

Fast Automatic Bayesian Cubature Using Sobol Sampling

R. Jagadeeswaran · Fred J. Hickernell

Received: date / Accepted: date

Abstract Automatic cubatures approximate integrals to user-specified error tolerances. For high dimensional problems, it is difficult to adaptively change the sampling pattern, but one can automatically determine the sample size, n , given a reasonable, fixed sampling pattern. We take this approach here using a Bayesian perspective. We postulate that the integrand is an instance of a Gaussian stochastic process parameterized by a constant mean and a covariance kernel defined by a scale parameter times a parameterized function specifying how the integrand values at two different points in the domain are related. These hyperparameters are inferred or integrated out using integrand values via one of three techniques: empirical Bayes, full Bayes, or generalized cross-validation. The sample size, n , is increased until the half-width of the credible interval for the Bayesian posterior mean is no greater than the error tolerance.

The process outlined above typically requires a computational cost of $O(N_{\text{opt}}n^3)$, where N_{opt} is the number of optimization steps required to identify the hyperparameters. Our innovation is to pair low discrepancy nodes with matching covariance kernels to lower the computational cost to $O(N_{\text{opt}}n \log n)$. This approach is demonstrated explicitly with rank-1 lattice sequences and shift-invariant kernels. Our algorithm is implemented in the Guaranteed Automatic Integration Library (GAIL).

Keywords Bayesian cubature · Fast automatic cubature · GAIL · Probabilistic numeric methods

R. Jagadeeswaran
Department of Applied Mathematics,
Illinois Institute of Technology
10 W. 32nd St., Room 208, Chicago IL 60616
E-mail: jrathin1@iit.edu

Fred J. Hickernell
Center for Interdisciplinary Scientific Computation and
Department of Applied Mathematics
Illinois Institute of Technology
10 W. 32nd St., Room 208, Chicago IL 60616
E-mail: hickernell@iit.edu

1 Introduction

JR: discuss prior work

Cubature is the problem of inferring a numerical value for a definite integral, $\mu := \int_{\mathbb{R}^d} g(\mathbf{x}) d\mathbf{x}$, where μ has no closed form analytic expression. Typically, g is accessible through a black-box function routine. Cubature means numerical multivariate integration and is a key component of many problems in scientific computing, finance Glasserman [2004], statistical modeling, imaging Keller [2013], uncertainty quantification, machine learning Goodfellow et al. [2016], etc.

The integral may often be expressed as

$$\mu := \mu(f) := \mathbb{E}[f(\mathbf{X})] = \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x}, \quad (1)$$

where $f : [0,1]^d \rightarrow \mathbb{R}$ is the integrand, and $\mathbf{X} \sim \mathcal{U}[0,1]^d$. The process of transforming the original integral into the form of (1) is addressed in Beckers and Haegemans [1992], Sidi [2008, 1993], Laurie [1996], Cristea et al. [2007]. The cubature may be an affine function of integrand values:

$$\hat{\mu} := \hat{\mu}(f) := w_0 + \sum_{i=1}^n f(\mathbf{x}_i)w_i, \quad \mathcal{P} := \{\mathbf{x}_i\}_{i=1}^n \subset [0,1]^d \quad (2)$$

where the weights, w_0 , and $\mathbf{w} = (w_i)_{i=1}^n \in \mathbb{R}^n$, and the nodes, \mathcal{P} , are chosen to make the error, $|\mu - \hat{\mu}|$, small. The integration domain $[0,1]^d$ is convenient for the low discrepancy node sets that we use. The nodes are assumed to be deterministic. This research focuses on multivariate numerical integrals where the computational cost is a bottleneck.

We construct a reliable stopping criterion that determines the number of integrand values required, n , to ensure that the error is no greater than a user-defined error tolerance denoted by ε , i.e.,

$$|\mu - \hat{\mu}| \leq \varepsilon. \quad (3)$$

Rather than relying on strong assumptions about the integrand, such as an upper bound on its variance or total variation, we construct a stopping criterion that is based on a credible interval arising from a Bayesian approach to the problem. We build upon the work of Briol et al. Briol et al. [2019], Diaconis Diaconis [1988], O'Hagan O'Hagan [1991], Ritter Ritter [2000], Rasmussen and Ghahramani Rasmussen and Williams [2003], and others. Our algorithm is an example of *probabilistic numerics*. To study numerical algorithms from a statistical point of view, where uncertainty is formally due to the presence of an unknown numerical error, is the goal of probabilistic numerics.

Our primary contribution in this research is to demonstrate how the choice of a family of covariance kernels that match the low discrepancy sampling nodes facilitates fast computation of the cubature and the data-driven stopping criterion. Our Bayesian cubature requires a computational cost of

$$\mathcal{O}(n\$(f) + N_{\text{opt}}[n\$(C) + n \log(n)]), \quad (4)$$

where $\$(f)$ is the cost of one integrand value, $\$(C)$ is the cost of a single covariance kernel value, $\mathcal{O}(n \log(n))$ is the cost of a fast Bayesian transform, and N_{opt} is an upper bound on the number of optimization steps required to choose the hyperparameters. If function evaluation is expensive, e.g., the output of a computationally intensive simulation, or if $\$(f) = \mathcal{O}(d)$ for large d , then $\$(f)$ might be similar

in magnitude to $N_{\text{opt}} \log(n)$ in practice. Typically, $\$(C) = \mathcal{O}(d)$. Note that the $\mathcal{O}(n \log(n))$ contribution is d independent.

In contrast to our fast algorithm, the typical computational cost for Bayesian cubature is

$$\mathcal{O}(n\$(f) + N_{\text{opt}}[n^2\$(C) + n^3]), \quad (5)$$

which is explained in Jagadeeswaran and Hickernell [2019]. Note that apart from evaluating the integrand, the computational cost in (5) is much larger than that in (4).

Low discrepancy points are characterized by how uniformly the points are distributed, which is measured by the *discrepancy*. The goal is to have maximum uniform space filling. The discrepancy is defined as below. Let \mathcal{M} be the set of all intervals of the form $\prod_{\ell=1}^d [a_\ell, b_\ell] = \{\mathbf{x} \in \mathbb{R}^d : a_\ell \leq x_\ell \leq b_\ell, 0 \leq a_\ell \leq b_\ell \leq 1\}$. Then, the discrepancy of a point set \mathcal{P} is,

$$D(\mathcal{P}) := \sup_{M \in \mathcal{M}} \left| \frac{|M \cap \mathcal{P}|}{|\mathcal{P}|} - \lambda_L(M) \right|,$$

where $|\mathcal{P}|$ is the cardinality of the set \mathcal{P} , and λ_L is the Lebesgue measure. The *low discrepancy points* satisfy $D(\mathcal{P}) = \mathcal{O}((\log n)^d/n)$. In this work we experiment with two most popular low discrepancy point sets, 1) lattice points, and 2) Sobol' points.

Hickernell [2018] compares different approaches to cubature error analysis depending on whether the rule is deterministic or random and whether the integrand is assumed to be deterministic or random. Error analysis that assumes a deterministic integrand lying in a Banach space leads to an error bound that is typically impractical for deciding how large n must be to satisfy (3). The deterministic error bound includes a (semi-)norm of the integrand, which is often more complex to compute than the original integral.

Hickernell and Jiménez-Rugama [2016], Hickernell and Jiménez Rugama [2016], Jiménez Rugama and Hickernell [2016] have developed stopping criteria for cubature rules based on low discrepancy nodes by tracking the decay of the discrete Fourier coefficients of the integrand. The algorithms proposed here also rely on discrete Fourier coefficients, but in a different way. We only discuss automatic Bayesian cubature for absolute error tolerances in this thesis. The recent work by Hickernell, Jiménez-Rugama, and Li [2018] suggests how one might accommodate more general error criteria, such as relative error tolerances which has been adapted in the MATLAB implementation of our algorithms.

Section 2 explains the Bayesian approach to calculate the posterior cubature error and defines our automatic Bayesian cubature. Although much of this material is known, it is included for completeness. We end Chapter 2 by demonstrating why Bayesian cubature is typically computationally expensive. Chapter 6 introduces the concept of covariance kernels that match the nodes and expedite the computations required by our automatic Bayesian cubature. Chapter ?? implements this concept for shift invariant kernels and rank-1 lattice nodes. It also develops approaches to build shift-invariant kernels of continuous valued kernel order rather than fixing the kernel order to integer values. Chapter 7 demonstrates another implementation of matching nodes and kernel using Sobol' points and Walsh kernels. It also shows that the fast Walsh Hadamard as the fast Bayesian transform for this case. Chapter ?? describes how to avoid cancellation error for kernels of product form. It also covers some of the additional techniques used in the implementation

of our Bayesian Cubature algorithms. Numerical examples are provided in Chapter 14 to demonstrate the performance and advantages of our new algorithms. We conclude with a brief discussion and potential future work in Chapter ??.

We use the terms integrand or function interchangeably to denote the function f being considered for the numerical integration. Also, we use the terms, nodes, points, node-sets, designs, and data-sites interchangeably to denote the points \mathcal{P} used in the cubature.

2 Bayesian Cubature

JR: Briefly explain and cite our BayesLattice paper

1. automatic Bayesian cubature
2. fast automatic Bayesian cubature

The Bayesian approach for numerical analysis was popularized by Diaconis [1988]. The earliest reference for such kind of approach dates back to Poincaré, where, the theory of interpolation was discussed. Diaconis motivates the reader by interpreting the most well known numerical methods, 1) trapezoidal rule and 2) splines, from the statistical point of view with whatever is known about the integrand as prior information. For example, the trapezoidal rule can be interpreted as a Bayesian method with prior information being modeled as a Brownian motion in the sample space $\mathcal{C}[0, 1]$, the space of continuous functions.

This research is focused on the Bayesian approach for numerical integration that is known as Bayesian cubature as introduced by O'Hagan [1991]. Bayesian cubature returns a probability distribution, that expresses belief about the true value of integral, $\mu(f)$. This posterior probability distribution is based on a prior that depends on f , which is computed via Bayes' rule using the *data* contained in the function evaluations Briol et al. [2019]. The distribution in general captures numerical uncertainty due to the fact that we have only used a finite number of function values to evaluate the integral.

2.1 Bayesian Posterior Error

We assume the integrand, f , is an instance of a stochastic Gaussian process, i.e., $f \sim \mathcal{GP}(m, s^2 C_\theta)$. Specifically, f is a real-valued random function with constant mean m and covariance function $s^2 C_\theta$, where s is a positive scale factor, and $C_\theta : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$ is a symmetric, positive-definite function and, parameterized by θ :

$$C^T = C, \quad \mathbf{a}^T C \mathbf{a} > 0, \quad \text{where } C = (C_\theta(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n, \\ \text{for all } \mathbf{a} \neq 0, n \in \mathbb{N}, \text{ distinct } \mathbf{x}_1, \dots, \mathbf{x}_n \in [0, 1]^d. \quad (6)$$

The covariance function, C , and the Gram matrix, C , depend implicitly on θ , but the notation may omit this for simplicity's sake. Procedures for estimating or integrating out the hyperparameters m , s , and θ are explained later in this section.

For a Gaussian process, all vectors of linear functionals of f have a multivariate Gaussian distribution. For any deterministic sampling scheme with distinct nodes,

$\{\mathbf{x}_i\}_{i=1}^n$, and defining $\mathbf{f} := (f(\mathbf{x}_i))_{i=1}^n$ as the multivariate Gaussian vector of function values, it follows from the definition of a Gaussian process that

$$\mathbf{f} \sim \mathcal{N}(m\mathbf{1}, s^2\mathbf{C}), \quad (7a)$$

$$\mu \sim \mathcal{N}(m, s^2c_0), \quad (7b)$$

$$\text{where } c_0 := \int_{[0,1]^d \times [0,1]^d} C_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{t}) \, d\mathbf{x} \, d\mathbf{t}, \quad (7c)$$

$$\text{cov}(\mathbf{f}, \mu) = \left(\int_{[0,1]^d} C(\mathbf{t}, \mathbf{x}_i) \, d\mathbf{t} \right)_{i=1}^n =: \mathbf{c}. \quad (7d)$$

Here, c_0 and \mathbf{c} depend implicitly on $\boldsymbol{\theta}$. We assume the covariance function C is simple enough that the integrals in these definitions can be computed analytically. We need the following lemma to derive the posterior error of our cubature.

Lemma 1 [Rasmussen and Williams, 2006, (A.6), (A.11–13)] *If $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2)^T \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$, where \mathbf{Y}_1 and \mathbf{Y}_2 are random vectors of arbitrary length, and*

$$\mathbf{m} = \begin{pmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \end{pmatrix} = \begin{pmatrix} \mathbb{E}(\mathbf{Y}_1) \\ \mathbb{E}(\mathbf{Y}_2) \end{pmatrix},$$

$$\mathbf{C} = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{C}_{21}^T \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{pmatrix} = \begin{pmatrix} \text{var}(\mathbf{Y}_1) & \text{cov}(\mathbf{Y}_1, \mathbf{Y}_2) \\ \text{cov}(\mathbf{Y}_2, \mathbf{Y}_1) & \text{var}(\mathbf{Y}_2) \end{pmatrix}$$

then

$$\mathbf{Y}_1 | \mathbf{Y}_2 \sim \mathcal{N}(\mathbf{m}_1 + \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} (\mathbf{Y}_2 - \mathbf{m}_2), \quad \mathbf{C}_{11} - \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{21}).$$

Moreover, the inverse of the matrix \mathbf{C} may be partitioned as

$$\mathbf{C}^{-1} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{21}^T \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix},$$

$$\mathbf{A}_{11} = (\mathbf{C}_{11} - \mathbf{C}_{12} \mathbf{C}_{22}^{-1} \mathbf{C}_{21})^{-1}, \quad \mathbf{A}_{21} = -\mathbf{C}_{22}^{-1} \mathbf{C}_{21} \mathbf{A}_{11},$$

$$\mathbf{A}_{22} = \mathbf{C}_{22}^{-1} + \mathbf{C}_{22}^{-1} \mathbf{C}_{21} \mathbf{A}_{11} \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1}.$$

It follows from Lemma 1 that the *conditional* distribution of the integral given observed function values, $\mathbf{f} = \mathbf{y}$ is also Gaussian:

$$\mu | (\mathbf{f} = \mathbf{y}) \sim \mathcal{N}(m(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}, \quad s^2(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})). \quad (8)$$

The natural choice for the cubature is the posterior mean of the integral, namely,

$$\hat{\mu} | (\mathbf{f} = \mathbf{y}) = m(1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}) + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}, \quad (9)$$

which takes the form of (2). Under this definition, the cubature error has zero mean and a variance depending on the choice of nodes:

$$(\mu - \hat{\mu}) | (\mathbf{f} = \mathbf{y}) \sim \mathcal{N}\left(0, \quad s^2(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})\right).$$

A credible interval for the integral is given by

$$\mathbb{P}_f[|\mu - \hat{\mu}| \leq \text{err}_{\text{CI}}] = 99\%, \quad (10a)$$

$$\text{err}_{\text{CI}} = 2.58s\sqrt{c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}}. \quad (10b)$$

Naturally, 2.58 and 99% can be replaced by other quantiles and credible levels.

3 Hyperparameter search

JR: This section explains hyperparameter estimation using gradient descent

The credible interval in (10) suggests how our automatic Bayesian cubature proceeds. Integrand data is accumulated until the width of the credible interval,

err_{CI} , is no greater than the error tolerance. As n increases, one expects $c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}$ to decrease for well-chosen nodes, $\{\mathbf{x}_i\}_{i=1}^n$. Please note that the credible interval depends on the parameters m, s , and $\boldsymbol{\theta}$

Note that err_{CI} shows no explicit dependence on the integrand values, even though one would intuitively expect that a larger integrand should imply a larger err_{CI} . This dependence will become explicit once the hyperparameters, m, s , and $\boldsymbol{\theta}$, are inferred from integrand data. After inferring the hyperparameters, err_{CI} *does reflect the size* of the integrand values. Three approaches for hyperparameter estimation were described in Jagadeeswaran and Hickernell [2019] but no explicit parameter search was discussed. Here we provide one such technique using steepest descent.

The methods developed for hyperparameter estimation from the previous sections are summarized as a theorem below:

Theorem 1 *There are at least three approaches to estimating or integrating out the hyperparameters defining the Gaussian process from which the integrand is drawn: empirical Bayes, full Bayes, and generalized cross-validation. Under these three approaches, we have the following:*

$$m_{\text{EB}} = \frac{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}}, \quad m_{\text{GCV}} = \frac{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{y}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1}}, \quad (11)$$

$$s_{\text{EB}}^2 = \frac{1}{n} \mathbf{y}^T \left[\mathbf{C}_{\boldsymbol{\theta}}^{-1} - \frac{\mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}} \right] \mathbf{y}, \quad (12)$$

$$\begin{aligned} \hat{\sigma}_{\text{full}}^2 &= \frac{1}{n-1} \mathbf{y}^T \left[\mathbf{C}_{\boldsymbol{\theta}}^{-1} - \frac{\mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}} \right] \mathbf{y} \\ &\times \left[\frac{(1 - \mathbf{c}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{c}) \right], \end{aligned} \quad (13)$$

$$s_{\text{GCV}}^2 = \mathbf{y}^T \left[\mathbf{C}_{\boldsymbol{\theta}}^{-2} - \frac{\mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1} \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1}} \right] \mathbf{y} \left[\text{trace}(\mathbf{C}_{\boldsymbol{\theta}}^{-1}) \right]^{-1},$$

$$\boldsymbol{\theta}_{\text{EB}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \left\{ \log \left(\mathbf{y}^T \left[\mathbf{C}_{\boldsymbol{\theta}}^{-1} - \frac{\mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}} \right] \mathbf{y} \right) + \frac{1}{n} \log(\det(\mathbf{C}_{\boldsymbol{\theta}})) \right\}, \quad (14)$$

$$\boldsymbol{\theta}_{\text{GCV}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \left\{ \log \left(\mathbf{y}^T \left[\mathbf{C}_{\boldsymbol{\theta}}^{-2} - \frac{\mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1} \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1}} \right] \mathbf{y} \right) - \log \left(\text{trace}(\mathbf{C}_{\boldsymbol{\theta}}^{-2}) \right) \right\}, \quad (15)$$

$$\hat{\mu}_{\text{EB}} = \hat{\mu}_{\text{full}} = \left(\frac{(1 - \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{c}) \mathbf{1}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}} + \mathbf{c} \right)^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{y}, \quad (16)$$

$$\hat{\mu}_{\text{GCV}} = \left(\frac{(1 - \mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{c}) \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{1}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1}} + \mathbf{c} \right)^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{y}. \quad (17)$$

The credible intervals widths, err_{CI} , are given by

$$\text{err}_{\mathbf{x}} = 2.58 s_{\mathbf{x}} \sqrt{c_0 - \mathbf{c}^T \mathbf{C}_{\boldsymbol{\theta}}^{-1} \mathbf{c}}, \quad \mathbf{x} \in \{\text{EB}, \text{GCV}\}, \quad (18)$$

$$\text{err}_{\text{full}} = t_{n-1, 0.995} \hat{\sigma}_{\text{full}} > \text{err}_{\text{EB}}. \quad (19)$$

The resulting credible intervals are then

$$\mathbb{P}_f [|\mu - \hat{\mu}_{\mathbf{x}}| \leq \text{err}_{\mathbf{x}}] = 99\%, \quad \mathbf{x} \in \{\text{EB}, \text{full}, \text{GCV}\}. \quad (20)$$

Here $t_{n-1,0.995}$ denotes the 99.5 percentile of a standard Student's t -distribution with $n-1$ degrees of freedom. In the formulas above, $\boldsymbol{\theta}$ is assumed to take on the values $\boldsymbol{\theta}_{\text{EB}}$ or $\boldsymbol{\theta}_{\text{GCV}}$ as appropriate.

