

Adaptive Bayesian Sum of Trees Model for Covariate Dependent Spectral Analysis

Scott A. Bruce¹

Department of Statistics
Texas A&M University

Yakun Wang, Department of Statistics, George Mason University
Zeda Li, Paul H. Chook Department of Information System and Statistics,
Baruch College, The City University of New York

May 24, 2022

¹Research is supported by the National Institute Of General Medical Sciences of the NIH under Award Number R01GM140476.

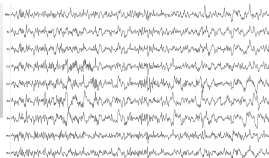
Examples of Biomedical Time Series Data

Clinicians and researchers collect a variety of [time series](#) data whose [oscillatory patterns](#) are of interest.

Examples of Biomedical Time Series Data

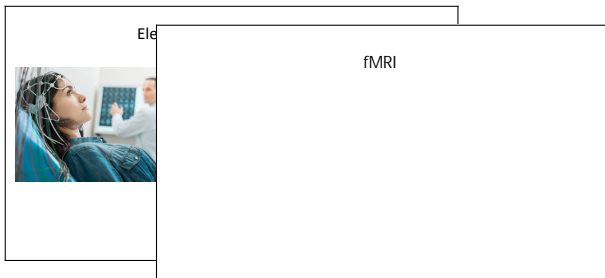
Clinicians and researchers collect a variety of **time series** data whose **oscillatory patterns** are of interest.

Electroencephalography (EEG)



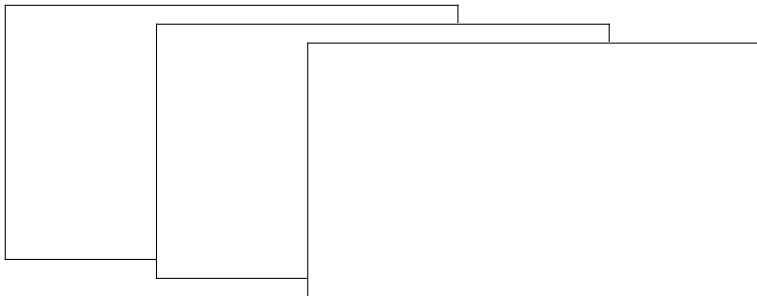
Examples of Biomedical Time Series Data

Clinicians and researchers collect a variety of **time series** data whose **oscillatory patterns** are of interest.



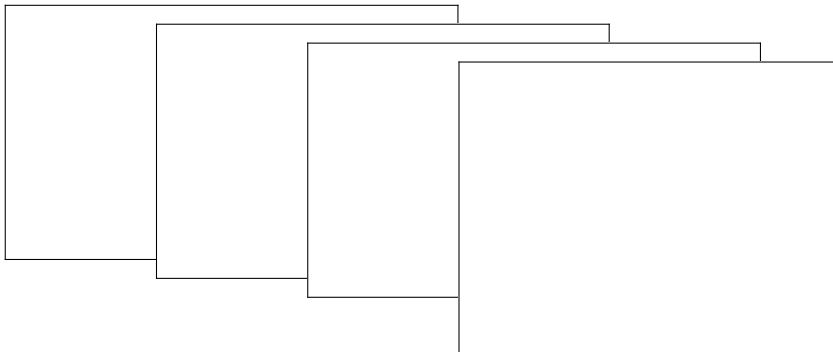
Examples of Biomedical Time Series Data

Clinicians and researchers collect a variety of **time series** data whose **oscillatory patterns** are of interest.



Examples of Biomedical Time Series Data

Clinicians and researchers collect a variety of **time series** data whose **oscillatory patterns** are of interest.



Gait Maturation Study [Hausdor et al. (1999)]

Maturation of gait dynamics

Immature gait in very young children results in unsteady walking patterns and **frequent falls**.

Gait is relatively mature by age 3. However, **neuromuscular control** continues to develop well beyond age 3.

