



# Divergence measures and a general framework for local variational approximation

Kazuho Watanabe<sup>a,\*</sup>, Masato Okada<sup>b</sup>, Kazushi Ikeda<sup>a</sup>

<sup>a</sup> Graduate School of Information Science, Nara Institute of Science and Technology, 8916-5, Takayama-cho, Ikoma, Nara, 630-0192, Japan

<sup>b</sup> Graduate School of Frontier Sciences, The University of Tokyo, 5-1-5, Kashiwanoha, Kashiwa, 277-8561, Japan

## ARTICLE INFO

### Article history:

Received 8 November 2010

Received in revised form 17 May 2011

Accepted 6 June 2011

### Keywords:

Local variational approximation

Bayesian learning

Divergence measure

Variational Bayes

## ABSTRACT

The local variational method is a technique to approximate an intractable posterior distribution in Bayesian learning. This article formulates a general framework for local variational approximation and shows that its objective function is decomposable into the sum of the Kullback information and the expected Bregman divergence from the approximating posterior distribution to the Bayesian posterior distribution. Based on a geometrical argument in the space of approximating posteriors, we propose an efficient method to evaluate an upper bound of the marginal likelihood. Moreover, we demonstrate that the variational Bayesian approach for the latent variable models can be viewed as a special case of this general framework.

© 2011 Elsevier Ltd. All rights reserved.

## 1. Introduction

Bayesian learning is in wide use in many applied data-modeling problems, and it is often accompanied by approximation schemes since it requires intractable computation of posterior distributions. Local variational approximation (LVA), also known as direct site bounding, is a technique to approximate the Bayesian posterior distribution by an analytically computable one. Bishop (2006) provides a tutorial on this scheme in terms of convex duality theory (Jordan, Ghahramani, Jaakkola, & Saul, 1999). Representative applications of this scheme include logistic regression (Bishop, 2006; Jaakkola & Jordan, 2000), the Gaussian process classifier (Gibbs & MacKay, 2000), mixture of experts (Bishop & Svensen, 2003), and sparse linear models (Seeger, 2008, 2009; Wipf & Rao, 2007).

LVA forms lower and/or upper bounds of the unnormalized posterior distribution. The approximation is optimized so as to maximize (minimize) the lower (upper) bound of the normalizing factor, also known as the marginal likelihood of the Bayesian posterior distribution. However, learning algorithms based on LVA have been derived on a case-by-case basis since LVA lacks a general framework and the principle behind it has yet to be described.

In this article, we investigate the relationship between information divergences and LVA to provide a general framework for this scheme. More specifically, we first show that the discrepancy between the log-marginal likelihood, also known as the free energy, and its bound by LVA is expressed as the sum of the Kullback information and the expected Bregman divergence. Decomposing the

bounds of the free energy provides a general interpretation of the optimization of LVA in terms of the minimization of the respective divergences. Moreover, it provides an efficient method for computing the lower bound of the free energy by using only the result of the upper bound minimization. We demonstrate how these results apply to practical models by taking two examples. One is for the kernelized logistic regression model, also known as the relevance vector machine. Another is for latent variable models such as the Gaussian mixture model and the hidden Markov model. We show that the so-called variational Bayesian method for the latent variable models can be viewed as a special case of LVA.

The rest of the paper is organized as follows. Section 2 summarizes the general framework of LVA. Section 3 derives equalities for the information divergences related to LVA. Section 4 elaborates on optimization of the approximation and describes an efficient method to combine the lower and upper bounds of the free energy. In Section 5, we demonstrate two applications of LVA to concrete models. Section 6 provides discussions on the comparison of general LVA with other approximation schemes. Section 7 concludes this paper.

## 2. Local variational approximation

Assume that we are given training examples or observations  $\mathbf{t} = \{t_1, t_2, \dots, t_n\}$ , where each observation  $t_i$  is defined in some domain. Let  $\mathbf{w} \in R^d$  be the parameter vector, and consider Bayesian learning for a model  $p(\mathbf{t}|\mathbf{w})$ .<sup>1</sup> By using the prior distribution  $p_0(\mathbf{w})$ ,

<sup>1</sup> The formulation in this paper also applies to discriminative or regression models by simply replacing  $p(\mathbf{t}|\mathbf{w})$  with the conditional distribution  $p(\mathbf{t}|\mathbf{x}, \mathbf{w})$  of the outputs  $\mathbf{t} = \{t_1, \dots, t_n\}$  given the inputs  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$ . Section 5.1 provides an example of this case.

\* Corresponding author.

