

A Subject-Independent Brain-Computer Interface based on Smoothed, Second-Order Baseline

Boris Reuderink

Jason Farquhar

Mannes Poel

Anton Nijholt

Abstract—A brain-computer interface (BCI) enables direct communication from the brain to devices, bypassing the traditional pathway of peripheral nerves and muscles. Traditional approaches to BCIs require the user to train for weeks or even months to learn to control the BCI. In contrast, BCIs based on machine learning only require a calibration session of less than an hour before the system can be used, since the machine adapts to the user's existing brain signals. However, this calibration session has to be repeated before each use of the BCI due to inter-session variability, which makes using a BCI still a time-consuming and an error-prone enterprise. In this work, we present a second-order baselining procedure that reduces these variations, and enables the creation of a BCI that can be applied to new subjects without such a calibration session. The method was validated with a motor-imagery classification task performed by 109 subjects. Results showed that our subject-independent BCI without calibration performs as well as the popular common spatial patterns (CSP)-based BCI that does use a calibration session.

I. INTRODUCTION

A brain-computer interface (BCI) automatically interprets electrical signals emitted by the brain, enabling users to interact with a computer without the need of efferent nerve signals. The traditional approach to BCIs is to provide a device that is controlled through a fixed function of the brain signals, and to train users to reliably modify their brain signals — a process that takes weeks, or even months [1]. An alternative, more user friendly approach is to adapt the BCI to the user's naturally occurring brain signals with machine learning (ML) methods (e.g. [2]). This reduces the investment of time necessary for the first use of a BCI from weeks to minutes.

Due to subject-related differences in the electroencephalography (EEG) signals, ML-based BCIs currently still rely on a calibration session in which the user performs a known series of mental tasks before the BCI can be used. Examples of the brain signals associated with these mental tasks are used to automatically train a classifier that is optimized to discriminate between the tasks. But due to the variable nature of the signals, the quality of the discrimination often degrades with time, requiring adaptation or retraining of this subject-dependent (SD) BCI.

Boris Reuderink is with the Human Media Interaction, Faculty of EEMCS, University of Twente, The Netherlands. Email: b.reuderink@ewi.utwente.nl.

Jason Farquhar is with Cognitive Artificial Intelligence, NICI, Radboud University Nijmegen, The Netherlands. Email: j.farquhar@nici.ru.nl.

Mannes Poel is with the Human Media Interaction, Faculty of EEMCS, University of Twente, The Netherlands. Email: m.poel@ewi.utwente.nl.

Anton Nijholt is with the Human Media Interaction, Faculty of EEMCS, University of Twente, The Netherlands. Email: anijholt@cs.utwente.nl.

Obviously, a BCI that does not require this frequent recalibration and can be used immediately is highly desirable for patients and other users alike. Removing this calibration session might even be necessary for large scale adoption. Furthermore, removing the need for calibration implies that the same BCI can discern types of brain signals *independently* of the subject from which the signal is recorded. Such a BCI might provide insights into the invariance characteristics of known neural correlates.

II. PREVIOUS WORK

Recently, progress has been made to make the ML based BCIs generalize to new sessions and new users. The zero training method described in [3] is one of the first attempts to extend the applicability of the popular common spatial patterns (CSP) algorithm to generalize from one session to another session. The method attempts to find prototypical spatial filters from past sessions of a specific subject, and uses a small number of trials of the current session to update the BCI classifier. Using these prototypical filters a performance similar to CSP performance was obtained. Although this result is a promising step towards zero training, historical EEG data and a minimal calibration session are still required.

To overcome these limitations, an ensemble method [4] was developed that selects a sparse set of SD spatio-spectral filters derived from a large database with the recordings of 45 subjects. With a wide-band frequency filter (as used in our study), a subject-independent (SI) classifier performed almost as well (68% correct) as the average SD CSP classifier (70% correct). However, the SI classifier's predictions were post-processed with a non-causal bias-correction, which prevents online application. Without post-processing the best SI classifier still scores much lower than the SD classifiers with 63% of the trials correctly classified.

Combinations of different feature extraction methods and different classifiers were compared on their ability to discriminate between classes of imaginary movement in unseen subjects in [5]. Of all tested combinations, a filter-bank CSP classifier that used frequency filters with different bandwidths had the best SI performance (71% correct). This is slightly above the SI performance of naive log band-power features (68%), and far below the best SD classifier (82%).

