

Fast Automatic Bayesian Cubature using Matching Kernels and Designs

Jagadeeswaran R.

Adviser: Prof. Fred J Hickernell

Department of Applied Mathematics, Illinois Institute of Technology
`jrathin1@iit.edu`

Thesis Defense • Oct 22, 2019

Contents

1 Introduction

2 Bayesian Cubature

3 Faster

4 Lattice Nodes

5 Sobol' Nets

6 Demonstration

7 Conclusion



Numerical Integration

A fundamental problem in various fields, including finance, machine learning and statistics,

$$\mu = \int_{\mathbb{R}^d} g(\mathbf{x}) \, d\mathbf{x} = \int_{[0,1]^d} f(\mathbf{x}) \, d\mathbf{x} = \mathbb{E}[f(\mathbf{X})], \quad \text{where } \mathbf{X} \sim \mathcal{U}[0,1]^d \quad (1)$$

by a cubature rule $\hat{\mu}_n := w_0 + \sum_{j=1}^n f(\mathbf{x}_j)w_j$

using points $\{\mathbf{x}_j\}_{j=1}^n$ and associated weights w_j .

The goal of this work is to

- Develop an automatic algorithm for integration
- Use an **extensible** point-set and an algorithm that allows extending points
- Determine ***n*** such that, given ***ε***, $|\mu - \hat{\mu}_n| \leq \epsilon$



Bayesian Cubature

The Bayesian approach for numerical analysis was popularized by Diaconis (1988). Diaconis interpreted the most well known numerical methods, such as 1) trapezoidal rule and 2) splines, from the statistical point of view.

Bayesian approach for numerical integration that is known as Bayesian cubature as introduced by O'Hagan (1991).

- Assume f is drawn from a Gaussian process
 - Need to estimate the mean and Covariance kernel
- Parameter estimation (Empirical, Cross validation) is expensive in general
 - Use points and kernel for which it is cheap

Problem:

- How to choose $\{\mathbf{x}_i\}_{i=1}^n$, and $\{\mathbf{w}_i\}_{i=1}^n$ to make $|\mu - \hat{\mu}_n|$ small? what is err_{CI} ? (Bayesian posterior error)
- How to find n such that $|\mu - \hat{\mu}_n| \leq \text{err}_{\text{CI}} \leq \epsilon$? (automatic cubature)



Bayesian posterior error

Assume random $f \sim \mathcal{GP}(m, s^2 C_\theta)$, a **Gaussian process** with mean m and covariance kernel $s^2 C_\theta$, $C_\theta : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$.

Define $c_0 = \int_{[0,1] \times [0,1]} C_\theta(x, t) dx dt$,

$$\mathbf{c} = \left(\int_{[0,1]} C_\theta(x_i, t) dt \right)_{i=1}^n, \quad \mathbf{C}_\theta = \left(C_\theta(x_i, x_j) \right)_{i,j=1}^n$$

$$\mu - \hat{\mu}_n | \mathbf{y} \sim \mathcal{N}\left(0, s^2(c_0 - \mathbf{c}^T \mathbf{C}_\theta^{-1} \mathbf{c})\right) \quad \text{where } \mathbf{y} = (f(x_i))_{i=1}^n.$$

Choosing $\hat{\mu}_n = \underbrace{m(1 - \mathbf{1}^T \mathbf{C}_\theta^{-1} \mathbf{c})}_{w_0} + (\underbrace{\mathbf{C}_\theta^{-1} \mathbf{c}}_w)^T \mathbf{y}$ makes error unbiased.

If $m = 0$, fixed, choosing $w = \mathbf{C}_\theta^{-1} \mathbf{c}$ makes error unbiased. Please note m, s and θ needs to be inferred.

Diaconis (1988), O'Hagan (1991), Ritter (2000), Rasmussen and Williams (2003), Briol et al. (2018+), Traub et al. (1988) and others



Parameter estimation - Empirical Bayes

Maximize the log-likelihood of the parameters given the data $\mathbf{y} = (f(\mathbf{x}_i))_{i=1}^n$ w.r.t. m then s^2 , then with θ :

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

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

$$\theta_{\text{EB}} = \underset{\theta}{\operatorname{argmin}} \log \left(\frac{1}{2n} \log(\det \mathbf{C}_{\theta}) + \log(s_{\text{EB}}) \right) \quad (\text{numeric})$$

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

The credible interval width, err_{EB} , is given by

$$\text{err}_{\text{EB}} = 2.58 s_{\text{EB}} \sqrt{c_0 - \mathbf{c}^T \mathbf{C}_{\theta}^{-1} \mathbf{c}}$$



1: **procedure** AUTOBAYESCUBATURE(f, ϵ)

Require: a generator for the sequence x_1, x_2, \dots ; a black-box function, f ; an absolute error tolerance, $\epsilon > 0$; the positive initial sample size, n_0 ; the maximum sample size n_{\max}

2: $n \leftarrow n_0, n' \leftarrow 0, \text{err}_n \leftarrow \infty$

3: **while** $\text{err}_{\text{CI}} > \epsilon$ and $n \leq n_{\max}$ **do**

4: Generate $\{x_i\}_{i=n'+1}^n$ and sample $\{f(x_i)\}_{i=n'+1}^n$,

5: Compute parameters, compute error bound err_{CI}

6: $n' \leftarrow n, n \leftarrow 2 \times n'$

7: **end while**

8: Sample size to compute $\hat{\mu}, n \leftarrow n'$

9: Compute approximate $\hat{\mu}_n$, the approximate integral

10: **return** $\hat{\mu}_n$ ▷ Integral estimate $\hat{\mu}_n$

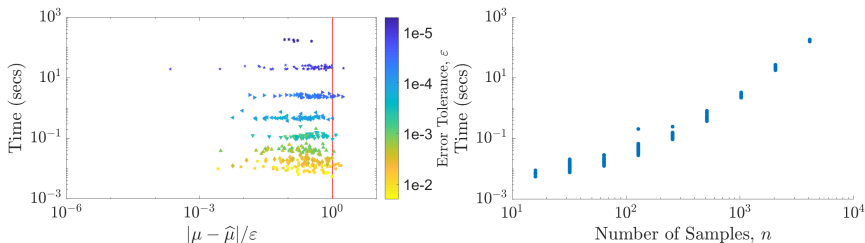
11: **end procedure**



Multivariate Gaussian integration with Matérn kernel

$$\text{Matérn covariance kernel: } C_{\theta}(x, t) = \prod_{\ell=1}^d \exp(-\theta|x_{\ell} - t_{\ell}|)(1 + \theta|x_{\ell} - t_{\ell}|) \quad (2)$$

