



Technical Report No. 120

Support Vector Channel Selection in BCI

Thomas Navin Lal¹, Michael Schröder²,
Thilo Hinterberger³, Jason Weston¹,
Martin Bogdan², Niels Birbaumer³, and
Bernhard Schölkopf¹

December 2003

¹ Max-Planck-Institut for Biological Cybernetics, Tübingen, Germany.
{navin|jason.weston|bs}@tuebingen.mpg.de

² Eberhard Karls University Tübingen, Department of Computer Engineering, Tübingen, Germany.
{schroedm|bogdan}@informatik.uni-tuebingen.de

³ Eberhard Karls University Tübingen, Institute of Medical Psychology and Behavioral Neurobiology,
Tübingen, Germany. {thilo.hinterberger|niels.birbaumer}@uni-tuebingen.de

Support Vector Channel Selection in BCI

*Thomas Navin Lal, Michael Schröder,
Thilo Hinterberger, Jason Weston,
Martin Bogdan, Niels Birbaumer, and
Bernhard Schölkopf*

Abstract. Designing a Brain Computer Interface (BCI) system one can choose from a variety of features that may be useful for classifying brain activity during a mental task. For the special case of classifying EEG signals we propose the usage of the state of the art feature selection algorithms Recursive Feature Elimination [1] and Zero-Norm Optimization [2] which are based on the training of Support Vector Machines (SVM) [3]. These algorithms can provide more accurate solutions than standard filter methods for feature selection [4].

We adapt the methods for the purpose of selecting EEG channels. For a motor imagery paradigm we show that the number of used channels can be reduced significantly without increasing the classification error. The resulting best channels agree well with the expected underlying cortical activity patterns during the mental tasks.

Furthermore we show how time dependent task specific information can be visualized.

1 Introduction

Most Brain Computer Interfaces (BCIs) make use of mental tasks that lead to distinguishable EEG signals of two or more classes. For some tasks the relevant recording positions are known, especially when the tasks comprise motor imagery, e.g. the imagination of limb movements, or the overall activity of large parts of the cortex that occurs during intentions or states of preparation and relaxation.

For the development of new paradigms whose neural correlates are not known in such detail, finding optimal recording positions for use in BCIs is challenging. New paradigms can become necessary when motor cortex areas show lesions, for the increase of the information rate of BCI systems or for robust multi-class BCIs. If good recording positions are not known, a simple approach is to use data from as many as possible EEG electrodes for signal classification. The drawback of this approach is that the extend to which fea-

ture selection and classification algorithms overfit to noise increases with the number of task-irrelevant features, especially when the ratio of training points and number of features is small. In addition it is difficult to understand which part of the brain generates the class relevant activity.

We show that the selection of recording positions can be done robustly in the absence of prior knowledge about the spatial distribution of brain activity of a mental task. Specifically we adapt the state of the art feature selection methods *Zero-Norm Optimization* (l0-Opt) and *Recursive Feature Elimination* (RFE) to the problem of channel selection and demonstrate the usefulness of these methods on the well known paradigm of motor imagery.

The paper is structured as follows: section 2 contains the experimental setup, the task, and the basic data preprocessing. In section 3 the feature selection methods and the classification algorithm are described. Results are given in section 4 and the final section concludes.

2 Data acquisition

2.1 Experimental setup and mental task

We recorded EEG signals from eight untrained right handed male subjects using 39 silver chloride electrodes (see figure 1). The reference electrodes were positioned at TP9 and TP10. The two electrodes Fp2 and 1cm lateral of the right eye (EOG) were used to record possible EOG artifacts and eye blinks while two fronto-temporal and two occipital electrodes were positioned to detect possible muscle activity during the experiment. Before sampling the data at 256 Hz an analog bandpass filter with cutoff frequencies 0.1 Hz and 40 Hz was applied.

The subjects were seated in an armchair at 1m distance in front of a computer screen. Following the experimental setup of [5] the subjects were asked to imagine left versus right hand movements during each trial. With every subject, we recorded 400 trials during one single session. The total length of each trial was 9 seconds. Additional inter-trial intervals for relaxation varied randomly between 2 and 4 seconds. No outlier detection was performed and no trials were removed during the data processing at any stage.

Each trial started with a blank screen. A small fixation cross was displayed in the center of the screen from second 2 to 9. A cue in the form of a small arrow pointing to the right or left side was visible for half a second starting with second 3. In order to avoid event related signals in later processing stages only data from seconds 4 to 9 of each trial was considered for further analysis. Feedback was not provided at any time.

