

# 重回帰モデル

N組のデータ ( $y_i, x_{i1}, x_{i2}, \dots, x_{ip}$ ) ;  $i = 1, \dots, N$ 、において、 $y_i$ を、次の1次式で表わすことを考えます。

$$y_i \approx b_1 x_{i1} + b_2 x_{i2} + \dots + b_p x_{ip}$$

記号  $\approx$  は、近似式であることを表わします。

上式には定数項が明記されていませんが、これは、例えば、

$$x_{ip} = 1$$

とすることによって、定数項をもつ場合も含まれていることがわかります。

近似式の、データ全体についての精度を、誤差の2乗和 $Q_e$ で表わします。すなわち、

$$Q_e = \sum_{i=1}^N \{y_i - (b_1 x_{i1} + \dots + b_p x_{ip})\}^2$$

$Q_e$ を最小にする $b_j$ を求めるために、 $Q_e$ の $b_j$ に関する偏導関数を求めます。

$$\begin{aligned} \frac{\partial Q_e}{\partial b_j} &= \sum_{i=1}^N 2\{y_i - (b_1 x_{i1} + \dots + b_p x_{ip})\}(-x_{ij}) \\ &= -2 \sum_{i=1}^N y_i x_{ij} + 2 \sum_{i=1}^N (b_1 x_{i1} x_{ij} + \dots + b_p x_{ip} x_{ij}) \end{aligned}$$

上の偏導関数を0とおくと、次式が得られます。

$$\begin{aligned} \sum_{i=1}^N y_i x_{ij} &= \sum_{i=1}^N (b_1 x_{i1} x_{ij} + \dots + b_p x_{ip} x_{ij}) \\ &= \sum_{i=1}^N \sum_{k=1}^p x_{ij} x_{ik} b_k \end{aligned} \quad (1)$$

上式を行列を使った式で表わすために、

$$y = [y_1, \dots, y_N]^T$$

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ x_{21} & \cdots & x_{2p} \\ \vdots & & \vdots \\ x_{N1} & \cdots & x_{Np} \end{bmatrix}$$

$$b = [b_1, \cdots, b_p]'$$

とおきます。

なお、ベクトルあるいは行列の右肩に  $'$  または  $'$  を付けたものは、転置行列を表わします。

このとき、式 (1) は次のように表わせます。

$$X'y = X'Xb \quad (2)$$

$X'X$  が正則である、すなわち逆行列  $(X'X)^{-1}$  が存在するとき、

$$b = (X'X)^{-1}X'y$$

となります。

なお、最初の 1 次式を、行列を使って表わすと、

$$y \approx Xb$$

となります。

いま、 $X$  の第  $p$  列目がすべて 1 であるとします (cf. 佐和, 1979)。

$$X = [X_1, \mathbf{1}_N]$$

と  $X$  を分割して表わします。 $\mathbf{1}_N$  は、1 を  $N$  個並べた列ベクトルを表わします。

$X$  に合わせて、 $b$  も

$$b = \begin{bmatrix} b_0 \\ b_p \end{bmatrix}$$

と分割します。 $b_0$  は  $(p - 1)$  次元列ベクトル、 $b_p$  は  $b$  の  $p$  番目の要素 (スカラー) です。

$X$  が上の形のときは、 $y$  の成分  $y_i$  についてみると

$$y_i \approx b_1 x_{i1} + b_2 x_{i2} + \cdots + b_{p-1} x_{i,p-1} + b_p$$

となっています。これは、1 次式が定数項  $b_p$  をもつ形です。

このとき、式 (4.2) は次のようになります。

$$\begin{aligned} [X_1, \mathbf{1}_N]' y &= [X_1, \mathbf{1}_N]' \begin{bmatrix} b_0 \\ b_p \end{bmatrix} \\ &= \begin{bmatrix} X_1' X_1 & X_1' \mathbf{1}_N \\ \mathbf{1}_N' X_1 & \mathbf{1}_N' \mathbf{1}_N \end{bmatrix} \begin{bmatrix} b_0 \\ b_p \end{bmatrix} \end{aligned}$$

すなわち、

$$\begin{cases} X_1' y = X_1' X_1 b_0 + X_1' \mathbf{1}_N b_p \\ \mathbf{1}_N' y = \mathbf{1}_N' X_1 b_0 + \mathbf{1}_N' \mathbf{1}_N b_p \end{cases} \quad (3)$$

