

Covariances, Robustness, and Variational Bayes

Ryan Giordano

*Department of Statistics, UC Berkeley
367 Evans Hall, UC Berkeley
Berkeley, CA 94720*

RGIORDANO@BERKELEY.EDU

Tamara Broderick

*Department of EECS, MIT
77 Massachusetts Ave., 38-401
Cambridge, MA 02139*

TBRODERICK@CSAIL.MIT.EDU

Michael I. Jordan

*Department of Statistics and EECS, UC Berkeley
367 Evans Hall, UC Berkeley
Berkeley, CA 94720*

JORDAN@CS.BERKELEY.EDU

Editor: Mohammad Emtiyaz Khan

Abstract

Mean-field Variational Bayes (MFVB) is an approximate Bayesian posterior inference technique that is increasingly popular due to its fast runtimes on large-scale data sets. However, even when MFVB provides accurate posterior means for certain parameters, it often mis-estimates variances and covariances. Furthermore, prior robustness measures have remained undeveloped for MFVB. By deriving a simple formula for the effect of infinitesimal model perturbations on MFVB posterior means, we provide both improved covariance estimates and local robustness measures for MFVB, thus greatly expanding the practical usefulness of MFVB posterior approximations. The estimates for MFVB posterior covariances rely on a result from the classical Bayesian robustness literature that relates derivatives of posterior expectations to posterior covariances and includes the Laplace approximation as a special case. Our key condition is that the MFVB approximation provides good estimates of a select subset of posterior means—an assumption that has been shown to hold in many practical settings. In our experiments, we demonstrate that our methods are simple, general, and fast, providing accurate posterior uncertainty estimates and robustness measures with runtimes that can be an order of magnitude faster than MCMC.

Keywords: Variational Bayes; Bayesian robustness; Mean field approximation; Linear response theory; Laplace approximation; Automatic differentiation

1. Introduction

Most Bayesian posteriors cannot be calculated analytically, so in practice we turn to approximations. Variational Bayes (VB) casts posterior approximation as an optimization problem in which the objective to be minimized is the divergence, among a sub-class of tractable distributions, from the exact posterior. For example, one widely-used and relatively simple flavor of VB is “mean field variational Bayes” (MFVB), which employs Kullback-Leibler (KL) divergence and a factorizing exponential family approximation for the tractable sub-class of posteriors (Wainwright and Jordan, 2008). MFVB has been increasingly popular as an alternative to Markov Chain Monte Carlo (MCMC) in part due to its fast runtimes on large-scale data sets. Although MFVB does not come with any general accuracy guarantees (except asymptotic ones in special cases (Westling and McCormick, 2015; Wang and Blei, 2017)), MFVB produces posterior mean estimates of certain parameters that are accurate enough to be useful in a number of real-world applications (Blei et al., 2016). Despite this ability to produce useful point estimates for large-scale data sets, MFVB is limited as an inferential tool; in particular, MFVB typically underestimates marginal variances (MacKay, 2003; Wang

and Titterton, 2004; Turner and Sahani, 2011). Moreover, to the best of our knowledge, techniques for assessing Bayesian robustness have not yet been developed for MFVB. It is these inferential issues that are the focus of the current paper.

Unlike the optimization approach of VB, an MCMC posterior estimate is an empirical distribution formed with posterior draws. MCMC draws lend themselves naturally to the approximate calculation of posterior moments, such as those required for covariances. In contrast, VB approximations lend themselves naturally to sensitivity analysis, since we can analytically differentiate the optima with respect to perturbations. However, as has long been known in the Bayesian robustness literature, the contrast between derivatives and moments is not so stark since, under mild regularity conditions that allow the exchange of integration and differentiation, there is a direct correspondence between derivatives and covariance (Gustafson, 1996b; Basu et al., 1996; Efron, 2015, Section 2.2 below).

Thus, in order to calculate local sensitivity to model hyperparameters, the Bayesian robustness literature re-casts derivatives with respect to hyperparameters as posterior covariances that can be calculated with MCMC. In order to provide covariance estimates for MFVB, we turn this idea on its head and use the sensitivity of MFVB posterior expectations to estimate their covariances. These sensitivity-based covariance estimates are referred to as “linear response” estimates in the statistical mechanics literature (Opper and Saad, 2001), so we refer to them here as *linear response variational Bayes* (LRVB) covariances. Additionally, we derive straightforward MFVB versions of hyperparameter sensitivity measures from the Bayesian robustness literature. Under the assumption that the posterior means of interest are well-estimated by MFVB for all the perturbations of interest, we establish that LRVB provides a good estimate of local sensitivities. In our experiments, we compare LRVB estimates to MCMC, MFVB, and Laplace posterior approximations. We find that the LRVB covariances, unlike the MFVB and Laplace approximations, match the MCMC approximations closely while still being computed over an order of magnitude more quickly than MCMC.

In Section 2 we first discuss the general relationship between Bayesian sensitivity and posterior covariance and then define local robustness and sensitivity. Next, in Section 3, we introduce VB and derive the linear system for the MFVB local sensitivity estimates. In Section 4, we show how to use the MFVB local sensitivity results to estimate covariances and calculate canonical Bayesian hyperparameter sensitivity measures. Finally, in Section 5, we demonstrate the speed and effectiveness of our methods with simple simulated data, an application of automatic differentiation variational inference (ADVI), and a large-scale industry data set.

2. Bayesian Covariances and Sensitivity

2.1 Local Sensitivity and Robustness

Denote an unknown model parameter by the vector $\theta \in \mathbb{R}^K$, assume a dominating measure for θ on \mathbb{R}^K given by λ , and denote observed data by x . Suppose that we have a vector-valued hyperparameter $\alpha \in \mathcal{A} \subseteq \mathbb{R}^D$ that parameterizes some aspects of our model. For example, α might represent prior parameters, in which case we would write the prior density with respect to λ as $p(\theta|\alpha)$, or it might parameterize a class of likelihoods, in which case we could write the likelihood as $p(x|\theta, \alpha)$. Without loss of generality, we will include α in the definition of both the prior and likelihood. For the moment, let $p_\alpha(\theta)$ denote the posterior density of θ given x and α , as given by Bayes’ Theorem (this definition of $p_\alpha(\theta)$ will be a special case of the more general Definition 2 below):

$$p_\alpha(\theta) := p(\theta|x, \alpha) = \frac{p(x|\theta, \alpha) p(\theta|\alpha)}{\int p(x|\theta', \alpha) p(\theta'|\alpha) \lambda(d\theta')} = \frac{p(x|\theta, \alpha) p(\theta|\alpha)}{p(x|\alpha)}.$$

We will assume that we are interested in a posterior expectation of some function $g(\theta)$ (e.g., a parameter mean, a posterior predictive value, or squared loss): $\mathbb{E}_{p_\alpha}[g(\theta)]$. In the current work, we will quantify the uncertainty of $g(\theta)$ by the posterior variance, $\text{Var}_{p_\alpha}(g(\theta))$. Other measures of central tendency (e.g., posterior medians) or uncertainty (e.g., posterior quantiles) may also be good choices but are beyond the scope of the current work.

Note the dependence of $\mathbb{E}_{p_\alpha} [g(\theta)]$ on both the likelihood and prior, and hence on α , through Bayes' Theorem. The choice of a prior and choice of a likelihood are made by the modeler and are almost invariably a simplified representation of the real world. The choices are therefore to some extent subjective, and so one hopes that the salient aspects of the posterior would not vary under reasonable variation in either choice. Consider the prior, for example. The process of prior elicitation may be prohibitively time-consuming; two practitioners may have irreconcilable subjective prior beliefs, or the model may be so complex and high-dimensional that humans cannot reasonably express their prior beliefs as formal distributions. All of these circumstances might give rise to a range of reasonable prior choices. A posterior quantity is "robust" to the prior to the extent that it does not change much when calculated under these different prior choices.

Quantifying the sensitivity of the posterior to variation in the likelihood and prior is one of the central concerns of the field of robust Bayes (Berger et al., 2000). (We will not discuss the other central concern, which is the selection of priors and likelihoods that lead to robust estimators.) Suppose that we have determined that the hyperparameter α belongs to some open set \mathcal{A} , perhaps after expert prior elicitation. Ideally, we would calculate the extrema of $\mathbb{E}_{p_\alpha} [g(\theta)]$ as α ranges over all of \mathcal{A} . These extrema are a measure of *global robustness*, and their calculation is intractable or difficult except in special cases (Moreno, 2000; Huber, 2011, Chapter 15). A more practical alternative is to examine how much $\mathbb{E}_{p_\alpha} [g(\theta)]$ changes locally in response to small perturbations in the value of α near some tentative guess, $\alpha_0 \in \mathcal{A}$. To this end we define the *local sensitivity at α_0* (Gustafson, 2000).

Definition 1 *The local sensitivity of $\mathbb{E}_{p_\alpha} [g(\theta)]$ to hyperparameter α at α_0 is given by*

$$\mathbf{S}_{\alpha_0} := \left. \frac{d\mathbb{E}_{p_\alpha} [g(\theta)]}{d\alpha} \right|_{\alpha_0}. \quad (1)$$

\mathbf{S}_{α_0} , the local sensitivity, can be considered a measure of *local robustness* (Gustafson, 2000). Throughout the paper we will distinguish between sensitivity, which comprises objectively defined quantities such as \mathbf{S}_{α_0} , and robustness, which we treat as a more subjective concept that may be informed by the sensitivity as well as other considerations. For example, even if one knows \mathbf{S}_{α_0} precisely, how much posterior change is too much change and how much prior variation is reasonable remain decisions to be made by the modeler. For a more in-depth discussion of how we use the terms sensitivity and robustness, see Appendix C.

The quantity \mathbf{S}_{α_0} can be interpreted as measuring sensitivity to hyperparameters within a small region near $\alpha = \alpha_0$ where the posterior dependence on α is approximately linear. Then local sensitivity provides an approximation to global sensitivity in the sense that, to first order,

$$\mathbb{E}_{p_\alpha} [g(\theta)] \approx \mathbb{E}_{p_{\alpha_0}} [g(\theta)] + \mathbf{S}_{\alpha_0}^\top (\alpha - \alpha_0).$$

Generally, the dependence of $\mathbb{E}_{p_\alpha} [g(\theta)]$ on α is not given in any closed form that is easy to differentiate. However, as we will now see, the derivative \mathbf{S}_{α_0} is equal, under mild regularity conditions, to a particular posterior covariance that can easily be estimated with MCMC draws.

2.2 Covariances and Sensitivity

We will first state a general result relating sensitivity and covariance and then apply it to our specific cases of interest as they arise throughout the paper, beginning with the calculation of \mathbf{S}_{α_0} from Section 2.1. Consider a general base density $p_0(\theta)$ defined relative to λ and define $\rho(\theta, \alpha)$ to be a λ -measurable log perturbation function that depends on $\alpha \in \mathcal{A} \subseteq \mathbb{R}^D$. We will require the following mild technical assumption:

Assumption 1 *For all $\alpha \in \mathcal{A}$, $\rho(\theta, \alpha)$ is continuously differentiable with respect to α , and, for a given λ -measurable $g(\theta)$ there exist λ -integrable functions $f_0(\theta)$ and $f_1(\theta)$ such that $|p_0(\theta) \exp(\rho(\theta, \alpha)) g(\theta)| < f_0(\theta)$ and $|p_0(\theta) \exp(\rho(\theta, \alpha))| < f_1(\theta)$.*

Under Assumption 1 we can normalize the log-perturbed quantity $p_0(\theta) \exp(\rho(\theta, \alpha))$ to get a density in θ with respect to λ .

Definition 2 Denote by $p_\alpha(\theta)$ the normalized posterior given α :

$$p_\alpha(\theta) := \frac{p_0(\theta) \exp(\rho(\theta, \alpha))}{\int p_0(\theta') \exp(\rho(\theta', \alpha)) \lambda(d\theta')}. \quad (2)$$

For example, $p_\alpha(\theta)$ defined in Section 2.1 is equivalent to taking $p_0(\theta) = p(\theta|x, \alpha_0)$ and $\rho(\theta, \alpha) = \log p(x|\theta, \alpha) + \log p(\theta|\alpha) - \log p(x|\theta, \alpha_0) - \log p(\theta|\alpha_0)$.

For a λ -measurable function $g(\theta)$, consider differentiating the expectation $\mathbb{E}_{p_\alpha}[g(\theta)]$ with respect to α :

$$\frac{d\mathbb{E}_{p_\alpha}[g(\theta)]}{d\alpha^\top} := \frac{d}{d\alpha} \int p_\alpha(\theta) g(\theta) \lambda(d\theta). \quad (3)$$

When evaluated at some $\alpha_0 \in \mathcal{A}$, this derivative measures the local sensitivity of $\mathbb{E}_{p_\alpha}[g(\theta)]$ to the index α at α_0 . Define $\mathcal{A}_0 \subseteq \mathcal{A}$ to be an open ball containing α_0 . Under Assumption 1 we assume without loss of generality that $\rho(\theta, \alpha_0) \equiv 0$ so that $p_0(\theta) = p_{\alpha_0}(\theta)$; if $\rho(\theta, \alpha_0)$ is non-zero, we can simply incorporate it into the definition of $p_0(\theta)$. Then, under Assumption 1, the derivative in Eq. (3) is equivalent to a particular posterior covariance.

Theorem 1 Under Assumption 1,

$$\left. \frac{d\mathbb{E}_{p_\alpha}[g(\theta)]}{d\alpha^\top} \right|_{\alpha_0} = \text{Cov}_{p_0} \left(g(\theta), \left. \frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \right|_{\alpha_0} \right). \quad (4)$$

Theorem 1 is a straightforward consequence of the Lebesgue dominated convergence theorem; see Appendix A for a detailed proof. Versions of Theorem 1 have appeared many times before; e.g., Diaconis and Freedman (1986); Basu et al. (1996); Gustafson (1996b); Pérez et al. (2006) have contributed variants of this result to the robustness literature.

By using MCMC draws from $p_0(\theta)$ to calculate the covariance on the right-hand side of Eq. (4), one can form an estimate of $d\mathbb{E}_{p_\alpha}[g(\theta)]/d\alpha^\top$ at $\alpha = \alpha_0$. One might also approach the problem of calculating $d\mathbb{E}_{p_\alpha}[g(\theta)]/d\alpha^\top$ using importance sampling as follows (Owen, 2013, Chapter 9). First, an importance sampling estimate of the dependence of $\mathbb{E}_{p_\alpha}[g(\theta)]$ on α can be constructed with weights that depend on α . Then, differentiating the weights with respect to α provides a sample-based estimate of $d\mathbb{E}_{p_\alpha}[g(\theta)]/d\alpha^\top$. We show in Appendix B that this importance sampling approach is equivalent to using MCMC samples to estimate the covariance in Theorem 1.

An immediate corollary of Theorem 1 allows us to calculate \mathbf{S}_{α_0} as a covariance.

Corollary 1 Suppose that Assumption 1 holds for some $\alpha_0 \in \mathcal{A}$, some $g(\theta)$, and for

$$\rho(\theta, \alpha) = \log p(x|\theta, \alpha) + \log p(\theta|\alpha) - \log p(x|\theta, \alpha_0) - \log p(\theta|\alpha_0).$$

Then Theorem 1 implies that

$$\mathbf{S}_{\alpha_0} = \text{Cov}_{p_0} \left(g(\theta), \left. \frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \right|_{\alpha_0} \right). \quad (5)$$

Corollary 1 can be found in Basu et al. (1996), in which a version of Corollary 1 is stated in the proof of their Theorem 1, as well as in Pérez et al. (2006) and Efron (2015). Note that the definition of $\rho(\theta, \alpha)$ does not contain any normalizing constants and so can typically be easily calculated. Given N_s MCMC draws $\{\theta_n\}_{n=1}^{N_s}$ from a chain that we assume to have reached equilibrium at the stationary distribution $p_0(\theta)$, one can calculate an estimate of \mathbf{S}_{α_0} using the sample covariance version of Eq. (4):

$$\hat{\mathbf{S}}_{\alpha_0} := \frac{1}{N_s} \sum_{n=1}^{N_s} g(\theta_n) \left. \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} \right|_{\alpha_0} - \left(\frac{1}{N_s} \sum_{n=1}^{N_s} g(\theta_n) \right) \left(\frac{1}{N_s} \sum_{n=1}^{N_s} \left. \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} \right|_{\alpha_0} \right) \quad (6)$$

for $\theta_n \sim p_0(\theta)$, where $n = 1, \dots, N_s$.

3. Variational Bayesian Covariances and Sensitivity

3.1 Variational Bayes

We briefly review variational Bayes and state our key assumptions about its accuracy. We wish to find an approximate distribution, in some class \mathcal{Q} of tractable distributions, selected to minimize the Kullback-Leibler divergence (KL divergence) between $q \in \mathcal{Q}$ and the exact log-perturbed posterior p_α . We assume that distributions in \mathcal{Q} are parameterized by a finite-dimensional parameter η in some feasible set $\Omega_\eta \subseteq \mathbb{R}^{K_\eta}$.

Definition 3 *The approximating variational family is given by*

$$\mathcal{Q} := \{q : q = q(\theta; \eta) \text{ for } \eta \in \Omega_\eta\}. \quad (7)$$

Given \mathcal{Q} , we define the optimal $q \in \mathcal{Q}$, which we call $q_\alpha(\theta)$, as the distribution that minimizes the KL divergence $KL(q(\theta; \eta) || p_\alpha(\theta))$ from $p_\alpha(\theta)$. We denote the corresponding optimal variational parameters as η^* .

Definition 4 *The variational approximation $q_\alpha(\theta)$ to $p_\alpha(\theta)$ is defined by*

$$q_\alpha(\theta) := q(\theta; \eta^*) := \operatorname{argmin}_{q \in \mathcal{Q}} \{KL(q(\theta; \eta) || p_\alpha(\theta))\}, \quad (8)$$

where

$$KL(q(\theta; \eta) || p_\alpha(\theta)) = \mathbb{E}_{q(\theta; \eta)} [\log q(\theta; \eta) - \log p_\alpha(\theta)].$$

In the KL divergence, the (generally intractable) normalizing constant for $p_\alpha(\theta)$ does not depend on $q(\theta)$ and so can be neglected when optimizing. In order for the KL divergence to be well defined, we assume that both $p_0(\theta)$ and $q(\theta)$ are given with respect to the same base measure, λ , and that the support of $q(\theta)$ is contained in the support of $p_\alpha(\theta)$. We will require some additional mild regularity conditions in Section 3.2 below.

A common choice for the approximating family \mathcal{Q} in Eq. (7) is the “mean field family” (Wainwright and Jordan, 2008; Blei et al., 2016),

$$\mathcal{Q}_{mf} := \left\{ q(\theta) : q(\theta) = \prod_k q(\theta_k; \eta_k) \right\}, \quad (9)$$

where k indexes a partition of the full vector θ and of the parameter vector η . That is, \mathcal{Q}_{mf} approximates the posterior $p_\alpha(\theta)$ as a distribution that factorizes across sub-components of θ . This approximation is commonly referred to as “MFVB,” for “mean field variational Bayes.” Note that, in general, each function $q(\theta_k; \eta_k)$ in the product is different. For notational convenience we write $q(\theta_k; \eta_k)$ instead of $q_k(\theta_k; \eta_k)$ when the arguments make it clear which function we are referring to, much as the same symbol p is used to refer to many different probability distributions without additional indexing.

One may additionally assume that the components $q(\theta_k; \eta_k)$ are in a convenient exponential family. Although the exponential family assumption does not in general follow from a factorizing assumption, for compactness we will refer to both the factorization and the exponential family assumption as MFVB.

In an MFVB approximation, Ω_η could be a stacked vector of the natural parameters of the exponential families, or the moment parameterization, or perhaps a transformation of these parameters into an unconstrained space (e.g., the entries of the log-Cholesky decomposition of a positive definite information matrix). For more concrete examples, see Section 5. Although all of our experiments and much of our motivating intuition will use MFVB, our results extend to other choices of \mathcal{Q} that satisfy the necessary assumptions.

3.2 Variational Bayes sensitivity

Just as MCMC approximations lend themselves to moment calculations, the variational form of VB approximations lends itself to sensitivity calculations. In this section we derive the sensitivity of VB posterior means to generic perturbations—a VB analogue of Theorem 1. In Section 4 we will choose particular perturbations to calculate VB prior sensitivity and, through Theorem 1, posterior covariances.

In Definition 4, the variational approximation is a function of α through the optimal parameters $\eta^*(\alpha)$, i.e., $q_\alpha(\theta) = q(\theta, \eta^*(\alpha))$. In turn, the posterior expectation $\mathbb{E}_{q_\alpha}[g(\theta)]$ is also a function of α , and its derivative at α_0 —the local sensitivity of the variational approximation to α —has a closed form under the following mild technical conditions. As with p_0 , define $q_0 := q_{\alpha_0}$, and define $\eta_0^* := \eta^*(\alpha_0)$.

All the following assumptions are intended to hold for a given $p_\alpha(\theta)$, approximating class \mathcal{Q} , λ -measurable function $g(\theta)$, and to hold for all $\alpha \in \mathcal{A}_0$ and all η in an open neighborhood of η_0^* .

Assumption 2 *The KL divergence at $KL(q(\theta; \eta) || p_0(\theta))$ and expected log perturbation $\mathbb{E}_{q(\theta; \eta)}[\rho(\theta, \alpha)]$ are twice continuously differentiable in η and α .*

Assumption 3 *There exists a strict local minimum, $\eta^*(\alpha)$, of $KL(q(\theta; \eta) || p_\alpha(\theta))$ in Eq. (8) such that $\eta^*(\alpha)$ is interior to Ω_η .*

Assumption 4 *The expectation $\mathbb{E}_{q(\theta; \eta)}[g(\theta)]$ is a continuously differentiable function of η .*

We define the following quantities for notational convenience.

Definition 5 *Define the following derivatives of variational expectations evaluated at the optimal parameters:*

$$\mathbf{H}_{\eta\eta} := \left. \frac{\partial^2 KL(q(\theta; \eta) || p_0(\theta))}{\partial \eta \partial \eta^\top} \right|_{\eta=\eta_0^*} \quad \mathbf{f}_{\alpha\eta} := \left. \frac{\partial^2 \mathbb{E}_{q(\theta; \eta)}[\rho(\theta, \alpha)]}{\partial \alpha \partial \eta^\top} \right|_{\eta=\eta_0^*, \alpha=\alpha_0} \quad \mathbf{g}_\eta := \left. \frac{\partial \mathbb{E}_{q(\theta; \eta)}[g(\theta)]}{\partial \eta^\top} \right|_{\eta=\eta_0^*}.$$

Since $g(\theta)$, α , and η are all vectors, the quantities $\mathbf{H}_{\eta\eta}$, $\mathbf{f}_{\alpha\eta}$, and \mathbf{g}_η are matrices. We are now ready to state a VB analogue of Theorem 1.

