

回 帰 分 析

重回帰分析のための確率モデルについて説明します。

Gauss-Markoff setup (Timm, 1975) と呼ばれているモデル

$$y = X\mathbf{b} + \mathbf{e}$$

において、 X と \mathbf{b} は定数部分で、それらの積 $X\mathbf{b}$ からのランダムな変動分 \mathbf{e} を加えたものとして、 y が表わされています。

$$\mathbf{e} = [\mathbf{e}_1, \dots, \mathbf{e}_N]'$$

において、各 \mathbf{e}_i は、平均 0、分散 \mathbf{s}^2 の正規分布に従い、お互いに独立であるとしています。

X が N 行 q 列の行列で階数が r であるとしています。このとき、残差平方和 Q_e を \mathbf{s}^2 で割ったものは、自由度 $N - r$ のカイ 2 乗分布に従います。

$$\frac{Q_e}{\mathbf{s}^2} \sim \mathbf{c}_{N-p}^2 \quad (1)$$

ここで、 Q_e は次式

$$\begin{aligned} Q_e &= \min_{\mathbf{b}} (\mathbf{y} - X\mathbf{b})'(\mathbf{y} - X\mathbf{b}) \\ &= (\mathbf{y} - X\mathbf{b})'(\mathbf{y} - X\mathbf{b}) \end{aligned}$$

で与えられ、 \mathbf{b} は次式

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

で与えられる \mathbf{b} の推定値です。ここで、 $(X'X)^{-}$ は $X'X$ の一般逆行列を表しています。

Q_e を求めるときは、 $(\mathbf{y} - X\mathbf{b})'(\mathbf{y} - X\mathbf{b})$ を最小にするという最小 2 乗基準で \mathbf{b} の推定を行っています。 $\mathbf{b} = [\mathbf{b}_1, \dots, \mathbf{b}_p]'$ に対するある制約、例えば $\mathbf{b}_1 = 0$ というような制約のもとで \mathbf{b} の推定を行うときは、その制約条件を満たす範囲で \mathbf{b} の値を求めることになります。

制約条件として、次の 1 次式で表わされるものについて考えます。

$$c_1 \mathbf{b}_1 + \cdots + c_p \mathbf{b}_p = 0$$

例えば、 $\mathbf{b}_1 = 0$ は

$$1 \cdot \mathbf{b}_1 + 0 \cdot \mathbf{b}_2 + \cdots + 0 \cdot \mathbf{b}_p = 0$$

という形で表わすことができます。

$$\mathbf{c} = [c_1, \cdots, c_p]$$

とおき、 \mathbf{c} をコントラストと呼びます。コントラストはお互いに独立なものを行として複数個まとめて行列

$$\mathbf{C} = \begin{bmatrix} c_{11} & \cdots & c_{1p} \\ \vdots & & \\ c_{g1} & \cdots & c_{gp} \end{bmatrix}$$

の形で与えることができます。

制約条件

$$\mathbf{C}\mathbf{b} = \mathbf{0}$$

もとでの $(\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b})$ の最小値を、 Q_t で表わします。

$$Q_t = \min_{\mathbf{C}\mathbf{b}=\mathbf{0}} (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b})$$

このとき、次式が成り立ちます。

$$Q_h = Q_t - Q_e = (\mathbf{C}\mathbf{b})'[C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\mathbf{C}\mathbf{b})$$

Q_h は、 $\mathbf{C}\mathbf{b} = \mathbf{0}$ という制約を付けたために増えた、残差平方和 Q_e の増加分です。この Q_h を \mathbf{s}^2 で割ったものは、仮説「 $\mathbf{C}\mathbf{b} = \mathbf{0}$ 」が真であるとき、自由度 g のカイ 2 乗分布に従い、自由度 Q_e / \mathbf{s}^2 と独立です。

$$\frac{Q_h}{\mathbf{s}^2} \sim \mathbf{c}_g^2 \quad (2)$$

(1) 式と (2) 式より、帰無仮説「 $\mathbf{C}\mathbf{b} = \mathbf{0}$ 」のもとで、

$$F = \frac{Q_h / s^2}{(Q_e / s^2) / (N - p)}$$

$$= \frac{Q_h}{Q_e / (N - p)}$$

は、自由度 (g、N - r) の F 分布に従います。

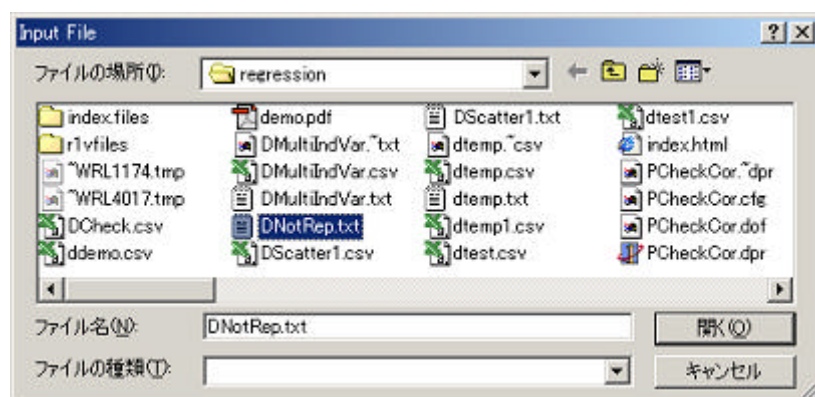
Q_h は、母数 \mathbf{b} の真の値に対して、 $C\mathbf{b}$ が 0 から離れた値であるほど大きい値になりやすいので、F 値が有意水準 α で設定された値 F_α より大きいとき、帰無仮説「 $C\mathbf{b} = 0$ 」は棄却されます。 F_α は、自由度 (g、N - r) の F 分布に従う確率変数を Z とするとき、

