

Probabilistic Models for Common Spatial Patterns: Parameter-Expanded EM and Variational Bayes

Hyohyeong Kang¹ and Seungjin Choi^{1,2,3}

¹ Department of Computer Science and Engineering,

² Division of IT Convergence Engineering,

³ Department of Creative IT Excellence Engineering,

Pohang University of Science and Technology,

San 31, Hyoja-dong, Nam-gu, Pohang 790-784, Korea

{paanguin,seungjin}@postech.ac.kr

Abstract

Common spatial patterns (CSP) is a popular feature extraction method for discriminating between positive and negative classes in electroencephalography (EEG) data. Two probabilistic models for CSP were recently developed: probabilistic CSP (PCSP), which is trained by expectation maximization (EM), and variational Bayesian CSP (VBCSP) which is learned by variational approximation. Parameter expansion methods use auxiliary parameters to speed up the convergence of EM or the deterministic approximation of the target distribution in variational inference. In this paper, we describe the development of parameter-expanded algorithms for PCSP and VBCSP, leading to PCSP-PX and VBCSP-PX, whose convergence speed-up and high performance are emphasized. The convergence speed-up in PCSP-PX and VBCSP-PX is a direct consequence of parameter expansion methods. The contribution of this study is the performance improvement in the case of CSP, which is a novel development. Numerical experiments on the BCI competition datasets, III IV a and IV 2a demonstrate the high performance and fast convergence of PCSP-PX and VBCSP-PX, as compared to PCSP and VBCSP.

Introduction

Electroencephalography (EEG) is the recording of electrical potentials at multiple sensors placed on the scalp, leading to multivariate time series data reflecting brain activities. EEG classification is a crucial part of non-invasive brain computer interface (BCI) systems, enabling computers to translate a subject's intention or mind into control signals for a device such as a computer, wheelchair, or neuroprosthesis (Wolpaw et al. 2002; Ebrahimi, Vesin, and Garcia 2003; Cichocki et al. 2008).

Common spatial patterns (CSP) is a widely-used discriminative EEG feature extraction method (Blankertz et al. 2008; Koles 1991; Müller-Gerking, Pfurtscheller, and Flyvbjerg 1999; Kang, Nam, and Choi 2009), also known as the Fukunaga-Koontz transform (Fukunaga and Koontz 1970), where we seek a discriminative subspace such that the variance for one class is maximized while the variance for the

other class is minimized. CSP was recently cast into a probabilistic framework (Wu et al. 2009), where a linear Gaussian model for each of the positive/negative classes was considered and the maximum likelihood estimate of the basis matrix shared across two models (positive and negative class models) was shown to yield the same solution as CSP. Bayesian models were also proposed for CSP (Wu et al. 2010; Kang and Choi 2011), where posterior distributions over variables of interest are estimated by variational approximation.

We revisit two probabilistic models for CSP. One is probabilistic CSP (PCSP) (Wu et al. 2009) where the maximum likelihood estimate is determined by the expectation maximization (EM) optimization and the other is variational Bayesian CSP (VBCSP) (Kang and Choi 2011) where the posterior distributions over variables in the model are computed by variational inference in the framework of Bayesian multi-task learning (Heskes 2000). EM and variational inference, while successful, often suffer from slow convergence to the solution. Parameter eXpanded-EM (PX-EM) (Liu, Rubin, and Wu 1998) is a method for accelerating EM, using the over-parameterization of the model. The underlying idea in PX-EM is to use a covariance adjustment to correct the analysis of the M step, thereby exploiting extra information captured in the imputed complete data. Similarly, Parameter eXpanded-VB (PX-VB) (Qi and Jaakkola 2007) expands a model with auxiliary parameters to reduce the coupling between variables in the original model, so that it accelerates the deterministic approximation of the target distribution in variational Bayesian inference.

In this study, we employ the parameter-expansion methods of (Liu, Rubin, and Wu 1998; Qi and Jaakkola 2007; Luttinen and Ilin 2010) in order to develop parameter-expanded algorithms for PCSP and VBCSP, leading to PCSP-PX and VBCSP-PX. By capitalizing on the convergence speed-up by parameter-expansion methods, we show that the expanded models, PCSP-PX and VBCSP-PX, converge to solutions faster than PCSP and VBCSP. In addition, we show that the generalization performance of PCSP-PX and VBCSP-PX is better than that of PCSP and VBCSP. In PCSP and VBCSP, feature vectors are constructed using only variances of the expected latent variables so that the information on covariances is neglected. In contrast, the auxiliary parameters in PCSP-PX and VBCSP-PX re-

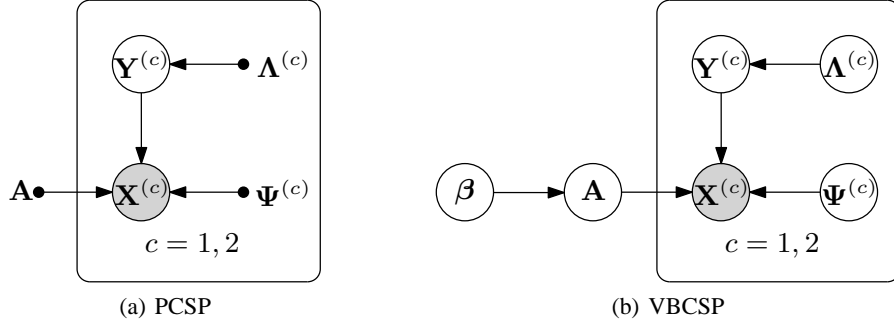


Figure 1: Graphical representation of PCSP and VBCSP models.

duce such information loss by simultaneous diagonalization of the second-order moments of the latent variables. This is found to improve the generalization performance in the classification. Numerical experiments on the BCI competition datasets, III IVa and IV 2a, confirmed that PCSP-PX and VBCSP-PX not only speed-up the computations but also improve the classification performances of the feature vectors, as compared to PCSP and VBCSP.

Probabilistic Models for CSP

We briefly review two probabilistic models, PCSP (Wu et al. 2009) and VBCSP (Kang and Choi 2011). The graphical representations of these models are shown in Fig. 1(a) and 1(b), respectively.