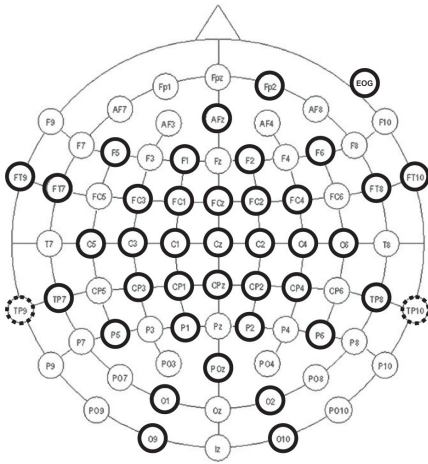


Figure 1: The position of 39 EEG electrodes used for data acquisition are marked in solid black circles. The two referencing electrodes are marked in dotted circles.

2.2 Pre analysis

As Pfurtscheller and da Silva have reported [6] that movement related desynchronization of the μ -rhythm (8-12 Hz) is not equally strong in subjects and might even fail for various reasons (e.g. because of too short inter-trial intervals that prevent a proper re-synchronization) we performed a pre analysis in order to identify and exclude subjects that did not show significant μ -activity at all.

For seven of the eight subjects the μ -band was only slightly differing from the 8-12 Hz usually given in the EEG literature. Only one subject showed scarcely any activity in this frequency range but instead a recognizable movement related desynchronization in the 16-20Hz band.

Restricted to only the 17 EEG channels that were located over or close to the motor cortex we calculated the maximum energy of the μ -band using the Welch method [7] for each subject. This feature extraction resulted in one parameter per trial and channel and explicitly incorporated prior knowledge about the task.

The eight data sets consisting of the Welch-features were classified with linear SVMs (see below) including individual model selection for each subject. Generalization errors were estimated by 10-fold cross validation. As for three subjects the pre analysis showed very poor error rates close to chance level their data sets were excluded from further analysis.

2.3 Data preprocessing

For the remaining five subjects the recorded 5s windows of each trial resulted in a time series of 1280 sample points per channel. We fitted an autoregres-

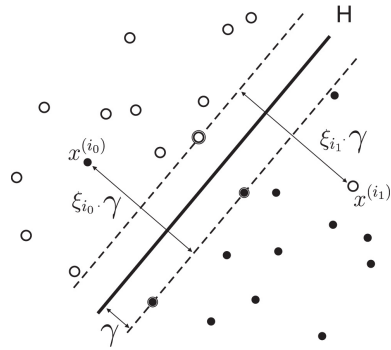


Figure 2: Linear SVM. For non separable data sets, slack variables ξ_i are introduced. The thick points on the dashed lines are called support vectors (SVs). The solution for the hyperplane H can be written in terms of the SVs. For more detail see section 3.1.

sive (AR) model of order 3 to the time series¹ of all 39 channels using forward backward linear prediction [8]. The three resulting coefficients per channel and trial formed the new representation of the data.

The extraction of the features did not explicitly incorporate prior knowledge although autoregressive models have successfully been used for motor related tasks (e.g. [5]). However, they are not directly linked to the μ -rhythm.

2.4 Notation

Let n denote the number of training vectors (trials) of the data sets ($n = 400$ for all five data sets) and let d denote the data dimension ($d = 3 \cdot 39 = 117$ for all five data sets). The training data for a classifier is denoted as $X = (x^{(1)}, \dots, x^{(n)}) \in \mathbb{R}^{n \times d}$ with labels $Y = (y_1, \dots, y_n) \in \{-1, 1\}^n$. For the task used in this paper $y = -1$ denotes imagined left hand movement, $y = 1$ denotes imagined right hand movement. The terms *dimension* and *feature* are used synonymously. For $l \in \mathbb{N}$, $l > 1$ the set $M^{-j} \subset \mathbb{R}^{l-1}$ is obtained from a set $M \subset \mathbb{R}^l$ by removing the dimension j from every point $m \in M$ (canonical projection).

3 Feature selection and classification methods

Feature selection algorithms can be characterized as either filter or wrapper methods [9]. They select or omit dimensions of the data depending on a performance measure.

¹For this choice we compared different model orders. For a given order we fitted an AR-model to each EEG sequence. After proper model selection a Support Vector Machine with 10-fold cross validation (CV) was trained on the coefficients. Model order 3 resulted in the best mean CV error.

