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

#1 ロジスティック回帰分析

2022年10月12日



Copyright © SAS Institute Inc. All rights reserved.

### Agenda

- ・ グラフによるデータ観察
  - 散布図、散布図行列の活用方法
  - 箱ひげ図の活用方法
- ・ ロジスティック回帰分析
  - 教師あり学習と分類手法
  - ロジスティック回帰の基本とその他手法との違い
  - マーケティングデータを用いてロジスティック回帰により契約確度の高い顧客を予測する
  - 回帰係数により契約の影響因子を探索する

ビッグデータとは



- ビッグデータとは、「4V」\* の特徴を持つデータ \*Veracity を除いて 3V と呼ばれることもある
- つまり、「大量かつ、様々な種類のデータが、次々と蓄積されていること」





Copyright © SAS Institute Inc. All rights reserved.

### ビッグデータからの新知識獲得=「ビッグデータ分析」

- ビッグデータ分析は、膨大なデータ(=事実、結果)に潜むデータ間の関係性や特徴(=新知識)をデータドリブンで抽出する試みである
- ・データから「一見関連のなさそうな新たな洞察を得る」手法と言える



マイニング = 採掘 と言う意味合いから、「データマイニング」 とも呼ばれる (mining)







CRISP-DM: データマイニング方法論





### 参考: SASの方法論 "SAS Analytics Lifecycle"

• SASの分析フレームワーク "SAS Analytics Lifecycle" においても、大きな流れは同じであり、 データ準備の重要性や、反復的に進めることの必要性を説いている

SAS Analytics Lifecycle





Copyright © SAS Institute Inc. All rights reserved

### 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





Copyright © SAS Institute Inc. All rights reserved.

## 教師あり学習のイメージ(数値予測と分類)

#### ・各顧客レコードに対して数値 or カテゴリー値(クラス)の解答を与え関係性を学習



参考:説明変数と目的変数

- ・原因となる変数のことを「説明変数」、その原因を受けて変動する結果系の変数のことを「目的 変数」と呼ぶ
- 書籍や講師、ツールなどによって呼称が異なる場合もあるため、慣れておく必要がある





### ロジスティック回帰分析の基本

・ロジスティック回帰は、回帰させる関数にS字カーブ状のロジット関数を用いることで、 ニクラスの分類(購入する/しない、良品/不良品、生/死など)を可能とする回帰



### 参考:重回帰モデル

- ・目的変数Yとその影響因子である説明変数Xとで関数 (Y=f(X)) に当てはめることを回帰という
- ・特にf(X)が線形関数の場合、線形回帰と呼び、1変数では単回帰、多変数では重回帰という



### 参考:教師あり分類手法における決定境界の比較

- 分類手法は、データの空間内で、いかに異なるクラスを綺麗に分割できるかで性能が決まる。
   この分割の境目を「決定境界」と呼ぶ
- ・ **ロジスティック回帰は、直線的な決定境界**を描くため、非線形に入り組んだ分類には向かない





### 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





- 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 smallest datasets are provided to test more computationally demanding machine learning algorithms (e.g., SVM).

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: (ユー	均残高 -□)			最終連續 会話時間	洛時の  (秒)	キャンペーン 連絡回	v中の 最終 数 糸	≷連絡からの 圣過日数	キャンペーン 連絡回	ン前の 前回= 数 の	ドャンペーン )結果
年前	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	≣¥8	日亦数	14	may	341	1	330	2	other 😑	的恋类
	39	technician	married	secondary	no	147	yes		JESA	6	may	151	2	-1	0	unkno	- 1 36 3
	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	······				101			11 1			100	2	4	^	1	
	56	technician	married	secondary	no	4073	no	no	cellular	27	aug	239	5	-1	0	unknow	BII ? 4
	37	admin.	single	tertiary	no	2317	ye	鶪(分	術)対象	20	apr	114	1	152	2	failure	no / J
	25	blue-collar	single	primary	no	-221	朝日	INTA	thme	23	may	250	1	-1	0	unknow	い対
	31	services	married	secondary	no	132	no	no	cellular	7	jul	148	1	152	1	other	no
	38	management	divorced	unknown	no	0	yes	no	cellular	18	nov	96	2	-1	0	unknowp	10
	42	management	divorced	tertiary	no	16	no	no	cellular	19	nov	140	3	-1	0	unknow	52
	44	services	single	secondary	no	106	C <b>∂<sub>l</sub>9</b> /right	t 🛛 🛛 AS I	nsunknown	rights re <b>1</b> 2	jun	109	2	-1	0	unknown	no
	4.4		the second second			0.2				-	4.1	105	2	4	0		