Theorem 2 *Consider a variational approximation $q_\alpha(\theta)$ to $p_\alpha(\theta)$ as given in Definition 4 and a λ -measurable function $g(\theta)$. Then, under Assumptions 1–4, using the definitions given in Definition 5, we have*

$$\left. \frac{d\mathbb{E}_{q_\alpha}[g(\theta)]}{d\alpha^\top} \right|_{\alpha_0} = \mathbf{g}_\eta \mathbf{H}_{\eta\eta}^{-1} \mathbf{f}_{\alpha\eta}^\top. \quad (10)$$

A proof of Theorem 2 is given in Appendix D. As with Theorem 1, by choosing the appropriate $\rho(\theta, \alpha)$ and evaluating $\mathbf{f}_{\alpha\eta}$, we can use Theorem 2 to calculate the exact sensitivity of VB solutions to any arbitrary local perturbations that satisfy the regularity conditions. Assumptions 1–4 are typically not hard to verify. For an example, see Appendix E, where we establish Assumptions 1–4 for a multivariate normal target distribution and a mean-field approximation.

Eq. (10) is formally similar to frequentist sensitivity estimates. For example, the pioneering paper of Cook (1986) contains a formula for assessing the curvature of a marginal likelihood surface (Cook, 1986, Equation 15) that, like our Theorem 2, represents the sensitivity as a linear system involving the Hessian of an objective function at its optimum. The geometric interpretation of local robustness suggested by Cook (1986) has been extended to Bayesian settings (see, for example, Zhu et al. (2007, 2011)). In addition to generality, one attractive aspect of their geometric approach is its invariance to parameterization. Investigating geometric interpretations of the present work may be an interesting avenue for future research.

3.3 Approximating with Variational Bayes

Recall that we are ultimately interested in $\mathbb{E}_{p_\alpha}[g(\theta)]$. Variational approximations and their sensitivity measures will be useful to the extent that both the variational means and sensitivities are close to the exact means and sensitivities. We formalize these desiderata as follows.

Condition 1 Under Assumptions 1–4 and the quantities defined therein, we additionally have, for all $\alpha \in \mathcal{A}$,

$$\mathbb{E}_{q_\alpha} [g(\theta)] \approx \mathbb{E}_{p_\alpha} [g(\theta)] \quad \text{and} \quad (11)$$

$$\left. \frac{d\mathbb{E}_{q_\alpha} [g(\theta)]}{d\alpha^\top} \right|_{\alpha_0} \approx \left. \frac{d\mathbb{E}_{p_\alpha} [g(\theta)]}{d\alpha^\top} \right|_{\alpha_0} \quad (12)$$

We will not attempt to be precise about what we mean by the ‘‘approximately equal’’ sign, since we are not aware of any practical tools for evaluating quantitatively whether Condition 1 holds other than running both VB and MCMC (or some other slow but accurate posterior approximation) and comparing the results. However, VB has been useful in practice to the extent that Condition 1 holds true for at least some parameters of interest. We provide some intuition for when Condition 1 might hold in Section 5.1, and will evaluate Condition 1 in each of our experiments below by comparing the VB and MCMC posterior approximate means and sensitivities.

Since Condition 1 holds only for a particular choice of $g(\theta)$, it is weaker than the assumption that q_α is close to p_α in KL divergence, or even that all the posterior means are accurately estimated. For example, as discussed in Appendix B of Giordano et al. (2015) and in Section 10.1.2 of Bishop (2006), a mean-field approximation to a multivariate normal posterior produces inaccurate covariances and may have an arbitrarily bad KL divergence from p_α , but Condition 1 holds exactly for the location parameters. We discuss the multivariate normal example further in Section 4.1 and Section 5.1 below.

4. Calculation and Uses of Sensitivity

In this section, we discuss two applications of Theorem 1 and Theorem 2: calculating improved covariance estimates and prior sensitivity measures for MFVB. Throughout this section, we will assume that we can apply Theorem 1 and Theorem 2 unless stated otherwise.

4.1 Covariances for Variational Bayes

Consider the mean field approximating family, \mathcal{Q}_{mf} , from Section 3.1 and a fixed exact posterior $p_0(\theta)$. It is well known that the resulting marginal variances also tend to be under-estimated even when parameters means are well-estimated (see, e.g., (MacKay, 2003; Wang and Titterton, 2004; Turner and Sahani, 2011; Bishop, 2006, Chapter 10)). Even more obviously, any $q \in \mathcal{Q}_{mf}$ yields zero as its estimate of the covariance between sub-components of θ that are in different factors of the mean field approximating family. It is therefore unreasonable to expect that $\text{Cov}_{q_0}(g(\theta)) \approx \text{Cov}_{p_0}(g(\theta))$. However, if Condition 1 holds, we may expect the sensitivity of MFVB means to certain perturbations to be accurate by Condition 1, and, by Theorem 1, we expect the corresponding covariances to be accurately estimated by the MFVB sensitivity. In particular, by taking $\rho(\theta, \alpha) = \alpha^\top g(\theta)$ and $\alpha_0 = 0$, we have by Condition 1 that

$$\left. \frac{d\mathbb{E}_{q_\alpha} [g(\theta)]}{d\alpha^\top} \right|_{\alpha=0} \approx \left. \frac{d\mathbb{E}_{p_\alpha} [g(\theta)]}{d\alpha^\top} \right|_{\alpha=0} = \text{Cov}_{p_0}(g(\theta)). \quad (13)$$

We can consequently use Theorem 2 to provide an estimate of $\text{Cov}_{p_0}(g(\theta))$ that may be superior to $\text{Cov}_{q_0}(g(\theta))$. With this motivation in mind, we make the following definition.

Definition 6 The linear response variational Bayes (LRVB) approximation, $\text{Cov}_{q_0}^{LR}(g(\theta))$, is given by

$$\text{Cov}_{q_0}^{LR}(g(\theta)) := \mathbf{g}_\eta \mathbf{H}_{\eta\eta}^{-1} \mathbf{g}_\eta^\top. \quad (14)$$

Corollary 2 For a given $p_0(\theta)$, class \mathcal{Q} , and function $g(\theta)$, when Assumptions 1–4 and Condition 1 hold for $\rho(\theta, \alpha) = \alpha^\top g(\theta)$ and $\alpha_0 = 0$, then

$$\text{Cov}_{q_0}^{LR}(g(\theta)) \approx \text{Cov}_{p_0}(g(\theta)).$$

The strict optimality of η_0^* in Assumption 3 guarantees that $\mathbf{H}_{\eta\eta}$ will be positive definite and symmetric, and, as desired, the covariance estimate $\text{Cov}_{q_0}^{LR}(g(\theta))$ will be positive semidefinite and symmetric. Since the optimal value of every component of $\mathbb{E}_{q_\alpha}[g(\theta)]$ may be affected by the log perturbation $\alpha^\top g(\theta)$, $\text{Cov}_{q_0}^{LR}(g(\theta))$ can estimate non-zero covariances between elements of $g(\theta)$ even when they have been partitioned into separate factors of the mean field approximation.

Note that $\text{Cov}_{q_0}^{LR}(g(\theta))$ and $\text{Cov}_{q_0}(g(\theta))$ differ only when there are at least some moments of p_0 that q_0 fails to accurately estimate. In particular, if q_α provided a good approximation to p_α for both the first and second moments of $g(\theta)$, then we would have $\text{Cov}_{q_0}^{LR}(g(\theta)) \approx \text{Cov}_{q_0}(g(\theta))$ since, for q_0 and p_0 ,

$$\begin{aligned} \mathbb{E}_{q_0}[g(\theta)] &\approx \mathbb{E}_{p_0}[g(\theta)] \text{ and} \\ \mathbb{E}_{q_0}[g(\theta)g(\theta)^\top] &\approx \mathbb{E}_{p_0}[g(\theta)g(\theta)^\top] \Rightarrow \\ \text{Cov}_{q_0}(g(\theta)) &\approx \text{Cov}_{p_0}(g(\theta)), \end{aligned}$$

and, for q_α and p_α ,

$$\begin{aligned} \mathbb{E}_{q_\alpha}[g(\theta)] &\approx \mathbb{E}_{p_\alpha}[g(\theta)] \Rightarrow \\ \text{Cov}_{q_0}^{LR}(g(\theta)) &\approx \text{Cov}_{p_0}(g(\theta)). \end{aligned}$$

Putting these two approximate equalities together, we see that, when the first and second moments of q_α approximately match those of p_α ,

$$\text{Cov}_{q_0}(g(\theta)) \approx \text{Cov}_{q_0}^{LR}(g(\theta)).$$

However, in general, $\text{Cov}_{q_0}^{LR}(g(\theta)) \neq \text{Cov}_{q_0}(g(\theta))$. In this sense, any discrepancy between $\text{Cov}_{q_0}^{LR}(g(\theta))$ and $\text{Cov}_{q_0}(g(\theta))$ indicates an inadequacy of the variational approximation for at least the second moments of $g(\theta)$.

Let us consider a simple concrete illustrative example which will demonstrate both how $\text{Cov}_{q_0}(g(\theta))$ can be a poor approximation to $\text{Cov}_{p_0}(g(\theta))$ and how $\text{Cov}_{q_0}^{LR}(g(\theta))$ can improve the approximation for some moments but not others. Suppose that the exact posterior is a bivariate normal,

$$p_0(\theta) = \mathcal{N}(\theta|\mu, \Sigma), \tag{15}$$

where $\theta = (\theta_1, \theta_2)^\top$, $\mu = (\mu_1, \mu_2)^\top$, Σ is invertible, and $\Lambda := \Sigma^{-1}$. One may think of μ and Σ as known functions of x via Bayes' theorem, for example, as given by a normal-normal conjugate model. Suppose we use the MFVB approximating family

$$\mathcal{Q}_{mf} = \{q(\theta) : q(\theta) = q(\theta_1)q(\theta_2)\}.$$

One can show (see Appendix E) that the optimal MFVB approximation to p_α in the family \mathcal{Q}_{mf} is given by

$$\begin{aligned} q_0(\theta_1) &= \mathcal{N}(\theta_1|\mu_1, \Lambda_{11}^{-1}) \\ q_0(\theta_2) &= \mathcal{N}(\theta_2|\mu_2, \Lambda_{22}^{-1}). \end{aligned}$$

Note that the posterior mean of θ_1 is exactly estimated by the MFVB procedure:

$$\mathbb{E}_{q_0}[\theta_1] = \mu_1 = \mathbb{E}_{p_0}[\theta_1].$$

However, if $\Sigma_{12} \neq 0$, then $\Lambda_{11}^{-1} < \Sigma_{11}$, and the variance of θ_1 is underestimated. It follows that the expectation of θ_1^2 is *not* correctly estimated by the MFVB procedure:

$$\mathbb{E}_{q_\alpha}[\theta_1^2] = \mu_1^2 + \Lambda_{11}^{-1} < \mu_1^2 + \Sigma_{11} = \mathbb{E}_{p_\alpha}[\theta_1^2].$$

An analogous statement holds for θ_2 . Of course, the covariance is also mis-estimated if $\Sigma_{12} \neq 0$ since, by construction of the MFVB approximation,

$$\text{Cov}_{q_0}(\theta_1, \theta_2) = 0 \neq \Sigma_{12} = \text{Cov}_{p_0}(\theta_1, \theta_2).$$

Now let us take the log perturbation $\rho(\theta, \alpha) = \theta_1 \alpha_1 + \theta_2 \alpha_2$. For all α in a neighborhood of zero, the log-perturbed posterior given by Eq. (2) remains multivariate normal, so it remains the case that, as a function of α , $\mathbb{E}_{q_\alpha}[\theta_1] = \mathbb{E}_{p_\alpha}[\theta_1]$ and $\mathbb{E}_{q_\alpha}[\theta_2] = \mathbb{E}_{p_\alpha}[\theta_2]$. Again, see Appendix E for a detailed proof. Consequently, Condition 1 holds with equality (not approximate equality) when $g(\theta) = \theta$. However, since the second moments are not accurate (irrespective of α), Condition 1 does not hold exactly when $g(\theta) = (\theta_1^2, \theta_2^2)^\top$, nor when $g(\theta) = \theta_1 \theta_2$. (Condition 1 may still hold approximately for second moments when Σ_{12} is small.) The fact that Condition 1 holds with equality for $g(\theta) = \theta$ allows us to use Theorem 1 and Theorem 2 to calculate $\text{Cov}_{q_0}^{LR}(g(\theta)) = \text{Cov}_{p_0}(g(\theta))$, even though $\mathbb{E}_{p_0}[\theta_1 \theta_2]$ and $\mathbb{E}_{p_0}[(\theta_1^2, \theta_2^2)^\top]$ are mis-estimated.

In fact, when Condition 1 holds with equality for some θ_i , then the estimated covariance in Eq. (14) for all terms involving θ_i will be exact as well. Condition 1 holds with equality for the means of θ_i in the bivariate normal model above, and in fact holds for the general multivariate normal case, as described in Appendix E. Below, in Section 5, in addition to robustness measures, we will also report the accuracy of Eq. (14) for estimating posterior covariances. We find that, for most parameters of interest, particularly location parameters, $\text{Cov}_{q_0}^{LR}(g(\theta))$ provides a good approximation to $\text{Cov}_{p_0}(g(\theta))$.

4.2 Linear Response Covariances in Previous Literature

The application of sensitivity measures to VB problems for the purpose of improving covariance estimates has a long history under the name “linear response methods.” These methods originated in the statistical physics literature (Tanaka, 2000; Oppen and Saad, 2001) and have been applied to various statistical and machine learning problems (Kappen and Rodriguez, 1998; Tanaka, 1998; Welling and Teh, 2004; Oppen and Winther, 2004). The current paper, which builds on this line of work and on our earlier work (Giordano et al., 2015), represents a simplification and generalization of classical linear response methods and serves to elucidate the relationship between these methods and the local robustness literature. In particular, while Giordano et al. (2015) focused on moment-parameterized exponential families, we derive linear-response covariances for generic variational approximations and connect the linear-response methodology to the Bayesian robustness literature.

A very reasonable approach to address the inadequacy of MFVB covariances is simply to increase the expressiveness of the model class \mathcal{Q} —although, as noted by Turner and Sahani (2011), increased expressiveness does not necessarily lead to better posterior moment estimates. This approach is taken by much of the recent VB literature (e.g., Tran et al., 2015a,b; Ranganath et al., 2016; Rezende and Mohamed, 2015; Liu and Wang, 2016). Though this research direction remains lively and promising, the use of a more complex class \mathcal{Q} sometimes sacrifices the speed and simplicity that made VB attractive in the first place, and often without the relatively well-understood convergence guarantees of MCMC. We also stress that the current work is not necessarily at odds with the approach of increasing expressiveness. Sensitivity methods can be a supplement to any VB approximation for which our estimators, which require solving a linear system involving the Hessian of the KL divergence, are tractable.

4.3 The Laplace Approximation and Linear Response Covariances

In this section, we briefly compare linear response covariances to the Laplace approximation (Gelman et al., 2014, Chapter 13). The Laplace approximation to $p_0(\theta)$ is formed by first finding the “maximum *a posteriori*” (MAP) estimate,

$$\hat{\theta}_{Lap} := \underset{\theta}{\operatorname{argmax}} p_0(\theta), \quad (16)$$

and then forming the multivariate normal posterior approximation

$$\mathbf{H}_{Lap} := - \left. \frac{\partial^2 p_0(\theta)}{\partial \theta \partial \theta^\top} \right|_{\hat{\theta}_{Lap}} \quad (17)$$

$$\begin{aligned} \text{Cov}_{q_{Lap}^{Lap}}(\theta) &:= \mathbf{H}_{Lap}^{-1} \\ q_{Lap}(\theta) &:= \mathcal{N}\left(\theta \mid \hat{\theta}_{Lap}, \text{Cov}_{q_{Lap}^{Lap}}(\theta)\right). \end{aligned} \quad (18)$$

Since both LRVB and the Laplace approximation require the solution of an optimization problem (Eq. (8) and Eq. (16) respectively) and the estimation of covariances via an inverse Hessian of the optimization objective (Eq. (14) and Eq. (17) respectively), it will be instructive to compare the two approaches.

Following Neal and Hinton (1998), we can, in fact, view the MAP estimator as a special variational approximation, where we define

$$\mathcal{Q}_{Lap} := \left\{ q(\theta; \theta_0) : \int q(\theta; \theta_0) \log p_0(\theta) \lambda(d\theta) = \log p_0(\theta_0) \text{ and} \right. \\ \left. \int q(\theta; \theta_0) \log q(\theta; \theta_0) \lambda(d\theta) = \text{Constant} \right\},$$

where the *Constant* term is constant in θ_0 . That is, \mathcal{Q}_{Lap} consists of “point masses” at θ_0 with constant entropy. Generally such point masses may not be defined as densities with respect to λ , and the *KL* divergence in Eq. (8) may not be formally defined for $q \in \mathcal{Q}_{Lap}$. However, if \mathcal{Q}_{Lap} can be approximated arbitrarily well by well-defined densities (e.g., normal distributions with variance fixed at an arbitrarily small number), then we can use \mathcal{Q}_{Lap} as a heuristic tool for understanding the MAP estimator.

Since \mathcal{Q}_{Lap} contains only point masses, the covariance of the variational approximation is the zero matrix: $\text{Cov}_{q_{Lap}}(\theta) = 0$. Thus, as when one uses the mean field assumption, $\text{Cov}_{q_{Lap}}(\theta)$ underestimates the marginal variances and magnitudes of the covariances of $\text{Cov}_{p_0}(\theta)$. Of course, the standard Laplace approximation uses $\text{Cov}_{q_{Lap}^{Lap}}(\theta)$, not $\text{Cov}_{q_{Lap}}(\theta)$, to approximate $\text{Cov}_{p_0}(\theta)$. In fact, $\text{Cov}_{q_{Lap}^{Lap}}(\theta)$ is equivalent to a linear response covariance matrix calculated for the approximating family \mathcal{Q}_{Lap} :

$$\begin{aligned} KL(q(\theta; \theta_0) \parallel p_0(\theta)) &= -\log p_0(\theta_0) - \text{Constant} \Rightarrow \\ \hat{\theta}_{Lap} &= \underset{\theta}{\text{argmax}} p_0(\theta) = \underset{\theta_0}{\text{argmin}} KL(q(\theta; \theta_0) \parallel p_0(\theta)) = \theta_0^* \\ \mathbf{H}_{Lap} &= - \left. \frac{\partial^2 p_0(\theta)}{\partial \theta \partial \theta^\top} \right|_{\hat{\theta}_{Lap}} = - \left. \frac{\partial^2 KL(q(\theta; \theta_0) \parallel p_0(\theta))}{\partial \theta_0 \partial \theta_0^\top} \right|_{\theta_0^*} = \mathbf{H}_{\eta\eta}. \end{aligned}$$

So $\hat{\theta}_{Lap} = \theta_0^*$, $\mathbf{H}_{Lap} = \mathbf{H}_{\eta\eta}$, and $\text{Cov}_{q_{Lap}^{Lap}}(\theta) = \text{Cov}_{q_0^{LR}}(\theta)$ for the approximating family \mathcal{Q}_{Lap} .

From this perspective, the accuracy of the Laplace approximation depends precisely on the extent to which Condition 1 holds for the family of point masses \mathcal{Q}_{Lap} . Typically, VB approximations use a \mathcal{Q} that is more expressive than \mathcal{Q}_{Lap} , and we might expect Condition 1 to be more likely to apply for a more expressive family. It follows that we might expect the LRVB covariance estimate $\text{Cov}_{q_0^{LR}}$ for general \mathcal{Q} to be more accurate than the Laplace covariance approximation $\text{Cov}_{q_{Lap}^{Lap}}$. We demonstrate the validity of this intuition in the experiments of Section 5.

4.4 Local Prior Sensitivity for MFVB

We now turn to estimating prior sensitivities for MFVB estimates—the variational analogues of \mathbf{S}_{α_0} in Definition 1. First, we define the variational local sensitivity.

Definition 7 *The local sensitivity of $\mathbb{E}_{q_\alpha}[g(\theta)]$ to prior parameter α at α_0 is given by*

$$\mathbf{S}_{\alpha_0}^q := \left. \frac{d\mathbb{E}_{q_\alpha}[g(\theta)]}{d\alpha} \right|_{\alpha_0}.$$

Corollary 3 *Suppose that Assumptions 1–4 and Condition 1 hold for some $\alpha_0 \in \mathcal{A}$ and for*

$$\rho(\theta, \alpha) = \log p(x|\theta, \alpha) + \log p(\theta|\alpha) - \log p(x|\theta, \alpha_0) - \log p(\theta|\alpha).$$

Then $\mathbf{S}_{\alpha_0}^q \approx \mathbf{S}_{\alpha_0}$.

Corollary 3 states that, as with the covariance approximations in Section 4.1, $\mathbf{S}_{\alpha_0}^q$ is a useful approximation to \mathbf{S}_{α_0} to the extent that Condition 1 holds—that is, to the extent that the MFVB means are good approximations to the exact means for the prior perturbations $\alpha \in \mathcal{A}_0$.

Under the $\rho(\theta, \alpha)$ given in Corollary 3, Theorem 2 gives the following formula for the variational local sensitivity:

$$\mathbf{S}_{\alpha_0}^q = \mathbf{g}_\eta \mathbf{H}_{\eta\eta}^{-1} \frac{\partial}{\partial \eta^\top} \mathbb{E}_{q(\theta;\eta)} \left[\left. \frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \right|_{\alpha_0} \right] \Big|_{\eta_0^*}. \quad (19)$$

We now use Eq. (19) to reproduce MFVB versions of some standard robustness measures found in the existing literature. A simple case is when the prior $p(\theta|\alpha)$ is believed to be in a given parametric family, and we are simply interested in the effect of varying the parametric family’s parameters (Basu et al., 1996; Giordano et al., 2016). For illustration, we first consider a simple example where $p(\theta|\alpha)$ is in the exponential family, with natural sufficient statistic θ and log normalizer $A(\alpha)$, and we take $g(\theta) = \theta$. In this case,

$$\begin{aligned} \log p(\theta|\alpha) &= \alpha^\top \theta - A(\alpha) \\ \mathbf{f}_{\alpha\eta} &= \frac{\partial}{\partial \eta^\top} \mathbb{E}_{q(\theta;\eta)} \left[\left. \frac{\partial}{\partial \alpha} (\alpha^\top \theta - A(\alpha)) \right|_{\alpha_0} \right] \Big|_{\eta_0^*} \\ &= \left(\frac{\partial}{\partial \eta^\top} \mathbb{E}_{q(\theta;\eta)} [\theta] - \frac{\partial}{\partial \eta^\top} \left. \frac{\partial A(\alpha)}{\partial \alpha} \right|_{\alpha_0} \right) \Big|_{\eta_0^*} \\ &= \frac{\partial}{\partial \eta^\top} \mathbb{E}_{q(\theta;\eta)} [\theta] \Big|_{\eta_0^*} \\ &= \mathbf{g}_\eta. \end{aligned}$$

Note that when $\mathbf{f}_{\alpha\eta} = \mathbf{g}_\eta$, Eq. (19) is equivalent to Eq. (14). So we see that

$$\mathbf{S}_{\alpha_0}^q = \text{Cov}_{q_0}^{LR}(\theta).$$

In this case, the sensitivity is simply the linear response covariance estimate of the covariance, $\text{Cov}_{q_0}^{LR}(\theta)$. By the same reasoning, the exact posterior sensitivity is given by

$$\mathbf{S}_{\alpha_0} = \text{Cov}_{p_0}(\theta).$$

Thus, $\mathbf{S}_{\alpha_0}^q \approx \mathbf{S}_{\alpha_0}$ to the extent that $\text{Cov}_{q_0}^{LR}(\theta) \approx \text{Cov}_{p_0}(\theta)$, which again holds to the extent that Condition 1 holds. Note that if we had used a mean field assumption and had tried to use the direct, uncorrected response covariance $\text{Cov}_{q_0}(\theta)$ to try to evaluate $\mathbf{S}_{\alpha_0}^q$, we would have erroneously concluded that the prior on one component, θ_{k_1} , would not affect the posterior mean of some other component, θ_{k_2} , for $k_2 \neq k_1$.

