

## 第4章 回帰分析

### 4.1 準備

#### 4.1.1 重要な公式

$$1. \sum_{i=1}^n X_i = n\bar{X}$$

$$2. \sum_{i=1}^n (X_i - \bar{X}) = 0$$

$$3. \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n X_i^2 - n\bar{X}^2$$

$$4. \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y} = \sum_{i=1}^n (X_i - \bar{X})Y_i = \sum_{i=1}^n (Y_i - \bar{Y})X_i$$

$$5. 2 \times 2 \text{ 行列の逆行列の公式: } \begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

#### 4.1.2 データについて

1. タイム・シリーズ(時系列)・データ：添え字  $i$  が時間を表す(第  $i$  期)。  $t$  を添え字に使う場合も多い。
2. クロス・セクション(横断面)・データ：添え字  $i$  が個人や企業を表す(第  $i$  番目の家計, 第  $i$  番目の企業)。

## 4.2 最小二乗法について：単回帰モデル

最小二乗法とは、線型モデルの係数の値をデータから求める時に用いられる手法である。

### 4.2.1 最小二乗法と回帰直線

$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$  のように  $n$  組のデータがあり、 $X_i$  と  $Y_i$  との間に以下の線型関係を想定する。

$$Y_i = \alpha + \beta X_i,$$

$X_i$  は説明変数、 $Y_i$  は被説明変数、 $\alpha, \beta$  はパラメータとそれぞれ呼ばれる。

上の式は回帰モデル（または、回帰式）と呼ばれる。切片  $\alpha$  と傾き  $\beta$  をデータ  $\{(X_i, Y_i), i = 1, 2, \dots, n\}$  から推定することを考える。

ある基準の下で、 $\alpha$  と  $\beta$  の推定値が求められたとしよう。それぞれ、 $\hat{\alpha}$  と  $\hat{\beta}$  とする。データ  $\{(X_i, Y_i), i = 1, 2, \dots, n\}$  と直線との関係は、

$$Y_i = \hat{\alpha} + \hat{\beta} X_i + \hat{u}_i,$$

となる。すなわち、実際のデータ  $Y_i$  と直線上の値  $\hat{\alpha} + \hat{\beta} X_i$  との間には、誤差  $\hat{u}_i$ （残差と呼ばれる）が生じる。

### 4.2.2 切片 $\alpha$ と傾き $\beta$ の求め方

$\alpha, \beta$  のある推定値を  $\hat{\alpha}, \hat{\beta}$  としよう。次のような関数  $S(\hat{\alpha}, \hat{\beta})$  を定義する。

$$S(\hat{\alpha}, \hat{\beta}) = \sum_{i=1}^n \hat{u}_i^2 = \sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta} X_i)^2$$

これは残差平方和と呼ばれる。

このとき、

$$\min_{\hat{\alpha}, \hat{\beta}} S(\hat{\alpha}, \hat{\beta})$$

となるような  $\hat{\alpha}, \hat{\beta}$  を求める（最小自乗法）。  
最小化のためには、

$$\frac{\partial S(\hat{\alpha}, \hat{\beta})}{\partial \hat{\alpha}} = 0, \quad \frac{\partial S(\hat{\alpha}, \hat{\beta})}{\partial \hat{\beta}} = 0$$

を満たす  $\hat{\alpha}, \hat{\beta}$  を求める。

すなわち、 $\hat{\alpha}, \hat{\beta}$  は、

$$\sum_{i=1}^n (Y_i - \hat{\alpha} - \hat{\beta}X_i) = 0, \quad (4.1)$$

$$\sum_{i=1}^n X_i(Y_i - \hat{\alpha} - \hat{\beta}X_i) = 0, \quad (4.2)$$

を満たす。

さらに、

$$\sum_{i=1}^n Y_i = n\hat{\alpha} + \hat{\beta} \sum_{i=1}^n X_i \quad (4.3)$$

$$\sum_{i=1}^n X_i Y_i = \hat{\alpha} \sum_{i=1}^n X_i + \hat{\beta} \sum_{i=1}^n X_i^2 \quad (4.4)$$

(4.3) 式の辺々を  $n$  で割って、

$$\frac{1}{n} \sum_{i=1}^n Y_i = \hat{\alpha} + \hat{\beta} \frac{1}{n} \sum_{i=1}^n X_i$$

すなわち、

$$\bar{Y} = \hat{\alpha} + \hat{\beta} \bar{X} \quad (4.5)$$

を得る。ただし、

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i,$$

とする。

さらに、 $\sum_{i=1}^n X_i = n\bar{X}$  と (4.5) 式を利用して、 $\hat{\alpha}$  を消去すると、

$$\sum_{i=1}^n X_i Y_i = (\bar{Y} - \hat{\beta}\bar{X})n\bar{X} + \hat{\beta} \sum_{i=1}^n X_i^2$$

$\hat{\beta}$  で整理して、

$$\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{S_{XY}}{S_X^2} \quad (4.6)$$

