

Fast Gaussian Process Regression with Derivative Information

SIAM UQ 2024

Aleksei G. Sorokin & Fred J. Hickernell

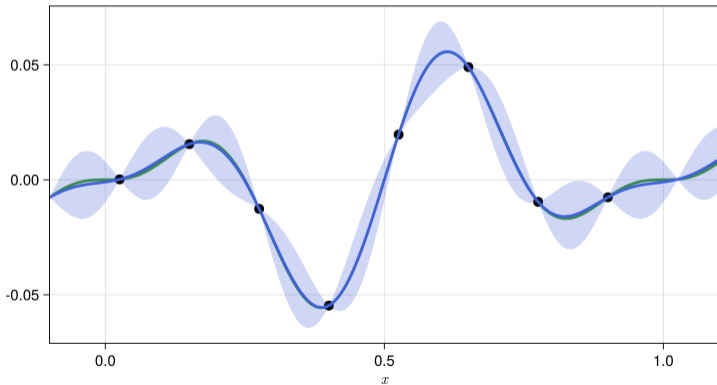
Illinois Institute of Technology, Department of Applied Mathematics

February 27, 2024

Why Gaussian Process Regression (GPR)?

- Encode simulation knowledge into model via kernel e.g. smoothness or periodicity
- Provides a distribution over simulations i.e. quantifies uncertainty in predictions

— $f^{(0)}(x)$ ● $(y_i^{(0)})_{i=1}^s$ — $m_n^{(0)}(x)$ ■ $m_n^{(0)}(x) \pm 1.96 \sigma_n^{(p)}(x)$



Motivation

Challenge: A GPR model costs $\mathcal{O}(n^3)$ to fit

- Applications often require GPR for $n > 10,000$ nodes which is very costly

Solution: Matching nodes and kernel reduces costs to $\mathcal{O}(n \log n)$

- ★ Require control over design of experiments

Observation: Derivative information can enhance GPR

- Derivatives available for free e.g. simulation is the numerical solution of a PDE
- Derivatives may be available at a nominal cost e.g. via automatic differentiation
- Derivatives may be the primary information source e.g. GPR for solving non-linear PDEs [Chen et al., 2021]

New Challenge: With m derivatives, can we improve the $\mathcal{O}(n^3 m^3)$ fitting cost?

New Solution: Exploit additional structure to reduce cost to $\mathcal{O}(m^2 n \log n + m^3 n)$

Outline and Related Work

1. **GPR**: follows book on GP for Machine Learning [Rasmussen et al., 2006]
2. **Fast GPR**: follows fast Bayesian cubature of Hickernell and Jagadeeswaran
 - Flavor #1: lattice sequence designs [Jagadeeswaran and Hickernell, 2019]
 - Flavor #2: digital sequence designs [Jagadeeswaran and Hickernell, 2022]
 - Unifying thesis of Jagadeeswaran Rathinavel [Rathinavel, 2019]
 - Adjacent work in a RKHS with lattice sequences [Kaarnioja et al., 2022]
 - Application to surrogate for PDE with random coefficients [Sorokin et al., 2023]
3. **GPR with derivative information**: incorporated gradients in [Solak et al., 2002]
4. **Fast GPR with derivative information**: [our novel contribution!](#)

Gaussian Process Regression

- Given *simulation* $f : [0, 1]^s \rightarrow \mathbb{R}$
- Assume simulation an instance of a *Gaussian process*, $f \sim GP(0, K)$
 - Assume prior mean is zero (not necessary but simplifies presentation)
 - *Prior covariance kernel* $K : [0, 1]^s \times [0, 1]^s \rightarrow \mathbb{R}$ is symmetric positive definite

$$K(\mathbf{x}, \mathbf{x}') = \text{Cov}[f(\mathbf{x}), f(\mathbf{x}')]]$$

- *Sampling sequence* $\mathbf{X} = (\mathbf{x}_i)_{i=1}^n \in [0, 1]^{n \times s}$
- *Observations* $\mathbf{y} = (y_i)_{i=1}^n = (f(\mathbf{x}_i) + \varepsilon_i)_{i=1}^n \in \mathbb{R}^{n \times 1}$ with *noise* $\varepsilon_1, \dots, \varepsilon_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \zeta)$
- *kernel (Gram) matrix* $\mathbf{K} = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n \in \mathbb{R}^{n \times n}$
- *kernel vector* $\mathbf{k}_{\mathbf{X}}(\mathbf{x}) = (K(\mathbf{x}, \mathbf{x}_i))_{i=1}^n \in \mathbb{R}^{n \times 1}$