Suppose that EEG signals involving two different mental tasks ($c \in \{1, 2\}$) recorded at D electrodes over multiple trials constitute a D -dimensional vector $\mathbf{x}_t^{(c)}$ for $t = 1, \dots, T_c$, where T_c represents the number of samples obtained over multiple trials. We denote the EEG data matrix by $\mathbf{X}^{(c)} = [\mathbf{x}_1^{(c)}, \dots, \mathbf{x}_{T_c}^{(c)}] \in \mathbb{R}^{D \times T_c}$. The probabilistic model for CSP assumes that data matrices $\mathbf{X}^{(c)}$ are generated by

$$\mathbf{X}^{(c)} = \mathbf{A}\mathbf{Y}^{(c)} + \mathbf{E}^{(c)}, \quad (1)$$

where $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_M] \in \mathbb{R}^{D \times M}$ is the *basis matrix* shared across two classes ($c = 1, 2$),

$$\mathbf{Y}^{(c)} = [\mathbf{y}_1^{(c)}, \dots, \mathbf{y}_{T_c}^{(c)}] \in \mathbb{R}^{M \times T_c},$$

is the *encoding matrix* (corresponding to latent variables), and $\mathbf{E}^{(c)} = [\boldsymbol{\epsilon}_1^{(c)}, \dots, \boldsymbol{\epsilon}_{T_c}^{(c)}] \in \mathbb{R}^{D \times T_c}$ is the *noise matrix*. Latent variables and noise are assumed to follow zero-mean Gaussian distributions,

$$\begin{aligned} \mathbf{y}_t^{(c)} &\sim \mathcal{N}\left(\mathbf{y}_t^{(c)} \mid 0, [\boldsymbol{\Lambda}^{(c)}]^{-1}\right), \\ \boldsymbol{\epsilon}_t^{(c)} &\sim \mathcal{N}\left(\boldsymbol{\epsilon}_t^{(c)} \mid 0, [\boldsymbol{\Psi}^{(c)}]^{-1}\right), \end{aligned}$$

where $\boldsymbol{\Lambda}^{(c)} = \text{diag}(\lambda_1^{(c)}, \dots, \lambda_M^{(c)}) \in \mathbb{R}^{M \times M}$ and $\boldsymbol{\Psi}^{(c)} = \text{diag}(\psi_1^{(c)}, \dots, \psi_D^{(c)}) \in \mathbb{R}^{D \times D}$ are precision matrices for $c = 1, 2$.

PCSP

In PCSP, the basis matrix \mathbf{A} is treated as a matrix of parameters, and their maximum likelihood estimate $\hat{\mathbf{A}}_{ML}$ is determined by EM, where the E-step involves computing the expectation of the complete-data log-likelihood $\log p(\{\mathbf{X}^{(c)}\}, \{\mathbf{Y}^{(c)}\} \mid \mathbf{A}, \{\boldsymbol{\Lambda}^{(c)}\}, \{\boldsymbol{\Psi}^{(c)}\})$ with respect to the posterior distribution $p(\{\mathbf{Y}^{(c)}\} \mid \{\mathbf{X}^{(c)}\})$ and the M-step re-estimates \mathbf{A} as well as other model parameters $\{\boldsymbol{\Lambda}^{(c)}, \boldsymbol{\Psi}^{(c)}\}$. It was shown in (Wu et al. 2009) that $\hat{\mathbf{A}}_{ML}^{-\top}$ is equal to the linear transformation matrix computed in CSP, in the case of zero noise limit and when \mathbf{A} is a square matrix. CSP feature vectors are constructed by taking logarithms of top- n variances of the projected variables for each class within each trial. In the case of PCSP, CSP features are computed using logarithms of top- n variances for each class of posterior means over latent variables in each trial.

VBCSP

In VBCSP, the basis matrix \mathbf{A} is treated as a matrix of random variables, and the automatic relevance determination (ARD) prior is applied to it, i.e.,

$$p(\mathbf{A} \mid \mathbf{D}_\beta) = \prod_{m=1}^M \mathcal{N}([\mathbf{A}]_{:,m} \mid 0, \beta_m^{-1} \mathbf{I}_D), \quad (2)$$

where $\mathbf{I}_D \in \mathbb{R}^{D \times D}$ is the identity matrix, $\mathbf{D}_\beta \in \mathbb{R}^{M \times M} = \text{diag}(\beta_1, \dots, \beta_M)$, and the precision hyperparameters β_m are assumed to follow Gamma distribution:

$$\beta_m \sim \text{Gam}(a_0^\beta, b_0^\beta). \quad (3)$$

Inferring posterior distributions over β_m leads us to predict an appropriate number of columns in \mathbf{A} . ARD priors are also applied to $\{\lambda_m^{(c)}\}$ and $\{\psi_d^{(c)}\}$ (which are diagonal entries of precision matrices $\boldsymbol{\Lambda}^{(c)}$ and $\boldsymbol{\Psi}^{(c)}$):

$$\lambda_m^{(c)} \sim \text{Gam}(a_0^\lambda, b_0^\lambda), \quad \psi_d^{(c)} \sim \text{Gam}(a_0^\psi, b_0^\psi). \quad (4)$$

Variational posterior distributions are determined by variational Bayesian inference (Kang and Choi 2011). Again, logarithms of top- n variances for each class of the variational posterior means of latent variables within each trial are used as the CSP features.

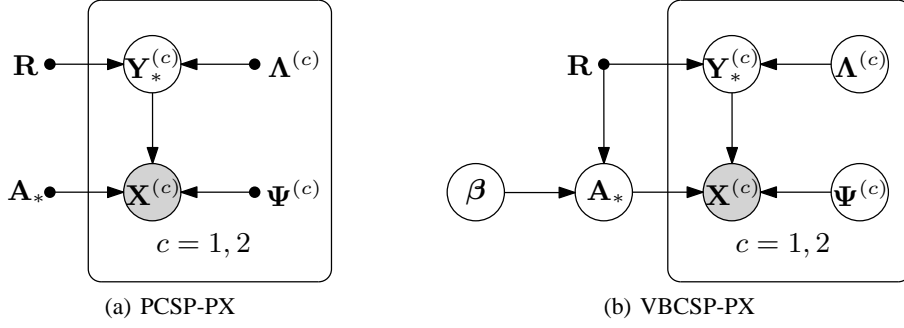


Figure 2: Graphical representation of PCSP-PX and VBCSP-PX models.

