

Adaptive Probabilistic Numerical Methods Using Fast Transforms

Fred J. Hickernell & Jagadees Rathinavel

Department of Applied Mathematics

Center for Interdisciplinary Scientific Computation

Illinois Institute of Technology

hickernell@iit.edu mypages.iit.edu/~hickernell

Thanks to the organizers, the GAIL team, NSF-DMS-1522687 and NSF-DMS-1638521 (SAMSI)

SAMSI-Lloyds-Turing Workshop on Probabilistic Numerical Methods, April 11, 2018





When Do We Stop?

Compute the solution to a **linear** problem:

$$\text{sol}(f) = \begin{cases} \int_{\mathbb{R}^d} f(\mathbf{x}) \varrho(\mathbf{x}) d\mathbf{x} & \text{Bayesian inference, financial risk, statistical physics, ...} \\ f(\mathbf{x}) & \text{surrogate modeling for computer experiments, ...} \\ \vdots \end{cases}$$

Desired solution: An adaptive algorithm, $\text{app}(\cdot, \cdot)$ of the form

$$\text{app}(f, \varepsilon) = w_{0,n} + \sum_{i=1}^n w_{i,n} f(\mathbf{x}_i),$$

where n , $\{\mathbf{x}_i\}_{i=1}^\infty$, $w_{0,n}$, and $\mathbf{w} = (w_{i,n})_{i=1}^n$ are chosen to **guarantee**

$$|\text{sol}(f) - \text{app}(f, \varepsilon)| \leq \varepsilon \text{ with high probability} \quad \forall \varepsilon > 0, \text{ reasonable } f$$



The Probabilistic Numerics Approach

Assume $f \sim \mathcal{GP}(m, s^2 C_\theta)$, a sample from a Gaussian process. Defining

$$c = \text{sol}(\text{sol}(C_\theta(\cdot, \cdot))), \quad \mathbf{c} = (\text{sol}(C_\theta(\cdot, \mathbf{x}_1)), \dots, \text{sol}(C_\theta(\cdot, \mathbf{x}_n)))^T, \quad \mathbf{C} = (C_\theta(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n$$

and choosing the **weights** as

$$w_0 = m[\text{sol}(1) - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}], \quad \mathbf{w} = \mathbf{C}^{-1} \mathbf{c}, \quad \text{app}(f, \varepsilon) = w_0 + \mathbf{w}^T \mathbf{f}, \quad \mathbf{f} = (f(\mathbf{x}_i))_{i=1}^n.$$

yields an unbiased approximation:

$$\text{sol}(f) - \text{app}(f, \varepsilon) \mid \mathbf{f} = \mathbf{y} \sim \mathcal{N}\left(0, s^2(c - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})\right)$$

If n is chosen large enough to make

$$2.58s\sqrt{c - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}} \leq \varepsilon,$$

then we are **assured** that

$$\mathbb{P}_f [|\text{sol}(f) - \text{app}(f, \varepsilon)| \leq \varepsilon] \geq 99\%.$$



The Probabilistic Numerics Approach

Assume $f \sim \mathcal{GP}(m, s^2 C_\theta)$, a sample from a Gaussian process. Defining

$$c = \text{sol}(\text{sol}(C_\theta(\cdot, \cdot))), \quad \mathbf{c} = (\text{sol}(C_\theta(\cdot, \mathbf{x}_1)), \dots, \text{sol}(C_\theta(\cdot, \mathbf{x}_n)))^T, \quad \mathbf{C} = (C_\theta(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n$$

and choosing the **weights** as

$$w_0 = m[\text{sol}(1) - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{1}], \quad \mathbf{w} = \mathbf{C}^{-1} \mathbf{c}, \quad \text{app}(f, \varepsilon) = w_0 + \mathbf{w}^T \mathbf{f}, \quad \mathbf{f} = (f(\mathbf{x}_i))_{i=1}^n.$$

yields an unbiased approximation:

$$\text{sol}(f) - \text{app}(f, \varepsilon) \mid \mathbf{f} = \mathbf{y} \sim \mathcal{N}\left(0, s^2(c - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})\right)$$

If n is chosen large enough to make

$$2.58s\sqrt{c - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c}} \leq \varepsilon,$$

then we are **assured** that

$$\mathbb{P}_f [|\text{sol}(f) - \text{app}(f, \varepsilon)| \leq \varepsilon] \geq 99\%.$$

There are issues requiring attention.



