

Variational Bayesian Approximation of Inverse Problems using Sparse Precision Matrices

Jan Povala^{a,c,1,*}, Ieva Kazlauskaitė^{b,1}, Eky Febrianto^{b,c}, Fehmi Cirak^{b,c}, Mark Girolami^{b,c}

^a*Department of Mathematics, Imperial College London, London, SW7 2AZ, UK*

^b*Department of Engineering, University of Cambridge, Cambridge, CB2 1PZ, UK*

^c*The Alan Turing Institute, London, NW1 2DB, UK*

Abstract

Inverse problems involving partial differential equations are widely used in science and engineering. Although such problems are generally ill-posed, different regularization approaches have been developed to ameliorate this problem. Among them is the Bayesian approach, where a prior probability measure is placed on the quantity of interest. The resulting posterior probability measure is usually analytically intractable, and approximation techniques must be used. The Markov Chain Monte Carlo (MCMC) method has been the go-to method for sampling from those posterior measures. This method, although asymptotically exact, is computationally infeasible for large-scale problems that arise in engineering practice. Variational Bayes (VB) was proposed as a more computationally tractable method for Bayesian inference, approximating a Bayesian posterior distribution with a simpler distribution by solving an optimization problem. In this work, we argue, through an empirical assessment, that VB methods are a flexible, fast, and scalable alternative to MCMC methods for this class of problems. We propose a natural choice of a family of trial distributions parametrised by precision matrices, thus taking advantage of the sparse structure encoded in the discretization of the problem. We utilize stochastic optimization to efficiently estimate the variational objective and assess not only the expected error in the mean solution but also the ability to quantify the uncertainty of the estimate. We test this on PDEs based on the Poisson equation in 1D and 2D, and we will make our Tensorflow implementation publically available.

Keywords: Inverse problems, Bayesian inference, variational Bayes, precision matrix, uncertainty quantification

1. Introduction

The increased availability of measurements from engineering systems allows for the development of new and the improvement of existing computational models, which are usually formulated as partial differential equations. Inferring model parameters (*e.g.*, material properties) from observations (*e.g.*, the strains) of the physical entities is termed the *inverse problem* (Tarantola, 2005; Kaipio and Somersalo, 2005; Stuart, 2010). In this work, we focus on the inverse problem where the quantities of interest (for example, some material properties) and the observations (*e.g.*, the displacement field) are related through elliptic PDEs. Most inverse problems are ill-posed, meaning that the existence, uniqueness, and/or stability (continuous dependence on the parameters) of the solution are violated

*Corresponding author

Email address: jan.povala@gmail.com (Jan Povala)

¹Equal contribution.

(Stuart, 2010; Kaipio and Somersalo, 2005). These issues are often alleviated through some regularization, like Tikhonov regularization (Tikhonov and Arsenin, 1977), that imposes assumptions on the regularity of the solution. Alternatively, the specification of the prior in the Bayesian formulation of inverse problems provides a natural choice for regularization, and any given regularization can be interpreted as a specific choice of priors in the Bayesian setting (Bishop, 2006). Furthermore, the Bayesian formulation provides not only a qualitative but also a quantitative estimate of the uncertainty in the solution. In particular, the mean of the posterior probability distribution corresponds to the point estimate of the solution while the credible intervals capture the range of the parameters consistent with the observed measurements and prior assumptions. For these reasons, Bayesian methods have gained popularity in computational mechanics for experimental design and inverse problems with uncertainty quantification (Abdulle and Garegnani, 2021; Pandita et al., 2021; Pyrialakos et al., 2021; Ni et al., 2021; Sabater et al., 2021; Huang et al., 2021; Ibrahimbegovic et al., 2020; Tarakanov and Elsheikh, 2020; Michelén Ströfer et al., 2020; Carlon et al., 2020; Wu et al., 2020; Uribe et al., 2020; Rizzi et al., 2019; Arnst and Soize, 2019; Beck et al., 2018; Betz et al., 2018; Chen et al., 2017; Asaadi and Heyns, 2017; Huang et al., 2017; Karathanasopoulos et al., 2017; Babuška et al., 2016)

The Bayesian formulation of inverse problems is the focal point of probabilistic machine learning, and in recent years significant progress has been made in adapting and scaling machine learning approaches to complex large-scale problems (Lu and Tang, 2015; Solin et al., 2018). One of the leading models for Bayesian inverse problems are Gaussian processes (GPs) which define probability distributions over functions and allow for ways to generalize from observed data. Given that most posterior distributions in Bayesian inference are analytically intractable, approximation methods need to be resorted to. Two classical approximation schemes are Markov Chain Monte Carlo (MCMC) and the Laplace approximation (LA). The MCMC algorithm proceeds by creating a Markov Chain whose stationary distribution is the desired posterior distribution. Although MCMC provides asymptotic convergence in distribution, devising an efficient, finite-time sampling scheme is challenging, especially in higher dimensions (Gelman et al., 2013). Application-specific techniques such as parameter space reduction and state space reduction have been proposed in the literature to help scale up MCMC methods, but these low-rank approximations are not specific to MCMC methods only (Cui et al., 2016). Due to the asymptotic correctness of MCMC, we use it as a benchmark for the experimental studies in this paper. Meanwhile, the Laplace approximation finds a Gaussian density centred around the mode of the true posterior, utilizing the negative Hessian of the unnormalised posterior log-density (Bishop, 2006). The Hessian is a large dense matrix, where forming each column requires multiple PDE solves; to make such calculation feasible, low-rank approximations are typically used (Villa et al., 2021; Bui-Thanh et al., 2013). Evidently, the Laplace approximation is not suitable for multi-modal posterior distributions due to the uni-modality of the Gaussian distribution.

1.1. Related work

In recent years, advances in variational Bayes (VB) methods have allowed for Bayesian inference to be successfully applied to large data sets. Variational Bayes translates a sampling problem that arises from applying the Bayes rule into an optimization problem (Jordan et al., 1999; Blei et al., 2017; Jordan and Wainwright, 2007). The method finds a solution that minimizes the Kullback-Leibler (KL) divergence between the true posterior distribution and a trial distribution from a chosen family of distributions, for instance, multivariate Gaussian distributions with a specific covariance structure. The strong appeal of VB is that one can explicitly choose the complexity of the trial distribution such that the resulting optimization problem is computationally tractable, and the approximate posterior adequately captures important aspects of the true posterior.

Further scalability of VB methods is due to advancements in sparse approximations and approximate inference. For instance, sparse GP methods such as Nyström approximation or fully independent training conditional method (FITC) rely on lower-dimensional representations that are defined by a smaller set of so-called inducing points to represent the full GP (Williams and Seeger, 2001; Csató and Oppé, 2002; Seeger et al., 2003; Quiñonero-Candela and Rasmussen, 2005; Snelson and Ghahramani, 2006; Titsias, 2009, 2008). Using this approximation for a data set of size N , algorithmic complexity is reduced from $\mathcal{O}(N^3)$ to $\mathcal{O}(NM^2)$, while storage demands go down from $\mathcal{O}(N^2)$ to $\mathcal{O}(NM)$, where M is a user selected number of inducing variables. To widen the applicability of VB to large datasets and non-conjugate models (combinations of prior distributions and likelihoods that do not result in a closed-form solution), *stochastic variational inference* (SVI) was proposed (Hensman et al., 2012; Hoffman et al., 2013; Hensman et al., 2013). Sub-sampling the original data and Monte Carlo estimation of the optimization objective and its gradients, allows for calibrating complex models using large amounts of data. Multiple further extensions to the sparse SVI framework were proposed, leveraging the Hilbert space formulation of VB (Cheng and Boots, 2017), introducing parametric approximations (Jankowiak et al., 2020), applying the Lanczos algorithm to efficiently factorize the covariance matrix (Pleiss et al., 2018), transforming to an orthogonal basis (Salimbeni et al., 2018; Shi et al., 2020), and adapting to compositional models (Salimbeni and Deisenroth, 2017).

The choice of prior is a central task in designing Bayesian models. If the prior is obtained from a domain expert, it is not necessarily less valuable than the data itself; one way of thinking about a prior is by considering how many observations one would be prepared to trade for a prior from an expert – if the expert is very knowledgeable, then one might be prepared to exchange a large part of a dataset to get access to that prior. Translating the expert knowledge into a prior probability distribution is a challenging task, and due to practical considerations, certain choices of priors are preferred for their simplicity and analytic tractability. When inferring values of parameters over a spatial domain, as is typically the case in finite elements, GP priors offer a natural way to incorporate the information about the smoothness and other known properties of the solution. We note that while other Bayesian models, such as Bayesian neural networks are gaining interest, it is very difficult to impose functional priors in such models, challenging the effective use of expert knowledge and leading to unrealistic uncertainty estimates (Sun et al., 2019; Burt et al., 2021).

1.2. Contributions

In this work, we advocate for the use of GP priors with stochastic variational inference as a principled and scalable way to solve the inverse problems arising in computational mechanics. We show, through an extensive empirical study, that variational Bayes methods provide a flexible, fast, and scalable alternative to MCMC methods in the context of Bayesian inverse problems based on elliptic PDEs while retaining the ability to quantify uncertainty. While similar directions have been explored in previous work, the focus there is on specific applications, such as parameter estimation problems in models of contamination (Tsilifis et al., 2016) or proof-of-concept on particular 1D inverse problems (Barajas-Solano and Tartakovsky, 2019).

We extend the previous works in multiple aspects, focusing on improving the utility of VB in inverse problems arising from elliptic PDEs and providing a thorough discussion of the empirical results that can be used by practitioners to guide their use of VB in applications. Specifically, we argue that the efficiency of the VB algorithms for PDE based inverse problems can be improved by taking in to account the structure of the problem, as encoded in the FEM discretization of the PDE. Motivated by previous uses of precision matrices as a way of describing conditional independence (Tan and Nott, 2018; Durrande et al., 2019), we leverage the sparse structure of the problems to impose conditional independence in the approximating posterior distribution. This choice of parametrisation

results in sparse matrices, which improve the computational and the memory cost of the resulting algorithms. Such parametrisation, combined with stochastic optimisation techniques, allows the method to be scaled up to large problems on 2D domains. Through extensive empirical comparisons, we demonstrate that VB provides high quality point estimates and uncertainty quantification comparable to the estimates attained by MCMC algorithms but with significant computational gains. Finally, we describe how the proposed framework can be seamlessly combined with existing solvers and optimization algorithms in the finite element implementations.

The main concern related to VB in statistics stems from the fact that it is constrained by the chosen family of trial distributions, which may not approximate the true posterior distribution well. If the choice of the trial distributions is too restrictive, the estimate of the posterior mean is biased while the uncertainty may be underestimated (MacKay, 2003; Wang and Titterton, 2005; Turner and Sahani, 2011). Furthermore, as noted in previous work, the commonly used mean-field factorization of the trial distributions does not come with general guarantees on accuracy (Giordano et al., 2018). However, VB has been demonstrated to work well in practice in a variety of settings (Kingma and Welling, 2014; Damianou et al., 2016; Blei et al., 2017; Zhang et al., 2019). Recent work on VB has provided some tools for assessing the robustness of the VB estimates (Giordano et al., 2018).

1.3. Overview

The rest of the paper is structured as follows. In Section 2, we define Bayesian inverse problems and detail some inference challenges related to their ill-posedness. In Section 3, we give a presentation of the variational inference framework, with strong focus on sparse parametrisation resulting from conditional independence. We give details of the experiments and the evaluation criteria, and discuss obtained results for each experiment in Section 4. Lastly, Section 5 concludes the paper and discusses some promising directions for future work.

2. Bayesian Formulation of Inverse Problems

2.1. Forward map and observation model

We closely follow the formulation by Stuart (2010). We are interested in finding $\kappa \in \mathcal{K}$, an input to a model, given $y \in \mathcal{Y}$, a noisy observation of the solution of the model, where \mathcal{K}, \mathcal{Y} are Banach spaces². We write the mapping as

$$y = \mathcal{G}(\kappa) + \eta, \quad (1)$$

where $\mathcal{G} : \mathcal{K} \rightarrow \mathcal{Y}$, $\eta \in \mathcal{Y}$ is additive observational noise. We focus on problems where \mathcal{G} maps solutions of elliptic partial differential equations with input $\kappa \in \mathcal{K}$ into the observation space \mathcal{Y} . For a suitable Hilbert space \mathcal{U} which we make concrete later, let $\mathcal{A} : \mathcal{K} \rightarrow \mathcal{U}$ be a possibly non-linear solution operator of the PDE. For a particular $\kappa \in \mathcal{K}$, the solution is

$$u = \mathcal{A}(\kappa). \quad (2)$$

To obtain observations y , we define a projection operator $\mathcal{P} : \mathcal{U} \rightarrow \mathcal{Y}$. Consequently, Eq. (1) can be written out in full as

$$y = \mathcal{P}(\mathcal{A}(\kappa)) + \eta, \quad (3)$$

²Respective norms for Banach spaces \mathcal{K}, \mathcal{Y} are $\|\cdot\|_{\mathcal{K}}$ and $\|\cdot\|_{\mathcal{Y}}$.

2.2. Inference

We solve Eq. (1) for κ by finding κ such that the data misfit norm, $\|y - \mathcal{G}(\kappa)\|_{\mathcal{Y}}$, is minimised. As already mentioned in the introduction, this is typically an ill-posed problem: there may be no solution, it may not be unique, and it may depend sensitively on y . To proceed, we choose the Bayesian framework of regularising the problem to make it amenable to analysis and practical implementation. We describe our prior knowledge about κ in terms of a prior probability measure μ_0 on the subspace of \mathcal{K} and use Bayes' formula to calculate the posterior probability measure, μ^y , for κ given y . The relationship between the posterior and prior is expressed as

$$\frac{d\mu^y}{d\mu_0}(\kappa) = \frac{1}{Z(y)} \exp(-\Phi(\kappa; y)), \quad (4)$$

where $\frac{d\mu^y}{d\mu_0}$ is the Radon-Nikodym derivative of μ^y with respect to μ_0 , and Φ is the potential function which is determined by the forward problem (1), specifically \mathcal{G} and η , and $Z(y) = \int_{\mathcal{K}} \exp(\kappa; y) d\mu_0(\kappa)$ so that μ^y is a valid probability measure.

From here on, we assume that $(\mathcal{Y}, \|\cdot\|_{\mathcal{Y}} = (\mathbb{R}^{n_y}, \|\cdot\|))$, where $\|\cdot\|$ is the Euclidean norm, and we treat data y and η as vectors, i.e. \mathbf{y} and $\boldsymbol{\eta}$. We specify the additive noise vector $\boldsymbol{\eta}$ as Gaussian such that

$$\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \Gamma = \sigma_y^2 \mathbb{I}),$$

where σ_y is the standard deviation of the measurement noise and \mathbb{I} is the identity matrix. We can write Φ conveniently as

$$\Phi(\kappa; \mathbf{y}) = \frac{1}{2} \|\mathcal{G}(\kappa) - \mathbf{y}\|_{\Gamma^{-1}}, \quad (5)$$

where $\|\cdot\|_{\Gamma^{-1}}$ is the norm induced by the weighted inner product³.

We restrict the space of solutions \mathcal{K} to be a Hilbert space and place a centred Gaussian prior measure on κ with covariance operator \mathcal{C}_κ :

$$\mu_0(\kappa) \sim \mathcal{N}(0, \mathcal{C}_\kappa), \quad (6)$$

which can be written out in the weighted norm form as follows:

$$d\mu_0(\kappa) \propto \exp \left\{ -\frac{1}{2} \|\kappa\|_{\mathcal{C}_\kappa^{-1}} \right\}. \quad (7)$$

Applying Bayes' rule from Eq. (4), the posterior is given as

$$d\mu^y \propto \exp \left\{ -\frac{1}{2} \|\mathcal{G}(\kappa) - \mathbf{y}\|_{\Gamma^{-1}} - \frac{1}{2} \|\kappa\|_{\mathcal{C}_\kappa^{-1}} \right\}. \quad (8)$$

For detailed assumptions on μ_0 , \mathcal{G} , and η that are required for deriving the posterior probability measure, we refer the reader to [Stuart \(2010, Sec. 2.4\)](#).

³For any self-adjoint positive operator \mathcal{T} , weighted inner product is $\langle \cdot, \cdot \rangle_{\mathcal{T}} = \langle \mathcal{T}^{-1/2} \cdot, \mathcal{T}^{-1/2} \cdot \rangle$, and the induced norm is $\|\cdot\|_{\mathcal{T}} = \|\mathcal{T}^{-1/2} \cdot\|$

2.2.1. Algorithms

The objective is to find the posterior measure μ^y conditioned on the observations, as dictated by Bayes's rule. The forward map (1) and the respective functions must be discretised. In Bayesian inference there are two approaches for discretisation: 1) apply the Bayesian methodology first, discretize afterwards, or 2) discretize the first, then apply the Bayesian methodology (Stuart, 2010).

