

An Introduction to Probabilistic Numerical Methods

Chris. J. Oates

School of Mathematics, Statistics and Physics @ Newcastle University
Programme on Data-Centric Engineering @ Alan Turing Institute

October 2017 @ Turing

- optimisation
- integration
- linear algebra
- solution of differential equations
- ...

What is the fuss all about?

The goal:

Numerical Task \implies Finite Computation \implies Distribution on Output

Optimisation

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

- Well-defined:
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$.

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

- **Well-defined:**

- $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .

- **Well-posed:**

- Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
- Aim to minimise $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$.

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

- Well-defined:
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$.

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

- Well-defined:
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$.

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

- Well-defined:
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$.

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

- Well-defined:
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2$.

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = \mathbf{x}_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(x_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = x_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(x_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = \mathbf{x}_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(x_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = x_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(x_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = \mathbf{x}_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(\mathbf{x}_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = \mathbf{x}_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(\mathbf{x}_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = \mathbf{x}_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(\mathbf{x}_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

$$\mathbf{x}^* = \arg \max f(\mathbf{x})$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. gradient ascent with estimated gradients?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{\mathbf{x}}^*$.
 - The empirical maximum $\hat{\mathbf{x}}^* = \mathbf{x}_{i^*}$ where $i^* = \arg \max_{i=1, \dots, n} f(\mathbf{x}_i)$?
 - Something better?
- Key idea: Estimator uncertainty quantification!

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$$

Bayesian linear regression onto a basis $\{\phi_i\}_{i=1}^m$:

$$f(\mathbf{x}) = \beta_1\phi_1(\mathbf{x}) + \cdots + \beta_m\phi_m(\mathbf{x})$$

with $n \leq m \in \mathbb{N} \cup \{\infty\}$.

- Prior $p(\beta_1, \dots, \beta_m)$
- Likelihood $\prod_{i=1}^n \delta(f(\mathbf{x}_i) - \beta_1\phi_1(\mathbf{x}_i) - \cdots - \beta_m\phi_m(\mathbf{x}_i))$
- Posterior $p(\beta_1, \dots, \beta_m | \mathcal{D})$
- Posterior marginal $p(x^* | \mathcal{D})$

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$$

Bayesian linear regression onto a basis $\{\phi_i\}_{i=1}^m$:

$$f(\mathbf{x}) = \beta_1\phi_1(\mathbf{x}) + \cdots + \beta_m\phi_m(\mathbf{x})$$

with $n \leq m \in \mathbb{N} \cup \{\infty\}$.

- Prior $p(\beta_1, \dots, \beta_m)$
- Likelihood $\prod_{i=1}^n \delta(f(\mathbf{x}_i) - \beta_1\phi_1(\mathbf{x}_i) - \cdots - \beta_m\phi_m(\mathbf{x}_i))$
- Posterior $p(\beta_1, \dots, \beta_m | \mathcal{D})$
- Posterior marginal $p(\mathbf{x}^* | \mathcal{D})$

Calculations for the conjugate set-up:

•

•

•

•

•

•

•

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$

-

-

-

-

-

-

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$

- Prior:

-

-

-

-

-

Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta}|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.

Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta} | \lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:



Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta}|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\boldsymbol{\beta}|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \boldsymbol{\Phi}^\top (\mathbf{I} + \boldsymbol{\Phi} \boldsymbol{\Phi}^\top)^{-1} \boldsymbol{\Phi} \mathbf{f}) (\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\boldsymbol{\Phi}]_{ij} = \phi_j(\mathbf{x}_i)$.

•

•
•
•

Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta}|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\boldsymbol{\beta}|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \boldsymbol{\Phi}^\top (\mathbf{I} + \boldsymbol{\Phi} \boldsymbol{\Phi}^\top)^{-1} \boldsymbol{\Phi} \mathbf{f}) (\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\boldsymbol{\Phi}]_{ij} = \phi_j(\mathbf{x}_i)$.

- Posterior marginal: $\mathbf{x}^*|\mathcal{D} \sim ?$

-
-
-

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda I)$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\beta|\mathcal{D} \sim \text{MVT} \left((I + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) (I + \Phi^\top \Phi)^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\Phi]_{ij} = \phi_j(\mathbf{x}_i)$.

- Posterior marginal: $\mathbf{x}^*|\mathcal{D} \sim ?$
 - Draw β from $\beta|\mathcal{D}$
 - Evaluate $\mathbf{x}^* = \arg \max \beta_1 \phi_1(\mathbf{x}) + \dots + \beta_m \phi_m(\mathbf{x})$
 - Repeat.

