

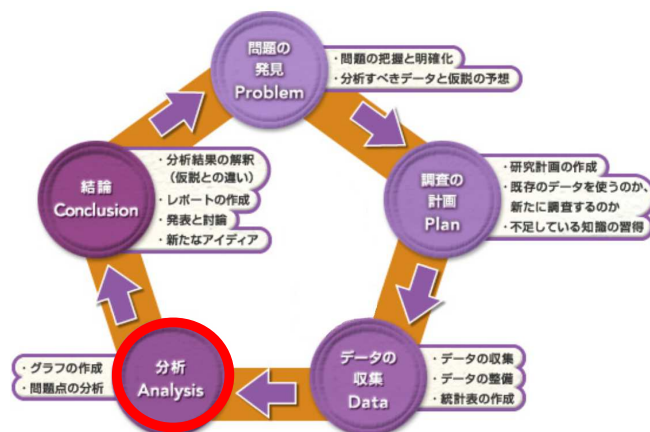
数理・データサイエンス・AI入門

第9回 データサイエンス実践(4)

Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

0

クロス集計を使った 分析と解釈



Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

1

クロス集計の分析

- この表のどこに着目すると良いか一目で分かりますか？

	好き	普通	嫌い	合計
50代以上	5	1	5	11
30代～40代	1	1	10	12
20代以下	2	3	2	7
合計	8	5	17	30

Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

2

クロス集計の分析手順

- まずは総数で全体を割ります

	好き	普通	嫌い	合計
50代以上	5	1	5	11
30代～40代	1	1	10	12
20代以下	2	3	2	7
合計	8	5	17	30

	好き	普通	嫌い	合計
50代以上	0.17	0.03	0.17	0.37
30代～40代	0.03	0.03	0.33	0.40
20代以下	0.07	0.10	0.07	0.23
合計	0.27	0.17	0.57	1.00

