

Functional Data Analysis of Cognitive Events in EEG

Jean-Philippe Thivierge, Ph.D.

Abstract—This paper investigates a data-driven approach to the detection of statistically reliable fluctuations in electroencephalographic (EEG) brain waves. The main advantage of this approach is the ability to automatically identify meaningful sources of variation in signals dominated by large spatiotemporal fluctuations due to noise. The approach is tested on artificial data as well as on EEG signals obtained in a cognitive experiment. Results not only show that the approach can robustly identify meaningful variations in brain activity, but also open possible applications to brain-computer interfaces based on single-trial analyses.

I. INTRODUCTION

NEUROIMAGING approaches such as EEG (electroencephalogram) offer the promise of understanding the links between neural activity and cognitive processes [1]. However, a major challenge lies in extracting meaningful information from neuroelectrical responses, typically characterized by poor signal-to-noise ratios. In EEG analysis, a well-established strategy to address this challenge is to average signals over a large number of trials [2,3]. The resulting waveform, termed “event-related potential”, can then be analyzed to identify certain signal fluctuations, termed “components”, that are statistically reliable as well as cognitively meaningful.

While several data analysis techniques are available, most possess one (or more) of the following shortcomings. First, many techniques assume that potentials can be analyzed using linear statistics (e.g., [4]), while in fact it is clear that a faithful approximation of EEG signals should rely on nonlinear methods [5]. Second, most extant approaches require some prior knowledge on the spatial location and time-course of components (c.f., [6]). Such knowledge is sometimes available, but is often imprecise and greatly limits the scope of analyses performed. Third, in currently available techniques, the detection of statistically significant components is not automated, thus introducing experimenter biases and considerably lengthening analysis times [7]. Finally, because of poor signal-to-noise ratios, several techniques are not suitable for analyzing data from single trials.

Because of these shortcomings, many techniques would not be adequate for applications based on brain-computer interfaces (BCIs). The goal of these applications is to rely on

cognitively-meaningful fluctuations in brain activity to control hardware or software devices [8]. Ideally, BCIs should be performed in real-time (i.e., on a single-trial basis) and achieve a high degrees of reliability. This is a particularly difficult task in cases where the recorded brain signals reflect highly abstract cognitive processes whose waveforms vary in sometimes unpredictable fashions.

The approach introduced here aims to address the above points. It combines nonlinear smoothing and Functional Data Analysis [9] to autonomously identify statistically significant fluctuations in EEG signals without requiring any a priori knowledge on the spatiotemporal distribution of these fluctuations.

We first provide an overview of the approach, followed by results on an artificial data set. Then, we apply the technique to a large population of signals obtained from a previously-published cognitive experiment [10]. Results demonstrate the usefulness of the technique in identifying well-documented EEG components. In both data sets, we discuss some preliminary data where the proposed technique is applied to single-trial EEGs, of potential importance for real-time applications, including BCIs.

II. METHODS

The proposed analysis of EEG signals can be broken down into three distinct steps. These general steps are identical regardless of whether the signals analyzed are obtained from single trials or averages over many trials. First, raw EEG data is fitted by a temporally continuous function that is smoother (i.e., exhibits less fluctuations) than the original data. Second, first- and second-order derivatives of this function are used to detect peaks (i.e., points characterized by large signal fluctuations). Third, a statistical test determines which, among all peaks identified, represent statistically reliable fluctuations, and are unlikely to be caused solely by noise. We now describe each of the above three steps in more detail. Further information and Matlab software is available at the following website: <http://www.psych.mcgill.ca/faculty/ramsay/ramsay.html>.

A. B-spline smoothing

From raw data y_j representing signal amplitude sampled at discrete time steps $t_j, j=1, \dots, n$, a smooth approximation $x(t)$ is obtained through b-spline smoothing (other smoothing approaches, not examined here, could also be explored). The goal of this procedure is to reduce noise levels by a function that does not overfit the raw data, yet represents its main components of variations. A linear summation

Manuscript received March 15, 2007. This work was supported in part by funding from the Fonds Québécois de Recherche sur les Natures et Technologies (FQRNT).

J.P.T. is with the Département de Physiologie, Université de Montréal, QC Canada. Phone: 514-342-5759; fax: 514-343-6113; e-mail: jean-philippe.thivierge@umontreal.ca.

