



**HAL**  
open science

# Detecting single-trial EEG evoked potential using a wavelet domain linear mixed model: application to error potentials classification

J Spinnato, M-C Roubaud, Boris Burle, Bruno Torr sani

► **To cite this version:**

J Spinnato, M-C Roubaud, Boris Burle, Bruno Torr sani. Detecting single-trial EEG evoked potential using a wavelet domain linear mixed model: application to error potentials classification. *Journal of Neural Engineering*, 2015, 12 (3), pp.036013. 10.1088/1741-2560/12/3/036013 . hal-01161911

**HAL Id: hal-01161911**

**<https://hal.science/hal-01161911v1>**

Submitted on 10 Jun 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destin e au d p t et   la diffusion de documents scientifiques de niveau recherche, publi s ou non,  manant des  tablissements d'enseignement et de recherche fran ais ou  trangers, des laboratoires publics ou priv s.

# Detecting Single-trial EEG Evoked Potential using a Wavelet Domain Linear Mixed Model: *Application to Error Potentials classification*

J Spinnato<sup>1,2</sup>, M-C Roubaud<sup>1</sup>, B Burle<sup>2</sup> and B Torr sani<sup>1</sup>

<sup>1</sup> Aix-Marseille Universit , CNRS, Centrale Marseille, I2M, UMR 7373, 13453 Marseille, France

<sup>2</sup> Aix-Marseille Universit , CNRS, LNC, UMR 7291, 13331 Marseille, France

E-mail: juliette.spinnato[AT]univ-amu.fr

## **Abstract.**

*Objective.* The main goal of this work is to develop a model for multi-sensor signals such as MEG or EEG signals, that accounts for the inter-trial variability, suitable for corresponding binary classification problems. An important constraint is that the model be simple enough to handle small size and unbalanced datasets, as often encountered in BCI type experiments.

*Approach.* The method involves linear mixed effects statistical model, wavelet transform and spatial filtering, and aims at the characterization of localized discriminant features in multi-sensor signals. After discrete wavelet transform and spatial filtering, a projection onto the relevant wavelet and spatial channels subspaces is used for dimension reduction. The projected signals are then decomposed as the sum of a signal of interest (i.e. discriminant) and background noise, using a very simple Gaussian linear mixed model.

*Main results.* Thanks to the simplicity of the model, the corresponding parameter estimation problem is simplified. Robust estimates of class-covariance matrices are obtained from small sample sizes and an effective Bayes plug-in classifier is derived.

The approach is applied to the detection of error potentials in multichannel EEG data, in a very unbalanced situation (detection of rare events). Classification results prove the relevance of the proposed approach in such a context.

*Significance.* The combination of linear mixed model, wavelet transform and spatial filtering for EEG classification is, to the best of our knowledge, an original approach, which is proven to be effective. This paper improves on earlier results on similar problems, and the three main ingredients all play an important role.

## 1. Introduction

Electro- and Magneto-encephalography (respectively, EEG and MEG) are of the rare techniques allowing non-invasive brain investigation with an excellent temporal resolution and, under some conditions, a fairly good spatial one. One main limitation, however, is that extracting the brain activity of interest from background activity and noise usually requires averaging a large number of repetitions of the signal recorded in the “same” condition (for example, averaging several epochs of signal following the same repeated stimulus). Such an averaging, however, distorts the signal and prevents precise investigation of the dynamics of the underlying processes. Indeed, it has long been known that averaging has non-linear impact on the latencies estimation [12, 30, 11]. For example, the latency of the onset of an activity measured on the average largely underestimates the real mean of the individual onsets [52, 41, 32], making its use problematic for direct comparison with chronometric variables such as Reaction Time (RT) [11, 41]. The averaging process also prevents from analyzing learning and/or adaptation effects across trials repetitions [46]. More generally averaging trials eliminates the signal of interest variability, although the latter contains important information that can be useful in various contexts, from signal interpretation to classification.

For these reasons, methods for single-trial EEG analysis and classification have been developed during the last decades. To reduce the impact of noise and background activity in single-trial analysis, common strategies have been to extract more elementary parts concentrating the signal of interest, either in time, in frequency, or in space. Such approaches include, among others, selection of time domains, frequency filtering, wavelet decomposition [45, 56], adaptive basis selection [54, 4], matching pursuit [15, 5] or blind source separation [29]. In those approaches, one aims at simply getting rid of the variability induced by the background activity and the noise, hoping that the remaining variability will only be attributable to the signal of interest.

Alternative strategies can be based on explicit signal modelling. For example the Linear Discriminant Analysis (LDA) which is one of the most popular classifiers for EEG single-trials detection (and for brain computer interface - BCI - applications, see [33]) can be interpreted in terms of very simple Gaussian mixture model. In this setting LDA assumes a class-independent covariance matrix. In a recent overview [7], Blankertz and coworkers proposed a regularized Gaussian mixture model for event related potentials (ERPs). In the latter, the single-trial is written as the sum of a signal of interest which is approximated as constant over trials and a random background activity modelled as a Gaussian noise. The inter-trial variability is not taken into account. The class-covariance matrices are assumed to be equal, which leads to a very simple detection algorithm that turns out to be very efficient for binary classification tasks.

However when this equality assumption is not valid, the problem becomes more complex, especially when the two classes are unbalanced. In that context, standard classifiers, such as LDA, may fail as the estimated common covariance matrix is largely

determined by the majority class [24, 61]. Taking into account the difference of the class-covariance matrices leads to a quadratic classification rule (QDA) and requires the estimation of a covariance matrix for each class. This leads to a robustness problem in the estimates, the more so when the size of the minority class is small.

In BCI and more generally in EEG experiments, such an unbalanced situation is not unfrequent. For example in BCI, the P300 speller protocol naturally generates two unbalanced classes [18]. Classical P300 spellers, with a  $6 \times 6$  matrix containing all letters and characters, yield an unbalanced datasets composed of 1/6 of ERP signals and of 5/6 of noERP. Another example of unbalanced classes can be found in RT tasks in which the participant has to respond as fast as possible to the appearance of predetermined stimulation. For example, standard experimental psychology protocols make use of biased probabilities across experimental conditions (see e.g. [43]). Besides manipulated factors, biased probability might also be the result of participants behavior, such as errors which are typically much lower than correct trials.

In this paper, we target more specifically situations where the class-covariance matrices differ and available datasets are unbalanced and of limited size. We propose a wavelet domain Gaussian linear mixed model (termed LMM in the following) for the binary signal classification problem. The model expresses each single-trial as a sum of a class-dependent signal of interest and a background activity as in [28, 27]. The signal of interest is modelled as a multivariate Gaussian vector whose covariance matrix, that describes the inter-trial variability, depends on a user-specified design matrix. Both mean and design matrix are class-dependent. The background signal is modelled as a class-independent Gaussian white noise. The resulting model is characterized by a remarkably small number of parameters.

The application context of the current paper is the detection of evoked potentials in M/EEG signals. The above described model turns out to be relevant for such signals when suitable preprocessings are performed, namely wavelet based time decorrelation, spatial filtering and corresponding dimension reductions. The procedure is more specifically applied to the problem of error negativity (ErrP) detection and analysis. Corresponding classification results compare very favorably with standard linear classifiers for small sample size datasets.

One of the main contributions of this work is a model for subject specific M/EEG signals, that belongs to the family of linear mixed models. Mixed models offer a rich and flexible framework for describing and quantifying variability and produce robust estimates of class-covariance matrices from small datasets. Such models have been considered in EEG applications in many different contexts. Let us for example mention [3] where an introduction to mixed model is provided and possible applications in psychology and neuroimaging are discussed. Mixed models have also recently been used by other authors in bayesian framework for electrophysiology applications. In [14], a functional mixed model is introduced for the purpose of regression analysis of ERP's data and where the model is combined with discrete wavelet transform and sparsity-inducing prior distribution. In [19] a subject-independent classifier is proposed using a

$\ell_1$ -penalized linear mixed model. Within-subject and between-subject variabilities are modelled through random and fixed terms and the method is applied to BCI achieving subject-to-subject transfer. The present work, inspired on the work of Huang and coworkers [28, 27], focuses on a subject specific model taking into account the variability across repetitions of the same experiment. The model proposed here is constructed directly in wavelet domain instead of time domain, and also differs from [28, 27] in both fixed and random parts. In our case the fixed part corresponds to the class-average whereas they used projection onto the first principal components. The modelling of inter-trial variability is different too. All together this results in a simplified Gaussian linear mixed model with less parameters to estimate. Moreover the application of the model to the binary classification problem is different. In our work we directly exploit the estimates of class-covariance matrices in a Bayes classifier while Huang *et al* [28, 27] use a likelihood test based on the predicted single-trial signals. Let us stress that we only address and model here the inter-trial variability. Even though the dataset used for validation involves several subjects, they are all treated independently and we do not attempt to model the inter-subject variability.

This paper is organized as follows. A short overview for linear and quadratic discriminant methods for unbalanced data is given in section 2. The model and the statistical procedure are described in section 3 while section 4 is devoted to the application to ErrP data. Discussion and conclusions are given in section 5 and section 6 respectively.

## 2. A short overview of linear and quadratic discriminant analysis for unbalanced data

We briefly discuss here the main issues encountered when using linear and quadratic discriminant analysis for unbalanced and small size datasets, and describe some classical solutions that will be of interest for this work.

### 2.1. Linear and Quadratic discriminant analysis

Let us take the probabilistic point of view, and derive Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA) from the decision problem in the case of a Gaussian mixture, see e.g. [25]. Observations are considered as multivariate Gaussian draws with respective prior probabilities  $p^c$  (with  $\sum_c p^c = 1$ ) and probability density function (pdf):

$$f^c(x) = \frac{1}{(2\pi)^{d/2}|\Sigma^c|^{1/2}} \exp \left[ -\frac{1}{2}(x - \mu^c)'(\Sigma^c)^{-1}(x - \mu^c) \right] , \quad (1)$$

where  $\mu^c \in \mathbb{R}^d$  is the class-mean,  $\Sigma^c \in \mathbb{R}^{d \times d}$  is the class-covariance matrix (assumed to be invertible) and  $|\Sigma^c|$  its determinant.

This leads to the discriminant function, defined for each class  $c$ :

$$\delta^c(x) = -\frac{1}{2} \log |\Sigma^c| - \frac{1}{2}(x - \mu^c)'(\Sigma^c)^{-1}(x - \mu^c) + \log p^c . \quad (2)$$

Let  $x_i$  denote the  $i$ -th sample to classify.  $x_i$  a  $d$ -dimensional vector. The decision rule is given by assigning  $x_i$  to the class with the maximal  $\delta^c(x_i)$ . In the general case, the class-covariances  $\Sigma^c$  are different, and the decision should be based on quadratic functions of  $x_i$  (thus the name QDA). When the class-covariances are assumed to be equal ( $\Sigma^c = \Sigma$  for all  $c$ ) the decision is based on linear functions of  $x_i$  (LDA).

In practice the parameters of the Gaussian distribution in each class are unknown and must be estimated, which is problematic in the case of small size and/or high dimensional datasets. We address this problem below, focusing on binary classification ( $c \in \{0, 1\}$ ) and unbalanced datasets. We will denote by  $N^0$  and  $N^1$  the corresponding sample sizes of datasets from which the estimation is done, and set  $N = N^0 + N^1$ .  $\hat{\Sigma}^c$  and  $\hat{\Sigma}$  will denote respectively the sample estimates of class-covariance matrix  $\Sigma^c$  and common covariance matrix  $\Sigma$  defined as follows:

$$\hat{\Sigma}^c = \frac{1}{N^c - 1} \sum_{i=1}^{N^c} (x_i - \bar{x}^c)(x_i - \bar{x}^c)', \quad \text{where} \quad \bar{x}^c = \frac{1}{N^c} \sum_{i=1}^{N^c} x_i, \quad (3)$$

$$\hat{\Sigma} = \frac{N^0}{N} \hat{\Sigma}^0 + \frac{N^1}{N} \hat{\Sigma}^1. \quad (4)$$

Two cases are to be considered.

