

# Bayesian Multi-Task Learning for Common Spatial Patterns

Hyohyeong Kang  
 Department of Computer Science  
 Pohang University of Science and Technology  
 Pohang, Korea  
 Email: paanguin@postech.ac.kr

Seungjin Choi  
 Department of Computer Science  
 Division of IT Convergence Engineering  
 Pohang University of Science and Technology  
 Pohang, Korea  
 Email: seungjin@postech.ac.kr

**Abstract**—Common spatial pattern (CSP) is a widely-used feature extraction method for electroencephalogram (EEG) classification and corresponding probabilistic models were recently developed, adopting a linear generative model for each class. These models are trained on a subject-by-subject basis so that inter-subject information is neglected. Moreover when only a few training samples are available for each subject, the performance is degraded. In this paper we employ Bayesian multi-task learning so that subject-to-subject information is transferred in learning the model for a subject of interest. We present two probabilistic models where precision parameters of multivariate or matrix-variate Gaussian prior for the dictionary are shared across subjects. Numerical experiments on the BCI competition IV 2a dataset confirm that our methods improve classification performance over the standard CSP (on a subject-by-subject basis), especially in the case of subjects with fewer number of training samples.

**Keywords**—Bayesian multi-task learning; brain computer interface; common spatial patterns;

## I. INTRODUCTION

Electroencephalogram (EEG) is multivariate brain wave signals recorded at multiple electrodes placed on scalp, reflecting electrical potentials induced by subjects' brain activities. EEG classification allows computers to translate a subject's intention or mind into a control signal for a device, which is important for brain computer interface (BCI) [1].

Common spatial pattern (CSP) is a popular discriminative EEG feature extraction methods, which is proven to be a useful subject-specific spatial filter [4]. Recently a probabilistic model was proposed for CSP, where two linear Gaussian generative models with the dictionary shared are jointly learned to infer common spatial patterns [6]. CSP is a subject-specific spatial filter, which does not consider other subjects' information involving the same task as the subject of interest. In the case of a subject with much fewer training samples, the performance of CSP is much deteriorated. Thus it is desirable to transfer useful information of subjects involving the same task to the subject of interest with fewer training samples, which is known as *subject-to-subject transfer* [3], [5]. In this paper we present two Bayesian multi-task extensions of probabilistic CSP, where precision parameters of multivariate or matrix-variate Gaussian prior for the dictionary are shared across subjects.

## II. METHODS

We denote by  $\mathbf{X}^{sc} = [\mathbf{x}_1^{sc}, \dots, \mathbf{x}_{N_{sc}}^{sc}] \in \mathbb{R}^{D \times N_{sc}}$  a collection of EEG signals measured from  $D$  electrodes over trials ( $N_{sc}$  is the number of samples) for subject  $s \in \{1, \dots, S\}$  who undergoes the mental task involving class  $c \in \{1, 2\}$ . The multi-task extension of the probabilistic CSP model assumes that  $\mathbf{X}^{sc}$  is generated by

$$\mathbf{X}^{sc} = \mathbf{W}^s \mathbf{Z}^{sc} + \mathbf{E}^{sc}, \quad (1)$$

where  $\mathbf{W}^s = [\mathbf{w}_1^s, \dots, \mathbf{w}_m^s] \in \mathbb{R}^{D \times M}$  is the *dictionary* for subject 's', containing  $M$  basis vectors, which is shared across two classes,  $\mathbf{Z}^{sc} = [\mathbf{z}_1^{sc}, \dots, \mathbf{z}_{N_c}^{sc}] \in \mathbb{R}^{M \times N_c}$  the *coefficient matrix*, and  $\mathbf{E}^{sc} = [\boldsymbol{\epsilon}_1^{sc}, \dots, \boldsymbol{\epsilon}_{N_c}^{sc}] \in \mathbb{R}^{D \times N_c}$  is the noise matrix. Each row of  $\mathbf{X}^{sc}$  is already centered (zero mean).

