

INVERTING RANDOM FUNCTIONS II: EXPLICIT BOUNDS FOR DISCRETE MAXIMUM LIKELIHOOD ESTIMATION, WITH APPLICATIONS*

MICHAEL A. STEEL[†] AND LÁSZLÓ A. SZÉKELY[‡]

Abstract. In this paper we study inverting random functions under the maximum likelihood estimation (MLE) criterion in the discrete setting. In particular, we consider how many independent evaluations of the random function at a particular element of the domain are needed for reliable reconstruction of that element. We provide explicit upper and lower bounds for MLE, both in the nonparametric and parametric setting, and give applications to coin-tossing and phylogenetic tree reconstruction.

Key words. random function, maximum likelihood estimation, Kullback–Leibler distance, phylogeny reconstruction

AMS subject classifications. 62B10, 62C05, 62F10, 92D15

PII. S089548010138790X

1. Review of random functions. This paper is a sequel of our earlier paper [12]. We assume that the reader is familiar with that paper; however, we repeat the most important definitions.

For two finite sets, A and U , let us be given a U -valued random variable ξ_a for every $a \in A$. We call the vector of random variables $(\xi_a : a \in A)$ a *random function* $\Xi : A \rightarrow U$. Ordinary functions are specific instances of random functions. It is easy to see [12] that an equivalent definition of random functions is obtained by picking one of the $|U|^{|A|}$ ordinary functions from A to U according to some distribution.

Given another random function, Γ , from U to V , we can speak about the composition of Γ and Ξ , $\Gamma \circ \Xi : A \rightarrow V$, which is the vector variable $(\gamma_{\xi_a} : a \in A)$. In this paper we are concerned with inverting random functions. In other words, we look for random functions $\Gamma : U \rightarrow A$ in order to obtain the best approximations of the identity function $\iota : A \rightarrow A$ by $\Gamma \circ \Xi$. *We always assume that Ξ and Γ are independent.* This assumption holds for free if either Ξ or Γ is a deterministic function.

Our motivation for the study of random functions came from phylogeny reconstruction. Stochastic models define how biomolecular sequences are generated at the leaves of a binary tree. If all possible binary trees on n leaves come equipped with a model for generating biomolecular sequences of length k , then we have a random function from the set of binary trees with n leaves to the ordered n -tuples of biomolecular sequences of length k . *Phylogeny reconstruction* is a random function from the set of ordered n -tuples of biomolecular sequences of length k to the set of binary trees with n leaves. It is a natural assumption that random mutations in the past are independent from any random choices in the phylogeny reconstruction algorithm. Criteria for phylogeny reconstruction may differ according to what one wishes to optimize.

*Received by the editors April 12, 2001; accepted for publication (in revised form) June 16, 2002; published electronically September 10, 2002.

<http://www.siam.org/journals/sidma/15-4/38790.html>

[†]Biomathematics Research Centre, University of Canterbury, Christchurch, New Zealand (m.steel@math.canterbury.ac.nz). This author was supported by the New Zealand Marsden Fund.

[‡]Department of Mathematics, University of South Carolina, Columbia, SC (szekely@math.sc.edu). This author was supported by NSF grants DMS 9701211 and 0072187, by a grant of SC CHE, and Hungarian NSF grant T 032455. This research was started when this author visited the University of Canterbury under the support of the New Zealand Marsden Fund.

Consider the probability of returning a from a by the composition of two random functions; that is, $r_a = \mathbb{P}[\gamma_{\xi_a} = a]$. The assumption on the independence of Ξ and Γ immediately implies

$$(1.1) \quad r_a = \sum_{u \in U} \mathbb{P}[\xi_a = u] \cdot \mathbb{P}[\gamma_u = a].$$

A natural criterion is to find Γ for a given Ξ in order to maximize $\sum_a r_a$. More generally, we may have a weight function $w : A \rightarrow \mathbb{R}^+$, and we may wish to maximize $\sum_a r_a w(a)$. This can happen if we give preference to returning certain a 's, or if we have a prior probability distribution on A and we want to maximize the expected return probability for a random element of A selected according to the prior distribution. A random function $\Gamma^* : U \rightarrow A$ can be defined in the following way: for any fixed $u \in U$,

$$(1.2) \quad \gamma_u^* = a^* \text{ for sure if } \forall a \in A, \mathbb{P}[\xi_{a^*} = u]w(a^*) \geq \mathbb{P}[\xi_a = u]w(a).$$

In case there is more than one element a^* that satisfies (1.2), we may select uniformly at random from the set of such elements. This function Γ^* is called the *maximum a posteriori estimator* (MAP) in the literature [7]. The special case when the weight function w is constant is known as the *maximum likelihood estimation* (MLE) [2, 7]. The MAP estimator Γ^* maximizes $\sum_a r_a w(a)$ for any given Ξ ; i.e., MAP is best on average. This result appears as Theorem 17.2 of [8], but an equivalent formulation, in the context of decision theory, is given by Theorem 10.3.1 of [2]; a further formulation, using a different proof, appears as Theorem 3.1 of [12]. However, it is at least as natural to look at a more conservative criterion: maximize the *smallest* value of r_a for $a \in A$; i.e., do the worst case the best. For this criterion MAP or MLE is, in general, not optimal. It is surprising, but little is known about the performance of MAP or MLE under this more conservative criterion.

Our paper [12] introduced a new abstract model for phylogeny reconstruction: inverting parametric random functions. Most of the work done on the mathematics of phylogeny reconstruction can be discussed in this context. This model is more structured than random functions, and hence is better suited to describe details of models of phylogeny and the evolution of biomolecular sequences. The approach is likely to be applicable in other areas where “nuisance” parameters are involved.

Assume that for a finite set A , for every $a \in A$, an (arbitrary, finite, or infinite) set $\Theta(a) \neq \emptyset$ is assigned, and, moreover, $\Theta(a) \cap \Theta(b) = \emptyset$ for $a \neq b$. Set $B = \{(a, \theta) : a \in A, \theta \in \Theta(a)\}$ and let π_1 denote the natural projection from B to A . A *parametric random function* is the collection Ξ of random variables such that

(i) for $a \in A$ and $\theta \in \Theta(a)$ there is a (unique) U -valued random variable $\xi_{(a,\theta)}$ in Ξ .

