

主成分分析

京都大学大学院工学研究科化学工学専攻
プロセスシステム工学研究室
加納 学

1997年1月 第1版作成

2002年5月 第2版作成

Copyright ©1997-2002 by Manabu Kano. All rights reserved.

[注意事項]

自由に利用させていただいて結構ですが、著作権は一切放棄していません。また、本資料の間違いなどによって生じた不利益などに対して、著者は一切責任を負いません。勿論、間違いの指摘やアドバイスは歓迎します。

目次

1	主成分分析とは	3
2	主成分の導出	4
2.1	準備	4
2.2	第1主成分の導出	4
2.3	第 m 主成分の導出	7
2.4	データの標準化	8
3	寄与率と因子負荷量	8
3.1	寄与率	8
3.2	因子負荷量	9
4	主成分分析の行列表現	10
4.1	主成分得点の計算	10
4.2	データの再構築	12
4.3	主成分分析と特異値分解	13

1 主成分分析とは

身体検査で身長、体重という2種類の変数が10人について測定されており、このデータをプロットするとFig.1のようになったとする。この図から、身長が大きくなるほど体重も大きくなるという傾向が把握できる。この傾向を表現するためには、新たに図中 z_1 のような軸を考えればよいであろう。この z_1 軸は「体の大きさ」を表していると解釈できる。しかし、「体の大きさ」だけで各人の特徴をすべて表現できるわけではない。10人の中には太っている人も痩せている人もいるのである。そこで、 z_1 軸と直交する z_2 軸を考える。この z_2 軸は「肥満の程度」を表していると解釈できる。このように、「身長」や「体重」といった変数を独立に扱うのではなく、「体の大きさ」や「肥満の程度」といった総合的な指標を導入することによって、データに含まれる変数間の関係や特徴が容易に把握できるようになる。このような総合的な指標を統計的に設定し、変数間の関係を把握するための手法が、主成分分析(PCA; Principal Component Analysis)と呼ばれるものである。主成分とは総合的な指標のことであり、身体検査の例では、「体の大きさ」を表す z_1 が第1主成分、「肥満の程度」を表す z_2 が第2主成分と呼ばれる。

「体の大きさ」を表す z_1 によってデータの持つ特徴(あるいは傾向)は概ね表現できているので、 z_1 のみでデータを代表させることもできる。この場合、2種類の変数によって表現されていたデータを1種類の変数で表現することになる。すなわち、データの持つ特徴を主成分を用いて表すことにより、情報の損失を最小限に抑えながら、2次元のデータを1次元のデータに変換したわけである。このような観点から、主成分分析をデータの低次元化を行うための手法とみることもできる。

身体検査の結果を「体の大きさ」で代表させる場合、「肥満の程度」に関する情報は失われることになる。このようにデータの低次元化は情報の損失を伴うが、そのときの情報の損失量は各データ点から第1主成分 z_1 へおろした垂線の長さで表される。データの持つ特徴を最もよく表現するためには、情報損失量をできる限り小さくする必要がある。実際、主成分分析では、情報損失量を最小化するという条件の下で主成分が決定される。

情報損失量を最小にするということは、得られる情報量を最大にすることの裏返しである。これは、データのばらつきを最もよく表す方向、すなわち分散が最大となる方向に主成分を設定することによって実現できる。このことをFig.1を用いて示そう。いま、点 $A(x_{11}, x_{12})$ に着目すると、第1主成分 z_1 のみでデータを代表させる場合の情報損失量は点 A から直線 z_1 へ下ろした垂線の長さ \overline{AB} となり、 z_1 を用いて得られる情報量は \overline{OB} で与えられる。ここで、点 O は z_1 軸の原点であり、データの重心である。この新しい情報量 \overline{OB} は第1主成分得点と呼ばれる。情報損失量と新しい情報量(主成分得点)の間には、明らかに

$$\overline{OA}^2 = \overline{OB}^2 + \overline{AB}^2 \quad (1)$$

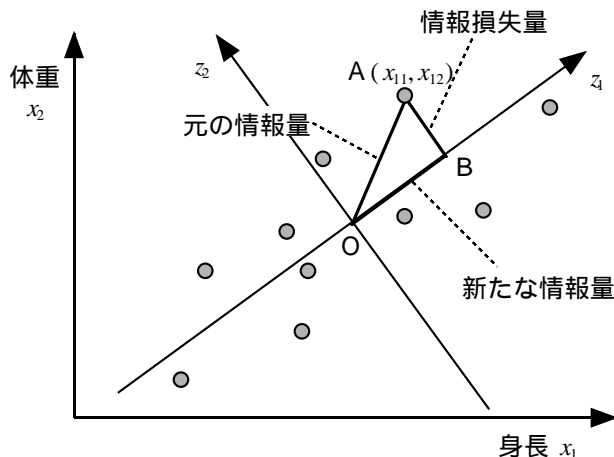


Fig. 1 身体検査の結果

という関係が成り立つ。 \overline{OA} を元の情報量と呼ぶことにすると、10 点のデータ全体について、

$$(\text{元の情報量の二乗和}) = (\text{新しい情報量の二乗和}) + (\text{情報損失量の二乗和})$$

という関係が成り立つ。ここで、元の情報量の二乗和はある定数であるから、情報損失量の二乗和を最小にすることと新しい情報量の二乗和を最大にすることが等価であることがわかる。

