BIONUMERICS Tutorial:

# Clustering Omnilog carbon source oxidation data

## 1  Aim

Cluster analysis is a collective noun for a variety of algorithms that have the common feature of visualizing the hierarchical relatedness between samples by grouping them in a dendrogram or tree. In this tutorial we will create a dendrogram based on trend data. We will also see how to alter the layout of the dendrogram and how to export the cluster analysis to use it in a publication, presentation, etc.

## 2  Example data

1. Import the Omnilog .csv trend data files as described in the tutorial: "Importing carbon source oxidation data from Omnilog csv files".

Each csv file contains information about the utilization of carbon substrates of a certain strain.

## 3  Comparison window

1. In the *Database entries* panel of the *Main* window, select all entries in the database for which Omnilog trend curves are present: use the **Ctrl-** key to select the entries, or alternatively right-click on the **Omnilog** column in the *Experiment presence* panel and select **Select entries with experiment**.

2. Highlight the *Comparisons* panel in the *Main* window and select **Edit** > **Create new object...** ( + ) to create a new comparison for the selected entries.

3. Click on the ⊚ next to the experiment name **Omnilog** in the *Experiments* panel to display the defined parameter(s) in the *Experiment data* panel (see Figure 1).

4. Select **TrendData** > **Show parameter values colors** to display the values of the parameter together with the color as defined in the *Trend type* window.

5. Select a parameter in the *Experiment data* panel and select **TrendData** > **Sort entries by parameter value** ( ⇓ ).

The entries are sorted according to increasing value of the selected parameter.

6. A tab-delimited text file of the entries and trend data values contained in the comparison can be exported with **TrendData** > **Export character table**.
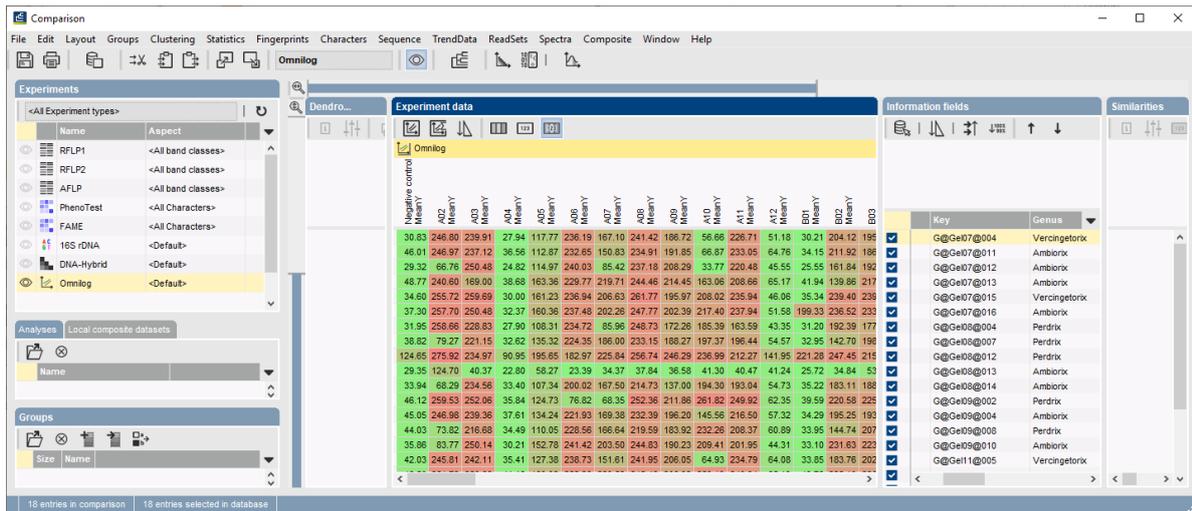
**Figure 1:** The *Comparison* window.

# 4 Cluster analysis

Cluster analysis is a two-step process. First, all pairwise similarity values are calculated with a **similarity coefficient**. Then, the resulting similarity matrix is converted into a dendrogram with a **clustering algorithm**. Although in practice these steps are performed together, they each require their own comparison settings.

1. Make sure **Omnilog** is selected in the *Experiments* panel and select ***Clustering*** > ***Calculate*** > ***Cluster analysis (similarity matrix)...***.

The first step deals with the similarity coefficient for the calculation of the similarity matrix (see Figure 2).

In case of trend data, two groups of coefficients can be applied for the calculation of the similarity matrix:

- Curve based coefficients: provide similarities based upon the original data points of the curves.

- Parameter based coefficient: measures the similarity by comparing the values of the parameter(s), defined in the *Trend type* window.

2. Select a coefficient from the list, e.g. ***Parameter similarity*** and press <***Next***>.

In step two the options related to the clustering algorithms are grouped. Under ***Method***, the clustering algorithm to be applied on the similarity matrix can be selected. A ***Dendrogram name*** can be entered in the corresponding text box. By default, the name of the experiment type will be used.