The problem of how to rate the relevance of a feature if nonlinear interactions between features are present is not trivial, especially since the overall accuracy might not be monotonic in the number of features used. Some feature selection methods try to overcome this problem by optimizing the feature selection for subgroups of fixed sizes (plus-1 take-away-r search) or by implementing floating strategies (e.g. floating forward search) [9]. Only few algorithms like e.g. genetic algorithms can choose subgroups of arbitrary size during the feature selection process. They have successfully been used for the selection of spatial features [10] in BCI applications but are computationally demanding.

For the application of EEG channel selection, it is necessary to treat a certain kind of grouped features homogeneously: numerical values belonging to one and the same EEG channel have to be dealt with in a congeneric way so that a spatial interpretation of the solution becomes possible. We adapted the state of the art feature selection methods *Zero-Norm Optimization* and *Recursive Feature Elimination* (RFE) as well as the Fisher Correlation to implement these specific requirements. The first two algorithms are closely related to Support Vector Machines (SVM).

3.1 Support Vector Machines (SVMs)

The Support Vector Machine is a relatively new classification technique developed by V. Vapnik [3] which has shown to perform strongly in a number of real-world problems, including BCI [11]. The central idea is to separate data $X \subset \mathbb{R}^d$ from two classes by finding a weight vector $w \in \mathbb{R}^d$ and an offset $b \in \mathbb{R}$ of a hyperplane

$$\begin{aligned} H : \mathbb{R}^d &\rightarrow \{-1, 1\} \\ x &\mapsto \text{sign}(w \cdot x + b) \end{aligned}$$

with the largest possible margin², which apart from being an intuitive idea has been shown to provide theoretical guaranties in terms of generalization ability [3]. One variant of the algorithm consists of solving the following optimization problem:

$$\begin{aligned} \min_{w \in \mathbb{R}^d} \quad & \|w\|_2 + C \sum_{i=1}^n \xi_i^2 \\ \text{s.t.} \quad & y_i(w \cdot x^{(i)} + b) \geq 1 - \xi_i \quad (i = 1, \dots, n) \end{aligned} \quad (1)$$

The parameters ξ_i are called slack variables and ensure that the problem has a solution in case the data is not linear separable³ (see figure 2). The margin is defined as $\gamma(X, Y, C) = 1/\|w\|_2$. In practice one has to

²Is X linear separable the margin of a hyperplane is the distance of the hyperplane to the closest point $x \in X$.

³Is the data linear separable the slack variables can improve the generalization ability of the solutions.

trade-off between a low training error, e.g. $\sum \xi_i^2$, and a large margin γ . This trade-off is controlled by the regularization parameter C . Finding a good value for C is part of the model selection procedure. If no prior knowledge is available C has to be estimated from the training data, e.g. by using cross validation. The value $2/C$ is also referred to as the *ridge*. For a detailed discussion please refer to [12].

3.2 Fisher Criterion (FC)

The Fisher Criterion determines how strongly a feature is correlated with the labels [13]. For a set $T = \{t^{(1)}, \dots, t^{(|T|)}\} \subset \mathbb{R}^d$ define the mean $\mu_j(T) = \frac{1}{|T|} \sum_{i=1}^{|T|} t_j^{(i)}$ and the variance $V_j(T) = \frac{1}{|T|} \sum_{i=1}^{|T|} (t_j^{(i)} - \mu_j(T))^2$ ($j = 1, \dots, d$). The score R_j of feature j is then given by:

$$R_j(X) = \frac{(\mu_j(X^+) - \mu_j(X^-))^2}{V_j(X^+) + V_j(X^-)}, \quad (2)$$

with $X^+ := \{x_i \in X \mid y_i = 1\}$ and X^- similarly. The rank of a channel is simply set to the mean score of the corresponding features.

3.3 Zero-Norm Optimization (l0-Opt)

Weston *et. al.* [2] recently suggested to minimize the zero-norm⁴ $\|w\|_0 := \text{cardinality}(\{w_j : w_j \neq 0\})$ instead of minimizing the l_1 -norm or l_2 -norm as in standard SVMs (cp. equation (1)):

$$\begin{aligned} \min_{w \in \mathbb{R}^d} \quad & \|w\|_0 + C \|\xi\|_0 \\ \text{s.t.} \quad & y_i(w \cdot x^{(i)} + b) \geq 1 - \xi_i \quad (i = 1, \dots, n). \end{aligned} \quad (3)$$

The solution of this optimization problem is usually much sparser than the solution of problem (1). Thus feature selection is done implicitly. Unfortunately the problem has shown to be NP-hard but the authors developed an iterative method to approximate the solution. In case the solution w^* has less than the desired number of zero entries, the remaining features $\{j\}$ can be ranked according to w_j^* (as in one iteration step of RFE).

