Multinomial Dirichlet Gaussian Process Model for Classification of Multidimensional Data

Wanhyun Cho, Soonja Kang, Sangkyoon Kim, Soonyoung Park

Abstract—We present probabilistic multinomial Dirichlet classification model for multidimensional data and Gaussian process priors. Here, we have considered efficient computational method that can be used to obtain the approximate posteriors for latent variables and parameters needed to define the multiclass Gaussian process classification model. We first investigated the process of inducing a posterior distribution for various parameters and latent function by using the variational Bayesian approximations and important sampling method, and next we derived a predictive distribution of latent function needed to classify new samples. The proposed model is applied to classify the synthetic multivariate dataset in order to verify the performance of our model. Experiment result shows that our model is more accurate than the other approximation methods.

Keywords—Multinomial dirichlet classification model, Gaussian process priors, variational Bayesian approximation, Importance sampling, approximate posterior distribution, Marginal likelihood evidence.

I. INTRODUCTION

GAUSSIAN PROCESSES (GPs) are natural generations of multivariate Gaussian random variables to infinite index sets. Gaussian process is also a process that is constructed from classical statistical models by replacing latent functions of parametric form by random processes with Gaussian prior. GPs have been applied in a large number of fields in machine learning. One of these fields is the Gaussian regression and classification problem. In the case of regression with Gaussian noise, inference can be done simply in closed from, since the posterior is also a GP. However, in the case of classification, exact inference is analytically intractable because non-Gaussian likelihood function is assumed [1]-[3].

One prolific line of attack is based on approximating the non-Gaussian posterior with a tractable Gaussian distribution [4]-[6]. The most popular technique for treating intractable probabilistic models in these areas is the variational Bayesian approximation [7]. Thus, method approximated an intractable probability distribution by the closest distribution within a tractable family, where closeness is defined by the Kullback -Leibler divergence. In most applications, the tractable families contain distributions which factorize in all or in tractable subgroups of random variables. Hence, the method neglects

Wanhyun Cho (professor) is with the Chonnam National University, Gwangju, 61186 South Korea (corresponding author, phone: +82-62-530-3443; fax +82-62-530-3449; e-mail whcho@ chnnam.ac.kr).

Soonja Kang (professor) is with the Chonnam National University, Gwangju, and 61186 South Korea (e-mail: sjkang@chonnam.ac.kr).

Sangkyoon Kim (researcher) and Soonyoung Park (professor) are with the Electrical Engineering Department, Mokpo National University, Chonnam, 58554 South Korea (e-mail: narciss76@moipo.ac.kr, sypark@mokpo.ac.kr).

correlations between variables which may be crucial in the learning of the hyper-parameters. However, despite these problems, empirical comparisons with exact analysis via MCMC and Laplace approximations illustrate the utility of the variational approximation as a computationally economic alternative to full MCMC and it is shown to be more accurate than the Laplace approximation.

Here, we will demonstrate that the variational Bayesian multinomial Dirichlet Gaussian process approach is efficient to classify multivariate data with Gaussian process priors. Our results have also motivated by the inclusion of the variational Gaussian approach within a larger study of different methods for classification with Gaussian processes.

II. GAUSSIAN PROCESS CLASSIFICATION MODEL

A. Our Model

Here, we are going to consider the multinomial Dirichlet Gaussian Process Classification Model (MDGPCM) defined such as. First, we have the data matrix as $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ which has dimension $(N \times D)$. Here, we represent the $(N \times K)$ dimensional matrix of Gaussian process latent variables as \mathbf{F} . And we also represent the $(N \times 1)$ dimensional columns of \mathbf{F} and the $(K \times 1)$ dimensional rows of \mathbf{F} as $\mathbf{f}^k, k = 1, \dots, K$ and $\mathbf{f}_n, n = 1, \dots, N$ respectively. Then, we assume that a GP prior for the *k* -th latent vector function $\mathbf{f}^k(\mathbf{x})$ is defined as the Gaussian distribution with zero mean vector and the $(N \times N)$ dimensional covariance matrix \mathbf{K}^k . That is, $\mathbf{f}^k \sim G(\mathbf{0}, \mathbf{K}^k)$. Here, the (i, j)-th element k_{ij}^k of the covariance matrix \mathbf{K}^k will be defined by:

$$\mathbf{K}(\mathbf{X}, \mathbf{\phi}^{k}) = (k_{ij}^{k}), \ k_{ij}^{k} = \exp(-\frac{1}{2}\sum_{d=1}^{D}\varphi_{d}^{k}(x_{id} - x_{jd})^{2})$$
(1)
$$i, j = 1, \cdots, n$$