Parameter-Expanded Algorithms

In this section, we present the main contribution of this paper, parameter-expanded algorithms for both PCSP and VBCSP. First, we introduce the parameter-expanded models (PCSP-PX and VBCSP-PX), as shown in Fig. 2(a) and 2(b), inspired by (Luttinen and Ilin 2010). Then we develop a parameter-expanded EM algorithm for PCSP-PX and parameter-expanded variational Bayesian inference for VBCSP-PX.

We introduce an invertible matrix $\mathbf{R} \in \mathbb{R}^{M \times M}$. We define $\mathbf{A}_* = \mathbf{A}\mathbf{R}$ and $\mathbf{Y}_*^{(c)} = \mathbf{R}^{-1}\mathbf{Y}^{(c)}$. Then, we can write the model (1) as

$$\mathbf{X}^{(c)} = \mathbf{A}_* \mathbf{Y}_*^{(c)} + \mathbf{E}^{(c)}, \quad (5)$$

for $c = 1, 2$, since $\mathbf{A}_* \mathbf{Y}_*^{(c)} = \mathbf{A}\mathbf{R}\mathbf{R}^{-1}\mathbf{Y}^{(c)}$. The invertible matrix \mathbf{R} introduces a transformation of the encoding matrix $\{\mathbf{Y}^{(c)}\}$ while preserving the conditional distribution $p(\mathbf{X}^{(c)} | \mathbf{A}, \mathbf{Y}^{(c)})$. We optimize the auxiliary parameters \mathbf{R} such that the expected complete-data log-likelihood (for PCSP-PX) or the variational lower-bound (for VBCSP-PX) is maximized.

PCSP-PX

We present a parameter-expanded EM algorithm to estimate the model parameters $\{\mathbf{A}_*, \Lambda^{(c)}, \Psi^{(c)}\}$ as well as auxiliary parameters \mathbf{R} . The basis matrix \mathbf{A} in the original model (1) is recovered by $\mathbf{A}_*\mathbf{R}^{-1}$. To this end, we consider the complete-data likelihood in the expanded model (5) (shown in Fig. 2(a)):

$$\begin{aligned} p(\{\mathbf{X}^{(c)}\}, \{\mathbf{Y}_*^{(c)}\} | \mathbf{A}_*, \{\Lambda^{(c)}\}, \{\Psi^{(c)}\}, \mathbf{R}) \\ = \prod_{c=1}^2 p(\mathbf{X}^{(c)} | \mathbf{Y}_*^{(c)}, \mathbf{A}_*, \Psi^{(c)}) p(\mathbf{Y}_*^{(c)} | \Lambda^{(c)}, \mathbf{R}), \end{aligned}$$

where

$$\begin{aligned} p(\mathbf{X}^{(c)} | \mathbf{Y}_*^{(c)}, \mathbf{A}_*, \Psi^{(c)}) \\ = \prod_{t=1}^{T_c} \mathcal{N}(\mathbf{x}_t^{(c)} | \mathbf{A}_* \mathbf{y}_{*t}^{(c)}, [\Psi^{(c)}]^{-1}), \\ p(\mathbf{Y}_*^{(c)} | \Lambda^{(c)}, \mathbf{R}) \\ = \prod_{t=1}^{T_c} \mathcal{N}(\mathbf{y}_{*t}^{(c)} | 0, [\mathbf{R}^\top \Lambda^{(c)} \mathbf{R}]^{-1}). \end{aligned}$$

In the E-step, we compute the expected complete-data log-likelihood $\langle \mathcal{L}_c \rangle$

$$\sum_{c=1}^2 \langle \log p(\mathbf{X}^{(c)}, \mathbf{Y}_*^{(c)} | \Theta) \rangle,$$

where the expectation $\langle \cdot \rangle$ is taken with respect to the posterior distribution over latent variables $p(\mathbf{Y}_*^{(c)} | \mathbf{X}^{(c)})$ given the current estimate of parameters $\Theta = \{\mathbf{A}_*, \Lambda^{(c)}, \Psi^{(c)}, \mathbf{R}\}$. In the M-step, we re-estimate parameters Θ that maximize $\langle \mathcal{L}_c \rangle$ computed in the E-step. The EM iteration for PCSP-PX, which is summarized in Algorithm 1, alternates between the E-step and M-step until convergence.

In contrast to PCSP, the auxiliary parameters \mathbf{R} should also be optimized. The stationary point equation for \mathbf{R} is given by

$$\begin{aligned} \frac{\partial \langle \mathcal{L}_c \rangle}{\partial \mathbf{R}} &= \left(\sum_{c=1}^2 T_c \right) \mathbf{R}^{-\top} - \sum_{c=1}^2 \Lambda^{(c)} \mathbf{R} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \\ &= 0, \end{aligned} \quad (6)$$

leading to

$$\sum_{c=1}^2 \Lambda^{(c)} \mathbf{R} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \mathbf{R}^\top = \left(\sum_{c=1}^2 T_c \right) \mathbf{I}_M, \quad (7)$$

which is solved for \mathbf{R} by simultaneous diagonalization of the second-order moments of encodings $\langle \mathbf{Y}_*^{(1)} \mathbf{Y}_*^{(1)\top} \rangle$ and $\langle \mathbf{Y}_*^{(2)} \mathbf{Y}_*^{(2)\top} \rangle$, followed by re-scaling.

