

因子分析

Ver.2022.09.04

愛媛大学

データサイエンスセンター (CDSE)

理工学研究科 / 理学部

まつうら まさや

松浦 真也

2020年4月設立



CDSE

Center for Data Science, Ehime University



回帰分析と主成分分析（再考）



CDSE

Center for Data Science, Ehime University



重回帰分析（再考）

体脂肪率のデータ 出典：The Data And Story LibraryのBodyfat
<https://dasl.datadescription.com/datafile/bodyfat/>

人物 i	体重 x_{i1}	身長 x_{i2}	体脂肪率 y_i
1	154.25	67.75	12.3
2	173.25	72.25	6.1
\vdots	\vdots	\vdots	\vdots
n	207.5	70	31.9

重回帰モデル

人物に依存しない未知の定数（これを推定）

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + e_i$$

人物 i の
体脂肪率

人物 i の
体重

人物 i の
身長

残差

$(1 \leq i \leq n)$

残差平方和 $S_e = \sum_{i=1}^n e_i^2$ が最小になるように $\beta_0, \beta_1, \beta_2$ を定める

重回帰分析（再考）

回帰モデル $y_i = \beta_0 + \sum_{j=1}^r \beta_j x_{ij} + e_i \quad (1 \leq i \leq n)$

規格化 $z_i = \sum_{j=1}^r b_j f_{ij} + e_i$ 平均0、分散1

$$\begin{cases} y_i & \rightarrow & z_i \\ x_{ij} & \rightarrow & f_{ij} \end{cases}$$

目的変数を m 種類に $z_{ik} = \sum_{j=1}^r b_{jk} f_{ij} + e_{ik} \quad (1 \leq k \leq m)$

$$\begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} = \begin{bmatrix} f_{11} & \cdots & f_{1r} \\ \vdots & \cdots & \vdots \\ f_{n1} & \cdots & f_{nr} \end{bmatrix} \begin{bmatrix} b_{11} & \cdots & b_{1m} \\ \vdots & \cdots & \vdots \\ b_{r1} & \cdots & b_{rm} \end{bmatrix} + \begin{bmatrix} e_{11} & \cdots & e_{1m} \\ \vdots & \cdots & \vdots \\ e_{n1} & \cdots & e_{nm} \end{bmatrix}$$

重回帰分析（再考）

$$\begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} = \begin{bmatrix} f_{11} & \cdots & f_{1r} \\ \vdots & \cdots & \vdots \\ f_{n1} & \cdots & f_{nr} \end{bmatrix} \begin{bmatrix} b_{11} & \cdots & b_{1m} \\ \vdots & \cdots & \vdots \\ b_{r1} & \cdots & b_{rm} \end{bmatrix} + \begin{bmatrix} e_{11} & \cdots & e_{1m} \\ \vdots & \cdots & \vdots \\ e_{n1} & \cdots & e_{nm} \end{bmatrix}$$

目的変数

残差

説明変数

$$Z = \begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix}$$

$$E = \begin{bmatrix} e_{11} & \cdots & e_{1m} \\ \vdots & \cdots & \vdots \\ e_{n1} & \cdots & e_{nm} \end{bmatrix}$$

$$F = \begin{bmatrix} f_{11} & \cdots & f_{1r} \\ \vdots & \cdots & \vdots \\ f_{n1} & \cdots & f_{nr} \end{bmatrix}$$

n 個(人)
の個体

m 種類の目的変数

m 種類の目的変数

r 種類の説明変数

係数

$$B = \begin{bmatrix} b_{11} & \cdots & b_{1m} \\ \vdots & \cdots & \vdots \\ b_{r1} & \cdots & b_{rm} \end{bmatrix}$$

(個体に依存しない)

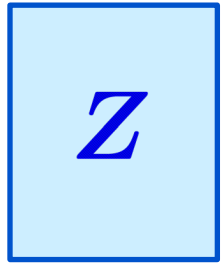
回帰式の行列表現

$$Z = FB + E$$

重回帰分析（再考）

$$\begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} = \begin{bmatrix} f_{11} & \cdots & f_{1r} \\ \vdots & \cdots & \vdots \\ f_{n1} & \cdots & f_{nr} \end{bmatrix} \begin{bmatrix} b_{11} & \cdots & b_{1m} \\ \vdots & \cdots & \vdots \\ b_{r1} & \cdots & b_{rm} \end{bmatrix} + \begin{bmatrix} e_{11} & \cdots & e_{1m} \\ \vdots & \cdots & \vdots \\ e_{n1} & \cdots & e_{nm} \end{bmatrix}$$

$n \times m$

A light blue square representing matrix Z.

