



## Variational Bayes With Intractable Likelihood

Minh-Ngoc Tran, David J. Nott & Robert Kohn

To cite this article: Minh-Ngoc Tran, David J. Nott & Robert Kohn (2017) Variational Bayes With Intractable Likelihood, Journal of Computational and Graphical Statistics, 26:4, 873-882, DOI: [10.1080/10618600.2017.1330205](https://doi.org/10.1080/10618600.2017.1330205)

To link to this article: <https://doi.org/10.1080/10618600.2017.1330205>



View supplementary material [↗](#)



Accepted author version posted online: 19 May 2017.  
Published online: 11 Oct 2017.



Submit your article to this journal [↗](#)



Article views: 774



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 11 View citing articles [↗](#)



## Variational Bayes With Intractable Likelihood

Minh-Ngoc Tran<sup>a</sup>, David J. Nott<sup>b</sup>, and Robert Kohn<sup>c</sup>

<sup>a</sup>University of Sydney Business School, University of Sydney, NSW, Australia; <sup>b</sup>Department of Statistics and Applied Probability, National University of Singapore, Singapore; <sup>c</sup>UNSW Business School, University of New South Wales, Sydney NSW, Australia

### ABSTRACT

Variational Bayes (VB) is rapidly becoming a popular tool for Bayesian inference in statistical modeling. However, the existing VB algorithms are restricted to cases where the likelihood is tractable, which precludes their use in many interesting situations such as in state–space models and in approximate Bayesian computation (ABC), where application of VB methods was previously impossible. This article extends the scope of application of VB to cases where the likelihood is intractable, but can be estimated unbiasedly. The proposed VB method therefore makes it possible to carry out Bayesian inference in many statistical applications, including state–space models and ABC. The method is generic in the sense that it can be applied to almost all statistical models without requiring too much model-based derivation, which is a drawback of many existing VB algorithms. We also show how the proposed method can be used to obtain highly accurate VB approximations of marginal posterior distributions. Supplementary material for this article is available online.

### ARTICLE HISTORY

Received January 2016  
Revised March 2017

### KEYWORDS

Approximate Bayesian computation; Marginal likelihood; Natural gradient; Quasi-Monte Carlo; State–space models; Stochastic optimization

### 1. Introduction

Let  $y$  be the data and  $\theta \in \Theta$  the vector of model parameters. Bayesian inference about  $\theta$  is based on the posterior distribution  $\pi(\theta) = p(\theta|y) \propto p(\theta)p(y|\theta)$ , with  $p(\theta)$  the prior and  $p(y|\theta)$  the likelihood function. In this article, we are interested in variational Bayes (VB), which is widely used as a computationally effective method for approximating the posterior distribution  $\pi(\theta)$  (Bishop 2006). VB approximates the posterior by a distribution  $q(\theta)$  within some tractable class, such as an exponential family, chosen to minimize the Kullback–Leibler divergence between  $q(\theta)$  and  $\pi(\theta)$ . Most of the VB algorithms in the literature require that the likelihood  $p(y|\theta)$  can be computed analytically for any  $\theta$ .

In many applications, however, the likelihood  $p(y|\theta)$  is intractable in the sense that it is infeasible to compute  $p(y|\theta)$  exactly at each value of  $\theta$ , which makes it difficult to use VB for inference. For example, in state–space models (Durbin and Koopman 2001), which are widely used in economics, finance, and engineering, the likelihood is a high-dimensional integral over the state variables governed by a Markov process. Ghahramani and Hinton (2000) were the first to use VB for inference in state–space models. However, they only considered the special case in which the time series is segmented into regimes with each regime assumed to follow a linear-Gaussian state–space model. For general state–space models, it is still a challenging problem to do inference with VB. Turner and Sahani (2011) discussed some of the difficulties in applying VB methods to time series models. Another example where implementing VB is difficult is ABC (Tavare et al. 1997; Peters, Sisson, and Fan 2012). ABC methods provide a way of

approximating the posterior  $\pi(\theta)$  when the likelihood is difficult to compute but it is possible to simulate data from the model. We are not aware of any work that uses VB for inference in ABC. Our article proposes a generic VB algorithm that approximates  $\pi(\theta)$  when the likelihood is intractable. The only requirement is that the intractable likelihood can be estimated unbiasedly. The proposed algorithm therefore makes it possible to carry out variational Bayes inference in many statistical models with an intractable likelihood, where this was previously impossible.

In many models, by introducing a latent variable  $\alpha$ , the joint density  $p(y, \alpha|\theta)$  is tractable. This makes it much easier to work with the joint posterior  $p(\theta, \alpha|y) \propto p(\theta)p(y, \alpha|\theta)$  rather than the marginal posterior of interest  $\pi(\theta)$  itself. In this situation, many VB algorithms in the literature approximate the joint posterior  $p(\theta, \alpha|y)$  by a factorized distribution  $q(\theta)q(\alpha)$ , and then use  $q(\theta)$  as an approximation to  $\pi(\theta)$ . The main drawback of this approach is that the (usually high) posterior dependence between  $\theta$  and  $\alpha$  is ignored, which might lead to a poor VB approximation (Neville, Ormerod, and Wand 2014). Our VB algorithm approximates  $\pi(\theta)$  directly with the latent variable  $\alpha$  integrated out and thus overcomes this drawback; see the example in Section 5.1.

Section 2 presents our approach, which we call variational Bayes with intractable likelihood (VBIL), when the likelihood can be estimated unbiasedly. VBIL transforms the problem of approximating the posterior  $\pi(\theta)$  into a stochastic optimization problem using a noisy gradient. It is essential for the success of stochastic optimization algorithms to have a gradient estimator with a sufficiently small variance. Section 3 describes

several techniques, including control variate and quasi-Monte Carlo, for reducing the variance of the estimated gradient. This section also discusses the importance of the natural gradient (Amari 1998), which takes into account the geometry of the variational distribution  $q(\theta)$  being learned.

Unlike many VB algorithms that are derived on a model-by-model basis and require analytical computation of some model-based expectations, one of the main advantages of VBIL is that it can be applied to almost all statistical models without requiring an analytical solution to model-based expectations. The only requirement is that we are able to estimate the intractable likelihood unbiasedly. The VBIL methodology is therefore more generic and widely applicable. As a by-product, VBIL provides an estimate of the marginal likelihood, which is useful for model choice.

There are several lines of work related to ours in terms of working with an intractable likelihood. Beaumont (2003) and Andrieu and Roberts (2009) showed that Markov chain Monte Carlo simulation based on an unbiased estimator of the likelihood is still able to generate samples from the posterior. This method is known in the literature as the pseudo-marginal approach (PM). More efficient variants of PM, called correlated PM and blocking PM, have been proposed recently (Deligianidis, Doucet, and Pitt 2016; Tran et al. 2016). Tran et al. (2013) showed that importance sampling with the likelihood replaced by its unbiased estimator is still valid for estimating expectations with respect to the posterior, and name their method as importance sampling squared (IS<sup>2</sup>). The main advantage of VBIL is that it is several orders of magnitude faster than these competitors.

Section 4 studies the link between the precision of the likelihood estimator to the variance of the VBIL estimator. This helps to understand how much accuracy is lost when working with an estimated likelihood compared to the case when the likelihood is available. In this spirit, Pitt et al. (2012) and Tran et al. (2013) showed that the asymptotic variance of PM and IS<sup>2</sup> estimators increases exponentially with the variance of the log of the likelihood estimator. Therefore, it is critical for these methods to have a likelihood estimator that is accurate enough. They show that the variance of the log of the estimated likelihood should be around 1 to minimize the computing time that is needed for the variance of PM and IS<sup>2</sup> estimators to have a fixed precision. For VBIL, we show that the asymptotic variance of VBIL estimators increases linearly with the variance of the log-likelihood estimator. The proposed methodology is therefore useful when only highly variable estimates of the likelihood are available. We discuss such a situation in Section 5.1 where VBIL works well while its competitors fail.

Several interesting applications of VBIL are presented in Section 5. Section 5.1 shows the use of VBIL for generalized linear mixed models and demonstrates the high accuracy of VBIL compared to the existing VB algorithms. Section 5.2 applies VBIL to Bayesian inference in state-space models and Section 5.3 shows how VBIL can be used for ABC. To the best of our knowledge, our article is the first to use a VB method in the most general way for Bayesian inference in state-space models and ABC. Another interesting application of VBIL is presented in Section 5.4, in which we illustrate that VBIL provides an attractive way to improve the accuracy of VB approximations

of marginal posteriors. Proof and technical details can be found in an online Appendix.

