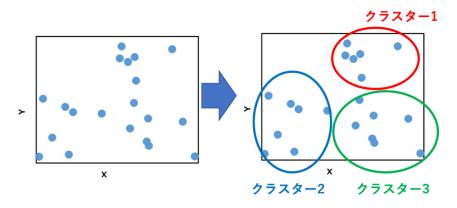
クラスター分析できる

クラスタ―分析を<mark>最短距離法、最長距離法、群平均法</mark>の3手法の違いを理解しながらマスターできる!

【1】クラスター分析とは

(1) データをクラスターで分ける

データ群をある規模のクラスターで分類することですね。下図のようなイメージが簡単にできますよね!



ここで問題になるのが、<mark>どうやってクラスターに分類するの?</mark>

(2) クラスター分析の主な2つの手法

よくあるのが、

- 1. 階層的方法(手計算で考えて解く方法)←これを解説!
- 2. 非階層的方法(計算機で解く方法)

本来は、非階層的方法で、計算機とプログラムを使って解きたいですが、<mark>何を解いているかがわからないので、手計算で理解できる階層的方法を使ってクラスター分析を理解しましょう。</mark>

【2】最短距離法、最長距離法、群平均法とは

階層的方法はさらに3つの方法に分類できます。<mark>比較しながら3手法をマスターしましょう!</mark>

- 1. 最短距離法(最も基本的)
- 2. 最長距離法
- 3. 群平均法

まずは、最短距離法でクラスター分類して、結果を可視化して納得いけばいいですが、 結果がいまいちな場合は、最長距離法、群平均法を使っていきます。

(1) 最短距離法

クラスターに含まれる対象の対の中で、<mark>最短距離なもの</mark>を選びます。式で書くと

$$d(C_i \cup C_j, C_k) = min(d(C_i, C_k), d(C_j, C_k))$$

「min」から最短とわかればOKです。

(2) 最長距離法

クラスターに含まれる対象の対の中で、<mark>最長距離なもの</mark>を選びます。式で書くと

$$d(C_i \cup C_j, C_k) = max(d(C_i, C_k), d(C_j, C_k))$$

「max」から最長とわかればOKです。

(3) 群平均法

最短でも最長でもなく、平均的な値で定義したい場合に使います。式で書くと

$$d(C_i \cup C_j, C_k) = \frac{n_i imes d(C_i, C_k) + n_j imes d(C_j, C_k)}{n_i + n_j}$$

「平均」を計算しているとわかればOKです。

【3】クラスター分析の解法

- (1) 共通の解き方(最短距離法、最長距離法、群平均法)
- 3つのステップがあります。
- 1. 全手法とも、最初は<mark>最短距離なペア</mark>でクラスターを作る
- 2. 手法別にクラスター間距離を計算
- 3. クラスターを合体
- の3ステップを全データが分類し終わるまで繰り返します。



特に注意が必要なのは、

|最長距離法、群平均法でも、最初は最短距離なペアを見つける点に注意しましょう。

(2) データ事例

【事例】5つのデータがあり、それぞれの距離がわかっている。

- (1)最短距離法
- (2)最長距離法
- (3)群平均法
- を使って、それぞれクラスター分析せよ。

_	А	В	С	D	E
А	_	_	_	_	_
В	31.6	_	_	_	_
С	20	51	_	_	_
D	31.6	28.3	42.4	_	_
Е	31.6	63.2	14.1	56.6	_

【4】最短距離法、最長距離法、群平均法を比較しながら解く

分類は3回実施しますので、丁寧に解説します。

(1) 分類 1 回目

①1 回目 step1

最短距離なペアを見つけましょう。CとEの14.1が最短ですね。見ればわかる!

