

回帰分析（単回帰モデル）

N 個のデータの組、 $(y_i, x_i); i=1 \cdots N$ 、が与えられている場合について考えてみます。

例えば、 y_i が国語の点数、 x_i が英語の点数とします。このとき、 y_i を x_i の一次式

$$ax_i + b$$

で表わすことを考えます。

y_i を x_i の一次式で表わしたときの誤差を e_i とおけば、

$$y_i = ax_i + b + e_i$$

と書けます。

上の一次式は回帰式、 a は回帰係数と呼ばれています。グラフ上で回帰式を表わしたものは、回帰直線と呼ばれています。

一次式における定数 a と b の値を、誤差の2乗 e_i^2 の和が最小になるように求めることを考えます。

2乗和を SSE とおけば、

$$SSE = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - ax_i - b)^2$$

となります。

SSE を最小にする a と b を求めるために、 SSE を a と b で偏微分します。

$$\frac{\partial SSE}{\partial a} = \sum_{i=1}^N 2 \cdot (y_i - ax_i - b) \cdot (-x_i)$$

$$\frac{\partial SSE}{\partial b} = \sum_{i=1}^N 2 \cdot (y_i - ax_i - b) \cdot (-1)$$

$\partial SSE / \partial a = 0$ において、式を変形します。

$$\sum_{i=1}^N (y_i - ax_i - b)x_i = 0 \quad (1)$$

$$a \sum_{i=1}^N x_i^2 + b \sum_{i=1}^N x_i = \sum_{i=1}^N x_i y_i \quad (2)$$

$\sum_{i=1}^N (y_i - ax_i - b) = 0$ において、同じように変形します。

$$\sum_{i=1}^N (y_i - ax_i - b) = 0 \quad (3)$$

$$a \sum_{i=1}^N x_i + bN = \sum_{i=1}^N y_i \quad (4)$$

$(1.5.2) \times N - (1.5.4) \times \sum_{i=1}^N x_i$ より、

$$a \left\{ N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2 \right\} = N \sum_{i=1}^N x_i y_i - \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right)$$

$$a = \frac{N \sum_{i=1}^N x_i y_i - \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right)}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} \quad (5)$$

となります。

式(2) $\times \sum_{i=1}^N x_i -$ 式(4) $\times \sum_{i=1}^N x_i^2$ より、

$$b \left\{ \left(\sum_{i=1}^N x_i \right)^2 - N \sum_{i=1}^N x_i^2 \right\} = \left(\sum_{i=1}^N x_i y_i \right) \left(\sum_{i=1}^N x_i \right) - \left(\sum_{i=1}^N y_i \right) \left(\sum_{i=1}^N x_i^2 \right)$$

$$b = \frac{\left(\sum_{i=1}^N x_i^2 \right) \left(\sum_{i=1}^N y_i \right) - \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N x_i y_i \right)}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2} \quad (6)$$

となります。

次に、一次式 $ax_i + b$ の係数 a および b を、式(5) および (6) で与えたときの誤差の2乗

和 $SSE = \sum \{y_i - (ax_i + b)\}^2$ の大きさについて考えます。

まず、 $ax_i + b$ の平均値を求めてみます。

式(4)の両辺を N で割って

$$a \frac{1}{N} \sum_{i=1}^N x_i + b = \frac{1}{N} \sum_{i=1}^N y_i$$

よって、

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N (ax_i + b) &= a \cdot \frac{1}{N} \left(\sum_{i=1}^N x_i \right) + b \\ &= \frac{1}{N} \sum_{i=1}^N y_i \\ &= \bar{y} \end{aligned}$$

となります。

すなわち

$$\frac{1}{N} \sum_{i=1}^N (ax_i + b) = a\bar{x} + b = \bar{y}$$

です。

但し、

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

とおきます。

次に、

$$s_x^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2,$$

$$s_y^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2,$$

および

$$s_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

とおきます。

s_x^2 、 s_y^2 および s_{xy} は、次のように変形できます。

$$\begin{aligned} s_x^2 &= \frac{1}{N} \left\{ \sum_{i=1}^N x_i^2 - 2\bar{x} \left(\sum_{i=1}^N x_i \right) + N\bar{x}^2 \right\} \\ &= \frac{1}{N} \sum_{i=1}^N x_i^2 - \frac{1}{N^2} \left(\sum_{i=1}^N x_i \right)^2 \end{aligned}$$