Sometimes it is easy to evaluate the derivative of the log prior even when it is not easy to normalize it. As an example, we will show how to calculate the local sensitivity to the concentration parameter of an LKJ prior (Lewandowski et al., 2009) under an inverse Wishart variational approximation. The LKJ prior is defined as follows. Let Σ (as part of θ) be an unknown $K \times K$ covariance matrix. Define the $K \times K$ scale matrix \mathbf{M} such that

$$\mathbf{M}_{ij} = \begin{cases} \sqrt{\Sigma_{ij}} & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

Using \mathbf{M} , define the correlation matrix \mathbf{R} as

$$\mathbf{R} = \mathbf{M}^{-1} \boldsymbol{\Sigma} \mathbf{M}^{-1}.$$

The LKJ prior on the covariance matrix \mathbf{R} with concentration parameter $\alpha > 0$ is given by:

$$p_{\text{LKJ}}(\mathbf{R}|\alpha) \propto |\mathbf{R}|^{\alpha-1}.$$

The Stan manual recommends the use of p_{LKJ} , together with an independent prior on the diagonal entries of the scaling matrix \mathbf{M} , for the prior on a covariance matrix that appears in a hierarchical model (Stan Team, 2015, Chapter 9.13).

Suppose that we have chosen the variational approximation

$$q(\boldsymbol{\Sigma}) := \text{InverseWishart}(\boldsymbol{\Sigma}|\boldsymbol{\Psi}, \nu),$$

where $\boldsymbol{\Psi}$ is a positive definite scale matrix and ν is the number of degrees of freedom. In this case, the variational parameters are $\eta = (\boldsymbol{\Psi}, \nu)$. We write η with the understanding that we have stacked only the upper-diagonal elements of $\boldsymbol{\Psi}$ since $\boldsymbol{\Psi}$ is constrained to be symmetric and η^* must be interior. As we show in Appendix G,

$$\mathbb{E}_q[\log p_{\text{LKJ}}(\mathbf{R}|\alpha)] = (\alpha - 1) \left(\log |\boldsymbol{\Psi}| - \psi_K\left(\frac{\nu}{2}\right) - \sum_{k=1}^K \log\left(\frac{1}{2}\boldsymbol{\Psi}_{kk}\right) + K\psi\left(\frac{\nu - K + 1}{2}\right) \right) + \text{Constant},$$

where *Constant* contains terms that do not depend on α , and where ψ_K denotes the multivariate digamma function. Consequently, we can evaluate

$$\begin{aligned} \mathbf{f}_{\alpha\eta} &= \frac{\partial}{\partial \eta^\top} \mathbb{E}_{q(\theta;\eta)} \left[\frac{\partial}{\partial \alpha} \log p(\boldsymbol{\Sigma}|\alpha) \right] \Bigg|_{\eta=\eta_0^*, \alpha=\alpha_0} \\ &= \frac{\partial}{\partial \eta^\top} \left(\log |\boldsymbol{\Psi}| - \psi_K\left(\frac{n}{2}\right) - \sum_{k=1}^K \log\left(\frac{1}{2}\boldsymbol{\Psi}_{kk}\right) + K\psi\left(\frac{n - K + 1}{2}\right) \right) \Bigg|_{\eta_0^*}. \end{aligned} \quad (20)$$

This derivative has a closed form, but the bookkeeping required to represent an unconstrained parameterization of the matrix $\boldsymbol{\Psi}$ within η would be tedious. In practice, we evaluate terms like $\mathbf{f}_{\alpha\eta}$ using automatic differentiation tools (Baydin et al., 2018).

Finally, in cases where we cannot evaluate $\mathbb{E}_{q(\theta;\eta)}[\log p(\theta|\alpha)]$ in closed form as a function of η , we can use numerical techniques as described in Section 4.5. We thus view $\mathbf{S}_{\alpha_0}^q$ as the exact sensitivity to an approximate KL divergence.

4.5 Practical Considerations when Computing the Sensitivity of Variational Approximations

We briefly discuss practical issues in the computation of Eq. (10), which requires calculating the product $\mathbf{g}_\eta \mathbf{H}_{\eta\eta}^{-1}$ (or, equivalently, $\mathbf{H}_{\eta\eta}^{-1} \mathbf{g}_\eta^\top$ since $\mathbf{H}_{\eta\eta}$ is symmetric). Calculating $\mathbf{H}_{\eta\eta}$ and solving this linear system can be the most computationally intensive part of computing Eq. (10).

We first note that it can be difficult and time consuming in practice to manually derive and implement second-order derivatives. Even a small programming error can lead to large errors in Theorem 2. To ensure accuracy and save analyst time, we evaluated all the requisite derivatives using the Python `autograd` automatic differentiation library (Maclaurin et al., 2015) and the Stan math automatic differentiation library (Carpenter et al., 2015).

Note that the dimension of $\mathbf{H}_{\eta\eta}$ is as large as that of η , the parameters that specify the variational distribution $q(\theta; \eta)$. Many applications of MFVB employ many latent variables, the number of which may even scale with the amount of data—including several of the cases that we examine in Section 5. However, these applications typically have special structure that render $\mathbf{H}_{\eta\eta}$ sparse, allowing the practitioner to calculate $\mathbf{g}_\eta \mathbf{H}_{\eta\eta}^{-1}$

quickly. Consider, for example, a model with “global” parameters, θ_{glob} , that are shared by all the individual datapoint likelihoods, and “local” parameters, $\theta_{loc,n}$, associated with likelihood of a single datapoint indexed by n . By “global” and “local” we mean the likelihood and assumed variational distribution factorize as

$$p(x, \theta_{glob}, \theta_{loc,1}, \dots, \theta_{loc,N}) = p(\theta_{glob}) \prod_{n=1}^N p(x|\theta_{loc,n}, \theta_{glob}) p(\theta_{loc,n}|\theta_{glob}) \quad (21)$$

$$q(\theta; \eta) = q(\theta_{glob}; \eta_{glob}) \prod_{n=1}^N q(\theta_{loc,n}; \eta_n) \text{ for all } q(\theta; \eta) \in \mathcal{Q}.$$

In this case, the second derivatives of the variational objective between the parameters for local variables vanish:

$$\text{for all } n \neq m, \frac{\partial^2 KL(q(\theta; \eta) || p_0(\theta))}{\partial \eta_{loc,n} \partial \eta_{loc,m}^\top} = 0.$$

The model in Section 5.3 has such a global / local structure; see Section 5.3.2 for more details. Additional discussion, including the use of Schur complements to take advantage of sparsity in the log likelihood, can be found in Giordano et al. (2015).

When even calculating or instantiating $\mathbf{H}_{\eta\eta}$ is prohibitively time-consuming, one can use conjugate gradient algorithms to approximately compute $\mathbf{H}_{\eta\eta}^{-1} \mathbf{g}_\eta^\top$ (Wright and Nocedal, 1999, Chapter 5). The advantage of conjugate gradient algorithms is that they approximate $\mathbf{H}_{\eta\eta}^{-1} \mathbf{g}_\eta^\top$ using only the Hessian-vector product $\mathbf{H}_{\eta\eta} \mathbf{g}_\eta^\top$, which can be computed efficiently using automatic differentiation without ever forming the full Hessian $\mathbf{H}_{\eta\eta}$. See, for example, the `hessian_vector_product` method of the Python `autograd` package (Maclaurin et al., 2015). Note that a separate conjugate gradient problem must be solved for each column of \mathbf{g}_η^\top , so if the parameter of interest $g(\theta)$ is high-dimensional it may be faster to pay the price for computing and inverting the entire matrix $\mathbf{H}_{\eta\eta}$. See 5.3.2 for more discussion of a specific example.

In Theorem 2, we require η_0^* to be at a true local optimum. Otherwise the estimated sensitivities may not be reliable (e.g., the covariance implied by Eq. (14) may not be positive definite). We find that the classical MFVB coordinate ascent algorithms (Blei et al. (2016, Section 2.4)) and even quasi-second order methods, such as BFGS (e.g., Regier et al., 2015), may not actually find a local optimum unless run for a long time with very stringent convergence criteria. Consequently, we recommend fitting models using second-order Newton trust region methods. When the Hessian is slow to compute directly, as in Section 5, one can use the conjugate gradient trust region method of Wright and Nocedal (1999, Chapter 7), which takes advantage of fast automatic differentiation Hessian-vector products without forming or inverting the full Hessian.

5. Experiments

We now demonstrate the speed and effectiveness of linear response methods on a number of simulated and real data sets. We begin with simple simulated data to provide intuition for how linear response methods can improve estimates of covariance relative to MFVB and the Laplace approximation. We then develop linear response covariance estimates for ADVI and apply them to four real-world models and data sets taken from the Stan examples library (Stan Team, 2017). Finally, we calculate both linear response covariances and prior sensitivity measures for a large-scale industry data set. In each case, we compare linear response methods with ordinary MFVB, the Laplace approximation, and MCMC. We show that linear response methods provide the best approximation to MCMC while still retaining the speed of approximate methods. Code and instructions to reproduce the results of this section can be found in the git repository `rgiordan/CovariancesRobustnessVBPaper`.

5.1 Simple Expository Examples

In this section we provide a sequence of simple examples comparing MFVB and LRVB with Laplace approximations. These examples provide intuition for the covariance estimate $\text{Cov}_{q_0}^{LR}(g(\theta))$ and illustrate how

the sensitivity analysis motivating $\text{Cov}_{q_0}^{LR}(g(\theta))$ differs from the local posterior approximation motivating $\text{Cov}_{q_{Lap}}^{Lap}(g(\theta))$.

For each example, we will explicitly specify the target posterior $p_0(\theta)$ using a mixture of normals. This will allow us to define known target distributions with varying degrees of skewness, over-dispersion, or correlation and compare the truth with a variational approximation. Formally, for some fixed K_z , component indicators z_k , $k = 1, \dots, K_z$, component probabilities π_k , locations μ_k , and covariances Σ_k , we set

$$p(z) = \prod_{k=1}^{K_z} \pi_k^{z_k}$$

$$p_0(\theta) = \sum_z p(z) p(\theta|z) = \sum_z p(z) \prod_{k=1}^{K_z} \mathcal{N}(\theta; m_k, \Sigma_k)^{z_k}.$$

The values π , m and Σ will be chosen to achieve the desired shape for each example using up to $K_z = 3$ components. There will be no need to state the precise values of π , m , and Σ ; rather, we will show plots of the target density and report the marginal means and variances, calculated by Monte Carlo.¹

We will be interested in estimating the mean and variance of the first component, so we take $g(\theta) = \theta_1$. Consequently, in order to calculate $\text{Cov}_{q_0}^{LR}(\theta_1)$, we will be considering the perturbation $\rho(\theta, \alpha) = \alpha\theta_1$ with scalar α and $\alpha_0 = 0$.

For the variational approximations, we will use a factorizing normal approximation:

$$\mathcal{Q}_{mf} = \left\{ q(\theta) : q(\theta) = \prod_{k=1}^K \mathcal{N}(\theta_k; \mu_k, \sigma_k^2) \right\}.$$

In terms of Eq. (7), we take $\eta = (\mu_1, \dots, \mu_K, \log \sigma_1, \dots, \log \sigma_K)^\top$. Thus $\mathbb{E}_{q(\theta; \eta)}[g(\theta)] = \mathbb{E}_{q(\theta; \eta)}[\theta_1] = \mu_1$. In the examples below, we will use multiple distinct components in the definition of $p_0(\theta)$, so that $p_0(\theta)$ is non-normal and $p_0(\theta) \notin \mathcal{Q}_{mf}$.

Since the expectation $\mathbb{E}_{q(\theta; \eta)}[\log p(\theta)]$ is intractable, we replace the exact KL divergence with a Monte Carlo approximation using the “re-parameterization trick” (Kingma and Welling, 2013; Rezende et al., 2014; Titsias and Lázaro-Gredilla, 2014). Let \circ denote the Hadamard (component-wise) product. Let $\xi_m \stackrel{iid}{\sim} \mathcal{N}(0, I_K)$ for $m = 1, \dots, M$. We define

$$\theta_m := \sigma \circ \xi_m + \mu$$

$$KL_{approx}(q(\theta; \eta) || p_0(\theta)) := -\frac{1}{M} \sum_{m=1}^M \log p_0(\theta_m) - \sum_{k=1}^K \log \sigma_k,$$

which is a Monte Carlo estimate of $KL(q(\theta; \eta) || p_0(\theta))$. We found $M = 10000$ to be more than adequate for our present purposes of illustration. Note that we used the same draws ξ_m for both optimization and for the calculation of $\mathbf{H}_{\eta\eta}$ in order to ensure that the η_0^* at which $\mathbf{H}_{\eta\eta}$ was evaluated was in fact an optimum. This approach is similar to our treatment of ADVI; see Section 5.2 for a more detailed discussion.

5.1.1 MULTIVARIATE NORMAL TARGETS

If we take only a single component in the definition of $p_0(\theta)$ ($K_z = 1$), then $p_\alpha(\theta)$ is a multivariate normal distribution for all α , and the Laplace approximation $q_{Lap}(\theta)$ is equal to $p_\alpha(\theta)$ for all α . Furthermore, as discussed in Section 4.1 and Appendix E, the variational means $\mathbb{E}_{q_\alpha}[\theta] = \mu$ are exactly equal to the exact posterior mean $\mathbb{E}_{p_\alpha}[\theta] = m_1$ for all α (even though in general $\text{Cov}_{q_0}(\theta) \neq \Sigma_1$). Consequently, for all α ,

1. MFVB is often used to approximate the posterior when the Bayesian generative model for data x is a mixture model (e.g., Blei et al. (2003)). By contrast, we note for clarity that we are *not* using the mixture model as a generative model for x here. E.g., z is not one of the parameters composing θ , and we are not approximating the distribution of z in the variational distribution $q(\theta)$. Rather, we are using mixtures as a way of flexibly defining skewed and over-dispersed targets, $p(\theta)$.

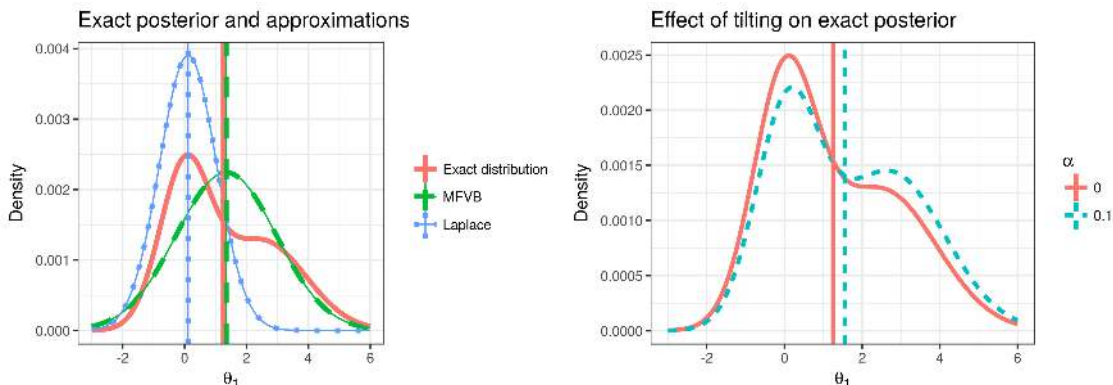


Figure 1: A univariate skewed distribution. Vertical lines show the location of the means.

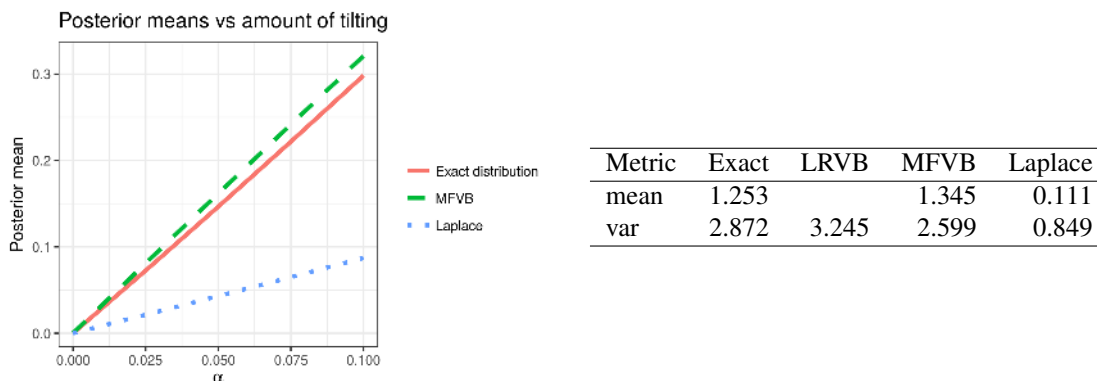


Figure 2: Effect of tilting on a univariate skew distribution.

the variational approximation, the Laplace approximation, and the exact $p_0(\theta)$ all coincide in their estimates of $\mathbb{E}[\theta]$, and by Corollary 2, $\Sigma = \text{Cov}_{p_0}(\theta) = \text{Cov}_{q_0}^{LR}(\theta) = \text{Cov}_{q_{Lap}}^{Lap}(\theta)$. Of course, if Σ is not diagonal, $\text{Cov}_{q_0}(\theta) \neq \Sigma$ because of the mean field assumption. Since this argument holds for the whole vector θ , it holds *a fortiori* for our quantity of interest, the first component $g(\theta) = \theta_1$.

In other words, the Laplace approximation will differ only from the LRVB approximation when $p_0(\theta)$ is not multivariate normal, a situation that we will now bring about by adding new components to the mixture; i.e., by increasing K_z .

5.1.2 A UNIVARIATE SKEWED DISTRIBUTION

If we add a second component ($K_z = 2$), then we can make $p_0(\theta)$ skewed, as shown (with the approximations) in Fig. 1. In this case, we expect $\mathbb{E}_{q_\alpha}[\theta_1]$ to be more accurate than the Laplace approximation $\mathbb{E}_{q_{Lap}}[\theta_1]$ because Q_{mf} is more expressive than Q_{Lap} . This intuition is born out in the left panel of Fig. 1. Since $\hat{\theta}_{Lap}$ uses only information at the mode, it fails to take into account the mass to the right of the mode, and the Laplace approximation’s mean is too far to the left. The MFVB approximation, in contrast, is quite accurate for the posterior mean of θ_1 , even though it gets the overall shape of the distribution wrong.

This example also shows why, in general, one cannot naively form a “Laplace approximation” to the posterior centered at the variational mean rather than at the MAP. As shown in the left panel of Fig. 1, in this case the posterior distribution is actually convex at the MFVB mean. Consequently, a naive second-order approximation to the log posterior centered at the MFVB mean would imply a negative variance.

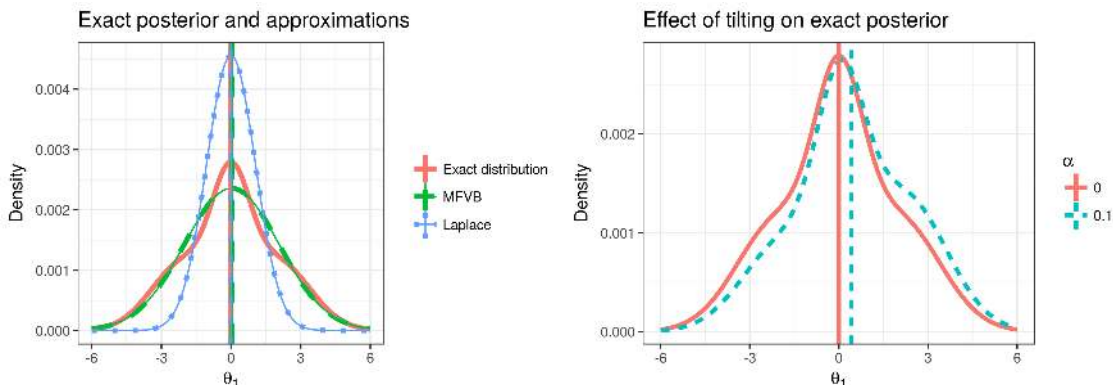


Figure 3: A univariate over-dispersed distribution. Vertical lines show the location of the means.

The perturbation $\rho(\theta, \alpha) = \alpha\theta_1$ is sometimes also described as a “tilting,” and the right panel of Fig. 1 shows the effect of tilting on this posterior approximation. Tilting increases skew, but the MFVB approximation remains accurate, as shown in Fig. 2. Since local sensitivity of the expectation of θ_1 to α is the variance of θ_1 (see Eq. (13)), we have in Fig. 2 that:

- The slope of the exact distribution’s line is $\text{Cov}_{p_0}(\theta_1)$;
- The slope of the MFVB line is the LRVB variance $\text{Cov}_{q_0}^{LR}(\theta_1)$; and
- The slope of the Laplace line is $\text{Cov}_{q_{Lap}}^{Lap}(\theta_1)$.

Since the MFVB and exact lines nearly coincide, we expect the LRVB variance estimate to be quite accurate for this example. Similarly, since the slope of the Laplace approximation line is lower, we expect the Laplace variance to underestimate the exact variance. This outcome, which can be seen visually in the left-hand panel of Fig. 2, is shown quantitatively in the corresponding table in the right-hand panel. The columns of the table contain information for the exact distribution and the three approximations. The first row, labeled “mean,” shows $\mathbb{E}[\theta_1]$ and the second row, labeled “var,” shows $\text{Cov}(\theta_1)$. (The “LRVB” entry for the mean is blank because LRVB differs from MFVB only in covariance estimates.) We conclude that, in this case, Condition 1 holds for \mathcal{Q}_{mf} but not for \mathcal{Q}_{Lap} .