Algorithm 1 EM for PCSP-PX

Input: EEG data $\{\mathbf{X}^{(c)}\}$.**Output:** estimate of parameters $\Theta = \{\mathbf{A}_*, \Lambda^{(c)}, \Psi^{(c)}, \mathbf{R}\}$ **initialize** $\Theta = \{\mathbf{A}_*, \Lambda^{(c)}, \Psi^{(c)}, \mathbf{R}\}$.**repeat****E-step** Calculate the posterior distribution over latent variables $p(\mathbf{y}_{*t}^{(c)} | \mathbf{x}_t^{(c)})$:

$$\begin{aligned} p(\mathbf{y}_{*t}^{(c)} | \mathbf{x}_t^{(c)}) &= \mathcal{N}(\mathbf{y}_{*t}^{(c)} | \boldsymbol{\mu}_t^{(c)}, \boldsymbol{\Sigma}^{(c)}), \\ \boldsymbol{\mu}_t^{(c)} &= \boldsymbol{\Sigma}^{(c)} \mathbf{A}_*^\top \Psi^{(c)} \mathbf{x}_t^{(c)}, \\ [\boldsymbol{\Sigma}^{(c)}]^{-1} &= \mathbf{R}^\top \Lambda^{(c)} \mathbf{R} + \mathbf{A}_*^\top \Psi^{(c)} \mathbf{A}_*. \end{aligned}$$

M-step Re-estimate Θ :- Update $\{\mathbf{A}_*, \Lambda^{(c)}, \Psi^{(c)}\}$:

$$\begin{aligned} [\mathbf{A}_*]_{d,:} &= \left(\sum_{c=1}^2 \psi_d^{(c)} [\mathbf{X}^{(c)}]_{d,:} \langle \mathbf{Y}_*^{(c)\top} \rangle \right) \\ &\quad \left(\sum_{c=1}^2 \psi_d^{(c)} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \right)^{-1}, \\ [\psi_d^{(c)}]^{-1} &= \frac{1}{T_c} \left[\mathbf{X}^c \mathbf{X}^{(c)\top} - 2 \mathbf{X}^{(c)} \langle \mathbf{Y}_*^{(c)\top} \rangle \mathbf{A}_*^\top \right. \\ &\quad \left. + \mathbf{A}_* \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \mathbf{A}_*^\top \right]_{d,d}, \\ [\lambda_m^{(c)}]^{-1} &= \frac{1}{T_c} \left[\mathbf{R} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \mathbf{R}^\top \right]_{m,m}. \end{aligned}$$

- Solve (7) for \mathbf{R} by simultaneous diagonalization of $\langle \mathbf{Y}_*^{(1)} \mathbf{Y}_*^{(1)\top} \rangle$ and $\langle \mathbf{Y}_*^{(2)} \mathbf{Y}_*^{(2)\top} \rangle$ to update \mathbf{R} .**until** convergence.

The posterior distribution over latent variables in the original model is easily computed by

$$p(\mathbf{y}_t^{(c)} | \mathbf{x}_t^{(c)}) = \mathcal{N}(\mathbf{y}_t^{(c)} | \mathbf{R} \boldsymbol{\mu}_t^{(c)}, \mathbf{R} \boldsymbol{\Sigma}^{(c)} \mathbf{R}^\top),$$

where $\boldsymbol{\mu}_t^{(c)}$ and $\boldsymbol{\Sigma}^{(c)}$ are calculated in the E-step in Algorithm 1. We compute CSP features using posterior means $\mathbf{R} \boldsymbol{\mu}_t^{(c)} = \mathbf{R} \boldsymbol{\Sigma}^{(c)} \mathbf{R}^\top \mathbf{A}_*^\top \Psi^{(c)} \mathbf{x}_t^{(c)}$ that correspond to projected variables in CSP. Given test trial data $\mathbf{X} \in \mathbb{R}^{D \times T}$, we first compute T M -dimensional posterior mean vectors over latent variables, constructing posterior mean matrices $\bar{\mathbf{Y}}^{(c)} \in \mathbb{R}^{M \times T}$ for $c = 1, 2$,

$$\bar{\mathbf{Y}}^{(c)} = \mathbf{R} \boldsymbol{\Sigma}^{(c)} \mathbf{R}^\top \mathbf{A}_*^\top \Psi^{(c)} \mathbf{X}.$$

To model the project variables in CSP, we average $\bar{\mathbf{Y}}^{(c)}$ of the two classes considering the class prior probability as $p(\mathbf{X} \in (c)) = \frac{T_c}{T_1 + T_2}$,

$$\bar{\mathbf{Y}} = \sum_{c=1}^2 \frac{T_c}{T_1 + T_2} \bar{\mathbf{Y}}^{(c)}. \quad (8)$$

Then, we build a vector $\mathbf{z} \in \mathbb{R}^M$, the m -th entry of which is computed by

$$[\mathbf{z}]_m = \log \left(\frac{1}{T} [\bar{\mathbf{Y}} \bar{\mathbf{Y}}^\top]_{m,m} - \left(\frac{1}{T} [\bar{\mathbf{Y}} \mathbf{1}_T]_m \right)^2 \right), \quad (9)$$

where $\mathbf{1}_T \in \mathbb{R}^T$ is the vector of all ones. We choose $2n$ $[\mathbf{z}]_m$'s with m corresponding to n largest and n smallest values of the ratio $\lambda_m^{(1)} / \lambda_m^{(2)}$ to construct the CSP feature vector $\mathbf{f} \in \mathbb{R}^{2n}$.

VBCSP-PX

We present a parameter-expanded variational Bayesian inference to speed up the deterministic approximation of the posterior distributions over variables $\mathcal{Z} = \{\mathbf{Y}_*^{(c)}, \mathbf{A}_*, \beta_m, \Psi^{(c)}, \Lambda^{(c)}\}$ with auxiliary parameters \mathbf{R} .

