

EVOLUTIONARY DISTANCES INFERRING PHYLOGENIES

Luca Bortolussi¹

¹Dipartimento di Matematica ed Informatica
Università degli studi di Trieste
`luca@dmf.units.it`

Trieste, 28th November 2007

OUTLINE

- 1 INFERRING PHYLOGENIES
- 2 OPTIMIZATION PROBLEMS ON TREES
- 3 LEAST SQUARE METHODS
- 4 CLUSTERING METHODS

OUTLINE

- 1 INFERRING PHYLOGENIES
- 2 OPTIMIZATION PROBLEMS ON TREES
- 3 LEAST SQUARE METHODS
- 4 CLUSTERING METHODS

RECONSTRUCTING HISTORY OF LIFE

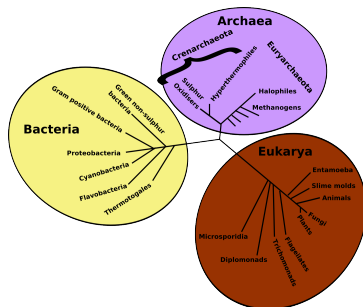
WHAT MEANS “PHYLOGENETIC INFERENCE”?

All species on Earth come from a common ancestor.

If we have data from a pool of species, we wish to reconstruct the history of speciation events that lead to their emergence:

We want to find the phylogenetic tree giving this information!

This is an hard task, because data is often **incomplete** (we lack information about most of the ancestor species) and **noisy**.



METHODS TO INFER PHYLOGENY

APPROACHES TO PHYLOGENY

- Distance-based methods
- Parsimony methods
- Likelihood methods
- Bayesian inference methods

DISTANCE-BASED METHODS

Given a matrix of pairwise distances, find the tree that explains it better. Several algorithms:

- UPGMA (clustering methods)
- Neighbor Joining
- Fitch-Margolias (sum of squares methods)

AN EXAMPLE: PRIMATES

DNA FROM PRIMATES

Tarsius	AAGTTTCATTGGAGCCACCACTCTTATAATTGCCCATGGCCTCACCTCCT...
Lemur	AAGCTTCATAGGAGCAACCATTCTAATAATCGCACATGGCCTTACATCAT...
Homo Sapiens	AAGCTTCACCGGCGCAGTCATTCTCATAATCGCCACGGGCTTACATCCT...
Chimp	AAGCTTCACCGGCGCAATTATCCTCATAATCGCCACGGACTTACATCCT...
Gorilla	AAGCTTCACCGGCGCAGTTGTTCTTATAATTGCCACGGACTTACATCAT...
Pongo	AAGCTTCACCGGCGCAACCACCCTCATGATTGCCCATGGACTCACATCCT...
Hylobates	AAGCTTACAGGTGCAACCGTCCTCATAATCGCCACGGACTAACCTCTT...
Macaco Fuscata	AAGCTTTTCGGGCGCAACCATCCTTATGATCGCTCACGGACTCACCTCTT...



DISTANCE MATRIX

0.00	0.29	0.40	0.39	0.38	0.34	0.38	0.37	Tarsius
0.29	0.00	0.37	0.38	0.35	0.33	0.36	0.34	Lemur
0.40	0.37	0.00	0.10	0.11	0.15	0.21	0.24	Homo Sapiens
0.39	0.38	0.10	0.00	0.12	0.17	0.21	0.24	Chimp
0.38	0.35	0.11	0.12	0.00	0.16	0.21	0.26	Gorilla
0.34	0.33	0.15	0.17	0.16	0.00	0.22	0.24	Pongo
0.38	0.36	0.21	0.21	0.21	0.22	0.00	0.26	Hylobates
0.37	0.34	0.24	0.24	0.26	0.24	0.26	0.00	Macaco Fuscata



OUTLINE

- 1 INFERRING PHYLOGENIES
- 2 OPTIMIZATION PROBLEMS ON TREES
- 3 LEAST SQUARE METHODS
- 4 CLUSTERING METHODS

WHAT IS AN OPTIMIZATION PROBLEM?

TWO INGREDIENTS

- 1 A **search space** S (possibly constrained)
- 2 A function $f : S \rightarrow \mathbb{R}$ to optimize:
find \bar{x} such that $f(\bar{x}) = \max_{x \in S} f(x)$.

- If S is discrete (integers, graphs, **trees**), then we talk of combinatorial optimization.
- If S is continuous (\mathbb{R}), then we talk of continuous optimization.

BAD NEWS... (OR GOOD ONES?)

“Interesting” combinatorial optimization problems are usually \mathcal{NP} -hard.