In the theorem above, note that if the original covariance kernel, C , is replaced by bC for some positive constant b , the cubature, $\hat{\mu}$, the estimates of $\boldsymbol{\theta}$, and the credible interval half-widths, err_x for $x \in \{\text{EB}, \text{full}, \text{GCV}\}$, all remain unchanged. The estimates of s^2 are multiplied by b^{-1} , as would be expected.

3.1 Gradient descent to find optimal shape parameter

The equation specifying $\boldsymbol{\theta}_{\text{EB}}$ as defined in (14) does not say how the parameter search can be done. There exist empirical algorithms Brent [1973], Forsythe et al. [1976] that one could use to accomplish the same. Since the objective function is known we could compute the gradient. Using the gradient of $l(s, m, \boldsymbol{\theta}|\mathbf{y})$, one can apply optimization techniques such as gradient descent to find the optimal value faster. Let us define the objective function for the same purpose by excluding the negative sign, which modifies the problem to become a minimization of

$$\mathcal{L}(\boldsymbol{\theta}|\mathbf{y}) := \frac{1}{n} \log(\det C) + \log \left((\mathbf{y} - m_{\text{EB}}\mathbf{1})^T C^{-1} (\mathbf{y} - m_{\text{EB}}\mathbf{1}) \right) + \text{constants}.$$

Taking derivative with respect to θ_ℓ , for $\ell = 1, \dots, d$

$$\begin{aligned} \frac{\partial}{\partial \theta_\ell} \mathcal{L}(\boldsymbol{\theta}|\mathbf{y}) &= \frac{1}{n} \frac{\partial}{\partial \theta_\ell} \log(\det C) + \frac{\partial}{\partial \theta_\ell} \log \left((\mathbf{y} - m_{\text{EB}}\mathbf{1})^T C^{-1} (\mathbf{y} - m_{\text{EB}}\mathbf{1}) \right) \\ &= \frac{1}{n} \text{trace} \left(C^{-1} \frac{\partial C}{\partial \theta_\ell} \right) - \frac{((\mathbf{y} - m_{\text{EB}}\mathbf{1})^T C^{-1})^T \left(\frac{\partial C}{\partial \theta_\ell} \right) ((\mathbf{y} - m_{\text{EB}}\mathbf{1})^T C^{-1})}{(\mathbf{y} - m_{\text{EB}}\mathbf{1})^T C^{-1} (\mathbf{y} - m_{\text{EB}}\mathbf{1})} \end{aligned}$$

where we used some of the results from Dong et al. [2017]. This can be used with gradient descent as follows,

$$\theta_\ell^{(j+1)} = \theta_\ell^{(j)} - \nu_\ell \frac{\partial}{\partial \theta_\ell} \mathcal{L}(\boldsymbol{\theta}|\mathbf{y}), \quad j = 0, 1, \dots \quad (21)$$

where ν_ℓ is the step size for the gradient descent.

4 Cone of Functions and the Credible interval

In this research we assume that the integrand belongs to a cone of well-behaved functions, \mathcal{C} , to make the computations bounded in terms of function data. The concept of cone in general for cubature error analysis can be stated using the error bound definition. Suppose that

$$|\mu(f) - \hat{\mu}_n(f)| \leq \text{err}_{\text{CI}}(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) \quad (22)$$

for some f , which it is 99% of the time under our hypothesis. Also note that our err_{CI} (18) (19) are positively homogeneous functions, meaning,

$$\text{err}_{\text{CI}}(ay_1, \dots, ay_n) = |a| \text{err}_{\text{CI}}(y_1, \dots, y_n).$$

One can verify the homogeneity of (18) and (19) easily. Thus if f satisfies (22), then

$$\begin{aligned} |\mu(af) - \hat{\mu}_n(af)| &= |a| |\mu(f) - \hat{\mu}_n(f)| \\ &\leq |a| \text{err}_{\text{CI}}(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) \\ &= \text{err}_{\text{CI}}(af(\mathbf{x}_1), \dots, af(\mathbf{x}_n)) \end{aligned}$$

for all real a . Thus the set of all f satisfying (22) is a *cone*, \mathcal{C} . Cones of functions satisfy the property that if $f \in \mathcal{C}$ then $af \in \mathcal{C}$.

In the context of Bayesian cubature, one can explain the cone concept beginning with the definition of credible interval (10). Let $f \sim \mathcal{GP}$, be an instance of a Gaussian stochastic process:

$$\mathbb{P}_f [|\mu(f) - \hat{\mu}_n(f)| \leq \text{err}_{\text{CI}}(f)] \geq 99\%.$$

This can be interpreted as $|\mu(f) - \hat{\mu}_n(f)| \leq \text{err}_{\text{CI}}(f)$ with 99% confidence. If f is in the 99% middle of the sample space with $f(\mathbf{x}_i) = y_i$ then af is also in the middle 99% of the sample space with $af(\mathbf{x}_i) = ay_i$.

We demonstrate the credible interval using the following example. For this purpose, choose a smooth and periodic integrand $f_{\text{smooth}}(\mathbf{x}) = \exp(\sum_{\ell=1}^d \cos(2\pi x_\ell))$ and another integrand $f_{\text{peaky}}(\mathbf{x}) = f_{\text{smooth}} + a_{\text{peaky}} f_{\text{noise}}$ where $a_{\text{peaky}} \in \mathbb{R}$. Here $f_{\text{noise}}(\mathbf{x}) = (1 - \exp(2\pi\sqrt{-1}\mathbf{x}^T\boldsymbol{\zeta}))$, $\boldsymbol{\zeta} \in \mathbb{R}^d$ is some d -dimensional vector belonging to the dual space of the lattice nodes for some $\{\mathbf{x}_i\}_{i=1}^n$. The $\boldsymbol{\zeta}$ in the dual space of lattice nodes implies that $f_{\text{noise}}(\mathbf{x}_i) = 0$ at the sampling nodes $\{\mathbf{x}_i\}_{i=1}^n$. The f_{nice} is obtained by kernel interpolation of the n samples of f_{smooth} at $\{\mathbf{x}_i\}_{i=1}^n$. We chose the Matérn kernel (??) for the interpolation. Please note that $f_{\text{peaky}}(\mathbf{x}_i) = f_{\text{nice}}(\mathbf{x}_i) = f_{\text{smooth}}(\mathbf{x}_i)$ for $i = 1, \dots, n$.

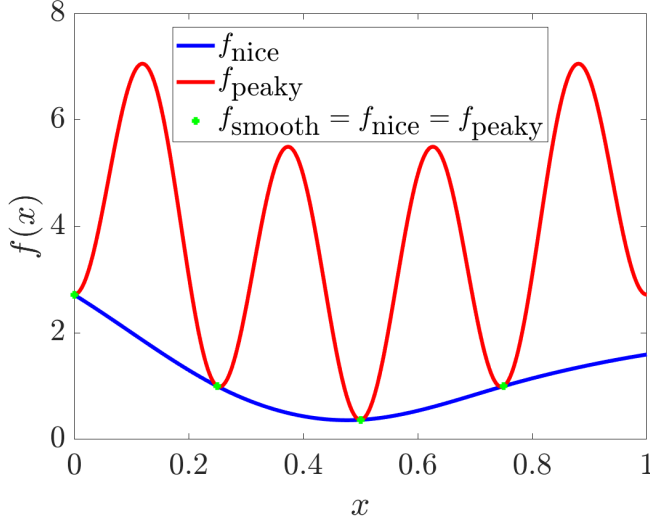


Fig. 1: Example integrands 1) f_{nice} , a smooth function, 2) f_{peaky} , a peaky function. The function values $f_{\text{peaky}}(\mathbf{x}_i) = f_{\text{nice}}(\mathbf{x}_i) = f_{\text{smooth}}(\mathbf{x}_i)$ for $i = 1, \dots, n$. This plot can be conditionally reproduced using `DemoCone.m`

In Figure 1, the sampled function values are shown as dots. One can imagine these samples were obtained from f_{nice} , a moderately smoother function or from f_{peaky} , a highly oscillating function. In this example, we used $a_{\text{peaky}} = 2$.

When using $n = 16$ rank-1 lattice points, and $r = 1$ shift-invariant kernel, we get the posterior distribution of μ as shown in Figure 2. The true integral

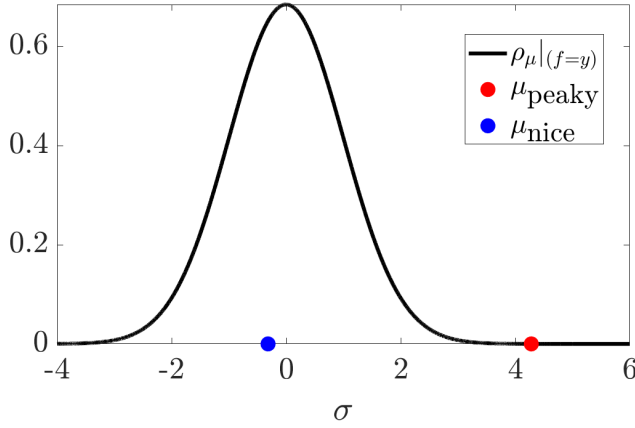


Fig. 2: Probability distributions showing the relative integral position of a smooth and a peaky function. f_{nice} lies within the center 99% of the confidence interval, and f_{peaky} lies on the outside of 99% of the confidence interval. This plot can be conditionally reproduced using `DemoCone.m`

value is shown as μ_{smooth} which is at the center of the plot. The integral of the peaky function f_{peaky} lies outside of the 99% of the credible interval given by (??), whereas the μ_{nice} falls within.

Our Bayesian cubature algorithms compute the approximate integral using only the samples of the integrand. Estimated integral value of our algorithm closely matches the integral of a smooth function that falls within the middle of the confidence interval. If the true integrand were to resemble the smooth approximate function then the estimated integral will be accurate.

5 The Automatic Bayesian Cubature Algorithm

The previous section presents three credible intervals, (??), (??), and (??), for the μ , the desired integral. Each credible interval is based on different assumptions about the hyperparameters m , s , and θ . We stress that one must estimate these hyperparameters or assume a prior distribution on them because the credible intervals are used as stopping criteria for our cubature rule. Since a credible interval makes a statement about a typical function—not an outlier—one must try to ensure that the integrand is a typical draw from the assumed Gaussian process.

Our Bayesian cubature algorithm increases the sample size until the width of the credible interval is small enough. This is accomplished through successively doubling the sample size. The steps are detailed in Algorithm 1.

We recognize that multiple applications of our credible intervals in one run of the algorithm is not strictly justified. However, if our integrand comes from the middle of the sample space and not the extremes, we expect our automatic Bayesian cubature to approximate the integral within the desired error tolerance with high probability. The example in the next section and the examples in Sec-

tion 14 support that expectation. We also believe that an important factor contributing to the occasional failure of our algorithm is unreasonable parameterizations of the stochastic process from which the integrand is hypothesized to be drawn. Overcoming this latter challenge is a topic for future research.

Algorithm 1 Automatic Bayesian Cubature

Require: a generator for the sequence $\mathbf{x}_1, \mathbf{x}_2, \dots$; a black-box function, f ; an absolute error tolerance, $\varepsilon > 0$; the positive initial sample size, n_0 ; the maximum sample size n_{\max}

- 1: $n \leftarrow n_0, n' \leftarrow 0, \text{err} \leftarrow \infty$
- 2: **while** $\text{err} > \varepsilon$ and $n \leq n_{\max}$ **do**
- 3: Generate $\{\mathbf{x}_i\}_{i=n'+1}^n$ and sample $\{f(\mathbf{x}_i)\}_{i=n'+1}^n$
- 4: Compute $\boldsymbol{\theta}$ by (14) or (15)
- 5: Compute err according to (18), (19), or (??)
- 6: $n' \leftarrow n, n \leftarrow 2n'$
- 7: **end while**
- 8: Sample size to compute $\hat{\mu}, n \leftarrow n'$
- 9: Compute $\hat{\mu}$, the approximate integral, according to (16) or (17)
- 10: **return** $\hat{\mu}, n$ and err

As described above, the computational cost of Algorithm 1 is the sum of the following:

- $\mathcal{O}(n\$(f))$ for the integrand data, where $\$(f)$ is the computational cost of a single $f(\mathbf{x})$; $\$(f)$ may be large if it is the result of an expensive simulation; $\$(f)$ is typically proportional to d ;
- $\mathcal{O}(N_{\text{opt}}n^2\$(C_{\boldsymbol{\theta}}))$ for the evaluation of the Gram matrix $\mathbf{C}_{\boldsymbol{\theta}}$, N_{opt} is the number of optimization steps required, and $\$(C_{\boldsymbol{\theta}})$ is the computational cost of a single $C_{\boldsymbol{\theta}}(\mathbf{t}, \mathbf{x})$; $\$(C_{\boldsymbol{\theta}})$ is typically proportional to d ; and
- $\mathcal{O}(N_{\text{opt}}n^3)$ for the matrix inversions and determinant calculations; this cost is independent of d .

As we see in the example in the next section, the cost increases quickly as the n required to meet the error tolerance increases. This motivates the fast Bayesian cubature algorithm presented in Section 6.

6 Fast Bayesian Transform Kernel

We introduce some further notation

$$\begin{aligned}
\mathbf{C} &= \mathbf{C}_{\boldsymbol{\theta}} = \left(C_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{x}_j) \right)_{i,j=1}^n = (\mathbf{C}_1, \dots, \mathbf{C}_n) \\
&= \frac{1}{n} \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^H, \quad \mathbf{V}^H = n \mathbf{V}^{-1}, \\
\mathbf{V} &= (\mathbf{v}_1, \dots, \mathbf{v}_n)^T = (\mathbf{V}_1, \dots, \mathbf{V}_n) \\
\mathbf{C}^p &= \frac{1}{n} \mathbf{V} \boldsymbol{\Lambda}^p \mathbf{V}^H, \quad \forall p \in \mathbb{Z},
\end{aligned} \tag{23}$$

where \mathbf{V}^H is the Hermitian of \mathbf{V} , $\mathbf{C}_1, \dots, \mathbf{C}_n$ are columns of \mathbf{C} , $\mathbf{V}_1, \dots, \mathbf{V}_n$ are columns of \mathbf{V} , and $\mathbf{v}_1, \dots, \mathbf{v}_n$ are rows of \mathbf{V} . The columns of matrix \mathbf{V} are eigenvectors of \mathbf{C} , and $\boldsymbol{\Lambda}$ is a diagonal matrix of eigenvalues of \mathbf{C} . In this and later sections, we drop the $\boldsymbol{\theta}$ dependence of various quantities for simplicity of notation.

The normalization of \mathbf{V} assumed in (23) conveniently allows the first eigenvector, \mathbf{V}_1 , to be the vector of ones in (24b) below. For any $n \times 1$ vector \mathbf{b} , define the notation $\tilde{\mathbf{b}} := \mathbf{V}^H \mathbf{b}$.

We make three assumptions that allow the fast computation:

$$\mathbf{V} \text{ may be identified analytically,} \quad (24a)$$

$$\mathbf{v}_1 = \mathbf{V}_1 = \mathbf{1}, \quad (24b)$$

$$\text{Computing } \mathbf{V}^H \mathbf{b} \text{ requires only } \mathcal{O}(n \log n) \text{ operations } \forall \mathbf{b}. \quad (24c)$$

We call the transformation $\mathbf{b} \mapsto \mathbf{V}^H \mathbf{b}$ a *fast Bayesian transform* and C_θ a *fast Bayesian transform kernel* for the matching nodes $\{\mathbf{x}_i\}_{i=1}^\infty$.

The covariance kernel used in practice also may satisfy an additional assumption:

$$\int_{[0,1]^d} C(\mathbf{t}, \mathbf{x}) d\mathbf{t} = 1 \quad \forall \mathbf{x} \in [0,1]^d, \quad (25)$$

which implies that $c_{0\theta} = 1$ and $\mathbf{c}_\theta = \mathbf{1}$. Under (25), the expressions above may be further simplified:

The assumptions and results in this chapter lead to the following theorem.