- 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)

		8	SAS <sup>®</sup> Studio	🕗 🗁 🤀 ビジュアルプログラマ - 🖨 ? サインアウ
A Dury Construction () / 1 Construction () / 1 C	None Label Type Long Store & de Canador II 201		<ul> <li>サーバーファイルとフォルダ</li> <li>● 面 上 平 目 い</li> <li>● godaws01-apse1</li> <li>□ フォルダショートカット</li> <li>□ ファイル(ホーム)</li> <li>● my_content</li> <li>● my_shared_file_links</li> <li>● moguchi0</li> <li>▶ ■ sasuser v94</li> </ul>	

### 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 <sup>®</sup> Studio	
<ul> <li>サーバーファイルとフォルダ</li> </ul>	8
は→ 竜 圭 平 目 い	_
⊿ 🚰 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_ma	arketing.xlsx			
ソースの場所: /home/u	62013505			
ワークシート名・				
第1日-カシート				
pp 1 2 - 2 2 - 1-				
7-8 04	25 M (15)	<u>=_</u>		
	和朱 1/J:	7-9		
╔ Ҏ ╔ ╩   ≞   ┍	7 📾			
<ul> <li>目次</li> </ul>		· · · · · · · · · · · · · · · · · · ·		
		CONTENTS 711997		
	データセット名	WORK.IMPORT1	オブザベーション数	4521
	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	2022/08/08 09:34:47	オブザベーションのバッファ長	120
	更新日時	2022/08/08 09:34:47	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)		
		エンジン/ホスト間連情報		
データセットのページ	ジサイズ 131072			
データセットのページ	ノ数 5			
データページの先頭	1			
	SIX 1090			
ページごとの最大OB				





### 読み込んだデータの確認

#### データ概要の確認

新し		うご	フザタブ	で 厚	く						
[a ,	コート 日次 日次	2   2	а а								
						CONT	ENTS ZON	Jur	1		
		データセット名 WORKIMPORT1 オブザペーション数 4521									
					CON	TENTS プロ	シジャ				
	データt	マット名	WORK.IMPORT1					オブザベーション数	4521		
	メンバ-	タイプ	DATA					変数の数	17		
	エンジン	/	V9					107007			
	作成日期	ŧ	2022/08/08 09:34:4	7				オブザベーションのパップ	7 7 長 120		
	更新日時	ŧ	2022/08/08 09:34:4	7				911272とイエ27	を確認し		
	保護							二龍河の	NU		
	データt	マットタイ	プ					(ビッグデータケ	)析の基本)		
	ラベル							(_))) //.	, 1, 1 • 2 · 1 · 1 · 1		
	データま	表現	SOLARIS_X86_64,	LINUX_	K86_64	, ALPHA_T	RU64, LINU	IX_IA64			
	エンコー	-ド	utf-8 Unicode (UTF	-8)							
					T 1/3	シッカフトリ	語と南本語書記				
データセットのペーシ	バサイズ	131072				200201	SAE IN TA				
データセットのペーシ	烫	5									
データページの先頭		1									
ページごとの最大OB	S数	1090	)90								
先頭ページのOBS数		1052									
データセットの修復数	\$	0									
ファイル名		/saswork/S	AS_work71F80001F3F	A_odaws	01-aps	e1-2.oda.sa	as.com/SAS	_workC7860001F3FA_odaws01-ap	se1-2.oda.sas.com/import1.sas7bdat		
作成したリリース		9.0401M6									
作成したホスト		Linux									
ノード番号		33850									
アクセス権限		rw-rr									
所有者名		u6201350	5								
ファイルサイス	k7 -	768KB		au 4	- 71	t==					
7712512 (111	谷(	710	)ナータ!	ギャ	て们	隹訟					
				恋勤	と尾性	リスト(アル	ファベット	(E)			
		#	変数	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	ク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	住宅ローンの有無	文字	3	\$3.	\$3.	住宅ローンの有無			
		8	個人ローンの有無	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	学歴	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	年間平均残高(ユーロ	数值	8	BEST.		年間平均残高(ユーロ)			
		1	年齢	数值	8	BEST.		年齢			
		14	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	最終連絡月	文字	3	\$3.	\$3.	最終連絡月			
		3	結婚歷	文字	8	\$8.	\$8.	結婚歷			
		2	職業	文字	13	\$13.	\$13.	職業			
		9	連絡手段	文字	9	\$9.	\$9.	連絡手段			

