

Quasi-Monte Carlo finite element methods for elliptic PDEs with lognormal random coefficients

I. G. Graham¹, F. Y. Kuo², J. A. Nichols², R. Scheichl¹, Ch. Schwab³ and I. H. Sloan²

¹ Dept of Mathematical Sciences, University of Bath, Bath BA2 7AY UK
I.G.Graham@bath.ac.uk, R.Scheichl@bath.ac.uk

² School of Mathematics and Statistics, University of NSW, Sydney NSW 2052, Australia.
f.kuo@unsw.edu.au, i.sloan@unsw.edu.au

³ Seminar für Angewandte Mathematik, ETH Zürich, Rämistrasse 101, 8092 Zürich, Switzerland.
christoph.schwab@sam.math.ethz.ch

Abstract

In this paper we analyze the numerical approximation of diffusion problems over polyhedral domains in \mathbb{R}^d ($d = 1, 2, 3$), with diffusion coefficient $a(\mathbf{x}, \omega)$ given as a lognormal random field, i.e., $a(\mathbf{x}, \omega) = \exp(Z(\mathbf{x}, \omega))$ where \mathbf{x} is the spatial variable and $Z(\mathbf{x}, \cdot)$ is a Gaussian random field. The analysis presents particular challenges since the corresponding bilinear form is not uniformly bounded away from 0 or ∞ over all possible realizations of a . Focusing on the problem of computing the expected value of linear functionals of the solution of the diffusion problem, we give a rigorous error analysis for methods constructed from (i) standard continuous and piecewise linear finite element approximation in physical space; (ii) truncated Karhunen-Loève expansion for computing realizations of a (leading to a possibly high-dimensional parametrized deterministic diffusion problem); and (iii) lattice-based quasi-Monte Carlo (QMC) quadrature rules for computing integrals over parameter space which define the expected values. The paper contains novel error analysis which accounts for the effect of all three types of approximation. The QMC analysis is based on a recent result on randomly shifted lattice rules for high-dimensional integrals over the unbounded domain of Euclidean space, which shows that (under suitable conditions) the quadrature error decays with $\mathcal{O}(n^{-1+\delta})$ with respect to the number of quadrature points n , where $\delta > 0$ is arbitrarily small and where the implied constant in the asymptotic error bound is independent of the dimension of the domain of integration.

1 Introduction

In this paper we propose and analyze a class of numerical methods for the diffusion problem with random coefficient:

$$-\nabla \cdot (a(\mathbf{x}, \omega) \nabla u(\mathbf{x}, \omega)) = f(\mathbf{x}), \quad \text{for almost all } \omega \in \Omega \text{ and } \mathbf{x} \in D, \quad (1.1)$$

subject to the homogeneous Dirichlet condition $u = 0$ on ∂D . Here, D is a bounded (spatial) domain in \mathbb{R}^d and $(\Omega, \mathcal{A}, \mathbb{P})$ is a probability space (clarified below). We focus on the lognormal case, assuming that

$$a(\mathbf{x}, \omega) = a_*(\mathbf{x}) + a_0(\mathbf{x}) \exp(Z(\mathbf{x}, \omega)), \quad (1.2)$$

where Z is a zero-mean Gaussian random field, and a_*, a_0 are given functions that are continuous on \overline{D} with a_* non-negative and a_0 strictly positive on \overline{D} . This is (slightly) more general than the classical lognormal case where $a_* \equiv 0$. The latter is commonly used in many applications, for example in hydrology (see, e.g., [25, 26] and the references there). Models where realizations of a are not smooth (e.g. Hölder continuous, almost surely) are regularly of interest. Thus we need to study problem (1.1) in its weak form: Seek $u(\cdot, \omega) \in H_0^1(D)$ such that

$$\mathcal{A}(\omega; u, v) = \langle f, v \rangle, \quad \text{for all } v \in H_0^1(D) \text{ and for almost all } \omega \in \Omega, \quad (1.3)$$

where $\langle \cdot, \cdot \rangle$ denotes the duality pairing between $H^s(D)$ and $(H^s(D))'$, and

$$\mathcal{A}(\omega; w, v) := \int_D a(\mathbf{x}, \omega) \nabla w(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x}, \quad w, v \in H^1(D),$$

and we assume that $f \in (H_0^1(D))'$.

For each $\mathbf{x} \in D$, $Z(\mathbf{x}, \cdot)$ is a Gaussian random variable, and thus $0 < a(\mathbf{x}, \omega) < \infty$ for any $\omega \in \Omega$. However, for any $\varepsilon > 0$ we have $\mathbb{P}[a(\mathbf{x}, \cdot) > \varepsilon^{-1}] > 0$ so that problem (1.3) is not uniformly bounded over all possible realizations of a . If $a_*(\mathbf{x}) = 0$, then we also have $\mathbb{P}[a(\mathbf{x}, \cdot) < \varepsilon] > 0$ so that (1.3) is not uniformly elliptic either. This loss of ellipticity and boundedness is one of the main difficulties in the (numerical) analysis of (1.3).

Motivated by applications in uncertainty quantification, we will be interested in expected values of linear functionals of the solution of (1.3). That is, if $\mathcal{G} \in (H_0^1(D))'$, we will be interested in the expected value $\mathbb{E}[\mathcal{G}(u)]$ of the random variable $\mathcal{G}(u(\cdot, \omega))$. We will use sampling methods for the computation of $\mathbb{E}[\mathcal{G}(u)]$. That is, we will compute realizations of $a(\mathbf{x}, \omega)$, which yield realizations of $u(\mathbf{x}, \omega)$, via the solution of the elliptic problem (1.3), and from these we shall compute an approximation of $\mathbb{E}[\mathcal{G}(u)]$ by an appropriate averaging. However, in contrast to standard Monte Carlo (MC) methods, we will sample $a(\mathbf{x}, \omega)$ using quasi-Monte Carlo (QMC) methods. One principal aim of the paper is to prove mathematically that, under suitable assumptions, QMC methods are faster than MC methods for this class of problems. This fact has already been demonstrated computationally in [15] for a closely related method, but a convergence analysis was missing in that paper. (In [15] we used a different type of QMC points as well as different methods to sample from a and to discretize (1.1).)

In this work we assume that the Gaussian random field Z which appears in (1.2) is given in terms of a Karhunen-Loève expansion

$$Z(\mathbf{x}, \omega) = \sum_{j=1}^{\infty} \sqrt{\mu_j} \xi_j(\mathbf{x}) Y_j(\omega), \quad \mathbf{x} \in D. \quad (1.4)$$

Here, the sequence $\{Y_j\}_{j \geq 1}$ denotes the i.i.d. $\mathcal{N}(0, 1)$ distributed random variables Y_j , and the sequence $\{(\mu_j, \xi_j)\}_{j \geq 1}$ denotes the real eigenvalues and eigenfunctions of the *covariance* integral operator

$$(\mathcal{C}v)(\mathbf{x}) := \int_D c(\mathbf{x}, \mathbf{x}') v(\mathbf{x}') \, d\mathbf{x}', \quad \mathbf{x} \in D. \quad (1.5)$$

The kernel function $c(\cdot, \cdot)$ is assumed to be continuous on $\overline{D} \times \overline{D}$. It represents the covariance function of Z , i.e., $c(\mathbf{x}, \mathbf{x}') = \mathbb{E}[Z(\mathbf{x}, \cdot)Z(\mathbf{x}', \cdot)]$ for $\mathbf{x}, \mathbf{x}' \in D$, and thus $c(\mathbf{x}, \mathbf{x}') = c(\mathbf{x}', \mathbf{x})$. These facts imply that the covariance operator is a compact and self-adjoint operator from $L^2(D)$ to $L^2(D)$. Throughout, we shall assume that the covariance operator of Z is non-degenerate so that the sum in (1.4) is infinite, and that the eigenfunctions are orthonormal in $L^2(D)$, i.e., $\int_D \xi_i(\mathbf{x})\xi_j(\mathbf{x}) \, d\mathbf{x} = \delta_{ij}$. Consequently, the eigenvalues μ_j are real and positive and $\{\mu_j\}_{j \geq 1} \in \ell^1(\mathbb{N})$. We assume that μ_j are enumerated in non-increasing magnitude.

Our methods approximate $\mathbb{E}[\mathcal{G}(u)]$ in a three-stage process: The first stage comprises the approximation of (1.3) with fixed ω via the finite element method. Let $V_h \subset H_0^1(D)$ denote the space of piecewise linear functions on a family of conforming shape regular triangulations of D (parametrized by maximum mesh diameter h). As usual, the standard finite element solution of problem (1.3) is denoted $u_h(\mathbf{x}, \omega)$. In the second stage, to sample the random field a , the infinite sum in (1.4) is truncated to s terms. The resulting approximation to Z is substituted into (1.2) to obtain an approximate field a^s and the resulting finite element solution to (1.3) (with a replaced by a^s) is then denoted $u_h^s(\mathbf{x}, \omega)$. The corresponding approximation of $\mathbb{E}[\mathcal{G}(u)]$ is then taken to be the expected value of the random variable $\mathcal{G}(u_h^s(\cdot, \omega))$, written $\mathbb{E}[\mathcal{G}(u_h^s)]$.

In fact, since u_h^s is a random field derived from Z and thus from the s i.i.d. $\mathcal{N}(0, 1)$ random variables Y_1, \dots, Y_s , we have the concrete formula

$$\mathbb{E}[\mathcal{G}(u_h^s)] = \int_{\mathbb{R}^s} \mathcal{G}(u_h^s(\cdot, \mathbf{y})) \prod_{j=1}^s \phi(y_j) d\mathbf{y}, \quad (1.6)$$

where $\phi(y) = \exp(-y^2/2)/\sqrt{2\pi}$ is the Gaussian normal probability density.

The computation of the (possibly high dimensional) integral (1.6) by suitable quadrature rules leads us to the third stage of the approximation process. By introducing a change of variables $\mathbf{y} = \Phi_s^{-1}(\mathbf{v})$, where $\Phi_s^{-1}(\mathbf{v})$ denotes the inverse cumulative normal applied to each entry of $\mathbf{v} \in \mathbb{R}^s$, and writing $F(\mathbf{y}) = \mathcal{G}(u_h^s(\cdot, \mathbf{y}))$, we obtain

$$\mathbb{E}[\mathcal{G}(u_h^s)] = \mathbb{E}[F] = \int_{(0,1)^s} F(\Phi_s^{-1}(\mathbf{v})) d\mathbf{v}. \quad (1.7)$$

We will use “randomly shifted lattice rules” for (1.7), leading to approximations of the form

$$\mathcal{Q}_{s,h,n}(\Delta) := \frac{1}{n} \sum_{i=1}^n F\left(\Phi_s^{-1}\left(\text{frac}\left(\frac{i\mathbf{z}}{n} + \Delta\right)\right)\right), \quad i = 1, \dots, n, \quad (1.8)$$

where $\mathbf{z} \in \mathbb{N}^s$ is a *generating vector*, $\Delta \in [0, 1]^s$ is a *random shift* which is uniformly distributed over $[0, 1]^s$ and “frac” denotes the fractional part function, applied componentwise. A key point here is that a suitable generating vector \mathbf{z} can be very efficiently computed via the “component-by-component” procedure. We give a discussion of this in §4 and provide further references there. Applying this approximation to (1.7) yields a computable approximation to $\mathbb{E}[\mathcal{G}(u)]$ which depends on s, h and n and also on the random shift Δ , and is here denoted $\mathcal{Q}_{s,h,n}(\Delta)$.

The principal result of this paper is a bound for the root-mean-square error

$$e_{s,h,n} := \sqrt{\mathbb{E}^\Delta \left[(\mathbb{E}[\mathcal{G}(u)] - \mathcal{Q}_{s,h,n}(\Delta))^2 \right]},$$

where \mathbb{E}^Δ denotes expectation with respect to the random shift Δ . Our main result is Theorem 22, where estimates from §2 and §4 are collected to obtain the overall bound

$$e_{s,h,n} \leq C \left(h^{2\tau} + s^{-\chi} + n^{-r} \right). \quad (1.9)$$

(Here and throughout the paper C denotes a generic, positive constant independent of s, h and n .) The parameter τ depends on the regularity of (realizations of) a and on the smoothness of D , while the parameters χ and r depend on the asymptotics of the Karhunen–Loève eigenvalues and eigenvectors. Broadly speaking, faster decay yields stronger bounds in (1.9). (For example when the kernel c is analytic on $\overline{D} \times \overline{D}$, we have $\tau = 1$ and χ and r can be taken arbitrarily close to ∞ and 1, respectively.) The key novel result of this paper is that for a range of covariance functions, convergence for the QMC quadrature rule of order arbitrarily close to $\mathcal{O}(n^{-1})$ is attained with an asymptotic constant independent of the truncation dimension s in (1.6) of the Gaussian random field Z in (1.4). This should be compared with the $\mathcal{O}(n^{-1/2})$ rate which is attained by standard MC methods and cannot be improved in general. To our knowledge, MC and QMC are currently the only quadrature rules which afford asymptotic error bounds with constants that are independent of dimension. Since, due to the slow convergence of Karhunen–Loève expansions, high truncation dimensions s are often encountered in applications in hydrology, the results have considerable practical significance, as outlined in [15].

The starting point for our analysis is the observation that since the random diffusion coefficient $a(\mathbf{x}, \omega)$ in (1.1) and the random shift Δ in the QMC rule are statistically independent, we can write

$$e_{s,h,n}^2 = (\mathbb{E}[\mathcal{G}(u) - \mathcal{G}(u_h^s)])^2 + \mathbb{E}^\Delta [(\mathbb{E}[\mathcal{G}(u_h^s)] - \mathcal{Q}_{s,h,n}(\cdot))^2]. \quad (1.10)$$

This expression will form the basis for our error analysis. The first term on the right hand side of (1.10) can be bounded using the results in [3, 5, 36, 31, 14, 17] and we will outline this in §2. The second term on the right hand side of (1.10) is the QMC quadrature error; this is estimated in §4 and there particular emphasis is put on obtaining a rate of convergence which is close to $\mathcal{O}(n^{-1})$, with an asymptotic constant which is independent of dimension s . To apply QMC quadrature rules and to obtain error bounds that are independent of the truncation dimension s , it is necessary to reformulate the random problem (1.3) as a parametric, deterministic problem on an infinite-dimensional parameter space and to study both (i) the effect of dimension truncation to finite QMC integration dimension s and (ii) the regularity of the solution with respect to the parameters. This is the subject of §3. There we prove the well-posedness of these parametric problems (pointwise in a set of full Gaussian measure) and establish measurability and integrability of the parametric, deterministic solution.

For the QMC analysis we will in part follow the recent paper [19]. There, QMC integration was applied to a simplified PDE problem, in which the coefficient $a(\mathbf{x}, \mathbf{y})$ in (1.1) was assumed to be a linear function of \mathbf{y} , with \mathbf{y} a uniform random vector from the bounded domain $[-\frac{1}{2}, \frac{1}{2}]^N$. Here, as in [19], we introduce *weight parameters* $(\gamma_{\mathbf{u}})_{\mathbf{u} \subset \mathbb{N}}$ to control the relative importance of various subsets of the variables $\mathbf{y}_{\mathbf{u}} = \{y_j : j \in \mathbf{u}\}$. In [19] the PDE solutions as functions of \mathbf{y} were, after truncation to s dimensions, continuous functions on $[-\frac{1}{2}, \frac{1}{2}]^s$ with square integrable mixed first derivatives, a standard setting for QMC analysis. However, in the present work the function space setting we need is non-standard since, in general, the integrand $F \circ \Phi_s^{-1}$ in (1.7) may *not* have square integrable mixed first derivatives. This is due to the presence of the inverse cumulative distribution function Φ_s^{-1} , which is unbounded near the boundary of the unit cube. We will therefore consider a function space setting which uses a sequence of *weight functions* ψ_j to counteract the growth of the mixed first derivatives of $F(\mathbf{y}) = \mathcal{G}(u_h^s(\cdot, \mathbf{y}))$ as the components of \mathbf{y} go to $\pm\infty$, and we will make use of the corresponding error analysis for randomly shifted lattice rules from [21, 28].

The precise details on randomly shifted lattice rules and on the function space setting are given in §4. Note that the weight functions ψ_j and weight parameters $\gamma_{\mathbf{u}}$ appear in the definition of the norm - see (4.3), and they are initially free. A crucial step in the analysis is to establish bounds on some spatial norm (in \mathbf{x}) of the mixed first derivatives of $u_h^s(\mathbf{x}, \mathbf{y})$ with respect to the parametric variables in \mathbf{y} ; this is done in §3.2. Based on these bounds, we are then able to choose suitable weight functions ψ_j that allow $F(\mathbf{y}) = \mathcal{G}(u_h^s(\cdot, \mathbf{y}))$ to be in the function space. Then, with appropriately chosen weight parameters $\gamma_{\mathbf{u}}$, the QMC convergence rate can be arbitrarily close to order n^{-1} , with a constant that is independent of the truncation dimension s . In our analysis, as in [19], we choose the weight parameters $\gamma_{\mathbf{u}}$ so as to minimize a certain upper bound on the QMC quadrature error (i.e., the second term in (1.10)), and then show that the constant is indeed independent of the dimension s . Moreover, the resulting weight parameters $\gamma_{\mathbf{u}}$, again as in [19], are of a special form called POD weights (which stands for “product and order dependent weights”). This POD structure of the weight parameters $\gamma_{\mathbf{u}}$ allows in turn for a fast algorithm to construct lattice rules that are tailored to our problem (for details see [28]). In the penultimate section we summarize the overall error. In the final section we present numerical results.

2 Discretization and dimension truncation

The main purpose of this section is to summarise results on approximating solutions of (1.3) by the finite element method and on the effect of truncating the Karhunen–Loève expansion in (1.4) to s terms. As described in §1, the two approximations, with and without truncation, are denoted u_h^s and u_h , and we are concerned in this section with estimating the first term in (1.10). We give mainly a summary here, and we make reference to proofs in [3, 5, 36, 31, 14, 17], and the references therein, giving details only where necessary. However, a crucial difference between the present treatment and these references is that here, to estimate the first term in (1.10), we write

$$\mathcal{G}(u) - \mathcal{G}(u_h^s) = (\mathcal{G}(u) - \mathcal{G}(u_h)) + (\mathcal{G}(u_h) - \mathcal{G}(u_h^s)) \quad (2.1)$$

and estimate the expectation for each of these two terms separately. The final result is in Corollary 11. We remark that previous analyses introduced $\mathcal{G}(u^s)$ and estimated $\mathcal{G}(u) - \mathcal{G}(u^s)$ and $\mathcal{G}(u^s) - \mathcal{G}(u_h^s)$ instead. As we will see, the analysis here allows for a weaker assumption on the Karhunen–Loève eigenvalues and eigenfunctions than that made in [3, 5, 36]. Let us start with some notation.