=

$n \times r$

A light blue vertical rectangle representing matrix F.

$r \times m$

A light pink horizontal rectangle representing matrix B.

+

$n \times m$

A light pink square representing matrix E.

目的変数 (既知)

説明変数 (既知)

残差 (最小化)



$n \times m$

A light blue square representing matrix Z.

\approx

$n \times r$

A light blue vertical rectangle representing matrix F.

$r \times m$

A light pink horizontal rectangle representing matrix B.

※ rank $FB \leq r$ に注意

主成分分析（再考）

アヤメのデータ

出典：R. A. Fisher, "The use of multiple measurements in taxonomic problems", Annals of Eugenics, Vol. 7, No. 2, 179–188, 1936.

No i	萼片長 x_{i1}	萼片幅 x_{i2}	...	花弁幅 x_{im}	品種
1	5.1	3.5	...	0.2	setosa
2	4.9	3.0	...	0.2	setosa
⋮	⋮	⋮	⋮	⋮	⋮
n	5.9	3.0	...	1.8	virginica

$$k = 1, 2, \dots, m$$

標本平均

$$\bar{x}_k = \frac{1}{n} \sum_{i=1}^n x_{ik}$$

標本標準偏差

$$s_k = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}$$

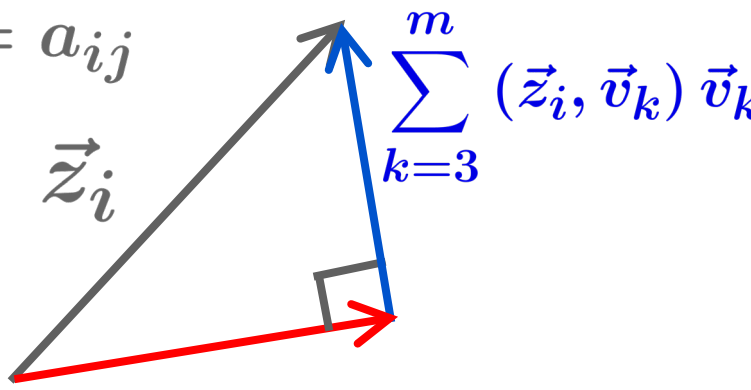
規格化

$$z_{ik} = \frac{x_{ik} - \bar{x}_k}{s_k}$$

(平均 0、分散 1)

主成分分析 (再考)

$$\vec{z}_i = \begin{bmatrix} z_{i1} \\ \vdots \\ z_{im} \end{bmatrix} \quad \vec{v}_1, \dots, \vec{v}_m \text{ を } \mathbb{R}^m \text{ の正規直交基底とすると}$$
$$\vec{z}_i = \sum_{k=1}^m a_{ik} \vec{v}_k \quad (a_{ik} \in \mathbb{R}) \text{ と書ける}$$

$$(\vec{z}_i, \vec{v}_j) = \sum_{k=1}^m a_{ik} (\vec{v}_k, \vec{v}_j) = a_{ij}$$


$$\vec{z}_i = \sum_{k=1}^m (\vec{z}_i, \vec{v}_k) \vec{v}_k$$

$$(\vec{z}_i, \vec{v}_1) \vec{v}_1 + (\vec{z}_i, \vec{v}_2) \vec{v}_2$$

$$|\vec{z}_i|^2 = (\vec{z}_i, \vec{v}_1)^2 + (\vec{z}_i, \vec{v}_2)^2 + \sum_{k=3}^m (\vec{z}_i, \vec{v}_k)^2$$