#### 🔎 🖻 🤀 SAS プログラマ・ 🖨 ? サインアウト 😰 プログラム1 × 🗊 \*bank\_marketing × 設定 コード/結果 分割 🏒 🔒 😡 🚼 目ログ 選コード ▼ 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: 第1ワークシート 出力データ SAS Server: SASApp データセット名: IMPORT1 WORK ライブラリ: 「出力データ」画面より、 変更 オプション 取り込んだ生のデータを確認 ファイルの種類: デフォルト (ファイル拡張子に基: 出力データ コード ログ 結果 ビュー: 列名 ・ 🗈 💄 😘 🔳 🍞 フィルタ: (なし) テーブル: WORK.IMPORT1 -合計行数: 4521 合計列数: 17 列 0 年齢 職業 結婚歴 学歴 クレジット... 年間平均残高 (ユーロ 住宅口 ✓ すべて選択 1787 no 30 unemployed married primary no 🜌 🔞 年齢 4789 yes 33 services married secondary no ☑ 🛕 職業 35 management 1350 yes single tertiary no ☑ 🛕 結婚歴 1476 yes 30 management married tertiary no 🜌 🛕 学歴 59 blue-collar 0 yes secondary no married 🜌 🛕 クレジットカード債務 747 no 35 management single tertiary no 🜌 🛛 毎間平均残高(ユーロ 307 yes 36 self-employed married tertiary no 🗹 🛕 住宅ローンの有無 39 technician 147 yes married secondary no 🜌 🛕 個人ローンの有無 41 entrepreneur married tertiary no 221 yes 🗹 🛕 連絡手段 10 43 services married -88 yes primary no 11 プロバティ 39 services secondary no 9374 yes 値 married 12 43 admin. 264 yes ラベル married secondary no 名前 13 36 technician married tertiary no 1109 no 14 20 student 502 no 長さ single secondary no 15 31 blue-collar secondary no 360 ves 種類 married (i) メッセージ: 4 ユーザー: u62013505

生データの確認





### 作成したプロセスフローの保存



#### 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

★プロセスフロー1 ×
- プロセスフロー1> "bank_marketing.xlsx" のインポー
設定 コード/結果 分割 🖌 😡 🎦
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: <b>/home/u62013505</b>

#### 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





データの特徴の捉え方

・ビッグデータでは個々のデータをくまなく見るのは難しいため、グラフ(ヒストグラムや散布図)や 要約統計量(平均値や標準偏差)を用いて全体傾向を把握する







ビッグデータでは個々のデータをくまなく見るのは難しいため、グラフ(ヒストグラムや散布図)や
 要約統計量(平均値や標準偏差)を用いて全体傾向を把握する



### ヒストグラム(度数分布図)

 ヒストグラムは、データを一定間隔(階級)ごとに頻度集計(度数)をとり、 横軸に階級、縦軸に度数をとって柱状(ビン)の集合で表したグラフ

• データの分布状況を可視化して、直感的にデータの特徴を捉えやすくする目的がある

元データ

レストランAの レビューアーごとの評点データ

#	評点	7	#	評点
1	3.5	1	1	3.0
2	3.5	1	2	3.5
3	4.5	1	3	3.0
4	3.5	1	4	4.0
5	3.5	1	5	4.0
6	2.5	1	6	4.0
7	3.0	1	7	2.5
8	3.5	1	8	4.5
9	4.0	1	9	1.0
10	3.0	2	0	5.0