In the original version of the method the features are multiplied with a scaling factor during each iteration. Once a scaling factor is zero, the corresponding feature is removed. We adapt this method in the following way: the scaling factors of the features corresponding to a channel are substituted by their mean. Thus all features of one channel are either removed completely (the channel is removed) or all features remain. As in the case of SVM and RFE, the parameter C has to be estimated from the training data in case prior knowledge is not available.

⁴The zero-norm of a vector v is equal to number of nonzero entries of v .

3.4 Recursive Feature Elimination (RFE)

This feature selection method was proposed by Guyon et al. [14] and is based on the concept of margin maximization. The importance of a dimension is determined by the influence it has on the margin of a trained SVM. Let W be the inverse of the margin:

$$W(X, Y, C) := \frac{1}{\gamma(X, Y, C)} = \|w\|_2$$

At each iteration one SVM is trained and the features \hat{j} which minimize $|W(X, Y, C) - W(X^{-j}, Y^{-j}, C)|$ are removed (typically that is one feature only); this is equivalent to removing the dimensions \hat{j} that correspond to the smallest $|w_j|$. We adapt this method for channel selection in the following way:

Let $F_k \subset \{1, \dots, d\}$ denote the features from channel k . Similar to the reformulation of the Fisher Criterion and the Zero-Norm-Optimization we define for each channel k the score $s_k := \frac{1}{|F_k|} \sum_{l \in F_k} |w_l|$. At each iteration step we remove the channel with the lowest score. The parameter C has to be estimated from the training data, if no prior knowledge is available.

For the remainder of the paper we refer to the adapted feature selection methods as channel selection methods. Furthermore we will also refer to the adapted RFE as *Recursive Channel Elimination*.

3.5 Generalization Error Estimation

For model selection purposes we estimated the generalization error of classifiers via 10-fold cross validation.

If the generalization error of a channel selection method had to be estimated, a somewhat more elaborated procedure was used. An illustration of this procedure is given in figure 3.

The whole data set is split up into 10 folds ($F1$ to $F10$) as for usual cross validation. In each fold F , the channel selection (CS in figure 3) is performed based on the train set of F only, leading to a specific ranking of the 39 EEG channels. For each fold F , 39 classifiers C_F^h , $h = 1, \dots, 39$ are trained as follows: C_F^h is trained on the h best⁵ channels, respectively, of the train set of F and tested on the corresponding channels of the test set of F . For each fold, this results in 39 test errors (E_F^1 to E_F^{39}).

During the last step, the corresponding test errors are averaged over all folds. This leads to an estimate of the generalization error for every number of selected channels.

⁵In this context, *best* means according to the calculated ranking of that fold.

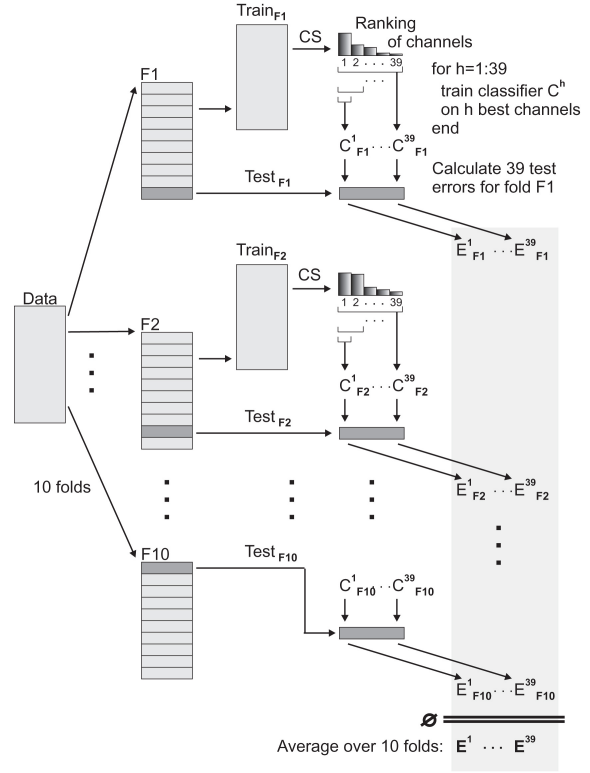


Figure 3: Illustration of the procedure for channel selection and error estimation using cross validation.