*Case 1: Equal class-covariance matrices.* When  $\Sigma^0 = \Sigma^1 = \Sigma$ , the decision is based upon the sign of  $w'x_i$ , where  $w = \Sigma^{-1}(\mu^0 - \mu^1)$ . The quality of the decision relies heavily on the quality of the estimation of the covariance matrix  $\Sigma$  and its inverse, which may turn out to be poor for small datasets and/or in high-dimensional situations. In addition, low quality estimate for  $\Sigma$  may lead to invertibility problems, which makes the classifier ill-defined. Solutions to such problems are discussed below.

Let us note that for equal class-covariance matrices, unbalancedness of the datasets does not introduce additional difficulty as  $\Sigma$  is estimated from the complete dataset.

*Case 2: Unequal class-covariance matrices.* In this case, two covariance matrices must be estimated, which amplifies the above mentioned difficulty, particularly when the dataset is small and unbalanced enough so that at least one of the two covariance matrices is poorly estimated. This is why LDA is generally preferred, even when the true class-covariances are different. However, when the datasets are unbalanced, the estimation tends to focus on the prevalent class and to ignore the rare one: the covariance estimate  $\hat{\Sigma}$  given in (4) is dominated by the majority class, so that  $\hat{\Sigma} \approx \hat{\Sigma}^0$  (in the case where  $N^1 \ll N^0$ ).

In such a situation, a common solution consists in re-balancing datasets using sampling methods [61, 62, 57]. Over-sampling increases the size of the minority class (without information gain) while under-sampling reduces the size of the majority class (which may result in prejudicial information loss).

## 2.2. The case of diagonally dominant covariance matrices

We stick here to binary classification problem, in the small sample size ( $d \gg N$ ) and unbalanced ( $N^1 \ll N^0$ ) situation. In such case, estimated covariance matrices are often non-invertible, and some additional assumptions or regularization are needed. We review a few solutions that will be of interest for us, that ensure invertibility by enforcing diagonal dominance of the covariance matrices.

- (i) Diagonal covariance. If data are assumed to be decorrelated, the covariance matrix is diagonal and only the  $d$  variances are to be estimated (leading to Diagonal LDA, DLDA for short). Under the multivariate Gaussian law, this assumption corresponds to independence assumption. However it is known that even for correlated data (which is the general situation) better classification results are often obtained when correlations are ignored, thus replacing  $\hat{\Sigma}$  by its diagonal (leading to the so-called naive Bayes classifier). This is supported by asymptotic arguments ( $d \gg n$ ) as well as experimental results (see e.g. [6] and references therein). In particular if the off-diagonal elements of the covariance matrix are expected to be nearly zero, it is usual to estimate them by zero. Otherwise in some situations, a decorrelating transformation can be found that enforces the diagonal dominance of the covariance matrix (see section 3.1).

The diagonal assumption as well as the above considerations can be extended to the quadratic situation (leading to the diagonal QDA, DQDA for short).

- (ii) Diagonal dominance enforcing regularization. Singularity issues can be overcome by shrinking the sample covariance matrix towards a multiple of the identity matrix [21]. A particular trace preserving shrinkage approach has been proposed in [7] in the context of EEG single-trial analysis. The corresponding estimate is a weighted average of the sample covariance matrix and the identity matrix

$$\hat{\Sigma}(\gamma) = (1 - \gamma)\hat{\Sigma} + \gamma \frac{\text{tr}(\hat{\Sigma})}{d} \mathbf{I}_d, \quad (5)$$

where  $\gamma \in [0, 1]$  is a shrinkage parameter and  $\text{tr}(\hat{\Sigma})$  is the trace of the sample covariance matrix, i.e. the sum of its diagonal elements. The standard practice is to find the optimal  $\gamma$  by cross-validation. Blankertz *et al* give a simple heuristic estimate with satisfactory results on EEG data [7]. This method will be called Regularized LDA (RDA for short).

- (iii) Explicit covariance modelling. When prior information on data is available, explicit models for covariance matrices can sometimes be proposed to reduce the number of parameters to estimate. We describe below an example where a suitable transformation combined with a linear mixed model leads to diagonal dominant class-covariance matrices characterized by  $(2 + 2 \times d)$  parameters only.

### 3. Method

We propose a classification method for multisensor signals to discriminate between two experimental conditions in view of application to EEG and BCI. The signals are recorded over  $M$  sensors in a fixed time-period with  $T$  samples. For a given participant, we obtain  $N^c$  trials per class,  $c \in \{0, 1\}$  being the class label, with each trial taking the form of  $M$  time series of  $T$  samples each.

Our procedure is based on three steps. The first step is preprocessing: transformations are applied to obtain diagonally dominant sample class-covariance matrices and the temporal and spatial dimensions are reduced. The second step models variability using a Gaussian linear mixed model to estimate the class-covariance matrices. The third step is classification.

#### 3.1. Decorrelation and dimension reduction

Raw recorded signals are known to be strongly correlated both across sensors and time points. Data are thus preprocessed in order to reduce dimensions and enforce the diagonal dominance of class-covariance matrices.

*3.1.1. Time-domain decorrelation using DWT.* Signals are recorded as time-series, with a given sampling frequency. However, samples are expected to be highly correlated, should they correspond to the background activity or to the signal of interest. Therefore, it is important to reduce the correlations, which we do through a linear transform, by replacing the time samples with the coefficients of the signals expansion with respect to a suitably chosen basis. In this work, we limit ourselves to a wavelet basis, generated by shifting and rescaling a generic waveform  $\psi$ . More precisely, in the framework of the so-called *multiresolution analysis*, it is possible to find pairs of functions  $(\phi, \psi)$  such that suitably rescaled and shifted copies of these form orthonormal bases of signal spaces of interest. Without going into mathematical details, for which we refer to e.g. [13, 36, 55] for detailed accounts, and also to [10] for a pedestrian introduction, this yields multiscale signal decompositions as follows. Given a reference (integer) scale index  $m_0$ , any signal in the signal space can be uniquely expanded as time-shifted copies of the scaling function at scale  $2^{m_0}$  and time-shifted copies of the wavelet, rescaled at smaller scales  $2^m$ ,  $m \leq m_0$  (also called *levels*). This is expressed mathematically in the form

$$f(t) = \sum_n s_{m_0,n} 2^{-m_0/2} \phi(2^{-m_0}t - n) + \sum_{m \leq m_0} \sum_n d_{m,n} 2^{-m/2} \psi(2^{-m}t - n) .$$

The coefficients that enter such an expansion are called respectively *scaling coefficients* (for the coefficients  $s_{m_0,n}$ ) and *wavelet coefficients* (coefficients  $d_{m,n}$ ) and are easily computed as inner products of the signal  $f$  with the corresponding basis functions. The numerical computation of these coefficients is performed using dedicated fast algorithms (see [55]), which leads to the so-called DWT (for discrete wavelet transform). In what follows, we will gather scaling and wavelet coefficients of a signal in a unique vector of *multiscale coefficients*.



Wavelets are introduced here because of their ability to decorrelate signals. It has been found empirically in many application domains (such as image processing) that the nonzero (or significant) multiscale coefficients of a correlated signal are often far less correlated than the signal, and that the number of significant coefficients is therefore much smaller than the signal length. This turns out to be the case for the EEG signal considered here. Therefore, a (time-domain) dimension reduction can be performed by moving to the multiscale domain using DWT, and setting to zero irrelevant multiscale coefficients, namely those coefficients that are numerically negligible for all channels. The resulting coefficients are therefore weakly correlated (and the corresponding covariance matrix is strongly diagonal dominant). It is worth noticing that wavelets only provide approximate decorrelation. Alternatives to wavelets have been proposed for improving the decorrelation by seeking optimally decorrelating basis (see [59] for detailed account of the best basis approach and [54, 4] for the adaptation to EEG signals). These techniques are adaptive and, as a consequence, the decomposition basis is signal dependent, which is not suitable for the model we are using here.

In addition to dimension reduction, a level-dependent rescaling is performed on the retained multiscale coefficients to correct for variability differences across scales. More precisely, at each scale, coefficients are normalized so that they all have the same variance.

As a result of this time (or time-scale) domain processing, each trial  $i$  gives rise to a multiscale coefficient matrix  $X_i^c \in \mathbb{R}^{M \times K}$  ( $c$  being the class label), where  $K$  is the number of retained multiscale coefficients and  $M$  being the number of sensors.

*3.1.2. Spatial filtering.* The considered signals are multisensor signals, with low spatial resolution and large spatial correlations. These can be reduced using spatial filtering. In single-trial analysis, spatial filters are used to improve signal-to-noise ratio while reducing both EEG data complexity and spatial dimension. Various spatial filter constructions have been proposed, depending on EEG application (see [42, 38] for reviews).

In the present study, we rely on a matrix-based technique inspired by Fisher's linear discriminant method [20], as done in [27]. The main purpose is to identify projections in the sensor space that keep the classes as separated as possible while minimizing the variance within classes. We introduce the between-class and the within-class covariance matrices  $S_B$  and  $S_W$ , respectively

$$S_B = \frac{1}{NT} \sum_{c=0}^1 N^c (\bar{X}^c - \bar{X}) (\bar{X}^c - \bar{X})', \quad (6)$$

$$S_W = \frac{1}{NT} \sum_{c=0}^1 \sum_{i=1}^{N^c} (X_i^c - \bar{X}^c) (X_i^c - \bar{X}^c)', \quad (7)$$

where  $\bar{X}^c \in \mathbb{R}^{M \times K}$  and  $\bar{X} \in \mathbb{R}^{M \times K}$  are the trial-average in class  $c$  and the total trial-average matrices. Here,  $N = N^0 + N^1$  is the total number of trials. Notice that  $S_B$  and  $S_W$  are both  $M \times M$  non-negative definite matrices.

The method seeks to identify the most discriminant linear combinations of sensors by optimizing the following Fisher criterion, sequentially under orthogonality constraints between vectors:

$$\max_{u \in \mathbb{R}^M} \frac{u' S_B u}{u' S_W u}. \quad (8)$$

The solution is obtained by computing the eigenvectors decomposition of the matrix  $S_W^{-1} S_B$  (assuming that  $S_W$  is invertible). The eigenvectors are sorted by decreasing discriminating power (eigenvalue), out of which the first  $J$  are selected. The number  $J$  is chosen according to a criterion based on the cumulative percentage of eigenvalues. In most situations,  $J \ll M$ , which yields a significant spatial dimension reduction. The  $J$  uncorrelated linear combinations of the  $M$  sensors will be called *channels* in the following.

Let  $U = [u_1, \dots, u_J]$  denote the matrix whose columns are the selected eigenvectors. Each single-trial signal  $X_i^c$  is projected onto the selected subspace spanned by these  $J$  eigenvectors as follows

$$Y_i^c = U' X_i^c, \quad (9)$$

where each row  $j$  of the matrix  $Y_i^c \in \mathbb{R}^{J \times K}$  consists in the wavelet and scaling function coefficients of the  $j$ -th channel for single-trial  $i$ . Each column of  $U$ , denoted  $u_j$ , can be referred to as *spatial filters* (see [38, 8, 7, 42]).

### 3.2. Model setup

In this section a single-trial analysis is proposed, in which the variability is modelled through a Gaussian linear mixed model [47, 37].