ここまで 2 変数の場合について述べてきたが、多変数の場合にも本質的な違いはない。すなわち、主成分分析とは、 P 個の変数 $\{x_p\} (p = 1, 2, \dots, P)$ の持つ情報を、情報の損失を最小限に抑えながら、 $\{x_p\}$ の一次結合として与えられる互いに独立な $M (M \leq P)$ 個の主成分 (総合的指標) $\{z_m\}$

$$z_m = \sum_{p=1}^P w_{pm} x_p \quad (m = 1, 2, \dots, M) \quad (2)$$

を用いて表現する手法である。なお、 z_m は第 m 主成分と呼ばれ、その結合係数 $\{w_{pm}\} (p = 1, 2, \dots, P; m = 1, 2, \dots, M)$ は以下の条件を満足するように決定される。

<条件>

第 1 主成分 z_1 の分散は $\{x_p\} (p = 1, 2, \dots, P)$ のあらゆる 1 次式の持つ分散の中で最大であり、第 m 主成分 $\{z_m\} (m = 2, \dots, M)$ の分散は $\{z_{m'}\} (m' = 1, 2, \dots, m-1)$ の全てと無相関な 1 次式の持つ分散の中で最大である。ただし、

$$\sum_{p=1}^P w_{pm}^2 = 1 \quad (m = 1, 2, \dots, M) \quad (3)$$

とする。

2 主成分の導出

2.1 準備

Sec.1 で述べた条件に従い、主成分の分散が最大となるように主成分を決定する方法について述べる。いま、 P 個の変数について N 個のサンプルがある場合を考え、測定値を $\{x_{np}^*\} (n = 1, 2, \dots, N; p = 1, 2, \dots, P)$ とする。以下の議論を簡単にするために、各変数についてその平均値 $\{\bar{x}_p\} (p = 1, 2, \dots, P)$ からの偏差 $\{x_{np}\}$ を導入する。すなわち、

$$x_{np} = x_{np}^* - \bar{x}_p \quad (n = 1, 2, \dots, N; p = 1, 2, \dots, P) \quad (4)$$

とする。このとき、測定データ全体は次の行列 \mathbf{X} で与えられる。

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1P} \\ x_{21} & x_{22} & \cdots & x_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NP} \end{pmatrix} \quad (5)$$

2.2 第 1 主成分の導出

第 1 主成分 z_1 は Eq.(2) で与えられるので、その結合係数を

$$\mathbf{w}_1 = \begin{pmatrix} w_{11} \\ w_{21} \\ \vdots \\ w_{P1} \end{pmatrix} \quad (6)$$

とすると， n 番目のサンプル

$$\mathbf{x}_n = \begin{pmatrix} x_{n1} & x_{n2} & \cdots & x_{nP} \end{pmatrix} \quad (7)$$

に対応する第 1 主成分 z_1 の値 t_{n1} は

$$\begin{aligned} t_{n1} &= \sum_{p=1}^P w_{p1} x_{np} \\ &= \mathbf{x}_n \mathbf{w}_1 \end{aligned} \quad (8)$$

となる．この第 1 主成分 z_1 の値 t_{n1} を第 1 主成分得点と呼ぶ． N 個のサンプルに対応する第 1 主成分得点を 1 つのベクトルにまとめ，

$$\mathbf{t}_1 = \begin{pmatrix} t_{11} \\ t_{21} \\ \vdots \\ t_{N1} \end{pmatrix} \quad (9)$$

とおくと，

$$\mathbf{t}_1 = \mathbf{X} \mathbf{w}_1 \quad (10)$$

が成り立つ．第 1 主成分得点の平均値 \bar{t}_1 は

$$\begin{aligned} \bar{t}_1 &= \frac{1}{N} \sum_{n=1}^N t_{n1} \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{w}_1 \\ &= \frac{1}{N} \sum_{n=1}^N \left(\sum_{p=1}^P w_{p1} x_{np} \right) \\ &= \frac{1}{N} \sum_{p=1}^P w_{p1} \left(\sum_{n=1}^N x_{np} \right) \\ &= 0 \end{aligned} \quad (11)$$

であるから，第 1 主成分 z_1 の分散 $\sigma_{z_1}^2$ は

$$\begin{aligned} \sigma_{z_1}^2 &= \frac{1}{N-1} \mathbf{t}_1^T \mathbf{t}_1 \\ &= \frac{1}{N-1} (\mathbf{X} \mathbf{w}_1)^T (\mathbf{X} \mathbf{w}_1) \\ &= \mathbf{w}_1^T \mathbf{V} \mathbf{w}_1 \\ &\geq 0 \end{aligned} \quad (12)$$

