

Maximum Likelihood Methods in Phylogenetics

ML trees and likelihood-based tree topology testing

Heiko A. Schmidt

Center for Integrative Bioinformatics Vienna (CIBIV)
 Max F. Perutz Laboratories (MFPL)
 Vienna, Austria
 heiko.schmidt@univie.ac.at

September 2007

Data	Method	Evaluation Criterion
Characters (Alignment)	Maximum Parsimony	Parsimony
	Statistical Approaches: Likelihood, Bayesian	Evolutionary Models
Distances	Distance Methods	

What is the Maximum Likelihood (ML) Approach?

Having the probabilistic process of evolution and its parameters, we could compute the probability of any outcoming sequence data.

$$p(\text{Data} \mid \text{Parameter set } \theta) \equiv L(\text{Parameter set } \theta \mid \text{Data})$$

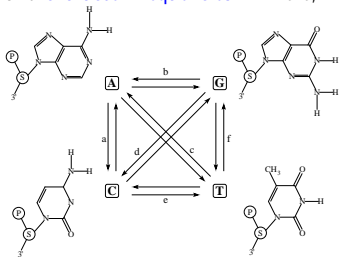
But here we are interested in the process and parameters themselves.

Hence, "likelihood flips the probability around."

Aim: The ML approach searches for that parameter set θ for the process (i.e., evolution) which maximizes the probability of our given dataset.

Substitution Models

Evolutionary models are often described using a substitution rate matrix R and character frequencies Π . Here, 4×4 matrix for DNA models:



$$R = \begin{pmatrix} & A & C & G & T \\ - & a & b & c \\ a & - & d & e \\ b & d & - & f \\ c & e & f & - \end{pmatrix}$$

$$\Pi = (\pi_A, \pi_C, \pi_G, \pi_T)$$

R and Π are then combined to a substitution probability matrix $P(t)$ which allows us to compute the probability $P_{ij}(t)$ of a change $i \rightarrow j$ over a time t .

Protein Models

Generally this is the same for protein sequences, but with 20×20 matrices. Some protein models are:

- Poisson model ("JC69" for proteins)
- Dayhoff (Dayhoff *et al.*, 1978)
- JTT (Jones *et al.*, 1992)
- mtREV (Adachi & Hasegawa, 1996)
- cpREV (Adachi *et al.*, 2000)
- VT (Müller & Vingron, 2000)
- WAG (Whelan & Goldman, 2000)
- BLOSUM-62 (Henikoff & Henikoff, 1992)

Problem: parameter sets

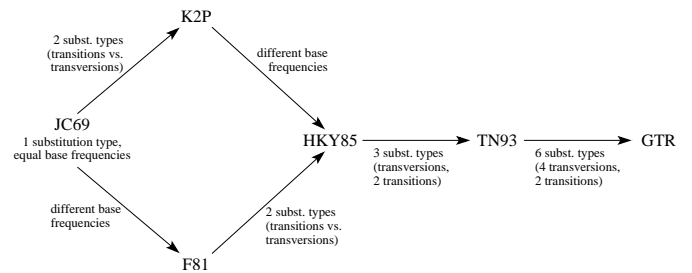
Problem: In phylogenetic analysis, the parameter set θ comprises:

- evolutionary model
- its parameters
- tree topology
- its branch lengths

This makes ML a high dimensional optimization problem that usually cannot be solved in one go.

Hence, some parameters like the substitution model and the model parameters are often determined/set separately from the tree.

DNA substitution models



There are further submodels (see Modeltest) and extensions like models assuming rate heterogeneity and codon models.

