

Enclosing the Maximum Likelihood of the Simplest DNA Model Evolving on Fixed Topologies: Towards a Rigorous Framework for Phylogenetic Inference

Raazesh Sainudiin

Cornell University, Ithaca, USA

Summary. An interval extension of the recursive formulation for the likelihood function of the simplest Markov model of DNA evolution on unrooted phylogenetic trees with a fixed topology is used to obtain rigorous enclosure(s) of the global maximum likelihood. Validated global maximizer(s) inside any compact set of the parameter space which is the set of all branch lengths of the tree are thus obtained. The algorithm is an adaptation of a widely applied global optimization method using interval analysis for the phylogenetic context. The method is applied to enclose the maximizer(s) and the global maximum for the simplest DNA model evolving on trees with 2, 3, and 4 taxa. The method is also applicable to a wide class of inclusion isotonic likelihood functions.

1. INTRODUCTION

When one is given a homologous set of distinct deoxyribo-nucleic acid (DNA) sequences of length v from s species and asked for an estimate of their inter-relationships back through time under some model of DNA evolution, a phylogenetic tree estimation problem arises. This problem is two-fold. First, one has to estimate the shape or topology of the tree, which captures the set of “who is related to whom and in what order? and whose ancestors are related to whose and in what order?” questions. Second, one has to estimate the lengths of the branches when given a particular topology. The branch lengths of a tree usually represent a scaled product of mutation rate and number of generations between the nodes. The s extant species are represented by the external nodes or leaves and their ancestors are represented by the internal nodes of the tree. A rooted tree always has a bifurcation at the root, typically the most recent common ancestor of all s leaves, where as, an unrooted tree has m -furcations at all internal nodes with $m \geq 3$. This work focusses on the second problem, namely, estimating the branch lengths for a given topology in a maximum likelihood framework.

When statistical inference is conducted in a maximum likelihood framework, one is interested in the global maximum of the likelihood function over the parameter space. Explicit analytical solutions for the maximum likelihood estimates of the branch lengths for a specified unrooted topology with more than 2 leaves are not available even for the simplest model of DNA evolution due to Jukes and Cantor (1969) without assuming a molecular clock. See 5.(b) of Yang (2000) for results on clocked 3-leaved rooted trees. Results are known for models with two character states superimposed on 3-leaved trees (Yang, 2000), as well as for specific observations on 4-leaved trees (Chor, 2000). In practice one settles for a local optimization algorithm to numerically approximate the global solution.

However, statistical inference procedures that rely on having found some global optimum through any numerical approach may suffer from at least five major sources of errors. To fully appreciate the sources of errors one needs some understanding of a number screen. Computers only support a finite set of numbers that are usually represented in a semi-logarithmic manner as a set of fixed length floating-point numbers of the form

$$x = \pm m \cdot b^e = \pm 0.m_1 m_2 \cdots m_p \cdot b^e$$

where, m is the signed mantissa of precision p , b is the base (usually 2) and e , bounded between \underline{e} and \bar{e} , is the exponent. When $b = 2$, the digits of the mantissa $m_1 = 1$ and $m_i \in \{0, 1\}$, $\forall i, 1 < i \leq p$, and the smallest and largest machine-representable numbers in absolute value are $\underline{x} = 0.10 \cdots 0 \cdot 2^{\underline{e}}$ and $\bar{x} = 0.11 \cdots 1 \cdot 2^{\bar{e}}$, respectively. For a detailed description of the binary floating-point system see P754 (1985). Thus, the binary floating-point system $\mathcal{R} = \mathcal{R}(2, p, \underline{e}, \bar{e})$ is said to form a screen of the real numbers in the interval $[-\bar{x}, +\bar{x}]$ with 0 uniquely represented by $0.00 \cdots 0 \cdot 2^{\underline{e}}$.

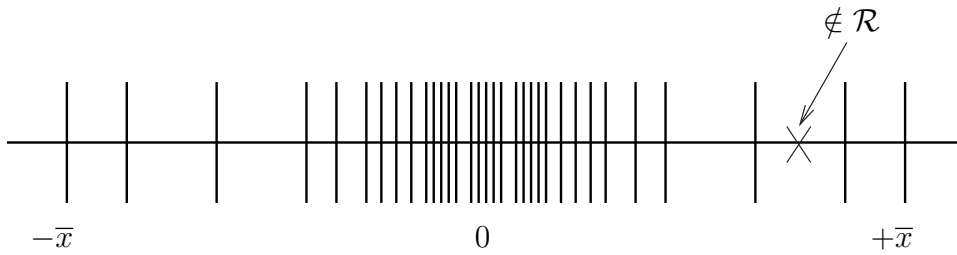


Fig. 1. A Number Screen \mathcal{R} for the interval $[-\bar{x}, +\bar{x}] \subset \mathbb{R}$

Arithmetic on a machine is performed with such a screen and thus cannot be exact. The computed result of an arithmetic operation is only an approximation to the true result, and this difference between the computed and the actual result is known as the *roundoff error*. Such errors can accumulate catastrophically in a long sequence of imprecise operations (Cuyt et al., 2001; Loh and Walster, 2002). The next major source of error comes from truncating various mathematical expressions, such as, derivatives, integrals, and transcendental functions, that are only defined in terms of the limit of some infinite sequence of operations. Since, only finitely many operations can be performed on a machine, the actual limit is only approximated, and the difference between the actual limit and the computed approximation to it is known as the *truncation error*. Another source of error arises in the *conversion* of constants represented in the decimal format to a binary floating-point format used in a machine. For example, the real number 0.1 has $0.000\overline{1100}_2$ as its unique, but infinite, binary representation. Thus, 0.1 does not have an exact binary representation in any binary floating-point system with a fixed length mantissa, and the truncation of the mantissa to any finite length in order to keep 0.1 from slipping through the screen \mathcal{R} results in a *conversion error*. Errors in statistical decision may also result from an *ill-posed statistical experiment or model*. For instance, when one has not proved that the model is identifiable, there may exist unknown nonidentifiable parameter subspaces that need to be rectified. Finally, no experimental procedure is immune to the physical limits on empirical resolution and therefore *measurement error* in the data cannot be ignored, especially for highly nonlinear models.

Furthermore, traditional nonlinear programming techniques that use local information, such as, clustering methods, generalized descent methods, and other stochastic search methods, including simulated annealing, start from some approximate trial point(s) and iterate by sampling only finitely many points. Therefore, they can neither validate that the objective function has not plunged between the sampled points, nor guarantee escape from a local minimum, albeit they can be made to increase the probability of such desired events. Figure 2 shows the trajectories (shaded circles) of two local searches with random initial conditions (white circles) that are attracted to some fixed point (black circle) in a local valley on two different functions. The first function has a sharp valley that is not visited by the local search trajectories. The second function is highly multi-modal with different basins of attraction for each valley. In both cases, the global minimum is missed by the search. Methods that use local information at finitely many points and do not account for all five major sources of errors, cannot be expected to yield anything more than a possible, approximate, and local solution.

Controversy exists over the nonrigorous nature of such a numerically-based statistical inference procedure, especially in parameter-rich models that have not been shown to be identifiable and/or are adorned with multiple local optima. In some problems, a local approximate solution may be sufficient, but in others one may base statistical decisions that address a real biological problem on an incorrect solution. Unfortunately, it currently seems impossible to know this difference *a priori*. This paper shows an existing method toward such knowledge and applies it to enclose the maximum likelihood value as well as the estimate of the most likely unrooted multifurcating four taxa tree for any given data set of four homologous DNA sequences that are assumed to evolve according to the Jukes and Cantor model (Jukes and Cantor, 1969). The global optimization method sketched below rigorously encloses the global maximum of the likelihood function through interval analysis. Such interval methods evaluate the likelihood function over a continuum of points including those that are not machine-representable and account for the five sources of errors described earlier. Thus, in contrast to local search methods, interval methods can enclose the global optimum with guaranteed

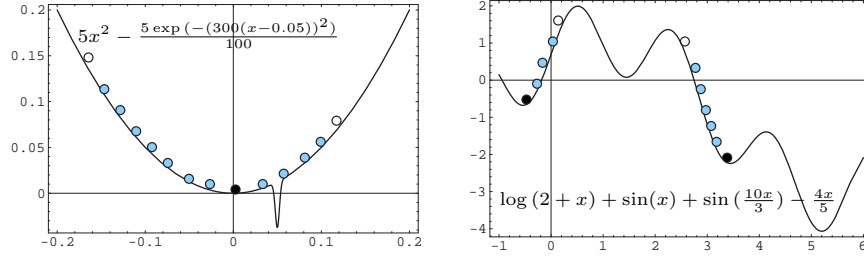


Fig. 2. Two randomly initialized local searches on a sharp-valleyed and a multi-modal function.

accuracy by exhaustive search within any compact set of the parameter space.

Section 2 contains a brief introduction to various enclosure arithmetics. For a recent introduction to such arithmetics see (Kulisch et al., 2001). Section 3 describes a problem that arises in phylogenetic inference through maximum likelihood. Section 4 gives the basic global optimization algorithm based on Hansen's method (Hansen, 1980, 1992) with Ratz's modifications (Ratz, 1992) as implemented in Hammer et al. (1995) with further extensions that account for non-stationary maxima at the boundaries and increase computational efficiency. Assuming the simplest model of DNA evolution, the method is applied in section 5 to rigorously estimate trees, with two, three, and four leaves, based on primate mitochondrial DNA sequences (Brown et al., 1982).