また、

$$\begin{aligned} \mathbf{1}_N' X b &= \mathbf{1}_N' [X_1, \mathbf{1}_N] \begin{bmatrix} b_0 \\ b_p \end{bmatrix} \\ &= \mathbf{1}_N' (X_1 b_0 + \mathbf{1}_N b_p) \\ &= \mathbf{1}_N' y \end{aligned} \quad (\text{式 (3) より})$$

すなわち、 $y$  の予測式

$$\hat{y} = Xb$$

の平均  $\bar{m}_y$  が、 $y$  の平均  $\bar{m}_y$  に等しい、という次式が成り立ちます。

$$\bar{m}_y = \frac{1}{N} \mathbf{1}_N' X b = \frac{1}{N} \mathbf{1}_N' y = \bar{m}_y \quad (4)$$

いま、 $y$  の総変動分 (平方和)

$$Q_y = (y - \bar{m}_y \mathbf{1}_N)' (y - \bar{m}_y \mathbf{1}_N)$$

を、次のように分解します。

$$Q_y = (y - \bar{m}_y \mathbf{1}_N)' (y - \bar{m}_y \mathbf{1}_N)$$

$$\begin{aligned}
 &= (y - Xb + Xb - \mathbf{m}_y \mathbf{1}_N)'(y - Xb + Xb - \mathbf{m}_y \mathbf{1}_N) \\
 &= (y - Xb)'(y - Xb) + (y - Xb)'(Xb - \mathbf{m}_y \mathbf{1}_N) \\
 &\quad + (Xb - \mathbf{m}_y \mathbf{1}_N)'(y - Xb) + (Xb - \mathbf{m}_y \mathbf{1}_N)'(Xb - \mathbf{m}_y \mathbf{1}_N)
 \end{aligned}$$

ここで、

$$(y - Xb)'(Xb - \mathbf{m}_y \mathbf{1}_N) = y'Xb - y'\mathbf{m}_y \mathbf{1}_N - b'X'Xb + b'X'\mathbf{m}_y \mathbf{1}_N$$

上式に、

$$b = (X'X)^{-1}X'y$$

を代入すると、

$$\begin{aligned}
 \text{上式} &= y'Xb - y'\mathbf{m}_y \mathbf{1}_N - y'X(X'X)^{-1}X'Xb + b'X'\mathbf{m}_y \mathbf{1}_N \\
 &= -y'\mathbf{m}_y \mathbf{1}_N + b'X'\mathbf{m}_y \mathbf{1}_N \\
 &= -y'\mathbf{m}_y \mathbf{1}_N + \mathbf{m}_y(1'_N Xb)' \\
 &= -y'\mathbf{m}_y \mathbf{1}_N + \mathbf{m}_y(1'_N y)' \quad (\text{式(4)より}) \\
 &= 0
 \end{aligned}$$

同様に、

$$(Xb - \mathbf{m}_y \mathbf{1}_N)'(y - Xb) = 0$$

となります。

したがって、 $Q_y$  は次のように分解できます。

$$\begin{aligned}
 Q_y &= (y - Xb)'(y - Xb) + (Xb - \mathbf{m}_y \mathbf{1}_N)'(Xb - \mathbf{m}_y \mathbf{1}_N) \\
 &= (y - Xb)'(y - Xb) + (Xb - \mathbf{m}_y \mathbf{1}_N)'(Xb - \mathbf{m}_y \mathbf{1}_N)
 \end{aligned}$$

すなわち、総変動平方和

$$Q_y = (y - \mathbf{m}_y \mathbf{1}_N)'(y - \mathbf{m}_y \mathbf{1}_N)$$

は、残差平方和

$$Q_e = (y - Xb)'(y - Xb) = (y - \hat{y})'(y - \hat{y})$$

と、1 次式  $\hat{y} = Xb$  で説明される変動平方和

$$Q_{\hat{y}} = (Xb - \mathbf{m}_{\hat{y}} \mathbf{1}_N)'(Xb - \mathbf{m}_{\hat{y}} \mathbf{1}_N) = (\hat{y} - \mathbf{m}_{\hat{y}} \mathbf{1}_N)'(\hat{y} - \mathbf{m}_{\hat{y}} \mathbf{1}_N)$$

