

IV 回帰分析入門

1) 2変量データの記述

1. 散布図の描画

【課題 17】 次の表は、日本の実質家計可処分所得と実質家計最終消費支出のデータである。このデータを入力し、散布図を描いてみよう。

☞ グラフの種類を「散布図」として、データ範囲を **B2:C19** とすればよい。

また、また、軸の書式設定で、

縦軸 最小値：**160** 最大値：**300** 目盛間隔：**20**

横軸 最小値：**200** 最大値：**340** 目盛間隔：**20**

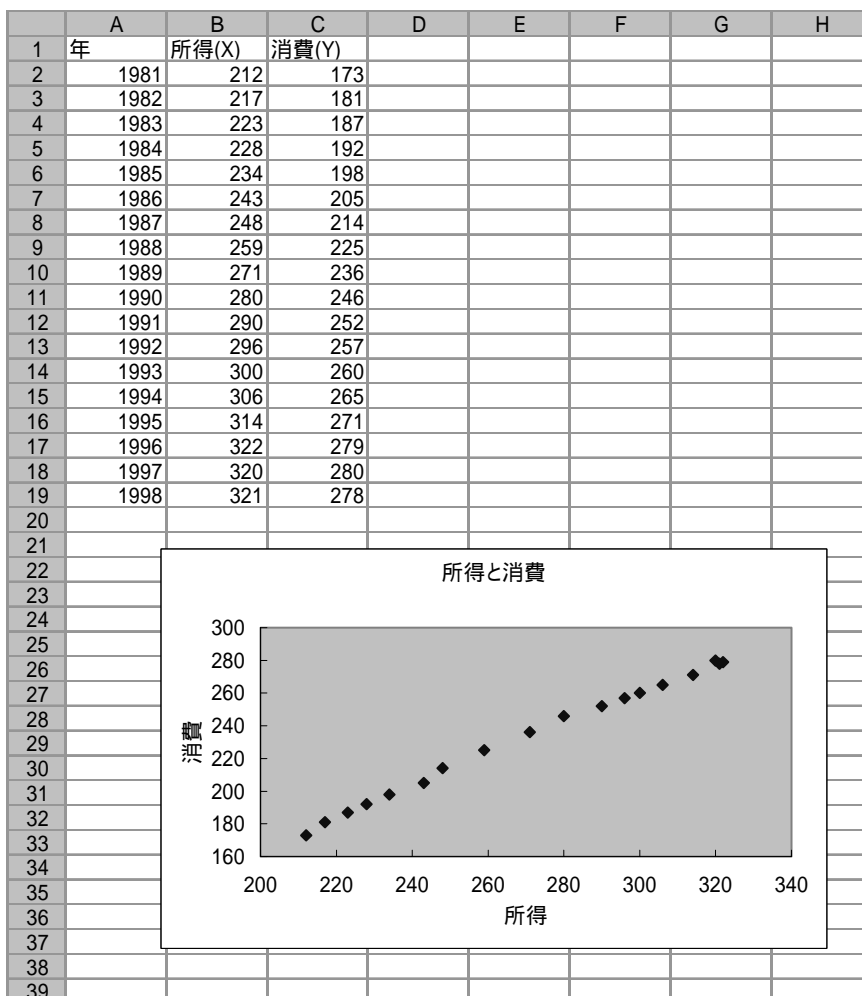
としたのがこのグラフである。

年	実質家計 可処分所得	実質家計 最終消費支出	年	実質家計 可処分所得	実質家計 最終消費支出
1981	212	173	1990	280	246
1982	217	181	1991	290	252
1983	223	187	1992	296	257
1984	228	192	1993	300	260
1985	234	198	1994	306	265
1986	243	205	1995	314	271
1987	248	214	1996	322	279
1988	259	225	1997	320	280
1989	271	236	1998	321	278

単位：兆円

出典：国民経済計算年報

<作成見本>



2. 相関係数の導出

【課題 18】 所得と消費のデータについて相関係数を求めてみよう。

相関係数の計算式は次のような式である。

$$R = \frac{n\sum XY - (\sum X)(\sum Y)}{\sqrt{\{n\sum X^2 - (\sum X)^2\}\{n\sum Y^2 - (\sum Y)^2\}}}$$

