

# Uncertainty Modeling and Propagation for Groundwater Flow: A Comparative Study of Surrogates

Dedicated to the memory of K. Andrew Cliffe (1953–2014)

Oliver G. Ernst<sup>1\*</sup>, Björn Sprungk<sup>2†</sup> and Chao Zhang<sup>1†</sup>

<sup>1\*</sup>Department of Mathematics, TU Chemnitz, Germany .

<sup>2\*</sup>Faculty of Mathematics and Computer Science, TU  
Bergakademie Freiberg, Freiberg, Germany .

\*Corresponding author(s). E-mail(s):

[uernst@math.tu-chemnitz.de](mailto:uernst@math.tu-chemnitz.de);

Contributing authors: [bjoern.sprungk@math.tu-freiberg.de](mailto:bjoern.sprungk@math.tu-freiberg.de);

[chao.zhang@math.tu-chemnitz.de](mailto:chao.zhang@math.tu-chemnitz.de);

†These authors contributed equally to this work.

## Abstract

We compare sparse grid stochastic collocation and Gaussian process emulation as surrogates for the parameter-to-observation map of a groundwater flow problem related to the Waste Isolation Pilot Plant in Carlsbad, NM. The goal is the computation of the probability distribution of a contaminant particle travel time resulting from uncertain knowledge about the transmissivity field. The latter is modelled as a lognormal random field which is fitted by restricted maximum likelihood estimation and universal kriging to observational data as well as geological information including site-specific trend regression functions obtained from technical documentation. The resulting random transmissivity field leads to a random groundwater flow and particle transport problem which is solved realization-wise using a mixed finite element discretization. Computational surrogates, once constructed, allow sampling the quantities of interest in the uncertainty analysis at substantially reduced computational cost. Special emphasis is placed on explaining the differences between the

001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046

047 two surrogates in terms of computational realization and interpreta-  
048 tion of the results. Numerical experiments are given for illustration.

049 **Keywords:** sparse grid stochastic collocation, Gaussian process emulation,  
050 uncertainty propagation, kriging, Darcy flow, mixed finite elements

051 **MSC Classification:** 60G60 , 60H35 , 62P12 , 62M30 , 65C05 , 65D12 , 65C30 ,  
052 65N75 ,

054

055

056

## 057 1 Introduction

058

059 By their very nature, the Earth Sciences have had to cope with uncertainty  
060 from early on, and scientists from this field such as Harold Jeffreys and Albert  
061 Tarantola have had foundational and lasting impact on how uncertainty is  
062 modeled and merged with physical models in the interdisciplinary field now  
063 known as *uncertainty quantification (UQ)*. A current account of uncertainty  
064 quantification in subsurface hydrology can be found in [Linde et al \(2017\)](#). Many  
065 UQ studies involve a system governed by a partial differential equation (PDE)  
066 in which one or more input quantities are uncertain. When this uncertainty  
067 is described in probabilistic terms, we arrive at a PDE with random data, or  
068 *random PDE* for short, the solutions of which are stochastic processes, also  
069 referred to in this context as random functions or *random fields*. The task of  
070 determining the probability distribution of the solution of a random PDE, or  
071 of *quantities of interest* derived from such solutions, is known as *uncertainty*  
072 *propagation* or *forward UQ* (cf. [Ernst et al \(2022\)](#)). Approximation methods  
073 for random fields and their incorporation into PDE solution methods have  
074 been actively developed in the engineering and numerical analysis communities  
075 in the past two decades, and excellent surveys can be found in [Ghanem and](#)  
076 [Spanos \(1991\)](#); [Babuška et al \(2010\)](#); [Schwab and Gittelsohn \(2011\)](#); [Gunzburger](#)  
077 [et al \(2014\)](#). At the same time, sampling-based simulation techniques known as  
078 *Gaussian process emulators* have gained popularity in the statistics community  
079 for solving similar problems, cf. [Sacks et al \(1989\)](#); [Currin et al \(1991\)](#); [Kennedy](#)  
080 [and O’Hagan \(2001\)](#); [O’Hagan \(2006\)](#). Our objective in this work is the direct  
081 comparison of these two approaches using Monte Carlo sampling as a reference  
082 in a case study on the hydrogeological transport of radionuclides within the  
083 site assessment for a nuclear waste repository. In doing so, we place particular  
084 emphasis on the careful construction of a stochastic model using geostatistical  
085 techniques.

086

087

088

089

090

091

092

The uncertainty propagation techniques we shall consider are based on gener-  
ating realizations (samples) of the uncertain input parameters, solving the  
PDE for each realization and then determining the statistical properties of the  
quantities of interest in a post-processing step. As each PDE solution typically  
requires considerable computational resources, the mapping of random input  
parameters to quantities of interest is often substituted by *surrogate models*,  
which are considerably less costly to evaluate, thus speeding up the uncertainty

propagation analysis. The two surrogates we shall compare, *sparse polynomial collocation* and Gaussian process emulation are interesting in that they were developed in different fields (numerical analysis and statistics), display different performance characteristics, and also differ in the interpretations of the surrogates they produce. Our work is closest in spirit to Owen et al (2017), where Gaussian process emulation is compared with polynomial chaos expansion surrogates for two black-box computer simulators. Although different in construction, polynomial chaos surrogates yield a multivariate polynomial approximation of the input-output map realized by the computer simulator as does stochastic collocation, whereas the latter is considerably easier to integrate into PDE solvers. In place of a small number of discrete parameters in the models considered in (Owen et al, 2017), the random input in our groundwater model is a random field, i.e., its realizations are functions, which can be considered as parameterized by a countably infinite number of parameters. A comparison between uncertainty propagation techniques in aerodynamic modeling can be found in Liu et al (2017).

The remainder of the paper is organized as follows: Section 2 presents the problem of predicting the travel or *exit time* of radionuclides transported by groundwater flow through a horizontal layer above the Waste Isolation Pilot Plant, an operational underground disposal site for nuclear waste, in a scenario where a hypothetical future accidental breach leads to the release of radioactive material. The physical as well as the probabilistic model are presented as well as how observational data of hydraulic transmissivity is incorporated, leading to the generation of samples of the exit time quantity of interest. Section 3 describes the computational realization for solving the Darcy flow equations, the construction of the truncated Karhunen-Loève representation of the random transmissivity field as well as the estimation of the cumulative distribution function of the exit time quantity of interest. Section 4 gives detailed description of the two surrogate types to be compared, Gaussian process emulation and sparse polynomial collocation, emphasising their differences with respect to construction, computation and interpretation. In Section 5, we present the results of numerical computations with both surrogates using original data from the WIPP site, and present our conclusions in Section 6.

## 2 Uncertainty Propagation for a Groundwater Flow Problem

In this section we introduce the application setting, physical model, UQ task as well as the probabilistic model with which this is addressed.

### 2.1 The Waste Isolation Pilot Plant (WIPP)

The Waste Isolation Pilot Plant (WIPP) in Carlsbad, NM is a long-term deep geologic storage facility for transuranic waste operated by the U.S. Department of Energy since 1999. One of the issues investigated in the course of an

139 extensive performance assessment for WIPP was the risk of hazardous mate-  
 140 rials escaping to the biosphere in the event of a future accidental breach of the  
 141 enclosure system. As the most likely pathway for such contaminants is trans-  
 142 port through the subsurface via groundwater, we are led to the objective of  
 143 predicting the groundwater flow and transport of contaminants released from  
 144 the storage site. In case of the WIPP, the disposal area lies within in the *Salado*  
 145 bedded salt formation. The Salado itself as well as the overlying formations are  
 146 essentially impermeable to groundwater with the exception of a laterally exten-  
 147 sive but narrow layer of rock known as the *Culebra Dolomite*. Details of the  
 148 geological site characterization can be found in the extensive documentation<sup>1</sup>  
 149 in the WIPP certification and recertification applications (U.S. Department  
 150 of Energy (DOE), 2004, 2014) which are produced every five years. Figure 1,  
 151 taken from (U.S. Department of Energy (DOE), 2014), shows the location of  
 152 the WIPP site within the UTM coordinate system, the location of boreholes  
 153 where measurements of transmissivity and hydraulic head were obtained as  
 154 well as the boundaries of areas with distinct geological features.

155 One of the most relevant quantities of interest is the travel or *exit time*  
 156 of radionuclides after release from a point within the Culebra layer above the  
 157 site to reach the boundary of the repository area, the computation of which  
 158 requires simulating the groundwater flow and transport in the Culebra. As the  
 159 precise transmissivity properties of the rock are uncertain, the same applies to  
 160 the exit time. In the remainder of this section we describe a model for ground-  
 161 water flow and contaminant transport in which the uncertain transmissivity  
 162 is modeled stochastically, incorporating geological background information,  
 163 standard geostatistical assumptions as well as available measurement data.

164

## 165 **2.2 Darcy Flow and Particle Transport**

166

167 We model the flow of groundwater through the Culebra dolomite geological  
 168 unit by stationary single-phase Darcy flow. Denoting by  $p$  the *hydraulic head*  
 169 (pressure) and by  $K$  the (scalar) *hydraulic conductivity*, the *volumetric flux*  
 170 (Darcy flux)  $\mathbf{q}$  is given by

$$171 \mathbf{q} = -K\nabla p. \quad (1)$$

172

173 If  $\mathbf{u}$  denotes the pore velocity of the groundwater, which is related to the Darcy  
 174 flux in terms of the *porosity*  $\phi$  as  $\mathbf{q} = \phi\mathbf{u}$ , conservation of mass in the absence  
 of sources and sinks leads to the divergence-free condition

$$175 \nabla \cdot \mathbf{u} = 0. \quad (2)$$

176

177 Since the aquifer under consideration is essentially horizontal with a much  
 178 larger lateral than vertical extent, we model the flow as two-dimensional and  
 179 consider the hydraulic *transmissivity*  $T = bK$  in place of conductivity, where  
 180  $b$  denotes the aquifer thickness.

181

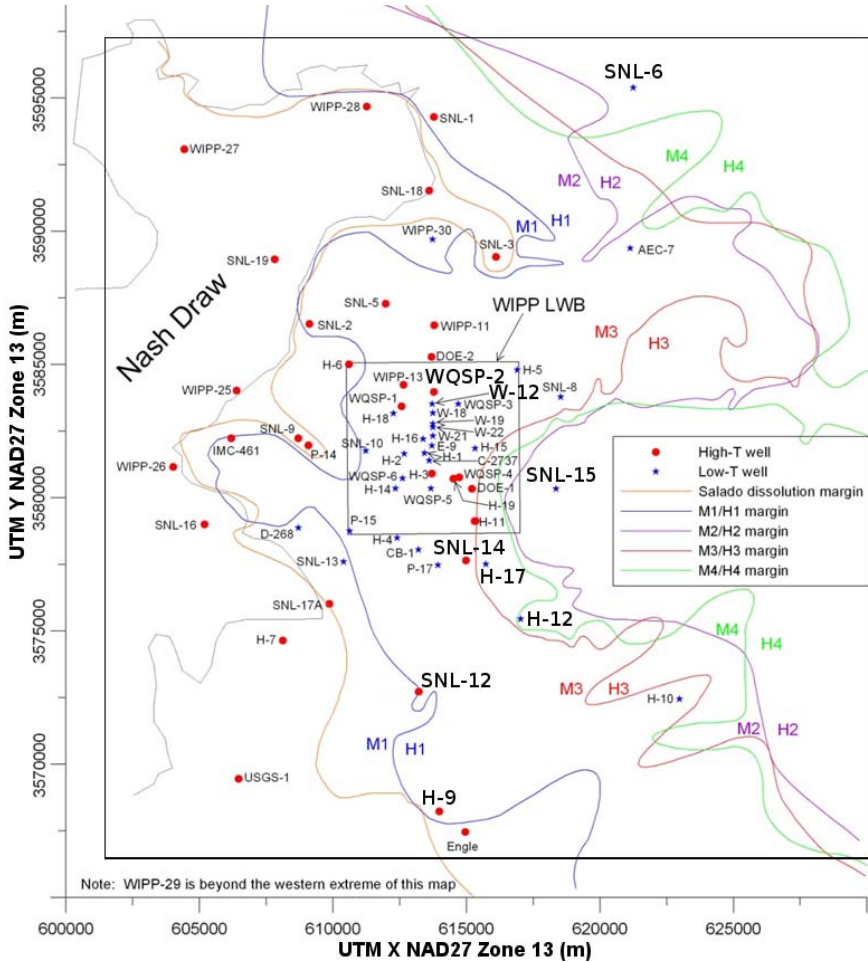
182 On the boundary  $\partial D$  of the bounded computational domain  $D$ , we dis-  
 183 tinguish impermeable segments  $\Gamma_N$  along which the normal flux vanishes and

184

---

<sup>1</sup>These can be found at <https://wipp.energy.gov/epa-certification-documents.asp>.





**Fig. 1** Horizontal location of WIPP repository (dashed lines), observation boreholes as well as boundaries of distinct geological features. Source: (U.S. Department of Energy (DOE), 2014).

their complement  $\Gamma_D = \partial D \setminus \Gamma_N$ , where we prescribe the value of the hydraulic head  $p$ . Denoting by  $\mathbf{n}$  the exterior unit normal vector along  $\Gamma_N$  and by  $g$  the prescribed head data along  $\Gamma_D$ , this leads to the boundary conditions

$$\mathbf{n} \cdot \mathbf{u} = 0 \quad \text{on } \Gamma_N, \quad p = g \quad \text{on } \Gamma_D. \quad (3)$$

The computational domain  $D$  as well as the boundary segments  $\Gamma_N$  and  $\Gamma_D$  are displayed in the left panel in Figure 2. The Dirichlet data  $g$  is obtained by evaluating a kriging interpolant (cf. Section 2.4.3) of observational hydraulic head data taken from (U.S. Department of Energy (DOE), 2014). As the flux

variable  $\mathbf{u}$  is of primary interest in view of the subsequent transport calculation we employ the usual mixed formulation of the boundary value problem presented by (1), (2) and (3). The associated variational formulation consists in finding the pair  $(\mathbf{u}, p) \in \mathcal{V} \times \mathcal{W}$  such that

$$\left( \frac{\phi b}{T} \mathbf{u}, \mathbf{v} \right) - (p, \nabla \cdot \mathbf{v}) = -\langle g, \mathbf{n} \cdot \mathbf{v} \rangle_{\Gamma_D} \quad \forall \mathbf{v} \in \mathcal{V}, \quad (4a)$$

$$(\nabla \cdot \mathbf{u}, q) = 0 \quad \forall q \in \mathcal{W} \quad (4b)$$

with suitable boundary data  $g \in H^{1/2}(\Gamma_D)$ . Here  $(\cdot, \cdot)$  denotes the  $L^2(D)$  inner product, the variational spaces are given by

$$\mathcal{V} = \{ \mathbf{v} \in \mathbf{H}(\text{div}; D), \mathbf{n} \cdot \mathbf{v}|_{\Gamma_N} = 0 \}, \quad \mathcal{W} = L^2(D)$$

and  $\langle \cdot, \cdot \rangle_{\Gamma_D}$  denotes the duality pairing of  $H^{1/2}(\Gamma_D) \times H^{-1/2}(\Gamma_D)$ . Given the flux solution  $\mathbf{u}$  of (4), the trajectory of a particle from a release point  $\mathbf{x}_0 \in D$  neglecting hydraulic dispersion is found as the solution of the initial value problem

$$\dot{\mathbf{x}}(t) = \mathbf{u}(\mathbf{x}(t)), \quad t \geq 0, \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (5)$$

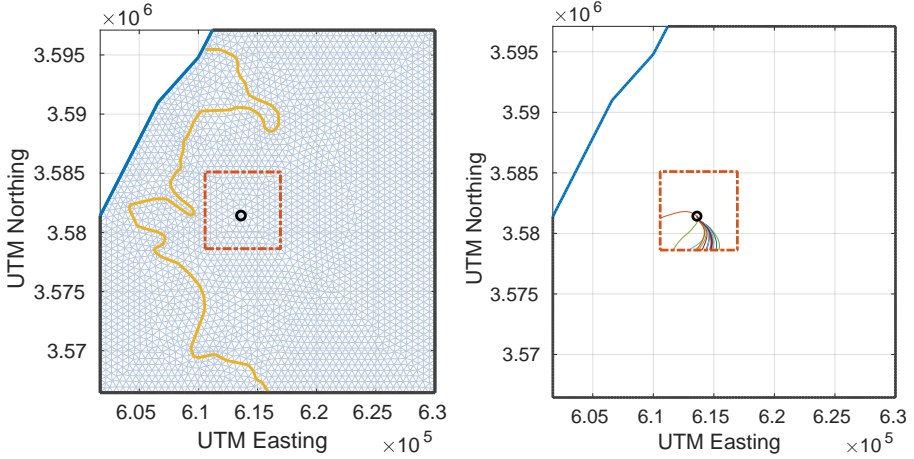
A discussion of the regularity requirements for the Darcy flow problem (4) needed to ensure existence and uniqueness of the particle trajectory (5) can be found in (Graham et al, 2016, Section 5.3). As we shall see below, for the probabilistic model of transmissivity with finite-dimensional noise, which we shall employ in our calculations, these requirements are satisfied. As a *quantity of interest* derived from the solution of the random Darcy flow equations, we choose the logarithm of the travel or exit time of a particle released at a location  $\mathbf{x}_0$  inside the Culebra layer above the WIPP repository until it reaches the boundary of the subdomain  $D_0 \subset D$  marking the edge of the WIPP site projected vertically up to the Culebra layer within the surrounding computational domain  $D$ ,

$$f_{\text{exit}} := \log \min\{t > 0 : \mathbf{x}(t) \notin D_0, \mathbf{x}_0 \in D_0\}.$$

The location of the release point  $\mathbf{x}_0$ , the perimeter of the WIPP site  $D_0$  as well as a number of particle trajectory realizations from  $\mathbf{x}_0$  to  $\partial D_0$  are displayed in Figure 2.