が得られ、 $\hat{\alpha}$  は (4.5) 式から、

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X} \quad (4.7)$$

となる。ただし、

$$S_{XY} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}), \quad S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

とする。

または、行列を用いて解くこともできる。行列表示によって、

$$\begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n X_i Y_i \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i & \sum_{i=1}^n X_i^2 \end{pmatrix} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix},$$

$\hat{\alpha}$ ,  $\hat{\beta}$  について、まとめて、

$$\begin{aligned} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} &= \begin{pmatrix} n & \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i & \sum_{i=1}^n X_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n X_i Y_i \end{pmatrix} \\ &= \frac{1}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \begin{pmatrix} \sum_{i=1}^n X_i^2 & -\sum_{i=1}^n X_i \\ -\sum_{i=1}^n X_i & n \end{pmatrix} \begin{pmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n X_i Y_i \end{pmatrix} \end{aligned}$$

さらに、 $\hat{\beta}$  について解くと、

$$\begin{aligned} \hat{\beta} &= \frac{n \sum_{i=1}^n X_i Y_i - (\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2} \\ &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \end{aligned}$$

$\hat{\alpha}$  については,

$$\begin{aligned}\hat{\alpha} &= \frac{(\sum_{i=1}^n X_i^2)(\sum_{i=1}^n Y_i) - (\sum_{i=1}^n X_i)(\sum_{i=1}^n X_i Y_i)}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} = \frac{\bar{Y} \sum_{i=1}^n X_i^2 - \bar{X} \sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2 - n\bar{X}^2} \\ &= \frac{\bar{Y}(\sum_{i=1}^n X_i^2 - n\bar{X}^2) - \bar{X}(\sum_{i=1}^n X_i Y_i - n\bar{Y}\bar{X})}{\sum_{i=1}^n X_i^2 - n\bar{X}^2} = \bar{Y} - \frac{\sum_{i=1}^n X_i Y_i - n\bar{Y}\bar{X}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2} \\ &= \bar{Y} - \hat{\beta}\bar{X}\end{aligned}$$

となる。

回帰直線は,

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i,$$

として与えられる。 $\hat{Y}_i$  は、 $X_i$  を与えたときの  $Y_i$  の予測値と解釈される。

数値例： 以下の数値例を使って、回帰式  $Y_i = \alpha + \beta X_i$  の  $\alpha$ ,  $\beta$  の推定値  $\hat{\alpha}$ ,  $\hat{\beta}$  を求める。

$i$	$X_i$	$Y_i$
1	5	4
2	1	1
3	3	1
4	2	3
5	4	4

$\hat{\alpha}$ ,  $\hat{\beta}$  を求めるための公式は,

$$\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2}, \quad \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X},$$

なので、必要なものは  $\bar{X}$ ,  $\bar{Y}$ ,  $\sum_{i=1}^n X_i^2$ ,  $\sum_{i=1}^n X_i Y_i$  である。

$i$	$X_i$	$Y_i$	$X_i^2$	$X_i Y_i$
1	5	4	25	20
2	1	1	1	1
3	3	1	9	3
4	2	3	4	6
5	4	4	16	16
合計	$\sum X_i$ 15	$\sum Y_i$ 13	$\sum X_i^2$ 55	$\sum X_i Y_i$ 46
平均	$\bar{X}$ 3	$\bar{Y}$ 2.6		

表中では,  $\sum_{i=1}^n$  を  $\Sigma$  と省略して表記している。

よって,

$$\hat{\beta} = \frac{46 - 5 \times 3 \times 2.6}{55 - 5 \times 3^2} = \frac{7}{10} = 0.7, \quad \hat{\alpha} = 2.6 - 0.7 \times 3 = 0.5,$$

となる。

注意事項 :

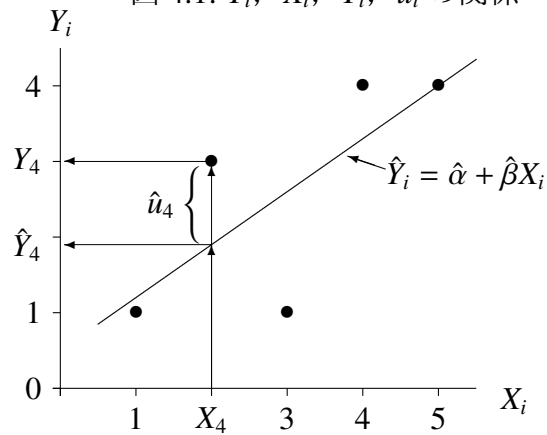
1.  $\alpha, \beta$  は真の値で未知である。
2.  $\hat{\alpha}, \hat{\beta}$  は  $\alpha, \beta$  の推定値でデータから計算される。

回帰直線は,  $\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i$  であり, 上の数値例では,

$$\hat{Y}_i = 0.5 + 0.7X_i,$$

となる。 $\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_5$  として, 次の表のように計算される。 $Y_i, X_i, \hat{Y}_i, \hat{u}_i$  の関係が図 4.1 に描かれている。

