

Fast automatic Bayesian Cubature with Lattice sampling

Jagadeeswaran R.

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

• Feb 27, 2019 •



Numerical Integration

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

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

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

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



Numerical Integration

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

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

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

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

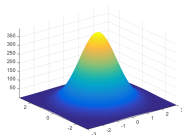
The goal of this work is to

- Develop an automatic algorithm for integration
- Assume f is drawn from a Gaussian process
 - Need to estimate the mean and Covariance kernel
- MLE is expensive in general
 - Use points and kernel for which MLE is cheap
- Use an **extensible** point-set and an algorithm that allows to add more points if needed
- Determine n such that, given ϵ , $|\mu - \hat{\mu}_n| \leq \epsilon$



Motivating Examples

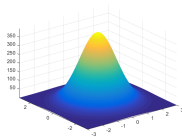
$$\text{Gaussian probability} = \int_{[a,b]} \frac{e^{-x^T \Sigma^{-1} x / 2}}{(2\pi)^{d/2} |\Sigma|^{1/2}} dx, \text{ (Genz, 1993)}$$





Motivating Examples

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



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

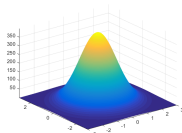
$$\text{where } \text{payoff}(x) = e^{-rT} \max \left(\frac{1}{d} \sum_{k=1}^d S_k(x_k) - K, 0 \right)$$

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



Motivating Examples

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



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

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

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

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



1: **procedure** AUTOCUBATURE(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** $\text{err}_n > \epsilon$ 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 θ

6: Compute error bound err_n

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

8: **end while**

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

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

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

12: **end procedure**

Problem:

■ How to choose $\{x_i\}_{i=1}^n$, and $\{w_i\}_{i=1}^n$ to make $|\mu - \hat{\mu}_n|$ small? what is err_n ?
(Bayesian posterior error)

■ How to find n such that $|\mu - \hat{\mu}_n| \leq \text{err}_n \leq \epsilon$? (automatic cubature)

1: **procedure** AUTOCUBATURE(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** $\text{err}_n > \epsilon$ 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 θ

6: Compute error bound err_n

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

8: **end while**

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

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

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

12: **end procedure**

Problem:

- How to choose $\{x_i\}_{i=1}^n$, and $\{w_i\}_{i=1}^n$ to make $|\mu - \hat{\mu}_n|$ small? what is err_n ?
(Bayesian posterior error)

- How to find n such that $|\mu - \hat{\mu}_n| \leq \text{err}_n \leq \epsilon$? (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, $s^2 C_\theta$, $C_\theta : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$.

Lets define $c_0 = \int_{[0,1] \times [0,1]} C_\theta(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$,

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



Bayesian posterior error

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

Lets define $c_0 = \int_{[0,1] \times [0,1]} C_\theta(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$,

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

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

where $\mathbf{y} = (f(\mathbf{x}_i))_{i=1}^n$. Moreover m, s and θ needs to be inferred.

$$\hat{\mu}_n = w_0 + \sum_{i=1}^n w_i f(\mathbf{x}_i) = w_0 + \mathbf{w}^T \mathbf{y}$$



Bayesian posterior error

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

Lets define $c_0 = \int_{[0,1] \times [0,1]} C_\theta(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$,

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

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

where $\mathbf{y} = (f(\mathbf{x}_i))_{i=1}^n$. Moreover m, s and θ needs to be inferred.

$$\hat{\mu}_n = w_0 + \sum_{i=1}^n w_i f(\mathbf{x}_i) = w_0 + \mathbf{w}^T \mathbf{y}$$

In general choosing $w_0 = m(1 - \mathbf{1}^T \mathbf{C}^{-1} \mathbf{c})$, $\mathbf{w} = \mathbf{C}^{-1} \mathbf{c}$, makes error unbiased

If $m = 0$ fixed, choosing $\mathbf{w} = \mathbf{C}^{-1} \mathbf{c}$, makes error unbiased



Parameter estimation - Maximum likelihood

The log-likelihood of the parameters given the data $\mathbf{y} = (f(x_i))_{i=1}^n$ is :

$$l(s, \boldsymbol{\theta} | \mathbf{y}) = \log \left[\frac{1}{\sqrt{(2\pi)^n \det(s^2 \mathbf{C})}} \exp \left(-\frac{1}{2} s^{-2} (\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m\mathbf{1}) \right) \right]$$



