## NOVEL IMAGE CLASSIFICATION BASED ON INTEGRATION OF EEG AND VISUAL FEATURES VIA MSLPCCA

Takuya Kawakami, Takahiro Ogawa and Miki Haseyama\*

Graduate School of Information Science and Technology, Hokkaido University N-14, W-9, Kita-ku, Sapporo, Hokkaido, 060-0814, Japan E-mail: {kawakami, ogawa}@lmd.ist.hokudai.ac.jp, miki@ist.hokudai.ac.jp

## ABSTRACT

This paper presents a novel image classification method based on integration of EEG and visual features. In the proposed method, we obtain classification results by separately using EEG and visual features. Furthermore, we merge the above classification results based on a kernelized version of Supervised learning from multiple experts and obtain the final classification result. In order to generate feature vectors used for the final image classification, we apply Multiset supervised locality preserving canonical correlation analysis (MSLPCCA), which is newly derived in the proposed method, to EEG and visual features. Our method realizes successful multimodal classification of images by the object categories that they contain based on MSLPCCA-based feature integration.

*Index Terms*— electroencephalogram (EEG), image classification, multimodal scheme, decision-level fusion, canonical correlation analysis.

## 1. INTRODUCTION

Image classification is an important task for image semantic analysis. Thus, various methods which classify images according to object categories that these images contain have intensively been proposed [1–3]. In order to classify images automatically, these methods utilize visual features extracted from each image. Although the classification accuracy has been improved by using several visual features [4–7], the improvement of the classification performance based on the discovery of the new visual features tends to be saturated. Therefore, it is necessary to introduce a new idea such as solving the problem by using alternative features.

In order to realize image classification based on this approach, we have proposed an image classification method [8] which utilizes both EEG features extracted from EEG signals recorded while a user stares at the images and their visual features. However, this method utilizes feature vectors generated without considering relationships between EEG features and visual features. Therefore, the performance improvement is expected by using vectors based on feature integration considering their relationships.

In this paper, we propose a novel image classification method based on integration of EEG and visual features. The proposed method consists of two stages. In the first stage, our method calculates EEG features and visual features, and classifies images based on Support vector machine (SVM) [9] by using each feature independently. Then, in the second stage, our method performs merging the above classification results, *i.e.*, decision-level fusion. In order

to perform the decision-level fusion, we employ a kernelized version of Supervised learning from multiple experts (KSLME) [8]. In this method, we utilize feature vectors generated by integration of EEG and visual features based on a new approach, i.e., the biggest contribution of this paper. Specifically, in order to integrate EEG and visual features with considering relationships between them, we newly derive Multiset supervised locality preserving canonical correlation analysis (MSLPCCA). MSLPCCA enables to apply Supervised locality preserving canonical correlation analysis [10] to more than three variables. Since the conventional method [8] utilizes EEG features and three types of visual features, it is necessary to derive MSLPCCA in order to integrate at least four types of features. Then MSLPCCA enables the feature integration with preserving the locality structure of each variable and using class labels which are generally effective for the classification problem. Consequently, successful image classification based on EEG and visual features becomes feasible by our method.

## 2. IMAGE CLASSIFICATION VIA MSLPCCA-BASED FEATURE INTEGRATION

The proposed method consists of two stages. In the first stage, we extract the EEG features from EEG signals recorded while a user stares at images, and the visual features are computed from these images. Then we perform image classification based on SVM [9] by inputting EEG and visual features into the classifiers, separately. Thus, multiple classification results are obtained for test data based on each feature. Furthermore, in the second stage, we employ the kernelized decision-level fusion approach, *i.e.*, merging the above classification results, considering their classification accuracy. In this stage, we utilize feature vectors generated by integrating EEG and visual features based on MSLPCCA which is our original method. The details of our method are described below.

## 2.1. Feature Extraction and Single Feature-based Image Classification

In this subsection, we explain the first stage of our method. Specifically, we explain the EEG features and the visual features used in the proposed method, and the single feature-based image classification method is shown.

### **EEG Feature Extraction**

First, segmentation of each channel's EEG signal is performed at fixed intervals with an overlapped Hamming window. In this paper,  $f_t$  ( $t = 1, 2, \dots, F$ ; F is the total number of EEG segments) denote EEG segments. Next, we compute the EEG features shown in Table 1 from each EEG segment. Note that C and P denote the number of channels of EEG signals and the number of symmetric electrode pairs placed on the scalp, respectively. Thus, the dimension of EEG features becomes 6C + 10P. In this table, we calculate

<sup>\*</sup>This work was partly supported by Grant-in-Aid for Scientific Research (B) 25280036 from JSPS.