図 4.1:  $Y_i, X_i, \hat{Y}_i, \hat{u}_i$  の関係



$i$	$X_i$	$Y_i$	$X_i^2$	$X_i Y_i$	$\hat{Y}_i$
1	5	4	25	20	4.0
2	1	1	1	1	1.2
3	3	1	9	3	2.6
4	2	3	4	6	1.9
5	4	4	16	16	3.3
合計	$\sum X_i$	$\sum Y_i$	$\sum X_i^2$	$\sum X_i Y_i$	$\sum \hat{Y}_i$
	15	13	55	46	13
平均	$\bar{X}$	$\bar{Y}$			
	3	2.6			

$\hat{Y}_i$  を実績値  $Y_i$  の予測値または理論値と呼ぶ。

$$\hat{u}_i = Y_i - \hat{Y}_i,$$

$\hat{u}_i$  を残差と呼ぶ。 $Y_i, \hat{Y}_i, \hat{u}_i$  の関係,  $\hat{Y}_i, X_i, \hat{\alpha}, \hat{\beta}$  の関係は,

$$Y_i = \hat{Y}_i + \hat{u}_i = \hat{\alpha} + \hat{\beta}X_i + \hat{u}_i,$$

の式でまとめられる。

### 4.2.3 残差 $\hat{u}_i$ の性質について

$\hat{u}_i = Y_i - \hat{\alpha} - \hat{\beta}X_i$  に注意すると, (4.1) 式, (4.2) 式から,

$$\sum_{i=1}^n \hat{u}_i = 0, \quad \sum_{i=1}^n X_i \hat{u}_i = 0,$$

を得る。また,  $\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i$  から,

$$\sum_{i=1}^n \hat{Y}_i \hat{u}_i = 0,$$

が得られる。なぜなら,

$$\sum_{i=1}^n \hat{Y}_i \hat{u}_i = \sum_{i=1}^n (\hat{\alpha} + \hat{\beta}X_i) \hat{u}_i = \hat{\alpha} \sum_{i=1}^n \hat{u}_i + \hat{\beta} \sum_{i=1}^n X_i \hat{u}_i = 0$$

となるからである。

数値例で確認してみよう。

$i$	$X_i$	$Y_i$	$X_i^2$	$X_i Y_i$	$\hat{Y}_i$	$\hat{u}_i$	$X_i \hat{u}_i$	$\hat{Y}_i \hat{u}_i$
1	5	4	25	20	4.0	0.0	0.0	0.00
2	1	1	1	1	1.2	-0.2	-0.2	-0.24
3	3	1	9	3	2.6	-1.6	-4.8	-4.16
4	2	3	4	6	1.9	1.1	2.2	2.09
5	4	4	16	16	3.3	0.7	2.8	2.31
合計	$\sum X_i$ 15	$\sum Y_i$ 13	$\sum X_i^2$ 55	$\sum X_i Y_i$ 46	$\sum \hat{Y}_i$ 13	$\sum \hat{u}_i$ 0.0	$\sum X_i \hat{u}_i$ 0.0	$\sum \hat{Y}_i \hat{u}_i$ 0.0
平均	$\bar{X}$ 3	$\bar{Y}$ 2.6						

### 4.2.4 決定係数 $R^2$ について

$Y_i, \hat{Y}_i, \hat{u}_i$  の関係は,

$$Y_i = \hat{Y}_i + \hat{u}_i,$$



であった。 $\bar{Y}$ を両辺から引くと、

$$(Y_i - \bar{Y}) = (\hat{Y}_i - \bar{Y}) + \hat{u}_i,$$

が得られる。さらに、両辺を二乗して、総和すると、

$$\begin{aligned} \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \sum_{i=1}^n ((\hat{Y}_i - \bar{Y}) + \hat{u}_i)^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + 2 \sum_{i=1}^n (\hat{Y}_i - \bar{Y})\hat{u}_i + \sum_{i=1}^n \hat{u}_i^2 \\ &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{u}_i^2 \end{aligned}$$

となる。二つ目の等式の右辺第二項では、 $\sum_{i=1}^n \hat{Y}_i \hat{u}_i = \bar{Y} \sum_{i=1}^n \hat{u}_i = 0$ が使われている。まとめると、

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{u}_i^2$$

を得る。さらに、両辺を左辺で割ると、

$$1 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} + \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2},$$

が得られる。それぞれの項は、

1.  $\sum_{i=1}^n (Y_i - \bar{Y})^2 \rightarrow Y_i$  の全変動
2.  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 \rightarrow \hat{Y}_i$  (回帰直線) で説明される部分
3.  $\sum_{i=1}^n \hat{u}_i^2 \rightarrow \hat{Y}_i$  (回帰直線) で説明されない部分

