

# Gaussian tree constraints applied to acoustic linguistic functional data



Nathaniel Shiers<sup>a</sup>, John A.D. Aston<sup>b,\*</sup>, Jim Q. Smith<sup>a</sup>, John S. Coleman<sup>c</sup>

<sup>a</sup> Department of Statistics, University of Warwick, Coventry, UK

<sup>b</sup> Statistical Laboratory, University of Cambridge, Cambridge, UK

<sup>c</sup> Phonetics Laboratory, University of Oxford, Oxford, UK

## ARTICLE INFO

### Article history:

Received 26 September 2015

Available online 11 October 2016

### AMS subject classifications:

primary 62P99

secondary 62G05

13P25

### Keywords:

Gaussian tree constraints

Canonical variate analysis

Functional data analysis

Phonetics

Separable covariance

## ABSTRACT

Evolutionary models of languages are usually considered to take the form of trees. With the development of so-called tree constraints the plausibility of the tree model assumptions can be assessed by checking whether the moments of observed variables lie within regions consistent with Gaussian latent tree models. In our linguistic application, the data set comprises acoustic samples (audio recordings) from speakers of five Romance languages or dialects. The aim is to assess these functional data for compatibility with a hereditary tree model at the language level. A novel combination of canonical function analysis (CFA) with a separable covariance structure produces a representative basis for the data. The separable-CFA basis is formed of components which emphasize language differences whilst maintaining the integrity of the observational language-groupings. A previously unexploited Gaussian tree constraint is then applied to component-by-component projections of the data to investigate adherence to an evolutionary tree. The results highlight some aspects of Romance language speech that appear compatible with an evolutionary tree model but indicate that it would be inappropriate to model all features as such.

© 2016 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Functional data analysis (FDA) is the statistical analysis of intrinsically infinite dimensional objects which can often be described in terms of curves or (hyper-)surfaces. The range of applications (e.g., brain imaging [69], climatology [7], and medical research [61]) for FDA is considerable, and this has resulted in considerable theoretical developments as well [35,60]. The use of FDA in statistical phonetics has recently attracted attention (e.g., [40,51]) as many phonological processes can essentially be seen as continuous, whether that be sound waves or derived objects such as spectrograms, where, in the tradition of FDA, the discretization of the signal is more properly considered as being inherently arbitrary and a result of the measurement rather than a property of the object under consideration. Analyses which involve acoustic functional data have provided particularly promising and interesting results in a diverse range of settings. For example, Grabe et al. [29] use a polynomial basis expansion to examine pitch variation in English, while Aston et al. [4] investigate Qiang, a Sino-Tibetan language, and find previously unidentified gender differences amongst speakers via a functional principal components based modeling approach.

\* Corresponding author.

E-mail addresses: [n.l.shiers@warwick.ac.uk](mailto:n.l.shiers@warwick.ac.uk) (N. Shiers), [j.aston@statslab.cam.ac.uk](mailto:j.aston@statslab.cam.ac.uk) (J.A.D. Aston), [j.q.smith@warwick.ac.uk](mailto:j.q.smith@warwick.ac.uk) (J.Q. Smith), [john.coleman@phon.ox.ac.uk](mailto:john.coleman@phon.ox.ac.uk) (J.S. Coleman).

<http://dx.doi.org/10.1016/j.jmva.2016.09.015>

0047-259X/© 2016 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

The acoustic structure of spoken words can be used to investigate areas of linguistic interest in a similar way that discrete (alphabetic) representations of speech have been utilized to make cross-language comparisons, and more recently inferences regarding proto languages [11,33]. The differences and similarities between spoken languages suggest that any meaningful functional observations taken across languages are unlikely to be independently, identically distributed. As such, it is probable that the language relationships form a tree or network structure, which may be informative about possible historical developments of these languages. If this alternative (acoustic) approach can be used to corroborate known and uncontroversial language relationships, then our methods offer great potential for less certain language relationships. For instance, this would be useful for languages where there are few historical records but in which inference of a family tree is reasonably supported by the contemporary data (e.g., African language families), or alternatively, in cases where reconstruction of a family tree is disputed, such as Greenberg's classification of Native American languages [10].

Relationships between languages have long been described as phylogenetic trees constructed using linguistic factors (e.g., Schleicher [64]) where all non-leaf variables are unobserved and represent features of the past languages before their divergence. Greenberg [30] developed some of the first quantitative methods which were used to investigate evolutionary relationships between languages. More recently there have been large scale attempts to reconstruct trees or networks of languages (e.g., Nakhleh et al. [52] for the Indo-European language family, Nicholls and Ryder [54] for the Semitic language family). Some researchers have shifted away from describing the evolutionary language relationships via trees toward using networks (for example, [26,53]). However, trees have a somewhat more natural interpretation in terms of evolutionary structure, and assessing the suitability of a tree model for language data would therefore be of interest to researchers in linguistics. The focus of this paper is to examine functional acoustic data from speakers of five Romance languages (French, Italian, Portuguese, American Spanish, and Iberian Spanish) to provide insight at an exploratory level as to whether a tree may be adequate for describing certain features of these language relationships. We will examine this through a specific acoustic representation (the spectrogram), which captures considerable acoustic information for each word.

To address questions of whether data is compatible with a latent tree model (tree-amenability), we appeal to the notion of tree constraints. The theory of tree constraints is embedded in the area of algebraic statistics, a field that has a significant recent literature related to phylogenetics (e.g., [2,70]). It has been known for some time that covariance functions of data on observed variables respecting an evolutionary tree must obey particular algebraic and semi-algebraic constraints, e.g., [65]. Recently these have become much better understood (for example [1,3,19]) and fully characterized in some cases (e.g., the binary case [75,76]). In this paper, building on developments in Shiers and Smith [67], a Gaussian analogue of the binary tree constraints is applied to the covariances of component-by-component projections of the data. By considering the data component-wise, a more realistic and nuanced analysis can be performed which permits some observed features of linguistic data to be tree-amenable and others not.

The overall aim of this paper is to present a methodology for assessing whether particular features of spoken languages may be suitably modeled by Gaussian trees with latent interior variables. To this end, it is necessary to identify phonetic features that effectively distinguish languages. This is achieved by projecting high dimensional spectrograms to a novel separable-canonical variate basis, to obtain a meaningful low dimensional representation of data. The features highlighted by this projection can then be assessed for compatibility with evolutionary trees. Throughout, a Romance language data set is used to illustrate the methodology as a proof of concept for phonetic data analysis.

In more detail, Section 2 describes the data set and the preprocessing in preparing an audio recording for FDA. In a similar spirit to the work on object data by Wang and Marron [72], for the application to Romance languages presented here, the observation units of interest are two-dimensional functional data objects known as spectrograms (time–frequency descriptions of the data). These are formed from transformations of audio recordings of people speaking single words. When regarded as functional data, observations are in fact stored as high dimensional objects. Therefore, in Section 3 tools from FDA are employed to transform high dimensional speaker data to a lower dimension. This is achieved through the novel approach of using between-language covariance as described in canonical function analysis (CFA) and combining it with a tensor decomposable covariance structure.

Having achieved the required reduction in dimension, Section 4 provides a brief summary of tree constraints. A fundamental but yet to be exploited constraint for use with Gaussian data is then introduced. Statistics associated with the violation of or adherence to the constraint are then constructed from the acoustic language data to answer the question of tree suitability. In Section 4.3, in preparation for use with this Gaussian tree constraint, we describe the construction of a between-language cross-covariance matrix using the scores (projections) of the acoustic data. In Section 4.4, the general effectiveness of the Gaussian tree constraint is investigated via simulations to assess its ability to correctly accept or reject tree-amenability. Section 4.5 entails further simulation, tailored to mimic aspects of the acoustic data, so the results are more immediately relevant to the setting. Section 5 addresses whether any of the isolated features of the spoken languages can be described by a Gaussian evolutionary tree model. These tree constraints are then used to explore tree-amenability for each of the components of the chosen basis. It is found that a subset of the components adhere to the tree constraint. This suggests that some features of the acoustic linguistic data which distinguish between languages could have evolved in a tree-like manner whilst others have not. This preliminary result is consistent with the current understanding that the development of Romance languages has been complex, involving much cross-language interaction [32] in addition to a historically well-documented common origin from Latin. Thus, attempting to fit a tree model to the entire data set would be misguided based on the empirical evidence presented here. More appropriately, a different model can be fitted for each component, using the tree constraints to indicate whether to restrict the space of models to trees.

## 2. Acoustic functional data set

The data set of interest comprises audio recordings originating from speakers of one of five different Romance languages: French, Italian, Portuguese, Spanish (American), and Spanish (Iberian)—while two dialects of Spanish are being used in this study, they are treated as different spoken languages in this analysis as the interest is in pronunciation rather than textual representation, the difference between ‘dialect’ and ‘language’ being a matter of degree of difference rather than an absolute qualitative difference. Each recording is of some individual saying an integer from ‘one’ to ‘ten’ in their particular language. Recordings were processed at a rate of 16 000 samples per second and a resolution of 16 bits, though there is some variation in some of the original recordings. In total there are 219 word recordings and each can be classified by the language, the gender of the speaker and the number being spoken (see S.2 in supplementary material for more details). Observations of the same word being spoken in different languages are treated as sharing the same word attribute. For example the word ‘four’ includes recordings of ‘quatre’ (French) and ‘quattro’ (Italian) as well as the word ‘four’ in other languages. Integers were chosen because these have no ambiguity in terms of translation making comparison of their use across languages straightforward. Furthermore, the cardinals ‘one’ to ‘ten’ of Romance languages (among many other words) stem from shared Latin forms [59]. This suggests that these words might also be suitable when comparing languages acoustically.

In this paper, the observations are modeled as functional data as is becoming increasingly common in studies involving sound recordings (e.g., [34]). Such models make the reasonable assumption that the data have been obtained by observing an underlying function at finitely many discrete points along a continuum, and that this underlying function is smooth (i.e., a certain number of derivatives exist).

The overall duration of a word can vary significantly per speaker as can the timings of intra-word elements (for instance syllables). Thus, to adjust for these differences all observations within a word grouping undergo registration (also known as alignment or warping, see [47,60], with a modified version of Tang and Müller [71] being used for the actual alignment). As has been noted in [48], it is very important in acoustic processing that distortion of frequencies does not occur. Thus, a short-time (10 ms window) Fourier transform is taken of each audio recording to produce a spectrogram, a spectrogram being a two-dimensional representation of audio signal energy intensity in frequency–time space [27]. Registration then took place in the spectrogram space (purely on the time dimension) using the modified Tang and Müller [71] method, where the modification accounted for the two-dimensional nature of the object under registration (see Pigoli et al. [58] for more details). Spectrograms are a natural choice for representing power with functional data [34,48], though approaches such as Mel-frequency cepstra can provide possible alternative representations [17]. In Holan et al. [34], spectrograms of mating calls are used as predictors of mating success of treehoppers. Martinez et al. [48] investigate regional differences in bat chirps by considering a functional mixed model with spectrograms as the image response. In contrast, the emphasis of this paper is not to seek a model that acts as a data generating process. Instead it is to identify meaningful low dimensional representations of spectrograms that highlight differences between languages, and subsequently assess whether these distinctive acoustic features are compatible with the class of Gaussian latent tree models.