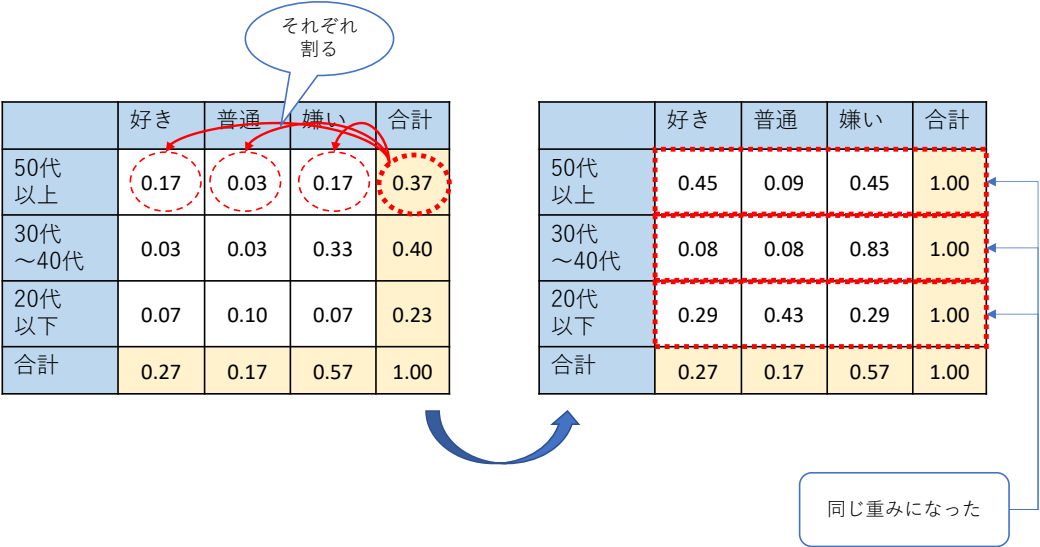
重みがバラバラで比較できない

Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

3

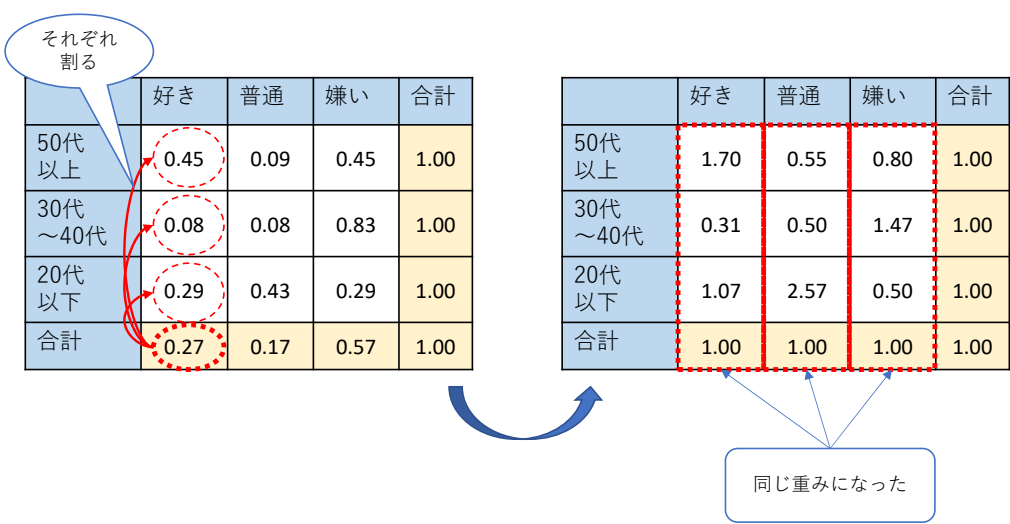
クロス集計の分析手順

- 横方向に割り算して縦項目同士の重みを揃える



クロス集計の分析手順

- 同様に縦方向に割り算して横項目同士の重みを揃える



クロス集計の分析手順

- 最終的に計算された値を見て
数値が特に大きい／小さい箇所にもフォーカスして考察する
- 実数では見えにくい課題を抽出することができる

	好き	普通	嫌い	合計
50代以上	5	1	5	11
30代～40代	1	1	10	12
20代以下	2	3	2	7
合計	8	5	17	30

実数では特に多い／少ないと言えるのかどうか判断できない

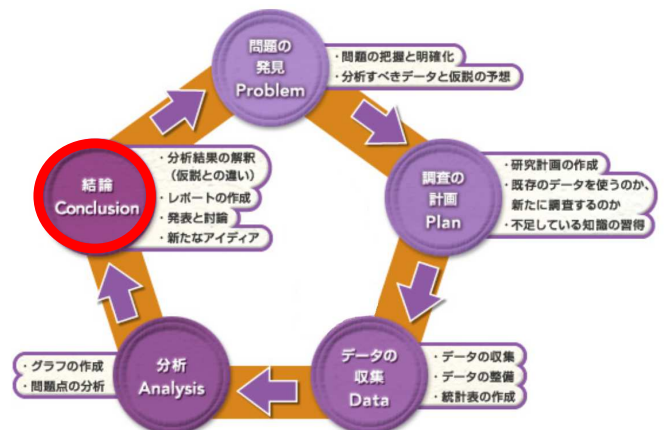
	好き	普通	嫌い	合計
50代以上	1.70	0.55	0.80	1.00
30代～40代	0.31	0.50	1.47	1.00
20代以下	1.07	2.57	0.50	1.00
合計	1.00	1.00	1.00	1.00

