



Cite this article: Hadjipantelis PZ, Jones NS, Moriarty J, Springate DA, Knight CG. 2013 Function-valued traits in evolution. *J R Soc Interface* 10: 20121032.
<http://dx.doi.org/10.1098/rsif.2012.1032>

Received: 15 December 2012

Accepted: 28 January 2013

Subject Areas:

biomathematics, bioinformatics

Keywords:

comparative analysis, Ornstein–Uhlenbeck process, non-parametric Bayesian inference, functional phylogenetics, ancestral reconstruction, functional Gaussian process regression

Author for correspondence:

Pantelis Z. Hadjipantelis

e-mail: p.z.hadjipantelis@warwick.ac.uk

Electronic supplementary material is available at <http://dx.doi.org/10.1098/rsif.2012.1032> or via <http://rsif.royalsocietypublishing.org>.

Function-valued traits in evolution

Pantelis Z. Hadjipantelis¹, Nick S. Jones², John Moriarty³, David A. Springate⁴ and Christopher G. Knight⁴

¹Centre for Complexity Science and Department of Statistics, University of Warwick, Coventry CV4 7AL, UK

²Department of Mathematics, Imperial College London, London SW7 2AZ, UK

³School of Mathematics, and ⁴Faculty of Life Sciences, University of Manchester, Oxford Road, Manchester M13 9PL, UK

Many biological characteristics of evolutionary interest are not scalar variables but continuous functions. Given a dataset of function-valued traits generated by evolution, we develop a practical, statistical approach to infer ancestral function-valued traits, and estimate the generative evolutionary process. We do this by combining dimension reduction and phylogenetic Gaussian process regression, a non-parametric procedure that explicitly accounts for known phylogenetic relationships. We test the performance of methods on simulated, function-valued data generated from a stochastic evolutionary model. The methods are applied assuming that only the phylogeny, and the function-valued traits of taxa at its tips are known. Our method is robust and applicable to a wide range of function-valued data, and also offers a phylogenetically aware method for estimating the autocorrelation of function-valued traits.

1. Introduction

The number, reliability and coverage of evolutionary trees are growing rapidly [1,2]. However, knowing organisms' evolutionary relationships through phylogenetics is only one step in understanding the evolution of their characteristics [3]. Three issues are particularly challenging. The first is limited information: empirical information is typically available only for extant taxa, represented by tips of a phylogenetic tree, whereas evolutionary questions frequently concern unobserved ancestors deeper in the tree. The second is dependence: the available information for different organisms in a phylogeny is independent because a phylogeny describes a complex pattern of non-independence; observed variation is a mixture of this inherited variation and specific variation [4]. The third is high dimensionality: the emerging literature on function-valued traits [5–7] recognizes that many characteristics of living organisms are best represented as a continuous function rather than a single factor or a small number of correlated factors. Such characteristics include growth or mortality curves [8], reaction norms [9] and distributions [10], where the increasing ease of genome sequencing has greatly expanded the range of species in which distributions of gene [11] or predicted protein [12] properties are available. Therefore, a function-valued trait is defined as a phenotypic trait that can be represented by a continuous mathematical function [9].

Previous work [13] proposed an evolutionary model for function-valued data d related by a phylogeny T . The data are regarded as observations of a phylogenetic Gaussian process (PGP) at the tips of T . That work shows that a PGP can be expressed as a stochastic linear operator X on a fixed set ϕ of basis functions (independent components of variation), so that

$$d = X^T \phi. \quad (1.1)$$

However, the study does not address the linear inverse problem of obtaining estimates $\hat{\phi}$ and \hat{X} of ϕ and X : our first contribution in this paper is to provide an approach to this problem in §2.2 via independent principal component analysis (IPCA; [14]).

We refer to X as the *mixing matrix*, and to the (i, j) th entry of X as the *mixing coefficient* of the i th basis function at the j th taxon. It is these mixing coefficients that we model as evolving. For each fixed value of i , the X_{ij} are correlated (owing to phylogeny) as j varies over the taxa; the basis functions themselves do not evolve in our model.

In §2.3, we address the problem of estimating the statistical structure of the mixing coefficients by performing phylogenetic Gaussian process regression (PGPR) on each of the rows of \hat{X} separately. This corresponds to assuming independence between the rows (i.e. that the coefficients of the different basis functions evolve independently). It is commonly argued in the quantitative genetics literature [15] that evolutionary processes can be modelled as Ornstein–Uhlenbeck (OU) processes. Under these assumptions, the estimation of the forward operator reduces to the estimation of a small vector γ of parameters [13]. In §2.1, we clarify the interpretation of these parameters in evolutionary contexts. The explicit PGPR posterior likelihood function is then used to obtain maximum-likelihood estimates (MLEs) for γ . The estimation of γ is known to be a challenging statistical problem [16]. We suggest an approach based on the principle of *bagging* [17] in §2.4.

Our final contribution (§2.5) addresses the problem of estimating the function-valued traits of ancestral taxa. The earlier-mentioned PGPR step also returns a posterior distribution for the mixing coefficient of each basis function at each ancestral taxon in the phylogeny. At any particular ancestor, the estimated basis functions may be combined statistically, using the posterior distributions of their respective mixing coefficients, to provide a posterior distribution for the function-valued trait. Because the univariate posterior distributions are Gaussian, and the mixing is linear, the posterior for the function-valued trait has a closed form representation as a GP (equation (2.6)) that provides a major analytical and computational advantage for the approach. We can verify the methods proposed by using a PGP as a stochastic generative model. This simulates correlated function-valued traits across the taxa of T . Given only the phylogeny and the function-valued traits of taxa at its tips, our estimates for $\hat{\phi}$ and the ancestral functions are then compared with the simulation.