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



Problem: Computation time (in seconds) increases rapidly, so it is not practical to use more than 4000 points in the cubature.



Fast Bayesian Transform

Choose the kernel C_θ and $\{x_i\}_{i=1}^n$, so the Gram matrix $C_\theta = (C_\theta(x_i, x_j))_{i,j=1}^n$ has:

$$C_\theta = (C_1, \dots, C_n) = \frac{1}{n} V \Lambda V^H, \quad V = (V_1, \dots, V_n) = (v_1, \dots, v_n)^T$$

C_θ is a fast Bayesian transform kernel, if

$$V \text{ may be identified analytically,} \quad (3a)$$

$$v_1 = V_1 = \mathbf{1}, \quad (3b)$$

$$\text{Computing } \tilde{\mathbf{b}} = V^H \mathbf{b} \text{ requires only } \mathcal{O}(n \log n) \text{ operations } \forall \mathbf{b}. \quad (3c)$$

The covariance kernel may also be normalized

$$\int_{[0,1]^d} C_\theta(\mathbf{t}, \mathbf{x}) d\mathbf{t} = 1 \quad \forall \mathbf{x} \in [0,1]^d, \text{ leading to } c_0 = 1 \text{ and } \mathbf{c} = \mathbf{1}. \quad (4)$$

Please note that

$$\Lambda = \text{diag}(\lambda), \quad \lambda = (\lambda_1, \dots, \lambda_n) = V^H C_1.$$



Parameter estimation - Full Bayes

Treat m and s as hyper-parameters with a non-informative, conjugate prior, namely $\rho_{m,s^2}(\xi, \lambda) \propto 1/\lambda$. Then the posterior density for the integral μ given the data is

$$\begin{aligned} \rho_{\mu}(z|\mathbf{f} = \mathbf{y}) &\propto \int_0^{\infty} \int_{-\infty}^{\infty} \rho_{\mu}(z|\mathbf{f} = \mathbf{y}, m = \xi, s^2 = \lambda) \rho_f(\mathbf{y}|\xi, \lambda) \rho_{m,s^2}(\xi, \lambda) d\xi d\lambda \\ &\propto \left(1 + \frac{1}{n-1} \frac{(z - \mu_{\text{full}})^2}{\hat{\sigma}_{\text{full}}^2}\right)^{-n/2} \end{aligned}$$

Where

$$\mu_{\text{full}} = \mu_{\text{EB}}$$

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

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

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



Parameter estimation - Generalized Cross validation

Let $\tilde{y}_i = \mathbb{E}[f(x_i) | f_{-i} = \mathbf{y}_{-i}]$. The cross-validation criterion, which is to be minimized, is sum of squares of the difference between these conditional expectations and the observed values :

$$\text{CV} = \sum_{i=1}^n (y_i - \tilde{y}_i)^2 = \sum_{i=1}^n \left(\frac{\zeta_i}{a_{ii}} \right)^2, \quad \text{where } \zeta = \mathbf{C}_{\theta}^{-1}(\mathbf{y} - m\mathbf{1}),$$

$$a_{ii} \text{ are diagonal elems of } \mathbf{C}_{\theta}^{-1} = \begin{pmatrix} a_{ii} & \mathbf{A}_{-i,i}^T \\ \mathbf{A}_{-i,i} & \mathbf{A}_{-i,-i} \end{pmatrix}$$

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

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

The credible interval width, err_{EB} , is given by

$$\text{err}_{\text{GCV}} = 2.58 s_{\text{GCV}} \sqrt{c_0 - \mathbf{c}^T \mathbf{C}_{\theta}^{-1} \mathbf{c}}$$



Computing the θ , err_{CI} , and $\hat{\mu}$ faster

$$\theta_{\text{EB}} = \underset{\theta}{\text{argmin}} \left[\log \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right) + \frac{1}{n} \sum_{i=1}^n \log(\lambda_i) \right] \quad (5a)$$

$$\theta_{\text{GCV}} = \underset{\theta}{\text{argmin}} \left[\log \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right) - 2 \log \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right) \right] \quad (5b)$$

$$\text{err}_{\text{EB}} = \frac{2.58}{n} \left\{ \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left(1 - \frac{n}{\lambda_1} \right) \right\}^{1/2} \quad (6a)$$

$$\text{err}_{\text{full}} = t_{n-1, 0.995} \left\{ \frac{1}{n(n-1)} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left(\frac{\lambda_1}{n} - 1 \right) \right\}^{1/2} \quad (6b)$$

$$\text{err}_{\text{GCV}} = \frac{2.58}{n} \left\{ \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{\lambda_i} \right]^{-1} \times \left(1 - \frac{n}{\lambda_1} \right) \right\}^{1/2} \quad (6c)$$

Similarly, $\hat{\mu}$ can be computed faster $\hat{\mu}_{\text{EB}} = \hat{\mu}_{\text{full}} = \hat{\mu}_{\text{GCV}} = \frac{1}{n} \sum_{i=1}^n y_i$