Coefficients and noise are assumed to be zero-mean Gaussians:

$$\begin{aligned} \mathbf{z}_t^{sc} &\sim \mathcal{N}(\mathbf{z}_t^{sc} | \mathbf{0}, (\boldsymbol{\Lambda}^{sc})^{-1}), \\ \boldsymbol{\epsilon}_t^{sc} &\sim \mathcal{N}(\boldsymbol{\epsilon}_t^{sc} | \mathbf{0}, (\boldsymbol{\Psi}^{sc})^{-1}), \end{aligned}$$

where  $\boldsymbol{\Lambda}^{sc} = \text{diag}(\lambda_1^{sc}, \dots, \lambda_M^{sc}) \in \mathbb{R}^{M \times M}$  and  $\boldsymbol{\Psi}^{sc} = \text{diag}(\psi_1^{sc}, \dots, \psi_D^{sc}) \in \mathbb{R}^{D \times D}$  are diagonal precision matrices for  $s = 1, \dots, S$  and  $c = 1, 2$ , which are assumed to follow Gamma distributions:

$$\begin{aligned} p(\boldsymbol{\Lambda}^{sc}) &= \prod_{m=1}^M \text{Gamma}(\lambda_m^{sc} | a_0^\lambda, b_0^\lambda), \\ p(\boldsymbol{\Psi}^{sc}) &= \prod_{d=1}^D \text{Gamma}(\psi_d^{sc} | a_0^\psi, b_0^\psi). \end{aligned}$$

In the case of  $S = 1$  (subject-specific model), the model (1) reduces to the probabilistic CSP model in [6] (see Fig. 1(a)). We describe two multi-task learning methods for probabilistic CSP in the following two subsections, where multivariate Gaussian prior (in MTL-CSP1) or matrix-variate Gaussian prior (in MTL-CSP2) for  $\{\mathbf{W}^s\}$  is defined, with some hyperparameters shared.

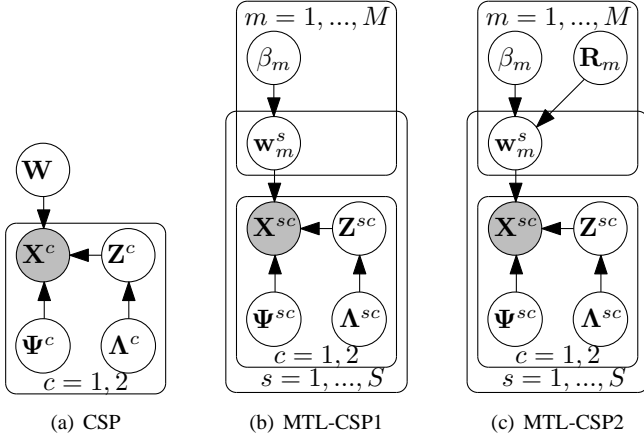


Figure 1. Graphical representations for three probabilistic CSP models: (a) probabilistic CSP on a subject-by-subject basis [6]; (b) multi-task learned CSP (proposed method 1); (c) multi-task learned CSP (proposed method 2).

#### A. MTL-CSP1: Sharing Multivariate Gaussian Prior

We consider the Gaussian prior distributions for the dictionaries  $\{\mathbf{W}^s\}$  which are common across subjects:

$$p(\mathbf{W}^s | \beta) = \prod_{m=1}^M \mathcal{N}(\mathbf{w}_m^s | \mathbf{0}, \beta_m^{-1} \mathbf{I}_D),$$

$$p(\beta_m) = \text{Gamma}(\beta_m | a_0^\beta, b_0^\beta).$$

Each basis vector  $\mathbf{w}_m^s$  follows the isotropic Gaussian prior with the hyperparameter  $\beta_m$  that is shared across subjects, leading to an indirect sharing of information across subjects (see Fig. 1(b)). The indirect information transfer by sharing hyperparameters in Bayesian framework was motivated by the empirical Bayesian multi-task learning method [2]. Hyperparameters are estimated by the empirical variational Bayesian method, which is explained in Section II-C.