同様にして

$$\begin{aligned} s_y^2 &= \frac{1}{N} \sum_{i=1}^N y_i^2 - \frac{1}{N^2} \left(\sum_{i=1}^N y_i \right)^2 \\ s_{xy} &= \frac{1}{N} \left\{ \sum_{i=1}^N x_i y_i - \bar{x} \left(\sum_{i=1}^N y_i \right) - \bar{y} \left(\sum_{i=1}^N x_i \right) + N\bar{x}\bar{y} \right\} \\ &= \frac{1}{N} \left(\sum_{i=1}^N x_i y_i \right) - \frac{1}{N^2} \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right) \end{aligned}$$

従って、式(5)は、次のように書けます。

$$\begin{aligned} a &= \frac{\frac{1}{N} \sum_{i=1}^N x_i y_i - \frac{1}{N^2} \left(\sum_{i=1}^N x_i \right) \left(\sum_{i=1}^N y_i \right)}{\frac{1}{N} \sum_{i=1}^N x_i^2 - \frac{1}{N^2} \left(\sum_{i=1}^N x_i \right)^2} \\ &= \frac{s_{xy}}{s_x^2} \end{aligned} \tag{7}$$

以上より、誤差の2乗和は、次のように計算できます。

$$\begin{aligned} \sum_{i=1}^N \{y_i - (ax_i + b)\}^2 &= \sum_{i=1}^N \{(y_i - \bar{y}) - (ax_i + b - \bar{y})\}^2 \\ &= \sum_{i=1}^N (y_i - \bar{y})^2 - 2 \sum_{i=1}^N (y_i - \bar{y})(ax_i + b - \bar{y}) + \sum_{i=1}^N (ax_i + b - \bar{y})^2 \end{aligned}$$

$$\begin{aligned}
 &= Ns_y^2 - 2 \sum_{i=1}^N (y_i - \bar{y})(ax_i + b - \bar{ax} - b) + \sum_{i=1}^N (ax_i + b - \bar{ax} - b)^2 \\
 &= Ns_y^2 - 2 \sum_{i=1}^N (y_i - \bar{y})a(x_i - \bar{x}) + \sum_{i=1}^N a^2(x_i - \bar{x})^2 \\
 &= Ns_y^2 - 2a \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}) + a^2 \sum_{i=1}^N (x_i - \bar{x})^2 \\
 &= Ns_y^2 - 2 \frac{s_{xy}}{s_x^2} Ns_{xy} + \left(\frac{s_{xy}}{s_x^2} \right)^2 Ns_x^2 \\
 &= Ns_y^2 - N \frac{s_{xy}^2}{s_x^2}
 \end{aligned}$$

従って、誤差の2乗和の平均値の、 y_i の分散 s_y^2 に対する比は、次のようになります。

$$\begin{aligned}
 \frac{\frac{1}{N} \sum_{i=1}^N \{y_i - (ax_i + b)\}^2}{s_y^2} &= \frac{s_y^2 - \frac{s_{xy}^2}{s_x^2}}{s_y^2} \\
 &= 1 - \left(\frac{s_{xy}}{s_x s_y} \right)^2 \\
 &= 1 - r^2
 \end{aligned} \tag{8}$$

ここで、 r は相関係数で、次式で与えられます。

$$\begin{aligned}
 r &= \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}} \\
 &= \frac{s_{xy}}{s_x s_y}
 \end{aligned}$$

式(8)より、 r の絶対値が1に近いほど、誤差の小さいことがわかります。式(8)で与えられる $1 - r^2$ は、非決定係数と呼ばれています。これに対して、 r の2乗 r^2 は、直線回帰の決定係数と呼ばれています。

いま、 y_i の分散 s_y^2 を、次のように分解します。

$$\begin{aligned} Ns_y^2 &= \sum_{i=1}^N (y_i - \bar{y})^2 \\ &= \sum_{i=1}^N \{y_i - (ax_i + b) + (ax_i + b) - \bar{y}\}^2 \\ &= \sum_{i=1}^N \{y_i - (ax_i + b)\}^2 + 2 \sum_{i=1}^N \{y_i - (ax_i + b)\} \{(ax_i + b) - \bar{y}\} + \sum_{i=1}^N \{(ax_i + b) - \bar{y}\}^2 \end{aligned}$$

