

PAPER

Mean Polynomial Kernel and Its Application to Vector Sequence Recognition

Raissa RELATOR^{†a)}, Yoshihiro HIROHASHI^{†*}, Eisuke ITO[†], *Nonmembers*, and Tsuyoshi KATO^{†b)}, *Member*

SUMMARY Classification tasks in computer vision and brain-computer interface research have presented several applications such as biometrics and cognitive training. However, like in any other discipline, determining suitable representation of data has been challenging, and recent approaches have deviated from the familiar form of one vector for each data sample. This paper considers a kernel between vector sets, the mean polynomial kernel, motivated by recent studies where data are approximated by linear subspaces, in particular, methods that were formulated on Grassmann manifolds. This kernel takes a more general approach given that it can also support input data that can be modeled as a vector sequence, and not necessarily requiring it to be a linear subspace. We discuss how the kernel can be associated with the Projection kernel, a Grassmann kernel. Experimental results using face image sequences and physiological signal data show that the mean polynomial kernel surpasses existing subspace-based methods on Grassmann manifolds in terms of predictive performance and efficiency.

key words: *kernel methods, support vector machines, Grassmann distance and kernels, face recognition, brain-computer interface, vector sequence, binary classification*

1. Introduction

Among currently trending fields, research efforts particularly related to computer vision and brain-computer interface (BCI) have been aimed at modeling data either as a low-dimensional subspace or a sequence of vectors. There have been studies in these areas dedicated to algorithms for such type of input [1]–[9]. For computer vision, this approach may be motivated by the presence of abundant material derived from videos and sets of image sequences [1]–[7] such as in Fig. 1 (a). Each video image extracted is represented by a vector, while the whole vector sequence, concatenated as a matrix, approximates the video for a given time frame. As for BCI adopting a similar approach, this may be induced by the nature of the data, which is commonly time series, such as electroencephalography (EEG) signals collected while subjects perform motor tasks or during induction of visual stimuli [8], [9]. EEG data is generated by placing several sensors accordingly on the head of the subject as shown in Fig. 2, and each sensor records neural activity depicted by the signals. The vector sequence

illustrated in Fig. 1 (b) corresponds to the signals collected from all sensors.

Appropriate data representation has been considered as one of the most important challenges in dealing with classification tasks. Vector form may be the simplest and most common representation of samples in existing literatures, especially when using popular techniques such as support vector machines (SVM) and kernel methods. However, this may not be the best representation to encompass significant, if not all, attributes and information useful for discrimination. To address this problem, new modes of data representation are constantly being explored [5], [10]–[24]. Along with this, various feature extraction techniques and discrimination methods are also being investigated, and several studies have proven kernel methods to be a flexible technique in supporting various data structures, such as graphs [10]–[15], strings [16]–[19], and even subspaces and sets of vectors [5], [20]–[24].

Kernel-based algorithms [25], [26] have come a long way since their introduction. Aside from the fact that kernel functions have provided algorithms a bridge between linearity and nonlinearity, their performance have been proven comparable to, if not better than, existing algorithms in various areas where they have been applied. Moreover, applying the so-called ‘kernel trick’ is very straightforward and new kernels can be easily derived from old kernels. Compared to other methods, the dimension of the feature space can also be treated more lightly since the technique can be reduced to simply performing inner products between data images on the space, thus making the algorithm computationally inexpensive.

In this paper, we focus on data represented as sets of vectors. Different algorithms have been formulated in such a way that data are approximated by low-dimensional linear subspaces [1]–[7], [22]. However, as previously pointed out [2], the task of appropriately handling data has become an issue, such as inconsistency in strategy when feature extraction is done in a Euclidean space while non-Euclidean metrics are used. For this purpose, they proposed a unified framework for subspace-based approaches by formulating the problem on the Grassmann manifold, a space of linear subspaces with a fixed dimension. On the other hand, these methods involve dimension reduction, and even with the use of the usual dimension reduction techniques such as Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA), there is always a possibility of information loss. This makes the selection of the subspace di-

Manuscript received November 6, 2013.

Manuscript revised April 1, 2014.

[†]The authors are with the Department of Computer Science, Graduate School of Science and Engineering, Gunma University, Kiryu-shi, 376–8515 Japan.

^{*}Presently, with the Electrical Energy Systems, Graduate School of Engineering, Tohoku University, Sendai-shi, 980–8579 Japan.

a) E-mail: relator-raissa@kato-lab.cs.gunma-u.ac.jp

b) E-mail: katotsu@cs.gunma-u.ac.jp

DOI: 10.1587/transinf.E97.D.1855

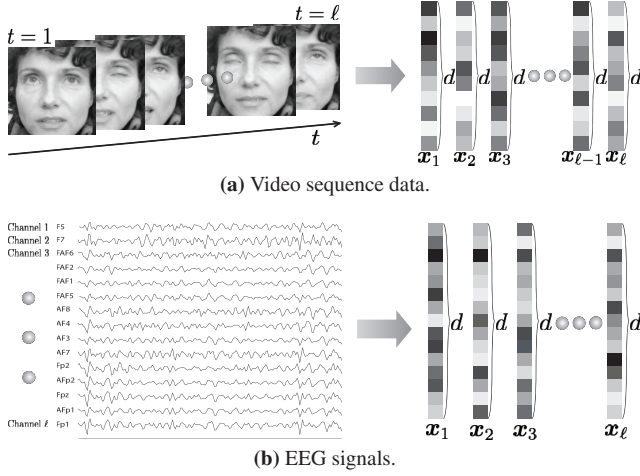


Fig. 1 Examples of data modeled as vector sequences. (a) For video sequences, each image frame extracted is represented as a vector of pixel intensity values. The vector sequence is usually concatenated to represent the vector set input. (b) For BCI, EEG signals are recorded over a certain time interval using several channels or sensors. Each vector in the sequence corresponds to a channel used in the procedure, and vector entry represents an instantaneous signal intensity.

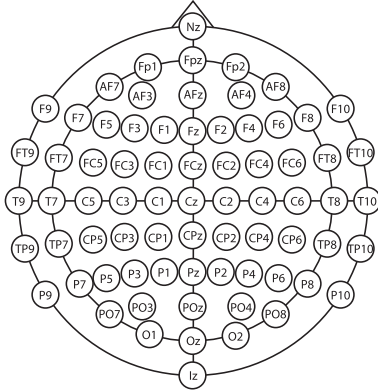


Fig. 2 Sample position map of sensors for EEG.

mension a crucial step. Furthermore, methods such as PCA and LDA usually employ eigendecomposition, and hence, may be very time consuming especially for high dimensional data.

