう。

$$\bar{\mathbf{x}}' = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p) = \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{X}' \quad (1)$$

$$\mathbf{S} = (s_{ij}) = \frac{1}{n-1} (\mathbf{X} - \bar{\mathbf{X}})(\mathbf{X} - \bar{\mathbf{X}})' \quad (2)$$

$$\mathbf{R} = \mathbf{D}_{1/\sqrt{s_{ii}}} \cdot \mathbf{S} \cdot \mathbf{D}_{1/\sqrt{s_{ii}}} \quad (3)$$

ただし、 $\mathbf{D}_{1/\sqrt{s_{ii}}}$ は i 対角要素が $1/\sqrt{s_{ii}}$ ($i=1, 2, \dots, p$) に等しい対角行列である。

このとき、主成分を見つけることは代数的には標本分散行列 \mathbf{S} の固有値と固有ベクトルを見つける問題に帰着する。つまりベクトル $\mathbf{a}' = (a_1, a_2, \dots, a_p)$ と \mathbf{x} との線形結合 $\mathbf{a}' \cdot \mathbf{x}$ を求め

$$\mathbf{a}' \cdot \mathbf{x} \quad (4)$$

式(4)を $\mathbf{a}' \cdot \mathbf{a} = 1$ の条件下で第1主成分が最大の分散を持つように \mathbf{a} を求める。このためには(4)式と条件 $\mathbf{a}' \cdot \mathbf{a} = 1$ より、ラグランジュ乗法により

$$\phi = \mathbf{a}' \cdot \mathbf{x} - \lambda(\mathbf{a}' \cdot \mathbf{a} - 1) \quad \text{subject to} \quad \mathbf{a}' \cdot \mathbf{a} = 1 \quad (5)$$

と書き改め、

$$\partial \phi / \partial \mathbf{a} = 0 \quad (6)$$

とすれば、分散 $V(\mathbf{a}' \cdot \mathbf{x}) = \mathbf{a}' \cdot \mathbf{S} \cdot \mathbf{a}$ より

$$\begin{aligned} \partial \phi / \partial \mathbf{a} &= 2 \cdot \mathbf{S} \cdot \mathbf{a} - 2 \cdot \lambda \cdot \mathbf{a} = 0 \\ (\mathbf{S} - \lambda \cdot \mathbf{I}) \cdot \mathbf{a} &= 0 \end{aligned} \quad (7)$$

が得られる。したがって、(7)式が $\mathbf{a} = 0$ 以外の解を持つのは

$$\begin{aligned} \text{行列式 } |\mathbf{S} - \lambda \cdot \mathbf{I}| \text{ が} \\ |\mathbf{S} - \lambda \cdot \mathbf{I}| &= 0 \end{aligned} \quad (8)$$

となることである。(8)を解くには分散行列 \mathbf{S} の固有値を求めればよいので $\mathbf{a}' \cdot \mathbf{a} = 1$ より

$$\begin{aligned} \mathbf{a}'(\mathbf{S} - \lambda \cdot \mathbf{I})\mathbf{a} &= \mathbf{a}'\mathbf{S}\mathbf{a} - \lambda = 0 \\ \therefore \mathbf{a}'\mathbf{S}\mathbf{a} &= \lambda \end{aligned}$$

でかつ、 $\text{Max}(\mathbf{a}'\mathbf{S}\mathbf{a})$ とするには $\text{Max}\{\lambda_i; i=1, 2, \dots, p\}$ となる $\lambda = \lambda_i$ を選択し、それに対応する固有ベクトルを求めれば良い。

しかしこの手順にも欠点がある。それは変量 x_i の測定単位が全て同じの場合には都合が良いが測定単位を変更する場合には、分散、共分散行列の両側に対角行列 ($D_{1/\sqrt{s_{ii}}}$) を掛けねばならない(3式)。このような場合には変量の尺度不変性は維持されず固有値、固有ベクトルに尺度交換がどのようなことを及ぼすかは複雑となる(固有値、固有ベクトルが全く異ってしまう場合がある。*)。したがって、主成分分析を行う場合に通常相関行列 \mathbf{R} からスタートする方がよりスムーズとされている。

しかし、相関行列から出発するには変量相互間は互いに線型関係を持っていることが前提とされており、その有意性の検定を行わなければならないので変量相互間に非線型の関係がみられる場合にはどのようにすべきなのかという問題に行きつく*²。一般的には変量相互間に線型でない関係がみられるときには通常相関係数は使用せず相関比により2つの関係の度合 ($\eta_{x_i x_j}$) を測ることができるが $\eta_{x_i x_j}$ の行列を作っても相関行列と異なり非対称行列*³になるためにこのような接近法は考えられない。