30・40代に好まれていない
また、20代は普通評価が特に多いと言える

Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

6

分析結果の読み取りと解釈



Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

7

結果を読み解く際のポイント

□ 結果を考察する

- 結果は他の人にも共有し、他者の視点を入れるようにしましょう
- 新たな観点や見逃していた課題が出て来たら、次のテーマにしましょう

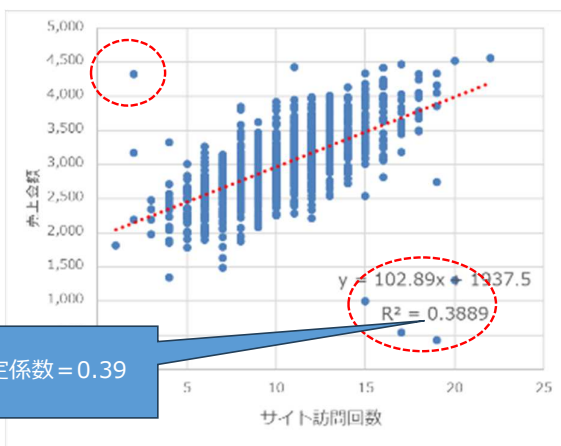
□ 結果を当初の目的と接合させる

- 結果（当初の仮説）ありきで解釈を歪めることが無いようにしましょう
- 仮説と異なれば理由をしっかりと考えましょう

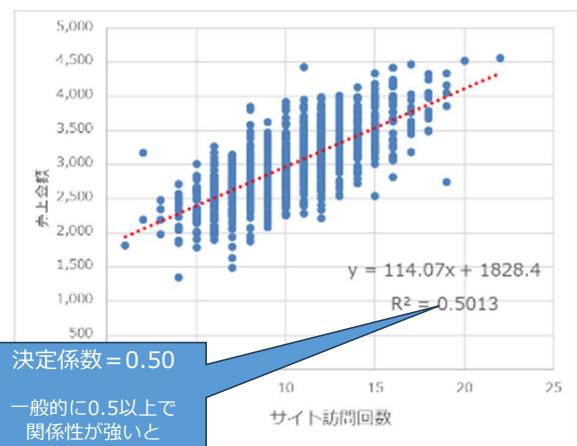
外れ値による影響

□ 関係性が強く出ない場合は外れ値の存在を疑ってみる

外れ値がある結果

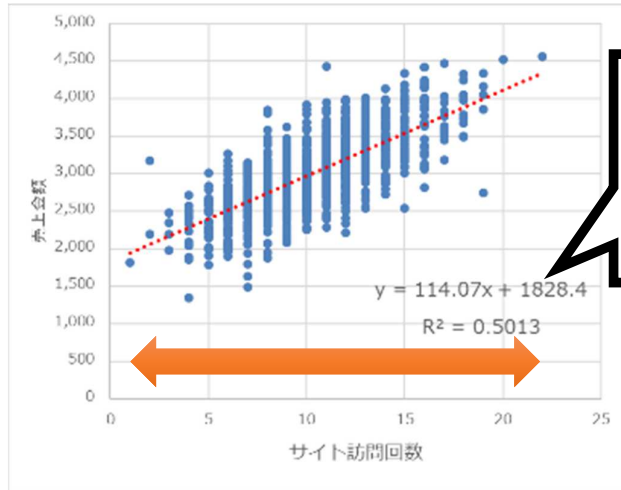


外れ値を取り除いた結果



データのカバー範囲

- 予測対象が元データの範囲外になっていないか確認が必要



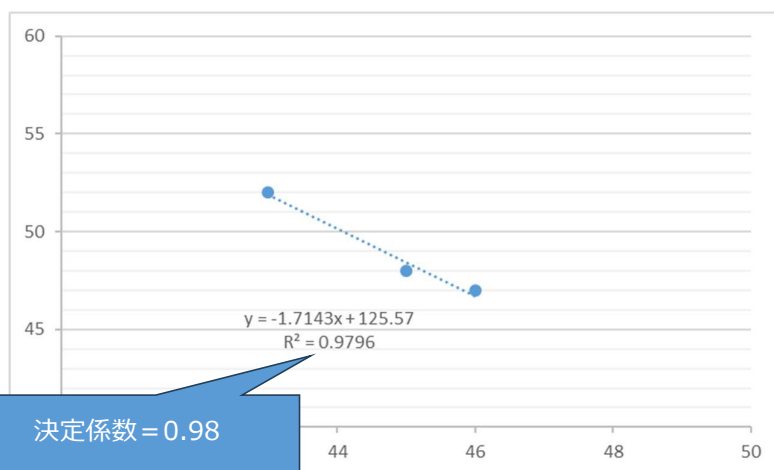
このモデル式なら
サイト訪問回数(x=) 5 回なら売り上げ(y=)約2,500円
サイト訪問回数(x=) 20 回なら売り上げ(y=)約4,000円
になると予測ができる