となる．なお，行列 \mathbf{V} は共分散行列と呼ばれる非負定値行列であり，

$$\mathbf{V} = \frac{1}{N-1} \mathbf{X}^T \mathbf{X} \quad (13)$$

で与えられ，その (i, j) 要素 v_{ij} は

$$v_{ij} = \frac{1}{N-1} \sum_{n=1}^N x_{ni} x_{nj} \quad (14)$$

$$= \frac{1}{N-1} \sum_{n=1}^N (x_{ni}^* - \bar{x}_i) (x_{nj}^* - \bar{x}_j) \quad (15)$$

である．また，明らかに， $v_{ij} = v_{ji}$ すなわち $V = V^T$ が成り立つ．

前節の条件より，第 1 主成分 z_1 は Eq.(3) の下でその分散 $\sigma_{z_1}^2$ が最大となるように決定されなければならない．この最適化問題は Lagrange 乗数法を用いて容易に解くことができる．すなわち，Lagrange 乗数 λ を導入して

$$J_1 = \mathbf{w}_1^T \mathbf{V} \mathbf{w}_1 - \lambda (\mathbf{w}_1^T \mathbf{w}_1 - 1) \quad (16)$$

とおき， J_1 を最大にするような結合係数 \mathbf{w}_1 を求めればよい．そこで， J_1 を \mathbf{w}_1 で偏微分して 0 とおくと，

$$\begin{aligned} \frac{\partial J_1}{\partial \mathbf{w}_1} &= \begin{pmatrix} \frac{\partial J_1}{\partial w_{11}} \\ \frac{\partial J_1}{\partial w_{21}} \\ \vdots \\ \frac{\partial J_1}{\partial w_{P1}} \end{pmatrix} \\ &= \begin{pmatrix} 2 \sum_{p=1}^P v_{1p} w_{p1} \\ 2 \sum_{p=1}^P v_{2p} w_{p1} \\ \vdots \\ 2 \sum_{p=1}^P v_{Pp} w_{p1} \end{pmatrix} - 2\lambda \begin{pmatrix} w_{11} \\ w_{21} \\ \vdots \\ w_{P1} \end{pmatrix} \\ &= 2\mathbf{V} \mathbf{w}_1 - 2\lambda \mathbf{w}_1 \\ &= \mathbf{0} \end{aligned} \quad (17)$$

となることから

$$(\mathbf{V} - \lambda \mathbf{I}) \mathbf{w}_1 = \mathbf{0} \quad (18)$$

という条件式が得られる．これは固有値問題に他ならず，Lagrange 乗数 λ が満たすべき条件は固有方程式

$$\det |\mathbf{V} - \lambda \mathbf{I}| = 0 \quad (19)$$

を用いて表わすこともできる．以上より，Lagrange 乗数 λ および第 1 主成分 z_1 の結合係数 \mathbf{w}_1 はそれぞれ共分散行列 V の固有値および固有ベクトルとして与えられることがわかる．

さて，共分散行列 V は P 次元正方行列であるため，その固有値 λ は P 個存在するが，その中のいずれの固有値が第 1 主成分 z_1 の分散 $\sigma_{z_1}^2$ を最大にする結合係数 (固有ベクトル) \mathbf{w}_1 に対応しているのだろうか．この問題について考えてみよう．第 1 主成分 z_1 の分散 $\sigma_{z_1}^2$ は Eq.(12) で与えられ，その結合係数 \mathbf{w}_1 は条件 Eq.(18) を満足しなければならない．そこで，Eq.(12) に Eq.(18) を代入し， $\mathbf{w}_1^T \mathbf{w}_1 = 1$ に注意すると，

$$\begin{aligned} \sigma_{z_1}^2 &= \mathbf{w}_1^T \mathbf{V} \mathbf{w}_1 \\ &= \mathbf{w}_1^T \lambda \mathbf{w}_1 \\ &= \lambda \end{aligned} \quad (20)$$

を得る．これより，第 1 主成分 z_1 の分散 $\sigma_{z_1}^2$ は共分散行列 V の固有値 λ に等しいことがわかる．従って，最大にすべき第 1 主成分 z_1 の分散 $\sigma_{z_1}^2$ は共分散行列 V の最大固有値に等しくなり，その結合係数 \mathbf{w}_1 は最大固有値に対応する固有ベクトルとして求めることができる．

2.3 第 m 主成分の導出

第2主成分以下の結合係数 $\{\mathbf{w}_m\} (m = 2, 3, \dots, M)$ も第1主成分と同様な手順で求めることができる。ただし、第 m 主成分 $z_m (m = 2, \dots, M)$ の分散 $\sigma_{z_m}^2$ が主成分 $\{z_{m'}\} (m' = 1, 2, \dots, m-1)$ の全てと無相関な1次式の持つ分散の中で最大となるように、結合係数 \mathbf{w}_m を決定しなければならない。当然ながら、Eq.(3)も満たされなければならない。