As part of the data preprocessing, the standardization of word duration results in the time dimension being measured in generic time units (i.e.,  $\mathcal{T} = [0, 1]$ ). The value stored at a frequency–time point is a function of the power (or amplitude), with frequencies binned every 100 Hz up to the Nyquist frequency of 8000 Hz (i.e.,  $\mathcal{F} = (0, 8000]$ ). The resulting spectrograms are stored as matrices of 81 frequency by 100 time points. Fig. 1 is the spectrogram of a female French speaker saying the word ‘quatre’. Broadly, this interpolated plot indicates that there is greater power in the lower frequencies, and that the beginning and the end portions of the standardized time period are quieter.

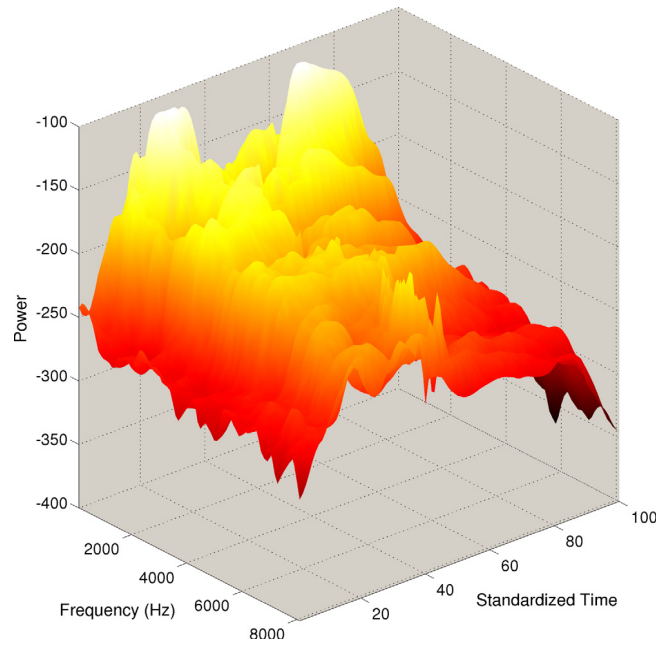
### 2.1. Notation

The underlying function of each spectrogram is denoted  $x_{\ell,m}^{d,g}(f, t)$  with the two dimensions  $f$  and  $t$  referring to frequency and time respectively. Recall that each spectrogram is derived from a spoken word—the subscripts and superscripts encode observational information:  $\ell = 1, \dots, n_\ell$  denotes the language being spoken;  $d = 1, \dots, n_d$  indicates the word being spoken;  $m = 1, \dots, m_{\ell d}$  is a counter where  $m_{\ell d}$  is the number of observations of word  $d$  from language  $\ell$ ;  $g$  refers to the gender of the speaker.

It is well documented that there are differences in the acoustics of male and female speakers which go beyond a simple shift in the spoken frequencies (for instance [55,57]). Parris and Carey [56] present a statistical method for discriminating between speaker gender of short acoustic recordings. In their analysis of seven Indo-European languages (of which Romance is a subset), gender was correctly identified on average 98% of the time. This suggests that there are commonalities in acoustic gender differences across Indo-European languages. In light of this result, it is judged that gender should be adjusted for at the macro level:

$$x_{\ell,m}^d(f, t) = x_{\ell,m}^{d,g}(f, t) + \tilde{x}^g(f, t)$$

where  $\tilde{x}^g$  is the difference between the mean of all samples with gender  $g$  and the mean of all samples. Henceforth it will be the gender adjusted function that will be the object of interest in this paper.



**Fig. 1.** Post-registration spectrogram of female French speaker saying 'quatre'. It can be seen that there is greater power in the lower frequencies, and that the very beginning and end of the word are unsurprisingly two of the quietest regions.

The mean spectrograms for language  $\ell$ , word  $d$  are defined in (1), for language  $\ell$  in (2), and the grand mean spectrogram in (3).

$$\bar{x}_{\ell}^d(f, t) = \frac{1}{m_{\ell d}} \sum_{m=1}^{m_{\ell d}} x_{\ell, m}^d(f, t) \quad (1)$$

$$\bar{x}_{\ell}(f, t) = \frac{1}{m_{\ell}} \sum_{d=1}^{n_d} m_{\ell d} \bar{x}_{\ell}^d(f, t) \quad (2)$$

$$\bar{x}(f, t) = \frac{1}{m_{..}} \sum_{\ell=1}^{n_{\ell}} m_{\ell} \bar{x}_{\ell}(f, t) \quad (3)$$

where  $m_{\ell} = \sum_{d=1}^{n_d} m_{\ell d}$ ,  $m_{..} = n = \sum_{\ell=1}^{n_{\ell}} m_{\ell}$ , and for  $t \in \mathcal{T}$ ,  $f \in \mathcal{F}$ . The parameters  $m_{..}$  and  $n$  will be used interchangeably depending on whether summation is being emphasized.

### 3. Group-based projections of functional data

Dimension reduction is a well-studied area of statistics with tools such as principal component analysis (PCA) and multidimensional scaling (e.g., see [16,38]) having widespread use. Functional counterparts of such techniques have also been formulated, for example functional principal component analysis (FPCA) [14,62,74], and functional multidimensional scaling [50]. In this paper, we further investigate a somewhat less well studied area of dimension reduction, canonical variates analysis (CVA) and its functional equivalent canonical function analysis (CFA) which explicitly base their dimension reductions on discriminatory power, with very strong connections to the theory of linear discriminant analysis [25], a widely used classification tool.

The acoustic data presented in this paper benefit from a dimension reduction in two main ways. First and foremost, dimension reduction provides a route to feature extraction whilst also reducing unwanted noise. Second, if subsequent to the reduction it is found that  $n \geq p$  then techniques which make use of inverse covariances can be implemented straightforwardly. If this is not so, standard estimates will produce singular sample covariance matrices. Of course these benefits must be balanced against potential information loss from the data reduction. One approach to feature extraction that mitigates against this loss is to find an ordered basis which prioritizes one or more characteristics of interest. Thus by projecting data onto the first few components of such a basis the most prominent aspects of the data are retained whilst what remains is treated as noise. In the cases of PCA and FPCA, the dimension reduction is optimized so as to efficiently capture modes of variation. Such techniques are often used in linguistic and semantic analyses, for example [44,73]. These techniques could also be used in the subsequent analysis below, with little difficulty.

It would be possible to investigate many appropriate measures of (dis-)similarity for our tree-based analysis which will be given in Section 4. However, unless a specific evolutionary model is chosen it is unlikely that there exists a *best* measure for this purpose. As our focus is on macro-language comparisons, we argue that the feature of interest is the between- to within-language covariance, and it is this which should directly inform the method selected to construct a basis. When data is known a priori to be grouped then CFA and its multivariate analogue CVA are standard techniques implemented to select variables to discriminate between groups. These tools are therefore the starting points for our analysis.

### 3.1. Canonical function analysis

Here we present CFA as a tool for FDA to produce a basis which maximizes between- to within-group variation (subject to the basis component functions being uncorrelated) with the intention of achieving an efficient dimension reduction. The finer details of CFA can be found in Kiiveri [39]. To illustrate CFA, consider a set of one-dimensional functional data denoting the between-group covariance function as  $B(u_1, u_2)$  and within-group covariance function as  $W(u_1, u_2)$  (with  $u_1, u_2 \in \mathcal{U}$ ,  $\mathcal{U} \subset \mathbb{R}$ ) where both functions are considered to be bounded and piecewise continuous. The aim is to identify canonical functions  $h_q(u)$  such that between-group variation is maximized relative to within-group variation under the restriction that each canonical function is uncorrelated to every other. This is expressed as a generalized eigen-equation which simplifies to:

$$\int_{\mathcal{U}} \{B(u_1, u_2) - \lambda_q W(u_1, u_2)\} h_q(u_2) du_2 = 0 \quad (4)$$

with eigenvalues  $\lambda_q$  and eigenfunctions  $h_q$ . There may be countably infinite solutions to Eq. (4) but in discretized estimation, only a maximum of  $s$  (say) will have non-zero  $\lambda_q$ . Pairs of canonical functions and real numbers  $(h_1(u), \lambda_1), \dots, (h_s(u), \lambda_s)$  can be found by solving (4) numerically, where  $\lambda_1, \dots, \lambda_s$  is a monotone decreasing sequence. An  $r$ -dimensional projection of the data is obtained using the first  $r$  canonical functions, and this projection is such that the between- to within-group covariance is maximally retained.

### 3.2. Separable covariance functions

Recall that the Romance data set comprises spectrograms which have time and frequency directions. A straightforward model for describing how these directions interact is that of separable covariance. The assumption underpinning this model can be encapsulated as there being no dependency between the (standardized) time and frequency of the data. This is, of course, a significant simplification of the likely underlying model, but comes with significant computational savings. It is linguistically justifiable over larger time scales (because consonants and vowels can be combined in very many different ways, both within a language and in different languages, meaning that frequency variation in a language is largely independent of temporal position), while from a statistical point of view is an intermediate position between a deterministic basis (such as a wavelet or Fourier basis) with its computational advantages, and a fully non-parametric basis with its inherent data driven nature, making it an appealing basis to work with. While we will define the basis through separable covariance structures (see below), it has recently been shown that in many other circumstances a basis defined by marginal covariances (an equivalent representation) is appropriate even in cases where separability does not hold [15].

Recall that a covariance function  $C$  is said to be separable if

$$C\{(f_1, t_1), (f_2, t_2)\} = C_f(f_1, f_2)C_t(t_1, t_2)$$

where  $C_f$  and  $C_t$  are functions only of their arguments. The factored covariances provide an understanding of how frequency or time dimensions of the spectrograms vary when the other has been averaged out. However, the main purpose of the assumption becomes apparent for the Romance data set subsequently (as described in Section 3.3.1) when use of the separable model of covariance overcomes the challenge of covariance rank deficiency.

Under the separable covariance assumption for two-dimensional data the CFA optimality equation equivalent to (4) is:

$$\int_{\mathcal{F}} \int_{\mathcal{T}} \{B_f(f_1, f_2)B_t(t_1, t_2) - \lambda_q W_f(f_1, f_2)W_t(t_1, t_2)\} h_q(f_2, t_2) dt_2 df_2 = 0. \quad (5)$$

It can be easily shown that the solutions to this equation can be obtained as the product of the solutions to two CFAs performed on the frequency and time covariances separately. Thus given any canonical function pairs  $\{h_{qf}(f_2), \lambda_{qf}\}$  and  $\{h_{qt}(t_2), \lambda_{qt}\}$  from a frequency and time CFA respectively, the products provide a solution to (5):  $h_q(f_2, t_2) = h_{qf}(f_2)h_{qt}(t_2)$  and  $\lambda_q = \lambda_{qf}\lambda_{qt}$ . Moreover, any solution to (5) can be obtained from such products. This result is useful when proceeding to obtain numerical solutions to (5).

