

# STATISTICALLY-CONSISTENT K-MER METHODS FOR PHYLOGENETIC TREE RECONSTRUCTION

ELIZABETH S. ALLMAN, JOHN A. RHODES, AND SETH SULLIVANT

ABSTRACT. Frequencies of  $k$ -mers in sequences are sometimes used as a basis for inferring phylogenetic trees without first obtaining a multiple sequence alignment. We show that a standard approach of using the squared-Euclidean distance between  $k$ -mer vectors to approximate a tree metric can be statistically inconsistent. To remedy this, we derive model-based distance corrections for orthologous sequences without gaps, which lead to consistent tree inference. The identifiability of model parameters from  $k$ -mer frequencies is also studied. Finally, we report simulations showing the corrected distance out-performs many other  $k$ -mer methods, even when sequences are generated with an insertion and deletion process. These results have implications for multiple sequence alignment as well, since  $k$ -mer methods are usually the first step in constructing a guide tree for such algorithms.

## 1. INTRODUCTION

The first step in most approaches to inference of a phylogenetic tree from sequence data is to construct an alignment of the sequences, intended to identify orthologous sites. When many sequences are considered at once, a full search over all possible sequence alignments is infeasible, so most algorithms reduce the range of possible alignments considered by constructing multiple alignments on subcollections of the sequences and then merging these together, using heuristic, rather than model-based, schemes. Deciding which subcollections of the sequences to align follows a *guide tree*, a rough tree approximating the evolutionary histories of all the sequences. This means that sequence alignment and phylogenetic tree construction are circularly entangled: finding a tree depends on knowing a multiple sequence alignment, and obtaining a sequence alignment requires knowing a tree.

To get around this “chicken-and-egg” problem of alignment and phylogeny several methods have been proposed. The most theoretically appealing methods are simultaneous alignment and phylogeny algorithms, built upon statistical models of insertion and deletions (indels) of bases as well as base substitutions [Thorne et al., 1991, 1992]. Unfortunately, such methods are computationally intensive, and do not scale well for large phylogenies. Alternatively, methods have been developed that iteratively compute alignments and phylogenies many times, using the output from one procedure as the input to the next [Liu et al., 2009, 2012]. These last investigations underscored that poor alignments can be a significant source of error in trees, and that better guide trees can lead to better tree inference.

If one is interested primarily in the phylogeny, an alternate strategy is to develop methods for inferring trees that do not require having a sequence alignment in hand. Current fully alignment-free phylogenetic methods were not developed with stochastic models of sequence evolution in mind, and are not widely accepted in the phylogenetics community. However, the construction of initial guide-trees for producing alignments generally follows an alignment-free approach. For example MUSCLE [Edgar, 2004a,b] uses  $k$ -mer distances with UPGMA or Neighbor-Joining to produce guide trees, whereas Clustal Omega [Sievers et al., 2011] uses a low dimensional geometric embedding based on  $k$ -mers [Blackshields et al., 2010] and  $k$ -means or UPGMA as a clustering algorithm. Thus even though tree inference is typically performed with model-based statistical methods, the initial step is built on heuristic ideas, with no evolutionary model in use.

As exemplified by these alignment algorithms, most common alignment-free methods are based on  $k$ -mers, contiguous subsequences of length  $k$ . To a sequence of length  $n$  for any natural number  $k \leq n$  we associate the vector of counts of its distinct  $k$ -mers. For a DNA sequence the  $k$ -mer count vector has  $4^k$  entries and sums to  $n - k + 1$ . Distance between two sequences might be calculated by measuring the (squared Euclidean) distance between their (suitably normalized)  $k$ -mer count vectors. In this way one obtains pairwise distances between all sequences, and can apply a standard distance-based method (e.g. Neighbor joining) to construct a phylogenetic tree.

Such  $k$ -mer methods are sometimes described as non-parametric, in that they do not depend on any underlying statistical model describing the generation of the sequences. For phylogenetic purposes, where an evolutionary model will be assumed in later stages of an analysis, it is hard to view this as desirable. As we will show in Section 3, if we do assume that data is produced according to a standard probabilistic model of sequence evolution, then a naive  $k$ -mer method is statistically inconsistent. That is, over a rather large range of metric trees, it will not recover the correct tree from sequence data, even with arbitrarily long sequences. The statistical inconsistency of such a  $k$ -mer method is similar to the ones seen for parsimony, in the “Felsenstein zone” [Felsenstein, 1978].

Our main result, presented in Section 2, is the derivation of a statistically consistent model-based  $k$ -mer distance under standard phylogenetic models with no indel process. It would, of course, be preferable to work with a model including indels, as only in that situation is an alignment-free method of real value. At this time, however, we are only able to offer a reasonable heuristic extension of our method for sequences evolving with a mild indel process. This appears in Section 5. We view this as only a first step towards developing rigorously-justified model-based  $k$ -mer methods for indel models; solid theoretical development of such methods is a project for the future.

Section 4 presents more detailed results on identifiability of model parameters from  $k$ -mer count vectors. While one of these plays a role in establishing the results of Section 2, they are of interest in their own right. Technical proofs for Sections 2 and 4 are deferred to the Appendices.

In Section 6 we report results from simulation studies on sequence data generated from models with and without an indel process, comparing  $k$ -mer methods with and without the model-based corrections. As expected, the  $k$ -mer methods with the model-based

corrections outperform both the uncorrected  $k$ -mer methods and a more traditional distance method based on first computing pairwise alignments of sequences. The simulation studies also illustrate the statistical consistency of the model-based methods, and the inconsistency of the standard  $k$ -mer method.

**Comparison to Prior Work on Alignment-Free Phylogenetic Algorithms.** There have been a number of papers in recent years developing alignment-free methods for phylogenetic tree reconstruction [Daskalakis and Roch, 2013, Reyes-Prieto et al., 2011, Yang and Zhang, 2008, Chan et al., 2014] or for clustering metagenomic data [Reinert et al., 2009, Shen et al., 2014]. Of these only one [Daskalakis and Roch, 2013] appears to be based on common phylogenetic modeling assumptions, but its focus is theory rather than practice. Others [Chan et al., 2014, Reinert et al., 2009] are model-based but the underlying model is not evolutionary in nature. Some are primarily simulations studies of the application of a method on larger trees than those we focus on here.

In our simulations, we follow the framework suggested by Huelsenbeck [1995], which allows us to graphically display performance on an important slice of tree space for 4-taxon trees. One then readily sees the effect on performance of varying branch length, and the strength of the common “long branch attraction” phenomenon. In comparison, the simulations in [Reyes-Prieto et al., 2011, Yang and Zhang, 2008, Chan et al., 2014] use trees that have more leaves but the range of branch lengths explored is significantly reduced. We believe following Huelsenbeck’s plan provides more fundamental insights into a methodology’s value.

Daskalakis and Roch [2013] derived a statistically consistent alignment-free method for a model with indels, although it appears to have not yet been tested, even on simulated data. Their method is based on computing the base distribution (i.e., the 1-mer distribution) in sub-blocks of the sequences, and motivated the similar approach we take here. In addition to restricting to 1-mers, their approach requires *a priori* knowledge of the value of certain model parameters, e.g., the proportion of gaps in a sequence, and several parameters defining the base substitution process. As our theoretical results involve no indel process and allow arbitrary  $k$ , the two works are not directly comparable. However, we are able to obtain stronger results on the identifiability of parameters of the base substitution model, and our simulations show that using  $k > 1$  can result in improved performance.

For advancing data analysis, it is highly desirable to develop theoretically-justified model-based  $k$ -mer methods that both account for indels and require few assumptions on model parameters. Neither Daskalakis and Roch [2013] nor we provide such methods; both of our works represent first steps, in slightly different directions, but pointing towards the same goal.

## 2. $k$ -MER FORMULAS FOR INDEL-FREE SEQUENCES

In this section we present formulas for model-based corrections to distances based on  $k$ -mer frequency counts. Technical proofs appear in Appendix A. Our main result, Theorem 2.1, is quite general, applying to arbitrary pairwise distributions that are at stationarity. We use this result to derive corrected distance calculations for the Juke-Cantor model

and the Kimura 2- and 3-parameter models. These corrections yield statistically consistent estimates of evolutionary times between extant taxa. Coupled with a statistically consistent method for constructing a tree from distances (for example, Neighbor Joining [Saitou and Nei, 1987]), this produces a statistically consistent method for reconstructing phylogenetic trees from  $k$ -mer counts.

Let  $S$  be a sequence on an  $L$ -letter alphabet,  $[L] := \{1, 2, \dots, L\}$ . For a natural number  $k$ , let  $X$  denote the vector of  $k$ -mer counts extracted from  $S$ . That is, for each  $W = w_1 w_2 \dots w_k \in [L]^k$  the coordinate  $X^W$  records the number of times that  $W$  occurs as a contiguous substring in  $S$ . A standard  $k$ -mer method computes a distance between two sequences  $S_1$  and  $S_2$  of lengths  $n_1$  and  $n_2$  by first computing their respective  $k$ -mer vectors  $X_1$  and  $X_2$  and then computing the squared-Euclidean distance

$$\|X_1 - X_2\|_2^2 = \sum_{W \in [L]^k} (X_1^W - X_2^W)^2.$$

Consider two sequences descended from a common ancestor while undergoing a base-substitution process described by standard phylogenetic modeling assumptions. More specifically, we may assume one of the sequences,  $S_1$ , is ancestral to the other,  $S_2$ , and its sites are assigned states in  $[L]$  according to an i.i.d. process with state probability vector  $\pi = (\pi^w)_{w \in [L]}$ . Additionally,  $\pi$  is the stationary distribution of an  $L \times L$  Markov matrix  $M$  describing the single-site state change process from sequence  $S_1$  to sequence  $S_2$ . For continuous-time models, with rate matrix  $Q$  and time (or branch length)  $t$ , one has  $M = \exp(Qt)$ . The probability of a  $k$ -mer  $W = w_1 w_2 \dots w_k \in [L]^k$  in any  $k$  consecutive sites of either single sequence is then  $\pi^W = \prod_{j=1}^k \pi^{w_j}$ . The  $k$ -mer vectors  $X_1$  and  $X_2$  are random variables which summarize  $S_1$  and  $S_2$ .