となる。

回帰式の当てはまりの良さを示す指標として、決定係数  $R^2$  が、

$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad (4.8)$$

のように定義される。 $R^2$  は  $Y_i$  のうち  $\hat{Y}_i$  (または,  $X_i$ ) で説明できる比率を意味する。または,

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad (4.9)$$

として書き換えることもできる。

**$R^2$  の取り得る範囲:** さらに,  $R^2$  の取り得る範囲を求める。(4.8) 式の右辺の分子と分母は共に正なので,  $R^2 \geq 0$  となる。(4.9) 式の右辺では 1 から第二項の正の値 (分子分母共に正) を差し引いているので,  $R^2 \leq 1$  となることが分かる。すなわち,  $R^2$  の取り得る範囲は,

$$0 \leq R^2 \leq 1,$$

となる。

$R^2 = 1$  となる場合はすべての  $i$  について  $\hat{u}_i = 0$  となり, 観測されたデータ  $(X_i, Y_i)$  は一直線上に並んでいる状態となる。

$R^2 = 0$  となる場合は二通りが考えられる。一つは,  $Y_i$  が  $X_i$  に影響されないときで,  $\hat{\beta} = 0$  の状態, すなわち, データが横軸に平行に一直線上に並んでいる状態となる。もう一つは, データが円状に散布していて, どこにも直線が引けない状態である (ちなみに, データが楕円上に散布している場合は, 直線が引ける状態である)。

実際のデータを用いた場合は  $R^2 = 0$  や  $R^2 = 1$  という状況はあり得ない。 $R^2$  が 1 に近づけば回帰式の当てはまりは良い,  $R^2$  が 0 に近づけば回帰式の当てはまりは悪いと言える。しかし, 「どの値よりも大きくなるべき」といった基準はない。慣習的には, メドとして 0.9 以上が当てはまりが良いと判断する。

データと  $R^2$  との関係は, 後述の 4.2.5 節で, 数値例を挙げながら解説する。

**$R^2$  の別の解釈:**  $R^2$  のもう一つの解釈をするために,  $R^2$  の右辺の分子を,

$$\begin{aligned} \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \bar{Y} - \hat{u}_i) = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \bar{Y}) - \sum_{i=1}^n (\hat{Y}_i - \bar{Y})\hat{u}_i \\ &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \bar{Y}), \end{aligned}$$

と書き換える。最初の等式では、括弧二乗の一つに  $\hat{Y}_i = Y_i - \hat{u}_i$  が用いられている。 $R^2$  は、

$$\begin{aligned} R^2 &= \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{(\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2)^2}{(\sum_{i=1}^n (Y_i - \bar{Y})^2)(\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2)} \\ &= \left( \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2} \sqrt{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}} \right)^2, \end{aligned}$$

と書き換えられる。この式では、 $R^2$  が  $Y_i$  と  $\hat{Y}_i$  の相関係数の二乗と解釈されることを意味する。なお、二つ目の等号の右式では、分子と分母に  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$  を掛けていることに注意せよ。

特に、単回帰の場合、 $\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i$  と  $\bar{Y} = \hat{\alpha} + \hat{\beta}\bar{X}$  を用いて、

$$\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \hat{\beta}^2 \sum_{i=1}^n (X_i - \bar{X})^2 = \hat{\beta} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}),$$

を利用すると、

$$\begin{aligned} R^2 &= \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{\hat{\beta}^2 \sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \left( \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2} \sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}} \right)^2 \\ &= \frac{S_{XY}^2}{S_X^2 S_Y^2}, \end{aligned}$$

としても書き換えられる。すなわち、単回帰の場合、決定係数は説明変数  $X_i$  と被説明変数  $Y_i$  との相関係数の二乗となる。

数値例： 決定係数の計算には以下の公式を用いる。

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n Y_i^2 - n\bar{Y}^2}$$

計算に必要なものは、 $\sum_{i=1}^n \hat{u}_i^2$ 、 $\bar{Y}$ 、 $\sum_{i=1}^n Y_i^2$  である。

$i$	$X_i$	$Y_i$	$X_i^2$	$X_i Y_i$	$\hat{Y}_i$	$\hat{u}_i$	$X_i \hat{u}_i$	$\hat{Y}_i \hat{u}_i$	$\hat{u}_i^2$	$Y_i^2$
1	5	4	25	20	4.0	0.0	0.0	0.00	0.00	16
2	1	1	1	1	1.2	-0.2	-0.2	-0.24	0.04	1
3	3	1	9	3	2.6	-1.6	-4.8	-4.16	2.56	1
4	2	3	4	6	1.9	1.1	2.2	2.09	1.21	9
5	4	4	16	16	3.3	0.7	2.8	2.31	0.49	16
合計	$\sum X_i$ 15	$\sum Y_i$ 13	$\sum X_i^2$ 55	$\sum X_i Y_i$ 46	$\sum \hat{Y}_i$ 13	$\sum \hat{u}_i$ 0.0	$\sum X_i \hat{u}_i$ 0.0	$\sum \hat{Y}_i \hat{u}_i$ 0.0	$\sum \hat{u}_i^2$ 4.3	$\sum Y_i^2$ 43
平均	$\bar{X}$ 3	$\bar{Y}$ 2.6								

