RESEARCH ARTICLE

# Model based generalization analysis of common spatial pattern in brain computer interfaces

**Gan Huang · Guangquan Liu · Jianjun Meng · Dingguo Zhang · Xiangyang Zhu**

**Abstract** In the motor imagery based Brain Computer Interface (BCI) research, Common Spatial Pattern (CSP) algorithm is used widely as a spatial filter on multi-channel electroencephalogram (EEG) recordings. Recently the overfitting effect of CSP has been gradually noticed, but what influence the overfitting is still unclear. In this work, the generalization of CSP is investigated by a simple linear mixing model. Several factors in this model are discussed, and the simulation results indicate that channel numbers and the correlation between signals influence the generalization of CSP significantly. A larger number of training trials and a longer time length of the trial would prevent overfitting. The experiments on real data also verify our conclusion.

**Keywords** Brain Computer Interfaces · Common Spatial Pattern · Generalization

## Introduction

Brain Computer Interface (BCI) brings human being direct ways to communication with computer programs or technical devices by intend alone [1]. This communicated way may help the patients with severe motor impairments like late-stage amyotrophic lateral sclerosis (ALS), severe cerebral palsy, head trauma, and spinal injuries in the life. It is well reported that both actual movements and motor imagery would lead to a short-lasting and circumscribed attenuation known as event related desynchronization (ERD) of rhythmic EEG components in the alpha and beta

G. Huang (✉) · G. Liu · J. Meng · D. Zhang · X. Zhu
State Key Laboratory of Mechanical System and Vibration,
Shanghai Jiao Tong University, 200240 Shanghai, China
e-mail: huanggan1982@gmail.com

band [2]. In contrast, foot movement results in an enhancement or event-related synchronization (ERS) of mu rhythms over the hand area [3].

Due to the poor spatial resolution of the raw EEG scalp potentials [4], a better spatial filter could help the BCI system improve its accuracy and information transform rate (ITR) effectively. Common spatial pattern (CSP) algorithm is a good alternative as a kind of spatial filter. By solving the generalized eigenvalue problem, CSP filter removes the signal's strong correlation among the original axes, and the distributions are maximally dissimilar along the new axes. In Ref. [5], a 2-D toy example was used to make a good understanding on CSP. It was reported that CSP algorithm was used more popular than others such as standard ear-reference, common average reference (CAR), small Laplacian (3 cm to set of surrounding electrodes), large Laplacian (6 cm to set of surrounding electrodes), Principal Commonents Analysis (PCA) and Independent Components Analysis (ICA) [6].

The CSP algorithm was firstly proposed by Fukunaga [7]. And Koles [8] used it to extract the abnormal components from EEG. Later, CSP was applied by Ramoser et al. [9] to create features for classification of ERD in motor imagery EEG. Based on the CSP algorithm, several variants have gradually appeared. Common Spectral Spatial Pattern (CSSP) was proposed as a improvement over CSP firstly [10], which simultaneously embedded a first-order finite impulse response (FIR) temporal filter into the CSP procedure. Considering the limitation in flexibility of first-order FIR filters, Dornhege et al. [11] presented the Common Sparse Spectral Spatial Pattern (CSSSP) algorithm. To obtain the optimizing spectral filters, some iterative algorithms were also introduced into CSP such as the SPECtrally Weighted CSP (SPEC-CSP) [12, 13] and Iterative Spatio-Spectral Patterns Learning (ISSPL) [14].

With the development of CSP algorithm, the issue of generalization has been gradually noticed [15, 16]. The generalization of CSP refers to how far the common spatial filter performs on the testing data from the training data. In Ref. [15], Hill et al. compared the performance of ICA and CSP algorithm in EEG, ECoG (electrocorticography) and MEG (magnetoencephalography) signals, and found that spatial filtering does not help in MEG, helps a little in ECoG, and improves performance a great deal in EEG. It was found that the supervised CSP algorithm suffers from poor generalization performance due to overfitting in ECoG and MEG.

Reuderink et al. [16] discussed the generalization of CSP over the amounts of training data, time and subjects. However, as Reuderink et al. pointed out, it is still unclear what influences the overfitting observed with CSP algorithm.