We denote the usual scale of Sobolev spaces by $H^s(D)$, $s \geq 0$. Then $H_0^1(D)$ denotes, as usual, the subspace of functions in $H^1(D)$ with vanishing trace on ∂D . The space of continuous functions in \bar{D} is denoted $C^0(\bar{D})$. The Hölder space $C^{0,t}(\bar{D})$, with $0 < t \leq 1$, is the space of all functions $v \in C^0(\bar{D})$ for which $\|v\|_{C^{0,t}(\bar{D})} := \|v\|_{C^0(\bar{D})} + |v|_{C^{0,t}(\bar{D})} < \infty$, with seminorm

$$|v|_{C^{0,t}(\bar{D})} := \sup_{\mathbf{x}, \mathbf{x}' \in \bar{D}: \mathbf{x} \neq \mathbf{x}'} \frac{|w(\mathbf{x}) - w(\mathbf{x}')|}{|\mathbf{x} - \mathbf{x}'|^t},$$

where $|\mathbf{x}|$ denotes the Euclidean norm in \mathbb{R}^d . If $t = 0$, we adopt the convention $C^{0,0}(\bar{D}) = C^0(\bar{D})$.

We will also require spaces of Bochner integrable functions, that is, for any Banach space X with norm $\|\cdot\|_X$ and for $1 \leq q < \infty$, we denote by $L^q(\Omega, \mathbb{P}; X)$ the space of all strongly \mathbb{P} -measurable mappings v from (Ω, \mathcal{A}) to $(X, \mathcal{B}(X))$ (where $\mathcal{B}(X)$ denotes the Borel sigma algebra over X), for which the Bochner integral

$$\|v\|_{L^q(\Omega, \mathbb{P}; X)} = \begin{cases} \left(\int_{\Omega} \|v\|_X^q d\mathbb{P} \right)^{1/q}, & \text{for } 1 \leq q < \infty, \\ \text{esssup}_{\omega \in \Omega} \|v\|_X, & \text{for } q = \infty, \end{cases}$$

is finite. When there is no ambiguity about the measure, we shall denote this space by $L^q(\Omega; X)$. In the particular case $X = \mathbb{R}$, we shall simply write $L^q(\Omega)$ in place of $L^q(\Omega; \mathbb{R})$.

2.1 Spatial regularity

It is a classical result (see, e.g. [1]) that the regularity of the coefficient $a(\mathbf{x}, \omega)$ in (1.2), when considered as a function of $\mathbf{x} \in \bar{D}$, depends on the spectrum of the covariance operator \mathcal{C} in (1.5). We now specify several assumptions on the covariance of the Gaussian field Z in (1.2) which will be used in the paper. Firstly, we assume that the Gaussian process $Z(\cdot, \omega)$ is *stationary*, i.e., the covariance function $c(\mathbf{x}, \mathbf{x}')$ in (1.5) depends only on the single argument $\mathbf{x} - \mathbf{x}'$. In addition to stationarity, we assume for most of the paper that the random field Z (1.2) is *isotropic*, i.e.,

$$c(\mathbf{x}, \mathbf{x}') = \rho(|\mathbf{x} - \mathbf{x}'|), \quad \mathbf{x}, \mathbf{x}' \in \bar{D} \quad (2.2)$$

for some continuous function $\rho : [0, \infty) \rightarrow [0, \infty)$. Recall that $c(\cdot, \cdot)$ was also assumed to be continuous on $\bar{D} \times \bar{D}$. Thirdly we assume that there exist constants $C, \beta > 0$ such that $\mathbb{E}[(Z(\mathbf{x}, \cdot) - Z(\mathbf{x}', \cdot))]^2 \leq C|\mathbf{x} - \mathbf{x}'|^{2\beta}$, or equivalently that

$$|\rho(|\mathbf{x} - \mathbf{x}'|) - \rho(\mathbf{0})| \leq C|\mathbf{x} - \mathbf{x}'|^{2\beta}, \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \bar{D}. \quad (2.3)$$

From these assumptions it follows – e.g. from [1, §3.3] – that the Karhunen–Loève expansion (1.4) exists and converges almost surely.

For stationary, lognormal fields (1.2) with isotropic covariance function of the form (2.2), the following result (which is referred to sometimes as *Kolmogorov’s theorem*) gives sufficient conditions for the almost sure Hölder regularity of realizations of this field.

Proposition 1 *Assume that the Gaussian random field Z in (1.4) satisfies conditions (2.2) and (2.3) for some $\beta \in (0, 1]$. Then realizations of $Z(\cdot, \omega)$ are in $\mathbf{C}^{0,t}(\overline{D})$, \mathbb{P} -almost surely, for any $0 \leq t < \beta \leq 1$. If in addition a_* , $a_0 \in \mathbf{C}^{0,t}(\overline{D})$, then $a(\mathbf{x}, \omega)$ defined by (1.2) also satisfies $a(\cdot, \omega) \in \mathbf{C}^{0,t}(\overline{D})$ \mathbb{P} -almost surely.*

Proof. A proof for Z is given in [1, Theorem 8.3.2]. The result for $a(\mathbf{x}, \omega)$ follows from the assumptions on a_* and a_0 and the smoothness of $\exp(\cdot)$. \square

In §2.2 we will use Proposition 1 to infer Hölder regularity of realizations for a particular example which, in turn, will entail \mathbb{P} -almost sure Hölder regularity of solutions $u(\cdot, \omega)$ of (1.1) and thereby \mathbb{P} -almost sure rates of convergence of finite element discretizations of (1.1).

Consider the weak form of (1.1) as defined in (1.3). To prove well-posedness of this variational problem, we define, for \mathbb{P} -almost every $\omega \in \Omega$,

$$\check{a}(\omega) := \min_{\mathbf{x} \in \overline{D}} a(\mathbf{x}, \omega) \quad \text{and} \quad \hat{a}(\omega) := \max_{\mathbf{x} \in \overline{D}} a(\mathbf{x}, \omega). \quad (2.4)$$

Under the assumptions of Proposition 1, for almost all $\omega \in \Omega$, $Z(\cdot, \omega)$ is a continuous function on \overline{D} and hence attains its (finite) maximum on \overline{D} . Thus the quantities \check{a} and \hat{a} defined in (2.4) are \mathbb{P} -measurable and, hence, random variables which satisfy $\check{a}(\omega) > 0$ and $\hat{a}(\omega) < \infty$ \mathbb{P} -almost surely. Therefore, we may apply the Lax–Milgram Lemma “pathwise” to infer the existence of a unique solution $u(\cdot, \omega)$ of (1.1), for \mathbb{P} -almost every realization $a(\mathbf{x}, \omega)$ of the coefficient function in (1.1). The Lipschitz continuity of the data-to-solution correspondence for (1.1), between $\mathbf{C}^{0,t}(\overline{D})$ and $H_0^1(D)$, guarantees \mathbb{P} -measurability and hence, u is a random field taking values in the separable Hilbert space $H_0^1(D)$ on the probability space $(\Omega, \mathcal{A}, \mathbb{P})$. Finally, an application of Fernique’s Theorem (see, e.g. [7]) allows us to extend \mathbb{P} -almost sure bounds on $u(\cdot, \omega)$ to infer boundedness of $\|u\|_{L^q(\Omega; H_0^1(D))}$ for any $0 < q < \infty$ (we refer to [3, Section 2] for details).

From now on we adopt the notation $V = H_0^1(D)$.

Theorem 2 *Assume that $a_*, a_0 \in \mathbf{C}^0(\overline{D})$ in (1.2). Then, for all q in the range $1 \leq q < \infty$, $1/\check{a} \in L^q(\Omega)$ and $\hat{a} \in L^q(\Omega)$, and for every $f \in V'$ the problem (1.3) admits a unique solution $u \in L^q(\Omega; V)$ that satisfies*

$$\|u\|_{L^q(\Omega; V)} \leq \|1/\check{a}\|_{L^q(\Omega)} \|f\|_{V'}.$$

Proof. See [3, Prop. 2.3 & 2.4]. \square

As usual, to quantify the rate of convergence of finite element solutions of (1.3), additional regularity of the solution u is required. We formulate this as an assumption here, but we will indicate immediately that this assumption is indeed satisfied in a wide range of cases.

Assumption A1 There exists some $\tau > 0$ such that $u \in L^q(\Omega; H^{1+\tau}(D))$, for all $1 \leq q < \infty$.

The following theorem from [36, §5] (see also [5]) characterizes τ in the 2D polygonal case. It depends on the \mathbb{P} -almost sure Hölder regularity¹ of realizations of a and on the largest interior angle of ∂D .

¹The assumption on a made in [5, 36] can be weakened in the boundary case $t = 1$ from $a(\cdot, \omega) \in \mathbf{C}^1(\overline{D})$ to $a(\cdot, \omega) \in \mathbf{C}^{0,1}(\overline{D})$, since this implies $\nabla a(\cdot, \omega) \in L^\infty(D)$ and then the crucial Lemma A.2 in [5] can be replaced by the simple inequality $\|bv\|_{H^1(D)} \leq C(\|\nabla b\|_{L^\infty(D)} \|v\|_{L^2(D)} + \|b\|_{L^\infty(D)} \|v\|_{H^1(D)})$, for all $b \in \mathbf{C}^{0,1}(\overline{D})$ and $v \in H^1(D)$.

Theorem 3 *Let D be a polygon in \mathbb{R}^2 such that the largest interior angle θ_{max} of all the corners is in $(0, 2\pi)$, and suppose that $a(\cdot, \omega) \in \mathbf{C}^{0,t}(\overline{D})$ \mathbb{P} -almost surely. Then Assumption A1 holds for any $\tau < \min(t, \frac{\pi}{\theta_{max}})$. If D is convex, i.e. $\theta_{max} \leq \pi$ and if $a(\cdot, \omega) \in \mathbf{C}^{0,1}(\overline{D})$ almost surely, then Assumption A1 also holds for $\tau = 1$.*

Similar results hold also in three dimensions (see [36] for details).

2.2 Matérn Covariances

The foregoing abstract conditions are satisfied by Gaussian random fields Z with Matérn covariances which are commonly used in practice. Isotropic Matérn covariance functions (2.2) are given by

$$\rho(r) = \rho_\nu(r) := \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} (r/\tilde{\lambda})^\nu K_\nu(r/\tilde{\lambda}), \quad (2.5)$$

with $\tilde{\lambda} = \lambda_C/(2\sqrt{\nu})$. Here Γ is the gamma function and K_ν is the modified Bessel function of the second kind. The parameter $\nu > 1/2$ is a smoothness parameter, σ^2 is the variance and λ_C is a length scale parameter. Using the asymptotics of the modified Bessel function it is possible to show that (2.3) holds for this covariance with $\beta = \nu$, when $\nu \in (1/2, 1)$, and so it follows that Proposition 1 holds in this case with $0 < t < \nu$.

Remark 4 By increasing the parameter ν in (2.5) it appears possible to generate random fields whose sample paths have higher Hölder regularity. In particular when $\nu \in (1, 2)$, the second derivative

$$\left(\frac{\partial^2}{\partial x_i \partial x'_i} \right) \rho_\nu(|\mathbf{x} - \mathbf{x}'|) \quad (2.6)$$

exists and is finite at $\mathbf{x}' = \mathbf{x}$. By [1, Theorem 2.2.2] it follows that $\partial Z/\partial x_i$ exists (in the mean square sense), and is itself a random field with covariance given by (2.6). Again, by examining the asymptotics of K_ν , it can be shown that condition (2.3) holds for Z replaced by $\partial Z/\partial x_i$, with $\beta = \nu - 1$. Hence $\partial Z/\partial x_i$ is \mathbb{P} -almost surely Hölder continuous with exponent t where $0 < t < \nu - 1 < 1$. In particular, this implies that, for any $\nu > 1$, the random field Z is continuously differentiable in quadratic mean with respect to \mathbb{P} and so is $a(\cdot, \omega)$ in (1.2) provided that a_* , a_0 in (1.2) are continuously differentiable in \overline{D} .

It is instructive to consider the limiting cases $\nu \rightarrow 1/2$ and $\nu \rightarrow \infty$ separately. By evaluating the Matérn class at $\nu = 1/2$, we obtain the *exponential covariance*

$$\rho_{1/2}(r) = \sigma^2 \exp(-r/\tilde{\lambda}), \quad (2.7)$$

with $\tilde{\lambda} = \lambda_C/\sqrt{2}$, and in this case Proposition 1 holds with $0 < t < 1/2$.

In the theory of isotropic and stationary covariances, a key role is played by the Fourier transform. In the case of the Matérn covariance, it is given by

$$\widehat{\rho}_\nu(\boldsymbol{\xi}) = \left(\frac{1}{2\pi} \right)^d \int_{\mathbb{R}^d} \exp(-i\boldsymbol{\xi} \cdot \mathbf{x}) \rho_\nu(|\mathbf{x}|) d\mathbf{x} =: \phi_\nu(|\boldsymbol{\xi}|),$$

where

$$\phi_\nu(r) = \sigma^2 \left(\frac{1}{\pi} \right)^{d/2} \left[\frac{\Gamma(d/2 + \nu)}{\Gamma(\nu)} \right] \left[\frac{\alpha^{2\nu}}{(r^2 + \alpha^2)^{\nu+d/2}} \right]. \quad (2.8)$$

The limit of this as $\nu \rightarrow \infty$ coincides with the Fourier transform of the *Gaussian covariance*

$$\rho_\infty(r) = \sigma^2 \exp(-r^2/\lambda_C^2). \quad (2.9)$$

In this case, the covariance function $c(\cdot, \cdot)$ is analytic in $\overline{D} \times \overline{D}$ and so the samples $a(\cdot, \omega)$ are also analytic, for \mathbb{P} -almost every $\omega \in \Omega$.

As stated in the introduction, the speed of decay of the Karhunen–Loève eigenvalues $\{\mu_j\}_{j \geq 1}$ plays an important role in our error analysis. In the case of the Matérn class this can be determined using classical results on the analysis of integral operators with difference type kernels by H. Widom [39]. If μ_j denotes the j th eigenvalue of the integral operator (1.5), where c is given by (2.2) and $\rho = \rho_\nu$ (given in (2.5)), then Widom’s results imply that μ_j^d has the same asymptotic rate of decay as $j \rightarrow \infty$ as does the Fourier transform $\phi_\nu(r)$ as $r \rightarrow \infty$. This leads to the following corollary to the Widom theory (see, e.g., [23]).

Corollary 5 *There exists $C > 0$ such that the j th largest eigenvalue of the Matérn covariance operator satisfies, for every $j \geq 1$, the bound*

$$\mu_j \leq Cj^{-(1+2\nu/d)}.$$

When $d = 1$ and $\nu = 1/2$ (the 1D exponential case), we get the decay $\mathcal{O}(j^{-2})$ which corresponds to the decay in the analytic formula for the Karhunen–Loève eigenvalues for this problem given in [13]. For (2.9), i.e., when $\nu = \infty$ in (2.5), the μ_j decays at least exponentially (see e.g. [32]).

2.3 Finite element discretization error

To discretize (1.3) in the physical domain D we consider now finite element approximations with standard, continuous, piecewise linear finite elements. We denote by $\{\mathcal{T}_h\}_{h>0}$ a shape-regular family of simplicial triangulations of the domain D , parametrized by the mesh width $h := \max_{T \in \mathcal{T}_h} \text{diam}(T)$. Associated with each triangulation \mathcal{T}_h we define the space $V_h \subset V$ of piecewise linear, continuous functions on this mesh, which vanish on ∂D . For \mathbb{P} -almost every $\omega \in \Omega$, we denote by $u_h(\omega, \cdot) \in V_h$ the solution of

$$\mathcal{A}(\omega; u_h(\cdot, \omega), v_h) = \langle f, v_h \rangle, \quad \text{for all } v_h \in V_h. \quad (2.10)$$

As in Theorem 2, for every h and for \mathbb{P} -almost every realization $a(\cdot, \omega)$, the finite element solution $u_h(\cdot, \omega) \in V_h$ exists, is unique and (like the exact solution $u(\cdot, \omega)$) satisfies the a priori bound

$$\|u_h\|_{L^q(\Omega; V)} \leq \|1/\check{a}\|_{L^q(\Omega)} \|f\|_{V'}, \quad \text{for all } 1 \leq q < \infty, \quad (2.11)$$

We are now in a position to bound the first term in the overall error bound (2.1) for our method. A proof can be found in [36].

Theorem 6 *Let $0 < t < 1$ be as in Proposition 1 and let Assumption A1 hold for some $0 < \tau < t$. Suppose $\mathcal{G}(\cdot)$ is a continuous linear functional on $H^{1-\tau}(D)$, i.e. there exists a constant $C_{\mathcal{G}}$ such that $|\mathcal{G}(v)| \leq C_{\mathcal{G}} \|v\|_{H^{1-\tau}(D)}$ for all $v \in H^{1-\tau}(D)$. Then*

$$|\mathbb{E}[\mathcal{G}(u) - \mathcal{G}(u_h)]| \leq Ch^{2\tau}. \quad (2.12)$$

If $a(\cdot, \omega) \in C^{0,1}(\overline{D})$ \mathbb{P} -almost surely and if Assumption A1 holds for $\tau = 1$, then (2.12) holds with $\tau = 1$.

Remark 7 Theorem 6 can be generalized to the case where the functional \mathcal{G} is random, i.e. for each $\omega \in \Omega$, $\mathcal{G} = \mathcal{G}_\omega \in (H^{1-\tau}(D))'$, where $|\mathcal{G}_\omega(v)| \leq C_{\mathcal{G}}(\omega) \|v\|_{H^{1-\tau}(D)}$. Then, (2.12) still holds, provided $C_{\mathcal{G}} \in L^q(\Omega)$ for some $q > 1$ ([36, §3]).

2.4 Dimension truncation error

In practice, in order to use (1.4) to numerically sample the Gaussian random field Z , it is of course necessary to truncate the infinite series expansion (1.4) and to control the resulting error. To analyze the truncation error we need to make some assumptions on the regularity and the decay of the Karhunen–Loève eigenvalues and eigenfunctions (μ_j, ξ_j) as $j \rightarrow \infty$. These assumptions can be verified rigorously for particular covariances, such as for isotropic Matérn covariances. As mentioned above, we follow closely [3, 5, 36] except that (via the error decomposition (2.1)), we need only consider the effect of dimension truncation for the finite element solution (see the last term in (2.1)), and not for u itself.

