実務直結! 分析カ向上ウェビナーシリーズ 機械学習によるビッグデータ分析の手法

#3 クラスター分析による分類(2) クラスター分析の応用 と 階層的クラスタリング

2022年12月15日



Agenda

・クラスター分析の応用(他の分析手法との組み合わせ)

- 主成分分析により説明変数を要約する
- 主成分軸でクラスター分析を行う

・クラスター分析による分類(2):階層的クラスタリング

- 階層的クラスタリング(群平均法、重心法、Ward法)のしくみ
- 樹形図(デンドログラム)とクラスタ数の検討
- 都道府県データを用いて階層的クラスタリングにより類似地域を分析する

・今後のデータサイエンス学習に向けたスキルアップ

- データサイエンティストに求められるスキル
- SAS内サンプルデータの紹介と使い方
- オープンデータの紹介





代表的な機械学習手法

- ・ 機械学習手法は、教師あり、教師なし、強化学習に大別される
- ・なかでも、教師あり分類、教師なし分類は極めて基本的かつ頻用される手法である



Agenda

・クラスター分析の応用(他の分析手法との組み合わせ)

- 主成分分析により説明変数を要約する
- 主成分軸でクラスター分析を行う
- ・クラスター分析による分類(2):階層的クラスタリング
 - 階層的クラスタリング(群平均法、重心法、Ward法)のしくみ
 - 樹形図(デンドログラム)とクラスタ数の検討
 - 都道府県データを用いて階層的クラスタリングにより類似地域を分析する
- ・今後のデータサイエンス学習に向けたスキルアップ
 - データサイエンティストに求められるスキル
 - SAS内サンプルデータの紹介と使い方
 - オープンデータの紹介





教師あり学習と教師なし学習



教師あり学習





教師なし学習のイメージ (クラスタリング)

- 各データ間の距離に基づき、近接データ(=類似度が高いデータ)同士のグループ(クラスタ)を作り、 データを分類する手法
- ・ 学習データなしでデータを大きく層別したい場合に有効





主成分分析の概要

- ・主成分分析は、多数の説明変数が存在する場合に、(それらの分散構造を考慮して)変数 を合成していくことで、より少ない変数(=主成分)でデータを説明しようとするアプローチ
- ・アンケート調査や官能評価でよく用いられるほか、分析前の次元削減としても活用される



▼主成分分析のイメージ

▼官能評価における活用例



Source: https://www.nodai.ac.jp/research/teacher-column/22913/





- 主成分分析
- 主成分に対するk-meansクラスタリング





使用データ

- UCI Machine Learning Repositoryでは様々な分野のデータが公開
- ・今回は、銀行のマーケティングデータを活用し、分析を行う



Bank Marketing Data Set

Download: Data Folder, Data Set Description

Abstract: The data is related with direct marketing campaigns (phone calls) of a Portuguese banking institution. The classification goal is to predict if the client will subscribe a term deposit (variable y).

Data Set Characteristics:	Multivariate	Number of Instances:	45211	Area:	Business
Attribute Characteristics:	Real	Number of Attributes:	17	Date Donated	2012-02-14
Associated Tasks:	Classification	Missing Values?	N/A	Number of Web Hits:	1577437

Source:

[Moro et al., 2014] S. Moro, P. Cortez and P. Rita. A Data-Driven Approach to Predict the Success of Bank Telemarketing. Decision Support Systems, Elsevier, 62:22-31, June 2014

Data Set Information:

The data is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls. Often, more than one contact to the same client was require ('yes') or not ('no') subscribed.

There are four datasets:

bank-additional-full.csv with all examples (41188) and 20 inputs, ordered by date (from May 2008 to November 2010), very close to the data analyzed in [Moro et al., 2014]
 bank-additional.csv with 10% of the examples (4119), randomly selected from 1), and 20 inputs.
 bank-full.csv with all examples and 17 inputs, ordered by date (older version of this dataset with less inputs).
 bank.csv with 10% of the examples and 17 inputs, randomly selected from 3 (older version of this dataset with less inputs).
 bank.csv with 10% of the examples and 17 inputs, randomly selected from 3 (older version of this dataset with less inputs).

The classification goal is to predict if the client will subscribe (yes/no) a term deposit (variable y).

https://archive.ics.uci.edu/ml/datasets/bank+marketing





4,521人



• 4,521人分の顧客について、顧客情報や営業アプローチ状況、最終的な狙いである「定期預金の契約有無」に関する情報(計17列)が格納されている

※クラウド型のSAS Studio (SAS OnDemand for Academics) において 列名を日本語にする場合、

			クレジット 債務不履行	カード テの有無	年間平5 (ユー	均残高 -□)			最終連絡 会話時間	各時の (秒)	キャンペーン 連絡回望	·中の 最終 数 約	経連絡からの 経過日数	キャンペーン 連絡回	v前の 前回キ 数 の	キャンペーン)結果
年齢	職業	結婚歴	学歴	クレカ債務	年間平均 残高	住宅 ローン	個人 ローン	連絡手段	最終連 絡日	最終連 絡月	最終会話 時間	CP中連絡 回数	最終連絡 日数	CP前連絡 回数	前回CP結果	定期預金 契約
30	unemployed	married	primary	no	1787	no	no	cellular	19	oct	79	1	-1	0	unknown	no
33	services	married	secondary	no	4789	yes	yes	cellular	11	may	220	1	339	4	failure	no
35	management	single	tertiary	no	1350	yes	no	cellular	16	apr	185	1	330	1	failure	no
30	management	married	tertiary	no	1476	yes	yes	unknown	3	jun	199	4	-1	0	unknown	no
59	blue-collar	married	secondary	no	0	yes	no	unknown	5	may	226	1	-1	0	unknown	no
35	management	single	tertiary	no	747	no	P		23	feb	141	2	176	3	failure	
36	self-employed	married	tertiary	no	307	yes	説即	[変数]	14	may	341	1	330	2	other 😑	的変数
39	technician	married	secondary	no	147	yes	L/0-7		6	may	151	2	-1	0	unkno	
41	entrepreneur	married	tertiary	no	221	yes	no	unknown	14	may	57	2	-1	0	unknown	no
43	services	married	primary	no	-88	yes	yes	cellular	17	apr	313	1	147	2	failure	no
39	services	married	secondary	no	9374	yes	no	unknown	20	may	273	1	-1	0	unknown	no
43	admin.	married	secondary	no	264	yes	no	cellular	17	apr	113	2	-1	0	unknown	no
36	technician	married	tertiary	no	1109	no	no	cellular	13	aug	328	2	-1	0	unknown	no
20	student	single	secondary	no	502	no	no	cellular	30	apr	261	1	-1	0	unknown	yes
31	blue-collar	married	secondary	no	360	yes	yes	cellular	29	jan	89	1	241	1	failure	no
10			1		101			11 1	20		100	2		0	1	

予測(分析)対象を 説明するための変数



10 S.Sas



- 今年のウェビナーでは、SAS Studio でデモを行います。
- SAS Studio はすべてのSAS製品に付帯しているGUI で、今回は学習用に自宅でもお使い 頂けるクラウド型無償版 SAS OnDemand for Academics を使っています。
 (※無償版の登録については、SAS からの申込完了メールをご参照ください)
- ・なお、SAS Studio起動時はコード入力画面となっていますが、画面右上の「SASプログラマ」を「ビジュアルプログラマ」に変更するとデモと同様の入力画面となります。

▼SAS Studio 画面イメージ

▼GUI画面への変更方法 (ビジュアルプログラマ)





参考:SAS Studio 起動方法

- SAS OnDemamd for Academics にログイン後、Dashboard より SAS Studio を起動
- ・ 起動後、前頁の通り、右上メニューより「ビジュアルプログラマ」を選択





データの読み込み (1/2)

① 左パネル内の 「アップロード」アイコン をクリック



②「ファイルの選択」ボタンをクリックし、ファイル選択画面で
 "bank_marketing.xlsx"を選択し、OKボタン
 ③「アップロード」ボタンをクリック

ファイルのアップロード	
ファイルのアップロード先: /home/u62013505	
ファイルの選択	
選択済みファイル:	
1 XLSX bank_marketing.xlsx	371.1 kb
	アップロード キャンセル

④左パネル内にファイルがアップロードされていることを確認

SAS®	² Studio	
• サ-	ーバーファイルとフォルダ	
Ľ	* 竜 圭 平 国 55	_
49	🖥 odaws02-apse1-2	*
	🔁 フォルダショートカット	1
	🖌 📮 ファイル (ホーム)	
_	sasuser.v94	
	🔀 bank_marketing.xlsx	