The $(D \times 1)$ vector of covariance kernel parameters for each class is denoted by $\boldsymbol{\varphi}^k$ and associated hyper-parameters $\boldsymbol{\psi}^k$ and hyper-hyper parameters $(\boldsymbol{\tau}^k, \boldsymbol{\zeta}^k)$ complete the model.

Moreover, we define the $(NK \times 1)$ dimensional vector $\mathbf{f}(\mathbf{x})$ of the latent variables for *K* classes classification as:

$$\mathbf{f}(\mathbf{x}) = (\mathbf{f}^1, \cdots, \mathbf{f}^k, \cdots, \mathbf{f}^K)^T,$$

$$\mathbf{f}^k = (\mathbf{f}^k_1, \cdots, \mathbf{f}^k_n, \cdots, \mathbf{f}^k_N)^T, k = 1, \cdots, K, \qquad (2)$$

where the superscript k denotes a particular class and the subscript n denotes the observation number. Then, the GP prior of latent vector $\mathbf{f}(\mathbf{x})$ for K classes classification has usually been chosen to have only intra-class correlations. The covariance matrix for the prior of latent function $\mathbf{f}(\mathbf{x})$ is defined as:

$$\mathbf{K} = \begin{pmatrix} \mathbf{K}(\mathbf{X}, \boldsymbol{\varphi}^{1}) & \cdots & \mathbf{0} \\ \vdots & \mathbf{K}(\mathbf{X}, \boldsymbol{\varphi}^{k}) & \vdots \\ \mathbf{0} & \cdots & \mathbf{K}(\mathbf{X}, \boldsymbol{\varphi}^{k}) \end{pmatrix}$$
(3)

Second, we define the $(N \times K)$ dimensional response matrix of associated target values as $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_N)^T$, where each component y_n^k of $(K \times 1)$ dimensional row vector $\mathbf{y}_n = (y_n^1, \dots, y_n^K)^T$ takes zero or one, and satisfies the condition $\sum_{k=1}^{k} y_{n}^{k} = 1 \text{ for all } n \text{ . And we assume that a } (K \times 1) \text{ dimensional}$ row vector $\mathbf{y}_n = (y_n^1, \dots, y_n^K)^T$ is distributed with one observation multinomial distribution with parameter vector $\boldsymbol{\pi}_n = (\pi_n^1, \cdots, \pi_n^K)^T$. That is, the $\pi_n^k, k = 1, \cdots, K$ denotes the probability that the n-th observation \mathbf{x}_n belongs to the k-th class. Moreover, we define the all of auxiliary random variables $\pi_n^k, k = 1, \dots, K, n = 1, \dots, N$ as a $(N \times K)$ dimensional random matrix $\boldsymbol{\pi} = (\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_N)^T$ with same size of response matrix **Y**. We also assume that the prior distribution of n-th parameter vector $\boldsymbol{\pi}_n = (\boldsymbol{\pi}_n^1, \cdots, \boldsymbol{\pi}_n^K)^T$ of the $(N \times K)$ dimensional random matrix $\boldsymbol{\pi} = (\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_N)^T$ is given as the Dirichlet distribution with K parameter vector $\boldsymbol{\alpha}_n = (\alpha_n^1, \cdots, \alpha_n^K)$.

