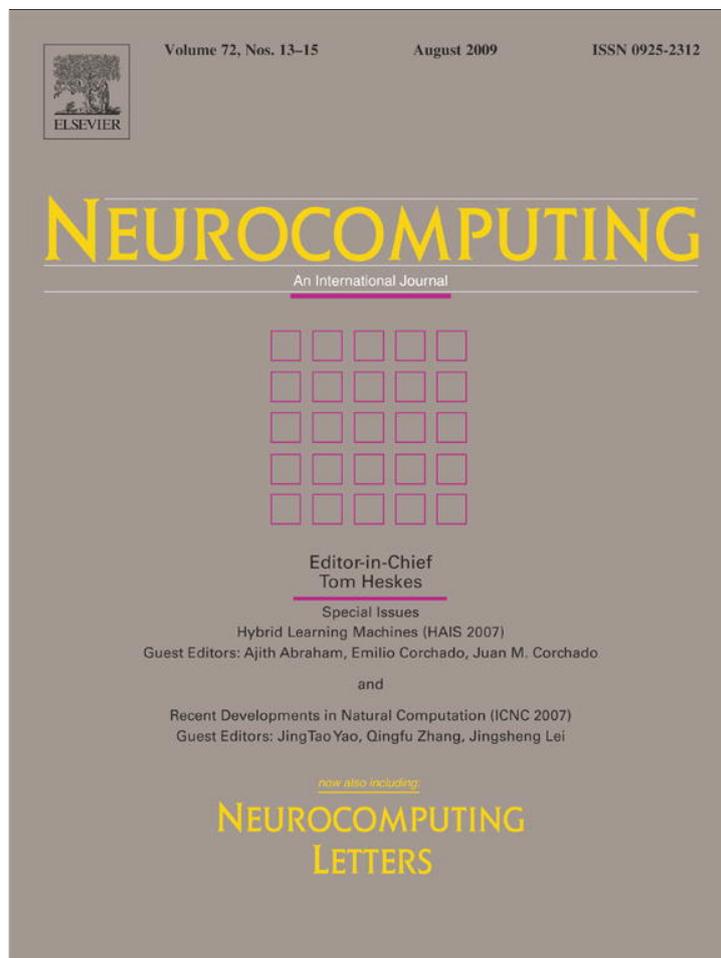


Provided for non-commercial research and education use.
Not for reproduction, distribution or commercial use.

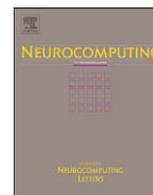


This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Letters

A novel view of the variational Bayesian clustering

Takashi Takekawa^{a,*}, Tomoki Fukai^{a,b}^a Laboratory for Neural Circuit Theory, RIKEN Brain Science Institute (BSI), 2-1 Hirosawa, Wako, Saitama 351-0198, Japan^b Department of Complexity Science and Engineering, University of Tokyo, 5-1-5 Kashiwanoha, Kashiwa, Chiba 277-8561, Japan

ARTICLE INFO

Article history:

Received 30 June 2008

Accepted 16 April 2009

Communicated by L.C. Jain

Available online 10 May 2009

Keywords:

Unsupervised learning

Bayesian estimation

Variational approximation

Model selection

Robust variational Bayes

ABSTRACT

We prove that the evaluation function of variational Bayesian (VB) clustering algorithms can be described as the log likelihood of given data minus the Kullback–Leibler (KL) divergence between the prior and the posterior of model parameters. In this novel formalism of VB, the evaluation functions can be explicitly interpreted as information criteria for model selection and the KL divergence imposes a heavy penalty on the posterior far from the prior. We derive the update process of the variational Bayesian clustering with finite mixture Student's *t*-distribution, taking the penalty term for the degree of freedoms into account.

© 2009 Elsevier B.V. All rights reserved.

1. Introduction

Data observed in experiments or statistical researches can often be categorized into clusters of data points with similar features. Such a clustering, however, is not trivially easy for various reasons: the features used for it may not be optimal, or the data set may include noisy elements. Here, we reformulate the variational Bayesian (VB) method to solve difficult clustering problems.