The complexity of EEG signals has restricted a further exploring for the overfitting of CSP. The complexity mainly comes from three aspects, the non-stationarity of the signals, the unclear internal mechanisms, and the large variance of the classification results. First, the non-stationarity of the EEG signals refers to the non-stationary in the single trial and between trials. Second, in BCI application, the EEG signals are often seen as the output of a black-box. Due to the complexity of the brain, it is still not very clear how the signals are produced from a well-defined paradigm (e.g. motor imagery). Third, the classification results are of high variation. Since the discovery of ERD, it has been found that the classification accuracy over subjects is quite different [17]. But what influences the results is still unknown. To get a thorough understanding for some methods, a large number of data sets with hundreds of trials from tens of subjects is needed, which increases the cost of investigation.

In this paper, a simple linear model is employed to investigate the generalization performance of CSP algorithm. By the method of modeling, all these difficulties mentioned above can be avoided. The signals produced by this model are locally stationary. It is not a black box any more, and we could adjust the parameters and compare the results as we want. For more convenience, it doesn't take us a lot of time and energy on data collection. Although there still exists a gap between real data and our model, they share some similar qualitative characteristics. The simple model can bring us a further understanding on the generalization of CSP in a controllable way.

In the next section, the CSP algorithm and the linear mixing model is introduced. In Sect. 3, the overfitting of CSP is analyzed based on the simulation of artificial data. The analysis of real data is presented in Sect. 4 and the final section is the discussion and conclusion.

## Methods

### CSP Algorithm

CSP algorithm is described as follows: the raw EEG signal of a single trial $k$ is represented as $X_k$ with its dimensions $ch \times len$, where $ch$ is the number of channels (i.e. recording electrodes) and $len$ is the time length of the trials. The normalized spatial covariance of the EEG can be obtained from

$$C_k = \frac{X_k X_k^T}{\text{trace}\left(X_k X_k^T\right)}$$

where $x^T$ denotes the transpose of the matrix $x$ and $\text{trace}(x)$ is the sum of the diagonal elements of the matrix $x$. Let

$$C_l = \sum_{k \in I_l} C_k, \qquad C_r = \sum_{k \in I_r} C_k;$$

which mean the spatial covariance of the EEG signals of each group. $I_l$, $I_r$ are the two index sets of the separated classes (i.e. the left and right hands motor imagery). Denote the composite spatial covariance $C = C_l + C_r$, and $C$ is decomposed as

$$C = U_C \Sigma U_C^T,$$

where $\Sigma$ is the diagonal matrix of the eigenvalues, and $U_C$ is the matrix of the corresponding eigenvectors.

Using whitening transformation

$$P = \sqrt{\Sigma^{-1}} U_C^T,$$

the spatial covariances $C_l$, $C_r$ can be transformed as

$$S_l = P C_l P^T = U \Sigma_l U^T,$$
$$S_r = P C_r P^T = U \Sigma_r U^T;$$

where $S_l$ and $S_r$ share common eigenvectors $U$, and

$$\Sigma_l + \Sigma_r = I,$$

$I$ is the identity matrix. Here, the whitening transform are used for simultaneous diagonalization.

Denote the projection matrix $W = U^T P$, which is named as spatial filters, and the columns of $A = W^{-1} \in R^{ch \times ch}$ are common spatial patterns. To the $k$-th trial, the filtered signal

$$Z_k = W X_k$$

is uncorrelated. By whitening transformation, the two distributions are maximally dissimilar along the new axes. In this work, the variance of three first and last rows of $Z_k$ are used as feature vector for classifier Linear Discriminant Analysis (LDA).

## Linear mixing model

At the interest frequencies (<100Hz), tissue is primarily treated as a resistive medium governed by Ohm's law and the capacitive effects can be neglected. Hence the multiple channels of EEG signals were usually considered as a linear model [18],

$$x(t) = As(t).$$

This model was also used in Ref. [5] to introduce the CSP method, where $s(t)$ were assumed to be two different distributions for two classes respectively.

Here, the following linear mixing model with local stationary sources are considerd:

$$\begin{cases} x_l = [A \ \tau I] \begin{bmatrix} s_l \\ s_n \end{bmatrix} = As_l + \tau s_n \\ x_r = [A \ \tau I] \begin{bmatrix} s_r \\ s_n \end{bmatrix} = As_r + \tau s_n \end{cases} \tag{1}$$

where $s_l$, $s_r$ are the signal sources corresponding to two class of imagery movements. $s_n$ represents the noise, and $\tau$ is the intensity of the noise. $s_l$, $s_r$, $s_n$ are assumed to be uncorrelated Gaussian distributions with their distributions