したがって相関係数を導出するためには、 X 、 Y 、 XY 、 X^2 、 Y^2 をまず求める必要がある。 $\sum X$ 、 $\sum Y$ は、**B** 列、**C** 列の和を求めればよいが、 $\sum XY$ 、 $\sum X^2$ 、 $\sum Y^2$ を求めるためには、交差積(XY)と 2 乗(X^2 、 Y^2)を **D** 列、**E** 列、**F** 列に計算したうえで、その和を求めることになる。

手順は次のようになる。

📖 手順

- ① **D** 列に X と Y の交差積を求める。**D2** セルに **=B2*D2** と入力し、これをコピーすればよい。
- ② **E** 列に X の 2 乗を、**F** 列に Y の 2 乗を求める。2 乗を表す演算子は '^' であり、**E2** セルに **=B2^2** と入力し、これをコピーする。**F** 列も同様である。
- ③ **B21** セルから **F21** セルに各列の合計を求める。これらのセルがそれぞれ $\sum X$ 、 $\sum Y$ 、 $\sum XY$ 、 $\sum X^2$ 、 $\sum Y^2$ である。
- ④ **C23** セルに **=(18 * D21 - B21 * C21)/SQRT((18 * E21 - B21^2)*(18 * F21 - C21^2))** と入力する。この式と計算式とを見比べてみよ。

< 作成見本 >

	A	B	C	D	E	F
1	年	所得(X)	消費(Y)	XY	X^2	Y^2
2	1981	212	173	36676	44944	29929
3	1982	217	181	39277	47089	32761
4	1983	223	187	41701	49729	34969
5	1984	228	192	43776	51984	36864
6	1985	234	198	46332	54756	39204
7	1986	243	205	49815	59049	42025
8	1987	248	214	53072	61504	45796
9	1988	259	225	58275	67081	50625
10	1989	271	236	63956	73441	55696
11	1990	280	246	68880	78400	60516
12	1991	290	252	73080	84100	63504
13	1992	296	257	76072	87616	66049
14	1993	300	260	78000	90000	67600
15	1994	306	265	81090	93636	70225
16	1995	314	271	85094	98596	73441
17	1996	322	279	89838	103684	77841
18	1997	320	280	89600	102400	78400
19	1998	321	278	89238	103041	77284
20						
21	合計	4884	4199	1163772	1351050	1002729
22						
23		相関係数	0.998068			
24						

2) 単回帰モデル(その 1)

1. 回帰直線の導出

【課題 19】 所得と消費のデータについて $Y = \alpha + \beta X$ という 1 次式をあてはめ、回帰係数 α 、 β の推定値を求めよ。

回帰係数の推定値を求める式は次のようなものである。

$$b = \frac{n \sum XY - \sum X \sum Y}{n \sum X^2 - (\sum X)^2}$$
$$a = \frac{\sum X^2 \sum Y - \sum X \sum XY}{n \sum X^2 - (\sum X)^2}$$

この式に相関係数の導出の際に求めた $\sum X$ 、 $\sum Y$ 、 $\sum XY$ 、 $\sum X^2$ 、 $\sum Y^2$ を代入すれば回帰係数の推定値がそれぞれ計算できる。

2. 予測値と残差の計算

【課題 20】 所得と消費のデータについて各年の X のデータに対する予測値 \hat{Y} と残差を求めよ。

📖 手順

- ① G 列に予測値を求める。予測値 \hat{Y} は各 X_i について $a + b X_i$ を計算すればよいので、G2 セルに 1981 年の X(B2 セル)に対応する予測値を求めるなら $=C\$25+\$C\$24*B2$ とし、これをコピーすればよい。ここでは、コピーの際に絶対参照をするので、'S'がついている。
- ② H 列に残差を求める。残差は Y から予測値 \hat{Y} を引いたものなので、H2 セルに $=C2-G2$ とし、これをコピーすればよい。