2. PRELIMINARIES

2.1. Interval arithmetic

Lower case letters denote real numbers, e.g. $x \in \mathbb{R}$, the set of real numbers. Upper case letters represent bounded and closed (compact) real intervals, e.g. $X = [\underline{x}, \bar{x}] = [\inf(X), \sup(X)]$. Any compact interval $X \in \mathbb{IR} := \{[a, b] : a \leq b, a, b \in \mathbb{R}\}$, the set of all compact real intervals. The diameter and the midpoint of X are $d(X) := \bar{x} - \underline{x}$ and $m(X) := (\underline{x} + \bar{x})/2$, respectively. The smallest and largest absolute value of an interval X are real numbers given by $\langle X \rangle := \min\{|x| : x \in X\} = \min\{|\underline{x}|, |\bar{x}|\}$ and $|X| := \max\{|x| : x \in X\} = \max\{|\underline{x}|, |\bar{x}|\}$, respectively, while the absolute value of an interval X is $|X|_{[\cdot]} := \{|x| : x \in X\} = [\langle X \rangle, |X|]$. The relative diameter of an interval X , denoted by d_{rel} is the diameter $d(X)$ itself if $0 \in X$, and $d(X)/\langle X \rangle$, otherwise. An interval X with zero diameter is called a thin interval with $\underline{x} = \bar{x} = x$. The hull of two intervals is $X \sqcup Y := [\min\{\underline{x}, \underline{y}\}, \max\{\bar{x}, \bar{y}\}]$. By the notation $X \subseteq Y$, it is meant that X is completely contained in Y , i.e., $\underline{x} > \underline{y}$ and $\bar{x} < \bar{y}$. No notational distinction is made between a real number $x \in \mathbb{R}$ and a real vector $x = (x_1, \dots, x_n)^T \in \mathbb{R}^n$ and between a real interval X and a real interval vector or box $X = (X_1, \dots, X_n)^T \in \mathbb{IR}^n$, i.e. $X_i = [\underline{x}_i, \bar{x}_i] = [\inf(X_i), \sup(X_i)] \in \mathbb{IR}$, where, $i = 1, \dots, n$. The dimension n should be clear from the context. For an interval vector X , the diameter, relative diameter, midpoint, and hull operations are defined component-wise to yield vectors, while the maximum over its components is taken to obtain the maximal diameter and the maximal relative diameter, $d_\infty(X) = \max_i d(X_i)$ and $d_{rel, \infty}(X) = \max_i d_{rel}(X_i)$, respectively. It can be seen that \mathbb{IR} under the metric \mathfrak{h} , given by,

$$\mathfrak{h}(X, Y) := \max\{|\underline{x} - \underline{y}|, |\bar{x} - \bar{y}|\},$$

is a complete metric space. Convergence of a sequence of intervals $\{X^{(i)}\}$ to an interval X under the metric \mathfrak{h} is equivalent to the sequence $\mathfrak{h}(X^{(i)}, X)$ approaching 0 as i approaches ∞ , which in turn is equivalent to both $\underline{x}^{(i)} \rightarrow \underline{x}$ and $\bar{x}^{(i)} \rightarrow \bar{x}$. Continuity and differentiability of a function $f : \mathbb{IR}^n \rightarrow \mathbb{IR}^k$ are defined in the usual way.

A real arithmetic operation $x \circ y$, where $\circ \in \{+, -, \cdot, /\}$ and $x, y \in \mathbb{R}$, is a continuous function $x \circ y := \circ(x, y) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, except when $y = 0$ under $/$ operation. An interval arithmetic operation $X \circ Y := \{x \circ y : x \in X, y \in Y\}$ thus yields the set that contains the result of the operation done for every real pair $(x, y) \in (X, Y)$. This definition of interval operation leads to the property of inclusion isotony which stipulates that $X \circ Y$ contain $V \circ W$ provided $V \subseteq X$ and $W \subseteq Y$. Since X and Y are simply connected compact intervals, so is their product $X \times Y$. On such a domain $X \times Y$, the continuity of $\circ(x, y)$ (except when $\circ = /$ and $0 \in Y$) ensures the attainment of a minimum, a maximum and all intermediate values. In

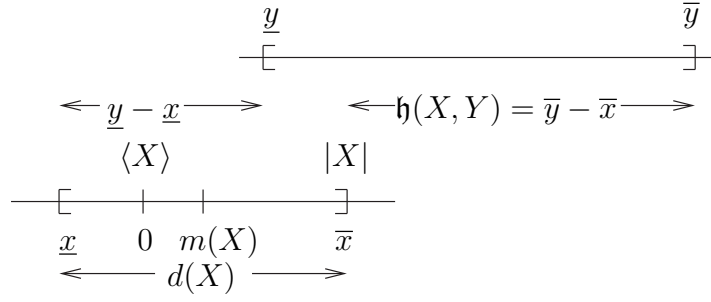


Fig. 3. Features of intervals

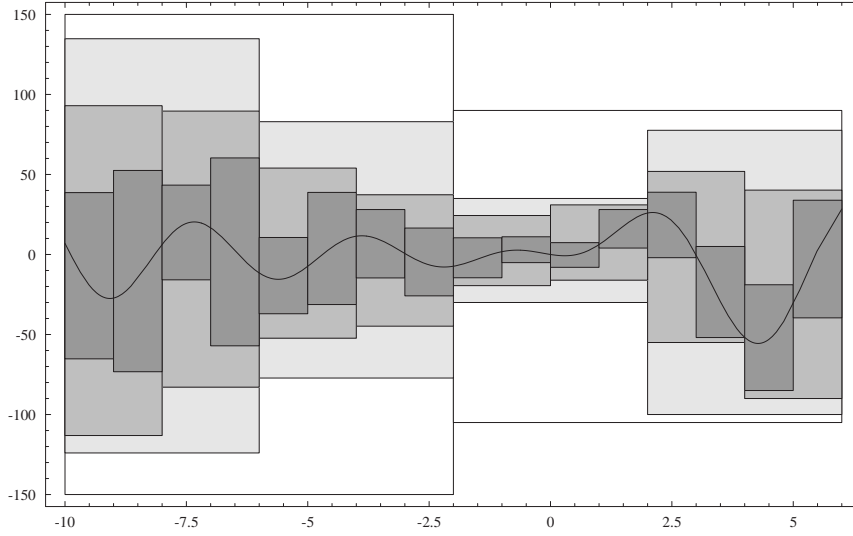


Fig. 4. Range enclosure of the natural interval extension of $-\sum_{k=1}^5 kx \sin(\frac{k(x-3)}{3})$ linearly tightens with the mesh

other words, with the exception of the case when $\circ = /$ and $0 \in Y$, the range $X \circ Y$ has an interval form $[\min(x \circ y), \max(x \circ y)]$, where the min and max are taken over all pairs $(x, y) \in X \times Y$. The particular forms of $X \circ Y$ for the elementary operations are,

$$\begin{aligned} X + Y &= [\underline{x} + \underline{y}, \bar{x} + \bar{y}], & X \cdot Y &= [\min\{\underline{x}\underline{y}, \underline{x}\bar{y}, \bar{x}\underline{y}, \bar{x}\bar{y}\}, \max\{\underline{x}\underline{y}, \underline{x}\bar{y}, \bar{x}\underline{y}, \bar{x}\bar{y}\}], \\ X - Y &= [\underline{x} - \bar{y}, \bar{x} - \underline{y}], \text{ and} & X/Y &= X \cdot [1/\bar{y}, 1/\underline{y}], \quad 0 \notin Y. \end{aligned}$$

The identity elements of $+$ and \cdot are the thin intervals $[0, 0]$ and $[1, 1]$, respectively. Multiplicative and additive inverses do not exist except when X is also thin, since $[0, 0] \subseteq X - X$, and $[1, 1] \subseteq X/X$. Although the commutative and associative laws are satisfied by $+$ and \cdot , only a weaker notion of distributivity called subdistributivity, i.e., $X \cdot (Y + Z) \subseteq (X \cdot Y) + (X \cdot Z)$, is satisfied.

For any real function $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ and some box $X \in \mathbb{IR}^n$, let the range of f over X be denoted by $f(X) := \{f(x) : x \in X\}$. Inclusion isotony also holds for interval evaluations that are compositions of arithmetic expressions and the elementary functions. When real constants, variables, and operations in f are replaced by their interval counterparts one obtains $F(X) : \mathbb{R} \rightarrow \mathbb{R}$, the natural interval extension of f . Guaranteed enclosures of the range $f(X)$ are obtained by $F(X)$, since inclusion isotony holds for F , i.e, if $X \subseteq Y$, then $F(X) \subseteq F(Y)$, and in particular, the inclusion property that $x \in X \implies f(x) \in F(X)$ holds. The natural interval extension $F(X)$ often overestimates $f(X)$, but can be shown under mild conditions to linearly approach the range as the maximal diameter of the box X goes to zero, i.e., $\mathfrak{h}(F(X), f(X)) \leq \alpha \cdot d_\infty(X)$ for some $\alpha \geq 0$. This implies that a partition of X into smaller boxes $\{X^{(1)}, \dots, X^{(m)}\}$ gives better enclosures of $f(X)$ through the union $\bigcup_{i=1}^m F(X^{(i)})$ as illustrated in Figure 4.

Let $\nabla F(x)$ and $\nabla^2 F(x)$ represent the interval extensions of $\nabla f(x)$ and $\nabla^2 f(x)$, the gradient and Hessian

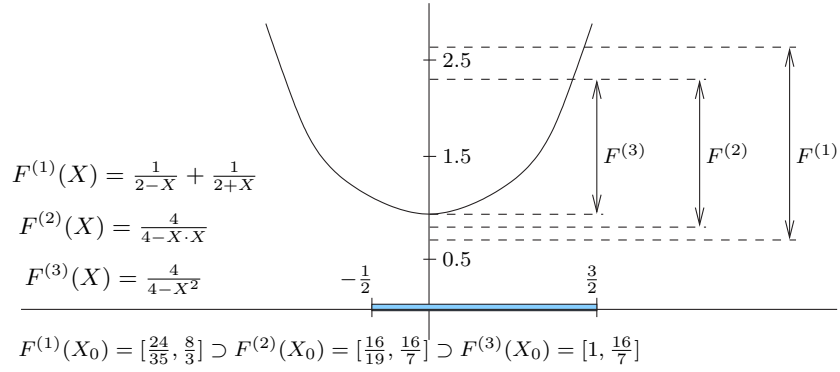


Fig. 5. Extension-specific dependence of range enclosures

of f . A better enclosure of $f(X)$ is possible for an f with the centered form,

$$f(x) = f(c) + \nabla f(b) \cdot (x - c) \in f(c) + \nabla f(X) \cdot (x - c) \subseteq F_c(X) := f(c) + \nabla F(X) \cdot (X - c),$$

where, $b, c, x \in X$ with b between c and x . $F_c(X)$ is the interval extension of the centered form of f with center c and decays quadratically to $f(X)$ as the maximal diameter of $X \rightarrow 0$. Finally, some interval extensions of f are better at enclosing the true range than others. Although the three functions shown in Figure 5 are equivalent, their interval extensions yield different range enclosures. The interval extension $F^{(3)}$ is better than $F^{(1)}$ and $F^{(2)}$ as depicted in Figure 5. Note that $F^{(3)} \subset F^{(2)}$ since $X^2 \subset X \cdot X$ in interval arithmetic. If X appears only once in the expression and all parameters are thin intervals, then it was shown by Moore (1979) that the natural interval extension does indeed yield a tight enclosure. In general, one can obtain tighter enclosures by minimizing the occurrence of X in the expression.

2.2. Differentiation Arithmetic

When it becomes too cumbersome or impossible to explicitly compute the derivative of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, or when f itself is only available as an algorithm, one may employ a differentiation arithmetic, often known as automatic differentiation (see for e.g. Griewank and Corliss (1991)) to obtain any $\nabla^k f$, the k th-order derivative of f . This approach circumvents the computation of a formal expression for f by defining a differentiation arithmetic on the ordered k -tuples $(f(x), \nabla f(x), \nabla^2 f(x), \dots, \nabla^k f(x))$ (Berz, 1991). A brief sketch of such an arithmetic is given for the case when $k = 2$ as it will be used in section 4.