*3.2.1. Background and definition.* Given multisensor and multi-trial signals  $f_{i,s}^c(t)$ , with  $i$  the trial index and  $s$  the sensor index, that are labelled by a class index  $c$ , we assume that such signals can be modelled as a sum of two independent components:

- A background activity, assumed to be stationary (in the sense of weakly stationary random processes, namely the two first order moments are translation invariant) and correlated (see e.g. [31]),
- Event related components (also called signals of interest), thus intrinsically non-stationary, characterized by a *class index* which labels the cerebral response to the events. The signal of interest is further modelled as the sum of a fixed component, common to all the trials in the class, and a random trial dependent component (the so-called *trial random effect*).

In what follows, the signals that will be modelled in such a way are the multi-channel multiscale (wavelet and scaling) coefficient arrays  $Y_i^c \in \mathbb{R}^{J \times K}$  resulting from the preprocessing. We will denote by  $y_i^c \in \mathbb{R}^{JK}$  the vector obtained by concatenation of all columns of  $Y_i^c$ .

3.2.2. *The statistical model.* For a given participant, we consider  $N^c$  trials per class,  $c \in \{0, 1\}$  being the class label. After preprocessing, we represent each trial  $i$  as a vector  $y_i^c \in \mathbb{R}^{JK}$ , where  $J$  is the number of channels and  $K$  the number of retained wavelet and scaling coefficients. In a given class  $c$  and according to the additivity assumption (see section 3.2.1 above),  $y_i^c$  ( $i = 1, \dots, N^c$ ) is therefore written as

$$y_i^c = \mu^c + \Gamma^c b_i + \varepsilon_i, \quad (10)$$

where

- $\mu^c \in \mathbb{R}^{JK}$  is the mean vector of class  $c$ ,
- $b_i \sim \mathcal{N}(0, \tau^2)$  is a zero-mean Gaussian random variable with variance  $\tau^2$ , modelling the trial random effect,
- $\Gamma^c \in \mathbb{R}^{JK}$  is a class-dependent coefficient vector which modulates the impact of the trial on each sensor and sample,
- $\varepsilon_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{JK})$  is the residual part modelling the background activity where  $\mathbf{I}_{JK}$  denotes the identity matrix of size  $JK$ . We further assume that the residual vector  $\varepsilon_i$  and the trial random effect  $b_i$  are independent.

After specification of the coefficient vector  $\Gamma^c$  (see section 3.2.3 below), the model given in (10) becomes a Linear Mixed Model (LMM), such that

$$y_i^c | \Gamma^c \sim \mathcal{N}(\mu^c, V^c), \quad \text{where} \quad V^c = \tau^2 \Gamma^c (\Gamma^c)' + \sigma^2 \mathbf{I}_{JK}, \quad (11)$$

(here  $(\Gamma^c)'$  denotes the matrix transpose of  $\Gamma^c$ ).

In summary, each single-trial is decomposed into a fixed part, corresponding to the class-mean, a trial random effect  $b_i$  modulated by the coefficient vector  $\Gamma^c$ , and a random term accounting for the background activity.

Let us note that the simplicity of the proposed model in wavelet domain permits to obtain precise parameter estimates over a small number of trials. Indeed, after specifying the coefficient vectors  $\Gamma^c$ , we only have to estimate four parameters: the two mean vectors  $\mu^c$ ,  $c \in \{0, 1\}$ , and the variances  $\tau^2$  and  $\sigma^2$ .

**Remark 1** (i) Since the class-mean is used as the fixed part, the model described in (10) reduces to a mixed-effects ANOVA model [51], where EEG-class is the fixed effect factor and trial is the random effect factor. This allows grounding the model on very solid foundations, and yields a more straightforward interpretation of the results.

- (ii) In the case of the between-trial variability is negligible, there is no random effect  $\Gamma^c b_i$  in (10) and the model corresponds to a simple linear one, as described in [7].
- (iii) A related approach, that actually inspired the present work, was proposed by Huang *et al* in [28, 27]. Our modelling, however, differs in three major points. First, we consider a linear mixed model on wavelet domain more adapted to the assumption of the noise decorrelation. Second, the fixed part is set to the class-mean and does not involve any projection. Third, the covariance matrix depends on

channel, sample and class while Huang *et al* propose dependence on trial, sample and class. We stress that our model is significantly simpler, avoids unnecessary approximations and requires less parameter estimations.

- (iv) Sticking to such a simple model is motivated by use cases where few trials are available. This is the case under some experimental conditions where the number of trials of interest is low, either because of the experimental design (e.g. probability bias), or because rare behaviour of the participants (e.g. errors in RT tasks). Another example is BCI type data, for which the training set has to be limited to reasonable size, and to which high complexity models therefore can hardly be fitted.

*3.2.3. Random part coefficient vector.* To completely specify the model in (10), the (class-dependent) coefficient vectors  $\Gamma^c \in \mathbb{R}^{JK}$  have to be chosen. Let us note that the components of these vectors modulate the random effects  $b_i$  in the same way for all trials in the same class.  $\Gamma^c$  is a design parameter to be fixed *a priori*. The model itself doesn't assume any specific form, the choice of  $\Gamma^c$  is problem dependent, and usually relies on prior information. In the considered situation, such prior information is not available, the choice of  $\Gamma^c$  was guided by preliminary exploration of the data. We refer to section 4.2.2 for details.

### 3.3. Classification procedure

Each single-trial  $y_i \in \mathbb{R}^{JK}$  to classify is assigned to a class  $c$  according to the maximum *a posteriori* calculated using the Bayes formula:

$$\max_{c \in \{0,1\}} \mathbb{P}(\text{Class} = c | y_i). \quad (12)$$

Based on the model proposed in (10),  $y_i$  is assumed to be distributed in each class  $c$  according to the multivariate Gaussian distribution  $\mathcal{N}(\mu^c, V^c)$  where  $\mu^c$  is the class-mean and  $V^c$  is the class covariance matrix such that  $V^c = \tau^2 \Gamma^c (\Gamma^c)' + \sigma^2 \mathbf{I}_{JK}$ .

Under these assumptions, the optimal classification is obtained with the following quadratic discriminant function for each class  $c$  :

$$g_c(y_i) = (y_i - \mu^c)' (V^c)^{-1} (y_i - \mu^c) - 2 \ln(p^c) + \ln |V^c|, \quad (13)$$

which yields to the following binary classifier :

$$\delta(y_i) = g_0(y_i) - g_1(y_i). \quad (14)$$

The trial  $y_i$  is assigned to class 0 if  $\delta(y_i) < 0$ , to class 1 otherwise.

**Remark 2** Let us note that the resulting classifier is a particular case of the QDA classifier presented in section 2.1. In the situation considered here, the class-covariance matrix  $\Sigma^c$  is decomposed in the following simple form  $\tau^2 \Gamma^c (\Gamma^c)' + \sigma^2 \mathbf{I}_{JK}$ . Therefore, after specifying the coefficient vector  $\Gamma^c \in \mathbb{R}^{JK}$  for each class, only  $\tau^2$  and  $\sigma^2$  have to be estimated, which makes QDA tractable in this situation.

### 3.4. A specific methodology for EEG single-trial classification

We shortly summarize the main steps of the LMM-based classification procedure (hereafter termed LMMC), which consists in two stages: a training step to estimate the model parameters and the test step to classify the EEG single-trials. For the reasons exposed above, in many contexts as BCI, the training set shall be as small as possible to get sufficiently precise estimates to correctly classify test set.

#### 3.4.1. Training Step.

(i) **Preprocessing.**

First whitening and time-domain reduction are performed through DWT. The  $K$  selected multiscale coefficients are scaled to set the coefficient variability independent of the decomposition level. Then a spatial dimension reduction is performed. Spatial filters are estimated using a matrix-based Fisher's linear discriminant and only the  $J$  most discriminant channels are selected.

(ii) **Wavelet-based modelling.**

Single-trials in the multiscale domain are modelled as in (10), with the class-mean as a fixed component and the random component depending on class, channel and trial. In the latter, the class and channel dependence is given by the coefficient vectors  $\Gamma^0$  and  $\Gamma^1$  that must be specified a priori, according to (18) in this model.

(iii) **LMM parameter estimation.**

After fixing  $\Gamma^c$  a priori, the complete parameter  $\theta = (\mu^0, \mu^1, \tau^2, \sigma^2) \in \mathbb{R}^{2JK+2}$  of the LMM can be estimated by likelihood-based methods using the assumption that  $b$  and  $\varepsilon$  are independent and normally distributed. We use Restricted Maximum Likelihood (REML) method that provides unbiased variance estimates (see e.g. [51]). The estimates will be denoted by  $\hat{\theta} = (\hat{\mu}^0, \hat{\mu}^1, \hat{\tau}^2, \hat{\sigma}^2)$ .

#### 3.4.2. Test Step.

(i) **Preprocessing.**

The transformations are the same as in the training step.

(ii) **Classification rule.**

The classifier  $\hat{\delta}(y_i)$  is obtained from the plug-in estimates of the discriminant functions given in (13). In the expressions (11) to (14) the parameters  $\mu^0, \mu^1, \tau^2, \sigma^2$  are replaced by their estimates computed during the training step.

From the variance component estimates  $\hat{\tau}^2$  and  $\hat{\sigma}^2$ , we obtain the estimation of the class-covariance matrix  $V^c$  described in the model (11) :

$$\hat{V}^c = \hat{\tau}^2 \Gamma^c (\Gamma^c)' + \hat{\sigma}^2 \mathbf{I}_{JK} . \quad (15)$$

### 3.5. Single-trial reconstruction

In addition, an estimate  $\hat{b}_i$  can be obtained for each single-trial  $i$  in class  $c$  through the Henderson mixed model equations (see e.g. [51]) :

$$\hat{b}_i = \hat{\tau}^2(\Gamma^c)'(\hat{V}^c)^{-1}(y_i - \hat{\mu}^c). \quad (16)$$

This, in turns, yields an estimate for the signal of interest as follow :

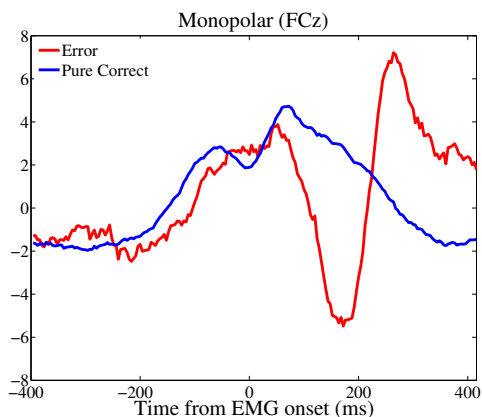
$$\hat{y}_i^c = \hat{\mu}^c + \Gamma^c \hat{b}_i. \quad (17)$$

Notice that unlike Huang *et al* [28, 27], this estimation is a sub-product of the model allowing single-trial visualization and does not play any role in the classification procedure.

## 4. Application to error potentials detection

### 4.1. ErrP Data and preprocessing

We now consider the problem of modelling and detecting error potentials using the approach developed above on the dataset reported in [49]. We start with a brief description of main aspects of the experiment that are relevant for the present study. Details on recordings and EEG data preprocessing are not exposed here since they had been largely detailed in [49] and are beyond the scope of the present paper.



**Figure 1.** Grand averages (FCz) of *error* (red) and *correct* (blue) trials for monopolar data.

*4.1.1. Experiment and signal acquisition.* The dataset includes 10 participants performing an Eriksen’s flanker task [16]. Each trial consists in the identification of a central letter (target) embedded in a set of three letters. Participants were asked to respond to the target letter by pressing a button with either the left hand or right hand. The signal was recorded with 64 scalp electrodes and after preprocessing (sampling to 1024 Hz and artifacts removal), data were downsampled to 256 Hz.