$$s_l \sim \quad \mathcal{N}(0, \Lambda^{(l)})$$
$$s_r \sim \quad \mathcal{N}(0, \Lambda^{(r)})$$
$$s_n \sim \quad \mathcal{N}(0, I)$$

and

$$\Lambda^{(l)} = \begin{bmatrix} \lambda & 0 \\ 0 & 1-\lambda \end{bmatrix}, \qquad \Lambda^{(r)} = \begin{bmatrix} 1-\lambda & 0 \\ 0 & \lambda \end{bmatrix}.$$

where $\lambda$ varies from 0 to 0.5, and $I$ is the identity matrix, which means the noises among channels are independent to each other. Without loss of generality, the signal sources are constructed in two-dimensional space for simplicity. The channel $A$ is defined as follows,

$$\begin{cases} A_{i1} = \mu, \quad A_{i2} = \sqrt{1-\mu^2}; \quad i = 1,3,5, \\ A_{i1} = \sqrt{1-\mu^2}, \quad A_{i2} = \mu; \quad i = 2,4,6, \end{cases}$$

and $A_{i1}^2 + A_{i2}^2 = 1$. $\mu$ varies between 0 and $\sqrt{2}/2$

In the mixing matrix $A$, $\mu$ represents the correlation between signal sources. $\lambda$ is used to define the separability of the signals. The time length of the signals is measured by *len*, which indicates the number of the sampling points. Take the channel number equal to 2 for instance. The spatial covariance matrix is represented as

$$x_r x_r^T = As_r s_r^T A^T + \tau^2 s_n s_n^T$$
$$= \left( A \begin{bmatrix} 1-\lambda & 0 \\ 0 & \lambda \end{bmatrix} A^T + \tau^2 I \right) \times len$$
$$= \begin{bmatrix} \lambda + (1-2\lambda)\mu^2 + \tau^2 & \mu\sqrt{(1-\mu^2)} \\ \mu\sqrt{(1-\mu^2)} & (1-\lambda) + (1-2\lambda)\mu^2 + \tau^2 \end{bmatrix} \times len$$

where the non-diagonal elements are the function of $\mu$. If $\mu = 0$, then $\mu\sqrt{(1-\mu^2)} = 0$, and the signals between the two channels are uncorrelated. With increasing of the value of $\mu$, the signals become more correlative. Eventually the correlation reaches the maximum as $\mu = \sqrt{2}/2$. If $\mu = 0$ and $\tau = 0$, the matrix $x_r x_r^T$ are proportional to $\Lambda^{(r)}$, similarly $x_l x_l^T$ are proportional to $\Lambda^{(l)}$. If $\lambda = 0$, the covariance matrices are maximally dissimilar from each other. As the value of $\lambda$ is increased to 0.5, the covariance matrices are just the same and can not be separated any more.

*Remark 1* In fact, there are multiple sources for real EEG signals. Here, for simplicity, the event related signals in this model are limited to two-dimensional, so that only two parameters $\lambda$ and $\mu$ are used to generate the signals $x_l$ and $x_r$.

*Remark 2* In the CSP algorithm, only the second order statistics are used. The algorithm doesn't require signals to be Gaussian or not. In the simulation, the signals are all assumed to be Gaussian distributed.

## Results

In the linear mixing model, the performance of CSP may be affected by the following factors listed in Table 1.
   We expect,

1. as more spatial information is given, a larger number of channels would lead to a better classification result;
2. with the increasing number of training trials, the performance should be improved;
3. the longer the time length of trials is, the higher the accuracy will be achieved;
4. as the noise intensity $\tau$ is enhanced, the correction rate would be lower;
5. a small value of $\lambda$ would make the signals from two classes more separable, so a higher correction rate would be guaranteed with a smaller $\lambda$;
6. the CSP algorithm has a perfect ability to make the signals de-correlated, but it is hard to guess how the correlation $\mu$ could influence the performance of CSP.

**Table 1** Factors in the linear mixing model

| | |
|---|---|
| $ch$ | The number of channels |
| $n$ | The number of training trials |
| $len$ | The time length of trials |
| $\tau$ | The noise intensity |
| $\lambda$ | The separability of the signals |
| $\mu$ | The correlation between signal sources |

For overfitting, it is reported that small training set has poor generalization [16]. The overfitting effect is worse when there are a larger number of channels related to the number of trials [15]. We have no idea about how the other factors can influence the generalization.