**Table 1**. Features used for EEG signals in our method. Note that C denotes the number of channels of EEG signals and P shows the number of symmetric electrode pairs placed on the scalp.

The Type of EEG	Num. of Dimension	
Zero Crossing	С	
Content percentage of the power spectrum	θ wave (4-7Hz)	С
	slow-α wave (7-9Hz)	С
	mid-a wave (9-11Hz)	С
	slow-α wave (11-13Hz)	С
	β wave (13Hz-)	С
	$\theta$ wave (4-7Hz)	2P
Power spectrum of the hemispheric asymmetry [13]	slow-α wave (7-9Hz)	2P
	mid-a wave (9-11Hz)	2P
	slow-a wave (11-13Hz)	2P
	β wave (13Hz-)	2P

Zero crossing rate [11] in the time domain, and the other features are computed in the frequency domain by applying short-time Fourier transform (STFT) to each channel's EEG signal. The details of EEG features in our method are shown in [12].

## **Visual Feature Extraction**

We utilize four kinds of visual features: Scale invariant feature transform (SIFT) [4], Pyramid histogram of oriented gradients (PHOG) [6], GIST descriptor [7] and the Intensity histogram (IHIST). We calculate feature vectors:  $\boldsymbol{x}^{V_S} \in \mathbb{R}^{300}$ ),  $\boldsymbol{x}^{V_P} \in \mathbb{R}^{3400}$ ),  $\boldsymbol{x}^{V_G} \in \mathbb{R}^{512}$ ) and  $\boldsymbol{x}^{V_I} \in \mathbb{R}^{256}$ ) based on each method. Note that we obtain  $\boldsymbol{x}^{V_S}$  by applying BoF approach [14] to extracted 128dimensional SIFT descriptors. Due to the limitation of pages, we only show the overview of the visual features. The details of each visual feature are shown in [4, 6, 7].

#### Single Feature-based Image Classification

We explain the method to classify images based on each feature in the first stage. First, since relationships between "stimuli to human beings from the outside" and "which parts of the human brain are affected by these stimuli" are not well-known, we employ the feature selection in order to obtain EEG feature vectors. This means we reduce the dimension of the features shown in Table 1 to select only features useful for the classification. In order to perform the dimensionality reduction, we apply Max-relevance and min-redundancy (mRMR) feature selection algorithm [15] to the EEG features calculated from each segment and obtain an efficient feature set for the image classification. After this procedure,  $\boldsymbol{x}_{i}^{f_{i}} \in \mathbb{R}^{d^{f_{i}}}$   $(i = 1, 2, \cdots, N; N \text{ is the number of images included}$ in training data;  $d^{f_t}$  is the number of the selected features based on mRMR algorithm for EEG segment  $f_t$ ) are obtained as EEG feature vectors for each EEG segment  $f_i$  ( $t = 1, 2, \dots, F$ ). As for visual feature vectors, we directly use the vectors  $\boldsymbol{x}_i^{\text{Vs}}, \boldsymbol{x}_i^{\text{Vp}}, \boldsymbol{x}_i^{\text{Vg}}$  and  $\boldsymbol{x}_i^{\text{VI}}$ , separately.

In the first stage of our method, we employ SVM as the classifier. Although SVM is a two class classifier, image classification is generally a multi-class problem. Fortunately, since the two class classification can be easily expanded into multi-class classification based on one vs. one approach [16] or one vs. all approach [17], we focus on the improvement of the two class classification performance. We train classifiers by separately using EEG feature vectors calculated from each EEG segment and visual feature vectors. This means multiple classifiers (F + 4 classifiers) are respectively obtained based on EEG features  $x_i^{V_S}$ ,  $x_i^{V_F}$ ,  $x_i^{V_G}$ ,  $x_i^{V_1}$  by using each feature vector for training. Therefore, we can classify images based on EEG and visual features by inputting feature vectors extracted from test data into each trained classifier. Finally, F + 4 kinds of classification results are obtained.

#### 2.2. Multiple Feature-based Image Classification

In this subsection, we explain the method to obtain the final classification result in the second stage. In the proposed method, we merge F + 4 classification results obtained in the first stage based on KSLME proposed in [8] to determine the final classification result. In the second stage of our method, we newly derive the feature integration method called "MSLPCCA". We utilize this method to generate feature vectors from EEG and visual features and input them into the KLSME. In this subsection, we merge  $1, 2, \dots, F$  EEG segments' classification results and 4 classification results based on visual features. In the proposed method, we regard the F + 4 classifiers based on EEG features extracted from each EEG segment and visual features as F + 4 annotators. In order to merge multiple classification results, we focus on the classification accuracy of each annotator and assign higher weights to classification results of annotators which have higher classification accuracy. The details of the second stage are shown below.