$\bar{Y} = 2.6$ ,  $\sum_{i=1}^n \hat{u}_i^2 = 4.3$ ,  $\sum_{i=1}^n Y_i^2 = 43$  なので,

$$R^2 = 1 - \frac{4.3}{43 - 5 \times 2.6^2} = \frac{4.9}{9.2} = 0.5326$$

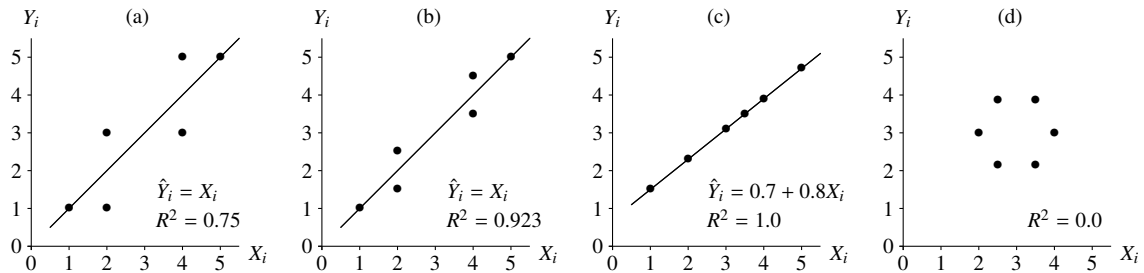
#### 4.2.5 決定係数の比較

次の数値例を用いて、決定係数の比較を行おう。 $X$ と $Y$ のプロットしたものが図4.2(a)~(d)である。

$i$	(a)		(b)		(c)		(d)	
	$X_i$	$Y_i$	$X_i$	$Y_i$	$X_i$	$Y_i$	$X_i$	$Y_i$
1	1	1	1	1	1	1.5	1	3
2	2	1	2	1.5	2	2.3	2.5	2.134
3	2	3	2	2.5	3	3.1	2.5	3.866
4	4	3	4	3.5	3.5	3.5	3.5	2.134
5	4	5	4	4.5	4	3.9	3.5	3.866
6	5	5	5	5	5	4.7	4	3

(a)と(b)のどちらの場合も、切片・傾きの値は $\hat{\alpha} = 0$ ,  $\hat{\beta} = 1$ として計算されるが、決定係数について、(a)は0.75, (b)は0.923となる(読者はチェック

図 4.2: 決定係数の比較



すること)。データのプロットと回帰直線は図4.2の(a)と(b)に描かれている。 $X_i$ はどちらも同じ数値とした。横軸 $X$ が2, 4のケースについて、(b)が(a)より直線に近くなるように、 $Y$ の値を変えてみた。(b)のデータの方が(a)より直線に近いために、決定係数が0.923と1に近い値となっているのが分かる。

(c)はデータが一直線上に並んでいる場合で、決定係数が1となる。決定係数がゼロとなるのは(d)の場合で、 $X$ と $Y$ との関係を表す直線が描けない場合である。(d)の数値例では、 $X$ と $Y$ との関係が円としているが、満遍なく散布している状態と考えてもらえれば良い。

#### 4.2.6 まとめ

$\hat{\alpha}$ ,  $\hat{\beta}$ を求めるための公式は

$$\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

なので、必要なものは $\bar{X}$ ,  $\bar{Y}$ ,  $\sum_{i=1}^n X_i^2$ ,  $\sum_{i=1}^n X_i Y_i$ である。

決定係数の計算には以下の公式を用いる。

$$R^2 = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n Y_i^2 - n\bar{Y}^2}$$

ただし、 $\hat{u}_i = Y_i - \hat{\alpha} - \hat{\beta}X_i$ である。計算に必要なものは、 $\sum_{i=1}^n \hat{u}_i^2$ ,  $\bar{Y}$ ,  $\sum_{i=1}^n Y_i^2$ である。

### 4.3 最小二乗法について：重回帰モデル

$k$ 変数の多重回帰モデルを考える。

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki}$$

$X_{ji}$ は  $j$ 番目の説明変数の第  $i$ 番目の観測値を表す。 $\beta_1, \beta_2, \dots, \beta_k$ は推定されるべきパラメータである。すべての  $i$ について、 $X_{1i} = 1$ とすれば、 $\beta_1$ は定数項として表される。 $n$ 組のデータ  $(Y_i, X_{1i}, X_{2i}, \dots, X_{ki})$ ,  $i = 1, 2, \dots, n$ を用いて、 $\beta_1, \beta_2, \dots, \beta_k$ を求める。

ある基準の下で、 $\beta_1, \beta_2, \dots, \beta_k$ の解を  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ としよう。データ  $\{(X_i, Y_i), i = 1, 2, \dots, n\}$ と直線との関係は、

$$Y_i = \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki} + \hat{u}_i = \hat{Y}_i + \hat{u}_i,$$