5.1.3 A UNIVARIATE OVER-DISPersed DISTRIBUTION

Having seen how MFVB can outperform the Laplace approximation for a univariate skewed distribution, we now apply that intuition to see why the linear response covariance can be superior to the Laplace approximation covariance for over-dispersed but symmetric distributions. Such a symmetric but over-dispersed distribution, formed with $K_z = 3$ components, is shown in Fig. 3 together with its approximations. By symmetry, both the MFVB and Laplace means are exactly correct (up to Monte Carlo error), as can be seen in the left panel of Fig. 3.

However, the right panel of Fig. 3 shows that symmetry is not maintained as the distribution is tilted. For $\alpha > 0$, the distribution becomes skewed to the right. Thus, by the intuition from the previous section, we expect the MFVB mean to be more accurate as the distribution is tilted and α increases from zero. In particular, we expect that the Laplace approximation’s mean will not shift enough as α varies, i.e., that the Laplace approximation variance will be underestimated. Fig. 4 shows that this is indeed the case. The slopes in the left panel once again correspond to the estimated variances shown in the table, and, as expected the LRVB variance estimate is superior to the Laplace approximation variance.

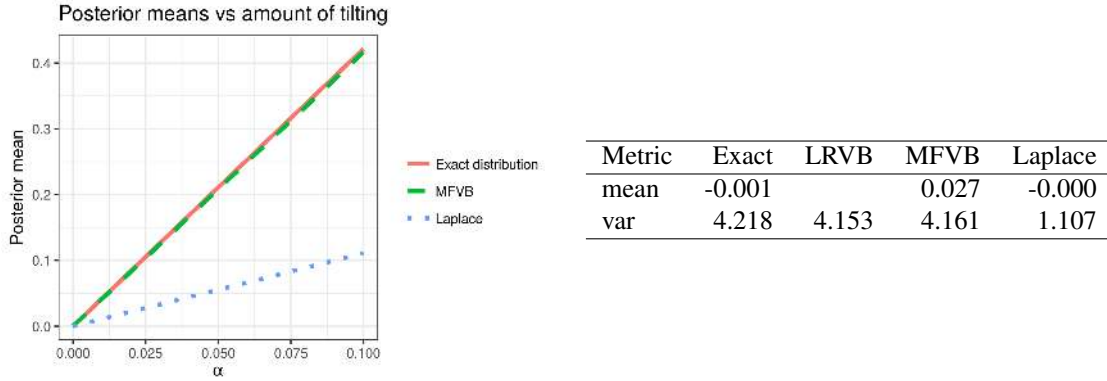


Figure 4: Effect of tilting on a univariate over-dispersed distribution.

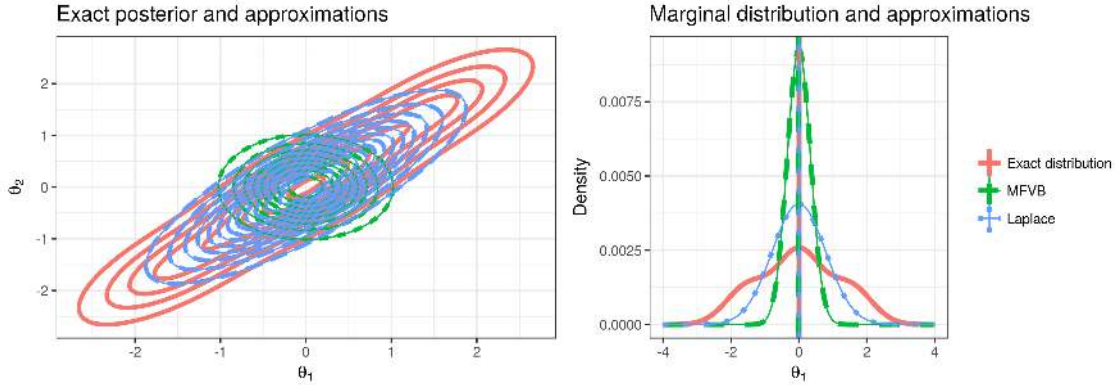


Figure 5: A bivariate over-dispersed distribution.

In this case, Condition 1 holds for \mathcal{Q}_{mf} . For the Laplace approximation, $\mathbb{E}_{q_{Lap}}[g(\theta)] = \mathbb{E}_{p_0}[g(\theta)]$ for $\alpha = 0$, so \mathcal{Q}_{Lap} satisfies Eq. (11) of Condition 1 for α near zero, the derivatives of the two expectations with respect to α are quite different, so Eq. (12) of Condition 1 does not hold for \mathcal{Q}_{Lap} .

5.1.4 A BIVARIATE OVER-DISPersed DISTRIBUTION

In the previous two examples the mean field approximation in \mathcal{Q} did not matter, since the examples were one-dimensional. The only reason that the variational approximation was different from the exact $p_0(\theta)$ was the normal assumption in \mathcal{Q}_{mf} . Indeed, the tables in Fig. 2 and Fig. 4 show that the MFVB variance estimate is also reasonably close to the exact variance. In order to demonstrate why the LRVB variance can be better than both the Laplace approximation and the MFVB approximation, we turn to a bivariate, correlated, over-dispersed $p_0(\theta)$. For this we use $K_z = 3$ correlated normal distributions, shown in the left panel of Fig. 5. The right panel of Fig. 5 shows the marginal distribution of θ_1 , in which the over-dispersion can be seen clearly. As Fig. 5 shows, unlike in the previous two examples, the mean field approximation causes $q_0(\theta)$ to dramatically underestimate the marginal variance of θ_1 . Consequently, the MFVB means will also be under-responsive to the skew introduced by tilting with α . Though the Laplace approximation has a larger marginal variance, it remains unable to take skewness into account. Consequently, as seen in Fig. 6, the LRVB variance, while not exactly equal to the correct variance, is still an improvement over the Laplace covariance, and a marked improvement on the badly under-estimated MFVB variance.

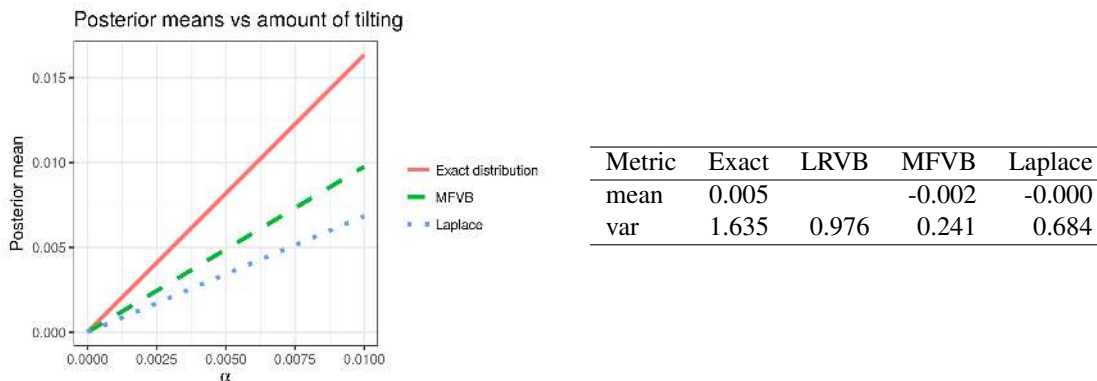


Figure 6: Effect of tilting on a bivariate over-dispersed distribution.

One might say, in this case, that Condition 1 does not hold for either Q_{mf} or Q_{Lap} , or, if it does, it is with a liberal interpretation of the “approximately equals” sign. However, the expressiveness of Q_{mf} allows LRVB to improve on the Laplace approximation, and the linear response allows it to improve over the MFVB approximation, and so LRVB gives the best of both worlds.

Thinking about problems in terms of these three simple models can provide intuition about when and whether Condition 1 might be expected to hold in a sense that is practically useful.

5.2 Automatic Differentiation Variational Inference (ADVI)

In this section we apply our methods to automatic differentiation variational inference (ADVI) (Kucukelbir et al., 2017). ADVI is a “black-box” variational approximation and optimization procedure that requires only that the user provide the log posterior, $\log p_0(\theta)$, up to a constant that does not depend on θ . To achieve this generality, ADVI employs:

- A factorizing normal variational approximation,²
- An unconstraining parameterization,
- The “re-parameterization trick,” and
- Stochastic gradient descent.

ADVI uses a family employing the factorizing normal approximation

$$Q_{ad} := \left\{ q(\theta) : q(\theta) = \prod_{k=1}^K \mathcal{N}(\theta_k | \mu_k, \exp(2\zeta_k)) \right\}.$$

That is, Q_{ad} is a fully factorizing normal family with means μ_k and log standard deviations ζ_k . Because we are including exponential family assumptions in the definition of MFVB (as described in Section 3.1), Q_{ad} is an instance of a mean-field family Q_{mf} . In the notation of Eq. (7),

$$\eta = (\mu_1, \dots, \mu_K, \zeta_1, \dots, \zeta_K)^\top, \quad (22)$$

2. Kucukelbir et al. (2017) describe a non-factorizing version of ADVI, which is called “fullrank” ADVI in Stan. The factorizing version that we describe here is called “meanfield” ADVI in Stan. On the examples we describe, in the current Stan implementation, we found that fullrank ADVI provided much worse approximations to the MCMC posterior means than the meanfield version, and so we do not consider it further.

$\Omega_\eta = \mathbb{R}^{2K}$, λ is the Lebesgue measure, and the objective function Eq. (8) is

$$KL(q(\theta; \eta) || p_0(\theta)) = - \int \mathcal{N}(\theta_k | \mu_k, \exp(2\zeta_k)) \log p_0(\theta) \lambda(d\theta) - \sum_{k=1}^K \zeta_k,$$

where we have used the form of the univariate normal entropy up to a constant.

The unconstraining parameterization is required because the use of a normal variational approximation dictates that the base measure on the parameters $\theta \in \mathbb{R}^K$ be supported on all of \mathbb{R}^K . Although many parameters of interest, such as covariance matrices, are not supported on \mathbb{R}^K , there typically exist differentiable maps from an unconstrained parameterization supported on \mathbb{R}^K to the parameter of interest. Software packages such as Stan automatically provide such transforms for a broad set of parameter types. In our notation, we will take these constraining maps to be the function of interest, $g(\theta)$, and take θ to be unconstrained. Note that, under this convention, the prior $p(\theta|\alpha)$ must be a density in the unconstrained space. In practice (e.g., in the Stan software package), one usually specifies the prior density in the constrained space and converts it to a density $p(\theta|\alpha)$ in the unconstrained space using the determinant of the Jacobian of the constraining transform $g(\cdot)$.

The re-parameterization trick allows easy approximation of derivatives of the (generally intractable) objective $KL(q(\theta; \eta) || p_0(\theta))$. By defining z_k using the change of variable

$$z_k := (\theta_k - \mu_k) / \exp(\zeta_k), \quad (23)$$

$KL(q(\theta; \eta) || p_0(\theta))$ can be re-written as an expectation with respect to a standard normal distribution. We write $\theta = \exp(\zeta) \circ z + \mu$ by using the component-wise Hadamard product \circ . Then

$$KL(q(\theta; \eta) || p_0(\theta)) = -\mathbb{E}_z [\log p_0(\exp(\zeta) \circ z + \mu)] - \sum_{k=1}^K \zeta_k + Constant.$$

The expectation is still typically intractable, but it can be approximated using Monte Carlo and draws from a K -dimensional standard normal distribution. For a fixed number M of draws z_1, \dots, z_M from a standard K -dimensional normal, we can define the approximate KL divergence

$$\widehat{KL}(\eta) := -\frac{1}{M} \sum_{m=1}^M \log p_0(\exp(\zeta) \circ z_m + \mu) - \sum_{k=1}^K \zeta_k + Constant. \quad (24)$$

For any fixed M ,

$$\mathbb{E} \left[\frac{\partial}{\partial \eta} \widehat{KL}(\eta) \right] = \frac{\partial}{\partial \eta} KL(q(\theta; \eta) || p_0(\theta)),$$

so gradients of $\widehat{KL}(\eta)$ are unbiased for gradients of the exact KL divergence. Furthermore, for fixed draws z_1, \dots, z_M , $\widehat{KL}(\eta)$ can be easily differentiated (using, again, the re-parameterization trick). Standard ADVI uses this fact to optimize $KL(q(\theta; \eta) || p_0(\theta))$ using the unbiased gradient draws $\frac{\partial}{\partial \eta} \widehat{KL}(\eta)$ and a stochastic gradient optimization method, where the stochasticity comes from draws of the standard normal random variable z . Note that stochastic gradient methods typically use a new draw of z at every gradient step.

5.2.1 LINEAR RESPONSE FOR ADVI (LR-ADVI)

Since ADVI uses a factorizing normal approximation, the intuition from Section 5.1 may be expected to apply. In particular, we might expect that the ADVI means $\hat{\mu}$ might be a good approximation to $\mathbb{E}_{p_0}[\theta]$, that the ADVI variances $\exp(2\hat{\zeta})$ would be under-estimates of the posterior variance $\text{Cov}_{p_0}(\theta)$, so that using $\text{Cov}_{q_0}^{LR}(\theta)$ could improve the approximations to the posterior variance. We refer to LRVB covariances calculated using an ADVI approximation as LR-ADVI.

To apply linear response to an ADVI approximation, we need to be able to approximate the Hessian of $KL(q(\theta; \eta) || p_0(\theta))$ and to be assured that we have found an optimal η_0^* . But, by using a stochastic gradient method, ADVI avoids ever actually calculating the expectation in $KL(q(\theta; \eta) || p_0(\theta))$. Furthermore even if a stochastic gradient method finds an point that is close to the optimal value of $KL(q(\theta; \eta) || p_0(\theta))$ it may not be close to an optimum of $\widehat{KL}(\eta)$ for a particular finite M . Indeed, we found that, even for very large M , the optimum found by ADVI’s stochastic gradient method is typically not close enough to an optimum of the approximate $\widehat{KL}(\eta)$ for sensitivity calculations to be useful. Sensitivity calculations are based on differentiating the fixed point equation given by the gradient being zero (see the proof in Appendix D), and do not apply at points for which the gradient is not zero either in theory nor in practice.

Consequently, in order to calculate the local sensitivity, we simply eschew the stochastic gradient method and directly optimize $\widehat{KL}(\eta)$ for a particular choice of M . (We will discuss shortly how to choose M .) We can then use $\widehat{KL}(\eta)$ in Eq. (10) rather than the exact KL divergence. Directly optimizing $\widehat{KL}(\eta)$ both frees us to use second-order optimization methods, which we found to converge more quickly to a high-quality optimum than first-order methods, and guarantees that we are evaluating the Hessian $\mathbf{H}_{\eta\eta}$ at an optimum of the objective function used to calculate Eq. (10).

As M approaches infinity, we expect the optimum of $\widehat{KL}(\eta)$ to approach the optimum of $KL(q(\theta; \eta) || p_0(\theta))$ by the standard frequentist theory of estimating equations (Keener, 2010, Chapter 9). In practice we must fix a particular finite M , with larger M providing better approximations of the true KL divergence but at increased computational cost. We can inform this tradeoff between accuracy and computation by considering the frequentist variability of η_0^* when randomly sampling M draws of the random variable z used to approximate the intractable integral in $\widehat{KL}(\eta)$. Denoting this frequentist variability by $\text{Cov}_z(\eta_0^*)$, standard results (Keener, 2010, Chapter 9) give that

$$\text{Cov}_z(\eta_0^*) \approx \mathbf{H}_{\eta\eta}^{-1} \text{Cov}_z \left(\left. \frac{\partial}{\partial \eta} \widehat{KL}(\eta) \right|_{\eta_0^*} \right) \mathbf{H}_{\eta\eta}^{-1}. \quad (25)$$

A sufficiently large M will be one for which $\text{Cov}_z(\eta_0^*)$ is adequately small. One notion of “adequately small” might be that the ADVI means found with $\widehat{KL}(\eta)$ are within some fraction of a posterior standard deviation of the optimum of $KL(q(\theta; \eta) || p_0(\theta))$. Having chosen a particular M , we can calculate the frequentist variability of μ^* using $\text{Cov}_{q_0}^{LR}(g(\theta))$ and estimate the posterior standard deviation using Eq. (14). If we find that each μ^* is probably within 0.5 standard deviations of the optimum of $KL(q(\theta; \eta) || p_0(\theta))$, we can keep the results; otherwise, we increase M and try again. In the examples we consider here, we found that the relatively modest $M = 10$ satisfies this condition and provides sufficiently accurate results.

Finally, we note a minor departure from Eq. (14) when calculating $\text{Cov}_{q_0}^{LR}(g(\theta))$ from $\mathbf{H}_{\eta\eta}$. Recall that, in this case, we are taking $g(\cdot)$ to be ADVI’s constraining transform, and that Eq. (14) requires the Jacobian, \mathbf{g}_η , of this transform. At the time of writing, the design of the Stan software package did not readily support automatic calculation of \mathbf{g}_η , though it did support rapid evaluation of $g(\theta)$ at particular values of θ . Consequently, we used linear response to estimate $\text{Cov}_{q_0}^{LR}(\theta)$, drew a large number N_s of Monte Carlo draws from $\theta_n \sim \mathcal{N}(\mu, \text{Cov}_{q_0}^{LR}(\theta))$ for $n = 1, \dots, N_s$, and then used these draws to form a Monte Carlo estimate of the sample covariance of $g(\theta)$. Noting that $\mathbb{E}_{q_\alpha}[\theta] = \mu$, and recalling the definition of η for ADVI in Eq. (22), by Eq. (14) we have

$$\text{Cov}_{q_0}^{LR}(\theta) = \frac{\partial \mathbb{E}_{q_\alpha}[\theta]}{\partial \eta^\top} \mathbf{H}_{\eta\eta}^{-1} \frac{\partial \mathbb{E}_{q_\alpha}[\theta^\top]}{\partial \eta} = \begin{pmatrix} I_K & 0 \\ 0 & 0 \end{pmatrix} \mathbf{H}_{\eta\eta}^{-1} \begin{pmatrix} I_K & 0 \\ 0 & 0 \end{pmatrix},$$

which is the upper-left quarter of the matrix $\mathbf{H}_{\eta\eta}^{-1}$. In addition to obviating the need for \mathbf{g}_η , this approach also allowed us to take into account possible nonlinearities in $g(\cdot)$ at little additional computational cost.

5.2.2 RESULTS

We present results from four models taken from the Stan example set, namely the models `election88` (“Election model”), `sesame_street1` (“Sesame Street model”), `radon_vary_intercept_floor` (“Radon

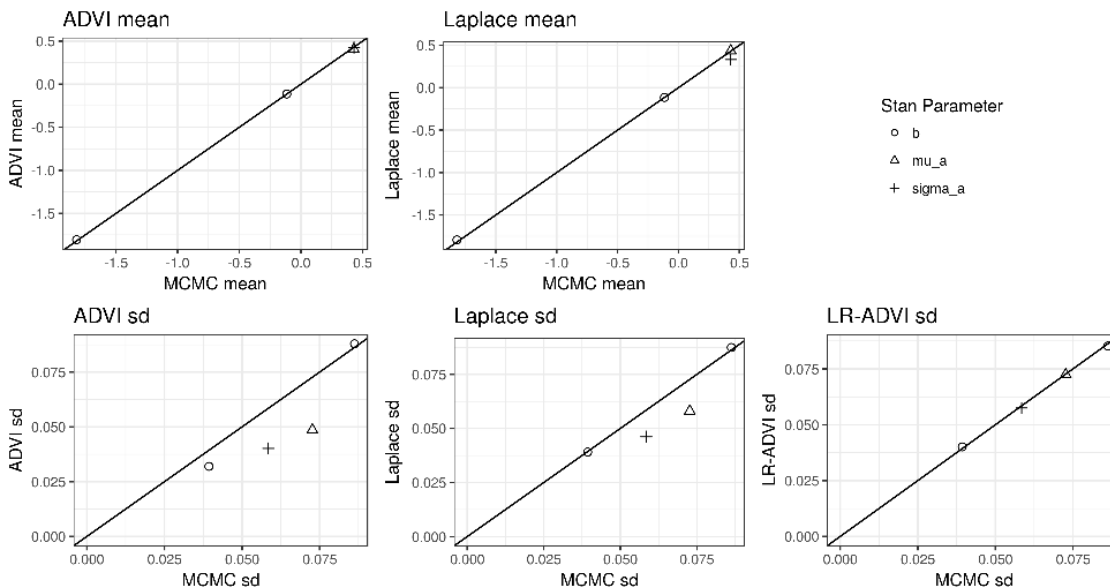


Figure 7: Election model

model”), and `cjs_cov_ranef` (“Ecology model”). We experimented with many models from the Stan examples and selected these four as representative of the type of model where LR-ADVI can be expected to provide a benefit—specifically, they are models of a moderate size. For very small models, MCMC runs quickly enough in Stan that fast approximations are not necessary, and for very large models (with thousands of parameters) the relative advantages of LR-ADVI and the Laplace approximation diminish due to the need to calculate $\mathbf{H}_{\eta\eta}$ or \mathbf{H}_{Lap} using automatic differentiation.³ The size of the data and size of the parameter space for our four chosen models are shown in Fig. 11. We also eliminated from consideration models where Stan’s MCMC algorithm reported divergent transitions or where Stan’s ADVI algorithm returned wildly inaccurate posterior mean estimates.

For brevity, we do not attempt to describe the models or data in any detail here; rather, we point to the relevant literature in their respective sections. The data and Stan implementations themselves can be found on the Stan website (Stan Team, 2017) as well as in Appendix F.

To assess the accuracy of each model, we report means and standard deviations for each of Stan’s model parameters as calculated by Stan’s MCMC and ADVI algorithms and a Laplace approximation, and we report the standard deviations as calculated by $\text{Cov}_{q_0}^{LR}(g(\theta))$. Recall that, in our notation, $g(\cdot)$ is the (generally nonlinear) map from the unconstrained latent ADVI parameters to the constrained space of the parameters of interest. The performance of ADVI and Laplace vary, and only LR-ADVI provides a consistently good approximation to the MCMC standard deviations. LR-ADVI was somewhat slower than a Laplace approximation or ADVI alone, but it was typically about five times faster than MCMC; see Section 5.2.7 for detailed timing results.

