BIONUMERICS Tutorial:

# Cluster analysis of sequences

## 1 Aim

Similarity and distance-based trees (e.g. UPGMA and Neighbor joining) and phylogenetic trees (e.g. Maximum likelihood and Maximum parsimony trees) can be calculated in the *Comparison* window in BIONUMERICS, reflecting the relationships between the analyzed sequences. Different groups can be defined for the data and these groups are visualized in the trees, giving a nice overview of the relationships, the clusters present in the database and the outliers. In this tutorial we will calculate several types of trees on a set of ribosomal sequences.

## 2 Example data

As an example we will import sequences from the FASTA file `FastaSeqCL.txt` into a new or existing BIONUMERICS database. This FASTA file contains DNA sequences of the 16S rRNA gene for a set of 14 strains. The example file can be found on the download page on our website (https://www.applied-maths.com/download/sample-data, "FASTA 16S sequences").

## 3 Preparing the database

1. Create a new database (see tutorial "Creating a new database") or open an existing database.

2. Select *File* > *Import...* ( , **Ctrl+I**) to open the *Import* dialog box.

3. Choose the option *Import FASTA sequences from text files* under the *Sequence type data* item in the tree and click <*Import*> (see Figure 1).

4. The import wizard allows you to browse for one or more text files as data source. Press <*Browse*>, navigate to the folder, select the `FastaSeqCL.txt` file and press <*Open*>.

5. With the option *Preview sequences* checked, press <*Next*>.

The import wizard now displays a preview of the sequence data in the FASTA file (see Figure 2). From this preview, it is clear that the first FASTA field contains the unique strain number.
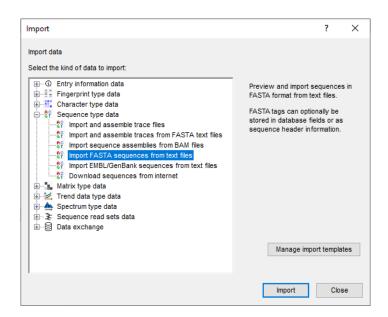
6. Press <*Next*>.

**Figure 1:** The Import tree.



**Figure 2:** Sequence preview.

The next step of the import wizard lists the templates that are present to import sequence information in the database. As this is the first time we import FASTA formatted sequences in the database, we need to create a new import template by specifying **Import rules**.

7. Click <**Create new**> to create a new import template.

8. Select "Field 1" in the list and click <**Edit destination**> or simply double-click on "Field 1". Select **Key** and press <**OK**>.

9. Double-click on "Field 4". Select **Create new** under **Entry info field** and click <**OK**>. Enter "Accession number" as name for the new information field, press <**OK**> and confirm the creation of the new field with <**Yes**>.

The grid is updated (see Figure 3).

**Figure 3:** Import rules.

10. Optionally, you can press <**Preview**> to obtain a preview of the data you are about to import (see Figure 4).



**Figure 4:** Preview of the rules.

11. Click <**Next**>.

12. Make sure **Key** is selected as **Entry link field**. Press <**Finish**>.

13. Specify a template name, e.g. **FASTA**, and optionally enter a description. Press <**OK**>.

14. Highlight the newly created template and select "Create new" as **Experiment type** (see Figure 5).

**Figure 5:** Import template.

15. Press <***Next***>.

16. Specify a sequence type name (e.g. **Ribosomal**) and press <***OK***> and confirm the action.

The *Database links* wizard page will indicate that 14 new entries will be created during import.

17. Press <***Finish***> to start the import into the database.

For 14 strains, strain information and sequences are imported in the database (see Figure 6).



**Figure 6:** The *Main* window after import of the sequences.

# 4 Comparison window

1. In the *Database entries* panel of the *Main* window, select all entries that have an associated **Ribosomal** experiment. To select all entries at once, use the **CTRL+A** shortcut combination.

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 **Ribosomal** in the *Experiments* panel to display the sequences in the *Experiment data* panel.

