

2008年7月8日

問題発見技法

6. クラスタ分析

情報学部 堀田敬介

クラスタ分析

Contents

1. クラスタ分析概要
2. 類似度の測定
3. クラスタ間の近さの決定
4. クラスタ分析の実施[SPSS, 手計算]
5. クラスター分析実施上の注意点
6. 演習: やってみよう!

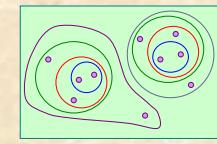
1. クラスタ分析概要

- クラスタ分析とは?
 - ・複数の対象(もの、変数など)を、
類似度(similarity)を定義し、
 - ・均質な**集団(cluster)**に分類する
 - 方法の総称

1. クラスタ分析概要

● 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(or決まっていて)、
クラスタリングを行う



1. クラスタ分析概要

- 例: x_1

どうやって
クラスタ間の近さ
を決めるのか
例) クラスタ(G,B)とクラ
スタ(D)の近さ?

どうやって
類似度(距離)を
測るのか
例) CとEの類似度?

補足: 類似度は相関で測る場合もある
(距離: 近いほうが類似している)
(相関: 高いほうが類似している)

2. 類似度の測定

● 距離, 間隔尺度

- 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
- ミンコフスキイ距離
- マハラノビス汎距離

● 相関, 間隔尺度

- Pearsonの積率相関係数
- ベクトル内積

● 相関, 順序尺度

- Spearmanの順位相関係数
- Kendallの順位相関係数

● 距離, 名義尺度 [0, 1]

- 類似比
- 一致係数
- Russel-Rao係数
- Rogers-Tanimoto係数
- Hamann係数
- ファイ係数

● 变量間類似度, 名義尺度

- 平均平方根一致係数
- グッドマン・クラスクアルの係数

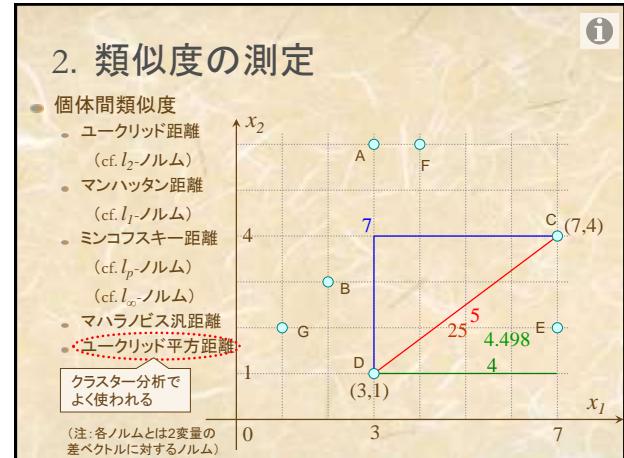


2. 類似度の測定

- データと尺度

学籍番号	氏名	性別	生年月日	身長	体重	問題発見技法成績	...
1	文教太郎	男	1987.5.6	175cm	69kg	B	...
2	湘南花子	女	1988.1.4	163cm	48kg	AA	
3	

 - 比率尺度**: 比に意味がある(絶対原点が存在)
例) 身長 180cmのAさんは息子(100cm)の1.8倍高い
 - 間隔尺度**: 差に意味がある
例) 気温20°Cより30°Cの方が10°C高い
 - 順序尺度**: 順序関係がある
例) 成績評価 (A > B > C > D)
 - 名義尺度**: 単なる分類
例) 名前、性別



2. 類似度の測定

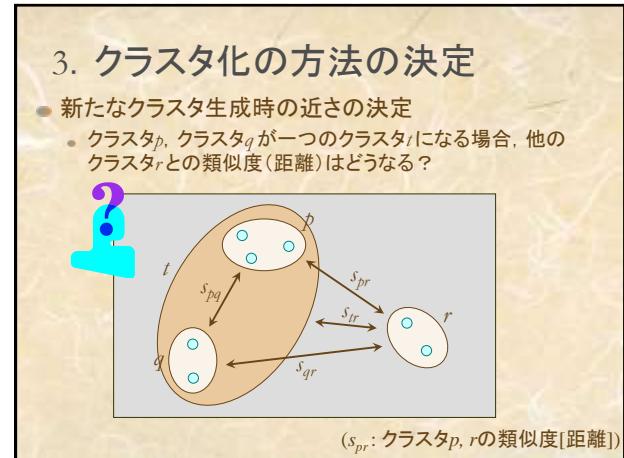
- 個体間類似度
 - ユークリッド距離
(cf. L_2 -ノルム)
 - マンハッタン距離
(cf. L_1 -ノルム)
 - ミンコフスキードイツ距離
(cf. L_p -ノルム)
(cf. L_∞ -ノルム)
 - マハラノビス汎距離
マハラノビス汎距離(2変量 x_1, x_2 版)