上式の第2項は、次のように変形できます。

$$\begin{aligned} &\sum_{i=1}^N \{y_i - (ax_i + b)\} \{(ax_i + b) - \bar{y}\} \\ &= \sum_{i=1}^N \{y_i - (ax_i + b)\} ax_i + \sum_{i=1}^N \{y_i - (ax_i + b)\} (b - \bar{y}) \end{aligned} \quad (9)$$

式(9)の第1項は、式(1)を用いると、

$$\begin{aligned} \sum_{i=1}^N \{y_i - (ax_i + b)\} ax_i &= a \sum_{i=1}^N (y_i - ax_i - b)x_i \\ &= 0 \end{aligned}$$

となります。

式(9)の第2項は、式(3)を用いると、

$$\begin{aligned} \sum_{i=1}^N \{y_i - (ax_i + b)\} (b - \bar{y}) &= (b - \bar{y}) \sum_{i=1}^N (y_i - ax_i - b) \\ &= 0 \end{aligned}$$

となります。

したがって、

$$\begin{aligned} Ns_y^2 &= \sum_{i=1}^N \{y_i - (ax_i + b)\}^2 + \sum_{i=1}^N \{(ax_i + b) - \bar{y}\}^2 \\ &= Ns_y^2(1 - r^2) + \sum_{i=1}^N \{(ax_i + b) - \bar{y}\}^2 \quad (\text{式(8)より}) \end{aligned}$$

ゆえに、

$$\frac{\frac{1}{N} \sum_{i=1}^N \{(ax_i + b) - \bar{y}\}^2}{s_y^2} = \frac{\frac{1}{N} \sum_{i=1}^N \{(ax_i + b) - (a\bar{x} + b)\}^2}{s_y^2} = r^2$$

となります。

上式より、決定係数 r^2 は、回帰直線による予測値の分散の、データ値の分散 s_y^2 に対する比になっていることがわかります。

回帰直線を求めるプログラム P0neIndVar.dpr を用意しました。

このプログラムのフォームは、図1のように用意されています。

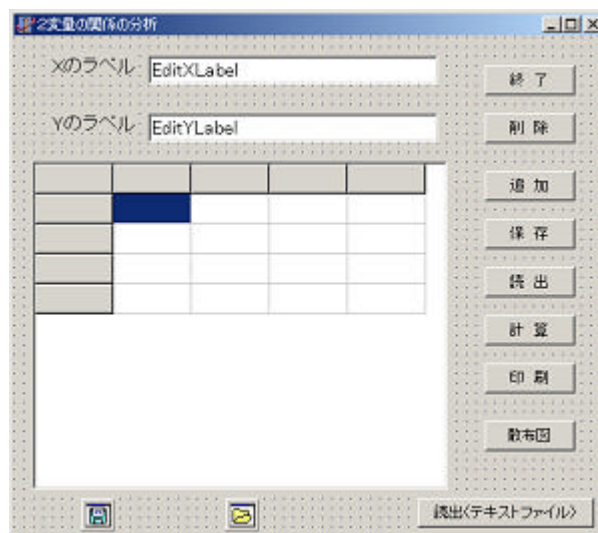


図1 フォームの準備

散布図や回帰直線は、実行時において図1の「散布図」ボタンをクリックしたときに生成・表示されるフォームの Image コンポーネントの Canvas に描きます。このフォームはプログラミング時には次のようになっています。