Researchers are interested in determining if stride-to-stride dynamics continue to become more steady and regular beyond age 3.

Gait Maturation Study [Hausdor et al. (1999)]

Stride-to-stride time series data

$N = 50$ healthy children ages 3-14.

$T = 256$ stride times recorded after removing stride times in the first 60 seconds and last 5 seconds.

Age, gender, height, weight, leg length, and gait speed are also collected for each child.

Goal: To better understand the maturation of gait dynamics with age in the presence of other related covariates.

Data From Three Subjects

Stationary Time Series

Consider a zero-mean **stationary** time series X_t .

Cramer Representation [Cramer (1942)]:

$$X_t = \int_{-1/2}^{1/2} A(\omega) \exp(2i\omega t) dZ(\omega);$$

Power spectrum: $f(\omega) = |A(\omega)|^2$.

The power spectrum represents a decomposition of variance over **frequencies**.

$$\text{Var}(X_t) = \int_{-1/2}^{1/2} f(\omega) d\omega.$$

Two Simulated Examples

Two More Simulated Examples

Periodogram Estimator

Periodogram from $X = (X_1; \dots; X_T)^T$:

$$I(k) = \frac{1}{T} \sum_{t=1}^T X_t \exp(-2i!_k t) : \quad 2$$

$k = k=T, k = 1; \dots; n = bT = 2c \quad 1.$

Unbiased but noisy estimates of $f(k)$.

Approximately distributed as **scaled**² to provide the **Whittle likelihood**:

$$p(x|f) \quad (2) \quad \prod_{k=1}^n \exp \left[-\frac{1}{2} [\log f(k) + I(k) - f(k)] \right]$$

Smoothing

Periodogram can be **smoothed** to obtain a consistent estimate.

One approach Bayesian penalized linear spline [Wahba (1990)]:

$$\log f(\omega) + \sum_{s=1}^S \lambda_s \cos(2s\omega)$$

Priors [Rosen, Wood, and Stoer (2012)]

$$N(0; \sigma^2)$$

$$N(0; \sigma^2 D_S); \text{ where } D_S = \text{diag}(\frac{1}{2}, \frac{1}{2}, \dots, \frac{1}{2})$$

half-t

Sampling via Metropolis-Hastings and Gibbs steps.

Covariate-dependent Power Spectrum

Covariate-dependent Cramer Representation

$$X_{\cdot t} = \sum_{i=1}^P A(i; \cdot) \exp(2it) dZ_{\cdot}(i);$$

where $i = (i_1; \dots; i_P)^0$ is a P -dimensional covariate vector, and $\cdot = 1; \dots; L$ independent subjects.

Covariate-dependent power spectrum :

$$f(i; \cdot) = jA(i; \cdot)j^2.$$

Goal: Develop an **adaptive** method that can capture both **smooth and abrupt** changes in power spectra across multiple covariates and provide a tool for **variable selection**.

One Option: Tree-based Approach

Regression tree illustration

Tree-based models provide a flexible and parsimonious approach for partitioning multiple covariates.

For **scalar** responses, Bayesian Additive Regression Tree (BART) model [Chipman et al. (2010)]

Adaptive Bayesian Sum of Trees Model

Idea: Develop a Bayesian sum-of-trees model for $f(\mathbf{x}; \boldsymbol{\beta})$

$$\log f(\mathbf{x}; \boldsymbol{\beta}) = \sum_{j=1}^M \sum_{b=1}^{B_j} \beta_{bj} \log f_{bj}(\mathbf{x}; U_j; b);$$

M is the number of trees

U_j represents the j th tree that has B_j terminal nodes

β_{bj} is a function that identifies terminal node membership such that $\beta_{bj}(\mathbf{x}; U_j; b) = 1$ if the i th observation falls into the b th terminal node and $\beta_{bj}(\mathbf{x}; U_j; b) = 0$ otherwise.