Overall, our three methods (in §§2.2, 2.4, 2.5) appropriately combine developments in functional data analysis with the evolutionary dynamics of quantitative phenotypic traits, allowing non-parametric Bayesian inference from phylogenetically correlated function-valued traits. An outline of the framework presented in this study can be found in figure 1.

2. Methods and implementation

2.1. Artificial evolution of function-valued traits

We begin by generating a random phylogenetic tree T with 128 tips, shown in figure 2. This fixes the experimental design for our simulation and inference, but further simulations given in the electronic supplementary material confirm that the statistical performance of our methods is consistent across a range of choices for T . Branch length distributions are surprisingly consistent across organisms [18]; branch lengths were drawn from the empirical branch length distribution (see electronic supplementary material, section S1) extracted from TREEFAM v. 8.0 [2].

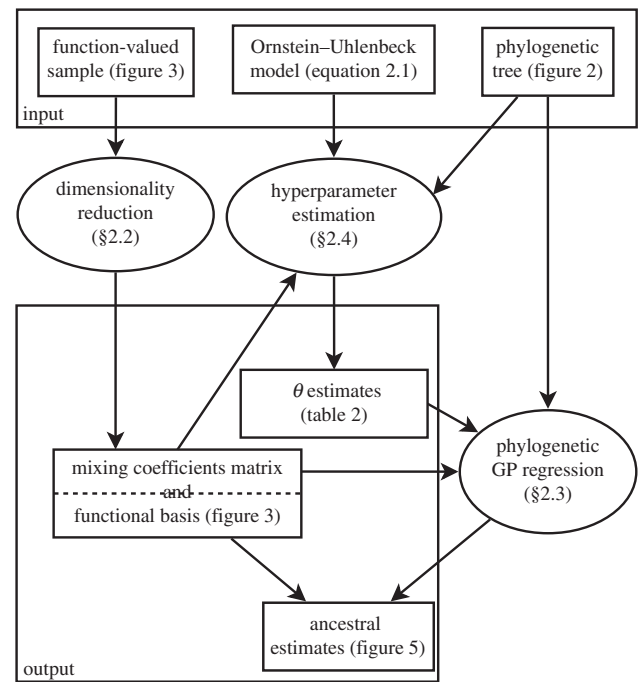


Figure 1. The three methods presented in this paper (ovals) and their interrelationships.

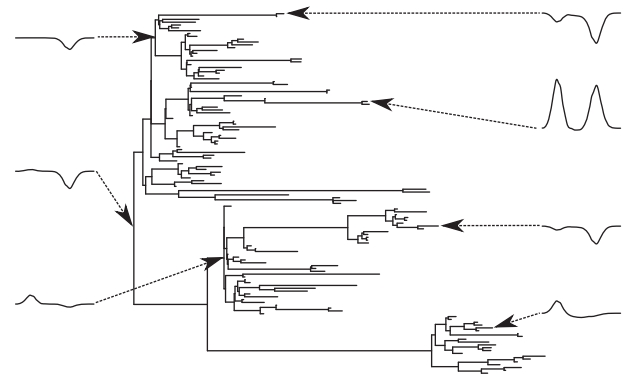


Figure 2. The random phylogenetic tree used and examples of the function-valued traits shown at the tips (extant taxa) and the internal nodes (ancestral taxa). A subset of these is used in figure 5.

Second, we chose a basis ϕ in equation (1.1). We have no reason *a priori* to suppose that this basis is orthogonal and, in general, there is no reason for our inference procedure to be sensitive to the particular shape of the basis functions. The three simple non-orthogonal, unimodal functions shown in figure 3 were therefore chosen as examples. For computational purposes, each basis function was stored numerically as a vector of length 1024, so that the basis matrix ϕ was of size 3×1024 and its i th row stored the i th basis function.

Third, different mixing coefficients were generated by a phylogenetic OU process for each basis function and stored in the respective row of X . Our modelling assumption is that the mixing coefficients for distinct basis functions ϕ_1 , ϕ_2 , ϕ_3 are statistically independent of each other: in equation (1.1), this means that the rows of X are independent. It is therefore sufficient to describe the stochastic process generating $X_{i\cdot}$, the i th row of X with $i \in \{1, 2, 3\}$. We calculated the mixing matrix at the 128 tip taxa so X is of size 3×128 . The ‘true’ ancestral values were established by generating phylogenetic OU processes over the whole phylogeny. The values of this process at tip taxa were stored in a row