In the following, six groups of simulations are made to study these factors respectively. In each simulation, one factor varies with the other factors fixed. The signals are generated by the linear mixing model (1) with the sampling rate of 100 Hz. And they are filtered by IIR butterworth bandpass filter with the band of 8–30 Hz, which is the frequency of mu-rhythm. All results are averaged by 400 runs in order to make the results smoother.

Number of channels

Training accuracies and testing accuracies are compared in Fig. 1 with different channel numbers. As the channel number rises from 10 to 100, the training accuracy continues to go up from 92 to 100%. In testing, the accuracy grows rapidly from 10 to 30 channels, and it gradually declines.

In fact this tendency is also observed in other experiments [19]. Lv and Liu attributed it to the artifacts imported by part of channels. Due to the different locations of the electrodes, some researchers also deduced it to that some channels are more correlated with movement imagery than others.

In this linear mixing model, we use the uniform $\mu$ for all channels, and the noise intensity equals to each other. Hence the channels have the same qualities in the simulation. It is observed that the increasing channel number would degrade the generalization performance. As the
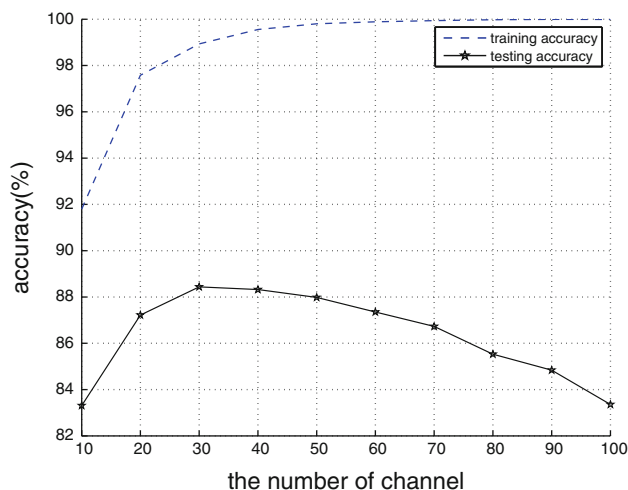
channel number increases, CSP algorithm will have more spatial information to do spatial filtering. Meanwhile, a large number of channels also bring CSP algorithm into a high-dimensional feature space, which leads to a poor generalization. Hence the peak in the curve could be seen as the best trade off between enough spatial information and overfitting.

In the following simulation, the generalization of CSP is measured by the accuracies varying with channel numbers. If the accuracy grows with the channel number, then CSP under the corresponding parameter set shows good generalization. If the peak occurs with a small channel number, then CSP has a poor generalization in this condition.

Number of training trials

Figure 2 shows the performance of CSP with different number of training trials. The same number of trials is used in test. As we expected, small training set gets a low accuracy, and tends to overfitting.

Time length of the trials

As mentioned in Ref. [5], shorter segments result in more responsive but more noisy feedback signal. Longer segments give a smoother control signal, but the delay from intention to control becomes longer. In the simulation (Fig. 3), the data segment varies from 100 ms to 800 ms. It could be found that longer segments get preferable accuracies above 90%, and the performance declines very slightly or tends to be stable when the channel numbers grows.
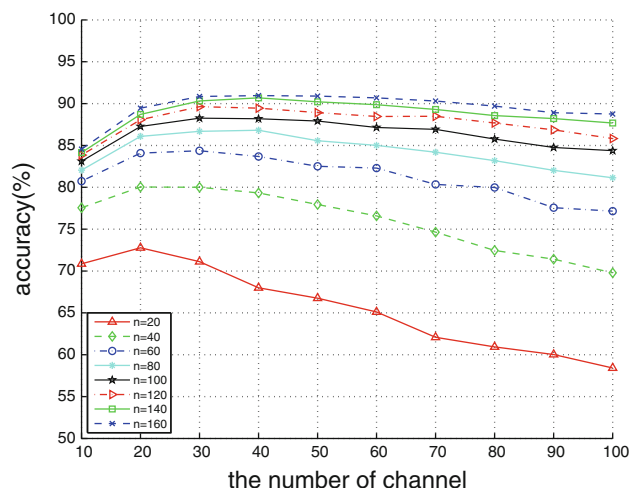


**Fig. 1** Mean training errors and testing errors depending on the number of channels, where $n = 100$, $len = 400$ ms, $\tau = 1$, $\lambda = 0.46$ and $\mu = \cos(2\pi/36)$