Model specification for $\log f_{bj}(\mathbf{x})$ then follows directly from the Bayesian penalized linear spline introduced previously.

Tree Structure Priors

A regularization prior is applied to encourage each tree to be a **weak learner**:

$$\Pr(\text{SPLIT}) = (1 + d)^{-\alpha} ; \quad \alpha \in (0; 1); \quad \alpha = 2 [0; 1) ;$$

d is the depth of a tree, $\alpha = 0.95$ and $\alpha = 2$ as default.
[Chipman et al. (2010)]

Terminal node parameters and trees are assumed to be **independent** a priori.

Uniform priors on split variables and cut points.

Sparsity-inducing Dirichlet prior on split variables can also be used for improved **variable selection**. [Linero, 2018]

Sampling Scheme

Back tting Markov chain Monte Carlo (MCMC) on `residual' of periodogram

$$R_j(k) = \log L(k) - \sum_{i \in j} \sum_{b=1}^B X_i^{(b)} (\theta; U_i; b) \log f_{bi}(\cdot)$$

allows for updating each individual tree structure in turn.

Reversible jump MCMC

Birth: splitting a terminal node into two child nodes

Death: dropping two terminal child nodes belonging to the same internal node

Change: modifying the variable and cut point associated with an internal node with two terminal child nodes

Overview of Proposed Approach

Adaptively **partition covariate space** using tree structures.

Bayesian penalized spline model for **local spectra estimation** within each terminal node.

Bayesian Back tting MCMC and **Reversible jump MCMC** techniques to sample from posterior of the trees

Inference averaged over distribution of trees.

Run times

Mean run times [for a single tree update](#) over 100 replicates of the three simulation settings with $M = 5$ trees.

Simulated Abrupt+Smooth Example

$$\text{AR}(1): x_t = \rho x_{t-1} + \epsilon_t$$

$$\rho =$$

$$\left(\begin{array}{l} 0.7 + 1.4I_2 \text{ for } 0 \leq I_1 < 0.5 \\ 0.9 - 1.8I_2 \text{ for } 0.5 \leq I_1 \leq 1; \end{array} \right.$$

$$\rho = 1; \quad ; L = 100 \text{ subjects}$$

$$t = 1; \quad ; T = 250$$

$$I_1; I_2 \stackrel{i.i.d.}{\sim} U(0; 1)$$

$$\epsilon_t \stackrel{i.i.d.}{\sim} N(0; 1)$$

Estimation Accuracy for Abrupt+Smooth Simulation

Simulated Latent Variable Example

AR(2) :

$$x_t = z_1 x_{t-1} + z_2 x_{t-2} + \epsilon_t \quad \text{Latent variable mapping}$$

$$(z_1; z_2) = \begin{matrix} \text{8} \\ \text{WWW} \\ \text{WWW} \\ \text{WWW} \\ \text{WWW} \end{matrix} \begin{matrix} (1.5; 0.75); z = 1 \\ (0.8; 0); z = 2 \\ (1.5; 0.75); z = 3 \\ (0.2; 0); z = 4 \end{matrix}$$

$\epsilon_t = 1; \quad ; L = 100$ subjects

$t = 1; \quad ; T = 250$

$\epsilon_{1i}; \epsilon_{2i} \stackrel{i.i.d.}{\sim} U(0; 1)$

$\epsilon_{ti} \stackrel{i.i.d.}{\sim} N(0; 1)$

Simulated Latent Variable Example

Red line: true log power spectra

Gray points: log periodogram
ordinates

Blue line: estimated log power
spectra using the proposed
Bayesian sum of trees model

Green line: estimated log power
spectra using the competing
smooth model

Gait Maturation Analysis

Data:

$N = 50$; $T = 256$ stride-to-stride time series.

Ages 3-14 years old.