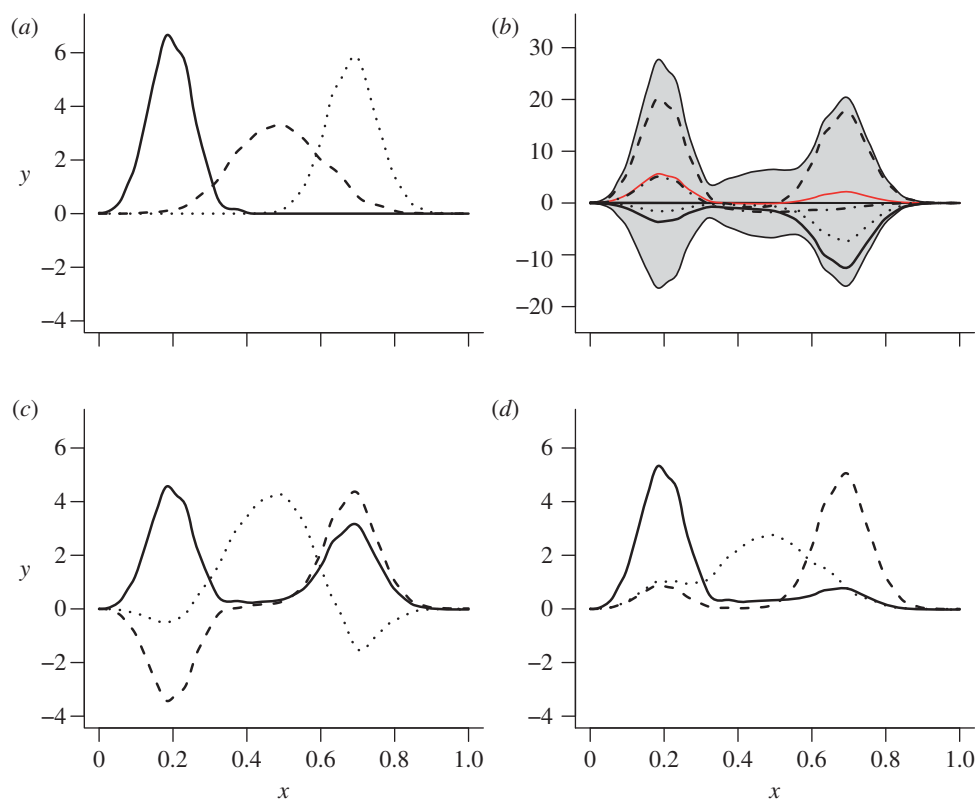


Figure 3. (a) original basis signals, ϕ ; (b) mixed sample at the tips, d (four individual function-valued traits are shown; red line and grey band show, respectively, the mean and 2 s.d. for all 128 function-valued data at the tips); (d) IPCA basis, $\hat{\phi}$; (c) PCA basis. (Online version in colour.)

vector \bar{X}_i (\bar{X}_i is a simulation of the tip taxa mixing coefficients X_i excluding the non-phylogenetic variation), and its values at internal taxa were stored in a row vector W_i for performance analyses in §2.5. To simulate the additional effect of non-phylogenetic variation (e.g. due to measurement error or environmental effects), independent (i.e. non-phylogenetic) variation was added to each entry of \bar{X}_i :

$$\bar{X}_i = \bar{X}_i + \epsilon_i,$$

where ϵ_i is a 1×128 vector of independent Gaussian errors with mean 0 and standard deviation σ_n^i , and, finally, the matrix multiplication in equation (1.1) was performed to obtain the simulated data d . The ‘extant’ function-valued trait at tip taxon j is thus $\sum_{i=1}^3 X_{ij} \phi_i$ (a vector of length 1024), whereas the ancestral function-valued trait at internal taxon g is $\sum_{i=1}^3 W_{ig} \phi_i$. The ancestral function-valued traits therefore exhibit only the phylogenetic part of simulated variation, whereas the extant function-valued traits exhibit both phylogenetic and non-phylogenetic variation. Of course, it is not possible to reconstruct non-phylogenetic variation using phylogenetic methods: we simulate non-phylogenetic variation only to demonstrate that it does not prevent the reconstruction of the phylogenetic part of variation for ancestral taxa in §§2.2–2.5.

We now comment on the specific parameters chosen for the phylogenetic OU processes mentioned earlier. As in Hansen [19], we refer to the *strength of selection parameter* α and the *random genetic drift* σ : we add superscripts to these parameters to distinguish between the three different OU processes. With this notation, the mixing coefficients for the row X_i have the following covariance function:

$$\begin{aligned} K_T^i(\mathbf{t}_1, \mathbf{t}_2) &= E[X_{ij} X_{ig}] \\ &= (\sigma_f^i)^2 \exp(-2\alpha^i D_T(\mathbf{t}_j, \mathbf{t}_g)) + (\sigma_n^i)^2 \delta_{\mathbf{t}_j, \mathbf{t}_g}^e, \end{aligned} \quad (2.1)$$

Table 1. The fixed values used for the parameters in equation (2.1) to generate the mixing coefficients X_{ij} . Each row constitutes a value of γ^i . 6.17 and 2.06 correspond to 0.75 and 0.25 of the tree’s ℓ_{\max} , respectively. When $i = 2$, ℓ^i is not applicable, because there is no phylogenetic variation in the sample.

i	σ_f^i	ℓ^i	σ_n^i
1	2.5	6.17	0.5
2	0	n.a.	1.0
3	1.5	2.06	0.5

where $\sigma_f^i = \sqrt{(\sigma^i)^2 / 2\alpha^i}$, $D_T(\mathbf{t}_j, \mathbf{t}_g)$ denotes the phylogenetic or patristic distance (i.e. the distance in \mathbf{T}) between the j th and g th tip taxa, σ_n is defined as earlier, and

$$\delta_{\mathbf{t}_j, \mathbf{t}_g}^e = \begin{cases} 1, & \text{if } t_j = t_g \text{ and } t_j \text{ is a tip taxon,} \\ 0, & \text{otherwise} \end{cases}$$

adds non-phylogenetic variation to extant taxa as discussed earlier, i.e. δ^e evaluates to 1 only for extant taxa, thus σ_n quantifies within-species genetic or environmental effects and measurement error in the i th mixing coefficient. We see from equation (2.1) that the proportion of variation in the row X_i attributable to the phylogeny is $(\sigma_f^i)^2 / ((\sigma_f^i)^2 + (\sigma_n^i)^2)$.

In the Gaussian process regression (GPR) literature in machine learning, $1/2\alpha$ is equivalent to ℓ , the characteristic length-scale [20] of decay in the correlation function and in the following we work with the latter. For all of the OU processes, we used characteristic length-scales relative to 8.22, the maximum patristic distance (ℓ_{\max}) between two extant taxa for our simulated tree (figure 2). The values we used are given in table 1. In particular, $\sigma_f^i = 0$ when $i = 2$

and it follows that the characteristic length-scale ℓ plays no role for this OU process, and equally we do not define the strength-of-selection parameter α^i when $i = 2$.

2.2. Dimensionality reduction and source separation for function-valued traits

Given a dataset d of function-valued traits, we would like to find appropriate estimates \hat{X} and $\hat{\phi}$ of the mixing matrix X and the basis set ϕ , respectively. The first task is to identify a good linear subspace S of the space of all continuous functions by choosing basis functions appropriately. The purpose is to work, not with the function-valued data directly, but with their projections in S . We may say that the chosen subspace S is good, if the projected data approximate the original data well, whereas the number of basis functions is not unnecessarily large so that S has the ‘effective’ dimension of the data.

