データ分析の基礎-3

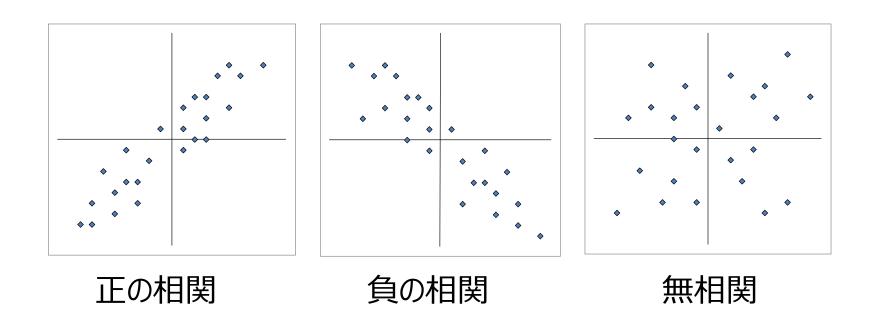
2022年10月20日



本日の内容

- ◇相関関係 散布図 積率相関係数 偏相関係数
- ◇回帰分析 回帰式 偏回帰係数、t検定、決定係数

相関関係



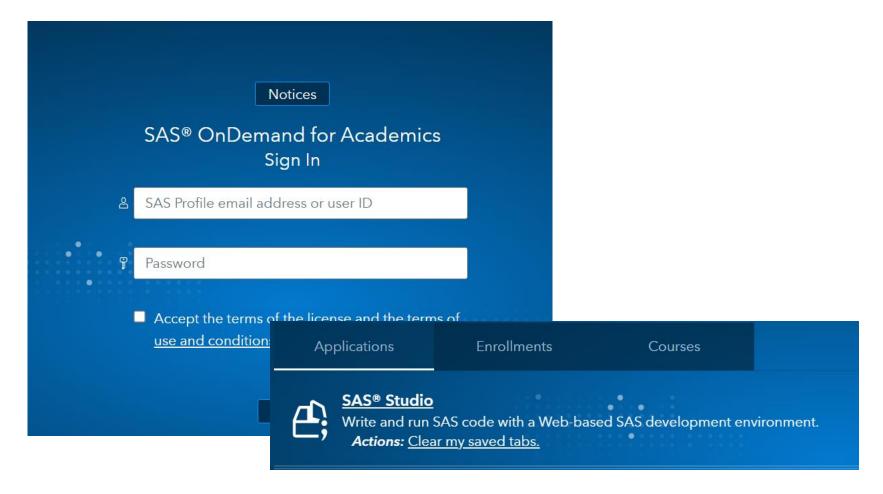
相関関係には正の相関、負の相関、無相関。点の集中度が関係の強さを測定する手がかり。

◇支店別広告費と売上高

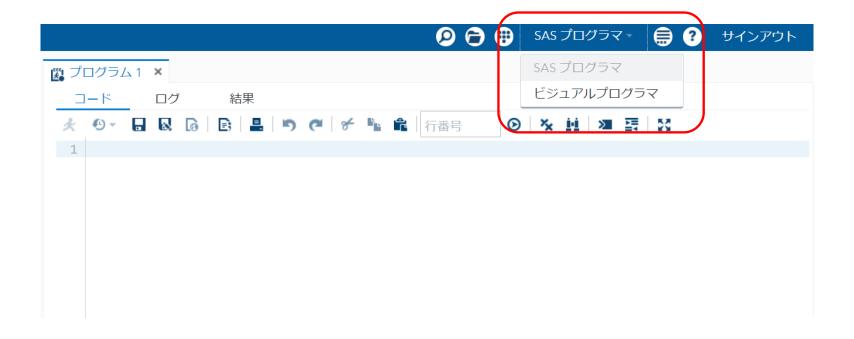
支店	広告費	売上高
北海道	92	44
東北	93	102
関東	332	288
北陸	78	54
中部	181	118
近畿	108	138
中国	113	138
四国	72	86
九州	243	152
沖縄	13	22

散布図の作成法 (SAS Studio)

