SASお客様向けwebinar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) 階層的クラスタリング \_\_\_\_\_

2021年12月15日



## **Agenda**

- 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
  - 各クラスタの解釈、予測への適用方法
  - 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て
  - データサイエンティストに求められるスキル
  - 学習リソースの紹介
  - SAS内サンプルデータの紹介と使い方



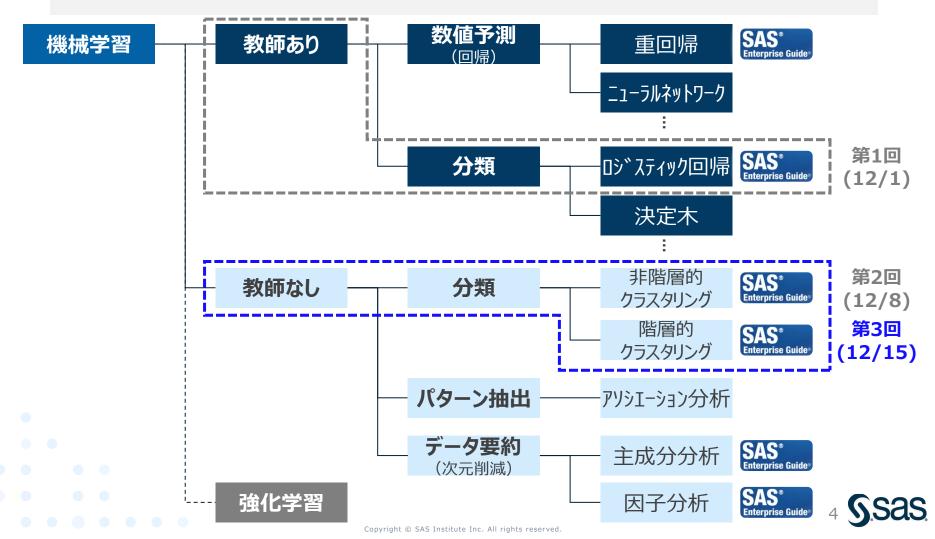
### **Agenda**

- 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
  - 各クラスタの解釈、予測への適用方法
  - 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て
  - データサイエンティストに求められるスキル
  - 学習リソースの紹介
  - SAS内サンプルデータの紹介と使い方



#### 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

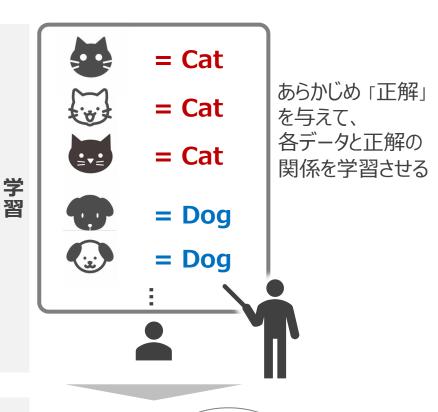


推論

(نعن)

## 教師あり学習と教師なし学習

## 教師あり学習



Dog!

## 教師なし学習



「正解」を与えずに、 各データのパターン (距離の近さ、頻出の 組み合わせなど)を 学習する



推論

学習



※分類されたグループの 意味づけは人が行う



## 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

クラスタリング

• **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

#### データ例

顧客ID	名前	年齢	年収	購入額	購入有無	
0001	XX	25	300万	35,000	購入	
0002	XX	35	600万	68,000	購入	
0003	XX	18	120万	0	非購入	
0004	XX	42	820万	85,000	購入	
:	:	1	:	1	:	

説明変数

※目的変数は無し



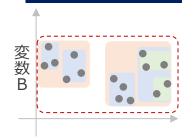
#### 非階層的クラスタリング



#### 主な手法

- k-means法 (k平均法)
- ・混合ガウス

#### 階層的クラスタリング



#### 主な手法

- 最短距離法
- 最長距離法
- 群平均法
- ウォード法



### クラスタリング手法の種類



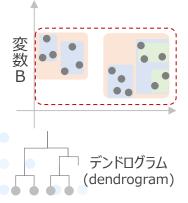
- クラスタリング手法は、「**非階層的**」と「**階層的**」に大別される
- 階層的クラスタリングはさらに 凝集型 と 分割型 があり、凝集型が用いられるのが一般的

#### 手法の分類

#### 非階層的クラスタリング



#### 階層的クラスタリング



#### 手法

• 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<sup>®</sup> Enterprise Guide

■ その他 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)

• 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ん	SAS* Enterprise Guide*
■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
■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
■重心法	クラスタ重心からの距離に基づき併合	SAS® Enterprise Guide®
■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	∰⊜ SAS°

■ その他 メディアン法、可変法

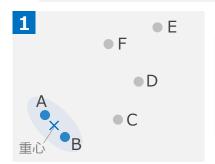
似ていないデータ/クラスタ同士を逐次分離させる(トップタウンアプローチ)

■ Diana法

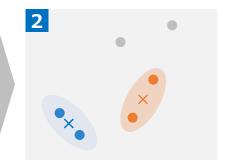


#### 代表的な**階層的**クラスタリング: 凝集型階層クラスタリング

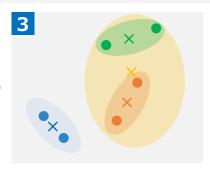
- 凝集型階層クラスタリングは、距離に応じて小さいクラスタを束ねて階層的に分類する手法
- クラスタ数は自動的に決定してくれる他、分類過程を可視化した**樹形図 (デンドログラム)** も同時 に出力されるので、結果の解釈やクラスタ数の決定に役立つ



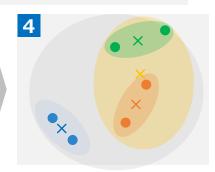
最も距離の近い点を 同一クラスタとしてまとめ、 重心を計算



同様に、重心と残りの点の中で 最も距離の近い点同士をまとめ、 重心を計算



②を繰り返し、クラスタの 重心同士が最近傍となった 場合は、クラスタ同士をまとめる



最後にクラスタ全体をまとめる

#### 凝集型階層的クラスタリング

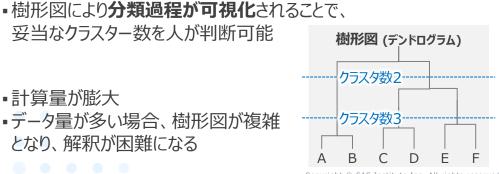
(agglomerative hierarchical clustering)

クラスタ数は自動決定

## ■計算量が膨大

■データ量が多い場合、樹形図が複雑 となり、解釈が困難になる

妥当なクラスター数を人が判断可能



#### Copyright © SAS Institute Inc. All rights reserved

## Iteration m-3 Builds up a sequence of clusters ("hierarchical") Dendrogram:

In matlab: "linkage" function (stats toolbox)

Algorithmic Complexity: O(m2 log m) + (m-3)\*O(m log m) + 出典:https://youtu.be/OcoE7JlbXvY