ならば、訪問回数100回(x=100)のときに
売り上げ(y)は約13,000円になるといえるか？
→データが存在する範囲についての説明しかできません

Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

データ量不足による影響

- 決定係数だけ見ると非常に良い分析結果に見える



決定係数 = 0.98

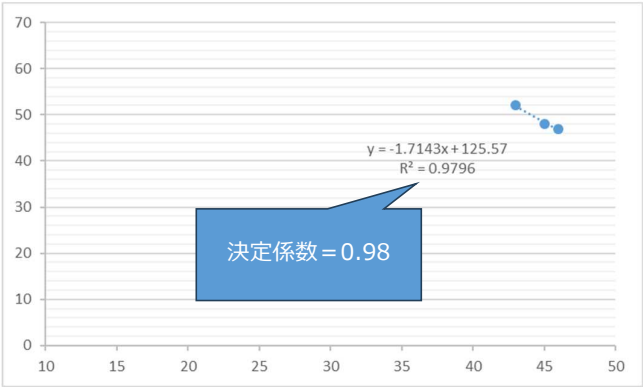
一般的に0.5以上で
関係性が強いと説明されます

Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

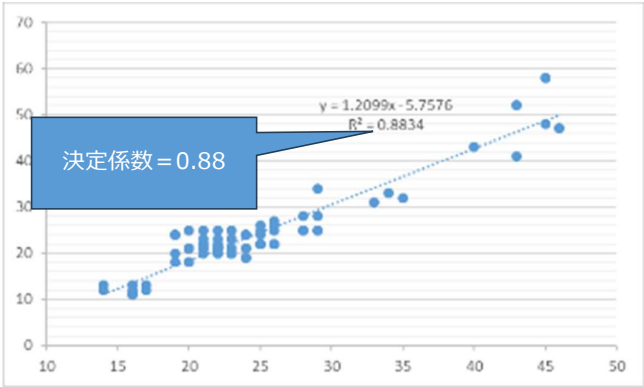
データ量不足による影響

元データを恣意的に選択することで意図的な結果にもできる

データが少ないときの分析結果

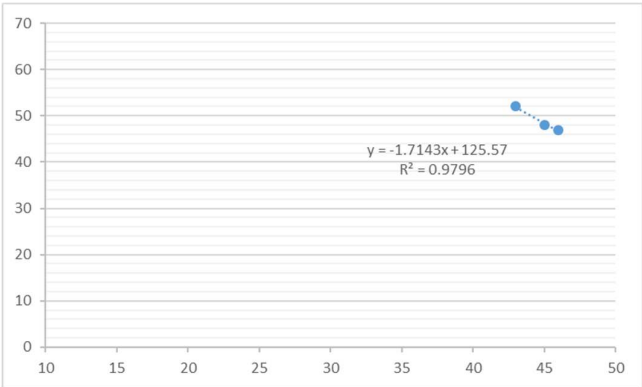


データ量を十分に増やした分析結果



データ量不足による影響

データ量が十分でない場合、モデルの信頼性の指標が著しく悪化



概要

回帰統計	
重相関 R	0.989743319
重決定 R2	0.979591837
補正 R2	0.959183673
標準誤差	0.3086067
観測数	3

分散分析表

	自由度	変動	分散	観測された分散比	有意 F
回帰	1	4.571428571	4.571428571	48	0.091257897
残差	1	0.095238095	0.095238095		
合計	2	4.666666667			

Excelで回帰分析した場合は
有意Fがモデルの信頼性の指標

仮説検定

仮説検定

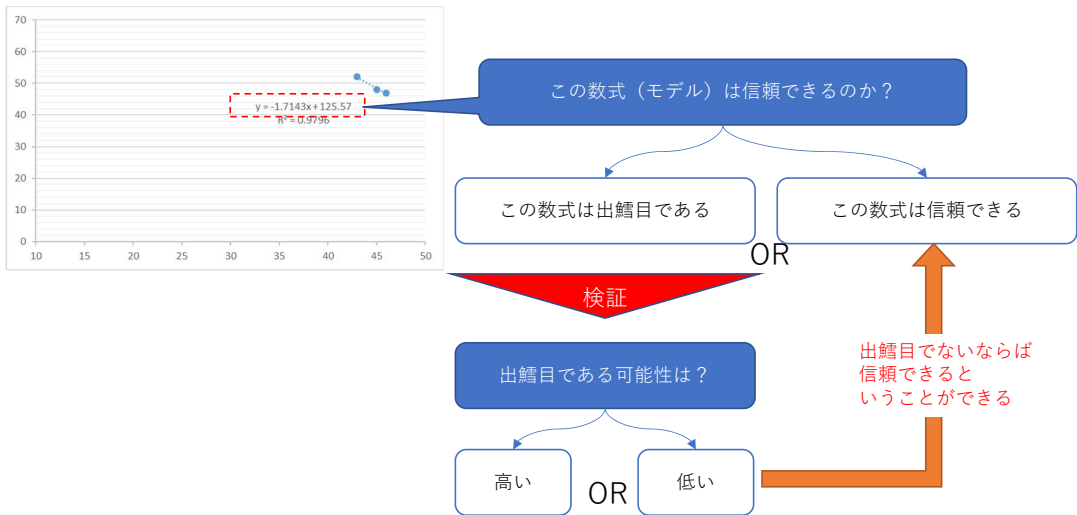
- 仮説検定とは、各種の統計分析を行う上で算出した統計量がどの程度信頼できるのかを判定する手法
 - 分析結果の信頼性を検証するためには仮説検定の考え方を 사용합니다