- 1.SAS Studio にログインする。
- 2.SAS®Studio をクリックする。



3. 「SASプログラマ」をクリックし、「ビジュアルプログラマ」を選択する。



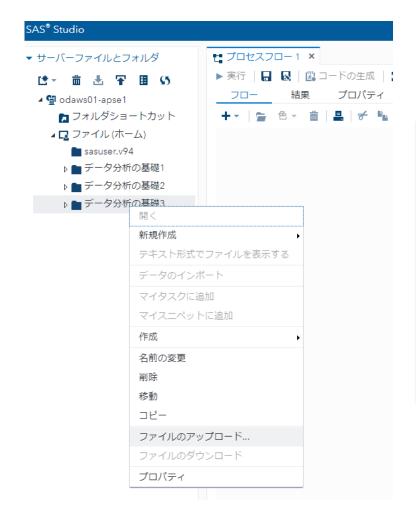
4. 新規フォルダを作成する。

「ファイル(ホーム)」を右クリックー「新規作成」ー「フォルダ」をクリック、 「新規フォルダ名(データ分析の基礎3)」を入力し、「保存」をクリックする。



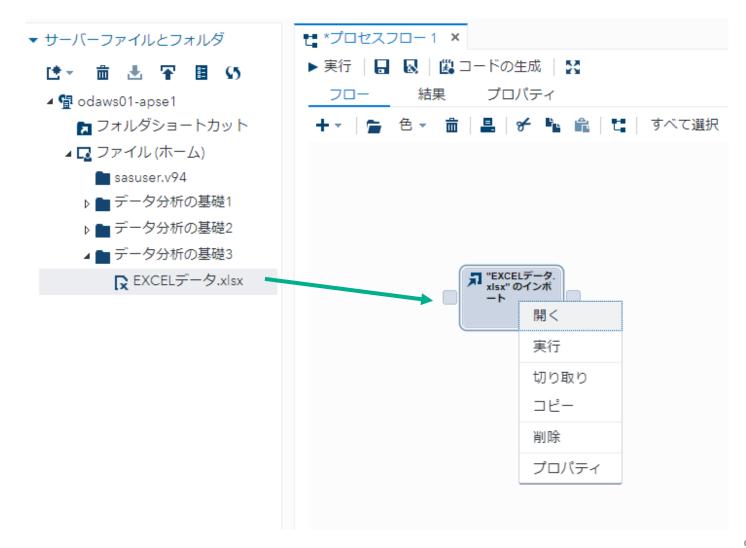
5. ファイルをアップロードする。

「データ分析の基礎3」を右クリック、「ファイルのアップロード」をクリックし、「ファイルの選択」から、ファイルを選択し「アップロード」をクリックする。





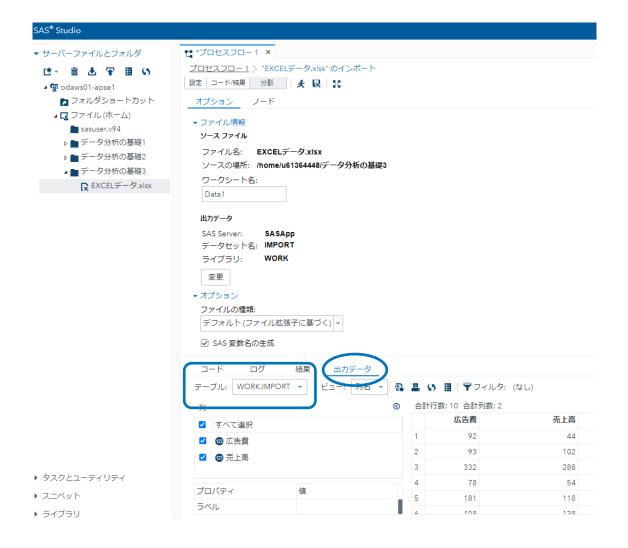
6. 「データ分析の基礎3」フォルダを開きファイル(EXCELデータ)を 右側のプロセスフロー画面にドラッグし、右クリックして「開く」を選択する。



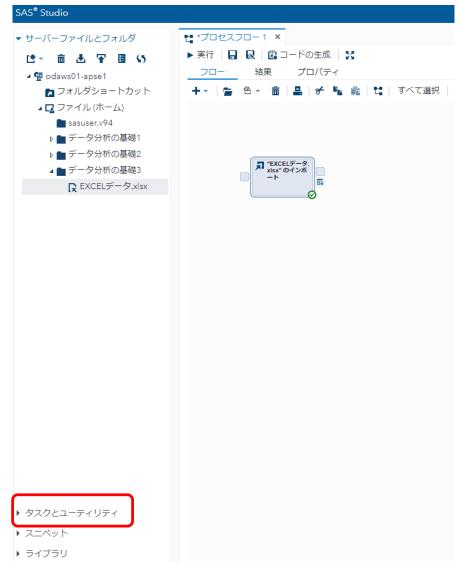
7. 「ワークシート名(Data1)」を入力し、実行ボタンをクリックする。



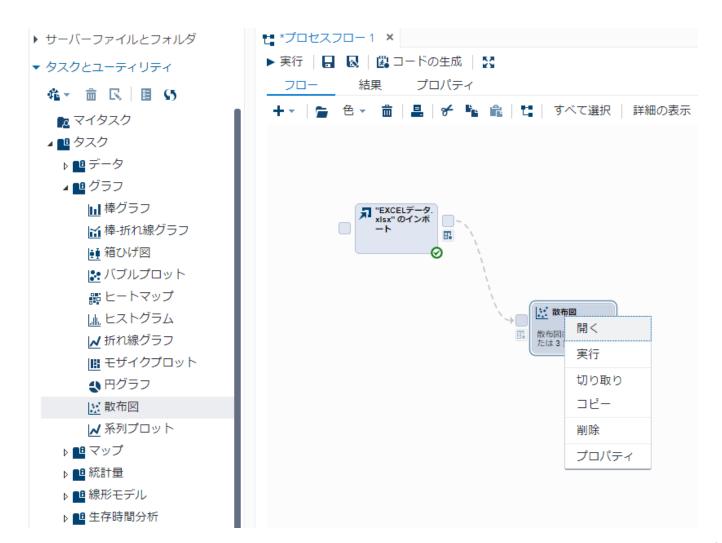
8. 「出力データ」をクリックし、「テーブル(データセット名)」、データの内容を確認する。



9. 「プロセスフロー」をクリックしてフロー画面に戻り、「タスクとユーティリティ」を開く。



10. 「タスクとユーティリティ」→「タスク」→「グラフ」の「散布図」をフロー画面に ドラッグし、「EXCELデータ・・」と結合、右クリック-「開く」をクリックする。



