

1 **PROBABILISTIC GRADIENTS FOR FAST CALIBRATION OF**
2 **DIFFERENTIAL EQUATION MODELS** *

3 JON COCKAYNE[†] AND ANDREW B. DUNCAN[‡]

4 **Abstract.** Calibration of large-scale differential equation models to observational or experimen-
5 tal data is a widespread challenge throughout applied sciences and engineering. A crucial bottleneck
6 in state-of-the art calibration methods is the calculation of local sensitivities, i.e. derivatives of the
7 loss function with respect to the estimated parameters, which often necessitates several numerical
8 solves of the underlying system of partial or ordinary differential equations. In this paper we present
9 a new probabilistic approach to computing local sensitivities. The proposed method has several advan-
10 tages over classical methods. Firstly, it operates within a constrained computational budget and
11 provides a probabilistic quantification of uncertainty incurred in the sensitivities from this constraint.
12 Secondly, information from previous sensitivity estimates can be recycled in subsequent computa-
13 tions, reducing the overall computational effort for iterative gradient-based calibration methods. The
14 methodology presented is applied to two challenging test problems and compared against classical
15 methods.

16 **Key words.** PDE constrained optimisation, sensitivity analysis, probabilistic numerics

17 **AMS subject classifications.** 62G86, 62P30, 35Q93, 35R30

18 **1. Introduction.** Complex systems arising in applied sciences and engineering
19 are often modelled by systems of coupled ordinary or partial differential equations
20 (ODEs or PDEs) derived from the underlying physical principles. Typically, the spe-
21 cific model behaviour will depend on a vector of parameters which must be *calibrated*
22 to observations of system. A major challenge in calibration is the high computational
23 cost associated with numerically solving the mathematical model for a given value
24 of the parameters. This is particularly relevant for large-scale models incorporat-
25 ing multi-physics and multiscale behaviour, as arise in the context of digital twins
26 [49]. This high cost often precludes the use of many iterative methods for calibration,
27 including both optimisation methods and Bayesian approaches that use sampling al-
28 gorithms such as Markov chain Monte-Carlo (MCMC). Each of these requires at least
29 one solve of the governing equations per iteration of the algorithm. In practice MCMC
30 often requires of $\mathcal{O}(10^5)$ model evaluations [20].

31 The calibration of differential equation models to observed data can be formulated
32 as a constrained optimisation problem [7, 22, 28], which is solved using deterministic
33 or stochastic optimisation methods. Most fundamental optimisation methods¹ either
34 require or are accelerated by access to derivatives of the functional to be minimised,
35 so that the solver for the underlying equations must be augmented with a routine that
36 provides the derivative of the solution with respect to model parameters. Employing
37 an approximation of the gradient, such as a finite-difference approximation, may seem
38 attractive due to ease of implementation, but obtaining accurate approximations can
39 be challenging and, when the parameter dimension is large, computationally expen-

*Submitted to the editors September 3, 2020.

Funding: JC was supported by Wave 1 of The UKRI Strategic Priorities Fund under the EPSRC Grant EP/T001569/1, particularly the “Digital Twins for Complex Engineering Systems” theme within that grant, and The Alan Turing Institute. ABD was supported by the Lloyds Register Foundation Programme on Data Centric Engineering and by The Alan Turing Institute under the EPSRC grant [EP/N510129/1].

[†]The Alan Turing Institute (jcockayne@turing.ac.uk)

[‡]Imperial College London and The Alan Turing Institute (a.duncan@imperial.ac.uk)

¹i.e. any method in the Newton family of optimisation methods.

40 sive. Thus it is usually preferable to obtain derivatives using *first-order sensitivity*
 41 *analysis*, which expresses the derivatives as the solution of an auxiliary system of dif-
 42 ferential equations known as the *sensitivity equations*. While the sensitivity equations
 43 are linear, they depend on the solution of the underlying equations and so must typi-
 44 cally be solved numerically. Thus, computing the sensitivities is at least as expensive
 45 as solving the system itself. Further, sensitivities must be computed for every pa-
 46 rameter value at which a gradient evaluation is required, making them prohibitively
 47 expensive for use in optimisation methods, where gradients are typically required over
 48 a large sequence of parameter values.

49 In the context of model calibration and uncertainty quantification, Gaussian pro-
 50 cesses (GPs) are often used as surrogate models for the solution of the underlying
 51 equations with the aim of making the calibration of such models tractable [27, 52].
 52 This approach is advantageous as derivatives of the GP posterior mean can usually be
 53 computed explicitly, permitting the use of gradient based optimisation methods and
 54 sampling. While GP surrogate models do provide an effective approach to calibrating
 55 black-box computer codes where little is known about the structure of the underly-
 56 ing model, this comes at the price of *data-efficiency*. Information about the gradient
 57 can only be obtained from multiple function evaluations near to the location of the
 58 required gradient, so again numerous evaluations may be required, particularly if the
 59 dimension of the parameter space is high.

60 These highlighted issues motivate the novel approach to computing sensitivities
 61 for optimisation problems presented in this paper. Our proposed approach is able to
 62 bridge the gap between the classical approach of numerically solving the sensitivity
 63 equations and the purely data-driven surrogate model approach. This is achieved by
 64 introducing a nonparametric Gaussian process model for the solution of the sensitivity
 65 equations that is defined over the entire parameter space. The output is a posterior
 66 distribution on the space of vector fields in parameter space, whose mean can be
 67 interpreted as an estimate of the local sensitivity across multiple parameter locations
 68 and whose variance controls the error in this estimate under regularity assumptions.

69 This approach offers various advantages to the state-of-the-art approaches: firstly
 70 the computational cost of the method can be carefully controlled by the user, either
 71 to attain a desired level of accuracy as measured by the “width” of the posterior
 72 distribution or to fit within a given fixed computational budget. Secondly, estimates
 73 of gradients at multiple parameter locations are able to share information between
 74 them to provide accurate gradient approximations without necessitating additional
 75 numerical solves of the underlying PDE model. Thirdly, the posterior distribution can
 76 be efficiently updated when a gradient evaluation at a new parameter value is required.
 77 These three advantages are particularly pertinent to model calibration methods which
 78 require multiple gradient evaluations along a trajectory.

79 Besides the immediate application to calibration of PDE models, the efficient
 80 approximation of sensitivities for large scale PDE models is of independent interest,
 81 with wide ranging applications including model order reduction [42], shape optimi-
 82 sation [33] and uncertainty quantification [2]. The probability distribution output
 83 from our new approach has a rigorous Bayesian interpretation, allowing it to be com-
 84 posed within inference and computation pipelines in a coherent manner to enable
 85 propagation of uncertainty.

86 **1.1. Related Work.** ODE- or PDE-constrained optimisation problems are a
 87 class of control problem in which the cost function involves the solution of a partial
 88 differential equation posed on a domain $D \subseteq \mathbb{R}^d$. Classically, such optimisation

89 problems arise in the context of design and control of engineering systems, for example
 90 in optimal topological design, shape design and optimal control of dynamic systems.
 91 See [26] for a review of optimisation algorithms for use in this context. Further,
 92 these problems arise naturally in the context of Bayesian inverse problems and model
 93 calibration. In particular, variational approaches to data assimilation for weather
 94 prediction can be naturally rephrased as PDE-constrained optimisation problems [18].

95 Sensitivity analysis seeks to quantify the dependence of a function $g(p)$ on pertur-
 96 bations of the problem data or parameters $p \in P$. Broadly speaking, we distinguish
 97 between *global* sensitivity analysis, which quantifies how input variability influences
 98 output variability of a model, and *local* sensitivity analysis which assesses the influence
 99 of infinitesimal input perturbations on model output. The former is typically assessed
 100 in terms of variance, classically using variants of Sobol’ indices [51]. By contrast, local
 101 sensitivity analysis involves the calculation of partial derivatives of function outputs
 102 with respect to parameters. Local sensitivity analysis plays a fundamental role in
 103 the context of ODE- or PDE-constrained optimisation [9]. In this setting, let $g(u, p)$
 104 denote the real-valued objective function for the optimisation problem, that depends
 105 on the solution $u(p)$ of a differential equation for a given parameter value $p \in P$. Then
 106 we seek to compute the total derivative $\frac{dg}{dp}(p)$, which constitutes the local sensitivities.

107 Generally speaking there are two approaches to computing such derivatives: the
 108 *forward* or *direct* method and the *adjoint* method. In the forward method, supposing
 109 that $p \subseteq \mathbb{R}^m$, the underlying equations are differentiated with respect to p_1, \dots, p_m to
 110 obtain a system of m equations for the sensitivities. The adjoint method originates in
 111 the theory of Lagrange multipliers in optimisation, and involves solving an auxiliary
 112 adjoint equation for the Lagrange multiplier λ from which the sensitivities can be
 113 directly computed. Given that the forward approaches involves solving a system of
 114 m equations while the adjoint approach involves solving only a single equation, the
 115 latter approach can be significantly more efficient for large m [48].

116 The computational cost of solving optimisation problems involving large-scale
 117 ODE or PDE models has motivated the use of *surrogate models*; approximations of
 118 the underlying model that have significantly lower computational overhead. Proposed
 119 approaches include using reduced order modelling based on reduced basis methods or
 120 proper orthogonal decompositions [4, 5]. These surrogate approaches are motivated
 121 by the fact that the adjoints, and therefore the gradients, of the low-dimensional sur-
 122rogate model can obtained efficiently. Recent efforts involve combining neural network
 123 models with low-dimensional physical models to obtain efficient and accurate surro-
 124gate models [16, 46, 47, 24, 50]. Again, these methods exploit the fact that gradients
 125 of neural network models can be obtained efficiently through back-propagation.

126 Gaussian processes (GPs) have been widely used to provide black-box emulation
 127 of computationally expensive codes [45], with [30] providing a mature Bayesian for-
 128mulation to the methodology. Emulation methods based on GPs are now widespread
 129 and find uses in numerous applications ranging from computer code calibration [27],
 130 uncertainty analysis [35] and MCMC [31, 12]. Among the first papers to consider
 131 application of emulation within sensitivity analysis was [36], which extended the work
 132 of [30] to computation of variance-based global sensitivities. Subsequent work by [29]
 133 considered a similar approach that exploited a tensor-product kernel to simplify the
 134 integration problems required, though this work did not consider the posterior co-
 135variance in their estimator. See [10] for a more extensive review of emulation-based
 136 global sensitivity analysis techniques, and [21, 3, 44] for a survey of applications of
 137 such approaches. One could envisage an analogous emulation strategy for local sen-
 138sitivity analysis of computationally expensive models that involves first constructing

139 an emulator \hat{g} of the objective function g and then evaluating the derivative $\frac{d\hat{g}}{dp}(p^*)$
 140 which, assuming a conducive emulator, can be computed at a lower cost than the
 141 derivative of g itself. A notable disadvantage of this approach is that to approximate
 142 local sensitivities in this way would require *global* information about g , since unless
 143 a highly structured surrogate model is used little information can be obtained about
 144 $\frac{dg}{dp}(p)$ from the single evaluation $g(p)$.

145 The method proposed in this paper aims to bridge the gap between classical nu-
 146 merical approaches and emulation-based approaches to calculating sensitivities within
 147 optimisation problems. The proposed method can be interpreted as a Bayesian prob-
 148 abilistic numerical method [14] for the solution of the forward or adjoint sensitivity
 149 equations over $D \times P$. It is similar to the probabilistic meshless methods for solutions
 150 of PDEs presented in [13, Chapter 5], but extended across parameter space. This
 151 formalism presents several advantages. Firstly it permits a high level of adaptivity,
 152 in that the solution can be refined over both P and D to increase accuracy either
 153 globally over parameter space P , or locally for particular value of the parameters
 154 $p \in P$. Secondly, subject to regularity assumptions, estimates of the gradient at a
 155 parameter p may re-use information from nearby gradient evaluations, exploiting the
 156 smoothness of the sensitivity equations to reduce the computational effort required
 157 for accurate gradient estimates at p when nearby gradients have already been evalu-
 158 ated. Thirdly, gradient estimates can be updated efficiently, allowing the adaptivity
 159 and smoothness properties mentioned to be exploited within algorithms that depend
 160 upon local sensitivities, such as gradient-based optimisation algorithms.

161 **1.2. Contributions.** The main contributions of the paper are as follows:

- 162 • We develop a probabilistic framework for computing gradients for differential
 163 equation models.
- 164 • We study the theoretical properties of this method, in particular its robustness
 165 to discretisation error.
- 166 • We demonstrate how the inferred gradients can be leveraged in optimisation
 167 problems.
- 168 • The results are demonstrated on a number of model problems to analyse its
 169 performance in comparison to classical approaches.

170 **1.3. Structure of the Paper.** The paper proceeds as follows. In [section 2](#)
 171 the classical approach to computing local sensitivities is formulated with examples
 172 of application to the problem of computing sensitivities for a simple PDE. [section 3](#)
 173 presents the novel probabilistic approaches and provides theoretical results relating to
 174 their accuracy and stability. [section 4](#) discusses the use of the probabilistic methods
 175 introduced in optimisation problems, and the empirical performance of these meth-
 176 ods is assessed in [section 5](#). We conclude with some discussion in [section 6](#). The
 177 supplementary material contains the proofs required for the paper in [section S1](#).

178 **1.4. Notation.** Let $W^{k,p}(D)$ denote the Sobolev space in which each function
 179 has k weak derivatives with finite $L^p(D)$ norm. We will use the notation $H^k(D) =$
 180 $W^{k,2}(D)$. Further let $H_0^k(D)$ denote the subset of $H^k(D)$ for which all $f \in H_0^k(D)$
 181 have $f = 0$ on ∂D and $H^{-k}(D)$ to be the dual of $H_0^k(D)$. For two normed spaces $\mathcal{U},$
 182 \mathcal{V} we will use the notation $\mathcal{L}(\mathcal{U}, \mathcal{V})$ to denote the set of all bounded linear operators
 183 from \mathcal{U} to \mathcal{V} . For the set of all bounded linear functionals on \mathcal{U} we will use the
 184 notation $\mathcal{U}^* = \mathcal{L}(\mathcal{U}, \mathbb{R})$. When \mathcal{U} is a set of functions on some domain D we will
 185 use the notation $\delta[x]$ to denote the evaluation functional for the point $x \in D$, i.e.
 186 $\delta[x](u) = u(x)$.