データの読み込み (2/2)

① 左パネル内の "bank_marketing.xlsx" を選択し、 画面右側のプログラムエリアにドラッグ & ドロップ



③詳細設定画面が開くので、実行ボタンをクリック (特に各設定は変更不要)



②右側のプロセスフローにノードが生成されるので、 当該ノードをダブルクリック



④「結果」のタブ画面に読み込んだデータの概要が出力

ファイル名: bank_mark	ting.xlsx			
ソースの場所: /home/u620	13505			
ワークシート名: 第1ワークシート				
	^{吉果} :力: 2	F-9		
• 目次		CONTENTS プロシジャ		
5	ータセット名	WORK.IMPORT1	オブザベーション数	4521
,	ンバータイプ	DATA	変数の数	17
3	ンジン	V9	インデックス数	0
1	成日時	2022/08/08 09:34:47	オブザベーションのバッファ長	120
3	新日時	2022/08/08 09:34:47	削除済みオブザベーション数	0
9	2 0		圧縮済み	NO
5	ータセットタイプ		ソート済み	NO
-	ベル			
5	ータ表現	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64		
3	ンコード	utf-8 Unicode (UTF-8)		
		エンジン/ホスト関連情報		
データセットのページサ	131072			
	5			
データセットのページ数				
データセットのページ数 データページの先頭	1			
データセットのページ数 データページの先頭 ページごとの最大OBS数	1090			





読み込んだデータの確認

データ概要の確認

新し		ブラウ	゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙゙	で 厚	剥く				
6 •	日 日 日 日 日 日 日 日 日 日 日 日 日 日 日 日 日 日 日	4 4	rt 2	,					
						CONT	ENTS 70 4	U ++	1
			データセット名 V	VORK.IMPO	ORT1	Contr	Lino y Li y	オブザベーション数	4521
					CON	TENTS プロ	シジャ		
	データ・	ヤットタ	WORK IMPORTS					オブザベーション数	4521
	*2/5-	-917	DATA					安敷の数	17
	エンジ	2	V9					36,50,77,50	
	作成日	38	2022/08/08 09:34:4	7				オブザベーションのパッファ	120-
	更新日		2022/08/08 09:34:4	7				石川米ケレイニ米ケオ	「な」の「
	保護		Lolli co co colo ini						「「一日日日」
	データ	セットタイフ	,					(ビッグデータイン)	后の其木)
	ラベル								川の至平)
	データ	表現	SOLARIS X86 64.	LINUX	X86 64	ALPHA T	RU64, LINU	JX_IA64	
	בעב	-8	utf-8 Unicode (UTF	-8)		_			
-					エンシ	ジン/ホスト	司連情報		
データセットのページ	191X	131072							
テータセットのペーン	/90X	5							
デージストの見古の町	c #/r	1							
作用なージの086数	590	1090							
光明ページのOBS数		1052							
テージセットの修復数	K	U (an an unadul C	0			-1 0 - d			and and any forward and The dat
ノアイル石		/saswork/5/	45_WORK/1F600001F3F	A_odaws	su i -aps	e1-2.00a.sa	IS.COM/SAS	_workC7660001F3FA_odaws01-apse1-	2.00a.sas.com/import1.sas/bdat
作成したまでト		Linux							
日本のためた		22950							
アクセス権限		55650 rw-rr							
所有考え		062013505							
ファイルサイズ		768KB							
ファイルサイズバイ	5	FILT	<u> </u>		T	たきす			
	合	, טניצ	ワーク	£7	ደው	目記			
				変数	と属性	リスト (アル	ファベット	順)	
		# 3	乏数	タイプ	長さ	出力形式	入力形式	ラベル	
		13 =	キャンペーン中の連絡	数值	8	BEST.		キャンペーン中の連絡回数	
		15 =	キャンペーン前の連絡	数值	8	BEST.		キャンペーン前の連絡回数	
		5	フレジットカード債務	文字	3	\$3.	\$3.	クレジットカード債務不履行有無	
		7 1	主宅ローンの有無	文字	3	\$3.	\$3.	住宅ローンの有無	
		8 1	固人ローンの有無	文字	3	\$3.	\$3.	個人ローンの有無	
		16	前回キャンペーンの結	文字	7	\$7.	\$7.	前回キャンペーンの結果	
		4 🕾	許歷	文字	9	\$9.	\$9.	学歴	
		17 🎗	定期預金契約有無	文字	3	\$3.	\$3.	定期預金契約有無	
		6 1	∓間平均残高(ユーロ	数值	8	BEST.		年間平均残高(ユーロ)	
		1 4	手給 令	数值	8	BEST.		年齢	
		14 3	最終連絡からの経過日	数值	8	BEST.		最終連絡からの経過日数	
		10 🛔	最終連絡日	数值	8	BEST.		最終連絡日	
		12 🗄	最終連絡時の会話時間	数值	8	BEST.		最終連絡時の会話時間(秒)	
		11 3	最終連絡月	文字	3	\$3.	\$3.	最終連絡月	
		3 8	吉婚歷	文字	8	\$8.	\$8.	結婚歴	
		2 1	11年	文字	13	\$13.	\$13.	職業	
		9 3	主格手段	文字	9	\$9.	\$9.	連絡手段	