We are interested in random functions $\Gamma : U \rightarrow A$ independent from Ξ so that $\gamma_{\xi_{(a,\theta)}}$ best approximates π_1 under certain criteria. Call $R_{(a,\theta)}$ the probability $\mathbb{P}[\gamma_{\xi_{(a,\theta)}} = a]$. MLE, as it is used in the practice of phylogeny reconstruction, would take the Γ' , for which, for every fixed u , $\gamma'_u = a'$ for sure, if

$$(1.3) \quad \forall (a, \theta) \in B \quad \exists \theta' \in \Theta(a') \quad \mathbb{P}[\xi_{(a',\theta')} = u] \geq \mathbb{P}[\xi_{(a,\theta)} = u].$$

In case there is more than one element a' that satisfies (1.3), we may select uniformly at random from the set of such elements. (We avoided using the more natural looking quantification $\exists \theta' \in \Theta(a')$ for all $(a, \theta) \in B$, since $\mathbb{P}[\xi_{(a',\theta')} = u]$ may not take a

maximum value!) We denote by $R'_{(a,\theta)}$ the probability that from the pair (a, θ) the MLE Γ' returns a , i.e.,

$$(1.4) \quad R'_{(a,\theta)} = \mathbb{P}[\gamma'_{\xi_{(a,\theta)}} = a].$$

In [12] we made further assumptions on parametric random functions that we do not make in this paper:

(ii) There is a measure space $(\Theta(a), \mu_a(\cdot))$ defined on every $\Theta(a)$ such that $\mu_a(\Theta(a)) < \infty$.

(iii) For all $u \in U$, and for all $a \in A$, $\mathbb{P}[\xi_{(a,\theta)} = u] \in L^1(\Theta(a), \mu_a(\cdot))$.

Under these additional conditions we showed in [12] that in the model of parametric random functions the MLE criterion has to be modified to ensure the property that Γ' maximizes:

$$(1.5) \quad \sum_{a \in A} \int R_{(a,\theta)} d\mu_a(\theta).$$

This criterion is natural, since if $\sum_{a \in A} \int d\mu_a(\theta) = 1$, the formula (1.5) can be interpreted as the expected probability of return of elements of A , given a prior distribution on A .

The purpose of this paper is to place explicit upper and lower bounds on the probability that MLE correctly reconstructs elements of A , in both the parametric and nonparametric settings. Our primary interest is in the situation where k independent experiments are carried out, and we wish to determine how large k needs to be in order to correctly recover the underlying element of A with high probability. To emphasize the role of k we will let $[r^{(k)}]_a^*$ (resp., $[R^{(k)}]_{(a,\theta)}'$) denote the probability that MLE correctly reconstructs a in the nonparametric (resp., parametric) setting. We illustrate our bounds in the nonparametric setting by applications to coin-tossing and phylogeny reconstruction.

For the parametric setting, we first show, by way of an example, that the nonparametric upper bound on k does not extend in the way one might hope or expect. Nevertheless, we provide (in Theorem 5.1) an explicit upper bound on the number k of experiments required for MLE to reconstruct elements of A accurately. This result can be regarded as an extension of a discrete version of Wald's theorem [15]. We describe some implications of this result for phylogeny reconstruction in the remarks following Theorem 5.1.

Most of present paper can be considered as an attempt to analyze the worst case behavior of MLE. This is a very natural question in situations where a prior distribution is not given on A , or the inverting of the random function is to be carried out only once. Such a situation arises in phylogeny reconstruction, where we do not have a prior distribution on alternative evolutionary scenarios, and the reconstruction is not going to be repeated—there is only one “tree of life” that we want to know. However, the results in this paper are not restricted to the phylogeny setting and may be relevant to several other areas where MLE estimation is employed.

Our approach is information-theoretic; we focus on the possibility or impossibility of inverting random functions, and not on the computational complexity issues. Our results can also be restated in the language of decision theory, by talking about “loss functions” and “risk function” associated with the ML decision rule. Although some of our consequences or applications (described in section 4) may be derivable from existing theory, as far as we are aware the main results in this paper are not special cases of published results in either information theory or statistical decision theory.

2. Distances between distributions. For $a, b \in A, \Xi : A \rightarrow U$, let

$$(2.1) \quad d(a, b) = \sum_{u \in U} |\mathbb{P}[\xi_a = u] - \mathbb{P}[\xi_b = u]|.$$

We will refer to $d(a, b)$ as the *variational distance* of the random variables ξ_a and ξ_b . We also use the *Hellinger distance* of the random variables ξ_a and ξ_b , defined by

$$(2.2) \quad d_H(a, b) = \sqrt{\sum_{u \in U} \left(\sqrt{\mathbb{P}[\xi_a = u]} - \sqrt{\mathbb{P}[\xi_b = u]} \right)^2}.$$

These measures sometimes appear with slightly different definitions, terminology, and normalization constants. (For example, $\frac{1}{2}d(a, b)$ is sometimes referred to as the “variation distance.”) It is well known (see p. 25 in [9]) that $0 \leq d(a, b) \leq 2$ and

$$(2.3) \quad d_H^2(a, b) \leq d(a, b) \leq 2d_H(a, b).$$

We are going to use a well known and elegant multiplicative property of the Hellinger distance. For any $\Xi : A \rightarrow U$ random function define the $\Xi^{(k)} : A \rightarrow U^k$ random function as a sequence of k independent trials of Ξ . Let $d_H^{(k)}(a, b)$ denote the Hellinger distance of the random variables $\xi_a^{(k)}$ and $\xi_b^{(k)}$. Then independence immediately implies the identity

$$(2.4) \quad 1 - \frac{1}{2} \left(d_H^{(k)}(a, b) \right)^2 = \left(1 - \frac{1}{2} d_H^2(a, b) \right)^k,$$

by virtue of the formula

$$(2.5) \quad \sum_{u \in U} \left(\sqrt{\mathbb{P}[\xi_a = u]} - \sqrt{\mathbb{P}[\xi_b = u]} \right)^2 = 2 - 2 \sum_{u \in U} \sqrt{\mathbb{P}[\xi_a = u]} \sqrt{\mathbb{P}[\xi_b = u]}.$$