4 Results

4.1 Channel Selection

We applied the three channel selection methods Fisher Criterion, Recursive Feature Elimination and Zero-Norm Optimization introduced in section 3 to the five data sets. As the experimental paradigm is well known we could examine the results concerning their physiological plausibility. Therefore we investigated whether the best ranked channels are those situated over or close to motor areas. Furthermore we analyzed if the number of channels can be reduced without a loss of accuracy in terms of cross validation error.

Initial to the channel selection and individually for each subject s , the regularization parameter C_s for later SVM trainings was estimated via 10-fold cross validation from the training data sets⁶.

The estimation of the generalization error for all 39 stages of the channel selection process⁷ was carried out using linear SVMs as classifiers with parameters C_s previously determined. Details about the 10-fold

⁶Estimating the parameter for each number of channels in the process of channel selection might improve the accuracy. However the chance of overfitting increases.

⁷In fact, methods RFE and 10-Opt perform rather a channel *removal* than a channel selection.

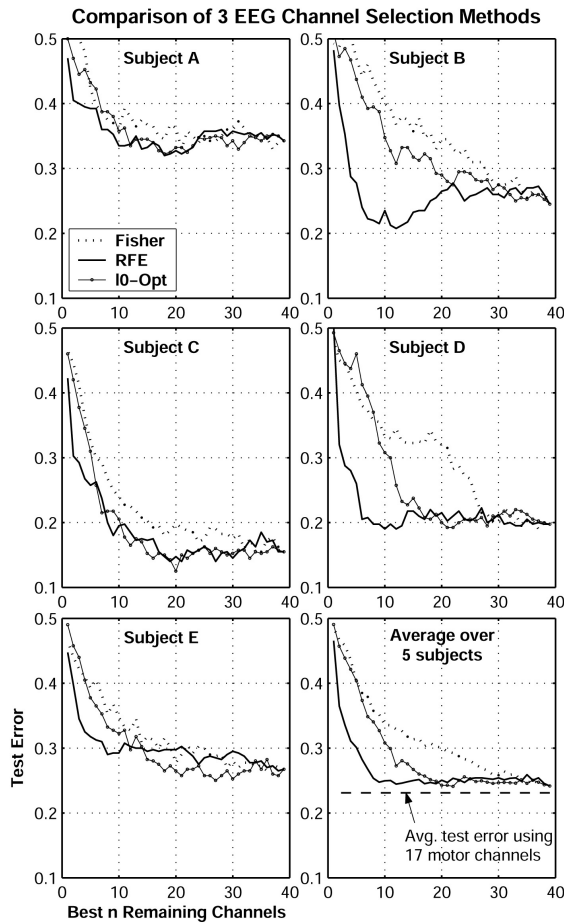


Figure 4: Comparison of the three channel selection methods *Fisher Score*, *RFE* and *l0-Opt* individually for five subjects and averaged over the subjects. Method *RFE* allows the strongest reduction of number of channels for all subjects.

cross validation during the estimation procedure are described in section 3.5 and figure 3.

The estimation results are depicted in figure 4. The first five plots show the individual generalization error for the five subjects against the different numbers of channels chosen by the three channel selection methods. The sixth plot in the bottom right corner shows the generalization error of the three methods averaged over the five subjects.

Recursive Feature Elimination and Zero-Norm Optimization prove to be capable of selecting relevant channels, whereas the Fisher Criterion fails for some subjects. Especially for small numbers of channels *RFE* is slightly superior over the Fisher Criterion and Zero-Norm Optimization. For larger numbers of channels the performance of *l0-Opt* is comparable to *RFE*. As can be seen in figure 4 it is possible to reduce the number of EEG channels significantly using the *RFE* method - for the investigated experimental paradigm

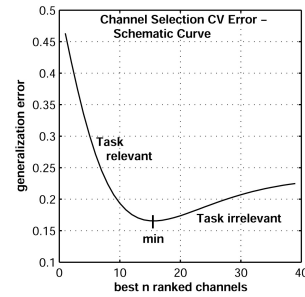


Figure 5: Idealized generalization error curve using a channel selection method in the presence of irrelevant channels. When removing channels iteratively the classification error decreases slightly until all irrelevant channels are removed. Removing more channels results in an increase of error.

this can be done without a loss of classification accuracy. E.g. using 8 channels for subject *D* yields the same error as the error obtained using all channels. On the data set of subject *B* the cross validation error of 24.5% can be reduced to 20.75% using 12 channels only.

It is not tractable to test all ($\approx 10^{11}$) possible combinations of channels to find the best combination. In this light the 17 channels located over or close to the motor cortex can be considered a very good solution that is close to the optimal one. For rating the overall accuracy of the *RFE* method we thus trained a classifier using these 17 channels. The result averaged over the five subjects is plotted as a baseline in the last figure. The average error rate (taken over all subjects) of 24% using 12 channels is very close to the error of the baseline which is 23%.