## 2. Variational Bayes With an Intractable Likelihood

This section describes the basic form of the proposed VBIL algorithm, where an unbiased estimator of the likelihood is available. Suppose  $\hat{p}_N(y|\theta)$  is an unbiased estimator of the likelihood  $p(y|\theta)$ , where  $N$  is an algorithmic parameter relating to the precision in estimating the likelihood, such as the number of samples if the likelihood is estimated by importance sampling or the number of particles if the likelihood in state-space models is estimated by a particle filter. Using the terminology in Pitt et al. (2012), we refer to  $N$  as the number of particles. Let  $z = \log \hat{p}_N(y|\theta) - \log p(y|\theta)$ , so that  $\hat{p}_N(y|\theta) = p(y|\theta)e^z$ , and denote by  $g_N(z|\theta)$  the density of  $z$ . Note that  $z$  is unknown as we do not know  $\log p(y|\theta)$ . As will become clear shortly, we never require  $z$  in practice. We sometimes write  $\hat{p}_N(y|\theta)$  as  $\hat{p}_N(y|\theta, z)$ . We note that  $\int e^z g_N(z|\theta) dz = 1$  because of the unbiasedness of the estimator  $\hat{p}_N(y|\theta)$ . Define the following density on the extended space  $\Theta \times \mathbb{R}$

$$\pi_N(\theta, z) = \frac{p(\theta)p(y|\theta)e^z g_N(z|\theta)}{p(y)} = \pi(\theta)e^z g_N(z|\theta).$$

This augmented density admits the posterior of interest  $\pi(\theta)$  as its marginal. It is useful to work with  $\pi_N(\theta, z)$  as the high-dimensional vector of random variables involved in estimating the likelihood is transformed into the scalar  $z$ . A direct approximation of  $\pi_N(\theta, z)$  is  $\tilde{q}_{\lambda,N}(\theta, z) = q_\lambda(\theta)e^z g_N(z|\theta)$ , where  $q_\lambda(\theta)$  is the variational distribution with the variational parameter  $\lambda$  to be estimated, and then  $q_\lambda(\theta)$  can be used as an approximation of  $\pi(\theta)$ . However, it turns out that it is impossible to estimate the gradient of the Kullback–Leibler divergence between  $\tilde{q}_{\lambda,N}(\theta, z)$  and  $\pi_N(\theta, z)$  as this requires knowing  $z$ .

We propose instead to approximate  $\pi_N(\theta, z)$  by  $q_{\lambda,N}(\theta, z) = q_\lambda(\theta)g_N(z|\theta)$ . This augmented density has the attractive features that  $q_\lambda(\theta)$  is its marginal for  $\theta$  and it is possible to estimate the gradient of the Kullback–Leibler divergence  $\text{KL}(\lambda)$  between  $q_{\lambda,N}(\theta, z)$  and  $\pi_N(\theta, z)$  (see Equation (2)). Although  $q_{\lambda,N}(\theta, z)$  does not provide a good approximation of the posterior marginal of  $z$ , the latter is not of interest to us. Furthermore, under Assumptions 1 and 2 given in Section 4, the minimization of  $\text{KL}(\lambda)$  is equivalent to the minimization of the KL divergence between  $q_\lambda(\theta)$  and  $\pi(\theta)$ .

The Kullback–Leibler divergence between  $q_{\lambda,N}(\theta, z)$  and  $\pi_N(\theta, z)$  is

$$\text{KL}(\lambda) = \int q_\lambda(\theta)g_N(z|\theta) \log \frac{q_\lambda(\theta)g_N(z|\theta)}{\pi_N(\theta, z)} dz d\theta, \quad (1)$$

where we omit to indicate dependence on  $N$  for notational convenience. The gradient of  $\text{KL}(\lambda)$  is

$$\begin{aligned} \nabla_\lambda \text{KL}(\lambda) &= \int (q_\lambda(\theta)g_N(z|\theta) \nabla_\lambda [\log q_\lambda(\theta)] \\ &\quad \times (\log q_\lambda(\theta) - \log(p(\theta)\hat{p}_N(y|\theta, z))) dz d\theta \\ &= \mathbb{E}_{\theta \sim q_\lambda(\theta), z \sim g_N(z|\theta)} (\nabla_\lambda [\log q_\lambda(\theta)] \\ &\quad \times (\log q_\lambda(\theta) - \log(p(\theta)\hat{p}_N(y|\theta, z))). \end{aligned} \quad (2)$$

Here, we have used the facts that  $\nabla_{\lambda}[q_{\lambda}(\theta)] = q_{\lambda}(\theta)\nabla_{\lambda}[\log q_{\lambda}(\theta)]$  and that  $\mathbb{E}(\nabla_{\lambda}[\log q_{\lambda}(\theta)]) = 0$ . It follows from (2) that, by generating  $\theta \sim q_{\lambda}(\theta)$  and  $z \sim g_N(z|\theta)$ , it is straightforward to obtain an unbiased estimator  $\widehat{\nabla_{\lambda}\text{KL}}(\lambda)$  of the gradient  $\nabla_{\lambda}\text{KL}(\lambda)$ . Therefore, we can use stochastic optimization to optimize  $\text{KL}(\lambda)$ . We note that  $z$  is never dealt with explicitly and it only plays a theoretical role in the mathematical derivations. The basic algorithm is as follows

**Algorithm 1.**

- Initialize  $\lambda^{(0)}$  and stop the following iteration if the stopping criterion is met.
- For  $t = 0, 1, \dots$ , compute  $\lambda^{(t+1)} = \lambda^{(t)} - a_t \widehat{\nabla_{\lambda}\text{KL}}(\lambda^{(t)})$ .

We will refer to this algorithm as variational Bayes with intractable likelihood (VBIL). The sequence  $\{a_t\}$  should satisfy  $a_t > 0$ ,  $\sum_t a_t = \infty$  and  $\sum_t a_t^2 < \infty$ . We choose  $a_t = 1/(1+t)$  in this article, but it is also possible to train  $a_t$  adaptively. It is important to note that each iteration is parallelizable, as the gradient  $\nabla_{\lambda}\text{KL}(\lambda)$  is estimated by independent samples from  $q_{\lambda,N}(\theta, z)$ .

## 2.1. Stopping Criterion and Marginal Likelihood Estimation

An easy-to-implement stopping rule is to stop the updating procedure if the change between  $\lambda^{(t+1)}$  and  $\lambda^{(t)}$ , for example, in terms of the Euclidean distance, is less than some threshold  $\epsilon$  (Ranganath, Gerrish, and Blei 2014). However, it is difficult to select  $\epsilon$  as such a distance depends on the scales and the length of the vector  $\lambda$ . It is easy to show that  $\log p(y) \geq \text{LB}(\lambda)$ , where

$$\text{LB}(\lambda) = \mathbb{E}_{\theta,z}[\log p(\theta) - \log q_{\lambda}(\theta) + \log \hat{p}_N(y|\theta, z)] \quad (3)$$

is the lower bound on the log of the marginal likelihood  $\log p(y)$ . This lower bound after convergence can be used as an approximation to  $\log p(y)$ , which is useful for model selection purposes. The expectation of the first two terms in (3) can be computed analytically, while the last term can be estimated unbiasedly by samples from  $q_{\lambda,N}(\theta, z)$ . However, in our experience, estimating the entire expectation (3) based on samples from  $q_{\lambda,N}(\theta, z)$  leads to a smaller variance. Denote by  $\widehat{\text{LB}}(\lambda)$  the resulting unbiased estimate of  $\text{LB}(\lambda)$ . Although

$\text{LB}(\lambda)$  is strictly nondecreasing over iterations, its sample estimate  $\widehat{\text{LB}}(\lambda)$  might not be. To account for this, we suggest to stop the updating procedure if the change in an averaged value of the lower bounds over a window of  $M$  iterations,  $\overline{\text{LB}}(\lambda_t) = (1/M) \sum_{k=1}^M \widehat{\text{LB}}(\lambda_{t-k+1})$ , is less than some threshold  $\epsilon$ . At convergence, the values  $\widehat{\text{LB}}(\lambda_t)$  stay the same, therefore  $\overline{\text{LB}}(\lambda_t)$  will average out the noise in  $\widehat{\text{LB}}(\lambda_t)$  and is stable. Furthermore, we suggest replacing  $\widehat{\text{LB}}(\lambda)$  by a scaled version of it,  $\widehat{\text{LB}}(\lambda)/n$  with  $n$  the size of the dataset such as the number of observations. The scaled lower bound is more or less independent of the size of the dataset (see Figure 1). We set  $M = 5$  and  $\epsilon = 10^{-5}$  in this article.

## 3. Variance Reduction and Natural Gradient

As is typical of stochastic optimization algorithms, the performance of Algorithm 1 depends greatly on the variance of the noisy gradient. This section describes several techniques for variance reduction.