Variational posterior distributions over \mathbf{A} and $\mathbf{Y}^{(c)}$ in the original model are easily recovered by variable transformation: $\mathbf{A} = \mathbf{A}_* \mathbf{R}^{-1}$ and $\mathbf{Y}^{(c)} = \mathbf{R} \mathbf{Y}_*^{(c)}$. We write the joint distribution over $\{\mathbf{X}^{(c)}\}$ and \mathcal{Z} in the expanded model (shown in Fig. 2(b)) with prior distributions defined in (2), (3), and (4) as

$$\begin{aligned} p(\{\mathbf{X}^{(c)}\}, \{\mathbf{Y}_*^{(c)}\}, \mathbf{A}_*, \{\beta_m\}, \{\Lambda^{(c)}\}, \{\Psi^{(c)}\} | \mathbf{R}) \\ = \prod_{c=1}^2 p(\mathbf{X}^{(c)} | \mathbf{Y}_*^{(c)}, \mathbf{A}_*, \Psi^{(c)}) p(\mathbf{Y}_*^{(c)} | \Lambda^{(c)}, \mathbf{R}) \\ p(\Lambda^{(c)}) p(\Psi^{(c)}) p(\mathbf{A}_* | \mathbf{D}_\beta, \mathbf{R}) \prod_{m=1}^M p(\beta_m), \end{aligned}$$

where

$$p(\mathbf{Y}_*^{(c)} | \Lambda^{(c)}, \mathbf{R}) = \prod_{t=1}^{T_c} \mathcal{N}(\mathbf{y}_{*t}^{(c)} | 0, [\mathbf{R}^\top \Lambda^{(c)} \mathbf{R}]^{-1}).$$

The prior distribution over $\mathbf{A}_* \in \mathbb{R}^{D \times M}$ is assumed to be *matrix-variate Gaussian* since it is not column-wise independent in contrast to (2). Thus we assume

$$p(\mathbf{A}_* | \mathbf{D}_\beta, \mathbf{R}) = \mathcal{N}_{D \times M}(\mathbf{A}_* | 0, \mathbf{I}_D \otimes \mathbf{R}^\top \mathbf{D}_\beta^{-1} \mathbf{R}),$$

where matrix-variate Gaussian distribution for a random matrix $\mathbf{B} \in \mathbb{R}^{D \times M}$ with mean matrix $\mathbf{M} \in \mathbb{R}^{D \times M}$ and covariance matrix $\boldsymbol{\Omega}_D \otimes \boldsymbol{\Omega}_M$ ($\boldsymbol{\Omega}_D \in \mathbb{R}^{D \times D}$ and $\boldsymbol{\Omega}_M \in \mathbb{R}^{M \times M}$) takes the form

$$\begin{aligned} \mathcal{N}_{D \times M}(\mathbf{B} | \mathbf{M}, \boldsymbol{\Omega}_D \otimes \boldsymbol{\Omega}_M) \\ = (2\pi)^{-\frac{DM}{2}} |\boldsymbol{\Omega}_D|^{-\frac{M}{2}} |\boldsymbol{\Omega}_M|^{-\frac{D}{2}} \\ \exp \left[-\frac{1}{2} \text{tr} \left(\boldsymbol{\Omega}_D^{-1} (\mathbf{B} - \mathbf{M}) \boldsymbol{\Omega}_M^{-1} (\mathbf{B} - \mathbf{M})^\top \right) \right]. \end{aligned}$$

Note that the prior distribution over \mathbf{A} in the original model, given in (2), can also be written as

$$p(\mathbf{A} | \mathbf{D}_\beta) = \mathcal{N}_{D \times M}(\mathbf{A} | 0, \mathbf{I}_D \otimes \mathbf{D}_\beta^{-1}).$$

The variational inference involves the maximization of a lower-bound on the *marginal log-likelihood* given by

$$\begin{aligned}
& \log p\left(\left\{\mathbf{X}^{(c)}\right\} \mid \mathbf{R}\right) \\
&= \log \int p\left(\left\{\mathbf{X}^{(c)}\right\}, \mathcal{Z} \mid \mathbf{R}\right) d\mathcal{Z} \\
&\geq \int q(\mathcal{Z}) \log \frac{p\left(\left\{\mathbf{X}^{(c)}\right\}, \mathcal{Z} \mid \mathbf{R}\right)}{q(\mathcal{Z})} d\mathcal{Z} \\
&\equiv \mathcal{F}(q \mid \mathbf{R}), \tag{10}
\end{aligned}$$

where variational posterior distribution $q(\mathcal{Z})$ is assumed to factorize as

$$q(\mathcal{Z}) = q(\mathbf{A}_*) q\left(\left\{\mathbf{y}_{*t}^{(c)}\right\}\right) q\left(\left\{\beta_m\right\}\right) q\left(\left\{\psi_d^{(c)}\right\}\right) q\left(\left\{\lambda_m^{(c)}\right\}\right).$$

Variational posterior distributions over each variable are alternatively updated such that the lower-bound $\mathcal{F}(q \mid \mathbf{R})$ is maximized. Updating equations are summarized in Algorithm 2. In addition, auxiliary parameters \mathbf{R} are optimized by the maximization of $\mathcal{F}(q \mid \mathbf{R})$, given the variational posterior distribution $q(\mathcal{Z})$. We consider the terms involving \mathbf{R} in the lower-bound $\mathcal{F}(q \mid \mathbf{R})$, given by

$$\begin{aligned}
& \frac{\sum_{c=1}^2 T_c - D}{2} \log(|\mathbf{R}|^2) \\
& - \frac{1}{2} \sum_{c=1}^2 \text{tr}\left(\langle \mathbf{\Lambda}^{(c)} \rangle \mathbf{R} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \mathbf{R}^\top\right) \\
& - \frac{1}{2} \text{tr}\left(\langle \mathbf{D}_\beta \rangle \mathbf{R}^{-\top} \langle \mathbf{A}_*^\top \mathbf{A}_* \rangle \mathbf{R}^{-1}\right). \tag{11}
\end{aligned}$$

