

Fast Automatic Bayesian Cubature Using Lattice Sampling

R. Jagadeeswaran · Fred. J. Hickernell

Received: date / Accepted: date

Abstract Automatic cubatures provide approximations to multidimensional integrals that satisfy user-specified error tolerances. For multidimensional problems, the sampling density is fixed, but the sample size, n , is determined automatically. Bayesian cubature postulates that the integrand is an instance of a stochastic process. Prior information about mean and covariance of this process is used to form data-driven error bounds. However, the process of inferring the mean and covariance governing the stochastic process from n integrand values involves computing matrix inverses and determinants, which are in general time-consuming $O(n^3)$ operations. Our work employs low discrepancy data sites and matching kernels that lower the computational cost to $O(n \log n)$. The confidence interval for the Bayesian posterior error is used to choose n automatically to satisfy the user-defined error tolerance. This approach is demonstrated using rank-1 lattice sequences and shift-invariant kernels.

Keywords Bayesian cubature · Probabilistic numeric methods · GAIL

1 Introduction

Cubature is the problem of inferring a numerical value for the integral $\mu = \int_{\mathbb{R}^d} g(\mathbf{x}) d\mathbf{x}$, of a multi-dimensional function g where μ has no closed form analytic expression. Typically, g is only accessible in the form of a

black-box function routine. Cubature is a key component of many problems in scientific computing, finance, statistical modeling, and machine learning.

The integral μ may often be expressed in the form
$$\mu := \mu(f) := \mathbb{E}[f(\mathbf{X})] = \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x}, \quad (1)$$
 where $f : [0,1]^d \rightarrow \mathbb{R}$ is the integrand, and $\mathbf{X} \sim \mathcal{U}[0,1]^d$. The problem of transforming the original integral into the above form is not addressed here. The cubature algorithm may take the form

$$\hat{\mu} := \hat{\mu}(f) := w_0 + \sum_{i=1}^n f(\mathbf{x}_i) w_i, \quad (2)$$

where the weights, w_0 , and $\mathbf{w} = (w_i)_{i=1}^n \in \mathbb{R}^n$, and the nodes, $\{\mathbf{x}_i\}_{i=1}^n \subset [0,1]^d$, are chosen to make the error, $|\mu - \hat{\mu}|$, small. In this article the nodes is assumed to be deterministic.

Our concern is constructing 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, i.e.,

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

Rather than relying on strong assumptions about the integrand, such as its variance or total variation, we construct a stopping criterion that is data-driven, and based on a credible interval arising from a Bayesian approach to the problem. We build upon the ideas of Diaconis[3], O’Hagan[11], Ritter[14], Rasmussen and Ghahramani[12], Briol et al.[1], and others. Thus, our algorithm is an example of *probabilistic numerics*.

Our primary contribution here is to demonstrate how the choice of a family of kernels that match the low discrepancy sampling nodes facilitates fast computation of the data-driven stopping criterion. If n function values, at a cost of $\$(f)$ each, are obtained for cubature purposes, then $\mathcal{O}(n \log(n))$ additional oper-

F. Author
first address
Tel.: +123-45-678910
Fax: +123-45-678910
E-mail: fauthor@example.com

S. Author
second address

ations are required to check whether the error tolerance is satisfied. The total cost of our algorithm is then $\mathcal{O}(\$(f)n + n \log(n))$. This is significantly fewer operations than the $\mathcal{O}(n^3)$ typically required for Bayesian cubature.

Hickernell [7] compares and contrasts 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[8,9] have developed stopping criteria for cubature rules based on low discrepancy nodes by tracking the discrete Fourier coefficients of the integrand. The algorithm proposed here also depends on the discrete Fourier coefficients of the integrand, but in a different way.

Section 2 explains the Bayesian approach to estimate the posterior cubature error and defines our automatic Bayesian cubature. Although much of this material is not new, it is included for completeness. We end Section 2 by demonstrating why Bayesian cubature is typically quite computationally expensive. Section 3 introduces the concept of covariance kernels that match well the nodes and expedites the computations required by our automatic Bayesian cubature. Section 4 implements this idea for shift invariant kernels and rank-1 lattice nodes. Section 5 covers further enhancements to the algorithm to avoid cancellation error and make it faster. Finally, numerical examples are shown using the faster algorithm developed.

2 Bayesian cubature

2.1 Bayesian posterior error

Suppose that the integrand, f , is drawn from a Gaussian process, i.e., $f \sim \mathcal{GP}(m, s^2 C_\theta)$. Specifically, f has real-valued mean m and covariance function $s^2 C_\theta$, where s is a non-negative scale factor $s C_\theta : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$ is a symmetric, positive-definite function and parameterized by θ :

$$C^T = C, \quad \mathbf{a}^T C \mathbf{a} > 0 \quad \forall \mathbf{a} \neq \mathbf{0},$$

where $C = (C_\theta(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n$,

$$\forall n \in \mathbb{N}, \mathbf{t}, \mathbf{x}, \mathbf{x}_1, \dots, \mathbf{x}_n \in [0, 1]^d. \quad (4)$$

The covariance function, C , and the Gram matrix, C depend implicitly on θ , but sometimes this is omitted in the notation for simplicity's sake.

For a Gaussian process, all vectors of linear functionals of f have a multivariate Gaussian distribution. Defining $\mathbf{f} := (f(\mathbf{x}_i))_{i=1}^n$ as the multivariate normal vector of function values, it follows that

$$\mathbf{f} \sim \mathcal{N}(m\mathbf{1}, s^2 C), \quad (5a)$$

$$\mu \sim \mathcal{N}(m, s^2 c_0), \quad (5b)$$

$$\text{where } c_0 = \int_{[0,1]^{2d}} C_\theta(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}, \quad (5c)$$

$$\text{cov}(\mathbf{f}, \mu) = \left(\int_{[0,1]^d} C(\mathbf{t}, \mathbf{x}_i) d\mathbf{t} \right)_{i=1}^n =: \mathbf{c}. \quad (5d)$$

Therefore, it follows from [13] Appendix A.2 Gaussian Identities that the *conditional* distribution of the integral given observed function values, $\mathbf{f} = \mathbf{y}$ is also Gaussian:

$$\mu | (\mathbf{f} = \mathbf{y}) \sim \mathcal{N}(m(1 - \mathbf{c}^T C^{-1} \mathbf{1}) + \mathbf{c}^T C^{-1} \mathbf{y}, s^2(c_0 - \mathbf{c}^T C^{-1} \mathbf{c})). \quad (6)$$

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

$$\hat{\mu} | (\mathbf{f} = \mathbf{y}) = m(1 - \mathbf{1}^T C^{-1} \mathbf{c}) + \mathbf{c}^T C^{-1} \mathbf{y}. \quad (7)$$

Under this definition, the cubature error has zero mean and a variance depending on the choice of nodes:

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

A credible interval for the integral is given given by

$$\mathbb{P}_f \left[|\mu - \hat{\mu}| \leq 2.58s \sqrt{c_0 - \mathbf{c}^T C^{-1} \mathbf{c}} \right] = 99\%. \quad (8)$$

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

2.2 Parameter estimation

The credible interval in (8) suggests how our automatic Bayesian cubature proceeds. Function data is accumulated until we can be confident that the error is no greater than the error tolerance. However, the credible interval still depends on the parameters m, s , and θ . These should not be assumed a priori but be inferred from the data. This section describes three approaches.

2.2.1 Empirical Bayes

One approach to estimate the parameters is via maximum likelihood estimation (MLE). The log-likelihood function of the parameters given the function data \mathbf{y} is:

$$l(s, m, \theta | \mathbf{y}) = -\frac{1}{2} s^{-2} (\mathbf{y} - m\mathbf{1})^T C^{-1} (\mathbf{y} - m\mathbf{1}) - \frac{1}{2} \log(\det C) - \frac{n}{2} \log(s^2) + \text{constants}. \quad (9)$$

Maximizing the log-likelihood first with respect to m , then with respect to s , and finally with respect to θ

yields

$$m_{\text{MLE}} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \quad (10)$$

$$\begin{aligned} s_{\text{MLE}}^2 &= \frac{1}{n} (\mathbf{y} - m_{\text{MLE}} \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{MLE}} \mathbf{1}) \\ &= \frac{1}{n} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}, \end{aligned} \quad (11)$$

$$\begin{aligned} \boldsymbol{\theta}_{\text{MLE}} &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left\{ \frac{1}{n} \log(\det \mathbf{C}) \right. \\ &\quad \left. + \log \left(\mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \right) \right\}. \end{aligned} \quad (12)$$

Under these estimates of the parameters, the cubature (7) and the credible interval (8) simplify to

$$\hat{\mu}_{\text{MLE}} = \left(\frac{(1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}) \mathbf{1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + \mathbf{c} \right)^T \mathbf{C}^{-1} \mathbf{y}, \quad (14)$$

$$\begin{aligned} \text{err}_{\text{MLE}}^2 &:= \frac{2.58^2}{n} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \\ &\quad \times (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}), \end{aligned} \quad (15)$$

$$\mathbb{P}_f [|\mu - \hat{\mu}_{\text{MLE}}| \leq \text{err}_{\text{MLE}}] = 99\%. \quad (16)$$

Here c_0 , \mathbf{c} , and \mathbf{C} are assumed implicitly to be based on $\boldsymbol{\theta} = \boldsymbol{\theta}_{\text{MLE}}$.