The selected trials were segmented into epochs of 800 ms from  $-400$  ms to  $+400$  ms, where zero corresponds to the electromyogram (EMG) onset that triggered the response. Each dataset in the training step corresponds to a *three-dimensional array* as we consider  $N^c$  trials in both classes  $c = 0$  (correct trials) and  $c = 1$  (errors) over 64 electrodes and 204 time-points.

*4.1.2. Error/Correct trials classification: an unbalanced dataset.* For each participant, we consider two categories of trials: *error* and *correct*. In this current application we do not take into account the third class discussed in [49] which corresponds to partial errors. In figure 1, differences on monopolar data between *error* and *correct* are observable on grand averages (over trials and subjects). For errors, a clear negativity is observable around 150 ms after the EMG onset (zero of time) whereas on correct trials a large positive wave occurs just after EMG onset, but no negativity is visible ‡.

The number of trials per class differs greatly for each participant. In table 1, the total number of trials for error and correct responses is given. The percentage of error trials approximately ranges from 2% to 12%, the dataset is therefore highly unbalanced, errors forming the minority class. The main goal here is to achieve *single-trial classification* on this unbalanced dataset.

For an illustrative purpose throughout the section, results are presented for participant A. Results concerning the other participants are given in Supplementary Data.

**Table 1.** Total number of trials per participant for error and correct responses.

	Participant									
	A	B	C	D	E	F	G	H	I	J
Error	130	105	43	39	18	28	94	145	100	63
Correct	1376	1575	1467	1735	690	759	1844	1238	1167	907

*4.1.3. Preprocessing.* The following two-step preprocessing for spatial and temporal decorrelation and dimension reduction is applied to the ErrP data.

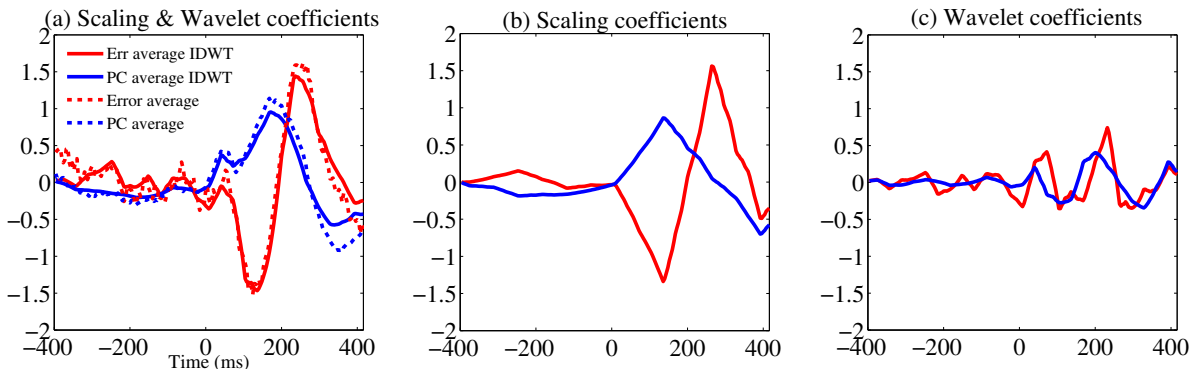
*Decorrelation and time-domain reduction.* DWT was performed using a Daubechies filter D6 (see [13, 36, 55] for details) with 5 decomposition levels. Given the sampling frequency, those 5 levels correspond to the following frequency bands:  $64 - 128Hz$ ,  $32 - 64Hz$ ,  $16 - 32Hz$ ,  $8 - 16Hz$  and  $4 - 8Hz$  respectively, corresponding to wavelet

‡ On correct trials, Current Source Density (CSD) analysis has revealed that a negative activity similar to the one observed on errors exists, but with a much smaller amplitude [49]. In the present context, since discrimination was the main goal, we did not resort to CSD analysis for maximizing the difference between correct and errors. This is why such a small negative wave is not observable on correct trials in this case

coefficients, the lowest frequency band ( $0 - 4Hz$ ) being encoded in scaling coefficients (see section 3.1.1).

Prior to DWT, signals were first extended using zero-padding. Doing so, we obtain a set of signals with appropriate length (for the DWT implementation which we use [9], it is convenient that the signal length be a power of two). Zero-padding is also an extension method which permits to avoid misinterpretation about the behaviour of signals beyond boundaries.

**Remark 3** The choice of D6 is the result of a tradeoff between localization and smoothness for the resulting wavelet. Wavelets with larger support generate more important boundary effects (we recall that in the dataset under consideration, signals are epoched, so that boundary effects have to be accounted for). Shorter filters yield wavelets with lower smoothness, so that the adjusted signals can also lack smoothness. For example, adjusted signals obtained using the Haar wavelet are discontinuous, which we didn't find satisfactory, while adjusted signals obtained with D4 are non differentiable. Further tests were made using Daubechies filters with various lengths. Corresponding results (which can be found in Supplementary data) are fully consistent with the results presented here.



**Figure 2.** Illustration of the relevance of DWT and multiscale coefficients selection. Figures (a), (b) and (c) represent error and correct responses averages over FCz. (a) Error (red) and correct responses (blue) averages obtained from the 24 selected coefficients projected back to time domain and compared with time-courses averages (dashed lines). Partial reconstructions based on scaling (b) and wavelet coefficients (c) are also represented separately.

Second, dimension reduction was achieved by removing the first three decomposition levels, that contain low amplitude high frequency phenomena in the signals. The truncation of the EEG wavelet coefficients sequence to 2 decomposition levels amounts to a projection onto an appropriate subspace [1] that keeps the relevant information in the remaining wavelet and scaling coefficients.

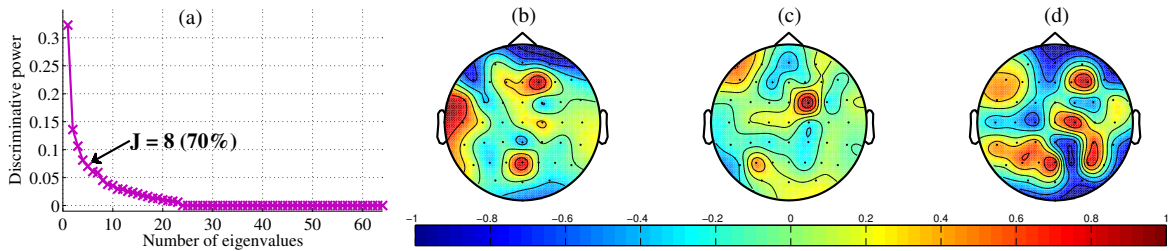
Finally, keeping in mind that the signals were zero-padded, boundary wavelet coefficients that originate from the zero-padding are excluded from the statistical analysis, as they turn out to exhibit a behaviour different from the relevant ones.



More precisely, coefficients likely to be sensitive to boundary effects were not taken into account, namely the 2 boundary scaling coefficients, the 2 boundary wavelet coefficients at coarsest scale, and the 4 boundary wavelet coefficients at finer scale. After this step,  $K = 24$  coefficients are selected: 12 + 6 wavelet coefficients for decomposition levels 4 and 5 respectively and 6 scaling coefficients.

Figure 2 illustrates the relevance of multiscale coefficients selection. The 24 coefficients remaining after DWT were projected back to the temporal space using an inverse DWT (IDWT). As seen in figure 2(a), errors and correct responses averages obtained from the selected coefficients are very similar to the 204 time-points averages. In addition, one can see that scaling coefficients (b) capture the general waveform while wavelet coefficients concentrate the details (c).

In the following, the selected multiscale coefficients are referred to as *multiscale features*.

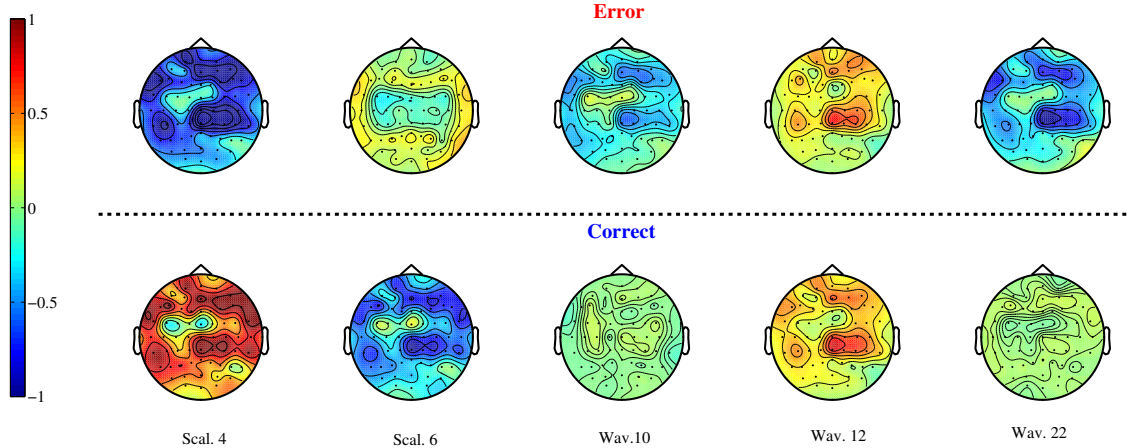


**Figure 3.** Spatial filtering results for participant A. (a) Using a cut-off of 70% of discriminative power,  $J = 8$  filters are selected based on the screeplot of eigenvalues. (b), (c), (d) display the three first selected filters, ordered by decreasing order of discriminative power.

*Spatial filtering.* Using the matrix-based Fisher’s linear discriminant proposed in section 3.1.2, spatial filtering consists in identifying the most discriminant filters to separate errors and correct responses. The number  $J$  of selected filters depends on the cumulative percentage of discriminative power, fixed to 70% in the present study for all subjects. This choice is somewhat arbitrary. We have chosen to stick to a very simple criterion, but other criteria could have been used as well. For example, Kaiser’s rule has been tested on our dataset, and results were not satisfactory in the sense that the number of selected channels was larger (and the computational load was increased), without improvement of the classification rates.

Figure 3(a) represents the screeplot, which plots the ordered eigenvalues, associated with the 64 electrodes and out of which only 8 filters are needed to concentrate 70% of the discriminative power. Figure 3 (b), (c), (d) correspond to topographical maps of the first three spatial filters i.e. the weights of the electrodes in the corresponding projections. The contributions of multiscale features in each class can be visualized using their spatial signature defined as their backprojection onto the sensor space: for each multiscale index, the pseudo-inverse of the spatial filter matrix is applied to the

vector of the corresponding coefficient values at all channels. For the sake of illustration, we display in figure 4 the spatial signatures of some multiscale features averaged across trials.



**Figure 4.** Topographical maps of spatial signatures of five multiscale features averaged across trials.

Here we have chosen to display five features corresponding to scaling coefficients 4 and 6 and to wavelet coefficient 10, 12 and 22. These specific multiscale coefficients were chosen for their relevance for *error*, *correct* or both classes: on average they appear as emergent, large magnitude coefficients (see figure 6 below). For example, scaling coefficient 4 has the largest average magnitude in both classes and was found to be highly discriminant. This behaviour is reflected by the corresponding spatial signatures (defined above) which are strongly negatively correlated. Scaling coefficient 6 is relevant only for the correct class and so is its spatial signature. Wavelet coefficients 10 and 22 are only relevant for error whereas the coefficient 12 is not discriminant. In agreement with a large literature on ErrP, the spatial signatures for the relevant multiscale features mainly loads on fronto-central sensors.