となる。すなわち、すべての  $i$ について、実際のデータ  $Y_i$ と直線上の値  $\hat{Y}_i = \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}$ が一致することはあり得ないので、残差  $\hat{u}_i$ の二乗和を考える。

次のような関数  $S(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)$ を定義する。

$$S(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k) = \sum_{i=1}^n u_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \cdots - \hat{\beta}_k X_{ki})^2$$

このとき、

$$\min_{\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k} S(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)$$

となるような  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ を求める。⇒ 最小自乗法  
最小化のためには、

$$\frac{\partial S(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)}{\partial \hat{\beta}_1} = 0, \quad \frac{\partial S(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)}{\partial \hat{\beta}_2} = 0, \quad \dots, \quad \frac{\partial S(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)}{\partial \hat{\beta}_k} = 0$$

を満たす  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$  となる。

すなわち,  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$  は,

$$\begin{aligned} \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki}) X_{1i} &= 0, \\ \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki}) X_{2i} &= 0, \\ &\vdots \\ \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki}) X_{ki} &= 0, \end{aligned}$$

を満たす。

さらに,

$$\begin{aligned} \sum_{i=1}^n X_{1i} Y_i &= \hat{\beta}_1 \sum_{i=1}^n X_{1i}^2 + \hat{\beta}_2 \sum_{i=1}^n X_{1i} X_{2i} + \dots + \hat{\beta}_k \sum_{i=1}^n X_{1i} X_{ki}, \\ \sum_{i=1}^n X_{2i} Y_i &= \hat{\beta}_1 \sum_{i=1}^n X_{1i} X_{2i} + \hat{\beta}_2 \sum_{i=1}^n X_{2i}^2 + \dots + \hat{\beta}_k \sum_{i=1}^n X_{2i} X_{ki}, \\ &\vdots \\ \sum_{i=1}^n X_{ki} Y_i &= \hat{\beta}_1 \sum_{i=1}^n X_{1i} X_{ki} + \hat{\beta}_2 \sum_{i=1}^n X_{2i} X_{ki} + \dots + \hat{\beta}_k \sum_{i=1}^n X_{ki}^2, \end{aligned}$$

行列表示によって,

$$\begin{pmatrix} \sum X_{1i} Y_i \\ \sum X_{2i} Y_i \\ \vdots \\ \sum X_{ki} Y_i \end{pmatrix} = \begin{pmatrix} \sum X_{1i}^2 & \sum X_{1i} X_{2i} & \cdots & \sum X_{1i} X_{ki} \\ \sum X_{1i} X_{2i} & \sum X_{2i}^2 & \cdots & \sum X_{2i} X_{ki} \\ \vdots & \vdots & \ddots & \vdots \\ \sum X_{1i} X_{ki} & \sum X_{2i} X_{ki} & \cdots & \sum X_{ki}^2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix}$$

が得られ,  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$  についてまとめると,

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} = \begin{pmatrix} \sum X_{1i}^2 & \sum X_{1i} X_{2i} & \cdots & \sum X_{1i} X_{ki} \\ \sum X_{1i} X_{2i} & \sum X_{2i}^2 & \cdots & \sum X_{2i} X_{ki} \\ \vdots & \vdots & \ddots & \vdots \\ \sum X_{1i} X_{ki} & \sum X_{2i} X_{ki} & \cdots & \sum X_{ki}^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum X_{1i} Y_i \\ \sum X_{2i} Y_i \\ \vdots \\ \sum X_{ki} Y_i \end{pmatrix}$$

を解くことになる。⇒ コンピュータによって計算

### 4.3.1 重回帰モデルにおける回帰係数の意味

結論：他の変数の影響を取り除いての被説明変数への影響を表す。

$k = 2$  の単純なモデル：

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad i = 1, 2, \dots, n$$

$\beta_1, \beta_2$  の最小二乗推定量は,

$$\min_{\beta_1, \beta_2} \sum_{i=1}^n (Y_i - \beta_1 X_{1i} - \beta_2 X_{2i})^2$$

を解いて、 $\hat{\beta}_1, \hat{\beta}_2$  が次のように得られる。

$$\begin{aligned} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} &= \begin{pmatrix} \sum X_{1i}^2 & \sum X_{1i}X_{2i} \\ \sum X_{1i}X_{2i} & \sum X_{2i}^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum X_{1i}Y_i \\ \sum X_{2i}Y_i \end{pmatrix} \\ &= \frac{1}{(\sum X_{1i}^2)(\sum X_{2i}^2) - (\sum X_{1i}X_{2i})^2} \begin{pmatrix} \sum X_{2i}^2 & -\sum X_{1i}X_{2i} \\ -\sum X_{1i}X_{2i} & \sum X_{1i}^2 \end{pmatrix} \begin{pmatrix} \sum X_{1i}Y_i \\ \sum X_{2i}Y_i \end{pmatrix} \\ &= \begin{pmatrix} \frac{(\sum X_{2i}^2)(\sum X_{1i}Y_i) - (\sum X_{1i}X_{2i})(\sum X_{2i}Y_i)}{(\sum X_{1i}^2)(\sum X_{2i}^2) - (\sum X_{1i}X_{2i})^2} \\ \frac{-(\sum X_{1i}X_{2i})(\sum X_{1i}Y_i) + (\sum X_{1i}^2)(\sum X_{2i}Y_i)}{(\sum X_{1i}^2)(\sum X_{2i}^2) - (\sum X_{1i}X_{2i})^2} \end{pmatrix} \end{aligned}$$