					9 🖨 🤅	SAS プログラ	र - 😑 🕐 🖽	ンアウト
コグラム1 × 🎜 *bank_marketin	g ×							
コード/結果 分割 🏒 🖬	R 20						箇 ログ 【	🕮 🗆 — H
イル情報								
スファイル								
イル名: bank_marketing.xlsx								
·スの場所: /home/u62013505								
・クシート名:								
ワークシート								
データ								
Server: SASApp								
タセット名: IMPORT1								
プラリ: WOHK	<u>.</u> »			~				
更	出力テ	-タ 画	面より	0.				
ション エ				- <u>-</u> T	ta = 33			
イルの種類:		テエの・	 -^	ノケヤ	住認			
	x.フル・/U		/	· · · ·				
7ォルト(ファイル拡張子に基づく)				۹ ک ۲				
フォルト (ファイル拡張子に基づく)			, , 		-			
フォルト (ファイル拡張子に基づく) ード ログ 結果 ブル: WOPK IMPOPT1					1 1			
**ルト (ファイル拡張子に基づく) ード ログ 結果 ブル: WORK.IMPORT1 -	出力データ - ・ ビュー: 列名 *		ア フィルタ: (オ					
オルト (ファイル拡張子に基プス) ード ログ 結果 ブル: WORK.IMPORT1 ▼ 8	出力データ ビュー:列名 × の 合計行数:452	小 こ 二 し し 、 、 、 、 、 、 、 、 、 、 、 、 、	 アフィルタ: (ボ 8業	C P こここ なし)		クレジット・	<u>たまた。それでのの</u> 全間平均残高(ユーロ	
オルト (ファイル拡張子に基プイ ード ログ 結果 ブル: WORKIMPORT1 ↓ ↓ t	出力データ	いて「「 ここ」の ここの 日 「 て 1 合計列数:17 年齢 調 30 m	フィルタ: (水 数業 nemploved	なし) 結婚歴 married	」 ↓ 学歴 primary	クレジット	注 ま ミュック 年間平均残高(ユーロ 1787	◆ → 住宅口
 オルト(ファイル拡張子に基づく) ド ログ 結果 ブル: WORKIMPORT1 ▼ 1 すべて選択 ● 年齢 	<u>出力データ</u> ビュー:列名 · の 合計行数:452	いてしていた。 ここのでは、 このででは、 このででは、 このででは、 このでで このでで このでで、 このでで このでで このでで このでで こので このでで このでで こので	マフィルタ: (水 電業 inemployed ervices	なし) 結婚歴 married married	学歴 primary secondary	クレジット no no	使 を 61100 年間平均残高 (ユーロ 1787 4789	◆ →1 住宅口 no ves
 オルト(ファイル拡張子に基づく) ード ログ 結果 ブル: WORKIMPORT1 * すべて選択 ●年約 ▲ 職業 	<u>出力データ</u> ビユー: 列名 · の 合計行数: 452 1 2 3	3 3 3 3 4 3 3 5 5 5 5 5 5 5 5 5 5 5 5 5	マフィルタ: (水 電業 nemployed ervices nanagement	なし) 結婚歴 married married single	学歴 primary secondary tertiary	クレジット no no	たまた1000 年間平均残高(ユーロ 1787 4789 1350	● ● 住宅口 no yes
 オルト(ファイル拡張子に基づく) ード ログ 結果 ブル: WORKIMPORT1 * すべて選択 ● 年齢 ▲ 職業 ▲ 結婚歴 	出力データ ビュー: 列名 の 合計行数:452 1 2 3 4	3. 上 5. 日 9 1 合計列数-17 年齢 題 30 ur 33 se 35 m 30 m	マイルタ: (オ 意業 inemployed ervices nanagement nanagement	なし) 結婚歴 married married single married	学歴 primary secondary tertiary	クレジット no no no	体 を 長 1000 年間平均残高 (ユーロ 1787 4789 1350 1476	中 中 no yes yes
 ォルト(ファイル拡張子に基づく) ドログ 結果 ブル: WORKIMPORTI ● すべて選択 ● 年齢 ● 年齢 ▲ 職業 ▲ 結婚歴 ▲ 学歴 	出力データ ● ● ビュー: 列名 ● 1 2 3 3 4 5	この目して、 この目して、 1 合計列数-17 年齢 署 30 ur 33 se 35 m 30 m 59 bl	マフィルタ: (オ 電業 inemployed ervices nanagement nanagement ilue-collar	なし) 結婚歴 married married single married married	学歴 primary secondary tertiary tertiary secondary	クレジット no no no no	体 会 G 1000 年間平均残高(ユーロ 1767 4769 1350 1476 0	中 住宅口 no yes yes yes yes
 ホルト(ファイル拡張子に基づく) ドログ 結果 ブル: WORKIMPORTI ● すべて選択 ● 年齢 職業 A 結婚歴 学歴 クレジットカード債務 	世力データ ビコー: 別名 の 合計行数:452 1 2 3 4 5 6	こ 4 (1) 日 (1) (1) (1) (1) (1) (1) (1) (1) (1) (1)	マフィルタ: (オ 事業 inemployed ervices nanagement nanagement ilue-collar nanagement	結婚歴 married married single married married single	学歴 primary secondary tertiary tertiary tertiary	クレジット no no no no no no	年間平均残高(ユーロ 1787 4789 1350 1476 0 0 747	中 中 中 中 中 中 中 中 中 中 中 中 中 中
 オルト(ファイル拡張子に基づく) ト ログ 結果 アル: WORKIMPORT1 ● すべて選択 ● 年齢 ▲ 職業 ▲ 結婚歴 今学歴 ヘクレジットカード債務 ● 年間平均残高(ユーロ 	 出力データ・ ビュー:列名・ 会社行数:452 1 2 3 4 5 6 7 	こ 2 (3) 目 「 「 (1 会計初数: 17 年齢 目 30 ur 33 se 35 m 30 m 59 bl 35 m 36 se	マイルタ: (水 車業 nemployed ervices nanagement anagement ilue-collar nanagement elf-employed	なし) 結婚歴 married married single married single married	♥ ♥ Primary secondary tertiary tertiary tertiary	クレジット no no no no no no no	年間平均残高(ユーロ 1787 4789 1350 1476 0 747 307	住宅口 no yes yes yes yes no yes
 オルト(ファイル拡張子に基づく) ード ログ 結果 ブル: WORKIMPORT1 → 1 すべて選択 ③ 年齢 ▲ 職業 ▲ 結婚歴 ▲ 学歴 ▲ クレジットカード債務 ③ 年間平均残高(ユーロ ▲ 住宅ローンの有無 	 出力データ ゴー:列名 会社行数:452 1 2 3 4 5 6 7 8 	 このののでは、 このののののののののののののののののののののののののののののののののののの	マイルタ: (化 意葉 nemployed ervices nanagement ilue-collar nanagement elf-employed achnician	なし) 結婚歴 married married single married married married married	学歴 primary secondary tertiary tertiary tertiary tertiary secondary	クレジット no no no no no no no	年間平均残高(ユーロ 1787 4789 1350 1476 0 747 307 147	住宅口 no yes yes yes yes no yes yes yes
 オルト(ファイル拡張子に基づく) ード ログ 結果 ブル: WORKIMPORT1 → 1 すべて選択 ● 年齢 ▲ 職業 ▲ 結婚歴 ◇ 学歴 ▲ クレジットカード債務 ● 年間平均残高(ユーロ ▲ 住宅ローンの有無 ▲ 個人ローンの有無 	出力データ - ・ ジュー: 列名 ● の 合計行数:452 1 2 3 4 5 6 7 8 9	Q 単 い 回 可 1 会社初数: 17 年齢 間 30 ur 33 se 35 m 30 m 59 bl 35 m 36 se 39 te 41 er	マイルタ: (水 意業 nemployed ervices nanagement lua-collar nanagement elf-employed echnician ntrepreneur	結婚歴 married married single married single married married married	Primary secondary tertiary secondary tertiary secondary tertiary secondary tertiary	クレジット ハロ ハロ ハ	体 を た 100 年間平均残高 (ユーロ 1787 4789 1350 1476 0 747 307 1477 221	dt宅口 no yes yes yes no yes yes yes yes yes yes yes
 オルト(ファイル拡張子に基づく) ド ログ 結果 ブル: WORKIMPORT1 ▼ 10 すべて選択 ● 年齢 ▲ 職業 ▲ 結婚歴 学歴 ◇ クレジットカード債務 ● 年間平均残高(ユーロ ▲ 住宅ローンの有無 ▲ 個人ローンの有無 ▲ 連絡手段 	出力データ ● ● ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ ○	日本部 第 1 会計初時:17 年齢 第 30 m 33 se 33 m 30 m 50 bi 35 m 36 se 37 te 41 er 41 se	マイルタ: (A 事業 nemployed ervices nanagement elue-collar nanagement elf-employed echnician ntrepreneur ervices	結婚歴 married married single married married married married married	Primary secondary tertiary tertiary secondary tertiary secondary tertiary primary	クレジット no no no no no no no no no no	生また1100 年間平均残高(ユーロ 1787 4789 1350 1476 0 747 307 147 307 147 - 8	住宅口 no yes
 オルト(ファイル拡張子に基づく) ドログ 結果 ブル: WORKIMPORT1 → 10 すべて選択 ● 年齢 ▲ 職業 ▲ 結婚歴 ▲ 学歴 ▲ クレジットカード債務 ● 年間平均残高(コーロ ▲ 住宅ローンの有無 ▲ 運絡手段 IVディ 値 	 出力データ ・ ・<	3 品 (3 画) 「「 (3 合計列数-17 年齢 = 30 ur 33 sc 35 m 30 ur 33 sc 35 m 30 m 30 m 30 ur 33 sc 35 m 30 ur 33 sc 35 m 30 sc 35 m 36 sc 39 tc 41 cc 43 sc 39 sc 43 sc 39 sc	Remployed ervices nanagement nanagement elf-employed echnician ntrepreneur ervices	なし) 結婚歴 married married married married married married married married	学歴 primary secondary tertiary tertiary tertiary tertiary secondary tertiary secondary tertiary	クレジット PO PO PO PO PO PO PO PO PO PO	生 を G 1000 年間平均残高(ユーロ 1787 4789 1350 1476 0 747 307 1476 - 0 747 307 1476 - 8 934	A defer A
 オルト(ファイル拡張子に基づく) ード ログ 結果 ブル: WORKIMPORT1 → 1 すべて選択 ② 年齢 ▲ 結婚歴 ▲ 学歴 ▲ クレジットカード債務 ③ 年間平均残高(ユーロ ▲ 住宅ローンの有無 ▲ 遠絶手段 バティ 値 バティ 値 	 出力データ ・ ・<	この日本の目的では、1000000000000000000000000000000000000	マイルタ: (パ 事業 nemployed ervices nanagement ilua-collar elf-employed achnician ntrepreneur ervices ervices dmin.	なし) 結婚歴 married married married married married married married married married	₽ ₽ ₽ primary secondary tertiary secondary tertiary tertiary secondary tertiary secondary secondary secondary	クレジット PO	体 を 日 1000 年間平均残高(ユーロ 1787 4789 1350 1476 0 747 307 747 307 147 307 147 307 147 307 307 147 307 147 221 -88 335 24 35 24 35 35 35 35 35 35 35 35 35 35	the second
x ルト (ファイル拡張子に基づく) ード ログ 結果 ブル: WORKIMPORT1 → 1 すべて選択 ③ 年齢 ▲ 糖婚歴 ▲ 学歴 ▲ クレジットカード債務 ③ 年間平均残高(ユーロ ▲ 住宅ローンの有無 ▲ 個人ローンの有無 ▲ 連絡手段 1パディ 値 パル 」	 出力データ ・ ・<	こ 2 (3) 日 (1) (1) 合計列数: 17 年齢 日 30 ur 33 se 35 m 30 m 59 bl 35 m 36 se 39 te 41 er 43 se 39 se 41 ar 43 se 39 se 43 ac 36 se 39 te 43 se 39 se	マイルタ: (ボ 事業 nemployed ervices nanagement ilue-collar nanagement ervices ervices ervices ervices ervices ervices	結婚證 married married single married single married married married married married married	Primary Primary secondary tertiary secondary tertiary tertiary secondary tertiary secondary tertiary primary secondary tertiary	クレジット	年間平均残高(ユーロ 1787 4789 1350 1476 0 0 747 307 147 221 - 688 9374 244 244 1109	・ ・ ・ ・ ・ ・ ・ ・ ・
マオルト(ファイル拡張子に基づく) ード ログ あまい マイン ブル: WORKIMPORT1 * 「すべて選択 ③ 年齢 ▲ 職業 ▲ 結婚歴 ◇ 学歴 ◇ クレジットカード債務 ③ 年間平均残高(ユーロ ▲ 住宅ローンの有無 ▲ 個人ローンの有無 > 運絡手段 1/(ディ) 値 1/(ディ) 値	 出力データ ゴージ 列名 会社行数・452 1 2 3 4 5 6 7 8 9 10 11 12 13 14 	C	T71149: (1 memployed ervices nanagement nanagement elf-employed echnician ntrepreneur ervices ervices dmin. echnician tudent	結婚歴 married married single married single married married married married married married single	₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽ ₽	クレジット	生また100 年間平均残高(ユーロ 1787 4789 1350 1476 0 0 747 307 747 307 747 307 1476 201 4 4 4 307 4 4 4 307 4 4 4 4 307 1476 307 4 4 5 2 2 14 5 6 374 4 307 5 2 14 5 6 374 30 5 7 8 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 30 5 374 375 375 374 377 377 377 377 377 377 377 377 377	