$\phi_k, k=1, \dots, K$ of b-spline bases is used to generate this function:

$$x(t) = \sum_k^K c_k \phi_k(t), \quad (1)$$

where c_k are smoothing coefficients. These coefficients are adjusted in order to minimize the following cost function, representing the sum of squared errors between raw data and weighted bases ϕ_k :

$$SSE_\lambda(y|c) = (y - \Phi c)'(y - \Phi c) + \lambda \times PEN(x), \quad (2)$$

where the K -vector c contains the coefficients c_k . The term $PEN(x)$ controls the smoothness of the approximated function by penalizing the 2nd order derivatives of $x(t)$:

$$PEN(x) = \int [\partial^2 x(t) / (\partial t)^2]^2 dt. \quad (3)$$

The term λ controls the amount of penalty to be applied. As a heuristic rule (which may not apply to all EEG data), we set $\lambda = 2n$; through trial-and-error, this rule was found to generate a cognitively-plausible number of components in the data analyzed (further analyses of the influence of λ are discussed below).

B. Automatic Detection of Components

Once an approximated function $x(t)$ is obtained, temporal intervals that contain peaks are defined by two endpoints $[t_1, t_2]$ where:

$$\begin{aligned} dx(t_1)/dt_1 &= 0, \\ dx(t_2)/dt_2 &= 0, \text{ and} \\ d^2x(t)/(dt)^2 &< 0 \quad \forall x(t) \in [t_1, t_2], \end{aligned} \quad (4)$$

that is, where t_1 and t_2 correspond to zero-crossings in the 1st derivative, and where the 2nd derivative is negative (see Fig.1A-C).

C. Statistical reliability of EEG peaks

Components identified according to the criteria of (4) can be subjected to a statistical test. First, the 2nd derivative of $x(t)$ is surrounded by confidence bands determining upper ($b_{max}(t)$) and lower ($b_{min}(t)$) 95% limits on statistical reliability [9]. Second, the minima immediately preceding (t'_1) and following (t'_2) a peak (t_{peak}) are defined as follow:

$$\begin{aligned} d^2x(t'_1)/d^2t'_1 &= 0, \\ d^2x(t'_2)/d^2t'_2 &= 0, \end{aligned} \quad (5)$$

where

$$\begin{aligned} d^2x(t)/(dt)^2 &< 0 \quad \forall x(t) \in [t'_1, t_{peak}], \\ d^2x(t)/(dt)^2 &> 0 \quad \forall x(t) \in [t_{peak}, t'_2]. \end{aligned} \quad (6)$$

A peak in the signal is considered statistically significant if the joint probability of it being larger than the two adjacent

minima t'_1 and t'_2 is greater than 95%:

$$p[b_{max}(t'_1) > b_{min}(t_{peak})] \times p[b_{max}(t'_2) > b_{min}(t_{peak})] > 0.95.$$

Computationally, this is equivalent to checking that the lower confidence band of a peak is larger than both the upper confidence bands of the adjacent minima:

$$b_{min}(t_{peak}) - b_{max}(t'_1) > 0, \text{ and}$$

$$b_{min}(t_{peak}) - b_{max}(t'_2) > 0. \quad (7)$$

Finally, an important question pertains to determining the temporal latency of EEG peaks with respect to the onset of a particular stimulus. Confidence intervals can be employed towards this end. Starting from the temporal location of a peak, confidence bands describing both forward and backward intervals of time are obtained by locating the times (both before and after a peak) where the amplitude of the upper confidence band is closest to the amplitude of the signal at peak (Fig.1D).