### 3.3. Canonical variate analysis

Although less frequently implemented than PCA, the theory of CVA as a multivariate tool has been well developed in Krzanowski [41, Chapter 11]. However, beyond being a purely multivariate technique, CVA can also be used with functional data as an approximation to CFA as is presented in Kiiveri [39]. The technicalities of implementing CVA do not differ



whether in a functional or multivariate setting, although it is sometimes necessary to interpret their outputs differently, as is encountered with other tools (e.g., [24]).

In practice spectrograms are often discretized representations of underlying functions, and so each function  $x_{\ell,m}^d$  is instead given by a matrix  $\mathbf{X}_{\ell,m}^d$  with time–frequency dimensions  $n_f \times n_t$  (i.e., the number of sample points of the frequency and time). These finite approximations tend to be high dimensional and so the question of dimension reduction is pertinent. By concatenating the rows of these matrix representations of spectrograms the data corresponds to the vector description of CVA. As CVA considers each covariance entry independently of its adjacent values, this does not affect the implementation of CVA. The only notable downside of concatenation is that it can obscure visual representation and description of the data.

Recall that the aim of CVA is to find directions which discriminate the groups in the data, and does this by finding successive uncorrelated vectors  $\mathbf{a}$  that form linear combinations  $y = \mathbf{a}\mathbf{x}^\top$  (where  $\mathbf{x}$  is  $1 \times p$  data vector) that maximize the ratio of the between-groups covariance ( $\mathbf{B}$ ) to the within-groups covariance ( $\mathbf{W}$ ). In the context of the paper,  $\mathbf{B}$  describes the variation between the per-language mean spectrograms and the grand mean spectrogram, whereas the  $\mathbf{W}$  describes the variation between individual observations and the associated per-language mean spectrograms.

Finding the optimal  $\mathbf{a}$  is equivalent to solving  $(\mathbf{W}^{-1}\mathbf{B} - \lambda\mathbf{I})\mathbf{a}^\top = \mathbf{0}$  where  $\lambda \in \mathbb{R}$ . This reduces to performing an eigenanalysis on  $\mathbf{W}^{-1}\mathbf{B}$ ; the eigenvector corresponding to the largest eigenvalue is the optimal  $\mathbf{a}$ . As with CFA, canonical pairs  $(\mathbf{a}_r^\top, \lambda_r)$  are sought. These are found through a full eigenanalysis of  $\mathbf{W}^{-1}\mathbf{B}$  such that  $\lambda_1 > \dots > \lambda_s > 0$ , where  $s = \min(p, n_\ell - 1)$  is the number of non-zero eigenvalues of  $\mathbf{W}^{-1}\mathbf{B}$ . Thus  $(\mathbf{a}_r^\top, \lambda_r)$  produces the  $r$ th greatest ratio of between- to within-language variability. Hence, the optimal projection to  $r$  dimensions requires only  $(\mathbf{a}_1^\top, \lambda_1), \dots, (\mathbf{a}_r^\top, \lambda_r)$ . Theoretically, CVA is designed for use with Gaussian data and assumes that within-group covariance is equal across groups. If these assumptions hold true then CVA is optimal for identifying modes of variability that distinguish between groups.

### 3.3.1. Separable-CVA

As mentioned in Section 3.2, the overall solutions to a CFA optimality problem with a separable covariance structure can be found as the product of solutions to CFAs of the decomposed covariance functions. We propose combining a tensor decomposable covariance structure with CVA in order to obtain numerical solutions to the decomposition of the separable-CFAs. This, when taking products, also gives solutions to the overall CFA.

The main purpose of assuming a tensor-decomposable covariance structure is to overcome the obstacle of rank-deficient sample covariance matrices caused by the length of the observations exceeding the number of observations (i.e.,  $p > n$ ). This is not just a problem with the Romance speaker data set but is commonly encountered with functional data sets due to their often high-dimensionality (e.g., [46]). Rank deficiency obstructs using CVA to obtain numerical solutions to CFA. Theoretically in CFA an inverse function  $\mathbf{W}^{-1}$  is neither required nor is usually bounded, whereas in CVA  $\mathbf{W}^{-1}$  is needed for the eigenanalysis of  $\mathbf{W}^{-1}\mathbf{B}$  but cannot be obtained because in this case  $\mathbf{W}$  is singular.

In the observational matrix setting,  $\mathbf{C}$  is separable if:

$$\mathbf{C}\{(f_1, t_1), (f_2, t_2)\} = \mathbf{C}_f(f_1, f_2) \otimes \mathbf{C}_t(t_1, t_2)$$

where  $\otimes$  is the standard Kronecker product. Using known results of the Kronecker product (see [42] for example), the separability assumption in the multivariate setting implies:

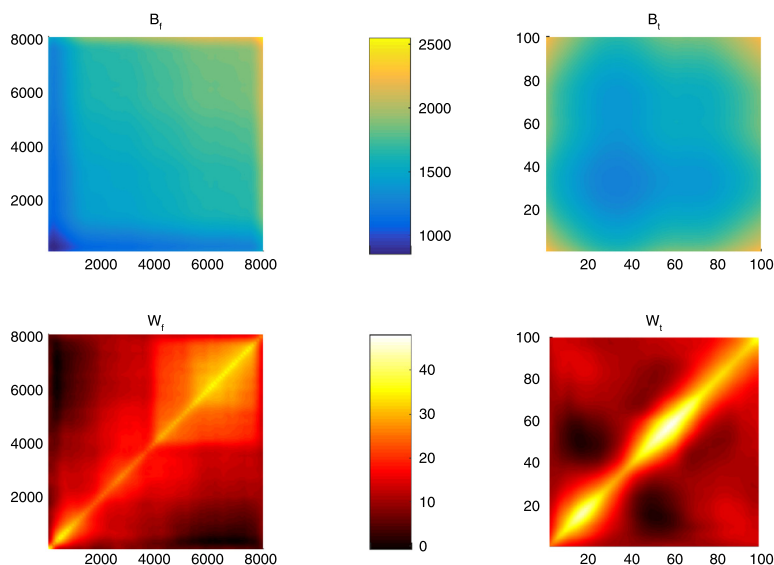
$$\mathbf{W}^{-1}\mathbf{B} = (\mathbf{W}_t^{-1} \otimes \mathbf{W}_f^{-1})(\mathbf{B}_t \otimes \mathbf{B}_f) = \mathbf{W}_t^{-1}\mathbf{B}_t \otimes \mathbf{W}_f^{-1}\mathbf{B}_f \quad (6)$$

where the estimates of separate within- and between-language covariance matrices in the frequency direction are:

$$\begin{aligned} \hat{\mathbf{B}}_f(f_1, f_2) &= \frac{1}{n_\ell - 1} \sum_{\ell=1}^{n_\ell} \frac{m_\ell}{n_t} \sum_{t=1}^{n_t} \tilde{\mathbf{X}}_\ell(f_1, t) \tilde{\mathbf{X}}_\ell(f_2, t) \\ \hat{\mathbf{W}}_f(f_1, f_2) &= \frac{1}{n - n_\ell} \sum_{\ell=1}^{n_\ell} \sum_{d=1}^{n_d} \sum_{m=1}^{m_{\ell d}} \frac{1}{n_t} \sum_{t=1}^{n_t} \tilde{\mathbf{X}}_{\ell,m}^d(f_1, t) \tilde{\mathbf{X}}_{\ell,m}^d(f_2, t) \end{aligned}$$

where  $\tilde{\mathbf{X}}_\ell(i, j) = \bar{\mathbf{X}}_\ell(i, j) - \bar{\mathbf{X}}(i, j)$  and  $\tilde{\mathbf{X}}_{\ell,m}^d(i, j) = \mathbf{X}_{\ell,m}^d(i, j) - \bar{\mathbf{X}}_\ell(i, j)$  with equivalent estimates for the time direction. Treating each frequency and time sample as a separate observation leads to the product covariance matrices  $\mathbf{W}$  and  $\mathbf{B}$  having higher ranks than previously. Explicitly, for  $\mathbf{W}^{-1} = (\mathbf{W}_f \otimes \mathbf{W}_t)^{-1}$  to be nonsingular, we need that  $nn_f \geq n_t$  and  $nn_t \geq n_f$ . This is equivalent to requiring  $n \geq \max(n_f, n_t) / \min(n_f, n_t)$ . This contrasts with the previous condition  $n \geq p = n_f n_t$ . So the new requirement is usually significantly more relaxed, and CVA can often then be implemented.

An eigenanalysis of  $\mathbf{W}_f^{-1}\mathbf{B}_f$  produces eigenvalues  $(\lambda_{f1}, \dots, \lambda_{fnf})$  and corresponding eigenvectors  $(\mathbf{c}_{f1}, \dots, \mathbf{c}_{fnf})$  with equivalent output for the time covariances  $\mathbf{W}_t^{-1}\mathbf{B}_t$ . Sorting decreasingly,  $(\lambda_{f1}, \dots, \lambda_{fnf}) \otimes (\lambda_{t1}, \dots, \lambda_{tn_t})$  produces a vector denoted  $(\lambda_1, \dots, \lambda_{nfn_t})$  and the Kronecker product of the corresponding eigenvectors results in matrices denoted  $(\mathbf{c}_1, \dots, \mathbf{c}_{nfn_t})$  of size  $n_f \times n_t$ , solving the overall CVA. It should be noted that while this basis defined is based on an assumption of separability, it nevertheless provides a complete basis of the space; see an analogous argument for separable PCA in Aston and Kirch [5].



**Fig. 2.** Sample between-language and within-language covariances of speech data for frequency and time directions. There is a clear ridge along the diagonal of the within-group covariances indicating that similar times and frequencies are highly positively correlated. The higher correlations in the high frequencies of the within-group frequency covariance are associated with recordings at lower audio sampling rates capturing fewer details in these ranges of frequencies. Rerunning the analyses performed in Section 5 whilst excluding these higher frequencies produces broadly similar results. Alternatively, excluding observations with low sampling rates has been examined in [66, Section 7.1] with both analyses leading to the same conclusion regarding tree-amenability.

### 3.4. Application of CVA as an approximation to CFA

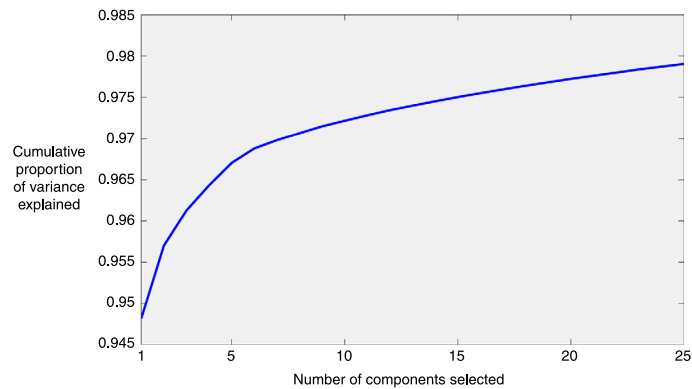
As motivated, the separable-CVA is used to approximate the separable-CFA of the Romance languages data to achieve a dimension reduction based on components which maximize between- to within-language variability. This is a suitable approximation to make as the functional spectrograms have been sufficiently densely sampled during discretization.