Combining the inequality $1 - (1 - x)^k \leq kx$ which holds for all $0 \leq x \leq 1$ and k positive integers, and (2.4), we obtain

$$(2.6) \quad \left(d_H^{(k)}(a, b) \right)^2 = 2 \left[1 - \left(1 - \frac{1}{2} d_H^2(a, b) \right)^k \right] \leq k d_H^2(a, b).$$

Using the notation $d^{(k)}(a, b)$ for the variational distance of the k independent trials, i.e., of the random variables $\xi_a^{(k)}$ and $\xi_b^{(k)}$, inequalities (2.3) and (2.6) imply

$$(2.7) \quad d^{(k)}(a, b) \leq 2\sqrt{k}d_H(a, b).$$

The nonsymmetric Kullback–Leibler distance (or relative entropy) of the random variables ξ_a and ξ_b is defined as

$$d_{KL}(a, b) = \sum_{u \in U} \mathbb{P}[\xi_a = u] \log \frac{\mathbb{P}[\xi_a = u]}{\mathbb{P}[\xi_b = u]}.$$

We will use the inequality [4]

$$(2.8) \quad d_{KL}(a, b) \geq \frac{1}{2}d^2(a, b).$$

3. MLE for inverting random functions. In this section we describe some lower and upper bounds on the probability that MLE correctly reconstructs elements of the set A . A classical upper bound on the average value of r_a over A —or more generally the value of $\sum_{a \in A} r_a w(a)$ for some probability distribution w on A —is given by “Fano’s inequality” (see, for example, [4]). Here we recall from [12] a different type of upper bound that applies also to r_a for any particular value of a and which is closely related to the variational distance.

THEOREM 3.1. *Assume that we have finite sets A and U and random functions $\Xi : A \rightarrow U$ and $\Gamma : U \rightarrow A$. Suppose that there is an element $b \in A$ and a subset N , $b \in N \subset A$ such that for all $a \in N$*

$$d(a, b) < \delta.$$

Then we have

$$\min_{a \in N} r_a \leq \frac{1}{|N|} + \delta \left(1 - \frac{1}{|N|} \right).$$

Now we can state the following lower bound for r_a in the setting of Theorem 3.1.

THEOREM 3.2. *Assume that we have finite sets A and U and a random function $\Xi : A \rightarrow U$. Assume that $\Gamma^* : U \rightarrow A$ is the MLE, and r_a^* is the return probability of $a \in A$ using Γ^* . Then we have*

$$(3.1) \quad r_a^* \geq 1 - \sum_{b \neq a} \left(1 - \frac{1}{2} d(a, b) \right).$$

If the MLE $\Gamma^* : U^k \rightarrow A$ is applied to invert the random function $\Xi^{(k)} : A \rightarrow U^k$, which is a sequence of k independent trials of Ξ , then

$$(3.2) \quad [r^{(k)}]_a^* \geq 1 - \sum_{b \neq a} \left(1 - \frac{1}{2} d_H^2(a, b) \right)^k.$$

Proof. For $y \in A$ let

$$U_y = \{u \in U \mid \forall x \in A, x \neq y, \mathbb{P}[\xi_y = u] > \mathbb{P}[\xi_x = u]\}$$

and similarly for V_y with “ \geq ” instead of “ $>$ ” in the definition. It is clear from independence (1.1) and the definition (1.2) that

$$(3.3) \quad r_a^* \geq \sum_{u \in U_a} \mathbb{P}[\xi_a = u].$$

For $x, y \in A$ set $p_y^x = \sum_{u \in V_y} \mathbb{P}[\xi_x = u]$. Now we claim

$$(3.4) \quad r_a^* \geq 1 - \sum_{y \neq a} p_y^a.$$

Note that

$$\sum_{y \neq a} p_y^a = \sum_{y \neq a} \sum_{u \in V_y} \mathbb{P}[\xi_a = u] \geq \mathbb{P}[\xi_a \notin U_a],$$

since the complement of U_a is a subset of $\cup_{y \neq a} V_a$, and

$$\mathbb{P}[\xi_a \notin U_a] = 1 - \mathbb{P}[\xi_a \in U_a] \geq 1 - r_a^*$$

by (3.3). This establishes (3.4). Finally, we have

$$\begin{aligned} d(a, y) &= \sum_{u \in U} |\mathbb{P}[\xi_a = u] - \mathbb{P}[\xi_y = u]| \\ &= \sum_{u \in V_y} (\mathbb{P}[\xi_y = u] - \mathbb{P}[\xi_a = u]) + \sum_{u \notin V_y} |\mathbb{P}[\xi_a = u] - \mathbb{P}[\xi_y = u]| \\ &\leq p_y^y - p_y^a + \sum_{u \notin V_y} (\mathbb{P}[\xi_a = u] + \mathbb{P}[\xi_y = u]) = p_y^y - p_y^a + (1 - p_y^a) + (1 - p_y^y) = 2 - 2p_y^a. \end{aligned}$$

Hence, $p_y^a \leq 1 - \frac{1}{2}d(a, y)$, and plugging this into (3.4) yields (3.1). To prove (3.2), apply (3.1) to $\Xi^{(k)}$ and invoke (2.4). \square

Remarks. First, note that (3.2) immediately implies that if $d_a = \min_{b \neq a} d_H(a, b)$, then $[r^{(k)}]_a^* > 1 - |A| \exp(-kd_a^2/2)$. Consequently, if

$$(3.5) \quad k > \frac{2}{d_a^2} \log \frac{|A|}{\epsilon},$$

then $[r^{(k)}]_a^* > 1 - \epsilon$. Second, note that an analogue of (3.2) also holds if, instead of k independent trials of Ξ , we take independent $A \rightarrow U$ random functions $\Xi_1, \Xi_2, \dots, \Xi_k$. Now the lower bound on $[r^{(k)}]_a^*$ is

$$1 - \sum_{b \neq a} \prod_{i=1}^k \left(1 - \frac{1}{2} d_H^2((\xi_i)_a, (\xi_i)_b)\right).$$