With the aforementioned issues in mind, the goal of this paper is to examine a kernel function, which we refer to as the mean polynomial kernel, that can retain data information while being computationally inexpensive. Also, as a more general approach than kernels for subspaces, we treat data as a common collection of vectors, instead of a linear subspace. The kernel is invariant of the permutation order of the vectors in the set. In addition, we present an interesting relationship between this kernel and the Projection kernel, which is a known Grassmann kernel. We give emphasis to face recognition and BCI applications posed as binary classification problems, which are of particular interest due to their practicality in various areas, biometrics and cognitive training and improvement, among others. Experimental re-

sults using real data modeled as vector sequences show that, aside from being computationally efficient, the performance of the mean polynomial kernel is comparable to methods employing kernels in the Grassmann manifold and subspace methods using Grassmann distances.

2. Preliminaries

Consider a set of data $\mathbf{x}_1, \dots, \mathbf{x}_\ell \in \mathbb{R}^d$, where ℓ is the number of data points. Let us denote the h th entry in the i th data point \mathbf{x}_i by $x_{h,i}$. A sample statistic,

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \prod_{h=1}^d x_{h,i}^{p_h},$$

is said to be the q th order moment if the d -dimensional vector $\mathbf{p} \in (\mathbb{N} \cup \{0\})^d$ satisfies $p_1 + \dots + p_d = q$. The *uncentered covariance matrix* defined by

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \mathbf{x}_i \mathbf{x}_i^{\top}$$

contains all the second order moments. Indeed, the (h, k) th entry in the uncentered covariance matrix is the second order moment with $\mathbf{p} = \mathbf{e}_h + \mathbf{e}_k$, where \mathbf{e}_h is a unit vector whose h th entry is one and the rest of the entries are zero. Let $\bar{\mathbf{x}} = [\bar{x}_1, \dots, \bar{x}_d]^{\top}$ be the mean vector of the data points. With the d -dimensional vector \mathbf{p} satisfying $p_1 + \dots + p_d = q$, the q th order central moment is defined as

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \prod_{h=1}^d (x_{h,i} - \bar{x}_h)^{p_h}.$$

Every second order central moment is included in the *central covariance matrix*

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^{\top},$$

which is usually referred to simply as the *covariance matrix*.

For succeeding sections, we refer to a matrix \mathbf{U} as *orthonormal* if $\mathbf{U}^{\top} \mathbf{U} = \mathbf{I}$, and define the *vectorization* of an $m \times n$ matrix \mathbf{A} as the column vector $\text{vec}(\mathbf{A}) = [a_{11}, a_{12}, \dots, a_{1n}, \dots, a_{m1}, a_{m2}, \dots, a_{mn}]^{\top}$.

3. Grassmann Kernels and Related Methods

We give a concise discussion of the Grassmann kernels [2], [22], [27], their analogy with the mean polynomial kernel, and some related methods.

A Grassmann manifold, or Grassmannian, is defined as a set of linear subspaces with a fixed number of dimensions, say, m . Several metrics used in literatures have been specified in this manifold, mostly incorporating principal angles or angles between subspaces in their characterization [1], [2], [4]–[7], [22], [27]. Moreover, kernels over these manifolds have also been introduced. In particular, we are interested in the following kernels:

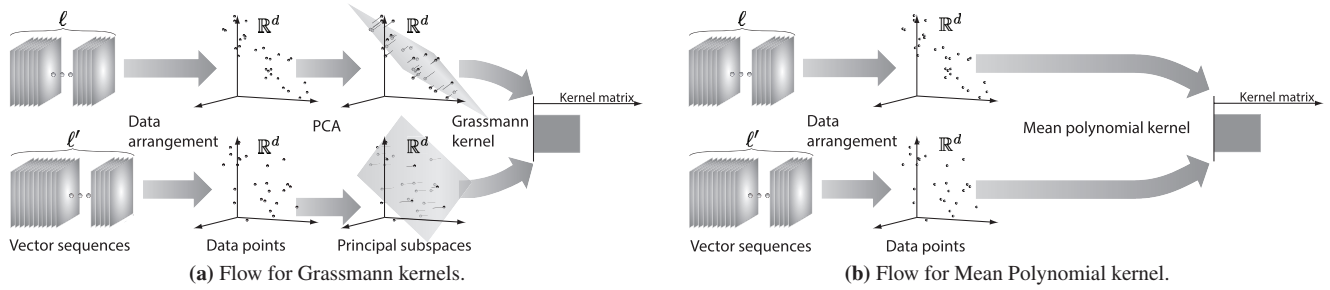


Fig. 3 Flow of methodology for computing values for Grassmann kernels and the mean polynomial kernel. Grassmann kernels are defined on a Grassmann manifold which is a set of linear subspaces. When employing these kernels, each vector sequence, represented by a set of data points on space, is approximated by a principal subspace obtained via PCA. However, this poses a threat of some degree of information loss, and is more likely to consume more time due to eigendecomposition. The mean polynomial kernel, on the other hand, can be directly applied to compute the kernel value between the sets of data points. It can avoid information loss while being more time efficient.

Definition 1. Let U_x and U_y be orthonormal matrices whose columns are bases of linear subspaces. The Projection kernel is defined as

$$k_{\text{PROJ}}(U_x, U_y) = \|U_x^T U_y\|_F^2,$$

where $\|\cdot\|_F$ denotes the Frobenius norm, and the Binet-Cauchy kernel is given by

$$k_{\text{BC}}(U_x, U_y) = (\det U_x^T U_y)^2 = \det U_x^T U_y U_y^T U_x.$$

Many existing problems can be realized on nonlinear manifolds such as the Grassmannian. This being said, various methods in the Grassmannian setting have been proposed. One such technique is the use of Grassmann kernels in conjunction with support vector machines (GK-SVM) [5]. This approach entails the computation of kernel matrices, which then proceed as the SVM input. Analogously, the mean polynomial kernel given in Sect. 4 is applied in this manner when SVM is the classifier. Figure 3 gives a general illustration of the flow of computation of the Grassmann kernels and the mean polynomial kernel, and also highlights the difference between the two kernels.