生データの確認





作成したプロセスフローの保存(別名で保存)

プロセスフローをクリックしてプロセスフロー画面に戻る

€ #1_ロジスティック回帰.cpf ×
- <u>#1_ロジスティック回帰</u> 〉 "bank_marketing.xlsx" のイン
設定 コード/結果 分割 🖌 <table-cell> 🔀</table-cell>
オプション ノード
▼ ファイル情報
ソース ファイル
ファイル名: bank_marketing.xlsx
ソースの場所: /home/u62013505
ワークシート名:

「名前を付けてプロセスフローを保存」 アイコンをクリックし、 保存場所、ファイル名を指定して保存ボタン







主成分分析 - 実行方法 (1/2) ノードの設置

①左パネルより、[タスクとユーティリティ]→[タスク] →[多変量解析]→[**主成分分析**]を選択

②右側のプロセスフロー内のインポートノードの 右端の四角 (コントロールポート)の上へドラッグ&ドロップ

③プロセスフロー上に 主成分分析ノードが 生成されるのでダブルクリックして詳細設定画面を開く









主成分分析 - 実行方法(2/2)説明変数・オプション・出力

コード/結果 分割 53 設定 * 🛛 オプション データ 出力 • データ -WORK.IMPORT1 **マ**フィルタ: (なし) データソースの設定確認 役割 *分析変数: ↑ ↓ 💼 🕇 23 年齢 🔞 年間平均残高 😰 最終連絡日 説明変数の設定 数值型変数 😰 最終会話時間 🔞 CP中連絡回数 🐵 最終連絡日数 ▶追加役割

[データ]の設定(説明変数)

設定 コード/結果 分割 🗶 民 🔀
データ オプション 出力 情報
▼手法
成分の数: すべて 🗸
▶ 詳細
▼ プロット
表示するプロットの選択:
デフォルトおよび追加プロット
✓ 固有値と成分 (スクリープロット)
□ 成分ペアのスコア
□ 成分スコア行列
□ 成分パターンプロファイル
✓ 成分ペアのパターン
▼ オプション
スコアとパターンプロットの成分の数:
3 🗸
[デフォルトおよび追加プロット]を選択し、
[固有値と成分(スクリープロット)]と
[成分ペアのパターン]にチェックを入れる

[オプション]の設定(各種出力)

山刀」の設定(分析結果の二次	(利用)
設定 コード/結果 分割 💉 😡 🚼	
データ オプション 出力 情報	
▼出力データセット	-
成方のスコアテータセットを作成する *データセット名:	
work.Princomp_scores	参照
□ 統計量データセットを作成する	i
*データセット名:	
work.Princomp_stats [成分のスコアデータセットを作成す	参照 「る]
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングに	参照 「 る] 活用可能
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングに	参照 「 る] 活用可能
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングにお	参照 「 る] 活用可能
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングに対	参照 「 る] 活用可能
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングに	参照 「 る] 活用可能
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングに)	参照 「 る] 活用可能
work.Princomp_stats [成分のスコアデータセットを作成す にチェックを入れる (これにより主成分分析結果をクラスタリングにき	参照 「 る] 活用可能



主成分分析 - 実行結果(主成分分析の出力)

- ・ 主成分分析では、まずスクリープロットと累積寄与率のグラフから、最適な主成分数を検討する
 → 今回の分析では、固有値の値と簡単のため、主成分数=2とする
- 各主成分に対する変数寄与度から、各主成分軸の意味を検討する



主成分分析 - 実行結果(主成分分析結果の可視化)

・「**散布図**」ノードを活用して、各主成分軸をX軸、Y軸にとり、目的変数で色分け表示することで、主成分軸における各データポイントの位置付けと、目的変数との関係性が観察できる

▼散布図の設定

設定 コード/結果 分割 🖌 😡 🔀
表示 情報 ノード
 データ
WORK.PRINCOMP_SCORES
 マィルタ: (なデータソースを設定 ・ 、役割 ・ 役割
*X 軸: (1 項目)
❷ Prin1 X軸=主成分①
*Y 軸: (1 項目)
Prin2 Y軸=主成分②
グループ: (1 項目)
▲ 定期預金契約 グループ=目的変数
凡例の場所: 外側(デフォルト) 🗸
▶追加役割
I
I

▼散布図の出力結果





主成分分析結果のクラスタリング – 実行方法 (1/2) ノードの設置

①前回と同様に、左パネルより、[タスクとユーティリティ]→[タスク] →[クラスター分析]→[K-Means クラスタリング]を選択

②右側のプロセスフロー内の**主成分分析の後の** 右端の四角 □ (コ>トロールポート)の上へドラッグ&ドロップ

③プロセスフロー上に K-Means クラスタリングノードが 生成されるのでダブルクリックして詳細設定画面を開く



主成分分析結果のクラスタリング – 実行方法 (2/2) 説明変数・オプション

[オプション]の設定(各種出力)

[データ]の設定(説明変数)



設定 コード/結果 分割			
ノード データ	オプション	出力	情報
▼手法			
▼標準化			
標準化法:			
範囲 (デフォルト)			-
最小値を引き、	範囲で割ります		
▼ クラスタリング			
次の2つの手法のい	いずれかを指定する	必要がありま	:す:
☑ 最大クラスタ	7一数		
*クラスター:	3 🗘		
□ 候補シードと	2既存シード間の最小	小距離	
 合オブザベーシ ド 	ョンのクラスター重	心法をアッフ	fu-
🗌 データセットの	クラスター重心法を	読み込む	
🗌 最大反復回数	[最大クラス	マー数]にチェックカ
▼ 統計量	入っている	ことを確	認し、
表示する統計量: デフォルト統計量	[クラスター	-数]を	3 に設定

	・出力データセット	<u> </u>
	✔ クラスター割り当てデータセットを作成する	
	*データセット名:	
L	work.Fastclus_scores	1000
	□ 統計量データセットを作成する	
[2	フラスタテ語り当てテータセットを作用	成する]
с .	チェックを入れるits	1W
	🗌 クラスター重心法データセットを作成する	
	*データセット名:	
	work.Fastclus_seeds	IWe

[出力]の設定(分析結果の二次利用)

オプション

* 8 8

出力

分割

設定 コード/結果

データ



主成分分析結果のクラスタリング – 実行結果 (クラスタリングの可視化)

・「散布図」ノードを活用して、各主成分軸をX軸、Y軸にとり、クラスタ番号で色分け表示することで、主成分軸における各データポイントの位置付けと、各クラスタとの関係性が観察できる