Maximum Likelihood Estimation

Minimize minus twice the log likelihood observed with $f = y$:

$$2n \log(s) + \log(\det(C)) + s^{-2}(y - m\mathbf{1})^T C^{-1}(y - m\mathbf{1})$$

first with respect to m , then s , then θ :

$$m_{MLE} = \frac{\mathbf{1}^T C^{-1} y}{\mathbf{1}^T C^{-1} \mathbf{1}}, \quad s_{MLE}^2 = \frac{1}{n} y^T \left[C^{-1} - \frac{C^{-1} \mathbf{1} \mathbf{1}^T C^{-1}}{\mathbf{1}^T C^{-1} \mathbf{1}} \right] y$$

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmin}} \left\{ n \log \left(y^T \left[C^{-1} - \frac{C^{-1} \mathbf{1} \mathbf{1}^T C^{-1}}{\mathbf{1}^T C^{-1} \mathbf{1}} \right] y \right) + \log(\det(C)) \right\}$$



Maximum Likelihood Estimation

Minimize minus twice the log likelihood observed with $f = y$:

$$2n \log(s) + \log(\det(C)) + s^{-2}(y - m\mathbf{1})^T C^{-1}(y - m\mathbf{1})$$

first with respect to m , then s , then θ :

$$m_{MLE} = \frac{\mathbf{1}^T C^{-1} y}{\mathbf{1}^T C^{-1} \mathbf{1}}, \quad s_{MLE}^2 = \frac{1}{n} y^T \left[C^{-1} - \frac{C^{-1} \mathbf{1} \mathbf{1}^T C^{-1}}{\mathbf{1}^T C^{-1} \mathbf{1}} \right] y$$

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmin}} \left\{ n \log \left(y^T \left[C^{-1} - \frac{C^{-1} \mathbf{1} \mathbf{1}^T C^{-1}}{\mathbf{1}^T C^{-1} \mathbf{1}} \right] y \right) + \log(\det(C)) \right\}$$

Stopping criterion becomes

$$2.58 \sqrt{\underbrace{\frac{(c - c^T C^{-1} c)}{n}}_{\text{depends on design}} \underbrace{y^T \left[C^{-1} - \frac{C^{-1} \mathbf{1} \mathbf{1}^T C^{-1}}{\mathbf{1}^T C^{-1} \mathbf{1}} \right] y}_{\text{depends on data}}} \leq \varepsilon,$$



Maximum Likelihood Estimation

$$m_{\text{MLE}} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}}, \quad s_{\text{MLE}}^2 = \frac{1}{n} \mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}$$

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

Stopping criterion becomes

$$2.58 \sqrt{\underbrace{\frac{(c - \mathbf{c}^T \mathbf{C}^{-1} \mathbf{c})}{n}}_{\text{depends on design}} \underbrace{\mathbf{y}^T \left[\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{1} \mathbf{1}^T \mathbf{C}^{-1}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} \right] \mathbf{y}}_{\text{depends on data}}} \leq \varepsilon,$$

Q1: How large a family of kernels, C_θ is needed in practice to be confident in the error bound?

Q2: Are there better ways of finding the right θ , say cross-validation?

Q3: Can we check out assumption that f really comes from a Gaussian process? If not, are there alternatives



Low Discrepancy Sampling

Suppose that the domain is $[0, 1]^d$. Low discrepancy sampling places the x_i more evenly than IID sampling¹

IID points



Sobol' points



Integration lattice points



¹Dick, J. et al. High dimensional integration — the Quasi-Monte Carlo way. *Acta Numer.* **22**, 133–288 (2013), H., F. J. et al. *SAMSI Program on Quasi-Monte Carlo and High-Dimensional Sampling Methods for Applied Mathematics*.

<https://www.samsi.info/programs-and-activities/year-long-research-programs/2017-18-program-quasi-monte-carlo-high-dimensional-sampling-methods-applied-mathematics-qmc/>.



