

# UCLA

## UCLA Previously Published Works

### Title

Hybrid principal components analysis for region-referenced longitudinal functional EEG data.

### Permalink

<https://escholarship.org/uc/item/5d05q498>

### Journal

Biostatistics, 21(1)

### ISSN

1465-4644

### Authors

Scheffler, Aaron  
Telesca, Donatello  
Li, Qian  
[et al.](#)

### Publication Date

2020

### DOI

10.1093/biostatistics/kxy034

Peer reviewed

# Hybrid principal components analysis for region-referenced longitudinal functional EEG data

AARON SCHEFFLER, DONATELLO TELESKA, QIAN LI

*Department of Biostatistics, University of California Los Angeles, 650 Charles E Young Drive, Los Angeles, CA, 90095, USA*

CATHERINE A. SUGAR

*Department of Biostatistics, University of California Los Angeles, 650 Charles E Young Drive, Los Angeles, CA, 90095, USA and Department of Psychiatry and Biobehavioral Sciences, University of California Los Angeles, 757 Westwood Plaza, Los Angeles, CA, 90095, USA*

CHARLOTTE DISTEFANO, SHAFALI JESTE

*Department of Psychiatry and Biobehavioral Sciences, University of California Los Angeles, 757 Westwood Plaza, Los Angeles, CA, 90095, USA*

DAMLA ŞENTÜRK\*

*Department of Biostatistics, University of California Los Angeles, 650 Charles E Young Drive, Los Angeles, CA, 90095, USA*  
dsenturk@ucla.edu

## SUMMARY

Electroencephalography (EEG) data possess a complex structure that includes regional, functional, and longitudinal dimensions. Our motivating example is a word segmentation paradigm in which typically developing (TD) children, and children with autism spectrum disorder (ASD) were exposed to a continuous speech stream. For each subject, continuous EEG signals recorded at each electrode were divided into one-second segments and projected into the frequency domain via fast Fourier transform. Following a spectral principal components analysis, the resulting data consist of region-referenced principal power indexed regionally by scalp location, functionally across frequencies, and longitudinally by one-second segments. Standard EEG power analyses often collapse information across the longitudinal and functional dimensions by averaging power across segments and concentrating on specific frequency bands. We propose a hybrid principal components analysis for region-referenced longitudinal functional EEG data, which utilizes both vector and functional principal components analyses and does not collapse information along any of the three dimensions of the data. The proposed decomposition only assumes weak separability of the higher-dimensional covariance process and utilizes a product of one dimensional eigenvectors and eigenfunctions, obtained from the regional, functional, and longitudinal marginal covariances, to represent the observed data, providing a computationally feasible non-parametric approach. A mixed

\*To whom correspondence should be addressed.

effects framework is proposed to estimate the model components coupled with a bootstrap test for group level inference, both geared towards sparse data applications. Analysis of the data from the word segmentation paradigm leads to valuable insights about group-region differences among the TD and verbal and minimally verbal children with ASD. Finite sample properties of the proposed estimation framework and bootstrap inference procedure are further studied via extensive simulations.

*Keywords:* Electroencephalography; Functional data analysis; Marginal covariances; Product functional principal components decomposition; Spectral principal components decomposition.

## 1. INTRODUCTION

Approximately 30% of children with autism spectrum disorder (ASD) never gain spoken language (referred to as “minimally verbal”) and the reasons are largely unknown (Tager-Flusberg and Kasari, 2013). A major barrier in conducting research with minimally verbal children is the limited availability of appropriate assessment techniques. The recording of electroencephalography (EEG) signals during our motivating study, involving a word segmentation paradigm, gave researchers a unique opportunity to compare and contrast neurocognitive processes involved in language and communication development among verbal ASD (vASD), minimally verbal ASD (mvASD), and typically developing (TD) children, without relying on the children’s ability to understand directions or provide an overt behavioral response. EEG is a popular non-invasive method for measuring voltage fluctuations across scalp regions in order to characterize neurocognitive processes and disorders. Children listened to a continuous speech stream, which contained four “made-up” words, each composed of three different phonemes or units of sound (Figure 1(a) and (b)). The four words were repeated 45 times in random order such that no word was used twice in a row, and there was no time gap between words. The full experiment took 144 s. Children were expected to segment the speech stream, i.e. identify boundaries between words, by recognizing the differential patterns in the phonemes (Scott-Van Zeeland and others, 2010).

EEG studies, including both event-related and resting state paradigms, create high-dimensional data with regional, functional, and longitudinal dimensions. Data from resting state paradigms are typically analyzed in the frequency domain, while event-related paradigms, where stimuli are applied repeatedly throughout the experiment, are analyzed either in the time or frequency domain. In our word segmentation paradigm, an event-related study, quantities considered of interest are in the frequency domain. Hence EEG signals, collected from an 128 electrode sensor net, were divided into one-second segments and projected into the frequency domain via fast Fourier transform (FFT). Given the fact that EEG signals have low spatial resolution and that neighboring electrodes have similar power spectra, spectral principal components analysis (PCA) has been proposed to combine information from EEG signals recorded at electrodes within a scalp region (Ombao and Moon-Ho, 2006). This pre-processing step produces region-referenced principal power, following a region-referenced longitudinal functional stochastic process. Specifically, the scalp locations represent the regional dimension, principal power obtained across frequencies represents the functional dimension, and the one-second EEG segments represent the longitudinal dimension. Similarly, if the quantities of interest in an event-related paradigm are in the time domain, event-related potentials (ERP) time-locked to each stimulus (potentially combined over electrodes within a scalp region) would represent the functional dimension, and repetitions of the stimuli throughout the experiment would represent the longitudinal dimension. Note that all three dimensions of the observed data carry distinct interpretations and that longitudinal time (captured through segments across the experiment) may play an important role, especially in learning paradigms in which the focus is on changes over experimental time as learning evolves.

Standard analysis of high-dimensional EEG data involves collapsing information along multiple dimensions. The longitudinal dimension is collapsed when power spectra are averaged over segments or ERP

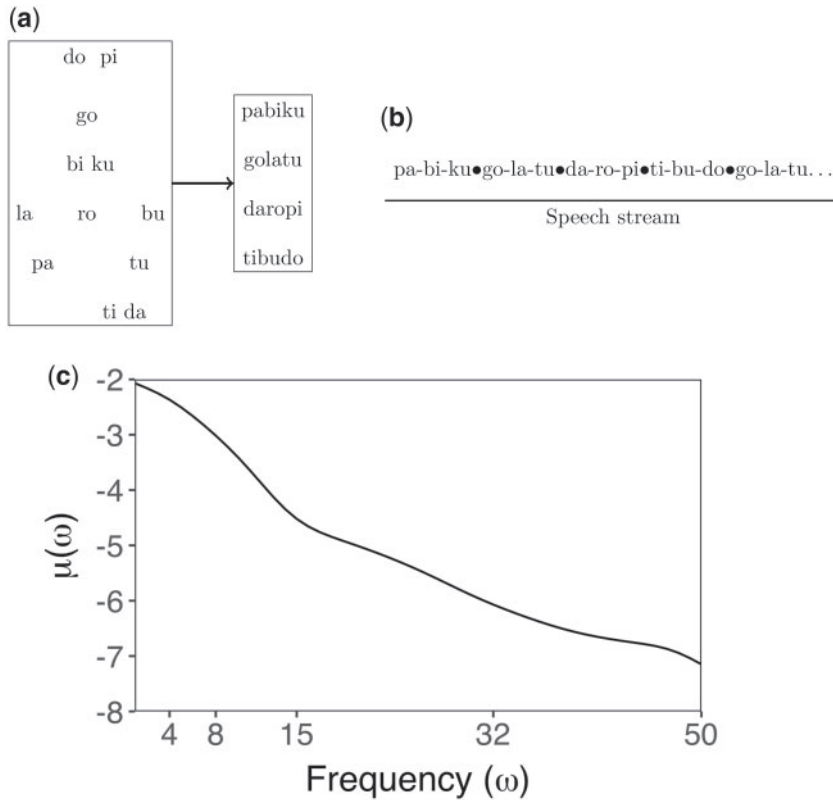


Fig. 1. (a) Four “made-up” words formed by concatenating three phonemes from a set of 12 phonemes without repetition in the word segmentation paradigm. (b) The artificial speech stream generated during the word segmentation paradigm. Breaks between phonemes are denoted by a dash and breaks between words are denoted by a dot. (c) The estimated mean log principal power  $\mu(\omega)$  for subjects pooled across the TD, vASD, and mvASD groups.

curves are averaged over stimuli. Similarly, analysis of spectral power from specific frequency bands or specific ERP curve features corresponds to collapsing of the functional dimension, while averages over scalp regions collapse the regional dimension of the data. We propose a hybrid principal components analysis (HPCA) for region-referenced longitudinal functional EEG data that does not collapse any of the three dimensions. We call the proposed decomposition hybrid, since it combines vector principal components analysis along the regional dimension (lacking a time order) and functional principal components analysis along the longitudinal and functional dimensions, providing an efficient non-parametric means of modeling high-dimensional EEG data. The HPCA decomposition involves a product of one-dimensional eigenvectors and eigenfunctions obtained from marginal covariances along the three dimensions of the data. A central assumption in this low dimensional, and hence computationally feasible, framework is the weak separability of the overall covariance process of the observed data. The concept of weak separability, recently proposed by [Chen and Lynch \(2017\)](#), refers to the idea that the covariance can be approximated by a weighted sum of separable covariance components and implies that the direction of variation (i.e. eigenvectors/eigenfunctions) along one of the three dimensions of the EEG data is the same across fixed slices of the other two dimensions. Note that this assumption is weaker than the commonly assumed strong separability of covariance surfaces in higher dimensions, which requires that the entire covariance

structure, not only the directions of variation, is the same up to a constant across fixed slices of the other dimensions.

