

Tropical Logistic Regression Model on Space of Phylogenetic Trees

Georgios Aliatimis Ruriko Yoshida Burak Boyaci
James A. Grant

Abstract

Motivation: Classification of gene trees is an important task both in analysis of multi-locus phylogenetic data, and assessment of the convergence of Markov Chain Monte Carlo (MCMC) analyses used in Bayesian phylogenetic tree reconstruction. The logistic regression model is one of the most popular classification models in statistical learning, thanks to its computational speed and interpretability. However, it is not appropriate to directly apply the logistic regression model to a set of phylogenetic trees with the same set of leaf labels, as the space of phylogenetic trees is not Euclidean.

Results: It is well-known in tropical geometry and phylogenetics that the space of phylogenetic trees is a tropical linear space in terms of the max-plus algebra. Therefore, in this paper, we propose an analogue approach of the logistic regression model in the setting of tropical geometry. In our proposed method, we consider two cases: where the numbers of the species trees are fixed as one and two, and we estimate the species tree(s) from a sample of gene trees distributed over the space of ultrametrics, which is a tropical linear space. We show that both models are statistically consistent and bounds on the generalization error of both models are derived. Finally, we conduct computational experiments on simulated data generated by the multi-species coalescent model and apply our model to African coelacanth genomes to infer the species tree.

1 Introduction

Phylogenomics is a new field that applies tools from phylogenetics to genome datasets. The multi-species coalescent model is often used to model the distribution of gene trees under a given species tree [15]. The first step in statistical analysis of phylogenomics is that evolutionary biologists, also known as systematists, analyze sequence alignments to determine whether their evolutionary histories are congruent with each other. In this step, systematists aim to identify genes with unusual evolutionary events, such as duplication, horizontal gene transfer, or hybridization [2]. To accomplish this, they compare multiple sets of *gene trees*, that is, phylogenetic trees reconstructed from alignments of genes.

The classification of gene trees into different categories is therefore important for analyzing multi-locus phylogenetic data, but also in assessing convergence of Markov Chain Monte Carlo (MCMC) analyses for Bayesian inference on phylogenetic tree reconstruction. Often we apply MCMC samplers to estimate the posterior distribution of a phylogenetic tree given an observed alignment. These samplers typically run multiple independent Markov chains on the same observed data set. The goal is to check whether these chains converge to the same distribution. This process is often done by comparing summary statistics computed from sampled trees, as there is no classification model over the *space of phylogenetic trees*, the set of all possible phylogenetic trees with leaves $[m] := \{1, \dots, m\}$. However, computing a summary statistic from a sample naturally loses information about the sample [17].

In Euclidean geometry, the logistic regression model is the simplest generalized linear model for classification. It is a supervised learning method that classifies data points by modeling the log-odds of having a response variable in a particular class as a linear combination of predictors. This model is highly popular in statistical learning due to its simplicity, computational speed and interpretability. However, directly applying classical supervised models to a set of sampled trees may be misleading, since the space of phylogenetic trees does not conform to Euclidean geometry.

The space of phylogenetic trees with labeled leaves $[m]$ is a union of lower dimensional polyhedral cones with dimension $m - 1$ over \mathbb{R}^e where $e = \binom{m}{2}$ [22, 11]. This space is not Euclidean and even lacks convexity [11]. In fact, [22] showed that the space of phylogenetic trees is a *tropicalization* of linear subspaces defined by a system of tropical linear equations [18] and is therefore a tropical linear space.

Consequently, many researchers have applied tools from tropical geometry to statistical learning methods in phylogenomics, such as principal component analysis over the space of phylogenetic trees with a given set of leaves $[m]$ [27, 18], kernel density estimation [25], MCMC sampling [24], and support vector machines [26]. Recently, [1] proposed a tropical linear regression over the tropical projective space as the best-fit tropical hyperplane.

In this paper, an analog of the logistic regression is developed over the tropical projective space, which is the quotient space $\mathbb{R}^e/\mathbb{R}\mathbf{1}$ where $\mathbf{1} := (1, 1, \dots, 1)$. Given a sample of observations within this space, the proposed model finds the “best-fit” tree representative $\omega_Y \in \mathbb{R}^e/\mathbb{R}\mathbf{1}$ of each class $Y \in \{0, 1\}$ and the “best-fit” deviation of the gene trees. This tree representative is a statistical parameter and can be interpreted as the corresponding species tree of the gene trees. It is established that the median tree, specifically the Fermat-Weber point, can asymptotically approximate the inferred tree representative of each class. The response variable $Y \in \{0, 1\}$ has conditional distribution $Y|X \sim \text{Bernoulli}(S(h(X)))$, where $h(x)$ is small when x is close to ω_0 and far away from ω_1 and vice versa.

In Section 2 an overview of tropical geometry and its connections to phylogenetics is presented. The one-species and two-species tropical logistic model is developed in Section 3. Theoretical results, including the

optimality of the proposed method over tropically distributed predictor trees, the distance distribution of those trees from their representative, the consistency of estimators and the generalization error of each model are stated in Section 4 and proved in Supplement A. Section 5 explains the benefit and suitability of using the Fermat-Weber point approximation for the inferred trees and a sufficient optimality condition is stated. Computational results are presented in Section 6 where a toy example is considered for illustration purposes. Additionally, a comparison study between classical, tropical and BHV logistic regression is conducted on data generated under the coalescent model and an implementation of the proposed method on the lungfish dataset is analysed. Finally, the paper concludes with a discussion in Section 7.

2 Tropical Geometry and Phylogenetic Trees

2.1 Tropical Basics

This section covers the basics of tropical geometry and provides the theoretical background for the model developed in later sections. For more details regarding the basic concepts of tropical geometry covered in this section, readers are recommended to consult [13].

2.1.1 Tropical Metric

A key tool from tropical geometry is the *tropical metric* also known as the *tropical distance* defined as follows:

Definition 2.1 (Tropical distance). *The tropical distance, more formally known as the Generalized Hilbert projective metric, between two vectors $v, w \in (\mathbb{R} \cup \{-\infty\})^e$ is defined as*

$$d_{\text{tr}}(v, w) := \|v - w\|_{\text{tr}} = \max_i \{v_i - w_i\} - \min_i \{v_i - w_i\}, \quad (1)$$

where $v = (v_1, \dots, v_e)$ and $w = (w_1, \dots, w_e)$.

Remark 1. *Consider two vectors $v = (c, \dots, c) = c\mathbf{1} \in \mathbb{R}^e$ and $w = \mathbf{0} \in \mathbb{R}^e$. It is easy to verify that $d_{\text{tr}}(v, w) = 0$ and as a result d_{tr} is not a metric in \mathbb{R}^e . The space in which d_{tr} is a metric treats all points in $\{c\mathbf{1} : c \in \mathbb{R}\} = \mathbb{R}\mathbf{1}$ as the same point. The quotient space $(\mathbb{R} \cup \{-\infty\})^e / \mathbb{R}\mathbf{1}$ achieves just that.*

Proposition 1. *The function d_{tr} is a well-defined metric on $(\mathbb{R} \cup \{-\infty\})^e / \mathbb{R}\mathbf{1}$, where $\mathbf{1} \in \mathbb{R}^e$ is the vector of all-ones.*

2.2 Equidistant Trees and Ultrametrics

Phylogenetic trees depict the evolutionary relationship between different taxa. For example, they may summarise the evolutionary history of certain species. The leaves of the tree correspond to the species studied, while internal nodes represent (often hypothetical) common ancestors of those

species and their ancestors. In this paper, only rooted phylogenetic trees are considered, with the common ancestor of all taxa based on the root of the tree. The branch lengths of these trees are measured in evolutionary units i.e. the amount of evolutionary change. Under the molecular clock hypothesis, the rate of genetic change between species is constant over time, which implies genetic equidistance and allows us to treat evolutionary units as proportional to time units. Consequently, phylogenetic trees of extant species are *equidistant trees*.

Definition 2.2 (Equidistant tree). *Let T be a rooted phylogenetic tree with leaf label set $[m]$, where $m \in \mathbb{N}$ is the number of leaves. If the distance from all leaves $i \in [m]$ to the root is the same, then T is an equidistant tree.*

It is noted that the molecular clock hypothesis has limitations and the rate of genetic change can in fact vary from one species to another. However, the assumption that gene trees are equidistant is not unusual in phylogenomics; the multispecies coalescent model makes that assumption in order to conduct inference on the species tree from a sample of gene trees [14]. The proposed classification method is not restricted to equidistant trees, but all coalescent model gene trees produced in Section 6.2. are equidistant.

In order to conduct any mathematical analysis, a vector representation of trees is needed. A common way is to use BHV coordinates [5] but in this paper *distance matrices* are used instead.

Definition 2.3 (Distance matrix). *Consider a phylogenetic tree T with leaf label set $[m]$. Its distance matrix $D \in \mathbb{R}^{m \times m}$ has components D_{ij} being the pairwise distance between a leaf $i \in [m]$ to a leaf $j \in [m]$. It follows that the matrix is symmetric with zeros on its diagonals. For equidistant trees, D_{ij} is equal to twice the difference between the current time and the latest time that the common ancestor of i and j was alive.*

To form a vector, the distance matrix D is mapped onto \mathbb{R}^e by vectorizing the strictly upper triangular part of D , i.e.

$$D \mapsto (D_{12}, \dots, D_{1m}, D_{23}, \dots, D_{2m}, \dots, D_{(m-1)m}) \in \mathbb{R}^e,$$

where the dimension of the resulting vector is equal to the number of all possible pairwise combinations of leaves in T . Hence the dimension of the phylogenetic tree space is $e = \binom{m}{2}$. In what follows, the connection between the space of phylogenetic trees and tropical linear spaces is established.

Definition 2.4 (Ultrametric). *Consider the distance matrix $D \in \mathbb{R}^{m \times m}$. Then if*

$$\max\{D_{ij}, D_{jk}, D_{ik}\}$$

is attained at least twice for any $i, j, k \in [m]$, D is an ultrametric. Note that the distance map $d(i, j) = D_{ij}$ forms a metric in $[m]$, with the strong triangular inequality satisfied. The space of ultrametrics is denoted as \mathcal{U}_m .

Theorem 2.5 (noted in [6]). *Suppose we have an equidistant tree T with a leaf label set $[m]$ and D as its distance matrix. Then, D is an ultrametric if and only if T is an equidistant tree.*

Using Theorem 2.5, if we wish to consider all possible equidistant trees, then it is equivalent to consider the space of ultrametrics as the space of phylogenetic trees on $[m]$. Here we define \mathcal{U}_m as the space of ultrametrics with a set of leaf labels $[m]$. Theorem B.1 in Supplement B establishes the connection between phylogenetic trees and tropical geometry by stating that the ultrametric space is a tropical linear space.

3 Method

Unlike tropical PCA developed by [27] which is an unsupervised learning to reduce its dimensionality, logistic regression is a supervised learning model with a categorical response variable associated with the input variable. In this paper, we only consider the simplest case, that is, a bivariate response variable $Y \in \{0, 1\}$ given with the explanatory variable $x \in \mathbb{R}^n$, where n is the number of covariates in the model. Under the logistic model, $Y \sim \text{Bernoulli}(p(x|\omega))$ where

$$p(x|\omega) = \mathbb{P}(Y = 1|x) = \frac{1}{1 + \exp(-h_\omega(x))} = S(h_\omega(x)), \quad (2)$$

where S is the logistic function and ω is the statistical/model parameter. The most intuitive and sensible classifier for this model is defined as

$$C(x) = \mathbb{I}(h_\omega(x) > 0) \in \{0, 1\}. \quad (3)$$

The log-likelihood function of logistic regression for N observation pairs $(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})$ is

$$l(\omega|x, y) = \frac{1}{N} \sum_{i=1}^N y^{(i)} \log p_\omega^{(i)} + (1 - y^{(i)}) \log(1 - p_\omega^{(i)}), \quad (4)$$

where $p_\omega^{(i)} = p(x^{(i)}|\omega)$. The training model seeks a statistical estimator $\hat{\omega}$ that maximizes this function.

Everything mentioned thus far follows from the classical logistic regression. The only difference is our choice of the function h . In fact, this function can be derived from the conditional distributions $Y|X$, as stated in Lemma 3.1. The classical and tropical logistic regression assume different distributions $Y|X$ and therefore get different functions h , as shown in Examples 3.2 and 3.3.

Lemma 3.1. *Let $Y \sim \text{Bernoulli}(r)$ and define the random vector $X \in \mathbb{R}^n$ with conditional distribution $X|Y \sim f_Y$, where f_0, f_1 are probability density functions defined in \mathbb{R}^n . Then, $Y|X \sim \text{Bernoulli}(p(X))$ with $p(x) = \sigma(h(x))$, where*

$$h(x) = \log \left(\frac{r f_1(x)}{(1-r) f_0(x)} \right). \quad (5)$$

Example 3.2 (Normal distribution and classical logistic regression). *Suppose that the two classes are equiprobable ($r = 1/2$) and that the covariate is multivariate normal*

$$X|Y \sim \mathcal{N}(\omega_Y, \sigma^2 I_n),$$

where n is covariate dimension and I_n is the identity matrix. Using Lemma 3.1, the optimal model has

$$h(x) = -\frac{\|x - \omega_1\|^2}{2\sigma^2} + \frac{\|x - \omega_0\|^2}{2\sigma^2} = \frac{(\omega_1 - \omega_0)^T}{\sigma^2}(x - \bar{\omega}), \quad (6)$$

where $\|\cdot\|$ is the Euclidean norm and $\bar{\omega} = (\omega_0 + \omega_1)/2$. This model is the classical logistic regression model with translated covariate $X - \bar{\omega}$ and $\omega = \sigma^{-2}(\omega_1 - \omega_0)$.

Example 3.3 (Tropical Laplace distribution). *It may be assumed that the covariates are distributed according to the tropical version of the Laplace distribution, as presented in [29], with mean ω_Y and probability density functions*

$$f_Y(x) = \frac{1}{\Lambda} \exp\left(-\frac{d_{\text{tr}}(x, \omega_Y)}{\sigma_Y}\right), \quad (7)$$

where Λ is the normalizing constant of the distribution.