Third, we have to consider the link function that specifies the relation between the latent function $\mathbf{f}(\mathbf{x})$ and the response mean vector $E(\mathbf{Y} | \mathbf{f})$. To drive this relationship, we have first considered the $(N \times K)$ dimensional random matrix \mathbf{M} of another auxiliary random variables $m_n^k, k = 1, \dots, K, n = 1, \dots, N$. Here, we have also assumed that each auxiliary random variable m_n^k is distributed with Gaussian $N(\mathbf{f}_n^k, 1)$. And we define *K*-dimensional parameter vector $\mathbf{a}_n = (\alpha_n^1, \dots, \alpha_n^K)$ of Dirichlet distribution $Dir(\mathbf{\pi}_n | \mathbf{a}_n)$ as *K* parameters $\mathbf{a}_n = (\exp(m_n^1), \dots, \exp(m_n^K))$. Then, the relationship between the response row vector $\mathbf{y}_n = (y_n^1, \dots, y_N^K)^T$ and the auxiliary latent vector $\mathbf{\pi}_n = (\mathbf{\pi}_n^1, \dots, \mathbf{\pi}_n^K)^T$ is adopted as the following form:

$$y_n^k = 1$$
 if $\pi_n^k = \max_{1 \le j \le K} \{\pi_n^j\}$ (4)

Moreover, by using an expectation property of Dirichlet distribution and the softmax method, the probability that a response variable y_n^k takes 1 can be defined by:

$$p(y_n^k = 1 | \mathbf{m}_n) = E(\pi_n^k | \mathbf{a}_n = \exp(\mathbf{m}_n)) = \frac{\exp(m_n^k)}{\sum_{l=1}^K \exp(m_n^l)}$$

$$k = 1, \dots, K, \ n = 1, \dots, N.$$
(5)

Therefore, since the random variables m_n^k is distributed with the Gaussian distribution $N(\mathbf{f}_n^k, \mathbf{l})$, the relation between the probability $p(y_n^k = 1 | \mathbf{m}_n)$ and latent function \mathbf{f}_n can be approximated as:

$$p(y_n^k = 1 | \mathbf{m}_n) = \frac{\exp(m_n^k)}{\sum_{l=1}^{K} \exp(m_n^l)} \approx \frac{\exp(f_n^k)}{\sum_{l=1}^{K} \exp(f_n^l)}$$

$$k = 1, \cdots, K, \ n = 1, \cdots, N.$$
(6)

Furthermore, an hierarchic prior on the covariance function parameters φ^k is employed such that each parameter has an independent exponential distribution $\varphi_d^k \sim Exp(\psi_d^k)$, $d = 1, \dots, D$ and a gamma distribution is placed on the mean value of the exponential $\psi_d^k \sim \Gamma(\sigma^k, \tau^k), d = 1, \dots, D$. Thus, they are forming a conjugate pair. The associated hyper-hyper-parameters $\omega = (\sigma^{k=1,\dots,K}, \tau^{k=1,\dots,K})$ can be set to reflect some prior knowledge of the data. Alternatively, vague priors can be employed such that, for example, each $\sigma^{k=1,\dots,K} = \tau^{k=1,\dots,K} = 1$.

Finally, if we define a set of the hidden variables as $\Theta = \{\pi, \mathbf{m}, \mathbf{f}\}$ and the parameters as $\Phi = \{\varphi^{k=1,\cdots,K}, \psi^{k=1,\cdots,K}\}$, the joint likelihood function for all hidden variables and parameters takes the following form:

$$p(\mathbf{Y}, \Theta, \Phi \mid \mathbf{X}, \omega) = \prod_{n=1}^{N} p(\mathbf{y}_n \mid \boldsymbol{\pi}_n) p(\boldsymbol{\pi}_n \mid \boldsymbol{\alpha}_n = \exp(\mathbf{m}_n))$$
$$\times \prod_{k=1}^{K} p(\mathbf{m}^k \mid \mathbf{f}^k) p(\mathbf{f}^k \mid \mathbf{K}_{(N \times N)}^k(\mathbf{X}, \boldsymbol{\varphi}^k)) \qquad (7)$$
$$\times p(\boldsymbol{\varphi}^k \mid \boldsymbol{\psi}^k) p(\boldsymbol{\psi}^k \mid \boldsymbol{\omega}^k)$$