### 3.1. Control Variate

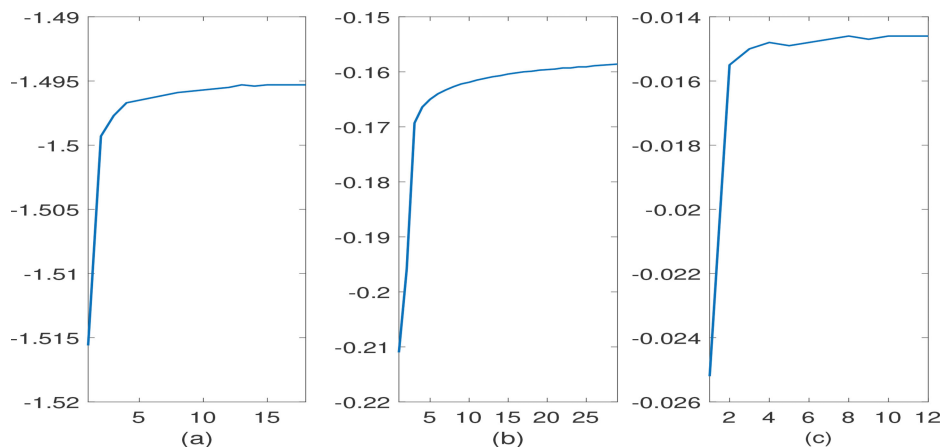
Denote  $\hat{h}(\theta, z) = \log(p(\theta)\hat{p}_N(y|\theta, z))$  for notational simplicity. Let  $\theta_s \sim q_{\lambda}(\theta)$  and  $z_s \sim g_N(z|\theta_s)$ ,  $s = 1, \dots, S$ , be  $S$  samples from the variational distribution  $q_{\lambda,N}(\theta, z)$ . A naive estimator of the  $i$ th element of  $\nabla_{\lambda}\text{KL}(\lambda)$  is

$$\widehat{\nabla_{\lambda_i}\text{KL}}(\lambda)^{\text{naive}} = \frac{1}{S} \sum_{s=1}^S \nabla_{\lambda_i}[\log q_{\lambda}(\theta_s)] (\log q_{\lambda}(\theta_s) - \hat{h}(\theta_s, z_s)), \quad (4)$$

whose variance is often too large to be useful. For any number  $c_i$ , consider

$$\widehat{\nabla_{\lambda_i}\text{KL}}(\lambda) = \frac{1}{S} \sum_{s=1}^S \nabla_{\lambda_i}[\log q_{\lambda}(\theta_s)] (\log q_{\lambda}(\theta_s) - \hat{h}(\theta_s, z_s) - c_i), \quad (5)$$

which is still an unbiased estimator of  $\nabla_{\lambda_i}\text{KL}(\lambda)$  since  $\mathbb{E}(\nabla_{\lambda}[\log q_{\lambda}(\theta)]) = 0$ , whose variance can be greatly reduced by an appropriate choice of  $c_i$ . Similar ideas are considered in the literature; see Paisley, Blei, and Jordan (2012), Nott et al. (2012), and Ranganath, Gerrish, and Blei (2014). The variance



**Figure 1.** Plots of scaled lower bounds over the iterations: (a) Six City example, (b) state-space example, (c) ABC example. The curves are smooth, which show that the update schemes converge quickly after a few iterations.

of  $\widehat{\nabla_{\lambda_i} \text{KL}(\lambda)}$  is

$$\frac{1}{S} \mathbb{V}(\nabla_{\lambda_i}[\log q_{\lambda}(\theta)](\log q_{\lambda}(\theta) - \widehat{h}(\theta, z))) + \frac{c_i^2}{S} \mathbb{V}(\nabla_{\lambda_i}[\log q_{\lambda}(\theta)]) - \frac{2c_i}{S} \text{cov}(\nabla_{\lambda_i}[\log q_{\lambda}(\theta)](\log q_{\lambda}(\theta) - \widehat{h}(\theta, z)), \nabla_{\lambda_i}[\log q_{\lambda}(\theta)]).$$

The optimal  $c_i$  that minimizes this variance is

$$c_i = \frac{\text{cov}(\nabla_{\lambda_i}[\log q_{\lambda}(\theta)](\log q_{\lambda}(\theta) - \widehat{h}(\theta, z)), \nabla_{\lambda_i}[\log q_{\lambda}(\theta)])}{\mathbb{V}(\nabla_{\lambda_i}[\log q_{\lambda}(\theta)])}. \quad (6)$$

Then  $\mathbb{V}(\widehat{\nabla_{\lambda_i} \text{KL}(\lambda)}) = \mathbb{V}(\widehat{\nabla_{\lambda_i} \text{KL}(\lambda)}^{\text{naive}})(1 - \rho_i^2) \leq \mathbb{V}(\widehat{\nabla_{\lambda_i} \text{KL}(\lambda)}^{\text{naive}})$ , where  $\rho_i$  is the correlation between  $\nabla_{\lambda_i}[\log q_{\lambda}(\theta)](\log q_{\lambda}(\theta) - \widehat{h}(\theta, z))$  and  $\nabla_{\lambda_i}[\log q_{\lambda}(\theta)]$ . Often,  $\rho_i^2$  is very close to 1.

We estimate the numbers  $c_i$  by samples  $(\theta_s, z_s) \sim q_{\lambda, N}(\theta, z)$  as in (6). To ensure the unbiasedness of the gradient estimator, the samples used to estimate  $c_i$  must be independent of the samples used to estimate the gradient. In practice, the  $c_i$  can be updated sequentially as follows. At iteration  $t$ , we use the  $c_i$  computed in the previous iteration  $t - 1$ , that is, based on the samples from  $q_{\lambda^{(t-1)}, N}(\theta, z)$ , to estimate the gradient  $\widehat{\nabla_{\lambda} \text{KL}(\lambda^{(t)})}$ , which is estimated using new samples from  $q_{\lambda^{(t)}, N}(\theta, z)$ . We then update the  $c_i$  using this new set of samples. By doing so, the unbiasedness is guaranteed while no extra samples are needed in updating the numbers  $c_i$ .

The gradient in the form of (2) can be written as a sum of two terms, where the first term  $\mathbb{E}_{\theta \sim q_{\lambda}(\theta)}(\nabla_{\lambda}[\log q_{\lambda}(\theta)] \log q_{\lambda}(\theta))$  can be in most cases computed analytically. However, as pointed out by a referee, this term should be estimated using the same samples of  $\theta$  as we do in (5). Doing so helps to reduce the noise in estimating the gradient. This is because the first term plays the role of a control variate. This issue was discussed in detail by Salimans and Knowles (2013).

### 3.2. Natural Gradient

Intuitively, a different learning rate should be used for each scale in the gradient vector. That is, the traditional gradient vector  $\nabla_{\lambda} \text{KL}(\lambda)$  should be multiplied by an appropriate scale matrix. It is well-known that the traditional gradient defined on the Euclidean space does not adequately capture the geometry of the variational distribution  $q_{\lambda}(\theta)$  (Amari 1998). A small Euclidean distance between  $\lambda$  and  $\lambda'$  does not necessarily mean a small Kullback–Leibler divergence between  $q_{\lambda}(\theta)$  and  $q_{\lambda'}(\theta)$ . Amari (1998) defined the natural gradient as

$$\nabla_{\lambda} \text{KL}(\lambda)^{\text{natural}} = I_F(\lambda)^{-1} \nabla_{\lambda} \text{KL}(\lambda), \quad (7)$$

with  $I_F(\lambda)$  the Fisher information matrix, and suggests using the natural gradient as an efficient alternative to the traditional gradient. See also Hoffman et al. (2013).

If the variational distribution  $q_{\lambda}(\theta)$  has the exponential family form

$$q_{\lambda}(\theta) = \exp(T(\theta)' \lambda - Z(\lambda)), \quad (8)$$

with  $T(\theta)$  the vector of sufficient statistics and  $\lambda$  the vector of natural parameters, then  $I_F(\lambda) = \text{cov}_{q_{\lambda}}(T(\theta), T(\theta))$  is computed analytically.

The use of the natural gradient in VB algorithms is considered, among others, by Honkela et al. (2010), Hoffman et al. (2013), and Salimans and Knowles (2013). A simple demonstration of the importance of the natural gradient can be found in the Appendix.

### 3.3. Factorized Variational Distribution

Often, the variational distribution  $q_{\lambda}(\theta)$  is factorized into  $K$  factors

$$q_{\lambda}(\theta) = q_{\lambda^{(1)}}(\theta^{(1)}) \dots q_{\lambda^{(K)}}(\theta^{(K)}). \quad (9)$$

Then, each factor  $q_{\lambda^{(k)}}(\theta^{(k)})$  is updated separately and the variance of the estimate of the corresponding gradient can be reduced. Salimans and Knowles (2013) and Ranganath, Gerrish, and Blei (2014) considered variance reduction using factorization. Denote by  $\widehat{h}_k(\theta, z)$  the terms in  $\widehat{h}(\theta, z)$  that involve only  $\theta^{(k)}$  and  $z$ . From (2), and noting that  $\mathbb{E}_{\theta, z}(\nabla_{\lambda^{(k)}}[\log q_{\lambda^{(k)}}(\theta^{(k)})]) = 0$ , the traditional gradient corresponding to factor  $k$  is

$$\nabla_{\lambda^{(k)}} \text{KL}(\lambda) = \mathbb{E}_{\theta, z}(\nabla_{\lambda^{(k)}}[\log q_{\lambda^{(k)}}(\theta^{(k)})] \times (\log q_{\lambda^{(k)}}(\theta^{(k)}) - \widehat{h}_k(\theta, z))). \quad (10)$$

In the case  $q_{\lambda^{(k)}}(\theta^{(k)}) = \exp(T_k(\theta^{(k)})' \lambda^{(k)} - Z_k(\lambda^{(k)}))$  belongs to an exponential family, the natural gradient corresponding to factor  $k$  is  $\nabla_{\lambda^{(k)}} \text{KL}(\lambda)^{\text{natural}} = I_{F, k}(\lambda^{(k)})^{-1} \nabla_{\lambda^{(k)}} \text{KL}(\lambda)$ , with  $I_{F, k}(\lambda^{(k)})$  the information matrix of distribution  $q_{\lambda^{(k)}}(\theta^{(k)})$ .