他方R. ニヤナデシカンは非線型の主成分分析を(*4)の中で展開している。その骨旨を(*4)に従って述べれば次のようになる。

例として $p=2$ の場合について説明している。 $y=f(x_1, x_2)$ の関数型 f を最初に特定化する。

* 1. 日科技連奥野忠一他, 多変量解析法 pp. 181-182

* 2. 非線型の成分分析の示唆については、竹内啓「計量経済学の研究」東洋経済新報社, p. 160 に記載されている。

* 3. 非対称行列の対称化の方法は存在しないであろう。

* 4. R. ニヤナデシカン, 丘本他訳, 統計的多変量データの解析 日科技連

2次多項式として f を特定化する。つまり

$$y = a_1 x_1 + a_2 x_2 + a_3 x_1 \cdot x_2 + a_4 \cdot x_1^2 + a_5 \cdot x_2^2 \quad (9)$$

と y を置いて、 y の分散を最大になるように係数を定める。

ここで、 $x_1 \cdot x_2$ を x_3 、 x_1^2 を x_4 、 x_2^2 を x_5 と表わせば

式(9)は

$$y^* = a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 + a_5 x_5 \quad (10)$$

となるので $\mathbf{a}^* = (a_1, a_2, a_3, a_4, a_5)$ で係数ベクトルを表わすならば問題は $\mathbf{a}^* \mathbf{a} = 1$ のもとで $\mathbf{a}^* \cdot \mathbf{y}^*$ の分散を最大にするように \mathbf{a}^* を決定することになる。また標本平均ベクトルと分散、共分散行列は次のように定義できる

$$\bar{\mathbf{x}}^* = (\bar{x}_1, \bar{x}_2, \bar{x}_3, \bar{x}_4, \bar{x}_5) \quad (11)$$

$$\mathbf{S}^* = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i^* - \bar{\mathbf{x}}^*)(\mathbf{x}_i^* - \bar{\mathbf{x}}^*)' \quad (12)$$

そして、通常の主成分分析（線型主成分分析）を実行する。また $n \times p$ コのデータ、マトリックスが与えられたときに2次の主成分分析を実行するには p コの変数に $p + [p(p-1)/2]$ コの変数を付け加えて

$$\left[2p + \frac{p(p-1)}{2} \right] \times \left[2p + \frac{p(p-1)}{2} \right]$$

の分散行列を使い、通常の主成分分析を行えば良いとしている。

次にこの手法による例*⁴を1つ上げその後の問題に触れよう。

(例) 表1-1は x_1 のデータを与えたときの、放物線 $x_2 = 2 + 4x_1 + 4x_1^2$ の軌跡をえがいている。このときの2次の主成分分析を実行した結果、得られた固有値と固有ベクトルが表1-1と表1-2に載っているものである。

次に、式(9)のバリエーションとしてデータに統計的な誤差を含めた場合の固有値と固有ベクトルが表2に載っている。この場合の統計的誤差とは各観測値の組に正規乱数（平均ゼロ、分散 $\frac{1}{6}$ ）を

表 1-1 2次の主成分分析の例に対する固有値解析

		固 有 値		
2,163.634	219.915	2.258	2.223	0.000009
		固 有 ベ ク ト ル		
-.002513	.246077	.508757	-.442445	.696310
.169321	.011882	.758811	.604229	-.164076
-.094253	.932909	-.212548	.274991	.0000004
.044843	-.243106	.319056	.593499	.696311
.980015	.099425	-.135640	-.106239	.0000003