#### 2.2.1. MSLPCCA-based Feature Integration

We explain the feature integration method called "MSLPCCA". We generate new integrated feature vectors by applying MSLPCCA to EEG features and visual features. First, we define 5 variables  $X^1, X^2, \dots, X^5$  as  $X^r = [x_1^{(1)}, x_2^{(1)}, \dots, x_{n_1}^{(l)}, x_2^{(0)}, \dots, x_{n_0}^{(l)}]$ , where  $x_i^{r(0)} \in \mathbb{R}^{d^r}$ ,  $n_1 + n_0 = N$ , and  $x_i^{r(1)}$  and  $x_i^{r(0)}$  are respectively feature vectors assigned to positive and negative classes. In the proposed method, 5 variables correspond to  $X^E, X^{V_S}, X^{V_P}, X^{V_G}$  and  $X^{V_1}$ , respectively. In particular,  $x_i^{E(0)}$  in a variable  $X^E$  are generated by  $1, 2, \dots, F$  EEG segments' features. We obtain these feature vectors by calculating the average and standard deviation of each EEG feature from  $1, 2, \dots, F$  EEG segments.

Next, we calculate similarity matrices  $S^{X^1}$ ,  $S^{X^2}$ ,  $\dots$ ,  $S^{X^5}$  in the same way as Locality preserving projection [18]. The (i, j) th component of each matrix is obtained as follows:

$$S_{ij}^{X^{r}} = \begin{cases} e^{-\left\|\boldsymbol{x}_{i}^{(k)} - \boldsymbol{x}_{j}^{(k)}\right\|^{2} / \sigma_{\boldsymbol{x}^{r}}^{2}} & \text{if } \boldsymbol{x}_{j}^{r^{(1)}} \in \boldsymbol{\Omega}_{\boldsymbol{x}_{i}^{r^{(1)}}}^{k} \text{ or } \boldsymbol{x}_{i}^{r^{(1)}} \in \boldsymbol{\Omega}_{\boldsymbol{x}_{j}^{r^{(1)}}}^{k} \\ 0 & \text{otherwise.} \end{cases}$$
(1)

In Eq. (1),  $\Omega_{x_i^{\ell}}^k$  is a set of *k* neighbors of  $x_i^{\ell^{(i)}}$ , and these neighbors are defined by the Euclidean distance. In addition,  $\sigma_{x'}^2$  is calculated by  $\sigma_{x'}^2 = \frac{\xi'}{N(N-1)}$ , where  $\xi'$  is the sum of the squared distances between all two vectors in  $\mathbf{X}^r$ . Furthermore, in the proposed method, if two vectors corresponding to (i, j) th component of  $\mathbf{S}^{X'}$  have different class labels, we replace the value of  $S_{ij}^{X'}$  as zero. This procedure enables to introduce supervised learning. After this procedure, we define similarity matrices as  $\tilde{\mathbf{S}}^{X'}$ .

We calculate weight matrices  $U^1, U^2, \dots, U^5$ , which maximize the following correlation between variables  $X^1, X^2, \dots, X^5$ , by using similarity matrices as follows:

Maximize 
$$\sum_{r=1}^{5} \sum_{s=1}^{5} \boldsymbol{U}^{r^{\top}} \boldsymbol{X}^{r} \boldsymbol{L}^{rs} \boldsymbol{X}^{s^{\top}} \boldsymbol{U}^{s},$$
  
subject to 
$$\sum_{r=1}^{5} \boldsymbol{U}^{r^{\top}} \boldsymbol{X}^{r} \boldsymbol{L}^{rr} \boldsymbol{X}^{r^{\top}} \boldsymbol{U}^{r} = 1.$$
 (2)

Weight matrices can be obtained by solving the generalized eigenvalue problem. Thus, each weight matrix is  $\hat{U}^r = [u_1^r, u_2^r, \cdots, u_{N_e}^r]$ , where  $u_l^r \in \mathbb{R}^{d^r}$  is the eigenvector, and  $N_e$  is the number of positive eigenvalues. In Eq. (2),  $L^{rr} = D^{rr} - \tilde{S}_{X^r} \circ \tilde{S}_{X}$  and  $L^{rs} = D^{rs} - \tilde{S}_{X^r} \circ \tilde{S}_{X^s}$ , respectively. Note that " $\circ$ " denotes the Hadamard product, and  $D^{rr} = \text{diag}(\sum_i (\tilde{S}_{i1}^{tr})^2, \sum_i (\tilde{S}_{2i}^{tr})^2, \cdots, \sum_i (\tilde{S}_{N_i}^{tr})^2)$  and

 $D^{rs} = \operatorname{diag}(\sum_{i} \tilde{S}_{1i}^{X^{r}} \tilde{S}_{1i}^{X^{s}}, \sum_{i} \tilde{S}_{2i}^{X^{r}} \tilde{S}_{2i}^{X^{s}}, \cdots, \sum_{i} \tilde{S}_{Ni}^{X^{r}} \tilde{S}_{Ni}^{X^{s}}),$  respectively.