Another comparable method is the Grassmann Distance Mutual Subspace Method (GD-MSM) [5]. This technique integrates the Grassmann metrics in the Mutual Subspace Method (MSM) [7]. Furthermore, the task of subspace classification can be approached in two ways. The first one, which is referred to as the *subject-wise dictionary*, is done by assuming that one subject or object corresponds to one principal subspace. During the training stage, the total of principal subspaces calculated is the same as the number of subjects. These serve as the bases to which the unlabeled principal subspaces of test subjects are compared to, and the subspace with the minimal Grassmann distance from the unlabeled subspace is determined. The second approach is done by assuming one principal subspace per class. The principal subspaces, which in this case is referred to as the *class-wise dictionary*, are derived from each class among the training data. This being said, we have only two principal subspaces in the case of binary classification, regardless of

the number of subjects. In the testing stage, unlabeled principal subspaces are classified according to which subspace they are closer to in terms of metric.

The score function can be considered for the two aforementioned mutual subspace methods. The SVM score can serve as a confidence level. Namely, a higher score may provide higher certainty of assigning the data to the positive class. For the *class-wise dictionary*, the difference between the distance to the subspace of the negative class, d_- , and the distance to the subspace of the positive class, d_+ , represents how confidently unknown labels are classified as positive. Hence, we define the score function as $d_- - d_+$. For the *subject-wise dictionary*, we define the score function by the difference between the minimal distance to negative class subspaces and the minimal distance to the positive class subspaces.

4. Mean Polynomial Kernel

In this section, we discuss the details of the *mean polynomial kernel*, which can be directly applied to data in the form of vector sets.

Consider two sets of vectors $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^{\ell}$ and $\mathcal{Y} = \{\mathbf{y}_j\}_{j=1}^{\ell'}$, where $\mathbf{x}_i, \mathbf{y}_j \in \mathbb{R}^d$. To define a kernel for such types of data, we introduce a notation of a set of vector sequences as $\mathcal{S} = \{\{\mathbf{z}_i\}_{i=1}^n \mid n \in \mathbb{N} \text{ and } \forall i \in \mathbb{N}_n, \mathbf{z}_i \in \mathbb{R}^d\}$, where \mathbb{N} is the set of natural numbers, and $\mathbb{N}_n = \{i \in \mathbb{N} \mid i \leq n\}$, such that \mathcal{S} is the input domain for the kernel defined as follows.

Definition 2. Let $k_q : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ such that

$$k_q(\mathcal{X}, \mathcal{Y}) = \frac{1}{\ell \ell'} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell'} \langle \mathbf{x}_i, \mathbf{y}_j \rangle^q,$$

where $\mathcal{X}, \mathcal{Y} \in \mathcal{S}$ and $q \in \mathbb{N}$. We shall refer to k_q as the q th order mean polynomial kernel.

It can be shown that this kernel is a special case of the multi-instance kernels [28] when instances involve linear kernels or polynomial kernels with constant $c = 0$. With regards to its characterization, we can easily confirm that

for the case $q = 2$, the covariance matrix is directly used as a feature vector. For instance, consider two matrices, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell]$ and $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{\ell'}]$, for the set of vectors \mathcal{X} and \mathcal{Y} , respectively. Then their respective uncentered covariance matrices are given by $\Sigma_x = \frac{1}{\ell} \sum_{i=1}^{\ell} \mathbf{x}_i \mathbf{x}_i^\top$ and $\Sigma_y = \frac{1}{\ell'} \sum_{j=1}^{\ell'} \mathbf{y}_j \mathbf{y}_j^\top$. By defining a feature map $\phi(\mathbf{X}) = \text{vec}(\Sigma_x)$, we have

$$\begin{aligned} \langle \phi(\mathbf{X}), \phi(\mathbf{Y}) \rangle &= \langle \text{vec}(\Sigma_x), \text{vec}(\Sigma_y) \rangle = \text{tr}(\Sigma_x \Sigma_y) \\ &= \frac{1}{\ell \ell'} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell'} \text{tr}(\mathbf{x}_i \mathbf{x}_i^\top \mathbf{y}_j \mathbf{y}_j^\top) = \frac{1}{\ell \ell'} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell'} \text{tr}(\mathbf{y}_j^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{y}_j) \\ &= \frac{1}{\ell \ell'} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell'} \langle \mathbf{x}_i, \mathbf{y}_j \rangle^2. \end{aligned} \quad (1)$$

Hence, the Euclidean inner product of vectorized covariance matrices is precisely the second order mean polynomial. Furthermore, all information contained within the uncentered matrices are preserved and can be exploited.

If we rewrite the definition of the kernel as

$$\bar{k}_q(\mathbf{X}, \mathbf{Y}) = \frac{1}{\ell \ell'} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell'} \langle \mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{y}_j - \bar{\mathbf{y}} \rangle^q, \quad (2)$$

where $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ are the mean vectors of \mathbf{X} and \mathbf{Y} , respectively, then the kernel is the inner product among centered covariance matrices when $q = 2$.

More generally, we can say that the q th order mean polynomial kernel contains all q th order moments as feature vectors. Indeed, if we let $\mathbb{P}_q = \{\mathbf{p} \in (\mathbb{N} \cup \{0\})^d | \mathbf{p}^\top \mathbf{1} = q\}$ and $x_{h,i}$ be the h th entry in \mathbf{x}_i , enumerating all q th order moments allows us to define

$$\phi_p(\mathbf{X}) = \frac{1}{\ell} \sqrt{\frac{q!}{p_1! \cdots p_d!}} \sum_{i=1}^{\ell} \prod_{h=1}^d x_{h,i}^{p_h}.$$

By using the feature map given by $\phi(\mathbf{X}) = [\phi_p(\mathbf{X})]_{p \in \mathbb{P}_q}$, we can derive the following equality

$$k_q(\mathbf{X}, \mathbf{Y}) = \langle \phi(\mathbf{X}), \phi(\mathbf{Y}) \rangle, \quad (3)$$

as given in Appendix A. Existence of a feature vector ensures the positive semidefiniteness of the mean polynomial kernel. Similarly for the centered version of the mean polynomial kernel, the features can be explicitly expressed as a set of all the q th order central moments (See Appendix C).