These three studies indicate that constructing an SI BCI classifier that generalizes to new subjects is quite challenging. With complex feature extraction as done in [5] and spatial filter matching as done in [3], [4], the performance can approach the SD CSP performance.

In this work, we present a simple second-order baseline (SOB) procedure that reduces inter-session and inter-subject variability, and results in features that can be used to generalize to new, unseen subjects with standard classification methods. In the next section, we will outline our baselining method and describe an off-line experiment used to assess the performance on unseen subjects. Then we will describe and discuss the results, and end with conclusions and remarks for future work.

III. METHODS

During the initiation of imaginary movement, an event-related desynchronization (ERD) in the μ -band is often observed in the motor cortices. For BCIs based on ERD, the spatially filtered EEG is used to compute band-power related features which are then classified with a linear classifier.

Unrelated changes in the EEG signal that manifest over time pose a problem for this classification scheme, because the power, and not the change in power is used to classify the individual trials. To counter this problem, a pre-trial baseline is often used in neuroscientific experiments; for example, the trial power spectrum is often divided by the power spectrum obtained from a pre-trial baseline to study the ERD. Surprisingly, baselining is rarely used in BCIs. In this work we propose a new feature domain for ERD classification that uses a pre-trial baseline to remove the covariance of the EEG channels. Before we describe our baselining approach in more detail, we will outline the CSP pipeline that serves as a control in this work.

A. CSP classification

The CSP algorithm [6], [7] is designed to find a set of m spatial filters that have a maximally different mean variance (band power) for two classes:

$$\Sigma_{WX} = I \quad (1)$$

$$\Sigma_{WX_+} = D, \quad (2)$$

where Σ_{WX} is the channel covariance of the band-pass filtered EEG signal X multiplied by the spatial filter matrix W , X_+ is the EEG signal generated during one specific task, I is the identity matrix and D is a diagonal matrix. The CSP transform can be decomposed into an unsupervised whitening transform to satisfy (1), and a class-specific (supervised) linear transformation to satisfy (2). Usually the $m = 6$ filters corresponding to the $\frac{m}{2}$ smallest and largest eigenvalues in D are used for classification, as extreme eigenvalues represent projections with the greatest mean difference in variance.

After projecting a trial to this $m \times n$ space, the logarithm of the variance of these m projections is typically used as a feature to automatically train a linear classifier. The combination of the feature extraction and a trained classifier results in the follow classification function:

$$f(X(i), \vec{w}, W) = \sum_m w_m \log(\text{diag}(\Sigma_{WX(i)})_m) + w_0 \quad (3)$$

with bias w_0 , feature weights \vec{w} , the m spatial filters W , and trial i 's band-pass filtered signals $X(i)$. The variance of the projections is expressed with the diagonal of the covariance matrix $\text{diag}(\Sigma)$. The logarithm is used to convert the band-power features to an approximately normal distribution as assumed by linear discriminant analysis (LDA) classifiers, but is not strictly necessary for classification.

B. Direct covariance classification

If we reformulate (3) to work in the channel covariance ($\Sigma_{X(i)}$) space and drop the logarithm:

$$f(\Sigma_{X(i)}, \vec{w}, W) = \sum_m w_m (W \Sigma_{X(i)} W^T)_{m,m} + w_0, \quad (4)$$

we can see that $f(\Sigma_{X(i)})$ is just a linear transformation of the vectorized (flattened) trial covariance matrix denoted by $\text{vec}(\Sigma_{X(i)})$:

$$f(\Sigma_{X(i)}, \vec{w}, W) = \sum_m \vec{w}_m (W_{m,\cdot} \Sigma_{X(i)} (W_{m,\cdot})^T) + w_0 \quad (5)$$

$$= \sum (W^T \text{diag}(\vec{w}) W) \circ \Sigma_{X(i)} + w_0 \quad (6)$$

$$= \vec{u}^T \text{vec}(\Sigma_{X(i)}) + w_0, \quad (7)$$

where \circ is the Hadamard (element wise) product, $W_{m,\cdot}$ is the m -th row of W , and $\text{diag}(\vec{w})$ is a diagonal matrix containing the values of \vec{w} on its diagonal. The combination of the spatial filters and the band-power feature weights $\vec{u} = \text{vec}(W^T \text{diag}(\vec{w}) W)$ can thus be modeled directly with a regularized linear classifier [8]. This simplification enables the classifier to learn spatial filters simultaneously with the projection's variance weighting in a single, supervised learning step.