ここでは、数学的帰納法を用いて、第 m 主成分の結合係数 \mathbf{w}_m を導出する方法を示す。いま、第 $m-1$ 主成分まで求められており、それらの結合係数 $\{\mathbf{w}_i\} (i = 1, 2, \dots, m-1)$ は条件

$$(\mathbf{V} - \lambda_i \mathbf{I}) \mathbf{w}_i = \mathbf{0} \quad (21)$$

$$\mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases} \quad (22)$$

を満たしているとして、第 m 主成分の結合係数 \mathbf{w}_m を Lagrange 乗数法を用いて求める。すなわち、

$$J_m = \mathbf{w}_m^T \mathbf{V} \mathbf{w}_m - \lambda_m (\mathbf{w}_m^T \mathbf{w}_m - 1) - \sum_{i=1}^{m-1} \mu_i \mathbf{w}_m^T \mathbf{w}_i \quad (23)$$

とにおいて、 J_m を最大にするような \mathbf{w}_m を求めることにする。ここで、右辺最終項が第 m 主成分 $z_m (m = 2, \dots, M)$ と主成分 $\{z_{m'}\} (m' = 1, 2, \dots, m-1)$ の全てが無相関となる条件を表している。 J_m を \mathbf{w}_m で偏微分して0とおくと、

$$\begin{aligned} \frac{\partial J_m}{\partial \mathbf{w}_m} &= 2\mathbf{V} \mathbf{w}_m - 2\lambda_m \mathbf{w}_m - \sum_{i=1}^{m-1} \mu_i \mathbf{w}_i \\ &= \mathbf{0} \end{aligned} \quad (24)$$

となる。 $\mathbf{w}_j^T (j = 1, 2, \dots, m-1)$ を左から掛けて、Eq.(22) を用いると、

$$\mathbf{w}_j^T \mathbf{V} \mathbf{w}_m - \mu_j = 0 \quad (j = 1, 2, \dots, m-1) \quad (25)$$

を得る。ここで、左辺第1項は、 $\mathbf{V} = \mathbf{V}^T$ および Eq.(21) より

$$\begin{aligned} \mathbf{w}_j^T \mathbf{V} \mathbf{w}_m &= \mathbf{w}_m^T \mathbf{V} \mathbf{w}_j \\ &= \mathbf{w}_m^T \lambda_j \mathbf{w}_j \\ &= 0 \quad (j = 1, 2, \dots, m-1) \end{aligned} \quad (26)$$

となることから、Eq.(25) は

$$\mu_j = 0 \quad (j = 1, 2, \dots, m-1) \quad (27)$$

となる。この結果を Eq.(24) に代入すると、最終的に

$$(\mathbf{V} - \lambda_m \mathbf{I}) \mathbf{w}_m = \mathbf{0} \quad (28)$$

を得る。これは Eq.(21) と同一の式であり、第 m 主成分 z_m の分散 $\sigma_{z_m}^2$ もまた共分散行列 \mathbf{V} の固有値に等しいことがわかる。ただし、大きい方から $m-1$ 個の固有値とそれに対応する固有ベクトルはすでに第 $m-1$ 番目までの主成分を表すのに使われているため、第 m 主成分 z_m の分散 $\sigma_{z_m}^2$ は共分散行列 \mathbf{V} の m 番目に大きい固有値に等しくなり、結合係数 \mathbf{w}_m はその固有値に対応する固有ベクトルとして求めることができる。

2.4 データの標準化

ここまででは、各変数の測定値をそのまま用いて、その共分散行列 V に基づいて主成分を決定する方法について述べた。しかし、各変数が異なる単位で測定されている場合には、単位の取り方によって異なる主成分が得られることになる。また、たとえ単位が同じであったとしても、大きく分散の異なる変数に対してそのまま主成分分析を適用すれば、その結果は分散の大きな変数の影響を強く受けることになり、変数間の関係を正しく把握することはできないであろう。従って、全ての変数を何らかの方法を用いて標準化する必要がある。

最も簡単で広く利用されている方法は、各変数を平均 1、分散 1 となるように標準化する方法である。具体的には、 P 個の変数について N 個のサンプルがある場合、その測定値 x_{np}^* ($n = 1, 2, \dots, N; p = 1, 2, \dots, P$) を用いる代わりに、

$$\tilde{x}_{np} = \frac{x_{np}^* - \bar{x}_p}{\sigma_{x_p}} \quad (29)$$

を用いる。ここで、 \bar{x}_p, σ_{x_p} はそれぞれ p 番目の変数 x_p の平均値および標準偏差である。この場合、標準化された変数間の共分散は相関係数に等しくなるので、標準化されたデータ行列 \tilde{X} の相関行列を

$$R = \frac{1}{N-1} \tilde{X}^T \tilde{X} \quad (30)$$

とおくと、第 m 主成分 z_m の分散 $\sigma_{z_m}^2$ は相関行列 R の m 番目に大きな固有値に等しくなり、第 m 主成分 z_m の結合係数 w_m はその固有値に対応する固有ベクトルとして求めることができる。

