

MultiDendrograms 2.1

Manual

Departament d'Enginyeria



**Informàtica i
Matemàtiques**



UNIVERSITAT
ROVIRA I VIRGILI

Developers: David Torres, Justo Montiel
Advisors: Sergio Gómez, Alberto Fernández
Universitat Rovira i Virgili, Tarragona (Spain)
<http://deim.urv.cat/~sgomez/multidendrograms.php>



Contents

Introduction	4
1. Input data	5
Matrix-like file format	5
List-like file	6
2. Loading data.....	7
3. Actions.....	9
4. Settings.....	10
Main data representation settings	11
Tree settings.....	12
Nodes settings.....	14
Axis settings	15
5. Analyzing and exporting results	18
APPENDIX A. Requirements, installation and execution	21
Requirements	21
Installation	21
Basic execution.....	21
Advanced execution	21
APPENDIX B. Preparing input data with Microsoft Excel.....	22
APPENDIX C. Languages	25
APPENDIX D. History of changes.....	26
APPENDIX E. License.....	27

MultiDendrograms - Manual

Introduction

MultiDendrograms is a simple yet powerful program to make the Hierarchical Clustering of real data, distributed under an Open Source license. Starting from a distances (or weights) matrix, *MultiDendrograms* calculates its dendrogram using the most common Agglomerative Hierarchical Clustering algorithms (e.g. Single Linkage, Complete Linkage and Unweighted Average), allows the tuning of many of the graphical representation parameters, and the results may be easily exported to file.

MultiDendrograms implements the variable-group algorithms in [1] to solve the non-uniqueness problem found in the standard pair-group algorithms and implementations. This problem arises when two or more minimum distances between different clusters are equal during the amalgamation process. The standard approach consists in choosing a pair, breaking the ties between distances, and proceeds in the same way until the final hierarchical classification is obtained. However, different clusterings are possible depending on the criterion used to break the ties (usually a pair is just chosen at random!).

The variable-group algorithms group more than two clusters at the same time when ties occur, given rise to a graphical representation called *multidendrogram*. Their main properties are:

- When there are no ties, the variable-group algorithms give the same results as the standard pair-group ones.
- They always give a uniquely determined solution.
- In the multidendrogram representation for the results one can explicitly observe the occurrence of ties during the agglomerative process. Furthermore, the height of any fusion interval (the *bands* in the program) indicates the degree of heterogeneity inside the corresponding cluster.

In this manual we show how to prepare the data file, load it into *MultiDendrograms*, the meaning of the parameters of the program, and how to export the results to file.

The main characteristics of *MultiDendrograms* are:

- Multiplatform, runs in Windows, Linux and MacOS.
- Graphical user interface.
- Implementation of variable-group algorithms for Agglomerative Hierarchical Clustering.
- Works with distance and weight matrices.
- Many parameters for the customization of the dendrogram layout.
- Navigation through the dendrogram information in a folder-like window.
- Calculation of corresponding ultrametric matrix.
- Calculation of deviation measures such as the cophenetic correlation coefficient.
- Save dendrogram details in text and Newick tree format.
- Save dendrogram image as JPG, PNG and EPS.

MultiDendrograms web page: <http://deim.urv.cat/~sgomez/multidendrograms.php>

Please cite [1] if you use *MultiDendrograms* in your publications.

[1] Alberto Fernández and Sergio Gómez, *Solving Non-Uniqueness in Agglomerative Hierarchical Clustering Using Multidendrograms*, *Journal of Classification* **25** (2008) 43-65

MultiDendrograms - Manual

1. Input data

The data file must represent a distances (or weights) matrix, like the one in the following table:

	a	b	c	d
a	0	2	4	7
b	2	0	2	5
c	4	2	0	3
d	7	5	3	0

There are two different arrangements these data can be stored in a text file such that *MultiDendrograms* may accept them: matrix and list formats.

Matrix-like file format

Each line in the text file contains a data matrix row. The characteristics of these files are:

- The matrix must be symmetric, and the diagonal elements must be zeros.
- Within each row, the elements are separated by: spaces (' '), tab character, semicolon (';'), comma (',') or vertical bar ('|').
- It is possible to include the names in an additional first row or column, but not in both.
- If present, the labels of the nodes can not contain any of the previous separators.

Some different representations for the previous matrix could be:

Node_a	Node_b	Node_c	Node_d
0.0	2.0	4.0	7.0
2.0	0.0	2.0	5.0
4.0	2.0	0.0	3.0
7.0	5.0	3.0	0.0

Matrix-like with node labels in first row, data separated by tabs

a	0.0	2.0	4.0	7.0
b	2.0	0.0	2.0	5.0
c	4.0	2.0	0.0	3.0
d	7.0	5.0	3.0	0.0

Matrix-like with node labels in the first column, data separated by spaces

0.0,2.0,4.0,7.0
2.0 0.0 2.0 5.0
4.0 2.0 0.0 3.0
7.0 5.0;3.0 0.0

Matrix-like file without node labels, data separated by all kind of separators

MultiDendrograms - Manual

List-like file

Each line in the text file contains three elements, which represent the labels of two nodes and the distance (or weight) between them. The characteristics of these files are:

- The separators between the three elements may be: spaces (' '), tab character, semicolon (';'), comma (',') or vertical bar ('|').
- The labels of the nodes can not contain any of the previous separators.
- Distances from an element to itself (e.g. "a a 0.0") must not be included.
- *MultiDendrograms* accepts either the presence or absence of the symmetric data elements, i.e. if the distance between nodes a and b is 2.0, it is possible to include in the list the line "a b 2.0", or both "a b 2.0" and "b a 2.0". If both are present, the program checks if they are equal.

For example, three list-like files for the previous matrix could be:

```
a b 2
a c 4
a d 7
b c 2
b d 5
c d 3
```

Simple list

```
a b 2
a c 4
a d 7
b a 2
b c 2
b d 5
c a 4
c b 2
c d 3
d a 7
d b 5
d c 3
```

Complete lists

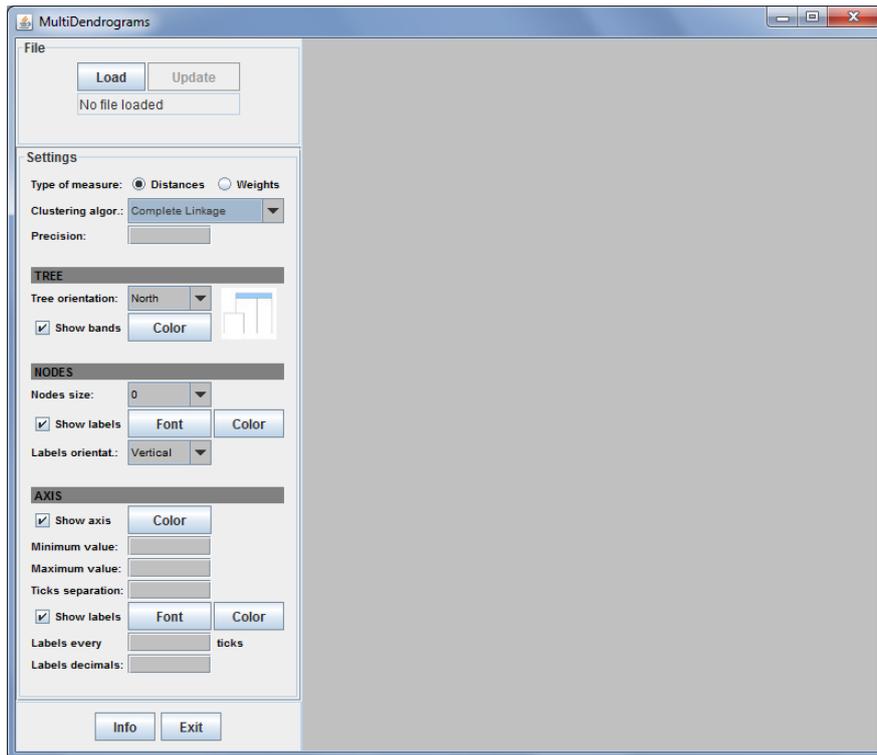
```
a b 2
a c 4
a d 7
b c 2
b d 5
c d 3
b a 2
c a 4
d a 7
c b 2
d b 5
d c 3
```