左側の対象内での、A-B間距離と右側の対象内でのA-B間距離が異なる！(ユークリッド距離などでは同じ)

$D = \sqrt{\frac{u_1^2 + u_2^2 - 2\rho u_1 u_2}{1 - \rho^2}}$

ただし、 u_1, u_2 は x_1, x_2 の標準化変量で、 $u_1 = \frac{x_1 - \mu_1}{\sigma_1}, u_2 = \frac{x_2 - \mu_2}{\sigma_2}$

また、 μ_1, μ_2 はそれぞれ、変量 x_1, x_2 の平均。
 σ_1, σ_2 は x_1, x_2 の標準偏差。 ρ は x_1, x_2 の相関係数

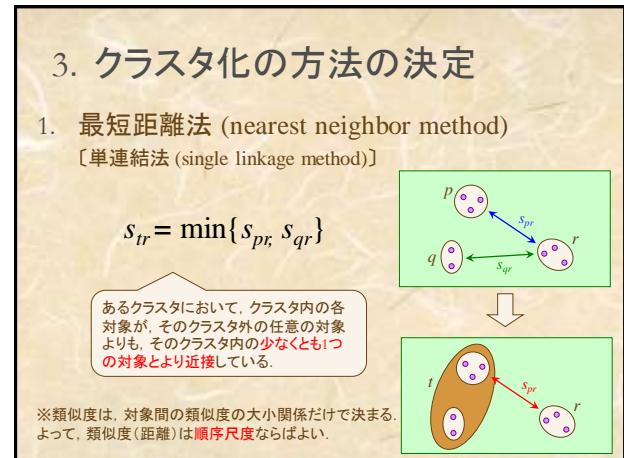


3. クラスタ化の方法の決定

- クラスタ間の近さ決定方法

(事前にクラスタ数を決める必要はない方法群)

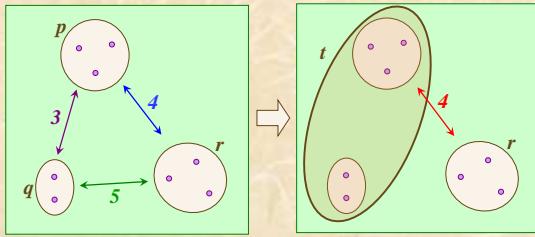
 - 最短距離法 (nearest neighbor method)
 - 最長距離法 (furthest neighbor method)
 - 群平均法 (group average method)
 - 重心法 (centroid method)
 - 中央値法 (median method)
 - ウォード法 (Ward method)



3. クラスタ化の方法の決定

1. 最短距離法

$$s_{tr} = \min\{s_{pr}, s_{qr}\}$$



3. クラスタ化の方法の決定

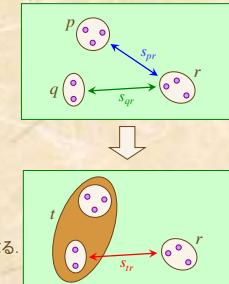
2. 最長距離法 (furthest neighbor method)

[完全連結法 (complete linkage method)]

$$s_{tr} = \max\{s_{pr}, s_{qr}\}$$

あるクラスタにおいて、クラスタ内の全ての対象が、そのクラスタ外の任意の対象との距離よりも常に近接している。

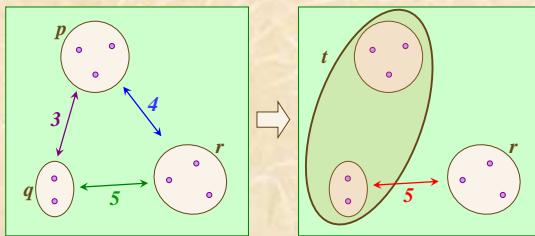
※類似度は、対象間の類似度の大小関係だけで決まる。よって、類似度(距離)は順序尺度ならばよい。



3. クラスタ化の方法の決定

2. 最長距離法

$$s_{tr} = \max\{s_{pr}, s_{qr}\}$$



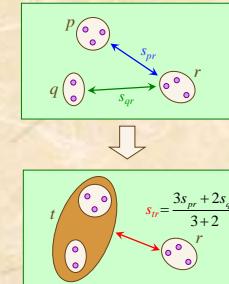
3. クラスタ化の方法の決定

3. 群平均法 (group average method)

$$s_{tr} = \frac{n_p s_{pr} + n_q s_{qr}}{n_p + n_q}$$

(n_p : クラスタ P に含まれる対象数)