The following theorem relates the expectation of an appropriately chosen norm of the difference of  $k$ -mer counts  $X_1 - X_2$  to the base-substitution model. Since the expectation can be estimated from  $k$ -mer data, this means that from  $k$ -mer data we can infer information on how much substitution has occurred.

**Theorem 2.1.** *Let  $S_1$  and  $S_2$  be two sequences of length  $n$  generated from an indel-free Markov model with transition matrix  $M$  and stationary distribution  $\pi$ , and let  $X_1$  and  $X_2$  be the resulting  $k$ -mer count vectors. Then*

$$(1) \quad \mathbb{E} \left[ \sum_{W \in [L]^k} \frac{1}{\pi^W} (X_1^W - X_2^W)^2 \right] = 2(n - k + 1)(L^k - (\text{tr } M)^k).$$

Since for each  $W$  the random variable  $X_1^W - X_2^W$  has mean 0, the expectation on the left of equation (1) can be viewed as a (weighted) variance of the  $k$ -mer count difference. Indeed this observation plays an important role in the proof, which appears in Appendix A.

We now derive consequences for the Jukes-Cantor model. In this setting the rate matrix  $Q$  has the form:

$$Q = \begin{pmatrix} -3\alpha & \alpha & \alpha & \alpha \\ \alpha & -3\alpha & \alpha & \alpha \\ \alpha & \alpha & -3\alpha & \alpha \\ \alpha & \alpha & \alpha & -3\alpha \end{pmatrix}.$$

In the Jukes-Cantor model, the rate parameter  $\alpha$  and the branch length  $t$  are confounded with only their product  $at$  identifiable. For simplicity we set  $\alpha = 1/3$  which gives the branch length  $t$  the interpretation of the expected number of substitutions per site. The stationary distribution is uniform. Theorem 2.1 then implies the following.

**Corollary 2.2.** *Let  $S_1$  and  $S_2$  be sequences of length  $n$  generated under the Jukes-Cantor model on an edge of length  $t$ . Let  $X_i$  be the  $k$ -mer count vector of  $S_i$  and let  $d = \mathbb{E} [\|X_1 - X_2\|_2^2]$  be the expected squared Euclidean distance between the  $k$ -mer counts. Then*

$$(2) \quad t = -\frac{3}{4} \ln \left( \frac{4}{3} \sqrt[k]{1 - \frac{d}{2(n-k+1)}} - \frac{1}{3} \right).$$

Equation (2) thus gives a model-corrected estimate of the branch length  $t$  under the Jukes-Cantor model, when in place of the true expected value  $d$  one uses an estimate obtained from data.

*Proof of Corollary 2.2.* To specialize Theorem 2.1 to the Jukes-Cantor model, take  $L = 4$ , and  $\pi^W = 4^{-k}$  for all  $W \in \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}^k$ . Dividing both sides of equation (1) by  $4^k$  we deduce that

$$(3) \quad d = 2(n-k+1) \left( 1 - (\text{tr } M/4)^k \right).$$

For the Jukes-Cantor model

$$M = \exp(Qt) = \begin{pmatrix} y & x & x & x \\ x & y & x & x \\ x & x & y & x \\ x & x & x & y \end{pmatrix}$$

with

$$(4) \quad x = \frac{1 - \exp(-4t/3)}{4}, \quad y = \frac{1 + 3 \exp(-4t/3)}{4},$$

so that  $\text{tr } M = 1 + 3 \exp(-4t/3)$ . Substituting this into equation (3) and solving for  $t$  yields the desired formula.  $\square$

Next we derive an analogous result for the Kimura 3-parameter model, with rate matrix

$$Q = \begin{pmatrix} * & \alpha & \beta & \gamma \\ \alpha & * & \gamma & \beta \\ \beta & \gamma & * & \alpha \\ \gamma & \beta & \alpha & * \end{pmatrix}.$$

**Corollary 2.3.** *Let  $S_1$  and  $S_2$  be two random sequences of length  $n$  generated under the Kimura 3-parameter model on an edge of length  $t$ . Let  $X_i$  be the  $k$ -mer count vector of  $S_i$ . Then*

$$\mathbb{E}[\|X_1 - X_2\|_2^2] = 2(n - k + 1) \left( 1 - \left( \frac{1 + e^{-2(\alpha+\beta)t} + e^{-2(\alpha+\gamma)t} + e^{-2(\beta+\gamma)t}}{4} \right)^k \right).$$

Note that the right side of this equation is strictly increasing as a function of  $t$ . Thus if  $\alpha, \beta, \gamma$  are known, and  $\mathbb{E}[\|X_1 - X_2\|_2^2]$  is estimated, it is straightforward to estimate  $t$  using a numerical root finding algorithm.

For general rate matrices  $Q$ , the matrix  $M = \exp(Qt)$  has trace

$$(5) \quad \text{tr } M = \sum_{i=1}^L e^{\lambda_i t}$$

where  $\lambda_1, \dots, \lambda_L$  are the eigenvalues of  $Q$ , counted with multiplicity. Since  $Q$  is a rate matrix, all these eigenvalues have nonpositive real part. If all the eigenvalues are real, then equation (5) shows  $\text{tr } M$  is a decreasing function of  $t$ . This means we can consistently estimate the branch length if we assume  $Q$  is known and we have an estimate for the expectation in equation (1). For instance, this argument shows that for any time-reversible rate matrix (i.e., from the general time-reversible model GTR) we can obtain statistically consistent estimates for the branch lengths.

### 3. JUKES-CANTOR CORRECTION

In this section, we give a detailed explanation of the statistical consistency for phylogenetic tree reconstruction using our Jukes-Cantor correction from Corollary 2.2. In particular, we explain that without this correction, even with arbitrary amounts of data generated from the model, the  $k$ -mer method based on the squared Euclidean distance is statistically inconsistent for every  $k$ .

Corollary 2.2 gives an estimate of branch lengths under the Jukes-Cantor model based on the value of  $d = \mathbb{E}[\|X_1 - X_2\|_2^2]$ . Applying the same formula to an empirical estimate  $\hat{d}$  of  $d$ , it can thus be viewed as giving a model-based distance correction to the naive distance estimate  $\hat{d}$ . This is similar to the usual Jukes-Cantor correction applied to the frequency  $\hat{p}$  of mismatches of bases in aligned sequences. When  $k = 1$ , equation (2) simplifies to

$$t = -\frac{3}{4} \ln \left( 1 - \frac{4}{3} \cdot \frac{d}{2n} \right)$$

which is clearly very similar to the usual Jukes-Cantor correction obtained from an alignment with  $\frac{d}{2n}$  playing the role of  $p$ .

That  $\frac{d}{2} = p$  for  $k = n = 1$  can be justified rigorously as follows: For a single aligned site in two sequences, the probability of a mismatch is  $p$  under the Jukes-Cantor model. The  $k$ -mer count vectors  $X_1$  and  $X_2$  are the elementary basis vectors  $X_1 = \mathbf{e}_i$  and  $X_2 = \mathbf{e}_j$ , and the quantity  $\|X_1 - X_2\|_2^2$  is 0 or 2 depending if  $i = j$  or  $i \neq j$ . Thus, the expected value  $d = \mathbb{E}[\|X_1 - X_2\|_2^2] = (1 - p) \cdot 0 + p \cdot 2 = 2p$ . It follows that our estimate for the

branch length  $t$  is exactly the Jukes-Cantor corrected estimate when 1-mer frequencies at each site are used to estimate  $d$ . Indeed, formula (2) gives a natural generalization of the pairwise corrected distance to the present context of  $k$ -mers.

To understand the potential impact of the correction of Corollary 2.2 we first work theoretically, by assuming we have the true expected value  $d$  in hand. Later, in Section 6, we use simulations to investigate the usefulness of the branch length estimate (2) with finite length sequences, to understand its practical impact.

We follow the framework suggested by Felsenstein [1978]. We consider an unrooted four-leaf tree with topology 12|34. Two branch lengths  $t_a$  and  $t_b$ , each ranging over the interval  $(0, \infty)$ , are used, with  $t_a$  on edges 2|134, 3|124, and 12|34 and  $t_b$  on the edges 1|234 and 4|123. This tree is depicted in Figure 3.1. The branch lengths are transformed to probabilities  $a$  and  $b$  in  $(0, .75)$ , probabilities that bases at a site differ at opposite ends of a branch.

Consider the naive  $k$ -mer method that uses  $d$  as a distance together with the 4-point condition (or equivalently, Neighbor Joining) to infer a tree topology. To analyze its behavior, we must first relate the expected values of

$$d = 2(n - k + 1) \left( 1 - \left( \frac{1 + 3 \exp(-4t/3)}{4} \right)^k \right),$$

for each taxon pair to the underlying branch parameters  $a$  and  $b$ . As  $a$  is the probability that some change from the current state is made along the edge of scaled length  $t_a$  ( $y$  from equation (4)), we have that the diagonal element from the associated Jukes-Cantor transition matrix is

$$(6) \quad 1 - a = \frac{1 + 3 \exp(-4t_a/3)}{4}$$

and thus

$$t_a = -\frac{3}{4} \ln \left( 1 - \frac{4}{3}a \right).$$

Similar arithmetic gives the formula for  $t_b$ .