Estimating the gradient using (10) has less variation than using (2). Intuitively, this is because the variation due to terms not involving  $\theta^{(k)}$  has been removed. Ranganath, Gerrish, and Blei (2014) explained this as a Rao–Blackwellization effect.

### 3.4. Randomized Quasi-Monte Carlo

Numerical integration using quasi-Monte Carlo (QMC) has been proven efficient in many applications. Instead of generating uniform random numbers  $U(0, 1)$  as in plain Monte Carlo methods, QMC generates deterministic sequences that are more evenly distributed in  $(0, 1)$  in the sense that they minimize the so-called star-discrepancy. Dick and Pillichshammer (2010) provided an extensive background on QMC. It is shown that, in many cases, QMC integration achieves a better convergence rate than Monte Carlo integration. In this article, we use randomized quasi-Monte Carlo (RQMC) as VBIL requires an unbiased estimator of the gradient. By introducing a random element into a QMC sequence, RQMC preserves the low-discrepancy property and, at the same time, leads to unbiased estimators (Owen 1997; Dick and Pillichshammer 2010).

Here, we use RQMC to sample  $\theta \sim q_{\lambda}(\theta)$ . This will help to reduce the variance of the noisy gradient if the dimension of  $\theta$  is not too high. Of course, one can also use RQMC in the likelihood estimation, but we do not pursue this idea in this article.



#### 4. The Effect of Estimating the Likelihood

This section studies the effect of the variance of the noisy likelihood on the VBIL estimators, and provides guidelines for selecting the number of particles  $N$ . A large  $N$  gives a precise likelihood estimate and therefore an accurate estimate of  $\lambda$ , but at a greater computational cost. A small  $N$  leads to a large variance of the likelihood estimator, so a larger number of iterations is needed for the procedure to settle down. It is therefore useful in practice to have some guidelines for selecting  $N$ .

To understand the effect of estimating the likelihood, we follow Pitt et al. (2012) and assume that

**Assumption 1.** There is a function  $\gamma^2(\theta) > 0$  such that  $\mathbb{E}(z|\theta) = -\frac{\gamma^2(\theta)}{2N}$  and  $\mathbb{V}(z|\theta) = \frac{\gamma^2(\theta)}{N}$ .

More precisely, Pitt et al. (2012) assumed further that  $z \sim \mathcal{N}(-\frac{\gamma^2(\theta)}{2N}, \frac{\gamma^2(\theta)}{N})$  to derive a theory for selecting an optimal  $N$ . This assumption is justified in Tran et al. (2013) and Doucet et al. (2015) making use of the unbiasedness of the likelihood estimate. The reason that the mean of  $z$  is  $-\frac{1}{2}$  times its variance is because  $\mathbb{E}(e^z) = 1$  in order for the likelihood estimator to be unbiased.

**Assumption 2.** For a given  $\sigma^2 > 0$ , let  $N$  be a function of  $\theta$  and  $\sigma^2$  such that  $\mathbb{V}(z|\theta) = \sigma^2$ , that is,  $N = N_{\sigma^2}(\theta) = \gamma^2(\theta)/\sigma^2$ . Then  $\mathbb{E}(z|\theta) = -\frac{\sigma^2}{2}$  and  $\mathbb{V}(z|\theta) = \sigma^2$ .

Suppose that the equation  $\nabla_{\lambda} \text{KL}(\lambda) = 0$ , with  $\text{KL}(\lambda)$  defined in (1), has the unique solution  $\lambda^*$ . Let  $\hat{\lambda}_M$  be the estimator of  $\lambda^*$  obtained by Algorithm 1 after  $M$  iterations, and  $\tilde{\lambda}_M$  be the corresponding estimator obtained when the exact likelihood is available. Denote  $\zeta_*(\theta) = \nabla_{\lambda} [\log q_{\lambda}(\theta)]|_{\lambda=\lambda^*}$  and denote by  $\mathbb{E}_*(\cdot)$  and  $\mathbb{V}_*(\cdot)$  the expectation and variance operators with respect to  $q_{\lambda^*}(\theta)$ . For simplicity, we consider the case that  $\lambda$  is scalar; the case with a multivariate  $\lambda$  can be obtained using Theorem 5 of Sacks (1958). We obtain the following results whose proof is in the Appendix.

**Theorem 1.** Suppose that Assumptions 1 and 2 are satisfied, and that the regularity conditions in Theorem 1 of Sacks (1958) hold.

(i) Then,

$$\sqrt{M}(\hat{\lambda}_M - \lambda^*) \xrightarrow{d} \mathcal{N}\left(0, c_{\lambda^*} \mathbb{V}(\nabla_{\lambda} \text{KL}(\lambda^*))\right), \quad \text{as } M \rightarrow \infty, \quad (11)$$

where  $c_{\lambda^*}$  is a positive constant that depends only on geometric properties of the function  $\nabla_{\lambda} \text{KL}(\lambda^*)$  and is independent of the random variables involving in estimating  $\nabla_{\lambda} \text{KL}(\lambda^*)$ , that is,  $c_{\lambda^*}$  is independent of  $\sigma^2$ .

(ii) Let  $\sigma_{\text{asym}}^2(\hat{\lambda}_M) = c_{\lambda^*} \mathbb{V}(\nabla_{\lambda} \text{KL}(\lambda^*))$  be the asymptotic variance of  $\hat{\lambda}_M$  as  $M \rightarrow \infty$ . Similarly, let  $\sigma_{\text{asym}}^2(\tilde{\lambda}_M)$  be the asymptotic variance of  $\tilde{\lambda}_M$ . Then,

$$\sigma_{\text{asym}}^2(\hat{\lambda}_M) = \sigma_{\text{asym}}^2(\tilde{\lambda}_M) + \sigma^2 \tau(\lambda^*, S), \quad (12)$$

where  $\tau(\lambda^*, S) = c_{\lambda^*} \mathbb{V}_*\{\zeta_*(\theta)\}/S$  if the noisy traditional gradient is used, and  $\tau(\lambda^*, S) = c_{\lambda^*} I_F(\lambda^*)^{-1} \mathbb{V}_*\{\zeta_*(\theta)\} I_F(\lambda^*)^{-1}/S$  if the noisy natural gradient in (7) is used.

These results show that the variance of VBIL estimators increases linearly with  $\sigma^2$ . For PM and IS<sup>2</sup> estimators, Pitt et al. (2012) and Tran et al. (2013) showed that their variances increase exponentially with  $\sigma^2$ . This means that VBIL is useful in cases where only a rough estimate of the likelihood is available, or it is expensive to obtain an accurate estimate of the likelihood.

We now discuss the issue of selecting  $\sigma^2$ . We note that under Assumption 2,  $N$  is tuned depending on  $\theta$  as  $N = N_{\sigma^2}(\theta) = \gamma^2(\theta)/\sigma^2$ , so the time to compute the likelihood estimate  $\hat{p}_N(y|\theta)$  is proportional to  $1/\sigma^2$ . Then, Pitt et al. (2012) and Tran et al. (2013) showed that, for the PM and IS<sup>2</sup> methods, the optimal  $\sigma^2$  that gives an optimal trade-off between the CPU time and the variance of the estimators is 1. For VBIL, the computing time can be defined as

$$\text{CT}(\sigma^2) = \frac{\sigma_{\text{asym}}^2(\hat{\lambda}_M)}{\sigma^2} = \frac{\sigma_{\text{asym}}^2(\tilde{\lambda}_M)}{\sigma^2} + \tau(\lambda^*, S), \quad (13)$$

where neither  $\sigma_{\text{asym}}^2(\tilde{\lambda}_M)$  nor  $\tau(\lambda^*, S)$  depends on  $\sigma^2$ . These results suggest that  $\sigma^2$  should be set to a large value, as long as it is not too large for the stochastic search procedure in Algorithm 1 to converge.

#### 5. Applications

##### 5.1. Application to Panel Data Models

Generalized linear mixed models (GLMM; see, e.g., Fitzmaurice, Laird, and Ware 2011), also known as panel data models, use a vector of random effects  $\alpha_i$  to account for the dependence between the observations  $y_i = \{y_{ij}, j = 1, \dots, n_i\}$  measured on the same individual  $i$ . Given the random effects  $\alpha_i$ , the conditional density  $p(y_i|\theta, \alpha_i)$  belongs to an exponential family. The joint likelihood function of the model parameters  $\theta$  and the random effects  $\alpha = (\alpha_1, \dots, \alpha_n)$ ,  $p(y, \alpha|\theta) = \prod_{i=1}^n p(\alpha_i|\theta) p(y_i|\theta, \alpha_i)$ , is tractable. Typically in the VB literature the joint posterior  $p(\theta, \alpha|y) \propto p(\theta) p(y, \alpha|\theta)$  is approximated by a variational distribution of the form  $q(\theta) q(\alpha)$ , and then  $q(\theta)$  is used as an approximation to the marginal posterior  $p(\theta|y)$ . For example, Tan and Nott (2013) took this approach but use partially noncentered parameterizations to reduce dependence between parameter blocks. Ormerod and Wand (2012) considered frequentist estimation of  $\theta$ , but using VB methods to integrate out  $\alpha$ . As discussed in the introduction, factorization of the VB distribution generally ignores the posterior dependence between  $\theta$  and  $\alpha$ , which often leads to underestimating the variance in the posterior distribution of  $\theta$ . In the following, we refer to such a VB method as classical VB.

