

## ② PCA (主成分分析)

★ 2変数を扱うために考える

$$X_1 = \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{pmatrix} \quad X_2 = \begin{pmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{n2} \end{pmatrix}$$

サンプルサイズ  $n$ , 特徴ベクトル 2 のベクトル

$$u_{1i} = \frac{x_{1i} - \bar{x}_1}{\sigma_1} \quad u_{2i} = \frac{x_{2i} - \bar{x}_2}{\sigma_2} \quad \text{※ 標準化}$$

ベクトル  $u_1, u_2$  の線形結合を用いたベクトル  $Z$  を定義する

$$Z_i = a_1 u_{1i} + a_2 u_{2i}$$

$$\begin{pmatrix} z_{1i} = a_1 u_{1i} + a_2 u_{2i} \end{pmatrix} \quad \dots (*)$$

★ 目的 ... データの情報を最大限残す

→  $Z$  の分散を最大化すればいい

$$V_Z = \frac{\sum_{i=1}^n (z_{1i} - \bar{z}_1)^2}{n}$$

$u_1, u_2$  が平均0と分散1のベクトル  $\bar{z}_1 = 0$  と同様

$$= \frac{\sum_{i=1}^n z_{1i}^2 - 2 \sum_{i=1}^n z_{1i} \bar{z}_1 + \sum_{i=1}^n \bar{z}_1^2}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n (a_1 u_{1i} + a_2 u_{2i})^2$$

$$= \frac{1}{n} \left\{ a_1^2 \sum_{i=1}^n u_{1i}^2 + 2a_1 a_2 \sum_{i=1}^n u_{1i} u_{2i} + a_2^2 \sum_{i=1}^n u_{2i}^2 \right\}$$

$$\text{※1 } \frac{\sigma_{u_1}^2}{\sigma_{u_1}^2} = \frac{\sum_{i=1}^n (u_{1i} - \bar{u}_1)^2}{n} = \frac{\sum_{i=1}^n (u_{1i}^2 - 2u_{1i}\bar{u}_1 + \bar{u}_1^2)}{n} = \frac{1}{n} \sum_{i=1}^n u_{1i}^2$$

$$\text{※2 } \frac{\sigma_{u_1 u_2}}{n} = \frac{\sum_{i=1}^n (u_{1i} - \bar{u}_1)(u_{2i} - \bar{u}_2)}{n} = \frac{\sum_{i=1}^n (u_{1i} u_{2i} - u_{1i} \bar{u}_2 - \bar{u}_1 u_{2i} + \bar{u}_1 \bar{u}_2)}{n}$$

$$= \frac{1}{n} \left\{ a_1^2 n + 2a_1 a_2 r_{u_1 u_2} n + a_2^2 n \right\}$$

$$= a_1^2 + a_2^2 + 2a_1 a_2 r_{u_1 u_2}$$

※2の2項和=1  
という制約を設ける

$$= 1 + 2a_1 a_2 r_{u_1 u_2} \quad \dots (***)$$

→  $(*)$  の最大問題をラグランジュ乗数法を用いて計算する

$$f(a_1, a_2, \lambda) = a_1^2 + a_2^2 + 2a_1 a_2 r_{u_1 u_2} - \lambda(a_1^2 + a_2^2 - 1)$$

$$\frac{\partial f}{\partial a_1} = 2a_1 + 2r_{u_1 u_2} a_2 - 2\lambda a_1 = 0$$

$$\frac{\partial f}{\partial a_2} = 2a_2 + 2r_{u_1 u_2} a_1 - 2\lambda a_2 = 0$$

$$\rightarrow \begin{pmatrix} 1 & r_{u_1 u_2} \\ r_{u_1 u_2} & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \lambda \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$$

※分散行列  $V$

→  $\lambda$  に対する固有方程式  $Va = \lambda a$  を解けば  
 $a_1, a_2$  の組がある固有ベクトルが求まる

Thm 固有方程式の固有値  $\lambda$  は元定義域  $Z$  の分散と等しい

[証明]

$Va = \lambda a$  に  $a^T$  を作用する

$$\rightarrow a^T Va = a^T \lambda a = \lambda(a_1^2 + a_2^2)$$

$$\rightarrow (a_1 \ a_2) \begin{pmatrix} 1 & r_{u_1 u_2} \\ r_{u_1 u_2} & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \lambda$$

$$\rightarrow \underbrace{a_1^2 + a_2^2 + 2a_1 a_2 r_{u_1 u_2}}_{V_Z} = \lambda$$

$$\rightarrow V_Z = \lambda \quad \text{証明終わり}$$

結局,  $u_1, u_2$  を用いて、生成した第1主成分は

$$Z_1 = a_1 u_1 + a_2 u_2$$

★ 実装 において

$$W = \begin{pmatrix} a_1 & a_2 \\ b_1 & b_2 \end{pmatrix} \quad X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix}$$

$$\begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix} \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \end{pmatrix}$$

$$= \begin{pmatrix} a_1 x_{11} + a_2 x_{12} & b_1 x_{11} + b_2 x_{12} \\ a_1 x_{21} + a_2 x_{22} & b_1 x_{21} + b_2 x_{22} \end{pmatrix}$$