Table 1 contains channel rankings, which are obtained by applying Recursive Channel Elimination to the data set of each subject⁸. As the *RFE* method has outperformed *CF* and *l0-Opt*, the rankings in table 1 were exclusively calculated by *RFE*.

To interpret the table it is useful to have a closer look at figure 5. It shows an idealized curve for an estimate of the generalization error when using a channel or feature selection method. As we have also seen in the experiments it is possible to reduce the number of channels without a loss of accuracy. For each subject we can obtain a heuristic estimate on the number of irrelevant channels from the generalization error curves in figure 4. We underlined one entry in each column of table 1. The row number of that entry is an estimate for the rank position that divides task relevant channels from task irrelevant ones. E.g. for subject *D* figure 4 shows a local minimum of the *RFE* generalization error curve at 10 channels. Thus the best 10 selected

⁸Please note that in this step cross validation was not applied.

Table 1: RFE Ranking of 39 EEG Channels

| Rank | Subjects | | | | |
|------|----------|------|------|------|------|
| | A | B | C | D | E |
| 1 | C4 | CP4 | CP4 | FC4 | CP4 |
| 2 | CP4 | C3 | CP3 | C4 | CPz |
| 3 | CP2 | C4 | C4 | CP2 | C2 |
| 4 | C2 | FC4 | C2 | CP1 | FC3 |
| 5 | Cz | FT9 | C1 | C3 | C4 |
| 6 | FC4 | FT10 | CPz | FC3 | C1 |
| 7 | FC2 | CP1 | CP2 | C2 | FCz |
| 8 | C3 | C1 | C3 | C1 | FC4 |
| 9 | CP3 | F6 | F1 | FC2 | C3 |
| 10 | F1 | Fp2 | FC1 | FC1 | POz |
| 11 | F2 | FC1 | FC2 | FT10 | P6 |
| 12 | C1 | AFz | C5 | FCz | O10 |
| 13 | FC3 | C2 | FT7 | F2 | FC1 |
| 14 | CPz | P6 | F2 | FT9 | C6 |
| 15 | CP1 | CP2 | FC3 | F1 | C5 |
| 16 | FCz | P1 | C6 | C5 | Cz |
| 17 | P2 | EOG | P1 | F5 | CP2 |
| 18 | P1 | FC3 | CP1 | C6 | O1 |
| 19 | C6 | Cz | O1 | POz | O9 |
| 20 | AFz | C6 | POz | AFz | TP8 |
| 21 | F5 | TP8 | TP7 | FT8 | CP1 |
| 22 | C5 | P2 | Fp2 | Fp2 | P1 |
| 23 | FT9 | POz | P5 | P2 | F1 |
| 24 | FC1 | F2 | P6 | P1 | F2 |
| 25 | FT7 | FC2 | FC4 | O10 | FT7 |
| 26 | POz | O10 | EOG | O9 | TP7 |
| 27 | O2 | O1 | FCz | P6 | P2 |
| 28 | P6 | CP3 | AFz | O1 | O2 |
| 29 | EOG | FCz | Cz | P5 | FT8 |
| 30 | P5 | P5 | FT10 | EOG | FT10 |
| 31 | FT10 | TP7 | F5 | Cz | F5 |
| 32 | Fp2 | O9 | TP8 | CPz | EOG |
| 33 | FT8 | CPz | P2 | F6 | P5 |
| 34 | O1 | O2 | O9 | O2 | CP3 |
| 35 | TP8 | F5 | O2 | TP7 | FC2 |
| 36 | O9 | FT7 | O10 | CP3 | FT9 |
| 37 | O10 | F1 | F6 | CP4 | Fp2 |
| 38 | F6 | FT8 | FT8 | FT7 | AFz |
| 39 | TP7 | C5 | FT9 | TP8 | F6 |

The ranking of the 39 EEG channels was calculated by the RFE method. The 17 channels over or close to motor areas of the cortex are marked with grey background for all five subjects. Underlined rank positions mark the estimated minimum of the RFE error curve for every subject from which on the error rate increases prominently (see figure 4 for the individual error curves).

channels can be used without increasing the error estimate.