By using generated weight matrices  $\hat{U}^1, \hat{U}^2, \dots, \hat{U}^5$ , we obtain feature vectors based on integration of EEG and visual features as follows:

$$\boldsymbol{x}_{i}^{\text{EV}} = \left[\boldsymbol{x}_{i}^{1^{(\circ)^{\top}}} \hat{\boldsymbol{U}}^{1}, \boldsymbol{x}_{i}^{2^{(\circ)^{\top}}} \hat{\boldsymbol{U}}^{2}, \cdots, \boldsymbol{x}_{i}^{5^{(\circ)^{\top}}} \hat{\boldsymbol{U}}^{5}\right]^{\top}.$$
(3)

MSLPCCA considers the class label of each vector and realizes feature integration with preserving a locality structure of each variable in new feature space. Therefore, we generate feature vectors, which are effective for image classification using both EEG and visual features, in Eq. (3) based on MSLPCCA.

# 2.2.2. Each annotator's classification accuracy and classification model

We explain the classification accuracy of each annotator and the classification model defined in our method. Let  $y^a \in \{0, 1\}$  be the label assigned to the feature vector  $\boldsymbol{x}^{\text{EV}}$  by annotator  $a \in \mathcal{A}$ , where  $\mathcal{A} = \{f_1, f_2, \cdots, f_r, V_S, V_P, V_G, V_I\}$  is a set of annotators, and f and V correspond to Frame (EEG segment) and Visual, respectively. Given the actual label  $y \in \{0, 1\}$ , *i.e.*, ground truth, the classification accuracy of each annotator,  $P_{se}^a$  (sensitivity) and  $P_{sp}^a$  (specificity) are respectively defined as follows:

$$P_{se}^{a} := \Pr[y^{a} = 1|y = 1], \tag{4}$$

$$P_{sp}^{a} := \Pr[y^{a} = 0|y = 0].$$
(5)

In our method, a classification model is specifically written as:

$$f_w(\boldsymbol{x}^{\rm EV}) = \boldsymbol{w}^{\top} \boldsymbol{\phi}(\boldsymbol{x}^{\rm EV}), \qquad (6)$$

where w is a weight.

In Eq. (6),  $\phi(\boldsymbol{x}^{\text{EV}})$  is obtained by mapping the feature vector  $\boldsymbol{x}^{\text{EV}}$  into a high-dimensional feature space. The final classification result  $\hat{\boldsymbol{y}}$  is obtained as follows:

$$\hat{y} = \begin{cases} 1 & f_w(\boldsymbol{x}^{\text{EV}}) \ge \tau \\ 0 & \text{otherwise,} \end{cases}$$
(7)

where  $\tau$  is a predetermined threshold. Given the training dataset  $\mathcal{D}$  consisting of N feature vectors  $\boldsymbol{x}_i^{\text{EV}}$   $(i = 1, 2, \dots, N)$ , a weight  $\boldsymbol{w}$  is specifically written as follows:

$$\boldsymbol{w} = \sum_{i=1}^{N} \alpha_i \boldsymbol{\phi}(\boldsymbol{x}_i^{\text{EV}}), \qquad (8)$$

where  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \cdots, \alpha_N]^{\mathsf{T}}$ . Therefore, by using  $\boldsymbol{\alpha}$  in Eq. (8), the discriminating function in Eq. (6) is rewritten as follows:  $f_w(\boldsymbol{x}^{\mathrm{EV}}) = \boldsymbol{w}^{\mathsf{T}} \boldsymbol{\phi}(\boldsymbol{x}^{\mathrm{EV}})$ 

$$\boldsymbol{\varphi}(\boldsymbol{x}^{\mathrm{EV}}) = \boldsymbol{w}^{\mathrm{EV}} \boldsymbol{\varphi}(\boldsymbol{x}^{\mathrm{EV}})$$
$$= \sum_{i=1}^{N} \alpha_{i} K \left( \boldsymbol{x}_{i}^{\mathrm{EV}}, \boldsymbol{x}^{\mathrm{EV}} \right), \tag{9}$$

where  $K(\cdot, \cdot)$  is a kernel function of  $\phi(\cdot)$ , and we specifically employ the Gaussian kernel. In order to determine the discriminating function  $f_{w}(\cdot)$ , we have to obtain the coefficients  $\alpha_i$   $(i = 1, 2, \dots, N)$  from training data by using each annotator's classification accuracy defined in Eqs. (4) and (5). The details are shown below.