MultiDendrograms - Manual

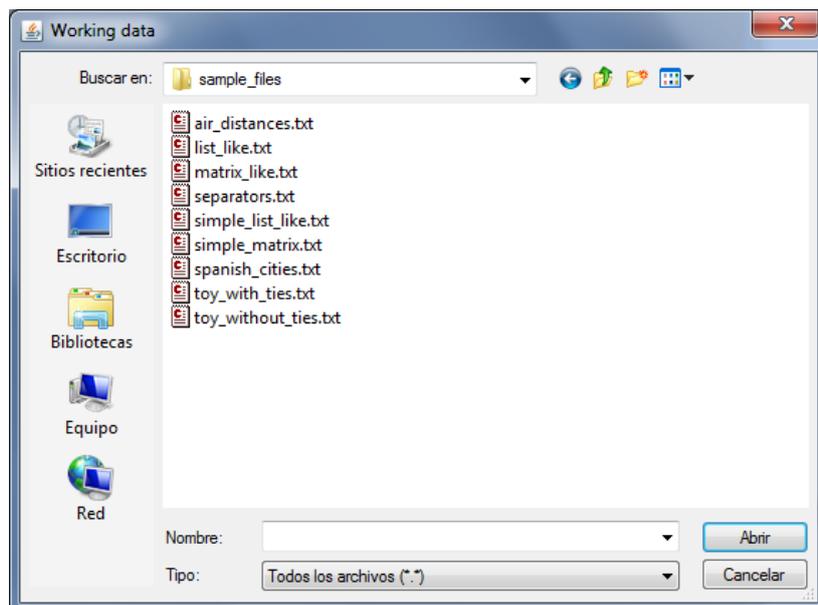
2. Loading data

Once we have our data in a compatible format, we can load them into *MultiDendrograms*.

1. Choose the desired settings, mainly the Type of measure and the Clustering algorithm. These settings will be explained in detail in the next sections.
2. Click on the 'Load button':