The likelihood,  $p(y|\theta) = \prod_{i=1}^n p(y_i|\theta)$ , with  $p(y_i|\theta) = \int p(y_i|\theta, \alpha_i) p(\alpha_i|\theta) d\alpha_i$ , is in most cases analytically intractable but can be easily estimated unbiasedly using importance sampling. Let  $h_i(\alpha_i|y, \theta)$  be an importance density for  $\alpha_i$ . The integral  $p(y_i|\theta)$  is estimated unbiasedly by

$$\begin{aligned} \hat{p}_{N_i}(y_i|\theta) &= \frac{1}{N_i} \sum_{j=1}^{N_i} w_i(\alpha_i^{(j)}, \theta), \quad w_i(\alpha_i^{(j)}, \theta) \\ &= \frac{p(y_i|\alpha_i^{(j)}, \theta) p(\alpha_i^{(j)}|\theta)}{h_i(\alpha_i^{(j)}|y, \theta)}, \quad \alpha_i^{(j)} \stackrel{\text{iid}}{\sim} h_i(\cdot|y, \theta). \end{aligned} \quad (14)$$

It is possible to use different  $N_i$  for each  $p(y_i|\theta)$ . Hence,  $\hat{p}_N(y|\theta) = \prod_{i=1}^n \hat{p}_{N_i}(y_i|\theta)$  is an unbiased estimator of the likelihood  $p(y|\theta)$ . The variance of  $z = \log \hat{p}_N(y|\theta) - \log p(y|\theta)$  is

$$\mathbb{V}(z|\theta) = \mathbb{V}(\log \hat{p}_N(y|\theta)) = \sum_{i=1}^n \mathbb{V}(\log \hat{p}_{N_i}(y_i|\theta)), \quad (15)$$

which can be estimated by  $\hat{\mathbb{V}}(z|\theta) = \sum_{i=1}^n \hat{\mathbb{V}}(\log \hat{p}_{N_i}(y_i|\theta))$  with

$$\hat{\mathbb{V}}(\log \hat{p}_{N_i}(y_i|\theta)) = \frac{\hat{\gamma}_i(\theta)}{N_i}, \quad \hat{\gamma}_i(\theta) = \frac{N_i \sum_{j=1}^{N_i} w_i(\alpha_i^{(j)}, \theta)^2}{\left(\sum_{j=1}^{N_i} w_i(\alpha_i^{(j)}, \theta)\right)^2} - 1. \quad (16)$$

Given a fixed  $\sigma^2$ , it is therefore straightforward to target  $\mathbb{V}(z|\theta) = \sigma^2$  by selecting  $N_i$  such that  $\hat{\mathbb{V}}(\log \hat{p}_{N_i}(y_i|\theta)) \approx \sigma^2/n$ .

*Six City data.* We now illustrate the VBIL algorithm using the Six City data in Fitzmaurice and Laird (1993). The data consist of binary responses  $y_{ij}$  which indicate the wheezing status (1 if wheezing, 0 if not wheezing) of the  $i$ th child at time-point  $j$ ,  $i = 1, \dots, 537$  and  $j = 1, \dots, 4$ . Covariates are the age of the child at time-point  $j$ , centered at 9 years, and the maternal smoking status (0 or 1). We consider the logistic regression model with a random intercept  $y_{ij}|\beta, \alpha_i \sim \text{Binomial}(1, p_{ij})$ , where  $\text{logit}(p_{ij}) = \beta_1 + \beta_2 \text{Age}_{ij} + \beta_3 \text{Smoke}_{ij} + \alpha_i$ ,  $\alpha_i \sim \mathcal{N}(0, \tau^2)$ . The model parameters are  $\theta = (\beta, \tau^2)$ . We use a normal prior  $\mathcal{N}(0, 50I_3)$  for  $\beta$  and a Gamma(1, 0.1) prior for  $\tau^2$ .

We use the variational distribution  $q_\lambda(\theta) = q(\beta)q(\tau^2)$ , where  $q(\beta)$  is a  $d = 3$ -variate normal  $\mathcal{N}(\mu, \Sigma)$  and  $q(\tau^2)$  is an inverse gamma distribution. We then run Algorithm 1, see the Appendix for details, with  $S = 1000$  samples to estimate the gradient. The likelihood is estimated as in (14) with the natural sampler  $h_i(\alpha_i|y, \theta) = p(\alpha_i|\theta)$ , which is the normal distribution  $\mathcal{N}(0, \tau^2)$  in this case. The  $\sigma^2$  in Section 4 is set to 4, which on average requires  $\tilde{N} = \sum N_i/n = 124$  particles. Using a larger  $\sigma^2$  leads to too small  $N_i$  that makes the estimate in (16) unreliable. Figure 1(a) plots the scaled lower bounds over the iterations.

We compare the performance of the classical VB and VBIL algorithms to the PM approach. For the PM approach, we set  $\sigma^2 = 1$  as suggested in Pitt et al. (2012). The MCMC chain,

based on the adaptive random walk Metropolis–Hastings algorithm in Haario, Saksman, and Tamminen (2001), consists of 20,000 iterates with another 10,000 iterates used as burn-in.

Figure 2 plots the estimated marginal posteriors. The MCMC density estimates are carried out using the kernel density estimation method based on the built-in Matlab function `ksdensity`. The figure shows that the VBIL estimates are very close to the MCMC estimates. The classical VB underestimates the posterior variance of  $\tau^2$  in this example. The clock times taken to run the VB, VBIL, and PM procedures are 4, 2.9, and 505 min, respectively. However, we note that the running time depends on many factors such as the programming language being used and the initialization of the procedures. We note that PM could be sped up greatly using the blocking PM as in Tran et al. (2016).

All the examples in this article are run on an Intel Core i7 3.2GHz desktop supported by the Matlab Parallel Toolbox with eight local processors. Obviously, the more processors are used, the faster the VBIL procedure will be.

*Simulation data.* One of the main advantages of VBIL is its scalability, that is, it is applicable in large data cases. Consider large data cases with a large number of panels  $n$ . From (15), for fixed  $N_i$ , the variance of the log-likelihood estimator  $\mathbb{V}(z|\theta)$  increases linearly with  $n$ . Therefore, when  $n$  is large enough, the PM and  $\text{IS}^2$  methods will not work in a practical sense, because  $\mathbb{V}(z|\theta)$  can be very large (Flury and Shephard 2011). In this GLMM setting, PM and  $\text{IS}^2$  do not work when  $\mathbb{V}(z|\theta)$  is as large as 6 or 7. One can decrease  $\mathbb{V}(z|\theta)$  by increasing  $N_i$ , but this can be too computationally expensive to be practical.

An alternative approach for sampling from  $p(\theta|y)$  is to sample from the joint posterior  $p(\theta, \alpha|y)$ . To study the accuracy of VBIL and its scalability, we compare VBIL against this approach. We use the RStan package `rstanarm` to sample from  $p(\theta, \alpha|y)$ . Stan (Stan Development Team 2016) is a programming language for inference and posterior analysis. Stan is known for its high efficiency because it uses C++ as the core and implements Hamiltonian MC for sampling. The package `rstanarm` is Stan's R interface (RStan), that performs posterior analysis for models with dependent data such as GLMMs.

We generate datasets from the logistic model with a random intercept as above, but with  $\text{logit}(p_{ij}) = \beta_1 + \beta_2 x_{ij} + \alpha_i$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, n_i = 5$ , where  $\beta = (-1.5, 2.5)'$ ,  $\tau^2 = 1.5$ ,  $x_{ij} \sim U(0, 1)$ . Table 1 summarizes the results for

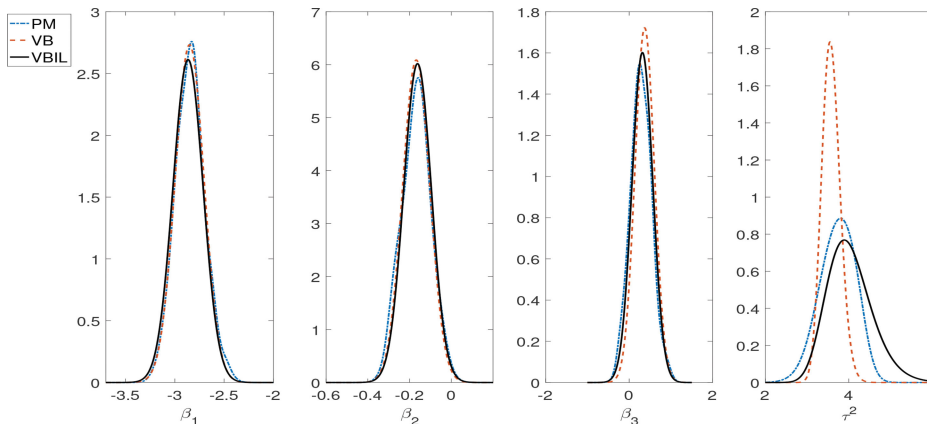


Figure 2. Six City data: plots of the classical VB estimates (dashed line), PM estimates (dotted line), and VBIL estimates (solid line) of the marginal posteriors  $p(\beta_i|y)$  and  $p(\tau^2|y)$ .