Parameter estimation - Maximum likelihood

The log-likelihood of the parameters given the data $\mathbf{y} = (f(x_i))_{i=1}^n$ is :

$$l(s, \theta | \mathbf{y}) = \log \left[\frac{1}{\sqrt{(2\pi)^n \det(s^2 \mathbf{C})}} \exp \left(-\frac{1}{2} s^{-2} (\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m\mathbf{1}) \right) \right]$$

Maximising w.r.t m and then s^2 , further with θ :

$$m_{\text{MLE}} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}, \quad s_{\text{MLE}}^2 = \frac{1}{n} (\mathbf{y} - m_{\text{MLE}} \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{MLE}} \mathbf{1}), \quad (\text{Explicit})$$

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

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



Parameter estimation - Maximum likelihood

The log-likelihood of the parameters given the data $\mathbf{y} = (f(x_i))_{i=1}^n$ is :

$$l(s, \theta | \mathbf{y}) = \log \left[\frac{1}{\sqrt{(2\pi)^n \det(s^2 \mathbf{C})}} \exp \left(-\frac{1}{2} s^{-2} (\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m\mathbf{1}) \right) \right]$$

Maximising w.r.t m and then s^2 , further with θ :

$$m_{\text{MLE}} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}, \quad s_{\text{MLE}}^2 = \frac{1}{n} (\mathbf{y} - m_{\text{MLE}} \mathbf{1})^T \mathbf{C}^{-1} (\mathbf{y} - m_{\text{MLE}} \mathbf{1}), \quad (\text{Explicit})$$

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

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

Why do we need θ_{MLE} ? Function space spanned by C_θ customized to contain the integrand function f .



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 :

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



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 :

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

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

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

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



Parameter estimation - Leave-one-out Cross validation

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

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

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

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



Parameter estimation - Leave-one-out 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 :

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

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

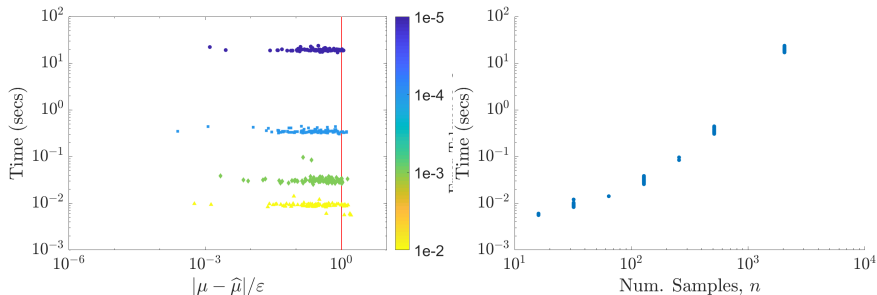
$$\text{GCV} = \frac{\sum_{i=1}^n \zeta_i^2}{\left(\frac{1}{n} \sum_{i=1}^n a_{ii}\right)^2} = \frac{(\mathbf{y} - m\mathbf{1})^T \mathbf{C}^{-2} (\mathbf{y} - m\mathbf{1})}{\left(\frac{1}{n} \text{trace}(\mathbf{C}^{-1})\right)^2}.$$

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

$$s_{\text{GCV}}^2 := \mathbf{y}^T \left[\mathbf{C}^{-2} - \frac{\mathbf{C}^{-2} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-2}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}} \right] \mathbf{y} [\text{trace}(\mathbf{C}^{-1})]^{-1}, \quad m_{\text{GCV}} := \frac{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-2} \mathbf{1}}.$$