#### 2.2.3. Training Phase

Given the training data  $\mathcal{D}$  consisting of N feature vectors with the classification results by F + 4 annotators and their actual labels,  $\mathcal{D} = \{y_i, \phi(\boldsymbol{x}_i^{EV}), \mathcal{Y}_i\}_{i=1}^N$ , where  $y_i$  is the actual label and  $\mathcal{Y}_i = \{y_i^{f_1}, y_i^{f_2}, \cdots, y_i^{f_r}, y_i^{V_s}, y_i^{V_b}, y_i^{V_c}, y_i^{V_l}\}$  is a set of classification results, the estimation target is the coefficients  $\alpha_i$  ( $i = 1, 2, \cdots, N$ ) in Eq. (9). From the training data  $\mathcal{D}$ , the likelihood of the coefficient vector  $\boldsymbol{\alpha}$  is defined as:

$$\Pr\left[\mathcal{D}|\boldsymbol{\alpha}\right] = \prod_{i=1}^{N} \Pr\left[\mathcal{Y}_{i}|\boldsymbol{\phi}(\boldsymbol{x}_{i}^{\text{EV}}),\boldsymbol{\alpha}\right].$$
(10)

By using a set of sensitivity  $\mathcal{P}_{se} = \{P_{se}^{f_1}, P_{se}^{f_2}, \cdots, P_{se}^{f_r}, P_{se}^{V_S}, P_{se}^{V_g}, P_{se}^{V_g}, P_{se}^{V_g}, P_{se}^{V_g}\}$  obtained from each annotator and that of specificity  $\mathcal{P}_{sp} = \{P_{sp}^{f_1}, P_{sp}^{f_2}, \cdots, P_{sp}^{f_r}, P_{sp}^{V_S}, P_{sp}^{V_g}, P_{sp}^{V_g}, P_{sp}^{V_g}\}$ , the above equation is rewritten as follows:

$$\Pr[\mathcal{D}|\boldsymbol{\alpha}] = \prod_{i=1}^{N} \left\{ \Pr[\mathcal{Y}_{i}|y_{i}=1,\mathcal{P}_{se}] \times \Pr[y_{i}=1|\boldsymbol{\phi}(\boldsymbol{x}_{i}^{\mathrm{EV}}),\boldsymbol{\alpha}] + \Pr[\mathcal{Y}_{i}|y_{i}=0,\mathcal{P}_{sp}] \times \Pr[y_{i}=0|\boldsymbol{\phi}(\boldsymbol{x}_{i}^{\mathrm{EV}}),\boldsymbol{\alpha}] \right\}.$$
(11)

If it is assumed that each annotator  $a \in \mathcal{A}$  is independent each other,  $\Pr[\mathcal{Y}_i|y_i = 1, \mathcal{P}_{se}]$  can be rewritten as follows:

$$\begin{bmatrix} \mathbf{\mathcal{Y}}_i | y_i = 1, \mathcal{P}_{se} \end{bmatrix} = \prod_{a \in \mathcal{A}} \Pr[y_i^a | y_i = 1, P_{se}^a]$$
$$= \prod_{a \in \mathcal{A}} [P_{se}^a]^{y_i^a} [1 - P_{se}^a]^{1 - y_i^a}.$$
(12)

Similarly,  $\Pr[\mathcal{Y}_i|y_i = 0, \mathcal{P}_{sp}]$  can be rewritten as follows:

$$\Pr[\mathcal{Y}_i|y_i = 0, \mathcal{P}_{sp}] = \prod_{a \in \mathcal{A}} \Pr[y_i^a|y_i = 0, P_{sp}^a]$$
$$= \prod_{a \in \mathcal{A}} [P_{sp}^a]^{1-y_i^a} [1 - P_{sp}^a]^{y_i^a}.$$
(13)