where the individual factors are respectively define as: $p(\mathbf{Y} | \boldsymbol{\pi})$ is one observation multinomial distribution with probability vector $\boldsymbol{\pi}_n$, $p(\boldsymbol{\pi} | \boldsymbol{\alpha})$ is the Dirichlet distribution with parameter vector $\boldsymbol{\alpha}_n = \exp(\mathbf{m}_n)$, $p(\mathbf{m} | \mathbf{f})$ is the Gaussian distribution with mean \mathbf{f}_n^k and variance 1, $p(\mathbf{f} | \mathbf{X}, \boldsymbol{\varphi})$ is a multivariate Gaussian distribution with zero mean vector and covariance matrix $\mathbf{K}_{(n\times n)}^k(\mathbf{X}, \boldsymbol{\varphi}^k)$, $p(\boldsymbol{\varphi} | \boldsymbol{\psi})$ is an exponential distribution with parameter $\boldsymbol{\psi}_d^k$, and $p(\boldsymbol{\psi} | \boldsymbol{\omega})$ is the gamma distribution with parameter $(\boldsymbol{\sigma}^k, \boldsymbol{\tau}^k)$.

B. The Variational Bayesian Approximation

We now consider an approximate variational Bayesian inference to drive the posterior distributions for latent function

and parameters in our model. This method is a particular variational method which aims to find some approximate joint distribution $q(\Theta) = q(\pi, \mathbf{m}, \mathbf{f})$ over latent variables $\Theta = \{\pi, \mathbf{m}, \mathbf{f}\}$ to approximate the true joint distribution $p(\Theta) = p(\pi, \mathbf{m}, \mathbf{f})$ by minimizing the KL-divergence $KL(q(\Theta) | p(\Theta)$ defined as:

$$KL(q(\Theta) \mid p(\Theta \mid \mathbf{Y}, \mathbf{X}, \Theta)) = -\int q(\Theta) \ln\left(\frac{p(\Theta \mid \mathbf{Y}, \mathbf{X}, \Phi, \omega)}{q(\Theta)}\right) d\Theta$$
(8)

Then, vatiational distribution $q(\Theta)$ is taken by minimizing the KL divergence, and it is equivalent to maximizing the free energy. By the way, the free energy is divided into two components as:

$$F(q,\Theta) = \int q(\Theta) \ln p(\mathbf{Y} | \mathbf{X}, \Phi, \omega) d\Theta - \int q(\Theta) \ln q(\Theta) d\Theta$$
(9)

Here, the variational distribution $q(\Theta)$ is usually assumed to factorize over some partition of latent variables $\Theta = \{\pi, \mathbf{m}, \mathbf{f}\}$ like as $q(\Theta) = q(\pi)q(\mathbf{m})q(\mathbf{f})$. It can be shown using the calculus of variations that the best distribution $q^*(\Theta_i)$ for each of the factors $q(\Theta_i), i = 1, 2, 3$ (in term of the distribution maximizing the free energy) can be expressed as:

$$q^*(\Theta_i) \approx \exp(E_{q(\Theta \setminus \Theta_i)}(\ln p(\mathbf{Y}, \Theta \mid \mathbf{X}, \Phi, \omega))), i = 1, 2, 3 \quad (10)$$

where $E_{q(\Theta \setminus \Theta_i)}(\ln p(\mathbf{Y}, \Theta \mid \mathbf{X}, \Phi, \boldsymbol{a}))$ is the expectation of the logarithm of the joint probability of the data \mathbf{Y} and latent variables Θ , taken over all variables not in the partition. In practice, we usually work in terms of logarithm of the variational distribution $q^*(\Theta_i)$, i.e.:

$$\ln q^{*}(\Theta_{i}) = E_{q(\Theta \setminus \Theta_{i})}(\ln p(\mathbf{Y}, \Theta \mid \mathbf{X}, \Phi, \boldsymbol{\omega})) + \text{constant}, i = 1, 2, 3$$
(11)

From a variational principle that we have considered so far, the variational Bayesian posterior distributions for latent function and all parameters can be summarized by the following iterations which, for all k and d, will optimize the bound on the marginal likelihood.

