

Bayesian kernel adaptive grouping learning for multi-dimensional datasets*

XIAOZHOU WANG[†] AND FANGLI DONG^{†‡}

With the development of information technology, a large number of datasets with complex structures, such as multi-dimensional datasets, need to be processed and analyzed. In this paper we propose a kernel-based statistical learning algorithm, Bayesian Kernel Adaptive Grouping Learning (BKAGL), to provide an innovative solution for the classification problem of multi-dimensional datasets. BKAGL can integrate information from different dimensions adaptively. Meanwhile, we adopt the Bayesian framework which can infer the approximate posterior distributions of parameters. The utilization of grouping features can help find which groups of features have more contributions to the response. Simulation results illustrate that BKAGL outperforms some classical classification methods and the corresponding ungrouped method. The analysis of the electrocardiogram (ECG) and electroencephalography (EEG) datasets shows that BKAGL has the highest classification accuracy and provides explanatory information.

AMS 2000 SUBJECT CLASSIFICATIONS: 62F15, 62H30.

KEYWORDS AND PHRASES: Classifier, Multi-dimensional dataset, Bayesian model, Adaptiveness, Kernel method.

1. INTRODUCTION

One of the main purposes of machine learning is to infer the relationship between the outcome Y and the variable X so as to improve human awareness [1, 7, 16]. However, with the explosive growth of the data, how to mine and integrate the information from the data is a challenge to the learning algorithms.

Support Vector Machine (SVM) [13] is a well-known machine learning method due to its excellent generalization properties. The decision function which is used for predicting the label y of one instance \mathbf{x} can be written as

$$(1) \quad f(\mathbf{x}) = \mathbf{a}^T \mathbf{k}(\mathbf{x}) + e,$$

*Supported by RGC Competitive Earmarked Research Grants, National Basic Research Program of China (973 Program, 2015CB856004) and National Natural Science Foundation of China (11531001).

[†]These authors contributed equally to this work.

[‡]Corresponding author.

where $\mathbf{a} = (a_1, \dots, a_N)^T$ is the sample weight, e is the bias, $\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}_1, \mathbf{x}), \dots, k(\mathbf{x}_N, \mathbf{x}))^T$ and the kernel function $k(\mathbf{x}_i, \mathbf{x})$ satisfies Mercer's condition [10]. The parameters in Eq. (1) can be obtained through a quadratic optimization process [27]. Furthermore, the parameters can also be optimized under a Bayesian framework [26], which is called Relevance Vector Machine (RVM). The sample weight parameter \mathbf{a} is viewed as a random vector. Both SVM and RVM are kernel-based methods which have a lot of advantages [14, 25]. Moreover, when we want to integrate data from different sources, the multiple kernel learning algorithm [12, 22, 24] is usually considered. For example, we use multiple omics data, including the high-throughput genomic, epigenomic, and transcriptomic data, to predict a person's phenotypic response [29]. From the viewpoint of data structure, multiple omics data correspond to multiple dimensions. For the aim of integrating these multi-dimensional data, multiple kernels can be used where each kernel is a function of each dimension. Then multiple kernels can integrate multi-dimensional data by a kernel-based decision function as follows:

$$(2) \quad f(\mathbf{x}) = \mathbf{a}^T \sum_{q=1}^Q \rho_q \mathbf{k}_q(\mathbf{x}) + e,$$

where Q is the number of kernels and also the number of dimensions, and $\boldsymbol{\rho} = (\rho_1, \dots, \rho_Q)^T$ is the kernel weight parameter. For the q -th kernel, $\mathbf{k}_q(\mathbf{x}) = (k_q(\mathbf{x}_1, \mathbf{x}), \dots, k_q(\mathbf{x}_N, \mathbf{x}))^T$. Through the estimation of kernel weight parameter, we can obtain which dimensions contribute more to the response. Moreover, for single omics data such as gene expression data, a pathway contains a group of genes. There are different pathways (groups) in gene expression data. If we want to decipher the correlation between the pathways and response, Eq. (2) can also be used through defining multiple kernels based on different pathways. However, a mass of parameters will need to be estimated. In this case, a novel model is necessary to integrate information from different groups in different dimensions. Naturally, the extension of Eq. (2) as below can be utilized.

$$(3) \quad f(\mathbf{x}) = \mathbf{a}^T \sum_{d=1}^D c_d \sum_{m=1}^{P_d} b_{md} \mathbf{k}_{dm}(\mathbf{x}) + e,$$

where D is the number of dimensions and $\mathbf{c} = (c_1, \dots, c_D)^T$

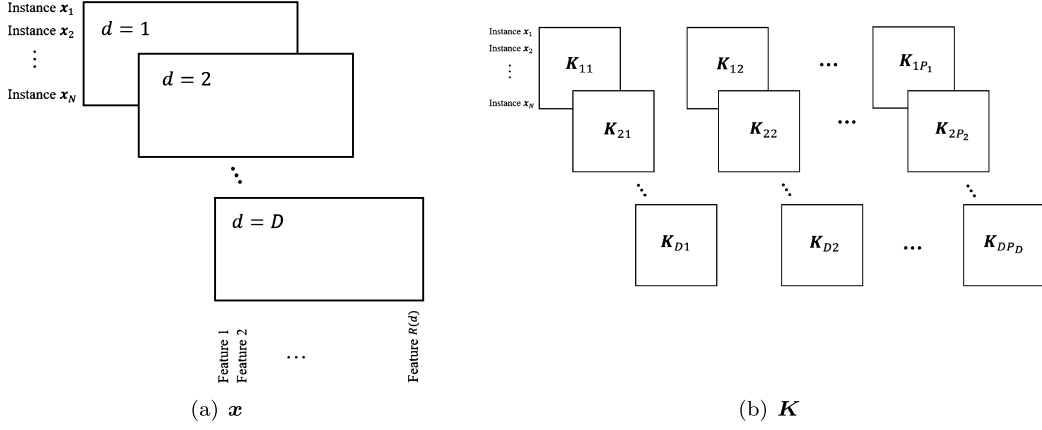


Figure 1. Structures of the dataset \mathbf{x} and kernel matrix set \mathbf{K} .

is the dimension weight parameter. P_d is the group number of features in the d -th dimension and \mathbf{b} is the two-dimensional group weight parameter. For the m -th group in the d -th dimension, $\mathbf{k}_{dm}(\mathbf{x}) = (k_{dm}(\mathbf{x}_1, \mathbf{x}), \dots, k_{dm}(\mathbf{x}_N, \mathbf{x}))^T$ and $k_{dm}(\mathbf{x}_i, \mathbf{x}_j)$ is the kernel function of \mathbf{x}_i and \mathbf{x}_j for the subjects i and j . Since the grouping result of each dimension may be different, the group number P_d and the parameter b_{md} are decided by the corresponding d -th dimension. As shown in Eq. (3), we propose a model which extracts information layer by layer. In particular, the proposed algorithm is based on a Bayesian framework [19, 21] and uses an approximate inference process [6]. BKAGL can simultaneously integrate information from different groups in different dimensions and then predict the outcome. At the same time, it has good interpretation ability, which enables us to find which dimensions and groups of features have higher contributions to the outcome. This is our concern.

The rest of the paper is organized as follows: Section 2 firstly introduces the notations used in the paper, and then explains the model. In Section 3, the process of parameter estimation and the prediction for a new point are provided. Section 4 examines the performance of our algorithm including accuracy and interpretability by the simulation and the experiment of the real ECG and EEG datasets. Conclusions and some extensions for further work are given in Section 5.