11. 「散布図」を右クリック、開き、「データ」、「X軸(広告費)」、「Y軸(売上高)」をセットする。



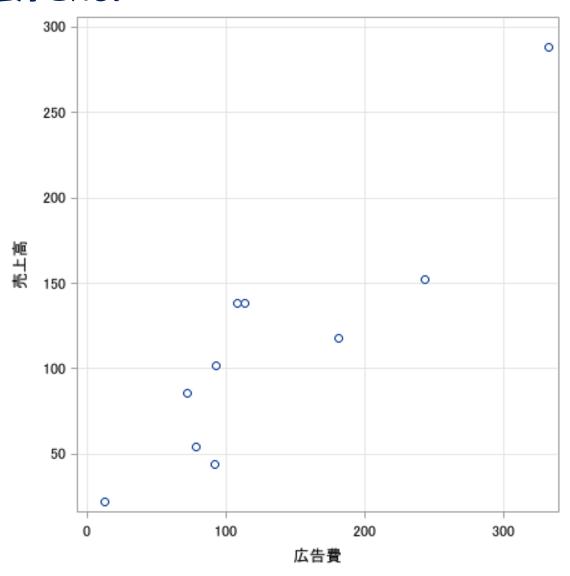
12. 「表示」をクリックし、「グラフサイズ」を幅「4.8」に変更し、出力する グラフを正方形にする。



13. 実行ボタンをクリックすると散布図が表示される。

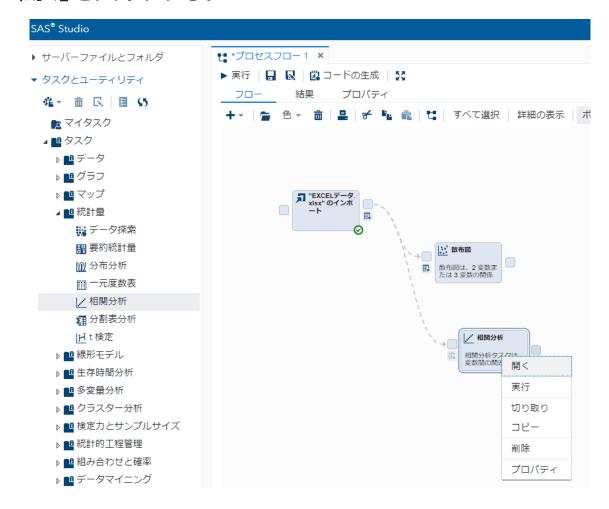


結果が表示される。



積率相関係数を求める(SAS Studio)

 「タスクとユーティリティ」→「タスク」→「統計量」の「相関分析」を フロー画面にドラッグし、「EXCELデータ・・」と結合、右クリック-「開く」をクリックする。



2. 「相関分析」を右クリック、開き、「データ」、「分析変数(広告費、売上高)」をセットする。



3. 実行ボタンをクリックする。



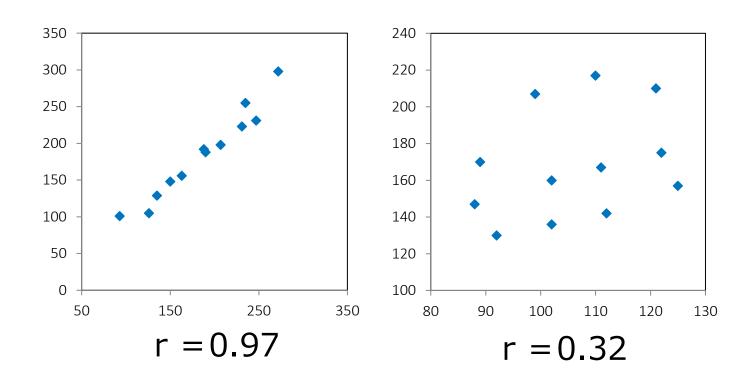
結果が表示される。



積率相関係数 (r) = 0.902

積率相関係数(r)

$-1 \leq r \leq 1$



積率相関係数 (r) は相関関係の強さ

積率相関係数(r)の解釈

0.32×0.32=0.1024 説明力は、10.24%

◇支店別広告費、売上高、人口

支店	広告費	売上高	人口
北海道	92	44	5506
東北	93	102	9335
関東	332	288	42604
北陸	78	54	5443
中部	181	118	18127
近畿	108	138	12912
中国	113	138	15554
四国	72	86	3976
九州	243	152	13204
沖縄	13	22	1393

3 変数: 広告費 売上高 人口

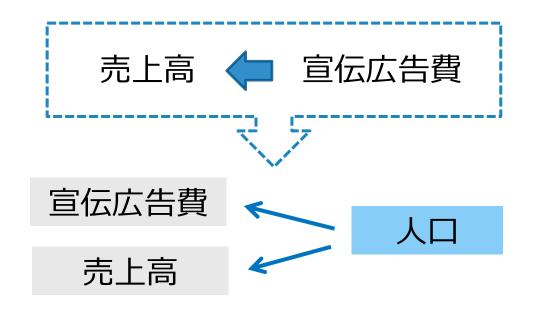
Pears	son の相関	関係数. N	= 10
	広告費	売上高	人口
広告費 広告費	1.00000	0.90236	0.89471
売上高 売上高	0.90236	1.00000	0.95093
煰	0.89471	0.95093	1,00000

広告費と売上高 人口と広告費 人口と売上高 r = 0.902

r = 0.895

r = 0.951

広告宣伝費は売上高に貢献?