Computing ML Distances Using $P_{ij}(t)$

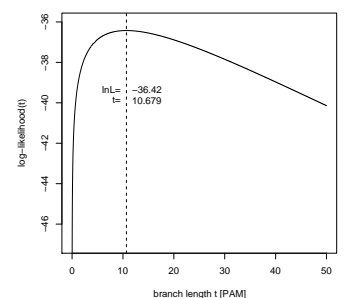
The Likelihood of sequence s evolving to s' in time t :

$$L(t|s \rightarrow s') = \prod_{i=1}^m (\Pi(s_i) \cdot P_{s_i s'_i}(t))$$

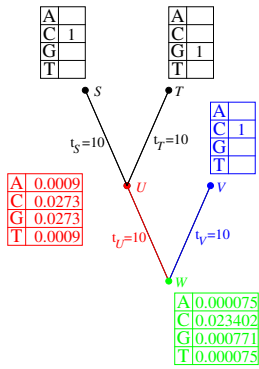
Likelihood surface for two sequences under JC69:

GATCCTGACAGAAATAAAC
 GTCCTGACAGAAATAAAC

Note: we do not compute the probability of the distance t but that of the data $D = \{s, s'\}$.



Likelihoods of Trees (Single column $\frac{C}{C}$, given tree)



Likelihoods of nucleotides at inner nodes:
 $L_U(i) = [P_{iC}(10) \cdot L(C)] \cdot [P_{iG}(10) \cdot L(G)]$

$$L_W(i) = \left[\sum_{u \in \{ACGT\}} P_{iu}(t_U) \cdot L_U(u) \right] \cdot \left[\sum_{v \in \{ACGT\}} P_{iv}(t_V) \cdot L_V(v) \right]$$

$$L^k = \sum_{i \in \{ACGT\}} \pi_i \cdot L_W(i) = 0.024323$$

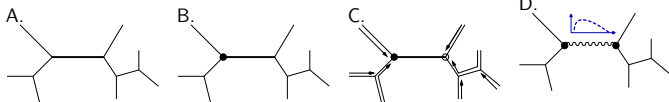
Site-Likelihood of an alignment column k :

$$L^k = \sum_{i \in \{ACGT\}} \pi_i \cdot L_W(i) = 0.024323$$

Adjusting Branch Lengths Step-By-Step

To compute optimal branch lengths do the following. Initialize the branch lengths.

Choose a branch (A). Move the virtual root to an adjacent node (B). Compute all partial likelihoods recursively (C). Adjust the branch length to maximize the likelihood value (D).



Repeat this for every branch until no better likelihood is gained.

Finding the ML Tree

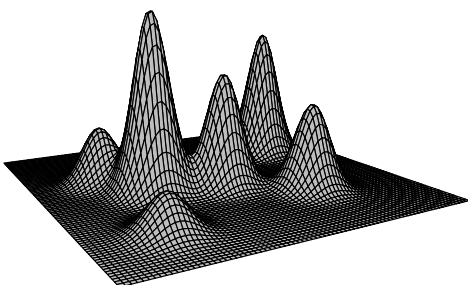
Exhaustive Search: guarantees to find the optimal tree, because all trees are evaluated, but not feasible for more than 10-12 taxa.

Branch and Bound: guarantees to find the optimal tree, without searching certain parts of the tree space – can run on more sequences, but often not for current-day datasets.

Heuristics: cannot guarantee to find the optimal tree, but are at least able to analyze large datasets.

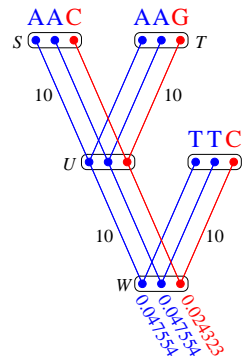
Local Maxima

What if we have **multiple maxima** in the likelihood surface?



Tree rearrangements to escape local maxima.