5.2.3 ELECTION MODEL ACCURACY

We begin with `election88`, which models binary responses in a 1988 poll using a Bernoulli hierarchical model with normally distributed random effects for state, ethnicity, and gender and a logit link. The model and data are described in detail in Gelman and Hill (2006, Chapter 14). Fig. 7 shows that both the Laplace

3. We calculated $\mathbf{H}_{\eta\eta}$ using a custom branch of Stan’s automatic differentiation software (Carpenter et al., 2015) that exposes Hessians and Hessian-vector products in the `Rstan::model_fit` class. When this custom branch is merged with the main branch of Stan, it will be possible to implement LR-ADVI for generic Stan models.

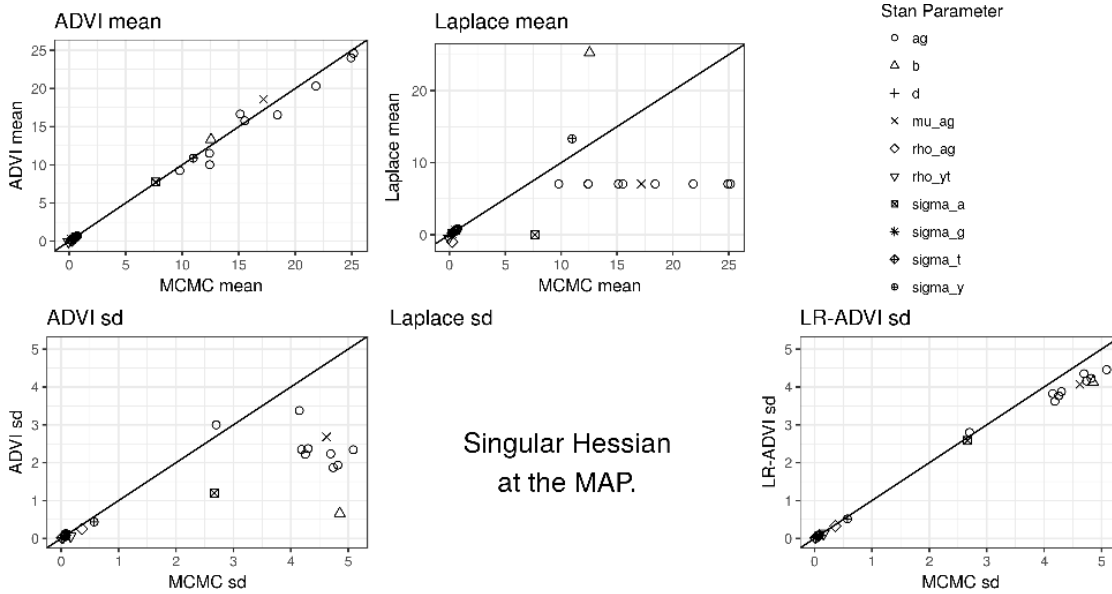


Figure 8: Sesame Street model

approximation and ADVI do a reasonable job of matching to MCMC, though LR-ADVI is slightly more accurate for standard deviations.

5.2.4 SESAME STREET MODEL ACCURACY

Next, we show results for `sesame_street1`, an analysis of a randomized controlled trial designed to estimate the causal effect of watching the television show Sesame Street on a letter-recognition test. To control for different conditions in the trials, a hierarchical model is used with correlated multivariate outcomes and unknown covariance structure. The model and data are described in detail in Gelman and Hill (2006, Chapter 23).

As can be seen in Fig. 8, the MAP under-estimates the variability of the random effects ag , and, in turn, under-estimates the variance parameter σ_a . Because the MAP estimate of σ_a is close to zero, the log posterior has a very high curvature with respect to the parameter ag at the MAP, and the Hessian used for the Laplace approximation is numerically singular. ADVI, which integrates out the uncertainty in the random effects, provides reasonably good estimates of the posterior means but underestimates the posterior standard deviations due to the mean-field assumption. Only LR-ADVI provides accurate estimates of posterior uncertainty.

5.2.5 RADON MODEL ACCURACY

We now turn to `radon_vary_intercept_floor`, a hierarchical model of radon levels in Minnesota homes described in Gelman and Hill (2006, Chapters 16 and 21). This model is relatively simple, with univariate normal observations and unknown variances. Nevertheless, the Laplace approximation again produces a numerically singular covariance matrix. The ADVI means are reasonably accurate, but the standard deviations are not. Only LR-ADVI produces an accurate approximation to the MCMC posterior standard deviations.

5.2.6 ECOLOGY MODEL ACCURACY

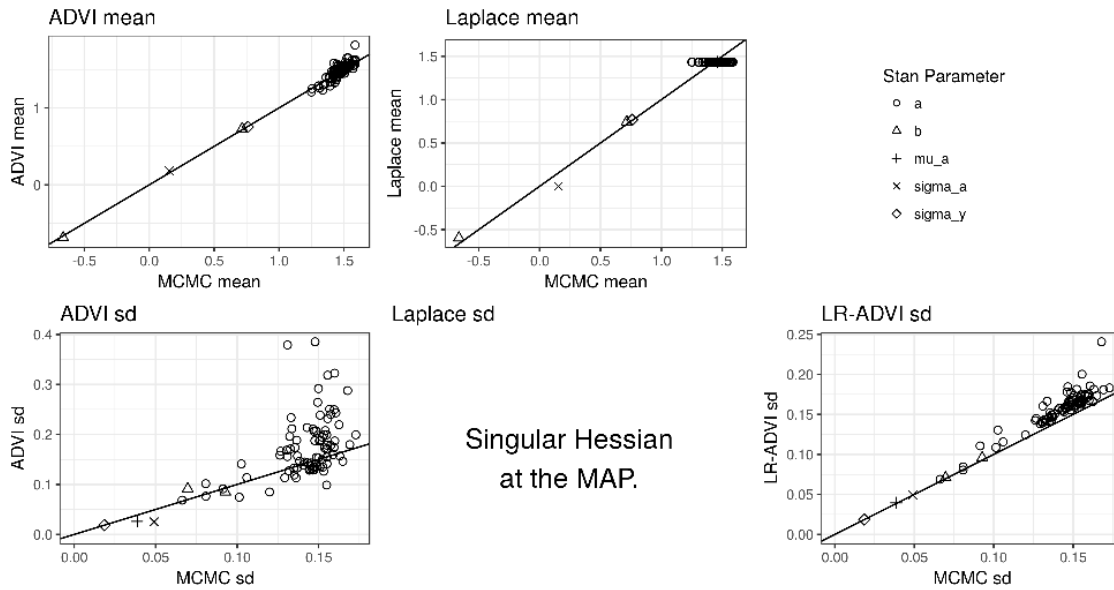


Figure 9: Radon model

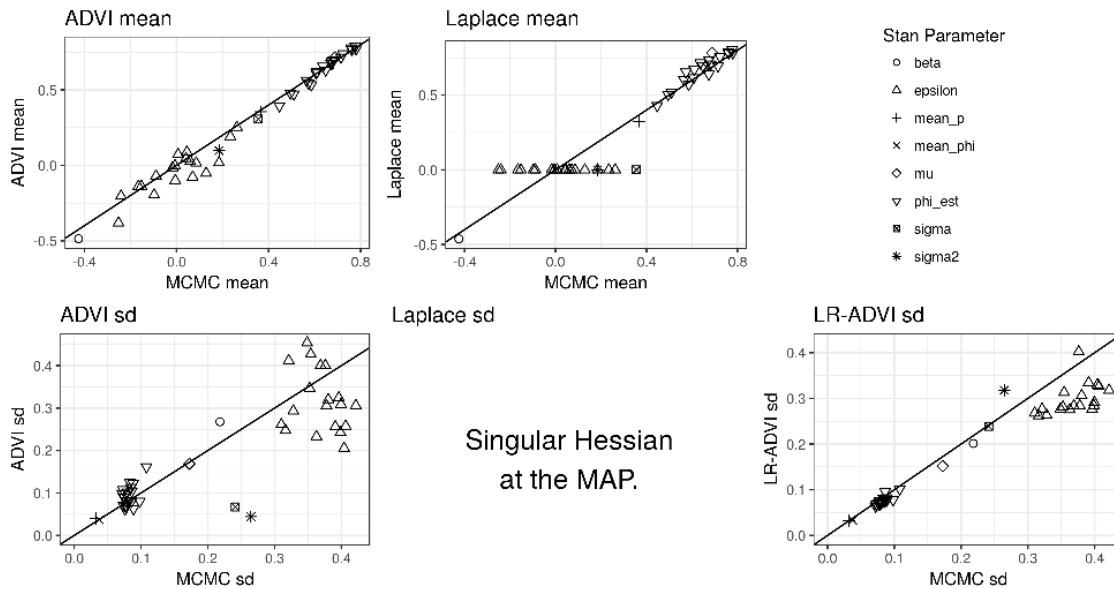


Figure 10: Ecology model

Finally, we consider a more complicated mark-recapture model from ecology known as the Cormack-Jolly-Seber (CJS) model. This model is described in detail in Kéry and Schaub (2011, Chapter 7), and discussion of the Stan implementation can be found in Stan Team (2015, Section 15.3).

The Laplace approximation is again degenerate, and the ADVI standard deviations again deviate considerably from MCMC. In this case, the ADVI means are also somewhat inaccurate, and some of the LR-ADVI standard deviations are mis-estimated in turn. However, LR-ADVI remains by far the most accurate method for approximating the MCMC standard errors.

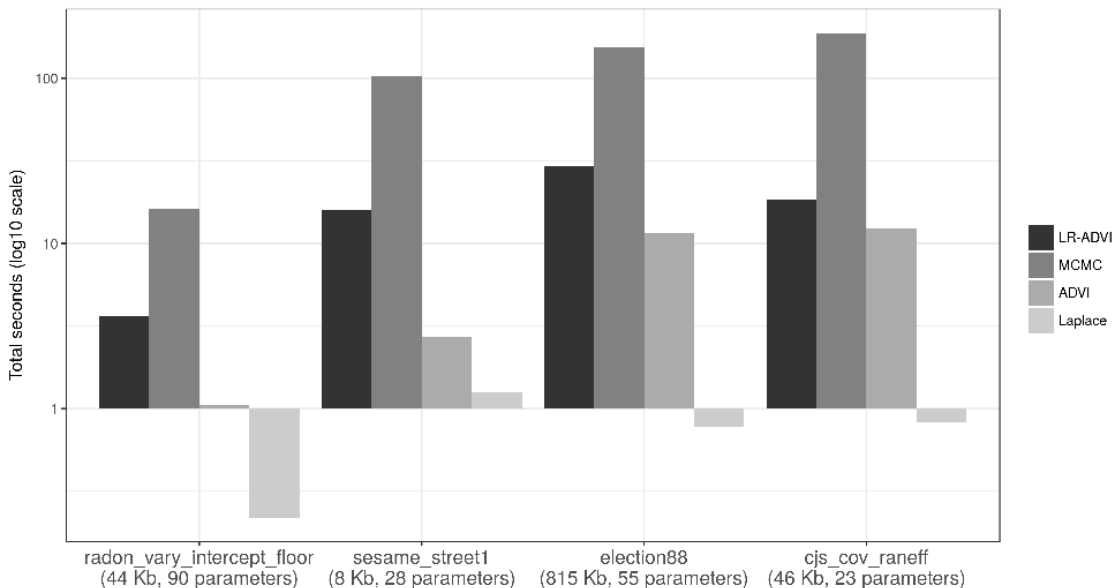


Figure 11: Comparison of timing in ADVI experiments

5.2.7 TIMING RESULTS

Detailed timing results for the ADVI experiments are shown in Fig. 11. Both the Laplace approximation and ADVI alone are faster than LR-ADVI, which in turn is about five times faster than MCMC. We achieved the best results optimizing $\widehat{KL}(\eta)$ by using the conjugate gradient Newton’s trust region method (`trust-ncg` of `scipy.optimize`), but the optimization procedure still accounted for an appreciable proportion of the time needed for LR-ADVI.

5.3 Criteo Dataset

We now apply our methods to a real-world data set using a logistic regression with random effects, which is an example of a generalized linear mixed model (GLMM) (Agresti and Kateri, 2011, Chapter 13). This data and model have several advantages as an illustration of our methods: the data set is large, the model contains a large number of imprecisely-estimated latent variables (the unknown random effects), the model exhibits the sparsity of $\mathbf{H}_{\eta\eta}$ that is typical in many MFVB applications, and the results exhibit the same shortcomings of the Laplace approximation seen above. For this model, we will evaluate both posterior covariances and prior sensitivities.

5.3.1 DATA AND MODEL

We investigated a custom subsample of the 2014 Criteo Labs conversion logs data set (Criteo Labs, 2014), which contains an obfuscated sample of advertising data collected by Criteo over a period of two months. Each row of the data set corresponds to a single user click on an online advertisement. For each click, the data set records a binary outcome variable representing whether or not the user subsequently “converted” (i.e., performed a desired task, such as purchasing a product or signing up for a mailing list). Each row contains two timestamps (which we ignore), eight numerical covariates, and nine factor-valued covariates. Of the eight numerical covariates, three contain 30% or more missing data, so we discarded them. We then applied a per-covariate normalizing transform to the distinct values of those remaining. Among the factor-valued covariates, we retained only the one with the largest number of unique values and discarded the others.

These data-cleaning decisions were made for convenience. The goal of the present paper is to demonstrate our inference methods, not to draw conclusions about online advertising.

Although the meaning of the covariates has been obfuscated, for the purpose of discussion we will imagine that the single retained factor-valued covariate represents the identity of the advertiser, and the numeric covariates represent salient features of the user and/or the advertiser (e.g., how often the user has clicked or converted in the past, a machine learning rating for the advertisement quality, etc.). As such, it makes sense to model the probability of each row’s binary outcome (whether or not the user converted) as a function of the five numeric covariates and the advertiser identity using a logistic GLMM. Specifically, we observe binary conversion outcomes, y_{it} , for click i on advertiser t , with probabilities given by observed numerical explanatory variables, x_{it} , each of which are vectors of length $K_x = 5$. Additionally, the outcomes within a given value of t are correlated through an unobserved random effect, u_t , which represents the “quality” of advertiser t , where the value of t for each observation is given by the factor-valued covariate. The random effects u_t are assumed to follow a normal distribution with unknown mean and variance. Formally,

$$\begin{aligned} y_{it}|p_{it} &\sim \text{Bernoulli}(p_{it}), \text{ for } t = 1, \dots, T \text{ and } i = 1, \dots, N_t \\ p_{it} &:= \frac{e^{\rho_{it}}}{1 + e^{\rho_{it}}} \quad \text{where } \rho_{it} := x_{it}^T \beta + u_t \\ u_t|\mu, \tau &\sim \mathcal{N}(\mu, \tau^{-1}). \end{aligned}$$

Consequently, the unknown parameters are $\theta = (\beta^T, \mu, \tau, u_1, \dots, u_T)^T$. We use the following priors:

$$\begin{aligned} \mu|\mu_0, \tau_\mu &\sim \mathcal{N}(\mu_0, \tau_\mu^{-1}) \\ \tau|\alpha_\tau, \beta_\tau &\sim \text{Gamma}(\alpha_\tau, \beta_\tau) \\ \beta|\beta_0, \tau_\beta, \gamma_\beta &\sim \mathcal{N}\left(\begin{pmatrix} \beta_0 \\ \vdots \\ \beta_0 \end{pmatrix}, \begin{pmatrix} \tau_\beta & \gamma_\beta & \gamma_\beta \\ \gamma_\beta & \ddots & \gamma_\beta \\ \gamma_\beta & \gamma_\beta & \tau_\beta \end{pmatrix}^{-1}\right). \end{aligned}$$

Note that we initially take $\gamma_\beta = 0$ so that the prior information matrix on β is diagonal. Nevertheless, by retaining γ_β as a hyperparameter we will be able to assess the sensitivity to the assumption of a diagonal prior in Section 5.3.6. The remaining prior values are given in Appendix H. It is reasonable to expect that a modeler would be interested both in the effect of the numerical covariates and in the quality of individual advertisers themselves, so we take the parameter of interest to be $g(\theta) = (\beta^T, u_1, \dots, u_T)^T$.

To produce a data set small enough to be amenable to MCMC but large and sparse enough to demonstrate our methods, we subsampled the data still further. We randomly chose 5000 distinct advertisers to analyze, and then subsampled each selected advertiser to contain no more than 20 rows each. The resulting data set had $N = 61895$ total rows. If we had more observations per advertiser, the “random effects” u_t would have been estimated quite precisely, and the nonlinear nature of the problem would not have been important; these changes would thus have obscured the benefits of using MFVB versus the Laplace approximation. In typical internet data sets a large amount of data comes from advertisers with few observations each, so our subsample is representative of practically interesting problems.

5.3.2 INFERENCE AND TIMING

We estimated the expectation and covariance of $g(\theta)$ using four techniques: MCMC, the Laplace approximation, MFVB, and linear response (LRVB) methods. For MCMC, we used Stan (Stan Team, 2015), and to calculate the MFVB, Laplace, and LRVB estimates we used our own Python code using `numpy`, `scipy`, and `autograd` (Jones et al., 2001; Maclaurin et al., 2015). As described in Section 5.3.3, the MAP estimator did not estimate $\mathbb{E}_{p_\theta}[g(\theta)]$ very well, so we do not report standard deviations or sensitivity measures for the Laplace approximations. The summary of the computation time for all these methods is shown in Table 1, with details below.

Method	Seconds
MAP (optimum only)	12
VB (optimum only)	57
VB (including sensitivity for β)	104
VB (including sensitivity for β and u)	553
MCMC (Stan)	21066

Table 1: Timing results

For the MCMC estimates, we used Stan to draw 5000 MCMC draws (not including warm-up), which took 351 minutes. We estimated all the prior sensitivities of Section 5.3.6 using the Monte Carlo version of the covariance in Eq. (5).

For the MFVB approximation, we use the following mean field exponential family approximations:

$$\begin{aligned}
q(\beta_k) &= \mathcal{N}(\beta_k; \eta_{\beta_k}), \text{ for } k = 1, \dots, K_x \\
q(u_t) &= \mathcal{N}(u_t; \eta_{u_t}), \text{ for } t = 1, \dots, T \\
q(\tau) &= \text{Gamma}(\tau; \eta_\tau) \\
q(\mu) &= \mathcal{N}(\mu; \eta_\mu) \\
q(\theta) &= q(\tau) q(\mu) \prod_{k=1}^{K_x} q(\beta_k) \prod_{t=1}^T q(u_t).
\end{aligned}$$

With these choices, evaluating the variational objective requires the following intractable univariate variational expectation:

$$\mathbb{E}_{q(\theta; \eta)} [\log(1 - p_{it})] = \mathbb{E}_{q(\theta; \eta)} \left[\log \left(1 - \frac{e^{\rho_{it}}}{1 + e^{\rho_{it}}} \right) \right].$$

We used the re-parameterization trick and four points of Gauss-Hermite quadrature to estimate this integral for each observation. See Appendix H for more details.

We optimized the variational objective using the conjugate gradient Newton’s trust region method, `trust-ncg`, of `scipy.optimize`. One advantage of `trust-ncg` is that it performs second-order optimization but requires only Hessian-vector products, which can be computed quickly by `autograd` without constructing the full Hessian. The MFVB fit took 57 seconds, roughly 370 times faster than MCMC with Stan.

With variational parameters for each random effect u_t , $\mathbf{H}_{\eta\eta}$ is a 10014×10014 dimensional matrix. Consequently, evaluating $\mathbf{H}_{\eta\eta}$ directly as a dense matrix using `autograd` would have been prohibitively time-consuming. Fortunately, our model can be decomposed into global and local parameters, and the Hessian term $\mathbf{H}_{\eta\eta}$ in Theorem 2 is extremely sparse. In the notation of Section 4.5, take $\theta_{glob} = (\beta^\top, \mu, \tau)^\top$, take $\theta_{loc,t} = u_t$, and stack the variational parameters as $\eta = (\eta_{glob}^\top, \eta_{loc,1}, \dots, \eta_{loc,T})^\top$. The cross terms in $\mathbf{H}_{\eta\eta}$ between the local variables vanish:

$$\frac{\partial^2 KL(q(\theta; \eta) || p_\alpha(\theta))}{\partial \eta_{loc,t_1} \partial \eta_{loc,t_2}} = 0 \text{ for all } t_1 \neq t_2.$$

Equivalently, note that the full likelihood in Appendix H, Eq. (31), has no cross terms between u_{t_1} and u_{t_2} for $t_1 \neq t_2$. As the dimension T of the data grows, so does the length of η . However, the dimension of η_{glob} remains constant, and $\mathbf{H}_{\eta\eta}$ remains easy to invert. We show an example of the sparsity pattern of the first few rows and columns of $\mathbf{H}_{\eta\eta}$ in Fig. 12.

Taking advantage of this sparsity pattern, we used `autograd` to calculate the Hessian of the KL divergence one group at a time and assembled the results in a sparse matrix using the `scipy.sparse` Python

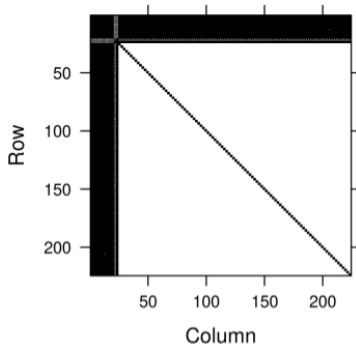


Figure 12: Sparsity pattern of top-left sub-matrix of $\mathbf{H}_{\eta\eta}$ for the logit GLMM model. The axis numbers represent indices within η , and black indicates non-zero entries of $\mathbf{H}_{\eta\eta}$.

package. Even so, calculating the entire sparse Hessian took 323 seconds, and solving the system $\mathbf{H}_{\eta\eta}^{-1}\mathbf{g}_{\eta}^{\top}$ using `scipy.sparse.linalg.spsolve` took an additional 173 seconds. These results show that the evaluation and inversion of $\mathbf{H}_{\eta\eta}$ was several times more costly than optimizing the variational objective itself. (Of course, the whole procedure remains much faster than running MCMC with Stan.)

We note, however, that instead of the direct approach to calculating $\mathbf{H}_{\eta\eta}^{-1}\mathbf{g}_{\eta}^{\top}$ one can use the conjugate gradient algorithm of `sp.sparse.linalg.cg` (Wright and Nocedal, 1999, Chapter 5) together with the fast Hessian-vector products of `autograd` to query one column at a time of $\mathbf{H}_{\eta\eta}^{-1}\mathbf{g}_{\eta}^{\top}$. On a typical column of $\mathbf{H}_{\eta\eta}^{-1}\mathbf{g}_{\eta}^{\top}$ in our experiment, calculating the conjugate gradient took only 9.4 seconds (corresponding to 81 Hessian-vector products in the conjugate gradient algorithm). Thus, for example, one could calculate the columns of $\mathbf{H}_{\eta\eta}^{-1}\mathbf{g}_{\eta}^{\top}$ corresponding to the expectations of the global variables β in only $9.4 \times K_x = 46.9$ seconds, which is much less time than it would take to compute the entire $\mathbf{H}_{\eta\eta}^{-1}\mathbf{g}_{\eta}^{\top}$ for both β and every random effect in u .