Rank-1 Lattice rules : low discrepancy point set

Given the **generating vector** \mathbf{h} , the construction of n - Rank-1 lattice points (Dick and Pillichshammer, 2010) is given by

$$\mathbf{L}_{n,\mathbf{h}} := \{\mathbf{x}_i := (\mathbf{h}\phi(i-1) + \Delta) \bmod 1; \ i = 1, \dots, n\} \quad (7)$$

where Δ is a random shift and \mathbf{h} is a *generalized Mahler integer* (∞ digit expression) (Hickernell and Niederreiter, 2003) also called **generating vector**. $\phi(i)$ is the Van der Corput sequence in base 2. Then the Lattice rule approximation is

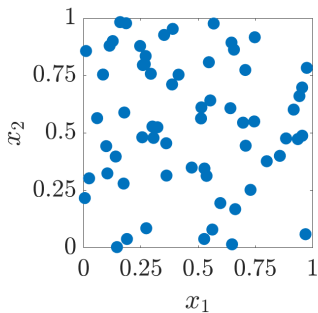
$$\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i).$$

Extensible integration lattices : The number of points in the node set can be increased while retaining the existing points.(Hickernell and Niederreiter, 2003)

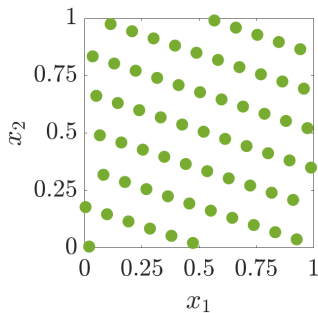


Rank-1 Lattice points in $d = 2$

An example of $n = 64$ IID points and Lattice nodes:



(a) IID



(b) Shifted Lattice

Shift invariant kernel + Lattice points = '*Symmetric circulant kernel*' matrix



Lattice nodes and shift invariant covariance kernels

$$C_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{t}) = \prod_{\ell=1}^d \left[1 - \eta_{\ell} \frac{(2\pi\sqrt{-1})^r}{r!} B_r(|x_{\ell} - t_{\ell}|) \right], \quad r \in 2\mathbb{N}, \quad \eta_{\ell} > 0, \quad \boldsymbol{\theta} = (r, \boldsymbol{\eta})$$

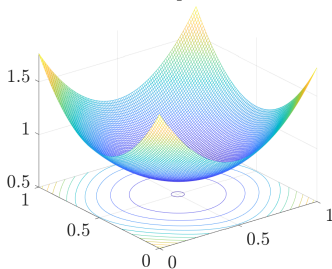
where B_r is Bernoulli polynomial of **order** r (Olver et al., 2013). We call $C_{\boldsymbol{\theta}}$, Fourier kernel. Also this kernel satisfies:

$$c_0 = \int_{[0,1]^2} C_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} = \mathbf{1}, \quad \mathbf{c} = \left(\int_{[0,1]} C_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{t}) d\mathbf{t} \right)_{i=1}^n = \mathbf{1}.$$

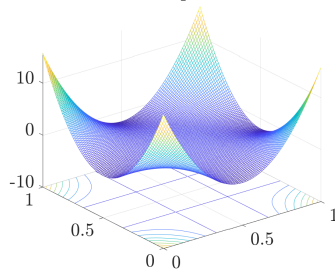


Fourier kernel

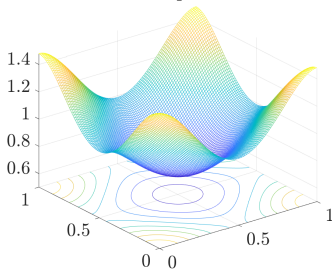
$r=2$ shape=0.10



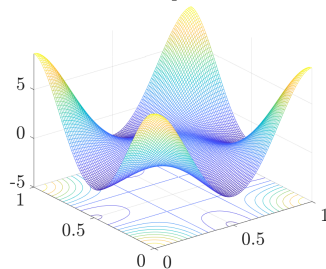
$r=2$ shape=0.90



$r=4$ shape=0.10



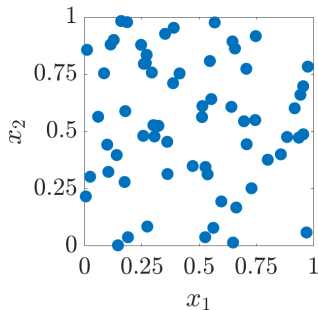
$r=4$ shape=0.90



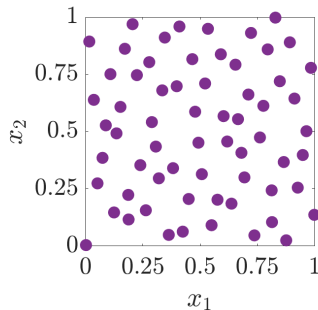


Sobol' Nets

An example of $n = 64$ IID points and Sobol' nets:



(a) IID



(b) Scrambled Sobol'



Walsh Kernels

The Walsh covariance kernels are of the form

$$C_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{t}) = K_{\boldsymbol{\theta}}(\boldsymbol{x} \ominus \boldsymbol{t}). \quad (8)$$

where

$$K_{\boldsymbol{\theta}}(\boldsymbol{x} \ominus \boldsymbol{t}) = \prod_{\ell=1}^d 1 + \eta_{\ell} \omega_r(x_{\ell} \ominus t_{\ell}), \quad \boldsymbol{\eta} = (\eta_1, \dots, \eta_d), \quad \boldsymbol{\theta} = (r, \boldsymbol{\eta}) \quad (9)$$

where r is the kernel order, $\boldsymbol{\eta}$ is the kernel shape parameter. For example, explicit expression is available for ω_r in the case of order $r = 1$ (Nuyens, 201308),

$$\omega_1(x) = 6 \left(\frac{1}{6} - 2^{\lfloor \log_2 x \rfloor - 1} \right). \quad (10)$$