The literature on functional data analysis has proliferated over the past two decades, with methodological developments motivated by the complex dependency structures of repeatedly measured curves. Most of the recent developments on functional principal components analysis (FPCA) consider either longitudinally or spatially repeated functional data but not both. For longitudinally repeated functional data, [Di and others \(2009\)](#) proposed multilevel ANOVA decompositions. [Greven and others \(2010\)](#) extended their work to linear longitudinal decompositions, and [Chen and Müller \(2012\)](#), [Park and Staicu \(2015\)](#), [Chen and others \(2016\)](#) and [Hasenstab and others \(2017\)](#) considered more flexible non-linear forms. For spatially repeated functional data, [Staicu and others \(2010\)](#), [Zhou and others \(2010\)](#) and [Liu and others \(2017\)](#) considered parametric forms, while [Huang and others \(2017\)](#) proposed a non-parametric decomposition. Of the proposed methods, only [Hasenstab and others \(2017\)](#) decomposed both longitudinal and regional sources of functional variation in three dimensions via a multi-dimensional FPCA procedure (MD-FPCA). MD-FPCA, motivated by the analysis of the high-dimensional event-related ERP data in the time domain (through ERPs), treated scalp regions as exchangeable. The proposed HPCA method relaxes this assumption and involves a much simpler and computationally efficient decomposition via the weak separability of the covariance process. Product FPCA of [Chen and others \(2016\)](#) also relies on weak separability and involves a product of one-dimensional eigenfunctions in the proposed decomposition; but their developments are obtained for two-dimensional functional data. HPCA extends product FPCA approach of [Chen and others \(2016\)](#) to higher dimensions targeting region-referenced longitudinal functional EEG data and combining vector and functional principal components analysis. In addition, while developments for product FPCA have only focused on densely measured longitudinally observed functional data, the estimation and inference procedures proposed for HPCA focus on sparse EEG data applications.

The outline of the article is as follows. Section 2 introduces spectral PCA as a pre-processing step with minimal loss of information that produces region-referenced longitudinal functional data. Section 3 introduces the HPCA decomposition, develops an innovative mixed effects framework for estimation of the model components, specifically geared towards sparse data applications, and outlines a bootstrap procedure for group-level inference. We highlight that the developments for sparse data applications are novel. Prediction of subject-specific scores based on sparse data have not yet been considered for decompositions based on weak separability of the covariance process, such as the product FPCA. The proposed mixed effects framework is also utilized to assess the weak separability assumption via the random effects correlation structure. Section 4 provides insights from the word segmentation paradigm including inference on group-region differences in spectral dynamics among TD, vASD, and mvASD children. We assess the proposed decomposition and the associated bootstrap test with an extensive simulation study summarized in Section 5 and conclude with a discussion in Section 6.

## 2. SPECTRAL PCA AND THE RESULTING REGION-REFERENCED LONGITUDINAL FUNCTIONAL EEG DATA

Given that EEG signals measured on neighboring electrodes are highly multi-collinear due to their spatial proximity, the analysis of EEG data collected from high density electrode arrays is often preceded by reduction of the electrode dimension to discard redundant information and facilitate interpretation. When analysis takes place in the frequency domain, dimensional reduction is often unsatisfactorily carried out by selecting spectra from a single electrode or averaging spectra within a scalp region. Alternatively, given that electrodes within a scalp region possess similar spectra, spectral PCA has been proposed to pool spectral information within a scalp region with minimal loss of information. Spectral PCA applications in the analysis of time series data date back to [Brillinger \(1981\)](#), but we follow a more recent application to EEG data by [Ombao and Moon-Ho \(2006\)](#). They utilize spectral PCA as an exploratory tool to consolidate power

spectra in a scalp region by utilizing overlapping segments of the continuous multi-channel time-series recorded at multiple electrodes in a seizure study. In contrast, we perform spectral PCA on non-overlapping EEG segments as a pre-processing step to be followed by scalp-wide analysis.

We highlight the outline of spectral PCA procedure here and defer details to Appendix A of the [supplementary material](#) available at *Biostatistics* online. Fourier coefficients at a fixed frequency are obtained via FFT for EEG signals measured from electrodes within the same scalp region and collected in a region-specific periodogram matrix. Following smoothing of each term of the periodogram matrices over frequencies, principal power is defined as the normalized leading eigenvalue of the smoothed periodogram matrix, representing the common variation in the fixed frequency across the electrodes (relative to variation in other frequencies) in a given scalp region along the direction of the leading eigenvector. The interpretation of principal power is closely tied to the goal of spectral PCA in combining signals across electrodes within a given scalp region. The assumption that electrodes within a scalp region have similar spectral densities implies that the region-specific periodogram matrix at a particular frequency would be of low rank. Hence, extracting the largest eigenvalue would serve as a reasonable summary of the spectral dynamics within a brain region. While our analysis focuses on the largest eigenvalue as principal power, note that second and third eigenvalues can also be modeled similarly via HPCA, allowing further analysis of the spectral dynamics among brain regions. Spectral PCA being applied at each segment and region for each subject, yields region-referenced longitudinal functional data, i.e. principal power as a function of region  $r$ , frequency  $\omega$ , and segment  $s$  denoted by  $Y_{di}(r, \omega, s)$ . If a given subject does not have valid data at a fixed segment then the principal power for that segment is considered missing. We model  $Y_{di}(r, \omega, s)$  as a summary measure of the power dynamics across the scalp.

### 3. HYBRID PRINCIPAL COMPONENTS ANALYSIS (HPCA)

#### 3.1. The HPCA decomposition

Let  $Y_{di}(r, \omega, s)$  denote the log principal power, which comprises region-referenced longitudinal functional data observed for subject  $i$ ,  $i = 1, \dots, n_d$ , from group  $d$ ,  $d = 1, \dots, D$ , in region  $r$ ,  $r = 1, \dots, R$ , at frequency  $\omega$ ,  $\omega \in \Omega$ , and segment  $s$ ,  $s \in \mathcal{S}$ . Here  $\Omega$  and  $\mathcal{S}$  represent the functional and longitudinal domains, respectively, and  $Y_{di}(r, \omega, s)$  is assumed to be square-integrable. Even though subjects may not be observed at all segments  $s \in \mathcal{S}$ , we use a common index set in the formulations below for notational ease. Note that the smoothing-based estimation procedure proposed in the next section, will readily extend to subject-specific sparse longitudinal domains. Further let  $Z_{di}(r, \omega, s) = Y_{di}(r, \omega, s) - \mu(\omega, s) - \eta_d(r, \omega, s) - \epsilon_{di}(r, \omega, s)$  denote a de-meaned and de-noised region-referenced stochastic process, where  $\mu(\omega, s)$  and  $\eta_d(r, \omega, s)$  denote the functional fixed effects that represent the overall mean function and group-region shifts, respectively, and  $\epsilon_{di}(r, \omega, s)$  denotes the measurement error with mean zero and variance  $\sigma_d^2$ .

The proposed HPCA decomposition provides a lower dimensional approximation of a stochastic process defined over regional, functional, and longitudinal dimensions in terms of an empirical orthonormal basis based on eigenvectors and eigenfunctions obtained from the marginal covariances in each dimension. A central assumption of HPCA is the weak separability of the overall three-dimensional covariance process, which implies that the direction of variation (i.e. eigenvectors/eigenfunctions) along any one of the three dimensions of the EEG data is the same across fixed slices of the other two dimensions. This assumption is less stringent than the strong separability commonly assumed in the analysis of spatio-temporal stochastic processes, which requires that the entire covariance process along one dimension, not only the direction of variation, is the same up to a constant across fixed slices of the other dimensions. Note that the eigenfunctions or eigenvectors being the same does not necessarily imply the same covariance surface at fixed slices of the other dimensions due to weighting through the eigenvalues. We refer readers to [Chen and Lynch \(2017\)](#) for a detailed comparison of weak versus strong separability

and note that we propose two separate checks for the weak separability assumption in Section 3.2 and Appendix D of the [supplementary material](#) available at *Biostatistics* online, through a test for the correlation structure of the random effects in the mixed effects modeling and through visualization of the data, respectively.

Under weak separability, the common eigenfunctions and eigenvectors along each of the three dimensions can be estimated using the marginal covariances. Let the functional and longitudinal marginal covariance surfaces be defined as

$$\begin{aligned}\Sigma_{d,\Omega}(\omega, \omega') &= \sum_r \int_S \text{cov}\{Z_{di}(r, \omega, s), Z_{di}(r, \omega', s)\} ds = \sum_{\ell=1}^{\infty} \tau_{d\ell,\Omega} \phi_{d\ell}(\omega) \phi_{d\ell}(\omega'), \\ \Sigma_{d,S}(s, s') &= \sum_r \int_{\Omega} \text{cov}\{Z_{di}(r, \omega, s), Z_{di}(r, \omega, s')\} d\omega = \sum_{m=1}^{\infty} \tau_{dm,S} \psi_{dm}(s) \psi_{dm}(s'),\end{aligned}$$

and let  $\Sigma_{d,\mathcal{R}}$  denote the regional marginal covariance matrix with  $(r, r')$ -th element equal to

$$(\Sigma_{d,\mathcal{R}})_{r,r'} = \int_S \int_{\Omega} \text{cov}\{Z_{di}(r, \omega, s), Z_{di}(r', \omega, s)\} d\omega ds = \sum_{k=1}^R \tau_{dk,\mathcal{R}} \mathbf{v}_{dk}(r) \mathbf{v}_{dk}(r'),$$

where  $\phi_{d\ell}(\omega)$  and  $\psi_{dm}(s)$  are the common eigenfunctions of the functional and longitudinal marginal covariance surfaces, respectively;  $\mathbf{v}_{dk}(r)$  are the common eigenvectors for the regional marginal covariance matrix; and  $\tau_{d\ell,\Omega}$ ,  $\tau_{dm,S}$  and  $\tau_{dk,\mathcal{R}}$  are the respective eigenvalues. While we estimate the regional marginal covariance matrix nonparametrically, we note that parametric approaches have been quite popular for modeling spatial covariances. An important difference of the current EEG application from typical environmental applications is that in the latter spatial data may typically be observed only once over the location grid at a fixed time point, while we observe the region-specific longitudinal functional EEG data repeatedly over subjects. Parametric assumptions to interpolate information across regions are thus not necessarily needed in modeling the spatial dependence in our application and we use a non-parametric region marginal covariance matrix, mimicking the non-parametric marginal functional and longitudinal covariance surfaces.