First, we consider the variational approximate posterior $q^*(\boldsymbol{\pi})$ for the parameter vector of classification probabilities $\boldsymbol{\pi}$.

$$q^{*}(\boldsymbol{\pi}) = \prod_{n=1}^{N} q^{*}(\boldsymbol{\pi}_{n}), \ q^{*}(\boldsymbol{\pi}_{n}) \sim \operatorname{Dir}(\boldsymbol{\pi}_{n} \mid \boldsymbol{\beta}_{n}^{*})$$
(12)
$$\boldsymbol{\beta}_{n}^{*} = (\boldsymbol{\beta}_{n}^{1*}, \cdots, \boldsymbol{\beta}_{n}^{K*}), \text{ where } \boldsymbol{\beta}_{n}^{k*} = y_{n}^{k} + E_{q(m_{n}^{k})}(\exp(m_{n}^{k})).$$

Second, we consider the approximate posterior $q^*(\mathbf{m})$ of normal auxiliary variables \mathbf{m} over Dirichlet auxiliary vector $\boldsymbol{\pi}$ and the latent variables \mathbf{f} . This is given as

$$q^{*}(\mathbf{m}) = \prod_{n=1}^{N} \prod_{k=1}^{K} q^{*}(m_{n}^{k}), \ q^{*}(m_{n}^{k}) = q_{1}^{*}(m_{n}^{k}) \times q_{2}^{*}(m_{n}^{k}) \quad (13)$$

where the distribution $q_1^*(m_n^k)$ is given by

$$q_{1}^{*}(m_{n}^{k}) = \frac{r(m_{n}^{k})}{\sum_{l=1}^{K} r(m_{n}^{l})} r(m_{n}^{k}) = \frac{\exp(\mathrm{E}_{q(\pi_{n})}[\ln \pi_{n}^{k}]\exp(m_{n}^{k}))}{\Gamma(\exp(m_{n}^{k}))}$$
(14)

and $\mathbb{E}_{q^*(\pi_n)}[\ln \pi_n^k] = \upsilon(\beta_n^{k^*}) - \upsilon(\sum_{l=1}^{K} \beta_n^{l^*})$, $\upsilon(x)$ is the digamma function, and the distribution $q_2^*(m_n^k)$ is given by the normal distribution $N(m_n^k | \mathbb{E}_{q(t_n^k)}(\mathbf{f}_n^k), \mathbf{1})$. Therefore, the required posterior expectation $E(m_n^k)$ can be computed as the following manner by using importance sampling method:

$$E(m_n^k) = \int m_n^k q^*(m_n^k) dm_n^k \approx \sum_{l=1}^L m_n^k w(m_n^{k(l)})$$
(15)

where each $m_n^{k(1)}, \dots, m_n^{k(L)}$ are random samples drawn from;

$$N(m_n^k \mid \mathsf{E}_{q(\mathbf{f}_n^k)}(\mathbf{f}_n^k), \mathbf{l}), \text{ and } w(m_n^{k(l)}) = \frac{q_1^*(m_n^{k(l)})}{\sum_{s=1}^L q_1^*(m_n^{k(s)})}$$
(16)

Third, we consider the approximate posterior $q^*(\mathbf{f})$ of latent variables \mathbf{f} over normal auxiliary variables \mathbf{m} and the parameters $\varphi^{k=1,\dots,K}$. We have obtained the approximate posterior $q^*(\mathbf{f})$ for latent function \mathbf{f} :

$$q^*(\mathbf{f}) = \prod_{k=1}^{K} N(\mathbf{f}^k \mid \mathrm{E}(\mathbf{f}^k), \Sigma^k), \qquad (17)$$

where $E(\mathbf{f}^k) = \Sigma^k E(\mathbf{m}^k)$, and

$$\Sigma^{k} = \mathbf{K}_{(N \times N)}^{k} (\mathbf{E}_{q(\boldsymbol{\varphi}^{k})}(\boldsymbol{\varphi}^{k})) (\mathbf{I} + \mathbf{K}_{(N \times N)}^{k} (\mathbf{E}_{q(\boldsymbol{\varphi}^{k})}(\boldsymbol{\varphi}^{k})))^{-1}$$