Walsh Kernels

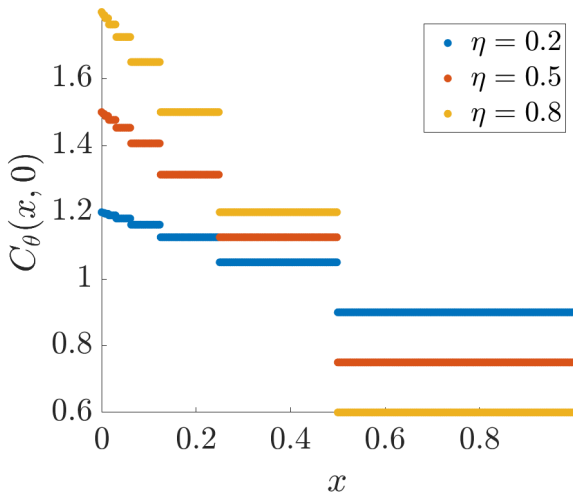


Figure: Walsh kernel of order $r = 1$ in dimension $d = 1$.



Sobol' Nets and Walsh Kernels

Walsh kernels + digital nets = 2×2 *block-Toeplitz* matrix



Fast Bayesian Transform

Theorem

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

$$H^{(m)} C_{\theta}^{(m)} = \Lambda^{(m)} H^{(m)}, \quad m \in \mathbb{N},$$

where (m) denotes the size of the matrix is $2^m \times 2^m$.

By this theorem

$$C_{\theta}^{(m)} = \frac{1}{n} H^{(m)} \Lambda^{(m)} H^{(m)}, \quad \text{where} \quad H^{(m)} = \underbrace{H^{(1)} \otimes \dots \otimes H^{(1)}}_{m \text{ times}}. \quad (11)$$



Cancellation error in err_{CI}

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

$$\text{Let } C_{\theta}(\mathbf{t}, \mathbf{x}) = \prod_{\ell=1}^d \left[1 + \mathring{C}_{\theta, \ell}(t_{\ell}, x_{\ell})\right], \quad \mathring{C}_{\theta, \ell} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}.$$

Direct computation of $\mathring{C}_{\theta}(\mathbf{t}, \mathbf{x}) = C_{\theta}(\mathbf{t}, \mathbf{x}) - 1$ introduces cancellation error if the \mathring{C}_{ℓ} are small. So, we employ the iteration,

$$\begin{aligned} \mathring{C}_{\theta}^{(1)}(\mathbf{t}, \mathbf{x}) &= \mathring{C}_{\theta, 1}(t_1, x_1), \\ \mathring{C}_{\theta}^{(\ell)}(\mathbf{t}, \mathbf{x}) &= \mathring{C}_{\theta}^{(\ell-1)}[1 + \mathring{C}_{\theta, \ell}(t_{\ell}, x_{\ell})] + \mathring{C}_{\theta, \ell}(t_{\ell}, x_{\ell}), \quad \ell = 2, \dots, d, \\ \mathring{C}_{\theta}(\mathbf{t}, \mathbf{x}) &= \mathring{C}_{\theta}^{(d)}(\mathbf{t}, \mathbf{x}). \end{aligned}$$

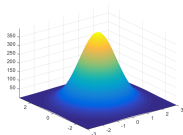
Eigenvalues of \mathring{C}_{θ} : $(\mathring{\lambda}_i)_{i=1}^n = \mathbf{V}^T \mathring{\mathbf{C}}_1$, $\mathring{\lambda}_1 = \lambda_1 - n, \lambda_2, \dots, \lambda_n$

$$\text{err}_{\text{EB}} = \frac{2.58}{n} \sqrt{\frac{\mathring{\lambda}_1}{\lambda_1} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i}}, \quad \theta_{\text{EB}} = \underset{\theta}{\operatorname{argmin}} \left[\log \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \right) + \frac{1}{n} \sum_{i=1}^n \log(\lambda_i) \right]$$



Example Integrands

$$\text{Gaussian probability} = \int_{[a,b]} \frac{e^{-x^T \Sigma^{-1} x / 2}}{(2\pi)^{d/2} |\Sigma|^{1/2}} dx, \text{ (Genz, 1993)}$$



$$\text{Option pricing} = \int_{\mathbb{R}^d} \text{payoff}(x) \underbrace{\frac{e^{-x^T \Sigma^{-1} x / 2}}{(2\pi)^{d/2} |\Sigma|^{1/2}}}_{\text{PDF of Brownian motion at } d \text{ times}} dx, \text{ (Glasserman, 2004)}$$

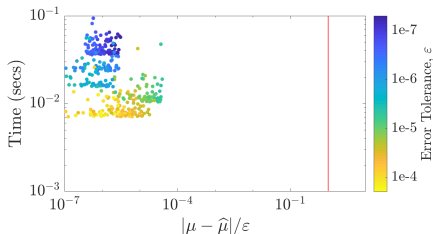
$$\text{where } \text{payoff}(x) = e^{-rT} \max \left(\frac{1}{d} \sum_{k=1}^d S_k(x_k) - K, 0 \right)$$

$$S_j(x_j) = S_0 e^{(r - \sigma^2/2)t_j + \sigma x_j} = \text{stock price at time } t_j = jT/d;$$

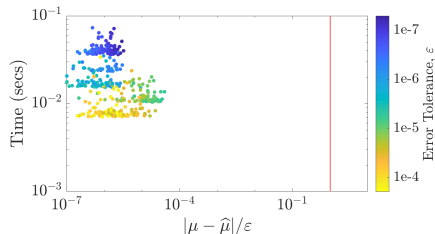
$$\text{Keister integral} = \int_{\mathbb{R}^d} \cos(\|x\|) \exp(-\|x\|^2) dx, \quad d = 1, 2, \dots \text{ (Keister, 1996)}$$