Suppose a_0^β and b_0^β are set to small values so that $\langle \beta_m \rangle$ can be approximated as

$$\begin{aligned}
\langle \beta_m \rangle &= \frac{a_m^\beta}{b_m^\beta} = \frac{a_0^\beta + D/2}{b_0^\beta + \frac{1}{2} \left[\mathbf{R}^{-\top} \langle \mathbf{A}_*^\top \mathbf{A}_* \rangle \mathbf{R}^{-1} \right]_{m,m}} \\
&\simeq \frac{D}{\left[\mathbf{R}^{-\top} \langle \mathbf{A}_*^\top \mathbf{A}_* \rangle \mathbf{R}^{-1} \right]_{m,m}},
\end{aligned}$$

implying that

$$\begin{aligned}
& \text{tr}\left(\langle \mathbf{D}_\beta \rangle \mathbf{R}^{-\top} \langle \mathbf{A}_*^\top \mathbf{A}_* \rangle \mathbf{R}^{-1}\right) \\
&= \sum_{m=1}^M \langle \beta_m \rangle \left[\mathbf{R}^{-\top} \langle \mathbf{A}_*^\top \mathbf{A}_* \rangle \mathbf{R}^{-1} \right]_{m,m} \\
&\simeq MD. \tag{12}
\end{aligned}$$

Applying the approximation (12) to (11), as in (Lutinen and Ilin 2010), we have the stationary point equation for \mathbf{R} given by

$$\begin{aligned}
& \frac{\partial \mathcal{F}(q \mid \mathbf{R})}{\partial \mathbf{R}} \\
&= \left(\sum_{c=1}^2 T_c - D \right) \mathbf{R}^{-\top} - \sum_{c=1}^2 \langle \mathbf{\Lambda}^{(c)} \rangle \mathbf{R} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \\
&= 0. \tag{13}
\end{aligned}$$

Again, we can solve the equation for \mathbf{R} by simultaneous diagonalization of the second-order moments of encodings $\langle \mathbf{Y}_*^{(1)} \mathbf{Y}_*^{(1)\top} \rangle$ and $\langle \mathbf{Y}_*^{(2)} \mathbf{Y}_*^{(2)\top} \rangle$, followed by re-scaling, as in PCSP-PX.

Algorithm 2 Variational Bayesian Inference for VBCSP-PX

Input: EEG data $\{\mathbf{X}^{(c)}\}$

Output: approximated posterior $q(\mathcal{Z})$ for variables $\mathcal{Z} = \{\mathbf{A}_*, \mathbf{Y}_*^{(c)}, \mathbf{\Lambda}^{(c)}, \mathbf{\Psi}^{(c)}\}$ and the auxiliary parameter \mathbf{R}
initialize $q(\mathcal{Z})$.

repeat

- Update $q(\mathbf{A}_*) = \prod_{d=1}^D \mathcal{N}([\mathbf{A}_*]_{d,:} \mid \bar{\nu}_d, \mathbf{\Omega}_d)$ by

$$\begin{aligned}
[\mathbf{\Omega}_d]^{-1} &= \mathbf{R}^{-1} \langle \mathbf{D}_\beta \rangle \mathbf{R}^{-\top} \\
&+ \sum_{c=1}^2 \langle \psi_d^{(c)} \rangle \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle, \\
\bar{\nu}_d &= \sum_{c=1}^2 \langle \psi_d^{(c)} \rangle [\mathbf{X}^{(c)}]_{d,:} \langle \mathbf{Y}_*^{(c)\top} \rangle \mathbf{\Omega}_d.
\end{aligned}$$

- Update $q(\mathbf{y}_{*t}^{(c)}) = \mathcal{N}(\boldsymbol{\mu}_t^{(c)}, \boldsymbol{\Sigma}^{(c)})$ by

$$\begin{aligned}
[\boldsymbol{\Sigma}^{(c)}]^{-1} &= \mathbf{R}^\top \mathbf{\Lambda}^{(c)} \mathbf{R} + \langle \mathbf{A}_*^\top \mathbf{\Psi}^{(c)} \mathbf{A}_* \rangle, \\
\boldsymbol{\mu}_t^{(c)} &= \boldsymbol{\Sigma}^{(c)} \langle \mathbf{A}_*^\top \rangle \langle \mathbf{\Psi}^{(c)} \rangle \mathbf{x}_t^{(c)}.
\end{aligned}$$

- Update $q(\beta_m) = \text{Gam}(a_m^\beta, b_m^\beta)$ by

$$\begin{aligned}
a_m^\beta &= a_0^\beta + D/2, \\
b_m^\beta &= b_0^\beta + \frac{1}{2} \left[\mathbf{R}^{-\top} \langle \mathbf{A}_*^\top \mathbf{A}_* \rangle \mathbf{R}^{-1} \right]_{m,m}.
\end{aligned}$$

- Update $q(\psi_d^{(c)}) = \text{Gam}(a_d^{\psi(c)}, b_d^{\psi(c)})$ by

$$\begin{aligned}
a_d^{\psi(c)} &= a_0^\psi + T_c/2, \\
b_d^{\psi(c)} &= b_0^\psi + \frac{1}{2} \left[\mathbf{X}^{(c)} \mathbf{X}^{(c)\top} \right. \\
&\quad \left. - 2 \langle \mathbf{A}_* \rangle \langle \mathbf{Y}_*^{(c)} \rangle \mathbf{X}^{(c)\top} \right. \\
&\quad \left. + \langle \mathbf{A}_* \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \mathbf{A}_*^\top \rangle \right]_{d,d}.
\end{aligned}$$

- Update $q(\lambda_m^{(c)}) = \text{Gam}(a_m^{\lambda(c)}, b_m^{\lambda(c)})$ by

$$\begin{aligned}
a_m^{\lambda(c)} &= a_0^\lambda + T_c/2, \\
b_m^{\lambda(c)} &= b_0^\lambda + \frac{1}{2} \left[\mathbf{R} \langle \mathbf{Y}_*^{(c)} \mathbf{Y}_*^{(c)\top} \rangle \mathbf{R}^\top \right]_{m,m}.
\end{aligned}$$

- Solve (13) for \mathbf{R} by simultaneous diagonalization of $\langle \mathbf{Y}_*^{(1)} \mathbf{Y}_*^{(1)\top} \rangle$ and $\langle \mathbf{Y}_*^{(2)} \mathbf{Y}_*^{(2)\top} \rangle$ to update \mathbf{R} .

until convergence.