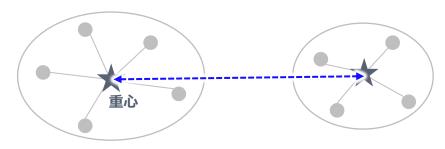


#### 「近い」の評価尺度バリエーション

- 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

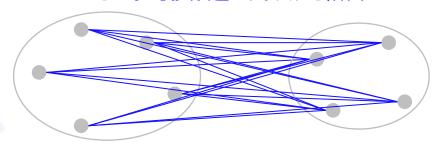
重心法

#### 重心間の距離が近いクラスタを結合



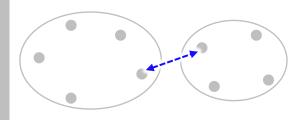
群平均法

クラスタ間で全データ間の距離を算出し、 その**平均値**が近いクラスタを結合



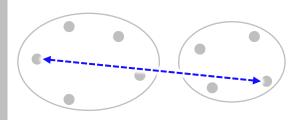
最短距離法

2つのクラスタ間で**最近傍**のデータを クラスタ間距離として採用



最長距離法

2つのクラスタ間で**最遠方**のデータを クラスタ間距離として採用



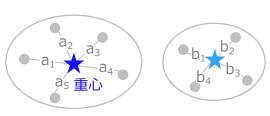


#### Ward法の考え方

- 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

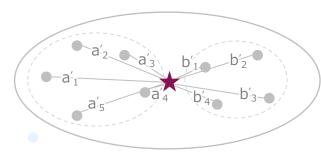
#### クラスタA

クラスタB





#### A, Bの結合を仮定した場合のクラスタAB



1

「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

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

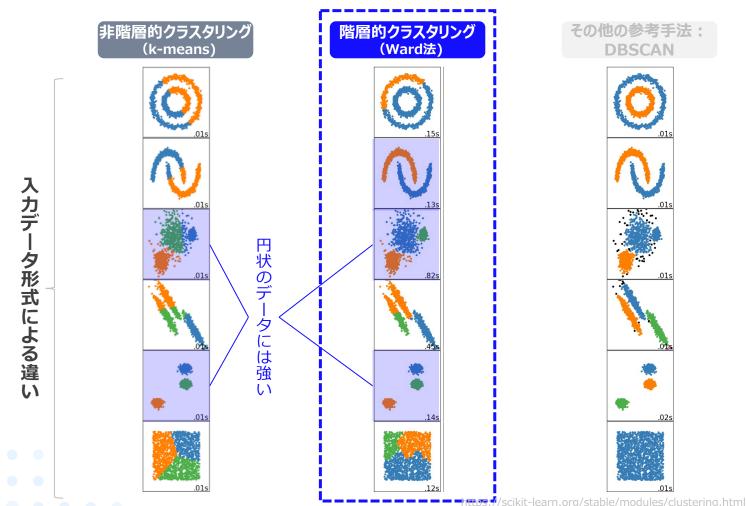
1 と 2 の差、つまり、AB - (A+B) が最小となるクラスター結合を採用 (結合前後でクラスタ内のばらつきに変化なし→統合してもOKと判定)

※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



#### 参考:クラスタリング手法における分類結果の比較

• 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



**S**.sas



# SAS Enterprise Guide での実装方法

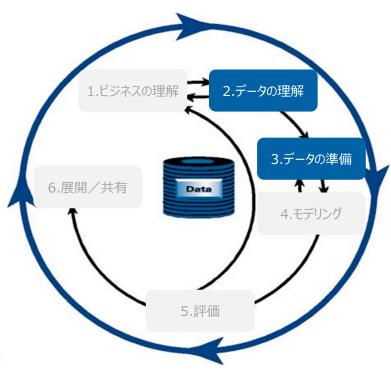
- データの読み込み
- 階層的クラスタリング(Ward法、重心法、群平均法)
- 標準化したクラスタリング



## 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



(CRoss Industry Standard Process for Data Mining)

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

## 使用データ (e-Statについて)

- 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



#### 統計で見る日本

e-Statは、日本の統計が閲覧できる政府統計ポータルサイトです

お問い合わせ | ヘルプ | English ログイン 新規登録

統計データを探す 統計データの活用 統計データの高度利用 統計関連情報 リンク集 ●統計データを探す (政府統計の調査結果を探します) その他の絞込 利用ガイド

1 すべて

政府統計一覧の中から探します

🗞 分野

17の統計分野から探します

1 組織

統計を作成した府省等から探します

キーワード検索: (例:国勢調査

統計データを活用する

主要指標をグラフで表示 (統計ダッシュボード)

時系列表

主要指標を時系列表で表示 (統計ダッシュボード)

地図

地図上に統計データを表示 (統計GIS)

地域

都道府県、市区町村の 主要データを表示

●統計データの高度利用

ミクロデータの利用 公的統計のミクロデータの利用案内

開発者向け

API、LODで統計データを取得

統計関連情報

統計分類・調査計画等

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



#### 使用データ

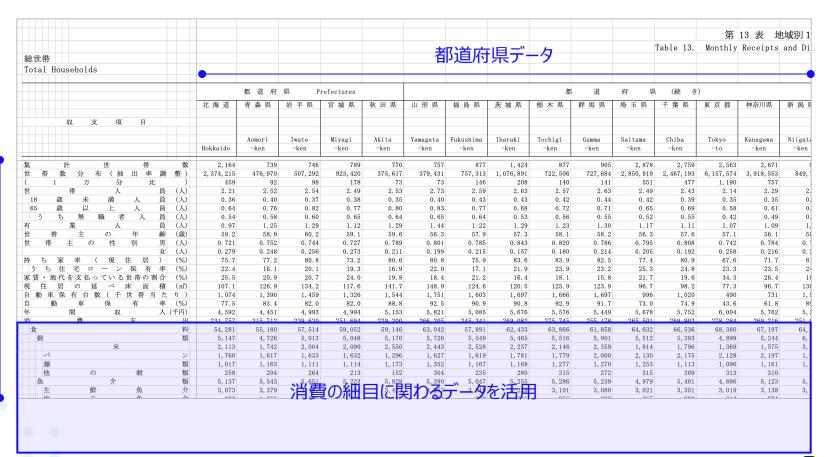
• 今回は、このうち、5年に1度実施している「全国消費実態調査」(現在の名称は「全国家計構造調査」) のデータを用いて、**都道府県別の消費動向から、類似の都道府県をグルーピング**すること を考える

					_									
政府統計名		全国家計構造調査(旧全国消費実態調査)												
提供統計名		平成26年全国消費実態調査												
提供分類1		全国												
提供分類2		家計収支に関する結果	家計収支に関する結果											
提供分類3		総世帯												
表番号		統計表	調査年月	公開 (更新) 日	表示・ダウンロー									
フロー編														
42	年間収入階級・年間収入十分位階級別1世帯当たり1か月間の収入と支出													
	総世帯		2014年	2015-12-16	<b>★</b> EXCEL → DB									
	勤労者世帯		2014年	2015-12-16	<b>★</b> EXCEL → DB									
43	世帯主の年齢階級別1世帯当たり1か月間の収入と支出													
	総世帯・勤労	<del>芍者世帯</del>	2014年	2015-12-16	<b>★</b> 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	<del>芍者世帯</del>	2014年	2015-12-16	<b>★</b> 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	<b>★</b> EXCEL → DB									
	勤労者世帯		2014年	2016-03-25	<b>★</b> 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	<b>★</b> 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	当たり1か月間の収入と支出												
	総世帯		2014年	2015-12-16	<b>★</b> EXCEL → DB									
	勤労者世帯		2014年	2015-12-16	<b>★</b> EXCEL → DB									

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

#### 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



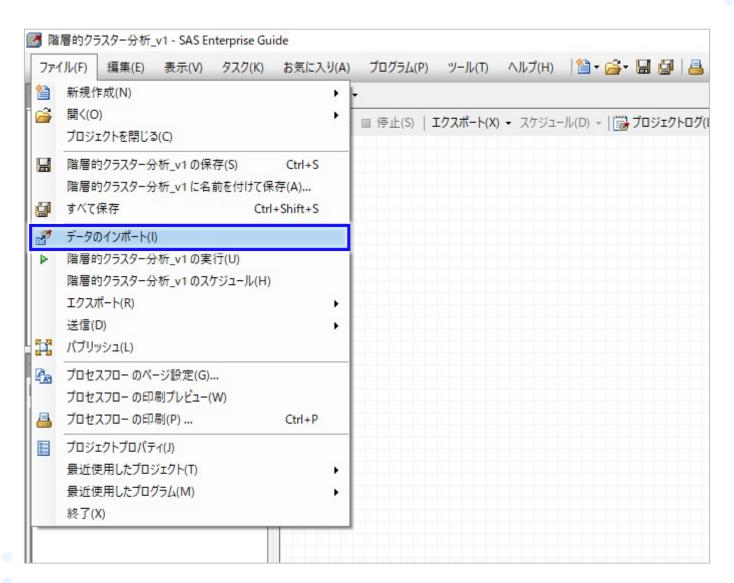
- e-Statより素データをダウンロードして開くと、開始行や開始列がずれていたり、空白行があった りと、加工が必要な形式であることがわかる
- 今回は、本データから都道府県別の消費細目データ部分を抽出し、加工済のデータを用いる