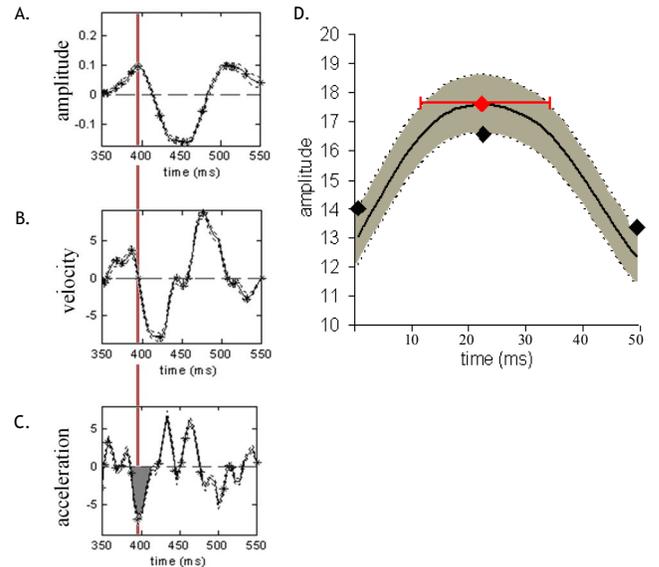


Fig. 1. Example of EEG signal (A-C). The red line identifies the temporal location of a statistically significant peak in the signal, as detected by a zero-crossing in the velocity, and negative peak in acceleration. The shaded area in the acceleration (C) marks the duration of the peak. (D) Example of the use of confidence bands to determine the statistical significance of a peak. The three points used in Equation 7 (corresponding to upper and lower markers on the confidence bands) are indicated by black diamonds. The red diamond identifies the highest point in the peak, and the red line indicates locations on the upper confidence band that delimit temporal intervals of confidence. Signal amplitude is measured in microvolts (μV).

III. ARTIFICIAL PROBLEM

In order to test the capabilities and limitations of the above-described approach, an artificial data set was devised. These data are characterized by a Gaussian function (centered at $t=300$ ms) combined with additive noise. All data were sampled at a rate of 1000Hz for a duration of 550ms. Different conditions varied the signal-to-noise ratio by

introducing different amounts of uniformly-distributed noise (Fig.2A). Each condition combined 1000 trials where noise was sampled independently. Signals from these trials were averaged together prior to running our analysis.

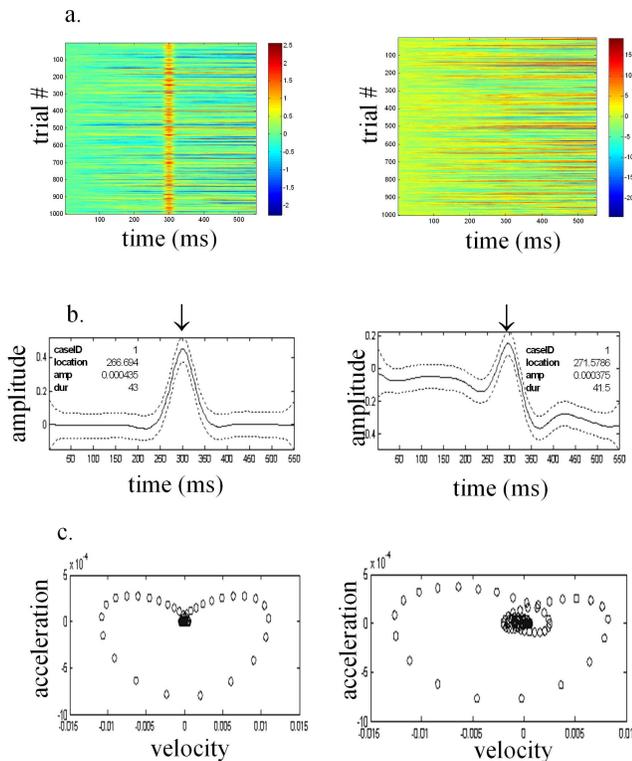


Fig. 2. Two conditions that varied signal-to-noise ratios by adding uniformly-distributed noise to a Gaussian signal centered at $t=300$ ms. Left column: 10:1 signal-to-noise ratio; right column: 1:10 signal-to-noise ratio. (A) Different trials that contained a Gaussian signal embedded in noise. (B) Population analysis of averaged trials. Dotted lines represent confidence intervals for statistical reliability, as described in (6). Arrows indicate the temporal location of statistically significant components. (C) Phase-plane plots showing the influence of noise on the main components of variation in the velocity and acceleration of averaged signals.