Covariance Kernels that Match the Design

Suppose that the **covariance kernel**, C_θ , and the **design**, $\{\mathbf{x}_i\}_{i=1}^n$, have special properties:

$$\mathbf{C} = \left(C_\theta(\mathbf{x}_i, \mathbf{x}_j) \right)_{i,j=1}^n = (\mathbf{C}_1, \dots, \mathbf{C}_n) = \frac{1}{n} \mathbf{V} \Lambda \mathbf{V}^H, \quad \mathbf{V}^H = n \mathbf{V}^{-1}, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) = \text{diag}(\boldsymbol{\lambda})$$

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

Suppose that $\mathbf{V}^T z$ is a **fast transform ($\mathcal{O}(n \log n)$ cost)** applied to z . Then it follows that

$$\boldsymbol{\lambda} = \mathbf{V}^T \mathbf{C}_1 \text{ is fast,} \quad \mathbf{C}^{-1} \mathbf{1} = \frac{\mathbf{1}}{\lambda_1}$$

Let \mathbf{y} be the observed function values. Recall $\mathbf{c} = \text{sol}(\text{sol}^\top(C_\theta(\cdot, \cdot)))$, and let

$$\hat{\mathbf{y}} = \mathbf{V}^T \mathbf{y}, \quad \hat{\mathbf{c}} = \mathbf{V}^T \mathbf{c}, \quad \mathbf{c} = (\text{sol}(C_\theta(\cdot, \mathbf{x}_1)), \dots, \text{sol}(C_\theta(\cdot, \mathbf{x}_n)))^T$$

Then using the MLE estimates, the approximate solution and the stopping criterion become:

$$\text{app}(f, \varepsilon) = \frac{\hat{y}_1 \text{sol}(1)}{n} + \sum_{i=2}^n \frac{\hat{y}_i^* \hat{c}_i}{\lambda_i}, \quad 2.58 \sqrt{\left(c - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i}} \leq \varepsilon$$



Covariance Kernels that Match the Design

$$c = \text{sol}^*(\text{sol}^*(C_\theta(\cdot, \cdot))), \quad C = \left(C_\theta(x_i, x_j) \right)_{i,j=1}^n = \frac{1}{n} V \Lambda V^H, \quad \Lambda = \text{diag}(\lambda), \quad V_1 = v_1 = 1$$

$$V^T z \text{ is } \mathcal{O}(n \log n), \quad \lambda = V^T C_1, \quad C^{-1} \mathbf{1} = \frac{1}{\lambda_1}$$

$$\hat{\mathbf{y}} = W^T \mathbf{y}, \quad \hat{\mathbf{c}} = W^T \mathbf{c}, \quad \mathbf{c} = (\text{sol}^*(C_\theta(\cdot, x_1)), \dots, \text{sol}^*(C_\theta(\cdot, x_n)))^T$$

$$\text{app}(f, \varepsilon) = \frac{\hat{y}_1 \text{sol}(1)}{n} + \sum_{i=2}^n \frac{\hat{y}_i^* \hat{c}_i}{\lambda_i}, \quad 2.58 \sqrt{\left(c - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i}} \leq \varepsilon$$

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

For integration with respect to a density and our special kernels, $\text{sol}(1) = 1$, $\mathbf{c} = 1$, and $c = 1$:

$$\text{app}(f, \varepsilon) = \frac{\hat{y}_1}{n}, \quad 2.58 \sqrt{\left(1 - \frac{n}{\lambda_1} \right) \frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i}} \leq \varepsilon$$



Covariance Kernels that Match the Design

$$c = \text{sol}^*(\text{sol}^*(C_\theta(\cdot, \cdot))), \quad C = \left(C_\theta(x_i, x_j) \right)_{i,j=1}^n = \frac{1}{n} V \Lambda V^H, \quad \Lambda = \text{diag}(\lambda), \quad V_1 = v_1 = 1$$

$$V^T z \text{ is } \mathcal{O}(n \log n), \quad \lambda = V^T C_1, \quad C^{-1} \mathbf{1} = \frac{1}{\lambda_1}$$

$$\hat{\mathbf{y}} = W^T \mathbf{y}, \quad \hat{\mathbf{c}} = W^T \mathbf{c}, \quad \mathbf{c} = (\text{sol}^*(C_\theta(\cdot, x_1)), \dots, \text{sol}^*(C_\theta(\cdot, x_n)))^T$$

$$\text{app}(f, \varepsilon) = \frac{\hat{y}_1 \text{sol}(1)}{n} + \sum_{i=2}^n \frac{\hat{y}_i^* \hat{c}_i}{\lambda_i}, \quad 2.58 \sqrt{\left(c - \frac{1}{n} \sum_{i=1}^n \frac{|\hat{c}_i|^2}{\lambda_i} \right) \frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i}} \leq \varepsilon$$

