

# 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 integral  $\mu = \int g(x) dx$  of a multi-dimensional function  $f$  where  $\mu$  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  $\mu$  is typically expressed in the form

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

where  $f : \mathcal{X} \rightarrow \mathbb{R}$  is the integrand and  $\nu$  is a probability measure defined on the measurable set  $\mathcal{X} \subseteq \mathbb{R}^d$ . The problem of choosing  $f$  and  $\nu$  well to correspond to the original integrand,  $g$ , is the matter of importance sampling, which we do not address here. The cubature algorithm is defined as

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

where the weights,  $w_0$ , and  $\mathbf{w} = \{w_i\}_{i=1}^n$ , and the nodes  $\{\mathbf{x}_i\}_{i=1}^n$  are chosen to make the error,  $|\mu(f) - \hat{\mu}(f)|$ , small.

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.,

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

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 by a Bayesian assumption on the problem. We build upon the ideas of Diaconis [Dia88], O'Hagan [O'H91], Ritter [Rit00], Rasmussen and Ghahramani [RG03], and others. Our 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 are obtained for cubature purposes, then  $\mathcal{O}(n \log(n))$  addi-

---

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

S. Author  
second address

tional operations are required to check whether the error tolerance is satisfied. This is significantly fewer operations than the  $\mathcal{O}(n^3)$  typically required for Bayesian cubature.

Traditional cubature methods assume the integrand to be a deterministic function. An alternative to this approach is to assume that the integrand is a *stochastic process*. Bayesian cubature (BC) methods assume the integrand is an instance of a Gaussian process. BC methods approximate the multivariate integrals that cannot be evaluated analytically, using weighted sums of integrand values at carefully chosen nodes over the domain  $\mathcal{X}$  as shown in eqn (2).

Traditional cubature methods may not provide any guaranteed error accuracy. The goal of this work is to develop a guaranteed, BC algorithm using the confidence intervals given by Bayesian posterior error. Our algorithm strives to find the minimum sample size ‘ $n$ ’ - *number of integrand values* required so that the error  $|\mu(f) - \hat{\mu}(f)|$  is less than the user-defined error threshold  $\varepsilon$ , i.e.,

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

But the true error  $|\mu(f) - \hat{\mu}(f)|$  will not be available for problems that do require cubature methods. So we compute an approximate error bound  $\text{err}_n$  using the confidence interval obtained from the posterior error which meets,

$$\text{err}_n \leq |\mu(f) - \hat{\mu}(f)|$$

By calculating a data-driven error bound, our algorithm adaptively determines ‘ $n$ ’. An added advantage of our approach is being data-driven, so it does not require any other parameter of the integrand like the ‘total variation’.

The “algorithms for numerical tasks” that return uncertainties in their calculations are called *probabilistic numeric methods*. Our BC algorithm can be categorized as a probabilistic numeric method, due to the nature of assumptions made and the confidence interval it provides.

Section 2 introduces the Bayesian approach to estimate the posterior error and develops the concepts to derive error bound. Then formulates the Automatic cubature algorithm using the error bound. Finally demonstrates why directly using Bayesian cubature algorithm is computationally very expensive. Section 3 Introduces the concept of Fast transform kernel and develops the concepts to make the Bayesian cubature faster which is the major contribution of this research. Section 4 Demonstrates the realization of Fast transform kernel using a shift-invariant kernel and Rank-1 Lattice points. 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.

## 1.1 Related work

## 2 Bayesian cubature

This section introduces the concepts of the Bayesian cubature and derives all the basic results which will be used in the next section for further speedup of the algorithm. *Bayesian posterior error* [Hic17] is used to estimate an error bound for the cubature.

### 2.1 Bayesian posterior error

Random  $f$  postulated by Diaconis [Dia88], O’Hagen [O’H91], Ritter [Rit00], Rasmussen and Ghahramani [RG03] and others:  $f \sim \mathcal{GP}(m, s^2 C_\theta)$ , an instance of a Gaussian process with mean  $m$ , non-negative scale factor  $s$  and covariance function  $s^2 C_\theta$ . Where  $C_\theta : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$  is a symmetric, positive-definite function and parameterized by  $\theta$ :

$$C(\mathbf{t}, \mathbf{x}) = C(\mathbf{x}, \mathbf{t}), \quad \sum_{i,j=1}^n a_i C(\mathbf{x}_i, \mathbf{x}_j) a_j > 0 \text{ when } \mathbf{a} \neq \mathbf{0},$$

$$\forall n \in \mathbb{N}, \mathbf{t}, \mathbf{x}, \mathbf{x}_1, \dots, \mathbf{x}_n \in [0, 1]^d, \mathbf{a} = (a_1, \dots, a_n) \in \mathbb{R}^n. \quad (4)$$

In this work, we use domain  $\mathcal{X} = [0, 1]^d$  and uniform measure. The scale parameter  $s^2$  and shape parameter  $\theta$  should be estimated.

For a Gaussian process, all vectors of linear functionals of  $f$  have a multivariate Gaussian distribution. Sometimes we drop  $\theta$  in the writings to simplify the notation. With this assumption,

$$\mu(f) \sim \mathcal{N}(m\mu(1), s^2 c_0), \quad \text{where } c_0 = \int_{[0,1]^2} C_\theta(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}, \quad (5a)$$

$$\mathbf{f} = (f(\mathbf{x}_i))_{i=1}^n \sim \mathcal{N}(m\mathbf{1}, s^2 \mathbf{C}), \quad \text{where } \mathbf{C} = (C_\theta(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n. \quad (5b)$$

Let us consider a fixed node set design  $\{\mathbf{x}_i\}_{i=1}^n$ , then we observe  $\mathbf{y} = \{y_i = f(\mathbf{x}_i)\}_{i=1}^n$ . The integral is correlated with the observed data as follows:

$$\text{cov}(\mathbf{f}, \mu) = \mathbb{E} \left\{ (\mathbf{f} - m\mathbf{1}) \left[ \int_{[0,1]^d} [f(\mathbf{x}) - m] d\mathbf{x} \right] \right\} = \mathbf{c},$$