2.2.2 Full Bayes

Rather than use maximum likelihood to determine m and s one treat them as hyper-parameters and may a non-informative, conjugate prior, namely $\rho_{m,s^2}(\xi, \lambda) \propto 1/\lambda$. Then the posterior density for the integral given the data using Bayes theorem is

$$\begin{aligned} \rho_\mu(z|\mathbf{f} = \mathbf{y}) &\propto \int_0^\infty \int_{-\infty}^\infty \rho_\mu(z|\mathbf{f} = \mathbf{y}, m = \xi, s^2 = \lambda) \\ &\quad \times \rho_{\mathbf{f}}(\mathbf{y}|\xi, \lambda) \rho_{m,s^2}(\xi, \lambda) d\xi d\lambda \\ &\propto \int_0^\infty \frac{1}{\lambda^{(n+3)/2}} \\ &\quad \times \int_{-\infty}^\infty \exp \left(-\frac{1}{2\lambda} \left\{ \frac{[z - \xi(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}]^2}{c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}} \right. \right. \\ &\quad \left. \left. + (\mathbf{y} - \xi \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - \xi \mathbf{1}) \right\} \right) d\xi d\lambda \\ &\quad \text{by (5a), (6) and } \rho_{m,s^2}(\xi, \lambda) \propto 1/\lambda \\ &\propto \int_0^\infty \frac{1}{\lambda^{(n+3)/2}} \cdots \\ &\quad \cdots \int_{-\infty}^\infty \exp \left(-\frac{\alpha \xi^2 - 2\beta \xi + \gamma}{2\lambda(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) d\xi d\lambda, \end{aligned}$$

where

$$\alpha = (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2 + \mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}),$$

$$\begin{aligned} \beta &= (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})(z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}) \\ &\quad + \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}), \end{aligned}$$

$$\gamma = (z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y})^2 + \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}).$$

In the derivation above and below, factors that are independent of ξ , λ , or z can be discarded since we only need to preserve the proportion, however, factors that depend on ξ , λ , or z must be kept. Completing the square $\alpha \xi^2 - 2\beta \xi + \gamma = \alpha(\xi - \beta/\alpha)^2 - (\beta^2/\alpha) + \gamma$, allows us to evaluate the integrals with respect to ξ and λ :

$$\begin{aligned} \rho_\mu(z|\mathbf{f} = \mathbf{y}) &\propto \int_0^\infty \frac{1}{\lambda^{(n+3)/2}} \exp \left(-\frac{\gamma - \beta^2/\alpha}{2\lambda(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) \cdots \\ &\quad \cdots \int_{-\infty}^\infty \exp \left(-\frac{\alpha(\xi - \beta/\alpha)^2}{2\lambda(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) d\xi d\lambda \\ &\propto \int_0^\infty \frac{1}{\lambda^{(n+2)/2}} \exp \left(-\frac{\gamma - \beta^2/\alpha}{2\lambda(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \right) d\lambda \\ &\propto \left(\gamma - \frac{\beta^2}{\alpha} \right)^{-n/2} \propto (\alpha\gamma - \beta^2)^{n/2}. \end{aligned}$$

Finally, we simplify the key term via straightforward calculations to the following:

$$\alpha\gamma - \beta^2 \propto 1 + \frac{(z - \hat{\mu}_{\text{MLE}})^2}{(n-1)\hat{\sigma}_{\text{full}}^2},$$

where

$$\begin{aligned} \hat{\sigma}_{\text{full}}^2 &:= \frac{1}{n-1} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \\ &\quad \times \left[\frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \right]. \end{aligned} \quad (17)$$

This means that $\mu(\mathbf{f} = \mathbf{y})$, properly centered and scaled, has a Student's t -distribution with $n-1$ degrees of freedom. The estimated integral is the same as in the empirical Bayes case, but the confidence interval is wider:

$$\mathbb{P}_f [|\mu - \hat{\mu}_{\text{MLE}}| \leq \text{err}_{\text{full}}] = 99\%, \quad (18)$$

where

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

Here $t_{n-1, 0.995}$ denotes the 99.5 percentile of a standard Student's t -distribution with $n-1$ degrees of freedom.

Because the shape parameter, $\boldsymbol{\theta}$, enters the definition of the covariance kernel in a non-trivial way, the only way to treat it as a hyperparameter and assign a tractable prior would be for the prior to be discrete. We believe in practice that choosing such a prior involves more guesswork than using the empirical Bayes estimate of $\boldsymbol{\theta}$ in (12) or the cross-validation approach described next.

2.2.3 Cross-validation

A third parameter optimization is *leave-one-out cross-validation*. Let $\tilde{y}_i = \mathbb{E}[f(\mathbf{x}_i) | \mathbf{f}_{-i} = \mathbf{y}_{-i}]$, where the subscript $-i$ denotes the vector excluding the i^{th} component. This is the conditional expectation of $f(\mathbf{x}_i)$ given all data but the function value at \mathbf{x}_i . The cross-validation criterion, which is to be minimized, is sum of squares of the difference between these conditional expectations and the observed values:

$$\text{CV} = \sum_{i=1}^n (y_i - \tilde{y}_i)^2. \quad (20)$$

Let $\mathbf{A} = \mathbf{C}^{-1}$, let $\boldsymbol{\zeta} = \mathbf{A}(\mathbf{y} - m\mathbf{1})$, and partition \mathbf{C} , \mathbf{A} , and $\boldsymbol{\zeta}$ as

$$\mathbf{C} = \begin{pmatrix} c_{ii} & \mathbf{C}_{-i,i}^T \\ \mathbf{C}_{-i,i} & \mathbf{C}_{-i,-i} \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} a_{ii} & \mathbf{A}_{-i,i}^T \\ \mathbf{A}_{-i,i} & \mathbf{A}_{-i,-i} \end{pmatrix},$$

$$\boldsymbol{\zeta} = \begin{pmatrix} \zeta_i \\ \boldsymbol{\zeta}_{-i} \end{pmatrix},$$

where the subscript i denotes the i^{th} row or column, and the subscript $-i$ denotes all rows or columns except the i^{th} . Following this notation, Lemma 1 implies that

$$\begin{aligned} \tilde{y}_i &= m + \mathbf{C}_{-i,i}^T \mathbf{C}_{-i,-i}^{-1} (\mathbf{y}_{-i} - m\mathbf{1}) \\ \zeta_i &= a_{ii}(y_i - m) + \mathbf{A}_{-i,i}^T (\mathbf{y}_{-i} - m\mathbf{1}) \\ &= a_{ii}[(y_i - m) - \mathbf{C}_{-i,i}^T \mathbf{C}_{-i,-i}^{-1} (\mathbf{y}_{-i} - m\mathbf{1})] \\ &= a_{ii}(y_i - \tilde{y}_i). \end{aligned}$$

Thus, (20) may be re-written as

$$\text{CV} = \sum_{i=1}^n \left(\frac{\zeta_i}{a_{ii}} \right)^2, \quad \boldsymbol{\zeta} = \mathbf{C}^{-1}(\mathbf{y} - m\mathbf{1}).$$

Wahba ?? recommended a generalized cross-validation criterion, which replaces the i^{th} diagonal element of \mathbf{A} in the denominator by the average diagonal element of \mathbf{A} :

$$\begin{aligned} \text{GCV} &= \frac{\sum_{i=1}^n \zeta_i^2}{\left(\frac{1}{n} \sum_{i=1}^n A_{ii} \right)^2} \\ &= \frac{(\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-2} (\mathbf{y} - m\mathbf{1})}{\left(\frac{1}{n} \text{trace}(\mathbf{C}^{-1}) \right)^2}. \end{aligned} \quad (21)$$

The loss function GCV depends on m and $\boldsymbol{\theta}$, but not on s . Minimizing the GCV yields

$$m_{\text{GCV}} = \frac{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}}, \quad (22)$$

$$\boldsymbol{\theta}_{\text{GCV}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \log \left(\mathbf{y}^T \left[\mathbf{C}^{-2} - \frac{\mathbf{C}^{-2} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-2}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}} \right] \mathbf{y} \right) - 2 \log (\text{trace}(\mathbf{C}^{-1})). \quad (23)$$

Plugging this value of m into (7) yields

$$\hat{\mu}_{\text{GCV}} = \left(\frac{(1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}) \mathbf{C}^{-1} \mathbf{1}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}} + \mathbf{c} \right)^T \mathbf{C}^{-1} \mathbf{y}. \quad (24)$$

An estimate for s may be obtained by noting that by Lemma 1,

$$\text{var}[f(\mathbf{x}_i) | \mathbf{f}_{-i} = \mathbf{y}_{-i}] = s^2 a_{ii},$$

Thus, we may estimate s using a similar idea

$$s^2 = \frac{\text{var}[f(\mathbf{x}_i) | \mathbf{f}_{-i} = \mathbf{y}_{-i}]}{a_{ii}} \quad (25)$$

$$\approx \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2 \quad (26)$$

2.3 Formulating the automatic Bayesian cubature algorithm

Let the function to integrate be $f \sim \mathcal{GP}(m, s^2 \mathbf{C}_{\boldsymbol{\theta}})$. Our goal is to compute $\hat{\mu}$ within the error tolerance ε , i.e., $|\mu - \hat{\mu}| \leq \varepsilon$.

It is not possible to know true error $|\mu - \hat{\mu}|$ for real problems. So, the error-bound err_n is used as the approximate error estimate. Basic principle of our algorithm is to keep adding more function values till the stopping criterion is met, i.e., $\text{err}_n \leq \varepsilon$. Once the stopping criterion is met, the approximation of μ can be computed:

$$\hat{\mu}_n = w_0 + \mathbf{w}^T \mathbf{y}, \quad (27)$$

Where the suffix n is used to imply the number of integrand-samples used. In every iteration of the *automatic cubature algorithm (1)*, we only need to numerically estimate the MLE of $\boldsymbol{\theta}$ and then use it to compute \mathbf{C} and err_n .

The following pseudo-code briefly explains the working of the automatic cubature algorithm.

Algorithm 1 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 initial sample size, n_0

- 1: $n \leftarrow n_0, n' \leftarrow 0$
- 2: Generate $\{\mathbf{x}_i\}_{i=n'+1}^n$ and sample $\{f(\mathbf{x}_i)\}_{i=n'+1}^n$
- 3: Compute $\boldsymbol{\theta}$ by (??)
- 4: Compute error bound err_n $\triangleright \text{err}_n$ data driven error bound
- 5: **if** $\text{err}_n \leq \varepsilon$ **then break**
- 6: **end if**
- 7: $n' \leftarrow n, n \leftarrow 2 \times n'$
- 8: Compute cubature weights $\{w_i\}_{i=1}^n$
- 9: Compute approximate integral $\hat{\mu}_n$
- 10: **return** $\hat{\mu}_n$ \triangleright Integral estimate $\hat{\mu}_n$

As shown, the algorithm continues the iteration loop till the stopping-criterion is met, i.e., err_n is smaller than error threshold ε . At the end of every iteration, if the computed error bound err_n is higher than the required ε , algorithm doubles the number of points n and repeats. When the error tolerance ε is met, exits the loop. Finally using the n and the MLE estimated parameters, computes the $\hat{\mu}_n$.

Accurately computing the data-driven error bound err_n which closely matches the true error is essential for the effectiveness of our algorithm. So, accurate and faster computation of err_n and $\hat{\mu}_n$ are the main objective of the following sections.

2.4 Example with Matern kernel

We would like to demonstrate and test numerical accuracy and computational cost of the automatic Bayesian cubature algorithm using the results obtained so far. For this example, we use the Matern kernel:

$$C_{\theta}(\mathbf{x}, \mathbf{t}) = \prod_{k=1}^d \exp(-\theta_k |\mathbf{x}_k - \mathbf{t}_k|) (1 + \theta_k |\mathbf{x}_k - \mathbf{t}_k|) \quad (28)$$

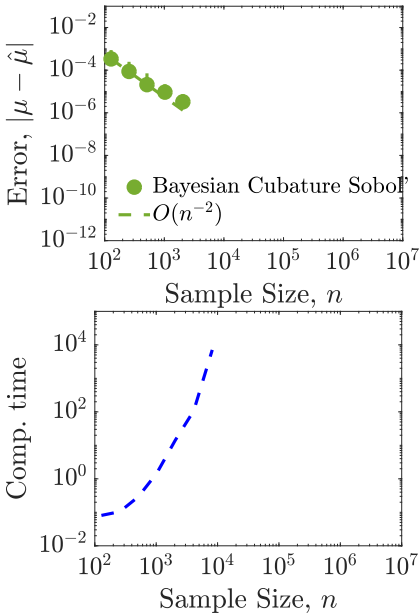


Fig. 1: MVN for $d=2$ with Matern kernel

On our test computer with Intel i7 3630QM and 16GB RAM memory, it took 118 minutes to compute $\hat{\mu}_n$ with $n = 2^{13}$. As shown in Figure 1, computation time increases rapidly with n . Especially, Maximum-likelihood-estimation of θ which needs the loss function, is the most time consuming of all. Because the loss functions need to be computed multiple times in every iteration to choose the optimal shape parameter till its minimum is found. Not only the complexity increases, also the kernel matrix becomes highly ill-conditioned with increasing n number of data-points. We must use alternative techniques to overcome this problem. So, our algorithm in the current form is not straightaway usable for any practical applications.