▼散布図の出力結果

主成分①、②が大きい範囲では、契約者はほとんどいなかったが、 クラスタリングではこの傾向をある程度捉えた分類が行われている





Agenda

- ・クラスター分析の応用(他の分析手法との組み合わせ)
 - 主成分分析により説明変数を要約する
 - 主成分軸でクラスター分析を行う

・クラスター分析による分類(2):階層的クラスタリング

- 階層的クラスタリング(群平均法、重心法、Ward法)のしくみ
- 樹形図(デンドログラム)とクラスタ数の検討
- 都道府県データを用いて階層的クラスタリングにより類似地域を分析する
- ・今後のデータサイエンス学習に向けたスキルアップ
 - データサイエンティストに求められるスキル
 - SAS内サンプルデータの紹介と使い方
 - オープンデータの紹介





クラスタリング手法の種類



- ・ クラスタリング手法は、「非階層的」と「階層的」に大別される
- ・ 階層的クラスタリングはさらに 凝集型 と 分割型 があり、凝集型が用いられるのが一般的

手法の分類		手法	
非階層的クラスタリング クラスタ1	•k-means法(k平均法)	クラスタ内データの平均値をクラスタ重心として、 距離に基づき、事前に設定したクラスタ数k個に分割	SAS [*] Studio
変 数 B		混合ガウス法、超体積法など	第2回で説明
	似ている(≒距離の近い)データ/	クラスタ同士を逐次まとめる (ボトムアップアプローチ)	
1	 ・ウォード法 	クラスタ内のデータの平方和を最小にするように併合	SAS [®] Studio
階層的クラスタリング	▪ 最短距離法(最近隣法)	距離の近いデータから順番に併合	本日
変	凝 集 ■最長距離法(最遠隣法)	距離の遠いデータから順番に併合	ご説明
	型 ● 重心法	クラスタ重心からの距離に基づき併合	SAS [®] Studio
•	■群平均法	各クラスタ同士で全データの距離の平均を基準に併	合 SAS [®] Studio
 デンドログラム (dendrogram) 	■その他	メディアン法、可変法	
	分 似ていないデータ/クラスタ同士	を逐次分離させる(トップダウンアプローチ)	
	Diana法 Copyright © SAS Institute Int	c. All rights reserved.	25 Sas

代表的な階層的クラスタリング: 凝集型階層クラスタリング

- ・凝集型階層クラスタリングは、距離に応じて小さいクラスタを束ねて階層的に分類する手法
- クラスタ数は自動的に決定してくれる他、分類過程を可視化した樹形図(デンドログラム)も同時 に出力されるので、結果の解釈やクラスタ数の決定に役立つ



「近い」の評価尺度バリエーション

- ・ クラスタ間の「近さ」を測る指標には様々あるが、一概にどれが良いとは言えないため、複数試し て比較するのが一般的である。ただし、一般には、群平均法やWard法 (次頁) が頻用される
- ・ 最短距離/最長距離法は、計算量が少なくて済む反面、1点の影響を大きく受けやすい



Ward法の考え方

• Ward法*は最もよく用いられる手法であり、計算量は多いが、各データ点とクラスタ重心との関係性まで評価しているため、他手法に比べ、分類感度が高いとされる *米国の統計学者Joe H. Ward, Jr.が1963年に発表した論文にちなむ



「クラスタ重心」と、「当該クラスタ内の各データ」との距離の総和(二乗和)を クラスタごとに算出

クラスタAの場合

 $\mathbf{A} = a_1^2 + a_2^2 + a_3^2 + a_4^2 + a_5^2$

クラスタBの場合

 $\mathbf{B} = b_1^2 + b_2^2 + b_3^2 + b_4^2$

注目する2つのクラスタを結合した場合を仮定し、「結合後のクラスタ重心」と 「当該クラスタ内の各データ」との距離の総和(二乗和)を算出

 $\mathbf{AB} = a'_{1}{}^{2} + a'_{2}{}^{2} + a'_{3}{}^{2} + a'_{4}{}^{2} + a'_{5}{}^{2}$ $+ b'_{1}{}^{2} + b'_{2}{}^{2} + b'_{3}{}^{2} + b'_{4}{}^{2}$

1 と 2 の差、つまり、 AB - (A+B) が最小となるクラスター結合を採用 (結合前後でクラスタ内のばらつきに変化なし→統合してもOKと判定)

※近くにあり、ばらつきの小さいクラスタ同士が結合しやすい



Copyright © SAS Institute Inc. All rights reserved.



再揭

クラスタリング手法によって得意なデータパターンは異なり、様々な手法を試しながら、最適な手法を選択することが望ましい。中でも、k-meansは「重心からの距離」を用いて分類するため、円状のデータには強いが、楕円状や曲線状のデータは苦手

29





ビッグデータ分析の進め方

・データマイニングの進め方に関する方法論「CRISP-DM」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

CRISP-DM: データマイニング方法論



1.ビジネスの理解	・ビジネス、データマイニング目標の決定 ・プロジェクトの立ち上げ
2.データの理解	・データの収集 ・データの調査 ・データ品質の検証
3.データの準備	・データの選択や除外 ・データのクリーニング ・データの構築や統合
4.モデル作成	 ●モデリング手法の選択 ●モデルの作成 ●モデルの評価
5.評価	 ・データマイニングの結果の評価 ・プロセスの見直し ・実行可能なアクションリストの作成
6.展開/共有	•業務への導入計画 •モニタリング、メンテナンスの計画

30 **S**Sc



- 政府が公開する政府統計のオープンデータ "e-Stat" のデータを活用する
- 今回扱うデータの他にも、様々な統計データが公開さているので、企業内のデータと組み合わせることで、さらなる付加価値を生む可能性がある



Copyright © SAS Institute Inc. All rights reserved.



今回は、このうち、5年に1度実施している「全国消費実態調査」(現在の名称は「全国家計構造調査」)のデータを用いて、都道府県別の消費動向から、類似の都道府県をグルーピングすることを考える

政府統計名		全国家計構造調査(旧全国消費実態調査)				詳細
提供統計名		平成26年全国消費実態調査				
		全国				
提供分類2		家計収支に関する結果				
提供分類3		総世帯				
表番号		統議	調査年月	公開(更新)日	表示・ダウンロー	۴
フロー編						
42	平間収入階級	・中間収入十分位階級別1世帯当たり1か月間の収入と文出				
	総世帯		2014年	2015-12-16	L EXCEL → DB	
	勤労者世帯		2014年	2015-12-16	🛃 EXCEL 🛛 🌩 DB	
43 世帯主の年齢 総世帯・劃		階級別1世帯当たり1か月間の収入と支出				
		労者世帯	2014年	2015-12-16	🛓 EXCEL 🛛 🌩 DB	
44	住居の所有関係別1世帯当たり1か月間の収入と支出					
	総世帯・勤	労者世帯	2014年	2015-12-16	EXCEL DB	
45	資産の種類・	資産額階級別1世帯当たり1か月間の収入と支出(純資産)				
	総世帯		2014年	2016-03-25	EXCEL DB	
	勤労者世帯		2014年	2016-03-25	EXCEL DB	
	資産の種類・	資産額階級別1世帯当たり1か月間の収入と支出(総資産)				
	総世帯		2014年	2016-03-25	EXCEL → DB	
	勤労者世帯		2014年	2016-03-25	📩 EXCEL 🛛 🌩 DB	
地域編						
13	地域別1世帯	当たり1か月間の収入と支出				
	総世帯		2014年	2015-12-16	★ EXCEL → DB	
	勤労者世帯		2014年	2015-12-16	🛃 EXCEL 🛛 🔶 DB	

Source: https://www.e-stat.go.jp/stat-

search/files?page=1&layout=datalist&toukei=00200564&tstat=000001073908&cycle=0&tclass1=000001073965 &tclass2=000001074840&tclass3=000001077457&tclass4val=0



データの概要(加工前)

- e-Statより素データをダウンロードして開くと、開始行や開始列がずれていたり、空白行があったりと、加工が必要な形式であることがわかる
- ・ 今回は、本データから都道府県別の消費細目データ部分を抽出し、加工済のデータを用いる





データの概要(加工後)

都道府県

- e-Statより素データをダウンロードして開くと、開始行や開始列がずれていたり、空白行があったりと、加工が必要な形式であることがわかる
- 今回は、本データから都道府県別の消費細目データ部分を抽出し、加工済のデータを用いる