Theorem 2 *Under assumptions (24), the parameters and credible interval half-widths in Theorem 1 may be expressed in terms of the fast Bayesian transforms of the integrand data, the first column of the Gram matrix, c_0 , and \mathbf{c} as follows:*

$$m_{\text{EB}} = m_{\text{full}} = m_{\text{GCV}} = \frac{\tilde{y}_1}{n} = \frac{1}{n} \sum_{i=1}^n y_i,$$

$$s_{\text{EB}}^2 = \frac{1}{n^2} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i},$$

$$\hat{\sigma}_{\text{full}}^2 = \frac{1}{n(n-1)} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left[\frac{\lambda_1}{n} \left(1 - \frac{\tilde{c}_1}{\lambda_1} \right)^2 + \left(c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\tilde{c}_i|^2}{\lambda_i} \right) \right],$$

$$s_{\text{GCV}}^2 = \frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \left[\sum_{i=1}^n \frac{1}{\lambda_i} \right]^{-1},$$

$$\boldsymbol{\theta}_{\text{EB}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\log \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right) + \frac{1}{n} \sum_{i=1}^n \log(\lambda_i) \right], \quad (26a)$$

$$\boldsymbol{\theta}_{\text{GCV}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\log \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right) \right], \quad (26b)$$

$$\hat{\mu}_{\text{EB}} = \hat{\mu}_{\text{full}} = \hat{\mu}_{\text{GCV}} = \frac{\tilde{y}_1}{n} + \frac{1}{n} \sum_{i=2}^n \frac{\tilde{c}_i^* \tilde{y}_i}{\lambda_i},$$

$$\text{err}_{\text{EB}} = \frac{2.58}{n} \sqrt{\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left(c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\tilde{c}_i|^2}{\lambda_i} \right)},$$

$$\text{err}_{\text{full}} = t_{n-1, 0.995} \hat{\sigma}_{\text{full}},$$

$$\text{err}_{\text{GCV}} = \frac{2.58}{n} \left\{ \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i} \right]^{-1} \left(c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\tilde{c}_i|^2}{\lambda_i} \right) \right\}^{1/2}.$$

Under the further assumption (25), it follows that

$$\hat{\mu}_{\text{EB}} = \hat{\mu}_{\text{full}} = \hat{\mu}_{\text{GCV}} = \frac{\tilde{y}_1}{n} = \frac{1}{n} \sum_{i=1}^n y_i, \quad (27)$$

and so $\hat{\mu}$ is simply the sample mean. Also, under assumption (25), the credible interval half-widths simplify to

$$\text{err}_{\text{EB}} = \frac{2.58}{n} \sqrt{\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left(1 - \frac{n}{\lambda_1}\right)}, \quad (28a)$$

$$\text{err}_{\text{full}} = t_{n-1, 0.995} \sqrt{\frac{1}{n(n-1)} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left(\frac{\lambda_1}{n} - 1\right)}, \quad (28b)$$

$$\text{err}_{\text{GCV}} = \frac{2.58}{n} \left\{ \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i} \right]^{-1} \left(1 - \frac{n}{\lambda_1}\right) \right\}^{1/2}. \quad (28c)$$

In the formulas for the credible interval half-widths and λ depends on θ , and θ is assumed to take on the values θ_{EB} or θ_{GCV} as appropriate.

6.1 Empirical Bayes: Gradient of the objective function using fast Bayesian transform

We refer back to Section 3.1, where we discuss about using gradient descent for hyperparameter search but the computational cost is of $\mathcal{O}(N_{\text{opt}} n^3)$. Here we develop a techniques to speed up the computation. If \mathbf{V} does not depend on θ then one can fast compute the derivative of Gram matrix \mathbf{C} . Starting from the definition (23) and taking derivative w.r.t. θ_ℓ ,

$$\frac{\partial \mathbf{C}}{\partial \theta_\ell} = \frac{1}{n} \mathbf{V} \frac{\partial \bar{\Lambda}}{\partial \theta_\ell} \mathbf{V}^H = \frac{1}{n} \mathbf{V} \bar{\Lambda}_{(\ell)} \mathbf{V}^H,$$

where $\bar{\Lambda}_{(\ell)} = \text{diag}(\bar{\lambda}_{(\ell)})$, and

$$\bar{\lambda}_{(\ell)} = \frac{\partial \lambda}{\partial \theta_\ell} = \left(\frac{\partial \lambda_i}{\partial \theta_\ell} \right)_{i=1}^n = \left(\frac{\partial}{\partial \theta_\ell} \mathbf{V}^H \mathbf{C}_1 \right) = \mathbf{V}^H \left(\frac{\partial}{\partial \theta_\ell} \mathbf{C}_\theta(\mathbf{x}_1, \mathbf{x}_i) \right)_{i=1}^n, \quad (29)$$

where we used the fast Bayesian transform property $\lambda = \mathbf{V}^H \mathbf{C}_1$. We use the notation $\bar{\lambda}_{(\ell)} = \mathbf{V}^H \bar{\mathbf{C}}_{1(\ell)}$ to denote the derivative of the eigenvalue $\lambda_{(\ell)}$, where $\bar{\mathbf{C}}_{1(\ell)}$ denotes the first row of the gram matrix after taking the derivative in the ℓ th variable, i.e.

$$\bar{\mathbf{C}}_{1(\ell)} = \left(\frac{\partial}{\partial \theta_\ell} \mathbf{C}_\theta(\mathbf{x}_1, \mathbf{x}_i) \right)_{i=1}^n.$$

The goal is to compute the derivative of the objective function faster. First, let's rewrite the objective function from (26a) in two parts,

$$\begin{aligned} \mathcal{L}_{\text{EB}}(\theta|\mathbf{y}) &= \underbrace{\frac{1}{n} \log(\det \mathbf{C})}_{\mathcal{L}_{|\mathbf{C}|}} + \underbrace{\log \left((\mathbf{y} - m_{\text{EB}} \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{EB}} \mathbf{1}) \right)}_{\mathcal{L}_{\mathbf{y}}}, \\ &=: \mathcal{L}_{|\mathbf{C}|} + \mathcal{L}_{\mathbf{y}}. \end{aligned}$$

Now, take the derivative:

$$\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\text{EB}}(\theta|\mathbf{y}) = \frac{\partial}{\partial \theta_\ell} \mathcal{L}_{|\mathbf{C}|} + \frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\mathbf{y}}.$$

Now we tackle the individual terms,

$$\begin{aligned}
\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{|\mathbf{C}|} &= \frac{\partial}{\partial \theta_\ell} \frac{1}{n} \log(\det \mathbf{C}) \\
&= \frac{1}{n} \text{trace} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_\ell} \right) = \frac{1}{n} \text{trace} \left(\mathbf{V} \Lambda^{-1} \mathbf{V}^H \frac{1}{n} \mathbf{V} \bar{\Lambda}_{(\ell)} \mathbf{V}^H \right) \\
&= \frac{1}{n} \text{trace} \left(\mathbf{V} \Lambda^{-1} \bar{\Lambda}_{(\ell)} \mathbf{V}^H \right), \quad \text{where we used } \mathbf{V}^H \mathbf{V} = \mathbf{I}, \\
&= \frac{1}{n} \text{trace} \left(\mathbf{V} \text{diag} \left(\frac{\bar{\lambda}_{i(\ell)}}{\lambda_i} \right)_{i=1}^n \mathbf{V}^H \right) = \frac{1}{n} \sum_{i=1}^n \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i},
\end{aligned}$$

where we used the fact from Higham [2008],

$$\log(\det \mathbf{C}) = \text{trace}(\log(\mathbf{C})).$$

Part of the $\mathcal{L}_{\mathbf{y}}$ was already simplified using the fast Bayesian transform,

$$(\mathbf{y} - m_{\text{EB}} \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{EB}} \mathbf{1}) = \frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i}.$$

Using the above result,

$$\begin{aligned}
\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\mathbf{y}} &= \frac{\partial}{\partial \theta_\ell} \log \left(\frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right) \\
&= \left(\frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right)^{-1} \frac{\partial}{\partial \theta_\ell} \left(\frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right) \\
&= \left(\frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right)^{-1} \frac{1}{n} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \left(-\frac{\partial \lambda_i}{\partial \theta_\ell} \right) \\
&= - \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right)^{-1} \left(\sum_{i=2}^n |\tilde{y}_i|^2 \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i^2} \right).
\end{aligned}$$

Finally, using the above results,

$$\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\text{EB}}(\boldsymbol{\theta}|\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i} - \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2 \bar{\lambda}_{i(\ell)}}{\lambda_i^2} \right) \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right)^{-1}, \quad (30)$$

where $\bar{\lambda}_{i(\ell)}$ is the derivative of the i th eigenvalue of \mathbf{C} in the ℓ th variable. Please recollect the gradient descent proposed in (21) can be computed faster in $\mathcal{O}(n \log n)$ using the result (30). A technique to compute this faster is discussed in Section 6.3.

6.2 Generalized Cross-Validation: Gradient of the objective function using fast Bayesian transform

Using the results obtained from the Section 6.1 with empirical Bayes, one can reduce the computational cost of the derivative of the objective function in (26b),

$$\mathcal{L}_{\text{GCV}}(\boldsymbol{\theta}|\mathbf{y}) = \log \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right).$$

Using the similar techniques from Section 6.1, the derivative of the objective function w.r.t θ_ℓ :

$$\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\text{GCV}}(\boldsymbol{\theta}|\mathbf{y})$$

$$\begin{aligned}
&= \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right)^{-1} \frac{\partial}{\partial \theta_\ell} \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right) - 2 \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right)^{-1} \frac{\partial}{\partial \theta_\ell} \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right) \\
&= \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right)^{-1} \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^3} (-2) \frac{\partial \lambda_i}{\partial \theta_\ell} \right) \\
&\quad - 2 \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right)^{-1} \left(\sum_{i=1}^n \frac{1}{\lambda_i^2} (-1) \frac{\partial \lambda_i}{\partial \theta_\ell} \right) \\
&= -2 \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right)^{-1} \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2 \bar{\lambda}_{i(\ell)}}{\lambda_i^3} \right) + 2 \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right)^{-1} \left(\sum_{i=1}^n \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i^2} \right).
\end{aligned}$$

Thus,

$$\begin{aligned}
\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\text{GCV}}(\boldsymbol{\theta} | \mathbf{y}) &= -2 \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right)^{-1} \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2 \bar{\lambda}_{i(\ell)}}{\lambda_i^3} \right) \\
&\quad + 2 \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right)^{-1} \left(\sum_{i=1}^n \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i^2} \right), \quad (31)
\end{aligned}$$

where $\bar{\lambda}_{i(\ell)}$ is the derivative of the i th eigenvalue of the Gram matrix, \mathbf{C} , in the ℓ th variable. We discuss a technique to compute $\bar{\lambda}_{i(\ell)}$ in the next section below.

6.3 Product Kernels

In this research, we use product kernels in the demonstrations and numerical implementations. They got nice properties which are helpful to obtain analytical results easily. Product kernels in d dimensions are of the form,

$$C_{\boldsymbol{\theta}}(\mathbf{t}, \mathbf{x}) = \prod_{\ell=1}^d \left[1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell) \right] \quad (32)$$

where η_ℓ is called shape parameter in the ℓ th variable for $\ell = 1, \dots, d$, and \mathfrak{C} is chosen such that to ensure $C_{\boldsymbol{\theta}}$ is symmetric and positive definite. Our goal is to compute $\bar{\lambda}_{i(\ell)}$ for which the kernel derivative is necessary. The derivative of the product kernels can be obtained easily. Please note that $\boldsymbol{\theta}$ denotes all the hyper parameters of the kernel C where η is one of them and called the shape parameter.

6.3.1 Derivative of the product kernel when $\eta_1 = \dots = \eta_d = \eta$

It was suggested to use gradient descent to find optimal shape parameter in Section 3.1. In this section, we compute the gradient for product kernels. When the $\eta_1 = \dots = \eta_d = \eta$, the derivative of a product kernel w.r.t. η can be obtained as below,

$$\begin{aligned}
\frac{\partial}{\partial \eta} C_{\boldsymbol{\theta}}(\mathbf{t}, \mathbf{x}) &= \frac{\partial}{\partial \eta} \prod_{j=1}^d \left[1 - \eta \mathfrak{C}(x_j, t_j) \right] \\
&= \sum_{\ell=1}^d \prod_{j=1, j \neq \ell}^d \left[1 - \eta \mathfrak{C}(x_j, t_j) \right] \left(-\mathfrak{C}(x_\ell, t_\ell) \right)
\end{aligned}$$

$$\begin{aligned}
&= \prod_{j=1}^d \left[1 - \eta \mathfrak{C}(x_j, t_j) \right] \sum_{\ell=1}^d \frac{\left(-\mathfrak{C}(x_\ell, t_\ell) \right)}{1 - \eta \mathfrak{C}(x_\ell, t_\ell)} \\
&= C_\theta(\mathbf{t}, \mathbf{x}) \frac{1}{\eta} \sum_{\ell=1}^d \frac{\left(1 - \eta \mathfrak{C}(x_\ell, t_\ell) - 1 \right)}{1 - \eta \mathfrak{C}(x_\ell, t_\ell)} \\
&= C_\theta(\mathbf{t}, \mathbf{x}) \frac{1}{\eta} \sum_{\ell=1}^d \left(1 - \frac{1}{1 - \eta \mathfrak{C}(x_\ell, t_\ell)} \right) \\
&= (d/\eta) \underbrace{\left(\prod_{j=1}^d \left[1 - \eta \mathfrak{C}(x_j, t_j) \right] \right)}_{C_\theta(\mathbf{t}, \mathbf{x})} \left(1 - \frac{1}{d} \sum_{\ell=1}^d \frac{1}{1 - \eta \mathfrak{C}(x_\ell, t_\ell)} \right).
\end{aligned}$$

Thus,

$$\frac{\partial}{\partial \eta} C_\theta(\mathbf{t}, \mathbf{x}) = (d/\eta) C_\theta(\mathbf{t}, \mathbf{x}) \left(1 - \frac{1}{d} \sum_{\ell=1}^d \frac{1}{1 - \eta \mathfrak{C}(x_\ell, t_\ell)} \right).$$

6.3.2 When η_ℓ is different for each $\ell = 1, \dots, d$

In this case, we will have a vector of length d shape parameters. Derivative of the kernel, $C_\theta(\mathbf{t}, \mathbf{x})$ (32), with respect to η_ℓ is,

$$\begin{aligned}
\frac{\partial}{\partial \eta_\ell} C_\theta(\mathbf{t}, \mathbf{x}) &= \frac{\partial}{\partial \eta_\ell} \prod_{j=1}^d \left[1 - \eta_j \mathfrak{C}(x_j, t_j) \right], \quad \text{where } \ell = 1, \dots, d \\
&= \prod_{j=1, j \neq \ell}^d \left[1 - \eta_j \mathfrak{C}(x_j, t_j) \right] \left(-\mathfrak{C}(x_\ell, t_\ell) \right) \\
&= \prod_{j=1}^d \left[1 - \eta_j \mathfrak{C}(x_j, t_j) \right] \frac{\left(-\mathfrak{C}(x_\ell, t_\ell) \right)}{1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell)} \\
&= C_\theta(\mathbf{t}, \mathbf{x}) \frac{1}{\eta_\ell} \frac{\left(1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell) - 1 \right)}{1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell)} \\
&= C_\theta(\mathbf{t}, \mathbf{x}) \frac{1}{\eta_\ell} \left(1 - \frac{1}{1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell)} \right) \\
&= \frac{1}{\eta_\ell} \underbrace{\left(\prod_{j=1}^d \left[1 - \eta_j \mathfrak{C}(x_j, t_j) \right] \right)}_{C_\theta(\mathbf{t}, \mathbf{x})} \left(1 - \frac{1}{1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell)} \right).
\end{aligned}$$

Thus,

$$\frac{\partial}{\partial \eta_\ell} C_\theta(\mathbf{t}, \mathbf{x}) = \frac{1}{\eta_\ell} C_\theta(\mathbf{t}, \mathbf{x}) \left(1 - \frac{1}{1 - \eta_\ell \mathfrak{C}(x_\ell, t_\ell)} \right).$$

Please note that the above derivatives do not depend on $\mathfrak{C}(x, t)$ and most importantly these computations are applicable to any product kernel of the form (32).

The $\bar{\lambda}_{i(\ell)}$ can be computed now using (29) with the computed kernel derivative, $\frac{\partial}{\partial \eta_\ell} C_\theta$.

6.4 Shape parameter search using steepest descent

Using the obtained derivative of the eigenvalues, $\bar{\lambda}_{i(\ell)}$, one can easily compute the gradient of the objective function (30) or (31). This can be further used to implement the steepest descent search as introduced in Section 3.1

$$\eta_\ell^{(j+1)} = \eta_\ell^{(j)} - \nu \frac{\partial}{\partial \eta_\ell} \mathcal{L}(\theta|\mathbf{y}), \quad j = 0, 1, \dots, \quad \ell = 1, \dots, d$$

where ν is the step size for the gradient descent, j is the iteration index, and $\frac{\partial}{\partial \eta_\ell} \mathcal{L}(\theta|\mathbf{y})$ is either (30) or (31) depending on the choice of the hyperparameter search method. The parameter η_ℓ is usually searched in the whole \mathbb{R} by using the simple domain transformation as explained in Section 12.

7 Sobol' Nets and Walsh Kernels

The previous section shows an automatic Bayesian cubature algorithm using rank-1 lattice nodes and shift-invariant kernels. In this chapter, we demonstrate a second approach to formulate fast Bayesian transform using matching kernel and point sets. Scrambled Sobol' nets and Walsh kernels are paired to achieve $\mathcal{O}(n^{-1+\epsilon})$ order error convergence where n is the sample size. Sobol' nets [Sobol' 1967] are low discrepancy points, used extensively in numerical integration, simulation, and optimization. The results of this chapter can be summarized as a theorem,

Theorem 3 *Any symmetric, positive definite, digital shift-invariant covariance kernel of the form (37) scaled to satisfy (25), when matched with digital net data-sites, satisfies assumptions (24). The fast Walsh-Hadamard transform (FWHT) can be used to expedite the estimates of θ in (26) and the credible interval widths (28) in $\mathcal{O}(n \log n)$ operations. The cubature, $\hat{\mu}$, is just the sample mean.*

We introduce the necessary concepts and prove this theorem in the remaining of this chapter.