$$\text{Posterior Mean: } m_n(\mathbf{x}) = \mathbf{k}_{\mathbf{X}}^{\top}(\mathbf{x})(\mathbf{K} + \zeta \mathbf{I})^{-1} \mathbf{y}$$

$$\text{Posterior Covariance: } K_n(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}, \mathbf{x}') - \mathbf{k}_{\mathbf{X}}^{\top}(\mathbf{x})(\mathbf{K} + \zeta \mathbf{I})^{-1} \mathbf{k}_{\mathbf{X}}(\mathbf{x}')$$

Posterior Mean: $m_n(\mathbf{x}) = \mathbf{k}_X^\top(\mathbf{x})(\mathbf{K} + \zeta\mathbf{I})^{-1}\mathbf{y}$

Posterior Covariance: $K_n(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}, \mathbf{x}') - \mathbf{k}_X^\top(\mathbf{x})(\mathbf{K} + \zeta\mathbf{I})^{-1}\mathbf{k}_X(\mathbf{x}')$

Key is to solve systems of the form

$$(\mathbf{K} + \zeta\mathbf{I})\mathbf{a} = \mathbf{b}$$

for $\mathbf{a} \in \mathbb{C}^n$ where $\mathbf{b} \in \mathbb{R}^n$

- $(\mathbf{K} + \zeta\mathbf{I})^{-1}\mathbf{y}$ precomputed during fitting, typically costs $\mathcal{O}(n^3)$
- $(\mathbf{K} + \zeta\mathbf{I})^{-1}\mathbf{k}_X(\mathbf{x}')$ computed when evaluating uncertainty, typically costs $\mathcal{O}(n^2)$ after precomputing factorization of $\mathbf{K} + \zeta\mathbf{I}$

Fast Gaussian Process Regression

What? Induce structure in $K + \zeta I$ so solving $(K + \zeta I)\mathbf{a} = \mathbf{b}$ for \mathbf{a} costs $\mathcal{O}(n \log n)$

How? Match quasi-random sequences with structured kernels [Rathinavel, 2019]

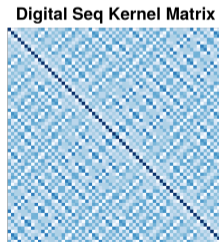
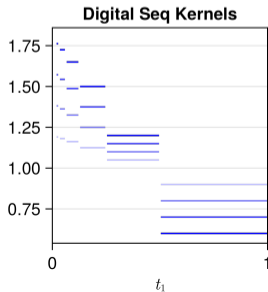
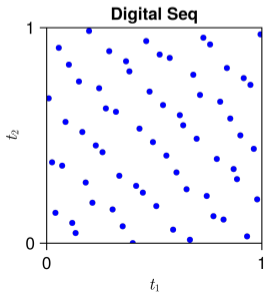
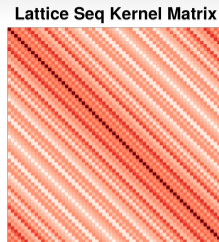
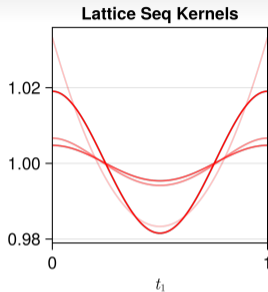
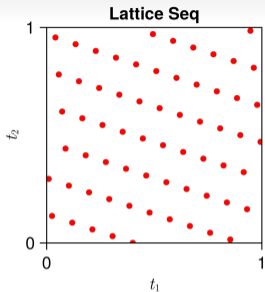
- K circulant with [lattice sequence](#) X and shift invariant kernel

$$K(\mathbf{x}, \mathbf{x}') = K((\mathbf{x} - \mathbf{x}') \bmod 1)$$

- K block-Toeplitz with [digital sequence](#) X and digitally shift invariant kernel

$$K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x} \ominus \mathbf{x}')$$