**Fig. 2** Mean accuracies depending on the number of training trials, where $len = 400$ ms, $\tau = 1$, $\lambda = 0.46$ and $\mu = \cos(2\pi/36)$
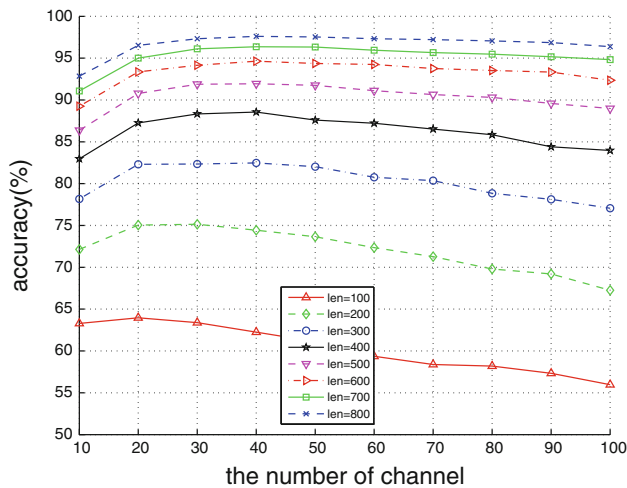
**Fig. 3** Mean accuracies depending on the time length of the trials, where $n = 100$, $\tau = 1$, $\lambda = 0.46$ and $\mu = \cos(2\pi/36)$
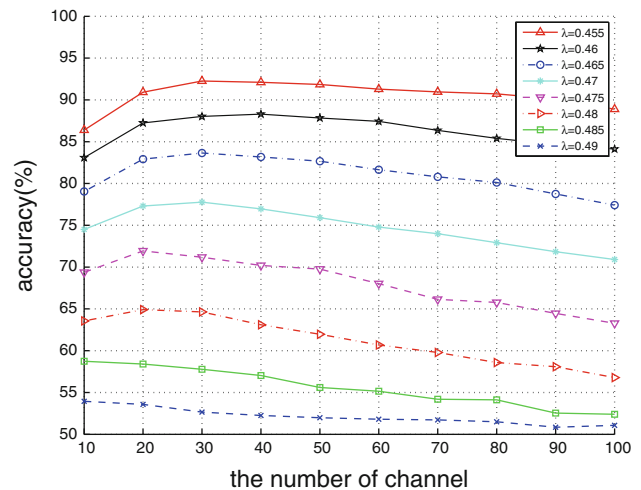


**Fig. 5** Mean accuracies depending on the separability of the signals, where $n = 100$, $len = 400$ ms, $\tau = 1$ and $\mu = \cos(2\pi/36)$
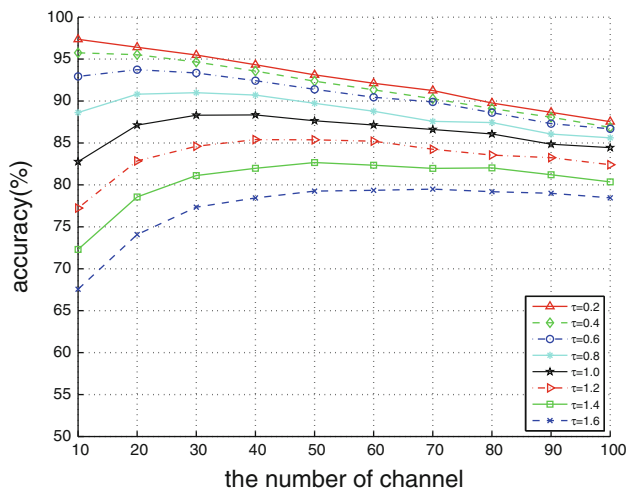


**Fig. 4** Mean accuracies depending on the intensity of the noise, where $n = 100$, $len = 400$ ms, $\lambda = 0.46$ and $\mu = \cos(2\pi/36)$
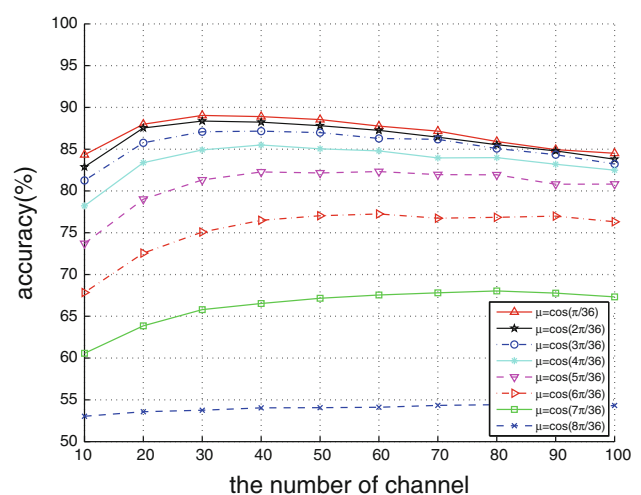