主成分分析 (再考)

$$\vec{z}_i = \begin{bmatrix} z_{i1} \\ \vdots \\ z_{im} \end{bmatrix} \quad \vec{v}_1, \dots, \vec{v}_m \text{ を } \mathbb{R}^m \text{ の正規直交基底とすると}$$

$$\vec{z}_i = \sum_{k=1}^m a_{ik} \vec{v}_k \quad (a_{ik} \in \mathbb{R})$$

$$= \sum_{k=1}^r a_{ik} \vec{v}_k + \sum_{k=r+1}^m a_{ik} \vec{v}_k$$

r 次元

$m-r$ 次元

m 次元 \rightarrow r 次元

$$\begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix}$$

$$\begin{bmatrix} {}^t \vec{z}_1 \\ \vdots \\ {}^t \vec{z}_n \end{bmatrix} = \begin{bmatrix} a_{11} & \cdots & a_{1r} \\ \vdots & \cdots & \vdots \\ a_{n1} & \cdots & a_{nr} \end{bmatrix} \begin{bmatrix} {}^t \vec{v}_1 \\ \vdots \\ {}^t \vec{v}_r \end{bmatrix} + \begin{bmatrix} a_{1\ r+1} & \cdots & a_{1m} \\ \vdots & \cdots & \vdots \\ a_{n\ r+1} & \cdots & a_{nm} \end{bmatrix} \begin{bmatrix} {}^t \vec{v}_{r+1} \\ \vdots \\ {}^t \vec{v}_m \end{bmatrix}$$

主成分分析 (再考)

$$\begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} = \begin{bmatrix} a_{11} & \cdots & a_{1r} \\ \vdots & \cdots & \vdots \\ a_{n1} & \cdots & a_{nr} \end{bmatrix} \begin{bmatrix} {}^t\vec{v}_1 \\ \vdots \\ {}^t\vec{v}_r \end{bmatrix} + \begin{bmatrix} a_{1\ r+1} & \cdots & a_{1m} \\ \vdots & \cdots & \vdots \\ a_{n\ r+1} & \cdots & a_{nm} \end{bmatrix} \begin{bmatrix} {}^t\vec{v}_{r+1} \\ \vdots \\ {}^t\vec{v}_m \end{bmatrix}$$

◆ 用語・記号を回帰分析に合わせる

「目的変数」

$$Z = \begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix}$$

「残差」

$$E = \begin{bmatrix} a_{1\ r+1} & \cdots & a_{1m} \\ \vdots & \cdots & \vdots \\ a_{n\ r+1} & \cdots & a_{nm} \end{bmatrix} \begin{bmatrix} {}^t\vec{v}_{r+1} \\ \vdots \\ {}^t\vec{v}_m \end{bmatrix}$$

「説明変数」

$$F = \begin{bmatrix} a_{11} & \cdots & a_{1r} \\ \vdots & \cdots & \vdots \\ a_{n1} & \cdots & a_{nr} \end{bmatrix}$$

「係数」(個体に依存しない)

$$B = \begin{bmatrix} {}^t\vec{v}_1 \\ \vdots \\ {}^t\vec{v}_r \end{bmatrix}$$

主成分分析の行列表現

$$Z = FB + E$$

主成分分析 (再考)

$$\begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} = \begin{bmatrix} a_{11} & \cdots & a_{1r} \\ \vdots & \cdots & \vdots \\ a_{n1} & \cdots & a_{nr} \end{bmatrix} \begin{bmatrix} {}^t\vec{v}_1 \\ \vdots \\ {}^t\vec{v}_r \end{bmatrix} + \begin{bmatrix} a_{1r+1} & \cdots & a_{1m} \\ \vdots & \cdots & \vdots \\ a_{nr+1} & \cdots & a_{nm} \end{bmatrix} \begin{bmatrix} {}^t\vec{v}_{r+1} \\ \vdots \\ {}^t\vec{v}_m \end{bmatrix}$$

$n \times m$



$n \times r$

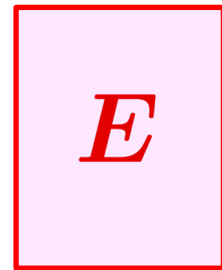


$r \times m$



係数(推測)