Multivariate normal integration with Matern kernel



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



Fast transform kernel

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

$$C = (C_\theta(x_i, x_j))_{i,j=1}^n = (C_1, \dots, C_n) = \frac{1}{n} V \Lambda V^H$$

$$V := (v_1, \dots, v_n)^T = (V_1, \dots, V_n), \quad V_1 = v_1 = \mathbf{1},$$

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

Then



Fast transform kernel

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

$$C = (C_\theta(x_i, x_j))_{i,j=1}^n = (C_1, \dots, C_n) = \frac{1}{n} V \Lambda V^H$$

$$V := (v_1, \dots, v_n)^T = (V_1, \dots, V_n), \quad V_1 = v_1 = \mathbf{1},$$

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

Then

$$V^H C_1 = V^H \left(\frac{1}{n} V \Lambda v_1^* \right) = \Lambda \mathbf{1} = (\lambda_1, \dots, \lambda_n)^T = \lambda$$

C_θ is a fast transform kernel, if the transform $\hat{z} = V^H z$ for arbitrary z , can be done in $\mathcal{O}(n \log n)$. Using the fast transform,

$$a^T C^p b = \frac{1}{n} a^T V \Lambda^p V^H b = \frac{1}{n} \tilde{a}^H \Lambda^p \tilde{b} = \frac{1}{n} \sum_{i=1}^n \lambda_i^p \tilde{a}_i^* \tilde{b}_i,$$

The covariance kernel used in practice also may be normalized

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



Faster parameters estimation

MLE and GCV estimates of θ made faster by using the properties of the fast transform kernel:

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{argmin}} \left[\log \left(\sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \right) + \frac{1}{n} \sum_{i=1}^n \log(\lambda_i) \right],$$

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

Also,

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

$$\hat{\sigma}_{\text{full}}^2 = \frac{1}{n(n-1)} \sum_{i=2}^n \frac{|\tilde{y}_i|^2}{\lambda_i} \left(\frac{\lambda_1}{n} - 1 \right), \quad \text{where}$$

$$\hat{\mathbf{y}} = (\hat{y}_i)_{i=1}^n = \mathbf{V}^T \mathbf{y}, \quad \boldsymbol{\lambda} = (\lambda_i)_{i=1}^n = \mathbf{V}^T \mathbf{C}_1, \quad \text{where } \mathbf{C}_1 = (C(\mathbf{x}_i, \mathbf{x}_1))_{i=1}^n$$

$\mathcal{O}(n \log n)$ operations to compute $\hat{\mathbf{y}}$ and $\hat{\boldsymbol{\lambda}}$, So the θ_{MLE}



Computing the error bound err and $\hat{\mu}$ faster

Using the properties of the **fast transform kernel**, the error bound err_n can be computed faster

$$\text{err}_{\text{MLE}} = \frac{2.58}{n} \left\{ \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i} \left(1 - \frac{n}{\lambda_1} \right) \right\}^{1/2}$$

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

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

similarly, $\hat{\mu}$ can be computed faster

$$\hat{\mu}_{\text{MLE}} = \hat{\mu}_{\text{full}} = \hat{\mu}_{\text{GCV}} = \mathbf{w}^T \mathbf{y} = \sum_{i=1}^n \frac{y_i}{n}$$

where

$$\hat{\mathbf{y}} = \mathbf{V}^T \mathbf{y}, \quad \boldsymbol{\lambda} = \mathbf{V}^T \mathbf{C}_1, \quad \text{where } \mathbf{C}_1 = (\mathbf{C}(\mathbf{x}_i, \mathbf{x}_1))_{i=1}^n$$

$\mathcal{O}(n \log n)$ operations to compute the err . $\mathcal{O}(n)$ operations to compute the $\hat{\mu}$



Special shift invariant covariance kernel

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



Special shift invariant covariance kernel

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

where B_r is Bernoulli polynomial of **order** r (Olver et al., 2013). We call C_{θ} , Fourier kernel. Also this kernel has:

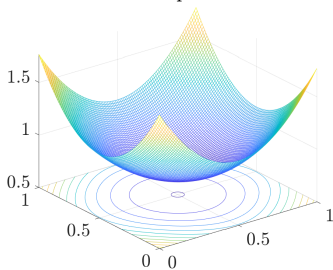
$$c_0 = \int_{[0,1]^2} C_{\theta}(x, t) dx dt = \mathbf{1}, \quad \mathbf{c} = \left(\int_{[0,1]} C_{\theta}(x_i, t) dt \right)_{i=1}^n = \mathbf{1}.$$