Compute $\mathbf{x}^* = \arg \max f(\mathbf{x})$:

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(x^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(x^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(x^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

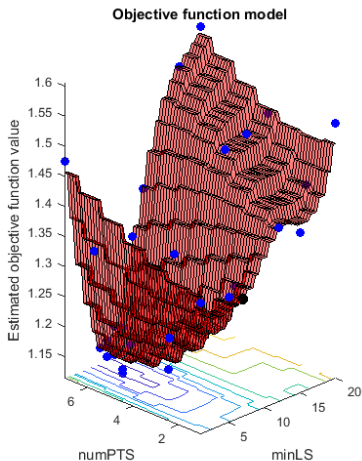
- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(x^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(\mathbf{x}^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(\mathbf{x}^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(\mathbf{x}^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.

- Close connection between statistics and design of numerical optimisation methods.
- Similar to “Bayesian optimisation” (Mockus, 1989).
- Kernel trick maps to Gaussian processes.
- The distributional output $p(\mathbf{x}^*|\mathcal{D})$ provides uncertainty quantification.
- Propagation and the Bayesian mantra of Dawid.
- Numerical analysts want to consider order of convergence and constants (of the point estimator).
- Similar considerations relevant to posterior contraction.



<https://uk.mathworks.com/help/stats/tune-random-forest-using-quantile-error-and-bayesian-optimization.html>
(but the MATLAB function doesn't provide uncertainty quantification!)

Integration

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:
 - $f \in L_2(\pi)$?
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathcal{X} \subset \mathbb{R}^d$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:

- $f \in L_2(\pi)$?
- $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathcal{X} \subset \mathbb{R}^d$ a compact subset of \mathbb{R}^d .

- Well-posed:

- Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
- Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:
 - $f \in L_2(\pi)$?
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathcal{X} \subset \mathbb{R}^d$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:
 - $f \in L_2(\pi)$?
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:
 - $f \in L_2(\pi)$?
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:
 - $f \in L_2(\pi)$?
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

- Well-defined:
 - $f \in L_2(\pi)$?
 - $f \in C^\alpha(\mathcal{X})$ for some $\alpha \geq 0$ and $\mathbf{x} \in \mathcal{X}$ a compact subset of \mathbb{R}^d .
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs which you can select.
 - Aim to minimise $|\hat{I} - I|$.

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(x_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(x_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(x_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(x_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(x_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(x_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$I = \int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$$

Two distinct requirements:

- A method to select the integrand evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n$.
 - Uniform grid over \mathcal{X} ?
 - Adaptive selection, e.g. based on local error indicators?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{I}$.
 - The arithmetic mean $\hat{I} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$$

Bayesian linear regression onto a basis $\{\phi_i\}_{i=1}^m$:

$$f(\mathbf{x}) = \beta_1\phi_1(\mathbf{x}) + \cdots + \beta_m\phi_m(\mathbf{x})$$

with $n \leq m \in \mathbb{N} \cup \{\infty\}$.

- Prior $p(\beta_1, \dots, \beta_m)$
- Likelihood $\prod_{i=1}^n \delta(f(\mathbf{x}_i) - \beta_1\phi_1(\mathbf{x}_i) - \cdots - \beta_m\phi_m(\mathbf{x}_i))$
- Posterior $p(\beta_1, \dots, \beta_n | \mathcal{D})$
- Posterior marginal $p(l | \mathcal{D})$

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$$

Bayesian linear regression onto a basis $\{\phi_i\}_{i=1}^m$:

$$f(\mathbf{x}) = \beta_1\phi_1(\mathbf{x}) + \cdots + \beta_m\phi_m(\mathbf{x})$$

with $n \leq m \in \mathbb{N} \cup \{\infty\}$.

- Prior $p(\beta_1, \dots, \beta_m)$
- Likelihood $\prod_{i=1}^n \delta(f(\mathbf{x}_i) - \beta_1\phi_1(\mathbf{x}_i) - \cdots - \beta_m\phi_m(\mathbf{x}_i))$
- Posterior $p(\beta_1, \dots, \beta_n | \mathcal{D})$
- Posterior marginal $p(I | \mathcal{D})$

Calculations for the conjugate set-up:

•

•

•

•

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$



Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior:

•

•

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta|\lambda \sim N(\mathbf{0}, \lambda I)$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta | \lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta}|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\boldsymbol{\beta}|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \boldsymbol{\Phi}^\top (\mathbf{I} + \boldsymbol{\Phi} \boldsymbol{\Phi}^\top)^{-1} \boldsymbol{\Phi} \mathbf{f}) (\mathbf{I} + \boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\boldsymbol{\Phi}]_{ij} = \phi_j(\mathbf{x}_i)$.