Consider a twice-continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with the gradient vector and Hessian matrix given by $\nabla f(x) := (\partial f(x)/\partial x_1, \dots, \partial f(x)/\partial x_n)^T \in \mathbb{R}^n$, and $\nabla^2 f(x) := ((\partial^2 f(x)/\partial x_i \partial x_j))_{i,j=\{1,\dots,n\}} \in \mathbb{R}^{n \times n}$, respectively. For every, $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, consider its corresponding ordered triple $(f(x), \nabla f(x), \nabla^2 f(x))$. The ordered triples corresponding to a constant function, $c(x) = c : \mathbb{R}^n \rightarrow \mathbb{R}$, and a component identifying function (or variable), $I_j(x) = x_j : \mathbb{R}^n \rightarrow \mathbb{R}$, are $(c, 0, 0)$ and $(x_j, e^{(j)}, 0)$, respectively, where, $e^{(j)}$ is the j -th unit vector and the 0's are additive identities in their appropriate spaces. To perform an elementary operation $\circ \in \{+, -, \cdot, /\}$ with a pair of such triples to obtain another, the rules of calculus apply as follows:

$$\begin{aligned}
& (h(x), \nabla h(x), \nabla^2 h(x)) \\
& := (f(x), \nabla f(x), \nabla^2 f(x)) \circ (g(x), \nabla g(x), \nabla^2 g(x)) \\
& = (f(x) \circ g(x), \nabla f(x) \circ \nabla g(x), \nabla^2 f(x) \circ \nabla^2 g(x)), & \text{if } \circ \in \{+, -\} \\
& = (f(x) \cdot g(x), f(x) \cdot \nabla g(x) + g(x) \cdot \nabla f(x), \\
& \quad g(x) \cdot \nabla^2 f(x) + \nabla f(x) \cdot \nabla g(x)^T + \nabla g(x) \cdot \nabla f(x)^T + f(x) \cdot \nabla^2 g(x)), & \text{if } \circ = \cdot \\
& = (f(x)/g(x), 1/g(x) \cdot \{\nabla f(x) - h(x) \cdot \nabla g(x)\}, \\
& \quad 1/g(x) \cdot \{\nabla^2 f(x) \cdot \nabla h(x) \cdot \nabla h(x)^T - \nabla g(x) \cdot \nabla h(x)^T - h(x) \cdot \nabla^2 h(x)\}), & \text{if } \circ = /, g(x) \neq 0
\end{aligned}$$

The arithmetic for composition of functions, such as, $h(x) = r(f(x)) : \mathbb{R} \rightarrow \mathbb{R}$, with the first and second derivative of r given by r' and r'' , on their corresponding triples, given by,

$$(h(x), \nabla h(x), \nabla^2 h(x)) = (r(f(x)), r'(f(x)) \cdot \nabla f(x), r''(f(x)) \cdot \nabla f(x) \cdot \nabla f(x)^T + r'(f(x)) \cdot \nabla^2 f(x)),$$

yields the triples for the elementary functions, $\exp(x)$ and $\ln(x)$, which are used to compute the likelihood in section 3.

For dyadic reasons, the differentiation arithmetic has been explained above only in terms of reals. By replacing the real x 's above by interval X 's and performing all operations in the real interval arithmetic with the interval extension F of f , as discussed in section 2.1, one can rigorously enclose the components of the interval triple $(F(X), \nabla F(X), \nabla^2 F(X))$ through interval differentiation arithmetic, such that, for every $x \in X \in \mathbb{IR}^n$, $f(x) \in F(X) \in \mathbb{IR}$, $\nabla f(x) \in \nabla F(X) \in \mathbb{IR}^n$, and $\nabla^2 f(x) \in \nabla^2 F(X) \in \mathbb{IR}^{n \times n}$.

2.3. Interval Newton method and its extension

Newton's method linearly approximates the differentiable real function $f(x)$ in the neighborhood of an initial value $x^{(0)}$ by the tangent,

$$t(x) = f(x^{(0)}) + f'(x^{(0)})(x - x^{(0)}),$$

where, $f'(x)$ is the first derivative of $f(x)$. This tangent equation can be used to solve for an approximation to the zero of $f(x)$, by means of the following discrete dynamical system known as Newton's method:

$$x^{(j+1)} = x^{(j)} - \frac{f(x^{(j)})}{f'(x^{(j)})}, \quad j = 0, 1, 2, \dots$$

Thus, the zero of the tangent to $f(x)$ at the current approximate value $x^{(j)}$ gives the next approximate value $x^{(j+1)}$. The following geometric interpretation shown in Figure 6 is useful. During each iteration of the Newton's method, a light beam is shone upon the domain from the point $(x^{(j)}, f(x^{(j)}))$ along the tangent to $f(x)$ at $x^{(j)}$. The intersection of this beam (white line in Figure 6) with the domain provides $x^{(j+1)}$, which is where the next iteration is resumed. If x^* is the only root of $f(x)$ in the search interval X containing all iterates $x^{(j)}$, $f(x)$ is twice continuously differentiable, and $x^{(0)}$ is sufficiently close to x^* , then it is well known that Newton's method converges quadratically fast to x^* . Otherwise, it may well diverge or oscillate.

The interval version of Newton's method computes an enclosure of the zero x^* of a continuously differentiable function $f(x)$ in the interval X through the following dynamical system in \mathbb{IR} :

$$X^{(j+1)} = \left(m(X^{(j)}) - \frac{f(m(X^{(j)}))}{F'(X^{(j)})} \right) \cap X^{(j)}, \quad j = 0, 1, 2, \dots$$

Here, $X^{(0)} = X$, $F'(X^{(j)})$ is the enclosure of $f'(x)$ over $X^{(j)}$, and $m(X^{(j)})$ is the mid-point of $X^{(j)}$. Provided, $0 \notin F'(X^{(0)})$ or equivalently a unique zero of f lies in $X^{(0)}$, the interval Newton method will never diverge. Moore (1967) derived the interval Newton method. Under natural conditions on f , the sequence of compact sets $X^{(0)} \supseteq X^{(1)} \supseteq X^{(2)} \dots$ can be shown to converge quadratically to x^* (Alefeld and Herzberger, 1983). One can derive the above dynamical system in \mathbb{IR} via the mean value theorem. Let $f(x)$ be continuously differentiable and $f'(x) \neq 0$ for all $x \in X$ such that x^* is the only zero of f in X . Then, by the mean value theorem, for every x , there exists a $c \in (x, x^*)$, such that, $f(x) - f(x^*) = f'(c)(x - x^*)$. Since, $f'(c) \neq 0$, by assumption, and $f(x^*) = 0$, it follows that:

$$x^* = x - \frac{f(x)}{f'(c)} \in x - \frac{f(x)}{F'(X)} =: N(X), \quad \forall x \in X$$

$N(X)$ is called the Newton operator and it contains x^* . Since our root of interest lies in X , $x^* \in N(X) \cap X$. Note that the above dynamical system in \mathbb{IR} is obtained by replacing x with $m(X)$ and X with $X^{(j)}$ in the previous expression. The interval Newton method can also be interpreted geometrically. At the j^{th} iteration, a set of light beams are shone from the point $(x^{(j)}, f(x^{(j)}))$ along the directions of all the tangents to $f(x)$ on the entire interval X . The intersection of these beams (gray floodlight of Figure 6) with the domain is $N(X^{(j)})$. The iteration is resumed with the new interval $X^{(j+1)} = N(X^{(j)}) \cap X^{(j)}$. Next we extend the interval Newton method in order to allow $F'(X)$ to contain 0.

By including two ideal points $+\infty$ and $-\infty$ to \mathbb{R} , it becomes possible to extend interval arithmetic to $\mathbb{IR}^* := \mathbb{IR} \cup \{(-\infty, \bar{x}] : \bar{x} \in \mathbb{R}\} \cup \{[\underline{x}, +\infty) : \underline{x} \in \mathbb{R}\} \cup (-\infty, +\infty)$, the set of intervals with end points in the complete lattice $\mathbb{R}^* := \mathbb{R} \cup \{+\infty\} \cup \{-\infty\}$, with respect to the ordering relation \leq . Since, division is the inverse operation of multiplication, obtaining any $x/y \in X/Y := \{x/y : x \in X, y \in Y\}$ is equivalent to

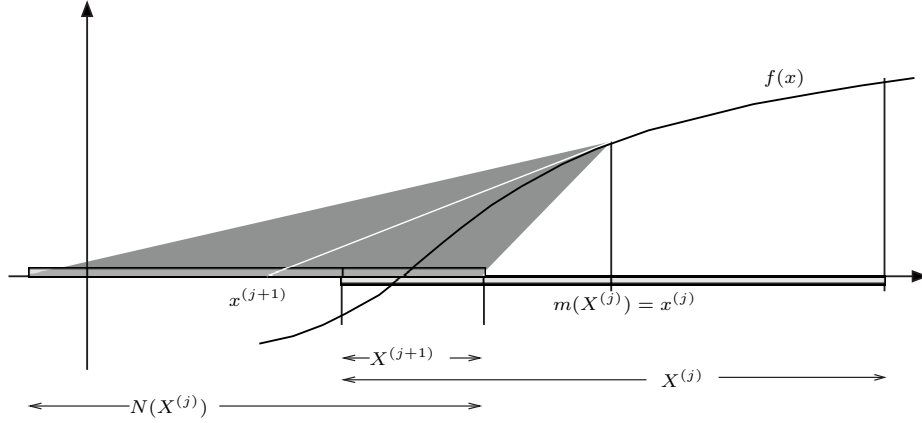


Fig. 6. Geometric interpretation of the interval Newton method

solving the equation $y \cdot s = x$ for s , i.e., $X/Y := \{s : y \cdot s = x, x \in X, y \in Y\}$. Let $[\]$ denote the empty interval. With the following rules, division by intervals containing 0 becomes possible.

$$X/Y := \begin{cases} (-\infty, +\infty) & \text{if } 0 \in X, \text{ or } Y = [0, 0] \\ [\] & \text{if } 0 \notin X, \text{ and } Y = [0, 0] \\ [\bar{x}/\underline{y}, +\infty) & \text{if } \bar{x} \leq 0, \text{ and } \bar{y} = 0 \\ [\underline{x}/\bar{y}, +\infty) & \text{if } 0 \leq \underline{x}, \text{ and } 0 = \underline{y} < \bar{y} \\ (-\infty, \bar{x}/\bar{y}] & \text{if } \bar{x} \leq 0, \text{ and } 0 = \underline{y} < \bar{y} \\ (-\infty, \underline{x}/\underline{y}] & \text{if } 0 \leq \underline{x}, \text{ and } \underline{y} < \bar{y} = 0 \\ (-\infty, \bar{x}/\bar{y}] \cup [\bar{x}/\underline{y}, +\infty) & \text{if } \bar{x} \leq 0, \text{ and } [0, 0] \in Y \\ (-\infty, \underline{x}/\underline{y}] \cup [\underline{x}/\bar{y}, +\infty) & \text{if } 0 \leq \underline{x}, \text{ and } [0, 0] \in Y \end{cases}$$