Interpolated plots of the between- and within-language covariances for both the time and frequency directions of the spectrograms are given in Fig. 2. The  $B_f$  plot displays positive covariance which increases as the frequency increases. The time covariance  $B_t$  is positive throughout with an internal minimum and the highest covariances corresponding to beginning and end time points. These particularly high time corner covariances could indicate genuine similarities at the beginnings and ends of spoken words or could be an artifact of the spectrogram registration. The covariances  $W_t$  and  $W_f$  both have a ridge along the diagonal which reassuringly suggests that similar time and frequency points have strong covariances within languages. The within-language covariances were each tested for equality with the remaining data within-language covariance using Box's M statistic [12]. None of the covariances were found to be significantly different at the 0.01 level and thus the CVA assumption of homogeneous within-group covariances is not rejected.

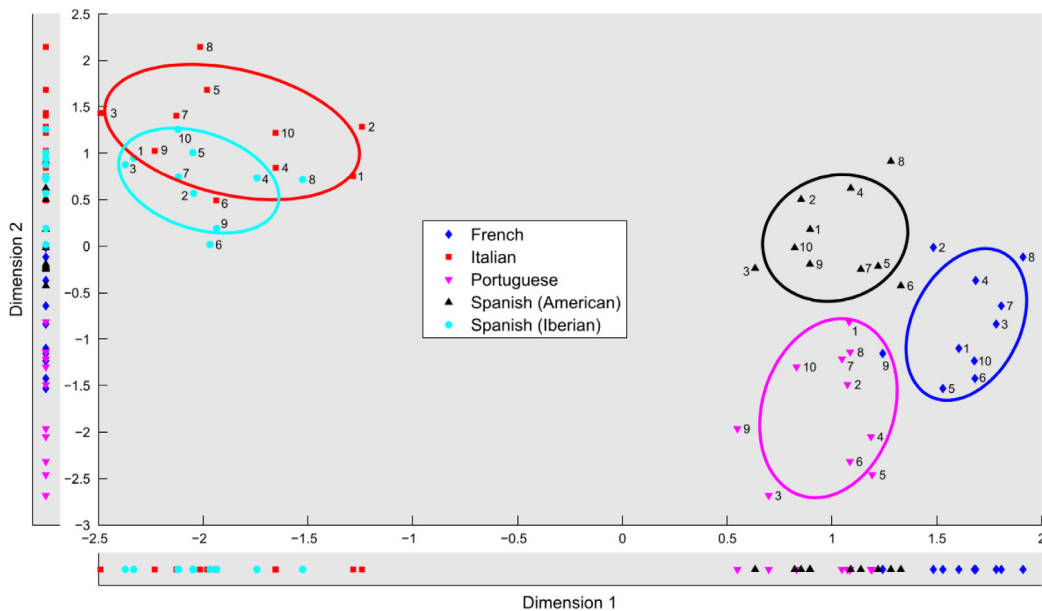
When selecting a dimension  $r$  to project to, it is unusual to have anything but an arbitrary albeit sensible method for selecting  $r$ . However, in some acoustic contexts (e.g., [31]) thresholds can be proposed based on sounds which are audible to humans. Otherwise, equivalent techniques to those employed with PCA [38] can be used. For this linguistic study Fig. 3 shows the cumulative variation explained by selecting particular numbers of components. This indicates that almost 95% of the between- to within-language variance can be explained by a single component, and an additional two components take this figure to over 96%.

Once a dimension  $r$  has been selected, each observation from the language data can be projected into  $r$ -dimensional space:  $\mathbf{Y} = \mathbf{A}\tilde{\mathbf{X}}^T$ , where  $\mathbf{A}$  is  $r \times p$  with columns  $\mathbf{c}_1, \dots, \mathbf{c}_r$  and where  $\tilde{\mathbf{X}}$  is  $1 \times p$  and formed by concatenating the  $n_f$  rows of length  $n_t$  of the observation  $\mathbf{X}$ . The sub and superscripts of the observations are omitted solely for notational clarity.

It is clear from the description that separable-CVA is not implemented on the belief that it reflects some underlying data generating process. Nonetheless, separable-CVA is a practical tool that produces a meaningful representation of the original data in lower dimension. To demonstrate the effectiveness of projection of the spectrograms in even two dimensions, the projections of the means of the word observations are plotted in Fig. 4. The results of the separable-CVA are encouraging as there are clear groupings of the projected word means from the same languages. Furthermore, the first two dimensions appear to reflect aspects known about the languages. The second dimension seems to distinguish the four distinct languages, and then the first dimension is able to separate the Spanish dialects (while also providing clear distinctions between some of the other languages). This is certainly not to imply interpretation of the dimensions as dialect and language related, but rather that even dialects do have discriminatory properties in the data set. Given that the acoustic data are undoubtedly noisy, the figure indicates the effectiveness of separable-CVA at selecting components which discriminate on a group basis. Note that whilst CVA operates on all languages simultaneously rather than in a pairwise manner, this does not necessarily



**Fig. 3.** Cumulative variation explained by number of components. The explanatory power of the first component in terms of between- to within-language variability is over 94%.



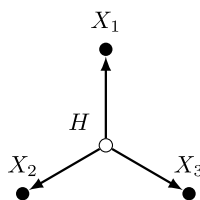
**Fig. 4.** Two dimensional separable-CVA projection of means of word observations with ellipses representing the standard deviation around the mean of plotted points. The second dimension is effective at distinguishing at the language level. The first dimension then discriminates between the Spanish dialects. As each component can relate to a different set of speech qualities, the proximities of the languages can differ depending on the dimension projection. For example, in the second dimension, the nasality of the speech appears to contribute to the separation of Italian and Portuguese. This is further explored in Section 5.

imply languages in close proximity post-projection share particular acoustic features. However, this can be examined in more detail, for example through studying the Hadamard matrices of each projection as is illustrated in Fig. 8 (with further components given in the supplementary material (see Appendix A)).

#### 4. Constraint diagnostics for tree-amenability

The aim of this study is to examine the suitability of a tree for acoustic functional data. It is of particular interest to know whether there are certain features of these spoken Romance languages which could have developed over time in the manner of an evolutionary tree. CFA (or CVA) is compatible with this aim as it effectively identifies components with features that distinguish between languages. Thus a projection to  $r$  dimensions provides  $r$  different combinations of characteristics which could potentially be modeled as a tree. In order to help assess the plausibility of these  $r$  hypotheses the topic of tree constraints is introduced. These are algebraic and semi-algebraic constraints which mathematically must be obeyed if a data set corresponds to the leaves of a Bayesian network which is a tree. Note here that the non-leaf vertices correspond to various unobserved ancestral languages. While there has been considerable recent interest in understanding graphical structures, including trees, in very general settings [45], these approaches are designed for settings when all nodes (both internal and leaves) are observed. Here the concern is with determining constraints when the internal nodes are unobserved.





**Fig. 5.** Tripod tree. A strictly trivalent tree with  $n_\ell$  leaves contains  $\binom{n_\ell}{3}$  conditional independence relationships which can each be expressed as a tripod tree.

Examining whether these constraints are respected for a given data set determines whether that data set adheres to a tree structure. If it respects the constraints of interest the data set is said to be tree-amenable. Provided with  $r$  component-by-component projections, a diagnostic for tree-amenable can be applied to each (see e.g., [67] for a binary example). This can then be used as an exploratory tool to give insight as to whether the distinguishing characteristics captured by a component could have an evolutionary tree structure. To begin with a brief overview is provided, detailing the concept of tree constraints and how they form a natural choice for use with data sets such as the acoustic recordings in this study.

#### 4.1. Constraints on observed covariances respecting latent tree models

Tree constraints are useful for studies investigating evolutionary relationships between both observable and unobservable variables. These evolutionary relationships can be expressed graphically as phylogenetic trees where traditionally the objects of interest have been biological species (e.g., [23]), however, this idea naturally extends to other fields such as linguistics (e.g., [20]). Phylogenetic trees are of particular statistical interest when formally considered as Bayesian networks describing conditional independence relationships (e.g., [21]). In the case of discrete graphs in which all variables are assumed to be observed, Loh and Wainwright [45] make use of precision matrices to assess graph structure. However, typically Bayesian networks have both observed and hidden variables and these are much more complex to analyze. Furthermore, the data set in question is functional and thus the analyses are embedded in Gaussian assumptions and observations treated as Gaussian processes as opposed to individual points. This is where Gaussian tree constraints can provide useful insight into the difficult question of tree-amenable.

A graph  $T$  is a tree if there exists a unique path between any two connected vertices. The trees of interest are those with interior hidden nodes  $H \in \mathcal{H}$  and manifest leaf nodes  $X \in \mathcal{X}$  (denoted by white and black circles respectively). The observed leaf variables can be thought of as contemporary languages, and the latent interior variables as past versions of languages. Attention is restricted to binary trees, i.e., trees where each interior node has three adjacent nodes. Any tree from this class with more than three leaves can be expressed as a bifurcating tree, a common model in phylogenetics (e.g., [22]). The tripod tree (shown in Fig. 5) describes any such relationship between three observed variables. To see this, simply note that any three connected leaves of a strictly trivalent tree share a unique interior node, and that the conditional independence of these four nodes is minimally described by the tripod tree. Hence, the tripod tree is a fundamental component for analyses involving tree constraints.

The distributions associated with moments of the observed variables of such tree models are known for some settings. The case of univariate binary random variables was expanded in Settini and Smith [65] and has recently attracted considerable interest (for example [1,76]). These advances have encouraged some authors to proceed to use the known geometry of these spaces to support inference and learning over the space of tree models (see [19]). These focus on the polynomial constraints that are implicit in these models. Zwiernik and Smith [75] fully characterized the space for binary trees by extending the understanding beyond algebraic relationships to also include the active semi-algebraic constraints, whilst Allman et al. [3] broadened the results to a  $k$ -state  $n$ -pod tree. Here the interest is in the use of the constraints which define these geometries for assessing whether data could have originated from a tree model, in particular whether the spoken languages in the acoustic data set could be modeled as a phylogenetic tree.

#### 4.2. A constraint on the covariance of Gaussian latent tree models

Most of the results in the literature of tree constraints are only applicable in settings with binary random variables. Alternative constraints specific to Gaussian tree models are required in order to perform equivalent tree-amenable diagnostics for Gaussian functional data such as the Romance data set. Here we describe a fundamental univariate Gaussian tree constraint which has not yet been exploited for the purpose of investigating tree-amenable.

