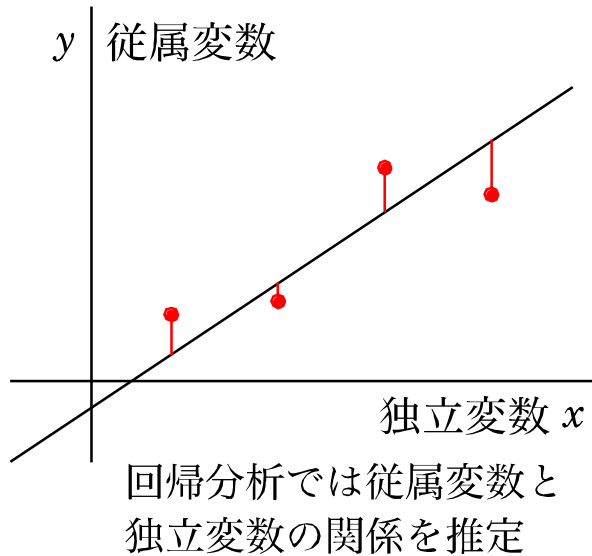


数值计算 (7)

最小二乘法

主成分分析

最小二乗法



データを直線で近似 (回帰直線)
一次回帰曲線

$$y = ax + b$$

データ (x_i, y_i) に対する残差

$$y_i - (ax_i + b)$$

残差平方和 \rightarrow 最小

$$S = \sum_{i=1}^N \{y_i - (ax_i + b)\}^2$$

2乗和を最小 \Rightarrow 偏微分 = 0

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^N (y_i - ax_i - b) x_i = 0$$

$$\frac{\partial S}{\partial b} = -2 \sum_{i=1}^N (y_i - ax_i - b) = 0$$

$$\left\{ \begin{array}{l} \left(\sum_{i=1}^N x_i \right) a + \left(\sum_{i=1}^N 1 \right) b = \sum_{i=1}^N y_i \\ \left(\sum_{i=1}^N x_i^2 \right) a + \left(\sum_{i=1}^N x_i \right) b = \sum_{i=1}^N x_i y_i \end{array} \right.$$

数值計算

| | | | | | | |
|-----|------|------|------|------|------|------|
| x | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 |
| y | 10.2 | 12.0 | 15.7 | 17.0 | 20.5 | 22.4 |

$$\left(\sum_{i=1}^6 x_i\right) = 15 \quad \left(\sum_{i=1}^6 x_i^2\right) = 55 \quad \left(\sum_{i=1}^6 y_i\right) = 97.8 \quad \left(\sum_{i=1}^6 x_i y_i\right) = 288.4$$

$$\begin{cases} 15a + 6b = 97.8 \\ 55a + 15b = 288.4 \end{cases} \quad \begin{cases} a = 2.51 \\ b = 10.03 \end{cases}$$

$$y = 2.51x + 10.03$$

多項式近似 (回帰曲線)

$$y = P_n(x) = a_0 + a_1x + a_2x^2 + \cdots + a_nx^n$$

残差 $y_i - P_n(x_i)$ 残差平方和 $S = \sum_{i=1}^N \{ y_i - P_n(x_i) \}^2$

a_j で偏微分 = 0

$$\frac{\partial S}{\partial a_j} = -2 \sum_{i=1}^N (y_i - P_n(x_i)) x_i^j = 0, \quad \sum_{i=1}^N P_n(x_i) x_i^j = \sum_{i=1}^N x_i^j y_i$$

$$\begin{array}{l} j=0 \quad \left(\sum_{i=1}^N 1 \right) a_0 + \left(\sum_{i=1}^N x_i \right) a_1 + \cdots + \left(\sum_{i=1}^N x_i^n \right) a_n = \sum_{i=1}^N y_i \\ j=1 \quad \left(\sum_{i=1}^N x_i \right) a_0 + \left(\sum_{i=1}^N x_i^2 \right) a_1 + \cdots + \left(\sum_{i=1}^N x_i^{n+1} \right) a_n = \sum_{i=1}^N x_i y_i \\ \vdots \\ j=n \quad \left(\sum_{i=1}^N x_i^n \right) a_0 + \left(\sum_{i=1}^N x_i^{n+1} \right) a_1 + \cdots + \left(\sum_{i=1}^N x_i^{2n} \right) a_n = \sum_{i=1}^N x_i^n y_i \end{array}$$

