

# Phylogeny of carnivorous plants

## comparison of basic algorithms on real data

K. Tucek

January 25, 2017

# Goals

This project was made as a part of fundamental course of bio-algorithms. The results are by no means to be interpreted as scientifically relevant! Our aim is to:

- Try basic DNA-sequence similarity algorithms on a real problem and real data.
- To see how these simple approaches behave compared to actual scientific results.
- To compare the behaviour of some basic and sometimes naive algorithms.
- To learn something about carnivorous plants.

## Topic choice.

We have chosen to try the construction of phylogeny trees of some representants of carnivorous plants. The class of carnivorous plants was picked quite arbitrarily.

This problem seems interesting since the carnivorous properties of these plants are somehow exotic and structurally interesting, since hypotheses about structural relationships being involved in phylogeny of carnivorous plants come as natural.

# Classification of Carnivorous Plants

We have picked one representants from every genus of carnivorous plants. We have picked multiple representants from the Drosera genera. The following slides sum up the order/family/genera hierarchy (nonbranching families are ommited):

- Caryophyllales
- Ericales
- Lamiales
- Oxalidales
- Poales

# Classification of Carnivorous Plants

- Caryophyllales
  - Aldrovanda (1 sp, Everywhere) (rare underwater plant)
  - Droseraceae
    - Dionaea (1 sp, NAm) (pitcher flytrap)
    - Drosera (233+ sp, Ev) (sticky flypaper)
    - Drosophyllum (1 sp, Spain) (suffocating flypaper)
  - Nepenthes (131+ sp, SEAs) (classical mawed pitcher plant)
  - Triphyophyllum (1 sp, Africa) (climbing vine flypaper)
- Ericales
  - Sarraceniaceae
    - Darlingtonia (1 sp, NAm) (pitfall pitcher)
    - Heliamphora (25+ sp, SAm) (heli-amphora)
    - Sarracenia (15 sp, NAm) (trumpet pitchers)
  - Roridula (2 sp, SAf) (non-digestive resin (!) flypaper)

# Classification of Carnivorous Plants

- Lamiales
  - Byblis (8+ sp, Au) (carnivorous via hosted insects)
  - Philcoxia (3 sp, SAm) (sand-underground flypaper)
  - Lentibulariaceae
    - Genlisea (27+ sp, SAm, SAf) (two types of pitfall traps)
    - Pinguicula (102+ sp, Ev) (flypaper with efficient digestion)
    - Utricularia (235+ sp, Ev) (underwater bladder traps)
- Oxalidales
  - Cephalotus (1 sp, Au) (another pitcher with actual 'teeth')
- Poales
  - Brocchinia (2 sp, SAm) (minimalistic pitfall pitcher)
  - Catopsis (1 sp, Am) (leafed pitfall trap)

# Trap Types

Basic trap types are the following:

- snapping traps - such as the fly-trap
- sticky 'flypaper' traps - plants with leaves covered by a sticky extract
- pitcher traps - traps relying on the insect falling into a liquid-filled pitch and drowning inside
- other pitfall traps - other traps which rely on insect's getting in and not being able to get out
- special underwater traps

## Methods overview

- We compare the organisms by the so-called 'rbcL' plastid gene (ribulose-1,5-bisphosphate carboxylase/oxygenase), which is one of substantive photosynthesis genes. (As used in paper by V. A. Albert, S. E. Williams and M. W. Chase)
- Distance matrix is constructed using various distance functions.
- Phylogeny tree is computed by the biopython's upgma algorithm.



# Compared Algorithms

- GenCompress compression distance.
- Simple levenstein dinstance using gap-aware rating.
- Uni/Bi - directional histogram.
- Uni/Bi - directional Dotplot area with 2/3 filtering window ratio.

# Compared Algorithms

## Bi/Uni directionality

- In order to take into account the possibility of complement-strand crossovers, we implement the histogram and dotplot algorithms in bidirectional variants.
- The bi-directional versions consider complement's patterns as if these were found in the forward algorithm.
- This consideration is mostly theoretical since We assume that crossovers with the complement strands are not of high interes inside of gene sequences. Also, this penalty-less rating in reverse strand is naive.

# Compared Algorithms

## Dotplot vs histogram algorithms

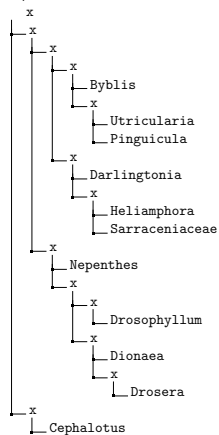
These two algorithms perform basically the same thing. However, unlike simple histogram algorithm, the dotplot surface algorithm may take into account:

- Locality of patterns.
- Mutations.
- Longer window, wince the results do not need to be explicitly stored.

On the other hand, the dotplot-surface takes reccuring patterns into account  $O(n^2)$  times unlike the simple histogram algorithm. Also, dotplot algorithm is much more prone to badly tuned parameters.