Utilizing the eigenfunctions and eigenvectors of the marginal covariances, the HPCA decomposition of  $Y_{di}(r, \omega, s)$  is given as

$$\begin{aligned}Y_{di}(r, \omega, s) &= \mu(\omega, s) + \eta_d(r, \omega, s) + Z_{di}(r, \omega, s) + \epsilon_{di}(r, \omega, s) \\ &= \mu(\omega, s) + \eta_d(r, \omega, s) + \sum_{k=1}^R \sum_{\ell=1}^{\infty} \sum_{m=1}^{\infty} \xi_{di,k\ell m} \mathbf{v}_{dk}(r) \phi_{d\ell}(\omega) \psi_{dm}(s) + \epsilon_{di}(r, \omega, s).\end{aligned}\quad (3.1)$$

In (3.1), the subject-specific scores  $\xi_{di,k\ell m}$  are defined through the projection,  $\langle Z_{di}(r, \omega, s), \mathbf{v}_{dk}(r) \phi_{d\ell}(\omega) \psi_{dm}(s) \rangle = \sum_{r=1}^R \int \int Z_{di}(r, \omega, s) \mathbf{v}_{dk}(r) \phi_{d\ell}(\omega) \psi_{dm}(s) d\omega ds$ , of the de-meaned and de-noised stochastic process,  $Z_{di}(r, \omega, s)$ , onto the orthonormal bases  $\mathbf{v}_{dk}(r) \phi_{d\ell}(\omega) \psi_{dm}(s)$  defined as the product of the one-dimensional eigenfunctions and eigenvectors of the marginal covariances. Note that the set of subject-specific scores  $(\xi_{di,k\ell m})$  are uncorrelated over regions, frequencies and segments under weak separability. Hence, the proposed HPCA expansion also leads to a decomposition of the total covariance,

$\Sigma_{d,T}\{(r, \omega, s), (r', \omega', s')\}$ , of  $Y_{di}(r, \omega, s)$ , as follows,

$$\begin{aligned} \Sigma_{d,T}\{(r, \omega, s), (r', \omega', s')\} &= \text{cov}\{Z_{di}(r, \omega, s), Z_{di}(r', \omega', s')\} + \sigma_d^2 \delta\{(r, \omega, s), (r', \omega', s')\} \\ &= \sum_{k=1}^R \sum_{\ell=1}^{\infty} \sum_{m=1}^{\infty} \tau_{d,k\ell m} \mathbf{v}_{dk}(r) \phi_{d\ell}(\omega) \psi_{dm}(s) \mathbf{v}_{dk}(r') \phi_{d\ell}(\omega') \psi_{dm}(s') + \sigma_d^2 \delta\{(r, \omega, s), (r', \omega', s')\}, \end{aligned}$$

where  $\tau_{d,k\ell m} = \text{var}(\xi_{di,k\ell m})$  and  $\delta\{(r, \omega, s), (r', \omega', s')\}$  denotes the indicator for  $\{(r, \omega, s) = (r', \omega', s')\}$ . Note that the total covariance is written as a weighted sum of separable regional, functional and longitudinal covariances. One way of assessing the weak separability assumption will be to examine the correlation structure of the subject-specific decomposition scores  $\xi_{di,k\ell m}$  via the mixed effects modeling framework proposed in Section 3.2.

In practice, the HPCA decomposition is truncated to include only  $K$ ,  $L$ , and  $M$  eigencomponents for the regional, functional, and longitudinal marginal covariances in the expansion, respectively, with truncation based on the fraction of variance explained (FVE). A general guideline is to initially include marginal eigencomponents in the HPCA expansion that explain approximately 90% of variation in their respective marginal covariances. Some of these components may be eliminated after subject-specific scores and their associated variance components are estimated via the proposed mixed effects modeling framework of Section 3.2, which provide an overall estimate of FVE, not only for the separate marginal covariances, but for the covariance based on the entire data. Details on the selection of the number of eigencomponents are deferred to Section 3.2.

Note that the three-dimensional (3D) HPCA introduced in (3.1) reduces to a two-dimensional (2D) HPCA with the regional and functional dimensions when the longitudinal dimension may not be of interest or may not exhibit change. Given that the remaining indices and arguments are unchanged, let  $Y_{di}(r, \omega)$  denote the region-referenced functional data with a weakly separable covariance process. Utilizing the eigenfunctions and eigenvectors of the marginal covariances, the 2D HPCA decomposition of  $Y_{di}(r, \omega)$  can be given as,

$$\begin{aligned} Y_{di}(r, \omega) &= \mu(\omega) + \eta_d(r, \omega) + Z_{di}(r, \omega) + \epsilon_{di}(r, \omega) \\ &= \mu(\omega) + \eta_d(r, \omega) + \sum_{k=1}^R \sum_{\ell=1}^{\infty} \xi_{di,k\ell} \mathbf{v}_{dk}(r) \phi_{d\ell}(\omega) + \epsilon_{di}(r, \omega), \end{aligned}$$

where model components and the decomposition of the total variance are defined as in the 3D HPCA by omitting the longitudinal argument  $s$ . The functional dimension can similarly be collapsed leading to the 2D HPCA with only the regional and longitudinal dimensions. The discussion will continue to center on the 3D HPCA with the understanding that extensions to 2D HPCA are available by omitting one of the continuous arguments.

Motivated by the high-dimensional EEG data, both 3D and 2D HPCA extend the product FPCA of [Chen and others \(2016\)](#) for longitudinally observed functional data by the addition of a regional dimension. Moreover, HPCA involves a hybrid decomposition for the region-referenced longitudinal functional EEG data, combining vector and functional principal components analysis under the assumption of weak separability. Another important divergence from the product FPCA formulation is in estimation. Motivated by the longitudinally sparse EEG data, we next propose a novel mixed effects procedure framework for estimation of the model components, specifically geared towards sparse data applications (with low number of repetitions and irregular spacing in observations over the longitudinal dimension). The estimation and testing procedures proposed for the product FPCA largely depend on projection techniques, which are applicable only to densely measured longitudinal functional data ([Chen and others, 2016](#); [Chen and Lynch, 2017](#)).



### 3.2. Estimation of model components

The section below outlines the estimation of all the model components, including functional fixed effects, marginal covariances, and eigencomponents, a novel mixed effects framework for estimation of subject-specific decomposition scores and associated variance components, and a recommendation to select the number of eigencomponents included in the proposed HPCA. We begin by introducing the HPCA estimation algorithm.

---

**Algorithm:** *HPCA estimation procedure*

---

- (1) Estimation of fixed effects
    - i. Calculate  $\hat{\mu}(\omega, s) = \sum_{d=1}^D \hat{\mu}_d(\omega, s)$  by applying a bivariate penalized spline smoother to all observed data  $\{\omega, s, Y_{di}(r, \omega, s) : i = 1, \dots, n_d; r = 1, \dots, R; \omega \in \Omega; s \in \mathcal{S}\}$ .
    - ii. Calculate  $\hat{\eta}_d(r, \omega, s)$  by applying a bivariate penalized spline smoother to all observed data  $\{\omega, s, \hat{Y}_{di}(r, \omega, s) - \hat{\mu}(\omega, s) : i = 1, \dots, n_d; \omega \in \Omega; s \in \mathcal{S}\}$ .
  - (2) Estimation of marginal covariances and measurement error variance
    - i. Calculate  $\tilde{\Sigma}_{d,\Omega}(\omega, \omega')$  and  $\tilde{\Sigma}_{d,\mathcal{S}}(s, s')$  by applying bivariate penalized spline smoothers to the pooled covariances,  $\hat{\Sigma}_{d,\Omega}(\omega, \omega')$  and  $\hat{\Sigma}_{d,\mathcal{S}}(s, s')$ , respectively.
    - ii. Calculate  $\hat{\sigma}_d^2$  by averaging the measurement error variance estimates  $\hat{\sigma}_{d,\Omega}^2$  and  $\hat{\sigma}_{d,\mathcal{S}}^2$ .
    - iii. Calculate  $\tilde{\Sigma}_{d,\mathcal{R}}$  by removing the estimated measurement error variance  $\hat{\sigma}_d^2$  from the diagonal entries of the pooled covariance  $\hat{\Sigma}_{d,\mathcal{R}}$ .
  - (3) Estimation of marginal eigencomponents
    - i. Employ FPCA on  $\tilde{\Sigma}_{d,\Omega}(\omega, \omega')$  and  $\tilde{\Sigma}_{d,\mathcal{S}}(s, s')$  to estimate the eigenvalue, eigenfunction pairs,  $\{\tau_{d\ell,\Omega}, \phi_{d\ell}(\omega) : \ell = 1, \dots, L\}$  and  $\{\tau_{dm,\mathcal{S}}, \psi_{dm}(s) : m = 1, \dots, M\}$ , respectively.
    - ii. Employ PCA on  $\tilde{\Sigma}_{d,\mathcal{R}}$  to estimate the eigenvalue, eigenvector pairs  $\{\tau_{dk,\mathcal{R}}, v_{dk}(r) : k = 1, \dots, K\}$ .
  - (4) Estimation of variance components and subject-specific scores via linear mixed effects models
    - i. Calculate  $\hat{k}_{dg}$  and  $\hat{\sigma}_d^2$  by fitting the proposed linear mixed effects model.
    - ii. Calculate  $\hat{\zeta}_{dig}$  as the BLUP  $\hat{\zeta}_{dig} = E(\zeta_{dig} | \mathbf{Y}_{di})$ .
    - iii. Select  $G'$  such that  $FVE_{dG'} > 0.8$  for  $d = 1, \dots, D$  and form predictions  $\hat{Y}_{di}(r, \omega, s)$ .
- 