Then the likelihood in Eq. (11) is rewritten as

Pr

$$\Pr[\mathcal{D}|\boldsymbol{\alpha}] = \prod_{i=1}^{N} [\gamma_i \rho_i + \delta_i (1 - \rho_i)], \qquad (14)$$

where

and

$$\gamma_{i} = \prod_{a \in \mathcal{A}} [P_{se}^{a}]^{y_{i}^{a}} [1 - P_{se}^{a}]^{1 - y_{i}^{a}},$$
(15)

$$\delta_{i} = \prod_{a \in \mathcal{A}} [P^{a}_{sp}]^{1-y^{a}_{i}} [1 - P^{a}_{sp}]^{y^{a}_{i}},$$
(16)

$$\rho_i = \Pr[y_i = 1 | \phi(\boldsymbol{x}_i^{\text{EV}}), \boldsymbol{\alpha}]$$
$$= \frac{1}{1 + \exp(-\boldsymbol{\alpha}^\top \boldsymbol{\kappa}_i)}, \qquad (17)$$

where  $\kappa_i = \left[ K(\boldsymbol{x}_1^{\text{EV}}, \boldsymbol{x}_i^{\text{EV}}), K(\boldsymbol{x}_2^{\text{EV}}, \boldsymbol{x}_i^{\text{EV}}), \cdots, K(\boldsymbol{x}_N^{\text{EV}}, \boldsymbol{x}_i^{\text{EV}}) \right]^{\prime}$ . The maximum-likelihood estimator is found by maximizing the following log-likelihood:

$$\hat{\boldsymbol{\alpha}}_{ML} = \arg \max \{ \ln \Pr[\mathcal{D}|\boldsymbol{\alpha}] \}.$$
(18)

Let  $\mathcal{L} = \{y_1, y_2, \dots, y_N\}$  be the set of the actual labels, and the complete data log-likelihood can be written as

$$\ln \Pr[\mathcal{D}, \mathcal{L}|\boldsymbol{\alpha}] = \sum_{i=1}^{N} \{ y_i \ln \rho_i \gamma_i + (1 - y_i) \ln(1 - \rho_i) \delta_i \}.$$
 (19)

In order to maximize this likelihood, the following Expectation-Maximization (EM) algorithm [19] is adopted. E-step

In the E-step, when the training data  $\mathcal{D}$  and the current estimate of the coefficient vector  $\boldsymbol{\alpha}$  are given, the conditional expected value of log-likelihood is computed as follows:

$$\mathbf{E}\left\{\ln\Pr[\mathcal{D},\mathcal{L}|\boldsymbol{\alpha}]\right\} = \sum_{i=1}^{N} \left\{\mu_{i}\ln\rho_{i}\gamma_{i} + (1-\mu_{i})\ln(1-\rho_{i})\delta_{i}\right\},$$
(20)

where  $\mu_i$  is computed as follows:

$$\mu_{i} \propto \Pr[\mathcal{Y}_{i}|y_{i} = 1, \alpha] \times \Pr[y_{i} = 1|\phi(\boldsymbol{x}_{i}^{\mathrm{EV}}), \alpha]$$
$$= \frac{\gamma_{i}\rho_{i}}{\gamma_{i}\rho_{i} + \delta_{i}(1-\rho_{i})}.$$
(21)

	Conventional Method [8]		Only Visu	al Features	Only EEG Features	Proposed Method	
		SIFT	PHOG	GIST	IHIST		
subject A	$79 \pm 0.10\%$					$59 \pm 0.30\%$	$82\pm0.06\%$
subject B	$87\pm0.05\%$					$77 \pm 0.05\%$	$87\pm0.05\%$
subject C	$88\pm0.06\%$	$68\pm0.07\%$	$67 \pm 0.07\%$	$67 \pm 0.09\%$	$71\pm0.01\%$	$72 \pm 0.09\%$	$88\pm0.02\%$
subject D	$94\pm0.05\%$					$69 \pm 0.02\%$	$93 \pm 0.04\%$
subject E	$93 \pm 0.04\%$					$71 \pm 0.08\%$	$95\pm0.02\%$

Table 2. Image Classification Accuracy: These values are the average and standard deviation over all target image categories.

#### M-step

In the M-step, the coefficient vector  $\boldsymbol{\alpha}$  is estimated based on the current estimate  $\mu_i$  and the training data  $\mathcal{D}$  by maximizing the conditional expected value in Eq. (20). Specifically, by solving equation  $\frac{\partial}{\partial \alpha} \{\ln \Pr[\mathcal{D}, \mathcal{L} | \boldsymbol{\alpha}]\} = 0$ , we obtain the estimated coefficient vector  $\boldsymbol{\alpha}$  as follows:

$$\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} - \eta \boldsymbol{H}^{-1} \boldsymbol{g}. \tag{22}$$

In Eq. (22),  $\boldsymbol{g}$  is a gradient vector,  $\boldsymbol{H}$  is a Hessian matrix and  $\eta$  is a step length. The gradient vector  $\boldsymbol{g}$  and the Hessian matrix  $\boldsymbol{H}$  are respectively computed as follows:  $\boldsymbol{g} = \sum_{i=1}^{N} [\mu_i - \sigma(\alpha^\top \kappa_i)] \boldsymbol{\kappa}_i$  and  $\boldsymbol{H} = -\sum_{i=1}^{N} [\sigma(\alpha^\top \kappa_i)] [1 - \sigma(\alpha^\top \kappa_i)] \boldsymbol{\kappa}_i \boldsymbol{\kappa}_i^\top$ , where  $\sigma(\alpha^\top \kappa_i) = \frac{1}{1 + \exp(-\alpha^\top \kappa_i)}$ .

