

# Quasi-Monte Carlo methods for PDE uncertainty quantification

Vesa Kaarnioja  
[vesa.kaarnioja@fu-berlin.de](mailto:vesa.kaarnioja@fu-berlin.de)

Winter School of SIAM/GAMM Student Chapter Berlin

February 21–22, 2024

# Practical matters

- Lecture notes and MATLAB programs are available at  
<https://www.iki.fi/vesakaar/winterschool24>
- Schedule:

	<b>Wednesday</b>	<b>Thursday</b>
09:10-09:30	Coffee/Opening	Coffee
09:30-10:50	V. Kaarnioja	D. Walter
10:50-11:10	Coffee Break	Coffee Break
11:10-12:30	V. Kaarnioja	D. Walter
12:30-13:30	Lunch Break	Lunch Break
13:30-14:50	D. Walter	V. Kaarnioja
14:50-15:10	Coffee Break	Coffee Break
15:10-16:30	D. Walter	V. Kaarnioja

# Uncertainty in groundwater flow

Risk analysis of radwaste disposal or CO<sub>2</sub> sequestration.

Darcy's law

$$\mathbf{q}(\mathbf{x}) + \mathbf{a}(\mathbf{x}) \nabla p(\mathbf{x}) = \mathbf{f}(\mathbf{x})$$

mass conservation law

$$\nabla \cdot \mathbf{q}(\mathbf{x}) = 0$$

in  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$

together with boundary conditions

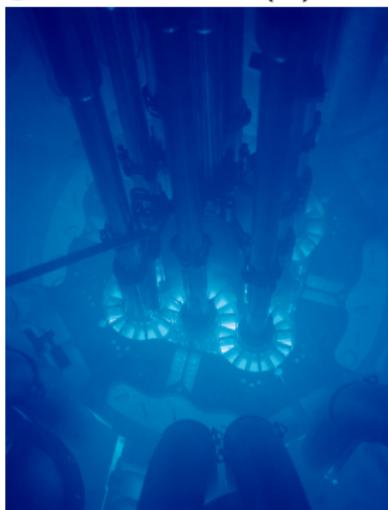


Uncertainty in  $\mathbf{a}(\mathbf{x}, \omega)$  leads to uncertainty in  $\mathbf{q}(\mathbf{x}, \omega)$  and  $p(\mathbf{x}, \omega)$

## Criticality problem for nuclear reactors

$$-\nabla \cdot (\underbrace{a(\mathbf{x})}_{\text{diffusion}} \nabla u(\mathbf{x})) + \underbrace{b(\mathbf{x})}_{\text{absorption}} u(\mathbf{x}) = \lambda \underbrace{c(\mathbf{x})}_{\text{fission}} u(\mathbf{x})$$

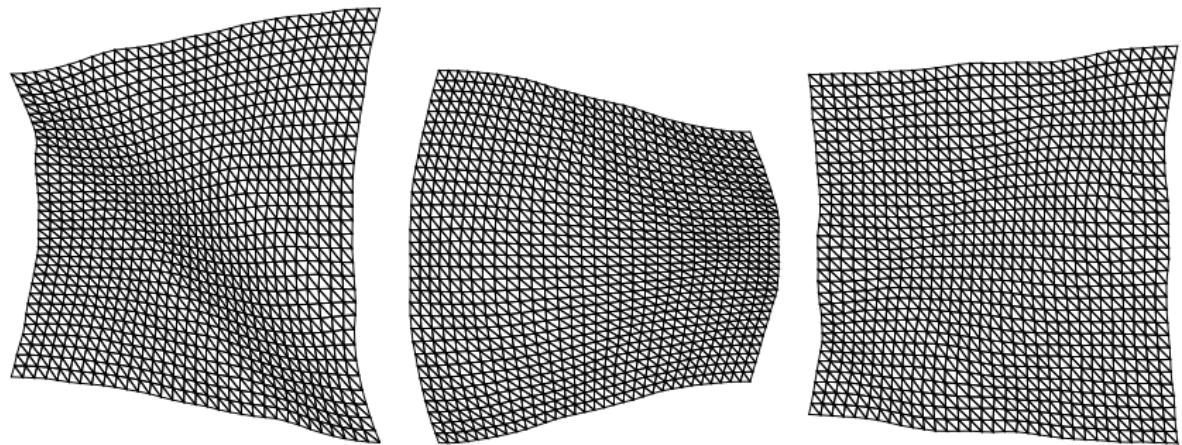
- The smallest eigenvalue  $\lambda_1 \in \mathbb{R}$  measures *criticality* of a reactor.
- Eigenfunction  $u_1(\mathbf{x})$  is the *neutron flux* at the point  $\mathbf{x}$ .



- $\lambda_1 \approx 1 \Rightarrow$  operating efficiently
- $\lambda_1 > 1 \Rightarrow$  not self-sustaining
- $\lambda_1 < 1 \Rightarrow$  supercritical

Source: Argonne National  
Laboratory on Flickr

# Domain uncertainty quantification

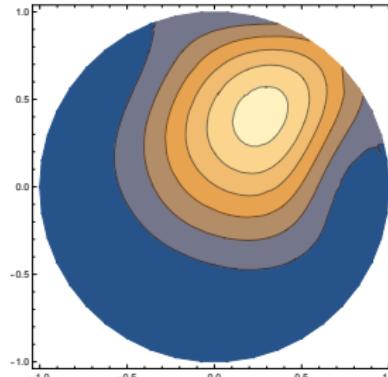
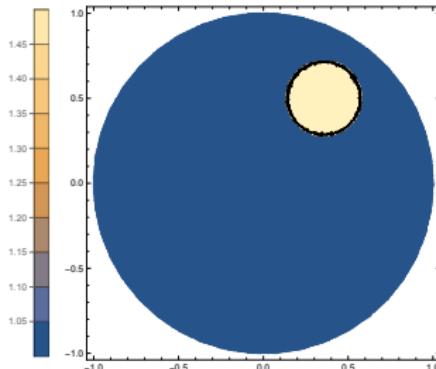
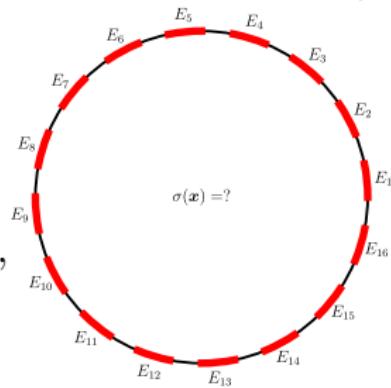


Three realizations of a random spatial domain

# Electrical impedance tomography

Use measurements of current and voltage collected at electrodes covering part of the boundary to infer the interior conductivity of an object/body.

$$\begin{cases} \nabla \cdot (\sigma \nabla u) = 0 & \text{in } D, \\ \sigma \frac{\partial u}{\partial \mathbf{n}} = 0 & \text{on } \partial D \setminus \bigcup_{k=1}^L \overline{E_k}, \\ u + z_k \sigma \frac{\partial u}{\partial \mathbf{n}} = U_k & \text{on } E_k, \ k \in \{1, \dots, L\}, \\ \int_{E_k} \sigma \frac{\partial u}{\partial \mathbf{n}} dS = I_k, & k \in \{1, \dots, L\}, \end{cases}$$



Consider the elliptic PDE problem:

$$\begin{cases} -\nabla \cdot (a(\mathbf{x}) \nabla u(\mathbf{x})) = f(\mathbf{x}) & \text{for } \mathbf{x} \in D, \\ +\text{boundary conditions.} \end{cases}$$

In practice, one or several of the material/system parameters may be uncertain or incompletely known and modeled as random fields:

- PDE coefficient  $a$  may be uncertain;
- Source term  $f$  may be uncertain;
- Boundary conditions may be uncertain;
- The domain  $D$  itself may be uncertain.

In forward uncertainty quantification, one is interested in assessing how uncertainties in the inputs of a mathematical model affect the output.

⇒ If the uncertain inputs are modeled as random fields, then the output of the PDE is also a random field. One may be interested in assessing the statistical response of the system, for example, the expectation or variance of the PDE solution (or some other quantity of interest thereof).

# High-dimensional numerical integration

$$\int_{[0,1]^s} f(\mathbf{y}) \, d\mathbf{y} \approx \sum_{i=1}^n w_i f(\mathbf{t}_i)$$

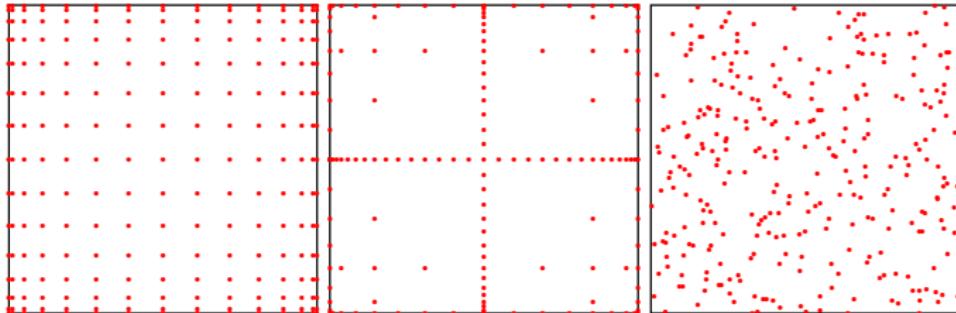


Figure: Tensor product grid, sparse grid, Monte Carlo nodes (not QMC rules)

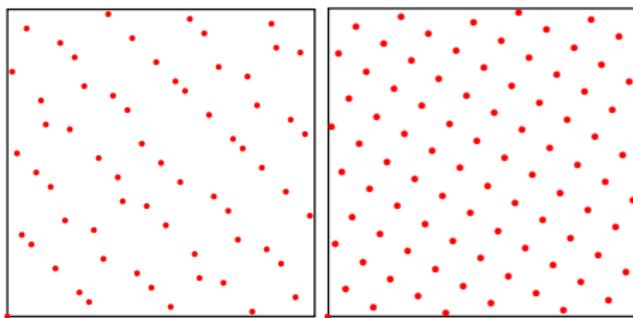


Figure: Sobol' points, lattice rule (examples of QMC rules)

*Quasi-Monte Carlo (QMC) methods* are a class of *equal weight* cubature rules

$$\int_{[0,1]^s} f(\mathbf{y}) d\mathbf{y} \approx \frac{1}{n} \sum_{i=1}^n f(\mathbf{t}_i),$$

where  $(\mathbf{t}_i)_{i=1}^n$  is an ensemble of *deterministic* nodes in  $[0, 1]^s$ .

The nodes  $(\mathbf{t}_i)_{i=1}^n$  are NOT random! Instead, they are *deterministically chosen*.

QMC methods exploit the smoothness and anisotropy of an integrand in order to achieve better-than-Monte Carlo rates.

# Table of contents

1. Preliminaries
  - elliptic partial differential equations (PDEs), Galerkin method, modeling random fields, elliptic PDEs with random coefficients
2. Quasi-Monte Carlo (QMC) methods
  - randomly shifted rank-1 lattice rules, weighted Sobolev spaces, error analysis
3. Constructing lattice rules
  - Naïve component-by-component (CBC) construction
  - Fast CBC algorithm
4. QMC methods for forward and inverse uncertainty quantification of elliptic PDEs with random coefficients
  - affine and uniform setting
  - lognormal setting (briefly)
  - some perspectives on applying QMC for Bayesian inverse problems

General reference on QMC:

-  J. Dick, F. Y. Kuo, and I. H. Sloan. High-dimensional integration: The quasi-Monte Carlo way. *Acta Numer.* **22**:133–288, 2013.  
<https://doi.org/10.1017/S0962492913000044>

QMC for PDE uncertainty quantification:

-  F. Y. Kuo and D. Nuyens. Application of quasi-Monte Carlo methods to elliptic PDEs with random diffusion coefficients: a survey of analysis and implementation. *Found Comput. Math.* **16**:1631–1696, 2016. <https://doi.org/10.1007/s10208-016-9329-5>
-  F. Y. Kuo and D. Nuyens. Application of quasi-Monte Carlo methods to PDEs with random coefficients – an overview and tutorial. In: A. B. Owen, P. W. Glynn (eds.), *Monte Carlo and Quasi-Monte Carlo Methods 2016*, 53–71, 2018.  
[https://doi.org/10.1007/978-3-319-91436-7\\_3](https://doi.org/10.1007/978-3-319-91436-7_3)

## QMC for Bayesian inverse problems:

-  R. Scheichl and A. M. Stuart and A. L. Teckentrup. Quasi-Monte Carlo and multilevel Monte Carlo methods for computing posterior expectations in elliptic inverse problems. *SIAM/ASA J. Uncertain. Quantif.* **5**(1):493–518, 2017.  
<https://doi.org/10.1137/16M1061692>
-  J. Dick, R. N. Gantner, Q. T. Le Gia, and Ch. Schwab. Higher order quasi-Monte Carlo integration for Bayesian PDE inversion. *Comput. Math. Appl.* **77**(1):144–172, 2019.  
<https://doi.org/10.1016/j.camwa.2018.09.019>
-  R. N. Gantner. *Computational Higher-Order Quasi-Monte Carlo for Random Partial Differential Equations*. PhD thesis, ETH Zürich, 2017.  
<https://doi.org/10.3929/ethz-b-000182695>
-  L. Herrmann, M. Keller, and Ch. Schwab. Quasi-Monte Carlo Bayesian estimation under Besov priors in elliptic inverse problems. *Math. Comp.* **90**:1831–1860, 2021. <https://www.ams.org/journals/mcom/2021-90-330/S0025-5718-2021-03615-7>

## 1. Preliminaries

## Elliptic PDE

Many physical phenomena can be modeled using elliptic partial differential equations of the form

$$\begin{cases} -\nabla \cdot (a(\mathbf{x}) \nabla u(\mathbf{x})) = f(\mathbf{x}), & \mathbf{x} \in D, \\ +\text{boundary conditions} \end{cases}$$

Uncertainties can appear in the material parameter  $a$ , source term  $f$ , boundary conditions, or the domain  $D$ .

- For the purposes of analysis, we consider the weak formulation of the PDE. Under certain conditions, the solution to the weak formulation can be shown to exist and be uniquely defined.
- When we solve the PDE numerically using the finite element method, we are actually approximating the solution to the *the weak formulation* of the PDE problem.
- Under suitably strong regularity assumptions ( $D$  convex Lipschitz domain,  $f \in L^2(D)$ , and  $a$  Lipschitz), the weak solution satisfies  $-\nabla \cdot (a(\mathbf{x}) \nabla u(\mathbf{x})) = f(\mathbf{x})$  for a.e.  $\mathbf{x} \in D$  with  $u|_{\partial D} = 0$ .

Let  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be a nonempty open set.

$$L^2(D) := \{v: D \rightarrow \mathbb{R} \mid v \text{ is measurable}, \|v\|_{L^2(D)} := (\int_D |v(x)|^2 dx)^{1/2} < \infty\},$$

$$H^1(D) := \{v \in L^2(D) \mid \partial_j v \in L^2(D) \text{ for all } j \in \{1, \dots, d\}\},$$

$$\text{with } \|v\|_{H^1(D)} := (\|v\|_{L^2(D)}^2 + \|\nabla v\|_{L^2(D)}^2)^{1/2},$$

$$C_0^\infty(D) := \{v \in C^\infty(D) \mid \text{supp}(v) \subset D \text{ is a compact set}\},$$

$$\text{where } \text{supp}(v) := \overline{\{x \in D \mid v(x) \neq 0\}},$$

$$H_0^1(D) := \text{cl}_{H^1(D)}(C_0^\infty(D)),$$

$$H^{-1}(D) := H_0^1(D)' := \{A: H_0^1(D) \rightarrow \mathbb{R} \mid A \text{ linear and bounded}\}.$$

The spaces  $L^2(D)$ ,  $H^1(D)$ ,  $H_0^1(D)$ , and  $H^{-1}(D)$  are Hilbert spaces.

The duality pairing  $\langle f, v \rangle_{H^{-1}(D), H_0^1(D)}$  represents the bounded, linear functional  $v \mapsto f(v)$  for  $f \in H^{-1}(D)$  and  $v \in H_0^1(D)$ .

**Poincaré's inequality:** if  $D \subset \mathbb{R}^d$  is a bounded domain, then there exists a constant  $C_P > 0$  (depending on the domain  $D$ ) such that

$$\|v\|_{L^2(D)} \leq C_P \|\nabla v\|_{L^2(D)} \quad \text{for all } v \in H_0^1(D).$$

Therefore, we can define an equivalent norm in  $H_0^1(D)$  by setting

$$\|v\|_{H_0^1(D)} := \|\nabla v\|_{L^2(D)}.$$

This induces exactly the same topology in  $H_0^1(D)$  as the usual Sobolev norm  $\|\cdot\|_{H^1(D)}$ .

## Trace theorem and boundary values

**Lipschitz domain:** A nonempty domain  $D \subset \mathbb{R}^d$  is said to be a Lipschitz domain if, for every  $\mathbf{x} \in \partial D$ , there exists a rigid transformation (i.e., a rotation plus a translation)  $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , a radius  $r > 0$ , and a Lipschitz function  $\xi: \mathbb{R}^{d-1} \rightarrow \mathbb{R}$  such that

$$f(D) \cap B(f(\mathbf{x}), r) = \{(y_1, \dots, y_d) \in \mathbb{R}^d \mid y_d < \xi(y_1, \dots, y_{d-1})\} \cap B(f(\mathbf{x}), r),$$

where  $B(\mathbf{x}, r)$  denotes the open ball of radius  $r$  centered at  $\mathbf{x}$ .

**Trace theorem:** Let  $D$  be a bounded Lipschitz domain. Then the trace operator

$$\gamma: C^\infty(\overline{D}) \rightarrow C^\infty(\partial D), \quad \gamma u = u|_{\partial D},$$

has a unique extension to a bounded linear operator  $\gamma: H^1(D) \rightarrow L^2(\partial D)$ .

This means that even though  $u \in H^1(D)$  is not well-defined over a set of measure zero, we can interpret its restriction to the boundary of the domain  $D$  as the trace  $\gamma u \in L^2(\partial D)$ .

Especially, Sobolev functions  $u \in H^1(D)$  with zero trace are precisely the elements of  $H_0^1(D)$ :

$$u \in H_0^1(D) \iff \gamma u = 0: \partial D \rightarrow \mathbb{R}.$$

**Q:** How to solve such PDE problems in practice?

**A:** We consider the weak formulation of the PDE problem: given  $f \in H^{-1}(D)$ , find  $u \in H_0^1(D)$  such that

$$\underbrace{\int_D a(\mathbf{x}) \nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x}}_{=:B(u,v)} = \langle f, v \rangle_{H^{-1}(D), H_0^1(D)} \quad \text{for all } v \in H_0^1(D). \quad (1)$$

If there exist  $a_{\min}, a_{\max} > 0$  s.t.  $0 < a_{\min} \leq a(\mathbf{x}) \leq a_{\max} < \infty$  for all  $\mathbf{x} \in D$ , then the bilinear form  $B: H_0^1(D) \times H_0^1(D) \rightarrow \mathbb{R}$  is bounded, i.e.,

$$|B(u, v)| = \left| \int_D a(\mathbf{x}) \nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x} \right| \leq a_{\max} \|u\|_{H_0^1(D)} \|v\|_{H_0^1(D)}$$

for all  $u, v \in H_0^1(D)$ , and coercive, i.e.,

$$B(u, u) = \left| \int_D a(\mathbf{x}) \nabla u(\mathbf{x}) \cdot \nabla u(\mathbf{x}) \, d\mathbf{x} \right| \geq a_{\min} \|u\|_{H_0^1(D)}^2 \quad \text{for all } u \in H_0^1(D),$$

then the *Lax–Milgram lemma* ensures that there exists a unique solution  $u \in H_0^1(D)$  to (1).

## Galerkin method

To solve the system approximately, let  $V_m \subset H_0^1(D)$  be a finite-dimensional subspace of the solution space  $H_0^1(D)$ .

The *Galerkin solution*  $u_m \in V_m$  of the system (1) is the unique solution such that

$$\int_D a(\mathbf{x}) \nabla u_m(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x} = \langle f, v \rangle_{H^{-1}(D), H_0^1(D)} \quad \text{for all } v \in V_m.$$

Let  $V_m$  be spanned by  $\phi_1, \dots, \phi_m$ . We can write the solution as  $u_m = \sum_{i=1}^m c_i \phi_i$ . The above system reduces to the linear system of equations

$$\begin{bmatrix} \int_D a(\mathbf{x}) \nabla \phi_1(\mathbf{x}) \cdot \nabla \phi_1(\mathbf{x}) \, d\mathbf{x} & \cdots & \int_D a(\mathbf{x}) \nabla \phi_1(\mathbf{x}) \cdot \nabla \phi_m(\mathbf{x}) \, d\mathbf{x} \\ \vdots & \ddots & \vdots \\ \int_D a(\mathbf{x}) \nabla \phi_m(\mathbf{x}) \cdot \nabla \phi_1(\mathbf{x}) \, d\mathbf{x} & \cdots & \int_D a(\mathbf{x}) \nabla \phi_m(\mathbf{x}) \cdot \nabla \phi_m(\mathbf{x}) \, d\mathbf{x} \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix} = \begin{bmatrix} \langle f, \phi_1 \rangle_{H^{-1}(D), H_0^1(D)} \\ \vdots \\ \langle f, \phi_m \rangle_{H^{-1}(D), H_0^1(D)} \end{bmatrix}.$$

Solving this system and plugging the expansion coefficients back into the expression for  $u_m$  yields the Galerkin solution.

## Céa's lemma

The solution to the Galerkin system is quasi-optimal in the following sense:

$$\|u - u_m\|_{H_0^1(D)} \leq \frac{a_{\max}}{a_{\min}} \inf_{v_m \in V_m} \|u - v_m\|_{H_0^1(D)}.$$

That is, the  $H_0^1(D)$  error between the true PDE solution  $u$  and the Galerkin approximation  $u_m$  differs from the *optimal approximation* in  $V_m$  up to a constant factor.

## Finite element method

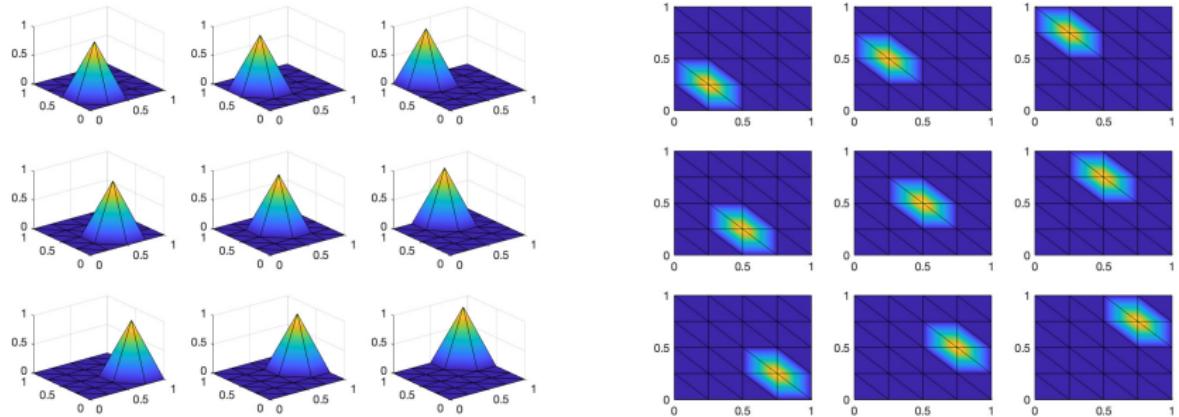
The finite element method is a particular method of constructing the finite-dimensional subspaces  $V_m$  of the solution space  $H_0^1(D)$ .