Proposition 2. *In distribution (7), the normalizing factor is $\Lambda = e! \sigma_Y^{e-1}$.*

Proof. See Supplement A. □

Combining the result of the proposition above with equations (5) and (7) yields

$$h_{\omega_0, \omega_1}(x) = \frac{d_{\text{tr}}(x, \omega_0)}{\sigma_0} - \frac{d_{\text{tr}}(x, \omega_1)}{\sigma_1} + (e-1) \log\left(\frac{\sigma_0}{\sigma_1}\right). \quad (8)$$

In its most general form, the model parameters are $(\omega_0, \omega_1, \sigma_0, \sigma_1)$ so the parameter space is a subset of $(\mathbb{R}^e / \mathbb{R}\mathbf{1})^2 \times \mathbb{R}_+^2$ with dimension $2e$. There are two notable special cases that are discussed in more detail than the general model in the subsequent sections; the one-species and two-species model.

For the one-species model, it is assumed that $\omega_0 = \omega_1$ and $\sigma_0 \neq \sigma_1$. If, without loss of generality, $\sigma_1 > \sigma_0$, equation (8) becomes

$$h_{\omega}(x) = \lambda (d_{\text{tr}}(x, \omega) - c), \quad (9)$$

where $\lambda = (\sigma_0^{-1} - \sigma_1^{-1})$ and $\lambda c = \log(\sigma_1/\sigma_0)$. Symbolically, the expression in equation (9) can be considered to be a scaled tropical inner product, whose direct analogue in classical logistic regression is the classical inner product $h_{\omega}(x) = \omega^T x$. See Section C in the supplement for more details. The classifier is $C(x) = \mathbb{I}(d_{\text{tr}}(x, \hat{\omega}) > c)$, where $\hat{\omega}$ is the inferred estimator of ω^* . Note that the classification threshold and the probability contours ($p(x)$) are tropical circles, illustrated in Figure 1.

For the two-species-tree model, it is assumed that $\sigma_0 = \sigma_1$. Equation (8) reduces to

$$h_{\omega_0, \omega_1}(x) = \sigma^{-1} (d_{\text{tr}}(x, \omega_0) - d_{\text{tr}}(x, \omega_1)), \quad (10)$$

with a classifier $C(x) = \mathbb{I}(d_{\text{tr}}(x, \hat{\omega}_0) > d_{\text{tr}}(x, \hat{\omega}_1))$, where $\hat{\omega}_y$ is the inferred tree for class $y \in \{0, 1\}$. The classification boundary is the tropical bisector which is extensively studied in [7] between the estimators $\hat{\omega}_0$ and $\hat{\omega}_1$ and

the probability contours are tropical hyperbolae with $\hat{\omega}_0$ and $\hat{\omega}_1$ as foci, as shown in Figure 3(right).

The one-species model is appropriate when the gene trees of both classes are concentrated around the same species tree ω with potentially different concentration rates. When the gene trees of each class come from distributions centered at different species trees the two-species model is preferred.

4 Theoretical Properties

In this section, we show some theoretical results on the model, including that it is statistically consistent. The proofs of these results can be found Section A in the supplement.

From Lemma 3.1, we have that given the distributions of Y and $X|Y$, the distribution of $Y|X$ follows. Proposition 3 states that from all possible distributions $Y|X$, the one that fits the data best is the true distribution as given by the lemma. Therefore, the best model to fit data that have been generated by tropical Laplace distribution (7) is the tropical logistic regression.

Proposition 3. *Let $Y \sim \text{Bernoulli}(r)$ and define the random vector $X \in \mathbb{R}^n$ with conditional distribution $X|Y \sim f_Y$, where f_0, f_1 are probability density functions defined in \mathbb{R}^n . The functional p that maximises the expected log-likelihood as given by equation (4) is $p(x) = \sigma(h(x))$, with h defined as in equation (5) of Lemma 3.1.*

Corollary 1, which is based on Proposition 4 allows us to derive the distribution of the radius $d(X, \omega)$ for both the Euclidean and the tropical metric. However, the arguments used in the proof of Corollary 1 do not work for distributions defined on the space of ultrametric trees \mathcal{U}_m , because these spaces are not translation invariant. For a similar reason, the corollary does not apply to the BHV metric.

Proposition 4. *Consider a function $d : \mathbb{R}^n \rightarrow \mathbb{R}$ with $\alpha d(x) = d(\alpha x)$, for all $\alpha \geq 0$. If $X \sim f$ with $f(x) \propto \exp(-d^i(x)/(i\sigma^i))$ being a valid probability density function, for some $i \in \mathbb{N}$, $\sigma > 0$. Then, $d^i(X) \sim i\sigma^i \text{Gamma}(n/i)$.*

Corollary 1. *If $X \in \mathbb{R}^e$ with $X \sim f \propto \exp(-d^i(x, \omega^*)/(i\sigma^i))$, where d is the Euclidean metric, then $d^i(X, \omega^*) \sim i\sigma^i \text{Gamma}(e/i)$. If $X \in \mathbb{R}^e/\mathbb{R}\mathbf{1}$ with $X \sim f \propto \exp(-d^i(x, \omega^*)/(i\sigma^i))$, where d is the tropical metric, then $d^i(X, \omega^*) \sim i\sigma^i \text{Gamma}((e-1)/i)$.*

The main result of this section is Theorem 4.1, which states that the general tropical logistic regression method is statistically consistent. The proof can also be adapted for the special cases of one-species and two-species model.

Theorem 4.1. *The estimator $(\hat{\omega}, \hat{\sigma}) = (\hat{\omega}_0, \hat{\omega}_1, \hat{\sigma}_0, \hat{\sigma}_1) \in \Omega^2 \times \Sigma^2$ of the parameter $(\omega^*, \sigma^*) = (\omega_0^*, \omega_1^*, \sigma_0^*, \sigma_1^*) \in \Omega^2 \times \Sigma^2$ is defined as the maximizer of the logistic likelihood function, where $\Omega \subset \mathbb{R}^e/\mathbb{R}\mathbf{1}$ and $\Sigma \subset \mathbb{R}_+$ are compact sets. Moreover, it is assumed that the covariate-response pairs $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ are independent and identically distributed*

with $X_i \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$, $d_{\text{tr}}(X, \omega_Y)$ being integrable and square-integrable and $Y_i \sim \text{Bernoulli}(S(h(X_i, (\omega^*, \sigma^*))))$. Then,

$$(\hat{\omega}, \hat{\sigma}) \xrightarrow{P} (\omega^*, \sigma^*) \text{ as } n \rightarrow \infty.$$

In other words, the model parameter estimator is consistent.

Finally, Proposition 5 and 6 provide generalization error bounds for the one-species and two species model respectively. In both cases the error bounds are getting worse as the estimation error ϵ increases. It is worth mentioning that in the case of exact estimation, the generalization error of the one-species model can be computed explicitly by equation (11). Moreover, there is a higher misclassification rate from the more dispersed class (inequality (12)).

Proposition 5. Consider the one-species model where $\omega = \omega_0 = \omega_1 \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$ and without loss of generality $\sigma_0 < \sigma_1$. The classifier is $C(x) = \mathbb{I}(h_{\hat{\omega}}(x) \geq 0)$, where h is defined in equation (9) and $\hat{\omega}$ is the estimate for ω^* . Define the covariate-response joint random variable (X, Y) with $Z = \sigma_Y^{-1} d_{\text{tr}}(X, \omega_Y^*)$ drawn from the same distribution with cumulative density function F . Then,

$$\begin{aligned} \mathbb{P}(C(X) = 1 | Y = 0) &\in [1 - F(\sigma_1(\alpha + \epsilon)), 1 - F(\sigma_1(\alpha - \epsilon))], \\ \mathbb{P}(C(X) = 0 | Y = 1) &\in [F(\sigma_0(\alpha - \epsilon)), F(\sigma_0(\alpha + \epsilon))], \text{ where} \\ \alpha &= \frac{\log \frac{\sigma_1}{\sigma_0}}{\sigma_1 - \sigma_0}, \text{ and } \epsilon = (e - 1) \frac{d_{\text{tr}}(\hat{\omega}, \omega^*)}{\sigma_1 \sigma_0}. \end{aligned}$$

The generalization error defined as $\mathbb{P}(C(X) \neq Y)$ lies in the average of the two intervals above. In particular, note that if $\hat{\omega} = \omega^*$, then $\epsilon = 0$ and the intervals shrink to a single point, so the misclassification probabilities and generalization error can be computed explicitly.

$$\mathbb{P}(C(X) \neq Y) = \frac{1}{2} (1 - F(\sigma_1 \alpha) + F(\sigma_0 \alpha)) \quad (11)$$

Moreover, if $\hat{\omega} = \omega_*$ and $Z \sim \text{Gamma}(e - 1, 1)$, then

$$\mathbb{P}(C(X) = 1 | Y = 0) < \mathbb{P}(C(X) = 0 | Y = 1). \quad (12)$$

Proposition 6. Consider the random vector $X \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$ with response $Y \in \{0, 1\}$ and the random variable $Z = d_{\text{tr}}(X, \omega_Y^*)$. Assuming that the probability density function is $f_X(x) \propto f_Z(d_{\text{tr}}(x, \omega_Y^*))$, the generalization error satisfies the following upper bound

$$\mathbb{P}(C(X) \neq Y) \leq \frac{1}{2} F_Z^C(\Delta_\epsilon) + h(\epsilon), \quad (13)$$

where $\epsilon = d_{\text{tr}}(\hat{\omega}_1, \omega_1^*) + d_{\text{tr}}(\hat{\omega}_0, \omega_0^*)$, $2\Delta_\epsilon = (d_{\text{tr}}(\omega_1^*, \omega_0^*) - \epsilon)$, F_Z^C is the complementary cumulative distribution of Z , and $h(\epsilon)$ is an increasing function of ϵ with $2h(\epsilon) \leq F_Z^C(\Delta_\epsilon)$ and $h(0) = 0$ assuming that $\mathbb{P}(d_{\text{tr}}(X, \omega_1^*) = d_{\text{tr}}(X, \omega_{-1}^*)) = 0$.

Observe that the upper bound is a strictly increasing function of ϵ .

Example 4.2. *The complementary cumulative distribution of $\text{Gamma}(n, \sigma)$ is $F^C(x) = \Gamma(n, x/\sigma)/\Gamma(n, 0)$, where Γ is the upper incomplete gamma function and $\Gamma(n, 0) = \Gamma(n)$ is the regular Gamma function. Therefore, the tropical distribution given in equation (7) yields the following upper bound for the generalization error*

$$\frac{\Gamma\left(e - 1, \frac{d_{\text{tr}}(\omega_0^*, \omega_1^*)}{2\sigma}\right)}{2\Gamma(e - 1)}, \quad (14)$$

under the assumptions of Proposition 6 and assuming that the estimators coincide with the theoretical parameters. This assumption is reasonable for large sample sizes and it follows from Theorem 4.1.

In Section 6, these theoretical properties are applied. Bounds on the generalization error from Propositions 5 and 6 are computed and the suitability of Euclidean and tropical distributions, and as a result of classical and tropical logistic regards, using the distance distribution of Proposition 4.

5 Optimization

As in the classical logistic regression, the parameter vectors $(\hat{\omega}, \hat{\sigma})$ maximising the log-likelihood (4), are chosen as statistical estimators. Identifying these requires the implementation of a continuous optimization routine. While root-finding algorithms typically work well for the classical logistic regression where the log-likelihood is concave, they are unsuitable here. The gradients of the log-likelihood under the proposed tropical logistic models are only piecewise continuous, with the number of discontinuities increasing along with the sample size. Furthermore, even if a parameter is found, it may merely be a local optimum. In light of this, the tropical Fermat-Weber problem of [12] is revisited.

5.1 Fermat-Weber Point

A Fermat-Weber point or geometric mean $\hat{\omega}_n$ of the sample set (X_1, \dots, X_n) is a point that minimizes the sum of distances from to sample points, i.e.

$$\hat{\omega}_n \in \arg \min_{\omega} \sum_{i=1}^n d_{\text{tr}}(X_i, \omega). \quad (15)$$

This point is rarely unique for finite n [12]. However, the proposition below gives conditions for asymptotic convergence.

Proposition 7. *Let $X_i \stackrel{\text{iid}}{\sim} f$, where f is a distribution that is symmetric around its center $\omega^* \in \mathbb{R}^e/\mathbb{R}\mathbf{1}$ i.e. $f(\omega^* + \delta) = f(\omega^* - \delta)$ for all $\delta \in \mathbb{R}^e/\mathbb{R}\mathbf{1}$. Let $\tilde{\omega}_n$ be any Fermat-Weber point as defined in equation (15). Then, $\tilde{\omega}_n \xrightarrow{P} \omega^*$ as $n \rightarrow \infty$.*

The significance of Proposition 4.1 is twofold. It proves that the Fermat-Weber sets of points sampled from symmetric distributions tend to a unique point. This is a novel result and ensures that for sufficiently

large sample sizes the topology of any Fermat-Weber point is fixed. Additionally, using Theorem 4.1 and Proposition 4.1, $\hat{\omega}_n - \tilde{\omega}_n \xrightarrow{P} 0$ as $n \rightarrow \infty$. As a result, for a sufficiently large sample size we may use the Fermat-Weber point as an approximation for the MLE vector, which is a simpler problem especially for the two-species model. Instead of having a single optimization problem with $2e - 1$ degrees of freedom, three simpler problems are considered; finding the Fermat-Weber point of each of the two clusters, which has $e - 1$ degrees of freedom and then finding the optimal σ which is a one dimensional root finding problem. The algorithms of our implementation for both model can be found in Supplement D.

There is also another yet another benefit of using Fermat-Weber points. In [12], Fermat-Weber points are computed by means of linear programming, which is computationally expensive. Employing a gradient-based method is much faster, but there is no guarantee of convergence. Nevertheless, if the gradient, which is an integer vector, vanishes, then it is guaranteed, as below, that the algorithm has reached a Fermat-Weber point.

Proposition 8. *Let $X_1, \dots, X_n \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$, $\omega \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$ and define the function*

$$f(\omega) = \sum_{i=1}^n d_{\text{tr}}(X_i, \omega).$$

- i. The gradient vector of f is defined at ω if and only if the vectors $\omega - X_i$ have unique maximum and minimum components for all $i \in [n]$.*
- ii. If the gradient of f at ω is well-defined and zero, then ω is a Fermat-Weber point.*

Proposition 8 provides a sufficient optimality condition that the MLE lacks, since a vanishing gradient in the log likelihood function merely shows that there is a local optimum.

6 Results

In this section, tropical logistic regression is applied in three different scenarios. The first and simplest one is an illustration that considers datapoints generated from the tropical Laplace distribution. Secondly, the coalescent model is employed to generate gene trees from a species tree, and finally a real dataset of 1290 gene trees is considered. The models' performance on these datasets is examined.

