

1.2. データの回帰分析

$y$  目的変数

$\mathbf{x} = (x_1, \dots, x_p)'$  : 説明変数

データ :  $(\mathbf{x}_i, y_i), i = 1, \dots, n$

(標本)  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})', (i = 1, \dots, n)$

行列記号で、 $X = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

統計量 :  $\left\{ \begin{array}{l} \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \\ \bar{y} \\ S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})': \text{標本分散行列} \\ s^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\ \mathbf{s} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(y_i - \bar{y}) : \text{標本共分散ベクトル} \end{array} \right.$

$y$	$x_1$	$\cdots$	$x_p$	$\mathbf{x}'$
$y_1$	$x_{11}$	$\cdots$	$x_{1p}$	$\mathbf{x}'_1$
$\vdots$	$\vdots$		$\vdots$	$\vdots$
$y_n$	$x_{n1}$	$\cdots$	$x_{np}$	$\mathbf{x}'_n$
$\bar{y}$	$\bar{x}_1$	$\cdots$	$\bar{x}_p$	$\bar{\mathbf{x}}'$

仮定 :  $\left\{ \begin{array}{l} S : \text{正值} \\ s > 0 \end{array} \right.$

[注意] :  $\text{rank } S \leq \min(n-1, p)$

より、 $S : \text{正值}$ の仮定は、 $n \geq p + 1$  を要求する。

理由 :  $X_R := X - \mathbf{1}_n \bar{\mathbf{x}}' = \begin{bmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_n \end{bmatrix} - \begin{bmatrix} \bar{\mathbf{x}}' \\ \vdots \\ \bar{\mathbf{x}}' \end{bmatrix} = \begin{bmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})' \\ \vdots \\ (\mathbf{x}_n - \bar{\mathbf{x}})' \end{bmatrix}$  とおくと、

$S = \frac{1}{n} X'_R X_R, p = \text{rank } S = \text{rank } X_R \leq n - 1$  より分かる。

**Prop.1.2.1.**  $Q(\alpha, \beta) := \sum_{i=1}^n (y_i - \alpha - \beta' \mathbf{x}_i)^2$

データにおいて  $y$  を  $\alpha + \beta' \mathbf{x}$  で近似したときの誤差

問題  $\left\{ \begin{array}{l} \min_{\alpha, \beta} Q(\alpha, \beta) \\ \text{条件 } \alpha \in \mathbf{R}^1, \beta : p \times 1 \end{array} \right.$

$\implies$

解  $\left\{ \begin{array}{l} \alpha = \hat{\alpha} := \bar{y} - \hat{\beta}' \bar{\mathbf{x}} \\ \beta = \hat{\beta} := S^{-1} \mathbf{s} \end{array} \right.$

最小値 =  $n(s^2 - \mathbf{s}' S^{-1} \mathbf{s})$

証明：r.v.  $\mathbf{X}$ ,  $Y$  を、次のように定義する。

$$Pr\{\mathbf{X} = \mathbf{x}_i, Y = y_i\} = \frac{1}{n} \quad (i = 1, \dots, n)$$

$$\text{このとき、} \boldsymbol{\mu} = E(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \bar{\mathbf{x}}, \nu = E(Y) = \bar{y},$$

$$\Sigma = \text{Var}(\mathbf{X}) = E\{\mathbf{X} - E(\mathbf{X})\}\{\mathbf{X} - E(\mathbf{X})\}' = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'$$

$$\boldsymbol{\sigma} = \text{Cov}(\mathbf{X}, Y) = \mathbf{s}, \sigma^2 = \text{Var}(Y) = s^2$$

より、

$$Q_0(\alpha, \boldsymbol{\beta}) = E(Y - \alpha - \boldsymbol{\beta}'\mathbf{X})^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \alpha - \boldsymbol{\beta}'\mathbf{x}_i)^2 = \frac{1}{n} Q(\alpha, \boldsymbol{\beta})$$