- Construct a triangulation for the computational domain  $D$ .
- The space  $V_m$  is spanned by piecewise linear functions  $\phi_1, \dots, \phi_m$  which are constructed to satisfy

$$\phi_i(\mathbf{n}_j) = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise,} \end{cases}$$

where  $\mathbf{n}_1, \dots, \mathbf{n}_m$  are the *interior* nodes of the triangulation.

- The finite element solution can be written as  $u_m(\mathbf{x}) = \sum_{i=1}^m c_i \phi_i(\mathbf{x}) \in V_m$ , where the expansion coefficients are solved from the Galerkin system. Note that  $u_m(\mathbf{n}_j) = c_j$ .
- If  $v_m(\mathbf{x}) = \sum_{i=1}^m c_i \phi_i(\mathbf{x}) \in V_m$ , then, e.g.,  $\|v_h\|_{L^2(D)} = \sqrt{\mathbf{c}^T M \mathbf{c}}$ , where  $\mathbf{c} := [c_1, \dots, c_m]^T$  and  $M = [M_{i,j}]_{i,j=1}^m$  is the mass matrix defined elementwise by  $M_{i,j} := \int_D \phi_i(\mathbf{x}) \phi_j(\mathbf{x}) d\mathbf{x}$ ,  $i, j \in \{1, \dots, m\}$ .



**Figure:** Left: An illustration of global, piecewise linear FE basis functions spanning  $V_m$  over a regular, uniform triangulation of  $(0, 1)^2$ . Right: Bird's-eye view of the same global FE basis functions.

# Random field

## Definition

Let  $D \subset \mathbb{R}^d$  and let  $(\Omega, \mathcal{F}, \mu)$  be a probability space. A function  $A: D \times \Omega \rightarrow X$  is called a *random field* if  $A(x, \cdot)$  is an  $X$ -valued random variable for all  $x \in D$ .

## Definition

We call a random field  $A: D \times \Omega \rightarrow X$  square-integrable if

$$\int_{\Omega} |A(x, \omega)|^2 \mu(d\omega) < \infty \quad \text{for all } x \in D.$$

Our goal will be to model (infinite-dimensional) input random fields using finite-dimensional expansions with  $s$  variables.

*Comment on notation:* In what follows,  $s$  will always refer to the “stochastic dimension” (dimension of the stochastic/parametric space) while  $d$  will refer to the “spatial dimension” (dimension of the spatial Lipschitz domain  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ ).

## Mercer's theorem

Let  $a(x, \omega)$  be a square-integrable random field with mean

$$\bar{a}(x) = \int_{\Omega} a(x, \omega) \mu(d\omega), \quad x \in D,$$

and a continuous, symmetric, positive definite<sup>†</sup> covariance

$$K(x, x') = \int_{\Omega} (a(x, \omega) - \bar{a}(x))(a(x', \omega) - \bar{a}(x')) \mu(d\omega).$$

**Mercer's theorem:** the covariance operator  $\mathcal{C}: L^2(D) \rightarrow L^2(D)$ ,

$$(\mathcal{C}u)(x) = \int_D K(x, x') u(x') dx', \quad x \in D,$$

has a countable sequence of eigenvalues  $\{\lambda_k\}_{k \geq 1}$  and corresponding eigenfunctions  $\{\psi_k\}_{k \geq 1}$  satisfying  $\mathcal{C}\psi_k = \lambda_k \psi_k$  such that  $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$  and  $\lambda_k \rightarrow 0$  and the eigenfunctions form an orthonormal basis for  $L^2(D)$ .

Note that this means:

$$\int_D K(x, x') \psi_k(x') dx' = \lambda_k \psi_k(x), \quad \int_D \psi_k(x) \psi_\ell(x) dx = \delta_{k,\ell}.$$

---

<sup>†</sup>In this context, positive definite means: for all choices of finitely many points  $x_1, \dots, x_k \in D$ ,  $k \in \mathbb{N}$ , the Gram matrix  $G := [K(x_i, x_j)]_{i,j=1}^k$  is positive semidefinite.

# The Karhunen–Loève (KL) expansion of a random field

## Theorem

Let  $(\Omega, \mathcal{F}, \mu)$  be a probability space, let  $D \subset \mathbb{R}^d$  be closed and bounded, and let  $a: D \times \Omega \rightarrow \mathbb{R}$  be a square-integrable random field with continuous, symmetric, positive definite covariance

$K(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(a(\mathbf{x}, \cdot) - \bar{a}(\mathbf{x}))(a(\mathbf{x}', \cdot) - \bar{a}(\mathbf{x}'))]$ . Then the eigensystem  $(\lambda_k, \psi_k) \in \mathbb{R}_+ \times L^2(D)$  of the covariance operator  $\mathcal{C}: L^2(D) \rightarrow L^2(D)$ , as described on the previous slide, can be used to write

$$a(\mathbf{x}, \omega) = \bar{a}(\mathbf{x}) + \sum_{k=1}^{\infty} \sqrt{\lambda_k} \xi_k(\omega) \psi_k(\mathbf{x}),$$

$$\text{where } \xi_k(\omega) = \frac{1}{\sqrt{\lambda_k}} \int_D (a(\mathbf{x}, \omega) - \bar{a}(\mathbf{x})) \psi_k(\mathbf{x}) d\mathbf{x},$$

where the convergence is in  $L^2$  w.r.t. the stochastic parameter and uniform in  $\mathbf{x}$ . Furthermore, the random variables  $\xi_k$  are zero-mean uncorrelated random variables with unit variance, i.e.,

$$\mathbb{E}[\xi_k] = 0 \quad \text{and} \quad \mathbb{E}[\xi_k \xi_\ell] = \delta_{k,\ell}.$$

The Karhunen–Loève (KL) expansion of random field  $a(\mathbf{x}, \omega)$  can be written as

$$a(\mathbf{x}, \omega) = \bar{a}(\mathbf{x}) + \sum_{k=1}^{\infty} \sqrt{\lambda_k} \xi_k(\omega) \psi_k(\mathbf{x}).$$

- The KL expansion minimizes the mean-square truncation error:  
$$\|a(\mathbf{x}, \omega) - \bar{a}(\mathbf{x}) - \sum_{k=1}^s \sqrt{\lambda_k} \xi_k(\omega) \psi_k(\mathbf{x})\|_{L^2(\Omega, \mu; L^2(D))} = \left( \sum_{k=s+1}^{\infty} \lambda_k \right)^{1/2}.$$
- The random variables  $\xi_k$  are centered and uncorrelated, but not necessarily independent.
- If the random field  $a(\mathbf{x}, \omega)$  is Gaussian – by definition, this means that  $(a(\mathbf{x}_1, \omega), \dots, a(\mathbf{x}_k, \omega))$  is a multivariate Gaussian random variable for all  $\mathbf{x}_1, \dots, \mathbf{x}_k \in D$ ,  $k \in \mathbb{N}$  – then the random variables  $\xi_k$  are independent.
- The KL expansion is an effective method of representing *input* random fields when their covariance structure is known. If the (infinite-dimensional) input random field has a known covariance (which satisfies the conditions of Mercer's theorem), then we can use the KL expansion to find a finite-dimensional approximation, optimal in the mean-square error sense.

## Modeling assumptions

In engineering and practical applications, the idea is that we have some *a priori* knowledge/belief that the unknown input random field is distributed according to some known probability distribution with a known covariance.

- If the input random field is Gaussian with a known, nice covariance function<sup>†</sup>, then we use the KL expansion to find a reasonable finite-dimensional approximation of the true input. Since the KL expansion decorrelates the stochastic variables, and uncorrelated jointly Gaussian random variables are independent, we can exploit the independence of the stochastic variables to parameterize the model problem.
- If the input random field is *not Gaussian*, then the stochastic variables in the KL expansion are uncorrelated *but not necessarily independent*. For the purposes of mathematical analysis, we typically assume that the random variables in the input random field are independent so that we can parameterize the model problem. (Transforming dependent random variables into independent random variables can be done using, e.g., the Rosenblatt transformation, but this is computationally expensive.)

Note especially that in the Gaussian setting we do not need to make any “extra” effort to ensure the independence of the stochastic variables in the KL expansion.

---

<sup>†</sup>Matérn covariance is an especially popular choice.

## Example (Lognormal input random field)

Let  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be a Lipschitz domain and consider the PDE problem

$$\begin{cases} -\nabla \cdot (a(\mathbf{x}, \omega) \nabla u(\mathbf{x}, \omega)) = f(\mathbf{x}) & \text{for } \mathbf{x} \in D, \\ u(\cdot, \omega)|_{\partial D} = 0, \end{cases}$$

where  $f: D \rightarrow \mathbb{R}$  is a fixed (deterministic) source term. We can represent a lognormally distributed random diffusion coefficient  $a: D \times \Omega \rightarrow \mathbb{R}$  using the KL expansion, e.g., as

$$a(\mathbf{x}, \omega) = a_0(\mathbf{x}) \exp \left( \sum_{k=1}^{\infty} y_k(\omega) \psi_k(\mathbf{x}) \right), \quad y_k \sim \mathcal{N}(0, 1),$$

where  $a_0 \in L^\infty(D)$  is such that  $a_0(\mathbf{x}) > 0$  and the random variables  $y_k$  are uncorrelated (and thus independent in the Gaussian case).

Due to the independence, we can consider the above as a *parametric PDE* with  $a(\mathbf{x}, \mathbf{y}) \equiv a(\mathbf{x}, \mathbf{y}(\omega))$  and  $u(\mathbf{x}, \mathbf{y}) \equiv u(\mathbf{x}, \mathbf{y}(\omega))$ , where (formally)  $\mathbf{y} \in \mathbb{R}^N$  is a *parametric vector* endowed with a product Gaussian measure.

## Example (Uniform and affine input random field)

Let  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be a Lipschitz domain,  $f: D \rightarrow \mathbb{R}$  is a fixed (deterministic) source term, and consider the PDE problem

$$\begin{cases} -\nabla \cdot (a(\mathbf{x}, \omega) \nabla u(\mathbf{x}, \omega)) = f(\mathbf{x}) & \text{for } \mathbf{x} \in D, \\ u(\cdot, \omega)|_{\partial D} = 0. \end{cases}$$

We can represent a uniformly distributed random diffusion coefficient  $a: D \times \Omega \rightarrow \mathbb{R}$  using the KL expansion, e.g., as

$$a(\mathbf{x}, \omega) = a_0(\mathbf{x}) + \sum_{k=1}^{\infty} y_k(\omega) \psi_k(\mathbf{x}), \quad y_k \sim \mathcal{U}(-\frac{1}{2}, \frac{1}{2}),$$

where the random variables  $y_k$  are uncorrelated. *Unlike the Gaussian setting, the random variables  $y_k$  are generally not independent!*

In numerical analysis, the random variables  $y_k$  are often **assumed** to be independent – this allows us to consider the above as a parametric PDE with  $a(\mathbf{x}, \mathbf{y}) \equiv a(\mathbf{x}, \mathbf{y}(\omega))$  and  $u(\mathbf{x}, \mathbf{y}) \equiv u(\mathbf{x}, \mathbf{y}(\omega))$ , where  $\mathbf{y} \in [-\frac{1}{2}, \frac{1}{2}]^{\mathbb{N}}$  is a *parametric vector* endowed with a uniform probability measure.

To estimate the statistical response, note that in the *lognormal model* the expected value of the PDE solution is (formally) given by

$$\mathbb{E}[u(\mathbf{x}, \cdot)] = \lim_{s \rightarrow \infty} \int_{\mathbb{R}^s} u(\mathbf{x}, \mathbf{y}) \prod_{j=1}^s \frac{e^{-\frac{1}{2}y_j^2}}{\sqrt{2\pi}} d\mathbf{y}$$

while in the *affine and uniform model* the expected value of the PDE solution is (formally) given by

$$\mathbb{E}[u(\mathbf{x}, \cdot)] = \lim_{s \rightarrow \infty} \int_{[-1/2, 1/2]^s} u(\mathbf{x}, \mathbf{y}) d\mathbf{y}.$$

- In practice, we need to truncate these infinite-dimensional integrals into finite-dimensional ones, incurring the so-called *dimension truncation error*. Since the PDE is solved numerically using the finite element method, this also incurs a *finite element discretization error*.
- To compute the resulting high-dimensional integrals for the dimensionally-truncated, finite element discretized PDE solution we use a *quasi-Monte Carlo (QMC) method*.

## 2. Quasi-Monte Carlo (QMC) methods

# Lattice rules

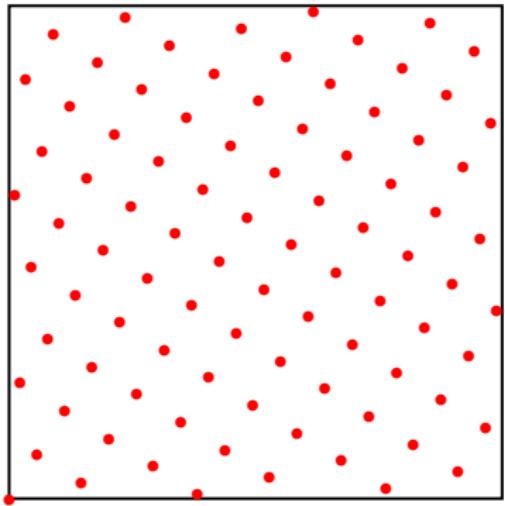
Rank-1 lattice rules

$$Q_{n,s}(f) = \frac{1}{n} \sum_{i=1}^n f(\mathbf{t}_i)$$

have the points

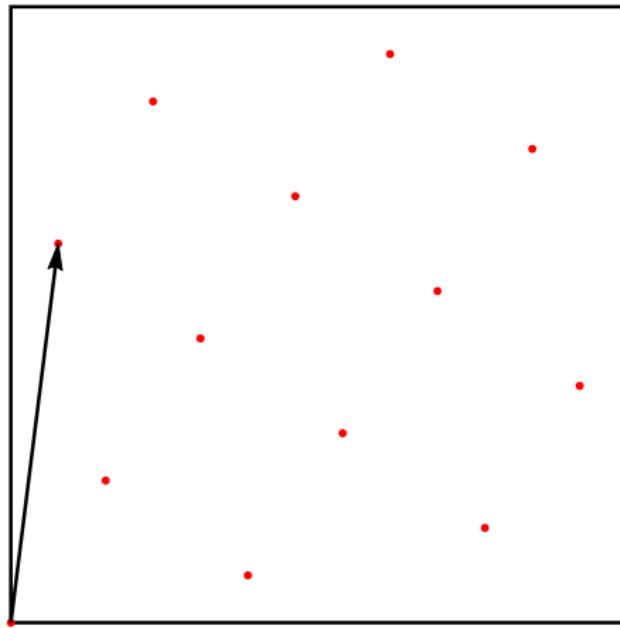
$$\mathbf{t}_i = \left\{ \frac{i\mathbf{z}}{n} \right\} = \text{mod}\left( \frac{i\mathbf{z}}{n}, 1 \right), \quad i \in \{1, \dots, n\},$$

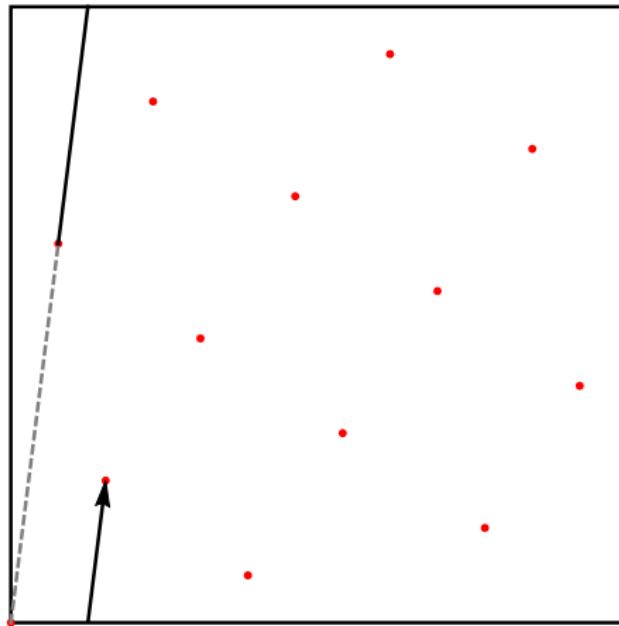
where the entire point set is determined by the *generating vector*  $\mathbf{z} \in \mathbb{N}^s$ , with all components *coprime* to  $n$ . The braces  $\{\cdot\}$  denote the componentwise fractional part of a vector.

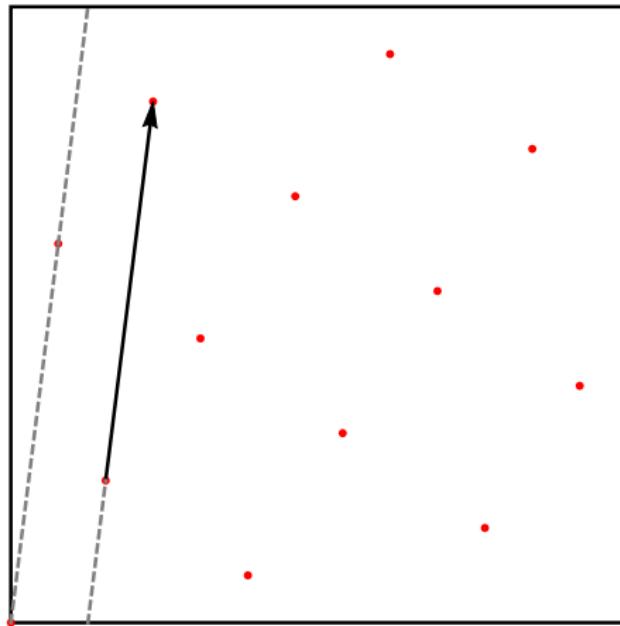


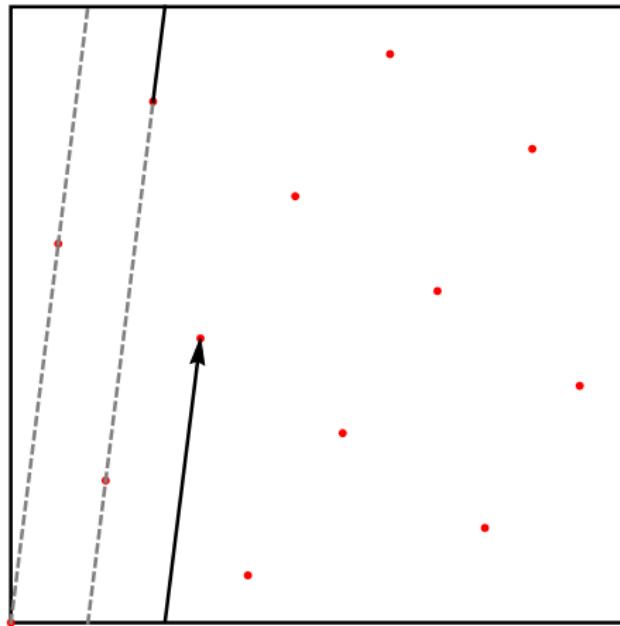
Lattice rule with  $\mathbf{z} = (1, 55)$  and  $n = 89$   
nodes in  $[0, 1]^2$

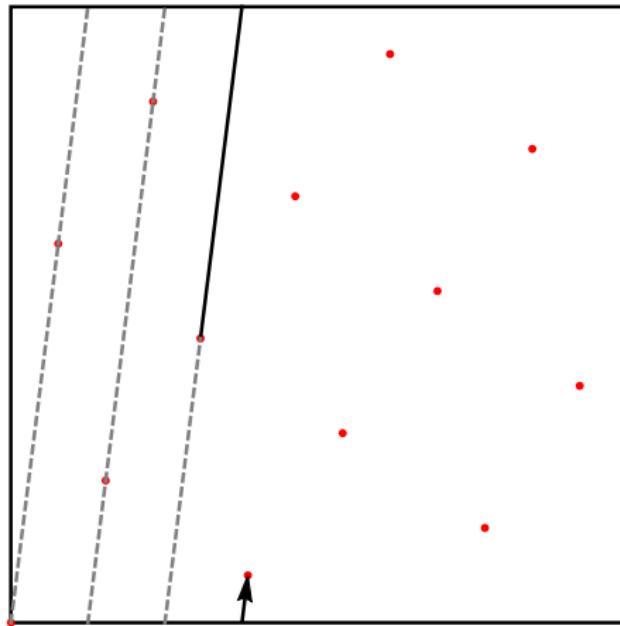
The quality of the lattice rule is determined by the choice of  $\mathbf{z}$ .