#### B. MTL-CSP2: Sharing Matrix-Variate Gaussian Prior

The MTL-CSP1 model does not consider correlations between subjects so that the negative information transfer (from subjects with negative correlations) might degrade the generalization performance. To overcome this negative effect, we present the MTL-CSP2 (see Fig. 1(c)) model which also takes inter-subject correlations into account. To this end, we consider the matrix-variate Gaussian prior distribution. The matrix-variate Gaussian distribution for a random matrix  $\mathbf{Y} \in \mathbb{R}^{M \times N}$ , denoted by  $\mathcal{N}_{M,N}(\mathbf{Y} | \mathbf{M}, \mathbf{\Sigma}, \mathbf{\Omega})$ , is defined by the probability density function which takes the form:

$$\mathcal{N}_{M,N}(\mathbf{Y} | \mathbf{M}, \mathbf{\Sigma}, \mathbf{\Omega}) = (2\pi)^{-\frac{MN}{2}} |\mathbf{\Sigma}|^{-\frac{N}{2}} |\mathbf{\Omega}|^{-\frac{M}{2}}$$

$$\exp \left\{ -\frac{1}{2} \text{tr} \left( \mathbf{\Omega}^{-1} (\mathbf{Y} - \mathbf{M})^\top \mathbf{\Sigma}^{-1} (\mathbf{Y} - \mathbf{M}) \right) \right\},$$

where  $|\cdot|$  denotes the determinant and  $\text{tr}(\cdot)$  denotes the trace operator,  $\mathbf{M} \in \mathbb{R}^{M \times N}$  is the mean matrix,  $\mathbf{\Sigma} \in \mathbb{R}^{M \times M}$  and  $\mathbf{\Omega} \in \mathbb{R}^{N \times N}$  are covariance matrices.

In MTL-CSP2 model, we consider the matrix-variate Gaussian distribution for  $\overline{\mathbf{W}}_m = [\mathbf{w}_m^1, \dots, \mathbf{w}_m^S] \in \mathbb{R}^{D \times S}$ , in order to accommodate correlations between  $\mathbf{w}_m^s$  and  $\mathbf{w}_m^{s'}$  for  $s, s' = 1, \dots, S$ :

$$p(\overline{\mathbf{W}}_m | \mathbf{R}_m, \beta_m) = \mathcal{N}_{D,S}(\overline{\mathbf{W}}_m | \mathbf{0}, \beta_m^{-1} \mathbf{I}_D, \mathbf{R}_m^{-1}),$$

where each  $\mathbf{w}_m^s$  obeys the isotropic Gaussian as in the MTL-CSP1 model and correlations between  $\mathbf{w}_m^s$  and  $\mathbf{w}_m^{s'}$  are reflected in  $\mathbf{R}_m \in \mathbb{R}^{S \times S}$  which is assumed to follow the Wishart distribution,  $p(\mathbf{R}_m) = \mathcal{W}(\mathbf{R}_m | \mathbf{W}_0, \nu_0)$ . Thus the prior distribution for  $\{\mathbf{W}^s\}$  is given by

$$p(\{\mathbf{W}^s\} | \{\mathbf{R}_m\}, \{\beta_m\}) = \prod_{m=1}^M p(\overline{\mathbf{W}}_m | \mathbf{R}_m, \beta_m).$$

Note that when  $\mathbf{R}_m = \mathbf{I}_S$ , MTL-CSP2 reduces to MTL-CSP1. Learning  $\mathbf{R}_m$  from training set enables us to capture positive or negative correlations between subjects so that the effect of negative transfer is reduced.