上記の標準化を施すことによって、単位に依存しない分析を行うことが可能になる。しかし、実際の測定データは大きさや性質の異なる様々な誤差を含むため、誤差の影響を強く受けている変数を用いた分析は誤った結果を導く恐れがある。そこで、誤差の影響の程度に応じて各変数に異なる重みを付けるという方法が考えられる。重みの付け方としてはいろいろ考えられるが、誤差の影響を全く受けていない変数は全て分散が 1 となるように、誤差のみによって変動している変数についてはその分散が 0 となるようにするのが直感的にも妥当であろう。このような重みの付け方として、測定値 x_{np}^* ($n = 1, 2, \dots, N; p = 1, 2, \dots, P$) の代わりに、

$$\tilde{x}'_{np} = \frac{x_{np}^* - \bar{x}_p}{\sigma_{x_p}} \frac{\sigma_{x_p} - \sigma_{x_{pe}}}{\sigma_{x_p}} \quad (31)$$

を用いる方法がある。ここで、 $\sigma_{x_{pe}}$ は p 番目の変数 x_p に含まれる誤差の標準偏差である。現実には、誤差の標準偏差を正確に把握することは困難であるが、何らかの方法でその推定値が与えられている場合には、ここで述べたような重み付けが有効となる。

3 寄与率と因子負荷量

3.1 寄与率

主成分分析とは少数の総合的指標 (主成分) を用いて変数間の関係や特徴を把握するための統計的手法である。従って、各主成分が元のデータに含まれる特徴をどの程度表現しているのか、あるいは何個の主成分を採用すれば元のデータに含まれる特徴を十分に表現できるのかを知ることが必要になる。このための指標として、寄与率および累積寄与率がある。

合計 P 個の変数の分散の和は、その共分散行列を V とすると、 V の (p, p) 要素 v_{pp} が変数 x_p の分散に等しくなるので、

$$\begin{aligned} \sum_{p=1}^P \sigma_{x_p}^2 &= \sum_{p=1}^P v_{pp} \\ &= \text{tr}(V) \end{aligned} \quad (32)$$

で与えられる。一方、第 m 主成分の分散 $\sigma_{z_m}^2$ は共分散行列 V の m 番目に大きな固有値 λ_m に等しいので、合計 P 個の主成分の分散の和は

$$\begin{aligned} \sum_{p=1}^P \sigma_{z_p}^2 &= \sum_{p=1}^P \lambda_p \\ &= \text{tr}(V) \end{aligned} \quad (33)$$

で与えられる。すなわち、変数の分散の総和と主成分の分散の総和とは等しくなる。そこで、第 m 主成分が元のデータに含まれる特徴をどの程度表現しているかを示す指標として、第 m 主成分の分散が分散の総和に占める割合

$$C_m = \frac{\lambda_m}{\sum_{p=1}^P \lambda_p} = \frac{\lambda_m}{\text{tr}(V)} \quad (34)$$

を利用することができる。 C_m は寄与率 (proportion) と呼ばれ、通常%を用いて表される。また、第 m 主成分までの分散の和が分散の総和に占める割合

$$P_m = \sum_{i=1}^m C_i = \frac{\sum_{i=1}^m \lambda_i}{\text{tr}(V)} \quad (35)$$

は累積寄与率 (accumulated proportion, cumulative proportion) と呼ばれる。なお、各変数を予め分散 1 に標準化してある場合には、共分散行列 V は相関行列 R に等しくなるため、相関行列を用いて寄与率および累積寄与率を求めることもできる。

上述の累積寄与率を用いて採用する主成分数を決定する場合、通常、累積寄与率が 80%程度となるように主成分数が決められる。しかし、80%という数値に深い意味はなく、常にこの基準が妥当であるとは限らない。従って、累積寄与率を参考にしながら、主成分分析の目的や適用対象に応じて採用する主成分数を適切に決定する必要がある。

各変数の分散を予め 1 に標準化する場合、すなわち相関行列に基づいて主成分分析を行う場合には、元の変数が持つ分散よりも大きな分散を持つ主成分のみを採用するという考え方に従い、大きさが 1 以上の固有値に対応する主成分を全て採用するという基準が利用されることもある。この場合、採用すべき主成分の個数は大きさが 1 以上の固有値の数として自動的に与えられることになる。