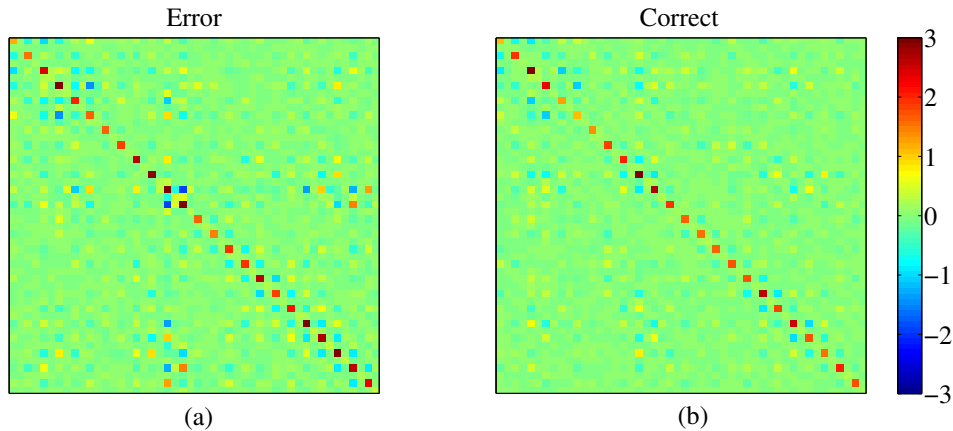
#### 4.2. Modelling ErrP single-trials

After applying spatial and temporal filters, the extracted features from ErrP data are decorrelated and dimensions are significantly reduced. We now consider multichannel trials in the wavelet domain, such that  $Y_i^c \in \mathbb{R}^{K \times J}$ , with  $K = 24$  and  $J \ll M$  (for the training set example  $J = 8$ , see figure 3). Modelling ErrP single-trials is based on the columnwise vectorized trial  $y_i^c$ .

*4.2.1. Testing class-covariance matrices inequality.* We display in figure 5 the sample covariance matrix for each class (*error* and *correct*), calculated after preprocessing. For simplicity the figure is limited to the first two channels. In addition for the relevance of the visual comparison, the two matrices have been computed on datasets of equal

size (actually maximal possible size, to ensure maximum precision). Visual inspection reveals different diagonals and noticeable non-diagonal terms in the sample covariance matrix of the *error* class. The significance of the difference between these two matrices is evaluated quantitatively using Box’s M test [2]. The latter clearly rejects the class-covariance matrices equality null hypothesis (p-value  $\ll 10^{-6}$ ).

Importantly, this conclusion differs from the claim of [7], who did not see any noticeable difference between sample class-covariance matrices (on a different EEG dataset). Although such a finding may be dataset dependent, we stress that for the dataset considered here, the significant difference results from time and space dimension reductions, which yield an important denoising. In the absence of dimension reductions (for example raw data), the sample covariance matrices are generally singular and the Box’s M test cannot be performed.

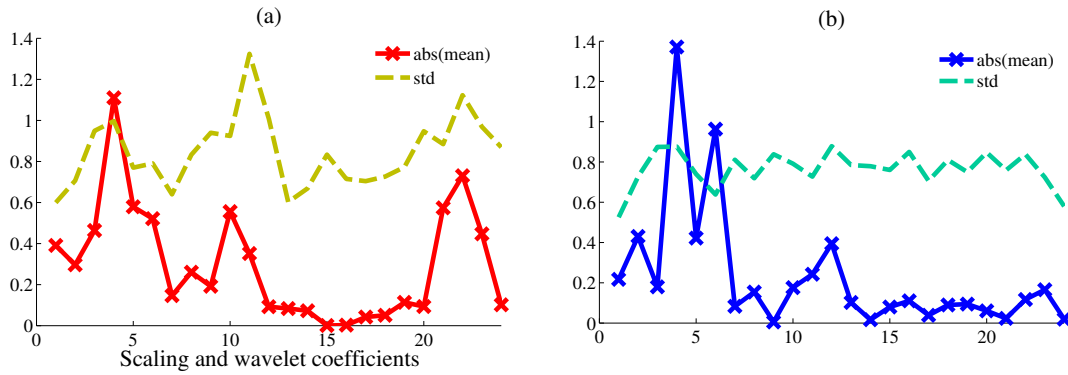


**Figure 5.** Representation of  $48 \times 48$  sample class-covariance matrices for 24 multiscale coefficients and 2 channels : *error* class (a) and *correct* class (b).

*4.2.2. Choice of the random part coefficient vector  $\Gamma^c$ .* As mentioned in section 3.2.3, the choice of  $\Gamma^c$  is problem dependent. In our application no prior information was available that could help specifying the form of  $\Gamma^c$ , we thus relied on preliminary data exploration. We plot in figure 6 the mean in absolute value and the standard deviation computed over all trials for participant A on the first channel in the two classes. This reveals a monotonic relationship between the amplitude and variability of the signal of interest, with an additional offset effect. The latter can be interpreted as originating from random noise  $\varepsilon$  in (10). Similar results are obtained for the other participants and channels.

In the present work, we take the simplest relationship, namely a linear one, by setting the vector  $\Gamma^c$  to the average over trials in class  $c$ :

$$\Gamma^c = \frac{1}{N^c} \sum_{i=1}^{N^c} y_i^c. \quad (18)$$



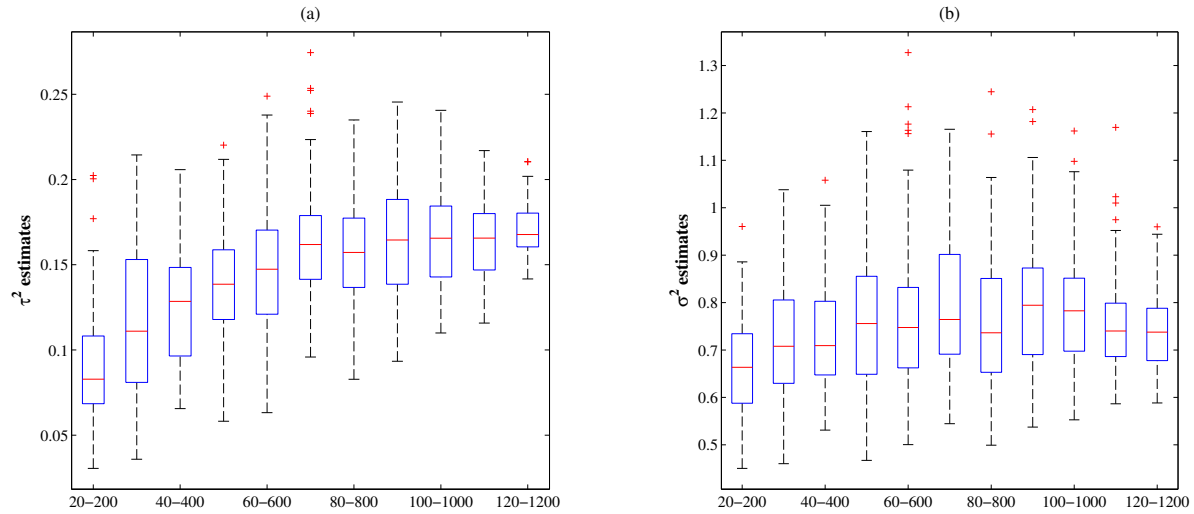
**Figure 6.** Illustration of the monotonic relationship between amplitude and variability in ErrP-related EEG signals: mean in absolute value ( $abs(mean)$ ) and to the standard deviation ( $std$ ) both calculated across all trials of the first channel. (a): *error* ; (b): *correct*.

*4.2.3. Variance components estimation.* We now explore one of the main assumptions of our study, namely the single-trial signal of interest can be written as the sum of two components: a class-dependent fixed term (the class-mean) plus a trial-dependent random term (the so-called random effect). The hypothesis to be tested is the significance of the random effect variance  $\tau^2$ . A negative answer would indicate that a random effect is not needed. In that case, a classical linear model as proposed in [7] would be more appropriate.

The implementation of the mixed model has been done using the MATLAB function *mixed.m* written by Witkovský [60] (more details can be found in Supplementary Data). REML estimates for variance components  $\tau^2$  and  $\sigma^2$  are displayed in figure 7. Each boxplot corresponds to a different sample size (10% error trials, 90% correct trials) and summarizes the distribution of estimates performed on 100 different training sets. We note that the magnitude of the residual variance estimates  $\hat{\sigma}^2$  is larger than that of  $\hat{\tau}^2$ . Globally the dispersion of the estimates decreases as a function of the sample size, nevertheless the figure shows that good estimates can be obtained even with small sample sizes. In figure 7 (a) the quartiles of the  $\hat{\tau}^2$  distribution tend to increase slightly with the sample size and minimum estimates rise above zero. Using a z-test we test that the random effect variance is strictly positive in average for each training sample size (p-value  $< 10^{-3}$ ).

### 4.3. Classification results

*4.3.1. Performance evaluation.* In the case of unbalanced classes, the performance evaluation methodology has to be chosen carefully as the majority class may shadow the behaviour of the minority class. This problem has been discussed in several works (see e.g. [58, 57, 40] and references therein). In this paper, we use various performance measures for comparing unbalanced EEG dataset classifiers. In particular, we focus on class-dependent evaluations.



**Figure 7.** Covariance components estimation. (a) Random effect variance  $\tau^2$  estimates. (b) Residual variance  $\sigma^2$  estimates. Each boxplot displays the series of 100 parameter estimates for different training sample sizes (10% error trials, 90% correct trials). *The central mark of boxplot is the median (second quartile), the bottom and top of the box are the first and third quartiles, the whiskers extend to the most extreme values not considered outliers, and outliers are plotted individually.*

- *Confusion matrix.* In a binary classification, given a classifier and a test sample, four outcomes are possible: true positive, false positive, true negative and false negative. For instance, if the test sample is positive and it is classified as negative, it is counted as a false negative. For a test dataset, the number of occurrences of the four outcomes are respectively denoted by (TP), (FP), (TN) and (FN), and form the  $2 \times 2$  confusion matrix given in table 2.

**Table 2.** Confusion matrix for the ErrP classification problem.

		predicted	
		<i>corret</i>	<i>error</i>
actual	<i>correct</i>	TN	FP
	<i>error</i>	FN	TP

- *Good classification rate.* In the particular case of unbalanced data, the good classification rate must be calculated for each class independently. Indeed, a global good classification rate may lead to misleading results. For example, in our problem, where the errors only represent 10% of the data, the naive classification strategy allocating all testing data to the majority class would automatically achieve a good classification rate of 90%. However the classifier is obviously irrelevant for detecting *error* trials.
- *Pierce's Skill Score.* Given the confusion matrix, the good prevision rate  $H = TP/(FP + TP)$  and the false alarm rate  $F = FP/(TN + FP)$  for the *error* class

lead to the so-called Pierce score

$$PSS = H - F. \quad (19)$$

By construction,  $PSS \in [-1; +1]$ . In our context the closer  $PSS$  is to +1, the better the classifier is for the detection of *error* trials.

*4.3.2. ErrP detection.* For the considered dataset the number of *error* trials ( $N^1$ ) is much smaller than the number of *correct* trials ( $N^0$ ). For performance evaluation, unbalanced training sets of increasing size are generated, all based on the *a priori* 10% of *errors*. For each participant and each randomly drawn unbalanced training set (namely,  $\{N^1 = 20, N^0 = 200\}, \{N^1 = 30, N^0 = 300\}, \dots$ ), the test set is composed of all remaining data, therefore its size differs. For each participant this splitting is performed 100 times and classifiers performance are recorded.