3 Fast automatic Bayesian cubature

Automatic Bayesian cubature algorithm described (Section 2.3) is slow due to its computationally intensive nature. This could be a hindrance in real world applications when a larger number of data points are needed to achieve the desired accuracy. This section and further sections explore techniques to make it faster.

Bayesian cubature algorithm uses error-bound err_n to decide when to stop, and uses MLE loss function (12) to find the optimal shape parameter. These computations involve covariance matrix inversions and multiplications, which are time-consuming and prone to numerical error. When the algorithm tries to adapt by increasing n to get better accuracy, the covariance matrix size also grows, along with it, the matrix's condition number also grows significantly causing numerical error in matrix operations. This phenomenon is called ill-conditioning.

One approach to overcome these issues could be to use an efficient method to compute the whole equation without incurring the ill-conditioning, which involves, for example, using *stable computation* [5]. But it is going to be still slower and time-consuming.

Another approach for stable and faster computation is, using a carefully chosen kernel with some special properties which leads to faster matrix operations, like faster matrix inversion and multiplication. This approach can be designed to be faster in computation while avoiding the numerical error due to ill-conditioning. The later approach is the major objective of this research work. Our approach is especially to choose a special form of kernel, which is called *fast transform kernel* along with a suitable pointset, that will avoid expensive matrix operations.

3.1 Fast transform kernel

Suppose the domain is $\mathcal{X} = [0, 1]^d$ and the probability measure is uniform. If the kernel $C_{\theta}(\mathbf{x}, \mathbf{t})$ is chosen carefully along with appropriate point set $\{\mathbf{x}_i\}_{i=1}^n$, have the special properties, so the resulting Gram matrix has factorization of the form:

$$\begin{aligned} \mathbf{C} &= \left(C_{\theta}(\mathbf{x}_i, \mathbf{x}_j) \right)_{i,j=1}^n = (\mathbf{C}_1, \dots, \mathbf{C}_n) \\ &= \frac{1}{n} \mathbf{V} \mathbf{\Lambda} \mathbf{V}^H = \frac{1}{n} \mathbf{V}^* \mathbf{\Lambda} \mathbf{V}^T, \quad \mathbf{V}^H = n \mathbf{V}^{-1}, \end{aligned} \quad (29)$$

Where \mathbf{V}^H is the Hermitian of \mathbf{V} ,

$\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)^T = (\mathbf{V}_1, \dots, \mathbf{V}_n)$, also $\mathbf{v}_1 = \mathbf{V}_1 = \mathbf{1}$

$\mathbf{\Lambda} = \text{diag}(\boldsymbol{\lambda})$, $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$,

$$\mathbf{V}^{-1} = \frac{1}{n} \mathbf{V}^H, \quad \mathbf{C}^{-1} = \frac{1}{n} \mathbf{V} \mathbf{\Lambda}^{-1} \mathbf{V}^H = \frac{1}{n} \mathbf{V}^* \mathbf{\Lambda}^{-1} \mathbf{V}^T.$$

Suppose the transform $\hat{\mathbf{z}} = \mathbf{V}^T \mathbf{z}$ for any arbitrary \mathbf{z} can be computed within the computational cost of $\mathcal{O}(n \log n)$, we call it a *Fast transform*. Consequently \mathbf{C}_θ is called a *fast transform kernel*. We use the notation $\hat{\mathbf{z}}$ to indicate the transform of \mathbf{z} . The properties of fast transform can be leveraged to make the computations in Bayesian cubature algorithm faster. To begin with, $\boldsymbol{\Lambda}$ can be computed faster by simply,

$$\begin{aligned} \mathbf{V}^T \mathbf{C}_1 &= \mathbf{V}^T \left(\frac{1}{n} \mathbf{V}^* \boldsymbol{\Lambda} \mathbf{V}_1 \right) \\ &= \underbrace{\left(\frac{1}{n} \mathbf{V}^T \mathbf{V}^* \right)}_{\mathbf{I}} \boldsymbol{\Lambda} \mathbf{V}_1 = \boldsymbol{\Lambda} \mathbf{1} = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{pmatrix} = \boldsymbol{\lambda} \quad (30) \end{aligned}$$

Most of the computations in the algorithm are of the form $\mathbf{a}^T \mathbf{C}^{-1} \mathbf{b}$, can be simplified easily:

$$\mathbf{C} \mathbf{V}_1 = \lambda_1 \mathbf{V}_1, \quad \text{by definition}$$

$$\text{So, } \mathbf{C}^{-1} \mathbf{1} = \mathbf{C}^{-1} \mathbf{V}_1 = \frac{\mathbf{V}_1}{\lambda_1} = \frac{\mathbf{1}}{\lambda_1},$$

Similarly,

$$\mathbf{a}^T \mathbf{C}^{-1} \mathbf{b} = \frac{1}{n} \mathbf{a}^T \mathbf{V} \boldsymbol{\Lambda}^{-1} \mathbf{V}^H \mathbf{b} = \frac{1}{n} \hat{\mathbf{a}}^T \boldsymbol{\Lambda}^{-1} \hat{\mathbf{b}}^* = \frac{1}{n} \sum_{i=1}^n \frac{\hat{a}_i \hat{b}_i^*}{\lambda_i},$$

$$\mathbf{a}^T \mathbf{C}^{-2} \mathbf{b} = \frac{1}{n} \mathbf{a}^T \mathbf{V} \boldsymbol{\Lambda}^{-2} \mathbf{V}^H \mathbf{b} = \frac{1}{n} \hat{\mathbf{a}}^T \boldsymbol{\Lambda}^{-2} \hat{\mathbf{b}}^* = \frac{1}{n} \sum_{i=1}^n \frac{\hat{a}_i \hat{b}_i^*}{\lambda_i^2},$$

where $\hat{\mathbf{a}} = \mathbf{V}^T \mathbf{a}$ and $\hat{\mathbf{b}} = \mathbf{V}^T \mathbf{b}$. By using this, some of the recurring terms in the algorithm are simplified,

$$\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} = \mathbf{1}^T \left(\frac{\mathbf{1}}{\lambda_1} \right) = \frac{n}{\lambda_1}$$

$$\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y} = \mathbf{1}^T \left(\frac{\mathbf{y}}{\lambda_1} \right) = \frac{\sum_{i=1}^n y_i}{\lambda_1} = \frac{\hat{y}_1}{\lambda_1}$$

$$\mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} = \frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i}$$

$$\mathbf{c}^T \mathbf{C}^{-1} \mathbf{1} = \mathbf{c}^T \left(\frac{\mathbf{1}}{\lambda_1} \right) = \frac{\hat{c}_1}{\lambda_1}$$

$$\mathbf{c}^T \mathbf{C}^{-1} \mathbf{c} = \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i}$$

$$\mathbf{y}^T \mathbf{C}^{-2} \mathbf{y} = \frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i^2}$$

Then the MLE estimates can be simplified using these results:

$$m_{\text{MLE}} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} = \frac{\hat{y}_1}{n} = \frac{1}{n} \sum_{i=1}^n y_i,$$

$$\begin{aligned} s_{\text{MLE}}^2 &= \frac{1}{n} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \\ &= \frac{1}{n} \left[\frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i} - \frac{|\hat{y}_1|^2 / \lambda_1^2}{n / \lambda_1} \right] = \frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \end{aligned}$$

3.2 Application of fast transform kernel

Using the properties of the fast transform, matrix multiplications and inversions can be avoided, so the overall computation cost is $\mathcal{O}(n \log n)$. To begin with, MLE estimate of $\boldsymbol{\theta}$ can be made faster:

3.2.1 Empirical Bayes

$$\begin{aligned} \boldsymbol{\theta}_{\text{MLE}} &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\frac{1}{n} \log(\det \mathbf{C}) + \log(s_{\text{MLE}}^2) \right] \\ &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\frac{1}{n} \sum_{i=1}^n \log(\lambda_i) + \log \left(\frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \right) \right] \quad (31) \end{aligned}$$

$$\text{where } \hat{\mathbf{y}} = (\hat{y}_i)_{i=1}^n = \mathbf{V}^T \mathbf{y}.$$

Similarly, the error bound err_n computation can be made faster:

$$\begin{aligned} \text{err}_n &= 2.58 \sqrt{\frac{1}{n} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \\ &= 2.58 \sqrt{\frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \left(c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right)}, \end{aligned}$$

$$\text{where } \hat{\mathbf{y}} = \mathbf{V}^T \mathbf{y}, \hat{\mathbf{c}} = \mathbf{V}^T \mathbf{c}.$$

Finally, the computation of $\hat{\mu}_n$ is made faster:

$$\begin{aligned} \hat{\mu}_{\text{MLE}}(f) &= m_{\text{MLE}}[1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}] + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y} \\ &= \frac{\hat{y}_1}{n} \left[1 - \frac{\hat{c}_1}{\lambda_1} \right] + \frac{1}{n} \sum_{i=1}^n \frac{\hat{c}_i \hat{y}_i^*}{\lambda_i} \\ &= \frac{\hat{y}_1}{n} + \frac{1}{n} \sum_{i=2}^n \frac{\hat{c}_i \hat{y}_i^*}{\lambda_i} \quad (32) \end{aligned}$$

It is interesting to note that, in general with no assumption on m , $\hat{\mu}_{\text{MLE}}$ is simply the sample mean. Moreover the choice of kernel or any MLE estimated parameters do not affect the $\hat{\mu}_n$ but influences the accuracy of the error bound err_n .