Likelihoods of Trees (multiple columns)



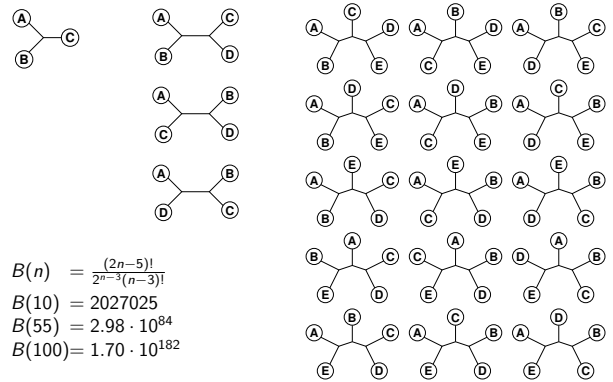
Considering this tree with $n = 3$ sequences of length $m = 3$ the tree likelihood of this tree is

$$\mathcal{L}(T) = \prod_{k=1}^m L^k = 0.047554^2 \cdot 0.024323 = 0.000055$$

or the log-likelihood

$$\ln \mathcal{L}(T) = \sum_{k=1}^m \ln L^k = -9.80811$$

Number of Trees to Examine...



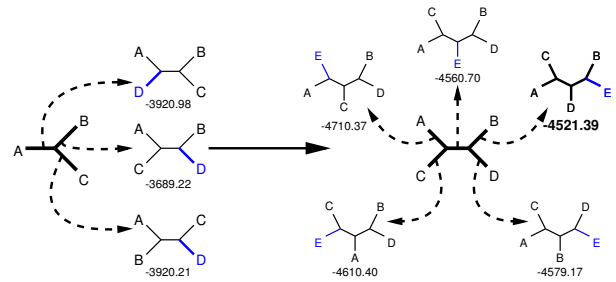
$$B(n) = \frac{(2n-5)!}{2^{n-3}(n-3)!}$$

$$B(10) = 2027025$$

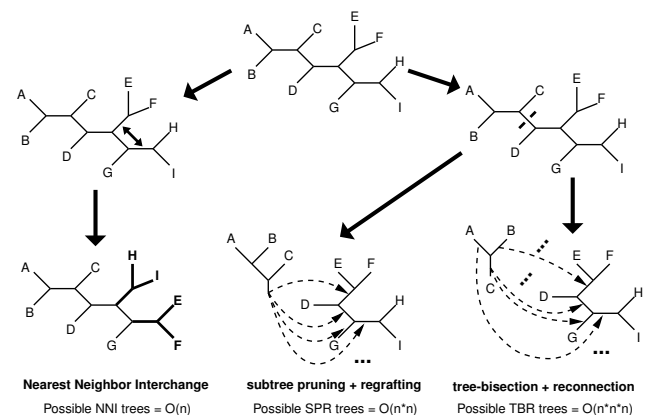
$$B(55) = 2.98 \cdot 10^{84}$$

$$B(100) = 1.70 \cdot 10^{182}$$

Build up a tree: Stepwise Insertion



Tree Rearrangements: Scanning a Tree's Neighborhood



Nearest Neighbor Interchange
Possible NNI trees = $O(n)$

subtree pruning + regrafting
Possible SPR trees = $O(n^2)$

tree-bisection + reconnection
Possible TBR trees = $O(n^2 \cdot n^2)$

Concept: Stepwise insertion + NNI/SPR

- 1 Build tree with **stepwise insertion**
 - (a) after each insertion optimize using NNI/local rearrangement (default, but user-adjustable gradually up to SPR; only fastDNAML)
 - (b) repeat (a) rearrangements until no better tree found.
- 2 after the last insertion optimize using SPR/global rearrangement (in DNAML; in fastDNAML user-adjustable gradually down to NNI)
- 3 repeat (2) rearrangements until no better tree found.

Pro: Evaluating large neighborhood with SPR.

Con: Slow.

Note: To save time, in other methods steps (1) and (2) are usually substituted by swiftly computed trees (e.g., BioNJ).

ML programs: PHYML

Concept: BioNJ tree + fastNNI

- 1 Start with BioNJ tree.
- 2 Do fastNNIs to optimize trees, i.e., evaluate all NNIs simultaneously and then merge all best ones which are non-conflicting.
- 3 Repeat (1) until no better tree found anymore.