When X is a thin interval with $x = \underline{x} = \bar{x}$ and Y has $+\infty$ or $-\infty$ as one of its bounds, then extended interval subtraction is also necessary for the extended interval Newton algorithm, and is defined as follows:

$$[\underline{x}, \bar{x}] - Y := \begin{cases} (-\infty, +\infty) & \text{if } Y = (-\infty, +\infty) \\ (-\infty, x - \underline{y}] & \text{if } Y = (\underline{y}, +\infty) \\ [x - \bar{y}, +\infty) & \text{if } Y = (-\infty, \bar{y}] \end{cases}$$

The extended interval Newton method sketched below uses the extended interval arithmetic described above and is a variant of the method based on Hansen and Sengupta (1981) with Ratz's modifications (Ratz, 1992) as implemented in Hammer et al. (1995). It can be used to enclose the roots of a continuously differentiable $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ in a given box $X \in \mathbb{I}\mathbb{R}^n$. Let $J_f(x) := ((\partial f_i(x)/\partial x_j))_{i,j=\{1,\dots,n\}} \in \mathbb{R}^{n \times n}$ denote the Jacobian matrix of f at x . Let $J_F(X) \supset J_f(X)$ denote the Jacobian of the interval extension of f . The Jacobian can be computed via automatic differentiation of section 2.2 by computing the gradient of each component f_i of f . By the mean value theorem, $f(m(X)) - f(x^*) = J_f(w) \cdot (m(X) - x^*)$, for some $x^* \in X, w = (w_1, w_2, \dots, w_n)$, where $w_i \in X, \forall i \in \{1, 2, \dots, n\}$. Interest in x^* with $f(x^*) = 0$, yields the following relation, provided $\forall x \in X, J_F(x)$ is invertible.

$$\begin{aligned} f(m(X)) &= J_f(w) \cdot (m(X) - x^*) \\ x^* &= m(X) - (J_f(w))^{-1} \cdot f(m(X)) \\ &\in m(X) - (J_f(X))^{-1} \cdot f(m(X)) \\ &\subseteq m(X) - (J_F(X))^{-1} \cdot F(m(X)) =: \mathcal{N}(X) \\ &\subseteq \mathcal{N}(X) \cap X \end{aligned}$$

An iteration scheme $X^{(j+1)} := \mathcal{N}(X^{(j)}) \cap X^{(j)}$, where $j = 0, 1, \dots$, and $X^{(0)} := X$, will enclose the zeros of f contained in X . To relax the assumption that every matrix in $J_F(X)$ be invertible, the inverse of the midpoint of $J_F(X)$, i.e., $(m(J_F(X)))^{-1} =: p \in \mathbb{R}^{n \times n}$, is used as a matrix preconditioner. The extended

interval Gauss-Seidel iteration, which is also applicable to singular systems (Neumaier, 1990), is used to solve the preconditioned interval linear equation,

$$\begin{aligned} p \cdot F(m(X)) &= p \cdot J_F(X) \cdot (m(X) - x^*) \\ a &= G \cdot (c - x^*), \end{aligned}$$

where, $a \in A := p \cdot F(m(X))$, $G := p \cdot J_F(X)$, and, $c := m(X)$. Thus, the solution set $\mathbf{S} := \{x \in X : g \cdot (c - x) = a, \forall g \in G\}$ of the interval linear equation $a = G \cdot (c - x)$ has the component-wise solution set $\mathbf{S}_i = \{x_i \in X_i : \sum_{j=1}^n (g_{i,j} \cdot (c_j - x_j)) = a_i, \forall g \in G\}$, $\forall i \in \{1, \dots, n\}$. Now, set $Y = X$, and solve the i th equation for the i th variable, iteratively for each i , as follows:

$$y_i = c_i - \frac{1}{g_{i,i}} \left(a_i + \sum_{j=1, j \neq i}^n (g_{i,j} \cdot (y_j - c_j)) \right) \in \left(c_i - \frac{1}{G_{i,i}} \left(A_i + \sum_{j=1, j \neq i}^n (G_{i,j} \cdot (Y_j - c_j)) \right) \right) \cap Y_i$$

The interval vector(s) Y obtained at the end of such an iteration is the set, $\mathcal{N}_{GS}(X)$, resulting from one extended interval Newton Gauss-Seidel step, such that, $\mathbf{S} \subseteq \mathcal{N}_{GS}(X) \subseteq X$. Thus, the roots of f are enclosed by the discrete dynamical system $X^{(j)} = \mathcal{N}_{GS}(X^{(j)})$ in \mathbb{R}^n . Every 0 of f that lies in X also lies in $\mathcal{N}_{GS}(X)$. If $\mathcal{N}_{GS}(X) = []$, the empty interval, then f has no solution in X . If $\mathcal{N}_{GS}(X) \Subset X$, then f has a unique solution in X . For proofs of the above three statements see Hansen (1992). When $G_{ii} \supset 0$, the method is applicable with extended interval arithmetic that allows for division by 0. In such cases, one may obtain upto two disjoint compact intervals for Y_i subsequent to extended interval arithmetic and intersection with the previous compact interval X_i . In such cases, the iteration is applied to each resulting sub-interval.

2.4. Machine interval arithmetic

All interval arithmetic was done above with real intervals. However, there are only finitely many floating-point numbers available on a computing machine. Let \mathcal{R} be this set of floating-point numbers. A machine interval is a real interval with floating-point bounds. Thus, on a computer, one works with $\mathbb{IR} := \{X \in \mathbb{IR} : \underline{x}, \bar{x} \in \mathcal{R}\}$, the set of all machine intervals. In spite of the finiteness of \mathbb{IR} , the strength of interval arithmetic lies in a machine interval X being able to enclose the entire continuum of reals between its machine-representable boundaries. Through rounding controlled floating-point arithmetic provided by the IEEE arithmetic standard, operations with real intervals can be tightly enclosed by the rounding directed operations with the smallest machine intervals containing them. See Kulisch (2001) for a recent description of machine interval arithmetic. The errors resulting from converting a decimal number, usually a constant or input data, which in general does not have a finite binary representation, to a binary floating-point number is controlled by first passing the decimal number as a string and then enclosing it with the smallest machine interval by proper outward rounding. The program is written in C++ using the C-XSC class libraries. The differentiation arithmetic of section 2.2 is implemented using the *hess_ari* module provided in Hammer et al. (1995).

2.5. Interval experiment

DEFINITION 1. The (statistical) experiment $\mathcal{E}_{\mathcal{P}} := (\mathcal{X}, \mathcal{F}, \mathcal{P})$, which is a measurable space $(\mathcal{X}, \mathcal{F})$ endowed with a family $\mathcal{P} = \{P_{\theta} : \theta \in \Theta\}$ of probability measures that are indexed by the compact set Θ , is said to be identifiable if and only if

$$\theta \neq \vartheta \Rightarrow P_{\theta} \neq P_{\vartheta}, \forall \theta, \vartheta \in \Theta$$

DEFINITION 2. The interval experiment $\mathcal{E}_{\mathcal{IP}} := (\mathcal{X}, \mathcal{F}, \mathcal{IP})$, is an extension of the compact index set Θ of the usual experiment $\mathcal{E}_{\mathcal{P}}$, to $\mathbb{I}\Theta$. This extended index set $\mathbb{I}\Theta$ indexes a collection \mathcal{IP} of compact sets of measures from \mathcal{P} ,

$$\mathcal{IP} = \{P_{\Theta} : \Theta \in \mathbb{I}\Theta\}, \quad \mathbb{I}\Theta := \{\Theta = [\underline{\theta}, \bar{\theta}] : \underline{\theta} \leq \bar{\theta}, \text{ and } \underline{\theta}, \bar{\theta} \in \Theta\}$$

DEFINITION 3. The experiment $\mathcal{E}_{\mathcal{IP}}$ is said to be identifiable if and only if

$$\Theta \neq \Theta \Rightarrow P_{\Theta} \neq P_{\Theta} \Leftrightarrow P_{\Theta} \Delta P_{\Theta} := \{P_{\Theta} \setminus P_{\Theta}\} \cup \{P_{\Theta} \setminus P_{\Theta}\} \neq \emptyset, \forall \Theta, \Theta \in \Theta$$

LEMMA 1. *The interval extension $\mathcal{E}_{\mathcal{IP}}$ of an identifiable experiment $\mathcal{E}_{\mathcal{P}}$ is identifiable, provided the mapping $P_{\Theta} : \mathbb{I}\Theta \rightarrow \mathcal{IP}$ is inclusion isotonic.*

PROOF. Let $\mathcal{E}_{\mathcal{IP}}$ be the interval extension of the experiment $\mathcal{E}_{\mathcal{P}}$ with an inclusion isotonic mapping $P_{\Theta} : \mathbb{I}\Theta \rightarrow \mathcal{IP}$. Suppose $\mathcal{E}_{\mathcal{P}}$ is identifiable. We need to show that $\mathcal{E}_{\mathcal{IP}}$ is also identifiable. Let Θ and Θ' be distinct non-empty elements of $\mathbb{I}\Theta$. Without loss of generality, it follows that $\{\Theta \setminus \Theta'\} \neq \emptyset$, since, $\Theta \neq \Theta' \Leftrightarrow \Theta \Delta \Theta' \neq \emptyset \Leftrightarrow \{\Theta \setminus \Theta'\} \cup \{\Theta' \setminus \Theta\} \neq \emptyset$. Thus, $\exists \theta' \in \{\Theta \setminus \Theta'\} \neq \emptyset$, and $P_{\theta'} \in P_{\Theta}$ by inclusion isotony. We will prove the identifiability of $\mathcal{E}_{\mathcal{IP}}$ by contradiction. Now, suppose to the contrary, that $P_{\Theta} \Delta P_{\Theta'} = \emptyset$ or equivalently that $P_{\Theta} = P_{\Theta'}$. Then, $P_{\theta'} \in P_{\Theta} \Rightarrow P_{\theta'} \in P_{\Theta'}$. By inclusion isotony $\exists \theta'' \in \Theta'$ such that $P_{\theta''} = P_{\theta'}$. But, $\theta' \in \{\Theta \setminus \Theta'\} \Rightarrow \theta' \neq \theta''$. Thus, we have that $\theta' \neq \theta''$ and yet $P_{\theta'} = P_{\theta''}$. This contradicts the assumed identifiability of $\mathcal{E}_{\mathcal{P}}$ and concludes the proof.