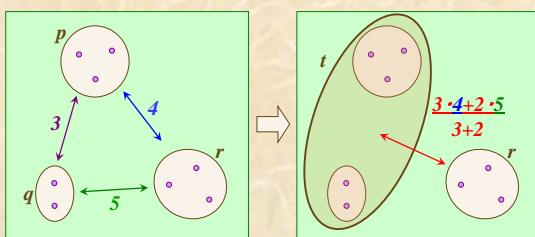
※類似度は、間隔尺度ならばよい。



3. クラスタ化の方法の決定

3. 群平均法

$$s_{tr} = \frac{n_p s_{pr} + n_q s_{qr}}{n_p + n_q}$$



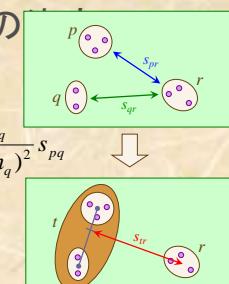
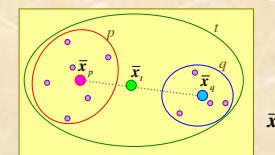
3. クラスタ化の方法の決定

4. 重心法 (centroid method)

$$s_{tr} = \frac{n_p}{n_p + n_q} s_{pr} + \frac{n_q}{n_p + n_q} s_{qr} - \frac{n_p n_q}{(n_p + n_q)^2} s_{pq}$$

(n_p : クラスタ P に含まれる対象数)

※導出過程 (tex-file参照) より、類似度 s_{tr} は ユークリッド平方距離の時のみ妥当。



3. クラスタ化の方法の決定

4. 重心法 $s_{tr} = \frac{n_p}{n_p + n_q} s_{pr} + \frac{n_q}{n_p + n_q} s_{qr} - \frac{n_p n_q}{(n_p + n_q)^2} s_{pq}$

3. クラスタ化の方法の決定

5. 中央値法 (median method)

$$s_{tr} = \frac{1}{2} s_{pr} + \frac{1}{2} s_{qr} - \frac{1}{4} s_{pq}$$

(重心法の簡易版、重心ではなく中央値を取る。
よって、重心法で $n_p := 1, n_q := 1$ に相当する)
※導出過程(重心法参照)より、類似度 S_{tr} は ユークリッド平方距離の時のみ妥当。

3. クラスタ化の方法の決定

5. 中央値法 $s_{tr} = \frac{1}{2} s_{pr} + \frac{1}{2} s_{qr} - \frac{1}{4} s_{pq}$

3. クラスタ化の方法の決定

6. ウオード法 (Ward method)

$$s_{tr} = \frac{n_p + n_r}{n_t + n_r} s_{pr} + \frac{n_q + n_r}{n_t + n_r} s_{qr} - \frac{n_r}{n_t + n_r} s_{pq}$$