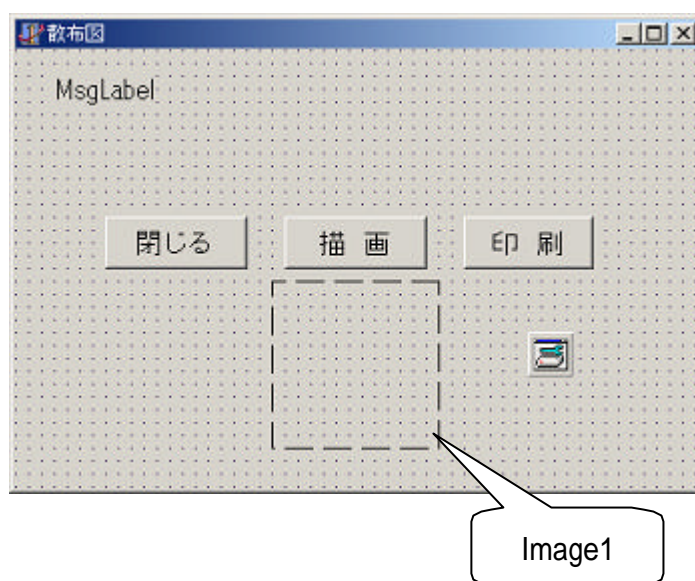


Image1 の Canvas において、第 2 変数を表わす Y 軸のラベルを縦方向に描くために（図 3 における「国語」）、Image コンポーネント Image2 を実行中に生成して、この Canvas に描いたラベル（文字列）を Image1 の Canvas に 90 度回転してコピーすることになります。

プログラムの実行開始時のフォームは、図 2 のようになっています。

	X =	Y =
1番目		

図 2 実行開始時のフォーム

「追加」ボタンを押して行を必要なだけ追加しデータを設定します。「Xのラベル」「Yのラベル」の右側の Edit コンポーネントにはそれぞれの変数の適当なラベル（名前）を設

定します。行を削除するときは、その行のセルをクリックしてアクティブにしてから「削除」ボタンをクリックします。

データはリスト 1 のように用意されているファイルから読み込むことも出来ます。

リスト 1 入力データファイル例

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

第 1 行目に第 1 変数のラベル、第 2 行目に第 2 変数のラベルを置きます。

次に、第 3 行目に、どのデータ値よりも小さい適当な値を、基準値として置きます。リスト 1 では -100 がおかれています。数値の後ろは、データとしては無視されるので、1 つ以上の空白を置いた後、「英語」、「国語」などの文字列がコメントとして置かれています。

4 行目から、各行に 1 組ずつ、データが x_i 、 y_i の順に並べられています。

最後のデータが並べられている行の次の行は、3 行目においた数値より小さい値を置きます。リスト 1 では -1000 が置かれています。

リスト 1 の形式で用意されたテキストファイルを読み込むときは「読出 (テキストファイ

ル)」ボタンをクリックします。下図はリスト1のデータを読み込んだ状態です。

2変量の関係の分析

Xのラベル 英語

Yのラベル 国語

	X =	Y =
1番目	66	72
2番目	46	41
3番目	87	93
4番目	82	82
5番目	42	41
6番目	72	64
7番目	31	27
8番目	69	68
9番目	84	79
10番目	68	67