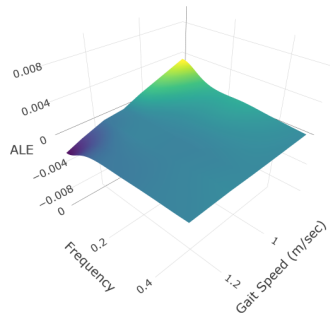
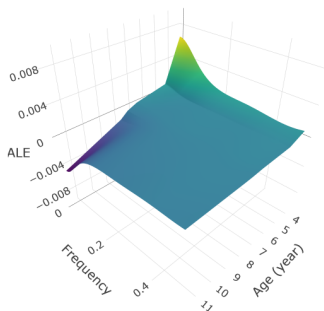
Age, gender and gait speed as covariates.

Low frequencies (LF) (0.05-0.25 stride¹) represent
uctuations over a longer-term scale (immature gait).

High frequencies (HF) (0.25-0.5 stride¹) represent
uctuations over a shorter-term scale (mature gait).

Covariate Effects on Power Spectrum

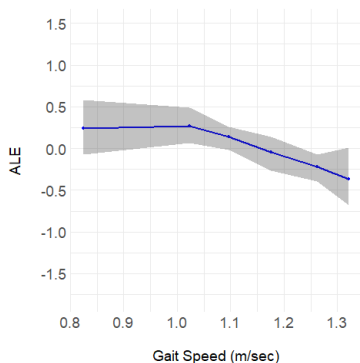
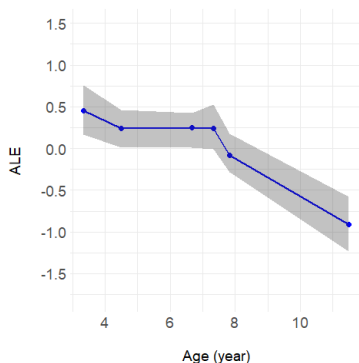
Accumulated local effects (ALE) is a method for evaluating covariate effects.



- 1 Power over all frequencies decreases as age increase.
- 2 Power in low frequencies decreases much more with age relative to higher frequencies.

Covariate Effects on Power Spectrum

Low-to-high frequency ratio $\widehat{\frac{LF}{HF}}(!) = \frac{\sum_k \mathcal{Z}(0.05; 0.25) \hat{f}(!; k)}{\sum_k \mathcal{Z}[0.25; 0.5] \hat{f}(!; k)}$.



Significant decreases for ages above 7 years and for speeds above 1 m/sec.

Current and Future Work

Current work

Proposed a nonparametric adaptive Bayesian sum of trees model for covariate-dependent spectral analysis.

Captures both **abrupt** and **smooth** changes.

Handle complex nonlinear and interaction effects.

Future work

Extend to time- and covariate-dependent time series.

Apply alternative partitioning frameworks such as Voronoi tessellations. [Payne et al., 2020]

Current and Future Work

Current work

Proposed a nonparametric adaptive Bayesian sum of trees model for covariate-dependent spectral analysis.

Captures both **abrupt** and **smooth** changes.

Handle complex nonlinear and interaction effects.

Future work

Extend to time- and covariate-dependent time series.

Apply alternative partitioning frameworks such as Voronoi tessellations. [Payne et al., 2020]

THANK YOU!

Gait Maturation Analysis

Accumulated local effects (ALE) is a method for evaluating covariate effects

ALE for $!_j = x$ on the power spectrum is

$$f_{j,ALE}(X;) = \int_{z_{0;j}}^{z_x} E_{!_{nj}!_j} \frac{f(!_j;)}{!_j} !_j = z_j dz_j \quad \text{constant}$$

$! = (!_j; !_{nj})$ where $!_j$ denotes the j th covariate and $!_{nj}$ denotes all covariates other than the j th covariate

$\mathbb{Z}_j = \{z_{0;j}, \dots, z_{H;j}\}$ is a collection of $H + 1$ partition points over the effective support of $!_j$

The constant is a value to vertically center the plot