**Table 1.** Simulation data: The table reports (i) the posterior mean and the posterior standard deviation (in brackets) for each parameter, (ii) the clock time (in minutes) for each method. For each dataset, two Markov chains, one with 5000 and the other with 10,000 iterations, are produced by RStan.

$n$	Methods	$\beta_1$	$\beta_2$	$\tau^2$	Clock time
500	RStan (5000)	-1.351(0.114)	2.446(0.182)	1.331(0.195)	2.4
	RStan (10,000)	-1.348(0.115)	2.446(0.184)	1.329(0.200)	4.4
	VBIL	-1.389(0.120)	2.513(0.187)	1.474(0.200)	1.4
1000	RStan (5000)	-1.473(0.085)	2.477(0.132)	1.417(0.147)	5.5
	RStan (10,000)	-1.472(0.085)	2.477(0.132)	1.416(0.146)	9.5
	VBIL	-1.507(0.087)	2.511(0.130)	1.551(0.144)	2.1
3000	VBIL	-1.509(0.062)	2.561(0.096)	1.506(0.102)	8.1
5000	VBIL	-1.509(0.048)	2.501(0.073)	1.506(0.092)	14.2

various  $n$ . The results show that VBIL enjoys a high accuracy compared to the “gold-standard” MCMC estimates. It is important to note that, unlike classical VB that ignores the posterior dependence of  $\theta$  and  $\alpha$ , VBIL does not underestimate posterior variances.

For the case with  $n = 3000$ , it takes, on average across different  $\theta$ , 30 sec to carry out each likelihood estimation with the numbers of particles  $N_i$  tuned to target  $\mathbb{V}(z|\theta) = 1$ . So if an optimal PM procedure was run on our computer to generate a chain of 10,000 iterations, it would take 3.5 days. Because the RStan method generates and stores a Markov chain of size  $(n + 3) \times \ell$  with  $\ell$  the number of iterations, this method becomes computationally infeasible for the case  $n = 3000$  and  $n = 5000$ .

## 5.2. Application to State-Space Models

In state-space models, the observations  $y_t$  are observed in time order. At time  $t$ , the distribution of  $y_t$  conditional on a state variable  $x_t$  is independently distributed as  $y_t|x_t \sim g_t(y_t|x_t, \theta)$ , and the state variables  $\{x_t\}_{t \geq 1}$  are a Markov chain with  $x_1 \sim \mu_\theta(\cdot)$  and  $x_t|x_{t-1} \sim f_t(x_t|x_{t-1}, \theta)$ . The likelihood of the data  $y = y_{1:T}$  is  $p(y|\theta) = \int p(y|x, \theta)p(x|\theta)dx$  with  $x = x_{1:T}$  and

$$p(x|\theta) = \mu_\theta(x_1) \prod_{t=2}^T f_t(x_t|x_{t-1}, \theta), \quad p(y|x, \theta) = \prod_{t=1}^T g_t(y_t|x_t, \theta).$$

Given a value of  $\theta$ , the likelihood  $p(y|\theta)$  can be unbiasedly estimated by an importance sampling estimator (Shephard and Pitt 1997; Durbin and Koopman 1997) or by a particle filter estimator (Pitt et al. 2012),  $\hat{p}_N(y|\theta)$ , with  $N$  the number of particles.

An important example of state-space models is the stochastic volatility (SV) model. The time series data  $y_t = \exp(x_t/2)w_t$ ,

where  $w_t \sim \mathcal{N}(0, 1)$  and

$$x_t = \mu + \phi(x_{t-1} - \mu) + \sigma v_t, \quad x_1 \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{1 - \phi^2}\right), \\ v_t \sim \mathcal{N}(0, 1),$$

with  $\mu \in \mathbb{R}$ ,  $\phi \in (-1, 1)$  and  $\sigma^2 > 0$ . Let  $\tau = (1 + \phi)/2 \in (0, 1)$ ; we will estimate  $\tau$  but report results for  $\phi$ . The model parameters are  $\theta = (\mu, \tau, \sigma^2)$ . We follow Kim, Shephard, and Chib (1998) and use a normal prior  $\mathcal{N}(0, 10)$  for  $\mu$ , a Beta prior  $B(20, 1.5)$  for  $\tau$ , and an inverse gamma  $IG(2.5, 0.025)$  for  $\sigma^2$ .

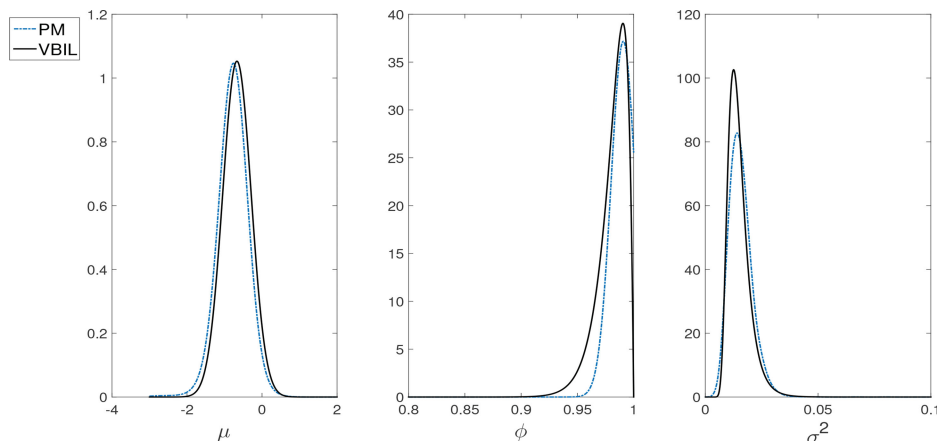
To illustrate the VBIL algorithm for state-space models, we analyze the weekday close exchange rates  $r_t$  for the Australian Dollar/U.S. Dollar from 5/1/2010 to 31/12/2013. The data are available from the Reserve Bank of Australia. The data  $y_t$  are

$$y_t = 100 \left( \log \frac{r_{t+1}}{r_t} - \frac{1}{T} \sum_{i=1}^T \log \frac{r_{i+1}}{r_i} \right), \quad t = 1, \dots, T = 1001.$$

We use the variational distribution  $q_\lambda(\theta) = q(\mu)q(\tau)q(\sigma^2)$ , where  $q(\mu)$  is  $\mathcal{N}(\mu_\mu, \sigma_\mu^2)$ ,  $q(\tau)$  is  $\text{Beta}(\alpha_\tau, \beta_\tau)$ , and  $q(\sigma^2)$  is inverse gamma  $IG(\alpha_{\sigma^2}, \beta_{\sigma^2})$ . We employ the constraint  $\alpha_\tau > 1$  and  $\beta_\tau > 1$  to make sure that  $q(\tau)$  has a mode. The likelihood estimator  $\hat{p}_N(y|\theta)$  is computed by the standard particle filter. We then run the VBIL algorithm with  $S = 1000$  samples, starting with  $\mu_\mu = 0$ ,  $\sigma_\mu^2 = 0.3$ ,  $\alpha_\tau = 95$ ,  $\beta_\tau = 5$ ,  $\alpha_{\sigma^2} = 11$ ,  $\beta_{\sigma^2} = 1$ . This initial point is set so that the initial mean values of  $\mu$ ,  $\phi$ , and  $\sigma^2$  are 0, 0.9, and 0.1, respectively, which is pretty far away from the posterior means; see Figure 3. The VBIL algorithm stops after 28 iterations. Figure 1(b) plots the scaled lower bounds over the iterations.

The VBIL is compared to pseudo-marginal simulation, based on an adaptive random walk Metropolis–Hastings algorithm, with 100,000 iterations starting from the same values  $\mu = 0$ ,  $\tau = 0.95$ , and  $\sigma^2 = 0.1$ . The number of particles used in PM is fixed at  $N = 300$ , so that  $\mathbb{V}(\hat{p}_N(y|\bar{\theta})) \approx 1$  at the initial value  $\bar{\theta} = (0, 0.95, 0.1)$ . The number of particles used in VBIL is fixed at  $N = 100$  as the use of randomized QMC for generating  $\theta$  helps reduce greatly the variance in estimating the gradient. We fix  $N$  in this example as it is difficult to estimate the variance of log-likelihood estimates obtained by the particle filter.

Figure 3 plots the PM estimates (dotted line) and the VBIL estimates (solid line) of the marginal posteriors. The figure shows that the VBIL estimates are close to the PM estimates but



**Figure 3.** Exchange rate data: plots of the PM (dashed line) and VBIL (solid line) estimates of the marginal posteriors. The VBIL estimates are similar to the PM estimates and do not underestimate the posterior variance.



require significantly less computational time. The CPU times taken to run the VBIL and PM procedures are 0.7 and 28 min, respectively.