•

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\beta|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) (\mathbf{I} + \Phi^\top \Phi)^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\Phi]_{ij} = \phi_j(\mathbf{x}_i)$.

- Posterior marginal:

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\beta|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) (\mathbf{I} + \Phi^\top \Phi)^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\Phi]_{ij} = \phi_j(\mathbf{x}_i)$.

- Posterior marginal:

$$\begin{aligned} I|\mathcal{D} \sim \text{Student-T} & \left(\Psi^\top (\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \right. \\ & \left. \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (\mathbf{I} + \Phi^\top \Phi)^{-1} \Psi, n \right) \end{aligned}$$

where $[\Psi]_j = \int \phi_j(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}$.

Compute $\int f(\mathbf{x})\pi(\mathbf{x})d\mathbf{x}$:

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{x_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{x_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select x_n to minimise (*) based on $\{x_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{x_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{x_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select x_n to minimise (*) based on $\{x_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_p(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_p(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{x_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n} (f^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi f) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{x_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select x_n to minimise (*) based on $\{x_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{x_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{x_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select x_n to minimise (*) based on $\{x_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_p(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_p(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{x_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{x_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select x_n to minimise (*) based on $\{x_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_p(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_p(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

- Similar to “Bayesian quadrature” (O’Hagan, 1991).
- Kernel trick maps to GPs.
- Theoretical results were provided in Briol *et al.* 2016:
 - Posterior mean \hat{I} satisfies $\hat{I} - I = O_P(n^{-\alpha/d+\epsilon})$.
 - Stronger assumptions on f , such as $f \in H^\alpha(0, 1) \otimes \dots \otimes H^\alpha(0, 1)$, lead to $\hat{I} - I = O_P(n^{-\alpha+\epsilon})$ for an appropriately sparse basis $\{\phi_i\}_{i=1}^n$.
 - Posterior is concentrated on \hat{I} , so rates of contraction to I can also be established.
- Posterior mean often coincides with standard quadrature methods (Särkkä *et al.*, 2016).
- How to select the $\{\mathbf{x}_i\}_{i=1}^n$?
 - Recall the posterior scale was determined by

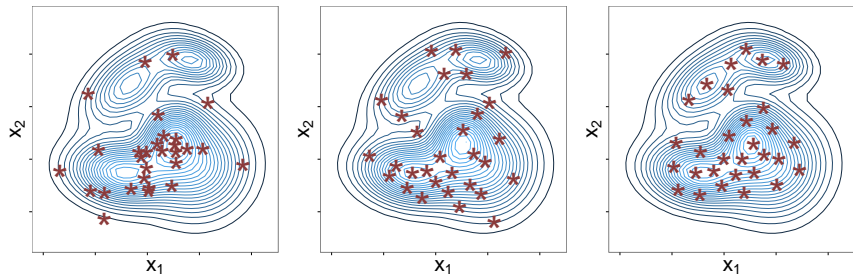
$$\frac{1}{n}(\mathbf{f}^\top \Phi^\top (I + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi$$

- Select $\{\mathbf{x}_i\}_{i=1}^n$ to minimise

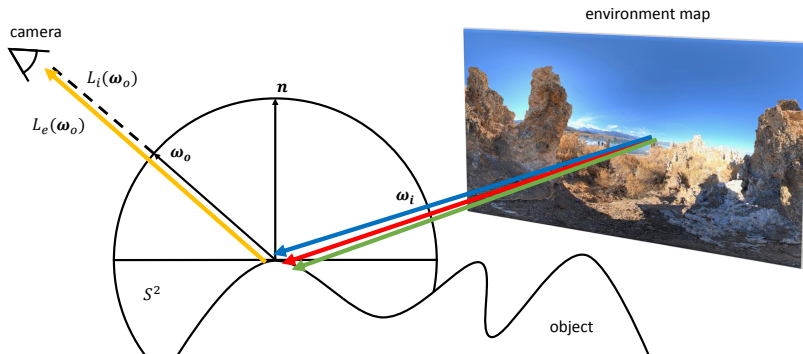
$$(*) \quad \Psi^\top (I + \Phi^\top \Phi)^{-1} \Psi ?$$

- Select \mathbf{x}_n to minimise (*) based on $\{\mathbf{x}_i\}_{i=1}^{n-1}$ fixed? (Similar to “sequential Bayesian quadrature”.)

From Briol *et al.*, 2016:



One of these is (a variant on) sequential Bayesian quadrature - which one?



From Briol *et al.*, 2016.