Fourth, if we also consider the set of hyper-parameters $\Phi = \{ \varphi^{k=1, \dots, K}, \psi^{k=1, \dots, K} \}$, in this variational treatment, then the expectation of the covariance kernel hyper-parameters $\varphi^{k=1, \dots, K}$ under the variational posterior distribution $q(\varphi^k)$ can be approximated by drawing *S* samples such that each $\varphi_s^{kd} \sim Exp(E_{q(w^{kd})}(\psi^{kd}))$ and so;

(18)

where,

$$w(\mathbf{\phi}_{s}^{k}) = \frac{N(E_{q^{*}(\mathbf{f}^{k})}(\mathbf{f}^{k}) | \mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{\phi}_{s}^{k}))}{\sum_{v=1}^{S} N(E_{q^{*}(\mathbf{f}^{k})}(\mathbf{f}^{k}) | \mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{\phi}_{v}^{k}))}$$

 $\mathbf{E}_{q(\boldsymbol{\varphi}^k)}[\boldsymbol{\varphi}^k] \approx \sum_{s=1}^{S} \boldsymbol{\varphi}_s^k w(\boldsymbol{\varphi}_s^k),$

Finally, the approximate posterior for hyper-parameters ψ_d^k is given as the gamma distribution with parameters $\sigma^k + 1$ and $\tau^k + E_{q(\varphi_d^k)}(\varphi_d^k)$. And the required posterior mean values are given as;

$$E_{q(\psi_d^k)}(\psi_d^k) = \frac{\sigma^k + 1}{\tau^k + \mathcal{E}_{q(\varphi_d^k)}(\varphi_d^k)}$$
(19)

C. Predictive Distribution for New Observation

Here, we have first to consider a predictive distribution $q(\mathbf{f}_{new} | \mathbf{X}, \mathbf{x}_*, \Theta, \mathbf{y})$ of latent function $\mathbf{f}_{new} = (\mathbf{f}_{new}^1, \cdots, \mathbf{f}_{new}^K)^T$ corresponding to new observation \mathbf{x}_{new} . Here, we define $\mathbf{K}^k(\mathbf{X} | \mathbf{E}(\varphi^k))$ as the $(N \times N)$ covariance matrix containing the covariance function values given by training data points \mathbf{X} , and we also define $\mathbf{k}_{new}^k(\mathbf{x}_{new}, \mathbf{X} | \mathbf{E}(\varphi^k))$ as the $(N \times 1)$ covariance function values given by training data points \mathbf{X} , and we also define $\mathbf{k}_{new}^k(\mathbf{x}_{new}, \mathbf{X} | \mathbf{E}(\varphi^k))$ as the $(N \times 1)$ covariance vector containing the covariance function values between the new point \mathbf{x}_{new} and those contained in training data \mathbf{X} . Moreover $k_{new}^k(\mathbf{x}_{new}, \mathbf{x}_{new} | \mathbf{E}(\varphi^k))$ denotes the covariance function value for the new point and itself.

Here, the predictive distribution $q(\mathbf{f}_{new} | \mathbf{X}, \mathbf{x}_*, \Theta, \mathbf{y})$ will be Gaussian. Hence, its mean vector and covariance matrix are given as follows. The predictive mean vector for latent function \mathbf{f}_{new} is given by:

$$E_{q}(\mathbf{f}_{new} | \mathbf{X}, \mathbf{x}_{*}, \Theta, \mathbf{y}) = Q(\mathbf{X}, \mathbf{x}_{new} | E(\varphi))^{T} \mathbf{K}(\mathbf{X}, E(\varphi))^{-1} \boldsymbol{\mu}_{\mathbf{f}}$$

$$= Q(\mathbf{X}, \mathbf{x}_{new} | E(\varphi))^{T} (\mathbf{I} + \mathbf{K}(\mathbf{X}, E(\varphi)))^{-1} E(\mathbf{m})$$
(20)

and the covariance matrix for latent function \mathbf{f}_{new} is given by:

$$Cov_{q}(\mathbf{f}_{new} | \mathbf{X}, \mathbf{y}, \Theta, \mathbf{x}_{new}) = \operatorname{diag}(k_{new}^{1}, \cdots, k_{new}^{K}) -Q(\mathbf{X}, \mathbf{x}_{new} | \mathbf{E}(\varphi))^{T} (\mathbf{I} + \mathbf{K}(\mathbf{X}, \mathbf{E}(\varphi)))^{-1} Q(\mathbf{X}, \mathbf{x}_{new} | \mathbf{E}(\varphi))$$
(21)