3. 回帰直線のグラフへの書き入れ

散布図に回帰直線を書き入れる場合、Excel では各 X に対応する予測値をグラフに書き入れ、それを直線でつなぐという手順をとる。

【課題 21】 所得と消費のデータについて散布図に回帰直線を書き入れよ。

📖 手順

- ① **グラフをアクティブ** (グラフエリアの枠の四隅および 4 辺に ■ という印が現れて入る状態。散布図のグラフエリアの白い部分をクリックすればこの状態になる。) にし、メニューバーから「グラフ」-「データの追加」を選ぶ。すると、「グラフの追加」というウィンドウが開くので、G2 セルから G19 セルまでをドラッグし、OK ボタンを押す。
- ② ①の操作で散布図上にピンク色の点が見れたはずである。これを直線で結ぶ。ピンク色の点のひとつをダブルクリックすると、「データ系列の書式設定」ウィンドウが開く。¹そこで「パターン」のタグにおいて、線を「指定」をチェックし、色を黒に変え、マーカーを「なし」にする。グラフエリアの外をクリックすると回帰直線が引けたことがわかるはずである。

¹ この操作はうまくいかないことが多々ある。その場合にはグラフエリアの外を一回クリックした後で、もう一度この操作を繰り返かえてみると良い。

4. 決定係数の導出

決定係数は回帰における当てはまりの尺度であり、全変動のうち回帰モデルによって説明される変動の割合を示すものである。決定係数は0と1の間の値をとるが、決定係数が0.2や0.3などの小さい値であるということは、あまり関係のないXとYの間に因果関係を想定し、分析を行っているということを意味し、モデルの再検討が必要となる。

決定係数は、相関係数との間に、

$$(\text{決定係数 } R^2) = (\text{相関係数})^2$$

という関係がある。

5. 残差の表示

残差 e_i は従属変数の個々の観測データと回帰直線との間のズレの大きさをあらわすものであった。この残差の状態を調べることで、回帰直線の当てはまり具合など、さまざまな情報入手することができる。残差を出発点としてモデルの設定やデータ間の関係を検討する分析を、残差分析(**residual analysis**)という。ここでは、残差を求めてそれをグラフに表示してみよう。

グラフを描くには、残差 e_i を縦軸にとり、横軸には、

① i

② X_i

③ \hat{Y}_i

などを用いる。

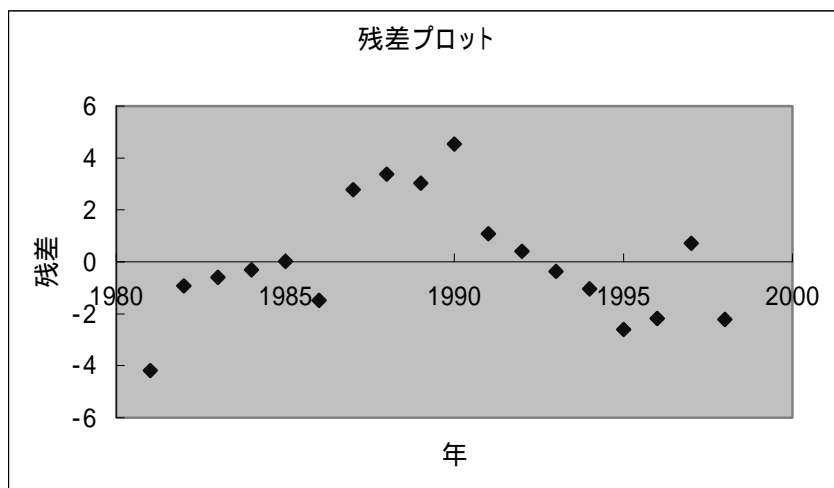
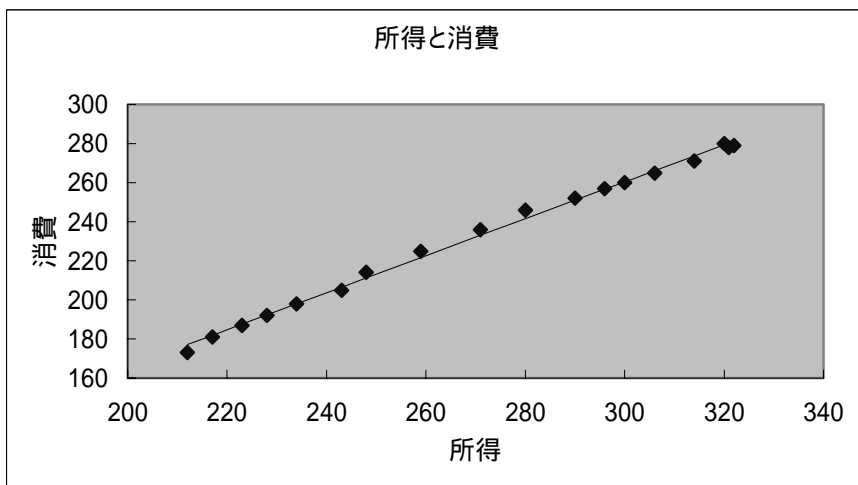
☒ **演習問題 3** : 消費と所得のデータについて、【課題 21】までに加え、決定係数を計算し、残差プロット(横軸は年とする)を描いたものを、A4用紙縦1枚に納まるようにレイアウトして、印刷してみよ。
次ページに見本を示す。