	#	都道府県	食料	住居	光熱・水道	家具·家事用品	被服及び履物	保健医療	交通・通信	教育	教養娯楽	読書	聴視•観覧	旅行	スポーツ	月謝	会費・つきあい費
• [1 1	北海道	54281	17491	19520	9089	10208	9234	35627	5630	23323	3894	2106	6506	2492	1430	1011
Τľ	2 7	青森県	55180	14357	22420	9162	8972	9936	33039	4880	16564	3256	2084	3839	1153	777	1079
ľ	3 1	岩手県	57514	14782	21267	8973	8288	10273	42912	5805	20278	3534	2428	6326	1447	1219	1458
ľ	4 7	宮城県	59052	16353	20331	9700	10640	10851	40742	7331	23394	3241	2322	7734	2172	1587	1168
ľ	5 ₹	秋田県	59146	12875	22394	9108	7467	9214	37645	4472	21037	3349	2403	6672	1617	979	1191
Ĩ	6 L	山形県	63042	13186	24030	11771	9486	10407	40365	6667	24577	3706	2227	5438	1659	1333	2230
Ĩ	7 7	福島県	57891	12579	20102	9633	8949	9613	44195	4614	21216	3637	1921	6701	1743	1432	1027
Ĩ	8 3	茨城県	62433	17292	20330	9186	10645	10668	44764	10113	26592	3638	2451	7254	2639	2041	1134
Ĩ	9 柞	栃木県	63866	18994	19997	9914	10252	12956	47208	6889	27561	3564	2723	7996	2908	2400	996
Ĩ	10 #	鮮馬県	61858	15629	18305	9605	9682	11270	43782	7539	25659	3767	2473	6800	2837	1965	931
~	11 İ	奇玉県	64632	20131	17747	8544	11403	11616	40152	12953	29055	3800	3115	8887	2699	2705	632
~	12 -	千葉県	66536	17887	18039	8820	11826	11859	39048	12165	30385	4186	3512	10407	2889	2800	684
~	13 3	東京都	68380	33295	16315	8691	12404	12151	33118	11060	32038	4180	3882	14361	3061	2964	808
~	14 7	神奈川県	67197	22708	16957	8783	11591	11443	38440	11004	31833	4159	3607	10790	2874	2894	889
~	15 ¥	新潟県	64400	15713	21881	8900	9077	10628	38983	6736	23878	3605	2535	7327	1955	1715	1032
1	16 2	富山県	67635	14518	21894	10624	9387	10776	51532	7879	27246	4003	4101	5869	2139	2433	1373
۲ ۲	17 7	石川県	66478	17678	18423	8733	9512	11135	42087	7993	27548	4220	3207	8516	1898	1966	1402
۲ ۲	18 7	福井県	67429	12168	20741	9034	10204	11287	45576	9585	27984	3482	4031	8052	1863	2080	1539
۲ ۲	19 L	山梨県	57641	17234	18209	7890	9429	10280	39392	9066	25849	3531	3699	6404	1877	2212	1132
~	20 1	長野県	62406	21145	21350	9866	9375	11987	42846	8047	27147	4035	2910	7265	2537	1830	1224
~	21 🛙	岐阜県	61939	12754	19952	9042	9942	10463	41580	704	25777	3627	2610	6754	2280	2056	1306
~	22 #	静岡県	62396	15048	18407	9012	9985	11488			28082	3714	2936	8355	2350	2259	1048
~	23	愛知県	64248	21485	17573	9010	11051	11880	=1	6 D D 7 D 7 W	28967	4041	3305	8547	2948	2655	818
~	24	三重県	63275	12856	19237	9036	11444	12889	—————————————————————————————————————	即经等	28462	3556	3578	8754	2773	2385	1173
~	25 %	茲賀県	63385	16479	18807	9587	10034	10818			26609	3280	2845	8377	1848	2052	1486
1	26 3	京都府	65337	13829	17928	8409	12630	9239	36645		26012	3984	3109	8260	1855	2174	985
~	27 2	大阪府	62386	18778	16292	7230	9898	10782	31046	10348	25016	3744	3264	7492	2465	2250	658
1	28 4	兵庫県	63620	19262	16725	8281	10712	10926	36040	9806	27000	3827	3023	8397	2546	2548	823
~	29 3	奈良県	66408	17630	19784	9875	11068	12405	42593	14481	27121	3849	3065	9467	2170	2609	986
~	30 1	和歌山県	58010	10696	17125	8152	9250	8326	36333	6001	23890	3376	2656	4995	2098	1787	888
~	31 /	鳥取県	58027	13626	18488	8143	9050	10320	41570	4966	24212	3198	3787	7077	1775	1865	904
~	32 8	島根県	59223	11926	19494	8915	8767	11814	40722	3866	23678	3446	3538	6335	1619	1512	1656
~	33 🕅	岡山県	58368	13776	18306	8286	9846	10347	38978	8451	24914	3052	2796	6682	2430	1789	931
ľ	34 /	広島県	58058	17721	17128	9180	9622	11195	38580	8773	24978	3308	2660	7944	1918	2004	997
ľ	35 L	山口県	55832	18576	16610	9381	8003	10961	35524	5193	23931	3557	2873	6234	1514	1798	840
Ĭ	36 1	徳島県	55896	16389	18015	8680	9656	10261	38507	6659	23923	3439	3762	7064	2039	1810	997
Ĩ	37 1	香川県	57352	15438	17319	8338	8754	11070	40876	6059	25565	3476	2814	6197	2689	2176	942
Ĩ	38	愛媛県	55531	13489	17201	8171	8284	9224	32679	7901	19353	2872	2450	5539	1732	1774	937
Ĩ	39 7	高知県	54971	14463	16479	7609	7510	10329	32613	5206	20184	3273	2447	5228	1722	1167	867
Ĩ	40 7	福岡県	54633	18999	16314	8029	9823	10405	36057	7360	24134	3073	2790	9478	2489	1809	784
	41 1	佐賀県	57104	13214	17556	8682	9647	11281	40406	6975	24864	3354	2813	6326	2333	1889	1189
	42]	長崎県	51798	18624	16853	7291	7934	10115	34480	6345	19631	2687	2421	7528	1383	1517	1013
1	43	熊本県	55006	11286	16802	8254	10041	11155	34633	6967	21046	2889	2309	5544	1688	2016	682
1	44 5	大分県	53558	14707	15685	8558	10853	11677	36458	3243	22105	3021	3223	5354	2139	1327	1153
1	45 3	宮崎県	53347	15963	15828	8228	8386	9476	36294	6276	21314	2565	3061	6208	2410	1413	1210
1	46 /	鹿児島県	50294	14792	15496	7800	7857	10022	39992	5063	18721	2593	2002	5533	2053	1160	1488
L	47 3	中縄県	48770	22616	17251	6750	5010	8088	28055	5169	16217	2500	1492	3913	1767	1700	970

予測(分析)対象を説明するための変数



34



- データの読み込み
- ・ 階層的クラスタリング (Ward法)
- •標準化したクラスタリング



新規プロセスフローの作成と保存

①左上メニューの <u>・</u>アイコンをクリックし、 [プロセスフロー] を選択

②新規のプロセスフローが作成されるので、 「名前を付けてプロセスフローを保存」アイコンをクリックし、 保存場所、ファイル名を指定して保存ボタン





データの読み込み (1/2)

① 左パネル内の 「アップロード」アイコン をクリック



②「ファイルの選択」ボタンをクリックし、ファイル選択画面で "Japan_expenditure_2014.csv"を選択し、OKボタン ③「アップロード」ボタンをクリック

ſ	ファイルのアップロード		×
	ファイルのアップロード先: /home/	u62013505	
	フ 選択済みファイル:	アイルの選択	
	CSV Japan_expenditure_201	4.csv 4.6 kb	C2
		アップロードキ・	ャンセル
l			
$\wedge +$			いファレムでない
4)/ _	ハイノレハリンノノノイノレノ サーバーファイルとフォルダ		
	は~ 竜 土 平 国 い	► 実行 □	
	⊿ 📲 odaws02-apse1-2		結果
	🔁 フォルダショートカット	+* 🗎 🖆	色 🔹 💼
	🛛 🔽 ファイル (ホーム)		
	sasuser.v94		

は~ 茴 土 平 国 い		100 L L L L L L L L L L L L L L L L L L
▲ 🚰 odaws02-apse1-2		和木
🔁 フォルダショートカット	+ - 🖆	色 💼
🔺 🔽 ファイル (ホーム)		
sasuser.v94		
₹ #1_ロジスティック回帰.cpf		
₶ #2_k-meansクラスタリング.cpf		
₹ #3_階層的クラスタリング.cpf		
bank_marketing.xlsx	_	
🕞 Japan_expenditure_2014.csv		
🚓 K-Means クラスタリング.ctk		

37

データの読み込み (2/2)

 ① 左パネル内の "Japan_expenditure_2014.csv" を選択し、右側のプログラムエリアにドラッグ & ドロップ



③詳細設定画面が開くので、実行ボタンをクリック (特に各設定は変更不要)