6.1 Toy Example

In this example, a set of data points is generated from the tropical normal distribution as defined in Equation (7) using rejection sampling.

The data points are defined in the tropical projective torus $\mathbb{R}^e / \mathbb{R}\mathbf{1}$, which is isomorphic to \mathbb{R}^{e-1} . To map $x \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$ to \mathbb{R}^{e-1} , simply set the last component of x to 0, or in other words $x \mapsto (x_1 - x_e, x_2 - x_e, \dots, x_{e-1} - x_e)$. For illustration purposes, it is desirable to plot points

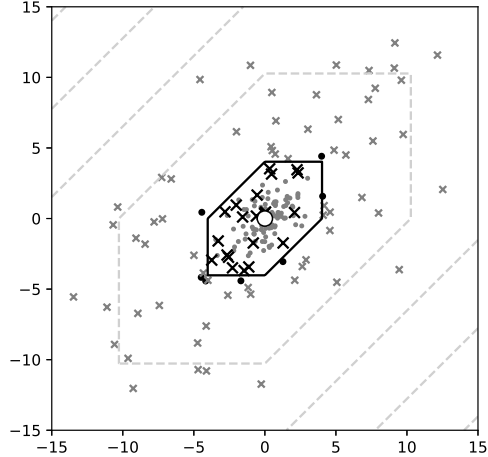


Figure 1: Scatterplot of 200 points - 100 dots for class 0 and 100 Xs for class 1, black for misclassified and grey otherwise - imposed upon a contour plot of the probability of inclusion in class 0, where the black contour is the classification threshold. The deviation parameters used in data generation were $\sigma_0 = 1, \sigma_1 = 5$ and the centre of the distribution (white-filled point) is the origin. The centres of the two distributions are $\omega_0 = \omega_1$.

in \mathbb{R}^2 , so $e = 3$ which corresponds to phylogenetic trees with 3 leaves. Both the one-species model and the two-species model are examined.

In the case of the former, $\omega = \omega_0 = \omega_1$ and $\sigma_0 \neq \sigma_1$. The classification boundary in this case is a tropical circle. If $\sigma_0 < \sigma_1$, the algorithm classifies points close to the inferred centre to class 0 and those that are more dispersed away from the centre as class 1. For simplicity, the centre is set to be the origin $\omega = (0, 0, 0)$ and no inference is performed. In Fig. 1 a scatterplot of the two classes is shown, where misclassified points are highlighted. As anticipated from Proposition 5 there are more misclassified points from the more dispersed class (class 1). Out of 100 points for each class, there are 7 and 21 misclassified points from class 0 and 1 respectively, while the theoretical probabilities calculated from equation (11) of Proposition 5 are 9% and 19% respectively.

Varying the deviation ratio σ_1/σ_0 in the data generation process allows exploration of its effect on the generalization error in the one-species model. The closer this ratio is to unity, the higher the generalization error. For $\sigma_0 = \sigma_1$ the classes are indistinguishable and hence any model is as good as a random guess i.e. the generalization error is $1/2$. The estimate of the generalization error for every value of that ratio is the proportion of misclassified points in both classes. Assuming an inferred ω that differs from the true parameter, Fig. 2(left) verifies the bounds of Proposition 5.

For the two-species model, tropical logistic regression is directly com-

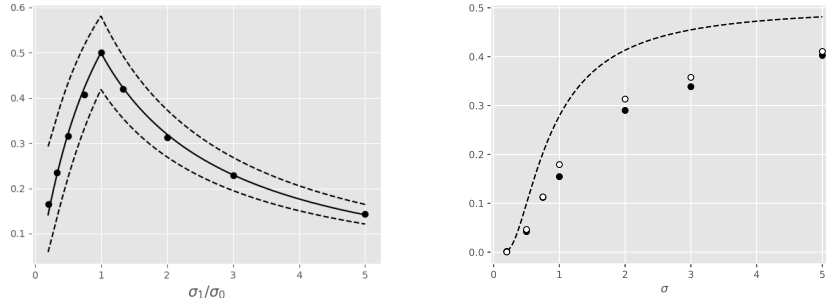


Figure 2: (left) Generalization error for 9 different deviation ratios. The estimator $\hat{\omega} = (0.3, 0, 3)$ differs from the true parameter $\omega = (0, 0)$. The upper and lower bounds of Proposition 5 are plotted in dashed lines and the generalization error for the correct estimator $\hat{\omega} = \omega^*$ plotted in solid line. The dots represent the proportion of misclassified points from a set of 2000 points in each experiment, 1000 points for each class. (right) Generalization errors for 7 different dispersion parameters with black markers for the two-species tropical logistic regression and white markers for the classical logistic regression. The upper bound (14) of Proposition 6 is plotted in dashed line.

pared to classical logistic regression. Data is generated using different centres $\omega_0 = (0, 0, 0)$, $\omega_1 = (3, 2, 0)$ but the same $\sigma = 0.5$. The classifier is $C(x) = \mathbb{I}(h(x) > 0)$ for both methods, using h as defined in equations (6) and (10) for the classical and tropical logistic regression respectively. Fig. 3 compares contours and classification thresholds of the classical (left) and tropical (right) logistic regression by overlaying them on top of the same data. Out of $100 + 100$ points there are $5 + 4$ and $4 + 3$ misclassifications in classical and tropical logistic regression respectively. Fig. 2(right) visualizes the misclassification rates of the two logistic regression methods for different values of dispersion σ , showing the tropical logistic regression to have consistently lower generalization error than the classical, even in this simple toy problem.

6.2 Coalescent Model

In Bayesian inference of phylogeny via MCMC, it is important to be able to recognise whether the chain of trees has converged. MCMC convergence can be tested using any classification algorithm. In Bayesian inference of phylogenetic trees generated from some distribution conditional on the species tree, the chain has likely converged if it is hard to distinguish it from another independent chain of the same length, generated from the same distribution [9, 8]. On the other hand, if it is easy to classify the generated trees coming from two different chains, then the chains have not converged yet.

However, instead of using trees from MCMC, the data that have been used in our simulations were generated under the multispecies coalescent

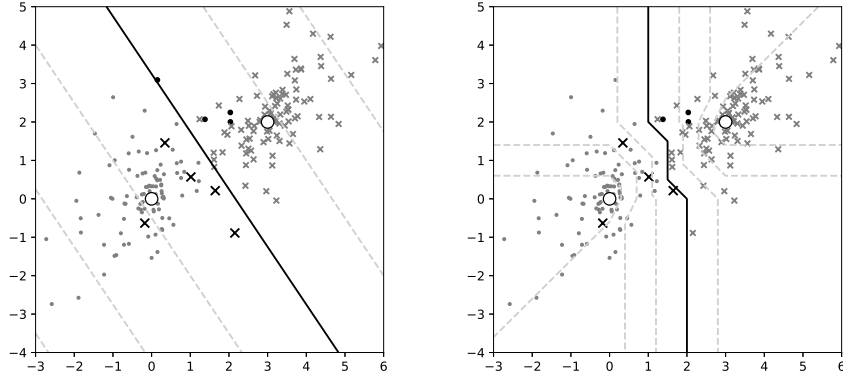


Figure 3: Scatterplot of points - dots for class 0 and X for class 1, black for misclassified according to (left) **classical logistic regression** or (right) **tropical logistic regression**, and grey otherwise - alongside a contour plot of the probabilities, where the black contour is the classification threshold. The centres, drawn as big white dots, are $\omega_0 = (0, 0, 0)$, $\omega_1 = (3, 2, 0)$ and $\sigma = 0.5$.

model, using the python library `dendropy` [23]. The classification method we propose is the two-species model because two distinct species tree have been used to generate gene tree data for each class.

Two distinct species trees are used, which are randomly generated under a Yule model. Then, using `dendropy`, 1000 gene trees are randomly generated for each of the two species. The trees have 10 leaves and so the number of the model variables is $\binom{10}{2} = 45$. They are labelled according to the species tree they are generated from. The tree generation is under the coalescent model for specific model parameters.

Since the species trees are known, we conduct a comparative analysis between classical, tropical and BHV ([5]) logistic regression. In the supplement, we show an approximation analog of our model to the BHV metric. The comparative analysis includes the distribution fitting of distances and the misclassification rates for different metrics.

In Fig. 4, the distribution of the radius $d(X, \omega)$ as given by Proposition 4, is fitted to the histograms of the Euclidean and tropical distances of gene trees to their corresponding species tree. According to Proposition 4, for both the classical and tropical Laplace distributed covariates, $d(X, \omega^*) \sim \sigma \text{Gamma}(n)$, shown in solid lines in Fig. 4, where $n = e = 45$ and $n = e - 1 = 44$ for the classical and tropical case respectively. Similarly, for normally distributed covariates, $d(X, \omega^*) \sim \sigma \sqrt{\chi_n^2}$, shown in dashed lines. It is clear that Laplacian distributions produce better fits in both geometries and that the tropical Laplacian fits the data best. This is reflected in the values of the average log-likelihood summarised in Table 1. As discussed in Section 4, the same analysis can not be applied to the BHV metric, because the condition of Proposition 4 does not hold.

	Laplace	Normal
Euclidean	-1.9	-2.5
Tropical	-1.1	-1.3

Table 1: Average log-likelihoods for distributions fitted in Fig. 4. The tropical Laplace distribution on which our model is based on, fits the data best.

Species depth SD is the time since the speciation event between the species and *effective population size* N quantifies genetic variation in the species. Datasets have been generated for a range of values $R := SD/N$ by varying species depth. For low values of R , speciation happens very recently and so the gene trees look very much alike. Hence, classification is hard for datasets with low values of R and vice versa, because the gene deviation σ_R is a decreasing function of R . We expect classification to improve in line with R . Fig. 9 and Fig. 8 in Supplement H confirm that, by showing that as R increases the receiver operating characteristic (ROC) curves are improving and the Robinson-Foulds and tropical distances of inferred (Fermat-Weber point) trees are decreasing. In addition, Fig. 5 shows that as R increases, AUCs increase (left) and misclassification rates decrease (right). It also shows that tropical logistic regression produces higher AUCs than classical logistic regression and lower misclassification rates than both the BHV and classical logistic regression. Finally, note that the generalization error upper bound as given in equation (14) is satisfied but it not very tight (dashed line in Fig. 5).

6.3 Lungfish Dataset

In this last application, we consider an empirical dataset of 1290 gene (loci) alignments of 10 species from [10] and reconstructed 1290 corresponding gene trees in [28]. Two different methods have been applied for tree reconstruction; the standard maximum likelihood estimator (MLE) method and the neighbor-joining (NJ) method [20]. After dimensionality reduction using principal component analysis, in [28] Yoshida et al. applied three different clustering methods and observed that normalised cuts is the best performing method. Using the two clusters as labels, our model is implemented to test how easy it is to differentiate the two clusters from each other. Failing to differentiate the clusters is indicative of poor performance in clustering. Five clustering cases have been considered for the MLE and NJ reconstructed trees, for each dimensionality reduction technique. The authors in [28] observed that the clustering algorithms performed better on the NJ gene trees than the MLE gene trees. Our model confirms that by comparing AUC values.

The one-species model is not applicable in this case, as the two clusters have different centers. The AUC values obtained for the one-species model are quite poor, around 60%. However, applying the two-species model results in significantly better performance. For the five clusterings of NJ reconstructed trees, the range of AUCs is between 96% and 98%, while the range for the MLE trees is much lower, ranging from 67% to 77%. These results indicate that the clustering methods implemented in [28]

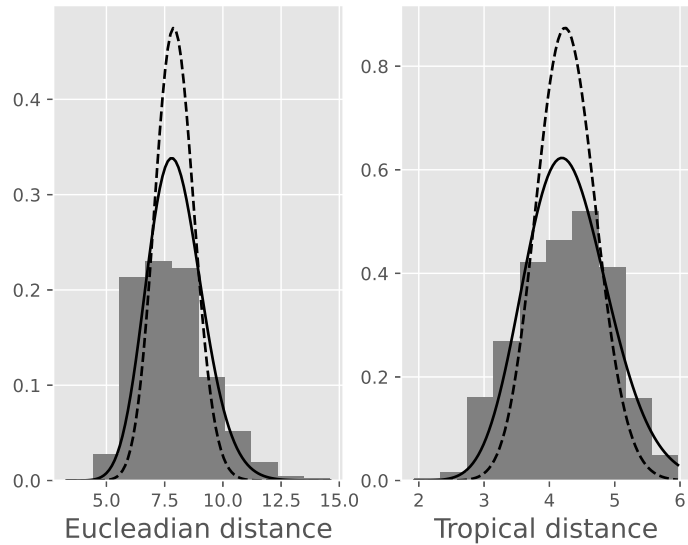


Figure 4: Histograms of the distances of 1000 gene trees from the species trees that generated them under the coalescent model with $R = 0.7$. The solid and dashed lines are fitted distributions $\sigma\text{Gamma}(n)$ and $\sigma\sqrt{\chi_n^2}$ respectively; σ is chosen to be the MLE, derived in the supplement. (left) Euclidean metric has worse fit than (right) the tropical metric. In both cases, log-likelihoods with quadratic distances (dashed) perform worse than linear distances (solid).

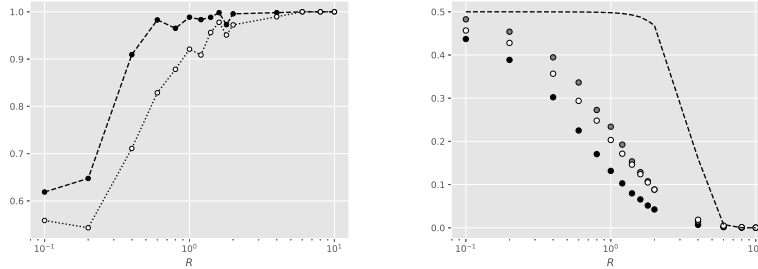


Figure 5: (left) A dashed line is for AUCs of the tropical logistic regression under two species-tree model and a grey dotted line for AUCs of the classical logistic regression model with simulated dataets generated under the multi-species coalescent model. It is important to state that if R is less than 1, it is very difficult to classify and the tropical logistic regression works significantly better than the classical logistic regression model. (right) the x-axis represents the ratio R and the y-axis represents misclassification rates. Black circles represent the tropical logistic regression, white circles represent the classical logistic regression, grey points represent the logistic regression with BHV metric, and the dashed line represents the theoretical generalization error shown in Proposition 6.

perform better on the former (NJ reconstructed trees) compared to the latter (MLE trees).