2. THE FRAMEWORK

2.1 Notations

Suppose that there are N i.i.d. instances $\{(\mathbf{x}_i, y_i) : i = 1, \dots, N\}$. Each instance \mathbf{x}_i has D dimensions and each dimension is a vector with $R(d)$ features. More precisely, x_{idr} represents the r -th feature of the d -th dimension of the i -th instance. The response y_i is the class label of the i -th instance \mathbf{x}_i . Assume that all the features of each dimension have been divided into different groups in a certain way. Denote P_d as the number of the groups of features in the d -th dimension. Define

- Kernel matrix set $\mathbf{K} = \{\mathbf{K}_{dm} : d = 1, \dots, D, m = 1, \dots, P_d\}$, where \mathbf{K}_{dm} is a $N \times N$ kernel matrix of the m -th group of features in the d -th dimension. $(\mathbf{K}_{dm})_{\cdot i}$ is the i -th column of \mathbf{K}_{dm} and $(\mathbf{K}_{dm})_j$ is the j -th row of \mathbf{K}_{dm} .
- $\mathbf{G} = \{\mathbf{G}_d : d = 1, \dots, D\}$, where \mathbf{G}_d is a $P_d \times N$ matrix. $\mathbf{G}_{d \cdot i}$ is the i -th column of \mathbf{G}_d and \mathbf{G}_{dm} is the m -th row of \mathbf{G}_d .
- \mathbf{L} is a $D \times N$ matrix. $\mathbf{L}_{\cdot i}$ is the i -th column of \mathbf{L} and \mathbf{L}_d is the d -th row of \mathbf{L} .
- \mathbf{f} is a vector with N elements. In fact, \mathbf{f} is a latent variable connecting the second layer \mathbf{L} of intermediate output and the class label vector \mathbf{y} of N instances. More details about \mathbf{f} are in Section 2.2.

Figure 1 shows the structures of dataset \mathbf{x} and kernel matrix set \mathbf{K} in order to help understand the data structure.

2.2 Model

The main aim of BKAGL is to add the interpretability of the model through grouping features and make the number of estimated parameters under control at the same time compared with Eq. (2). The proposed BKAGL algorithm extracts information from layer to layer. After getting the kernel matrix set \mathbf{K} through the input training set, we calculate the first layer \mathbf{G} of intermediate output by combining different kernels for different groups in each dimension as

$$(4) \quad \mathbf{G}_d = \begin{pmatrix} \mathbf{a}^T \mathbf{K}_{d1} \\ \vdots \\ \mathbf{a}^T \mathbf{K}_{dP_d} \end{pmatrix}, \quad d = 1, \dots, D.$$

After obtaining the first layer \mathbf{G} of intermediate output, the second layer \mathbf{L} of intermediate output can be calculated.

$$(5) \quad \mathbf{L} = \begin{pmatrix} \mathbf{b}_1^T \mathbf{G}_1 \\ \vdots \\ \mathbf{b}_D^T \mathbf{G}_D \end{pmatrix},$$

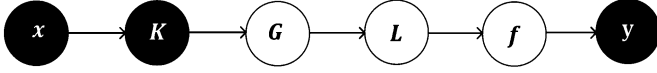


Figure 2. Classification of multi-dimensional datasets with BKAGL.

where $\mathbf{b}_{\cdot d} = (b_{1d}, \dots, b_{P_d d})^T$ for $d = 1, 2, \dots, D$. Then through combining Eqs. (4) and (5), f_i , $i = 1, \dots, N$, takes the following form:

$$\begin{aligned}
f_i &= \mathbf{c}^T \mathbf{L}_{\cdot i} + e = \sum_{d=1}^D c_d L_{di} + e \\
&= \sum_{d=1}^D c_d \mathbf{b}_{\cdot d}^T \mathbf{G}_{d \cdot i} + e = \sum_{d=1}^D c_d \sum_{m=1}^{P_d} b_{md} G_{dmi} + e \\
&= \sum_{d=1}^D c_d \sum_{m=1}^{P_d} b_{md} \mathbf{a}^T (\mathbf{K}_{dm})_{\cdot i} + e \\
&= \sum_{d=1}^D c_d \sum_{m=1}^{P_d} b_{md} \sum_{j=1}^N a_j (\mathbf{K}_{dm})_{ji} + e \\
&= \mathbf{a}^T \sum_{d=1}^D \sum_{m=1}^{P_d} c_d b_{md} (\mathbf{K}_{dm})_{\cdot i} + e \triangleq \mathbf{a}^T \mathcal{K}_i + e.
\end{aligned}$$

This formula indicates that the similarity measure \mathcal{K}_i is a linear combination of kernels defined in different groups in different dimensions. Through learning the parameters \mathbf{b} and \mathbf{c} , we can find the more informative groups and dimensions for the response. The complete learning process is summarized in Figure 2.

For computational convenience, conjugate prior distributions [11] are applied in BKAGL. Denote all hyper-parameters $\Xi = \{\alpha_\lambda, \beta_\lambda, \alpha_\eta, \beta_\eta, \alpha_\gamma, \beta_\gamma, \alpha_\omega, \beta_\omega\}$ and the prior set $\Theta = \{\lambda, \eta, \gamma, \omega\}$. Denote weight and bias parameters by the set $\Lambda = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, e\}$, where vector \mathbf{a} is the sample weight parameter, two-dimensional parameter \mathbf{b} is the group weight parameter, vector \mathbf{c} is the dimension weight parameter and e is the bias. Assume that the distributions of the related random variables in our proposed model (Eq. (3)) satisfy the following assumption (A1).

Assumption (A1): Note that $i = 1, \dots, N$, $d = 1, \dots, D$ and $m = 1, \dots, P_d$. For parameters in Θ :

$$\begin{aligned}
\lambda_i &\sim \text{Gamma}(\lambda_i; \alpha_\lambda, \beta_\lambda), \\
\eta_{md} &\sim \text{Gamma}(\eta_{md}; \alpha_\eta, \beta_\eta), \\
\gamma_d &\sim \text{Gamma}(\gamma_d; \alpha_\gamma, \beta_\gamma), \\
\omega &\sim \text{Gamma}(\omega; \alpha_\omega, \beta_\omega).
\end{aligned}$$

For parameters in Λ :

$$\begin{aligned}
a_i | \lambda_i &\sim \mathcal{N}(a_i; 0, \lambda_i^{-1}), \\
b_{md} | \eta_{md} &\sim \mathcal{N}(b_{md}; 0, \eta_{md}^{-1}), \\
c_d | \gamma_d &\sim \mathcal{N}(c_d; 0, \gamma_d^{-1}), \\
e | \omega &\sim \mathcal{N}(e; 0, \omega^{-1}).
\end{aligned}$$

For intermediate output variables:

$$\begin{aligned}
G_{dmi} | \mathbf{a}, (\mathbf{K}_{dm})_{\cdot i} &\sim \mathcal{N}(G_{dmi}; \mathbf{a}^T (\mathbf{K}_{dm})_{\cdot i}, 1), \\
L_{di} | \mathbf{b}_{\cdot d}, \mathbf{G}_{d \cdot i} &\sim \mathcal{N}(L_{di}; \mathbf{b}_{\cdot d}^T \mathbf{G}_{d \cdot i}, 1).
\end{aligned}$$

For the latent variable and the class label:

$$\begin{aligned}
f_i | \mathbf{c}, e, \mathbf{L}_{\cdot i} &\sim \mathcal{N}(f_i; \mathbf{c}^T \mathbf{L}_{\cdot i} + e, 1), \\
y_i | f_i &\sim \delta(f_i y_i > \tau),
\end{aligned}$$

where τ is a given margin parameter relevant to the low density area of the distribution [15], $\text{Gamma}(\cdot; \alpha, \beta)$ denotes the gamma distribution with the mean $\alpha\beta$ and the variance $\alpha\beta^2$, $\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ represents the normal distribution, and $\delta(\cdot)$ represents the Kronecker delta function that equals 1 if the condition is satisfied and 0 otherwise. In the process of training the model, we can tune the hyper-parameters $(\alpha_\lambda, \beta_\lambda)$, $(\alpha_\eta, \beta_\eta)$, and $(\alpha_\gamma, \beta_\gamma)$ to control the sparsity of the sample, group and dimension weights, respectively.