3.2.2 Full Bayes

Generalized cross validation criterion:

$$\begin{aligned} \boldsymbol{\theta}_{\text{GCV}} &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\log \left(\frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} - \frac{\hat{y}_1^2 / \lambda_1^4}{n / \lambda_1^2} \right) \right. \\ &\quad \left. - 2 \log \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \\ &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[\log \left(\frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \end{aligned}$$

The cubature rule:

$$\begin{aligned}\hat{\mu}_{\text{full}}(\mathbf{f} = \mathbf{y}) &= \left[\frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) \mathbf{1}^T}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + \mathbf{c}^T \right] \mathbf{C}^{-1} \mathbf{y} \\ &= \frac{\hat{y}_1}{n} + \frac{1}{n} \sum_{i=2}^n \frac{\hat{c}_i \hat{y}_i^*}{\lambda_i},\end{aligned}\quad (33)$$

$$\begin{aligned}\hat{\sigma}_{\text{full}}^2(\mathbf{f} = \mathbf{y}) &= \frac{1}{n-1} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \\ &\quad \times \left[\frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \right] \\ &= \frac{1}{n-1} \frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \\ &\quad \times \left[\frac{\lambda_1}{n} \left(1 - \frac{\hat{c}_1}{\lambda_1} \right)^2 + \left(c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \right]\end{aligned}$$

The stopping criterion for the full Bayes case,

$$\begin{aligned}n = \min \left\{ n_j \in \mathbb{Z}_{>0} : \frac{t_{n_j-1,0.995}^2}{n_j(n_j-1)} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \right. \\ \left. \times \left[\frac{\lambda_1}{n_j} \left(1 - \frac{\hat{c}_1}{\lambda_1} \right)^2 + \left(c_0 - \frac{1}{n_j} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \right] \leq \varepsilon^2 \right\}.\end{aligned}\quad (34)$$

These simplified computations involve no matrix inversion or multiplications. It just uses scalar divisions and multiplications so the computational cost is $\mathcal{O}(n)$. But computation of fast transform $\hat{\mathbf{y}}$ is $\mathcal{O}(n \log n)$. Consequently the overall computation cost is $\mathcal{O}(n \log n)$.

4 Implementation 1: Using a shift invariant kernel for Bayesian cubature

We have established the concept of fast transform and showed how it can make the computations faster. But we assumed there exist a kernel that meets the requirements. Now we are going to show a example. The following shift invariant kernel satisfies the requirements of the fast transform kernel when combined with Rank-1 lattice points:

$$C_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{t}) := \sum_{\mathbf{k} \in \mathbb{Z}^d} \alpha_{\mathbf{k}, \boldsymbol{\theta}} e^{2\pi \sqrt{-1} \mathbf{k}^T \mathbf{x}} e^{-2\pi \sqrt{-1} \mathbf{k}^T \mathbf{t}}, \quad (35)$$

where $\alpha_{\mathbf{k}, \boldsymbol{\theta}} := \prod_{l=1}^d \frac{1}{\max(\frac{|k_l|}{\theta_l}, 1)^r}$, with $\alpha_{\mathbf{0}, \boldsymbol{\theta}} = 1$,

where d is number of dimensions and $\alpha_{\mathbf{k}}$ is a scalar. The Gram matrix formed by this kernel is a Hermitian matrix. With some more simplifications:

$$C_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{t}) = \prod_{l=1}^d \left[1 + \sum_{k_l=1}^{\infty} \left| \frac{\theta_l}{k_l} \right|^r 2 \cos(2\pi \sqrt{-1} k_l (x_l - t_l)) \right].$$

The *shape parameter* $\boldsymbol{\theta}$ is used to make the kernel customizable. To be specific, the shape parameter

tweaks the kernel, so that the function space spanned by the kernel closely resembles the space where the integrand function belongs.

This form of the kernel is very convenient to use in any analytical derivations or proofs, but not suitable for use with finite precision computers as this kernel involves infinite sum. It is preferred to have a simpler expression of the kernel without infinite sum for practical computations.

4.1 Using Lattice points

Along with the kernel (35), Rank-1 lattice points $\mathbf{x}_i \in \mathbb{L}_{n, \mathbf{h}}$ are used to get the *fast transform kernel*. In this work, we use the lattice points as defined in [15]:

$$\mathbb{L}_{n, \mathbf{h}} := \{ \mathbf{x}_i := \mathbf{h} \frac{i-1}{n} \pmod{1}; i = 1, \dots, n \},$$

Where \mathbf{h} is the generating vector. Its dual lattice is defined as:

$$\mathbb{L}_{n, \mathbf{h}}^{\perp} := \{ \mathbf{k} \in \mathbb{Z}^d : \mathbf{h}^T \mathbf{k} \equiv 0 \pmod{n} \},$$

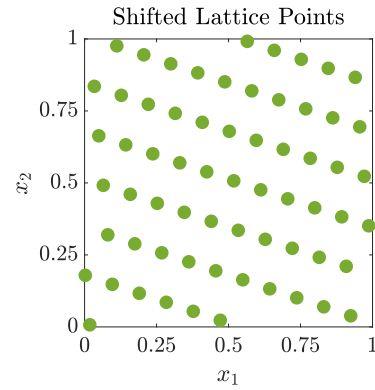


Fig. 2: Randomly shifted Lattice points in d=2

With lattice points, the kernel is written as:

$$\begin{aligned}C(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{\mathbf{k} \in \mathbb{Z}^d} \alpha_{\mathbf{k}} e^{2\pi \sqrt{-1} \mathbf{k}^T \frac{\mathbf{h}_i}{n}} e^{-2\pi \sqrt{-1} \mathbf{k}^T \frac{\mathbf{h}_j}{n}} \\ &= \sum_{\mathbf{k} \in \mathbb{Z}^d} \alpha_{\mathbf{k}} e^{2\pi \sqrt{-1} \mathbf{k}^T \mathbf{h} \frac{(i-j)}{n}}\end{aligned}$$

Please note that ‘ $\pmod{1}$ ’ need not be used explicitly in $C(\mathbf{x}_i, \mathbf{x}_j)$ since $e^{\pm 2\pi \sqrt{-1}} = 1$. The kernel (35), when used with Rank-1 Lattice points, leads to symmetric and circulant Gram matrix \mathbf{C} . We can demonstrate the resulting Gram matrix satisfies the conditions of a fast transform.

4.2 Computing the kernel using Bernoulli polynomials

The shift-invariant-kernel (35) cannot be computed directly due to it’s infinite sum. If the coefficients $\alpha_{\mathbf{k}}$ were

chosen appropriately, then there exist a direct expression for the kernel in terms of Bernoulli polynomial. We can use the Fourier series expansion properties [4] of Bernoulli polynomial B_r to find the appropriate α_k . This provides a closed form expression for the kernel without the infinite sum. It also allows to choose the smoothness of kernel, i.e, how fast the coefficients α_k in (35) decay. The following very useful expansion is referenced from [4] under eqn (24.8.3):

$$B_r(x) = \frac{-r!}{(2\pi\sqrt{-1})^r} \sum_{\substack{k \neq 0, \\ k=-\infty}}^{\infty} \frac{e^{2\pi\sqrt{-1}kx}}{k^r} \quad \begin{cases} \text{for } r = 1, & 0 < x < 1 \\ \text{for } r = 2, 3, \dots & 0 \leq x \leq 1 \end{cases} \quad (36)$$

Our goal is to replace the infinite sum of the kernel using Bernoulli polynomials. This would allow the computations to be carried out in any software like Matlab. For 1-D ($d = 1$) rewriting (36):

$$\begin{aligned} \sum_{k \neq 0, k=-\infty}^{\infty} \frac{e^{2\pi\sqrt{-1}k|x|}}{\left(\frac{k}{\theta}\right)^r} \\ = -\theta^r B_r(|x|) \frac{(2\pi\sqrt{-1})^r}{r!} \text{ for } -1 \leq x \leq 1 \end{aligned}$$

Comparing this expansion with (35), one can deduce, for $d = 1$:

$$C(x, t) = \alpha_0 + \sum_{k \neq 0, k=-\infty}^{\infty} \alpha_k e^{2\pi\sqrt{-1}k|x-t|}, \quad -1 \leq x, t \leq 1,$$

Where the kernel coefficients α_k can be expressed explicitly:

$$\alpha_k = \begin{cases} 1, & k = 0 \\ \frac{1}{k^r}, & \text{otherwise} \end{cases},$$

In general for $d > 1$, the coefficients α_k for $C_{\theta}(x, t)$ is computed using:

$$\alpha_{k, \theta} = \frac{1}{\prod_{l=1}^d \max\left(\frac{|k_l|}{\theta_l}, 1\right)^r},$$

Where θ is the shape parameter and ' r ' is the Bernoulli polynomial. The value of ' r ' will be chosen specific to the integrand when using in the cubature algorithm. We keep the value of ' r ' to be even positive integer in this work. Using the above-found results, we get the simplified form of the kernel:

$$C_{\theta}(x, t) = \prod_{l=1}^d 1 - \theta_l^r \frac{(2\pi\sqrt{-1})^r}{r!} B_r(|x_l - t_l|), \quad \text{where } 0 \leq |x_l - t_l| \leq 1, \quad (37)$$

where $\theta \in (0, 1]^d$ is the shape parameter, ' r ' is the order of the Bernoulli polynomial B_r .

We call (37) a Fourier kernel as the resulting Fast transform is the *Fast Fourier transform* (FFT). Plots of this kernel in $d = 2$ is shown in Figure 3.

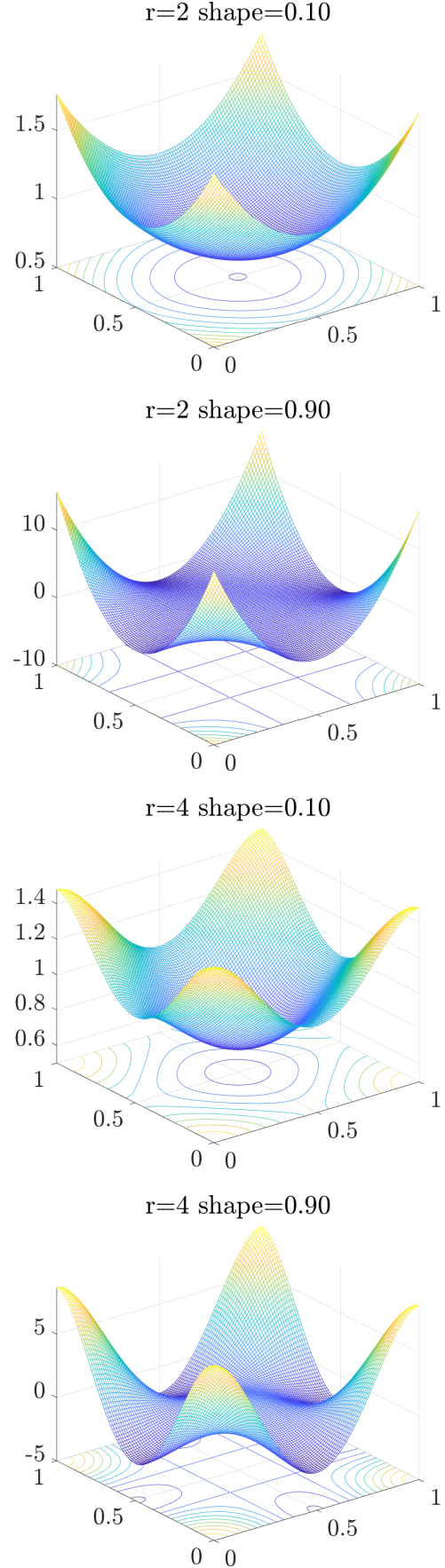


Fig. 3: Fourier kernel (37) with $r = 2, 4$ and $\theta = 0.1, 0.9$

4.3 Eigenvalues and Eigenvectors of \mathbf{C}

The kernel's Gram matrix \mathbf{C} is circulant. It is generated by the vector:

$$(C_{\theta}(\mathbf{x}_1, \mathbf{x}_1), C_{\theta}(\mathbf{x}_2, \mathbf{x}_1), \dots, C_{\theta}(\mathbf{x}_n, \mathbf{x}_1))^T =: \mathbf{C}_1, \quad (38)$$