予測(分析)対象を説明するための変数

# 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 北海道	54281	17491	19520	9089	10208	9234	35627	5630	23323	3894	2106	6506	2492	1430	1011
2 青森県	55180	14357	22420	9162	8972	9936	33039	4880	16564	3256	2084	3839	1153	777	1079
3 岩手県	57514	14782	21267	8973	8288	10273	42912	5805	20278	3534	2428	6326	1447	1219	1458
1 宮城県	59052	16353	20331	9700	10640	10851	40742	7331	23394	3241	2322	7734	2172	1587	1168
5 秋田県	59146	12875	22394	9108	7467	9214	37645	4472	21037	3349	2403	6672	1617	979	1191
5 山形県	63042	13186	24030	11771	9486	10407	40365	6667	24577	3706	2227	5438	1659	1333	2230
7 福島県	57891	12579	20102	9633	8949	9613	44195	4614	21216	3637	1921	6701	1743	1432	1027
8 茨城県	62433	17292	20330	9186	10645	10668	44764	10113	26592	3638	2451	7254	2639	2041	1134
9 栃木県	63866	18994	19997	9914	10252	12956	47208	6889	27561	3564	2723	7996	2908	2400	996
10 群馬県	61858	15629	18305	9605	9682	11270	43782	7539	25659	3767	2473	6800	2837	1965	931
11 埼玉県	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 東京都	68380	33295	16315	8691	12404	12151	33118	11060	32038	4180	3882	14361	3061	2964	808
14 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
.6 富山県	67635	14518	21894	10624	9387	10776	51532	7879	27246	4003	4101	5869	2139	2433	1373
7 石川県	66478	17678	18423	8733	9512	11135	42087	7993	27548	4220	3207	8516	1898	1966	1402
18 福井県	67429	12168	20741	9034	10204	11287	45576	9585	27984	3482	4031	8052	1863	2080	1539
19 山梨県	57641	17234	18209	7890	9429	10280	39392	9066	25849	3531	3699	6404	1877	2212	1132
0 長野県	62406	21145	21350	9866	9375	11987	42846	8047	27147	4035	2910	7265	2537	1830	1224
1 岐阜県	61939	12754	19952	9042	9942	10463	41580	-0.4	25777	3627	2610	6754	2280	2056	1306
2 静岡県	62396	15048	18407	9012	9985	11488	0		28082	3714	2936	8355	2350	2259	1048
3 愛知県	64248	21485	17573	9010	11051	11880	=14	ロロガに坐	28967	4041	3305	8547	2948	2655	818
4 三重県	63275	12856	19237	9036	11444	12889	<b>一</b>	明変数	28462	3556	3578	8754	2773	2385	1173
滋賀県	63385	16479	18807	9587	10034	10818			26609	3280	2845	8377	1848	2052	1486
5 京都府	65337	13829	17928	8409	12630	9239	36645		26012	3984	3109	8260	1855	2174	985
7 大阪府	62386	18778	16292	7230	9898	10782	31046	10348	25016	3744	3264	7492	2465	2250	658
8 兵庫県	63620	19262	16725	8281	10712	10926	36040	9806	27000	3827	3023	8397	2546	2548	823
29 奈良県	66408	17630	19784	9875	11068	12405	42593	14481	27121	3849	3065	9467	2170	2609	986
30 和歌山県	58010	10696	17125	8152	9250	8326	36333	6001	23890	3376	2656	4995	2098	1787	888
1 鳥取県	58027	13626	18488	8143	9050	10320	41570	4966	24212	3198	3787	7077	1775	1865	904
2 島根県	59223	11926	19494	8915	8767	11814	40722	3866	23678	3446	3538	6335	1619	1512	1656
3 岡山県	58368	13776	18306	8286	9846	10347	38978	8451	24914	3052	2796	6682	2430	1789	931
84 広島県	58058	17721	17128	9180	9622	11195	38580	8773	24978	3308	2660	7944	1918	2004	997
35 山口県	55832	18576	16610	9381	8003	10961	35524	5193	23931	3557	2873	6234	1514	1798	840
36 徳島県	55896	16389	18015	8680	9656	10261	38507	6659	23923	3439	3762	7064	2039	1810	997
37 香川県	57352	15438	17319	8338	8754	11070	40876	6059	25565	3476	2814	6197	2689	2176	942
38 愛媛県	55531	13489	17201	8171	8284	9224	32679	7901	19353	2872	2450	5539	1732	1774	937
39 高知県	54971	14463	16479	7609	7510	10329	32613	5206	20184	3273	2447	5228	1722	1167	867
40 福岡県	54633	18999	16314	8029	9823	10405	36057	7360	24134	3073	2790	9478	2489	1809	784
41 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
43 熊本県	55006	11286	16802	8254	10041	11155	34633	6967	21046	2889	2309	5544	1688	2016	682
44 大分県	53558	14707	15685	8558	10853	11677	36458	3243	22105	3021	3223	5354	2139	1327	1153
45 宮崎県	53347	15963	15828	8228	8386	9476	36294	6276	21314	2565	3061	6208	2410	1413	1210
46 鹿児島県	50294	14792	15496	7800	7857	10022	39992	5063	18721	2593	2002	5533	2053	1160	1488
47 沖縄県	48770	22616	17251	6750	5010	8088	28055	5169	16217	2500	1492	3913	1767	1700	970

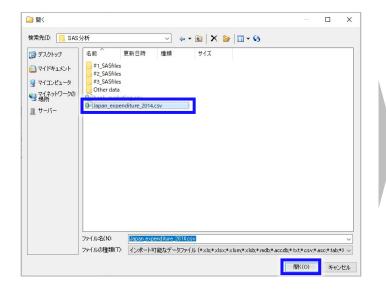


## データの読み込み (1/2)





## データの読み込み (2/2)



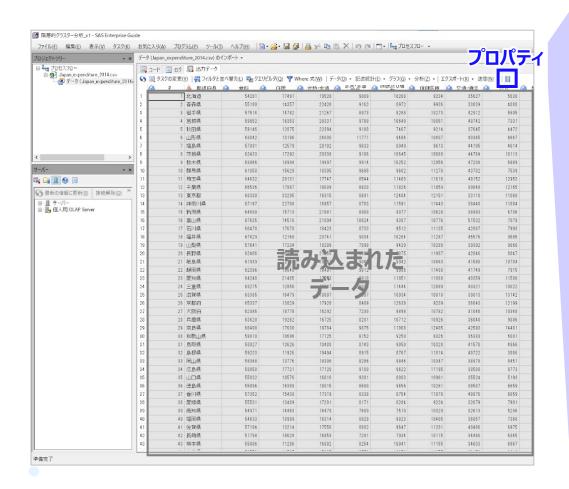


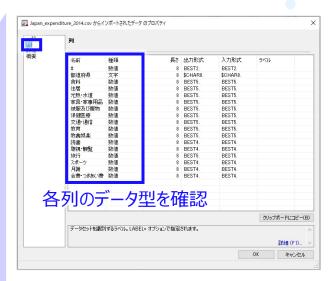


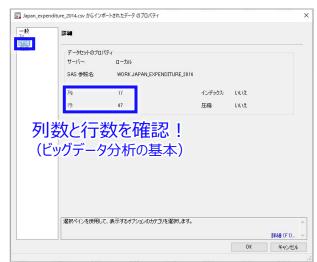




#### 読み込んだデータの確認





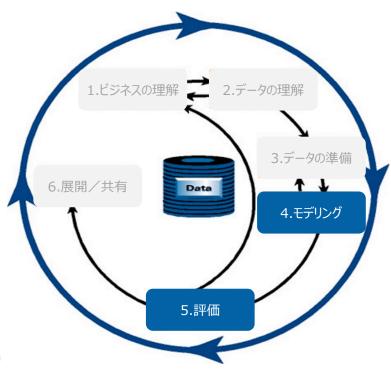




## ビッグデータ分析の進め方

• 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



(CRoss Industry Standard Process for Data Mining)