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

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

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

## (6) 2.2.1 Empirical Bayes

Similarly, the conditional probability density of cubature error  $\mu - \hat{\mu}$  given observed data  $\mathbf{y}$ :

$$\mu - \hat{\mu} | \mathbf{y} \sim \mathcal{N} \left( \begin{matrix} -w_0 + m(\mu(1) - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}) + \mathbf{y}^T (\mathbf{C}^{-1} \mathbf{c} - \mathbf{w}) \\ s^2(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \end{matrix} \right) \quad (7)$$

Properly choosing the weights:

$$w_0 = m(\mu(1) - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}), \quad \mathbf{w} = \mathbf{C}^{-1} \mathbf{c},$$

we can force the posterior error  $\mu - \hat{\mu} | \mathbf{y}$  to have zero mean:

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

This leads to the cubature rule as unbiased solution:

$$\hat{\mu}(\mathbf{f} = \mathbf{y}) = w_0 + \mathbf{w}^T \mathbf{y} = m(\mu(1) - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}) + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}, \quad (8)$$

which is also the posterior mean:

$$\mathbb{E}[\mu | (\mathbf{f} = \mathbf{y})] = m(\mu(1) - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c}) + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}. \quad (9)$$

This cubature rule is an affine function of the data.

If we can assume the observed data comes from the source of Gaussian process with zero mean  $m = 0$  then  $w_0 = 0$ ; This fact can be used where applicable. For the integration problem as defined eqn (1), which is the focus of this work,  $\mu(1) = 1$ . So, further in this work, we use the actual value instead.

Subsequently, since the posterior error  $\mu - \hat{\mu} | \mathbf{y}$  is a Normal random variable, we can obtain a confidence interval or inference using integrand-samples and estimated-parameters. If  $n$  is chosen large enough to make

$$\text{err}_n := 2.58 \sqrt{s^2(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})} \leq \varepsilon \quad (10)$$

Then

$$\mathbb{P}_f[|\mu - \hat{\mu}| \leq \varepsilon] \geq 99\% \quad (11)$$

where the constant 2.58 comes from the fact that 99% standard Normal distribution is comprised within 2.58 times the standard deviation. We call  $\text{err}_n$ , error bound and  $\text{err}_n \leq \varepsilon$ , stopping criterion.

As mentioned above, our algorithm assumes that the sample size belongs a sequence of positive integers,  $\mathbb{Z}_{>0}$ . If  $n$  is chosen so that

$$n = \min\{n_j \in \mathbb{Z}_{>0} : 2.58^2 s^2 [c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}] \leq \varepsilon^2, \quad j \in \mathbb{N}\}, \quad (12)$$

then one can be confident that (3) is satisfied. In the above stopping criterion it is understood that  $c_0$ ,  $\mathbf{c}$ , and  $\mathbf{C}$ , all depend on  $n_j$

## 2.2 Parameters estimation

The covariance scale parameter  $s^2$ , mean  $m$  and kernel shape parameter  $\boldsymbol{\theta}$  must be estimated, since the cubature (9) and the credible interval in (11) involve these parameters.

One approach to do this estimation through Maximum likelihood estimation (MLE), by using the observed integrand values for the purpose of estimating the integral. The log-likelihood function of the parameters given the data  $\mathbf{y} = \{f(\mathbf{x}_i)\}_{i=1}^n$  is:

$$\begin{aligned} l(s, m, \boldsymbol{\theta} | \mathbf{y}) &= \log \left( \frac{\exp(-\frac{1}{2} s^{-2} (\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m\mathbf{1}))}{\sqrt{(2\pi)^n \det(s^2 \mathbf{C})}} \right) \\ &= -\frac{1}{2} s^{-2} (\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m\mathbf{1}) \\ &\quad - \frac{1}{2} \log(\det \mathbf{C}) - \frac{n}{2} \log(s^2) + \text{constants}. \end{aligned} \quad (13)$$

Maximizing  $l(s, m, \boldsymbol{\theta} | \mathbf{y})$  with respect to  $m$ , provides the MLE of  $m$ :

$$m_{\text{MLE}} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \quad (14)$$

Then using this result,

$$\begin{aligned} (\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m\mathbf{1}) &= \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} - \frac{\mathbf{y}^T \mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \\ &= \mathbf{y}^T \left[ \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}. \end{aligned}$$

Maximizing  $l(s, m, \boldsymbol{\theta} | \mathbf{y})$  with respect to 's' provides the MLE of  $s$ :

$$\begin{aligned} s_{\text{MLE}}^2 &= \frac{1}{n} (\mathbf{y} - m_{\text{MLE}} \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{MLE}} \mathbf{1}) \\ &= \frac{1}{n} \mathbf{y}^T \left[ \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}. \end{aligned} \quad (15)$$

Plug in the results of  $m_{\text{MLE}}$  and  $s_{\text{MLE}}$  to simplify the log likelihood:

$$l(s, m, \boldsymbol{\theta} | \mathbf{y}) = -\frac{1}{2n} - \frac{1}{2} \log(\det \mathbf{C}) - \frac{n}{2} \log(s_{\text{MLE}}^2) + \text{const}$$

By minimizing the negative of  $l(s, m, \boldsymbol{\theta} | \mathbf{y})$  with respect to ' $\boldsymbol{\theta}$ ', the MLE of  $\boldsymbol{\theta}$  can be obtained. This is usually done by numerically searching for the minimum:

$$\boldsymbol{\theta}_{\text{MLE}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \left[ \frac{1}{n} \log(\det \mathbf{C}) + \log(s_{\text{MLE}}^2) \right]. \quad (16)$$

Using those two MLE results of  $m$  and  $s$ , the simplified stopping criterion to choose  $n$  becomes:

$$\begin{aligned} n &= \min \left\{ n_j \in \mathbb{Z}_{>0} : \right. \\ &\quad \left. \frac{2.58^2}{n_j} \mathbf{y}^T \left[ \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \times (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \leq \varepsilon^2 \right\}. \end{aligned} \quad (17)$$