4. Applications.

4.1. Solving biased coin-tossing with MLE. We want to show an example where our upper and lower bounds for reconstructing random functions are nearly tight. Assume that $U = \{T, H\}$; i.e., we are tossing coins. Let a set A consist of $n + 1$ biased coins, denoted by $0, 1, 2, \dots, n$. Define the random function Ξ as follows: coin i shows H with probability i/n and shows T with probability $1 - i/n$. We show the following: there is a constant c_1 such that for $k = c_1 n^2$, for k independent trials of Ξ , $\Xi^{(k)}$, $[r^{(k)}]_i$ cannot be uniformly close to 1, no matter which method is used for inverting $\Xi^{(k)}$. However, there is a constant c_2 such that for $k = c_2 n^2$, using MLE, we find $[r^{(k)}]_i^*$ uniformly close to 1.

For simplicity we assume that n is odd. We are going to use Theorem 3.1 in the following setting: $b = \frac{n-1}{2}$, $N = \{\frac{n-3}{2}, \frac{n-1}{2}, \frac{n+1}{2}\}$. Then,

$$\min_{a \in N} [r^{(k)}]_a \leq \frac{1}{3} + \frac{2}{3}\delta,$$

where δ is the largest variational distance for $\Xi^{(k)}$ among b and the elements of N . Observe that for Ξ , by formula (2.5), we have

$$(4.1) \quad d_H^2(i, j) = 2 \left(1 - \frac{\sqrt{ij}}{n} - \frac{\sqrt{(n-i)(n-j)}}{n}\right).$$

It is easy to see that, for $i = b, j \in N$, (4.1) is maximized by $j_0 = \frac{n+1}{2}$ at the value $2(1 - \sqrt{1 - \frac{1}{n^2}}) \leq 2/n^2$. By (2.7), $d^{(k)}(b, x) \leq 2\sqrt{k}d_H(b, x)$, for every x , and therefore $\delta \leq 2\sqrt{k}d_H(b, j_0) \leq 4\sqrt{k}/n$. Any choice of $c_1 < 1/4$ suffices to keep either $r_{\frac{n-3}{2}}$ or $r_{\frac{n+1}{2}}$ separated from 1.

In the other direction we use Theorem 3.2. By (3.2) and (4.1) we have

$$(4.2) \quad [r^{(k)}]_i^* \geq 1 - \sum_{\substack{j=0 \\ j \neq i}}^n \left(1 - \frac{1}{2}d_H^2(i, j)\right)^k = 1 - \sum_{\substack{j=0 \\ j \neq i}}^n \left(\frac{\sqrt{ij}}{n} + \frac{\sqrt{(n-i)(n-j)}}{n}\right)^k.$$

By the classical inequality

$$\sqrt{a_1 a_2} + \sqrt{b_1 b_2} \leq \sqrt{a_1 + b_1} \cdot \sqrt{a_2 + b_2},$$

the generic subtracted term in the summation (4.2) is estimated from above by

$$\left(1 - \frac{(j-i)^2}{n^2}\right)^{k/2}.$$

Hence,

$$(4.3) \quad [r^{(k)}]_i^* \geq 1 - 2 \sum_{m=1}^n \left(1 - \frac{m^2}{n^2}\right)^{k/2}.$$

Now observe that

$$(4.4) \quad \sum_{m=1}^n \left(1 - \frac{m^2}{n^2}\right)^{k/2} \leq \left(1 - \frac{1}{n^2}\right)^{k/2} + n \int_{1/n}^1 (1 - x^2)^{k/2} dx$$

and

$$n \int_{1/n}^1 (1 - x^2)^{k/2} dx \leq n \int_{1/n}^1 e^{-k \cdot x^2/2} dx \leq \frac{n}{\sqrt{k}} \int_{\sqrt{k}/n}^{\sqrt{k}} e^{-t^2/2} dt \leq \frac{n}{\sqrt{k}} \cdot \frac{\sqrt{2\pi}}{2}.$$

Therefore, for a sufficiently large c_2 , selecting $k = c_2 n^2$, both terms in the right-hand side of (4.4) will be as small as wanted, and hence in (4.3) $[r^{(k)}]_i^*$ will be as close to 1 as wanted.

4.2. Phylogeny reconstruction. As a second application, we consider a problem arising in phylogenetic analysis. In this setting we have a model for generating sequences at the leaves of a tree, and the question is how long such sequences need to be in order to correctly reconstruct the tree—with high probability—from just the generated sequences.

The simplest stochastic model, for two-state sequences, is the symmetric model, due to Neyman [10], which we call the Neyman-2 model. (Related models also arise in statistical physics and in the theory of noisy communication—see, for example, [6].) Let $\{0, 1\}$ denote the two states. Let us be given a binary tree T (a tree in which each vertex has degree 1 or 3) with n labeled leaves. We describe how a single site in the sequence develops on T , and then we assume that the sites are independently and identically distributed (i.i.d.).

For each edge e of T we have an associated *transition probability*, which lies strictly between 0 and 0.5. Let $p : E(T) \rightarrow (0, 0.5)$ denote the associated map. Select one of

the leaves¹ and assign it state 0 or state 1 with probability 0.5. Direct all edges away from this leaf and recursively assign random states to the vertices of T as follows: if $e = \{u, v\}$ is directed from u to v , and u (but not v) has a state assignment, then v is assigned the same state as u with probability $1 - p_e$ or the other state with probability p_e . (In this latter case, we say there is a *transition on e* .) It is assumed that all assignments are made independently, and so the pair (T, p) determines the joint probability of any assignment of states to the vertices of T and thereby the marginal probability of any assignment of states to the leaves of T . If we independently generate k such assignments of states to the leaves of T , we obtain n sequences of length k . For this model, upper bounds on the sequence length k required to reconstruct the underlying tree were given in [5, 12]. These papers showed that, for accurate tree reconstruction, k needs to grow only quadratically in $1/f$, where f is the smallest transition probability in the tree, when other parameters are fixed. We now show that this rate of growth is not only sufficient but is also necessary.

