

# NAG アドインとクラスター分析

このドキュメントでは、「Excel NAG 統計解析アドイン」（以下「NAG アドイン」と略称）を用いて行うことができる（階層的な）クラスター分析の方法を簡単に説明し、分析の結果として樹形図を作成するまでの具体的な操作例を示します。

## 目 次

1	クラスター分析とは？ .....	1
2	距離の定義.....	2
2.1	個体間の距離.....	2
2.2	クラスター間の距離.....	3
3	クラスター分析の流れ .....	4
4	「NAG アドイン」を用いた分析例.....	5
4.1	クラスター分析（CLUSTER） .....	6
4.2	樹形図の作成（NAG_DendoPlotData） .....	11
4.3	クラスター分け（GROUPS_FROM_CLUSTER） .....	14
4.4	結果の検討 .....	18
5	参考文献.....	19

## 1 クラスター分析とは？

クラスター分析とは、対象（個体の場合もあるし、変数の場合もある）の集合に対して、対象間の測度（非類似度（距離など）又は類似度（相関係数など））※を定義することで、類似した対象が同じ部分集合（クラスターと呼ぶ）になるように分類する方法の総称です。方法を大別すると、階層的な方法と非階層的な方法に分けられます。

※ 非類似度はその値が小さいほど類似性が高く、類似度はその値が大きいほど類似性が高いと考える。

このドキュメントでは、前者（階層的な方法）を取り上げます。以下「クラスター分析」と言った場合は「階層的な方法」のことを意味し、対象としては個体を、測度としては距離を想定します。また、一つの個体から成るクラスターを単に「個体」と呼ぶこともあります（文脈上混乱は生じないと思います）。

階層的な方法は、結果として樹形図（樹状図、デンドログラムとも言う）が得られる方法で、事前にクラスター数を定めることなく、個体間の階層的構造を求めることができます。この樹形図を適当な高さ（距離）で切断することにより個体をいくつかのクラスターに分け、そして各クラスターに含まれる個体を調べることでよりクラスター（クラスター間）の特徴を把握することが分析の目的です。

（この「いくつかのクラスターに分けるか」及び「クラスターの特徴は何か」については、数学的に定まった方法があるわけではなく、分析者自身の判断に委ねられます。）

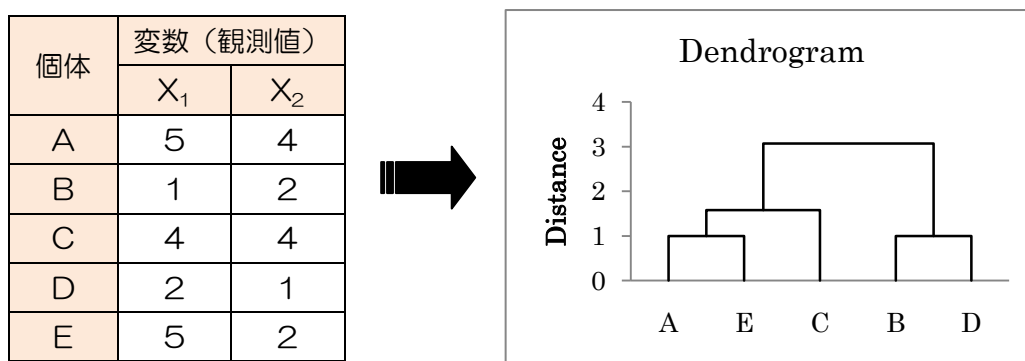


図1 （階層的な）クラスター分析のイメージ

## 2 距離の定義

個体間（個体と個体の間）及びクラスター間（クラスター（特に個体の場合もある）とクラスターの間）の距離の定義を説明します。

### 2.1 個体間の距離

$n$  個の各個体が  $p$  変数の観測値  $\{x_{ji}, j = 1, \dots, n, i = 1, \dots, p\}$  を持っているとし  
ます。個体  $j$  と個体  $k$  の間の距離を  $d_{jk}$  と書いて、以下の様に定義します。

$$d_{jk} = \left\{ \sum_{i=1}^p D(x_{ji}/s_i, x_{ki}/s_i) \right\}^{\alpha}$$

ここで、 $s_i$  は  $i$  番目の変数の標準化、 $D$  は適当な関数です。

「NAG アドイン」では、以下の3つの距離の定義（関数  $D$ ）が利用できます。

- (a) ユークリッド距離：  $D(u, v) = (u - v)^2, \alpha = \frac{1}{2}$
- (b) ユークリッド平方距離：  $D(u, v) = (u - v)^2, \alpha = 1$
- (c) 絶対距離（city block 距離）：  $D(u, v) = |u - v|, \alpha = 1$

また、以下の3つの標準化が利用できます。

(a) 標準偏差：  $s_i = \sqrt{\sum_{j=1}^n (x_{ji} - \bar{x})^2 / (n - 1)}$

(b) 範囲（レンジ）：  $s_i = \max(x_{1i}, x_{2i}, \dots, x_{ni}) - \min(x_{1i}, x_{2i}, \dots, x_{ni})$