Consider the precision matrix  $\Sigma^{-1}$  related to the tripod tree in Fig. 5 with one latent and three manifest Gaussian random variables. It is well known in the Gaussian setting that if  $X_i$  is independent of  $X_j$  conditional on all other observed variables then the corresponding entry  $\Sigma_{ij}^{-1} = 0$  (see [43, Chapter 5] for example). So the precision matrix for the univariate Gaussian

tripod tree has the form:

$$\Sigma^{-1} = \left( \begin{array}{ccc|c} \sigma_1^{-1} & 0 & 0 & \sigma_{1H}^{-1} \\ 0 & \sigma_2^{-1} & 0 & \sigma_{2H}^{-1} \\ 0 & 0 & \sigma_3^{-1} & \sigma_{3H}^{-1} \\ \hline \sigma_{1H}^{-1} & \sigma_{2H}^{-1} & \sigma_{3H}^{-1} & \sigma_H^{-1} \end{array} \right)$$

where the final row/column relates to the interior hidden variable. The covariance matrix resulting from taking the inverse of the precision matrix can be expressed algebraically in terms of entries of  $\Sigma^{-1}$ . Using the resulting entries of the covariance  $\Sigma = (\sigma_{ij})$  it is then straightforward to calculate that a necessary condition for the leaves of this tree to be the margin of the tripod tree given in Fig. 5 is that

$$\forall_{i < j < k} \quad \sigma_{ij}\sigma_{ik}\sigma_{jk} \geq 0. \quad (7)$$

This constraint shall be denoted positivity constraint. For a strictly trivalent tree with  $n_\ell$  observed leaf nodes, there are  $\binom{n_\ell}{3}$  tripod trees that must be valid – one for each triple – and thus  $\binom{n_\ell}{3}$  such positivity constraints. Whilst there are further constraints imposed on the observed moment space of trees (see Shiers et al. [68]), the derivations are considerably more involved and not the emphasis of this methodological paper. Given that applications of Gaussian tree-constraints are not apparent in the literature, there is still much to explore focusing solely on this fundamental Gaussian tree constraint.

These constraints can be used with the linguistic data set as a diagnostic for assessing tree-amenability. More precisely, for the  $r$  component-by-component projections, each can be assessed for tree-amenability. Thus some components may be found to violate the tree constraints whereas others may satisfy them. This may suggest which attributes of the spoken languages have evolved in a tree-like manner and which have not. Note that although we have functional data, the formulation of this positivity constraint does not assume this, and thus the constraint is equally valid for multivariate Gaussian data.

#### 4.3. Constructing a suitable covariance statistic

In pursuit of assessing tree-amenability of the  $r$  projections of the Romance data using the positivity constraint  $\sigma_{ij}\sigma_{ik}\sigma_{jk} \geq 0$ , it is clear that a sample covariance of the scores must be constructed. Recall that the relationships of interest in this study are at the language level and thus between-language covariances (each  $5 \times 5$ ) are the appropriate statistics to produce, one for each of the  $r$  components. One approach to calculating the entries of these matrices is to treat the mean score of each word in a language as an observation and then measure the distance from the overall word mean projection. Then using appropriate weights, a between-language covariance matrix can be estimated as follows. Let

$$\bar{y}_d^i = \frac{1}{m_d} \sum_{\ell=1}^{n_\ell} m_{\ell d} \bar{y}_{\ell d}^i, \quad m_d = \sum_{\ell} m_{\ell d},$$

where recall  $m_{\ell d}$  is the number of samples of word  $d$  in language  $\ell$ , and  $\bar{y}_{\ell d}^i = \mathbf{c}_i \bar{x}_{\ell d}$  the projection of the mean of word  $d$  of language  $\ell$  using component  $\mathbf{c}_i$ . Then for component  $i$  the between-groups cross-covariance for the projected data has the following form

$$\Sigma_{\mathbf{y}_i} = (\sigma_{\ell, \ell'}^i) \quad \text{where} \quad \sigma_{\ell, \ell'}^i = \sum_{d=1}^{n_d} \frac{\sqrt{m_{\ell d}} \sqrt{m_{\ell' d}} (\bar{y}_{\ell d}^i - \bar{y}_d^i) (\bar{y}_{\ell' d}^i - \bar{y}_d^i)}{n_d - 1} \quad (8)$$

where recall  $n_d$  is the number of unique words. Note that this between-group covariance differs from that used in the CVA—this is of the projected data and the word means are used to provide an observational summary of the data. This is a valid construction in the sense that (8) is an inner product (see [37] for instance). Furthermore, for the cross-covariance to be meaningful, equivalent statistics must be compared, in this case per language word means. The sample matrices  $\hat{\Sigma}_{\mathbf{y}_i}$  will be rank deficient if  $n_\ell \geq n_d$ . Also, observe that if for at least one word  $d$  the number of observations is unequal across languages then the weighted word mean  $\bar{y}_d^i$  differs from the unweighted version. This relaxes a zero-sum condition on the rows or columns of  $\hat{\Sigma}_{\mathbf{y}_i}$  permitting the covariance matrix to be full rank. In the alternate case of a balanced observational design, full rank can be achieved through an alternative construction (for example adding the unweighted word means back to each language-word mean).

Now component-by-component covariances can be used to indicate adherence to a Gaussian tree model using the tripod tree positivity constraint on all  $\binom{n_\ell}{3}$  selections of languages. Each component captures a different combination of variability. Thus it is not unexpected that some components may show violations of the constraint whereas others may indicate tree-amenability.

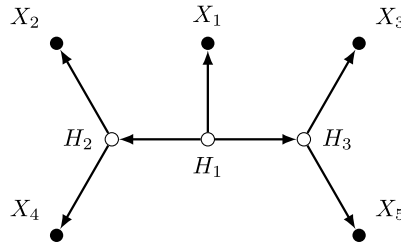


Fig. 6. Five leaf tree.

Table 1

Percentage of simulations that satisfy the positivity constraint for the sample size (left headings) and for the four components for each data set  $D$  and  $D^*$  (above headings), for parameter values  $a = 5$ ,  $b = 6$  in the structural equations.

Sample size	$D$				$D^*$			
	$c_1$	$c_2$	$c_3$	$c_4$	$c_1$	$c_2$	$c_3$	$c_4$
50,000	100	100	97	61	11	9	10	10
5,000	100	99	85	41	13	7	9	11
500	98	87	58	28	13	8	9	12
50	88	61	27	27	12	9	11	11

#### 4.4. Simulation

Before analyzing the Romance language data, a short simulation is performed. The purpose of the simulation is to demonstrate the effectiveness of the tripod positivity constraint at identifying tree-amenability. Also, by varying the simulation sample size the robustness of the positivity constraint is examined.

In order to mimic the linguistic data set, a tree with five leaves is used to generate data (see Fig. 6). Standard structural equations are used [8] adapted for a Gaussian tree to model the dependencies between nodes. The interior node  $H_1$  is the root and this variable is simulated from a Gaussian distribution. Data on nodes adjacent to the root are generated as linear combinations of this root simulation with the addition of Gaussian noise. In a  $p$ -dimensional scenario, the root node is simulated as follows:

$$H_1 \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V}_{h_1})$$

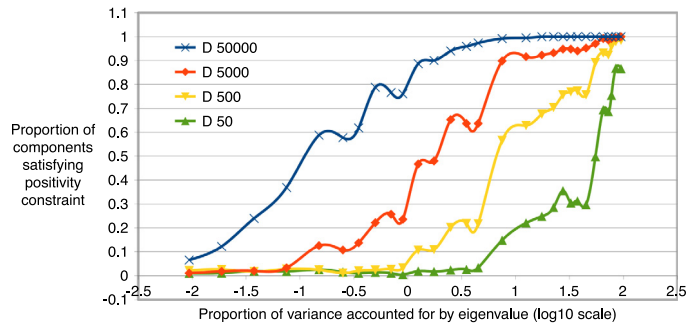
where  $\mathbf{V}_{h_1}$  is a  $p \times p$  symmetric positive definite matrix, constructed as  $\mathbf{A} + \mathbf{A}^\top + p\mathbf{I}_p$  with entries of  $\mathbf{A}$  generated independently from the continuous uniform distribution  $\mathcal{U}(0, 1)$ .  $\mathbf{0}_p$  is a  $p$ -vector of zeros, and  $\mathbf{I}_p$  is the  $p \times p$  identity matrix. Subsequent nodes are simulated from previous ones, for example:

$$X_1 = \lambda_{h_1x_1}H_1 + \epsilon_{x_1}.$$

Each entry of  $\lambda_{h_1x_1} \sim \mathcal{U}(-a, a)$ ,  $a \in (0, \infty)$ . Within a repetition, all entries of  $\lambda_{h_1x_1}$  are fixed for all observations, although do differ depending on the node transition being considered. Additionally,  $\epsilon_{x_1} \sim \mathcal{N}_p(\mathbf{0}_p, b^2\mathbf{I}_p)$ ,  $b \in (0, \infty)$ , and  $\epsilon_{x_1}$  is pairwise independent with all other random variables. Following the directions of the arrows, the remaining variables are simulated equivalently.

A four-dimensional data set  $D$  is generated for all eight nodes of the tree though only the observed leaf nodes are of interest and are utilized with the tree constraint. Another data set  $D^*$  is also generated using a corrupted version of the tree where each transition between nodes provides an opportunity for sign-reversal of entries of the data. Given this corruption can occur on any dimension of any of the data, the simulation is no longer from a Gaussian tree model. Although higher dimensional data could be generated, under standard CVA only a maximum of four ( $n_\ell - 1$ ) eigenvalues would be non-zero (the use of separable CVA relaxes this assumption in our analysis). Thus there would not be much information gained from the increased dimension but the computational resource required would certainly be greater. Four sample sizes are investigated: 50,000, 5,000, 500, and 50. These are average sample sizes as the simulation is designed such that the number of observations differs across languages allowing the use of (8) for the reasons relating to rank noted in Section 4.3. Each simulated data set undergoes the same basis change using CVA to approximate CFA as described in Section 3.3. For a full analysis, all four dimensions are retained for projecting the data, and thus for each data set four covariance matrices are calculated using the construction given in (8). The resulting covariance matrices are then used to check the positivity constraint  $\sigma_{ij}\sigma_{ik}\sigma_{jk} \geq 0$ . Having repeated the simulation 1000 times for each sample size the results are summarized in Table 1. For the projections using components  $c_1, \dots, c_4$  (ordered by explanatory power), the percentage of samples satisfying the positivity constraint is recorded.

There are three notable features highlighted by Table 1. First, the data sets  $D$  are found to be tree-amenable consistently more than  $D^*$ . Second, the performance of the positivity constraint on  $D$  decreases as the sample size decreases, particularly in the third and fourth components, whereas for  $D^*$  there is little difference. Third, it is the highly explanatory components



**Fig. 7.** For each sample size, interpolated plots are given of the estimated relationship between the explanatory power of components known to be tree-amenable and the chances that the components are correctly identified as being tree-amenable. This broadly suggests that the larger the sample size and/or the higher the explanatory power of a component, the higher the chances of a component being correctly identified as tree-amenable.

of  $D$  that are most effective at correctly satisfying the positivity constraint. Further investigation into this last feature can give a guide as to how much variance a component should account for before the positivity constraint becomes a reliable diagnostic. Using the  $D$  simulations, the proportion of times tree-amenable components are correctly identified as such is estimated over small ranges of the explanatory power; the interpolated estimates are plotted as the midpoints of the ranges. Although the results displayed are for when the structural equation parameter values are  $a = 5$ ,  $b = 6$ , similar results are found when these parameters are varied. In general, as  $a$  increases and  $b$  decreases the performance of the constraint  $D$  over  $D^*$  is more notable.