Finally, using the MLE estimated parameters  $m_{\text{MLE}}$ ,  $s_{\text{MLE}}$  and  $\boldsymbol{\theta}_{\text{MLE}}$ , the cubature rule:

$$\hat{\mu}_{\text{MLE}}(f) = \left( \frac{(1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c})}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \mathbf{C}^{-1} \mathbf{1} + \mathbf{C}^{-1} \mathbf{c} \right)^T \mathbf{y} \quad (18)$$

### 2.2.2 Full Bayes

Full Bayes approach includes  $m, s$  as hyper-parameters while estimating the distribution of cubature error. Then assumes conjugate for the hyper-parameters namely  $\rho_{m,s^2}(\xi, \lambda) \propto 1/\lambda$ . Then the posterior density for the integral given the data using the Bayes theorem is

$$\begin{aligned} \rho_\mu(z|\mathbf{f} = \mathbf{y}) &= \int_0^\infty \int_{-\infty}^\infty \rho_\mu(z|\mathbf{f} = \mathbf{y}, m = \xi, s^2 = \lambda) \rho_{m,s^2}(\xi, \lambda|\mathbf{f} = \mathbf{y}) d\xi d\lambda \\ &\quad \text{by the properties of conditional probability} \\ &\propto \int_0^\infty \int_{-\infty}^\infty \rho_\mu(z|\mathbf{f} = \mathbf{y}, m = \xi, s^2 = \lambda) \rho_{\mathbf{f}}(\mathbf{y}|\xi, \lambda) \rho_{m,s^2}(\xi, \lambda) d\xi d\lambda \\ &\quad \text{by Bayes' Theorem} \\ &\propto \int_0^\infty \frac{1}{\lambda^{(n+3)/2}} \int_{-\infty}^\infty \exp\left(-\frac{1}{2\lambda} \left\{ \frac{[z - \xi(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}]^T \mathbf{C}^{-1} [z - \xi(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) + \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}]}{[c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}]} \right. \right. \\ &\quad \left. \left. + (\mathbf{y} - \xi \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - \xi \mathbf{1}) \right\} \right) d\xi d\lambda \\ &\quad \text{by (5b), (6) and } \rho_{m,s^2}(\xi, \lambda) \propto 1/\lambda \\ &\propto \int_0^\infty \frac{1}{\lambda^{(n+3)/2}} \int_{-\infty}^\infty \exp\left(-\frac{\alpha \xi^2 - 2\beta \xi + \gamma}{2\lambda(c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})}\right) d\xi d\lambda, \end{aligned}$$

where

$$\begin{aligned} \alpha &= (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2 + \mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}), \\ \beta &= (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})(z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}) + \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}), \\ \gamma &= (z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y})^2 + \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}). \end{aligned}$$

In the derivation above and below, factors that are independent of  $\xi, \lambda$ , or  $z$  can be discarded since we only need to preserve the proportion, however, factors that depend on  $\xi, \lambda$ , 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 compute the integral with respect to  $\xi$ :

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

Finally, we simplify the key term:

$$\begin{aligned} \alpha\gamma - \beta^2 &= \mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) (z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y})^2 \\ &\quad - 2\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) (z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y}) \\ &\quad + (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2 \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \\ &\quad + [\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} - (\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y})^2] (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})^2 \\ &\propto \mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} \left( z - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{y} - \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right)^2 \\ &\quad - \frac{[(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}]^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2 \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} \\ &\quad (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) [\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} - (\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y})^2] \\ &\propto \left( z - \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) \mathbf{1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + \mathbf{c} \right]^T \mathbf{C}^{-1} \mathbf{y} \right)^2 \\ &\quad + \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \right] \times \mathbf{y}^T \left[ \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}. \end{aligned}$$

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 cubature rule is the same as in the empirical Bayes case in (18):

$$\begin{aligned} \mu_{\text{full}}|(\mathbf{f} = \mathbf{y}) &= \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) \mathbf{1}^T}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + \mathbf{c}^T \right] \mathbf{C}^{-1} \mathbf{y}, \quad (19) \\ \hat{\sigma}_{\text{full}}^2|(\mathbf{f} = \mathbf{y}) &= \frac{1}{n - 1} \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \right] \\ &\quad \times \mathbf{y}^T \left[ \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}. \end{aligned} \quad (20)$$

But the variance is different. So, the stopping criterion for the full Bayes case,

$$\begin{aligned} n &= \min \left\{ n_j \in \mathbb{Z}_{>0} : \frac{t_{n_j-1,0.995}^2}{n_j - 1} \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \right] \right. \\ &\quad \left. \times \mathbf{y}^T \left[ \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y} \leq \varepsilon^2 \right\}. \quad (21) \end{aligned}$$

is more conservative than that in the empirical Bayes case, (17). 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,  $\theta$ , enters the definition of the covariance kernel in a non-trivial way, it typically cannot be practically considered as a hyper-parameter in a full Bayesian model. The empirical Bayes approach for determining  $\theta$  in (16) is one possibility

### 2.2.3 Cross-validation

Another alternative to determining  $m, s$ , and  $\theta$  is *leave-one-out cross-validation*. Let  $\hat{y}_i$  denote the expected

value of  $f(\mathbf{x}_i)$  conditioned on  $\mathbf{f}_{-i} = \mathbf{y}_{-i}$ , where the subscript  $-i$  denotes vector excluding the  $i^{\text{th}}$  component. The cross-validation criterion, which is to be minimized, is

$$\text{CV} = \sum_{i=1}^n (y_i - \tilde{y}_i)^2. \quad (22)$$

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

$$\mathbf{C} = \begin{pmatrix} C_{ii} & \mathbf{C}_{-i,i}^T \\ \mathbf{C}_{-i,i} & \mathbf{C}_{-i,-i} \end{pmatrix}, \quad \mathbf{A} = \begin{pmatrix} A_{ii} & \mathbf{A}_{-i,i}^T \\ \mathbf{A}_{-i,i} & \mathbf{A}_{-i,-i} \end{pmatrix}, \quad \boldsymbol{\zeta} = \begin{pmatrix} \zeta_i \\ \boldsymbol{\zeta}_{-i} \end{pmatrix}$$