#### C. Variational Inference

We apply the variational inference method to compute the approximate posterior distribution over a set of variables,  $\theta = (\{\mathbf{W}^s\}, \{\Psi^{sc}\}, \{\mathbf{Z}^{sc}\}, \{\Lambda^{sc}\}, \{\beta_m\})$  for MTL-CSP1 or  $\theta = (\{\mathbf{W}^s\}, \{\Psi^{sc}\}, \{\mathbf{Z}^{sc}\}, \{\Lambda^{sc}\}, \{\beta_m\}, \{\mathbf{R}_m\})$  for MTL-CSP2. In MTL-CSP1, we assume that the variational distribution  $q(\theta)$  factorizes over variables:

$$q(\theta) = q(\{\mathbf{W}^s\})q(\{\Psi^{sc}\})q(\{\mathbf{Z}^{sc}\})q(\{\Lambda^{sc}\})q(\{\beta_m\}),$$

where

$$q(\{\mathbf{W}^s\}) = \prod_{s=1}^S \prod_{d=1}^D \mathcal{N}(\overline{\mathbf{w}}_d^{s\top} | \nu_d^s, \Phi_d^s),$$

$$q(\{\Psi^{sc}\}) = \prod_{s=1}^S \prod_{c=1}^2 \prod_{d=1}^D \text{Gamma}(\psi_d^{sc} | a_d^{\psi^{sc}}, b_d^{\psi^{sc}}),$$

$$q(\{\mathbf{Z}^{sc}\}) = \prod_{s=1}^S \prod_{c=1}^2 \prod_{t=1}^{N_{sc}} \mathcal{N}(z_t^{sc} | \mu_t^{sc}, \Sigma^{sc}),$$

$$q(\{\Lambda^{sc}\}) = \prod_{s=1}^S \prod_{c=1}^2 \prod_{m=1}^M \text{Gamma}(\lambda_m^{sc} | a_m^{\lambda^{sc}}, b_m^{\lambda^{sc}}),$$

$$q(\{\beta_m\}) = \prod_{m=1}^M \text{Gamma}(\beta_m | a_m^\beta, b_m^\beta),$$

where  $\overline{\mathbf{w}}_d^s$  represents the  $d$ -th row of  $\mathbf{W}^s$ . Variational parameters are estimated by maximizing the lower bound on

the marginal likelihood,  $\mathcal{L}(q) = \mathbb{E}_q [\log (p(\mathbf{X}, \boldsymbol{\theta}) / q(\boldsymbol{\theta}))]$ ,

$$\begin{aligned}
(\Phi_d^s)^{-1} &= \mathbb{E}[D_\beta] + \sum_{c=1}^2 \mathbb{E}[\psi_d^{sc}] \sum_{t=1}^{N_{sc}} \mathbb{E}[\mathbf{z}_t^{sc} \mathbf{z}_t^{sc\top}], \\
\boldsymbol{\nu}_d^s &= \Phi_d^s \sum_{c=1}^2 \mathbb{E}[\psi_d^{sc}] \sum_{t=1}^{N_{sc}} \mathbf{x}_t^{sc}(d) \mathbb{E}[\mathbf{z}_t^{sc}], \\
\Sigma^{sc} &= \left( \mathbb{E}[\Lambda^{sc}] + \mathbb{E}[\mathbf{W}^{s\top} \Psi^{sc} \mathbf{W}^s] \right)^{-1}, \\
\boldsymbol{\mu}_t^{sc} &= \Sigma^{sc} \mathbb{E}[\mathbf{W}^{s\top}] \mathbb{E}[\Psi^{sc}] \mathbf{x}_t^{sc}, \\
a_d^{\psi sc} &= a_0^{\psi sc} + \frac{N_{sc}}{2}, \quad a_m^{\lambda sc} = a_0^{\lambda sc} + \frac{N_{sc}}{2}, \\
b_d^{\psi sc} &= b_0^{\psi sc} + \frac{1}{2} \sum_{t=1}^{N_{sc}} \mathbb{E}[(\mathbf{x}_t^{sc}(d) - \bar{\mathbf{w}}_d^s \mathbf{z}_t^{sc})^2], \\
b_m^{\lambda sc} &= b_0^{\lambda sc} + \frac{1}{2} \sum_{t=1}^{N_{sc}} \mathbb{E}[(\mathbf{z}_t^{sc}(m))^2], \\
a_m^\beta &= a_0^\beta + \frac{DS}{2}, \quad b_m^\beta = b_0^\beta + \frac{1}{2} \sum_{s=1}^S \mathbb{E}[\mathbf{w}_m^{s\top} \mathbf{w}_m^s],
\end{aligned}$$