Fig. 7 displays graphically the property indicated in Table 1—that the lower the explanatory power of a component and the lower the sample size, the less reliable the positivity constraint is as an indicator of tree-amenability. However, the performance of the positivity constraint is still relatively high; even in the lowest sample size, the first component is effective 88% of the time. Furthermore, the simulations suggest the reliability of the constraint is not symmetric, particularly for low sample sizes; if a component satisfies the positivity constraint then this is good evidence of true tree-amenability, but if a component does not satisfy the positivity constraint then underlying tree-amenability should not be ruled out.

#### 4.5. Further simulation based on characteristics of observed data

In order to better assess the fundamental Gaussian tree constraint when applied to the data, a further simulation is performed using characteristics of the acoustic data set (e.g., sample size, eigendecomposition, sample variance). This is achieved by generating data from cross-covariance matrices  $\Sigma_Y$  (as in (8)) for which tree-amenability status is known. Then by reversing the eigenbasis projection as performed for the Romance data set, new observations are obtained. Once again, a lower dimensional representation is produced using a truncated eigenbasis projection. Then using the fundamental Gaussian tree constraint, an assessment can be made of whether tree-amenability of the source  $\Sigma_Y$  has an affect on the new sample being deemed tree-amenable.

Three scenarios are considered: (A) cross-covariance  $\Sigma_Y$  which is tree-amenable; (B)  $\Sigma_Y$  which is not tree-amenable; (C)  $\Sigma_Y$  is tree-amenable 50% of the time independently for each component and replication. For (A) and (B) the  $\Sigma_{Y_i}$  calculated from the Romance data set are used (for the analysis presented in Section 5). Each of these 15  $\Sigma_{Y_i}$  are thus known to be tree-amenable or not for the fundamental Gaussian tree constraint, and furthermore, the number of  $\binom{5}{3} = 10$  constraints satisfied is known and can be denoted  $T(\Sigma_{Y_i})$ . Recall,  $\Sigma_{Y_i}$  is only deemed tree-amenable if  $T(\Sigma_{Y_i}) = 10$ . That is, it satisfies the 10 positivity constraints comprising products of covariances relating to the 10 possible choices of three languages. Note that it is not possible for  $T(\Sigma_{Y_i})$  to be exactly 1, 2, 8 or 9 due to the appearance of each off-diagonal entry of  $\Sigma_{Y_i}$  in three of the ten combinations. For the first 15  $\Sigma_{Y_i}$  from the actual data set, the percentage of tree-amenable samples out 1000 generated is recorded against  $T(\Sigma_{Y_i})$ . The four  $\Sigma_{Y_i}$  which are tree-amenable provide higher proportions of tree-amenable samples than the remaining covariance matrices. Considering the full results (see Table S.1 in the supplementary material), a positive correlation is seen between  $T(\Sigma_{Y_i})$  and the percentage of samples simulated from  $\Sigma_{Y_i}$  which are then found to be tree-amenable. Considering scenario (C), a mix of the best and worst performing of the actual  $\Sigma_{Y_i}$  is used, having respectively 35% and 7% sample tree-amenability. The simulation results in an overall 18% tree amenability across all sampled components, suggesting that proportionately a similar level of tree-amenability is retained even when tree and non-tree components are combined.

This second simulation approach closely resembles the linguistic data set both in structure and parameter values. Thus these simulations provide comfort that even for the relatively small sample size the analysis is able to produce reasonable results. Scenarios (A) and (B) provide evidence that components that truly satisfy a low number of constraints are unlikely to give a false positive in terms of tree-amenability. Usefully they also identify that this risk grows as the number of constraints satisfied increases, and this can thus be considered when assessing the results. Scenario (C) demonstrates that tree-amenable components can be recovered even when combined with non-tree components, and furthermore, that these appear to be retained in the output almost proportionally to the input. The sampling in all these simulation scenarios may also provide

a proxy for distributional results of tree-amenability, where the higher the proportion of tree-amenability in a sample the more robust the conclusion. This hypothesis is an open problem which we look to address in a subsequent paper through derivation of explicit moments of relevant parameters.

## 5. Assessing tree-amenability of Romance languages

As illustrated in Fig. 4, an effective projection of the Romance data can be performed in even two dimensions. However, to account for a higher proportion of between- to within-language variability further components can be included for the projection. A threshold of cumulative between-language variability is set at 97.5% which for these data equates to including all components which account for at least 0.05% explanatory power. Application of this threshold provides a dimension reduction from 8100 to 15. Each one of these 15 components  $\mathbf{c}_1, \dots, \mathbf{c}_{15}$  accounts for some mode of variability between languages. Although the earlier components have high explanatory power, the latter components may isolate directions of variability which are of more interest from a linguistic perspective.

Applying the positivity constraint to each of the 15 component covariance matrices results in four of the components ( $\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_4, \mathbf{c}_6$ ) adhering to the positivity constraint. Recall that the simulations in Section 4.4 hint that for less explanatory components the constraint appears more likely to identify true tree components than false positives. This suggests that the identified tree-amenable components are more likely to truly be so, whereas the rejected components (e.g.,  $\mathbf{c}_3, \mathbf{c}_5$ ) may or may not be truly tree-amenable.

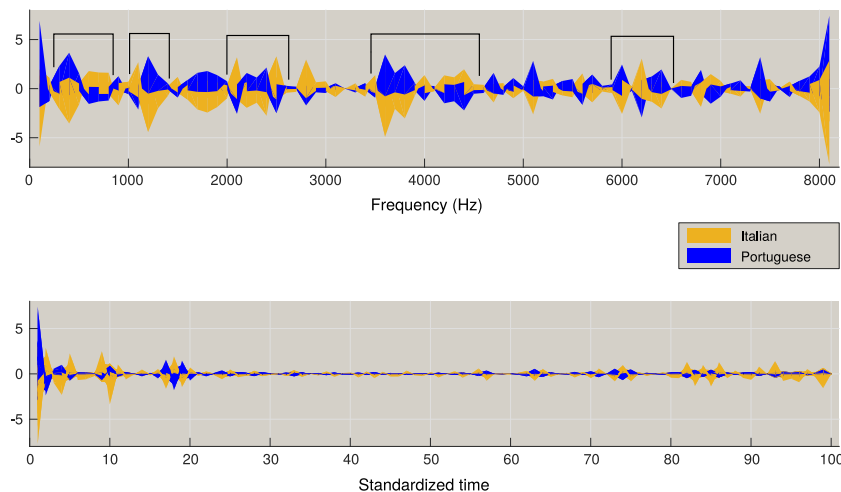
To develop an insight into which aspects of the languages are being identified by the separable-CVA, the Hadamard (entrywise) product of each component with each concatenated mean language spectrogram is calculated. Restoring the  $81 \times 100$  dimensions, the resulting matrices indicate the contribution of each frequency–time point to the overall co-ordinates produced by the component projections. This is useful for highlighting particular ranges of standardized time and frequency which distinguish between languages. For illustration, the time and frequency perspectives are plotted for the second component for both Italian and Portuguese in Fig. 8. These two languages are selected as they are clearly separated in both projected co-ordinates (Fig. 4), although the plots for the remaining languages are similar. Contrasting the plots in Fig. 8, immediately it can be seen that there are many frequency points which contribute to the overall projection whereas in the time dimension most points do not appear to be integral to distinguishing between languages in the second component. Considering the frequency perspectives, there is some symmetry in power between Italian and Portuguese. This is exhibited well in frequency ranges 300–800, 1000–1500, 2000–2500, 3500–4500, and 6000–6500 Hz which show reflections along the line of zero power and are indicated by brackets in Fig. 8. Some of these ranges are suggestive of possible differences though cannot be said to be definitive. The frequency ranges can be considered in relation to different phonetic features. For example, the 300–800 Hz range likely relates to the first formant F1 being different, due to vowel differences. The 1000–1500 Hz range likely relates to nasality (Portuguese is more ‘nasal’ than Italian, and therefore has less energy in this portion of the spectrum than Italian). The 3500–4500 Hz range is in the region of the third formant and could correspond to differences in lip rounding between speakers of the languages. The variation at the highest frequencies, around 6000 Hz, are likely to be due to idiosyncratic differences in speakers (since humans cannot readily control speech frequencies in that range) or in the recordings (equipment or recording location). Examining a smaller range of the data that excludes these high frequencies does not affect the main results of the analyses (data not shown).

Considering the time perspective, it is clear that the interesting time ranges are approximately 1–20 and 70–100. This suggests that the differences in the earliest and latest portions may be particularly effective at separating the languages. These results appear robust to trimming of the data set as well as to standardization of the covariance (data not shown), which suggests that the analysis is identifying more than just a feature of the data registration.

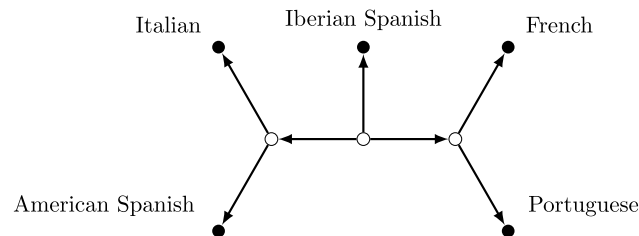
These projections are particularly effective at indicating graphically the dominant features of components which are often obscured when displayed numerically. Of course more detailed analyses are required to determine whether these are general features of the languages or simply of the particular data set studied. However, in either case, this type of exploratory data analysis is clearly a helpful one to perform before any more detailed analysis has taken place, especially if based on the firm assumption that the data originate from a tree.

The exploratory analysis can go one step further by also proposing a preliminary tree for tree-amenable components. Consider the correlation matrix corresponding to the covariance matrix  $\Sigma_{\mathbf{Y}_i}$ . For each entry  $\rho_{jk}$  make the transformation  $d_{jk} = -\ln\{(\rho_{jk} + 1)/2\}$  and use with an existing tree reconstruction method. For example, considering the second component once more, the UPGMA algorithm [49] produces a tree with topology as shown in Fig. 9. Observe the similarity to the projection in Fig. 4. While the proposed tree does differ from an accepted overall evolutionary tree for the Romance languages (it is expected that the two Spanish dialects would be closest, with these then close to Portuguese), it does allow the possibility of producing data extracted plausible evolutionary paths of very specific phonetic components, which might individually be different from the overall globally expected paths for the languages as a whole. Of course, the set of 10 words for numerals that we use in each language is admittedly a very small subset of the entire vocabulary, and that acoustic similarity alone is not all there is to linguistic relatedness. With a larger sample of the vocabulary that more reasonably approaches what we ordinarily understand by ‘a language’, we would hope that such trees yielded by these methods would be more in line with established phylogenetic trees for Romance. This application demonstrates a method for isolating and identifying distinguishing aspects of variability in acoustic functional data which may be of evolutionary interest. It shows that it is possible to identify prominent features which render particular components effective for distinguishing the





**Fig. 8.** Cross-sections from the frequency and pseudo-time perspectives of the interpolated second component Hadamard matrices. For the frequency plots, the points with power of opposite sign (for instance points denoted 3500–4500 Hz) indicate the ranges which are separating Italian and Portuguese in the second co-ordinate. Some of possible frequency regions of interest are bracketed in the figure. The time points with higher opposite powers are between time units 1–20 and 70–100. Thus it is the non-central time points which are useful for distinguishing Italian and Portuguese in the second co-ordinate.