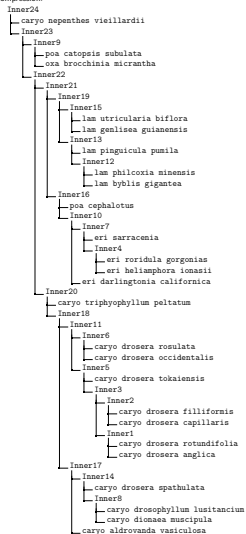
# Expected results

expected:

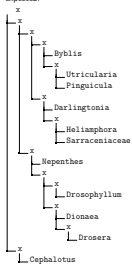


# GenCompress distance

compression:

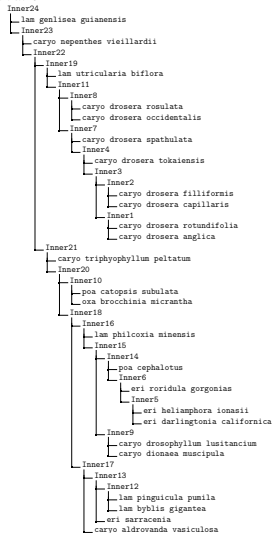


expected:

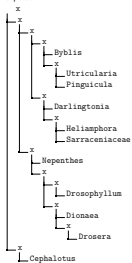


# Levenstein distance

levensteinUniDir:

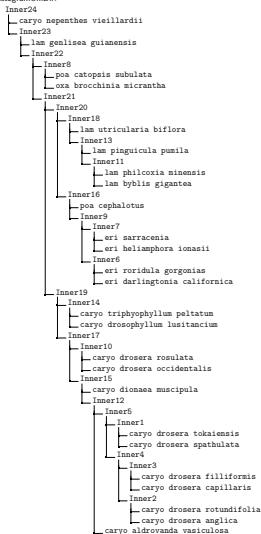


expected:

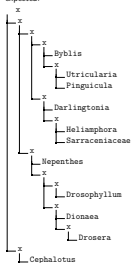


# Unidirectional histogram

histogramUniDir:

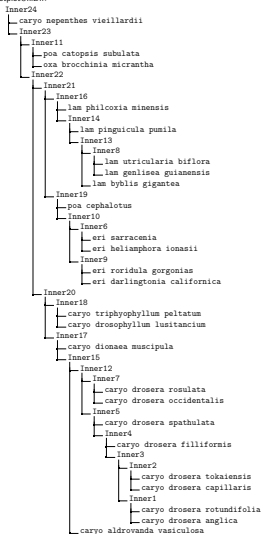


expected:

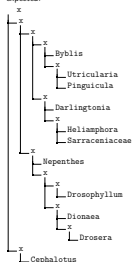


# Unidirectional dotplot surface area

dotplotUnDir:



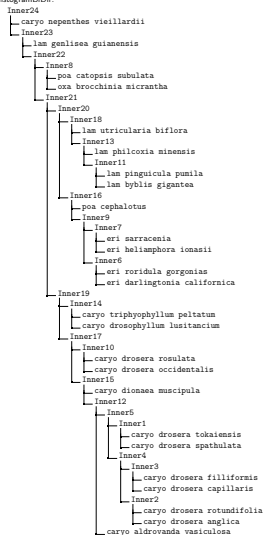
expected:



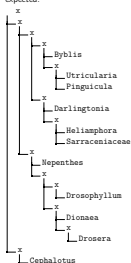


# Bidirectional histogram

histogramBiDir:



expected:













## Data analysis

- It turns out that manual analysis of the tree similarity is almost impossible.
- The best results are given by the compression algorithm and then the surface dotplot area algorithm. The dotplot algorithm actually seems to give a bit more consistent results.
- Dotplot area algorithm and histogram give quite similar results. The results of histogram are a bit worse.
- Levenstein distance gives quite garbled trees.
- Graphical dotplots show that similarity of samples is exceptional (approximately 70% of DNA has 100% correspondence with window of length 20).
- As the consequence, uni and bi directional versions of the histogram give exactly the same results.
- Nepenthes family results turn out to be crucial for result analysis. Unfortunately, the only present sample does not match its expected position in any of the trees.

## Data analysis

- The phylogenetic trees mostly correspond to those the original paper proposes.
- Also, the sample corresponds to the standard taxonomy.
- The data sample is too small to draw other conclusion with respect to evolution of trap types.
- Our sample contains some genes which were not available to the original study. On the other hand, the original study contains much more sophisticated system of representants of the Caryophyllales order. This results in an effective impossibility of drawing exact results.
- GenCompress is unable digest *Drosera Spathulata* for some reason.



- The original study:  
[http://www.botany.wisc.edu/courses/botany\\_400/papers/Carnivorou](http://www.botany.wisc.edu/courses/botany_400/papers/Carnivorou)
- <http://www.carnivorousplants.org/cp/WhatAreCPs.php>
- <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.00>
- <https://www.czplants.com>
- <http://www.masozravky.com/>
- [http://botany.org/Carnivorous\\_Plants/](http://botany.org/Carnivorous_Plants/)