**Fig. 6** Mean accuracies depending on the correlation between channels, where $n = 100$, $len = 400$ ms, $\tau = 1$ and $\lambda = 0.46$

Noise intensity

In the simulation, as illustrated in Fig. 4, the noise intensity $\tau$ varies from 0.2 to 1.6. It is observed that the signals under less noisy circumstance have high accuracies but poor generalization. On the other hand, the signals with large $\tau$ would need more channels for CSP to get better performance.

Separability of the signals

As $\lambda$ decreasing from 0.49 to 0.455 in Fig. 5, the signals are more separable. The correction rates are improved remarkably, while the generalization is enhanced slightly.

Correlation between signal sources

For correlation, with the increasing of $\mu$ from $\cos(8\pi/36)$ to $\cos(\pi/36)$ in Fig. 6, the accuracies are improved significantly, while the generalization is weaken.

In summary, the generalization of CSP tends to degrade as the channel numbers increases. Large number of training trials and long segments would help CSP algorithm improving the accuracy greatly and preventing overfitting slightly. Furthermore, the noise intensity $\tau$ and the correlation $\mu$ influence the generalization of CSP severely. In general, both $\tau$ and $\mu$ can be seen as the correlation. $\tau$ is the correlation between the signal sources and noises, and $\mu$ represents the correlation between two different signal

sources. Under less correlation circumstances, CSP would get better performance with fewer channels. Finally, if the signals are more separable in nature, the accuracies would be higher.
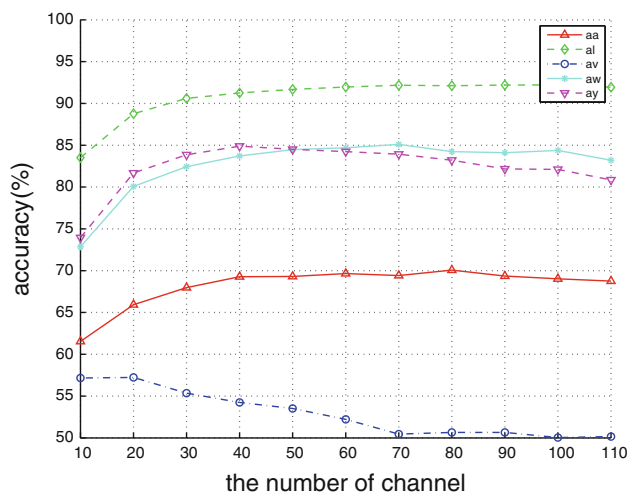


**Fig. 7** Mean training errors and testing errors depending on the number of channels. All results are averaged by 400 runs, and in every run, the channels are selected randomly in the 118 channels

## Real data analysis

Data set IVa of the BCI Competition III [6] is used for our experiment. This data set contains 280 trials for each of the 5 subjects. Cued motor imagery with 2 classes (right hand, foot) were recorded. The recording was made using BrainAmp amplifiers and a 128 channel Ag/AgCl electrode cap from BCI. 118 EEG channels were measured at positions of the extended international 10/20-system. Signals were band-pass filtered between 0.05 and 200 Hz and then digitized at 1,000 Hz with 16 bit (0.1 uV) accuracy. This data set was selected because it contained enough trials with large number of channels. In our experiment, the signals are down-sampled to 100 Hz. The bandpass filter is made by butterworth IIR-filter from 8 to 30 Hz. The trials are separated into 3.5 s.