数值計算 $y = P_3(x) = a_0 + a_1x + a_2x^2 + a_3x^3$

| | | | | | | |
|-----|------|------|------|------|------|------|
| x | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 |
| y | 10.2 | 14.0 | 15.7 | 16.0 | 18.5 | 22.4 |

$$\left(\sum_{i=1}^6 1\right) = 6, \quad \left(\sum_{i=1}^6 x_i\right) = 15, \quad \left(\sum_{i=1}^6 x_i^2\right) = 55, \quad \left(\sum_{i=1}^6 x_i^3\right) = 255,$$

$$\left(\sum_{i=1}^6 x_i^4\right) = 979, \quad \left(\sum_{i=1}^6 x_i^5\right) = 4425, \quad \left(\sum_{i=1}^6 x_i^6\right) = 20515,$$

$$\sum_{i=1}^6 y_i = 96.8, \quad \sum_{i=1}^6 x_i y_i = 279.4, \quad \sum_{i=1}^6 x_i^2 y_i = 1076.8, \quad \sum_{i=1}^6 x_i^3 y_i = 4555.6$$

$$\begin{cases} 6a_0 + 15a_1 + 55a_2 + 255a_3 = 96.8 \\ 15a_0 + 55a_1 + 255a_2 + 979a_3 = 279.4 \\ 55a_0 + 255a_1 + 979a_2 + 4425a_3 = 1076.8 \\ 255a_0 + 979a_1 + 4425a_2 + 20515a_3 = 4555.6 \end{cases}$$

多項式近似 行列版

$$\begin{pmatrix} \sum_{i=1}^N 1 & \sum_{i=1}^N x_i & \cdots & \sum_{i=1}^N x_i^n \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 & \cdots & \sum_{i=1}^N x_i^{n+1} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_i^n & \sum_{i=1}^N x_i^{n+1} & \cdots & \sum_{i=1}^N x_i^{2n} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N x_i y_i \\ \vdots \\ \sum_{i=1}^N x_i^n y_i \end{pmatrix}$$

ここで,

$$\mathbf{A} = \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^n \\ 1 & x_2 & x_2^2 & \cdots & x_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \cdots & x_N^n \end{pmatrix} \quad \text{とすると,} \quad \mathbf{A}^t \mathbf{A} \mathbf{a} = \mathbf{A}^t \mathbf{y}$$

線形最小二乗法

$$y = a_1 f_1(x) + a_2 f_2(x) + \cdots + a_n f_n(x) = \sum_{k=1}^n a_k f_k(x)$$

残差平方和 $S = \sum_{i=1}^N \left(y_i - \sum_{k=1}^n a_k f_k(x_i) \right)^2$

$$\frac{\partial S}{\partial a_j} = -2 \sum_{i=1}^N \left(y_i - \sum_{k=1}^n a_k f_k(x_i) \right) f_j(x_i) = 0, \quad \sum_{k=1}^n \sum_{i=1}^N f_j(x_i) f_k(x_i) a_k = \sum_{i=1}^N f_j(x_i) y_i$$

ここで,

$$\mathbf{A} = \begin{pmatrix} f_1(x_1) & f_2(x_1) & \cdots & f_n(x_1) \\ f_1(x_2) & f_2(x_2) & \cdots & f_n(x_2) \\ \vdots & \vdots & & \vdots \\ f_1(x_N) & f_2(x_N) & \cdots & f_n(x_N) \end{pmatrix} \text{とすると,} \quad \mathbf{A}^t \mathbf{A} \mathbf{a} = \mathbf{A}^t \mathbf{y}$$

数値計算

$$f_1(x) = 1, \quad f_2(x) = x, \quad f_3(x) = \sin\left(\frac{2\pi}{5}x\right)$$

| | | | | | | |
|-----|------|------|------|------|------|------|
| x | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 |
| y | 10.2 | 14.0 | 15.7 | 16.0 | 18.5 | 22.4 |

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0.00 \\ 1 & 1 & 0.95 \\ 1 & 2 & 0.58 \\ 1 & 3 & -0.58 \\ 1 & 4 & -0.95 \\ 1 & 5 & 0.00 \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 10.2 \\ 14.0 \\ 15.7 \\ 16.0 \\ 18.5 \\ 22.4 \end{pmatrix}$$

$$\mathbf{A}^t \mathbf{A} \mathbf{a} = \mathbf{A}^t \mathbf{y} \quad \text{を解くと,} \quad \begin{cases} a_0 = 10.01 \\ a_1 = 2.45 \\ a_2 = 1.59 \end{cases}$$

よって,

$$y = 10.01 + 2.45x + 1.59\sin\left(\frac{2\pi}{5}x\right)$$

多次元の場合

独立変数 x_1, x_2, \dots, x_n と従属変数 y

j 番目のデータを $(x_{j1}, x_{j2}, \dots, x_{jn}, y_j)$ と表すと ($1 \leq j \leq N$)