3. INFERENCE

3.1 Inference using variational approximation

In this section, we apply variational approximation [6] to the estimation of parameters. There are some advantages. For instance, there are closed-form iterative formulas and therefore the computational efficiency of variational approximation is superior to that of Markov Chain Monte Carlo (MCMC). Denote $\Phi = \Lambda \cup \{\mathbf{G}, \mathbf{L}, \mathbf{f}\}$. The exact posterior distribution is $p(\Theta, \Phi | \mathbf{y}, \mathbf{K})$ and has the factorable ensemble approximation $q(\Theta, \Phi)$ [5] as follows:

$$\begin{aligned}
p(\Theta, \Phi | \mathbf{y}, \mathbf{K}) &\approx q(\Theta, \Phi) = q(\lambda)q(\mathbf{a})q(\mathbf{G})q(\eta)q(\mathbf{b})q(\mathbf{L})q(\gamma) \\
&\quad \cdot q(\omega)q(e, \mathbf{c})q(\mathbf{f}).
\end{aligned}$$

This factored form of $q(\Theta, \Phi)$ corresponds to mean field theory [20]. Then, the log marginal likelihood function can be calculated as

$$\begin{aligned}
\log p(\mathbf{y} | \mathbf{K}) &= \mathbb{E}_{q(\Theta, \Phi)} [\log p(\mathbf{y} | \mathbf{K})] \\
&= \int q(\Theta, \Phi) \left[\log \frac{p(\mathbf{y}, \Theta, \Phi | \mathbf{K})}{q(\Theta, \Phi)} \right. \\
&\quad \left. - \log \frac{p(\Theta, \Phi | \mathbf{y}, \mathbf{K})}{q(\Theta, \Phi)} \right] d\Theta d\Phi \\
&= \mathcal{L}(\Theta, \Phi) + KL(q || p),
\end{aligned}$$

where

$$\mathcal{L}(\Theta, \Phi) \triangleq \int q(\Theta, \Phi) \log \frac{p(\mathbf{y}, \Theta, \Phi | \mathbf{K})}{q(\Theta, \Phi)} d\Theta d\Phi$$

and

$$KL(q||p) \triangleq - \int q(\Theta, \Phi) \log \frac{p(\Theta, \Phi | \mathbf{y}, \mathbf{K})}{q(\Theta, \Phi)} d\Theta d\Phi.$$

Note that the Kullback-Leibler divergence $KL(q||p) \geq 0$, therefore the log marginal likelihood function satisfies

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{K}) &\geq \mathcal{L}(\Theta, \Phi) = \mathbb{E}_{q(\Theta, \Phi)}[\log p(\mathbf{y}, \Theta, \Phi | \mathbf{K})] \\ &\quad - \mathbb{E}_{q(\Theta, \Phi)}[\log q(\Theta, \Phi)]. \end{aligned}$$

Afterwards optimize the lower bound $\mathcal{L}(\Theta, \Phi)$ of the log marginal likelihood $\log p(\mathbf{y} | \mathbf{K})$. That means we need to solve an optimization problem as

$$\arg \max_{\Theta, \Phi} \mathcal{L}(\Theta, \Phi).$$

Then based on the calculus of variations [2], we know that each factor of $q(\Theta, \Phi)$ can be calculated from the following formula:

$$(6) \quad q(\cdot) \propto \exp(\mathbb{E}_{q(\Theta, \Phi)}[\log p(\mathbf{y}, \Theta, \Phi | \mathbf{K})]).$$

Under the conjugate assumptions, we can compute the approximate posterior distribution of each parameter. First update the parameters in Θ and Λ , then the intermediate output variables $\{\mathbf{G}, \mathbf{L}\}$ and finally the class label vector \mathbf{y} . The results are shown in Proposition 1.

Proposition 1. *Under the assumption (A1) and formula 6, the calculation results of the approximate posterior distributions are given below.*

$$\begin{aligned} \lambda &\sim \prod_{i=1}^N \text{Gamma}(\lambda_i; \alpha(\lambda), \beta(\lambda_i)), \\ \eta &\sim \prod_{d=1}^D \prod_{m=1}^{P_d} \text{Gamma}(\eta_{md}; \alpha(\eta), \beta(\eta_{md})), \\ \gamma &\sim \prod_{d=1}^D \text{Gamma}(\gamma_d; \alpha(\gamma), \beta(\gamma_d)), \\ \omega &\sim \text{Gamma}(\omega; \alpha(\omega), \beta(\omega)), \\ \mathbf{a} &\sim \mathcal{N}(\mathbf{a}; \boldsymbol{\mu}(\mathbf{a}), \boldsymbol{\Sigma}(\mathbf{a})), \\ \mathbf{b} &\sim \prod_{d=1}^D \mathcal{N}(\mathbf{b}_d; \boldsymbol{\mu}(\mathbf{b}_d), \boldsymbol{\Sigma}(\mathbf{b}_d)), \\ \begin{pmatrix} e \\ \mathbf{c} \end{pmatrix} &\sim \mathcal{N}\left(\begin{pmatrix} e \\ \mathbf{c} \end{pmatrix}; \boldsymbol{\mu}(e, \mathbf{c}), \boldsymbol{\Sigma}(e, \mathbf{c})\right), \\ \mathbf{G} &\sim \prod_{d=1}^D \prod_{i=1}^N \mathcal{N}(\mathbf{G}_{d,i}; \boldsymbol{\mu}(\mathbf{G}_{d,i}), \boldsymbol{\Sigma}(\mathbf{G}_{d,i})), \end{aligned}$$

$$\begin{aligned} \mathbf{L} &\sim \prod_{i=1}^N \mathcal{N}(\mathbf{L}_i; \boldsymbol{\mu}(\mathbf{L}_i), \boldsymbol{\Sigma}(\mathbf{L}_i)), \\ \mathbf{f} &\sim \prod_{i=1}^N \mathcal{TN}\left(f_i; \widetilde{\mathbf{c}}^T \widetilde{\mathbf{L}}_i + \tilde{e}, 1, f_i y_i > \tau\right), \end{aligned}$$

where

$$\begin{aligned} \alpha(\lambda) &= \alpha_\lambda + \frac{1}{2}, & \beta(\lambda_i) &= \left(\beta_\lambda^{-1} + \frac{\tilde{a}_i^2}{2}\right)^{-1}, \\ \alpha(\eta) &= \alpha_\eta + \frac{1}{2}, & \beta(\eta_{md}) &= \left(\beta_\eta^{-1} + \frac{\widetilde{b_{md}^2}}{2}\right)^{-1}, \\ \alpha(\gamma) &= \alpha_\gamma + \frac{1}{2}, & \beta(\gamma_d) &= \left(\beta_\gamma^{-1} + \frac{\widetilde{c_d^2}}{2}\right)^{-1}, \\ \alpha(\omega) &= \alpha_\omega + \frac{1}{2}, & \beta(\omega) &= \left(\beta_\omega^{-1} + \frac{\tilde{e}^2}{2}\right)^{-1}, \\ \boldsymbol{\mu}(\mathbf{a}) &= \boldsymbol{\Sigma}(\mathbf{a}) \sum_{d=1}^D \sum_{m=1}^{P_d} (\mathbf{K}_{dm}) \cdot \widetilde{\mathbf{G}}_{dm}^T, \\ \boldsymbol{\Sigma}(\mathbf{a}) &= \left(\text{diag}(\tilde{\lambda}) + \sum_{d=1}^D \sum_{m=1}^{P_d} (\mathbf{K}_{dm}) \cdot (\mathbf{K}_{dm})^T\right)^{-1}, \\ \boldsymbol{\mu}(\mathbf{b}_d) &= \boldsymbol{\Sigma}(\mathbf{b}_d) \widetilde{\mathbf{G}}_d \cdot \widetilde{\mathbf{L}}_d^T, \\ \boldsymbol{\Sigma}(\mathbf{b}_d) &= \left(\text{diag}(\tilde{\eta}_d) + \widetilde{\mathbf{G}}_d \cdot \widetilde{\mathbf{G}}_d^T\right)^{-1}, \\ \boldsymbol{\mu}(e, \mathbf{c}) &= \boldsymbol{\Sigma}(e, \mathbf{c}) \begin{pmatrix} \mathbb{1}^T \tilde{\mathbf{f}} \\ \widetilde{\mathbf{L}} \tilde{\mathbf{f}} \end{pmatrix}, \\ \boldsymbol{\Sigma}(e, \mathbf{c}) &= \begin{pmatrix} \tilde{\omega} + N & \mathbb{1}^T \widetilde{\mathbf{L}}^T \\ \widetilde{\mathbf{L}} \mathbb{1} & \text{diag}(\tilde{\gamma}) + \widetilde{\mathbf{L}} \widetilde{\mathbf{L}}^T \end{pmatrix}^{-1}, \\ \boldsymbol{\mu}(\mathbf{G}_{d,i}) &= \boldsymbol{\Sigma}(\mathbf{G}_{d,i}) \begin{pmatrix} (\mathbf{K}_{d1})_i \\ \vdots \\ (\mathbf{K}_{dP_d})_i \end{pmatrix} \tilde{\mathbf{a}} + \widetilde{\mathbf{b}}_d \widetilde{\mathbf{L}}_{di}, \\ \boldsymbol{\Sigma}(\mathbf{G}_{d,i}) &= \left(\mathbf{I} + \widetilde{\mathbf{b}}_d \widetilde{\mathbf{b}}_d^T\right)^{-1}, \\ \boldsymbol{\mu}(\mathbf{L}_i) &= \boldsymbol{\Sigma}(\mathbf{L}_i) (\tilde{\mathcal{G}} + \widetilde{\mathbf{c}} \tilde{f}_i - \tilde{c} \tilde{e}), \\ \boldsymbol{\Sigma}(\mathbf{L}_i) &= \left(\mathbf{I} + \widetilde{\mathbf{c}} \widetilde{\mathbf{c}}^T\right)^{-1}, \\ \tilde{\mathcal{G}} &\triangleq \begin{pmatrix} \widetilde{\mathbf{b}}_{1\cdot}^T \widetilde{\mathbf{G}}_{1\cdot i} \\ \vdots \\ \widetilde{\mathbf{b}}_{D\cdot}^T \widetilde{\mathbf{G}}_{D\cdot i} \end{pmatrix}, \end{aligned}$$

where $\widetilde{(\cdot)}$ represents $\mathbb{E}(\cdot)$, the expectation operator. $\mathbb{1}$ is the all-one vector $(1, 1, \dots, 1)^T$, \mathbf{I} represents identity matrix and \mathcal{TN} is the truncated normal distribution.

