

Embeddability of Kimura 3ST Markov matrices

Jordi Roca-Lacostena and Jesús Fernández-Sánchez

February 25, 2019

Abstract

In this note, we characterize the embeddability of generic Kimura 3ST Markov matrices in terms of their eigenvalues. As a consequence, we are able to compute the volume of such matrices relative to the volume of all Markov matrices within the model. We also provide examples showing that, in general, mutation rates are not identifiable from substitution probabilities. These examples also illustrate that symmetries between mutation probabilities do not necessarily arise from symmetries between the corresponding mutation rates.

Keywords: Markov matrix; Markov generator; eigenvalues; evolutionary model; embeddability

1 Introduction

Genomic data expressed by means of sequence alignments is widely used to infer phylogenetic relationships between species. Substitution models are used to describe the evolutionary process that leads from one DNA sequence to another. These models are usually given in terms of a family of Markov matrices with a prescribed structure. The entries of these matrices represent the conditional probabilities of nucleotide substitution between one sequence and the other, and can be obtained either by counting the relative frequencies of these substitutions or fitting the parameters of the model using maximum likelihood. Usually, the structure imposed by the model is motivated by some biological / biochemical properties observed (e.g. the Kimura 3ST model [Kim81]) or some computational / mathematical convenient assumptions to deal with the model (e.g. the GTR model [Tav86] or Lie Markov models [SFSJ12]). Moreover, evolution is usually modelled by means of Markov chains, together with the additional assumption that all sites in the sequences evolve independently and according to the same probabilities.

A general approach in modelling evolution corresponds to regarding time as a continuous variable where substitution events always happen at the same rate, which remains constant throughout the whole evolutionary process. This leads to the homogeneous continuous-time substitution models, where only Markov matrices that are the exponential of a rate matrix are considered. Clearly this is used as an approximation to biological reality where it is well known that transition rates vary over time [HPCD05, HSP⁺07] and also among the different branches of the phylogenetic tree [LSB⁺98]. However, given the bias / variance compensation of the statistical analysis [BA02], modelling phylogenetic evolution as a non-homogeneous process is not statistically feasible in practice (cf. [SFSJ12]).

A different approach appears when one regards the evolutionary process as a whole and only takes into account the conditional probabilities between the original and the final sequences, without caring about rates of mutation. When these probabilities are taken as the parameters of the model, we deal with the so-called *algebraic* models¹. Algebraic models have been used in a number of theoretical papers, including [AR08, SS05, DK08, CFS10].

If one attempts to connect both approaches, a natural question is to decide whether a given Markov matrix is the exponential of some rate matrix, whose entries would be some kind of average of the rates involved throughout the evolutionary process. In this case, we say the matrix is *embeddable* and this question is known in the literature as the *embedding problem* for Markov matrices. An easier version of

¹Here, “algebraic” refers to the fact that the probabilities of pattern observation at the leaves of a phylogenetic tree evolving under these models are given by algebraic expressions (only sums and products) in terms of the parameters of the model.

this problem is to decide whether the rate matrices associated to the embeddable matrices of a particular (algebraic) model \mathcal{M} should keep the same symmetries as the model (\mathcal{M} -embeddability, see definition in Section 2.2). The embedding problem is relevant even if restricted to continuous-time models since it is not true in general that the product of embeddable matrices is necessarily embeddable (indeed, the Baker-Campbell-Hausdorff formula [Cam97] leads to ask whether some series of matrices is convergent or not, which is not always true [BC04]). These questions are closely related to the problem of the multiplicative closure of continuous-time models, namely whether the product of matrices $e^{Q_1}e^{Q_2}$ where Q_1 and Q_2 are rate matrices in one particular (continuous-time) model can be obtained as some e^Q for some rate matrix Q in *the same model*. After [SFSJ12, Sum17], it is known that there are popular models which are not multiplicatively closed, notably including the GTR model and the HKY model.

The reader is referred to [Dav10] for a nice overview of the embedding problem from a mathematical point of view. In a more biological and applied setting, the paper by Verbyla *et al.* [VYP⁺13] deals with the possible consequences for phylogenetic inference. Also, the paper [SJFS⁺12] and the more recent paper [WSL⁺17] deal with the incidental question of how the lack of (multiplicative) closure in substitution models have consequences for the phylogenetic analysis of data.

In this paper, we deal with the embedding problem from a theoretical perspective. The main goal is to obtain a characterization for the embeddability of generic matrices of the Kimura 3ST model [Kim81]. From our results, we will be able to compute the whole volume of embeddable Kimura 3ST matrices and compare it with the volume of the whole space of Kimura 3ST Markov matrices. At the same time, we provide a number of examples showing matrices that are embeddable but for which the mutation rates are not identifiable or do not keep the same structure of the model. The recent paper [KK17] deals with the similar question of characterizing embeddable matrices of symmetric group-based phylogenetic models, but focusing on the existence of rate matrices strictly in the model.

The organization of the paper is as follows. In section 2, we recall some definitions and basic facts concerning the embedding problem and the Kimura 3 parameter model. Here, we also show that any embeddable matrix is biologically relevant since it can be seen as the transition matrix of a concatenation of *realistic* evolutionary processes (“realistic” here means a process whose transition matrix is close to the identity matrix, see Theorem 2.2). In section 3, we prove the main theorem which characterizes under the (generic) assumption of having different eigenvalues the Kimura 3ST embeddable matrices in terms of inequalities to be satisfied by the eigenvalues. We devote as well some attention to the case of matrices with repeated eigenvalues as they present certain situations that may be interesting from a theoretical and applied point of view. Namely, these matrices show that the identifiability of the mutation rates is not a generic property for the Kimura 2ST model or the Jukes-Cantor model, as well as that there are embeddable matrices with rate matrices that do not keep the same symmetries of the model (see Theorem 3.9). As a consequence of the characterization mentioned above, in section 4 we are able to compute the volume of embeddable matrices and compare it to the volume of all Kimura 3ST Markov matrices. Finally, Section 5 discusses implications and possibilities for future work.