*4 統計的多変量データ解析，ニヤナデンカン，丘本他訳，日科技連 pp. 49-50

表 1-2 データおよび最大，最小固有値から得られた座標

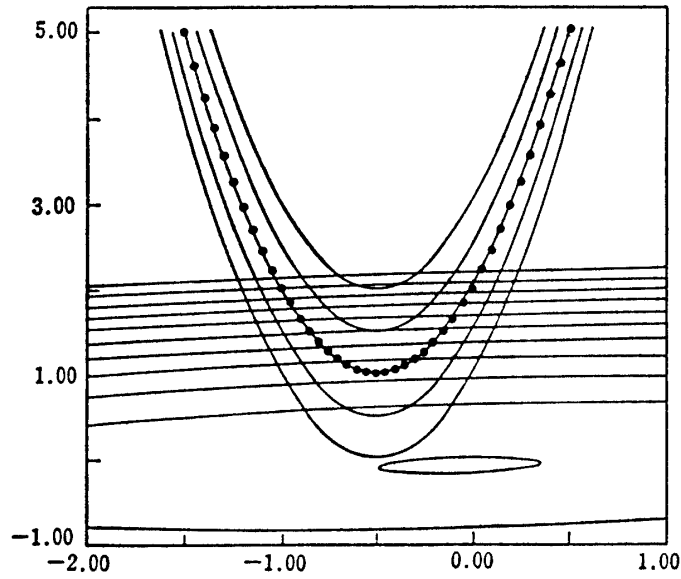
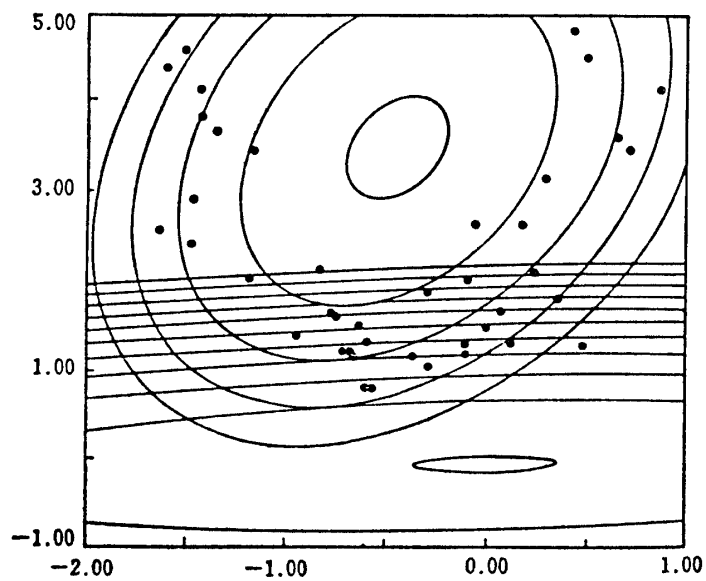


表 2-1 2次の主成分分析の例に対する固有値解析

		固 有 値				
1969.563	285.938	5.717	2.034	1.081		
		固 有 ベ ク ト ル				
-.010915	.260105	-.269383	.610264	.698023		
.172863	.012363	.189273	.761768	-.594853		
-.106665	.930317	.301394	-.146763	-.103707		
.062983	-.230187	.892766	.062053	.377047		
.977064	.117118	-.061142	-.147978	.077413		

表 2-2 データおよび最大，最小固有値から得られた座標



加えたものである。この2つの結果をみると、誤差が、固有値、固有ベクトルに影響を与えていることもわかる。

しかし、ここまでのところで変量が誤差を含む場合のモデルについて疑問点が生じよう。

先づ第1にノイズが入った場合について式(9)をそのまま適用しても良いのであろうか。通常、変量に誤差が含まれる場合、変量 x_1, x_2 は

$$x_1 = \hat{x}_1 + \varepsilon_1 \quad (13)$$

$$x_2 = \hat{x}_2 + \varepsilon_2$$

ただし、 \hat{x}_1, \hat{x}_2 は真の値、 $\varepsilon_1, \varepsilon_2$ は互いに $\varepsilon_1 \sim N\left(0, \frac{1}{6}\right)$, $\varepsilon_2 \sim N\left(0, \frac{1}{6}\right)$ とする。

と表わすことができるので、(9)式は

$$y = a_1 \hat{x}_1 + a_2 \hat{x}_2 + a_3 \hat{x}_1 \hat{x}_2 + a_4 \hat{x}_1^2 + a_5 \hat{x}_2^2 + \zeta_i \quad (14)$$

ただし、 $\zeta_i = a_1 \varepsilon_1 + a_2 \varepsilon_2 + a_3 \varepsilon_1 \varepsilon_2 + a_4 \varepsilon_1^2 + a_5 \varepsilon_2^2 + a_3 \hat{x}_1 \varepsilon_2 + a_4 \hat{x}_2 \varepsilon_1 + 2 a_4 \hat{x}_1 \varepsilon_1 + 2 a_5 \hat{x}_2 \varepsilon_2$

と書き改めることができよう。したがい、データに統計的誤差が含まれるときには、誤差 ζ_i をうまく考慮して後に(14)式を処理すべきではなかろうか。

第2は $y = \sum_i a_i f_i(x_1, x_2, \dots, x_p)$ の関数型の特定化であろう。 f の役割は p 次元空間的に散布している n 個の点をいかにうまくその曲面上にのせるかにあるからである。