### 2.3 Probabilistic Modeling of Uncertain Transmissivity

The primary source of uncertainty in the modeling of flow and transport in the Culebra dolomite is the spatial variation of hydraulic conductivity, or, in our horizontal two-dimensional setting, transmissivity  $T$ . The prevailing mathematical description of uncertainty is probabilistic, i.e., the quantities in question are modeled as random variables following a given probability distribution. The randomness thus introduced is an expression of uncertainty due



**Fig. 2** Left: Computational domain  $D$  with Neumann boundary  $\Gamma_N$  (blue) and Dirichlet boundary  $\Gamma_D$  (black) as well as the perimeter of the WIPP site  $D_0$  (red dashed), location of particle release point  $\mathbf{x}_0$  (black circle), and boundary of the Salado dissolution zone  $D_1$  (yellow), cf. Section 2.4.1 below, respected by the triangular finite element mesh. Right: Simulation of several realizations of random particle trajectories from  $\mathbf{x}_0$  to  $\partial D_0$ .

to lack of knowledge of the precise spatial variation of transmissivity throughout the domain  $D$  in the sense that some realizations of transmissivity across the domain are more likely than others. Rather than a deterministic value  $T = T(\mathbf{x})$ , transmissivity at a point  $\mathbf{x} \in D$  (scaled by porosity and thickness) is thus expressed as a random variable  $T(\mathbf{x}, \omega)$  governed by a probability measure  $\mathbf{P}$  defined on a probability space  $(\Omega, \mathfrak{A}, \mathbf{P})$  with elementary outcome set  $\Omega$  carrying a  $\sigma$ -algebra  $\mathfrak{A}$  on which a probability measure  $\mathbf{P}$  is defined. The collection of all such random variables  $\{T(\mathbf{x}, \omega) : \mathbf{x} \in D\}$  is known as a *random field*, i.e., a stochastic process for which the index variable  $\mathbf{x}$  is a spatial coordinate. The most well-established probabilistic model for transmissivity in the hydrology literature assumes that  $T(\mathbf{x}, \cdot)$  follows a *lognormal* distribution, i.e., that  $Z(\mathbf{x}, \cdot) := \log T(\mathbf{x}, \cdot)$  is a Gaussian random field (cf. Freeze (1975); Hoeksema and Kitanidis (1985) and (de Marsily, 1986, Chapter 11)). By consequence, realizations of  $T = \exp(Z)$  are always positive. Such a Gaussian random field  $Z$  is completely specified by its mean and covariance function

$$\begin{aligned} \bar{Z}(\mathbf{x}) &= \mathbf{E}[Z(\mathbf{x}, \cdot)], & \mathbf{x} \in D, \\ \text{and } c(\mathbf{x}, \mathbf{y}) &= \mathbf{E}[(Z(\mathbf{x}) - \bar{Z}(\mathbf{x}))(Z(\mathbf{y}) - \bar{Z}(\mathbf{y}))], & \mathbf{x}, \mathbf{y} \in D, \end{aligned}$$

respectively, where  $\mathbf{E}[\cdot]$  denotes mathematical expectation with respect to  $\mathbf{P}$ .

We assume throughout that the covariance function of  $Z = \log T$  is *isotropic* and that the fluctuation  $Z - \bar{Z}$  is *wide-sense stationary* such that we have  $c(\mathbf{x}, \mathbf{y}) = c(|\mathbf{x} - \mathbf{y}|)$ , i.e., the covariance depends only on the (Euclidean) separation distance  $r = |\mathbf{x} - \mathbf{y}|$ . Moreover, we assume  $c(r)$  to belong to the

277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322

323 *Matérn* family of covariance models

324

$$325 \quad c(r) = \frac{\sigma^2}{2^{\nu-1} \Gamma(\nu)} \left( \frac{2\sqrt{\nu} r}{\rho} \right)^\nu K_\nu \left( \frac{2\sqrt{\nu} r}{\rho} \right), \quad r = |\mathbf{x} - \mathbf{y}|, \quad (6)$$

326

327 where  $K_\nu$  denotes the modified Bessel function of order  $\nu > 0$ . The quantity  
 328  $\nu$  is called the *smoothness parameter*,  $\sigma^2 = c(0) = \mathbf{Var} Z(\mathbf{x})$  is the (marginal)  
 329 *variance* (constant in  $\mathbf{x}$ ) and the parameter  $\rho > 0$  is called the *correlation*  
 330 *length*, a measure of how quickly the covariance decays with separation dis-  
 331 tance. An extensive justification for using the Matérn model as well as its  
 332 properties and scaling variants can be found in (Stein, 1999, pp. 48). For the  
 333 particular scaling (6), the Matérn covariance coincides with the exponential  
 334 covariance for  $\nu = \frac{1}{2}$ , the Bessel covariance for  $\nu = 1$  and the squared expo-  
 335 nential covariance in the limit  $\nu \rightarrow \infty$ . The smoothness of the realizations of  
 336  $Z$  increases with  $\nu$ , and the spatial scale of variation is described by  $\rho$ . We  
 337 determine the values of the *hyperparameters*  $(\sigma, \rho, \nu)$  by statistical estimation  
 338 based on data published in the WIPP Compliance Recertification Assessment  
 339 U.S. Department of Energy (DOE) (2014) documents, which contain measure-  
 340 ments of transmissivity in the Culebra dolomite at 62 boreholes throughout  
 341 the assessment site (cf. Figure 1). Figure 3 displays realizations of a Gaus-  
 342 sian random field describing  $Z = \log T$  throughout the computational domain  
 343  $D$  representing the Culebra flow domain. It can be seen that larger values  
 344 of  $\nu$  result in realizations that are smoother, and smaller values of  $\rho$  lead to  
 345 structures which decorrelate faster with separation distance.

346

## 347 2.4 Statistical Estimation of Transmissivity Field

348

349 As described in Section 2.3, we model the uncertain hydraulic transmissivity  
 350  $T$  as a lognormal random field on the bounded simulation domain  $D \subset \mathbb{R}^2$ ,  
 351 i.e., the random field

352

$$353 \quad Z := \log T = \bar{Z}(\mathbf{x}) + \tilde{Z}(\mathbf{x}, \omega) \quad (7)$$

354

355 is Gaussian with (deterministic) mean  $\bar{Z}$  and (centered) residual field  $\tilde{Z}$ . Due  
 356 to the complexity and irregular features of geological structures, it is crucial  
 357 to merge stochastic models with available measurement data in a transparent  
 358 fashion. Below we summarize the statistical techniques by which available data  
 359 is incorporated into the stochastic model of uncertain transmissivity.

360

361

362

363

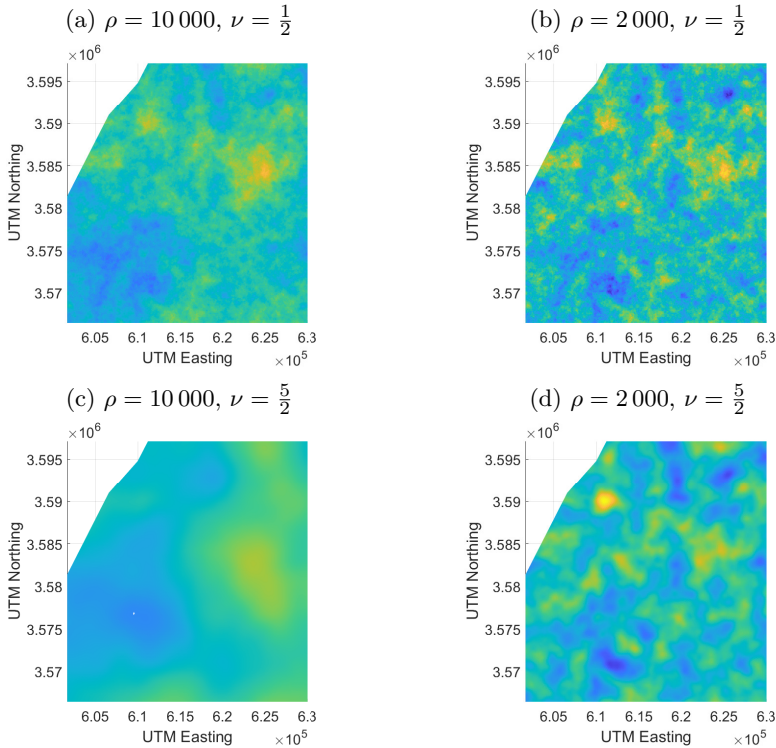
364

365

366

367

368



**Fig. 3** Realizations of mean-zero Gaussian random fields with Matérn covariance function for different values of  $\rho$  and  $\nu$ . All plots use the same color map and  $\sigma^2$  was set to 1 in each case.

### 2.4.1 Regression Modeling of Mean Transmissivity

The deterministic mean  $\bar{Z}$  of the log-transmissivity field is constructed as a linear regression model

$$\bar{Z}(\mathbf{x}) = \sum_{j=1}^k \beta_j h_j(\mathbf{x}) = \mathbf{h}(\mathbf{x})^\top \boldsymbol{\beta}, \quad \mathbf{h}(\mathbf{x}) = \begin{bmatrix} h_1(\mathbf{x}) \\ \vdots \\ h_k(\mathbf{x}) \end{bmatrix}, \quad (8)$$

in which the  $k$  components of  $\mathbf{h}$  contain *regression functions* from which the trend behavior of  $Z$  can be obtained by linear combination. Known geological features of the area under study can be incorporated by choosing the regression functions as, e.g., indicator functions of subdomains possessing distinguishing characteristics, linear or polynomial trends to be fitted as well as the variation of available quantities known or believed to affect the transmissivity field. Based on the available WIPP technical documents, model comparison was

369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414

415 made using the five regression functions

$$\begin{aligned}
 416 & \\
 417 & \quad h_1(\mathbf{x}) \equiv 1 \quad (\text{constant}), \quad h_4(\mathbf{x}) = d(\mathbf{x}) \quad (\text{overburden}), \\
 418 & \quad h_2(\mathbf{x}) = x_1 \quad (\text{linear in } x_1), \quad h_5(\mathbf{x}) = \mathbf{1}_{D_1}(\mathbf{x}) \quad (\text{zone indicator}). \quad (9) \\
 419 & \quad h_3(\mathbf{x}) = x_2 \quad (\text{linear in } x_2), \\
 420 &
 \end{aligned}$$

421 The first three regression functions allow to fit a basic affine trend. The *over-*  
 422 *burden*  $d(\mathbf{x})$  denotes the vertical distance between the ground surface and the  
 423 top of the Culebra layer above location  $\mathbf{x}$ . This is an indication of the extent  
 424 to which erosion has led to stress relief on the underlying Culebra layer, possi-  
 425 bly causing new fracturing or the opening of pre-existing fractures and thereby  
 426 enhancing transmissivity. Regression function  $h_5$  is the indicator function of a  
 427 subdomain  $D_1 \subset D$  to the north, south and west of the WIPP site, where dis-  
 428 solution of the upper Salado formation has led to strain in the rock overlying  
 429 the Salado, including the Culebra, leading to larger apertures in existing frac-  
 430 tures, collapse and brecciation and thus to a generally higher transmissivity  
 431 (cf. [U.S. Department of Energy \(DOE\) \(2004\)](#)).  
 432

### 433 2.4.2 Hyperparameter Estimation and Model Selection

434 For all combinations of the regression functions (9), a *restricted maximum*  
 435 *likelihood* (RML) estimation procedure detailed in [Appendix A](#) was used to  
 436 determine the hyperparameters  $\sigma, \rho$  and  $\nu$  of the Matérn covariance model  
 437 (6) based on the 62 transmissivity observations published in [U.S. Department](#)  
 438 [of Energy \(DOE\) \(2014\)](#). Based on this calibrated covariance structure, a  
 439 model comparison was carried out following a procedure proposed in [Kitanidis](#)  
 440 [\(1997b\)](#), in which a significance test is used to determine whether adding fur-  
 441 ther regression functions to a model better explains the data. In this way, we  
 442 arrived at a trend model (8) consisting of the regression functions  $\{h_1, h_2, h_5\}$   
 443 from (9). In the following we refer to this parametrization of the mean as the  
 444 *best model* and to that containing only the constant trend function  $h_1$  as the  
 445 *constant model*. The resulting estimates of the hyperparameters  $\sigma, \rho$  and  $\nu$  for  
 446 both models are given in [Table 1](#). Note that we have fixed  $\nu = 0.5$  in both cases  
 447 since the estimates for  $\nu$  were sufficiently close to this value<sup>2</sup>, which also allows  
 448 a more efficient evaluation of the associated covariance function. The regres-  
 449 sion model estimated by the (weighted) least-squares method for the mean is  
 450 then  
 451

$$452 \quad \bar{Z}(\mathbf{x}) = 143.98 - 2.55 \cdot 10^{-4} x_1 + 3.31 \mathbf{1}_{D_1}(\mathbf{x}).$$

453 Note that the values for  $x_1$  (UTM Easting coordinates) are of order  $6 \cdot 10^5$  for  
 454 the WIPP computational domain  $D$ .  
 455

---

456 <sup>2</sup>If we do not fix  $\nu = 0.5$  but estimate it as well the RML results are  $\hat{\sigma}^2 = 6.14$ ,  $\hat{\rho} = 2005.2$ ,  
 457 and  $\hat{\nu} = 0.48$ .  
 458  
 459  
 460



Trend model	Sill $\sigma^2$	Range $\rho$	Smoothness $\nu$
$h_1$	17.12	6509.8	0.5
$h_1, h_2, h_5$	6.15	1948.0	0.5

**Table 1** Restricted maximum likelihood estimation of hyperparameters  $\sigma^2$  (variance or *sill*) and  $\rho$  (correlation length or *range*) for two trend models based on the 64 observations of transmissivity. The smoothness parameter was fixed at  $\nu = 1/2$ , which corresponds to the exponential covariance kernel.

### 2.4.3 Conditioning on Transmissivity Data

Once the mean and covariance functions of the Gaussian random field  $Z = \log T$  have been determined, the log transmissivity measurements  $\{z(\mathbf{x}_j)\}_{j=1}^N$  may be used to further calibrate the stochastic model to fit the observations in a statistical sense using the technique known as *kriging* (cf. Cressie (1991); Kitanidis (1997a); Stein (1999)). Kriging refers to *best linear unbiased prediction* (BLUP) in which the value of the random field  $Z$  at an arbitrary location  $\mathbf{x} \in D$  is estimated as an affine combination

$$\hat{Z} = \hat{Z}(\mathbf{x}, \omega) = \lambda_0(\mathbf{x}) + \boldsymbol{\lambda}(\mathbf{x})^\top \mathbf{Z}(\omega) \quad (10)$$

of the (random) realizations  $\mathbf{Z} = (Z(\mathbf{x}_1), \dots, Z(\mathbf{x}_N))^\top$ , with spatially varying coefficients  $\lambda_0 : D \rightarrow \mathbb{R}$  and  $\boldsymbol{\lambda} = (\lambda_1(\mathbf{x}), \dots, \lambda_N(\mathbf{x})) : D \rightarrow \mathbb{R}^N$  chosen to make the estimator *unbiased* and *mean square optimal*, which requires that, for all  $\mathbf{x} \in D$ , we have

$$\mathbf{E} [\hat{Z}(\mathbf{x})] = \mathbf{E} [Z(\mathbf{x})] \quad \text{and} \quad \mathbf{E} \left[ |Z(\mathbf{x}) - \hat{Z}(\mathbf{x})|^2 \right] \rightarrow \min!_{\lambda_0, \boldsymbol{\lambda}}$$

For a known mean function  $\bar{Z}$  the solution is given by the (*simple*) *kriging prediction* or *interpolation*

$$\hat{Z}(\mathbf{x}) = \bar{Z}(\mathbf{x}) + \mathbf{c}(\mathbf{x})^\top \mathbf{C}^{-1} (\mathbf{Z} - \bar{\mathbf{Z}}),$$

where  $\bar{\mathbf{Z}} := [\bar{Z}(\mathbf{x}_1), \dots, \bar{Z}(\mathbf{x}_N)]^\top$ ,  $\mathbf{c}(\mathbf{x}) := (c(\mathbf{x}, \mathbf{x}_1), \dots, c(\mathbf{x}, \mathbf{x}_N))^\top$  and  $\mathbf{C} := (c(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1,\dots,N} \in \mathbb{R}^{N \times N}$ , with mean square error given via the *kriging (error) covariance*

$$\mathbf{E} \left[ |Z(\mathbf{x}) - \hat{Z}(\mathbf{x})|^2 \right] = \hat{c}(\mathbf{x}, \mathbf{x}), \quad \hat{c}(\mathbf{x}, \mathbf{y}) := c(\mathbf{x}, \mathbf{y}) - \mathbf{c}(\mathbf{x})^\top \mathbf{C}^{-1} \mathbf{c}(\mathbf{y}).$$

Note that for a Gaussian random field  $Z$  the kriging prediction  $\hat{Z}$  is again Gaussian and coincides with the conditioned random field  $Z(\mathbf{x}) | \mathbf{Z} = \mathbf{z}$ , where  $\mathbf{z} = (z_1, \dots, z_N)^\top$ , so that  $\hat{Z}(\mathbf{x}) \sim \mathbf{N}(\bar{Z}(\mathbf{x}) + \mathbf{c}(\mathbf{x})^\top \mathbf{C}^{-1} (\mathbf{z} - \bar{\mathbf{Z}}), \hat{c}(\mathbf{x}, \cdot))$ . It is easily verified that at the observation sites  $\{\mathbf{x}_j\}_{j=1}^N$  we have  $\hat{Z}(\mathbf{x}_j) = z(\mathbf{x}_j)$  and  $\hat{c}(\mathbf{x}_j, \mathbf{x}_j) = 0$ , hence the kriging estimate  $\hat{Z}$  of the random field  $Z$  interpolates the measurements.



507 In the variant called *universal kriging*, the mean  $\bar{Z}$  is not assumed known  
 508 and instead modelled as in (8). Forming the least squares estimate  $\hat{\beta}$  of  $\beta$  and  
 509 proceeding as above with  $\bar{Z}(\mathbf{x}) = \mathbf{h}(\mathbf{x})^\top \hat{\beta}$  would fail to account for uncer-  
 510 tainty in this estimate. Instead, we require that unbiasedness of the kriging  
 511 estimate (10) hold for all  $\beta \in \mathbb{R}^k$ , resp. for all possible mean functions. Apply-  
 512 ing unbiasedness as a constraint in the pointwise minimization over  $\lambda_0, \boldsymbol{\lambda}$  via  
 513 Lagrange multipliers yields the *universal kriging prediction* or interpolation

$$514 \hat{Z}(\mathbf{x}) = \begin{bmatrix} \mathbf{c}(\mathbf{x}) \\ \mathbf{h}(\mathbf{x}) \end{bmatrix}^\top \begin{bmatrix} \mathbf{C} & \mathbf{H} \\ \mathbf{H}^\top & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{Z} \\ 0 \end{bmatrix}, \quad (11)$$

518 where

$$519 \mathbf{H} = \begin{bmatrix} h_1(\mathbf{x}_1) & \dots & h_k(\mathbf{x}_1) \\ \vdots & & \vdots \\ h_1(\mathbf{x}_N) & \dots & h_k(\mathbf{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times k},$$

522 or, equivalently,

$$524 \hat{Z}(\mathbf{x}) = \mathbf{h}(\mathbf{x})^\top \hat{\beta} + \mathbf{c}(\mathbf{x})^\top \mathbf{C}^{-1} (\mathbf{Z} - \mathbf{H}\hat{\beta}), \quad (12)$$

527 where  $\hat{\beta} = (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}^{-1} \mathbf{Z}$ , with mean square error  
 528  $\mathbf{E} \left[ |\mathbf{Z}(\mathbf{x}) - \hat{Z}(\mathbf{x})|^2 \right] = \hat{c}(\mathbf{x}, \mathbf{x})$  given in this case by the *universal kriging*  
 529 (*error*) *covariance*

$$531 \hat{c}(\mathbf{x}, \mathbf{y}) := \mathbf{c}(\mathbf{x}, \mathbf{y}) - \mathbf{c}(\mathbf{x})^\top \mathbf{C}^{-1} \mathbf{c}(\mathbf{y}) + \boldsymbol{\gamma}(\mathbf{x})^\top \mathbf{V} \boldsymbol{\gamma}(\mathbf{y}), \quad (13)$$

533 where  $\boldsymbol{\gamma} = \mathbf{h}(\mathbf{x}) - \mathbf{H}^\top \mathbf{C}^{-1} \mathbf{c}(\mathbf{x})$  and  $\mathbf{V} = (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1}$ . Thus, the univer-  
 534 sal kriging prediction (12) consists in obtaining the mean as the least squares  
 535 estimate  $\mathbf{h}(\mathbf{x})^\top \hat{\beta}$  and proceeding as in simple kriging. However, the univer-  
 536 sal kriging mean square error contains the additional term  $\boldsymbol{\gamma}(\mathbf{x})^\top \mathbf{V} \boldsymbol{\gamma}(\mathbf{x}) \geq 0$   
 537 compared to that of simple kriging, which accounts for the additional uncer-  
 538 tainty present in the estimated mean and  $\beta$ , respectively. Note further that,  
 539 even for Gaussian  $Z$ , the universal kriging mean and (co)variance do not, in  
 540 general, possess an interpretation as those of a conditioned Gaussian random  
 541 field as is the case with simple kriging.