- 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



## 階層的クラスタリング(群平均法) - 実行方法(1/3)

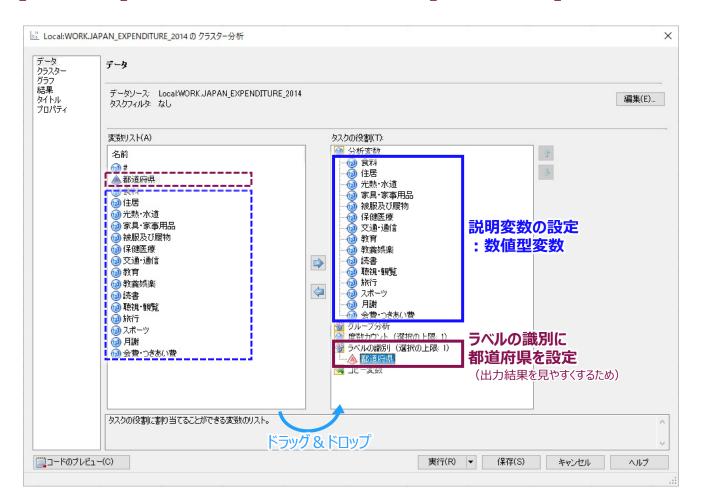
#### [分析] → [多変量解析] → [クラスター分析] をクリック

津料	(3) f	据	6 光熱·水道	19 家品.歌事	⑩ (放版及U)腹 物	16	分散分析	f(A)	- F	123	教育	(123)	教養娯楽	(123)	読書	(123)	聴視·観覧	(123)	旅行	(23)	スポーツ
54281		17491	19520	9089	De la constant de la		▼ 回帰分析		-,		5630	_	233	_	3894		2106		6506		2492
55180		14357	22420	9162	8972		多変量解			1/	相関分析( <u>C</u>	1		i4	3256		2084		3839		1153
57514		14782	21267	8973	8288		生存時間			11 13	正準相関分			78	3534		2428		6326		1447
59052		16353	20331	9700	10640									94	3241		2322		7734		2172
59146		12875	22394	9108	7467		工程能力	1分析( <u>Y</u> )	•	逐.	主成分分析			37	3349		2403		6672		1617
63042		13186	24030	11771	9486		管理図(0	)	•	ha.	因子分析( <u>F</u>			77	3706		2227		5438		1659
57891		12579	20102	9633	8949	lim	パレート図	l( <u>P</u> )	Ī	Ŀ	クラスター分	沂( <u>L</u> )		16	3637		1921		6701		1743
62433		17292	20330	9186	10645		時系列分	/析(T)		×	判別分析( <u>D</u>	)		92	3638		2451		7254		2639
63866		18994	19997	9914	10252		データマイ				6889		275	61	3564		2723		7996		2908
61858		15629	18305	9605	9682	_	7-2 (4	_//( <u>IN</u> )			7539		256	59	3767		2473		6800		2837
64632		20131	17747	8544	11403		11616		40152		12953		290	55	3800		3115		8887		2699
66536		17887	18039	8820	11826		11859		39048		12165		3038	35	4186		3512		10407		2889
68380		33295	16315	8691	12404		12151		33118		11060		320	38	4180		3882		14361		3061
67197		22708	16957	8783	11591		11443		38440		11004		3183	33	4159		3607		10790		2874



### 階層的クラスタリング(群平均法) - 実行方法(2/3)

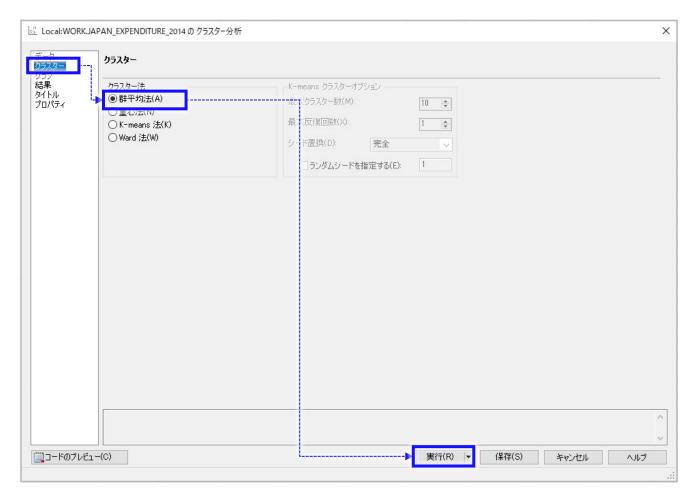
- ・ 左側の [変数リスト] より、数値型の変数を選択し、右側の [分析変数] にドラッグ&ドロップ
- ・ 左側の [変数リスト] より、都道府県を選択し、右側の [ラベルの識別] にドラッグ&ドロップ





## 階層的クラスタリング(群平均法) - 実行方法(3/3)

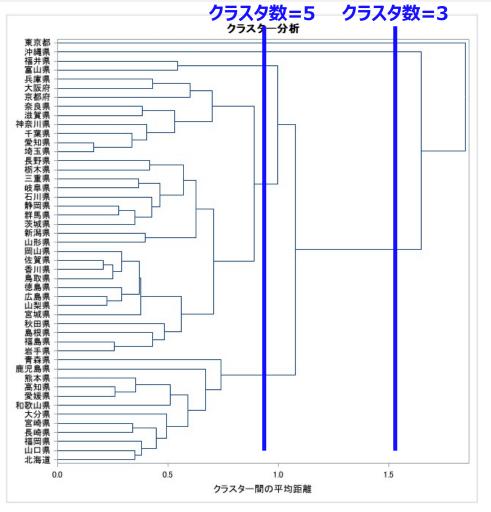
- ・ 左パネルから [クラスター]メニュー をクリック
- ・ クラスター法から[群平均法]を選択し、実行ボタン





#### 階層的クラスタリング(群平均法) - 実行結果

- 群平均法の結果、47都道府県が階層的にクラスタリングされた
- 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

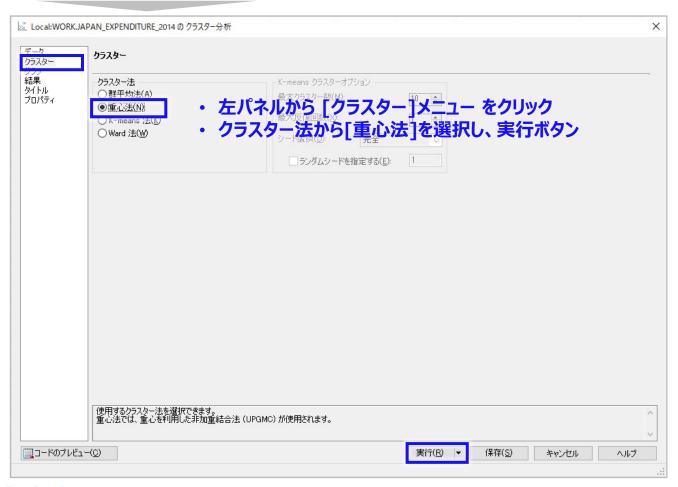




## 階層的クラスタリング (重心法):実行方法



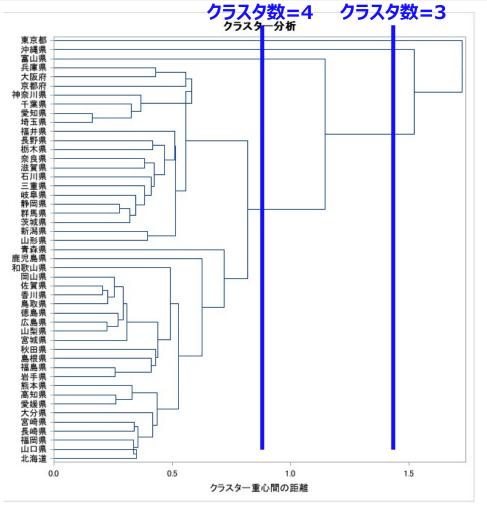
クラスタリング分析結果の画面上において 上部メニュー [タスクの変更] をクリック





#### 階層的クラスタリング(重心法) - 実行結果

- 重心法の結果、47都道府県が階層的にクラスタリングされた。群平均法と類似の傾向を示すが、部分的に、分割結果が複雑化している
- 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