This yields the following formulas for the expected distance in terms of the parameters  $a, b$ :

$$\begin{aligned} d_{12} = d_{34} &= 2(n - k + 1) \left( 1 - \left( \frac{1 + 3(1 - \frac{4}{3}a)(1 - \frac{4}{3}b)}{4} \right)^k \right), \\ d_{13} = d_{24} &= 2(n - k + 1) \left( 1 - \left( \frac{1 + 3(1 - \frac{4}{3}a)^2(1 - \frac{4}{3}b)}{4} \right)^k \right), \\ d_{14} &= 2(n - k + 1) \left( 1 - \left( \frac{1 + 3(1 - \frac{4}{3}a)(1 - \frac{4}{3}b)^2}{4} \right)^k \right), \\ d_{23} &= 2(n - k + 1) \left( 1 - \left( \frac{1 + 3(1 - \frac{4}{3}a)^3}{4} \right)^k \right). \end{aligned}$$



To construct correctly the unique true tree 12|34 using the 4-point condition or Neighbor Joining requires that these distances satisfy

$$d_{12} + d_{34} < \min(d_{13} + d_{24}, d_{14} + d_{23}).$$

Note that  $d_{12} + d_{34} < d_{13} + d_{24}$  for all  $a > 0$ , so we focus on the condition

$$d_{12} + d_{34} < d_{14} + d_{23}.$$

Using the formulas above, this becomes:

$$2 \left(1 + 3\left(1 - \frac{4}{3}a\right)\left(1 - \frac{4}{3}b\right)\right)^k \geq \left(1 + 3\left(1 - \frac{4}{3}a\right)\left(1 - \frac{4}{3}b\right)^2\right)^k + \left(1 + 3\left(1 - \frac{4}{3}a\right)^3\right)^k.$$

The values of  $a$ ,  $b$  for which this is satisfied are shown by the white regions in Figure 3.1, for  $k = 1, 3, 5$ . As  $k$  increases the white regions change; when  $k = 1$  the boundary curve is a circle, and as  $k \rightarrow \infty$  it approaches a parabola with vertex in the upper right corner, passing through the lower left. Note that the white region indicates where the naive  $k$ -mer distance inference behaves well provided one knows  $d$  exactly — in practice one only has an estimate of  $d$  and should not expect even this good behavior.

In contrast, using the corrected Jukes-Cantor  $k$ -mer distance from equation (2) to make diagrams analogous to those of Figure 3.1 would show the entire square white. If  $d$  were known exactly, inference would be perfect. The corrected distances lead to statistically consistent distance methods on 4-taxon trees. More generally, our argument in Section 2 shows that we can use Theorem 2.1 to derive statistically consistent estimates for the evolutionary time between species when we have a known time-reversible rate matrix  $Q$ .

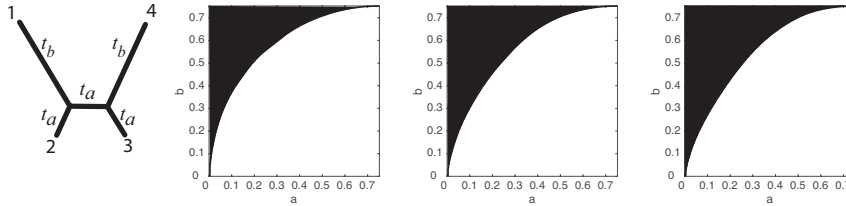


FIGURE 3.1. White region is the zone of consistency for tree inference using the naive  $k$ -mer distance combined with the 4-point condition. From left to right,  $k = 1, 3$ , and  $5$ . The usual “Felsenstein Zone” is in the upper left.

#### 4. IDENTIFIABILITY OF INDEL-FREE MODEL PARAMETERS

The results in Section 2 prove that, with knowledge of the stationary base frequency  $\pi$  of an unknown Markov matrix  $M$  describing base substitutions from one sequence to another,  $\text{tr} M$  is identifiable from the joint distribution of  $k$ -mer counts in the two sequences. If one assumes a continuous time model with  $M = \exp(Qt)$  and  $Q$  a known time-reversible rate matrix, then this is sufficient to identify lengths  $t$  between taxa on the tree. As a consequence, with  $Q$  known the metric phylogenetic tree relating many taxa is identifiable.

In fact, more is true:  $\pi$  and  $M$  are identifiable from 1-mer count distributions as well. This is the result in the next proposition, which in addition to being interesting in its



own right, plays a role in the proof of Theorem 2.1. Its proof appears in Appendix B. Note that in this result we do not assume that base frequency distributions  $\pi_i$  are the stationary vectors of the Markov matrix.

**Proposition 4.1.** *From the joint distribution of 1-mer count vectors  $X_1$  and  $X_2$  of two sequences  $S_1$  and  $S_2$  of length  $n$ , one can identify the distributions  $\pi_1$  and  $\pi_2$  of bases in each sequence, and the joint distribution  $P = \text{diag}(\pi_1)M$  of bases at a single site in the two sequences. Specifically,  $\pi_i = \frac{1}{n}\mathbb{E}[X_i]$ , and for  $w, u \in [L]$ ,*

$$P_{wu} = \frac{1}{2} \left( \pi_1^w + \pi_2^u - \frac{\pi_1^w \pi_2^u}{n} \mathbb{E} \left[ \left( \frac{X_1^w}{\pi_1^w} - \frac{X_2^u}{\pi_2^u} \right)^2 \right] \right).$$

The formula for  $P_{wu}$  in this proposition ultimately underlies our suggested practical inference method. However, there is a simpler formula, applying for any  $k$ , showing that from a joint  $k$ -mer count vector distribution one can identify the joint probabilities  $P_{wu}$ : For sequences of length  $n$  and the particular  $k$ -mers  $W = www \dots w$  and  $U = uuu \dots u$ ,

$$\text{Prob}(X_1^W = n - k + 1, X_2^U = n - k + 1) = (P_{wu})^n.$$

Of course the method of estimation suggested by this approach is useless in practice, since it is based on events that are rarely, if ever, observed.

Nonetheless, since  $P$  and  $\pi_1$  can be found from the joint distribution of  $X_1$  and  $X_2$  for any  $k$ , the transition matrix  $M = \text{diag}(\pi_1)^{-1}P$  is also identifiable. In the continuous-time model setting, where  $M = \exp(Qt)$ ,  $Q$  can be found, first up to a scalar multiple, and then normalized. Putting this together yields the following.

**Theorem 4.2.** *For an indel-free GTR model, all parameters, both numerical ones and tree topology, can be identified from pairwise joint  $k$ -mer count distributions.*

If we consider sequences three-at-a-time, rather than pairwise, we obtain an analog of Proposition 4.1, again without assuming stationarity. This new result is based on third moments, rather than second, and its proof is given in Appendix B.

**Proposition 4.3.** *For a 3-leaf tree, the joint distribution  $P = (P_{uvw})$  of site patterns is identifiable from the joint 1-mer count vector distributions of the 3 taxa. Specifically, define a random variable*

$$Y_{uvw} = \alpha X_1^u + \beta X_2^v + \gamma X_3^w,$$

where  $\alpha, \beta, \gamma$  are constants chosen so

$$\alpha \pi_1^u + \beta \pi_2^v + \gamma \pi_3^w = 0.$$

Then

$$(7) \quad P_{uvw} = \frac{1}{6\alpha\beta\gamma n} \mathbb{E}(Y_{uvw}^3) + \frac{1}{2} \left( \frac{\alpha + \beta}{\gamma} P_{uv+} + \frac{\alpha + \gamma}{\beta} P_{u+v} + \frac{\beta + \gamma}{\alpha} P_{+vw} \right) - \frac{1}{6} \left( \frac{\alpha^2}{\beta\gamma} \pi_1^u + \frac{\beta^2}{\alpha\gamma} \pi_2^v + \frac{\gamma^2}{\alpha\beta} \pi_3^w \right).$$

where the pairwise marginal distributions  $P_{uv+}$ ,  $P_{u+v}$ , and  $P_{+vw}$  in equation (7) are identifiable by Proposition 4.1.

Proposition 4.3 is significant in that it establishes that the distribution of 1-mer counts contains enough information to identify parameters of more general models than our preceding arguments allow. Recall, for instance, that parameters for the General Markov (GM) model, in which the base substitution process on each edge of the tree can be specified by a different Markov matrix, are identifiable from the marginalization of the site pattern distribution to 3-taxon sets [Chang, 1996], but are not identifiable from pairwise marginalizations. In the present context of  $k$ -mers, we obtain the following.

**Corollary 4.4.** *For an indel-free GM model, all parameters, both numerical ones and tree topology, are identifiable from the joint 1-mer count vector distributions on  $n$  taxa.*

## 5. PRACTICAL $k$ -MER DISTANCES BETWEEN SEQUENCES

In this section, we apply the results of Section 2 to develop practical methods for estimating pairwise distances between sequences. Those derivations were made under the assumption that sequences evolved in the absence of an indel process, and thus that sequences could be unambiguously aligned. In practice, however, we desire a method of distance estimation that can be applied in the presence of a mild indel process, without a precise alignment. Although this violates our model assumptions, in Section 6 we use simulations to investigate how robust our resulting method is to such a violation.

Assuming a Jukes-Cantor process of site substitution and no indel process, formula (2) of Corollary 2.2 suggests a natural definition for a distance, provided we have a good method of approximating  $d = \mathbb{E}[\|X_1 - X_2\|_2^2]$ . If the observed values of the random variables  $X_1$  and  $X_2$  are denoted in lower case, so  $x_1$  and  $x_2$  are observed  $k$ -mer count vectors, then one could simply compute

$$(8) \quad \|x_1 - x_2\|_2^2 = \sum_{W \in [L]^k} (x_1^W - x_2^W)^2$$

as a point estimate for  $d$ . This is a very poor estimate for the expected value, however, since only one sample ( $\|x_1 - x_2\|_2^2$ ) is used to estimate a mean. Indeed, this estimate has large variance. Moreover, naively increasing sequence length (number of  $k$ -mers) would do nothing to address the fundamental problem of needing more samples to estimate a mean well.

To obtain a better estimate of  $d$ , with smaller variance, we instead subdivide the two sequences into a fixed number  $B$  of contiguous blocks. Assuming for  $1 \leq i \leq B$  that the  $i$ th blocks of the two sequences are at least roughly orthologous, we compute the  $k$ -mer frequencies  $x_{j,i}$  for each block  $i$  in sequence  $j$ . Then the values of  $\|x_{1,i} - x_{2,i}\|_2^2$  for the  $B$  blocks can be averaged to estimate  $d$ . In this framework, we are adopting the approach introduced by Daskalakis and Roch [2013].