An identifiable interval experiment $\mathcal{E}_{\mathcal{IP}}$ allows for the notion of statistical consistency in the complete metric space of the index set $\mathbb{I}\Theta$ under the Hausdorff metric \mathfrak{h} . An estimator $\widehat{\Theta}_n$ of Θ from n observations is asymptotically consistent if $\mathfrak{h}(\widehat{\Theta}_n - \Theta) \xrightarrow{P} 0$. Observe that $\mathbb{I}\Theta$ is a complete lattice, *ie*, every subset has an infimum and a supremum and its machine counterpart is a complete sublattice w.r.t. \leq and \subseteq . The machine interval experiment is an interval experiment indexed by Θ 's with floating-point bounds in $\mathcal{R} = \mathcal{R}(2, p, \underline{e}, \bar{e})$. Thus, in a machine interval experiment, the estimator $\widehat{\Theta}_n$ may not converge in probability to a thin box even if the true measure Θ is thin and all observations are measurable with infinite precision, since $p < \infty$. The diameter $d(\widehat{\Theta}_n)$ reflects the limit of numerical resolution caused by computations that account for all numerical errors using the machine's finite number screen.

3. A PHYLOGENETIC PROBLEM

Let D denote a homologous set of distinct DNA sequences of length v from s species. The objective of this paper is to find the maximum likelihood estimates of branch lengths for the best tree under a particular topology. Recall that the branch lengths usually represent a scaled product of mutation rate and number of generations. Let b denote the number of branches and n denote the number of nodes of a tree with topology τ . For an s -leaved unrooted tree of a given topology, there are at most $2s - 3$ branches, *i.e.*, $b \leq 2s - 3$. Since the number of topologies for multifurcating unrooted trees grows as a factorial of the number of leaves, it is difficult to exhaustively search through all possible topologies to find the most likely tree even when the number of leaves is reasonably small. For example, there are 12,818,912 topologies when $s = 10$. See chapter 3 of Felsenstein (2003) to appreciate this problem. However, one can find the most likely tree among a small set of specified topologies, by first computing the most likely tree under each topology of interest and then choosing the tree with the highest likelihood.

Thus, for a given unrooted topology τ with s leaves and b branches, the unknown parameter $\theta = (\theta_1, \dots, \theta_b)$ is the real vector of branch lengths in the positive orthant, where each positive branch length $\theta_q \in [\theta_{\delta}, R] \rightarrow \mathbb{R}_+$, as $\theta_{\delta} \rightarrow 0$ and $R \rightarrow \infty$. An explicit model of DNA evolution is needed to construct the likelihood function which gives the probability of observing data D as a function of the parameter θ . The simplest such continuous time Markov chain model (JC69) on the state space $\mathcal{S} := \{A, G, C, T\}$ is due to Jukes and Cantor (1969) with the rate of mutation between nucleotides i and j given by $q_{i,j} = 1/3$, if $i \neq j$, and -1 , otherwise. Its stationary distribution $\pi = (1/4, 1/4, 1/4, 1/4)$, and $P_{i,j}(t)$, the probability of transition from i to j in time t is, $1/4 + 3/4 \exp(-4t/3)$ if $i = j$, and $1/4 - 1/4 \exp(-4t/3)$, otherwise. Felsenstein's algorithm (Felsenstein, 1981) to compute $\ell^{(k)}(\theta)$, the likelihood at site $k \in \{1, \dots, v\}$, is the following postorder traversal:

- (a) Associate with each node $q \in \{1, \dots, n\}$ a real vector $\ell_q := (\ell_q^A, \ell_q^C, \ell_q^G, \ell_q^T) \in \mathbb{R}^4$, and let the length of the branch leading to its ancestor be θ_q .
- (b) For a leaf node q with nucleotide i , set $\ell_q^i = 1$ and $\ell_q^j = 0$ for all $j \neq i$. For any internal node q , set $\ell_q := (1, 1, 1, 1)$.
- (c) For an internal node q with descendants s_1, s_2, \dots, s_m ,

$$\ell_q^i = \sum_{j_1, \dots, j_m \in \mathcal{S}} \{ \ell_{s_1}^{j_1} \cdot P_{i,j_1}(\theta_{s_1}) \cdot \ell_{s_2}^{j_2} \cdot P_{i,j_2}(\theta_{s_2}) \cdot \dots \cdot \ell_{s_m}^{j_m} \cdot P_{i,j_m}(\theta_{s_m}) \}$$

- (d) Compute ℓ_q for each sub-terminal node q , then those of their ancestors recursively to finally compute ℓ_r for the root node r and obtain the likelihood,

$$\ell^{(k)}(\theta) = \sum_{i \in \mathcal{S}} (\pi_i \cdot \ell_r^i)$$

for each site k .

Assuming independence across sites one obtains the likelihood function for the entire sequence by multiplying the site-specific likelihoods together. The problem of finding the global maximum of this likelihood function is equivalent to finding the global minimum of $l(\theta)$, the negative of the natural logarithm of the likelihood function given by,

$$l(\theta) = - \sum_{k=1}^v \ln \ell^{(k)}(\theta).$$

$l(\theta)$ is of interest because algorithms in the optimization literature are usually addressed in terms of minimization. Replacing θ , a positive real vector of branch lengths, in the above algorithm by a positive real interval vector or box Θ and all real operations by their interval counterparts, yields $L(\Theta)$, the natural interval extension of the negative log likelihood function $l(\theta)$ over Θ . Since $\nabla L(\Theta)$ and $\nabla^2 L(\Theta)$, the enclosures of the gradient and the Hessian of $l(\theta)$ over Θ , respectively, are needed in Section 4, one may use the constant triples, $(C, 0, 0)$, variable triples, $(\Theta_j, e^{(j)}, 0)$, appropriate triples for the elementary functions, exp and ln, and perform all operations in the interval differentiation arithmetic of Section 2.2, in order to obtain the negative log likelihood triple $(L(\Theta), \nabla L(\Theta), \nabla^2 L(\Theta))$.

4. GLOBAL OPTIMIZATION

4.1. Branch-and-bound

The most basic strategy in global optimization through enclosure methods is to employ rigorous branch-and-bound techniques. Such techniques recursively partition (branch) the original compact space of interest into compact subspaces and discard (bound) those subspaces that are guaranteed to not contain the global optimizer(s). For the real scalar-valued multi-dimensional objective function $l(\theta)$, the interval branch-and-bound technique can be applied to its natural interval extension $L(\Theta)$ to obtain an interval enclosure L^* of the global minimum value l^* as well as the set of minimizer(s) to a specified accuracy ϵ . Note that this set of minimizer(s) of $L(\theta)$ is the set of maximizer(s) of the likelihood function for the observed data D . The strength of such methods arises from the algorithmic ability to discard large sub-boxes from the original search region,

$$\Theta^{(0)} = (\Theta_1^{(0)}, \dots, \Theta_b^{(0)}) := ([\underline{\theta}_1^{(0)}, \bar{\theta}_1^{(0)}], \dots, [\underline{\theta}_b^{(0)}, \bar{\theta}_b^{(0)}]) \subset \mathbb{R}^b,$$

that are not candidates for global minimizer(s). Four tests that help discard sub-regions are described below. Let \mathfrak{L} denote a list of ordered pairs of the form $(\Theta^{(i)}, \underline{L}_{\Theta^{(i)}})$, where, $\Theta^{(i)} \subseteq \Theta^{(0)}$, and $\underline{L}_{\Theta^{(i)}} := \min(L(\Theta^{(i)}))$ is a lower bound for the range of the negative log likelihood function l over $\Theta^{(i)}$. Let \tilde{l} be an upper bound for l^* and $\nabla L(\Theta^{(i)})_k$ denote the k -th interval of the gradient box $\nabla L(\Theta^{(i)})$. If no information is available for \tilde{l} , then $\tilde{l} = \infty$.

- (a) *Midpoint Cut-off test*: The basic idea of the *midpoint cut-off test* is to discard sub-boxes of the search space $\Theta^{(0)}$ with the lower bound for their range enclosures above \tilde{l} , the current best estimate of an upper bound for l^* . Figure 7 shows a multi-modal l as a function of a scalar valued θ over $\Theta^{(0)} = \cup_{i=1}^{16} \Theta^{(i)}$. For this illustrative example, \tilde{l} is set as the upper bound of the range enclosure of l over the smallest machine interval containing the midpoint of $\Theta^{(15)}$, the interval with the smallest lower bound of its range enclosure. The shaded rectangles show the range enclosures over intervals that lie strictly above \tilde{l} . In this example the *midpoint cut-off test* would discard all other intervals except $\Theta^{(1)}$, $\Theta^{(2)}$, and $\Theta^{(4)}$.

- Given a list \mathfrak{L} and \tilde{l}
- Choose an element j of \mathfrak{L} , such that, $j = \underset{i}{\operatorname{argmin}} \underline{L}_{\Theta^{(i)}}$, since $\Theta^{(j)}$ is likely to contain a minimizer.