< 作成見本 >

わが国の実質家計可処分所得と実質家計消費支出 E35-000 徳山太郎

年	所得(X)	消費(Y)	XY	X^2	Y^2	予測値	残差
1981	212	173	36676	44944	29929	177.1905	-4.19051
1982	217	181	39277	47089	32761	181.917	-0.91697
1983	223	187	41701	49729	34969	187.5887	-0.58872
1984	228	192	43776	51984	36864	192.3152	-0.31517
1985	234	198	46332	54756	39204	197.9869	0.013084
1986	243	205	49815	59049	42025	206.4945	-1.49453
1987	248	214	53072	61504	45796	211.221	2.779011
1988	259	225	58275	67081	50625	221.6192	3.380811
1989	271	236	63956	73441	55696	232.9627	3.037319
1990	280	246	68880	78400	60516	241.4703	4.529701
1991	290	252	73080	84100	63504	250.9232	1.076791
1992	296	257	76072	87616	66049	256.595	0.405046
1993	300	260	78000	90000	67600	260.3761	-0.37612
1994	306	265	81090	93636	70225	266.0479	-1.04786
1995	314	271	85094	98596	73441	273.6102	-2.61019
1996	322	279	89838	103684	77841	281.1725	-2.17252
1997	320	280	89600	102400	78400	279.2819	0.718063
1998	321	278	89238	103041	77284	280.2272	-2.22723
合計	4884	4199	1163772	1351050	1002729		

相関係数 0.9980678
 b 0.9452909
 a -23.21117
 決定係数 0.9961394



3) 単回帰モデル(その2)

前章では、相関係数と回帰係数の推定値を交差積和($\sum XY$)、2乗和($\sum X$, $\sum Y$)を求め、それを計算式に代入することによって求めた。

しかしExcelによって相関係数と回帰直線を求めるには、以下に説明するようなExcelが備えている関数を用いることもできる。ここではregの例について、統計関数を用いた方法についても行ってみることにする。その際には、データを入力したセルの範囲に名前をつけておくと便利である。まず、実習の準備として、regのSheet1から、年次、所得X、消費Yの部分(A3:C21)をSheet2のA3:C21にコピーし、所得のデータに_X、消費のデータに_Yという名前を定義しておく。

1. 統計関数による相関係数と回帰直線の導出

Excelが備えている関数を用いた相関係数と回帰直線の導出を行ってみることにする。regのSheet2のA3:C21に、年次、所得X、消費Yのデータが入力されているものとする。

(1) 関数 PEARSON (CORREL), RSQ

相関係数を求めるには、関数 **PEARSON**(引数1, 引数2) を用いる。**PEARSON** は相関係数を最初に導出した **Karl Pearson (イギリス;1851-1936)** にちなんでつけられた名前である。または、**CORREL** という名前の関数もあるが、どちらも全く同じものである。引数は2個あり、それぞれがデータの範囲(名前でもよい)である。戻り値は r_{xy} である。

= PEARSON(B4:B21, C4:C21)
= PEARSON(_X , _Y)