The proposed method (named LMM hereafter) is compared with several classifiers: classical linear discriminant analysis (LDA), regularized method (regularized LDA - RDA - as proposed in [7]) and diagonal discriminant analysis (diagonal LDA (DLDA) and diagonal QDA (DQDA)). In addition, LMM is also compared to LDA, RDA and DLDA classifiers after sub-sampling the majority class to balance training sets. The classic QDA classifier cannot be used, because the sample covariance matrix in *error* class is very often singular in the considered range of sample sizes. Corresponding results are displayed in figure 8. Let us stress that the error bars represent the standard deviation of the good classification rates distribution. Since the training and test sets are complementary, when the size of training set increases, the number of single-trial signals used to evaluate the good classification rates decreases, and consequently the error bar increases. These error bars only measure the dispersion of the estimates and cannot be used for quantitative pairwise comparison (see below).

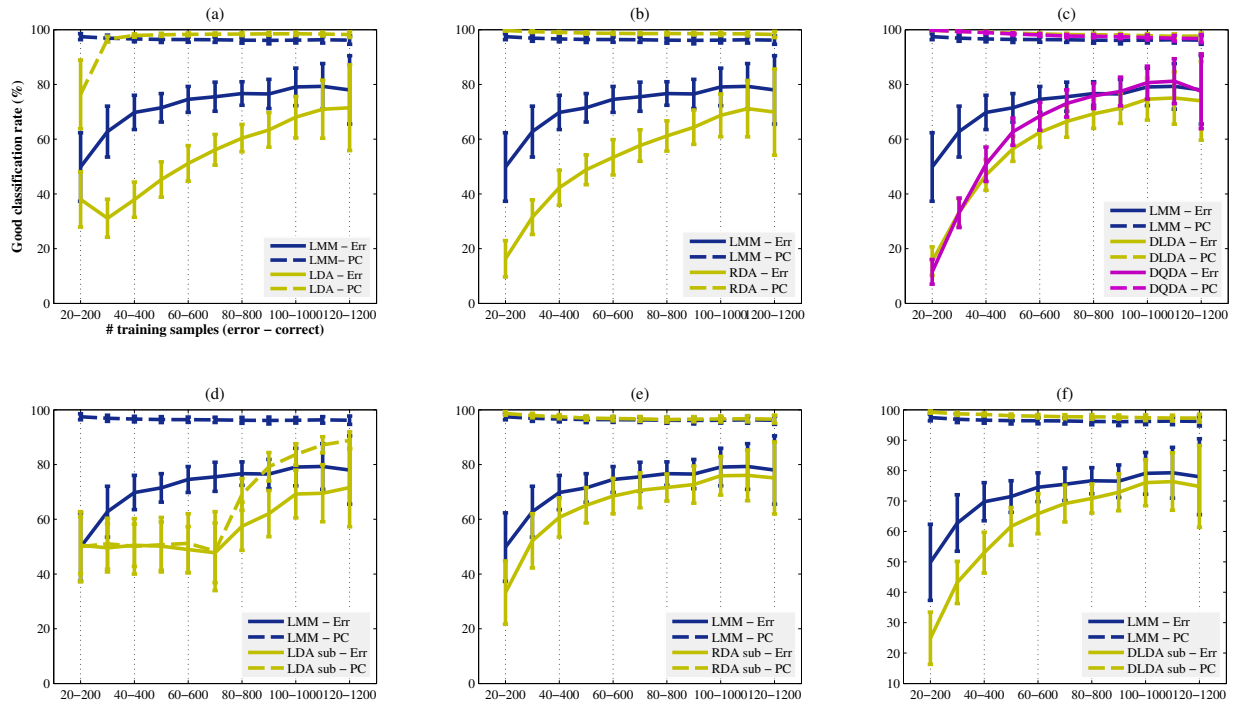
Clearly LMM greatly outperforms each method in terms of *error* detection in all situations. This result is particularly spectacular for small training set which we interpret as follows. The simplicity of the underlying model (which involves very few parameters) and its relevance for the considered dataset yield accurate and robust parameter estimates even for small sample size. In addition, our method performs equivalently to other methods to detect *correct* trials which correspond to the majority class. Even with small sample size, classifiers detect *correct* trials with a good classification rate close to 100%.

We provide in table 3 quantitative classifiers comparison based on Pierce's Skill Score (PSS) defined in (19). For each training set, PSS is computed for each classifier. LMM is compared with each other classifier (OC) using Wilcoxon signed rank test (a paired non-parametric statistical test, see [23, 26]) based on the average difference  $PSS_{LMM} - PSS_{OC}$ .

For participant A, LMM outperforms all the other considered classifiers very significantly except DQDA from the sample size (80 – 800). Indeed the DQDA performances are equivalent to or significantly better than the LMM ones for the largest

sample sizes. Quantitative results for all participants can be found in Supplementary Data. The conclusions drawn for participant *A* are confirmed in most configurations but should be modulated as follows. In most situations LMM performs consistently better than the other classifiers. This in particular true for small training sets namely when the number of errors is in the range (20, 30, 40, 50), where the superiority is supported by statistical tests (with only two exceptions where results are equivalent). These results can be explained by the preprocessing which enforces the data decorrelation and are coherent with the results of Box’s M test (see section 4.2.1).

The proposed method is therefore a valuable alternative to linear discriminant techniques which involves quadratic decision rule avoiding the difficulty of usual shortcomings of classical QDA.



**Figure 8.** Good classification rates for different training sample sizes (*error–correct*). Results correspond to good classification rates averaged over 100 iterations for each training sample size. Vertical bars represent standard deviations (std). Results are given for *error* (plain lines) and *correct* (dotted lines) test trials. LMM classifier is compared with three different classifiers on the unbalanced and re-balanced datasets. (a), (b), (c) display results in the unbalanced case (10% *error* trials) for classic LDA (LDA), Regularized LDA (RDA) and Diagonal LDA and QDA (DLDA & DQDA) respectively. (d), (e), (f) display results of the LDA, RLDA and DLDA classifiers after re-balancing datasets by subsampling the *correct* class.

**Remark 4** As mentioned above, the correct classification rate depends heavily on the quality of the estimation of the covariance matrices. In our case, the latter are strongly simplified, and the issue is the quality of the estimation of the mean vectors  $\mu^c$  and the variances  $\tau^2$  and  $\sigma^2$ . Thanks to pre-processing and dimension reduction, variances

**Table 3.** Participant A: Classifiers comparison based on PSS. LMM performances are compared with other classifiers using the Wilcoxon signed rank test. For each pair of classifiers LMM-other classifier, the difference between the two mean PSS is given and p-value order of magnitude is provided between parentheses.

$N^1-N^0$	Classifiers comparison			
	LMM-LDA	LMM-RDA	LMM-DLDA	LMM-DQDA
20-200	0.14 ( $10^{-12}$ )	0.36 ( $10^{-18}$ )	0.37 ( $10^{-18}$ )	0.41 ( $10^{-18}$ )
30-300	0.34 ( $10^{-18}$ )	0.34 ( $10^{-18}$ )	0.32 ( $10^{-18}$ )	0.32 ( $10^{-18}$ )
40-400	0.35 ( $10^{-18}$ )	0.30 ( $10^{-18}$ )	0.25 ( $10^{-18}$ )	0.21 ( $10^{-18}$ )
50-500	0.28 ( $10^{-18}$ )	0.25 ( $10^{-18}$ )	0.16 ( $10^{-18}$ )	0.09 ( $10^{-16}$ )
60-600	0.25 ( $10^{-18}$ )	0.23 ( $10^{-18}$ )	0.13 ( $10^{-18}$ )	0.07 ( $10^{-13}$ )
70-700	0.21 ( $10^{-18}$ )	0.19 ( $10^{-18}$ )	0.10 ( $10^{-18}$ )	0.03 ( $10^{-4}$ )
80-800	0.18 ( $10^{-18}$ )	0.17 ( $10^{-18}$ )	0.08 ( $10^{-17}$ )	0.009 (0.43)
90-900	0.14 ( $10^{-18}$ )	0.13 ( $10^{-18}$ )	0.06 ( $10^{-13}$ )	-0.01 ( $10^{-3}$ )
100-1000	0.12 ( $10^{-17}$ )	0.11 ( $10^{-17}$ )	0.05 ( $10^{-10}$ )	-0.02 ( $10^{-4}$ )
110-1100	0.09 ( $10^{-12}$ )	0.09 ( $10^{-13}$ )	0.04 ( $10^{-7}$ )	-0.02 ( $10^{-4}$ )
120-1200	0.07 ( $10^{-4}$ )	0.08 ( $10^{-6}$ )	0.04 ( $10^{-2}$ )	0.005 (0.23)

are estimated from  $KJ(N^0 + N^1)$  samples, while mean vectors  $\mu^c$  are estimated from  $N^c$  samples. Hence the quality of the classification is mainly driven by the quality of the estimation of the mean vectors. Since we consider unbalanced situations, the difficult class is the minority class (in our case the *error* class) for which approximately 30 samples or more are needed to get more than 60% correct classifications. The other class being ten times bigger, the corresponding classification rate is close to 100%, and the global classification rate is satisfactory.

Similar results using different Daubechies filters are presented in Supplementary Data. Whatever the filter, the proposed method consistently outperforms concurrent approaches in terms of classification, best results being obtained when the Haar filter is used.

#### 4.4. Back to time courses: adjusted single-trials

As described in subsection 3.5, for each single-trial an estimate of the signal of interest can be obtained according to (17). After back projection onto sensor space (using a Moore-Penrose pseudo-inverse) and inverse DWT we obtain time courses called *adjusted single-trials*. Examples of adjusted single-trials (electrode FCz) are displayed in figure 9. These plots are obtained using the D6 Daubechies filter. Adjusted single-trials obtained using other Daubechies filters are displayed in Supplementary Data.

The left hand plot (a) represents the raw signals averaged over trials (within each class) together with the averages of adjusted single-trials. As can be seen for each class the two averages are very close to each other for participant A. For other participants, this similarity is less striking but still present. This shows that neither dimension reductions nor modelling have strongly affected the average information. Although

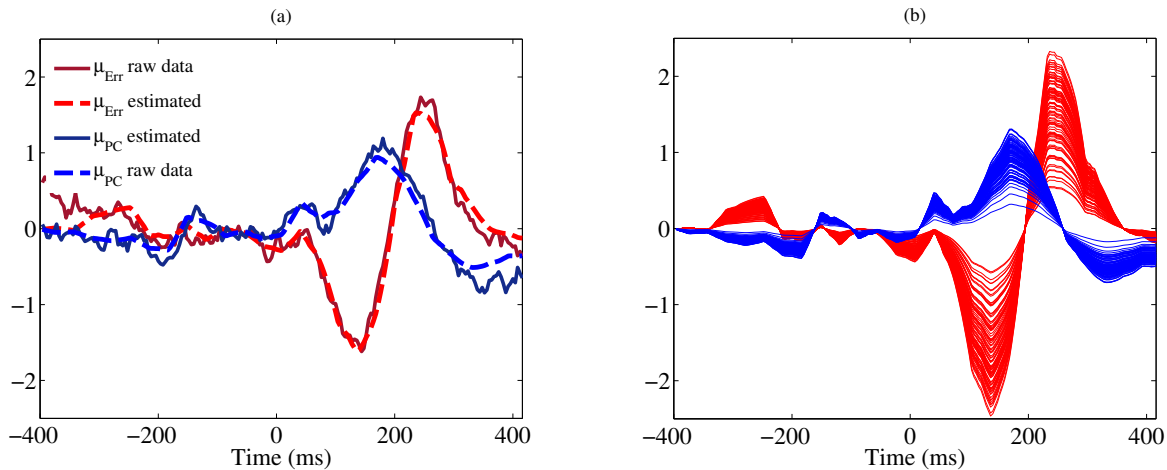


not surprising and even expected, the recovery of the global shape of the average is a first necessary step before looking at single-trial reconstruction. It appears that the modelling indeed captures the signal of interest and filters out the background activity.