$\sum_{i=1}^n X_{ji}X_{li}$ ,  $\sum_{i=1}^n X_{ji}Y_i$  をそれぞれ  $\sum X_{ji}X_{li}$ ,  $\sum X_{ji}Y_i$  と表記する。

ただし、 $j = 1, 2$ ,  $l = 1, 2$  とする。

一方、次の2つの回帰式を考える。

$$Y_i = \alpha_1 X_{2i} + v_i$$

$$X_{1i} = \alpha_2 X_{2i} + w_i$$

$\alpha_1, \alpha_2$  のそれぞれの最小二乗推定量を求めると、

$$\hat{\alpha}_1 = \frac{\sum X_{2i}Y_i}{\sum X_{2i}^2}, \quad \hat{\alpha}_2 = \frac{\sum X_{2i}X_{1i}}{\sum X_{2i}^2}$$



となる。

$\hat{\alpha}_1, \hat{\alpha}_2$  を用いて、残差  $\hat{v}_i, \hat{w}_i$  を下記のようにそれぞれ求める。

$$\hat{v}_i = Y_i - \hat{\alpha}_1 X_{2i}, \quad \hat{w}_i = X_{1i} - \hat{\alpha}_2 X_{2i}$$

$\hat{v}_i, \hat{w}_i$  は  $Y_i, X_{1i}$  から  $X_{2i}$  の影響を取り除いたものと解釈できる。

更に、次の回帰式を考える。

$$\hat{v}_i = \gamma \hat{w}_i + \epsilon_i$$

$\gamma$  の最小二乗推定量  $\hat{\gamma}$  は  $\hat{\beta}_1$  に一致することを示す。

$$\begin{aligned} \hat{\gamma} &= \frac{\sum \hat{w}_i \hat{v}_i}{\sum \hat{w}_i^2} = \frac{\sum (X_{1i} - \hat{\alpha}_2 X_{2i})(Y_i - \hat{\alpha}_1 X_{2i})}{\sum (X_{1i} - \hat{\alpha}_2 X_{2i})^2} \\ &= \frac{\sum X_{1i} Y_i - \hat{\alpha}_1 \sum X_{1i} X_{2i} - \hat{\alpha}_2 \sum X_{2i} Y_i + \hat{\alpha}_1 \hat{\alpha}_2 \sum X_{2i}^2}{\sum X_{1i}^2 - 2\hat{\alpha}_2 \sum X_{1i} X_{2i} + \hat{\alpha}_2^2 \sum X_{2i}^2} \\ &= \frac{\sum X_{1i} Y_i - \frac{(\sum X_{2i} Y_i)(\sum X_{1i} X_{2i})}{\sum X_{2i}^2}}{\sum X_{1i}^2 - \frac{(\sum X_{1i} X_{2i})^2}{\sum X_{2i}^2}} \\ &= \frac{(\sum X_{2i}^2)(\sum X_{1i} Y_i) - (\sum X_{1i} X_{2i})(\sum X_{2i} Y_i)}{(\sum X_{1i}^2)(\sum X_{2i}^2) - (\sum X_{1i} X_{2i})^2} = \hat{\beta}_1, \end{aligned}$$

「 $Y_i$  から  $X_{2i}$  の影響を取り除いた変数」を被説明変数、「 $X_{1i}$  から  $X_{2i}$  の影響を取り除いた変数」を説明変数とした回帰係数が  $\beta_1$  に等しい。

一般化： 次の回帰モデルを考える。

$$Y_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki}$$

$j$  番目の回帰係数  $\beta_j$  の意味は、「 $Y_i$  から  $X_{1i}, \dots, X_{j-1,i}, X_{j+1,i}, \dots, X_{ki}$  (すなわち、 $X_{ji}$  以外の説明変数) の影響を取り除いた変数」を被説明変数、「 $X_{ji}$  から  $X_{1i}, \dots, X_{j-1,i}, X_{j+1,i}, \dots, X_{ki}$  (すなわち、 $X_{ji}$  以外の説明変数) の影響を取り除いた変数」を説明変数とした回帰係数となる。

### 4.3.2 決定係数 $R^2$ と自由度修正済み決定係数 $\bar{R}^2$ について

また、決定係数  $R^2$  についても同様に表される。

$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

ただし、 $\hat{Y}_i = \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}$ 、 $Y_i = \hat{Y}_i + \hat{u}_i$  である。

$R^2$  は、説明変数を増やすことによって、必ず大きくなる。なぜなら、説明変数が増えることによって、 $\sum_{i=1}^n \hat{u}_i^2$  が必ず減少するからである。