5. Mean Polynomial Kernel and Projection Kernel Relationship

We aim to establish a relationship between the mean polynomial kernel and the Projection kernel. In principle, Grassmann kernels are considered as kernel functions for principal subspaces. Eigendecomposition of two symmetric matrices Σ_x and Σ_y is essential for the computation of the Projection kernel value between two vector sequences \mathcal{X} and

\mathcal{Y} . Moreover, it can be shown that the bases of the principal subspaces are exactly the m major eigenvectors. To obtain the value of the Projection kernel between two subspaces \mathcal{X} and \mathcal{Y} , the m eigenvectors are initially stored in the matrices \mathbf{U}_x and \mathbf{U}_y . Let us define $\Sigma'_x = \mathbf{U}_x \mathbf{U}_x^\top$ and $\Sigma'_y = \mathbf{U}_y \mathbf{U}_y^\top$, the uncentered covariance matrices where the m major eigenvalues are replaced with ones and the rest of the eigenvalues are disregarded. Then the two kernels are related by the equality

$$k_{\text{PROJ}}(\mathbf{U}_x, \mathbf{U}_y) = \langle \text{vec}(\Sigma'_x), \text{vec}(\Sigma'_y) \rangle, \quad (4)$$

with the derivation given in Appendix B.

An assessment of both Eqs. (1) and (4) suggests that while the second order mean polynomial kernel preserves every bit of information in the uncentered covariance matrices, the Projection kernel possesses the possibility to disregard and lose information of each dimension of the principal subspaces, and all the information on their orthogonal complements. A similar case can be said for the centered version of the mean polynomial kernel (2) versus the Projection kernel, by using the centered covariance matrices. Although the first dilemma of the Projection kernel has been addressed by Hamm and Lee [22] by extending the kernel, resulting to the scaling of information of each dimension in linear subspaces and their preservation, data on the orthogonal complement are still overlooked. As with all dimension reduction techniques, there is always a risk of losing information when employing the Grassmann kernel. Though the hope is to retain the dimensions that are most discriminant, dimension number selection must be done with care and has become a critical stage in the implementation process. Furthermore, implementation via eigenvalue decomposition adds to the computational cost of k_{PROJ} , and also k_{BC} , giving the mean polynomial kernel an efficiency advantage, especially when presented with very high dimensional data.

6. Experiments and Discussion

We evaluate the performance of the mean polynomial kernel in binary classification tasks using data with underlying subspace structures. Techniques using the Grassmann kernels and Grassmann Distance Mutual Subspace method (GD-MSM) were also performed for comparison.

6.1 Face Membership Authentication

An important application of face recognition is face membership verification. The goal of this operation is to determine whether a subject is a 'member' or not. Moreover, we can also extend this to determining whether the given query is the authorized user or owner, which are common situations in accessing secured buildings or offices, logging on to computers, unlocking mobile phones, availing of online services, and other access control systems. The task can easily be modeled as a binary classification problem. For this purpose, we attempt to classify image sequences extracted from videos. The data was from the MOBIO database [29], and contains video data taken from 152 persons, each having

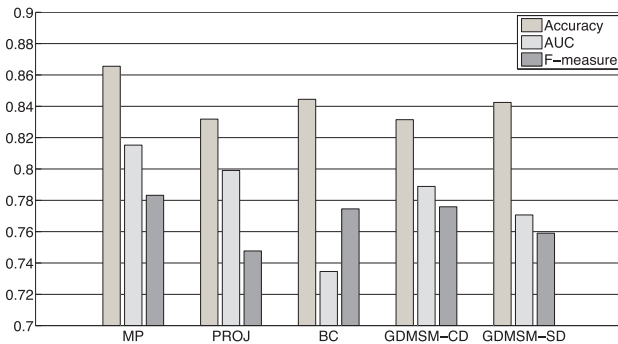


Fig. 4 Average performance of all methods for the face membership authentication task. The bar plot represents the average accuracy, average AUC, and average F-measure values computed.

12 video sessions divided into two: 6 sessions for Phase 1, and 6 sessions for Phase 2. Only data from 25 subjects and the 6 sessions from Phase 1 were used for the face membership verification task. Each session contains 21 image sequences of varying length. For the experiments, we set the sequence length to 25 images, where each image is a cropped face image of the subject, obtained using a face detection program, transformed to gray scale and resized to 25×25 pixels. Among the 25 subjects, 10 were randomly selected and labeled as ‘member’ (+1), and the remaining 15 as ‘nonmember’ (−1).

Two methods were employed: one using kernels with SVM and the other one using GD-MSM. For the first method, three types of kernel functions were utilized: the Grassmann kernels, Projection (PROJ) and Binet-Cauchy (BC) kernels, and the mean polynomial kernel (MP). For the GD-MSM, eight metrics were used for comparison: average distance, Binet-Cauchy metric, Geodesic distance, maximum correlation, minimum correlation, Frobenius norm based Procrustes distance, 2-norm based Procrustes distance, and Projection metric, as defined in [5]. For the SVM setting, 6-fold cross-validation was employed to evaluate the performance of the kernels such that one session per subject is used as test data while the remaining five sessions are used for training. On the other hand, class-wise (GDMSM-CD) and subject-wise (GDMSM-SD) dictionaries were implemented for the GD-MSM, as described in Sect. 3.

As for the parameters of the kernel methods, the value of q for the MP kernel was varied from 1 to 5, while the dimension of the subspace, m , was varied from 1 to 10. The regularization parameter C for SVM was varied over the set $\{10^0, 10^1, 10^2, 10^3, 10^4, 10^5\}$. To optimize the tuning of the said parameters, we implemented a 3-fold cross-validation grid search of the pairs (q, C) and (m, C) on the training data, for each cross-validation set. Values of the pairs were chosen such that the highest accuracy value is obtained. Variation and selection of the value of m for GDMSM was also done in a similar manner. The area under the ROC curve (AUC), accuracy, and F-measure values were considered for evaluating the performance of each method.