Introduction
○○Bayesian Cubature
○○○○Faster
○○○○Lattice Nodes
○○○○Sobol' Nets
○○○○○Demonstration
○○●○○○○Conclusion
○○○○

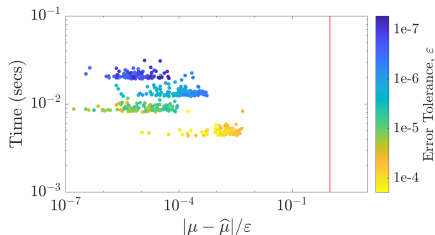
Multivariate normal probability: Lattice



(a) Empirical Bayes



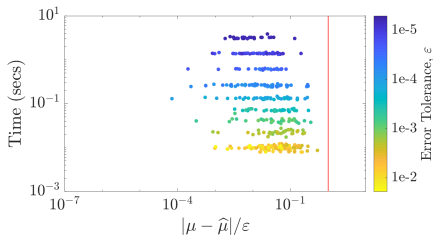
(b) Full Bayes



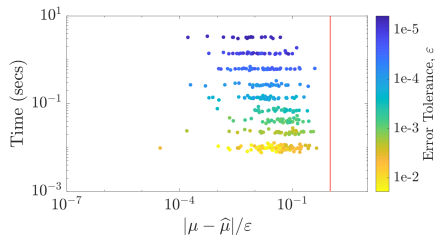
(c) GCV

Introduction
○○Bayesian Cubature
○○○○Faster
○○○○Lattice Nodes
○○○○Sobol' Nets
○○○○○Demonstration
○○○●○○○Conclusion
○○○○

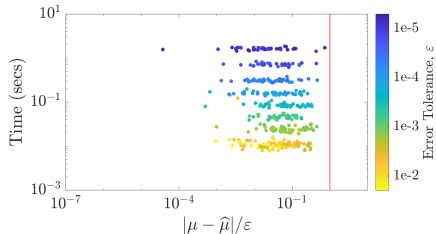
Multivariate normal probability: Sobol'



(a) Empirical Bayes



(b) Full Bayes



(c) GCV



Introduction
○○

Bayesian Cubature
○○○○

Faster
○○○○

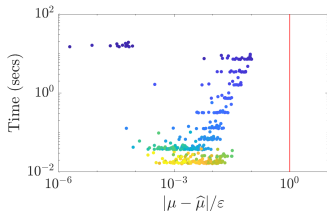
Lattice Nodes
○○○○

Sobol' Nets
○○○○○

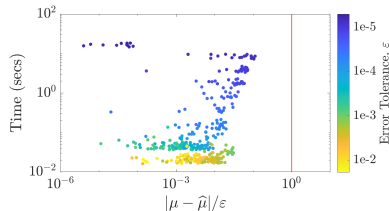
Demonstration
○○○○●○○○

Conclusion
○○○○

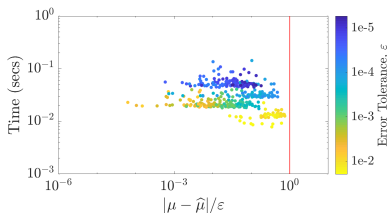
Keister Integral: Lattice



(a) Empirical Bayes



(b) Full Bayes



(c) GCV



Introduction
○○

Bayesian Cubature
○○○○

Faster
○○○○

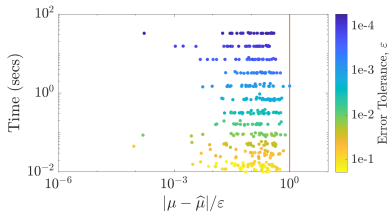
Lattice Nodes
○○○○

Sobol' Nets
○○○○○

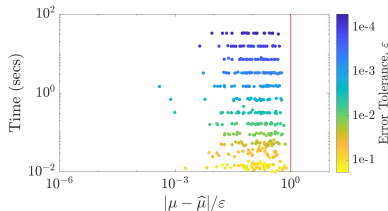
Demonstration
○○○○○●○○

Conclusion
○○○○

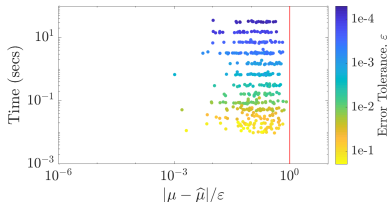
Keister Integral: Sobol'



(a) Empirical Bayes



(b) Full Bayes



(c) GCV



Introduction
○○

Bayesian Cubature
○○○○

Faster
○○○○

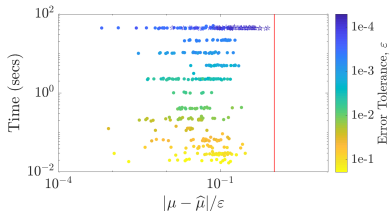
Lattice Nodes
○○○○

Sobol' Nets
○○○○○

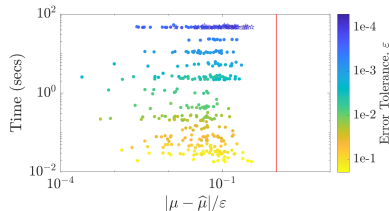
Demonstration
○○○○○●○

Conclusion
○○○○

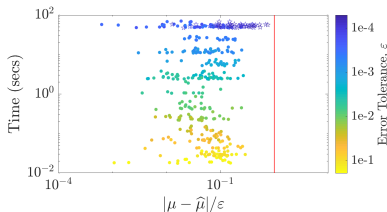
Option pricing: Lattice



(a) Empirical Bayes



(b) Full Bayes



(c) GCV



Introduction

○○

Bayesian Cubature

○○○○

Faster

○○○○

Lattice Nodes

○○○○

Sobol' Nets

○○○○○

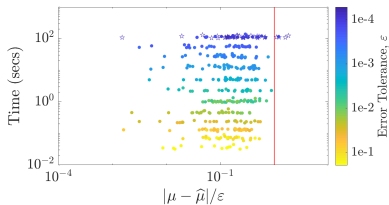
Demonstration

○○○○○○●

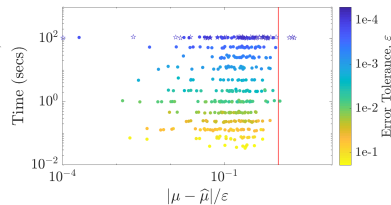
Conclusion

○○○○

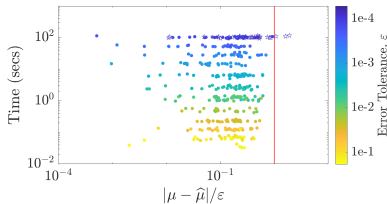
Option pricing: Sobol'