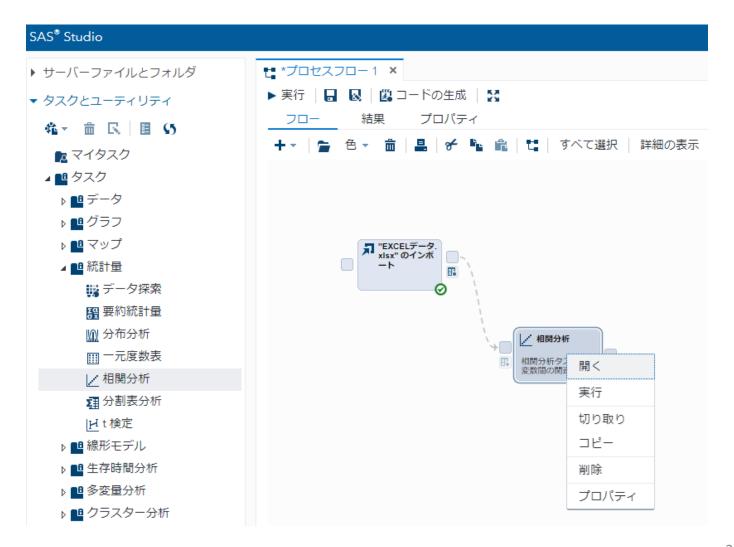
人口の影響を除いたときの広告費と売上高の 積率相関係数 → 偏相関係数

偏相関係数(SAS Studio)

1.「ワークシート名(Data3)」を入力し、実行ボタンをクリックする。



2. 「相関分析」を右クリック-「開く」をクリックする。



3. 「相関分析」を右クリック、開き、「データ」、「分析変数(広告費、売上高)」、「部分変数(人口)」をセットする。



4. 実行ボタンをクリックする。



結果が表示される。



偏相関係数=0.373

◇事例

都道府県別コンビニ件数、人口、甲子園の勝率

1		コンビニ数	人口	甲子園勝率
2	東京都	6847	13513.7	0.543
3	神奈川県	3431	9127.3	0.623
4	大阪府	3654	8838.9	0.616
5	愛知県	3576	7484.1	0.602
6	埼玉県	2597	7261.3	0.538
7	千葉県	2421	6224.0	0.554
8	兵庫県	1852	5537.0	0.563
9	北海道	1807	5383.6	0.335
10	福岡県	2040	5102.9	0.475
11	静岡県	1685	3701.2	0.493

Pearson の相関係数, N = 47			
	コンビニ数	人口	甲子園勝率
コンビニ数 コンビニ数	1.00000	0.98272	0.40987
人口	0.98272	1.00000	0.43197
甲子園勝率 甲子園勝率	0.40987	0.43197	1.00000

・人口とコンビニ数: 0.983

人口の多い都道府県はコンビニ件数が多い。

・人口と甲子園勝率:0.432

人口が多いと高校生の数も多く、結果として野球のレベルも上がる。

·コンビニ数と甲子園勝率: 0.410?



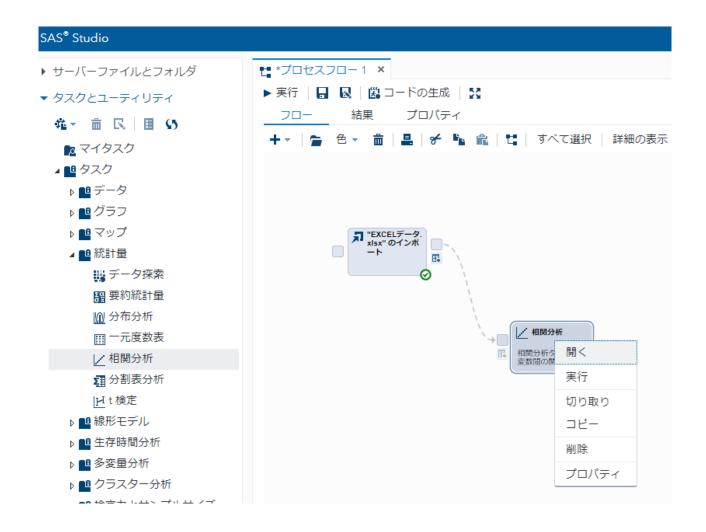
→ 人口の影響を除いた積率相関係数(偏相関係数)

偏相関係数を求める(SAS Studio)

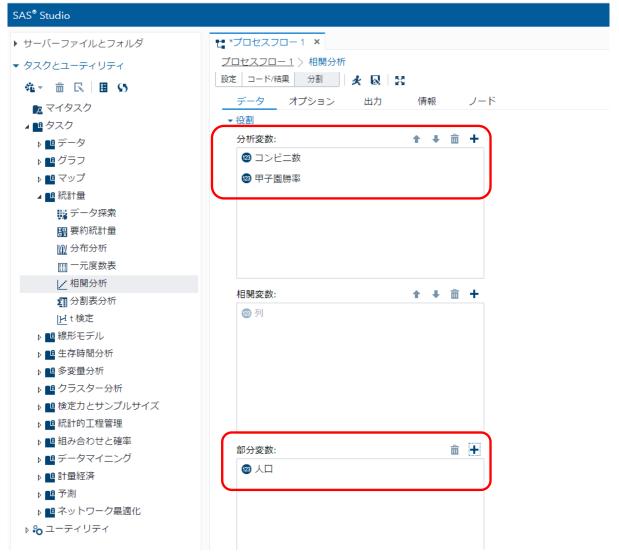
1.「ワークシート名(Data5)」を入力し、実行ボタンをクリックする。



2. 「相関分析」を右クリック-「開く」をクリックする。



3. 「相関分析」を右クリック、開き、「データ」、「分析変数(コンビニ数、甲子園勝率)」、「部分変数(人口)」をセットする。



結果が表示される。

1 Partial 変数	女:	人口		
2 変数:		コンビコ	— * — #	女 甲子園勝率
Pearson	の偏	相関係	赦	, N = 47
	⊐ 2	ノビニ数		甲子園勝率
コンビニ数		1.00000		-0.08766
コンビニ数				
甲子園勝率	-	0.08766		1.00000
甲子園勝率			_	