The first approach develops the solution of the inference problem in the function space before discretising it. A widely used algorithm of this form is the pre-conditioned Crank-Nicholson (pCN) Markov chain Monte Carlo scheme, where proposals are based on the measure μ_0 and the current state of the Markov chain. The pCN method is a standard choice for high-dimensional sampling problems, as its implementation is well-defined and is invariant to mesh refinement (Cotter et al., 2013; Pinski et al., 2015). Since we will use this algorithm as one of the baselines, we provide a summary of the algorithm in Sec. Appendix C.1. Recently, infinite-dimensional MCMC schemes that leverage the geometry of the posterior to improve the efficiency have been proposed, see Beskos et al. (2017). Other than MCMC schemes, some variational Bayes formulations in function space have been proposed (for example, Minh (2017); Burt et al. (2021)), though currently they do not offer a viable computational alternative to the finite-dimensional formulation of variational inference.

The second approach proceeds by first discretizing the problem and then deriving the solution of the inference method. This approach forms the basis of almost all inference procedures developed in engineering: MCMC algorithms such as Metropolis-Hastings (Metropolis et al., 1953) or Hamiltonian Monte Carlo (HMC) (Duane et al., 1987), the Laplace approximation, or variational Bayes (Jordan et al., 1998, 1999) are used to approximate the posterior. In the discretised formulation, HMC has achieved recognition as the *gold standard* for its good convergence properties, favourable performance on high-dimensional and poorly conditioned problems, and universality of implementation that enables its generic use in many applications through probabilistic programming languages (*e.g.*, Stan (Carpenter et al., 2017)). Therefore, along with the pCN scheme mentioned above, our baseline for inference methods includes the HMC method, and we provide a summary of the HMC scheme in Sec. Appendix C.2.

For the rest of the exposition in this paper, we will focus on algorithms in the finite-dimensional case, where we discretise κ to a vector $\boldsymbol{\kappa}$. In finite dimensions, probability densities with respect to the Lebesgue measure can be defined, thus leading to a more familiar form of the Bayes's rule:

$$p(\boldsymbol{\kappa} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\kappa}) p(\boldsymbol{\kappa})}{p(\mathbf{y})} \propto p(\mathbf{y} | \boldsymbol{\kappa}) p(\boldsymbol{\kappa}), \quad (9)$$

where $p(\boldsymbol{\kappa} | \mathbf{y})$ is the posterior density, $p(\mathbf{y} | \boldsymbol{\kappa})$ is the likelihood of the observed data for a given discretised $\boldsymbol{\kappa}$ and is determined by the discretised forward problem (1) and noise $\boldsymbol{\eta}$. The prior density for $\boldsymbol{\kappa}$, which itself may depend on some (hyper-) parameters $\boldsymbol{\psi}$, is denoted by $p(\boldsymbol{\kappa})$. Next two sections focus on discussing $p(\mathbf{y} | \boldsymbol{\kappa})$ and $p(\boldsymbol{\kappa})$, respectively.

2.3. Poisson Equation and likelihood

Let us consider a specific forward problem where u is the solution to the Poisson equation with prescribed boundary conditions:

$$-\nabla \cdot (\exp(\kappa(\mathbf{x})) \nabla u(\mathbf{x})) = f(\mathbf{x}), \quad (10)$$

where $\mathbf{x} \in \Omega \subset \mathbb{R}^d$, with $d \in \{1, 2, 3\}$, $\kappa(\mathbf{x}) \in \mathbb{R}$ is the log-diffusion coefficient, $u(\mathbf{x}) \in \mathbb{R}$ is the unknown, and $f(\mathbf{x}) \in \mathbb{R}$ is a deterministic forcing term. We are given the $\mathbf{y} \in \mathbb{R}^{n_y}$ noisy observations of the solution u at a finite set of points, $\{\mathbf{x}_i\}_{i=1}^{n_y}$. The observation points $\mathbf{x}_i \in \Omega$ are collected in the matrix $\mathbf{X} \in \mathbb{R}^{n_y \times d}$. Although this PDE is linear in u for a given κ , the methodology in this

paper applies to non-linear cases and also for time-dependent cases such as the inverse problem of inferring initial conditions of a system given observations of the system at a later time.

We discretize the weak form of the Poisson (10) with a standard finite element approach. Specifically, the domain of interest Ω is subdivided into a set $\{\omega_e\}$ of non-overlapping elements of size $h = \max_e \text{diam}(\omega_e)$ such that:

$$\Omega = \bigcup_{e=1}^{n_e} \omega_e. \quad (11)$$

The unknown field $u(x)$ is approximated with Lagrange basis functions $\phi_i(x)$ and the respective nodal coefficients $\mathbf{u} = [\mathbf{u}_1, \dots, \mathbf{u}_{n_u}]^\top$ of the n_u non-Dirichlet boundary mesh nodes by

$$u_h(\mathbf{x}) = \sum_{i=1}^{n_u} \phi_i(\mathbf{x}) \mathbf{u}_i. \quad (12)$$

The discretization of the weak form of the Poisson equation yields the linear system of equations

$$\mathbf{A}(\boldsymbol{\kappa}) \mathbf{u} = \mathbf{f}, \quad (13)$$

where $\mathbf{A}(\boldsymbol{\kappa}) \in \mathbb{R}^{n_u \times n_u}$ is the stiffness matrix, $\boldsymbol{\kappa} \in \mathbb{R}^{n_\kappa}$ is the vector of log-diffusion coefficients, $\mathbf{f} \in \mathbb{R}^{n_u}$ is the nodal source vector. The stiffness matrix of an element with the label e is given by

$$A_{ij}^e(\boldsymbol{\kappa}_e) = \int_{\omega_e} \exp(\kappa_e) \frac{\partial \phi_i(\mathbf{x})}{\partial \mathbf{x}} \cdot \frac{\partial \phi_j(\mathbf{x})}{\partial \mathbf{x}} d\mathbf{x}, \quad (14)$$

where the log-diffusion coefficient κ_e of the element is assumed to be *constant* within the element. The source vector is discretised as:

$$\mathbf{f}_i(\kappa_e) = \int_{\Omega} f(\mathbf{x}) \phi_i(\mathbf{x}) d\mathbf{x}. \quad (15)$$

Hence, according to the observation model (5) the likelihood is given by

$$p(\mathbf{y} \mid \boldsymbol{\kappa}) = p(\mathbf{y} \mid \mathbf{u}_{\boldsymbol{\kappa}}) = \mathcal{N}(\mathbf{P}\mathbf{A}(\boldsymbol{\kappa})^{-1}\mathbf{f}, \sigma_y^2 \mathbb{I}), \quad (16)$$

where the matrix \mathbf{P} represents the discretisation of the observation operators \mathcal{P} .

Then the mapping from the coefficients $\boldsymbol{\kappa}$ to the solution \mathbf{u} is $\mathbf{u} = \mathbf{A}(\boldsymbol{\kappa})^{-1}\mathbf{f}$. The unconditional distribution of \mathbf{u} is:

$$p(\mathbf{u}) = \int p(\mathbf{u} \mid \boldsymbol{\kappa}) p(\boldsymbol{\kappa}) d\boldsymbol{\kappa}, \quad (17)$$

where $p(\mathbf{u} \mid \boldsymbol{\kappa})$ is deterministic as defined in (13) but $\boldsymbol{\kappa}$ appears in it non-linearly, implying that the inference is not analytically tractable.

Throughout the experiments in the later sections, we either set Dirichlet (essential) boundary conditions everywhere (for example $u(\mathbf{x}) = 0$ on $\partial\Omega$), or assume Neumann (natural) boundary conditions on parts of the boundary. The choice will be made explicit in each experiment. To compute the likelihood, we solve (10) for $u(\mathbf{x})$ using the finite element method (FEM).

2.4. Prior

As discussed above, we place a centred Gaussian measure on κ , $\mu_0(\kappa) \sim \mathcal{N}(0, \mathcal{C}_\kappa)$. Properties of sample paths depend on the mean function which we assume to be zero and on the spectral

properties of \mathcal{C}_κ . We restrict the space of prior functions to $L^2(\Omega; \mathbb{R})$. Then, covariance operator \mathcal{C}_κ can be constructed from the covariance function, $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[\kappa(\mathbf{x})\kappa(\mathbf{x}')] as:$

$$(\mathcal{C}_\kappa \phi)(\mathbf{x}) = \int_{\Omega} k(\mathbf{x}, \mathbf{x}') \phi(\mathbf{x}') d\mathbf{x}'. \quad (18)$$

This formulation is what is commonly referred to as a Gaussian process (GP) with mean function $m(\cdot)$, which we assume to be zero, and covariance function $k(\cdot, \cdot)$:

$$\kappa \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot)). \quad (19)$$

The index set of the two functions is infinite which allows GPs to be interpreted as non-parametric priors over the space of functions. However, even though the process is infinite-dimensional, an instantiation of the process is finite and reduces to a multivariate Gaussian distribution by definition. The covariance function is typically parametrised by a set of hyperparameters ψ . One popular option, which satisfies assumptions about μ_0 as per [Stuart \(2010\)](#), is the squared exponential kernel (also called the exponentiated quadratic or the radial basis function (RBF) kernel):

$$k_{\text{SE}}(\mathbf{x}, \mathbf{x}') = \sigma_\kappa^2 \exp\left(-\frac{r^2}{2\ell_\kappa^2}\right), \quad (20)$$

where $r = \|\mathbf{x} - \mathbf{x}'\|_2$ is the Euclidean distance between the inputs. It depends on two parameters $\psi = \{\sigma_\kappa, \ell_\kappa\}$, the scaling parameter σ_κ , and the length-scale ℓ_κ . Note that, $k_{\text{SE}}(\cdot, \cdot)$ is an infinitely smooth function, which implies that so is $\kappa(\cdot)$. The RBF kernel imposes smoothness and stationarity assumptions on the solution; in addition, such choice of kernel offers a way to regularise the resulting optimisation problem. However, depending on the expert knowledge of the true solution, other kernels may be used to impose other assumptions such as periodicity.

Both conditioning and marginalization can be done in closed form. In particular, consider the joint model of the values κ at training locations \mathbf{X} and the test values κ^* at test locations \mathbf{X}^* :

$$\begin{bmatrix} \kappa \\ \kappa^* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_\psi(\mathbf{X}, \mathbf{X}) & \mathbf{K}_\psi(\mathbf{X}, \mathbf{X}^*) \\ \mathbf{K}_\psi(\mathbf{X}^*, \mathbf{X}) & \mathbf{K}_\psi(\mathbf{X}^*, \mathbf{X}^*) \end{bmatrix}\right), \quad (21)$$

where $\mathbf{K}_\psi(\mathbf{X}, \mathbf{X}^*)$ is the matrix resulting from evaluating $k(\cdot, \cdot)$ at all pairs of training and test points. The conditional distribution of the function values κ^* given the values κ at \mathbf{X} is:

$$\begin{aligned} \kappa^* | \kappa &\sim \mathcal{N}\left(\tilde{\kappa}^*, \tilde{\mathbf{K}}\right), \text{ where} \\ \tilde{\kappa}^* &= \mathbf{K}(\mathbf{X}^*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X})]^{-1} \kappa \\ \tilde{\mathbf{K}} &= \mathbf{K}(\mathbf{X}^*, \mathbf{X}^*) - \mathbf{K}(\mathbf{X}^*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X})]^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}^*). \end{aligned} \quad (22)$$

The marginal distribution can be recovered by finding the relevant part of the covariance matrix; for example, the marginal of κ given \mathbf{X} is $\kappa \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_\psi(\mathbf{X}, \mathbf{X}))$.

In this work, we place a Gaussian process prior on $\kappa(\mathbf{x})$ and assume the squared exponential kernel with length-scale ℓ_κ and fixed variance $\sigma_\kappa^2 = 1$. As mentioned in the previous section, we assume that $\kappa(\mathbf{x})$ is constant on each element of the mesh (we use the same mesh as for discretizing $u(\mathbf{x})$ and $f(\mathbf{x})$). We place the prior on κ so that the centroids of the elements are the training points of the GP:

$$p(\kappa) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_\psi(\mathbf{X}, \mathbf{X})). \quad (23)$$

3. Variational Bayes Approximation

3.1. Variational Bayes

We assume that any hyper-parameters ψ of the prior are fixed, and are only interested in the posterior distribution of $\boldsymbol{\kappa}$. The variational approach proceeds by approximating the true posterior $p(\boldsymbol{\kappa} \mid \mathbf{y})$ according to (9) with a density $q(\boldsymbol{\kappa})$, which is the minimizer of the discrepancy between a chosen family of trial densities \mathcal{D}_q and the true posterior distribution $p(\boldsymbol{\kappa} \mid \mathbf{y})$ (Jordan et al., 1999; Jordan and Wainwright, 2007). A typical choice for the measure of discrepancy between distributions is the Kullback-Leibler (KL) divergence (which due to the lack of symmetry is not a metric). To find the approximate posterior distribution we have:

$$q^*(\boldsymbol{\kappa}) = \arg \min_{q(\boldsymbol{\kappa}) \in \mathcal{D}_q} \text{KL}(q(\boldsymbol{\kappa}) \parallel p(\boldsymbol{\kappa} \mid \mathbf{y})). \quad (24)$$

Expanding the KL divergence term we obtain

$$\begin{aligned} \text{KL}(q(\boldsymbol{\kappa}) \parallel p(\boldsymbol{\kappa} \mid \mathbf{y})) &= \int q(\boldsymbol{\kappa}) \log \frac{q(\boldsymbol{\kappa})}{p(\boldsymbol{\kappa} \mid \mathbf{y})} d(\boldsymbol{\kappa}) \\ &= \mathbb{E}_q [\log q(\boldsymbol{\kappa})] - \mathbb{E}_q [\log p(\boldsymbol{\kappa} \mid \mathbf{y})] \\ &= \mathbb{E}_q [\log q(\boldsymbol{\kappa})] - \mathbb{E}_q \left[\log \frac{p(\mathbf{y}, \boldsymbol{\kappa})}{p(\mathbf{y})} \right] \\ &= \mathbb{E}_q [\log q(\boldsymbol{\kappa})] - \mathbb{E}_q [\log p(\mathbf{y}, \boldsymbol{\kappa})] + \log p(\mathbf{y}) \end{aligned} \quad (25)$$

The last term of the KL divergence, the log-marginal likelihood, is usually not known. However, we use the fact that the KL divergence is non-negative to obtain the bound

$$\log p(\mathbf{y}) \geq \mathbb{E}_q [\log p(\mathbf{y}, \boldsymbol{\kappa})] - \mathbb{E}_q [\log q(\boldsymbol{\kappa})]. \quad (26)$$

This inequality becomes an equality when the trial density $q(\boldsymbol{\kappa})$ and the posterior $p(\boldsymbol{\kappa} \mid \mathbf{y})$ are equal. To minimize the KL divergence, it is sufficient to maximise $\mathbb{E}_q [\log p(\mathbf{y}, \boldsymbol{\kappa})] - \mathbb{E}_q [\log q(\boldsymbol{\kappa})]$, which is commonly referred to as the evidence lower bound (ELBO). The ELBO term can be rewritten as

$$\begin{aligned} \text{ELBO}(q) &= \mathbb{E}_q [\log p(\mathbf{y} \mid \boldsymbol{\kappa}) + \log p(\boldsymbol{\kappa})] - \mathbb{E}_q [\log q(\boldsymbol{\kappa})] \\ &= \mathbb{E}_q [\log p(\mathbf{y} \mid \boldsymbol{\kappa})] - \text{KL}(q(\boldsymbol{\kappa}) \parallel p(\boldsymbol{\kappa})). \end{aligned} \quad (27)$$

To summarize, the task now becomes:

$$q^*(\boldsymbol{\kappa}) = \arg \max_{q(\boldsymbol{\kappa}) \in \mathcal{D}_q} \mathbb{E}_q [\log p(\mathbf{y} \mid \boldsymbol{\kappa})] - \text{KL}(q(\boldsymbol{\kappa}) \parallel p(\boldsymbol{\kappa})). \quad (28)$$

To maximize the ELBO with a gradient-based optimiser, we need to be able to evaluate it and its gradients with respect to the parameters of $q(\boldsymbol{\kappa})$. The ELBO is usually non-convex which makes its optimization particularly challenging. Although the KL divergence term of the ELBO is often available in closed form, $\mathbb{E}_q [\log p(\mathbf{y} \mid \boldsymbol{\kappa})]$ involving the likelihood is generally not available. It can be approximated using a Monte Carlo approximation with N_{SVI} samples from the trial density $q(\boldsymbol{\kappa})$ as follows:

$$\mathbb{E}_q [\log p(\mathbf{y} \mid \boldsymbol{\kappa})] \approx \frac{1}{N_{\text{SVI}}} \sum_{i=1}^{N_{\text{SVI}}} \log p(\mathbf{y} \mid \boldsymbol{\kappa}^{(i)}), \quad (29)$$

where $\boldsymbol{\kappa}^{(i)}$ is the i -th sample from $q(\boldsymbol{\kappa})$. Our empirical tests show that the value of N_{SVI} in the range of 2–5 provides fast convergence of the optimisation, agreeing with the previous literature (Kingma and Welling, 2014). This approach is often referred to as stochastic variational inference (SVI). The Monte Carlo approximation is in line with the work in Barajas-Solano and Tartakovsky (2019) but in contrast with the analytic approximation based on the Hessian calculations proposed in Tsilifis et al. (2016).