### 5.3. Application to ABC

ABC approximates the intractable likelihood by  $p_{\text{LF}}(y|\theta) = \int K_\epsilon(S(y'), S(y))p(y'|\theta)dy'$ , where  $K_\epsilon(\cdot, \cdot)$  is a kernel with the bandwidth  $\epsilon$  and  $S(\cdot)$  is a vector of summary statistics. Inference is then based on the approximate posterior  $p_{\text{ABC}}(\theta|y) \propto p(\theta)p_{\text{LF}}(y|\theta)$ . Because the likelihood-free function  $p_{\text{LF}}(y|\theta)$  can be unbiasedly estimated by

$$\hat{p}_N^{\text{LF}}(y|\theta) = \frac{1}{N} \sum_{i=1}^N K_\epsilon(S(y^{[i]}), S(y)), \quad y^{[i]} \stackrel{\text{iid}}{\sim} p(\cdot|\theta),$$

it is straightforward to use the VBIL algorithm to approximate  $p_{\text{ABC}}(\theta|y)$ .

We illustrate the application of the VBIL algorithm to ABC by using it to fit an  $\alpha$ -stable distribution.  $\alpha$ -stable distributions (Nolan 2007) are a class of heavy-tailed distributions used in many statistical applications. An  $\alpha$ -stable distribution  $\mathcal{S}(\alpha, \beta, \gamma, \delta)$  is parameterized by the stability parameter  $\alpha \in (0, 2)$ , skewness  $\beta \in (-1, 1)$ , scale  $\gamma > 0$ , and location  $\delta \in \mathbb{R}$ . The main difficulty when working with  $\alpha$ -stable distributions is that they do not have closed-form densities, which makes it difficult to do inference. However, ABC can be used for Bayesian inference as it is easy to sample from an  $\alpha$ -stable distribution (Peters, Sisson, and Fan 2012). We illustrate in this example that VBIL provides an efficient approach for fitting  $\alpha$ -stable distributions.

We generate a dataset  $y$  with  $n = 500$  observations from a univariate  $\alpha$ -stable distribution  $\mathcal{S}(1.5, 0.5, 1, 0)$ . We use the summary statistics as in Peters, Sisson, and Fan (2012); the interested reader is referred to Peters, Sisson, and Fan (2012) for the details. As the parameterization is discontinuous at  $\alpha = 1$ , resulting in poor estimates of the summary statistics, we consider the case with  $\alpha > 1$  and restrict the support of  $\alpha$  to the interval  $(1.1, 2)$  as in Peters, Sisson, and Fan (2012).

We reparameterize  $\tilde{\alpha} = \log(\frac{\alpha-1.1}{2-\alpha})$ ,  $\tilde{\beta} = \log(\frac{\beta+1}{1-\beta}) \in$ ,  $\tilde{\gamma} = \log(\gamma)$ ,  $\tilde{\delta} = \delta$ , and estimate  $\tilde{\theta} = (\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\delta})$ , but report the results for  $(\alpha, \beta, \gamma, \delta)$ . We use a normal prior  $\tilde{\theta} \sim \mathcal{N}(0, 100I_4)$  and approximate the posterior  $p(\tilde{\theta}|y)$  by a normal variational distribution  $q_\lambda(\tilde{\theta}) = \mathcal{N}(\mu_{\tilde{\theta}}, \Sigma_{\tilde{\theta}})$ . One can work with the original parameterization  $(\alpha, \beta, \gamma, \delta)$  and use some form of factorization  $q(\alpha)q(\beta)q(\gamma)q(\delta)$ . We choose to work with  $\tilde{\theta}$  to account for the posterior dependence between the parameters. This also illustrates the flexibility of the VBIL method in the sense that it can be used without requiring factorization.

We use the Gaussian kernel with covariance matrix  $0.01I_4$  for the kernel  $K_\epsilon$ . VBIL is compared to PM with 20,000 iterations after 5000 burn-in iterations. For the standard PM (Andrieu and Roberts 2009), the number of pseudo-datasets  $N = 20$ , which was set after some experimentation to obtain a well mixing chain. Efficient versions of PM have been proposed recently, which are more tolerant of noise in the likelihood estimates. Here, we compare VBIL to the blocking PM method of Tran et al. (2016). For the blocking PM, we set  $N = 5$ . We also use this value of  $N$  in VBIL. Table 2 shows the VBIL and MCMC

**Table 2.** ABC example: Standard PM, block PM, and VBIL estimates of the posterior means and posterior standard deviations. The VBIL method is roughly 104 and 63 times faster than the standard and block PM approach, respectively.

	True	Standard PM	Blocking PM	VBIL
$\alpha$	1.5	1.57 (0.15)	1.58 (0.14)	1.57 (0.11)
$\beta$	0.5	0.46 (0.21)	0.45 (0.21)	0.48 (0.16)
$\gamma$	1	1.04 (0.12)	1.04 (0.12)	1.02 (0.12)
$\delta$	0	-0.08 (0.21)	-0.09 (0.18)	-0.08 (0.14)
CPU time (min)		12.56	7.62	0.12

estimates, and the CPU times. As shown, VBIL is orders of magnitude faster than MCMC in this example. Figure 1(c) plots the scaled lower bounds over the iterations.

### 5.4. Using VBIL to Improve Estimates of the Marginal Posteriors

A drawback of VB methods in general is that the factorization assumption as in (9) ignores the posterior dependence between the factors, which might lead to poor approximations of the posterior variances (Neville, Ormerod, and Wand 2014). We now show how the VBIL algorithm can be used to help overcome this problem.

Suppose that we would like to have a highly accurate VB approximation to the marginal posterior  $p(\theta^{(j)}|y)$ . We restrict ourselves to the case with a tractable likelihood for simplicity, but the following discussion also applies when the likelihood is intractable. The likelihood of  $\theta^{(j)}$ ,

$$p(y|\theta^{(j)}) = \int p(\theta^{(\setminus j)}|\theta^{(j)}) p(y|\theta^{(1)}, \dots, \theta^{(K)}) d\theta^{(\setminus j)}, \quad (17)$$

with  $\theta^{(\setminus j)} = (\theta^{(1)}, \dots, \theta^{(j-1)}, \theta^{(j+1)}, \dots, \theta^{(K)})$ , is in general intractable but can be estimated unbiasedly. Let  $q(\theta^{(\setminus j)})$  be an approximation to the marginal posterior  $p(\theta^{(\setminus j)}|y)$  resulting from a classical VB method that uses the factorization (9). The integral in (17) can be estimated unbiasedly using importance sampling with the proposal density  $q(\theta^{(\setminus j)})$  or a tail-flattened version of it. This is accurate enough in practice because VBIL does not require a very accurate estimate of  $p(y|\theta^{(j)})$  as discussed in Section 4. The VBIL algorithm can then be used to approximate the marginal posterior  $p(\theta^{(j)}|y)$  directly with  $\theta^{(\setminus j)}$  integrated out. The resulting approximation is often highly accurate as the dependence between  $\theta^{(j)}$  and  $\theta^{(\setminus j)}$  is taken into account.

A formal justification is as follows. We use the notation as in (9) and write  $\lambda = (\lambda^{(j)}, \lambda^{(\setminus j)})$ . Suppose that we estimate the marginal posterior of  $\lambda^{(j)}$  by  $q_{\lambda^{(\setminus j)}}(\theta^{(j)})$ , which belongs to a family  $\mathcal{F} = \{q_{\lambda^{(\setminus j)}}(\theta^{(j)}), \lambda^{(\setminus j)} \in \Lambda\}$ . VBIL proceeds by minimizing

$$\text{KL}_j(\lambda^{(j)}) = \int q_{\lambda^{(\setminus j)}}(\theta^{(j)}) \log \frac{q_{\lambda^{(\setminus j)}}(\theta^{(j)})}{p(\theta^{(j)}|y)} d\theta^{(j)}$$

over  $\lambda^{(j)} \in \Lambda$ . Let  $\lambda_*^{(j)}$  be the VBIL estimator. Under Assumptions 1 and 2, or when the number of samples  $N$  used to estimate (17) is large enough,  $\lambda_*^{(j)}$  is guaranteed to be a minimizer of  $\text{KL}_j(\lambda^{(j)})$ . If  $\text{KL}_j(\lambda^{(j)})$  is also convex, then

$$\text{KL}_j(\lambda_*^{(j)}) \leq \text{KL}_j(\lambda^{(j)}) \quad \text{for all } \lambda^{(j)} \in \Lambda. \quad (18)$$