Initially, the sequences are not aligned and no similarity matrix exists.

# 5   Pairwise and multiple alignment

1. The similarity matrix is calculated with **Clustering** > **Calculate** > **Cluster analysis (similarity matrix)...** or by pressing ▤ and selecting **Calculate cluster analysis**.

2. Choose **Fast algorithm** under **Pairwise alignment** with default settings and press <**Next**>.

3. In *Cluster analysis* wizard page, choose **UPGMA** as the dendrogram type and press <**Finish**>.

A dendrogram is now calculated, but the sequences are still unaligned.

4. Select **Sequence** > **Multiple alignment...** ( ▦ ).

5. Leave the settings as they are (default) and press <**OK**> to start the multiple alignment. When the calculations are done, the sequences are aligned in the *Experiment data* panel (see Figure 7).



**Figure 7:** Multiple alignment.

In order to facilitate visual interpretation of global alignments there are three methods to highlight homologous regions. The first one is the *Neighbor match* representation:

6. Select **Sequence** > **Block type** > **Neighbor blocks** to show the Neighbor match representation (see Figure 8).

This representation shows bases as blocks (highlighted) if at least one of the neighboring sequences has the same base at the corresponding position. Between two different groups of consensus, a small line is drawn.

The second visualization, the *Consensus match*, first requires a consensus sequence to be present.

7. Select the root and select **Sequence** > **Create consensus of branch**.

8. Leave the default setting and press <**OK**>.

A consensus sequence of the root is shown on the header of the image panel. Bases for which there is a consensus in more than 50% of the sequences, are named in the consensus, the other
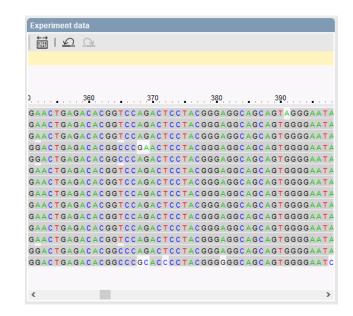
**Figure 8:** Neighbor match representation.

bases are unnamed (N).

9. Select ***Sequence*** > ***Block type*** > ***Consensus blocks*** to show the consensus match representation (see Figure 9).



**Figure 9:** Consensus match representation.

A third method, the *Consensus difference*, displays the consensus sequence in the editor caption, and only shows bases that differ from the consensus. Bases that are the same as the consensus are shown as "|" (see Figure 10).

10. Select ***Sequence*** > ***Block type*** > ***Consensus difference***.

A multiple alignment can be edited manually and is saved along with the comparison.

11. Select ***File*** > ***Save*** (🖫, **Ctrl+S**) to save the multiple alignment.

**Figure 10:** Consensus difference representation.

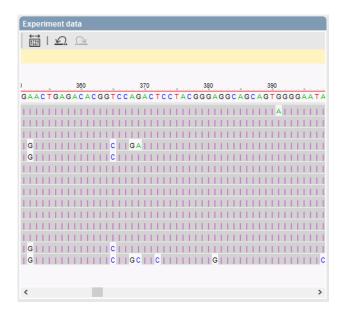12. A window pops up asking you to enter a name for the comparison. Enter "SeqAlign" and press <***OK***>.

# 6 Cluster analysis

## 6.1 Similarity or distance-based clustering

The mutual similarities between all the sequences are calculated from the aligned sequences as present in the multiple alignment.

1. Select ***Clustering*** > ***Calculate*** > ***Cluster analysis (similarity matrix)...*** or press 🗲 and select ***Calculate cluster analysis***.

2. Choose ***Multiple alignment based*** with the default settings and press <***Next***>.

3. In the second step of the wizard, select ***Neighbor Joining*** as dendrogram type and click <***Finish***>.



**Figure 11:** Neighbor joining based on multiple alignment.

4. Save the comparison with **File** > **Save** (🖫, **Ctrl+S**).

## 6.2 Phylogenetic clustering

### 6.2.1 Introduction

The most widely used phylogenetic clustering methods are maximum parsimony (Fitch, 1971) and maximum likelihood (Felsenstein, 1981). In both methods, an evolution is reconstructed based upon sequence data by optimizing a certain criterion.