の和に分解できます。

$$Q_y = Q_e + Q_{\hat{y}}$$

$Q_{\hat{y}}$  と  $Q_y$  の比

$$R^2 = \frac{Q_{\hat{y}}}{Q_y} = 1 - \frac{Q_e}{Q_y}$$

は決定係数と呼ばれているもので、その正の平方根  $R$  は重相関係数と呼ばれています(佐和、1979)。

$R^2$  が 1 に近いほど、 $Q_{\hat{y}}$  は  $Q_y$  に近い、すなわち、誤差  $Q_e$  が  $Q_y$  に比べて相対的に小さい値であることを表わします。

プログラム PMultiIndVar.dpr は、上の重回帰モデルによる分析を行うものです。

このプログラムのフォームは、プログラミング時には図 1 のように用意されています。

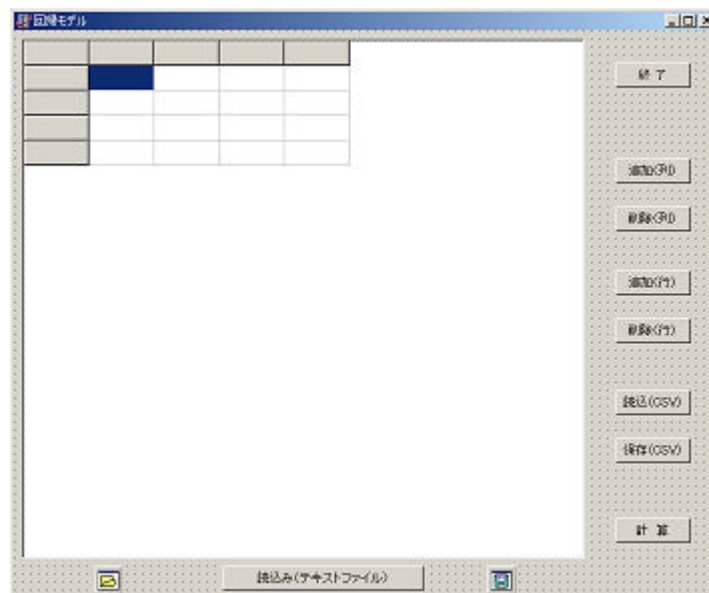


図 1 プログラミング時のフォーム

プログラムを実行すると、図 2 のようなフォームになります。

図2 実行開始時のフォーム

「追加」「削除」ボタンをクリックして、行および列を必要な数に調整します。変数 0 の欄に  $y_i$ 、変数 1 の欄に  $x_{i1}$ 、・・・、変数 p の欄に  $x_{ip}$  を書きます。i 番目のデータ、 $y_i$ 、 $x_{i1}$ 、・・・、 $x_{ip}$  は「i 番目」の行に書きます。「変数名」の行には変数のラベルを書きますが、変数名は空白でもかまいません。

「読み込み(テキストファイル)」ボタンをクリックすると、リスト 1 のような形式で用意されたテキストファイルデータを読み込むことができます。

#### リスト 1 入力データ例(テキストファイル)

-100			
2			
国語			
英語			
切片			
	72	66	1
	41	46	1
	93	87	1
	82	82	1
	41	42	1
	64	72	1
	27	31	1
	68	69	1
	79	84	1
	67	68	1
	57	50	1
	90	90	1
	87	85	1
	29	29	1
	75	76	1
	66	69	1
	78	84	1
	31	35	1
-1000			

上のデータは、次式

$$\text{国語} = a \times \text{英語} + b + \text{誤差項}$$

を、

$$y_i = \text{国語}、\quad x_{i1} = \text{英語}、\quad x_{i2} = 1$$

とにおいて、

$$y_i = \text{国語} = [\text{英語}, 1] \begin{bmatrix} a \\ b \end{bmatrix} + \text{誤差項} = [x_{i1}, x_{i2}] \begin{bmatrix} a \\ b \end{bmatrix} + \text{誤差項}$$

の形として、データ値

$$\text{国語の点数}、\quad \text{英語の点数}、\quad \text{定数 } 1$$

を並べたものです。

すなわち、リスト1の入力データは次の形式で用意します。

まず、第1行目にデータ読み終了の判定に用いる基準値を書きます。行の先頭の数値がこの1行目の値より小さいところで読みを終了します。リスト1では、1行目に100が書かれていて、最後の1000がデータの終わりを示しています。2行目には、独立変数