We have in mind two scenarios for using this approach on data, which are displayed in Figure 5.1. The first is under the assumption that if indels occurred, they were distributed evenly over the sequences. Then if the blocks are defined as a fixed fraction of the full sequence lengths, most of the sites in the  $i$ th blocks of the two sequences will be orthologous. The second is that the blocks arise naturally in the data; for instance if a dataset consists of multiple genes, then each gene can be treated as a block. In this case,

the point estimates for each gene would be averaged over all genes, making appropriate adjustments for their varying lengths.

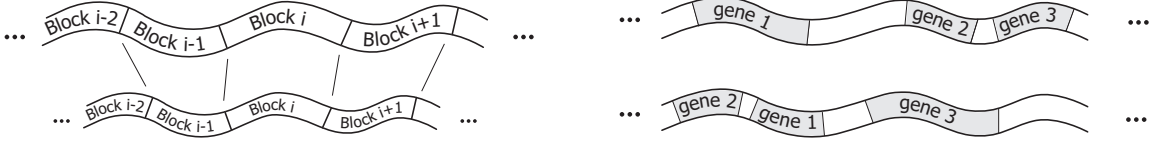


FIGURE 5.1. On the left, two sequences in which blocks  $i$  are roughly orthologous, perhaps due to a uniform indel process. On the right, two genomes in which genes serve as blocks for data analysis.

To be precise, in addition to specifying  $k$ , under the first scenario we must also specify a number  $B$  of blocks to be used in our calculations. To subdivide a sequence  $S_j$  of length  $n_j$  as uniformly as possible, each block will have length  $n_{j,i} = n_j/B$ , suitably rounded for  $1 \leq i \leq B$ , so block lengths for a single sequence can differ at most by one. Under the second scenario, using natural blocks like genes, the length  $n_{j,i}$  is specified by the data, and will vary more widely.

Now for block  $i$  in sequence  $j$ , let  $x_{j,i}$  be the  $k$ -mer count vector and  $\mu_{j,i} = (n_{j,i}-k+1)/4^k$  the mean  $k$ -mer count under the Jukes-Cantor model. We define

$$(9) \quad \tilde{d} = \frac{1}{B} \sum_{i=1}^B \left\| \frac{(x_{1,i} - \mu_{1,i})}{\sqrt{n_{1,i} - k + 1}} - \frac{(x_{2,i} - \mu_{2,i})}{\sqrt{n_{2,i} - k + 1}} \right\|_2^2.$$

Note that in this formula both the centering of  $x_{j,i}$  by subtracting  $\mu_{j,i}$  and the normalization by dividing by the square root of the number of  $k$ -mers depend upon the length  $n_{j,i}$ . In the special situation where  $n_{j,i} = n$  for all  $i, j$ , and hence  $\mu_{j,i} = \mu$ , this reduces to

$$\tilde{d} = \left( \frac{1}{n - k + 1} \right) \frac{1}{B} \sum_{i=1}^B \|x_{1,i} - x_{2,i}\|_2^2 \approx \frac{d}{n - k + 1}.$$

Comparing this estimate for  $\tilde{d}$  to equation (2), it is natural to define a Jukes-Cantor  $k$ -mer distance  $d_{JC}^{k,B}$ , dependent on  $k$  and  $B$ , by

$$(10) \quad d_{JC}^{k,B} = -\frac{3}{4} \ln \left( \frac{4}{3} \sqrt[3]{1 - \frac{\tilde{d}}{2} - \frac{1}{3}} \right).$$

We use this formula extensively in the simulations whose results are presented in the next section.

In examining (10), it is unclear *a priori* which values of  $k$  and  $B$  will yield the best estimate for  $d_{JC}^{k,B}$ . In the particular case that sequences evolved without an indel process, the lowest variance estimate of  $d_{JC}^{k,B}$  is obtained by taking the largest number  $B$  of samples, i.e. each block has length  $k$  (the smallest possible length which allows  $k$ -mers to be counted). However, in the presence of an evolutionary indel process a true alignment of sequences would contain gaps, and such short block sizes would give poor results. For

good performance, we need the  $i$ th blocks in the two sequences to be composed mostly of orthologous sites. If the block size is small, this is unlikely to be true, as even a mild indel process might result in orthologs residing in different blocks. The art is to find the right compromise between a large number  $B$  of blocks and a large enough length  $n_{j,i}$  for each block to ensure many orthologs. Results of simulation studies in the next section confirm this trade-off.

Using 1-mer distributions and taking into account a particular model of the indel process, Daskalakis and Roch [2013] give a detailed analysis of a distance method along the lines described here. Their results suggest that the block sizes should be of size roughly the square root of total sequence length. While the approach of Daskalakis and Roch inspired our results, since our approach to a  $k$ -mer distance is based on a model without indels, and our extension to a distance formula for sequence evolution in the presence of indels is heuristic, we can offer no such guidance. A fruitful direction for future research is to explore  $k$ -mer distances under some explicit model of sequence evolution with indels.

## 6. SIMULATION STUDIES

**Methods.** We performed extensive simulations to attempt to understand how the distance formula in (10) might work in practice and to compare distance methods with  $d_{JC}^{k,B}$  to other alignment-free methods for reconstructing phylogenetic trees from sequence data. Data was simulated using the sequence evolution simulator INDELible [Fletcher and Yang, 2009], which produces sequence data under standard base substitution models with or without an additional insertion and deletion process.

All of our simulations use the Jukes-Cantor (JC) substitution model on 4-taxon tree. We consider only trees in which two branch lengths occur,  $t_a$  and  $t_b$ , as shown in Figure 6.1. This allows us to investigate performance over an important range of parameter space, yet still display the success of an algorithm in an easily-interpretable 2-dimensional display, as introduced by Huelsenbeck [1995].

The two branch lengths  $t_a$  and  $t_b$  each range over the interval  $(0, \infty)$ , but are transformed to probabilities  $a$  and  $b$  in range  $(0, .75)$ , probabilities that bases at a site differ at opposite ends of a branch (see equation 6). In this interval, we sampled points from .01 to .73, with increments of .02, to get a  $37 \times 37$  grid of transformed branch lengths. For each choice of branch lengths we generated 100 sets of four sequences, used a specific method to recover the tree topology, and recorded the frequency the method under study reconstructs the correct tree 12|34 from the simulated data.

The middle and the right plots in Figure 6.1 show typical Huelsenbeck diagrams presenting results from such simulations. The dark red regions correspond to regions where the method reconstructs the true tree topology, with split 12|34, close to 100% of the time. Dark blue regions are regions where inference is strongly biased against the correct tree, reconstructing it close to 0% of the time. Light blue corresponds to a method constructing the true tree correctly about 33% of the time; that is, the method is indistinguishable from the process of randomly picking the tree topology to return. For any phylogenetic method applied to simulated sequence data, one typically sees light blue in the upper right of these figures ( $a \approx b \gg 0$ ), darker blue in the upper left ( $b \gg a$ ) in the “long-branch

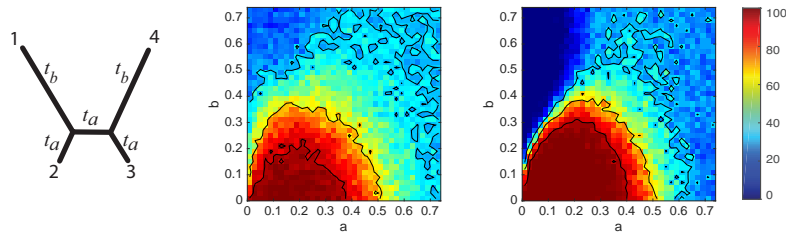


FIGURE 6.1. Figure on the left displays the model tree used for simulations. The middle and right figures are representative Huelsenbeck diagrams for some (unspecified) methods of inference. The horizontal axis is labeled by  $a$  and the vertical one with  $b$ , both in the range  $(0, .75)$  after transformation. Contour lines are drawn at levels .95, .67, and .33. The figure on the right suggests significant long branch attraction is present, as witnessed by the strong bias against the correct tree (much less than 33% correct) along the upper left side.

attraction” zone where the tree with split 14|23 tends to be inferred, and red where  $a \geq b$  are of small to moderate size.

In our simulation studies parameters other than branch lengths were also varied. Several of these govern the details of the model of sequence evolution:

- (1) Sequence length
- (2) The rates of insertions and deletions
- (3) Parameters for the distribution of the size of indels

Other parameters control the specifics of implementing our  $k$ -mer method:

- (4)  $k$ , the  $k$ -mer length
- (5)  $B$ , the number of blocks

For simulations that combine a site substitution process with an indel process, one must specify the location of a root in the tree, since indels change the sequence length; we chose the midpoint of the interior branch to root the tree. For initial sequence length at this root, we chose the lengths  $L = 1000, 10000$ . INDELible requires users to choose a rate of insertion events and a rate of deletion events, specified relative to the substitution rate; we set these equal and denote the common value  $\mu$ . In assuming that insertions and deletions are rare relative to base substitutions, we varied this parameter over the values  $\mu = .01, .05, .1$ . We used the Lavalette distribution as implemented in INDELible for determining the lengths of inserted and deleted segments: For parameters  $(a, M)$ , this is the distribution on  $S = \{1, 2, \dots, M\}$  such that for  $G \in S$ ,  $Pr(G) \propto \left(\frac{GM}{M-G+1}\right)^{-a}$ . Large  $M$  and small  $a$  tend to produce longer insertion and deletion events. Fletcher and Yang [2009] suggest that values of  $a \in [1.5, 2]$  with a large  $M$  give a reasonable match with data. We tried values  $a = 1.1$  (as used in [Chan et al., 2014]), 1.5, 1.8, and  $M = 100$ .

For testing our  $k$ -mer methods on simulated data, we varied  $k = 1, 3, 5, 7$  and the number of blocks  $B$  ranged over 1, 5, 25, 100, 250, 500, provided this allowed a block size at least  $k$ .

**Performance on simulated sequence data.** As presentation of all simulation results would require considerable space, here we present only representative examples to illustrate key points. The supplementary materials [Allman et al., 2015] contain results of other simulations.