一定間隔(階級)ごとに レビュー件数(度数)を集計

度数分布表

階級	度数
0.0~0.5	0
0.5 <b>~1.0</b>	0
1.0 <b>~1.5</b>	1
1.5 <b>~2.0</b>	0
2.0 <b>~2.5</b>	2
2.5 <b>~3.0</b>	4
3.0 <b>~3.5</b>	6
3.5 <b>~4.0</b>	4
4.0~4.5	2
4.5 <b>~5.0</b>	1

※境界値は小さい側の階級に含める

ヒストグラム(度数分布図)

横軸に階級、縦軸に度数を取り、 柱(ビン)の集合で可視化









 分布形状の把握は極めて重要であるが、主にヒストグラムは1変数、散布図は2変数での分布 可視化に適している。一方、箱ひげ図は、分布の概形を多変数間で比較するのに適している



26 **Sas** 

散布図と相関

- 二つの値(変数)間において、一方が上がれば他方も上がる(or下がる)ような関係性のことを「相関」と呼ぶ
- 各変数を各軸にとってグラフ化した「散布図」を描くことで、相関関係の視覚的な把握が可能



Copyright © SAS Institute Inc. All rights reserved







28 **S.Sas** 



ビッグデータでは全データをくまなく見るのは難しいため、「要約した値(=要約統計量)」を用いて全体傾向を把握するのが一般的



#### 相関行列 ※第2回で取り扱う予定

- 事前に各変数間の相関係数を総当たりで調べておくと、後々の結果解釈に役立つ(相関行列)
- ・また、共線性が高い変数 (相関の高い) が複数混ざっていると、その変数の影響を強く受け、 偏った分析結果になることがある。この場合、共線性が高い変数は除外することが有効







- 要約統計量
- ヒストグラム/箱ひげ図
- 散布図/層別散布図



Copyright © SAS Institute Inc. All rights reserved.



要約統計量/ヒストグラム/箱ひげ図 – 実行方法 (1/3)

## ①左パネルより、[タスクとユーティリティ]→[タスク] →[統計量]→[要約統計量]を選択

#### 

#### SAS<sup>®</sup> Studio ₱■ \*#1\_ロジスティック回帰.cpf × サーバーファイルとフォルダ ▶ 実行 🛛 🛃 🔣 🛛 🛣 コードの生成 ▼ タスクとユーティリティ フロー 結果 プロパティ 维→ 前 民 目 い 色 🗸 + - 🖌 🖕 侖 💻 🖌 🖳 👧 マイタスク 💷 タスク コントロールポート **⊐** "bank\_marketi ng.xlsx" のイ ンポート ▶ 💷 データ ▶ 👥 グラフ - 👥 マップ 🗈 統計量 🔡 データ探索 要約統計量 🕅 分布分析 🔪 Ⅲ 一元度数表 ドラッグ&ドロップ |∕ 相関分析 ☑ 分割表分析 ┣ t 検定 ▶ 💶 線形モデル

#### ③プロセスフロー上に要約統計量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



32 Sas



要約統計量/ヒストグラム/箱ひげ図 – 実行方法 (2/3)



Copyright © SAS Institute Inc. All rights reserved.