2 Preliminaries

2.1 Embedding problem of Markov matrices

We denote by $M_k(\mathbb{K})$ the space of all square k -matrices with entries in a field \mathbb{K} , where \mathbb{K} is \mathbb{R} or \mathbb{C} . Given a matrix $A \in M_k(\mathbb{K})$, we say that $B \in M_k(\mathbb{K})$ is a *logarithm* of A if $e^B = A$, where the exponential of a matrix is defined as

$$e^X = \sum_{n \geq 0} \frac{X^n}{n!}.$$

A classical result states that $\det(e^X) = e^{\text{tr}(X)}$, so the determinant of any matrix of the form e^X is never 0. Given a non-negative complex number $x \in \mathbb{C} \setminus \mathbb{R}^-$, we will denote by $\log(x)$ its *principal logarithm*, that is, the only logarithm of x that lies in the strip $\{z \mid -\pi < \text{Im}(z) < \pi\}$. Although the exponential map of matrices is not injective, it is known that if A is a matrix with no negative eigenvalues, there is a unique logarithm X of A all of whose eigenvalues are given by the principal logarithm of the eigenvalues

of A (Theorem 1.31 of [Hig08]). We will refer to this as the *principal logarithm of A* and we will denote it by $\text{Log}(A)$. In the particular case where the matrix A is diagonalizable, $A = SDS^{-1}$ then $\text{Log}(A) = S\text{Log}(D)S^{-1}$, where $\text{Log}(D)$ is the diagonal matrix with diagonal entries equal to the principal logarithm of the eigenvalues of A .

Definition 2.1. A matrix $M \in M_k(\mathbb{R})$ is said to be a *Markov matrix* if all the entries are non-negative and the rows sum to one. A matrix $Q \in M_k(\mathbb{R})$ is said to be a *rate matrix* if all the non-diagonal entries are non-negative and the rows sum to zero.

If Q is a rate matrix, it is well-known that $e^{tQ} = \sum_{n \geq 0} \frac{t^n Q^n}{n!}$ is a Markov matrix for all $t \geq 0$. That is why rate matrices are also referred as *Markov generators* [Dav10]. However, not every Markov matrix can be obtained in this way. A Markov matrix M is said to be *embeddable* if $M = e^Q$ for some rate matrix Q . The *embedding problem* attempts to decide which (Markov) matrices are embeddable, that is, which matrices can be written as $M = e^Q$, where Q is a rate matrix. We would like to point out that every embeddable matrix can be obtained as the substitution matrix of a long-running biologically realistic Markov process. Namely,

Theorem 2.2. *Every embeddable matrix is the product of embeddable matrices close to the identity matrix.*

Proof. Assume that M is an embeddable Markov matrix: $M = e^Q$. Clearly, $Q_n := \frac{1}{n}Q$ is still a rate matrix for any $n \geq 1$, so $M = (e^{Q_n})^n$ appears as the n -th power of a Markov matrix. Moreover, since

$$\lim_{n \rightarrow \infty} e^{Q_n} = e^{\lim_{n \rightarrow \infty} Q_n} = e^{(0)} = Id,$$

we can take n big enough so that e^{Q_n} is as close to Id as wanted. □

2.2 Kimura models

In this work we deal with the substitution model introduced by Kimura in [Kim81]. The Kimura 3ST model assigns three parameters to different type of substitutions: one parameter for transitions, i.e. substitutions between purines ($A \leftrightarrow G$) or pyrimidines ($C \leftrightarrow T$), and two parameters for transversions, i.e. substitutions that change the type of nucleotide: from purine to pyrimidine or vice versa. Ordering the set of nucleotides as A, G, C, T , the Markov matrices within the model are described by the following structure:

Definition 2.3. A matrix $M \in M_4(\mathbb{C})$ is *Kimura 3ST* (is K3 or has K3 form, for short) if it has the following structure:

$$M = \begin{pmatrix} a & b & c & d \\ b & a & d & c \\ c & d & a & b \\ d & c & b & a \end{pmatrix}. \tag{1}$$

For ease of reading we will use the notation $M = K(a, b, c, d)$ to denote a matrix with the structure in (1).

When M is a Markov matrix, the structure above describes the symmetry between the substitution probabilities of the Kimura 3ST model. Keeping the order of the nucleotides, if $i, j = A, G, C, T$, the (i, j) -entry corresponds to the probability of nucleotide i being replaced by nucleotide j . The submodels of Kimura 3ST model, namely Kimura 2ST [Kim80] and Jukes-Cantor [JC69], appear when more symmetries are considered: if $c = d$, we say that the matrix M is *Kimura 2ST* (K2, for short); if $b = c = d$, we say that M is *Jukes-Cantor* (JC, for short).

If one restricts the embedding problem to one particular model \mathcal{M} given by some equalities between the entries of the Markov matrices (such as in Kimura 3ST, Kimura 2ST, Jukes-Cantor), it is natural to ask whether these matrices have Markov generators fulfilling the same equalities. In this case, we say

that these matrices are \mathcal{M} -embeddable. Example 3.6 of the next section shows that this is not true in general.

The following lemma is fundamental for our study. It essentially claims that all K3 matrices are diagonalized through the following Hadamard matrix:

$$S := \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}.$$

Note that $S^2 = 4 \cdot Id$, thus $S^{-1} = \frac{1}{4}S$.

Lemma 2.4. *A matrix is K3 if and only if it can be diagonalized through S . In this case, $K(a, b, c, d) = S D S^{-1}$, where $D = \text{diag}(a + b + c + d, a + b - c - d, a - b + c - d, a - b - c + d)$. In particular, a K3 matrix is real if and only if its eigenvalues are all real numbers.*

Proof. The proof is straightforward and follows by direct computation. \square

Since the rows of K3 Markov matrices sum to $a + b + c + d = 1$, the first eigenvalue must be equal to one. Hence, we derive that every K3 Markov matrix is determined by the set of the other eigenvalues:

$$x := a + b - c - d, \quad y := a - b + c - d, \quad z := a - b - c + d. \quad (2)$$

Moreover, it is immediate to check that a matrix M as in (1) is K2 (resp. JC) if and only if $y = z$ (resp. $x = y = z$).