We defer details on steps 1–3 to Appendix B of the [supplementary material](#) available at *Biostatistics* online in which we refer readers to previous works on well-established mean, covariance, and eigencomponent estimation. However, we briefly highlight two novel estimation procedures found in step 2 for the measurement error variance,  $\sigma_d^2$ , and regional marginal covariance,  $\Sigma_{d,\mathcal{R}}$ . While previous authors obtain estimates of the measurement error variance using smoothing techniques on the raw covariance from a single dimension (Yao and others, 2005; Park and Staicu, 2015), we adapt this method to high dimensional settings by pooling information across both the functional and longitudinal marginal covariances. We then use this pooled estimate to remove the measurement error variance from the diagonals of the raw regional marginal covariance, which as a matrix is not amenable to smoothing techniques. Thus, we are able to leverage information from both the functional and longitudinal dimensions to obtain a decontaminated estimate of the regional marginal covariance.

In step 4, we utilize the estimated functional fixed effects and marginal eigencomponents to propose a linear mixed effects framework for modeling sparsely observed region-referenced longitudinal functional EEG data. In addition to allowing estimation of subject-specific scores under the assumption of their joint normality with the measurement error, the proposed mixed effects framework also provides final estimates for the corresponding variance components and the measurement error variance. The variance

components estimates associated with the subject-specific scores are utilized in selection of the number of eigenvectors included in the HPCA decomposition via estimation of the proportion of variance explained, as well as in the construction of a hypothesis testing procedure for group-level inference via the bootstrap. Finally, the proposed mixed effects framework provides an opportunity to check the weak separability assumption via examining the correlation structure of the random effects.

For ease of notation, we replace the triple index  $k\ell m$  in HPCA truncated at  $K$ ,  $L$ , and  $M$  with a single index  $g = (k - 1) + K(\ell - 1) + KL(m - 1) + 1$ ,

$$Y_{di}(r, \omega, s) = \mu(\omega, s) + \eta_d(r, \omega, s) + \sum_{g=1}^G \zeta_{dig} \varphi_{dg}(r, \omega, s) + \epsilon_{di}(r, \omega, s),$$

where  $\varphi_{dg}(r, \omega, s) = v_{dk}(r) \phi_{d\ell}(\omega) \psi_{dm}(s)$ ,  $\zeta_{dig} = \langle Z_{di}(r, \omega, s), \varphi_{dg}(r, \omega, s) \rangle$ ,  $\kappa_{dg} = \text{cov}(\zeta_{dig})$  and  $G = KLM$ . Denote by  $\mathbf{Y}_{di}$  the vectorized form of  $Y_{di}(r, \omega, s)$  over the subject-specific region, frequency and segment grid for subject  $i$ ,  $i = 1, \dots, n_d$ . In our EEG application, while the region and frequency grids are the same for all subjects, the segment grid is subject-specific due to data quality issues. Similar subject-specific vectorized forms for the functional fixed effects,  $\mu(\omega, s)$  and  $\eta_d(r, \omega, s)$ , the region-referenced stochastic process  $Z_{di}(r, \omega, s)$ , the measurement error  $\epsilon_{di}(r, \omega, s)$ , and the multidimensional orthonormal basis  $\varphi_{dg}(r, \omega, s)$  are denoted by  $\boldsymbol{\mu}_i$ ,  $\boldsymbol{\eta}_{di}$ ,  $\mathbf{Z}_{di}$ ,  $\boldsymbol{\epsilon}_{di}$  and  $\boldsymbol{\varphi}_{dig}$ , respectively. Note that the mean vectors  $\boldsymbol{\mu}_i$ ,  $\boldsymbol{\eta}_{di}$  are indexed by subject since they are defined over the subject-specific region, frequency, and segment grids. Under the assumption that  $\zeta_{dig}$  and  $\boldsymbol{\epsilon}_{di}$  are jointly Gaussian, the proposed linear mixed effects model is given as

$$\mathbf{Y}_{di} = \boldsymbol{\mu}_i + \boldsymbol{\eta}_{di} + \mathbf{Z}_{di} + \boldsymbol{\epsilon}_{di} = \boldsymbol{\mu}_i + \boldsymbol{\eta}_{di} + \sum_{g=1}^G \zeta_{dig} \boldsymbol{\varphi}_{dig} + \boldsymbol{\epsilon}_{di}, \quad \text{for } i = 1, \dots, n_d. \quad (3.2)$$

The model can be fit separately in each group,  $d = 1, \dots, D$ , with both  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\eta}_{di}$  previously obtained by smoothing. The functional, longitudinal, and regional dependencies in  $\mathbf{Y}_{di}$  are induced through the subject-specific random effects  $\zeta_{dig}$  in (3.2). Given estimates for  $\boldsymbol{\mu}_i$ ,  $\boldsymbol{\eta}_{di}$ , and  $\boldsymbol{\varphi}_{dig}$ , estimates of the variance components,  $\kappa_{dg}$  and  $\sigma_d^2$  are obtained using maximum likelihood.

Following Yao *and others* (2005) in using conditional expectations to estimate subject-specific scores for sparse functional data, the  $\zeta_{dig}$  are estimated using best linear unbiased prediction (BLUP),  $\hat{\zeta}_{dig} = E(\zeta_{dig} | \mathbf{Y}_{di}) = \hat{\kappa}_{dg} \hat{\boldsymbol{\varphi}}_{dig} \hat{\boldsymbol{\Sigma}}_{\mathbf{Y}_{di}}^{-1} (\mathbf{Y}_{di} - \hat{\boldsymbol{\mu}}_i - \hat{\boldsymbol{\eta}}_{di})$ , where  $\hat{\boldsymbol{\Sigma}}_{\mathbf{Y}_{di}} = \sum_g \hat{\kappa}_{dg} \hat{\boldsymbol{\varphi}}_{dig} \hat{\boldsymbol{\varphi}}_{dig}' + \hat{\sigma}_d^2 \mathbf{I}_i$  with  $\mathbf{I}_i$  denoting the identity matrix of the same dimension as the length of the vectorized response  $\mathbf{Y}_{di}$ . Compared with the projection-based estimator of the subject-specific random effects in Chen *and others* (2016) and Chen and Lynch (2017), which is only applicable for densely measured longitudinal functional two-dimensional process observed without measurement error, the proposed approach via mixed effects modeling is specifically geared towards sparse region-referenced longitudinal functional EEG data observed with measurement error. It also allows for assessing the weak separability assumption via a likelihood ratio test for the independence of the random effects (for details see Appendix D of the [supplementary material](#) available at *Biostatistics* online).

The subject-specific scores and variance components estimated via the proposed mixed effects model are used to obtain predicted subject-specific trajectories and to choose the number of eigenvectors included in the HPCA decomposition. Using subject-specific scores estimated from the mixed effects model, subject-specific trajectories can be predicted via  $\hat{Y}_{di}(r, \omega, s) = \hat{\mu}(r, \omega, s) + \hat{\eta}_d(r, \omega, s) + \sum_{g=1}^{G'} \hat{\zeta}_{dig} \hat{\varphi}_{dig}(r, \omega, s)$ , where  $G'$  denotes a set of eigenvectors such that the total fraction of variance explained ( $FVE_{dG'}$ ) is greater than 0.8 in all groups  $d = 1, \dots, D$ . We recommend starting with a

larger number  $G = KLM$  of eigencomponents in the mixed effects modeling used for the estimation of  $(\kappa_{dg} : g = 1, \dots, G)$ . In order to estimate the group-specific fraction of total variance explained via the  $G$  eigencomponents, we consider the quantity,  $FVE_{dG} = \{n_d \sum_{g=1}^G \hat{\kappa}_{dg}\} / [\sum_{i=1}^{n_d} \{ \|Y_{di}(r, \omega, s) - \hat{\mu}(\omega, s) - \hat{\eta}_d(r, \omega, s)\|^2 - R|\Omega||\mathcal{S}|\hat{\sigma}_d^2\}]$ , where  $\|f(r, \omega, s)\|^2 = \sum_{r=1}^R \int \int f(r, \omega, s)^2 d\omega ds$ . Note that the above formulation utilizes variance components estimates  $\hat{\kappa}_{dg}$  and  $\hat{\sigma}_d^2$  obtained from the proposed mixed effects model and considers the ratio of the variance in the  $G$  eigencomponents to the total variation in the observed data  $Y_{di}(r, \omega, s)$  without measurement error. The denominator of  $FVE_{dG}$  does not use variation in a large number of eigencomponents to estimate the total variation in the observed data due to computational costs in fitting the proposed mixed effects model, but instead uses the three-dimensional norm of the de-measured data, similar to the approach by [Chen and others \(2016\)](#). As a result, a limitation of  $FVE_{dG}$  is that when measurement error variance is overestimated and scaled by a factor of  $R|\Omega||\mathcal{S}|$ ,  $FVE_{dG}$  may exceed 1.