187 When both \mathcal{U} and \mathcal{V} are Hilbert spaces, for an operator $A \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ let $A^\dagger \in$
 188 $\mathcal{L}(\mathcal{V}, \mathcal{U})$ denote the adjoint of A . For $A \in \mathcal{L}(\mathcal{U}, \mathcal{U})$, recall that the trace of A is defined
 189 as $\text{trace}(A) = \sum_{i=1}^{\infty} \langle Ae_i, e_i \rangle$ where $(e_i)_{i=1}^{\infty}$ is an arbitrary orthonormal basis of \mathcal{U} .

190 Several operator norms will be required. For an operator $A : \mathcal{U} \rightarrow \mathcal{V}$ we will
 191 denote the operator norm by $\|A\|_{\mathcal{U} \rightarrow \mathcal{V}} = \sup_{u \in \mathcal{U}} \|Au\|_{\mathcal{V}} / \|u\|_{\mathcal{U}}$. When $\mathcal{U} = \mathcal{V}$ we will
 192 simply use the notation $\|A\|_{\mathcal{U}}$. The trace norm is given by $\|A\|_{\text{tr}} = \text{trace}([A^\dagger A]^{1/2})$ while
 193 the Hilbert-Schmidt norm is given by $\|A\|_{\text{HS}} = \text{trace}(A^\dagger A)^{1/2}$. Recall that $\|A\|_{\mathcal{U} \rightarrow \mathcal{V}} \leq$
 194 $\|A\|_{\text{HS}} \leq \|A\|_{\text{tr}}$.

195 **1.4.1. Fréchet Derivatives.** Of central importance to the paper is the concept
 196 of a Fréchet derivative. Let \mathcal{U} and \mathcal{V} each be normed spaces and consider a function
 197 $f : \mathcal{U} \rightarrow \mathcal{V}$. When it exists, Fréchet derivative of f at $u \in \mathcal{U}$ is defined to be the
 198 operator $\frac{df}{du}[u] \in \mathcal{L}(\mathcal{U}, \mathcal{V})$ that satisfies

$$199 \quad (1.1) \quad \lim_{\|h\| \rightarrow 0} \frac{\|f(u+h) - f(u) - \frac{df}{du}[u]h\|}{\|h\|} = 0$$

200 where the notation $\|h\| \rightarrow 0$ is a shorthand for the requirement that the limit exist
 201 uniformly across sequences (h_n) in \mathcal{U} such that $\|h_n\| \rightarrow 0$ as $n \rightarrow \infty$. It is important
 202 to observe that $\frac{df}{du}[u]$ is a linear operator that depends upon u , so that $\frac{\partial f}{\partial u}[u](v)$ is
 203 the Fréchet derivative at the location $u \in \mathcal{U}$ in the direction $v \in \mathcal{U}$. For a function
 204 $f : \mathcal{U} \times \mathcal{V} \rightarrow \mathcal{W}$ the partial Fréchet derivative is defined analogously to be the operator
 205 $\frac{\partial f}{\partial u} \in \mathcal{L}(\mathcal{U} \times \mathcal{V}, \mathcal{W})$ that satisfies

$$206 \quad \lim_{\|h\| \rightarrow 0} \frac{\|f(u+h, v) - f(u, v) - \frac{\partial f}{\partial u}[u, v](h)\|}{\|h\|} = 0$$

207 whenever the above limit exists. Finally, consider the case where the function u
 208 depends on v . Let $f : \mathcal{U} \times \mathcal{V} \rightarrow \mathcal{W}$, suppose that \mathcal{U} is a space of functions with
 209 domain \mathcal{V} . Then, when it exists, the Fréchet derivative of f with-respect-to v is the
 210 operator $\frac{df}{dv} \in \mathcal{L}(\mathcal{U} \times \mathcal{V}, \mathcal{W})$ that satisfies

$$211 \quad (1.2) \quad \lim_{\|h\| \rightarrow 0} \frac{\|f(u(v+h), v+h) - f(u, v) - \frac{df}{dv}[u(v), v](h)\|}{\|h\|} = 0.$$

212 This will sometimes be referred to as the *total* Fréchet derivative of f .

213 **2. Background.** In this section a formal presentation of local sensitivity analysis
 214 is provided. In [subsection 2.1](#) the problem is introduced, while [subsections 2.2](#) and [2.3](#)
 215 present forward and adjoint sensitivity analysis, respectively. Lastly in [subsection 2.4](#)
 216 we will briefly discuss probabilistic numerical methods for the solution of PDEs, and
 217 discuss their similarity to this work.

218 **2.1. Local Sensitivity Analysis.** We begin by introducing the relevant spaces
 219 for the problem. Let \mathcal{U} , P , \mathcal{F} and \mathcal{G} each be real-valued Banach spaces. In this paper
 220 it will be assumed that \mathcal{U} and \mathcal{F} are infinite-dimensional spaces of functions defined
 221 on spatial domain D , with \mathcal{U} referred to as the *solution space* and \mathcal{F} as the *constraint*
 222 *space*. Define \mathcal{U}_P to be a space of real-valued functions on $D \times P$ with the property
 223 that $u(\cdot, p) \in \mathcal{U}$ for all $p \in P$, and let $\mathcal{U}_{\partial P} = \left\{ \frac{\partial u}{\partial p} : u \in \mathcal{U}_P \right\}$. The set \mathcal{F}_P is defined
 224 analogously. The *parameter space* P may be finite- or infinite-dimensional. The space
 225 \mathcal{G} will be referred to as the *quantity of interest (QoI) space*, and will be assumed to

226 be finite-dimensional. In particular it will often be the case that $\dim(\mathcal{G}) = 1$, though
 227 we note that this is not required for the presentation below.

228 Two functions define the problem. The function $F : \mathcal{U} \times P \rightarrow \mathcal{F}$ is referred to as
 229 the *constraint function*, and loosely speaking this encapsulates all of the constraints
 230 that must be satisfied in order for a pair $(u, p) \in \mathcal{U} \times P$ to constitute a solution to
 231 the PDE. The function $g : \mathcal{U} \times P \rightarrow \mathcal{G}$ is referred to as the *QoI function*, and this
 232 describes a typically low-dimensional quantity of interest derived from the solution;
 233 in the context of optimisation problems this will generally be the objective function
 234 whose minimiser is sought.

235 More formally, F is such that for each $p \in P$ there is a unique $u^\dagger \in \mathcal{U}_P$ that
 236 satisfies $F(u^\dagger(\cdot, p), p) = 0$ for each $p \in P$. For convenience, define the *parameter-to-*
 237 *solution map* $U : P \rightarrow \mathcal{U}$ which provides the solution to the underlying differential
 238 equation for a particular value of the parameter, i.e. $U(p) = u^\dagger(\cdot, p)$. As a result, the
 239 equation $F(U(p), p) = 0$ is automatically satisfied for all $p \in P$.

240 It will be assumed the partial Fréchet derivatives of F and g with-respect-to both
 241 u and p exist and are tractably computable for all pairs $(u, p) \in \mathcal{U} \times P$. It will also be
 242 assumed that the derivative of U with-respect-to p exists but is not tractable. Note
 243 that this implies the existence of the total derivatives $\frac{dF}{dp}$ and $\frac{dg}{dp}$. Lastly we assume
 244 that $\frac{\partial F}{\partial u}[U(p), p]$ is nonsingular for each $p \in P$.

245 The objective is to estimate the value of the Fréchet derivative

$$246 \quad \frac{dg}{dp}[U(p), p] \in \mathcal{L}(P, \mathcal{G})$$

247 for a pair $(U(p), p)$. Note that since the location at which the derivative is taken is
 248 $U(p)$, this should be interpreted as a total Fréchet derivative in the form of (1.2). To
 249 fix ideas we consider the following simple parameter sensitivity problem.

250 **EXAMPLE 2.1** (Partial Differential Equation). *Let $P \subseteq \mathbb{R}^n$ be an open set.*
 251 *Consider the following parametrised steady state conductivity model:*

$$252 \quad \begin{aligned} -\nabla \cdot (\kappa(x; p) \nabla u(x)) &= f(x) & x \in D \\ u(x) &= 0 & x \in \partial D \end{aligned}$$

253 where $f \in H^{-1}(D)$ and $\kappa : D \times P \rightarrow \mathbb{R}^{d \times d}$ satisfies $\lambda_p |e|^2 \leq e \cdot \kappa(x, p) e \leq \Lambda_p |e|^2$
 254 for all $x \in D$ and $e \in \mathbb{R}^d$ for some $0 < \lambda_p < \Lambda_p < \infty$ for all $p \in P$. Standard
 255 existence theory for elliptic PDEs [17, Section 6.2, Theorem 3] states that a weak
 256 solution $u \in H_0^1(D)$ exists for every $p \in P$. For convenience we will suppose that the
 257 boundary conditions are implicitly satisfied, i.e. $\mathcal{U} = H_0^1(D)$. The constraint equation
 258 is given by $F(u, p) = -\nabla \cdot (\kappa(x; p) \nabla u(x)) - f(x)$ so that $\mathcal{F} = H^{-1}(D)$. Suppose that
 259 the quantity-of-interest is $g(x) = \|u\|_2 = \left(\int_D u^2(x) dx\right)^{\frac{1}{2}}$.
 260

261 Both forward and adjoint sensitivities are computed by first observing that the
 262 total derivative of interest, $\frac{dg}{dp}$ satisfies the following identity:
 263

$$264 \quad (2.1) \quad \frac{dg}{dp}[U(p), p] = \frac{\partial g}{\partial u}[U(p), p] \frac{dU}{dp}[p] + \frac{\partial g}{\partial p}[U(p), p].$$

265 Since it is assumed that $\frac{\partial g}{\partial u}$ and $\frac{\partial g}{\partial p}$ are each analytically tractable, the only remaining
 266 quantity that must be computed is $\frac{dU}{dp}$. The challenge is that since the parameter-to-
 267 solution map $U(p)$ is typically inaccessible and must be approximated independently
 268 for each $p \in P$, $\frac{dU}{dp}$ is also difficult to compute. The forward and adjoint approaches
 269 handle this intractability in different ways, which will now be presented.

270 **2.2. Forward Sensitivity Analysis.** In forward sensitivity analysis we seek to
 271 calculate $\frac{dU}{dp}$ directly. Note that we have

$$272 \quad \frac{\partial F}{\partial p}[U(p), p] = \frac{\partial F}{\partial u}[U(p), p] \frac{dU}{dp}(p) + \frac{\partial F}{\partial p}[U(p), p], \quad p \in P.$$

273 Further, since by construction $F(U(p), p) = 0$, we also have that $\frac{dF}{dp}[U(p), p] = 0$.
 274 This gives the forward sensitivity equation

$$275 \quad (2.2) \quad \frac{\partial F}{\partial u}[U(p), p] \frac{dU}{dp}[U(p), p] = -\frac{\partial F}{\partial p}[U(p), p], \quad p \in P$$

276 which is a linear system whose solution can be computed to determined $\frac{dU}{dp}$, since
 277 $\frac{\partial F}{\partial u}$ is assumed to be invertible. This solution can then be substituted into (2.1) to
 278 compute $\frac{dq}{dp}$.

279 Note that both the operator $\frac{\partial F}{\partial u}[U(p), p]$ and the right-hand-side $-\frac{\partial F}{\partial p}[U(p), p]$
 280 depend both on the parameter value p and the solution $U(p)$. This has two important
 281 consequences. Firstly, if sensitivities are required at another point $q \neq p$ then the
 282 solution $U(q)$ must be recomputed and the forward sensitivity equation (2.2) must
 283 be solved anew to determine $\frac{dU}{dp}[U(q), q]$. Secondly, for most problems of interest
 284 $U(p)$ will not be available explicitly and one must substitute an approximate solution
 285 $\hat{U}(p) \approx U(p)$. This may induce further numerical error, the impact of which must
 286 in turn be analysed, but also means that even though (2.2) is linear, its solution is
 287 unlikely to be available in closed-form owing to its dependence on $\hat{U}(p)$. We now
 288 consider the computation of the forward sensitivities for [Example 2.1](#).