Finally, the Fermat-Weber points of the MLE reconstructed trees are computed for both the one-species and two-species model and projected onto the space of ultrametrics via complete-linkage hierarchical clustering. The tree topologies of the resulting trees matched (shown in Fig. 6) and are almost entirely in agreement with that of the inferred tree from [10]. However, lungfish and coelacanth seem to be equally away from tetrapods. Nonetheless, this tree topology is also proposed by many researchers [30, 21].

7 Discussion

In this paper we developed an analogue of the classical logistic regression model and considered two special cases; the one species-tree model and two species-tree model. Even if the former was not suitable for the datasets considered in this paper, it could still be useful in other settings. The main benefit is having the same number of parameters as the number of predictors, unlike the two-species model which has almost twice as many. Therefore, it fits the standard definition of a generalized linear model and could even generalize to a stack of GLMs to produce a "tropical" neural network.

The two-species model implemented on data generated under the coalescent model outperformed classical and BHV logistic regression models on misclassification rates, AUCs and fitness of the distribution of distances

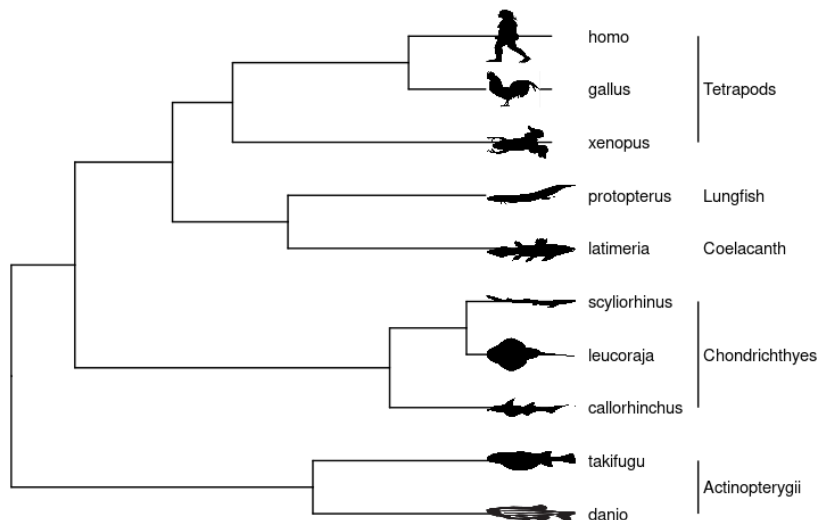


Figure 6: Inferred species tree using the tropical one-species model. Using the Fermat-Weber point of all gene trees as a distance matrix, complete-linkage clustering produces a dendrogram, which is the equidistant tree above.

from their centre. It was also observed that Laplacian distributions were more fit than Gaussians, for both geometries. Further research on the generalization error for the two-species model would provide tighter bounds.

Our model was applied to the lungfish dataset and verified the claim that the clustering method of [28] performed better for NJ trees than MLE trees. The inferred species tree of the MLE trees agrees with literature. However, this is not the case with NJ trees, whose species trees were topologically far from each other. Perhaps this incongruence has to do with the fact that NJ trees are unrooted.

Finally, the Fermat-Weber point is not always an equidistant tree and has to be projected onto the space of ultrametric trees. This is the case even when all datapoints are ultrametric, as in the case of coalescent model data. This projection was performed via hierarchical clustering and it was observed that complete-linkage clustering produced the smallest Robinson-Foulds distances between the projected tree and the true species tree. An interesting extension would be to consider properties of Fermat-Weber points when the sample is ultrametric and to investigate the best way of projecting them onto the space of ultrametrics.

Funding

RY is partially funded by NSF Division of Mathematical Sciences: Statistics Program DMS 1916037. GA is funded by EPSRC through the STOR-i Centre for Doctoral Training under grant EP/L015692/1.

References

- [1] M. Akian, S. Gaubert, Y. Qi, and O. Saadi. Tropical linear regression and mean payoff games: or, how to measure the distance to equilibria, 2021. <https://arxiv.org/abs/2106.01930>.
- [2] C Ané, B Larget, DA Baum, SD Smith, and A Rokas. Bayesian estimation of concordance among gene trees. *Mol Biol Evol.*, 24(2):412–26, 2007.
- [3] F. Ardila and C. J. Klivans. The bergman complex of a matroid and phylogenetic trees. *journal of combinatorial theory. Series B*, 96(1):38–49, 2006.
- [4] Herman J Bierens. *Topics in advanced econometrics: estimation, testing, and specification of cross-section and time series models*. Cambridge University Press, 1996.
- [5] L.J. Billera, S.P. Holmes, and K. Vogtmann. Geometry of the space of phylogenetic trees. *Adv Appl Math*, 27(4):733–767, 2001.
- [6] P. Buneman. A note on the metric properties of trees. *J. Combinatorial Theory Ser. B.*, 17:48–50, 1974.
- [7] Francisco Criado, Michael Joswig, and Francisco Santos. Tropical bisectors and voronoi diagrams. *Foundations of Computational Mathematics*, pages 1–38, 2021.
- [8] L. Guimarães Fabreti and S. Höhna. Convergence assessment for bayesian phylogenetic analysis using mcmc simulation. *Methods in Ecological Society*, pages <https://doi.org/10.1111/2041-210X.13727>, 2021.
- [9] D.C. Haws, P. Huggins, E.M. O’Neill, D.W. Weisrock, and R. Yoshida. A support vector machine based test for incongruence between sets of trees in tree space. *BMC Bioinformatics*, 13(210), 2012. <https://doi.org/10.1186/1471-2105-13-210>.
- [10] Dan Liang, Xing Xing Shen, and Peng Zhang. One thousand two hundred ninety nuclear genes from a genome-wide survey support lungfishes as the sister group of tetrapods. *Molecular biology and evolution*, 30(8):1803–1807, 2013.
- [11] B. Lin, B. Sturmfels, X. Tang, and R. Yoshida. Convexity in tree spaces. *SIAM Discrete Math*, 3:2015–2038, 2017.
- [12] Bo Lin and Ruriko Yoshida. Tropical fermat–weber points. *SIAM Journal on Discrete Mathematics*, 32(2):1229–1245, 2018.
- [13] D. Maclagan and B. Sturmfels. *Introduction to Tropical Geometry*, volume 161 of *Graduate Studies in Mathematics*. Graduate Studies in Mathematics, 161, American Mathematical Society, Providence, RI, 2015.
- [14] W. P. Maddison and D.R. Maddison. Mesquite: a modular system for evolutionary analysis. version 2.72, 2009. Available at <http://mesquiteproject.org>.
- [15] Wayne P Maddison. Mesquite: a modular system for evolutionary analysis. *Evolution*, 62:1103–1118, 2008.

- [16] Whitney K Newey and Daniel McFadden. Large sample estimation and hypothesis testing. *Handbook of econometrics*, 4:2111–2245, 1994.
- [17] JA Nylander, JC Wilgenbusch, DL Warren, and DL Swofford. AWTY (are we there yet?): a system for graphical exploration of mcmc convergence in bayesian phylogenetics. *Bioinformatics*, 24(4):581–3, 2008.
- [18] Robert Page, Ruriko Yoshida, and Leon Zhang. Tropical principal component analysis on the space of phylogenetic trees. *Bioinformatics*, 36(17):4590–4598, 06 2020.
- [19] Jean-Eric Pin. Tropical semirings, 1998.
- [20] Naruya Saitou and Masatoshi Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular biology and evolution*, 4(4):406–425, 1987.
- [21] Yunfeng Shan and Robin Gras. 43 genes support the lungfish-coelacanth grouping related to the closest living relative of tetrapods with the bayesian method under the coalescence model. *BMC Research Notes*, 4(1):1–6, 2011.
- [22] D. Speyer and B. Sturmfels. Tropical mathematics. *Mathematics Magazine*, 82:163–173, 2009.
- [23] Jeet Sukumaran and Mark T Holder. Dendropy: a python library for phylogenetic computing. *Bioinformatics*, 26(12):1569–1571, 2010.
- [24] R. Yoshida, K. Miura, and D. Barnhill. Hit and run sampling from tropically convex sets, 2022. <https://arxiv.org/abs/2209.15045>.
- [25] R. Yoshida, K. Miura, D. Barnhill, and D. Howe. Tropical density estimation of phylogenetic trees, 2022. <https://arxiv.org/abs/2206.04206>.
- [26] R. Yoshida, M. Takamori, H. Matsumoto, and K. Miura. Tropical support vector machines: Evaluations and extension to function spaces, 2021. <https://arxiv.org/abs/2101.11531>.
- [27] R. Yoshida, L. Zhang, and X. Zhang. Tropical principal component analysis and its application to phylogenetics. *Bulletin of Mathematical Biology*, 81:568–597, 2019.
- [28] Ruriko Yoshida, Kenji Fukumizu, and Chrysafis Vogiatzis. Multilocus phylogenetic analysis with gene tree clustering. *Annals of Operations Research*, 276:293–313, 2019.
- [29] Ruriko Yoshida, Keiji Miura, and David Barnhill. Hit and run sampling from tropically convex sets. *arXiv preprint arXiv:2209.15045*, 2022.
- [30] Min Zhu, Xiaobo Yu, and Per E Ahlberg. A primitive sarcopterygian fish with an eyestalk. *Nature*, 410(6824):81–84, 2001.

A Proofs

Proof of Lemma 3.1. A simple application of the Bayes rule for continuous random variables yields

$$\begin{aligned} p(x) = \mathbb{P}(Y = 1|X = x) &= \frac{f_1(x)\mathbb{P}(Y = 1)}{f_0(x)\mathbb{P}(Y = 0) + f_1(x)\mathbb{P}(Y = 1)} \\ &= \frac{1}{1 + \frac{f_1(x)(1-r)}{f_0(x)r}} = S(h(x)). \end{aligned}$$

□

Proof of Proposition 3. The expected log-likelihood is expressed as

$$\begin{aligned} \mathbb{E}(l) &= \mathbb{E}(Y \log(p(X)) + (1 - Y) \log(1 - p(X))) \\ &= \mathbb{P}(Y = 1) \int_{\mathbb{R}^n} f_1(x) \log(p(x)) dx \\ &\quad + \mathbb{P}(Y = 0) \int_{\mathbb{R}^n} f_0(x) \log(1 - p(x)) dx \\ &= \int_{\mathbb{R}^n} L(x, p(x)) dx, \end{aligned}$$

where $L(x, p) = r f_1(x) \log(p) + (1 - r) f_0(x) \log(1 - p)$ is treated as the Lagrangian. The Euler-Lagrange equation can be generalized to a several variables (in our case there are n variables). Since there are no derivatives of p , the stationary functional satisfies $\partial_p L = 0$, which yields the desired result. □

Proof of Proposition 4. The pdf of X is

$$f_\omega(x) = \frac{1}{C_\alpha} \exp\left(-\alpha^i \frac{d^i(x)}{i}\right), x \in \mathbb{R}^n$$

where $\alpha = \sigma^{-1}$ is the precision. Using the variable transformation $y = \alpha x$ with Jacobian $1/\alpha^n$ and remembering that $\alpha d(x) = d(y)$,

$$C_\alpha = \int_{\mathbb{R}^n} \exp\left(-\alpha^i \frac{d^i(x)}{i}\right) dx = \int_{\mathbb{R}^n} \exp\left(-\frac{d^i(x)}{i}\right) \frac{dy}{\alpha^n} = \frac{C_1}{\alpha^n}.$$

The moment generating function of $d^i(X)$ is

$$\begin{aligned} M_{d^i(X)} &= \int_{\mathbb{R}^n} \exp\left(z d^i(x)\right) \frac{\exp\left(-\alpha^i \frac{d^i(x)}{i}\right)}{C_\alpha} dx \\ &= \frac{C_{\sqrt{\alpha^i/i-z}}}{C_\alpha} = \frac{1}{(\sqrt{1-i\sigma^i z})^n}, \end{aligned}$$

which coincides with the MGF of $\Gamma(n/i, i\sigma^i)$. □

Proof of Proposition 2. From the proof of Proposition 4, it was established that the normalizing constant is $C_{\sigma_Y} = C_1 \sigma_Y^{e-1}$ for the tropical projective torus, whose dimension is $n = e - 1$.

The volume of a unit tropical sphere in the tropical projective torus $\mathbb{R}^e/\mathbb{R}\mathbf{1}$ is equal to e . If the tropical radius is r , then the volume is er^{e-1} and hence the surface area is $e(e-1)r^{e-2}$. Therefore,

$$\begin{aligned} C_1 &= \int_{\mathbb{R}^e/\mathbb{R}\mathbf{1}} \exp(-d_{\text{tr}}(x, \mathbf{0})) dx \\ &= \int_0^\infty e(e-1)r^{e-2} \exp(-r) dr \\ &= e(e-1)\Gamma(e-1) = e! \end{aligned}$$

It follows that the normalizing constant is $C_{\sigma_Y} = e! \sigma_Y^{e-1}$. \square

Proof of Corollary 1. Suppose that X comes from the Laplace or the Normal distribution, whose pdf is proportional to $\exp(-d^i(x, \omega^*)/(i\sigma^i))$ for $i = 1$ and 2 respectively, for all $x \in \mathbb{R}^n$ where d is the Euclidean metric. Then, $X - \omega^*$ has a distribution proportional to $\exp(-d^i(x, \mathbf{0})/(i\sigma^i))$. Clearly, $\alpha d(x, \mathbf{0}) = d(\alpha x, \mathbf{0})$ and so from Proposition 4, it follows that $d^i(X - \omega^*, \mathbf{0}) = d^i(X, \omega^*) \sim i\sigma^i \text{Gamma}(n/i)$. Note that for the normal distribution ($i = 2$), $d^i(X, \omega^*) \sim \sigma^2 \chi_{n/2}$. The same argument applies for tropical Laplace and tropical Normal distributions, where the metric is tropical ($d = d_{\text{tr}}$), the distribution is defined on $\mathbb{R}^e/\mathbb{R}\mathbf{1} \cong \mathbb{R}^{e-1}$ and the dimension is hence $n = e - 1$. \square

Prerequisites for proof of Theorem 4.1

Theorem A.1. (Theorem 4.2.1 in [4]) Let $(Q_n(\theta))$ be a sequence of random functions on a compact set $\Theta \subset \mathbb{R}^m$ such that for a continuous real function $Q(\theta)$ on Θ ,

$$\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| \xrightarrow{P} 0 \text{ as } n \rightarrow \infty.$$

Let θ_n be any random vector in Θ satisfying $Q_n(\theta_n) = \inf_{\theta \in \Theta} Q_n(\theta)$ and let θ_0 be a unique point in Θ such that $Q(\theta_0) = \inf_{\theta \in \Theta} Q(\theta)$. Then $\theta_n \xrightarrow{P} \theta_0$.