よって、Prop.1.1.1. より、

解  $\alpha = \alpha_0 = \nu - \boldsymbol{\beta}'\boldsymbol{\mu} = \bar{y} - \boldsymbol{\beta}_0'\bar{\mathbf{x}} = \hat{\alpha}$  とする。

$$\boldsymbol{\beta} = \boldsymbol{\beta}_0 = \Sigma^{-1}\boldsymbol{\sigma} = S^{-1}\mathbf{s} = \hat{\boldsymbol{\beta}}$$

$$Q \text{ の最小値は、} n\left(\sigma^2 - \boldsymbol{\sigma}'\Sigma^{-1}\boldsymbol{\sigma}\right) = n\left(s^2 - \mathbf{s}'S^{-1}\mathbf{s}\right)$$

証明終

**Def.1.2.1.**  $y = \hat{\alpha} + \hat{\boldsymbol{\beta}}'\mathbf{x}$  標本回帰関数 (標本回帰平面)

$$\hat{y}_i := \hat{\alpha} + \hat{\boldsymbol{\beta}}'\mathbf{x}_i \quad \text{回帰}$$

$$e_i := y_i - \hat{y}_i \quad \text{残差}$$

$$R := \sqrt{\mathbf{s}'S^{-1}\mathbf{s}}/s \quad \text{標本重相関係数}$$

$$R^2 \quad \text{寄与率 (標本決定係数)}$$

$$0 \leq R \leq 1$$

$$\left( \text{Min } Q = \sigma^2(1 - R^2) \geq 0 \right)$$

**Prop.1.2.2.** (i)  $\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}$  : 標本平均,  $\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = s^2 R^2$  : 標本分散

$$(ii) \frac{1}{n} \sum_{i=1}^n e_i = 0, \quad \frac{1}{n} \sum_{i=1}^n e_i^2 = s^2(1 - R^2)$$

$$(iii) \sum_{i=1}^n (\hat{y}_i - \bar{y})e_i = 0 : \text{標本共分散}$$

証明：(i)  $\hat{y}_i = \bar{y} + \hat{\boldsymbol{\beta}}'(\mathbf{x}_i - \bar{\mathbf{x}})$

$$\text{ゆえに、} \frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 &= \frac{1}{n} \sum_{i=1}^n \left\{ \hat{\beta}' (\mathbf{x}_i - \bar{\mathbf{x}}) \right\}^2 \\ &= \frac{1}{n} \hat{\beta}' \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' \hat{\beta} = \hat{\beta}' S \hat{\beta} = \mathbf{s}' S^{-1} \mathbf{s} = s^2 R^2 \end{aligned}$$

$$(ii) e_i = (y_i - \bar{y}) - \hat{\beta}' (\mathbf{x}_i - \bar{\mathbf{x}})$$

$$\text{ゆえに、} \frac{1}{n} \sum_{i=1}^n e_i = 0$$

$$\frac{1}{n} \sum_{i=1}^n e_i^2 = Q(\hat{\alpha}, \hat{\beta}) = s^2(1 - R^2)$$

$$(iii) \sum_{i=1}^n (\hat{y}_i - \bar{y}) e_i = \sum_{i=1}^n \hat{\beta}' (\mathbf{x}_i - \bar{\mathbf{x}}) \left[ (y_i - \bar{y}) - \hat{\beta}' (\mathbf{x}_i - \bar{\mathbf{x}}) \right] = \hat{\beta}' n(\mathbf{s} - S\hat{\beta}) = 0 \quad \text{証明終}$$

[注意] 要因の分解

変量	$y_i =$	$\hat{y}_i$	+	$e_i$
標本平均	$\bar{y} =$	$\bar{y}$	+	0
標本分散	$s^2 =$	$s^2 R^2$	+	$s^2(1 - R^2)$
		回帰		残差

**Prop.1.2.3.**  $\beta' \mathbf{x}_i$  ( $i = 1, \dots, n$ ) と  $y_i$  ( $i = 1, \dots, n$ ) の標本共分散を  $\text{Cov}(\beta' \mathbf{x}, y)$  で表すことにする。  $\text{Var}(\beta' \mathbf{x})$ ,  $\text{Var}(y)$  も同様。

$$\text{cor}(\beta' \mathbf{x}, y) = \text{Cov}(\beta' \mathbf{x}, y) / \left( \sqrt{\text{Var}(\beta' \mathbf{x})} \sqrt{\text{Var}(y)} \right).$$

$$\text{問題} \begin{cases} \max_{\beta} \text{cor}(\beta' \mathbf{x}, y) \\ \text{条件 } \beta : p \times 1 \end{cases}$$