Consider binary trees having four labeled leaves and two unlabeled interior vertices. There are three such trees (up to equivalence), and we will denote them as a, b, c . Each tree has four leaf edges (an edge incident to a leaf) and one interior edge. Take $A = \{a, b, c\}$, and let U be the set of binary functions defined on the four leaves. Assume that a, b, c are Neyman-2 trees with transition probability f on the interior edge. (We do not care what the other transition probabilities are.) Let Ξ denote the state assignment of the leaves of a, b, c under the Neyman-2 model.

THEOREM 4.1. *For the three Neyman-2 binary trees a, b, c on four leaves (as described above), and state assignment Ξ , under any method for inverting random function Ξ from k independent trials (i.e., from binary sequences of length k associated with the leaves) with success probability near 1 for all three trees, $k = \Omega(\frac{1}{f^2})$.*

Proof. We are going to prove that for f sufficiently close to 0, for some constant $C > 0$,

$$(4.5) \quad d_H(a, b) \leq Cf.$$

Now (2.7) and (4.5) imply $d^{(k)}(a, b) \leq 2Cf\sqrt{k}$, and one similarly obtains $d^{(k)}(c, b) \leq 2Cf\sqrt{k}$. So if we apply Theorem 3.1 with $N = \{a, b, c\}$,

$$(4.6) \quad \min\{r_a, r_c\} \leq \frac{1}{3} + \frac{4}{3}Cf\sqrt{k},$$

and the right-hand side of (4.6) is well separated from 1 as k is a small constant over f^2 .

To complete the proof, we have to verify (4.5). Assume that a is the tree in which the interior edge separates leaves 1, 2 from leaves 3, 4; and b is the tree in which the interior edge separates leaves 1, 3 from leaves 2, 4. By (2.2)

$$(4.7) \quad d_H(a, b)^2 = \sum_{u \in U} \left(\sqrt{\mathbb{P}[\xi_a = u]} - \sqrt{\mathbb{P}[\xi_b = u]} \right)^2,$$

where the summation goes for 16 terms which correspond to the 16 elements of U : functions with domain $\{1, 2, 3, 4\}$ and range $\{0, 1\}$. We are going to condition on the

¹One assumes that mutations of an ancestral sequence happen in this way. However, it can be shown that the selection of the special leaf has no effect on the reconstruction in the model considered here.

event Φ , denoting that there is transition on the interior edge of the tree and also for the complement of this event. For $x = a, b$ define

$$\begin{aligned} A(x, u) &= \mathbb{P}[\xi_x = u \mid \neg\Phi], \\ B(x, u) &= \mathbb{P}[\xi_x = u \mid \Phi] - \mathbb{P}[\xi_x = u \mid \neg\Phi]. \end{aligned}$$

Notice that $A(x, u)$ and $B(x, u)$ are constants that do not depend on f . Also, observe that

$$\begin{aligned} \mathbb{P}[\xi_x = u] &= \mathbb{P}[\xi_x = u \mid \neg\Phi] \cdot (1 - f) + \mathbb{P}[\xi_x = u \mid \Phi] \cdot f \\ &= \mathbb{P}[\xi_x = u \mid \neg\Phi] + f \cdot (\mathbb{P}[\xi_x = u \mid \Phi] - \mathbb{P}[\xi_x = u \mid \neg\Phi]) \\ &= A(x, u) + fB(x, u). \end{aligned}$$

It easily follows from the geometry of the trees a and b that $A(a, u) = A(b, u)$. Furthermore, it is easily seen that $A(a, u) \neq 0$ for all values of u , which ensures (below) that we may divide expressions by $A(a, u)$. Hence, by the Taylor expansion of the square root function, we have

$$\begin{aligned} \sqrt{\mathbb{P}[\xi_a = u]} - \sqrt{\mathbb{P}[\xi_b = u]} &= \sqrt{A(a, u)} \left(\sqrt{1 + \frac{fB(a, u)}{A(a, u)}} - \sqrt{1 + \frac{fB(b, u)}{A(a, u)}} \right) \\ (4.8) \qquad \qquad \qquad &= f \frac{B(a, u) - B(b, u)}{2\sqrt{A(a, u)}} + O(f^2), \end{aligned}$$

and summing up 16 terms like (4.8) we obtain

$$d_H^2(a, b) = f^2 \sum_{u \in U} \frac{(B(a, u) - B(b, u))^2}{4A(a, u)} + O(f^3),$$

and this proves (4.5) for all

$$C > \sqrt{\sum_{u \in U} \frac{(B(a, u) - B(b, u))^2}{4A(a, u)}}. \quad \square$$

5. MLE for inverting parametric random functions. We start with an example showing that for parametric MLE there is no counterpart of (3.2); that is, there is no function $f = f(\delta, k)$ such that, for all $\delta > 0$, $\lim_{k \rightarrow \infty} f(\delta, k) = 0$ and

$$(5.1) \qquad [R^{(k)}]_{(a, \theta)}' \geq 1 - \sum_{b \neq a} f(\delta((a, \theta), b), k),$$

where

$$\delta((a, \theta), b) = \inf_{\theta' \in \Theta(b)} d_H((a, \theta), (b, \theta')).$$

Take $A = \{a_1, a_2\}$, $U = \{u_1, u_2, \dots, u_{2k^2}\}$, $\Theta(a_1) = \Theta(a_2) = U^k$. We denote a generic element of U^k by \mathbf{u} , and $\text{supp}(\mathbf{u})$ denotes the set of elements of U which occur as coordinates in \mathbf{u} . Let $B = (\{a_1\} \times \Theta(a_1)) \cup (\{a_2\} \times \Theta(a_2))$. Define the parametric random function $\Xi : B \rightarrow U$ as follows. Set $\mathbb{P}[\xi_{(a_1, \mathbf{u})} = v] = 1/|U|$ for each $v \in U$. For