Proof. The likelihood function and the approximate posterior distribution are

$$\begin{aligned} p(\mathbf{y}, \Theta, \Phi | \mathbf{K}) &= p(\boldsymbol{\lambda}) p(\mathbf{a} | \boldsymbol{\lambda}) p(\mathbf{G} | \mathbf{a}, \mathbf{K}) p(\boldsymbol{\eta}) p(\mathbf{b} | \boldsymbol{\eta}) p(\mathbf{L} | \mathbf{b}, \mathbf{G}) \\ &\quad \cdot p(\boldsymbol{\gamma}) p(\mathbf{c} | \boldsymbol{\gamma}) p(\boldsymbol{\omega}) p(e | \boldsymbol{\omega}) p(\mathbf{f} | \mathbf{c}, e, \mathbf{L}) p(\mathbf{y} | \mathbf{f}), \\ q(\Theta, \Phi) &= q(\boldsymbol{\lambda}) q(\mathbf{a}) q(\mathbf{G}) q(\boldsymbol{\eta}) q(\mathbf{b}) q(\mathbf{L}) q(\boldsymbol{\gamma}) q(\boldsymbol{\omega}) q(e, \mathbf{c}) q(\mathbf{f}). \end{aligned}$$

and each factor $q(\cdot)$ of the $q(\Theta, \Phi)$ can be calculated by formula 6. So we can first compute the approximate posterior distribution of the prior $\boldsymbol{\lambda}$ as below.

$$\begin{aligned} q(\boldsymbol{\lambda}) &\propto \exp\{E_{q(\Theta, \Phi | \boldsymbol{\lambda})} [\log p(\mathbf{y}, \Theta, \Phi | \mathbf{K})]\} \\ &\propto \exp\{\log p(\boldsymbol{\lambda}) + E_{q(\mathbf{a})} [\log p(\mathbf{a} | \boldsymbol{\lambda})]\} \\ &\propto \prod_{i=1}^N \left\{ \lambda_i^{(\alpha_\lambda + \frac{1}{2}) - 1} e^{-\left(\beta_\lambda^{-1} + \frac{E(a_i^2)}{2}\right) \lambda_i} \right\}. \end{aligned}$$

So far we obtain the approximate posterior distribution

$$q(\boldsymbol{\lambda}) = \prod_{i=1}^N \text{Gamma} \left(\lambda_i; \alpha_\lambda + \frac{1}{2}, \left(\frac{1}{\beta_\lambda} + \frac{E a_i^2}{2} \right)^{-1} \right).$$

Similarly, we can obtain $q(\boldsymbol{\eta})$, $q(\boldsymbol{\gamma})$ and $q(\boldsymbol{\omega})$.

Second, the approximate posterior distribution of the sample weight parameter can be calculated as

$$\begin{aligned} q(\mathbf{a}) &\propto \exp\{E_{q(\Theta, \Phi | \mathbf{a})} [\log P(\mathbf{a} | \boldsymbol{\lambda}) p(\mathbf{G} | \mathbf{a}, \mathbf{K})]\} \\ &\propto \exp \left\{ -\frac{1}{2} \left[\mathbf{a}^T \left(\text{diag}(\tilde{\boldsymbol{\lambda}}) + \sum_{d=1}^D \sum_{m=1}^{P_d} (\mathbf{K}_{dm}) \cdot (\mathbf{K}_{dm})^T \right) \right. \right. \\ &\quad \left. \left. \mathbf{a} - \mathbf{a}^T \sum_{d=1}^D \sum_{m=1}^{P_d} (\mathbf{K}_{dm}) \cdot \widetilde{\mathbf{G}}_{dm}^T \right. \right. \\ &\quad \left. \left. - \sum_{d=1}^D \sum_{m=1}^{P_d} \widetilde{\mathbf{G}}_{dm} \cdot (\mathbf{K}_{dm})^T \mathbf{a} \right] \right\}. \end{aligned}$$

From the formula of $q(\mathbf{a})$ we know that the calculation process can update the mean and covariance matrix simultaneously which can reveal the correlation information of the samples. Through similar calculation, the approximate posterior distributions of \mathbf{b} and (e, \mathbf{c}) can be obtained as

$$\begin{aligned} q(\mathbf{b}_{\cdot d}) &\propto \exp \left\{ -\frac{1}{2} \left[\mathbf{b}_{\cdot d}^T \left(\text{diag}(\widetilde{\boldsymbol{\eta}}_{\cdot d}) + \widetilde{\mathbf{G}}_{d \cdot} \cdot \widetilde{\mathbf{G}}_{d \cdot}^T \right) \mathbf{b}_{\cdot d} \right. \right. \\ &\quad \left. \left. - \mathbf{b}_{\cdot d}^T \widetilde{\mathbf{G}}_{d \cdot} \cdot \widetilde{\mathbf{L}}_{d \cdot}^T - \widetilde{\mathbf{L}}_{d \cdot} \cdot \widetilde{\mathbf{G}}_{d \cdot}^T \mathbf{b}_{\cdot d} \right] \right\}, \quad d = 1, \dots, D, \end{aligned}$$

$$\begin{aligned} q(e, \mathbf{c}) &\propto \exp \left\{ -\frac{1}{2} \left[\begin{pmatrix} e \\ \mathbf{c} \end{pmatrix}^T \right. \right. \\ &\quad \left. \left(\begin{array}{cc} \tilde{\omega} + N & \mathbf{1}^T \widetilde{\mathbf{L}}^T \\ \widetilde{\mathbf{L}} \mathbf{1} & \text{diag}(\tilde{\boldsymbol{\gamma}}) + \widetilde{\mathbf{L}} \widetilde{\mathbf{L}}^T \end{array} \right) \begin{pmatrix} e \\ \mathbf{c} \end{pmatrix} \right. \\ &\quad \left. \left. - \begin{pmatrix} e \\ \mathbf{c} \end{pmatrix}^T \left(\begin{array}{c} \mathbf{1}^T \tilde{\mathbf{f}} \\ \widetilde{\mathbf{L}} \tilde{\mathbf{f}} \end{array} \right) - \left(\begin{array}{c} \mathbf{1}^T \tilde{\mathbf{f}} \\ \widetilde{\mathbf{L}} \tilde{\mathbf{f}} \end{array} \right)^T \begin{pmatrix} e \\ \mathbf{c} \end{pmatrix} \right] \right\}. \end{aligned}$$

Third, the approximate posterior distributions of the intermediate outputs \mathbf{G} and \mathbf{L} are

$$\begin{aligned} q(\mathbf{G}) &= \prod_{d=1}^D \prod_{i=1}^N \mathcal{N} \left(\mathbf{G}_{d \cdot i}; \boldsymbol{\Sigma}(\mathbf{G}_{d \cdot i}) \left(\begin{pmatrix} (\mathbf{K}_{d1})_{i \cdot} \\ \vdots \\ (\mathbf{K}_{dP_d})_{i \cdot} \end{pmatrix} \tilde{\mathbf{a}} + \widetilde{\mathbf{b}}_{\cdot d} \widetilde{\mathbf{L}}_{di} \right), \right. \\ &\quad \left. \left(\mathbf{I} + \widetilde{\mathbf{b}}_{\cdot d} \widetilde{\mathbf{b}}_{\cdot d}^T \right)^{-1} \right), \end{aligned}$$

$$q(\mathbf{L}) = \prod_{i=1}^N \mathcal{N} \left(\mathbf{L}_{\cdot i}; \boldsymbol{\Sigma}(\mathbf{L}_{\cdot i}) (\tilde{\mathbf{G}} + \tilde{\mathbf{c}} \tilde{\mathbf{f}}_i - \tilde{\mathbf{c}} e), \left(\mathbf{I} + \widetilde{\mathbf{c}} \widetilde{\mathbf{c}}^T \right)^{-1} \right).$$

Through the update of the parameters above, we obtain the following approximate posterior distribution of \mathbf{f} .

$$q(\mathbf{f}) = \prod_{i=1}^N \mathcal{TN} \left(f_i; \widetilde{\mathbf{c}}^T \widetilde{\mathbf{L}}_{\cdot i} + \tilde{e}, 1, f_i y_i > \tau \right). \quad \square$$

Under the framework of approximate inference [6], the convergence criterion of the algorithm is summarized in Proposition 2.