(a) Empirical Bayes



(b) Full Bayes



(c) GCV



Summary

- Developed a *technique* for a **Fast Bayesian transform**
- Developed two **fast automatic Bayesian cubature** algorithms with $\mathcal{O}(n \log n)$ complexity
- Having the advantages of a kernel method and the low computation cost of Quasi Monte carlo
- Scalable based on the complexity of the Integrand i.e, Kernel order and Lattice-points can be chosen to suit the smoothness of the integrand
- Conditioning problem if the kernel C is very smooth
- Source code : https://github.com/GailGithub/GAIL_Dev/tree/feature/BayesianCubature.
- A part of this work was published as a paper "Fast Automatic Bayesian Cubature Using Lattice Sampling by R. Jagadeeswaran, Fred J. Hickernell" (<https://arxiv.org/abs/1809.09803>)



Future work

- Choosing the **kernel order r** and **periodization** transform automatically
- Diagnostics for Gaussian process assumption
- Broaden the choice of numerical examples
- Better handling of **conditioning** problem and numerical errors
- Higher order nets and Walsh kernels could be used to achieve higher order or accuracy.



Future work : More applications

- **Control variates** : We would like to approximate a function of the form $(f - \beta_1 g_1 - \dots - \beta_p g_p)$, then

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

- **Function approximation** : consider approximating a function of the form

$$\int_{[0,1]^d} \underbrace{f(\boldsymbol{\Phi}(t))}_{g(t)} \cdot \left| \frac{\partial \boldsymbol{\Phi}}{\partial t} \right| dt, \quad \text{where } \left| \frac{\partial \boldsymbol{\Phi}}{\partial t} \right| \text{ is Jacobian, then}$$

$$g(\boldsymbol{\Psi}(x)) = f(\underbrace{\boldsymbol{\Phi}(\boldsymbol{\Psi}(x))}_x) \cdot \left| \frac{\partial \boldsymbol{\Phi}}{\partial t} \right|(\boldsymbol{\Psi}(x)), \quad f(x) = g(\boldsymbol{\Psi}(x)) \cdot \frac{1}{\left| \frac{\partial \boldsymbol{\Phi}}{\partial t} \right|(\boldsymbol{\Psi}(x))}$$

Finally, the function approximation is

$$\tilde{f}(x) = \tilde{g}(\boldsymbol{\Psi}(x)) = \sum w_i C(.,.)$$

Thank you!



Parameter estimation - Full Bayes - General prior

What if $\rho_{m,s^2}(\xi, \lambda) \propto g(1/\lambda)$?

$$\begin{aligned} \rho_{\mu}(z|\mathbf{f} = \mathbf{y}) &\propto \mathcal{L}\mathcal{T}\{g(1/\cdot)\}^{(\frac{n-4}{2})}(\chi) \\ &\propto \mathcal{L}\mathcal{T}\{g(1/\cdot)\}^{(\frac{n-4}{2})}\left(1 + \frac{(z - \hat{\mu}_{\text{full}})^2}{(n-1)\hat{\sigma}_{\text{full}}^2}\right) \end{aligned}$$

Thus, $\rho_{\mu}(z|\mathbf{f} = \mathbf{y})$ is proportional to $\left(\frac{n-4}{2}\right)$ th derivative of the Laplace transform of $g(1/\cdot)$ evaluated at χ , where $\chi \propto 1 + \frac{(z - \hat{\mu}_{\text{full}})^2}{(n-1)\hat{\sigma}_{\text{full}}^2}$.

Our motivation to experiment with the general prior was to show that it may be possible to infer the prior from the integrand samples.



Periodization Transforms

Suppose the original integral is

$$\mu := \int_{(a,b)^d} g(\mathbf{t}) \, d\mathbf{t}, \quad \text{where } g \text{ is smooth, not periodic.}$$

The Baker's transform, the tent transform,

$$\Psi : x \mapsto (\Psi(x_1), \dots, \Psi(x_d)), \quad \Psi(x) = 1 - 2|x - 1/2|, \quad f(x) = g(\Psi(x)).$$

A family of smoother variable transforms:

$$\Psi : x \mapsto (\Psi(x_1), \dots, \Psi(x_d)), \quad \Psi : [0, 1] \mapsto [0, 1], \quad f(x) = g(\Psi(x)) \prod_{\ell=1}^d \Psi'(x_\ell).$$

Example:

$$C^1 : \Psi(x) = x^3(10 - 15x + 6x^2), \quad \Psi'(x) = 30x^2(1 - x)^2,$$

$$\text{Sidi's } C^1 : \Psi(x) = x - \frac{\sin(2\pi x)}{2\pi}, \quad \Psi'(x) = 1 - \cos(2\pi x),$$

when it holds $\Psi \in C^{r+1}[0, 1]$, $\lim_{x \downarrow 0} x^{-r-1} \Psi'(x) = \lim_{x \uparrow 1} (1-x)^{-r-1} \Psi'(x) = 0$, and

$g \in C^{(r, \dots, r)}[0, 1]^d$, for $r \in \mathbb{N}_0$.