When the mean $E(f_{new}^k)$ and the variance $Var(f_{new}^k)$ of latent function f_{new}^k for each class $k = 1, \dots, K$ are given, we extract the *S* samples $f_{new}^{k(1)}, \dots, f_{new}^{k(S)}$ of latent variable f_{new}^k from a normal distribution $N(f_{new}^k | E(f_{new}^k), Var(f_{new}^k))$ having this mean and variance. And using a similar method, we also extract the *S* samples $m_{new}^{k(1)}, \dots, m_{new}^{k(S)}$ of the auxiliary variables m_{new}^k for each class $k = 1, \dots, K$ from the normal distribution $N(m_{new}^k | f_{new}^k, 1)$ having a mean f_{new}^k and variance 1. Next, we define the parameters $\boldsymbol{\alpha}_{new}^{s} = (\exp(m_{new}^{1(s)}), \cdots, \exp(m_{new}^{K(s)}))$ of the Dirichlet probability distribution using the extracted samples $(m_{new}^{1(s)}, \cdots, m_{new}^{K(s)})$ and we extract again the classification probabilities vector $\boldsymbol{\pi}_{new}^{(s)} = (\boldsymbol{\pi}_{new}^{1(s)}, \cdots, \boldsymbol{\pi}_{new}^{K(s)})$ from the Dirichlet distribution.

$$\boldsymbol{\pi}_{new}^{(s)} \sim Dir((\pi_{new}^{1(s)}, \cdots, \pi_{new}^{K(s)}) | (\exp(m_{new}^{1(s)}), \cdots, \exp(m_{new}^{K(s)})) .$$
(22)

Repeating this procedure S times, we have generated a total S number of classification probabilities vector $\boldsymbol{\pi}_{new}^{(1)}, \dots, \boldsymbol{\pi}_{new}^{(S)}$. And using them, we calculate the mean of the classification probabilities vectors as:

$$\overline{\boldsymbol{\pi}}_{new} = \frac{1}{S} (\boldsymbol{\pi}_{new}^{(1)} + \dots + \boldsymbol{\pi}_{new}^{(S)})$$
(23)

III. EXPERIMENTAL RESULTS

To estimate the performance of the proposed algorithm, we will consider four partially overlapping Gaussian sources of data in two dimensions. First, in order to train a model, we generate four classes' bivariate Gaussian random samples. One hundred twenty data points were generated by the four bivariate normal distributions with the mean vectors and covariance matrices described in Table I.

TABLE I MEAN VECTOR AND COVARIANCE MATRIX FOR EACH CLASS S

class	Mean vector	Covariance matrix
Class1	(1.75,-1.0)	$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
Class2	(-1.75,1.0)	$\begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$
Class3	(2,2)	$\begin{pmatrix}1 & -0.5\\-0.5 & 1\end{pmatrix}$
Class4	(-2,-2)	$\begin{pmatrix}1 & -0.5\\-0.5 & 1\end{pmatrix}$

One hundred and twenty draws were made from Table I and the test samples are used in the proposed algorithm with a further 4,000 points. Each of the sample values were sampled uniformly thus creating a balance of samples drawn from the four target classes. Fig. 1 (a) shows a plotting of training data points on a two-dimensional space. Second, in order to verify the performance of the model, we generate different four classes of bivariate Gaussian random sample. Four thousands data points were generated by the bivariate normal distribution.

Fig. 2 (a) shows the results that were monitored for covariance parameters during each step and as would be expected, it shows us a steady convergence in the improvements of the parameters. From Fig. 2 (b), we can also see that the development of the predictive performance follows that achieving a predictive performance of 92.2% at convergence. Therefore, it knows that the proposed method can entirely classify well data points.