Note that in this notation, Lemma 1 implies that

$$\tilde{y}_i = m + \mathbf{C}_{-i,i}^T \mathbf{C}_{-i,-i}^{-1} (\mathbf{y}_{-i} - m\mathbf{1})$$

$$\zeta_i = A_{ii}(y_i - m) + \mathbf{A}_{-i,i}^T (\mathbf{y}_{-i} - m\mathbf{1})$$

$$= A_{ii}[(y_i - m) - \mathbf{C}_{-i,i}^T \mathbf{C}_{-i,-i}^{-1} (\mathbf{y}_{-i} - m\mathbf{1})]$$

$$= A_{ii}(y_i - \tilde{y}_i).$$

Thus, (24) may be re-written as

$$\text{CV} = \sum_{i=1}^n \left( \frac{\zeta_i}{A_{ii}} \right)^2, \quad \boldsymbol{\zeta} = \mathbf{C}^{-1}(\mathbf{y} - m\mathbf{1}) \quad (23)$$

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

$$\text{GCV} = \frac{\sum_{i=1}^n \zeta_i^2}{\left( \frac{1}{n} \sum_{i=1}^n A_{ii} \right)^2} = \frac{(\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-2} (\mathbf{y} - m\mathbf{1})}{\left( \frac{1}{n} \text{trace}(\mathbf{C}^{-1}) \right)^2}. \quad (24)$$

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

$$m_{\text{GCV}} = \frac{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}}, \quad (25)$$

$$\boldsymbol{\theta}_{\text{GCV}} = \underset{\boldsymbol{\theta}}{\text{argmin}} \log \left( \mathbf{y}^T \left[ \mathbf{C}^{-2} - \frac{\mathbf{C}^{-2} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-2}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}} \right] \mathbf{y} \right) - 2 \log (\text{trace}(\mathbf{C}^{-1})) \quad (26)$$

### 2.3 Formulating the automatic Bayesian cubature algorithm

Let the function to integrate be  $f \sim \mathcal{GP}(m, s^2 C_{\boldsymbol{\theta}})$ . Our goal is to compute  $\hat{\mu}$  within the error tolerance  $\varepsilon$ , 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  $\text{err}_n$  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  $\mu$  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 *auto-*

*matic cubature algorithm* (1), we only need to numerically estimate the MLE of  $\boldsymbol{\theta}$  and then use it to compute  $\mathbf{C}$  and  $\text{err}_n$ .

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

---

#### Algorithm 1

---

```

1: procedure AUTOCUBATURE( $f, \varepsilon$ )  $\triangleright$  Integrate within the
    $\zeta_i$  error tolerance
2:  $n_0 \leftarrow 2^8$   $\triangleright$  start with minimum number of points
3:  $n \leftarrow n_0, n' \leftarrow 0$ 
4: while true do  $\triangleright$  Iterate till error tolerance is met
5:   Generate  $\{\mathbf{x}_i\}_{i=n'+1}^n$  and sample  $\{f(\mathbf{x}_i)\}_{i=n'+1}^n$ 
6:   Compute error bound  $\text{err}_n$   $\triangleright$   $\text{err}_n$  data driven
   error bound
7:   if  $\text{err}_n \leq \varepsilon$  then break
8:   end if
9:    $n' \leftarrow n, n \leftarrow 2 \times n'$ 
10: end while
11: Compute cubature weights  $\{w_i\}_{i=1}^n$ 
12: Compute approximate integral  $\hat{\mu}_n$ 
13: return  $\hat{\mu}_n$   $\triangleright$  Integral estimate  $\hat{\mu}_n$ 
14: end procedure

```

---

As shown, the algorithm continues the iteration loop till the stopping-criterion is met, i.e.,  $\text{err}_n$  is smaller than error threshold  $\varepsilon$ . At the end of every iteration, if the computed error bound  $\text{err}_n$  is higher than the required  $\varepsilon$ , algorithm doubles the number of points  $n$  and repeats. When the error tolerance  $\varepsilon$  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  $\text{err}_n$  which closely matches the true error is essential for the effectiveness of our algorithm. So, accurate and faster computation of  $\text{err}_n$  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:

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

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  $\boldsymbol{\theta}$  which needs the loss function, is the most time consuming of all. Because the loss

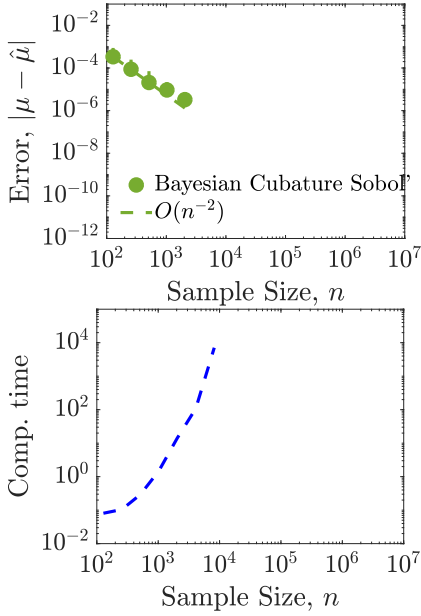


Fig. 1: MVN for  $d=2$  with Matern kernel

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  $\text{err}_n$  to decide when to stop, and uses MLE loss function (16) 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* [Fas07]. 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  $C_\theta$  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,

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

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

$$\begin{aligned} \mathbf{C} \mathbf{V}_1 &= \lambda_1 \mathbf{V}_1, \quad \text{by definition} \\ \text{So, } \mathbf{C}^{-1} \mathbf{1} &= \mathbf{C}^{-1} \mathbf{V}_1 = \frac{\mathbf{V}_1}{\lambda_1} = \frac{\mathbf{1}}{\lambda_1}, \end{aligned}$$

$$\begin{aligned} \text{Similarly, } \mathbf{a}^T \mathbf{C}^{-1} \mathbf{b} &= \frac{1}{n} \mathbf{a}^T \mathbf{V} \mathbf{\Lambda}^{-1} \mathbf{V}^H \mathbf{b} = \frac{1}{n} \hat{\mathbf{a}}^T \mathbf{\Lambda}^{-1} \hat{\mathbf{b}}^* = \frac{1}{n} \sum_{i=1}^n \frac{\hat{a}_i \hat{b}_i^*}{\lambda_i}, \\ \mathbf{a}^T \mathbf{C}^{-2} \mathbf{b} &= \frac{1}{n} \mathbf{a}^T \mathbf{V} \mathbf{\Lambda}^{-2} \mathbf{V}^H \mathbf{b} = \frac{1}{n} \hat{\mathbf{a}}^T \mathbf{\Lambda}^{-2} \hat{\mathbf{b}}^* = \frac{1}{n} \sum_{i=1}^n \frac{\hat{a}_i \hat{b}_i^*}{\lambda_i^2}, \end{aligned} \quad (31)$$