where \ominus is XOR (exclusive or) of base 2 digits



For circulant or block-Toeplitz \mathbf{K} we have

- $\mathbf{K} + \zeta\mathbf{I}$ inherits same structure as \mathbf{K}
- Eigendecomposition $\mathbf{K} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\dagger$ with $\mathbf{V}^{-1} = \mathbf{V}^\dagger =$ Hermitian of \mathbf{V}
- $\mathcal{F}(\mathbf{a}) := \mathbf{V}^\dagger\mathbf{a}$ and $\mathcal{F}^{-1}(\mathbf{b}) := \mathbf{V}\mathbf{b}$ can be computed in $\mathcal{O}(n\log n)$
 - Circulant \mathbf{K} means $\mathcal{F}(\mathbf{a})$ is the fast Fourier transform of \mathbf{a}
 - Block-Toeplitz \mathbf{K} means $\mathcal{F}(\mathbf{a})$ is the fast Walsh-Hadamard transform of \mathbf{a}
- First column of \mathbf{V} is $\mathbf{1}/\sqrt{n}$

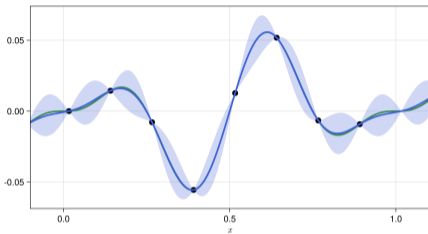
Solve $(\mathbf{K} + \zeta\mathbf{I})\mathbf{a} = \mathbf{b}$ for \mathbf{a} at cost $\mathcal{O}(n\log n)$ with

$$\mathbf{a} = \mathcal{F}^{-1}\left(\frac{\mathcal{F}(\mathbf{b})}{\boldsymbol{\lambda} + \zeta}\right)$$

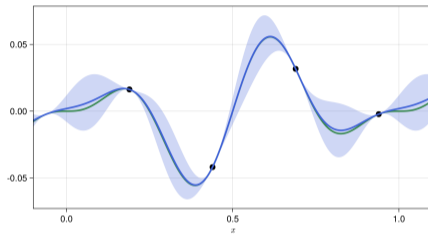
where $\boldsymbol{\lambda} = \text{diag}(\mathbf{\Lambda}) = \sqrt{n}\mathcal{F}(\mathbf{k}_\chi(\mathbf{x}_1))$ and the division is done elementwise

Derivative Informed Gaussian Process Regression

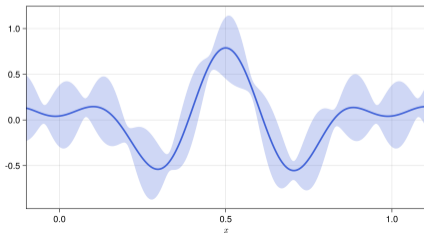
— $f^{(0)}(x)$ ● $(y^{(0)})_{t-1}^*$ — $m_t^{(0)}(x)$ ■ $m_t^{(0)}(x) \pm 1.96 \sigma_t^{(0)}(x)$



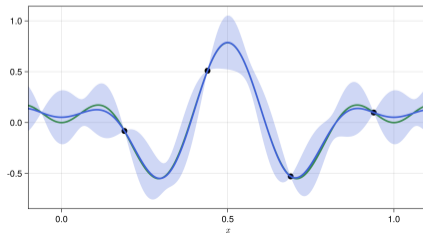
— $f^{(1)}(x)$ ● $(y^{(1)})_{t-1}^*$ — $m_t^{(1)}(x)$ ■ $m_t^{(1)}(x) \pm 1.96 \sigma_t^{(1)}(x)$



— $m_t^{(2)}(x)$ ■ $m_t^{(2)}(x) \pm 1.96 \sigma_t^{(2)}(x)$



— $f^{(2)}(x)$ ● $(y^{(2)})_{t-1}^*$ — $m_t^{(2)}(x)$ ■ $m_t^{(2)}(x) \pm 1.96 \sigma_t^{(2)}(x)$



Linear functional of a Gaussian process is still a Gaussian process

$$f^{(\beta)}(\mathbf{x}) := \frac{\partial^{|\beta|}}{\partial \mathbf{x}^\beta} f(\mathbf{x}) := \frac{\partial^{|\beta|}}{\partial x_1^{\beta_1} \dots \partial x_s^{\beta_s}} f(\mathbf{x})$$

$$\text{Cov}[f^{(\beta)}(\mathbf{x}), f^{(\beta')}(\mathbf{x}')] = \frac{\partial^{|\beta|}}{\partial \mathbf{x}^\beta} \frac{\partial^{|\beta'|}}{\partial \mathbf{x}'^{\beta'}} \text{Cov}[f(\mathbf{x}), f(\mathbf{x}')] =: K^{(\beta, \beta')}(\mathbf{x}, \mathbf{x}')$$