289 **EXAMPLE 2.2 (Elliptic PDE: Forward Sensitivity Analysis).** *We begin by de-*
 290 *rivng $\frac{\partial F}{\partial p}$. Assume that κ is once-differentiable in each coordinate of p and that*
 291 *$\sup_{x \in D} |\partial_{p_i} \kappa(x; p)| < \infty$. The Frechét derivative of F with respect to p at $(U(p), p)$ is*
 292 *defined by*

$$293 \quad (2.3) \quad \frac{\partial F}{\partial p}[U(p), p]q = -\sum_{i=1}^m \nabla \cdot \left(\frac{\partial \kappa}{\partial p_i}(x; p) \nabla U(p)(x) \right) q_i, \quad q \in P.$$

294 *From energy estimates for weak solutions of elliptic PDEs, $\nabla U(p)(x) \in L^2(D)$. For*
 295 *illustration, it is straightforward to show that the RHS of (2.3) lies in $H^{-1}(D)$. The*
 296 *derivative $\frac{\partial F}{\partial u}$ is given by*

$$297 \quad (2.4) \quad \frac{\partial F}{\partial u}[U(p), p](v) = -\nabla \cdot (\kappa(x; p) \nabla v(x)) \quad v \in \mathcal{U}.$$

298 *so that clearly $\frac{\partial F}{\partial u}[U(p), p] \in \mathcal{F}$ since in this case, owing to the linearity of the PDE*
 299 *operator, $\frac{dF}{du}$ is identical to this operator and independent of both p and $U(p)$, though*
 300 *for general nonlinear problems this will not be the case. The sensitivities of U with*
 301 *respect to the p_i are therefore defined by the following system of PDEs*

$$302 \quad (2.5) \quad -\nabla \cdot \left(\kappa(x; p) \nabla \frac{dU(p)}{dp_i}(x) \right) = -\nabla \cdot \left(\frac{\partial \kappa}{\partial p_i}(x; p) \nabla U(p)(x) \right), \quad (x, p) \in D \times P.$$

303 *For fixed p system of equations is well-posed, guaranteeing the existence of unique*
 304 *solutions $\frac{dU}{dp_i} \in H_0^1(D)$, $i = 1, \dots, m$.*

305 *Once these m PDEs have been solved, the computed solutions can be substituted*
 306 *into (2.1) to determine $\frac{dg}{dp}$. To accomplish this we are required to compute the de-*
 307 *rivatives $\frac{\partial g}{\partial u}$ and $\frac{\partial g}{\partial p}$. Note that in this case g is independent of p , and it is further*
 308 *straightforward to show that*

$$309 \quad \frac{\partial g}{\partial u}[u](v) = \frac{dg}{du}[u](v) = \frac{\langle u, v \rangle_2}{\|u\|_2}.$$

310 *Once again, note that this is a linear operator in v , but is nonlinear in u . We therefore*
 311 *have that*

$$312 \quad \frac{dg}{dp_i}[U(p), p] = \frac{1}{\|U(p)\|_2} \left\langle U(p), \frac{dU}{dp_i} \right\rangle_2$$

313 *for the derivatives $\frac{dU}{dp_i}$ identified by solution of (2.5).*

314 The central challenge with the forward approach, which motivates the adjoint
 315 approach that will be presented in the next section, is the dependence of the forward
 316 sensitivity equation (2.2) on the dimension of the parameter space: solving for $\frac{dU}{dp}$ re-
 317 quires the solution of $\dim(P)$ PDEs. In many practical problems the parameter space
 318 is extremely large; thus, a method for computing the sensitivities that is independent
 319 of the dimension of the parameter space is also of interest.

320 **2.3. Adjoint Sensitivity Analysis.** Adjoint sensitivity analysis begins by in-
 321 troducing the operator $\lambda \in \mathcal{L}(\mathcal{F}, \mathcal{G})$. Supposing that $\dim(\mathcal{G}) = n$, we can express g
 322 as (g_1, \dots, g_n) and consequently $\lambda = (\lambda_1, \dots, \lambda_n)$ where $\lambda_i \in \mathcal{F}^*$ for $i = 1, \dots, n$. For
 323 fixed p , the auxiliary term λ is selected to solve

$$324 \quad (2.6) \quad \lambda_i \frac{\partial F}{\partial u}[U(p), p] = \frac{\partial g_i}{\partial u}[U(p), p], \quad i = 1, \dots, n.$$

325 Assuming this is a unique solution λ exists, one can then recover the sensitivity of the
 326 quantity of interest g as follows

$$327 \quad (2.7) \quad \frac{dg_i}{dp} = -\lambda_i \frac{\partial F}{\partial p} + \frac{\partial g_i}{\partial p}, \quad i = 1, \dots, n.$$

328 which provides a computable expression for the local sensitivities.

329 We note that compared to subsection 2.2 which, in the finite-dimensional case,
 330 necessitates $m = \dim(P)$ solutions of the forward sensitivity equation, the adjoint
 331 system requires $n = \dim(G)$ solutions of the adjoint sensitivity equation. In typical
 332 situations where $n \ll m$ then there is a clear computational benefit to this approach.

333 **EXAMPLE 2.3** (Elliptic PDE: Adjoint Sensitivity Analysis). *Recalling $\frac{\partial F}{\partial u}$ and*
 334 *$\frac{\partial g}{\partial u}$ as derived in Example 2.2, the problem that must be solved to identify $\lambda \in H_0^1(D)$*
 335 *such that*

$$336 \quad (2.8) \quad \nabla \cdot (\kappa^\top(x; p) \nabla \lambda(x)) = \frac{U(p)}{\|U(p)\|_2}.$$

337 *Once λ has been determined, referring again to the derivation in Example 2.2 we*
 338 *have that*

$$339 \quad \frac{\partial g}{\partial p_i} = - \int \nabla \lambda(x) \cdot \frac{\partial \kappa}{\partial p_i}(x; p) \nabla U(p)(x) dx$$

340 *which is real-valued, as required. Again note that in the equation that determines λ ,*
 341 *$U(p)$ appears on the right-hand-side, so for each value of p for which sensitivities*
 342 *are required the PDE must be solved. Nevertheless the fact that in this example only*
 343 *a single system needs to be solved for each p makes the adjoint method significantly*
 344 *cheaper to apply when $m = \dim(P)$ is large.*

345 In the next section we will describe the new probabilistic approaches to both
 346 forward and adjoint sensitivity analysis, each of which operates with a constrained
 347 computational budget.

348 **2.4. Probabilistic Numerical Methods for PDEs.** When applied to PDEs,
 349 there is a marked similarity between this work and probabilistic numerical methods²
 350 applied to linear PDEs. In this section we will discuss these methods, and the simi-
 351 larity to the present approach. Broadly speaking these methods begin by placing a
 352 Gaussian prior on the function space occupied by the solution to the PDE. Finite-
 353 dimensional information about the unknown solution is then produced by projecting
 354 the linear PDE through a set of d functionals, referred to as *information functionals*
 355 in this work. The conjugacy of Gaussian distributions with linear projections can then
 356 be exploited to write down the posterior distribution in closed-form. For a detailed
 357 introduction to this perspective see [13, Chapter 5], in which it is referred to as the
 358 *probabilistic meshless method* (PMM).

359 This approach is equivalent to symmetric collocation with radial basis functions
 360 [55, 11], in that it is possible to construct the prior such that the posterior mean from
 361 PMM coincides with the estimator for the solution of the PDE produced in symmetric
 362 collocation. To our knowledge this approach was first presented in [55, Chapter 16],
 363 and extended in [11] to refine the error analysis, as well as explore applications in
 364 stochastic PDEs. In symmetric collocation the posterior distribution itself is not of
 365 interest, but the error analysis that appears in those works is relevant here as it
 366 provides an important interpretation for the posterior covariance. Specifically, the
 367 bound that appears in [55] connects the error to an object referred to as the *power*
 368 *function*, which can be shown to be directly connected to the posterior covariance
 369 that appears in the PMM.

370 In addition to the PMM, other works that could be interpreted as probabilistic
 371 numerical methods for PDEs include a series of papers that introduced *gamblets* for
 372 the solution of PDEs with rough coefficients [39, 41, 40]. These papers construct a
 373 probabilistic solution to the PDE in a broadly similar way to [13], but with several
 374 distinct differences. Firstly, the probability model is motivated by a game theoretic
 375 argument rather than Bayesian reasoning, though the ultimate conditioning procedure
 376 arrived at is equivalent. Secondly, the information about the solution is constructed
 377 in a distinctly different way, by projecting the defining equations of the PDE against a
 378 hierarchical basis formed by a nested partitioning of the domain, whereas in the PMM
 379 and in symmetric collocation it is obtained by evaluating those equations at a set of
 380 points referred to as *collocation points*. However this results in a very different error
 381 analysis, since collocation methods typically bound the estimation error in terms of
 382 the fill distance of these collocation points, whereas in gamblet-based methods, since
 383 there is no analogue of these points, a different approach must be adopted.

384 The chief similarities of these approaches to the approach presented in this paper
 385 is that, when the system defined by F is a PDE, the sensitivity equations will involve
 386 solving a system of PDEs. In this setting the approach that we describe is similar in

²See [25] for a high-level introduction, and [37] for a thorough literature review.

principal to the approaches we describe above, in that for a particular choice of prior and information, the method we employ will be equivalent to these methods. There are several distinct differences however. Firstly, it is possible that the system described by F is *not* a PDE, and indeed in this work we will explore sensitivity analysis for ODEs in addition to PDEs. While there exist probabilistic numerical methods for solving ODEs, they typically make approximations to account for nonlinearity which are not required in this work, as the systems which must be solved in sensitivity analysis are linear. Secondly, in the PDE case we do not make specific assumptions on the form of the information functionals, as these will typically be problem specific. Thirdly, by formulating the sensitivity equations as a single (degenerate) PDE on the joint space $D \times P$, the continuity of the sensitivities with respect to p is exploited to permit implicit interpolation of the sensitivities across different values of p . And lastly, the focus of this paper is on computing sensitivities, *not* on the solution of the PDE itself, which is assumed to be obtained by some classical numerical solver.

3. Probabilistic Approaches. In this section we will present two probabilistic approaches to computing parameter sensitivities. Each allows a user to restrict the amount of computational effort expended and still obtain an estimate of the sensitivities, while also providing an estimate of the error incurred as a result. Familiarity with Gaussian processes is assumed for this section; we refer the unfamiliar reader to the introduction given in [43]; see also [8] for a more mathematical treatment.

We will assume that there exist reproducing kernel Hilbert spaces (RKHSs) \mathcal{U}'_P , \mathcal{F}'_P such that \mathcal{U}'_P is dense in \mathcal{U}_P and \mathcal{F}'_P is dense in \mathcal{F}_P . Let $\mathcal{U}_{\partial P} = \left\{ \frac{\partial u}{\partial p} : u \in \mathcal{U}_P \right\}$ and let $\mathcal{U}'_{\partial P}$ be defined analogously for \mathcal{U}'_P . It will also be assumed that g is a functional, so that $\mathcal{G} = \mathbb{R}$; this last assumption can readily be generalised, and is made to simplify the presentation.

3.1. Probabilistic Forward Sensitivity Analysis. We first consider forward sensitivity analysis. We begin by modelling prior uncertainty about $\frac{\partial U}{\partial p}$ with the random variable X_F , distributed as $X_F \sim \mu_F = \mathcal{N}(a_F, C_F)$, where $a_F \in \mathcal{U}'_{\partial P}$ and $C_F : \mathcal{U}'_{\partial P} \rightarrow \mathcal{U}'_{\partial P}$ is a positive-definite covariance operator. It will be assumed that $\mu_F(\mathcal{U}_{\partial P}) = 1$. When $\dim(P) < \infty$ this prior takes the form of a vector-valued Gaussian process prior [1]. In the infinite-dimensional setting, we note that a discretisation of the parameter space will nevertheless be required for computational purposes, resulting in a parameter space that is effectively finite-dimensional, though a finite-dimensional parameter space is not strictly required for the theoretical results presented herein.

To obtain a posterior belief over the forward sensitivities, this prior will be conditioned on observations of (2.2). Let $\tilde{\mathcal{I}}_{F,1}, \dots, \tilde{\mathcal{I}}_{F,d}$ be such that $\tilde{\mathcal{I}}_{F,j} \in (\mathcal{F}^m)^*$ for $j = 1, \dots, d$ and let $\{p_1, \dots, p_d\} \subset P$. Let X_F be a random variable with law μ_F . Note that the prior distribution μ_F implies a prior distribution over $\frac{dg}{dp}$ by projecting through the linear map given in (2.1); this will be denoted ν_F . By applying each operator $\tilde{\mathcal{I}}_{F,j}$ to (2.2) we obtain

$$(3.1) \quad \tilde{\mathcal{I}}_{F,j} \frac{\partial F}{\partial u}[U(p_j), p_j] X_F = -\tilde{\mathcal{I}}_{F,j} \frac{\partial F}{\partial p}[U(p_j), p_j]$$

which, under the assumptions made at the start of this section, yields the information

$$f_{F,j} = -\tilde{\mathcal{I}}_{F,j} \frac{\partial F}{\partial p}[U(p_j), p_j]$$

431 where $f_{F,j} \in \mathbb{R}$. Let $f_F \in \mathbb{R}^d$ be the vector with $[f_F]_j = f_{F,j}$.

432 It is more mathematically convenient to think of the $\tilde{\mathcal{I}}_{F,j}$ in terms of functionals
 433 defined on $\mathcal{U}_{\partial P}$. To this end, let $\mathcal{I}_{F,j} \in \mathcal{U}_{\partial P}^*$ be defined by

$$434 \quad \mathcal{I}_{F,j} \frac{\partial u}{\partial p} = \tilde{\mathcal{I}}_{F,j} \frac{\partial F}{\partial u} [U(p_j), p_j] \frac{\partial u}{\partial p}(\cdot, p_j).$$

435 We refer to $\mathcal{I}_{F,1}, \dots, \mathcal{I}_{F,d}$ as the *information functionals*, and will assume that the
 436 information functionals are linearly independent.

437 The posterior is obtained by conditioning the prior on the information functionals.
 438 First, introduce the operator $\mathcal{I}_F : \mathcal{U}_P \rightarrow \mathbb{R}^d$, given by

$$439 \quad \mathcal{I}_F \frac{\partial u}{\partial p} = \begin{bmatrix} \mathcal{I}_{F,1} \frac{\partial u}{\partial p} \\ \vdots \\ \mathcal{I}_{F,d} \frac{\partial u}{\partial p} \end{bmatrix}.$$

440 Then we seek to compute $X_F | \mathcal{I}_F X_F = f_F$. Owing to the linearity of \mathcal{I}_F , the resulting
 441 posterior distribution is again Gaussian and is given in the following proposition.

442 **PROPOSITION 3.1** (Probabilistic Forward Sensitivity Analysis). *The posterior*
 443 $X_F | \mathcal{I}_F X_F = f_F$ has law $\bar{\mu}_F$ given by

$$444 \quad \bar{\mu}_F = \mathcal{N}(\bar{a}_F, \bar{C}_F)$$

$$445 \quad \bar{a}_F = a_F + C_F \mathcal{I}_F^\dagger [\mathcal{I}_F C_F \mathcal{I}_F^\dagger]^{-1} (f_F - \mathcal{I}_F a_F)$$

$$446 \quad \bar{C}_F = C_F - C_F \mathcal{I}_F^\dagger [\mathcal{I}_F C_F \mathcal{I}_F^\dagger]^{-1} \mathcal{I}_F C_F.$$

448 The implied posterior distribution over $\frac{dg}{dp}$, denoted $\bar{\nu}_F$, is given by

$$449 \quad \bar{\nu}_F = \mathcal{N}(\bar{g}_F, \bar{G}_F)$$

$$450 \quad \bar{g}_F(p) = \frac{\partial g}{\partial u} [U(p), p] (\bar{a}_F(\cdot, p)) + \frac{\partial g}{\partial p} [U(p), p]$$

$$451 \quad \bar{G}_F(p, p') = \frac{\partial g}{\partial u} [U(p), p] \delta[\cdot, p] \bar{C}_F \delta[\cdot, p']^\dagger \frac{\partial g}{\partial u} [U(p'), p']^\dagger.$$

453 An important note is that even when underlying system described by F is nonlin-
 454 ear, the posterior distribution remains Gaussian owing to the linearity of the Fréchet
 455 derivatives. Choice of prior mean and covariance is highly problem specific, and will
 456 be discussed for the specific examples considered in this paper in [section 5](#). Next we
 457 turn to the adjoint approach.