where  $D_\beta \in \mathbb{R}^{M \times M}$  is the diagonal matrix with diagonal entries  $\beta_m$  for  $m = 1, \dots, M$ .

For MTL-CSP2, we also assume that the variational distribution  $q(\boldsymbol{\theta})$  factorizes over variables:

$$q(\boldsymbol{\theta}) = q(\{\mathbf{W}^s\})q(\{\Psi^{sc}\})q(\{\mathbf{Z}^{sc}\})q(\{\Lambda^{sc}\})q(\{\beta_m\})q(\{\mathbf{R}_m\}),$$

where  $q(\{\mathbf{R}_m\}) = \prod_{m=1}^M \mathcal{W}(\mathbf{R}_m | \mathbf{W}_m, \nu_m)$ , and variational distributions over other variables are the same as MTL-CSP1. Variational parameters are updated:

$$\begin{aligned}
(\Phi_d^s)^{-1} &= \mathbb{E}[D_\beta \mathbf{R}(s, s)] + \sum_{c=1}^2 \sum_{t=1}^{N_{sc}} \mathbb{E}[\psi_d^{sc}] \mathbb{E}[\mathbf{z}_t^{sc} \mathbf{z}_t^{sc\top}], \\
\boldsymbol{\nu}_d^s &= \Phi_d^s \left( \sum_{c=1}^2 \mathbb{E}[\psi_d^{sc}] \sum_{t=1}^{N_{sc}} \mathbf{x}_t^{sc}(d) \mathbb{E}[\mathbf{z}_t^{sc}] \right. \\
&\quad \left. - \sum_{j \neq s} \mathbb{E}[D_\beta \mathbf{R}(s, j)] \mathbb{E}[\bar{\mathbf{w}}_d^{j\top}] \right), \\
\nu_m &= \nu_0 + D, \\
(\mathbf{W}_m)^{-1} &= \mathbf{W}_0^{-1} + \mathbb{E}[\beta_m] \mathbb{E}[\bar{\mathbf{W}}_m^\top \bar{\mathbf{W}}_m],
\end{aligned}$$

where  $\mathbf{R}(i, j)$  is the  $M \times M$  diagonal matrix with diagonal entries  $[\mathbf{R}_m]_{ij}$ , ( $m = 1, \dots, M$ ). The expectations  $\mathbb{E}$  appearing in updating equations for variational parameters are taken with respect to the variational distribution  $q(\cdot)$ , the values of which are left out due to the space limit. Hyperparameters are also estimated by maximizing the lower-bound  $\mathcal{L}(q)$ .

### III. NUMERICAL EXPERIMENTS

We tested our models on BCI competition IV<sup>1</sup> 2a dataset. The dataset contains EEG measurements of 9 subjects with

<sup>1</sup><http://www.bbci.de/competition/iv/index.html>

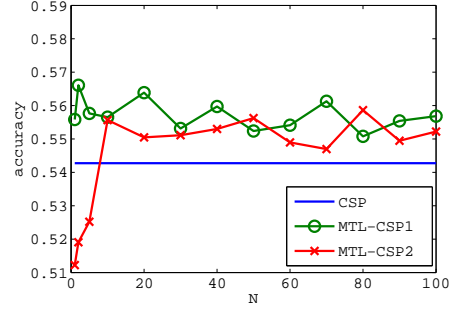


Figure 3. The classification performance averaged over every subject.