The approximated posterior over latent variables in the

original model is computed by

$$q\left(\mathbf{y}_t^{(c)}\right)=\mathcal{N}\left(\mathbf{y}_t^{(c)}\mid\mathbf{R}\boldsymbol{\mu}_t,\mathbf{R}\boldsymbol{\Sigma}^{(c)}\mathbf{R}^\top\right),$$

where $\boldsymbol{\mu}_t$ and $\boldsymbol{\Sigma}^{(c)}$ are calculated as in Algorithm 2. We compute a CSP feature vector \mathbf{f} for a test trial $\mathbf{X}\in\mathbb{R}^{D\times T}$ using the posterior means $\mathbf{R}\boldsymbol{\mu}_t^{(c)}=\mathbf{R}\boldsymbol{\Sigma}^{(c)}\langle\mathbf{A}_*^\top\rangle\langle\boldsymbol{\Psi}^{(c)}\rangle\mathbf{x}_t^{(c)}$,

$$\overline{\mathbf{Y}}^{(c)}=\mathbf{R}\boldsymbol{\Sigma}^{(c)}\mathbf{R}^\top\langle\mathbf{A}^\top\rangle\langle\boldsymbol{\Psi}^{(c)}\rangle\mathbf{X}.$$

We compute $\overline{\mathbf{Y}}$ and $[z]_m$ as in (8) and (9). Then we select $[z]_m$'s with m corresponding to n -largest and n -smallest values of the ratio of $\langle\lambda_m^{(1)}\rangle/\langle\lambda_m^{(2)}\rangle$ to construct the CSP feature vector $\mathbf{f}\in\mathbb{R}^{2n}$.

Numerical Experiments

We compared the performances of PCSP, VBCSP, PCSP-PX, and VBCSP-PX on the BCI competition datasets, III IVa (Blankertz et al. 2006)¹ and IV 2a². Both datasets consist of the EEG measurements of several subjects during motor imagery tasks. The EEG data was pre-processed by band-pass filtering, to emphasize important frequency bands for recognizing the motor imagery tasks. Every trial was divided into the same number of time intervals after each visual cue, which contains EEG variation caused by the imagination of the subject.

We extracted feature vectors $\mathbf{f}\in\mathbb{R}^{2n}$ using PCSP, VBCSP, PCSP-PX, and VBCSP-PX, and we applied the linear discriminant analysis (LDA) to transform these feature vectors down to scalar values which are fed into a minimum distance classifier. We set $D=M$ and $n=3$ for every model. The classification performance of each model is represented by the prediction accuracy of the LDA classifier on the test trials. The accuracy was calculated as the ratio of the number of correctly classified trials to the total number of test trials. We repeated the experiments 10 times, varying the number of training trials, while the number of test trials was fixed. We selected half of the trials in each data as the test trials, and randomly selected some of the remaining trials as the training trials. The classes were strictly balanced by selecting the same number of trials from each class.

BCI competition III IVa dataset was collected from five subjects using 118 electrodes ($D=118$) during the imagery movements of the right hand and right foot. The trials were separated by up to 3.5s after each cue, and we used the down-sampled version (100 Hz) of the data. 140 trials were conducted for each subject and each class. BCI competition IV 2a dataset contains 4 motor imagery tasks of 9 subjects, recorded using 22 electrodes ($D=22$). We considered only the binary classification problem so that we selected the imagery left/right hand movement classes. The trials were separated by from 3.5s to 5.5s after each cue, and the sampling rate was 250 Hz. 144 trials were conducted for each subject and each class.

¹<http://www.bbci.de/competition/iii/>

²<http://www.bbci.de/competition/iv/>

Compared to PCSP and VBCSP, parameter-expanded algorithms PCSP-PX and VBCSP-PX perform additional computation to optimize \mathbf{R} at every iteration. However, PCSP-PX and VBCSP-PX converge in a smaller number of iterations; hence, they are faster than PCSP and VBCSP, respectively (Table 1). In general, the classification performance of PCSP-PX and VBCSP-PX was also higher than that of PCSP and VBCSP (Fig. 3).

Conclusions

We have presented two new parameter-expanded algorithms for PCSP and VBCSP, leading to PCSP-PX and VBCSP-PX, where we expanded the models using auxiliary parameters \mathbf{R} to speed up the convergence as well as to improve the performance. The auxiliary parameters \mathbf{R} were estimated by simultaneous diagonalization of $\langle\mathbf{Y}_*^{(1)}\mathbf{Y}_*^{(1)\top}\rangle$ and $\langle\mathbf{Y}_*^{(2)}\mathbf{Y}_*^{(2)\top}\rangle$, reducing the coupling so that the convergence was accelerated and the performance was improved, while CSP features determined by PCSP or VBCSP ignored off-diagonal entries of the empirical second-order moment matrix of posterior mean vectors. Numerical experiments on the BCI competition datasets, III IVa and IV 2a, demonstrated the high performance of PCSP-PX and VBCSP-PX, as compared to their counterparts PCSP and VBCSP.

Acknowledgments: This work was supported by National Research Foundation (NRF) of Korea (2011-0018283, 2011-0018284), MEST Converging Research Center Program (2011K000673), NIPA Program of Software Engineering Technologies Development and Experts Education, MKE and NIPA "IT Consilience Creative Program" (C1515-1121-0003), and NRF World Class University Program (R31-10100).

References

- Blankertz, B.; Müller, K. R.; Krusierski, D. J.; Schalk, G.; Wolpaw, J. R.; Schlögl, A.; Pfurtscheller, G.; and Birbaumer, N. 2006. The BCI competition III: Validating alternative approaches to actual BCI problems. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 14:153–159.
- Blankertz, B.; Tomioka, R.; Lemm, S.; Kawanabe, M.; and Müller, K. R. 2008. Optimizing spatial filters for robust EEG single-trial analysis. *IEEE Signal Processing Magazine* 41–56.
- Cichocki, A.; Washizawa, Y.; Rutkowski, T.; Bakardjian, H.; Phan, A. H.; Choi, S.; Lee, H.; Zhao, Q.; Zhang, L.; and Li, Y. 2008. Noninvasive BCIs: Multiway signal-processing array decompositions. *IEEE Computer* 41(10):34–42.
- Ebrahimi, T.; Vesin, J. F.; and Garcia, G. 2003. Brain-computer interface in multimedia communication. *IEEE Signal Processing Magazine* 20(1):14–24.
- Fukunaga, K., and Koontz, W. L. G. 1970. Application of the Karhunen-Loève expansion to feature selection and ordering. *IEEE Transactions on Computers* 19(4):311–318.