As results clearly demonstrate, our technique successfully identified the peak of the Gaussian function along with its location, duration, and amplitude, even in a condition where noise fluctuations were ten-fold greater than signal fluctuations (i.e., 1:10 signal-to-noise ratio). (A full description of how the locations, durations, and amplitudes of signals are obtained is beyond the scope of the current paper.) Despite the fact that noise disturbed the average waveform (c.f., Fig.2B, left vs. right columns), our analysis robustly detected the embedded Gaussian function. Phase-plane plots (Fig.2C) demonstrate that the main components of signal variation are not drastically affected by noise. Because noise seems to have its greatest effect on small fluctuations in signal amplitude, it is efficiently smoothed out by our technique.

Interestingly, the proposed approach may also potentially be applied to single-trial analyses, albeit with a certain cost in

reliability as noise levels are increased (see Fig.3); in fact, detection rates for true positives (successful detection when a Gaussian signal was present) and false positives (detection when a Gaussian signal was not present) are directly related to noise levels. However, it is encouraging that true positives outweigh false positives by a factor of two even in the case where signal-to-noise ratios are 10:1.

Of course, the artificial problem described here represents an idealized scenario in many respects. For instance, the use of a uniform noise distribution does not capture phase variations in the peak amplitudes of components, which is a typical characteristic of EEG data [11]. In order to examine more closely whether our proposed approach can extract meaningful information from EEG data, the next section turns to a more realistic data set.

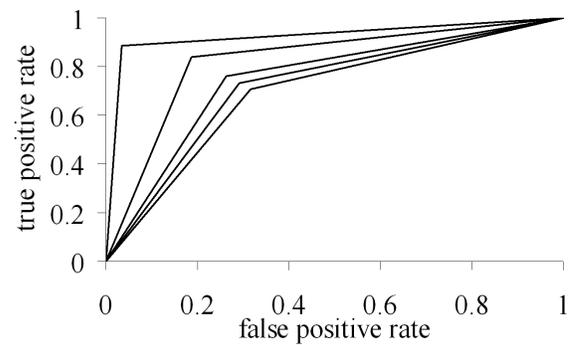


Fig. 3. Receiver-operant characteristic (ROC) curve of trials in which a statistically reliable peak was detected. Increases in the proportion of noise (from top curve, signal-to-noise ratio of 10:1, to bottom curve, signal-to-noise ratio of 1:10) deteriorate detection performances; as a result of increased noise, less components were detected in trials where they were present (true positives), and more components were detected in trials where they were absent (false positives).

IV. COGNITIVE DATA SET

A second data set, comprised of EEG signals collected during a cognitive task, was analyzed using our proposed approach. These data were obtained from an experiment where participants were presented with sentences (e.g., “*The train is never on time*”) followed by different probe words that were unrelated to the sentence (e.g., “*table*”) [10]. The data used here were pre-processed using digital filtering and manual rejection of trials containing artifacts. The time window of recording was 550ms following the probe word; sampling was performed at a rate of 1000Hz.

An example of possible analysis (many analyses were performed, comparing subject groups, electrodes, and various experimental conditions; these are beyond the scope of the current work) is presented in Fig.4. This analysis was performed over a signal averaged over trials, and restricted to the P4 electrode (located over the parietal cortex). The proposed approach successfully removed a large portion of the fluctuations present in the original signal (Fig.4A), yet retained the main components of variation (Fig.4B). A

phase-plane plot relating signal velocity and acceleration (Fig.4C) shows that there are indeed only a few large components of variation in the data. Among these components, three were found to be statistically reliable (located around 200ms, 300ms, and 400ms). These components can be roughly associated with cognitive processes known to occur around these times, including stimulus discrimination (N200), memory updating (P300), and semantic analysis (N400) [12]. In particular, the semantic component detected here (i.e., N400, a negative peak in the EEG waveform appearing 400ms after presentation of the probe word) is confirmed by previous analyses using the same set of data [10], and relates to the processing of implausible probe words.