(c) 標準化しない：  $s_i = 1$

## 2.2 クラスタ間の距離

それぞれ  $n_i, n_j, n_k$  個の個体から成る  $i, j, k$  の3つのクラスターがあり、各クラスター間の距離は  $d_{ij}, d_{ik}, d_{jk}$  として与えられているとします。

(特に  $i, j, k$  が個体の場合、 $n_i = 1, n_j = 1, n_k = 1$  とします。)

$j$  クラスタと  $k$  クラスタが統合されて  $jk$  クラスタと成った時に、 $i$  クラスタと  $jk$  クラスタの間の距離を  $d_{i,jk}$  と書くとする、**「NAG アドイン」** では、以下の6つの距離の定義が利用できます。

(このクラスタ間の距離の定義それぞれがクラスタ分析の方法に対応します。)

(a) 最短距離法:  $d_{i,jk} = \min(d_{ij}, d_{ik})$

(b) 最長距離法:  $d_{i,jk} = \max(d_{ij}, d_{ik})$

(c) 群平均法:  $d_{i,jk} = \frac{n_j}{n_j+n_k} d_{ij} + \frac{n_k}{n_j+n_k} d_{ik}$

(d) 重心法:  $d_{i,jk} = \frac{n_j}{n_j+n_k} d_{ij} + \frac{n_k}{n_j+n_k} d_{ik} - \frac{n_j n_k}{(n_j+n_k)^2} d_{jk}$

(e) メディアン法:  $d_{i,jk} = \frac{1}{2} d_{ij} + \frac{1}{2} d_{ik} - \frac{1}{4} d_{jk}$

(f) 最小分散法:  $d_{i,jk} = \{(n_i + n_j) d_{ij} + (n_i + n_k) d_{ik} - n_i d_{jk}\} / (n_i + n_j + n_k)$

### 3 クラスター分析の流れ

クラスター分析は、以下の様なステップで行います。

- ① 個体間の距離の定義、及び、クラスター間の距離の定義（クラスター分析の方法）を選択する。
- ② 個体間（全ての組合せ）の距離を計算する。  
（この時点で、個体それぞれを（最小の）クラスターと見なす。）
- ③ 距離が最小となる個体と個体を統合して1つのクラスターとする。
- ④ 新しく形成されたクラスターと既存のクラスター（特に個体の場合もある）の間の距離を計算し、（全ての組合せから）距離が最小のものを統合して新たなクラスターとする。  
これを全てのクラスターが統合され単一のクラスター（全体集合）になるまで繰り返す。
- ⑤ クラスターの統合過程に基づいて樹形図を作成する。
- ⑥ 樹形図を適当な高さ（距離）で切断し、個体をいくつかのクラスターに分ける。
- ⑦ 各クラスターに含まれる個体を調べ、クラスター（クラスター間）の特徴を把握する。

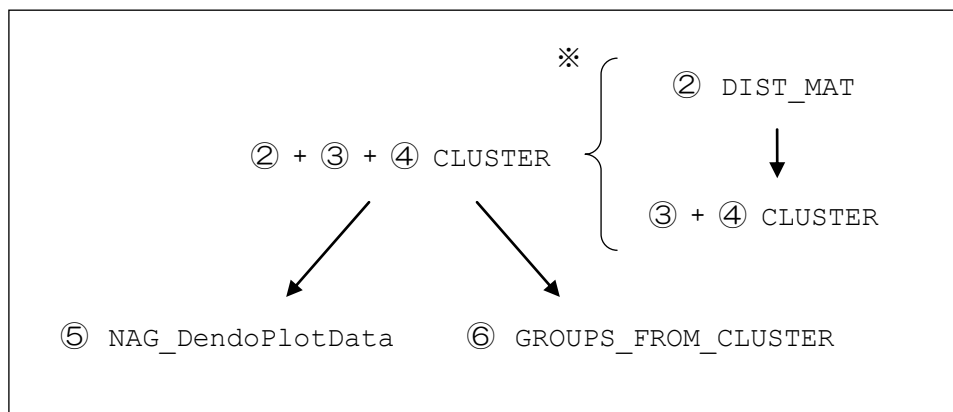


図2 「NAG アドイン」の関数との関係

※ 関数 `DIST_MAT` を用いて、個体間（全ての組合せ）の距離（距離行列）を別途計算することもでき、その結果（距離行列）を入力（引数）として関数 `CLUSTER` でクラスター分析を行うことも可能です。（このドキュメントでは、関数 `DIST_MAT` に関する操作例は割愛します。）

## 4 「NAG アドイン」を用いた分析例

以下の表1のデータ（5つの個体がそれぞれ2変数を持っている）に対して、実際に「NAG アドイン」を用いたクラスター分析の例を示します。

（本章の説明は Excel 2007 に「NAG アドイン」がインストールされていることが前提となります。また、Excel 2003 をご利用の方は適宜読み替えてください。）