We then face a linear inverse problem: given the dataset d of function-valued traits, the task is to generate estimates \hat{X} and $\hat{\phi}$ (equation (1.1)). This task is also known as *source separation* [21], which has a variety of implementations making different assumptions about the basis ϕ and mixing coefficients X . One widely used approach is principal components analysis (PCA) [22], which returns orthogonal sets of basis functions to explain the greatest possible variation. PCA has been extended to take account of phylogenetic relationships [23], however, if a sample of functions is generated by mixing non-orthogonal basis functions, the PCs of the sample (whether or not they account for phylogeny) will not equal the basis curves, due to the assumption of orthogonality (figure 3). In the independent component analysis (ICA), the alternative assumption is made that the rows X_i of X are statistically independent. This assumption fits more naturally with our modelling assumptions, because we assume that the rows X_i are mutually independent [21]. ICA has proved fruitful in other biological applications [24] as has passing the results of PCA to ICA, which has been termed IPCA [14].

PCA is an appropriate tool for identifying the effective dimension of a high-dimensional dataset [25]. Therefore, to achieve both dimension reduction and source separation, we first applied PCA to the dataset d (the 128 function-valued traits at the tips of T) to determine the appropriate number of basis functions. The PCs were then passed to the *CubICA* implementation of ICA [26]. *CubICA* returned a new set of basis functions (figure 3*d*) that were taken as the estimated basis $\hat{\phi}$.

2.3. Phylogenetic Gaussian process regression

ICA also returns the estimated mixing coefficients at tip taxa, \hat{X} . Our next step was to perform PGPR [13] separately on each row \hat{X}_i , assuming knowledge of the phylogeny T , in order to obtain posterior distributions for all mixing coefficients throughout the tree T .

GPR [20] is a flexible Bayesian technique in which prior distributions are placed on continuous functions. Its range of priors includes the Brownian motion and OU processes, which are by far the most commonly used models of character evolution [15,27]. Its implementation is particularly straightforward, because the posterior distributions are also GPs and have closed forms. We now give a brief exposition of GPR, using notation standard in the machine learning literature [20].

A GP may be specified by its mean surface and its covariance function $K(\gamma)$, where γ is a vector of parameters. Because the components of γ parametrize the prior distribution, they are referred to as *hyperparameters*. The GP prior distribution is denoted

$$f \sim \mathcal{N}(0, K(\gamma)).$$

If x^* is a set of unobserved coordinates and x is a set of observed coordinates, then the posterior distribution of the vector $f(x^*)$ given the observations $f(x)$ is

$$f(x^*)|f(x) \sim \mathcal{N}(A, B), \quad (2.2)$$

where

$$A = K(x^*, x, \gamma)K(x, x, \gamma)^{-1}f(x), \quad (2.3)$$

and

$$B = K(x^*, x^*, \gamma) - K(x^*, x, \gamma)K(x, x, \gamma)^{-1}K(x^*, x, \gamma)^T \quad (2.4)$$

and $K(x^*, x, \gamma)$ denotes the $|x^*| \times |x|$ matrix of the covariance function K evaluated at all pairs $x_i^* \in X^*, x_j \in X$. Equations (2.3) and (2.4) convey that the posterior mean estimate will be a linear combination of the given data and that the posterior variance will be equal to the prior variance minus the amount that can be explained by the data. Additionally, the log-likelihood of the sample $f(x)$ is

$$\begin{aligned} \log p(f(x)|\gamma) &= -\frac{1}{2}f(x)^T K(x, x, \gamma)^{-1}f(x) \\ &\quad -\frac{1}{2}\log(\det(K(x, x, \gamma))) - \frac{|x|}{2}\log 2\pi. \end{aligned} \quad (2.5)$$

It can be seen from equation (2.5) that the MLE is subject both to the fit it delivers (the first term) and the model complexity (the second term). Thus, GPR is non-parametric in the sense that no assumption is made about the structure of the model: the more data gathered, the longer the vector $f(x)$, and the more intricate the posterior model for $f(x^*)$.

PGPR extends the applicability of GPR to evolved function-valued traits. A *PGP* is a GP indexed by a phylogeny T , where the function-valued traits at each pair of taxa are conditionally independent given the function-valued traits of their common ancestors. When the evolutionary process has the same covariance function along any branch of T beginning at its root (called the *marginal covariance function*), these assumptions are sufficient to uniquely specify the covariance function of the PGP, K_T . As we assume that T is known in our inverse problem, the only remaining modelling choice is therefore the marginal covariance function. As can be seen from equation (2.1), K is a function of patristic distances on the tree rather than Euclidean distances as standard in spatial GPR.

In comparative studies, where one has observations at the tips of T , the covariance function K_T may be used to construct a GP prior for the function-valued traits, allowing functional regression. In the model that we use, this is equivalent to specifying a Gaussian prior distribution for the mixing coefficients Y_{ij} and X_{ij} . This may be carried out by regarding the row vectors Y_i and X_i as observations of a univariate PGP. As noted in Jones & Moriarty [13], if we assume that the evolutionary process is Markovian and stationary, then the modelling choice vanishes, and the marginal covariance function is specified uniquely: it is the stationary OU covariance function. If we also add explicit modelling of non-phylogenetically related variation at the tip taxa, the

