2024 | 11 | 19 Jue.

問題解決技法入門

3. Data Analysis 3. Cluster Analysis using R

堀田 敬介

クラスター分析とは

どれとどれが似てる?

(同じクラスター?)

- クラスタ分析とは?
 - 複数の対象(もの,変数など)を、その
 属性によって類似度(similarity)をはか
 - り、均質な

 集団(cluster)に

 分類する方

 法の総称



クラスター分析とは

クラスタ分析の種類

- 階層的方法
 - 樹形図(デンドログラム)を作成
 - 目的により高さを決めてクラス
 タリング

- 非階層的方法
 - 予めクラスタ数を決めて
 - (or 決まっていて) クラスタリングを行う





例:3つのクラスタに分類







		3	1	2	3	4	6	6
<i>x</i> ₁	<i>x</i> ₂	1	2	3	5	5	5	3
3	1							
1	2							
2	3							
3	5							
4	5							
6	5							
6	3							

2. 類似度の測定

- <u>距離【間隔尺度】</u>
 - ユークリッド距離
 - ユークリッド平方距離
 - 重み付きユークリッド距離
 - マンハッタン距離
 - ミンコフスキー距離
 - マハラノビス汎距離
- <u>相関【間隔尺度】</u>
 - Pearsonの積率相関係数
 - ベクトル内積
- <u>相関【順序尺度】</u>
 - Spearmanの順位相関係数
 - Kendallの順位相関係数

類似度は尺度により距離や相関で測る (距離:近いほうが類似) (相関:高いほうが類似)

- <u>距離【名義尺度 [0, 1]</u>】
 - 類似比
 - 一致係数
 - Russel-Rao係数
 - Rogers-Tanimoto係数
 - Hamann係数
 - ファイ係数
- <u>変量間類似度【名義尺度】</u>
 - - グッドマン・クラスカルのλ



2. 類似度の測定

氏名

比率尺度

間隔尺度

順序尺度

名義尺度

データと尺度

名義尺度

学籍番号

2

3

量的データ (数値データ)

質的データ

(カテゴリデータ)



比に意味がある(絶対原点が存在する) 例)身長 180cmのAさんは息子(100cm)の1.8倍背が高い

差に意味がある 例)温度 気温20℃より30℃の方が10℃高い

順序関係がある(順序に意味がある) 例) 成績評価 (A>B>C>D)

単なる分類,区別ができる 例) 名前, 性別





2. 類似度の測定

- 個体間類似度
 - ユークリッド距離
 - (cf. *l₂-ノ*ルム)
 - マンハッタン距離
 - (cf. *l₁-ノ*ルム)
 - ミンコフスキー距離
 - (cf. l_p -ノルム) (cf. l_∞ -ノルム)



- マハラノビス汎距離 2変量版 x=(x₁, x₂) $D = \sqrt{\frac{u_1^2 + u_2^2 - 2\rho u_1 u_2}{1 - \rho^2}}$ $u_1, u_2 \downarrow x_1, x_2 O 標準化変量で, u_1 = \frac{x_1 - \mu_1}{\sigma_2}, u_2 = \frac{x_2 - \mu_2}{\sigma_2}$ $\int \mu_1, \mu_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot x_1, x_2 O \mp b \uparrow \sigma_1, \sigma_2 \downarrow \xi + \Lambda \xi h \cdot \xi$

- 2. 類似度の測定
- どうやって<u>類似度を測る</u>か?
 - ・ 例:ユークリッド平方距離



		3	1	2	3	4	6	6
x_1	<i>x</i> ₂	1	2	3	5	5	5	3
3	1		5	5	16	17	25	13
1	2			2	<i>13</i>	18	34	26
2	3				5	8	20	16
3	5					1	9	13
4	5						4	8
6	5							4
6	3							





		3	1	2	3	4	6	6
x_1	<i>x</i> ₂	1	2	3	5	5	5	3
3	1		5	5	16	17	25	13
1	2			2	13	18	34	26
2	3				5	8	20	16
3	5				•	1	9	<i>13</i>
4	5						4	8
6	5							4
6	3							





			3	1	2	3,4	6	6
	x_1	<i>x</i> ₂	1	2	3	5,5	5	3
	3	1		5	5	16,17	25	<i>13</i>
	1	2			2	13,18	34	26
	2	3				5,8	20	16
	3,4	5,5				1	9,4	13,8
	6	5						4
	6	3						

3. クラスタ化の方法

- ・新たなクラスタ生成時の類似度の更新方法
 - クラスタp, クラスタq が一つのクラスタt になる場合, 他のクラスタr との類似度をどう更新する?