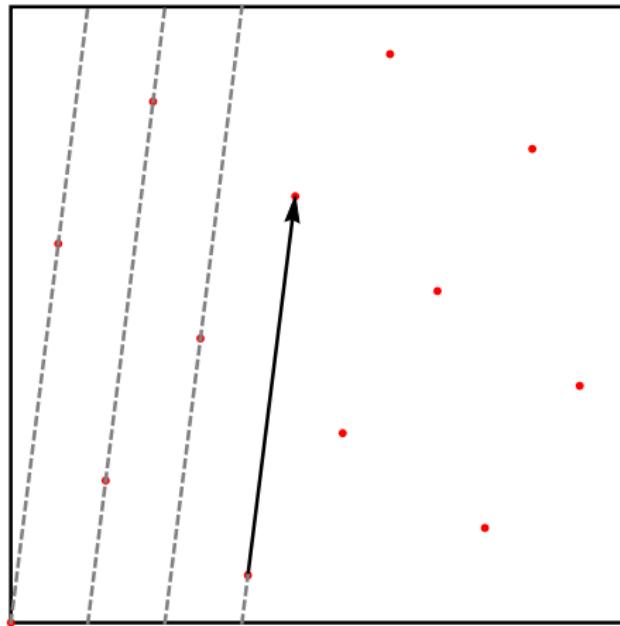


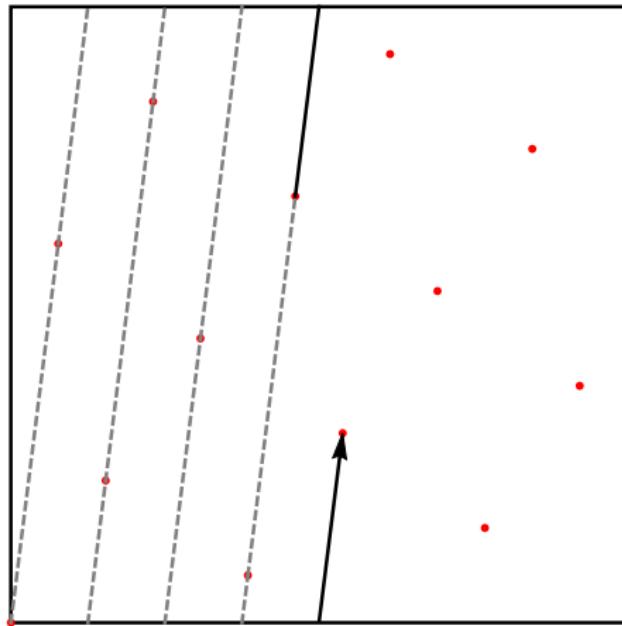


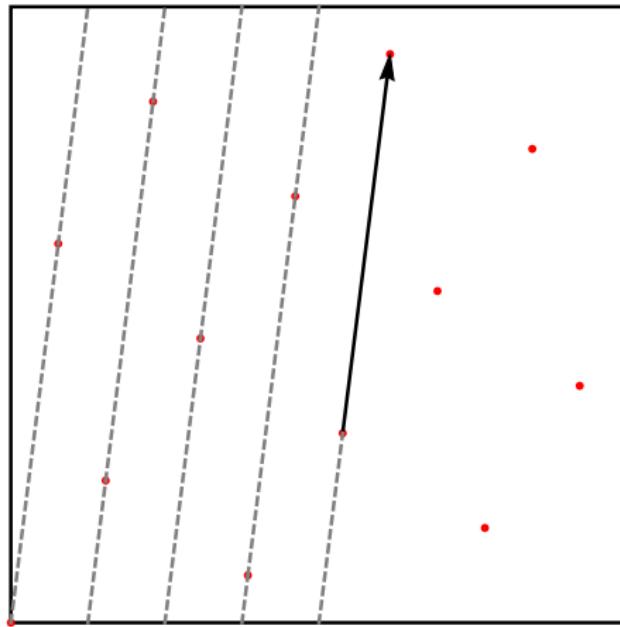


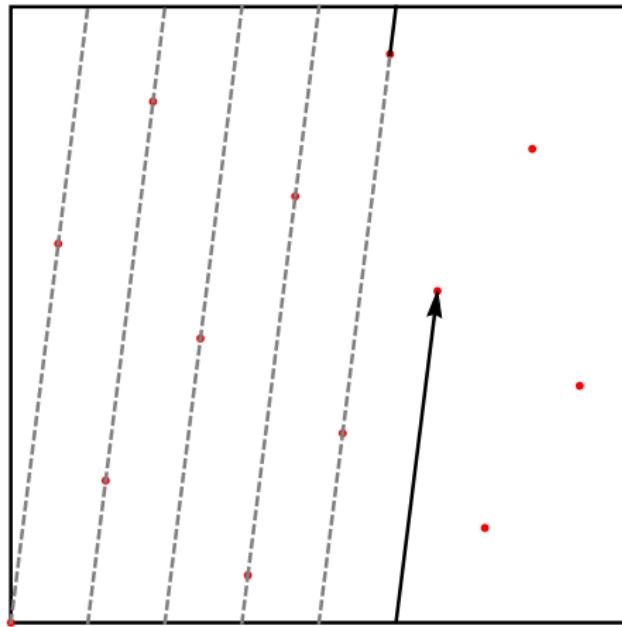


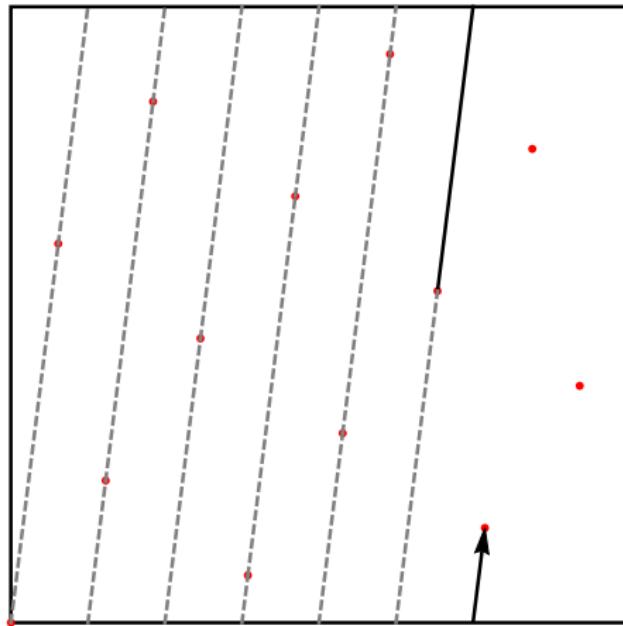


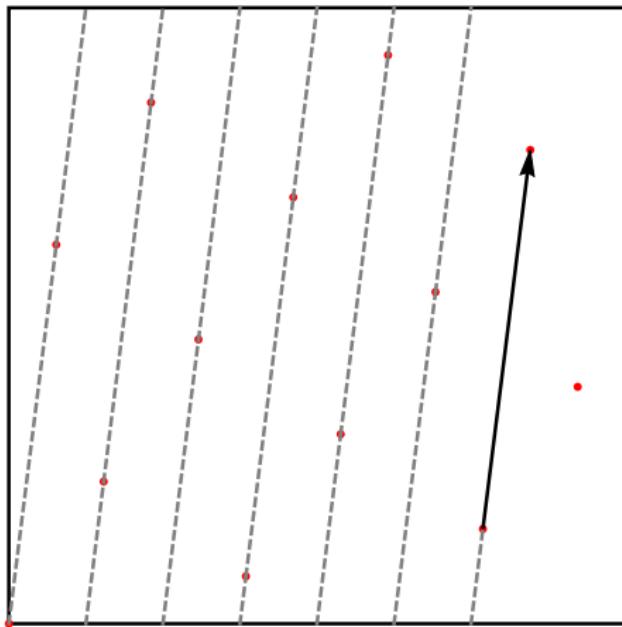


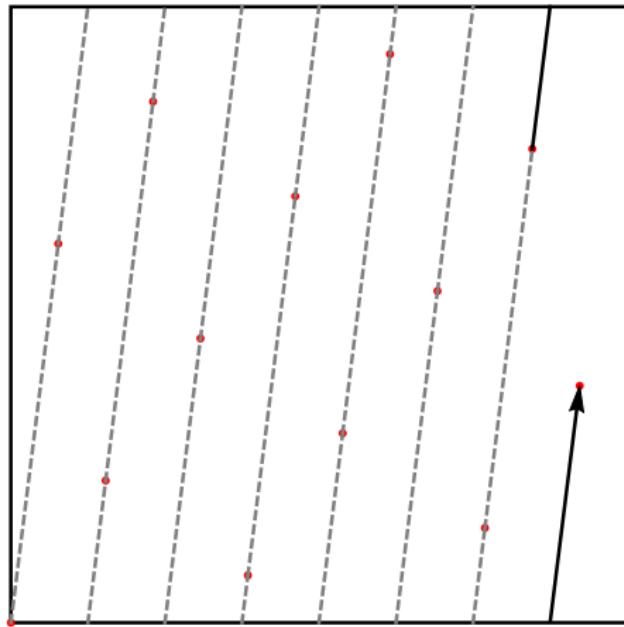












## Worst-case error

In the classical study of quadrature and cubature rules, we usually consider the so-called *worst-case error*. Suppose that  $f \in H$ , where  $H$  is a Hilbert space continuously embedded in  $C([0, 1]^s)$ . Let  $I_s : H \rightarrow \mathbb{R}$  be an integral operator

$$I_s f := \int_{[0,1]^s} f(\mathbf{x}) d\mathbf{x}$$

and let  $Q_{n,s} : H \rightarrow \mathbb{R}$  be a QMC rule

$$Q_{n,s} f := \frac{1}{n} \sum_{i=0}^{n-1} f(\mathbf{t}_i),$$

where  $P := \{\mathbf{t}_i \in [0, 1]^s \mid 0 \leq i \leq n - 1\}$  is a collection of cubature nodes. The worst-case error of cubature rule  $Q_{n,s}$  in  $H$  is defined by

$$e_{n,s}(P; H) := \sup_{\substack{f \in H \\ \|f\|_H \leq 1}} |I_s f - Q_{n,s} f|.$$

Note that this is precisely the operator norm of  $\|I_s - Q_{n,s}\|_{H \rightarrow \mathbb{R}}$ .

Since the worst-case error is just the operator norm of  $I_s - Q_{n,s}$ , we can express the cubature error as

$$|I_s f - Q_{n,s} f| \leq e_{n,s}(P; H) \|f\|_H.$$

Worst-case errors are in general hard to compute – except for the special case, when  $H$  is a *reproducing kernel Hilbert space* (RKHS).

Our strategy will be to *choose* the Hilbert space  $H$  (where our integrand  $f$  lives) to be such that it is possible to write down the expression for  $e_{n,s}(P; H)$  *explicitly* given a family of QMC rules. This allows us to analyze the dependence of the cubature error w.r.t.  $n$  and  $s$ .

We will end up taking  $H$  as an *unanchored, weighted Sobolev space* since this choice turns out to be “compatible” with the family of (randomly shifted) lattice rules!

## Reproducing kernel Hilbert space (RKHS)

Let  $H$  be a Hilbert space of functions on  $D \subseteq \mathbb{R}^s$ , with the property that *every point evaluation is a bounded linear functional*. That is, for any  $\mathbf{y} \in D$ , let

$$T_{\mathbf{y}}(f) := f(\mathbf{y}) \quad \text{for all } f \in H.$$

Then, since  $T_{\mathbf{y}}$  is a bounded linear functional, by Riesz representation theorem there exists a unique representer  $a_{\mathbf{y}} := K(\cdot, \mathbf{y}) \in H$  such that

$$T_{\mathbf{y}}(f) = \langle f, a_{\mathbf{y}} \rangle = \langle f, K(\cdot, \mathbf{y}) \rangle \quad \text{for all } f \in H.$$

The function  $K(\mathbf{x}, \mathbf{y})$  is known as the *reproducing kernel* of  $H$ .

### Definition (Reproducing kernel)

A *reproducing kernel* of a Hilbert space  $H$  of functions on  $D \subseteq \mathbb{R}^s$  is a function  $K: D \times D \rightarrow \mathbb{R}$  which satisfies

$$K(\cdot, \mathbf{y}) \in H \quad \text{for all } \mathbf{y} \in D$$

$$\text{and } f(\mathbf{y}) = \langle f, K(\cdot, \mathbf{y}) \rangle \quad \text{for all } f \in H \text{ and } \mathbf{y} \in D.$$

The latter property is known as the *reproducing property*.

## Remarks

- A *reproducing kernel Hilbert space* (RKHS) is a Hilbert space equipped with a reproducing kernel, or equivalently, it is a Hilbert space in which *every point evaluation is a bounded linear functional.*
- For any other bounded linear functional  $A: H \rightarrow \mathbb{R}$ , its representer  $a \in H$  satisfying  $A(f) = \langle f, a \rangle$  for all  $f \in H$  is given by

$$a(\mathbf{y}) = \langle a, K(\cdot, \mathbf{y}) \rangle = \langle K(\cdot, \mathbf{y}), a \rangle = A(K(\cdot, \mathbf{y})) \quad \text{for all } \mathbf{y} \in D.$$

- Any reproducing kernel  $K(\mathbf{x}, \mathbf{y})$  is symmetric in its arguments:

$$K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x}) \quad \text{for all } \mathbf{x}, \mathbf{y} \in D.$$

*Proof.* For fixed  $\mathbf{y} \in D$ , apply the reproducing property to the function  $f = K(\cdot, \mathbf{y})$  to get

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= f(\mathbf{x}) = \langle f, K(\cdot, \mathbf{x}) \rangle = \langle K(\cdot, \mathbf{y}), K(\cdot, \mathbf{x}) \rangle \\ &= \langle K(\cdot, \mathbf{x}), K(\cdot, \mathbf{y}) \rangle = K(\mathbf{y}, \mathbf{x}). \quad \square \end{aligned}$$

## Example

Suppose that we have a Hilbert space containing continuous functions on  $[0, 1]$  with square-integrable first order derivatives, equipped with the inner product

$$\langle f, g \rangle = \left( \int_0^1 f(x) dx \right) \left( \int_0^1 g(x) dx \right) + \int_0^1 f'(x)g'(x) dx.$$

Then this space has the *reproducing kernel*

$$K(x, y) = 1 + \eta(x, y), \quad \eta(x, y) = \frac{1}{2}B_2(|x - y|) + (x - \frac{1}{2})(y - \frac{1}{2}),$$

where  $B_2(x) := x^2 - x + \frac{1}{6}$  denotes the *Bernoulli polynomial of degree 2*.

That is, we claim that

$$\langle f, K(\cdot, y) \rangle = f(y) \quad \text{for all } y \in [0, 1].$$

### Example (continued)

By observing that

$$\int_0^1 K(x, y) dx = 1 \quad \text{and} \quad \frac{\partial}{\partial x} K(x, y) = x - \frac{1}{2} - \frac{1}{2} \text{sign}(x - y),$$

there holds

$$\begin{aligned}\langle f, K(\cdot, y) \rangle &= \left( \int_0^1 f(x) dx \right) \underbrace{\left( \int_0^1 K(x, y) dx \right)}_{=1} + \int_0^1 f'(x) \left( x - \frac{1}{2} - \frac{1}{2} \text{sign}(x - y) \right) dx \\ &= \int_0^1 f(x) dx + \int_0^1 f'(x)x dx - \frac{1}{2} \int_0^1 f'(x) dx + \frac{1}{2} \int_0^y f'(x) dx - \frac{1}{2} \int_y^1 f'(x) dx \\ &= \cancel{\int_0^1 f(x) dx} + \cancel{f(1)} - \cancel{\int_0^1 f(x) dx} - \frac{1}{2} \cancel{f(1)} + \frac{1}{2} \cancel{f(0)} + \frac{1}{2} f(y) - \frac{1}{2} \cancel{f(0)} - \frac{1}{2} \cancel{f(1)} + \frac{1}{2} f(y) \\ &= f(y)\end{aligned}$$

for all  $y \in [0, 1]$ , as desired.

## Theorem

Let  $H := H_s(K)$  be an RKHS and let  $K: [0, 1]^s \times [0, 1]^s \rightarrow \mathbb{R}$  be a reproducing kernel that satisfies

$$\int_{[0,1]^s} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} < \infty.$$

Then

$$\begin{aligned} e_{n,s}^2(P; H_s(K)) &= \int_{[0,1]^s} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} - \frac{2}{n} \sum_{i=0}^{n-1} \int_{[0,1]^s} K(\mathbf{t}_i, \mathbf{y}) d\mathbf{y} \\ &\quad + \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} K(\mathbf{t}_i, \mathbf{t}_j). \end{aligned} \tag{2}$$

*Proof.* For  $f \in H$ , we apply the reproducing property  $f(\mathbf{t}_k) = \langle f, K(\cdot, \mathbf{t}_k) \rangle_H$  and average the results to obtain

$$Q_{n,s}f = \frac{1}{n} \sum_{k=0}^{n-1} f(\mathbf{t}_k) = \frac{1}{n} \sum_{k=0}^{n-1} \langle f, K(\cdot, \mathbf{t}_k) \rangle_H = \left\langle f, \frac{1}{n} \sum_{k=0}^{n-1} K(\cdot, \mathbf{t}_k) \right\rangle_H. \quad (3)$$

Similarly, we find that

$$I_s f = \int_{[0,1]^s} f(\mathbf{x}) d\mathbf{x} = \int_{[0,1]^s} \langle f, K(\cdot, \mathbf{x}) \rangle_H d\mathbf{x} = \left\langle f, \int_{[0,1]^s} K(\cdot, \mathbf{x}) d\mathbf{x} \right\rangle_H, \quad (4)$$

which holds provided that  $\int_{[0,1]^s} K(\cdot, \mathbf{x}) d\mathbf{x} \in H$ . However, this is guaranteed by our assumption since

$$\begin{aligned} \left\| \int_{[0,1]^s} K(\cdot, \mathbf{x}) d\mathbf{x} \right\|_H^2 &= \int_{[0,1]^s} \int_{[0,1]^s} \langle K(\cdot, \mathbf{x}), K(\cdot, \mathbf{y}) \rangle_H d\mathbf{x} d\mathbf{y} \\ &= \int_{[0,1]^s} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} < \infty, \end{aligned}$$

which will hold for all the kernels we shall consider.

Taking the difference of (3) and (4) yields

$$I_s f - Q_{n,s} f = \left\langle f, \int_{[0,1]^s} K(\cdot, \mathbf{x}) d\mathbf{x} - \frac{1}{n} \sum_{i=0}^{n-1} K(\cdot, \mathbf{t}_i) \right\rangle_H = \langle f, \xi \rangle_H,$$

where

$$\xi(\mathbf{y}) := \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} - \frac{1}{n} \sum_{i=0}^{n-1} K(\mathbf{y}, \mathbf{t}_i), \quad \mathbf{y} \in [0, 1]^s,$$

is called the *representer* of the integration error since

$$e_{n,s}(P; H) = \sup_{\|f\| \leq 1} |\langle f, \xi \rangle_H| = \|\xi\|_H.$$

Especially, the supremum is attained by  $f = \xi / \|\xi\|_H \in H$  and we obtain

$$\begin{aligned} e_{n,s}^2(P; H) &= \left\| \int_{[0,1]^s} K(\cdot, \mathbf{x}) d\mathbf{x} - \frac{1}{n} \sum_{i=0}^{n-1} K(\mathbf{x}, \mathbf{t}_i) \right\|^2 \\ &= \int_{[0,1]^s} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} - \frac{2}{n} \sum_{i=0}^{n-1} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{t}_i) d\mathbf{x} + \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} K(\mathbf{t}_i, \mathbf{t}_j), \end{aligned}$$

as desired. □



## Randomly shifted rank-1 lattice points

In what follows, we will discuss randomly shifted QMC rules.

Consider the rank-1 lattice point set  $\mathbf{t}_k := \left\{ \frac{kz}{n} \right\}$  for some generating vector  $z \in \mathbb{N}^s$  and fixed  $n \in \mathbb{N}$ . Given a vector  $\Delta \in [0, 1]^s$ , known as the *shift*, the  $\Delta$ -shift of the QMC points  $\mathbf{t}_0, \dots, \mathbf{t}_{n-1}$  is defined as the point set

$$\{\mathbf{t}_k + \Delta\}, \quad k = 0, \dots, n - 1.$$

Shifting preserves the lattice structure. In practice, we will generate a number of independent random shifts  $\Delta_0, \dots, \Delta_{R-1}$  from  $\mathcal{U}([0, 1]^s)$  and take the average of  $\Delta_0, \dots, \Delta_{R-1}$ -shifted QMC rules as our approximation of  $I_s$ .

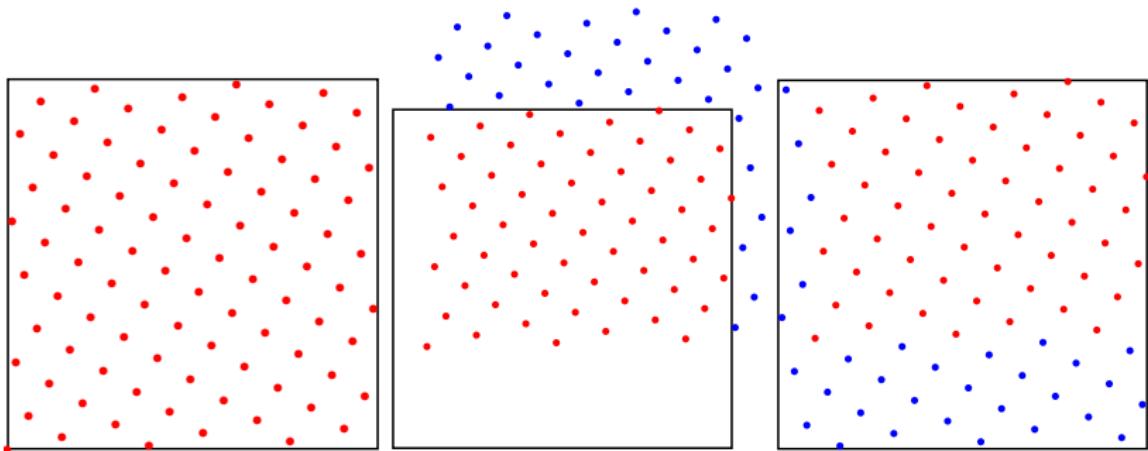
Advantages:

- Leads to a shift-invariant kernel (advantageous for high-dimensional computation).
- Randomization yields an unbiased estimator of the integral.
- Randomization provides a practical error estimate.

Shifted rank-1 lattice rules have points

$$\left\{ \frac{k\mathbf{z}}{n} + \boldsymbol{\Delta} \right\}, \quad k = 0, \dots, n-1.$$

*Use a number of random shifts for error estimation.*



Lattice rule shifted by  $\boldsymbol{\Delta} = (0.1, 0.3)$ .

## Randomization in practice

- Generate  $R$  independent random shifts  $\Delta_0, \dots, \Delta_{R-1}$  from  $\mathcal{U}([0, 1]^s)$ .
- For a given QMC rule with points  $(\mathbf{t}_i)_{i=0}^{n-1} \subset [0, 1]^s$ , form the approximations  $Q_{n,s}^{(0)} f, \dots, Q_{n,s}^{(R-1)} f$ , where

$$Q_{n,s}^{\Delta_r} f = \frac{1}{n} \sum_{i=0}^{n-1} f(\{\mathbf{t}_i + \Delta_r\}), \quad r = 0, \dots, R-1,$$

is the approximation of the integral using a  $\Delta_r$ -shift of the original QMC rule.