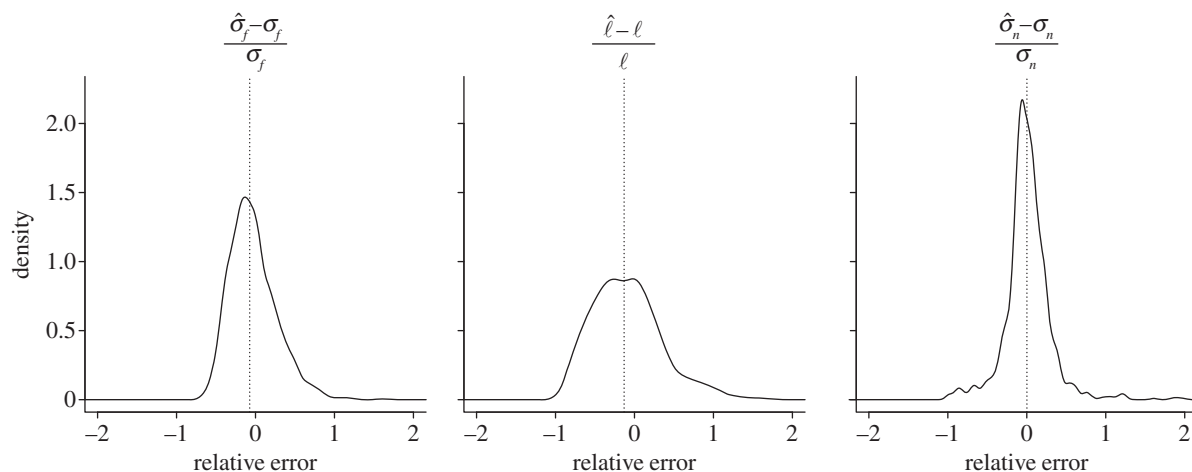


Figure 4. Kernel density estimates of the relative errors in 1024 runs of the γ estimation procedure, each time for a different tree, a different set of mixing coefficients and a different set of parameters in γ ; no components of γ are assumed to be known beforehand. Estimation results are commented on in §3. The median values shown by the dotted line are $(-0.073, -0.131$ and $0.001)$, respectively.

univariate prior covariance function has the unique functional form presented in equation (2.1). We do not assume knowledge of the parameters of equation (2.1), however, their estimation is the subject of §2.4.

2.4. Hyperparameter estimation

Because the posterior distributions returned by PGPR depend on the hyperparameter vector γ , we must estimate γ in order to reconstruct ancestral function-valued traits, and the estimation procedure should correct for the dependence owing to phylogeny. MLE of the phylogenetic variation, non-phylogenetic variation and characteristic length-scale hyperparameters σ_f^i , σ_n^i and ℓ^i , respectively, may be attempted numerically using the explicit prior likelihood function (equation (2.5)). Because estimating σ_f^i and ℓ^i alone is challenging [16] (although the estimation improves significantly with increased sample size), and we have further increased the challenge by introducing non-phylogenetic variation, we propose an improved estimation procedure using the machine learning technique *bagging* [17], which a member of the *boosting* framework [22]. We show that these estimates may be further improved if one knows the value of the ratio $(\sigma_f^i)^2/(\sigma_n^i)^2$, which is closely related to Pagel's λ [28].

Bagging (bootstrap aggregating) seeks to reduce the variance of an estimator by generating multiple estimates and averaging. It is simple to implement given an existing estimation procedure: one adds a loop front end that selects a bootstrap sample and sends it to the estimation procedure and a back end that aggregates the resulting estimates [17]. We generated 100 (sub)trees of 100 taxa by sampling without replacement our original 128 taxa tree, obtained the MLE for γ on each subtree, and averaged these estimates to obtain the aggregated estimate $\hat{\gamma}$. Our results are shown in table 2: for $i = 1$ and $i = 3$, given our moderate sample size (128 taxa), the accuracy of these results is at least in line with the state of the art [16] despite the additional challenge posed by non-phylogenetic variation. For $i = 2$, where phylogenetic variation is absent from the generative model ($\sigma_f^i = 0$), our estimation procedure indicates its absence by returning estimates for ℓ^i whose magnitude is unrealistically small for the examined tree (less than the first percentile of the tree's

Table 2. The bagging estimates for the hyperparameters in equation (2.1) (standard deviations of bagging estimates in parentheses). Each row corresponds to a given estimate of the vector $\hat{\gamma}^i$. These estimates provide the maximum-likelihood value for equation (2.5) and are comparable with the original ones from table 1.

i	$\hat{\sigma}_f^i$	$\hat{\ell}^i$	$\hat{\sigma}_n^i$
1	3.41 (0.62)	2.83 (0.47)	0.78 (0.47)
2	0.55 (0.33)	0.05 (0.02)	0.84 (0.34)
3	2.83 (0.33)	2.06 (0.50)	0.73 (0.29)

patristic distances). Commenting further on this matter, exceptionally *small* characteristic length-scales relative to the tree patristic distances, as seen here, practically suggest taxa-specific phylogenetic variation, i.e. non-phylogenetic variation. This holds also in its reverse: exceptionally *large* characteristic length-scales suggest a stable, non-decaying variation across the examined taxa that is indifferent to their patristic distances, again suggesting the absence of phylogenetic variance among the nodes.

To assess the robustness of this hyperparameter estimation method, we performed 1024 simulations, randomly regenerating the tree and parameter vector γ each time (see electronic supplementary material, section S2). The accuracy of these estimates is shown in figure 4. Improved results when the ratio $(\sigma_f^i)^2/(\sigma_n^i)^2$ is known *a priori* (e.g. through knowledge of Pagel's λ) are also given in the electronic supplementary material, sections S2 and S3. Our ultimate aim is ancestor reconstruction rather than hyperparameter estimation *per se*, and this is the subject of §2.5.