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,

$$\begin{aligned} \mathbf{1}^T \mathbf{C}^{-1} \mathbf{1} &= \mathbf{1}^T \left( \frac{\mathbf{1}}{\lambda_1} \right) = \frac{n}{\lambda_1} \\ \mathbf{1}^T \mathbf{C}^{-1} \mathbf{y} &= \mathbf{1}^T \left( \frac{\mathbf{y}}{\lambda_1} \right) = \frac{\sum_{i=1}^n y_i}{\lambda_1} = \frac{\hat{y}_1}{\lambda_1} \\ \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} &= \frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i} \\ \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1} &= \mathbf{c}^T \left( \frac{\mathbf{1}}{\lambda_1} \right) = \frac{\hat{c}_1}{\lambda_1} \\ \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c} &= \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \\ \mathbf{y}^T \mathbf{C}^{-2} \mathbf{y} &= \frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} \end{aligned}$$

Then the MLE estimates can be simplified using these results:

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

### 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

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

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

Similarly, the error bound  $\text{err}_n$  computation can be made faster:

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

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:

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

It is interesting to note that, with  $m \neq 0$  assumption,  $\hat{\mu}_{\text{MLE}}$  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  $\text{err}_n$ .

#### 3.2.2 Full Bayes

$$\begin{aligned} \boldsymbol{\theta}_{\text{GCV}} &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[ \log \left( \frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} - \frac{\hat{y}_1^2 / \lambda_1^4}{n / \lambda_1^2} \right) - 2 \log \left( \sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \\ &= \underset{\boldsymbol{\theta}}{\text{argmin}} \left[ \log \left( \frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left( \sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \end{aligned} \quad (33)$$

$$\hat{\mu}_{\text{full}}(\mathbf{f} = \mathbf{y}) = \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}) \mathbf{1}^T}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + \mathbf{c}^T \right] \mathbf{C}^{-1} \mathbf{y} = \frac{\hat{y}_1}{n} + \frac{1}{n} \sum_{i=2}^n \frac{\hat{c}_i \hat{y}_i^*}{\lambda_i}, \quad (34)$$

$$\hat{\sigma}_{\text{full}}^2(\mathbf{f} = \mathbf{y}) = \frac{1}{n-1} \left[ \frac{(1 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1})^2}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} + (c_0 - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}) \right] \quad (35)$$

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

The stopping criterion for the full Bayes case,

$$n = \min \left\{ n_j \in \mathbb{Z}_{>0} : \frac{t_{n_j-1,0.995}^2}{n_j(n_j-1)} \left[ \frac{\lambda_1}{n_j} \left( 1 - \frac{\hat{c}_1}{\lambda_1} \right)^2 + \left( c_0 - \frac{1}{n_j} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \times \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \leq \varepsilon^2 \right] \right\}. \quad (36)$$

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:

$$C_{\theta}(\mathbf{x}, \mathbf{t}) := \sum_{\mathbf{k} \in \mathbb{Z}^d} \alpha_{\mathbf{k}, \theta} e^{2\pi\sqrt{-1}\mathbf{k}^T \mathbf{x}} e^{-2\pi\sqrt{-1}\mathbf{k}^T \mathbf{t}}, \quad (37)$$

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

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:

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

The *shape parameter*  $\theta$  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 (37), 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 [SD06]:

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

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

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

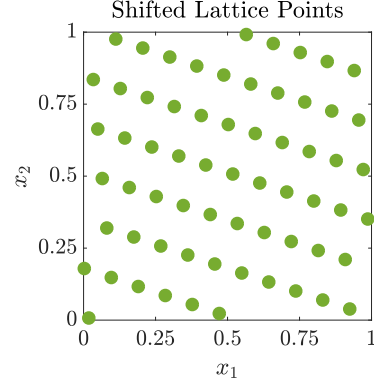


Fig. 2: Randomly shifted Lattice points in d=2

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 (37), 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 (37) 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 [DLM12] of Bernoulli polynomial  $B_r$  to find the appropriate  $\alpha_{\mathbf{k}}$ . 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  $\alpha_{\mathbf{k}}$  in (37) decay. The following very useful expansion is referenced from [DLM12] under eqn (24.8.3):

$$B_r(x) = \frac{-r!}{(2\pi\sqrt{-1})^r} \sum_{\substack{k \neq 0, \\ k=-\infty}}^{\infty} \frac{e^{2\pi\sqrt{-1}kx}}{k^r} \begin{cases} \text{for } r = 1, & 0 < x < 1 \\ \text{for } r = 2, 3, \dots & 0 \leq x \leq 1 \end{cases} \quad (38)$$

Our goal is to replace the infinite sum of the kernel using Bernoulli polynomials. This would allows the



computations to be carried out in any software like Matlab. For 1-D ( $d = 1$ ) rewriting (38):

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

Comparing this expansion with (37), 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  $\alpha_k$  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  $\alpha_k$  for  $C_{\theta}(\mathbf{x}, \mathbf{t})$  is computed using:

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

Where  $\theta$  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:

$$C_{\theta}(\mathbf{x}, \mathbf{t}) = \prod_{l=1}^d 1 - \theta_l^r \frac{(2\pi\sqrt{-1})^r}{r!} B_r(|x_l - t_l|), \quad 0 \leq |x_l - t_l| \leq 1, \quad (39)$$

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

### 4.3 Eigenvalues and Eigenvectors of C

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

$$(C_{\theta}(\mathbf{x}_1, \mathbf{x}_1), C_{\theta}(\mathbf{x}_2, \mathbf{x}_1), \dots, C_{\theta}(\mathbf{x}_n, \mathbf{x}_1))^T =: \mathbf{C}_1, \quad (40)$$