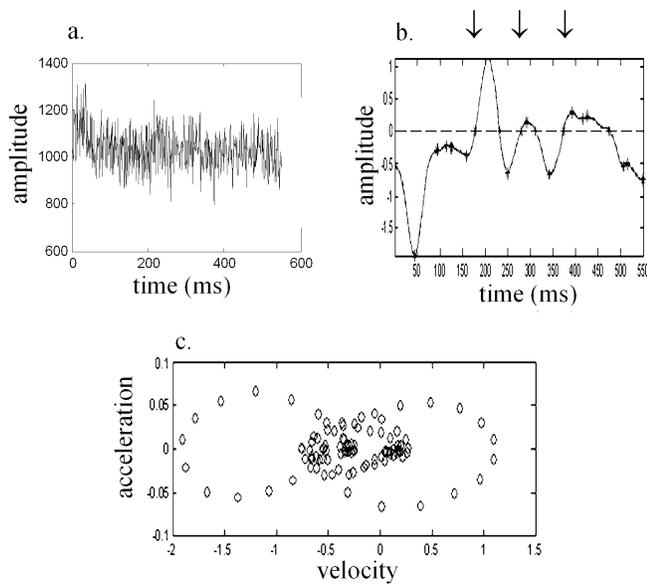


Fig. 4. Population analysis of EEG data from [10]. (A) Raw data of the P4 electrode averaged over participants (N=19) and trials (N=150). (B) B-spline smoothing of the signal in a; arrows indicate the temporal location of statistically significant components. (C) Phase-plane plot relating variations in velocity and acceleration.

Of course, the above results are sensitive to the particular degree of smoothing that is applied to the original signal, as controlled by the choice of λ parameter. A relatively low value of λ (e.g., $\lambda=10$, Fig.5A) results in a larger number of peaks detected, and broad frequency spectrum. Increasing λ to a larger value diminishes the number of peaks detected, and limits the frequency spectrum to lower bands (Fig.5B-D, summarized in Fig.5E). The choice of an adequate value of λ requires some *a priori* expectations of the approximate number of peaks expected in a given interval of time following stimulus onset. In the psycholinguistic literature, three main peaks of interest are associated with cognitive processes within 550ms following stimulus onset, reflecting a shift in attention (N100), a surprise to unexpected stimuli (P300), and a response to semantically implausible words

(N400) [2,3]. A value of $\lambda=10^5$ yields a number of peaks that approximates these *a priori* expectations.

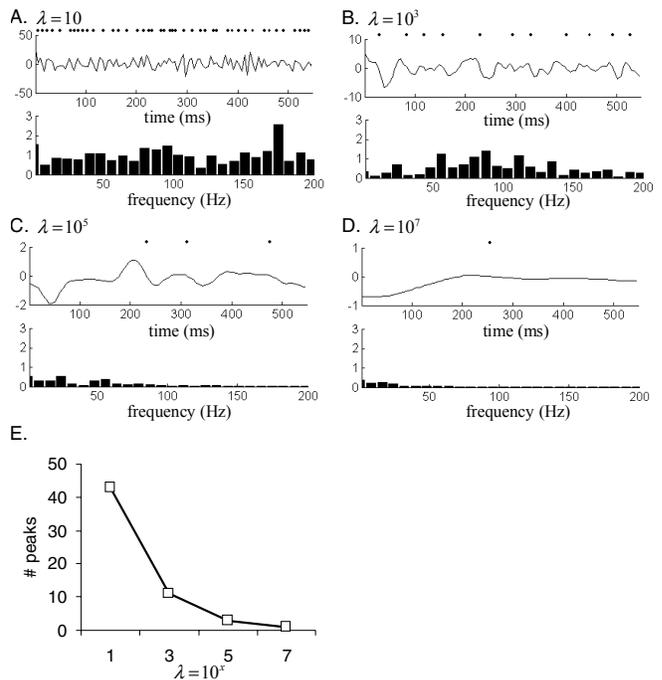


Fig. 5. Different settings of the λ parameter affect the degree of smoothing performed, and alter the number of significant peaks detected. (A-D) Top figures show signal amplitudes following stimulus onset (dots above figures show locations of significant peaks). Bottom figures show frequency decomposition performed by Fourier transform. E. Number of peaks detected decreasing as a function of the λ smoothing parameter.