2.5. Ancestor reconstruction

Having generated function-valued data (§2.1), extracted mixing coefficients \hat{X} (§2.2) and performed hyperparameter estimation (§2.4), we may now perform PGPR (§2.3) on each row \hat{X}_i , to obtain the univariate Gaussian posterior distribution for the mixing coefficient W_{it^*} at any internal taxon t^* . As discussed in §2.3, the GP prior distribution has covariance function (equation (2.1)). We have assessed the accuracy

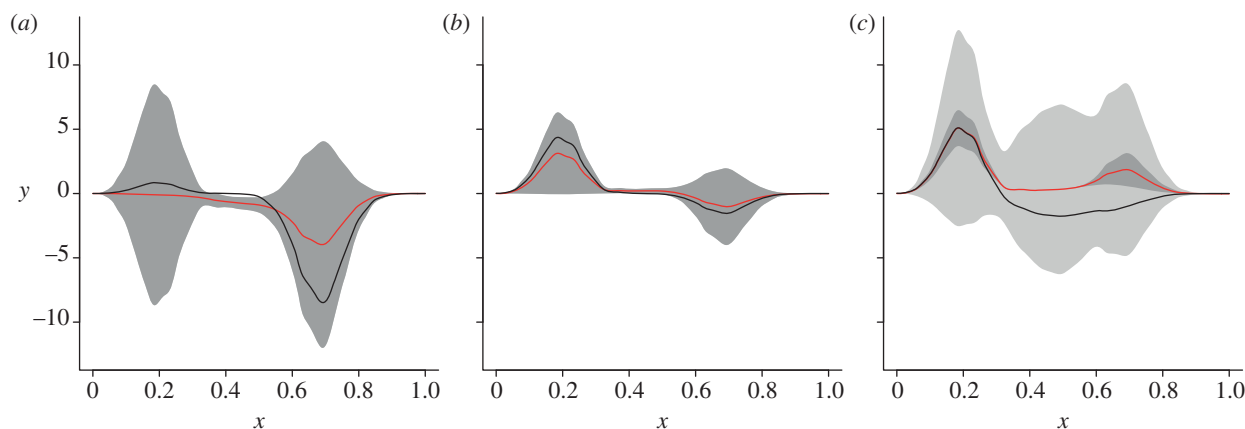


Figure 5. Posterior distributions at three points in the phylogeny using the estimated $\hat{\phi}$ and $\hat{\gamma}$. The prediction made by the regression analysis is shown via the posterior mean (red line), the component of posterior variance due to phylogenetic variation (2 s.d., dark grey band) and non-phylogenetic variation (2 s.d., light grey band). The black line shows the simulated data enabling visual validation of the ancestral predictions. (c), the black line is the training data at a tip taxon the red line and dark grey band represent the posterior distribution of its phylogenetic component, whereas the light grey band represents the estimated magnitude of non-phylogenetic variation. The root and internal taxon here are the same as those indicated in figures 2 and 3, and the tip is the second from bottom on the same figure. (a) root estimate; (b) node estimate; and (c) tip estimate. (Online version in colour.)

of our bagging estimate $\hat{\gamma}$ in §2.4 and we now substitute $\hat{\gamma}^i$ into equation (2.1). Taking a simple and direct approach, our estimate $\hat{\phi}$ obtained in §2.2 may then be substituted into equation (1.1) to obtain the function-valued posterior distribution $f_{\mathbf{t}^*}$ for the function-valued trait at taxon \mathbf{t}^* . Because our estimated basis functions are stored numerically as vectors of length 1024, this gives the same discretion for the ancestral traits.

Conditioning on our estimated mixing coefficients \hat{X}_i for the tip taxa, the posterior distribution of $W_{i\mathbf{t}^*}$ is

$$W_{i\mathbf{t}^*} \sim \mathcal{N}(\hat{A}_i, \hat{B}_i),$$

where the vector \hat{A}_i and matrix \hat{B}_i are obtained from equations (2.3) and (2.4), taking $x = \hat{X}_i$, $x^* = W_{i\mathbf{t}^*}$ and $\hat{\gamma}^i$, respectively, for our observation coordinates, estimation coordinates and hyperparameter vector. Because our prior assumption is that the rows of X are statistically independent of each other, it follows from equation (1.1) that

$$f_{\mathbf{t}^*} \sim \mathcal{N}\left(\sum_{i=1}^k \hat{A}_i \hat{\phi}_i, \sum_{i=1}^k \hat{\phi}_i^T \hat{B}_i \hat{\phi}_i\right). \quad (2.6)$$

The marginal distributions of this representation (mean and standard deviation) are shown in figure 5.

Figure 5 compares the function-valued estimates $\hat{f}_{\mathbf{t}^*}$ to the simulated function-valued traits at the root (figure 5a), an internal node (figure 5b), and at a tip (figure 5c). In figure 5a,b, the simulated function-valued data are shown in black, and can be seen typically to lie within two posterior standard deviations. In figure 5c, the black line is the observed function-valued trait at that tip: the red line and dark grey band represent the posterior distribution of its phylogenetic component, and the light grey band represents the estimated magnitude of the additional non-phylogenetic variation. Uncertainty over the phylogenetic part of variation (dark grey band) decreases from root to tip, as all observations are at the extant tip taxa. We note that the posterior distributions, even at the root, put clear statistical constraints on the phylogenetic part of ancestral function-valued data: in

this (admittedly simulated and highly controlled) setting, we can reason effectively about ancestral function-valued traits.

3. Discussion

In §2.1, we have appealed to equation (1.1) in the setting of mathematical inverse problems where, given data d , the challenge is to infer a forward operator G and model ϕ such that

$$d = G(\phi), \quad (3.1)$$

and such problems are typically under-determined and require additional modelling assumptions [29]. Given a phylogeny \mathbf{T} and function-valued data d at its tips, we wish to infer the forward operator $G_{\mathbf{T}}$ and model ϕ such that

$$d = G_{\mathbf{T}}(\phi). \quad (3.2)$$

When the data d are a small number of correlated factors per tip taxon, a variety of statistical approaches are available [30,31]. When the data are functions, the PGPs [13,32] have been proposed as the forward operator and this is the approach we have taken in this work.