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:

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

The matrix constructed with these eigenvectors is:

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

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

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:

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

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

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

This gives simplified error bound:

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

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

And the cubature:

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

#### 4.4.1 Full Bayes

$$\hat{\mu}_{\text{full}}(\mathbf{f} = \mathbf{y}) = \frac{\hat{y}_1}{n}, \quad (41)$$

$$\hat{\theta}_{\text{GCV}} = \underset{\theta}{\text{argmin}} \left[ \log \left( \frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left( \sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \quad (42)$$

$$\hat{\sigma}_{\text{full}}^2(\mathbf{f} = \mathbf{y})$$

$$\begin{aligned} &= \frac{1}{n-1} \left[ \frac{\lambda_1}{n} \left( 1 - \frac{\hat{c}_1}{\lambda_1} \right)^2 + \left( c_0 - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \right] \times \frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \\ &= \frac{1}{n-1} \left( 1 - \frac{n}{\lambda_1} \right) \left[ \frac{\lambda_1}{n} \left( 1 - \frac{n}{\lambda_1} \right) + 1 \right] \times \frac{1}{n} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \\ &= \frac{1}{n(n-1)} \left( \frac{\lambda_1}{n} - 1 \right) \times \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \end{aligned}$$

The stopping criterion for the full Bayes case,

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

#### 4.5 Variable transforms

Fast transform kernel technique discussed with shift invariant kernel eqn(37) 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.

$$\text{Baker's : } \tilde{f}(\mathbf{t}) = f \left( \left( 1 - 2 \left| t_j - \frac{1}{2} \right| \right)_{j=1}^d \right) \quad (44)$$

$$\text{C0 : } \tilde{f}(\mathbf{t}) = f(\mathbf{g}_0) \prod_{j=1}^d g'_0(t_j), \quad \mathbf{g}_0 = (g_0(t_j))_{j=1}^d \quad (45)$$

$$\text{C1 : } \tilde{f}(\mathbf{t}) = f(\mathbf{g}_1) \prod_{j=1}^d g'_1(t_j), \quad \mathbf{g}_1 = (g_1(t_j))_{j=1}^d \quad (46)$$

$$\text{Sidi's C1 : } \tilde{f}(\mathbf{t}) = f(\boldsymbol{\psi}_2) \prod_{j=1}^d \psi'_2(t_j), \quad \boldsymbol{\psi}_2 = (\psi_2(t_j))_{j=1}^d \quad (47)$$

$$\text{Sidi's C2 : } \tilde{f}(\mathbf{t}) = f(\boldsymbol{\psi}_3) \prod_{j=1}^d \psi'_3(t_j), \quad \boldsymbol{\psi}_3 = (\psi_3(t_j))_{j=1}^d \quad (48)$$

where

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

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

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

$$\psi_3(t) = \frac{1}{16} (8 - 9 \cos(\pi t) + \cos(3\pi t)), \quad \psi'_3(t) = \frac{1}{16} (9 \sin(\pi t) - \sin(3\pi t))$$

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  $j = 2l + 1$ :

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

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:

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

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  $\lambda_1$  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  $\tilde{\lambda}_1 = \lambda_1 - n$ . In fact for the rest of the values are:

$$\tilde{\lambda}_j = \lambda_j, \forall j = 1, \dots, n-1. \quad (49)$$

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. (49).

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

Thus proven. Using this result, cancellation error in the computation of  $\text{err}_n$  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:

$$d = 1: C_1 = 1 + \theta B(\{x_{i_1} - x_{j_1}\}) = 1 + \tilde{C}_1$$

$$\begin{aligned} d = 2: C_2 &= [1 + \theta B(\{x_{i_2} - x_{j_2}\})][1 + \tilde{C}_1] \\ &= 1 + \theta B(\{x_{i_2} - x_{j_2}\})[1 + \tilde{C}_1] + \tilde{C}_1 \\ &= 1 + \underbrace{\theta B(\{x_{i_2} - x_{j_2}\})C_1}_{\tilde{C}_2} + \tilde{C}_1 \end{aligned}$$

$$\begin{aligned} d = 3: C_3 &= [1 + \theta B(\{x_{i_3} - x_{j_3}\})][1 + \tilde{C}_2] \\ &= 1 + \theta B(\{x_{i_3} - x_{j_3}\})[1 + \tilde{C}_2] + \tilde{C}_2 \\ &= 1 + \theta B(\{x_{i_3} - x_{j_3}\})C_2 + \tilde{C}_2 \end{aligned}$$

$\vdots$

$$\begin{aligned} \forall d > 2: C_d &= [1 + \theta B(\{x_{i_d} - x_{j_d}\})][1 + \tilde{C}_{d-1}] \\ &= 1 + \theta B(\{x_{i_d} - x_{j_d}\})[1 + \tilde{C}_{d-1}] + \tilde{C}_{d-1} \\ &= 1 + \theta B(\{x_{i_d} - x_{j_d}\})C_{d-1} + \tilde{C}_{d-1} \end{aligned}$$

## 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  $m, s, \theta$  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,

$$\mathbf{f} = (f(\mathbf{x}_i))_{i=1}^n \sim \mathcal{N}(m\mathbf{1}, \mathbf{C}), \quad \text{where } \mathbf{C} = \frac{1}{n} \mathbf{V} \mathbf{V}^H, \quad \mathbf{V}^H \mathbf{V} = n$$

Then,

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

And,

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

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

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

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.

