

# 最小二乗法による単回帰の回帰係数の推定

クリスカ\*

2011年7月27日

## 1 回帰分析とは

身長と体重の間には、一方が増加すれば他方も増加するという関係があることが分かります。このような関係が成り立っているとき、この間の関係式を推定しようというのが回帰分析です。回帰分析は複数種類のデータを分析するための1つであり、あるデータを他のデータの関数によって説明するものです。

## 2 データの統計量

2種類のデータ  $x, y$  について  $n$  個の観測値  $(x_i, y_i)$  ( $i = 1, 2, \dots, n$ ) が与えられていたとします。

$$\text{データの平均: } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{データの分散: } (s_x)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - (\bar{x})^2 = \overline{x_i^2} - (\bar{x})^2$$

$$\text{データの共分散: } s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} = \overline{xy} - \bar{x} \bar{y}$$

$$\text{データの相関係数: } r_{xy} = \frac{s_{xy}}{s_x s_y}$$

Remark

データから平均を差し引き標準偏差で割って、位置、尺度の調整をすると、平均は0、標準偏差は1に揃ったこととなります。この操作をデータの標準化といいます。

$x_i$  を  $z_i = \frac{x_i - \bar{x}}{s_x}$ ,  $y_i$  を  $w_i = \frac{y_i - \bar{y}}{s_y}$  と標準化したデータ  $z_i, w_i$  を考えると、定義より、

$$r_{xy} = \frac{1}{n} \sum_{i=1}^n z_i w_i$$

となります。これは標準化したデータ  $z_i, w_i$  の共分散に等しいことが分かります。

---

\* k-nakagawa@h6.dion.ne.jp

このことから、 $-1 \leq r_{xy} \leq 1$  と言えます。証明は以下の通りです。

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n (z_i \pm w_i)^2 &= \frac{1}{n} \sum_{i=1}^n (z_i^2 \pm 2z_i w_i + w_i^2) \\ &= \frac{1}{n} \sum_{i=1}^n z_i^2 \pm \frac{2}{n} \sum_{i=1}^n z_i w_i + \frac{1}{n} \sum_{i=1}^n w_i^2 \\ &= 1 \pm 2r_{xy} + 1 \\ &= 2(1 \pm r_{xy})\end{aligned}$$

左辺は常に非負であるので  $1 \pm r_{xy} \geq 0$  より  $-1 \leq r_{xy} \leq 1$  と言えます。

また、 $x' = ax + b$ ,  $y' = cx + d$  ( $ac > 0$ ) のように線型変換した相関係数を  $r'_{xy}$  とすると線型不変性  $r'_{xy} = r_{xy}$  も言えます。

### 3 回帰分析

相関係数が 1 (もしくは  $-1$ ) に近く、正の相関 (負の相関) が認められるとき、一方の変数を他方の変数の一次関数として表すことが考えられます。

つまり、 $y$  を  $x$  の関数として考えるならば、 $y = \alpha + \beta x$  というモデルになりますが、このモデルを (線型) 回帰モデル、このモデルを用いた分析を 回帰分析 といいます。また、パラメータ  $\alpha, \beta$  を 回帰係数 といいます。

上記のように  $y$  を  $x$  の関数として考えるとき、 $x$  を 説明変数 (または 独立変数)、 $y$  を 被説明変数 (または 従属変数) といい、 $y$  を  $x$  によって説明することを「 $y$  を  $x$  に回帰する」といいます。

ここでは説明変数が 1 つの場合を考えていますが、観測データが  $(x_{1i}, x_{2i}, \dots, x_{mi}, y_i)$  ( $i = 1, 2, \dots, n$ ) として  $m$  種類あり、 $y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$  というモデル ( $y$  を  $x_1, x_2, \dots, x_m$  に回帰するモデル) を考えることもできます。このように説明変数が複数ある場合を 重回帰 と呼び、これに対して説明変数が 1 つである場合を 単回帰 と呼びます。