Our dimensionality-reduction methodology in §2.2 can be easily varied or extended. For example, any suitable implementation of PCA may be used to perform the initial dimension reduction step: in particular, if the data have an irregular design (as happens frequently with function-valued data), the method of Yao *et al.* [33] may be applied to account for this; the ICA step then proceeds unchanged. We also note that while we find the *CubiICA* implementation of ICA to be the most successful in our signal separation task, other implementations such as *FastICA* [21] or *JADE* [34] can also be used. In general, ICA gives rows \hat{X}_i of the estimated mixing matrix that are maximally independent under a particular measure of independence involving, for example, higher sample moments or mutual information, in order to approximate the solution of the inverse problem in equation (1.1) under our assumption of independence between the rows of X . PCA and ICA have different purposes (respectively, orthogonal decomposition of variation and separation

of independently mixed signals) and we use them sequentially in IPCA. IPCA is non-parametric and, in particular, both distributionally and phylogenetically agnostic. This means that unlike PCA, IPCA is robust to non-Gaussianity in the data and, unlike phylogenetically corrected PCA, IPCA is robust to mis-specification of the phylogeny and to mixed phylogenetic and non-phylogenetic variation in the data: any of these can be features of biological data.

It can be seen in figure 4 that the estimation of ℓ is more challenging than the estimation of σ_n or σ_f , having greater bias and variance. This corresponds to the documented difficulty of estimating the parameter α in the OU model, particularly for smaller sample sizes. Our work on hyperparameter estimation in §2.4 mitigates these difficulties due to small sample size [16,35] by using bagging in order to bootstrap our sample. Somewhat unintuitively, bagging ‘works’ exactly because the subsample \hat{y} estimates are variable and thus we avoid overfitted final estimates (see electronic supplementary material, section S2). Conceptually, our work on hyperparameter estimation, when taken together with §2.2, relates to the character process models of Pletcher & Geyer [8] and orthogonal polynomial methods of Kirkpatrick & Heckman [5], which give estimates for the autocovariance of function-valued traits. Writing out equation (1.1) for a single function-valued trait (at the j th tip taxon, say), our model may be viewed as

$$f(x) = \sum_{i=1}^3 g_{ij} \phi_i(x) + \sum_{i=1}^3 e_{ij} \phi_i(x), \quad (3.3)$$

where the mixing coefficient X_{ij} has been expressed as the sum of g_{ij} , the genetic (i.e. phylogenetic) part of variation, plus e_{ij} , the non-phylogenetic (e.g. environmental) part of variation, just as in these references. Then, the autocorrelation

of the function-valued trait is

$$E[f(x_1)f(x_2)] = \sum_{i=1}^3 ((\sigma_i^f)^2 + (\sigma_i^n)^2) \phi_i(x_1)\phi_i(x_2). \quad (3.4)$$

The estimates of σ_i^f and σ_i^n obtained in §2.4 may be substituted into equation (3.4) to obtain an estimate of the autocovariance of the function-valued traits under study. This estimate has the attractions both of being positive definite (by construction) and of taking phylogeny into account.

Various frameworks exist that could be used to generalize the method presented in §2.4, to model heterogeneity of evolutionary rates along the branches of a phylogeny [36] or for multiple fixed [15] or randomly evolving [16,37] local optima of the mixing coefficients. For the stationary OU process, the optimum trait value appears only in the mean, and not in the covariance function, and so does not play a role as a parameter in GPR [20]. We have not implemented such extensions here, effectively assuming that a single fixed optimum is adequate for each mixing coefficient. Nonetheless, our framework is readily extensible to include such effects, either implicitly through branch-length transformations [38], or explicitly by replacing the OU model with the more general Hansen model [37].

R code for the IPCA, ancestral reconstruction and hyperparameter estimation is available from <https://github.com/fpgpr/>.

P.Z.H. and D.A.S. were supported by an EPSRC Institutional Sponsorship award to the University of Manchester and the EPSRC Mathematics Platform Engagement activity grant. C.G.K. was supported by a fellowship from the Wellcome Trust. The authors express their gratitude to John A. D. Aston for his encouragement and scientific advice.