### 2.2.4. Testing Phase

Given the test data, the final classification result is obtained as follows. In the previous phase, we essentially solved a regular logistic regression problem with probabilistic labels  $\mu_i$ . Thus, we obtain the final classification result by applying a threshold to  $\mu$  calculated from a test data { $\phi(\mathbf{x}^{\text{EV}}), y^{f_1}, y^{f_2}, \cdots, y^{f_r}, y^{\text{Vs}}, y^{\text{Vg}}, y^{\text{Vg}}, y^{\text{Vl}}$ }, where its label is unknown, instead of directly using  $\hat{\alpha}_{ML}$ . The value of  $\mu$  is computed by using the estimated coefficient vector  $\hat{\alpha}_{ML}$  and  $\gamma$ ,  $\delta$  calculated from the training data. Specifically,  $\rho = \frac{1}{1 + \exp(-\hat{\alpha}_{ML}^{\text{T}}\kappa)}$  is calculated, where  $\kappa = \left[K(\mathbf{x}_1^{\text{EV}}, \mathbf{x}^{\text{EV}}), K(\mathbf{x}_2^{\text{EV}}, \mathbf{x}^{\text{EV}}), \cdots, K(\mathbf{x}_N^{\text{EV}}, \mathbf{x}^{\text{EV}})\right]^{\text{T}}$ . Then,  $\gamma = \prod_{a \in \mathcal{A}} [P_{ae}^a]^{1-p^a} [1 - P_{ae}^a]^{1-p^a}$  and  $\delta = \prod_{a \in \mathcal{A}} [P_{ap}^a]^{1-p^a} [1 - P_{ap}^a]^{a}$  are obtained, where  $P_{se}^a$  and  $P_{sp}^a$  are accuracy of annotator *a* calculated from training data and  $\gamma^a$  is classification result of the test data. Therefore, we obtain the final classification result of by using  $\rho$ ,  $\gamma$  and  $\delta$ . Finally, we obtain the final classification result as follows:

$$\hat{y} = \begin{cases} 1 & \mu \ge \mu' \\ 0 & \text{otherwise,} \end{cases}$$
(23)

where  $\mu'$  is a predetermined threshold. The value of  $\mu$  is the posterior probability.

## 3. EXPERIMENTAL RESULTS

In this section, we show experimental results to verify the effectiveness of the proposed method. First, we explain image dataset and EEG signal collection. In this experiment, we utilized Caltech101 dataset [20]. Specifically, we used the images included in "panda", "soccer ball" and "strawberry" in the database for image classification, and the number of images was 35 per category. These 105 images were randomly selected in advance. In this paper, we define images used for image classification as target images. We also used images included in "airplane", "elephant", "joshua tree", "pyramid" and "stapler" in the same database for the non-target images.

Next, we explain how to collect EEG signals in this experiment. In this study, five healthy subjects (subject A, B, C, D and E) participated, and EEG recordings were conducted during staring at images. The age of each subject was 22 or 23 years old. We recorded EEG signals from 12 channels (Fp1, Fp2, F7, F8, T3, T4, C3, C4, P3, P4, O1 and O2) according to the international 10-20 system. Since EEG signals are weak, we amplified these signals by using an amplifier (MEG-6116M, NIHON KOHDEN). We also applied a bandpass filter to recorded EEG signals to avoid artifacts, and set the filter bandwidth to 0.04-30Hz. In addition, the EEG signals were sampled at 2kHz. In this experiment, we collected single-trial EEG signals for each target image by the same experimental procedure as those shown in our previous work [8]. Note that the time length of staring at each image was two seconds (subject A and B) and three seconds (subject C, D and E). In this experiment, we perform image classification by using first one second of EEG signals recorded while a subject stared at target images.

Furthermore, we present the experimental condition. In this experiment, we performed the three class image classification (panda, soccer ball and strawberry) by the object categories that images contain based on one vs. all approach [17]. Therefore, the final classification was determined according to the posterior probability obtained from the testing phase (2.2.4). We followed [2, 3] for our experimental setup. Specifically, we randomly selected 30 training images per class and tested on the remaining images. Then we calculated the classification accuracy which was normalized according to the number of test images per class. We repeated the random selection 10 times and show the average classification accuracy over all classes.