Recalling (1.2) and (1.4), the approximation of a obtained by the dimensionally truncated Karhunen–Loève expansion of Z is

$$a^s(\mathbf{x}, \omega) := a_*(\mathbf{x}) + a_0(\mathbf{x}) \exp\left(\sum_{j=1}^s \sqrt{\mu_j} \xi_j(\mathbf{x}) Y_j(\omega)\right), \quad \text{for some } s \in \mathbb{N}. \quad (2.13)$$

The number of terms s is the dimension of the parameter domain for QMC integration in §4.

For \mathbb{P} -almost every $\omega \in \Omega$, we can now define $u_h^s(\cdot, \omega) \in V_h$ to be the solution of the dimensionally truncated, discretized boundary value problem

$$\mathcal{A}^s(\omega; u_h^s(\cdot, \omega), v_h) = \langle f, v_h \rangle, \quad \text{for all } v_h \in V_h, \quad (2.14)$$

where

$$\mathcal{A}^s(\omega; w, v) := \int_D a^s(\mathbf{x}, \omega) \nabla w(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x}, \quad \text{for any } v, w \in V.$$

For simplicity, we work here under the assumption that, for any $v_h, w_h \in V_h$, we evaluate the integrals in $\mathcal{A}^s(\omega; w_h, v_h)$ exactly. It is possible to also include quadrature errors in the analysis (see [5, §3.3] for details). Existence and uniqueness for $u_h^s(\cdot, \omega)$ for \mathbb{P} -almost every $\omega \in \Omega$ follows again by the Lax–Milgram Lemma.

To obtain a bound on $|\mathbb{E}[\mathcal{G}(u_h) - \mathcal{G}(u_h^s)]|$ we apply the truncation error analysis in [3, 5]. It is well known that the integral operator \mathcal{C} maps $L^2(D)$ to $L^\infty(D)$ and so $\xi_j \in L^\infty(D)$ for all $j \geq 1$. In what follows, we draw on some results in [3]. Therefore, we make the following assumptions on the Karhunen–Loève eigenfunctions ξ_j and eigenvalues μ_j .

Assumption A2 (a) There exist $C > 0$ and $\Theta > 1$ such that $\mu_j \leq Cj^{-\Theta}$ for $j \geq 1$.

(b) The Karhunen–Loève eigenfunctions ξ_j are continuously differentiable and there exist $C > 0$ and $\varepsilon \in [0, \frac{\Theta-1}{2\Theta})$ such that $\|\xi_j\|_{C^0(\bar{D})} + \mu_j \|\nabla \xi_j\|_{C^0(\bar{D})} \leq C\mu_j^{-\varepsilon}$ for $j \geq 1$.

Theorem 8 *Let Assumption A2 hold. Then $\|1/\tilde{a}^s\|_{L^q(\Omega)}$ is bounded independently of s , for all $1 \leq q < \infty$. Suppose further that $\mathcal{G} \in V'$. Then*

$$|\mathbb{E}[\mathcal{G}(u_h) - \mathcal{G}(u_h^s)]| \leq C_\chi s^{-\chi}, \quad \text{for all } 0 < \chi < (\tfrac{1}{2} - \varepsilon)\Theta - \tfrac{1}{2}. \quad (2.15)$$

Proof. Note that Assumption A2 implies

$$\begin{aligned} \sum_{j \geq 1} \mu_j \|\xi_j\|_{C^0(\bar{D})}^2 &\leq C \sum_{j \geq 1} \mu_j^{1-2\varepsilon} \leq C \sum_{j \geq 1} j^{-(1-2\varepsilon)\Theta} < \infty, \quad \text{and} \\ \sum_{j \geq 1} \mu_j \|\xi_j\|_{C^0(\bar{D})}^{2(1-\alpha)} \|\nabla \xi_j\|_{C^0(\bar{D})}^{2\alpha} &\leq C \sum_{j \geq 1} \mu_j^{1-2\varepsilon-2\alpha} \leq C \sum_{j \geq 1} j^{-(1-2\varepsilon-2\alpha)\Theta} < \infty, \end{aligned} \quad (2.16)$$

for arbitrary $\alpha \in (0, \frac{\Theta-1}{2\Theta} - \varepsilon)$. Thus Assumption 3.1 of [3] holds. The required result that $\|1/\tilde{a}^s\|_{L^q(\Omega)}$ is bounded for $0 < q < \infty$, independently of s , then follows from [3, Proposition 3.10]. Moreover, (2.16) implies

$$\max \left(\sum_{j>s} \mu_j \|\xi_j\|_{C^0(\bar{D})}^2, \sum_{j>s} \mu_j \|\xi_j\|_{C^0(\bar{D})}^{2(1-\alpha)} \|\nabla \xi_j\|_{C^0(\bar{D})}^{2\alpha} \right) \leq C \sum_{j>s} j^{-(1-2\varepsilon-2\alpha)\Theta} \leq C s^{-(1-2\varepsilon-2\alpha)\Theta+1}$$

and hence Assumption 3.5 of [3] (which requires that the quantity on the left hand side, raised to the p_0 th power, be summable with respect to s for some $p_0 > 0$) holds for arbitrary $p_0 > ((1-2\varepsilon-2\alpha)\Theta - 1)^{-1}$. In turn this allows us to use [3, Theorem 4.2] to obtain

$$\|u - u^s\|_{L^q(\Omega;V)} \leq C_{q,\chi} s^{-\chi}, \quad (2.17)$$

where $\chi = (1/2 - \varepsilon - \alpha)\Theta - 1/2$, and u^s is the solution of the dimensionally truncated problem

$$\mathcal{A}^s(\omega; u^s(\cdot, \omega), v) = \langle f, v \rangle, \quad \text{for all } v \in V.$$

Finally, since the finite element solution u_h satisfies the same a priori bound (2.11) as the exact solution u (in Theorem 2) and since the right hand sides in (2.10) and in (2.14) are identical, the bound (2.17) holds also for $\|u_h - u_h^s\|_{L^q(\Omega;V)}$ with constant $C_{q,\chi} > 0$ being independent of s and h . This follows immediately from the proofs of [3, Theorems 4.1 and 4.2], since all the identities and bounds involving $u - u^s$ there, hold equally for $u_h - u_h^s$. The final result (2.15) then follows upon taking $q = 1$ and from the fact that $\mathcal{G} \in V'$. \square

Note that it follows from [12, 32] that Assumption A2 is satisfied in the case of the Gaussian covariance kernel for any $\Theta > 1$ with $\varepsilon = 0$. It is in fact possible in that case to obtain (see [3, §7.2] and [12, 32]), instead of (2.15) the exponential convergence estimate:

$$|\mathbb{E}[\mathcal{G}(u_h) - \mathcal{G}(u_h^s)]| \leq C \exp(-c_1 s^{1/d}), \quad \text{for some } c_1 > 0.$$

Another example for which we know that Assumption A2 is satisfied with $\Theta = 2$ and $\varepsilon = 0$, is the exponential covariance kernel (2.7) in one dimension or on rectangular domains, when the Euclidean norm $|\mathbf{x}| = \|\mathbf{x}\|_2$ in the definition of the kernel c in (2.2) is replaced by the 1-norm $\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i|$, since in that case the Karhunen–Loève eigenvalues and eigenfunctions are known explicitly [3, §7.1]. Here, the rate of convergence in Theorem 8 is $\chi < \frac{1}{2}$.

For the Matérn class we know from Corollary 5 that Assumption A2(a) is satisfied with $\Theta = 1 + 2\nu/d$. The numerical experiments in §6 suggest that, at least in one dimension, Assumption A2(b) is also satisfied for all $\nu > 1/2$ with $\varepsilon = 0$. This would lead to a convergence rate of $\chi < \frac{\nu}{d}$. We are so far only able to prove this rigorously for $\varepsilon > \frac{1}{2\Theta}$ and $\chi < \frac{\nu}{d} - \frac{1}{2}$ using the Sobolev embedding theorem and an interpolation argument, as the following result shows.

Proposition 9 *Consider the Matérn covariance given by (2.2) and (2.5) with $\nu > \frac{d}{2}$. Then Assumption A2 holds with $\Theta = 1 + \frac{2\nu}{d}$ and $\varepsilon \in (\frac{1}{2\Theta}, \frac{\Theta-1}{2\Theta})$, and the truncation error bound (2.15) in Theorem 8 holds for all $0 < \chi < \frac{\nu}{d} - \frac{1}{2}$.*

Proof. It follows immediately from Corollary 5 that Assumption A2(a) holds with $\Theta = 1 + 2\nu/d$. To show A2(b), note that $\mathcal{C}v(\mathbf{x}) = \int_{\mathbb{R}^d} \rho_\nu(|\mathbf{x} - \mathbf{x}'|) \tilde{v}(\mathbf{x}') d\mathbf{x}'$, where \tilde{v} is the extension of v by zero. Considering this for all $\mathbf{x} \in \mathbb{R}^d$, it follows from (2.8) and by the definition of fractional Sobolev norms via Fourier transforms, the convolution theorem and Plancherel's theorem (e.g. [24, p.75]), that \mathcal{C} is bounded from $L^2(D)$ to $H^r(D)$, provided $r < d + 2\nu = d\Theta$. Moreover,

$$\|\xi_j\|_{H^r(D)} = \frac{1}{\mu_j} \|\mathcal{C}\xi_j\|_{H^r(D)} \leq \frac{1}{\mu_j} \|\mathcal{C}\|_{L^2 \rightarrow H^r} \|\xi_j\|_{L^2(D)}.$$

Now, using the fact that $\|\xi_j\|_{L^2(D)} = 1$ and interpolation between $L^2(D)$ and $H^r(D)$ we get

$$\|\xi_j\|_{H^{\tilde{r}}(D)} \leq \mu_j^{-\tilde{r}/r} \|\mathcal{C}\|_{L^2 \rightarrow H^r}^{\tilde{r}/r}, \quad \text{for } 0 \leq \tilde{r} \leq r \text{ and for every } j \geq 1.$$

By choosing $\tilde{r} > d/2$, it follows from the Sobolev embedding theorem that $\|\xi_j\|_{C^0(\bar{D})} \leq C\mu_j^{-\varepsilon}$, for any $\varepsilon > \frac{1}{2\Theta}$. Moreover, noting that $d + 2\nu > d/2 + 1$, we can also choose \tilde{r} in the range $d/2 + 1 < \tilde{r} < r < d + 2\nu$, allowing us to infer that $\mu_j \|\nabla \xi_j\|_{C^0(\bar{D})} \leq \mu_j \|\xi_j\|_{H^{\tilde{r}}(D)} \leq \mu_j^{1-\tilde{r}/r}$, which is bounded as $j \rightarrow \infty$. These two estimates allow us to conclude that A2(b) holds for all $\varepsilon \in (\frac{1}{2\Theta}, \frac{\Theta-1}{2\Theta})$. This interval is guaranteed to be non-empty by the requirement that $\Theta > 2$ which is equivalent to $\nu > \frac{d}{2}$.

To establish the final part, note that for any $\delta > 0$ sufficiently small, there exists an ε such that $(\frac{1}{2} - \varepsilon)\Theta - \frac{1}{2} = (\frac{1}{2} - \frac{1}{2\Theta})\Theta - \frac{1}{2} - \delta = \frac{\nu}{d} - \frac{1}{2} - \delta$, which is positive when $\nu > \frac{d}{2}$. \square

Remark 10 Assumption A2 is weaker than the assumptions in [3, 5, 36] because (due to the new splitting (2.1)) we have separated the truncation error analysis fully from the discretization error analysis. It is in fact possible to further weaken Assumption A2(b) by requiring only Hölder continuity with index $0 < t < 1$ (as in Proposition 1) for the eigenfunctions (see also [4]). Since this does not improve the result in Proposition 9 and would only further complicate the presentation, we did not do this. The results of [4] suggest that under additional assumptions it should be possible to strengthen the $s^{-\chi}$ term in the estimate of Theorem 8 to $s^{-2\chi}$. This rate can be observed numerically (cf. [4, 36]).

Combining Theorems 6 and 8, we obtain the following estimate of the first term in (1.10).

Corollary 11 *Under Assumptions A1 and A2 and with $\tau, \chi > 0$ as defined in Theorems 6 and 8 we have*

$$|\mathbb{E}[\mathcal{G}(u) - \mathcal{G}(u_h^s)]| \leq C (h^{2\tau} + s^{-\chi}).$$

This result extends straightforwardly to Fréchet differentiable nonlinear functionals (cf. [36]).

3 Parametric deterministic problem

In this section, we prepare the analysis of the second term in (1.10), which is the error in approximating the expectation $\mathbb{E}[\mathcal{G}(u_h^s)]$ by a suitable randomly shifted QMC quadrature approximation. Whereas standard Monte Carlo sampling of the random diffusion problem (1.1) requires merely a ‘sampler’ for the random permeability $a(\mathbf{x}, \omega)$, QMC quadrature requires “integration coordinates” (see (1.6) and (1.7)). We introduce these coordinates here via the Karhunen–Loève expansion (1.4) of the Gaussian random field Z .

3.1 Parametric, deterministic variational formulation

By (1.4) and (1.2), the coefficient $a(\mathbf{x}, \omega)$ of the problem (1.1) is parametrized by a vector $\mathbf{y}(\omega) = (Y_1(\omega), Y_2(\omega), \dots) \in \mathbb{R}^{\mathbb{N}}$ of i.i.d. random variables $Y_j \sim \mathcal{N}(0, 1)$. The law of the random vectors \mathbf{y} is defined on the product probability space $(\mathbb{R}^{\mathbb{N}}, \mathcal{B}(\mathbb{R}^{\mathbb{N}}), \bar{\mu}_G)$, where $\mathcal{B}(\mathbb{R}^{\mathbb{N}})$ denotes the sigma-algebra generated by the cylinder sets, i.e., by (countable) products of intervals $I \in \mathcal{B}(\mathbb{R}^1)$, and $\bar{\mu}_G$ is the product Gaussian measure (see, e.g., [2])

$$\bar{\mu}_G = \bigotimes_{j=1}^{\infty} \mathcal{N}(0, 1). \quad (3.1)$$

In our subsequent QMC error analysis, we will need the following assumption.

Assumption A3 The sequence $\mathbf{b} = \{b_j\}_{j \geq 1}$ defined by

$$b_j := \sqrt{\mu_j} \|\xi_j\|_{C^0(\bar{D})}, \quad j \geq 1, \quad (3.2)$$

satisfies

$$\sum_{j \geq 1} b_j^p < \infty, \quad \text{for some } p \in (0, 1].$$

Assumption A3 follows immediately from Assumption A2 under the significantly more restrictive condition $\Theta > \frac{2}{1-2\varepsilon} p^{-1}$. Note however, that in the Matérn case, due to Corollary 5 we can again (as in Proposition 9) show theoretically that Assumption A3 holds for all $\nu > d$. If $\|\xi_j\|_{C^0(\bar{D})}$ is uniformly bounded, as seems to be the case numerically in the Matérn case for $d = 1$ (see Figure 1), then Assumption A3 holds for all $\nu > d/2$.

Under Assumption A3 we can now define the *admissible* parameter set

$$U_{\mathbf{b}} := \left\{ \mathbf{y} \in \mathbb{R}^{\mathbb{N}} : \sum_{j=1}^{\infty} b_j |y_j| < \infty \right\} \subset \mathbb{R}^{\mathbb{N}}.$$

The set $U_{\mathbf{b}} \subset \mathbb{R}^{\mathbb{N}}$ is not a countable product of subsets of \mathbb{R} , but, as we show in the following lemma (cf. [31, Lemma 2.28]), it is $\bar{\mu}_G$ -measurable and of full (Gaussian) measure.

Lemma 12 *If Assumption A3 holds for some $0 < p < 1$ then $U_{\mathbf{b}} \in \mathcal{B}(\mathbb{R}^{\mathbb{N}})$ and $\bar{\mu}_G(U_{\mathbf{b}}) = 1$.*

With a slight abuse of notation, we identify the *stochastic* coefficient $a(\mathbf{x}, \omega)$ defined in (1.2) with its parametric representation $a(\mathbf{x}, \mathbf{y}(\omega))$, that is, for each $\mathbf{x} \in D$ and $\mathbf{y} \in U_{\mathbf{b}}$, we define the *deterministic, parametric* coefficient and its s -term truncation (cf. (2.13)) by

$$a(\mathbf{x}, \mathbf{y}) = a_*(\mathbf{x}) + a_0(\mathbf{x}) \exp \left(\sum_{j=1}^{\infty} \sqrt{\mu_j} \xi_j(\mathbf{x}) y_j \right), \quad (3.3)$$

$$a^s(\mathbf{x}, \mathbf{y}) = a_*(\mathbf{x}) + a_0(\mathbf{x}) \exp \left(\sum_{j=1}^s \sqrt{\mu_j} \xi_j(\mathbf{x}) y_j \right). \quad (3.4)$$

Note that $a^s(\mathbf{x}, \mathbf{y})$ can be considered as the exact coefficient $a(\mathbf{x}, \mathbf{y})$ evaluated at the particular vector $\mathbf{y} = (y_1, \dots, y_s, 0, 0, \dots)$. More generally, denoting by $\mathbf{u} \subset \mathbb{N}$ any set of “active” coordinates, we denote by $(\mathbf{y}_{\mathbf{u}}; \mathbf{0})$ the vectors $\mathbf{y} \in U_{\mathbf{b}}$ with the constraint that $y_j = 0$ if $j \notin \mathbf{u}$.