The positions of the 17 channels over or close to the motor cortex were marked with a grey background. Except for very few of them, these channels have a high rank. For four of the subjects only few other (non-motor) channels were ranked above the marked minimum-error positions (see underlined ranks). For subject *B* channels FT9, FT10, and FP2 are relevant according to the ranking. To verify this observation we

- estimated the classification error using the seventeen motor channels and compared it to the error using the the motor channels plus FT9, FT10, FP2, and EOG. Indeed by adding artefact channels the error could be reduced from 24% to 21%.
- trained an SVM based on these artefact channels only. The performance was poor: only 0.55% accuracy could be reached in a 10-fold CV SVM training⁹.

That means that although feedback was not provided this subject showed task relevant muscle activity. However his performance was only supported by this muscle activity. The other four subjects did not accompany the left/right tasks with corresponding muscle movements¹⁰.

We conclude that the RFE method was capable of estimating physiologically meaningful EEG channels for the imagined left/right hand paradigm.

4.2 Visualization

The visualization of channel scores can support the analysis of BCI experiments, reveal activation patterns or channels carrying misleading artifacts and ease the choice of channel subgroups.

For visualization purposes we assigned a score calculated by RFE to each channel. The channels below the underlined entries of table 1 receive a score of 0. The ones above the underlined entries are mapped to the grey value scale according to their rank. Figures 6 and 7 show the task relevant channels for the five subjects. Black regions in both plots mark channels irrelevant for the classification task whereas white regions mark relevant ones.

For all subjects the informative regions are located close to the motor cortex. Subject *D* shows a clear and symmetrical concentration of important channels. The second column of figure 6 also shows, that subject *B* has additional important channels outside the

⁹The ridge was explicitly optimized for this test.

¹⁰This observation was supported by visual inspection and frequency analysis of the raw EEG signal - only very little muscle activity or other forms of artifacts could be detected.

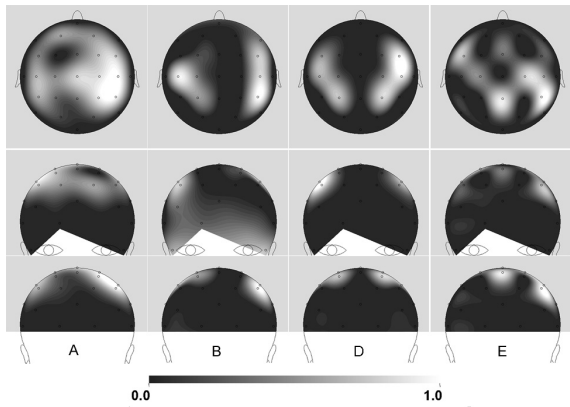


Figure 6: Visualization of task relevant regions for subjects *A, B, D* and *E* (one subject per column) during imagined hand movements. The score for each channel was obtained by using Recursive Feature Elimination (RFE) method and is based on the full duration of 5s. The top row depicts the view from above, the second and third row show the frontal view and view from the back. Please see also the left column of figure 7 for the corresponding mapping of subject *C*.

motor area probably resulting from muscle activity (as discussed above).

As the generalization error was minimal for the data of subject *C* we performed a closer examination of this data. Columns 2 to 4 of figure 7 visualize the spatial distribution of task specific information *over time*. We split the training data into three overlapping windows each of 2.5 seconds length. For every time window, we applied channel selection via RFE separately. It can be observed that the three resulting score patterns vary from window to window. This could be due to an unstable channel selection. Another reason might be that the task related activation pattern changes over time. Both issues will be addressed in future experiments.

5 Conclusion

We adapted two state of the art feature selection algorithms Recursive Feature Elimination (RFE) and Zero-Norm Optimization (l₀-Opt) as well as the Fisher Criterion for the special case of EEG channel selection for BCI applications.

The methods were applied to the paradigm of motor imagery. We showed that both RFE and l₀-Opt are capable of significantly reducing the number of channels needed for a robust classification without an increase of error. In our experiments, the Fisher Criterion failed to discover satisfying channel rankings.

The reason for the decrease in performance of the l₀-Opt compared to the RFE for smaller numbers of channels might be that on average the recursive l₀-Opt algorithm could not decrease the number of chosen channels to less than 25 before the recursion con-

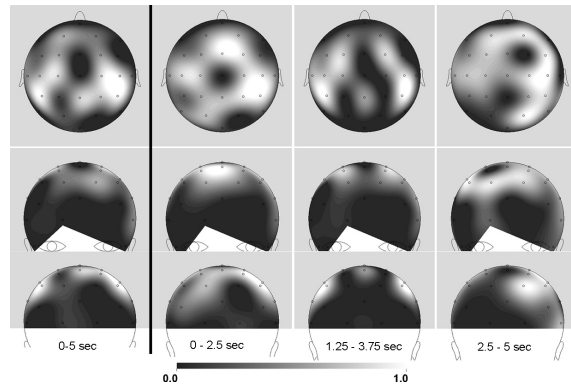


Figure 7: Visualization of task relevant regions for subject *C* (top, front and back view). The leftmost column shows the scores obtained by RFE based on the complete duration of 5s. The remaining three columns show the development of the scores over time. The rankings were obtained by applying the RFE method separately on the three shorter, overlapping time windows.