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

For integration with respect to a density and our special kernels, $\text{sol}(1) = 1$, $\mathbf{c} = 1$, and $c = 1$:

$$\text{app}(f, \varepsilon) = \frac{\hat{y}_1}{n}, \quad 2.58 \sqrt{\left(1 - \frac{n}{\lambda_1} \right) \frac{1}{n^2} \sum_{i=2}^n \frac{|\hat{y}_i|^2}{\lambda_i}} \leq \varepsilon$$

Q4: How do we avoid round-off in evaluating $1 - n/\lambda_1$?



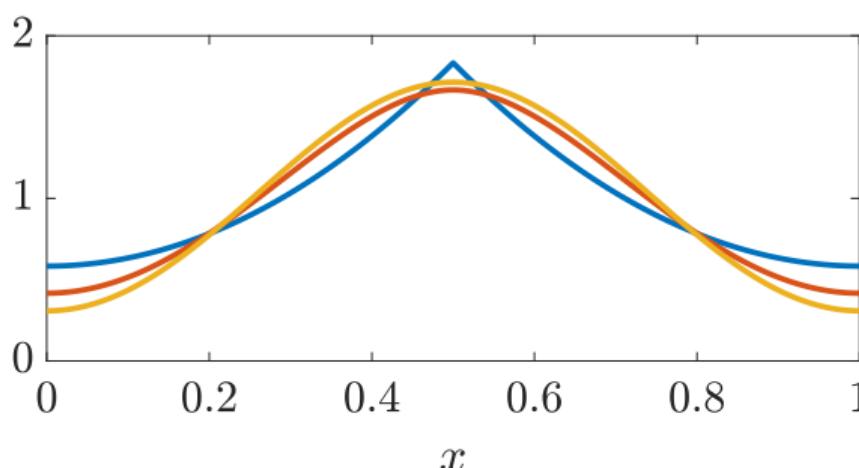
Form of Matching Covariance Kernels

Typically the domain of f is $[0, 1)^d$, and

$$C(\mathbf{x}, \mathbf{t}) = \begin{cases} \tilde{C}(\mathbf{x} - \mathbf{t} \bmod 1) & \text{integration lattices} \\ \tilde{C}(\mathbf{x} \oplus \mathbf{t}) & \text{Sobol' sequences, } \oplus \text{ means digitwise addition modulo 2} \end{cases}$$

E.g., for integration lattices

$$C(\mathbf{x}, \mathbf{t}) = \prod_{k=1}^d [1 - \theta_1(-1)^{\theta_2} B_{2\theta_2}(x_k - t_k \bmod 1)], \quad \theta_1 > 0, \theta_2 \in \mathbb{N}$$





Form of Matching Covariance Kernels

Typically the domain of f is $[0, 1)^d$, and

$$C(\mathbf{x}, \mathbf{t}) = \begin{cases} \tilde{C}(\mathbf{x} - \mathbf{t} \bmod 1) & \text{integration lattices} \\ \tilde{C}(\mathbf{x} \oplus \mathbf{t}) & \text{Sobol' sequences, } \oplus \text{ means digitwise addition modulo 2} \end{cases}$$

E.g., for integration lattices

$$C(\mathbf{x}, \mathbf{t}) = \prod_{k=1}^d [1 - \theta_1(-1)^{\theta_2} B_{2\theta_2}(x_k - t_k \bmod 1)], \quad \theta_1 > 0, \theta_2 \in \mathbb{N}$$

Q5: How do we periodize f to take advantage of integration lattices and smoother covariance kernels?
Does this even work for function approximation?

Q6: May we get higher order convergence with higher order nets and smoother kernels? What do those kernels look like?

Q7: Are low discrepancy designs also good for function approximation?

Q8: Are there other kernel/design combinations that expedite vector-matrix operations?