距離の定義には、個体間の距離に「ユークリッド距離」、クラスター間の距離（クラスター分析の方法）に「最短距離法」を用いることにします。

表1 5段階評価による成績

生徒	教科	
	国語	数学
A	5	4
B	1	2
C	4	4
D	2	1
E	5	2

この分析例では、以下の3つの関数を使用します。

### 《 CLUSTER 》

表1の様な行列データを入力（引数）としてクラスター分析を行います。

各種距離の定義を指定でき、計算結果としてクラスターの統合過程と樹形図を描く為の情報出力されます。

### 《 NAG\_DendoPlotData 》

関数 CLUSTER の計算結果を入力（引数）として樹形図を作成します。

結果は、Excel のグラフとして出力されます。

### 《 GROUPS\_FROM\_CLUSTER 》

関数 CLUSTER の計算結果を入力（引数）としてクラスター分けを行います。

クラスター数、又は、クラスター距離（高さ）のどちらかを指定して、どの個体が同じクラスターに属しているかを示すインデックス（指標）を出力します。

## 4.1 クラスタ分析 (CLUSTER)

Excel の表に上記の表 1 のデータを入力します。

	A	B	C	D	E
1					
2		生徒	教科		
3			国語	数学	
4		A	5	4	
5		B	1	2	
6		C	4	4	
7		D	2	1	
8		E	5	2	
9					

適当なセル（例えば F2）を選択し、関数 CLUSTER を挿入します。  
（Excel の標準関数と同様に「数式 | 関数の挿入」から行います。）

The screenshot shows an Excel spreadsheet with the data table from the previous block. Cell F2 is selected, and the formula bar shows '='. The '関数の挿入' (Insert Function) dialog box is open, displaying a list of functions. 'CLUSTER' is highlighted in the list. Below the list, there is a description: 'CLUSTER(手法,距離行列,距離タイプ,標準化方法,名前)' and 'NAG:階層的クラス分析を行います。' (NAG: Hierarchical cluster analysis is performed). The dialog box has buttons for '検索開始(G)' (Search), 'この関数のヘルプ' (Help for this function), 'OK', and 'キャンセル' (Cancel).

関数の引数を入力するウィンドウが表示されます。

各引数を以下の様に設定し、「OK」ボタンを押します。

（各引数の詳細は「この関数のヘルプ」をご参照ください。）

「手法」⇒ 最短距離法 ⇒ 1

「距離行列」⇒ 今回は距離行列を使わない ⇒ 設定なし

「X」⇒ 変数データ（成績）⇒ C4:D8

「距離タイプ」⇒ ユークリッド距離 ⇒ E

「標準化方法」⇒ 今回は標準化を行わない ⇒ U

「名前」⇒ 個体の名前（生徒）⇒ B4:B8

	A	B	C	D	E	F	G	H	I	J	K	L	M
1													
2													
3		生徒		教科									
4				国語									
5		A	5	4									
6		B	1	2									
7		C	4	4									
8		D	2	1									
9		E	5	2									
10													
11													
12													
13													
14													
15													
16													
17													
18													
19													
20													
21													

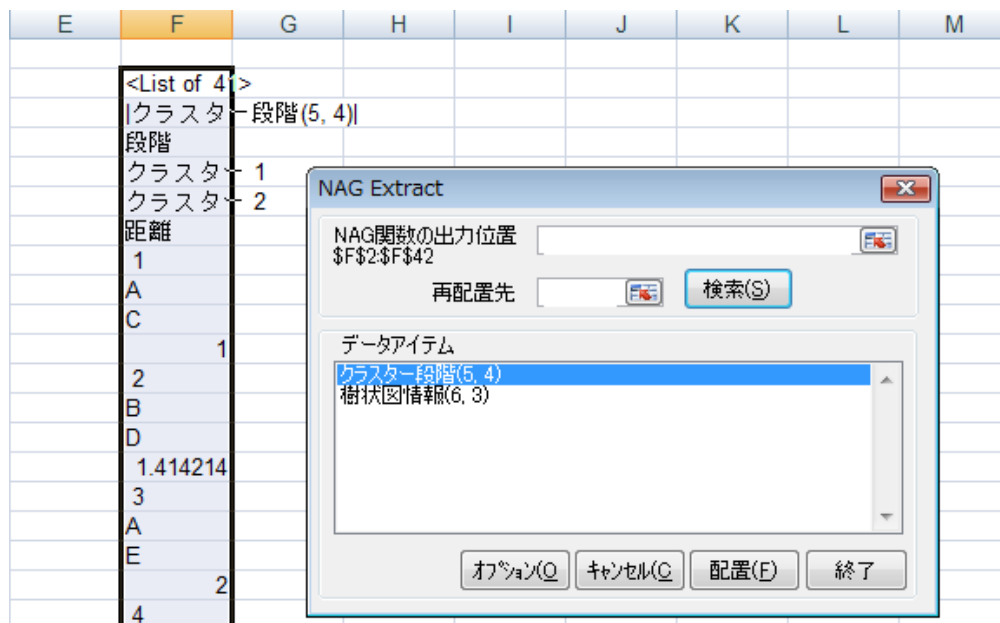
関数を挿入したセル（F2）に <List of 41> などと出力されれば計算は完了です。