Theorem A.2. (Lemma 2.4 in [16]) If the data z_1, \dots, z_n are independent and identically distributed, the parameter space Θ is compact, $f(z_i, \theta)$ is continuous at each $\theta \in \Theta$ almost surely and there is $r(z) \geq |f(z, \theta)|$ for all $\theta \in \Theta$ and $\mathbb{E}(r(z)) < \infty$, then $\mathbb{E}(f(z, \theta))$ is continuous and

$$\sup_{\theta \in \Theta} \left| n^{-1} \sum_{i=1}^n f(z_i, \theta) - \mathbb{E}(f(z, \theta)) \right| \xrightarrow{P} 0.$$

Lemma A.3. Consider two points $x, y \in \mathbb{R}^e/\mathbb{R}\mathbf{1}$. There exists $\eta > 0$ such that

$$d_{\text{tr}}(x + \epsilon E_i, y) = d_{\text{tr}}(x, y) + \epsilon \phi_i(x - y), \forall \epsilon \in [0, \eta], \forall i \in [e], \text{ where}$$

$$\phi_i(v) = \begin{cases} 1, & \text{if } v_i \geq v_j \forall j \in [e] \\ -1, & v_i < v_j \forall j \in [e] \setminus \{i\}, \\ 0, & \text{otherwise} \end{cases} \quad (16)$$

and $E_i \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$ is a vector with 1 in the i -th coordinate and 0 elsewhere.

Proof. By setting $v := x - y$, $M := \max_{j \in [e]} \{v_j\}$ and $m := \min_{j \in [e]} \{v_j\}$,

$$\begin{aligned} d_{\text{tr}}(x, y) &= M - m \\ d_{\text{tr}}(x + \epsilon E_i, y) &= \max_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} - \min_{j \in [e]} \{v_j + \epsilon \delta_{ij}\}, \end{aligned}$$

where $\epsilon \geq 0$, and $\delta_{ij} = \mathbb{I}(i = j)$ with \mathbb{I} being the indicator function. Three separate cases are considered.

i. If $v_i = M$, then

$$\max_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} = v_i + \epsilon = M + \epsilon, \quad (17)$$

$$\min_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} = m, \quad (18)$$

and so $d_{\text{tr}}(x + \epsilon E_i, y) = d_{\text{tr}}(x, y) + \epsilon$. Note that equations (17) and (18) hold for all $\epsilon > 0$.

ii. If $v_i = m$ **and** $v_i < v_k$ for all $k \neq i$, i.e. if v_i is the **unique** minimum component of vector v , then

$$\max_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} = M, \text{ for all } \epsilon \leq M - m \quad (19)$$

$$\min_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} = v_i + \epsilon = m + \epsilon, \text{ for all } \epsilon \leq m' - m, \quad (20)$$

where $m' := \min_{j: v_j > m} \{v_j\} > m$ is well-defined unless $v_j = m$ for all $j \in [e]$ i.e. for $v = m \cdot (1, \dots, 1) = \mathbf{0}$, which falls under the first case. Clearly, $M \geq m'$, so for all $\epsilon \in [0, m' - m]$ equations (19) and (20) are satisfied and thence $d_{\text{tr}}(x + \epsilon E_i, y) = d_{\text{tr}}(x, y) - \epsilon$.

iii. Otherwise, if none of the first two cases hold then $\exists k \neq i$ such that $m = v_k \leq v_i < M$ and so

$$\min_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} = v_k = m, \text{ for all } \epsilon > 0 \quad (21)$$

$$\max_{j \in [e]} \{v_j + \epsilon \delta_{ij}\} = M, \text{ if } \epsilon \leq M - v_i \quad (22)$$

Define $M' := \max_{j: v_j < M} \{v_j\} < M$ which is well-defined for all $v \neq \mathbf{0}$ (first case). Since $v_i < M$, it follows by definition that $v_i \leq M'$ and so $M - v_i \geq M - M' > 0$. As a result, for all $\epsilon \in [0, M - M']$, equations (21) and (22) are satisfied and thence $d_{\text{tr}}(x + \epsilon E_i, y) = d_{\text{tr}}(x, y)$.

If $v = \mathbf{0}$, set $\eta = +\infty$. Otherwise, for $v \neq \mathbf{0}$, with m', M' being well-defined, set

$$\eta = \min(m' - m, M - M') > 0.$$

In all three cases and for all $\epsilon \in [0, \eta]$ the desired result is satisfied. \square

Lemma A.4. Consider the function $q : \mathbb{R}^e / \mathbb{R}\mathbf{1} \rightarrow \mathbb{R}$,

$$q(x) = \lambda_\alpha d_{\text{tr}}(x, \alpha) - \lambda_\beta d_{\text{tr}}(x, \beta) - \lambda_\gamma d_{\text{tr}}(x, \gamma) + \lambda_\delta d_{\text{tr}}(x, \delta) \\ + \log\left(\frac{\lambda_\beta}{\lambda_\alpha}\right) - \log\left(\frac{\lambda_\delta}{\lambda_\gamma}\right),$$

where $\alpha, \beta, \gamma, \delta \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$, $\lambda_\alpha, \lambda_\beta, \lambda_\gamma, \lambda_\delta > 0$ and $(\alpha, \lambda_\alpha) \neq (\beta, \lambda_\beta)$. A set \mathcal{X} contains neighbourhoods of $\alpha, \beta, \gamma, \delta$. If $q(x) = 0, \forall x \in \mathcal{X}$ then $(\alpha, \lambda_\alpha) = (\gamma, \lambda_\gamma)$ and $(\beta, \lambda_\beta) = (\delta, \lambda_\delta)$.

Proof. According Lemma A.3, there exists $\eta_1 > 0$ such that for all $\epsilon \in [0, \eta_1]$

$$d_{\text{tr}}(x + \epsilon E_i, y) = d_{\text{tr}}(x, y) + \epsilon \phi_i(x - y). \quad (23)$$

Moreover, $d_{\text{tr}}(x - \epsilon E_i, y) = d_{\text{tr}}(y, x - \epsilon E_i) = d_{\text{tr}}(y + \epsilon E_i, x)$ and so using Lemma A.3 again (but with x and y swapped), there exists $\eta_2 > 0$ such that for all $\epsilon \in [0, \eta_2]$

$$d_{\text{tr}}(x - \epsilon E_i, y) = d_{\text{tr}}(x, y) + \epsilon \phi_i(y - x), \quad (24)$$

for all $\epsilon \in [0, \epsilon_0(y - x)]$. For all $\epsilon \in [0, \eta]$ where $\eta := \min(\eta_1, \eta_2)$, equations (23), (24) are satisfied and so

$$q(x + \epsilon E_i) = q(x) + \\ \epsilon (\lambda_\alpha \phi_i(x - \alpha) - \lambda_\beta \phi_i(x - \beta) - \lambda_\gamma \phi_i(x - \gamma) + \lambda_\delta \phi_i(x - \delta)), \\ q(x - \epsilon E_i) = q(x) + \\ \epsilon (\lambda_\alpha \phi_i(\alpha - x) - \lambda_\beta \phi_i(\beta - x) - \lambda_\gamma \phi_i(\gamma - x) + \lambda_\delta \phi_i(\delta - x)).$$

Consequently, for all $\epsilon \in [0, \eta]$,

$$q(x + \epsilon E_i) + q(x - \epsilon E_i) - q(x) = 0 \quad (25) \\ = \epsilon (\lambda_\alpha s_i(x - \alpha) - \lambda_\beta s_i(x - \beta) - \lambda_\gamma s_i(x - \gamma) + \lambda_\delta s_i(x - \delta)),$$

where

$$s_i(v) := \phi_i(v) + \phi_i(-v) = \quad (26) \\ \begin{cases} 2, & \text{if } v = \mathbf{0} \\ 1, & \text{if } v \neq \mathbf{0} \text{ and } v_i \text{ is the non-unique maximizer or minimizer of } v \\ 0, & \text{otherwise} \end{cases}$$

By summing equation (25) over $i \in [e]$ and defining $s(v) = \sum_{i=1}^e s_i(v)$,

$$\lambda_\alpha s(x - \alpha) - \lambda_\beta s(x - \beta) - \lambda_\gamma s(x - \gamma) + \lambda_\delta s(x - \delta) = 0, \quad (27)$$

$\forall x \in \mathcal{X}$.

Here we try to prove by contradiction that $\mathcal{S} := \{\alpha, \delta\} \cap \{\gamma, \beta\}$ is not empty. Suppose that $\mathcal{S} := \{\alpha, \delta\} \cap \{\gamma, \beta\} = \emptyset$. Then, setting $x = \alpha$ in equation (27) and noting that $s(0) = 2e$ and $0 \leq s(v) \leq e$ for $v \neq 0$, we

get $2e\lambda_\alpha \leq e\lambda_\beta + e\lambda_\gamma$, since $\beta, \gamma \neq \alpha$. Applying the same argument to $x = \beta, \gamma, \delta$, the following system of inequalities holds

$$\begin{aligned} 2\lambda_\alpha &\leq \lambda_\beta + \lambda_\gamma \\ 2\lambda_\beta &\leq \lambda_\alpha + \lambda_\delta \\ 2\lambda_\gamma &\leq \lambda_\alpha + \lambda_\delta \\ 2\lambda_\delta &\leq \lambda_\beta + \lambda_\gamma. \end{aligned}$$

It follows that $\lambda_\alpha = \lambda_\beta = \lambda_\gamma = \lambda_\delta$. Then, rewrite equation (27) as

$$s(x - \alpha) - s(x - \beta) - s(x - \gamma) + s(x - \delta) = 0, \quad (28)$$

Note now equation (28) can only hold at $x = \alpha$ iff $s(\alpha - \gamma) = s(\alpha - \beta) = e$ and $s(\alpha - \delta) = 0$. But $s(v) = e$ if and only if all the components of v are non-unique minimizers and maximizers or $\{v_i : i \in [e]\} = \{\zeta, \kappa\}$, where $\zeta < \kappa$ and $|\{i : v_i = \zeta\}| = n_\zeta, |\{i : v_i = \kappa\}| = n_\kappa$, such that $n_\zeta + n_\kappa = e$ and $n_\zeta, n_\kappa \geq 2$.

Consider $z = v + \epsilon E_i$, where $v_i = \zeta$ and $0 < \epsilon < \kappa - \zeta$. The minimum and maximum components of z are ζ and κ , and $\{z_i : i \in [e]\} = \{\zeta, \zeta + \epsilon, \kappa\}$ with $|\{i : z_i = \zeta\}| = n_\zeta - 1, |\{i : z_i = \kappa\}| = n_\kappa$. It follows that,

$$s(z) = |\{i : z_i = \zeta\}| + |\{i : z_i = \kappa\}| = e - 1.$$

Now consider $z = v + \epsilon E_i$ where $v_i = \kappa$. The maximum is no longer unique, but the n_ζ minima are still unique. Therefore, $s(z) = n_\zeta \geq 2$. Combining the two cases, it is concluded that $s(v + \epsilon E_i) \geq 2$ for all $i \in [e]$.

Set $x = \alpha + \epsilon E_i$, where $\alpha_i - \beta_i = \min_k \{\alpha_k - \beta_k\}$. Then,

$$s(x - \alpha) = s(\epsilon E_i) = e - 1, \quad (29)$$

since there is a unique maximizer, but all the other $e - 1$ components are 0, which is the minimum. Furthermore,

$$s(x - \beta) = s(\alpha - \beta + \epsilon E_i) = e - 1, \quad (30)$$

since for $v = \alpha - \beta$ with $s(v) = e$, it corresponds to the first case examined. It is assumed that $\epsilon < \kappa - \zeta = d_{\text{tr}}(\alpha - \beta)$. Moreover,

$$s(x - \gamma) = s(\alpha - \gamma + \epsilon E_i) \geq 2, \quad (31)$$

for $v = \alpha - \gamma$ with $s(v) = e$. Finally, since $s(\alpha - \delta) = 0$ and so the components of $\alpha - \delta$ have a unique minimum and a unique maximum, there exists a neighborhood around $x = \alpha$ such that $x - \alpha$ still has that property, i.e.

$$s(x - \delta) = s(\alpha - \delta + \epsilon E_i) = 0 \quad (32)$$

for all $\epsilon < \eta$ for some $\eta > 0$.

From equations (29) – (32), it is concluded that

$$s(x - \alpha) - s(x - \beta) - s(x - \gamma) + s(x - \delta) \leq -2, \quad (33)$$

which contradicts equation (28). Therefore $\mathcal{S} = \{\alpha, \delta\} \cap \{\gamma, \beta\} \neq \emptyset$.

Define another set $\mathcal{T} = \{\alpha, \beta, \gamma, \delta\}$. Since $\mathcal{S} \neq \emptyset$, $|\mathcal{T}| \leq 3$. Suppose that $|\mathcal{T}| = 3$ with $\mathcal{T} = \{\tau, \nu, \phi\}$. Then, without loss of generality equation (27) becomes

$$\lambda_\tau s(x - \tau) + \lambda_\nu s(x - \nu) - \lambda_\phi s(x - \phi) = 0 \quad (34)$$

Similarly to before, setting $x = \tau, \nu, \phi$ yields,

$$\begin{aligned} 2\lambda_\tau &\leq \lambda_\phi \\ 2\lambda_\nu &\leq \lambda_\phi \\ 2\lambda_\phi &\leq \lambda_\tau + \lambda_\nu, \end{aligned}$$

which is contradictory since $\lambda_\tau + \lambda_\nu > 0$. Therefore, $|\mathcal{T}| \leq 2$. There are 4 cases to consider

- i. $\alpha = \delta \neq \beta = \gamma$, but then $\mathcal{S} = \emptyset$,
- ii. $\alpha = \beta \neq \gamma = \delta$, but then equation (27) can only be satisfied $x = \alpha, \gamma$ if $\lambda_\alpha = \lambda_\beta$ and $\lambda_\gamma = \lambda_\delta$ which violates the statement that $(\alpha, \lambda_\alpha) \neq (\beta, \lambda_\beta)$,
- iii. $\alpha = \gamma \neq \beta = \delta$ and from equation (27) at $x = \alpha, \gamma$ it follows that $\lambda_\alpha = \lambda_\gamma, \lambda_\beta = \lambda_\delta$ and hence the desired result,
- iv. $\alpha = \beta = \gamma = \delta$, in which case

$$q(x) = (\lambda_\alpha - \lambda_\beta - \lambda_\gamma + \lambda_\delta) d_{\text{tr}}(x, \alpha) + \log\left(\frac{\lambda_\beta}{\lambda_\alpha}\right) - \log\left(\frac{\lambda_\delta}{\lambda_\gamma}\right),$$

which can only be uniformly 0 at \mathcal{X} if and only if $\lambda_\alpha + \lambda_\delta = \lambda_\beta + \lambda_\gamma$. Observe that $(\lambda_\alpha, \lambda_\delta)$ and $(\lambda_\beta, \lambda_\gamma)$ are the two roots of the same quadratic $z^2 - (\lambda_\alpha + \lambda_\delta)z + \lambda_\alpha \lambda_\delta$ and noting that in this case $\lambda_\alpha \neq \lambda_\beta$, it follows that $\lambda_\alpha = \lambda_\gamma$ and $\lambda_\beta = \lambda_\delta$.