### 要約統計量/ヒストグラム/箱ひげ図 – 実行方法 (3/3)

#### ④[オプション]の画面で、出力の設定



#### ⑤実行ボタンをクリックすると、結果が出力される





#### 要約統計量/ヒストグラム/箱ひげ図 – 実行結果 (要約統計量)

定期預金契約	Obs 数	変数	ラベル	平均	標準偏差	最小値	最大値	中央値	Ν	欠損値の数
no	4000	年齢	年齢	40.9980000	10.1883977	19.000000	86.0000000	39.0000000	4000	0
		年間平均残高	年間平均残高	1403.21	3075.35	-3313.00	71188.00	419.5000000	4000	0
		最終連絡日	最終連絡日	15.9487500	8.2497356	1.0000000	31.0000000	16.0000000	4000	0
		最終会話時間	最終会話時間	226.3475000	210.3136306	4.0000000	3025.00	167.0000000	4000	0
		CP中連絡回数	CP中連絡回数	2.8622500	3.2126088	1.0000000	50.0000000	2.0000000	4000	0
	1	最終連絡日数	最終連絡日数	36.0060000	96.2976572	-1.0000000	871.0000000	-1.0000000	4000	0
		CP前連絡回数	CP前連絡回数	0.4712500	1.6273707	0	25.0000000	0	4000	0
yes	21	年齢	年齢	42.4913628	13.1157723	19.000000	87.0000000	40.000000	521	0
	È.	年間平均残高	年間平均残高	1571.96	2444.40	-1206.00	26965.00	710.0000000	521	0
	$\sim 10^{-1}$	最終連絡日	最終連絡日	15.6583493	8.2351482	1.0000000	31.0000000	15.0000000	521	0
		最終会話時間	最終会話時間	552.7428023	390.3258046	30.0000000	2769.00	442.0000000	521	0
	1 N	CP中連絡回数	CP中連絡回数	2.2667946	2.0920709	1.0000000	24.0000000	2.0000000	521	0
	` <u>`</u>	最終連絡日数	最終連絡日数	68.6391555	121.9630631	-1.0000000	804.0000000	-1.0000000	521	0
		CP前連絡回数	CP前連絡回数	1.0902111	2.0553682	0	14.0000000	0	521	0

#### 目的変数である「定期預金の契約」の有無で、傾向が異なりそうな変数





#### 要約統計量/ヒストグラム/箱ひげ図 – 実行結果 (ヒストグラム/箱ひげ図)

#### ※定期預金の契約有無で差が大きい変数







#### 参考:要約統計量/ヒストグラム/箱ひげ図-実行結果(ヒストグラム/箱ひげ図)

※定期預金の契約有無であまり差が認められない変数





Copyright © SAS Institute Inc. All rights reserved



### 参考:ヒストグラム・箱ひげ図のみ描画したい場合







### 散布図の描画 - 実行方法 (1/2)

## ①左パネルより、[タスクとユーティリティ]→[タスク] →[グラフ]→[**散布図**]を選択

#### 



#### ③プロセスフロー上に散布図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







### 散布図の描画 – 実行方法 (2/2)

#### ④[オプション]の画面で、X軸、Y軸の設定をし、実行ボタン

サーバーファイルとフォルダ	⊗ *#1_ロジスティック回帰.cpf ×	
▼ タスクとユーティリティ	#1 ロジスティック回帰 実行ボタン	
维→ 亩 民 目 55	設定 コード/結果 分割 🔀 🔀	
■ マイタスク	データ 表示 情報 ノード	<u>コード</u> ログ 結
▲ <u>■</u> タスク	<b>・</b> データ	👩   💽   🚢   🐂   行番号
▶ 🛄 データ	WORK.IMPORT1	1 /*
<b>▲ 🛄</b> グラフ	マングングリクランクション	2 * 3 * SAS Studio 3.8
<ul> <li>■ 棒グラフ</li> <li>● 棒-折れ線グラフ</li> <li>● 箱ひげ図</li> <li>● バブルプロット</li> <li>■ ヒートマップ</li> <li>□ ヒストグラム</li> <li>● 折れ線グラフ</li> <li>■ モザイクプロット</li> <li>■ モザイクプロット</li> <li>■ 円グラフ</li> </ul>	・役割 *X軸:(1項目) @ CP中連絡回数 *Y軸:(1項目) @ 未 @ 最終会話時間 //X軸の変数、Y軸の変数を設定 ★ + 》 列	4 * 5 * 生成日 '2022/08 6 * 生成者 'u620135 7 * 生成に使用された* 8 * 生成に使用された 9 * 生成に使用された 10 * 生成に使用された 11 * 生成に使用された 12 * 13 */ 14 15 ods graphics / re 16
<ul> <li>※ 散布図</li> <li>▲ 系列プロット</li> <li>● 2 マップ</li> <li>● 統計量</li> <li>● データ探索</li> <li>● 要約統計量</li> <li>● 公布分析</li> <li>● 一元度数表</li> </ul>	▶ 追加役割	<pre>17 proc sgplot data= 18 scatter x='CP 19 xaxis grid; 20 yaxis grid; 21 run; 22 23 ods graphics / re</pre>