引き続き、この選択状態のままで Ctrl + E を押します。

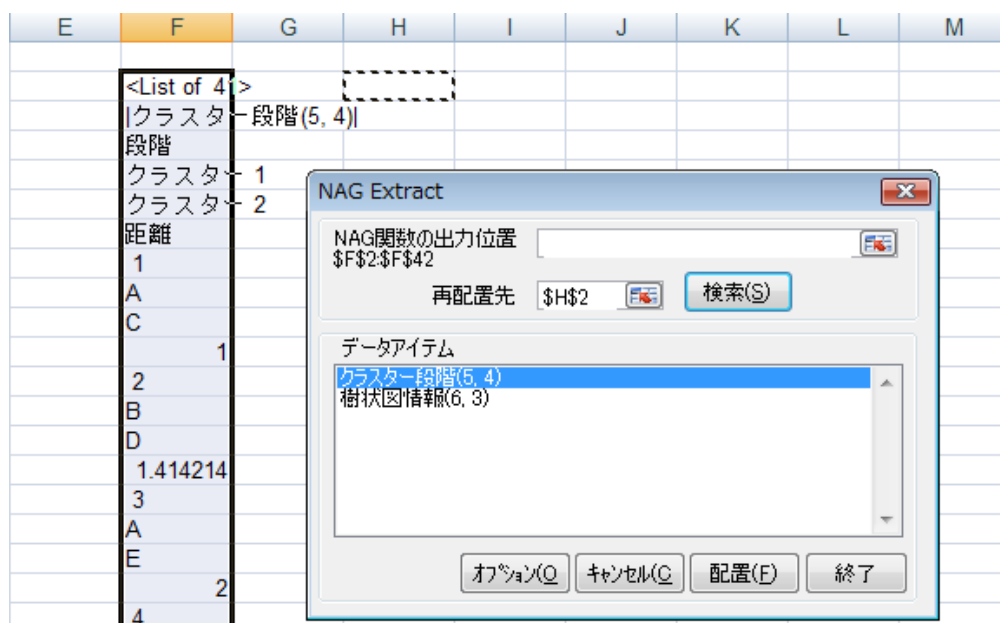
	A	B	C	D	E	F	G
1							
2		生徒	教科			<List of 41>	
3			国語	数学			
4		A	5	4			
5		B	1	2			
6		C	4	4			
7		D	2	1			
8		E	5	2			
9							



関数を挿入したセル（F2）の下方向に出力結果が展開され、NAG Extract（出力結果再配置ツール）ウィンドウが開きます。



「データアイテム」から「クラスター段階」選択し、「再配置先」として適当なセル（例えば H2）をクリックし、「終了」ボタンを押します。



出力結果（この場合は「クラスター段階」）が整理された状態で配置されます。

H	I	J	K
<b>クラスター段階</b>			
段階	クラスター 1	クラスター 2	距離
1	A	C	1
2	B	D	1.414213562
3	A	E	2
4	A	B	3.16227766

他のデータアイテムの配置も同様です。関数を挿入したセル（F2）を選択し、Ctrl + E を押し NAG Extract ウィンドウを開きます。

「データアイテム」で今度は「樹状図情報」選択し、「再配置先」として適当なセル（例えば H9）をクリックし、「終了」ボタンを押します。

出力結果（この場合は「樹状図情報」）が整理された状態で配置されます。

H	I	J	K
<b>クラスター段階</b>			
段階	クラスター 1	クラスター 2	距離
1	A	C	1
2	B	D	1.414213562
3	A	E	2
4	A	B	3.16227766
<b>樹状図情報</b>			
指標	オブジェクト	距離	
1	A	1	
3	C	2	
5	E	3.16227766	
2	B	1.414213562	
4	D	3.16227766	

最後に、出力結果を簡単に説明します。

「クラスター段階」とは、クラスターの統合過程を示します。

段階 1 で個体 A と個体 C が統合されます。

（このクラスターは改めて A と名づけられます。）

段階 2 で個体 B と個体 D が統合されます。

（このクラスターは改めて B と名づけられます。）

段階 3 でクラスター A と個体 E が統合されます。

（このクラスターは改めて A と名づけられます。）

段階 4 でクラスター A とクラスター B が統合されます。

（これで全てのクラスターが統合されたので終了です。）

「樹状図情報」は樹形図（樹状図、デンドログラムとも言う）を描く為の情報です。  
この情報から樹形図を描くルールは以下のようになります。

- ① 「オブジェクト」にある個体名（A, C, E, B, D）を左から右に並べる。
- ② その個体名の上に「距離」で示される長さの垂直な線を引く。
- ③ 左から順に線の最上部から右に向かって水平線を引き次の線にぶつかった所で止める。