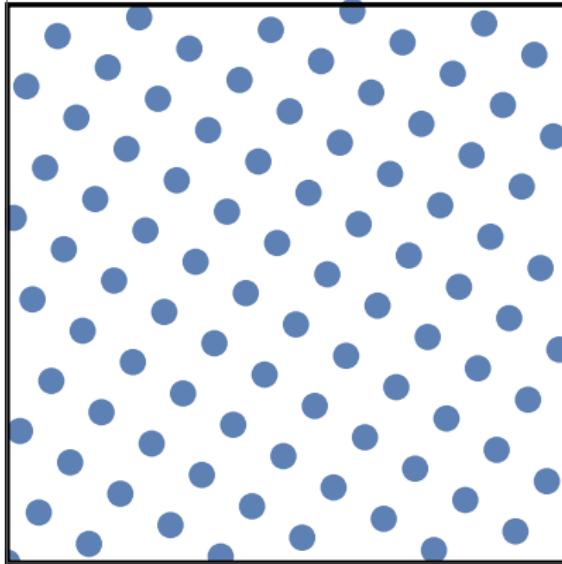
- We take the *average*

$$\overline{Q}_{n,s,R} f = \frac{1}{R} \sum_{r=0}^{R-1} Q_{n,s}^{\Delta_r} f$$

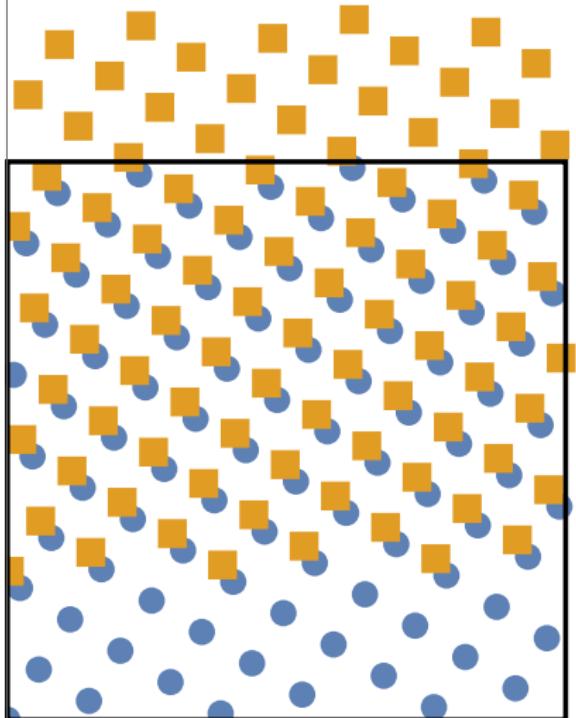
as our *final* approximation of the integral.

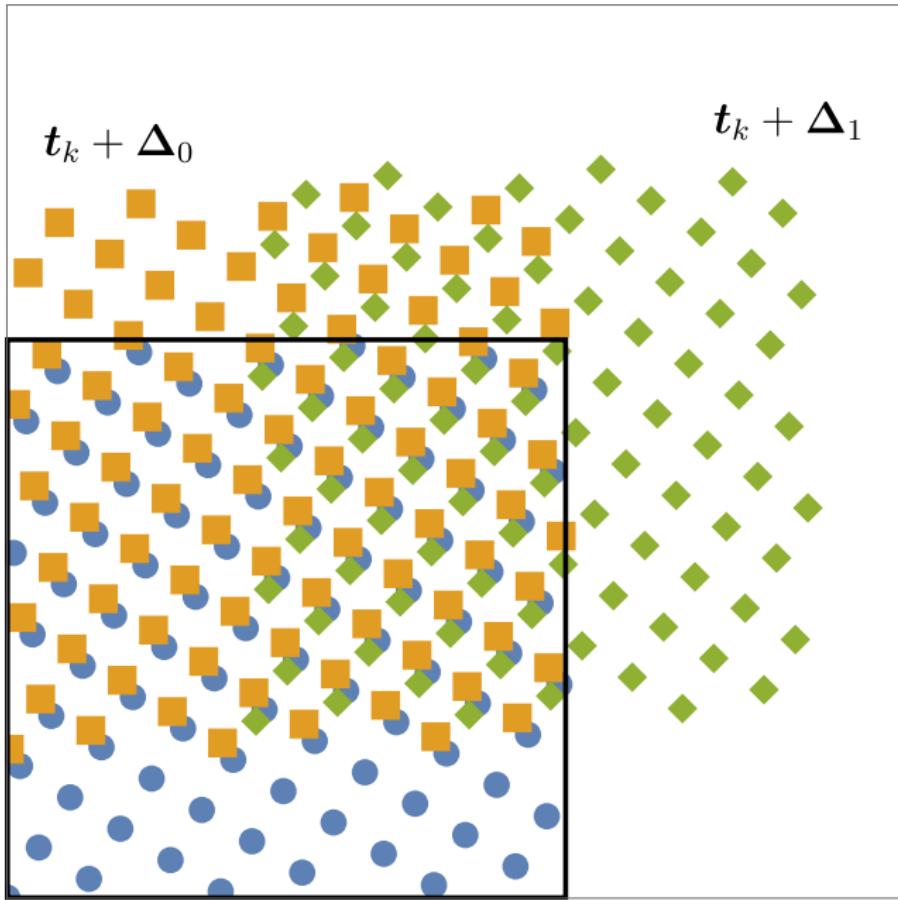
- An *unbiased* estimate for the mean-square error of  $\overline{Q}_{n,s,R} f$  is given by

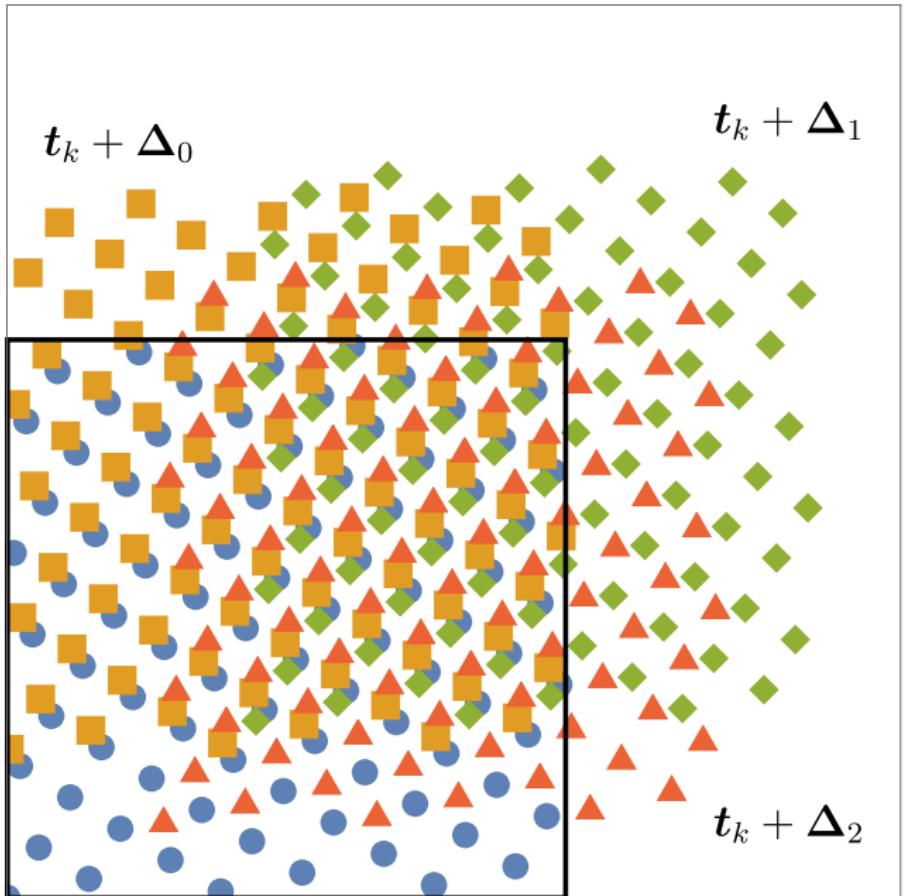
$$\mathbb{E}_{\Delta} |I_s f - Q_{n,s}^{\Delta} f|^2 \approx \frac{1}{R(R-1)} \sum_{r=0}^{R-1} (Q_{n,s}^{\Delta_r} f - \overline{Q}_{n,s,R} f)^2.$$



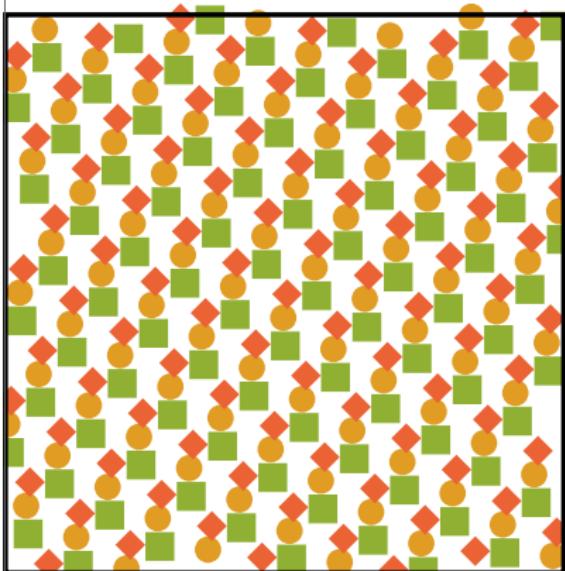
$$t_k + \Delta_0$$







$$\{\{t_k + \Delta_0\}, \{t_k + \Delta_1\}, \{t_k + \Delta_2\}\}$$



$$Q_{n,s}^{\Delta_0} f = \frac{1}{n} \sum_{i=0}^{n-1} f(\{t_i + \Delta_0\}), \quad Q_{n,s}^{\Delta_1} f = \frac{1}{n} \sum_{i=0}^{n-1} f(\{t_i + \Delta_1\}), \quad Q_{n,s}^{\Delta_2} f = \frac{1}{n} \sum_{i=0}^{n-1} f(\{t_i + \Delta_2\})$$

$$\text{QMC approximation with 3 random shifts: } \overline{Q}_{n,s,3} f = \frac{Q_{n,s}^{\Delta_0} f + Q_{n,s}^{\Delta_1} f + Q_{n,s}^{\Delta_2} f}{3}.$$

## Shift-averaged worst-case error

For any QMC point set  $P = \{\mathbf{t}_0, \dots, \mathbf{t}_{n-1}\}$  and any shift  $\Delta \in [0, 1]^s$ , let

$$P + \Delta := \{\{\mathbf{t}_i + \Delta\} \mid i = 0, 1, \dots, n-1\}$$

denote the *shifted QMC point set*, and let  $Q_{n,s}^\Delta f$  denote a corresponding shifted QMC rule (over the point set  $P + \Delta$ ). For any integrand  $f \in H$ , it follows from the definition of the worst-case error that

$$|I_s f - Q_{n,s}(\Delta; f)| \leq e_{n,s}(P + \Delta; H) \|f\|_H,$$

where  $e_{n,s}(P + \Delta; H) := \sup_{\|f\|_H \leq 1} |I_s(f) - Q_{n,s}^\Delta f|$ . We deduce a bound for the *root-mean-square* error

$$\sqrt{\mathbb{E}_\Delta |I_s f - Q_{n,s}^\Delta f|^2} \leq e_{n,s}^{\text{sh}}(P; H) \|f\|_H,$$

where the expected value  $\mathbb{E}_\Delta$  is taken over the random shift  $\Delta$  which is uniformly distributed over  $[0, 1]^s$  and the quantity

$$e_{n,s}^{\text{sh}}(P; H) := \sqrt{\int_{[0,1]^s} e_{n,s}^2(P + \Delta; H) d\Delta}$$

is called the *shift-averaged worst-case error*.

## Theorem (Formula for the shift-averaged worst-case error)

$$[e_{n,s}^{\text{sh}}(P; H_s(K))]^2 = - \int_{[0,1]^s} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} + \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} K^{\text{sh}}(\mathbf{t}_i, \mathbf{t}_j),$$

where

$$K^{\text{sh}}(\mathbf{x}, \mathbf{y}) := \int_{[0,1]^s} K(\{\mathbf{x} + \boldsymbol{\Delta}\}, \{\mathbf{y} + \boldsymbol{\Delta}\}) d\boldsymbol{\Delta}, \quad \mathbf{x}, \mathbf{y} \in [0, 1]^s.$$

*Proof.* The definition of shift-averaged WCE and (2) imply

$$\begin{aligned} [e_{n,s}^{\text{sh}}(P; H_s(K))]^2 &= \int_{[0,1]^s} e_{n,s}^2(P + \boldsymbol{\Delta}; H) d\boldsymbol{\Delta} \\ &= \int_{[0,1]^s} \int_{[0,1]^s} K(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} - \frac{2}{n} \sum_{i=0}^{n-1} \int_{[0,1]^s} \int_{[0,1]^s} K(\{\mathbf{t}_i + \boldsymbol{\Delta}\}, \mathbf{y}) d\boldsymbol{\Delta} d\mathbf{y} \\ &\quad + \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \int_{[0,1]^s} K(\{\mathbf{t}_i + \boldsymbol{\Delta}\}, \{\mathbf{t}_j + \boldsymbol{\Delta}\}) d\boldsymbol{\Delta}. \end{aligned}$$

The result follows by a change of variables  $\mathbf{x} = \{\mathbf{t}_i + \boldsymbol{\Delta}\}$  in the second term. □



## Remarks

$$K^{\text{sh}}(\mathbf{x}, \mathbf{y}) := \int_{[0,1]^s} K(\{\mathbf{x} + \boldsymbol{\Delta}\}, \{\mathbf{y} + \boldsymbol{\Delta}\}) d\boldsymbol{\Delta}, \quad \mathbf{x}, \mathbf{y} \in [0, 1]^s.$$

- The function  $K^{\text{sh}}$  is actually a reproducing kernel, with the *shift-invariant property*

$$K^{\text{sh}}(\mathbf{x}, \mathbf{y}) = K^{\text{sh}}(\{\mathbf{x} + \boldsymbol{\Delta}\}, \{\mathbf{y} + \boldsymbol{\Delta}\}) \quad \text{for all } \mathbf{x}, \mathbf{y}, \boldsymbol{\Delta} \in [0, 1].$$

Equivalently,

$$K^{\text{sh}}(\mathbf{x}, \mathbf{y}) = K^{\text{sh}}(\{\mathbf{x} - \mathbf{y}\}, \mathbf{0}) \quad \text{for all } \mathbf{x}, \mathbf{y} \in [0, 1].$$

- The function  $K^{\text{sh}}$  is called the *shift-invariant kernel associated with  $K$* .

## Weighted Sobolev spaces

## Unanchored, weighted Sobolev space

For our purposes, the relevant function space setting will be the *unanchored, weighted Sobolev space*. For any given collection  $(\gamma_u)_{u \subseteq \{1:s\}}$  of positive numbers (called *weights*), we associate a space  $H_{s,\gamma}$  containing continuous functions on  $[0, 1]^s$  whose *mixed first partial derivatives are square-integrable*. It is defined by the reproducing kernel

$$K_{s,\gamma}(\mathbf{x}, \mathbf{y}) = \sum_{u \subseteq \{1:s\}} \gamma_u \prod_{j \in u} \eta(x_j, y_j), \quad \eta(x, y) := \frac{1}{2}B_2(|x-y|) + (x - \frac{1}{2})(y - \frac{1}{2}),$$

where  $B_2(x) := x^2 - x + \frac{1}{6}$  is the Bernoulli polynomial of degree 2 and we use the notation  $\{1 : s\} := \{1, \dots, s\}$ .

The norm  $\|f\|_{s,\gamma} = \sqrt{\langle f, f \rangle_{s,\gamma}}$  is induced by the inner product

$$\begin{aligned} \langle f, g \rangle_{s,\gamma} &= \sum_{u \subseteq \{1:s\}} \frac{1}{\gamma_u} \int_{[0,1]^{|u|}} \left( \int_{[0,1]^{s-|u|}} \frac{\partial^{|u|}}{\partial \mathbf{x}_u} f(\mathbf{x}) d\mathbf{x}_{-u} \right) \\ &\quad \times \left( \int_{[0,1]^{s-|u|}} \frac{\partial^{|u|}}{\partial \mathbf{x}_u} g(\mathbf{x}) d\mathbf{x}_{-u} \right) d\mathbf{x}_u, \end{aligned}$$

where  $d\mathbf{x}_u := \prod_{j \in u} dx_j$  and  $d\mathbf{x}_{-u} := \prod_{j \in \{1:s\} \setminus u} dx_j$ .

## Remarks

- We sum over all  $2^s$  possible subsets of the indices  $\{1 : s\}$ . By convention, an empty product is 1.
- Each term of the sum corresponds to a subset of variables  $\mathbf{x}_{\mathfrak{u}} = \{x_j \mid j \in \mathfrak{u}\}$ . We refer to these as the “active” variables, and denote the remaining “inactive” variables by  $\mathbf{x}_{-\mathfrak{u}}$ .
- The cardinality  $|\mathfrak{u}|$  of the set  $\mathfrak{u}$  is referred to as the “order” of the subset of variables  $\mathbf{x}_{\mathfrak{u}}$ . There is a *weight* parameter  $\gamma_{\mathfrak{u}}$  associated with every subset of variables  $\mathbf{x}_{\mathfrak{u}}$ . The weights together model the relative importance between different subsets of variables. A small weight  $\gamma_{\mathfrak{u}}$  means that the  $L^2$  norm of  $\frac{\partial^{|\mathfrak{u}|} f}{\partial \mathbf{x}_{\mathfrak{u}}}$  must also be small.
- Note that  $\|\cdot\|_{s,\gamma}$  and  $\|\cdot\|_{s,c\gamma}$  are equivalent norms for any  $c > 0$ .<sup>†</sup> Therefore we do not lose any generality by assuming that the weights have been normalized s.t.  $\gamma_{\emptyset} = 1$ . WLOG, we will always use the convention that  $\gamma_{\emptyset} := 1$ .

---

<sup>†</sup>Here,  $c\gamma = (c\gamma_{\mathfrak{u}})_{\mathfrak{u} \subseteq \{1:s\}}$ .

## Special forms of weights

- *Product weights*: we have a sequence of numbers satisfying  $\gamma_1 \geq \gamma_2 \geq \dots$  and we take

$$\gamma_{\mathbf{u}} = \prod_{j \in \mathbf{u}} \gamma_j.$$

In this case, the reproducing kernel is given by the product

$$K_{s,\gamma}(\mathbf{x}, \mathbf{y}) = \prod_{j=1}^s \left( 1 + \gamma_j \eta(x_j, y_j) \right).$$

- *Finite order weights*: there exists  $q \in \mathbb{N}$  s.t.  $\gamma_{\mathbf{u}} = 0$  for all  $|\mathbf{u}| > q$ .
- *Order dependent weights*: we have a sequence of numbers  $\Gamma_1, \Gamma_2, \dots$ , and take

$$\gamma_{\mathbf{u}} = \Gamma_{|\mathbf{u}|}.$$

- *Product-and-order dependent (POD) weights*: we have two sequences  $\gamma_1, \gamma_2, \dots$  and  $\Gamma_1, \Gamma_2, \dots$ , and take

$$\gamma_{\mathbf{u}} = \Gamma_{|\mathbf{u}|} \prod_{j \in \mathbf{u}} \gamma_j.$$

## Why weighted spaces are interesting

Theorem (Sloan and Woźniakowski 1998)

Consider  $H_{s,\gamma}$  equipped with product weights  $\gamma_u = \prod_{j \in u} \gamma_j$ . Then there exist point sets  $P_n \subset [0, 1]^s$  for  $n = 1, 2, \dots$  such that the worst-case error  $e_{n,s}(P_n; H_{s,\gamma})$  is bounded independently of  $s$  if and only if

$$\sum_{j=1}^{\infty} \gamma_j < \infty. \tag{5}$$

---

To be more precise, the result has two parts:

- If condition (5) does *not* hold, then no matter how the points are chosen, the worst-case error is unbounded as  $s \rightarrow \infty$ .
- However, if (5) holds, then “good points” exist (although the result does not say how to find them).

Recall that  $H_{s,\gamma}$  is defined via the reproducing kernel

$$K_{s,\gamma}(\mathbf{x}, \mathbf{y}) = \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{j \in \mathfrak{u}} \eta(x_j, y_j), \quad \eta(x, y) := \frac{1}{2}B_2(|x-y|) + (x - \frac{1}{2})(y - \frac{1}{2}),$$

where  $B_2(x) := x^2 - x + \frac{1}{6}$  is the Bernoulli polynomial of degree 2.

### Lemma

$$\int_{[0,1]^s} K_{s,\gamma}(\mathbf{x}, \mathbf{y}) d\mathbf{y} = 1,$$

$$\int_{[0,1]^s} \int_{[0,1]^s} K_{s,\gamma}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = 1,$$

$$\int_{[0,1]^s} K_{s,\gamma}(\mathbf{x}, \mathbf{x}) d\mathbf{x} = \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} (\frac{1}{6})^{|\mathfrak{u}|}.$$

Recall that  $H_{s,\gamma}$  is defined via the reproducing kernel

$$K_{s,\gamma}(\mathbf{x}, \mathbf{y}) = \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{j \in \mathfrak{u}} \eta(x_j, y_j), \quad \eta(x, y) := \frac{1}{2}B_2(|x-y|) + (x - \frac{1}{2})(y - \frac{1}{2}),$$

where  $B_2(x) := x^2 - x + \frac{1}{6}$  is the Bernoulli polynomial of degree 2.

For our analysis, we will need the shift-invariant kernel associated with  $K_{s,\gamma}$ .

### Lemma

$$\begin{aligned} K_{s,\gamma}^{\text{sh}}(\mathbf{x}, \mathbf{y}) &:= \int_{[0,1]^s} K_{s,\gamma}(\{\mathbf{x} + \boldsymbol{\Delta}\}, \{\mathbf{y} + \boldsymbol{\Delta}\}) d\boldsymbol{\Delta} \\ &= \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{j \in \mathfrak{u}} B_2(|x_j - y_j|). \end{aligned}$$

*Proof.* This is an immediate consequence of

$$\int_0^1 \eta(\{x + \Delta\}, \{y + \Delta\}) d\Delta = B_2(|x - y|). \quad \square$$

Let

$$P = \left\{ \left\{ \frac{k\mathbf{z}}{n} \right\} \mid k = 0, \dots, n-1 \right\}$$

be a rank-1 lattice point set corresponding to generating vector  $\mathbf{z} \in \mathbb{N}^s$  and  $n \in \mathbb{N}$ .