458 **3.2. Probabilistic Adjoint Sensitivity Analysis.** For the adjoint problem,
 459 the system that must be solved is now (2.6). Since \mathcal{F}'_P is assumed to be an RKHS,
 460 due to the representer theorem (see e.g. [6, Section 4.4]) we have $\lambda f = \langle f, \beta \rangle_{\mathcal{F}}$, where
 461 $f, \beta \in \mathcal{F}'_P$,

462 The proposed approach is as in the previous section. We model uncertainty in β
 463 with the random variable X_A , whose law is $\mu_A = \mathcal{N}(a_A, C_A)$, where $a_A \in \mathcal{F}'_P$ and
 464 $C_A : \mathcal{F}'_P \rightarrow \mathcal{F}'_P$ is a positive-definite covariance operator. Note that this again implies
 465 a distribution ν_A over $\frac{dg}{dp}$ by projecting through the linear map

$$466 \quad \mathcal{J}[p](\beta) = \left\langle \frac{\partial F}{\partial p} [U(p), p], \beta(\cdot, p) \right\rangle_{\mathcal{F}}.$$

467 An important remark, however, is that unless μ_A and μ_F are chosen carefully, the
468 implied distributions ν_A and ν_F will not be equal.

469 To define the information functionals let $\{(e_1, p_1), \dots, (e_d, p_d)\} \subset \mathcal{U} \times P$. Then

$$470 \quad (3.2) \quad \mathcal{I}_{A,j}\beta = \left\langle \frac{\partial F}{\partial u}[U(p_j), p_j](e_j), \beta(\cdot; p_j) \right\rangle_{\mathcal{F}},$$

471 so that $\mathcal{I}_{A,j} \in \mathcal{F}_P^*$. Furthermore note that the information $f_{A,j} := \frac{\partial g}{\partial u}[U(p_j), p_j](e_j)$ is
472 clearly computable. Let f_A and \mathcal{I}_A be defined analogously to previous sections; then
473 the posterior on β is given in the following proposition.

474 **PROPOSITION 3.2** (Probabilistic Adjoint Sensitivity Analysis). *The posterior*
475 *distribution $\beta|f_A \sim \bar{\mu}_A$ is given by*

$$476 \quad \bar{\mu}_A = \mathcal{N}(\bar{a}_A, \bar{C}_A)$$

$$477 \quad \bar{a}_A = a_A + C_A \mathcal{I}_A^\dagger (\mathcal{I}_A C_A \mathcal{I}_A^\dagger)^{-1} (f_A - \mathcal{I}_A a_A)$$

$$478 \quad \bar{C}_A = C_A - C_A \mathcal{I}_A^\dagger (\mathcal{I}_A C_A \mathcal{I}_A^\dagger)^{-1} \mathcal{I}_A C_A.$$

480 *The implied posterior distribution $\bar{\nu}_A$ is given by*

$$481 \quad \bar{\nu}_A = \mathcal{N}(\bar{g}_A, \bar{G}_A)$$

$$482 \quad \bar{g}_A(p) = -\mathcal{J}[p](\bar{a}_A) + \frac{\partial g}{\partial p}[U(p), p]$$

$$483 \quad \bar{G}_A(p, p') = \mathcal{J}[p] C_A \mathcal{J}^\dagger[p']$$

485 Note that the form of the posterior over β is essentially identical to the form
486 of the posterior from [Proposition 3.1](#), modulo the choice of information functionals
487 and prior. Next we will present some theoretical analysis of the forward and adjoint
488 methods.

489 **3.3. Theoretical Analysis.** Our first theoretical result concerns a local error
490 bound for the posterior mean in terms of the posterior covariance. This result is
491 a general result about conditional distributions of Gaussian process, and so is not
492 specific to either the forward or adjoint method; as a result we adopt generic notation.

493 **PROPOSITION 3.3** (Local error bound). *Let $\mu = \mathcal{N}(a, C)$ be the prior, for $a \in$
494 \mathcal{H}_C , and let $\bar{\mu} = \mathcal{N}(\bar{a}, \bar{C})$ be the posterior measure based on observations $\mathcal{I}u^\dagger = f$
495 where $u^\dagger \in \mathcal{H}_C$, $\mathcal{I} \in (\mathcal{H}_C^*)^d$ and $f \in \mathbb{R}^d$. Then we have that, for each $\mathcal{L} \in \mathcal{H}_C^*$*

$$496 \quad |\mathcal{L}\bar{a} - \mathcal{L}u^\dagger| \leq (\mathcal{L}\bar{C}\mathcal{L}^\dagger)^{\frac{1}{2}} \|a - u^\dagger\|_{C^{-1}}.$$

497 The result from [Proposition 3.3](#) is similar to results on error bounds in scattered
498 data approximation with radial basis functions, such as in [\[55\]](#). The term $(\mathcal{L}\bar{C}\mathcal{L}^\dagger)^{\frac{1}{2}}$
499 is analogous to the *power function* [\[55, Section 11.1\]](#), but the focus in that work is on
500 the case when both \mathcal{L} and \mathcal{I}_j are evaluation functionals. In [\[55, Chapter 16\]](#) each of
501 these restrictions is relaxed, however the form of the power function derived in this
502 case is more abstract than presented here.

503 Similar bounds appear in the literature on solution of PDEs by symmetric col-
504 location with radial basis functions (see e.g. [\[55, Section 16.3\]](#), [\[13, 11\]](#)). In these
505 cases it is typically assumed that the $\tilde{\mathcal{I}}_j$ are evaluation functionals, so that the obser-
506 vations are point evaluations of the right-hand-side of the PDE, and that \mathcal{L} is again
507 an evaluation functional. It is then possible to bound $(\mathcal{L}\bar{C}\mathcal{L}^\dagger)^{\frac{1}{2}}$ in terms of the fill

508 distance in the interior and on the boundary of the domain. We have opted to make
 509 minimal assumptions on the form of the information operators and test functions in
 510 [Proposition 3.3](#), to avoid tying the result to a particular numerical method. Further
 511 note that the cited results only apply for fixed p when performing sensitivity analysis
 512 for an elliptic PDE; as a global function of (x, p) the sensitivity analysis equations
 513 may not be elliptic even when for fixed p the underlying PDE is elliptic.

514 The next proposition provides theoretical guarantees for the setting when the
 515 solution $U(p)$ cannot be accessed directly, and instead a numerical estimate is provided
 516 by the map $\hat{U} : P \rightarrow \mathcal{U}$. The natural way to provide such guarantees is by bounding
 517 the distance between the measure conditioned based on $U(p)$ to that based on $\hat{U}(p)$.
 518 This is closely related to results that appear in [[53](#), e.g. Theorem 4.6], though the
 519 results presented therein assume that the two measures have a common dominating
 520 measure, which is not the case in the present setting. A consequence of this is that
 521 the Hellinger metric, which is commonly used to measure distance in the space of
 522 probability measures in the field of uncertainty quantification, is not suitable here.

523 To proceed we introduce the 2-Wasserstein metric, which is suitable for measures
 524 that are mutually singular. Perhaps the most common way to define this metric is
 525 in terms of couplings of probability measures. Let μ_1 and μ_2 be measures on some
 526 abstract normed space \mathcal{V} . Let $\Gamma(\mu_1, \mu_2)$ be the set of couplings of μ_1 and μ_2 , that
 527 is, the set of all Borel probability measures $\pi \in \mathcal{P}(\mathcal{V} \times \mathcal{V})$ with the property that
 528 $\pi(A \times \mathcal{V}) = \mu_1(A)$ and $\pi(\mathcal{V} \times A) = \mu_2(A)$ for each Borel set $A \subset \mathcal{V}$. Then the
 529 2-Wasserstein metric [[54](#), Definition 6.1] is given by

$$530 \quad (3.3) \quad W_2(\mu_1, \mu_2) = \left(\inf_{\pi \in \Gamma(\mu_1, \mu_2)} \int_{\mathcal{V} \times \mathcal{V}} \|v - v'\|^2 d\pi(v, v') \right)^{\frac{1}{2}}.$$

531 We now proceed to state a generic result concerning robustness to approximation
 532 error, which will then be applied to the methods described in [Proposition 3.1](#) and
 533 [Proposition 3.2](#).

534 **PROPOSITION 3.4 (Robustness to Numerical Error).** *Let $\mu = \mathcal{N}(a, C)$ be a*
 535 *Gaussian distribution with associated RKHS \mathcal{H} , for $a \in \mathcal{H}$ and $C : \mathcal{H} \rightarrow \mathcal{H}$ positive-*
 536 *definite. Assume that $\mathcal{I}, \hat{\mathcal{I}}$ are each bounded linear operators from \mathcal{H} to \mathbb{R}^d . Let*
 537 *$\bar{\mu} = \mathcal{N}(\bar{a}, \bar{C})$ be the posterior measure based on observations $\mathcal{I}u^\dagger = f$ where $u^\dagger \in \mathcal{H}$.*
 538 *Let $\hat{\mu} = \mathcal{N}(\hat{a}, \hat{C})$ be the same prior conditioned on observations $\hat{\mathcal{I}}u^\dagger = \hat{f}$. Then it*
 539 *holds that*

$$540 \quad W_2(\bar{\mu}, \hat{\mu}) \leq (C_{\mathcal{I},1} + C_{\mathcal{I},2}) \|\mathcal{I} - \hat{\mathcal{I}}\|_{\mathcal{H} \rightarrow \mathbb{R}^d} + C_f \|f - \hat{f}\|_{\mathbb{R}^d} + \mathcal{O} \left(\|\mathcal{I} - \hat{\mathcal{I}}\|_{\mathcal{H} \rightarrow \mathbb{R}^d}^2 \right)$$

541 where

$$543 \quad C_{\mathcal{I},1} = \|C\|_{\mathcal{H}} \left[\left(\|a\|_{\mathcal{H}} + \|C\|_{\mathbb{H}\mathbb{S}}^{\frac{1}{2}} \right) \left(\|G^{-1}\mathcal{I}\|_{\mathcal{H} \rightarrow \mathbb{R}^d} + \|\hat{G}^{-1}\hat{\mathcal{I}}\|_{\mathcal{H} \rightarrow \mathbb{R}^d} \right) + \|G^{-1}f\|_{\mathbb{R}^d} \right]$$

$$544 \quad C_{\mathcal{I},2} = \alpha \|G^{-1}\|_{\mathbb{R}^d}^2 \|C\|_{\mathcal{H}} \left(\|G^{-1}f\|_{\mathbb{R}^d} + \|\mathcal{I}\|_{\mathcal{H} \rightarrow \mathbb{R}^d} \|\hat{\mathcal{I}}\|_{\mathcal{H} \rightarrow \mathbb{R}^d} \left(\|a\|_{\mathcal{H}} + \|C\|_{\mathbb{H}\mathbb{S}}^{\frac{1}{2}} \right) \right)$$

$$545 \quad C_f = \|C\|_{\mathcal{H}} \|G^{-1}\hat{\mathcal{I}}\|_{\mathcal{H} \rightarrow \mathbb{R}^d}$$

$$546 \quad \alpha = \|\mathcal{I}C\|_{\mathcal{H} \rightarrow \mathbb{R}^d} + \|\hat{\mathcal{I}}C\|_{\mathcal{H} \rightarrow \mathbb{R}^d}$$

$$548 \quad \text{and } G = \mathcal{I}C\mathcal{I}^\dagger, \hat{G} = \hat{\mathcal{I}}C\hat{\mathcal{I}}^\dagger.$$

549 We next prove a corollary of this result which establishes a bound for the error in
 550 the posterior distribution for both forward and adjoint sensitivity analysis as a result
 551 of the need to use $\hat{U}(p)$ rather than having access to $U(p)$ directly.

COROLLARY 3.5. Assume that for each $p \in P$ there exists $\epsilon > 0$ such that

$$\left\| \frac{\partial F}{\partial u}[U(p), p] - \frac{\partial F}{\partial u}[\hat{U}(p), p] \right\|_{\mathcal{U} \rightarrow \mathcal{F}} \leq \epsilon,$$

$$\left\| \frac{\partial F}{\partial p}[U(p), p] - \frac{\partial F}{\partial p}[\hat{U}(p), p] \right\|_{P \rightarrow \mathcal{F}} \leq \epsilon, \quad \text{and}$$

$$\left\| \frac{\partial g}{\partial u}[U(p), p] - \frac{\partial g}{\partial u}[\hat{U}(p), p] \right\|_{\mathcal{U} \rightarrow \mathcal{G}} \leq \epsilon.$$

Further assume that the $\tilde{\mathcal{I}}_{F,j}$ and $\tilde{\mathcal{I}}_{A,j}$ are such that, for all $j = 1, \dots, d$

$$\|\tilde{\mathcal{I}}_{F,j}\|_{\mathcal{F} \rightarrow \mathbb{R}} < M < \infty,$$

$$\|\tilde{\mathcal{I}}_{A,j}\|_{\mathcal{F} \rightarrow \mathbb{R}} < M < \infty.$$

Lastly assume that $\|\cdot\|_{\mathbb{R}^d} = \|\cdot\|_2$.