この作図に関しましては、Excel のグラフ機能を用いた樹形図を描く為の関数  
NAG\_DendoPlotData を弊社から別途無償にて提供しています。  
詳細は次節 4.2 をご参照ください。

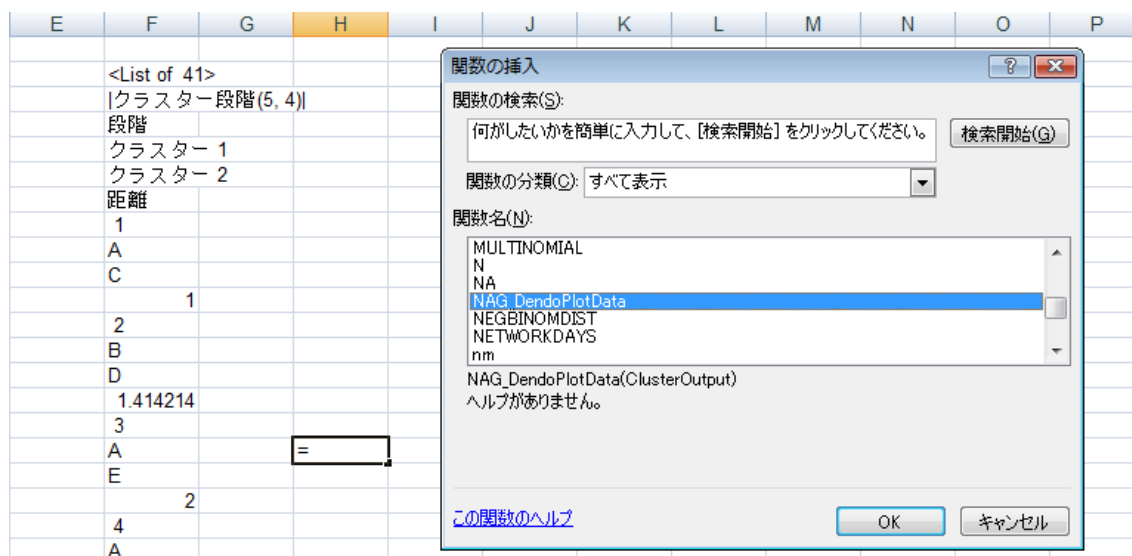
以上で関数 CLUSTER によるクラスター分析は終了です。

## 4.2 樹形図の作成 (NAG\_DendoPlotData)

この節では 4.1 節の操作が完了し「樹状図情報」が得られていることが前提となります。

適当なセル（例えば H17）を選択し、関数 NAG\_DendoPlotData を挿入します。

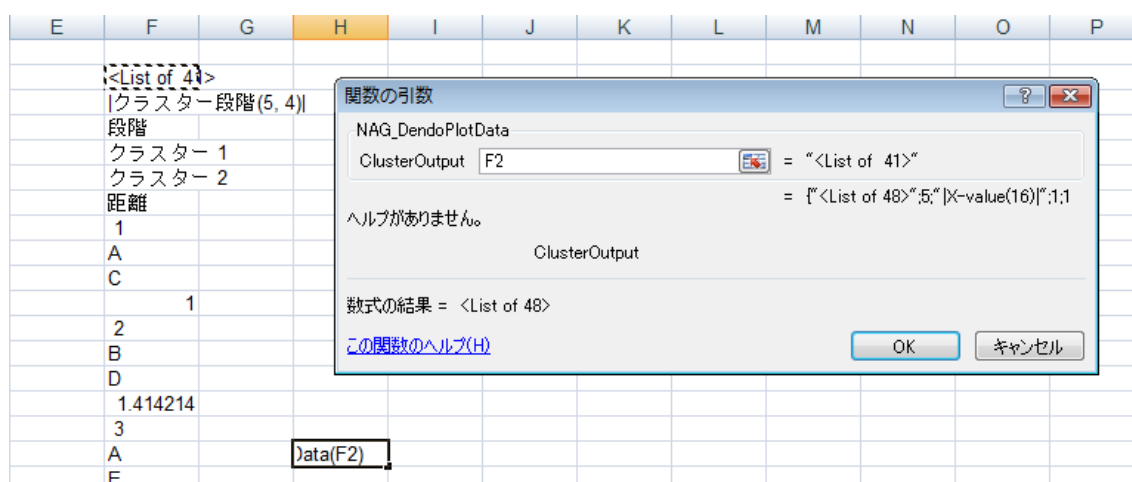
