# Inferring phylogenetic trees

#### László Kozma Lkozma@cc.hut.fi

Helsinki University of Technology T-61.6070 Special course in bioinformatics: Modeling of proteomics data

2. 4. 2008

<ロ> <日> <日> <日> < 三> < 三> < 三> 三 のへで

#### References

- Notes on phylogenetics, School of Computer Science, Tel-Aviv University http://www.cs.tau.ac.il/ rshamir/algmb/00/scribe00/html/lec08/lec08.html
- Kim, Warnow: Tutorial on Phylogenetic Tree Estimation http://kim.bio.upenn.edu/ jkim/media/ISMBtutorial.pdf
- Lecture notes on phylogenetic tree construction, Frank Olken, Berkeley PGA http://pga.lbl.gov/Workshop/April2002/lectures/Olken.pdf
- 4 L. A. Salter, Algorithms for Phylogenetic Tree Reconstruction (2000)
- **5** J. Felsenstein: Inferring Phylogenies (2003, Sinauer Associates)
- 6 Google, Wikipedia, ...

# **Phylogenetics**

- Fundamental idea in biology: *All organisms living on Earth today share common ancestors* (P.-L. de Maupertuis 1745; I. Kant 1790; Ch. Darwin, 1859)
- Evolution happened through:
  - mutation: DNA sequence changed due to single base changes, deletion/insertion
  - selection bias
  - speciation: physical separation into groups
- Phylogenetics: study the evolutionary relations between organisms
- Goal: recover order of speciation (evolution history)
- Evolution history can be (roughly) depicted as a tree

# **Phylogenetics**

- Taxon: organism, species, gene, protein, etc.
- Infer tree from observed features of taxa:
  - morphological characters (used in classical taxonomies): number of legs, size, color, swims/flies, etc.
  - DNA sequences (alignment of selected sequences that are conserved in many organisms, example: Ribosomal RNA 16S)
  - protein sequences
  - metabolic pathways data, etc.
- Often different data produce similar trees (despite the fact that most genetic variation does not cause different external morphology) ⇒indirect evidence for evolution

 $\mathcal{O} \mathcal{Q} \mathcal{O}$ 



Figure: Phylogenetic tree.

- leaves: taxa, species observable today
- internal nodes: hypothetical common ancestors

#### • Evidence of extinct common ancestors:

- fossil data:
  - rarely available
  - mostly recent
  - unreliable
- Recapitulation theory (Ernst Haeckel, 19th century): "ontogeny recapitulates phylogeny"
- Only data from current taxa used to infer tree
- Other data can still be used for verifying predictions

# Phylogenetics

Why do phylogenetics:

- "Tree of Life"
- understand lineage of species
- understand how a function evolved
- design drugs, vaccines
- assist epidemiology
- study of relationships between parasite/host

theory applicable to other fields, for ex. natural languages

# **Phylogenetics**

Example applications:

- Human, Chimpanzee, Gorilla tree
- analysis of mitochondrial DNA of 182 people, one common ancestor found who lived in Africa
- HIV origins
- Influenza virus phylogeny, matching frozen ancient samples

Sac



Figure: Rooted/Unrooted tree.

- In rooted tree one node given special significance (common ancestor)
- Can be easily picked if molecular clock hypothesis holds
- Edge length can indicate genetic distance ( =? time)

Difficulties with data:

- Noise
- Hybridisation between species that were not nearest neighbors
- Recombination
- Convergent evolution
- · Gene duplication, horizontal gene transfer
- Conserved sequences
- Extinct species

Methods for inferring the tree (binary/bifurcating):

- Distance-based methods
  - UPGMA (agglomerative clustering)
  - Neighbor Joining
- Character-based methods
  - Maximum Parsimony
  - Maximum Likelihood
  - Bayesian methods, etc.

- ロ > - 4 回 > - 4 回 > - 4 回 > - 9 9 9 9 9

Inferring phylogenetic trees is computationally expensive:

- Most of the criteria lead to NP-hard problems, therefore heuristics/approximations are used.
- Optimal sequence alignment itself NP-hard
- Exhaustive search:

$$\frac{(2n-3)!}{2^{n-2}(n-2)!}$$
 rooted

 $\frac{(2n-5)!}{2^{n-3}(n-3)!}$  unrooted

trees are possible with n given leaves.