Pro: Fast

Con: Prone to get stuck on local optima due to NNI-only. (SPR-based version PhyML-SPR has not been released yet.)

ML programs: Other strategies

- Genetic Algorithms (GARLI, GAML, MetaPIGA)
- Simulated Annealing (SSA, RAxML-SA)
- Quartet-based trees (TREE-PUZZLE, Qstar)
- ...

Note: The first two are also based on NNI/SPR/TBR.

Branch Support

- We can now reconstruct ML trees, but how comparable are the likelihoods, how reliable the groupings?
- Branch reliability can be checked, support values computed using:
 - Randomizing input orders in stepwise insertions (e.g., TREE-PUZZLE).
 - Jackknifing alignment columns + consensus.
 - Bootstrapping alignment columns + consensus.
 - Trees from Bayesian MCMC sampling + consensus.

Concept: MP tree + LSR

Descendant on fastDNAML, but ...

- 1 Starting with MP tree.
- 2 Uses *lazy subtree rearrangements* (only the 3 insertion branches are optimized), collecting candidates.
- 3 Candidates are evaluated.
- 4 Iterating (2)-(4).

Pro: Fast,

smart algorithmic and numerical optimized ML computation.

Con: Only few trees fully evaluated trees.

ML programs: IQPNNI

Concept: BioNJ tree + randomization + fastNNI

- 1 Start with BioNJ tree.
- 2 Do fastNNIs to optimize trees, i.e., evaluate all NNIs simultaneously and then accept all best ones which are non-conflicting. (after first round, identical to PHYML).
- 3 Remove randomly a certain amount of taxa and re-insert them by a fast and rough quartet-based method. (some randomization)
- 4 Repeat (2)-(3) until stop criterion is met.

Pro: Can evade local optima, offers automatic stopping criterion, hints when search didn't run enough, numerically optimized ML computation, offers codon models

Con: slower than PhyML/RAxML

How reliable is the reconstructed tree:

- Usually programs deliver a single tree, but without confidence values for the subtrees.
- How can we assess reliability for the subtree?

Are two evolutionary trees/models different?

Given sequence alignments and substitution models, we can reconstruct tree and compute their likelihoods.

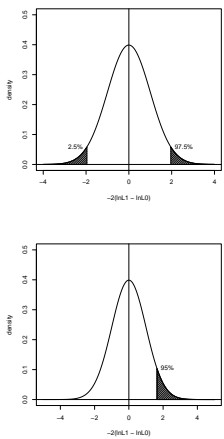
But can we decide from the likelihood

- which is the best substitution model to use?
 - yes, using (hierarchical) LRTs
- which tree is better?
 - only if likelihoods are computed from identical datasets/taxa.
- whether two tree likelihoods are significantly different?
 - yes, for identical datasets/taxa
- whether one tree likelihood is significantly better/worse?
 - yes, for identical datasets/taxa

These questions can be assessed by hypothesis testing.

- What question do I want to answer?
 - Say should I use the JC model or the GTR model?
 - Or perhaps, is tree T_a statistically different from tree T_b ?
- It is important to note that you should know the null hypothesis/hypotheses **before** you "collect" the data.

Usual Null-Hypotheses:



First the Null hypothesis has to be stated, for example:

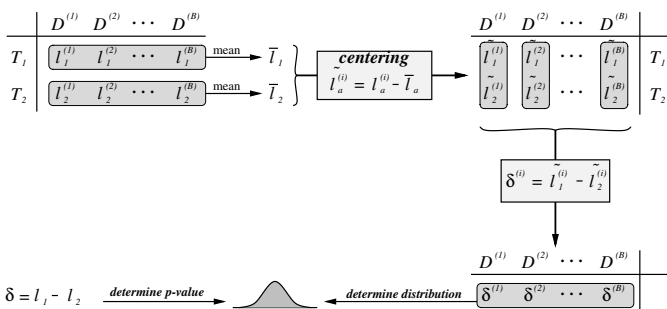
- **top:** The two likelihood are not significantly different – i.e. their expected difference $E(\ln L_1 - \ln L_0) = 0$.
- **bottom:** The 2nd likelihood is not significantly worse – i.e. their expected difference $E(\ln L_1 - \ln L_0) \leq 0$.