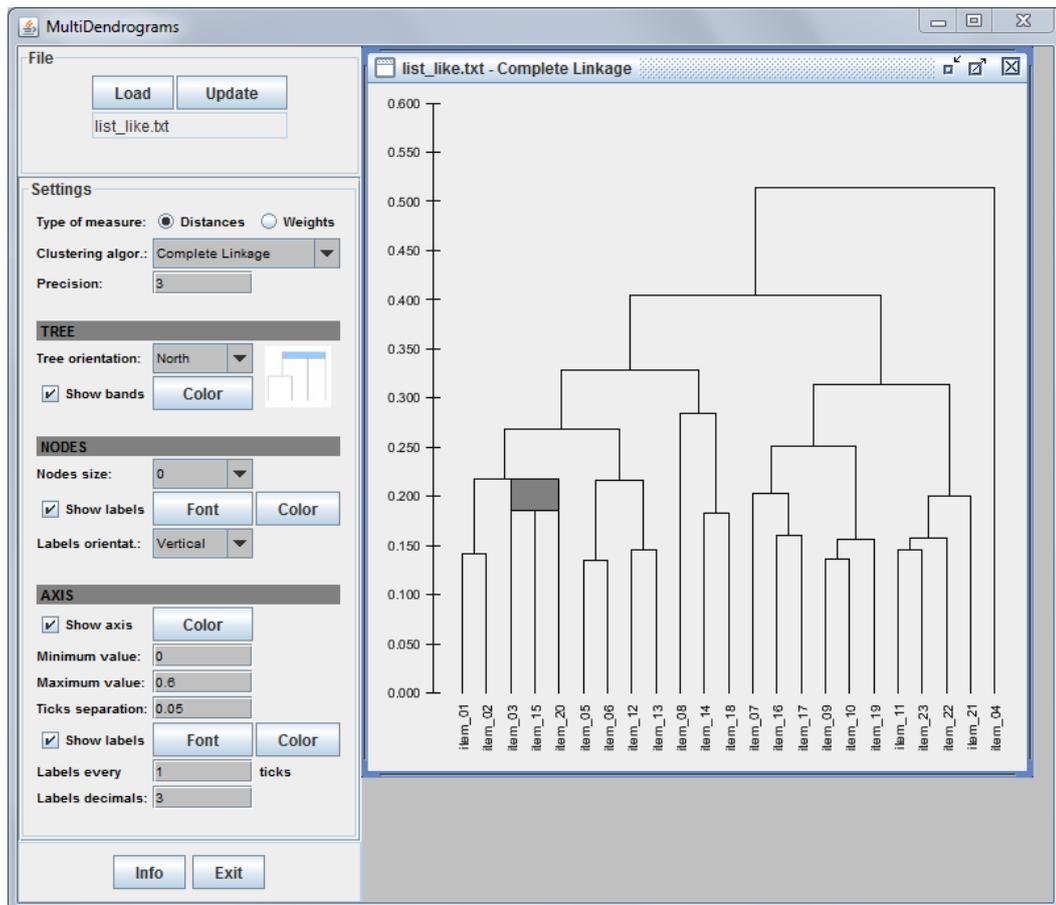


3. Select the file to open and then click on the 'Open' button:



MultiDendrograms - Manual

4. Now the data is loaded and its dendrogram representation is shown:



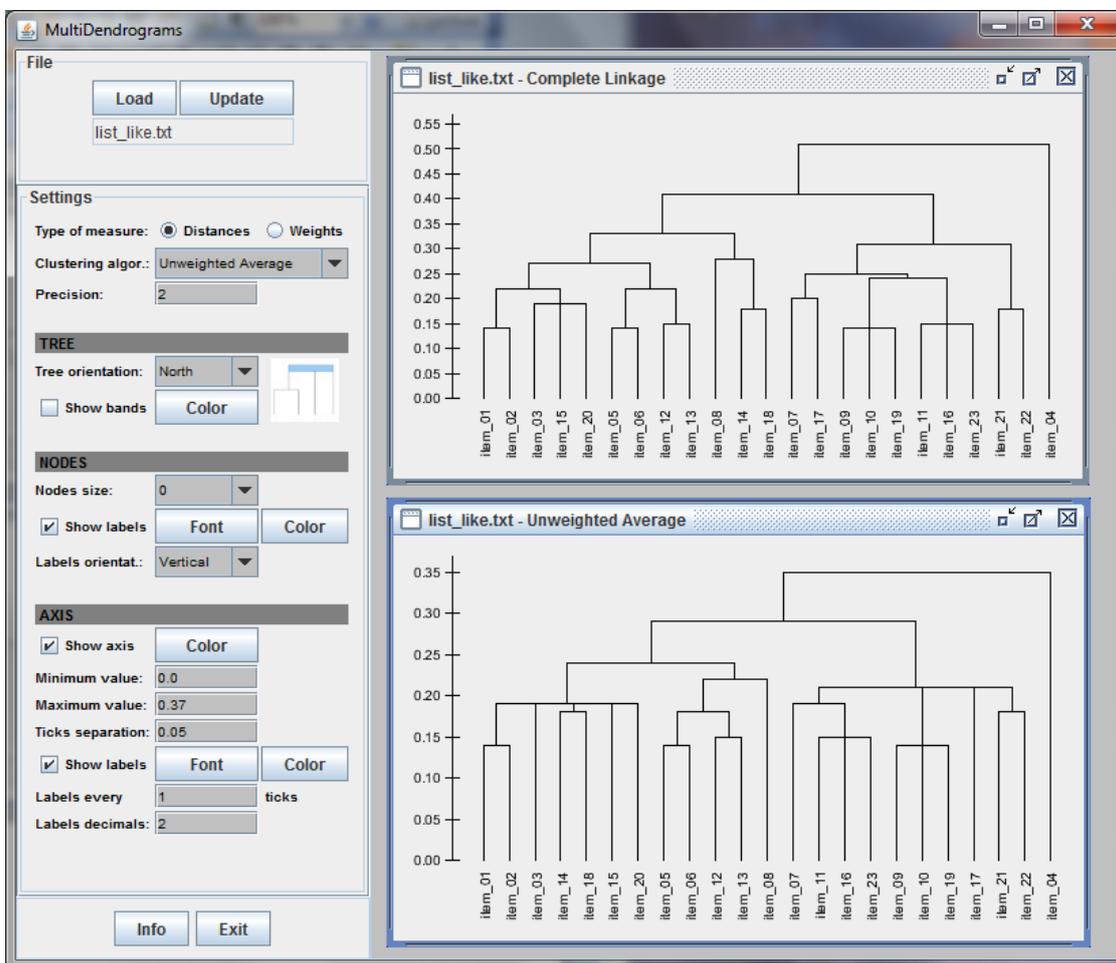
MultiDendrograms - Manual

3. Actions

MultiDendrograms only has two action buttons, **Load** and **Update**. **Load** is used to read data from a file and create a new window for the dendrogram, using the current values of the parameters, while **Update** is needed for the actualization of the active dendrogram when one or more parameters are changed. Below these buttons it is shown the name of the data file of the active dendrogram.



It is possible to load the same data file several times, in order to compare the dendrogram appearance for different parameters settings.



The parameters shown always correspond to the active (selected) dendrogram window.

Finally, there are two additional buttons, **Info** to show the information of the program, and **Exit** to quit the application.



MultiDendrograms - Manual

4. Settings

The program automatically applies default values to the parameters depending on the data loaded, which should be adjusted as desired. The following figure shows the settings tab, with four different areas corresponding to the main data representation, tree, nodes and axis settings respectively:

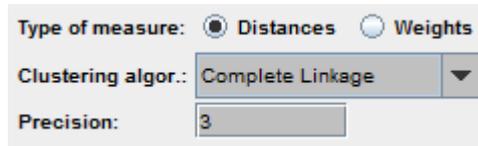
The screenshot shows the 'Settings' dialog box with the following configuration:

- Type of measure:** Distances Weights
- Clustering algor.:** Unweighted Average (dropdown)
- Precision:** 2 (input field)
- TREE:**
 - Tree orientation:** North (dropdown)
 - Show bands
 - Color (button)
- NODES:**
 - Nodes size:** 0 (dropdown)
 - Show labels
 - Font (button)
 - Color (button)
 - Labels orientat.:** Vertical (dropdown)
- AXIS:**
 - Show axis
 - Color (button)
 - Minimum value:** 0.0 (input field)
 - Maximum value:** 0.37 (input field)
 - Ticks separation:** 0.05 (input field)
 - Show labels
 - Font (button)
 - Color (button)
 - Labels every:** 1 (input field) ticks
 - Labels decimals:** 2 (input field)

Changes in the main data representation parameters affect the structure of the dendrogram tree, thus it needs to be fully recalculated, operation which may take several seconds, even minutes (depending of the data size and the computer speed). On the other hand, changes in the tree, nodes and axis settings only modify the visual representation of the dendrogram, which are much faster to update.

MultiDendrograms - Manual

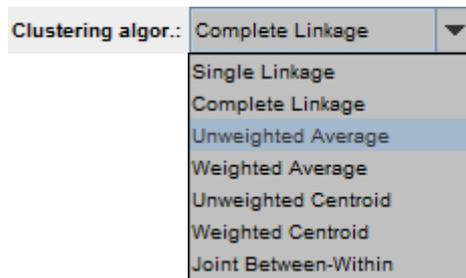
Main data representation settings



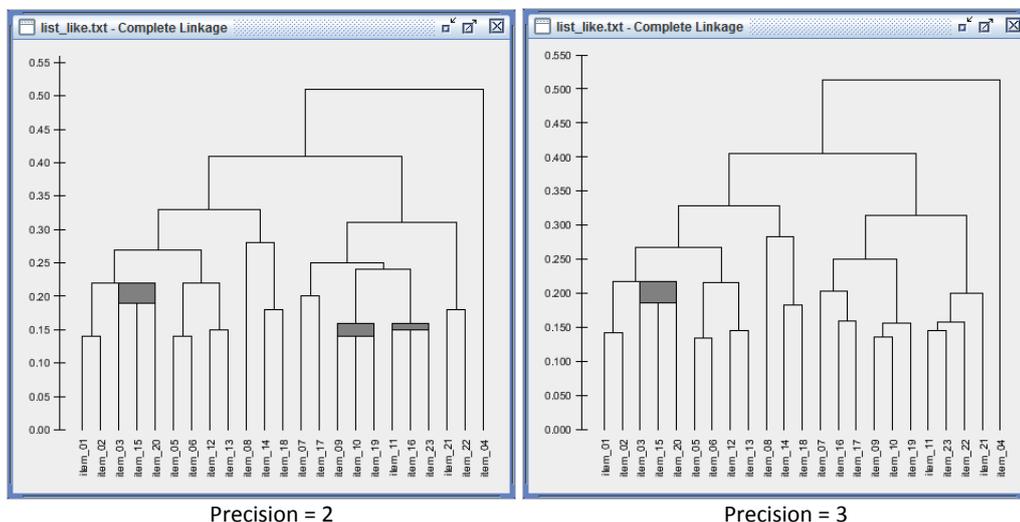
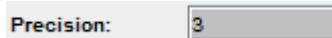
- **Type of measure:** It allows choosing between two kinds of measures, *distances* and *weights*. Choose between them according to the meaning of the loaded data. With *distances*, the closer the elements the lower their distance. On the contrary, with *weights*, the closer the elements the larger their weight. By default, *distances* is selected.



- **Clustering algorithm:** Seven clustering algorithms are available, *single linkage*, *complete linkage*, *unweighted average*, *weighted average*, *unweighted centroid*, *weighted centroid* and *joint between-within*. By default, *unweighted average* is selected.