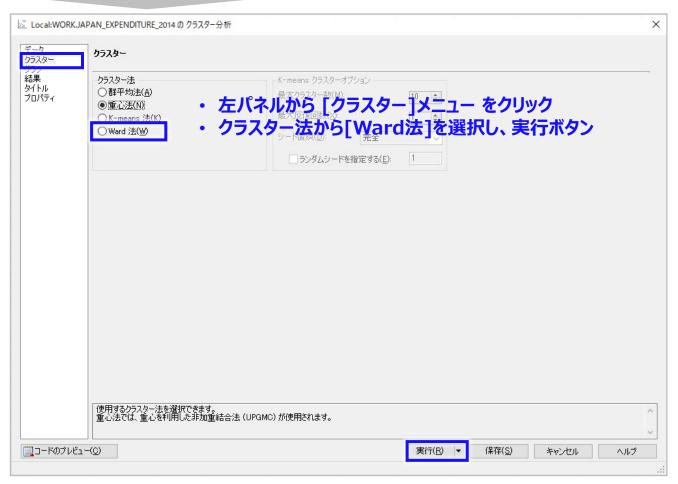




#### 階層的クラスタリング (Ward法):実行方法



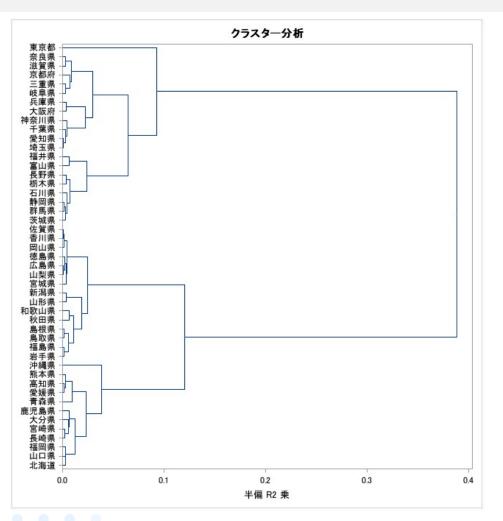
クラスタリング分析結果の画面上において 上部メニュー [タスクの変更] をクリック



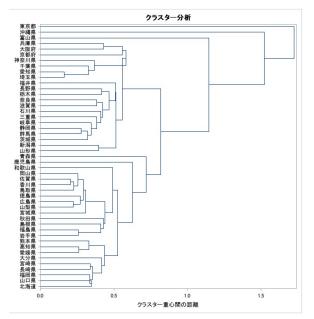


## 階層的クラスタリング (Ward法) - 実行結果

• 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



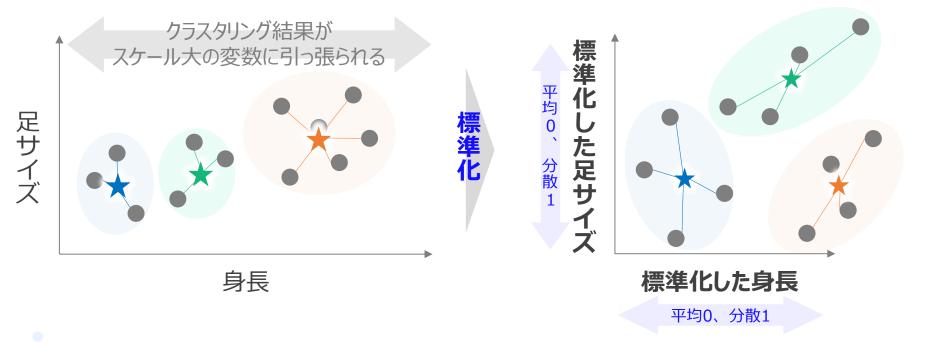
参考:重心法の結果





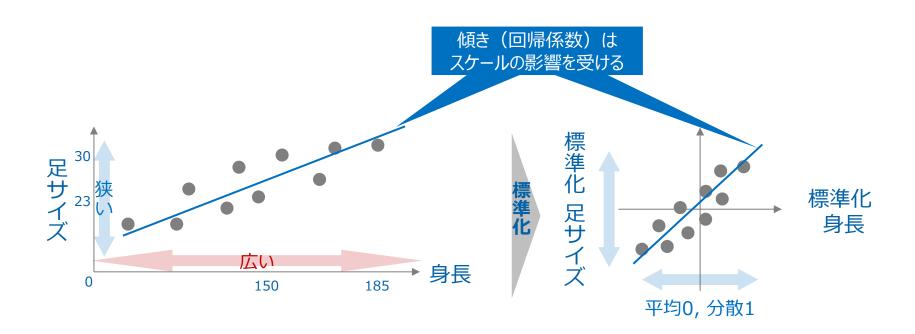
### クラスタリングにおける変数スケールの影響と標準化

• 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

- 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





## データの標準化 - 実行方法 (1/2)

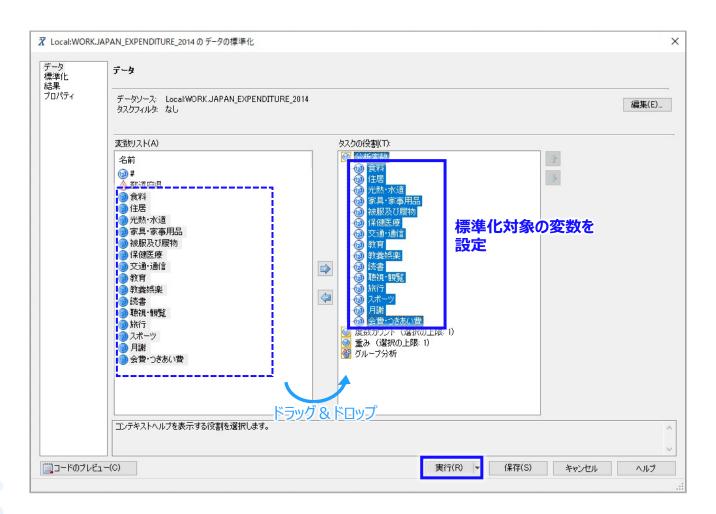
- ・ [入力データ]タブ をクリック
- ・ [データ]→ [データの標準化] をクリック

🔣 入力データ			🔛 出力データ 🕍					
77 フィルタと並/	<替え(L	== ##	(C) Two tou	デー	タ(D) ▼ 記述統計(E) ▼ グラフ(G) ▼ 分析(Z)	<ul><li>▼   エクスポート</li></ul>	·(X) · 送信(N) ·	=
130 #	4	▲ 都道府県	⑩ 食料 ⑩		テーブルの追加(B)	放肥ない腹 物	⑩ 保健医療	⑩ 交道
1	1 4	比海道	54281	T.	データの並べ替え(S)	10208	9234	
2	2 7	<b>手森県</b>	55180	Sw.	出力形式の作成(C)	8972	9936	
3	3 7	台手県	57514	Sw.	データセットから出力形式を作成(M)	8288	10273	
4	4 3	官城県	59052	7.0	転置(T)	10640	10851	
5	5 ₹	火田県	59146			7467	9214	
6	6 L	L形県	63042		列の分割(P)	9486	10407	
7	7 7	副島県	57891	<b>=</b>	列の積み上げ(K)	8949	9613	
8	8 7	<b>茨城県</b>	62433	§+8	ランダムサンプル(R)	10645	10668	
9	9 A	<b>厉木県</b>	63866	1=	▼ランク(A)	10252	12956	
10	10 君		61858	$\overline{X}$	データの標準化(Z)	9682	11270	
11	11 ±	奇玉県	64632	闘	データセットの属性(D)	11403	11616	
12	12 =	f葉県	66536	12	データの比較(O)	11826	11859	
13	13 3	東京都	68380			12404	12151	
14	14 🕏	<b>•奈川県</b>	67197	SX.	データセットと出力形式の削除(N)	11591	11443	
15	15 ≇	消湯県	64400	5	LASR にアップロード(U)	9077	10628	
16	16 2	山県	67635	5	CAS にアップロード(U)	9387	10776	
17	17 7	5川県	66478		データファイルをサーバーにアップロード(F)	9512	11135	



### データの標準化 - 実行方法 (2/2)

• 左側の[変数リスト]から、全変数を選択\*し、右側の[分析変数]ヘドラッグ&ドロップにより移動 \*全変数を選択してドラッグ&ドロップをすれば、自動的に数値型変数のみがセットされる