Let m be the number of clusters and $\theta = \{\alpha_1, \alpha_2, \dots, \alpha_m, \beta_1, \beta_2, \dots, \beta_m\}$ be the set of parameters, where α and β represent the size and other parameters of cluster k , respectively. When m and θ are given, we can define a finite mixture model by the following conditional probability $p(x_n, z_n = k | \theta, m)$ that the n -th data takes a value of $x_n \in \mathcal{R}^D$ and belongs to the k -th cluster with probability α_k :

$$p(x_n, z_n = k | \theta, m) = \alpha_k \mathcal{A}(x_n | \beta_k). \quad (1)$$

Here, $\mathcal{A}(x_n | \beta_k)$ is some probability distribution to generate x_n , and variable $z_n \in \{1, 2, \dots, m\}$ specifies the cluster from which x_n is taken.

Conversely, we want to estimate the optimal values of m , θ and $\mathbf{z} = \{z_1, z_2, \dots, z_N\}$ when data points are given as $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$. If the prior distribution is set as $p(\theta, m) = p(m)p(\theta|m)$, we can in principle estimate those quantities by calculating the

marginal likelihood

$$p(\mathbf{x}) = \sum_m p(m) \int p(\theta|m) \sum_n \sum_k p(x_n, z_n = k | \theta, m) d\theta \quad (2)$$

and the posterior distributions:

$$p(m|\mathbf{x}) = \frac{p(m) \int p(\theta|m) \sum_n \sum_k p(x_n, z_n = k | \theta, m) d\theta}{p(\mathbf{x})}, \quad (3)$$

$$p(z_n = k | \mathbf{x}, m) = \frac{\int p(\theta|m) p(x_n, z_n = k | \theta, m) d\theta}{p(\mathbf{x})}. \quad (4)$$

In practice, however, the above high-dimensional integrals are difficult to calculate. To overcome this difficulty, we may employ Markov Chain Monte Carlo (MCMC) methods for performing the above integrations (see [6]). In MCMC, however, the convergence to a correct solution is generally slow and is ensured only on a strict condition.

VB gives a deterministic algorithm to find approximate solutions to the posterior distributions [3]. The algorithm of VB is fast and resembles expectation–maximization (EM) algorithm for the maximum likelihood estimation [4]. In particular, optimizing the resultant clusters does not use information criteria, such as AIC [1] or BIC [7]. Rather, the model selection in VB is automatically accomplished by maximizing an estimation function. VB often shows a better generalization ability than EM when the number of data points is small. Clustering by VB can also be improved by using mixtures of principal component analyzers [5]. However, the formalism of VB is not transparent since the estimation function does not provide a simple information-theoretic interpretation.

* Corresponding author.

E-mail address: takekawa@riken.jp (T. Takekawa).

In this report, we show that the estimation function can be expressed as a sum of the log likelihood and a penalty term that has a simple geometrical meaning. Namely, the penalty term can be expressed as the Kullback–Leibler divergence of the prior and the posterior of model parameters. This result achieves a deeper insight into the theoretical framework of VB. Moreover, the result provides a more sophisticated software implementation than the previous ones. To show this explicitly, we derive a VB algorithm by using a mixture of Student's t -distributions. Such an algorithm is known to be more robust against noisy data points (i.e., outliers) than the algorithm with Gaussian mixtures [2,8].

The previous algorithm of VB with Student's t -mixtures assumed a uniform distribution for the priori of the degree of freedom ν [2]. This prior, however, makes the calculations of the penalty term of ν difficult, so the term was simply neglected in the formula. Here, we derive a natural expression of the penalty term by employing an exponential distribution for the prior of ν . Preliminary results obtained in the analysis of multi-unit recording data [9] revealed that our algorithm significantly improves the speed and precision of spike sorting. The results of spike sorting will be reported elsewhere.

2. Preparations

Below we explain our framework of variational Bayesian formulated with a mixture of Student's t -distributions.

2.1. Approximation

The Student's t -distribution can be described as follows:

$$\mathcal{F}(x|\nu, \mu, \Sigma) = \int_0^\infty \mathcal{N}(x|\mu, u^{-1}\Sigma) \mathcal{G}\left(u\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) du, \quad (5)$$

where the normal distribution \mathcal{N} and the gamma distribution \mathcal{G} are given as