- 用語の理解
 - 対立仮説 = 本来主張したい内容
 - 帰無仮説 = 主張の逆の内容

仮説検定

モデルの信頼性を仮説検定で検証

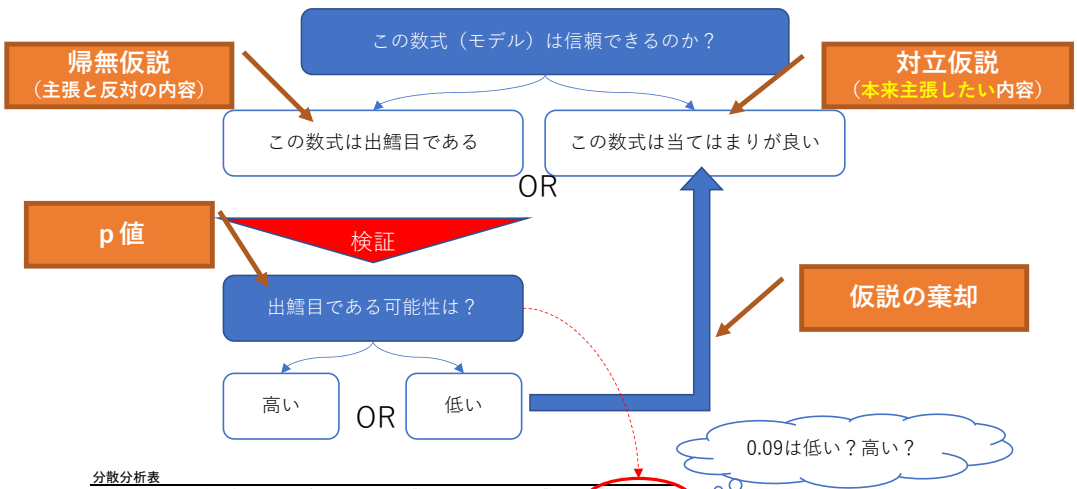
- ・ 検証したい回帰式が「出鱈目である」という仮説（帰無仮説）を立てます
- ・ その可能性が低い場合「当てはまりが良い」と判断します



Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

回帰分析における仮説検定

モデル全体の信頼度は「有意 F」の数値で判断



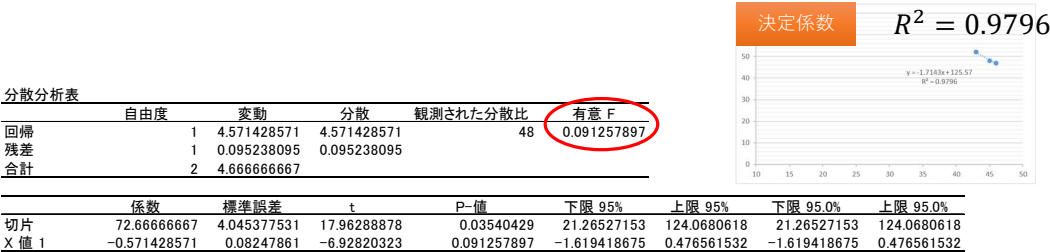
分散分析表								
	自由度	変動	分散	観測された分散比	有意 F			
回帰	1	4.571428571	4.571428571	48	0.091257897			
残差	1	0.095238095	0.095238095					
合計	2	4.666666667						

	係数	標準誤差	t	P-値	下限 95%	上限 95%	下限 95.0%	上限 95.0%
切片	72.66666667	4.045377531	17.96288878	0.03540429	21.26527153	124.0680618	21.26527153	124.0680618
X 値 1	-0.571428571	0.08247861	-6.92820323	0.091257897	-1.619418675	0.476561532	-1.619418675	0.476561532