**Fig. 9.** Topology of UPGMA generated tree for the second component.

language groups. However, it also highlights the challenge of precise physical interpretation of particular components, a task which appears notably more complex due to having both a time and a frequency dimension. It would be of interest to express these differences back in the sound domain, although given the difficulties in inverting spectrograms to sound, this is not a trivial task. However, it is the subject of ongoing work, including experiments with other parametric acoustic representations that are more easily inverted.

## 6. Discussion

This paper presents a method for assessing tree-amenability of Gaussian functional data via Gaussian tree constraints. Through simulation the capability of the positivity constraint and its level of robustness to sample size have been demonstrated. Application of this fundamental tree constraint has indicated that the implementation is practically feasible. Moreover, it has shown that meaningful (albeit high-level) interpretations can be made making this particularly useful as an exploratory data tool.

The application to a linguistic data set comprising recordings from Romance language speakers provides several interesting preliminary results. First, it suggests that at least four components of the new basis are tree-amenable, and hence, the linguistic features these particular components represent may have a tree structure. Second, examining projections of the two tree-amenable components with highest explanatory power indicates that the second component is sufficient to distinguish languages, and including the first dimension separates dialects. Finally, closer study of the first two tree components indicates that broadly it is the beginning and end of utterances that are identified as effective language discriminants. Whilst these findings are somewhat speculative and the scope of the analysis is limited to a small set of carefully selected words, it does provide a starting point for describing interesting tree and non-tree like features of grouped data which may otherwise be obscured. Importantly, the methodology has been demonstrated as a proof of concept for extension to larger data sets. Furthermore, if applied on a larger scale to languages lacking established historical pathways, these techniques may offer fresh indications of the plausibility of different hypotheses concerning their historical development.

In the process of preparing data for tree constraint diagnostics, the concept of decomposable covariance structure has been combined with CVA to produce separable-CVA. Although not central to the aim of the paper, it is worth underlining

its usefulness. As with standard CVA it identifies modes of variability which effectively distinguish grouped data, but furthermore, in many circumstances it overcomes the common functional data problem of variable dimension exceeding sample size. This makes separable-CVA an appealing tool to use in practice and one which worked well for separating these languages (Fig. 4).

Whilst the range of techniques applied in this study can be carried across to other functional data sets, it is important to consider the assumptions which underpin the methods employed. Although no data set will truly obey all of the underlying assumptions, it is important nonetheless to understand the purpose of and the consequences of not satisfying each.

The theory of CVA assumes that the data are Gaussian and that the within-language covariances are equal. Considering the acoustic data, Gaussianity is unlikely to strictly hold. However, it may be sufficiently close to doing so and the first and second order moments may still answer questions about language relationships. If there is information available that the Gaussianity assumption is violated (as may well be evident in larger studies than the one presented here), then applying a copula transformation (see [28]) will provide us with transformed variables that are marginally Gaussian (a necessary condition for joint Gaussianity). Rerunning the analysis after applying univariate copula transforms, the data set passes the Royston  $H$  test for multivariate normality at the 1% level [63]. In terms of tree-amenability based on the positivity constraint, the first two components are unchanged and the third and fourth differ by just one violation. Some of the lower weighted components do differ more and so this suggests that the more dominant components may be more robust to slight transformations. The use of such copula techniques broadens the scope of these diagnostics. For more information see the supplementary material (see [Appendix A](#)).

Although for this study the within-language covariances were not significantly different this may not always be the case. In such instances, pooling the covariances will nevertheless emphasize shared commonality of covariances, and pooling may produce a more stable estimate than otherwise would be possible. Importantly, in circumstances where either assumption is violated CVA still produces a valid basis, the downside being it will not necessarily be as efficient at capturing the variability thus hindering the dimension reduction.

Covariance separability may also be unrealistic as there is often some correlation between the two separated dimensions (e.g., frequency and time), but in practice this is not a problem. CVA is an optimality tool and thus deviations from the separability assumption only decrease the efficacy of the basis obtained. That is, if covariance separability holds perfectly CVA determines an ordered basis which linearly maximizes between- to within-language variation as intended. Otherwise the ordered basis is suboptimal, the consequence being more dimensions may be required to capture a suitable proportion of the variance.

Therefore, although the three stated assumptions are not guaranteed to hold, infringement of any of them may impair performance but is unlikely to completely negate such an exploratory analysis. This permits a greater range of data sets for which a similar study could be applicable.

Aside from improving performance by selecting a data set that better meets the assumptions outlined above, the application to the Romance languages could instead be improved through a larger data set with greater breadth of words, and indeed a greater number and diversity of speakers. Currently, the between-covariances of the projected data are only calculated using the means of the ten unique word observations. Furthermore, some languages have more recordings than others, meaning that the weights differ by group. We reran the analysis using a flat weighting structure, but this resulted in little to no difference in the inferred results (data not shown). By increasing the number of recordings for all groups, whilst also including a wider variety of words, the outcomes of the analyses are likely to be more robust, with the caveat that inclusion of new words must be carefully considered to ensure they are linguistically suitable across all languages.

In this paper only some of the conditions were checked that are necessary for consistency with a phylogenetic tree. However, it is possible to enhance the power of such diagnostics to include all such checks of possible violations of phylogenetic tree constraints [68]. One method is to utilize the 4-point conditions based on the metric properties noted in Buneman [13]. If these conditions are satisfied then the tree structure can be uniquely identified. For one particular constraint based on tetrads, the work of Bollen and Ting [9] gives a basis for possible test statistics and Drton et al. [18] provides some useful distributional results for particular moments. An exploratory approach could use an extension of the binary graphical tree diagnostics described in Shiers and Smith [67] to the Gaussian domain and is currently under investigation. Of course to fully and formally evaluate these methods it will be necessary to evaluate the probabilistic properties of our methods. A Bayesian approach is currently under investigation where replicates are generated from Wishart covariance matrices, with properties derived from the projected data. These then provide posterior distributions across a graphical model space [6]. For graphs with high posterior probability, trees with hidden variables could then be induced. As with any formal testing framework, it will be important to consider potential effects of multiple comparisons, and to adjust for them appropriately (see [36, Chapter 12]).

Each of these methods could enrich the analysis, as the proposed trees could be compared to existing knowledge about the relationships between the languages. Contrasting results could even suggest new directions of research using traditional linguistic methods.

There is much potential for applications in fields beyond linguistics to make use of additional algebraic and semi-algebraic constraints. However, even existing constraints are currently underutilized. For example, despite the relatively simple derivation of the fundamental Gaussian tree constraint, this study marks its first use as a diagnostic for tree-amenability. The inclusion of further tree constraints, both in univariate and multivariate settings could provide further insight into

phylogenetic relationships in applications. However, derivation and application of such constraints is non-trivial and is the subject of ongoing work.

## Acknowledgments

NS acknowledges the support of Economics and Social Science Research Council grant ES/I90427/1. JADA acknowledges the support of UK Engineering and Physical Sciences Research Council grant EP/K021672/2. JSC acknowledges the support of UK Arts and Humanities Research Council grant AH/M002993/1. The authors wish to thank Piotr Zwiernik for very helpful discussions. They are also very much indebted to Pantelis Hadjipantelis for preprocessing the Romance language data set, and to Davide Pigoli for helpful discussions. The authors are also indebted to the Editor, the AE and two anonymous referees for their helpful comments, which led to considerable improvements of the paper.

## Appendix A. Supplementary material

Supplementary material related to this article can be found online at <http://dx.doi.org/10.1016/j.jmva.2016.09.015>.