The *maximum parsimony* method tries to find a tree that explains the sequence diversity with a minimum number of total mutations needed (i.e. the most parsimonious tree). The branch lengths of the tree reflect the number of mutations along the branches.

The *maximum likelihood* method is based on a probabilistic model for base substitution. A tree is searched for that has the highest likelihood, i.e. the probability that the given sequences are the result of an evolution along that tree, following the assumed probabilistic module. The branch lengths of the tree correspond to evolutionary time.

The Neighbor joining method can also be classified under phylogenetic clustering methods, but in this tutorial it is covered in 6.1 since it uses a distance matrix as input.

### 6.2.2 Maximum parsimony

5. In the *Comparison* window with the multiple alignment calculated press the <**F4**> to clear the selection.

6. Select CL007, CL008, CL013 using the check boxes and choose **Groups** > **Create new group from selection** ( 📑 , **Ctrl+G**). Specify a group name (e.g. **Group 1**) and press <**OK**> to confirm the creation of the group.

7. Press <**F4**> again to clear the selection.

8. Select CL002, CL003, CL005, CL009, CL011, CL012, and choose **Groups** > **Create new group from selection** ( 📑 , **Ctrl+G**). Specify a group name (e.g. **Group 2**) and press <**OK**> to confirm the creation of the group.

9. Press <**F4**> again to clear the selection.

10. Select CL010 and add the entry to a third group (e.g. **Group 3**).

11. Select **Clustering** > **Calculate** > **Advanced cluster analysis...**.

12. Name the cluster analysis **Max parsimony**, choose **Ribosomal** as experiment and **Maximum parsimony tree** as predefined template (see Figure 13).

To view and modify the settings of the selected template check the option **Modify template settings for new analysis**.

13. Press <**Next**> to start the calculations.

The resulting tree is displayed in *Advanced cluster analysis* window (see Figure 14). The parsimony (the total number of base conversions over the tree) is given in the *Statistics* tab. The entries are represented in the colors of the groups we have defined earlier.

14. Select **Display** > **Display settings**, check **Show node labels** in the *Node labels and sizes* tab, and check **Show branch labels** in the *Branch labels and sizes* tab. Press <**OK**>.

**Figure 12:** The *Comparison* window with groups defined.



**Figure 13:** Maximum parsimony template.

15. Select **Display** > **Zoom to fit** to optimize the view of the tree in the current window.

16. By adjusting the zoom bar, the node sizes can be increased/decreased.

The right bottom information panel displays the cluster analysis settings, the branch properties and the total parsimony of the tree in the statistics tab.

17. A node or branch can be selected by clicking on them, or several nodes/branches by holding the **Shift**-key while clicking.

We will now perform a bootstrap analysis on the Maximum Parsimony tree.

18. Go to **Edit** > **Compute statistics** to call the *Statistics* dialog.

19. Check **Perform resampling**, set the resampling strategy on **Bootstrap resampling** with 100 samples and click <**OK**>.
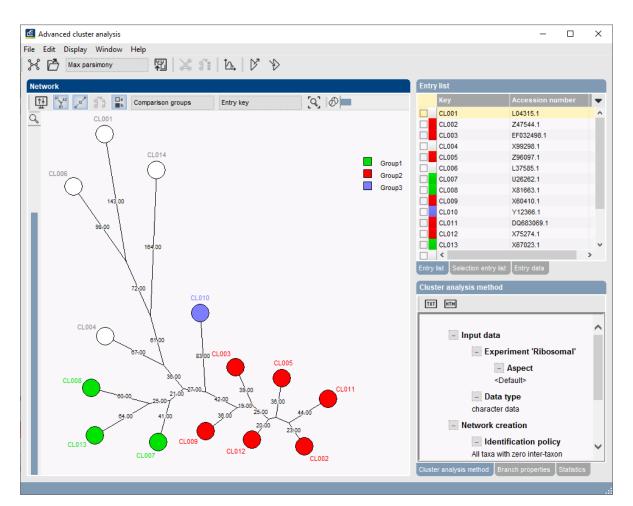
**Figure 14:** Maximum parsimony tree.