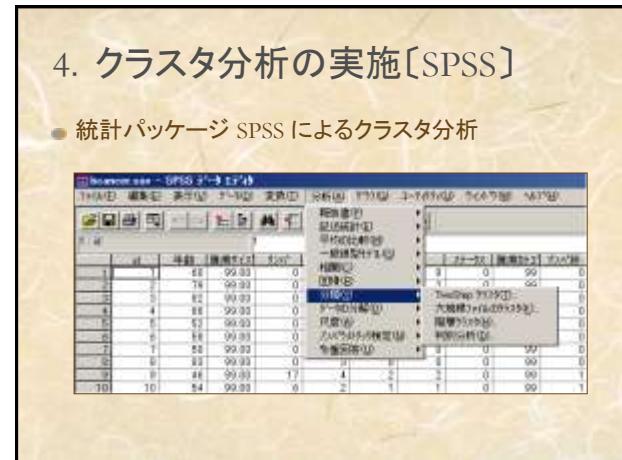
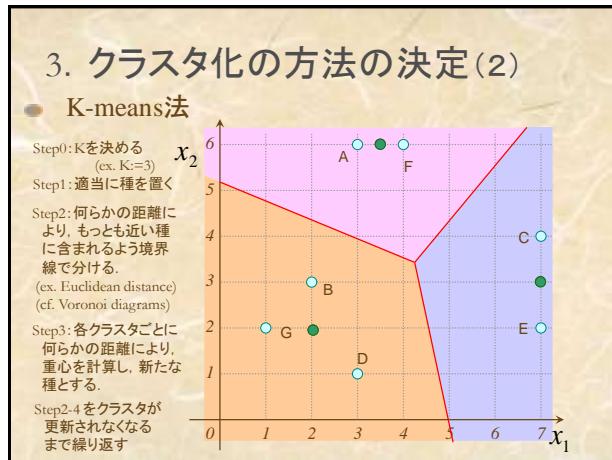
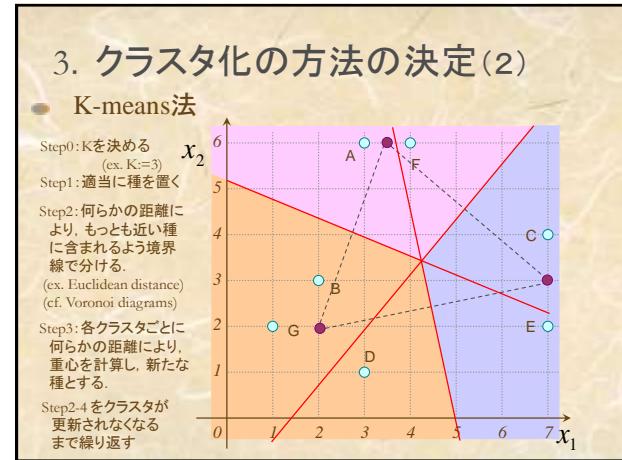
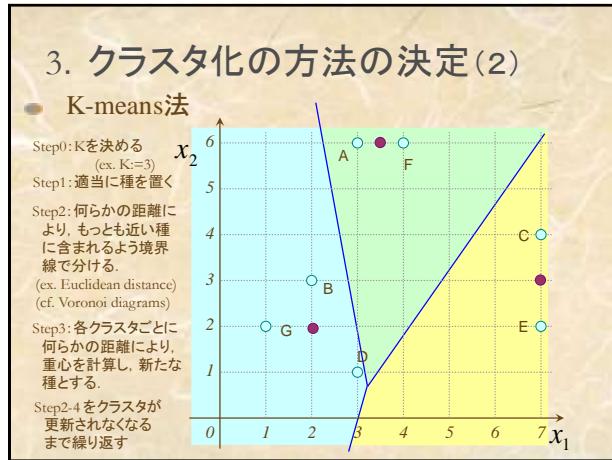
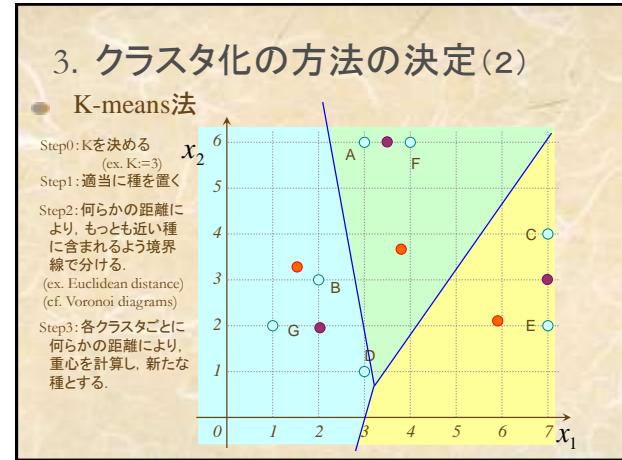
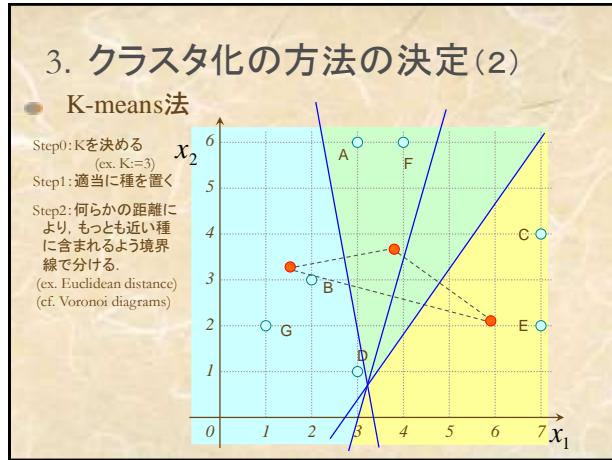
※導出過程(tex-file参照)より、類似度 S_{tr} は ユークリッド平方距離の時のみ妥当。