Let $\hat{\mu}_F$ be the posterior distribution from Proposition 3.1, with $\hat{U}(\cdot)$ substituted for $U(\cdot)$. Likewise let $\hat{\mu}_A$ be the posterior from Proposition 3.2 with the same substitution. Then we have

$$W_2(\bar{\mu}_F, \hat{\mu}_F) \leq (C_{\mathcal{I},1} + C_{\mathcal{I},2} + C_f)M\epsilon\sqrt{d} + \mathcal{O}(\epsilon^2)$$

$$W_2^2(\bar{\mu}_A, \hat{\mu}_A) \leq (C_{\mathcal{I},1} + C_{\mathcal{I},2} + C_f)M\epsilon\sqrt{d} + \mathcal{O}(\epsilon^2)$$

3.4. Comparison of Forward and Adjoint Approaches. We conclude this section with a brief discussion of the relative merits of the forward and adjoint approaches, compared to the classical approach.

Choice of Method. The forward approach requires the user to specify a prior on the parameter space; this is a space of dimension $\dim(P)$. While the space in which the prior is placed for the adjoint problem is less directly connected to the derivative of interest, which might make eliciting a prior more challenging, in the finite-dimensional case reasoning about the correlation structure between the components of $\frac{\partial u}{\partial p}$ for the forward problem may also be challenging. As a result, much as in classical sensitivity analysis, we are inclined to recommend the adjoint approach whenever $\dim(\mathcal{G}) < \dim(P)$, as will often be the case. However if the user has strong prior information about the correlation structure between these components, the forward approach may still perform well. Indeed, in the infinite-dimensional case such information is provided by knowledge about the smoothness of the function p .

Experimental Design. Propositions 3.1 and 3.2 each allow the user to construct a global model for the required derivatives. However in order to perform inference globally, one requires a set of points in P with which to construct the posterior. Both the forward and the adjoint approach suffer from the curse of dimensionality in this respect, since Gaussian processes typically require such designs to be “space-filling”³, and if P is high-dimensional constructing a space filling design will be equally prohibitive in either mode. However in the present paper we focus on application of these methods within iterative optimisation algorithms, so that rather than requiring a space-filling design we only require good estimates of the gradient along the path in parameter space followed by the optimiser. This will be discussed in detail in the next section.

³Since, typically, the rate of convergence of Gaussian processes with this type of information depends on the “fill distance”, i.e. the maximum distance of any point in the space to a design point. See e.g. [55], or [11] in the context of PDEs.

593 **4. Optimisation and Probabilistic Sensitivity Analysis.** We now explore
 594 a potential application of probabilistic local sensitivity analysis, as a way to provide
 595 approximations of gradients in optimisation algorithms. As a starting point we will
 596 consider the most fundamental of gradient-based optimisation algorithms, gradient
 597 descent [15]. In subsection 4.1 we briefly recall the GD algorithm. In subsection 4.2
 598 we describe how probabilistic gradients can be incorporated into the algorithm. Then,
 599 in section 5 we explore the use of this approach in two applications.

600 **4.1. Gradient Descent.** We now describe the GD algorithm. GD is in many
 601 respects a prototypical gradient-based optimisation method, making it a natural start-
 602 ing point for studying the integration of probabilistic gradients into such algorithms.
 603 In GD the goal is to compute a (local) minimiser p^* of a function $g(p)$. To accomplish
 604 this a sequence of points (p^n) , $p^n \in P, n \in \mathbb{N}$ is generated iteratively starting from
 605 some user-defined initial point p^0 and advancing according to

$$606 \quad p^{n+1} = p^n - \gamma^n \frac{dg}{dp}(p^n)$$

607 where γ^n is a parameter of the method known as the *step size* or *learning rate*. Under
 608 specific conditions on f and γ^n it can be shown that $p^n \rightarrow p^*$ (again, a local minimiser)
 609 as $n \rightarrow \infty$; see [34, Section 3.2] for further details. GD is presented as an algorithm
 610 in Algorithm 1 in the supplement.

611 There are various methods for choosing the parameter γ_n . Since the focus of
 612 this work is on the performance when $\frac{dg}{dp}$ is replaced by the probabilistic gradients
 613 introduced in section 3, we will use a probabilistic version of the backtracking line
 614 search method described in [34, Algorithm 3.1], based on the method described in
 615 [32].

616 **4.2. Gradient Descent with Probabilistic Gradients.** We now discuss a
 617 probabilistic modification of GD. Heuristically the approach followed is to replace the
 618 computation of $\frac{dg}{dp}$ with a probabilistic gradient obtained from either Proposition 3.1
 619 or Proposition 3.2; to simplify the exposition we will describe the former, but the ap-
 620 proach is essentially identical in the latter. The approach is presented as an algorithm
 621 in Algorithm 4.1. Essentially, we begin with a prior μ_F which is projected to ν_F as
 622 described in Proposition 3.1. We then construct a sequence of random variables (X_F^n) ,
 623 where X_F^0 has law ν_F , by sequentially updating this prior with information collected
 624 over the course of the optimisation. This provides a posterior distribution over the
 625 gradient which is used in place of $\frac{\partial g}{\partial p}$ in GD. The principal advantages, illustrated in
 626 section 5, are that (i) for each value of p^n , one can often obtain an approximation of
 627 $\frac{dg}{dp}$ that is sufficiently accurate for the purposes of taking a gradient step, at a lower
 628 cost than that of computing $\frac{dg}{dp}$ directly, and (ii) since the posterior is defined over the
 629 entire parameter space, for some values of p^n no inversion problem must be solved to
 630 advance the gradient descent.

631 There are two main issues to address. The first is that that it is well-established
 632 in the literature on stochastic gradient descent that line-search algorithms such as the
 633 BLS routine are not robust to inaccurate gradients. This is discussed in [32]. Since
 634 the gradients we propose to use in this work are also inaccurate, an alternative line-
 635 search strategy for selecting the step sizes γ^n must be adopted in the probabilistic case.
 636 Borrowing from the literature on stochastic gradient descent, our proposed approach
 637 incorporates ideas from the probabilistic line search of [32] into the backtracking line
 638 search from [34, Section 3.2]. The PLS routine is described in Algorithm 4.2.

639 A second issue is that if X_F^n is not sufficiently accurate, the step size γ^n found by
 640 the probabilistic line search will be selected to be below the tolerance ϵ , causing the
 641 algorithm to terminate. To address this we propose to couple the computation of γ^n
 642 with the calculation of the gradient, as described in PROBJAC within [Algorithm 4.1](#).
 643 Once the tolerance has been achieved, we calculate the step size γ according to a
 644 probabilistic version of backtracking line search that will be described presently. If
 645 γ is above the tolerance the procedure returns the current gradient estimate, along
 646 with the posterior distribution and the step size; otherwise, the tolerance δ is reduced
 647 and the conditioning procedure is repeated. This continues until delta is below some
 648 minimum value δ_{\min} , at which point convergence is accepted.

Algorithm 4.1 Probabilistic version of gradient descent. The routines METRIC and INFO are problem specific and must be supplied by the user, with the former assessing the distribution of the currently computed posterior distribution to determine whether it is sufficiently narrow to accept it as a valid gradient and the latter supplying information, iteratively, based on the current distribution and location. The routine CONDITION implements [Proposition 3.1](#). PLS is the probabilistic version of the Armijo line search, and is given in [Algorithm 4.2](#). Of the new parameters, δ reflects how much accuracy is demanded of the posterior at each iteration, δ_{\min} specifies a maximum level of accuracy to protect against numerical instabilities resulting from large Gram matrices in CONDITION, and τ describes how rapidly δ is reduced when a valid descent direction cannot be found.

```

1: procedure PGD( $p^0, g, \mu_F^0, \epsilon, \delta, \delta_{\min}, \tau_1$ )
2:   Compute  $\nu_F^0$  from  $\mu_F^0$  and let  $X_F^0$  be the random variable with law  $\nu_F^0$ 
3:   for  $n = 1, 2, \dots$  do
4:      $s_F^n, X_F^n, \gamma_n \leftarrow$  PROBJAC( $X_F^{n-1}, g, p^{n-1}, \epsilon, \delta, \delta_{\min}$ )
5:     if  $\gamma^n < \epsilon$  then
6:       return  $p^{n-1}$ 
7:     end if
8:      $p^n \leftarrow p^{n-1} + \gamma^n s_F^n$ 
9:   end for
10: end procedure
11: procedure PROBJAC( $X, g, p, \epsilon, \delta, \delta_{\min}, \tau_1$ )
12:   while  $\delta > \delta_{\min}$  do
13:     while METRIC( $X$ )  $> \delta$  do
14:        $\mathcal{I}, f \leftarrow$  INFO( $X, p$ )
15:        $X \leftarrow$  CONDITION( $X, \mathcal{I}, f$ )
16:        $s \leftarrow -\mathbb{E}(X(p)) / \|\mathbb{E}(X(p))\|_2$ 
17:        $\gamma \leftarrow$  PLS( $p, g, X$ )
18:       if  $\gamma < \epsilon$  then
19:          $\delta \leftarrow \tau_1 \delta$ 
20:       else
21:         return  $s, X, \gamma$ 
22:       end if
23:     end while
24:   end while
25: end procedure

```

Algorithm 4.2 Probabilistic line search algorithm. This is essentially a modification of the backtracking line search described in [34, Algorithm 3.1] to account for the fact that the gradient is a random variable rather than a constant. The parameters p, g and X are the parameter value, objective function and current posterior, respectively. The remaining parameters control the behaviour of the algorithm; we have specified sensible defaults for these and assume those defaults are used throughout the text. τ_2 controls how rapidly γ is decreased, while c controls how large a reduction in the objective function is required when a step is taken in the chosen direction and P^{crit} is the probability with which this reduction must be achieved. γ and γ^{min} control the initial and minimum values of γ respectively.

```

procedure PLS( $p, g, X; \tau_2 = 0.5, c = 0.5, P^{\text{crit}}, \gamma = 1, \gamma^{\text{min}} = 10^{-10}$ )
   $s \leftarrow -\mathbb{E}(X)/\|\mathbb{E}(X)\|_2$ 
  while  $\gamma > \gamma^{\text{min}}$  do
     $p_\gamma \leftarrow p + \gamma s$ 
    if  $g(p_\gamma) > g(p)$  then
      continue
    end if
     $Z \leftarrow -c\gamma X^\top s$ 
    if  $\mathbb{P}(Z > g(p_\gamma) - g(p)) < P^{\text{crit}}$  then
      return  $\gamma$ 
    end if
     $\gamma \leftarrow \tau_2 \gamma$ 
  end while
end procedure

```

649 **4.2.1. Discussion.** We now provide some important remarks about the algo-
 650 rithm presented above.

651 *Choice of Direction.* The direction chosen in Algorithm 4.1 at each iteration is the
 652 posterior mean. A natural alternative would be to instead *sample* a direction from the
 653 posterior distribution. This requires only minor modification of the above algorithm,
 654 but empirically was found to perform slightly worse in general; consequently we have
 655 opted to use the posterior mean as the descent direction.

656 *Recycling Information.* Note that the gradient here is computed based on infor-
 657 mation collected at all points p_F^1, \dots, p_F^n , i.e. based on a *global* model for the gradient
 658 as a function of p . Since the sequence (p_F^n) will increasingly concentrate in a region
 659 of p_* as n increases, one expects that the prior $\bar{\mu}_F^{n-1}$ will be an increasingly accu-
 660 rate predictor for the gradient $\frac{dg}{dp}(p_F^n)$ as n increases. This means that once some
 661 computational effort has been expended to obtain a relatively accurate gradient, it
 662 is possible for PROBJAC to perform many further iterations based on this gradient
 663 without needing calls to CONDITION, as we shall see in section 5.

664 *Linearly Independent Information.* A global model introduces some additional
 665 burden to ensure that \mathcal{I}_F^n is linearly independent of $\mathcal{I}_F^1, \dots, \mathcal{I}_F^{n-1}$, both to maximise
 666 the amount of new information obtained at each p^n and to ensure that the linear
 667 system that must be solved to compute the posterior does not become singular. Thus,
 668 INFO must be carefully designed to ensure that the information returned is not too
 669 highly correlated with information already observed.

670 *Computational Cost.* To compute the posterior distributions from Proposition 3.1
 671 and Proposition 3.2, it is necessary to compute the matrix $M = (ICT^\dagger)^{-1}IC$ by solv-

ing the linear system $\mathcal{I}CT^\dagger M = \mathcal{I}C$. To accomplish this one typically computes a Cholesky factorisation of $\mathcal{I}CT^\dagger$, which becomes computationally intensive once many information functionals have been collected. However, we note that the sequential nature of the algorithm proposed is such that, rather than recomputing the full factorisation at each iteration of PROBJAC, one can use an updating formula for the factorisation such as presented in [38, Appendix B]; this is described in detail in [section S2](#). In brief, one must only compute the Cholesky factorisation of a smaller matrix, whose dimension is only the same size as the dimension of the *new* information, which naturally dramatically reduces the cost of computing the probabilistic gradients.

The other factor that influences the cost is the size of $\mathcal{I}C$, and since this defines *how many* linear systems must be solved, it may be that ultimately the cost of assembling the posterior $\bar{\mu}_F^n$ exceeds that of simply computing $\frac{dg}{dp}(p_F^n)$ despite the efficient updating formula for the factorisation. Thus in practise we propose that the PROBJAC is used only to perform the initial iterations, and that when the method is determined to be close to the truth, or the cost of constructing the posterior is too great, we revert to classical GD to complete the optimisation. In [section 5](#) we adopt the crude rule of thumb that PROBJAC is terminated when the dimension of $\mathcal{I}CT^\dagger$ exceeds 10,000, though this is never exceeded in practise for one of the two examples examined. In future work more sophisticated switching schemes will be explored.

Choice of Metric. The routine METRIC must assess whether the posterior distribution at a particular iteration is sufficiently accurate for the probabilistic gradient to be accepted as a valid direction for the gradient descent. To determine this we focus on the width of the posterior covariance, and in this work we exclusively use the square-root of the trace of the posterior covariance, $\sqrt{\text{trace}(G_F)}$ as a proxy for the width. An exploration of other choices is not expected to affect the performance of the algorithm dramatically, and is reserved for future work.