We will now display a rooted tree from the unrooted tree.

20. Select the CL010 entry to act as root, go to **Edit** > **Determine root** and choose **Use selected node as root position**.

Once a rooted tree has been calculated, it can be visualized in the *Comparison* window:

21. Select **File** > **Show dendrogram in comparison** and confirm.

In a phylogenetic context, a dendrogram is sometimes represented as a real tree with a stem and branches:

22. This can be achieved by selecting **Clustering** > **Display rendered tree** in the *Comparison* window and pressing <**OK**> (see Figure 15).

23. Close the rendered tree representation.

### 6.2.3 Maximum likelihood

To reduce the calculation time of the maximum likelihood tree we will delete the 4 entries that do not belong to a group:

24. Press the F4 key in the *Comparison* window to clear the selection. Select entry CL004, CL006, CL001 and CL014 using the check boxes and select **Edit** > **Cut selection** (⇥✗, **Ctrl+X**). Confirm the deletion of the four entries from the comparison.
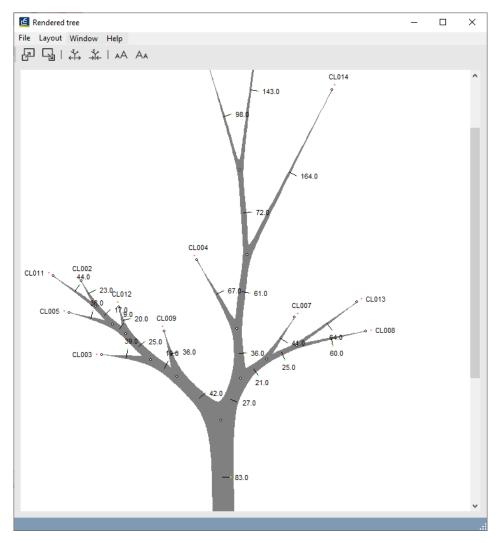
**Figure 15:** Rendered tree representation.

25. Select *Clustering* > *Calculate* > *Advanced cluster analysis...*.

26. Name the cluster analysis **Max likelihood**, choose *Ribosomal* as experiment and *Maximum likelihood tree* as predefined template (see Figure 16).

✎ To view and modify the settings of the selected template check the option *Modify template settings for new analysis*.

27. Press <*Next*> to start the calculations.

The resulting tree is displayed in *Advanced cluster analysis* window (see Figure 17). The entries are represented in the colors of the groups we have defined earlier.

28. Select *Display* > *Display settings*, check **Show node labels** in the *Node labels and sizes* tab. Press <*OK*>.

29. Select *Display* > *Zoom to fit* to optimize the view of the tree in the current window.

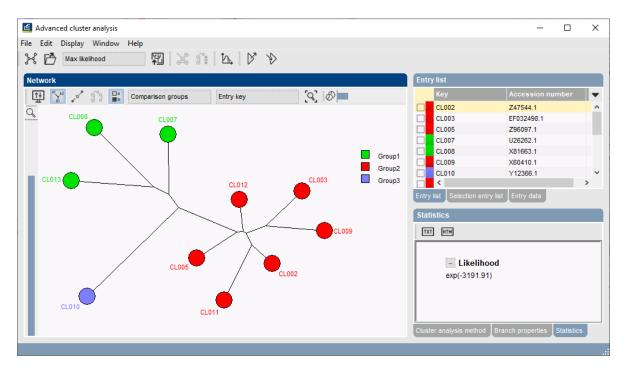30. By adjusting the zoom bar, the node sizes can be increased/decreased.

**Figure 16:** Maximum likelihood template.



**Figure 17:** Maximum likelihood tree.