			3	1	2	3:4	6	6
	<i>x</i> ₁	<i>x</i> ₂	1	2	3	5:5	5	3
	3	1		5	5	16:17	25	<i>13</i>
	1	2			2	13:18	34	26
	2	3				5:8	20	16
	3:4	5:5				1	9:4	<i>13</i> :8
	6	5						4
	6	3						

3. クラスタ化の方法









						0	
	<i>x</i> ₁	<i>x</i> ₂					
	•		5	5	<i>21.7</i>	25	<i>13</i>
			•	2	20.3	34	26
	2				8. 3	20	16
						8. 3	13.7
0000							4

3. クラスタ化の方法



					.	
	_					
	<i>x</i> ₁	<i>x</i> ₂				
			5:5	<i>21.7</i>	25	<i>13</i>
			2	20.3:8.3	34:20	26:16
					8. 3	13.7
0						4

3. クラスタ化の方法





3. クラスタ化の方法



	_					
	<i>x</i> ₁	<i>x</i> ₂				
			6	<i>21.7</i>	25	<i>13</i>
				20.5	35.3	27.3
					8. 3	13.7
e						4

3. クラスタ化の方法



<i>x</i> ₁	<i>x</i> ₂			
		6	<i>21.7</i>	25:13
			20.5	35.3:27.3
				8.3:13.7
				4

3. クラスタ化の方法









_				
<i>x</i> ₁	<i>x</i> ₂			
		6	<i>21.7</i>	24
			20.5	45
				14.5













3. クラスタ化の方法 14.5





3. クラスタ化の方法 14.5





3. クラスタ化の方法 14.5





 $\frac{2+3}{2+2+3}27 + \frac{2+3}{2+2+3}48 - \frac{3}{2+2+3}14.5$

3. クラスタ化の方法





4. R commanderでクラスター分析

① データの準備:csv ファイル

		算数	理科	国語	英語	社会
	太郎	90	100	70	90	30
	次郎	80	60	70	70	20
CSV	三郎	100	40	30	70	80
	四郎	60	30	40	80	80
	花子	30	60	80	90	90
	寒子	50	60	40	30	60
	湘子	90	100	90	80	70

data-seiseki.csv

② Rの起動:「R x64 X.X.X」を選択

- 注) x64 = 64bit用のプログラム(アプリ), X.X.X = Rのバージョン
- 注)32bit PCの場合は、「R i386 X.X.X」を選択
- 注)起動すると「R Console」が開く、コマンドプロンプト(>)で 「library(Rcmdr)」と打ち[Enter] → R commander が起動



4. R commanderでクラスター分析

③ データの読込(読み込んだファイルの確認)

- ▶ [データセットを表示]ボタンをクリックし, 内容を確認
 - 注1)氏名の項目名が「X」であることを確認(もとのファイルに項目名がな いデータは自動的に「X」となる)
 - ▶ 注2)ケース名(左端)が通し番号(1,2,...,7)となっていることを確認

(ℝ R 2マンダー)	_		\times
ファイル 編集 データ 統計量 グラフ モデル 分布 ツール ヘルプ			
マテータセット: TDataset / データセットの編集 (2) データセットを表示 モデル:	<u>Σ</u> <アクラ	ティブモデル	なし>
Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマークダウン Rスクリプト Rマーク			_
Dataset <- read.table("C:/Users/ header=TRUE, sep=",", na.str n 1 太郎 90 100 70 90 3 RUE) 2 次郎 80 60 70 70 2 3 三郎 100 40 30 70 8 4 四郎 60 30 40 80 8 5 花子 30 60 80 90 9 6 寒子 50 60 40 30 6	3eiseki.	csv″,	~ ~
<			>
出力		実行	
			~

注3)確認後は、必ず「Dataset」の「×」をクリックして「閉じる」こと





 \times



ファイル 編集 データ 統計量 グラフ モデル 分布 ツール ヘルプ	
マデータセット: Dataset アータセットの編集 データセットを表示 モデル: アクティブモデルなし>]
Rスクリプト Rマークダウン R Dat ロ × Dataset <- read.table("C:/Users/ header=TRUE, sep=",", na.strin row.names(Dataset) <- as.charact	
< > >	
出力	