With m derivative multi-indices $\beta_1, \dots, \beta_m \in \mathbb{N}_0^s$ the kernel (Gram) matrix becomes

$$\mathbf{K} = \begin{pmatrix} K^{(\beta_1, \beta_1)} & \dots & K^{(\beta_1, \beta_m)} \\ \vdots & \ddots & \vdots \\ K^{(\beta_m, \beta_1)} & \dots & K^{(\beta_m, \beta_m)} \end{pmatrix} \in \mathbb{R}^{nm \times nm}, \quad K^{(\beta_k, \beta_l)} = \left(K^{(\beta_k, \beta_l)}(\mathbf{x}_i, \mathbf{x}_j) \right)_{i,j=1}^n$$

so solving $(\mathbf{K} + \zeta \mathbf{I})\mathbf{a} = \mathbf{b}$ for $\mathbf{a} \in \mathbb{C}^{mn}$ where $\mathbf{b} \in \mathbb{R}^{mn}$ costs $\mathcal{O}(m^3 n^3)$ in general

Fast Gaussian Process Regression with Derivative Information

$K(\beta_k, \beta_l)$ retains structure of $K^{(0,0)}$ e.g. circulant or block Toeplitz

$$\begin{pmatrix} K(\beta_1, \beta_1) & \dots & K(\beta_1, \beta_m) \\ \vdots & \ddots & \vdots \\ K(\beta_m, \beta_1) & \dots & K(\beta_m, \beta_m) \end{pmatrix} = \begin{pmatrix} V & & \\ & \ddots & \\ & & V \end{pmatrix} \underbrace{\begin{pmatrix} \Lambda(\beta_1, \beta_1) & \dots & \Lambda(\beta_1, \beta_m) \\ \vdots & \ddots & \vdots \\ \Lambda(\beta_m, \beta_1) & \dots & \Lambda(\beta_m, \beta_m) \end{pmatrix}}_{\Lambda \in \mathbb{R}^{nm \times nm}} \begin{pmatrix} V^\dagger & & \\ & \ddots & \\ & & V^\dagger \end{pmatrix}$$

Let \otimes be the Kronecker product so

$$K + \zeta I = (I \otimes V)(\Lambda + \zeta I)(I \otimes V^\dagger)$$

$$K + \zeta I = (I \otimes V)(\Lambda + \zeta I)(I \otimes V^\dagger)$$

Since Λ is a diagonal block (striped) matrix, there is some permutation matrix P with

$$P^T(\Lambda + \zeta I)P = \Upsilon + \zeta I$$

where

$$\Upsilon = \begin{pmatrix} \Upsilon_1 & & \\ & \ddots & \\ & & \Upsilon_n \end{pmatrix}$$

is block diagonal with $\Upsilon_{i,kl} = \lambda_i^{(\beta_k, \beta_l)}$. Then

$$K + \zeta I = (I \otimes V)P(\Upsilon + \zeta I)P^T(I \otimes V^\dagger)$$

Cost of solving $(\mathbf{K} + \zeta\mathbf{I})\mathbf{a} = \mathbf{b}$ for \mathbf{a} with structured $\mathbf{K} + \zeta\mathbf{I}$

$$\mathbf{K} + \zeta\mathbf{I} = (\mathbf{I} \otimes \mathbf{V})\mathbf{P}(\Upsilon + \zeta\mathbf{I})\mathbf{P}^\top(\mathbf{I} \otimes \mathbf{V}^\dagger)$$

Reduce cost from $\mathcal{O}(m^3n^3)$ to $\mathcal{O}(m^2n \log n + m^3n)$ with the following algorithm