$$P(Z \geq F_\alpha) = \alpha$$

を満たす値として求められます。

プログラム PRegU.dpr は、上で説明した重回帰分析を行うものです。

このプログラムを実行すると、下図の入力データ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。



入力データファイルはテキストファイルとしてリスト 1 のように用意します。

リスト1 入力データファイル例

岡本安晴「Delphi で学ぶデータ分析法」C Q 出版社,1998				
表 2.2.4.1.1 より				
/				
3				
83	1	0	0	文学部
84	1	0	0	
74	1	0	0	
58	1	0	0	
73	1	0	0	
74	1	0	0	
66	1	0	0	
85	1	0	0	
61	1	0	0	
77	1	0	0	
65	0	1	0	法学部
80	0	1	0	
81	0	1	0	
81	0	1	0	
80	0	1	0	
72	0	1	0	
56	0	1	0	
64	0	1	0	
70	0	1	0	
66	0	1	0	
77	0	0	1	経済学部
72	0	0	1	
70	0	0	1	
75	0	0	1	
65	0	0	1	
53	0	0	1	
67	0	0	1	
75	0	0	1	
60	0	0	1	
63	0	0	1	
/				
/				
2				
1	-1	0		
1	0	-1		
/				
-1				

リスト1のデータは次の形式に従っています。

行の先頭が/で始まる行を区切り行として用います。

最初の区切り行の次の行に独立変数 x_1 、 \dots 、 x_p の数 p を書きます。リスト1では3が書かれています。独立変数の数を書いた行の次の行からデータ y 、 x_1 、 \dots 、 x_p を1行に1組ずつ書きます。データを各行の先頭は空白文字（半角）を1つ以上置きます。行の

先頭が/で始まる区切り行が現れた前の行でデータは終わります。

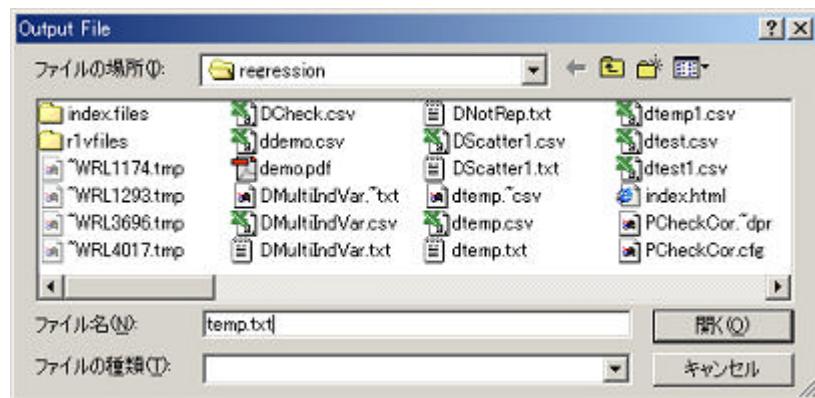
次の区切り行に続けてコントラスト行列を書きます。コントラスト行列は、先ず行数 g を書いてから行単位で書きます。リスト 1 では行数 2 のコントラスト行列

$$C = \begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{bmatrix}$$

が書かれています。

区切り行の次に負の整数が書かれたところでデータファイルの終了とみなされます。

入力データファイル名を設定した後「開く」ボタンをクリックすると、計算結果を書き出す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出力されるので、プログラムの実行終了後適当なエディタで開いて見ることができます。

出力ファイル名の設定後、「開く」ボタンをクリックすると次のフォームが表示されます。

「OK」ボタンのクリックで計算が始まり、計算が終了すると次のフォームになります。

「Exit」ボタンのクリックでプログラムの実行終了となります。計算結果は、リスト1のデータの場合理スト2となっています。

リスト2 計算結果の例

入力データファイル = D:\yasuharu\MyHomePage\jwu\openwww\regression\NotRep.txt		
83.000	1	0 0
84.000	1	0 0
74.000	1	0 0
58.000	1	0 0
73.000	1	0 0
74.000	1	0 0
66.000	1	0 0
85.000	1	0 0
61.000	1	0 0
77.000	1	0 0
65.000	0	1 0
80.000	0	1 0
81.000	0	1 0
81.000	0	1 0
80.000	0	1 0
72.000	0	1 0

56.000	0 1 0
64.000	0 1 0
70.000	0 1 0
66.000	0 1 0
77.000	0 0 1
72.000	0 0 1
70.000	0 0 1
75.000	0 0 1
65.000	0 0 1
53.000	0 0 1
67.000	0 0 1
75.000	0 0 1
60.000	0 0 1
63.000	0 0 1
B =	
73.5	
71.5	
67.7	
Qe	= 2017.1
Qy	= 2190.7
1-Qe/Qy	= 0.0792440772
sqrt(1-Qe/Qy)	= 0.281503246
コントラスト	
1.0	-1.0 0.0
1.0	0.0 -1.0
F = 1.16186605 df1 = 2 df2 = 27	

コントラストの次に、そのコントラストに対する F 値が書き出されています。リスト 2 の場合、「F = 1.16 df1 = 2 df2 = 27」は、リスト 1 の表す 1 要因 3 水準のデータに対する分散分析の結果と一致しています。

参考文献

- (1) Tim, N. H. (1975) *Multivariate Analysis with Applications in Education and Psychology*. Brooks/Cole Publishing Company, Inc.
- (2) 岡本安晴 (1998) Delphi で学ぶデータ分析法 . CQ 出版社 .