3.2 因子負荷量

主成分分析の結果を解釈するという事は、主成分 (総合的指標) の持つ意味を解釈するという事である。主成分は各変数の線形結合として与えられるので、その解釈のためには、主成分と各変数との相関を把握することによって主成分に強く影響を及ぼす変数を特定することが有効である。このための指標が因子負荷量 (factor loading) であり、主成分と変数との相関係数として定義される。すなわち、第 m 主成分 z_m と p 番目の変数 x_p との間の因子負荷量は

$$r_{z_m x_p} = \frac{\sigma_{z_m x_p}^2}{\sigma_{z_m} \sigma_{x_p}} \quad (36)$$

で与えられる．ここで， $\sigma_{z_m}, \sigma_{x_p}$ はそれぞれ z_m, x_p の標準偏差であり， $\sigma_{z_m x_p}^2$ は共分散を表す．いま，サンプル数が N であるとする， z_m, x_p の要素はそれぞれ

$$z_m : \begin{pmatrix} t_{1m} \\ t_{2m} \\ \vdots \\ t_{Nm} \end{pmatrix} = \mathbf{X} \mathbf{w}_m \quad (37)$$

$$x_p : \begin{pmatrix} x_{1p} \\ x_{2p} \\ \vdots \\ x_{Np} \end{pmatrix} = \mathbf{X} \mathbf{c}^{(p)} \quad (38)$$

で与えられる．ここで， $\mathbf{c}^{(p)}$ は第 p 列のみを取り出すためのベクトルであり，

$$\mathbf{c}^{(p)} = \begin{pmatrix} 0^{(1)} \\ \vdots \\ 1^{(p)} \\ \vdots \\ 0^{(P)} \end{pmatrix} \quad (39)$$

とする．ただし， $\mathbf{c}^{(p)}$ の各要素の右肩にある () 内の番号は行番号を表している．このとき， z_m と x_p の共分散は

$$\begin{aligned} \sigma_{z_m x_p}^2 &= \frac{1}{N-1} (\mathbf{X} \mathbf{w}_m)^T \mathbf{X} \mathbf{c}^{(p)} \\ &= \frac{1}{N-1} \mathbf{w}_m^T \mathbf{X}^T \mathbf{X} \mathbf{c}^{(p)} \\ &= (\mathbf{V} \mathbf{w}_m)^T \mathbf{c}^{(p)} \\ &= \lambda_m \mathbf{w}_m^T \mathbf{c}^{(p)} \\ &= \lambda_m w_{pm} \end{aligned} \quad (40)$$

となる．従って，因子負荷量は， $\sigma_{z_m}^2 = \lambda_m$ より

$$\begin{aligned} r_{z_m x_p} &= \frac{\lambda_m w_{pm}}{\sqrt{\lambda_m} \sigma_{x_p}} \\ &= \frac{\sqrt{\lambda_m} w_{pm}}{\sigma_{x_p}} \end{aligned} \quad (41)$$

で与えられる．なお，各変数の分散が予め 1 に標準化されている場合には，因子負荷量は

$$r_{z_m x_p} = \sqrt{\lambda_m} w_{pm} \quad (42)$$

となる．

4 主成分分析の行列表現

4.1 主成分得点の計算

P 個の変数について N 個のサンプルを含む，適当に標準化されたデータ行列 \mathbf{X} があり，その共分散行列の固有値が大きい方から順に $\lambda_i (i = 1, 2, \dots, P)$ で与えられているとする．共分散行列は非負定値行列で

あるから，その固有値は必ずゼロ以上の値をとる．いま，

$$\lambda_i \begin{cases} > 0, & \text{for } i = 1, 2, \dots, r \\ = 0, & \text{for } i = r + 1, r + 2, \dots, P \end{cases} \quad (43)$$

であるとする． $r + 1$ 番目以降の固有値はいずれもゼロであるから，それらに対応する主成分の分散もゼロとなる．すなわち， $r + 1$ 番目以降の主成分の方向にはデータは分散を持たない．これは，元々 P 次元空間に含まれているサンプル $\{x_n\} (n = 1, 2, \dots, N)$ が実は高々 r 次元部分空間に含まれることを意味する．従って，実際に得られる主成分の数はゼロでない固有値の数 (r 個) に等しくなり，ゼロでない固有値に対応する規格化された固有ベクトル $\{w_i\} (i = 1, 2, \dots, r)$ が r 次元部分空間の正規直交基底となる．

この場合， i 番目に大きな固有値に対応する固有ベクトル w_i が第 i 主成分の結合係数となるので，第 i 主成分 z_i は

$$z_i = \sum_{p=1}^P w_{pi} x_p \quad (44)$$

となり，第 i 主成分得点 t_i は

$$t_i = X w_i \quad (45)$$

で与えられる．ここで，第 $M (M \leq P)$ 主成分までを考慮し，主成分得点からなる行列 (score matrix) を

$$T = \begin{pmatrix} t_1 & t_2 & \dots & t_M \end{pmatrix} \quad (46)$$

とおき，さらに，主成分の結合係数からなる行列 (loading matrix) を

$$P = \begin{pmatrix} w_1 & w_2 & \dots & w_M \end{pmatrix} \quad (47)$$

とおくと，

$$T = X P \quad (48)$$

を得る．

新しく得られたサンプルが既に得られているサンプルと比較してどのような特徴を持つかを知りたい場合，既存のサンプルから得られる主成分に基づいて新しいサンプルの主成分得点を求め，各サンプルの主成分得点を比較することが有効である．この際，主成分得点の計算は以下の手順で行われる．