3.2. Specification of trial distribution

The specification of the approximating family of distributions determines how much structure of the true posterior distribution is captured by the variational approximation. To model complex relationships between the components of the posterior, a more complex approximating family of distributions is needed. As the richer family of distributions is likely to require more parameters, the optimization of the usually non-convex ELBO becomes harder. A balance must be struck in this trade-off: the family should be rich enough, but the optimization task should still be computationally tractable.

A practical and widely used variational family is the multivariate Gaussian distribution, parametrised by the mean vector and the covariance matrix. One of the key benefits of this choice is that the KL divergence term of the ELBO in (27) is available in closed form for a GP prior. The choice of the parametrisation of the covariance matrix determines how much structure, other than the mean estimate, is captured by the variational family. We discuss this in more detail in the next section.

Numerous approaches have been proposed to extend the trial distribution beyond the Gaussian family. A standard approach in situations when the true posterior distribution is likely to be multimodal is to consider mixtures of variational densities (Bishop et al., 1998). A more recent development is embedding parameters of a mean-field approximation in a hierarchical model to induce variational dependencies between latent variables (Tran et al., 2015; Ranganath et al., 2016).

3.2.1. Gaussian trial distribution

Choosing the trial distribution $q(\boldsymbol{\kappa})$ as a multivariate Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ requires optimization over the mean $\boldsymbol{\mu}$ (free-form) and the covariance matrix $\boldsymbol{\Sigma}$. The flexibility in choosing how we specify both of these parameters, especially the covariance matrix, enables us to balance the trade-off between the expressiveness of the approximating distribution and the computational efficiency.

The richest specification corresponds to parametrising the covariance matrix $\boldsymbol{\Sigma}$ using its full Cholesky factor \mathbf{L} , i.e.,

$$q(\boldsymbol{\kappa}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^\top). \quad (30)$$

This choice results in a dense covariance matrix that may be able to capture the full covariance structure between the inputs (*i.e.* each input may be correlated with every other input). Parametrising the components of \mathbf{L} automatically ensures that the covariance matrix $\boldsymbol{\Sigma}$ is positive definite as necessary. The number of parameters to optimize grows as $\mathcal{O}(n_{\boldsymbol{\kappa}}^2)$ and this leads to a difficult optimization task that needs to be carefully initialized and parametrised. We refer to this parametrisation as full-covariance variational Bayes (FCVB).

A much more efficient choice is a diagonal covariance matrix, which is often referred to as mean-field variational Bayes (MFVB). By limiting the number of parameters that need to be optimized, the optimization task becomes simpler and the number of parameters grows only as $\mathcal{O}(n_{\boldsymbol{\kappa}})$. While more computationally efficient and easier to initialize, MFVB ignores much of the dependence structure of the posterior distribution.

3.3. Conditional Independence, Sparse Precision Matrices

Instead of parametrising the covariance matrix $\boldsymbol{\Sigma}$, or its Cholesky decomposition \mathbf{L} in physical systems it is often advantageous to parametrise the precision matrix, \mathbf{Q} , where $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$. While a component of the covariance matrix $\boldsymbol{\Sigma}$ expresses *marginal* dependence between the two corresponding random variables, the elements of the precision matrix reflect their *conditional independence* (Rue et al., 2009). Or, more specifically, for two components κ_i and κ_j of the random vector $\boldsymbol{\kappa}$ we note

$$p(\kappa_i, \kappa_j) = p(\kappa_i)p(\kappa_j) \quad \Leftrightarrow \quad \Sigma_{ij} = 0, \quad (31)$$

where Σ_{ij} denotes the respective component of Σ . Furthermore, defining the vector $\kappa_{-\{i,j\}}$ from the random vector κ by removing its i -th and j -th component, we note

$$p(\kappa_i, \kappa_j \mid \kappa_{-\{i,j\}}) = p(\kappa_i \mid \kappa_{-\{i,j\}})p(\kappa_j \mid \kappa_{-\{i,j\}}) \Leftrightarrow Q_{ij} = 0. \quad (32)$$

That is, $Q_{ij} = 0$ if and only if κ_i is independent from κ_j , *conditional* on all other components of κ .

A succinct way to represent conditional independence is using an undirected graph whose nodes correspond to the random variables (Bishop, 2006). A graph edge is present between two graph vertices i and j if the corresponding random variables are *not* conditionally independent from each other, given all the other random variables. Or, expressed differently, the edges between the graph vertices correspond to non-zeros in the precision matrix. In our context, each graph vertex represents a finite element and graph edges are introduced according to geometric adjacency of the finite elements as determined by the mesh. To this end, we define the 1-neighbourhood of a finite element as the union of the element itself and of elements sharing a node with the element. The n -neighbourhood is defined recursively as the union of all 1-neighbourhoods of all the elements in the $(n - 1)$ -neighbourhood. We introduce an edge between two graph vertices when the respective elements are in the same n -neighbourhood.

Figure 1 shows examples of adjacency graphs and the structure of the corresponding precision matrices \mathbf{Q} for 5 random variables resulting from a discretization of a 1D domain with 5 finite elements. In the considered examples the random variables represent the constant log-diffusion coefficient in the elements. As shown in Figures 1b and 1c choosing a larger n -neighbourhood for graph construction leads to a denser precision matrix. For instance, from the structure of the precision matrix in Figure 1b, which assumes a 1-neighbourhood structure, we can read for the log-diffusion coefficient of element j the following conditional independence relationship:

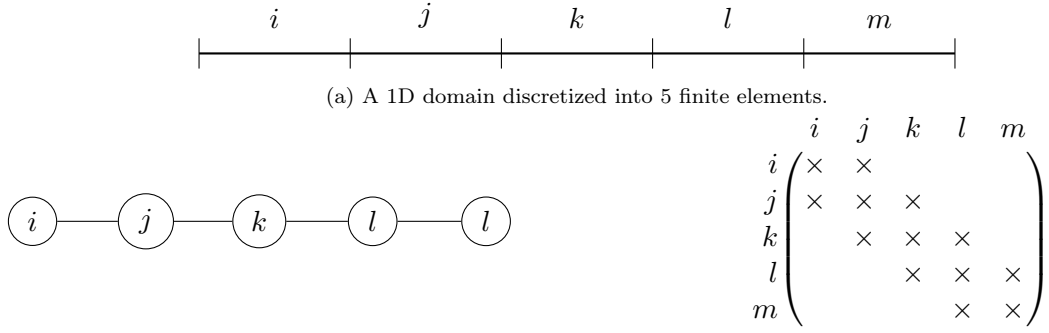
$$Q_{ik} = 0 \wedge Q_{il} = 0 \wedge Q_{im} = 0 \Rightarrow p(\kappa_i \mid \kappa_j, \kappa_k, \kappa_l, \kappa_m) = p(\kappa_i \mid \kappa_j), \quad (33)$$

where \wedge is the logical *and* operator. That is, when the coefficient of element j is given, the coefficient of the neighbouring element i is independent from all the remaining coefficients. This is intuitively plausible and in line with physical observations. Clearly, the covariance matrices corresponding to the given sparse precision matrices are dense. Hence, in the considered case the coefficient of element i may still be correlated to the coefficient of element m , i.e. $p(\kappa_i \mid \kappa_m) \neq p(\kappa_i)$. This correlation will most likely be relatively weak given the large distance between the two elements, but knowing the coefficient of element m will certainly restrict the range of possible values for the coefficient of element i .

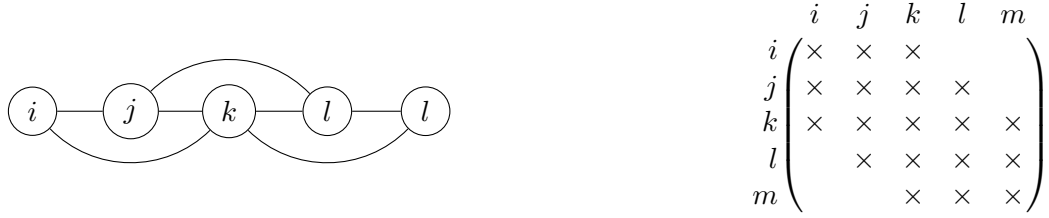
After obtaining the structure of the precision matrix, which is sparse but, in general, not banded, one can reorder the numbering of the elements in the finite element mesh to reduce its bandwidth. This allows for efficient linear algebra operations. See Cuthill and McKee (1969) for an example of a reordering algorithm. Once a minimum bandwidth ordering with b_{\min} has been established, we use the property that the bandwidth of the Cholesky factor \mathbf{L}_Q of matrix \mathbf{Q} is less than or equal to the bandwidth of \mathbf{Q} (Rue and Held, 2005). Finally, the parameters we optimize are the components of the lower band of size b_{\min} of matrix \mathbf{L}_Q , so that the approximating distribution reads

$$q(\kappa) \sim \mathcal{N}\left(\mu, (\mathbf{L}_Q \mathbf{L}_Q^\top)^{-1}\right). \quad (34)$$

This process of devising a parametrisation for the precision matrix for a more complex mesh in 2D is shown in Fig. 2. This approach is computationally efficient – the number of parameters grows as $\mathcal{O}(n_\kappa)$ – but also captures dependencies between all the random variables.



(b) Adjacency graph and the corresponding adjacency matrix based on 1-neighbourhood structure: there is an edge between two graph vertices if the corresponding elements share a node.



(c) Adjacency graph and the corresponding adjacency matrix based on 2-neighbourhood structure: there is an edge between two vertices if the corresponding elements are in each others 2-neighbourhoods.

Figure 1: An example of a 1D bar discretized into 5 elements and two different ways of assuming conditional independence.

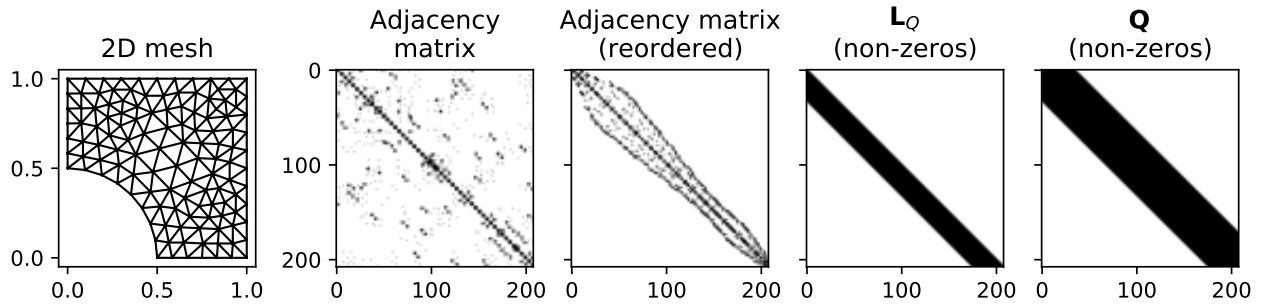


Figure 2: An illustration of how the sparse precision matrix parametrisation in Eq. (34) is derived. The adjacency matrix encodes the 2-neighborhood structure of the elements in the mesh. By reordering the nodes one can obtain a banded adjacency matrix, which is then used to parametrise the Cholesky factor of the precision matrix, as described in section 3.2.1.

3.4. Stochastic optimization

To maximize the ELBO in (28), we use the ADAM algorithm (Kingma and Ba, 2015). ADAM is a member of a larger class of stochastic optimization methods that have become popular as tools for maximizing non-convex cost functions. These methods construct a stochastic estimate of the gradient to perform gradient descent-based optimization. In the most basic form, the simplest method is stochastic gradient descent (SGD) which uses an unbiased gradient to perform standard gradient descent (Robbins and Monro, 1951). Due to their usability in the stochastic optimization of neural networks, there have been significant advances in the SGD methodology, utilizing momentum, preconditioning, and various other components to accelerate the convergence of the optimization procedure. ADAM, a stochastic gradient descent algorithm with an adaptive step size is one popular algorithm that exhibits a stable behaviour on many problems and is easy to use without significant tuning. The algorithm uses a per-parameter step size which is based on the first two moments of the estimate of the gradient for each parameter. Specifically, the step size is proportional to the ratio of the exponential moving average of the 1st moment to the square root of the exponential moving average of the non-centred 2nd moment. At any point, the exponential moving average is computed with decay parameters β_1 and β_2 for the 1st and 2nd moment, respectively. We adopt the parameter values suggested by the authors: $\beta_1 = 0.9$ and $\beta_2 = 0.99$. The speed of convergence is further controlled by the learning parameter α which is used to regulate the step size for all parameters in the same way. In our experiments, we set it to 0.01 and let it decay exponentially every 2,500 steps (1,000 for MFVB), with the decay rate of 0.96. While the ADAM algorithm performs well on a variety of problems, it has been shown that the convergence of this algorithm is poor on some problems (Reddi et al., 2018). We discuss alternative approaches as potential future work in Sec. 5.

To monitor convergence, we use a rule that tracks an exponentially weighted moving average of the decrease in the loss values between successive steps, and stops when that average drops below a threshold. The use of such an adaptive rule gives us a way to track the convergence of the algorithm and provides a conservative estimate for the time it takes for the optimization to converge. This rule can be adapted based on the available computational budget.

3.5. The algorithm

The maximization of the ELBO in (27) involves finding the parameters of the trial distribution $q(\boldsymbol{\kappa})$, i.e. its mean $\boldsymbol{\mu}$ and Cholesky factor \mathbf{L}_Q , that minimize KL between $q(\boldsymbol{\kappa})$ and the posterior $p(\boldsymbol{\kappa}|\mathbf{y})$. Algorithm 1 shows the required steps to compute the ELBO and its gradients with respect to the parameters of the trial distribution. Different from the discussion so far, in Algorithm 1 it is assumed that there are multiple independent observation vectors \mathbf{y}_i with $i \in \{1, 2, \dots, N_{\mathbf{y}}\}$.

4. Examples

We evaluate the efficacy of variational inference on (a) a 1D and 2D Poisson equation examples and (b) a benchmark proposed by Aristoff and Bangerth (2021). We discretize the examples with a standard finite element method using linear Lagrange basis functions as described in the previous section. We compare against two sampling-based inference schemes, Hamiltonian Monte Carlo (HMC) and pre-conditioned Crank-Nicholson Markov Chain Monte Carlo (pCN); these are known to be asymptotically correct as the number of samples increases. The evaluation criteria we use focus on three aspects of an inference scheme: the accuracy with respect to capturing the mean and the variance of the solution; propagation of uncertainty in derived quantities of interest; and the time until convergence of the solution.

Algorithm 1: ELBO estimation and its gradient with respect to the parameters of the trial distribution.

Input: Current parameters μ and L_Q of $q(\kappa)$
Output: ELBO and its gradients with respect to the parameters of $q(\kappa)$

- 1 Sample $[\kappa^{(1)}, \kappa^{(2)}, \dots, \kappa^{(N_{\text{SVI}})}]$ from $q(\kappa)$
- 2 **for each** $\kappa^{(i)}$ **do**
- 3 Solve for $\mathbf{u}(\kappa^{(i)})$ and obtain gradients with respect to κ using the FEM
- 4 $p(y | \kappa^{(i)}) \leftarrow \prod_{j=1}^{N_y} p(\mathbf{y}_j; \mathbf{u}(\kappa^{(i)}), \sigma_y^2)$ and propagate its gradient with respect to $\kappa^{(i)}$
- 5 ELBO $\leftarrow N_{\text{SVI}}^{-1} \sum_{i=1}^{N_{\text{SVI}}} \log p(\mathbf{y} | \kappa^{(i)}) + \text{KL}(q(\kappa) \parallel p(\kappa))$ and propagate the gradient with respect to the parameters of $q(\kappa)$ using the reparametrisation trick (see section [Appendix B.1](#) and [Kingma and Welling \(2014\)](#))
- 6 **return** ELBO, ∇ELBO

4.1. Quantity of Interest

To assess the propagation of uncertainty, we consider a summary quantity for which a point estimate alone may not be informative enough for downstream tasks. In particular, we compute the log of total boundary flux through the boundary Γ_b :

$$r(\kappa) = \log \int_{\Gamma_b} e^{\kappa(s)} \nabla u(s) \cdot \mathbf{n} \, ds, \quad (35)$$

where \mathbf{n} is a unit vector normal to the boundary Γ_b . The quantity $\nabla u(s) \cdot \mathbf{n}$ represents the magnitude of the change in $u(s)$ at the boundary point s . For example, if u represented the temperature of an object, it would represent the loss of energy flowing out of the object.