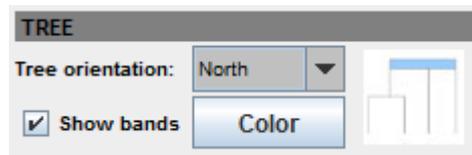


- **Precision:** Number of significant digits of the data and for the calculations. This is a very important parameter, since equal distances at a certain precision may become different by increasing its value. Thus, it may be responsible of the existence of tied distances. The rule should be not to use a precision larger than the resolution given by the experimental setup which has generated the data. By default, the precision is set to that of the data value with the largest number of significant digits.

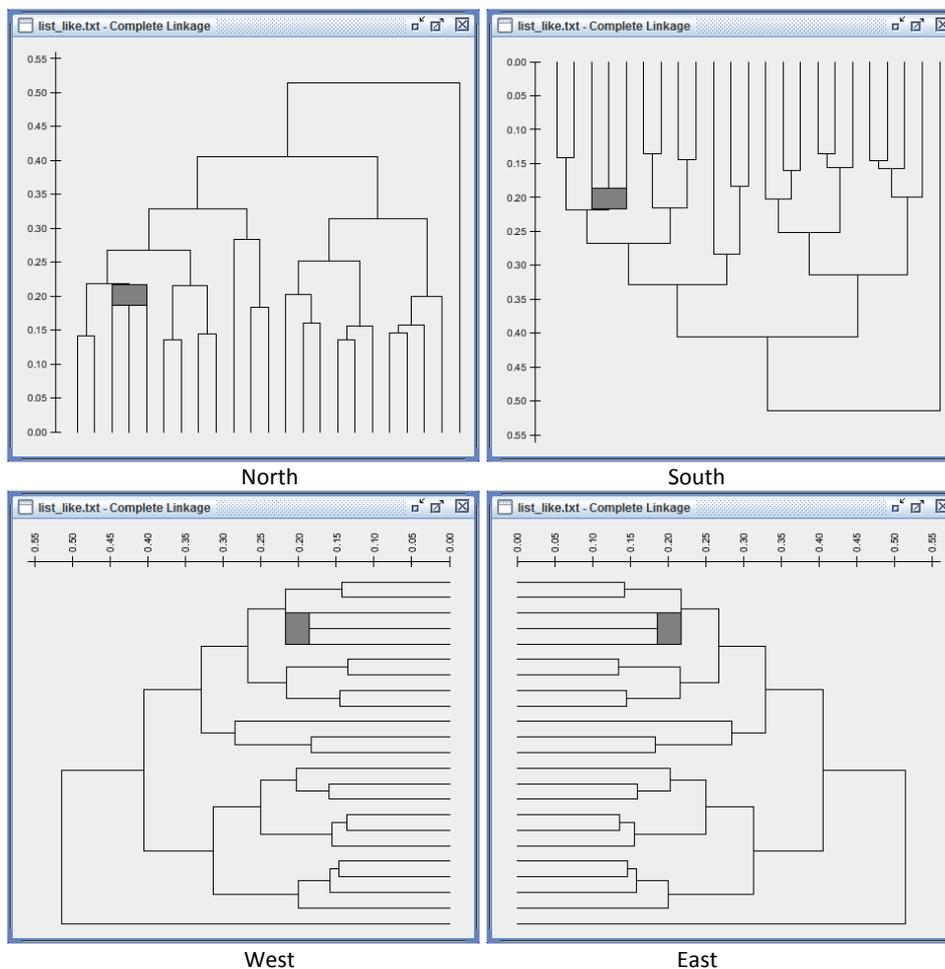


MultiDendrograms - Manual

Tree settings

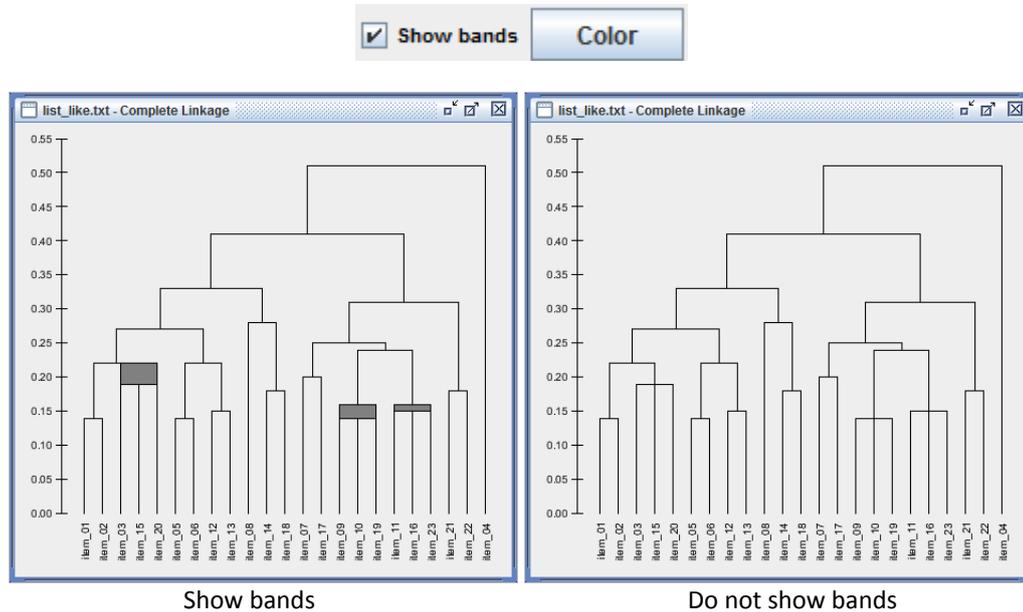


- **Tree orientation:** Four orientations are available, *north*, *south*, *east* and *west*, which refer to the relative position of the root of the tree. By default, *north* is selected.



MultiDendrograms - Manual

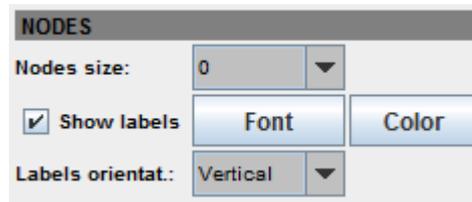
- **Show bands:** It allows showing a band or not in case of tied minimum distances between three or more elements, and selecting the color of the band. If selected, the bands show the heterogeneity of all the distances between the clustered elements. Otherwise, the elements are grouped at their minimum distance. By default, **show bands** is selected, and its default color is **light gray**.



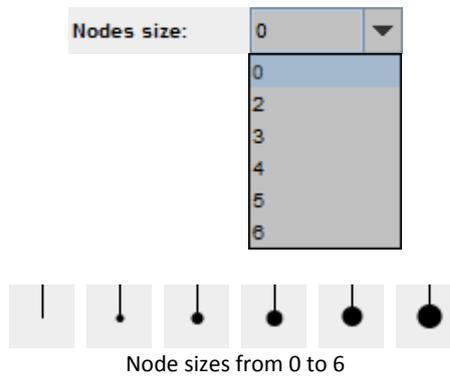
Let us explain the meaning of the bands. In *MultiDendrograms*, if several pairs of elements share the same minimal distance, they are clustered together in one step. For instance, suppose that the minimal distance is 0.4, and that they correspond to the tied pairs (A,B) and (B,C). *MultiDendrograms* puts them together in the same cluster (A,B,C) at height 0.4. However, if the distance (A,C) is 0.5, it is possible to represent the cluster (A,B,C) as a rectangle which spans between heights 0.4 and 0.5, thus showing the heterogeneity of the clustered elements.