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{|2\pi\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} \text{tr} \Sigma^{-1}(x - \mu)(x - \mu)^\top\right\}, \quad (6)$$

$$\mathcal{G}(u|a, b) = \frac{b^a}{\Gamma(a)} u^{a-1} \exp(-bu). \quad (7)$$

We define a probability distribution involving unobserved latent variable u_{nk} as

$$p(x_n, z_n = k, u_{nk}|\theta, m) = \alpha_k \mathcal{N}(x_n|\mu_k, u_{nk}^{-1}\Sigma_k) \mathcal{G}\left(u_{nk}\left|\frac{\nu_k}{2}, \frac{\nu_k}{2}\right.\right), \quad (8)$$

where ν_k , μ_k and Σ_k represent the degree of freedom, mean and variance matrix of cluster k , respectively. Then, we introduce a functional of test function $q(\mathbf{z}, \mathbf{u}, \theta, m)$, which should approximate the posterior distribution $p(\mathbf{z}, \mathbf{u}, \theta, m|\mathbf{x})$, as

$$F[q(\mathbf{z}, \mathbf{u}, \theta, m)] = \sum_m \sum_z \int \int q(\mathbf{z}, \mathbf{u}, \theta, m) \log \frac{p(\mathbf{x}, \mathbf{z}, \mathbf{u}, \theta, m)}{q(\mathbf{z}, \mathbf{u}, \theta, m)} d\theta d\mathbf{u} = \log p(\mathbf{x}) - \text{KL}[q(\mathbf{z}, \mathbf{u}, \theta, m), p(\mathbf{z}, \mathbf{u}, \theta, m|\mathbf{x})]. \quad (9)$$

Here, Kullback–Leibler divergence

$$\text{KL}[q(\mathbf{z}, \mathbf{u}, \theta, m), p(\mathbf{z}, \mathbf{u}, \theta, m|\mathbf{x})] = \sum_m \sum_z \int \int q(\mathbf{z}, \mathbf{u}, \theta, m) \log \frac{q(\mathbf{z}, \mathbf{u}, \theta, m)}{p(\mathbf{z}, \mathbf{u}, \theta, m|\mathbf{x})} d\theta d\mathbf{u} \quad (10)$$

vanishes when the test distribution coincides with the posterior. Since $p(\mathbf{x})$ takes a constant for given \mathbf{x} , minimizing $\text{KL}[q(\mathbf{z}, \mathbf{u}, \theta, m), p(\mathbf{z}, \mathbf{u}, \theta, m|\mathbf{x})]$ is equivalent to maximizing F and hence allows $q(\mathbf{z}, \mathbf{u}, \theta, m)$ to approach $p(\mathbf{z}, \mathbf{u}, \theta, m|\mathbf{x})$.

2.2. Model selection

To maximize F , we introduce the following function:

$$F_m[q(\mathbf{z}, \mathbf{u}, \theta|m)] = \sum_z \int \int q(\mathbf{z}, \mathbf{u}, \theta|m) \log \frac{p(\mathbf{x}, \mathbf{u}, \mathbf{z}, \theta|m)}{q(\mathbf{z}, \mathbf{u}, \theta|m)} d\theta d\mathbf{u}. \quad (11)$$

By using the above function and $q(\mathbf{z}, \mathbf{u}, \theta, m) = q(m)q(\mathbf{z}, \mathbf{u}, \theta|m)$, we can rewrite F as

$$F = \sum_m \left\{ q(m)F_m - q(m) \log \frac{q(m)}{p(m)} \right\}. \quad (12)$$

By maximizing F with respect to $q(m)$, we obtain

$$q(m) \propto p(m) \exp F_m. \quad (13)$$

Since no prior information is available for m , we may regard $p(m)$ as constant. Therefore, we should select the value of m that maximizes F_m as the optimal number of clusters. Thus, we have to solve the maximization problem of F_m .

3. Implementations

Hereafter, we explain our algorithm to maximize F_m . For brevity, we do not explicitly show the dependence of m in the argument and write $q(a|m)$, $p(a|m)$ as $q(a)$, $p(a)$ in the abbreviated form. In addition, we assume that the following factorization holds: $q(\mathbf{z}, \mathbf{u}, \theta) = q(\mathbf{z}, \mathbf{u})q(\theta)$. With this assumption, Eq. (11) can be transformed into

$$F_m[q(\mathbf{z}, \mathbf{u}), q(\theta)] = \sum_z \int \int q(\mathbf{z}, \mathbf{u})q(\theta) \log \frac{p(\mathbf{x}, \mathbf{z}, \mathbf{u}|\theta)p(\theta)}{q(\mathbf{z}, \mathbf{u})q(\theta)} d\theta d\mathbf{u}. \quad (14)$$