For direct covariance classification, we have chosen to explicitly decorrelate the channels with a symmetric whitening transform P that is estimated on the training trials to satisfy (1):

$$P = \Sigma_X^{-\frac{1}{2}} = U \Lambda^{-\frac{1}{2}} U^T \quad \text{with} \quad U \Lambda U^T = \Sigma_X, \quad (8)$$

and only learn the rotational spatial filters and their associated weights implicitly as in (7) with a linear support vector machine (SVM). Without this whitening transform, the SVM's ℓ_2 regulariser strongly biases the classifier to focus on high-powered sources that are probably not task-relevant.

In summary, we used the covariance of whitened trials as features to directly train a linear classifier that is, except for the log transform, almost equivalent to the commonly used CSP pipeline.

C. Covariance classification with second order baseline

Likewise, the proposed second-order baseline (SOB) method learns the spatial filters implicitly, but instead of the static whitening transform (8) a pre-trial baseline is used to *adaptively* normalize ongoing second order covariance statistics.

We estimated a whitening transform $P(i)$ for each trial i based on past pre-trial baselines, and applied this transform to the data of the trial itself. Without a change in brain activity, the covariance during normalized trial $P(i)X(i)$ would approximate the identity matrix, even during (slow) sensor covariance changes. But when the coupling between or the power in certain brain regions changes, a perturbation appears in the sensor covariance during the normalized trial. This perturbation can be used for classification. The specific whitening transform (8) has the crucial property that it removes correlations, but at the same time maximizes the correspondence between the projection and the original signals. This property preserves the task-relevant topography, which is needed to have consistent features over time and over subjects. A similar normalization procedure was outlined in [9] to adapt the session covariance matrix in order to reduce the influence of non-stationarities. The main differences between our method and [9] is that in our method each trial is normalized differently based on the pre-trial baseline, and that this baseline period is used to estimate the resting state covariance instead of the global session covariance.

Estimating the symmetrical whitening transform from the pre-trial baseline covariance $\Sigma_{B(i)}$ for trial i is difficult, since there is a large number of parameters to estimate from a few independent samples for the s sensors. To improve the robustness, we used the regularized Ledoit-Wolf covariance estimator [10], and used an exponentially weighted moving average (EMWA) to combine the baseline covariance of past trials into a covariance estimate $\hat{\Sigma}_{B(i)}$ for the baseline of trial i :

$$\hat{\Sigma}_{B(i)} = \alpha S_{B(i)}^* + (1 - \alpha) \hat{\Sigma}_{B(i-1)}, \quad (9)$$

where $S_{B(i)}^*$ is the Ledoit-Wolf covariance estimate of the baseline before trial i , and α is known as the forgetting factor that determines the rate of adaptation. With $\alpha = 1 - \sqrt[n]{\frac{1}{2}}$, the forgetting factor α is then associated with a decay that halves in n trials.

Specifically, we calculate a whitening transform $P(i)$ for each of the trials $X(i)$ based on $\hat{\Sigma}_{B(i)}$, and apply this $P(i)$ to the Ledoit-Wolf covariance estimate of the current trial $S_{X(i)}^*$:

$$\tilde{X}(i) = P(i) S_{X(i)}^* P(i)^T \quad \text{with} \quad P(i) = \Sigma_{B(i)}^{-\frac{1}{2}}. \quad (10)$$

The new features $\text{vec}(\tilde{X}(i))$ are more robust over time and subject related variations, but are still sensitive to task-related (co)variance changes.

D. Dataset

To evaluate the performance of the invariant features we used the movement imagery dataset¹ from the experiment described in [11] contributed to Physiobank [12]. This dataset contains sessions of 109 different subjects with trials for actual and imagined movement. We have chosen to use the

blocks where the subjects had to imagine either movement of both feet or movement with both hands, as a preliminary experiment indicated that BCI classification performance above chance level can be obtained with small training sets.

E. Preprocessing

To pre-process the data, we applied a 6th-order Butterworth notch filter at 60 Hz, applied a 6th-order Butterworth filter between 8–30 Hz and extracted trials for movement imagery of both hands or of both feet in the interval from $[-2, 4]$ s after the stimulus. All evaluated methods used the same interval $[0.5, 4]$ s after the stimulus presentation for classification. For the SOB method, the interval $[-2, 0]$ s was used to estimate the pre-trial baseline. The same pre-processing was used for all BCI classifiers.