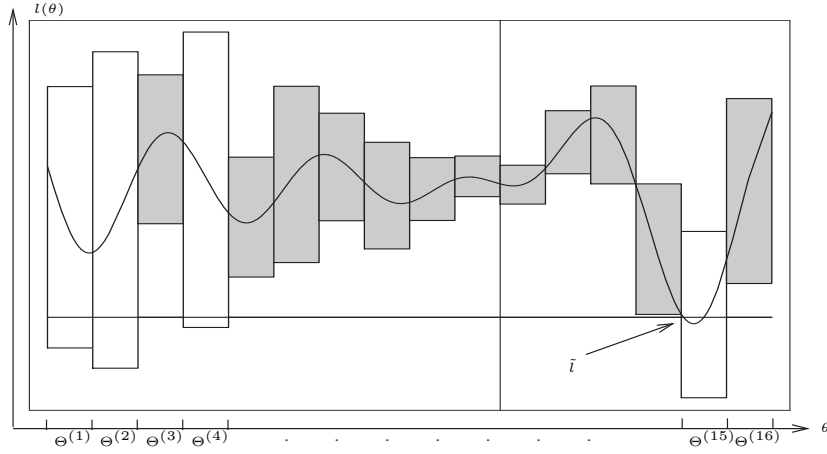


Fig. 7. Midpoint Cut-off test

- Find its midpoint $c = m(\Theta^{(j)})$ and let C be the smallest machine interval containing c .
 - Compute a possibly improved $\tilde{l} = \min \{\tilde{l}, \bar{L}_C\}$, where, $\bar{L}_C := \max(L(C))$
 - Discard any i -th element of \mathfrak{L} for which $\underline{L}_{\Theta^{(i)}} > \tilde{l} \geq l^*$
- (b) *Monotonicity test*: For a continuously differentiable function $l(\theta)$, the *monotonicity test* determines whether $l(\theta)$ is strictly monotone over an entire sub-box $\Theta^{(i)} \subset \Theta^{(0)}$. If l is strictly monotone over $\Theta^{(i)}$, then a global minimizer cannot lie in the interior of $\Theta^{(i)}$. Therefore, $\Theta^{(i)}$ can only contain a global minimizer as a boundary point if this point also lies in the boundary of $\Theta^{(0)}$. Figure 8 illustrates the *monotonicity test* for the one-dimensional case. In this example the search space of interest, $\Theta^{(0)} = [\underline{\theta}^{(0)}, \bar{\theta}^{(0)}] = \cup_{i=1}^8 \Theta^{(i)}$, can be reduced considerably. In the interior of $\Theta^{(0)}$, one may delete $\Theta^{(2)}$, $\Theta^{(5)}$, and $\Theta^{(7)}$, since $l(\theta)$ is monotone over them as indicated by the enclosure of the derivative $l'(\theta)$ being bounded away from 0. Since $l(\theta)$ is monotonically decreasing over $\Theta^{(1)}$ one can also delete it since we are only interested in minimization. $\Theta^{(8)}$ may be pruned to its right boundary point $\theta^{(8)} = \bar{\theta}^{(8)} = \bar{\theta}^{(0)}$ due to the strictly decreasing nature of $l(\theta)$ over it. Thus, the *monotonicity test* has pruned $\Theta^{(0)}$ to the smaller candidate set $\{\bar{\theta}^{(0)}, \Theta^{(3)}, \Theta^{(4)}, \Theta^{(6)}\}$ for a global minimizer.
- Given $\Theta^{(0)}$, $\Theta^{(i)}$, and $\nabla L(\Theta^{(i)})$
 - Iterate for $k = 1, \dots, b$
 - If $0 \in \nabla L(\Theta^{(i)})_k$, then leave $\Theta_k^{(i)}$ unchanged, as it may contain a stationary point of l .
 - Otherwise, $0 \notin \nabla L(\Theta^{(i)})_k$. This implies that $\Theta^{(i)}$ can be pruned, since $l^* \notin \Theta^{(i)}$ except possibly at the boundary points, as follows:
 - (i) if $\min(\nabla L(\Theta^{(i)})_k) > 0$ and $\underline{\theta}_k^{(0)} = \underline{\theta}_k^{(i)}$, then $\Theta_k^{(i)} = [\underline{\theta}_k^{(i)}, \underline{\theta}_k^{(i)}]$,
 - (ii) Else if $\max(\nabla L(\Theta^{(i)})_k) < 0$ and $\bar{\theta}_k^{(0)} = \bar{\theta}_k^{(i)}$, then $\Theta_k^{(i)} = [\bar{\theta}_k^{(i)}, \bar{\theta}_k^{(i)}]$.
 - (iii) Else, delete the i -th element of \mathfrak{L} and stop the iteration.
- (c) *Concavity test*: Given $\Theta^{(i)} \in \Theta^{(0)}$, and the diagonal elements $(\nabla^2 L(\Theta^{(i)}))_{kk}$ of $\nabla^2 L(\Theta^{(i)})$, note that if $\min((\nabla^2 L(\Theta^{(i)}))_{kk}) < 0$ for some k , then, $\nabla^2 L(\Theta^{(i)})$ cannot be positive semidefinite, and therefore $l(\theta)$ cannot be convex over $\Theta^{(i)}$ and thus cannot contain a minimum in its interior. In the one-dimensional example shown in Figure 8, an application of the *concavity test* to the candidate set $\{\underline{\theta}^{(0)}, \Theta^{(4)}, \Theta^{(6)}\}$ for a global minimizer returned by the *monotonicity test*, would result in the deletion of $\Theta^{(6)}$ due to the concavity of $l(\theta)$ over it.
- Given $\Theta^{(i)} \in \Theta^{(0)}$ and $\nabla^2 L(\Theta^{(i)})$
 - If $\min((\nabla^2 L(\Theta^{(i)}))_{kk}) < 0$ for any $k \in \{1, \dots, b\}$, then delete the i -th element of \mathfrak{L} .

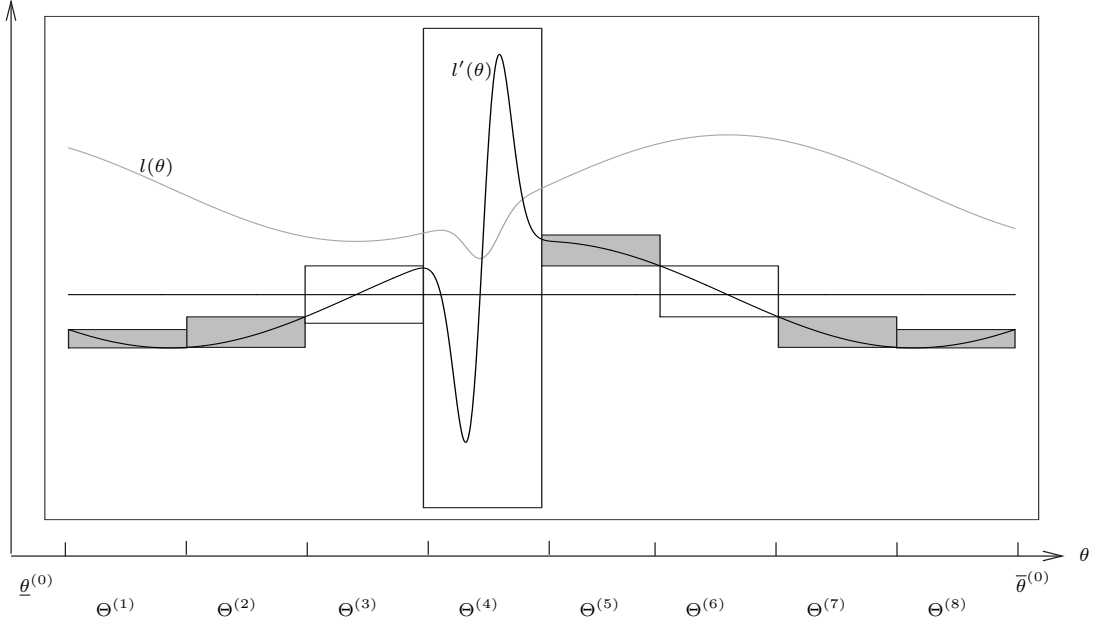


Fig. 8. Monotonicity test

(d) *Interval Newton test:* Given $\Theta^{(i)} \in \Theta^{(0)}$, and $\nabla L(\Theta^{(i)})$, attempt to solve the system, $\nabla L(\theta) = 0$, in terms of $\theta \in \Theta^{(i)}$.

- Apply one extended interval Newton Gauss-Seidel step of Section 2.3 to the linear interval equation $a = G \cdot (c - \theta)$, where, $a := p \cdot L(m(\Theta^{(i)}))$, $G := p \cdot \nabla^2 L(\Theta^{(i)})$, $c := m(\Theta^{(i)})$, and $p := (m(\nabla^2 F(X)))^{-1}$, in order to obtain $\mathcal{N}'_{GS}(\Theta^{(i)})$.
- One of the following can happen,
 - (i) If $\mathcal{N}'_{GS}(\Theta^{(i)})$ is empty, then discard $\Theta^{(i)}$.
 - (ii) If $\mathcal{N}'_{GS}(\Theta^{(i)}) \subseteq \Theta^{(i)}$, then replace $\Theta^{(i)}$ by the contraction $\mathcal{N}'_{GS}(\Theta^{(i)}) \cap \Theta^{(i)}$.
 - (iii) If $0 \in G_{jj}$, and the extended interval division splits $\Theta_j^{(i)}$ into a non-empty union of $\Theta_j^{(i),1}$ and $\Theta_j^{(i),2}$, then the iteration is continued on $\Theta_j^{(i),1}$, while $\Theta_j^{(i),2}$, if non-empty, is stored in \mathcal{L} for future processing. Thus, one extended interval Newton Gauss-Seidel step can add at most $b + 1$ sub-boxes to \mathcal{L} .

4.2. Verification

Given a collection of sub-boxes, $\{\Theta^{(1)}, \dots, \Theta^{(n)}\}$, each of width $\leq \epsilon$, that could not be discarded by the tests in Section 4.1, one can attempt to verify the existence and uniqueness of a local minimizer within each sub-box $\theta^{(i)}$ by checking whether the conditions of the following two theorems are satisfied. For proof of these two theorems see Hansen (1992) and Ratz (1992).

- (a) If $\mathcal{N}'_{GS}(\Theta^{(i)}) \subseteq \Theta^{(i)}$, then there exists a unique stationary point of L , i.e., a unique zero of ∇L exists in $\Theta^{(i)}$.
- (b) If $(I + \frac{1}{\kappa} \cdot (\nabla^2 L(\Theta^{(i)}))) \cdot Z \subseteq Z$, where $(\nabla^2 L(\Theta^{(i)}))_{d,\infty} \leq \kappa \in \mathbb{R}$, for some $Z \in \mathbb{IR}^n$, then, the spectral radius $\rho(s) < 1$ for all $s \in (I - \frac{1}{\kappa} \cdot (\nabla^2 L(\Theta^{(i)})))$, and all symmetric matrices in $\nabla^2 L(\Theta^{(i)})$ are positive definite.

If the conditions of the above two theorems are satisfied by some $\Theta^{(i)}$, then a unique stationary point exists in $\Theta^{(i)}$ and this stationary point is a local minimizer. Therefore, if exactly one candidate sub-box for minimizer(s) remained after pruning the search box $\Theta^{(0)}$ with the tests in Section 4.1, and if this sub-box satisfies the above two conditions for the existence of a unique local minimizer within it, then one has

rigorously enclosed the global minimizer in the search interval. On the other hand, if there are two or more sub-boxes in our candidate list for minimizer(s) that satisfy the above two conditions, then one may conclude that each sub-box contains a candidate for a global minimizer which may not necessarily be unique (disconnected sub-boxes, for example). Observe that failure to verify the uniqueness of a local minimizer in a sub-box can occur if it contains two or more points or even a continuum of points that are stationary (nonidentifiable manifolds in the sub-box, for example).