散布図の描画 – 実行結果





### 層別散布図 – 実行方法と実行結果

6

カテゴリー値で層別化した散布図を描くことで、当初見えなかった関係性が見えてくる可能性がある。分析前における変数間の関係把握、仮説検討を行う上で極めて重要なプロセス

#### [グループ]に層別したい変数を設定し、再度実行

€ *#1_ロジスティック回帰.	cpf ×	
- <u>#1 ロジスティック回帰</u> 🍃	ミ行ボタン	,
設定 コード/結果 分割	* 2	53
データ 表示	情報	ノード
<ul> <li>データ</li> </ul>		
WORK.IMPORT1		
▼フィルタ: (なし)		
▼ 役割		
*X 軸: (1 項目)		<b>≙</b> +
2 CP中連絡回数		
		÷ +
*** 1 <i>項日)</i>		ш т
₩ 取松云品时间		
グループ: (1 項目)		<u> </u>
▲ 定期預金契約		
「グルニアジョドに層別	りしたい変	こ数を設定
(まずは目的変数を設	定してみる	のが効果的)
		L
		L







## 参考:グループごとに散布図を分けて出力

6

#### [追加役割]→[グループ**分析**] に層別したい変数を設定し、再度実行 € \*#1\_ロジスティック回帰.cpf × #1 ロジスティック回帰 実行ボタン 設定 コード/結果 大区 50 分割 データ 情報 ノード 表示 • データ :::C WORK.IMPORT1 -**マ**フィルタ: (なし) - 役割 **命** + \*X 軸: (1 項目) 🐵 CP中連絡回数 盫 \*Y 軸: (1 項目) - +-🐵 最終会話時間 グループ: (1 項目) 亩 + 💊 列 追加役割 **亩** + グループ分析: (1項目) 定期預金契約 「グループ分析」 に層別したい変数を設定 (まずは目的変数を設定してみるのが効果的)









#### **層別散布図行列 – 実行方法** (1/2)

①これまでと同様に、[タスクとユーティリティ]→[タスク]→[統計量]→[データ探索] を選択し、 データインポートノードのコントロールポートに ドラッグ&ドロップ

②生成された [データ探索] ノードをダブルクリックして、詳細設定画面を開く



Copyright © SAS Institute Inc. All rights reserved.



### **層別散布図行列 – 実行方法** (2/2)

[データ]の設定



#### [プロット]の設定







### 層別散布図行列 – 実行結果



先ほど単一の散布図 で確認した2変数

46 **S.Sas** 

Copyright © SAS Institute Inc. All rights reserved.



### 参考:カテゴリー変数の可視化方法

- ・カテゴリー変数における可視化の基本は、「頻度集計」もしくは、 前頁までで活用した「層別化変数」としての活用である
- ここでは、棒グラフを用いた頻度集計の可視化方法について紹介する

▼棒グラフノードの追加



#### ▼ノードの詳細設定









### 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









・ロジスティック回帰は、回帰させる関数にS字カーブ状のロジット関数を用いることで、 ニクラスの分類(購入する/しない、良品/不良品、生/死など)を可能とする回帰



#### ロジスティック回帰のバリエーション

再揭

- 順序ロジスティック回帰: 順序尺度 (1,2,3 / 大,中,小 など)の予測に用いる
  - 多項ロジスティック回帰: 3クラス以上の名義尺度 (犬,猫,人 / 病気の種類 など)の予測に用いる



### 分類モデルの評価指標

- ・モデル評価は、分析対象や目的に応じて最適な指標を選択したり、複数の指標を用いて行う







• ロジスティック回帰



Copyright © SAS Institute Inc. All rights reserved.