Figure 4 illustrates the average accuracy, AUC, and F-

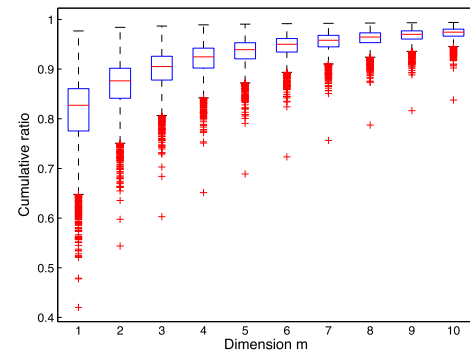


Fig. 5 Cumulative ratio distribution of the eigenvalues of the face video sequences as dimension m is varied.

measure values of each method for all 6 cross-validation sets. From the graph, it is evident that the MP kernel outperforms the other methods on all three benchmarks (with accuracy, AUC and F-measure values of 81.5%, 0.866, and 0.783, respectively). Meanwhile, PROJ, BC, and GDMSM-CD obtained the second best accuracy (79.9%), AUC (0.845), and F-measure (0.776), respectively. The values presented here for the two GDMSM’s are the highest obtained among all eight metrics used, which, interestingly, is the maximum correlation. We can therefore conclude that the method employing the MP kernel plus SVM is better than the GD-MSM regardless of the selected metric.

The cumulative ratio distributions of the eigenvalues of the image sequences as the dimension parameter m of the Grassmann methods is varied are presented as box plots in Fig. 5. As evident in the figure, even when dimension $m = 1$, median of the cumulative ratio is acceptably high (0.827). This gradually increases to 0.974 when $m = 10$, which is the maximum dimension we considered for the cross-validation in setting the parameter m . Values of the cumulative sum, in general, are relatively high, which may give decent representation of data using only the leading components. However, even with such decent representation, we experience a significant difference in the variance of the cumulative ratio as the dimension increases, and a significant difference in the performance of the MP kernel and Grassmann kernels. One factor affecting the performance of the Projection kernel may be due to the several outliers that can be observed below the minimum of the first quartile of each box plot which are quite low in value (at least 0.420), having a large discrepancy from the maximum (at most 0.557 difference).

6.2 EEG Signal Task Classification

We also compared the performances of MP, PROJ, BC, GDMSM-CD and GDMSM-SD on the BCI competition III-IVa dataset [8]. The data contains recorded measurements of five subjects (aa, al, av, aw, and ay) during motor imagery tasks (right hand and right foot movement) using 118 channels of electrodes. The EEG signals were recorded for 3.5 seconds with 1000 Hz sampling rate for each trial. However, we used the available downsampled version (at 100 Hz) of

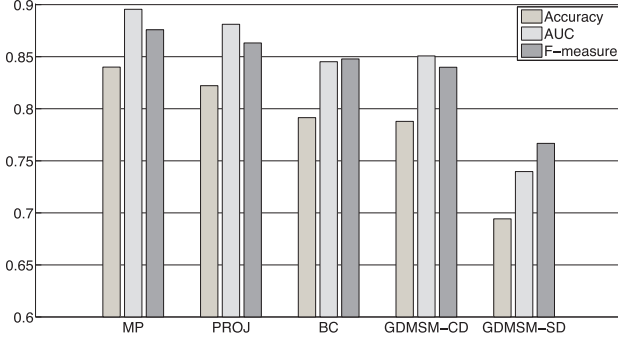


Fig. 6 Average performance of all methods for the EEG signal task classification. The bar plot represents the average accuracy, average AUC, and average F-measure values computed.

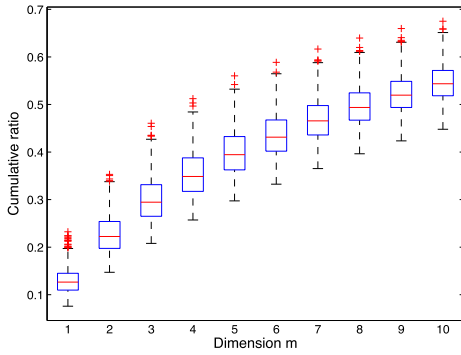


Fig. 7 Cumulative ratio distribution of the eigenvalues of the EEG signal sequences as dimension m is varied.

the data, and utilized the 0.5 to 3.5-second interval from the visual cues for each trial, resulting to a time range of 3.0 sec per trial. For data preprocessing, frequency band selection was done, and data was filtered between frequencies of 10 to 35 Hz. For each subject, 140 trials were conducted for each task, for a total of 280 trials per subject. Settings similar to the previous application were applied to the experiments using BCI data, including the approach on parameter selection.

The average values of the performance evaluators for all subjects over the 5 cross-validation sets are given in Fig. 6. As expected, the MP kernel bests the other approaches, with accuracy, AUC, and F-measure values of 84.0%, 0.896, and 0.876, respectively. This is followed by the PROJ method, with values 82.2%, 0.881, and 0.863, respectively. The GD-MSM results are also of the best performing metric, which in this case is also the maximum correlation. Hence, in a parallel logic to the previous experiments, we also conclude that the proposed method surpasses the GD-MSM approach for this task, irrespective of the metric used.

Figure 7 shows the box plots of the cumulative ratio distribution of the EEG signal sequences for each dimension $m = 1$ to 10. In contrast to the ratios presented in Fig. 5, the values for the EEG data are very low. At minimum, the median is 0.127 ($m = 1$), and maximum is 0.544 ($m = 10$). The

Table 1 Time complexity comparison of the kernels.

| | Training Stage (For kernel matrix computation) | Testing Stage (For prediction of a single sequence) |
|------|---|---|
| MP | $O(n_{\text{tra}}^2 \ell^2 d \log_2 q)$ | $O(n_{\text{sv}} \ell^2 d \log_2 q)$ |
| PROJ | For covariance matrix computation $O(d^2 \ell n_{\text{tra}})$ | $O(d^2 \ell)$ |
| | Eigendecomposition $O(k^3 n_{\text{tra}})$ | $O(k^3)$ |
| | Kernel value computation $O(dm^2 n_{\text{tra}}^2)$ | $O(d^2 m n_{\text{sv}})$ |
| BC | For covariance matrix computation $O(d^2 \ell n_{\text{tra}})$ | $O(d^2 \ell)$ |
| | Eigendecomposition $O(k^3 n_{\text{tra}})$ | $O(k^3)$ |
| | Kernel value computation $O(m^3 n_{\text{tra}}^2)$ | $O(d^2 m n_{\text{sv}})$ |

cumulative ratios also vary significantly as the dimension changes. Outliers can be observed above the fourth quartile, but not as much as in the face video data. Moreover, the difference between the maximum (outlier) and the minimum in each respective box plot is at least 0.157 and at most 0.263, which are significantly lower than those in the previous dataset. This may explain why the Projection kernel and most of the other Grassmann-based methods perform better in terms of AUC and F-measure values on this data set.