22 electrodes. Although there were 4 target classes in the data set, we only used the trials that correspond to the imagery left hand and right hand classes. For every trial, we picked the interval from 3.5 seconds to 5.5 seconds after each cue. Each subject performed 144 trials per class so that the dataset has 288 trials per subject.

We randomly selected only one trial from each class per subject as a training data. We extracted features and trained classifiers using the training data. The remaining trials are considered as the test data. We calculated log-variances projected by 6 CSP filters that corresponds to 3 largest and 3 smallest eigenvalues. Compared to original CSP, our MTL-CSP models can exploit every subject's data. When we extracting features for a subject  $s$  with MTL-CSP models, we randomly chose  $N$  trials of each class from the other subjects, and used them as an additional source of information. After the models were inferred, we calculated the log value of the variance of  $\mathbb{E}[\mathbf{z}_t]$  within each trial to obtain features. By the using  $\mathbb{E}[\lambda_m^{s1}] / \mathbb{E}[\lambda_m^{s2}]$ , we selected 3 largest-valued and 3 smallest-valued features. In every case, we trained the logistic regression classifier, and predicted the labels of test data using these features. We repeated the experiment for each value of  $N$  for 10 times, and averaged the results. The classification accuracy was estimated by the ratio of the number of accurately classified test trials compared to the total number of test trials (Fig. 2, 3).

### IV. CONCLUSIONS

We have presented Bayesian multi-task learning extensions of the probabilistic CSP method, allowing subject-to-subject transfer by sharing parameters of multivariate or matrix-variate Gaussian prior for the dictionary. We have presented two models, MTL-CSP1 and MTL-CSP2, developing variational inference algorithms to compute posterior distributions over variables. Numerical experiments on BCI competition IV 2a dataset demonstrated that our models improved EEG classification performance over the standard probabilistic CSP which is subject-specific. There were improvements for 5 subjects out of the total 9 subjects (Fig. 2). Although the performances for remaining subjects were