3.1. Prior

The prior of the model parameters can be decomposed as $p(\theta) = p(\boldsymbol{\alpha})\prod_k\{p(\nu_k)p(\Sigma_k)p(\mu_k|\Sigma_k)\}$. We use the following distributions in this study:

$$p(\boldsymbol{\alpha}) = \mathcal{D}(\{\alpha_1, \alpha_2, \dots, \alpha_m\}|\{\kappa_0, \kappa_0, \dots, \kappa_0\}), \quad (15)$$

$$p(\nu_k) = \xi_0 \exp(-\xi_0 \nu_k), \quad (16)$$

$$p(\Sigma_k) = \mathcal{W}^{-1}(\Sigma_k|\gamma_0, \gamma_0 \Sigma_0), \quad (17)$$

$$p(\mu_k|\Sigma_k) = \mathcal{N}(\mu_k|\mu_0, \eta_0^{-1}\Sigma_k), \quad (18)$$

where the Dirichlet distribution \mathcal{D} and the inverse Wishart distribution \mathcal{W}^{-1} are given as

$$\mathcal{D}(\{\alpha_1, \dots, \alpha_m\}|\{\kappa_1, \dots, \kappa_m\}) = \frac{\Gamma(\sum_k \kappa_k)}{\prod_k \Gamma(\kappa_k)} \prod_k \alpha_k^{\kappa_k - 1}, \quad (19)$$

$$\mathcal{W}^{-1}(\Sigma|\gamma, \Delta) = \frac{|\frac{1}{2}\Delta|^{\gamma/2} \exp(-\frac{1}{2}\text{tr} \Delta \Sigma^{-1})}{\Gamma_D(\gamma/2)|\Sigma|^{1/2(\gamma+D+1)}} \quad (20)$$

and Γ_D is the multivariate gamma function with dimension D . The conventional distributions are used for $\boldsymbol{\alpha}$, μ_k and Σ_k . As mentioned before, we adopted an exponential distribution for ν_k rather than the uniform distribution chosen previously [2].

3.2. Initial conditions

VB algorithm alternately performs an E step and an M step. An E step computes the expectation of the likelihood as if the latent variables were observable, and an M step computes the maximum likelihood estimates of the parameters by maximizing the expected likelihood found on the preceding E step. The parameters found on the M step are then used to begin another E step.

The value of optimization function F_m in Eq. (14) is evaluated after every E step, and the process is repeated until the value of F_m converges to a steady value.

The test distributions can be obtained in the derivation of the M step. Below, hyper-parameters of the test distribution are shown with tilde. As in the case of the prior of the model parameters $p(\theta)$, the test distribution of the model parameters can also be decomposed as $q(\theta) = q(\boldsymbol{\alpha}) \prod_k \{q(v_k)q(\Sigma_k)q(\mu_k|\Sigma_k)\}$. We use the following distributions in this study:

$$q(\boldsymbol{\alpha}) = \mathcal{D}(\{\alpha_1, \alpha_2, \dots, \alpha_m\} | \{\tilde{\kappa}_1, \tilde{\kappa}_2, \dots, \tilde{\kappa}_m\}), \quad (21)$$

$$q(v_k) = \frac{(v_k/2)^{v_k/2} \exp(-\tilde{\zeta}_k v_k)}{C_v(\tilde{\zeta}_k) \Gamma(v_k/2)}, \quad (22)$$

$$q(\Sigma_k) = \mathcal{W}^{-1}(\Sigma_k | \tilde{\gamma}_k, \tilde{\gamma}_k \tilde{\Sigma}_k), \quad (23)$$

$$q(\mu_k|\Sigma_k) = \mathcal{N}(\mu_k | \tilde{\mu}_k, \tilde{\eta}_k^{-1} \Sigma_k), \quad (24)$$

where Γ is the gamma function. For the convenience of practical computations, we may replace $q(v_k)$ with $\delta(v_k - \tilde{v}_k)$ by MAP approximation, where \tilde{v}_k is the value of v that maximizes $q(v_k)$. Thus, \tilde{v}_k is a solution to $dq(v_k)/dv_k = 0$. Normalization constant $C_\mu(\tilde{\zeta})$ is determined as

$$C_\mu(\tilde{\zeta}) = \int \frac{(v/2)^{v/2} \exp(-\tilde{\zeta} v)}{\Gamma(v/2)} dv. \quad (25)$$