THE SPACE OF TREES

The search space is the **space of trees**, usually with **branch lengths**. Branch lengths can be optimized for a given tree topology in an easy way. **The bottleneck is the identification of the correct tree topology!**

TREE TOPOLOGIES

- Rooted vz Unrooted
- Labeled vz Unlabeled (leaves)
- Bifurcating vz Multifurcating

HOW MANY TOPOLOGIES ARE THERE?

We count the **rooted labeled bifurcating trees**.

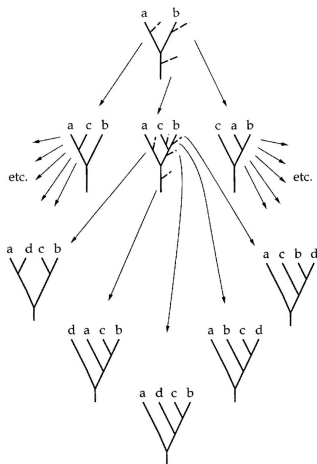
KEY PROPERTY

Each tree with n tips (leaves) can be obtained in a **unique** way from a tree with k leaves **adding in sequence** the remaining $n - k$ leaves.

In fact, reversing the sequence and removing the $n - k$ leaves, one obtains a unique tree with k leaves.

HOW MANY TOPOLOGIES ARE THERE?

We count the **rooted labeled bifurcating trees**.



TREES WITH n LEAVES

$$NT_{r,l,b}(n) = 3 \cdot 5 \cdots (2n-3) = \frac{(2n-3)!}{2^{n-2}(n-2)!}$$

HOW MANY TOPOLOGIES ARE THERE?

We count the **rooted labeled bifurcating trees**.

Species	Number of Trees
3	3
4	15
5	105
6	945
7	10,395
8	135,135
9	2,027,025
10	34,459,425
15	213,458,046,676,875
20	8,200,794,532,637,891,559,375
30	$4,9518 \times 10^{38}$
40	$1,00985 \times 10^{57}$
50	2.75292×10^{76}

TREES WITH n LEAVES

$$NT_{r,l,b}(n) = \frac{(2n-3)!}{2^{n-2}(n-2)!}$$

SEARCHING THE TREE SPACE: HEURISTIC METHODS

- The \mathcal{NP} -hardness of optimization problems on the tree space calls for heuristic solutions.
- **Heuristic optimization algorithms** usually search the space by following a **trajectory** that hopefully will hit good solutions, even the optimal one (trajectory-based heuristics).
- The next point of the trajectory is chosen in the **neighbor** of the current one.

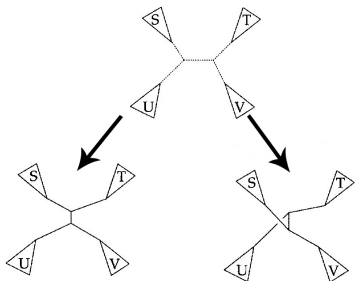
NEIGHBORHOODS FOR THE TREE SPACE

- Nearest-neighbor interchanges
- Subtree pruning and regrafting

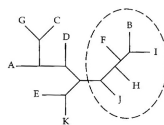
NEIGHBORHOODS FOR THE TREE SPACE

Subtree pruning and regrafting

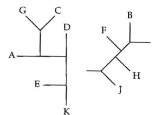
Nearest-neighbor interchanges



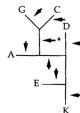
42 Chapter 4



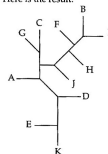
Break a branch, remove a subtree



Add it in, attaching it to one (*) of the other branches

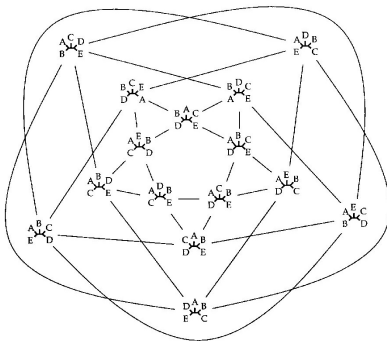
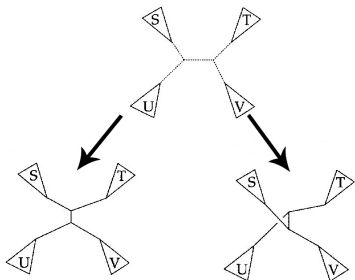


Here is the result:



NEIGHBORHOODS FOR THE TREE SPACE

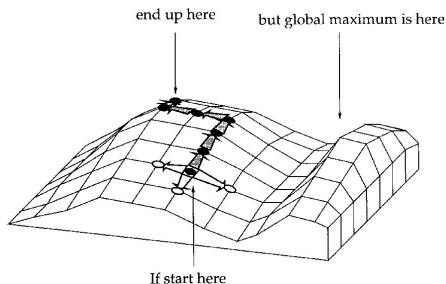
Nearest-neighbor interchanges



HEURISTIC ALGORITHMS

TRAJECTORY METHODS

- Local Optimization
- Greedy Search
- Simulated Annealing
- GRASP
- Taboo Search