4.2. Evaluation Criteria

To assess the inference of κ , we obtain S samples from the posterior distribution of κ , $\{\kappa^{(s)}\}_{s=1}^S$. For synthetic experiments, where we know the true κ which generated the observations, we compute the *expected error norm*. The computation is the average Euclidean norm of the error:

$$\text{Expected error} = \frac{1}{S} \sum_{s=1}^S \|\kappa^{(s)} - \kappa\|_2. \quad (36)$$

4.3. Poisson 1D

For this experiment, we assume the unit-line domain, which is discretized into 32 equal-length elements. We impose Dirichlet boundary conditions on both boundaries, specifically we set $u(0) = u(1) = 0$; we set the forcing term to be constant everywhere $f(x) = 1$. Unless specified otherwise, all experiments in this section use $n = 5$ observations per sensor and the sensor noise $\sigma_y = 0.01$. Sensors are located on each of the discretization nodes. We compare the results for three specifications of the prior length-scale, $\ell_\kappa = 0.1, 0.2, 0.3$. The length-scale used to generate the data is $\ell_\kappa = 0.2$. For inferences made using a shorter length-scale to generate the data, see [Appendix A](#).

4.3.1. Mean of variational posterior provides a good point estimate

Fig. 3 shows the expected error norm of the posterior estimates from Hamiltonian Monte Carlo (HMC), pre-conditioned Crank-Nicholson MCMC (pCN), as well as VB inference with different parametrisations of the covariance/precision matrix. It is evident that for prior length-scales $\ell_\kappa =$

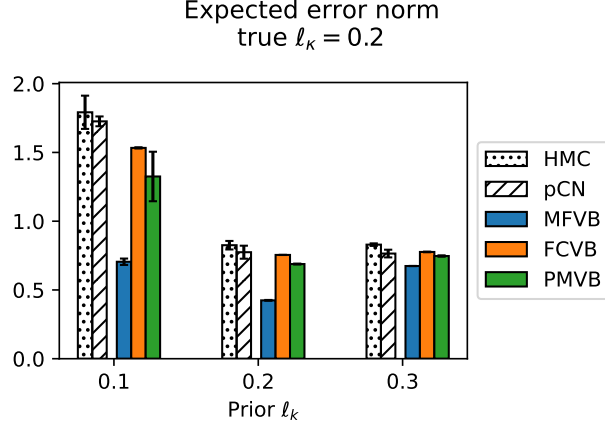


Figure 3: Expected error norm for the Poisson 1D problem, estimated using 10,000 samples from the inferred posterior distribution. Quantitatively, the sampling methods (HMC and pCN) and VB produce comparable results in terms of the expected error on the mean of the estimate. For a qualitative comparison, see Fig. 4.

0.2, 0.3, the posterior means computed by the variational Bayes methods are very close to the estimates from HMC and pCN. For prior $\ell_\kappa = 0.1$, MFVB error norm is lower than other VB methods and MCMC methods. This is most likely due to MFVB being a much easier optimisation task compared to other VB methods with more optimisation parameters that capture dependencies. While MCMC methods are asymptotically correct, in practice, devising efficient samplers for high-dimensional problems within a limited compute budget is still a challenging task and requires substantial hand-tuning. To affirm that all the VB methods provide a good estimate of the mean as compared to MCMC methods is better demonstrated by inspecting Fig. 4 which shows not only the mean but also the posterior uncertainty, which we discuss next.

4.3.2. VB adequately estimates posterior variance

Fig. 4 shows the true values of κ (red), the posterior means (black) and the corresponding estimate of the posterior variance (blue shaded regions) estimated by HMC, pCN, and variational inference with mean-field (MFVB), full covariance (FCVB), and precision matrix (PMVB) parametrisations for different values of true and prior length-scales. We observe that the posterior variance estimates computed by HMC, pCN, and full covariance VB are qualitatively very similar, with the estimated uncertainty increasing with increasing distance from the fixed boundary. However, the mean-field solution greatly underestimates posterior variance while computing a reasonable estimate of the posterior mean. For the PMVB parametrisation, the uncertainty is underestimated to a much lesser extent.

The observations above are further confirmed by the density plot of our quantity of interest: the log of the total flux on the boundary, shown in Fig. 5. For this example, we compute the flux on the left boundary at $x = 0$ and show the posterior distribution of this quantity. For longer prior length-scales, FCVB and PMVB agree with the estimates obtained from pCN and HMC, whereas mean-field VB underestimates the uncertainty. For the short prior length-scale ($\ell_\kappa = 0.1$), both PMVB and MFVB underestimate the uncertainty as compared with HMC, pCN, and FCVB schemes. The posterior distribution of FCVB approximately agrees with the MCMC schemes.

For the results obtained using the PMVB scheme, we used the 10-neighborhood structure to define the adjacency matrix and the non-zero elements of the precision matrix, \mathbf{Q} (see Sec. 3.3). The order of the neighbourhood structure, which corresponds to the precision matrix bandwidth, determines how much dependence within κ is captured by the approximating posterior distribution.

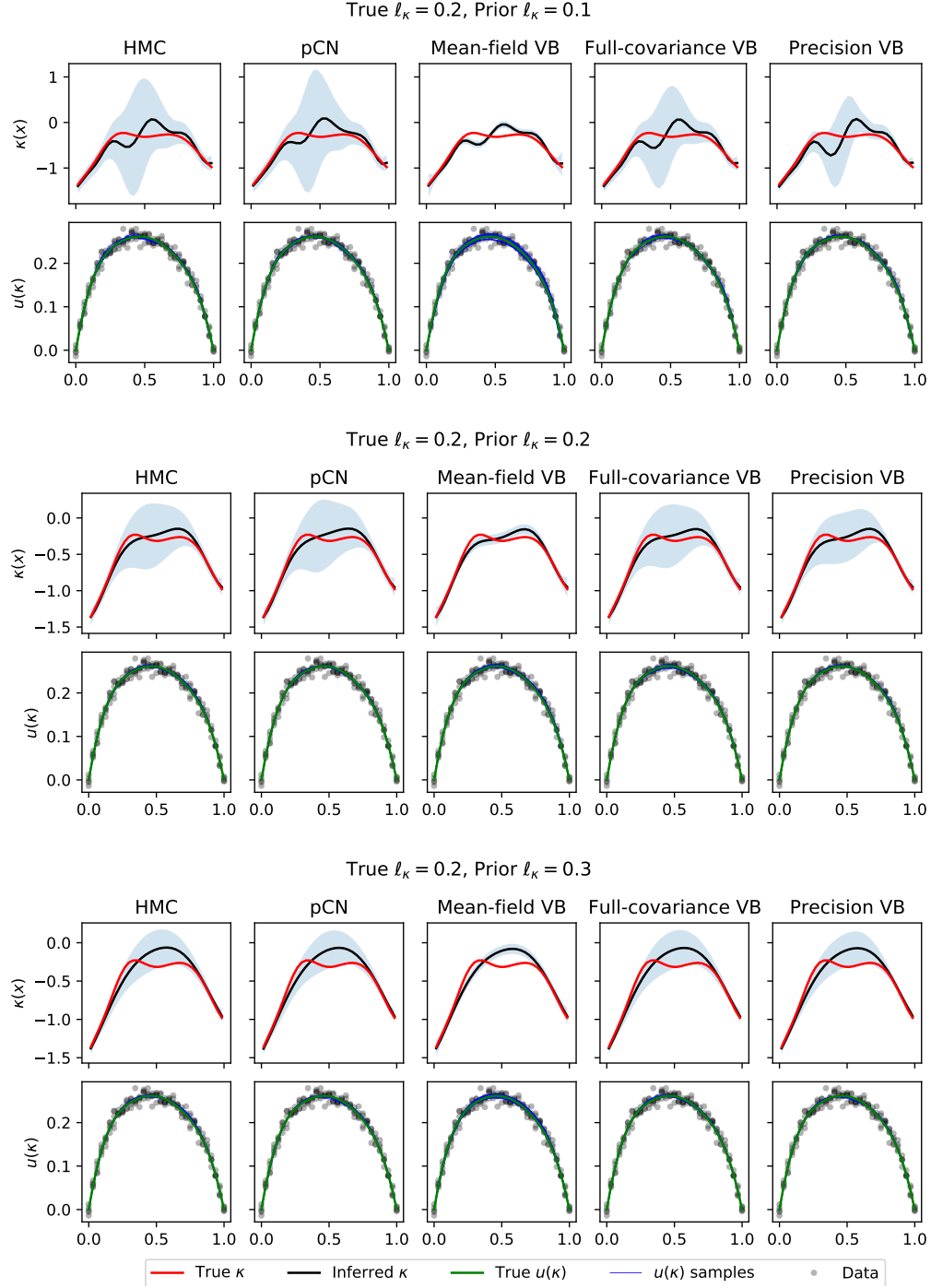


Figure 4: Top row in each of the three panels show true values of κ (red), posterior means (black) and posterior variances (blue shaded regions) for HMC, pCN, and VB variants for different values of prior length-scales ℓ_κ . The bottom rows show the data (black), true solution \mathbf{u} (green), solutions for different samples of κ (blue). For the PMVB estimate, the bandwidth is set to 10.

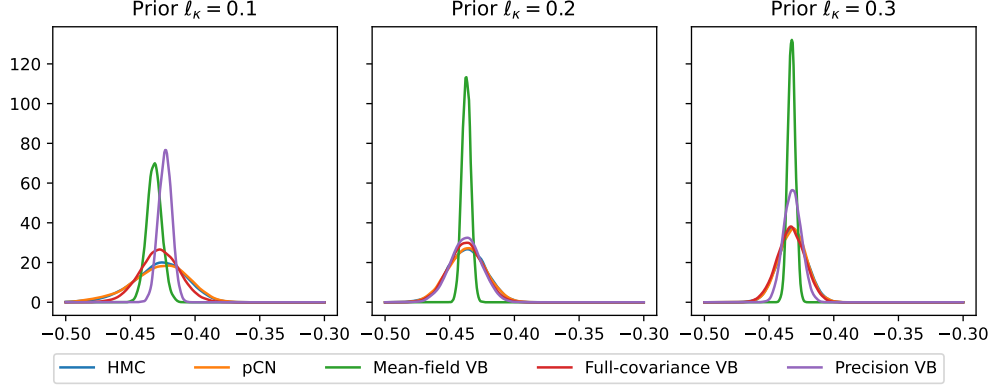


Figure 5: Log of the boundary flux at the left boundary node ($x = 0$) for the 1D Poisson example. For PMVB, the precision matrix bandwidth of 10 is used.

In Fig. 6, we show how the estimate of the mean and the variance of κ changes for different orders of neighbourhood structure. As expected, with the increasing bandwidth, the posterior estimate of κ gets closer to the estimate of FCVB, HMC, and pCN (shown in Fig. 4). While there is a significant change in the uncertainty estimate when we increase the bandwidth from 2 to 10, it is less pronounced when we change it from 10 to 20. For this reason, we choose the value of 10 for the PMVB parametrisation in 1D.

4.3.3. VB estimates improve with more observations and decreasing observational noise

The consistency of the posterior refers to the contraction of the posterior distribution to the truth as the data quality increases, *i.e.* either the number of observations increases or observation noise tends to zero. A recent line of work (Abraham and Nickl, 2020; Monard et al., 2020; Giordano and Nickl, 2020) showed the posterior consistency for the estimates obtained using popular MCMC schemes such as pCN or unadjusted discretized Langevin algorithm for Bayesian inverse problems based on PDE forward mappings. While similar results are not available for VB methods in infinite-dimensional case, consistency and Bernstein-von Mises type results have been shown for the finite-dimensional case, including Bayesian inverse problems (Wang and Blei, 2019; Lu et al., 2017). Empirically, our experiments show that for the given family of trial distributions the VB posterior distribution contracts to the true κ .

Firstly, we show that increasing the number of observations, N_y , results in a more accurate estimate. Given that the observations, $\{\mathbf{y}_i\}_{i=1}^{N_y}$, are independent of each other, the likelihood term of the ELBO (see Eq. (27)) is the product of the individual likelihood terms:

$$p(\mathbf{y}_1, \dots, \mathbf{y}_{N_y} \mid \kappa) = \prod_i^{N_y} p(\mathbf{y}_i \mid \kappa). \quad (37)$$

Secondly, by decreasing the observational noise, σ_y , we expect the posterior distribution to get closer to the ground truth and with lower uncertainty. Fig. 7 shows the true values of κ (red), the posterior mean estimates (black), and the corresponding estimates of the posterior variance (blue shaded regions) obtained by different variants of variational Bayes for varying numbers of observations (top panel) and different values of observational noise (bottom panel). We can see that MFVB underestimates the posterior variances and these estimates do not depend on the number of observations (top panel in Fig. 7) or the amount of observational noise (bottom panel in Fig. 7). However, the FCVB and PMVB uncertainty estimates get narrower with increasing number of

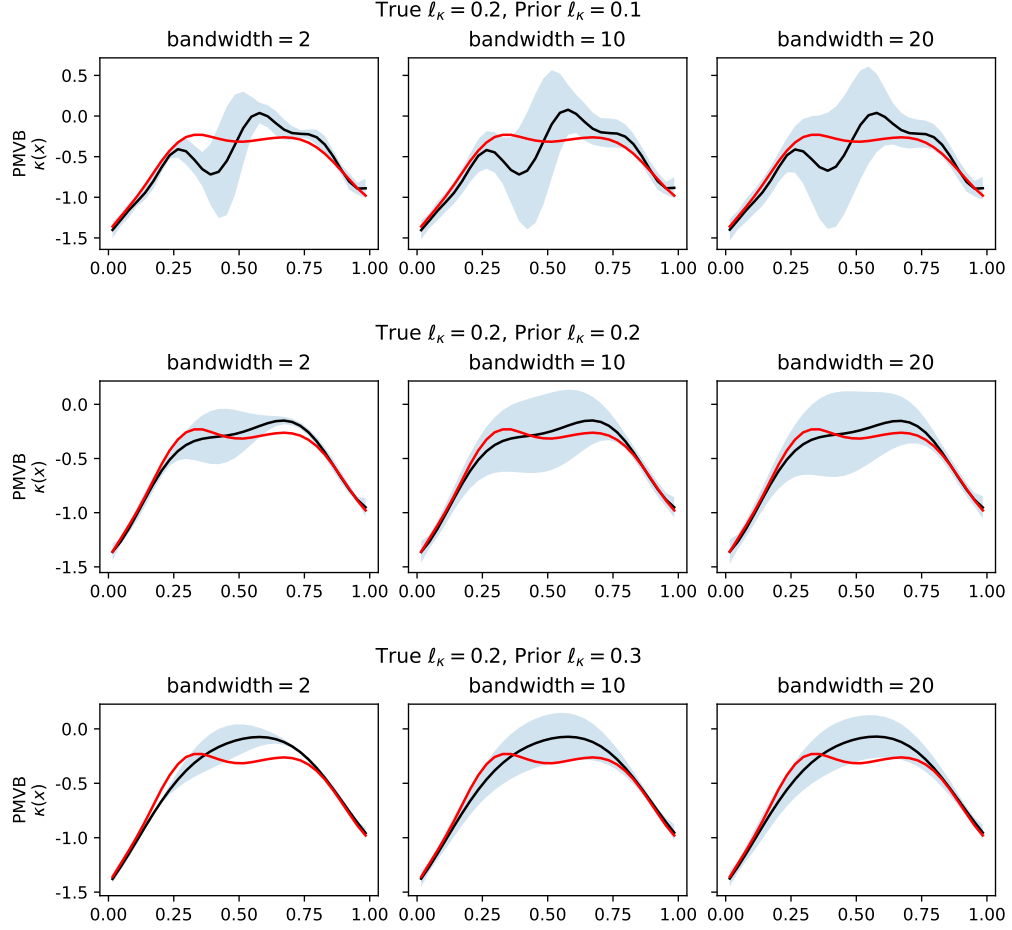


Figure 6: True values of κ (red), posterior means (black), and posterior variances (blue shaded region) for different matrix bandwidths of the precision matrix parametrisation of VB. Bandwidth corresponds to the order of neighbourhood structure considered when parametrising Q .

true ℓ_κ	prior ℓ_κ	Time (hours)				
		HMC		MFVB	FCVB	PMVB
0.1	0.1	15.2	(871–3244)	1.1	3.6	2.1
	0.2	11.1	(1043–4006)	0.7	2.7	2.1
	0.3	7.2	(1130–5408)	0.6	2.3	2.0
0.2	0.1	15.2	(1600–4700)	0.6	2.2	1.8
	0.2	10.4	(1067–3468)	0.6	2.3	2.0
	0.3	7.0	(1487–3969)	0.5	1.7	1.8

Table 1: Run-times for different inference schemes in hours for the Poisson 1D problem. For VB methods, $N_{\text{SVI}} = 3$. The column for HMC includes the range of effective sample sizes (ESS) across different components of κ .

observations and with decreasing observational noise, which is a desirable behaviour that should be exhibited by any consistent uncertainty estimation method. We can also see that the true solution is contained within the uncertainty bounds for all numbers of observations and noise levels for the full covariance parametrisation. This is not the case for the mean-field VB, providing another indication of uncertainty underestimation for this parametrisation.