Choice of Information Functionals. Lastly, we note that we have not yet discussed the selection of information functionals in INFO. We expect this to be highly problem dependent. We make a proposal in the next section that appears to be well adapted to the two examples presented therein, but do not expect that there exists a unique optimal choice of information for all settings.

5. Applications. In this section we apply [Algorithm 4.1](#) to compute the maximum *a-posteriori* (MAP) point in Bayesian inversion problems for two problems. In [subsection 5.1](#) we seek to infer a small number of parameters of an ODE using the forward approach, and in [subsection 5.2](#) inference of a larger number of parameters of a challenging PDE using the adjoint approach.

5.1. FitzHugh—Nagumo Model. As a first example we examine the problem of inferring the parameters for the Fitzhugh—Nagumo model [19], a nonlinear oscillatory ODE. Since this problem has four parameters, we use the forward approach from [subsection 3.1](#).

5.1.1. Problem Definition. The equations that define the FitzHugh—Nagumo model are

$$\frac{dv}{dt} = v - \frac{v^3}{3} - w + I \qquad \frac{dw}{dt} = \frac{v + a - bw}{\tau}$$

where $a, b, I, \tau \in \mathbb{R}_+$ are parameters of the model. We concatenate the parameters as $p = [I, a, b, \tau]^\top \in \mathbb{R}^4 = P$. The solution to this system of ODEs for

719 $p^* = [0.5, 0.8, 0.7, 12.5]^\top$ is shown in the supplement in [Figure S1a](#), while sensitiv-
 720 ities are displayed in [Figures S1b](#) and [S1c](#). The solution space \mathcal{U} is a space of once
 721 differentiable functions $u : D \rightarrow \mathbb{R}^2$, where $D = [0, T]$ for some $T > 0$. A refor-
 722 mulation of this problem in terms of the constraint function $F(u, p)$ can be found in
 723 [section S4](#), along with the form of its derivatives $\frac{\partial F}{\partial u}, \frac{\partial F}{\partial p}$.

724 To set up the inference problem we generated data for true parameter values
 725 p^* by evaluating the $v(t_i^{\text{data}}; p^*)$ at times $t_i^{\text{data}} = i, i = 1, \dots, 20$. These locations
 726 are distinguished as dashed gray lines in [Figure S1](#) in the supplement. Observations
 727 were then corrupted with centred Gaussian noise with standard deviation 10^{-2} , i.e.
 728 $y_i = v(t_i^{\text{data}}; p^*) + \xi_i$ where $\xi_i \sim \mathcal{N}(0, 10^{-2})$ IID. The prior distribution over the
 729 parameters was set to be log-Gaussian with mean $m_p = [1, 1, 1, 10]^\top$ and covariance
 730 I . The objective function is twice the negative logarithm of the likelihood multiplied
 731 by the prior, and is thus given by

$$732 \quad g(p) := \frac{1}{\gamma^2} \sum_{i=1}^M (v(t_i^{\text{data}}; p) - y_i)^2 + (\log(p) - \mu)^\top \Sigma^{-1} (\log(p) - \mu)$$

733 **5.1.2. Probabilistic Gradient Descent.** To apply the probabilistic gradient
 734 descent algorithm from [Algorithm 4.1](#) we must first specify the prior over $\frac{dU}{dp}$. Since
 735 the parameter space is four-dimensional and \mathcal{U} is a space of vector-valued functions,
 736 formally $\frac{dU}{dp}$ is $\mathbb{R}^{2 \times 4}$ -valued. For convenience, we place a prior on $X : D \times P \rightarrow \mathbb{R}^8$,
 737 and form $\frac{dU}{dp}$ as

$$738 \quad \frac{dU}{dp} = \begin{bmatrix} X_{1:4}^\top \\ X_{5:8}^\top \end{bmatrix}$$

739 where $X_{i:j}$ denotes components i to j of X . Noting that the posterior covariance
 740 is independent of the data, we assume an independent and identical prior over each
 741 column of $\frac{dU}{dp}$, so that the inference is identical but for the distinct right-hand-side
 742 for each component of p in the posterior mean of [Proposition 3.1](#).

743 Since the initial condition is independent of p , this prior was taken to be $X \sim$
 744 $\mathcal{N}(\mathbf{0}, k)$ where

$$745 \quad k((t, p), (t', p')) = Cq(t)q(t')k_{5/2}([t, p]^\top, [t', p']^\top \sigma, L)$$

$$746 \quad (5.1) \quad k_{5/2}(r, r'; \sigma, L) = \sigma^2 \left(1 + \sqrt{5}d(r, r'; L) + \frac{5}{3}d(r, r'; L)^2 \right) \exp \left(-\sqrt{5}d(r, r'; L) \right)$$

$$747 \quad d(r, r'; L) = \sqrt{r^\top L^{-1} r'}$$

$$748 \quad C = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

$$749 \quad q(t) = t.$$

751 Multiplication by the linear functions $q(t)$ ensures that there is no uncertainty at
 752 $t = 0$, where the sensitivity is known to be zero.

753 The kernel $k_{5/2}$ in [\(5.1\)](#) is a member of the Matérn family [[43](#), Section 4.2] and is
 754 the covariance kernel for a prior over functions with at least two continuous derivatives.
 755 To ease computation the length-scale matrix L was selected to be diagonal, $L =$
 756 $\text{diag}(\ell)$ for $\ell \in \mathbb{R}^6$. This parameter was further restricted to $\ell = [\ell_x \mathbf{1}_2, \ell_p \mathbf{1}_4]$ where
 757 $\ell_x, \ell_p \in \mathbb{R}$. The scalars σ, ℓ_x and ℓ_p were then selected by maximising the marginal
 758 likelihood of an initial candidate design (see e.g. [[43](#), Section 5.4]). This was obtained

759 by sampling a set of candidate parameters p_i^{calib} , $i = 1, \dots, 5$ from the prior over
 760 the parameters and defining the corresponding evaluation functionals $\tilde{I}_{ij} = \delta[i]$, $i =$
 761 $1, \dots, 20$ (i.e. using equally spaced points inside the spatial domain). The parameter
 762 ρ , which describes the degree of prior covariance between the components $u_1 = v$ and
 763 $u_2 = w$, was fixed to 0.5.

764 For this problem it was convenient to restrict the information functionals to be
 765 evaluation functionals, i.e. $\tilde{I}_i = \delta[t_i^{\text{info}}]$. The points t_i^{info} were restricted to a fine
 766 grid of 1000 points in $(0, T]$, denoted $t_1^{\text{info}}, \dots, t_{1000}^{\text{info}}$. To choose the next information
 767 functionals at iteration n within the function INFO in [Algorithm 4.1](#), we choose new
 768 conditioning locations within this set by attempting to minimise a heuristic based on
 769 the fill distance which often appears as an upper bound in Gaussian process regression
 770 problems. To be specific, we begin by constructing an augmented point set:

$$771 \quad z_{ij} = \begin{bmatrix} t_i^{\text{info}} \\ p^j \end{bmatrix}.$$

772 for $j = 1, \dots, n$ denoting the iteration number in PROBJAC and p^j the corresponding
 773 parameter value for that iteration. The information functionals were then selected
 774 to be the \tilde{I}_j for which the distance between z_{in} and $z_{i'j}$, $i, i' = 1, \dots, 1000$, $j =$
 775 $1, \dots, n - 1$, is maximised.

776 **5.1.3. Results.** The paths taken by the probabilistic optimiser are contrasted
 777 with classical gradient descent in [Figure 5.1](#). [Figure 5.1a](#) shows the value of $g(p^n)$,
 778 while [Figure 5.1b](#) shows the distance from the minimum obtained by gradient descent.
 779 All of the methods were started from the initial parameter value $p_0 = m_p$, and the
 780 GD tolerance was set to $\epsilon = 10^{-6}$. The threshold δ was varied from 1 (representing a
 781 high level of allowed error in the posterior gradient estimate) to 0.001 (representing a
 782 low level of allowed error). In each case δ_{\min} was set to 10^{-6} . For $\delta = 0.9, 0.5, 0.1$ the
 783 performance of the probabilistic approach is initially worse, as expected, though as
 784 the iterates near p^* the performance of the probabilistic approaches improves. Inter-
 785 estingly, for $\delta = 1$ and $\delta = 0.01$ the probabilistic approach actually seems to initially
 786 converge *faster* than the classical approach. This should not generally be expected,
 787 though we note that since the GD directions have no particular optimality properties
 788 nothing prevents an approximate method from achieving faster convergence.

789 [Figure 5.1c](#) tracks the amount of data collected (i.e. the size of f_F^n) as a function
 790 of the iteration number. This exhibits the expected behaviour of increasing inversely
 791 proportional to δ . However it is noteworthy that even in the strictest case, $\delta = 0.001$,
 792 only 3000 evaluations of $\frac{\partial F}{\partial p}$ are required over the course of 9840 iterations to perform
 793 almost as well as as gradient descent. For context computing the gradient $\frac{dg}{dp}$ using
 794 the DOP853 algorithm [[23](#), Section II] method as implemented in `scipy` required an
 795 average of 781 evaluations of $\frac{\partial F}{\partial p}$ *per iteration* of gradient descent, with a total of over
 796 1.5 million evaluations over the course of the 2013 iterations performed with exact
 797 gradients. While $\frac{\partial F}{\partial p}$ is cheap to evaluate in this example, in a setting in which this
 798 was a bottleneck it is clear that the probabilistic method would be preferable. Further
 799 note that while 3000 evaluations of $\frac{\partial F}{\partial p}$ were required, as noted in [subsection 4.2.1](#) this
 800 does not translate directly to inversion of a 3000×3000 Gram matrix, as the updating
 801 formula for Cholesky factorisations was exploited.

802 **5.2. Groundwater Flow Model.** We now consider a linear PDE that describes
 803 the steady-state flow of fluid through a porous medium. In this section the parameter
 804 is formally function-valued. Since after discretisation its dimension can be large, the

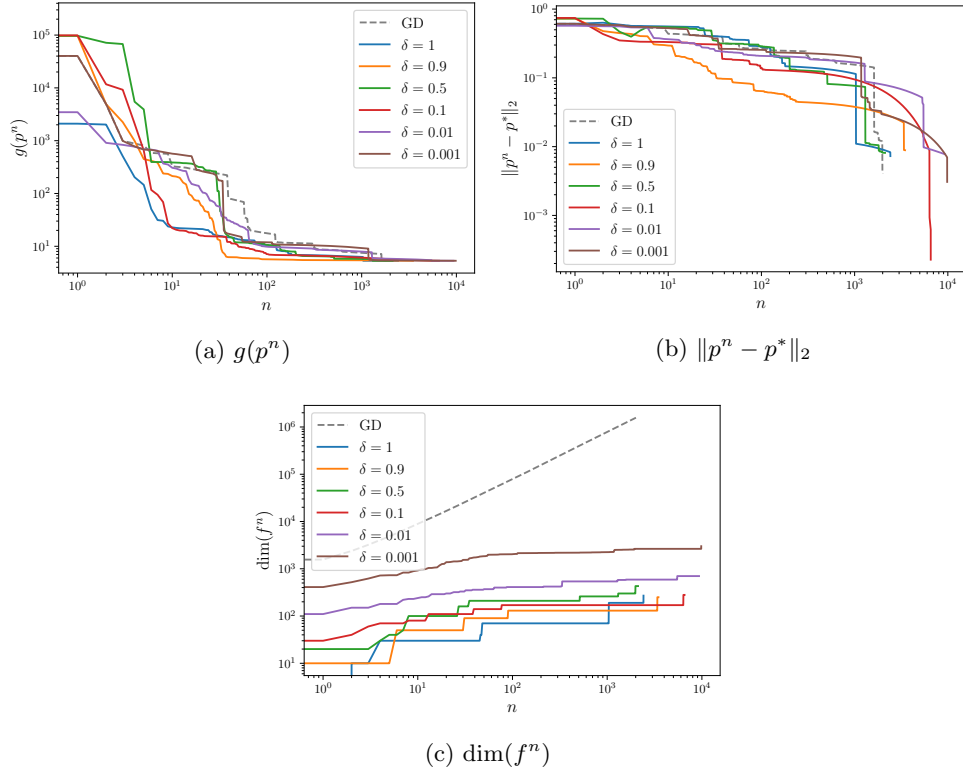


Fig. 5.1: Performance of the probabilistic gradient descent algorithm on the FitzHugh—Nagumo model described in subsection 5.1 as the parameter δ , which roughly controlling the accuracy demanded of the probabilistic gradient estimate, is varied. Here n is the iteration number. Figure 5.1c shows the value of the quantity of interest g , in this case the value of the negative log-target in a Bayesian inference problem described in subsection 5.1. Figure 5.1b shows the distance from the parameter at iteration n to the true MAP point. Figure 5.1c shows the dimension of the matrix inversion problem that was solved in order to compute the posterior distribution.

805 adjoint approach is adopted.

806 **5.2.1. Problem Definition.** For a fixed value of the parameter p , the forward
807 model is given by

$$\begin{aligned}
 808 \quad & -\nabla \cdot (p(x)\nabla u(x)) = 0 & x \in D \\
 809 \quad & u(x) = x & x_2 = 0 \\
 810 \quad & u(x) = 1 - x & x_2 = 1 \\
 811 \quad & \frac{\partial u}{\partial x_1} = 0 & x_1 = 0 \text{ or } 1 \\
 812
 \end{aligned}$$

813 Here the domain $D = [0, 1]^2$ and $p : D \rightarrow \mathbb{R}$. We assume that $p(x) > 0$ for all $x \in D$.

814 The solution $u(x)$ was obtained by discretising the domain above with FEM on

815 a fine triangular meshing of the unit square based on a grid of 32×32 points using
 816 piecewise-linear basis functions. The mesh is depicted in [Figure S2a](#) in the sup-
 817 plement, and the discretisation results in a finite-dimensional approximation of the
 818 solution $u(x)$ in with 1089 degrees of freedom. The solution to the PDE above for the
 819 parameter value $p = 1$ is depicted in [Figure S2b](#), again found in the supplement.