$$\mathbf{v} = \left(e^{2\pi n \sqrt{-1} \phi(i-1) \phi(j-1)} \right)_{i=1}^n$$

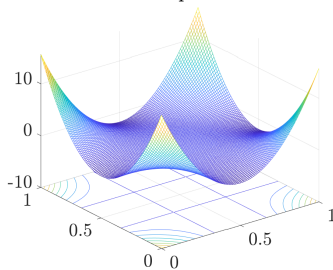


Fourier kernel

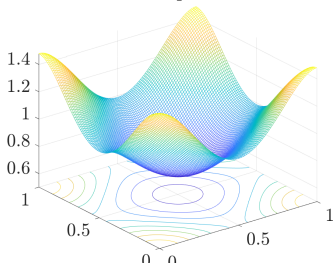
$r=2$ shape=0.10



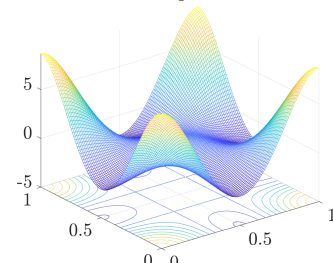
$r=2$ shape=0.90



$r=4$ shape=0.10



$r=4$ shape=0.90





Rank-1 Lattice rules : low discrepancy point set

Given the “generating vector” h , the construction of n - Rank-1 lattice points (Dick and Pillichshammer, 2010) is given by

$$L_{n,h} := \{x_i := h\phi(i-1) \bmod 1; \ i = 1, \dots, n\} \quad (1)$$

where 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

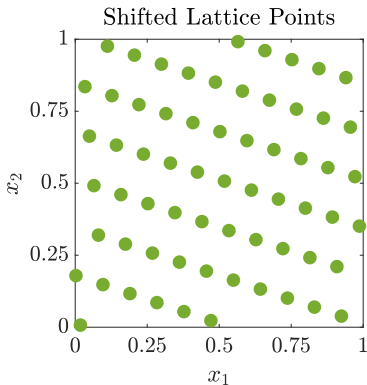
$$\frac{1}{n} \sum_{k=1}^n f \left(\left\{ \frac{kh}{n} + \Delta \right\}_1 \right)$$

where $\{.\}_1$ the fractional part, i.e, *modulo 1* operator and Δ a random shift.

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 $d = 2$



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



The shift invariant kernel with rank-1 Lattice points

- Satisfies all the requirements to be a **fast transform kernel**
- Fast transform = fast Fourier transform
- Complexity of fast Fourier transform is $\mathcal{O}(n \log n)$
- No need to form the kernel matrix C explicitly, so $\mathcal{O}(n^2)$ memory not required
- There are **no** matrix inversions, **no** matrix multiplications
- Factorization of matrix C does not need any computations.

where V is just the Fourier coefficient matrix: $V = \left(e^{2\pi n \sqrt{-1}(i-1)(j-1)} \right)_{i=1}^n$



Iterative DFT

We can avoid recomputing the whole Fourier transform for function values $y = (y_i = f(x_i))_{i=1}^n$ in every iteration. Discrete Fourier transform is defined as

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



Iterative DFT

We can avoid recomputing the whole Fourier transform for function values $y = (y_i = f(x_i))_{i=1}^n$ in every iteration. Discrete Fourier transform is defined as

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

Rearrange sum into even indexed $j = 2l$ and odd indexed $j = 2l + 1$.

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

we use this concept along with *extensible point set*, to avoid recomputing the DFT of y in every iteration.



Cancellation error in err

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



Cancellation error in err

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

Let $\tilde{C}(\mathbf{x}, t) = C(\mathbf{x}, t) - 1$, then $\tilde{C} = C - \mathbf{1}\mathbf{1}^T$, and $\tilde{C} = V\tilde{\Lambda}V^H$

where

$$\tilde{\Lambda} = \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_n), \text{ to compute } (\tilde{\lambda}_i)_{i=1}^n = V^T \tilde{C}_1$$

$$\tilde{\lambda}_1 = \lambda_1 - n, \quad \tilde{\lambda}_j = \lambda_j, \quad \forall j = 2, \dots, n$$