For the Laplace approximation, we calculated the MAP estimator and \mathbf{H}_{Lap} using Python code similar to that used for the MFVB estimates. We observe that the MFVB approximation to posterior means would be expected to improve on the MAP estimator only in cases when there is both substantial uncertainty in some parameters and when this uncertainty, through nonlinear dependence between parameters, affects the values of posterior means. These circumstances obtain in the logistic GLMM model with sparse per-advertiser data since the random effects u_t will be quite uncertain and the other posterior means depend on them through the nonlinear logistic function.

5.3.3 POSTERIOR APPROXIMATION RESULTS

In this section, we assess the accuracy of the MFVB, Laplace, and LRVB methods as approximations to $\mathbb{E}_{p_0}[g(\theta)]$ and $\text{Cov}_{p_0}(g(\theta))$. We take the MCMC estimates as ground truth. Although, as discussed in Section 5.3, we are principally interested in the parameters $g(\theta) = (\beta^{\top}, u_1, \dots, u_T)^{\top}$, we will report the results for all parameters for completeness. For readability, the tables and graphs show results for a random selection of the components of the random effects u .

5.3.4 POSTERIOR MEANS

We begin by comparing the posterior means in Table 2, Fig. 13, and Fig. 14. We first note that, despite the long running time for MCMC, the β_1 and μ parameters did not mix well in the MCMC sample, as is reflected in the MCMC standard error and effective number of draws columns of Table 2. The x_{it} data corresponding to β_1 contained fewer distinct values than the other columns of x , which perhaps led to some co-linearity between β_1 and μ in the posterior. This co-linearity could have caused both poor MCMC mixing

Parameter	MCMC	MFVB	MAP	MCMC std. err.	Eff. # of MCMC draws
β_1	1.454	1.447	1.899	0.02067	33
β_2	0.031	0.033	0.198	0.00025	5000
β_3	0.110	0.110	0.103	0.00028	5000
β_4	-0.172	-0.173	-0.173	0.00016	5000
β_5	0.273	0.273	0.280	0.00042	5000
μ	2.041	2.041	3.701	0.04208	28
τ	0.892	0.823	827.724	0.00051	1232
u_{1431}	1.752	1.757	3.700	0.00937	5000
u_{4150}	1.217	1.240	3.699	0.01022	5000
u_{4575}	2.427	2.413	3.702	0.00936	5000
u_{4685}	3.650	3.633	3.706	0.00862	5000

Table 2: Results for the estimation of the posterior means

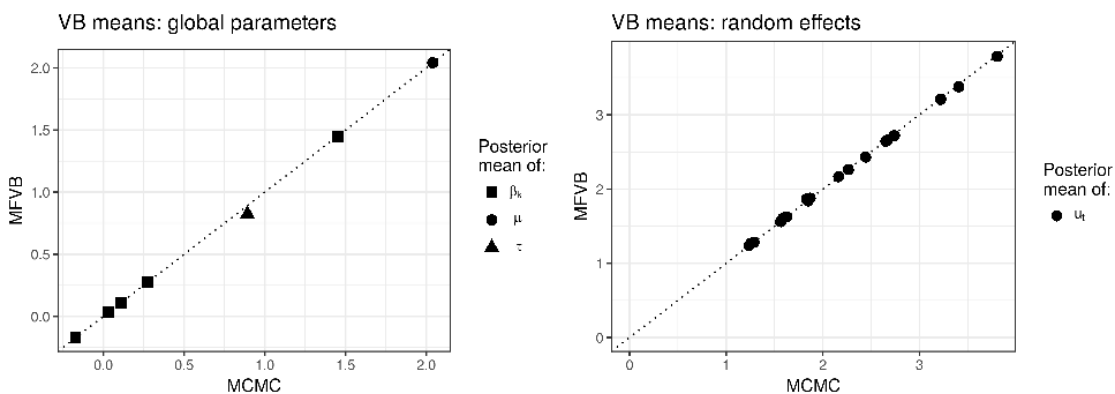


Figure 13: Comparison of MCMC and MFVB means

and, perhaps, excessive measured prior sensitivity, as discussed below in Section 5.3.6. Although we will report the results for both β_1 and μ without further comment, the reader should bear in mind that the MCMC “ground truth” for these two parameters is somewhat suspect.

The results in Table 2 and Fig. 13 show that MFVB does an excellent job of approximating the posterior means in this particular case, even for the random effects u and the related parameters μ and τ . In contrast, the MAP estimator does reasonably well only for certain components of β and does extremely poorly for the random effects parameters. As can be seen in Fig. 14, the MAP estimate dramatically overestimates the information τ of the random effect distribution (that is, it underestimates the variance). As a consequence, it estimates all the random effects to have essentially the same value, leading to mis-estimation of some location parameters, including both μ and some components of β . Since the MAP estimator performed so poorly at estimating the random effect means, we will not consider it any further.

5.3.5 POSTERIOR COVARIANCES

We now assess the accuracy of our estimates of $\text{Cov}_{p_0}(g(\theta))$. The results for the marginal standard deviations are shown in Table 3 and Fig. 15. We refer to the standard deviations of $\text{Cov}_{q_0}(g(\theta))$ as the “uncorrected MFVB” estimate, and of $\text{Cov}_{q_0}^{LR}(g(\theta))$ as the “LRVB” estimate. The uncorrected MFVB variance estimates of β are particularly inaccurate, but the LRVB variances match the exact posterior closely.

In Fig. 16, we compare the off-diagonal elements of $\text{Cov}_{q_0}^{LR}(g(\theta))$ and $\text{Cov}_{p_0}(g(\theta))$. These covariances are zero, by definition, in the uncorrected MFVB estimates $\text{Cov}_{q_0}(g(\theta))$. The left panel of Fig. 16 shows the

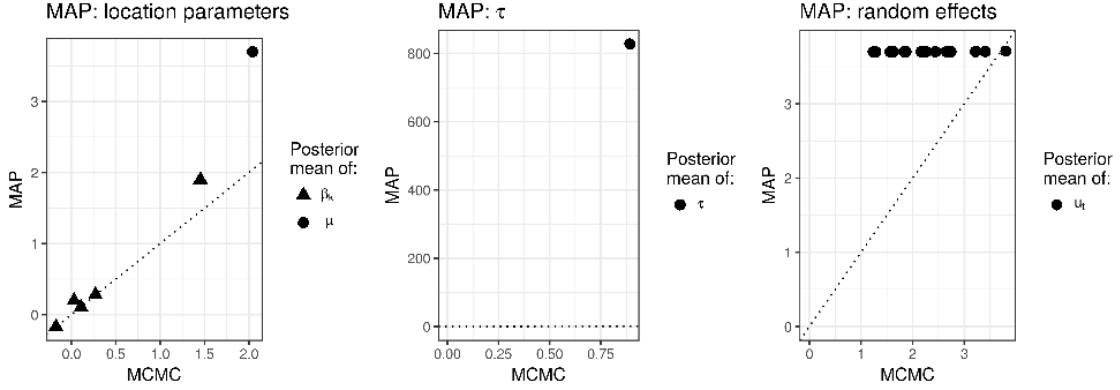


Figure 14: Comparison of MCMC and Laplace means

Parameter	MCMC	LRVB	Uncorrected MFVB
β_1	0.118	0.103	0.005
β_2	0.018	0.018	0.004
β_3	0.020	0.020	0.004
β_4	0.012	0.012	0.004
β_5	0.029	0.030	0.004
μ	0.223	0.192	0.016
τ	0.018	0.033	0.016
u_{1431}	0.663	0.649	0.605
u_{4150}	0.723	0.707	0.662
u_{4575}	0.662	0.649	0.615
u_{4685}	0.610	0.607	0.579

Table 3: Standard deviation results

estimated covariances between the global parameters and all other parameters, including the random effects, and the right panel shows only the covariances amongst the random effects. The LRVB covariances are quite accurate, particularly when we recall that the MCMC draws of μ may be inaccurate due to poor mixing.

5.3.6 PARAMETRIC SENSITIVITY RESULTS

Finally, we compare the MFVB prior sensitivity measures of Section 4.4 to the covariance-based MCMC sensitivity measures of Section 2.1. Since sensitivity is of practical interest only when it is of comparable order to the posterior uncertainty, we report sensitivities normalized by the appropriate standard deviation. That is, we report $\hat{S}_{\alpha_0} / \sqrt{\text{diag}(\hat{Cov}_{p_0}(g(\theta)))}$, and $S_{\alpha_0}^q / \sqrt{\text{diag}(Cov_{q_0}^{LR}(g(\theta)))}$, etc., where $\text{diag}(\cdot)$ denotes the diagonal vector of a matrix, and the division is element-wise. Note that we use the sensitivity-based variance estimates $Cov_{q_0}^{LR}$, not the uncorrected MFVB estimates Cov_{q_0} , to normalize the variational sensitivities. We refer to a sensitivity divided by a standard deviation as a “normalized” sensitivity.

The comparison between the MCMC and MFVB sensitivity measures is shown in Fig. 17. The MFVB and MCMC sensitivities correspond very closely, though the MFVB means appear to be slightly more sensitive to the prior parameters than the MCMC means. This close correspondence should not be surprising. As shown in Section 5.3.3, the MFVB and MCMC posterior means match quite closely. If we assume, reasonably, that they continue to match to first order in a neighborhood of our original prior parameters, then Condition 1 will hold and we would expect $\hat{S}_{\alpha_0} \approx S_{\alpha_0}^q$.

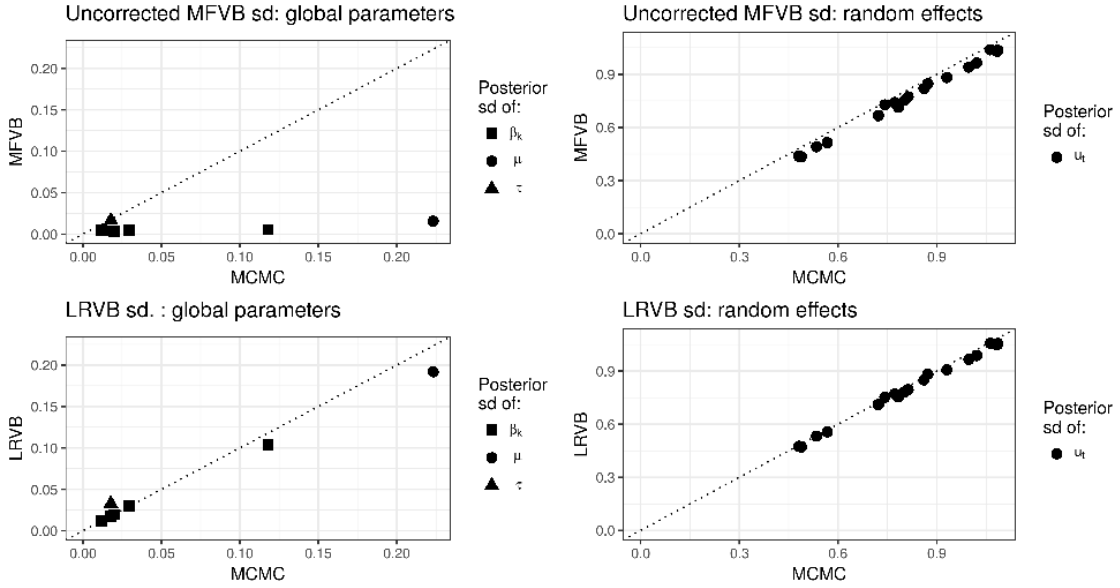


Figure 15: Comparison of MCMC, MFVB, and LRVB standard deviations

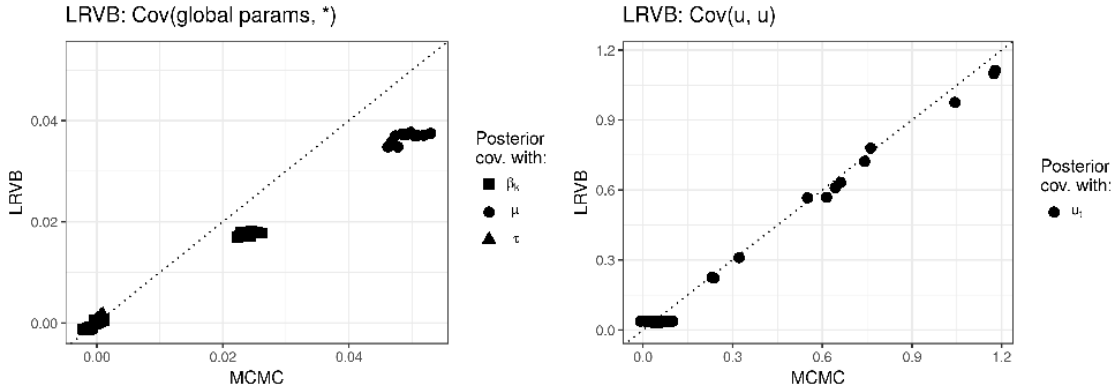


Figure 16: Comparison of MCMC and LRVB off-diagonal covariances

Table 4 shows the detailed MFVB normalized sensitivity results. Each entry is the sensitivity of the MFVB mean of the row’s parameter to the column’s prior parameter. One can see that several parameters are quite sensitive to the information parameter prior τ_μ . In particular, $\mathbb{E}_{p_\alpha}[\mu]$ and $\mathbb{E}_{p_\alpha}[\beta_1]$ are expected to change approximately -0.39 and -0.35 standard deviations, respectively, for every unit change in τ_μ . This size of change could be practically significant (assuming that such a change in τ_μ is subjectively plausible). To investigate this sensitivity further, we re-fit the MFVB model at a range of values of the prior parameter τ_μ , assessing the accuracy of the linear approximation to the sensitivity. The results are shown in Fig. 18. Even for very large changes in τ_μ —resulting in changes to $\mathbb{E}_{p_\alpha}[\mu]$ and $\mathbb{E}_{p_\alpha}[\beta_1]$ far in excess of two standard deviations—the linear approximation holds up reasonably well. Fig. 18 also shows a (randomly selected) random effect to be quite sensitive, though not to a practically important degree relative to its posterior standard deviation. The insensitivity of $\mathbb{E}_{p_\alpha}[\beta_2]$ is also confirmed. Of course, the accuracy of the linear approximation cannot be guaranteed to hold as well in general as it does in this particular case, and the quick

	β_0	τ_β	γ_β	μ_0	τ_μ	α_τ	β_τ
μ	0.0094	-0.1333	-0.0510	0.0019	-0.3920	0.0058	-0.0048
τ	0.0009	-0.0086	-0.0142	0.0003	-0.0575	0.0398	-0.0328
β_1	0.0089	-0.1464	-0.0095	0.0017	-0.3503	0.0022	-0.0018
β_2	0.0012	-0.0143	-0.0113	0.0003	-0.0516	0.0062	-0.0051
β_3	-0.0035	0.0627	-0.0081	-0.0006	0.1218	-0.0003	0.0002
β_4	0.0018	-0.0037	-0.0540	0.0004	-0.0835	0.0002	-0.0002
β_5	0.0002	0.0308	-0.0695	0.0002	-0.0383	0.0011	-0.0009
u_{1431}	0.0028	-0.0397	-0.0159	0.0006	-0.1169	0.0018	-0.0015
u_{4150}	0.0026	-0.0368	-0.0146	0.0005	-0.1083	0.0022	-0.0018
u_{4575}	0.0028	-0.0406	-0.0138	0.0006	-0.1153	0.0011	-0.0009
u_{4685}	0.0028	-0.0409	-0.0142	0.0006	-0.1163	0.0003	-0.0002

Table 4: MFVB normalized prior sensitivity results

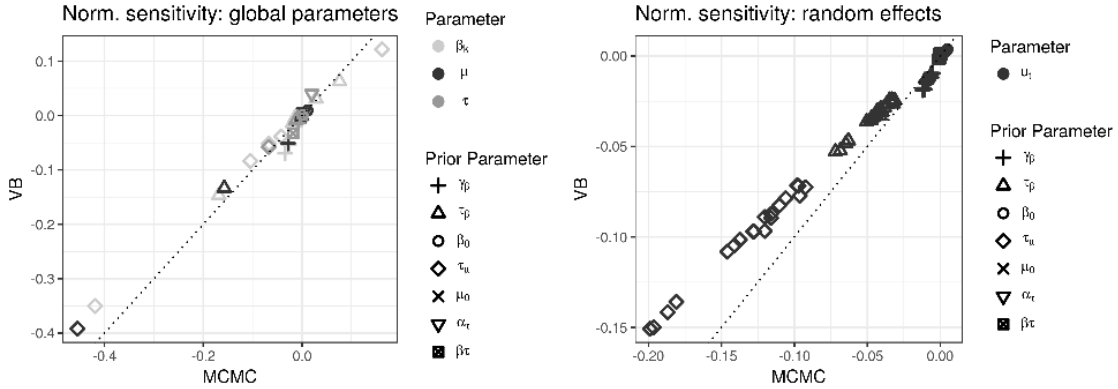


Figure 17: Comparison of MCMC and MFVB normalized parametric sensitivity results

and reliable evaluation of the linearity assumption without re-fitting the model remains interesting future work.

Since we started the MFVB optimization close to the new, perturbed optimum, each new MFVB fit took only 27.2 seconds on average. Re-estimating the MCMC posterior so many times would have been extremely time-consuming. (Note that importance sampling would be useless for prior parameter changes that moved the posterior so far from the original draws.) The considerable sensitivity of this model to a particular prior parameter, which is perhaps surprising on such a large data set, illustrates the value of having fast, general tools for discovering and evaluating prior sensitivity. Our framework provides just such a set of tools.

6. Conclusion

By calculating the sensitivity of MFVB posterior means to model perturbations, we are able to provide two important practical tools for MFVB posterior approximations: improved variance estimates and measures of prior robustness. When MFVB models are implemented in software that supports automatic differentiation, our methods are fast, scalable, and require little additional coding beyond the MFVB objective itself. In our experiments, we were able to calculate accurate posterior means, covariances, and prior sensitivity measures orders of magnitude more quickly than MCMC.

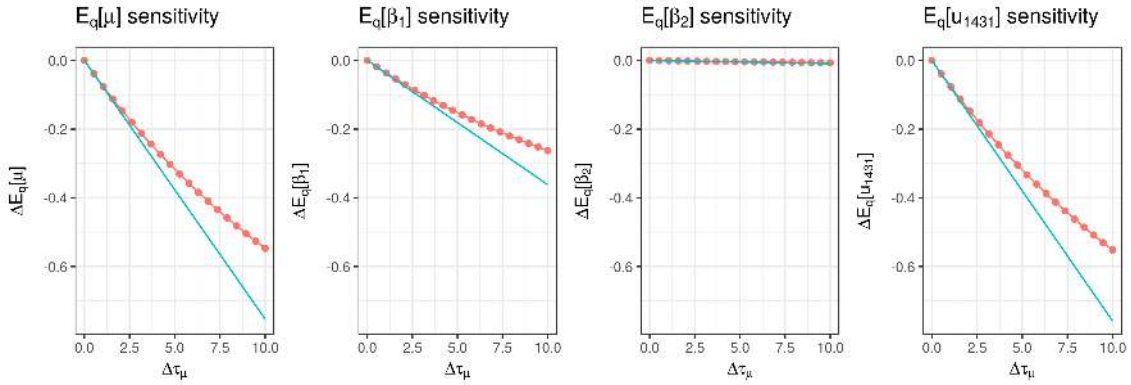


Figure 18: MFVB sensitivity as measured both by linear approximation (blue) and re-fitting (red)

Acknowledgments

We are grateful to the anonymous reviewers for their insightful comments and suggestions. Ryan Giordano’s research was funded in part by the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract number DE-AC02-05CH11231, and in part by the Gordon and Betty Moore Foundation through Grant GBMF3834 and by the Alfred P. Sloan Foundation through Grant 2013-10-27 to the University of California, Berkeley. Tamara Broderick’s research was supported in part by an NSF CAREER Award, an ARO YIP Award, and a Google Faculty Research Award. This work was also supported by the DARPA program on Lifelong Learning Machines, the Office of Naval Research under contract/grant number N00014-17-1-2072, and the Army Research Office under grant number W911NF-17-1-0304.

Appendices

Appendix A. Proof of Theorem 1

In this section we prove Theorem 1.