6.3 Efficiency Comparison

We investigate the time complexity of the MP kernel, and compare it with the Grassmann kernels. Suppose we are given n_{tra} number of training samples, and n_{sv} number of support vectors. For simplicity, we will assume that every (feature) vector sequence has length ℓ , and that each vector has length d . Moreover, we denote the dimension of the principal subspace as m for the Grassmann kernels, and let $k = \min(\ell, d)$. In Table 1, we give the computation time for each step in the calculation of the kernels. From this table, we conclude that the MP kernel is not only better in terms of performance, but it is also more efficient in terms of computational cost compared to the Grassmann kernels. This was confirmed empirically, as the average CPU time recorded for the MP kernel, for any value of q , is around 383 seconds for the MOBIO data, and 58 sec for the EEG data. On the other hand, computation of both Grassmann kernel matrices is around 1.21×10^4 sec when $m = 5$, and 1.24×10^4 sec for PROJ, and 1.25×10^4 sec for BC when $m = 10$, using the MOBIO data. On the EEG data, CPU time of PROJ is about 1.20×10^3 for any m , while the BC takes 1.22×10^3 and 1.23×10^3 when $m = 5$ and $m = 10$, respectively. It is also worth mentioning that should the number of features d increase, the computational time for the Grassmann kernels will drastically increase, whereas the increase with the MP kernel is only linear.

6.4 Discussion

We conclude this section by considering an extension of

the mean polynomial kernel. There are many possible extensions, one of which is by replacing the sample mean $\langle \mathbf{x}_i, \mathbf{y}_j \rangle^q$ with the expected value with respect to a probabilistic distribution: $k'_q(X, Y) = \mathbb{E}(\langle \mathbf{x}, \mathbf{y} \rangle^q)$. From this, the mean polynomial kernel can be derived as a special case when $p_x(\mathbf{x}) = \frac{1}{\ell} \sum_{i=1}^{\ell} \delta(\mathbf{x} - \mathbf{x}_i)$ and $p_y(\mathbf{y}) = \frac{1}{\ell'} \sum_{i=1}^{\ell'} \delta(\mathbf{y} - \mathbf{y}_i)$, where $\delta(\cdot)$ is the Dirac delta function.

Another choice of a probabilistic distribution can be Gaussian mixture. Suppose we are given two Gaussian mixtures

$$p_x(\mathbf{x}) = \sum_{i=1}^{\ell} \pi_{x,i} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{x,i}, \boldsymbol{\Sigma}_{x,i})$$

and

$$p_y(\mathbf{y}) = \sum_{j=1}^{\ell'} \pi_{y,j} \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_{y,j}, \boldsymbol{\Sigma}_{y,j}),$$

where ℓ and ℓ' are the number of Gaussian components for the two probabilistic distributions p_x and p_y , respectively, $\pi_{z,i}$ is the mixing coefficient satisfying $\sum_{i=1}^n \pi_{z,i} = 1$, and $\boldsymbol{\mu}_{z,i}$ and $\boldsymbol{\Sigma}_{z,i}$ are the mean vector and covariance matrix of the i th Gaussian component, respectively. The second order mean polynomial kernel can be readily computed as

$$\begin{aligned} k_2(p_x, p_y) = & \sum_{i=1}^{\ell} \sum_{j=1}^{\ell'} \pi_{x,i} \pi_{y,j} \left((\boldsymbol{\mu}_{x,i}^T \boldsymbol{\mu}_{y,j})^2 + \right. \\ & \text{tr}(\boldsymbol{\Sigma}_{x,i} \boldsymbol{\Sigma}_{y,j}) + \boldsymbol{\mu}_{x,i}^T \boldsymbol{\Sigma}_{y,j} \boldsymbol{\mu}_{x,i} + \\ & \left. \boldsymbol{\mu}_{y,j}^T \boldsymbol{\Sigma}_{x,i} \boldsymbol{\mu}_{y,j} \right). \end{aligned}$$

This example includes the original definition of the mean polynomial kernel in Definition 2, which can be shown by letting

$$\begin{aligned} \pi_{x,i} &= 1/\ell, & \boldsymbol{\mu}_{x,i} &= \mathbf{x}_i, & \boldsymbol{\Sigma}_{x,i} &= \sigma_{x,i}^2 \mathbf{I}, \\ \pi_{y,j} &= 1/\ell', & \boldsymbol{\mu}_{y,j} &= \mathbf{y}_j, & \boldsymbol{\Sigma}_{y,j} &= \sigma_{y,j}^2 \mathbf{I}, \end{aligned}$$

for all $i \in \mathbb{N}_{\ell}$ and $j \in \mathbb{N}_{\ell'}$, and taking the limit as $\sigma^2 \rightarrow 0$. When one wishes to weight each frame in image sequences, the weights can be set to $\pi_{x,i}$ or $\pi_{y,j}$. Positive $\sigma_{x,i}^2$ or positive $\sigma_{y,j}^2$ can be used to represent uncertainties in observations,

Similar to the original mean polynomial kernel, we can explicitly represent features of the extended mean polynomial kernel, as given in Appendix C.

7. Conclusion

We have examined the mean polynomial kernel as a kernel for binary classification of data modeled as vector sets or sequences. Analogy and connection to related methods, Grassmann Projection kernel in particular, have also been drawn. The effectiveness of the MP kernel was empirically supported using data of face image sequences, and motor imagery EEG recordings. Furthermore, we present a comparison of computational costs between methods, and some

interesting extensions of the MP kernel by considering the probabilistic distribution of the data. In brief, the mean polynomial kernel excels known methods from literature, both in performance and efficiency. In addition to the performed experiments, application to data vector sets in a multi classification problem setting may prove to be an interesting direction.

Acknowledgments

The work of RR is supported by the Ministry of Education, Culture, Sports, Science and Technology of Japan. TK is supported by MEXT KAKENHI Grant number 23500373.