F. Evaluation

We used two CSP based pipelines and a log band-power (logBP) based pipeline as a comparison method in both an SD and an SI BCI classification scheme. One CSP pipeline was based on CSP projected log band-power features classified with a LDA, the other CSP classifiers used band-power features without the log transform, classified with a linear soft-margin SVM [13]. The logBP classifiers simply used the variance of each band-pass filtered channel as a feature for a linear SVM. The whitened covariance features and SOB normalized covariance features were also classified with a linear SVM. The SVM's c -parameter was always estimated on the training set using a sequential (chronological) 5-fold cross-validation procedure with a logarithmic step size for the c -values.

To evaluate these classifiers in an SD context, the first half of the session was used for training, and the second half was used for evaluation. This simulates a calibration and application session, respectively. Chronological separation of training and test set is needed since random splits lead to overly optimistic performance estimates. Because the dataset also contains blocks with other mental tasks, we have only about 22 trials in total for training and 22 trials for evaluation per subject.

The performance of SI application was assessed by training an SI classifier on the first 50 users, and then applying the classifier to the test set formed from the remaining 51 subjects (we removed 8 subjects from the test pool because they had fewer trials). The final SI performance was calculated on the second half of the predictions for these 51 test subjects to allow for a paired comparison with the SD classifiers.

IV. RESULTS

A. Subject-dependent classification

The performance of the various control features is shown in Table I, as well as the performance on the newly proposed SOB covariance features. For subject-dependent (SD) classification, the LDA classification of CSP log-variance features had the highest mean accuracy. However, there was no significant difference between the performance of this

¹<http://www.physionet.org/pn4/eegmidb/>

TABLE I

THE SUBJECT-DEPENDENT ACCURACY OF THE DIFFERENT PIPELINES ON THE LAST 22 TRIALS OF EACH SESSION FOR THE 51 TEST SUBJECTS.

Pipeline	Half life	Mean (std.)	Median
logBP SVM		62.2 (11.5)	63.6
CSP logvar LDA		69.5 (14.6)	68.2
CSP var SVM		68.6 (15.0)	68.2
whcov SVM		68.9 (15.2)	72.7
SOB cov SVM (best)	13.8 trials	67.1 (13.3)	68.2
SOB cov SVM (worst)	1 trial	64.8 (13.1)	63.6

TABLE II

THE SUBJECT-INDEPENDENT ACCURACY OF THE DIFFERENT PIPELINES ON THE LAST 22 TRIALS OF EACH SESSION FOR THE 51 TEST SUBJECTS.

Pipeline	Half life	Mean (std.)	Median
logBP SVM		58.1 (11.1)	54.5
CSP logvar LDA		59.3 (11.8)	54.5
CSP var SVM		56.4 (9.7)	54.5
whcov SVM		59.1 (10.8)	54.5
SOB cov SVM (best)	4.0 trials	67.3 (13.4)	68.2
SOB cov SVM (worst)	18.9 trials	64.9 (13.0)	63.6

CSP pipeline and direct covariance classification (whcov SVM), and the latter had a higher median performance (72.7% accuracy). This demonstrates the feasibility of direct covariance classification. The logBP features performed much worse than the spatially filtered alternatives. The SOB features worked almost as well as the CSP features when low α -values were used; with faster adaption rates the SOB did not perform as well with SD application.

B. Subject-independent classification

The performance of the various control methods severely degraded due to inter-subject variability (Table II) when these classifiers were applied in a subject-independent (SI) fashion. The CSP based classifiers, which did outperform the naive logBP classifier with subject-dependent training, now performed at the same level as logBP classifier with SI

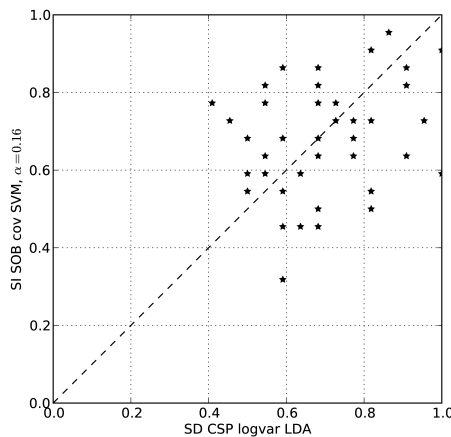


Fig. 1. The accuracy of a subject-dependent CSP pipeline versus the performance of a subject-independent SOB pipeline. There is no significant difference between the classifiers, despite the fact that the SOB pipeline was not trained on the subject.