When dealing with the shift-invariant kernel corresponding to the unanchored, weighted Sobolev space  $H_{s,\gamma}$ , we use the shorthand notation

$$e_{n,s}^{\text{sh}}(\mathbf{z}) := e_{n,s}^{\text{sh}}(P; H_{s,\gamma}).$$

## Lemma

The shift-averaged worst-case error for a rank-1 lattice rule in the weighted unanchored Sobolev space satisfies

$$[e_{n,s}^{\text{sh}}(\mathbf{z})]^2 = \frac{1}{n} \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \sum_{k=0}^{n-1} \prod_{j \in \mathfrak{u}} B_2 \left( \left\{ \frac{k z_j}{n} \right\} \right).$$

*Proof.* Let  $\mathbf{t}_j = \left\{ \frac{j \mathbf{z}}{n} \right\}$ . We have the kernel

$$K_{s,\gamma}(\mathbf{x}, \mathbf{y}) = \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{j \in \mathfrak{u}} \eta(x_j, y_j), \quad \eta(x, y) := \frac{1}{2} B_2(|x-y|) + (x - \frac{1}{2})(y - \frac{1}{2}),$$

which satisfies  $\int_{[0,1]^s} \int_{[0,1]^s} K_{s,\gamma}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = 1$ . We showed that the shift-invariant kernel related to  $K$  is given by

$$K_{s,\gamma}^{\text{sh}}(\mathbf{x}, \mathbf{y}) = \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{k \in \mathfrak{u}} B_2(|x_k - y_k|).$$

Moreover, we showed that the shift-averaged WCE is given by

$$[e_{n,s}^{\text{sh}}(\mathbf{z})]^2 = - \int_{[0,1]^s} \int_{[0,1]^s} K_{s,\gamma}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} + \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} K_{s,\gamma}^{\text{sh}}(\mathbf{t}_i, \mathbf{t}_j).$$

Making the obvious substitutions, we arrive at

$$\begin{aligned}
 [e_{n,s}^{\text{sh}}(\mathbf{z})]^2 &= -1 + \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \sum_{\mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{k \in \mathfrak{u}} B_2 \left( \left\{ \frac{(i-j)z_k}{n} \right\} \right) \quad (\gamma_{\emptyset} := 1) \\
 &= \frac{1}{n^2} \sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{k \in \mathfrak{u}} B_2 \left( \left\{ \frac{\text{mod}(i-j, n)z_k}{n} \right\} \right) \\
 &= \frac{1}{n^2} \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \underbrace{\sum_{i=0}^{n-1} \sum_{j=0}^{n-1} \prod_{k \in \mathfrak{u}} B_2 \left( \left\{ \frac{\text{mod}(i-j, n)z_k}{n} \right\} \right)}_{= n \sum_{\ell=0}^{n-1} \prod_{k \in \mathfrak{u}} B_2 \left( \left\{ \frac{\ell z_k}{n} \right\} \right)}.
 \end{aligned}$$

Final step: as  $i$  and  $j$  range from 0 to  $n-1$ , the values of  $\text{mod}(i-j, n)$  are just  $0, \dots, n-1$  in a different order (see next slide for illustration), with each value occurring  $n$  times. Thus

$$[e_{n,s}^{\text{sh}}(\mathbf{z})]^2 = \frac{1}{n} \sum_{\ell=0}^{n-1} \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}} \prod_{k \in \mathfrak{u}} B_2 \left( \left\{ \frac{\ell z_k}{n} \right\} \right),$$

as desired. □



## An illustration of the counting argument used on the previous slide

$i/j$	0	1	2	3	4	$\dots$	$n - 1$
0	$f(0)$	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$\dots$	$f(n - 1)$
1	$f(n - 1)$	$f(0)$	$f(1)$	$f(2)$	$f(3)$	$\dots$	$f(n - 2)$
2	$f(n - 2)$	$f(n - 1)$	$f(0)$	$f(1)$	$f(2)$	$\dots$	$f(n - 3)$
3	$f(n - 3)$	$f(n - 2)$	$f(n - 1)$	$f(0)$	$f(1)$	$\dots$	$f(n - 4)$
4	$f(n - 4)$	$f(n - 3)$	$f(n - 2)$	$f(n - 1)$	$f(0)$	$\dots$	$f(n - 5)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
$n - 1$	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$\dots$	$f(0)$

Table of the values  $f(\text{mod}(i - j, n))$ , when  $i, j \in \{0, 1, \dots, n - 1\}$ .

By a simple counting argument we can write

$$\sum_{i=0}^{n-1} \sum_{j=0}^{n-1} f(\text{mod}(i - j, n)) = n \sum_{\ell=0}^{n-1} f(\ell)$$

for any function  $f: \{0, 1, \dots, n - 1\} \rightarrow \mathbb{R}$ .

## Component-by-component construction

The components of the generating vector  $\mathbf{z}$  can be restricted to the set

$$\mathbb{U}_n := \{z \in \mathbb{Z} \mid 1 \leq z \leq n-1 \text{ and } \gcd(z, n) = 1\},$$

whose cardinality is given by the Euler totient function  $\varphi(n) := |\mathbb{U}_n|$ . When  $n$  is prime,  $\varphi(n)$  takes its largest value  $n - 1$ .

We know that for  $f \in H_{s,\gamma}$ , there holds

$$\sqrt{\mathbb{E}_{\Delta} |I_s f - Q_{n,s}^{\Delta} f|^2} \leq e_{n,s}^{\text{sh}}(\mathbf{z}) \|f\|_{s,\gamma}.$$

Finding  $\mathbf{z}^* = \arg \min_{\mathbf{z} \in \mathbb{U}_n} e_{n,s}^{\text{sh}}(\mathbf{z})$  is not computationally feasible: the search space contains altogether up to  $(n-1)^s$  possible choices for  $\mathbf{z}$ . However, the *component-by-component (CBC) construction* provides a feasible way to obtain good lattice generating vectors.

# CBC construction

**CBC construction.** Given  $n, s$ , and weights  $(\gamma_u)_{u \subseteq \{1:s\}}$ .

1. Set  $z_1 = 1$ .
2. For  $k = 2, 3, \dots, s$ , choose  $z_k \in \mathbb{U}_n$  to minimize  $[e_{n,k}^{\text{sh}}(z_1, \dots, z_k)]^2$ .

Remarks:

- Note that we have the (in principle computable) expression

$$[e_{n,k}^{\text{sh}}(\mathbf{z})]^2 = \frac{1}{n} \sum_{\emptyset \neq u \subseteq \{1:k\}} \gamma_u \sum_{\ell=0}^{n-1} \prod_{j \in u} B_2 \left( \left\{ \frac{\ell z_j}{n} \right\} \right). \quad (6)$$

- We will show that when the weights  $(\gamma_u)_{u \subseteq \{1:s\}}$  are so-called *product-and-order dependent (POD)* weights, i.e., they can be written in the form

$$\gamma_u := \Gamma_{|u|} \prod_{j \in u} \gamma_j, \quad u \subseteq \{1 : s\},$$

where  $\gamma_\emptyset := 1$ ,  $(\Gamma_k)_{k=1}^\infty$  and  $(\gamma_j)_{j=1}^\infty$  are sequences of positive numbers, then the value of (6) can be obtained in  $\mathcal{O}(s n \log n + s^2 n)$  time using the so-called *fast CBC algorithm*. This is quadratic, not exponential, w.r.t. the dimension  $s$ .

- The CBC algorithm is a greedy algorithm: in general, it will **not** produce a generating vector which minimizes  $e_{n,s}^{\text{sh}}(\mathbf{z})$ . Regardless, we **can** produce an error estimate for the *QMC rule based on a generating vector constructed by the CBC algorithm!*

## Theorem (CBC error bound)

The generating vector  $\mathbf{z} \in \mathbb{U}_n^s$  constructed by the CBC algorithm, minimizing the squared shift-averaged worst-case error  $[e_{n,s}^{\text{sh}}(\mathbf{z})]^2$  for the weighted unanchored Sobolev space in each step, satisfies

$$[e_{n,s}^{\text{sh}}(\mathbf{z})]^2 \leq \left( \frac{1}{\varphi(n)} \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}}^{\lambda} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^{\lambda}} \right)^{|\mathfrak{u}|} \right)^{1/\lambda} \quad \text{for all } \lambda \in (1/2, 1],$$

where  $\zeta(x) := \sum_{k=1}^{\infty} k^{-x}$  denotes the Riemann zeta function for  $x > 1$ .

**Significance:** Suppose that  $f \in H_{s,\gamma}$  for all  $\gamma = (\gamma_u)_{u \subseteq \{1:s\}}$ . Then for any given sequence of weights  $\gamma$ , we can use the CBC algorithm to obtain a generating vector satisfying the error bound

$$\sqrt{\mathbb{E}_\Delta |I_s f - Q_{n,s}^\Delta f|^2} \leq \left( \frac{1}{\varphi(n)} \sum_{\emptyset \neq u \subseteq \{1:s\}} \gamma_u^\lambda \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{|u|} \right)^{1/(2\lambda)} \|f\|_{s,\gamma} \quad (7)$$

for all  $\lambda \in (1/2, 1]$ . We can use the following strategy:

- For a given integrand  $f$ , estimate the norm  $\|f\|_{s,\gamma}$ .
- Find weights  $\gamma$  which *minimize* the error bound (7).
- Using the optimized weights  $\gamma$  as input, use the CBC algorithm to find a generating vector which *satisfies* the error bound (7).

### Remarks:

- If  $n$  is prime, then  $\frac{1}{\varphi(n)} = \frac{1}{n-1}$ . If  $n = 2^k$ , then  $\frac{1}{\varphi(n)} = \frac{2}{n}$ . For general (composite)  $n \geq 3$ ,  $\frac{1}{\varphi(n)} \leq \frac{e^\gamma \log \log n + \frac{3}{\log \log n}}{n}$ , where  $\gamma = 0.57721566\dots$  (Euler–Mascheroni constant).
- The optimal convergence rate close to  $\mathcal{O}(n^{-1})$  is obtained with  $\lambda \rightarrow 1/2$ , but note that  $\lambda = 1/2$  is not permitted since  $\zeta(2\lambda) \rightarrow \infty$  as  $\lambda \rightarrow 1/2$ .

### 3. Constructing lattice rules

## Naïve CBC construction

We write the error criterion as

$$\begin{aligned}
 [e_{n,d}^{\text{sh}}(z_1, \dots, z_d)]^2 &= \frac{1}{n} \sum_{k=0}^{n-1} \sum_{\emptyset \neq u \subseteq \{1:d\}} \gamma_u \prod_{j \in u} B_2 \left( \left\{ \frac{kz_j}{n} \right\} \right) \\
 &= \frac{1}{n} \sum_{k=0}^{n-1} \sum_{\ell=1}^d \underbrace{\sum_{\substack{|u|=\ell \\ u \subseteq \{1:d\}}} \gamma_u \prod_{j \in u} B_2 \left( \left\{ \frac{kz_j}{n} \right\} \right)}_{=: p_{d,\ell}(k)}.
 \end{aligned}$$

By plugging in the POD weights  $\gamma_u := \Gamma_{|u|} \prod_{j \in u} \gamma_j$ , note that we have the following recursion (we split the sum over  $u$  in two parts depending on whether  $d \in u$ ):

$$\begin{aligned}
 p_{d,\ell}(k) &= \sum_{\substack{|u|=\ell \\ u \subseteq \{1:d\}}} \Gamma_\ell \left( \prod_{j \in u} \gamma_j B_2 \left( \left\{ \frac{kz_j}{n} \right\} \right) \right) \\
 &= \sum_{\substack{|u|=\ell \\ u \subseteq \{1:d-1\}}} \Gamma_\ell \left( \prod_{j \in u} \gamma_j B_2 \left( \left\{ \frac{kz_j}{n} \right\} \right) \right) \\
 &\quad + \sum_{\substack{|u|=\ell-1 \\ u \subseteq \{1:d-1\}}} \Gamma_\ell \gamma_d B_2 \left( \left\{ \frac{kz_d}{n} \right\} \right) \left( \prod_{j \in u} \gamma_j B_2 \left( \left\{ \frac{kz_j}{n} \right\} \right) \right) \\
 &= p_{d-1,\ell}(k) + \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d B_2 \left( \left\{ \frac{kz_d}{n} \right\} \right) p_{d-1,\ell-1}(k).
 \end{aligned}$$

Plugging the recurrence

$$p_{d,\ell}(k) = p_{d-1,\ell}(k) + \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d B_2\left(\left\{ \frac{kz_d}{n} \right\}\right) p_{d-1,\ell-1}(k)$$

into the expression for the squared shift-averaged WCE yields

$$\begin{aligned} [e_{n,d}^{\text{sh}}(z_1, \dots, z_d)]^2 &= \frac{1}{n} \sum_{k=0}^{n-1} \sum_{\ell=1}^d p_{d,\ell}(k) \\ &= \frac{1}{n} \sum_{k=0}^{n-1} \sum_{\ell=1}^d p_{d-1,\ell}(k) + \frac{1}{n} \sum_{k=0}^{n-1} \sum_{\ell=1}^d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d B_2\left(\left\{ \frac{kz_d}{n} \right\}\right) p_{d-1,\ell-1}(k) \\ &= [e_{n,d-1}^{\text{sh}}(z_1, \dots, z_{d-1})]^2 + \frac{1}{n} \sum_{k=0}^{n-1} B_2\left(\left\{ \frac{kz_d}{n} \right\}\right) \sum_{\ell=1}^d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d p_{d-1,\ell-1}(k). \end{aligned}$$

Recall that in the  $d^{\text{th}}$  step of the CBC algorithm, the components  $z_1, \dots, z_{d-1}$  are fixed and it is therefore sufficient to find  $z_d \in \mathbb{U}_n$  which minimizes the expression  $\sum_{k=0}^{n-1} B_2\left(\left\{ \frac{kz_d}{n} \right\}\right) \sum_{\ell=1}^d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d p_{d-1,\ell-1}(k)$ .

Let us introduce the matrix  $\Omega_n := \left[ B_2\left(\left\{ \frac{kz}{n} \right\} \right) \right]_{k \in \{0, \dots, n-1\}, z \in \mathbb{U}_n}$  and define a set of  $n$ -vectors recursively via

$$\mathbf{p}_{d,\ell} = \mathbf{p}_{d-1,\ell} + \gamma_d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \Omega_n(z_d, :) * \mathbf{p}_{d-1,\ell-1}$$

starting from the initial values

$$\mathbf{p}_{d,0} = \mathbf{1}_n \quad \text{for all } d \geq 1,$$

$$\mathbf{p}_{d,\ell} = \mathbf{0}_n \quad \text{for all } d \geq 1 \text{ and } \ell > d,$$

with  $.*$  denoting the componentwise product between two vectors.

Then the value of  $\sum_{k=0}^{n-1} B_2\left(\left\{ \frac{kz_d}{n} \right\} \right) \sum_{\ell=1}^d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d \mathbf{p}_{d-1,\ell-1}(k)$  in the  $d^{\text{th}}$  step of the CBC algorithm can be obtained for all  $z_d \in \mathbb{U}_n$  via

$$\Omega_n \mathbf{x}, \quad \text{where } \mathbf{x} = \sum_{\ell=1}^d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d \mathbf{p}_{d-1,\ell-1}.$$

## CBC algorithm – naïve version

1. Define the matrix  $\Omega_n := \left[ B_2\left(\left\{\frac{kz}{n}\right\}\right) \right]_{k \in \{0, \dots, n-1\}, z \in \mathbb{U}_n}$  and initialize the  $n$ -vectors

$$\mathbf{p}_{d,0} = \mathbf{1}_n \quad \text{for all } d \geq 1,$$

$$\mathbf{p}_{d,\ell} = \mathbf{0}_n \quad \text{for all } d \geq 1 \text{ and } \ell > d.$$

**for**  $d = 1, \dots, s$ , **do**

2. Pick the value  $z_d \in \{1, \dots, n - 1\}$  corresponding to the smallest entry in the matrix-vector product

$$\Omega_n \mathbf{x}, \quad \text{where } \mathbf{x} = \sum_{\ell=1}^d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \gamma_d \mathbf{p}_{d-1,\ell-1}. \quad (8)$$

3. Update  $\mathbf{p}_{d,\ell} = \mathbf{p}_{d-1,\ell} + \gamma_d \frac{\Gamma_\ell}{\Gamma_{\ell-1}} \Omega_n(z_d, :) * \mathbf{p}_{d-1,\ell-1}$ .

**end for**

**Remarks:** We only need the ratio  $a_\ell := \frac{\Gamma_\ell}{\Gamma_{\ell-1}}$  for the implementation, e.g., for  $\Gamma_\ell = \ell!$  this is  $a_\ell = \ell$ . The computational bottleneck is the dense matrix-vector product  $\Omega_n \mathbf{x}$  in (8), which has complexity  $\mathcal{O}(n^2)$ . The *fast CBC algorithm* reduces this product down to  $\mathcal{O}(n \log n)$  complexity.

Fast CBC algorithm

## What makes fast CBC fast?

The matrix-vector product  $\Omega_n \mathbf{x}$  has time complexity  $\mathcal{O}(n^2)$ , which is too slow if  $n$  is, say, of the order of a million or more. (Not to mention the problem of storing a dense matrix of such size!)

However, the matrix  $\Omega_n$  has a lot of structure. It turns out that we can implement the matrix-vector product  $\Omega_n \mathbf{x}$  in  $\mathcal{O}(n \log n)$  time using some sophisticated mathematical tools.

In a nutshell, we let  $n \geq 3$  be *prime* and do the following:

- Using some natural symmetries of  $\Omega_n$ , we can ignore the first column (since it corresponds to shifting the objective functional in the CBC minimization step by a constant value) and it will be sufficient to consider only the top-left block  $\Omega'_n := \Omega_n(1 : m, 2 : m + 1)$ , where  $m := (n - 1)/2$ .
- For *prime*  $n$ , we can find a *generator*  $g$  (primitive root modulo  $n$ ) and use this to permute  $\Omega'_n$  into a circulant matrix.
- A circulant matrix implements a circular convolution, so a matrix-vector product (in the permuted indexing) can be implemented in  $\mathcal{O}(n \log n)$  time using the fast Fourier transform (FFT).

Before getting to the implementational details of fast CBC, we will need to

- discuss an algorithm to find a primitive root modulo  $n$ ;
- discuss how to compute a circulant matrix-vector product using FFT.

# Primitive root modulo $n$

## Definition

Let  $g, n \in \mathbb{N}$ . The number  $g$  is called a *primitive root modulo  $n$*  if for any integer  $a \in \mathbb{N}$  such that  $\gcd(a, n) = 1$ , there exists an integer  $k$  (called the *index*) such that

$$g^k \equiv a \pmod{n}.$$

Such a number  $g$  is the *generator* of the multiplicative group of integers modulo  $n$ , i.e.,  $(\mathbb{Z}/n\mathbb{Z})^\times$ .

## Theorem (Gauss 1801)

A primitive root modulo  $n$  exists if and only if

- $n$  is 1, 2, 4, or
- $n = p^k$ , where  $p \geq 3$  is a prime and  $k \in \mathbb{N}$ , or
- $n = 2p^k$ , where  $p \geq 3$  is a prime and  $k \in \mathbb{N}$ .

Note especially that a primitive root modulo  $n$  exists whenever  $n$  is prime.

Recall that the Euler totient function is defined by

$\varphi(n) := |\{k \in \mathbb{N} \mid 1 \leq k \leq n-1, \gcd(k, n) = 1\}|$ . We have the following.

### Proposition

*The number  $g$  is a primitive root modulo  $n$  if and only if the smallest positive integer  $k$  for which  $g^k \equiv 1 \pmod{n}$  is precisely  $k = \varphi(n)$ .*

**Lagrange's theorem:** the smallest  $k$  satisfying  $g^k \equiv 1 \pmod{n}$  divides  $\varphi(n)$ . Therefore, it is enough to check for all proper divisors  $d|\varphi(n)$  that  $g^d \not\equiv 1 \pmod{n}$ .

*However, we can do even better!*

Find the prime number factorization  $\varphi(n) = p_1^{a_1} \cdots p_\ell^{a_\ell}$ . It turns out that it is enough to check that  $g^d \not\equiv 1 \pmod{n}$  for all  $d \in \left\{ \frac{\varphi(n)}{p_1}, \dots, \frac{\varphi(n)}{p_\ell} \right\}$ . To see this, let  $d$  be any proper divisor of  $\varphi(n)$ . Then there exists  $j$  such that  $d \mid \frac{\varphi(n)}{p_j}$ , meaning that  $dk = \frac{\varphi(n)}{p_j}$  for some  $k \in \mathbb{N}$ . However, if  $g^d \equiv 1 \pmod{n}$ , we would get

$$g^{\frac{\varphi(n)}{p_j}} \equiv g^{dk} \equiv (g^d)^k \equiv 1^k \equiv 1 \pmod{n}.$$

That is, if  $g$  was not a primitive root, then one could find a number of the form  $\frac{\varphi(n)}{p_j}$  for which  $g^{\frac{\varphi(n)}{p_j}} \equiv 1 \pmod{n}$ .

$\therefore$  It is enough to check that  $g^{\frac{\varphi(n)}{p_j}} \not\equiv 1 \pmod{n}$  for all  $j \in \{1, \dots, \ell\}$ .

## Algorithm for finding a primitive root modulo $n$

1. Find the prime number factorization  $\varphi(n) = p_1^{a_1} \cdots p_\ell^{a_\ell}$ .

Iterate through all numbers  $g = 1, 2, \dots, n - 1$  and, for each number, check whether it is a primitive root by doing the following:

2. Calculate  $\text{mod}(g^{\frac{\varphi(n)}{p_j}}, n)$  for all  $j \in \{1, \dots, \ell\}$ .
3. If all the calculated values are different from 1, then  $g$  is a primitive root.

**Remark:** In MATLAB, the quantities in step 2 can be computed, e.g., via `powermod(g,eulerPhi(n)/pj,n)`

## Discrete and fast Fourier transform

The *discrete Fourier transform* of (complex) vector  $\mathbf{x} := (x_j)_{j=1}^n$  is defined as the vector  $\mathbf{y} := (y_j)_{j=1}^n$  with

$$y_j = \sum_{k=1}^n x_k e^{-2\pi i(j-1)(k-1)/n}, \quad j \in \{1, \dots, n\},$$

and the *inverse discrete Fourier transform* is given by

$$x_j = \frac{1}{n} \sum_{k=1}^n y_k e^{2\pi i(j-1)(k-1)/n}, \quad j \in \{1, \dots, n\}.$$

The *fast Fourier transform (FFT)* can be used to carry out these operations in  $\mathcal{O}(n \log n)$  time. In MATLAB, one has  $\mathbf{y} = \text{fft}(\mathbf{x})$  and  $\mathbf{x} = \text{ifft}(\mathbf{y})$ .

## Circular convolution

Let  $\mathbf{x} := (x_i)_{i=1}^n$  and  $\mathbf{y} := (y_i)_{i=1}^n$  be (complex) vectors. Then the sequence  $\mathbf{z} := (z_i)_{i=1}^n$  defined by

$$z_i = \sum_{k=1}^n x_k y_{\text{mod}(i-k, n)+1}, \quad i \in \{1, \dots, n\},$$

is called the *circular convolution* of  $\mathbf{x}$  and  $\mathbf{y}$  and we denote it by  $\mathbf{z} := \mathbf{x} \star \mathbf{y}$ .

Similarly to the continuous convolution, we have the following identity using discrete/fast Fourier transform:

$$\text{fft}(\mathbf{x} \star \mathbf{y}) = \text{fft}(\mathbf{x}).*\text{fft}(\mathbf{y}),$$

where  $\mathbf{x}.*\mathbf{y} := (x_i y_i)_{i=1}^n$  is the pointwise product of two vectors.