*Simulations with no indel process.* We begin by discussing simulations in which no indel process occurs. This is the situation in which our theoretical results were derived, and these runs investigate solely the effect of having simulated sequence data of finite length. These trials are, of course, somewhat artificial in that in the absence of an indel process we have exact alignments of sequences, and there is no reason to use an alignment-free phylogenetic method. Nonetheless, they represent a measuring rod for evaluating the performance of the new methods presented here.

We set the sequence length to 1000, and for comparison to traditional approaches, produce Hulsenbeck diagrams in Figure 6.2 using (i) the standard JC pairwise distance formula for the sequences with the true alignment as produced by INDELible together with Neighbor Joining (NJ), and (ii) the standard JC distance formula after a pairwise alignment, followed by NJ. Alignment in (ii) was performed by the Needleman-Wunsch algorithm implemented in MATLAB’s Bioinformatics Toolbox, but with scoring parameters set to NCBI defaults: match= 2, mismatch= -3, gap existence= -5, gap extension= -2. Simulation (i) represents a standard that would be desirable, but probably impossible, to match, as  $k$ -mer methods make no use of the alignment itself and a true alignment is never known in practice. Simulation (ii) offers a more realistic setting with results we might hope to match or beat, in which large amounts of substitution results in quite dissimilar sequences, and the introduction of gaps in the alignment process. The distance estimates computed with these ‘gappy’ alignments can be quite far from the true pairwise distances underlying the simulated data.

In Figure 6.2 (ii), for the simulation in which sequences were aligned before distances were computed, there is a rather pronounced region of parameter space to the upper left displaying the phenomenon of long branch attraction. In addition, surrounding the red area where the true tree is reliably constructed, we see a halo of darkish blue, illustrating another region of parameter space with a weaker bias against the correct tree. Comparing (i) and (ii), it is clear that the alignment process markedly degrades performance of the inference procedure.

In Figure 6.3, we present results using the same simulated sequences (JC and no indels) as in the previous figure, but use the distance  $d_{JC}^{5,B}$  with NJ. With  $k = 5$  held fixed, we vary  $B = 1, 5, 25, 100$ . This sequence of diagrams, in which the red area increases with  $B$ , illustrates that in the absence of indels and with  $k$  held constant, increasing the number of blocks is advantageous, as was anticipated in Section 5.

Comparing Figure 6.3 with Figure 6.2 (ii) suggests that when data sequences are quite dissimilar, and a researcher might be inclined to align sequences before a phylogenetic analysis, that our  $k$ -mer method can outperform the traditional approach (alignment +  $d_{JC}$  + NJ). In particular, using  $d_{JC}^{k,B}$  the red region of good performance is enlarged, and the phenomenon of long branch attraction is significantly lessened. (Further simulations below will return to this issue when there is a mild indel process.)



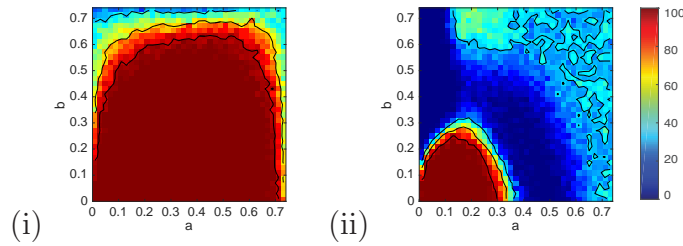


FIGURE 6.2. Figures illustrating the accuracy of inference of tree topology on simulated data with no gaps, using the Jukes-Cantor distance and Neighbor Joining. Simulated sequences have length 1000 bp with no indel process. In (i) the correct alignment is used, and in (ii) pairwise alignments are found before the JC distance is computed.

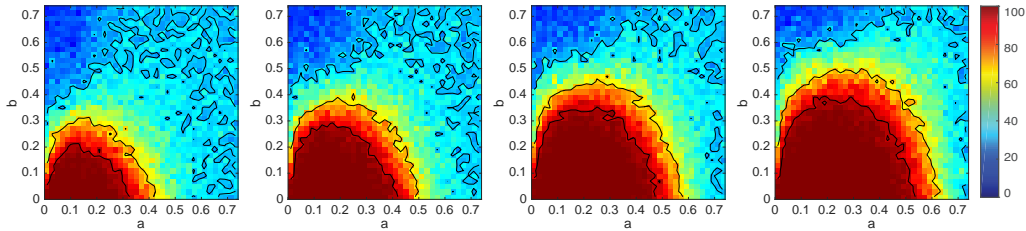


FIGURE 6.3. Figures illustrating the accuracy of inference of tree topology on simulated data with no indels, using a 5-mer distance  $d_{JC}^{5,B}$  and Neighbor Joining. Simulated sequences have length 1000 bp with no indel process. From left to right,  $B = 1, 5, 25, 100$ .

Now fixing the number of blocks  $B = 25$ , but varying  $k$  in  $d_{JC}^{k,25}$ , with NJ we produce Figure 6.4. Notice here that with a fixed number of blocks, both too small and too large a value of  $k$  reduces performance.

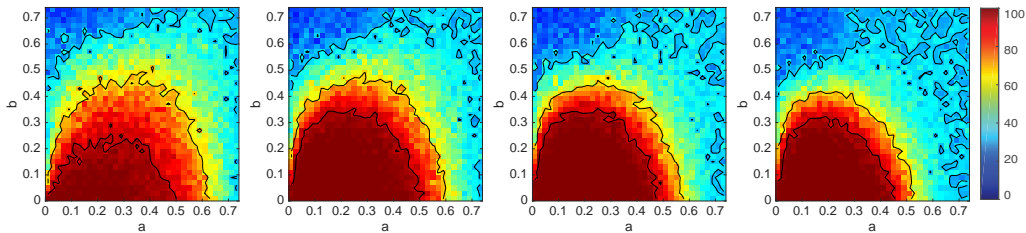


FIGURE 6.4. Figures illustrating the accuracy of inference of tree topology on data with no indels, using a  $k$ -mer distance  $d_{JC}^{k,25}$  and Neighbor Joining. Simulated sequences have length 1000 bp with no indel process. From left to right,  $k = 1, 3, 5, 7$ .

In summary, while no performance of our  $k$ -mer distance comes close to the ideal of Figure 6.2 (i) (true alignment+ $d_{JC}$ +NJ), the  $k$ -mer methods often perform better



than (alignment+ $d_{JC}$ +NJ) as shown in Figure 6.2 (ii). Computing erroneous pairwise alignments results in a large region of parameter space in which long branch attraction is pronounced, but such biased inference is almost absent when  $d_{JC}^{k,B}$  is used. When the sequence length and number of blocks  $B$  are fixed, the choice of  $k$  can affect performance, with either too large or too small a  $k$  causing degradation. It is unclear how to determine an “optimal” choice of  $k$  except through simulation.

*Simulations with an indel process.* With a length of 1000 bp for the sequence at the root of the tree, we now introduce an indel process with rate  $\mu = .05$  and Lavalette parameters  $a = 1.8$ ,  $M = 100$ . This means on average one insertion event and one deletion event occurs for every 20 base substitutions. Repeating reconstruction methods (i) and (ii) of Figure 6.2 on these datasets with indels, we obtain Figure 6.5.

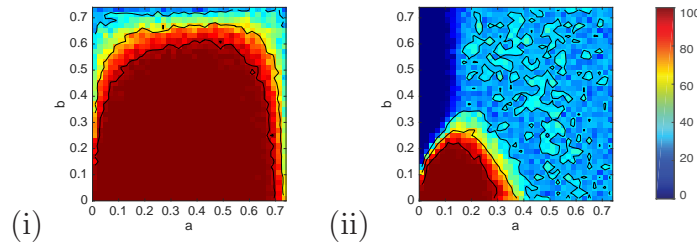


FIGURE 6.5. Figures illustrating the accuracy of inference of tree topology on simulated data with indels using the Jukes-Cantor distance and Neighbor Joining. The root sequence is 1000 bp. The indel process is determined by  $\mu = .05$  and Lavalette parameters  $a = 1.8$ ,  $M = 100$ . In (i) the true alignment is used, and in (ii) pairwise alignments are found before the JC distance is computed.

While Figure 6.5 (i) shows excellent performance, it assumes the correct alignment (including gaps) is known, which is unrealistic in any empirical study. Analysis (ii) is one that could be performed on real data, and should be compared to Figure 6.2 (ii) above. For sequence data with indels the region of good performance is similarly shaped, but smaller, than that for data without indels. This is to be expected, since even when few substitutions occur, indels could lead to erroneous alignment. In both Figure 6.2 (ii) and Figure 6.5 (ii), long branch attraction is present in the upper left corner of parameter space. In contrast, however, in Figure 6.5 (ii) the area to the upper right surrounding the area of good reconstruction does not display a bias against correct reconstruction, but rather a uniform randomness in selection of the tree.

Setting  $k = 5$  and  $B = 1, 5, 25, 100$ , and using  $d_{JC}^{5,B}$ +NJ on the sequence data with indels produces Figure 6.6. Note that increasing the number of blocks first improves performance, but then degrades it. This is explained by a large number of blocks producing a small block size, which increases the chance that corresponding blocks in two sequences share few homologous sites, as was discussed in Section 5. This phenomenon is only seen on data simulated with an indel process.

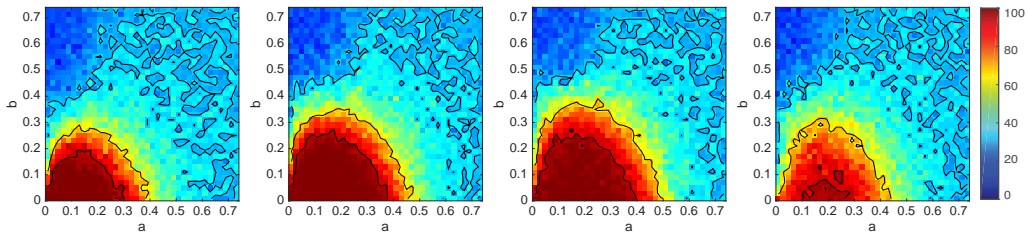


FIGURE 6.6. Figures illustrating the accuracy of inference of tree topology on simulated data with indels, using a 5-mer distance  $d_{JC}^{5,B}$  and Neighbor Joining. The root sequence is 1000 bp. The indel process is determined by  $\mu = .05$  and Lavalette parameters  $a = 1.8$ ,  $M = 100$ . From left to right,  $B = 1, 5, 25, 100$ .