1. Constructing Υ from eigenvalues $\lambda^{(\beta_k, \beta_l)} = \mathcal{F}\left(\mathbf{k}_X^{(\beta_k, \beta_l)}(\mathbf{x}_1)\right)$ costs $\mathcal{O}(m^2n \log n)$
2. $\check{\mathbf{b}} := \mathbf{P}^\top(\mathbf{I} \otimes \mathbf{V}^\dagger)\mathbf{b}$ can be computed at cost $\mathcal{O}(mn \log n)$
3. $\check{\mathbf{a}} := (\Upsilon + \zeta\mathbf{I})^{-1}\check{\mathbf{b}}$ can be computed at cost $\mathcal{O}(m^3n)$
4. $\mathbf{a} = (\mathbf{I} \otimes \mathbf{V})\mathbf{P}\check{\mathbf{a}}$ can be computed at cost $\mathcal{O}(mn \log n)$

Fast Kernel Parameter Optimization

K often depends on parameters $\boldsymbol{\theta}$ e.g. scaling factor, lengthscales, noise variance ζ
 $\boldsymbol{\theta}$ which maximizes the marginal log likelihood is

$$\begin{aligned}\operatorname{argmin}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}|\mathbf{y}) &= \operatorname{argmin}_{\boldsymbol{\theta}} \left[\log \det(\mathbf{K} + \zeta\mathbf{I}) + \mathbf{y}^\top (\mathbf{K} + \zeta\mathbf{I})^{-1} \mathbf{y} \right] \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{i=1}^n \left[\log \det(\Upsilon_i + \zeta\mathbf{I}) + \check{\mathbf{y}}_i^\dagger (\Upsilon_i + \zeta\mathbf{I})^{-1} \check{\mathbf{y}}_i \right]\end{aligned}$$

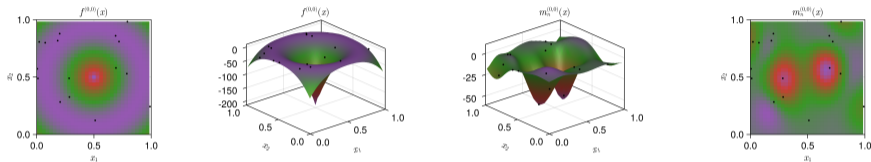
where

$$\check{\mathbf{y}} := \begin{pmatrix} \check{\mathbf{y}}_1 \\ \vdots \\ \check{\mathbf{y}}_n \end{pmatrix} := \mathbf{P}^\top (\mathbf{I} \otimes \mathbf{V}^\dagger) \mathbf{y}.$$

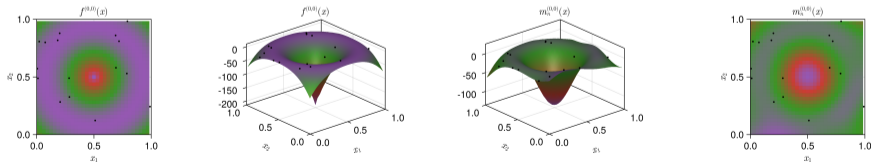
Both $L(\boldsymbol{\theta}|\mathbf{y})$ and $\partial_{\theta_j} L(\boldsymbol{\theta}|\mathbf{y})$ can still be computed in $\mathcal{O}(m^2 n \log n + m^3 n)$

Analytic Donut¹ in UMBridge [Seelinger et al., 2023]

IID Points, SE Kernel: No Gradient Information

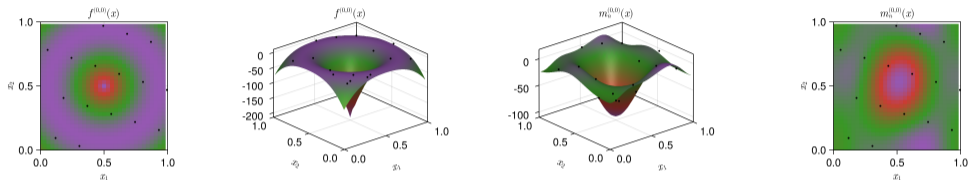


IID Points with SE Kernel: With Gradient Information

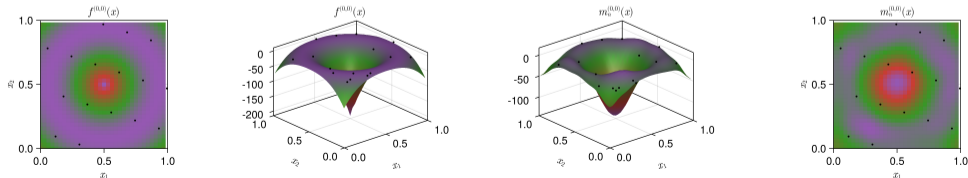


¹<https://um-bridge-benchmarks.readthedocs.io/en/docs/inverse-benchmarks/analytic-donut.html>