注2)確認後は、必ず「Dataset」の「×」をクリックして「閉じる」こと









▶ 名前をつけて保存する

5. クラスター分析実施上の注意点

- ・クラスター分析の長所
 - 探索的手法なので、データ構造を事前に知らなくてよい
 - あらゆる種類のデータに適用可能:数値・カテゴリー

適用が
 簡単

- ・クラスター分析の短所
 - どんな属性値を選んだらいいのか?
 - どの類似度(距離)測定法を選んだらいいのか?
 - どのクラスタ化更新法を選んだらいいのか?
 - データのスケーリング
 - 結果の解釈が困難な可能性がある



迷ったらとりあえず

「ユークリッド平方距離」

で



- ◆ 田中豊・脇本和昌『多変量統計解析法』現代数学社(1983)
- 河口至商『多変量解析入門Ⅱ』森北出版(1978,2005)
- ◆ 青木繁伸『Rによる統計解析』オーム社(2009)
- ◆ 荒木孝治 『RとRコマンダーではじめる多変量解析』日科技連(2007)
- ◆ 金明哲 『Rによるデータサイエンス』 森北出版(2007)
- 新納浩幸『Rで学ぶクラスタ解析』オーム社(2007)

もっと知りたい人へ

- 関連する経営学科の授業
 - 「基礎統計」(1/2セメ)
 - 「基礎統計演習」(3/4セメ)
 - -「**データ処理応用**」(2/3セメ)
 - 「統計モデル分析」(5セメ)
 - 「**ビッグデータ・AI演習**」(6セメ)

6. Rでクラスター分析

 Rを起動、csvファイルをデータとして読込み –「マイドキュメント(Y:)」の「R」フォルダに保存

		算数	理科	国語	英語	社会
	太郎	90	100	70	90	30
	次郎	80	60	70	70	20
data-seiseki.csv	三郎	100	40	30	70	80
	四郎	60	30	40	80	80
	花子	30	60	80	90	90
	寒子	50	60	40	30	60
	湘子	90	100	90	80	70

csvファイルを読み込み,変数seisekiに代入

> seiseki <- read.csv("Y:/R/data-seiseki.csv", header=T, row.names=1)</p>
※読み込むファイル名 ※1行目にheaderあり ※各行1列目は名前

6. Rでクラスター分析

関数 dist() で距離を計算し, seiseki.dに代入

> seiseki.d <- dist(seiseki, "manhattan")</pre>

※マンハッタン距離("manhattan")を用いて距離を計算している 他の距離を使いたいときは"manhattan"を以下に変更

"euclidean" =ユークリッド距離 "minkowski", p=3 = p=3のミンコフスキー距離 "maximum" = l_∞ノルム (える むげんだい のるむ)

階層クラスター分析をし、結果をseiseki.hcに代入

> seiseki.hc <- hclust(seiseki.d, "ward.D2")</pre>

※ウォード法("ward.D2")を用いてクラスター分析を実施している 他の方法を使いたいときは、"ward.D2"を以下に変更

"single"=最短距離法,	"complete"=最長距離法	
"average"=群平均法,	"centroid"=重心法,	"median"=中央值法

6. Rでクラスター分析

結果をデンドログラム(樹形図)で描画①

> plot(seiseki.hc, hang=-1)

・結果をデンドログラム(樹形図)で描画②

> plot(seiseki.hc)





6. Rでクラスター分析

Tips! 結果(樹形図)をレーダーチャートと比較

- > install.packages("fmsb")
- > library(fmsb)
- > radarchart(seiseki, axistype=2, オプション指定勉強せよ)



6. Rでクラスター分析

【練習】距離とクラスター化の方法,分割数を以下の設定に 従ってクラスター分析をし,樹形図を描き,比較せよ

	距離	クラスター化の方法	分割数	※たくさん計算する
1	ユークリッド距離(euclidean)	最短距離法(single)	4	変数は整理して使
2	ユークリッド距離(euclidean)	最長距離法(complete)	4	例えば、距離は
3	ユークリッド距離(euclidean)	群平均法(average)	4	seiseki.euc <- dist(.
4	ユークリッド距離(euclidean)	重心法(centroid)	4	などとし、ユークリッ
5	ユークリッド距離(euclidean)	中央值法(median)	4	離とマンハッタン距
6	ユークリッド距離(euclidean)	ウォード法(ward.D2)	4	計算した結果を、れ
				易い名則の別変数
	距離	クラスター化の方法	分割数	
1	マンハッタン距離(manhattan)	最短距離法(single)	4	_{西西} 万法] 17加で区。 seiseki.e si <-hclus
2	マンハッタン距離(manhattan)	最長距離法(complete)	4	seiseki.e_cp <-hclu
3	マンハッタン距離(manhattan)	群平均法(average)	4	seiseki.e_av <-hclu
4	マンハッタン距離(manhattan)	重心法(centroid)	4	seiseki.e_ce <-hclu
		由山庙注(modian)	Λ	seiseki.e_m <-hclus
(5)	マンハッタン距離(mannattan)	中大直云(meulan)	4	sojsoki m si c holu