#### Align relevant parts of DNA sequences:

Species	A:	ATGCTCAGT
Species	в:	ATGCCCAGG
Species	C:	ATGCTCAGC
Species	D:	CGCCTCAGA
Species	Е:	GTGCTCAGT

#### Similarity matrix:

	А	В	С	D	Е
А	-	0.3	0.4	0.55	0.7
В		-	0.24	0.45	0.5
С			-	0.4	0.2
D				-	0.5
Е					-

- Usual axioms for distance metrics must hold
- Possible distance measure: observed percent difference
- Extension: weighted using different transition probabilities
- Similar approaches for protein sequences
- Other distance measures: edit distance, more complex measures

Unweighted Pair Group Method with Arithmetic Mean (UPGMA):

- let d<sub>ij</sub> be the distance between two nodes
- $C_i$  a cluster of nodes, with  $n_i$  elements
- the distance between two clusters  $C_i$ ,  $C_j$  is:

$$D_{ij} = \frac{1}{n_i n_j} \sum_{p \in C_i, q \in C_j} d_{pq} \tag{1}$$

(ロ) < 同) < 三) < 三) < 三) 三 の(()</p>

## **UPGMA**

(ロ) < 同) < 三) < 三) < 三) 三 の(()</p>

Steps for building the tree:

- 1 Construct distance matrix between taxa.
- 2 Add each taxa to the tree as a leaf.
- 3 Assign each taxa to its own cluster.
- Find the pair of clusters with the shortest distance  $D_{ij}$ .
- **6** In the matrix replace clusters i and j with new cluster k.
- **6** Calculate distances from k to all other clusters as  $D_{kl} = \frac{n_i D_{il} + n_j D_{jl}}{n_i + n_j}$
- Add branches in the tree from k to i and from k to j such that distance to leaves is  $D_{ij}/2$ .
- 8 Repeat from 4. until only root node left.

# **UPGMA**

		А	В	С	D	Е	F
Example	Α						
	В	2					
	С	4	4				
	D	6	6	6			
	Е	6	6	6	4		
	F	8	8	8	8	8	



イロン イロン イヨン イヨン

 $\mathcal{O} \mathcal{Q} \mathcal{O}$ 

€



Figure: Resulting tree.

Sac

#### UPGMA

- Simple to implement, fast
- Not very accurate
- Estimates branch length
- Assumes that rate of evolution is constant
- Assumes distances are additive
- No statistical evolutionary model
- Does not give quality of tree

## Character-based methods

- Input: n species, m characters/features (an n\*m matrix)
- Can be the aligned DNA/protein sequences as seen before
- Output: A phylogenetic tree that maximizes some target function
- Possible other input: relative importance (weight) of characters, distribution of mutations, etc.
- We make two rather strong assumptions:
  - Characters mutually independent
  - After two species diverge, they continue to evolve independently

## Character-based methods

Maximum Parsimony

- Predicts evolutionary tree with minimum number of evolutionary steps required to generate observed variation
- Parsimony, Occam's Razor, MDL, "most defensible theory, least number of assumptions"
- Not necessarily meaningful biologically
- Naive algorithm:
  - Construct all possible trees
  - Assign scores to all trees
  - Choose tree with best score

• Parsimony score: 
$$S(T) = \sum_{(v,u)\in E(T)} |\{j: v_j \neq u_j\}|$$

< ロ > < 同 > < 三 > < 三 >

Jac.



Figure: Two possible trees.

Parsimony suggests that the first tree is better, as it has fewer mutations.

Parsimony algorithms have two components:

- **1** Search through the space of trees
- Assign minimum number of changes to a given tree topology
  - First problem NP-complete:
    - use heuristics, hill-climbing, local search
    - branch-and-bound
  - Second problem: Fitch's algorithm
    - any state can convert to any state
    - locations are uniform
    - each type of transition has same cost
  - Weighted variants exist

Fitch's algorithm:

- Traverse tree from leaves to root, determining possible states for each internal node
- 2 Traverse tree from root to leaves, picking ancestral states for internal nodes

Run separately for each character location.

State  $r_j$  of node j with parent i:  $r_j = r_i, ifr_i \in R_j$ arbitrary state  $\in R_j, otherwise$ 

( ) < </p>

∍

 $\mathcal{O} \mathcal{Q} \mathcal{O}$ 





Figure: Resulting tree.

- Advantages:
  - Models evolutionary history
  - Gives score of tree
- Disadvantages:
  - Does not estimate branch length.
  - Computationally expensive
  - Biased tree under some conditions.

Algorithm for finding best tree:

- for each possible tree
  - for each character location in the alignment
    - compute the likelihood of the tree, given data and evolution model

Choose tree that has highest likelihood.  $L(tree|data) \propto P(data|tree)$ 

explanation from:

http://noble.gs.washington.edu/ noble/genome373/lectures

< ロ > < 回 > < 三 > < 三 > < 三 > へ 回 > < 回 > < < つ へ ()

Questions left open:

- Heuristics to speed up the algorithm...
- How to find optimal branch lengths?

- Advantages:
  - Takes evolutionary model into account
  - More accurate than the other methods
  - All the sequence information is used
  - Evaluates all possible trees
- Disadvantages:
  - Very slow.
  - Impractical for analyzing large data sets.

# Summary

- What is phylogenetics?
- What are phylogenetic trees?
- Information used for constructing phylogenetic trees
- Methods for building phylogenetic trees
  - Distance based: UPGMA
  - Character based: Parsimony, Maximum Likelihood