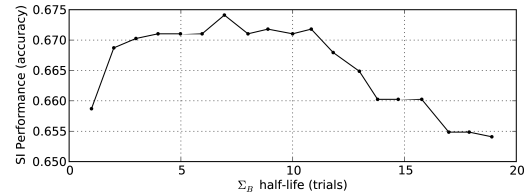


Fig. 2. The mean SI accuracy of the SOB as a function of the half-life of the baseline estimate $\Sigma_{B(i)}$ of the predictions on the last 22 trials of all the 59 test subjects.

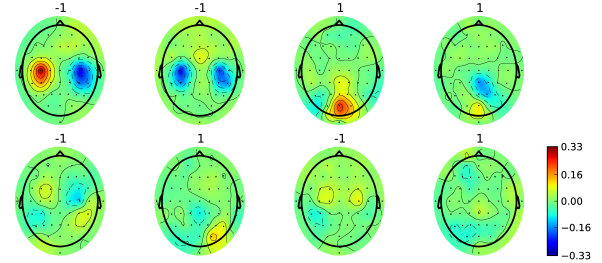


Fig. 3. The most influential spatial filters W_i , for the best SI SOB classifier scaled by the magnitude of their weight w_i . The number above the plot is the sign of the weight w_i (i.e. a positive sign indicates filters with a response that corresponds to imagery of foot movement, a negative sign indicates imagined hand movement). Most of the contribution seems to originate from the motor areas, the central parietal regions and the visual cortex.

application — the advantage that spatial filtering provided in the SD training disappeared with SI application. The performance of the SOB based predictions however, was not affected at all. The accuracy of the best subject-independent SOB-based predictions was not even significantly different from the best subject-dependent (CSP log-variance LDA based) predictions ($p=0.16$ with a Wilcoxon signed-rank test on 51 paired observations, see Fig. 1).

The best SI performance was obtained with a volatile baseline with a half life of 4 trials, while with SD application the best results were obtained with long half-lives (low α 's). Note that even the worst performing SOB classifier outperformed all of the control classifiers. Fig. 2 displays the fraction of correctly classified trials as a function of the amount of smoothing of the pre-trial baseline covariance. The best performance was obtained with a baseline half-life between 2–11 trials.

The most contributing spatial filters that were learned implicitly by the SI SOB covariance classifier are shown in Fig. 3. These filters can be easily extracted with an eigenvalue decomposition of the covariance matrix, see (6). The most relevant features originated from the motor cortex region, but also occipital and central parietal features contributed to the classifiers predictions. There was no apparent contribution of muscle or eye movement artifacts to the classification.

V. DISCUSSION

The results indicate that the new second-order baseline covariance features provide a robust alternative to CSP features for classification of motor imagery, and generalize

to new, unseen subjects without additional calibration or training. Apparently, the normalization performed with the SOB removes enough inter-subject variability to generalize to new subjects. However, the dataset used in this research contains rather few trials, hence the SD CSP performance might have suffered from insufficient training data. Nevertheless, recording more trials is not always an option, and the SD performance obtained in our study is similar to scores presented in [4] where much longer sessions were used.

The SOB's α parameter seems of some importance for generalization over subjects. While for SD classification a long half-life was preferred, α 's with a short half-life are preferred for SI classification. Presumably, slow adaptation is preferred for SD classification because it allows the classifier to model and exploit session-specific variations, such as for example bad channels and EEG artifacts. For SI classifications modeling these variations is generally not helpful as they are not consistent over subjects. Shorter half-lives reduce these variations, and are thus preferred for SI classification.

It is noteworthy to mention that the best α -value for SOB-covariance features was selected based on the performance on the test subjects. This might slightly overestimate the true performance. Usually these hyper-parameters (e.g. the SVM's c -parameter) are set based on performance estimates obtained with cross-validation. The half life constant α could be chosen with cross-validation, but since the SOB is a preprocessing method the time and space complexity is often prohibitive. Current BCI pipelines have several preprocessing hyper-parameters that are fixed a priori (e.g. the cut-off values in the band-pass filter, or the $m = 6$ spatial filters). Given that even the worst α performs better than the alternatives in SI classification, the performance gain seems fairly robust for a wide range of α -values (Fig. 2). Therefore, we expect that choosing $\alpha = 0.16$ (a half-life of 4 trials) a priori will be adequate in practise.