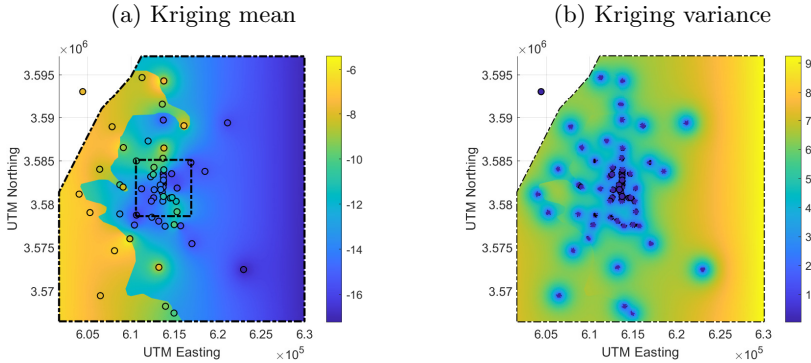
542 We now use the *universal kriged Gaussian random field*  $\hat{Z}$  obtained from  
 543 the available log transmissivity measurements  $\mathbf{z} = \{z(\mathbf{x}_j)\}_{j=1}^N$  as our final  
 544 stochastic model for the uncertain transmissivity field, i.e.,

$$546 \hat{Z}(\mathbf{x}) \sim \mathbf{N}(\hat{z}(\mathbf{x}), \hat{c}(\mathbf{x}, \cdot))$$

548 with  $\hat{c}$  given in (13) and  $\hat{z}$  resulting by inserting the realization  $\mathbf{Z} = \mathbf{z}$  in (11).  
 549 The resulting kriged mean  $\hat{z}$  and pointwise variance  $\hat{c}$  are displayed in Figure 4.

550 In summary, our approach for the stochastic modeling of the transmissivity  
 551 field  $T = T(\mathbf{x}, \omega)$  is characterized by

- (1) the assumptions that  $T$  has a lognormal distribution and that the covariance function of  $\log T$  belongs to the Matérn class; 553
- (2) obtaining the parameters  $\sigma$ ,  $\nu$  and  $\rho$  in the Matérn covariance function using *restricted maximum likelihood estimation (RML)*; 554
- (3) conditioning the thus obtained lognormal field on the available observations of transmissivity at the WIPP site. 555



**Fig. 4** Universal kriging prediction of  $Z = \log T$  based on 62 available transmissivity observations. Left: kriged mean field  $\hat{z}(\mathbf{x})$ . Right: pointwise kriging variance  $\hat{c}(\mathbf{x}, \mathbf{x})$ . The circular markers indicate the locations (and values) of the observational log transmissivity data. The interpolation property of  $\hat{z}(\mathbf{x})$  is apparent. 560

## 2.5 Uncertainty Propagation for the Quantity of Interest 561

For a random transmissivity field  $T(\omega) = T(\cdot, \omega)$ ,  $\omega \in \Omega$ , we consider individual realizations of the associated random boundary value problem in its mixed formulation (4), i.e., 562

$$\left( \frac{\phi b}{T(\omega)} \mathbf{u}(\omega), \mathbf{v} \right) - (p(\omega), \nabla \cdot \mathbf{v}) = -\langle g, \mathbf{n} \cdot \mathbf{v} \rangle_{\Gamma_D} \quad \forall \mathbf{v} \in \mathcal{V}, \quad (14a) \quad 563$$

$$(\nabla \cdot \mathbf{u}(\omega), q) = 0 \quad \forall q \in \mathcal{W}, \quad (14b) \quad 564$$

with random solution pair  $(\mathbf{u}(\omega), p(\omega)) \in \mathcal{V} \times \mathcal{W}$ . The equations (14) are now understood as holding  $\mathbf{P}$ -almost surely. Under suitable assumptions (cf. Babuška et al (2007)) we have for the random solution pair, that  $(\mathbf{u}, p) \in L^2_{\mathbf{P}}(\mathcal{V} \times \mathcal{W})$ , i.e., the norm of the solution is square integrable w.r.t. the probability measure  $\mathbf{P}$ . 565

For the quantity of interest under consideration, the exit time for particle trajectories, each realization of the random flux yields a realization of the associated random initial value problem 566

$$\dot{\mathbf{x}}(t, \omega) = \mathbf{u}(\mathbf{x}(t, \omega), \omega), \quad t \geq 0, \quad \mathbf{x}(0, \omega) = \mathbf{x}_0. \quad (15) \quad 567$$

599 **P**-almost surely, and hence, the quantity of interest becomes a random variable

$$600 \quad f_{\text{exit}}(\omega) := \log \min\{t > 0 : \mathbf{x}(t, \omega) \notin D_0, \mathbf{x}_0 \in D_0\}. \quad (16)$$

602 A complete characterization of the uncertainty in  $f$  is given by its cumulative  
603 distribution function (CDF)

$$604 \quad F(s) := \mathbf{P}(f_{\text{exit}} \leq s), \quad F: \mathbb{R} \rightarrow [0, 1].$$

605  
606 Due to the complexity of the problem,  $F$  cannot be given in analytic form and  
607 has to be approximated. We comment on the computational aspects in the  
608 next section.

## 609 **3 Computational Realization**

610 In this section we describe (i) the spatial discretization used for solving the  
611 Darcy flow equations (4) or (14), respectively, given a realization of the trans-  
612 missivity field  $T$ , (ii) a discrete representation of the random model for the  
613 transmissivity field  $T$  as well as (iii) a Monte Carlo approach for approximating  
614 the CDF of the quantity of interest.

### 615 **3.1 Finite Element Solution of Darcy Flow Problem**

616 We solve the Darcy flow equations (4) – or individual realizations of their  
617 random form (14) – using a mixed finite element discretization consisting of the  
618 lowest order Raviart-Thomas space  $\mathcal{V}_h \subset \mathcal{V}$  for the flux variable and piecewise  
619 constant space  $\mathcal{W}_h \subset \mathcal{W}$  for the hydraulic head with respect to a triangulation  
620  $\mathcal{T}_h$  of the domain  $D$ , where  $h > 0$  is a measure of mesh resolution. This  
621 discretization is known to be inf-sup-stable (cf. (Boffi et al, 2013, Chapter 7),  
622 (Ern and Guermond, 2021, Chapter 51)).

623 We choose a fixed triangulation of the two-dimensional computational  
624 domain with mesh width  $h$  chosen such that at least 10 elements correspond to  
625 the correlation length of the random transmissivity field, resulting in a mesh  
626 consisting of 28 993 triangles with the associated finite element spaces contain-  
627 ing 72 705 degrees of freedom (43 712 for flux and 28 993 for hydraulic head).  
628 Note that a coarser mesh is depicted in Figure 2 for illustration purposes. The  
629 particle tracking is performed by solving the ordinary differential equation  
630 (15) for the given realization. For the lowest-order Raviart-Thomas discretiza-  
631 tion, the constraint of zero divergence results in an elementwise constant flux,  
632 making this computation trivial and incurring no additional discretization  
633 error.

### 634 **3.2 Conditioned Karhunen-Loève Expansion**

635 Various methods exist to generate realizations of random fields, like turning  
636 bands, circulant embedding and Karhunen-Loève expansion, see Lord et al  
637  
638  
639  
640  
641  
642  
643  
644

(2014). In this work, we generate approximate realizations of the Gaussian log transmissivity field by truncating its Karhunen-Loève expansion, an orthogonal expansion of a random field based on the spectral decomposition of its covariance operator

$$C : L^2(D) \rightarrow L^2(D), \quad u \mapsto Cu, \quad (Cu)(\mathbf{x}) = \int_D c(\mathbf{x}, \mathbf{y})u(\mathbf{y}) \, d\mathbf{y}, \quad (17)$$

which for continuous covariance functions is compact and selfadjoint, positive definite and hence possesses a system of orthonormal eigenfunctions  $(z_m)_{m=1}^\infty$  which are complete in  $L^2(D)$ . Denoting by  $\lambda_m \geq 0$  the eigenvalue associated with eigenfunction  $z_m$  (ordered descending), a second-order random field  $Z$  on  $D$  with mean  $\bar{Z}$  possesses the expansion

$$Z(\mathbf{x}) = \bar{Z}(\mathbf{x}) + \sum_{m=1}^{\infty} \sqrt{\lambda_m} z_m(\mathbf{x}) \xi_m, \quad \mathbf{x} \in D, \quad (18)$$

converging in  $L^2$ , where  $(\xi_m)_{m \in \mathbb{N}}$  is a sequence of pairwise uncorrelated random variables and  $(\lambda_m)_{m \in \mathbb{N}}$  is square summable. In the present setting, the log transmissivity field  $Z$  is Gaussian, as stated in Section 2.3, therefore we have  $\xi_m \sim \mathbf{N}(0, 1)$  for all  $m$ .

An approximation suited for computation is obtained by truncating the infinite expansion in (17) after a finite number  $M$  of terms, hence the resulting approximation

$$Z(\mathbf{x}) \approx \bar{Z}(\mathbf{x}) + \sum_{m=1}^M \sqrt{\lambda_m} z_m(\mathbf{x}) \xi_m \quad (19)$$

for fixed  $M$  will depend on the decay rate of the eigenvalues.

Once a truncation index  $M$  has been fixed, the random field can be regarded as parameterized by the uncorrelated  $M$ -variate normal random vector  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_M)^\top \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$ , which takes values in  $\mathbb{R}^M$ . We may thus consider all random quantities in (14), i.e., the transmissivity field  $T$  and the solution  $(\mathbf{u}, p)$  of the Darcy flow equations as well as the particle trajectories (15) and exit time  $f_{\text{exit}}$  in (16) as parameterized by realizations of this single random vector.

Explicit closed-form solutions to the eigenvalue problem (17) are known only for a small number of special cases, hence we approximate the eigenpairs numerically. We approximate the covariance operator  $C$ , where the covariance kernel is obtained from the universal kriging covariance  $\hat{c}$  in (13), by Galerkin projection into a finite-dimensional subspace  $\mathcal{W}_h$  of  $L^2(D)$  consisting of piecewise constant functions with respect to a triangulation of the domain  $D$ , which we assume to be polygonal for simplicity<sup>3</sup>. Denoting by  $\{\phi_1, \dots, \phi_N\}$  a basis

---

<sup>3</sup>We use the same finite element space as for the piecewise constant discretization of the hydraulic head  $p$  for convenience.

of  $\mathcal{W}_h$ , we represent functions in  $\mathcal{W}_h$  as

$$u(\mathbf{x}) = \sum_{i=1}^N u_i \phi_i(\mathbf{x}) \quad (20)$$

with coefficient vector  $\mathbf{u} = (u_1, \dots, u_N)^\top$ . Substituting (20) into (17), multiplying it by test functions  $\phi_j$  and integrating over  $D$  we arrive at the discrete generalized eigenvalue problem

$$\mathbf{C}\mathbf{u} = \lambda\mathbf{M}\mathbf{u}, \quad (21)$$

where  $\mathbf{C}$  is a symmetric positive semi-definite matrix with entries

$$[\mathbf{C}]_{i,j} = (C\phi_i, \phi_j)_{L^2(D)} = \int_D \phi_j(\mathbf{x}) \int_D c(\mathbf{x}, \mathbf{y}) \phi_i(\mathbf{y}) \, d\mathbf{y} \, d\mathbf{x} \quad (22)$$

and  $\mathbf{M}$  is the symmetric positive definite Gram matrix of the piecewise constant basis with entries

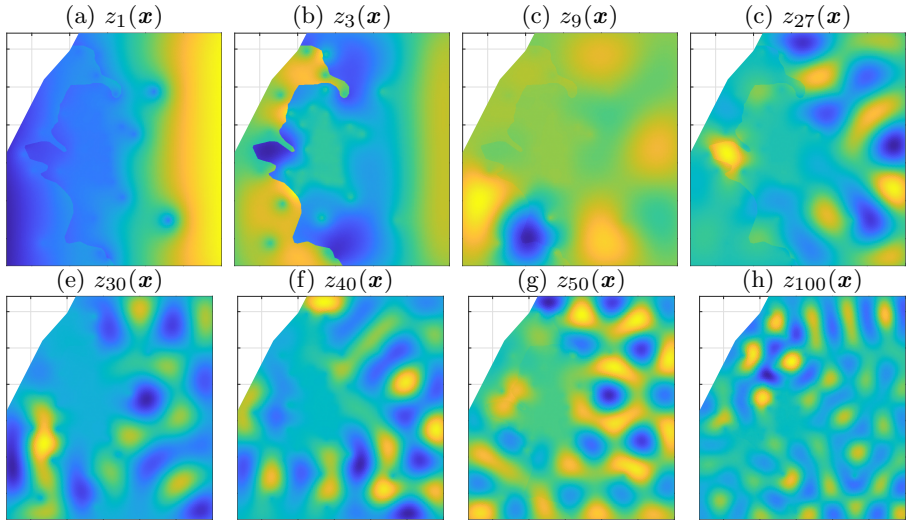
$$[\mathbf{M}]_{i,j} = (\phi_j, \phi_i)_{L^2(D)} = \int_D \phi_j(\mathbf{x}) \phi_i(\mathbf{x}) \, d\mathbf{x}. \quad (23)$$

An immediate difficulty with solving (21) is that  $\mathbf{C}$  is a dense matrix due to the nonlocal nature of the integral operator  $C$ , hence generating and storing  $\mathbf{C}$  is computationally expensive already for problems on two-dimensional domains, and even more so in three dimensions. Note that  $\mathbf{M}$  is diagonal due to the disjoint supports of the  $\phi_i$ . Moreover, even if generating and storing  $\mathbf{C}$  were feasible, solving a dense eigenvalue problem by the standard symmetric QR algorithm results in excessive computation costs. We address this problem by first using an iterative method for approximating only the dominant  $M$  eigenvalues of  $\mathbf{C}$  using a variant of the *thick-restart-Lanczos method* of Wu and Simon (2000), which requires only matrix vector products with  $\mathbf{C}$  in the course of the iteration. Second, we represent  $\mathbf{C}$  in *hierarchical matrix format* (cf. Hackbusch (2015)), which brings the cost of generating, storing and multiplying  $\mathbf{C}$  by a vector from  $\mathcal{O}(N^2)$  to a complexity  $\mathcal{O}(N \log N)$ . Further details on using hierarchical matrices in the context of random field generation with the Galerkin method can be found in Eiermann et al (2007) and Khoromskij et al (2009).

Figure 5 shows a few computed eigenfunctions  $z_m$  for the kriging covariance function  $\hat{c}$  in (13) displayed in Figure 4.

### 3.3 Empirical Estimation of the CDF

A common and convenient way to approximate the CDF  $F$  of the random quantity of interest  $f_{\text{exit}}(\boldsymbol{\xi}) := \log \min\{t > 0 : \mathbf{x}(t, \boldsymbol{\xi}) \notin D_0, \mathbf{x}_0 \in D_0\}$  is by generating  $n$  samples  $f_1, \dots, f_n$  of the random  $f_{\text{exit}}$  by sampling  $n$  different



**Fig. 5** Computed eigenfunctions of the kriging covariance function  $\hat{c}$  in (13), cf. Figure 4

realizations  $\xi_1, \dots, \xi_n$  of the random coefficient vector  $\xi$  in the KL expansion of the random field  $\log T$  and solving the corresponding  $n$  boundary and initial value problems to obtain  $f_i = f_{\text{exit}}(\xi_i)$ . We then determine the empirical CDF (ECDF)

$$F_n(s) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}_{(-\infty, f_j]}(s).$$

The ECDF  $F_n$  is a random approximation to the CDF  $F$  of the quantity of interest  $f_{\text{exit}}$  due to the randomly drawn samples  $f_1, \dots, f_n$ . We denote the error between the (random) ECDF and the true CDF by

$$D_n := \sup_{s \in \mathbb{R}} |F(s) - F_n(s)|. \quad (24)$$

For i.i.d. samples a classical result known as Donsker's theorem ([Athreya and Lahiri, 2006](#), Corollary 11.4.13) states

$$\sqrt{n}D_n \xrightarrow[n \rightarrow \infty]{d} \sup_{t \in [0,1]} |B(t)|,$$

where  $B$  denotes a standard Brownian bridge on the unit interval  $[0, 1]$ . This theoretical result can be employed to compute the necessary minimal sample size  $n$  for a desired error criterion, which we fix here by requiring

$$\mathbf{P}(D_n > 0.01) \leq 0.05. \quad (25)$$

Using the asymptotic result provided by Donsker's theorem as well as  $\mathbf{P}(\|B\|_{C[0,1]} > 1.36) \approx 0.05$ , see ([Williams, 2004](#), p. 343), we obtain for

737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782

783  $n \approx 20\,000$  that  $\mathbf{P}(D_n > 0.01) \approx 0.05$ . Hence, in the present setting this means  
 784 that, for this level of accuracy in approximating the CDF of the quantity of  
 785 interest, we need to solve  $n = 20\,000$  Darcy flow equations and compute the  
 786 associated particle trajectories. Thus, the question arises whether we could  
 787 save computational work by employing surrogates for the mapping from the  
 788 random parameter vector  $\boldsymbol{\xi}$  to the solution of the random PDE or the quantity  
 789 of interest  $f_{\text{exit}}$  itself.