The performance of five subjects is illustrated in Fig. 7. As the channel numbers is increased, the accuracy rises for subjects *aa*, *al*, *aw*, *ay*, but declines for subject *av*. Furthermore, the covariance matrices $C$ for the five subjects are plotted in Fig. 8. From the covariance matrices, we find the signals near the area of POZ (channel 106) are strong for all subjects except *av*, and the correlation for this signal to the other channels decreases by distance. In the area of left and right primary motor cortex (localized over C3 for channel 52 and C4 for channel 56) the correlation is still
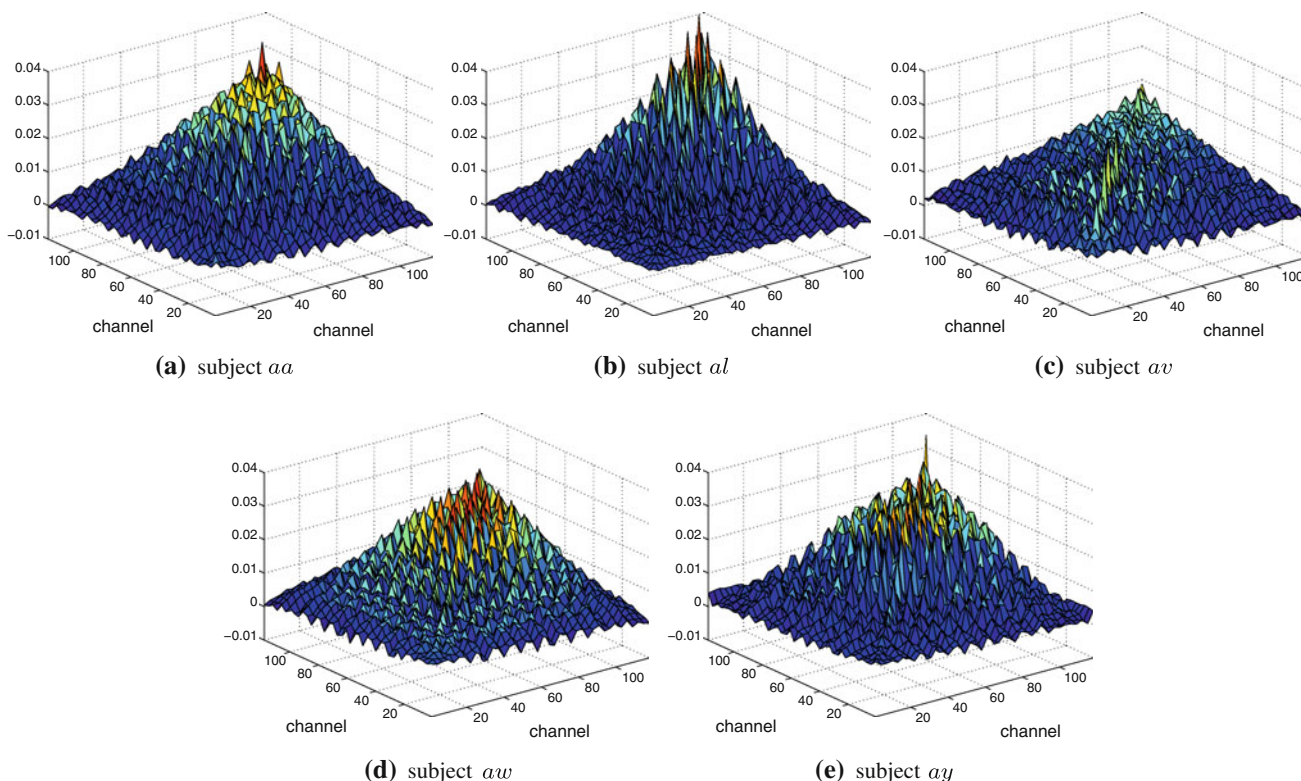


**(a)** subject *aa*



**(b)** subject *al*



**(c)** subject *av*



**(d)** subject *aw*



**(e)** subject *ay*

**Fig. 8** The spatial covariance of the five subjects. **a** subject *aa*. **b** subject *al*. **c** subject *av*. **d** subject *aw*. **e** subject *ay*

large, which could be seen as a strong noise. Under the noisy circumstance, the simulation results in Sect. 3 tell us that larger channel numbers can help improving the accuracies, which may interpret the results in Fig. 7.

**Remark 3** In this experiment, the channels are selected randomly, and all results are averaged by 400 runs. So the quality of different channels will not affect the analysis. Here, we focus on the generalization of CSP and channel selection method is not needed.

## Discussion and conclusion

By employing the linear mixing model, we explored the generalization performance of CSP on different channel numbers, number of training trials and data length. Some intrinsic factors like the noise intensity, the separability of the signals and the correlation between signal sources are also investigated. From the simulation, we find that larger number of training trials and longer time length of trials would improve the accuracies and prevent overfitting. Furthermore, channel numbers and the correlation between signals affect the overfitting of CSP significantly. The generalization performance would degrade when more channels are used.