4.3.4. VB is an order of magnitude faster than HMC

For HMC estimates, we obtain 200,000 samples out of which the first 100,000 are used to calibrate the sampling scheme and are subsequently discarded. The VB algorithm follows as described in Sec. 3. Table 1 provides the run-times for HMC, MFVB, FCVB, and PMVB. For the HMC column, we also report (shown in brackets) the range of effective sample sizes (ESS) across different components of κ . For details on ESS, we refer the reader to (Gelman et al., 2013, Ch. 11). Even with a conservative convergence criteria (described in Sec. 3.4), the computational cost of VB algorithms is up to 25 times lower than that of HMC. To emphasize the computational efficiency of the variational inference, we show the posterior estimates for different number of Monte Carlo samples in the estimation of ELBO. Fig. 8 shows that on a qualitative level, a low number of samples is sufficient to obtain a good estimate. In particular, even with 2 Monte Carlo samples, the estimates are very similar to the case where $N_{\text{SVI}} = 20$. However, a lower number of samples may result in slower convergence of the optimization scheme. Fig. 9 shows that for the FCVB and PMVB parametrisations, where the number of optimized parameters is larger than for MFVB, increasing the number of SVI samples may speed up the convergence of the optimization. The effect is not as strong for the MFVB parametrisation.

4.4. Poisson 2D

We consider a 2D Poisson problem on the unit-square domain with a circular hole as shown in Fig. 10, with boundary conditions as indicated in the same figure. The weak form of the problem is discretized with the finite element mesh consisting of 208 standard linear triangular elements and 125 nodes. The forcing term is assumed to be constant throughout the domain, $f(\mathbf{x}) = 1$. Unless specified otherwise, all experiments in this section use $n = 5$ observations per sensor and the sensor noise $\sigma_y = 0.001$ (note that for the 1D example we used $\sigma_y = 0.01$). The sensors are located at each node of the mesh.

Firstly, the results for the expected error norm in Fig. 11 show that the error of VB methods is very similar to the sampling methods (pCN and HMC). The results also show that the error is lowest when the prior ℓ_κ matches the length-scale used to generate the data.

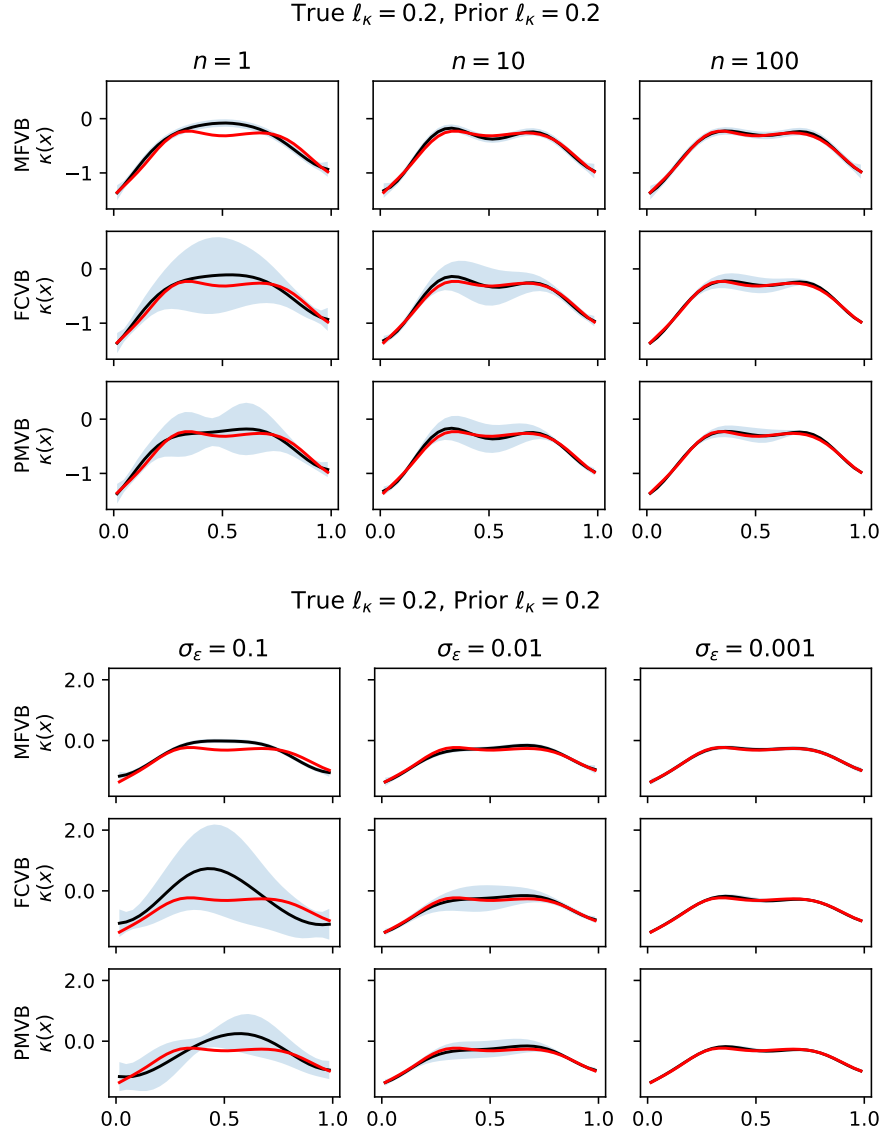


Figure 7: True values of $\kappa(x)$ (red), posterior means (black) and posterior variances (blue shaded regions) for VB with different parametrisations for different number of observations per sensor, $n = 1, 10, 100$ (top panel), and for different values of sensor noise $\sigma_\epsilon = 0.1, 0.01, 0.001$ (bottom panel).

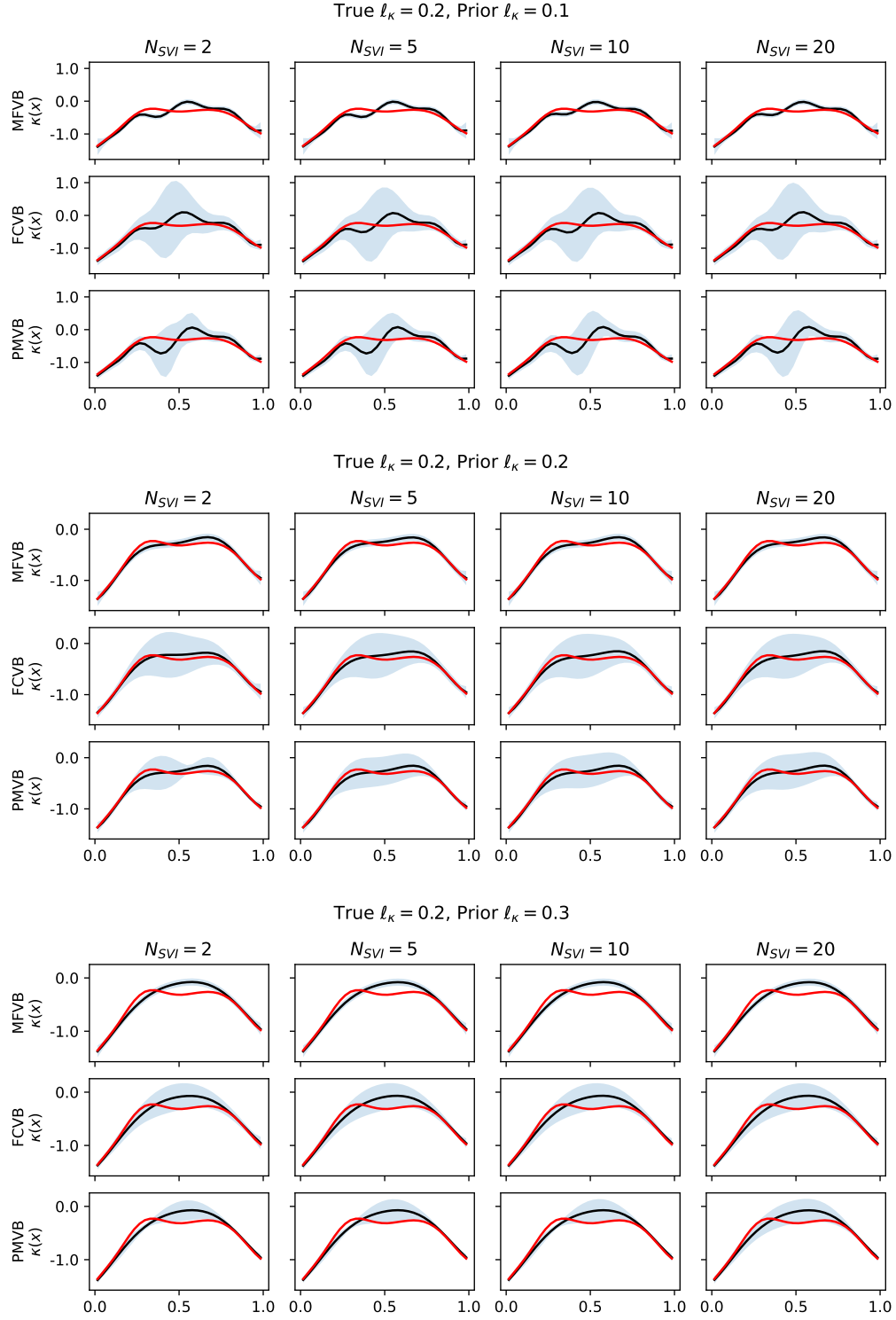


Figure 8: True values of $\kappa(x)$ (red), posterior means (black) and posterior variances (blue shaded regions) of VB with different parametrisations for varying number of Monte Carlo samples when computing ELBO. Three different length-scales for the prior are shown: 0.1, 0.2, 0.3.

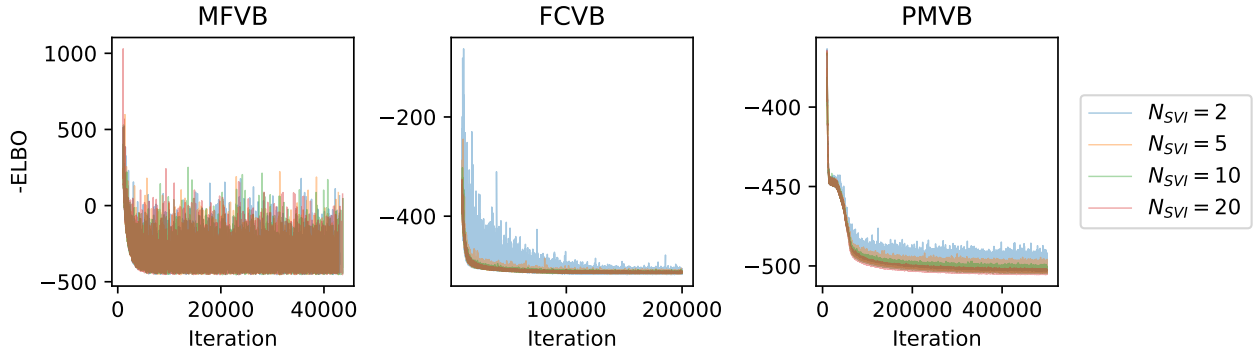


Figure 9: Negative ELBO trace plot for both MFVB and FCVB for different values of N_{SVI} . For this example, true $\ell_\kappa = 0.2$ and prior $\ell_\kappa = 0.1$.

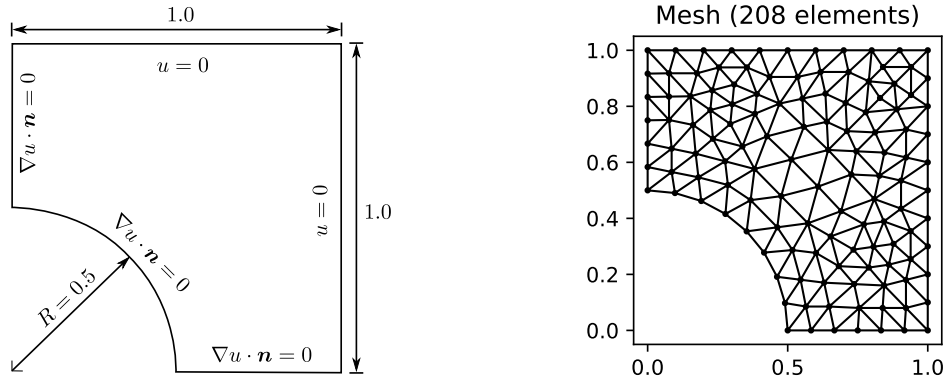


Figure 10: Left: Specification of the domain for the 2D Poisson problem. Note that we impose Dirichlet boundary conditions $u(x, y) = 0$ when $x = 1$ or $y = 1$. We impose Neumann boundary conditions on the rest of the boundary. Right: a triangular discretization of the domain.

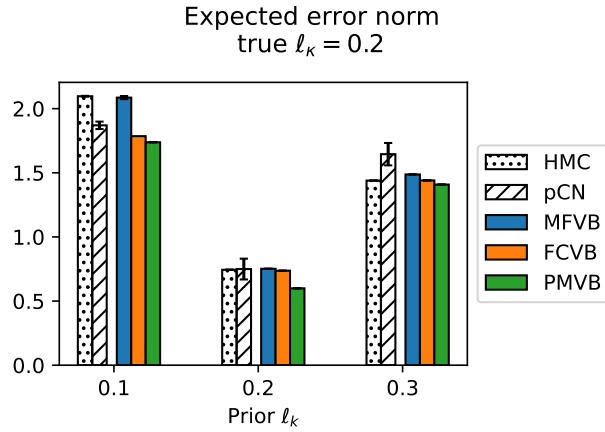


Figure 11: Expected error norm for the Poisson equation 2D example.

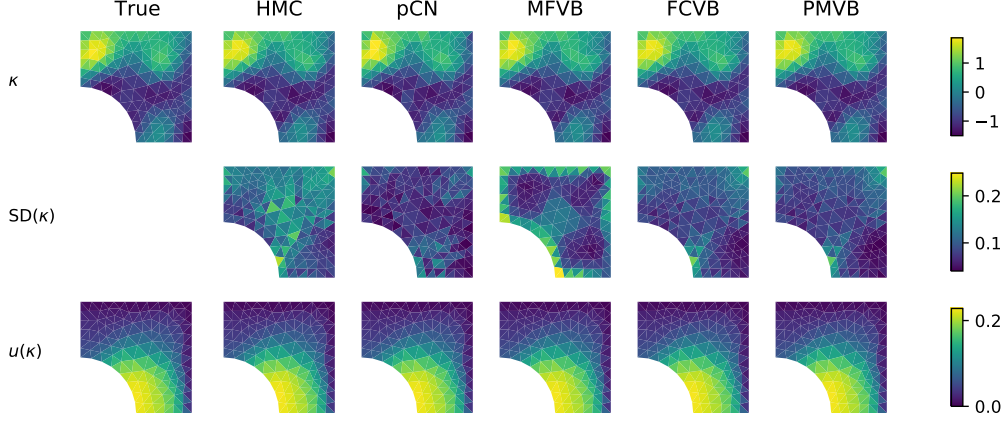


Figure 12: Posterior mean and variance for κ and corresponding \mathbf{u} for 2D Poisson example with prior length-scale $\ell_\kappa = 0.1$.