820 The parameter is defined to be piecewise constant over supersets of the cells of
 821 this mesh, defined by grouping the cells based on a subdivision of the domain into
 822 squares. For a parameter $N > 1$ these are obtained by placing down a regular grid
 823 of N^2 points, with $N + 1$ equispaced points along each axis. The points of this grid
 824 form the vertices of the N^2 parameter cells. In [Figure S2a](#), the parameter cells for
 825 $N = 4$ are surrounded by green lines.

826 To construct the inverse problem, we use a Gaussian prior $p \sim \mathcal{N}(\mu, \Sigma) =: \pi(p)$
 827 with $\mu = 5\mathbf{1}_{N^2}$, where $\mathbf{1}_{N^2}$ here denotes the vector of ones in \mathbb{R}^{N^2} . Letting x_i^{param}
 828 denote the centroid of cell i according to some arbitrary ordering of the cells, $i =$
 829 $1, \dots, N$, the covariance is given by $\Sigma_{ij} = k(x_i^{\text{param}}, x_j^{\text{param}})$, where k is the Matérn
 830 $5/2$ kernel given in (5.1), with amplitude and length-scale each set to 1. The data-
 831 generating parameter p^* was sampled from the prior over p . To define the likelihood,
 832 data was obtained by taking direct measurements of the solution $u(x; p^*)$ at locations
 833 $x_1^{\text{data}}, \dots, x_M^{\text{data}}$ where $M = 25$ the x_j^{data} are the nearest mesh points to points on a
 834 regular 5×5 grid starting at $(0.1, 0.1)$ and ending at $(0.9, 0.9)$. The points of this
 835 grid are shown in [Figure S2b](#) in the supplement. Let $\tilde{y} \in \mathbb{R}^M$ be the vector with
 836 $\tilde{y}_j = u(x_j^{\text{data}}; p^*)$. These points were corrupted with IID Gaussian noise $\xi_j \sim \mathcal{N}(0, \gamma)$,
 837 $\gamma = 0.01$ $j = 1, \dots, M$ to obtain data $y = \tilde{y} + \xi$. Denoting the likelihood by $\pi(p|y, u)$
 838 with dependence on u emphasised, the QoI for gradient descent was then given by
 839 $g(u, p) = -2 \log \pi(p|y, u)\pi(p)$, i.e.

$$840 \quad (5.2) \quad g(p) := \frac{1}{\gamma^2} \sum_{i=1}^M (u(x_i^{\text{data}}; p) - y_i)^2 + (p - \mu)^\top \Sigma^{-1} (p - \mu)$$

841 **5.2.2. Probabilistic Gradient Descent.** To test the algorithm described in
 842 [subsection 4.2](#) we attempt to compute the MAP point of the posterior distribution
 843 for the inverse problem described above. Owing to the potentially high dimension
 844 of the problem to be solved, the adjoint approach was used. For the prior we used
 845 $\beta \sim \mathcal{N}(0, k)$, where k is given by

$$846 \quad k((x, p), (x', p')) = k_{52}((x, p), (x', p'); \sigma, \ell) q(x_2) q(x'_2).$$

847 Here $q(x) = 1 - (2x - 1)^2$, so that $q(0) = q(1) = 0$, ensuring that the relevant
 848 boundary condition is encoded in the prior since we note that the boundary conditions
 849 do not depend upon p . Thus, the prior is formally over functions from \mathbb{R}^{N^2+2} to \mathbb{R} ,
 850 though since the problem has been discretised with the finite-element method the
 851 discretised prior is finite-dimensional. Strictly speaking to project the prior into the
 852 finite-element space requires computing integrals of the form $\int k(x, x') \phi_j(x) dx$ for
 853 $j = 1, \dots, 1089$, however since these integrals do not generally have a closed-form we
 854 opt to approximate them as $\int k(x, x') \phi_j(x) dx \approx k(x_j, x')$ where x_j is the nodal point
 855 corresponding to the basis function ϕ_j .

856 For the parameters of the prior, a separate constant length-scale was assigned
 857 to the spatial variables and the parameters, denoted ℓ_x and ℓ_p respectively, i.e.
 858 $\ell = [\ell_x \mathbf{1}_2, \ell_p \mathbf{1}_{N^2}]$. The amplitude σ and the length-scale ℓ_p were again selected
 859 by maximising the marginal likelihood of these parameters given a candidate de-
 860 sign obtained again by sampling a set of candidate parameters p_i^{calib} , $i = 1, \dots, 10$,

861 from the prior over parameters, and choosing corresponding information functionals
 862 $\tilde{\mathcal{I}}_{ij}u = \int_D u(x)\phi_j(x)dx$. Here the ϕ_j are the finite element basis functions correspond-
 863 ing to the nearest mesh points to a regular 10×10 grid of points within D , with basis
 864 functions on the top and bottom boundaries excluded.

865 For the remaining parameter, ℓ_x , we note that since in (5.2) g depends only on
 866 the value of u at the points x_i^{data} , we therefore have that $\frac{\partial g}{\partial u}$ is zero everywhere but
 867 at these locations. Since this function is so rough, it is impossible to infer the spatial
 868 length-scale ℓ_x from evaluations of it. As a result, we opted to fix $\ell_x = 0.2$, based on
 869 the observed smoothness of the solution to the adjoint equations.

870 For the information functionals we selected $\tilde{\mathcal{I}}_j u = \int u(x)\phi_j(x)dx$, i.e. projection
 871 against the j^{th} finite element basis function. This is straightforward to implement
 872 since after discretisation it is simply projection against the canonical basis vector e_j^\top .
 873 The function INFO was implemented similarly to in subsection 5.1, with the fine grid
 874 of points now consisting of the mesh locations which the basis functions correspond
 875 to, again excluding points on the top and bottom boundaries. However, to ensure
 876 that the information f is nonzero, we enforce that when $\text{METRIC}(X^n) > \delta$, the first
 877 locations to be conditioned upon are those basis functions corresponding to x_i^{data} .

878 **5.2.3. Results.** The results of the optimisation are displayed in Figure 5.2. As
 879 in subsection 5.1 one can clearly see the behaviour of the method reverting to that
 880 of gradient descent as the size of δ is decreased. Further, performance appears to
 881 be broadly similar as the parameter dimension increases, reflecting that only a single
 882 function $\beta(x)$ must be learned, rather than $\frac{dU}{dp_i}$ for $i = 1, \dots, \dim(P)$ as would be
 883 required in the forward approach. Thus, the output dimension of the inferred function
 884 is independent of the parameter dimension. While the *input* dimension does grow
 885 with $\dim(P)$, for the purposes of the gradient descent algorithm, at iteration n only
 886 the quality of inferences at and in the region of p^n is relevant. Since these points
 887 concentrate near p^* , performance does not appear to decay as the input dimension
 888 grows.

889 Figure 5.3 compares the cost of the probabilistic approach with that of the clas-
 890 sical approach, for $N = 2$, $\dim(P) = 4$, by plotting the size of the matrix whose
 891 Cholesky factorisation that must be computed at each iteration in order to update
 892 the Cholesky factorisation of the Gram matrix with novel information, as discussed in
 893 subsection 4.2.1. We note that in general more information seems to be required than
 894 for the Fitzhugh-Nagumo example, so that the $10,000 \times 10,000$ limit on the size of
 895 the Gram matrix discussed in subsection 4.2.1 is generally what causes the algorithm
 896 to terminate, though from Figure 5.2 it is clear that nevertheless PROBJAC is *close*
 897 to convergence when this occurs. The higher cost is perhaps due to the fact that the
 898 right-hand side, $\frac{\partial g}{\partial u}$, is highly localised in this example. It is nevertheless the case that
 899 throughout the gradient descent procedure, the size of the inversion problem that must
 900 be computed with the probabilistic approach is significantly smaller than that which
 901 must be computed with the classical approach, though since the matrix inverted in
 902 the classical approach is sparse the costs are not directly comparable. Furthermore, as
 903 in subsection 5.1, for larger values of δ the approach shows the behaviour of being able
 904 to conduct a large number of iterations without needing to collect *any* evaluations of
 905 the right-hand-side, due to the fact that the model is global over parameter space.

906 **6. Conclusion.** In this paper we have presented a probabilistic approach to
 907 computing local sensitivities of differential equation models in both the forward and
 908 adjoint modes. We presented an approach for incorporating these probabilistic gradi-
 909 ents into a gradient descent algorithm, and examined the properties of this algorithm

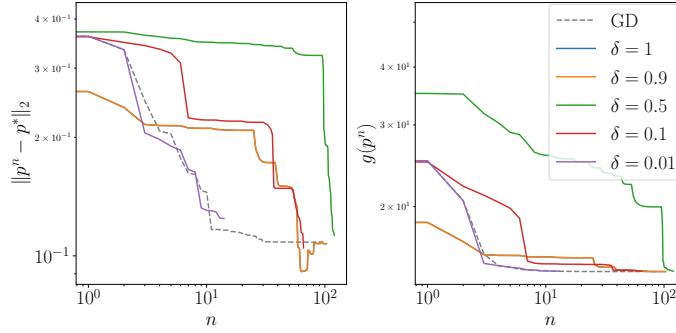
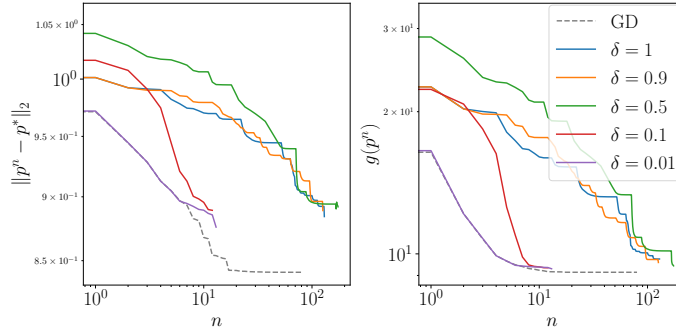
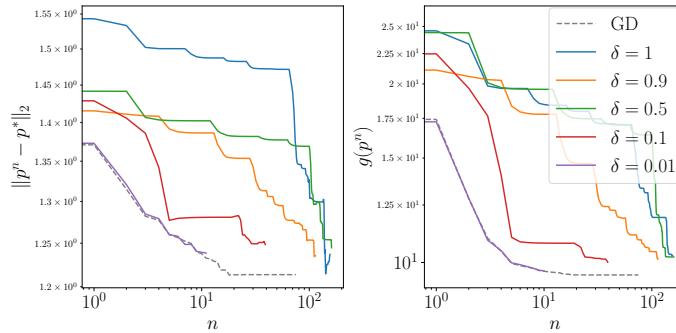
(a) $N = 2$; $\dim(P) = 4$ (b) $N = 4$; $\dim(P) = 16$ (c) $N = 8$; $\dim(P) = 64$

Fig. 5.2: Results for the groundwater flow example from subsection 5.2, for a variety of parameter dimensions. In each row, the left-hand plot shows the distance between the parameter found at iteration n and the true value p^* of the parameter. The right-hand plot shows the value of the objective function, the negative log-likelihood in the Bayesian inference problem.

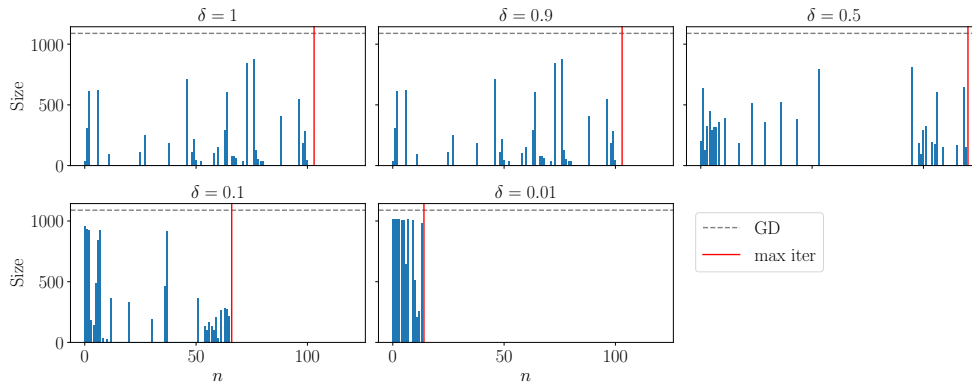


Fig. 5.3: Size of the matrix whose Cholesky factorisation must be computed for each iteration of gradient descent in the groundwater flow example from subsection 5.2, as the parameter δ is varied. The figure is for $N = 2$, $\dim(P) = 4$, but results for other parameter dimensions are similar. The dashed gray line shows the size of the (sparse) problem that must be solved for the classical approach, while the red line indicates the iteration number at which convergence was achieved.

910 on two challenging applied problems with favourable results compared to classical
 911 approaches. The chief advantages of the approach are that (i) gradients can be cal-
 912 culated at a lower cost than in classical approaches, (ii) that a global model for the
 913 gradient across parameter space is constructed, allowing for re-use of computational
 914 effort from previous iterations of gradient descent and (iii) that a full probability
 915 model is output, providing an error indicator that we used both to determine when
 916 to refine the approximation and to perform line searches.

917 Several possible avenues for future work present themselves. The first would be
 918 continuing to develop applications of this algorithm within optimisation, either by de-
 919 veloping versions of more sophisticated gradient-based optimisation algorithms which
 920 exploit probabilistic gradients, or by extending the framework to obtain higher order
 921 information to accelerate the optimisation. Another would be to explore the use of
 922 probabilistic gradients in other applications. In particular, we note that while com-
 923 puting the MAP point is an important problem in Bayesian inference, sophisticated
 924 Markov-chain Monte-Carlo algorithms for sampling the posterior also make use of this
 925 information, and the posterior distribution over the gradient presented herein could
 926 straightforwardly be incorporated into such algorithms.