### 3.3. Group-level inference via bootstrap

To test the null hypothesis that all groups have equal means in the scalp region  $r$ , i.e.  $H_0 : \eta_d(r, \omega, s) = \eta(r, \omega, s)$  for  $d = 1, \dots, D$ , we propose a parametric bootstrap procedure based on the HPCA decomposition. The proposed parametric bootstrap generates outcomes based on the estimated model components under the null hypothesis as  $Y_{di}^b(r, \omega, s) = \hat{\mu}(\omega, s) + \hat{\eta}(r, \omega, s) + Z_{di}^b(r, \omega, s) + \epsilon_{di}^b(r, \omega, s) = \hat{\mu}(\omega, s) + \hat{\eta}(r, \omega, s) + \sum_{g=1}^{G'} \zeta_{dig}^b \hat{\varphi}_{dig}(r, \omega, s) + \epsilon_{di}^b(r, \omega, s)$  in region  $r$  and as  $Y_{di}^b(r, \omega, s) = \hat{\mu}(\omega, s) + \hat{\eta}_d(r, \omega, s) + Z_{di}^b(r, \omega, s) + \epsilon_{di}^b(r, \omega, s) = \hat{\mu}(\omega, s) + \hat{\eta}_d(r, \omega, s) + \sum_{g=1}^{G'} \zeta_{dig}^b \hat{\varphi}_{dig}(r, \omega, s) + \epsilon_{di}^b(r, \omega, s)$  in the other regions, where subject-specific scores and measurement error are sampled from  $\zeta_{dig}^b \sim \mathcal{N}(0, \hat{\kappa}_{dg})$  and  $\epsilon_{di}^b(r, \omega, s) \sim \mathcal{N}(0, \hat{\sigma}_d^2)$ . The proposed test statistic  $T_r = [\sum_{d=1}^D \int \int \{\hat{\eta}_d(r, \omega, s) - \hat{\eta}(r, \omega, s)\}^2 d\omega ds]^{1/2}$  is based on the norm of the sum of the departures of the estimated group-region shifts  $\hat{\eta}_d(r, \omega, s)$  from the estimate of the common shift across groups,  $\hat{\eta}(r, \omega, s)$ . The common region shift estimate  $\hat{\eta}(r, \omega, s)$ , under the null, is set to the point-wise average of the group-region shift estimates,  $\hat{\eta}_d(r, \omega, s)$ ,  $d = 1, \dots, D$ . We utilize the proposed parametric bootstrap to estimate the distribution of the test statistic  $T_r$ . The proposed procedure can be extended to test for equal means from specific frequency bands (i.e. subsets of  $\Omega$ ). We defer steps of the bootstrap algorithm to Appendix C of the [supplementary material](#) available at *Biostatistics* online.

## 4. APPLICATION TO THE WORD SEGMENTATION DATA

### 4.1. Data structure and methods

In our motivating word segmentation study, EEG data were recorded for 144 s using an 128 electrode HydroCel Geodesic Sensor Net for 9 TD, 13 vASD, and 19 mvASD children ranging between 4 and 12 years of age. The EEG data is divided into non-overlapping segments of 1.024 seconds, producing a maximum of 140 observable segments for each subject at each electrode. Descriptions on the pre-processing steps and the final study sample are deferred to Appendix D of the [supplementary material](#) available at *Biostatistics* online. We consider 11 regions made up of 4–7 electrodes; left and right for the temporal region (LT and RT) and left, right, and middle for the frontal, central, and posterior regions (LF, RF, MF, LC, RC, MC, LP, RP, and MP, respectively). We employ the spectral PCA described in Section 2 and Appendix A [supplementary material](#) available at *Biostatistics* online as a pre-processing procedure to reduce the spectra within each brain region to its corresponding log transformed principal power. The functional domain ranges from 0 to 50 Hz, to include the clinically defined frequency bands of delta (0–4 Hz), theta (4–8 Hz), alpha (8–15 Hz), beta (15–32 Hz), and gamma (32–50 Hz). Even though HPCA captures power dynamics across the total frequency domain, we note that the gamma band was of

Table 1. FVE of the marginal covariances for the selected eigencomponents in each diagnostic group in the 2D HPCA decomposition

TD		vASD		mvASD	
$\mathcal{R}$	$\Omega$	$\mathcal{R}$	$\Omega$	$\mathcal{R}$	$\Omega$
0.652	0.698	0.706	0.653	0.583	0.656
0.113	0.159	0.112	0.249	0.133	0.231
0.084	0.091	0.083	—	0.113	0.048
0.058	—	—	—	0.062	—
—	—	—	—	0.042	—

The number of eigencomponents are chosen to explain at least 90% FVE.

particular interest in the word segmentation study since a higher gamma power is associated with better performance in cognitive processes.

We employ a 3D HPCA to model log principal power as a function of region, frequency, and segment. Based on the 3D HPCA decomposition, we observe minimal variability in the segment dimension in both the functional fixed effects (Figures S1(b) and S5 of the [supplementary material](#) available at *Biostatistics* online) and leading marginal eigenfunctions (Figure S2(c) of the [supplementary material](#) available at *Biostatistics* online), accounting for more than 85% of the marginal segment variation in each group (Table S1 of the [supplementary material](#) available at *Biostatistics* online). Collectively, these two observations suggest that log principal power dynamics do not substantially change in the segment dimension both within subjects and among groups. Therefore, we collapse the segment dimension by averaging log principal power across segments within regions and employ a 2D HPCA decomposition to model the resulting average log principal power  $Y_{di}(r, \omega)$  as a function of region and frequency. Thus, we utilize the 3D HPCA decomposition to justify the collapse of the segment dimension allowing for a more interpretable analysis based on the 2D HPCA decomposition. Finally, we illustrate the benefit of modeling the unreduced frequency dimension by integrating the average log principal power  $Y_{di}(r, \omega)$  over clinically defined frequency bands and comparing separate linear mixed models (LMMs) of the resulting region-referenced log principal power bands with the 2D and 3D HPCA. In the LMMs, group-region dynamics are captured through group-region interactions while within-subject region variation is modeled using a subject-specific random intercept. LMMs were fit using `nlme` (Pinheiro and others, 2017). For the 2D and 3D HPCA decompositions, the smoothing parameters for the functional fixed effects and marginal covariances were selected by generalized cross-validation/restricted maximum likelihood.

#### 4.2. Data analysis results

We present full results from the 2D HPCA decomposition but defer details from the 3D HPCA decomposition, including detailed checks of the weak separability assumption on the 2D and 3D covariance processes, to Appendix D of the [supplementary material](#) available at *Biostatistics* online. Our main focus is inference on group-region differences but we will briefly discuss the estimated model components from the 2D HPCA decomposition. Table 1 displays the eigencomponents for the regional and functional marginal covariances that explain at least 90% marginal FVE in all three diagnostic groups. The leading four, three, and five regional marginal eigenvector and three, two, and three functional marginal eigenfunctions are collectively found to explain 0.998, 1.000, and 0.999 of the total FVE ( $FVE_{dG}$ ) in the TD, vASD, and mvASD groups, respectively.

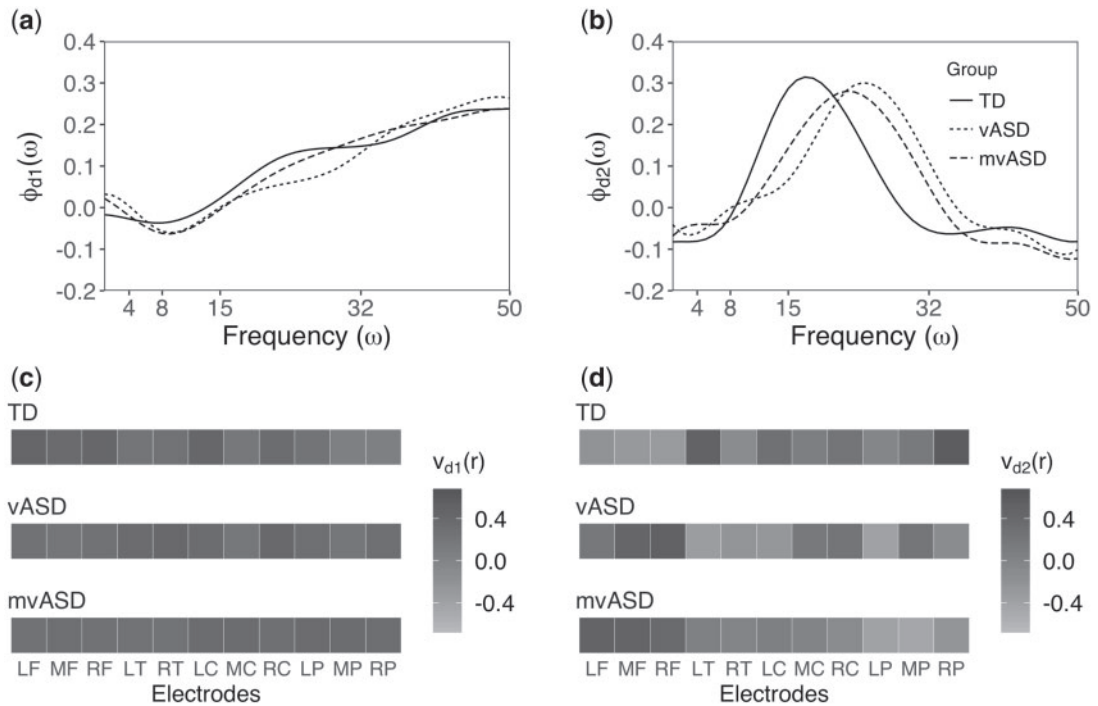


Fig. 2. (a, b) Estimated first and second leading functional and longitudinal marginal eigenfunctions  $\phi_{d1}(\omega)$  and  $\phi_{d2}(\omega)$ . (c, d) Estimated first and second leading regional marginal eigenvectors  $v_{d1}(r)$  and  $v_{d2}(r)$ .