Which is the first row or column of \mathbf{C} (since it is a symmetric matrix). By the properties of circulant matrix, the normalized eigenvectors are:

$$\left(1, e^{-2\pi\sqrt{-1}\frac{j}{n}}, e^{-2\pi\sqrt{-1}\frac{2j}{n}}, \dots, e^{-2\pi\sqrt{-1}\frac{j(n-1)}{n}}\right)^T, \\ \text{for } j = 0, 1, \dots, n-1.$$

The matrix constructed with these eigenvectors is:

$$\mathbf{V} := \left(e^{-2\pi\sqrt{-1}\frac{ij}{n}}\right)_{i,j=0}^{n-1}, \quad \mathbf{V}^H := \frac{1}{n} \left(e^{2\pi\sqrt{-1}\frac{ij}{n}}\right)_{i,j=0}^{n-1},$$

Where the first column of \mathbf{V} , $\mathbf{V}_1 = \mathbf{1}$. The columns of \mathbf{V} are complex exponential vectors independent of the kernel values. But their corresponding eigenvalues are:

$$\Lambda = (\lambda_j)_{j=1}^n = \mathbf{V}^T \mathbf{C}_1 := \mathcal{DFT}\{\mathbf{C}_1\},$$

i.e., eigenvalues of the matrix \mathbf{C} are computed by applying DFT over \mathbf{C}_1 . for this kernel, the corresponding fast transform is *Fast Fourier Transform* (FFT). The kernel's Gram matrix \mathbf{C} can be written in the factorized form:

$$\mathbf{C} = \frac{1}{n} \mathbf{V} \Lambda \mathbf{V}^H,$$

This is the exact factorization used for the construction of fast transform. Another requirement $\mathbf{V}_1 = \mathbf{1}$ is also met. Additionally, computational cost of FFT is $\mathcal{O}(n \log n)$. Thus we can conclude the shift invariant kernel in eqn (35), obeys all the requirements to be considered as a fast transform kernel. So the operation $\mathbf{V}^T \mathbf{z}$ is a fast transform. Numerical experiments using this kernel will be shown in the section (4.9).

4.4 Additional assumptions

The shift invariant kernel (35) has some more desirable properties, leading to $c_0 = 1$ and $\mathbf{c} = \mathbf{1}$, which will help to make the computations even faster:

$$c_0 = \int_{[0,1]^d \times [0,1]^d} C_{\theta}(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} = 1, \\ \mathbf{c} = \left(\int_{[0,1]^d} C_{\theta}(\mathbf{x}_i, \mathbf{t}) d\mathbf{t} \right)_{i=1}^n = \mathbf{1}.$$

This gives simplified error bound:

$$\text{err}_n = 2.58 \sqrt{\frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \left(c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right)}, \\ = 2.58 \sqrt{\frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \left(1 - \frac{n}{\lambda_1} \right)},$$

$$\text{where } \hat{\mathbf{c}} = (n, 0, \dots, 0)^T, \quad \hat{\mathbf{y}} = \mathbf{V}^T \mathbf{y}, \hat{\mathbf{c}} = \mathbf{V}^T \mathbf{c}.$$

And the cubature:

$$\hat{\mu}_n(f) = w_0 + \mathbf{w}^T \mathbf{y} = \frac{\hat{y}_1}{n} + \frac{1}{n} \sum_{i=2}^n \frac{\hat{c}_i \hat{y}_i^*}{\lambda_i} \\ = \frac{\hat{y}_1}{n}, \quad \text{since } \hat{\mathbf{c}}_i = 0 \quad \forall i \neq 0$$

4.4.1 Full Bayes

Final version of full Bayes equations optimized using the properties of the shift invariant kernel are:

$$\hat{\mu}_{\text{full}}(\mathbf{f} = \mathbf{y}) = \frac{\hat{y}_1}{n}, \quad (39)$$

$$\theta_{\text{GCV}} = \underset{\theta}{\text{argmin}} \left[\log \left(\frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \quad (40)$$

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

The stopping criterion for the full Bayes case,

$$n = \min \left\{ n_j \in \mathbb{Z}_{>0} : \right. \\ \left. \frac{t_{n_j-1, 0.995}^2}{n_j(n_j-1)} \left(\frac{\lambda_1}{n_j} - 1 \right) \times \sum_{i=2}^{n_j} \frac{|\hat{y}_i|^2}{\lambda_i} \leq \varepsilon^2 \right\}. \quad (41)$$

4.5 Variable transforms

Fast transform kernel technique discussed with shift invariant kernel eqn(35) and lattice points so far, assumes the integrand is periodic and has continuous derivatives on the boundaries of the domain $[0, 1]^d$. Non-periodic functions do not live in the space spanned by the kernel. Variable transformation or periodization transform techniques are typically used to improve the accuracy in multi-dimensional numerical integrations where boundary conditions needs to be enforced. These transformations could be either polynomial, exponential and also

trigonometric nature.

$$\text{Baker's : } \tilde{f}(t) = f\left(\left(1 - 2\left|t_j - \frac{1}{2}\right|\right)_{j=1}^d\right)$$

$$\text{C0 : } \tilde{f}(t) = f(\mathbf{g}_0) \prod_{j=1}^d g'_0(t_j), \mathbf{g}_0 = (g_0(t_j))_{j=1}^d$$

$$\text{C1 : } \tilde{f}(t) = f(\mathbf{g}_1) \prod_{j=1}^d g'_1(t_j), \mathbf{g}_1 = (g_1(t_j))_{j=1}^d$$

$$\text{Sidi's C1 : } \tilde{f}(t) = f(\boldsymbol{\psi}_2) \prod_{j=1}^d \psi'_2(t_j), \boldsymbol{\psi}_2 = (\psi_2(t_j))_{j=1}^d$$

$$\text{Sidi's C2 : } \tilde{f}(t) = f(\boldsymbol{\psi}_3) \prod_{j=1}^d \psi'_3(t_j), \boldsymbol{\psi}_3 = (\psi_3(t_j))_{j=1}^d$$

where

$$g_0(t) = 3t^2 - 2t^3, g'_0(t) = 6t(1 - t)$$

$$g_1(t) = t^3(10 - 15t + 6t^2), g'_1(t) = 30t^2(1 - t)^2$$

$$\psi_2(t) = \left(t - \frac{1}{2\pi} \sin(2\pi t)\right), \psi'_2(t) = (1 - \cos(2\pi t))$$

$$\psi_3(t) = \frac{1}{16} (8 - 9 \cos(\pi t) + \cos(3\pi t)),$$

$$\psi'_3(t) = \frac{1}{16} (9 \sin(\pi t)\pi - \sin(3\pi t)3\pi)$$

These transforms vary in terms of computational complexity and accuracy, shall be chosen on a need basis. Such as:

1. Baker's : Baker's transform or called tent map in each coordinate. It preserves only continuity but it is easier to compute.
2. C0 : Polynomial transformation only ensures periodicity of function.
3. C1 : Polynomial transformation preserving the first derivative.
4. C1sin : Sidi's transform with Sine, preserving the first derivative. This is, in general, a better option than 'C1'.
5. C2sin : Sidi's transform with Sine, preserving upto second derivative. We use this when smoothness of 'C1sin' is not sufficient and need to preserve upto second derivative.

4.6 Iterative Discrete Fourier transform for function values

Automatic cubature algorithms needs to compute Discrete Fourier transform of function values $\mathbf{y} = (f(\mathbf{x}_i))_{i=1}^n$ in every iteration with newly added points. Recomputing the whole Fourier transform in every iteration can be avoided when using Rank-1 Lattice points by using structural properties of the Lattice points. Discrete

Fourier transform is defined as :

$$\mathcal{DFT}\{y\} := \hat{y} = \left(\sum_{j=1}^n y_j e^{-\frac{2\pi\sqrt{-1}}{n}(j-1)(i-1)}\right)_{i=1}^n$$

In essence:

$$\hat{y}_i = \sum_{j=1}^n y_j e^{-\frac{2\pi\sqrt{-1}}{n}(j-1)(i-1)}$$

We could rearrange the sum into even indexed $j = 2l$ and odd indexed $j = 2l + 1$:

$$\begin{aligned} \hat{y}_i = & \underbrace{\sum_{l=1}^{n/2} y_{2l} e^{-\frac{2\pi\sqrt{-1}}{n/2}(l-1)(i-1)}}_{\text{DFT of even-indexed part of } y_i} \\ & + e^{-\frac{2\pi\sqrt{-1}}{n}(i-1)} \underbrace{\sum_{l=1}^{n/2} y_{2l+1} e^{-\frac{2\pi\sqrt{-1}}{n/2}(l-1)(i-1)}}_{\text{DFT of odd-indexed part of } y_i}, \end{aligned}$$

Which shows two separately computed DFTs can be combined to produce single output. For example, the odd indexed were the existing DFT and the even indexed are from the new halve of samples, algorithm wants to add to improve accuracy. We use this concept to avoid recomputing the full length DFT of \mathbf{y} in every iteration. In other words, DFT is computed only for the newly added samples in every iteration.