3.3. E step

On an E step, $F_m[q(\mathbf{z}, \mathbf{u}), q(\theta)]$ is maximized with respect to $q(\mathbf{z}, \mathbf{u})$ while $q(\theta)$ is fixed. We can find such a $q(\mathbf{z}, \mathbf{u})$ as

$$q(z_n = k, u_{nk}) = \frac{\zeta_{nk}}{\sum_{k'} \rho_{nk'}} \quad (26)$$

by using the Lagrange multiplier method, where

$$\begin{aligned} \zeta_{nk} &= \exp \int q(\theta) \log p(x_n, z_n = k, u_{nk} | \theta) d\theta \\ &= \frac{\hat{\alpha}_k (\tilde{v}_k/2)^{\tilde{v}_k/2}}{(2\pi)^{D/2} \hat{\Sigma}^{1/2} \Gamma(\tilde{v}_k/2)} u_{nk}^{(\hat{\alpha}_k - 1)} \exp(-b_{nk} u_{nk}), \end{aligned} \quad (27)$$

$$\rho_{nk} = \int \zeta_{nk} du_{nk} = \frac{\hat{\alpha}_k (\tilde{v}_k/2)^{\tilde{v}_k/2}}{(2\pi)^{D/2} \hat{\Sigma}^{1/2} \Gamma(\tilde{v}_k/2)} \Gamma(\hat{\alpha}_k) b_{nk}^{-\hat{\alpha}_k}. \quad (28)$$

The parameters appearing in the above equations are defined in Appendix A. Then, using

$$q(u_{nk} | z_n = k) = \frac{\zeta_{nk}}{\rho_{nk}} = \mathcal{G}(u_{nk} | a_{nk}, b_{nk}) \quad (29)$$

we can derive the following quantities:

$$\bar{z}_{nk} = q(z_n = k) = \frac{\rho_{nk}}{\sum_{k'} \rho_{nk'}}, \quad (30)$$

$$\bar{u}_{nk} = \int q(u_{nk} | z_n = k) u_{nk} du_{nk} = \frac{a_{nk}}{b_{nk}}, \quad (31)$$

$$\begin{aligned} \log \hat{u}_{nk} &= \int q(u_{nk} | z_n = k) \log u_{nk} du_{nk} \\ &= \Psi(a_{nk}) - \log b_{nk}, \end{aligned} \quad (32)$$

which are used on the succeeding M step, where Ψ is the digamma function.

3.4. Evaluation

The output of each E step is evaluated by the value of $F_m[q(\mathbf{z}, \mathbf{u}), q(\theta)]$ before a succeeding M step is processed. We first

rewrite Eq. (14) as

$$\begin{aligned} F_m &= \sum_{\mathbf{z}} \int q(\mathbf{z}, \mathbf{u}) q(\theta) \log \frac{p(\mathbf{x}, \mathbf{z}, \mathbf{u} | \theta)}{q(\mathbf{z}, \mathbf{u})} d\theta d\mathbf{u} \\ &\quad - \sum_{\mathbf{z}} \int q(\mathbf{z}, \mathbf{u}) d\mathbf{u} \int q(\theta) \log \frac{q(\theta)}{p(\theta)} d\theta. \end{aligned} \quad (33)$$

Since $\sum_{\mathbf{z}} \int q(\mathbf{z}, \mathbf{u}) d\mathbf{u} = 1$ holds, the second term coincides with $\text{KL}[q(\theta), p(\theta)]$.

Then, using (27), we can transform the first term into

$$\sum_n \sum_k \int q(z_n = k, u_{nk}) \{ \log \zeta_{nk} - \log q(z_n = k, u_{nk}) \} du_{nk}, \quad (34)$$

which we can further rewrite as

$$\begin{aligned} &\sum_n \sum_k \int q(z_n = k, u_{nk}) \left\{ \log \zeta_{nk} - \log \frac{\zeta_{nk}}{\sum_{k'} \rho_{nk'}} \right\} du_{nk} \\ &= \sum_n \sum_k \int q(z_n = k, u_{nk}) du_{nk} \log \sum_{k'} \rho_{nk'} \\ &= \sum_n \log \sum_{k'} \rho_{nk'} \end{aligned} \quad (35)$$

by using (26) and $\sum_k \int q(z_n = k, u_{nk}) du_{nk} = 1$. Thus, we finally obtain the following expression of the estimation function:

$$F_m = \sum_n \log \sum_{k'} \rho_{nk'} - \text{KL}[q(\theta), p(\theta)]. \quad (36)$$

The above results reveal that the estimation function consists of the log likelihood given by $\sum_k \rho_{nk}$, which has already been calculated at the preceding E step, and the penalty term calculated from the test function $q(\theta)$. Thus, VB method maximizes the log likelihood calculated from the test function with a penalty given as the deviation of the test function from the prior. Note that our evaluation of F_m no longer requires heavy calculations of the multi-dimensional integrals in Eq. (33), whereas the conventional evaluation step calculates the integrals directly.