Lattice Points, Matching Kernel: No Gradient Information



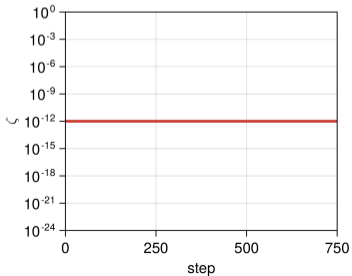
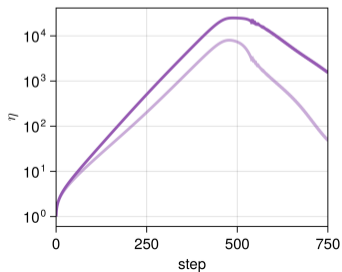
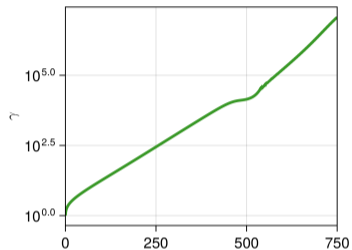
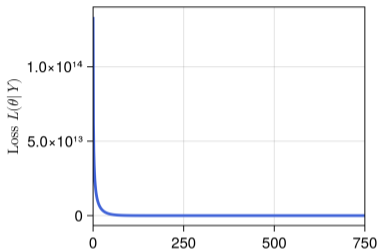
Lattice Points, Matching Kernel: With Gradient Information



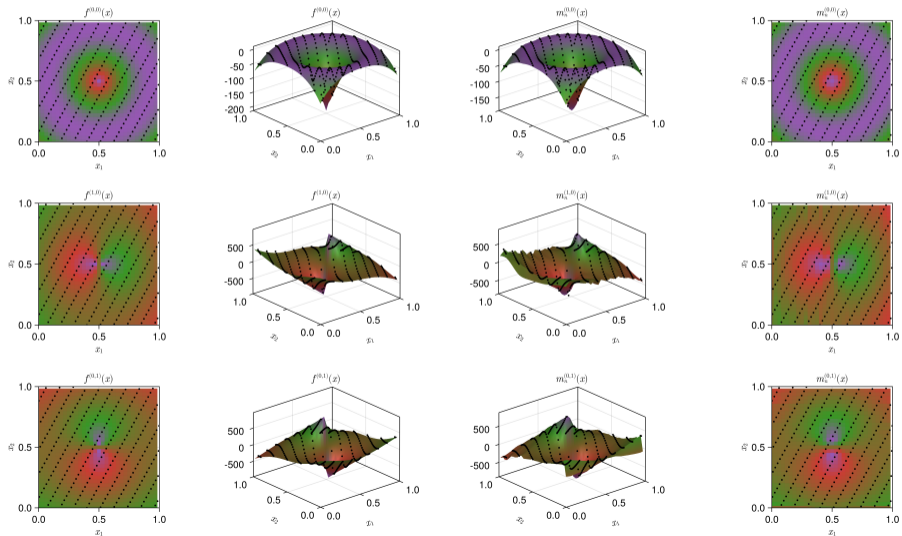
Lattice GP with Gradients for Donut Example

```
1 import FastGaussianProcesses; import QMCGenerators; import UMBridge
2
3 # docker run -it -p 4243:4243 linusseelinger/benchmark-analytic-donut
4 model = UMBridge.HTTPModel("posterior", "http://localhost:4243")
5 β = [0 0; 1 0; 0 1] # observe  $f^{\{(0,0)\}}$ ,  $f^{\{(1,0)\}}$ ,  $f^{\{(0,1)\}}$ 
6
7 function f(x::Vector{Float64})
8     z = [6*x[1]-3, 6*x[2]-3] # both componenets between -3 and 3
9     [    UMBridge.evaluate(model, [z], Dict())[1][1], #  $f^{\{(0,0)\}}$ 
10      (UMBridge.gradient(model, 0, 0, [z], [1]) .* [6, 6])...] #  $f^{\{(1,0)\}}$ ,  $f^{\{(0,1)\}}$ 
11 end
12
13 seq = QMCGenerators.RandomShift(QMCGenerators.LatticeSeqB2(2), 1, 7) # s=2, seed=7
14 gp = FastGaussianProcesses.FastGaussianProcess(f, seq, 2^8; β=β, optim_steps=750) # n=2^10
15 FastGaussianProcesses.plot_gp_optimization(gp, figpath=joinpath(@_DIR_, "optim.png"))
16 FastGaussianProcesses.plot_gp_2s(gp; f=f, β=β, figpath=joinpath(@_DIR_, "gp.png"))
```