終了

削除

追加

保存

読出

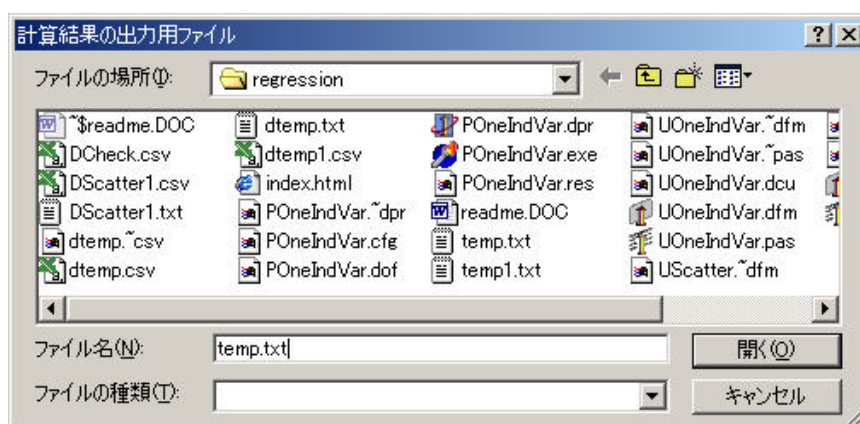
計算

印刷

散布図

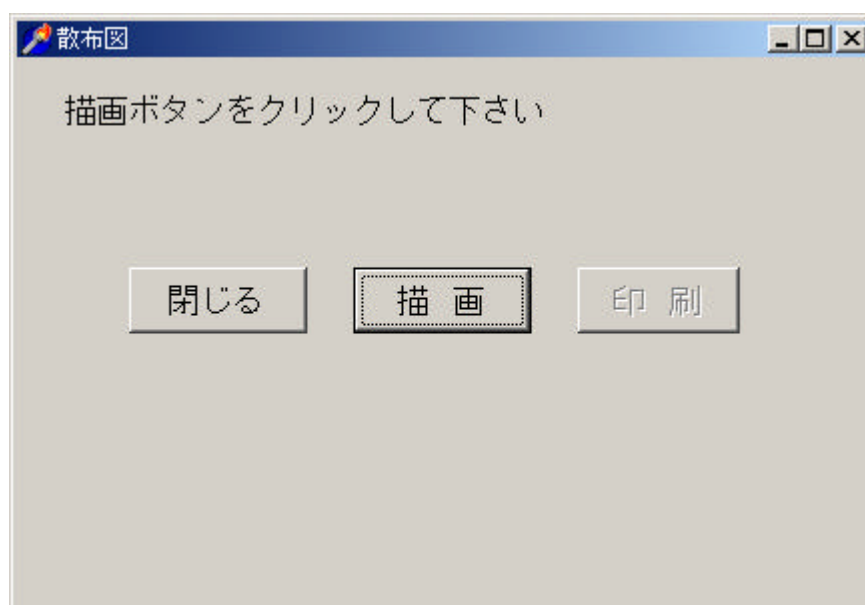
読出(テキストファイル)

データの設定後、「散布図」あるいは「計算」ボタンをクリックすると計算が始まります。「計算」ボタンのクリックでは、各変数の平均値、標準偏差、相関係数などが計算されます。このボタンをクリックすると、まず計算結果出力用のファイル名を求めるダイアログボックスが表示されます。

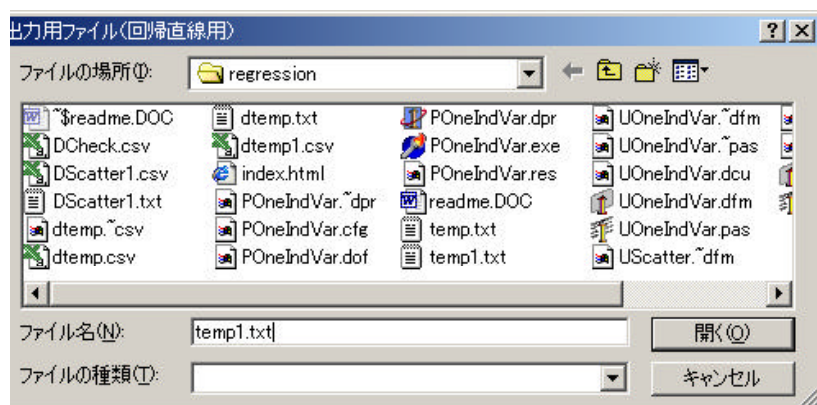


名前の設定後、「開く」ボタンのクリックで計算が始まり、計算結果が設定したファイル名のファイルに書き出されます。書き出されたファイルはテキストファイルなので、プログラムの実行終了後、適当なエディタで開いて見ることが出来ます。

「散布図」ボタンをクリックすると、まず下図のフォームが表示されます。



「描画」ボタンのクリックで回帰曲線の描画が行われますが、その前に計算結果を出力するテキスト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。



ファイル名の設定後、「開く」ボタンをクリックすると計算が始まります。 $x[i]$ 、 $y[i]$ と、直線回帰による値 $a \cdot x[i] + b$ 、誤差 $y[i] - (a \cdot x[i] + b)$ が出力ファイルに書き出され、続いて a と b の値も書き出されます。