In the functional dimension, the first leading marginal eigenfunction  $\phi_{d1}(\omega)$  (Figure 2(a)) displays increasing variation with increasing frequency for all diagnostic groups, with the peak observed in the beta and gamma bands (15–50 Hz). The second leading marginal eigenfunction  $\phi_{d2}(\omega)$  (Figure 2(b)) displays peak variation mostly in the beta band (15–32 Hz). The first two eigenfunctions together explain at least 85% of the variation in the functional marginal covariance in all three diagnostic groups. In the regional dimension, the weights of the first leading marginal eigenvector  $v_{d1}(r)$  (Figure 2(c)) are uniform across scalp locations in all the diagnostic groups, implying equal variation, while the weights of the second leading marginal eigenvector  $v_{d2}(r)$  (Figure 2(d)) are highest for the LT and RP regions, and MF and RF regions for the TD and vASD groups, respectively. In the mvASD group, the leading components signal a contrast between LF and MP regions. The first two regional marginal eigenvectors together explain at least 70% of the variation in the regional marginal covariance in all three diagnostic groups.

The estimated overall mean log principal power  $\mu(\omega)$  curve, given in Figure 1(c), follows the well known trend of decreasing power with increasing frequency. In order to test for differences in the group-region means among the three diagnostic groups, we utilize the bootstrap test proposed in Section 3.3 originally for the 3D HPCA decomposition, which can be extended to the 2D HPCA decomposition via the test statistic  $T_r = [\sum_{d=1}^D \int \{\hat{\eta}_d(r, \omega) - \hat{\eta}(r, \omega)\}^2 d\omega]^{1/2}$ . For each scalp region  $r$ , we test the null hypothesis that the three diagnostic groups share a common mean, which takes the form  $H_0 : \eta_d(r, \omega, s) = \eta(r, \omega, s)$  and  $H_0 : \eta_d(r, \omega) = \eta(r, \omega)$ ,  $d = 1, 2, 3$ , for 3D and 2D bootstrap procedures, respectively. The 2D and 3D bootstrap tests find significant differences among the group-region means for the three frontal regions: LR, RF, and MF (Figures 3(a, c, e) and S5(a–f) of the [supplementary material](#) available at *Biostatistics* online,  $P < 0.05$ ) across the full frequency domain. The 3D bootstrap procedure also identifies a significant

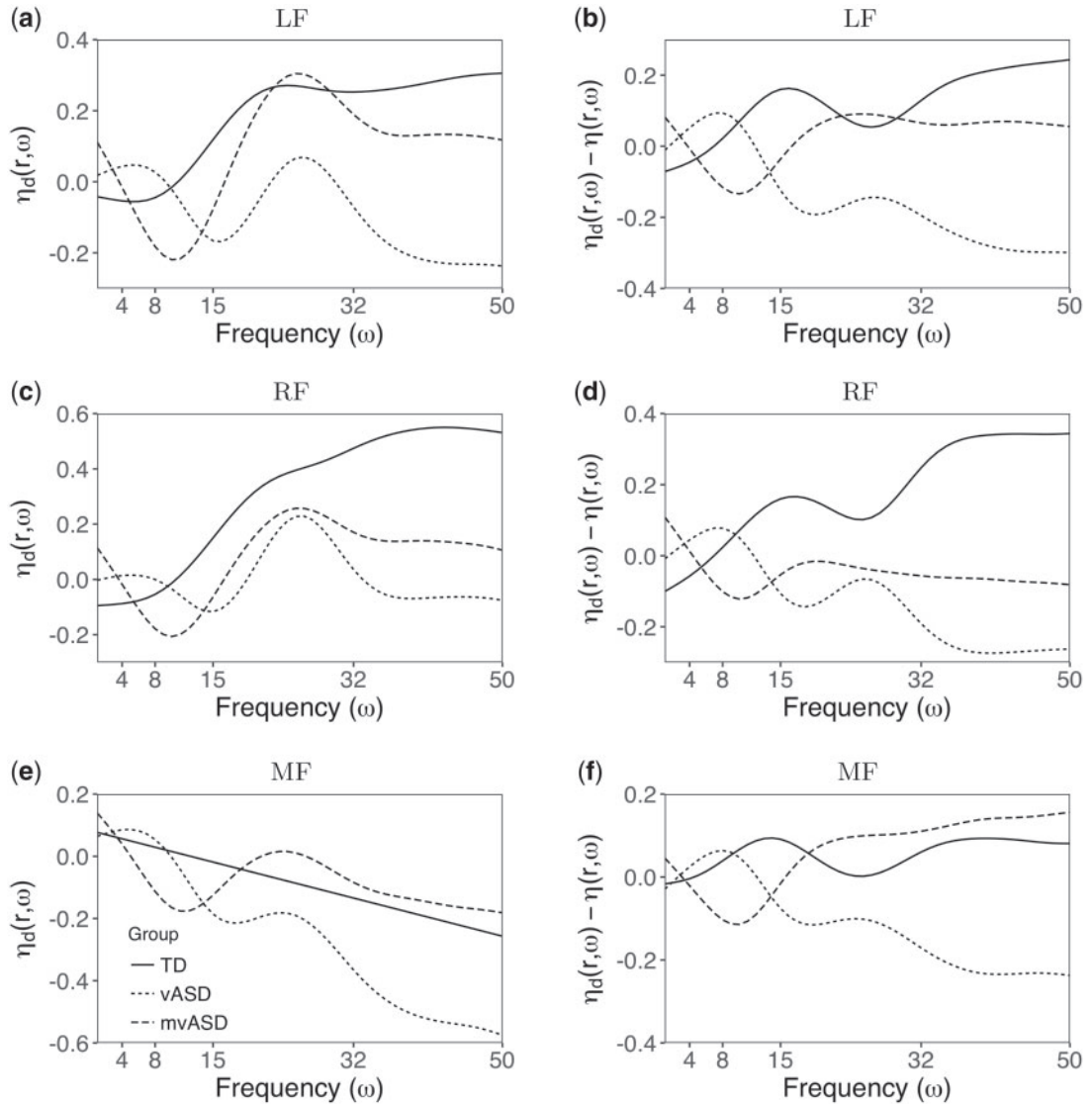


Fig. 3. (a, c, e) The estimated group-region shifts  $\eta_d(r, \omega)$  in the left, right, and middle frontal regions in the TD, vASD, and mvASD groups. (b, d, f) The differences of the estimated group-region shifts  $\eta_d(r, \omega)$  from group-region averages in the left, right, and middle frontal regions in the TD, vASD, and mvASD groups. Note, the quantity  $\eta_d(r, \omega) - \eta(r, \omega)$  forms the basis of the proposed bootstrap test statistic.

difference among the group-region means for the total frequency domain in the LT, MC, and RP, although for the LT and RP regions this may be ascribed to edge effects in the segment dimension inflating the observed test statistic and for the MC region the 2D bootstrap test is nearly significant ( $P=0.05$ ). While the 2D and 3D bootstrap tests provide insight into group-level dynamics for the full frequency domain, we also employ their band-specific extensions to enhance interpretation and enable comparisons with band-specific LMMs.

Table 2 displays the results of hypothesis tests for all scalp regions and frequency bands from the three separate models, the 2D and 3D HPCA decompositions and a set of band-specific LMMs. The greatest variation in group-region means from the 2D HPCA decomposition are observed in LF and RF regions for the gamma band (Figure 3(a, c),  $P < 0.05$ ), with the highest gamma principal power observed in the TD group, followed by the mvASD and vASD groups regions as evidenced by their relative difference from the group-region averages (Figure 3(b, d)). The mvASD group appears to have higher gamma principal power than the vASD group in the LF and LR regions, contradicting the expectation that the ordering of verbal impairment would be mirrored in group-region shifts in gamma activity, which is thought to signal cognitive processes. One reason could be that the three diagnostic groups were not age-matched. The age distribution of vASD group had minimal overlap with those of the TD and mvASD groups and was over 20 months younger on average than the other diagnostic groups, which may explain its lower gamma principal power. Further evidence of this age imbalance may be observed in LF and RF regions for the theta band in which the vASD group displays higher activity than the TD and mvASD groups across brain regions, consistent with the expected trend that theta activity is higher in younger children (Figure 3(a, c)). Finally, for each brain region the TD group followed by the vASD group have the highest alpha activity, which is thought to be associated with relaxation, suggesting that mvASD children are not as relaxed as their verbally able peers.

The 2D HPCA decomposition enhances the analysis of principal power by not only capturing the whole frequency domain but also by detecting significant differences among group-region means that are missed when the frequency domain is collapsed into specific bands and modeled via band-specific LMMs. In the MF region for the gamma band, the TD and mvASD groups display higher principal power than the vASD group (Figure 3(e, f)) and the null hypothesis of a common group-region mean is rejected by the 2D bootstrap test but not by the band-specific LMM. In addition, the 2D bootstrap procedure finds significant differences in theta band dynamics among groups in all regions but the RC region while the LMM finds no significant differences among group-region means. By collapsing the frequency dimension prior to modeling, analysis methods such as the LMM cannot capture dynamics within frequency bands among groups (e.g. two signals crossing in a given interval) that may be modeled by maintaining the full frequency dimension.