3. Select ***UPGMA*** and select ***Cophenetic correlation*** from the ***Branch quality*** list (see Figure 3).

The ***Cophenetic Correlation*** is a parameter that expresses the consistency of a cluster. This method calculates the correlation between the dendrogram-derived similarities and the matrix similarities. The value is calculated for each cluster thus estimating the faithfulness of each sub-cluster of the dendrogram.
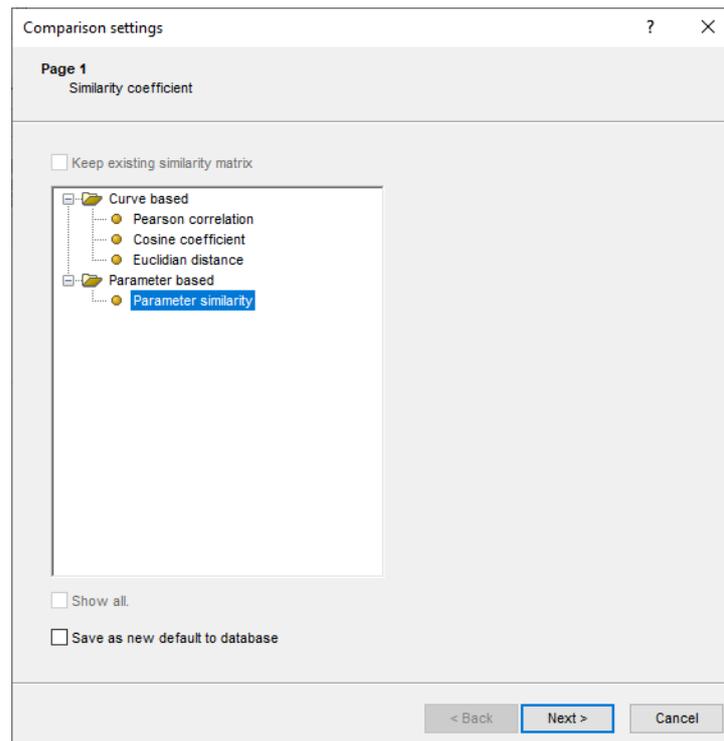
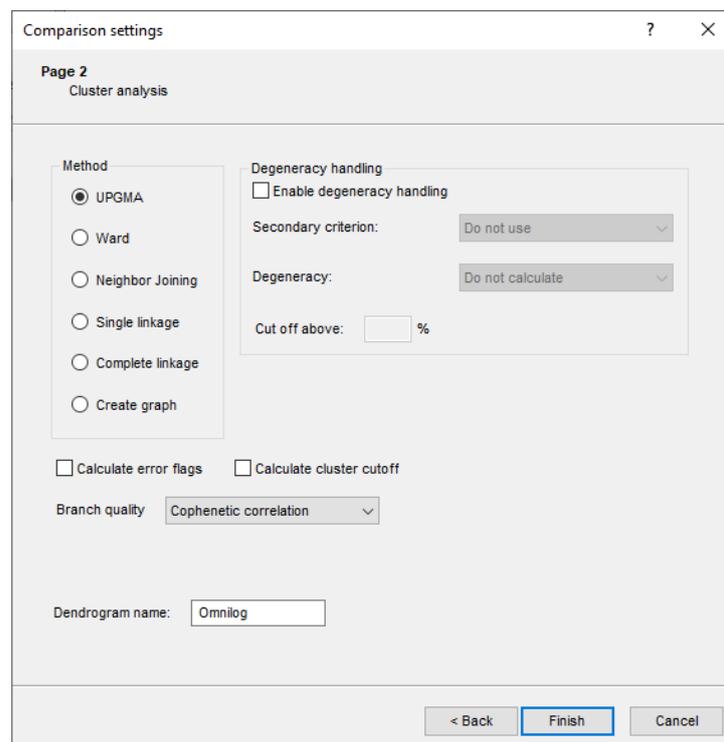**Figure 2:** Select similarity coefficient.



**Figure 3:** Select clustering algorithm.

4. Press <***Finish***> to start the cluster analysis.

During the calculations, the program shows the progress in the *Comparison* window's caption (as a percentage), and there is a green progress bar in the bottom of the window.