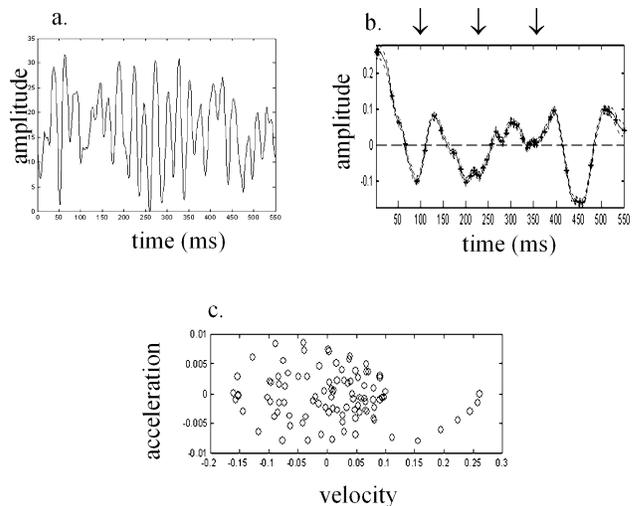


Fig. 6. A representative example of single-trial analysis. (A) Raw data from a single electrode (P4), single subject, and single trial. (B) B-spline smoothing of the signal in a; arrows indicate the temporal location of statistically significant components. (C) Phase-plane plot relating variations in velocity and acceleration.

Our analysis of artificial data (Section III) opened the possibility of single-trial analysis; we now explore this possibility further using the EEG data obtained from the cognitive task. As shown in Fig.6, the proposed approach retrieved statistically reliable components from single-trial EEG data. While in general it is not expected that reliability will be sufficient for BCI applications (c.f., Fig.3), where signal-to-noise ratios can reach beyond 1:10, it is nonetheless a step in the right direction; basing BCI applications on single-trial analyses of highly abstract cognitive processes is likely to remain one of the most important challenges of the field in the years to come.

V. CONCLUSIONS

This paper introduced a novel approach to the analysis of EEG signals, and presented a proof of principle of its application to the extraction of meaningful components in EEG data.

The proposed method overcomes several shortcomings of related approaches. First, it fully automates the detection of statistically reliable peaks in the EEG waveform, thus reducing potential experimenter biases. This is a clear advantage that many approaches are incapable of at present [13]. Second, it requires no prior knowledge of the temporal location and spatial distribution of components, an advantage characteristic of only a handful of other methods [6,14]. Finally, it treats EEG signals as temporally continuous and nonlinear, thus remaining more faithful to the data analyzed. At the current time, it is not clear how this last point advantages our technique compared to other statistical approaches available (e.g., [15-19]); more direct comparisons will be required to investigate this question.

Other areas of future exploration include the incorporation of blind source separation [20] in order to identify the topological epicenters of EEG components. In its current form, our approach focuses on the temporal evolution of EEG waveforms; a wealth of additional information could be gained by analysis of their spatial evolution.

Finally, further developments are required to provide reliable BCIs based on real-time analyses of EEG signals. Our results show some promises in this respect: some of the peaks detected in signals averaged across trials and subjects (e.g., ~400ms after stimulus onset, Fig.4) are also detected in single-trial EEGs (Fig.6). This result suggests that some of the sources of fluctuation detected in analyses of averaged signals, and reflecting fundamental cognitive processes such as the integration of semantic information, can also be uncovered with the presentation of a single stimulus. While this is a necessary step towards real-time BCIs, a major challenge still to be met is to reduce the computational burden required by our method. Indeed, the analysis of a single signal (regardless of whether it is averaged over trials or not) requires about 25 seconds of processing time on an Intel 1.7GHz PC. The solution to this challenge may reside in the development of more efficient algorithms both for smoothing

raw signals and detecting significant peaks.

In closing, we emphasize that the proposed method is not a stand-alone solution to all problems in EEG analysis (and certainly not in BCIs). The complex nature of the relationship between neural activity and cognition may be better tackled by an approach that combines several techniques, each focusing on a different aspect of this relationship.

ACKNOWLEDGMENT

This work benefited from ongoing collaborations with Frederic Dandudand, Thomas R. Shultz, and Natalie Phillips. The author is grateful to Jim Ramsay for useful discussions.