4.3. Algorithm

- *Initialization:*

- Let the search region be a single box $\Theta^{(0)}$ or a collection of not necessarily connected, but pair-wise disjoint boxes, $\Theta^{(i)}$, $i \in \{1, \dots, r\}$.
- Initialize the list \mathfrak{L} which may just contain one element $(\Theta^{(0)}, \underline{L}_{\Theta^{(0)}})$ or several elements

$$\{(\Theta^{(1)}, \underline{L}_{\Theta^{(1)}}), (\Theta^{(2)}, \underline{L}_{\Theta^{(2)}}), \dots, (\Theta^{(r)}, \underline{L}_{\Theta^{(r)}})\}.$$

- Let ϵ be a specified tolerance.
- Let $\max_{\mathfrak{L}}$ be the maximal length allowed for list \mathfrak{L} .
- Set the noninformative lower bound for l^* , i.e., $\tilde{l} = \infty$

- *Iteration:*

- Improve $\tilde{l} = \min\{\tilde{l}, \max(L(m(\Theta^{(j)})))\}$, where $j = \underset{i}{\operatorname{argmin}}\{\underline{L}_{\Theta^{(i)}}\}$.
 - Perform the *midpoint cut-off test* to \mathfrak{L} .
 - Set $L^* = [\underline{L}_{\Theta^{(j)}}, \tilde{l}]$.
- Bisect $\Theta^{(j)}$ along its longest side k , i.e., $d(\Theta_k^{(j)}) = d_{\infty}(\Theta^{(j)})$, to obtain sub-boxes $\Theta^{(j_q)}$, $q \in \{1, 2\}$.
- For each sub-box $\Theta^{(j_q)}$, evaluate its triple $(L(\Theta^{(j_q)}), \nabla L(\Theta^{(j_q)}), \nabla^2 L(\Theta^{(j_q)}))$, and do the following:

- Perform *monotonicity test* to possibly discard $\Theta^{(j_q)}$.
- Centered form cut-off test:*
Improve the range enclosure of $L(\Theta^{(j_q)})$ by replacing it with its centered form $L_c(\Theta^{(j_q)})$,

$$L_c(\Theta^{(j_q)}) := \{L(m(\Theta^{(j_q)})) + \nabla L(\Theta^{(j_q)}) \cdot (\Theta^{(j_q)} - m(\Theta^{(j_q)}))\} \cap L(\Theta^{(j_q)}),$$

and then discarding $\Theta^{(j_q)}$, if $\tilde{l} < \underline{L}_{\Theta^{(j_q)}}$.

- Perform *concavity test* to possibly discard $\Theta^{(j_q)}$.
- Apply an *extended interval Newton Gauss-Seidel step* to $\Theta^{(j_q)}$, in order to either entirely discard it or shrink it into v sub-sub-boxes, where v is at most $2s - 2$.
- For each one of these sub-sub-boxes $\Theta^{(j_q, u)}$, $u \in \{1, \dots, v\}$
 - Perform *monotonicity test* to possibly discard $\Theta^{(j_q, u)}$.
 - Try to discard $\Theta^{(j_q, u)}$ by applying the *centered form cut-off test* in cii to it.
 - Append $(\Theta^{(j_q, u)}, \underline{L}_{\Theta^{(j_q, u)}})$ to \mathfrak{L} if $\Theta^{(j_q, u)}$ could not be discarded by steps c(v)A and c(v)B.

- *Termination:*

- Terminate iteration if $d_{rel, \infty}(\Theta^{(j)}) < \epsilon$, or $d_{rel, \infty}(L^*) < \epsilon$, or \mathfrak{L} is empty, or $\text{Length}(\mathfrak{L}) > \max_{\mathfrak{L}}$
- Verify uniqueness of minimizer(s) in the final list \mathfrak{L} by applying algorithm of section 4.2 to each of its elements.

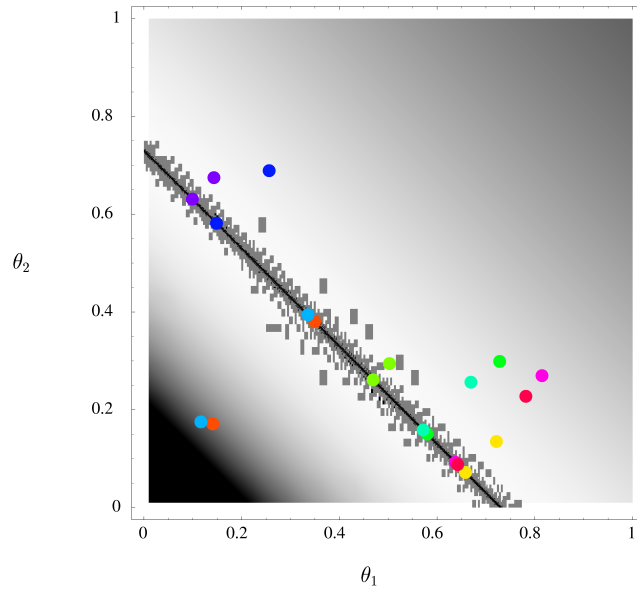


Fig. 9. For a pair of homologous sequences of 600 nucleotides out of which 280 sites are polymorphic, the nonidentifiable subspace of minimizers $\theta_1 + \theta_2 = \frac{3}{4} \log(45/17) = 0.730087$ of the negative log likelihood function under the JC69 model evolving on a rooted two-leaved tree is enclosed by a union of upto 30,000 boxes. The larger grey, and smaller black boxes have tolerances of $\epsilon = 1.0 \times 10^{-4}$ and $\epsilon = 1.0 \times 10^{-6}$, respectively. The 10 pairs of colored circles are the initial and final points of 10 BFGS searches with random initializations.

5. APPLICATIONS

5.1. Enclosing nonidentifiable subspaces

For time reversible Markov chains, such as JC69, evolving on a rooted tree, only the sum of the branch lengths emanating from the root is identifiable. Identifiability is a prerequisite for statistical consistency of estimators. To demonstrate the ability of interval methods, unlike the local search methods, to enclose the nonidentifiable ridge along $\theta_1 + \theta_2$, in the simplest case of a 2-leaved tree, a nonidentifiable negative log likelihood function $l(\theta)$ is formulated and its global minimizers along $\theta_1 + \theta_2 = \frac{3}{4} \log(45/17) = 0.730087$ are enclosed as shown in Figure 9 for a fictitious dataset for which 280 out of 600 sites were polymorphic. Observe that the basin of attraction for each point on $\theta_1 + \theta_2 = 0.730087$ under the BFGS local search algorithm is the line running orthogonal to it. This trivial example is only chosen for pedantic reasons. Enclosing possibly nonidentifiable submanifolds, that may not even be simply connected, within any compact subset of higher dimensional parameter spaces, may be accomplished, at least partly, by studying the rates of decay of the hyper-volume of the union of all pending boxes as the algorithm progresses, for instance.

5.2. Unrooted 3-leaved Tree

The global maximum of the log likelihood function for the JC69 model of DNA evolution on the three taxa unrooted tree with data from the mitochondria of Chimpanzee, Gorilla, and Orangutan (Brown et al., 1982) is enclosed. There is only one unrooted multifurcating topology for three species with all three branches emanating from the root like a star. The data set for this problem can be summarized by the following 29 data patterns:

PATTERN COUNTS :

232 71 229 168 13 31 16 18 9 20 1 8 22 3 10 8 7 1 9 2 4 2 1 2 1 1 2 1 3

PATTERNS:

agctatcaccatctgcccgtactaagcgt

agctgttatcaacacgcaaatccggtat

agctaccgttcccataataataaagcgca

Table 1. Enclosures of the Maximum log likelihood and their corresponding parameter estimates for 3 taxa tree relating Chimpanzee, Gorilla, and Orangutan.

Tree	$\Theta^{(0)}$	$\Theta^* \supset \theta^*$	$-L(\Theta^*) \supset -l(\theta^*)$
star	$[1.0 \times 10^{-11}, 1.0 \times 10^9]^{\otimes 3}$	$5.9816221384_0^2 \times 10^{-2}$ $5.4167416794_0^2 \times 10^{-2}$ $1.3299089685_8^9 \times 10^{-1}$	$-2.150318065856_6^5 \times 10^3$

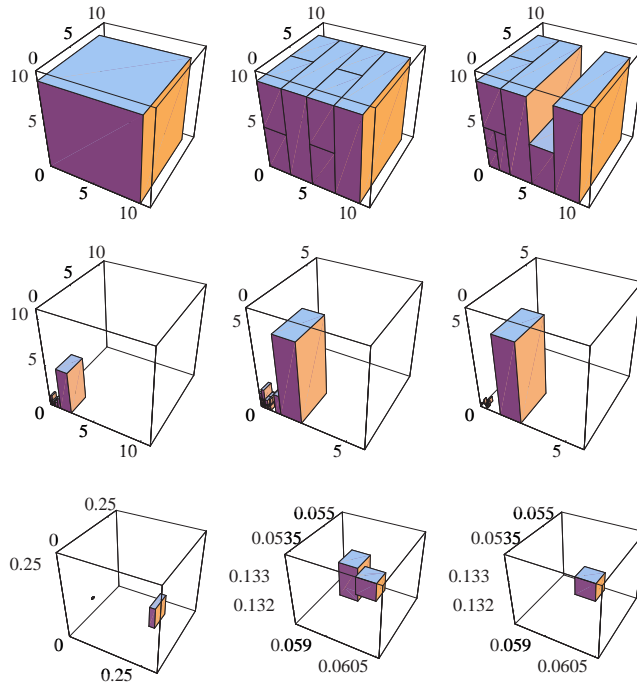


Fig. 10. Progress of the algorithm as it prunes $[0.001, 10.0]^{\otimes 3}$.

The parameter space is three dimensional corresponding to the three branch lengths of the 3-leaved tree. The algorithm is given a large search box $\Theta^{(0)}$. The results are summarized in Table 1. The notation x_a^b means the interval $[x_a, x_b]$, for e.g. $5.9816221384_0^2 \times 10^{-2} = [5.98162213840 \times 10^{-2}, 5.98162213842 \times 10^{-2}]$. Figure 10 shows the the parameter space being rigorously pruned as the algorithm progresses according to section 4.

5.3. Four Unrooted 4-leaved Trees

By adding the homologous mitochondrial sequence from Gibbon to the previous problem, one has the simplest phylogeny estimation problem with the following 61 data patterns:

PATTERN COUNTS :

209 71 192 157 28 5 11 20 2 10 10 15 5 15 1 5 1 2 15 3 14 3 4 2 5 4 5
 4 9 2 4 5 1 1 7 6 3 1 1 4 2 3 3 1 1 2 3 1 1 1 4 1 1 1 1 1 1 1 1 2 1

PATTERNS:

agctccatcacctacaatccatctgtcgggattaccctatttacgcgtcgcttc
 agctccgtgttatcctaaaaacatacacctgccaagaatactttcccatgctattacctt
 agctccacaccgttcaccaacatctccctatataaagattaccaaacatgtcgcagattaa
 agcttaacgtcaccactcgtcatgcatgctcctaagtgaatacaaaaactgaaagcaata

Four topologies are considered for a tree with four leaves (Figure 11). The star tree τ_1 has all four lineages

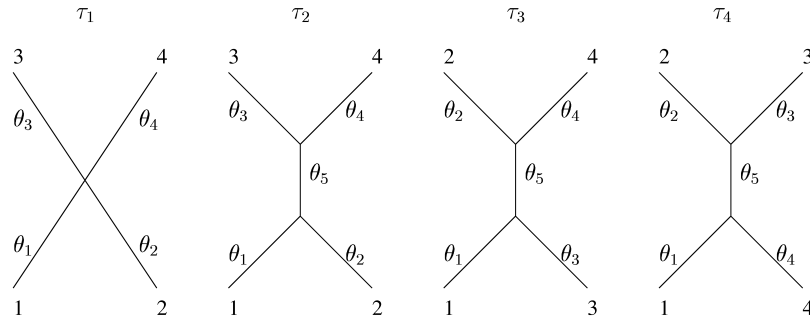


Fig. 11. The four different topologies, $\tau_1, \tau_2, \tau_3,$ and τ_4 with 4 leaves. The four leaves 1, 2, 3, and 4 denote the four primates Chimpanzee, Gorilla, Orangutan, and Gibbon, respectively.

coalescing at the same time, while the other three trees have an additional parameter θ_5 representing the only internal branch length. They differ due to the order in which the leaves relate to one another as shown in Figure 11. The parameter space is four dimensional for the star topology τ_1 , and five dimensional for each of the unrooted topologies $\tau_2, \tau_3,$ and τ_4 .