## 4.9 Numerical Experiments

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

### 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,

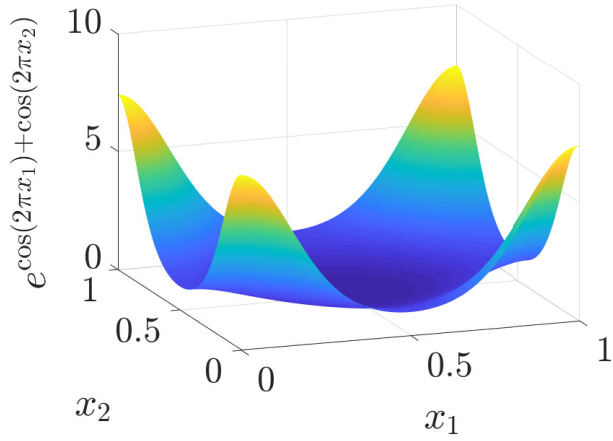
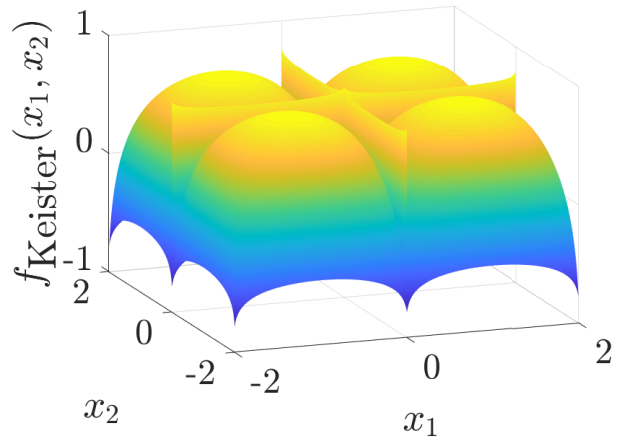
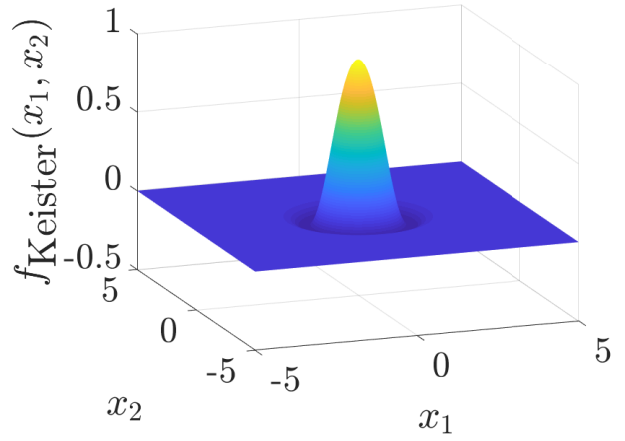


Fig. 3: Exp(Cos) in 2 dimensions

Fig. 4: Keister function in original form and its transform to  $[0, 1]^2$ 

where the function is defined as

$$f(\mathbf{x}) = e^{\sum_{i=1}^d \cos(2\pi x_i)},$$

$$\int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} = \text{Bessell}(0,1)^d$$

where ‘Bessell’ 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 [Kei96], inspired by a physics application.

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

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

where

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

$$\text{sininteg}(1) = \int_{x=0}^{\infty} \exp(-\mathbf{x}^T \mathbf{x}) \sin(\mathbf{x}) d\mathbf{x} = 0.424431835029215$$

$$\text{cosinteg}(2) = \frac{1 - \text{sininteg}(1)}{2}$$

$$\text{sininteg}(2) = \frac{\text{cosinteg}(1)}{2}$$

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

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

where  $\text{gamma}(\cdot) := \text{gamma function}$

### 3. Multivariate Normal:

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.

$$\mu = \int_{[\mathbf{a}, \mathbf{b}] \in \mathbb{R}^d} \frac{\exp(-\frac{1}{2} \mathbf{t}^T \Sigma^{-1} \mathbf{t})}{\sqrt{(2\pi)^d \det(\Sigma)}} d\mathbf{t} \stackrel{[\text{Gen93}]}{=} \int_{[0,1]^{d-1}} f_{\text{Genz}}(\mathbf{x}) d\mathbf{x}$$

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

$$\boldsymbol{\alpha}_1 = \Phi(a_1), \quad \boldsymbol{\beta}_1 = \Phi(b_1)$$

$$\boldsymbol{\alpha}_j(x_1, \dots, x_{j-1}) = \Phi\left(\frac{1}{l_{jj}} \left(a_j - \sum_{k=1}^{j-1} l_{jk} \Phi^{-1}(\boldsymbol{\alpha}_k + x_k(\boldsymbol{\beta}_k - \boldsymbol{\alpha}_k))\right)\right)$$

$$\boldsymbol{\beta}_j(x_1, \dots, x_{j-1}) = \Phi\left(\frac{1}{l_{jj}} \left(b_j - \sum_{k=1}^{j-1} l_{jk} \Phi^{-1}(\boldsymbol{\alpha}_k + x_k(\boldsymbol{\beta}_k - \boldsymbol{\alpha}_k))\right)\right)$$

$$f_{\text{Genz}}(\mathbf{x}) = \prod_{j=1}^d [\boldsymbol{\beta}_j(\mathbf{x}) - \boldsymbol{\alpha}_j(\mathbf{x})]$$

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  $\text{payoff}(\cdot)$  describes 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

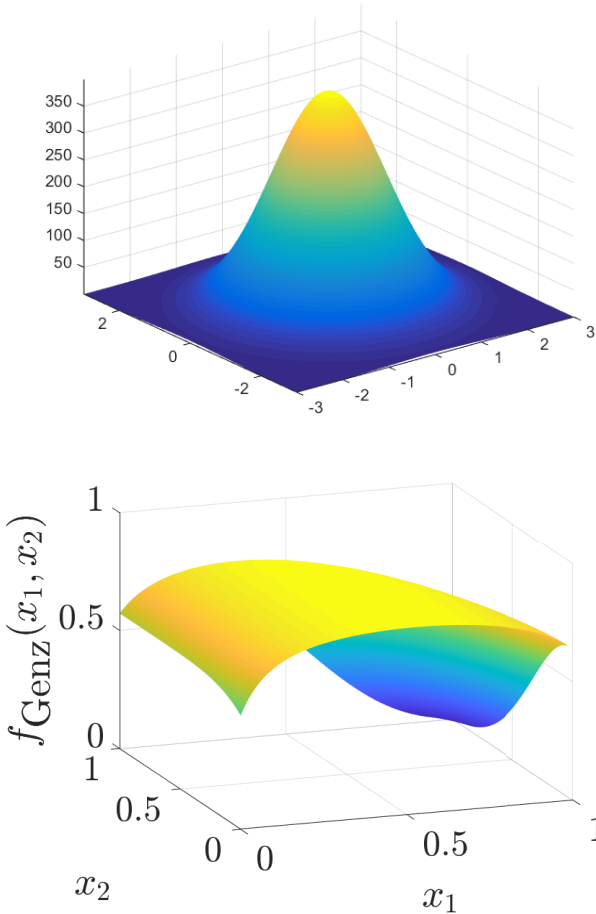


Fig. 5: Multivariate Normal and Genz function

$$\text{payoff}(z) = \max \left( \frac{1}{d} \sum_{j=1}^d S_j - K, 0 \right) e^{-rt}, \quad S_j = S_0 \exp((r - \beta_1 g_1 - \dots - \beta_p g_p) t)$$