### ロジスティック回帰 – 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] →[線形モデル]→[**二項ロジスティック回帰分析**]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二項ロジスティック回帰分析ノードが 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







### ロジスティック回帰 – 実行方法(2/2)説明/目的変数・オプションの設定





### ロジスティック回帰 – 実行結果 (1/4) 基本情報

- 分析結果画面の一番初めには、読み込んだデータの件数、目的変数における各クラスの件数 内訳が出力される
- ・ 本データは、未契約:4,000件 vs 契約:521件の不均衡性の高いデータであることがわかる



モデルの情報						
データセット	WORK.IMPORT1					
応答変数	定期預金契約	定期預金契約				
応答の水準数	2					
モデル	binary logit					
最適化の手法	Fisher's scoring					

読み込んだオブザベーション数 4521使用されたオブザベーション数 4521

	反応プロファ			
順番	定期預金契約	度数の合計		
1	no	4000	未契約:	88.5%
2	yes	521	契約:	11.5%





#### ロジスティック回帰 - 実行結果 (2/4) 精度

- ROC曲線では、左上に凸形状をしており、曲線下面積 (AUC) も0.9程度とまずまずである
- 最もバランスの良い閾値は0.1前後であり、感度 (再現率) 81.8%, 特異度80.8% である



閾値を増加

### 参考: ROC曲線によるモデル評価

- ROC曲線\*は、再現率 (真陽性率) と偽陽性率のバランスの観点から 「良いモデル」を評価する 指標 \*<u>R</u>eceiver Operating Characteristic curve; 受信者動作特性曲線・・・元々は軍事のレーダー性能評価に用いられた指標
- ・ AUC (Area Under the Curve; 曲線下面積) の算出と併せて用いる





#### ロジスティック回帰 – 実行結果 (3/4) 回帰係数

- 標準化回帰係数を見ると、学歴=tertiary (高等教育)や、個人ローン=無、連絡手段=電話、 前回CPの結果=success、最終会話時間で、有意に高い値を示し、関係性が示唆される
- •一方で、最終連絡月が夏(7,8月)や5月・11月になると契約率が下がる可能性が示唆される

#### ▼各説明変数の回帰係数

パラメータ		自由度	推定値	標準誤差	Wald カイ 2 乗値	Pr > ChiSq	標準化した推定値
Intercept		1	-3.8928	0.9183	17.9716	<.0001	
職業	admin.	1	-0.5206	0.5853	0.7914	0.3737	-0.0883
職業	blue-collar	1	-0.9131	0.5810	2.4695	0.1161	-0.2048
職業	entrepreneur	1	-0.7704	0.6460	1.4223	0.2330	-0.0804
職業	housemaid	1	-0.8736	0.6636	1.7333	0.1880	-0.0749
職業	management	1	-0.5937	0.5707	1.0822	0.2982	-0.1343
職業	retired	1	0.1109	0.5985	0.0343	0.8531	0.0134
職業	self-employed	1	-0.7018	0.6294	1.2433	0.2648	-0.0763
職業	services	1	-0.6663	0.5986	1.2388	0.2657	-0.1063
職業	student	1	-0.1422	0.6464	0.0484	0.8259	-0.0106
職業	technician	1	-0.7133	0.5743	1.5427	0.2142	-0.1477
職業	unemployed	1	-1.1602	0.6711	2.9891	0.0838	-0.1061
職業	unknown	0	0				
結婚歷	divorced	1	0.3051	0.2038	2.2414	0.1344	0.0540
結婚歷	married	1	-0.1644	0.1478	1.2385	0.2658	-0.0440
結婚歷	single	0	0				
学歴	primary	1	0.4210	0.3572	1.3890	0.2386	0.0829
学歴	secondary	1	0.5011	0.3245	2.3846	0.1225	0.1381
学歴	tertiary	1	0.7418	0.3361	4.8711	0.0273	0.1872
学歴	unknown	0	0				
クレカ債務	no	1	-0.5446	0.4315	1.5935	0.2068	-0.0386
クレカ債務	yes	0	0				
住宅ローン	no	1	0.2600	0.1381	3.5463	0.0597	0.0710
住宅ローン	yes	0	0				
個人ローン	no	1	0.6296	0.2000	9.9137	0.0016	0.1249
個人ローン	yes	0	0				