3.5. M step

An M step completes one cycle of EM algorithm. M step is complementary to E step, that is, $F_m[q(\mathbf{z}, \mathbf{u}), q(\theta)]$ is maximized with respect to $q(\theta)$ while $q(\mathbf{z}, \mathbf{u})$ is fixed. Such a $q(\theta)$ can be expressed in a fashion similar to Eqs. (21)–(24), and the hyper-parameters are derived as

$$\tilde{\kappa}_k = \kappa_0 + \bar{N}_k, \quad (37)$$

$$\tilde{\gamma}_k = \gamma_0 + \bar{N}_k, \quad (38)$$

$$\tilde{\eta}_k = \eta_0 + \bar{M}_k, \quad (39)$$

$$\tilde{\zeta}_k = \zeta_0 + \frac{\bar{M}_k - \hat{M}_k}{2\bar{N}}, \quad (40)$$

$$\tilde{\mu}_k = \frac{\eta_0 \mu_0 + \bar{M}_k \bar{\mu}_k}{\eta_0 + \bar{M}_k}, \quad (41)$$

$$\tilde{\Sigma}_k = \frac{1}{\gamma_0 + \bar{N}_k} \left\{ \gamma_0 \Gamma_0 + \bar{M}_k \tilde{\Sigma}_k + \frac{\eta_0 \bar{M}_k}{\eta_0 + \bar{M}_k} (\bar{\mu}_k - \mu_0)(\bar{\mu}_k - \mu_0)^T \right\}. \quad (42)$$

The parameters appearing in the above equations are defined in Appendix A.

4. Discussion

If we take the limit of vanishing ζ_0 in Eq. (40), the prior approaches a uniform distribution and the renewal rule coincides

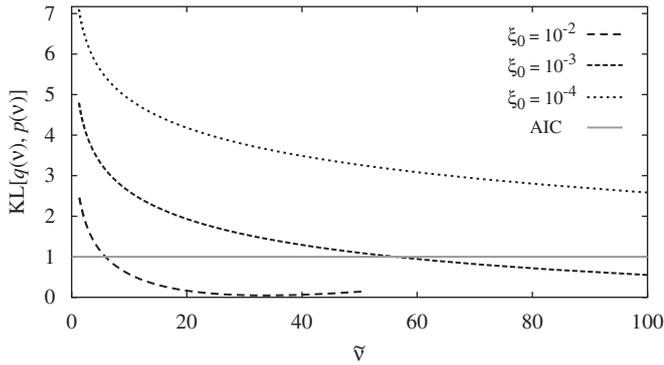


Fig. 1. Penalty terms of ν for various values of ξ_0 .

with the previous one [2], in which the penalty term diverges, as shown below. By contrast, too large ξ_0 implies too strong influences of the prior. Thus, an appropriate value range of ξ_0 should exist. This point was studied in an example.

In Fig. 1, we plot the values of the penalty term of ν (i.e., $\text{KL}[q(\nu), p(\nu)]$) for the estimated values of $\tilde{\nu}$. Three values of ξ_0 were tested. As a reference, we calculated the penalty term that should be added to AIC, which is 1, by the inclusion of ν . Note that at $\xi_0 = 10^{-2}$ the estimated value of $\tilde{\nu}$ cannot be larger than 50. When ξ_0 is smaller than 10^{-4} , the prior distribution is very close to a uniform distribution. However, the penalty term is rapidly increased for such a small value of ξ_0 . We suggest that an appropriate value of the hyper-parameter of the prior distribution is $\xi_0 \approx 10^{-3}$.

5. Conclusions

In this letter, we have derived a simple expression for the estimation function of VB method. This expression is mathematically transparent and is easy to calculate numerically. We applied the obtained framework to derive the penalty term and the renewal rule for the degrees of freedom ν of the robust VB method.