Walsh Transform

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

$$H^{(1)} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix},$$

$$H^{(2)} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix},$$

$$\vdots$$

$$H^{(m)} = \begin{pmatrix} H^{(m-1)} & H^{(m-1)} \\ H^{(m-1)} & -H^{(m-1)} \end{pmatrix} = \underbrace{H^{(1)} \otimes \dots \otimes H^{(1)}}_{m \text{ times}} = H^{(1)} \otimes H^{(m-1)} \quad (12)$$

where \otimes is Kronecker product.



Sobol' Nets and Walsh Kernels

Theorem

Any symmetric, positive definite, digital shift-invariant covariance kernel of the form (9) scaled to satisfy (4), when matched with digital net data-sites, satisfies assumptions (3). The fast Walsh-Hadamard transform (FWHT) can be used to expedite the estimates of θ in (5) and the credible interval widths (6) in $\mathcal{O}(n \log n)$ operations. The cubature, $\hat{\mu}$, is just the sample mean.

Walsh kernels + digital nets = 2×2 block-Toeplitz matrix



Eigenvectors of C_θ

The columns of Walsh-Hadamard matrix are the eigenvectors of C_θ , i.e., $V := H$

Theorem

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

$$C_\theta = (C(x_i, x_j))_{i,j=0}^{n-1} = (K(x_i \ominus x_j))_{i,j=0}^{n-1},$$

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



Iterative Computation of Walsh Transform

Let $\tilde{\mathbf{y}} = \mathbf{H}^{(m+1)} \mathbf{y}$ for some arbitrary $\mathbf{y} \in \mathbb{R}^{2^n}$, $n = 2^m$. Define,

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

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

Then,

$$\begin{aligned} \tilde{\mathbf{y}} &= \mathbf{H}^{(m+1)} \mathbf{y} = \begin{pmatrix} \mathbf{H}^{(m)} & \mathbf{H}^{(m)} \\ \mathbf{H}^{(m)} & -\mathbf{H}^{(m)} \end{pmatrix} \begin{pmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \end{pmatrix}, \quad \text{by (12)} \\ &= (\mathbf{H}^{(m)} \mathbf{y}^{(1)} + \mathbf{H}^{(m)} \mathbf{y}^{(2)}, \mathbf{H}^{(m)} \mathbf{y}^{(1)} - \mathbf{H}^{(m)} \mathbf{y}^{(2)}) = \begin{pmatrix} \tilde{\mathbf{y}}^{(1)} + \tilde{\mathbf{y}}^{(2)} \\ \tilde{\mathbf{y}}^{(1)} - \tilde{\mathbf{y}}^{(2)} \end{pmatrix} =: \tilde{\mathbf{y}} \end{aligned}$$



Cone of functions and the Credible interval

In this research we assume that the integrand belongs to a cone of well-behaved functions, \mathcal{C} . Suppose that

$$|\mu(f) - \hat{\mu}_n(f)| \leq \text{err}_{\text{CI}}(f(x_1), \dots, f(x_n)) \quad (13)$$

for some f , which it is 99% of the time under our hypothesis. Also note that our $\text{err}_{\text{CI}}, \text{CI} \in \{\text{EB}, \text{GCV}\}$ are positively homogeneous functions, meaning,

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

Thus if f satisfies (13), then

$$\begin{aligned} |\mu(af) - \hat{\mu}_n(af)| &= |a| |\mu(f) - \hat{\mu}_n(f)| \\ &\leq |a| \text{err}_{\text{CI}}(f(x_1), \dots, f(x_n)) = \text{err}_{\text{CI}}(af(x_1), \dots, af(x_n)) \end{aligned}$$

for all real a . Thus the set of all f satisfying (13) is a *cone*, \mathcal{C} . Cones of functions satisfy the property that if $f \in \mathcal{C}$ then $af \in \mathcal{C}$.

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



Let $f_{\text{TRUE}}(\mathbf{x}) = \exp(\sum_{\ell=1}^d \cos(2\pi x_{\ell}))$ and the peaky integrand

$f_{\text{PEAKY}}(\mathbf{x}) = f_{\text{TRUE}} + a_{\text{PEAKY}} f_{\text{NOISE}}$, $f_{\text{NOISE}}(\mathbf{x}) = (1 - \exp(2\pi\sqrt{-1}\mathbf{x}^T \boldsymbol{\zeta}))$, $\boldsymbol{\zeta} \in \mathbb{R}^d$ is some d -dimensional vector belonging to the dual space of the lattice nodes. The f_{NICE} is obtained by kernel interpolation.

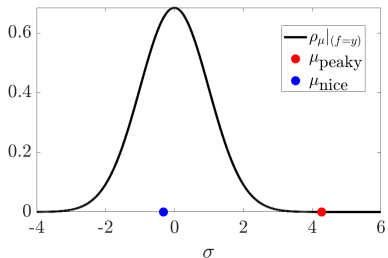
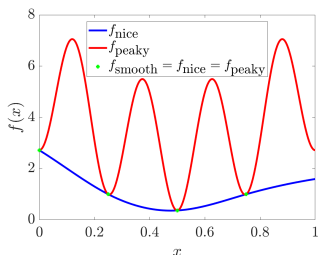


Figure: Left: Example integrands 1) f_{NICE} , a smooth function, 2) f_{PEAKY} , a peaky function. The function values $f_{\text{PEAKY}}(\mathbf{x}_i) = f_{\text{NICE}}(\mathbf{x}_i) = f_{\text{TRUE}}(\mathbf{x}_i)$ for $i = 1, \dots, n$. Right: Probability distributions showing the relative integral position of a smooth and a peaky function. f_{NICE} lies within the center 99% of the confidence interval, and f_{PEAKY} lies on the outside of 99% of the confidence interval.