パラメータ		自由度	推定值	標準誤差	Wald カイ 2 乗値	Pr > ChiSq	標準化した推定値
連絡手段	cellular	1	1.4161	0.2277	38.6781	<.0001	0.3747
連絡手段	telephone	1	1.3459	0.3119	18.6232	<.0001	0.1850
連絡手段	unknown	0	0				
最終連絡月	apr	1	-0.6572	0.4115	2.5504	0.1103	-0.0892
最終連絡月	aug	1	-0.9653	0.3995	5.8382	0.0157	-0.1847
最終連絡月	dec	1	-0.5428	0.7216	0.5657	0.4520	-0.0199
最終連絡月	feb	1	-0.4550	0.4233	1.1552	0.2825	-0.0542
最終連絡月	jan	1	-1.7805	0.4998	12.6885	0.0004	-0.1747
最終連絡月	jul	1	-1.4087	0.4113	11.7294	0.0006	-0.2820
最終連絡月	jun	1	-0.1030	0.4252	0.0587	0.8086	-0.0183
最終連絡月	mar	1	0.8413	0.4959	2.8783	0.0898	0.0480
最終連絡月	may	1	-1.1472	0.3981	8.3049	0.0040	-0.2924
最終連絡月	nov	1	-1.5002	0.4228	12.5876	0.0004	-0.2320
最終連絡月	oct	1	0.7038	0.4559	2.3832	0.1226	0.0512
最終連絡月	sep	0	0				
前回CP結果	failure	1	0.1216	0.3199	0.1445	0.7038	0.0208
前回CP結果	other	1	0.6128	0.3495	3.0749	0.0795	0.0690
前回CP結果	success	1	2.5665	0.3148	66.4610	<.0001	0.2356
前回CP結果	unknown	0	0				
年齡		1	-0.00423	0.00713	0.3528	0.5525	-0.0247
年間平均残高		1	-3.91E-6	0.000017	0.0500	0.8230	-0.00649
最終連絡日		1	0.0164	0.00816	4.0427	0.0444	0.0746
最終会話時間		1	0.00423	0.000202	437.3055	<.0001	0.6053
CP中連絡回数		1	-0.0704	0.0282	6.2314	0.0126	-0.1207
最終連絡日数		1	-0.00010	0.000996	0.0097	0.9217	-0.00540
CP前連絡回数		1	-0.00551	0.0382	0.0208	0.8853	-0.00514



### 参考:回帰係数とデータの標準化

- ・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





### ロジスティック回帰 - 実行結果 (4/4) オッズ比

• 2群間の事象の起こりやすさを表す「オッズ比」では、関係性が明白である「前回キャンペーン 成功」以外にも、学歴や電話連絡、3月の連絡など、意外な影響因子が浮かび上がった



#### まとめ

- ・グラフによるデータ観察
  - 散布図、層別散布図、散布図行列を用いて、変数間の関係性を俯瞰的に把握した
  - ヒストグラム、箱ひげ図を用いて、各変数の分布や外れ値を確認した
- ・ ロジスティック回帰分析
  - 教師あり学習の分類手法には、ロジスティック回帰の他、 決定木やサポートベクターマシンなどがある
  - ロジスティック回帰の特徴は、決定境界が直線的であること、0/1の予測だけでなく発生確率 も出力できることである
  - ロジスティック回帰を用いて、契約確度の高い顧客を予測し、精度の評価を行なった
  - 回帰係数の比較や、オッズ比の確認を行うことで、影響因子の探索を行った



### アンケートのお願い・ご質問

### <u>10月12日 機械学習によるビッグデータ分析の手法-1</u>

今後の参考にさせていただくため、ぜひともアンケートにご協力を お願いします。

・無記名
 ・所要時間目安: 1~3分



https://sas.qualtrics.com/jfe/form/SV\_1HeF5SyxqC9FqtM



- ・本日のアーカイブは、2022年10月17日~2023年3月31日迄 視聴できます。
- 本日の内容に関するご質問は、以下宛にご連絡ください。
   que@datascience.co.jp
- ご視聴ありがとうございました。

# **End of File**





Copyright © SAS Institute Inc. All rights reserved.