Proposition 2. *The update of parameters converges by achieving the local maximum value of the lower bound as below.*

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{K}) &\geq E_{q(\Theta, \Phi)} [\log p(\mathbf{y}, \Theta, \Phi | \mathbf{K})] - E_{q(\Theta, \Phi)} [\log q(\Theta, \Phi)]. \end{aligned}$$

Algorithm 1 summarizes the calculation process of updating approximate posterior distributions.

3.2 Classification

Given a new data point \mathbf{x}_* , denote the kernel matrix set as \mathbf{K}^* . According to the training parameters, the predictive distribution of the first layer $\mathbf{G}_{d \cdot *}$ of intermediate output can be obtained by replacing the posterior distribution of \mathbf{a} with its approximation $q(\mathbf{a})$ as follows:

$$P(\mathbf{G}_{d \cdot *} | \mathbf{K}^*, \mathbf{a}, \mathbf{K}, \mathbf{y}) = \prod_{m=1}^{P_d} \mathcal{N}(\mathbf{G}_{dm*}; \boldsymbol{\mu}(\mathbf{G}_{dm*}), \boldsymbol{\Sigma}(\mathbf{G}_{dm*})),$$

Algorithm 1 Calculation of approximate posterior distributions in BKAGL

Input: the kernel matrix set \mathbf{K} and the label vector \mathbf{y} of the training set, the hyper-parameter set Ξ , the iteration number $iter$;

- 1: initial parameters of distributions of $\{\Theta, \Phi\}$;
- 2: **for** $t = 1, \dots, iter$ **do**
- 3: compute $\beta^t(\lambda_i)$ and related parameters;
- 4: compute $\boldsymbol{\mu}^t(\mathbf{a})$, $\boldsymbol{\Sigma}^t(\mathbf{a})$ and related parameters;
- 5: compute $\boldsymbol{\mu}^t(\mathbf{G}_{d.i})$, $\boldsymbol{\Sigma}^t(\mathbf{G}_{d.i})$ and related parameters;
- 6: compute $\beta^t(\eta_{md})$ and related parameters;
- 7: compute $\boldsymbol{\mu}^t(\mathbf{b}_d)$, $\boldsymbol{\Sigma}^t(\mathbf{b}_d)$ and related parameters;
- 8: compute $\boldsymbol{\mu}^t(\mathbf{L}_{.i})$, $\boldsymbol{\Sigma}^t(\mathbf{L}_{.i})$ and related parameters;
- 9: compute $\beta^t(\gamma_d)$ and related parameters;
- 10: compute $\beta^t(\omega)$ and related parameters;
- 11: compute $\boldsymbol{\mu}^t(e, c)$, $\boldsymbol{\Sigma}^t(e, c)$ and related parameters;
- 12: compute $\boldsymbol{\mu}^t(\mathbf{f})$ and related parameters;
- 13: **end for**

Output: parameters of the approximate posterior distributions of $\{\Theta, \Phi\}$.

where

$$\boldsymbol{\mu}(G_{dm*}) = \boldsymbol{\mu}^T(\mathbf{a})(\mathbf{K}_{dm})_{.*}$$

and

$$\boldsymbol{\Sigma}(G_{dm*}) = 1 + (\mathbf{K}_{dm})_{.*}^T \boldsymbol{\Sigma}(\mathbf{a})(\mathbf{K}_{dm})_{.*}$$

Similarly, the predictive distribution of the second layer L_{d*} of intermediate output is

$$\begin{aligned} P(L_{d*} | \mathbf{G}_{d*}, \mathbf{b}_d) \\ = \mathcal{N}(L_{d*}; \boldsymbol{\mu}^T(\mathbf{b}_d) \mathbf{G}_{d*}, 1 + \mathbf{G}_{d*}^T \boldsymbol{\Sigma}(\mathbf{b}_d) \mathbf{G}_{d*}). \end{aligned}$$

Then the predictive distribution of f_* can be formulated as

$$P(f_* | \mathbf{L}_{.*}, c, e) = \mathcal{N}(f_*; \boldsymbol{\mu}(f_*), \boldsymbol{\Sigma}(f_*)),$$

where

$$\boldsymbol{\mu}(f_*) = \boldsymbol{\mu}^T(e, c) \begin{pmatrix} 1 \\ \mathbf{L}_{.*} \end{pmatrix}$$

and

$$\boldsymbol{\Sigma}(f_*) = 1 + (1 \quad \mathbf{L}_{.*}^T) \boldsymbol{\Sigma}(e, c) \begin{pmatrix} 1 \\ \mathbf{L}_{.*} \end{pmatrix}.$$

Moreover, given the margin parameter τ , we can obtain the predictive distribution of label y_* as given in Proposition 3.

Proposition 3. *The predictive distribution of y_* for \mathbf{x}_* is*

$$\begin{aligned} P(y_* = +1 | \mathbf{K}^*, \mathbf{K}, \mathbf{y}) &= Z_*^{-1} \Phi \left(\frac{\boldsymbol{\mu}(f_*) - \tau}{\boldsymbol{\Sigma}(f_*)} \right), \\ P(y_* = -1 | \mathbf{K}^*, \mathbf{K}, \mathbf{y}) &= Z_*^{-1} \Phi \left(\frac{-\boldsymbol{\mu}(f_*) - \tau}{\boldsymbol{\Sigma}(f_*)} \right), \end{aligned}$$

Algorithm 2 Prediction of the label y_* for a new point \mathbf{x}_*

Input: the kernel vector set \mathbf{K}^* of the new point, the training parameters obtained from Algorithm 1;

- 1: compute $p(\mathbf{G}_{d*} | \mathbf{K}^*, \mathbf{a}, \mathbf{K}, \mathbf{y})$;
- 2: compute $p(L_{d*} | \mathbf{G}_{d*}, \mathbf{b})$;
- 3: compute $p(f_* | \mathbf{L}_{.*}, c, e)$;
- 4: compute $p(y_* = +1 | \mathbf{K}^*, \mathbf{K}, \mathbf{y})$ and $p(y_* = -1 | \mathbf{K}^*, \mathbf{K}, \mathbf{y})$;
- 5: **if** $p(y_* = +1 | \mathbf{K}^*, \mathbf{K}, \mathbf{y}) \geq 0.5$ **then**
- 6: label = +1;
- 7: **else**
- 8: label = -1;
- 9: **end if**

Output: the label of \mathbf{x}_* .

where $Z_* = \Phi(\frac{\boldsymbol{\mu}(f_*) - \tau}{\boldsymbol{\Sigma}(f_*)}) + \Phi(\frac{-\boldsymbol{\mu}(f_*) - \tau}{\boldsymbol{\Sigma}(f_*)})$ is the normalization coefficient and $\Phi(\cdot)$ represents the cumulative distribution function of the standard normal distribution.

According to Proposition 3, we can calculate the classification probability of the test instance by Algorithm 2. Note that Algorithm 2 distinguishes the label by 0.5. That means 0.5 acts as the cut point. In fact, other cut points can also be selected through cross-validation techniques if needed.