We show the results of image classification in Table 2. As for the proposed method, we vary the number of neighbors k in Eq. (1) and present the best results. In this table, we also show the results of the conventional method [8]. From the obtained results, the proposed method realizes more accurate classification than the conventional method [8]. Therefore, the effectiveness of our method can be verified. Thus, the feature integration based on MSLPCCA is effective for the image classification using both EEG and visual features. In Table 2, we also present the average image classification accuracy by using either EEG or visual features based on SVM. From this table, although accuracy obtained by using each feature separately is not satisfactory, our method realizes successful classification based on collaborative use of all features.

## 4. CONCLUSION

In this paper, we have proposed a novel image classification method based on MSLPCCA-based feature integration. MSLPCCA enables the feature integration with preserving the locality structure of each variable and using class labels which are generally effective for the classification problem. In our method, this method is applied to EEG and visual features in order to generate feature vectors for the decision-level fusion of the image classification. Experimental results show the effectiveness of the proposed method.

#### 5. REFERENCES

- M. Haseyama, T. Ogawa, and N. Yagi, "A review of video retrieval based on image and video semantic understanding," *ITE Transactions on Media Technology and Applications*, vol. 1, no. 1, pp. 2–9, 2013.
- [2] J. Zhang, M. Marszałek, S. Lazebnik, and C. Schmid, "Local features and kernels for classification of texture and object categories: A comprehensive study," *International journal of computer vision*, vol. 73, no. 2, pp. 213–238, 2007.
- [3] K. Grauman and T. Darrell, "The pyramid match kernel: Discriminative classification with sets of image features," in in*Proceedings of IEEE International Conference on Computer Vision (ICCV)*, 2005, vol. 2, pp. 1458–1465.
- [4] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [5] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (*CVPR*), 2005, vol. 1, pp. 886–893.
- [6] Anna Bosch, Andrew Zisserman, and Xavier Munoz, "Representing shape with a spatial pyramid kernel," in *Proceedings* of the 6th ACM international conference on Image and video retrieval, 2007, pp. 401–408.
- [7] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *International journal of computer vision*, vol. 42, no. 3, pp. 145–175, 2001.
- [8] T. Kawakami, T. Ogawa, and M. Haseyama, "Novel image classification based on decision-level fusion of EEG and visual features," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5915– 5919, 2014.
- [9] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, pp. 273 –297, 1995.
- [10] Jinfeng Yang and Xu Zhang, "Feature-level fusion of fingerprint and finger-vein for personal identification," *Pattern Recognition Letters*, vol. 33, no. 5, pp. 623–628, 2012.
- [11] Alexander A. Borbely and HansUlrich Neuhaus, "Sleepdeprivation: Effects on sleep and EEG in the rat," *Journal* of comparative physiology, vol. 133, no. 1, pp. 71–87, 1979.
- [12] T. Kawakami, T. Ogawa, and M. Haseyama, "Vocal segment estimation in music pieces based on collaborative use of EEG and audio features," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing* (ICASSP), 2013, pp. 1197–1201.
- [13] Y. P. Lin, C. H. Wang, T. P. Jung, T. L. Wu, S. K. Jeng, J. R. Duann, and J. H. Chen, "EEG-based emotion recognition in music listening," *IEEE Transactions on Biomedical Engineering*, vol. 57, no. 7, pp. 1798–1806, 2010.
- [14] G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray, "Visual categorization with bags of keypoints," in *Proceedings* of European Conference on Computer Vision, 2004, vol. 1, pp. 1–22.
- [15] H. Peng, F. Long, and C. Ding, "Feature selection based on mutual information criteria of max-dependency, maxrelevance, and min-redundancy," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 8, pp. 1226 –1238, 2005.

- [16] T. Hastie and R. Tibshirani, "Classification by pairwise coupling," *The annals of statistics*, vol. 26, no. 2, pp. 451–471, 1998.
- [17] R. Rifkin and A. Klautau, "In defense of one-vs-all classification," *The Journal of Machine Learning Research*, vol. 5, pp. 101–141, 2004.
- [18] X. He and P. Niyogi, "Locality preserving projections," in *NIPS*, 2003, vol. 16, pp. 234–241.
- [19] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Jour*nal of the Royal Statistical Society, pp. 1–38, 1977.
- [20] L. Fei-Fei, R. Fergus, and P. Perona, "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 2004, pp. 178–178.