## Circular convolution and circulant matrices

A matrix  $A \in \mathbb{R}^{n \times n}$  is called *circulant* if it has the form

$$A = \begin{bmatrix} a_0 & a_{n-1} & \cdots & a_2 & a_1 \\ a_1 & a_0 & a_{n-1} & & a_2 \\ \vdots & a_1 & a_0 & \ddots & \vdots \\ a_{n-2} & & \ddots & \ddots & a_{n-1} \\ a_{n-1} & a_{n-2} & \cdots & a_1 & a_0 \end{bmatrix}.$$

- Each row is equal to the row above shifted to the right by one (wrapping around the edge in a periodic way).
- The first column/row contains all information about the matrix.
- A circulant matrix implements a circular convolution:

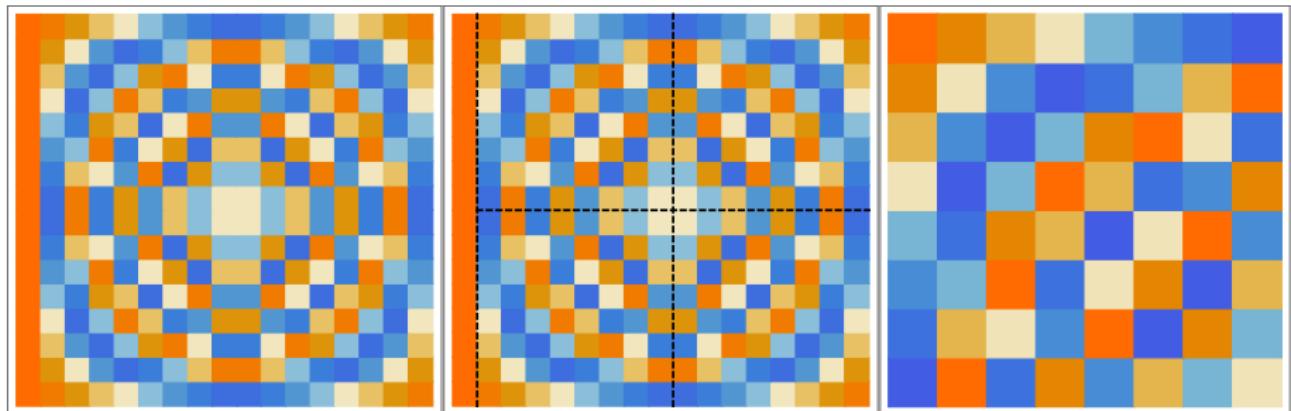
$$Ax = \mathbf{a} \star \mathbf{x}, \tag{9}$$

where  $\mathbf{a} := [a_0, a_1, \dots, a_{n-1}]^T$  is the first column of matrix  $A$ .

- The identity (9) implies that a circulant matrix-vector product can be implemented in  $\mathcal{O}(n \log n)$  time as  $A\mathbf{x} = \text{ifft}(\text{fft}(\mathbf{a}) \cdot \text{fft}(\mathbf{x}))$ .

## Putting it all together

The matrix-vector product  $\Omega_n \mathbf{x}$  in the CBC loop costs  $\mathcal{O}(n^2)$  operations. However, it was shown by Kuo, Nuyens, and Cools (2006) that the blocks of  $\Omega_n$  can be permuted into circulant form → the matrix-vector product can be implemented in  $\mathcal{O}(n \log n)$  operations using FFT.

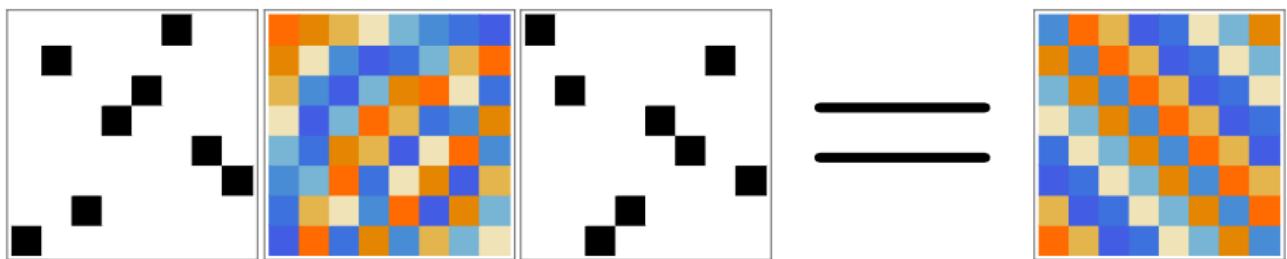


**Figure:** Example with  $\Omega_{17}$ . Note that the first column is a constant and can be left out (the components of  $\Omega_n \mathbf{x}$  are shifted by a constant → the smallest component stays invariant). Noting the obvious symmetries in the remaining four blocks, we can focus on the top left block.

When  $n$  is prime, it is possible to use the so-called Rader transformation to permute the top-left  $m \times m$  matrix  $\Omega'$  into circulant form:

$$\Omega_n^g(i, j) = \Omega'_n(g^i, (g^{-1})^j), \quad i, j \in \{1, \dots, m\},$$

where  $g$  is the primitive root modulo  $n$ . Here,  $g^{-1}$  denotes the modular multiplicative inverse  $gg^{-1} \equiv 1 \pmod{n}$ .



**Figure:** The original block matrix is multiplied from both sides by Rader permutation matrices (the black elements indicate the value 1 and white elements indicate the value 0) to obtain a circulant matrix.

## Example with $n = 1009$

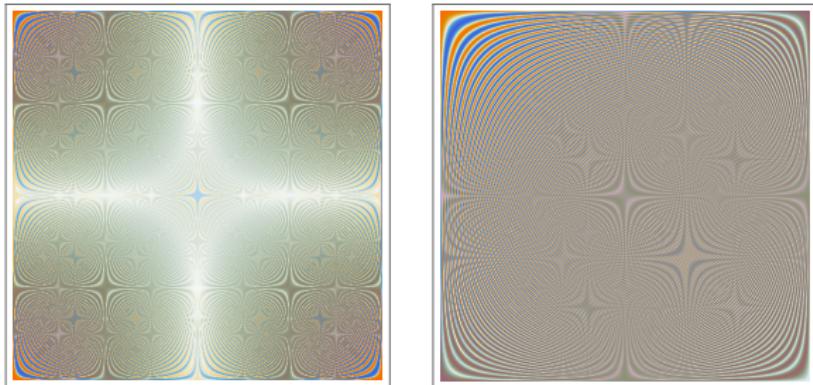


Figure: LHS: Original  $\Omega_{1009}$ . RHS: top left block of  $\Omega_{1009}$  (sans first column).

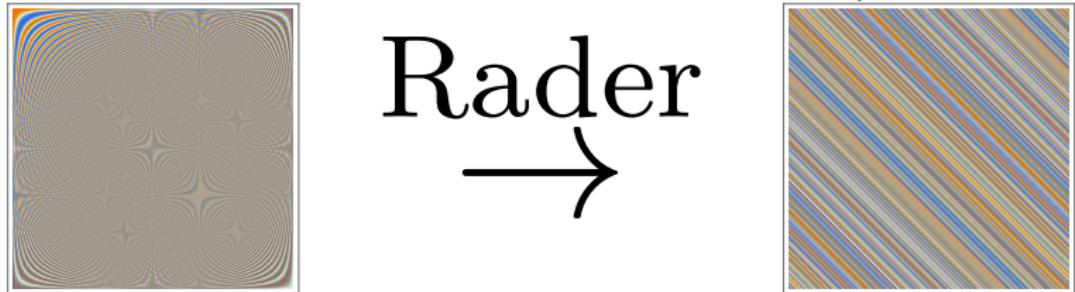


Figure: Rader transformation turns the top left block matrix circulant.

- The overall cost of the CBC algorithm with POD weights is  $\mathcal{O}(s n \log n + s^2 n)$ .
- For simplicity, we considered only the case where  $n$  is prime. An extension for composite  $n$  was discussed by Nuyens and Cools (J. Complexity 2006). The idea for composite  $n$  is that the complete matrix  $\Omega_n$  can be partitioned in blocks which have a circulant or block-circulant structure. The special case of  $n$  being a power of 2 has been discussed by Cools, Kuo, and Nuyens (SIAM J. Sci. Comput. 2006).
- There also exist freely available software implementing the fast CBC construction, cf., e.g.,

<https://people.cs.kuleuven.be/~dirk.nuyens/qmc4pde/>,  
<https://people.cs.kuleuven.be/~dirk.nuyens/fast-cbc/>,  
<https://qmcpy.org/>, ...

#### 4. QMC methods for forward and inverse uncertainty quantification of elliptic PDEs with random coefficients

Uniform and affine model

*Uniform and affine model:* let  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be a bounded Lipschitz domain, let  $f \in H^{-1}(D)$ , and let

$U := [-1/2, 1/2]^{\mathbb{N}} := \{(a_j)_{j \geq 1} : -1/2 \leq a_j \leq 1/2\}$  be a set of parameters.

Consider the problem of finding, for all  $\mathbf{y} \in U$ ,  $u(\cdot, \mathbf{y}) \in H_0^1(D)$  such that

$$\int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} = \langle f, v \rangle_{H^{-1}(D), H_0^1(D)} \quad \text{for all } v \in H_0^1(D),$$

where the diffusion coefficient has the parameterization

$$a(\mathbf{x}, \mathbf{y}) := a_0(\mathbf{x}) + \sum_{j=1}^{\infty} y_j \psi_j(\mathbf{x}), \quad \mathbf{x} \in D, \quad \mathbf{y} \in U,$$

where  $a_0 \in L^\infty(D)$ , there exist  $a_{\min}, a_{\max} > 0$

s.t.  $0 < a_{\min} \leq a(\mathbf{x}, \mathbf{y}) \leq a_{\max} < \infty$  for all  $\mathbf{x} \in D$  and  $\mathbf{y} \in U$ , and the *stochastic fluctuations*  $\psi_j: D \rightarrow \mathbb{R}$  are functions of the spatial variable such that

- $\psi_j \in L^\infty(D)$  for all  $j \in \mathbb{N}$ ,
- $\sum_{j=1}^{\infty} \|\psi_j\|_{L^\infty(D)} < \infty$ ,
- $\sum_{j=1}^{\infty} \|\psi_j\|_{L^\infty(D)}^p < \infty$  for some  $p \in (0, 1)$ .

## Total error decomposition

In practice, we need to truncate the infinite-dimensional parametric vector  $\mathbf{y} \in [-1/2, 1/2]^{\mathbb{N}}$  to a finite number of terms. Moreover, the PDE needs to be discretized spatially using, e.g., the finite element method.

Let  $u_s(\cdot, \mathbf{y}) := u_s(y_1, \dots, y_s, 0, 0, \dots)$  denote the dimensionally-truncated PDE solution for  $\mathbf{y} \in [-1/2, 1/2]^s$ , and let  $u_{s,h}(\cdot, \mathbf{y}) \in V_h$  denote the dimensionally-truncated FE solution in the FE subspace spanned by piecewise linear FE basis functions. Furthermore, let

$(\mathbf{t}_i)_{i=1}^n = (\{\frac{i\mathbf{z}}{n}\} - \frac{1}{2})_{i=1}^n$  be a QMC point set in  $[-1/2, 1/2]^s$ .

## Total error decomposition

For simplicity, let us consider the problem of computing  $\mathbb{E}[G(u)]$ , where  $u(\cdot, \mathbf{y}) \in H_0^1(D)$  is the PDE solution for  $\mathbf{y} \in [-1/2, 1/2]^{\mathbb{N}}$  and  $G: H_0^1(D) \rightarrow \mathbb{R}$  is a linear functional (quantity of interest). We decompose the total error as

$$\begin{aligned}& \int_{[-1/2,1/2]^{\mathbb{N}}} G(u(\cdot, \mathbf{y})) d\mathbf{y} - \frac{1}{n} \sum_{i=1}^n G(u_{s,h}(\cdot, \mathbf{t}_i)) \\&= \int_{[-1/2,1/2]^{\mathbb{N}}} (G(u(\cdot, \mathbf{y}) - u_s(\cdot, \mathbf{y}_{\leq s}))) d\mathbf{y} \\&+ \int_{[-1/2,1/2]^s} G(u_s(\cdot, \mathbf{y}) - u_{s,h}(\cdot, \mathbf{y})) d\mathbf{y} \\&+ \int_{[-1/2,1/2]^s} G(u_{s,h}(\cdot, \mathbf{y})) d\mathbf{y} - \frac{1}{n} \sum_{i=1}^n G(u_{s,h}(\cdot, \mathbf{t}_i)).\end{aligned}$$

Using the triangle inequality, we are left with the total error decomposition

$$\begin{aligned}
 & \left| \int_{[-1/2,1/2]^{\mathbb{N}}} G(u(\cdot, \mathbf{y})) d\mathbf{y} - \frac{1}{n} \sum_{i=1}^n G(u_{s,h}(\cdot, \mathbf{t}_i)) \right| \\
 & \leq \left| \int_{[-1/2,1/2]^{\mathbb{N}}} (G(u(\cdot, \mathbf{y}) - u_s(\cdot, \mathbf{y}_{\leq s})) d\mathbf{y} \right| \quad (\text{dimension-truncation error}) \\
 & + \left| \int_{[-1/2,1/2]^s} G(u_s(\cdot, \mathbf{y}) - u_{s,h}(\cdot, \mathbf{y})) d\mathbf{y} \right| \quad (\text{finite element error}) \\
 & + \left| \int_{[-1/2,1/2]^s} G(u_{s,h}(\cdot, \mathbf{y})) d\mathbf{y} - \frac{1}{n} \sum_{i=1}^n G(u_{s,h}(\cdot, \mathbf{t}_i)) \right|. \quad (\text{cubature error})
 \end{aligned}$$

We focus on the cubature error.

- If  $\|\psi_1\|_{L^\infty(D)} \geq \|\psi_2\|_{L^\infty(D)} \geq \dots$ , then the dimension truncation error decays like  $\mathcal{O}(s^{-2/p+1})$  [Gantner (2018)].
- If  $D$  is a convex polygon (2d)/polyhedron (3d) and we have additional regularity, e.g.,  $f, G \in L^2(D)$ ,  $a$  is Lipschitz, and the family  $\{V_h\}_h$  of first-order finite element spaces, indexed by the mesh size  $h > 0$ , is a sequence of regular, simplicial meshes in  $D$  obtained from an initial, regular triangulation of  $D$  by recursive, uniform bisection of simplices, then the  $L^2$  finite element error satisfies  $\mathcal{O}(h^2)$  as  $h \rightarrow 0$  independently of  $s$  [Kuo, Schwab, Sloan (2012)].

## Multi-index notation

We introduce the set of *finitely-supported multi-indices*

$$\mathcal{F} := \{\boldsymbol{\nu} \in \mathbb{N}_0^{\mathbb{N}} : |\text{supp}(\boldsymbol{\nu})| < \infty\},$$

where the *support* of a multi-index  $\boldsymbol{\nu}$  is defined as the set

$$\text{supp}(\boldsymbol{\nu}) := \{i \in \mathbb{N} : \nu_i \neq 0\}.$$

As before, the *order* of a multi-index is defined as

$$|\boldsymbol{\nu}| := \sum_{j \geq 1} \nu_j$$

and we use the special multi-index notations

$$\partial^{\boldsymbol{\nu}} := \partial_{\mathbf{y}}^{\boldsymbol{\nu}} := \prod_{j \in \text{supp}(\boldsymbol{\nu})} \frac{\partial^{\nu_j}}{\partial y_j^{\nu_j}}, \quad \mathbf{x}^{\boldsymbol{\nu}} := \prod_{j \in \text{supp}(\boldsymbol{\nu})} x_j^{\nu_j}, \quad \binom{\boldsymbol{\nu}}{\mathbf{m}} := \prod_{j \in \text{supp}(\boldsymbol{\nu})} \binom{\nu_j}{m_j}.$$

## Recursive bound

Consider the weak formulation

$$\int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} = \langle f, v \rangle_{H^{-1}(D), H_0^1(D)}. \quad (10)$$

Noting that

$$\partial^\nu a(\mathbf{x}, \mathbf{y}) = \begin{cases} a(\mathbf{x}, \mathbf{y}) & \text{if } \nu = \mathbf{0}, \\ \psi_j(\mathbf{x}) & \text{if } \nu = \mathbf{e}_j, \\ 0 & \text{otherwise,} \end{cases}$$

we let  $\nu \in \mathcal{F} \setminus \{\mathbf{0}\}$  and differentiate (10) on both sides with  $\partial^\nu$  and use the Leibniz product rule<sup>†</sup> to obtain

$$\partial^\nu \int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} = 0$$

$$\Leftrightarrow \sum_{\mathbf{m} \leq \nu} \binom{\nu}{\mathbf{m}} \int_D \partial^\mathbf{m} a(\mathbf{x}) \nabla \partial^{\nu-\mathbf{m}} u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} = 0$$

$$\Leftrightarrow \int_D a(\mathbf{x}, \mathbf{y}) \nabla \partial^\nu u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x} = - \sum_{j \in \text{supp}(\nu)} \nu_j \int_D \psi_j(\mathbf{x}) \nabla \partial^{\nu - \mathbf{e}_j} u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) d\mathbf{x}.$$

---

<sup>†</sup> $\partial^\nu(fg) = \sum_{\mathbf{m} \leq \nu} \binom{\nu}{\mathbf{m}} \partial^\mathbf{m} f \partial^{\nu-\mathbf{m}} g$

Testing this against  $v = \partial^\nu u(\mathbf{x}, \mathbf{y})$  yields

$$\begin{aligned} & a_{\min} \|\partial^\nu u(\cdot, \mathbf{y})\|_{H_0^1(D)}^2 \\ & \leq \int_D a(\mathbf{x}, \mathbf{y}) \|\nabla \partial^\nu u(\mathbf{x}, \mathbf{y})\|^2 d\mathbf{x} \\ & \leq \sum_{j \in \text{supp}(\nu)} \nu_j \|\psi_j\|_{L^\infty(D)} \|\partial^{\nu - e_j} u(\cdot, \mathbf{y})\|_{H_0^1(D)} \|\partial^\nu u(\cdot, \mathbf{y})\|_{H_0^1(D)} \end{aligned}$$

Thus we obtain the recursive relation

$$\|\partial^\nu u(\cdot, \mathbf{y})\|_{H_0^1(D)} \leq \sum_{j \in \text{supp}(\nu)} \nu_j \underbrace{\frac{\|\psi_j\|_{L^\infty(D)}}{a_{\min}}}_{=: b_j} \|\partial^{\nu - e_j} u(\cdot, \mathbf{y})\|_{H_0^1(D)}.$$

For later convenience, we introduce here the sequence  $\mathbf{b} := (b_j)_{j \geq 1}$  defined by  $b_j := \frac{\|\psi_j\|_{L^\infty(D)}}{a_{\min}}$ . Recall that by the assumptions we placed on the uniform and affine model, there holds  $\mathbf{b} \in \ell^p$  for some  $p \in (0, 1)$ .

# Parametric regularity

## Proposition

For all  $\mathbf{y} \in [-1/2, 1/2]^{\mathbb{N}}$  and  $\nu \in \mathcal{F}$ , there holds

$$\|\partial^\nu u(\cdot, \mathbf{y})\|_{H_0^1(D)} \leq \frac{\|f\|_{H^{-1}(D)}}{a_{\min}} b^\nu |\nu|!.$$

*Proof.* By induction w.r.t. the order of the multi-index  $\nu \in \mathcal{F}$ . If  $\nu = \mathbf{0}$ , then this is the ordinary Lax–Milgram *a priori* bound

$$\begin{aligned} a_{\min} \underbrace{\int_D |\nabla u(\mathbf{x}, \mathbf{y})|^2 d\mathbf{x}}_{= \|u(\cdot, \mathbf{y})\|_{H_0^1(D)}^2} &\leq \int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla u(\mathbf{x}, \mathbf{y}) d\mathbf{x} = \langle f, u(\cdot, \mathbf{y}) \rangle_{H^{-1}(D), H_0^1(D)} \\ &\leq \|f\|_{H^{-1}(D)} \|u(\cdot, \mathbf{y})\|_{H_0^1(D)} \end{aligned}$$

whence

$$\|u(\cdot, \mathbf{y})\|_{H_0^1(D)} \leq \frac{\|f\|_{H^{-1}(D)}}{a_{\min}}.$$

Next, let  $\nu \in \mathcal{F} \setminus \{\mathbf{0}\}$  and suppose that the claim has been proved for all multi-indices with order less than  $|\nu|$ . Then using the recursive relation we derived previously, we obtain

$$\begin{aligned} \|\partial^\nu u(\cdot, \mathbf{y})\|_{H_0^1(D)} &\leq \sum_{j \in \text{supp}(\nu)} \nu_j b_j \|\partial^{\nu - e_j} u(\cdot, \mathbf{y})\|_{H_0^1(D)} \\ &\leq \frac{\|f\|_{H^{-1}(D)}}{a_{\min}} \sum_{j \in \text{supp}(\nu)} \nu_j b_j |\nu - e_j|! b^{\nu - e_j} \\ &= \frac{\|f\|_{H^{-1}(D)}}{a_{\min}} b^\nu (|\nu| - 1)! \sum_{j \geq 1} \nu_j \\ &= \frac{\|f\|_{H^{-1}(D)}}{a_{\min}} b^\nu |\nu|!, \end{aligned}$$

as desired. □

**Remark.** Note that the same regularity bound holds for the dimensionally-truncated FE solution  $u_{s,h}$  as long as a (conforming) Galerkin FE discretization has been used to construct the FE approximation. This is due to the fact that the weak formulation of the Galerkin discretization is exactly the same (only the function space differs).

Now that we know the regularity of the PDE problem, we can analyze the QMC cubature error! Let  $G: H_0^1(D) \rightarrow \mathbb{R}$  be a linear and bounded functional,  $u_{s,h}$  the dimensionally-truncated FE solution, and define  $F(\mathbf{y}) := G(u_{s,h}(\cdot, \mathbf{y} - \frac{1}{2}))$  for  $\mathbf{y} \in [0, 1]^s$ . Let  $\gamma = (\gamma_u)_{u \subseteq \{1:s\}}$  be a sequence of positive weights. Then we know that the generating vector obtained by the CBC algorithm satisfies the error bound

$$\sqrt{\mathbb{E}_{\Delta} |I_s F - Q_{n,s}^{\Delta} F|^2} \leq \left( \frac{1}{\varphi(n)} \sum_{\emptyset \neq u \subseteq \{1:s\}} \gamma_u^{\lambda} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^{\lambda}} \right)^{|u|} \right)^{1/(2\lambda)} \|F\|_{s,\gamma}$$

for all  $\lambda \in (1/2, 1]$ , where

$$\begin{aligned} \|F\|_{s,\gamma}^2 &= \sum_{u \subseteq \{1:s\}} \frac{1}{\gamma_u} \int_{[0,1]^{|u|}} \left( \int_{[0,1]^{s-|u|}} \frac{\partial^{|u|}}{\partial \mathbf{x}_u} F(\mathbf{y}) d\mathbf{y}_{-u} \right)^2 d\mathbf{y}_u \\ &\leq \left( \frac{\|G\|_{H^{-1}(D)} \|f\|_{H^{-1}(D)}}{a_{\min}} \right)^2 \sum_{u \subseteq \{1:s\}} \frac{1}{\gamma_u} (|u|!)^2 \prod_{j \in u} b_j^2. \end{aligned}$$

Plugging this norm bound back into the QMC error bound yields...

$$\begin{aligned} \sqrt{\mathbb{E}_{\Delta} |I_s F - Q_{n,s}^{\Delta} F|^2} &\lesssim \left( \frac{1}{\varphi(n)} \right)^{1/(2\lambda)} \left( \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}}^{\lambda} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^{\lambda}} \right)^{|\mathfrak{u}|} \right)^{1/(2\lambda)} \\ &\quad \times \left( \sum_{\mathfrak{u} \subseteq \{1:s\}} \frac{1}{\gamma_{\mathfrak{u}}} (|\mathfrak{u}|!)^2 \prod_{j \in \mathfrak{u}} b_j^2 \right)^{1/2}. \end{aligned}$$

The upper bound can be *minimized* by choosing the *POD weights*

$$\gamma_{\mathfrak{u}} := \left( |\mathfrak{u}|! \prod_{j \in \mathfrak{u}} \frac{b_j}{\sqrt{\frac{2\zeta(2\lambda)}{(2\pi^2)^{\lambda}}}} \right)^{2/(1+\lambda)},$$

as explained by the following lemma.