偏相関係数 = -0.088

回帰分析

例)売上高と売上高に影響を与える要因との関係

売上高←広告宣伝費、人口、セールスマン数、・・・

従属変数←説明変数(独立変数)

説明変数が1つ:単回帰

説明変数が2つ以上:重回帰

目的1:回帰式を求め、予測する。

回帰式(y=a+bx+...)を求める。



a,b,••• 📥 偏回帰係数

◇駅前コンビニの売上高と乗降客数

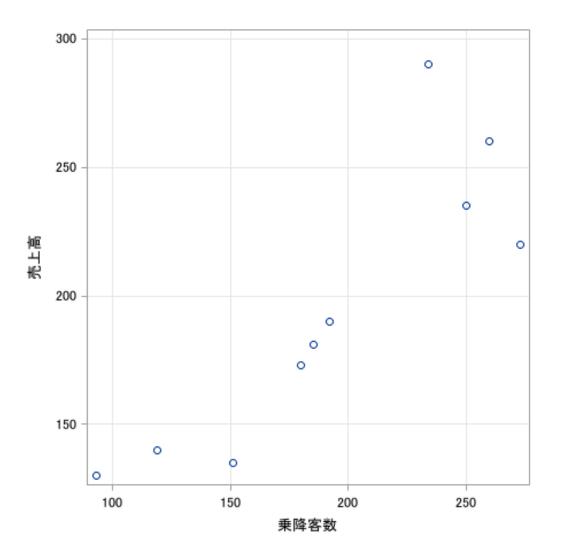
売上高←乗降客数

	売上高	乗降客数
	(百万円/月)	(百人/日)
1	130	93
2	290	234
3	235	250
4	260	260
5	140	119
6	173	180
7	135	151
8	190	192
9	220	273
10	181	185

2 変数: 売上高 乗降容数

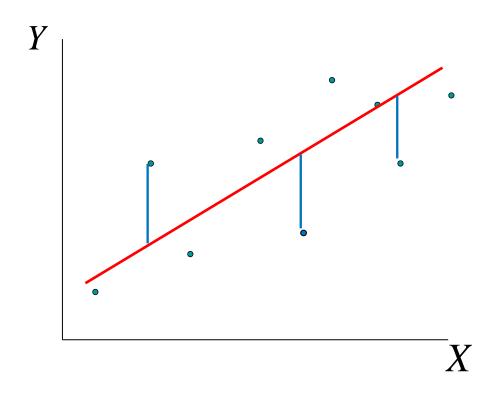
Pearson の相関係数, N = 10					
売上高 乗降客数					
売上高 売上高	1.00000	0.86747			
乗降容数 乗降容数	0.86747	1,00000			

積率相関係数 = 0.867



回帰式(売上高 = a + b×乗降客数)を求めて予測する。

回帰直線の求め方(最小二乗法)



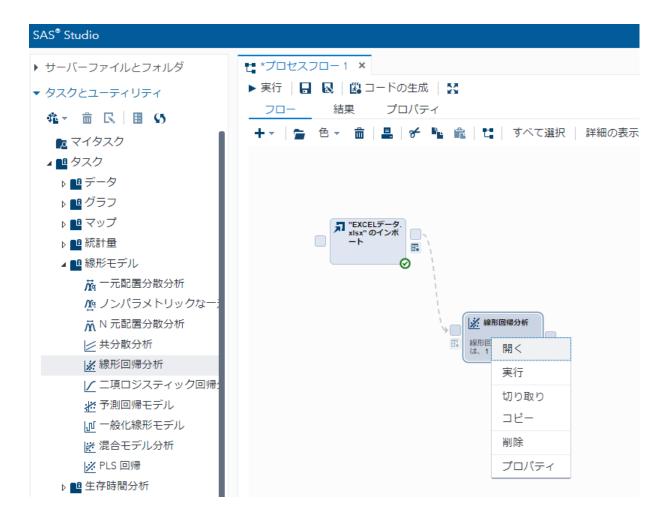
各データと回帰直線との垂直距離の2乗和を最小にする

回帰分析(SAS Studio)

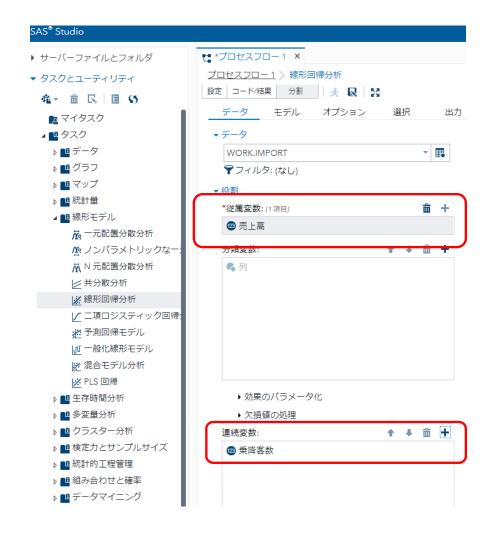
1. 「ワークシート名 (Data 10) 」を入力し、実行ボタンをクリックする。