## 5. SIMULATION STUDY

We studied the finite sample properties of the proposed HPCA and the bootstrap test for group-level inference via extensive simulations. While the results of the simulations are summarized here briefly, we defer details including data generation, discussion of the total and marginal FVEs, and further details on the bootstrap test to Appendix E of the [supplementary material](#) available at *Biostatistics* online. We conducted simulations for two sample sizes ( $n_d = 15$  and 50), two signal-to-noise ratios (SNRs= 2.5 and 10), and two data sparsity levels (complete, partial in the longitudinal domain), for a total of eight settings. The lower sample size and sparsity levels were chosen to mimic the word segmentation data. To assess the performance of the proposed estimation algorithm in targeting the components of HPCA, we utilize normalized mean squared errors (MSE) and relative squared errors (RSE), based on the norms of the deviations of the estimate from the targeted quantities. In addition, we report the total and marginal FVE along the regional, functional, and longitudinal dimensions,  $FVE_{dK, \mathcal{R}}$ ,  $FVE_{dL, \Omega}$ , and  $FVE_{dM, \mathcal{S}}$ , based on the  $K$ ,  $L$ , and  $M$  marginal eigenvectors included in the decomposition, respectively.

Figures 4 and S7 of the [supplementary material](#) available at *Biostatistics* online display the estimated model components based on 200 Monte Carlo runs from the sparse simulation set-up with  $n_d = 15$  and high SNR. The estimated overall mean function and group-region shift with the median RSE values (Figures S7(b),(d) of the [supplementary material](#) available at *Biostatistics* online) track the corresponding

Table 2. From left to right within each column grouping, results from the total spectral domain and band-specific hypothesis tests for all scalp regions for the 3D bootstrap procedure, 2D bootstrap procedure, and linear mixed models (LMMs), respectively

Region	3D HPCA/2D HPCA/LMM																		
	Delta (0–4 Hz)	Theta (4–8 Hz)	Alpha (8–15 Hz)	Beta (15–32 Hz)	Gamma (32–50 Hz)	Total (0–50 Hz)	Delta (0–4 Hz)	Theta (4–8 Hz)	Alpha (8–15 Hz)	Beta (15–32 Hz)	Gamma (32–50 Hz)	Total (0–50 Hz)							
LT	<b>0.005</b>	<b>0.000</b>	<b>0.001</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.174	<b>0.000</b>	<b>0.000</b>	0.080	0.215	0.413	<b>0.025</b>	0.145	0.217	<b>0.015</b>	0.095	-	
RT	<b>0.030</b>	<b>0.000</b>	<b>0.001</b>	<b>0.000</b>	<b>0.005</b>	<b>0.001</b>	0.442	<b>0.005</b>	<b>0.000</b>	0.210	0.410	0.777	0.240	0.465	0.517	0.125	0.350	-	
LF	0.235	0.075	0.435	<b>0.000</b>	<b>0.014</b>	0.080	<b>0.014</b>	0.080	<b>0.005</b>	<b>0.026</b>	<b>0.035</b>	0.065	0.073	<b>0.000</b>	<b>0.030</b>	<b>0.038</b>	<b>0.005</b>	<b>0.020</b>	-
RF	0.125	<b>0.005</b>	0.074	<b>0.000</b>	<b>0.005</b>	0.090	0.062	0.090	<b>0.000</b>	<b>0.034</b>	<b>0.040</b>	0.105	0.179	<b>0.000</b>	<b>0.000</b>	<b>0.013</b>	<b>0.000</b>	<b>0.000</b>	-
MF	0.615	0.550	0.679	<b>0.010</b>	<b>0.000</b>	0.085	0.198	0.085	<b>0.000</b>	<b>0.033</b>	0.170	0.125	0.218	<b>0.030</b>	0.088	<b>0.025</b>	<b>0.035</b>	-	
LC	<b>0.025</b>	<b>0.020</b>	<b>0.005</b>	<b>0.010</b>	<b>0.000</b>	<b>0.020</b>	0.163	<b>0.005</b>	<b>0.002</b>	0.130	0.160	0.387	0.185	0.275	0.262	0.110	0.175	-	
MC	0.085	<b>0.045</b>	0.272	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.232	<b>0.000</b>	<b>0.029</b>	0.075	0.150	0.532	<b>0.040</b>	0.105	0.362	<b>0.025</b>	0.050	-	
RC	<b>0.000</b>	<b>0.035</b>	<b>0.001</b>	<b>0.010</b>	<b>0.010</b>	<b>0.000</b>	0.624	<b>0.000</b>	<b>0.003</b>	0.260	0.370	0.748	0.130	0.345	0.396	0.100	0.320	-	
LP	<b>0.045</b>	<b>0.015</b>	<b>0.004</b>	<b>0.005</b>	<b>0.045</b>	<b>0.000</b>	0.450	<b>0.000</b>	<b>0.002</b>	0.120	0.205	0.542	0.310	0.795	0.822	0.110	0.350	-	
MP	<b>0.025</b>	<b>0.010</b>	<b>0.040</b>	<b>0.000</b>	<b>0.015</b>	<b>0.000</b>	0.253	<b>0.000</b>	<b>0.015</b>	0.115	0.180	0.490	0.310	0.545	0.646	0.130	0.265	-	
RP	<b>0.035</b>	0.130	<b>0.029</b>	<b>0.000</b>	<b>0.010</b>	<b>0.000</b>	0.162	<b>0.000</b>	<b>0.002</b>	0.135	0.260	0.644	0.065	0.220	0.386	<b>0.020</b>	0.155	-	

*P*-values less than 0.05 are displayed in bold.



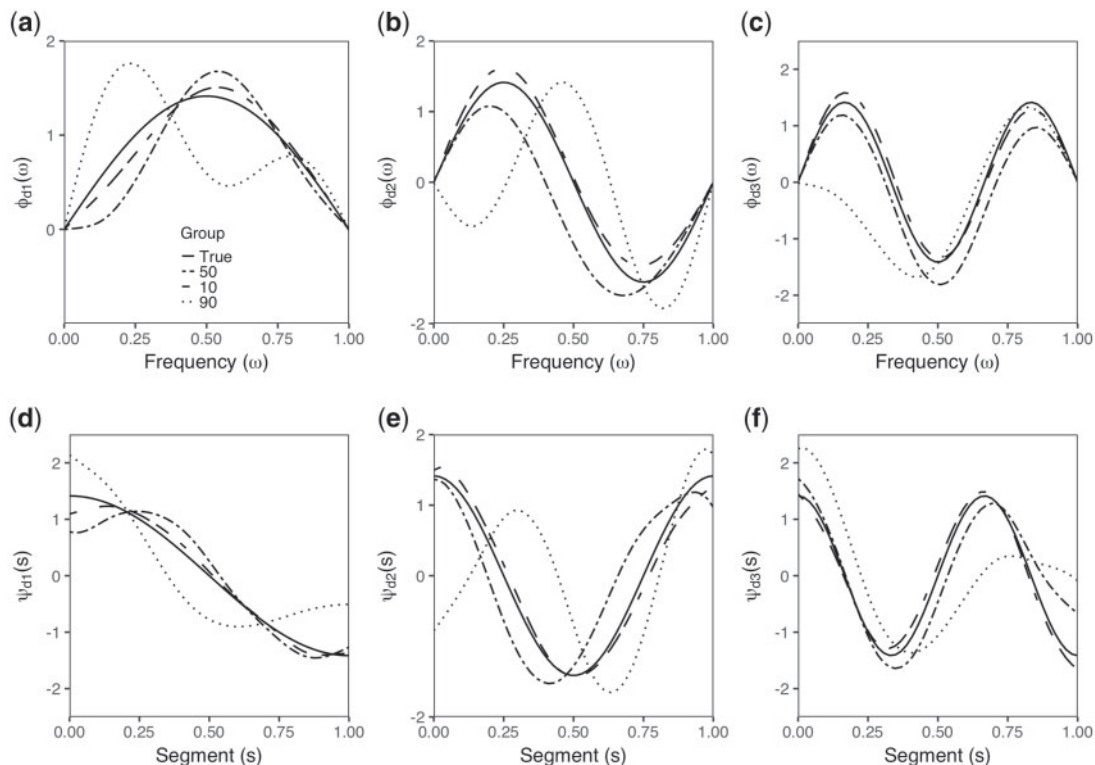


Fig. 4. The true and estimated functional (first row) and longitudinal (second row) marginal eigenfunctions corresponding to the 10th, 50th, and 90th percentile RSE values across groups based on 200 Monte Carlo runs from the sparse simulation design at  $n_d = 15$  and high SNR.

true surfaces (Figure S7(a) and (c) of the [supplementary material](#) available at *Biostatistics* online). The estimated functional and longitudinal marginal eigenfunctions (Figure 4) are displayed from runs with RSE values at the 10th, 50th, and 90th percentiles, overlaid by their true quantities. Even with a small sample size, HPCA captures the periodicity and magnitude of the true components; patterns of estimated components from the dense case are similar and are deferred to Figures S8 and S9 of the [supplementary material](#) available at *Biostatistics* online. Tables 3 and S2 of the [supplementary material](#) available at *Biostatistics* online display median, 10th, and 90th percentile RSE, normalized MSE values, and both total and marginal FVEs based on 200 Monte Carlo runs corresponding to the estimated HPCA components from all eight simulation settings. In general, the RSEs for all model components decrease with higher sample size and lower level of sparsity in the data. The fitted surfaces  $Y_{di}(r, \omega, s)$  are the most susceptible to changes in SNR. The RSEs associated with the marginal eigenvectors are not sensitive to changes in SNR, suggesting that the estimation procedure successfully corrects for measurement error when obtaining the marginal covariances. For simulation set-ups with  $n_d = 15$ , the 90th percentile RSE for the marginal eigenvectors and eigenfunctions can exceed 1 but we note that 15 subjects is small for PCA and FPCA decompositions and that for  $n_d = 50$  the comparable RSE values improve dramatically. The estimated level is on target and the power increases faster as one moves away from the null for the larger sample size, as expected (Figure S10 of the [supplementary material](#) available at *Biostatistics* online).