REFERENCES

- [1] S. Karakas, O. U. Erzen, and E. Basar, "A new strategy involving multiple cognitive paradigms demonstrates that ERP components are determined by the superposition of oscillatory responses," *Clin Neurophysiol*, vol. 111, pp. 1719-32, 2000.
- [2] T. Picton, *Human event-related potentials. Handbook of electroencephalography and clinical neurophysiology*. Amsterdam: Elsevier, 1988.
- [3] F. Lopes da Silva, "Event-related potentials: methodology and quantification," in *Electroencephalography*, E. Neidermeyer and F. Lopes da Silva, Eds., 4th ed. Philadelphia: Lippincott Williams & Wilkins, 1999.
- [4] F. Carbonell, L. Galan, P. Valdes, K. Worsley, R. J. Biscay, L. Diaz-Comas, M. A. Bobes, and M. Parra, "Random field-union intersection tests for EEG/MEG imaging," *Neuroimage*, vol. 22, pp. 268-76, 2004.
- [5] S. Carrubba, C. Frilot, A. Chession, and A. A. Marino, "Detection of nonlinear event-related potentials," *J Neurosci Methods*, vol. 157, pp.39-47, 2006.
- [6] N. J. Lobaugh, R. West, and A. R. McIntosh, "Spatiotemporal analysis of experimental differences in event-related potential data with partial least squares," *Psychophysiology*, vol. 38, pp. 517-30, 2001.
- [7] R. A. Dobie, "Objective response detection," *Ear Hear*, vol. 14, pp. 31-5, 1993.
- [8] C. W. Chen, C. C. Lin, and M. S. Ju, "Detecting movement-related EEG change by wavelet decomposition-based neural networks trained with single thumb movement," *Clin Neurophysiol*, 2007.
- [9] J. Ramsay and B. Silverman, *Functional data analysis, 2nd Edition*. New York: Springer-Verlag, 2005.
- [10] N. A. Phillips and D. Lesperance, "Breaking the waves: age differences in electrical brain activity when reading text with distractors," *Psychol Aging*, vol. 18, pp. 126-39, 2003.
- [11] J. Mocks, T. Gasser, and T. Pham Dinh, "Variability of single visual evoked potentials evaluated by two new statistical tests," *Electroencephalogr Clin Neurophysiol*, vol. 57, pp. 571-80, 1984.
- [12] C. L. Yang, C. A. Perfetti, and F. Schmalhofer, "Event-related potential indicators of text integration across sentence boundaries," *J Exp Psychol Learn Mem Cogn*, vol. 33, pp. 55-89, 2007.
- [13] D. Melkonian, T. D. Blumenthal, and R. Meares, "High-resolution fragmentary decomposition--a model-based method of non-stationary electrophysiological signal analysis," *J Neurosci Methods*, vol. 131, pp. 149-59, 2003.
- [14] E. Duzel, R. Habib, B. Schott, A. Schoenfeld, N. Lobaugh, A. R. McIntosh, M. Scholz, and H. J. Heinze, "A multivariate, spatiotemporal analysis of electromagnetic time-frequency data of recognition memory," *Neuroimage*, vol. 18, pp. 185-97, 2003.
- [15] J. Yordanova, V. Kolev, O. A. Rosso, M. Schurmann, O. W. Sakowitz, M. Ozgoren, and E. Basar, "Wavelet entropy analysis of event-related potentials indicates modality-independent theta dominance," *J Neurosci Methods*, vol. 117, pp. 99-109, 2002.
- [16] R. C. Blair and W. Karniski, "An alternative method for significance testing of waveform difference potentials," *Psychophysiology*, vol. 30, pp. 518-24, 1993.
- [17] K. J. Friston, K. M. Stephan, J. D. Heather, C. D. Frith, A. A. Ioannides, L. C. Liu, M. D. Rugg, J. Vieth, H. Keber, K. Hunter, and

- R.S. Frackowiak, "A multivariate analysis of evoked responses in EEG and MEG data," *Neuroimage*, vol. 3, pp. 167-74, 1996.
- [18] L. Galan, R. Biscay, J. L. Rodriguez, M. C. Perez-Abalo, and R. Rodriguez, "Testing topographic differences between event related brain potentials by using non-parametric combinations of permutation tests," *Electroencephalogr Clin Neurophysiol*, vol. 102, pp. 240-7, 1997.
- [19] H. Utku, O. U. Erzenin, E. D. Cakmak, and S. Karakas, "Discrimination of brain's neuroelectric responses by a decision-making function," *J Neurosci Methods*, vol. 114, pp. 25-31, 2002.
- [20] A. Tang, M. Sutherland, and Y. Wang, "Contrasting single-trial ERPs between experimental manipulations: improving differentiability by blind source separation," *Neuroimage*, vol. 29, pp. 335-46, 2006.