$n \times m$



=

+

目的変数(既知)

説明変数(推測)

残差(最小化)

観測変数

主成分得点

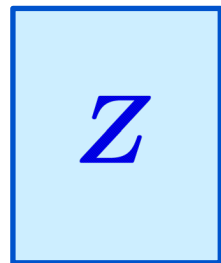
主成分



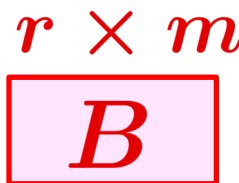
近似

$n \times m$

$n \times r$



\approx



$r \times m$

※ rank $BF \leq r$
に注意

因子分析



CDSE

Center for Data Science, Ehime University



因子分析（主成分分析との比較）

◆ 主成分分析

$$\begin{array}{c} n \times m \\ \boxed{Z} \end{array} = \begin{array}{c} n \times r \\ \boxed{F} \end{array} \begin{array}{c} r \times m \\ \boxed{B} \\ \text{係数(推測)} \end{array} + \begin{array}{c} n \times m \\ \boxed{E} \\ \text{残差(最小化)} \end{array}$$

目的変数(既知) 説明変数(推測)

◆ 因子分析

$$\begin{array}{c} n \times m \\ \boxed{Z} \end{array} = \begin{array}{c} n \times r \\ \boxed{F} \end{array} \begin{array}{c} r \times m \\ \boxed{B} \\ \text{係数(推測)} \end{array} + \begin{array}{c} n \times m \\ \boxed{E} \\ \text{残差(最小化)} \end{array} + \begin{array}{c} n \times m \\ \boxed{W} \\ \text{ノイズ} \\ \text{(無相関)} \end{array}$$

目的変数(既知) 説明変数(推測)

非構造的

※ ${}^t W W$ は対角行列

因子分析 (モデル)

※ よくある説明：下記の「モデル」から出発

※ ${}^t W W$ は対角行列

$n \times m$

$$Z$$

=

$n \times r$

$$F$$

$r \times m$

$$B$$

係数 (推測)

+

$n \times m$

$$W$$

目的変数 (既知)

説明変数 (推測)

ノイズ (無相関)

観測変数

因子得点

因子負荷量

独自因子

※ 実際は、上記の等号が厳密に成立する可能性は低い。
実際に解析する際は、下記の「残差」を切り捨て近似

切り捨て
近似

$n \times m$

$$Z$$

=

$n \times r$

$$F$$

$r \times m$

$$B$$

係数 (推測)

+

$n \times m$

$$E$$

+

$n \times m$

$$W$$

目的変数 (既知)

説明変数 (推測)

残差 (最小化)

ノイズ

因子分析（計算手順）

$$\mathbf{Z} = \mathbf{F} \mathbf{B} + \mathbf{E} + \mathbf{W}$$

観測変数 = 因子得点 × 因子負荷量 + 残差 (最小化) + 独自因子

※ ${}^t \mathbf{W} \mathbf{W}$ は対角行列

◆ 計算手順（イメージ）

- (1) 独自因子 \mathbf{W} を推測して、観測変数 \mathbf{Z} から除去
- (2) $\mathbf{Z} - \mathbf{W}$ に対し、主成分分析と同様の分解を実行
- (3) \mathbf{B} , \mathbf{F} を具体的に求める

※ 独自因子の推測が最大の難関（完璧な方法はない）

因子分析（独自因子の推定）

$$\mathbf{Z} = \mathbf{F} \mathbf{B} + \mathbf{E} + \mathbf{W}$$

観測変数 = 因子得点 × 因子負荷量 + 残差 (最小化) + 独自因子

◆ 独自因子の推定法の例（イメージ）

- (1) 観測変数 Z_i を、残りの観測変数 Z_k ($k \neq i$) に回帰
- (2) 回帰の残差を W_i とし、これを並べた行列を W とする
- (3) 一般に、 ${}^t W W$ は対角行列にならないが、非対角成分は無視