□

Lemma A.5. Consider a compact set $\Sigma \subseteq \mathbb{R}_+ = (0, \infty)$. Then the set $\Lambda = \{\sigma^{-1} : \sigma \in \Sigma\} \in \mathbb{R}_+$ is also compact.

Proof. In metric spaces, a set is compact iff it is sequentially compact. Therefore, for every sequence $\sigma_n \in \Sigma$, $\sigma_n \rightarrow \sigma \in \Sigma$. Every sequence in Λ can be expressed as $1/\sigma_n$, which tends to $1/\sigma \in \Lambda$. Therefore, Λ is sequentially compact and hence compact. □

Proof of Theorem 4.1. This proof has been written for precision estimators $\lambda = 1/\sigma$ instead of deviation estimators. For the rest of the proof consider $\lambda_y = \sigma_y^{-1}$ for $y = 0, 1$ and define the set

$$\Lambda = \{\sigma^{-1} : \sigma \in \Sigma\} \in \mathbb{R}_+.$$

According to Lemma A.5, Λ is also compact. Define the function f and h as

$$\begin{aligned} f &: \mathbb{R}^e / \mathbb{R}\mathbf{1} \times \{0, 1\} \times \Omega^2 \times \Lambda^2 \rightarrow \mathbb{R}, \\ f((x, y), (\omega, \lambda)) &= y \log S(h(x, (\omega, \lambda))) + (1 - y) \log S(-h(x, (\omega, \lambda))), \\ h &: \mathbb{R}^e / \mathbb{R}\mathbf{1} \times \Omega^2 \times \Lambda^2 \rightarrow \mathbb{R}, \\ h(x, (\omega, \lambda)) &= \lambda_0 d_{\text{tr}}(x, \omega_0) - \lambda_1 d_{\text{tr}}(x, \omega_1) + (e - 1) \log \frac{\lambda_1}{\lambda_0}, \end{aligned}$$

where S is the logistic function. Also denote the empirical (Q_n) and expected (Q) log-likelihood functions as

$$\begin{aligned} Q_n(\omega, \lambda) &= \frac{1}{n} \sum_{i=1}^n f((X_i, Y_i), (\omega, \lambda)) \quad \text{with} \\ Q_n(\hat{\omega}_n, \hat{\lambda}_n) &= \sup_{\omega \in \Omega^2, \lambda \in \Lambda^2} Q_n(\omega), \quad \text{and} \\ Q(\omega, \lambda) &= \mathbb{E}_{(X, Y)} (f((X, Y), (\omega, \lambda))) \\ &= \mathbb{E}_X (S(h(X, (\omega^*, \lambda^*))) \log(S(h(X, (\omega, \lambda)))) \\ &\quad + S(-h(X, (\omega^*, \lambda^*))) \log(S(-h(X, (\omega, \lambda))))). \end{aligned}$$

The last equation follows from conditioning on

$$Y \sim \text{Bernoulli}(S(h(X, (\omega^*, \lambda^*))))).$$

Before we move on, we need to prove that $f((X, Y), (\omega, \lambda))$ is integrable so that Q is well-defined. Without loss of generality assume that $\lambda_1 \geq \lambda_0$. It suffices to prove that $\mathbb{E}(f((X, Y), (\omega, \lambda)), Y = y)$ is integrable for both $y = 0, 1$. Observe that

$$\begin{aligned} h(X, (\omega, \lambda)) &\leq (\lambda_0 - \lambda_1) d_{\text{tr}}(X, \omega_0) + \lambda_1 d_{\text{tr}}(\omega_0, \omega_1) + \text{const} \\ &\leq \lambda_1 d_{\text{tr}}(\omega_0, \omega_1) + \text{const}. \end{aligned}$$

Since $h(X, (\omega, \lambda))$ is bounded above, $f((X, Y), (\omega, \lambda))$ is also bounded below on $Y = 0$ and is hence integral on $Y = 0$. Also, observe that

$$h(X, (\omega, \lambda)) \geq (\lambda_0 - \lambda_1) d_{\text{tr}}(X, \omega_1) - \lambda_0 d_{\text{tr}}(\omega_0, \omega_1) + \text{const}$$

and noting that $\log(S(x)) > x - 1$ for all $x < 0$

$$\log(S(h(X, (\omega, \lambda)))) \geq h(X, (\omega, \lambda)) - 1 \geq (\lambda_0 - \lambda_1) d_{\text{tr}}(X, \omega_1) + \text{const}.$$

Since $d_{\text{tr}}(X, \omega_1)$ is integrable on $Y = 1$, the LHS is integrable on $Y = 1$ too. It follows that $f(X, (\omega, \lambda))$ is integrable and hence Q is well-defined.

First, we prove that Q is maximised at $(\omega, \lambda) = (\omega^*, \lambda^*)$ and that this maximizer is unique. Consider the function

$$g : \mathbb{R} \rightarrow \mathbb{R}, \quad g(t) = S(\alpha) \log S(t) + S(-\alpha) \log S(-t),$$

where $\alpha \in \mathbb{R}$ is some constant. The function g is maximised at $t = \alpha$ and applying Taylor's theorem yields

$$g(x) = g(\alpha) - \frac{1}{2} S(\xi) S(-\xi) (x - \alpha)^2, \quad \text{for some } \xi \in (\alpha, x).$$

Setting $\alpha = h(X, (\omega^*, \lambda^*))$ and denoting ξ as a random variable

$$\xi(X) \in (h(X, (\omega^*, \lambda^*)), h(X, (\omega, \lambda)))$$

observe that

$$\begin{aligned} Q(\omega, \lambda) &= \mathbb{E}_X (g(h(X, (\omega, \lambda)))) \\ &= \mathbb{E}_X (g(h(X, (\omega^*, \lambda^*)))) - \\ &\quad \frac{1}{2} \mathbb{E}_X (S(\xi(X)) S(-\xi(X)) [h(X, (\omega, \lambda)) - h(X, (\omega^*, \lambda^*))]^2) \\ &\leq Q(\omega^*, \lambda^*), \end{aligned} \tag{35}$$

Hence, from the expression above it is deduced that (ω^*, λ^*) is a maximizer. Now, consider the function $q : \mathcal{X} \rightarrow \mathbb{R}$

$$q(x) = h(x, (\omega^*, \lambda^*)) - h(x, (\omega, \lambda)),$$

where $\Omega \subset \mathcal{X} \subset \mathbb{R}^e / \mathbb{R}\mathbf{1}$ such that for some $\zeta > 0$

$$\mathcal{X} = \{x \in \mathbb{R}^e / \mathbb{R}\mathbf{1} : \inf_{\omega \in \Omega} d_{\text{tr}}(x, \omega) < \zeta\},$$

so that for any $\omega \in \Omega$ there is a neighborhood of ω within \mathcal{X} . Note that \mathcal{X} is a bounded set since Ω is bounded too.

We will prove by contradiction that $q(x) = 0, \forall x \in \mathcal{X}$. Suppose there exists $x_0 \in \mathcal{X}$ such that $q(x_0) > 0$, then since q is continuous there exists a neighborhood U with $x_0 \in U$ such that $q(x) > 0$ for all $x \in U$ and so

$$\mathbb{E}(q^2(X)\mathbb{I}(X \in U)) > 0,$$

where \mathbb{I} is the indicator function. Since $h(x, (\omega, \lambda))$ is continuous with respect to x and \mathcal{X} is bounded, the function takes values on a bounded interval and hence $q(x)$ is bounded in \mathcal{X} i.e. there exists $\epsilon > 0$ such that $\mathbb{P}(S(\xi(X))S(-\xi(X)) > \epsilon | X \in U) = 1$ and so equation (35) becomes

$$Q(\omega, \lambda) \leq Q(\omega^*, \lambda^*) - \frac{\epsilon}{2} \mathbb{E}(q^2(X)\mathbb{I}(X \in U)) < Q(\omega^*, \lambda^*),$$

since $\mathbb{P}(X \in U) > 0$ (X has positive density everywhere). Therefore, for (ω, λ) to be a maximizer, $q(x) = 0$ for all $x \in \mathcal{X}$. Apply Lemma A.4 with $\omega^* = (\alpha, \beta)$, $\omega = (\gamma, \delta)$, $\lambda^* = (\lambda_\alpha, \lambda_\beta)$ and $\lambda = (\lambda_\gamma, \lambda_\delta)$ with the set \mathcal{X} containing neighbourhoods of $\alpha, \beta, \gamma, \delta$ and $q(x) = 0$ for all x in those neighbourhoods. It is concluded that $\omega = \omega^*$ and $\lambda = \lambda^*$, thus proving the uniqueness of the maximizer.

Theorem A.2 provides the uniform law of large numbers. The parameter space $\Omega^2 \times \Lambda^2$ is compact since Ω and Λ are compact. Moreover, $f((x, y), (\omega, \lambda))$ is clearly continuous at each $(\omega, \lambda) \in \Omega^2 \times \Lambda^2$. Finally, consider the function

$$r(z) = \sup_{\omega \in \Omega^2, \lambda \in \Lambda^2} \{|f(z, (\omega, \lambda))|\} = -f(z, \omega(z), \lambda(z)),$$

since f is non-positive. The function $\omega(z), \lambda(z)$ are chosen to minimize f . Using equation (35),

$$\begin{aligned} \mathbb{E}(r(X)) &\leq -Q(\omega^*, \lambda^*) + \\ &\quad \frac{1}{2} \mathbb{E}([h(X, (\omega(X), \lambda(X))) - h(X, (\omega^*, \lambda^*))]^2), \end{aligned}$$

since the sigmoid function is bounded by 1. Note that

$$\mathbb{E}((Z + W)^2) \leq 2(\mathbb{E}(Z^2) + \mathbb{E}(W^2)),$$

and set $W = \log(\lambda_1(X)/\lambda_0(X)) - \log(\lambda_1^*/\lambda_0^*)$. Since $\lambda_y(X) \in \Lambda \subseteq [a, b]$ for some $b \geq a > 0$, it follows that W^2 is integrable and so now we just have to prove that Z is integrable, where $Z = Z_1 + Z_2 + Z_3 + Z_4$ with the

four terms corresponding to tropical distance function $\lambda d_{\text{tr}}(X, \omega)$. It also holds

$$\mathbb{E}((Z_1 + Z_2 + Z_3 + Z_4)^2) \leq 2(\mathbb{E}(Z_1^2) + \mathbb{E}(Z_2^2) + \mathbb{E}(Z_3^2) + \mathbb{E}(Z_4^2))$$

and so $\mathbb{E}(Z^2)$ is bounded above by

$$\begin{aligned} & \mathbb{E} \left(\sum_{i=0}^1 \lambda_i^2 d_{\text{tr}}^2(X, \omega_i(X)) + (\lambda_i^*)^2 d_{\text{tr}}^2(X, \omega_i^*(X)) \right) \leq \\ & \mathbb{E}_Y \left[2 \left(\sum_{i=0}^1 \lambda_i^2 + (\lambda_i^*)^2 \right) \mathbb{E} (d_{\text{tr}}^2(X, \omega_Y^*) | Y) + \right. \\ & \left. 2 \left(\sum_{i=0}^1 \lambda_i^2 d_{\text{tr}}^2(\omega_i(X), \omega_Y^*) + (\lambda_i^*)^2 d_{\text{tr}}^2(\omega_i^*, \omega_Y^*) \right) \right], \end{aligned}$$

where the second inequality came from applying the triangular inequality four times in the form $d_{\text{tr}}(X, \tau) \leq d_{\text{tr}}(X, \omega_Y^*) + d_{\text{tr}}(\omega_Y^*, \tau)$. The final expression is finite because Ω is compact and hence $d_{\text{tr}}(\omega_i(X), \omega_Y^*)$ is finite, $d_{\text{tr}}(X, \omega_Y^*) | Y$ is square-integrable. Therefore, $\mathbb{E}(r(X))$ is finite.

All conditions of the theorem are satisfied and so

$$\begin{aligned} & \sup_{\omega \in \Omega^2} \left| \frac{1}{n} \sum_{i=1}^n f((X_i, Y_i), \omega) - \mathbb{E}(f((X, Y), \omega)) \right| = \\ & \sup_{\omega \in \Omega^2} |Q_n(\omega) - Q(\omega)| \xrightarrow{P} 0. \end{aligned}$$

Finally, using Theorem A.1 and combining the uniqueness of the maximizer with the uniform bound result, it is concluded that $\hat{\omega} \xrightarrow{P} \omega^*$. \square

Proof of Proposition 5. First, define $\Delta_0 = \{C(X) \neq 1 | Y = 0\}$. By definition of $C(X)$,

$$\begin{aligned} \Delta_0 &= \left\{ (\sigma_0^{-1} - \sigma_1^{-1}) d_{\text{tr}}(X, \hat{\omega}) - (e - 1) \log \left(\frac{\sigma_1}{\sigma_0} \right) \geq 0 \mid Y = 0 \right\} \\ &= \{d_{\text{tr}}(X, \hat{\omega}) \geq \alpha \sigma_0 \sigma_1 \mid Y = 0\}. \end{aligned}$$

Triangular inequality dictates that

$$d_{\text{tr}}(X, \omega^*) - d_{\text{tr}}(\omega^*, \hat{\omega}) \leq d_{\text{tr}}(X, \hat{\omega}) \leq d_{\text{tr}}(X, \omega^*) + d_{\text{tr}}(\omega^*, \hat{\omega}),$$

and so it follows that

$$\begin{aligned} \Delta_0 &\supseteq \{d_{\text{tr}}(X, \omega^*) \geq \sigma_0 \sigma_1 (\alpha + \epsilon) \mid Y = 0\} \\ \Delta_0 &\subseteq \{d_{\text{tr}}(X, \omega^*) \geq \sigma_0 \sigma_1 (\alpha - \epsilon) \mid Y = 0\}, \end{aligned}$$

and since $Z = \sigma_0^{-1} d_{\text{tr}}(X, \omega^*) | Y = 0 \sim F$,

$$\mathbb{P}(Z \geq \sigma_1(\alpha + \epsilon)) \leq \mathbb{P}(\Delta_0) \leq \mathbb{P}(Z \geq \sigma_1(\alpha - \epsilon)),$$

which yields the desired result.