4. EXPERIMENTS

In this section, we provide experiments on simulated data and real data to illustrate the performance of the proposed BKAGL method. We compare the proposed method with the k -nearest neighbors (KNN) algorithm, support vector machine (SVM) algorithm, the generalized linear model (GLM) with lasso, random forest and naive Bayes. For these contrastive classification methods, matrix data should be vectorized into a vector pattern before classification. In other words, a $D \times R$ matrix should be transformed into a DR -dimensional vector for each instance. Select $k = 4$ for KNN through cross-validation. The generalized linear model we use is the logistic regression and we use k -fold cross-validation for GLM with lasso. With regard to BKAGL, apply the k -medoids algorithm for grouping. To reflect the advantage of grouping, we also compare the proposed method with BKAGL without grouping. For all the kernel selection problems involved in this section, such as SVM, BKAGL and BKAGL without grouping, we choose the linear kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$ uniformly for a fair comparison. Other kernels, including the polynomial kernel and the Gaussian kernel, can be used, but parameter tuning is computationally expensive. Hence, we choose the linear kernel. Through the sensitivity analysis of the hyper-parameters, we find that the proposed method is not sensitive to the choice of the hyper-parameters. Therefore, we set the hyper-parameters $(\alpha_\lambda, \beta_\lambda, \alpha_\eta, \beta_\eta, \alpha_\gamma, \beta_\gamma, \alpha_\omega, \beta_\omega) = (1, 1, 1, 1, 1, 1, 1, 1)$ throughout the experiment study. We implement the proposed BKAGL method in R and the codes are available at <https://github.com/wangxz021/bkagl>.

4.1 Simulated data

In this part, we discuss the classification of two types of multi-dimensional datasets. We first introduce the generation models, and then give the comparison of the classification results. Finally, we analyze the interpretation of grouping and the learning parameters \mathbf{c} and \mathbf{b} .

Consider the classification problem of the following two models. Each data point \mathbf{X} has D dimensions and each dimension is a vector with R features.

- Model 1: the label $Y = +1$, \mathbf{X} is a $D \times R$ matrix:
 $X_{dr} \sim \mathcal{N}(0, 1), 1 \leq d \leq D, 1 \leq r \leq R.$
- Model 2: the label $Y = -1$, \mathbf{X} is a $D \times R$ matrix:
 $X_{1r} \sim \mathcal{N}(h, 1), (1-p)R + 1 \leq r \leq R, 0 \leq p \leq 1;$
 $X_{dr} \sim \mathcal{N}(0, 1), \textit{otherwise}.$

From the construction, the two generation models only differ in the last pR features of the first dimension. The parameters p and h decide the difference between the two models. The parameter p denotes the proportion of the difference and h determines the size of the difference. When p and h are large, classification is obviously easier. On the contrary, small p and h result in smaller differences, which makes classification more difficult. On the other hand, the classification method needs to be sensitive to the fluctuation due to the small sample size.

Generate N instances, among which half of them are from the first model and the rest are from the second model. In the classification process, 80% of instances are selected randomly from all of the instances as the training set and the remaining 20% form the test set. The probability of successful classification in the test set is used to measure the performance of each method. The results reported are the average of 100 independent runs.

First, take $N = 100$, $D = 2$, $R = 100$ and fix $p = 50\%$, then compare the seven algorithms on different h chosen from $\{0.1 \times i, 3 \leq i \leq 10\}$. For the sake of comparable performance measure, the classification accuracy of every method is calculated on the same training set and test set. Repeat the computing process and average the results. Figure 3 displays the relation between the classification accuracies of the seven methods and the difference parameter h . To clearly show the results and the trend, the X-axis is reversed. From Figure 3, one can see that BKAGL has the best classification performance. When the difference parameter h is large, BKAGL has higher classification accuracies than GLM with lasso, random forest and KNN. When h becomes smaller, naive Bayes, SVM and BKAGL without grouping exhibit a faster drop and the superiority of BKAGL becomes more obvious. Meanwhile, BKAGL has the advantage of interpretation ability provided by grouping and the learning parameters, which will be introduced in more detail later.

Then we consider the impact of the difference ratio p on the classification results. The value of p decides the proportion of differences. We fix $N = 100$, $D = 2$, $R = 100$,

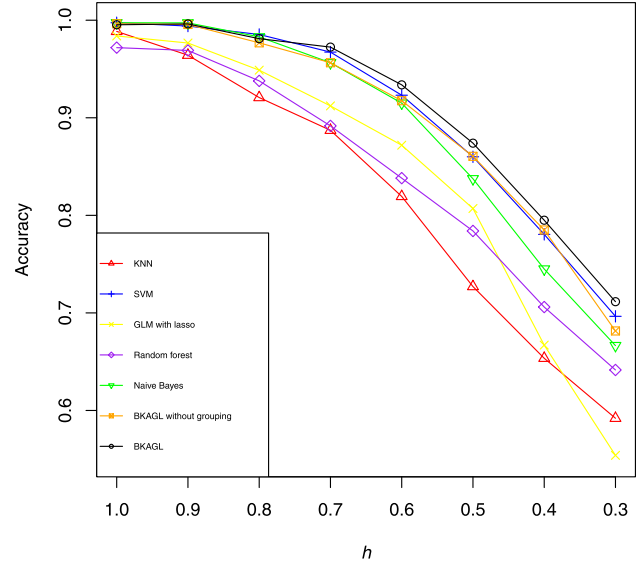


Figure 3. Classification accuracies of the seven methods under different h .

$h = 0.3$ and choose p from $\{10\%, 20\%, 30\%, 40\%, 50\%\}$. Repeat the calculation process and average the classification results. Table 1 shows the classification accuracies of the seven methods under different p . We observe from Table 1 that BKAGL performs constantly better than other methods under different p .

Next we discuss the learning parameters \mathbf{c} and \mathbf{b} . Choose $N = 100$, $D = 2$, $R = 100$, $h = 1$, $p = 50\%$. Through observing the learning parameters in each iteration, we find that the grouping and the analysis results are consistent. Therefore, we randomly select one classification process from the iteration to provide the detailed analysis. The results are recorded in Table 2. We observe that $c_1 = 0.5029$ and $c_2 = 0.0000$. The value of \mathbf{c} indicates that the classification result is only determined by the first dimension, which matches the simulation model. With regard to the learning parameter \mathbf{b} , whether all features are properly divided into several groups and the parameter \mathbf{b} of different groups in different dimensions correctly reflects the importance of partial features are two important things that need to be verified. From the construction of the simulation model, we know that the key to classification is the second half features of the first dimension. Therefore, the features of the first dimension should become a group. The grouping and learning parameters in Table 2 provide the explanatory results of the classification.

The grouping results of the first dimension successfully select most of the second half features to form Group 2. Compared to the unimportant Group 1, the learning parameter \mathbf{b} of Group 2 of the first dimension has a much higher value. The grouping of the second dimension looks random; nevertheless, the learning parameter \mathbf{b} of each group is 0,

Table 1. Test accuracies of classification under different p

p	10%	20%	30%	40%	50%
KNN	0.5115±0.1148	0.5338±0.1166	0.5390±0.1102	0.5907±0.0969	0.5920±0.1091
SVM	0.5310±0.1114	0.5981±0.1126	0.5910±0.0960	0.6648±0.0993	0.6965±0.1067
GLM with lasso	0.4635±0.1000	0.4644±0.1004	0.4830±0.1209	0.5456±0.1396	0.5540±0.1263
Random forest	0.5155±0.1275	0.5513±0.1046	0.5690±0.1134	0.6011±0.1103	0.6415±0.1078
Naive Bayes	0.5230±0.1401	0.5438±0.1192	0.5685±0.0926	0.6126±0.1297	0.6665±0.1113
BKAGL without grouping	0.5345±0.1121	0.5894±0.0986	0.6060±0.1018	0.6670±0.0981	0.6815±0.1012
BKAGL	0.5380±0.1146	0.6156±0.1174	0.6325±0.1106	0.7033±0.0933	0.7115±0.1054

Table 2. The grouping of features and learning parameters

d	c_d	$P_d = 2$	Real group	Experimental group	b_{md}
1	0.5029	Group 1	1-50	1-18,20-28,30-32,34-50,52,67,78,98	0.1339
		Group 2	51-100	19,29,33,51,53-66,68-77,79-97,99-100	0.4288
d	c_d	$P_d = 2$	Real group	Experimental group	b_{md}
2	0.0000	Group 1	None	1,4,8,11,12,13,14,16,17,18,20,21,22,24,25,28,29,30,32,33,34,36,37,38,39,41,42,45,52,53,54,57,59,60,63,65,67,70,71,72,74,75,77,78,80,82,85,87,88,89,94,95,97,98,99,100	0.0000
		Group 2	None	2,3,5,6,7,9,10,15,19,23,26,27,31,35,40,43,44,46,47,48,49,50,51,55,56,58,61,62,64,66,68,69,73,76,79,81,83,84,86,90,91,92,93,96	0.0000

which shows the lack of importance on the final classification result.