where  $d$  - number of dimensions and  $T, S, S_0, K, r, \sigma$  are the parameters to be specified

#### 4.10 Results and observation

### 5 Conclusion

We developed a generalized technique of *Fast transform kernel*. Using this technique, further developed 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  $\text{err}_n$  still depends on the kernel order. So when we use higher order  $r$ , the error bound  $\text{err}_n$  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  $\text{err}_n$  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  $f(\mathbf{x}) = \beta_0 + \beta_1 g_1 + \dots + \beta_p g_p$  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

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

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

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

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

Finally, the function approximation is

$$\begin{aligned} \tilde{f}(\mathbf{x}) &= \tilde{g}(\psi(\mathbf{x})) \\ &= \sum w_i C(\cdot, \cdot) \end{aligned}$$

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

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

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* : Note that

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

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

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

Let's denote the first partition  $\mathbf{Y}_1$  and the second  $\mathbf{Y}_2$ . 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

$$\mathbb{E}(\mathbf{Y}_1 | \mathbf{Y}_2) = \mathbf{m}_1 + \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} (\mathbf{y}_2 - \mathbf{m}_2)$$

$$\text{var}(\mathbf{Y}_1 | \mathbf{Y}_2) = \mathbf{C}_{11} - \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{21}.$$

Now define  $\mathbf{z} = \mathbf{Y}_1 + \mathbf{A} \mathbf{Y}_2$  where  $\mathbf{A} = -\mathbf{C}_{21}^T \mathbf{C}_{22}^{-1}$ .

We can write

$$\begin{aligned} \text{cov}(\mathbf{z}, \mathbf{Y}_2) &= \mathbf{C}_{21}^T + \text{cov}(\mathbf{A} \mathbf{Y}_2, \mathbf{Y}_2) \\ &= \mathbf{C}_{21}^T + \mathbf{A} \text{var}(\mathbf{Y}_2) \\ &= \mathbf{C}_{21}^T - \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} \mathbf{C}_{22} = 0 \end{aligned}$$

Therefore  $z$  and  $\mathbf{Y}_2$  are uncorrelated and, since they are jointly normal, they are independent. Now, clearly

$\mathbb{E}(z) = m_1 + A m_2$ , therefore it follows that

$$\begin{aligned}\mathbb{E}(\mathbf{Y}_1|\mathbf{Y}_2) &= \mathbb{E}(z - A\mathbf{Y}_2|\mathbf{Y}_2) \\ &= \mathbb{E}(z|\mathbf{Y}_2) - \mathbb{E}(A\mathbf{Y}_2|\mathbf{Y}_2) = \mathbb{E}(z) - A\mathbf{Y}_2 \\ &= m_1 + A(m_2 - \mathbf{Y}_2) = m_1 + \mathbf{C}_{21}^T \mathbf{C}_{22}^{-1} (\mathbf{Y}_2 - m_2)\end{aligned}$$

which proves the mean  $\mathbb{E}(\mathbf{Y}_1|\mathbf{Y}_2)$ . For the covariance matrix, note that

$$\begin{aligned}\text{var}(\mathbf{Y}_1|\mathbf{Y}_2) &= \text{var}(z - A\mathbf{Y}_2|\mathbf{Y}_2) \\ &= \text{var}(z|\mathbf{Y}_2) + \text{var}(A\mathbf{Y}_2|\mathbf{Y}_2) - A\text{cov}(z, \mathbf{Y}_2) - \text{cov}(z, \mathbf{Y}_2)A^T \\ &= \text{var}(z|\mathbf{Y}_2) \\ &= \text{var}(z)\end{aligned}$$

So, it is enough to show:

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

which proves the variance.

## References

- [Dia88] P. Diaconis. Bayesian numerical analysis. *Statistical decision theory and related topics IV, Papers from the 4th Purdue symp., West Lafayette, Indiana 1986,*, page 163–175, 1988.
- [DLM12] DLMF. Nist digital library of mathematical functions. <http://dlmf.nist.gov/>, Part 2:Release 1.0.5 of 2012–10–01, 2012.
- [Fas07] G. E. Fasshauer. *Meshfree Approximation Methods with Matlab*. World Scientific Publishing Co, 2007.
- [Gen93] A. Genz. Comparison of methods for the computation of multivariate normal probabilities. *Computing Science and Statistics*, 25:400 – 405, 1993.
- [Hic17] Fred J. Hickernell. Error analysis for quasi-monte carlo methods. *arXiv:1702.01487 [math.NA]*, 2017.
- [Kei96] B. D. Keister. Multidimensional quadrature algorithms. *Computers in Physics*, 10:119–122, 1996.
- [O’H91] A. O’Hagan. Bayes-hermite quadrature. *J. Statist. Plann. Inference*, 29:245–260, 1991.
- [RG03] C. E. Rasmussen and Z. Ghahramani. Bayesian monte carlo. *Advances in Neural Information Processing Systems*, pages 489–496, 2003.
- [Rit00] K. Ritter. Average-case analysis of numerical problems. *Lecture Notes in Mathematics, Springer-Verlag, Berlin*, vol. 1733:163–175, 2000.
- [SD06] Alexander Keller Sabrina Dammertz. Image synthesis by rank-1 lattices. *Monte Carlo and Quasi-Monte Carlo Methods, Part 2*:217–236, 2006.