The series in (3.3) converges in $C^0(\bar{D})$ for all $\mathbf{y} \in U_{\mathbf{b}} \subset \mathbb{R}^{\mathbb{N}}$. Thus, analogously to (2.4), setting

$$\check{a}(\mathbf{y}) := \min_{\mathbf{x} \in \bar{D}} a(\mathbf{x}, \mathbf{y}) \quad \text{and} \quad \hat{a}(\mathbf{y}) := \max_{\mathbf{x} \in \bar{D}} a(\mathbf{x}, \mathbf{y}),$$

and recalling that $a_*(\mathbf{x}) \geq 0$ and $a_0(\mathbf{x}) > 0$, for all $\mathbf{x} \in D$, we have

$$0 < \check{a}(\mathbf{y}) \leq a(\mathbf{x}, \mathbf{y}) \leq \hat{a}(\mathbf{y}) < \infty, \quad \text{for all } \mathbf{x} \in D \text{ and } \mathbf{y} \in U_{\mathbf{b}}. \quad (3.5)$$

Due to (3.5) and Lemma 12, we can consider $U_{\mathbf{b}}$ to be the parameter space instead of $\mathbb{R}^{\mathbb{N}}$. Even though $U_{\mathbf{b}}$ is not a product domain, we can define product measures such as the Gaussian measure $\bar{\mu}_G$ on $U_{\mathbf{b}}$ by restriction. For each $\mathbf{y} \in U_{\mathbf{b}}$, we can now consider the following *parametric, deterministic variational formulation* of the lognormal diffusion problem (1.1): Find $u(\cdot, \mathbf{y}) \in V$ such that

$$\mathcal{A}(\mathbf{y}; u, v) = \langle f, v \rangle, \quad \text{for all } v \in V, \quad (3.6)$$

where the *parametric, deterministic bilinear form* is defined as

$$\mathcal{A}(\mathbf{y}; w, v) := \int_D a(\mathbf{x}, \mathbf{y}) \nabla w(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x}, \quad w, v \in V.$$

The problem (3.6) is equivalent to its stochastic counterpart (1.3) (after substituting Gaussian random variables Y_j instead of parameter values y_j). For every $\mathbf{y} \in U_{\mathbf{b}}$, we can also define the finite element solution $u_h^s(\cdot, \mathbf{y}) \in V_h \subset V$ for the *s-term truncated parametric, deterministic* problem as the solution of

$$\mathcal{A}^s(\mathbf{y}; u_h^s, v_h) = \langle f, v_h \rangle, \quad \text{for all } v_h \in V_h, \quad (3.7)$$

where

$$\mathcal{A}^s(\mathbf{y}; w, v) := \int_D a^s(\mathbf{x}, \mathbf{y}) \nabla w(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x}, \quad w, v \in V.$$

Our aim now is the analysis of QMC approximations of integrals with respect to $\mathbf{y} \in U_{\mathbf{b}}$ of $\mathcal{G}(u_h^s(\cdot, \mathbf{y}))$, where \mathcal{G} is a continuous, linear functional on V . These integrals require pointwise evaluation of $u_h^s(\cdot, \mathbf{y})$, and so to develop the QMC error analysis we first establish that, for every $\mathbf{y} \in U_{\mathbf{b}}$, the two parametric, deterministic problems (3.6) and (3.7) admit unique weak solutions. Due to (3.5), this is again a direct consequence of the Lax-Milgram Lemma.

Theorem 13 *For every $\mathbf{y} \in U_{\mathbf{b}}$, $s \in \mathbb{N}$ and $h > 0$, the parametric, deterministic problems (3.6) and (3.7) admit unique solutions $u(\cdot, \mathbf{y}) \in V$ and $u_h^s(\cdot, \mathbf{y}) \in V_h$. Moreover,*

$$\|u(\cdot, \mathbf{y})\|_V \leq \frac{1}{\check{a}(\mathbf{y})} \|f\|_{V'}, \quad \text{for all } \mathbf{y} \in U_{\mathbf{b}},$$

with the same bound holding also for $\|u_h^s(\cdot, \mathbf{y})\|_V$.

3.2 Regularity with respect to the parametric variables

To estimate the second term in (1.10), it is crucial to bound the mixed first derivatives of $u_h^s(\cdot, \mathbf{y})$ with respect to \mathbf{y} . Here we state and prove a more general result which gives bounds also for higher order mixed derivatives. In fact we prove the result for $u(\cdot, \mathbf{y})$ and explain subsequently why the argument also applies (with constants that are independent of h and of s) to $u_h^s(\cdot, \mathbf{y})$. We introduce some additional notation. Let $\boldsymbol{\nu} = (\nu_j)_{j \in \mathbb{N}}$ denote a multi-index of non-negative integers, with finitely many nonzero elements, i.e. $|\boldsymbol{\nu}| := \sum_{j \geq 1} \nu_j < \infty$. As usual, the value of ν_j shall determine the number of derivatives to be taken with respect to y_j , and we write $\partial^{\boldsymbol{\nu}} u$ to denote the mixed derivative of u with respect to all variables specified by the multi-index $\boldsymbol{\nu}$.

It is a simple exercise to deduce from (3.3) that

$$(\partial^{\boldsymbol{\nu}} a)(\mathbf{x}, \mathbf{y}) = (a(\mathbf{x}, \mathbf{y}) - a_*(\mathbf{x})) \prod_{j \geq 1} (\sqrt{\mu_j} \xi_j(\mathbf{x}))^{\nu_j} \quad \text{for } \boldsymbol{\nu} \neq \mathbf{0}.$$

Since from (1.2) we have $0 \leq a_*(\mathbf{x}) \leq a(\mathbf{x}, \mathbf{y})$ for all $\mathbf{x} \in D$ and $\mathbf{y} \in U_{\mathbf{b}}$, it follows that

$$\left\| \frac{\partial^{\boldsymbol{\nu}} a(\cdot, \mathbf{y})}{a(\cdot, \mathbf{y})} \right\|_{L^\infty(D)} \leq \prod_{j \geq 1} b_j^{\nu_j} \quad \text{for all } \mathbf{y} \in U_{\mathbf{b}}, \quad (3.8)$$

where b_j is defined in (3.2). The same bound holds with $a(\cdot, \mathbf{y})$ replaced by the truncated parametric coefficient $a^s(\cdot, \mathbf{y})$, uniformly with respect to $s \in \mathbb{N}$. In this case if $\nu_j > 0$ for any $j > s$ then the left hand side of (3.8) vanishes and the bound holds trivially. This leads to the following regularity result with respect to the parameters.

Theorem 14 *For any $\mathbf{y} \in U_{\mathbf{b}}$, any $f \in V'$, and for any multi-index $\boldsymbol{\nu}$ with $|\boldsymbol{\nu}| := \sum_{j \geq 1} \nu_j < \infty$, the solution $u(\cdot, \mathbf{y})$ of the parametric weak problem (3.6) satisfies the a-priori estimate*

$$\|\partial^{\boldsymbol{\nu}} u(\cdot, \mathbf{y})\|_V \leq \frac{|\boldsymbol{\nu}|!}{(\ln 2)^{|\boldsymbol{\nu}|}} \left(\prod_{j \geq 1} b_j^{\nu_j} \right) \frac{\|f\|_{V'}}{\check{a}(\mathbf{y})}. \quad (3.9)$$

Moreover, the estimate (3.9) also holds with u replaced by u_h^s .

Proof. We only establish in detail the result for u as an identical argument will apply to u_h^s with all constants appearing in the bounds being independent of s and of h . We first prove by induction on $|\boldsymbol{\nu}|$ that, for any fixed $\mathbf{y} \in U_b$,

$$\left(\int_D a(\mathbf{x}, \mathbf{y}) |\nabla(\partial^\nu u)(\mathbf{x}, \mathbf{y})|^2 d\mathbf{x} \right)^{1/2} \leq \Lambda_{|\boldsymbol{\nu}|} \left(\prod_{j \geq 1} b_j^{\nu_j} \right) \left(\int_D a(\mathbf{x}, \mathbf{y}) |\nabla u(\mathbf{x}, \mathbf{y})|^2 d\mathbf{x} \right)^{1/2}, \quad (3.10)$$

where the sequence $(\Lambda_n)_{n \geq 0}$ is defined recursively by

$$\Lambda_0 := 1 \quad \text{and} \quad \Lambda_n := \sum_{i=0}^{n-1} \binom{n}{i} \Lambda_i \quad \text{for all } n \geq 1. \quad (3.11)$$

To simplify the presentation we suppress the arguments (\mathbf{x}, \mathbf{y}) where it is clear from the context.

The case $|\boldsymbol{\nu}| = 0$ holds trivially. Suppose now that (3.10) holds for all multi-indices $\boldsymbol{\nu}$ with $|\boldsymbol{\nu}| \leq n-1$ where $n \geq 1$. Given any multi-index $\boldsymbol{\nu}$ with $|\boldsymbol{\nu}| = n$, we have the Leibniz product rule

$$\partial^\nu(AB) = \sum_{\mathbf{m} \preceq \boldsymbol{\nu}} \binom{\boldsymbol{\nu}}{\mathbf{m}} (\partial^{\boldsymbol{\nu}-\mathbf{m}} A) (\partial^{\mathbf{m}} B).$$

Here, $\mathbf{m} \preceq \boldsymbol{\nu}$ means that the multi-index \mathbf{m} satisfies $m_j \leq \nu_j$ for all j , $\boldsymbol{\nu} - \mathbf{m}$ denotes a multi-index with the elements $\nu_j - m_j$, and $\binom{\boldsymbol{\nu}}{\mathbf{m}} := \prod_{j \geq 1} \binom{\nu_j}{m_j}$. Now, applying ∂^ν to the variational formulation (3.6), and recalling that f is independent of \mathbf{y} , we obtain the identity

$$\int_D \left(\sum_{\mathbf{m} \preceq \boldsymbol{\nu}} \binom{\boldsymbol{\nu}}{\mathbf{m}} (\partial^{\boldsymbol{\nu}-\mathbf{m}} a) \nabla(\partial^{\mathbf{m}} u) \cdot \nabla v \right) d\mathbf{x} = 0 \quad \text{for all } v \in V.$$

Taking $v = \partial^\nu u(\cdot, \mathbf{y})$, separating out the $\mathbf{m} = \boldsymbol{\nu}$ term, dividing and multiplying by a , and using the Cauchy-Schwarz inequality, we obtain

$$\begin{aligned} \int_D a |\nabla(\partial^\nu u)|^2 d\mathbf{x} &= - \sum_{\substack{\mathbf{m} \preceq \boldsymbol{\nu} \\ \mathbf{m} \neq \boldsymbol{\nu}}} \binom{\boldsymbol{\nu}}{\mathbf{m}} \int_D (\partial^{\boldsymbol{\nu}-\mathbf{m}} a) \nabla(\partial^{\mathbf{m}} u) \cdot \nabla(\partial^\nu u) d\mathbf{x} \\ &\leq \sum_{\substack{\mathbf{m} \preceq \boldsymbol{\nu} \\ \mathbf{m} \neq \boldsymbol{\nu}}} \binom{\boldsymbol{\nu}}{\mathbf{m}} \left\| \frac{(\partial^{\boldsymbol{\nu}-\mathbf{m}} a)(\cdot, \mathbf{y})}{a(\cdot, \mathbf{y})} \right\|_{L^\infty(D)} \left(\int_D a |\nabla(\partial^{\mathbf{m}} u)|^2 d\mathbf{x} \right)^{1/2} \left(\int_D a |\nabla(\partial^\nu u)|^2 d\mathbf{x} \right)^{1/2}. \end{aligned}$$

Canceling the common factor on both sides and using (3.8), we arrive at

$$\left(\int_D a |\nabla(\partial^\nu u)|^2 d\mathbf{x} \right)^{1/2} \leq \sum_{\substack{\mathbf{m} \preceq \boldsymbol{\nu} \\ \mathbf{m} \neq \boldsymbol{\nu}}} \binom{\boldsymbol{\nu}}{\mathbf{m}} \left(\prod_{j \geq 1} b_j^{\nu_j - m_j} \right) \left(\int_D a |\nabla(\partial^{\mathbf{m}} u)|^2 d\mathbf{x} \right)^{1/2}.$$

We now use the inductive hypothesis (that (3.10) holds when $|\boldsymbol{\nu}| \leq n-1$) in each of the terms on the right hand side to obtain

$$\begin{aligned} \left(\int_D a |\nabla(\partial^\nu u)|^2 d\mathbf{x} \right)^{1/2} &\leq \sum_{i=0}^{n-1} \sum_{\substack{\mathbf{m} \preceq \boldsymbol{\nu} \\ |\mathbf{m}|=i}} \binom{\boldsymbol{\nu}}{\mathbf{m}} \left(\prod_{j \geq 1} b_j^{\nu_j - m_j} \right) \Lambda_i \left(\prod_{j \geq 1} b_j^{m_j} \right) \left(\int_D a |\nabla u|^2 d\mathbf{x} \right)^{1/2} \\ &= \sum_{i=0}^{n-1} \binom{n}{i} \Lambda_i \left(\prod_{j \geq 1} b_j^{\nu_j} \right) \left(\int_D a |\nabla u|^2 d\mathbf{x} \right)^{1/2} = \Lambda_n \left(\prod_{j \geq 1} b_j^{\nu_j} \right) \left(\int_D a |\nabla u|^2 d\mathbf{x} \right)^{1/2}. \end{aligned}$$

Here, we also used the identity

$$\sum_{\substack{\mathbf{m} \preceq \boldsymbol{\nu} \\ |\mathbf{m}|=i}} \binom{\boldsymbol{\nu}}{\mathbf{m}} = \binom{|\boldsymbol{\nu}|}{i},$$

which follows from a simple counting argument (i.e., consider the number of ways to select i distinct balls from some baskets containing a total number of $|\boldsymbol{\nu}|$ distinct balls). This completes the proof of (3.10).

Next we prove by induction that

$$\Lambda_n \leq \frac{n!}{(\ln 2)^n} \quad \text{for all } n \geq 0. \quad (3.12)$$

Clearly the result holds for Λ_0 . Suppose the result holds for all Λ_i with $i \leq n-1$. Then we have

$$\Lambda_n \leq \sum_{i=0}^{n-1} \binom{n}{i} \frac{i!}{(\ln 2)^i} = \frac{n!}{(\ln 2)^n} \sum_{i=0}^{n-1} \frac{(\ln 2)^{n-i}}{(n-i)!} = \frac{n!}{(\ln 2)^n} \sum_{k=1}^n \frac{(\ln 2)^k}{k!} \leq \frac{n!}{(\ln 2)^n} (e^{\ln 2} - 1),$$

and so (3.12) holds for all n .

The desired bound (3.9) is obtained by inserting (3.12) into the right hand side of (3.10), and by noting that the left hand side of (3.10) can be bounded from below by $\sqrt{\tilde{a}(\mathbf{y})} \|\partial^{\boldsymbol{\nu}} u(\cdot, \mathbf{y})\|_V$. To estimate the last factor on the right hand side of (3.10), we take $v = u(\cdot, \mathbf{y})$ in the variational form (3.6) to obtain

$$\int_D a |\nabla u|^2 \, d\mathbf{x} \leq \|f\|_{V'} \|u(\cdot, \mathbf{y})\|_V \leq \frac{\|f\|_{V'}}{\sqrt{\tilde{a}(\mathbf{y})}} \left(\int_D a |\nabla u|^2 \, d\mathbf{x} \right)^{1/2},$$

and then cancel the common factor from both sides.

Our proof argument is based entirely on the parametric weak form which is satisfied also by $u_h^s(\cdot, \mathbf{y})$ if V is replaced by V_h , $\mathbf{y} \in U_b$ is such that $y_j = 0$ for $j > s$, and a is replaced by a^s . Thus the result holds also for the finite element solution $u_h^s(\cdot, \mathbf{y})$ of the dimensionally truncated problem, with all constants independent of s and of h . \square

3.3 Equivalence of integrals

Recall that the second term in (1.10) is the error in approximating the integral $\mathbb{E}[\mathcal{G}(u_h^s)]$ by applying an n -point QMC integration rule to the s -fold iterated integral with respect to the Gaussian measure on \mathbb{R}^s given by (1.6). We transformed (1.6) to the integral (1.7) on the unit cube $(0, 1)^s$ by changing the coordinates \mathbf{y} from \mathbb{R}^s to $(0, 1)^s$ using the inverse cumulative normal distribution. The ensuing convergence analysis thus deals with a double limit, as $s \rightarrow \infty$ and $n \rightarrow \infty$, in the transformed integrals. It is a consequence of a classical theorem of Kakutani that the limit is the same, *independently of the order in which we choose to do the dimension truncation and the integral transformation into the unit cube*, as we now explain.

Reinserting the countably many standard Gaussian variables $Y_j(\omega) \sim \mathcal{N}(0, 1)$, $j = 1, 2, \dots$, into the parametric solution $u(\cdot, \mathbf{y})$, we recover the random field

$$u(\cdot, \omega) = u(\cdot, \mathbf{y})|_{\mathbf{y}=(Y_j(\omega))_{j \geq 1}} \in V. \quad (3.13)$$

Furthermore, on recalling $\bar{\mu}_G$ defined in (3.1), we see that the mathematical expectation of the random field $u(\cdot, \omega)$ defined in (3.13) is well-defined as an element of V , and we may rewrite mathematical expectations $\mathbb{E}[\cdot]$ with respect to the measure \mathbb{P} as parametric, deterministic integrals with respect to the measure $\bar{\mu}_G$, i.e.

$$\mathbb{E}[u] = \int_{\mathbf{y} \in \mathbb{R}^{\mathbb{N}}} u(\cdot, \mathbf{y}) \bar{\mu}_G(d\mathbf{y}) \in V \quad \text{and} \quad \mathbb{E}[\mathcal{G}(u)] = \int_{\mathbf{y} \in \mathbb{R}^{\mathbb{N}}} \mathcal{G}(u(\mathbf{y})) \bar{\mu}_G(d\mathbf{y}) \in \mathbb{R} \quad (3.14)$$

exist for every $u \in L^1(U_{\mathbf{b}}, \bar{\mu}_G; V)$ and for every $\mathcal{G} \in V'$.

We may now use Kakutani's theorem on the equivalence of infinite product measures (see e.g. [2, 7]), to identify the image of the Gaussian measure $\bar{\mu}_G$ under the mapping

$$\mathbf{y} \in \mathbb{R}^{\mathbb{N}} \quad \mapsto \quad \Phi(\mathbf{y}) := (\Phi(y_1), \Phi(y_2), \dots) \in (0, 1)^{\mathbb{N}},$$

with $\Phi(y) := \int_{-\infty}^y \phi(t) dt$, with the uniform probability measure λ on $(0, 1)^{\mathbb{N}}$, i.e., with the countable product of Lebesgue measures on $(0, 1)$. Moreover, the image of $U_{\mathbf{b}}$ is of full measure in $(0, 1)^{\mathbb{N}}$ and so we can reparametrize the (countably) iterated integrals with respect to the Gaussian measures $\bar{\mu}_G$ in (3.14) as integrals over the unit cube $[0, 1]^{\mathbb{N}}$. Since $u_h^s \in V_h \subset V$, exactly the same statements are true for u_h^s .

Therefore, in these parametric integrals, we may interchange the order of dimension truncation and reparametrization without affecting the numerical values of the limit as $s \rightarrow \infty$.