MultiDendrograms - Manual

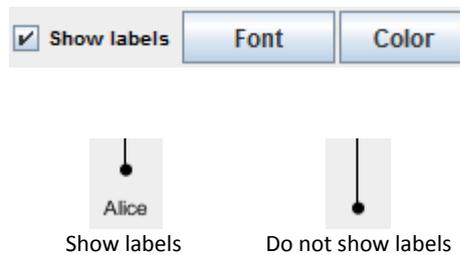
Nodes settings



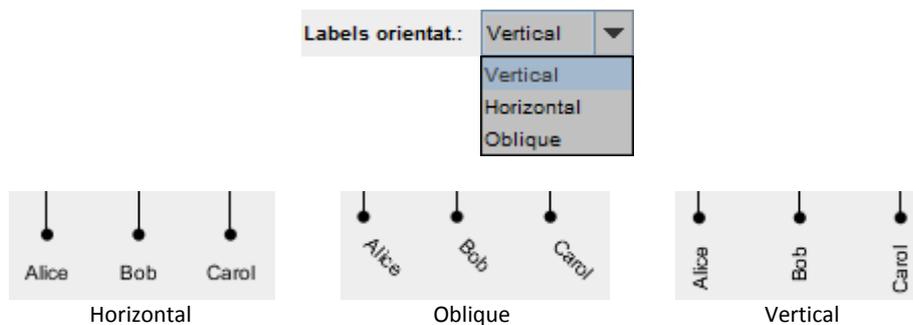
- **Nodes size:** Six different node sizes are available. By default, 0 is selected (i.e. nodes not shown):



- **Show labels:** It allows showing or not the labels of the nodes, and selecting their color and font. By default, **show labels** is selected, the font is **Arial** and the color is **black**:



- **Labels orientation:** Three orientations are available: **vertical**, **horizontal** and **oblique**. By default, **vertical** is selected:

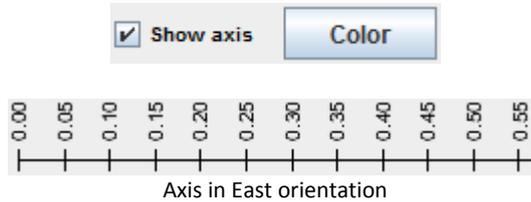


MultiDendrograms - Manual

Axis settings

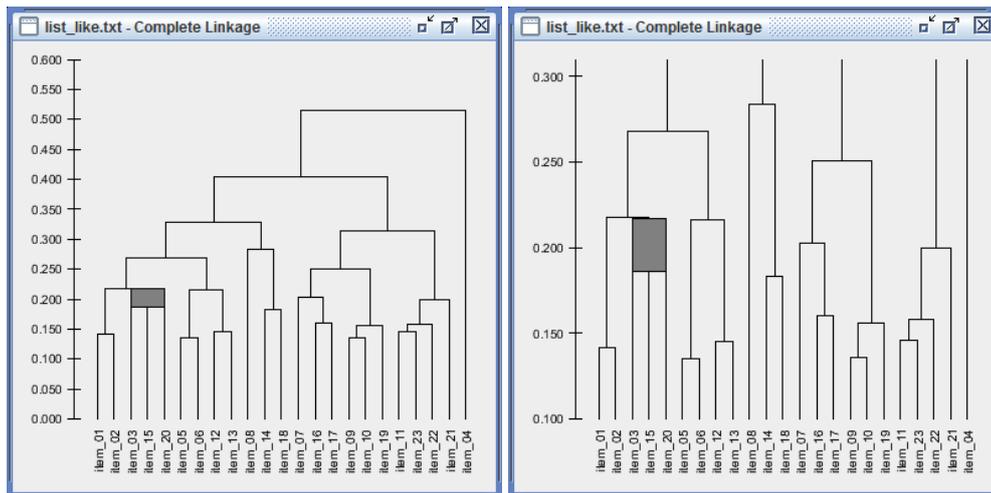
AXIS		
<input checked="" type="checkbox"/> Show axis	Color	
Minimum value:	0	
Maximum value:	0.570	
Ticks separation:	0.050	
<input checked="" type="checkbox"/> Show labels	Font	Color
Labels every	1	ticks
Labels decimals:	3	

- **Show axis:** It allows showing or not the axis, and selecting its color. By default, **show axis** is selected and the selected color is **black**.



- **Minimum value / Maximum value:** They allow choosing the minimum and maximum value of the axis, respectively. They also affect the view of the dendrogram. The default values are calculated from the data.

Minimum value:	0.10
Maximum value:	0.31

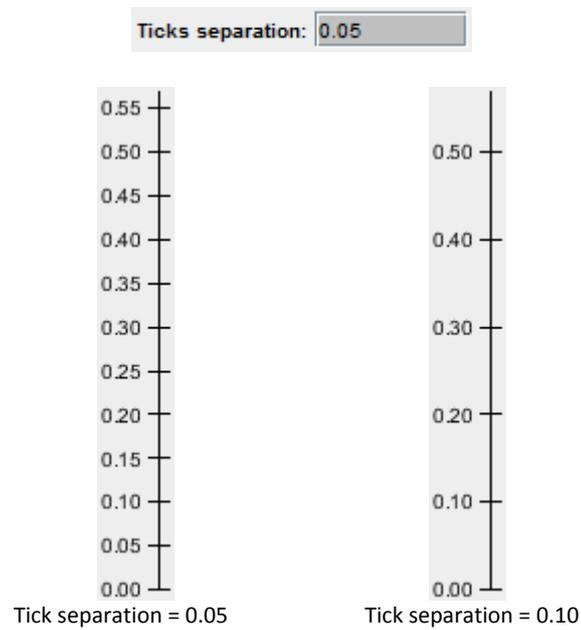


Minimum = 0, Maximum = 0.6 (default)

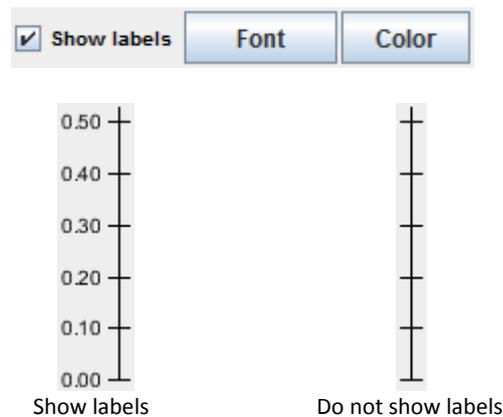
Minimum = 0.10, Maximum = 0.31

MultiDendrograms - Manual

- **Tick separation:** It allows choosing the separation between consecutive ticks of the axis. The default value is calculated from the data:

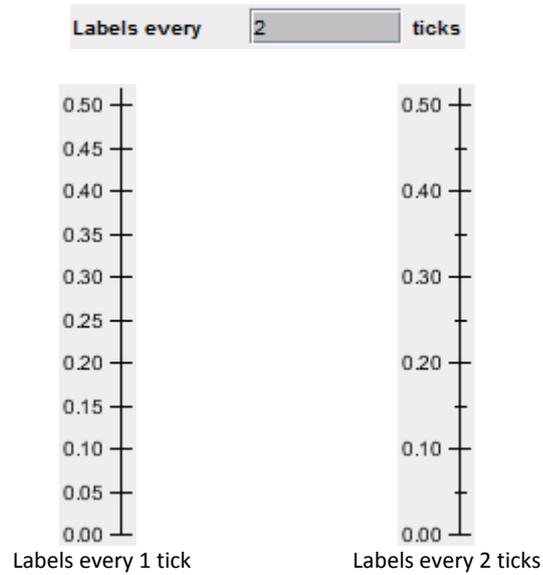


- **Show labels:** It allows showing or not the labels of the axis, and selecting their color and font. By default, *show labels* is selected, the font is *Arial* and the color is *black*.