As largely reported in the literature (see [17] for an overview), in the *error* class, the estimated averages as well as the adjusted single-trials show a clear negative wave around 150 ms after EMG onset followed by a positive one. Such a negative wave is not observable in the *correct* class, as commonly reported with monopolar recordings.

The added value of the model is the explicit description and the numerical quantification of the inter-trial variability. The latter appears clearly in the right hand plot of figure 9(b) where adjusted single-trials are displayed. The proposed method therefore allows one to reveal the amplitude variation of the ErrP: while some errors induce a very large activity, some others present virtually no response. The same figure is displayed for each participant in Supplementary Data.

Although speculative at this point, this could reveal differences in sensitivity to the errors across trials, an hypothesis that now becomes testable thanks to the proposed method.



**Figure 9.** Illustration of single-trial reconstruction. (a) *Error* and *correct* trial averages on FCz for participant A. Dashed curves: averages across raw trials in both classes; full curves: averages across adjusted trials. (b) Adjusted single-trials for *Err* (red) and *PC*(blue) classes on FCz.

## 5. Discussion

We first highlight the key points and main contributions of our method and then discuss other possible applications of this work.

### 5.1. The impact of time-domain decorrelation

Within the proposed procedure, we would like to emphasize that the introduction of DWT constitutes a key step for the analysis and the classification of single-trials.

Considering multiscale coefficients instead of time-point values first gives a simplified representation of data, but it also has important consequences for both modelling and classification. Indeed, decorrelation and dimension reduction lead to small size, diagonal dominant covariance matrices, easier to estimate. Consequently, in the wavelet domain single-trials can be modelled by simple linear mixed model with a small number of unknown parameters, allowing robust estimates from small unbalanced training set and good classification results.

### *5.2. Fine structure of the model*

Mixed model provides a general and flexible approach to analyze single-trial EEG signals. It allows one to include both fixed and random effects within the same analysis. In the case of EEG datasets it appears natural to consider trial effects as random effects since they can be seen as originating from a random selection from a much larger set. Furthermore, it is of direct relevance to separate inter and intra-trial variability. However for each problem and dataset of interest, the main difficulty lies in the choice of the dimension of the random effect vector, its class-dependence and the covariance structure.

For ErrP detection problem, the training sample size being small, we choose the simplest possible linear mixed model: class as fixed effect (class-average), one dimensional random effect vector (class-independent) modulated by a class-dependent coefficient vector  $\Gamma^c$  plus a white noise. The choice of  $\Gamma^c$  depends on the dataset and can be done using prior information or exploratory data analysis. In our application the observed relationship between amplitude and variability led us to set  $\Gamma^c$  equal to the average over trials in each class.

### *5.3. Relevant classification results*

In section 4.3, we provided a systematic comparison of related classification techniques, all based upon linear models, involving the same fixed part (class-average) and various approaches to model variability in the unbalanced situation. All methods were applied after the same preprocessing step, that enforces diagonal dominance for covariance matrices. Consequently, methods that can exploit such diagonal dominance properties are particularly relevant. As expected, results show that they significantly outperform classic LDA or QDA in terms of classification. Moreover our classification results highlight two important points. The first one is that our procedure is particularly efficient when the size of the training sample is small, emphasizing the relevance of the proposed variability modelling through the linear mixed model. Secondly, beyond a certain sample size, namely when the number of observations permits to estimate precisely the diagonal of the covariance matrix in the minority class, diagonal QDA outperforms all considered methods. That point is coherent with the results of the Box's M test of equality of class-covariance matrices. For all participants, when the number of trials is sufficient to calculate the test statistic, Box's M test rejects the

hypothesis of equality (p-value  $\leq 10^{-3}$ ). Unlike the balanced case where LDA generally outperforms QDA even when class-covariance matrices are different, diagonal QDA gives better results than diagonal LDA. In the unbalanced case, we stress that it is preferable to take into account the difference in variability of both classes. However it is well known that the classifier performance depends on the accuracy of the estimate of parameters involved in the procedure, and to obtain accurate parameter estimates the number of samples must be much larger than the number of parameters. Therefore when the sample size is small, modelling the intra-class variability is well suitable. In our model, in addition to the class average estimates, only two parameters are needed to estimate the two class-covariance matrices.

#### 5.4. Possible applications

Being able to exploit the single-trial M/EEG component would open a new window on the dynamics of brain processes, and the range of applications is wide. Although some might be obvious, let us focus on some areas where the proposed approach might be especially useful.

*5.4.1. BCI.* In the present application, we considered a discriminant framework. In this context, BCI naturally becomes a straightforward extension for the proposed LMMC procedure. Indeed, classification results, presented in section 4.3.2, illustrate the performance of single-trial variability modelling on classifying *error vs correct* trials. Importantly, the total variance explicit modelling, combined with spatial and temporal dimension reductions, allows robust parameters estimates, even when based on a limited training dataset. Sticking to a small training dataset is essential to keep the BCI calibration session as short as possible, a condition that can be reached thanks to the simplicity of the proposed model. The procedure, as described above, is not restricted to error potential estimate, however, and can easily be extended to other BCI protocols. Indeed, a strength of LMMC method is the modularity and interchangeability of the preprocessing steps. Concerning the spatial dimension reduction, well established spatial filtering methods have been developed for specific BCI tasks: for example Common Spatial Pattern (CSP) are especially suited spatial filters for motor imagery [50] whereas xDAWN has been developed for P300 data [48]. Those methods can easily be implemented in the LMMC procedure while leaving the modelling step almost unchanged. The same holds for temporal dimension reduction. In the current application, to get rid of temporal correlation and for temporal dimension reduction, we used wavelet transform which is well suited to extract transient phenomena such as ERP's (ErrP, P300, etc.). However, the choice of the basis can depend on the data to analyze, and for more oscillatory signals such as the ones recorded in motor imagery (like  $\mu$  and  $\beta$  rhythms, see [39]), other transforms such as time-frequency transforms (MDCT, STFT, see [35]) or adaptive transforms [54, 4] might be more suitable. Again, changing the transform leaves the variance modelling steps almost

unchanged, as long as the transform is invertible.

*5.4.2. Single-trial extraction in cognitive neurosciences.* In many experimental situations, being able to extract the single-trial amplitude of brain activity would open a new window on the dynamics of brain processes. One example is the temporal evolution of brain activity during learning: across repetitions of the learning situation, some brain regions might become more active, while some other might disengage from the task. Relying on averaging, at best, dramatically reduces the analysable temporal dynamic of the learning process. For example, while no clear learning effect could be evidenced on the raw Auditory Evoked Potentials (AEP) in rats, extracting the single-trial component (through denoising in this case) allowed to regress learning curves on the trial-by-trial responses, and show that the amplitude of some components of the AEP reduces with learning while other remain unchanged [46]. Averaging also prevents analyzing correlations between brain activities and behaviour across trials. For example, RT or psychophysical judgements are known to be highly variable even for the very same stimulation. The variability in RT have long be known to provide essential information to put into test different model of information processing [34]. Unfortunately, the equivalent variability is normally not accessible for EEG/MEG brain responses, preventing the use of similar constraints on models based on brain activity (see e.g. [22], for trial-to-trial variability to constraint RT models). The LMM approach furnishes a way to recover this variability (at least in amplitude) and may help to correlate brain activities to behaviour. Further applications will probe what type of constraint such variability can bring for model testing.

Another key interest of the proposed approach is the ability to better analyze unbalanced data, which is a quite common situation. Indeed, from a statistical point of view, LMM allows increasing the robustness of the estimates from small datasets and hence reducing the risk of type II error (false negative), which is a common problem in small and noisy dataset. It may also help to decrease the number of necessary trials to reach a given statistical power, which can be of tremendous interest on some populations (children, patients, etc.).

One may wonder, however, how general this methodology can be and whether it would apply with similar success on other brain activities which might not be as stable as the ErrP. Although we do not have any empirical response, formal analysis of the model allows speculating a bit. In the current modelling, the class average plays a central role. One can therefore anticipate that the quality of the modelling will largely be driven by the quality of the average. As a consequence, as long as the average emerges from the background noise, that is, as long as there is an ERP, it should be possible to apply the model. For example, it should be possible to apply the method on robust components, such as the standard visual evoked potentials (N1, P2, N2 etc.). This will be investigated on future research.

## 6. Conclusions and future works

In the present work we propose a single-trial classification procedure (LMMC) based on an explicit model of the signal variability in the wavelet domain. Each single-trial signal is assumed to be the sum of two components: a background activity and a signal of interest. The latter can again be decomposed into the sum of a fixed part (the class-mean) and a random part that models the deviation of each single-trial from the class-mean. This hypothesis is formalized mathematically as a Gaussian linear mixed model (LMM). This model generalizes the Gaussian linear model given in [7] which was adapted to situations where between-trial variability is negligible. An important step is the discrete wavelet transform in the preprocessing. On the one hand, this transformation allows one to reduce the time dimension and on the other hand, to consider a very simple model where class-covariance matrices are characterized by only few parameters (to be estimated). The advantage of such a model is to yield robust estimates of parameters from a small dataset.

This model is applied to a classification problem, namely separating correct from erroneous responses in a RT task. As expected, the proposed classification procedure is particularly efficient to detect EEG Error Potentials in the unbalanced case when the sample size of the minority class is small. Moreover these results validate the chosen modelling.

In addition, the proposed procedure allows the extraction of relevant and discriminant information from EEG signals. As a matter of fact model parameters estimates provide quantitative measures of the between-trial variability in each class in addition to the average behaviour. In the specific context of error potential detection, the proposed model allows one to recover single-trial signal of interest and more interestingly reveals the trial-by-trial amplitude variation of ErrP's.

Several directions might be interesting to investigate further and we list a few of these below.

First, the choice that was made for random part coefficient vector  $\Gamma^c$  explicitly assumes a linear relationship between the signal of interest's amplitude and variability. However, although the existence of a monotonic dependence between mean and variance seems a reasonable assumption, the relationship needs not be linear and is likely to be more complex (for example, class dependent). Further investigating and characterizing the link between mean and variance should improve the choice of the random component and hence the modelling performance.

Another improvement would be to better take into account the cubic structure of M/EEG data (*i.e* space  $\times$  time  $\times$  trials). Decoupling spatial and temporal features extraction [53] could lead to a more complex form for the random effects design matrix and hence to a spatio-temporal modelling of between-trial variability.

Comparing the EEG error potentials across participants, not only on this dataset but in the literature, see *e.g.* [49, figure 3], suggests large similarities between participants, but also some differences, both in terms of topography and time courses.

One might therefore consider incorporating also the between participant variability within a more sophisticated mixed model, which would then model the deviation of the participant (in time and in topography) to the mean behaviour. The model would include common settings and parameters specific to each subject. From a statistical point of view, the introduction of such common settings would presumably increase the robustness of the estimates from small datasets. This would in particular allow a more robust estimate of the single participant topography and time course. The question as to whether such differences in topography is related to fine anatomical differences [44] would become easier to test, especially when few trials are available.

### Acknowledgments

The authors are grateful to the anonymous reviewers, whose valuable questions and comments have contributed to significantly improve the quality of the manuscript. This work was supported by the ANR project CO-ADAPT (ANR-09-EMER-002-05). J. Spinnato's work is funded by a PhD grant from Région Provence-Alpes-Côte d'Azur, France. B. Burle is supported by a European Research Council under the European Community's Seventh Framework Program (FP7/2007-2013 Grant Agreement no. 241077). This work has been carried out in the framework of the Labex Archimède (ANR-11-LABX-0033) and of the A\*MIDEX project (ANR-11-IDEX-0001-02), funded by the "Investissements d'Avenir" French Government programme managed by the French National Research Agency (ANR).