$\alpha, \beta_1, \beta_2, \dots$  は観測値を用いて推定することになりますが、推定された  $\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2, \dots$  を用いた関係式を 回帰式 (もしくは 回帰直線) といいます。

Remark

相関係数  $r_{xy}$  は  $x$  と  $y$  の直線的関係を測る尺度であり、この値から非線形の関係を読みとることはできません。

### 4 最小二乗法による単回帰の回帰係数の推定

観測値  $(x_i, y_i)$  ( $i = 1, 2, \dots, n$ ) を用いた単回帰の回帰係数の推定について考えることにします。

この場合、 $n$  個の観測値に対して、 $y_i = \alpha + \beta x_i$  ( $i = 1, 2, \dots, n$ ) という回帰モデルをあてはめて回帰係数  $\alpha, \beta$  の推定を行うこととなります。しかし、観測値が 1 つの直線上に完全に乗ることはまずあり得ません。ですので、回帰式の誤差である  $y_i - (\alpha + \beta x_i)$  を小さくする  $\alpha, \beta$  を推定値とすることとなります。

最小二乗法 では、誤差の平方和である下式  $S$  を最小にする  $\alpha, \beta$  を推定値とします。

$$S = \sum_{i=1}^n \{y_i - (\alpha + \beta x_i)\}^2$$

$x_i, y_i$  は変数ではなく、既に確定した値をとっていることに注意すると、 $S$  は  $\alpha, \beta$  の関数だということが分かります。  $S$  は非負値の二次式であるので、これを最小化する  $\alpha, \beta$  は、 $\alpha$  と  $\beta$  についての偏導関数が 0 となる必要があります。 よって、次の連立方程式の解が求める推定量  $\hat{\alpha}, \hat{\beta}$  となります。

$$\begin{aligned} & \begin{cases} \frac{\partial S}{\partial \alpha} = 0 \\ \frac{\partial S}{\partial \beta} = 0 \end{cases} \\ & \Leftrightarrow \begin{cases} -2 \sum_{i=1}^n \{y_i - (\alpha + \beta x_i)\} = 0 \\ -2 \sum_{i=1}^n x_i \{y_i - (\alpha + \beta x_i)\} = 0 \end{cases} \\ & \Leftrightarrow \begin{cases} \sum_{i=1}^n \{y_i - (\hat{\alpha} + \hat{\beta} x_i)\} = 0 \\ \sum_{i=1}^n x_i \{y_i - (\hat{\alpha} + \hat{\beta} x_i)\} = 0 \end{cases} \\ & \Leftrightarrow \begin{cases} \hat{\alpha} n + \hat{\beta} \sum_{i=1}^n x_i & = \sum_{i=1}^n y_i \\ \hat{\alpha} \sum_{i=1}^n x_i + \hat{\beta} \sum_{i=1}^n x_i^2 & = \sum_{i=1}^n x_i y_i \end{cases} \end{aligned}$$

ここで、式変形によって得られた最後の  $\hat{\alpha}, \hat{\beta}$  に関する連立方程式を 正規方程式 といいます。 正規方程式の上式の両辺を  $n$  で割って整理すると

$$\begin{aligned} \hat{\alpha} + \hat{\beta} \frac{1}{n} \sum_{i=1}^n x_i &= \frac{1}{n} \sum_{i=1}^n y_i \Leftrightarrow \hat{\alpha} + \hat{\beta} \bar{x} = \bar{y} \\ &\Leftrightarrow \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \end{aligned}$$

となります。これを下式に代入して、両辺を  $n$  で割って整理すると

$$\begin{aligned} (\bar{y} - \hat{\beta} \bar{x}) \frac{1}{n} \sum_{i=1}^n x_i + \hat{\beta} \frac{1}{n} \sum_{i=1}^n x_i^2 &= \frac{1}{n} \sum_{i=1}^n x_i y_i \Leftrightarrow (\bar{y} - \hat{\beta} \bar{x}) \bar{x} + \hat{\beta} \bar{x}^2 = \overline{xy} \\ &\Leftrightarrow \hat{\beta} (\bar{x}^2 - (\bar{x})^2) = \overline{xy} - \bar{x} \bar{y} \\ &\Leftrightarrow \hat{\beta} (s_x)^2 = s_{xy} \\ &\Leftrightarrow \hat{\beta} = \frac{s_{xy}}{(s_x)^2} \\ &\Leftrightarrow \hat{\beta} = \frac{r_{xy} s_y}{s_x} \end{aligned}$$