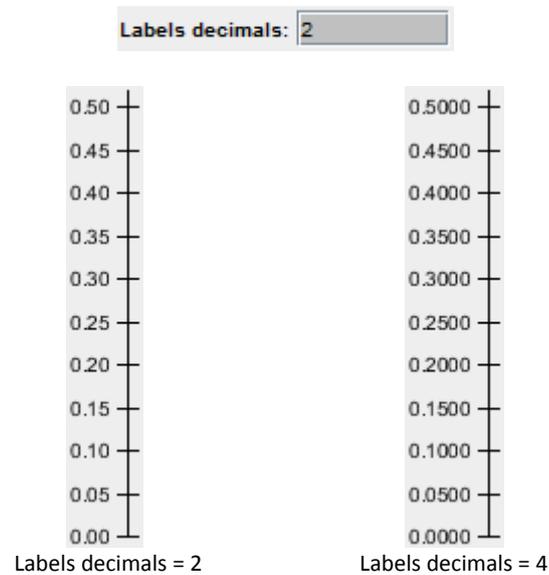


MultiDendrograms - Manual

- **Labels every ... ticks:** Number of consecutive ticks to find the next labeled tick. By default is set to 1.

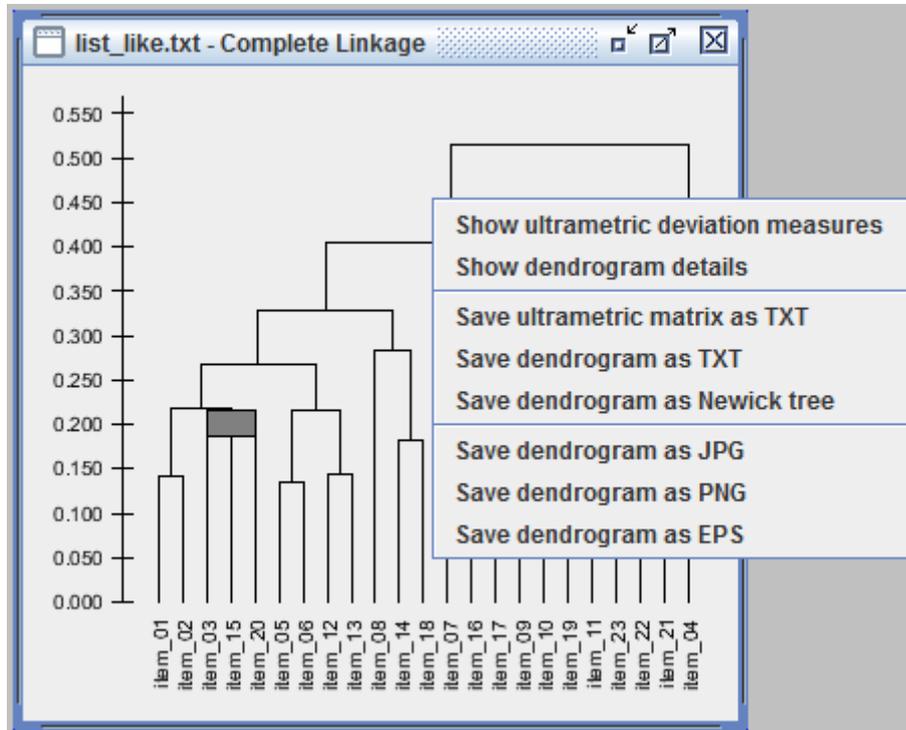


- **Labels decimals:** Number of decimal digits of the tick labels. By default it is set equal to the **precision** parameter.

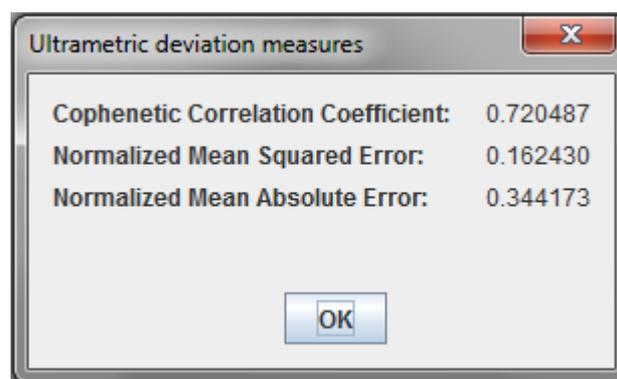


5. Analyzing and exporting results

The contextual menu, available by right-clicking the dendrogram windows, gives access to several options for analyzing and exporting the results to file.



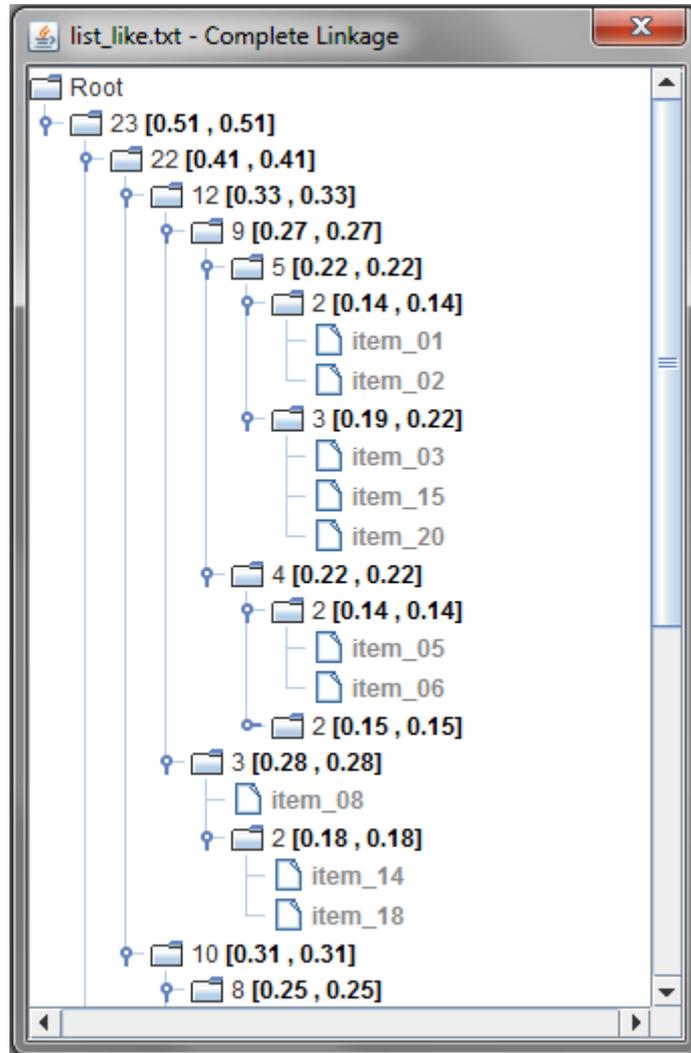
- **Show ultrametric deviation measures:** Calculates the ultrametric matrix corresponding to the active dendrogram and obtains three different deviation measures between the original and the ultrametric matrices: the **cophenetic correlation coefficient**, the **normalized mean squared error** and the **normalized mean absolute error**.