Option Pricing

$$\mu = \text{fair price} = \int_{\mathbb{R}^d} e^{-rT} \max \left(\frac{1}{d} \sum_{j=1}^d S_j - K, 0 \right) \frac{e^{-z^T z / 2}}{(2\pi)^{d/2}} dz \approx \$13.12$$

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

$$x = Az, \quad AA^T = \Sigma = \left(\min(i, j) T/d \right)_{i,j=1}^d, \quad T = 1/4, \quad d = 13 \text{ here}$$



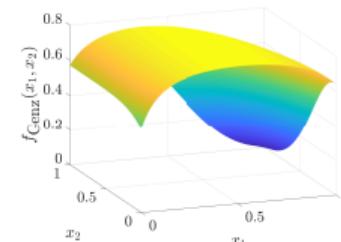
Abs. Error Tolerance	Method	Median Error	Accuracy	Worst 10% n	Worst 10% Time (s)
1E-2	IID	diff	2E-3	100%	6.1E7 33.000
1E-2	Scr. Sobol'	PCA	1E-3	100%	1.6E4 0.040
1E-2	Scr. Sob. cont. var.	PCA	2E-3	100%	4.1E3 0.019
1E-2	Bayes. Latt.	PCA	2E-3	99%	1.6E4 0.051



Gaussian Probability

$$\mu = \int_{[\mathbf{a}, \mathbf{b}]} \frac{\exp\left(-\frac{1}{2}\mathbf{t}^T \Sigma^{-1} \mathbf{t}\right)}{\sqrt{(2\pi)^d \det(\Sigma)}} d\mathbf{t} = \int_{[0,1]^{d-1}} f(\mathbf{x}) d\mathbf{x}$$

For some typical choice of $\mathbf{a}, \mathbf{b}, \Sigma, d = 3; \mu \approx 0.6763$



Tolerance	Method	Median		n	Worst 10% Time (s)
		Error	Accuracy		
1E-2	IID	5E-4	100%	8.1E4	0.020
1E-2	Scr. Sobol'	4E-5	100%	1.0E3	0.005
1E-2	Bayes. Latt.	5E-5	100%	4.1E3	0.023
1E-3	IID	9E-5	100%	2.0E6	0.400
1E-3	Scr. Sobol'	2E-5	100%	2.0E3	0.006
1E-3	Bayes. Latt.	3E-7	100%	6.6E4	0.076
1E-4	Scr. Sobol'	4E-7	100%	1.6E4	0.018
1E-4	Bayes. Latt.	6E-9	100%	5.2E5	0.580
1E-4	Bayes. Latt. Smth.	1E-7	100%	3.3E4	0.047



Recap of Questions

- Q1:** How large a family of kernels, C_Θ is needed in practice to be confident in the error bound?
- Q2:** Are there better ways of finding the right Θ , say cross-validation?
- Q3:** Can we check out assumption that f really comes from a Gaussian process? If not, are there alternatives
- Q4:** How do we avoid round-off in evaluating $1 - n/\lambda_1$?
- Q5:** How do we periodize f to take advantage of integration lattices and smoother covariance kernels? Does this even work for function approximation?
- Q6:** May we get higher order convergence with higher order nets and smoother kernels? What do those kernels look like?
- Q7:** Are low discrepancy designs also good for function approximation?
- Q8:** Are there other kernel/design combinations that expedite vector-matrix operations?
- Q9:** Is this adaptive Bayesian algorithm competitive with others?
- Q10:** What other problems are amenable to matched kernels and designs beyond cubature and function approximation?

Thank you

Slides available at [www.slideshare.net/fjhickernell/
probabilistic-numerics-2018-April-11-london](http://www.slideshare.net/fjhickernell/probabilistic-numerics-2018-April-11-london)



Dick, J., Kuo, F. & Sloan, I. H. High dimensional integration — the Quasi-Monte Carlo way. *Acta Numer.* **22**, 133–288 (2013).



H., F. J., Kuo, F. Y., L'Ecuyer, P. & Owen, A. B. *SAMSI Program on Quasi-Monte Carlo and High-Dimensional Sampling Methods for Applied Mathematics.*

<https://www.samsi.info/programs-and-activities/year-long-research-programs/2017-18-program-quasi-monte-carlo-high-dimensional-sampling-methods-applied-mathematics-qmc/>.