4 Analysis of the QMC integration error for $\mathcal{G}(u_h^s)$

In this section, we bound the QMC integration error, which is the second term on the right hand side of (1.10). Recalling (1.6), we address the efficient numerical evaluation, for large s , of integrals

$$I_s(F) := \int_{\mathbf{y} \in \mathbb{R}^s} F(\mathbf{y}) \prod_{j=1}^s \phi(y_j) d\mathbf{y}, \quad \text{with} \quad F(\mathbf{y}) := \mathcal{G}(u_h^s(\cdot, \mathbf{y})), \quad (4.1)$$

where $\phi(y) = e^{-y^2/2}/\sqrt{2\pi}$ is the standard normal probability density function. Let $\Phi(y) = \int_{-\infty}^y e^{-t^2/2}/\sqrt{2\pi} dt$ denote the cumulative normal distribution function and let Φ^{-1} denote its inverse. The integral $I_s(F)$ is transformed to the unit cube by applying Φ^{-1} component-wise, denoted by Φ_s^{-1} , as described in (1.7). We then approximate the resulting integral over the unit cube by *randomly shifted lattice rules*, leading to the formula (1.8), which we denote in this section by

$$Q_{s,n}(\Delta; F) := \frac{1}{n} \sum_{i=1}^n F \left(\Phi_s^{-1} \left(\text{frac} \left(\frac{i\mathbf{z}}{n} + \Delta \right) \right) \right). \quad (4.2)$$

We recall that $\mathbf{z} \in \mathbb{N}^s$ is the (deterministic) *generating vector* and $\Delta \in [0, 1]^s$ is the *random shift* which is uniformly distributed over $[0, 1]^s$. The quality of a randomly shifted lattice rule is determined by the choice of the generating vector \mathbf{z} . To find the best generating vector for the particular PDE problem, we need to identify a suitable weighted function space in which the integrand lies.

In §4.1 and §4.2 we write $I_s(F)$ and $Q_{s,n}(\Delta; F)$ for any general integrand F , not necessarily the one given in (4.1). In §4.3 and §4.4 we then focus on the integrand in (4.1). For simplicity we consider only linear functionals \mathcal{G} here.

4.1 A suitable weighted function space setting in \mathbb{R}^s

Most QMC methods are defined over the unit cube, and thus most QMC analyses in the literature are carried out for spaces of functions defined over the unit cube. The “standard” function spaces are *weighted Sobolev spaces* consisting of functions whose mixed first derivatives are square integrable, see e.g., [34, 35]. In particular, it is known from the standard theory that good randomly shifted lattice rules can be constructed to achieve the optimal rate of convergence close to $\mathcal{O}(n^{-1})$, provided that the integrand lies in such a weighted Sobolev space, see e.g., [33, 18, 8, 29, 30, 6, 11]; recent surveys can be found in [20, 9]. There are also higher order QMC methods that can achieve better than order one convergence for smooth integrands, see e.g., [10, 9].

For an integral of the form (4.1) over unbounded domain \mathbb{R}^s , the transformation to the unit cube yields the transformed integrand $F(\Phi_s^{-1}(\cdot))$ that may be unbounded near the boundary of the unit cube, and thus does *not* belong to the weighted Sobolev space. Consequently, the standard theory *cannot* be applied.

A suitable (but “non-standard”) function space setting for the integral (4.1) has been studied in [21, 28] (see also the earlier papers [37, 38, 16, 22]), and it is known that in this case randomly shifted lattice rules can still be constructed to achieve the optimal rate of convergence close to $\mathcal{O}(n^{-1})$. The norm in this case is given by

$$\|F\|_{\mathcal{W}_s}^2 := \sum_{\mathbf{u} \subseteq \{1:s\}} \frac{1}{\gamma_{\mathbf{u}}} \int_{\mathbb{R}^{|\mathbf{u}|}} \left(\int_{\mathbb{R}^{s-|\mathbf{u}|}} \frac{\partial^{|\mathbf{u}|} F}{\partial \mathbf{y}_{\mathbf{u}}}(\mathbf{y}_{\mathbf{u}}; \mathbf{y}_{\{1:s\} \setminus \mathbf{u}}) \prod_{j \in \{1:s\} \setminus \mathbf{u}} \phi(y_j) \, d\mathbf{y}_{\{1:s\} \setminus \mathbf{u}} \right)^2 \prod_{j \in \mathbf{u}} \psi_j^2(y_j) \, d\mathbf{y}_{\mathbf{u}}, \quad (4.3)$$

where $\{1:s\}$ is a shorthand notation for the set of indices $\{1, 2, \dots, s\}$, $\frac{\partial^{|\mathbf{u}|} F}{\partial \mathbf{y}_{\mathbf{u}}}$ denotes the mixed first derivative with respect to each of the “active” variables y_j with $j \in \mathbf{u}$, and $\mathbf{y}_{\{1:s\} \setminus \mathbf{u}}$ denotes the “inactive” variables y_j with $j \notin \mathbf{u}$. The norm (4.3) is said to be “unanchored” because the inactive variables are integrated out, as opposed to being “anchored” at some fixed value, say, 0. The unanchored norm (4.3) was first considered in [28] (which allowed also a general interval of integration instead of \mathbb{R}); an anchored norm was considered in [21].

For each $j \geq 1$, the function $\psi_j : \mathbb{R} \rightarrow \mathbb{R}^+$ in (4.3) is a positive and continuous *weight function* (not necessarily a probability density function, i.e., it does not need to integrate to 1), which is to be chosen to control the kind of functions F that are included in the space. Coordinate-dependent weight functions ψ_j were first considered in [28], while [21] used the same weight function for all coordinates, i.e. $\psi_j = \psi$ for all j . For the analysis from [21, 28] to hold, we need $\psi_j^2(y)$ to decay slower than the standard normal density $\phi(y)$ in (4.1) as $|y| \rightarrow \infty$. On the other hand, our later analysis of the integrand F in §4.3 indicates that ψ_j^2 must decay exponentially. This prompts us to restrict ourselves to the choice

$$\psi_j^2(y) = \exp(-2\alpha_j|y|) \quad \text{for some } \alpha_j > 0, \quad (4.4)$$

with the value of α_j to be specified later.

To every set $\mathbf{u} \subset \mathbb{N}$ of finite cardinality $|\mathbf{u}| < \infty$, we associate a *weight parameter* $\gamma_{\mathbf{u}} > 0$, which controls the relative importance of various (groups of) variables: small weights $\gamma_{\mathbf{u}}$ quantify “weak” dependence of the integrand F on the group of variables $\mathbf{y}_{\mathbf{u}} = \{y_j : j \in \mathbf{u}\}$. We write collectively $\boldsymbol{\gamma} = (\gamma_{\mathbf{u}})_{\mathbf{u} \subset \mathbb{N}}$, and we define $\gamma_{\emptyset} := 1$. In [21] only “product weights” were analyzed, that is, [21] assumed that there is a sequence $\gamma_1 \geq \gamma_2 \geq \dots > 0$, associated each γ_j with a single integration variable y_j , and then set $\gamma_{\mathbf{u}} := \prod_{j \in \mathbf{u}} \gamma_j$ for any nonempty subset \mathbf{u} of indices. The results of [21] were extended to more general weight parameters in [28]. (Note that [28] also allowed the weight parameters to depend on the dimension s , but we do not take this approach due to the infinite dimensional nature of our underlying PDE problem.)

The proper choice of the weight parameters $\gamma_{\mathbf{u}}$ is crucial to ensure that the constant in the error bound which we shall obtain below does not grow exponentially with increasing dimension. In the following we will consider a special form of weight parameters known as “POD weights”, which stands for “product and order dependent weights” (first seen in [19]): In this choice there are two sequences $\Gamma_0 = \Gamma_1 = 1, \Gamma_2, \dots$ and $\gamma_1 \geq \gamma_2 \geq \dots > 0$ such that $\gamma_{\mathbf{u}} := \Gamma_{|\mathbf{u}|} \prod_{j \in \mathbf{u}} \gamma_j$, where $|\mathbf{u}|$ denotes the cardinality, or the *order*, of the set \mathbf{u} .

4.2 Error analysis for randomly shifted lattice rules

To obtain a bound on the QMC integration error, we define the *worst case error* of the shifted lattice rule (4.2) with generating vector \mathbf{z} and shift $\mathbf{\Delta}$ by

$$e_{s,n}^{\text{wor}}(\mathbf{z}, \mathbf{\Delta}) := \sup_{\|F\|_{\mathcal{W}_s} \leq 1} |I_s(F) - Q_{s,n}(\mathbf{\Delta}, F)| .$$

Due to linearity of the exact and approximate integration problems, for any integrand $F \in \mathcal{W}_s$ we then have the error bound

$$|I_s(F) - Q_{s,n}(\mathbf{\Delta}, F)| \leq e_{s,n}^{\text{wor}}(\mathbf{z}, \mathbf{\Delta}) \|F\|_{\mathcal{W}_s} . \quad (4.5)$$

For randomly shifted lattice rules, we consider the *root-mean-square error*, i.e., on both sides of the inequality in (4.5) we take the square, take the expectation over the random shift $\mathbf{\Delta}$ which is uniformly distributed over $[0, 1]^s$, and then take the square root. This yields

$$\sqrt{\mathbb{E}^{\mathbf{\Delta}} |I_s(F) - Q_{s,n}(\cdot, F)|^2} \leq e_{s,n}^{\text{sh}}(\mathbf{z}) \|F\|_{\mathcal{W}_s} , \quad (4.6)$$

with

$$[e_{s,n}^{\text{sh}}(\mathbf{z})]^2 := \int_{[0,1]^s} [e_{s,n}^{\text{wor}}(\mathbf{z}, \mathbf{\Delta})]^2 d\mathbf{\Delta} .$$

The quantity $e_{s,n}^{\text{sh}}(\mathbf{z})$ is often referred to as the *shift averaged worst case error*. An upper bound of the form (4.6) conveniently decouples the error dependence on the generating vector \mathbf{z} from its dependence on the integrand F .

A generating vector \mathbf{z} for a randomly shifted lattice rule will be constructed using a *component-by-component algorithm* which determines in turn z_1, z_2, z_3 , and so on. The quantity $[e_{s,n}^{\text{sh}}(\mathbf{z})]^2$ will be used as the search criterion: assuming that the components z_1, \dots, z_j are already chosen and fixed, the component z_{j+1} is chosen from the set $\{1 \leq z \leq n-1 : \gcd(z, n) = 1\}$ of at most $n-1$ values to minimize $[e_{j+1,n}^{\text{sh}}(z_1, \dots, z_{j+1})]^2$. The precise formula for $[e_{s,n}^{\text{sh}}(\mathbf{z})]^2$ depends on the choices of weight functions ψ_j and weight parameters γ_u , see [28]:

$$[e_{s,n}^{\text{sh}}(\mathbf{z})]^2 = \sum_{\emptyset \neq \mathbf{u} \subseteq \{1:s\}} \frac{\gamma_{\mathbf{u}}}{n} \sum_{i=1}^n \prod_{j \in \mathbf{u}} \theta_j \left(\left\{ \frac{iz_j}{n} \right\} \right) ,$$

where

$$\theta_j(u) = \int_{\Phi^{-1}(u)}^{\infty} \frac{\Phi(t) - u}{\psi_j^2(t)} dt + \int_{\Phi^{-1}(1-u)}^{\infty} \frac{\Phi(t) - 1 + u}{\psi_j^2(t)} dt - \int_{-\infty}^{\infty} \frac{\Phi^2(t)}{\psi_j^2(t)} dt .$$

It is shown in [28] that, with POD weights γ_u , the total cost for constructing a lattice rule up to dimension s is $\mathcal{O}(n \log n s + ns^2)$ operations using FFT. It is also proved in [28] that for certain combinations of ϕ and weight functions ψ_j , we can obtain close to $\mathcal{O}(n^{-1})$ convergence for $e_{s,n}^{\text{sh}}(\mathbf{z})$. We present the relevant result from [28].

Theorem 15 *Let $F \in \mathcal{W}_s$. Given $s, n \in \mathbb{N}$, weight parameters $\gamma = (\gamma_{\mathbf{u}})_{\mathbf{u} \subseteq \mathbb{N}}$, standard normal density function ϕ , and weight functions ψ_j defined by (4.4), a randomly shifted lattice rule with n points in s dimensions can be constructed using a component-by-component algorithm such that, for all $\lambda \in (1/2, 1]$,*

$$\sqrt{\mathbb{E}^{\mathbf{\Delta}} |I_s(F) - Q_{s,n}(\cdot, F)|^2} \leq \left(\sum_{\emptyset \neq \mathbf{u} \subseteq \{1:s\}} \gamma_{\mathbf{u}}^{\lambda} \prod_{j \in \mathbf{u}} \varrho_j(\lambda) \right)^{1/(2\lambda)} [\varphi_{\text{tot}}(n)]^{-1/(2\lambda)} \|F\|_{\mathcal{W}_s} , \quad (4.7)$$

with

$$\varrho_j(\lambda) := 2 \left(\frac{\sqrt{2\pi} \exp(\alpha_j^2/\eta_*)}{\pi^{2-2\eta_*}(1-\eta_*)\eta_*} \right)^\lambda \zeta \left(\lambda + \frac{1}{2} \right) \quad \text{and} \quad \eta_* := \frac{2\lambda - 1}{4\lambda}, \quad (4.8)$$

where $\varphi_{\text{tot}}(n) := |\{1 \leq z \leq n-1 : \gcd(z, n) = 1\}|$ denotes the Euler totient function, and $\zeta(x) := \sum_{k=1}^{\infty} k^{-x}$ denotes the Riemann zeta function.

Proof. Theorem 8 in [28] yields the error bound (4.7) together with $\varrho_j(\lambda) = 2c_{2,j}^\lambda \zeta(2r_2\lambda)$, which holds for all $\lambda \in (1/(2r_2), 1]$, with the precise values of $c_{2,j}$ and r_2 depending on the particular combination of ϕ and ψ_j , which in [28] are general. For the choice of ϕ and ψ_j in this theorem, we have (see [21, Example 5])

$$c_{2,j} = \frac{\sqrt{2\pi} \exp(\alpha_j^2/\eta)}{\pi^{2-2\eta}(1-\eta)\eta} \quad \text{and} \quad r_2 = 1 - \eta \quad \text{for any} \quad \eta \in (0, 1/2).$$

Thus the bound holds for all $\eta \in (0, 1/2)$ and $\lambda \in (1/(2-2\eta), 1]$. Equivalently, the bound holds for all $\lambda \in (1/2, 1]$ and $\eta \in (0, 1 - 1/(2\lambda))$. We simplify the result by taking $\eta = \eta_*$ to be the mid-point of the latter interval, see (4.8). \square

Note that $\varphi_{\text{tot}}(n) = n - 1$ for n prime, and it can be verified that $1/\varphi_{\text{tot}}(n) < 9/n$ for all $n \leq 10^{30}$. Hence, from the practical point of view, we can replace the reciprocal of the Euler totient function by a constant factor times $1/n$.

4.3 Justifying the choice of the weight functions ψ_j

In this subsection we show in Theorem 16 that $\|F\|_{\mathcal{W}_s} < \infty$ for each s and for each choice of weight parameters γ_u (where the norm is defined in (4.3)). In the proof of this result we make crucial use of our specific choice of weight functions ψ_j in (4.4). This leads to Theorem 17, which gives an estimate for the root-mean-square error and shows that this attains a rate of convergence arbitrarily close to $\mathcal{O}(n^{-1})$, but with a possibly s -dependent asymptotic constant. Then in the following subsection we show that a careful choice of the weight parameters γ_u can be made so that the asymptotic constant in the convergence estimate is bounded uniformly with respect to s , leading to the main result Theorem 20.

We assume throughout the remainder of the paper that Assumption A3 holds for some $p \leq 1$. We assume also that the parameters α_j in (4.4) satisfy for some constants $0 < \alpha_{\min} < \alpha_{\max} < \infty$

$$\max(b_j, \alpha_{\min}) < \alpha_j \leq \alpha_{\max}, \quad j \in \mathbb{N}. \quad (4.9)$$

Theorem 16 *For each $j \geq 1$, let b_j be defined by (3.2) and ψ_j by (4.4) with parameters α_j satisfying (4.9). Then the norm (4.3) of the integrand F in (4.1) satisfies the bound*

$$\|F\|_{\mathcal{W}_s}^2 \leq (C^*)^2 \sum_{u \subseteq \{1:s\}} \frac{(|u|!)^2}{\gamma_u (\ln 2)^{2|u|}} \prod_{j \in u} \frac{\tilde{b}_j^2}{\alpha_j - b_j} \quad (4.10)$$

where

$$\tilde{b}_j^2 := \frac{b_j^2}{2 \exp(b_j^2/2) \Phi(b_j)}, \quad (4.11)$$

with $\Phi(\cdot)$ denoting the cumulative standard normal distribution function, and with

$$C^* := \frac{\|f\|_{V'} \|\mathcal{G}\|_{V'}}{\min_{\mathbf{x} \in \overline{D}} a_0(\mathbf{x})} \exp \left(\frac{1}{2} \sum_{j \geq 1} b_j^2 + \frac{2}{\sqrt{2\pi}} \sum_{j \geq 1} b_j \right). \quad (4.12)$$

Proof. To ease readability of the proof we introduce $K^* := \|f\|_{V'} \|\mathcal{G}\|_{V'}/\min_{\mathbf{x} \in \bar{D}} a_0(\mathbf{x})$. Now, for the integrand F from (4.1) and for any $\mathbf{y} \in \mathbb{R}^s$ (which we identify throughout this proof, with slight abuse of notation, with the sequence $\mathbf{y} \in \mathbb{R}^{\mathbb{N}}$ with $y_j = 0$ for $j > s$), we have from Theorem 14 with $\nu_j \in \{0, 1\}$ and the linearity of \mathcal{G} , that

$$\left| \frac{\partial^{|\mathbf{u}|} F}{\partial \mathbf{y}_{\mathbf{u}}}(\mathbf{y}) \right| \leq \|\mathcal{G}\|_{V'} \left\| \frac{\partial^{|\mathbf{u}|} u_h^s(\cdot, \mathbf{y})}{\partial \mathbf{y}_{\mathbf{u}}} \right\|_V \leq \|f\|_{V'} \|\mathcal{G}\|_{V'} \frac{|\mathbf{u}|!}{(\ln 2)^{|\mathbf{u}|}} \left(\prod_{j \in \mathbf{u}} b_j \right) \frac{1}{\tilde{a}(\mathbf{y})}.$$