Similarly, for $\Delta_1 = \{C(X) \neq 0 | Y = 1\} = \{d_{\text{tr}}(X, \hat{\omega}) \leq \sigma_0 \sigma_1 \alpha\}$,

$$\begin{aligned} \Delta_1 &\supseteq \{d_{\text{tr}}(X, \omega^*) \leq \sigma_0 \sigma_1 (\alpha - \epsilon) \mid Y = 1\} \\ \Delta_1 &\subseteq \{d_{\text{tr}}(X, \omega^*) \leq \sigma_0 \sigma_1 (\alpha + \epsilon) \mid Y = 1\}, \end{aligned}$$

and since $Z = \sigma_1^{-1} d_{\text{tr}}(X, \omega^*) | Y = 1 \sim F$,

$$\mathbb{P}(Z \leq \sigma_0(\alpha - \epsilon)) \leq \mathbb{P}(\Delta_1) \leq \mathbb{P}(Z \leq \sigma_0(\alpha + \epsilon)),$$

which is the desired interval.

For the second part of the proposition, $\hat{\omega} = \omega^*$ and so $\epsilon = 0$. Hence,

$$\begin{aligned} \mathbb{P}(\Delta_0) &= 1 - F(\sigma_1 \alpha) = 1 - F(xu(x)) \\ \mathbb{P}(\Delta_1) &= F(\sigma_0 \alpha) = F(u(x)), \text{ where} \\ x &= \frac{\sigma_1}{\sigma_0} \text{ and } u(x) = (e - 1) \frac{\log x}{x - 1} \end{aligned}$$

Consider the function

$$g(x) = 1 - F(xu(x)) - F(u(x))$$

Proving that $g(x) < 0$ for all $x > 1$ is equivalent to proving the desired result that $\mathbb{P}(\Delta_0) < \mathbb{P}(\Delta_1)$ for $\sigma_1 > \sigma_0$. First,

$$\lim_{x \rightarrow 1} u(x) = \lim_{x \rightarrow 1} xu(x) = e - 1,$$

and so $\lim_{x \rightarrow 1} g(x) = 1 - 2F(e - 1)$. It is a well-known fact that the median of the Gamma distribution is less than the mean. Hence, for $Z \sim \text{Gamma}(e - 1, 1)$ with mean $e - 1$, $F(e - 1) > \frac{1}{2}$ and so

$$\lim_{x \rightarrow 1} g(x) < 0. \quad (36)$$

Finally, the derivative of g is

$$g'(x) = -F'(u(x))u'(x) - F'(xu(x))(xu'(x) + u(x))$$

The following two inequalities

$$F'(u(x)) \geq F'(xu(x)), \quad (37)$$

$$u'(x) + xu'(x) + u(x) \geq 0, \quad (38)$$

imply that

$$g'(x) \leq -F'(xu(x))(u'(x) + xu'(x) + u(x)) \leq 0. \quad (39)$$

From (36) and (39) it follows that $g(x) < 0$ for all $x > 1$.

For inequality (37), remember that

$$F'(x) = \frac{x^{e-2} \exp(-x)}{\Gamma(e-1)}$$

and so

$$\begin{aligned} F'(u(x)) - F'(xu(x)) &= F'(u(x)) (1 - x^{e-2} \exp(-(x-1)u(x))) \\ &= F'(u(x)) (1 - x^{e-2} \exp(-(e-1)\log(x))) \\ &= F'(u(x))(1 - x^{-1}) > 0, \end{aligned}$$

for all $x > 1$.

For inequality (38),

$$u'(x) + xu'(x) + u(x) = \frac{e-1}{(x-1)^2} (x - x^{-1} - 2 \log x),$$

is a non-negative function for $x > 1$ iff v is a non-negative function, where

$$\begin{aligned} v(x) &= x - x^{-1} - 2 \log x, \text{ with} \\ v'(x) &= \frac{(x-1)^2}{x^2} \geq 0 \text{ and } v(1) = 0. \end{aligned}$$

Clearly, v is a non-negative function for $x > 1$, so inequality (38) is satisfied. \square

Proof of Proposition 6. For symbolic convenience, in this proof class 0 is referred to as class -1 and so $Y \in \{-1, 1\}$. Applying the triangular inequality twice,

$$\begin{aligned} D_X &= d_{\text{tr}}(X, \omega_Y^*) - d_{\text{tr}}(X, \omega_{-Y}^*) \\ &\geq (d_{\text{tr}}(X, \hat{\omega}_Y) - d_{\text{tr}}(\omega_Y^*, \hat{\omega}_Y)) \\ &\quad - (d_{\text{tr}}(X, \hat{\omega}_{-Y}) + d_{\text{tr}}(\omega_{-Y}^*, \hat{\omega}_{-Y})) \\ &= d_{\text{tr}}(X, \hat{\omega}_Y) - d_{\text{tr}}(X, \hat{\omega}_{-Y}) - \epsilon, \end{aligned}$$

it follows that

$$\{C(X) \neq Y\} = \{d_{\text{tr}}(X, \hat{\omega}_Y) - d_{\text{tr}}(X, \hat{\omega}_{-Y}) \geq 0\} \subseteq \{D_X \geq -\epsilon\}$$

and so the generalization error has the following upper bound

$$\mathbb{P}(C(X) \neq Y) \leq \mathbb{P}(D_X \geq -\epsilon).$$

Note that if $d_{\text{tr}}(X, \omega_Y^*) < \Delta_\epsilon$, then by the use of triangular inequality

$$\begin{aligned} D_X &= d_{\text{tr}}(X, \omega_Y^*) - d_{\text{tr}}(\omega_{-Y}^*, X) \\ &\leq d_{\text{tr}}(X, \omega_Y^*) - (d_{\text{tr}}(\omega_{-Y}^*, \omega_Y^*) - d_{\text{tr}}(\omega_Y^*, X)) \\ &< 2\Delta_\epsilon - d_{\text{tr}}(\omega_1^*, \omega_{-1}^*) = -\epsilon. \end{aligned}$$

Consequently,

$$\mathbb{P}(C(X) \neq Y) \leq \mathbb{P}(D_X \geq -\epsilon, Z_X \geq \Delta_\epsilon) \quad (40)$$

Since the distribution of X is symmetric around ω_Y^* , the random variable $2\omega_Y^* - X$ has the same distribution and so

$$\mathbb{P}(D_X \geq -\epsilon, Z_X \geq \Delta_\epsilon) = \mathbb{P}\left(D_{2\omega_Y^* - X} \geq -\epsilon, Z_{2\omega_Y^* - X} \geq \Delta_\epsilon\right). \quad (41)$$

It will be proved that

$$Z_{2\omega_Y^* - X} = Z_X, \quad (42)$$

$$D_X + D_{2\omega_Y^* - X} \leq 0, \quad (43)$$

and so $\{D_{2\omega_Y^* - X} \geq -\epsilon, Z_{2\omega_Y^* - X} \geq \Delta_\epsilon\} \subseteq \{D_X \leq \epsilon, Z_X \geq \Delta_\epsilon\}$. Then, using equation (41),

$$\mathbb{P}(D_X \geq -\epsilon, Z_X \geq \Delta_\epsilon) \leq \mathbb{P}(D_X \leq \epsilon, Z_X \geq \Delta_\epsilon),$$

and substituting it to inequality (40),

$$\begin{aligned} \mathbb{P}(C(X) \neq Y) &= \frac{1}{2}(\mathbb{P}(D_X \geq -\epsilon, Z_X \geq \Delta_\epsilon) \\ &\quad + \mathbb{P}(D_X \leq \epsilon, Z_X \geq \Delta_\epsilon)) \\ &= \mathbb{P}(Z_X \geq \Delta_\epsilon) + h(\epsilon) \end{aligned}$$

where $h(\epsilon) = \mathbb{P}(Z_X \geq \Delta_\epsilon, |D_X| \leq \epsilon)$ is an increasing function with respect to ϵ , which completes the proof.

Equation (42) follows from the observation that

$$d_{\text{tr}}(2\omega_Y^* - x, \omega_Y^*) = d_{\text{tr}}(x, \omega_Y^*).$$

For equation (43),

$$\begin{aligned} D_{2\omega_Y^* - X} + D_X &= Z_{2\omega_Y^* - X} - d_{\text{tr}}(2\omega_Y^* - X, \omega_{-Y}^*) \\ &\quad + Z_X - d_{\text{tr}}(X, \omega_{-Y}^*) \\ &\stackrel{(42)}{=} 2Z_{2\omega_Y^* - X} - d_{\text{tr}}(2\omega_Y^* - X, \omega_{-Y}^*) - d_{\text{tr}}(\omega_{-Y}^*, X) \\ &\leq 2Z_{2\omega_Y^* - X} - d_{\text{tr}}(2\omega_Y^* - X, X) = 0, \end{aligned}$$

where the last inequality comes from the triangular inequality. \square

Proof of Proposition 7. Consider the random variable $d_{\text{tr}}(X, \alpha)$. From the triangular inequality

$$d_{\text{tr}}(X, \alpha) \leq d_{\text{tr}}(X, \omega^*) + d_{\text{tr}}(\alpha, \omega^*),$$

it is deduced that $d_{\text{tr}}(X, \alpha)$ is integrable, bounded above by an integrable random variable.

Now consider the function $F : \mathbb{R}^e / \mathbb{R}\mathbf{1} \rightarrow \mathbb{R}$,

$$F(x) = d_{\text{tr}}(x, \omega) + d_{\text{tr}}(2\omega^* - x, \omega) - 2d_{\text{tr}}(x, \omega^*).$$

Noting that $d_{\text{tr}}(2\omega^* - x, \omega) = d_{\text{tr}}(x, 2\omega^* - \omega)$, it follows that $F(X)$ is integrable as the sum of integrable random variables.

From triangular inequality and the fact that $d_{\text{tr}}(2\omega^* - x, x) = 2d_{\text{tr}}(x, \omega^*)$ it follows that $F(x) \geq 0$ for all $x \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$. Furthermore, $F(\omega^*) > 0$ and since F is continuous, there exists a neighbourhood U that contains ω^* such that $F(x) > 0$ for all $x \in U$. Moreover, the function has positive density in a neighbourhood V that contains the centre ω^* . Therefore, there exists a neighbourhood $W = U \cap V$ such that $F(x) > 0$ for all $x \in W$ and $\mathbb{P}(X \in W) > 0$. Hence, since $F(X) \geq 0$,

$$\mathbb{E}(F(X)) \geq \mathbb{E}(F(X)|X \in W)\mathbb{P}(X \in W) > 0.$$

In other words,

$$\mathbb{E}(d_{\text{tr}}(X, \omega)) + \mathbb{E}(d_{\text{tr}}(2\omega^* - X, \omega)) > 2\mathbb{E}(d_{\text{tr}}(x, \omega^*)) \quad (44)$$

Moreover, consider the isometry $y = 2\omega^* - x$ and note that for symmetric probability density functions around ω^* , $f(\omega^* - \delta) = f(\omega^* + \delta)$ and so for $\delta = \omega^* - x$, we have $f(y) = f(x)$. Applying this transformation to the following integral yields

$$\begin{aligned}\mathbb{E}(d_{\text{tr}}(2\omega^* - X, \omega)) &= \int_{\mathbb{R}^e/\mathbb{R}\mathbf{1}} d_{\text{tr}}(2\omega^* - x, \omega) f(x) dx \\ &= \int_{\mathbb{R}^e/\mathbb{R}\mathbf{1}} d_{\text{tr}}(y, \omega) f(y) dy = \mathbb{E}(d_{\text{tr}}(X, \omega)).\end{aligned}\quad (45)$$

Combining equation (45) with inequality (44) shows that the function $Q(\omega) = \mathbb{E}(d_{\text{tr}}(X, \omega))$ has a global minimum at ω^* .

From Theorem A.2 (uniform law of large numbers), set $f(x, \omega) = d_{\text{tr}}(x, \omega)$ and observe that $f(x, \omega)$ is always continuous w.r.t. ω . Setting $r(x) = \sup_{\omega \in \Omega} d_{\text{tr}}(x, \omega)$, which is finite since Ω is compact, observe that

$$r(x) := \sup_{\omega \in \Omega} d_{\text{tr}}(x, \omega) \leq d_{\text{tr}}(x, \omega^*) + \sup_{\omega \in \Omega} d_{\text{tr}}(\omega, \omega^*).$$

Since Ω is compact, the second term is finite and hence $r(X)$ is integrable, since $d_{\text{tr}}(X, \omega^*)$ is integrable. All conditions of the theorem are satisfied so $Q(\omega) = \mathbb{E}(d_{\text{tr}}(x, \omega))$ is continuous with respect to ω and

$$\sup_{\omega \in \Omega} |Q_n(\omega) - Q(\omega)| \xrightarrow{P} 0 \text{ as } n \rightarrow \infty,$$

where $Q_n(\omega) = n^{-1} \sum_{i=1}^n d_{\text{tr}}(X_i, \omega)$. Since $Q(\omega)$ has a unique minimum at ω^* , all conditions of Theorem A.1 are satisfied and so $\tilde{\omega}_n \rightarrow \omega^*$ as $n \rightarrow \infty$. \square

Proof of Proposition 8. i. If $\omega - X_i$ has a unique maximum $M_i = \arg \max_j \{\omega_j - (X_i)_j\}$ and unique minimum $m_i = \arg \min_j \{\omega_j - (X_i)_j\}$, then the gradient is

$$(\nabla f(x))_j = |\{i : M_i = j\}| - |\{i : m_i = j\}|. \quad (46)$$

For the converse, assume that the gradient is well-defined. From equations (23)–(24) and following the first few sentences of Lemma A.4

$$d_{\text{tr}}(x + \epsilon E_j, y) + d_{\text{tr}}(x - \epsilon E_j, y) - 2d_{\text{tr}}(x, y) = \epsilon s_j(x - y),$$

where s_j is defined in equation (26) of Lemma A.4. Consequently,

$$f(x + \epsilon E_j) + f(x - \epsilon E_j) - 2f(x) = \epsilon \sum_{i=1}^n s_j(X_i - \omega_i)$$

Since f has a well-defined gradient, $\sum_{i=1}^n s_j(X_i - \omega) = 0$ i.e. $s_j(X_i - \omega) = 0$ for all $(i, j) \in [n] \times [e]$. This can only happen iff $X_i - \omega$ has unique maximum and minimum component for all $i \in [n]$.