It also follows that

Theorem 2.5. *The product of K3-embeddable matrices is K3-embeddable.*

Proof. By virtue of Lemma 2.4, we have

$$K(a, b, c, d) \cdot K(a', b', c', d') = (S D S^{-1}) \cdot (S D' S^{-1}) = S D D' S^{-1},$$

which is a K3 matrix. Since D and D' are diagonal matrices, we have that $D D' = D' D$ and from this, the product of K3 matrices is commutative. Therefore, if Q_1, Q_2 are K3 rate matrices, the Baker-Campbell-Hausdorff formula [Cam97] gives

$$e^{Q_1} e^{Q_2} = \exp\left(Q_1 + Q_2 + \frac{1}{2}[Q_1, Q_2] + \dots\right) = \exp(Q_1 + Q_2),$$

where $[A, B] = AB - BA$ is the Lie bracket. Since $Q_1 + Q_2$ is also a K3 rate matrix, the claim follows. \square

3 Embeddability of Kimura Markov matrices

The first result of this section describes how to compute the principal logarithm of a K3 Markov matrix and characterizes when it is a Markov generator.

First, we need a lemma which follows from the definition of the exponential matrix.

Lemma 3.1. *Let Q be a logarithm of M . If v is an eigenvector of Q with eigenvalue μ , then v is an eigenvector of M with eigenvalue e^μ .*

Remark 3.2. Note that the converse is not true in general. For instance, any vector of \mathbb{C}^2 is an eigenvector of $M := \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix}$ with eigenvalue -1 . However, the matrix $Q = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$ is a logarithm of M , with eigenvalues i and $-i$ and corresponding eigenspaces $[(i, 1)]$ and $[(-i, 1)]$. In particular, any vector not in these subspaces cannot be an eigenvector of Q .

The following result solves the K3-embeddability.

Theorem 3.3. *Let M be a K3 Markov matrix with eigenvalues $1, x, y, z$. Then,*

i) M has a real logarithm Q with K3 form if and only if $x, y, z > 0$. In this case, Q is necessarily the principal logarithm $\text{Log}(M) = S \text{diag}(0, \log(x), \log(y), \log(z)) S^{-1}$.

ii) M is K3-embeddable (i.e. $\text{Log}(M)$ is a Markov generator) if and only if

$$x \geq yz, \quad y \geq xz, \quad z \geq xy. \quad (3)$$

Proof. (i) First of all, if Q is a real logarithm of M with K3 form, then the eigenvalues of Q are real (Lemma 2.4) and, by Lemma 3.1, the eigenvalues of M have to be positive. Conversely, if $x, y, z > 0$ then we can take principal logarithms of x, y, z and define Q as the principal logarithm of M , i.e. $Q := \text{Log}(M) = S \text{diag}(0, \log(x), \log(y), \log(z)) S^{-1}$. Then, Q is a real matrix, $e^Q = M$ and it has K3 form because of Lemma 2.4.

We proceed to show that $\text{Log}(M)$ is the only possible real logarithm with K3 form. By Lemma 2.4, if Q' is any other real logarithm with K3 form, then it can be written as $Q' = S \text{diag}(r, s, u, v) S^{-1}$, for some $r, s, u, v \in \mathbb{R}$. Because of Lemma 3.1, these values are real logarithms of $1, x, y, z$. Necessarily, Q' must be the principal logarithm of M .

(ii) Using Lemma 2.4 we have $\text{Log}(M) = K(\alpha, \beta, \gamma, \delta)$, where $\alpha = \frac{1}{4}(\log(x) + \log(y) + \log(z))$, $\beta = \frac{1}{4}(\log(x) - \log(y) - \log(z))$, $\gamma = \frac{1}{4}(\log(y) - \log(x) - \log(z))$ and $\delta = \frac{1}{4}(\log(z) - \log(x) - \log(y))$. It is immediate that $\alpha + \beta + \gamma + \delta = 0$. Therefore, we only need to check that the non-diagonal entries β, γ and δ of $\text{Log}(M)$ are non-negative if and only if the inequalities (3) are satisfied. This is straightforward: for instance,

$$\beta \geq 0 \Leftrightarrow \frac{\log(x) - \log(y) - \log(z)}{4} \geq 0 \Leftrightarrow \log\left(\frac{x}{yz}\right) \geq 0 \Leftrightarrow x \geq yz.$$

The other inequalities are proved similarly. □

Lemma 3.4. *If M is a K3 matrix with no repeated eigenvalues and Q is a real logarithm of M , then Q is the principal logarithm of M . In particular, Q has K3 form.*

Proof. Because of Lemma 3.1, if the matrix M has no repeated eigenvalues, then so does the matrix Q . It follows that Q diagonalizes, and that M and Q have the same eigenvectors. In particular they both diagonalize through the matrix S and hence Q must be K3 (Lemma 2.4). Now, it is enough to apply Theorem 3.3. □

Corollary 3.5. *Let M be a K3 Markov matrix with no repeated eigenvalues. Then, the following are equivalent:*

(i) M is embeddable;

(ii) M is K3-embeddable;

(iii) the eigenvalues x, y, z of M are strictly positive, and satisfy

$$x \geq yz, \quad y \geq xz, \quad z \geq xy.$$

Proof. It follows directly from theorem 3.3 and Lemma 3.4. □

The case of repeated eigenvalues

After the previous result, it is natural to ask what can be said in the case of repeated eigenvalues. In this case, there are some theoretically interesting examples showing that:

1. There are embeddable K3 matrices with repeated *negative* eigenvalues. Remarkably, these matrices do not admit K3 Markov generators and their rates are not identifiable. See forthcoming Example 3.6.

2. Restricted to the case of repeated positive eigenvalues, there are K3 Markov matrices for which rates are not identifiable.

Although the matrices presented in the forthcoming examples are close to saturation and have a big mutation rate, they have still biological interest by virtue of Theorem 2.2.