Shape parameter search using gradient descent

Steepest descent search is defined as:

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

where ν is the step size for the gradient descent, j is the iteration index, and $\frac{\partial}{\partial \eta_{\ell}} \mathcal{L}(\boldsymbol{\theta}|\mathbf{y})$ is either (14) or (15) depending on the choice of the hyperparameter search method. The parameter η_{ℓ} is usually searched in the whole \mathbb{R} by using the simple domain transformation.



Computing the derivative of $\mathcal{L}_{\text{EB}}(\boldsymbol{\theta}|\mathbf{y})$

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

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

where we used some of the results from (Dong et al., 2017). After using the fast Bayesian transform properties

$$\frac{\partial}{\partial \theta_\ell} \mathcal{L}_{\text{EB}}(\boldsymbol{\theta}|\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i} - \left(\sum_{i=2}^n \frac{|\tilde{\mathbf{y}}_i|^2 \bar{\lambda}_{i(\ell)}}{\lambda_i^2} \right) \left(\sum_{i=2}^n \frac{|\tilde{\mathbf{y}}_i|^2}{\lambda_\ell} \right)^{-1} \quad (14)$$



Computing the derivative of $\mathcal{L}_{\text{GCV}}(\boldsymbol{\theta}|\mathbf{y})$

Similarly for the generalized cross-validation

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

$$\text{where } m_{\text{GCV}} = \frac{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{y}}{\mathbf{1}^T \mathbf{C}_{\boldsymbol{\theta}}^{-2} \mathbf{1}},$$

After using the fast Bayesian transform properties

$$\begin{aligned} \frac{\partial}{\partial \theta_{\ell}} \mathcal{L}_{\text{GCV}}(\boldsymbol{\theta}|\mathbf{y}) = & -2 \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i^2} \right)^{-1} \left(\sum_{i=2}^n \frac{|\tilde{y}_i|^2 \bar{\lambda}_{i(\ell)}}{\lambda_i^3} \right) \\ & + 2 \left(\sum_{i=1}^n \frac{1}{\lambda_i} \right)^{-1} \left(\sum_{i=1}^n \frac{\bar{\lambda}_{i(\ell)}}{\lambda_i^2} \right), \end{aligned} \quad (15)$$

where $\bar{\lambda}_{i(\ell)}$ is the derivative of the i th eigenvalue of the Gram matrix, \mathbf{C} , in the ℓ th variable.



Product Kernels

Product kernels in d dimensions are of the form,

$$C_{\theta}(\mathbf{t}, \mathbf{x}) = \prod_{\ell=1}^d \left[1 - \eta_{\ell} \mathfrak{C}(x_{\ell}, t_{\ell}) \right] \quad (16)$$

where η_{ℓ} is called shape parameter.

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

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

When η_{ℓ} is different for each $\ell = 1, \dots, d$

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



To compute $\bar{\lambda}_{i(\ell)}$

If V does not depend on θ then one can fast compute the derivative of Gram matrix C ,

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

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

$$\bar{\lambda}_{(\ell)} = \frac{\partial \lambda}{\partial \theta_\ell} = \left(\frac{\partial \lambda_i}{\partial \theta_\ell} \right)_{i=1}^n = \left(\frac{\partial}{\partial \theta_\ell} V^H C_1 \right) = V^H \left(\frac{\partial}{\partial \theta_\ell} C_\theta(x_1, x_i) \right)_{i=1}^n, \quad (17)$$

where we used the fast Bayesian transform property $\lambda = V^H C_1$.



References I

Briol, F-X, C. J. Oates, M. Griolami, M. A. Osborne, and D. Sejdinovic. 2018+. *Probabilistic integration: A role in statistical computation?*, Statist. Sci. to appear.

Diaconis, P. 1988. *Bayesian numerical analysis*, Statistical decision theory and related topics iv, papers from the 4th purdue symp., west lafayette/indiana 1986, pp. 163–175.

Dick, J. and F. Pillichshammer. 2010. *Digital nets and sequences: Discrepancy theory and quasi-Monte Carlo integration*, Cambridge University Press, Cambridge.

Dong, K., D. Eriksson, H. Nickisch, D. Bindel, and A. G. Wilson. 2017. *Scalable log determinants for gaussian process kernel learning*, NIPS. in press.

Genz, A. 1993. *Comparison of methods for the computation of multivariate normal probabilities*, Computing Science and Statistics **25**, 400–405.

Glasserman, P. 2004. *Monte Carlo methods in financial engineering*, Applications of Mathematics, vol. 53, Springer-Verlag, New York.

Hickernell, F. J. and H. Niederreiter. 2003. *The existence of good extensible rank-1 lattices*, J. Complexity **19**, 286–300.

Keister, B. D. 1996. *Multidimensional quadrature algorithms*, Computers in Physics **10**, 119–122.

Nuyens, Dirk. 201308. *The construction of good lattice rules and polynomial lattice rules*.



References II

O'Hagan, A. 1991. *Bayes-Hermite quadrature*, J. Statist. Plann. Inference **29**, 245–260.

Olver, F. W. J., D. W. Lozier, R. F. Boisvert, C. W. Clark, and A. B. O. Dalhuis. 2013. *Digital library of mathematical functions*.

Rasmussen, C. E. and C. Williams. 2003. *Bayesian Monte Carlo*, Advances in Neural Information Processing Systems, pp. 489–496.

Ritter, K. 2000. *Average-case analysis of numerical problems*, Lecture Notes in Mathematics, vol. 1733, Springer-Verlag, Berlin.

Sobol', I. M. 1976. *Uniformly distributed sequences with an additional uniformity property*, Zh. Vychisl. Mat. i Mat. Fiz. **16**, 1332–1337 (Russian).

Traub, J. F., G. W. Wasilkowski, and H. Woźniakowski. 1988. *Information-based complexity*, Academic Press, Boston.