※ ${}^t W W$ は対角行列

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \cdots & \mathbf{Z}_m \\ z_{11} & \cdots & z_{1m} \\ \vdots & \cdots & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} \quad \mathbf{W} = [W_1, \dots, W_m]$$

因子分析（回転不定性）

T を直交行列とすると、 $F^t T T B = F B$ なので、

$$\begin{array}{c} \boxed{Z} \\ \text{観測変数} \end{array} = \begin{array}{c} \overbrace{\boxed{F}}^{\tilde{F}} \\ \text{因子得点} \end{array} \begin{array}{c} \boxed{tT} \\ \boxed{T} \end{array} \begin{array}{c} \overbrace{\boxed{B}}^{\tilde{B}} \\ \text{因子負荷量} \end{array} + \begin{array}{c} \boxed{E} \\ \text{残差} \\ \text{(最小化)} \end{array} + \begin{array}{c} \boxed{W} \\ \text{独自因子} \end{array}$$

$$\begin{cases} F & \Rightarrow & \tilde{F} = F^t T \\ B & \Rightarrow & \tilde{B} = T B \end{cases} \quad \text{因子を変換（例えば回転）可能}$$

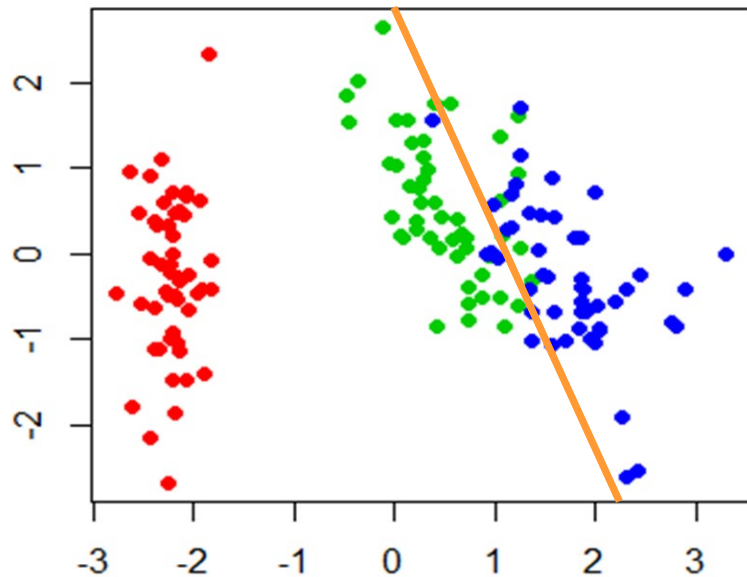
※ T を上手く選ぶことで、因子の解釈をしやすくする

主成分分析の回転不定性

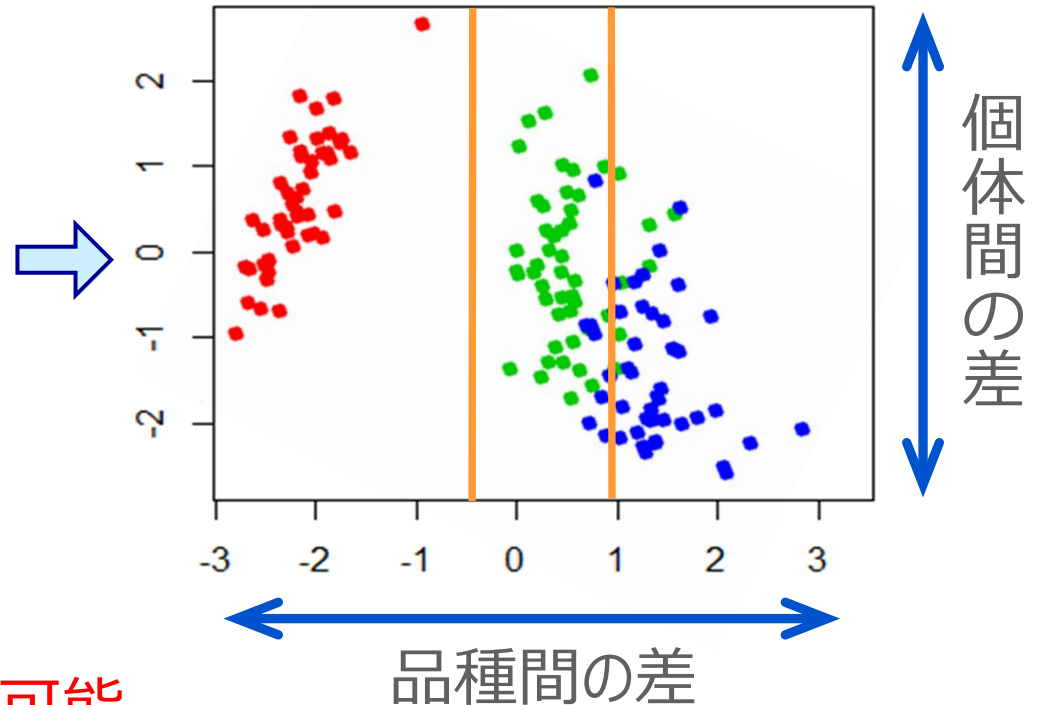
【例】アヤメの花のデータを要約したい (4次元 → 2次元)