If the observed value falls into the white area, the Null hypothesis cannot be rejected. If it falls into the grey area, this is interpreted as support for the alternative by rejecting the Null hypothesis.

Time Saving: REL

- The re-optimization to get the log-likelihood values $L_x^{(i)}$ is very time consuming.
- Hence, often the site-likelihoods are used fixed.
- During the bootstrap the already estimated site-log-likelihoods are sampled and added to produce $L_x^{(i)}$.
- The resampling of estimated log-likelihoods (RELL) has been shown to be often sufficient to produce the distribution of log-likelihood differences.

Kishino-Hasegawa test:



- **Only nested models** can be tested by ordinary LRT: One model (H_0 , Null-model, constraint model) is nested in another model (H_A , alternative, unconstraint model) if the model H_0 can be produced by restricting parameters in model H_A .
- two different topologies are not nested.
- Thus, LRT cannot be used on different topologies, because the assumption of the χ^2 distribution does not fit.
- Hence, other (bootstrap-bases) methods have been devised to determine the distribution of log-likelihood differences for testing (e.g., KH or SH test).

Basic Idea:

- Compute log-likelihood values L_1, \dots, L_N for your trees T_1, \dots, T_N .
- Draw bootstrap samples i from the alignment, re-estimate the log-likelihood values $L_x^{(i)}$ for each tree T_x and for each sample i .
- Adjust the log-likelihoods with the mean by setting $\tilde{L}_x^{(i)} = L_x^{(i)} - \bar{L}_x^{(i)}$ (Centering) **Centering is needed to correct make the bootstrap samples conforming to the Null model.**
- Use the differences between the $\tilde{L}_x^{(i)}$ to determine the distribution of differences $\delta^{(i)} = \tilde{L}_y^{(i)} - \tilde{L}_z^{(i)}$.
- Use the distribution of $\delta^{(i)}$ to test your trees.

Original Kishino and Hasegawa test (KH test)

- This test was devised to test whether two *a priori* chosen trees (e.g., from a Markov Chain) are equally well supported by the dataset.
- H_0 : the expected $\delta = L_1 - L_2 = 0$.
- H_A : the expected $\delta = L_1 - L_2 \neq 0$.
- KH assumes that the ML tree is not among the trees.

Mis-use of the Kishino and Hasegawa test (KH test)

- Often, instead two *a priori* chosen trees, one tree is tested against the ML tree T_{ML} .
- That means, $\delta = L_{ML} - L_1$ is rarely negative.
- Hence, δ has to be tested in a single-sided regime.
- H_0 : the expected $\delta = L_1 - L_2 = 0$. H_A : the expected $\delta = L_1 - L_2 > 0$.

Multiple trees (Shimodaira and Hasegawa test - SH test)

- The SH test offers a correct way to test a set of trees, which may be chosen *a posteriori* after ML analysis.
- H_0 : All trees including T_{ML} are equally supported. H_A : Some or all trees T_x are not equally well supported.
- The SH test assumes, that the ML tree T_{ML} is among the trees.