8 Sobol' Nets

Nets were developed to provide deterministic sample points for quasi-Monte Carlo rules Niederreiter [2005]. Nets are defined geometrically using elementary intervals, which are subintervals of the unit cube $[0, 1]^d$. The (t, m, d) -nets in base b , introduced by Niederreiter, whose quality is governed by t . Lower values of t correspond to (t, m, d) -nets of higher quality Baldeaux [2010].

Definition 1 *Let \mathcal{A} be the set of all elementary intervals $\mathcal{A} \subset [0, 1]^d$ where $\mathcal{A} = \prod_{\ell=1}^d [\alpha_\ell b^{-\gamma_\ell}, (\alpha_\ell + 1)b^{-\gamma_\ell})$, with $d, b, \gamma_\ell \in \mathbb{N}, b \geq 2$ and $b^{\gamma_\ell} > \alpha_\ell \geq 0$. For $m, t \in \mathbb{N}, m \geq t \geq 0$, the point set $\mathcal{P}_m \in [0, 1]^d$ with $n = b^m$ points is a (t, m, d) -net in base b if every \mathcal{A} with volume b^{t-m} contains b^t points of \mathcal{P}_m .*

Digital (t, m, d) -nets are a special case of (t, m, d) -nets, constructed using matrix-vector multiplications over finite fields. Digital sequences are infinite length digital nets, i.e., the first $n = b^m$ points of a digital sequence comprise a digital net for all integer $m \in \mathbb{N}_0$.

Definition 2 For any non-negative integer $i = \dots i_3 i_2 i_1$ (base b), define the $\infty \times 1$ vector \mathbf{i} as the vector of its digits, that is, $\mathbf{i} = (i_1, i_2, \dots)^T$. For any point $z = 0.z_1 z_2 \dots$ (base b) $\in [0, 1]$, define the $\infty \times 1$ vector of the digits of z , that is, $\mathbf{z} = (z_1, z_2, \dots)^T$. Let $\mathbf{G}_1, \dots, \mathbf{G}_d$ denote predetermined $\infty \times \infty$ generator matrices. The digital sequence in base b is $\{\mathbf{z}_0, \mathbf{z}_1, \mathbf{z}_2, \dots\}$, where each $\mathbf{z}_i = (z_{i1}, \dots, z_{id})^T \in [0, 1]^d$ is defined by

$$\mathbf{z}_{i\ell} = \mathbf{G}_\ell \mathbf{i}, \quad \ell = 1, \dots, d, \quad i = 0, 1, \dots$$

The value of t as mentioned in Definition 1 depends on the choice of \mathbf{G}_ℓ .

Digital nets have a group structure under digitwise addition, which is a very useful property exploited in our algorithm, especially to develop a fast Bayesian transform that speedups computations. Digitwise addition, \oplus , and subtraction \ominus , are defined in terms of b -ary expansions of points in $[0, 1]^d$,

$$\mathbf{z} \oplus \mathbf{y} = \left(\sum_{j=1}^{\infty} [z_{\ell j} + y_{\ell j} \bmod b] b^{-j} \bmod 1 \right)_{\ell=1}^d, \\ \mathbf{z} \ominus \mathbf{y} = \left(\sum_{j=1}^{\infty} [z_{\ell j} - y_{\ell j} \bmod b] b^{-j} \bmod 1 \right)_{\ell=1}^d,$$

where

$$\mathbf{z} = \left(\sum_{j=1}^{\infty} z_{\ell j} b^{-j} \right)_{\ell=1}^d, \quad \mathbf{y} = \left(\sum_{j=1}^{\infty} y_{\ell j} b^{-j} \right)_{\ell=1}^d, \quad z_{\ell j}, y_{\ell j} \in \{0, \dots, b-1\}.$$

Similarly for integer values in \mathbb{N}_0^d , the digitwise addition, \oplus , and subtraction \ominus , are defined in terms of their b -ary expansions,

$$\mathbf{k} \oplus \mathbf{l} = \left(\sum_{j=0}^{\infty} [k_{\ell j} + l_{\ell j} \bmod b] b^j \bmod 1 \right)_{\ell=1}^d, \\ \mathbf{k} \ominus \mathbf{l} = \left(\sum_{j=0}^{\infty} [k_{\ell j} - l_{\ell j} \bmod b] b^j \bmod 1 \right)_{\ell=1}^d,$$

where

$$\mathbf{k} = \left(\sum_{j=0}^{\infty} k_{\ell j} b^j \right)_{\ell=1}^d, \quad \mathbf{l} = \left(\sum_{j=0}^{\infty} l_{\ell j} b^j \right)_{\ell=1}^d, \quad k_{\ell j}, l_{\ell j} \in \{0, \dots, b-1\}.$$

Let $\{\mathbf{z}_i\}_{i=0}^{b^m-1}$ be a digital net. Then

$$\forall i_1, i_2 \in \{0, \dots, b^m-1\}, \quad \mathbf{z}_{i_1} \oplus \mathbf{z}_{i_2} = \mathbf{z}_{i_3}, \quad \text{for some } i_3 \in \{0, \dots, b^m-1\}.$$

The following very useful result, which will be further used to obtain the fast Bayesian transform, arises from the fundamental property of digital nets.

Lemma 2 Let $\{\mathbf{z}_i\}_{i=0}^{b^m-1}$ be the digital-net and the corresponding digitally shifted net be $\{\mathbf{x}_i\}_{i=0}^{b^m-1}$, i.e.,

$$\mathbf{x}_{i\ell} = \mathbf{z}_{i\ell} + \Delta_\ell \bmod 1,$$

where $\mathbf{x}_{i\ell}$ is the ℓ th component of i th digital net and Δ_ℓ is the digital shift for the ℓ th component. Then,

$$\mathbf{x}_i \ominus \mathbf{x}_j = \mathbf{z}_i \ominus \mathbf{z}_j = \mathbf{z}_{i \ominus j}, \quad \forall i, j \in \mathbb{N}_0. \quad (33)$$

Also the digital subtraction is symmetric,

$$\mathbf{x}_i \ominus \mathbf{x}_i = \mathbf{0}, \quad \mathbf{x}_i \ominus \mathbf{x}_j = \mathbf{x}_j \ominus \mathbf{x}_i, \quad \forall i, j \in \mathbb{N}_0. \quad (34)$$

Proof The proof can be obtained from the definition of digital nets which stated that the digital nets are obtained using generator matrices, $\mathbf{z}_{i\ell} = \mathbf{G}_\ell \mathbf{i} \bmod b$. Rewriting the subtraction using the generating matrix provides the result,

$$\begin{aligned} \mathbf{z}_{i\ell} - \mathbf{z}_{j\ell} \bmod b &= (\mathbf{G}_\ell \mathbf{i} \bmod b) - (\mathbf{G}_\ell \mathbf{j} \bmod b) \\ &= (\mathbf{G}_\ell \mathbf{i} - \mathbf{G}_\ell \mathbf{j}) \bmod b \\ &= \mathbf{G}_\ell (\mathbf{i} - \mathbf{j}) \bmod b \\ &= \mathbf{G}_\ell (\overline{i \ominus j}) \bmod b \\ &= \mathbf{z}_{i \ominus j \ell}. \end{aligned}$$

The rest of the lemma is obvious from the definition of digital nets.

We chose digitally shifted and scrambled nets Hickernell and Yue [2000] for our Bayesian cubature algorithm. Digital shifts help to avoid having nodes at the origin, similar to the random shift used with lattice nodes. Scrambling helps to eliminate bias while retaining the low-discrepancy properties. A proof that a scrambled net preserves the property of (t, m, d) -net almost surely can be found in Owen Owen [1995]. The scrambling method proposed by Matoušek Matoušek [1998] is preferred since it is more efficient than the Owen's scrambling.

Sobol' nets Sobol' [1976] are a special case of (t, m, d) -nets when base $b = 2$. An example of 64 Sobol' nets in $d = 2$ is given in Figure 3. The even coverage of the unit cube is ensured by a well chosen generating matrix. The choice of generating vector is typically done offline by computer search. See Kuo and Nuyens [2016] and Nuyens for more on generating matrices. We use randomly scrambled and digitally shifted Sobol' sequences in this research Hong and Hickernell [2003].

9 Walsh Kernels

Walsh kernels are product kernels based on the Walsh functions. We introduce the necessary concepts in this section.

9.1 Walsh functions

Like the Fourier transform used with lattice points (Jagadeeswaran and Hickernell [2019]), the Walsh-Hadamard transform, which we will simply call Walsh transform, is used for the digital nets. The Walsh transform is defined using Walsh functions. Recall $\mathbb{N}_0 := \{0, 1, 2, \dots\}$. The one-dimensional Walsh functions in base b are defined as

$$\text{wal}_{b,k}(x) := e^{2\pi\sqrt{-1}(x_1 k_0 + x_2 k_1 + \dots)/b} = e^{2\pi\sqrt{-1}\mathbf{k}^T \mathbf{x}/b}, \quad (35)$$

for $x \in [0, 1)$ and $k \in \mathbb{N}_0$ and the unique base b expansions $x = \sum_{j \geq 1} x_j b^{-j} = (0.x_1 x_2 \dots)_b$, $\mathbf{x} = (x_1, x_2, \dots)^T$, $k = \sum_{j \geq 0} k_j b^j = (\dots k_1 k_0)_b$, $\mathbf{k} = (k_0, k_1, \dots)^T$,

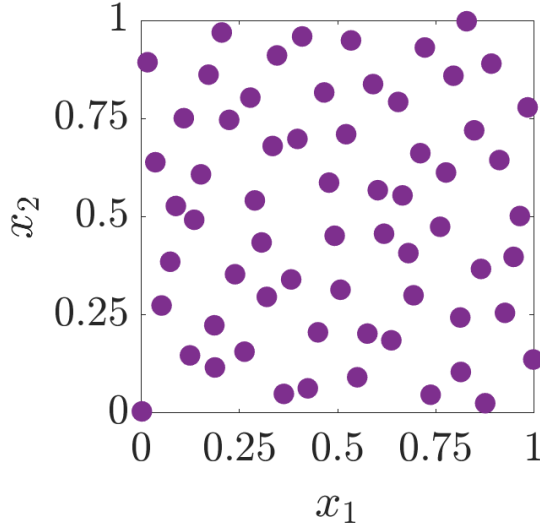


Fig. 3: Example of a scrambled Sobol' node set in $d = 2$. This plot can be reproduced using `PlotPoints.m`.

and $\mathbf{k}^T \mathbf{x} = x_1 k_0 + x_2 k_1 + \dots$ where the number of digits used in (35) are limited to the length required to represent x or k , i.e., $\max(\lceil -\log_b x \rceil, \lceil \log_b k \rceil)$. Multivariate Walsh functions are defined as the product of the one-dimensional Walsh functions,

$$\text{wal}_{b,\mathbf{k}}(\mathbf{x}) := \prod_{\ell=1}^d \text{wal}_{b,k_\ell}(x_\ell)$$

As shown in (35), for the case of $b = 2$, the Walsh functions only take the values in $\{1, -1\}$, i.e., $\text{wal}_{b,\mathbf{k}} : [0, 1)^d \rightarrow \{-1, 1\}$, $\mathbf{k} \in \mathbb{N}_0^d$. Walsh functions form an orthonormal basis of the Hilbert space $L^2[0, 1)^d$,

$$\int_{[0,1)^d} \text{wal}_{b,\mathbf{l}}(\mathbf{x}) \text{wal}_{b,\mathbf{k}}(\mathbf{x}) d\mathbf{x} = \delta_{\mathbf{l},\mathbf{k}}, \quad \forall \mathbf{l}, \mathbf{k} \in \mathbb{N}_0^d$$

Digital nets are designed to integrate certain Walsh functions without error. Thus our Bayesian cubature algorithm integrates linear combinations of certain Walsh functions without error. Functions that are well approximated by such linear combinations are then integrated with small errors.

In this research we use Sobol' nodes which are digital nets with base $b = 2$. So here afterwards base $b = 2$ is assumed. In this case, the Walsh function is simply

$$\text{wal}_{2,\mathbf{k}}(\mathbf{x}) = (-1)^{\mathbf{k}^T \mathbf{x}}.$$

9.2 Walsh kernels

Consider the covariance kernels of the form,

$$C_\theta(\mathbf{x}, \mathbf{t}) = K_\theta(\mathbf{x} \ominus \mathbf{t}) \tag{36}$$

where \ominus is bitwise subtraction. This is called a *digitally shift invariant kernel* because shifting both arguments of the covariance function by the same amount

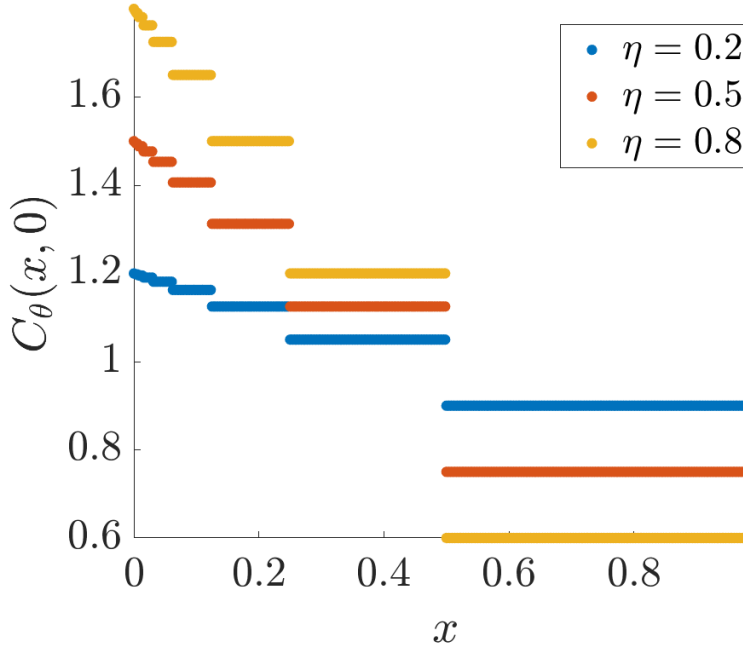


Fig. 4: Walsh kernel of order $r = 1$ in dimension $d = 1$. This figure can be reproduced using `plot_walsh_kernel.m`.

leaves the value unchanged. By a proper scaling of the function K_{θ} , it follows that assumption (25) is satisfied. The function K_{θ} must be of the form that ensures that C_{θ} is symmetric and positive definite, as assumed in (6). We drop the θ sometimes to make the notation simpler. The Walsh kernels are of the form,

$$K_{\theta}(\mathbf{x} \ominus \mathbf{t}) = \prod_{\ell=1}^d 1 + \eta_{\ell} \omega_r(x_{\ell} \ominus t_{\ell}), \quad \boldsymbol{\eta} = (\eta_1, \dots, \eta_d), \quad \boldsymbol{\theta} = (r, \boldsymbol{\eta}) \quad (37)$$

where r is the kernel order, $\boldsymbol{\eta}$ is the kernel shape parameter, and

$$\omega_r(x) = \sum_{k=1}^{\infty} \frac{\text{wal}_{2,k}(x)}{2^{2r \lfloor \log_2 k \rfloor}}.$$

Explicit expression is available for ω_r in the case of order $r = 1$ Nuyens [2013],

$$\omega_1(x) = 6 \left(\frac{1}{6} - 2^{\lfloor \log_2 x \rfloor - 1} \right). \quad (38)$$

The Figure 4 shows the Walsh kernel (37) of order $r = 1$ in the interval $[0, 1)$. Unlike the shift-invariant kernels used with lattice nodes, low order Walsh kernels are discontinuous and are only piecewise constant. Smaller η_{ℓ} implies lesser variation in the amplitude of the kernel. Also, the Walsh kernels are digitally shift invariant but not periodic.

10 Eigenvectors

We show the eigenvectors \mathbf{V} in (23) of the Gram matrix formed by the covariance kernel (37) and Sobol' nets are the columns of the Walsh-Hadamard matrix. First we introduce the necessary concepts.

10.1 Walsh transform

The Walsh-Hadamard transform (WHT) is a generalized class of discrete Fourier transform (DFT) and is much simpler to compute than the DFT. The WHT matrices are comprised of only ± 1 values, so the computation usually involves only ordinary additions and subtractions. Hence, the WHT is also sometimes called the integer transform. In comparison, the DFT that was used with lattice nodes, uses complex exponential functions and the computation involves complex, non-integer multiplications.

The WHT involves multiplications by $2^m \times 2^m$ Walsh-Hadamard matrices, which is constructed recursively, starting with $\mathbf{H}^{(0)} = 1$,

$$\begin{aligned} \mathbf{H}^{(1)} &= \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \\ \mathbf{H}^{(2)} &= \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}, \\ &\vdots \\ \mathbf{H}^{(m)} &= \begin{pmatrix} \mathbf{H}^{(m-1)} & \mathbf{H}^{(m-1)} \\ \mathbf{H}^{(m-1)} & -\mathbf{H}^{(m-1)} \end{pmatrix} = \underbrace{\mathbf{H}^{(1)} \otimes \dots \otimes \mathbf{H}^{(1)}}_{m \text{ times}} = \mathbf{H}^{(1)} \otimes \mathbf{H}^{(m-1)} \end{aligned} \quad (39)$$

where \otimes is Kronecker product. Alternatively for base $b = 2$, these matrices can be directly obtained by,

$$\mathbf{H}^{(m)} = \left((-1)^{(\mathbf{i}^T \mathbf{j})} \right)_{i,j=0}^{2^m-1},$$

where the notation $\mathbf{i}^T \mathbf{j}$ indicates the bitwise dot product.