出典 : R. A. Fisher, "The use of multiple measurements in taxonomic problems", Annals of Eugenics, Vol. 7, No. 2, 179-188, 1936.

主成分分析の結果



主成分分析の結果 (回転後)



主成分分析でも要約を回転可能

「因子分析は主成分分析と違い、回転不定性がある」は誤解の元

回帰・主成分・因子分析 (m 次元 \rightarrow r 次元)

回帰分析

$$\begin{matrix} n \times m \\ \boxed{Z} \\ \text{目的変数} \end{matrix} = \begin{matrix} n \times r \\ \boxed{F} \\ \text{説明変数} \end{matrix} \begin{matrix} r \times m \\ \boxed{B} \\ \text{係数} \end{matrix} + \begin{matrix} n \times m \\ \boxed{E} \\ \text{残差(最小化)} \end{matrix}$$

主成分分析

$$\begin{matrix} n \times m \\ \boxed{Z} \\ \text{観測変数} \end{matrix} = \begin{matrix} n \times r \\ \boxed{F} \\ \text{主成分得点} \end{matrix} \begin{matrix} r \times m \\ \boxed{B} \\ \text{主成分} \end{matrix} + \begin{matrix} n \times m \\ \boxed{E} \\ \text{残差(最小化)} \end{matrix}$$

因子分析

$$\begin{matrix} n \times m \\ \boxed{Z} \\ \text{観測変数} \end{matrix} = \begin{matrix} n \times r \\ \boxed{F} \\ \text{因子得点} \end{matrix} \begin{matrix} r \times m \\ \boxed{B} \\ \text{因子負荷量} \end{matrix} + \begin{matrix} n \times m \\ \boxed{E} \\ \text{残差(最小化)} \end{matrix} + \begin{matrix} n \times m \\ \boxed{W} \\ \text{独自因子} \end{matrix}$$

因子分析（解析例）

【例】アヤメの花のデータは1因子で説明可能？（4次元 → 1次元）

出典：R. A. Fisher, "The use of multiple measurements in taxonomic problems", *Annals of Eugenics*, Vol. 7, No. 2, 179–188, 1936.