Observe that τ_1 is really a special case of the other three trees, $\tau_2, \tau_3,$ and τ_4 , as their internal branch θ_5 vanishes. Since we assume that the branches are $\geq \theta_\delta > 0$, and let $\theta_\delta \rightarrow 0$ on a sequence of floating-point numbers, it is convenient to treat the star tree τ_1 separately. The algorithm is given a large search box $\Theta^{(0)}$ for each topology and the results are summarized in Table 2.

Within each one of the four topologies there exists a unique global maximum. However, the global maximizer over all five topologies falls under topology τ_2 with the global maximum $-l^*$ contained in the interval $-L^* = -2.656936470946_6^5 \times 10^3$.

6. CONCLUSIONS

A general procedure has been provided to rigorously enclose the maximum likelihood value, as well as the most likely branch lengths of a tree (with a specified topology) upon which the simplest Markov model of DNA evolution is superimposed. The global optimization algorithm is general and can be used by frequentists to conduct rigorous numerical inference in a likelihood framework when the statistical experiment is indexed by a compact subset of a finite dimensional continuum. This procedure is not susceptible to errors caused by *undirected rounding, catastrophic cancellation, discretization, conversion, and ill-posed model*. Modifications of this algorithm are also applicable to Markov models with unknown parameters. For models that have parameter constraints, several constrained global optimization algorithms already exist (Hansen, 1992). When analytical spectral decompositions are not available for more complicated Markov models, one may use one of several rigorous eigen system solvers (e.g. Mayer (1994)) to compute the transition probabilities.

The *extended interval Newton* method in combination with the *midpoint cut-off, monotonicity,* and *concavity* tests may be used to study the shape of the likelihood surface itself. For instance, one could rigorously enclose all the local maxima, or fish for nonidentifiable subspaces, above a given level-set of the likelihood function within any compact subset of the parameter space. Several efficiency increasing steps could be taken. Pre-enclosing the transition probabilities and accessing them through hash functions can save computational effort. Asynchronous parallelization of the algorithm across 6 processors is also observed to increase the rate of convergence to the global maximum. It also provides a natural framework to manage the memory requirements for larger trees through partial likelihood evaluations for non-overlapping subtrees in parallel prior to obtaining the full likelihood.

Finally, it is worth noting that inclusion isotony does indeed hold by the continuity of the likelihood function in the $CAT(0)$ space of trees (Billera et al., 2004), and thus in conjunction with interval analysis may be made to provide a rigorous numerical framework for global maximization of the likelihood over compact sets containing distinct topologies. Preliminary results indicate that efficiency increases when one starts with

Table 2. Enclosures of the Maximum log likelihood and their corresponding parameter estimates for 4 taxa trees relating Chimpanzee, Gorilla, Orangutan, and Gibbon.

Tree	$\Theta^{(0)}$	$\Theta^* \supset \theta^*$	$-L(\Theta^*) \supset -l(\theta^*)$
τ_1	$[1.0 \times 10^{-11}, 1.0 \times 10^9]^{\otimes 4}$	$6.57882493333_4^5 \times 10^{-2}$ $6.236162512403_2^8 \times 10^{-2}$ $1.324874902248_4^5 \times 10^{-1}$ $1.635912562476_3^4 \times 10^{-1}$	$-2.7027434501964_4^1 \times 10^3$
τ_2	$[1.0 \times 10^{-11}, 1.0 \times 10^9]^{\otimes 5}$	$4.96281934326_8^9 \times 10^{-2}$ $5.89926424690_7^8 \times 10^{-2}$ $5.51849077387_3^4 \times 10^{-2}$ $9.09714007596_2^3 \times 10^{-2}$ $1.231516018310_1^2 \times 10^{-1}$	$-2.656936470946_6^5 \times 10^3$
τ_3	$[1.0 \times 10^{-11}, 1.0 \times 10^9]^{\otimes 5}$	$9.0717704_6^7 \times 10^{-3}$ $6.14239111_3^4 \times 10^{-2}$ $1.296383822_4^5 \times 10^{-1}$ $5.650692181_0^3 \times 10^{-2}$ $1.60005431656_1^7 \times 10^{-1}$	$-2.69987813617_5^0 \times 10^3$
τ_4	$[1.0 \times 10^{-11}, 1.0 \times 10^9]^{\otimes 5}$	$1.149516430296_5^9 \times 10^{-2}$ $5.82580613431_7^8 \times 10^{-2}$ $1.588816609252_1^3 \times 10^{-1}$ $5.706958180199_2^9 \times 10^{-2}$ $1.293214169489_0^1 \times 10^{-1}$	$-2.6985586285405_5^5 \times 10^3$

a disjoint union of compact subsets of branch lengths from finitely many topologies, i.e, $\Theta^{(0)} = \cup_i \Theta^{(0, \tau_i)}$ and simultaneously prunes away sub-boxes from distinct $\Theta^{(0, \tau_i)}$ with a variant of the above algorithm that allows for compact sets contained in each $\Theta^{(0, \tau_i)}$ with its corresponding τ_i -specific post-order traversal to specify its topology-specific likelihood function. Thus, interval methods may be able to enclose the global maximum more efficiently when several topologies are considered simultaneously than when the global maximum is enclosed for each member of a finite set of topologies, one at a time, and finally compared.

Most statistical inference today is done on computing machines through numerical methods that do not rigorously account for the physical realities of such machines with finite memory. Possibly sub-optimal decisions may suffice for several decision problems. However, for others, such as, parameter estimation for nonlinear stochastic differential equations, or finding the equilibrium configuration of the n atoms in a folding protein molecule, or data fitting problems in training neural networks, one may want/have to guarantee the globally optimal decision. The enclosure methods provide statisticians with powerful tools for rigorous numerical inference.

7. Acknowledgements

The seeds for this work were sown in the inter-disciplinary environment of integrative graduate education and research traineeship program in complex non-linear systems funded by the NSF grant DGE-9870631 and completed with support from a joint NSF/NIGMS grant to Durrett, Aquadro, and Nielsen DMS 0201037. Many thanks to Dave Capella for insights into accelerations based on local search(es), Mike Steel for warning about the nonstationary optima at boundaries, Warwick Tucker for an introduction to interval analysis, Tandy Warnow for posing the problem, Wendy Shuk Wan Wong for sharing local search and data parsing codes, and Ziheng Yang for discussions about analytical results.

References

Alefeld, G. and J. Herzberger (1983). *An introduction to interval computations*. New York: Academic press.

- Berz, M. (1991). Forward algorithms for high orders and many variables with application to beam physics. See Griewank and Corliss (1991), pp. 147–156.
- Billera, L. J., S. Holmes, and K. Vogtmann (2004). Geometry of the space of phylogenetic trees. *Advances in Applied Math.*
- Brown, W. M., E. M. Prager, A. Wang, and A. C. Wilson (1982). Mitochondrial DNA sequences of primates, tempo and mode of evolution. *Jnl. Mol. Evol.* 18, 225–239.
- Chor, B. (2000). Multiple maxima of likelihood in phylogenetic trees: An analytic approach. *Mol. Biol. Evol.* 17, 1529–1541.
- Cuyt, A., B. Verdonk, S. Becuwe, and P. Kuterna (2001). A remarkable example of catastrophic cancellation unraveled. *Computing* 66, 309–320.
- Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood approach. *Jnl. Mol. Evol.* 17, 368–376.
- Felsenstein, J. (2003). *Inferring phylogenies*. Sunderland, MA: Sinauer Associates.
- Griewank, A. and G. Corliss (Eds.) (1991). *Automatic differentiation of algorithms: theory, implementation and applications*. Philadelphia: SIAM.
- Hammer, R., M. Hocks, U. Kulisch, and D. Ratz (1995). *C++ toolbox for verified computing: basic numerical problems*. Berlin: Springer-Verlag.
- Hansen, E. (1980). Global optimization using interval analysis - the multi-dimensional case. *Numerische Mathematik* 34, 247–270.
- Hansen, E. (1992). *Global optimization using interval analysis*. New York: Marcel Dekker.
- Hansen, E. and S. Sengupta (1981). Bounding solutions of systems of equations using interval analysis. *BIT* 21, 203–211.
- Jukes, T. H. and C. R. Cantor (1969). Evolution of protein molecules. In H. N. Munro (Ed.), *Mammalian protein metabolism*, Volume 3, pp. 21–123. New York: Academic Press.
- Kulisch, U. (2001). Advanced arithmetic for the digital computer, interval arithmetic revisited. See Kulisch et al. (2001), pp. 50–70.
- Kulisch, U., R. Lohner, and A. Facius (Eds.) (2001). *Perspectives on enclosure methods*. New York: Springer-Verlag.
- Loh, E. and G. W. Walster (2002). Rump’s example revisited. *Reliable Computing* 8, 245–248.
- Mayer, G. (1994). Result verification for eigenvectors and eigenvalues. In J. Herzberger (Ed.), *Topics in validated computations*, Volume 5 of *Studies in computational mathematics*, pp. 209–276. New York: North-Holland.
- Moore, R. E. (1967). *Interval analysis*. Englewood Cliffs, New Jersey: Prentice-Hall.
- Moore, R. E. (1979). *Methods and Applications of Interval analysis*. Philadelphia, Pennsylvania: SIAM.
- Neumaier, A. (1990). *Interval methods for systems of equations*. Cambridge: Cambridge university press.
- P754, I. T. (1985). *ANSI/IEEE 754-1985, Standard for Binary Floating-Point Arithmetic*. IEEE, New York. A preliminary draft was published in the January 1980 issue of IEEE Computer, together with several companion articles. Available from the IEEE Service Center, Piscataway, NJ, USA.
- Ratz, D. (1992). *Automatische ergebnisverifikation bei globalen optimierungsproblemen*. PhD dissertation, Universitat Karlsruhe, Karlsruhe.
- Yang, Z. (2000). Complexity of the simplest phylogenetic estimation problem. *Proceedings Royal Soc. London B Biol. Sci.* 267, 109–119.