4.7 Overcoming the cancellation error

During the numerical experiments, we noticed numerical cancellation error in the computation of the term, especially for the bigger values $n > 2^{15}$. Cancellation error happens because two almost equal values are subtracted, when they differ only in very high decimal values.

$$(c_{0,\theta} - \mathbf{c}_\theta^T \mathbf{C}^{-1} \mathbf{c}_\theta) = \left(1 - \frac{n}{\lambda_1}\right)$$

In this expression $\frac{n}{\lambda_1}$ almost close to 1. We would like to explore techniques to avoid the subtraction in this computation. Let's recollect the definition of the shift invariant kernel:

$$C(\mathbf{x}_i, \mathbf{x}_j) = \prod_{k=1}^d [1 + \theta B(\{x_{i_k} - x_{j_k}\})].$$

Let's define the modified kernel:

$$\tilde{C}(\mathbf{x}_i, \mathbf{x}_j) := C(\mathbf{x}_i, \mathbf{x}_j) - 1,$$

Then Its Gram matrix: $\tilde{\mathbf{C}} = \mathbf{C} - \mathbf{1}\mathbf{1}^T$,

Let $(\lambda_1, \dots, \lambda_n)$ be the eigenvalues of \mathbf{C} . Similarly let $(\tilde{\lambda}_1, \dots, \tilde{\lambda}_n)$ be the eigenvalues of $\tilde{\mathbf{C}}$. As per the definition of the Gram matrix \mathbf{C} , the eigenvector corresponding to λ_1 is a vector one $\mathbf{1}$. Then:

$$\lambda_1 \mathbf{1} = \mathbf{C}\mathbf{1} = (\tilde{\mathbf{C}} + \mathbf{1}\mathbf{1}^T)\mathbf{1} = \tilde{\lambda}_1 \mathbf{1} + n\mathbf{1},$$

This shows $\tilde{\lambda}_1 = \lambda - n$. In fact for the rest of the values are:

$$\tilde{\lambda}_j = \lambda_j, \quad \forall j = 1, \dots, n-1. \quad (42)$$

Let $\mathbf{v}_j, \forall j = 0, \dots, n-1$ be the eigenvectors of \mathbf{C} , similarly $\tilde{\mathbf{v}}_j, \forall j = 0, \dots, n-1$ be the eigenvectors of $\tilde{\mathbf{C}}$, $\mathbf{v}_j^T \mathbf{C} \mathbf{1} = \lambda_1 \mathbf{v}_j^T \mathbf{1} = \mathbf{1}^T \mathbf{C} \mathbf{v}_j = \lambda_j \mathbf{1}^T \mathbf{v}_j$. Since $\lambda_1 \neq \lambda_j$, the above statement implies $\mathbf{v}_j^T \mathbf{1} = 0$. This interesting property provides the proof of the eqn. (42).

$$\lambda_j \mathbf{v}_j = \mathbf{C} \mathbf{v}_j = \mathbf{1}_{n \times n} \mathbf{v}_j + \tilde{\mathbf{C}} \mathbf{v}_j = \tilde{\mathbf{C}} \mathbf{v}_j = \tilde{\lambda}_j \mathbf{v}_j,$$

Thus proven. Using this result, cancellation error in the computation of err_n can be avoided:

$$\begin{aligned} \left(1 - \frac{n}{\lambda_1}\right) &= \left(1 - \frac{n}{\tilde{\lambda}_1 + n}\right) \\ &= \left(\frac{\tilde{\lambda}_1 + n - n}{\tilde{\lambda}_1 + n}\right) = \left(\frac{\tilde{\lambda}_1}{\tilde{\lambda}_1 + n}\right). \end{aligned}$$

The following technique shows an iterative approach to compute $\tilde{C}(\mathbf{x}_i, \mathbf{x}_j)$ for $d > 1$. This iterative technique is developed using induction:

$$\begin{aligned} d=1: \quad C_1 &= 1 + \theta B(\{x_{i_1} - x_{j_1}\}) = 1 + \tilde{C}_1 \\ d=2: \quad C_2 &= [1 + \theta B(\{x_{i_2} - x_{j_2}\})][1 + \tilde{C}_1] \\ &= 1 + \theta B(\{x_{i_2} - x_{j_2}\})[1 + \tilde{C}_1] + \tilde{C}_1 \\ &= 1 + \underbrace{\theta B(\{x_{i_2} - x_{j_2}\})C_1}_{\tilde{C}_2} + \tilde{C}_1 \\ d=3 \quad C_3 &= [1 + \theta B(\{x_{i_3} - x_{j_3}\})][1 + \tilde{C}_2] \\ &= 1 + \theta B(\{x_{i_3} - x_{j_3}\})[1 + \tilde{C}_2] + \tilde{C}_2 \\ &= 1 + \theta B(\{x_{i_3} - x_{j_3}\})C_2 + \tilde{C}_2 \\ &\vdots \\ \forall d > 2 \quad C_d &= [1 + \theta B(\{x_{i_d} - x_{j_d}\})][1 + \tilde{C}_{d-1}] \\ &= 1 + \theta B(\{x_{i_d} - x_{j_d}\})[1 + \tilde{C}_{d-1}] + \tilde{C}_{d-1} \\ &= 1 + \theta B(\{x_{i_d} - x_{j_d}\})C_{d-1} + \tilde{C}_{d-1} \end{aligned}$$

4.8 Validating Gaussian process assumption

We begin developing the Bayesian cubature with the assumption that the integrand arises from a Gaussian process. How can we check if m, s, θ and \mathbf{C} are chosen appropriately so that the \mathbf{f} is a draw from the Gaussian process?. Here is an attempt to validate that assumption. Let,

$$\mathbf{f} = (f(\mathbf{x}_i))_{i=1}^n \sim \mathcal{N}(m\mathbf{1}, \mathbf{C}),$$

$$\text{where } \mathbf{C} = \frac{1}{n} \mathbf{V} \mathbf{\Lambda} \mathbf{V}^H, \quad \mathbf{V}^H \mathbf{V} = n$$

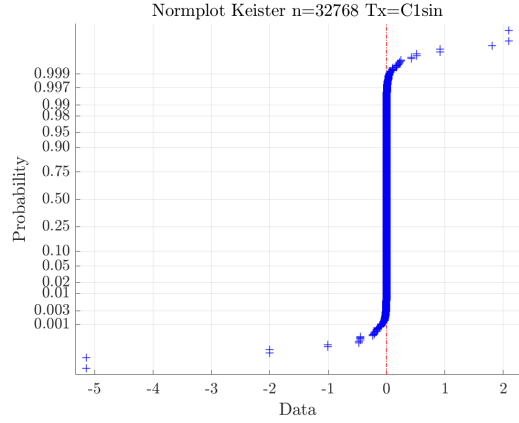


Fig. 4: Normal plot : Keister function with arbMean assumption

Then,

$$\begin{aligned} \mathbb{E} \left[\frac{1}{\sqrt{n}} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \mathbf{f} \right] &= \frac{1}{\sqrt{n}} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \mathbb{E}[\mathbf{f}] \\ &= \frac{1}{\sqrt{n}} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H m\mathbf{1} \\ &= m \sqrt{\frac{n}{\lambda_1}} \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \end{aligned}$$

$$\text{Let } \mathbf{f}' = \frac{1}{\sqrt{n}} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \mathbf{f},$$

$$\begin{aligned} \text{cov}[\mathbf{f}'] &= \frac{1}{n} \mathbb{E} \left[\mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H (\mathbf{f} - m\mathbf{1})(\mathbf{f} - m\mathbf{1})^T \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}} \right] \\ &= \frac{1}{n} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \mathbb{E}[(\mathbf{f} - m\mathbf{1})(\mathbf{f} - m\mathbf{1})^T] \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}} \\ &= \frac{1}{n} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{V}^H \frac{1}{n} \mathbf{V} \mathbf{\Lambda} \mathbf{V}^H \mathbf{V} \mathbf{\Lambda}^{-\frac{1}{2}} \\ &= \mathbf{I} \end{aligned}$$

Thus

$$\mathbf{f}' \sim \mathcal{N}(m' \mathbf{e}_1, \mathbf{I}),$$

Where $m' = m \sqrt{\frac{n}{\lambda_1}}$. If we can verify the sample distribution of \mathbf{f}' is approximately $\mathcal{N}(m' \mathbf{e}_1, \mathbf{I})$ by using Normal plots, could validate our assumption.

We plotted the Normal plots of \mathbf{f}' for Keister and Multivariate Normal functions as shown in Figure 5 and Figure 4. These plots show the transformed \mathbf{f}' follows approximately standard normal with a constant mean.

4.9 Numerical Experiments

Having all the tools and optimization done handy, now we shall run the numerical simulations to see the performance.

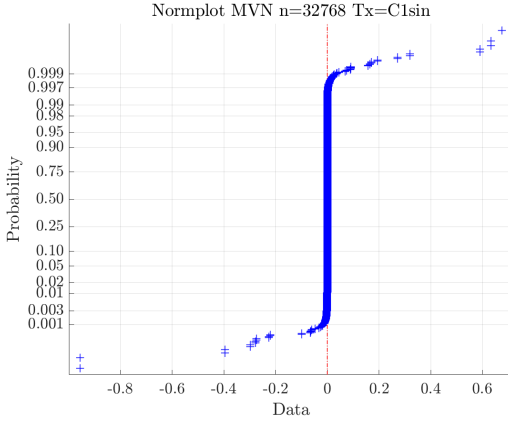


Fig. 5: Normal plot : MVN with arbMean assumption

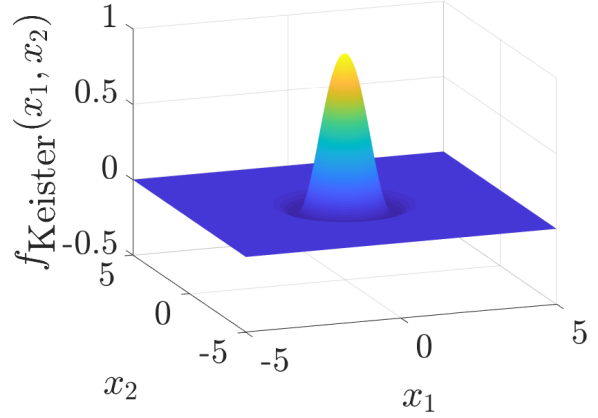


Fig. 7: Keister function in d=2

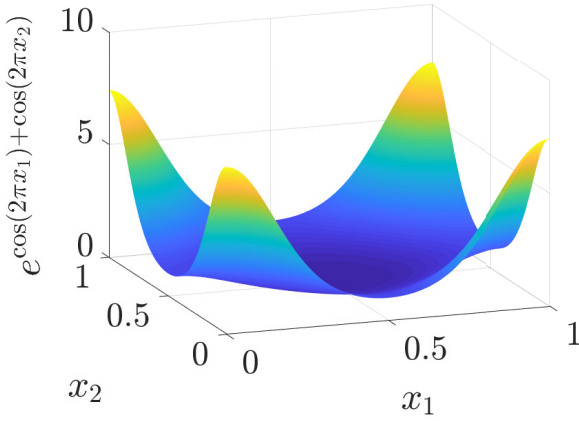
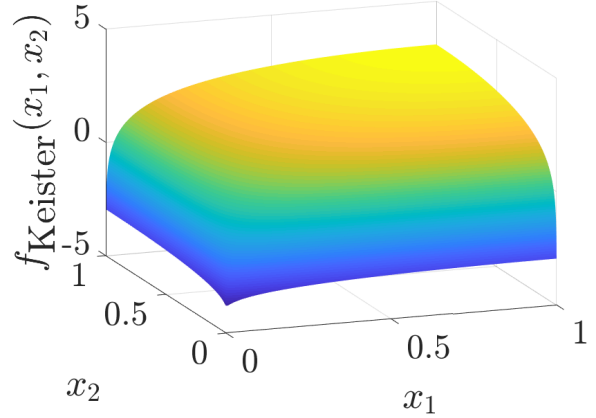


Fig. 6: Exp(Cos) in d=2

Fig. 8: Keister function transformed to $[0, 1]^2$

4.9.1 Test Functions

The following test functions were used to test the integration speed and accuracy of *Fast Bayesian cubature algorithm* that we developed so far

1. Exponential of Cosine:

This is a very simple function, serves as an example to check the working of the cubature, where the function is defined as

$$f(\mathbf{x}) = e^{\sum_{i=1}^d \cos(2\pi x_i)}, \quad \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} = \text{BesselI}(0,1)^d$$

where ‘BesselI’ is the *modified Bessel* function. Exp(Cos) function is periodic in $[0, 1]$, so we do not need to use any *transform* to make it periodic.