9	t #2_k-meansクラスタリング.cpf × t *#3_階層的クラスタリング.cpf × #3 階層的クラスタリング 〉 "Japan_expenditure_2014.csv" のインポート 酸定 コード/結果 分割 大 Q 緊
/	<u>オプション</u> *ファイル情報
7回帰.cpf スタリング.cpf りング.cpf lsx re_2014.csv リング.ctk	ソースファイル ファイル名: Japan_expenditure_2014.csv ソースの場所: /home/u62013505 行末の区切り記号: デフォルト ▼ 出力データ
	SAS Server: SASApp データセット名: IMPORT ライプラリ: WORK

②右側のプロセスフローにノードが生成されるので、 当該ノードをダブルクリック



④「結果」のタブ画面に読み込んだデータの概要が出力

P P &	8 K 8					
目次		· · · · · · · · · · · · · · · · · · ·				
		CONTENTS プロシジャ				
	データセット名	WORK.IMPORT	オブザベーション数	47		
	メンバータイプ	DATA	変数の数	17		
	エンジン	V9	インデックス数	0		
	作成日時	2022/10/01 22:53:29	オブザベーションのバッファ長	144		
更新日時		2022/10/01 22:53:29	削除済みオブザベーション数	0		
	保護		圧縮済み	NO		
データセットタイプ			ソート済み	NO		
	ラベル					
	データ表現	SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64				
	エンコード	utf-8 Unicode (UTF-8)				
		the second state of the second state of the second state of the				
	101070	エンシンパベスト関連情報				
ズ	131072					
データセットのペーシ	7数 1					
データページの先頭 1						
ページごとの最大OBS数 909						
先頭ページのOBS数 47						
データセットの修復要	k 0					
ファイル名 /saswork/SAS_work223C00006AE6_odaws01-apse1-2.oda.sas.com/SAS_workB87100006AE6_odaws01-apse1-						

38 5.2

Copyright © SAS Institute Inc. All rights reserved.



ビッグデータ分析の進め方

・データマイニングの進め方に関する方法論「CRISP-DM」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

CRISP-DM: データマイニング方法論







階層的クラスタリング (ward法) - 実行方法 (1/2) ノードの設置

①左パネルより、[タスクとユーティリティ]→[タスク] →[クラスター分析]→[オブザーベーションのクラスタリング] を選択

②右側のプロセスフロー内のインポートノードの 右端の四角 □ (コントロールポート)の上へドラッグ&ドロップ ③プロセスフロー上にオブザーベーションのクラスタリングノードが 生成されるのでダブルクリックして詳細設定画面を開く







階層的クラスタリング (ward法) - 実行方法 (2/2) 説明変数・オプション



[オプション]の設定(各種出力)





参考: 変数の尺度(名義尺度・順序尺度・間隔尺度・比例尺度)

• 変数の種類は大きく「質的データ」と「量的データ」に分けられ、それぞれの特性に合わせて扱う必要がある

種類	変数の尺度	概要	データの例	扱い方
				大小 差分 比率
質的データ	名義尺度	単にデータを区別するための分類ラベル。 演算不可で、順序も意味をなさない	 ●性別、血液型、顧客ID ●作業者、個品ID、 良品/不良品 	(A <b) (a="" (a-b)="" b)<br=""> べま計によるカウントのみ可能</b)>
(カテゴリーデータ)	順序尺度	順序(大小関係)にのみ意味がある尺度。 したがって、平均値は意味を持たないが、順 序統計量(最大・最小など)は算出可能	■顧客満足度、震度 ■不良レベル、工程順序	•
量的データ	間隔尺度	数値演算可能だが、 値の差 のみに意味が ある尺度。 0はあくまで相対的な位置関係でしかない	■年齡、西暦、偏差値 ■温度(℃)、製造日時	• • -
(数量データ)	比例尺度	数値演算可能で、値の差に加え、 値の比 にも意味がある尺度。 0が「何もない」という絶対的な意味を持つ	■身長、売上金額 ■寸法、圧力、作業時間、 絶対温度	• • •



クラスタリングにおける変数スケールの影響と標準化

k-means法などの「距離」に基づくによるクラスタリング手法は、データの「スケール」に大きく影響を受ける。このため、必要に応じて、「標準化」の処理を行なった上でクラスタリングを行う必要がある



(参考)回帰分析における標準化の有効性

再揭

- ・機械学習では、各変数間でスケール (値範囲) が大きく異なると、計算に時間がかかったり、 回帰係数などのパラメータの直接比較が困難になるため、スケールを揃えることが有効
- ・特に、各変数を平均0,分散1に変換する「標準化」を用いることが多い







階層的クラスタリング (Ward法) – 実行結果

- Ward法の結果、47都道府県が階層的にクラスタリングされた。Ward法では、全体的にバランスよく、分類が行われていることが確認できる
- ・ 任意の場所で区切ることで、最適なクラスタ数の検討が可能





Agenda

- ・クラスター分析の応用(他の分析手法との組み合わせ)
 - 主成分分析により説明変数を要約する
 - 主成分軸でクラスター分析を行う
- ・クラスター分析による分類(2):階層的クラスタリング
 - 階層的クラスタリング(群平均法、重心法、Ward法)のしくみ
 - 樹形図(デンドログラム)とクラスタ数の検討
 - 都道府県データを用いて階層的クラスタリングにより類似地域を分析する

・今後のデータサイエンス学習に向けたスキルアップ

- データサイエンティストに求められるスキル
- SAS内サンプルデータの紹介と使い方
- オープンデータの紹介





- データサイエンティスト協会が定義するデータサイエンティスト:
 「データサイエンティストとは、データサイエンスカ、データエンジニアリングカをベースにデータから 価値を創出し、ビジネス課題に答えを出すプロフェッショナル」
- ・これら3スキルはどれも不可欠で、分析フェーズによって中心となるスキルが変化する、としている



データサイエンティストの育成モデルの例

- 理想的には、すべてのスキルを持つ人材を育成できればよいが、元々のバックグラウンドや 経験値を生かした育成を行うのが現実的
- チームでスキルを補い合い、プロジェクトフェーズに応じて、役割分担や協業体制が必要



Copyright © SAS Institute Inc. All rights reserved.

参考:データサイエンティスト協会が定義するスキルレベル

バックグランドベース育成モデルを実現する上でも、各メンバーは、各スキルについて最低限のレベルを保有する必要がある。データサイエンティスト協会が定義する「独り立ちレベル」がまず目標

<u>レベル定義</u>	コンサル力 (ビジネス力)	統計力(データサイエンス力)	【T力(データエンジニアリング力)
業界代表	 組織や市場全体にインパウトを出せる 対象とする事業全体、産業領域における課題の切り分け、テーマ、論点の明確化ができる 	 新しいアルゴリズムや分析手法の開発ができる 複数のパラメータやアルゴリズムの選択など、適切な分析アプローチの設定ができる 	 複数のデータソースを統合したデータシステム、 もしくはデータプロダクトの構築、全体最適化が できる
棟梁 全てのスキル	 分析を通じオペレーション上の革新が実現できる 仮説や可視化された問題がない中で、適切に問題を定義し、解き、価値を見出すことができる 特定の課題領域において、課題と取組のテーマを構造的に整理し、見極めるべき論点をクリアにできる 組織全体を見渡して、必要なデータの当たりをでて強い立ちもベル」以上の習得が理想 	 多変量解析の概念を理解し、活用することができる 機械学習、自然言語、画像処理のアルゴリズムを理解し、適切に活用、問題解決することができる モデルを構築できる 	 分析に必要なデータフォーマナ、取得蓄積仕様等を設計できる 問題設定に応じた新規データマート設計ができる 構造化データノ非構造化データを問わず、分析システムを設計できる 構築したモデルを実装できる データ分析を作ったシステムを自身で構築できる
独り立ち	 仮説や既知の問題が与えられた中で、最適解 最大解を見出すことができる 扱っている課題領域で新規の課題を切り分け、 構造化できる 当該プロジェクト・サービスを超えて、必要なデータの当たりをつけることができる 	 SPSS/SAS/R等が使える。指示されなくてもサンプル抽出ができるとともに内容を確認できる データクレンジング、分布、単回帰やP値の概 念を理解し、活用することができる 	• 大規模のファイルや、データベースにアクセスし 大量の構造化データを処理することができる
見習い	 ビジネスにおける論理とデータの重要性を認識している 仮説や既知の問題が与えられた中で、必要な データに当たりをつけて、データを用いて改善することができる 扱っている課題領域における基本的な課題の 枠組みが理解できる 	 基本統計量(平均、中央値など)の知識を 有し、指示されればデータの抽出、グラフ作成 を正しく行うことができる 	 一般的なアクセス解析システムを使うことができる 抽出されたデータサブセットに対し、Excelや Access等の統合環境を用い、目的に応じた 処理をすることができる
 未経験 ● 	 ビジネスは勘と経験だけで回すものと思っている 課題を解決する際に、そもそも定量化する意識 がない 	 基本統計量の意味を正しく理解していない 指数を指数で割り算したりする 「平均年収」をそのまま鵜呑みする グラフ・チャートの使い方が不適切 	 レポートされてくる数値サマリに目は通すが、 特に記憶には残らない アクセス解析システムを使っていない ExcelやAccessは数字しか入れない

SASHELPデータ (サンプルデータ) の活用

SAS Studioでは、デフォルトで様々なサンプルデータ(SASHELP Data Sets)が格納されており、分析のトレーニングなどに有効活用できる

・左パネルより、 [ライブラリ]→[マイライブラリ]→[SASHELP]を選択

C.