Since a_* in (1.2) was assumed to be non-negative and since $\min_{\mathbf{x} \in \bar{D}} \sqrt{\mu_j} \xi_j(\mathbf{x}) \geq -b_j$, this implies

$$\left| \frac{\partial^{|\mathbf{u}|} F}{\partial \mathbf{y}_{\mathbf{u}}}(\mathbf{y}) \right| \leq K^* \frac{|\mathbf{u}|!}{(\ln 2)^{|\mathbf{u}|}} \left(\prod_{j \in \mathbf{u}} b_j \right) \left(\prod_{j \in \{1:s\} \setminus \mathbf{u}} \exp(b_j |y_j|) \right). \quad (4.13)$$

Since the final factor on the right hand side of (4.13) is a product, we can group the factors corresponding to $j \in \mathbf{u}$ and $j \in \{1:s\} \setminus \mathbf{u}$ separately, allowing us to estimate the norm (4.3) as

$$\begin{aligned} \|F\|_{\mathcal{W}_s}^2 &\leq (K^*)^2 \sum_{\mathbf{u} \subseteq \{1:s\}} \frac{1}{\gamma_{\mathbf{u}}} \frac{|\mathbf{u}|!^2}{(\ln 2)^{2|\mathbf{u}|}} \left(\prod_{j \in \mathbf{u}} b_j \right)^2 \left(\int_{\mathbb{R}^{s-|\mathbf{u}|}} \prod_{j \in \{1:s\} \setminus \mathbf{u}} \exp(b_j |y_j|) \phi(y_j) \, d\mathbf{y}_{\{1:s\} \setminus \mathbf{u}} \right)^2 \\ &\quad \times \int_{\mathbb{R}^{|\mathbf{u}|}} \prod_{j \in \mathbf{u}} \exp(2b_j |y_j|) \psi_j^2(y_j) \, d\mathbf{y}_{\mathbf{u}}. \end{aligned} \quad (4.14)$$

The integrals on the right hand side of (4.14) can be readily estimated. Firstly,

$$\begin{aligned} \int_{\mathbb{R}^{s-|\mathbf{u}|}} \prod_{j \in \{1:s\} \setminus \mathbf{u}} \exp(b_j |y_j|) \phi(y_j) \, d\mathbf{y}_{\{1:s\} \setminus \mathbf{u}} &= \prod_{j \in \{1:s\} \setminus \mathbf{u}} \left(\int_{-\infty}^{\infty} \exp(b_j |y|) \frac{\exp(-y^2/2)}{\sqrt{2\pi}} \, dy \right) \\ &= \prod_{j \in \{1:s\} \setminus \mathbf{u}} \left(2 \exp(b_j^2/2) \int_0^{\infty} \frac{\exp(-(y-b_j)^2/2)}{\sqrt{2\pi}} \, dy \right) \\ &= \prod_{j \in \{1:s\} \setminus \mathbf{u}} (2 \exp(b_j^2/2) \Phi(b_j)). \end{aligned} \quad (4.15)$$

Secondly,

$$\int_{\mathbb{R}^{|\mathbf{u}|}} \prod_{j \in \mathbf{u}} \exp(2b_j |y_j|) \psi_j^2(y_j) \, d\mathbf{y}_{\mathbf{u}} = \prod_{j \in \mathbf{u}} \left(\int_{-\infty}^{\infty} \exp(2b_j |y|) \psi_j^2(y) \, dy \right),$$

and from this we understand the requirement (explained in §4.1) that ψ_j^2 must decay exponentially. With ψ_j defined by (4.4) and using the condition (4.9), we obtain

$$\int_{\mathbb{R}^{|\mathbf{u}|}} \prod_{j \in \mathbf{u}} \exp(2b_j |y_j|) \psi_j^2(y_j) \, d\mathbf{y}_{\mathbf{u}} = \prod_{j \in \mathbf{u}} \frac{1}{\alpha_j - b_j}. \quad (4.16)$$

Combining (4.14) with (4.15) and (4.16), we obtain

$$\|F\|_{\mathcal{W}_s}^2 \leq (K^*)^2 \prod_{j \in \{1:s\}} (2 \exp(b_j^2/2) \Phi(b_j)) \sum_{\mathbf{u} \subseteq \{1:s\}} \left(\frac{1}{\gamma_{\mathbf{u}}} \frac{|\mathbf{u}|!^2}{(\ln 2)^{2|\mathbf{u}|}} \prod_{j \in \mathbf{u}} \frac{\tilde{b}_j^2}{\alpha_j - b_j} \right). \quad (4.17)$$

Now, to obtain the bound (4.10), it remains to bound the product in (4.17) independently of s . To do this we note that $2 \exp(b_j^2/2) \Phi(b_j) \geq 1$ and

$$\Phi(b_j) \leq \frac{1}{2} \left(1 + \frac{2b_j}{\sqrt{2\pi}} \right) \leq \frac{1}{2} \exp\left(\frac{2b_j}{\sqrt{2\pi}}\right) \quad \text{since } b_j \geq 0.$$

Thus we have $\prod_{j \in \{1:s\}} (2 \exp(b_j^2/2) \Phi(b_j)) \leq \prod_{j \geq 1} \exp(b_j^2/2 + 2b_j/\sqrt{2\pi})$ and the bound (4.10) then follows. \square

The root-mean-square error can now be estimated by combining Theorems 15 and 16.

Theorem 17 *Let F be the integrand in (4.1), and for each $j \geq 1$ let ψ_j be defined by (4.4) with α_j satisfying (4.9). Given $s, n \in \mathbb{N}$ with $n \leq 10^{30}$, weights $\gamma = (\gamma_u)_{u \subset \mathbb{N}}$, and standard normal density function ϕ , a randomly shifted lattice rule with n points in s dimensions can be constructed by a component-by-component algorithm such that, for all $\lambda \in (1/2, 1]$,*

$$\sqrt{\mathbb{E}^\Delta |I_s(F) - Q_{s,n}(\cdot; F)|^2} \leq 9C^* C_{\gamma,s}(\lambda) n^{-1/(2\lambda)}, \quad (4.18)$$

with

$$C_{\gamma,s}(\lambda) := \left(\sum_{\emptyset \neq u \subseteq \{1:s\}} \gamma_u^\lambda \prod_{j \in u} \varrho_j(\lambda) \right)^{1/(2\lambda)} \left(\sum_{u \subseteq \{1:s\}} \frac{(|u|!)^2}{\gamma_u (\ln 2)^{2|u|}} \prod_{j \in u} \frac{\tilde{b}_j^2}{\alpha_j - b_j} \right)^{1/2},$$

where b_j is defined in (3.2), \tilde{b}_j is defined in (4.11), C^* is defined in (4.12), and $\varrho_j(\lambda)$ is defined in (4.8) (making $\varrho_j(\lambda)$ depend on α_j).

Without a careful choice of the weight parameters γ_u , the quantity $C_{\gamma,s}(\lambda)$ might grow (even exponentially) with increasing s . To ensure that $C_{\gamma,s}(\lambda)$ is bounded independently of s , we choose the weight parameters to ensure that

$$C_\gamma(\lambda) := \left(\sum_{|u| < \infty} \gamma_u^\lambda \prod_{j \in u} \varrho_j(\lambda) \right)^{1/(2\lambda)} \left(\sum_{|u| < \infty} \frac{(|u|!)^2}{\gamma_u (\ln 2)^{2|u|}} \prod_{j \in u} \frac{\tilde{b}_j^2}{\alpha_j - b_j} \right)^{1/2} < \infty. \quad (4.19)$$

(Note that $\tilde{b}_j \leq b_j$ and that it tends to b_j rapidly as $j \rightarrow \infty$.) Provided (4.19) holds, it follows immediately that $C_{\gamma,s}(\lambda) \leq C_\gamma(\lambda) < \infty$ for all s , and so the asymptotic constant in the convergence estimate (4.18) is independent of the truncation dimension s .

4.4 Choosing the weight parameters γ_u

For any given $\lambda \in (1/2, 1]$, we now follow the strategy in [19] and choose the weight parameters γ_u to minimize the constant $C_\gamma(\lambda)$ given in (4.19). We shall see that the resulting minimal value of $C_\gamma(\lambda)$ is finite. To do this we will use the following two lemmas.

Lemma 18 ([19, Lemma 6.2]) *Let $m \in \mathbb{N}$, $\lambda > 0$, and $A_i, B_i > 0$ for all i . Then the function*

$$\left(\sum_{i=1}^m x_i^\lambda A_i \right)^{1/\lambda} \left(\sum_{i=1}^m \frac{B_i}{x_i} \right) \quad (4.20)$$

is minimized over all sequences $(x_i)_{1 \leq i \leq m}$ when

$$x_i = c \left(\frac{B_i}{A_i} \right)^{1/(1+\lambda)} \quad \text{for any } c > 0. \quad (4.21)$$

The function obtained by letting $m \rightarrow \infty$ in (4.20) is minimized when x_i is given by (4.21) for all i and has a finite value if and only if the series $\sum_{i=1}^{\infty} (A_i B_i^\lambda)^{1/(1+\lambda)}$ converges.

Lemma 19 ([19, Lemma 6.3]) *For all $A_j > 0$ with $\sum_{j \geq 1} A_j < 1$ we have*

$$\sum_{|u| < \infty} |u|! \prod_{j \in u} A_j \leq \sum_{k=0}^{\infty} \left(\sum_{j \geq 1} A_j \right)^k = \frac{1}{1 - \sum_{j \geq 1} A_j},$$

and for all $B_j > 0$ with $\sum_{j \geq 1} B_j < \infty$ we have

$$\sum_{|u| < \infty} \prod_{j \in u} B_j = \prod_{j \geq 1} (1 + B_j) = \exp \left(\sum_{j \geq 1} \log(1 + B_j) \right) \leq \exp \left(\sum_{j \geq 1} B_j \right).$$

Since the constant $C_{\gamma,s}(\lambda)$ in Theorem 17 and the uniform upper bound $C_\gamma(\lambda)$ in (4.19) are of the same general form as the function in Lemma 18, we obtain the formula (4.23) for the weights γ_u below. We then specify the parameter λ to obtain a good convergence rate, while ensuring that the constant $C_\gamma(\lambda)$ is indeed finite for our choice of weights (4.23).

Theorem 20 *For each $j \geq 1$, let ψ_j be defined by (4.4) with α_j satisfying (4.9). Suppose that Assumption A3 holds for some $p \leq 1$, and when $p = 1$ assume additionally that*

$$\sum_{j \geq 1} b_j < \ln 2 \sqrt{\frac{\mathcal{J}}{\varrho_{\max}(1)}}, \quad (4.22)$$

where $\mathcal{J} := \inf_{j \geq 1} (\alpha_j - b_j) > 0$ and $\varrho_{\max}(\lambda)$ is defined by replacing α_j in (4.8) by α_{\max} in (4.9). Then, for any given $\lambda \in (1/2, 1]$, the choice of weights

$$\gamma_u = \gamma_u^*(\lambda) := \left(\frac{(|u|!)^2}{(\ln 2)^{2|u|}} \prod_{j \in u} \frac{\tilde{b}_j^2}{(\alpha_j - b_j) \varrho_j(\lambda)} \right)^{1/(1+\lambda)} \quad (4.23)$$

minimizes $C_\gamma(\lambda)$ given in (4.19), if a finite minimum exists. If we furthermore choose

$$\lambda = \lambda_* := \begin{cases} \frac{1}{2-2\delta} & \text{for arbitrary } \delta \in (0, 1/2] \text{ when } p \in (0, 2/3], \\ \frac{p}{2-p} & \text{when } p \in (2/3, 1), \\ 1 & \text{when } p = 1, \end{cases} \quad (4.24)$$

and set $\gamma_u = \gamma_u^*(\lambda_*)$, then $C_\gamma(\lambda) < \infty$. Moreover, a randomly shifted lattice rule can be constructed for the approximation of the integral (4.1) such that

$$\sqrt{\mathbb{E}^\Delta |I_s(F) - Q_{s,n}(\cdot; F)|^2} = \begin{cases} \mathcal{O}(n^{-(1-\delta)}) & \text{when } p \in (0, 2/3], \\ \mathcal{O}(n^{-(1/p-1/2)}) & \text{when } p \in (2/3, 1), \\ \mathcal{O}(n^{-1/2}) & \text{when } p = 1, \end{cases}$$

with the implied constant independent of s , but depending on p and, when relevant, δ .

Proof. The fact that the choice of weights (4.23) minimizes $C_\gamma(\lambda)$ follows from Lemma 18, as in [19, Theorem 6.4], on observing that the finite subsets of \mathbb{N} in (4.19) can be ordered (i.e. are countable), and that the particular ordering is immaterial, as the convergence is absolute and hence unconditional. In the following, we show that $C_\gamma(\lambda)$ is indeed finite for the weights given by (4.23) and parameter λ given by (4.24). In the course of our derivation below we shall choose the value of λ according to the value of p , see (4.24), but until then λ is independent of p .

Let us define

$$S_\lambda := \sum_{|u| < \infty} (\gamma_u^*)^\lambda \prod_{j \in u} \varrho_j(\lambda) = \sum_{|u| < \infty} \left(\frac{(|u|!)^2}{(\ln 2)^{2|u|}} \prod_{j \in u} \frac{[\varrho_j(\lambda)]^{1/\lambda} \tilde{b}_j^2}{\alpha_j - b_j} \right)^{\lambda/(1+\lambda)}. \quad (4.25)$$

Then $S_\lambda^{1/(2\lambda)}$ is the first factor of $C_\gamma(\lambda)$ in (4.19) with the choice of weight parameters (4.23). Moreover, the second factor in $C_\gamma(\lambda)$ can also be shown to reduce to $S_\lambda^{1/2}$. Thus we have $C_\gamma(\lambda) = S_\lambda^{1/(2\lambda)+1/2}$. So, to prove $C_\gamma(\lambda)$ is finite it suffices to prove that S_λ is finite.

By definition we have $\alpha_j - b_j \geq \mathcal{J}$ and $\tilde{b}_j \leq b_j$ for all $j \leq s$. (Note that $\mathcal{J} > 0$, since b_j converges to 0 while $\alpha_j \geq \alpha_{\min} > 0$.) On the other hand, we see from (4.8) that, for fixed λ ,

$\varrho_j(\lambda)$ increases monotonically with α_j . Thus we have $\varrho_j(\lambda) \leq \varrho_{\max}(\lambda)$ for all $j \geq 1$. Applying these estimates to S_λ in (4.25) yields

$$S_\lambda \leq \sum_{|\mathbf{u}| < \infty} (|\mathbf{u}|!)^{2\lambda/(1+\lambda)} \prod_{j \in \mathbf{u}} \left(\frac{[\varrho_{\max}(\lambda)]^{1/\lambda}}{\mathcal{J}(\ln 2)^2} b_j^2 \right)^{\lambda/(1+\lambda)}. \quad (4.26)$$

In the following we consider the cases $\lambda \neq 1$ and $\lambda = 1$ separately.

For $\lambda \in (1/2, 1)$, we have $2\lambda/(1+\lambda) < 1$ and we further estimate S_λ as follows: we multiply and divide the terms on the right hand side of (4.26) by $\prod_{j \in \mathbf{u}} A_j^{2\lambda/(1+\lambda)}$, where $A_j > 0$ will be specified later, and then apply Hölder's inequality with conjugate exponents $(1+\lambda)/(2\lambda)$ and $(1+\lambda)/(1-\lambda)$, to obtain

$$\begin{aligned} S_\lambda &\leq \sum_{|\mathbf{u}| < \infty} (|\mathbf{u}|!)^{2\lambda/(1+\lambda)} \prod_{j \in \mathbf{u}} A_j^{2\lambda/(1+\lambda)} \prod_{j \in \mathbf{u}} \left(\frac{[\varrho_{\max}(\lambda)]^{1/\lambda}}{\mathcal{J}(\ln 2)^2} \frac{b_j^2}{A_j^2} \right)^{\lambda/(1+\lambda)} \\ &\leq \left(\sum_{|\mathbf{u}| < \infty} |\mathbf{u}|! \prod_{j \in \mathbf{u}} A_j \right)^{2\lambda/(1+\lambda)} \left(\sum_{|\mathbf{u}| < \infty} \prod_{j \in \mathbf{u}} \left(\frac{[\varrho_{\max}(\lambda)]^{1/\lambda}}{\mathcal{J}(\ln 2)^2} \frac{b_j^2}{A_j^2} \right)^{\lambda/(1-\lambda)} \right)^{(1-\lambda)/(1+\lambda)} \\ &\leq \left(\frac{1}{1 - \sum_{j \geq 1} A_j} \right)^{2\lambda/(1+\lambda)} \exp \left(\frac{1-\lambda}{1+\lambda} \left(\frac{[\varrho_{\max}(\lambda)]^{1/\lambda}}{\mathcal{J}(\ln 2)^2} \right)^{\lambda/(1-\lambda)} \sum_{j \geq 1} \left(\frac{b_j}{A_j} \right)^{2\lambda/(1-\lambda)} \right). \end{aligned}$$

In the last step we applied Lemma 19 which holds and guarantees that S_λ is finite, provided that

$$\sum_{j \geq 1} A_j < 1 \quad \text{and} \quad \sum_{j \geq 1} \left(\frac{b_j}{A_j} \right)^{2\lambda/(1-\lambda)} < \infty. \quad (4.27)$$

We now choose

$$A_j := \frac{b_j^p}{\varpi} \quad \text{for some } \varpi > \sum_{j \geq 1} b_j^p.$$

Then we have $\sum_{j \geq 1} A_j < 1$ due to Assumption A3. Noting that Assumption A3 also implies that $\sum_{j \geq 1} b_j^{p'} < \infty$ for all $p' \geq p$, we conclude that the second sum in (4.27) converges for

$$\frac{2\lambda}{1-\lambda}(1-p) \geq p \quad \iff \quad p \leq \frac{2\lambda}{1+\lambda} \quad \iff \quad \lambda \geq \frac{p}{2-p}.$$

Recall that λ must be strictly between $1/2$ and 1 for the argument above. When $p \in (0, 2/3]$, we choose $\lambda = 1/(2-2\delta)$ for some $\delta \in (0, 1/2)$. When $p \in (2/3, 1)$, we set $\lambda = p/(2-p)$.