With the number of blocks set at 25, but varying  $k$  in  $d_{JC}^{k,25}$ , we obtain Figure 6.7. Again we note that for a fixed number of blocks, too small or large a value of  $k$  degrades performance.

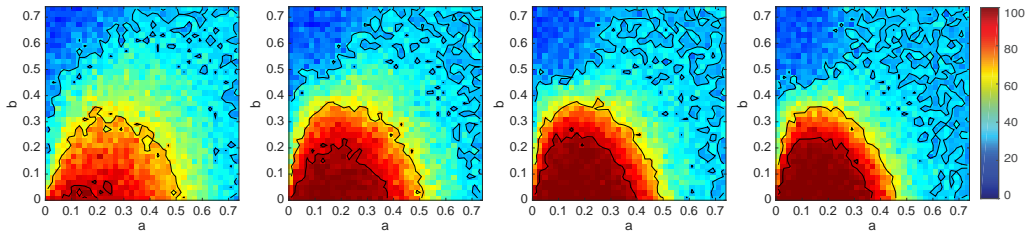


FIGURE 6.7. Figures illustrating the accuracy of inference of tree topology on simulated data with indels using a  $k$ -mer distance  $d_{JC}^{k,25}$  and Neighbor Joining. The root sequence is 1000 bp. The indel process is determined by  $\mu = .05$  and Lavalette parameters  $a = 1.8$ ,  $M = 100$ . From left to right,  $k = 1, 3, 5, 7$ .

These figures illustrate that even in the presence of a mild indel process, the  $k$ -mer method described here can perform as well as pairwise alignment with traditional distance methods in the regions of parameter space where those work well, yet greatly reduce the pronounced long-branch attraction problems that incorrect alignments introduce in other regions of parameter space. Although the  $k$ -mer distance  $d_{JC}^{k,B}$  was derived using a model with no indels, these simulations demonstrate its performance is somewhat robust to violation of that assumption.

*Other  $k$ -mer methods.* To conclude, in Figure 6.8 we display some diagrams that illustrate the performance of other  $k$ -mer distance methods [Vinga and Almeida, 2003, Chan et al., 2014, Reinert et al., 2009, Wan et al., 2010, Edgar, 2004b] on simulated data with indels. The datasets were the same ones used in producing Figures 6.5, 6.6, 6.7.

In the figure below, we use  $k$ -mer distances previously proposed: With  $x_1$  and  $x_2$  the observed  $k$ -mer count vectors, two of these distances are

$$(11) \quad L_2^2 = \|x_1 - x_2\|_2^2$$

$$(12) \quad \theta = \arccos(x_1 \cdot x_2 / \|x_1\|_2 \|x_2\|_2)$$

These have long been studied for sequence comparison [Vinga and Almeida, 2003], though primarily for non-phylogenetic applications. Yang and Zhang [2008] used a variation of the  $L_2^2$  distance based on replacing  $x_1$  and  $x_2$  with  $x_1/(n_1 - k + 1)$  and  $x_2/(n_2 - k + 1)$ , respectively, where  $n_i$  is the length of sequence  $i$ .

The next three have appeared in phylogenetic investigations of Chan et al. [2014], but are based on sequence comparison methods developed for other purposes, as reviewed by Song et al. [2014]. With  $\tilde{x}_i = x_i - \mathbb{E}(x_i)$  the centralized count vector, let

$$D_2(x_1, x_2) = x_1 \cdot x_2,$$

$$D_2^S(x_1, x_2) = \sum_W \frac{\tilde{x}_1^W \tilde{x}_2^W}{\sqrt{(\tilde{x}_1^W)^2 + (\tilde{x}_2^W)^2}},$$

$$D_2^*(x_1, x_2) = \sum_W \frac{\tilde{x}_1^W \tilde{x}_2^W}{\sqrt{\mathbb{E}(x_1^W) \mathbb{E}(x_2^W)}},$$

as did Reinert et al. [2009] and Wan et al. [2010]. Then define the distances

$$(13) \quad d_2 = \left| \ln \frac{D_2(x_1, x_2)}{\sqrt{D_2(x_1, x_1) D_2(x_2, x_2)}} \right|,$$

$$(14) \quad d_2^S = \left| \ln \frac{D_2^S(x_1, x_2)}{\sqrt{D_2^S(x_1, x_1) D_2^S(x_2, x_2)}} \right|,$$

$$(15) \quad d_2^* = \left| \ln \frac{D_2^*(x_1, x_2)}{\sqrt{D_2^*(x_1, x_1) D_2^*(x_2, x_2)}} \right|,$$

with the convention that the logarithm of a negative number is set to  $\infty$ . As the distances  $d_2$  and  $\theta$  differ from each other by the application of a monotone function, for 4-leaf trees they perform identically using UPGMA, and quite similarly with NJ. Thus, in Figure 6.8 the plot for the  $\theta$  distance is not shown.

Finally, for comparison purposes, we include the distance used in the initial step of the MUSCLE alignment algorithm [Edgar, 2004b],

$$(16) \quad m = 1 - \sum_W \frac{\min\{x_1^W, x_2^W\}}{(n - k + 1)},$$

where  $n = \min(n_1, n_2)$  is the length of the shorter of the two sequences. Since MUSCLE uses UPGMA as its default for tree building, we performed both NJ and UPGMA for all of these distances.

As is apparent in in Figure 6.8, with  $k = 5$  most of the distances in (11-16) exhibit long branch attraction bias which is generally quite pronounced, and fail to match the performance of the 5-mer distance derived here.

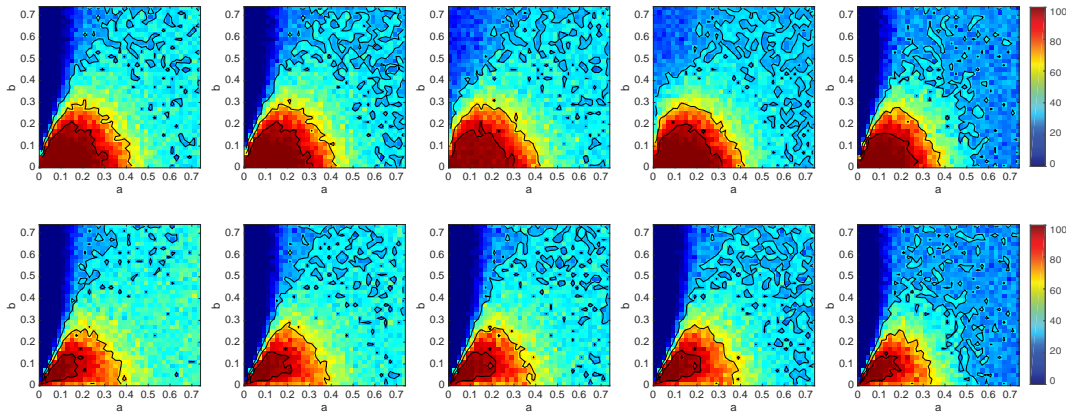


FIGURE 6.8. Figures illustrating the accuracy of inference of tree topology on simulated data with indels, using a variety of distances and Neighbor Joining and UPGMA. The root sequence is 1000 bp. The indel process is determined by  $\mu = .05$  and Lavalette parameters  $a = 1.8$ ,  $M = 100$ . Columns in the figure are, from left to right, obtained using the distances given in equations (11), (13), (14), (15), and (16), all with  $k = 5$ . The top row of figures uses Neighbor Joining, and the bottom UPGMA.

## 7. CONCLUSIONS

We have derived model-based distance corrections for the squared-Euclidean distance between  $k$ -mer count vectors of sequences. Our results show that the uncorrected use of the squared-Euclidean distance leads to statistically inconsistent estimation of the tree topology, with inherent long-branch attraction problems. This statistical inconsistency occurs even at short branch lengths, and is strongly manifested in simulations. Simulations show that our corrected distance outperforms previously proposed  $k$ -mer methods, and suggest that many of those are statistically inconsistent with long-branch attraction biases.

All our results have been derived under the assumption that there are no insertions or deletions in the evolution of sequences. Our simulations indicate that even if a mild indel process occurred, a simple extension of the corrected method still performs well. It remains to develop  $k$ -mer methods assuming an indel process, using the indel model structure to develop a more precise correction on the distance.

Daskalakis and Roch [2013] developed an alignment-free phylogenetic tree inference method for a model with a simple indel process. Their method can be seen as a 1-mer method. While we have not compared their method directly to any of ours, our simulations suggest that 1-mer methods perform poorly compared to  $k$ -mer methods with larger  $k$ . This suggests that a natural line for future research would be to combine the approach

of Daskalakis and Roch with ours to develop consistent  $k$ -mer methods that take into account the structure of an underlying indel model.

#### ACKNOWLEDGEMENT

Seth Sullivant was partially supported by the David and Lucille Packard Foundation and the US National Science Foundation (DMS 0954865).