HEURISTIC ALGORITHMS

TRAJECTORY METHODS

- Local Optimization
- Greedy Search
- Simulated Annealing
- GRASP
- Taboo Search

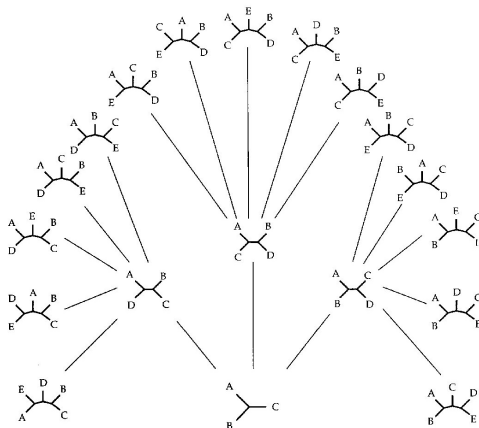
POPULATION METHODS

- Genetic Algorithms
- Ant Colony Optimization
- Particle swarm optimization

BRANCH AND BOUND SEARCH

- Branch and Bound is a technique that explores a **portion of the search space** that is **guaranteed to contain the optimal solution**. It is an **intelligent exhaustive search**.
- Trees are constructed **incrementally**, adding one species at time. If the current tree is known to lead to a worse solution to the best one found so far, the algorithm **backtracks** and reconsider previous choices.
- The algorithm needs to estimate **upper and lower bounds** on the function to be optimized from partial solutions.
- Its worst complexity is exponential, but it works well in practice.

BRANCH AND BOUND - SEARCH TREE



Branch and Bound can be seen as a visit of a **search tree**, **pruning** some subtrees (according to **bounds**) and **backtracking** before reaching leaves.

OUTLINE

- 1 INFERRING PHYLOGENIES
- 2 OPTIMIZATION PROBLEMS ON TREES
- 3 LEAST SQUARE METHODS
- 4 CLUSTERING METHODS

LEAST SQUARE METHOD

We have our **observed distance** matrix D_{ij} and a tree T with branch lengths predicting an **additive distance** matrix d_{ij} .

TARGET

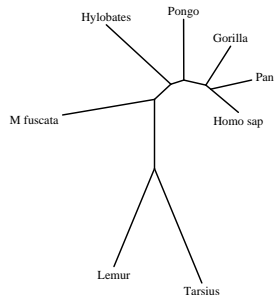
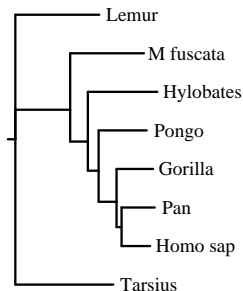
Find the tree T minimizing the error between d_{ij} and D_{ij} , i.e. the tree minimizing the weighted least square sum

$$S(T) = \sum_{i,j} w_{ij} (D_{ij} - d_{ij})^2$$

- Given a tree topology, the best branch lengths for S can be computed by solving a linear system.
- A least square algorithms needs to **search the tree space** for the best tree T : this is an **\mathcal{NP} -hard** problem.
- The search for the best tree can use branch and bound methods or heuristic state space explorations.
- This method gives the best explanation of the data

FITCH-MARGOLIAS ALGORITHM

- Letting $w_{ij} = \frac{1}{D_{ij}^2}$ in $S(T) = \sum_{i,j} w_{ij}(D_{ij} - d_{ij})^2$, we obtain the method of **Fitch-Margoliash**.
- The choice of $\frac{1}{D_{ij}^2}$ has statistical reasons: it takes into account the variance from the expected additive distances.



OUTLINE

- 1 INFERRING PHYLOGENIES
- 2 OPTIMIZATION PROBLEMS ON TREES
- 3 LEAST SQUARE METHODS
- 4 CLUSTERING METHODS

HEURISTIC METHODS: UPGMA

HIERARCHICAL CLUSTERING

- **Hierarchical clustering** works by iteratively merging the two closest clusters (sets of elements) in the current collection of clusters.
- It requires a matrix of distances among singletons.
- Different ways of computing **intercluster distances** give rise to different HC-algorithms.