Repeated negative eigenvalues It is well known that a matrix with negative eigenvalues has a real logarithm if and only if the negative eigenvalues have even multiplicity [Cul66]. Consider a K3 Markov matrix (actually, it is a K2 matrix) with eigenvalues 1, $e^{-\lambda}$ and $-e^{-\mu}$ with multiplicity 2, where $\lambda, \mu \geq 0$:

$$M = \frac{1}{4} \begin{pmatrix} 1 + e^{-\lambda} - 2e^{-\mu} & 1 + e^{-\lambda} + 2e^{-\mu} & 1 - e^{-\lambda} & 1 - e^{-\lambda} \\ 1 + e^{-\lambda} + 2e^{-\mu} & 1 + e^{-\lambda} - 2e^{-\mu} & 1 - e^{-\lambda} & 1 - e^{-\lambda} \\ 1 - e^{-\lambda} & 1 - e^{-\lambda} & 1 + e^{-\lambda} - 2e^{-\mu} & 1 + e^{-\lambda} + 2e^{-\mu} \\ 1 - e^{-\lambda} & 1 - e^{-\lambda} & 1 + e^{-\lambda} + 2e^{-\mu} & 1 + e^{-\lambda} - 2e^{-\mu} \end{pmatrix} \quad (4)$$

By virtue of Theorem 3.3, M has no real logarithm with K3 or K2 form.

According to Theorem 1.28 in [Hig08], distinct logarithms for the matrix M are obtained as $\overline{S}D\overline{S}^{-1}$, where D is a diagonal matrix with distinct determinations of the logarithm of the eigenvalues of M and the columns of \overline{S} form a basis of eigenvectors of M (see also Chap. VIII§8 in [Gan59]). Hence, if we choose a pair of conjugated eigenvectors v, \bar{v} of the eigenvalue $-e^{-\mu}$ as the third and fourth columns of the matrix \overline{S} , then all the matrices of the form

$$Q_k = \overline{S} \operatorname{diag}(0, -\lambda, -\mu + (2k+1)\pi i, -\mu - (2k+1)\pi i) \overline{S}^{-1}, \quad k \in \mathbb{Z}$$

are real logarithms of M . For example, we can take

$$\overline{S} = \begin{pmatrix} 1 & 1 & 1+i & 1-i \\ 1 & 1 & -1-i & -1+i \\ 1 & -1 & 1-i & 1+i \\ 1 & -1 & -1+i & -1-i \end{pmatrix}$$

to obtain

$$Q_k = \frac{1}{4} \begin{pmatrix} -\lambda - 2\mu & -\lambda + 2\mu & \lambda - 2\pi(2k+1) & \lambda + 2\pi(2k+1) \\ -\lambda + 2\mu & -\lambda - 2\mu & \lambda + 2\pi(2k+1) & \lambda - 2\pi(2k+1) \\ \lambda + 2\pi(2k+1) & \lambda - 2\pi(2k+1) & -\lambda - 2\mu & -\lambda + 2\mu \\ \lambda - 2\pi(2k+1) & \lambda + 2\pi(2k+1) & -\lambda + 2\mu & -\lambda - 2\mu \end{pmatrix}.$$

It is straightforward to check that this is a Markov generator if and only if $2\pi|2k+1| \leq \lambda$ and $\lambda \leq 2\mu$. In particular, we see that M is embeddable if

$$2\pi \leq \lambda \leq 2\mu. \quad (5)$$

Moreover, in this case, it is enough to take $k = 0$ and $k = -1$ to obtain a pair of Markov generators for M .

We illustrate this construction with a numerical example.

Example 3.6. Let us take $\lambda = 7$ and $\mu = 4$, so that $2\pi \leq \lambda \leq 2\mu$. Then the matrix M is (rounding off to 7 decimals)

$$\begin{pmatrix} 0.2410701 & 0.2593858 & 0.2497720 & 0.2497720 \\ 0.2593858 & 0.2410701 & 0.2497720 & 0.2497720 \\ 0.2497720 & 0.2497720 & 0.2410701 & 0.2593858 \\ 0.2497720 & 0.2497720 & 0.2593858 & 0.2410701 \end{pmatrix}. \quad (6)$$

As mentioned above, $\operatorname{Log}(M)$ is not a real matrix and hence is not a rate matrix either. In spite of that we are still able to find a pair of Markov generators for M by taking $k = 0$ and $k = -1$ respectively:

$$Q_0 = \frac{1}{4} \begin{pmatrix} -15 & 1 & 7 - 2\pi & 7 + 2\pi \\ 1 & -15 & 7 + 2\pi & 7 - 2\pi \\ 7 + 2\pi & 7 - 2\pi & -15 & 1 \\ 7 - 2\pi & 7 + 2\pi & 1 & -15 \end{pmatrix} \quad Q_{-1} = \frac{1}{4} \begin{pmatrix} -15 & 1 & 7 + 2\pi & 7 - 2\pi \\ 1 & -15 & 7 - 2\pi & 7 + 2\pi \\ 7 - 2\pi & 7 + 2\pi & -15 & 1 \\ 7 + 2\pi & 7 - 2\pi & 1 & -15 \end{pmatrix}$$

The previous example shows that a K3 Markov matrix M can be embeddable, even if it does not have any Markov generator with K3 form (see Theorem 3.3 (i)). Thus, embeddability of K3 matrices does not imply K3-embeddability, and is not determined by the principal logarithm (cf. [VYP⁺13, KK17]). This fact exhibits that the structure of the K3 model, which imposes certain symmetries between transitions and between transversions, is not always captured by the same symmetries between the mutation rates (cf. [Kim80, Kim81]). The reader may note that the expected number of substitutions for a process ruled by a matrix M as in (4) is $-\frac{1}{4}\text{tr}(Q_k) = \frac{1}{4}(\lambda + 2\mu) \geq \pi$.

Remark 3.7. After the previous example, we derive easily that the product of embeddable matrices within the Kimura 3ST model is not necessarily embeddable (cf. Theorem 2.5). Indeed, it is enough to consider the embeddable matrix M shown in (6) and any K3 matrix N with positive eigenvalues $1, x, y, z$ satisfying the inequalities (3) so that N is embeddable. The product of M and N is clearly a K3 Markov matrix, whose eigenvalues are the product of the eigenvalues of M and N . Thus, MN has two *different negative* eigenvalues and another positive eigenvalue. By virtue of [Cul66], MN has no real logarithm so, in particular, it cannot be embeddable.