次に、相関係数 r が求められて、決定係数 r^2 とともに書き出されてから、出力用ファイルが閉じられます。

出力用ファイルへの書き出しが終わると、散布図と回帰直線が図3のように描画されます。

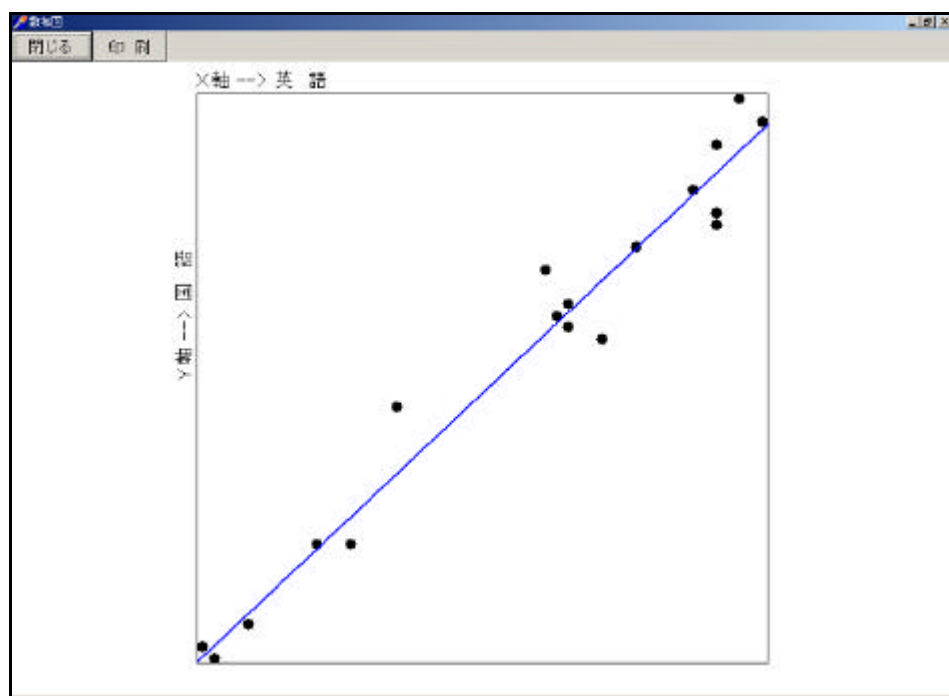


図3 回帰直線と散布図の描画面

図3において、「印刷」ボタンをクリックすると回帰直線と散布図がプリンタに出力され

ます。

出力ファイルは、リスト2のようになっています。

リスト2 出力ファイル

```

変量 X のラベル --> 英 語
変量 Y のラベル --> 国 語

X =      Y =      a*X+b =      誤差
66.000    72.000    65.026    6.974
46.000    41.000    44.618   -3.618
87.000    93.000    86.455    6.545
82.000    82.000    81.353    0.647
42.000    41.000    40.536    0.464
72.000    64.000    71.149   -7.149
31.000    27.000    29.311   -2.311
69.000    68.000    68.087   -0.087
84.000    79.000    83.394   -4.394
68.000    67.000    67.067   -0.067
50.000    57.000    48.699    8.301
90.000    90.000    89.516    0.484
85.000    87.000    84.414    2.586
29.000    29.000    27.270    1.730
76.000    75.000    75.230   -0.230
69.000    66.000    68.087   -2.087
84.000    78.000    83.394   -5.394
35.000    31.000    33.393   -2.393

a =          1.02042
b =         -2.32185

r =          0.98063
決定係数 = 0.96164

```

リスト2から、リスト1のデータの場合、

$$a \quad 1.0 \quad b \quad -2.3$$

であることがわかります。

$a \quad 1.0$ ということは、英語の点数1点分の変化量に対応する国語の点数の変化量が、約1点ということになります。

さらに、 $b \quad -2.3$ ですから、国語の点は、英語の点から約2.3点減じた値が直線回帰による値となります。

すなわち、

$$\text{国語} = 1 \times \text{英語} - 2.3$$

です。

決定係数は約 0.96 になっています。このことから、英語の点数によって、国語の点数の変動量の約 96% が予測できていると解釈できます。

プログラム POneIndVar.dpr では、図 4 のフォームにおいて

	X =	Y =
1番目	66	72
2番目	46	41
3番目	87	93
4番目	82	82
5番目	42	41
6番目	72	64
7番目	31	27
8番目	69	68
9番目	84	79
10番目	68	67

図 4 データの設定されたフォーム

「保存」ボタンを押すと、CSV の形式でデータが保存されます。このときは次のダイアログボックス

が表示されるので、拡張子が.csv であるファイル名を設定します。この形式で保存したファ

イルは図 2 のフォームにおいて「読出」ボタンをクリックすると読み込むことができます。

また、Excel でも読み込めます。Excel で読み込むと次図のようになります。

	A	B	C
1	英語	国語	
2	66	72	
3	46	41	
4	87	93	
5	82	82	
6	42	41	
7	72	64	
8	31	27	
9	69	68	
10	84	79	
11	68	67	
12	50	57	
13	90	90	
14	85	87	
15	29	29	
16	76	75	
17	69	66	
18	84	78	
19	35	31	
20			
21			

上図のように設定された Excel データを拡張子が.csv である CSV 形式で保存すると、プログラム POneIndVar.dpr から「読出」ボタンのクリックで読み込むことができます。