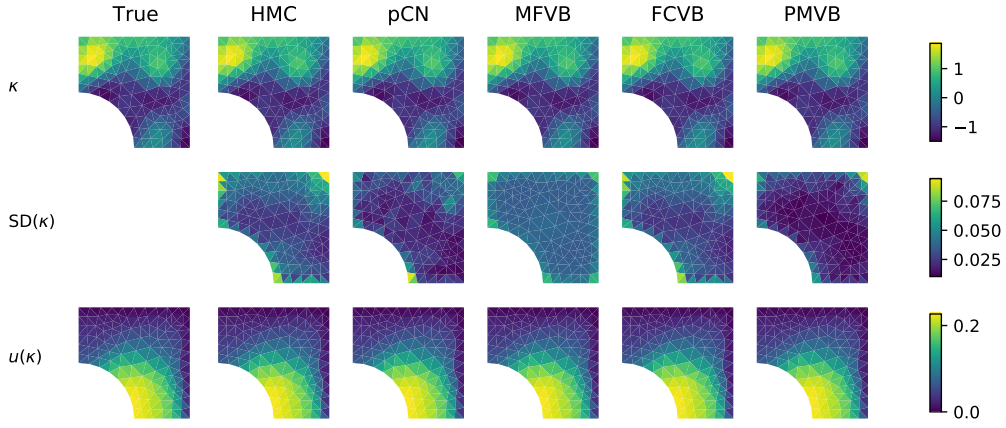


Figure 13: Posterior mean and variance for κ and corresponding \mathbf{u} for 2D Poisson example with prior length-scale $\ell_\kappa = 0.2$.

Fig. 12-14 show the results for the posterior mean and the variance of κ , and the solution \mathbf{u} corresponding to the mean of the posterior. We consider three configurations with prior length-scale $\ell_\kappa = 0.1, 0.2, 0.3$ where the length-scale used to generate the data is 0.2. In all cases, the estimates of the posterior mean of κ and the corresponding solutions \mathbf{u} are very close to the true values. Similarly to the 1D case discussed in Sec. 4.3, the variance estimates between HMC and FCVB are consistent, especially for longer prior length-scales. There seems to be a discrepancy between the estimates obtained using MFVB and those obtained by other methods. The estimates obtained using precision-matrix parametrisation are qualitatively very close to the FCVB and MCMC estimates.

For the quantity of interest, we compute the log of the total flux along the right boundary of the domain ($x = 1$), and the results are shown in Fig. 15. Unlike the 1D case, the posterior estimates of the boundary flux are approximately the same for all the considered methods, except for the mean-field estimate when prior $\ell_\kappa = 0.1$, where the MFVB estimate is biased as compared to the other methods.

The empirical computational cost for these experiments is given in Table 2. For the HMC experiments, we obtained 250,000 samples, out of which the first 125,000 were used to calibrate the sampling scheme and discarded afterwards. The timing results show that HMC takes an order of magnitude longer than variational Bayes, with some variation that depends on the parametrisation.

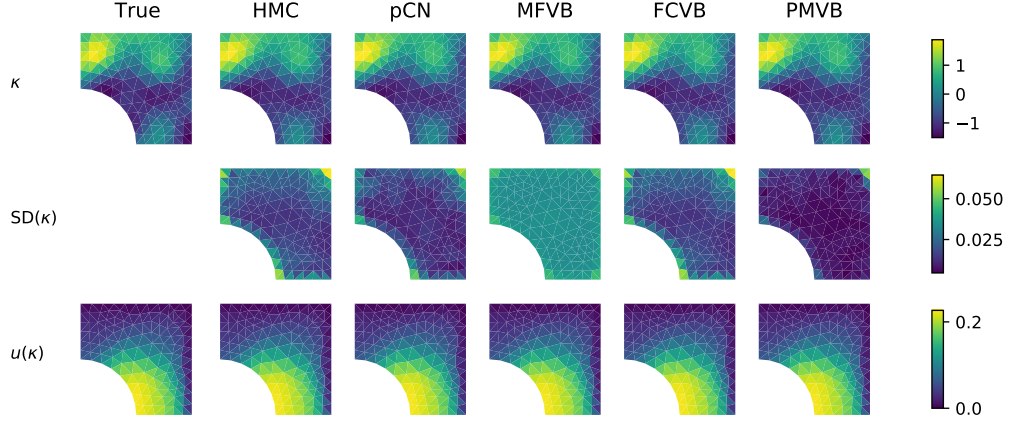


Figure 14: Posterior mean and variance for κ and corresponding \mathbf{u} for 2D Poisson example with prior length-scale $\ell_\kappa = 0.3$.

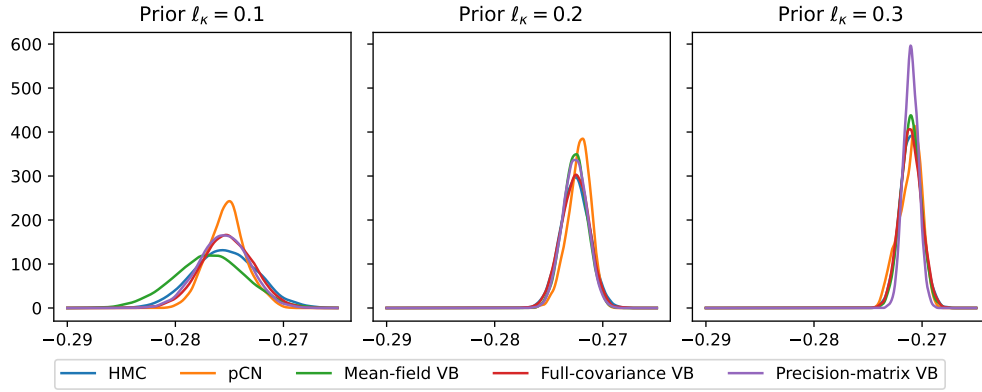


Figure 15: Log of the total flux computed along the right boundary ($x = 1$). For PMVB, the precision matrix is parametrised using the second-order neighbourhood structure, as shown in Fig. 2.

true ℓ_κ	prior ℓ_κ	Time (hours)				
		HMC	MFVB	FCVB	PMVB	
0.1	0.1	240.6	(930–11200)	6.4	29.6	28.1
	0.2	295.5	(1537–11067)	6.6	32.6	28.9
	0.3	242.0	(1057–6068)	7.3	27.3	30.6
0.2	0.1	242.7	(1102–18235)	6.2	34.3	27.2
	0.2	264.3	(1304–9848)	7.4	33.7	34.0
	0.3	221.9	(1192–6356)	7.8	31.3	34.0

Table 2: Run-times for different inference schemes in seconds. The number of Monte Carlo samples is $N_{\text{SVI}} = 5$ for all MFVB, FCVB, and PMVB. The column for HMC includes the range of effective sample sizes (ESS) across different components of κ .

4.5. Inverse Problem Benchmark

We evaluate the effectiveness of VB methods on a recently proposed benchmark for Bayesian inverse problems (Aristoff and Bangerth, 2021). The benchmark aims to provide a test case that reflects practical applications, but at the same time is easy to replicate. Like above, the test case is a Poisson equation inverse problem where the task is to recover log-diffusion, κ , from a finite set of noisy observations. The problem domain is a unit square, the forcing function $f(\mathbf{x}) = 10$ is constant throughout the domain, and the solution of the PDE is imposed to be zero on all four boundaries.

The benchmark discretizes κ using 64 quadrilateral elements, such that κ is constant for each individual element as shown in Fig. 16. The forward solution of the PDE is obtained after discretizing u using a 32×32 bilinear Lagrange rectangle elements. The locations where the solution is observed are placed on a uniform grid of 169 points (13×13). The measurements are corrupted by the Gaussian noise with standard deviation $\sigma_y = 0.05$. The authors of the benchmark provide the measurements as well as the true log-diffusion coefficient, κ , which generated the observations. The true log-diffusion coefficient, shown in Fig. 16, is zero throughout the domain, except two regions, where the value is $\log(10)$ and $\log(0.1)$. It is these two jumps that make it a non-trivial test case.

Unlike in the previous examples, we place a prior on κ which does not induce any spatial correlation between any of the κ coefficients. The role of the prior is to express our belief about the ranges of the coefficients, rather than any dependencies. Although authors place $\mathcal{N}(\mu = 4, \sigma^2 = 4)$ for each component of κ independently, we choose $\mathcal{N}(\mu = 0, \sigma^2 = 1)$ as most of the coefficients of the true κ are at the baseline level equal to zero, and the fact that the κ corresponds to the diffusion parameter on the log-scale, a priori we do not expect such high variance.

We performed the inference using HMC, MFVB, FCVB, and PMVB. The means and standard deviations of inferred log-diffusion coefficients, together with the PDE solutions corresponding to the inferred means, are shown in Fig. 16. The results suggest that the mean estimates of all three methods do capture the jumps and the overall structure of κ . Specifically, the FCVB estimate of the mean of κ is closest to the true value. As for uncertainty quantification, the MFVB and PMVB estimates are closer to the HMC estimate (our assumed ground truth for the uncertainty) than the FCVB estimate. The FCVB estimate seems to overestimate the uncertainty at a few locations. This is potentially due to being stuck in a local optimum during the optimization procedure, which for FCVB involves high-dimensional exploration.

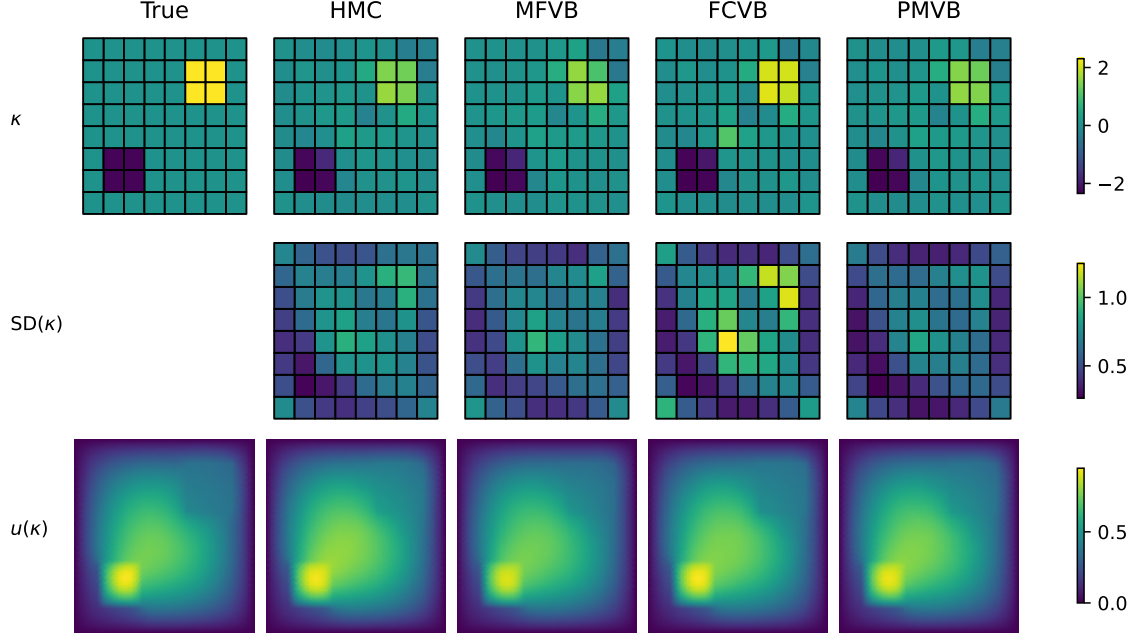


Figure 16: Independent prior for each coefficient of κ : $\kappa_i \sim \mathcal{N}(0, 1)$

5. Conclusions

In this paper, we have presented the variational inference framework for Bayesian inverse problems and investigated its efficacy on problems based on elliptic PDEs. Computationally, variational Bayes offers a tractable alternative to the intractable MCMC methods, and provides consistent mean and uncertainty estimates on the problems inspired by questions in computational mechanics. VB recasts the integration problem associated with Bayesian inference into an optimization problem. As such, it is naturally integrated with existing FEM solvers, using the gradient calculations from the FEM solvers to optimize the VB objective. Furthermore, the geometry of the problem encoded in the FEM mesh is utilized through the use of a sparse precision matrix that defines the conditional independence structure of the problem. Our results on the 1D and 2D Poisson problems support the claims of accuracy and scalability of VB. In particular, our results show that

- the mean of the variational posterior provides an accurate point estimate irrespective of the choice of the parametrisation of the covariance structure,
- the variational approximation with a full-covariance structure and the structured precision structure adequately estimates posterior variance when compared to HMC and pCN which are known to be asymptotically correct,
- parametrising the multivariate Gaussian distribution using a sparse precision matrix provides a way to balance the trade-off between computational complexity and the ability to capture dependencies in the posterior distribution,
- variational Bayes inference provides a good estimate for the mean and the variance of the posterior distribution in a time that is an order of magnitude faster than HMC or pCN,
- the multivariate Gaussian variational family is flexible enough to capture the true posterior distribution with high accuracy,

- the VB estimates may be used effectively in downstream tasks to estimate various quantities of interest,
- variational Bayes method is flexible enough to model multimodal posteriors, as illustrated on the truss example.

Our work may be extended in a number of natural ways that allows for greater adaptivity to the specific problems encountered in applications and integration within existing frameworks. Firstly, taking advantage of fast implementations of sparse linear algebra routines would further improve the scalability of VB with the structured precision matrix, as proposed in our work. Secondly, casting the inverse problem in a multi-level setting has potential to further improve computational efficiency (Nagel and Sudret, 2016). Thirdly, the results provided in this paper use standard off-the-shelf optimization routines; further computational improvements may be achieved using customized algorithms. As a further extension, in some applications it may be informative to consider the uncertainty in the forcing function so that the forward mapping is stochastic, as discussed in (Giolami et al., 2021). Finally, one of the aims of our work is to take advantage of the advances in Bayesian inference and adapt the novel algorithms to inverse problems in computational mechanics. As such, any further developments in VB as applied to machine learning and computational statistics problems may be directly applied using the framework proposed in this paper.

6. Implementation

Codes for performing all forms of variational Bayes inference presented in this paper are available on Github at <https://github.com/jp2011/bip-pde-vi>. The user must provide their own PDE solver which accepts κ as input parameter and computes $\log p(\mathbf{y} \mid \kappa)$, together with its gradient with respect to κ .

Acknowledgements

We would like to thank William Baldwin for providing us the necessary background for and the implementation of the bimodal example; Garoe Dorta for his help with Tensorflow debugging. JP, IK, and MG were supported by the EPSRC grant EP/P020720/2 for Inference, COmputation and Numerics for Insights into Cities (ICONIC, <https://iconicmath.org/>). JP was also supported by EPSRC (EP/L015129/1). MG was supported by EPSRC grants EP/T000414/1, EP/R018413/2, EP/R034710/1, EP/R004889/1, and a Royal Academy of Engineering Research Chair in Data Centric Engineering.

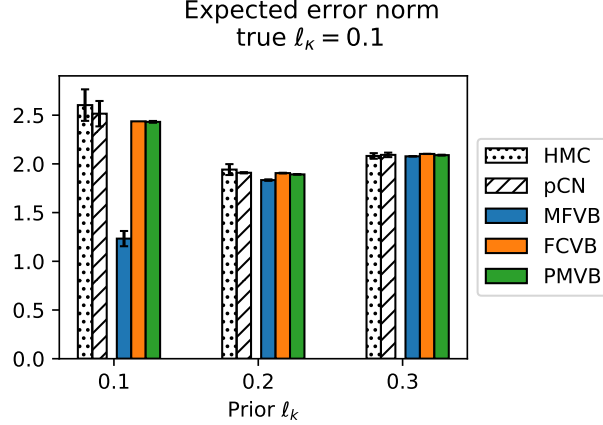


Figure A.17: Expected error norm for the Poisson 1D problem, estimated using 10,000 samples from the inferred posterior distribution.

Appendix A. Short length-scale results

Appendix B. Variational Inference

Appendix B.1. Reparametrisation trick

Reparametrisation trick allows computing the gradients of quantities derived from samples from a probability distribution with respect to the parameters ϕ of that probability distribution. This holds for probability distributions where samples can be obtained by a deterministic mapping, parametrised by ϕ , of other random variables.

Let ϵ be a set of random variables. We assume that samples of $\kappa \sim q(\kappa; \phi)$ are given by a deterministic mapping

$$\kappa = t(\phi, \epsilon). \quad (\text{B.1})$$

The KL divergence between approximating distribution $q(\kappa)$ and the prior $p(\kappa)$ is often available in closed form and so are its gradients with respect to ϕ . To estimate the gradients of the Monte Carlo estimate of the log-likelihood of the data,

$$\mathbb{E}_q [\log p(\mathbf{y} | \kappa)] \approx N_{\text{SVI}}^{-1} \sum_{i=1}^{N_{\text{SVI}}} \log p(\mathbf{y} | \kappa^{(i)}), \quad (\text{B.2})$$

we can use the chain rule of differentiation to obtain

$$\nabla_{\phi} N_{\text{SVI}}^{-1} \sum_{i=1}^{N_{\text{SVI}}} \log p(\mathbf{y} | \kappa^{(i)}) = N_{\text{SVI}}^{-1} \sum_{i=1}^{N_{\text{SVI}}} \nabla_{\kappa} \log p(\mathbf{y} | \kappa^{(i)}) \cdot \nabla_{\phi} t(\phi, \epsilon^{(i)}). \quad (\text{B.3})$$

Appendix C. Markov Chain Monte Carlo

Appendix C.1. Pre-Conditioned Crank-Nicholson Scheme

We consider the pre-conditioned Crank-Nicholson scheme proposed by [Cotter et al. \(2013\)](#). We summarise the procedure in Algorithm 2.