Cancellation error in err

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

$$\text{Let } \tilde{C}(x, t) = C(x, t) - 1, \quad \text{then } \tilde{C} = C - \mathbf{1}\mathbf{1}^T, \quad \text{and } \tilde{C} = V\tilde{\Lambda}V^H$$

where

$$\tilde{\Lambda} = \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_n), \quad \text{to compute } (\tilde{\lambda}_i)_{i=1}^n = V^T \tilde{C}_1$$

$$\tilde{\lambda}_1 = \lambda_1 - n, \quad \tilde{\lambda}_j = \lambda_j, \quad \forall j = 2, \dots, n$$

vector $\tilde{C}_1 = \tilde{C}_1^{(d)}$ computed iteratively

$$\tilde{C}_1^{(1)} = \theta (B(x_{i1} - x_{11}))_{i=1}^n, \quad C_1^{(1)} = \mathbf{1} + \tilde{C}_1^{(1)},$$

$$\forall 1 < k \leq d, \quad \tilde{C}_1^{(k)} = \theta C_1^{(k-1)} \circ (B(x_{ik} - x_{1k}))_{i=1}^n + \tilde{C}_{k-1}, \quad C_1^{(k)} = \mathbf{1} + \tilde{C}_1^{(k)}$$

where \circ is elementwise multiplication. MATLAB=.*

Using this to avoid cancellation error

$$1 - \frac{n}{\lambda_1} = 1 - \frac{n}{n + \tilde{\lambda}_1} = \frac{\tilde{\lambda}_1}{n + \tilde{\lambda}_1}$$



Periodization transforms

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

$$\text{C0 : } \tilde{f}(t) = f(\tilde{g}_0(t)) \prod_{j=1}^d g'_0(t_j), \quad g_0(t) = 3t^2 - 2t^3, \quad g'_0(t) = 6t(1 - t)$$

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

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

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

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

$$\text{Sidi's C2 : } \tilde{f}(t) = f(\tilde{\psi}_3(t)) \prod_{j=1}^d \psi'_3(t_j), \quad \psi_3(t) = \frac{1}{16} (8 - 9 \cos(\pi t) + \cos(3\pi t)),$$

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



Multivariate normal probability

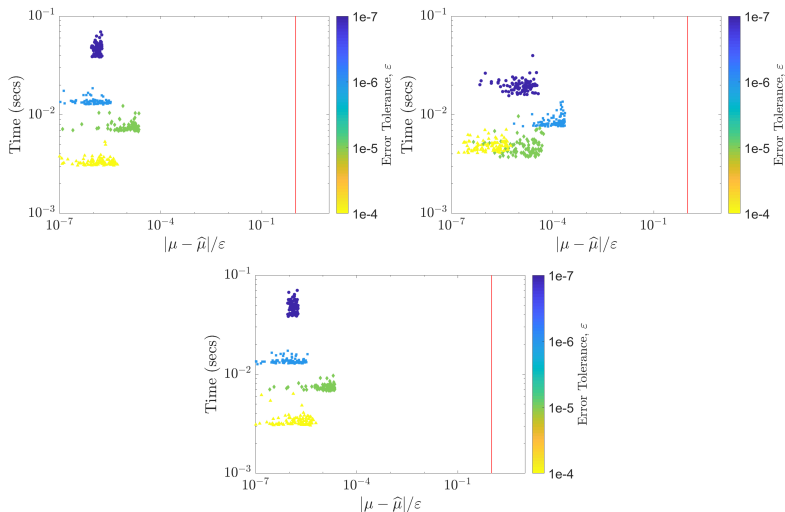


Figure: Multivariate normal probability example using 1) Empirical Bayes, 2) GCV, 3) Full Bayes stopping criterion