References

- Maddison DR, Schulz KS. 2007 *The tree of life web project*. See <http://tolweb.org>.
- Li H *et al.* 2006 TREEFAM: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res.* **34**, D572–D580. (doi:10.1093/nar/gkj118)
- Yang Z, Rannala B. 2012 Molecular phylogenetics: principles and practice. *Nat. Rev. Genet.* **13**, 303–314. (doi:10.1038/nrg3186)
- Cheverud JM, Dow MM, Leutenegger W. 1985 The quantitative assessment of phylogenetic constraints in comparative analyses: sexual dimorphism in body weight among primates. *Evolution* **39**, 1335–1351. (doi:10.2307/2408790)
- Kirkpatrick M, Heckman N. 1989 A quantitative genetic model for growth, shape, reaction norms, and other infinite-dimensional characters. *J. Math. Biol.* **27**, 429–450. (doi:10.1007/BF00290638)
- The Functional Phylogenies Group. 2012 Phylogenetic inference for function-valued traits: speech sound evolution. *Trends Ecol. Evol.* **27**, 160–166. (doi:10.1016/j.tree.2011.10.001)
- Stinchcombe JR, Kirkpatrick M, Function-valued traits working group 2012 Phylogenetic inference for function-valued traits: speech sound evolution. *Trends Ecol. Evol.* **27**, 637–647. (doi:10.1016/j.tree.2012.07.002)
- Pletcher SD, Geyer CJ. 1999 The genetic analysis of age-dependent traits: modeling the character process. *Genetics* **153**, 825–835.
- Kingsolver JG, Gomulkiewicz R, Carter PA. 2001 Variation, selection and evolution of functionvalued traits. *Genetica* **112–113**, 87–104. (doi:10.1023/A:1013323318612)
- Zhang Z, Müller HG. 2011 Functional density synchronization. *Comput. Stat. Data Anal.* **55**, 2234–2249. (doi:10.1016/j.csda.2011.01.007)
- Moss SP, Joyce DA, Humphries S, Tindall KJ, Lunt DH. 2011 Comparative analysis of teleost genome sequences reveals an ancient intron size expansion in the zebrafish lineage. *Genome Biol. Evol.* **3**, 1187–1196. (doi:10.1093/gbe/evr090)
- Knight CG, Kassen R, Hebestreit H, Rainey PB. 2004 Global analysis of predicted proteomes: functional adaptation of physical properties. *Proc. Natl Acad. Sci. USA* **101**, 8390–8395. (doi:10.1073/pnas.0307270101)
- Jones NS, Moriarty J. 2013 Evolutionary inference of function-valued traits: Gaussian process regression on phylogenies. *J. R. Soc. Interface* **10**, 20120616. (doi:10.1098/rsif.2012.0616)
- Yao F, Coquery J, Le Cao KA. 2012 Independent principal component analysis for biologically meaningful dimension reduction of large biological data sets. *BMC Bioinformatics* **13**, 13–24. (doi:10.1186/1471-2105-13-13)
- Butler MA, King AA. 2004 Phylogenetic comparative analysis: a modelling approach for adaptive evolution. *Am. Nat.* **164**, 683–695. (doi:10.1086/426002)
- Beaulieu JM, Jhwueng DC, Boettiger C, O’eara BC. 2012 Modeling stabilizing selection: expanding the Ornstein–Uhlenbeck model of adaptive evolution. *Evolution* **8**, 2369–2383. (doi:10.1111/j.1558-5646.2012.01619.x)
- Breiman L. 1996 Bagging predictors. *Mach. Learn.* **24**, 123–140.
- Venditti C, Meade A, Pagel M. 2010 Phylogenies reveal new interpretation of speciation and the red queen. *Nature* **463**, 349–352. (doi:10.1038/nature08630)
- Hansen T. 1997 Stabilizing selection and the comparative analysis of adaptation. *Evolution* **51**, 1341–1351. (doi:10.2307/2411186)

20. Rasmussen CE, Williams CKI. 2006 *Gaussian processes for machine learning*. Cambridge, MA: MIT Press.
21. Hyvärinen A, Oja E. 2000 Independent component analysis: algorithms and applications. *Neural Netw.* **13**, 411–430. (doi:10.1016/S0893-6080(00)00026-5)
22. Bishop C. 2006 *Pattern recognition and machine learning*. Berlin, Germany: Springer.
23. Revell LJ. 2009 Size-correction and principal components for interspecific comparative studies. *Evolution* **63**, 3258–3268. (doi:10.1111/j.1558-5646.2009.00804.x)
24. Scholz M, Gatzek S, Sterling A, Fiehn O, Selbig J. 2004 Metabolite fingerprinting: detecting biological features by independent component analysis. *Bioinformatics* **20**, 2447–2454. (doi:10.1093/bioinformatics/bth270)
25. Minka TP. 2000 Automatic choice of dimensionality for PCA. *Adv. Neural Inf. Process. Syst.* **13**, 514.
26. Blaschke T, Wiskott L. 2004 CuBICA: independent component analysis by simultaneous third- and fourth-order cumulant diagonalization. *IEEE Trans. Signal Process.* **52**, 1250–1256. (doi:10.1109/TSP.2004.826173)
27. Hansen T, Martins E. 1996 Translating between microevolutionary process and macroevolutionary patterns: the correlation structure of interspecific data. *Evolution* **50**, 1404–1417. (doi:10.2307/2410878)
28. Pagel M. 1997 Inferring evolutionary processes from phylogenies. *Zool. Scr.* **26**, 331–348. (doi:10.1111/j.1463-6409.1997.tb00423.x)
29. Jaynes ET. 1984 Prior information and ambiguity in inverse problems. *Inverse Probl.* **14**, 151–166.
30. Salamin N, Wuest RO, Lavergne S, Thuiller W, Pearman PB. 2010 Assessing rapid evolution in a changing environment. *Trends Ecol. Evol.* **25**, 692–698. (doi:10.1016/j.tree.2010.09.009)
31. Hadfield JD, Nakagawa S. 2010 General quantitative genetic methods for comparative biology: phylogenies, taxonomies and multi-trait models for continuous and categorical characters. *J. Evol. Biol.* **23**, 494–508. (doi:10.1111/j.1420-9101.2009.01915.x)
32. Kerr M. 2012 Evolutionary inference for functional data: using Gaussian processes on phylogenies of functional data objects. MSc Thesis, University of Glasgow, Glasgow, UK.
33. Yao F, Müller HG, Wang JL. 2005 Functional data analysis for sparse longitudinal data. *J. Am. Stat. Assoc.* **100**, 577–590. (doi:10.1198/016214504000001745)
34. Cardoso JF. 1999 High-order contrasts for independent component analysis. *Neural Comput.* **11**, 157–192. (doi:10.1162/089976699300016863)
35. Collar DC, O'Meara BC, Wainwright PC, Near TJ. 2009 Piscivory limits diversification of feeding morphology in centrarchid fishes. *Evolution* **63**, 1557–1573. (doi:10.1111/j.1558-5646.2009.00626.x)
36. Revell LJ, Collar DC. 2009 Phylogenetic analysis of the evolutionary coreally using likelihood. *Evolution* **63–64**, 1090–1100. (doi:10.1111/j.1558-5646.2009.00616.x)
37. Hansen T, Pienaar J, Orzack SH. 2008 A comparative method for studying adaptation to a randomly evolving environment. *Evolution* **8**, 1965–1977. (doi:10.1111/j.1558-5646.2008.00412.x)
38. Pagel M. 1999 Inferring the historical patterns of biological evolution. *Nature* **401**, 877–884. (doi:10.1038/44766)