2. 「タスクとユーティリティ」→「タスク」→「線形モデル」の「線形回帰分析」 をフロー画面にドラッグし、「EXCELデータ・・」と結合、右クリックー「開く」 をクリックする。



3. 「線形回帰分析」を右クリック、開き、「データ」、「従属変数(売上高)」、「連続変数(乗降客数)」をセットする。



4. 「モデル」をクリックし、「モデル効果」の「編集」をクリックする。



5.「変数(乗降客数)」→「単一効果(追加)」をセットする。 「切片」にチェックがつき、「乗降客数」が表示されていることを確認し、 OKボタンをクリックする。



7. 実行ボタンをクリックする。



結果が表示される。

Root MSE	28.91642	R2 乗	0.7525
従属変数の平均	195.40000	調整済み R2 乗	0.7216
変動係数	14.79858		

パラメータの 推定						
変数	ラベル	自由度	バラメータ 推定値	標準誤差	t 値	Pr > t
Intercept	Intercept	1	43.99116	32.03230	1.37	0.2069
乗降客数	乗降客数	1	0.78167	0.15849	4.93	0.0011

売上高=43.99+0.782×乗降客数

*乗降客数=180のときの売上高の予測 売上高=43.99+0.782×180=184.75

回帰モデルのチェック

◇自由度調整済み決定係数

 $R^2 = 0.7216$

約72.16%説明できる。

◇偏回帰係数のt検定

乗降客数:P値(P値) = 0.0011

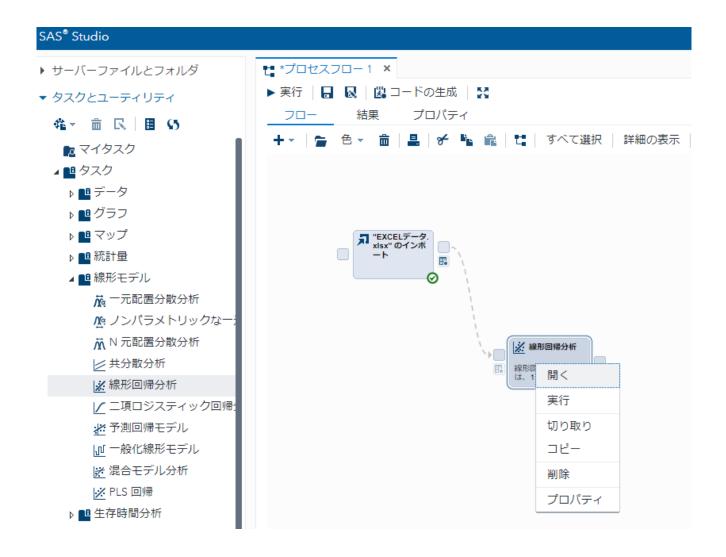
説明変数に「取扱品目数」を追加

	売上高	乗降客数	取扱品目数
	(百万円/月)	(百人/日)	(品)
_1	130	93	150
2	290	234	311
3	235	250	182
4	260	260	245
5	140	119	149
6	173	180	160
7	135	151	98
8	190	192	180
9	220	273	113
10	181	185	105

1.「ワークシート名(Data11)」を入力し、実行ボタンをクリックする。



2. 「線形回帰分析」を右クリック-「開く」をクリックする。



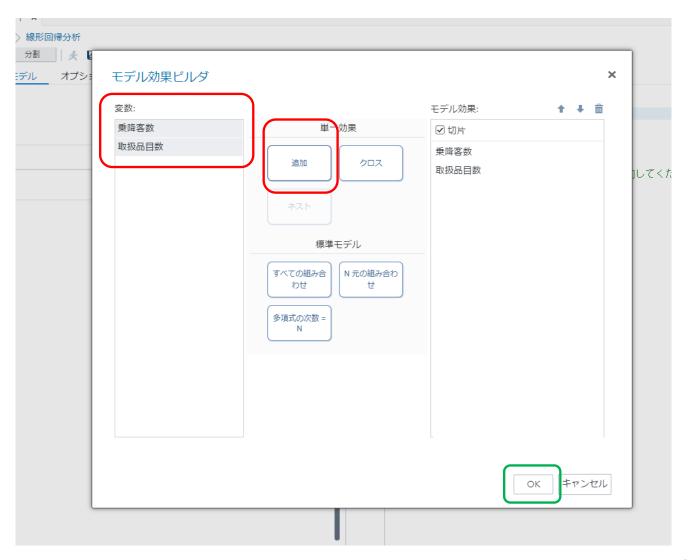
3. 「線形回帰分析」を右クリック、開き、「データ」、「従属変数(売上高)」、「連続変数(乗降客数、取扱品目数)」をセットする。



4. 「モデル」をクリックし、「モデル効果」の「編集」をクリックする。



5.「変数(乗降客数、取扱品目数)」→「単一効果(追加)」をセットし、 OKボタンをクリックする。



6. 実行ボタンをクリックする。



結果が表示される。

Root MSE	10.04926	R2 乗	0.9738
従属変数の平均	195.40000	調整済み R2 乗	0.9664
変動係数	5.14292		

	バラメータの 推定						
変数	ラベル	自由度	バラメータ 推定値	標準誤差	t 値	Pr > t	
Intercept	Intercept	1	7.15385	12.11739	0.59	0.5735	
乗降客数	乗降客数	1	0.60153	0.05985	10.05	<.0001	
取扱品目数	取扱品目数	1	0.42368	0.05505	7.70	0.0001	