The grouping and the learning parameters \mathbf{c} and \mathbf{b} accurately capture the characteristics of the simulation model. For practical applications, valuable information can be obtained from the grouping and the learning parameters.

4.2 Real data

In this section, we compare the performance of BKAGL with classical classification methods on two real datasets: ECG and EEG. Meanwhile, we will analyze the learning parameters \mathbf{b} and \mathbf{c} based on the practical meaning. The results reported are the average of 10 independent runs.

4.2.1 ECG dataset

The electrocardiogram (ECG) dataset [4] contains 200 instances, where 67 instances are abnormal (label +1) and 133 instances are normal (label -1). The researcher collected heartbeat data through two electrodes. Therefore, the heartbeat record for each instance is a multivariate time series. For instances in the raw data, the time series of each electrode have unequal lengths. To solve this problem, we choose to linearly interpolate the time series of each electrode so that every time series can have the same length as the longest length in the raw data. This is a common tool for data pre-processing [3, 9]. Some researchers used this form of pre-processing to show that for time series classification problems, the length of the time series is not an issue [23]. After the procedure, the length of the time series of each electrode for every instance is 152.

Table 3. Test accuracies of the ECG dataset

Method	Accuracy
KNN	0.6563±0.0182
SVM	0.7944±0.0474
GLM with lasso	0.6806±0.0526
Random forest	0.7456±0.0405
Naive Bayes	0.7150±0.0928
BKAGL without grouping	0.7738±0.0340
BKAGL	0.8056±0.0387

We note that the dimensions in the model in Section 2.2 are two electrodes and the features are 152 time points in this problem. We randomly select 40 instances as the training set and the remaining 160 instances form the test set. The classification accuracy of each method is shown in Table 3. As we can see, the proposed BKAGL outperforms other classification methods.

Next, we observe the learning parameters \mathbf{c} and \mathbf{b} . We find that the grouping and analysis results are consistent under each iteration by observing each group of learning parameters. Consequently, we randomly select one process to give further analysis and the results are recorded in Table 4. The learning parameters c_d for the two dimensions (electrodes) are 0.4755 and 0.0000, respectively, which illustrates that the records on the first electrode have the most influence on the final output. On the other hand, Table 4 records the learning parameter \mathbf{b} for the two dimensions. For each dimension, 152 time points are divided into two groups: earlier-stage (about the first two-thirds of time points) and later-stage (about the last third of time points). For the

first dimension, the earlier-stage period has a much higher absolute value than the later-stage period. For the second dimension, the learning parameters b_{m2} of two groups are very close to 0. According to the analysis based on the learning parameters c and b , we know that the earlier-stage records of the first electrode play the most important role in the ECG classification problem.

4.2.2 EEG dataset

Next we consider the multiple electrode time series electroencephalography (EEG) dataset [8] of alcoholism (<http://kdd.ics.uci.edu/databases/eeg/eeg.data.html>). The study consists of 77 alcoholic subjects and 45 control subjects and contains voltage values which are measured from 64 electrodes placed on each subject’s scalp. The 64 electrodes were located at standard sites (Standard Electrode Position Nomenclature, American Electroencephalographic Association 1990) [28]. The voltage values are recorded at 256 time points. Every subject had 120 trials exposed to three types of stimuli, which are single stimulus, two matched stimuli and two unmatched stimuli. The goal of the study is to explore the associations between alcoholism and voltage values over both channels and time.

We focus on the data exposed to the single stimulus and average all the trials under the single stimulus for each subject [17]. For intuitive cognition of the dataset, we average

the voltage values of all alcoholic and control subjects respectively. Figure 4 displays the average EEG recordings of two groups. From Figure 4, we can see that the differences of two groups mainly lie in the medium-term and final performance of the EEG readings. Consider the medium-term and final data to find the relation among alcoholism, channels and the medium-term and final performance of the EEG readings.

Now the dimensions in the model in Section 2.2 are replaced by the 64 channels and the features are time points here. In other words, there are 64 time series data for each subject. According to the above analysis of Figure 4, the differences of two groups are mainly concentrated in the medium-term and final data. Hence, choose the 51st to 70th time points together with the 237th to 256th time points to be a time chain with length 40. Randomly select 100 subjects as the training set and the remaining 22 subjects form the test set. Calculate the classification accuracy of each method. Repeat the process and average the results. The classification results are shown in Table 5.

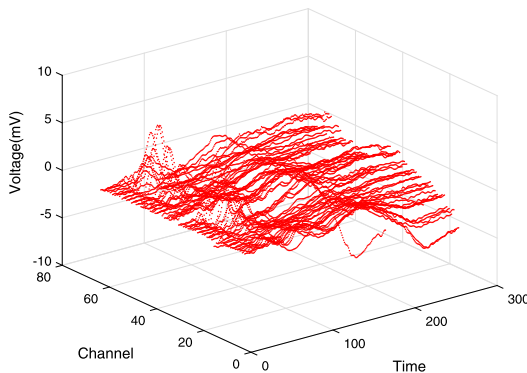
From Table 5 we know that BKAGL shows the best performance. Next we discuss the learning parameters c and b to find useful information. By observing each group of learning parameters under each iteration, we find the consistency of grouping and analysis results. Hence, randomly select one classification process from the iteration and give

Table 4. The grouping and learning parameters of the ECG dataset

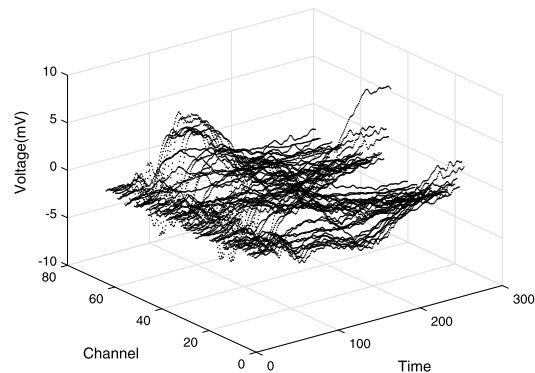
d	c_d	$P_d = 2$	Time points	b_{md}
1	0.4755	Group 1	1-99	-0.4115
		Group 2	100-152	-0.0013
d	c_d	$P_d = 2$	Time points	b_{md}
2	0.0000	Group 1	1-100	0.0000
		Group 2	101-152	0.0000

Table 5. Test accuracies of the EEG dataset

Method	Accuracy
KNN	0.6705±0.0776
SVM	0.7159±0.0227
GLM with lasso	0.6705±0.0937
Random forest	0.7386±0.0435
Naive Bayes	0.6591±0.0587
BKAGL without grouping	0.7159±0.0937
BKAGL	0.7614±0.0435



(a) The alcoholic group



(b) The control group

Figure 4. The average EEG recordings of two groups.

Table 6. The grouping and learning parameters of the EEG dataset

d	Channel name	c_d	$P_d = 3$	Time points	Period	b_{md}
55	PO7	1.1025	Group 1	51,52,53,54,55,56,57,58	Medium-term	0.2266
			Group 2	59,60,61,62,63,64,65,66,67,68,69,70	Medium-term	0.1974
			Group 3	237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256	Final-stage	0.1794

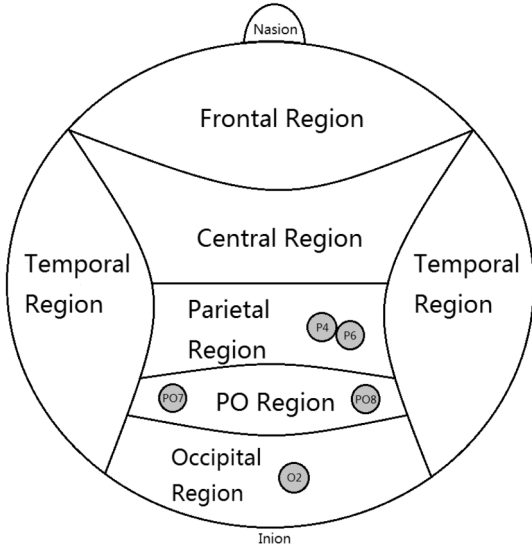


Figure 5. The locations of the top five channels on the scalp.

the following analysis. We obtain 64 c values for all channels. Different channels have different degrees of influence on the classification results. The order of the top five channels from high to low are PO7, P6, PO8, P4 and O2. According to electrode labelling, every placement site has a letter to identify the region of the brain. For example, P stands for Parietal and O stands for Occipital. PO means the intermediate electrode place between parietal and occipital. Odd and even numbers mean the left hemisphere and the right hemisphere, respectively. Each numeric value means a distance. Take PO7 as an example, number 7 is utilized for the left hemisphere to stand for 40% of the inion-to-nasion distance. We can observe from Figure 5 that the locations of the top five channels marked as grey circles are concentrated in the parietal region, occipital region and their intermediate region PO. Hence, the EEG readings of these regions are the key indices to classify the alcoholic group and the control group.