6. Rでクラスター分析

Tips! 画面を分割して, 複数の図を比較する

> par(mfrow=c(2,3))

- ※1画面を2×3に分割

※c(x,y)のx,yに分割したい 行数(x)と列数(y)を指定する

※この命令の後, plotなどで 図を描画すると, 左上から順 に描画されていく

※6個描かれた後,7個目を描 くと,画面がクリアされてまた 左上から順に描画される

※別の分割に変えたい場合 は,変えたい設定でもう一度 実行すれば良い(何度でも変 更可能)



6. Rでクラスター分析

Tips! 画面分割, 複数図描画(前ページの場合の実行例)



> seiseki.e_si <- hclust(seiseki.euc, "single")
> seiseki.e_cp <- hclust(seiseki.euc, "complete")
> seiseki.e_av <- hclust(seiseki.euc, "average")</pre>

> seiseki.e_ce <- hclust(seiseki.euc, "centroid")</pre>

> seiseki.e_m <- hclust(seiseki.euc,"median")</pre>

> seiseki.e_wa <- hclust(seiseki.euc,"ward.D2")</pre>

> par(mfrow=c(2,3))

> plot(seiseki.e_si,hang=-1)

> rect.hclust(seiseki.e_si,k=4,border="red")

> plot(seiseki.e_cp,hang=-1)

> rect.hclust(seiseki.e_cp,k=4,border="red")

> plot(seiseki.e_av,hang=-1)

> rect.hclust(seiseki.e_av,k=4,border="red")

> plot(seiseki.e_ce,hang=-1)

> rect.hclust(seiseki.e_ce,k=4,border="red")

> plot(seiseki.e_m,hang=-1)

> rect.hclust(seiseki.e_m,k=4,border="red")

> plot(seiseki.e_wa,hang=-1)

> rect.hclust(seiseki.e_wa,k=4,border="red")



6. Rでクラスター分析

Tips! たくさんの命令を打つのは大変だし間違えちゃう!

一度にまとめて命令したい!

- まとめて実行したい命令(右)を1つの ファイルに書く、制作には「TeraPad」 や「メモ帳」「秀丸」などのテキストエ ディタを使う
- ② ファイルの種類を「全てのファイル」にし、「ファイル名.R」で保存.このとき、ファイル名は半角アルファベットが良い(例:ファイル名「euc_clust.R」とし「Y:/R/」フォルダに保存)

③ R(R Studio)で以下を実行

> source("Y:/R/euc_clust.R")

※ソースコード「euc_clust.R」内に間違いが なければ全て順に実行される.間違いがあ る場合は、その場所でエラーが出て止まる

【演習】manhattan 距離で同様のファイル 「man_clust.R」をつくり実行しよう seiseki.euc <- dist(seiseki, "euclidean")</pre>

seiseki.e_si <- hclust(seiseki.euc, "single")
seiseki.e_cp <- hclust(seiseki.euc, "complete")
seiseki.e_av <- hclust(seiseki.euc, "average")
seiseki.e_ce <- hclust(seiseki.euc, "centroid")
seiseki.e_m <- hclust(seiseki.euc, "median")
seiseki.e_wa <- hclust(seiseki.euc, "ward.D2")</pre>

par(mfrow=c(2,3))

plot(seiseki.e_si,hang=-1)
rect.hclust(seiseki.e_si,k=4,border="red")
plot(seiseki.e_cp,hang=-1)
rect.hclust(seiseki.e_cp,k=4,border="red")
plot(seiseki.e_av,hang=-1)
rect.hclust(seiseki.e_av,k=4,border="red")
plot(seiseki.e_ce,hang=-1)
rect.hclust(seiseki.e_ce,k=4,border="red")
plot(seiseki.e_m,hang=-1)
rect.hclust(seiseki.e_m,k=4,border="red")
plot(seiseki.e_wa,hang=-1)
rect.hclust(seiseki.e_wa,k=4,border="red")