（Excel の標準関数と同様に「数式 | 関数の挿入」から行います。）



関数の引数を入力するウィンドウが表示されます。

引数を以下の様に設定し、「OK」ボタンを押します。

“ClusterOutput” ⇒ 関数 CLUSTER を配置したセルを選択する ⇒ F2

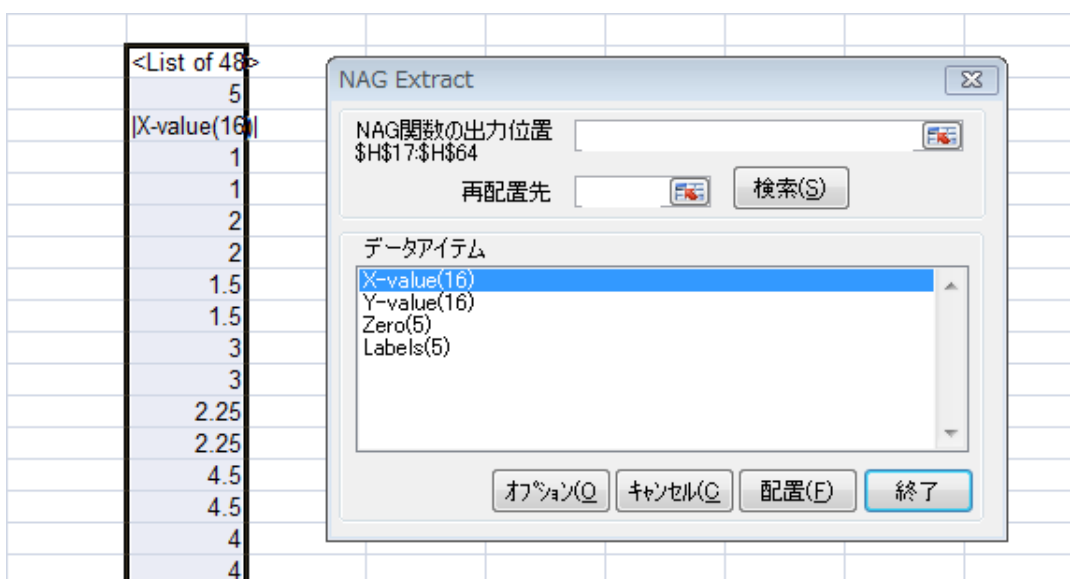


関数を挿入したセル（H17）に <List of 48> などと出力されれば計算は完了です。  
引き続き、この選択状態のままで Ctrl + E を押します。

	<List of 48>	

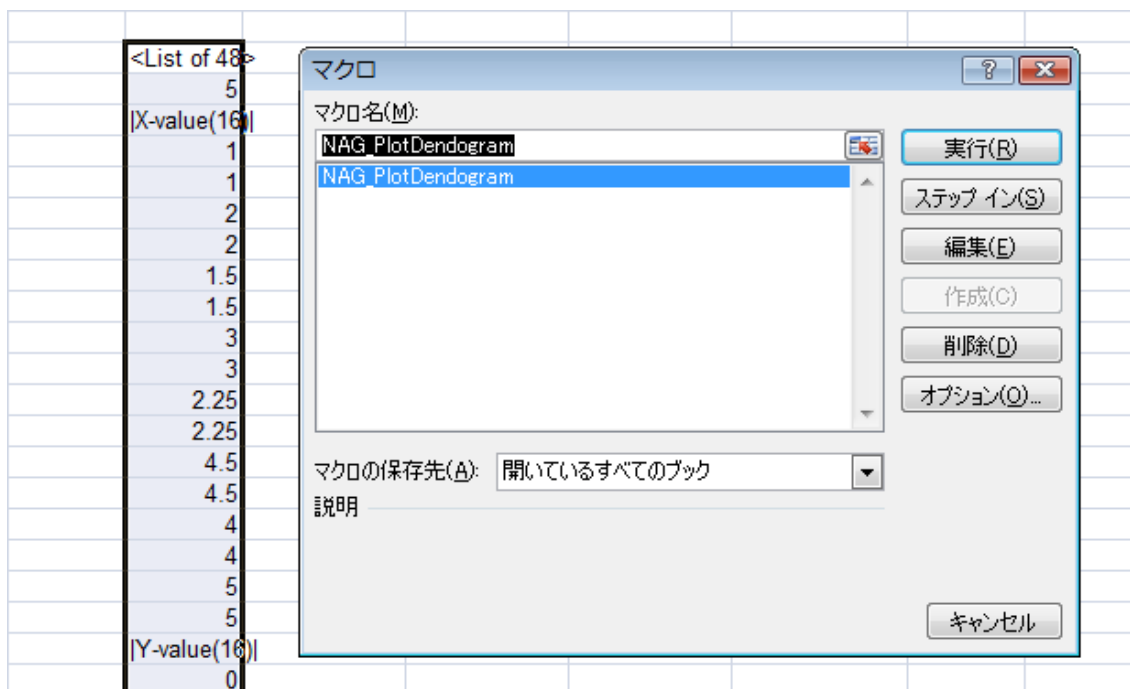
関数を挿入したセル（H17）の下方向に出力結果が展開され、NAG Extract（出力結果再配置ツール）ウィンドウが開きます。