◇偏回帰係数のt検定

乗降客数 P値 = < 0.0001

取扱品目数 P値 = 0.0001

売上高の予測

売上高=7.15+0.602×乗降客数+0.424×品目数

*乗降客数=200 取扱品目数が180のとき売上高の予測

売上高=7.15+0.602×200+0.424×180=203.87

◇自由度調整済決定係数

 $R^2 = 0.9664$ 約96.64%説明できる。

説明変数に世帯数をさらに追加

	売上高 (百万円/月)	乗降客数 (千人/日)	取扱品目数 (品)	世帯数
1	130	93.0	150	143
2	290	234.0	311	284
3	235	250.0	182	320
4	260	260.0	245	302
5	140	119.0	149	182
6	173	180.0	160	225
7	135	151.0	98	190
8	190	192.0	180	242
9	220	273.0	113	320
10	181	185.0	105	235

「ワークシート名(Data12)」を入力し、実行ボタンをクリックする。



重回帰分析 結果

Root MSE	10.55938	R2 乗	0.9748
従属変数の平均	195.40000	調整済み R2 乗	0.9622
変動係数	5,45516		

バラメータの推定						
	パラメータ					
変数	ラベル	自由度	推定値	標準誤差	t 値	Pr > t
Intercept	Intercept	1	- 1 82664	23,00272	-0.08	0.9393
乗降客数	乗降容数	1	0.42311	0.38429	1.10	0.3131
取扱品目数	取扱品目数	1	0.421 44	0.05858	7.19	0.0004
世帯数	世帯數	1	0.17978	0.38190	0.47	0.6544

◇偏回帰係数のt検定

乗降客数 P値 = 0.3131

取扱品目数 P値 = 0.0004

世帯数 P値 = 0.6544



売上高を乗降客数、取扱品目数、世帯数の3つの説明変数で説明する重回帰モデルは、不成立!

説明変数の検討

- ①乗降客数&取扱品目数
- ②乗降客数&世帯数
- ③取扱品目数&世帯数
- ①説明変数:乗降客数&取扱品目数

Root MSE	10.04926	R2 乗	0.9738
従属変数の平均	195,40000	調整済み R2 乗	0.9664
変動係数	5.14292		

バラメータの推定							
変数	ラベル	自由度	パラメータ 推定値	標準誤差	t値	Pr > [t]	
Intercept	Intercept	1	7.15385	12.11739	0.59	0.5735	
乗降客数	乗降容数	1	0.60153	0.05985	10.05	<.0001	
取扱品目数	取扱品目数	1	0.42368	0.05505	7.70	0.0001	

乗降客数 取扱品目数 P値 = <0.0001

P値 = 0.0001

②説明変数:乗降客数&世帯数

Root MSE	30.61661	R2 乗	0.7572
従属変数の平均	195,40000	調整済み R2 乗	0.6879
変動係数	15,66869		

バラメータの推定							
変数	ラベル	Pr > t					
Intercept	Intercept	1	23,40110	65.29769	0.36	0.7306	
乗降客数	乗降客数	1	0.37915	1.10365	0.34	0.7413	
世帯数	世帯数	1	0.40343	1.09329	0.37	0.7230	

乗降客数 P値 = 0.7413

世帯数 P値 = 0.7230

③説明変数:取扱品目数&世帯数

Root MSE	10.81976	R2 乗	0.9697
従属変数の平均	195,40000	調整済み R2 乗	0.9610
変動係数	5.53723		

バラメータの推定							
変数	パラメータ ラベル 自由度 推定値 標準誤差 t 値 Pr						
Intercept	Intercept	1	-21,00871	15.24568	-1.38	0.21 06	
世帯数	世帯数	1	0.59449	0.06403	9.28	<.0001	
取扱品目数	取扱品目数	1	0.42041	0.05946	7.07	0.0002	

世帯数 P値 = < 0.0001

取扱品目数 P値 = 0.0002

説明変数	t 検定
乗降客数、取扱品目数、世帯数	×
①乗降客数、取扱品目数	0
②乗降客数、世帯数	×
③取扱品目数、世帯数	0

積率相関係数

Pearson の相関係数, N = 10								
売上高 乗降客数 取扱品目数 世帯								
売上高 売上高	1.00000	0.86747	0.77224	0.86784				
乗降客数 乗降客数	0.86747	1,00000	0.391 08	0.98837				
取扱品目数 取扱品目数	0.77224	0.391 08	1,00000	0.39792				
世帯数 世帯数	0.86784	0.98837	0.39792	1.00000				



いずれも高い値⇒売上高を説明する説明変数として妥当

乗降客数 —— 取扱品目数 0.391

乗降客数 ——— 世帯数 0.988

取扱品目数 —— 世帯数 0.398

乗降客数と世帯数の値0.988は高い値

説明変数相互の積率相関係数は低い方が良い! 説明変数⇒独立変数

t検定におけるサンプルサイズの影響

売上高	乗降客数	間口の広さ
130	93	150
290	234	148
235	250	182
260	260	245
140	119	149
173	180	160
135	151	135
190	192	180
220	273	113
181	185	105