1. 既存のサンプルの loading matrix P を求める．
2. 既存のサンプルの主成分得点を Eq.(48) を用いて計算する．
3. 新しく得られたサンプル x_{new} の主成分得点 t_{new} を次式を用いて計算する．

$$t_{new} = x_{new} P \quad (49)$$

既存のサンプルの主成分得点に対する新しいサンプルの主成分得点の特徴 (大きい，平均値 0 に近いなど) を各主成分について検討することにより，新しいサンプルの特徴を把握することができる．例えば，Sec.1 で述べた身体検査の例では，第 1 主成分が「体の大きさ」を，第 2 主成分が「肥満の程度」を表しているので，新しいサンプルの第 1 主成分得点が大きく，第 2 主成分得点が小さければ，その人は体が大きいけど痩せていることがわかる．この例は高々 2 変数の場合であるので，主成分による特徴抽出の効果はほとんど無いに等しい．しかし，プロセスデータのような多変数データを扱う場合には，生データから特徴を把握することは一般に困難であり，主成分分析による特徴抽出が威力を発揮するであろう．実際に，運転状態の監視や異常検出，さらには異常診断といった分野に主成分分析が利用されている．

4.2 データの再構築

第 1 主成分得点は P 次元空間内の各サンプル $x_n (n = 1, 2, \dots, N)$ を第 1 主成分を表す z_1 軸に射影して得られるベクトルの長さであるから, あるサンプル x_n を第 1 主成分のみで代表させる場合, その予測値 \hat{x}_n は

$$\begin{aligned}\hat{x}_n &= t_{n1} \mathbf{w}_1^T \\ &= x_n \mathbf{w}_1 \mathbf{w}_1^T\end{aligned}\quad (50)$$

で与えられる. ここで, t_{n1} はサンプル x_n の第 1 主成分得点, \mathbf{w}_1 は第 1 主成分の結合係数である. サンプル x_n を第 M 主成分までの全てを用いて表現する場合には, その予測値 \hat{x}_n は各主成分方向の成分の線形結合として

$$\hat{x}_n = \sum_{m=1}^M t_{nm} \mathbf{w}_m^T \quad (51)$$

で与えられる. 簡単な例として, $P = 3, M = 2$ の場合にデータが再構築される様子を Fig.2 に示しておく. Eq.(51) は loading matrix P を用いて

$$\begin{aligned}\hat{x}_n &= \begin{pmatrix} t_{n1} & t_{n2} & \cdots & t_{nM} \end{pmatrix} P^T \\ &= x_n P P^T\end{aligned}\quad (52)$$

と表すこともできる. これより, N 個のサンプル全体については,

$$\hat{X} = X P P^T \quad (53)$$

が成り立つ. 主成分数 M が共分散行列 (あるいは相関行列) のゼロでない固有値の数に等しい場合には, \hat{X} は元のデータ行列 X と正確に一致する. しかし, それ以外の場合には \hat{X} は X と一致しない. このとき, 元データと再構築データの残差行列 E は

$$\begin{aligned}E &= X - \hat{X} \\ &= X - T P^T \\ &= X - X P P^T \\ &= X (I - P P^T)\end{aligned}\quad (54)$$

となる.

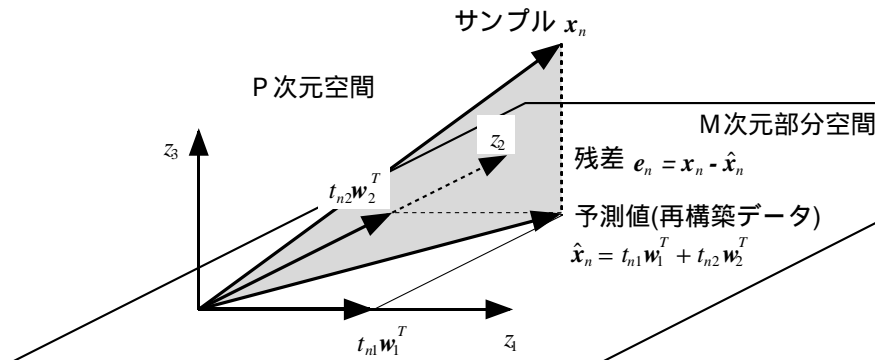


Fig. 2 データの再構築 ($P = 3, M = 2$ の場合)