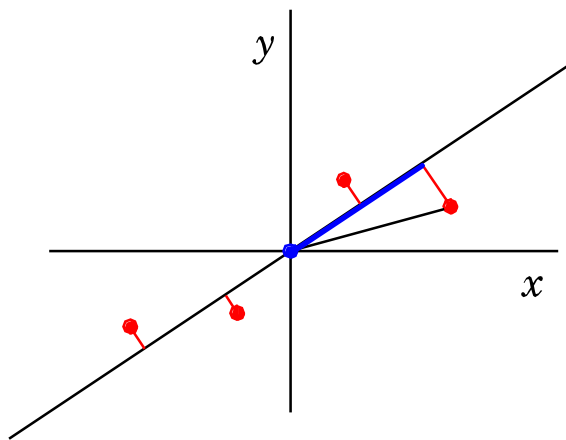
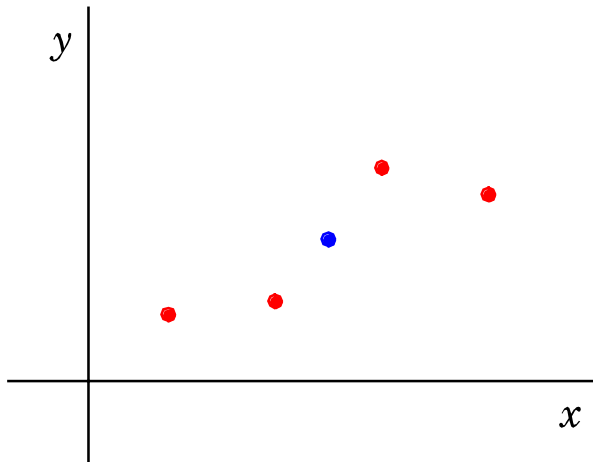
$x_1 = f_1(\xi), x_2 = f_2(\xi), \dots, x_n = f_n(\xi)$ と考える

$$y = a_1 x_1 + a_2 x_2 + \dots + a_n x_n = \sum_{k=1}^n a_k x_k$$

$f_k(\xi)$? 計算に必要なのは, データ $f_k(\xi_j) = x_{jk}$

$$\mathbf{A} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nn} \end{pmatrix} \quad \mathbf{A}^t \mathbf{A} \mathbf{a} = \mathbf{A}^t \mathbf{y}$$

主成分分析



生データ (x_i^*, y_i^*)

データの平均 (\bar{x}^*, \bar{y}^*)

$$\begin{aligned} x_i &= x_i^* - \bar{x}^* \\ y_i &= y_i^* - \bar{y}^* \end{aligned} \quad \mathbf{x}_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix}$$

データを最も良く表す直線

$$y = \frac{b}{a} x$$

データ (x_i, y_i) に対する距離最小

↓

\mathbf{x}_i の直線に対する射影長最大

$$\mathbf{u} = \begin{pmatrix} a \\ b \end{pmatrix} \quad (a^2 + b^2 = 1)$$

射影長は $\mathbf{x}_i^t \mathbf{u}$

$$\text{射影長の平方和} = \sum_{i=1}^N (\mathbf{x}_i^t \mathbf{u})^2$$

先ほどと同様に

$$\mathbf{A} = \begin{pmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_N & y_N \end{pmatrix} \quad \text{とおくと,} \quad \mathbf{A}\mathbf{u} = \begin{pmatrix} \mathbf{x}_1^t \mathbf{u} \\ \mathbf{x}_2^t \mathbf{u} \\ \vdots \\ \mathbf{x}_N^t \mathbf{u} \end{pmatrix}$$

$$\begin{aligned} (\mathbf{A}\mathbf{u})^t \mathbf{A}\mathbf{u} &= (\mathbf{x}_1^t \mathbf{u} \quad \mathbf{x}_2^t \mathbf{u} \quad \cdots \quad \mathbf{x}_N^t \mathbf{u}) \begin{pmatrix} \mathbf{x}_1^t \mathbf{u} \\ \mathbf{x}_2^t \mathbf{u} \\ \vdots \\ \mathbf{x}_N^t \mathbf{u} \end{pmatrix} = \sum_{i=1}^N (\mathbf{x}_i^t \mathbf{u})^2 \\ \parallel \\ \mathbf{u}^t \mathbf{A}^t \mathbf{A} \mathbf{u} &\leftarrow \text{これを最大にすれば良い} \end{aligned}$$

$\mathbf{u}^t (\mathbf{A}^t \mathbf{A} \mathbf{u})$ は, \mathbf{u} と $\mathbf{A}^t \mathbf{A} \mathbf{u}$ の内積

$\mathbf{A}^t \mathbf{A}$ により向きが変わらないベクトル (固有ベクトル) で,
変換倍率 (固有値) が最大のものを選べば良いのでは?