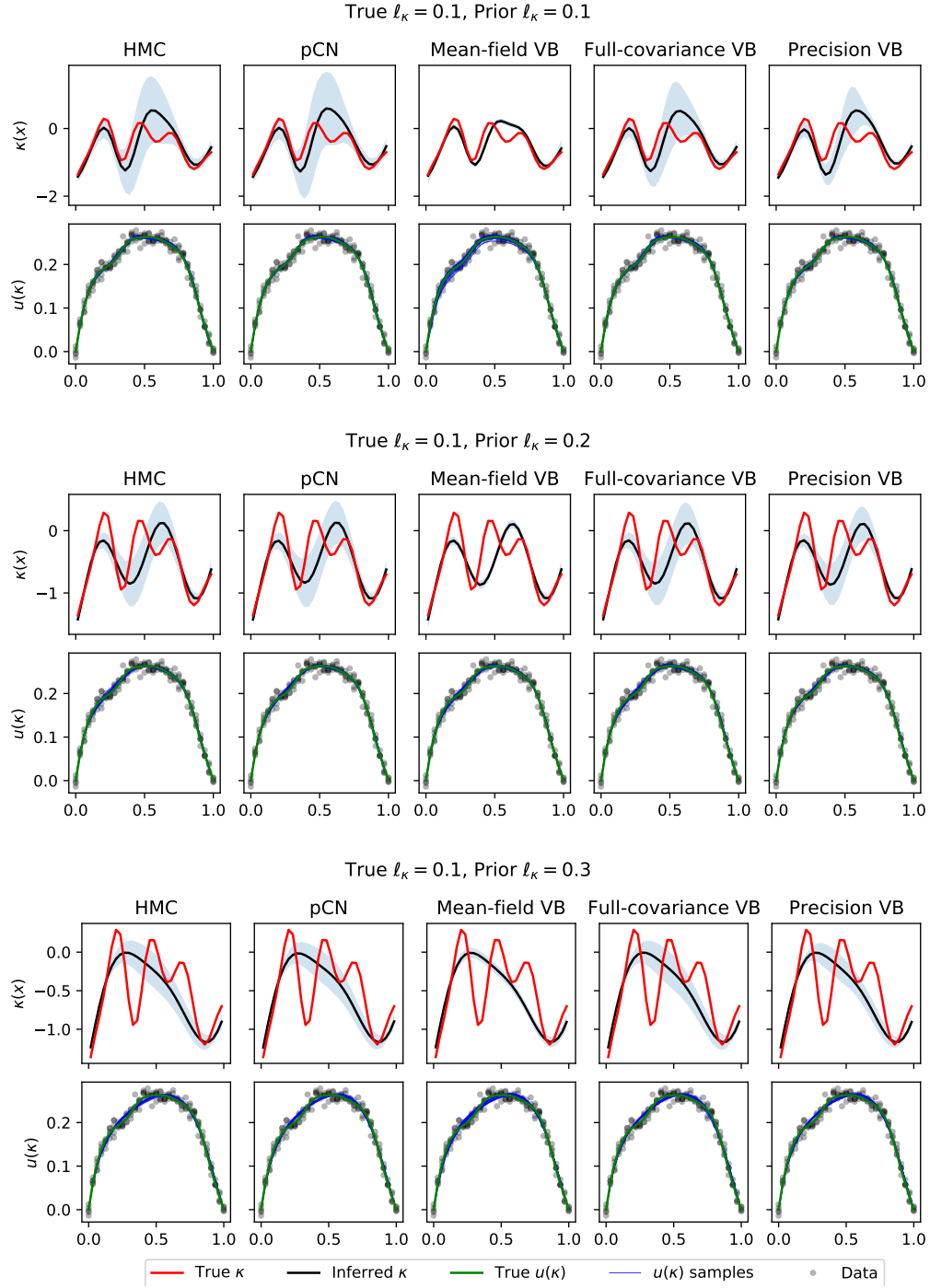


Figure A.18: Top row in each of the three panels show true values of $\kappa(x)$ (red), posterior means (black) and posterior variances (blue shaded regions) for HMC and VB variants for different values of prior length-scales ℓ_κ . The bottom rows show the data (black), true solution \mathbf{u} (green), solutions for different samples of κ (blue). For the PMVB estimate, the bandwidth is set to 10.

Algorithm 2: PRE-CONDITIONED CRANK-NICHOLSON MCMC (Cotter et al., 2013)

Input: $p(\mathbf{y} \mid \kappa)$: likelihood of the data, $\mu_0(\kappa)$: prior measures, β : corresponds to the amount of innovation in the proposal. If the value is small, there is little innovation and the proposed sample will be close to the previous sample.

Output: A list of samples from $\mu^y(\kappa)$.

- 1 **for** $t \leftarrow 1, 2, \dots$ **do**
- 2 Sample $\xi^{(t)} \sim \mu_0(\kappa)$
- 3 $v^{(t)} \leftarrow \sqrt{(1 - \beta^2)\kappa^{(t)}} + \beta\xi^{(t)}$
- 4 $\kappa^{(t+1)} \leftarrow \begin{cases} v^{(t)} & \text{with probability } \min\left(1, \exp(\Phi(\kappa^{(t)}; \mathbf{y}) - \Phi(v^{(t)}; \mathbf{y}))\right) \\ \kappa^{(t)} & \text{otherwise} \end{cases}$
- 5 **return** $[\kappa^{(1)}, \kappa^{(2)}, \dots]$

Appendix C.2. Hamiltonian Monte Carlo

Random walk Metropolis-Hastings algorithm may lead to inefficient explorations of the sample space, especially in the case of oddly-shaped densities and in higher dimensions. Improved proposal densities can alleviate this issue, but as the number of dimensions increases these measures become less effective. Hamiltonian Monte Carlo is a variant of Metropolis-Hastings which takes advantage of the gradients of the target distribution in the proposal, allowing for a more rapid exploration of the sample space, even in a high-dimensional target space. For each component κ_i of the target space, the scheme adds a ‘momentum’ variable ϕ_j (note that this is different from ϕ used in section Appendix B.1). Subsequently, κ and ϕ are updated jointly in a series of updates to propose a new sample (κ^*, ϕ^*) that is then accepted/rejected.

The proposal is largely driven by the momentum variable. The proposal step starts with drawing a new value of ϕ from $p(\phi)$ which needs to be specified. Then in a series of user-specified steps, L , the momentum variable ϕ is updated based on the gradient of the log of the target density, and κ is moved based on the momentum. Usually, the distribution of the momentum variable is $\mathcal{N}(0, M)$, where M is the so called ‘mass’ matrix. A diagonal matrix is often chosen to be able to efficiently sample from the momentum distribution. The full steps of the procedure are given in Algorithm 3.

The reason why HMC is suitable for high-dimensional problems is that to efficiently explore the target space, the algorithm exploits gradients of log of the target distribution.

The performance of the algorithm can be tuned in three ways: (i) choice of the momentum distribution $p(\phi)$, which in the version above requires specifying the mass matrix, (ii) adjusting the scaling factor of the leapfrog step, ϵ , and (iii) the number of leapfrog steps, L . Gelman et al. (2013) suggest setting ϵ and L so that $\epsilon L = 1$. They suggest tuning these so that the acceptance rate is about 65%. As for the mass matrix, the authors suggest that it should approximately scale with the inverse covariance matrix of the posterior distribution, $(\text{Cov}(\kappa \mid \mathbf{y}))^{-1}$. This can be achieved by a pre-run from which the empirical covariance matrix can be computed.

Appendix D. Bimodal example

One of the advantages of VB over Laplace approximation is the flexibility of the approximating distribution. To illustrate this, we consider an example with a one-dimensional truss which is fixed at one node and contains three degrees of freedom that correspond to the horizontal motion of the three nodes as shown in the left panel of Fig. D.19. We assume that the stiffness b_i of each member i is constant within each member and, furthermore, the stiffness of the member 1 is the average

Algorithm 3: HAMILTONIAN MONTE CARLO as presented in [Gelman et al. \(2013\)](#)

Input: $p(\boldsymbol{\kappa} \mid \mathbf{y})$: unnormalised target density, $p(\boldsymbol{\phi})$: momentum density and its mass matrix M , L : leapfrog steps, ϵ : scaling factor

Output: A list of samples from $p(\boldsymbol{\kappa} \mid \mathbf{y})$.

```

1 for  $t \leftarrow 1, 2, \dots$  do
2   Sample  $\boldsymbol{\phi}$  from  $p(\boldsymbol{\phi})$ 
3   for  $i \leftarrow 1$  to  $L$  do
4      $\boldsymbol{\phi} \leftarrow \boldsymbol{\phi} + \frac{1}{2}\epsilon \frac{d \log p(\boldsymbol{\kappa} \mid \mathbf{y})}{d \boldsymbol{\kappa}}$ 
5      $\boldsymbol{\kappa} \leftarrow \boldsymbol{\kappa} + \epsilon M^{-1} \boldsymbol{\phi}$ 
6      $\boldsymbol{\phi} \leftarrow \boldsymbol{\phi} + \frac{1}{2}\epsilon \frac{d \log p(\boldsymbol{\kappa} \mid \mathbf{y})}{d \boldsymbol{\kappa}}$ 
7      $r \leftarrow \frac{p(\boldsymbol{\kappa}^* \mid \mathbf{y}) p(\boldsymbol{\phi}^*)}{p(\boldsymbol{\kappa}^{t-1} \mid \mathbf{y}) p(\boldsymbol{\phi}^{t-1})}$ 
8      $\boldsymbol{\kappa}^t \leftarrow \begin{cases} \boldsymbol{\kappa}^* & \text{with probability } \min(r, 1) \\ \boldsymbol{\kappa}^{t-1} & \text{otherwise} \end{cases}$ 
9 return  $[\boldsymbol{\kappa}^1, \boldsymbol{\kappa}^2, \dots]$ 

```

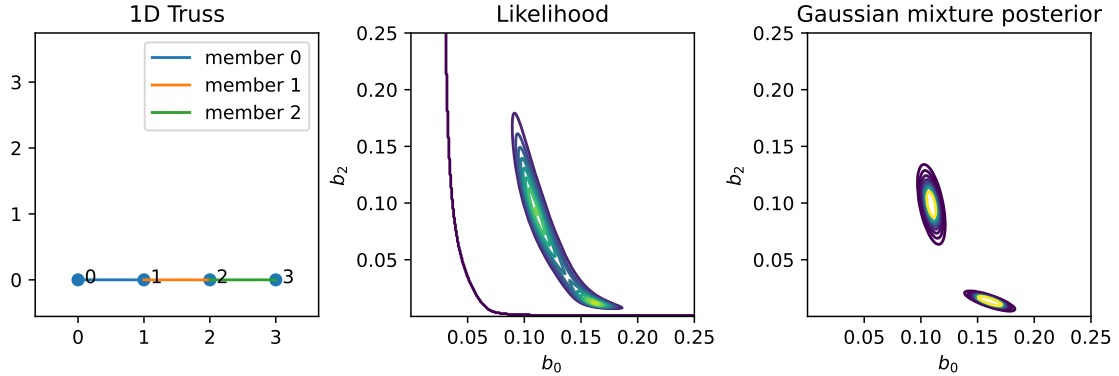


Figure D.19: One-dimensional truss discretized using three elements (shown in the left panel). The middle panel shows the likelihood surface for varying stiffness of the first and the third element. The right panel shows the posterior inference of the stiffness after displacements have been observed.

of the stiffness of the members at the ends, *i.e.* $b_1 = (b_0 + b_2)/2$. The inverse problem is then defined as follows. Given the displacement vector \mathbf{d} and boundary conditions, find the unknown stiffness parameters b_0 and b_1 . To prevent negative or small stiffness, constraints are imposed on the stiffness of each member, $b_i > 0.1$, and Neumann boundary conditions are set to $\mathbf{f} = (0, 1, 0.05)^T$. Due to the constraint on the stiffness of member 1, the image of the forward problem is a manifold with dimension 2 embedded in displacement space \mathbb{R}^3 . Due to the symmetry in this problem, the likelihood function, shown in centre panel of Fig. D.19, is bimodal.

We place a multivariate Gaussian as prior on the stiffness parameters, and use a bimodal trial distribution to infer the posterior distribution of the parameters b_0 and b_2 given observed displacement $\mathbf{d} = (0.1, 0.17, 0.23)^T$. Specifically, we consider a mixture of two multivariate Gaussians with equal fixed mixture weights. As there is no closed form expression for the KL divergence between a mixture of Gaussians and a single Gaussian, we estimate the KL divergence term in the ELBO using Monte Carlo sampling. As shown in the right panel of Fig. D.19, the resulting posterior distribution is bimodal and recovers the two modes present in the likelihood function. This illustrative example shows that when a proposed model exhibits multi-modality, the flexibility of variational Bayes methodology allows for specifying a family of trial distributions that can capture that property of the model.

References

- Abdulle, A. and Garegnani, G. (2021), ‘A probabilistic finite element method based on random meshes: A posteriori error estimators and Bayesian inverse problems’, *Computer Methods in Applied Mechanics and Engineering* **384**, 113961.
- Abraham, K. and Nickl, R. (2020), ‘On statistical Calderón problems’, *Mathematical Statistics and Learning* **2**(2), 165–216.
- Aristoff, D. and Bangerth, W. (2021), A benchmark for the bayesian inversion of coefficients in partial differential equations.
URL: <https://arxiv.org/abs/2102.07263>
- Arnst, M. and Soize, C. (2019), ‘Identification and sampling of Bayesian posteriors of high-dimensional symmetric positive-definite matrices for data-driven updating of computational models’, *Computer Methods in Applied Mechanics and Engineering* **352**, 300–323.
- Asaadi, E. and Heyns, P. S. (2017), ‘A computational framework for Bayesian inference in plasticity models characterisation’, *Computer Methods in Applied Mechanics and Engineering* **321**, 455–481.
- Babuška, I., Sawlan, Z., Scavino, M., Szabó, B. and Tempone, R. (2016), ‘Bayesian inference and model comparison for metallic fatigue data’, *Computer Methods in Applied Mechanics and Engineering* **304**, 171–196.
- Barajas-Solano, D. A. and Tartakovsky, A. M. (2019), ‘Approximate bayesian model inversion for PDEs with heterogeneous and state-dependent coefficients’, *Journal of Computational Physics* **395**, 247–262.
- Beck, J., Dia, B. M., Espath, L. F., Long, Q. and Tempone, R. (2018), ‘Fast Bayesian experimental design: Laplace-based importance sampling for the expected information gain’, *Computer Methods in Applied Mechanics and Engineering* **334**, 523–553.

- Beskos, A., Girolami, M., Lan, S., Farrell, P. E. and Stuart, A. M. (2017), ‘Geometric MCMC for infinite-dimensional inverse problems’, *Journal of Computational Physics* **335**, 327–351.
- Betz, W., Papaioannou, I., Beck, J. L. and Straub, D. (2018), ‘Bayesian inference with Subset Simulation: Strategies and improvements’, *Computer Methods in Applied Mechanics and Engineering* **331**, 72–93.
- Bishop, C., Lawrence, N., Jaakkola, T. and Jordan, M. (1998), Approximating posterior distributions in belief networks using mixtures, in M. Jordan, M. Kearns and S. Solla, eds, ‘Advances in Neural Information Processing Systems’, Vol. 10, MIT Press.
URL: <https://proceedings.neurips.cc/paper/1997/file/c0826819636026dd1f3674774f06c51d-Paper.pdf>
- Bishop, C. M. (2006), *Pattern Recognition and Machine Learning*, Information Science and Statistics, Springer, New York.
- Blei, D. M., Kucukelbir, A. and McAuliffe, J. D. (2017), ‘Variational inference: A review for statisticians’, *Journal of the American Statistical Association* **112**(518), 859–877.
- Bui-Thanh, T., Ghattas, O., Martin, J. and Stadler, G. (2013), ‘A computational framework for infinite-dimensional bayesian inverse problems part I: The linearized case, with application to global seismic inversion’, *SIAM Journal on Scientific Computing* **35**(6), A2494–A2523.
- Burt, D. R., Ober, S. W., Garriga-Alonso, A. and van der Wilk, M. (2021), Understanding variational inference in function-space, in ‘Third Symposium on Advances in Approximate Bayesian Inference’.
URL: <https://openreview.net/forum?id=7P9y3sRa5Mk>
- Carlton, A. G., Dia, B. M., Espath, L., Lopez, R. H. and Tempone, R. (2020), ‘Nesterov-aided stochastic gradient methods using Laplace approximation for Bayesian design optimization’, *Computer Methods in Applied Mechanics and Engineering* **363**, 112909.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M., Guo, J., Li, P. and Riddell, A. (2017), ‘Stan : A probabilistic programming language’, *Journal of Statistical Software* **76**(1), 1–32.
- Chen, P., Villa, U. and Ghattas, O. (2017), ‘Hessian-based adaptive sparse quadrature for infinite-dimensional Bayesian inverse problems’, *Computer Methods in Applied Mechanics and Engineering* **327**, 147–172.
- Cheng, C.-A. and Boots, B. (2017), Variational inference for gaussian process models with linear complexity, in I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, eds, ‘Advances in Neural Information Processing Systems’, Vol. 30, Curran Associates, Inc.
URL: <https://proceedings.neurips.cc/paper/2017/file/f8da71e562ff44a2bc7edf3578c593da-Paper.pdf>
- Cotter, S. L., Roberts, G. O., Stuart, A. M. and White, D. (2013), ‘MCMC methods for functions: Modifying old algorithms to make them faster’, *Statistical Science* **28**(3), 424–446.
- Csató, L. and Opper, M. (2002), ‘Sparse on-line gaussian processes’, *Neural Computation* **14**(3), 641–668.