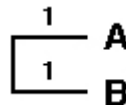
UPGMA

- UPGMA (Unweighted Pair Group Method with Arithmetic mean) computes the distance between two clusters as

$$d(A, B) = \frac{1}{|A||B|} \sum_{i \in A, j \in B} d_{ij}.$$
- When two clusters A and B are merged, their union is represented by their ancestor node in the tree.
- The distance between A and B is evenly split between the two branches entering in A and B

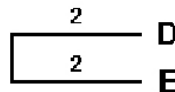
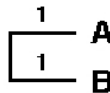
UPGMA - AN EXAMPLE

	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	2				
<i>C</i>	4	4			
<i>D</i>	6	6	6		
<i>E</i>	6	6	6	4	
<i>F</i>	8	8	8	8	8



$$d([A, B], C) = \frac{1}{2}(d(A, C) + d(B, C)) = 4$$

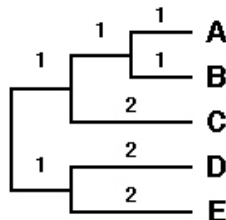
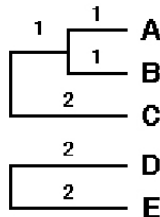
	<i>A, B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>C</i>	4			
<i>D</i>	6	6		
<i>E</i>	6	6	4	
<i>F</i>	8	8	8	8



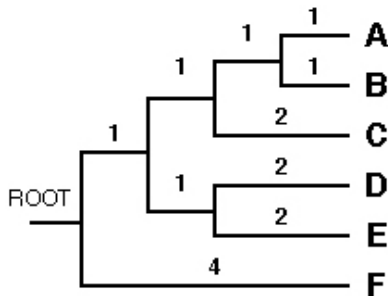
UPGMA - AN EXAMPLE

	<i>A, B</i>	<i>C</i>	<i>D, E</i>
<i>C</i>	4		
<i>D, E</i>	6	6	
<i>F</i>	8	8	8

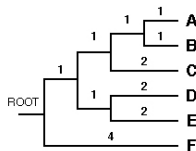
	<i>A, B, C</i>	<i>D, E</i>
<i>D, E</i>	6	
<i>F</i>	8	8



UPGMA - AN EXAMPLE

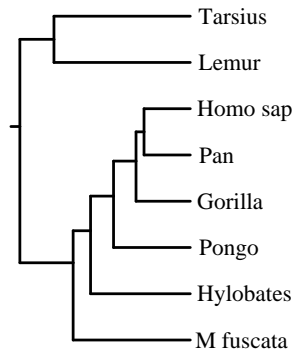


UPGMA - CONSIDERATIONS



HYPOTHESIS

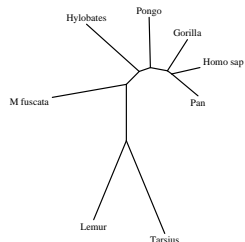
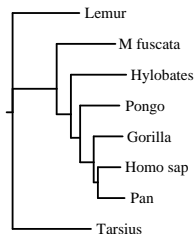
UPGMA reconstructs correctly the tree if the input distance is an **ultrametric** (**molecular clock**).



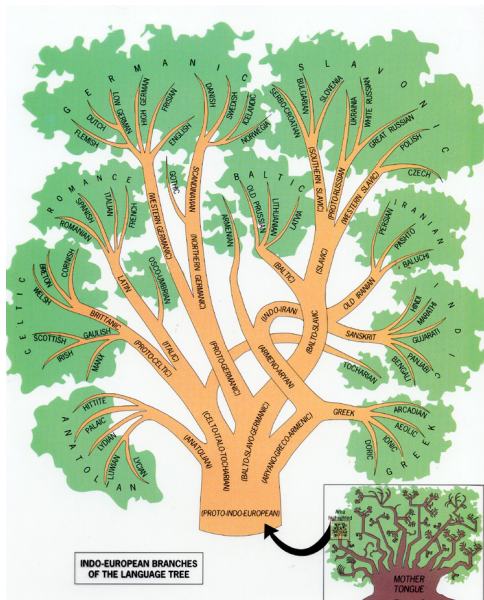
HEURISTIC METHODS: NEIGHBOR-JOINING

NEIGHBOR-JOINING

- Neighbor-Joining works similarly to UPGMA, but it merges together the two clusters minimizing $D_{ij} = d_{ij} - r_i - r_j$, where $r_i = \frac{1}{c-2} \sum_k d_{ik}$ is the average distance of i from all other nodes.
- When i and j are merged, their new ancestor x has distances from another node k equal to $d_{xk} = \frac{1}{2}(d_{ik} + d_{jk} - d_{ij})$
- The branch lengths are $d_{ix} = \frac{1}{2}(d_{ij} + r_i - r_j)$ and $d_{jx} = \frac{1}{2}(d_{ij} + r_j - r_i)$.
- NJ reconstructs the correct tree if the input distance is **additive**.



NOT ONLY DNA EVOLVE...



THE END

Thanks for the attention!