Keister Integral with arb mean m

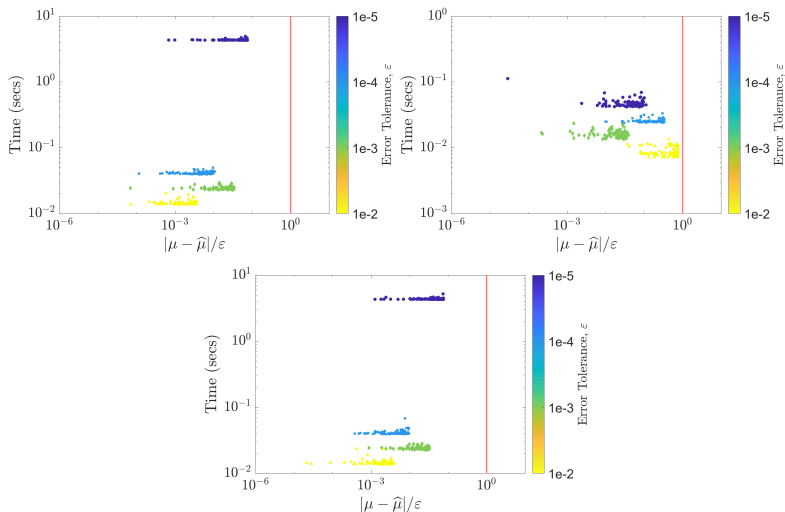


Figure: Integrating Keister function using 1) Empirical Bayes, 2) GCV, 3) Full Bayes stopping criterion



Option pricing

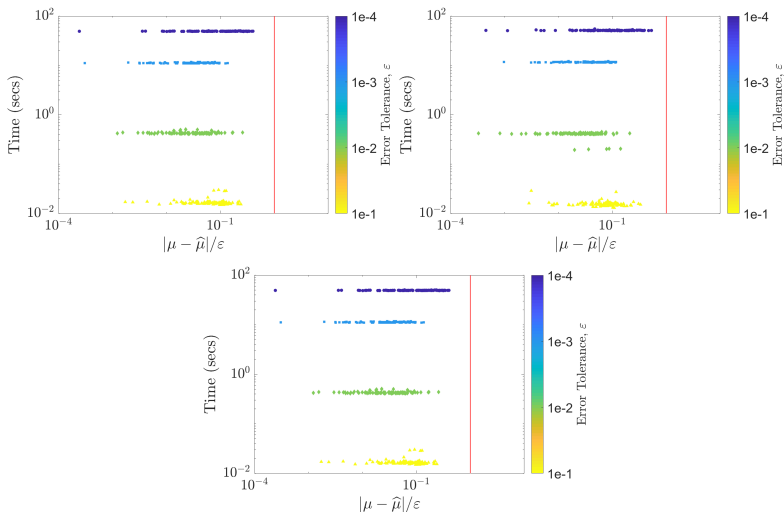


Figure: Option pricing using 1) Empirical Bayes, 2) GCV, 3) Full Bayes stopping criterion



Summary

- Developed a *general technique* for a **Fast transform kernel**
- Developed a **fast automatic Bayesian cubature** 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
- More about Guaranteed Automatic Algorithms (GAIL):
http://gailgithub.github.io/GAIL_Dev/



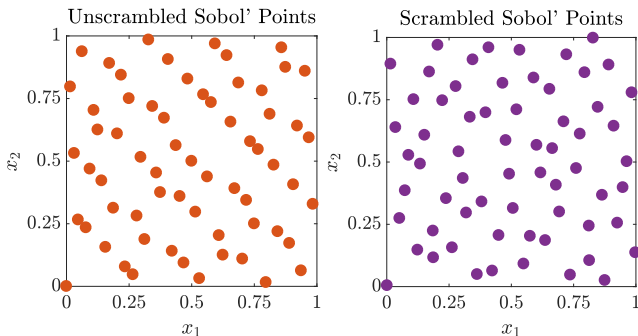
Future work

- Choosing the **kernel order r** and **periodization** transform automatically
- **Deterministic** interpretation of Bayesian cubature
- Broaden the choice of numerical examples
- Better handling of **conditioning** problem and numerical errors
- Sobol pointset and Fast Walsh Transform with smooth kernels ...



Future work : Sobol points and Fast walsh transform

- Using the established generalized theory for a *Fast transform kernel*, we could use other kernels with suitable point sets to achieve similar or better performance and accuracy.
- One such point sets to consider in future is, *Sobol points* and with appropriate choice of smooth kernel, should lead to *Fast Walsh Transform*.





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

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

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

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

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

Finally, the function approximation is

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

Thank you!



References I

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.

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.

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. 2003. *Bayesian Monte Carlo*, Advances in Neural Information Processing Systems, pp. 489–496.



References II

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