2. Keister function

The following multidimensional integral example is based on the paper [10], inspired by a physics application.

$$f(\mathbf{x}) = \cos(\|\mathbf{x}\|) \exp(-\|\mathbf{x}\|^2) d\mathbf{x},$$

$$\int_{\mathbb{R}^d} f(\mathbf{x}) d\mathbf{x} = \frac{2\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2})} \text{CI}(d), \quad d = 1, 2, 3, \dots$$

where, $\Gamma :=$ gamma function, and

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

$$\begin{aligned} \text{SI}(1) &= \int_{x=0}^{\infty} \exp(-\mathbf{x}^T \mathbf{x}) \sin(\mathbf{x}) d\mathbf{x} \\ &= 0.4244363835020225, \end{aligned}$$

$$\text{CI}(2) = \frac{1 - \text{SI}(1)}{2}, \quad \text{SI}(2) = \frac{\text{CI}(1)}{2}$$

$$\text{CI}(j) = \frac{(j-2)\text{CI}(j-2) - \text{SI}(j-1)}{2}, \quad j = 3, 4, \dots$$

$$\text{SI}(j) = \frac{(j-2)\text{SI}(j-2) - \text{CI}(j-1)}{2}, \quad j = 3, 4, \dots$$

3. Multivariate Normal:

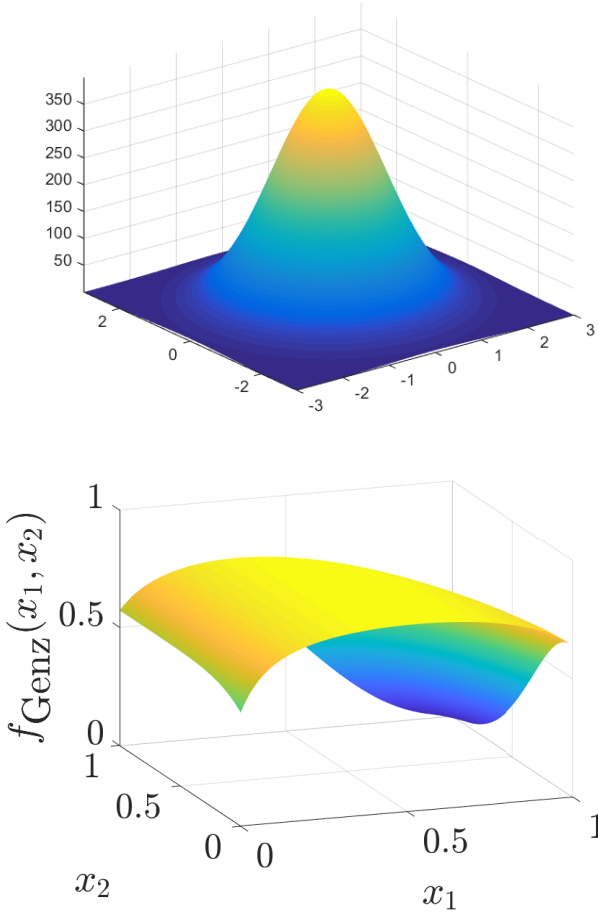


Fig. 9: Mutivariate Normal and Genz function

We use the Genz's method to compute Multivariate normal probability as explained below. This method reduces the original dimension of the problem by 1.

$$\mu = \int_{[\mathbf{a}, \mathbf{b}] \in \mathbb{R}^d} \frac{\exp(-\frac{1}{2} \mathbf{t}^T \Sigma^{-1} \mathbf{t})}{\sqrt{(2\pi)^d \det(\Sigma)}} d\mathbf{t} \quad \stackrel{[6]}{=} \int_{[0,1]^{d-1}} f_{\text{Genz}}(\mathbf{x}) d\mathbf{x}$$

where $\Sigma = \mathbf{L}\mathbf{L}^T$ is the Cholesky decomposition of the covariance matrix, $\mathbf{L} = (l_{jk})_{j,k=1}^d$ is a lower triangular matrix, and

$$\begin{aligned} \alpha_1 &= \Phi(a_1), & \beta_1 &= \Phi(b_1) \\ \alpha_j(x_1, \dots, x_{j-1}) &= \\ \Phi\left(\frac{1}{l_{jj}} \left(a_j - \sum_{k=1}^{j-1} l_{jk} \Phi^{-1}(\alpha_k + x_k(\beta_k - \alpha_k))\right)\right), \\ & j = 2, \dots, d, \\ \beta_j(x_1, \dots, x_{j-1}) &= \\ \Phi\left(\frac{1}{l_{jj}} \left(b_j - \sum_{k=1}^{j-1} l_{jk} \Phi^{-1}(\alpha_k + x_k(\beta_k - \alpha_k))\right)\right), \\ & j = 2, \dots, d, \\ f_{\text{Genz}}(\mathbf{x}) &= \prod_{j=1}^d [\beta_j(\mathbf{x}) - \alpha_j(\mathbf{x})]. \end{aligned}$$

As we see from the figure, Genz function is not periodic. So we need to make it periodic to get the best accuracy with Bayesian cubature.

We use the following parameter values for the numerical examples

	a	b	L
$d = 2$	$\begin{pmatrix} -6 \\ -2 \\ -2 \end{pmatrix}$	$\begin{pmatrix} 5 \\ 2 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 4 & 1 & 1 \\ 0 & 1 & 0.5 \\ 0 & 0 & 0.25 \end{pmatrix}$
$d = 3$	$\begin{pmatrix} -6 \\ -2 \\ -2 \\ 2 \end{pmatrix}$	$\begin{pmatrix} 5 \\ 2 \\ 1 \\ 2 \end{pmatrix}$	$\begin{pmatrix} 4 & 1 & 1 & 1 \\ 0 & 1 & 0.5 & 0.5 \\ 0 & 0 & 0.25 & 0.25 \\ 0 & 0 & 0 & 0.25 \end{pmatrix}$

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 Brownian motion, B , at time t_1, \dots, t_d , then $Z = (B(t_1), \dots, B(t_d)) \sim \mathcal{N}(\mathbf{0}, \Sigma)$, where $\Sigma = (\min(t_j, t_k))_{j,k=1}^d$, and 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 the function $\text{payoff}(\cdot)$ defines the discounted payoff of the option,

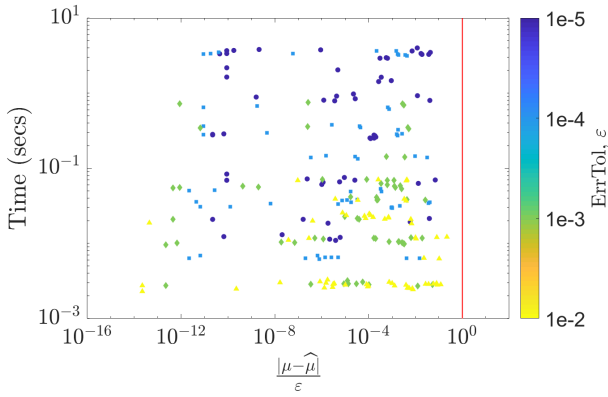
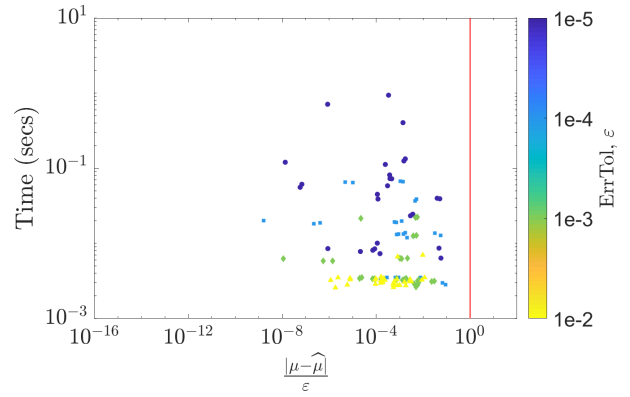
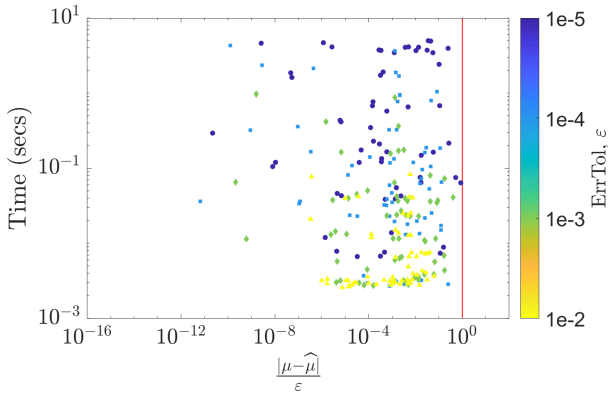
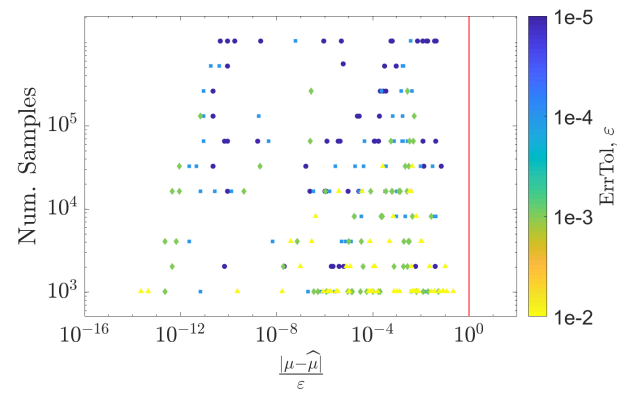
$$f(\mathbf{x}) = \text{payoff}(\mathbf{z}), \quad \mathbf{z} = \mathbf{L} \begin{pmatrix} \Phi^{-1}(x_1) \\ \vdots \\ \Phi^{-1}(x_d) \end{pmatrix},$$

and \mathbf{L} is any square matrix satisfying $\Sigma = \mathbf{L}\mathbf{L}^T$.

For the Asian arithmetic mean call option:

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

Where $S_j = S_0 \exp((r - \sigma^2/2)t_j + \sigma z_j)$, d - number of dimensions and T, S, S_0, K, r, σ are the parameters to be specified.