790

### 791 *Estimation of CDF based on surrogates*

792 Assuming now that we have an approximation  $\hat{f}_{\text{exit}}: \mathbb{R}^M \rightarrow \mathbb{R}$  to the quantity  
 793 of interest  $f$  seen as mapping from  $\boldsymbol{\xi} \in \mathbb{R}^M \rightarrow \mathbb{R}$ , the resulting approximate  
 794 ECDF  $\hat{F}_n(s)$  based on  $n$  samples  $\hat{f}_1, \dots, \hat{f}_n$  of  $\hat{f}_{\text{exit}}$  resulting from  $n$  samples  
 795  $\boldsymbol{\xi}_i$  of the random KL parameter  $\boldsymbol{\xi}$ , where  $\hat{f}_i = \hat{f}_{\text{exit}}(\boldsymbol{\xi}_i)$  is given by

796

797

798

799

800

$$\hat{F}_N(s) = \frac{1}{N} \sum_{j=1}^N \mathbf{1}_{(-\infty, \hat{f}_j]}(s).$$

801 The question we investigate in this work is whether, for common surrogate  
 802 constructions such as stochastic collocation and Gaussian process emulators,  
 803 the approximation error  $\|f_{\text{exit}} - \hat{f}_{\text{exit}}\|$  (measured in a suitable norm) can  
 804 be made smaller than the sampling error  $D_n$  in the empirical estimation of  
 805 the CDF. To this end, we evaluate the quality of the surrogate  $\hat{f}_{\text{exit}}$  by a  
 806 two-sample *Kolmogorow-Smirnov (KS)* test which is a well-known hypothe-  
 807 sis test for checking whether sets of two samples—in our case  $\hat{f}_1, \dots, \hat{f}_n$  and  
 808  $f_1, \dots, f_n$ —are likely to have been drawn from the same distribution. Specifi-  
 809 cally, in our case the KS test is passed at significance level  $\alpha = 0.05$  if the  
 810 KS-statistic  $K$  satisfies

811

812

813

814

$$K := \sup_{s \in \mathbb{R}} \left| \hat{F}_n(s) - F_n(s) \right| \leq 1.36 \frac{\sqrt{2}}{n},$$

815 cf. [Williams \(2004\)](#).

816

## 817 **4 Propagation Surrogates**

818

819 In the following, we recall *sparse grid polynomial collocation* and *Gaussian pro-*  
 820 *cess emulators* as surrogate techniques for approximating a function  $f: \Xi \rightarrow \mathcal{Y}$   
 821 of  $M$  (random or parametric) variables  $\boldsymbol{\xi} \in \mathbb{R}^M$  taking values either in  $\mathcal{Y} = \mathbb{R}$ ,  
 822 as for scalar quantities of interest such as the breakthrough time, or a function  
 823 space, e.g.,  $\mathcal{Y} = \mathcal{V} \times \mathcal{W}$ , as for the solution of the mixed formulation (14) of  
 824 the Darcy flow equations with random conductivity.

825 We begin by illustrating the basic principles of polynomial collocation and  
 826 Gaussian process emulation for the case of a single variable, i.e.,  $\Xi \subseteq \mathbb{R}$ , before  
 827 proceeding to the technical details for the multivariate case  $\Xi \subseteq \mathbb{R}^M$ , where  
 828 we assume  $\Xi$  to be of product form  $\Xi = \Xi^M$  with  $\Xi \subseteq \mathbb{R}$ .



## 4.1 Univariate Collocation and Emulation

As a simple example in the style of the GPE tutorial O'Hagan (2006), consider the function

$$y = f(\xi) := \xi + 3 \sin \frac{3\xi}{4}, \quad \xi \in \Xi := [0, 6].$$

The presence of *input uncertainty*, i.e., uncertainty with regard to the precise value of the independent variable  $\xi$ , is accounted for by modeling it as a random variable  $\xi \sim \text{U}[0, 6]$ . Suppose further that  $f$  is only accessible in the form of a finite number of point evaluations  $f(\xi)$ , as is the case for the breakthrough time in our WIPP case study, where each evaluation of the former requires solving the Darcy flow problem followed by particle tracking up to the exit boundary. The task is to construct a computationally inexpensive approximation  $\hat{f}: \Xi \rightarrow \mathbb{R}$  of  $f$  given  $n$  evaluations

$$y_j = f(\xi_j), \quad j = 1, \dots, n.$$

The points of evaluation  $\xi_j$  are often called *design points* in the emulator literature and *nodes* or *knots* for collocation. Their choice depends on the type of surrogate being constructed. We begin with an elementary numerical analysis procedure and then contrast this with an approach rooted in the statistics community.

### Polynomial Collocation

In the univariate case polynomial collocation simplifies to Lagrange interpolation by global polynomials, and the surrogate  $\hat{f}$  for  $f$  takes the familiar form

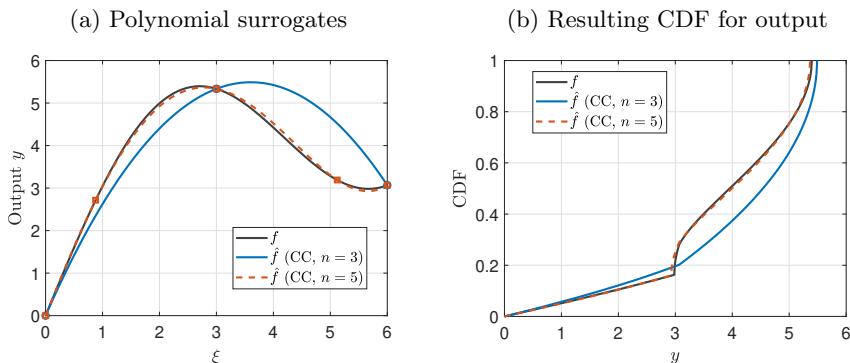
$$\hat{f}_n(\xi) := \sum_{j=1}^n f(\xi_j) \ell_j(\xi), \quad \ell_j(\xi) = \prod_{k \neq j} \frac{\xi - \xi_k}{\xi_j - \xi_k}$$

with  $\{\ell_j\}_{j=1}^n$  the Lagrange fundamental polynomials associated with the nodes  $\{\xi_1, \dots, \xi_n\}$ . Although this expression is well-defined for any set of distinct nodes, good approximation quality is only achieved if the points are chosen with care. A classical choice for bounded intervals is the family of *Clenshaw–Curtis nodes* (also called *Chebyshev nodes*). Scaled to the interval  $[0, 6]$ , the set of  $n$  Clenshaw–Curtis nodes is given by

$$\xi_j = 3 + 3 \cos \left( \frac{j-1}{n-1} \pi \right), \quad j = 1, \dots, n.$$

Other common choices, particularly for UQ applications, are the roots of the  $n$ -th orthogonal polynomial associated with the probability density of  $\xi$  on  $\Xi$ , e.g., Gauss–Legendre nodes for the uniform distribution or Gauss–Hermite nodes for the normal distribution, cf. Babuška et al (2010). For optimal convergence of the interpolants for smooth functions  $f$  it is well known that the spatial distribution of the nodes  $\xi_j \in \Xi$  should follow the equilibrium distribution in the sense of logarithmic potential theory, which for the standard interval

875  $\Xi = [-1, 1]$  is given by  $d\mu(\xi) = 1/\pi\sqrt{1 - \xi^2}$ , cf. (Trefethen, 2013, Chapter 12).  
 876 In particular, the nodes should cluster near the interval endpoints. Figure 6  
 877 shows two polynomial interpolation surrogates for  $f$  as well as the CDF of the  
 878 output  $f(\xi)$ .



891 **Fig. 6** The function  $f(\xi) = \xi + 3\sin(3\xi/4)$  on  $\Xi = [0, 6]$  and its Lagrange interpolation  
 892  $\hat{f}_n$  based on  $n = 3$  and  $n = 5$  Clenshaw-Curtis nodes (left) and the resulting CDF for the  
 893 output  $y = f(\xi)$  and  $\hat{y} = \hat{f}_n(\xi)$ , resp., if  $\xi \sim \mathcal{U}(\Xi)$ .  
 894

895  
 896 The approximation quality of polynomial interpolation depends not only  
 897 on the choice of interpolation nodes, but also on the *smoothness* of  $f$ . For  
 898 example, we have for  $f \in C^r(\Xi)$ ,  $r \in \mathbb{N}$ , that

$$900 \quad \|f - \hat{f}_n\|_\infty \leq c_r(f) n^{-r} (1 + \Lambda_{\xi_1, \dots, \xi_n})$$

901 where  $\|f - \hat{f}_n\|_\infty = \sup_{\xi \in \Xi} |f(\xi) - \hat{f}_n(\xi)|$ ,  $c_r(f)$  is a constant depending only  
 902 on  $r$  and  $f$ , and  $\Lambda_{\xi_1, \dots, \xi_n}$  denotes the *Lebesgue constant* of the nodes  $\xi_1, \dots, \xi_n$ .  
 903 Thus, we should choose nodes which have a small Lebesgue constant, and one  
 904 which grows only slowly with  $n$ . This is the case for Chebyshev and Clenshaw-  
 905 Curtis nodes, for which

$$906 \quad \Lambda_{\xi_1, \dots, \xi_n} \in \mathcal{O}(\log n).$$

907  
 908 Beside uniform convergence there are also classical results on convergence in  
 909 the  $L^p$  sense Nevai (1976, 1980, 1984), e.g., for Gauss-Legendre and Gauss-  
 910 Hermite nodes

$$911 \quad \lim_{n \rightarrow \infty} \|f - \hat{f}_n\|_{L_\mu^p} = 0, \quad \|f - \hat{f}_n\|_{L_\mu^p} = \left( \int_\Xi |f(\xi) - \hat{f}_n(\xi)|^p \mu(dx) \right)^{1/p},$$

912  
 913 where  $\mu = \mathcal{U}(\Xi)$  or  $\mu = \mathcal{N}(0, 1)$ , respectively. However, if  $f$  has low regularity  
 914 or is discontinuous, then convergence may fail or it may take a very large  
 915 number of nodes to approximate  $f$  with sufficient accuracy.

916  
 917 In summary, polynomial collocation constructs a (deterministic) interpolat-  
 918 ing polynomial as a surrogate for  $f$  based on evaluations of  $f$  at  $n$  judiciously  
 919  
 920

chosen nodes, for which the error decays with  $n$  at a rate depending on the smoothness of  $f$ .

### ***Gaussian Process Emulation***

The GPE approach consists in applying a method originating in geostatistics, namely the conditioning of Gaussian processes on observations (kriging), to the input-output map of a computer code, which for simplicity we assume to be represented by the scalar-valued function  $f: \Xi \rightarrow \mathbb{R}$ . Again, we are given a finite number of input-output pairs and the goal is to approximate function values at other locations. The uncertainty due to incomplete information at these locations is accounted for in the GPE approach by modeling  $f$  as an unknown realization of a Gaussian process indexed by  $\Xi$ , which is conditioned on the given observations  $y_j = f(\xi_j)$ —analogous to the conditioning of the Gaussian log transmissivity on observational values in Section 2.4.3. In this sense, GPE constructs not only a surrogate  $\hat{f}$  for  $f$ , but also a probabilistic representation of its pointwise deviation  $\hat{f}(\xi) - f(\xi)$ . In this sense the term *emulator* designates a statistical approximation of a function, which in this context is referred to as the *simulator* (cf. O’Hagan (2006)). Before we provide a more detailed discussion of this form of *output uncertainty quantification*, we briefly describe how a GPE surrogate is derived.

Analogously, to Section 2.3 we first choose a Gaussian process model  $G \sim \mathbf{N}(m, c)$  on  $\Xi$  with a (parametrized) mean function  $m: \Xi \rightarrow \mathbb{R}$ , e.g.,

$$m(\xi) = m(\xi; \boldsymbol{\beta}) = \sum_{k=1}^p \beta_k h_k(\xi), \quad \boldsymbol{\beta} \in \mathbb{R}^p,$$

and a (parametrized) covariance function  $c: \Xi \times \Xi \rightarrow \mathbb{R}$ , e.g., a Matérn covariance (6) or squared exponential covariance

$$c(\xi, \xi') = c(\xi, \xi'; \sigma^2, \rho) = \sigma^2 \exp(-(\xi - \xi')^2 / \rho), \quad \xi, \xi' \in \Xi. \quad (26)$$

In a true Bayesian approach, prior probability distributions are placed on the hyperparameters  $\boldsymbol{\beta}, \sigma^2, \rho$  of  $m$  and  $c$ . For now, however, we assume the covariance  $c$  to be fixed and  $m$  to be given as linear regression model—in analogy to Section 2.3. Conceptually, the Gaussian process describes our “prior beliefs” about the unknown  $f$  in the form of, e.g., characteristic dependencies reflected in the regression functions  $h_k$  in the mean model or smoothness properties encoded in the choice of  $c$ . Given evaluations  $f(\xi_j)$  of  $f$  at  $n$  design points  $\xi_j$ , we condition the Gaussian process  $G$  on this data and obtain  $\hat{G}_n \sim \mathbf{N}(\hat{m}_n, \hat{c}_n)$  with  $\hat{m}_n$  and  $\hat{c}_n$  determined by the relations for (simple or universal) kriging, see Section 2.4.3. The resulting *surrogate*  $\hat{f}_n$  is the conditional mean (or

967 kriging prediction) of  $\hat{G}_n$

968

969

970

971

972

973 where the coefficients  $\hat{\beta}_k$  and  $\hat{\gamma}_k$  depend on  $\xi_j$  and linearly on the  $f(\xi_j)$  and are  
 974 computed via universal kriging. We illustrate the GPE mean/surrogate for  $f$   
 975 as above and the resulting CDF for the output  $\hat{f}_n(\xi)$  if  $\xi \sim \text{U}[0, 6]$  in Figure 7.  
 976 Here we have used, similar to O'Hagan (2006),

977

$$978 \quad m(\xi; \boldsymbol{\beta}) = \beta_1 + \beta_2 \xi, \quad c(\xi, \xi') = \exp\left(-\frac{1}{4}(\xi - \xi')^2\right).$$

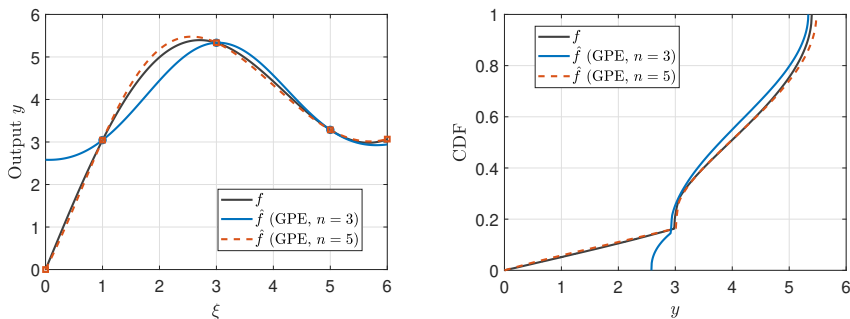
979

980

981

982 (a) GPE surrogates

(b) Resulting CDF for output



993 **Fig. 7** The function  $f(\xi) = \xi + 3 \sin(3\xi/4)$  on  $\Xi = [0, 6]$  and its GPE surrogates based on  
 994  $n = 3$  and  $n = 5$  design points  $\xi_j \in \{1, 3, 5\}$  and  $\xi_j \in \{0, 1, 3, 5, 6\}$  (left) and the resulting  
 995 CDF for the output  $y = f(\xi)$  and  $\hat{y} = \hat{f}_n(\xi)$ , resp., if  $\xi \sim \text{U}(\Xi)$ .

996

997 The choice of design points  $\xi_j$  for GPE follows different considerations  
 998 than for polynomial interpolation. It is well known that kriging coincides with  
 999 *kernel interpolation*, see Scheuerer et al (2013). If we assume for simplicity  
 1000 that  $m \equiv 0$  and  $c$  is given, then we can straightforwardly apply established  
 1001 approximation results from kernel interpolation theory by (Narcowich et al,  
 1002 2006, Proposition 3.2), (Wendland, 2004, Theorem 11.14), i.e., for  $f \in H^r(\Xi)$   
 1003 with  $r \geq 1$  and suitable<sup>4</sup> covariance functions  $c$  such as Matérn kernels (6)

1004

$$1005 \quad \|f - \hat{f}_n\|_\infty \leq C_r(f) D_{\xi_1, \dots, \xi_n}(\Xi)^{r - \frac{1}{2}}$$

1006

1007 where

$$1008 \quad D_{\xi_1, \dots, \xi_n}(\Xi) := \max_{\xi \in \Xi} \min_{j=1, \dots, n} |\xi - \xi_j|$$

1009

1010 

---

1011 <sup>4</sup>“Suitable” means here, that the *native or reproducing kernel Hilbert space* of  $c$  coincides with  
 1012  $H^r(\Xi)$ . For more details we refer to Scheuerer et al (2013); Wendland (2004).

denotes the *fill distance* of the node set  $\{\xi_1, \dots, \xi_n\}$ . For the Gaussian covariance function (26) we even obtain exponential convergence if the function  $f$  is *analytic*, see Wendland (2004),

$$\|f - \hat{f}_n\|_\infty \leq C(f) r^{D_{\xi_1, \dots, \xi_n}(\Xi)}, \quad r < 1.$$

Thus, for good approximation properties, GPE requires a *space filling* strategy for choosing design points, i.e., one which minimizes fill distance. In the univariate case that would be achieved by equispaced points, in stark contrast to the optimal equilibrium distribution for interpolation nodes.

As mentioned, a GPE not only provides a surrogate  $\hat{f}_n$  but also a probabilistic quantification of the remaining pointwise error  $f - \hat{f}_n$ , which represents another important difference to (polynomial) collocation. In order to understand this probabilistic error, recall that the conditioned Gaussian process  $\hat{G}_n$  can be seen as our “posterior belief” about the unknown  $f$  given  $n$  evaluations  $f(\xi_j)$ . Thus, as for the transmissivity field in subsurface flow (which is deterministic but unknown) we model our *uncertainty about the true output*  $f(\xi)$  at a *fixed input*  $\xi \in \Xi$  by  $\hat{G}_n(\xi) \sim \mathbf{N}(\hat{f}_n(\xi), \hat{c}_n(\xi))$ . This is called *code* or *output uncertainty* in the GPE literature, and is distinct from the *input uncertainty* modelled by *random*  $\xi$ : we have

input uncertainty:  $\xi$  random and  $\xi \mapsto f(\xi)$  fixed

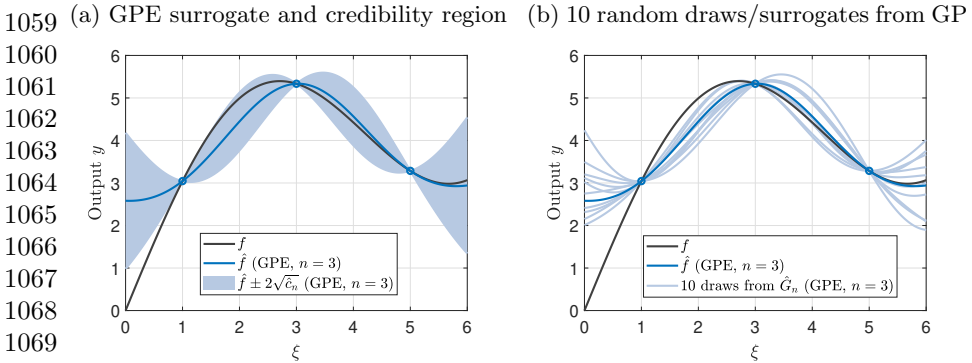
output uncertainty:  $\xi$  fixed and  $\xi \mapsto f(\xi)$  random

Of course, both uncertainty types can be superposed, as we shall see later. We illustrate the output uncertainty provided by the GPE in Figure 8: the left panel shows  $f$ ,  $\hat{f}_n$  as well as pointwise error estimates for  $f - \hat{f}_n$  given by two times the standard deviation of  $\hat{G}_n(\xi)$ , which can be also understood as the pointwise 95% credibility region for the unknown  $f(\xi)$ ; the right panel shows 10 realizations of the Gaussian process  $\hat{G}_n$ . Each of these could equally well be used as a surrogate  $\hat{f}_n$  in place of  $\hat{m}_n$ , since they are also valid (random) guesses for  $f$ . In this way,  $\hat{G}_n$  provides a *random* surrogate for  $f$ .