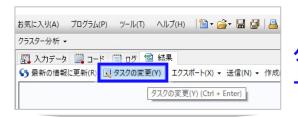
## データの標準化 - 実行結果

• 標準化を行うと、元データの右側に、標準化された後のデータが出力される

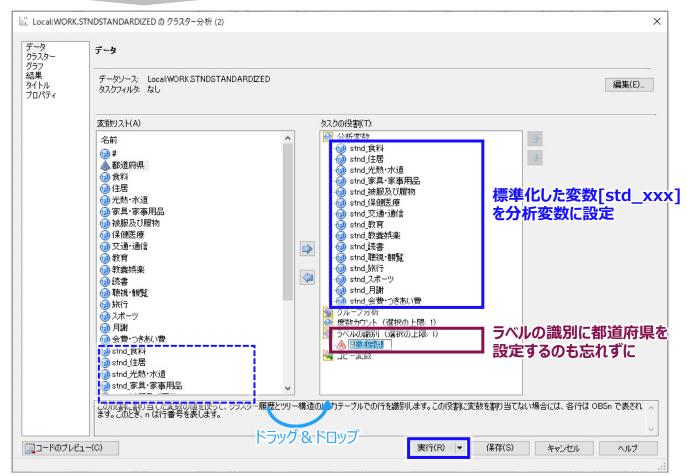
F-90	標準化 -																	
7	カデータ 📋 コー	r 📋 o5 🖁	3 出力データ															
5	タスクの変更(Y)	<b>ジ</b> フィルタと並	べ替え(L) 🛅 クエ!			データ(D) • 記述統	計(E) • グラフ(G)	→ 分析(Z) →										
6	車視·観覧 🚳	抗行于	□ スポーツ 6	月謝	⊕ 会費・つきあ	⊚ stnd_食料	⑩ stnd_住居	⊚ stnd 光熱・ 水道	⑩ stnd 家具・ 家事用品	動 stnd 被服及び履物	ஞ stnd_保健医 療	op stnd 交通 i順篇	⊚ stnd_教育	⑩ stnd_教養娯 ♀	stnd_読書	⊚ stnd 聴視・ 観覧	stnd_抗行	
1	2106	6506	2492	1430	101	-1.126440835	0.2824331478	0.4643504342		0.3941421031	-1.414683268	-0.774476972	-0.838208783	-0.370316738	0.9167736168	-1.300082128	-0.403623942	0.696109100
2	2084	3839	1153	777	1075	-0.946537808	-0.523179513	1.8855653575	0.403847784	-0.489142484	-0.753005689	-1.363677215	-1.115883954	-2.246876089	-0.53543841	-1.337029273	-1.851479295	-2.1063623
3	2428	6326	1447	1219	1458	-0.479470328	-0.413930827	1.8205099069	0.1885059348	-0.977950459	-0.435362749	0.8840716269	-0.77341791	-1.215726299	0.0978436828	-0.75931028	-0.501341963	-1.4910325
1	2322	7734	2172	1587	1168	-0.17169407	-0.010096274	0.8618005386	1.0168314604	0.7028629295	0.1094373082	0.3900358742	-0.208441495	-0.350604397	-0.569581326	-0.937328342	0.2630301167	0.02636237
5	2403	6672	1617	979	119	-0.152883298	-0.904136106	1.8728234306	0.3423215414	-1.564662955	-1.433534481	-0.315046488	-1.266939247	-1.004998596	-0.323752331	-0.801295672	-0.313506211	-1.1352296
6	2227	5438	1659	1333	2230	0.6267631901	-0.824191774	2.6745846769	3.3764768022	-0.121821871	-0.309059622	0.3042056997	-0.45427658	-0.022157924	0.4888490696	-1.096872831	-0.983417536	-1.04732535
7	1921	6701	1743	1432	1027	-0.404027128	-0.980224603	0.7495735671	0.9404933446	-0.505579009	-1.05745278	1.1761674199	-1.214366081	-0.955301285	0.331791656	-1.610774028	-0.297762752	-0.87151683
3	2451	7254	2639	2041	1134	0.5048933977	0.2312790574	0.8613104645	0.4311927807	0.7064360872	-0.063051291	1.3057095136	0.8215482727	0.5372838712	0.3340678504	-0.720683719	0.0024487258	1.00377400
3	2723	7996	2908	2400	991	0.7916576219	0.6687879115	0.6981157854	1.2606576812	0.425585891	2.0935274826	1.8621258176	-0.372084729	0.8063156822	0.1656294648	-0.263882655	0.4052641259	1.566779848
10	2473	6800	2837	1965	93	0.3898275018	-0.196204623	-0.131089611	0.9085908484	0.0182459117	0.5043702217	1.0821412605	-0.131432914	0.27824705	0.627696928	-0.683736574	-0.24401784	1.418179792
11	3115	8887	2699	2705	633	0.9449454067	0.9610602777	-0.404550965	-0.300285882	1.2481267965	0.8304962077	0.2557127894	1.873011587	1.2211077624	0.7028113432	0.8944482907	0.8889683328	1.129351518
12	3512	10407	2889	2800	684	1.3259636083	0.3842272173	-0.261449325	0.0141815807	1.550415939	1.0595384463	0.0043692543	1.5812675406	1.590367111	1.5814223815	1.0611763147	1.7141427872	1.527013636
3	3882	14361	3061	2964	808	1.6949749337	4.3449419205	-1.106337093	-0.132797777	1.9634729706	1.334766157	-1.345691581	1.1721594553	2.0493037298	1.5677652151	1.6825601155	3.8606819447	1.88700250
4	3607	10790	2874	2894	888	1.4582393599	1.6234928965	-0.791709514	-0.027975289	1.3824775265	0.6674332147	-0.134051823	1.1514263759	1.9923878152	1.5199651327	1.2207208041	1.922064972	1.49561925
5	2535	7327	1955	1715	1033	0.8985188191	-0.174611942	1.6214154114	0.1053315698	-0.414106172	-0.100753717	-0.010429052	-0.428730464	-0.216227311	0.2589534352	-0.579612802	0.0420788124	-0.42780963
16	4101	5869	2139	2433	1873	1.5458895555	-0.48179354	1.6277863748	2.0696138341	-0.192570393	0.0387452592	2.8465546632	-0.005553503	0.7188595207	1.1648788063	2.0503521488	-0.749437162	-0.0427052
7	3207	8516	1898	1966	140:	1.3143569614	0.3305025695	-0.073260866	-0.084944032	-0.103241451	0.3771245336	0.6962469744	0.0366531226	0.8027063803	1.6588129911	0.5489545331	0.6875606327	-0.54710826
8	4031	8052	1863	2080	1535	1.5046659475	-1.085874508	1.0627309243	0.2580078015	0.391283577	0.5203937528	1.490574505	0.6260649523	0.9237568134	-0.021018476	1.9327930514	0.4356652882	-0.6203618
9	3699	6404	1877	2212	1133	-0.454055773	0.216369825	-0.178136726	-1.045437042	-0.162555869	-0.428764825	0.0826864428	0.4339137339	0.3309983855	0.0905150496	1.3752270465	-0.458997487	-0.59106039
20	2910	7265	2537	1830	1224	0.4994903035	1.2217147889	1.3611860582	1.2059676877	-0.201145972	1.1801862099	0.8690456547	0.0566457849	0.6913732985	1.2377170271	0.0501680768	0.0084203827	0.790292234
21	2610	6754	2280	2056	1300	0.4060367845	-0.93523985	0.6760624503	0.2671228004	0.2040501128	-0.256276225	0.5808201879	1.0403563074	0.3110084057	0.309029712	-0.453656626	-0.268990223	0.252401892
22	2936	8355	2350	2259	104	0.4974891574	-0.345554003	-0.081102052	0.2329415545	0.2347792691	0.7098484441	0.6192957834	0.0077749048	0.9509653969	0.5070586248	0.0938328845	0.6001572912	0.398908989
3	3305	8547	2948	2655	818	0.8681013997	1.3091137374	-0.489823861	0.2306628048	0.9965764935	1.07933222	0.3028397022	1.3350622223	1.1966755649	1.2513741935	0.7135372695	0.7043898475	1.650498189
4	3578	8754	2773	2385	1173	0.6733898923	-0.909020165	0.3256594606	0.2602865512	1.2774266897	2.0303759189	0.2941883849	0.7878570186	1.056468068	0.1474199096	1.1720177495	0.8167655724	1.284230446
25	2845	8377	1848	2052	1480	0.6954024985	0.0222927481	0.1149275927	0.8880821009	0.2697962147	0.0783328066	0.1785339321	1.9429857801	0.5420037275	-0.480809745	-0.058993942	0.6121006049	-0.65175619
6	3109	8260	1855	2174	988	1.0860262009	-0.658904939	-0.315847551	-0.454101488	2.1249796992	-1.409970464	-0.543168066	1.5938554817	0.3762534786	1.1216311127	0.3843717967	0.5485838909	-0.63710548
7	3264	7492	2465	2250	658	0.4954880114	0.6132638736	-1.117608798	-1.797424452	0.1726063249	0.0444006231	-1.817416042	0.9085531596	0.0997254251	0.5753444568	0.6446812267	0.1316536655	0.639599220
8	3023	8397	2546	2548	823	0.7424294299	0.7376788474	-0.905406707	-0.599941471	0.7543164006	0.1801293571	-0.680450812	0.7078865693	0.6505604232	0.764268592	0.2399420484	0.6229581629	0.80912886
9	3065	9467	2170	2609	986	1.3003489393	0.3181638945	0.5937299996	1.2162220615	1.0087252297	1.5741765628	0.8114460946	2.4387284687	0.6841546947	0.8143448688	0.3104775069	1.2038375134	0.022176454
0	2656	4995	2098	1787	888	-0.380213485	-1.464260544	-0.709377063	-0.746920828	-0.290474915	-2.270528341	-0.613744602	-0.700852132	-0.212895647	-0.262295082	-0.376403505	-1.223912445	-0.128516
1	3787	7077	1775	1865	904	-0.376811537	-0.711087252	-0.041406049	-0.757175202	-0.433401223	-0.391062399	0.5785435255	-1.084043868	-0.123496016	-0.667457685	1.5230156261	-0.093640662	-0.8045421
2	3538	6335	1619	1512	1656	-0.137474474	-1.148081995	0.4516085073	0.1224221928	-0.63564195	1.0171232169	0.3854825493	-1.491300785	-0.271755032	-0.102961474	1.1048411224	-0.496456062	-1.1310436
3	2796	6682	2430	1789	93	-0.308572458	-0.672528893	-0.130599537	-0.594244596	0.1354454847	-0.365613261	-0.011567383	0.2062200937	0.0714062871	-0.999782068	-0.14128531	-0.308077432	0.56634567
4	2660	7944	1918	2004	99	-0.370607984	0.341555966	-0.70790684		-0.024631981	0.4336781728	-0.102178549	0.3254353004	0.089175158	-0.417076802	-0.369685843	0.3770344752	-0.5052490
5	2873	6234	1514	1798	841	-0.816063088	0.561338616	-0.96176523	0.653370879	-1.181620448	0.21311898	-0.797926595	-1.000000849	-0.201512465	0.149696104	-0.011970303	-0.55128678	-1.3508043
6	3762	7064	2039	1810	99	-0.803255753	-0.000842268	-0.273211104		-0.000334508	-0.446673477	-0.118798184	-0.457238448	-0.203733573	-0.118894835	1,4810302341	-0.100698075	-0.25200111
2	2014	6107	2600	9176	0.44	0.511000000	-0.045000060	0.014000000		-0.64400016	0.0150500011	0.4005401511	0.070070505	0.0501400000	0.004675640	-0.1110EE000	0.571070010	1 100401000