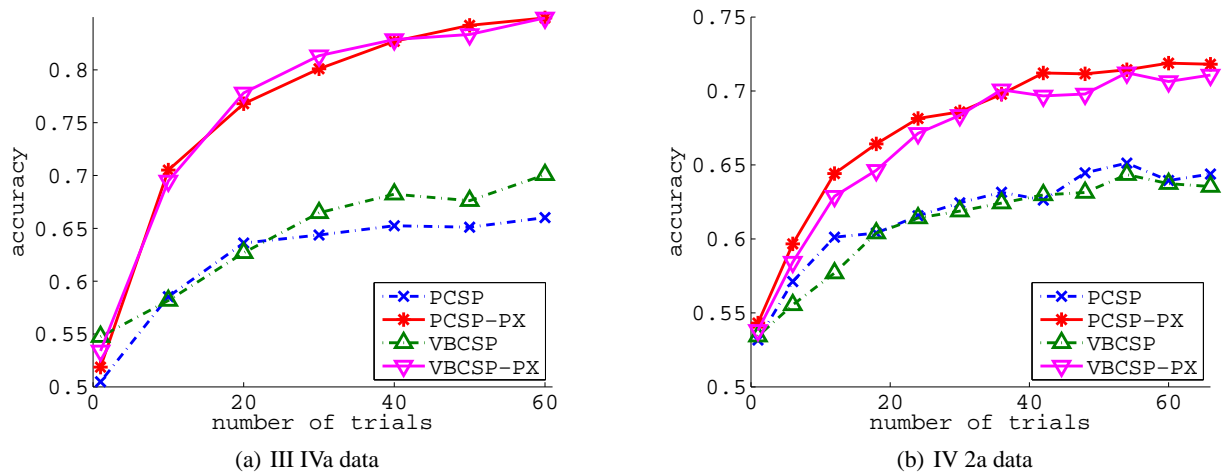


Figure 3: Classification performances of existing methods (PCSP and VBCSP) as well as our proposed methods (PCSP-PX and VBCSP-PX).

Table 1: Performance comparison in terms of number of iterations and run time. For each algorithm, the iterations stop when the variant of the expected complete-data log-likelihood (for PCSP and PCSP-PX) or the variational lower-bound (for VBCSP and VBCSP-PX) falls within a pre-defined value. The maximum number of iterations was set as 100 for each algorithm. The 'number of iterations', 'run time' and 'run time per iteration' were averaged over multiple runs with 10 different training sets.

data	measure	PCSP	PCSP-PX	VBCSP	VBCSP-PX
III IVa	number of iterations	86.6400 \pm 24.0742	15.7057 \pm 7.0065	54.9114 \pm 26.1892	13.8800 \pm 3.4429
	run time (sec)	14.3137 \pm 4.5311	3.0026 \pm 1.4251	37.1477 \pm 16.7673	12.2370 \pm 3.2140
	run time per iteration	0.1631 \pm 0.0166	0.1901 \pm 0.0246	0.6917 \pm 0.1013	0.8804 \pm 0.0638
IV 2a	number of iterations	58.0417 \pm 27.4807	34.2093 \pm 16.1989	18.5046 \pm 7.9813	13.5111 \pm 3.8341
	run time (sec)	0.1702 \pm 0.0841	0.1338 \pm 0.0668	0.1579 \pm 0.0740	0.1337 \pm 0.0427
	run time per iteration	0.0029 \pm 0.0003	0.0039 \pm 0.0004	0.0084 \pm 0.0009	0.0098 \pm 0.0011

Heskes, T. 2000. Empirical Bayes for learning to learn. In *Proceedings of the International Conference on Machine Learning (ICML)*.

Kang, H., and Choi, S. 2011. Bayesian multi-task learning for common spatial patterns. In *Proceedings of the IEEE International Workshop on Pattern Recognition in NeuroImaging (PRNI)*.

Kang, H.; Nam, Y.; and Choi, S. 2009. Composite common spatial pattern for subject-to-subject transfer. *IEEE Signal Processing Letters* 16(8):683–686.

Koles, Z. J. 1991. The quantitative extraction and topographic mapping of the abnormal components. *EEG and Clinical Neurophysiology* 79:440–447.

Liu, C.; Rubin, D. B.; and Wu, Y. N. 1998. Parameter expansion to accelerate EM: The PX-EM algorithm. *Biometrika* 85(4):755–770.

Luttinen, J., and Ilin, A. 2010. Transformations in variational Bayesian factor analysis to speed up learning. *Neurocomputing* 73:1093–1102.

Müller-Gerking, J.; Pfurtscheller, G.; and Flyvbjerg, H. 1999. Designing optimal spatial filters for single-trial EEG

classification in a movement task. *Clinical Neurophysiology* 110:787–798.

Qi, Y., and Jaakkola, T. S. 2007. Parameter expanded variational Bayesian methods. In *Advances in Neural Information Processing Systems (NIPS)*, volume 19. MIT Press.

Wolpaw, J. R.; Birbaumer, N.; McFarland, D. J.; Pfurtscheller, G.; and Vaughan, T. M. 2002. Brain-computer interfaces for communication and control. *Clinical Neurophysiology* 113:767–791.

Wu, W.; Chen, Z.; Gao, S.; and Brown, E. N. 2009. 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*.

Wu, W.; Chen, Z.; Gao, S.; and Brown, E. N. 2010. Hierarchical Bayesian modeling of inter-trial variability and variational Bayesian learning of common spatial patterns from multichannel EEG. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*.