$\mathbf{A}^t \mathbf{A}$ は分散共分散行列であり, 固有値はすべて非負である

数値計算

| | | | | | | |
|-----|------|------|------|------|------|------|
| x | 0.0 | 1.0 | 2.0 | 3.0 | 4.0 | 5.0 |
| y | 10.2 | 12.0 | 15.7 | 17.0 | 20.5 | 22.4 |

$$\bar{x} = 2.5$$

$$\bar{y} = 16.3$$

$$\mathbf{A} = \begin{pmatrix} -2.5 & -6.1 \\ -1.5 & -4.3 \\ -0.5 & -0.6 \\ 0.5 & 0.7 \\ 1.5 & 4.2 \\ 2.5 & 6.1 \end{pmatrix}$$

$$\mathbf{A}^t \mathbf{A} = \begin{pmatrix} 17.5 & 43.9 \\ 43.9 & 111.4 \end{pmatrix}$$

$\mathbf{A}^t \mathbf{A}$ 固有値, 固有ベクトルは,

$$\det(\mathbf{A}^t \mathbf{A} - \lambda \mathbf{E}) = 0$$

$$(17.5 - \lambda)(111.4 - \lambda) - 43.9^2 = 0 \text{ を解いて}$$

$$\lambda_1 = 128.73 \quad \lambda_2 = 0.17$$

$$\mathbf{u}_1 = \begin{pmatrix} 0.37 \\ 0.93 \end{pmatrix} \quad \mathbf{u}_2 = \begin{pmatrix} -0.93 \\ 0.37 \end{pmatrix}$$

$$y = \frac{0.93}{0.37} (x - 2.5) + 16.3$$

$$y = 2.51x + 9.67$$

多次元のとき

$$\mathbf{A} = \begin{matrix} \leftarrow n \text{次元} \rightarrow \\ \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nn} \end{pmatrix} \\ \begin{matrix} \uparrow \\ N \text{データ} \\ \downarrow \end{matrix} \end{matrix}$$

$\mathbf{A}^t \mathbf{A}$ の最大の固有値に対応する固有ベクトル \mathbf{u}_1

$$\mathbf{z}_1 = \mathbf{A} \mathbf{u}_1 \quad \begin{matrix} \text{主成分 } \mathbf{z}_1 \\ \downarrow \\ \begin{pmatrix} z_{11} \\ z_{21} \\ \vdots \\ z_{N1} \end{pmatrix} \end{matrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nn} \end{pmatrix} \begin{matrix} \text{固有ベクトル } \mathbf{u}_1 \\ \downarrow \\ \begin{pmatrix} u_{11} \\ u_{21} \\ \vdots \\ u_{n1} \end{pmatrix} \end{matrix}$$

n 次元 \rightarrow 1次元

$\mathbf{A}^t \mathbf{A}$ の大きい方から1~ k 番目の固有値に対応する固有ベクトル

$$\mathbf{u}_1 \quad \mathbf{u}_2 \quad \cdots \quad \mathbf{u}_k$$

$$\mathbf{Z}_k = \mathbf{A} \mathbf{U}_k \quad \begin{array}{c} \text{第1主成分 } \mathbf{z}_1 \quad \text{第}k\text{主成分 } \mathbf{z}_k \\ \downarrow \qquad \qquad \downarrow \\ \begin{pmatrix} z_{11} & z_{12} & \cdots & z_{1k} \\ z_{21} & z_{22} & \cdots & z_{2k} \\ \vdots & \vdots & & \vdots \\ z_{N1} & z_{N2} & \cdots & z_{Nk} \end{pmatrix} \end{array} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nn} \end{pmatrix} \begin{array}{c} \mathbf{u}_1 \qquad \qquad \mathbf{u}_k \\ \downarrow \qquad \qquad \downarrow \\ \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1k} \\ u_{21} & u_{22} & \cdots & u_{2k} \\ \vdots & \vdots & & \vdots \\ u_{n1} & u_{n2} & \cdots & u_{nk} \end{pmatrix} \end{array}$$

n 次元 \rightarrow k 次元 (次元削減, 次元圧縮)

$$\text{寄与率} = \frac{\lambda_k}{\sum_{i=1}^n \lambda_i} \qquad \text{累積寄与率} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i}$$