10.2 Eigenvectors of \mathbf{C} are columns of Walsh-Hadamard matrix

The Gram matrix \mathbf{C}_θ formed by Walsh kernels and Sobol' nodes have a special structure called block-Toeplitz, which can be used to construct the fast Bayesian transform. A Toeplitz matrix is a diagonal-constant matrix in which each descending diagonal from left to right is constant. A block Toeplitz matrix is a special block matrix, which contains blocks that are repeated down the diagonals of the matrix. We prove that the eigenvectors of \mathbf{C}_θ are columns of a Walsh-Hadamard matrix in two theorems.

Theorem 4 Let $(\mathbf{x}_i)_{i=0}^{n-1}$ be digitally shifted Sobol' nodes and K be any function, then the Gram matrix,

$$\mathbf{C}_\theta = (C(\mathbf{x}_i, \mathbf{x}_j))_{i,j=0}^{n-1} = (K(\mathbf{x}_i \ominus \mathbf{x}_j))_{i,j=0}^{n-1},$$

where $n = 2^m$, $C(\mathbf{x}, \mathbf{t}) = K(\mathbf{x} \ominus \mathbf{t})$, $\mathbf{x}, \mathbf{t} \in [0, 1]^d$, is a 2×2 block-Toeplitz matrix and all the sub-blocks and their sub-sub-blocks, etc. are also 2×2 block-Toeplitz.

Proof We prove this theorem by induction. Let $\mathbf{C}_\theta^{(m)}$ denote the Gram matrix of size $2^m \times 2^m$. The relation between sub-block matrices can be deciphered using the properties of digital nets. To help with the proof of block-Toeplitz structure, consider the digital net properties (33), (34), and notations,

$$\mathbf{K}^{(m)} := (K(\mathbf{z}_i \ominus \mathbf{z}_j))_{i,j=0}^{2^m-1} = (K(\mathbf{z}_{i \ominus j}))_{i,j=0}^{2^m-1}, \quad m = 1, 2, \dots,$$

$$\mathbf{K}^{(m,q)} := (K(\mathbf{z}_{i \ominus j + q2^m}))_{i,j=0}^{2^m-1}, \quad q = 0, 1, \dots$$

These two notations are related by $\mathbf{K}^{(m)} = \mathbf{K}^{(m,0)}$. Please note that $\mathbf{C}_\theta^{(m)} = \mathbf{K}^{(m,0)}$.

We will prove $\mathbf{K}^{(m,q)}$ is a 2×2 block-toeplitz matrix for all $m \in \mathbb{N}, q \in \mathbb{N}$.

As the first step, we verify the property holds for $m = 1$,

$$\mathbf{K}^{(1,q)} = \begin{pmatrix} K(\mathbf{z}_{0 \ominus 0 + q2^1}) & K(\mathbf{z}_{1 \ominus 0 + q2^1}) \\ K(\mathbf{z}_{0 \ominus 1 + q2^1}) & K(\mathbf{z}_{1 \ominus 1 + q2^1}) \end{pmatrix} = \begin{pmatrix} K(\mathbf{z}_{2q}) & K(\mathbf{z}_{1+2q}) \\ K(\mathbf{z}_{1+2q}) & K(\mathbf{z}_{2q}) \end{pmatrix}, \quad \text{by (33)}$$

has diagonal elements repeated. Thus by definition, it is a 2×2 block-Toeplitz.

Now assume that $\mathbf{K}^{(m,q)}$ is block-Toeplitz. We need to prove $\mathbf{K}^{(m+1,q)}$ is also a 2×2 block-Toeplitz. Let $n = 2^m$,

$$\begin{aligned} \mathbf{K}^{(m+1)} &= \begin{pmatrix} K(\mathbf{z}_{0 \ominus 0}) & \dots & K(\mathbf{z}_{0 \ominus n-1}) & K(\mathbf{z}_{0 \ominus n}) & \dots & K(\mathbf{z}_{0 \ominus 2n-1}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ K(\mathbf{z}_{n-1 \ominus 0}) & \dots & K(\mathbf{z}_{n-1 \ominus n-1}) & K(\mathbf{z}_{n-1 \ominus n}) & \dots & K(\mathbf{z}_{n-1 \ominus 2n-1}) \\ K(\mathbf{z}_{n \ominus 0}) & \dots & K(\mathbf{z}_{n \ominus n-1}) & K(\mathbf{z}_{n \ominus n}) & \dots & K(\mathbf{z}_{n \ominus 2n-1}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ K(\mathbf{z}_{2n-1 \ominus 0}) & \dots & K(\mathbf{z}_{2n-1 \ominus n-1}) & K(\mathbf{z}_{2n-1 \ominus n}) & \dots & K(\mathbf{z}_{2n-1 \ominus 2n-1}) \end{pmatrix} \\ &= \begin{pmatrix} \begin{pmatrix} K(\mathbf{z}_0) & \dots & K(\mathbf{z}_{n-1}) \\ \vdots & \vdots & \vdots \\ K(\mathbf{z}_{n-1}) & \dots & K(\mathbf{z}_0) \end{pmatrix} & \begin{pmatrix} K(\mathbf{z}_n) & \dots & K(\mathbf{z}_{2n-1}) \\ \vdots & \vdots & \vdots \\ K(\mathbf{z}_{2n-1}) & \dots & K(\mathbf{z}_n) \end{pmatrix} \\ \begin{pmatrix} K(\mathbf{z}_n) & \dots & K(\mathbf{z}_{2n-1}) \\ \vdots & \vdots & \vdots \\ K(\mathbf{z}_{2n-1}) & \dots & K(\mathbf{z}_n) \end{pmatrix} & \begin{pmatrix} K(\mathbf{z}_0) & \dots & K(\mathbf{z}_{n-1}) \\ \vdots & \vdots & \vdots \\ K(\mathbf{z}_{n-1}) & \dots & K(\mathbf{z}_0) \end{pmatrix} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{K}^{(m)} & \mathbf{K}^{(m,1)} \\ \mathbf{K}^{(m,1)} & \mathbf{K}^{(m)} \end{pmatrix} \end{aligned}$$

is a 2×2 block-Toeplitz, where we used the properties (33), (34) and facts $2n - 1 \ominus n = n - 1$, $2n - 1 \ominus n - 1 = n$, and $n \ominus n - 1 = 2n - 1$. Thus $\mathbf{K}^{(m+1)}$ is a 2×2 block-Toeplitz. Similarly

$$\mathbf{K}^{(m+1,q)} = \begin{pmatrix} \mathbf{K}^{(m,q)} & \mathbf{K}^{(m,q+1)} \\ \mathbf{K}^{(m,q+1)} & \mathbf{K}^{(m,q)} \end{pmatrix}$$

is a 2×2 block-Toeplitz. Thus $\mathbf{C}_\theta^{(m)}$ of size $2^m \times 2^m$, for $m \in \mathbb{N}$, is a 2×2 block-Toeplitz and every block and it's sub-blocks of size 2^p , $p \in \mathbb{N}$, $p \leq m$ are also 2×2 block-Toeplitz.

Theorem 5 The Walsh-Hadamard matrix $\mathbf{H}^{(m)}$ factorizes $\mathbf{C}_\theta^{(m)}$, so that the columns of Walsh-Hadamard matrix are the eigenvectors of $\mathbf{C}_\theta^{(m)}$, i.e.,

$$\mathbf{H}^{(m)} \mathbf{C}_\theta^{(m)} = \mathbf{\Lambda}^{(m)} \mathbf{H}^{(m)}, \quad m \in \mathbb{N}.$$

Proof Again, we use the proof-by-induction technique to show that the Walsh-Hadamard matrix factorizes $\mathbf{K}^{(m,q)}$. We can easily see the Hadamard matrix $\mathbf{H}^{(1)}$ diagonalizes $\mathbf{K}^{(1,q)}$,

$$\begin{aligned}\mathbf{H}^{(1)}\mathbf{K}^{(1,q)} &= \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} K(\mathbf{z}_{0+q2^1}) & K(\mathbf{z}_{1+q2^1}) \\ K(\mathbf{z}_{1+q2^1}) & K(\mathbf{z}_{0+q2^1}) \end{pmatrix}, \quad \text{by Theorem 4} \\ &= \begin{pmatrix} K(\mathbf{z}_{2q}) + K(\mathbf{z}_{2q+1}) & K(\mathbf{z}_{2q}) - K(\mathbf{z}_{2q+1}) \\ K(\mathbf{z}_{2q}) - K(\mathbf{z}_{2q+1}) & K(\mathbf{z}_{2q}) + K(\mathbf{z}_{2q+1}) \end{pmatrix} \\ &= \begin{pmatrix} K(\mathbf{z}_{2q}) + K(\mathbf{z}_{2q+1}) & 0 \\ 0 & K(\mathbf{z}_{2q}) - K(\mathbf{z}_{2q+1}) \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \\ &= \Lambda^{(1,q)}\mathbf{H}^{(1)},\end{aligned}$$

where $\Lambda^{(1,q)}$ is a diagonal matrix, thus $\mathbf{H}^{(1)}$ factorizes $\mathbf{K}^{(1,q)}$.

Now assume $\mathbf{H}^{(m)}$ factorizes $\mathbf{K}^{(m,q)}$, so $\mathbf{H}^{(m)}\mathbf{K}^{(m,q)} = \Lambda^{(m,q)}\mathbf{H}^{(m)}$ where $\Lambda^{(m,q)}$ is diagonal. We need to prove $\mathbf{H}^{(m+1)}$ factorizes $\mathbf{K}^{(m+1,q)}$,

$$\begin{aligned}\mathbf{H}^{(m+1)}\mathbf{K}^{(m+1,q)} &= \begin{pmatrix} \mathbf{H}^{(m)} & \mathbf{H}^{(m)} \\ \mathbf{H}^{(m)} & -\mathbf{H}^{(m)} \end{pmatrix} \begin{pmatrix} \mathbf{K}^{(m,q)} & \mathbf{K}^{(m,q+1)} \\ \mathbf{K}^{(m,q+1)} & \mathbf{K}^{(m,q)} \end{pmatrix}, \quad \text{by Theorem 4} \\ &= \begin{pmatrix} \mathbf{H}^{(m)}(\mathbf{K}^{(m,q)} + \mathbf{K}^{(m,q+1)}) & \mathbf{H}^{(m)}(\mathbf{K}^{(m,q)} - \mathbf{K}^{(m,q+1)}) \\ \mathbf{H}^{(m)}(\mathbf{K}^{(m,q)} - \mathbf{K}^{(m,q+1)}) & \mathbf{H}^{(m)}(\mathbf{K}^{(m,q+1)} + \mathbf{K}^{(m,q)}) \end{pmatrix} \\ &= \begin{pmatrix} (\Lambda^{(m,q)} + \Lambda^{(m,q+1)})\mathbf{H}^{(m)} & (\Lambda^{(m,q)} - \Lambda^{(m,q+1)})\mathbf{H}^{(m)} \\ (\Lambda^{(m,q)} - \Lambda^{(m,q+1)})\mathbf{H}^{(m)} & (\Lambda^{(m,q+1)} + \Lambda^{(m,q)})\mathbf{H}^{(m)} \end{pmatrix} \\ &= \begin{pmatrix} \Lambda^{(m,q)} + \Lambda^{(m,q+1)} & 0 \\ 0 & \Lambda^{(m,q)} - \Lambda^{(m,q+1)} \end{pmatrix} \begin{pmatrix} \mathbf{H}^{(m)} & \mathbf{H}^{(m)} \\ \mathbf{H}^{(m)} & -\mathbf{H}^{(m)} \end{pmatrix} \\ &= \Lambda^{(m+1,q)}\mathbf{H}^{(m+1)}.\end{aligned}$$

Thus, $\mathbf{H}^{(m+1)}$ factorizes $\mathbf{K}^{(m+1,q)}$ to a diagonal matrix $\Lambda^{(m+1,q)}$. This implies $\mathbf{H}^{(p)}$ factorizes $\mathbf{C}_\theta^{(p)}$ for $p \in \mathbb{N}$. Please recall $\mathbf{C}_\theta^{(p)} = \mathbf{K}^{(p,0)}$. Here we used the fact that both \mathbf{H} and \mathbf{K} are symmetric positive definite.

10.3 Fast Bayesian transform

We can easily show that the Walsh-Hadamard matrices satisfy the assumptions of fast Bayesian transform (24). As shown in Section 10.2 the columns of $\mathbf{H}^{(m)}$ are the eigenvectors. Since the Gram matrix \mathbf{C} is symmetric, the columns/rows of Walsh-Hadamard matrices are mutually orthogonal. Thus the Gram matrix can be written as

$$\mathbf{C}^{(m)} = \frac{1}{n}\mathbf{H}^{(m)}\Lambda^{(m)}\mathbf{H}^{(m)}, \quad \text{where} \quad \mathbf{H}^{(m)} = \underbrace{\mathbf{H}^{(1)} \otimes \dots \otimes \mathbf{H}^{(1)}}_{m \text{ times}}. \quad (40)$$

Assumption (24b) follows automatically by the fact that Walsh-Hadamard matrices can be constructed analytically. Assumption (24a) can also be verified as the first row/column are one vectors. Finally, assumption (24c) is satisfied due to the fact that fast Walsh transform can be computed in $\mathcal{O}(n \log n)$ operations using fast Walsh-Hadamard transform. Thus the Walsh-Hadamard transform is a fast Bayesian transform, $\mathbf{V} := \mathbf{H}$, as per (24).

We have implemented a fast adaptive Bayesian cubature algorithm using the kernel (37) with $r = 1$ and Sobol' points Bratley and Fox [1988] in MATLAB as

part of the Guaranteed Adaptive Integration Library (GAIL) Choi et al. [2013–2017] as `cubBayesNet.g`. The Sobol' points used in this algorithm are generated using MATLAB's builtin function `sobolset` and scrambled using MATLAB function `scramble` Hong and Hickernell [2003]. The fast Walsh-Hadamard transform (40) is computed using MATLAB's builtin function `fwht` with *hadamard* ordering.

10.4 Iterative Computation of Walsh Transform

In every iteration of our algorithm, we double the number of function values. Using the technique described here, we have to only compute the Walsh transform for the newly added function values. Similar to the lattice points, Sobol' points are extensible by definition. This property is used in our algorithm to improve the integration accuracy till the required error tolerance is met. Sobol' nodes can be combined with Hadamard matrices as demonstrated here for iterative computation. Let $\tilde{\mathbf{y}} = \mathbf{H}^{(m+1)}\mathbf{y}$ for some arbitrary $\mathbf{y} \in \mathbb{R}^{2n}$, $n = 2^m$. Define,

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_{2n} \end{pmatrix}, \quad \mathbf{y}^{(1)} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{y}^{(2)} = \begin{pmatrix} y_{n+1} \\ \vdots \\ y_{2n} \end{pmatrix},$$

$$\tilde{\mathbf{y}}^{(1)} = \mathbf{H}^{(m)}\mathbf{y}^{(1)} = \begin{pmatrix} \tilde{y}_1^{(1)} \\ \tilde{y}_2^{(1)} \\ \vdots \\ \tilde{y}_n^{(1)} \end{pmatrix}, \quad \tilde{\mathbf{y}}^{(2)} = \mathbf{H}^{(m)}\mathbf{y}^{(2)} = \begin{pmatrix} \tilde{y}_1^{(2)} \\ \tilde{y}_2^{(2)} \\ \vdots \\ \tilde{y}_n^{(2)} \end{pmatrix}.$$

Then,

$$\begin{aligned} \tilde{\mathbf{y}} &= \mathbf{H}^{(m+1)}\mathbf{y} \\ &= \begin{pmatrix} \mathbf{H}^{(m)} & \mathbf{H}^{(m)} \\ \mathbf{H}^{(m)} & -\mathbf{H}^{(m)} \end{pmatrix} \begin{pmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \end{pmatrix}, \quad \text{by (39)} \\ &= \begin{pmatrix} \mathbf{H}^{(m)}\mathbf{y}^{(1)} + \mathbf{H}^{(m)}\mathbf{y}^{(2)} \\ \mathbf{H}^{(m)}\mathbf{y}^{(1)} - \mathbf{H}^{(m)}\mathbf{y}^{(2)} \end{pmatrix} \\ &= \begin{pmatrix} \tilde{\mathbf{y}}^{(1)} + \tilde{\mathbf{y}}^{(2)} \\ \tilde{\mathbf{y}}^{(1)} - \tilde{\mathbf{y}}^{(2)} \end{pmatrix} =: \tilde{\mathbf{y}}. \end{aligned}$$

As before with the lattice nodes, the computational cost to compute $\mathbf{V}^{(m+1)H}\mathbf{y}$ is twice the cost of computing $\mathbf{V}^{(m)H}\mathbf{y}^{(1)}$ plus $2n$ additions, where $n = 2^m$. An inductive argument shows that for any $m \in \mathbb{N}$, $\mathbf{V}^{(m)H}\mathbf{y}$ requires only $\mathcal{O}(n \log n)$ operations. Usually the multiplications in $\mathbf{V}^{(m)H}\mathbf{y}^{(1)}$ are multiplications by -1 which are simply accomplished using sign change or negation, requiring no multiplications at all.

11 Higher Order Nets

Higher order digital nets are an extension of (t, m, d) -nets, introduced in Dick [2008]. They can be used to numerically integrate smoother functions which are not necessarily periodic, but have square integrable mixed partial derivatives of order α , at a rate of $\mathcal{O}(n^{-\alpha})$ multiplied by a power of a $\log n$ factor using rules

corresponding to the modified (t, m, d) -nets. We want to emphasize that quasi-Monte Carlo rules based on these point sets can achieve convergence rates faster than $\mathcal{O}(n^{-1})$. Higher order digital nets are constructed using matrix-vector multiplications over finite fields.

One could develop matching digitally shift invariant kernels to formulate the fast Bayesian cubature. Bayesian cubatures using higher order digital nets are a topic for future research.

12 Kernel Hyperparameters Search

JR: Explain the transformation used to make the search range positive, > 0 , etc.