ここでは何もせず「終了」ボタンを押します（出力結果を展開するのが目的です）。

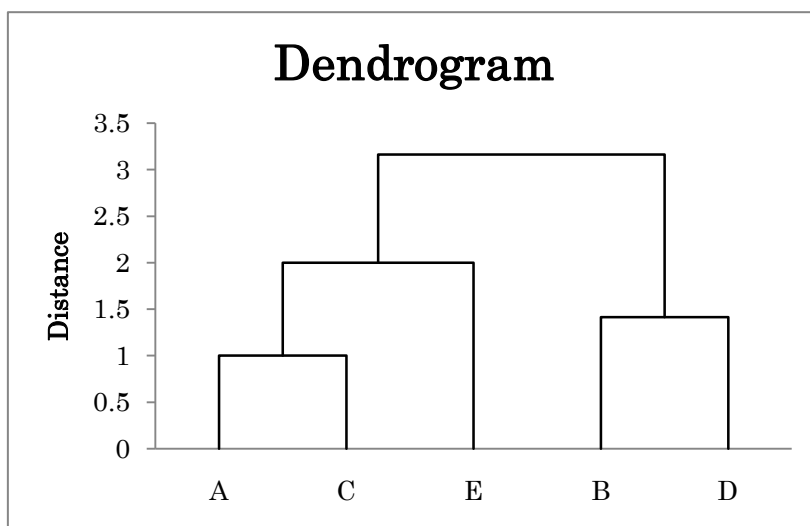


引き続き、関数を挿入したセル（H17）を選択した状態のままで Excel のメニュー「表示 | マクロ」からマクロウィンドウを開きます。

“NAG\_PlotDendogram” を選択し、「実行」ボタンを押します。



Excel のグラフとして樹形図が出力されます。



以上で関数 NAG\_DendoPlotData による樹形図の作成は終了です。

### 4.3 クラスター分け (GROUPS FROM CLUSTER)

この節では 4.1 節の操作が完了し「樹状図情報」が得られていることが前提となります。

適当なセル（例えば M2）を選択し、関数 **GROUPS FROM CLUSTER** を挿入します。

(Excel の標準関数と同様に「数式 | 関数の挿入」から行います。)

関数の引数を入力するウィンドウが表示されます。

各引数を以下の様に設定し、「OK」ボタンを押します。

(各引数の詳細は「この関数のヘルプ」をご参照ください。)

「物体名」 ⇒ 樹状図情報の「オブジェクト」 ⇒ I11:I15

「距離」 ⇒ 樹状図情報の「距離」 ⇒ J11:J15

「クラスター数」 ⇒ 今回は「クラスター距離」を用いる ⇒ 設定なし

「クラスター距離」 ⇒ 樹形図を高さ 1.6 で切断する ⇒ 1.6

「Index」 ⇒ 出力される指標は樹形図での個体の順番（デフォルト） ⇒ 設定なし

	H	I	J	K	L	M	N	O	P	Q	R
		クラスター段階									
段階		クラスター 1	クラスター 2	距離		115, 1.6					
1		A	C		1						
2		B	D		1.414213562						
3		A	E		2						
4		A	B		3.16227766						
		樹状図情報									
指標		オブジェクト	距離								
1		A			1						
3		C			2						
5		E			3.16227766						
2		B			1.414213562						
4		D			3.16227766						

関数の引数

GROUPS\_FROM\_CLUSTER

物体名 I11:I15 = {"A","C","E","B","D"}

距離 J11:J15 = {1;2;3.16227766016838;1.41421356...}

クラスター数 =

クラスター距離 1.6 = 1.6

Index =

NAG: クラスター指標変数の計算を行います。

距離

数式の結果 = <List of 11>

[この関数のヘルプ\(H\)](#) OK キャンセル



関数を挿入したセル（M2）に <List of 11> などと出力されれば計算は完了です。  
引き続き、この選択状態のままで Ctrl + E を押します。

	H	I	J	K	L	M	N
		クラスター段階				<List of 11>	
段階		クラスター 1	クラスター 2	距離			
1		A	C	1			
2		B	D	1.414213562			
3		A	E	2			
4		A	B	3.16227766			
		樹状図情報					
指標		オブジェクト	距離				
1		A	1				
3		C	2				
5		E	3.16227766				
2		B	1.414213562				
4		D	3.16227766				

関数を挿入したセル（M2）の下方方向に出力結果が展開され、NAG Extract（出力結果再配置ツール）ウィンドウが開きます。

The screenshot shows the NAG Extract dialog box open over a spreadsheet. The spreadsheet has columns L through T. Column M contains the following data from row 2 to 7:

- <List of 11>
- |クラスター数|
- 3
- |クラスター距離|
- 1.6
- |グループ指標(5)|
- 1
- 1
- 2
- 3
- 3