Random draws from  $\hat{G}_n$  can then be used to quantify the effect of the output uncertainty about the value  $f(\xi) \neq \hat{f}_n(\xi)$  within an uncertainty analysis for varying  $\xi$ , e.g., for estimating the CDF of  $f(\xi)$  when  $\xi \sim \mathbf{U}(\Xi)$ , see, e.g. Oakley and O’Hagan (2002). To explain this in more detail: Regarding the input uncertainty modelled by  $\xi \sim \mathbf{U}(\Xi)$  we would like to quantify its effect on the outcome by the CDF

$$F(y) = \mathbf{P}(f(\xi) < y).$$

This is a deterministic function for uncertainty analysis for random  $\xi$ . However, if we are not able to use  $f$  itself to compute  $F$  but rather use a GPE  $\hat{G}_n$  for  $f$ , we can, besides a deterministic approximation of  $F$  based on a deterministic



1070 **Fig. 8** The function  $f(\xi) = \xi + 3\sin(3\xi/4)$  on  $\Xi = [0, 6]$ , its GPE surrogate and the  
 1071 related 95% credibility region for  $f$  (left) as well as 10 paths (or surrogates) drawn from the  
 1072 conditioned GP  $\hat{G}_n$ .

1073  
 1074 surrogate  $\hat{f}_n$  for  $f$

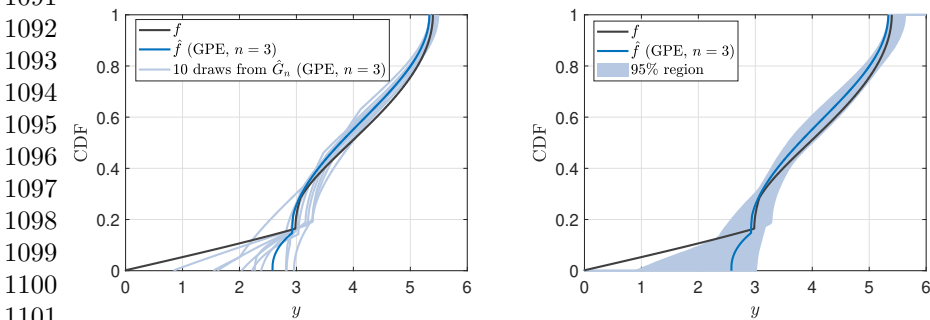
$$F(y) \approx \mathbf{P}_\xi(\hat{f}_n(\xi) < y),$$

1075  
 1076 also incorporate our remaining output uncertainty about  $f$  via the conditioned  
 1077 Gaussian process  $\hat{G}_n$  for  $f$ . This then yields a *random CDF*

$$\hat{F}_n(y) = \mathbf{P}_\xi(\hat{G}_n(\xi) < y),$$

1081 due to the random  $\hat{G}_n$  where we emphasize that the CDF is *only* w.r.t. ran-  
 1082 domness of the  $\xi$ . To illustrate this we show in Figure 9 the resulting CDFs  
 1083 for  $\hat{f}_n(\xi)$ ,  $\xi \sim \mathbf{U}(\Xi)$  using  $\hat{f}_n = \hat{m}_n$  as well as  $\hat{f}_n$  set to be each of the 10 draws  
 1084 from  $\hat{G}_n$  (left) as well as the 95% credibility region for the true (but unknown)  
 1085 CDF values  $F(y) = \mathbf{P}(f(\xi) < y)$  based on 10,000 draws from  $\hat{G}_n$ . The cred-  
 1086 ity region, thus, quantifies our uncertainty about the true CDF resulting  
 1087 from using a (random) surrogate instead of the true quantity of interest  $f$ .  
 1088

1089 (a) CDF for output resulting from GP draws (b) 95% credibility region for CDF based  
 1090 on GPE



1102 **Fig. 9** Resulting CDFs for the output  $\hat{y} = \hat{f}(\xi)$ ,  $\xi \sim \mathbf{U}(\Xi)$ , based on the mean and 10  
 1103 random draws from the GPE  $\hat{G}_n$  (left), and the resulting 95% credibility region for the CDF  
 1104 of  $y = f(\xi)$  derived from the GPE (right).

## Discussion

Polynomial collocation and Gaussian process emulators are well-established surrogate techniques using solely point evaluations of the underlying quantity of interest  $f$ , and both approaches rely on a certain smoothness of  $f$ . However, they also differ in several aspects. This includes the type of basis functions from which each surrogate is constructed (polynomials vs. kernel functions or radial basis functions) as well as the selection strategies for nodes  $\xi_j$  (equilibrium distribution in the sense of potential theory vs. space filling). Moreover, the GPE surrogate  $\hat{f}_n = \hat{m}_n$  is based on minimizing the *average error* w.r.t. an assumed probability distribution over a function space, whereas interpolation error bounds are obtained by a *worst-case error* analysis over a function class. We refer to Ritter (2000) for more details on these two contrasting approaches. In particular, for GPE we explicitly assume a probability distribution for the unknown function  $f$ , given by the prior Gaussian process model  $G$ , whereas for collocation we simply assume that  $f$  is sufficiently smooth. This prior probability distribution for  $f$  is then updated given the data  $f(\xi_j)$  in a Bayesian fashion. Thus, GPE can be related to *Bayesian numerical analysis*, see Diaconis (1988), or *probabilistic numerics*, see Hennig et al (2022), respectively, and be seen as a Bayesian approach to kernel interpolation. In particular, the conditioned (posterior) distribution for the unknown  $f$  provided by  $\hat{G}_n$  yields an indicator for the remaining (output) uncertainty about  $f$  after its evaluation at  $n$  nodes  $\xi_j$ . Of course, the assumption of Gaussianity for this computer output uncertainty is debatable. We refer to Bastos and O’Hagan (2009) for diagnostics to validate the GP ansatz as well as to Kracker et al (2010) for a performance study of GPE for “Gaussian” as well as “non-Gaussian”  $f$ .

## 4.2 Polynomial Sparse Grid Collocation

Polynomial collocation in the context of UQ or parametric problems can roughly be described as computing an  $M$ -variate polynomial approximation to  $f: \Xi \rightarrow \mathcal{Y}$ ,  $\Xi \subseteq \mathbb{R}^d$ , based on multivariate Lagrange interpolation. Sparse grid collocation uses sparse grids as multivariate interpolation node sets in order to mitigate the curse of dimensionality associated with straightforward tensor-product interpolation for high-dimensional parameter spaces.

While more sophisticated sparse grid techniques have been developed in recent years, in this work we consider a basic and simple construction known as (*Smolyak sparse grid collocation*) introduced for UQ settings, e.g., in Xiu and Hesthaven (2005); Nobile et al (2008). To this end, assume  $f \in C(\Xi; \mathcal{Y})$ , i.e., the mapping  $f$  is continuous, and denote by

$$\mathcal{P}_n(\Xi; \mathcal{Y}) = \left\{ \sum_{k=0}^n a_k \xi^k : a_k \in \mathcal{Y} \right\}$$

the space of all  $\mathcal{Y}$ -valued univariate polynomials of degree at most  $n$ . Then for a given sequence of univariate node sets  $\Xi_k := \{\xi_1^{(k)}, \dots, \xi_{n_k}^{(k)}\} \subseteq \Xi, k \geq 1$ , where we assume  $n_1 = 1$  and  $n_k < n_{k+1}$  throughout, we denote the associated



1151 univariate (Lagrange) interpolation operators by

1152

1153

$$1154 \quad \mathcal{I}_k: C(\Xi; \mathcal{Y}) \rightarrow \mathcal{P}_{n_k}(\Xi; \mathcal{Y}), \quad (\mathcal{I}_k f)(\xi) := \sum_{j=1}^{n_k} f\left(\xi_j^{(k)}\right) \ell_j^{(k)}(\xi), \quad \xi \in \Xi,$$

1155

1156

1157 with  $\ell_j^{(k)} \in \mathcal{P}_{n_k}(\Xi; \mathbb{R})$  the Lagrange fundamental polynomials associated  
 1158 with  $\Xi_k$ . The most immediate extension of the interpolation operator to the  
 1159  $M$ -dimensional parameter domain  $\Xi$  would be the multivariate interpolation  
 1160 operator  $\mathcal{I}_k: C(\Xi; \mathcal{Y}) \rightarrow \mathcal{P}_{n_k}(\Xi; \mathcal{Y})$  obtained by tensorization

1161

$$1162 \quad (\mathcal{I}_k f)(\boldsymbol{\xi}) := (\mathcal{I}_{k_1} \otimes \cdots \otimes \mathcal{I}_{k_M}) f(\boldsymbol{\xi}) = \sum_{j \leq n_k} f\left(\boldsymbol{\xi}_j^{(k)}\right) \ell_j^{(k)}(\boldsymbol{\xi}),$$

1163

1164

1165 with multi-indices  $\mathbf{j} = (j_1, \dots, j_M)$ ,  $\mathbf{n}_k = (n_{k_1}, \dots, n_{k_M}) \in \mathbb{N}^M$ , multivari-  
 1166 ate nodes  $\boldsymbol{\xi}_j^{(k)} = (\xi_{j_1}^{(k_1)}, \dots, \xi_{j_M}^{(k_M)}) \in \Xi_k := \Xi_{k_1} \times \cdots \times \Xi_{k_M}$ , and tensorized  
 1167 Lagrange fundamental polynomials  $\ell_j^{(k)}(\boldsymbol{\xi}) = \ell_{j_1}^{(k_1)}(\xi_1) \cdots \ell_{j_M}^{(k_M)}(\xi_M)$  for  $\boldsymbol{\xi} =$   
 1168  $(\xi_1, \dots, \xi_M) \in \Xi$ . However, this construction suffers heavily from the curse of  
 1169 dimensionality since the computational work for evaluating  $f$  at all points in  
 1170 the Cartesian product grid  $\Xi_k$  grows exponentially with dimension  $M$ .

1171

1172 Sparse grid constructions, which improve this to polynomial complexity in

1173

1174

$$1174 \quad \Delta_i = \mathcal{I}_i - \mathcal{I}_{i-1}, \quad i \geq 1, \quad \mathcal{I}_0 \equiv 0,$$

1175

1176

1177 so that  $\mathcal{I}_k = \sum_{i=1}^k \Delta_i$ , yielding the tensor product interpolation operator as

1178

$$1179 \quad \mathcal{I}_k f = \sum_{\mathbf{i} \leq \mathbf{k}} \Delta_{\mathbf{i}} f, \quad \Delta_{\mathbf{i}} = \Delta_{i_1} \otimes \cdots \otimes \Delta_{i_M}.$$

1180

1181

1182 By contrast, the (*Smolyak*) *sparse grid collocation operator* is defined by

1183

$$1184 \quad \mathcal{S}_{\ell, M} f := \sum_{|\mathbf{i}-\mathbf{1}|_1 \leq \ell} \Delta_{\mathbf{i}} f, \quad |\mathbf{i}-\mathbf{1}|_1 := \sum_{j=1}^M |i_j - 1|, \quad \ell \geq 0.$$

1185

1186

1187 By combinatorial arguments, one can obtain the equivalent *combination*  
 1188 *technique* representation

1189

1190

$$1191 \quad \mathcal{S}_{\ell, M} f = \sum_{\ell-M+1 \leq |\mathbf{i}-\mathbf{1}| \leq \ell} (-1)^{\ell+M-|\mathbf{i}|} \binom{M-1}{\ell+M-|\mathbf{i}|} \mathcal{I}_{\mathbf{i}} f,$$

1192

1193

1194

1195

1196

which expresses the Smolyak operator as a linear combination of selected  $M$ -variate tensor product interpolation operators. For the associated *sparse grid*

$$\Xi_{\ell, M} := \bigcup_{\ell-M+1 \leq |i-1| \leq \ell} \Xi_i$$

consisting of all multivariate nodes occurring in these representations, the cardinality  $|\Xi_{\ell, M}|$  grows only polynomially w.r.t.  $M$  (cf. Novak and Ritter (1999)), while the overall order of accuracy remains close to that of the full tensor product  $\mathcal{I}_{(\ell+1, \dots, \ell+1)}$ . In particular, it can be shown (Bäck et al, 2011, Proposition 1) that  $\mathcal{S}_{\ell, M}$  is a projection on

$$\mathcal{P}_{\ell, M}(\Xi; \mathcal{Y}) := \sum_{|i-1| \leq \ell} \mathcal{P}_{n_{i_1}}(\Xi; \mathcal{Y}) \otimes \dots \otimes \mathcal{P}_{n_{i_M}}(\Xi; \mathcal{Y}).$$

Note, however, that in general  $\mathcal{S}_{\ell, M}$  is *not* interpolatory unless the univariate nodes sets are *nested*  $\Xi_k \subset \Xi_{k+1}$  (Barthelmann et al, 2000, Proposition 6). The latter is the case for Clenshaw–Curtis nodes with the “doubling sequence”  $n_k = 2^k - 1$  ( $k \geq 1$ ), or (weighted) Leja nodes with linear growth  $n_k = k$  Ernst et al (2021). In the following, we shall use the non-nested nodal sequence of *Gauss-Hermite* nodes, i.e., the roots of Hermite polynomials. This choice is common for collocation applied to functions of Gaussian random variables, see Babuška et al (2007); Nobile et al (2008); Ernst and Sprungk (2014).

### Convergence and Application

If  $f$  is sufficiently smooth then  $\mathcal{S}_{\ell, M}f$  can be shown to converge to  $f$ , specifically

$$\|f - \mathcal{S}_{\ell, M}f\|_{L^2_{\mu}} \in \mathcal{O}(|\Xi_{\ell, M}|^{-r}),$$

for an  $r < 1$  using Gauss-Hermite nodes  $\xi_i^{(k)}$  with linear growth  $n_k = k$  or doubling growth  $n_k = 2^{k-1} + 1$  ( $k \geq 1$ ), see, e.g., Ernst and Sprungk (2014); Ernst et al (2018). The rate of convergence  $r$  w.r.t. the number of collocation points depends, of course, on the smoothness class of  $f$ .

It was shown in (Ernst and Sprungk, 2014, Section 3) that the solution  $(\mathbf{u}, p)$  of the random/parametric mixed variational problem (4) allows for a holomorphic extension into  $\mathbb{C}^M$  under suitable assumptions, which are satisfied by truncated KL expansions (19) of a lognormal transmissivity field. Thus, applying  $\mathcal{S}_{\ell, M}$  to approximate the solution map  $(\mathbf{u}, p): \Xi \rightarrow \mathcal{V} \times \mathcal{W}$  is justified. By contrast, the quantity of interest given by the breakthrough time  $f_{\text{exit}}$  may, in general, not even be a continuous function of the parameters  $\xi$ , as is immediate from considering the case of a particle grazing the exit boundary and returning into the domain for a particular parameter setting. Thus, applying  $\mathcal{S}_{\ell, M}$  to approximate  $f_{\text{exit}}$  directly may lead to inaccurate surrogate approximation or even divergence with increasing  $|\Xi_{\ell, M}|$ .

1243 However, a simple remedy is to use the surrogate

1244

$$1245 \hat{f}_{\text{exit},\ell} = G_{\text{exit}}(\mathcal{S}_{\ell,M}\mathbf{u})$$

1246

1247 where  $G_{\text{exit}}: \mathcal{V} \rightarrow \mathbb{R}$  denotes the mapping from a velocity field on  $D$  to the log  
 1248 breakthrough time of a particle following this field released at  $\mathbf{x}_0$  at time  $t = 0$ ,  
 1249 which is inexpensive to evaluate compared to solving the Darcy flow equations.  
 1250 Then, since  $L^2$ -convergence implies convergence in distribution, assuming that  
 1251 the set of points of discontinuity of the mapping  $G_{\text{exit}}$  has probability measure  
 1252 zero, we have by the continuous mapping theorem

1253

$$1254 \lim_{\ell \rightarrow \infty} \|F - \hat{F}_\ell\|_\infty = 0, \quad \hat{F}_\ell(s) := \mathbf{P}_{\boldsymbol{\xi} \sim \mu}(G_{\text{exit}}(\mathcal{S}_{\ell,M}\mathbf{u}(\boldsymbol{\xi})) \leq s),$$

1255

1256 where  $F$  denotes the true CDF of  $f_{\text{exit}}$ . Thus, we are assured convergence of  
 1257 the CDF based on the surrogate  $\mathcal{S}_{\ell,M}\mathbf{u}$  for the true velocity  $\mathbf{u}$  to the true  
 1258 CDF for the breakthrough time.