- ii. Using equation (46), the gradient of f vanishes at $x = \omega$ if and only if

$$|\{i : M_i = j\}| = |\{i : m_i = j\}|. \quad (47)$$

Moreover,

$$\begin{aligned} f(\omega + v) &= \sum_{i=1}^n \max_k \{\omega_k - (X_i)_k + v_k\} - \min_k \{\omega_k - (X_i)_k + v_k\} \\ &\geq \sum_{i=1}^n \omega_{M_i} - (X_i)_{M_i} + v_{M_i} - \omega_{m_i} + (X_i)_{m_i} - v_{m_i} \\ &= f(\omega) + \sum_{i=1}^n v_{M_i} - v_{m_i} \end{aligned}$$

Finally, note that because of equation (47),

$$\begin{aligned} \sum_{i=1}^n v_{M_i} &= \sum_{j=1}^e v_j |\{i \in [n] : M_i = j\}| \\ &\stackrel{(47)}{=} \sum_{j=1}^e v_j |\{i \in [n] : m_i = j\}| = \sum_{i=1}^n v_{m_i}, \end{aligned}$$

and so $f(\omega + v) \geq f(\omega)$ for all $v \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$. □

B Space of ultrametrics

Theorem B.1 (explained in [3, 18]). *Suppose we have a classical linear subspace $L_m \subset \mathbb{R}^e$ defined by the linear equations $x_{ij} - x_{ik} + x_{jk} = 0$ for $1 \leq i < j < k \leq m$. Let $\text{Trop}(L_m) \subseteq \mathbb{R}^e / \mathbb{R}\mathbf{1}$ be the tropicalization of the linear space $L_m \subset \mathbb{R}^e$, that is, classical operators are replaced by tropical ones (defined in Section C in the supplement) in the equations defining the linear subspace L_m , so that all points $(v_{12}, v_{13}, \dots, v_{m-1,m})$ in $\text{Trop}(L_m)$ satisfy the condition that*

$$\max_{i,j,k \in [m]} \{v_{ij}, v_{ik}, v_{jk}\}.$$

is attained at least twice. Then, the image of \mathcal{U}_m inside of the tropical projective torus $\mathbb{R}^e / \mathbb{R}\mathbf{1}$ is equal to $\text{Trop}(L_m)$.

C Tropical Arithmetics and Tropical Inner Product

In tropical geometry, addition and multiplication are different than regular arithmetic. The arithmetic operations are performed in the max-plus tropical semiring $(\mathbb{R} \cup \{-\infty\}, \oplus, \odot)$ as defined in [19].

Definition C.1 (Tropical Arithmetic Operations). *In the tropical semiring, the basic tropical arithmetic operations of addition and multiplication are defined as:*

$$a \oplus b := \max\{a, b\}, \quad a \odot b := a + b, \quad \text{where } a, b \in \mathbb{R} \cup \{-\infty\}.$$

The element $-\infty$ ought to be included as it is the identity element of tropical addition. Tropical subtraction is not well-defined and tropical division is classical subtraction.

The following definitions are necessary for the definition of the tropical inner product

Definition C.2 (Tropical Scalar Multiplication and Vector Addition). *For any scalars $a, b \in \mathbb{R} \cup \{-\infty\}$ and for any vectors $v, w \in (\mathbb{R} \cup \{-\infty\})^e$, where $e \in \mathbb{N}$,*

$$\begin{aligned} a \odot v &:= (a + v_1, \dots, a + v_e), \\ a \odot v \oplus b \odot w &:= (\max\{a + v_1, b + w_1\}, \dots, \max\{a + v_e, b + w_e\}). \end{aligned}$$

From the definitions above, it follows that the tropical inner product is $\omega^T \odot x = \max\{\omega + x\}$ for all vectors $\omega, x \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$. In classical logistic regression a linear function in the form of a classical inner product $h_\omega(x) = \omega^T x$, $\omega \in \mathbb{R}^n$ is used. The tropical symbolic equivalent is

$$h_\omega(x) = \omega^T \odot x = \max_{l \in [e]} \{\omega_l + x_l\}. \quad (48)$$

This expression is not well-defined, since the statistical parameter and covariate vectors $\omega, u \in \mathbb{R}^e / \mathbb{R}\mathbf{1}$ are only defined up to addition of a scalar multiple of the vector $(1, \dots, 1)$. To resolve this issue, we fix

$$-\min_{l \in [e]} \{\omega_l + x_l\} = c, \quad (49)$$

where $c \in \mathbb{R}$ is a constant for all observations. Combining equations (49), (48), and the definition of tropical distance (1),

$$h_\omega(x) = d_{\text{tr}}(x, -\omega) - c.$$

For simplicity, under the transformation $-\omega \rightarrow \omega$ the expression becomes

$$h_\omega(x) = d_{\text{tr}}(x, \omega) - c. \quad (50)$$

D Tropical Logistic Regression Algorithm

Algorithm 1: One-species tropical logistic regression

Input: distance matrix $D \in \mathbb{R}_+^{N \times e}$, labels $Y \in \{0, 1\}^N$
 $\tilde{\omega} = \text{FW_point}(D)$
 $\hat{\sigma}_0, \hat{\sigma}_1 = \arg \max_{\sigma_0, \sigma_1 > 0} l(\tilde{\omega}, \sigma_0, \sigma_1 | D, Y)$ with root solving.
Output: $(\tilde{\omega}, \hat{\sigma}_0, \hat{\sigma}_1)$

Algorithm 2: Two-species tropical logistic regression

Input: distance matrix $D \in \mathbb{R}_+^{N \times e}$, labels $Y \in \{0, 1\}^N$
 $\tilde{\omega}_0 = \text{FW_point}(D[Y == 0])$
 $\tilde{\omega}_1 = \text{FW_point}(D[Y == 1])$
 $\hat{\sigma} = \arg \max_{\sigma > 0} l(\tilde{\omega}_0, \tilde{\omega}_1, \sigma | D, Y)$ with root solving.
Output: $(\tilde{\omega}_0, \tilde{\omega}_1, \hat{\sigma})$

E Fermat-Weber Point Visualization

As noted in Section 5, the gradient method is much faster than linear programming. Unfortunately, there is no guarantee that it will guide us to a Fermat-Weber point. However, in practice, the gradient method tends to work well. Figure 7 illustrates just that. Given, ten datapoint $X_1, \dots, X_{10} \in \mathbb{R}^3 / \mathbb{R}\mathbf{1} \cong \mathbb{R}^2$, the Fermat-Weber set is found to be a trapezoid. This is in agreement with [12], which states that all Fermat-Weber sets are classical polytopes. The two-dimensional gradient vector, plotted as a vector field in Figure 7, always points towards the Fermat-Weber set. Therefore, the gradient algorithm should always guide us to a Fermat-Weber point.

F MLE Estimator for σ

If $Z_i \stackrel{\text{iid}}{\sim} \text{Gamma}(n, k)$, where n is constant and k is a statistical parameter, then it is well-known that the maximum likelihood estimator is

$$\hat{k} = \bar{Z}/n, \quad (51)$$

where \bar{Z} is the sample average. In our case $Z_i = d(X_i, \omega^*)$ and $k = i\sigma^i$. From Proposition 4, $Z_i \sim \text{Gamma}(n/i, i\sigma^i)$ and by substituting these parameters in equation 51, it follows that the MLE for σ is

$$\hat{\sigma}^i = \bar{Z}/n,$$

where \bar{Z} is the average distance of the covariates (gene trees) from their mean (species tree). This results holds for all $i \in \mathbb{N}$ and both Euclidean and tropical metrics. The only difference is that for Euclidean spaces $X \in \mathbb{R}^e$ and so $n = e$, while for the tropical projective torus $\mathbb{R}^e / \mathbb{R}\mathbf{1}$, $n = e - 1$.

G Approximate BHV Logistic Regression

Similar to the tropical Laplace distribution, in [5] the following distribution was considered

$$f_{\lambda, \omega}(x) = K_{\lambda, \omega} \exp(-\lambda d_{\text{BHV}}(x, \omega)),$$

where $\lambda = 1/\sigma$ is a concentration/precision parameter, d_{BHV} is the BHV metric and $K_{\lambda, \omega}$ is the normalization constant that depends on λ and

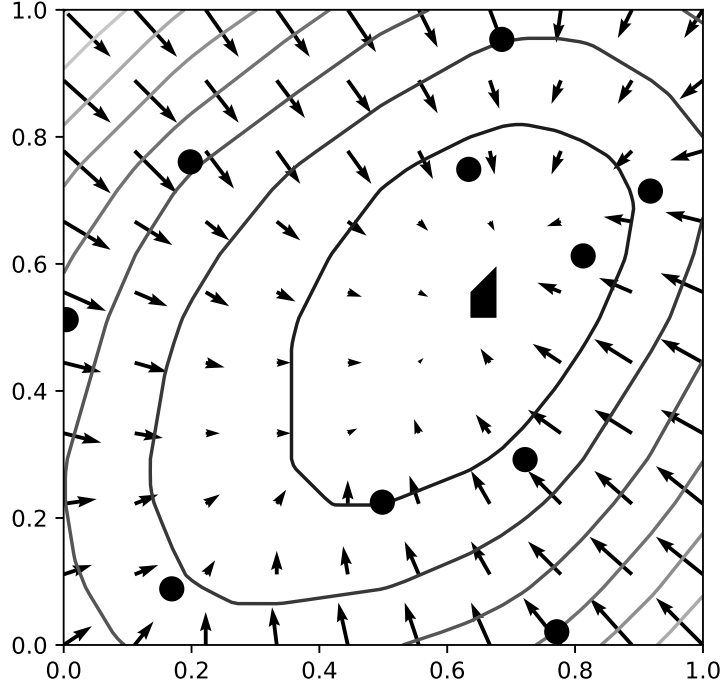


Figure 7: Visualization of the function $f(\omega) = \sum_{i=1}^{10} d_{\text{tr}}(X_i, \omega)$ for X_i . The black circles are the datapoints X_1, \dots, X_{10} , the solid lines are contours of f , the vector field is the gradient and the small black trapezoid at $(0.65, 0.55)$ is the Fermat-Weber set.

ω . We consider an adaptation of the two-species model for this metric, where the data from the two classes have the same concentration rate but different centre. If $X|Y \sim f_{\lambda, \omega_Y^*}$, then

$$h_{\omega_0, \omega_1}(x) = \lambda (d_{\text{BHV}}(x, \omega_0^*) - d_{\text{BHV}}(x, \omega_1^*)) + \log \frac{K_{\lambda, \omega_0^*}}{K_{\lambda, \omega_1^*}}. \quad (52)$$

Unlike in the tropical projective torus or the euclidean space, in the BHV space $K_{\lambda, \omega_0^*} \neq K_{\lambda, \omega_1^*}$, because the space is not translation-invariant. However, if we assume that the two centres are far away from trees with bordering topologies, it may be assumed that the trees are mostly distributed in the Euclidean space and as a result $K_{\lambda, \omega_0^*} \approx K_{\lambda, \omega_1^*}$. Under this assumption, equation (52) becomes

$$h_{\omega_0, \omega_1}(x) \approx \lambda (d_{\text{BHV}}(x, \omega_0^*) - d_{\text{BHV}}(x, \omega_1^*)).$$

Therefore, the classification/decision boundary for the BHV is the BHV bisector $d_{\text{BHV}}(x, \omega_0^*) = d_{\text{BHV}}(x, \omega_1^*)$ and the most sensible classifier is

$$C(x) = \mathbb{I}(d_{\text{BHV}}(x, \omega_0^*) > d_{\text{BHV}}(x, \omega_1^*)),$$

where \mathbb{I} is the indicator function.

H Graphs for Simulated Data under the Multi-Species Coalescent Model for different R

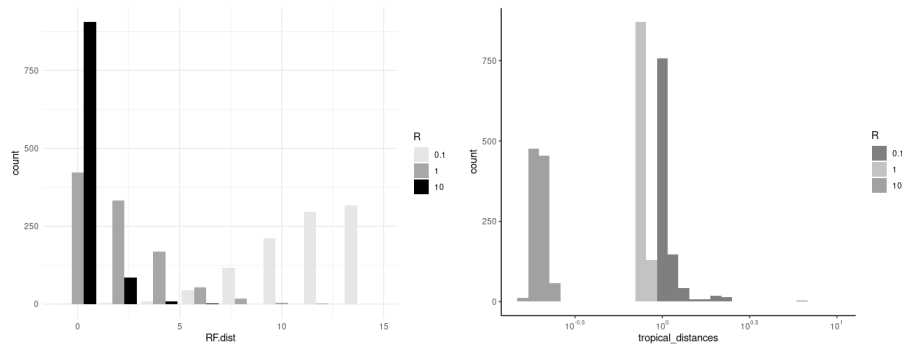


Figure 8: (LEFT) Robinson-Foulds distances and (RIGHT) tropical distances of inferred species trees $\hat{\omega}$ from the actual species trees ω^* for $R = 0.1, 1, 10$.

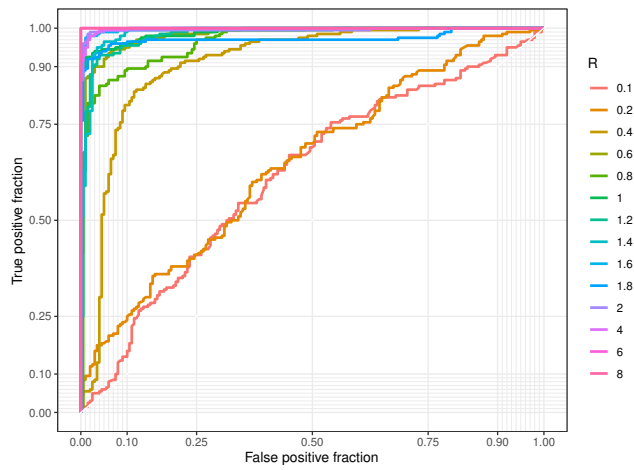


Figure 9: ROC curves for the tropical logistic regression with different values of R . Higher the value of R is the closer an estimated ROC curve for the tropical logistic regression model gets to the point $(0, 1)$.