VI. CONCLUSIONS AND FUTURE WORK

We presented an SOB procedure that reduces inter-subject and inter-session variability, and demonstrated that SOB-covariance features allow for cross-subject motor imagery classification without a loss of performance compared to within-subject classification with the popular CSP based BCI classifier. The advantage of the SOB based covariance features is that they are robust to inter-session and inter-subject variation, and that standard classifiers such as the SVM can be used without the need of adaptation or post-processing of the outputs, such as done with bias-adaptation. Furthermore, the online processing is simplified as it can be implemented as a stateless, fixed pipeline that does not handle the incoming data differently during a calibration or online application session.

In addition to the practical advantages of removing the need for the laborious calibration sessions, changing from subject-dependent to subject-independent BCIs also simplifies multi-disciplinary BCI research. It allows researchers to work with validated BCI classifiers that are known to

work with a certain probability on the target population, and focusses on the intended brain regions. The development of subject-independent BCIs can facilitate new applications for which collecting enough subject-specific training data before each session is not feasible, such as for example fatigue detection, screening of neurological disorders or classification of emotional states.

The method described uses a single, broad frequency band. For future work, the features can be extended to include multiple frequency bands as in [5], [8]. Another interesting research direction is to generalize to recordings with different electrode layouts. As the learned covariance weights were quite sparse, the correspondence of a few key sensor locations might be enough to generalize to new sensor configurations. Finally, although the presented method works causally, it should be validated in an online experiment with a user in the loop.

ACKNOWLEDGEMENTS

The authors gratefully acknowledge the support of the BrainGain Smart Mix Programme of the Netherlands Ministry of Economic Affairs and the Netherlands Ministry of Education, Culture and Science.

REFERENCES

- [1] J. R. Wolpaw, N. Birbaumer, D. J. McFarland, G. Pfurtscheller, and T. M. Vaughan, "Brain-computer interfaces for communication and control," *Clinical Neurophysiology*, vol. 113, no. 6, pp. 767–791, 2002.
- [2] B. Blankertz, R. Tomioka, S. Lemm, M. Kawanabe, and K.-R. Müller, "Optimizing spatial filters for robust EEG single-trial analysis," *Signal Processing Magazine*, vol. 25, no. 1, pp. 41–56, 2008.
- [3] M. Krauledat, M. Tangermann, B. Blankertz, and K.-R. Müller, "Towards zero training for brain-computer interfacing," *PLoS ONE*, vol. 3, p. e2967, 2008.
- [4] S. Fazli, F. Popescu, M. Danóczy, B. Blankertz, K.-R. Müller, and C. Grozea, "Subject-independent mental state classification in single trials," *Neural Networks*, vol. 22, pp. 1305–1312, 2009.
- [5] F. Lotte, C. Guan, and K. K. Ang, "Comparison of designs towards a subject-independent brain-computer interface based on motor imagery," in *Proceedings of the 31st Annual International Conference of the IEEE EMBS*, 2009, pp. 4543–4546.
- [6] Z. J. Koles, "The quantitative extraction and topographic mapping of the abnormal components in the clinical EEG," *Electroencephalography and Clinical Neurophysiology*, vol. 79, no. 6, pp. 440–447, December 1991.
- [7] H. Ramoser, J. Müller-Gerking, and G. Pfurtscheller, "Optimal spatial filtering of single trial EEG during imagined hand movement," *IEEE Transactions on Rehabilitation Engineering*, vol. 8, no. 4, pp. 441–446, 2000.
- [8] J. Farquhar, "A linear feature space for simultaneous learning of spatio-spectral filters in BCI," *Neural Networks*, vol. 22, pp. 1278–1285, 2009.
- [9] R. Tomioka, J. Hill, B. Blankertz, and K. Aihara, "Adapting spatial filtering methods for nonstationary BCIs," in *2006 Workshop on Information-Based Induction Sciences (IBIS2006)*, 2006.
- [10] O. Ledoit and M. Wolf, "A well-conditioned estimator for large-dimensional covariance matrices," *Journal of Multivariate Analysis*, vol. 88, no. 2, pp. 365–411, February 2004.
- [11] G. Schalk, D. J. McFarland, T. Hinterberger, N. Birbaumer, and J. R. Wolpaw, "BCI2000: a general-purpose brain-computer interface (BCI) system," *IEEE Transactions on Biomedical Engineering*, vol. 51, no. 6, pp. 1034–1043, 2004.
- [12] A. L. Goldberger, L. A. N. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, "PhysioBank, PhysioToolkit, and PhysioNet," *Circulation*, vol. 101, no. 23, pp. e215–e220, 2000.
- [13] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297.