### References

- [1] F. Abramovich, T. C. Bailye, and T. Sapatinas. Wavelet analysis and its statistical applications. Royal Statistical Society : Series D (The statistician), 49(1):1–29, 2000.
- [2] T. W. Anderson. An Introduction to Multivariate Statistical Analysis. John Wiley and sons, 3rd edition, 2003.
- [3] R. H. Baayen, D. J. Davidson, and D. M. Bates. Mixed-effect modeling with crossed random effects for subjects and items. Journal of Memory and Language, 59:390–412, March 2008.
- [4] S. Barbieri and B. Torrèsani. Optimal time-frequency bases for EEG signal classification in the context of bci. In Proceedings of GRETSI 2013 (Brest), september 2013.
- [5] C. G. Bénar, T. Papadopoulo, B. Torrèsani, and M. Clerc. Consensus matching pursuit for multi-trial EEG signals. Journal of Neuroscience Methods, 180:161–170, 2009.
- [6] P. J. Bickel and E. Levina. Some theory for Fisher's linear discriminant function, naive Bayes, and some alternatives when there are many more variables than observations. Bernoulli, 10:989–1010, 2004.

- [7] B. Blankertz, S. Lemm, M. Treder, S. Haufe, and K. Müller. Single-trial analysis and classification of ERP components - a tutorial. Neuroimage, 56:814–825, 2011.
- [8] B. Blankertz, S. Tomioka, R. and Lemm, M. Kawanabe, and K. Müller. Optimizing spatial filters for robust EEG single-trial analysis. IEEE Signal processing magazine, 20, 2008.
- [9] J. Buckheit, S. Chen, D. Donoho, I. Johnstone, and J. Scargle. About wavelab. Technical report, Stanford University, 2005.
- [10] B. Burke-Hubbard. The World According to Wavelets: The Story of a Mathematical Technique in the Making. A.K. Peters, 1998.
- [11] B. Burle, C. Roger, F. Vidal, and T. Hasbroucq. Spatio-temporal dynamics of information processing in the brain: Recent advances, current limitations and future challenges. International Journal of Bioelectromagnetism, 10:17–21, 2008.
- [12] E. Callaway, R. Halliday, H. Naylor, and D. Thouvenin. The latency of the average is not the average of the latencies. Psychophysiology, 21:571, 1984.
- [13] I. Daubechies. Ten lectures on wavelets. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1992.
- [14] D. J. Davidson. Functional mixed-effect models for electrophysiological responses. Neurophysiology, 41:79–87, January-February 2009.
- [15] P. Durka. Matching Pursuit and Unification in EEG Analysis. Engineering in Medicine & Biology. Artech House, 2007.
- [16] B. A. Eriksen and C. W. Eriksen. Effects of noise letters upon the identification of target letter in a non-search task. Perception and Psychophysics, 16:143–149, 1974.
- [17] M. Falkenstein, J. Hoormann, S. Christ, and J. Hohnsbein. ERP components on reaction errors and their functional significance: a tutorial. Biological Psychology. Special Issue: Error Processing and Adaptive Responding, 51:87–107, 2000.
- [18] L. Farwell and E. Donchin. Talking off the top of your head: toward a mental prosthesis utilizing event-related brain potentials. Electroencephalography and clinical Neurophysiology, 70:510–523, 1988.
- [19] S. Fazli, M. Danóczy, J. Schelldorfer, and K.-R. Müller.  $l_1$ -penalized linear mixed-effects models for high dimensional data with application to bci. NeuroImage, 56:2100–2108, April 2011.
- [20] R. A. Fisher. The use of multiple measurements in taxonomic problems. Annals of eugenics, 7:179–188, 1936.
- [21] J. H. Friedman. Regularized discriminant analysis. Journal of the American Statistical Association, 405:165–175, July 1988.
- [22] A. D. Gerson, L. C. Parra, and P. Sajda. Cortical origins of response time variability during rapid discrimination of visual objects. NeuroImage, 28:342–353, September 2005.

- [23] D. Gibbons and S. Chakraborti. Nonparametric Statistical Inference. Statistics: Textbooks and Monographs, fifth edition, 2010.
- [24] D. J. Hand and V. Vinciotti. Choosing  $k$  for two-class nearest neighbour classifiers with unbalanced classes. Pattern Recognition Letters, 24:1555–1562, 2003.
- [25] T. Hastie, R. Tibshirani, and J. Friedman. The elements of statistical learning. Springer, second edition, 2009.
- [26] M. Hollander, D. A. Wolfe, and E. Chicken. Nonparametric Statistical Methods. Wiley, third edition, 2014.
- [27] Y. Huang, D. Erdogmus, K. Hild II, M. Pavel, and S. Mathan. Mixed effects models for single-trial ERP detection in noninvasive brain computer interface design. In Recent Advances in Biomedical Signal Processing, pages 171–180. E-Book Preprint Bentham Science Publishers, October 2009.
- [28] Y. Huang, D. Erdogmus, and M. Pavel. Mixed effects models for EEG evoked response detection. In Machine Learning for Signal Processing, IEEE Workshop, 2008.
- [29] T.-P. Jung, S. Makeig, M. Westerfield, J. Townsend, E. Courchesne, and T. J. Sejnowski. Analysis and visualizations of single-trial event related potentials. Human Brain Mapping, 14:166–185, 2001.
- [30] K. H. Knuth, A. S. Shah, W. A. Truccolo, M. Ding, S. L. Bressler, and C. E. Schroeder. Differentially variable component analysis: Identifying multiple evoked components using trial-to-trial variability. Journal of Neurophysiology, 95:3257–3276, 2006.
- [31] L. H. Koopmans. The spectral analysis of time series. Probability and Mathematical Statistics. Academic Press, 1995.
- [32] M. Kukleta and M. Lamarche. Steep early negative slopes can be demonstrated in pre-movement Bereitschaftspotential. Clinical Neurophysiology, 112:1642–1649, 2001.
- [33] F. Lotte, M. Congedo, A. Lécuyer, F. Lamarche, and B. Arnaldi. A review of classification algorithms for EEG-based brain computer interfaces. Journal of Neural Engineering, 4, 2007.
- [34] R. D. Luce. Response times: Their roles in inferring mental organization. Oxford University Press, New York, 1986.
- [35] M. Mahanta, A. Aghaei, and K. Plataniotis. A Bayes Optimal Matrix-variate LDA for Extraction of Spatio-Spectral Features from EEG Signals. In Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC 2012), pages 3955–3958, 2012.
- [36] S. Mallat. A Wavelet Tour of Signal Processing, Third Edition: The Sparse Way. Academic Press, 3rd edition, 2008.
- [37] C. E. McCulloch, S. R. Searle, and J. M. Neuhaus. Generalized, Linear and Mixed Models. Wiley, second edition, 2008.



- [38] D. J. McFarland, L. M. McCane, S. V. David, and J. R. Wolpaw. Spatial filter selection for EEG-based communication. Electroencephalography and clinical Neurophysiology, 103:386–394, 1997.
- [39] D. J. McFarland, L. A. Miner, T. M. Vaughan, and J. R. Wolpaw. Mu and Beta Rythm Topographies During Motor Imagery and Actual Movements. Brain Topography, 12(3):177–186, 2000.
- [40] G. Menardi and N. Torelli. Training and assessing classification rules with unbalanced data. Technical report, DEAMS, 2010.
- [41] D. E. Meyer, A. Osman, D. E. Irwin, and S. Yantis. Modern mental chronometry. Biological Psychology, 26:3–67, 1988.
- [42] L. D. Parra, C. D. Spence, A. D. Gerson, and P. Sajda. Recipes for the linear analysis of EEG. Neuroimage, 28:326–341, 2005.
- [43] M. I. Posner. Orienting of attention. Q J Exp Psychol, 32(1):3–25, Feb 1980.
- [44] E. Procyk, C. R. E. Wilson, F. M. Stoll, M. C. M. Faraut, M. Petrides, and C. Amiez. Midcingulate motor map and feedback detection: Converging data from humans and monkeys. Cereb Cortex, Sep 2014.
- [45] Q. R. Quiroga. Obtaining single stimulus evoked potentials with wavelet denoising. Physica D, 145:278–292, 2000.
- [46] Q. R. Quiroga and E. L. J. M. van Luijtelaar. Habituation and sensitization in rat auditory evoked potentials: a single-trial analysis with wavelet denoising. International Journal of Psychophysiology, 43(2):141–153, Feb 2002.
- [47] C. R. Rao and J. Kleffe. Estimation of variance components and applications. Elsevier, Amsterdam, 1988.
- [48] B. Rivet, A. Souloumiac, V. Attina, and G. Gibert. xDAWN Algorithm to enhance evoked potentials: application to brain computer interface. In Biomedical engineering, IEEE Transactions, volume 56, pages 2035–2043, 2009.
- [49] C. Roger, C. G. Bénar, F. Vidal, T. Hasbroucq, and B. Burle. Rostral cingulate zone and correct response monitoring: ICA and source localization evidences for the unicity of correct- and error-negativities. NeuroImage, 51:391–403, February 2010.
- [50] C. Sanelli, C. Vidaurre, K.-R. Müller, and B. Blankertz. Common Spatial Patter Patches - an optimized filter ensemble for adaptive brain-computer interfaces. In 32nd Annual conference of the IEEE EMBS, pages 4351–4354, 2010.
- [51] S. R. Searle, G. Casella, and C. E. McCulloch. Variance components. Wiley, 1992.
- [52] F. T. Smulders, J. L. Kenemans, and A. Kok. Effects of task variables on measures of the mean onset latency of LRP depend on the scoring method. Psychophysiology, 33(2):194–205, Mar 1996.
- [53] J. Spinnato, M.-C. Roubaud, B. Burle, and B. Torrèsani. Finding EEG space-time-scale localized features using matrix-based penalized discriminant analysis.

- In International conference on Acoustics, Speech and Signal Processing (ICASSP 2014), 2014.
- [54] D. Vautrin, X. Artusi, and M.-F. Lucas. A novel criterion of wavelet packet best basis selection for signal classification with application to brain-computer interfaces. IEEE Transactions on Biomedical Engineering, 56(11):2734–2738, November 2009.
- [55] M. Vetterli and J. Kovačević. Wavelets and subband coding. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1995.
- [56] Z. Wang, A. Maier, D. A. Leopold, N. K. Logothetis, and H. Liang. Single-trial evoked potential estimation using wavelets. Computers in Biology and Medicine, 37:463–473, 2007.
- [57] G. M. Weiss. Mining with rarity: A unifying framework. ACM SIGKDD Explorations Newsletter, 6(1):7–19, 2004.
- [58] G. M. Weiss and F. Provost. The effect of class distribution on classifier learning: an empirical study. Technical report, Data, Inference, Analytics and Learning Lab, 2001.
- [59] M. V. Wickerhauser. Adapted Wavelet Analysis from Theory to Software. A.K. Peters, 1996.
- [60] V. Witkovský. Matlab algorithm mixed.m for solving Henderson’s mixed model equations. Technical report, Institute of Measurement Sciences, Slovak Academy of Sciences, 2002.
- [61] J. Xie and Z. Qiu. The effect of imbalanced data sets on LDA: A theoretical and empirical analysis. Pattern Recognition, 40:557–562, 2006.
- [62] J. H. Xue and M. D. Titterington. Do unbalanced data have a negative effect on LDA? Pattern Recognition, 41:1558–1571, 2008.