#### 階層的クラスタリング (Ward法):標準化 - 実行方法



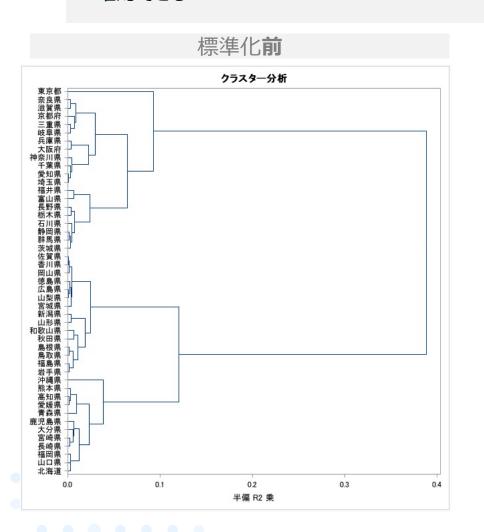
クラスタリング分析結果の画面上において 上部メニュー [タスクの変更] をクリック

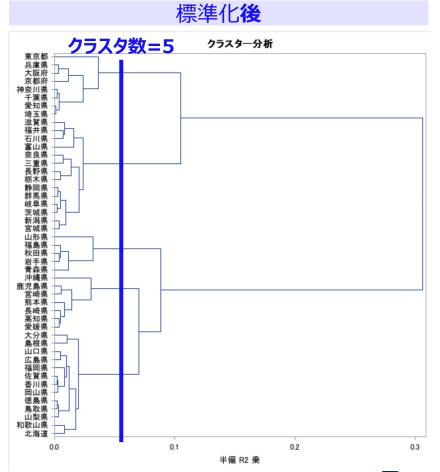




#### 階層的クラスタリング (Ward法):標準化 - 実行結果

標準化を行うことで、全ての変数がフェアに分類に寄与するようになり、分類結果が変わることが確認できる

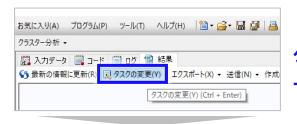




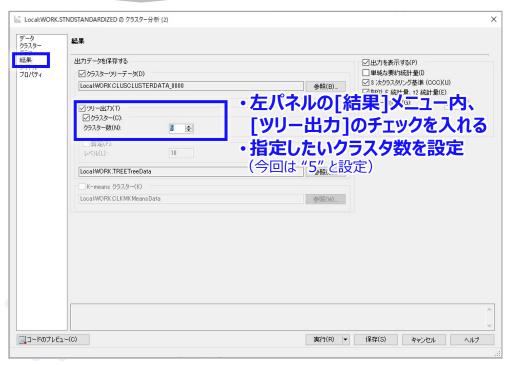


### 指定したクラスタ数の分類結果を出力する方法

- 階層的クラスタリングでは、分析の試行錯誤を重ね、結果の妥当な解釈を行いながら、最適な クラスタ数を決めていく
- 最適なクラスタ数の決定後は、明示的にクラスタ数を指定した分類結果を出力することが可能



クラスタリング分析結果の画面上において 上部メニュー [タスクの変更] をクリック

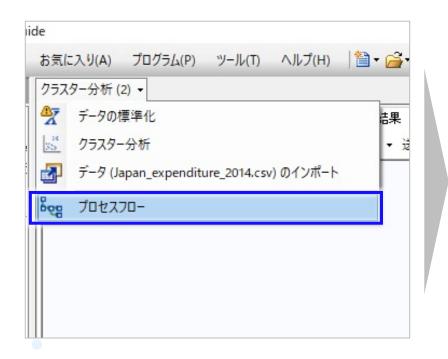


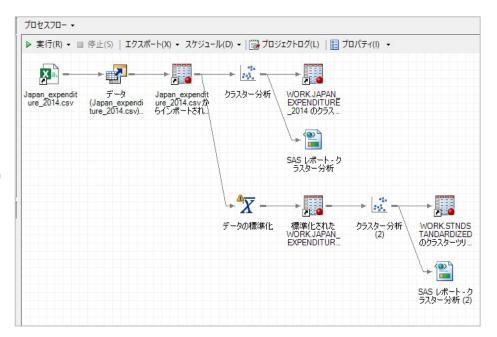
#### 指定したクラスタ数での 分類結果が出力される ⑤ タスクの変更(Y) | □ フィルタと並べ替え(L) □ クエリビルダ( 🔌 都道府県 🔞 CLUSTER 🔌 CLUSNAME CL5 愛知県 CL5 CL11 佐賀県 CL11 鳥取県 OL11 徳島県 OL11 千葉県 CL5 CL5 神奈川県 CL9 茶城田 10 岐阜県 CL9 11 岡山県 OL11 12 群馬県 CL9 13 静岡県 CL9 14 大阪府 CL5 CL5 15 兵庫県 16 山梨県 OL11 17 岩手県 CL6 CL6 18 秋田県 CL7 19 愛媛県 CL7 20 高知県 CL9 21 宮城県 22 新潟県 CL9 CL11 23 広島県 24 山口県 OL11 25 福島県 CL6 26 栃木県 CL9 27 長野県 OL9 28 長崎県 OL7 29 宮崎県 CL7 CL7 30 鹿児島県 31 石川県 CL9 32 福井県 CL9 CL9 33 三重県 CL9 34 奈良県