Fig. 10: Exp(Cos) integrated within ε and finite timeFig. 12: MVN integrated within ε and finite timeFig. 11: Keister function integrated within ε and finite timeFig. 13: Exp(Cos) integrated within ε using finite number of samples

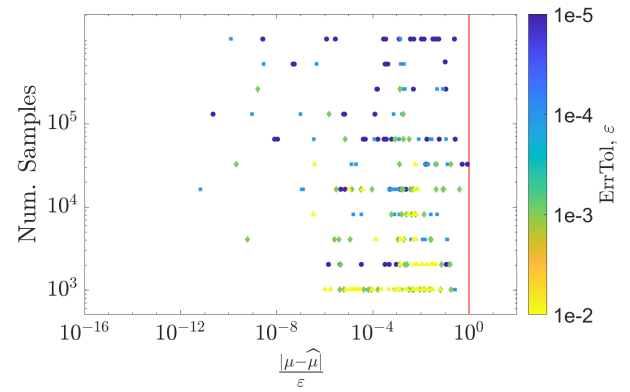
4.10 Results and observation

We tested our cubature to integrate the above specified functions. The following plots summarize the results. We used a random shift of $\delta \sim \mathcal{U}[0, 1]^d$ to randomly shift the rank-1 Lattice points. This provides randomness in the cubature's behavior. Then we run the cubature for 100 times and take the mean of it as the final result. We integrated with the error thresholds $\varepsilon \in \{1E-2, 1E-3, 1E-4, 1E-5\}$. We use $\frac{|\mu - \hat{\mu}|}{\varepsilon}$ as x-axis so that a single plot can include all these ε .

The figures Figure 10, Figure 11 and Figure 12 show the integration done in a finite time within the error tolerance ε . Similarly, the figures Figure 13, Figure 14 and Figure 15 show the integration done using a finite number of samples within the error tolerance ε .

5 Conclusion

We developed a generalized technique of *Fast transform kernel*. Using this technique, further developed

Fig. 14: Keister integrated within ε using finite number of samples

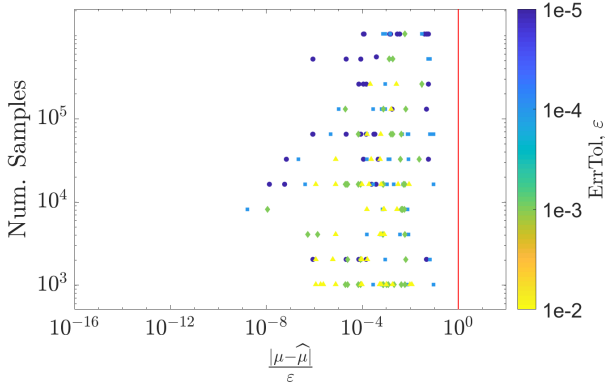


Fig. 15: MVN integrated within ε using finite number of samples

fast automatic Bayesian cubature algorithm that takes very low computational cost in the order of $\mathcal{O}(n \log n)$ comparing to direct Bayesian cubature of $\mathcal{O}(n^3)$, so it can be used in practical applications. By adjusting the Bernoulli order r of the kernel and using the appropriately smoother variable transformation, we could get higher order of error convergence when the integrand is assumed to have zero mean. In general case without any assumption on the integrand mean, the optimal $\hat{\mu}$ is just the sample mean and so the order of convergence is defined by how smoother the underlying function being integrated and the variable transformation being used. Though the optimal $\hat{\mu}$ is just the sample mean, the error bound err_n still depends on the kernel order. So when we use higher order r , the error bound err_n closely matches the actual error. We could use the algorithm to integrate upto 2^{23} data points on a 16GB of RAM memory and i7-3630QM computer within 5 minutes. Numerical results show that the theoretical error err_n closely matches the actual error but it could still be improved with more tighter error bound especially when the r is smaller.

6 Future work

As an extension to the ideas and techniques established in this work, the following are considered potential future work

1. Sobol points and Fast Walsh Transform (FWT)

We have shown one example of a special form of kernel with defined requirements to build a *Fast transform kernel*. Using the established generalized requirements for a fast-transform-kernel, we could use the same approach with other kernels with suitable point-sets to achieve similar or better performance and accuracy. One such point-sets that to

consider in future is, *Sobol points* and with appropriate choice of kernel, which should lead to *Fast Walsh Transform*.

2. Control variates

We would like to approximate a function of the form $(f - \beta_1 g_1 - \dots, -\beta_p g_p)$ than

$$f = \mathcal{N}(\beta_0 + \beta_1 g_1 + \dots, +\beta_p g_p, s^2 \mathbf{C})$$

3. Function approximation

Let us consider approximating a function of the form

$$\int_{[0,1]^d} \underbrace{f(\phi(t)) \cdot \left| \frac{\partial \phi}{\partial t} \right|}_{g(t)} dt$$

where $\left| \frac{\partial \phi}{\partial t} \right|$ is Jacobian and then

$$g(\psi(x)) = f(\underbrace{\phi(\psi(x))}_x) \cdot \left| \frac{\partial \phi}{\partial t} \right|(\psi(x))$$

$$f(x) = g(\psi(x)) \cdot \frac{1}{\left| \frac{\partial \phi}{\partial t} \right|(\psi(x))}$$

Finally, the function approximation is

$$\tilde{f}(x) = \tilde{g}(\psi(x))$$

$$= \sum w_i C(.,.)$$

4. Deterministic interpretation of Bayesian cubature

7 Appendix

7.1 Properties of Multivariate Normal Distributions

Lemma 1 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

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

$$\mathbf{C}_{11} - \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{21}).$$

Proof : The inverse of the matrix \mathbf{C} may be partitioned as [?]

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

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

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

Since \mathbf{Y}_1 and \mathbf{Y}_2 have a Gaussian distribution, the conditional distribution of $\mathbf{Y}_1 | \mathbf{Y}_2$ is also Gaussian. So, it is sufficient to prove that $\mathbb{E}(\mathbf{Y}_1 | \mathbf{Y}_2)$ and $\text{cov}(\mathbf{Y}_1 | \mathbf{Y}_2)$ are as hypothesized.

Define $\mathbf{Z} = \mathbf{Y}_1 - \mathbf{B}\mathbf{Y}_2$ where $\mathbf{B} = \mathbf{C}_{21}^T \mathbf{C}^{-1}_{22}$. Note that $\mathbb{E}(\mathbf{Z}) = \mathbf{m}_1 - \mathbf{B}\mathbf{m}_2$. Elementary calculations establish that \mathbf{Z} and \mathbf{Y}_2 are uncorrelated and therefore independent. Therefore it follows that the conditional mean is

$$\begin{aligned}\mathbb{E}(\mathbf{Y}_1|\mathbf{Y}_2) &= \mathbb{E}(\mathbf{Z} + \mathbf{B}\mathbf{Y}_2|\mathbf{Y}_2) \\ &= \mathbb{E}(\mathbf{Z}|\mathbf{Y}_2) + \mathbb{E}(\mathbf{B}\mathbf{Y}_2|\mathbf{Y}_2) = \mathbb{E}(\mathbf{Z}) + \mathbf{B}\mathbf{Y}_2 \\ &= \mathbf{m}_1 + \mathbf{B}(\mathbf{Y}_2 - \mathbf{m}_2) \\ &= \mathbf{m}_1 + \mathbf{C}_{21}^T \mathbf{C}^{-1}_{22}(\mathbf{Y}_2 - \mathbf{m}_2)\end{aligned}$$

Likewise the conditional covariance is

$$\begin{aligned}\text{cov}(\mathbf{Y}_1|\mathbf{Y}_2) &= \text{cov}(\mathbf{Z} + \mathbf{B}\mathbf{Y}_2|\mathbf{Y}_2) \\ &= \text{cov}(\mathbf{Z}) = \text{var}(\mathbf{Y}_1 - \mathbf{B}\mathbf{Y}_2) \\ &= \text{cov}(\mathbf{Y}_1) + \mathbf{B}\text{cov}(\mathbf{Y}_2)\mathbf{B}^T \\ &\quad - 2\text{cov}(\mathbf{Y}_1, \mathbf{Y}_2)\mathbf{B}^T \\ &= \mathbf{C}_{11} + \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{22} \mathbf{C}_{22}^{-1} \mathbf{C}_{21} - 2\mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{21} \\ &= \mathbf{C}_{11} - \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{21}.\end{aligned}$$

This completes the proof. \square

References

1. Briol, F.X., Oates, C.J., Griolami, M., Osborne, M.A., Sejdinovic, D.: Probabilistic integration: A role in statistical computation? *Statist. Sci.* (2018+). To appear
2. Cools, R., Nuyens, D. (eds.): Monte Carlo and Quasi-Monte Carlo Methods: MCQMC, Leuven, Belgium, April 2014, *Springer Proceedings in Mathematics and Statistics*, vol. 163. Springer-Verlag, Berlin (2016)
3. Diaconis, P.: Bayesian numerical analysis. Statistical decision theory and related topics IV, Papers from the 4th Purdue symp., West Lafayette, Indiana 1986, p. 163–175 (1988)
4. DLMF: Nist digital library of mathematical functions. <http://dlmf.nist.gov/> **Part 2**, Release 1.0.5 of 2012–10–01 (2012)
5. Fasshauer, G.E.: Meshfree Approximation Methods with Matlab. World Scientific Publishing Co (2007)
6. Genz, A.: Comparison of methods for the computation of multivariate normal probabilities. *Computing Science and Statistics* **25**, 400 – 405 (1993)
7. Hickernell, F.J.: The trio identity for quasi-Monte Carlo error analysis. In: P. Glynn, A. Owen (eds.) Monte Carlo and Quasi-Monte Carlo Methods: MCQMC, Stanford, USA, August 2016, *Springer Proceedings in Mathematics and Statistics*, pp. 13–37. Springer-Verlag, Berlin (2018). ArXiv:1702.01487
8. Hickernell, F.J., Jiménez Rugama, L.I.A.: Reliable adaptive cubature using digital sequences. In: Cools and Nuyens [2], pp. 367–383. ArXiv:1410.8615 [math.NA]
9. Jiménez Rugama, L.I.A., Hickernell, F.J.: Adaptive multidimensional integration based on rank-1 lattices. In: Cools and Nuyens [2], pp. 407–422. ArXiv:1411.1966
10. Keister, B.D.: Multidimensional quadrature algorithms. *Computers in Physics* **10**, 119–122 (1996)
11. O’Hagan, A.: Bayes-hermite quadrature. *J. Statist. Plann. Inference* **29**, 245–260 (1991)
12. Rasmussen, C.E., Ghahramani, Z.: Bayesian monte carlo. *Advances in Neural Information Processing Systems* pp. 489–496 (2003)
13. Rasmussen, C.E., Williams, C.: Gaussian Processes for Machine Learning. MIT Press, Cambridge, Massachusetts (2006). (online version at <http://www.gaussianprocess.org/gpml/>)
14. Ritter, K.: Average-case analysis of numerical problems. *Lecture Notes in Mathematics*, Springer-Verlag, Berlin **vol. 1733**, 163–175 (2000)
15. Sabrina Dammertz, A.K.: Image synthesis by rank-1 lattices. *Monte Carlo and Quasi-Monte Carlo Methods Part 2*, 217–236 (2006)