E-mail address: [wkazuho@is.naist.jp](mailto:wkazuho@is.naist.jp) (K. Watanabe).

the Bayesian posterior distribution of the parameter  $\mathbf{w}$  is defined by

$$p(\mathbf{w}|\mathbf{t}) = \frac{p(\mathbf{t}|\mathbf{w})p_0(\mathbf{w})}{\int p(\mathbf{t}|\mathbf{w})p_0(\mathbf{w})d\mathbf{w}} = \frac{p(\mathbf{w}, \mathbf{t})}{Z}. \quad (1)$$

Then Bayesian learning requires calculating the predictive distribution  $p(t|\mathbf{t}) = \int p(t|\mathbf{w})p(\mathbf{w}|\mathbf{t})d\mathbf{w}$  as an estimate of the underlying distribution or, for example, the posterior mean  $\int \mathbf{w}p(\mathbf{w}|\mathbf{t})d\mathbf{w}$  as an estimate of the unknown parameter  $\mathbf{w}$ .

However, as is often the case, the normalizing constant of the posterior distribution (1),

$$Z = p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{w})p_0(\mathbf{w})d\mathbf{w},$$

called the marginal likelihood, is analytically intractable, and so is the Bayesian posterior distribution (1). Examples include the case of the logistic regression model (the Bernoulli likelihood) with the Gaussian prior distribution (Jaakkola & Jordan, 2000) and that of the Gaussian model with the sparsity-inducing prior (Seeger, 2008).

**Example 1.** Consider the following simple example of the pair of the Bernoulli distribution for  $t \in \{0, 1\}$  and the Gaussian prior distribution for  $w \in R$ ,

$$p(t|w) = \exp\{tw - \log(1 + e^w)\} \quad (2)$$

and

$$p_0(w) = \frac{1}{\sqrt{2\pi}} e^{-w^2/2}.$$

The Bernoulli distribution has the probability of success ( $t = 1$ )  $\frac{e^w}{1+e^w}$ . Given the observation  $t$ , the integration

$$\int p(t|w)p_0(w)dw$$

is not analytically tractable, and hence prevents the calculation of the posterior distribution  $p(w|\mathbf{t})$ . This is mainly due to the existence of the term  $-\log(1 + e^w)$  in the exponent of Eq. (2). We will revisit this example in Sections 3 and 5.1.

LVA forms upper and/or lower bounds of the joint distribution  $p(\mathbf{w}, \mathbf{t})$ , denoted by  $\bar{p}_\lambda(\mathbf{w}, \mathbf{t})$  and  $\underline{p}_\xi(\mathbf{w}, \mathbf{t})$ , respectively. If the bounds satisfying

$$\bar{p}_\lambda(\mathbf{w}, \mathbf{t}) \geq p(\mathbf{w}, \mathbf{t}) \quad \text{and} \quad (3)$$

$$\underline{p}_\xi(\mathbf{w}, \mathbf{t}) \leq p(\mathbf{w}, \mathbf{t}), \quad (4)$$

for all  $\mathbf{w}$  and  $\mathbf{t}$ , are analytically integrable, then, by normalizing the bounds instead of  $p(\mathbf{w}, \mathbf{t})$ , LVA approximates the posterior distribution by

$$p_\lambda(\mathbf{w}|\mathbf{t}) = \frac{\bar{p}_\lambda(\mathbf{w}, \mathbf{t})}{\bar{Z}(\lambda)}, \quad \text{and} \quad (5)$$

$$p_\xi(\mathbf{w}|\mathbf{t}) = \frac{\underline{p}_\xi(\mathbf{w}, \mathbf{t})}{\underline{Z}(\xi)}, \quad (6)$$

respectively, where  $\bar{Z}(\lambda)$  and  $\underline{Z}(\xi)$  are the normalization constants defined by

$$\bar{Z}(\lambda) = \int \bar{p}_\lambda(\mathbf{w}, \mathbf{t})d\mathbf{w}$$

and

$$\underline{Z}(\xi) = \int \underline{p}_\xi(\mathbf{w}, \mathbf{t})d\mathbf{w}.$$

Here,  $\lambda$  and  $\xi$  are called the variational parameters, which are introduced to make the bounds adjustable.

The respective approximations are optimized by estimating the variational parameters,  $\xi$  and  $\lambda$ , so that  $\underline{Z}(\xi)$  is maximized and  $\bar{Z}(\lambda)$  is minimized, since the inequalities

$$\underline{Z}(\xi) \leq Z \leq \bar{Z}(\lambda) \quad (7)$$

hold by definition.