The various hyperparameters introduced and used by our algorithms need to be optimally chosen. The parameter search can be done in two major ways. Bounded minima search, if the search interval is known, else unbounded search. Most of the scenarios, the search interval is unknown. So the natural choice is to use unbounded search over the unbound domain such as `fminsearch` provided by MATLAB. However hyperparameters need to live in a domain that is bounded or semi-bounded. There are some simple domain transformations available to achieve this.

12.1 Positive kernel shape parameter

The following parameter map is used to ensure that the shape parameter values are positive real numbers. For $\eta > 0$ as introduced in Section 6.3.1, let

$$\eta(t_1) = e^{t_1}, \quad \eta : (-\infty, \infty) \rightarrow (0, \infty).$$

Instead of searching for $\eta \in (0, \infty)$, we may search for the optimal $t_1 = \log(\eta)$ over the whole real line \mathbb{R} . The optimal value $t_{1,\text{opt}}$ can be transformed back to the $(0, \infty)$ interval using

$$\eta_{\text{opt}} = e^{t_{1,\text{opt}}}.$$

12.2 Kernel order $1 < r < \infty$

The following map is used to ensure that the kernel order values are positive real number and greater than one, i.e., in the $(1, \infty)$ interval as required in Section ??,

$$r(t_2) = 1 + e^{t_2}, \quad r : (-\infty, \infty) \rightarrow (1, \infty).$$

So one may search for the optimal $t_2 = -\log(r - 1)$ in the whole real line \mathbb{R} . The optimal value $t_{2,\text{opt}}$ can be transformed back to the desired interval $(0, 1)$ using

$$r_{\text{opt}} = 1 + e^{t_{2,\text{opt}}}.$$

12.3 Kernel order $0 < q < 1$

The following multivariate map is used to ensure that the kernel order values are positive real and less than one, i.e., in the $(0, 1)$ interval to use with exponentially decaying kernel, as introduced in Section ??,

$$q(t_3) = \frac{1}{1 + e^{t_3}}, \quad q : (-\infty, \infty) \rightarrow (0, 1).$$

So one may search for the optimal $t_3 = \log(q^{-1} - 1)$ in the whole real line \mathbb{R} . The optimal value $t_{3,\text{opt}}$ can be transformed back to the desired interval $(0, 1)$ by using

$$q_{\text{opt}} = \frac{1}{1 + e^{t_{3,\text{opt}}}}.$$

12.4 Combined searching of kernel order r and shape parameter $\boldsymbol{\eta}$

Instead of searching $\boldsymbol{\eta}$ and r separately one would prefer to search them together so that the most optimal values can be obtained, where $\boldsymbol{\eta} = (\eta_1, \dots, \eta_d)$ such that $\eta_l \neq \eta_k$, for $l \neq k$. We can combine the parameter maps used above to ensure that the kernel order values in $(1, \infty)$ and shape parameter $\boldsymbol{\eta}$ in $(0, \infty)^d$ as required in Section ??,

$$\boldsymbol{\theta}(\mathbf{t}) = \begin{pmatrix} r(t_1) \\ \eta(t_2) \\ \vdots \\ \eta(t_{d+1}) \end{pmatrix} = \begin{pmatrix} 1 + e^{t_1} \\ e^{t_2} \\ \vdots \\ e^{t_{d+1}} \end{pmatrix}, \quad \boldsymbol{\theta} : \mathbb{R}^{d+1} \rightarrow (1, \infty) \times (0, \infty)^d.$$

So instead of searching for $\boldsymbol{\theta}_{\text{opt}}$ in $(1, \infty) \times (0, \infty)^d$, one may search for the optimal

$$\mathbf{t} = \mathbf{t}^{-1}(\boldsymbol{\theta}) = \begin{pmatrix} \log(r - 1) \\ \log(\eta_1) \\ \vdots \\ \log(\eta_d) \end{pmatrix}$$

in the whole real line \mathbb{R}^{d+1} . The optimal value \mathbf{t}_{opt} can be transformed back to the desired interval $(1, \infty) \times (0, \infty)^d$ using

$$\boldsymbol{\theta}_{\text{opt}} = \begin{pmatrix} 1 + e^{t_{1,\text{opt}}} \\ e^{t_{2,\text{opt}}} \\ \vdots \\ e^{t_{d+1,\text{opt}}} \end{pmatrix}.$$

Similarly one can map the kernel order $q \in (0, 1)$ Section ??, and $\boldsymbol{\eta}$ in to a multivariate hyperparameter search.

13 Diagnostics for the Gaussian Process Assumption

The starting point for our Bayesian cubature is the assumption that the integrand arises from a Gaussian process. This means that the function data, \mathbf{f} , shall satisfy a multivariate Gaussian distribution, as in (7). We propose a transform, $\mathbf{Z} = \frac{1}{n} \mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H (\mathbf{f} - m\mathbf{1})$, to the given function data \mathbf{f} to verify the assumption. We can verify that the transformed data, \mathbf{Z} , has zero mean and is also uncorrelated because

$$\begin{aligned} \mathbb{E}[\mathbf{Z}] &= \frac{1}{n} \mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H (\mathbb{E}[\mathbf{f}] - m\mathbf{1}) \\ &= 0 \\ \text{cov}(\mathbf{Z}) &= \frac{1}{n^2} \mathbb{E} \left[\mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H (\mathbf{f} - m\mathbf{1}) (\mathbf{f} - m\mathbf{1})^T \mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \right] \\ &= \frac{1}{n^2} \mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \mathbb{E} \left[(\mathbf{f} - m\mathbf{1}) (\mathbf{f} - m\mathbf{1})^T \right] \mathbf{V} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n^2} \mathbf{V} \Lambda^{-\frac{1}{2}} \mathbf{V}^H \frac{1}{n} \mathbf{V} \Lambda \mathbf{V}^H \mathbf{V} \Lambda^{-\frac{1}{2}} \mathbf{V}^H \\
&= \frac{1}{n^3} \mathbf{V} \Lambda^{-\frac{1}{2}} (n) \Lambda(n) \Lambda^{-\frac{1}{2}} \mathbf{V}^H = \mathbf{I}
\end{aligned}$$

Thus, the elements of \mathbf{Z} are IID standard Gaussian random variables. In practice, using the estimated m_{EB} further simplifies,

$$\begin{aligned}
\mathbf{Z} &= \frac{1}{n} \mathbf{V} \Lambda^{-\frac{1}{2}} \mathbf{V}^H (\mathbf{f} - m \mathbf{1}) \\
&= \frac{1}{n} \mathbf{V} \Lambda^{-\frac{1}{2}} (\mathbf{V}^H \mathbf{f} - \frac{\tilde{f}_1}{n} \mathbf{V}^H \mathbf{1}) \\
&= \frac{1}{n} \mathbf{V} \Lambda^{-\frac{1}{2}} \left(\mathbf{V}^H \mathbf{f} - \tilde{f}_1 (1, 0, \dots, 0)^T \right)
\end{aligned}$$

To verify this hypothesis, we build a toy function using the known Bernoulli polynomial series,

$$f_{\theta,r}(x) = \hat{f}_0 + \sqrt{\theta} \sum_{k>0} \frac{f_{c,k} \cos(2\pi kx) + f_{s,k} \sin(2\pi kx)}{k^r} \quad (41)$$

where $f_0 \sim \text{IIDN}(c, b^2)$, $f_{c,k}, f_{s,k} \sim \text{IIDN}(0, a^2)$ and a, b, c are constants. The series in (41) is truncated to an arbitrarily chosen length N . It has zero mean and covariance

$$\begin{aligned}
\text{cov}(f) &= \mathbb{E}[f_{\theta,r}(x)f_{\theta,r}(t)] \\
&= \text{cov}(f_0) + \theta \sum_{k=1}^N \left(\mathbb{E}(f_{c,k}^2) \frac{\cos(2\pi kx) \cos(2\pi kt)}{k^{2r}} + \mathbb{E}(f_{s,k}^2) \frac{\sin(2\pi kx) \sin(2\pi kt)}{k^{2r}} \right) \\
&= b^2 + a^2 \theta \sum_{k=1}^N \frac{\cos(2\pi k(x-t))}{k^{2r}}
\end{aligned}$$

Figures 5 and 6 are normal probability plots of the Z_i using empirical Bayes estimates of m and θ . **more goes here.**

13.1 Posterior

We compute posterior distribution to verify the Gaussian assumption.

13.1.1 Non-informative prior

Assume the prior is of the form $\rho_{m,s^2}(\xi, \lambda) = \frac{1}{\lambda}$

$$\begin{aligned}
\rho_{f(x)}(\mathbf{y}) &= \int_0^\infty \int_{-\infty}^\infty \rho_f(\mathbf{y}|m = \xi, s^2 = \lambda) \rho_{m,s^2}(\xi, \lambda) d\xi d\lambda \\
&= \int_0^\infty \int_{-\infty}^\infty \frac{1}{(\sqrt{2\pi})^n \sqrt{\det(\lambda \mathbf{C})}} \exp\left(-\frac{1}{2\lambda} (\mathbf{y} - \xi \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - \xi \mathbf{1})\right) \frac{1}{\lambda} d\xi d\lambda \\
&\propto \frac{1}{(\sqrt{2\pi})^n} \int_0^\infty \frac{1}{\lambda^{\frac{n}{2}+1}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} (\mathbf{y} - \xi \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - \xi \mathbf{1})\right) d\xi d\lambda \\
&\propto \int_0^\infty \frac{1}{\lambda^{\frac{n}{2}+1}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} \underbrace{(\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1})}_D \left[\xi^2 - 2\xi \underbrace{\frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}}_A + \underbrace{\frac{\mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}}_B \right] \right) d\xi d\lambda
\end{aligned}$$

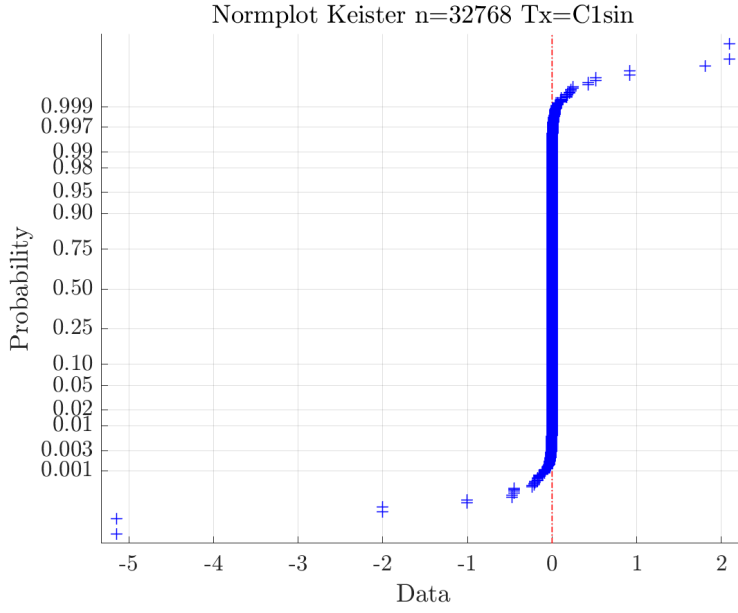


Fig. 5: Normal plot : Keister function with arbMean assumption

$$\begin{aligned}
& \propto \int_0^\infty \frac{1}{\lambda^{\frac{n}{2}+1}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} D \left[\xi^2 - 2\xi A + A^2 - A^2 + B \right]\right) d\xi d\lambda \\
& \propto \int_0^\infty \frac{1}{\lambda^{\frac{n}{2}+1}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} D \left[\xi^2 - 2\xi A + A^2 \right]\right) d\xi \exp\left(-\frac{1}{2\lambda} D[B - A^2]\right) d\lambda \\
& \propto \int_0^\infty \frac{1}{\lambda^{\frac{(n+1)}{2}}} \frac{1}{\sqrt{2\pi\lambda/D}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} D [\xi - A]^2\right) d\xi \exp\left(-\frac{1}{2\lambda} D[B - A^2]\right) d\lambda \\
& \propto \int_0^\infty \frac{1}{\lambda^{\frac{(n+1)}{2}}} \exp\left(-\frac{1}{2\lambda} D[B - A^2]\right) d\lambda \\
& \propto (D[B - A^2])^{-\frac{(n-1)}{2}} \\
& \propto \left(\mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} - \frac{(\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right)^{-\frac{(n-1)}{2}} \\
& = \left(\mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{(\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1})}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \right)^{-\frac{(n-1)}{2}} \\
& = \left((\mathbf{y} - m_{\text{EB}})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{EB}}) \right)^{-\frac{(n-1)}{2}}, \quad \text{where } m_{\text{EB}} = \frac{(\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y})}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}
\end{aligned}$$

13.1.2 TODO: How to interpret as Student's t -distribution?

13.1.3 Generic prior

Assume the prior is of the form $\rho_{m,s^2}(\xi, \lambda) = g(\lambda)$

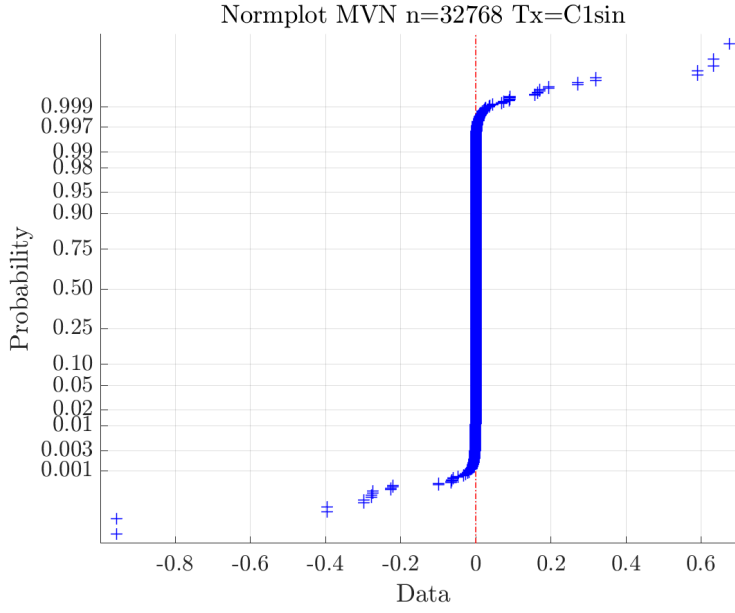


Fig. 6: Normal plot : MVN with arbMean assumption

$$\begin{aligned}
\rho_{f(\mathbf{x})}(\mathbf{y}) &= \int_0^\infty \int_{-\infty}^\infty \rho_f(\mathbf{y}|m = \xi, s^2 = \lambda) \rho_{m, s^2}(\xi, \lambda) d\xi d\lambda \\
&= \int_0^\infty \int_{-\infty}^\infty \frac{1}{(\sqrt{2\pi})^n \sqrt{\det(\lambda \mathbf{C})}} \exp\left(-\frac{1}{2\lambda} (\mathbf{y} - \xi \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - \xi \mathbf{1})\right) g(\lambda) d\xi d\lambda \\
&\propto \frac{1}{(\sqrt{2\pi})^n} \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{n}{2}}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} (\mathbf{y} - \xi \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - \xi \mathbf{1})\right) d\xi d\lambda \\
&\propto \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{n}{2}}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} \underbrace{(\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1})}_D \left[\xi^2 - 2\xi \underbrace{\frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}}_A + \underbrace{\frac{\mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}}_B \right] \right) d\xi d\lambda \\
&\propto \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{n}{2}}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} D \left[\xi^2 - 2\xi A + A^2 - A^2 + B \right] \right) d\xi d\lambda \\
&\propto \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{n}{2}}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} D \left[\xi^2 - 2\xi A + A^2 \right] \right) d\xi \exp\left(-\frac{1}{2\lambda} D[B - A^2] \right) d\lambda \\
&\propto \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{(n-1)}{2}}} \frac{1}{\sqrt{2\pi\lambda/D}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} D [\xi - A]^2 \right) d\xi \exp\left(-\frac{1}{2\lambda} D[B - A^2] \right) d\lambda \\
&\propto \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{(n-1)}{2}}} \exp\left(-\frac{1}{2\lambda} D[B - A^2] \right) d\lambda
\end{aligned}$$

This can be interpreted as Laplace transform of $g(\lambda)$. Let $\eta = \frac{1}{2}D[B - A^2]$ and $\lambda = \frac{1}{w}$, $d\lambda = -w^{-2}dw$

$$\begin{aligned}\rho_{f(\mathbf{x})}(\mathbf{y}) &\propto \int_0^\infty \frac{g(\lambda)}{\lambda^{\frac{(n-1)}{2}}} \exp\left(-\frac{1}{2\lambda}D[B - A^2]\right) d\lambda \\ &= \int_0^\infty \frac{g(1/w)}{w^{\frac{(n-1)}{2}}} \exp(-w\eta) (-w^{-2}) dw \\ &= \int_\infty^0 -g(1/w) w^{\frac{(n-1)}{2}-2} \exp(-w\eta) dw \\ &= \int_0^\infty g(1/w) w^{\frac{(n-5)}{2}} \exp(-w\eta) dw \\ &= \mathcal{LT}\{g(1/\eta)\}^{\left(\frac{n-5}{2}\right)'}\end{aligned}$$

where $\mathcal{LT}(\cdot)$ denotes the Laplace transform and $\left(\frac{n-5}{2}\right)'$ indicates the $\left(\frac{n-5}{2}\right)$ th derivative taken after the transform. Here we used frequency domain derivative property of the Laplace transform. $\rho_{f(\mathbf{x})}(\mathbf{y})$ is proportional to $\left(\frac{n-5}{2}\right)$ th derivative of the Laplace transform of $g(1/\eta)$.

$$\begin{aligned}D[B - A^2] &= \left(\mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} - \frac{(\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right) \\ &= \left((\mathbf{y} - m_{\text{EB}})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{EB}}) \right), \quad \text{where } m_{\text{EB}} = \frac{(\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y})}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}\end{aligned}$$

13.1.4 TODO: How to integrate further?

14 Numerical Experiments

JR: use uniformly randomly chosen ϵ instead 4 fixed

Fast Bayesian cubature algorithms developed in this research are demonstrated using three commonly used integration examples. These integrals are evaluated using our algorithm `cubBayesNet.g` and compared with `cubBayesLattice.g` which was developed in Jagadeeswaran and Hickernell [2019]. The `cubBayesLattice.g` is a fast automatic Bayesian cubature algorithm which uses rank-1 lattice nodes. The first example shows evaluating a multivariate Gaussian probability given the interval. The second example shows integrating the Keister's function, and the final example shows computing an Asian arithmetic option pricing.