Repeated positive eigenvalues A similar computation to that for negative eigenvalues can be done by assuming positive and repeated eigenvalues: $1, e^{-\lambda}$ and $e^{-\mu}$ with multiplicity 2. In this case, we can produce real logarithms by taking

$$\begin{aligned} Q_k &= \overline{S} \text{diag}(0, -\lambda, -\mu + 2k\pi i, -\mu - 2k\pi i) \overline{S}^{-1} = \\ &= \frac{1}{4} \begin{pmatrix} -\lambda - 2\mu & -\lambda + 2\mu & \lambda - 4\pi k & \lambda + 4\pi k \\ -\lambda + 2\mu & -\lambda - 2\mu & \lambda + 4\pi k & \lambda - 4\pi k \\ \lambda + 4\pi k & \lambda - 4\pi k & -\lambda - 2\mu & -\lambda + 2\mu \\ \lambda - 4\pi k & \lambda + 4\pi k & -\lambda + 2\mu & -\lambda - 2\mu \end{pmatrix}. \end{aligned}$$

Such a matrix is a Markov generator if and only if $4\pi|k| \leq \lambda$ and $\lambda \leq 2\mu$. In particular, we see that if $0 \leq \lambda \leq 2\mu$, then M is K2-embeddable since Q_0 is a Markov generator. If in addition,

$$4\pi \leq \lambda \leq 2\mu, \tag{7}$$

then we have (at least) three Markov generators, which correspond to $k = 0, \pm 1$.

Note that the Markov generator Q_0 is a K3 matrix (corresponding to the principal logarithm) while Q_1 and Q_{-1} are not. The expected number of substitutions for a process ruled by these matrices is $-\frac{1}{4}\text{tr}(Q_k) = \frac{1}{4}(\lambda + 2\mu) \geq 2\pi$.

Remark 3.8. Following the construction of Theorem 2.2, for any $n \geq 1$ the matrix $e^{(1/n)Q_0}$ is a K2 Markov matrix that approaches to Id when n grows. If a substitution process is ruled by such a matrix for some long time $t > 0$, the rates of the resulting Markov matrix $e^{t(Q_0/n)}$ become unidentifiable at some point. This situation cannot occur for K3 embeddable matrices with different eigenvalues (see Theorem 3.3).

The existence of two or more Markov generators for the same Markov matrix (both in the case of negative and positive eigenvalues) exhibit that, in general, mutation rates are not identifiable from the mutation probabilities. Even more, if we restrict to the Kimura 3ST submodels, such matrices do not appear as marginal cases. Indeed, we have seen that identifiability of rates of a K2 embeddable matrix M with eigenvalues $1, e^{-\lambda}, e^{-\mu}$ (with multiplicity 2) does not hold in the subspace defined by the inequalities (7), which has positive measure within the space of all K2 Markov matrices. Furthermore, we have also seen that those Markov matrices with a negative repeated eigenvalue $1, e^{-\lambda}, -e^{-\mu}$ (with multiplicity 2) satisfying (5) are embeddable but not K2-embeddable (nor K3-embeddable). The space of such matrices has also positive measure within space of K2 Markov matrices.

Similarly, for a Jukes-Cantor matrix with eigenvalues 1 and $e^{-\lambda}$ (with multiplicity 3), identifiability of rates does not hold in the subspace defined by $4\pi \leq \lambda$, which has positive measure within the space of all JC matrices. On the other hand, every embeddable matrix is JC-embeddable since by [Cul66], a

necessary condition for a Markov Jukes-Cantor matrix to be embeddable is that its eigenvalues $1, x$ are positive, and then Theorem 3.3 ensures that the principal logarithm (which is a JC matrix) is a Markov generator (it is enough to check that $x \geq x^2$, which is true since $x \in [0, 1]$).

To conclude this section, we state the following theorem which summarizes the consequences of the previous examples:

- Theorem 3.9.** *1. For the Kimura 3ST model, K3-embeddability and the identifiability of rates are generic properties of embeddable matrices.²*
- 2. For the Kimura 2ST model, K2-embeddability and the identifiability of rates are not generic properties of embeddable matrices.*
- 3. For the Jukes-Cantor model, every embeddable matrix is JC-embeddable but identifiability of rates is not a generic property of embeddable matrices.*

4 The volume of embeddable K3 matrices

Roughly speaking, the goal of this section is to measure how many K3 Markov matrices we are considering when the continuous-time approach is taken and compare this value with the corresponding value of K3 Markov matrices with no further restriction. These values are expressed in terms of the volume of the corresponding subspaces. At the same time, it is direct to obtain the volume of subspaces of K3 Markov matrices with some constraints to make them biologically realistic. To this aim, we proceed to represent K3 Markov matrices in a geometrical way as follows (cf. [CFS08]): keeping the notation introduced in (2), Lemma 2.4 allows us to identify the K3 Markov matrices with the coordinates (x, y, z) of a 3-dimensional space. Moreover, since every K3 matrix is a convex combination of the identity matrix and permutation matrices:

$$M = a \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + b \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix} + c \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} + d \begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}, \quad \begin{array}{l} a, b, c, d \geq 0, \\ a + b + c + d = 1 \end{array}$$

the space of all K3 Markov matrices describes the 3-dimensional simplex (a regular tetrahedron) with vertices given by the corresponding eigenvalues: $p_1 = (1, 1, 1)$, $p_2 = (1, -1, -1)$, $p_3 = (-1, 1, -1)$ and $p_4 = (-1, -1, 1)$ (see Figure 1). The centroid of this simplex has coordinates (eigenvalues) $O = (0, 0, 0)$ and corresponds to the matrix

$$M = \frac{1}{4} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}.$$

According to this representation, the Jukes-Cantor matrices [JC69] ($b = c = d$) correspond to the line determined by the identity vertex and the centroid of the simplex ($x = y = z$), while Kimura 2ST matrices [Kim80] ($c = d$) correspond to a plane section of the simplex ($y = z$).

We proceed to compute the volume of a number of subspaces of K3 Markov matrices that can be interesting from a biological perspective. First of all, we introduce some notation:

- Δ is the space of all K3 Markov matrices.
- Δ_* is the space of matrices for which $a \geq b, c, d$.
- Δ_+ is the subspace of matrices with only positive eigenvalues.
- Δ_d is the subspace of matrices that are diagonally dominant, i.e. $a \geq b + c + d$.
- Δ_ε is the subspace of embeddable K3 matrices.