To consider the respective LVAs in terms of information divergence in later sections, let us introduce the free energy,

$$F = -\log Z,$$

following statistical mechanical terminology, and its lower and upper bounds,  $\underline{F}(\lambda) = -\log \bar{Z}(\lambda)$  and  $\bar{F}(\xi) = -\log \underline{Z}(\xi)$ . By taking the negative logarithms on both sides of Eq. (7), we have

$$\underline{F}(\lambda) \leq F \leq \bar{F}(\xi). \quad (8)$$

Hereafter, we follow the measure of the free energy and adopt the phrases the lower bound maximization ( $\underline{F}(\lambda)$  maximization) and the upper bound minimization ( $\bar{F}(\xi)$  minimization) to signify the respective local variational approximations (5) and (6).

### 3. Divergence measures in LVA

In this section, we derive key equations relating LVA with information divergence. Most existing LVA techniques are based on the convexity of the log-likelihood function or the log-prior (Bishop, 2006; Seeger, 2008). We describe these cases by using general convex functions,  $\phi$  and  $\psi$ , and show that the objective functions

$$\bar{F}(\xi) - F = \log \frac{Z}{\underline{Z}(\xi)} \geq 0$$

and

$$F - \underline{F}(\lambda) = \log \frac{\bar{Z}(\lambda)}{Z} \geq 0$$

to be minimized in the approximations (5) and (6) are decomposable into the sum of the Kullback information and the expected Bregman divergence.

Let  $\phi$  and  $\psi$  be twice differentiable real-valued strictly convex functions, and denote by  $d_\phi$  the Bregman divergence associated with the function  $\phi$  (Banerjee, Merugu, Dhillon, & Ghosh, 2005),

$$d_\phi(\mathbf{v}_1, \mathbf{v}_2) = \phi(\mathbf{v}_1) - \phi(\mathbf{v}_2) - (\mathbf{v}_1 - \mathbf{v}_2) \cdot \nabla \phi(\mathbf{v}_2) \geq 0, \quad (9)$$

where  $\nabla \phi(\mathbf{v}_2)$  denotes the gradient vector of  $\phi$  at  $\mathbf{v}_2$ .

Let us consider the case when  $\phi$  and  $\psi$  are respectively used to form the following bounds of the joint distribution  $p(\mathbf{w}, \mathbf{t})$ ,

$$\underline{p}_\xi(\mathbf{w}, \mathbf{t}) = p(\mathbf{w}, \mathbf{t}) \exp\{-d_\phi(\mathbf{h}(\mathbf{w}), \mathbf{h}(\xi))\}, \quad (10)$$

$$\bar{p}_\lambda(\mathbf{w}, \mathbf{t}) = p(\mathbf{w}, \mathbf{t}) \exp\{d_\psi(\mathbf{g}(\mathbf{w}), \mathbf{g}(\lambda))\}, \quad (11)$$

where  $\mathbf{h}$  and  $\mathbf{g}$  are vector-valued functions of  $\mathbf{w}$ .<sup>2</sup>

Eq. (10) is interpreted as follows.  $\log p(\mathbf{w}, \mathbf{t})$  includes a term that prevents analytic integration of  $p(\mathbf{w}, \mathbf{t})$  with respect to  $\mathbf{w}$ . If such a term is expressed by the convex function  $\phi$  of some function  $\mathbf{h}$  transforming  $\mathbf{w}$ , it is replaced by the tangent hyperplane,  $\phi(\mathbf{h}(\xi)) + (\mathbf{h}(\mathbf{w}) - \mathbf{h}(\xi)) \cdot \nabla \phi(\mathbf{h}(\xi))$ , so that  $\log \underline{p}_\xi(\mathbf{w}, \mathbf{t})$  makes a simpler function of  $\mathbf{w}$ , such as a quadratic function. Remember that, if  $\log \underline{p}_\xi(\mathbf{w}, \mathbf{t})$  is quadratic with respect to  $\mathbf{w}$ ,  $\underline{p}_\xi(\mathbf{w}, \mathbf{t})$  is analytically integrable by the Gaussian integral.

<sup>2</sup> The functions  $\mathbf{g}$  and  $\mathbf{h}$  (also  $\psi$  and  $\phi$ ) can be dependent on  $\mathbf{t}$  in this discussion. However, we denote them as if they are independent of  $\mathbf{t}$  for simplicity. They are actually independent of  $\mathbf{t}$  in the example of Section 5.1 and in most applications (Bishop, 2006; Seeger, 2008).

Download English Version:

<https://daneshyari.com/en/article/404289>

Download Persian Version:

<https://daneshyari.com/article/404289>

[Daneshyari.com](https://daneshyari.com)