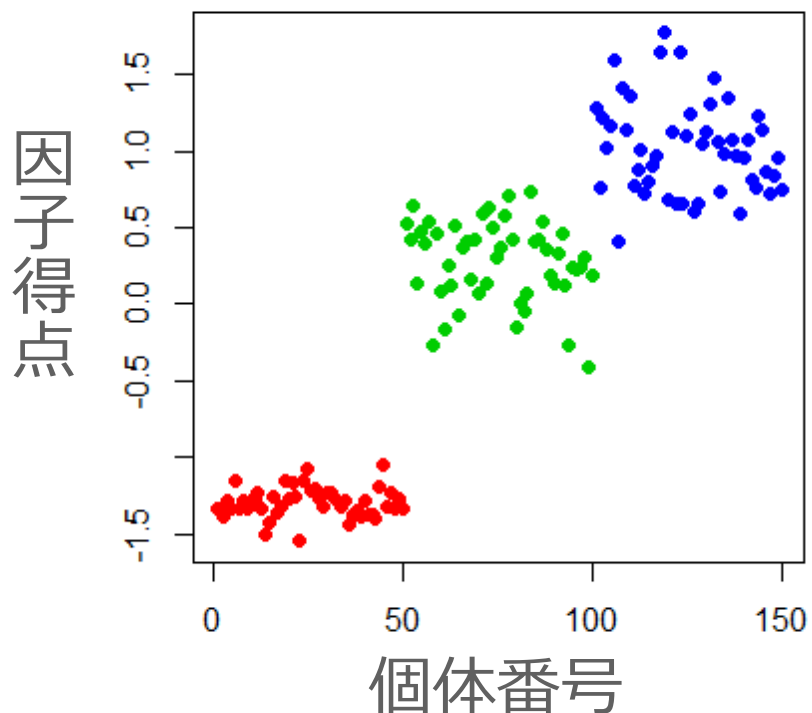
No	萼片長	萼片幅	花弁長	花弁幅	品種
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
⋮	⋮	⋮	⋮	⋮	⋮
51	7.0	3.2	4.7	1.4	versicolor
52	6.4	3.2	4.5	1.5	versicolor
⋮	⋮	⋮	⋮	⋮	⋮
101	6.3	3.3	6.0	2.5	virginica
102	5.8	2.7	5.1	1.9	virginica
⋮	⋮	⋮	⋮	⋮	⋮
150	5.9	3.0	5.1	1.8	virginica

因子分析（解析例）

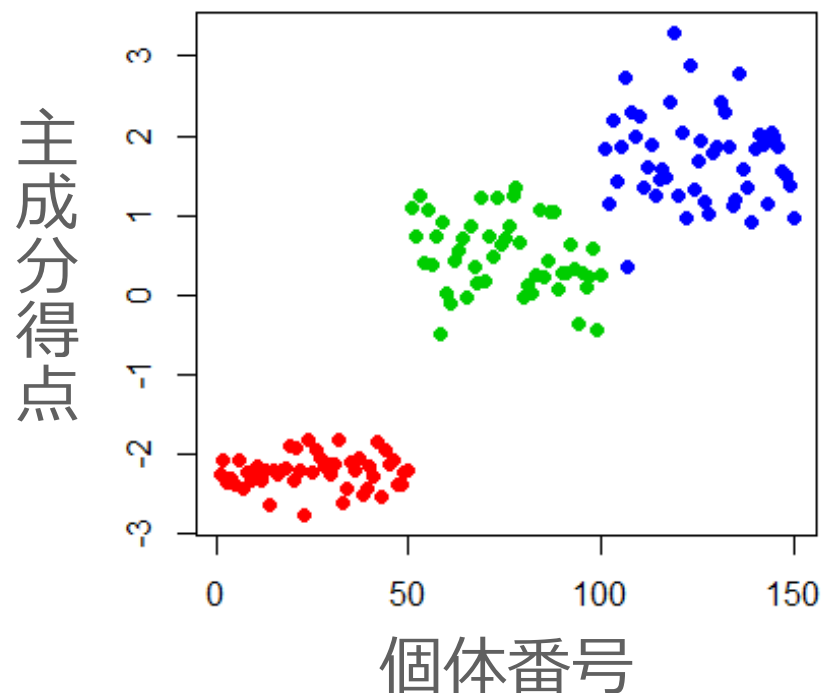
【例】アヤメの花のデータは1因子で説明可能？（4次元 → 1次元）

出典：R. A. Fisher, "The use of multiple measurements in taxonomic problems", Annals of Eugenics, Vol. 7, No. 2, 179-188, 1936.

因子分析の結果



主成分分析の結果



因子分析（参考文献）

参考文献

- ・ アヤメのデータ

R. A. Fisher, "The use of multiple measurements in taxonomic problems", *Annals of Eugenics*, Vol. 7, No. 2, 179–188, 1936.

- ・ 体脂肪率のデータ

The Data And Story Library, Bodyfat,
<https://dasl.datadescription.com/datafile/bodyfat/>

（最終閲覧日 2022年9月4日）