When finished, the dendrogram and the similarity matrix are displayed in their corresponding panels. The cluster analysis is listed in the *Analyses* panel of the *Comparison* window (see Figure 4).
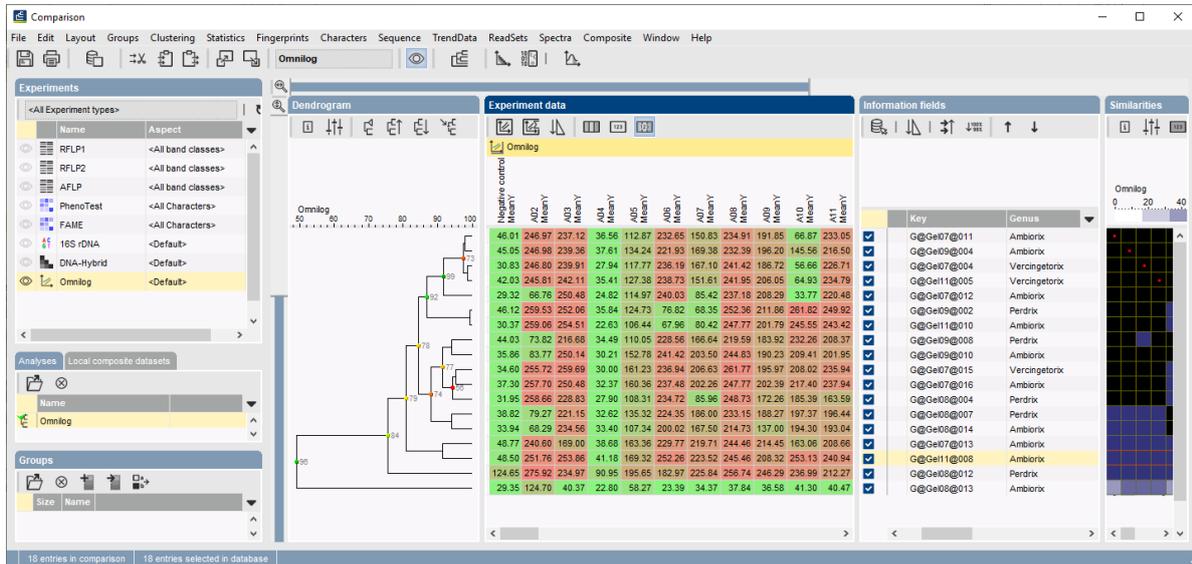


**Figure 4:** The *Comparison* window.

The ***Cophenetic correlation*** is shown at each branch, together with a colored dot, of which the color ranges between green-yellow-orange-red according to decreasing cophenetic correlation. This makes it easy to detect reliable and unreliable clusters at a glance.

5. Press the **F4** key to clear any selection in the database.

6. Left-click on the dendrogram to place the cursor on any node or tip (where a branch ends in an individual entry).

7. To select entries in a cluster, click on the node of the cluster while holding the **Ctrl-** button.

8. Press ***Edit*** > ***Cut selection*** ( ⇥X , **Ctrl+X**) to remove the selected entries from the cluster analysis. Confirm the action. The dendrogram is automatically updated.

9. Select ***Edit*** > ***Paste selection*** ( ⧉ , **Ctrl+V**). The cluster analysis is recalculated automatically, and the selected entries are placed back in the dendrogram.

A branch can be moved up or down to improve the layout of a dendrogram:

10. Click the branch which you want to move up in the dendrogram and select ***Clustering*** > ***Move branch up*** ( ⧉↑ ).

11. Click the branch which you want to move down in the dendrogram and select ***Clustering*** > ***Move branch down*** ( ⧉↓ ).

To simplify the representation of large and complex dendrograms, it is possible to simplify branches by abridging them as a triangle.

12. Select a cluster of closely related entries and select ***Clustering*** > ***Collapse/expand branch*** ( ⧉ ). Repeat this action to undo the abridge operation.

13. Select ***Clustering*** > ***Dendrogram display settings...*** ( ⭭↕ ) to call the *Dendrogram display settings* dialog box.

14. Uncheck ***Show branch quality*** and press <***OK***> to remove the cophenetic correlation from the tree.

15. Select **Clustering** > **Show information** ( 🗊 ) to display a report containing the comparison settings. Close the report.

The similarity values in the *Similarities* panel are represented by shades of blue.

16. To show the values in the matrix, select **Clustering** > **Similarity matrix** > **Show values** ( 🖽 ).

17. Save the comparison with the dendrogram by selecting **File** > **Save** ( 🖫, **Ctrl+S**). Specify a name (e.g. **Omnilog**) and press <**OK**>.

# 5 Exporting and printing a cluster analysis

BIONUMERICS can export the cluster analysis as it appears in the *Comparison* window.

1. Select **File** > **Print preview...** ( 🖶, **Ctrl+P**).

The *Comparison print preview* window now appears.

2. To scan through the pages that will be printed out, use **Edit** > **Previous page** ( ◀ , **Page Up**) and **Edit** > **Next page** ( ▶ , **Page Down**).

3. To zoom in or out, use **Edit** > **Zoom in** ( 🗗, **Ctrl+Page Up**) and **Edit** > **Zoom out** ( 🗗, **Ctrl+Page Down**) or use the zoom slider.

4. To enlarge or reduce the whole image, use **Layout** > **Enlarge image size** ( ᴀA ) or **Layout** > **Reduce image size** ( Aᴀ ).

5. If a similarity matrix is available, it can be included with **Layout** > **Show similarity matrix** ( 🖿 ).

6. On top of the page, there are a number of small yellow slider bars, which can be moved.

7. To preview and print the image in full color select **Layout** > **Use colors** ( 🌐 ).

8. Export the image to the clipboard with **File** > **Copy page to clipboard** ( 🗐 ) and selecting an appropriate format.

9. If a printer is available, use **File** > **Print this page** ( 🗐 ) or **File** > **Print all pages** ( 🗐 ) to print one or all pages.

10. Select **File** > **Exit** to close the *Comparison print preview* window.