$\bar{y} = \hat{\alpha} + \hat{\beta} \bar{x}$  より、回帰直線は観測データの平均  $(\bar{x}, \bar{y})$  を通ることが分かります。

$y = \hat{\alpha} + \hat{\beta} x$  を  $y$  を被説明変数 ( $x$  を説明変数) とする回帰直線といいます。

$y$  を  $x$  で回帰するともいいます。

また,

$$\begin{aligned}y = \hat{\alpha} + \hat{\beta}x &\iff y = (\bar{y} - \hat{\beta}\bar{x}) + \hat{\beta}x \\&\iff y - \bar{y} = \hat{\beta}(x - \bar{x}) \\&\iff y - \bar{y} = \frac{r_{xy}s_y}{s_x}(x - \bar{x}) \\&\iff \frac{y - \bar{y}}{s_y} = r_{xy} \times \frac{x - \bar{x}}{s_x}\end{aligned}$$

と変形できることより ( $y$  の標準化) = (相関係数)  $\times$  ( $x$  の標準化) という関係があることも分かります.

次に, 回帰式に  $x_i$  を代入して得られる値  $\hat{y}_i$  を  $y_i$  の内挿値といい,  $y_i$  と  $\hat{y}_i$  の差  $e_i = y_i - \hat{y}_i$  を残差といいます. これを用いると

$$\begin{cases} \sum_{i=1}^n \{y_i - (\hat{\alpha} + \hat{\beta}x_i)\} = 0 \\ \sum_{i=1}^n x_i \{y_i - (\hat{\alpha} + \hat{\beta}x_i)\} = 0 \end{cases} \iff \begin{cases} \sum_{i=1}^n e_i = 0 \\ \sum_{i=1}^n e_i x_i = 0 \end{cases}$$

が分かります. これは,  $e_i$  の平均が 0 であり, また,  $\{e_i\}$  と  $\{x_i\}$  はベクトルとして直交していることを示しています.

## 5 多重線型回帰と正規方程式

単回帰のときと同様な方法で求めます.

観測値  $(x_{1i}, x_{2i}, \dots, x_{mi}, y_i)$  ( $i = 1, 2, \dots, n$ ) を用いた多重線型回帰の回帰係数の推定について考えることにします.

この場合,  $n$  個の観測値に対して,  $y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_m x_{mi}$  ( $i = 1, 2, \dots, n$ ) という回帰モデルをあてはめて回帰係数  $\alpha, \beta_1, \beta_2, \dots, \beta_m$  の推定を行うこととなります. 回帰式の誤差である  $y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_m x_{mi})$  を小さくする  $\alpha, \beta_1, \beta_2, \dots, \beta_m$  を推定値とすることとなります.

最小二乗法では, 誤差の平方和である下式  $S$  を最小にする  $\alpha, \beta_1, \beta_2, \dots, \beta_m$  を推定値とします.

$$S = \sum_{i=1}^n \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_m x_{mi})\}^2$$

$S$  は非負値の二次式であるので, これを最小化する  $\alpha, \beta_1, \beta_2, \dots, \beta_m$  は, 各々の変数についての偏導関数が 0 となる必要があります. よって, 次の連立方程式の解が求める推定量  $\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_m$  となります.