In the case $p = 1$ we take $\lambda = 1$. Then, using Lemma 19, we obtain from (4.26) that

$$S_1 \leq \sum_{|\mathbf{u}| < \infty} |\mathbf{u}|! \prod_{j \in \mathbf{u}} \left(\frac{\varrho_{\max}(1)}{\mathcal{J}(\ln 2)^2} b_j^2 \right)^{1/2} \leq \left(1 - \sum_{j \geq 1} \sqrt{\frac{\varrho_{\max}(1)}{\mathcal{J}} \frac{b_j}{\ln 2}} \right)^{-1},$$

which is finite due to the assumption (4.22). This completes the proof. \square

Corollary 21 *Let $\lambda = \lambda_*$ and $\gamma_{\mathbf{u}} = \gamma_{\mathbf{u}}^*(\lambda_*)$, as defined in (4.24) and (4.23), respectively. Then the constant $C_{\gamma}(\lambda)$ in (4.19) is minimized by choosing*

$$\alpha_j = \frac{1}{2} \left(b_j + \sqrt{b_j^2 + 1 - \frac{1}{2\lambda_*}} \right), \quad \text{for all } j \geq 1. \quad (4.28)$$

Proof. Recall from the proof of Theorem 20 that $C_\gamma(\lambda) = S_\lambda^{1/(2\lambda)+1/2}$ with S_λ given by (4.25). Since all terms on the right hand side of (4.25) are positive, minimizing $C_\gamma(\lambda)$ with respect to the parameters $\{\alpha_j\}_{j \geq 1}$ is equivalent to minimizing each of the functions $[\varrho_j(\lambda)]^{1/\lambda}/(\alpha_j - b_j)$ individually with respect to α_j . But due to (4.8), $[\varrho_j(\lambda)]^{1/\lambda} = c \exp(\alpha_j^2/\eta_*)$, for some constant c independent of α_j and for $\eta_* = 1/2 - 1/(4\lambda)$, leading to the choice (4.28) for the minimizer. \square

Following the argument in the proof of [19, Theorem 6.5], we can prove that the alternative choice of weights

$$\gamma_{\mathbf{u}} = \gamma_{\mathbf{u}}^{**} := \left(|\mathbf{u}|! \prod_{j \in \mathbf{u}} (\kappa b_j) \right)^{2-p} \quad \text{for arbitrary } \kappa > 0,$$

while not minimizing $C_\gamma(\lambda)$, still ensures that $C_\gamma(\lambda) < \infty$ and yields the same convergence rates under the same conditions on b_j . This result might seem to indicate that the approximation is somewhat robust with respect to the scaling parameters κ . However, numerical experiments indicate that arbitrary choices of κ can lead to poor lattice rules (see [27] for details). Therefore, we recommend the choice of weight parameters (4.23) that minimizes the bound.

Similarly, although Theorem 20 holds for any choice of $\{\alpha_j\}_{j \geq 1}$ that satisfies (4.9), numerical experiments show that arbitrary choices, such as $\alpha_j = 2b_j$, can again lead to poor lattice rules.

5 Final result

We now summarize our theoretical results and state our combined bound for the root-mean-square error, which includes the finite element error, the dimension truncation error and the QMC quadrature error, estimated in Theorems 6, 8 and 20, respectively.

Theorem 22 *We consider approximations of the expected value of $\mathcal{G}(u)$ via quasi-Monte Carlo finite element methods. In particular, we apply a randomly shifted lattice rule $Q_{s,n}$ to $\mathcal{G}(u_h^s)$. Then, under the same assumptions and definitions as in Theorems 6, 8 and 20, the root-mean-square error with respect to the uniformly distributed shift $\Delta \in [0, 1]^s$ can be bounded by*

$$\sqrt{\mathbb{E}^\Delta \left[\left(\mathbb{E}[\mathcal{G}(u)] - Q_{s,n}(\cdot; \mathcal{G}(u_h^s)) \right)^2 \right]} \leq C (h^{2\tau} + s^{-\chi} + n^{-r}), \quad (5.1)$$

for some $0 < \tau \leq 1$ and $0 < \chi < (1/2 - \varepsilon)\Theta - 1/2$, and with $r = 1/p - 1/2$ for $p \in (2/3, 1]$ and $r = 1 - \delta$ for $p \leq 2/3$, with δ arbitrarily small.

The rate τ depends on the spatial regularity of u in Assumption A1, while the rates χ and r depend on the parameters ε , Θ and p which in turn depend on the asymptotics of the Karhunen–Loève eigenvalues and eigenvectors in Assumptions A2 and A3. The constant C is independent of h , s , and n .

In one spatial dimension in the case of the Matérn covariance ρ_ν from (2.5), if $\|\xi_j\|_{C^0(\overline{D})}$ is uniformly bounded and $\|\nabla \xi_j\|_{C^0(\overline{D})}$ grows no faster than μ_j^{-1} (as seems to be the case numerically; see Figure 1), then the result holds for any $\tau < \min(\nu, 1)$, $r < \min(\nu, 1)$, and for $\chi < \nu$, provided $\nu > 1/2$. For details see the discussions after each of the assumptions in Sections 2 and 3 above.

Note that the rate r is capped at 1 even for $p < 2/3$ because we are using only QMC methods of order one. With higher order QMC methods we might expect to have $r = 1/p - 1/2$ also for $p < 2/3$, and in one spatial dimension with the Matérn covariance ρ_ν , we might expect r to be close to ν , even for $\nu > 1$. Similarly, we could use higher order finite elements in space to ensure that τ is close to ν , for all $\nu > 1/2$, but this is classical. A recent result in [4] shows that under slightly stronger conditions on the data, the rate χ in the truncation error can also be increased to 2χ , which is what is observed numerically.

Finally, we briefly discuss the ε -cost of the algorithm, i.e., the cost to obtain an overall error bounded by ε . Treating the construction of the randomly shifted lattice rule as pre-computation and assuming for the moment that the Karhunen-Loève eigensystem is known, an evaluation of the approximate coefficient $a^s(\mathbf{x}, \omega)$ in (2.13) at all the grid points costs $\mathcal{O}(s h^{-d})$ operations. Therefore, the total cost of our algorithm is $\mathcal{O}(n s h^{-d})$, provided a linear-complexity finite element solver is available, i.e., provided each linear solve costs $\mathcal{O}(h^{-d})$ operations. Note that the cost of computing s terms of the Karhunen-Loève eigensystem is of lower order when n is large and h is small. Similarly, the pre-computation cost of the lattice rule is of lower order when h is small (cf. [28]). For a given error threshold of $\varepsilon > 0$, we can choose n, s, h such that each of the three error contributions in (5.1) is of order ε . The cost is then $\mathcal{O}(\varepsilon^{-[1/r+1/\chi+d/(2\tau)]})$ operations, which in the Matérn case in the limit as $\nu \rightarrow \infty$ approaches $\mathcal{O}(\varepsilon^{-(d/2+1)})$, as opposed to $\mathcal{O}(\varepsilon^{-(d/2+2)})$ operations for standard Monte Carlo, or in terms of the mesh size, $\mathcal{O}(h^{-(d+2)})$ instead of $\mathcal{O}(h^{-(d+4)})$ operations.

6 Numerical results

We present here a numerical study of the algorithm described above over a range of parameters. Theorem 22 provides us with a theoretical bound for the error in the method, and we examine here whether we see, in the numerics, the behavior predicted by the theory. We focus our experiments on the convergence of the QMC error, since this is the novel element of this paper.

We solve (3.6) with spatial dimension $d = 1$ on $D = [0, 1]$, with a forcing term $f(x) = 1$. The Karhunen-Loève expansion of the random field is truncated at $s = 400$, so that $\mathbf{y} \in \mathbb{R}^{400}$. The strong form of the problem we are solving is the parametrized ODE

$$-\frac{d}{dx} \left(a^s(x, \mathbf{y}) \frac{du^s(x, \mathbf{y})}{dx} \right) = 1, \quad (6.1)$$

with homogeneous Dirichlet boundary conditions, $u(0, \mathbf{y}) = u(1, \mathbf{y}) = 0$. We solve (6.1) using the piecewise linear finite element method with uniform meshes of diameter $h = 1/M$ to get the approximate solution $u_h^s(\cdot, \mathbf{y})$. The tridiagonal systems which arise are solved in $\mathcal{O}(M)$ time by the Thomas algorithm. In the numerical experiments that follow, we set $M = 1024$ and compute the entries of the tridiagonal system using the composite mid-point rule applied elementwise. Our choices of s and h ensure that both the dimension truncation error and finite element discretization error are sufficiently small compared to the QMC error.

The quantity of interest is here taken to be $\mathbb{E}[\mathcal{G}(u_h^s)]$, where the functional \mathcal{G} is taken to be point evaluation at $1/3$, i.e.

$$F(\mathbf{y}) = \mathcal{G}(u_h^s(\cdot, \mathbf{y})) = u_h^s(1/3, \mathbf{y}).$$

To specify a^s , we choose here the Matérn class of covariances defined in (2.2) and (2.5) with $a_* \equiv 0$ and $a_0 \equiv 1$ in (1.2). To compute a^s we use the formula (3.4) which requires computation of the eigenpairs (μ_j, ξ_j) for $1 \leq j \leq s$. We do this by discretizing the integral operator in (1.5) using the Nyström method based on Gauss-Legendre quadrature on $[0, 1]$ with 10,000 quadrature points and then solving the resulting algebraic eigenvalue problem.

To define the weighted space \mathcal{W}_s in (4.3) and to perform the component-by-component algorithm for calculating the generating vector \mathbf{z} , we must choose the weight parameters γ_u and weight function ψ_j . Thus we need to specify α_j in (4.4) and λ_* from (4.24). In principle, our weighted function space framework in §4 allows us to adjust the QMC rule to the integrand behavior with respect to every coordinate via the j -dependent parameters α_j in (4.4), and we can choose α_j according to (4.28), to minimize the constant in the theoretical error bound. However, as discussed in [27, 28], allowing a different value of α_j for each j would cause a substantial increase in the cost of the component-by-component (CBC) algorithm. To maintain

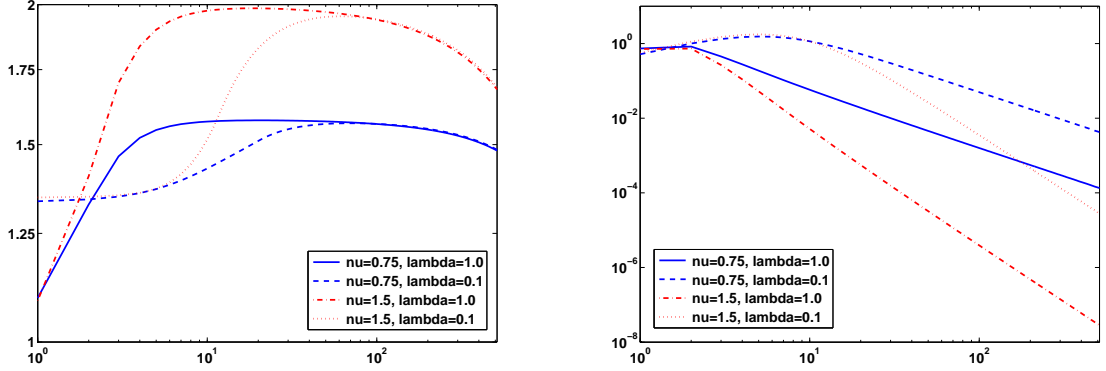


Figure 1: Log-log plots of $\|\xi_j\|_{\mathcal{C}^0(\overline{D})}$ (left) and $\mu_j \|\nabla \xi_j\|_{\mathcal{C}^0(\overline{D})}$ (right) against j for the Matérn covariance in one dimension for various ν and λ_C

the full efficiency of the currently available CBC construction algorithms, the α_j should be coordinate-independent, at least for large blocks of coordinates. In our numerical experiments we found that the use of a single value of α_j for all j led to unsatisfactory results, but that two values (chosen according to the prescription below) led to acceptable results. Noting that as a consequence of Assumption A3 $b_j \rightarrow 0$ as $j \rightarrow \infty$, we choose j_0 to be the smallest positive integer such that $b_j < b_*/2$ for all $j \geq j_0$, where $b_* := \max_{j \geq 1} b_j$. Then following (4.28), we define

$$\alpha_j = \begin{cases} \frac{1}{2} \left(b_* + \sqrt{b_*^2 + 1 - 1/(2\lambda_*)} \right) & \text{for } j < j_0 \\ \frac{1}{2} \left(b_{j_0} + \sqrt{b_{j_0}^2 + 1 - 1/(2\lambda_*)} \right) & \text{for } j \geq j_0 \end{cases},$$

with λ_* to be specified below, which satisfies the general condition (4.9).

The weight parameters γ_u are now defined as in (4.23), after first setting the parameter λ_* . By (4.24), this parameter λ_* is related to p , which in turn was introduced in Assumption A3 and is thus constrained by the choice of covariance function $c(x, x')$. For the Matérn covariance ρ_ν from (2.5), we can observe numerically the boundedness of both the sequences $\|\xi_j\|_{\mathcal{C}^0(\overline{D})}$ and $\mu_j \|\nabla \xi_j\|_{\mathcal{C}^0(\overline{D})}$ (see Figure 1). Using this together with Corollary 5 we can deduce empirically that Assumption A2 holds with $\Theta = 1 + 2\nu$ and $\varepsilon = 0$, and that $b_j = \mathcal{O}(j^{-(1+2\nu)/2})$, and hence that Assumption A3 holds for arbitrary $p \in (\frac{2}{1+2\nu}, 1]$. Using this relationship between ν and p , we find from (4.24) that for arbitrarily small $q > 0$ we can simplify the expression for λ_* to

$$\lambda_* = \begin{cases} \frac{1}{2\nu} + q, & \text{if } 1/2 < \nu < 1 \\ \frac{1}{2} + q, & \text{if } \nu \geq 1. \end{cases} \quad (6.2)$$

This choice of λ_* implies, see Theorem 17, that we obtain theoretical QMC convergence close to $\mathcal{O}(n^{-\min(\nu, 1)})$. Note that the choice of q involves a trade-off. Smaller values of q lead to a faster convergence, but also to a larger value for $C_\gamma(\lambda_*)$ in (4.19). In fact, for $\nu \geq 1$, we easily see that $q \rightarrow 0$ is equivalent to $\delta \rightarrow 0$ in (4.24). This in turn implies $C_\gamma(\lambda_*) \rightarrow \infty$, by way of (4.19) and (4.8). Whereas for $\nu \in (1/2, 1)$, we see that $q \rightarrow 0$ is equivalent to $p \rightarrow 2/(1 + 2\nu)$, so that the sum in Assumption A3 grows without bound, leading, as in the proof of Theorem 20, again to $C_\gamma(\lambda_*) \rightarrow \infty$. Here we choose $q = 0.05$.

Recalling (4.2), we write $Q_i = Q_{s,n}(\Delta_i; F)$, where Δ_i is the i -th independent random shift, uniformly distributed on $[0, 1]^s$. Denoting by \bar{Q} the mean of the Q_i , we have the following

Table 1: QMC standard errors for $\nu = 1.5$ (using POD weights γ_u as in (4.23))

n	$\sigma^2 = 0.25$		$\sigma^2 = 1.0$		$\sigma^2 = 4.0$	
	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$
8,009	2.69e-05	1.77e-05	1.90e-04	1.00e-04	1.12e-02	3.22e-03
16,001	1.38e-05	8.12e-06	1.02e-04	7.52e-05	5.47e-03	2.44e-03
32,003	8.85e-06	6.22e-06	6.79e-05	5.11e-05	3.83e-03	1.20e-03
64,007	4.49e-06	3.02e-06	3.33e-05	3.49e-05	2.36e-03	7.02e-04
120,011	2.66e-06	1.79e-06	2.46e-05	1.79e-05	3.18e-03	7.87e-04
240,007	1.43e-06	9.95e-07	1.48e-05	9.80e-06	1.74e-03	4.42e-04
480,013	7.82e-07	6.72e-07	9.17e-06	8.41e-06	7.68e-04	2.91e-04
Rate	0.86	0.80	0.73	0.66	0.55	0.58
90% Interval	[0.89, 0.83]	[0.87, 0.74]	[0.78, 0.69]	[0.75, 0.57]	[0.71, 0.40]	[0.68, 0.48]

Table 2: QMC standard errors for $\nu = 0.75$ (using POD weights γ_u as in (4.23))

n	$\sigma^2 = 0.25$		$\sigma^2 = 1.0$		$\sigma^2 = 4.0$	
	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$
8,009	2.80e-05	1.80e-05	1.76e-04	1.12e-04	8.97e-03	2.01e-03
16,001	1.37e-05	7.37e-06	1.25e-04	7.10e-05	7.25e-03	1.69e-03
32,003	8.37e-06	5.78e-06	5.72e-05	3.98e-05	2.42e-03	1.26e-03
64,007	4.36e-06	2.93e-06	3.39e-05	2.70e-05	1.72e-03	8.35e-04
120,011	2.58e-06	1.82e-06	2.00e-05	1.82e-05	1.43e-03	5.63e-04
240,007	1.32e-06	9.56e-07	1.14e-05	1.31e-05	1.57e-03	2.64e-04
480,013	7.06e-07	5.57e-07	6.31e-06	7.52e-06	5.60e-04	2.05e-04
Rate	0.89	0.82	0.83	0.64	0.63	0.60
90% Interval	[0.91, 0.86]	[0.89, 0.76]	[0.88, 0.79]	[0.68, 0.61]	[0.81, 0.44]	[0.70, 0.50]

unbiased estimator with R random shifts of the mean-square error (with respect to the shifts):

$$\frac{1}{R} \frac{1}{R-1} \sum_{i=1}^R (Q_i - \bar{Q})^2 \approx \mathbb{E}^\Delta |I_s(F) - Q_{s,n}(\cdot; F)|^2. \quad (6.3)$$

The square-root of the left hand side of (6.3) is an estimate of the ‘‘standard error’’. In the following experiments we estimate this standard error for the following selection of parameters,

$$\nu = 0.75, 1.5 \quad \sigma^2 = 0.25, 1.0, 4.0 \quad \lambda_c = 0.1, 1.0,$$

using $R = 32$ random shifts. Here, σ^2 and λ_c refer to the variance and length-scale parameters for the Matérn covariance in (2.5). Recall further that we have fixed the truncation dimension at $s = 400$ and the spatial resolution at $h = 1/1024$.

Tables 1 and 2 present results using the QMC quadrature analyzed in §4, along with estimated values of the rate-of-convergence parameter r in the error representation cn^{-r} . The rate r is estimated together with its 90% confidence interval by linear regression of the negative log of the standard error against $\log n$. Here we see a strong dependence on the variance σ^2 , but a weaker dependence on the choices of λ_c and ν . While Theorem 20 suggests that the asymptotic behavior of the root-mean-square error depends on p (and hence ν), in practice the observed rates of convergence bear little relation with that prediction. One explanation may be that with the range of n presented we are in a pre-asymptotic regime. This seems especially true for larger values of σ^2 , and hence may explain why we see our QMC quadrature performing similarly to standard Monte Carlo (MC) quadrature for $\sigma^2 = 4.0$. Another explanation is that our predicted rate of convergence is based on a rigorous upper bound which is likely not sharp.