14.1 Testing Methodology

Four hundred different error tolerances, ϵ , were randomly chosen from a fixed interval for each example. The intervals for error tolerance were chosen depending on the difficulty of the problem. The nodes used in `cubBayesNet.g` were the randomly scrambled and shifted Sobol' points supplied by MATLAB's Sobol' sequence generator.

For each integral example, and each stopping criteria—empirical Bayes, full Bayes, and generalized cross-validation—our algorithm is run with each randomly chosen error tolerance as mentioned above. For each test, the execution time is

plotted against $|\mu - \hat{\mu}|/\varepsilon$. We expect $|\mu - \hat{\mu}|/\varepsilon$ to be no greater than one, but hope that it is not too much smaller than one, which would indicate a stopping criterion that is too conservative.

Periodization variable transforms are not used in the examples with `cubBayesNet.g`. Since the `cubBayesNet.g` does not need this requirement, the integrands are used directly.

14.2 Multivariate Gaussian Probability

This example is introduced in Section ??, where we use the Matérn covariance kernel. We reuse f_{Genz} (??) and apply a periodization transform to obtain f_{GenzP} when required.

Here we use f_{Genz} (??) without any periodization, and chose $d = 3$ and $r = 1$. The simulation results for this example integrand are summarized in Figures 7, 8, and 9. In all cases, `cubBayesNet.g` returns an approximation within the prescribed error tolerance. We used the same setting as before with generic slow Bayesian cubature in Section ?? for comparison. For error threshold $\varepsilon = 10^{-5}$ with empirical stopping criterion, our fast algorithm takes about 2 seconds as shown in Figure ?? whereas the basic algorithm takes 30 seconds as shown in Figure ?. `cubBayesNet.g` uses fast Walsh transform which is slower in MATLAB due to the way it was implemented. This is reason it takes more longer the `cubBayesLattice.g`. But comparing the number of samples, n , used for integration provides more insight which directly relates to algorithm's computational cost. The `cubBayesLattice.g` used $n = 16384$ samples whereas `cubBayesNet.g` used $n = 32768$ samples even with $r = 1$ order kernel.

Amongst the three stopping criteria, GCV achieved the results faster than others but it is less conservative. One can also observe from the figures that the credible intervals are narrower than in Figure ?. This shows that `cubBayesNet.g` with $r = 1$ kernel more accurately approximates the integrand.

14.3 Keister's Example

This multidimensional integral function comes from Keister [1996] and is inspired by a physics application:

$$\begin{aligned}\mu &= \int_{\mathbb{R}^d} \cos(\|\mathbf{t}\|) \exp(-\|\mathbf{t}\|^2) d\mathbf{t} \\ &= \int_{[0,1]^d} f_{\text{Keister}}(\mathbf{x}) d\mathbf{x},\end{aligned}\tag{42}$$

where

$$f_{\text{Keister}}(\mathbf{x}) = \pi^{d/2} \cos\left(\left\|\Phi^{-1}(\mathbf{x})/2\right\|\right),$$

and Φ is the standard normal distribution. The true value of μ can be calculated iteratively in terms of a quadrature as follows:

$$\mu = \frac{2\pi^{d/2} I_c(d)}{\Gamma(d/2)}, \quad d = 1, 2, \dots$$

where Γ denotes the gamma function, and

$$I_c(1) = \frac{\sqrt{\pi}}{2 \exp(1/4)},$$

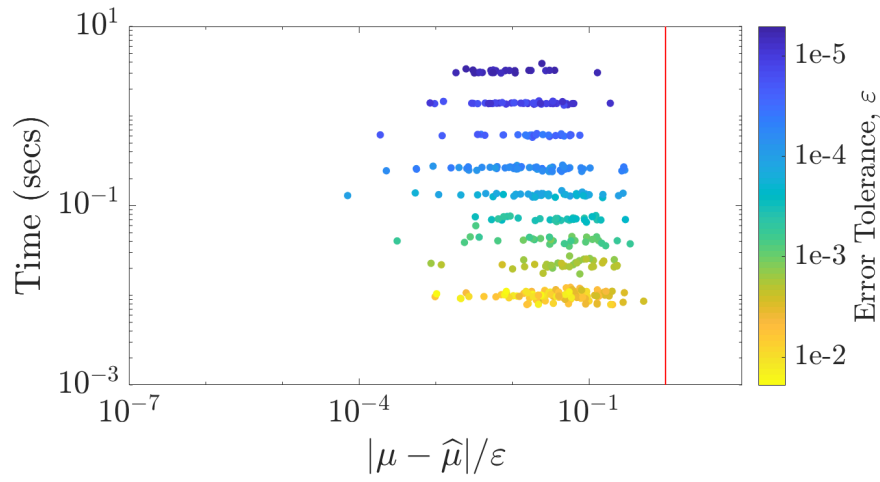


Fig. 7: Multivariate normal probability example with empirical Bayes stopping criterion.

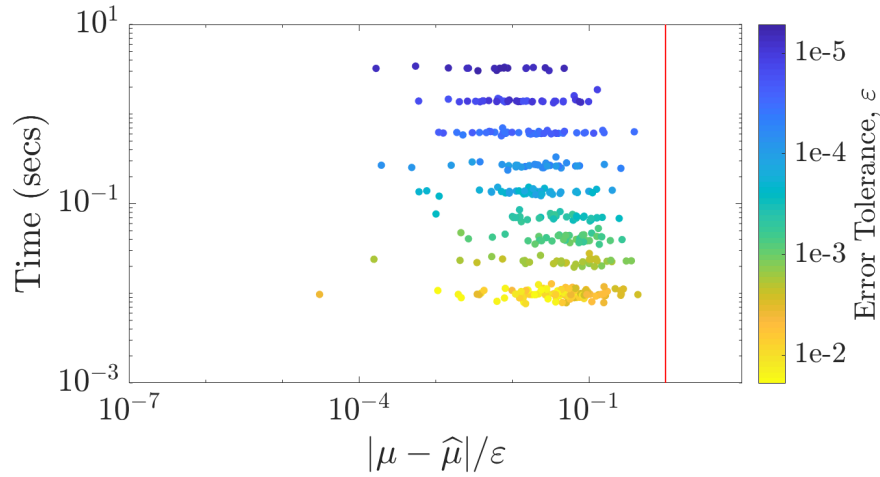


Fig. 8: Multivariate normal probability example with the full-Bayes stopping criterion.

$$\begin{aligned}
 I_s(1) &= \int_{x=0}^{\infty} \exp(-\mathbf{x}^T \mathbf{x}) \sin(\mathbf{x}) \, d\mathbf{x} \\
 &= 0.4244363835020225, \\
 I_c(2) &= \frac{1 - I_s(1)}{2}, \quad I_s(2) = \frac{I_c(1)}{2} \\
 I_c(j) &= \frac{(j-2)I_c(j-2) - I_s(j-1)}{2}, \quad j = 3, 4, \dots
 \end{aligned}$$

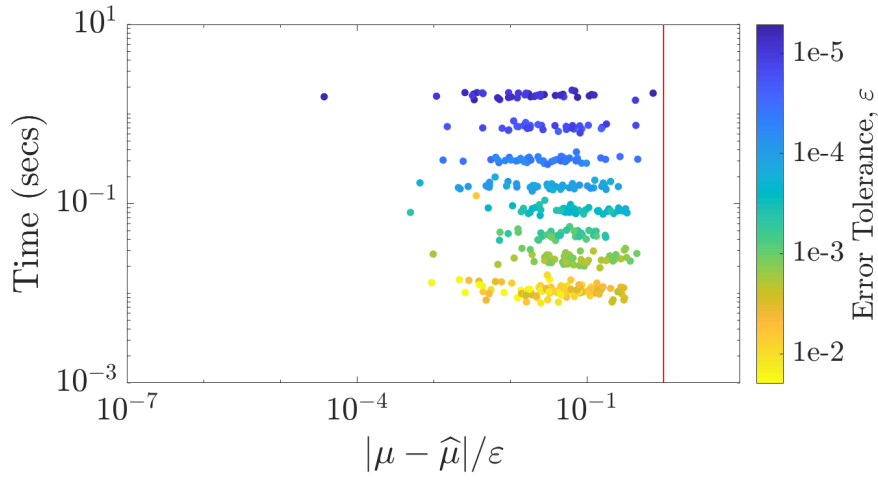


Fig. 9: Multivariate normal probability example with the GCV stopping criterion.

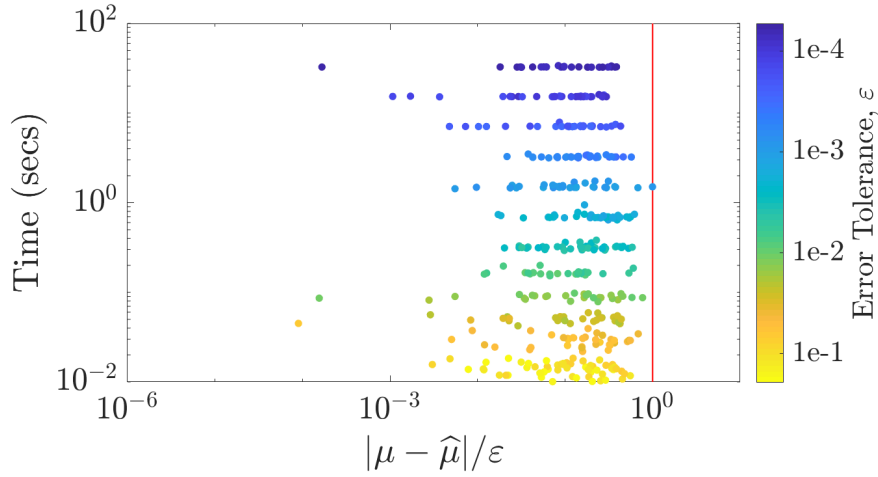


Fig. 10: Keister example using the empirical Bayes stopping criterion.

$$I_s(j) = \frac{(j-2)I_s(j-2) - I_c(j-1)}{2}, \quad j = 3, 4, \dots$$

Figures 10, 11 and 12 summarize the numerical tests for this case. We used dimension $d = 4$, and $r = 1$. No periodization transform was used as the integrand need not be periodic. In this example, we use $r = 1$ order kernel whereas in Section ??, $r = 2$ kernel was used. This necessitates `cubBayesNet_g` to use more samples for integration. As observed from the figures, the GCV stopping criterion achieved the results faster than the others but it is less conservative which is also the case with the multivariate Gaussian example.

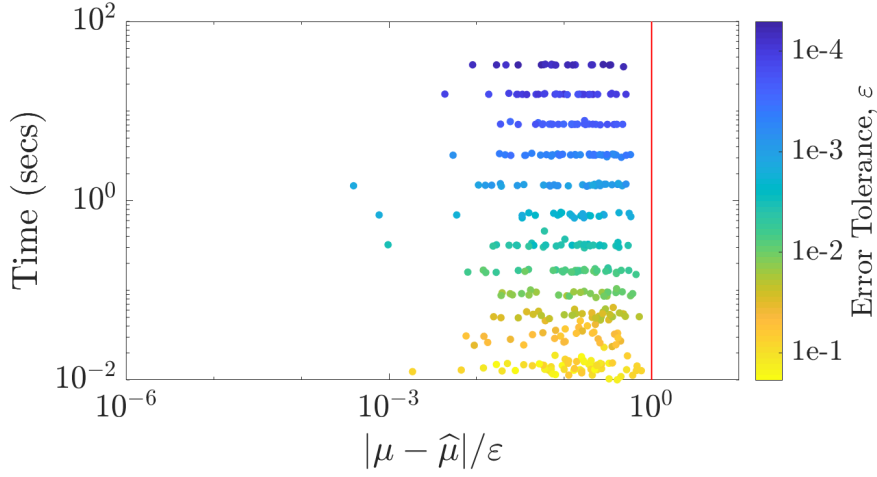


Fig. 11: Keister example using the full-Bayes stopping criterion.

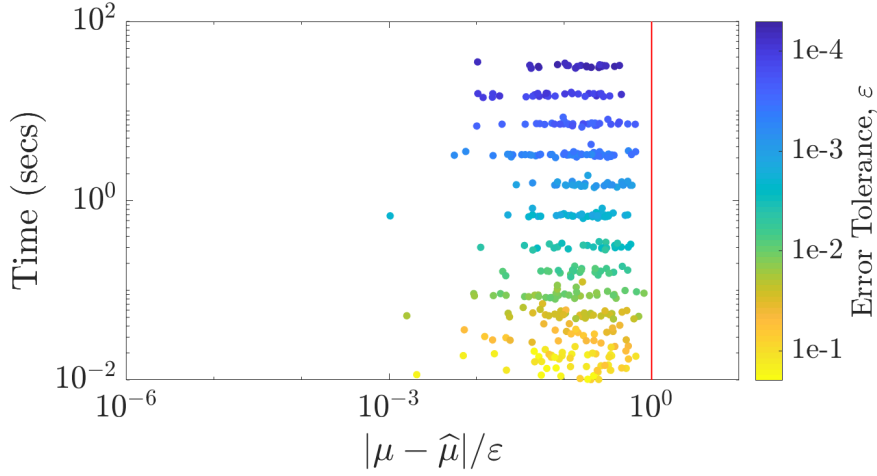


Fig. 12: Keister example using the GCV stopping criterion.

14.4 Option Pricing

The price of financial derivatives can often be modeled by high dimensional integrals. If the underlying asset is described in terms of a discretized geometric Brownian motion, then the fair price of the option is:

$$\mu = \int_{\mathbb{R}^d} \text{payoff}(\mathbf{z}) \frac{\exp(\frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z})}{\sqrt{(2\pi)^d \det(\Sigma)}} d\mathbf{z} = \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x},$$

where $\text{payoff}(\cdot)$ defines the discounted payoff of the option,

$$\Sigma = (T/d) (\min(j, k))_{j,k=1}^d = \mathbf{L} \mathbf{L}^T,$$

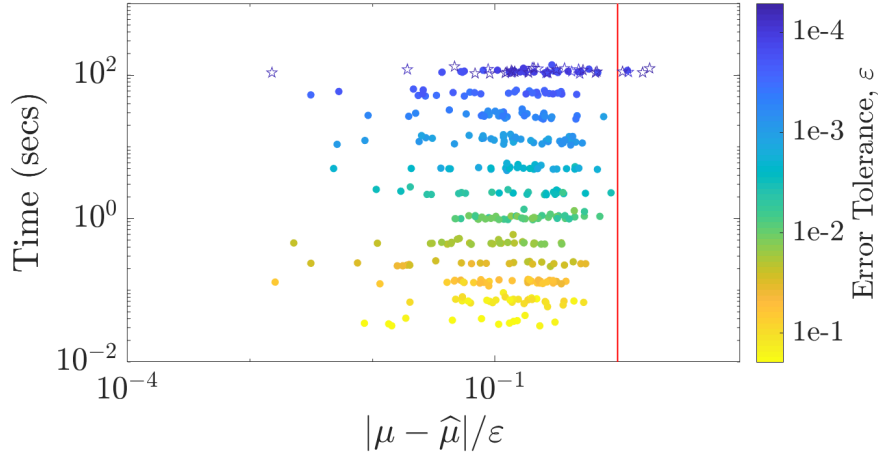


Fig. 13: Option pricing using the empirical Bayes stopping criterion. The hollow stars indicate the algorithm has not met the error threshold ϵ even with using maximum n .

$$f(\mathbf{x}) = \text{payoff} \left(\mathbf{L} \begin{pmatrix} \Phi^{-1}(x_1) \\ \vdots \\ \Phi^{-1}(x_d) \end{pmatrix} \right).$$

The Asian arithmetic mean call option has a payoff of the form

$$\text{payoff}(z) = \max \left(\frac{1}{d} \sum_{j=1}^d S_j(z) - K, 0 \right) e^{-rT},$$

$$S_j(z) = S_0 \exp \left((r - \sigma^2/2)jT/d + \sigma \sqrt{T/d} z_j \right).$$

Here, T denotes the time to maturity of the option, d the number of time steps, S_0 the initial price of the stock, r the interest rate, σ the volatility, and K the strike price.

The Figures 13, 14 and 15 summarize the numerical results for the option pricing example using the same values for, T , d , S_0 , r , σ , K , as in Section ?? . As mentioned before, this integrand has a kink caused by the max function, so, `cubBayesNet_g` could be more efficient than `cubBayesLattice_g`, as no periodization transform is required. This can be observed from the number of samples used for integration to meet the same error threshold. For the error tolerance, $\epsilon = 10^{-3}$, `cubBayesLattice_g` used $n = 2^{20}$ samples, whereas `cubBayesNet_g` used $n = 2^{17}$ samples.

14.5 Discussion

JR: Move this to end of chapter

As shown in Figures ?? to 15, both the algorithms computed the integral within user specified threshold most of the time except on a few occasions. This

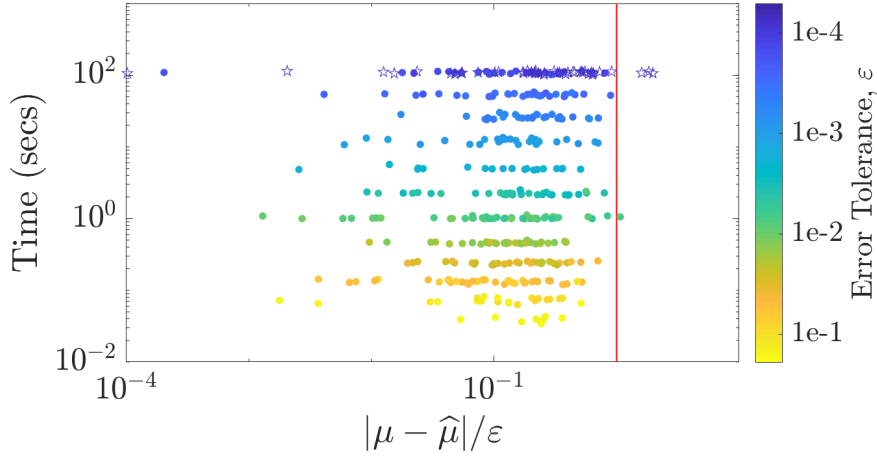


Fig. 14: Option pricing using the full-Bayes stopping criterion. The hollow stars indicate the algorithm has not met the error threshold ϵ even with using maximum n .