### Lemma

Let  $(\alpha_i)$  and  $(\beta_i)$  be sequences of positive real numbers. The expression

$$g(\gamma) := \left( \sum_i \alpha_i \gamma_i^{\lambda} \right)^{1/\lambda} \left( \sum_i \beta_i \gamma_i^{-1} \right)$$

is minimized by  $\gamma_i = c \left( \frac{\beta_i}{\alpha_i} \right)^{1/(1+\lambda)}$  for arbitrary  $c > 0$ .

*Proof.* Let us find out when the gradient vanishes:

$$0 = \partial_j g(\boldsymbol{\gamma}) = \frac{1}{\lambda} \left( \sum_i \alpha_i \gamma_i^\lambda \right)^{1/\lambda - 1} \lambda \alpha_j \gamma_j^{\lambda-1} \left( \sum_i \beta_i \gamma_i^{-1} \right) \\ - \beta_j \gamma_j^{-2} \left( \sum_i \alpha_i \gamma_i^\lambda \right)^{1/\lambda}.$$

After some trivial simplifications, we can see that this is equivalent to

$$\gamma_j^{\lambda+1} = \frac{\beta_j}{\alpha_j} \frac{\sum_i \alpha_i \gamma_i^\lambda}{\sum_i \beta_i \gamma_i^{-1}}.$$

Furthermore, this condition is satisfied if

$$\gamma_j = c \left( \frac{\beta_j}{\alpha_j} \right)^{1/(1+\lambda)},$$

where  $c > 0$  is arbitrary. □

Note that plugging  $\gamma_i = c \left( \frac{\beta_i}{\alpha_i} \right)^{1/(1+\lambda)}$  into  $(\sum_i \alpha_i \gamma_i^\lambda)^{1/(2\lambda)} (\sum_i \beta_i \gamma_i^{-1})^{1/2}$  yields the expression  $(\sum_i \alpha_i^{1/(1+\lambda)} \beta_i^{\lambda/(1+\lambda)})^{(1+\lambda)/(2\lambda)}$ . Thus, plugging the optimal POD weights into the QMC error bound results in

$$\sqrt{\mathbb{E}_\Delta |I_s F - Q_{n,s}^\Delta F|^2} \lesssim \left( \frac{1}{\varphi(n)} \right)^{1/(2\lambda)} C(s, \gamma, \lambda)^{(1+\lambda)/(2\lambda)},$$

where

$$C(s, \gamma, \lambda) := \sum_{\mathfrak{u} \subseteq \{1:s\}} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{|\mathfrak{u}|/(1+\lambda)} (|\mathfrak{u}|!)^{2\lambda/(1+\lambda)} \prod_{j \in \mathfrak{u}} b_j^{2\lambda/(1+\lambda)}.$$

This is the punchline:

Lemma

*By choosing*

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

*there exists a constant  $C(\gamma, \lambda) < \infty$  independently of  $s$  s.t.  $C(s, \gamma, \lambda) \leq C(\gamma, \lambda) < \infty$ .*

*Proof.* First observe that

$$\begin{aligned}
 C(s, \gamma, \lambda) &= \sum_{\mathfrak{u} \subseteq \{1:s\}} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{|\mathfrak{u}|/(1+\lambda)} (|\mathfrak{u}|!)^{2\lambda/(1+\lambda)} \prod_{j \in \mathfrak{u}} b_j^{2\lambda/(1+\lambda)} \\
 &= \sum_{\ell=0}^s \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{\ell/(1+\lambda)} (\ell!)^{2\lambda/(1+\lambda)} \sum_{\substack{|\mathfrak{u}|=\ell \\ \mathfrak{u} \subseteq \{1:s\}}} \prod_{j \in \mathfrak{u}} b_j^{2\lambda/(1+\lambda)} \\
 &\leq \sum_{\ell=0}^{\infty} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{\ell/(1+\lambda)} (\ell!)^{2\lambda/(1+\lambda)-1} \left( \sum_{j \geq 1} b_j^{2\lambda/(1+\lambda)} \right)^\ell
 \end{aligned}$$

where we used the inequality  $\sum_{|\mathfrak{u}|=\ell, \mathfrak{u} \subseteq \mathbb{Z}_+} \prod_{j \in \mathfrak{u}} c_j \leq \frac{1}{\ell!} \left( \sum_{j \geq 1} c_j \right)^\ell$ .

**Case 1:**  $p \in (2/3, 1)$ . We choose  $p = \frac{2\lambda}{1+\lambda} \Leftrightarrow \lambda = \frac{p}{2-p} \in (1/2, 1)$ , and

$$C(s, \gamma, \lambda) \leq \underbrace{\sum_{\ell=0}^{\infty} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{\ell/(1+\lambda)} (\ell!)^{p-1} \left( \sum_{j \geq 1} b_j^p \right)^\ell}_{=: a_\ell}$$

It is easy to see that  $\frac{a_{\ell+1}}{a_\ell} \xrightarrow{\ell \rightarrow \infty} 0$ . By the ratio test, this upper bound is finite independently of  $s$ .

**Case 2:**  $p \in (0, 2/3]$ . Let  $\delta \in (0, 1/2)$  be arbitrary. We choose  $\lambda = \frac{1}{2-2\delta} \in (1/2, 1)$ . Now  $\frac{2\lambda}{1+\lambda} = \frac{2}{3-2\delta} \in (2/3, 1)$ . Especially,  $\|\mathbf{b}\|_{\ell^{2\lambda/(1+\lambda)}} \leq \|\mathbf{b}\|_{\ell^p}$ , and we obtain from the estimate on the previous slide that

$$\begin{aligned} C(s, \gamma, \lambda) &\leq \sum_{\ell=0}^{\infty} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{\ell/(1+\lambda)} (\ell!)^{2\lambda/(1+\lambda)-1} \left( \sum_{j \geq 1} b_j^{2\lambda/(1+\lambda)} \right)^\ell \\ &\leq \underbrace{\sum_{\ell=0}^{\infty} \left( \frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda} \right)^{\ell/(1+\lambda)} (\ell!)^{2/(3-2\delta)-1} \left( \sum_{j \geq 1} b_j^p \right)^{2\ell/((3-2\delta)p)}}_{=: a_\ell} \end{aligned}$$

Again,  $\frac{a_{\ell+1}}{a_\ell} \xrightarrow{\ell \rightarrow \infty} 0$ , so by the ratio test this upper bound is finite independently of  $s$ . □

## Theorem

Let  $\delta \in (0, 1/2)$  be arbitrary. By choosing the POD weights

$$\gamma_{\mathfrak{u}} := \left( |\mathfrak{u}|! \prod_{j \in \mathfrak{u}} \frac{b_j}{\sqrt{\frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda}}} \right)^{2/(1+\lambda)}, \quad \lambda := \begin{cases} \frac{p}{2-p} & \text{if } p \in (2/3, 1), \\ \frac{1}{2-2\delta} & \text{if } p \in (0, 2/3], \end{cases}$$

then the QMC approximation for the expected value of the PDE problem satisfies

$$\text{R.M.S. error} \lesssim \begin{cases} \left(\frac{1}{\varphi(n)}\right)^{1/p-1/2} & \text{if } p \in (2/3, 1), \\ \left(\frac{1}{\varphi(n)}\right)^{1-\delta} & \text{if } p \in (0, 2/3], \end{cases}$$

where the implied coefficient is independent of the dimension  $s$ .

**Remark:** We have the following dimension-independent convergence rates:

- $n$  is prime  $\Rightarrow \frac{1}{\varphi(n)} = \frac{1}{n-1} \Rightarrow$  QMC rate  $\mathcal{O}(n^{\max\{-1/p+1/2, -1+\delta\}})$ .
- $n = 2^k \Rightarrow \frac{1}{\varphi(n)} = \frac{2}{n} \Rightarrow$  QMC rate  $\mathcal{O}(n^{\max\{-1/p+1/2, -1+\delta\}})$ .
- For general composite  $n$ , the dimension-independent QMC rate is at best essentially linear up to a double logarithmic factor of  $n$ .

## Remarks on implementation

Let  $G: H_0^1(D) \rightarrow \mathbb{R}$  be a bounded linear functional. Consider the problem of approximating

$$\mathbb{E}[G(u_{s,h})] = \int_{[-1/2,1/2]^s} G(u_{s,h}(\cdot, \mathbf{y})) d\mathbf{y},$$

where  $u_{s,h}$  is the dimensionally-truncated FE approximation to the elliptic PDE with a uniform and affine diffusion coefficient.

Our QMC approximation is guaranteed to satisfy the R.M.S. error bound from the previous slide if we plug the theoretically derived weights as input to the fast CBC algorithm. This produces a generating vector  $\mathbf{z} \in \mathbb{N}^s$ . The generating vector is designed to be used to compute the estimate

$$\overline{Q}_{n,s,R} G(u_{s,h}) := \frac{1}{R} \sum_{r=0}^{R-1} Q_{n,s}^{\Delta_r} G(u_{s,h}),$$

where  $Q_{n,s}^{\Delta_r} F := \frac{1}{n} \sum_{i=0}^{n-1} f(\{\mathbf{t}_i + \Delta_r\} - \frac{1}{2})$ ,  $\mathbf{t}_k := \{\frac{k\mathbf{z}}{n}\}$ , and  $\Delta_0, \dots, \Delta_{R-1}$  are independent random shifts drawn from  $\mathcal{U}([0, 1]^s)$ .

- Typically, the number of random shifts is taken to be rather small, e.g.,  $8 \leq R \leq 64$ .
- A practical estimate for the R.M.S. error is given by the formula

$$\sqrt{\mathbb{E}_{\Delta} |I_s F - Q_{n,s}^{\Delta} F|^2} \approx \sqrt{\frac{1}{R(R-1)} \sum_{r=0}^{R-1} (Q_{n,s}^{\Delta_r} F - \bar{Q}_{n,s,R} F)^2}.$$

- If we instead wish to estimate  $\mathbb{E}[u_{s,h}(x, \cdot)]$  (i.e., leave out the *quantity of interest*  $G: H_0^1(D) \rightarrow \mathbb{R}$ ), the same weights can be used as input to the CBC algorithm.

## Numerical example

Let us consider the PDE problem

$$-\nabla \cdot (a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y})) = x_1, \quad u(\cdot, \mathbf{y})|_{\partial D} = 0,$$

in the physical domain  $D = (0, 1)^2$  with the diffusion coefficient

$$a(\mathbf{x}, \mathbf{y}) = 1 + \sum_{j=1}^s y_j \psi_j(\mathbf{x}), \quad \mathbf{x} \in D, \quad y_j \in [-\frac{1}{2}, \frac{1}{2}],$$

where  $\psi_j(\mathbf{x}) = j^{-2} \sin(j\pi x_1) \sin(j\pi x_2)$ . We compute  $\mathbb{E}[G(u)]$  using QMC with  $R = 8$  random shifts, where  $G(v) = \int_D v(\mathbf{x}) d\mathbf{x}$ .

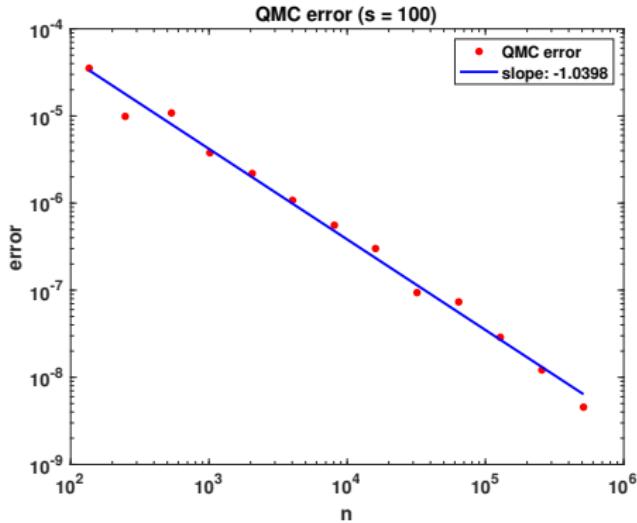


Figure: QMC with  $s = 100$  constructed using the weights

$$\gamma_{\mathfrak{u}} = \left( |\mathfrak{u}|! \prod_{j \in \mathfrak{u}} \frac{b_j}{\sqrt{2\zeta(2\lambda)/(2\pi^2)^\lambda}} \right)^{\frac{2}{1+\lambda}}, \quad \lambda = \frac{1}{2-2\delta}, \quad \delta = 0.05, \quad \text{for all } \mathfrak{u} \subseteq \{1, \dots, s\}.$$

Lognormal model (briefly)

*Lognormal model:* let  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be a bounded Lipschitz domain, and let  $f \in H^{-1}(D)$ . Let  $\psi_j \in L^\infty(D)$  and  $b_j := \|\psi_j\|_{L^\infty}$  for  $j \in \mathbb{N}$  such that  $\sum_{j=1}^{\infty} b_j < \infty$ , and set

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

Consider the problem of finding, for all  $\mathbf{y} \in U$ ,  $u(\cdot, \mathbf{y}) \in H_0^1(D)$  such that

$$\int_D a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x} = \langle f, v \rangle_{H^{-1}(D), H_0^1(D)} \quad \text{for all } v \in H_0^1(D),$$

where the diffusion coefficient is assumed to have the parameterization

$$a(\mathbf{x}, \mathbf{y}) := a_0(\mathbf{x}) \exp \left( \sum_{j=1}^{\infty} y_j \psi_j(\mathbf{x}) \right), \quad \mathbf{x} \in D, \quad \mathbf{y} \in U_b,$$

where  $a_0 \in L^\infty(D)$  is such that  $a_0(\mathbf{x}) > 0$ ,  $\mathbf{x} \in D$ .

## Standing assumptions for the lognormal model

- (B1) We have  $a_0 \in L^\infty(D)$  and  $\sum_{j=1}^{\infty} b_j < \infty$ .
- (B2) For every  $\mathbf{y} \in U_b$ , the expressions  $a_{\max}(\mathbf{y}) := \max_{\mathbf{x} \in \bar{D}} a(\mathbf{x}, \mathbf{y})$  and  $a_{\min}(\mathbf{y}) := \min_{\mathbf{x} \in \bar{D}} a(\mathbf{x}, \mathbf{y})$  are well-defined and satisfy  $0 < a_{\min}(\mathbf{y}) \leq a(\mathbf{x}, \mathbf{y}) \leq a_{\max}(\mathbf{y}) < \infty$ .
- (B3)  $\sum_{j=1}^{\infty} b_j^p < \infty$  for some  $p \in (0, 1)$ .

*Remark:* Note that in the lognormal case,  $a(\mathbf{x}, \mathbf{y})$  can take values which are arbitrarily close to 0 or arbitrarily large. Thus, the best we can do is to find  $\mathbf{y}$ -dependent lower and upper bounds  $a_{\min}(\mathbf{y})$  and  $a_{\max}(\mathbf{y})$ . This will lead to a  $\mathbf{y}$ -dependent *a priori* bound and, consequently,  $\mathbf{y}$ -dependent parametric regularity bounds. This will make the QMC analysis more involved, leading one to consider “special” weighted, unanchored Sobolev spaces.

In this setting, we have

$$I_s(F) := \int_{\mathbb{R}^s} F(\mathbf{y}) \prod_{j=1}^s \phi(y_j) d\mathbf{y} = \int_{(0,1)^s} F(\Phi^{-1}(\mathbf{w})) d\mathbf{w}.$$

where  $\phi(y) := \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y^2}$  is the probability density function of  $\mathcal{N}(0, 1)$  and  $\Phi^{-1}(\mathbf{w}) = [\Phi^{-1}(w_1), \dots, \Phi^{-1}(w_s)]^T$  denotes the corresponding (componentwise) inverse cumulative distribution function. We use the randomly shifted QMC rules

$$Q_{n,s}^{\Delta_r}(F) = \frac{1}{n} \sum_{k=1}^n F(\Phi^{-1}(\{\mathbf{t}_k + \Delta_r\})),$$

$$\overline{Q}_{n,R}(F) := \frac{1}{R} \sum_{r=1}^R Q_{n,s}^{\Delta_r}(F),$$

where we have  $R$  independent random shifts  $\Delta_1, \dots, \Delta_R$  drawn from  $\mathcal{U}([0, 1]^s)$ ,  $\mathbf{t}_k := \{\frac{k\mathbf{z}}{n}\}$ , with generating vector  $\mathbf{z} \in \mathbb{N}^s$ .

The appropriate function space for unbounded integrands is a “special” weighted, unanchored Sobolev space equipped with the norm

$$\|F\|_{s,\gamma} = \left[ \sum_{\mathfrak{u} \subseteq \{1:s\}} \frac{1}{\gamma_{\mathfrak{u}}} \int_{\mathbb{R}^{|\mathfrak{u}|}} \left( \int_{\mathbb{R}^{s-|\mathfrak{u}|}} \frac{\partial^{|\mathfrak{u}|}}{\partial \mathbf{y}_{\mathfrak{u}}} F(\mathbf{y}) \left( \prod_{j \in \{1:s\} \setminus \mathfrak{u}} \phi(y_j) \right) d\mathbf{y}_{-\mathfrak{u}} \right)^2 \times \left( \prod_{j \in \mathfrak{u}} \varpi_j^2(y_j) \right) d\mathbf{y}_{\mathfrak{u}} \right]^{1/2}$$

where we have the weights

$$\varpi_j^2(y) := \exp(-2\alpha_j|y_j|), \quad \alpha_j > 0.$$

## Theorem (Graham, Kuo, Nichols, Scheichl, Schwab, Sloan (2015))

Let  $F$  belong to the special weighted space over  $\mathbb{R}^s$  with weights  $\gamma$ , with  $\phi$  being the standard normal density, and the weight functions  $\varpi_j$  defined as above. A randomly shifted lattice rule in  $s$  dimensions with  $n$  being a prime power can be constructed by a CBC algorithm such that

$$\sqrt{\mathbb{E}_{\Delta}|I_s F - Q_{n,s}^{\Delta} F|^2} \leq \left( \frac{2}{n} \sum_{\emptyset \neq \mathfrak{u} \subseteq \{1:s\}} \gamma_{\mathfrak{u}}^{\lambda} \prod_{j \in \mathfrak{u}} \varrho_j(\lambda) \right)^{1/(2\lambda)} \|F\|_{s,\gamma},$$

where  $\lambda \in (1/2, 1]$  and

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

with  $\zeta(x) := \sum_{k=1}^{\infty} k^{-x}$  denoting the Riemann zeta function for  $x > 1$ .

The steps for QMC analysis are the same as in the uniform case: (1) estimate  $\|\cdot\|_{s,\gamma}$  for a given integrand (2) find weights  $\gamma$  which minimize the upper bound (3) plug the weights into the new error bound and estimate the constant (which ideally can be bounded independently of  $s$ ). 126

**Proposition (Parametric regularity bound for the lognormal model**  
**Graham, Kuo, Nichols, Scheichl, Schwab, Sloan (2015))**

For all  $\mathbf{y} \in U_b$  and  $\nu \in \mathcal{F}$ , there holds

$$\|\partial^\nu u(\cdot, \mathbf{y})\|_{H_0^1(D)} \leq \frac{\|f\|_{H^{-1}(D)}}{\min_{\mathbf{x} \in \bar{D}} a_0(\mathbf{x})} \frac{|\nu|!}{(\log 2)^{|\nu|}} b^\nu \prod_{j \geq 1} \exp(b_j |y_j|).$$

This parametric regularity bound is valid also for the dimensionally-truncated finite element solution  $u_{s,h}$ . If  $G: H_0^1(D) \rightarrow \mathbb{R}$  is a bounded linear functional and  $F(\mathbf{y}) := G(u_{s,h}(\cdot, \mathbf{y}))$  for  $\mathbf{y} \in \mathbb{R}^s$ , then

$$\|F\|_{s,\gamma}^2 \leq \sum_{\mathbf{u} \subseteq \{1:s\}} \frac{(|\mathbf{u}|!)^2}{\gamma_{\mathbf{u}}} \left( \prod_{j=1}^s 2 \exp(2b_j^2) \Phi(2b_j) \right) \left( \prod_{j \in \mathbf{u}} \frac{b_j^2}{2(\log 2)^2 \exp(2b_j^2) \Phi(2b_j) (\alpha_j - b_j)} \right).$$

By choosing  $\alpha_j = \frac{1}{2}(b_j + \sqrt{b_j^2 + 1 - \frac{1}{2\lambda}})$  and using the POD weights

$$\gamma_{\mathbf{u}} := \left( |\mathbf{u}|! \prod_{j \in \mathbf{u}} \frac{b_j}{2(\log 2) \exp(b_j^2/2) \Phi(b_j) \sqrt{(\alpha_j - b_j) \varrho_j(\lambda)}} \right)^{\frac{2}{1+\lambda}}, \quad \lambda := \begin{cases} \frac{p}{2-p} & \text{if } p \in (2/3, 1), \\ \frac{1}{2-2\delta} & \text{if } p \in (0, 2/3], \end{cases}$$

as inputs to the CBC algorithm yields a randomly shifted rank-1 lattice rule satisfying the R.M.S. error

$$\sqrt{\mathbb{E}_\Delta |I_s F - Q_{n,s}^\Delta F|^2} \lesssim n^{\max\{-1/p+1/2, -1+\delta\}},$$

where the constant is independent of the dimension.

Similarly to the uniform and affine setting, the truncation of the input random series and the finite element discretization incur a *dimension truncation error* and a *finite element discretization error*, respectively.