By Lemma 2.4, it is straightforward to check that $\Delta_d \subset \Delta_+ \subset \Delta_* \subset \Delta$. Note that the examples of the preceding section show that $\Delta_\varepsilon \not\subset \Delta_*$ (see the matrix in (6)). However, according to Corollary 3.5,

²A *generic property* is a property that holds almost everywhere, that is, except for a set of measure zero.

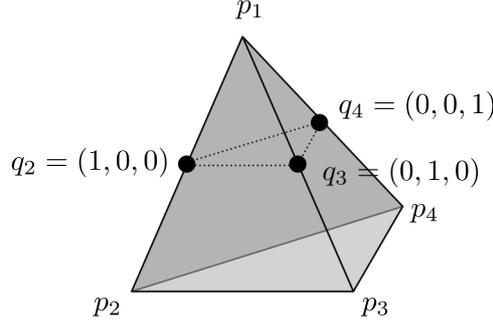


Figure 1: Simplex representing all K3 Markov matrices. Each matrix is represented by its eigenvalues.

all embeddable matrices that are not in Δ_+ correspond to the case of repeated eigenvalues, so they are a marginal case with measure (i.e. volume) 0 within the whole space of K3 matrices. This is because these matrices are constrained by nontrivial algebraic constraints, which make the dimension of the corresponding subspace necessarily smaller. Nevertheless, as shown in (e) and (f) of the next result, there are lots of embeddable matrices with no repeated eigenvalues that are not diagonal dominant.

Theorem 4.1. *We have the following:*

- (a) $V(\Delta_d) = 1/3$; (b) $V(\Delta_+) = 1/2$; (c) $V(\Delta_*) = 2/3$; (d) $V(\Delta) = 8/3$; (e) $V(\Delta_\varepsilon) = 1/4$;
(f) $V(\Delta_\varepsilon \cap \Delta_d) \simeq 0.20336$

Proof. (a) The space Δ_d of diagonal dominant matrices is defined by the inequality $a \geq b + c + d$. Since $a + b + c + d = 1$, this is equivalent to $a \geq 1/2$. Thus, Δ_d is the regular simplex with vertices p_1 , q_2 , q_3 and q_4 (see Figure 1), and its volume is given by the well known formula

$$V(\Delta_d) = \frac{1}{6} |\det(\overrightarrow{p_1q_2}, \overrightarrow{p_1q_3}, \overrightarrow{p_1q_4})| = 1/3.$$

(b) The space Δ_+ of matrices with positive eigenvalues is composed of Δ_d together with the simplex with vertices q_1 , q_2 , q_3 and the centroid O . The formula above gives that this last simplex has volume $1/6$. Therefore, $V(\Delta_+) = V(\Delta_d) + 1/6 = 1/2$.

(c) The three inequalities defining Δ_* are equivalent to $x + y \geq 0$, $x + z \geq 0$ and $y + z \geq 0$. If we denote by p_{ijk} the centroid of the triangle defined by p_i , p_j and p_k , it is straightforward to see that Δ_* is composed of Δ_+ together with the three simplices defined by $\{O, q_2, q_3, p_{123}\}$, $\{O, q_3, q_4, p_{134}\}$ and $\{O, q_4, q_2, p_{124}\}$. These three simplices have the same volume:

$$\frac{1}{6} |\det(\overrightarrow{Oq_i}, \overrightarrow{Oq_j}, \overrightarrow{Op_{1ij}})| = 1/18.$$

It follows that $V(\Delta_*) = 1/2 + 3(1/18) = 2/3$.

(d) follows similarly to the computation of (a).

(e) The computation of Δ_ε is more involved. First of all, as noted above, we can restrict the computation to matrices with no repeated eigenvalues. Therefore, it is enough to compute the volume of the space \mathcal{E} defined by the inequalities (3). By cutting the space \mathcal{E} with planes of the form $z = a$ with $a \in [0, 1]$, we obtain a shape like Figure 2. The computation follows by integrating the corresponding area for all values of a . Given $a \in [0, 1]$ and $x \in [0, a]$, the range of values for y is between ax and x/a , while for $x \in [a, 1]$, the value of y lies between ax and a/x . We are led to compute the following integral

$$V(\Delta_\varepsilon) = \int_0^1 \left(\int_0^z \int_{zx}^{x/z} dy dx + \int_z^1 \int_{xz}^{z/x} dy dx \right) dz$$

which can be easily shown to be equal to $1/4$.

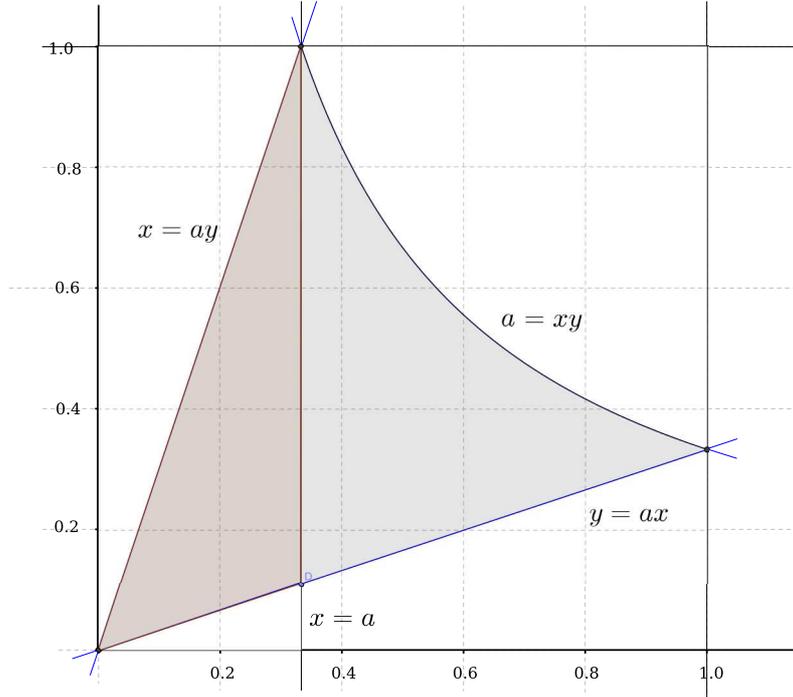


Figure 2: Plane section of the space ε of embeddable K3 matrices with the plane $z = a$, $a \in [0, 1]$.