$\implies$

$$\begin{aligned} \text{解} : \beta &\propto \hat{\beta} \\ \text{最大値} &= R \end{aligned}$$

理由 : Prop.1.2.1 の  $\mathbf{X}$ ,  $Y$  を使うと、

$$\text{Var}(\beta' \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n (\beta' \mathbf{x}_i - \beta' \bar{\mathbf{x}})^2 = \text{Var}(\beta' \mathbf{X})$$

$$\text{Var}(y) = \text{Var}(Y)$$

$$\text{Cov}(\beta' \mathbf{x}, y) = \text{Cov}(\beta' \mathbf{X}, Y)$$

$$\text{ゆえに、} \text{cor}(\beta' \mathbf{x}, y) = \text{cor}(\beta' \mathbf{X}, Y)$$

$$\text{よって、Prop.1.1.3 より、} \beta \propto \beta_0 = \Sigma^{-1} \sigma = S^{-1} \mathbf{s} = \hat{\beta}$$

理由終

これまでは、 $\alpha$  と  $\beta$  の扱いが非対称であった。

$$y_i \leftarrow \alpha + \beta' \mathbf{x}_i$$

( $\alpha$ : 定数項、 $\beta$ : 偏回帰係数)

これを、 $y_i \leftarrow \alpha \cdot 1 + \beta' \mathbf{x}_i$

と考え、 $\begin{bmatrix} \alpha \\ \beta \end{bmatrix}$ ,  $\begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix}$  を改めて、 $\beta$ ,  $\mathbf{x}_i$  と書けば、 $y_i \leftarrow \beta' \mathbf{x}_i$  となる。

逆に新しい記号において、 $\mathbf{x}_i$  の第1成分が恒等的に1ならば、前の場合になる。その場合を、「定数項のある場合」という。

データ:  $(\mathbf{x}_i, y_i)$ , ( $i = 1, \dots, n$ )  
行列記号で、 $X$ ,  $\mathbf{y}$

仮定:  $\text{rank } X = p$  このとき、 $X$  は **full rank** という。

[注意]: 定数項がある場合、これを除いた  $(p-1)$  個の変数についての標本分散行列を  $S$  とすると、

$$\text{full rank} \iff S: \text{正值}$$

**Prop.1.2.4.**  $Q(\beta) := \|\mathbf{y} - X\beta\|^2$  とする。このとき、以下の性質が成り立つ。

(i)  $Q(\beta) = Q(\hat{\beta}) + \|X(\beta - \hat{\beta})\|^2$

ここで、 $\hat{\beta} := (X'X)^{-1}X'\mathbf{y}$

(ii) 問題  $\begin{cases} \min_{\beta} Q(\beta) \\ \text{条件 } \beta: p \times 1 \end{cases}$

$\implies$

解:  $\beta = \hat{\beta}$

最小値 =  $\mathbf{y}'(I - \Pi_X)\mathbf{y}$

さらに、 $\Pi_X := X(X'X)^{-1}X': n \times n$ , 対称、中等、 $\text{rank } \Pi_X = p$

証明: (i)  $Q(\beta) = \|(\mathbf{y} - X\hat{\beta}) - X(\beta - \hat{\beta})\|^2$

$$= \|\mathbf{y} - X\hat{\beta}\|^2 + \|X(\beta - \hat{\beta})\|^2$$

なぜならば、積和 =  $2\{X(\beta - \hat{\beta})\}'(\mathbf{y} - X\hat{\beta}) = 2(\beta - \hat{\beta})' \underline{X'(\mathbf{y} - X\hat{\beta})} = 0$   
(下線部は、 $\mathbf{0}$ )

(ii)  $Q(\beta)$  を最小にする  $\beta$  は、 $\|X(\beta - \hat{\beta})\|^2 = 0$  の解

$$(\beta - \hat{\beta})' X' X (\beta - \hat{\beta}) = 0$$

ここで、 $X'X$  は正值であるから、 $\beta = \hat{\beta}$  に限る。

$$\mathbf{e} := \mathbf{y} - X\hat{\beta}$$

$$= \mathbf{y} - X(X'X)^{-1}X'\mathbf{y} = \{I_n - X(X'X)^{-1}X'\}\mathbf{y} = (I - \Pi_X)\mathbf{y}$$