Acknowledgments

This research was partially supported by RIKEN Special Postdoctoral Researchers Program (to T. Takekawa) and MEXT Grant-in-Aid for Scientific Research on Priority Areas 17022036 (to T. Fukai).

Appendix A

The parameters used in Eqs. (27)–(28) are defined as follows:

$$\log \hat{\alpha}_k = \int q(\boldsymbol{\alpha}) \log \alpha_k d\boldsymbol{\alpha} = \Psi(\tilde{\kappa}_k) - \Psi\left(\sum_k \tilde{\kappa}_k\right), \quad (43)$$

$$\begin{aligned} \log \hat{\Sigma}_k &= \int q(\Sigma_k) \log |\Sigma_k| d\Sigma_k \\ &= \log \left| \frac{\tilde{\gamma}_k}{2} \tilde{\Sigma}_k \right| - \sum_{i=0}^{D-1} \Psi\left(\frac{\tilde{\gamma}_k - i}{2}\right), \end{aligned} \quad (44)$$

$$a_{nk} = \frac{1}{2}(\tilde{\nu}_k + D), \quad (45)$$

$$b_{nk} = \frac{1}{2}(\tilde{\nu}_k + D/\tilde{\eta}_k + \text{tr} \tilde{\Sigma}_k^{-1} (x_n - \tilde{\mu}_k)(x_n - \tilde{\mu}_k)^\top). \quad (46)$$

The definitions of the parameters used in Eqs. (37)–(42) are listed as follows:

$$\tilde{N}_k = \sum_n \tilde{z}_{nk}, \quad (47)$$

$$\tilde{M}_k = \sum_n \tilde{z}_{nk} \tilde{u}_{nk}, \quad (48)$$

$$\hat{M}_k = \sum_n \tilde{z}_{nk} \log \hat{u}_{nk}, \quad (49)$$

$$\tilde{\mu}_k = \frac{1}{\tilde{M}_k} \sum_n \tilde{z}_{nk} \tilde{u}_{nk} x_n, \quad (50)$$

$$\tilde{\Sigma}_k = \frac{1}{\tilde{M}_k} \sum_n \tilde{z}_{nk} \tilde{u}_{nk} (x_n - \tilde{\mu}_k)(x_n - \tilde{\mu}_k)^\top. \quad (51)$$

References

- [1] H. Akaike, A new look at the statistical model identification, IEEE Transactions on Automatic Control 19 (1974) 716–723.
- [2] A. Archambeau, M. Verleysen, Robust Bayesian clustering, Neural Networks 20 (2007) 129–138.
- [3] H. Attias, Inferring parameters and structure of latent variable models by variational Bayes, in: Proceeding of the 15th Conference on Uncertainty in Artificial Intelligence, 1999, pp. 21–30.
- [4] A.P. Dempster, N.M. Laird, D.B. Rubin, Maximum-likelihood from incomplete data via the EM algorithm, Journal of the Royal Statistics Society Series B 39 (1977) 1–38.
- [5] S. Oba, M. Sato, S. Ishii, Variational Bayes method for mixture of principal component analyzers, in: Proceeding for 7th International Conference on Neural Information Processing, vol. 2, 2000, pp. 1416–1421.
- [6] S. Richardson, P. Green, On Bayesian analysis of mixtures with unknown number of components, Journal of the Royal Statistical Society Series B 59 (1997) 731–792.
- [7] G. Schwarz, Estimating the dimension of a model, Annals of Statistics 6 (1978) 461–464.
- [8] M. Svensén, C.M. Bishop, Robust Bayesian mixture modeling, Neurocomputing 64 (2005) 235–252.
- [9] T. Takekawa, S. Kang, Y. Isomura, T. Fukai, Robust and accurate spike sorting with matching pursuit and variational Bayesian clustering, in: Abstract of the 37th Annual Meeting of the Society for Neuroscience, 2007, 790.7.



Takashi Takekawa received the Ph.D. degree in Informatics from Kyoto University, Kyoto, in 2005. Currently, he is a Special Postdoctoral Researcher in RIKEN Brain Science Institute (BSI). His research interests include computational functions of neuronal dynamics. Since 2007, he has developed an accurate spike sorting method for more detailed analysis of experimental data.



Tomoki Fukai received a Ph.D degree in Physics from Waseda University, Tokyo, in 1985. Since 2005, he has been serving as the director of Theoretical Neuroscience Group of RIKEN Brain Science Institute. His main interests are in modeling dynamics and computational functions of neuronal networks.