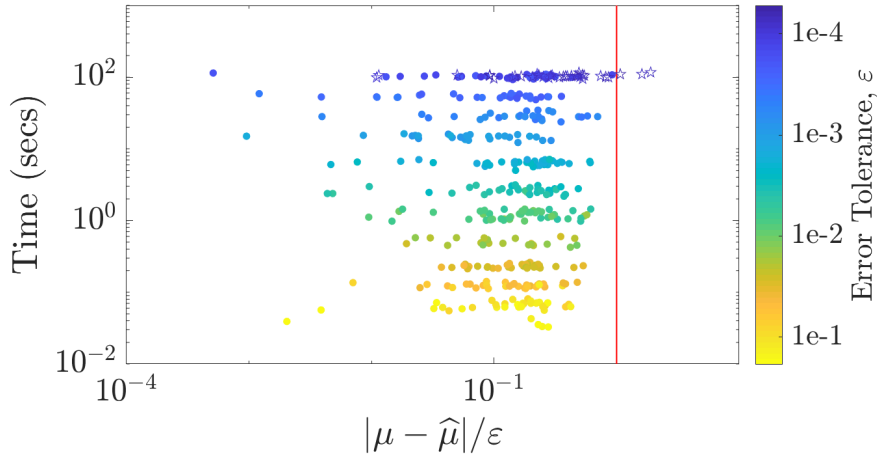


Fig. 15: Option pricing using the GCV stopping criterion. The hollow stars indicate the algorithm has not met the error threshold ϵ even with using maximum n .

is especially the case with option pricing example due to the complexity and high dimension of the integrand. Also notice that the `cubBayesLattice_g` algorithm finished within 10 seconds for Keister and multivariate Gaussian. Option pricing took closer to 70 seconds due to the complexity of the integrand.

Another noticeable aspect from the plots of `cubBayesLattice_g` is how much the error bounds differ from the true error. For option pricing example, the error bound is not as conservative as it is for the multivariate Gaussian and Keister

examples. A possible reason is that the latter integrands are significantly smoother than the covariance kernel. This is a matter for further investigation.

Most noticeable aspect from the plots of `cubBayesNet.g` is how closer the error bounds are to the true error. This shows that the `cubBayesNet.g`'s estimation of expected error in the stopping criterion is very accurate. Similar to `cubBayesLattice.g`, it missed meeting the given error threshold for the option pricing example, as marked by the hollow stars, for $\varepsilon = 10^{-4}$. The algorithm reached max allowed number of samples, $n = 2^{20}$ due to the complexity of the integrand.

14.6 Comparison with `cubMC.g`, `cubLattice.g` and `cubSobol.g`

GAIL library provides variety of numerical integration algorithms based on different theoretical foundations, We would like to compare how our algorithms perform relatively to these. We consider three GAIL algorithms 1) `cubMC.g`, a simple Monte-Carlo method for multi-dimensional integration, 2) `cubLattice.g`, a quasi-Monte-Carlo method using Lattice points, and 3) `cubSobol.g`, a quasi-Monte-Carlo method using Sobol points.

Keister integral The Table 1 summarizes the performance of the methods MC, Lattice, Sobol, BayesLat, and BayesSob—which refer to the GAIL cubatures, `cubMC.g`, `cubLattice.g`, `cubSobol.g`, `cubBayesLattice.g`, `cubBayesNet.g`, respectively for estimating Keister integral defined in (42). We conducted two simulations with $d = 3$ and 8. In the case of $d = 3$, all five methods succeeded completely, meaning, the absolute error is less than given tolerance, i.e., $|\mu - \hat{\mu}| \leq \varepsilon$, where $\hat{\mu}$ is a cubature's approximated value. The fastest method was `cubBayesLattice.g`. In the case of $d = 8$, `cubSobol.g` achieved 100% success rate and was the fastest. But `cubBayesLattice.g` was competitive and had the smallest average absolute error. `cubBayesNet.g` used lowest number of samples but was slower than `cubSobol.g`.

JR: avoid Exp notation for 2 decimals

Table 1: Comparison of average performance of cubatures for estimating the Keister integral (42) for 1000 independent runs. These results can be conditionally reproduced with the script, `KeisterCubatureExampleBayes.m`, in GAIL.

$d = 3, \varepsilon = 0.005$					
Method	MC	Lattice	Sobol	BayesLat	BayesSobol
Absolute Error	0.001 100	0.000 510	0.000 520	0.000 430	0.000 560
Tolerance Met	100%	100%	100%	100%	100%
n	2 500 000	4100	3900	1000	1900
Time (seconds)	0.1800	0.0069	0.0054	0.0029	0.0700
$d = 8, \varepsilon = 0.050$					
Method	MC	Lattice	Sobol	BayesLat	BayesSobol
Absolute Error	0.012 000	0.015 000	0.007 300	0.001 800	0.008 300
Tolerance Met	100%	99%	100%	100%	100%
n	7 400 000	15 000	16 000	66 000	8200
Time (seconds)	1.2000	0.0220	0.0160	0.2100	0.3500

Multivariate Gaussian Table 2 summarizes the performance of the methods MC, Lattice, Sobol, BayesLat, and BayesSob for estimating the multi-dimensional Gaussian probability $\mathbf{X} \sim \mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. This experiment demonstrates our algorithm's ability to handle high-dimensional integral.

We conducted two simulations with different $\boldsymbol{\Sigma}$ and estimation intervals (\mathbf{a}, \mathbf{b}) but fixed $\boldsymbol{\mu} = \mathbf{0}$ and required error threshold, $\varepsilon = 10^{-3}$. In the first case, all five methods succeeded completely. The fastest method was `cubBayesLattice.g` but `cubBayesNet.g` used the lowest number of samples. In the second case also, all five methods succeeded, but `cubLattice.g` was the fastest. The `cubBayesNet.g` was competitive and had the smallest average absolute error using lowest number of samples. The `cubBayesLattice.g` achieved the next lowest average error but was slower than `cubSobol.g`.

Table 2: Comparison of average performance of cubatures for estimating the $d = 20$ Multivariate Normal (??) for 1000 independent runs with $\varepsilon = 10^{-3}$. These results can be conditionally reproduced with the script, `MVNCubatureExampleBayes.m`, in GAIL.

$\Sigma = \mathbf{I}_d, \mathbf{b} = -\mathbf{a} = (3.5, \dots, 3.5)$						
Method	MC	Lattice	Sobol	BayesLat	BayesSobol	
Absolute Error	2.20×10^{-16}	2.70×10^{-14}	2.70×10^{-14}	2.20×10^{-16}	2.20×10^{-16}	
Tolerance Met	100%	100%	100%	100%	100%	
n	10 000	1000	1000	1000	260	
Time (seconds)	0.0410	0.0820	0.0710	0.0650	0.0790	
$\Sigma = 0.4 \mathbf{I}_d + 0.6 \mathbf{11}^T, \mathbf{a} = (-\infty, \dots, -\infty), \mathbf{b} = \sqrt{d}(U_1, \dots, U_d)$						
Method	MC	Lattice	Sobol	BayesLat	BayesSobol	
Absolute Error	2.30×10^{-4}	2.10×10^{-4}	4.40×10^{-4}	1.00×10^{-4}	4.80×10^{-5}	
Tolerance Met	100%	100%	100%	100%	100%	
n	10 000	1000	1000	1000	260	
Time (seconds)	0.0350	0.0120	0.0140	0.0150	0.0300	

Shape Parameter Fine-tuning **JR: Numerical examples for the case of shape parameter per dimension**

Allowing the kernel shape parameter to vary for each dimension could improve the accuracy of numerical integration when the integrand under consideration has only very low effective dimension as in the Option Pricing example we demonstrated. We demonstrate this advantage by integrating a function that is not symmetric across dimensions,

$$f(\mathbf{x}) = \sum_{j=1}^d v_j \sin(2\pi x_j^2) \quad (43)$$

which has known integral

$$\int_{[0,1]^d} f(\mathbf{x}) = \frac{1}{2} \mathbf{fresnels}(d) \sum_{j=1}^d v_j$$

where $\mathbf{fresnels}$ is the Fresnel Sine integral,

$$\mathbf{fresnels}(z) = \int_0^z \sin\left(\frac{\pi t^2}{2}\right) dt.$$

Table 3: Comparison of average performance of Bayesian Cubature with common shape parameter vs dimension specific shape parameter for estimating the $d = 3$ Fresnel Sine integral. These results can be conditionally reproduced with the script, `demoMultiTheta.m`, in GAIL.

Method	Fresnel Sine Integral in $d = 3$	
	OneTheta	MultiTheta
Absolute Error	0.000 23	0.063 00
n	4100	260
Time (seconds)	0.0270	0.0230

The results are summarized from the two different approaches in Table 3. The first method, called **OneTheta**, uses common shape parameter across all the dimensions, whereas the second method, called **MultiTheta**, allows the shape parameters to vary across the dimensions. In the **MultiTheta** method, the shape parameter search is multivariate, so the magnitude of shape parameter depends on the integrand’s magnitude in each dimension. We have chosen an integrand particularly to demonstrate this aspect (43) where we used $d = 3$ and the constants $\mathbf{v} = (10^{-4}, 1, 10^4)$. The choice of magnitude variations in constants \mathbf{v} allows to make the integrand varies significantly across dimensions.

We ran this test for 1000 times. In comparison, both the methods successfully computed the integral all the time but **MultiTheta** was slightly faster. The **MultiTheta** method used less number of samples but the integration error was bigger than the **OneTheta**. For the same number of samples, the **OneTheta** method will be much faster since the shape parameter search is faster. The **MultiTheta** method is useful in scenarios where we want to use smaller size, n , and the integrand varies significantly across dimensions.

15 Conclusion and future work

TBD

References

- J. F. Baldeaux. *Higher order nets and sequences*. PhD thesis, The School of Mathematics and Statistics at The University of New South Wales, June 2010.
- M. Beckers and A. Haegemans. Transformation of integrands for lattice rules. In T. O. Espelid and A. C. Genz, editors, *Numerical Integration: Recent Developments, Software and Applications*, pages 329–340. Kluwer Academic Publishers, Dordrecht, 1992.
- P. Bratley and B. L. Fox. Algorithm 659: Implementing Sobol’s quasirandom sequence generator. *ACM Trans. Math. Software*, 14:88–100, 1988.
- R. Brent. *Algorithms for Minimization Without Derivatives*. Prentice-Hall, 1973.
- F.-X. Briol, C. J. Oates, M. Girolami, M. A. Osborne, and D. Sejdinovic. Probabilistic integration: A role in statistical computation? *Statist. Sci.*, 2019. to appear.
- S.-C. T. Choi, Y. Ding, F. J. Hickernell, L. Jiang, Ll. A. Jiménez Rugama, D. Li, R. Jagadeeswaran, X. Tong, K. Zhang, Y. Zhang, and X. Zhou. GAIL: Guaranteed Automatic Integration Library (versions 1.0–2.2). MATLAB software, 2013–2017. URL http://gailgithub.github.io/GAIL_Dev/.
- R. Cools and D. Nuyens, editors. *Monte Carlo and Quasi-Monte Carlo Methods: MCQMC, Leuven, Belgium, April 2014*, volume 163 of *Springer Proceedings in Mathematics and Statistics*, 2016. Springer-Verlag, Berlin.
- L. L. Cristea, J. Dick, G. Leobacher, and F. Pillichshammer. The tent transformation can improve the convergence rate of quasi-Monte Carlo algorithms using digital nets. *Numer. Math.*, 105:413–455, 2007.
- P. Diaconis. Bayesian numerical analysis. In S. S. Gupta and J. O. Berger, editors, *Statistical Decision Theory and Related Topics IV, Papers from the 4th Purdue Symp., West Lafayette, Indiana 1986*, volume 1, pages 163–175. Springer-Verlag, New York, 1988.
- J. Dick. Walsh spaces containing smooth functions an quasi-Monte Carlo rules of arbitrary high order. *SIAM J. Numer. Anal.*, 46(1519–1553), 2008.
- K. Dong, D. Eriksson, H. Nickisch, D. Bindel, and A. G. Wilson. Scalable log determinants for gaussian process kernel learning. *NIPS*, 2017. in press.
- G. Forsythe, M. Malcolm, and C. Moler. *Computer methods for mathematical computations*. Prentice-Hall, 1976.
- P. Glasserman. *Monte Carlo Methods in Financial Engineering*, volume 53 of *Applications of Mathematics*. Springer-Verlag, New York, 2004.
- I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- F. J. Hickernell. The trio identity for quasi-Monte Carlo error analysis. In P. Glynn and A. Owen, editors, *Monte Carlo and Quasi-Monte Carlo Methods: MCQMC, Stanford, USA, August 2016*, Springer Proceedings in Mathematics and Statistics, pages 3–27. Springer-Verlag, Berlin, 2018. doi: 10.1007/978-3-319-91436-7.

- F. J. Hickernell and Ll. A. Jiménez Rugama. Reliable adaptive cubature using digital sequences. In Cools and Nuyens [2016], pages 367–383. arXiv:1410.8615 [math.NA].
- F. J. Hickernell and R. X. Yue. The mean square discrepancy of scrambled (t, s) -sequences. *SIAM J. Numer. Anal.*, 38:1089–1112, 2000. doi: 10.1137/S0036142999358019.
- F. J. Hickernell, Ll. A. Jiménez Rugama, and D. Li. Adaptive quasi-Monte Carlo methods for cubature. In J. Dick, F. Y. Kuo, and H. Woźniakowski, editors, *Contemporary Computational Mathematics — a celebration of the 80th birthday of Ian Sloan*, pages 597–619. Springer-Verlag, 2018. doi: 10.1007/978-3-319-72456-0.
- N. J. Higham. *Functions of matrices: theory and computation*. SIAM, 2008.
- H. S. Hong and F. J. Hickernell. Algorithm 823: Implementing scrambled digital nets. *ACM Trans. Math. Software*, 29:95–109, 2003. doi: 10.1145/779359.779360.
- R. Jagadeeswaran and F. J. Hickernell. Fast automatic Bayesian cubature using lattice sampling. *Statist. Comp.*, 29:1215–1229, 2019. doi: 10.1007/s11222-019-09895-9.
- Ll. A. Jiménez Rugama and F. J. Hickernell. Adaptive multidimensional integration based on rank-1 lattices. In Cools and Nuyens [2016], pages 407–422. arXiv:1411.1966.
- B. D. Keister. Multidimensional quadrature algorithms. *Computers in Physics*, 10:119–122, 1996. doi: 10.1063/1.168565.
- A. Keller. Quasi-Monte Carlo image synthesis in a nutshell. In J. Dick, F. Y. Kuo, G. W. Peters, and I. H. Sloan, editors, *Monte Carlo and Quasi-Monte Carlo Methods 2012*, volume 65 of *Springer Proceedings in Mathematics and Statistics*, pages 213–249. Springer Berlin Heidelberg, 2013.
- 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. *Foundations of Computational Mathematics*, 16(6):1631–1696, 2016.
- D. Laurie. Periodizing transformations for numerical integration. *J. Comput. Appl. Math.*, 66:337–344, 1996.
- J. Matoušek. On the L_2 -discrepancy for anchored boxes. *J. Complexity*, 14:527–556, 1998.
- H. Niederreiter. Constructions of (t, m, s) -nets and (t, s) -sequences. *Finite Fields Appl.*, 11:578–600, 2005.
- D. Nuyens. URL <https://people.cs.kuleuven.be/~dirk.nuyens/>.
- D. Nuyens. The construction of good lattice rules and polynomial lattice rules. Aug 2013.
- A. O’Hagan. Bayes-Hermite quadrature. *J. Statist. Plann. Inference*, 29:245–260, 1991. doi: 10.1016/0378-3758(91)90002-V.
- A. OHagan. Bayes-hermite quadrature. *Journal of Statistical Planning and Inference*, 29(3):245260, 1991.
- A. B. Owen. Randomly permuted (t, m, s) -nets and (t, s) -sequences. In H. Niederreiter and P. J.-S. Shiue, editors, *Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing*, volume 106 of *Lecture Notes in Statistics*, pages 299–317. Springer-Verlag, New York, 1995.
- C. E. Rasmussen and C. Williams. Bayesian Monte Carlo. In S. Thrun, L. K. Saul, and K. Obermayer, editors, *Advances in Neural Information Processing Systems*, volume 15, pages 489 – 496. MIT Press, 2003.

- C. E. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, Massachusetts, 2006. (online version available at <http://www.gaussianprocess.org/gpml/>).
- K. Ritter. *Average-Case Analysis of Numerical Problems*, volume 1733 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 2000.
- A. Sidi. A new variable transformation for numerical integration. In H. Brass and G. Hämmerlin, editors, *Numerical Integration IV*, number 112 in International Series of Numerical Mathematics, pages 359–373. Birkhäuser, Basel, 1993.
- A. Sidi. Further extension of a class of periodizing variable transformations for numerical integration. *J. Comput. Appl. Math.*, 221:132–149, 2008.
- I. M. Sobol'. The distribution of points in a cube and the approximate evaluation of integrals. *U.S.S.R. Comput. Math. and Math. Phys.*, 7:86–112, 1967.
- I. M. Sobol'. Uniformly distributed sequences with an additional uniformity property. *Zh. Vychisl. Mat. i Mat. Fiz.*, 16:1332–1337, 1976.