Proof Under Assumption 1, we can exchange differentiation and integration in $\frac{\partial}{\partial \alpha^\top} \int p_0(\theta) \exp(\rho(\theta, \alpha)) g(\theta) \lambda(d\theta)$ and $\frac{\partial}{\partial \alpha^\top} \int p_0(\theta) \exp(\rho(\theta, \alpha)) \lambda(d\theta)$ by Fleming (1965, Chapter 5-11, Theorem 18), which ultimately depends on the Lebesgue dominated convergence theorem. By Assumption 1, $\mathbb{E}_{p_\alpha}[g(\theta)]$ is well-defined for $\alpha \in \mathcal{A}_0$ and

$$\frac{\partial p_0(\theta) \exp(\rho(\theta, \alpha))}{\partial \alpha} = p_0(\theta) \exp(\rho(\theta, \alpha)) \frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \quad \lambda\text{-almost everywhere.}$$

Armed with these facts, we can directly compute

$$\begin{aligned} \left. \frac{d\mathbb{E}_{p_\alpha}[g(\theta)]}{d\alpha^\top} \right|_{\alpha_0} &= \left. \frac{d}{d\alpha^\top} \frac{\int g(\theta) p_0(\theta) \exp(\rho(\theta, \alpha)) \lambda(d\theta)}{\int p_0(\theta) \exp(\rho(\theta, \alpha)) \lambda(d\theta)} \right|_{\alpha_0} \\ &= \frac{\frac{\partial}{\partial \alpha^\top} \int g(\theta) p_0(\theta) \exp(\rho(\theta, \alpha)) \lambda(d\theta) \Big|_{\alpha_0}}{\int p_0(\theta) \exp(\rho(\theta, \alpha_0)) \lambda(d\theta)} - \mathbb{E}_{p_0}[g(\theta)] \frac{\frac{\partial}{\partial \alpha^\top} \int p_0(\theta) \exp(\rho(\theta, \alpha)) \lambda(d\theta) \Big|_{\alpha_0}}{\int p_0(\theta) \exp(\rho(\theta, \alpha_0)) \lambda(d\theta)} \\ &= \frac{\int g(\theta) p_0(\theta) \exp(\rho(\theta, \alpha)) \frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \Big|_{\alpha_0} \lambda(d\theta)}{\int p_0(\theta) \exp(\rho(\theta, \alpha_0)) \lambda(d\theta)} - \mathbb{E}_{p_0}[g(\theta)] \mathbb{E}_{p_0} \left[\frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \Big|_{\alpha_0} \right] \\ &= \text{Cov}_{p_0} \left(g(\theta), \frac{\partial \rho(\theta, \alpha)}{\partial \alpha} \Big|_{\alpha_0} \right). \end{aligned}$$

■

Appendix B. Comparison With MCMC Importance Sampling

In this section, we show that using importance sampling with MCMC samples to calculate the local sensitivity in Eq. (1) is precisely equivalent to using the same MCMC samples to estimate the covariance in Eq. (4) directly. For this section, will suppose that Assumption 1 holds. Further suppose, without loss of generality, we have samples θ_i drawn IID from $p_0(\theta)$:

$$\begin{aligned} \theta_n &\stackrel{iid}{\sim} p_0(\theta), \text{ for } n = 1, \dots, N_s \\ \mathbb{E}_{p_0}[g(\theta)] &\approx \frac{1}{N_s} \sum_{n=1}^{N_s} g(\theta_n). \end{aligned}$$

Typically we cannot compute the dependence of the normalizing constant $\int p(\theta') \exp(\rho(\theta', \alpha)) \lambda(d\theta')$ on α , so we use the following importance sampling estimate for $\mathbb{E}_{p_\alpha}[g(\theta)]$ (Owen, 2013, Chapter 9):

$$\begin{aligned} w_n &= \exp(\rho(\theta_n, \alpha) - \rho(\theta_n, \alpha_0)) \\ \tilde{w}_n &:= \frac{w_n}{\sum_{n'=1}^{N_s} w_{n'}} \\ \mathbb{E}_{p_\alpha}[g(\theta)] &\approx \sum_{n=1}^{N_s} \tilde{w}_n g(\theta_n). \end{aligned}$$

Note that $\tilde{w}_n|_{\alpha_0} = \frac{1}{N_s}$, so the importance sampling estimate recovers the ordinary sample mean at α_0 . The derivatives of the weights are given by

$$\begin{aligned} \frac{\partial w_n}{\partial \alpha} &= w_n \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} \\ \frac{\partial \tilde{w}_n}{\partial \alpha} &= \frac{\frac{\partial w_n}{\partial \alpha}}{\sum_{n'=1}^{N_s} w_{n'}} - \frac{w_n \sum_{n'=1}^{N_s} \frac{\partial w_{n'}}{\partial \alpha}}{\left(\sum_{n'=1}^{N_s} w_{n'}\right)^2} \\ &= \frac{w_n}{\sum_{n'=1}^{N_s} w_{n'}} \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} - \frac{w_n}{\sum_{n'=1}^{N_s} w_{n'}} \sum_{n'=1}^{N_s} \frac{w_{n'}}{\sum_{n'=1}^{N_s} w_{n'}} \frac{\partial \rho(\theta_{n'}, \alpha)}{\partial \alpha} \\ &= \tilde{w}_n \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} - \tilde{w}_n \sum_{n'=1}^{N_s} \tilde{w}_{n'} \frac{\partial \rho(\theta_{n'}, \alpha)}{\partial \alpha}. \end{aligned}$$

It follows that

$$\begin{aligned} \frac{\partial}{\partial \alpha} \sum_{n=1}^{N_s} \tilde{w}_n g(\theta_n) \Big|_{\alpha_0} &= \sum_{n=1}^{N_s} \left(\tilde{w}_n \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} - \tilde{w}_n \sum_{n'=1}^{N_s} \tilde{w}_{n'} \frac{\partial \rho(\theta_{n'}, \alpha)}{\partial \alpha} \right) \Big|_{\alpha_0} g(\theta_n) \\ &= \frac{1}{N_s} \sum_{n=1}^{N_s} \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} \Big|_{\alpha_0} g(\theta_n) - \left[\frac{1}{N_s} \sum_{n=1}^{N_s} \frac{\partial \rho(\theta_n, \alpha)}{\partial \alpha} \Big|_{\alpha_0} \right] \left[\frac{1}{N_s} \sum_{n=1}^{N_s} g(\theta_n) \right], \end{aligned}$$

which is precisely the sample version of the covariance in Theorem 1.

Appendix C. Our Use of the Terms ‘‘Sensitivity’’ and ‘‘Robustness’’

In this section we clarify our usage of the terms ‘‘robustness’’ and ‘‘sensitivity.’’ The quantity $\mathbf{S}_{\alpha_0}^T(\alpha - \alpha_0)$ measures the *sensitivity* of $\mathbb{E}_{p_\alpha}[g(\theta)]$ to perturbations in the direction $\Delta\alpha$. Intuitively, as sensitivity increases, robustness decreases, and, in this sense, sensitivity and robustness are opposites of one another. However, we emphasize that sensitivity is a clearly defined, measurable quantity and that robustness is a subjective judgment informed by sensitivity, but also by many other less objective considerations.

Suppose we have calculated \mathbf{S}_{α_0} from Eq. (1) and found that it has a particular value. To determine whether our model is robust, we must additionally decide

1. How large of a change in the prior, $|\alpha - \alpha_0|$, is plausible, and
2. How large of a change in $\mathbb{E}_{p_\alpha}[g(\theta)]$ is important.

The set of plausible prior values necessarily remains a subjective decision.⁴ Whether or not a particular change in $\mathbb{E}_{p_\alpha}[g(\theta)]$ is important depends on the ultimate use of the posterior mean. For example, the posterior standard deviation can be a guide: if the prior sensitivity is swamped by the posterior uncertainty then it can be neglected when reporting our subjective uncertainty about $g(\theta)$, and the model is robust. Similarly, even if the prior sensitivity is much larger than the posterior standard deviation but small enough that it would not affect any actionable decision made on the basis of the value of $\mathbb{E}_{p_\alpha}[g(\theta)]$, then the model is robust. Intermediate values remain a matter of judgment. An illustration of the relationship between sensitivity and robustness is shown in Fig. 19.

Finally, we note that if \mathcal{A} is small enough that $\mathbb{E}_{p_\alpha}[g(\theta)]$ is roughly linear in α for $\alpha \in \mathcal{A}$, then calculating Eq. (1) for all $\alpha \in \mathcal{A}$ and finding the worst case can be thought of as a first-order approximation to a global robustness estimate. Depending on the problem at hand, this linearity assumption may not be plausible except

4. This decision can be cast in a formal decision theoretic framework based on a partial ordering of subjective beliefs (Insua and Criado, 2000).

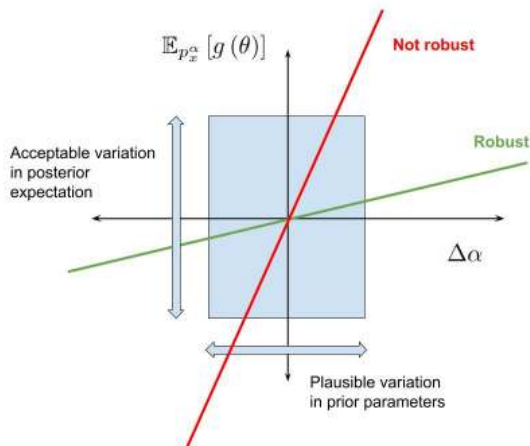


Figure 19: The relationship between robustness and sensitivity

for very small \mathcal{A} . This weakness is inherent to the local robustness approach. Nevertheless, even when the perturbations are valid only for a small \mathcal{A} , these easily-calculable measures may still provide valuable intuition about the potential modes of failure for a model.

If $g(\theta)$ is a scalar, it is natural to attempt to summarize the high-dimensional vector \mathbf{S}_{α_0} in a single easily reported number such as

$$\mathbf{S}_{\alpha_0}^{sup} := \sup_{\alpha: \|\alpha - \alpha_0\| \leq 1} |\mathbf{S}_{\alpha_0}^T (\alpha - \alpha_0)|.$$

For example, the calculation of $\mathbf{S}_{\alpha_0}^{sup}$ is the principal ambition of Basu et al. (1996). The use of such summaries is also particularly common in work that considers function-valued perturbations (e.g., Gustafson, 1996b; Roos et al., 2015). (Function-valued perturbations can be connected to the finite-dimensional perturbations of the present work through the notion of the Gateaux derivative (Huber, 2011, Chapter 2.5), the elaboration of which we leave to future work.) Although the summary $\mathbf{S}_{\alpha_0}^{sup}$ has obvious merits, in the present work we emphasize the calculation only of \mathbf{S}_{α_0} in the belief that its interpretation is likely to vary from application to application and require some critical thought and subjective judgment. For example, the unit ball $\|\alpha - \alpha_0\| \leq 1$ (as in Basu et al. (1996)) may not make sense as a subjective description of the range of plausible variability of $p(\theta|\alpha)$. Consider, e.g.: why should the off-diagonal term of a Wishart prior plausibly vary as widely as the mean of some other parameter, when the two might not even have the same units? This problem is easily remedied by choosing an appropriate scaling of the parameters and thereby making the unit ball an appropriate range for the problem at hand, but the right scaling will vary from problem to problem and necessarily be a somewhat subjective choice, so we refrain from taking a stand on this decision. As another example, the worst-case function-valued perturbations of Gustafson (1996a,b) require a choice of a metric ball in function space whose meaning may not be intuitively obvious, may provide worst-case perturbations that depend on the data to a subjectively implausible degree, and may exhibit interesting but perhaps counter-intuitive asymptotic behavior for different norms and perturbation dimensions. Consequently, we do not attempt to prescribe a particular one-size-fits-all summary measure. The local sensitivity \mathbf{S}_{α_0} is a well-defined mathematical quantity. Its relationship to robustness must remain a matter of judgment.

Appendix D. Proof of Theorem 2

In this section we prove Theorem 2.

Proof For notational convenience, we will define

$$KL(\eta, \alpha) := KL(q(\theta; \eta) || p_\alpha(\theta)).$$

By Assumption 3, $\eta^*(\alpha)$ is both optimal and interior for all $\alpha \in \mathcal{A}_0$, and by Assumption 2, $KL(\eta, \alpha)$ is continuously differentiable in η . Therefore, the first-order conditions of the optimization problem in Eq. (8) give:

$$\left. \frac{\partial KL(\eta, \alpha)}{\partial \eta} \right|_{\eta=\eta^*(\alpha)} = 0 \text{ for all } \alpha \in \mathcal{A}_0. \quad (26)$$

$\left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \eta^\top} \right|_{\alpha_0}$ is positive definite by the strict optimality of η^* in Assumption 3, and $\left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \alpha^\top} \right|_{\alpha_0}$ is continuous by Assumption 2. It follows that $\eta^*(\alpha)$ is a continuously differentiable function of α by application of the implicit function theorem to the first-order condition in Eq. (26) (Fleming, 1965, Chapter 4.6). So we can use the chain rule to take the total derivative of Eq. (26) with respect to α .

$$\begin{aligned} \frac{d}{d\alpha} \left(\left. \frac{\partial KL(\eta, \alpha)}{\partial \eta} \right|_{\eta=\eta^*(\alpha)} \right) &= 0 \text{ for all } \alpha \in \mathcal{A}_0 \Rightarrow \\ \left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \eta^\top} \right|_{\eta=\eta^*(\alpha)} \frac{d\eta^*(\alpha)}{d\alpha^\top} + \left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \alpha^\top} \right|_{\eta=\eta^*(\alpha)} &= 0 \text{ for all } \alpha \in \mathcal{A}_0. \end{aligned}$$

The strict optimality of $KL(\eta, \alpha)$ at $\eta^*(\alpha)$ in Assumption 3 requires that $\left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \eta^\top} \right|_{\eta=\eta^*(\alpha)}$ be invertible. So we can evaluate at $\alpha = \alpha_0$ and solve to find that

$$\left. \frac{d\eta^*(\alpha)}{d\alpha^\top} \right|_{\alpha_0} = - \left(\left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \eta^\top} \right|_{\eta=\eta_0^*, \alpha=\alpha_0} \right)^{-1} \left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \alpha^\top} \right|_{\eta=\eta_0^*, \alpha=\alpha_0}.$$

$\mathbb{E}_{q_\alpha}[g(\theta)]$ is a continuously differentiable function of $\eta^*(\alpha)$ by Assumption 4. So by the chain rule and Assumption 2, we have that

$$\left. \frac{d\mathbb{E}_{q(\theta; \eta)}[g(\theta)]}{d\alpha^\top} \right|_{\alpha_0} = \left. \frac{\partial \mathbb{E}_{q(\theta; \eta)}[g(\theta)]}{\partial \eta} \frac{d\eta^*(\alpha)}{d\alpha^\top} \right|_{\eta=\eta_0^*, \alpha=\alpha_0}.$$

Finally, we observe that

$$\begin{aligned} KL(\eta, \alpha) &= \mathbb{E}_{q(\theta; \eta)}[\log q(\theta; \eta) - \log p(\theta) - \rho(\theta, \alpha)] + Constant \Rightarrow \\ \left. \frac{\partial^2 KL(\eta, \alpha)}{\partial \eta \partial \alpha^\top} \right|_{\eta=\eta_0^*, \alpha=\alpha_0} &= - \left. \frac{\partial^2 \mathbb{E}_{q(\theta; \eta)}[\rho(\theta, \alpha)]}{\partial \eta \partial \alpha^\top} \right|_{\eta=\eta_0^*, \alpha=\alpha_0}. \end{aligned}$$

Here, the term *Constant* contains quantities that do not depend on η . Plugging in gives the desired result. ■

Appendix E. Exactness of Multivariate Normal Posterior Means

In this section, we show that the MFVB estimate of the posterior means of a multivariate normal with known covariance is exact and that, as an immediate consequence, the linear response covariance recovers the exact posterior covariance, i.e., $\text{Cov}_{q_0}^{LR}(\theta) = \text{Cov}_{p_0}(\theta)$.

Suppose we are using MFVB to approximate a non-degenerate multivariate normal posterior, i.e.,

$$p_0(\theta) = \mathcal{N}(\theta; \mu, \Sigma)$$

for full-rank Σ . This posterior arises, for instance, given a multivariate normal likelihood $p(x|\mu) = \prod_{n=1:N} \mathcal{N}(x_n|\theta, \Sigma_x)$ with known covariance Σ_x and a conjugate multivariate normal prior on the unknown mean parameter $\theta \in \mathbb{R}^K$. Additionally, even when the likelihood is non-normal or the prior is not conjugate, the posterior may be closely approximated by a multivariate normal distribution when a Bayesian central limit theorem can be applied (Le Cam and Yang, 2012, Chapter 8).

We will consider an MFVB approximation to $p_0(\theta)$. Specifically, let the elements of the vector θ be given by scalars θ_k for $k = 1, \dots, K$, and take the MFVB normal approximation with means m_k and variances v_k :

$$\mathcal{Q} = \left\{ q(\theta) : q(\theta) = \prod_{k=1}^K \mathcal{N}(\theta_k; m_k, v_k) \right\}.$$

In the notation of Eq. (9), we have $\eta_k = (m_k, v_k)^\top$ with $\Omega_\eta = \{\eta : v_k > 0, \forall k = 1, \dots, K\}$. The optimal variational parameters are given by $\eta_k^* = (m_k^*, v_k^*)^\top$.

Lemma 1 *Let $p_0(\theta) = \mathcal{N}(\theta; \mu, \Sigma)$ for full-rank Σ and let $\mathcal{Q} = \left\{ q(\theta) : q(\theta) = \prod_{k=1}^K \mathcal{N}(\theta_k; m_k, v_k) \right\}$ be the mean field approximating family. Then there exists an $\eta^* = (m^*, v^*)$ that solves*

$$\eta^* = \underset{\eta: q(\theta; \eta) \in \mathcal{Q}}{\operatorname{argmin}} KL(q(\theta; \eta) || p_\alpha(\theta))$$

with $m^* = \mu$.

Proof Let $\operatorname{diag}(v)$ denote the $K \times K$ matrix with the vector v on the diagonal and zero elsewhere. Using the fact that the entropy of a univariate normal distribution with variance v is $\frac{1}{2} \log v$ plus a constant, the variational objective in Eq. (8) is given by

$$\begin{aligned} KL(q(\theta; \eta) || p_\alpha(\theta)) &= \mathbb{E}_{q(\theta; \eta)} \left[\frac{1}{2} (\theta - \mu)^\top \Sigma^{-1} (\theta - \mu) \right] - \frac{1}{2} \sum_k \log v_k + \text{Constant} \\ &= \frac{1}{2} \operatorname{trace}(\Sigma^{-1} \mathbb{E}_{q(\theta; \eta)}[\theta \theta^\top]) - \mu^\top \Sigma^{-1} \mathbb{E}_{q(\theta; \eta)}[\theta] - \frac{1}{2} \sum_k \log v_k + \text{Constant} \\ &= \frac{1}{2} \operatorname{trace}(\Sigma^{-1} (m m^\top + \operatorname{diag}(v))) - \mu^\top \Sigma^{-1} m - \frac{1}{2} \sum_k \log v_k + \text{Constant} \\ &= \frac{1}{2} \operatorname{trace}(\Sigma^{-1} \operatorname{diag}(v)) + \frac{1}{2} m^\top \Sigma^{-1} m - \mu^\top \Sigma^{-1} m - \frac{1}{2} \sum_k \log v_k + \text{Constant}. \end{aligned} \tag{27}$$

The first-order condition for the optimal m^* is then

$$\begin{aligned} \frac{\partial KL(q(\theta; \eta) || p_0(\theta))}{\partial m} \Big|_{m=m^*, v=v^*} &= 0 \Rightarrow \\ \Sigma^{-1} m^* - \Sigma^{-1} \mu &= 0 \Rightarrow \\ m^* &= \mu. \end{aligned}$$

The optimal variances follow similarly:

$$\begin{aligned} \frac{\partial KL(q(\theta; \eta) || p_0(\theta))}{\partial v_k} \Big|_{m=m^*, v=v^*} &= 0 \Rightarrow \\ \frac{1}{2} (\Sigma^{-1})_{kk} - \frac{1}{2} \frac{1}{v_k^*} &= 0 \Rightarrow \\ v_k^* &= \frac{1}{(\Sigma^{-1})_{kk}}. \end{aligned}$$

Since $v_k^* > 0$, we have $\eta^* \in \Omega_\eta$.

Lemma 1 can be also be derived via the variational coordinate ascent updates (Bishop (2006, Section 10.1.2) and Giordano et al. (2015, Appendix B)). ■

Next, we show that Lemma 1 holds for all perturbations of the form $\rho(\theta, \alpha) = \alpha^\top \theta$ with $\alpha_0 = 0$ and that Assumptions 1–4 are satisfied for all finite α .

Lemma 2 *Under the conditions of Lemma 1, let $p_\alpha(\theta)$ be defined from Eq. (2) with $\rho(\theta, \alpha) = \alpha^\top \theta$ and $\alpha_0 = 0$. Take $g(\theta) = \theta$. Then, for all finite α , Assumptions 1–4 are satisfied, and Condition 1 is satisfied with equality.*

Proof Up to a constant that does not depend on θ , the log density of $p_\alpha(\theta)$ is

$$\begin{aligned} \log p_\alpha(\theta) &= -\frac{1}{2}(\theta - \mu)^\top \Sigma^{-1}(\theta - \mu) + \alpha^\top \theta + \text{Constant} \\ &= -\frac{1}{2}\theta^\top \Sigma^{-1}\theta - \frac{1}{2}\mu^\top \Sigma^{-1}\mu + (\mu^\top \Sigma^{-1} + \alpha^\top)\theta + \text{Constant}. \end{aligned}$$

Since θ is a natural sufficient statistic of the multivariate normal distribution and the corresponding natural parameter of $p_\alpha(\theta)$, $\Sigma^{-1}\mu + \alpha$, is interior when Σ is full-rank, $p_\alpha(\theta)$ is multivariate normal for any finite α . Assumption 1 follows immediately.

By inspection of Eq. (27), Assumption 2 is satisfied. Because Ω_η is an open set and Σ is positive definite, Assumption 3 is satisfied. Since $\mathbb{E}_{q(\theta;\eta)}[g(\theta)] = m$, Assumption 4 is satisfied. Finally, by Lemma 1, $\mathbb{E}_{q_\alpha}[\theta] = \mathbb{E}_{p_\alpha}[\theta]$, so Condition 1 is satisfied with equality. ■

It now follows immediately from Definition 6 that the linear response variational covariance exactly reproduces the exact posterior covariance for the multivariate normal distribution.

Corollary 4 *Under the conditions of Lemma 2, $\text{Cov}_{q_0}^{LR}(\theta) = \text{Cov}_{p_0}(\theta)$.*

Appendix F. ADVI Model Details

This section reports the Stan code for the models used in Section 5.2. For details on how to interpret the models as well as the unconstraining transforms, see the Stan manual (Stan Team, 2015). For the associated data, see the Stan example models wiki (Stan Team, 2017).