927

REFERENCES

- 928 [1] M. A. ÁLVAREZ, L. ROSASCO, AND N. D. LAWRENCE, *Kernels for vector-valued functions:*
 929 *a review*, Foundations and Trends in Machine Learning, 4 (2012), pp. 195–266, <https://doi.org/10.1561/22000000036>.
 930
 931 [2] L. ARRIOLA AND J. M. HYMAN, *Sensitivity analysis for uncertainty quantification in math-*
 932 *ematical models*, in Mathematical and statistical estimation approaches in epidemiology,
 933 Springer, 2009, pp. 195–247.
 934 [3] A. V. BEDDOWS, N. KITWIROON, M. L. WILLIAMS, AND S. D. BEEVERS, *Emulation and*
 935 *sensitivity analysis of the community multiscale air quality model for a UK ozone pol-*
 936 *lution episode*, Environmental Science & Technology, 51 (2017), pp. 6229–6236, <https://doi.org/10.1021/acs.est.6b05873>.
 937

- 938 [4] P. BENNER, S. GUGERCIN, AND K. WILLCOX, *A survey of projection-based model reduction*
939 *methods for parametric dynamical systems*, SIAM review, 57 (2015), pp. 483–531.
- 940 [5] P. BENNER, E. SACHS, AND S. VOLKWEIN, *Model order reduction for PDE constrained opti-*
941 *mization*, in Trends in PDE constrained optimization, Springer, 2014, pp. 303–326.
- 942 [6] A. BERLINET AND C. THOMAS-AGNAN, *Reproducing Kernel Hilbert Spaces in Probability and*
943 *Statistics*, Springer US, 2004, <https://doi.org/10.1007/978-1-4419-9096-9>.
- 944 [7] L. T. BIEGLER, O. GHATTAS, M. HEINKENSCHLOSS, D. KEYES, AND B. VAN BLOEMEN WAAN-
945 *ders*, *Real-time PDE-constrained Optimization*, SIAM, 2007.
- 946 [8] V. I. BOGACHEV, *Gaussian Measures*, vol. 62, American Mathematical Society Providence,
947 1998.
- 948 [9] J. F. BONNANS AND A. SHAPIRO, *Perturbation analysis of optimization problems*, Springer
949 Science & Business Media, 2013.
- 950 [10] K. CHENG, Z. LU, C. LING, AND S. ZHOU, *Surrogate-assisted global sensitivity analysis: an*
951 *overview*, Structural and Multidisciplinary Optimization, 61 (2020), pp. 1187–1213, <https://doi.org/10.1007/s00158-019-02413-5>.
- 952 [11] I. CIALENCO, G. E. FASSHAUER, AND Q. YE, *Approximation of stochastic partial differential*
953 *equations by a kernel-based collocation method*, Int. J. Comput. Math., 89 (2012), pp. 2543–
954 2561, <https://doi.org/10.1080/00207160.2012.688111>.
- 955 [12] E. CLEARY, A. GARBUNO-INIGO, S. LAN, T. SCHNEIDER, AND A. M. STUART, *Calibrate, Emu-*
956 *late, Sample*, Journal of Computational Physics, (2020), p. 109716, <https://doi.org/https://doi.org/10.1016/j.jcp.2020.109716>.
- 957 [13] J. COCKAYNE, *Bayesian Probabilistic Numerical Methods*, PhD thesis, University of Warwick,
958 2019.
- 959 [14] J. COCKAYNE, C. J. OATES, T. J. SULLIVAN, AND M. GIROLAMI, *Bayesian probabilistic numeri-*
960 *cal methods*, SIAM Review, 61 (2019), pp. 756–789, <https://doi.org/10.1137/17m1139357>.
- 961 [15] H. B. CURRY, *The method of steepest descent for non-linear minimization problems*, Q APPL
962 MATH, 2 (1944), pp. 258–261, <https://doi.org/10.1090/qam/10667>.
- 963 [16] M. DROHMANN AND K. CARLBERG, *The ROMES method for statistical modeling of reduced-*
964 *order-model error*, SIAM/ASA Journal on Uncertainty Quantification, 3 (2015), pp. 116–
965 145.
- 966 [17] L. EVANS, *Partial Differential Equations*, American Mathematical Society, Mar. 2010, <https://doi.org/10.1090/gsm/019>.
- 967 [18] M. FISHER, J. NOCEDAL, Y. TRÉMOLET, AND S. J. WRIGHT, *Data assimilation in weather*
968 *forecasting: a case study in PDE-constrained optimization*, Optimization and Engineering,
969 10 (2009), pp. 409–426.
- 970 [19] R. FITZHUGH, *Impulses and physiological states in theoretical models of nerve membrane*, Bio-
971 *physical Journal*, 1 (1961), pp. 445–466, [https://doi.org/10.1016/s0006-3495\(61\)86902-6](https://doi.org/10.1016/s0006-3495(61)86902-6).
- 972 [20] C. GEYER, *Introduction to markov chain monte carlo*, Handbook of markov chain monte carlo,
973 20116022 (2011), p. 45.
- 974 [21] S. GIRARD, V. MALLET, I. KORSAKISSOK, AND A. MATHIEU, *Emulation and Sobol' sensitivity*
975 *analysis of an atmospheric dispersion model applied to the Fukushima nuclear accident*,
976 *Journal of Geophysical Research: Atmospheres*, 121 (2016), pp. 3484–3496, <https://doi.org/10.1002/2015jd023993>.
- 977 [22] M. GUNZBURGER, *Perspective in flow control and optimization (2003)*, SIAM, Philadelphia,
978 2002.
- 979 [23] E. HAIRER, S. NØRSETT, AND G. WANNER, *Solving Ordinary Differential Equations I: Nonstiff*
980 *Problems*, Springer, 1993.
- 981 [24] D. HARTMAN AND L. K. MESTHA, *A deep learning framework for model reduction of dynamical*
982 *systems*, in 2017 IEEE Conference on Control Technology and Applications (CCTA), IEEE,
983 2017, pp. 1917–1922.
- 984 [25] P. HENNIG, M. A. OSBORNE, AND M. GIROLAMI, *Probabilistic numerics and uncertainty in*
985 *computations*, J. R. Stat. Soc. A Stat., 471 (2015), pp. 20150142, 17, <https://doi.org/10.1098/rspa.2015.0142>.
- 986 [26] R. HERZOG AND K. KUNISCH, *Algorithms for PDE-constrained optimization*, GAMM-
987 *Mitteilungen*, 33 (2010), pp. 163–176, <https://doi.org/10.1002/gamm.201010013>.
- 988 [27] D. HIGDON, M. KENNEDY, J. C. CAVENDISH, J. A. CAPEO, AND R. D. RYNE, *Combining field*
989 *data and computer simulations for calibration and prediction*, SIAM Journal on Scientific
990 *Computing*, 26 (2004), pp. 448–466.
- 991 [28] K. ITO AND K. KUNISCH, *Lagrange multiplier approach to variational problems and applications*,
992 SIAM, 2008.
- 993 [29] R. JIN, W. CHEN, AND A. SUDJANTO, *Analytical metamodel-based global sensitivity analy-*
994 *sis and uncertainty propagation for robust design*, in SAE Technical Paper Series, SAE
995

- International, Mar. 2004, <https://doi.org/10.4271/2004-01-0429>.
- 1001 [30] M. C. KENNEDY AND A. O'HAGAN, *Bayesian calibration of computer models*, Journal of the
1002 Royal Statistical Society: Series B (Statistical Methodology), 63 (2001), pp. 425–464, <https://doi.org/10.1111/1467-9868.00294>.
- 1003 [31] S. LAN, T. BUI-THANH, M. CHRISTIE, AND M. GIROLAMI, *Emulation of higher-order tensors in*
1004 *manifold Monte Carlo methods for Bayesian inverse problems*, Journal of Computational
1005 Physics, 308 (2016), pp. 81–101.
- 1006 [32] M. MAHSERECI AND P. HENNIG, *Probabilistic line searches for stochastic optimization*, in Ad-
1007 vances in Neural Information Processing Systems 28, C. Cortes, N. D. Lawrence, D. D.
1008 Lee, M. Sugiyama, and R. Garnett, eds., Curran Associates, Inc., 2015, pp. 181–189, <http://papers.nips.cc/paper/5753-probabilistic-line-searches-for-stochastic-optimization.pdf>.
- 1009 [33] J. C. NEWMAN III, A. C. TAYLOR III, R. W. BARNWELL, P. A. NEWMAN, AND G. J.-W.
1010 HOU, *Overview of sensitivity analysis and shape optimization for complex aerodynamic*
1011 *configurations*, Journal of Aircraft, 36 (1999), pp. 87–96.
- 1012 [34] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer New York, 2006, <https://doi.org/10.1007/978-0-387-40065-5>, <https://doi.org/10.1007/978-0-387-40065-5>.
- 1013 [35] J. OAKLEY AND A. O'HAGAN, *Bayesian inference for the uncertainty distribution of computer*
1014 *model outputs*, Biometrika, 89 (2002), pp. 769–784.
- 1015 [36] J. E. OAKLEY AND A. O'HAGAN, *Probabilistic sensitivity analysis of complex models: a*
1016 *Bayesian approach*, Journal of the Royal Statistical Society: Series B (Statistical Method-
1017 ology), 66 (2004), pp. 751–769, <https://doi.org/10.1111/j.1467-9868.2004.05304.x>.
- 1018 [37] C. J. OATES AND T. J. SULLIVAN, *A modern retrospective on probabilistic numerics*, Statistics
1019 and Computing, 29 (2019), pp. 1335–1351, <https://doi.org/10.1007/s11222-019-09902-z>.
- 1020 [38] M. OSBORNE, *Bayesian Gaussian Processes for Sequential Prediction, Optimisation and Quad-*
1021 *rateure*, PhD thesis, PhD thesis, University of Oxford, 2010.
- 1022 [39] H. OWHADI, *Bayesian numerical homogenization*, Multiscale Modeling & Simulation, 13 (2015),
1023 pp. 812–828, <https://doi.org/10.1137/140974596>.
- 1024 [40] H. OWHADI, *Multigrid with rough coefficients and multiresolution operator decomposition from*
1025 *hierarchical information games*, SIAM Review, 59 (2017), pp. 99–149, <https://doi.org/10.1137/15m1013894>.
- 1026 [41] H. OWHADI AND L. ZHANG, *Gamblets for opening the complexity-bottleneck of implicit schemes*
1027 *for hyperbolic and parabolic ODEs/PDEs with rough coefficients*, Journal of Computational
1028 Physics, 347 (2017), pp. 99–128, <https://doi.org/10.1016/j.jcp.2017.06.037>.
- 1029 [42] R. PULCH, E. J. W. TER MATEN, AND F. AUGUSTIN, *Sensitivity analysis and model order re-*
1030 *duction for random linear dynamical systems*, Mathematics and Computers in Simulation,
1031 111 (2015), pp. 80–95.
- 1032 [43] C. E. RASMUSSEN AND C. K. I. WILLIAMS, *Gaussian Processes for Machine Learning*, The
1033 MIT Press, 2005, <https://doi.org/10.7551/mitpress/3206.001.0001>.
- 1034 [44] M. RENARDY, T.-M. YI, D. XIU, AND C.-S. CHOU, *Parameter uncertainty quantification using*
1035 *surrogate models applied to a spatial model of yeast mating polarization*, PLOS Computa-
1036 tional Biology, 14 (2018), p. e1006181, <https://doi.org/10.1371/journal.pcbi.1006181>.
- 1037 [45] J. SACKS, W. J. WELCH, T. J. MITCHELL, AND H. P. WYNN, *Design and analysis of computer*
1038 *experiments*, Statistical science, (1989), pp. 409–423.
- 1039 [46] O. SAN AND R. MAULIK, *Extreme learning machine for reduced order modeling of turbulent*
1040 *geophysical flows*, Physical Review E, 97 (2018), p. 042322.
- 1041 [47] O. SAN AND R. MAULIK, *Neural network closures for nonlinear model order reduction*, Advances
1042 in Computational Mathematics, 44 (2018), pp. 1717–1750.
- 1043 [48] B. SENGUPTA, K. J. FRISTON, AND W. D. PENNY, *Efficient gradient computation for dynamical*
1044 *models*, NeuroImage, 98 (2014), pp. 521–527.
- 1045 [49] M. SHAFTO, M. CONROY, R. DOYLE, E. GLAESSGEN, C. KEMP, J. LEMOIGNE, AND L. WANG,
1046 *Modeling, simulation, information technology & processing roadmap*, National Aeronautics
1047 and Space Administration, (2012).
- 1048 [50] S. SHERIFFDEEN, J. C. RAGUSA, J. E. MOREL, M. L. ADAMS, AND T. BUI-THANH, *Accelerating*
1049 *PDE-constrained inverse solutions with Deep Learning and Reduced Order Models*, arXiv
1050 preprint arXiv:1912.08864, (2019).
- 1051 [51] I. SOBOL', *Global sensitivity indices for nonlinear mathematical models and their Monte Carlo*
1052 *estimates*, Mathematics and Computers in Simulation, 55 (2001), pp. 271–280, [https://doi.org/10.1016/s0378-4754\(00\)00270-6](https://doi.org/10.1016/s0378-4754(00)00270-6).
- 1053 [52] A. STUART AND A. TECKENTRUP, *Posterior consistency for Gaussian process approximations*
1054 *of Bayesian posterior distributions*, Mathematics of Computation, 87 (2018), pp. 721–753.
- 1055 [53] A. M. STUART, *Inverse problems: A Bayesian perspective*, Acta Numerica, 19 (2010), pp. 451–
1056 559, <https://doi.org/10.1017/s0962492910000061>.

- 1062 [54] C. VILLANI, *Optimal Transport*, Springer Berlin Heidelberg, 2009, <https://doi.org/10.1007/>
1063 [978-3-540-71050-9](https://doi.org/10.1007/978-3-540-71050-9).
- 1064 [55] H. WENDLAND, *Scattered Data Approximation*, Cambridge University Press, Dec. 2004, [https:](https://doi.org/10.1017/cbo9780511617539)
1065 [//doi.org/10.1017/cbo9780511617539](https://doi.org/10.1017/cbo9780511617539).