$x_{i1}$ 、 $\dots$ 、 $x_{ip}$  の数  $p$  を書きます。リスト 1 では 2 が書かれています。3 行目から  $(3 + p)$  行目までに、各変数  $y_i$ 、 $x_{i1}$ 、 $\dots$ 、 $x_{ip}$  のラベルを 1 行に 1 つずつ書きます。3 行目に  $y_i$  のラベル、4 行目に  $x_{i1}$ 、 $\dots$ 、 $(3 + p)$  行目に  $x_{ip}$  のラベルを書きます。リスト 1 では、「国語」、「英語」、「切片」が順番に書かれています。

$(4 + p)$  行目からデータの組  $\{ y_i, x_{i1}, \dots, x_{ip} \}$  次の形式

$$y_i \quad x_{i1} \quad \dots \quad x_{ip}$$

で、1 行に 1 組ずつ書きます。

最後の組

$$y_N \quad x_{N1} \quad \dots \quad x_{Np}$$

の次の行に、第 1 行に書いた数値より小さい数を書いてデータの終わりとします。リスト 1 では 1000 が書かれています。

リスト 1 のファイルを読み込むと図 3 のようになります。



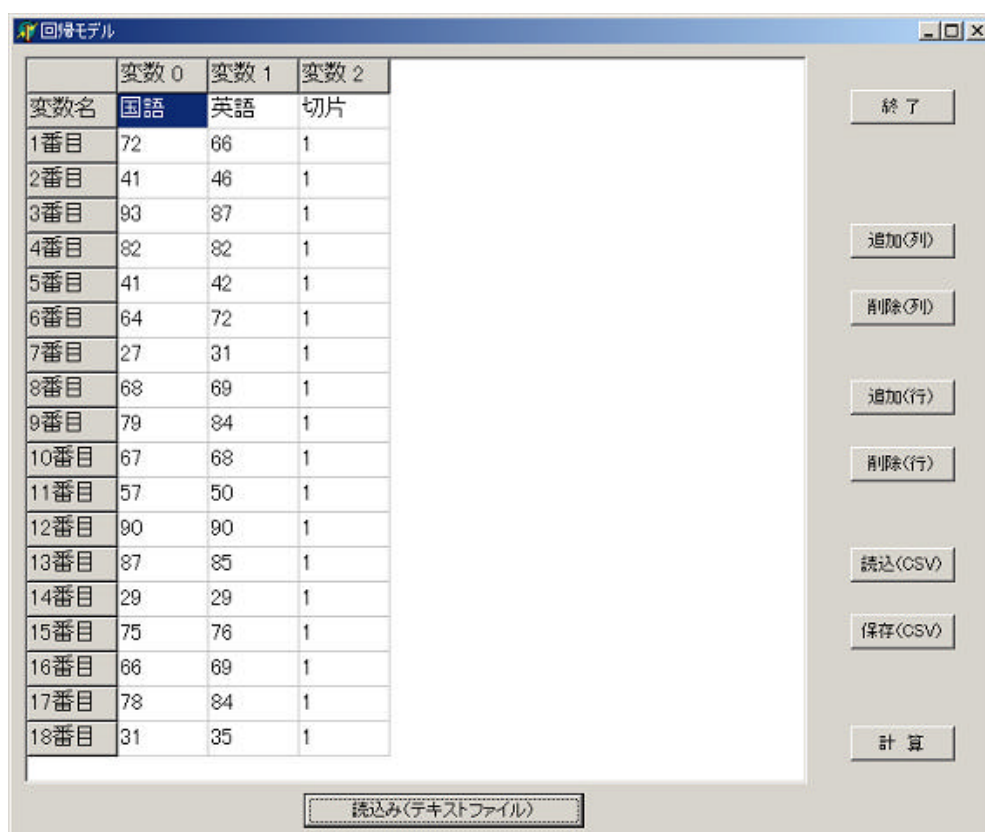
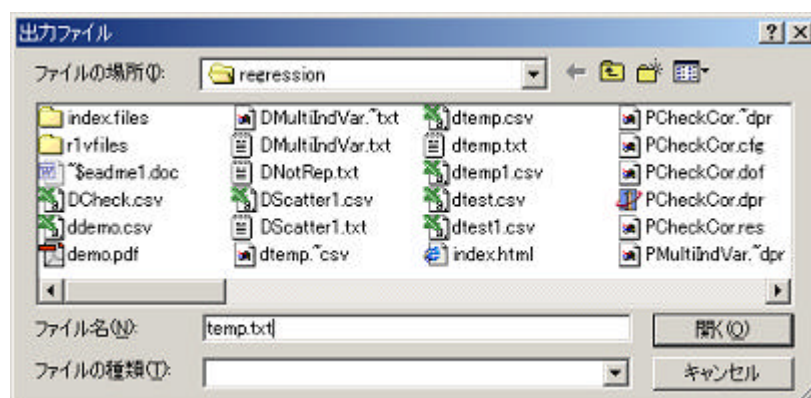