#### REFERENCES

- E.S. Allman, J.A. Rhodes, and S. Sullivant. Supplementary materials. <http://www.dms.uaf.edu/~eallman/Papers/kmerSupp.html>, 2015.
- G. Blackshields, F. Sievers, W. Shi, A. Wilm, and D.G. Higgins. Sequence embedding for fast construction of guide trees for multiple sequence alignment. *Algorithms Mol. Biol.*, 5:21, 2010.
- C.X. Chan, B. Guillaume, O. Poirion, J.M. Hogan, and M.A. Ragan. Inferring phylogenies of evolving sequences without multiple sequence alignment. *Scientific Reports*, 4:1–9, 2014.
- J.T. Chang. Full reconstruction of Markov models on evolutionary trees: Identifiability and consistency. *Math. Biosci.*, 137(1):51–73, 1996.
- C. Daskalakis and S. Roch. Alignment-free phylogenetic reconstruction: Sample complexity via a branching process analysis. *Annals of Applied Probability*, 23:693–721, 2013.
- R.C. Edgar. Muscle: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.*, 32:1792–1797, 2004a.
- R.C. Edgar. Muscle: A multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics*, 5:113, 2004b.
- J. Felsenstein. Cases in which parsimony or compatibility methods will be positively misleading. *Systematic Biology*, 27:401–410, 1978.
- W. Fletcher and Z. Yang. Indelible: A flexible simulator of biological sequence evolution. *Mol. Biol. and Evol.*, 26:1879–1888, 2009.
- J. Huelsenbeck. Performance of phylogenetic methods in simulation. *Syst. Biol.*, 44:17–48, 1995.
- K. Liu, S. Raghavan, S. Nelesen, C. R. Linder, and T. Warnow. Rapid and accurate large scale coestimation of sequence alignments and phylogenetic trees. *Science*, 324:1561–1564, 2009.
- K. Liu, T.J. Warnow, M.T. Holder, S. Nelesen, J. Yu, A. Stamatakis, and C.R. Linder. SATé-II: Very fast and accurate simultaneous estimation of multiple sequence alignments and phylogenetic trees. *Systematic Biology*, 61:90–106, 2012.
- G. Reinert, D. Chew, F. Sun, and M. Waterman. Alignment-free sequence comparison (I): Statistics and power. *Journal of Computational Biology*, 16:1615–1634, 2009.
- F. Reyes-Prieto, A.J. Garcia-Chequer, H. Jaimes-Diaz, J. Casique-Almazan, J. M Espinosa-Lara, R. Palma-Orozco, Alfonso Mendez-Tenorio, Rogelio Maldonado-Rodriguez, and Kenneth L Beattie. Lifeprint: A novel  $k$ -tuple distance method for construction of phylogenetic trees. *Adv Appl Bioinforma Chem.*, 4:13–27, 2011.

- N. Saitou and M. Nei. The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4:406–425, 1987.
- W. Shen, H. Wong, Q. Xiao, X. Guo, and S. Smale. Introduction to the peptide binding problem of computational immunology: New results. *Foundations of Computational Mathematics*, 14:951–984, 2014.
- F. Sievers, A. Wilm, D.G. Dineen, T.J. Gibson, K. Karplus, W. Li, R. Lopez, H. McWilliam, M. Remmert, J. Söding, J.D. Thompson, and D.G. Higgins. Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Molecular Systems Biology*, 7:539, 2011.
- K. Song, J. Ren, G. Reinert, M. Deng, M. Waterman, and F. Sun. New developments of alignment-free sequence comparison: Measures, statistics and next-generation sequencing. *Briefings in Bioinformatics*, 15:343–353, 2014.
- J.L. Thorne, H. Kishino, and J. Felsenstein. An evolutionary model for maximum likelihood alignment of DNA sequences. *Journal of Molecular Evolution*, 33:114–124, 1991.
- J.L. Thorne, H. Kishino, and J. Felsenstein. Inching toward reality: An improved likelihood model of sequence evolution. *Journal of Molecular Evolution*, 34:3–16, 1992.
- S. Vinga and J. Almeida. Alignment-free sequence comparison — a review. *Bioinformatics*, 19:513–523, 2003.
- L. Wan, G. Reinert, F. Sun, and M. Waterman. Alignment-free sequence comparison (II): Theoretical power of comparison statistics. *Journal of Computational Biology*, 17:1467–1490, 2010.
- K. Yang and L. Zhang. Performance comparison between k-tuple distance and four model-based distances in phylogenetic tree reconstruction. *Nucleic acids research*, 36:e33, 2008.

## APPENDIX A: PROOFS FOR SECTION 2

Here we establish Theorem 2.1.

Recalling notation from Section 2,  $\pi = (\pi^w)_{w \in L}$  is the stationary distribution for  $M$ , an  $L \times L$  Markov matrix describing the single-site state change process from sequence  $S_1$  to sequence  $S_2$ . The probability of a  $k$ -mer  $W = w_1 w_2 \dots w_k \in [L]^k$  in any  $k$  consecutive sites of either single sequence is  $\pi^W = \prod_{j=1}^k \pi^{w_j}$ .

Then  $P = \text{diag}(\pi)M$  is the joint distribution of states in aligned sites of the two sequences. We can alternately view the state changes from  $S_2$  to  $S_1$  as described by the Markov matrix  $N = \text{diag}(\pi)^{-1}P^T$ , where T denotes transpose. For future use, note that  $\text{tr } M = \text{tr } N$ , where tr denotes the trace.

Let  $X_{\ell i}^W$  be an indicator variable for the occurrence of a  $k$ -mer  $W$  in sequence  $\ell = 1, 2$  starting at position  $i$ . Then  $X_\ell^W = \sum_{i=1}^{n-k+1} X_{\ell i}^W$  is the count of occurrences of  $k$ -mer  $W$  in sequence  $\ell$ , and  $X_\ell = (X_\ell^W)_{W \in [L]^k}$  is the random vector of  $k$ -mer counts in the sequence.

Let  $Z_i^W = X_{1i}^W - X_{2i}^W$  and  $Z^W = \sum_{i=1}^{n-k+1} Z_i^W = X_1^W - X_2^W$ . These random variables have mean 0.

**Proposition 7.1.** *For  $i \neq j$ ,*

$$\sum_{W \in [L]^k} \frac{1}{\pi^W} \text{Cov}[Z_i^W, Z_j^W] = 0.$$

*Proof.* We may assume  $i < j$ . If  $j - i \geq k$ , the variables  $Z_i^W, Z_j^W$  depend on disjoint sets of sites, hence are independent. This is all that is needed for  $k = 1$ .

We now proceed by induction, assuming the result holds for  $(k - 1)$ -mers, and considering only cases with  $0 < j - i < k$ . Writing a  $k$ -mer  $W$  as a  $(k - 1)$ -mer  $W'$  followed by a 1-mer  $w$ , so that  $\pi^W = \pi^{W'}\pi^w$ , we have

$$\begin{aligned} \sum_{W \in [L]^k} \frac{1}{\pi^W} \text{Cov}[Z_i^W, Z_j^W] &= \\ & \sum_{W' \in [L]^{k-1}} \frac{1}{\pi^{W'}} \sum_{w \in [L]} \frac{1}{\pi^w} \mathbb{E}[Z_i^{W'w} Z_j^{W'w}] \\ (17) \quad &= \sum_{W' \in [L]^{k-1}} \frac{1}{\pi^{W'}} \sum_{w \in [L]} \frac{1}{\pi^w} \mathbb{E}[X_{1i}^{W'w} X_{1j}^{W'w} - X_{1i}^{W'w} X_{2j}^{W'w} - X_{2i}^{W'w} X_{1j}^{W'w} + X_{2i}^{W'w} X_{2j}^{W'w}] \end{aligned}$$

Now since  $j - i < k$ ,

$$X_{1i}^{W'w} X_{1j}^{W'w} = X_{1i}^{W'} X_{1j}^{W'} X_{1(i+k-1)}^w X_{1(j+k-1)}^w$$

and

$$\mathbb{E}[X_{1i}^{W'w} X_{1j}^{W'w}] = \mathbb{E}[X_{1i}^{W'} X_{1j}^{W'}] \delta(u, w) \pi^w,$$

where  $u = w_{k-j+i}$  is the  $(k - j + i)$ th letter in  $W$ , and  $\delta(u, w)$  is the Kronecker delta. Thus

$$\sum_{w \in [L]} \frac{1}{\pi^w} \mathbb{E}[X_{1i}^{W'w} X_{1j}^{W'w}] = \sum_{w \in [L]} \mathbb{E}[X_{1i}^{W'} X_{1j}^{W'}] \delta(u, w) = \mathbb{E}[X_{1i}^{W'} X_{1j}^{W'}].$$

Likewise,  $\sum_{w \in [L]} \frac{1}{\pi^w} \mathbb{E}[X_{2i}^{W'w} X_{2j}^{W'w}] = \mathbb{E}[X_{2i}^{W'} X_{2j}^{W'}]$ .

In a similar way we see

$$X_{2i}^{W'w} X_{1j}^{W'w} = X_{2i}^{W'} X_{1j}^{W'} X_{2(i+k-1)}^w X_{1(j+k-1)}^w$$

and

$$\mathbb{E}[X_{2i}^{W'w} X_{1j}^{W'w}] = \mathbb{E}[X_{2i}^{W'} X_{1j}^{W'}] M(u, w) \pi^w,$$

where  $u = w_{k-j+i}$  is the  $(k - j + i)$ th letter in  $W$ , and  $M$  is the Markov matrix describing the substitution process from sequence 1 to sequence 2. Thus

$$\sum_{w \in [L]} \frac{1}{\pi^w} \mathbb{E}[X_{2i}^{W'w} X_{1j}^{W'w}] = \sum_{w \in [L]} \mathbb{E}[X_{2i}^{W'} X_{1j}^{W'}] M(u, w) = \mathbb{E}[X_{2i}^{W'} X_{1j}^{W'}]$$

and, similarly,  $\sum_{w \in [L]} \frac{1}{\pi^w} \mathbb{E}[X_{1i}^{W'w} X_{2j}^{W'w}] = \mathbb{E}[X_{1i}^{W'} X_{2j}^{W'}]$ .