## References

- [1] E.S. Allman, C. Matias, J.A. Rhodes, Identifiability of parameters in latent structure models with many observed variables, *Ann. Statist.* 37 (2009) 3099–3132.
- [2] E.S. Allman, J.A. Rhodes, Phylogenetic ideals and varieties for the general Markov model, *Adv. in Appl. Math.* 40 (2008) 127–148.
- [3] E.S. Allman, J.A. Rhodes, A. Taylor, A semialgebraic description of the general Markov model on phylogenetic trees, *SIAM J. Discrete Math.* 28 (2014) 736–755.
- [4] J.A.D. Aston, J.-M. Chiou, J.P. Evans, Linguistic pitch analysis using functional principal component mixed effect models, *J. Roy. Statist. Soc. Ser. C* 59 (2010) 297–317.
- [5] J.A.D. Aston, C. Kirch, Evaluating stationarity via change-point alternatives with applications to fMRI data, *Ann. Appl. Stat.* 6 (2012) 1906–1948.
- [6] A. Atay-Kayis, H. Massam, A Monte Carlo method for computing the marginal likelihood in nondecomposable Gaussian graphical models, *Biometrika* 92 (2005) 317–335.
- [7] P.C. Besse, H. Cardot, D.B. Stephenson, Autoregressive forecasting of some functional climatic variations, *Scand. J. Stat.* 27 (2000) 673–687.
- [8] K.A. Bollen, *Structural Equations with Latent Variables*, Wiley, New York, 1989.
- [9] K.A. Bollen, K.-F. Ting, Confirmatory tetrad analysis, *Sociol. Methodol.* 23 (1993) 147–175.
- [10] D.A.W. Bolnick, B.A.S. Shook, L. Campbell, I. Goddard, Problematic use of Greenberg's linguistic classification of the Americas in studies of Native American genetic variation, *Am. J. Hum. Genet.* 75 (2004) 519.
- [11] A. Bouchard-Côté, D. Hall, T.L. Griffiths, D. Klein, Automated reconstruction of ancient languages using probabilistic models of sound change, *Proc. Natl. Acad. Sci.* 110 (2013) 4224–4229.
- [12] G.E.P. Box, A general distribution theory for a class of likelihood criteria, *Biometrika* 36 (1949) 317–346.
- [13] O.P. Buneman, The recovery of trees from measures of dissimilarity, in: D.G. Kendall, P. Tautu (Eds.), *Mathematics in the Archaeological and Historical Sciences*, Edinburgh University Press, 1971, pp. 387–395.
- [14] P. Castro, W. Lawton, E. Sylvestre, Principal modes of variation for processes with continuous sample curves, *Technometrics* 28 (1986) 329–337.
- [15] K. Chen, P. Delicado, H.-G. Müller, Modelling function-valued stochastic processes, with applications to fertility dynamics, *J. R. Stat. Soc. Ser. B Stat. Methodol.* (ISSN: 1467-9868) (2016) <http://dx.doi.org/10.1111/rssb.12160>, in press.
- [16] T.F. Cox, M.A. Cox, *Multidimensional Scaling*, CRC Press, London, 2010.
- [17] S. Davis, P. Mermelstein, Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences, *IEEE Trans. Acoust. Speech Signal Process.* 28 (1980) 357–366.
- [18] M. Drton, H. Massam, I. Olkin, Moments of minors of Wishart matrices, *Ann. Statist.* 36 (2008) 2261–2283.
- [19] M. Drton, S. Sullivant, Algebraic statistical models, *Statist. Sinica* 17 (2007) 1273–1297.
- [20] M. Dunn, A. Terrill, G. Reesink, R.A. Foley, S.C. Levinson, Structural phylogenetics and the reconstruction of ancient language history, *Science* 309 (2005) 2072–2075.
- [21] J. Dutkowski, J. Tiurny, Identification of functional modules from conserved ancestral protein–protein interactions, *Bioinformatics* 23 (2007) 149–158.
- [22] J. Felsenstein, The number of evolutionary trees, *Syst. Biol.* 27 (1978) 27–33.
- [23] J. Felsenstein, Statistical inference of phylogenies, *J. R. Stat. Soc. Ser. A* 146 (1983) 246–272.
- [24] G. Fervaha, G. Remington, Interpreting a multivariate analysis of functional neuroimaging data, *Front. Psychiatry* 3 (2012) 1–52.
- [25] R.A. Fisher, The use of multiple measurements in taxonomic problems, *Ann. Eugenics* 7 (1936) 179–188.
- [26] P. Forster, A. Toth, Toward a phylogenetic chronology of ancient Gaulish, Celtic, and Indo-European, *Proc. Natl. Acad. Sci.* 100 (2003) 9079–9084.
- [27] S. Fulop, *Speech Spectrum Analysis*, Springer, Berlin, 2011.
- [28] C. Genest, J. Nešlehová, Copulas and copula models, in: A.H. El-Shaarawi, W.W. Piegorsch (Eds.), *Encyclopedia of Environmetrics*, Vol. 2, second ed., Wiley, Chichester, 2012, pp. 541–553.
- [29] E. Grabe, G. Kochanski, J. Coleman, Connecting intonation labels to mathematical descriptions of fundamental frequency, *Lang. Speech* 50 (2007) 281–310.
- [30] J.H. Greenberg, A quantitative approach to the morphological typology of language, in: R.F. Spencer (Ed.), *Method and Perspective in Anthropology: Papers in Honor of Wilson D. Wallis*, University of Minnesota Press, 1954, pp. 192–220.
- [31] P.Z. Hadjipantelis, J.A.D. Aston, J.P. Evans, Characterizing fundamental frequency in Mandarin: A functional principal component approach utilizing mixed effect models, *J. Acoust. Soc. Am.* 131 (2012) 4651–4664.
- [32] M. Harris, N. Vincent, *The Romance Languages*, Taylor & Francis, London, 1988.
- [33] H.H. Hock, *Principles of Historical Linguistics*, Walter de Gruyter, Berlin, 1991.
- [34] S.H. Holan, C.K. Wikle, L.E. Sullivan-Beckers, R.B. Croft, Modeling complex phenotypes: Generalized linear models using spectrogram predictors of animal communication signals, *Biometrics* 66 (2010) 914–924.
- [35] L. Horváth, P. Kokoszka, *Inference for Functional Data with Applications*, Springer, New York, 2012.
- [36] D. Howell, *Statistical Methods for Psychology*, Cengage Learning, 2012.
- [37] V. Istratescu, *Inner Product Structures: Theory and Applications*, Springer, Berlin, 1987.
- [38] I.T. Jolliffe, *Principal Component Analysis*, second ed., Springer, New York, 2002.
- [39] H.T. Kiiveri, Canonical variate analysis of high-dimensional spectral data, *Technometrics* 34 (1992) 321–331.
- [40] L.L. Koenig, J.C. Lucero, E. Perlman, Speech production variability in fricatives of children and adults: Results of functional data analysis, *J. Acoust. Soc. Am.* 124 (2008) 3158–3170.

- [41] W.J. Krzanowski, Principles of Multivariate Analysis, in: Oxford Statistical Science Series, vol. 3, The Clarendon Press Oxford University Press, New York, 1990, A user's perspective, Corrected reprint of the 1988 edition, Oxford Science Publications.
- [42] P. Lancaster, M. Tismenetsky, The Theory of Matrices, second ed., Academic Press Inc., Orlando, FL, 1985.
- [43] S.L. Lauritzen, Graphical Models, Oxford University Press, 1996.
- [44] C. Lee, S. Narayanan, R. Pieraccini, Recognition of negative emotions from the speech signal, in: IEEE Workshop on Automatic Speech Recognition and Understanding, 2001. ASRU '01, 2001, pp. 240–243.
- [45] P.-L. Loh, M.J. Wainwright, Structure estimation for discrete graphical models: Generalized covariance matrices and their inverses, *Ann. Statist.* 41 (2013) 3022–3049.
- [46] C. Long, E. Brown, C. Triantafyllou, I. Aharon, L. Wald, V. Solo, Nonstationary noise estimation in functional MRI, *NeuroImage* 28 (2005) 890–903.
- [47] J.C. Lucero, L.L. Koenig, Time normalization of voice signals using functional data analysis, *J. Acoust. Soc. Am.* 108 (2000) 1408–1420.
- [48] J.G. Martinez, K.M. Bohn, R.J. Carroll, J.S. Morris, A study of Mexican free-tailed bat chirp syllables: Bayesian functional mixed models for nonstationary acoustic time series, *J. Amer. Statist. Assoc.* 108 (2013) 514–526.
- [49] C.D. Michener, R.R. Sokal, A quantitative approach to a problem in classification, *Evolution* 11 (1957) 130–162.
- [50] M. Mizuta, Graphical representation of functional clusters and MDS configurations, in: *Data Analysis, Classification and the Forward Search*, Springer, New York, 2006, pp. 31–37.
- [51] C. Mooshammer, Acoustic and laryngographic measures of the laryngeal reflexes of linguistic prominence and vocal effort in German, *J. Acoust. Soc. Am.* 127 (2010) 1047–1058.
- [52] L. Nakhleh, D.A. Ringe, T. Warnow, Perfect phylogenetic networks: a new methodology for reconstructing the evolutionary history of natural languages, *Language* 81 (2005) 382–420.
- [53] S. Nelson-Sathi, J.-M. List, H. Geisler, H. Fangerau, R.D. Gray, W. Martin, T. Dagan, Networks uncover hidden lexical borrowing in Indo-European language evolution, *Proc. R. Soc. B: Biol. Sci.* 278 (2011) 1794–1803.
- [54] G.K. Nicholls, R.J. Ryder, Phylogenetic models for semitic vocabulary, in: *Proceedings of the 26th International Workshop on Statistical Modelling*, 2011.
- [55] S. Nittrouer, R.S. McGowan, P.H. Milenkovic, D. Beehler, Acoustic measurements of men's and women's voices: A study of context effects and covariation, *J. Speech Hear. Res.* 33 (1990) 761–775.
- [56] E. Parris, M. Carey, Language independent gender identification, in: *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 1996. ICASSP-96 Conference Proceedings, 1996, Vol. 2, 1996, pp. 685–688.
- [57] E. Pépiot, Voice, speech and gender: Male–female acoustic differences and cross-language variation in English and French speakers. *Actes des Rencontres Jeunes Chercheurs de l'ED 268* 2011–2012, 2013.
- [58] D. Pigoli, P.Z. Hadjipantelis, J.S. Coleman, J.A.D. Aston, The analysis of acoustic phonetic data: Exploring differences in the spoken Romance languages, 2015. [arXiv:1507.07587](https://arxiv.org/abs/1507.07587).
- [59] G. Price, Romance, in: J. Gvozdanović (Ed.), *Indo-European Numerals*, Mathematical Research, Berlin, 1992, (Chapter 13).
- [60] J.O. Ramsay, B.W. Silverman, *Functional Data Analysis*, second ed., Springer, New York, 2005.
- [61] S.J. Ratcliffe, L.R. Leader, G.Z. Heller, Functional data analysis with application to periodically stimulated foetal heart rate data. I: Functional regression, *Stat. Med.* 21 (2002) 1103–1114.
- [62] J. Rice, B. Silverman, Estimating the mean and covariance structure nonparametrically when the data are curves, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 53 (1991) 233–243.
- [63] J. Royston, Some techniques for assessing multivariate normality based on the Shapiro–Wilk W, *J. Appl. Stat.* 32 (1983) 121–133.
- [64] A. Schleicher, *Die deutsche Sprache*, Cotta'scher Verlag, Stuttgart, 1860.
- [65] R. Settini, J.Q. Smith, Geometry, moments and conditional independence trees with hidden variables, *Ann. Statist.* 28 (2000) 1179–1205.
- [66] N. Shiers, *Gaussian Latent Tree Model Constraints for Linguistics and Other Applications* (Doctoral dissertation), University of Warwick, 2016.
- [67] N. Shiers, J.Q. Smith, Graphical inequality diagnostics for phylogenetic trees, in: *Proceedings of 6th European Workshop on Probabilistic Graphical Models*, Granada, Spain, 19–21 Sep 2012, 2012, pp. 291–298.
- [68] N. Shiers, P. Zwiernik, J.A.D. Aston, J.Q. Smith, The correlation space of Gaussian latent tree models and model selection without fitting, *Biometrika* 103 (3) (2016) 531–545. <http://dx.doi.org/10.1093/biomet/asw032>.
- [69] H. Sørensen, J. Goldsmith, L.M. Sangalli, An introduction with medical applications to functional data analysis, *Stat. Med.* 32 (2013) 5222–5240.
- [70] B. Sturmfels, S. Sullivant, Toric ideals of phylogenetic invariants, *J. Comput. Biol.* 12 (2005) 204–228.
- [71] P. Tang, H. Müller, Pairwise curve synchronization for high-dimensional data, *Biometrika* 95 (2008) 875–889.
- [72] H. Wang, J. Marron, Object oriented data analysis: Sets of trees, *Ann. Statist.* 35 (2007) 1849–1873.
- [73] L. Wenyin, X. Quan, M. Feng, B. Qiu, A short text modeling method combining semantic and statistical information, *Inform. Sci.* 180 (2010) 4031–4041.
- [74] F. Yao, T. Lee, Penalized spline models for functional principal component analysis, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 68 (2006) 3–25.
- [75] P. Zwiernik, J.Q. Smith, Implicit inequality constraints in a binary tree model, *Electron. J. Stat.* 5 (2011) 1276–1312.
- [76] P. Zwiernik, J.Q. Smith, Tree cumulants and the geometry of binary tree models, *Bernoulli* 18 (2012) 290–321.