(f) The computation of the volume of the intersection $\Delta_d \cap \Delta_\varepsilon$ is similar to the computation of $V(\Delta_\varepsilon)$ but for each plane section $z = a$, $a \in [0, 1]$ we have to remove the area of the space below the line $x + y = 1 - a$ (corresponding to non-embeddable matrices). This leads to two different situations: for $a \in [0, 3 - 2\sqrt{2}]$ the line cuts the hyperbola $xy = a$; while for $a \in [3 - 2\sqrt{2}, 1]$ line and hyperbola does not meet. The computation of the corresponding integrals is tedious and we do not include it here. The final value has been obtained using the mathematical software SAGE [Dev]. \square

The values of these volumes illustrate the relative size between the spaces of Markov matrices considered above. Table 1 shows these volumes and the relative volume of embeddable matrices in each of the above subspaces of K3 Markov matrices. These relative volumes are a measure of how many matrices are rejected when taking the continuous-time approach instead of considering other subspaces of matrices. The figures in the second row of the table show, in particular, that there is a big difference between these volumes. For example, embeddable matrices suppose half of the matrices with positive eigenvalues. Similarly, only three out of eight K3 Markov matrices satisfying $a \geq b, c, d$ are embeddable, while we observe that they represent more than the 60% of the diagonal dominant matrices. However, if we only consider diagonal dominant matrices we are rejecting a non-negligible number of embeddable matrices (the difference of volumes between embeddable and diagonal dominant embeddable is $1/4 - 0.20336 = 0.046641$, see Theorem 4.1).

	Δ_ε	Δ	Δ_*	Δ_+	Δ_d
$V(\cdot)$	1/4	8/3	2/3	1/2	1/3
relative vol. of embeddable	1	3/32	3/8	1/2	0.61008

Table 1: Volumes and relative volumes of the embeddable K3 matrices. The relative volumes of embeddable matrices within each spaces are shown in the second row of the table and are obtained as the quotients $V(\Delta_\varepsilon \cap \cdot)/V(\Delta)$.

Remark 4.2. From Theorem 3.9 (i) it follows that all the values of Table 1 remain the same if we only consider $K3$ -embeddable matrices (instead of embeddable matrices).

Since we have no characterization for embeddable matrices with repeated eigenvalues, and this set has positive measure within the Kimura 2ST model (Theorem 3.9), we are not able to compute volumes within this model.

Nevertheless, for the Jukes-Cantor model, Theorem 3.9 states that embeddability is equivalent to JC-embeddability, which has been characterized in Theorem 3.3. Adapting the notation used for the $K3$ matrices to the JC model, we get that

$$\begin{aligned}\Delta^{JC} &= \{(x, x, x) \in \Delta \mid x \in [-1/3, 1]\}; \\ \Delta_*^{JC} &= \{(x, x, x) \in \Delta \mid x \in [0, 1]\}; \\ \Delta_+^{JC} = \Delta_\varepsilon^{JC} &= \{(x, x, x) \in \Delta \mid x \in (0, 1]\}; \\ \Delta_d^{JC} &= \{(x, x, x) \in \Delta \mid x \in [1/3, 1]\}.\end{aligned}$$

A straightforward computation shows that the space of embeddable Jukes-Cantor matrices has volume (length) $\sqrt{3}$ while the space of all Markov Jukes-Cantor matrices has volume $\frac{4}{3}\sqrt{3}$. That is, three out of four Jukes-Cantor matrices are embeddable.

5 Discussion

As suggested in the introduction, there are four connected problems relative to the algebraic and continuous-time evolutionary models. Namely,

1. whether a given Markov matrix is e^Q for some rate matrix Q (embedding problem);
2. whether a given Markov matrix M within an algebraic model \mathcal{M} is e^Q for some rate matrix Q with the same symmetries as M (\mathcal{M} -embedding problem);
3. whether the product of embeddable matrices is embeddable, i.e. if Q_1 and Q_2 are rate matrices, then is it true that $e^{Q_1}e^{Q_2} = e^{Q_3}$ for some rate matrix Q_3 ? (multiplicative closure of embeddable matrices);
4. same question as in 3. but with Q_1, Q_2 and Q_3 within a particular continuous-time model (multiplicative closure of \mathcal{M} -embeddable matrices).

The first two questions are motivated by the connection between the algebraic and the continuous-time models. The last two are more intrinsic of the continuous-time approach. In this paper, we have discussed the embedding problem for Markov matrices within the Kimura 3-parameter model (**1.**). Under a generic assumption (that is, eigenvalues should be different), we have obtained a characterization of the embeddability of the matrices for this model in terms of some inequalities relative to their eigenvalues (Corollary 3.5). Moreover, Theorem 3.9 shows that in this case, rates are identifiable and keep the $K3$ symmetries, so that embeddability holds if and only if so does $K3$ -embeddability, which has been characterized in Theorem 3.3 (**2.**). As for the problem (**3.**), Remark 3.7 shows that the product of embeddable matrices within the Kimura 3ST is not embeddable in general. However, under the assumption of different eigenvalues again, problems (**3.**) and (**4.**) become equivalent and Theorem 2.5 gives an affirmative answer to both questions.

As a consequence of the characterization of Corollary 3.5, we have been able to compute and compare the volume of embeddable and $K3$ -embeddable matrices within some subspaces of $K3$ Markov matrices that may be regarded as a compromise solution between the continuous-time approach and the whole space of $K3$ Markov matrices (see Theorem 4.1 and Table 1).

Despite the fact that the evolutionary submodels of Kimura 3ST (Kimura 2ST and Jukes-Cantor) have repeated eigenvalues, the results obtained here also give information about them. While a characterization of the the embeddability of the Jukes-Cantor model follows easily from the general case of $K3$ matrices, the embeddability of Kimura 2ST model is more involved and remains an open question. Under this model,

we have provided examples of matrices with repeated eigenvalues showing that there are embeddable matrices within the K2ST model that have no Markov generator with K3 form (see Theorem 3.9) and although these matrices are close to saturation, Theorem 2.2 shows they are still biologically relevant.