If we use a VB procedure with a factorization of the form  $q_{\lambda}(\theta) = q_{\lambda^{(j)}}(\theta^{(j)})q_{\lambda^{(\setminus j)}}(\theta^{(\setminus j)})$  where  $q_{\lambda^{(j)}}(\theta^{(j)})$  belongs to the same family  $\mathcal{F}$ , then VB proceeds by minimizing the KL divergence  $\text{KL}(\lambda^{(j)}, \lambda^{(\setminus j)})$

$$\begin{aligned} & \int q_{\lambda^{(j)}}(\theta^{(j)}) q_{\lambda^{(\setminus j)}}(\theta^{(\setminus j)}) \log \frac{q_{\lambda^{(j)}}(\theta^{(j)}) q_{\lambda^{(\setminus j)}}(\theta^{(\setminus j)})}{p(\theta^{(j)}, \theta^{(\setminus j)} | y)} d\theta^{(j)} d\theta^{(\setminus j)} \\ &= \text{KL}_j(\lambda^{(j)}) + \int q_{\lambda^{(j)}}(\theta^{(j)}) \int q_{\lambda^{(\setminus j)}}(\theta^{(\setminus j)}) \\ & \quad \times \log \frac{q_{\lambda^{(\setminus j)}}(\theta^{(\setminus j)})}{p(\theta^{(\setminus j)} | \theta^{(j)}, y)} d\theta^{(\setminus j)} d\theta^{(j)}. \end{aligned} \quad (19)$$

Let  $(\tilde{\lambda}^{(j)}, \tilde{\lambda}^{(\setminus j)})$  be a minimizer of (19). Because of the decomposition in (19), the estimator  $\tilde{\lambda}^{(j)}$  is not necessarily the minimizer of  $\text{KL}_j(\lambda^{(j)})$ . From (18),

$$\text{KL}_j(\lambda_*^{(j)}) \leq \text{KL}_j(\tilde{\lambda}^{(j)}). \quad (20)$$

So the VBIL estimator  $\lambda_*^{(j)}$  is no worse than the factorization-based VB estimator  $\tilde{\lambda}^{(j)}$  in terms of KL divergence. An example to illustrate this can be found in the Appendix.

## 6. Conclusion

VBIL is a useful VB algorithm for Bayesian inference in statistical modeling where the likelihood is intractable. The method makes it possible to do inference in statistical models using VB in some situations that were previously impossible. The main advantage of VBIL over its competitors, such as PM and IS<sup>2</sup>, is its scalability. We show in the examples that VBIL is several orders of magnitude faster than these existing methods.

## Supplementary Materials

The online Appendix contains the proof of Theorem 1, a detailed derivation for Section 5.1, and an example to show the importance of the natural gradient.

## Acknowledgments

The research of Tran and Kohn was partially supported by the ARC COE grant CE140100049. Nott's research was supported by a Singapore Ministry of Education Academic Research Fund Tier 2 grant (R-155-000-143-112).

## References

- Amari, S. (1998), "Natural Gradient Works Efficiently in Learning," *Neural Computation*, 10, 251–276. [874,876]
- Andrieu, C., and Roberts, G. (2009), "The Pseudo-Marginal Approach for Efficient Monte Carlo Computations," *The Annals of Statistics*, 37, 697–725. [874,880]
- Beaumont, M. A. (2003), "Estimation of Population Growth or Decline in Genetically Monitored Populations," *Genetics*, 164, 1139–1160. [874]
- Bishop, C. M. (2006), *Pattern Recognition and Machine Learning*, New York: Springer. [873]
- Deligiannidis, G., Doucet, A., and Pitt, M. (2016), "The Correlated Pseudo-Marginal Method," Technical Report. arXiv:1511.04992v3. [874]
- Dick, J., and Pillichshammer, F. (2010), *Digital Nets and Sequences. Discrepancy Theory and Quasi-Monte Carlo Integration*, Cambridge: Cambridge University Press. [876]

- Doucet, A., Pitt, M. K., Deligiannidis, G., and Kohn, R. (2015), "Efficient Implementation of Markov Chain Monte Carlo When Using an Unbiased Likelihood Estimator," *Biometrika*, 102, 295–313. [877]
- Durbin, J., and Koopman, S. J. (1997), "Monte Carlo Maximum Likelihood Estimation for Non-Gaussian State Space Models," *Biometrika*, 84, 669–684. [879]
- (2001), *Time Series Analysis by State Space Methods*, Oxford: Oxford University Press. [873]
- Fitzmaurice, G. M., and Laird, N. M. (1993), "A Likelihood-Based Method for Analysing Longitudinal Binary Responses," *Biometrika*, 80, 141–151. [878]
- Fitzmaurice, G. M., Laird, N. M., and Ware, J. H. (2011), *Applied Longitudinal Analysis* (2nd ed.), Hoboken, NJ: Wiley. [877]
- Flury, T., and Shephard, N. (2011), "Bayesian Inference Based Only on Simulated Likelihood: Particle Filter Analysis of Dynamic Economic Models," *Econometric Theory*, 1, 1–24. [878]
- Ghahramani, Z., and Hinton, G. E. (2000), "Variational Learning for Switching State-Space Models," *Neural Computation*, 12, 831–864. [873]
- Haario, H., Saksman, E., and Tamminen, J. (2001), "An Adaptive Metropolis Algorithm," *Bernoulli*, 7, 223–242. [878]
- Hoffman, M. D., Blei, D. M., Wang, C., and Paisley, J. (2013), "Stochastic Variational Inference," *Journal of Machine Learning Research*, 14, 1303–1347. [876]
- Honkela, A., Raiko, T., Kuusela, M., Tornio, M., and Karhunen, J. (2010), "Approximate Riemannian Conjugate Gradient Learning for Fixed-Form Variational Bayes," *Journal of Machine Learning Research*, 11, 3235–3268. [876]
- Kim, S., Shephard, N., and Chib, S. (1998), "Stochastic Volatility: Likelihood Inference and Comparison With Arch Models," *The Review of Economic Studies*, 65, 361–393. [879]
- Neville, S. E., Ormerod, J. T., and Wand, M. P. (2014), "Mean Field Variational Bayes for Continuous Sparse Signal Shrinkage: Pitfalls and Remedies," *Electronic Journal of Statistics*, 8, 1113–1151. [873,880]
- Nolan, J. (2007), *Stable Distributions: Models for Heavy-Tailed Data*, Boston: Birkhauser. [880]
- Nott, D. J., Tan, S., Villani, M., and Kohn, R. (2012), "Regression Density Estimation With Variational Methods and Stochastic Approximation," *Journal of Computational and Graphical Statistics*, 21, 797–820. [875]
- Ormerod, J. T., and Wand, M. P. (2012), "Gaussian Variational Approximate Inference for Generalized Linear Mixed Models," *Journal of Computational and Graphical Statistics*, 21, 2–17. [877]
- Owen, A. (1997), "Monte Carlo Variance of Scrambled Net Quadrature," *SIAM Journal on Numerical Analysis*, 34, 1884–1910. [876]
- Paisley, J., Blei, D., and Jordan, M. (2012), "Variational Bayesian Inference With Stochastic Search," in *International Conference on Machine Learning*, Edinburgh, Scotland, UK, pp. 1367–1374. [875]
- Peters, G., Sisson, S., and Fan, Y. (2012), "Likelihood-Free Bayesian Inference for  $\alpha$ -Stable Models," *Computational Statistics & Data Analysis*, 56, 3743–3756. [873,880]
- Pitt, M. K., Silva, R. S., Giordani, P., and Kohn, R. (2012), "On Some Properties of Markov Chain Monte Carlo Simulation Methods Based on the Particle Filter," *Journal of Econometrics*, 171, 134–151. [874,877,878,879]
- Ranganath, R., Gerrish, S., and Blei, D. M. (2014), "Black Box Variational Inference," in *International Conference on Artificial Intelligence and Statistics* (Vol. 33), Iceland: Reykjavik, pp. 814–822. [875,876]
- Sacks, J. (1958), "Asymptotic Distribution of Stochastic Approximation Procedures," *The Annals of Mathematical Statistics*, 29, 373–405. [877]
- Salimans, T., and Knowles, D. A. (2013), "Fixed-Form Variational Posterior Approximation Through Stochastic Linear Regression," *Bayesian Analysis*, 8, 741–908. [876]
- Shephard, N., and Pitt, M. K. (1997), "Likelihood Analysis of Non-Gaussian Measurement Time Series," *Biometrika*, 84, 653–667. [879]
- Stan Development Team (2016), "The Stan C++ Library," Version 2.14.0. [878]

- Tan, L. S. L., and Nott, D. J. (2013), “Variational Inference for Generalized Linear Mixed Models Using Partially Noncentered Parametrizations,” *Statistical Science*, 28, 168–188. [877]
- Tavare, S., Balding, D. J., Griffiths, R. C., and Donnelly, P. (1997), “Inferring Coalescence Times From DNA Sequence Data,” *Genetics*, 145, 505–518. [873]
- Tran, M.-N., Kohn, R., Quiroz, M., and Villani, M. (2016), “Block-Wise Pseudo-Marginal Metropolis-Hastings,” Technical Report. arXiv:1603.02485. [874,878,880]
- Tran, M.-N., Scharth, M., Pitt, M. K., and Kohn, R. (2013), “Importance Sampling Squared for Bayesian Inference in Latent Variable Models,” available at <http://arxiv.org/abs/1309.3339>. [874,877]
- Turner, R. E., and Sahani, M. (2011), “Two Problems With Variational Expectation Maximisation for Time-Series Models,” in *Bayesian Time Series Models*, eds. D. Barber, A. T. Cemgil, and S. Chiappa, Cambridge: Cambridge University Press, pp. 109–130. [873]