F.1 Election Model (`election88.stan`)

Listing 1: `election88.stan`

```

1 data {
2   int<lower=0> N;
3   int<lower=0> n_state;
4   vector<lower=0, upper=1>[N] black;
5   vector<lower=0, upper=1>[N] female;
6   int<lower=1, upper=n_state> state[N];
7   int<lower=0, upper=1> y[N];
8 }
9 parameters {
10  vector[n_state] a;
11  vector[2] b;
12  real<lower=0, upper=100> sigma_a;
13  real mu_a;
14 }
```

```

15 transformed parameters {
16   vector[N] y_hat;
17
18   for (i in 1:N)
19     y_hat[i] <- b[1] * black[i] + b[2] * female[i] + a[state[i]];
20 }
21 model {
22   mu_a ~ normal(0, 1);
23   a ~ normal(mu_a, sigma_a);
24   b ~ normal(0, 100);
25   y ~ bernoulli_logit(y_hat);
26 }

```

F.2 Sesame Street Model (`sesame_street1`)

Listing 2: `sesame_street1.stan`

```

1 data {
2   int<lower=0> J;
3   int<lower=0> N;
4   int<lower=1,upper=J> siteset[N];
5   vector[2] yt[N];
6   vector[N] z;
7 }
8 parameters {
9   vector[2] ag[J];
10  real b;
11  real d;
12  real<lower=-1,upper=1> rho_ag;
13  real<lower=-1,upper=1> rho_yt;
14  vector[2] mu_ag;
15  real<lower=0,upper=100> sigma_a;
16  real<lower=0,upper=100> sigma_g;
17  real<lower=0,upper=100> sigma_t;
18  real<lower=0,upper=100> sigma_y;
19 }
20 model {
21  vector[J] a;
22  vector[J] g;
23  matrix[2,2] Sigma_ag;
24  matrix[2,2] Sigma_yt;
25  vector[2] yt_hat[N];
26
27  //data level
28  Sigma_yt[1,1] <- pow(sigma_y,2);
29  Sigma_yt[2,2] <- pow(sigma_t,2);
30  Sigma_yt[1,2] <- rho_yt*sigma_y*sigma_t;
31  Sigma_yt[2,1] <- Sigma_yt[1,2];
32
33  // group level
34  Sigma_ag[1,1] <- pow(sigma_a,2);
35  Sigma_ag[2,2] <- pow(sigma_g,2);

```

```

36 Sigma_ag[1,2] <- rho_ag*sigma_a*sigma_g;
37 Sigma_ag[2,1] <- Sigma_ag[1,2];
38
39 for (j in 1:J) {
40   a[j] <- ag[j,1];
41   g[j] <- ag[j,2];
42 }
43
44 for (i in 1:N) {
45   yt_hat[i,2] <- g[siteset[i]] + d * z[i];
46   yt_hat[i,1] <- a[siteset[i]] + b * yt[i,2];
47 }
48
49 //data level
50 sigma_y ~ uniform (0, 100);
51 sigma_t ~ uniform (0, 100);
52 rho_yt ~ uniform(-1, 1);
53 d ~ normal (0, 31.6);
54 b ~ normal (0, 31.6);
55
56 //group level
57 sigma_a ~ uniform (0, 100);
58 sigma_g ~ uniform (0, 100);
59 rho_ag ~ uniform(-1, 1);
60 mu_ag ~ normal (0, 31.6);
61
62 for (j in 1:J)
63   ag[j] ~ multi_normal(mu_ag, Sigma_ag);
64
65 //data model
66 for (i in 1:N)
67   yt[i] ~ multi_normal(yt_hat[i], Sigma_yt);
68
69 }

```

F.3 Radon Model (radon_vary_intercept_floor)

Listing 3: radon_vary_intercept_floor.stan

```

1 data {
2   int<lower=0> J;
3   int<lower=0> N;
4   int<lower=1,upper=J> county[N];
5   vector[N] u;
6   vector[N] x;
7   vector[N] y;
8 }
9 parameters {
10  vector[J] a;
11  vector[2] b;
12  real mu_a;
13  real<lower=0,upper=100> sigma_a;

```



```

14  real<lower=0,upper=100> sigma_y;
15  }
16  transformed parameters {
17    vector[N] y_hat;
18
19    for (i in 1:N)
20      y_hat[i] <- a[county[i]] + u[i] * b[1] + x[i] * b[2];
21  }
22  model {
23    mu_a ~ normal(0, 1);
24    a ~ normal(mu_a, sigma_a);
25    b ~ normal(0, 1);
26    y ~ normal(y_hat, sigma_y);
27  }

```

F.4 Ecology Model (cjs_cov_ranef)

Listing 4: cjs_cov_ranef.stan

```

1 // This models is derived from section 12.3 of "Stan Modeling Language
2 // User's Guide and Reference Manual"
3
4 functions {
5   int first_capture(int[] y_i) {
6     for (k in 1:size(y_i))
7       if (y_i[k])
8         return k;
9     return 0;
10  }
11
12  int last_capture(int[] y_i) {
13    for (k_rev in 0:(size(y_i) - 1)) {
14      // Compound declaration was enabled in Stan 2.13
15      int k = size(y_i) - k_rev;
16      //   int k;
17      //   k = size(y_i) - k_rev;
18      if (y_i[k])
19        return k;
20    }
21    return 0;
22  }
23
24  matrix prob_uncaptured(int nind, int n_occasions,
25                        matrix p, matrix phi) {
26    matrix[nind, n_occasions] chi;
27
28    for (i in 1:nind) {
29      chi[i, n_occasions] = 1.0;
30      for (t in 1:(n_occasions - 1)) {
31        // Compound declaration was enabled in Stan 2.13
32        int t_curr = n_occasions - t;
33        int t_next = t_curr + 1;

```

```

34     /*
35     int t_curr;
36     int t_next;
37
38     t_curr = n_occasions - t;
39     t_next = t_curr + 1;
40     */
41     chi[i, t_curr] = (1 - phi[i, t_curr])
42                   + phi[i, t_curr] * (1 - p[i, t_next - 1]) * chi[
43     i, t_next];
44   }
45   return chi;
46 }
47 }
48
49 data {
50   int<lower=0> nind;           // Number of individuals
51   int<lower=2> n_occasions;   // Number of capture occasions
52   int<lower=0,upper=1> y[nind, n_occasions]; // Capture-history
53   vector[n_occasions - 1] x; // Covariate
54 }
55
56 transformed data {
57   int n_occ_minus_1 = n_occasions - 1;
58   // int n_occ_minus_1;
59   int<lower=0,upper=n_occasions> first[nind];
60   int<lower=0,upper=n_occasions> last[nind];
61   vector<lower=0,upper=nind>[n_occasions] n_captured;
62
63   // n_occ_minus_1 = n_occasions - 1;
64   for (i in 1:nind)
65     first[i] = first_capture(y[i]);
66   for (i in 1:nind)
67     last[i] = last_capture(y[i]);
68   n_captured = rep_vector(0, n_occasions);
69   for (t in 1:n_occasions)
70     for (i in 1:nind)
71       if (y[i, t])
72         n_captured[t] = n_captured[t] + 1;
73 }
74
75 parameters {
76   real beta; // Slope parameter
77   real<lower=0,upper=1> mean_phi; // Mean survival
78   real<lower=0,upper=1> mean_p; // Mean recapture
79   vector[n_occ_minus_1] epsilon;
80   real<lower=0,upper=10> sigma;
81   // In case a weakly informative prior is used
82   // real<lower=0> sigma;
83 }

```

```

84
85 transformed parameters {
86   matrix<lower=0,upper=1>[nind, n_occ_minus_1] phi;
87   matrix<lower=0,upper=1>[nind, n_occ_minus_1] p;
88   matrix<lower=0,upper=1>[nind, n_occasions] chi;
89   // Compoud declaration was enabled in Stan 2.13
90   real mu = logit(mean_phi);
91   // real mu;
92
93   // mu = logit(mean_phi);
94   // Constraints
95   for (i in 1:nind) {
96     for (t in 1:(first[i] - 1)) {
97       phi[i, t] = 0;
98       p[i, t] = 0;
99     }
100    for (t in first[i]:n_occ_minus_1) {
101      phi[i, t] = inv_logit(mu + beta * x[t] + epsilon[t]);
102      p[i, t] = mean_p;
103    }
104  }
105
106  chi = prob_uncaptured(nind, n_occasions, p, phi);
107 }
108
109 model {
110   // Priors
111   // Uniform priors are implicitly defined.
112   // mean_phi ~ uniform(0, 1);
113   // mean_p ~ uniform(0, 1);
114   // sigma ~ uniform(0, 10);
115   // In case a weakly informative prior is used
116   // sigma ~ normal(5, 2.5);
117   beta ~ normal(0, 100);
118   epsilon ~ normal(0, sigma);
119
120   for (i in 1:nind) {
121     if (first[i] > 0) {
122       for (t in (first[i] + 1):last[i]) {
123         l ~ bernoulli(phi[i, t - 1]);
124         y[i, t] ~ bernoulli(p[i, t - 1]);
125       }
126       l ~ bernoulli(chi[i, last[i]]);
127     }
128   }
129 }
130
131 generated quantities {
132   real<lower=0> sigma2;
133   vector<lower=0,upper=1>[n_occ_minus_1] phi_est;
134

```

```

135 sigma2 = square(sigma);
136 // inv_logit was vectorized in Stan 2.13
137 phi_est = inv_logit(mu + beta * x + epsilon); // Yearly survival
138 /*
139 for (t in 1:n_occ_minus_1)
140   phi_est[t] = inv_logit(mu + beta * x[t] + epsilon[t]);
141 */
142 }

```

Appendix G. LKJ Priors for Covariance Matrices in Mean Field Variational Inference

In this section we briefly derive closed-form expressions for using an LKJ prior with a Wishart variational approximation.

Proposition 3 *Let Σ be a $K \times K$ positive definite covariance matrix. Define the $K \times K$ matrix \mathbf{M} such that*

$$\mathbf{M}_{ij} = \begin{cases} \sqrt{\Sigma_{ij}} & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

Define the correlation matrix \mathbf{R} as

$$\mathbf{R} = \mathbf{M}^{-1} \Sigma \mathbf{M}^{-1}.$$

Define the LKJ prior on \mathbf{R} with concentration parameter ξ (Lewandowski et al., 2009):

$$p_{\text{LKJ}}(\mathbf{R}|\xi) \propto |\mathbf{R}|^{\xi-1}.$$

Let $q(\Sigma|\mathbf{V}^{-1}, \nu)$ be an inverse Wishart distribution with matrix parameter \mathbf{V}^{-1} and ν degrees of freedom. Then

$$\mathbb{E}_q[\log |\mathbf{R}|] = \log |\mathbf{V}^{-1}| - \psi_K\left(\frac{\nu}{2}\right) - \sum_{k=1}^K \log((\mathbf{V}^{-1})_{kk}) + K\psi\left(\frac{\nu - K + 1}{2}\right) + \text{Constant}$$

$$\mathbb{E}_q[\log p_{\text{LKJ}}(\mathbf{R}|\xi)] = (\xi - 1) \mathbb{E}_q[\log |\mathbf{R}|] + \text{Constant},$$

where *Constant* does not depend on \mathbf{V} or ν . Here, ψ_K is the multivariate digamma function.

Proof First note that

$$\begin{aligned} \log |\Sigma| &= 2 \log |\mathbf{M}| + \log |\mathbf{R}| \\ &= 2 \sum_{k=1}^K \log \sqrt{\Sigma_{kk}} + \log |\mathbf{R}| \\ &= \sum_{k=1}^K \log \Sigma_{kk} + \log |\mathbf{R}| \Rightarrow \\ \log |\mathbf{R}| &= \log |\Sigma| - \sum_{k=1}^K \log \Sigma_{kk}. \end{aligned} \tag{28}$$

By Eq. B.81 in (Bishop, 2006), a property of the inverse Wishart distribution is the following relation.

$$E_q[\log |\Sigma|] = \log |\mathbf{V}^{-1}| - \psi_K\left(\frac{\nu}{2}\right) - K \log 2, \tag{29}$$

where ψ_K is the multivariate digamma function. By the marginalization property of the inverse Wishart distribution,

$$\begin{aligned} \Sigma_{kk} &\sim \text{InverseWishart} \left((\mathbf{V}^{-1})_{kk}, \nu - K + 1 \right) \Rightarrow \\ E_q [\log \Sigma_{kk}] &= \log \left((\mathbf{V}^{-1})_{kk} \right) - \psi \left(\frac{\nu - K + 1}{2} \right) - \log 2. \end{aligned} \quad (30)$$

Plugging Eq. (29) and Eq. (30) into Eq. (28) gives the desired result. \blacksquare

Appendix H. Logistic GLMM Model Details

In this section we include extra details about the model and analysis of Section 5. We will continue to use the notation defined therein. We use *Constant* to denote any constants that do not depend on the prior parameters, parameters, or data. The log likelihood is

$$\begin{aligned} \log p(y_{it}|u_t, \beta) &= y_{it} \log \left(\frac{p_{it}}{1 - p_{it}} \right) + \log(1 - p_{it}) \\ &= y_{it}\rho + \log(1 - p_{it}) + \text{Constant} \\ \log p(u|\mu, \tau) &= -\frac{1}{2}\tau \sum_{t=1}^T (u_t - \mu)^2 - \frac{1}{2}T \log \tau \\ &= -\frac{1}{2}\tau \sum_{t=1}^T (u_t^2 - \mu u_t + \mu^2) - \frac{1}{2}T \log \tau + \text{Constant} \\ \log p(\mu, \tau, \beta) &= -\frac{1}{2}\sigma_\mu^{-2} (\mu^2 + 2\mu\mu_0) + \\ &\quad (1 - \alpha_\tau) \tau + \beta_\tau \log \tau + \\ &\quad -\frac{1}{2} \left(\text{trace} \left(\Sigma_\beta^{-1} \beta \beta^T \right) + 2\text{trace} \left(\Sigma_\beta^{-1} \beta_0 \beta^T \right) \right). \end{aligned} \quad (31)$$

The prior parameters were taken to be

$$\begin{aligned} \mu_0 &= 0.000 \\ \sigma_\mu^{-2} &= 0.010 \\ \beta_0 &= 0.000 \\ \sigma_\beta^{-2} &= 0.100 \\ \alpha_\tau &= 3.000 \\ \beta_\tau &= 3.000. \end{aligned}$$

Under the variational approximation, ρ_{it} is normally distributed given x_{it} , with

$$\begin{aligned} \rho_{it} &= x_{it}^T \beta + u_t \\ \mathbb{E}_q [\rho_{it}] &= x_{it}^T \mathbb{E}_q [\beta] + \mathbb{E}_q [u_t] \\ \text{Var}_q (\rho_{it}) &= \mathbb{E}_q [\beta^T x_{it} x_{it}^T \beta] - \mathbb{E}_q [\beta]^T x_{it} x_{it}^T \mathbb{E}_q [\beta] + \text{Var}_q (u_t) \\ &= \mathbb{E}_q [\text{tr} (\beta^T x_{it} x_{it}^T \beta)] - \text{tr} \left(\mathbb{E}_q [\beta]^T x_{it} x_{it}^T \mathbb{E}_q [\beta] \right) + \text{Var}_q (u_t) \\ &= \text{tr} \left(x_{it} x_{it}^T \left(\mathbb{E}_q [\beta \beta^T] - \mathbb{E}_q [\beta] \mathbb{E}_q [\beta]^T \right) \right) + \text{Var}_q (u_t). \end{aligned}$$

We can thus use $n_{MC} = 4$ points of Gauss-Hermite quadrature to numerically estimate $\mathbb{E}_q \left[\log \left(1 - \frac{e^\rho}{1+e^\rho} \right) \right]$:

$$\rho_{it,s} := \sqrt{\text{Var}_q(\rho_{it})} z_s + \mathbb{E}_q[\rho_{it}]$$

$$\mathbb{E}_q \left[\log \left(1 - \frac{e^{\rho_{it}}}{1 + e^{\rho_{it}}} \right) \right] \approx \frac{1}{n_{MC}} \sum_{s=1}^{n_{MC}} \log \left(1 - \frac{e^{\rho_{it,s}}}{1 + e^{\rho_{it,s}}} \right)$$

We found that increasing the number of points used for the quadrature did not measurably change any of the results. The integration points and weights were calculated using the `numpy.polynomial.hermite` module in Python (Jones et al., 2001).

References

- A. Agresti and M. Kateri. *Categorical Data Analysis*. Springer, 2011.
- S. Basu, S. Rao Jammalamadaka, and W. Liu. Local posterior robustness with parametric priors: Maximum and average sensitivity. In *Maximum Entropy and Bayesian Methods*, pages 97–106. Springer, 1996.
- A. Baydin, B. Pearlmutter, A. Radul, and J. Siskind. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*, 18:1–43, 2018.
- J. O. Berger, D. R. Insua, and F. Ruggeri. Robust Bayesian analysis. In D. R. Insua and F. Ruggeri, editors, *Robust Bayesian Analysis*, volume 152. Springer Science & Business Media, 2000.
- C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, New York, 2006. Chapter 10.
- D. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3: 993–1022, 2003.
- D. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *arXiv preprint arXiv:1601.00670*, 2016.
- B. Carpenter, M. D. Hoffman, M. Brubaker, D. Lee, P. Li, and M. Betancourt. The Stan math library: Reverse-mode automatic differentiation in C++. *arXiv preprint arXiv:1509.07164*, 2015.
- R. D. Cook. Assessment of local influence. *Journal of the Royal Statistical Society: Series B*, 28(2):133–169, 1986.
- Criteo Labs. Criteo conversion logs dataset, 2014. URL <http://criteolabs.wpengine.com/downloads/2014-conversion-logs-dataset/>. Downloaded on July 27th, 2017.
- P. Diaconis and D. Freedman. On the consistency of Bayes estimates. *The Annals of Statistics*, 14(1):1–26, 1986.
- B. Efron. Frequentist accuracy of Bayesian estimates. *Journal of the Royal Statistical Society: Series B*, 77(3):617–646, 2015.
- W. H. Fleming. *Functions of Several Variables*. Addison-Wesley Publishing Company, Inc., 1965.
- A. Gelman and J. Hill. *Data analysis using regression and multilevel/hierarchical models*. Cambridge university press, 2006.
- A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*, volume 2. Chapman & Hall/CRC, 2014.
- R. J. Giordano, T. Broderick, and M. I. Jordan. Linear response methods for accurate covariance estimates from mean field variational Bayes. In *Advances in Neural Information Processing Systems*, pages 1441–1449, 2015.
- R. J. Giordano, T. Broderick, R. Meager, J. Huggins, and M. I. Jordan. Fast robustness quantification with variational Bayes. *arXiv preprint arXiv:1606.07153*, 2016.
- P. Gustafson. Local sensitivity of inferences to prior marginals. *Journal of the American Statistical Association*, 91(434):774–781, 1996a.
- P. Gustafson. Local sensitivity of posterior expectations. *The Annals of Statistics*, 24(1):174–195, 1996b.
- P. Gustafson. Local robustness in Bayesian analysis. In D. R. Insua and F. Ruggeri, editors, *Robust Bayesian Analysis*, volume 152. Springer Science & Business Media, 2000.

- P. J. Huber. *Robust Statistics*. Springer, 2011.
- D. R. Insua and R. Criado. Topics on the foundations of robust Bayesian analysis. In D. R. Insua and F. Ruggeri, editors, *Robust Bayesian Analysis*, volume 152. Springer Science & Business Media, 2000.
- E. Jones, T. Oliphant, P. Peterson, et al. SciPy: Open source scientific tools for Python, 2001. URL <http://www.scipy.org/>.
- H. J. Kappen and F. B. Rodriguez. Efficient learning in Boltzmann machines using linear response theory. *Neural Computation*, 10(5):1137–1156, 1998.
- R. W. Keener. *Theoretical Statistics: Topics for a Core Course*. Springer Science & Business Media, 2010.
- M. Kéry and M. Schaub. *Bayesian population analysis using WinBUGS: a hierarchical perspective*. Academic Press, 2011.
- D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. Blei. Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474, 2017.
- L. Le Cam and G. L. Yang. *Asymptotics in Statistics: Some Basic Concepts*. Springer Science & Business Media, 2012.
- D. Lewandowski, D. Kurowicka, and H. Joe. Generating random correlation matrices based on vines and extended onion method. *Journal of Multivariate Analysis*, 100(9):1989–2001, 2009.
- Q. Liu and D. Wang. Stein variational gradient descent: A general purpose Bayesian inference algorithm. In *Advances In Neural Information Processing Systems*, pages 2370–2378, 2016.
- D. J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003. Chapter 33.
- D. Maclaurin, D. Duvenaud, and R. P. Adams. Autograd: Effortless gradients in numpy. In *International Conference on Machine Learning 2015 AutoML Workshop*, 2015.
- E. Moreno. Global Bayesian robustness for some classes of prior distributions. In D. R. Insua and F. Ruggeri, editors, *Robust Bayesian Analysis*, volume 152. Springer Science & Business Media, 2000.
- R. M. Neal and G. E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In *Learning in Graphical Models*, pages 355–368. Springer, 1998.
- M. Opper and D. Saad. *Advanced Mean Field Methods: Theory and Practice*. MIT press, 2001.
- M. Opper and O. Winther. Variational linear response. In *Advances in Neural Information Processing Systems*, pages 1157–1164, 2004.
- A. B. Owen. *Monte Carlo Theory, Methods and Examples*. 2013. URL <http://statweb.stanford.edu/~owen/mc/>. Accessed November 23rd, 2016.
- C. J. Pérez, J. Martín, and M. J. Rufo. MCMC-based local parametric sensitivity estimations. *Computational Statistics & Data Analysis*, 51(2):823–835, 2006.
- R. Ranganath, D. Tran, and D. Blei. Hierarchical variational models. In *International Conference on Machine Learning*, pages 324–333, 2016.
- J. Regier, A. Miller, J. McAuliffe, R. Adams, M. Hoffman, D. Lang, D. Schlegel, and M. Prabhat. Celeste: Variational inference for a generative model of astronomical images. In *International Conference on Machine Learning*, pages 2095–2103, 2015.

- D. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, pages 1530–1538, 2015.
- M. Roos, T. G. Martins, L. Held, and H. Rue. Sensitivity analysis for Bayesian hierarchical models. *Bayesian Analysis*, 10(2):321–349, 2015.
- Stan Team. *Stan Modeling Language Users Guide and Reference Manual, Version 2.8.0*, 2015. URL <http://mc-stan.org/>.
- Stan Team. Stan example models wiki, 2017. URL <https://github.com/stan-dev/example-models/wiki>. Referenced on May 19th, 2017.
- T. Tanaka. Mean-field theory of Boltzmann machine learning. *Physical Review E*, 58(2):2302, 1998.
- T. Tanaka. Information geometry of mean-field approximation. *Neural Computation*, 12(8):1951–1968, 2000.
- Michalis Titsias and Miguel Lázaro-Gredilla. Doubly stochastic variational Bayes for non-conjugate inference. In *International Conference on Machine Learning*, pages 1971–1979, 2014.
- D. Tran, D. Blei, and E. M. Airoldi. Copula variational inference. In *Advances in Neural Information Processing Systems*, pages 3564–3572, 2015a.
- D. Tran, R. Ranganath, and D. Blei. The variational Gaussian process. *arXiv preprint arXiv:1511.06499*, 2015b.
- R. E. Turner and M. Sahani. Two problems with variational expectation maximisation for time-series models. In D. Barber, A. T. Cemgil, and S. Chiappa, editors, *Bayesian Time Series Models*. 2011.
- M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1-2):1–305, 2008.
- B. Wang and M. Titterton. Inadequacy of interval estimates corresponding to variational Bayesian approximations. In *Workshop on Artificial Intelligence and Statistics*, pages 373–380, 2004.
- Y. Wang and D. M. Blei. Frequentist consistency of variational Bayes. *arXiv preprint arXiv:1705.03439*, 2017.
- M. Welling and Y. W. Teh. Linear response algorithms for approximate inference in graphical models. *Neural Computation*, 16(1):197–221, 2004.
- T. Westling and T. H. McCormick. Establishing consistency and improving uncertainty estimates of variational inference through m-estimation. *arXiv preprint arXiv:1510.08151*, 2015.
- S. Wright and J. Nocedal. Numerical optimization. *Springer Science*, 35:67–68, 1999.
- H. Zhu, J. G. Ibrahim, S. Lee, and H. Zhang. Perturbation selection and influence measures in local influence analysis. *The Annals of Statistics*, 35(6):2565–2588, 2007.
- H. Zhu, J. G. Ibrahim, and N. Tang. Bayesian influence analysis: A geometric approach. *Biometrika*, 98(2):307–323, 2011.