verged. This means that all the remaining channels were ranked according to the solution of only one SVM. To overcome this shortcoming of l₀-Opt we suggest the following procedure: channels are reduced with l₀-Opt until the minimum l₀-norm for w is obtained. In a next step the remaining channels are ranked using an iterative method like RFE instead of relying on a single SVM solution. This combination method was not investigated in this paper but will be subject to future research.

Although we did not incorporate explicit prior knowledge of the mental task or its underlying neural substrates, channels that are well known to be important (from a physiological point of view) were consistently selected by RFE whereas task irrelevant channels were disregarded. Furthermore the method revealed the use of muscular activity for one subject.

We introduced a method to visualize the channel rankings. This method can also be used to visualize the spatial change of task relevant information over time.

The results suggest that the RFE method can be used for new experimental paradigms in future BCI research - especially if no a priori knowledge about the location of important channels is available.

Acknowledgment

The authors would like to thank Rebecca Rörig for her restless data processing as well as Bernd Battes and Prof. Dr. Kuno Kirschfeld for their help with the EEG recordings. Special thanks to Gökhan Bakır for fruitful discussion. Parts of this work have been supported by DFG und NIH.

References

- [1] I. Guyon and A. Elisseeff. Introduction to variable and feature selection. *Journal of Machine Learning, Special Issue on Variable and Feature Selection*, pages 1157–1182, 2003.
- [2] J. Weston, A. Elisseeff, B. Schölkopf, and M. Tipping. Use of the zero-norm with linear models and kernel methods. *Journal of Machine Learning Research*, 3:1439–1461, March 2003.
- [3] V. N. Vapnik. *Statistical Learning Theory*. John Wiley and Sons, New York, USA, 1998.
- [4] J. Weston, S. Mukherjee, O. Chapelle, M. Pontil, T. Poggio, and V. Vapnik. Feature selection for SVMs. In S.A. Solla, T.K. Leen, and K.-R. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12, pages 526–532, Cambridge, MA, USA, 2000. MIT Press.
- [5] G. Pfurtscheller., C. Neuper and A. Schlögl, and K. Lugger. Separability of EEG signals recorded during right and left motor imagery using adaptive autoregressive parameters. *IEEE Transactions on Rehabilitation Engineering*, 6(3):316–325, 1998.
- [6] G. Pfurtscheller and F.H. Lopes da Silva. Event-related EEG/MEG synchronization and desynchronization: basic principles. *Clinical Neurophysiology*, 110(11):1842–1857, November 1999.
- [7] P.D. Welch. The use of fast fourier transform for the estimation of power spectra: A method based on time averaging over short, modified periodograms. In *IEEE Trans. Audio Electroacoustics*, volume AU-15, pages 70–73, June 1967.
- [8] S. Haykin. *Adaptive Filter Theory*. Prentice-Hall International, Inc., Upper Saddle River, NJ, USA, 1996.
- [9] P. Pudil, F. Ferri, J. Novovicova, and J. Kittler. Floating search methods for feature selection with non-monotonic criterion functions. In *Proceedings of the IEEE International Conference on Pattern Recognition*, pages 279–283, 1994.
- [10] M. Schröder, M. Bogdan, W. Rosenstiel, T. Hinterberger, and N. Birbaumer. Automated eeg feature selection for brain computer interfaces. In *Proceedings of the 1st International IEEE EMBS Conference on Neural Engineering*, pages 626–629, March 2003.
- [11] B. Blankertz, G. Curio, and K. Müller. Classifying single trial EEG: Towards brain computer interfacing. In T.K. Leen, T.G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems*, volume 14, Cambridge, MA, USA, 2001. MIT Press.
- [12] B. Schölkopf and A. Smola. *Learning with Kernels*. MIT Press, Cambridge, USA, 2002.
- [13] C. Bishop. *Neural Networks for Pattern Recognition*. Oxford UP, Oxford, UK, 1995.
- [14] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik. Gene selection for cancer classification using support vector machines. *Journal of Machine Learning Research*, 3:1439–1461, March 2003.