3. クラスタ化の方法の決定

6. ウオード法 $s_{tr} = \frac{n_p + n_r}{n_t + n_r} s_{pr} + \frac{n_q + n_r}{n_t + n_r} s_{qr} - \frac{n_r}{n_t + n_r} s_{pq}$

3. クラスタ化の方法の決定(2)

- クラスタ間の近さ決定方法(非階層的方法)
 - K-means法
 - 事前にクラスタ数をKとしてクラスタリングを行う。



4. クラスタ分析の実施[SPSS]

- 統計パッケージ SPSS によるクラスタ分析
 - 「クラスタ化の方法」の選択

4. クラスタ分析の実施[SPSS]

- 統計パッケージ SPSS によるクラスタ分析
 - 「クラスタ間距離測定法」の選択

〔佐々木作成スライド(2003堀田研ゼミ)〕

4. クラスター分析の実施[例題]

ある部屋に次の8人(①~⑧)が図7-1の様に座った。
座り方は任意なので、似たもの同士、親しい同士ほど接近して座るものと思われる。座った位置のデータは表7-1である。
このデータをもとに8人を分類しデンドログラムをつくる。

図7-1

表7-1

No.	A	B	C	D	E	F	G
1		4		1			
2		2		4			
3		5		7			
4		5		2			
5		3		4			
6		6		5			
7		2		6			
8		7		2			

〔佐々木作成スライド(2003堀田研ゼミ)〕

8人相互間のユークリッド平方距離を表7-1から求める。

表7-2

	A	B	C	D	E	F	G
A		13	37	2	10	20	29
B	13		8	10	4	20	
C	25	5		10	29		
D		10	25				
E			17	10			
F				5			
G					10		
						20	
							29
							10

図7-1

たて座標
0
よこ座標

表7-3

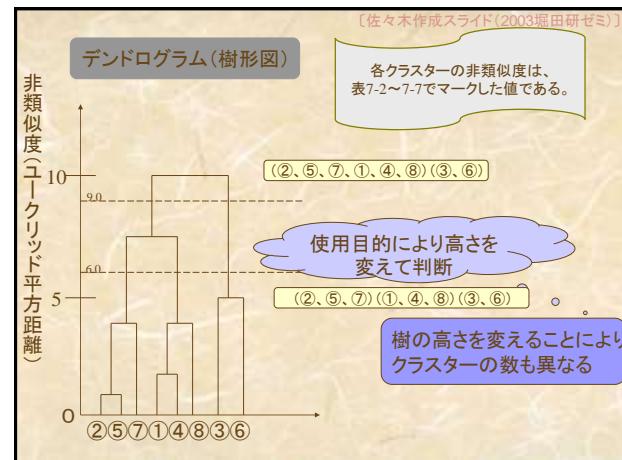
	A	B	C	D	E	F	G
A	10	37	2	20	29	10	
B	13		8	10	4	20	
C	25	5		10	29		
D		10					
E			17	10			
F				5			
G					10		
						20	
							29
							10

表7-4

	A	B	C	D	E	F	G
A		8	25	10	25	4	
B	13		10	4	20		
C		5		10	29		
D			17	10			
E				39			
F					10		
G						20	
							29
							10

図7-1

たて座標
0
よこ座標



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

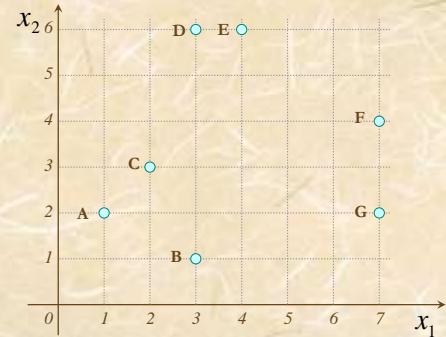
- 探索的手法：データの構造を事前に知らなくてよい
- あらゆる種類のデータに適用可能：数値・カテゴリー
- 適用が簡単

● クラスター分析の短所

- 類似度（距離）測定法の選択が困難の可能性
- クラスタ化法の選択が困難の可能性
- 非階層的手法の場合、事前に決定するクラスタ数の決定が困難の可能性
- 結果の解釈が困難の可能性

6. やってみよう！

- 演習：類似度をユークリッド平方距離、クラスタ化を最短距離法でクラスター分析しよう！



参考文献

- 田中豊ほか『多変量統計解析法』現代数学社
- 河口至商『多変量解析入門 II』森北出版
- 浅利英吉ほか『パソコンによるデータマイニング』日刊工業新聞社
- 東大教養学部統計学教室編『統計学入門』東京大学出版会
- M.J.A.ペリーほか『データマイニング手法』海文堂