Finally, we discuss the learning parameter b . We note that the 55th channel PO7 has the significantly largest c value that equals 1.1025, which means PO7 has the greatest impact on the classification. Therefore, we consider the learning parameter b of the 55th channel PO7, which is shown in Table 6. All 40 time points are divided into three groups where the 20 medium-term data points are split into

two groups and the 20 final data points become a group independently. The parameters b of the two medium-term groups have larger values than that of the final-stage group, which means greater impacts on the final classification. More specifically, the parameter b of Group 1 has the largest value, which implies the importance of the time points in this group. On the whole, the medium-term of EEG readings has a higher influence on classification than the final performance of the EEG readings, which coincides with the rough observation of Figure 4.

The grouping and the learning parameters c and b can provide valuable information for practical applications. The proposed method BKAGL can be applied to various practical classification problems. The classification process is similar to the above analysis.

5. DISCUSSION AND EXTENSIONS

This paper proposes a novel kernel-based algorithm to solve the binary classification problem for multi-dimensional datasets. The proposed method can integrate information from different groups in different dimensions. Through the estimation of parameters, BKAGL has good interpretative ability, which enables us to find the groups of features and dimensions that have a greater impact on the response. Meanwhile, the utilization of variational approximation inference and conjugate Bayesian models provides the closed-form iterative formulas. For the multiclass classification problem, methods such as one-to-all [18] can be adopted. In addition, the correlation information between dimensions can be considered. In fact, if there is a priori knowledge about the dimensions, we can take advantage of it by designing a more reasonable prior distribution to improve the classification accuracy.

Received 1 August 2018

REFERENCES

- [1] ALTIPARMAK, F., FERHATOSMANOGLU, H., ERDAL, S. and TROST, D. C. (2006). Information mining over heterogeneous and high-dimensional time-series data in clinical trials databases. *IEEE Transactions on Information Technology in Biomedicine*, **10**(2), 254–263.
- [2] AUBERT, G. and KORNPORST, P. (2006). *Mathematical problems in image processing: partial differential equations and the calculus of variations*. Springer Science & Business Media, New York. [MR2244145](#)

- [3] BAGNALL, A., LINES, J., BOSTROM, A., LARGE, J. and KEOGH, E. (2017). The great time series classification bake off: a review and experimental evaluation of recent algorithmic advances. *Data Mining and Knowledge Discovery*, **31**(3), 606–660. [MR3640659](#)
- [4] BAYDOGAN, M. G. (2015). Multivariate time series classification datasets. <http://www.mustafabaydogan.com>.
- [5] BEAL, M. J. (2003). Variational algorithms for approximate Bayesian inference. PhD thesis, The Gatsby Computational Neuroscience Unit, University College London, London, The UK.
- [6] BISHOP, C. M. (2006). *Pattern recognition and machine learning*. Springer, New York. [MR2247587](#)
- [7] CAMACHO, D. M., COLLINS, K. M., POWERS, R. K., COSTELLO, J. C. and COLLINS, J. J. (2018). Next-generation machine learning for biological networks. *Cell*, **173**(7), 1581–1592.
- [8] DUA, D. and GRAFF, C. (2019). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.
- [9] FAWAZ, H. I., FORESTIER, G., WEBER, J., IDOUMGHAR, L. and MULLER, P. A. (2019). Deep learning for time series classification: a review. *Data Mining and Knowledge Discovery*, **33**(4), 917–963. [MR3962039](#)
- [10] FERREIRA, J. C. and MENEGATTO, V. A. (2009). Eigenvalues of integral operators defined by smooth positive definite kernels. *Integral Equations & Operator Theory*, **64**(1), 61–81. [MR2501172](#)
- [11] GÖNEN, M. (2012). Bayesian efficient multiple kernel learning. *Proceedings of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, The UK.
- [12] GÖNEN, M. and ALPAYDIN, E. (2011). Multiple kernel learning algorithms. *Journal of Machine Learning Research*, **12**, 2211–2268. [MR2825425](#)
- [13] HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction, second edition*. Springer, New York. [MR2722294](#)
- [14] HOFMANN, T., SCHÖLKOPF, B. and SMOLA, A. J. (2008). Kernel methods in machine learning. *The Annals of Statistics*, **36**(3), 1171–1220. [MR2418654](#)
- [15] LAWRENCE, N. D. and JORDAN, M. I. (2005). Semi-supervised learning via Gaussian processes. *Advances in Neural Information Processing Systems*, **17**, 753–760.
- [16] LECUN, Y., BENGIO, Y. and HINTON, G. (2015). Deep learning. *Nature*, **521**, 436–444.
- [17] LI, B., KIM, M. K. and ALTMAN, N. (2010). On dimension folding of matrix-or array-valued statistical objects. *The Annals of Statistics*, **38**(2), 1094–1121. [MR2604706](#)
- [18] LIU, Y. and ZHENG, Y. F. (2005). One-against-all multi-class SVM classification using reliability measures. *2005 IEEE International Joint Conference on Neural Networks*, **2**, 849–854.
- [19] MURPHY, K. P. (2012). *Machine learning: a probabilistic perspective*. MIT Press, Cambridge, MA.
- [20] PARISI, G. (1988). *Statistical field theory*. Addison-Wesley, Redwood City.
- [21] PARK, T. and CASELLA, G. (2008). The Bayesian lasso. *Journal of the American Statistical Association*, **103**(482), 681–686. [MR2524001](#)
- [22] PENG, J., ZHU, X., WANG, Y., AN, L. and SHEN, D. (2019). Structured sparsity regularized multiple kernel learning for Alzheimer’s disease diagnosis. *Pattern Recognition*, **88**, 370–382.
- [23] RATANAMAHATANA, C. A. and KEOGH, E. (2005). Three myths about dynamic time warping data mining. *Proceedings of the 2005 SIAM International Conference on Data Mining*, 506–510.
- [24] SONNENBURG, S., RÄTSCHE, G., SCHÄFER, C. and SCHÖLKOPF, B. (2006). Large scale multiple kernel learning. *Journal of Machine Learning Research*, **7**, 1531–1565. [MR2274416](#)
- [25] TIPPING, M. E. (2000). The relevance vector machine. *Advances in Neural Information Processing Systems*, **12**, 652–658.
- [26] TIPPING, M. E. (2001). Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, **1**, 211–244. [MR1875838](#)
- [27] VAPNIK, V. N. (2000). *The nature of statistical learning theory, second edition*. Springer-Verlag, New York. [MR1719582](#)
- [28] ZHANG, X. L., BEGLEITER, H., PORJESZ, B., WANG, W. and LITKE, A. (1995). Event related potentials during object recognition tasks. *Brain Research Bulletin*, **38**(6), 531–538.
- [29] ZHU, B., SONG, N., SHEN, R., ARORA, A., MACHIELA, M.J., SONG, L., LANDI, M.T., GHOSH, D., CHATTERJEE, N., BALADANDAYUTHAPANI, V. and ZHAO, H. (2017). Integrating clinical and multiple omics data for prognostic assessment across human cancers. *Scientific Reports*, **7**(1), 1–13.

Xiaozhou Wang
 School of Mathematical Sciences
 Shanghai Jiao Tong University
 Shanghai 200240
 China
 E-mail address: wangxiaozhou@sjtu.edu.cn

Fangli Dong
 School of Mathematical Sciences
 SJTU-Yale Joint Center of Biostatistics and Data Science
 Shanghai Jiao Tong University
 Shanghai 200240
 China
 E-mail address: dongfl@sjtu.edu.cn