$\mathbf{u} \in U^k$ and $v \in U$, set $\mathbb{P}[\xi_{(a_2, \mathbf{u})} = v] = i/k$ if v occurs at $i = i(v)$ coordinates in \mathbf{u} . Now for any $\mathbf{w}, \mathbf{u} \in U^k$ we have

$$(5.2) \quad d((a_1, \mathbf{w}), (a_2, \mathbf{u})) \geq 2 - \frac{1}{k}$$

by the calculation

$$\sum_{v \in \text{supp}(\mathbf{u})} \left(\frac{i(v)}{k} - \frac{1}{|U|} \right) + \sum_{v \notin \text{supp}(\mathbf{u})} \frac{1}{|U|} = 2 - 2 \sum_{v \in \text{supp}(\mathbf{u})} \frac{1}{|U|} \geq 2 - \frac{2k}{|U|} = 2 - \frac{1}{k}.$$

Now consider k independent trials of $\Xi, \Xi^{(k)}$. We study inverting $\Xi^{(k)}$ with parametric MLE. Note that, for any $\mathbf{u} \in U^k$,

$$\mathbb{P}[\xi_{(a_2, \mathbf{u})} = \mathbf{u}] = \prod_{i=1}^k \mathbb{P}[\xi_{(a_2, \mathbf{u})} = u_i] \geq \left(\frac{1}{k} \right)^k ;$$

and for any $\mathbf{w} \in U^k$,

$$\mathbb{P}[\xi_{(a_1, \mathbf{w})} = \mathbf{u}] = \prod_{i=1}^k \mathbb{P}[\xi_{(a_1, \mathbf{w})} = u_i] = \left(\frac{1}{2k^2} \right)^k < \left(\frac{1}{k} \right)^k .$$

Therefore, one *always* has $[R^{(k)}]_{(a_1, \mathbf{w})}' = 0$ (see (1.4)), while by (5.2) and (2.3) the d_H distances between the random variables corresponding to a_1 and a_2 are well separated from zero. This establishes our claim at the start of this section regarding the nonexistence of an analogue of (3.2) from Theorem 3.2.

Intuitively, the reason this construction works is that we have selected range and parameter spaces whose size *depends* on the sequence length k . Note that we could have allowed $|U|$ to grow just linearly with k and still obtained the same conclusion. However, by allowing $|U|$ to grow more quickly with k our construction has a further notable property. Namely, the random variables corresponding to a_1 and a_2 become maximally distant under variation distance as $k \rightarrow \infty$, as inequality (5.2) reveals.

However, with mild extra conditions we can state a positive result. This positive result provides explicit bounds on the convergence of the MLE in the parametric setting.

THEOREM 5.1. *Assume $B = \{(a, \theta) : a \in A, \theta \in \Theta(a)\}$, and $\Xi : B \rightarrow U$ is a parametric random function, where A and U are finite sets. Assume that for a particular $(a, \theta) \in B$ there exists a $d_0 > 0$ such that for all $b \in A, b \neq a$, and $\theta' \in \Theta(b)$*

$$(5.3) \quad d((a, \theta), (b, \theta')) \geq d_0,$$

where d , as usual, denotes the variational distance. If the MLE is applied to invert the parametric random function $\Xi^{(k)} : A \rightarrow U^k$, which is a sequence of k independent trials of Ξ , then

$$(5.4) \quad \lim_{k \rightarrow \infty} [R^{(k)}]_{(a, \theta)}' = 1.$$

For a more precise result, set $U^+ = \{u \in U : \mathbb{P}[\xi_{(a, \theta)} = u] > 0\}$, and $\alpha = \min_{u \in U^+} \mathbb{P}[\xi_{(a, \theta)} = u]$. If

$$(5.5) \quad k > f(\alpha, d_0) \log \left(\frac{2|U^+|}{\epsilon} \right),$$

then MLE estimation returns a with probability at least $1 - \epsilon$, where

$$f(\alpha, d_0) = \max \left\{ \frac{16}{\alpha}, \frac{17 \log^2 \alpha (1 + \frac{2}{\alpha})^2}{\alpha d_0^4} \right\}.$$

Proof. For $u \in U$, define $p(u) = \mathbb{P}[\xi_{(a,\theta)} = u]$, and then $\alpha = \min_{u \in U^+} \{p(u)\} > 0$. Define $\hat{p}(u)$ as the corresponding relative frequency, i.e.,

$$(5.6) \quad \hat{p}(u) = \frac{1}{k} \#\{j : (\xi_j)_{(a,\theta)} = u\},$$

where ξ_j is the j th trial of the random function. Let $\delta = \frac{4}{\sqrt{17}}$, and let

$$\eta = \min \left\{ \frac{1}{2}, \frac{\delta d_0^2}{2|\log \alpha|(1 + \frac{2}{\alpha})} \right\}.$$

Then,

$$(5.7) \quad \eta |\log \alpha| + \frac{\eta |\log \alpha|}{\alpha(1 - \eta)} \leq \frac{\delta}{2} d_0^2.$$

By the large deviation inequality given in formula (14) of Appendix A in [1], we have

$$(5.8) \quad \mathbb{P}[|p(u) - \hat{p}(u)| > \eta p(u)] < 2e^{-c_\eta k p(u)},$$

where $c_\eta = \min\{-\log [e^\eta(1 + \eta)^{-(1+\eta)}], \frac{\eta^2}{2}\}$. Note that for $0 < \eta < 1/2$ we have $-\log[e^\eta(1 + \eta)^{-(1+\eta)}] \geq \frac{\eta^2(1-\eta)}{2}$ by Taylor expansion, and hence $c_\eta \geq \eta^2/4$. Therefore, formula (5.8) holds if we change c_η to $\eta^2/4$ in the exponent. Now suppose k satisfies inequality (5.5). Then,

$$k > \frac{4}{\alpha \eta^2} \log \left(\frac{2|U^+|}{\epsilon} \right)$$

by the definition of f and η . Consequently, $2|U^+|e^{-\eta^2 k \alpha/4} < \epsilon$, and so, with probability at least $1 - \epsilon$, we have

$$(5.9) \quad \forall u \in U \quad |p(u) - \hat{p}(u)| \leq \eta p(u).$$

(We also used the Bonferroni inequality, and the fact that, with probability 1, $p(u) = \hat{p}(u) = 0$ for all $u \in U \setminus U^+$.) For $x \in A, \omega \in \Theta(x)$, consider

$$(5.10) \quad L(x, \omega) = \sum_{u \in U} \hat{p}(u) \log \mathbb{P}[\xi_{x,\omega} = u].$$