Heiko A. Schmidt ML trees and topology testing

The SH procedure:

- Compute log-likelihood values and the differences $\delta_x = L_{ML} - L_x$ for your trees.
- Draw bootstrap samples i from the alignment (with REML) to gain log-likelihood values $L_x^{(i)}$ for each tree T_x .
- Adjust the log-likelihoods with the mean over the samples i by setting $\tilde{L}_x^{(i)} = L_x^{(i)} - \bar{L}_x^{(i)}$ (Centering)
- For each sample i , find $\tilde{L}_{ML}^{(i)}$ over all topologies T_x .
- and compute $\delta_x^{(i)} = \tilde{L}_{ML}^{(i)} - \tilde{L}_x^{(i)}$.
- For each tree T_x , test whether δ_x is a plausible sample from the distribution of $\delta_x^{(i)}$ (over all replicates i).
- We use a single sided test, since $\tilde{L}_{ML}^{(i)} \geq \tilde{L}_x^{(i)}$.

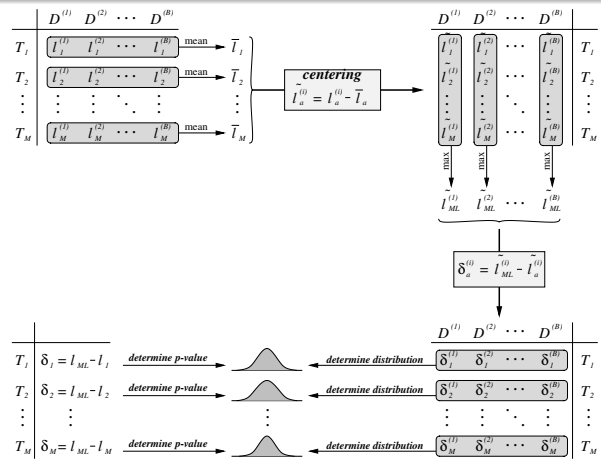
Heiko A. Schmidt ML trees and topology testing

Pros and Cons of various tests:

- Kishino-Hasegawa test (KH)** – usually mis-used if the trees are not chosen *a priori*.
- Shimodaira-Hasegawa (SH)** – test for multiple trees, affected by selection error, i.e., gets more conservative with the number of trees.
- Weighted Shimodaira-Hasegawa (wSH)** – SH test weighted with the variance of the likelihood difference (less conservative than SH).
- Expected likelihood weights (ELW)** – less conservative than SH, but the impact of model mis-specification unclear.
- Approximately unbiased test (AU)** – fixes the conservativeness issue of SH, but many similarly good trees can lead to artificial over-confidence.

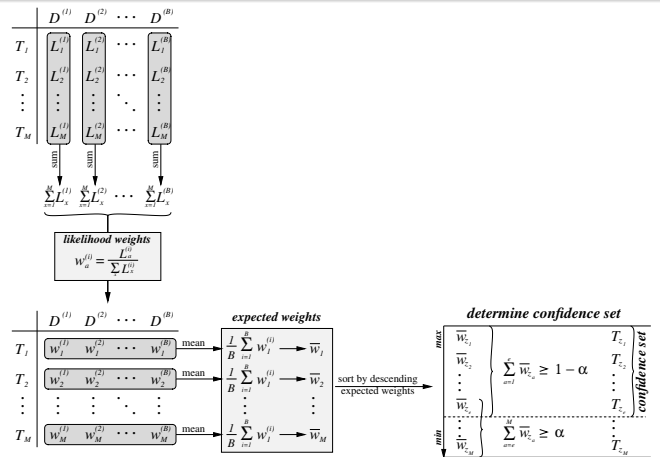
Heiko A. Schmidt ML trees and topology testing

Shimodaira-Hasegawa test:



Heiko A. Schmidt ML trees and topology testing

Expected Likelihood weights:



Heiko A. Schmidt ML trees and topology testing

Summary

- Likelihoods gives a strong statistical framework for hypothesis testing.
- Proper experimental design and proper use of tests is required.
- One should always be aware of the hypothesis a test assesses and should make sure that this answers the question asked.
- Testing tree topologies can be used to assess whether two competing hypotheses are really substantially different. If they are not, one cannot be preferred over the other.

Heiko A. Schmidt ML trees and topology testing