範囲B4:B21に名前_X、範囲C4:C21に名前_Yを付けてあれば、どちらの式でも結果は同じである。以下の説明では下式の書き方で示す。

関数 **RSQ**は相関係数の2乗(=**決定係数**)を求める関数であるが、引数は **PEARSON** と同じである。したがって、べき乗を求める演算子 ^ を用いれば **RSQ** は不要となる。

= RSQ(_X, _Y)
= (PEARSON(_X, _Y)) ^ 2

どちらも全く同じ結果を与える。

(2) 関数 SLOPE とINTERCEPT²

SLOPE は回帰直線の傾き(回帰係数) **b** を、**INTERCEPT** は切片(回帰定数) **a** を求める関数で、どちらも引数は2個あるが、最初の引数が従属変数の範囲で、2個目の引数が独立変数の範囲をとる。引数の順序に注意しなければならない。

= SLOPE(_Y, _X)
= INTERCEPT(_Y, _X)

² 回帰直線の傾きと切片を求める関数には、**LINEST** という関数がある。この関数は傾きと切片以外に分析結果に関する多くの情報量を与えてくれる、非常に便利な関数である反面、使用法および結果の解釈の仕方が難しい。LINEST 関数の説明はここでは省略する。

(3) 関数 FORECAST と TREND

予測値 \hat{Y} を求める関数には2種類のものが用意されている。関数**FORECAST** は引数を3個とり、**FORECAST(X_i , Y範囲, X範囲)** として用いる。戻り値は $a + bx$ として求められた数値1個である。

=FORECAST(B4, _Y, _X) セル**B4**の値を **x** としたときの $a + bx$ が戻り値

残りの **X** の値に対する予測値は、これをコピーして求めればよい。

あるいは、n個の予測値を書き込む範囲を指定しておき、配列数式とすることもできる。たとえば、**D4:D21**の範囲を指定して、

=FORECAST(_X, _Y, _X)

を入力して、**Ctrl** + **Shift** + **Enter** とする。

関数 **TREND** も同じ予測値を求めるものであるが、引数の数が **FORECAST** より1個多く、計4個となる。一般的な型式は **TREND(Y範囲, X範囲, X_i, 1)**となる。最後の引数は **0** か **1** で、**0** のときは、原点を通る直線 $Y = bX$ による予測値、**1** のときはこれまで通りの $Y = a + bX$ による予測値を戻り値として求める。第4引数を省略した場合は、**1** を指定したものとみなす。

TRENDの第3引数として、**a, b** の計算に用いない任意の数値を指定することもできる。たとえば、

=TREND(_Y, _X, 190, 1)

とすれば、 $\hat{Y} = a + b * 190$ を求めることになる。**reg**の **X** のデータの中には **190**という数値はなく、これによって求まる \hat{Y} は未知の **X** の値に対する**予測値** (外挿値) である。同様のことを、**FORECAST**を用いてもおこなうことができる。

=FORECAST(190, _Y, _X)

FORECASTでは **190** が第1引数となる点に注意されたい。また、**190** という数値を直接指定するのではなく、セル番地で指定することもできる。セル **N4** に**190**が書き込んであれば、

=TREND(_Y, _X, N4, 1)

=FORECAST(N4, _Y, _X)

とすればよい。また、**N4**から**N13**に $\hat{Y} = a + bx$ として求めたい **x** の値が連続して書き込まれていれば、

=TREND(_Y, _X, N4:N13, 1)

=FORECAST(N4:N13, _Y, _X)

とすればよい。範囲 **N4:N13** に名前を付けて、それを使用してもよい。

3. 分析ツールの利用

Excel には統計分析を行うためのいくつかの分析ツールが付属している。これらのツールを使えば一度に詳細な分析結果を得ることができる。

分析ツールを最初に使用する場合にはメニューバーの「ツール」－「アドイン」を選び、分析ツールをチェックすることによって、分析ツールをアドイン(有効にすること)しなくてはならない。

アドインを行った後で、再びメニューバーから「ツール」を選ぶと、下のほうに「分析ツール」と表示される。ここで分析ツールを選び、回帰分析を選べばよい。