(Here, as always in this kind of calculation, we use the convention $0 \times (-\infty) = 0$, which is supported by $\lim_{x \rightarrow 0^+} x \log x = 0$.) $L(x, \omega)$ is $\frac{1}{k}$ times the natural logarithm of the probability that the observed sequence of U -elements came from (x, ω) . Therefore $L(x, \omega) \leq 0$ is proportional to the log-likelihood of (x, ω) .

Now consider a fixed $b \in A, b \neq a$ and a fixed $\theta' \in \Theta(b)$. For $u \in U$, we use the notation $q(u) = \mathbb{P}[\xi_{(b,\theta')} = u]$.

We finish the proof conditional on the following event:

$$(5.11) \quad [(5.9) \text{ holds}] \text{ and } [u \notin U^+ \text{ implies } \hat{p}(u) = 0].$$

Note that the second part of the condition holds with probability 1, and so event (5.11) occurs with probability at least $1 - \epsilon$.

We distinguish two cases. In both cases we show

$$(5.12) \quad L(a, \theta) - L(b, \theta') > 0.$$

Since $L(a, \theta)$ (resp., $L(b, \theta')$) is the log-likelihood of getting the observed sequence from (a, θ) (resp., (b, θ')), (5.12) implies the correct reconstruction of a from the observed data by MLE by (1.4). Since this holds (with probability 1) for all θ' , conditional on event (5.11), and event (5.11) occurs with probability at least $1 - \epsilon$, the probability that MLE correctly reconstructs a will be at least $1 - \epsilon$, as required.

Case 1. There exists a $v \in U^+$ with $q(v) < \exp(\frac{\log \alpha}{\alpha(1-\eta)})$. In this case $L(b, \theta') \leq \hat{p}(v) \log q(v) < \log \alpha$, so $L(b, \theta') < \log \alpha$. On the other hand,

$$L(a, \theta) = \sum_{u \in U} \hat{p}(u) \log p(u) \geq \sum_{u \in U} \hat{p}(u) \log \alpha = \log \alpha.$$

Therefore, $L(a, \theta) > L(b, \theta')$.

Case 2. For all $u \in U^+$, $q(u) \geq \exp(\frac{\log \alpha}{\alpha(1-\eta)})$. We have, for all $u \in U^+$, $|\log q(u)| \leq \frac{|\log \alpha|}{\alpha(1-\eta)}$. Consider

$$(5.13) \quad \begin{aligned} L(a, \theta) - L(b, \theta') &= \sum_{u \in U} \hat{p}(u) \log \frac{p(u)}{q(u)} = \sum_{u \in U^+} \hat{p}(u) \log \frac{p(u)}{q(u)} \\ &= \sum_{u \in U} p(u) \log \frac{p(u)}{q(u)} + \sum_{u \in U^+} (\hat{p}(u) - p(u)) \log \frac{p(u)}{q(u)}. \end{aligned}$$

Notice that the first sum in (5.13) is exactly the Kullback–Leibler distance $d_{KL}((a, \theta), (b, \theta'))$. By formulae (2.8, 5.3) this first sum is at least $\frac{1}{2}d_0^2$. Since we condition on (5.9), $|\hat{p}(u) - p(u)| \leq \eta p(u)$. Hence, we can estimate the absolute value of the second sum in (5.13) by

$$(5.14) \quad \begin{aligned} \sum_{u \in U} \eta p(u) (|\log p(u)| + |\log q(u)|) &\leq \sum_{u \in U} \eta p(u) \left(|\log \alpha| + \frac{|\log \alpha|}{\alpha(1-\eta)} \right) \\ &= \eta |\log \alpha| + \frac{\eta |\log \alpha|}{\alpha(1-\eta)} \leq \frac{\delta}{2} d_0^2 \end{aligned}$$

by (5.7), and so $L(a, \theta) - L(b, \theta') > 0$. \square

Remarks.

1. Notice that, because $|U^+|\alpha \leq 1$, inequality (5.5) will hold whenever $k \geq f(\alpha, d_0) \log(\frac{2}{\alpha\epsilon})$. Notice that this bound on k (that suffices for parametric MLE to reconstruct a with probability at least $1 - \epsilon$) depends only on ϵ , d_0 , and α , and it is independent of the cardinality of A and U (cf. the bound we described for nonparametric MLE in the remark following Theorem 3.2).
2. Note also that the example described at the beginning of section 5 shows that one cannot strengthen Theorem 5.1 by simply dropping the role of α . That is, Theorem 5.1 fails if we replace (5.5) with the weaker condition that

$$k \geq f_1(d_0) \log \left(\frac{2|U^+|}{\epsilon} \right)$$

for some suitable function f_1 (that does not depend on α), since in the example described any such inequality will be satisfied for sufficiently large k ($|U|$ grows only quadratically with k), yet MLE fails to recover a_1 . A closer examination of this example shows that α converges to zero sufficiently fast with k for the bound in (5.5) to be violated.