Table 3: MC standard errors for $\nu = 1.5$

n	$\sigma^2 = 0.25$		$\sigma^2 = 1.0$		$\sigma^2 = 4.0$	
	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$
8,009	7.24e-04	4.19e-04	2.21e-03	1.11e-03	2.70e-02	6.62e-03
16,001	3.98e-04	2.58e-04	1.15e-03	7.22e-04	1.42e-02	4.98e-03
32,003	2.97e-04	1.52e-04	9.73e-04	4.45e-04	1.65e-02	3.56e-03
64,007	1.87e-04	1.07e-04	6.21e-04	3.08e-04	1.02e-02	2.43e-03
120,011	1.25e-04	7.59e-05	4.11e-04	2.17e-04	5.78e-03	1.65e-03
240,007	9.40e-05	6.19e-05	2.97e-04	1.50e-04	4.02e-03	8.78e-04
480,013	7.06e-05	4.16e-05	2.12e-04	9.75e-05	2.79e-03	5.06e-04
Rate	0.56	0.55	0.56	0.59	0.55	0.63
90% Interval	[0.62, 0.51]	[0.61, 0.49]	[0.61, 0.50]	[0.61, 0.57]	[0.65, 0.44]	[0.71, 0.55]

Table 4: MC standard errors for $\nu = 0.75$

n	$\sigma^2 = 0.25$		$\sigma^2 = 1.0$		$\sigma^2 = 4.0$	
	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$	$\lambda_c = 1.0$	$\lambda_c = 0.1$
8,009	6.89e-04	4.01e-04	2.05e-03	1.07e-03	2.30e-02	6.54e-03
16,001	3.82e-04	2.47e-04	1.08e-03	6.90e-04	1.21e-02	4.85e-03
32,003	2.81e-04	1.45e-04	9.05e-04	4.23e-04	1.35e-02	3.43e-03
64,007	1.78e-04	1.02e-04	5.76e-04	2.93e-04	8.51e-03	2.36e-03
120,011	1.20e-04	7.20e-05	3.85e-04	2.07e-04	5.00e-03	1.62e-03
240,007	9.91e-05	5.95e-05	3.11e-04	1.52e-04	3.84e-03	1.24e-03
480,013	6.92e-05	4.08e-05	2.11e-04	9.60e-05	3.85e-03	8.15e-04
Rate	0.55	0.55	0.54	0.58	0.45	0.51
90% Interval	[0.61, 0.49]	[0.61, 0.48]	[0.59, 0.48]	[0.61, 0.55]	[0.56, 0.34]	[0.53, 0.49]

Recall from the theory that we expect a convergence rate close to $\mathcal{O}(n^{-\min(\nu,1)})$. We see, however, that the results converge almost as well for $\nu = 0.75$ as for $\nu = 1.5$. This seems to indicate that our theory is not sharp, and that optimal (close to $\mathcal{O}(n^{-1})$) convergence could potentially be demonstrated for ν lower than our current cross-over point of $\nu = 1$ in (6.2). This is also indicated by the fact that our method sometimes converges faster than predicted by the theory, for example for $\nu = 0.75$, $\sigma^2 = 0.25$ and $\lambda_c = 1.0$, where the observed rate of convergence of approximately 0.89 is significantly larger than the predicted rate of 0.75 from Theorem 20.

Tables 3 and 4 present the same experiments as Tables 1 and 2 respectively, but for MC quadrature. The results agree with the usual behavior of MC methods where standard errors converge with approximately $\mathcal{O}(n^{-1/2})$. Figure 2 charts all the findings in Tables 1 to 4. They demonstrate that in all our test cases QMC always does better than MC, especially for small σ^2 , where QMC outperforms MC by up to two orders of magnitude.

As a final comparison, in Tables 5 and 6 we look at the standard errors for a generic lattice rule, which is not specifically designed to fit the problem. We choose here a lattice rule generated for the Sobolev space of mixed first-order derivatives on $[0, 1]^s$, with product weight parameters $\gamma_j = 1/j^2$. We see that these lattice rules still behave very well, attaining similar results to the lattice rules constructed for our specific problem.

7 Conclusion

We have been able to demonstrate good convergence of our QMC finite element method for this class of PDE problems both theoretically and numerically. The success of this project has required the analysis of QMC methods for integrals over an unbounded region, with the weight

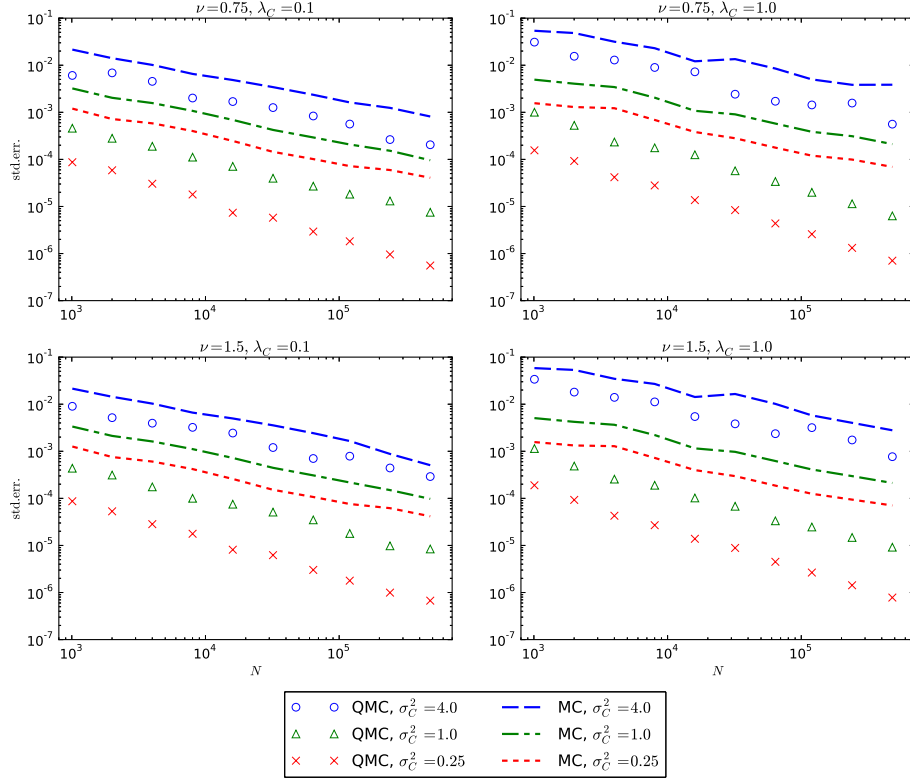


Figure 2: Standard errors from Tables 1 to 4 for QMC and MC plotted against n .

Table 5: QMC standard errors for $\nu = 1.5$ (using generic lattice rules)

n	$\sigma^2 = 0.25$		$\sigma^2 = 1.0$		$\sigma^2 = 4.0$	
	$\lambda_C = 1.0$	$\lambda_C = 0.1$	$\lambda_C = 1.0$	$\lambda_C = 0.1$	$\lambda_C = 1.0$	$\lambda_C = 0.1$
8,009	3.38e-05	1.90e-05	2.48e-04	1.24e-04	1.29e-02	2.81e-03
16,001	1.68e-05	9.91e-06	1.23e-04	7.41e-05	7.00e-03	1.78e-03
32,003	8.27e-06	6.48e-06	6.47e-05	5.56e-05	4.05e-03	1.19e-03
64,007	4.32e-06	4.60e-06	3.55e-05	4.05e-05	2.88e-03	9.36e-04
120,011	2.46e-06	2.01e-06	2.59e-05	2.08e-05	4.33e-03	6.41e-04
240,007	1.77e-06	1.41e-06	2.07e-05	1.31e-05	3.54e-03	3.76e-04
480,013	6.84e-07	6.32e-07	7.10e-06	7.26e-06	7.60e-04	2.27e-04
Rate	0.92	0.80	0.80	0.68	0.52	0.59
90% Interval	[0.99, 0.84]	[0.88, 0.72]	[0.91, 0.68]	[0.76, 0.61]	[0.78, 0.25]	[0.65, 0.54]

Table 6: QMC standard errors for $\nu = 0.75$ (using generic lattice rules)

n	$\sigma^2 = 0.25$		$\sigma^2 = 1.0$		$\sigma^2 = 4.0$	
	$\lambda_C = 1.0$	$\lambda_C = 0.1$	$\lambda_C = 1.0$	$\lambda_C = 0.1$	$\lambda_C = 1.0$	$\lambda_C = 0.1$
8,009	3.25e-05	1.80e-05	2.27e-04	1.21e-04	1.06e-02	2.82e-03
16,001	1.61e-05	9.42e-06	1.12e-04	6.77e-05	5.52e-03	1.66e-03
32,003	7.87e-06	5.92e-06	5.71e-05	5.13e-05	3.36e-03	1.16e-03
64,007	4.06e-06	4.21e-06	3.11e-05	3.64e-05	2.18e-03	8.49e-04
120,011	2.41e-06	1.81e-06	2.34e-05	1.84e-05	3.13e-03	5.77e-04
240,007	1.69e-06	1.32e-06	1.81e-05	1.25e-05	2.40e-03	3.71e-04
480,013	6.60e-07	5.96e-07	6.69e-06	7.19e-06	6.44e-04	2.40e-04
Rate	0.91	0.81	0.80	0.68	0.54	0.58
90% Interval	[0.99, 0.84]	[0.88, 0.73]	[0.91, 0.69]	[0.74, 0.61]	[0.76, 0.31]	[0.61, 0.55]

parameters γ_u and the weight functions ψ_j both tuned to the specific problem. This led to the extension of shifted lattice rule theory presented in [28].

Our numerical results demonstrate that QMC rules comfortably beat MC rules in most cases, or certainly perform no worse in the cases of large σ^2 . Furthermore we see that this is the case even for arbitrarily chosen lattice rules, as is demonstrated in Tables 5 and 6, despite the fact that the theory for these lattice rules does not apply to this problem.

It is important to note that in these experiments neither the MC nor the QMC rules are enhanced using any variance reduction techniques such as the use of antithetic variates. This way we have a fair comparison between two unflavored implementations. Evidently there is scope for further work to sharpen our error bounds, as demonstrated by our numerics. Nevertheless, the results show that QMC finite element methods provide an excellent solution to the lognormal porous flow problem, and present a marked improvement over MC methods.

Acknowledgements. The authors would like to thank Tony Shardlow and Helmut Harbrecht for fruitful discussions that lead to Corollary 5 and Proposition 9. We would also like to acknowledge the financial support of the Australian Research Council, of the UK Engineering and Physical Sciences Research Council and of the European Research Council under grant ERC AdG247277.

References

- [1] R.J. Adler, *The Geometry of Random Fields*, Wiley, London, 1981.
- [2] V.I. Bogachev, *Gaussian Measures*, AMS Monographs Vol. **62**, American Mathematical Society, R.I., USA (1998).
- [3] J. Charrier, Strong and weak error estimates for elliptic partial differential equations with random coefficients, *SIAM J. Numer. Anal.*, **50**, 216–246 (2012).
- [4] J. Charrier, and A. Debussche, Weak truncation error estimates for elliptic PDEs with lognormal coefficients, *Stoch. PDEs: Anal. Comput.*, **1**, 63–93 (2013).
- [5] J. Charrier, R. Scheichl, and A.L. Teckentrup, Finite element error analysis of elliptic PDEs with random coefficients and its application to multilevel Monte Carlo methods, *SIAM J. Numer. Anal.*, **51**, 322–352 (2013).
- [6] R. Cools, F.Y. Kuo, and D. Nuyens, Constructing embedded lattice rules for multivariate integration, *SIAM J. Sci. Comput.*, **28**, 2162–2188 (2006).
- [7] G. Da Prato and J. Zabczyk, *Stochastic Equations in Infinite Dimensions*, Vol. 44 of *Encyclopedia of Mathematics and its Applications*. Cambridge University Press, 1992.
- [8] J. Dick, On the convergence rate of the component-by-component construction of good lattice rules, *J. Complexity*, **20**, 493–522 (2004).
- [9] J. Dick, F.Y. Kuo, and I.H. Sloan, *High dimensional integration – the quasi-Monte Carlo way*, *Acta Numer.* **22**, 133–288 (2013).
- [10] J. Dick and F. Pillichshammer, *Digital Nets and Sequences. Discrepancy Theory and Quasi-Monte Carlo Integration*, Cambridge, Cambridge University Press, 2010.
- [11] J. Dick, F. Pillichshammer, and B.J. Waterhouse, The construction of good extensible rank-1 lattices, *Math. Comp.*, **77**, 2345–2374 (2008).
- [12] P. Frauenfelder, Ch. Schwab, and R.A. Todor, Finite elements for elliptic problems with stochastic coefficients, *Compu. Methods Appl. Mech. Engrg.*, **194**, 205–228 (2005).
- [13] R.G. Ghanem and P.D. Spanos, *Stochastic Finite Elements*, Dover, 1991.

- [14] C.J. Gittelsohn, Stochastic Galerkin discretization of the isotropic lognormal diffusion problem, *Math. Mod. Meth. Appl. Sci.*, **20**, 237–263 (2010).
- [15] I.G. Graham, F.Y. Kuo, D. Nuyens, R. Scheichl, and I.H. Sloan, Quasi-Monte Carlo methods for elliptic PDEs with random coefficients and applications, *J. Comput. Phys.*, **230**, 3668–3694 (2011).
- [16] F.J. Hickernell, I.H. Sloan, and G.W. Wasilkowski, On tractability of weighted integration for certain Banach spaces of functions, in *Monte Carlo and Quasi-Monte Carlo Methods 2002* (H. Niederreiter, ed.), Springer-Verlag, Berlin, 2004, pp. 51–71.
- [17] V.-H. Hoang and Ch. Schwab, N -term Wiener chaos approximation rates for elliptic PDEs with lognormal Gaussian random inputs, *Math. Mod. Meth. Appl. Sci.*, **24**, 797–826 (2014).
- [18] F.Y. Kuo, Component-by-component constructions achieve the optimal rate of convergence, *J. Complexity*, **19**, 301–320 (2003).
- [19] F.Y. Kuo, Ch. Schwab, and I. H. Sloan, Quasi-Monte Carlo finite element methods for a class of elliptic partial differential equations with random coefficient, *SIAM J. Numer. Anal.*, **50**, 3351–3374 (2012).
- [20] F.Y. Kuo, Ch. Schwab, and I.H. Sloan, Quasi-Monte Carlo methods for very high dimensional integration: the standard weighted-space setting and beyond, *ANZIAM J.*, **53**, 1–37 (2011).
- [21] F.Y. Kuo, I.H. Sloan, G.W. Wasilkowski and B.J. Waterhouse, Randomly shifted lattice rules with the optimal rate of convergence for unbounded integrands, *J. Complexity*, **26**, 135–160 (2010).
- [22] F.Y. Kuo, G.W. Wasilkowski, and B.J. Waterhouse, Randomly shifted lattice rules for unbounded integrals, *J. Complexity* **22**, 630–651 (2006).
- [23] G.J. Lord, C.E. Powell and T. Shardlow, *An Introduction to Computational Stochastic PDEs*, Cambridge University Press, Cambridge, 2014.
- [24] W. McLean, *Strongly Elliptic Systems and Boundary Integral Equations*, Cambridge University Press, Cambridge, 2000.
- [25] R.L. Naff, D.F. Haley, and E.A. Sudicky, High-resolution Monte Carlo simulation of flow and conservative transport in heterogeneous porous media 1. Methodology and flow results, *Water Resour. Res.*, **34**, 663–677 (1998).
- [26] R.L. Naff, D.F. Haley, and E.A. Sudicky, High-resolution Monte Carlo simulation of flow and conservative transport in heterogeneous porous media 2. Transport Results, *Water Resour. Res.*, **34**, 679–697 (1998).
- [27] J.A. Nichols, Quasi-Monte Carlo methods with applications to partial differential equations with random coefficients, PhD Thesis, University of New South Wales, 2014.
- [28] J.A. Nichols and F.Y. Kuo, Fast CBC construction of randomly shifted lattice rules achieving $\mathcal{O}(N^{-1+\delta})$ convergence for unbounded integrands in \mathbb{R}^s in weighted spaces with POD weights, *J. Complexity*, **30**, 444–468 (2014).
- [29] D. Nuyens and R. Cools, Fast algorithms for component-by-component construction of rank-1 lattice rules in shift-invariant reproducing kernel Hilbert spaces, *Math. Comp.*, **75**, 903–920 (2006).
- [30] D. Nuyens and R. Cools, Fast component-by-component construction of rank-1 lattice rules with a non-prime number of points, *J. Complexity*, **22**, 4–28 (2006).
- [31] Ch. Schwab and C.J. Gittelsohn, Sparse Tensor Discretizations of High Dimensional and Stochastic PDEs, *Acta Numerica*, **20**, 291–467 (2011).
- [32] Ch. Schwab and R.A. Todor, Karhunen-Loève approximation of random fields by generalized fast multipole methods, *J. Comput. Phys.*, **217**, 100–122 (2006).

- [33] I.H. Sloan, F.Y. Kuo, and S. Joe, Constructing randomly shifted lattice rules in weighted Sobolev spaces, *SIAM J. Numer. Anal.*, **40**, 1650–1665 (2002).
- [34] I.H. Sloan and H. Woźniakowski, When are quasi-Monte Carlo algorithms efficient for high-dimensional integrals?, *J. Complexity*, **14**, 1–33 (1998).
125, 569-600 (2013).
- [35] I.H. Sloan, X. Wang, and H. Woźniakowski, Finite-order weights imply tractability of multivariate integration, *J. Complexity*, **20**, 46–74 (2004).
- [36] A.L. Teckentrup, R. Scheichl, M.B. Giles, and E. Ullmann, Further analysis of multilevel Monte Carlo methods for elliptic PDEs with random coefficients, *Numer. Math.*, **125**, 569–600 (2013).
- [37] G.W. Wasilkowski and H. Woźniakowski, Complexity of weighted approximation over \mathbb{R}^1 , *J. Approx. Theory.*, **103**, 223–251 (2000).
- [38] G.W. Wasilkowski and H. Woźniakowski, Tractability of approximation and integration for weighted tensor product problems over unbounded domains, in *Monte Carlo and Quasi-Monte Carlo Methods 2000* (K.-T. Fang, F.J. Hickernell, H. Niederreiter, eds.), Springer, Berlin, 2002, pp. 497–522.
- [39] H. Widom, Asymptotic behaviour of the eigenvalues of certain integral operators, *Transactions of the American Mathematical Society*, **109**, 278–295 (1963).