Table 3. Percentiles 50% (10%, 90%) of the relative squared errors, normalized mean squared errors, and both total and marginal fraction of variance explained across groups for model components based on 200 Monte Carlo runs from the sparse simulation design at  $n_d = 15, 50$  for low and high SNR

	Low SNR				High SNR			
	$n_d = 15$		$n_d = 50$		$n_d = 15$		$n_d = 50$	
$\mu(\omega, s)$	0.000	(0.000, 0.000)	0.000	(0.000, 0.000)	0.000	(0.000, 0.000)	0.000	(0.000, 0.000)
$\eta_d(r, \omega, s)$	0.004	(0.003, 0.006)	0.001	(0.001, 0.002)	0.004	(0.003, 0.006)	0.001	(0.001, 0.002)
$Y_{di}(r, \omega, s)$	0.149	(0.072, 0.155)	0.077	(0.072, 0.155)	0.016	(0.015, 0.035)	0.016	(0.015, 0.035)
$v_{d1}(r)$	0.085	(0.019, 0.422)	0.025	(0.005, 0.095)	0.075	(0.012, 0.348)	0.023	(0.003, 0.082)
$v_{d2}(r)$	0.259	(0.063, 1.166)	0.076	(0.017, 0.350)	0.236	(0.039, 1.118)	0.057	(0.011, 0.250)
$v_{d3}(r)$	0.198	(0.042, 1.025)	0.059	(0.012, 0.332)	0.159	(0.024, 0.954)	0.039	(0.005, 0.189)
$\phi_{d1}(\omega)$	0.080	(0.015, 0.376)	0.022	(0.003, 0.078)	0.074	(0.009, 0.351)	0.021	(0.004, 0.083)
$\phi_{d2}(\omega)$	0.267	(0.033, 1.109)	0.057	(0.007, 0.232)	0.231	(0.029, 0.966)	0.058	(0.008, 0.243)
$\phi_{d3}(\omega)$	0.143	(0.023, 1.073)	0.038	(0.006, 0.196)	0.137	(0.018, 0.857)	0.037	(0.004, 0.211)
$\psi_{d1}(s)$	0.092	(0.011, 0.414)	0.026	(0.004, 0.088)	0.078	(0.012, 0.330)	0.027	(0.004, 0.087)
$\psi_{d2}(s)$	0.223	(0.025, 1.130)	0.068	(0.011, 0.325)	0.212	(0.023, 1.044)	0.060	(0.009, 0.247)
$\psi_{d3}(s)$	0.143	(0.022, 0.871)	0.059	(0.012, 0.297)	0.141	(0.021, 0.899)	0.051	(0.008, 0.221)
$FVE_{dK, \mathcal{R}}$	0.985	(0.974, 0.993)	0.989	(0.983, 0.994)	0.995	(0.992, 0.998)	0.997	(0.995, 0.998)
$FVE_{dL, \Omega}$	0.974	(0.962, 0.986)	0.980	(0.972, 0.988)	0.990	(0.984, 0.994)	0.993	(0.990, 0.996)
$FVE_{dM, S}$	0.938	(0.916, 0.957)	0.963	(0.948, 0.974)	0.986	(0.979, 0.992)	0.992	(0.988, 0.995)
$FVE_{dG'}$	1.009	(0.981, 1.038)	1.011	(0.994, 1.030)	1.007	(0.995, 1.019)	1.010	(1.003, 1.016)
$\tau_{d,klm}$	0.106	(0.004, 0.769)	0.031	(0.001, 0.232)	0.101	(0.003, 0.686)	0.032	(0.001, 0.226)
$\sigma_d^2$	0.025	(0.001, 0.131)	0.005	(0.000, 0.034)	0.001	(0.000, 0.005)	0.000	(0.000, 0.002)

Due to their small magnitude, MSE values are scaled by a factor of  $10^3$  for presentation.

## 6. DISCUSSION

We proposed a hybrid principal components analysis technique (HPCA) which combines tools from vector and functional principal components analysis to decompose three-dimensional region-referenced longitudinal functional EEG data in a computationally efficient manner through the product of the one-dimensional eigenvectors and eigenfunctions of marginal covariances. Hence, the proposed estimation procedure scales up well to large data sets since estimation of the covariances and eigencomponents are performed within the marginal dimensions. To ease the computational burden in fitting the proposed mixed effects model for large data applications, the size of the grid chosen along each marginal dimension affecting the length of the design matrices can be controlled. Note also that the number of subjects in most EEG studies are similar to those in our data application, hence HPCA would be applicable in most EEG paradigms.

The proposed estimation procedure centered around weak separability was developed to specifically handle realistic scenarios observed in EEG studies with potentially sparse data in the longitudinal dimension measured with noise. Note that similar ideas can be used to handle sparsity in the functional and regional dimensions as well. The HPCA decomposition paves the way for future work on regression analysis involving the high-dimensional EEG signals. A particular question of interest in autism research

centers around relating behavioral outcomes to information within the EEG signals collected during an experiment. HPCA is a promising dimension reduction tool to enable regression modeling involving high-dimensional EEG signals.

#### SUPPLEMENTARY MATERIAL

Supplementary material is available at <http://biostatistics.oxfordjournals.org>.

#### ACKNOWLEDGEMENTS

We thank the Editor, Associate Editor, and two referees for their helpful comments. *Conflict of Interest:* None declared.

#### FUNDING

This work was supported by National Institute of General Medical Sciences [R01 GM111378-01A1 to D.S., D.T., and C.S.].

#### REFERENCES

- BRILLINGER, D. R. (1981). *Time Series: Data Analysis and Theory*. Holden-Day Series in Time Series Analysis. Holden-Day, San Francisco.
- CHEN, K., DELICADO, P. AND MÜLLER, H.-G. (2016). Modelling function-valued stochastic processes, with applications to fertility dynamics. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **79**, 177–196.
- CHEN, K. AND LYNCH, B. (2017). Weak separability for two-way functional data: concept and test. *Cornell University Library*, arXiv:1703.10210.
- CHEN, K. AND MÜLLER, H.-G. (2012). Modeling repeated functional observations. *Journal of the American Statistical Association* **107**, 1599–1609.
- DI, C.-Z., CRAINICEANU, C. M., CAFFO, B. S. AND PUNJABI, N. M. (2009). Multilevel functional principal component analysis. *The Annals of Applied Statistics* **3**, 458–488.
- GREVEN, S., CRAINICEANU, C. M., CAFFO, B. S. AND REICH, D. S. (2010). Longitudinal functional principal component analysis. *Electronic Journal of Statistics* **4**, 1022–1054.
- HASENSTAB, K., SCHEFFLER, A., TELESKA, D., SUGAR, C. A., JESTE, S., DiSTEFANO, C. AND ŞENTÜRK, D. (2017). A multi-dimensional functional principal components analysis of eeg data. *Biometrics* **73**, 999–1009.
- HUANG, L., REISS, P.T., XIAO, L., ZIPUNNIKOV, V., LINDQUIST, M. A. AND CRAINICEANU, C.M. (2017). Two-way principal component analysis for matrix-variate data, with an application to functional magnetic resonance imaging data. *Biostatistics* **18**, 214.
- LIU, C., RAY, S. AND HOOKER, G. (2017). Functional principal component analysis of spatially correlated data. *Statistics and Computing* **27**, 1639–1654.
- OMBAO, H. AND MOON-HO, R. H. (2006). Time-dependent frequency domain principal components analysis of multichannel non-stationary signals. *Computational Statistics and Data Analysis* **50**, 2339–2360.
- PARK, S. Y. AND STAIU, A.-M. (2015). Longitudinal functional data analysis. *Stat* **4**, 212–226.
- PINHEIRO, J., BATES, D., DEBROY, S., SARKAR, D. AND R CORE TEAM. (2017). *nlme: Linear and Nonlinear Mixed Effects Models*. R package version 3.1-131.
- SCOTT-VAN ZEELAND, A. A., MCNEALY, K., WANG, A. T., SIGMAN, M., BOOKHEIMER, S. Y. AND DAPRETTO, M. (2010). No neural evidence of statistical learning during exposure to artificial languages in children with autism spectrum disorders. *Biological Psychiatry* **68**, 345–351.

- STAIKU, A.-M., CRAINICEANU, C. M. AND CARROLL, R. J. (2010). Fast methods for spatially correlated multilevel functional data. *Biostatistics* **11**, 177–194.
- TAGER-FLUSBERG, H. AND KASARI, C. (2013). Minimally verbal school-aged children with autism spectrum disorder: the neglected end of the spectrum. *Autism Research* **6**, 468–478.
- YAO, F., MÜLLER, H.-G. AND WANG, J.-L. (2005). Functional data analysis for sparse longitudinal data. *Journal of the American Statistical Association* **100**, 577–590.
- ZHOU, L., HUANG, J. Z., MARTINEZ, J. G., MAITY, A., BALADANDAYUTHAPANI, V. AND CARROLL, R. J. (2010). Reduced rank mixed effects models for spatially correlated hierarchical functional data. *Journal of the American Statistical Association* **105**, 390–400.

[Received June 20, 2017; revised January 25, 2018; accepted for publication June 11, 2018]