3. Suppose that for each $b \in A$ we have (i) the set $\Theta(b)$ is a compact topological space, and (ii) the mapping from $\Theta(b)$ to the interval $[0, 1]$ defined by $(b, \theta) \mapsto \mathbb{P}(\xi_{(b,\theta)} = u)$ is continuous for each element $u \in U$. Then the separation property (5.3) required in Theorem 5.1 becomes equivalent to the (in general weaker) condition that for all $b \in A$, $b \neq a$, and $\theta' \in \Theta(b)$

$$(5.15) \quad d((a, \theta), (b, \theta')) > 0.$$

For example, for most models in the phylogenetic setting, assumptions (i) and (ii) will apply, and so MLE will be statistically consistent (that is, satisfy (5.4)), provided the model satisfies (5.15). In particular, the detailed analysis and additional assumption required by Chang [3] in order to establish (for a general Markov model on trees) a strengthening of (5.15) to the case $b = a$, $\theta \neq \theta'$ is unnecessary if one wishes simply to establish the statistical consistency of MLE in the estimation of a binary tree (and not the associated transition matrices of the model). There are also other models in use that satisfy (5.15) and thereby justify the statistical consistency of MLE. For example, consider a model in which sites evolve i.i.d. on a binary tree according to a stationary, reversible Markov process (with an unknown rate matrix) and with a rate factor (constant across the tree) drawn from a distribution \mathcal{D} . Such models satisfy (5.15) if \mathcal{D} is known and therefore the same for each possible tree [14, section 3.3]; however, (5.15) may fail if \mathcal{D} is unknown [13]. We note that Theorem 5.1 also provides the first explicit upper bounds on the sequence length required for MLE to accurately reconstruct a binary tree in the phylogenetic setting.

6. Conclusion and open problems. It would be interesting to see how much the bound on k given by Theorem 5.1 might be improved. This question applies both for the general setting in which Theorem 5.1 is stated, and also for more particular settings, such as arises in phylogeny.

In the general setting, observe that our upper bound on k given by Theorem 5.1 grows at the rate d_0^{-4} . Yet, in the nonparametric setting, if we let $d_0 = \min\{d(a, b) : a, b \in A, a \neq b\}$, then the analogous upper bound on k grows at the rate d_0^{-2} (by inequalities (3.5) and (2.3)). An interesting question is whether this discrepancy is essential in moving from the nonparametric to the parametric setting, or whether it can be avoided by a different argument. There are other significant differences between our results for the nonparametric and parametric setting—for example, although $|A|$ (but not $|U|$) enters directly into our bound on k in the nonparametric setting (inequality (3.5)), in the parametric setting $|A|$ is not directly mentioned but $|U^+|$ is.

Regarding more particular parametric MLE settings (such as in phylogeny) it is quite likely that the additional structure present in these instances may yield tighter bounds than those given by Theorem 5.1. It would be particularly desirable to set matching lower and upper bounds on the sequence length (the number of samples k) required by MLE in phylogeny reconstruction. It is clear that, for certain choices of the parameter θ , MLE may require longer sequences than other methods to correctly

reconstruct a phylogenetic tree (as discussed in [11] and [12]). Indeed, the statistical consistency of MLE in phylogeny was established only in 1996 by [3] in a result that, like Wald's earlier result [15], is based on a compactness argument that does not give an explicit bound on k . The significance of Theorem 5.1 is that it gives the first such explicit bounds for MLE, both in the phylogenetic setting and beyond.

6.1. Correction. Theorem 2.3 in [12], cited as Theorem 3.1 in our current paper, tacitly assumed $b \in N$. This assumption has to be made explicit.

Acknowledgments. The authors are indebted to Éva Czabarka for her invaluable comments on the manuscript. MAS also thanks Joe Chang for suggesting the usefulness of the Hellinger distance in the proof of Theorem 4.1. The authors are indebted for a number of useful comments from two anonymous referees.

REFERENCES

- [1] N. ALON AND J. H. SPENCER, *The Probabilistic Method*, John Wiley and Sons, New York, 1992.
- [2] G. CASELLA AND R. L. BERGER, *Statistical Inference*, Duxbury Press, Belmont, CA, 1990.
- [3] J. T. CHANG, *Full reconstruction of Markov models on evolutionary trees: Identifiability and consistency*, *Math. Biosci.*, 137 (1996), pp. 51–73.
- [4] T. M. COVER AND J. A. THOMAS, *Elements of Information Theory*, John Wiley and Sons, New York, 1991.
- [5] P. L. ERDŐS, M. A. STEEL, L. A. SZÉKELY, AND T. WARNOW, *A few logs suffice to build (almost) all trees (I)*, *Random Structures Algorithms*, 14 (1999), pp. 153–184.
- [6] W. EVANS, C. KENYON, Y. PERES, AND L. J. SCHULMAN, *Broadcasting on trees and the Ising model*, *Ann. Appl. Probab.*, 10 (2000), pp. 410–433.
- [7] B. S. EVERITT, *The Cambridge Dictionary of Statistics*, Cambridge University Press, Cambridge, UK, 1998.
- [8] S. GUIASU, *Information Theory with Applications*, McGraw-Hill, New York, 1977.
- [9] L. LE CAM AND G. L. YANG, *Asymptotics in Statistics: Some Basic Concepts*, Springer-Verlag, New York, 1990.
- [10] J. NEYMAN, *Molecular studies of evolution: A source of novel statistical problems*, in *Statistical Decision Theory and Related Topics*, S. S. Gupta and J. Yackel, eds., Academic Press, New York, 1971, pp. 1–27.
- [11] M. STEEL AND D. PENNY, *Parsimony, likelihood, and the role of models in molecular phylogenetics*, *Mol. Biol. Evol.*, 17 (2000), pp. 839–850.
- [12] M. A. STEEL AND L. A. SZÉKELY, *Inverting random functions*, *Ann. Comb.*, 3 (1999), pp. 103–113.
- [13] M. A. STEEL, L. A. SZÉKELY, AND M. D. HENDY, *Reconstructing trees when sequence sites evolve at variable rates*, *J. Comput. Biol.*, 1 (1994), pp. 153–163.
- [14] C. TUFFLEY AND M. STEEL, *Modelling the covarion hypothesis of nucleotide substitution*, *Math. Biosci.*, 147 (1998), pp. 63–91.
- [15] A. WALD, *Note on the consistency of the maximum likelihood estimate*, *Ann. Math. Statistics*, 20 (1949), pp. 595–600.