The NAG Extract dialog box has the following fields and options:

- NAG関数の出力位置:** \$M\$2:\$M\$12
- 再配置先:** (empty field)
- 検索(S):** (button)
- データアイテム:**
  - クラスター数
  - クラスター距離
  - グループ指標(5)
- Buttons:** オプション(O), キャンセル(C), 配置(E), 終了

「データアイテム」から「グループ指標」選択し、「再配置先」として樹状図情報の右端（例えば K10）をクリックし、「終了」ボタンを押します。

段階	クラスター 1	クラスター 2	距離
1	A	C	1
2	B	D	1.414213562
3	A	E	2
4	A	B	3.16227766

指標	オブジェクト	距離	グループ指標
1	A	1	1
3	C	2	1
5	E	3.16227766	2
2	B	1.414213562	3
4	D	3.16227766	3

出力結果（この場合は「グループ指標」）が整理された状態で配置されます。

H	I	J	K
クラスター段階			
段階	クラスター 1	クラスター 2	距離
1	A	C	1
2	B	D	1.414213562
3	A	E	2
4	A	B	3.16227766
樹状図情報			
指標	オブジェクト	距離	グループ指標
1	A	1	1
3	C	2	1
5	E	3.16227766	2
2	B	1.414213562	3
4	D	3.16227766	3

「グループ指標」は、同じインデックス（指標）の個体（オブジェクト）が同じクラスターに属していることを意味します。

従って、高さ（距離）1.6 で樹形図を切断した場合、{A, C}, {E}, {B, D} の3つのクラスターに分類されることが分かります。

以上で関数 GROUPS\_FROM\_CLUSTER によるクラスター分けは終了です。

4.4 結果の検討

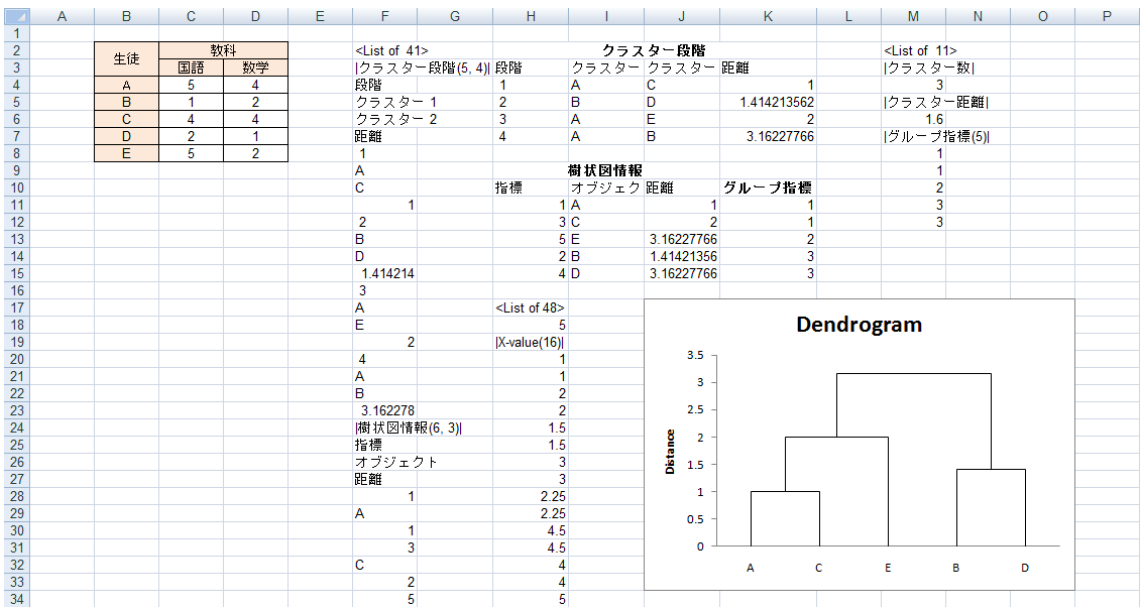


図3 クラスター分析の終了イメージ

最後に、各クラスターの特徴を調べてみます。

表2のデータ（表1の再録）を見ますと、クラスター {A, C} は国語も数学も得意な生徒、クラスター {E} は国語は得意だが数学は苦手な生徒、クラスター {B, D} は国語も数学も苦手な生徒、といった特徴を持っている様です。

表2 5段階評価による成績（表1の再録）

生徒	教科	
	国語	数学
A	5	4
B	1	2
C	4	4
D	2	1
E	5	2

## 5 参考文献

- [1] 田中 豊, 脇本和昌, 「多変量統計解析法」, 現代数学社, 1983.
- [2] 永田 靖, 棟近雅彦, 「多変量解析法入門」, サイエンス社, 2001.
- [3] Everitt B S, "Cluster Analysis", Heinemann, 1974.
- [4] Krzanowski W J, "Principles of Multivariate Analysis",  
Oxford University Press, 1990.