$R^2$  を基準にすると、被説明変数にとって意味のない変数でも、説明変数が多いほど、よりよいモデルということになる。この点を改善するために、自由度修正済み決定係数  $\bar{R}^2$  を用いる。

$$\bar{R}^2 = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2 / (n - k)}{\sum_{i=1}^n (Y_i - \bar{Y})^2 / (n - 1)},$$

$\sum_{i=1}^n \hat{u}_i^2 / (n - k)$  は  $u_i$  の分散  $\sigma^2$  の不偏推定量であり、 $\sum_{i=1}^n (Y_i - \bar{Y})^2 / (n - 1)$  は  $Y_i$  の分散の不偏推定量である。分散や不偏推定量の意味は、統計学の知識を必要とし、後述する。

$R^2$  と  $\bar{R}^2$  との関係は、

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - k},$$

となる。さらに、

$$\frac{1 - \bar{R}^2}{1 - R^2} = \frac{n - 1}{n - k} \geq 1,$$

という関係から、 $\bar{R}^2 \leq R^2$  という結果を得る。(  $k = 1$  のときのみ、等号が成り立つ。 )

数値例： 今までと同じ数値例で、 $\bar{R}^2$  を計算する。

$i$	$X_i$	$Y_i$	$X_i^2$	$X_i Y_i$	$\hat{Y}_i$	$\hat{u}_i$	$X_i \hat{u}_i$	$\hat{Y}_i \hat{u}_i$	$\hat{u}_i^2$	$Y_i^2$
1	5	4	25	20	4.0	0.0	0.0	0.00	0.00	16
2	1	1	1	1	1.2	-0.2	-0.2	-0.24	0.04	1
3	3	1	9	3	2.6	-1.6	-4.8	-4.16	2.56	1
4	2	3	4	6	1.9	1.1	2.2	2.09	1.21	9
5	4	4	16	16	3.3	0.7	2.8	2.31	0.49	16
合計	$\sum X_i$ 15	$\sum Y_i$ 13	$\sum X_i^2$ 55	$\sum X_i Y_i$ 46	$\sum \hat{Y}_i$ 13	$\sum \hat{u}_i$ 0.0	$\sum X_i \hat{u}_i$ 0.0	$\sum \hat{Y}_i \hat{u}_i$ 0.0	$\sum \hat{u}_i^2$ 4.3	$\sum Y_i^2$ 43
平均	$\bar{X}$ 3	$\bar{Y}$ 2.6								

$\bar{Y} = 2.6$ ,  $\sum_{i=1}^n \hat{u}_i^2 = 4.3$ ,  $\sum_{i=1}^n Y_i^2 = 43$  なので,

$$R^2 = 1 - \frac{\sum \hat{u}_i^2}{\sum Y_i^2 - n\bar{Y}^2} = 1 - \frac{4.3}{43 - 5 \times 2.6^2} = 1 - \frac{4.3}{9.2} = 0.5326$$

となり,  $\bar{R}^2$  は,

$$\bar{R}^2 = 1 - \frac{\sum \hat{u}_i^2 / (n - k)}{(\sum Y_i^2 - n\bar{Y}^2) / (n - 1)} = 1 - \frac{4.3 / (5 - 2)}{9.2 / (5 - 1)} = 0.3768$$

となる。

自由度について： 分子について、残差  $\hat{u}_i$  を求めるためには、 $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$  の  $k$  個の推定値を得なければならない。データ数  $n$  から推定値の数  $k$  を差し引いたものを自由度 (degree of freedom) と呼ぶ。

一方、分母については、 $X_{1i}$  が定数項だとして、 $Y_i$  が定数項を除く  $X_{2i}, X_{3i}, \dots, X_{ki}$  に依存しない場合を考える。この場合、 $\beta_2 = \beta_3 = \dots = \beta_k = 0$  とするので、 $\hat{u}_i = Y_i - \hat{\beta}_1$  となる。 $\hat{u}_i$  を得るためには  $\hat{\beta}_1$  だけを求めればよい。最小二乗法の考え方に沿って求めれば、 $\hat{\beta}_1 = \bar{Y}$  となる (読者は確認すること)。すなわち、自由度は「データ数 - 推定値の数 =  $n - 1$ 」ということになる。

このように、決定係数の第二項目の分子・分母をそれぞれの自由度で割ることによって、自由度修正済み決定係数が得られる。

注意：  $R^2$  や  $\bar{R}^2$  を比較する場合，被説明変数が同じであることが重要である。被説明変数が対数かまたはそのままの値であれば，決定係数・自由度修正済み決定係数の大小比較は意味をなさない。ただし，被説明変数が異なる場合であっても，被説明変数を上昇率とするかそのままの値を用いるかの比較では，決定係数・自由度修正済み決定係数の大小比較はできないが，誤差項  $u_i$  の標準誤差での比較は可能である（標準誤差の小さいモデルを採用する）。⇒ 関数型の選択