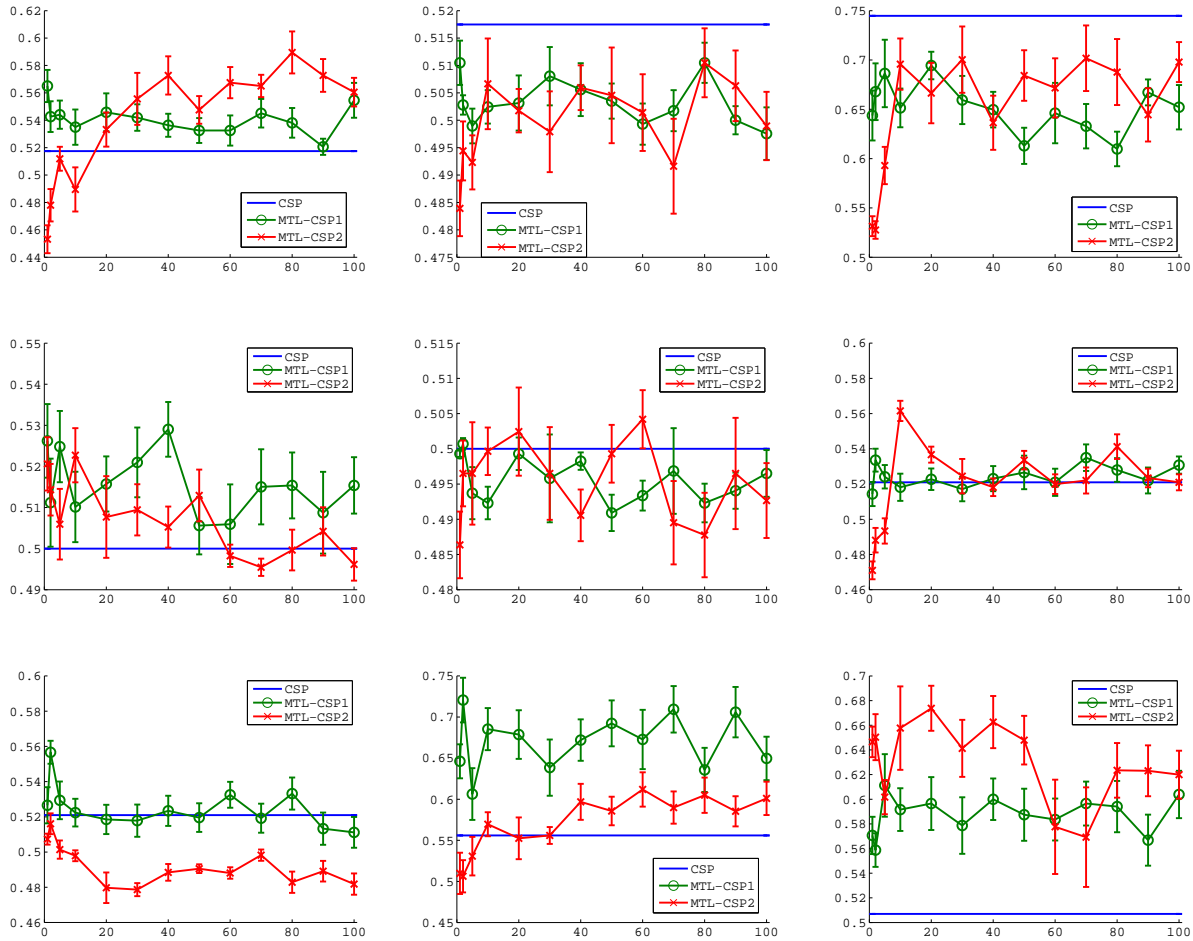


Figure 2. Classification accuracies of all subjects in the BCI competition IV 2a dataset (mean value  $\pm$  1 standard error). The first row shows the graphs for the subjects 1,2 and 3, starting from the left. The second row shows the graphs of the subject 4,5,6, and the third row shows that of the subject 7,8,9. For the subject 1, 4, 8, 6 and 9, there were some improvements by the information transfer, while other subjects suffer from negative side effects.

decreased, the sum of improvement were generally greater than that of the decrease. Thus the overall performance was increased, so that the experiment confirmed effectiveness of our methods (Fig. 3). Because collecting large number of training samples is often difficult in BCI applications, this result will be helpful for increasing recognition rate by better data utilization.

#### ACKNOWLEDGMENTS

This work was supported by National Research Foundation (NRF) of Korea (No. 2010-0018828 and 2010-0018829) and WCU Program (R31-2010-000-10100-0).

#### REFERENCES

- [1] A. Cichocki, Y. Washizawa, T. Rutkowski, H. Bakardjian, A. H. Phan, S. Choi, H. Lee, Q. Zhao, L. Zhang, and Y. Li, "Noninvasive BCIs: Multiway signal-processing array decompositions," *IEEE Computer*, vol. 41, no. 10, pp. 34–42, 2008.
- [2] T. Heskes, "Empirical Bayes for learning to learn," in *Proceedings of the International Conference on Machine Learning (ICML)*, San Francisco, CA, 2000.
- [3] H. Kang, Y. Nam, and S. Choi, "Composite common spatial pattern for subject-to-subject transfer," *IEEE Signal Processing Letters*, vol. 16, no. 8, pp. 683–686, 2009.
- [4] Z. J. Koles, "The quantitative extraction and topographic mapping of the abnormal components," *EEG and Clinical Neurophysiology*, vol. 79, pp. 440–447, 1991.
- [5] F. Lotte and G. Cuntai, "Learning from other subjects helps reducing brain-computer interface calibration time," in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Dallas, TX, 2010.
- [6] W. Wu, Z. Chen, S. Gao, and E. N. Brown, "A probabilistic framework for learning robust common spatial patterns," in *Proceedings of the 31st Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, Minneapolis, MN, 2009.