References

- [1] K. Fukui and O. Yamaguchi, "Face recognition using multi-viewpoint pattern for robot vision," *Proc. Int. Symp. Robotics Research*, pp.192–201, 2003.
- [2] J. Hamm and D. Lee, "Grassmann discriminant analysis: A unifying view on subspace-based learning," *Proc. 25th ICML*, pp.376–383, 2008.
- [3] T.K. Kim, J. Kittler, and R. Cipolla, "Discriminative learning and recognition of images set classes using canonical correlations," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol.29, pp.1005–1018, 2007.
- [4] H. Sakano and N. Mukawa, "Kernel mutual subspace method for robust facial image recognition," *Proc. Int. Conf. on Knowledge-Based Intell. Eng. Sys. And App. Tech.*, pp.245–248, 2000.
- [5] R. Shigenaka, B. Raychev, T. Tamaki, and K. Kaneda, "Face sequence recognition using Grassmann distances and Grassmann kernels," *Proc. IJCNN*, pp.2630–2636, 2012.
- [6] L. Wolf and A. Shashua, "Learning over sets using kernel principal angles," *J. Mach. Learn. Res.*, vol.4, pp.913–931, 2003.
- [7] O. Yamaguchi, K. Fukui, and K. Maeda, "Face recognition using temporal image sequence," *Proc. Int. Conf. Automatic Face and Gesture Recognition*, pp.318–323, 1998.
- [8] B. Blankertz, K. Müller, D. Krusienki, G. Schalk, J. Wolpaw, A. Schlögl, G. Pfurtscheller, J. Millán, M. Schröder, and N. Birbaumer, "The BCI competition III: Validating alternative approaches to actual BCI problems," *IEEE Trans. Neural Sys. Rehab. Eng.*, vol.14, no.2, pp.153–159, 2006.
- [9] J. Müller-Gerking, G. Pfurtscheller, and H. Flyvbjerg, "Designing optimal spatial filters for single-trial EEG classification in movement task," *Clinical Neurophysiology*, vol.110, no.5, pp.787–798, 1999.
- [10] H. Kashima, K. Tsuda, and A. Inokuchi, "Kernels for graphs," in *Kernels and Bioinformatics*, pp.155–170, MIT Press, Cambridge, MA, USA, 2004.
- [11] R. Kondor and J. Lafferty, "Diffusion kernels on graphs and other discrete input spaces," *Proc. 19th ICML*, pp.315–322, July 2002.
- [12] T. Kuboya, K. Hirata, and K. Aoiki-Kinoshita, "An efficient unordered tree kernel and its application to glycan classification," *Proc. PAKDD*, pp.184–195, 2008.
- [13] M. Neumann, N. Patricia, R. Garnett, and K. Kersting, "Efficient graph kernels by randomization," *Proc. ECML/PKDD*, pp.378–393, 2012.
- [14] Q. Shi, J. Petterson, G. Dror, J. Langford, A. Smola, and S. Vishwanathan, "Hash kernels for structured data," *J. Mach. Learn. Res.*, vol.10, pp.2615–2637, 2009.
- [15] S. Vishwanathan, N. Schraudolph, R. Kondor, and K. Borgwardt, "Graph kernels," *J. Mach. Learn. Res.*, vol.11, pp.1201–1242, 2010.
- [16] C. Leslie, E. Eskin, and W. Noble, "The spectrum kernel: A string kernel for SVM protein classification," *Proc. Pacific Symposium on Biocomputing*, pp.564–575, 2002.

- [17] H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins, "Text classification using string kernels," *J. Mach. Learn. Res.*, vol.2, pp.419–444, 2002.
- [18] S. Qiu, T. Lane, and L. Buturovic, "A randomized string kernel and its application to RNA interference," *Proc. AAAI*, pp.627–632, 2007.
- [19] S. Vishwanathan and A. Smola, "Fast kernels for string and tree matching," *Advances in Neural Info. Proc. Sys.*, pp.569–576, 2003.
- [20] F. Desobry, M. Davy, and W. Fitzgerald, "A class of kernels for sets of vectors," *Proc. 13th ESANN*, pp.461–466, 2005.
- [21] K. Grauman and T. Darrell, "The pyramid match kernel: Discriminative classification with sets of image features," *Proc. IEEE ICCV*, vol.2, pp.1458–1465, Beijing, China, Oct. 2005.
- [22] J. Hamm and D. Lee, "Extended Grassmann kernels for subspace-based learning," *Proc. NIPS*, pp.601–608, 2008.
- [23] T.K. Kim, J. Kittler, and R. Cippolla, "Discriminative learning and recognition of image set classes using canonical correlations," *IEEE Trans. Pattern Anal. Mach. Intell.*, pp.1005–1018, 2007.
- [24] R. Kondor and T. Jebara, "A kernel between sets of vectors," *Proc. 20th ICML*, pp.361–368, Aug. 2003.
- [25] B. Schölkopf and A. Smola, *Learning with kernels: Support vector machines, regularization, optimization, and beyond*, MIT Press, Cambridge, MA, USA, 2001.
- [26] J. Shawe-Taylor and N. Cristianini, *Kernel methods for pattern analysis*, Cambridge University Press New York, NY, USA, 2004.
- [27] S. Vishwanathan and A. Smola, "Binet-cauchy kernels," *Proc. NIPS*, pp.1441–1448, 2004.
- [28] T. Gärtner, P. Flach, A. Kowalczyk, and A. Smola, "Multi-instance kernels," *Proc. 19th ICML*, pp.179–186, June 2002.
- [29] C. McCool, S. Marcel, A. Hadid, M. Pietikainen, P. Matejka, J. Cernocky, N. Poh, J. Kittler, A. Larcher, C. Levy, D. Matrouf, J.F. Bonastre, P. Tresadern, and T. Cootes, "Bi-modal person recognition on a mobile phone: Using mobile phone data," *Proc. IEEE ICME Workshop on Hot Topics in Mobile Multimedia*, pp.635–640, 2012.

Appendix A: Derivation of Eq. (3)

Proposition 3. A mapping function of mean polynomial kernel is

$$\phi(X) = \left[\frac{1}{\ell} \sqrt{\frac{q!}{p_1! \cdots p_d!}} \sum_{i=1}^{\ell} \prod_{h=1}^d x_{h,i}^{p_h} \right],$$

where $\mathbf{p} \in (\mathbb{N} \cup \{0\})^d$ such that $\mathbf{p}^T \mathbf{1} = q$.