	4=		1201	н	>-	۱
最	뻤	ĭΕ	国	Ħ	77	

step1 クラスタ選び 最短								
	Α	В	С	D	Е			
Α	-	1	-	-	-			
В	31.6	1	-	-	-			
С	20	51	-	-	-			
D	31.6	28.3	42.4	-	-			
Е	31.6	63.2	14.1	56.6	-			

最長距離法

step1 クラスタ選び 最短									
	Α	В	С	D	E				
Α	-	-	1	-	-				
В	31.6	-	-	-	-				
С	20	51	-	-	-				
D	31.6	28.3	42.4	-	-				
Ε	31.6	63.2	14.1	56.6	-				

群平均法

step1 クラスタ選び 最短								
	Α	В	С	D	E			
Α	-	-	-	-	-			
В	31.6	-	-	-	-			
С	20	51	-	-	-			
D	31.6	28.3	42.4	-	-			
Е	31.6	63.2	14.1	56.6	-			

②1 回目 step2

CE が1つのクラスターになったので、

- ●AとCEクラスター
- ●BとCEクラスター
- ●Dと CE クラスター

との距離を最短距離法、最長距離法、群平均法で解きます。

最短距離法

step2 CEとの距離(最短)								
	Α	A B C D E						
Α	-	-	-	-	-			
В	31.6	-	-	-	-			
С	20	51	-	-	-			
D	31.6	28.3	42.4	-	-			
Е	31.6	63.2	14.1	56.6	-			

最長距離法

step2 CEとの距離(最長)								
	Α	В	С	D	E			
Α	-	-	-	-	-			
В	31.6	-	-	-	-			
С	20	51	-	-	-			
D	31.6	28.3	42.4	-	-			
Е	31.6	63.2	14.1	56.6	-			

群平均法

step2 CEとの距離(群平均)								
	Α	В	С	D	Е			
Α	-	-	-	-	-			
В	31.6	-	-	-	-			
С	25.8	57.1	-	-	-			
D	31.6	28.3	49.5	-	-			
E	25.8	57.1	49.5		-			

●最短距離法では、

- ・A と CE クラスター⇒ 黄色の 20 と 31.6 から <mark>20</mark> を選択
- ・BとCEクラスター⇒ 緑色の51と63.2から<mark>51</mark>を選択
- ・Dと CE クラスター⇒ 青色の 42.4 と 56.6 から <mark>42.4</mark> を選択

●最長距離法では、

- ・A と CE クラスター⇒ 黄色の 20 と 31.6 から <mark>31.6</mark> を選択
- ・BとCEクラスター⇒ 緑色の51と63.2から63.2を選択
- ・Dと CE クラスター⇒ 青色の 42.4 と 56.6 から <mark>56.6</mark> を選択

●群平均では、

- ・A と CE クラスター 黄色の 20 と 31.6 から平均 $1/2 \times (20+31.6) = \frac{25.8}{25.8}$ を選択
- ・BとCEクラスター⇒ 緑色の51と63.2から平均1/2×(51+63.2)=57.1を選択
- ・D と CE クラスター \Rightarrow 青色の 42.4 と 56.6 から平均 $1/2 \times (42.4+56.6) = 49.5$ を選択

③1 回目 step3

step2の計算結果を反映します。

最短距離法

step3 クラスター合体(最短)								
	Α	В	CE	D				
Α	-	-	-	-				
В	31.6	-	-	-				
CE	20	51	-	-				
D	31.6	28.3	42.4	-				

最長距離法

step3 クラスター合体(最長)									
	Α	В	CE	D					
Α	-	-	-	-					
В	31.6	-	-	-					
CE	31.6	63.2	-	-					
D	31.6	28.3	56.6	-					

群平均法

step3 クラスター合体(群平均)								
	Α	В	CE	D				
Α	-	-	-	-				
В	31.6	-	-	-				
CE	25.8	57.1	-	-				
D	31.6	28.3	49.5	-				

ここで1回目が終了です。3手法の違いが見えましたね。2回目も同様に解けます!

(2) 分類 2 回目

①2 回目 step1

最短距離なペアを見つけましょう。

- ●最短距離法では、A-CE 間の 20
- ●最長距離法では、B-D間の28.3
- ●群平均法では、A-CE 間の 25.8

が最短ですね。見ればわかるけど、候補と距離の数字が手法によって変わっていますね。

最短距離法

HXVIIIIIIIII								
step1 クラスタ選び 最短								
	Α	В	CE	D				
Α	-	-	-	-				
В	31.6	-	-	-				
CE	20	51	-	-				
D	31.6	28.3	42.4	-				

最長距離法

step1 クラスタ選び 最短						
	Α	В	CE	D		
Α	-	-	-	-		
В	31.6	-	-	-		
CE	31.6	63.2	-	-		
D	31.6	28.3	56.6	-		

群平均法

step1 クラスタ選び 最短						
	Α	В	CE	D		
Α	-	-	-	-		
В	31.6	-	-	-		
CE	25.8	57.1	-	-		
D	31.6	28.3	49.5	-		

② 2 回目 step2

- ●最短距離法では、ACE クラスターと B,D との距離
- ●最長距離法では、 $A \subset CE$ クラスターと BD クラスターとの距離
- ●群平均法では、ACE クラスターと B,D との距離

との距離を最短距離法、最長距離法、群平均法で解きます。ここが一番難しい所です!