図3 データの設定例

上図のようにデータを設定後、「計算」ボタンをクリックすると計算が始まります。「計算」ボタンをクリックすると、まず計算結果を出力するためのファイル名の設定を求めるダイアログボックスが表示されます。



計算結果はテキストファイルとして書き出されるので、プログラムの実行終了後エディタなどで開いて見ることができます。

ファイル名の設定後「開く」ボタンをクリックすると、計算は一瞬で終わります。「終了」ボタンをクリックして実行を終了します。リスト 1 のデータの場合、出力ファイルの内容はリスト 2 のようになっています。

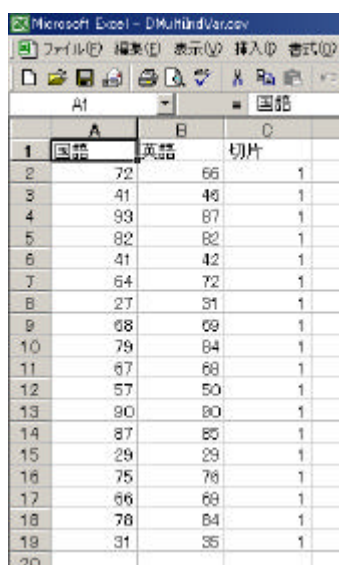
リスト 2 リスト 1 のデータの分析結果

国語	英語	切片
72.000	66.000	1.000
41.000	46.000	1.000
93.000	87.000	1.000
82.000	82.000	1.000
41.000	42.000	1.000
64.000	72.000	1.000
27.000	31.000	1.000
68.000	69.000	1.000
79.000	84.000	1.000
67.000	68.000	1.000
57.000	50.000	1.000
90.000	90.000	1.000
87.000	85.000	1.000
29.000	29.000	1.000
75.000	76.000	1.000
66.000	69.000	1.000
78.000	84.000	1.000
31.000	35.000	1.000
B =		
	1.02042344 <英語>	
	-2.32185061 <切片>	
Qe	= 298.997754	
Qy	= 7793.61111	
決定係数	= 0.961635531	
重相関係数	= 0.98063017	

「B =」の次行から、係数ベクトル  $\mathbf{b} = [b_1, \dots, b_p]$  の各要素の値が 1 行に一つずつ書かれています。リスト 2 に示されている値は、次の回帰式を示しています。

$$\text{国語} = 1.02 \times \text{英語} - 2.32 + \text{誤差項}$$

図 3 のように設定されたデータは、「保存 (CSV)」ボタンをクリックすると CSV 形式で保存することができます。この CSV 形式のデータを Excel で読み込むと図 4 のようになります。



	A	B	C
1	国語	英語	切片
2	72	66	1
3	41	48	1
4	93	87	1
5	82	82	1
6	41	42	1
7	64	72	1
8	27	31	1
9	68	69	1
10	79	84	1
11	67	68	1
12	57	50	1
13	90	80	1
14	87	85	1
15	29	29	1
16	75	76	1
17	66	69	1
18	78	84	1
19	31	35	1
20			

図4 Excelで読み込んだデータ

Excelで図4の形式で作成したデータは、CSV形式で保存すれば「読込（CSV）」ボタンをクリックすることによりプログラム PMultiIndVar.dpr に読み込むことができます。

## 参考文献

- (1) 佐和隆光 (1979) 回帰分析． 朝倉書店．
- (2) 岡本安晴 (1998) Delphiで学ぶデータ分析法． CQ 出版株式会社．