## Numerical example

Let us consider the PDE problem

$$-\nabla \cdot (a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y})) = x_1, \quad u(\cdot, \mathbf{y})|_{\partial D} = 0,$$

in the physical domain  $D = (0, 1)^2$  with the diffusion coefficient

$$a(\mathbf{x}, \mathbf{y}) = \exp \left( \sum_{j=1}^s y_j \psi_j(\mathbf{x}) \right), \quad \mathbf{x} \in D, \quad y_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1),$$

where  $\psi_j(\mathbf{x}) = j^{-2} \sin(j\pi x_1) \sin(j\pi x_2)$ . We compute  $\mathbb{E}[G(u)]$  using QMC with  $R = 8$  random shifts, where  $G(v) = \int_D v(\mathbf{x}) \, d\mathbf{x}$ .

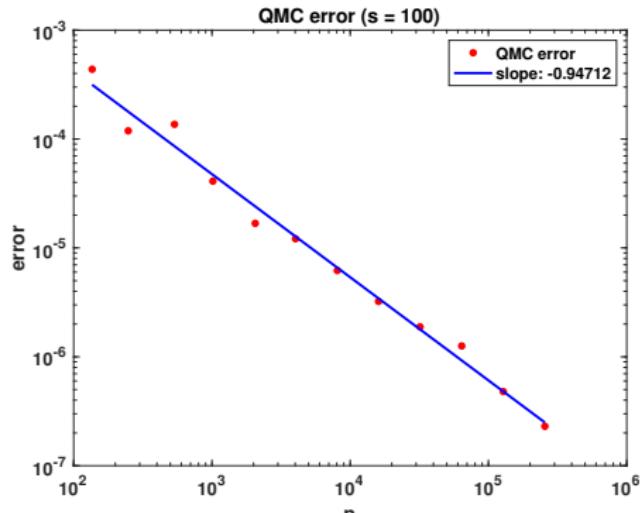


Figure: QMC with  $s = 100$  constructed using the weights

$$\gamma_{\mathfrak{u}} = \left( |\mathfrak{u}|! \prod_{j \in \mathfrak{u}} \frac{b_j}{2(\log 2) \exp(b_j^2/2) \Phi(b_j) \sqrt{(\alpha_j - b_j) \varrho_j(\lambda)}} \right)^{\frac{2}{1+\lambda}}, \quad \lambda = \frac{1}{2-2\delta}, \quad \delta = 0.05,$$
for all  $\mathfrak{u} \subseteq \{1, \dots, s\}$ .

## Computational implementation

Consider the task of approximating  $\int_{[0,1]^s} f(\mathbf{y}) d\mathbf{y}$  using a randomly shifted rank-1 lattice rule with  $R$  random shifts.

Once a generating vector  $\mathbf{z} \in \mathbb{N}^s$  has been obtained for a given number  $n$  of QMC nodes and dimension  $s$  (using, e.g., the CBC algorithm), then:

Remarks:

```
for r = 1, ..., R, do
    draw Δ(r) ~ U([0, 1]s);
    initialize Qr = 0;
    for i = 1, ..., n, do
        set ti = mod(iz/n + Δ(r), 1);
        set Qr = Qr + f(ti);
    end for
    set Qr = Qr/n;
end for
return Q̄ = Q1 + ... + QR/R;
(This is the QMC estimator
with R random shifts.)
```

- If integrating

$$\int_{\mathbb{R}^s} f(\mathbf{y}) \prod_{j=1}^s \frac{e^{-\frac{1}{2}y_j^2}}{\sqrt{2\pi}} d\mathbf{y}$$

then use  $t_i = \Phi^{-1}(\text{mod}(iz/n + \Delta^{(r)}, 1))$ , where  $\Phi^{-1}$  is the (componentwise) inverse cumulative distribution function of  $\mathcal{N}(0, 1)$ .

- The R.M.S. error can be estimated by

R.M.S. error

$$\approx \sqrt{\frac{1}{R(R-1)} \sum_{r=1}^R (\bar{Q} - Q_r)^2}.$$

Some perspectives on applying QMC for Bayesian inverse problems

Let  $U \subseteq \mathbb{R}^s$  be a nonempty set of parameters and let  $G: U \rightarrow \mathbb{R}^k$  be a forward mapping depending on some (unknown) parameter  $\mathbf{y} \in U$ .

Measurement model:

$$\boldsymbol{\delta} = G(\mathbf{y}) + \boldsymbol{\eta},$$

where  $\boldsymbol{\delta} \in \mathbb{R}^k$  is the measurement data and  $\boldsymbol{\eta} \in \mathbb{R}^k$  is Gaussian noise such that  $\boldsymbol{\eta} \sim \mathcal{N}(0, \Gamma)$ , where  $\Gamma \in \mathbb{R}^{k \times k}$  is a symmetric, positive definite covariance matrix.

If we endow  $\mathbf{y}$  with a *prior* density  $\pi_{\text{pr}}$  and  $\mathbf{y}$  and  $\boldsymbol{\eta}$  are independent, then Bayes' formula yields the *posterior distribution* with density

$$\pi(\mathbf{y}|\boldsymbol{\delta}) \propto \pi(\boldsymbol{\delta}|\mathbf{y})\pi_{\text{pr}}(\mathbf{y}),$$

where we have the *likelihood*  $\pi(\boldsymbol{\delta}|\mathbf{y}) \propto e^{-\frac{1}{2}(\boldsymbol{\delta}-G(\mathbf{y}))^T \Gamma^{-1} (\boldsymbol{\delta}-G(\mathbf{y}))}$ .

The *conditional mean (CM)* estimator of the unknown parameter is

$$\mathbf{y}_{\text{CM}} = \int_U \mathbf{y} \pi(\mathbf{y}|\boldsymbol{\delta}) d\mathbf{y} = \frac{\int_U \mathbf{y} e^{-\frac{1}{2}(\boldsymbol{\delta}-G(\mathbf{y}))^T \Gamma^{-1} (\boldsymbol{\delta}-G(\mathbf{y}))} \pi_{\text{pr}}(\mathbf{y}) d\mathbf{y}}{\int_U e^{-\frac{1}{2}(\boldsymbol{\delta}-G(\mathbf{y}))^T \Gamma^{-1} (\boldsymbol{\delta}-G(\mathbf{y}))} \pi_{\text{pr}}(\mathbf{y}) d\mathbf{y}}.$$

For simplicity, let us make the following standing assumptions:

We have  $U = U_s = [-\frac{1}{2}, \frac{1}{2}]^s$  and there is a constant  $C \geq 1$  and a sequence  $\mathbf{b} := (b_j)_{j \geq 1} \in \ell^p$  of nonnegative real numbers for some  $p \in (0, 1)$  such that

(A1) the forward model satisfies

$$\|\partial^\nu G(\mathbf{y})\|_{\mathbb{R}^k} \leq C|\nu|! \mathbf{b}^\nu \quad \text{for all } \nu \in \mathcal{F} \text{ and } \mathbf{y} \in U, \text{ and}$$

(A2)  $\pi_{\text{pr}}(\mathbf{y}) = 1$  for  $\mathbf{y} \in U$  and 0 otherwise.

(A3) The smallest eigenvalue of  $\Gamma$  is bounded from below by  $0 < \mu_{\min} \leq 1$ .

## Example

Let  $D \subset \mathbb{R}^d$ ,  $d \in \{1, 2, 3\}$ , be a nonempty, bounded Lipschitz domain and let  $f \in H^{-1}(D)$ . For each  $\mathbf{y} \in U$ , there exists a weak solution  $u(\cdot, \mathbf{y}) \in H_0^1(D)$  to

$$\begin{cases} -\nabla \cdot (a(\mathbf{x}, \mathbf{y}) \nabla u(\mathbf{x}, \mathbf{y})) = f(\mathbf{x}), & \mathbf{x} \in D, \mathbf{y} \in U, \\ u(\mathbf{x}, \mathbf{y}) = 0, & \mathbf{x} \in \partial D, \mathbf{y} \in U, \end{cases}$$

where we assume that  $\mathbf{y} = (y_j)_{j=1}^s$  are i.i.d. uniformly distributed in  $[-1/2, 1/2]$ , and

$$a(\mathbf{x}, \mathbf{y}) = a_0(\mathbf{x}) + \sum_{j=1}^s y_j \psi_j(\mathbf{x}), \quad \mathbf{x} \in D, \mathbf{y} \in [-1/2, 1/2]^s,$$

with  $a_0 \in L^\infty(D)$  and  $\psi_j \in L^\infty(D)$ ,  $j \geq 1$ , such that

$0 < a_{\min} \leq a(\mathbf{x}, \mathbf{y}) \leq a_{\max} < \infty$  for all  $\mathbf{x} \in D$ ,  $\mathbf{y} \in [-1/2, 1/2]^s$ .

Let  $\mathcal{O}_j: H_0^1(D) \rightarrow \mathbb{R}$  be linear, bounded observation operators for  $j = 1, \dots, k$ . Then  $G(\mathbf{y}) = [\mathcal{O}_j(u(\cdot, \mathbf{y}))]_{j=1}^k$  with

$$\|\partial^\nu G(\cdot, \mathbf{y})\|_{\mathbb{R}^k} \leq C|\nu|! b^\nu, \text{ where } C := \left( \sum_{j=1}^k \|\mathcal{O}_j\|_{H^{-1}(D)}^2 \right)^{1/2} \text{ and} \\ b_j := \frac{\|\psi_j\|_{L^\infty(D)}}{a_{\min}}.$$

We are interested in the CM estimate

$$\mathbf{y}_{\text{CM}} = \frac{1}{Z} \mathbf{Z}',$$

where

$$\begin{aligned}\mathbf{Z}' &= \int_U \mathbf{y} e^{-\frac{1}{2}(\delta - G(\mathbf{y}))^T \Gamma^{-1} (\delta - G(\mathbf{y}))} d\mathbf{y}, \\ Z &= \int_U e^{-\frac{1}{2}(\delta - G(\mathbf{y}))^T \Gamma^{-1} (\delta - G(\mathbf{y}))} d\mathbf{y}.\end{aligned}$$

It can be shown that

$$|\partial^\nu e^{-\frac{1}{2}(\delta - G(\mathbf{y}))^T \Gamma^{-1} (\delta - G(\mathbf{y}))}| \leq 3.82^k \cdot C^{|\nu|} 2^{|\nu|-1} \mu_{\min}^{-|\nu|/2} |\nu|! b^\nu \quad \text{for } \nu \neq \mathbf{0}$$

and

$$\begin{aligned}&\left| \frac{\partial^{|\mathbf{u}|}}{\partial \mathbf{y}_\mathbf{u}} y_\ell e^{-\frac{1}{2}(\delta - G(\mathbf{y}))^T \Gamma^{-1} (\delta - G(\mathbf{y}))} \right| \\ &\leq 3.82^k \cdot C^{|\mathbf{u}|} 2^{|\mathbf{u}|-1} \mu_{\min}^{-|\mathbf{u}|/2} |\mathbf{u}|! \left(1 + \frac{1}{b_s}\right) \prod_{j \in \mathbf{u}} b_j\end{aligned}$$

for  $\emptyset \neq \mathbf{u} \subseteq \{1 : s\}$ .

For the QMC approximation of both the numerator and the denominator, this suggests choosing the weights

$$\gamma_{\mathfrak{u}} = \left( |\mathfrak{u}|! \prod_{j \in \mathfrak{u}} \frac{\beta_j}{\sqrt{\frac{2\zeta(2\lambda)}{(2\pi^2)^\lambda}}} \right)^{2/(1+\lambda)}, \quad \lambda = \begin{cases} \frac{p}{2-p} & \text{if } p \in (2/3, 1), \\ \frac{1}{2-2\delta} & \text{if } p \in (0, 2/3], \end{cases}$$

where  $\delta \in (0, 1/2]$  is arbitrary and

$$\beta_j = 2C\mu_{\min}^{-1/2} b_j, \quad j = 1, \dots, s,$$

as inputs to the CBC algorithm. The QMC rule obtained in this way has

- a **dimension-independent** QMC convergence rate  $\mathcal{O}(n^{\max\{-1/p+1/2, -1+\delta\}})$  for the denominator when the number of QMC nodes  $n$  is a prime power.
- a **dimension-dependent** QMC convergence rate  $\mathcal{O}\left(\left(1 + \frac{1}{b_\ell}\right)n^{\max\{-1/p+1/2, -1+\delta\}}\right)$  for the  $\ell^{\text{th}}$  component of the vector  $\mathbf{Z}'$  when the number of QMC nodes  $n$  is a prime power.  
For example, if  $b_s \propto s^{-2}$ , then  $\left(1 + \frac{1}{b_s}\right) \propto s^2$ .

(Note that the constant in the error bounds also depends on  $k$ , the dimension of the measurement data!)

## Numerical example

Let us consider the PDE problem

$$-\nabla \cdot (a(\mathbf{x}) \nabla u(\mathbf{x})) = x_1, \quad u|_{\partial D} = 0,$$

in the physical domain  $D = (0, 1)^2$ , where the diffusion coefficient is assumed to be *unknown*.

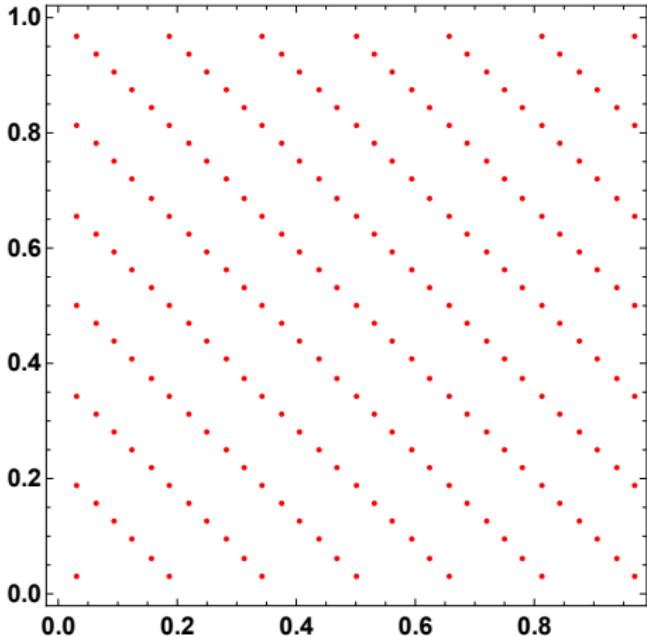
Given some noisy observations  $\mathbf{y} = [u_1, \dots, u_k]^T$  of the PDE solution  $G(u) = [\mathcal{O}_1(u), \dots, \mathcal{O}_k(u)]^T = [u(\mathbf{x}_1), \dots, u(\mathbf{x}_k)]^T$  over a point set  $\mathbf{x}_1, \dots, \mathbf{x}_k \in D$ , we wish to recover the diffusion coefficient which caused the observations.

We model the uncertain diffusion coefficient using the affine parameterization

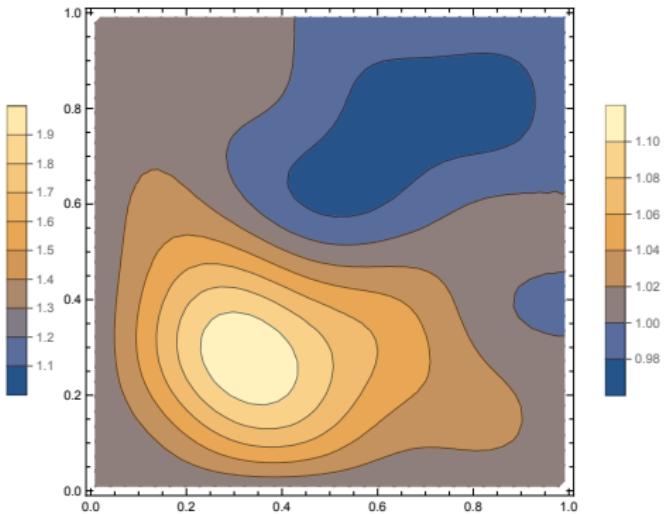
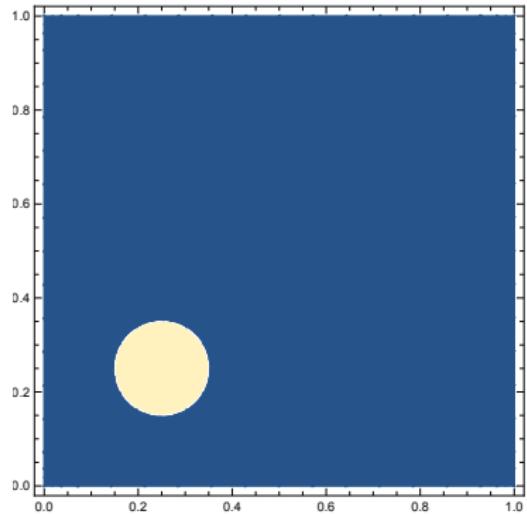
$$a(\mathbf{x}, \mathbf{y}) = 1 + \sum_{j=1}^{30} y_j \psi_j(\mathbf{x}), \quad y_j \in (-\frac{1}{2}, \frac{1}{2}),$$

where  $\psi_j(\mathbf{x}) = \frac{1}{(k_j^2 + \ell_j^2)^{1.1}} \sin(\pi k_j x_1) \sin(\pi \ell_j x_2)$  and the sequence  $(k_j, \ell_j)_{j \geq 1}$  is an ordering of the elements of  $\mathbb{N} \times \mathbb{N}$  so that the sequence  $(\|\psi_j\|_{L^\infty(D)})_{j \geq 1}$  is non-increasing. This implies that  $\|\psi_j\|_{L^\infty(D)} \sim j^{-1.1}$ .

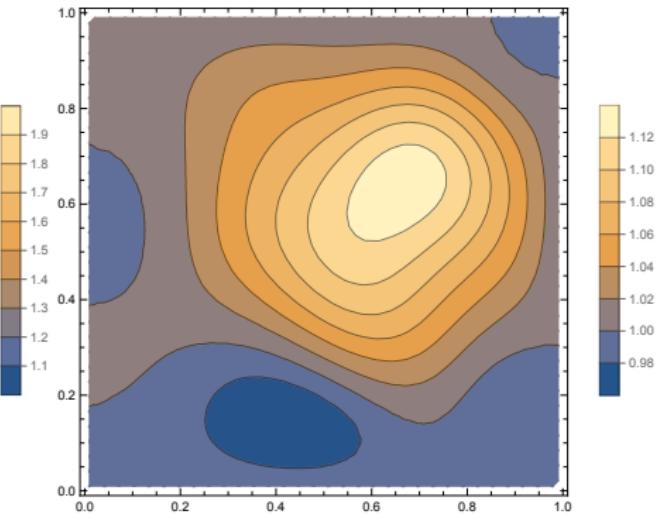
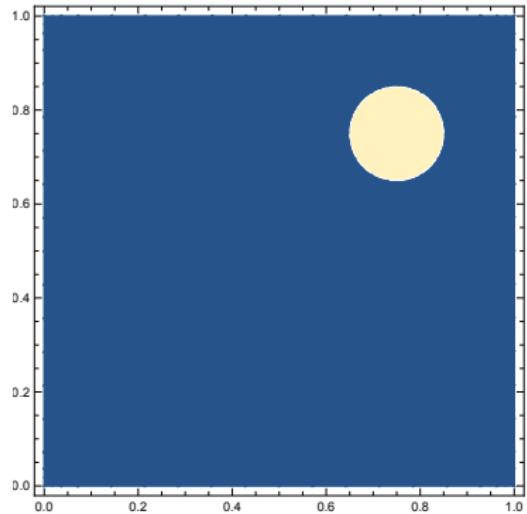
As the reconstruction, we compute the CM estimate  $\mathbf{y}_{\text{CM}} \in [-1/2, 1/2]^{30}$  of the parametric diffusion coefficient which fits the observations and plot  $a(\mathbf{x}, \mathbf{y}_{\text{CM}})$ . The observations were simulated using a FE mesh with  $h = 2^{-7}$  and contaminated with 1% relative white noise. The CM estimate was approximated using QMC with a single random shift and  $n = 2^{14}$  nodes, and the PDE was discretized using a coarser FE mesh  $h = 2^{-5}$  in order to avoid the so-called “inverse crime”.



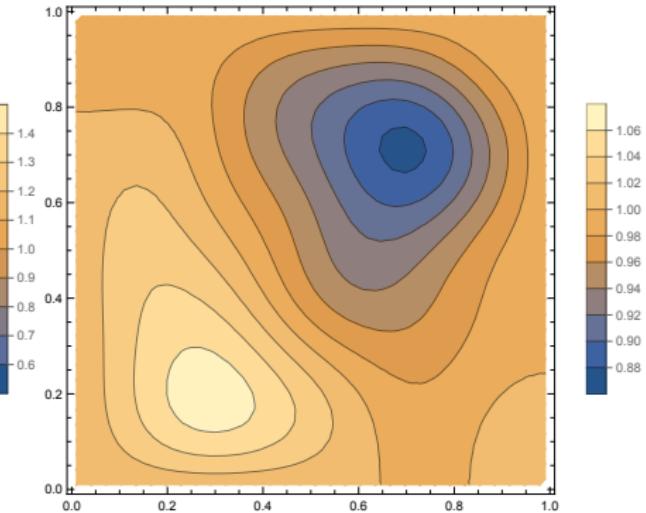
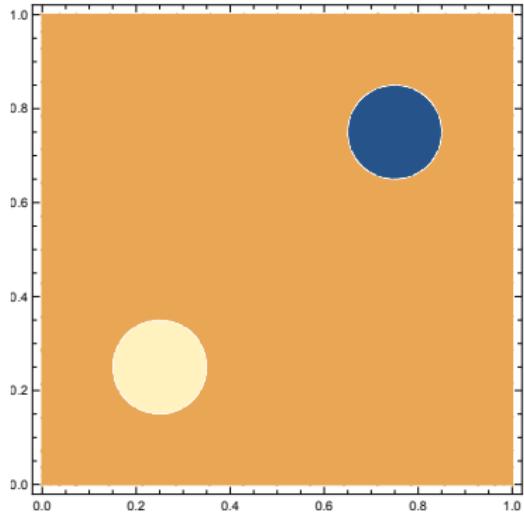
The point set  $\mathbf{x}_1, \dots, \mathbf{x}_k \in D$  in the spatial domain where the noisy measurements  $u_1, \dots, u_k$  were collected with  $k = 193$ .



Left: the ground truth diffusion with  $a_{\text{background}} \equiv 1$  and  $a_{\text{inclusion}} \equiv 2$ .  
Right: the reconstruction  $a(\mathbf{x}, \mathbf{y}_{\text{CM}})$ .



Left: the ground truth diffusion with  $a_{\text{background}} \equiv 1$  and  $a_{\text{inclusion}} \equiv 2$ .  
Right: the reconstruction  $a(\mathbf{x}, \mathbf{y}_{\text{CM}})$ .



Left: the ground truth diffusion with  $a_{\text{background}} \equiv 1$ ,  
 $a_{\text{inclusion, bottom left}} \equiv 1.5$ , and  $a_{\text{inclusion, top right}} \equiv 0.5$ .  
Right: the reconstruction  $a(\mathbf{x}, \mathbf{y}_{\text{CM}})$ .