Kernel Parameter Optimization



Gaussian Process and Gradient Visualization



Future Work

Theory

- Can we improve the $\mathcal{O}(m^2n \log n + m^3n)$ cost by relating $\lambda^{(\beta, \kappa)}$ to $\lambda^{(\beta', \kappa')}$?
- Link with RKHS setting
 - General GPR and RKHS kernel interpolation connections in [Kanagawa et al., 2018]
 - Optimize weights in [Kaarnioja et al., 2022] with GPR kernel parameter optimization
- Analogous developments for digital sequences

Practical Software

- QMCGenerators.jl²: Quasi-random sequence generators with randomizations
- FastGaussianProcesses.jl³: Fast GPR with derivatives (in development)
- QMCPy⁴ [Choi et al., 2022]
 - Quasi-random sequence generators with randomizations
 - Fast GPR cubature [Rathinavel, 2019]

²<https://github.com/alegresor/QMCGenerators.jl>

³<https://github.com/alegresor/FastGaussianProcesses.jl>

⁴<https://github.com/QMCSoftware/QMCSoftware>

References I

- Yifan Chen, Bamdad Hosseini, Houman Owhadi, and Andrew M. Stuart. Solving and learning nonlinear pdes with gaussian processes. *Journal of Computational Physics*, 447:110668, 2021. ISSN 0021-9991. doi:
<https://doi.org/10.1016/j.jcp.2021.110668>. URL <https://www.sciencedirect.com/science/article/pii/S0021999121005635>.
- Sou-Cheng T. Choi, Fred J. Hickernell, Rathinavel Jagadeeswaran, Michael J. McCourt, and Aleksei G. Sorokin. Quasi-monte carlo software. In Alexander Keller, editor, *Monte Carlo and Quasi-Monte Carlo Methods*, pages 23–47, Cham, 2022. Springer International Publishing. ISBN 978-3-030-98319-2.
- R. Jagadeeswaran and Fred J. Hickernell. Fast automatic bayesian cubature using lattice sampling. *Statistics and Computing*, 29(6):1215–1229, Sep 2019. ISSN 1573-1375. doi: 10.1007/s11222-019-09895-9. URL <http://dx.doi.org/10.1007/s11222-019-09895-9>.

References II

- Rathinavel Jagadeeswaran and Fred J Hickernell. Fast automatic bayesian cubature using sobol sampling. In *Advances in Modeling and Simulation: Festschrift for Pierre L'Ecuyer*, pages 301–318. Springer, 2022.
- Vesa Kaarnioja, Yoshihito Kazashi, Frances Y Kuo, Fabio Nobile, and Ian H Sloan. Fast approximation by periodic kernel-based lattice-point interpolation with application in uncertainty quantification. *Numerische Mathematik*, 150(1):33–77, 2022.
- Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*, 2018.
- Carl Edward Rasmussen, Christopher KI Williams, et al. *Gaussian processes for machine learning*, volume 1. Springer, 2006.
- Jagadeeswaran Rathinavel. *Fast automatic Bayesian cubature using matching kernels and designs*. Illinois Institute of Technology, 2019.

References III

- Linus Seelinger, Vivian Cheng-Seelinger, Andrew Davis, Matthew Parno, and Anne Reinartz. Um-bridge: Uncertainty quantification and modeling bridge. *Journal of Open Source Software*, 8(83):4748, 2023.
- Ercan Solak, Roderick Murray-Smith, WE Leithead, D Leith, and Carl Rasmussen. Derivative observations in gaussian process models of dynamic systems. *Advances in neural information processing systems*, 15, 2002.
- Aleksei G. Sorokin, Aleksandra Pachalieva, Daniel O'Malley, James M. Hyman, Fred J. Hickernell, and Nicolas W. Hengartner. Computationally efficient and error aware surrogate construction for numerical solutions of subsurface flow through porous media, 2023.