Deviation measures

MultiDendrograms - Manual

- **Show dendrogram details:** Opens a window which contains all the information of the dendrogram in a navigable folder-like structure:



The available information in the details window is:

- Number of data items (leaves of the tree) under each interior node of the dendrogram. The interior nodes in the dendrogram representation correspond to the clusters found during the agglomeration process.
- Maximum and minimum distances at which the sons of an interior node are joined to form a new cluster. These values may only be different in case of tied distances, which become a band in the multidendrogram representation.
- List of sons of each interior node, which may be either interior nodes or data items.

MultiDendrograms - Manual

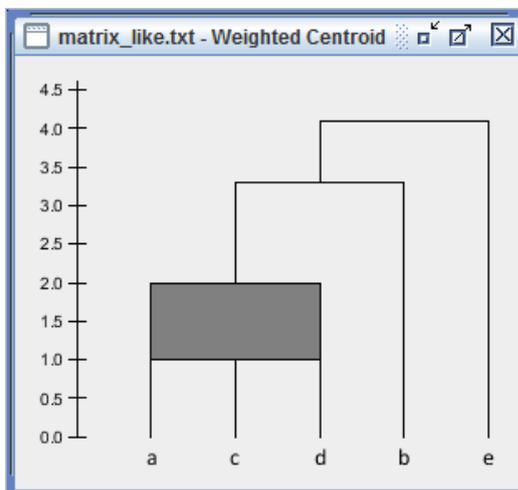
- **Save ultrametric matrix as TXT:** Calculates the ultrametric matrix corresponding to the loaded data and saves it to a text file represented as a matrix, not as a list, with the nodes labels in the first row. This text file can then be easily loaded into any text editor or spreadsheet application (e.g. Microsoft Excel).

a	b	c	d	e
0.0	1.2	1.0	2.0	5.0
1.2	0.0	3.0	7.0	4.0
1.0	3.0	0.0	1.0	8.0
2.0	7.0	1.0	0.0	6.0
5.0	4.0	8.0	6.0	0.0

Original distances matrix

a	b	c	d	e
0.0	3.3	1.0	1.0	4.1
3.3	0.0	3.3	3.3	4.1
1.0	3.3	0.0	1.0	4.1
1.0	3.3	1.0	0.0	4.1
4.1	4.1	4.1	4.1	0.0

Ultrametric matrix



- **Save dendrogram as TXT:** Saves the dendrogram details to a text file.

```
+ 5 [4.1, 4.1]
  + 4 [3.3, 3.3]
    + 3 [1.0, 2.0]
      * a
      * c
      * d
    * b
  * e
```

Dendrogram details file

- **Save dendrogram as Newick tree:** Saves the dendrogram details in Newick tree format (see http://en.wikipedia.org/wiki/Newick_format).

```
(( (a:1.0,c:1.0,d:1.0):2.3,b:3.3):0.8,e:4.1);
```

Dendrogram in Newick format

In this format the information given by the bands is lost, only the minimum distance is used. However, it has the advantage that is a standard format used in many other applications, thus allowing their use to generate other graphical representations.

- **Save dendrogram as JPG, PNG, EPS:** It is also possible to save the image of the dendrogram in three different formats (JPG, PNG and EPS) using their corresponding **Save dendrogram as ...** context menu items.

MultiDendrograms - Manual

APPENDIX A. Requirements, installation and execution

Requirements

To run *MultiDendrograms* it is necessary to have installed a recent version of the Java Runtime Environment (JRE):

- Java: <http://java.com>

You can check if Java is already in your computer following these steps:

1. Open a shell or command prompt (In Windows: Start -> Run -> type "cmd" -> Enter):
2. Type: java -version

If JRE is installed, you will get its version.

Installation

MultiDendrograms does not require installation. Just unpack the main ZIP file into a folder using any unzip program, e.g. 7-zip, iZarc, WinRAR or WinZip.

Basic execution

- Windows: double-click **multidendrograms.bat** or **multidendrograms.jar**
- Linux: run **multidendrograms.sh** or **multidendrograms.jar**
- MacOS: : double-click **multidendrograms.jar**

Advanced execution

In the command-line:

```
java -jar multidendrograms.jar [-h | -help] [-level LEVEL] [-XML | -TXT]
```

The program accepts these parameters:

-h, -help

Prints a list of options

-level LEVEL

Sets the detail level of the log file to one of the following LEVEL:

OFF (logging deactivated)

SEVERE

WARNING

INFO

CONFIG

FINE

FINER

FINEST

ALL (print all messages)

-XML

Prints the messages in the log file in XML format

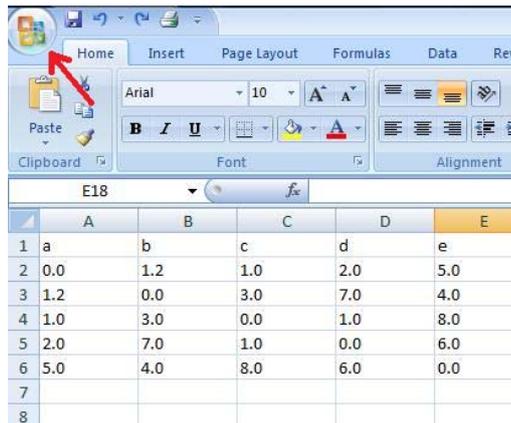
-TXT

Prints the messages in the log file in TXT format

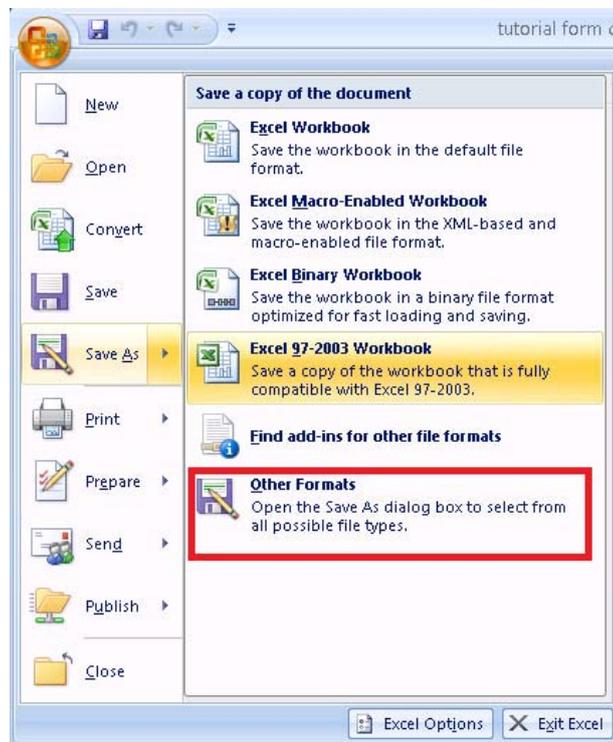
APPENDIX B. Preparing input data with Microsoft Excel

MultiDendrograms cannot load data directly from a Microsoft Excel (or similar) file, we first need to save our data in a compatible format. We will assume Microsoft Excel 2007, but similar procedures apply to other versions and similar programs.