※1-成分

※2-成分

0 第1主成分を考察。

$$Z_2 = b_1 u_1 + b_2 u_2 \text{ と } \langle u_1, u_2 \rangle \text{ を定義して求める。}$$

第1主成分を考察すると同様、 $b_1^2 + b_2^2 = 1$  という制約を用いる。

→  $Z_1$  と  $Z_2$  は 相関が 0 になる。

( $Z_1$  は 分散が最大になる  $Z_1, Z_2, \dots$  を定義する)

$$V[Z_1], V[Z_2] \neq 0 \text{ であるから, } \text{Cov}[Z_1, Z_2] = 0 \text{ である}$$

$$\rightarrow \frac{1}{n} \sum_{i=1}^n (Z_{1i} - \bar{Z}_1)(Z_{2i} - \bar{Z}_2) = 0$$

$$\rightarrow \frac{1}{n} \left\{ \sum_{i=1}^n Z_{1i} Z_{2i} - \underbrace{\bar{Z}_1 \sum_{i=1}^n Z_{2i}}_{0} - \underbrace{\sum_{i=1}^n Z_{1i} \bar{Z}_2}_{0} + \underbrace{\bar{Z}_1 \bar{Z}_2}_{0} \right\} = 0$$

$$\rightarrow \frac{1}{n} \sum_{i=1}^n Z_{1i} Z_{2i} = \frac{1}{n} \sum_{i=1}^n \{a_1 u_{1i} + a_2 u_{2i}\} \{b_1 u_{1i} + b_2 u_{2i}\} = 0$$

$$\rightarrow \frac{1}{n} \sum_{i=1}^n \left\{ a_1 b_1 \underbrace{u_{1i}^2}_{*1} + a_1 b_2 \underbrace{u_{1i} u_{2i}}_{*2} + a_2 b_1 \underbrace{u_{2i} u_{1i}}_{*2} + a_2 b_2 \underbrace{u_{2i}^2}_{*1} \right\} = 0$$

$$*1 \quad \frac{\sigma_{u_1}^2}{n} = \frac{\sum_{i=1}^n (u_{1i} - \bar{u}_1)^2}{n} = \frac{\sum_{i=1}^n (u_{1i}^2 - 2u_{1i}\bar{u}_1 + \bar{u}_1^2)}{n} = \frac{1}{n} \sum_{i=1}^n u_{1i}^2$$

$$*2 \quad \frac{\sigma_{u_1 u_2}}{n} = \frac{\sum_{i=1}^n (u_{1i} - \bar{u}_1)(u_{2i} - \bar{u}_2)}{n} = \frac{\sum_{i=1}^n (u_{1i} u_{2i} - u_{1i} \bar{u}_2 - \bar{u}_1 u_{2i} + \bar{u}_1 \bar{u}_2)}{n} = r_{u_1 u_2}$$

$$\rightarrow a_1 b_1 + a_1 b_2 r_{u_1 u_2} + a_2 b_1 r_{u_1 u_2} + a_2 b_2 = 0$$

$$\rightarrow (a_1, a_2) \begin{pmatrix} 1 & r_{u_1 u_2} \\ r_{u_1 u_2} & 1 \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = 0$$

$$\rightarrow a^T V b = 0 \quad * \text{ 第1主成分を求める際に得られる } V a = \lambda_1 a \text{ を用いる } \rightarrow$$

$$\rightarrow \lambda_1 a^T b = 0 \quad * \quad (V a)^T = \lambda_1 a^T \rightarrow a^T V = \lambda_1 a^T$$

つまり、制約条件  $a_1 b_1 + a_2 b_2 = 0$

$V_{Z_2} = b_1^2 + b_2^2 + 2r_{u_1 u_2} b_1 b_2$  の最大化は  $b_1^2 + b_2^2 = 1$  と  $a_1 b_1 + a_2 b_2 = 0$  を併せて行う

$$g = b_1^2 + b_2^2 + 2r_{u_1 u_2} b_1 b_2 - \lambda(b_1^2 + b_2^2 - 1) - \eta(a_1 b_1 + a_2 b_2) \text{ とする}$$

$$\frac{\partial g}{\partial b_1} = 2b_1 + 2r_{u_1 u_2} b_2 - 2\lambda b_1 - \eta a_1 = 0$$

$$\frac{\partial g}{\partial b_2} = 2b_2 + 2r_{u_1 u_2} b_1 - 2\lambda b_2 - \eta a_2 = 0$$

$$\rightarrow V b = \lambda b + \frac{\eta}{2} a \quad \left. \begin{array}{l} \\ \end{array} \right\} a^T \text{ は 圧力作用}$$

$$\rightarrow \underbrace{a^T V b}_0 = \underbrace{a^T \lambda b}_0 + \frac{\eta}{2} a^T a$$

$$\rightarrow \eta = 0$$

結局  $V b = \lambda b$  と 第1主成分  $a$  と同様