1259

### 1260 4.3 Gaussian Process Emulators

1261

1262 Having described basic GPE methodology in Section 4.1, we now turn to  
 1263 the construction of GPEs for multivariate scalar-valued functions  $f: \Xi \rightarrow \mathbb{R}$ .  
 1264 Again, the approach is similar to multivariate geostatistics. We shall con-  
 1265 sider the *full Bayesian* approach to GPE (cf. Kennedy and O'Hagan (2001);  
 1266 O'Hagan (2006)), which also entails specifying prior distributions for the  
 1267 hyperparameters contained in the mean and covariance functions which are  
 1268 also conditioned on the evaluations of  $f$  at the design points  $\boldsymbol{\xi}_j$ . As before, we  
 1269 start with a linear regression model for the mean

1270

$$1271 m: \Xi \rightarrow \mathbb{R}, \quad m(\boldsymbol{\xi}) = m(\boldsymbol{\xi}; \boldsymbol{\beta}) = \sum_{k=1}^p \beta_k h_k(\boldsymbol{\xi}), \quad \boldsymbol{\beta} \in \mathbb{R}^p,$$

1272

1273 with known regression functions  $\mathbf{h} = (h_1, \dots, h_p)$ ,  $h_k: \Xi \rightarrow \mathbb{R}$  ( $h_1 \equiv 1$  and  
 1274  $h_2(\boldsymbol{\xi}) = \boldsymbol{\xi}$  are common choices) and unknown coefficients  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ .  
 1275 For the emulator's covariance function  $c: \Xi \times \Xi \rightarrow \mathbb{R}$  we fix the squared  
 1276 exponential kernel

1277

$$1279 c(\boldsymbol{\xi}, \boldsymbol{\xi}') = c(\boldsymbol{\xi}, \boldsymbol{\xi}'; \sigma^2, B) = \sigma^2 \exp(-(\boldsymbol{\xi} - \boldsymbol{\xi}')^\top B(\boldsymbol{\xi} - \boldsymbol{\xi}')), \quad \boldsymbol{\xi}, \boldsymbol{\xi}' \in \Xi, \quad (27)$$

1280

1281 where  $\sigma^2 > 0$  is the marginal variance and  $B = \text{diag}(b_1, \dots, b_M) \in \mathbb{R}^{M \times M}$ ,  
 1282  $b_i > 0$  is a matrix of so-called *smoothness parameters*. For the squared exponen-  
 1283 tial covariance (27) and choices for  $h_1$  and  $h_2$  mentioned above, it is known that  
 1284 the realizations of the Gaussian process are almost surely analytic w.r.t.  $\boldsymbol{\xi}$ . For

1285

1286

1287

1288

other covariance functions, such as the family of Matérn kernels, one obtains Gaussian processes with realizations of different smoothness orders.<sup>5</sup>

Thus, for fixed given  $\boldsymbol{\beta}$ ,  $\sigma$ , and  $B$ , the (prior) Gaussian process model for the output of  $f$  for an arbitrary input  $\boldsymbol{\xi} \in \Xi$  is

$$f(\boldsymbol{\xi}) \sim \mathbf{N}(m(\boldsymbol{\xi}; \boldsymbol{\beta}), c(\boldsymbol{\xi}, \boldsymbol{\xi}; \sigma^2, B)).$$

Similarly, for fixed  $\boldsymbol{\beta}$ ,  $\sigma$ , and  $B$ , the vector  $\mathbf{f} = (f(\boldsymbol{\xi}_1), \dots, f(\boldsymbol{\xi}_n))^\top$  of values of the Gaussian process at a set of design points  $\{\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n\}$  has the  $n$ -variate Gaussian distribution

$$\mathbf{f} = (f(\boldsymbol{\xi}_1), \dots, f(\boldsymbol{\xi}_n))^\top \sim \mathbf{N}(\mathbf{H}\boldsymbol{\beta}, \mathbf{C}_{\sigma^2, B})$$

where  $\mathbf{H} = (h_k(\boldsymbol{\xi}_j)) \in \mathbb{R}^{n \times p}$  and  $\mathbf{C}_{\sigma^2, B} = (c(\boldsymbol{\xi}_i, \boldsymbol{\xi}_j; \sigma^2, B)) \in \mathbb{R}^{n \times n}$ . We denote the probability density of this random vector  $\mathbf{f} \in \mathbb{R}^n$  by

$$p(\mathbf{f} \mid \boldsymbol{\beta}, \sigma, B) \propto \exp\left(-\frac{1}{2}(\mathbf{f} - \mathbf{H}\boldsymbol{\beta})^\top \mathbf{C}_{\sigma^2, B}^{-1}(\mathbf{f} - \mathbf{H}\boldsymbol{\beta})\right).$$

Suitable values for the parameters  $\boldsymbol{\beta}$ ,  $\sigma$ , and  $B$  are usually not known a priori and should be inferred based on the evaluations  $\mathbf{f}$ . This is typically done in a Bayesian fashion, i.e., we choose hyperpriors for these parameters which are then conditioned on the data  $\mathbf{f} = (f(\boldsymbol{\xi}_1), \dots, f(\boldsymbol{\xi}_n))^\top$ . Common choices for  $(\boldsymbol{\beta}, \sigma^2)$  are a normal-inverse-gamma prior or a Jeffreys prior with density  $p(\boldsymbol{\beta}, \sigma^2) \propto \sigma^{-2}$  (cf. [Oakley and O'Hagan \(2002\)](#); [Stone \(2011\)](#)) since these allow for closed-form expressions for the resulting (marginal) posteriors. Given evaluations  $\mathbf{f}$ , the resulting posterior for the parameters  $(\boldsymbol{\beta}, \sigma^2)$  is then

$$p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{f}, B) \propto p(\mathbf{f} \mid \boldsymbol{\beta}, \sigma, B)p(\boldsymbol{\beta}, \sigma^2).$$

For the estimation of the smoothness parameters  $B$  a “full” Bayesian inference based on data  $\mathbf{f}$  would require Markov chain Monte Carlo simulations. Instead, one often simply computes a point estimate based on maximizing the marginal likelihood  $p(\mathbf{f} \mid B) \propto \int p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{f}, B)p(\boldsymbol{\beta}, \sigma^2) d\boldsymbol{\beta} d\sigma^2$  for which analytic formulas are available ([Stone, 2011](#), Section 2.3.4). This often yields competitive results to a full Bayesian inference [Kracker et al \(2010\)](#).

Given  $\mathbf{f}$ , the posterior density for the output  $f(\boldsymbol{\xi})$  at new location  $\boldsymbol{\xi}$  is then

$$p(f(\boldsymbol{\xi}) \mid \mathbf{f}, \boldsymbol{\beta}, \sigma, B) \propto p(\mathbf{f} \mid \boldsymbol{\beta}, \sigma, B)p(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{f}, B).$$

Marginalization by integrating out  $\boldsymbol{\beta}$  and  $\sigma^2$  can be done analytically for a normal-inverse-gamma or Jeffreys prior  $p(\boldsymbol{\beta}, \sigma^2)$  and results in a *Student-t*

---

<sup>5</sup>We have also explored other covariance models such as the Matérn kernels for GPE surrogates; however, the overall conclusions in the numerical experiments were about the same as for the squared exponential ([27](#)).

1335 process (cf. Shah et al (2014)) for the prediction of the output of  $f$ , i.e.,

$$1336 \quad f(\boldsymbol{\xi}) \mid \mathbf{f} \sim t_{n-p}(\hat{m}_n(\boldsymbol{\xi}), \hat{\sigma}^2 \hat{c}_n(\boldsymbol{\xi}, \boldsymbol{\xi})), \quad (28)$$

1338 where  $\hat{m}_n$  and  $\hat{c}_n$  are the mean and covariance obtained by universal kriging applied to  $f$  given the observations  $\mathbf{f}$  (see (??) and (??)) with  $\sigma^2 = 1$ , 1340 respectively, and where  $\hat{\sigma}^2$  is given by

$$1342 \quad \hat{\sigma}^2 = \frac{1}{n-p} \mathbf{f}^\top \mathbf{C}^{-1/2} \left( \mathbf{I} - \mathbf{C}^{-1/2} \mathbf{H} (\mathbf{H}^\top \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}^{-1/2} \right) \mathbf{C}^{-1/2} \mathbf{f}.$$

1344 For the prediction of  $f$  at multiple new points we obtain a multivariate Student- $t$ -distribution with mean vector given by the evaluation of  $\hat{m}_n$  at those points 1346 and covariance matrix given by evaluating  $\hat{\sigma}^2 \hat{c}_n$ .

1348 Regarding the choice of the design points for multivariate GPE we require again *space filling* designs. For compact  $\Xi \subset \mathbb{R}^M$  these are, e.g., Sobol' points 1350 (Owen et al, 2017) or Latin hypercube designs (Viana, 2015). The latter extend 1351 also to  $\Xi = \mathbb{R}^M$  w.r.t.  $\mu = \mathbf{N}(0, \mathbf{I})$  as we require for the WIPP problem. As 1352 for the appropriate number  $n \in \mathbb{N}$  of training points  $\Xi_n = \{\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n\} \subset \Xi$ , 1353 a common rule of thumb calls for  $n = cM$  (Loeppky et al, 2009) with a factor 1354  $c \geq 10$ .

### 1356 *Convergence and Application*

1357 Since the GPE surrogate  $\hat{f}_n = \hat{m}_n$  and its covariance  $\hat{c}_n$  are derived by uni- 1358 versal kriging, we can again exploit the relation between kriging and kernel 1359 interpolation Scheuerer et al (2013). Again, assume  $m \equiv 0$  for simplicity and 1360  $c$  fixed as in (27). Then for compact  $\Xi \subset \mathbb{R}^M$  and analytic  $f: \Xi \rightarrow \mathbb{R}$  we have

$$1362 \quad \|f - \hat{f}_n\|_\infty \leq C(f) r^{\mathbf{D}_{\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n}(\Xi)},$$

1364 for a  $0 < r < 1$  as well as

$$1366 \quad \hat{c}_n(x, x) \leq C r^{2\mathbf{D}_{\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n}(\Xi)}.$$

1368 Thus, besides uniform convergence of the surrogate  $\hat{f}_n \rightarrow f$ , we also have 1369 vanishing output uncertainty regarding  $f(\boldsymbol{\xi})$  as  $n \rightarrow \infty$ —which is a consistency 1370 statement for the posterior for  $f$  here given by the Gaussian or Student- $t$  1371 process  $\hat{G}_n$ . However, to our knowledge, no  $L^2$ -convergence statements are 1372 available for the case of unbounded  $\Xi = \mathbb{R}^M$ , as the setting of the WIPP 1373 problem would require.

1375 In the next section we will apply GPE to approximate the quantity of inter- 1376 est  $f_{\text{exit}}$  directly. Thus, for convergence with  $n \rightarrow \infty$ , we require  $f_{\text{exit}}$  to be 1377 sufficiently smooth (see above) which may not be the case in general. How- 1378 ever, it may well be that the surrogate  $\hat{f}_n$  and the related output uncertainty 1379 provided by the GPE for finite  $n = cM$  design points is sufficiently accurate 1380 for CDF estimation. We note that also vector-valued GPE are available, see

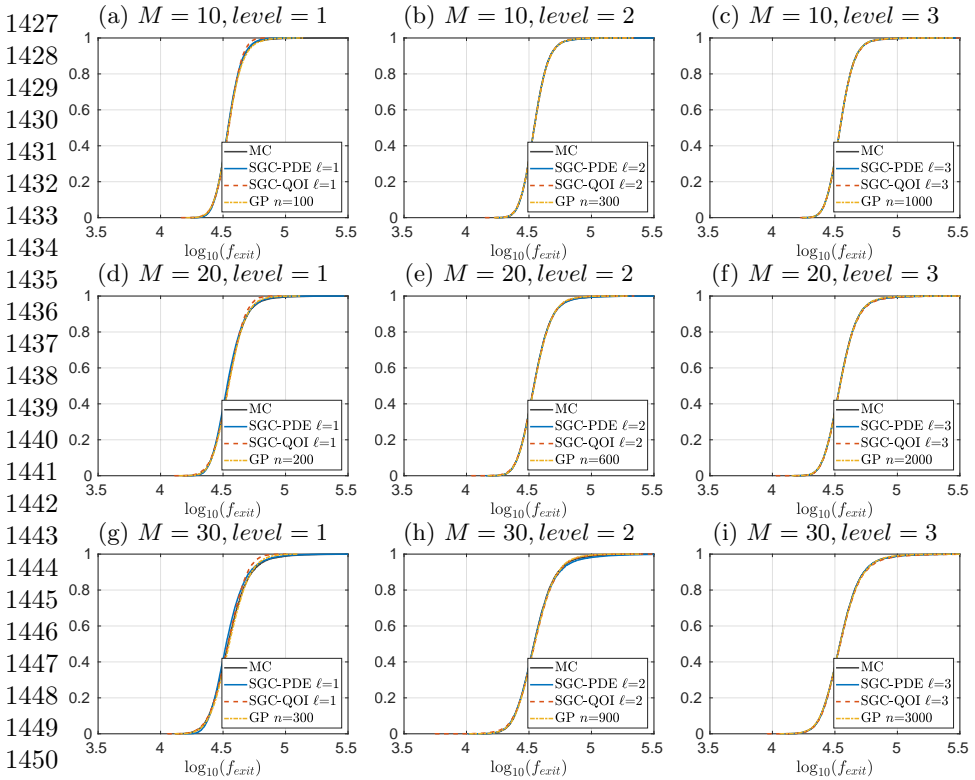
Álvarez et al (2012); Bilonis et al (2013); Cleary et al (2021); Higdon et al (2008). Hence, we could apply a GPE to approximate the FE solution of the random parametric variational problem (which depends analytically on  $\xi$ , see comment above) and proceed as for polynomial collocation to provide approximate samples of  $f_{\text{exit}}(\xi)$ . We do not consider this option in this work, since the FE space is very high dimensional (of order  $10^4$ ) and thus the GPE would involve too many parameters to estimate based on less than 20,000 design points.

## 5 Numerical Results

We now perform a numerical study comparing sparse grid polynomial collocation and Gaussian process emulators as surrogates for the task of approximating the CDF of the exit time  $f_{\text{exit}}(\xi)$  using  $M$  terms or coefficients  $\xi \sim \mathbf{N}(0, \mathbf{I})$  in the truncated KL expansion of the log transmissivity field  $Z = \log T$ . We vary  $M = 10, 20, 30$  and apply the following three surrogate approaches:

- **SGC-PDE:** We apply Smolyak sparse grid polynomial collocation  $\mathcal{S}_{\ell, M}$  to approximate the solution pair  $(\mathbf{u}, p)$  of the weak mixed formulation and then obtain approximate samples  $\hat{f}_{\text{exit}}(\xi_i)$  of the exit time by solving the particle transport given the approximate velocity field  $\mathcal{S}_{\ell, M}\mathbf{u}(\xi_i)$ , i.e.,  $\hat{f}_{\text{exit}}(\xi_i) = G_{\text{exit}}(\mathcal{S}_{\ell, M}\mathbf{u}(\xi_i))$  where  $\xi_i \sim \mathbf{N}(0, \mathbf{I})$ ,  $i = 1, \dots, N$  iid.
- **SGC-QoI:** We apply Smolyak sparse grid polynomial collocation  $\mathcal{S}_{\ell, M}$  directly to approximate the exit time  $f_{\text{exit}}(\xi_i)$  and in this way obtain approximate samples via  $\hat{f}_{\text{exit}}(\xi_i) = \mathcal{S}_{\ell, M}f_{\text{exit}}(\xi_i)$  where  $\xi_i \sim \mathbf{N}(0, \mathbf{I})$ ,  $i = 1, \dots, N$  iid.
- **GPE:** We apply Gaussian process emulation to approximate the exit time  $f_{\text{exit}}(\xi_i)$  and obtain approximate samples via  $\hat{f}_{\text{exit}}(\xi_i) = \hat{m}_n(\xi_i)$  where  $\xi_i \sim \mathbf{N}(0, \mathbf{I})$ ,  $i = 1, \dots, N$  iid and  $\hat{m}_n$  denotes the GPE mean.

For each surrogate we generate  $N = 20\,000$  approximate samples of the quantity of interest and compare these to  $N = 20\,000$  samples of the “true”  $f_{\text{exit}}$  evaluated by solving the Darcy flow equations and particle transport problem each time (denoted **MC** for Monte Carlo in the following). The number  $N = 20\,000$  of samples is derived from the error criterion outlined in Section 3.3. For SGC we use different levels  $\ell = 1, 2, 3$ , and for the GPE different numbers of design points  $n = cM$  with  $c = 10, 20, 30, 50, 100$ . We show the resulting empirical CDFs for the log exit time in Figure 10. It is apparent that, for each  $M = 10, 20, 30$  all surrogate methods yield a very good fit to the reference ECDF obtained by the plain Monte Carlo approach. Slight deviations can be seen for the lowest level  $\ell = 1$  for **SGC-QoI**, but at least for  $\ell \geq 2$  it is difficult to distinguish the four ECDFs. Therefore, we take a closer look at the performance of the surrogates in Table 2, where we report the resulting values of the KS statistic  $K = \sup_{s \in \mathbb{R}} \left| \hat{F}_n(s) - F_n(s) \right|$  of the empirical CDF  $F_n$  obtained by Monte Carlo sampling of  $f_{\text{exit}}$  and the empirical CDF  $\hat{F}_n$  obtained by Monte Carlo sampling of the surrogate  $\hat{f}_{\text{exit}}$ . Moreover,



1451 **Fig. 10** Empirical CDFs obtained by Monte Carlo, SGC and GPE surrogates for different  
 1452 lengths  $M$  of the KLE.

1453

1454 we indicate by an asterisk that the error  $K$  in the ECDFs is negligible, i.e.,  
 1455 that the Kolmogorov–Smirnov test is passed (at significance level  $\alpha = 0.05$ ),  
 1456 and hence there is no indication that the samples were drawn from different  
 1457 distributions. We make the following observations:

- 1458 • For  $M = 10, 20$  all three surrogates pass the KS-test at least for level  
 1459  $\ell \geq 2$  (SGC) or  $n \geq 30M$  design points (GPE). For  $M = 30$  this is  
 1460 also the case for SGC-PDE with  $\ell \geq 2$  and GPE with  $n = 100M$ . Thus,  
 1461 by employing the considered surrogates we can obtain an ECDF for the  
 1462 exit time which is essentially indistinguishable (for  $\alpha = 0.05$ ) from the  
 1463 “true” ECDF but which required just a fraction of the computational  
 1464 cost of the latter. Indeed, compared to  $N = 20\,000$  solutions of the Darcy  
 1465 flow equations, we require merely between  $\approx 200$  ( $M = 10$ ) and  $\approx 2000$   
 1466 ( $M = 30$ ) PDE solves when a surrogate is used.
- 1467 • For SGC-PDE we observe a steep increase in  $n$  with  $M$  but overall a  
 1468 robust and good performance.
- 1469 • For the SGC-QoI approach we observe a significantly worse performance  
 1470 for  $M = 30$  which may be due to insufficient (mixed) smoothness of  $f_{\text{exit}}$ .

1471

1472



Surrogate		M = 10		M = 20		M = 30	
		$n$	$K$	$n$	$K$	$n$	$K$
SGC-PDE	$\ell = 1$	21	0.0128*	41	0.0281	61	0.0495
SGC-PDE	$\ell = 2$	241	0.0028*	881	0.0045*	1921	0.0118*
SGC-PDE	$\ell = 3$	2001	0.0019*	13201	0.0023*	41601	0.0052*
SGC-QoI	$\ell = 1$	21	0.0271	41	0.0293	61	0.0435
SGC-QoI	$\ell = 2$	241	0.0065*	881	0.0088*	1921	0.0196
SGC-QoI	$\ell = 3$	2001	0.0048*	13201	0.0089*	41601	0.0138
GPE	$c = 10$	100	0.0136	200	0.0245	300	0.0309
GPE	$c = 20$	200	0.0092*	400	0.0191	600	0.0228
GPE	$c = 30$	300	0.0062*	600	0.0116*	900	0.0171
GPE	$c = 50$	500	0.0041*	1000	0.0070*	1500	0.0141
GPE	$c = 100$	1000	0.0031*	2000	0.0064*	3000	0.0087*

**Table 2** Performance of the SGC and GPE surrogates for different lengths  $M$  of the KL expansion measured by the value of the resulting KS statistic  $K$ . Here,  $n$  refers to the number of PDEs to be solved for building the surrogate and an asterisk denotes that the KS-test was passed at significance level  $\alpha = 0.05$ .

- For the GPE approach we observe deteriorating performance for increasing  $M$ , i.e., we require a larger factor  $c$  for the number of design points  $n = cM$  in order to pass the KS test and have small values of  $K$  ( $c = 20$  for  $M = 10$ ,  $c = 30$  for  $M = 20$  and  $c = 100$  for  $M = 30$ ). This may be due to the curse of dimensionality for kernel interpolation methods.