- Cui, T., Marzouk, Y. and Willcox, K. (2016), ‘Scalable posterior approximations for large-scale Bayesian inverse problems via likelihood-informed parameter and state reduction’, *Journal of Computational Physics* **315**, 363–387.
- Cuthill, E. and McKee, J. (1969), Reducing the bandwidth of sparse symmetric matrices, in ‘Proceedings of the 1969 24th National Conference’, ACM ’69, Association for Computing Machinery, New York, NY, USA, pp. 157–172.
- Damianou, A. C., Titsias, M. K. and Lawrence, N. D. (2016), ‘Variational inference for latent variables and uncertain inputs in gaussian processes’, *Journal of Machine Learning Research* **17**(42), 1–62.
- Duane, S., Kennedy, A. D., Pendleton, B. J. and Roweth, D. (1987), ‘Hybrid monte carlo’, *Physics Letters B* **195**(2), 216–222.
- Durrande, N., Adam, V., Bordeaux, L., Eleftheriadis, S. and Hensman, J. (2019), Banded Matrix Operators for Gaussian Markov Models in the Automatic Differentiation Era, in K. Chaudhuri and M. Sugiyama, eds, ‘Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics’, Vol. 89 of *Proceedings of Machine Learning Research*, PMLR, pp. 2780–2789.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A. and Rubin, D. B. (2013), *Bayesian Data Analysis*, Chapman and Hall/CRC.
- Giordano, M. and Nickl, R. (2020), ‘Consistency of Bayesian inference with Gaussian process priors in an elliptic inverse problem’, *Inverse Problems* **36**(8).
- Giordano, R., Broderick, T. and Jordan, M. I. (2018), ‘Covariances, robustness, and variational bayes’, *Journal of Machine Learning Research* **19**, 1981–2029.
- Girolami, M., Febrianto, E., Yin, G. and Cirak, F. (2021), ‘The statistical finite element method (statFEM) for coherent synthesis of observation data and model predictions’, *Computer Methods in Applied Mechanics and Engineering* **375**.
- Hensman, J., Fusi, N. and Lawrence, N. D. (2013), Gaussian processes for big data, in ‘Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence’, UAI’13, AUAI Press, Arlington, Virginia, USA, pp. 282–290.
- Hensman, J., Rattray, M. and Lawrence, N. D. (2012), Fast variational inference in the conjugate exponential family, in ‘Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 2’, NIPS’12, Curran Associates Inc., Red Hook, NY, USA, pp. 2888–2896.
- Hoffman, M. D., Blei, D. M., Wang, C. and Paisley, J. (2013), ‘Stochastic variational inference’, *Journal of Machine Learning Research* **14**, 1303–1347.
- Huang, Y., Beck, J. L. and Li, H. (2017), ‘Bayesian system identification based on hierarchical sparse Bayesian learning and Gibbs sampling with application to structural damage assessment’, *Computer Methods in Applied Mechanics and Engineering* **318**, 382–411.
- Huang, Y., Beck, J. L., Li, H. and Ren, Y. (2021), ‘Sequential sparse Bayesian learning with applications to system identification for damage assessment and recursive reconstruction of image sequences’, *Computer Methods in Applied Mechanics and Engineering* **373**, 113545.

- Ibrahimbegovic, A., Matthies, H. G. and Karavelić, E. (2020), ‘Reduced model of macro-scale stochastic plasticity identification by Bayesian inference: Application to quasi-brittle failure of concrete’, *Computer Methods in Applied Mechanics and Engineering* **372**, 113428.
- Jankowiak, M., Pleiss, G. and Gardner, J. (2020), Parametric gaussian process regressors, in H. D. III and A. Singh, eds, ‘Proceedings of the 37th International Conference on Machine Learning’, Vol. 119 of *Proceedings of Machine Learning Research*, PMLR, pp. 4702–4712.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S. and Saul, L. K. (1999), ‘An introduction to variational methods for graphical models’, *Machine Learning* **37**(2), 183–233.
- Jordan, M. I. and Wainwright, M. J. (2007), ‘Graphical models, exponential families, and variational inference’, *Foundations and Trends® in Machine Learning* **1**(1–2), 1–305.
- Jordan, R., Kinderlehrer, D. and Otto, F. (1998), ‘The variational formulation of the fokker-planck equation’, *Siam Journal on Mathematical Analysis* **29**(1), 1–17.
- Kaipio, J. and Somersalo, E. (2005), *Statistical and Computational Inverse Problems*, number v. 160 in ‘Applied Mathematical Sciences’, Springer, New York.
- Karathanasopoulos, N., Angelikopoulos, P., Papadimitriou, C. and Koumoutsakos, P. (2017), ‘Bayesian identification of the tendon fascicle’s structural composition using finite element models for helical geometries’, *Computer Methods in Applied Mechanics and Engineering* **313**, 744–758.
- Kingma, D. P. and Ba, J. (2015), Adam: A method for stochastic optimization, in Y. Bengio and Y. LeCun, eds, ‘International Conference on Learning Representations’, San Diego, CA, USA.
- Kingma, D. P. and Welling, M. (2014), Auto-encoding variational bayes, in ‘2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings’.
- Lu, C. and Tang, X. (2015), ‘Surpassing Human-Level Face Verification Performance on LFW with GaussianFace’, *Proceedings of the AAAI Conference on Artificial Intelligence* **29**(1).
- Lu, Y., Stuart, A. and Weber, H. (2017), ‘Gaussian approximations for probability measures on \mathbb{R}^d ’, *SIAM/ASA Journal on Uncertainty Quantification* **5**(1), 1136–1165.
- MacKay, D. J. C. (2003), *Information Theory, Inference, and Learning Algorithms*, Cambridge University Press.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. and Teller, E. (1953), ‘Equation of state calculations by fast computing machines’, *The Journal of Chemical Physics* **21**(6), 1087–1092.
- Michélen Ströfer, C. A., Zhang, X.-L., Xiao, H. and Coutier-Delgosha, O. (2020), ‘Enforcing boundary conditions on physical fields in Bayesian inversion’, *Computer Methods in Applied Mechanics and Engineering* **367**, 113097.
- Minh, H. Q. (2017), ‘Infinite-dimensional Log-Determinant divergences between positive definite trace class operators’, *Linear Algebra and its Applications* **528**, 331–383.
- Monard, F., Nickl, R. and Paternain, G. P. (2020), ‘Statistical guarantees for Bayesian uncertainty quantification in non-linear inverse problems with Gaussian process priors’.

- Nagel, J. B. and Sudret, B. (2016), ‘A unified framework for multilevel uncertainty quantification in Bayesian inverse problems’, *Probabilistic Engineering Mechanics* **43**, 68–84.
- Ni, P., Li, J., Hao, H., Han, Q. and Du, X. (2021), ‘Probabilistic model updating via variational Bayesian inference and adaptive Gaussian process modeling’, *Computer Methods in Applied Mechanics and Engineering* **383**, 113915.
- Pandita, P., Tsilifis, P., Awalgaonkar, N. M., Bilonis, I. and Panchal, J. (2021), ‘Surrogate-based sequential Bayesian experimental design using non-stationary Gaussian Processes’, *Computer Methods in Applied Mechanics and Engineering* **385**, 114007.
- Pinski, F. J., Simpson, G., Stuart, A. M. and Weber, H. (2015), ‘Algorithms for Kullback–Leibler approximation of probability measures in infinite dimensions’, *SIAM Journal on Scientific Computing* **37**(6), A2733–A2757.
- Pleiss, G., Gardner, J., Weinberger, K. and Wilson, A. G. (2018), Constant-time predictive distributions for gaussian processes, in J. Dy and A. Krause, eds, ‘Proceedings of the 35th International Conference on Machine Learning’, Vol. 80 of *Proceedings of Machine Learning Research*, PMLR, pp. 4114–4123.
URL: <http://proceedings.mlr.press/v80/pleiss18a.html>
- Pyrialakos, S., Kalogeris, I., Sotiropoulos, G. and Papadopoulos, V. (2021), ‘A neural network-aided Bayesian identification framework for multiscale modeling of nanocomposites’, *Computer Methods in Applied Mechanics and Engineering* **384**, 113937.
- Quiñonero-Candela, J. and Rasmussen, C. E. (2005), ‘A unifying view of sparse approximate gaussian process regression’, *Journal of Machine Learning Research* **6**(65), 1939–1959.
- Ranganath, R., Tran, D. and Blei, D. (2016), Hierarchical variational models, in M. F. Balcan and K. Q. Weinberger, eds, ‘Proceedings of the 33rd International Conference on Machine Learning’, Vol. 48 of *Proceedings of Machine Learning Research*, PMLR, New York, New York, USA, pp. 324–333.
URL: <https://proceedings.mlr.press/v48/ranganath16.html>
- Reddi, S. J., Kale, S. and Kumar, S. (2018), On the convergence of Adam and beyond, in ‘International Conference on Learning Representations’.
- Rizzi, F., Khalil, M., Jones, R., Templeton, J., Ostien, J. and Boyce, B. (2019), ‘Bayesian modeling of inconsistent plastic response due to material variability’, *Computer Methods in Applied Mechanics and Engineering* **353**, 183–200.
- Robbins, H. and Monro, S. (1951), ‘A stochastic approximation method’, *The Annals of Mathematical Statistics* **22**(3), 400–407.
- Rue, H. and Held, L. (2005), *Gaussian Markov Random Fields: Theory and Applications*, number 104 in ‘Monographs on Statistics and Applied Probability’, Chapman & Hall/CRC, Boca Raton.
- Rue, H., Martino, S. and Chopin, N. (2009), ‘Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **71**(2), 319–392.

- Sabater, C., Le Maître, O., Congedo, P. M. and Görtz, S. (2021), ‘A Bayesian approach for quantile optimization problems with high-dimensional uncertainty sources’, *Computer Methods in Applied Mechanics and Engineering* **376**, 113632.
- Salimbeni, H., Cheng, C.-A., Boots, B. and Deisenroth, M. (2018), Orthogonally decoupled variational gaussian processes, in S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi and R. Garnett, eds, ‘Advances in Neural Information Processing Systems’, Vol. 31, Curran Associates, Inc.
URL: <https://proceedings.neurips.cc/paper/2018/file/cc638784cf213986ec75983a4aa08cdb-Paper.pdf>
- Salimbeni, H. and Deisenroth, M. (2017), Doubly stochastic variational inference for deep gaussian processes, in I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, eds, ‘Advances in Neural Information Processing Systems’, Vol. 30, Curran Associates, Inc.
- Seeger, M. W., Williams, C. K. I. and Lawrence, N. D. (2003), Fast forward selection to speed up sparse gaussian process regression, in C. M. Bishop and B. J. Frey, eds, ‘Proceedings of the Ninth International Workshop on Artificial Intelligence and Statistics’, Vol. R4 of *Proceedings of Machine Learning Research*, PMLR, pp. 254–261.
- Shi, J., Titsias, M. and Mnih, A. (2020), Sparse orthogonal variational inference for gaussian processes, in S. Chiappa and R. Calandra, eds, ‘Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics’, Vol. 108 of *Proceedings of Machine Learning Research*, PMLR, pp. 1932–1942.
- Snelson, E. and Ghahramani, Z. (2006), Sparse gaussian processes using pseudo-inputs, in Y. Weiss, B. Schölkopf and J. Platt, eds, ‘Advances in Neural Information Processing Systems’, Vol. 18, MIT Press.
- Solin, A., Cortes, S., Rahtu, E. and Kannala, J. (2018), PIVO: Probabilistic Inertial-Visual Odometry for Occlusion-Robust Navigation, in ‘2018 IEEE Winter Conference on Applications of Computer Vision (WACV)’, IEEE, pp. 616–625.
- Stuart, A. M. (2010), ‘Inverse problems: A bayesian perspective’, *Acta Numerica* **19**, 451–559.
- Sun, S., Zhang, G., Shi, J. and Grosse, R. B. (2019), Functional variational bayesian neural networks, in ‘7th International Conference on Learning Representations’, New Orleans, LA, USA.
URL: <https://openreview.net/forum?id=rkxacs0qY7>
- Tan, L. S. L. and Nott, D. J. (2018), ‘Gaussian variational approximation with sparse precision matrices’, *Statistics and Computing* **28**(2), 259–275.
- Tarakanov, A. and Elsheikh, A. H. (2020), ‘Optimal Bayesian experimental design for subsurface flow problems’, *Computer Methods in Applied Mechanics and Engineering* **370**, 113208.
- Tarantola, A. (2005), *Inverse Problem Theory and Methods for Model Parameter Estimation*, Society for Industrial and Applied Mathematics, Philadelphia, PA.
- Tikhonov, A. N. and Arsenin, V. Y. (1977), *Solutions of Ill-Posed Problems*, Scripta Series in Mathematics, Winston ; distributed solely by Halsted Press, Washington : New York.

- Titsias, M. (2008), Variational model selection for sparse gaussian process regression, Technical report, School of Computer Science, University of Manchester.
- Titsias, M. (2009), Variational learning of inducing variables in sparse gaussian processes, *in* D. van Dyk and M. Welling, eds, ‘Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics’, Vol. 5 of *Proceedings of Machine Learning Research*, PMLR, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, pp. 567–574.
- Tran, D., Blei, D. and Airoldi, E. M. (2015), Copula variational inference, *in* C. Cortes, N. Lawrence, D. Lee, M. Sugiyama and R. Garnett, eds, ‘Advances in Neural Information Processing Systems’, Vol. 28, Curran Associates, Inc.
- Tsilifis, P., Bilonis, I., Katsounaros, I. and Zabararas, N. (2016), ‘Computationally efficient variational approximations for bayesian inverse problems’, *Journal of Verification, Validation and Uncertainty Quantification* **1**(031004).
- Turner, R. E. and Sahani, M. (2011), Two problems with variational expectation maximisation for time series models, *in* A. T. Cemgil, D. Barber and S. Chiappa, eds, ‘Bayesian Time Series Models’, Cambridge University Press, Cambridge, pp. 104–124.
- Uribe, F., Papaioannou, I., Betz, W. and Straub, D. (2020), ‘Bayesian inference of random fields represented with the Karhunen–Loève expansion’, *Computer Methods in Applied Mechanics and Engineering* **358**, 112632.
- Villa, U., Petra, N. and Ghattas, O. (2021), ‘hIPPYlib: An extensible software framework for large-scale inverse problems governed by PDEs: Part I: Deterministic inversion and linearized bayesian inference’, *ACM Transactions on Mathematical Software* **47**(2), 1–34.
- Wang, B. and Titterton, D. M. (2005), Inadequacy of interval estimates corresponding to variational Bayesian approximations, *in* R. G. Cowell and Z. Ghahramani, eds, ‘Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics’, Vol. R5 of *Proceedings of Machine Learning Research*, PMLR, pp. 373–380.
- Wang, Y. and Blei, D. M. (2019), ‘Frequentist consistency of variational bayes’, *Journal of the American Statistical Association* **114**(527), 1147–1161.
- Williams, C. and Seeger, M. (2001), Using the nyström method to speed up kernel machines, *in* T. Leen, T. Dietterich and V. Tresp, eds, ‘Advances in Neural Information Processing Systems’, Vol. 13, MIT Press.
- Wu, L., Zulueta, K., Major, Z., Arriaga, A. and Noels, L. (2020), ‘Bayesian inference of non-linear multiscale model parameters accelerated by a Deep Neural Network’, *Computer Methods in Applied Mechanics and Engineering* **360**, 112693.
- Zhang, C., Butepage, J., Kjellstrom, H. and Mandt, S. (2019), ‘Advances in Variational Inference’, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **41**(8), 2008–2026.