1. Click on the button with the Microsoft Office logo:

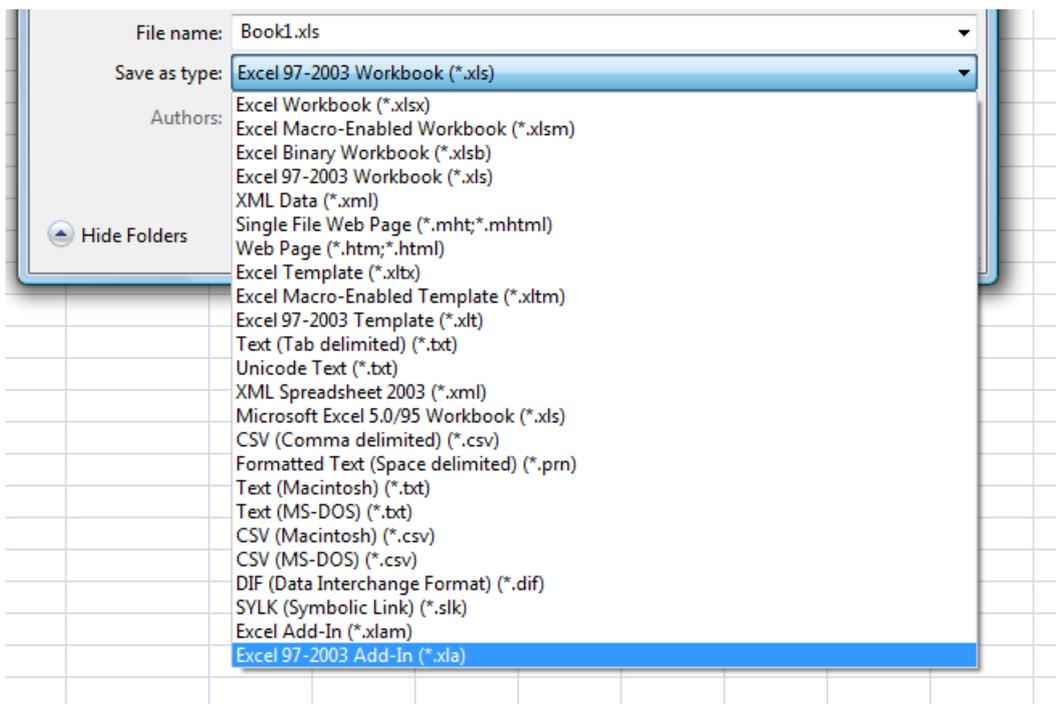


2. Select the option 'Save As' and then 'Other Formats':



MultiDendrograms - Manual

3. Now choose the file format to save the file:

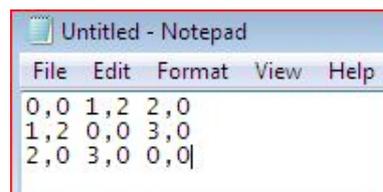


The compatible formats are the following:

- Text (Tab delimited)
- Unicode Text
- CSV (Comma delimited)
- Formatted text (Space delimited)
- Text (Macintosh)
- Text (MS-DOS)
- CSV (Macintosh)
- CSV (MS-DOS)

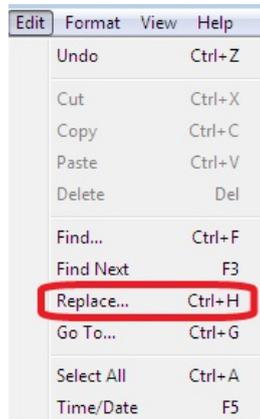
4. *MultiDendrograms* needs that decimal numbers use the character ‘.’ as the decimal symbol, e.g. “3.1416”. Unfortunately, some regional system configurations use different decimal symbols, e.g. in Spanish it is ‘,’ as in “3,1416”. In these cases, the previously exported file has to be edited to change the decimal symbol to ‘.’:

- Open the file with Notepad (or any other file editor):

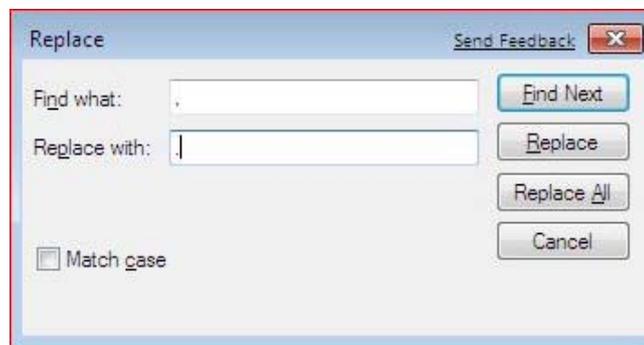


MultiDendrograms - Manual

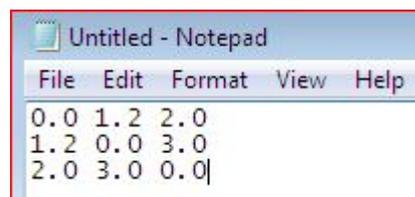
- Go to the 'Edit' menu and select 'Replace':



- Replace all appearances of ',' with '.':



- Save the modified file:



APPENDIX C. Languages

By default, the graphical user interface of *MultiDendrograms* is shown in English. Currently, it is possible to choose between the following languages:

- English
- Catalan
- Spanish

In the `ini` folder there is a language file (e.g. `english.l`) for each of the languages available. The selection of the language is made in the first lines of the configuration file `dendo.ini` found in the `ini` folder. To change the selected language, just open `dendo.ini` in an editor, uncomment the desired language and comment the rest (comments start with a '#' character). For example, to choose English the contents should be:

```
# Language
language = ini/english.l
#language = ini/catalan.l
#language = ini/spanish.l
```

To translate *MultiDendrograms* to other languages, just create a new language file `newlanguage.l`, containing the translation of all the lines in any other language file, and add the corresponding line to the configuration file `dendo.ini`. If you send us your new language file, we can include it in future versions of *MultiDendrograms*.

MultiDendrograms - Manual

APPENDIX D. History of changes

MultiDendrograms 2.1

- Export dendrograms to Newick format
- Show calculation progress
- Improved GUI
- Improved performance

MultiDendrograms 2.0

- Completely new multiplatform (Windows, Linux, MacOS, etc.) application
- Added Graphical User Interface (GUI)
- Control of the dendrogram appearance
- Navigation through the dendrogram details
- Accepts distance and similarity matrices
- Export dendrograms to JPG, PNG and EPS
- Calculation of ultrametric deviation measures

MultiDendrograms 1.0

- Windows command-line application to compute multidendrograms
- Windows command-line application to compute ultrametric matrices
- Windows command-line application to generate EPS plots

MultiDendrograms - Manual

APPENDIX E. License

MultiDendrograms is free software; you can redistribute it and/or modify it under the terms of the GNU Lesser General Public License as published by the Free Software Foundation; either version 2.1 of the License, or (at your option) any later version.

MultiDendrograms is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Lesser General Public License for more details.

You should have received a copy of the GNU Lesser General Public License along with *MultiDendrograms*; if not, see <http://www.gnu.org/licenses/lgpl.html>.