### Changing the trend model for $\log T$

Despite the overall positive observations for the employed surrogates made so far we report how the outcome may change if we simply use a different trend model for the mean of the log transmissivity field  $\log T$ . Instead of using the constant, linear in  $x_1$ , and zone indicator regression functions  $h_1$ ,  $h_2$ , and  $h_5$ , respectively, see (9), we only use the constant  $h_1$ . This leads to a different Matérn covariance function used for  $\log T$ , see Table 1 and thus also to different eigenvalues and eigenfunctions in the KL expansion. Moreover, the smoothness properties of the mapping  $\xi \mapsto f_{\text{exit}}(\xi)$  may change as well. In fact, in Table 3 we observe a much diminished performance of all three surrogate techniques: Now only SGC-PDE passes the KS test and only for the shorter KL truncation length  $M = 10, 20$ . However, SGC-PDE and GPE provide a visually acceptable fit to the reference ECDF in Figure 11, whereas we clearly see a deterioration for the SGC-QoI surrogate. This distinctly worse performance of SGC-QoI may be due to lacking smoothness of  $\xi \mapsto f_{\text{exit}}(\xi)$  in this case.

For the GPE surrogate we also evaluate to what extent the accompanying Gaussian model for this surrogate's output uncertainty covers the deviation from the reference CDF. To this end, we focus on the setting where the GPE surrogate performs worst, i.e.,  $M = 30$  using  $n = 300$  design points, and compute a 95% credibility region for the CDF based on 10 000 random draws of surrogates from the trained GPE. The results are reported in Figure 12 for both trend models. We observe that the Gaussian output uncertainty model appears overconfident in the case of the constant trend model. Thus, this

	Surrogate		M = 10		M = 20		M = 30	
			<i>n</i>	<i>K</i>	<i>n</i>	<i>K</i>	<i>n</i>	<i>K</i>
1519	SGC-PDE	$\ell = 1$	21	0.0537	41	0.0653	61	0.0621
1520	SGC-PDE	$\ell = 2$	241	0.0123*	881	0.0130*	1921	0.0146
1521	SGC-PDE	$\ell = 3$	2001	0.0121*	13201	0.0345	41601	0.0387
1522	SGC-QOI	$\ell = 1$	21	0.1099	41	0.1340	61	0.1301
1523	SGC-QOI	$\ell = 2$	241	0.0485	881	0.0798	1921	0.0697
1524	SGC-QOI	$\ell = 3$	2001	0.0369	13201	0.0577	41601	0.1711
1525	GPE	$c = 10$	100	0.0366	200	0.0546	300	0.0815
1526	GPE	$c = 20$	200	0.0373	400	0.0415	600	0.0591
1527	GPE	$c = 30$	300	0.0153	600	0.0368	900	0.0615
1528	GPE	$c = 50$	500	0.0188	1000	0.0405	1500	0.0415
1529	GPE	$c = 100$	1000	0.0192	2000	0.0258	3000	0.0422

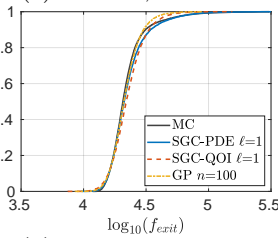
1529 **Table 3** Rerun of Table 2 but for constant mean for log  $T$ .

1530

1531

1532 experiment indicates that a sufficiently good performance of the surrogates for  
 1533 CDF estimation of exit times may depend on various aspects of the problem—  
 1534 such as the choice of the trend model for the log transmissivity field.

1535

1536 (a)  $M = 10, level = 1$ 

1537

1538

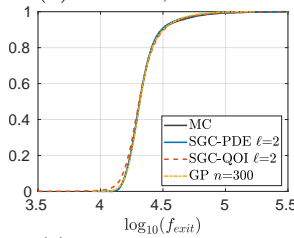
1539

1540

1541

1542

1543

1536 (b)  $M = 10, level = 2$ 

1544

1545

1546

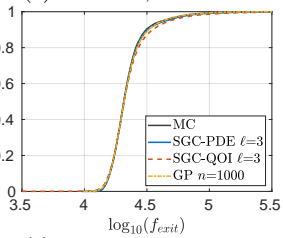
1547

1548

1549

1550

1551

1536 (c)  $M = 10, level = 3$ 

1544

1545

1546

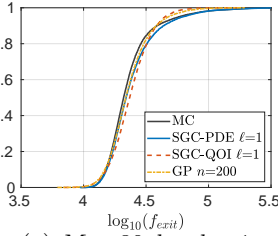
1547

1548

1549

1550

1551

1544 (d)  $M = 20, level = 1$ 

1544

1545

1546

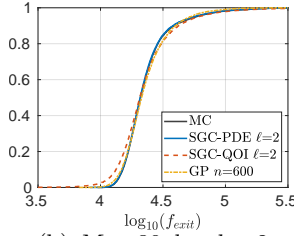
1547

1548

1549

1550

1551

1544 (e)  $M = 20, level = 2$ 

1544

1545

1546

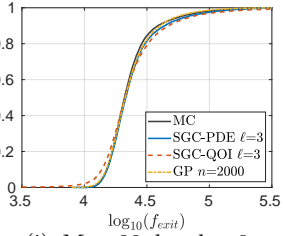
1547

1548

1549

1550

1551

1544 (f)  $M = 20, level = 3$ 

1544

1545

1546

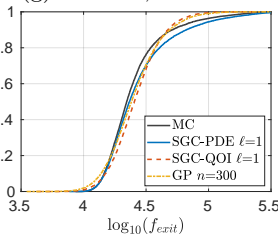
1547

1548

1549

1550

1551

1552 (g)  $M = 30, level = 1$ 

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

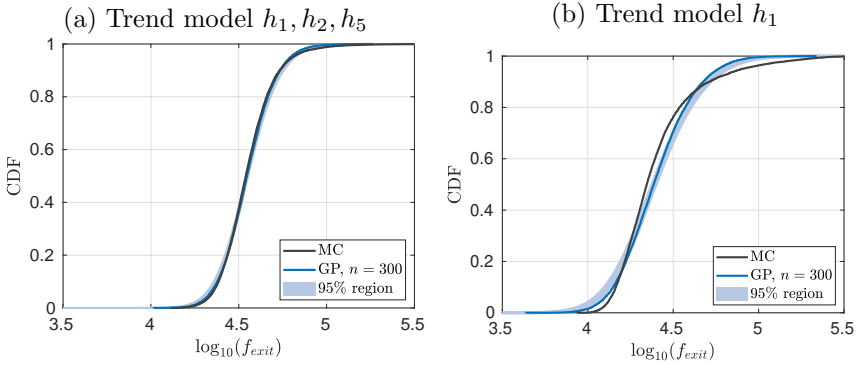
1561

1562

1563

1564

1560 **Fig. 11** Rerun of Figure 10 but for constant mean for log  $T$ .



**Fig. 12** 95% credibility region for CDF of breakthrough time based on GPE with  $n = 300$  for  $M = 30$  KL terms for different trend models

### Convergence Study

The negative results for the constant trend model raise the question whether we simply did not use enough design points  $n$  or sparse collocation level  $\ell$  for the GPE and SGC surrogates, respectively, or whether the quantity of interest is simply too rough to be approximated well by these methods. To this end, we perform a convergence study for both scenarios: constant trend model and “best” trend model using  $h_1, h_2$ , and  $h_5$  in (9). We report the associated  $L^2_\mu$ -errors of the SGC surrogates for the flux  $\mathbf{u}$  and the quantity of interest in Tables 4 and 5, respectively. We notice significantly larger errors for the constant trend model. In order to allow for a sufficiently high polynomial degree for SGC to observe a significant error decay, we restrict ourselves to the low-dimensional case of  $M = 2$  and  $M = 5$  KL terms. We report the resulting errors of the velocity and the quantity of interest in Figure 13. There we clearly observe a decaying error for increasing level  $\ell$  and number of sparse grid nodes  $|\Xi_{\ell, M}|$ , respectively. Moreover, we observe that the rate of convergence for both quantities is affected by the larger number of KL terms and the choice of the trend model. The former was already observed in Ernst and Sprungk (2014). The latter is also related to an observation made in Ernst and Sprungk (2014): since the constant trend model yields a larger estimated value for the variance  $\sigma^2$ , this in turn leads to a slower convergence rate of SGC.

Regarding the application of GPE to approximate the quantity of interest, we perform a similar study as for SGC using  $M = 2$  and  $M = 5$  KL terms. The results are displayed in Figure 14. We observe that the  $L^2_\mu$ -error (left panel) does not decay with increasing number of design points, at least not in the applied regime of up to  $n = 1000M$  design points. However, this is mainly due to high approximation errors in the tail regions of the distribution of  $\boldsymbol{\xi} \sim \mu = \mathbf{N}(0, \mathbf{I})$ . Despite this, we observe a decay of the KS test statistic value  $K$ , i.e., the  $L^\infty$ -error of the ECDF for the quantity of interest, except for  $M = 5$  and the constant trend model.

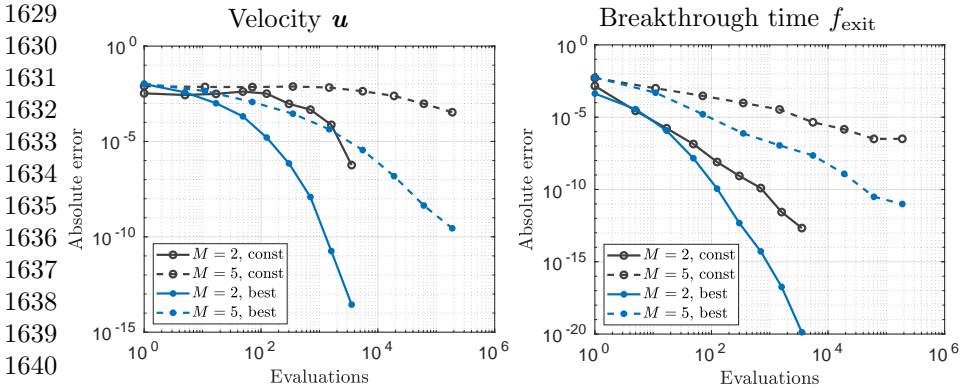
1565  
1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610

	Trend model	Surrogate	M = 10	M = 20	M = 30
1611	$h_1, h_2, h_5$	SGC $\ell = 1$	5.9897E-3	1.2933E-2	1.6810E-2
1612		SGC $\ell = 2$	2.1354E-3	6.3868E-3	9.3400E-3
1613		SGC $\ell = 3$	6.1168E-4	2.5686E-3	4.3738E-3
1614	$h_1$	SGC $\ell = 1$	4.0723E-2	1.1149E-1	1.7963E-1
1615		SGC $\ell = 2$	4.0331E-2	1.1113E-1	1.7329E-1
1616		SGC $\ell = 3$	3.9595E-2	1.0598E-1	1.6928E-1

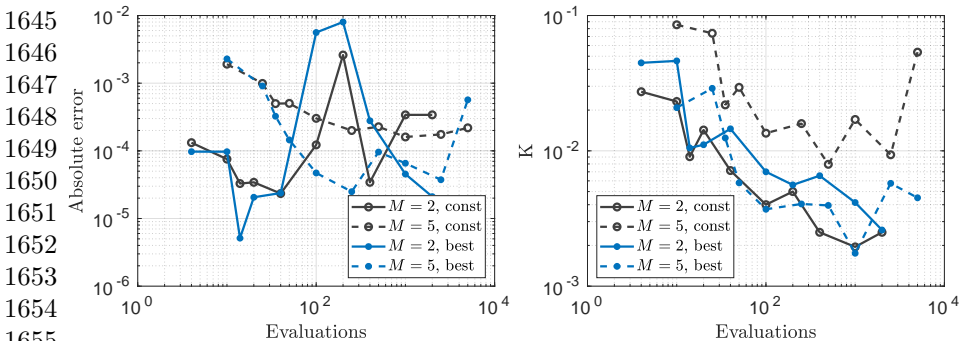
1617 **Table 4**  $L^2_\mu(\Xi, \mathbf{H}(\text{div}; D))$  error of SGC surrogates for the flux  $\mathbf{u}$  for the two different  
 1618 trend models.

	Trend model	Surrogate	M = 10	M = 20	M = 30
1620	$h_1, h_2, h_5$	SGC $\ell = 1$	1.2296E-3	2.7434E-3	6.0602E-3
1621		SGC $\ell = 2$	1.2699E-4	4.8917E-4	2.1426E-3
1622		SGC $\ell = 3$	2.0075E-5	9.6514E-5	4.4401E-4
1623	$h_1$	SGC $\ell = 1$	7.0990E-3	1.5259E-2	2.8396E-2
1624		SGC $\ell = 2$	2.9464E-3	7.7314E-3	1.5502E-2
1625		SGC $\ell = 3$	1.9730E-3	9.2632E-3	1.8001E-2

1626 **Table 5**  $L^2_\mu(\Xi, \mathbb{R})$  error of SGC surrogates for the exit time  $f_{\text{exit}}$  for the two different  
 1627 trend models.



1641 **Fig. 13**  $L^2_\mu$ -error of SGC surrogates for the velocity  $\mathbf{u}$  (left) and exit time  $f_{\text{exit}}$  (right).  
 1642 For the flux we used the norm in  $\mathbf{H}(\text{div}; D)$  to quantify the difference between  $\mathbf{u}(\xi)$  and  
 1643  $\mathcal{S}_{\ell, M} \mathbf{u}(\xi)$ .



1654 **Fig. 14**  $L^2_\mu$ -error (left) and K-S test value  $K$  (right) of GPE surrogates for exit time  $f_{\text{exit}}$ .  
 1655  
 1656

## 6 Conclusion

In this work we have presented a complete uncertainty propagation workflow for groundwater flow and particle transport simulations based on a real-world application related to the site performance assessment for a nuclear waste repository. We described in detail the construction of a stochastic model for an uncertain transmissivity field by geostatistical methods using the available observational data. Our main focus was the direct comparison of two established surrogate approaches for uncertainty propagation analysis: sparse grid stochastic collocation and Gaussian process emulation. Both methods originate from different communities, i.e., numerical analysis and computational statistics, respectively. Our purpose was to describe the fundamental ideas and principles underlying both methods and compare their performance for the UQ problem under consideration, specifically for CDF estimation of scalar quantities of interest, in this case the travel time of groundwater-borne radionuclides. The overall conclusion is that both methods can achieve significant reduction in computational cost, reducing the computational burden by a factor of 10 to even 100 in some cases considered. Moreover, we have observed that the GPE surrogate seems to be more adversely affected by the high dimensionality of the input space compared with sparse grid collocation, which is not surprising given the bad scaling of the filling distance with dimension. On the other hand, stochastic collocation must also be applied with care, since the quantity to be approximated has to depend sufficiently smoothly on the random inputs—such as the solution of the random PDE. However, the remarkable performance of both surrogates seems to be affected by modelling choices for the random log transmissivity field such as choice of the trend or regression model for the mean. Although this effect could be explained mathematically in our case, it does place limitations on the practical benefits of UQ surrogate methods for CDF estimation in groundwater flow applications.

**Acknowledgments.** The authors would like to thank Elisabeth Ullmann (TU München) as well as Gerald van den Boogaart and Silke Konsulke (HZDR Dresden-Rossendorf) for numerous contributions in an early phase of this work. The first two authors gratefully acknowledge their invitation to the 2018 Uncertainty Quantification Programme at the Isaac Newton Institute, where they participated in numerous clarifying discussions on surrogates.

**Funding and/or Conflicts of interests/Competing interests.** The first and third author gratefully acknowledge partial funding support from the German society for radioactive waste disposal BGE in the MeQUR project within the research cluster URS. The authors have no competing interests to declare that are relevant to the content of this article.

1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702

1703 **References**

- 1704  
 1705 Álvarez MA, Rosasco L, Lawrence ND (2012) Kernels for vector-valued  
 1706 functions: A review. *Foundations and Trends® in Machine Learning*  
 1707 4(3):195–266. <https://doi.org/10.1561/22000000036>
- 1708  
 1709 Athreya KB, Lahiri SN (2006) *Measure Theory and Probability Theory*.  
 1710 Springer-Verlag
- 1711  
 1712 Babuška I, Nobile F, Tempone R (2007) A stochastic collocation method  
 1713 for elliptical partial differential equations with random input data. *SIAM*  
 1714 *Journal on Numerical Analysis* 45(3):1005–1034
- 1715  
 1716 Babuška I, Nobile F, Tempone R (2010) A stochastic collocation method for  
 1717 elliptic partial differential equations with random input data. *SIAM Review*  
 1718 52(1):317–355
- 1719  
 1720 Bäck J, Nobile F, Tamellini L, et al (2011) Stochastic spectral Galerkin  
 1721 and collocation methods for PDEs with random coefficients: A numerical  
 1722 comparison. In: *Spectral and High Order Methods for Partial Differential*  
 1723 *Equations*. Springer Berlin Heidelberg, pp 43–62, [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-642-15337-2_3)  
 1724 [978-3-642-15337-2\\_3](https://doi.org/10.1007/978-3-642-15337-2_3)
- 1725  
 1726 Barthelmann V, Novak E, Ritter K (2000) High dimensional polynomial  
 1727 interpolation on sparse grids. *Advances in Computational Mathematics*  
 1728 12:273–288
- 1729  
 1730 Bastos LS, O’Hagan A (2009) Diagnostics for gaussian process emulators.  
 1731 *Technometrics* 51(4):524–438. URL <https://www.jstor.org/stable/40586652>
- 1732  
 1733 Bilonis I, Zabaras N, Konomi BA, et al (2013) Multi-output separable gaus-  
 1734 sian process: Towards an efficient, fully bayesian paradigm for uncertainty  
 1735 quantification. *Journal of Computational Physics* 241:212–239. [https://doi.](https://doi.org/10.1016/j.jcp.2013.01.011)  
 1736 [org/10.1016/j.jcp.2013.01.011](https://doi.org/10.1016/j.jcp.2013.01.011)
- 1737  
 1738 Boffi D, Brezzi F, Fortin M (2013) *Mixed Finite Element Methods and Appli-*  
 1739 *cations*, Springer Series in Computational Mathematics, vol 44. Springer  
 1740 Science & Business Media, <https://doi.org/10.1007/978-3-642-36519-5>
- 1741  
 1742 Cleary E, Garbuno-Inigo A, Lan S, et al (2021) Calibrate, emulate, sample.  
 1743 *Journal of Computational Physics* 424:109,716. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.jcp.2020.109716)  
 1744 [jcp.2020.109716](https://doi.org/10.1016/j.jcp.2020.109716)
- 1745  
 1746 Cressie NA (1991) *Statistics for Spatial Data*. Wiley-Interscience
- 1747  
 1748 Currin C, Mitchell T, Morris M, et al (1991) Bayesian prediction of deter-  
 ministic functions, with applications to the design and analysis of computer