### (参考) プロセスフローでの確認

• SAS Enterprise Guideでは、「プロセスフロー」によりデータ加工やグラフ観察、モデル構築などの一連の分析プロセスを俯瞰的/反復的に確認可能



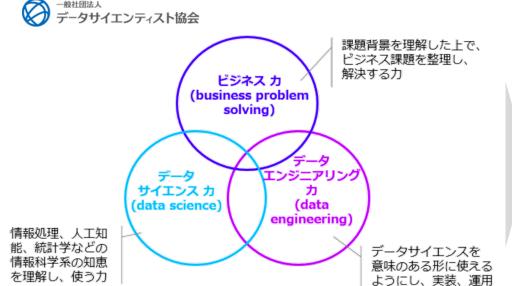


### **Agenda**

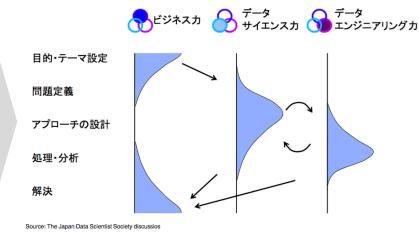
- クラスター分析による分類(2):階層的クラスタリング
  - 階層的クラスタリング(群平均法、重心法、Ward法)のしくみ
  - 各クラスタの解釈、予測への適用方法
  - 樹形図(デンドログラム)とクラスタ数の検討
  - 都道府県データを用いて階層的クラスタリングにより類似地域を分析する
- 今後のデータサイエンス学習に向けて
  - データサイエンティストに求められるスキル
  - 学習リソースの紹介
  - SAS内サンプルデータの紹介と使い方

### データサイエンティストに求められる知識・スキルセット

- 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



#### 課題解決の各フェーズで要求されるスキルセットのイメージ



出展:データサイエンティスト協会資料 http://www.datascientist.or.jp/news/2014/pdf/1210.pdf

できるようにする力

### データサイエンティストの育成モデルの例

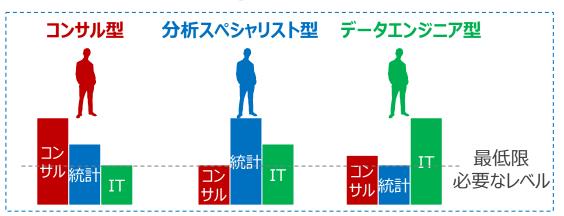
- 理想的には、すべてのスキルを持つ人材を育成できればよいが、元々のバックグラウンドや 経験値を生かした育成を行うのが現実的
- チームでスキルを補い合い、プロジェクトフェーズに応じて、役割分担や協業体制が必要

シニアアナリスト? アナリスト? コン統計 IT

全方位 育成モデル

#### 現実的には・・・

バックグランド ベース 育成モデル



データサイエンティスト?

チームでスキルを補い合う



### 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

# レベル定義 業界代表 棟梁

#### コンサルカ (ビジネスカ)

#### ■ 組織や市場全体にインパクトを出せる

- 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
- 分析を通じオペレーション上の革新が実現でき多変量解析の概念を理解し、活用することがで
- 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 • **モデルを構築**できる
- 組織全体を見渡して、必要なデータの当たりを

#### 統計力(データサイエンスカ)

- 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
- 機械学習、自然言語、画像処理のアルゴリズ △を理解し、適切に活用、問題解決することが できる

#### **IT力**(データTンジニアリングカ)

- 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
- 分析に必要なデータフォーマット、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
- 構造化データ/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

#### 全てのスキルで「独り立ちレベル」以上の習得が理想

#### 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
- ビジネスにおける論理とデータの重要性を認識し 基本統計量 (平均、中央値など) の知識を

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

- SPSS/SAS/R等が使える。指示されなくてもサ 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
- 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
- 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
- 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

#### 見習い

独り立ち

- 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

### 今後のデータサイエンス学習に向けて (スキルアップ編)

- おすすめのオンライン学習サイト
  - Gacco: <a href="https://gacco.org/">https://gacco.org/</a>
    - 国 (総務省) がやっている無料の講座集。データサイエンス学部がある滋賀大の講座が受講可能
  - Coursera: https://www.coursera.org/learn/machine-learning
    - 世界中でもっとも有名かつ人気と言っても過言ではない、スタンフォード大の機械学習講座。無料です
    - 日本語字幕もつけられます
  - 東大松尾研の無料講座: https://gci.t.u-tokyo.ac.jp/
    - 日本のデータサイエンスでもっとも有名とも言える東大松尾先生の無料講座
    - Pythonベースなので、Pythonの学習にも良いと思います
    - 年数回(不定期?)受講生を募集しているので、こまめにチェックしてください
- 今後の分析技術
  - テキスト分析、画像分析
  - スパースモデリング、ベイズ統計、xAI (eXplanable AI)



### 今後のデータサイエンス学習に向けて (腕試し/リアルデータ編)

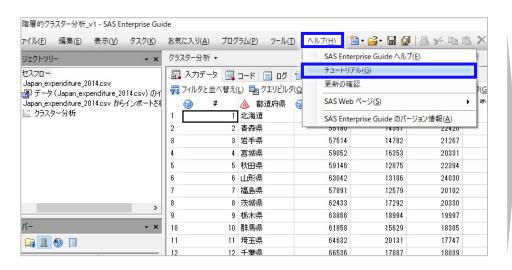
- データ分析コンペ
  - Kaggle: https://www.kaggle.com/
    - Kaggleは本場アメリカのコンペサイト
    - ここで上位に入り込めれば、AI系の有力企業からも注目を浴びると思います
  - **SIGNATE**: <a href="https://signate.jp/">https://signate.jp/</a>
    - "日本版Kaggle" とも言え、規模は小さいですがまずはこちらから挑戦いただくのが良いかもしれません
- プロジェクトベース
  - **AI Quest** (経産省主催、BCG運営の課題解決型データサイエンスコース) <a href="https://aiquest.meti.go.jp/2021/">https://aiquest.meti.go.jp/2021/</a>
- オープンデータを活用した独自分析
  - データカタログサイト: https://www.data.go.jp/
  - 観光統計データ: <a href="https://statistics.jnto.go.jp/">https://statistics.jnto.go.jp/</a>
  - 政府統計e-stat: <a href="https://www.e-stat.go.jp/">https://www.e-stat.go.jp/</a>
  - 医療: NDBオープンデータ: https://www.mhlw.go.jp/stf/seisakunitsuite/bunya/0000177182.html



### SAS Enterprise Guide チュートリアルの活用

• SAS Enterprise Guideの使い方などに困った場合は、オフィシャルのチュートリアルを閲覧すると有益な情報が得られることがある

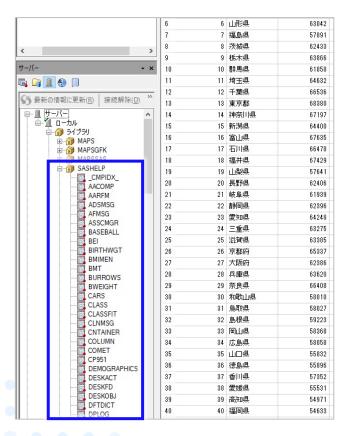
・ 上部メニュー[ヘルプ]から、[チュートリアル]をクリック





### SASHELPデータ (サンプルデータ) の活用

- SAS Enterprise Guide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
- ・ウインドウ左下の[サーバー]のパネルから [サーバー]→[ローカル]→[ライブラリ]→[SASHELP]を選択



#### ▼データセットの例

BASEBALL 1986年のメジャーリーガーの成績データ

BMIMEN 年齢とBMIに関するデータ

BMT 骨髄移植患者の生存期間データ

BIRTHWGT 2003年の乳児死亡率に関するデータ

FAILURE 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

JUNKMAIL 迷惑メールデータ

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



### まとめ

- 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
  - データを標準化することで、さらに精緻なクラスタリングを行うことができた
- 今後のデータサイエンス学習に向けて
  - 無料の有用な学習リソースはたくさん存在
  - 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
  - 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

### アンケートのお願い・ご質問

### 機械学習によるビッグデータ分析の手法

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

### アンケートURL

https://sas.qualtrics.com/jfe/form/SV\_e9byCaaw8q5xob4

・お客様講演会のアーカイブは、2022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

## **End of File**