Proof: Let $x_{h,i}$ and $y_{h,i}$ be the (h, i) th entries in \mathbf{X} and \mathbf{Y} , respectively.

Let d be the number of rows in \mathbf{X} and \mathbf{Y} ,

$$\begin{aligned} k_q(\mathbf{X}, \mathbf{Y}) &= \frac{1}{\ell \ell'} \sum_{i,j} \langle \mathbf{x}_i, \mathbf{y}_j \rangle^q \\ &= \frac{1}{\ell \ell'} \sum_{i,j} \left(\sum_{h=1}^d x_{h,i} y_{h,j} \right)^q. \end{aligned}$$

Using the multinomial theorem, we get

$$\begin{aligned} k_q(\mathbf{X}, \mathbf{Y}) &= \frac{1}{\ell \ell'} \sum_{i,j} \sum_{\mathbf{p}} \frac{q!}{p_1! \cdots p_d!} \prod_{h=1}^d x_{h,i}^{p_h} y_{h,j}^{p_h} \\ &= \sum_{\mathbf{p}} \left(\frac{1}{\ell} \sqrt{\frac{q!}{p_1! \cdots p_d!}} \sum_{i=1}^{\ell} \prod_{h=1}^d x_{h,i}^{p_h} \right) \times \end{aligned}$$

$$\begin{aligned} &\left(\frac{1}{\ell'} \sqrt{\frac{q!}{p_1! \cdots p_d!}} \sum_{j=1}^{\ell'} \prod_{h=1}^d y_{h,j}^{p_h} \right) \\ &= \langle \phi(\mathbf{X}), \phi(\mathbf{Y}) \rangle, \end{aligned}$$

where $\mathbf{p} \in (\mathbb{N} \cup \{0\})^d$ such that $\mathbf{p}^T \mathbf{1} = q$. □

Appendix B: Derivation of Eq. (4)

Suppose the transformed covariance matrices are given by $\Sigma'_x = \mathbf{U}_x \Lambda_x \mathbf{U}_x^T = \mathbf{U}_x \mathbf{U}_x^T$ and $\Sigma'_y = \mathbf{U}_y \Lambda_y \mathbf{U}_y^T = \mathbf{U}_y \mathbf{U}_y^T$, obtained via eigendecomposition of the covariance matrices Σ_x and Σ_y , and setting the major eigenvalues to one and the minor eigenvalues to zero. Then we can write

$$\begin{aligned} k_{\text{PROJ}}(\mathbf{U}_x, \mathbf{U}_y) &= \|\mathbf{U}_x^T \mathbf{U}_y\|_F^2 = \text{tr}(\mathbf{U}_x^T \mathbf{U}_y \mathbf{U}_y^T \mathbf{U}_x) \\ &= \text{tr}(\mathbf{U}_x \mathbf{U}_x^T \mathbf{U}_y \mathbf{U}_y^T) = \text{tr}(\Sigma'_x \Sigma'_y) = \langle \text{vec}(\Sigma'_x), \text{vec}(\Sigma'_y) \rangle. \end{aligned}$$

This concludes the derivation of Eq. (4).

Appendix C: Explicit Representation of Features

In Sect. 4, we have shown that the features of the mean polynomial kernel can be represented explicitly. Features of the centered mean polynomial kernel are represented by the q th central moments:

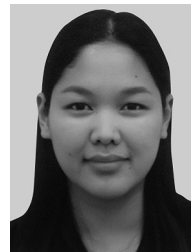
$$\bar{\phi}_{\mathbf{p}}(\mathbf{X}) = \frac{1}{\ell} \sqrt{\frac{q!}{p_1! \cdots p_d!}} \sum_{i=1}^{\ell} \prod_{h=1}^d (x_{h,i} - \bar{x}_h)^{p_h}.$$

The features that produce the extended mean polynomial kernel are given by

$$\phi_{\mathbf{p}}(p_x) = \sqrt{\frac{q!}{p_1! \cdots p_d!}} \mathbb{E} \left(\prod_{h=1}^d x_{h,i}^{p_h} \right)$$

for all $\mathbf{p} \in (\mathbb{N} \cup \{0\})^d$ such that $\mathbf{p}^T \mathbf{1} = q$, and the features for the extended centered mean polynomial kernel are given by

$$\bar{\phi}_{\mathbf{p}}(p_x) = \sqrt{\frac{q!}{p_1! \cdots p_d!}} \mathbb{E} \left(\prod_{h=1}^d (x_h - \mathbb{E}(x_h))^{p_h} \right).$$



Raissa Relator received the B.S. and M.S. degrees in Mathematics from the University of the Philippines Diliman in 2005 and 2008, respectively. She is currently pursuing a Ph.D. in Computer Science at the Graduate School of Science and Engineering, Gunma University. Her research interests include pattern recognition, machine learning, and bioinformatics.



Yoshihiro Hirohashi received the B.E. from Gunma University in 2013. He is currently pursuing a M.E. degree at the Graduate School of Engineering, Tohoku University. His research interests include biomedical engineering, machine learning and computer vision.



Eisuke Ito received the B.E. from Gunma University in 2012, and is currently pursuing his M.E. degree at the Graduate School of Science and Engineering, Gunma University. His research interests include neuroscience, machine learning and brain-computer interfaces. He is a student member of the JSAI.



Tsuyoshi Kato received the B.E., M.E., and Ph.D. degree from Tohoku University, Sendai, Japan, in 1998, 2000, and 2003, respectively. From 2003 to 2005, he was with the National Institute of Advanced Industrial Science Technology (AIST) as a postdoctoral fellow in the Computational Biology Research Center (CBRC) in Tokyo. From 2005 to 2008, he was an assistant professor at the Graduate School of Frontier Sciences, University of Tokyo. From 2008 to 2010, he was an associate professor at the Center

for Informational Biology, Ochanomizu University. He then moved back to Graduate School of Frontier Sciences, University of Tokyo, and as of present, he is an associate professor at the Graduate School of Science and Engineering, Gunma University. His current scientific interests include pattern recognition, computer vision and bioinformatics. He is a member of IEICEJ.