4.3 主成分分析と特異値分解

ここまで、主成分が固有値や固有ベクトルと深く関係していることをみてきた。その関係とは、適当に標準化されたデータ行列が与えられたとき、その共分散行列の固有ベクトルが主成分の結合係数になり、固有値が主成分の分散に等しくなるというものであった。

P 個の変数について N 個のサンプルからなるデータ行列 X が与えられている場合、その共分散行列の i 番目に大きな固有値を λ_i 、それに対応する規格化された固有ベクトルを w_i とすると、次式が成り立つ。

$$\frac{1}{N-1} X^T X w_i = \lambda_i w_i \quad (i = 1, 2, \dots, P) \quad (55)$$

ここで、 $(N-1)\lambda_i$ を改めて λ_i とおくと、

$$X^T X w_i = \lambda_i w_i \quad (i = 1, 2, \dots, P) \quad (56)$$

を得る。この場合、 λ_i は行列 $X^T X$ の i 番目に大きな固有値、 w_i はそれに対応する固有ベクトルとなる。ここで、loading matrix を

$$P = \begin{pmatrix} w_1 & w_2 & \dots & w_P \end{pmatrix} \quad (57)$$

とおくと、 P が正規直交行列であることから、

$$X = X P P^T \quad (58)$$

が成り立つ。さらに、

$$\sigma_i = \sqrt{\lambda_i} \quad (i = 1, 2, \dots, P) \quad (59)$$

$$\Sigma = \begin{pmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_P \end{pmatrix} \quad (60)$$

$$u_i = \frac{1}{\sigma_i} X w_i \quad (i = 1, 2, \dots, P) \quad (61)$$

$$U = \begin{pmatrix} u_1 & u_2 & \dots & u_P \end{pmatrix} \quad (62)$$

とおくと、

$$X P = U \Sigma \quad (63)$$

が成り立つので、これを Eq.(58) に代入すると、最終的に

$$X = U \Sigma P^T \quad (64)$$

を得る。ここで、行列 Σ は対角行列であり、行列 U および P は正規直交行列であることから、Eq.(64) は、行列 X の特異値分解に他ならないことがわかる。なお、行列 P が正規直交行列であることは、その列が規格化された固有ベクトルであり、Eq.(22) を満たしていることから明らかである。一方、行列 U が正規直交行列であることは、行列 $U^T U$ の各要素を考えることにより容易に示すことができる。すなわち、行列 $U^T U$ の (i, j) 要素を $(U^T U)_{ij}$ とすると、

$$\begin{aligned} (U^T U)_{ij} &= \frac{1}{\sigma_i} (X w_i)^T \cdot \frac{1}{\sigma_j} X w_j \\ &= \frac{1}{\sigma_i \sigma_j} w_i^T X^T X w_j \end{aligned}$$

$$\begin{aligned}
&= \frac{\lambda_j}{\sigma_i \sigma_j} \mathbf{w}_i^T \mathbf{w}_j \\
&= \begin{cases} 1 & (i \neq j) \\ 0 & (i = j) \end{cases}
\end{aligned} \tag{65}$$

となるので，行列 U は正規直交行列である．

行列 X の特異値分解である Eq.(64) は，ゼロでない固有値の数を r とすると，

$$\mathbf{X} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{w}_i^T \tag{66}$$

と表現することもできる．さらに，

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases} \tag{67}$$

$$\mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1 & (i = j) \\ 0 & (i \neq j) \end{cases} \tag{68}$$

であることから，

$$\mathbf{X}^T \mathbf{X} = \sum_{i=1}^r \lambda_i \mathbf{w}_i \mathbf{w}_i^T \tag{69}$$

$$\mathbf{X} \mathbf{X}^T = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^T \tag{70}$$

となる．これらの式はそれぞれ非負定値行列 $\mathbf{X}^T \mathbf{X}$ および $\mathbf{X} \mathbf{X}^T$ のスペクトル分解と呼ばれるものであり， $\mathbf{X}^T \mathbf{X}$ と $\mathbf{X} \mathbf{X}^T$ の固有値は等しくなることを示している．これらのスペクトル分解を行列を用いて表現すると，

$$\mathbf{X}^T \mathbf{X} = \mathbf{P} \mathbf{\Sigma}^2 \mathbf{P}^T \tag{71}$$

$$\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^T \tag{72}$$

となる．

以上より，主成分分析は特異値分解およびスペクトル分解と密接に関係していることがわかる．その関係を簡単にまとめると，

1. 特異値分解に基づく場合

第 i 主成分の結合係数 w_i は，データ行列 X の i 番目に大きな特異値 σ_i に対応する右特異ベクトルとして与えられる．

2. スペクトル分解に基づく場合

第 i 主成分の結合係数 w_i は，行列 $\mathbf{X}^T \mathbf{X}$ の i 番目に大きな固有値 λ_i に対応する固有ベクトルとして与えられる．

となる．