$$\begin{aligned}
\frac{\partial S}{\partial \alpha} = \frac{\partial S}{\partial \beta_1} = \frac{\partial S}{\partial \beta_2} = \cdots = \frac{\partial S}{\partial \beta_m} = 0 \\
\iff \left\{ \begin{array}{l} -2 \sum_{i=1}^n \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \\ -2 \sum_{i=1}^n x_{1i} \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \\ -2 \sum_{i=1}^n x_{2i} \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \\ \vdots \\ -2 \sum_{i=1}^n x_{mi} \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \end{array} \right. \\
\iff \left\{ \begin{array}{l} \sum_{i=1}^n \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \\ \sum_{i=1}^n x_{1i} \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \\ \sum_{i=1}^n x_{2i} \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \\ \vdots \\ \sum_{i=1}^n x_{mi} \{y_i - (\alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_m x_{mi})\} = 0 \end{array} \right. \\
\iff \left\{ \begin{array}{l} \hat{\alpha} n + \hat{\beta}_1 \sum_{i=1}^n x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i} + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{mi} = \sum_{i=1}^n y_i \\ \hat{\alpha} \sum_{i=1}^n x_{1i} + \hat{\beta}_1 \sum_{i=1}^n x_{1i}^2 + \hat{\beta}_2 \sum_{i=1}^n x_{1i} x_{2i} + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{1i} x_{mi} = \sum_{i=1}^n x_{1i} y_i \\ \hat{\alpha} \sum_{i=1}^n x_{2i} + \hat{\beta}_1 \sum_{i=1}^n x_{2i} x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i}^2 + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{2i} x_{mi} = \sum_{i=1}^n x_{2i} y_i \\ \vdots \\ \hat{\alpha} \sum_{i=1}^n x_{mi} + \hat{\beta}_1 \sum_{i=1}^n x_{mi} x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{mi} x_{2i} + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{mi}^2 = \sum_{i=1}^n x_{mi} y_i \end{array} \right.
\end{aligned}$$

ここで、式変形によって得られた最後の  $\hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_m$  に関する連立方程式を正規方程式といいます。正規方程式は式の数が少ないときは容易に解けますが、式の数が多いときは電子計算機が必要になります。

Remark

多重線形回帰の正規方程式は行列を用いることにより以下のように書き直すことができます。

$$\begin{cases} \hat{\alpha}n + \hat{\beta}_1 \sum_{i=1}^n x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i} + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{mi} & = \sum_{i=1}^n y_i \\ \hat{\alpha} \sum_{i=1}^n x_{1i} + \hat{\beta}_1 \sum_{i=1}^n x_{1i}^2 + \hat{\beta}_2 \sum_{i=1}^n x_{1i}x_{2i} + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{1i}x_{mi} & = \sum_{i=1}^n x_{1i}y_i \\ \hat{\alpha} \sum_{i=1}^n x_{2i} + \hat{\beta}_1 \sum_{i=1}^n x_{2i}x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{2i}^2 + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{2i}x_{mi} & = \sum_{i=1}^n x_{2i}y_i \\ \vdots & \\ \hat{\alpha} \sum_{i=1}^n x_{mi} + \hat{\beta}_1 \sum_{i=1}^n x_{mi}x_{1i} + \hat{\beta}_2 \sum_{i=1}^n x_{mi}x_{2i} + \cdots + \hat{\beta}_m \sum_{i=1}^n x_{mi}^2 & = \sum_{i=1}^n x_{mi}y_i \end{cases}$$

⇕

$$\begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix} \begin{pmatrix} 1 & x_{11} & x_{21} & \cdots & x_{m1} \\ 1 & x_{12} & x_{22} & \cdots & x_{m2} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{mn} \end{pmatrix} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_m \end{pmatrix} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

ここで、 $X = \begin{pmatrix} 1 & x_{11} & x_{21} & \cdots & x_{m1} \\ 1 & x_{12} & x_{22} & \cdots & x_{m2} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{mn} \end{pmatrix}$  とおくと上式は

$$X^T X \begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_m \end{pmatrix} = X^T \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

と変形できるので、最小二乗法による推定値は

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_m \end{pmatrix} = (X^T X)^{-1} X^T \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

となります。

単回帰のときのように、各パラメータについて明示的な式を示すことはできませんが、上式右辺の行列計算をした結果が最小二乗法による推定値となります。