Fig. 1 (a) Training data, (b) Testing data



Fig. 2 Experiments Results (a) Covariance Parameters ϕ_k , (b) Correct Prediction rate

IV. CONCLUSION

In this paper, we consider probabilistic multinomial Dirichlet classification model for multidimensional data with Gaussian process priors placed over the latent function. Variational Byesian algorithm is used to drive approximate posteriors for latent function as well as parameters needed to define the Multinomial dirichlect Gaussian process classification model. The proposed algorithm was performed in two steps: each of which is the training step and classification step. First, in the training step, using the variational Bayesian formula, we derived approximately the posterior distribution of the latent function and parameters on the basis of the learning data. Second, in the classification step, using a derived posterior distribution of the latent function, we estimate the classification probabilities to assign new sample into proper class. We assign this sample into a class that has the maximum probability.

ACKNOWLEDGMENT

This work was jointly supported by the National Research Foundation of Korea Government (2014R1A1A4A0109398) and the Research Foundation of Chonnam National University (2014-2256).

REFERENCES

- [1] C. E. Rasmussen, and C. K. I. Williams, *Gaussian Processes for Machine Learning*, MIT Press, 2006.
- [2] H. Nicklisch, and C. E. Rasmussen, "Approximation for Binary Gaussian process Classification", *Journal of Machine Learning Research*, vol 9, pp 2035-75, 2008.
- [3] C. K. I. Williams, and D. Barber, "Bayesian Classification with Gaussian Processes", *IEEE Tran. On PAMI*, vol 12, pp 1342-1351, 1998.
- [4] T. P. Minka, "Expectation Propagation for Approximate Bayesian Inference", *In UAI, Morgan Kaufmann*, pp 362-369, 2001.
- [5] M. Opper, and O. Winther, "Gaussian Processes for Classification: Mean Field Algorithms", *Neural Computation*, vol 12 pp 2655-2204, 2000.
 [6] L. Csato, E. Fokoue, M. Opper, and B. Schottky, "Efficient Approaches
- [6] L. Csato, E. Fokoue, M. Opper, and B. Schottky, "Efficient Approaches to Gaussian Process Classification", *In Neural Information Processing Systems*, vol 12, pp 251-257, 2000.
- [7] M. Girolami and S. Rogers, "Variational Bayesian Multinomial Probit Regression with Gaussian Process Priors", *Neural Computation*, vol 18 pp 1790-1817, 2006.

Wan-Hyun Cho received both B.S. degree and M.S. degree from the Department of Mathematics, Chonnam National University, Korea in 1977 and 1981, respectively and Ph.D. degree from the Department of Statistics, Korea University, Korea in 1988. He is now teaching in Chonnam National University. His research interests are statistical modeling, pattern recognition, image processing, and medical image processing.

Soon-Ja Kang received both B.S. degree and M.S. degree from the Department of Mathematics, Chonnam National University, Korea in 1979 and 1981, respectively and Ph.D. degree from the Department of Mathematics, Seogang University, Korea in 1988. She is now teaching in Chonnam National University. Her research fields are mathematical education, advanced calculus, and education for the gifted children.

Sang-Kyoon Kim received the B.S., M.S. and Ph.D. degrees in Electronics Engineering, Mokpo National University, Korea in 1998, 2000 and 2015 respectively. From 2011 to 2015, he was a Visiting Professor in the Department of Information & Electronics Engineering, Mokpo National University, Korea. His research interests include image processing, pattern recognition and computer vision.

Soon-Young Park received B.S. degree in Electronics Engineering from Yonsei University, Korea in 1982 and M.S and Ph.D. degrees in Electrical and Computer Engineering from State University of New York at Buffalo, in 1986 and 1989, respectively. From 1989 to 1990 he was a Postdoctoral Research Fellow in the department of Electrical and Computer Engineering at the State University of New York at Buffalo. Since 1990, he has been a Professor with Department of Electronics Engineering, Mokpo National University, Korea. His research interests include image and video processing, image protection and authentication and image retrieval techniques.