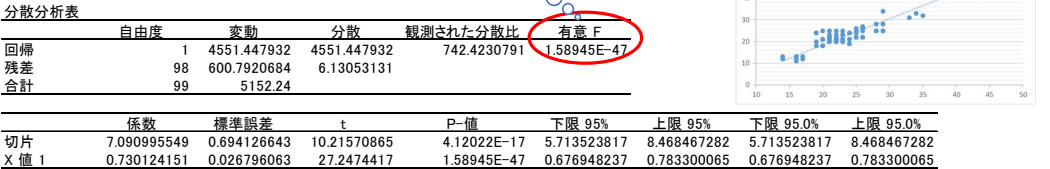
Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

回帰分析における仮説検定

- 2つのモデルを比較
 - 下のほうが相関係数はやや下がるものの信頼性は非常に高いといえます



小数以下47桁までゼロ
= 限りなくゼロ



Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved. 18

演習

演習内容

- SDGsの取り組みに課題がある自治体が、なぜ推進されないのかデータ分析を通じて原因を調査したいと思います
 - 今日のゴール
 - SDGsの取り組みに対して
 - 課題がない自治体を「0」
 - 課題がある自治体を「1」
- として、各自が考えた条件で分類した結果を「目的変数」の列に記入していきましょう

作業内容の確認

□ ここまでの作業のおさらい



Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

20

作業内容の確認

□ ここまでの作業のおさらい



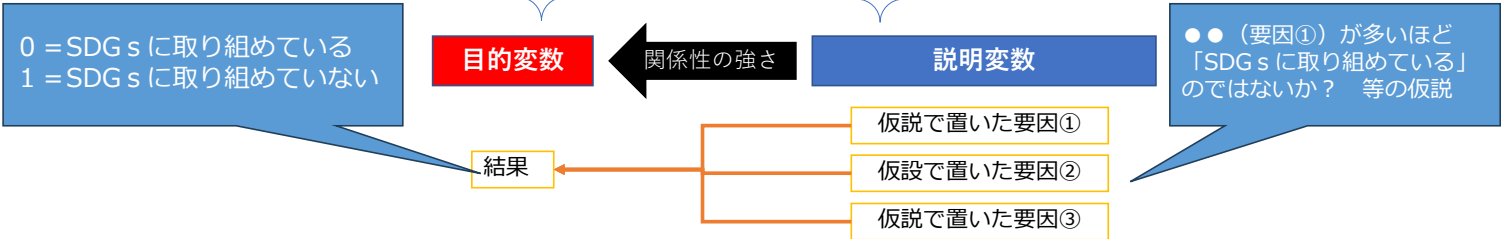
Copyright UCHIDA HUMAN DEVELOPMENT all rights reserved.

21

作業内容の確認

これから行うこと

自治体名	目的変数	説明変数1	説明変数2	説明変数3	説明変数4	説明変数5	説明変数6
北海道	1	(各自で用意したデータ)					
札幌市	1						
函館市	0						
小樽市	0						
旭川市	0						
室蘭市	0						



作業内容の確認

結果を記録

自治体名	目的変数	説明変数1	説明変数2	説明変数3	説明変数4	説明変数5	説明変数6
北海道	1						
札幌市							

この「列名」を記入します

分析結果の記録	説明変数	①	②	③	④	⑤	⑥	⑦	⑧	⑨	⑩	補正R2
列名を記入		婚姻有無	正社員か否か	パートか否か	最大連続利用月数	年間購入回数	サイト閲覧時間	サービスへの満足度	年齢	性別		
1回目		○	○	○	○	○	○	○	○	○		→ 0.438
2回目				○	○	○	○	○		○		→ 0.441
3回目							○	○				→ 0.443
4回目												→
5回目												→

分析の進め方（1回目の実行）

□ 最初はそのまま全部を説明変数にして分析してみます

説明変数 → 関係性の強さ → 目的変数

婚姻有無
正社員か否か
パートか否か
最大連続利用月数
年間購入回数
サイト閲覧時間
サービスへの満足度
年齢
性別

お得意様は否か

分析の進め方（2回目の実行）

□ 既婚と正社員と年齢を外して分析してみます

説明変数 → 関係性の強さ → 目的変数

婚姻有無
正社員か否か
パートか否か
最大連続利用月数
年間購入回数
サイト閲覧時間
サービスへの満足度
年齢
性別

お得意様は否か

分析の進め方（3回目の実行）

□ さらに説明変数を絞って分析してみます

説明変数 → 関係性の強さ → 目的変数

婚姻有無
正社員か否か
パートか否か
最大連続利用月数
年間購入回数
サイト閲覧時間
サービスへの満足度
年齢
性別

お得意様は否か