Combining these expected values with equation (17) we have

$$\begin{aligned} \sum_{W \in [L]^k} \frac{1}{\pi^W} \text{Cov}[Z_i^W, Z_j^W] &= \sum_{W' \in [L]^{k-1}} \frac{1}{\pi^{W'}} \mathbb{E}[X_{1i}^{W'} X_{1j}^{W'} - X_{1i}^{W'} X_{2j}^{W'} - X_{2i}^{W'} X_{1j}^{W'} + X_{2i}^{W'} X_{2j}^{W'}] \\ &= \sum_{W' \in [L]^{k-1}} \frac{1}{\pi^{W'}} \text{Cov}[Z_i^{W'}, Z_j^{W'}] = 0 \end{aligned}$$

by the inductive hypothesis.  $\square$

*Proof of Theorem 2.1.* For  $k = 1$ , using Proposition 4.1 we have

$$\begin{aligned} \mathbb{E} \left[ \sum_w \frac{1}{\pi^w} (X_1^w - X_2^w)^2 \right] &= \sum_w \pi^w \mathbb{E} \left[ \left( \frac{X_1^w}{\pi^w} - \frac{X_2^w}{\pi^w} \right)^2 \right] \\ &= \sum_w \pi^w \frac{n}{(\pi^w)^2} 2(\pi^w - P_{ww}) \\ &= \sum_w 2n(1 - M_{ww}) \\ &= 2n(L - \text{tr } M). \end{aligned}$$

Now inductively suppose the result holds for  $(k-1)$ -mers, and consider  $k$ -mers. Then, since  $Z^W$  has mean zero,

$$\begin{aligned} \mathbb{E} \left[ \sum_W \frac{1}{\pi^W} (X_1^W - X_2^W)^2 \right] &= \sum_W \frac{1}{\pi^W} \mathbb{E} [(Z^W)^2] \\ &= \sum_W \frac{1}{\pi^W} \text{Var} [Z^W] \\ &= \sum_W \frac{1}{\pi^W} \text{Var} \left[ \sum_i Z_i^W \right] \\ &= \sum_W \frac{1}{\pi^W} \left( \sum_i \text{Var} [Z_i^W] + \sum_{i \neq j} \text{Cov} [Z_i^W, Z_j^W] \right) \\ &= \sum_W \frac{1}{\pi^W} \left( \sum_i \text{Var} [Z_i^W] \right). \end{aligned}$$

Here Proposition 7.1 justifies the last equality. Now since  $(Z_i^W)^2$  is the indicator variable for when exactly one of  $X_{1i}^W, X_{2i}^W$  is 1,

$$\text{Var}[Z_i^W] = \mathbb{E}[(Z_i^W)^2] = \pi^W \left( 1 - \prod_{j=1}^k M(w_j, w_j) \right) + \pi^W \left( 1 - \prod_{j=1}^k N(w_j, w_j) \right).$$

Thus

$$\begin{aligned} \mathbb{E} \left[ \sum_W \frac{1}{\pi^W} (X_1^W - X_2^W)^2 \right] &= (n - k + 1) \sum_W \left( 2 - \prod_{j=1}^k M(w_j, w_j) - \prod_{j=1}^k N(w_j, w_j) \right) \\ &= (n - k + 1) (2L^k - (\text{tr } M)^k - (\text{tr } N)^k) \\ &= 2(n - k + 1) (L^k - (\text{tr } M)^k). \end{aligned}$$

□

## APPENDIX B: PROOFS FOR SECTION 4

We establish Proposition 4.1. Our proof is independent of earlier arguments, as the result is needed in Appendix A.

Sequences  $S_1$  and  $S_2$  each have i.i.d. sites with state probabilities given by  $\pi_1$  and  $\pi_2$ , and site transition probabilities from  $S_1$  to  $S_2$  are given by the matrix  $M$ . Note that the  $\pi_\ell$  need not be stationary vectors for  $M$ .

As in Appendix A, define random variables  $X_{\ell k}^w$  for  $w \in [L]$ ,  $\ell \in \{1, 2\}$ ,  $j \in [n]$  to be indicators of state  $w$  in sequence  $\ell$  at site  $j$ . The 1-mer distribution vector for the sequence  $S_\ell$  is then  $X_\ell$  with entries  $X_\ell^w = \sum_{j=1}^n X_{\ell j}^w$ .

*Proof of Proposition 4.1.* That  $\pi_\ell = \frac{1}{n} \mathbb{E}(X_\ell)$  is clear.

Since  $\mathbb{E} \left[ \frac{X_1^u}{\pi_1^u} - \frac{X_2^w}{\pi_2^w} \right] = 0$ ,

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{X_1^u}{\pi_1^u} - \frac{X_2^w}{\pi_2^w} \right)^2 \right] &= \text{Var} \left[ \frac{X_1^u}{\pi_1^u} - \frac{X_2^w}{\pi_2^w} \right] = \text{Var} \left[ \sum_{j=1}^n \left( \frac{X_{1j}^u}{\pi_1^u} - \frac{X_{2j}^w}{\pi_2^w} \right) \right] \\ (18) \qquad \qquad \qquad &= \sum_{j=1}^n \text{Var} \left[ \frac{X_{1j}^u}{\pi_1^u} - \frac{X_{2j}^w}{\pi_2^w} \right] = n \cdot \text{Var} \left[ \frac{X_{11}^u}{\pi_1^u} - \frac{X_{21}^w}{\pi_2^w} \right] \end{aligned}$$

by the i.i.d. assumption. But

$$\begin{aligned} \text{Var} \left[ \frac{X_{11}^u}{\pi_1^u} - \frac{X_{21}^w}{\pi_2^w} \right] &= \mathbb{E} \left[ \left( \frac{X_{11}^u}{\pi_1^u} - \frac{X_{21}^w}{\pi_2^w} \right)^2 \right] = \mathbb{E} \left[ \frac{(X_{11}^u)^2}{(\pi_1^u)^2} + \frac{(X_{21}^w)^2}{(\pi_2^w)^2} - 2 \frac{X_{11}^u X_{21}^w}{\pi_1^u \pi_2^w} \right] \\ &= \mathbb{E} \left[ \frac{X_{11}^u}{(\pi_1^u)^2} + \frac{X_{21}^w}{(\pi_2^w)^2} - 2 \frac{X_{11}^u X_{21}^w}{\pi_1^u \pi_2^w} \right]. \end{aligned}$$

Since  $\mathbb{E}[X_{\ell 1}^u] = \pi_\ell^u$  and  $\mathbb{E}[X_{11}^u X_{21}^w] = P_{uw}$ , this shows

$$(19) \qquad \qquad \qquad \text{Var} \left[ \frac{X_{11}^u}{\pi_1^u} - \frac{X_{21}^w}{\pi_2^w} \right] = \frac{1}{\pi_1^u} + \frac{1}{\pi_2^w} - 2 \frac{P_{uw}}{\pi_1^u \pi_2^w}.$$

Substituting equation (19) into equation (18) and solving for  $P_{uw}$  completes the proof. □

To establish Proposition 4.3, recall that for  $\ell = 1, 2, 3$  we consider sequences  $S_\ell$  with 1-mer count vectors  $X_\ell$  and base distribution vector  $\pi_\ell = \mathbb{E}(X_\ell)$ . Let

$$Y_{uvw} = \alpha X_1^u + \beta X_2^v + \gamma X_3^w,$$

where  $\alpha, \beta, \gamma$  are constants chosen so

$$\alpha\pi_1^u + \beta\pi_2^v + \gamma\pi_3^w = 0.$$

*Proof of Proposition 4.3.* Note  $\mathbb{E}(Y_{uvw}) = 0$ . Using the fact that the 3rd central moment is additive over independent variables, and that sites are identically distributed

$$\begin{aligned} \mathbb{E}(Y_{uvw}^3) &= \mathbb{E}\left(\left(\sum_{i=1}^n (\alpha X_{1i}^u + \beta X_{2i}^v + \gamma X_{3i}^w)\right)^3\right) \\ &= \sum_{i=1}^n \mathbb{E}\left((\alpha X_{1i}^u + \beta X_{2i}^v + \gamma X_{3i}^w)^3\right) = n \cdot \mathbb{E}\left((\alpha X_{11}^u + \beta X_{21}^v + \gamma X_{31}^w)^3\right), \end{aligned}$$

where  $n$  is the sequence length. But, since  $(X_{\ell 1}^u)^2 = X_{\ell 1}^u$ ,

$$\begin{aligned} \mathbb{E}\left((\alpha X_{1\ell}^u + \beta X_{2\ell}^v + \gamma X_{3\ell}^w)^3\right) &= \mathbb{E}\left(\alpha^3 X_{11}^u + \beta^3 X_{21}^v + \gamma^3 X_{31}^w \right. \\ &\quad \left. + 3(\alpha^2\beta + \alpha\beta^2)X_{11}^u X_{21}^v + 3(\alpha^2\gamma + \alpha\gamma^2)X_{11}^u X_{31}^w \right. \\ &\quad \left. + 3(\beta^2\gamma + \beta\gamma^2)X_{21}^v X_{31}^w + 6\alpha\beta\gamma X_{11}^u X_{21}^v X_{31}^w\right) \\ &= \alpha^3\pi_1^u + \beta^3\pi_2^v + \gamma^3\pi_3^w - 3\alpha\beta(\alpha + \beta)\mathbb{E}(X_{11}^u X_{21}^v) \\ &\quad - 3\alpha\gamma(\alpha + \gamma)\mathbb{E}(X_{11}^u X_{31}^w) - 3\beta\gamma(\beta + \gamma)\mathbb{E}(X_{21}^v X_{31}^w) \\ &\quad + 6\alpha\beta\gamma\mathbb{E}(X_{11}^u X_{21}^v X_{31}^w). \end{aligned}$$

Using  $\mathbb{E}(X_{11}^u X_{21}^v) = P_{uv+}$  and variants, and  $\mathbb{E}(X_{11}^u X_{21}^v X_{31}^w) = P_{uvw}$ , this shows

$$\begin{aligned} \mathbb{E}(Y_{uvw}^3) &= n(\alpha^3\pi_1^u + \beta^3\pi_2^v + \gamma^3\pi_3^w - 3\alpha\beta(\alpha + \beta)P_{uv+} - 3\alpha\gamma(\alpha + \gamma)P_{u+w} \\ &\quad - 3\beta\gamma(\beta + \gamma)P_{+vw} + 6\alpha\beta\gamma P_{uvw}), \end{aligned}$$

and the claim readily follows.  $\square$

DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF ALASKA FAIRBANKS, 99775  
E-mail address: e.allman@alaska.edu

DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF ALASKA FAIRBANKS, 99775  
E-mail address: j.rhodes@alaska.edu

DEPARTMENT OF MATHEMATICS, NORTH CAROLINA STATE UNIVERSITY, RALEIGH, NC, 27695  
E-mail address: smsulli2@ncsu.edu