Among the six factors studied in this work, $\tau$, $\lambda$ and $\mu$ are determined by the intrinsic factors of the signals and the technique for signal acquisition. Some methods like regularization can be used to suppress the effect of noise. The factors $n$ and $len$ influence the training time and the bit rate of the BCI applications directly. Hence, channel selection is the prevailing way to avoid the overfitting problem and improve the classification results. Some channel selection methods have been proposed and compared in Meng et al. [20]. However, a great number of channels unselected mean that a part of spatial information is not used by CSP method. A paper extending the current model to make full use of the spatial information and deal with overfitting of CSP algorithm is in preparation.

## References

1. Freeman W (2007) Definitions of state variables and state space for brain-computer interface. Cogn Neurodyn 1:3–14

2. Pfurtscheller G, Aranibar A (1979) Evaluation of event-related desynchronization (ERD) preceding and following voluntary self-paced movement. Electroencephalogr Clin Neurophysiol 46:138–46

3. Pfurtscheller G, Neuper C (1994) Event-related synchronization of mu rhythm in the EEG over the cortical hand area in man. Neurosci Lett 174:93–96

4. Nunez P, Srinivasan R, Westdorp A, Wijesinghe R, Tucker D, Silberstein R, Cadusch P (1997) EEG coherency I: statistics, reference electrode, volume conduction, Laplacians, cortical imaging, and interpretation at multiple scales. Electroencephalogr Clin Neurophysiol 103:499–515

5. Blankertz B, Tomioka R, Lemm S, Kawanabe M, Muller K (2008) Optimizing spatial filters for robust EEG single-trial analysis. Signal Process Mag IEEE 25:41–56

6. http://ida.first.fraunhofer.de/projects/bci/competitions

7. Fukunaga K (1990) Introduction to statistical pattern recognition. Academic Press, London

8. Koles Z (1991) The quantitative extraction and topographic mapping of the abnormal components in the clinical EEG. Electroencephalogr Clin Neurophysiol 79:440–447

9. Ramoser H, Muller-Gerking J, Pfurtscheller G (2000) Optimal spatial filtering of single trial EEG during imagined handmovement. IEEE Trans Rehabil Eng 8:441–446

10. Lemm S, Blankertz B, Curio G, Muller K (2005) Spatio-spectral filters for improving the classification of single trial EEG. IEEE Trans Biomed Eng 52:1541

11. Dornhege G, Blankertz B, Krauledat M, Losch F, Curio G, Muller K (2006) Combined optimization of spatial and temporal filters for improving brain-computer interfacing. IEEE Trans Biomed Eng 53:2274

12. Tomioka R, Dornhege G, Nolte G, Aihara K, Muller K (2006) Optimizing spectral filters for single trial EEG classification. Lect Notes Comput Sci 4174:414

13. Tomioka R, Dornhege G, Nolte G, Blankertz B, Aihara K, Müller K (2006) Spectrally weighted common spatial pattern algorithm for single trial eeg classification. Dept. Math. Eng., Univ. Tokyo, Tokyo, Japan, Tech. Rep, 40

14. Wu W, Gao X, Hong B, Gao S (2008) Classifying single-trial EEG during motor imagery by iterative spatio-spectral patterns learning (ISSPL). IEEE Trans Biomed Eng 55:1733–1743

15. Hill N, Lal T, Schroder M, Hinterberger T, Widman G, Elger C, Scholkopf B, Birbaumer N (2006) Classifying event-related desynchronization in EEG, ECoG and MEG signals. Lect Notes Comput Sci 4174:404

16. Reuderink B, Poel M (2008) Robustness of the common spatial patterns algorithm in the BCI-pipeline

17. Guger C, Edlinger G, Harkam W, Niedermayer I, Pfurtscheller G, OEG G, Graz A (2003) How many people are able to operate an EEG-based brain-computer interface (BCI)?. IEEE Trans Neural Syst Rehabil Eng 11:145–147

18. Parra L, Spence C, Gerson A, Sajda P (2005) Recipes for the linear analysis of EEG. NeuroImage 28:326–341

19. Lv J, Liu M (2008) Common spatial pattern and particle swarm optimization for channel selection in BCI. In: Innovative computing information and control, 2008. ICICIC'08. 3rd international conference on, pp 457–457

20. Meng J, Liu G, Huang G, Zhu X (2009) Automated selecting subset of channels based on CSP in motor imagery brain-computer interface system. In: 2009 IEEE International conference on robotics and biomimetics (ROBIO), pp 2290–2294