Another open last issue is the identifiability of rates for K3 embeddable matrices with repeated eigenvalues. In this case, the inequalities (5) for negative eigenvalues and (7) for positive eigenvalues are sufficient for a K3 Markov matrix to have different Markov generators. We believe that these inequalities are actually necessary, but a much more technical analysis is required. Moreover, we conjecture that under the Kimura 3ST model, the rates of an embeddable matrix that is not K3-embeddable are not identifiable. An affirmative answer to these questions would allow us to characterize the embeddability of any K3 Markov matrix (cf. Theorem 3.4), and show that for Markov matrices with positive eigenvalues, embeddability is equivalent to \mathcal{M} -embeddability for K3ST and any of its submodels. We defer such a study for a future publication.

6 Acknowledgement

The authors want to thank Marta Casanellas and Jeremy Sumner for conversations held on this topic. They also wish to thank the reviewers for useful comments and suggestions.

JRL and JFS are partially supported by Spanish government MTM2015-69135-P. Second author is also supported by Generalitat de Catalunya 2014SGR634.

References

- [AR08] Elizabeth S. Allman and John A. Rhodes. Phylogenetic ideals and varieties for the general markov model. *Advances in Applied Mathematics*, 40(2):127–148, 2008.
- [BA02] K. P. Burnham and D. Anderson. *Model Selection and Multi-Model Inference*. Springer-Verlag, 2002.
- [BC04] Sergio Blanes and Fernando Casas. On the convergence and optimization of the Baker-Campbell-Hausdorff formula. *Linear Algebra Appl.*, 378:135–158, 2004.
- [Cam97] J. E. Campbell. On a law of combination of operators (second paper). *Proc. London Math. Soc.*, 28:381390, 1897.
- [CFS08] M. Casanellas and J. Fernández-Sánchez. Geometry of the Kimura 3-parameter model. *Adv. in Appl. Math.*, 41(3):265–292, 2008.
- [CFS10] M. Casanellas and J. Fernández-Sánchez. Relevant phylogenetic invariants of evolutionary models. *Journal de Mathématiques Pures et Appliquées*, 96:207–229, 2010.
- [Cul66] Walter J. Culver. On the existence and uniqueness of the real logarithm of a matrix. *Proc. Amer. Math. Soc.*, 17:1146–1151, 1966.
- [Dav10] E. B. Davies. Embeddable Markov matrices. *Electron. J. Probab.*, 15:no. 47, 1474–1486, 2010.
- [Dev] The Sage Developers. *SageMath, the Sage Mathematics Software System (Version 6.1.1)*. <http://www.sagemath.org>.
- [DK08] Jan Draisma and Jochen Kuttler. On the ideals of equivariant tree models. *Mathematische Annalen*, 344:619–644, 2008.
- [Gan59] F.R. Gantmacher. *The theory of matrices Vol.1*. Chelsea Publishing Company, 1959.
- [Hig08] Nicholas J. Higham. *Functions of Matrices: Theory and Computation*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2008.

- [HPCD05] S. Y. W. Ho, M. J. Phillips, A. Cooper, and A. J. Drummond. Time dependency of molecular rate estimates and systematic overestimation of recent divergence times. *Mol. Biol. Evol.*, 22:1561–1568, 2005.
- [HSP⁺07] Simon Y W Ho, Beth Shapiro, Matthew J Phillips, Alan Cooper, and Alexei J Drummond. Evidence for time dependency of molecular rate estimates. *Syst Biol*, 56(3):515–522, 2007.
- [JC69] TH Jukes and CR Cantor. Evolution of protein molecules. *In Mammalian Protein Metabolism*, pages 21–132, 1969.
- [Kim80] M Kimura. A simple method for estimating evolutionary rates of base substitution through comparative studies of nucleotide sequences. *J. Mol. Evol.*, 16:111–120, 1980.
- [Kim81] M. Kimura. Estimation of evolutionary distances between homologous nucleotide sequences. *Proc. Natl. Acad. Sci.*, 78:1454–1458, 1981.
- [KK17] Dimitra Kosta and Kaie Kubjas. Geometry of symmetric group-based models. *ArXiv e-prints 1705.09228*, 2017.
- [LSB⁺98] P. J. Lockhart, M. A. Steel, A. C. Barbrook, D. H. Huson, and C. J. Howe. A covariotide model describes the evolution of oxygenic photosynthesis. *Mol. Biol. Evol.*, 15:1183–1188, 1998.
- [SFSJ12] J. G. Sumner, J. Fernández-Sánchez, and P. D. Jarvis. Lie Markov models. *J. Theor. Biol.*, 298:16–31, 2012.
- [SJFS⁺12] Jeremy G. Sumner, Peter D. Jarvis, Jesús Fernández-Sánchez, Bodie T. Kaine, Michael D. Woodhams, and Barbara R. Holland. Is the general time-reversible model bad for molecular phylogenetics? *Systematic Biology*, 61(6):1069–1074, 2012.
- [SS05] B. Sturmfels and S. Sullivant. Toric ideals of phylogenetic invariants. *J. Comput. Biol.*, 12:204–228, 2005.
- [Sum17] J. G. Sumner. Multiplicatively closed Markov models must form Lie algebras . *ArXiv e-prints 1704.01418*, 2017.
- [Tav86] S. Tavaré. Some probabilistic and statistical problems in the analysis of dna sequences. *Lectures on Mathematics in the Life Sciences (American Mathematical Society)*, 17:57–86, 1986.
- [VYP⁺13] Klara L. Verbyla, Von Bing Yap, Anuj Pahwa, Yunli Shao, and Gavin A. Huttley. The embedding problem for markov models of nucleotide substitution. *PLoS ONE*, 8:e69187, 7 2013.
- [WSL⁺17] Michael D. Woodhams, Jeremy G. Sumner, David A. Liberles, Michael A. Charleston, and Barbara R. Holland. Exploring the consequences of lack of closure in codon models. *ArXiv e-prints 1709.05079*, 2017.