データ10組のとき

「ワークシート名(Data15)」を入力し、実行ボタンをクリックする。



Root MSE	29.46589	R2 乗	0.7751
従属変数の平均	195.40000	調整済み R2 乗	0.71 09
変動係数	15,07978		

バラメータの推定							
変数パラメータ変数ラベル自由度推定値標準誤差t 値Pr 3						Pr > t	
Intercept	Intercept	1	17.72278	45.22169	0.39	0.7068	
乗降客数	乗降客数	1	0.74276	0.16802	4.42	0.0031	
間口の広さ	間口の広さ	1	0.21573	0.25704	0.84	0.4290	

間口の広さ: P値 = 0.4290

データ50組のとき

「ワークシート名(Data16)」を入力し、実行ボタンをクリックする。



Root MSE	25.42756	R2 乗	0.7751
従属変数の平均	195.40000	調整済み R2 乗	0.7656
変動係数	13.01308		

	バラメータの推定								
変数	ラベル	自由度	バラメータ 推定値	標準誤差	t値	Pr > [t]			
Intercept	Intercept	1	17.72278	17.45207	1.02	0.3151			
乗降客数	乗降客数	1	0.74276	0.06484	11.45	<.0001			
間口の広さ	間口の広さ	1	0.21573	0.09920	2.17	0.0347			

間口の広さ: P値 = 0.0347

t 検定のP値はサンプルサイズの影響を受ける。

◇満足度調査

男性19名 女性21名

<従属変数> 満足度

- <説明変数>
- •機能
- ・デザイン
- ·性別 (男性1 女性0)

満足度	機能	デザイン	性別	満足度	機能	デザイン	性別
5	5	4	1	5	2	5	0
4	4	2	1	5	2	4	0
4	4	3	1	5	4	4	0
4	3	5	1	5	3	3	0
4	3	2	1	4	5	3	0
3	4	3	1	4	2	5	0
3	3	5	1	4	3	4	0
3	3	2	1	4	5	3	0
3	3	4	1	4	3	5	0
3	3	3	1	3	2	4	0
3	3	3	1	3	1	4	0
3	3	5	1	3	5	2	0
3	3	3	1	3	5	3	0
2	2	4	1	3	3	2	0
2	2	3	1	2	3	3	0
2	2	3	1	2	2	3	0
2	2	3	1	2	3	2	0
1	1	3	1	1	1	3	0
1	1	4	1	1	3	2	0
5	5	5	0	1	4	1	0



「ワークシート名(Data20)」を入力し、実行ボタンをクリックする。

SAS® Studio

- ▼ サーバーファイルとフォルダ
 - **않→ 亩 赴 平 Ⅲ い**
 - Margaret Amage

 Marga
 - 🤼 フォルダショートカット
 - ▲ 🖳 ファイル (ホーム)
 - sasuser.v94
 - ▶ m データ分析の基礎1
 - ▶ データ分析の基礎2
 - データ分析の基礎3
 - EXCELデータ.xlsx



重回帰分析の結果

	Root MSE		0.9	1 407	R2 兼		0.4813	
	従属変数の ³	平均	3.1	00000	調整	ያሉ R2 ∰	0.4533	
	変動係数		29,4	18628				
)<=	·	の推定	•		
			,,,		メータ			
変数	ラベル	自由	康	推	定値	標準誤差	t値	Pr > t
変数 Intercep		自由	1度			標準誤差 0.67665	t値 -1.10	-
		自由	B度 1 1	-0	定値			0.2783 <.0001

- ・「機能」、「デザイン」の偏回帰係数のP値は小さい。
- ・自由度調整済み決定係数 = 0.4533

データ(男性のみ)の場合

「ワークシート名(Data21)」を入力し、実行ボタンをクリックする。



データ(男性のみ)の場合

Root MSE	0.42484	R2 乗	0.8541	
従属変数の平均	2.89474	調整済み R2 乗	0.8358	\supset
変動係数	14.67616			

		パラ	メータの推定			
変数	ラベル	自由度	パラメータ 推定値	標準誤差	t値	Pr > t
Intercept	Intercept	1	0.05533	0.47289	0.12	0.9083
機能	機能	1	0.95629	0.09883	9.58	<.0001
デザイン	チザイン	1	0.03608	0.10498	034 (0.7356

- ・男性は、機能重視
- ・自由度調整済み決定係数 = 0.8358

データ(女性のみ)の場合

「ワークシート名(Data22)」を入力し、実行ボタンをクリックする。



データ(女性のみ)の場合

	Root MSE		0.8	39023	R2 乗			0.6274	
	従属変数の ³	平均	3.2	28571	調整	賽办 R2	垂(05860	
	変動係数		27.0	09389					
			157	メータ	1の推定				
変数	ラベル	自由		パラ	!の推定 メータ 能定値	標準部	髊	t値	Pr > t
変数 Intercept		自由		パラ.	メータ			t 値 -1 39	
		éŧ		パラ. 指 -1	メータ 能定値	標準部	924		Pr > t 0.1822 0.0140

- ・女性は、機能、デザイン共に重視 機能 <デザイン
- ・自由度調整済み決定係数 = 0.5860

男性は、機能重視 女性は、機能、デザイン共に重視

性別によって異なる。



交互作用

まとめ

- ◇相関関係
 - •散布図、積率相関係数
 - •交絡要因(疑似相関)、偏相関係数
- ◇重回帰分析
 - •従属変数、説明変数(独立変数)
 - ・回帰モデルのチェック t 検定 P値(有意確率) 自由度調整済み決定係数
 - ・説明変数相互の積率相関係数
 - ·交互作用