最短距離法

step2 A,CEとの距離(最短)					
	Α	В	D		
Α	-	-	-	-	
В	31.6	-	-	-	
CE	20	51	-	-	
D	31.6	28.3	42.4	-	

最長距離法

step2 BDとの距離(最長)						
	Α	A B		D		
Α	-	-	-	-		
В	31.6	-	-	-		
CE	31.6	63.2	-	-		
D	31.6	28.3	56.6	-		

群平均法

step2 A,CEとの距離(群平均)						
	Α	A B CE				
Α	-	-	-	-		
В	48.6	-	-	-		
CE	25.8	48.6	-	-		
D	43.5	28.3	43.5	-		

●最短距離法では、

- ・ACE クラスターと B⇒ 橙色の 31.6 と 51 から <mark>31.6</mark> を選択
- ・ACE クラスターと D⇒ 緑色の 31.6 と 42.4 から <mark>31.6</mark> を選択

●最長距離法では、

- ・A と BD クラスター⇒ 灰色の 31.6 と 31.6 から 31.6 を選択
- ・BD クラスターと CE クラスター⇒ 紫色の 63.2 と 56.6 から 63.2 を選択

●最短距離法では、

- ・ACE クラスターと B⇒ 橙色から 1/3×31.6+2/3×57.1=48.6 を選択
- ・ACE クラスターと D \Rightarrow 緑色の $1/3 \times 31.6 + 2/3 \times 49.5 = 43.5$ を選択

③ 2 回目 step3

step2の計算結果を反映します。

最短距離法

step3 クラスター合体(最短)						
	ACE	В	D			
ACE	-	-	-			
В	31.6	-	-			
D	31.6	28.3	-			

最長距離法

AX						
step3 クラスター合体(最長)						
	Α	BD	CE			
Α	-	-	-			
BD	31.6	-	-			
CE	31.6	62.2	1			

群平均法

step3 クラスター合体(群平均)						
	ACE	В	D			
ACE	-	-	-			
В	48.6	-	-			
D	43.5	28.3	-			

ここで2回目が終了です。3手法の違いが見えましたね。3回目も同様に解けます!

(3) 分類 3 回目

①3 回目 step1

最短距離なペアを見つけましょう。

- ●最短距離法では、B-D 間の 28.3
- ●最長距離法では、A-CE 間の 31.6
- ●群平均法では、B-D 間の 28.3

が最短ですね。見ればわかるけど、候補と距離の数字が手法によって変わっていますね。

最短距離法

step1 クラスタ選び 最短						
	ACE	В	D			
ACE	-	-	-			
В	31.6	-	-			
D	31.6	28.3	-			
				_		

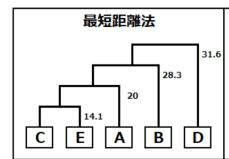
最長距離法

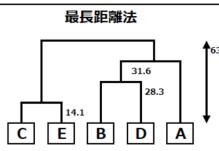
step1 クラスタ選び 最短						
	Α	BD	CE			
Α	-	ı	-			
BD	31.6	-	-			
CE	31.6	62.2	-			

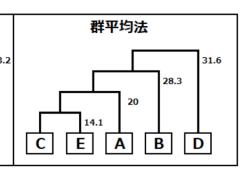
群平均法

	step1 クラスタ選び 最短						
	ACE	В	D				
ACE	-	-	-				
В	31.6	-	-				
D	31.6	28.3	-				

で、ここで、分類が完了したので、結果を比較すると







となりました。手法間で結果が異なりますが、実データと比較してどれを使うかを吟味すればOKです。

以上、「クラスター分析ができる」を解説しました。