- experiments. *Journal of the American Statistical Association* 86(416):953–963 1749  
1750  
1751
- de Marsily G (1986) *Quantitative Hydrogeology: Groundwater Hydrology for Engineers*. Academic Press 1752  
1753  
1754
- Diaconis P (1988) Bayesian numerical analysis. In: Gupta SS, Berger JO (eds) *Statistical Decision Theory and Related Topics IV*, vol 1. Springer, New York, NY, pp 163–175 1755  
1756  
1757
- Eiermann M, Ernst OG, Ullmann E (2007) Computational aspects of the stochastic finite element method. *Computing and visualization in science* 10(1):3–15 1758  
1759  
1760  
1761
- Ern A, Guermond JL (2021) *Finite Elements II: Galerkin Approximation, Elliptic and Mixed PDEs*, Texts in Applied Mathematics, vol 73. Springer Nature Switzerland AG, <https://doi.org/10.1007/978-3-030-56923-5> 1762  
1763  
1764  
1765
- Ernst OG, Sprungk B (2014) Stochastic collocation for elliptic PDEs with random data: The lognormal case. In: Garcke J, Pflüger D (eds) *Sparse Grids and Applications – Munich 2012*, LNCSE, vol 97. Springer International Publishing, p 29–53 1766  
1767  
1768  
1769  
1770
- Ernst OG, Sprungk B, Tamellini L (2018) Convergence of sparse collocation for functions of countably many Gaussian random variables (with application to elliptic PDEs). *SIAM Journal on Numerical Analysis* 56(2):877–905. <https://doi.org/10.1137/17m1123079>, URL <https://doi.org/10.1137/17m1123079> 1771  
1772  
1773  
1774  
1775
- Ernst OG, Sprungk B, Tamellini L (2021) On expansions and nodes for sparse grid collocation of lognormal elliptic PDEs. In: Bungartz HJ, Garcke J, Pflüger D (eds) *Sparse Grids and Applications – Munich 2018*, LNCSE, vol 144. Springer International Publishing, p 1–31, [https://doi.org/10.1007/978-3-030-81362-8\\_1](https://doi.org/10.1007/978-3-030-81362-8_1) 1776  
1777  
1778  
1779  
1780
- Ernst OG, Pichler A, Sprungk B (2022) Wasserstein sensitivity of risk and uncertainty propagation. *SIAM/ASA Journal on Uncertainty Quantification* 10(3):915–948. <https://doi.org/10.1137/20M1325459> 1781  
1782  
1783  
1784  
1785
- Freeze RA (1975) A stochastic-conceptual analysis of one-dimensional groundwater flow in nonuniform homogeneous media. *Water Resources Research* 11(5):725–741. <https://doi.org/10.1029/WR011i005p00725> 1786  
1787  
1788
- Ghanem RG, Spanos PD (1991) *Stochastic Finite Elements: A Spectral Approach*. Springer-Verlag, New York 1789  
1790  
1791
- Graham IG, Scheichl R, Ullmann E (2016) Mixed finite element analysis of lognormal diffusion and multilevel Monte Carlo methods. *Stoch PDE Anal* 1792  
1793  
1794



1795 Comp 4:41–75

1796

1797 Gunzburger MD, Webster CG, Zhang G (2014) Stochastic finite element meth-  
1798 ods for partial differential equations with random input data. *Acta Numerica*  
1799 23:521–650. <https://doi.org/10.1017/S0962492914000075>

1800

1801 Hackbusch W (2015) *Hierarchical matrices: algorithms and analysis*, vol 49.  
1802 Springer

1803

1804 Hennig P, Osborne MA, Kersting HP (2022) *Probabilistic Numerics - Compu-*  
1805 *tation as Machine Learning*. Cambridge University Press, [https://doi.org/](https://doi.org/10.1017/9781316681411)  
1806 [10.1017/9781316681411](https://doi.org/10.1017/9781316681411)

1807

1808 Higdon D, Gattiker J, Williams B, et al (2008) Computer model calibra-  
1809 tion using high-dimensional output. *Journal of the American Statistical*  
1810 *Association* 103(482):570–583

1811

1812 Hoeksema RJ, Kitanidis PK (1985) Analysis of the spatial structure of prop-  
1813 erties of selected aquifers. *Water Resources Research* 21(4):563–572. [https://doi.org/](https://doi.org/10.1029/WR021i004p00563)  
1814 [10.1029/WR021i004p00563](https://doi.org/10.1029/WR021i004p00563)

1815

1816 Kennedy MC, O’Hagan A (2001) Bayesian calibration of computer models. *J*  
1817 *R Statistical Society Part B* 63(Part 3):425–464

1818

1819 Khoromskij BN, Litvinenko A, Matthies HG (2009) Application of hierarchical  
1820 matrices for computing the Karhunen-loève expansion. *Computing* 84(1-  
1821 2):49–67

1822

1823 Kitanidis PK (1987) Parametric estimation of covariances of regionalized  
1824 variables. *Water Resources Bulletin* 23(4):557–567

1825

1826 Kitanidis PK (1997a) *Introduction to Geostatistics: Applications to Hydroge-*  
1827 *ology*. Cambridge University Press

1828

1829 Kitanidis PK (1997b) A variance-ratio test for supporting a variable mean  
1830 in kriging. *Mathematical Geology* 29(3):335–348. [https://doi.org/10.1007/](https://doi.org/10.1007/BF02769639)  
1831 [BF02769639](https://doi.org/10.1007/BF02769639)

1832

1833 Kracker H, Bornkamp B, Kuhnt S, et al (2010) Uncertainty in gaussian  
1834 process interpolation. In: Devroye L, Karasözen B, Kohler M, et al (eds)  
1835 *Recent Developments in Applied Probability and Statistics: Dedicated to*  
1836 *the Memory of Jürgen Lehn*. Physica-Verlag HD, Heidelberg, pp 79–102,  
[https://doi.org/10.1007/978-3-7908-2598-5\\_4](https://doi.org/10.1007/978-3-7908-2598-5_4)

1837

1838 Linde N, Ginsbourger D, Irving J, et al (2017) On uncertainty quantification in  
1839 hydrogeology and hydrogeophysics. *Advances in Water Resources* 110:166–  
1840 181. <https://doi.org/10.1016/j.advwatres.2017.10.014>

- Liu D, Litvinenko A, Schillings C, et al (2017) Quantification of airfoil geometry-induced aerodynamic uncertainties—comparison of approaches. *SIAM/ASA Journal on Uncertainty Quantification* 5(1):334–352. <https://doi.org/10.1137/15M1050239>
- Loeppky JL, Sacks J, Welch WJ (2009) Choosing the sample size of a computer experiment: A practical guide. *Technometrics* 51(4):366–376
- Lord GJ, Powell CE, Shardlow T (2014) *An Introduction to Computational Stochastic PDEs*, Cambridge Texts in Applied Mathematics, vol 50. Cambridge University Press
- Narcowich F, Ward J, Wendland H (2006) Sobolev error estimates and a Bernstein inequality for scattered data interpolation via radial basis functions. *Constr Approx* 24:175–186. <https://doi.org/10.1007/s00365-005-0624-7>
- Nevai GP (1976) Mean convergence of lagrange interpolation, i. *Journal of Approximation Theory* 18(4):363–377. [https://doi.org/10.1016/0021-9045\(76\)90008-3](https://doi.org/10.1016/0021-9045(76)90008-3)
- Nevai P (1984) Mean convergence of lagrange interpolation. III. *Trans Amer Math Soc* 282(2):669–698. <https://doi.org/10.2307/1999259>
- Nevai PG (1980) Mean convergence of lagrange interpolation, ii. *Journal of Approximation Theory* 30(4):263–276. [https://doi.org/10.1016/0021-9045\(80\)90030-1](https://doi.org/10.1016/0021-9045(80)90030-1)
- Nobile F, Tempone R, Webster CG (2008) A sparse grid stochastic collocation method for partial differential equations with random input data. *SIAM J Numer Anal* 46(5):2309–2345
- Novak E, Ritter K (1999) Simple cubature formulas with high polynomial exactness. *Constructive Approximation* 15:499–522
- Oakley J, O’Hagan A (2002) Bayesian inference for the uncertainty distribution of computer model outputs. *Biometrika* 89(4):769–784
- O’Hagan A (2006) Bayesian analysis of computer code outputs: A tutorial. *Reliability Engineering and System Safety* 91:1290–1300. <https://doi.org/10.1016/j.res.2005.11.025>
- Owen NE, Challenor P, Menon PP, et al (2017) Comparison of surrogate-based uncertainty quantification methods for computationally expensive simulators. *SIAM/ASA J Uncertainty Quantification* 5:403–436. <https://doi.org/10.1137/15M1046812>

- 1887 Ritter K (2000) Average-Case Analysis of Numerical Problems. No. 1733 in  
 1888 Lecture Notes in Mathematics, Springer-Verlag, <https://doi.org/10.1007/>  
 1889 [BFb0103934](https://doi.org/10.1007/BFb0103934)
- 1890
- 1891 Sacks J, Welch WT, Mitchell TJ, et al (1989) Design and analysis of computer  
 1892 experiments. *Statistical Science* 4(4):409–423
- 1893
- 1894 Scheuerer M, Schaback R, Schlather M (2013) Interpolation of spatial data:  
 1895 A stochastic or a deterministic problem? *European Journal of Applied*  
 1896 *Mathematics* 24(4):601–629. <https://doi.org/10.1017/S0956792513000016>
- 1897
- 1898 Schwab C, Gittelsohn CJ (2011) Sparse tensor discretizations of high-  
 1899 dimensional parametric and stochastic PDEs. *Acta Numerica* 20:291–467
- 1900
- 1901 Shah A, Wilson A, Ghahramani Z (2014) Student-t Processes as Alternatives  
 1902 to Gaussian Processes. In: Kaski S, Corander J (eds) *Proceedings of the*  
 1903 *Seventeenth International Conference on Artificial Intelligence and Statis-*  
 1904 *tics, Proceedings of Machine Learning Research*, vol 33. PMLR, Reykjavik,  
 1905 Iceland, pp 877–885, URL <https://proceedings.mlr.press/v33/shah14.html>
- 1906
- 1907 Stein M (1999) *Interpolation of Spatial Data: Some Theory for Kriging*.  
 1908 Springer, New York, [https://doi.org/10.1007/978-1-4612-1494-6\\_1](https://doi.org/10.1007/978-1-4612-1494-6_1)
- 1909
- 1909 Stone N (2011) *Gaussian process emulators for uncertainty analysis in*  
 1910 *groundwater flow*. PhD thesis, University of Nottingham
- 1911
- 1912 Trefethen LN (2013) *Approximation Theory and Approximation Practice*.  
 1913 SIAM, Philadelphia, <https://doi.org/10.1137/1.9781611975949>
- 1914
- 1915 U.S. Department of Energy (DOE) (2004) Title 40 CFR Part 191 Subparts  
 1916 B and C. Compliance Recertification Application 2004 for the Waste Isola-  
 1917 tion Pilot Plant Appendix TFIELD-2004 Transmissivity Fields. Tech. Rep.  
 1918 DOE/WIPP 2004/3231, Carlsbad Field Office, Carlsbad, NM
- 1919
- 1920 U.S. Department of Energy (DOE) (2014) Title 40 CFR Part 191 Subparts  
 1921 B and C. Compliance Recertification Application 2014 for the Waste Isola-  
 1922 tion Pilot Plant Appendix TFIELD-2014 Transmissivity Fields. Tech. Rep.  
 1923 DOE/WIPP 14-3503, Carlsbad Field Office, Carlsbad, NM
- 1924
- 1925 Viana FAC (2015) A tutorial on latin hypercube design of experiments. *Quality*  
 1926 *and Reliability Engineering International* 32:1975–1985. [https://doi.org/](https://doi.org/10.1002/qre.1924)  
 1927 [10.1002/qre.1924](https://doi.org/10.1002/qre.1924)
- 1928
- 1928 Wendland H (2004) *Scattered Data Approximation*. Cambridge University  
 1929 Press, <https://doi.org/10.1017/cbo9780511617539>
- 1930
- 1931
- 1932

- Williams D (2004) Weighing the Odds, 2nd edn. Cambridge University Press, Cambridge 1933  
1934  
1935
- Wu K, Simon H (2000) Thick-restart Lanczos method for large symmetric eigenvalue problems. SIAM Journal on Matrix Analysis and Applications 22(2):602–616 1936  
1937  
1938  
1939
- Xiu D, Hesthaven JS (2005) High-order collocation methods differential equations with random inputs. SIAM Journal on Scientific Computing 37(3):1118–1139 1940  
1941  
1942  
1943

## A Restricted Maximum Likelihood Estimation 1944

Under the models for the mean (8) and covariance structure (6), the Gaussian log-transmissivity field (7) has the form 1945  
1946

$$Z(\mathbf{x}, \omega) = \mathbf{h}(\mathbf{x})^\top \boldsymbol{\beta} + \tilde{Z}(\mathbf{x}, \omega) \quad (29) \quad 1947$$

with  $\mathbf{h}(\mathbf{x})^\top = [h_1(\mathbf{x}), \dots, h_k(\mathbf{x})]$  denoting the vector of regression functions evaluated at  $\mathbf{x} \in D$  and  $\boldsymbol{\beta} \in \mathbb{R}^k$  the vector of regression coefficients. The residual field  $\tilde{Z}$  is Gaussian with zero mean. The covariance function  $c_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y})$  of  $Z$  (and  $\tilde{Z}$ ) is given by (6), where  $\boldsymbol{\theta} = (\sigma, \rho, \nu)$  denotes the triplet of parameters consisting of standard deviation  $\sigma$ , correlation length  $\rho$  and smoothness parameter  $\nu$ . The specification of the probabilistic model for the random field  $Z$  consists in determining the vector  $\boldsymbol{\beta}$  of regression coefficients and the covariance parameter vector  $\boldsymbol{\theta}$ . It is desired that estimation techniques for these based on observations be *unbiased*, i.e., that the average estimation error be zero, and that this error be optimal in a least-squares sense. Another desirable property is *consistency* in the sense that the estimates converge to the true values as more and more observations are added. 1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962

The restriction of  $Z$  to a finite set of observation points  $\{\mathbf{x}_j\}_{j=1}^n \subset D$  is a multivariate Gaussian random vector, which we denote by 1963  
1964

$$\mathbf{Z} = \mathbf{Z}(\omega) = \begin{bmatrix} Z(\mathbf{x}_1, \omega) \\ \vdots \\ Z(\mathbf{x}_n, \omega) \end{bmatrix} : \Omega \rightarrow \mathbb{R}^n. \quad (30) \quad 1965  
1966  
1967  
1968  
1969$$

In view of (29), its expectation is 1970

$$\mathbf{E}[\mathbf{Z}] = \mathbf{H}\boldsymbol{\beta}, \quad [\mathbf{H}]_{i,j} = h_j(\mathbf{x}_i), \quad i = 1, \dots, n, \quad j = 1, \dots, k, \quad 1971  
1972$$

and its joint probability density function is given by 1973  
1974

$$p(\boldsymbol{\xi}; \boldsymbol{\beta}, \boldsymbol{\theta}) = \frac{1}{\sqrt{(2\pi)^n \det(\mathbf{C}_{\boldsymbol{\theta}})}} \exp\left(-\frac{1}{2}(\boldsymbol{\xi} - \mathbf{H}\boldsymbol{\beta})^\top \mathbf{C}_{\boldsymbol{\theta}}^{-1}(\boldsymbol{\xi} - \mathbf{H}\boldsymbol{\beta})\right), \quad \boldsymbol{\xi} \in \mathbb{R}^n, \quad 1975  
1977  
(31) \quad 1978$$

1979 in which  $\mathbf{C}_\theta$  denotes the covariance matrix

1980

$$1981 \quad \mathbf{C}_\theta = \mathbf{E} [\mathbf{Z}\mathbf{Z}^\top] = [c_\theta(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^n \in \mathbb{R}^{n \times n}$$

1982

1983 of the random vector  $\mathbf{Z}$ .

1984 When the covariance parameters  $\theta$  are known, an unbiased, consistent and  
1985 optimal estimate of  $\beta$ , given a vector of observations  $\zeta \in \mathbb{R}^n$ , is obtained by  
1986 minimizing the weighted least-squares functional

1987

$$1988 \quad \|\zeta - \mathbf{H}\beta\|_{\mathbf{C}_\theta^{-1}}^2 := (\zeta - \mathbf{H}\beta)^\top \mathbf{C}_\theta^{-1} (\zeta - \mathbf{H}\beta),$$

1989

1990 resulting in the estimate

1991

$$1992 \quad \hat{\beta} = (\mathbf{H}^\top \mathbf{C}_\theta^{-1} \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{C}_\theta^{-1} \zeta.$$

1993

1994 In *maximum likelihood (ML)* estimation, the joint probability density function  
1995 (31) is maximized for the given observation vector  $\xi = \zeta$  as a function of  
1996 the parameters  $\beta$  and  $\theta$ . To solve this nonlinear optimization problem one  
1997 usually *minimizes* the negative logarithm  $\ell(\zeta; \beta, \theta) := -\log p(\zeta; \beta, \theta)$  of the  
1998 likelihood given by

1999

$$2000 \quad \ell(\zeta; \beta, \theta) = \frac{1}{2} [n \log(2\pi) + \log \det \mathbf{C}_\theta + (\zeta - \mathbf{H}\beta)^\top \mathbf{C}_\theta^{-1} (\zeta - \mathbf{H}\beta)]. \quad (32)$$

2001

2002

2003 As is argued, e.g., in [Kitanidis \(1987\)](#), when random field hydrogeological  
2004 parameters are estimated based on data from a finite region where the sep-  
2005 aration distance of the measurements is of the same order as the correlation  
2006 length, the use of fitted means may introduce a bias in the estimation of the  
2007 covariance parameters, resulting typically in an underestimation of both the  
2008 variance and correlation length parameters. This bias is the result of strong  
2009 correlations in the observations, preventing the estimation error from entering  
2010 the asymptotic regime as more observations are added, since the number of  
2011 independent measurements does not increase due to these strong correlations.

2012 A remedy known as *restricted maximum likelihood estimation (RML)* is to  
2013 apply a transformation to the data which filters out the mean. In the case of  
2014 the linear model (8) for the mean, we consider the random vector  $\mathbf{Z}'$  obtained  
2015 by projecting  $\mathbf{Z}$  orthogonally onto the orthogonal complement of the range  
2016 of  $\mathbf{H}$ , hence removing any effect of the estimated regression coefficients  $\beta$   
2017 on the estimation of the covariance parameters. Indeed, if the columns of  
2018  $\mathbf{Q} \in \mathbb{R}^{n \times (n-k)}$  form an orthonormal basis of  $\text{range}(\mathbf{H})$ , then  $\mathbf{Q}^\top \mathbf{H} = \mathbf{O}$  and  
2019 therefore the random vector

2020

$$2021 \quad \mathbf{Z}' := \mathbf{Q}\mathbf{Q}^\top \mathbf{Z}$$

2022

2023

2024

has expectation

$$\mathbf{E}[\mathbf{Z}'] = \mathbf{E}[\mathbf{Q}\mathbf{Q}^\top(\mathbf{H}\boldsymbol{\beta} + \tilde{\mathbf{Z}})] = \mathbf{E}[\tilde{\mathbf{Z}}] = \mathbf{0}$$

regardless of the value of  $\boldsymbol{\beta}$ . Here  $\tilde{\mathbf{Z}}$  denotes the random vector obtained by restricting the residual random field  $\tilde{\mathbf{Z}}$  to the observation points. RML now maximizes the likelihood of the transformed random vector  $\mathbf{Z}'$ , which has an  $(n - k)$ -dimensional multivariate normal distribution with zero mean and covariance matrix  $\mathbf{Q}^\top \mathbf{C}_\theta \mathbf{Q} \in \mathbb{R}^{(n-k) \times (n-k)}$ .

2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051  
2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070