影

SAS	5° Studio	
۱.t	ナーバーファイルとフォルダ	
• 5	マスクとユーティリティ	
• 2	スニペット	
	ライブラロ	
é	5) 67 💼 🗏 ()	
4	🔐 マイライブラリ	
	🖻 靜 MAPS	
	🖻 靜 MAPSGFK	
	🛚 🔐 MAPSSAS	
	🛚 🗃 SASDATA	
	🔺 🔐 SASHELP	
	▷ 📰 _CMPIDX_	
	Þ 📰 AACOMP	
	Þ 🌇 AARFM	
	🖻 🌇 ADSMSG	
	🖻 🌇 AFMSG	
	▷ 📰 AIR	
	▷ 📰 AIRLINE	
	Þ 📰 APPLIANC	
	ASSCMGR	
	🖻 🌇 BASEBALL	
	Þ 🔡 BEI	
	BIRTHWGT	

▼データセットの例 1986年のメジャーリーガーの成績データ BASEBALL 年齢とBMIに関するデータ BMIMEN 骨髄移植患者の生存期間データ **BMT** 2003年の乳児死亡率に関するデータ BIRTHWGT **FATI URF** 機械の不具合に関するデータ DEMOGRAPHICS 各国の人口などに関するデータ 迷惑メールデータ JUNKMATI ORSALES 売上に関するデータ

データセットの一覧と詳細は、下記リンクを参照(英文) https://support.sas.com/documentation/tools/sashelpug.pdf



オープンデータの候補リスト (2022年4月現在)

★…オススメもしくは、よく使われている

カテゴリ	#	サイト名	概要	データ数	マーケ	製造	医療	URL
政府系	1	政府統計e-stat	各府省の公表データを1つにまとめたサイト。 また、これを分析-readyな形式に加工して扱いやす くした、「教育用標準データセット (SSDSE)」 というサ イトもある。	約700件		※統計 調査のみ		https://www.e-stat.go.jp/ ▼教育用標準データセットssDSE https://www.nstac.go.jp/ use/literacy/ssdse/
	2	データカタログサイト	二次利用が可能な公共データの横断的検索が可能 なデータカタログサイト。ただし、PDFやHTMLなどの未 整形なデータがほとんど。	約2.5万件				https://www.data.go.jp/
	3	観光統計データ	日本政府観光局が運営する日本の観光統計データ。	不明				https://statistics.jnto.go.jp/
大学系	7 4	UCI Machine Learning Repository	米カルフォルニア大学アーバイン校による公開データ セットで、非常に有名で、利用者が多い。	約600件				http://archive.ics.uci.edu/ml/i ndex.php
	5	Harvard Dataverse	米ハーバード大学による公開データセット。主に論文 で公開された、様々な分野のデータが揃っている。	約15万件				https://dataverse.harvard.ed u/
民間系	6	AWS パブリックデータセット	米Amazon Web Service社による公開データセッ ト。画像、ゲノム、テキストなど非構造化データが多い。	約300件				https://registry.opendata.aws /
	7	Tableau Public Sample Data Set	米Tableau社による公開データセット。 アメリカ国内の公共系データが多い。	約30件				https://public.tableau.com/s/ <u>resources</u>
	8	Google Dataset Search	米Google社が提供している、データセットの検索エン ジン。一部、日本語でも検索可能。 ただし独自データではなく、単なるWeb上の寄せ集め。	不明		•		https://datasetsearch.researc h.google.com/ ▼使い方を解説しているサイト https://atmarkit.itmedia.co.jp/ait/ articles/2007/15/news021.html
	9	Microsoft Research Open Data	米Microsoft社が提供しているデータセットで、同社 の研究部門が研究に用いたデータを公開。 テキスト、画像系などの非構造化データが多い。	約100件				https://msropendata.com/
	10	Yahoo Webscope Datasets	米Yahoo社が提供しているデータセットで、同社のサ イトなどで収集されたマーケティングデータが中心。 ※ただし、非営利団体の研究目的でのみ利用可能	約70件				https://webscope.sandbox.ya hoo.com/
	11	日経平均プロファイル	日本経済新聞が公開するデータセットで、日経平均、 日経アジア指数などが利用可能。	約30件				https://indexes.nikkei.co.jp/n kave/index?type=download



オープンデータの候補リスト (2022年4月現在)

★…オススメもしくは、よく使われている

カテゴリ	#	サイト名	概要	データ数	マーケ	製造	医療	URL
分析 コンペ サイト	12	Kaggle	言わずと知れた、米国の分析コンペティションサイト。 コンペ用の様々なデータセットが公開されており、中に は民間企業から提供受けたリアルなデータもある。 (コンペ終了後、非公開にされているデータもあり) 分析コードも公開されているため、自学習におすすめ。	不明		•		https://www.kaggle.com/
	13	SIGNATE	日本版Kaggleともいうべき、日本の分析コンペティ ションサイト。企業から提供を受けたリアルなデータが 大半。 ※ただし、基本的にはコンペ目的以外での利用を禁 止しているため、要注意。	不明				<u>https://signate.jp</u>
学術 研究 向け	14	NDBオープンデータ	厚労省が提供する、匿名化したレセプト情報・特定 健診等情報に関するデータセット。	約300件				https://www.mhlw.go.jp/stf/s eisakunitsuite/bunya/000017 7182.html
	15	国立情報学研究データ リポジトリ	国立情報学研究所(NII)が公開するデータセット。 民間企業や大学等から提供を受けたリアルデータ (楽天の購買データ、アットホームの不動産データ、 Yahoo!知恵袋データ、など)が公開されている。 ※ただし、非営利団体の研究目的でのみ利用可能	約20件				https://www.nii.ac.jp/dsc/idr/ datalist.html
その他	16	Python scikit-learn 内のデータセット	Pythonのscikit-learnライブラリに付属しているデー タセット。ボストンの住宅価格や、アヤメ品種、糖尿病 患者や乳がんのデータなど。 ▼参考:Pythonでの読込み例 from sklearn.datasets import load_boston boston = load_boston() df = pd.DataFrame(boston.data,columns=boston.feature_names)	14件	•		•	▼わかりやすくまとめているサイト https://zenn.dev/nekoallergy/ articles/scikit-learn-datasets
(付属データ/ 個人サイト等)	17	松原望先生(東大名誉 教授)の個人サイト	松原先生が公開する様々なデータ(Webから収集?)。 少し古いデータが多いものの、Excelによる分析手法 とともに掲載されているので、自学習用に向いている。	約90件				https://www.bayesco.org/top /datasite
	18	データで学ぶ!統計活用 授業のための教材サイト	統計教育推進委員会が公開する様々なデータ。や や古めで、項目数もそれほど多くないため、手元での 簡単な分析学習に使うイメージ。	26件				<u>https://estat.sci.kagoshima-</u> u.ac.jp/data/



まとめ

 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •
 •

- 主成分分析とクラスター分析の組み合わせ
 - 主成分分析を行うことで、多数の説明変数を要約することができた
 - 主成分分析結果に対してクラスター分析を適用することで、別観点の知見を得ることができた
- ・ 階層的クラスタリングによるデータ分類
 - 階層的クラスタリング(群平均法、重心法、Ward法)のしくみについて学習した
 - 各手法を都道府県データに適用し、類似の都道府県をグルーピングすることができた
 - デンドログラムを観察することで、最適なクラスタ数を検討することができた
- ・今後のデータサイエンス学習に向けて
 - 実践的なスキルを鍛えるには、座学だけでなく、実際のデータを触ってみることが一番
 - 様々なオープンデータを活用して、スキルアップを目指す
 - データサイエンスの知識だけでなく、「ビジネスカ」「ITカ」も極めて重要

End of File