Linear Algebra

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

- Well-defined:
 - \mathbf{A} is a $N \times N$ symmetric positive definite matrix.
- Well-posed:
 - Allowed $n \ll N$ matrix-vector multiplications.
 - Represented as $[\mathbf{s}_i^\top \mathbf{A}] \mathbf{x} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$.
 - You are allowed to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Aim to minimise $\|\hat{\mathbf{x}} - \mathbf{x}\|_{\mathbf{A}}$ where $\|\mathbf{z}\|_{\mathbf{A}} = \sqrt{\mathbf{z}^\top \mathbf{A} \mathbf{z}}$.

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $\mathbf{s}_i^\top \mathbf{A} \hat{\mathbf{x}} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(s_i^\top \mathbf{A}, s_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $s_i^\top \mathbf{A} \hat{\mathbf{x}} = s_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(s_i^\top \mathbf{A}, s_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $s_i^\top \mathbf{A} \hat{\mathbf{x}} = s_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(s_i^\top \mathbf{A}, s_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $s_i^\top \mathbf{A} \hat{\mathbf{x}} = s_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $\mathbf{s}_i^\top \mathbf{A} \hat{\mathbf{x}} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $\mathbf{s}_i^\top \mathbf{A} \hat{\mathbf{x}} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $\mathbf{s}_i^\top \mathbf{A} \hat{\mathbf{x}} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\mathbf{Ax} = \mathbf{b}$$

Two distinct requirements:

- A method to select the directions $\mathbf{s}_1, \dots, \mathbf{s}_n$.
 - Random projections?
 - Sequential selection, e.g. gradient descent or conjugate gradient?
- An estimator $\{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n \mapsto \hat{\mathbf{x}}$.
 - A minimal $\|\cdot\|_2$ norm vector that satisfies $\mathbf{s}_i^\top \mathbf{A} \hat{\mathbf{x}} = \mathbf{s}_i^\top \mathbf{b}$ for $i = 1, \dots, n$?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

Recap: Conjugate Gradient Method

Let $\langle \cdot, \cdot \rangle_{\mathbf{A}}$ be the inner-product induced by $\| \cdot \|_{\mathbf{A}}$. (i.e. $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = \mathbf{z}^{\top} \mathbf{A} \tilde{\mathbf{z}}$.)



Recap: Conjugate Gradient Method

Let $\langle \cdot, \cdot \rangle_{\mathbf{A}}$ be the inner-product induced by $\| \cdot \|_{\mathbf{A}}$. (i.e. $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = \mathbf{z}^{\top} \mathbf{A} \tilde{\mathbf{z}}$.)

Call $\mathbf{z}, \tilde{\mathbf{z}} \in \mathbb{R}^N$ conjugate (w.r.t. \mathbf{A}) if $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = 0$.



Recap: Conjugate Gradient Method

Let $\langle \cdot, \cdot \rangle_{\mathbf{A}}$ be the inner-product induced by $\| \cdot \|_{\mathbf{A}}$. (i.e. $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = \mathbf{z}^{\top} \mathbf{A} \tilde{\mathbf{z}}$.)

Call $\mathbf{z}, \tilde{\mathbf{z}} \in \mathbb{R}^N$ conjugate (w.r.t. \mathbf{A}) if $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = 0$.

Suppose that we have a conjugate basis $\{\mathbf{s}_1, \dots, \mathbf{s}_N\}$ for \mathbb{R}^N (i.e. $\langle \mathbf{s}_i, \mathbf{s}_j \rangle_{\mathbf{A}} = 0$ for all $i \neq j$).



Recap: Conjugate Gradient Method

Let $\langle \cdot, \cdot \rangle_{\mathbf{A}}$ be the inner-product induced by $\| \cdot \|_{\mathbf{A}}$. (i.e. $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = \mathbf{z}^{\top} \mathbf{A} \tilde{\mathbf{z}}$.)

Call $\mathbf{z}, \tilde{\mathbf{z}} \in \mathbb{R}^N$ conjugate (w.r.t. \mathbf{A}) if $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = 0$.

Suppose that we have a conjugate basis $\{\mathbf{s}_1, \dots, \mathbf{s}_N\}$ for \mathbb{R}^N (i.e. $\langle \mathbf{s}_i, \mathbf{s}_j \rangle_{\mathbf{A}} = 0$ for all $i \neq j$).

Consider the natural sequence of approximations

$$\hat{\mathbf{x}}_n = \sum_{i=1}^n \alpha_i \mathbf{s}_i$$

where each

$$\alpha_i = \frac{\langle \mathbf{s}_i, \mathbf{x} \rangle_{\mathbf{A}}}{\langle \mathbf{s}_i, \mathbf{s}_i \rangle_{\mathbf{A}}} = \frac{\mathbf{s}_i^{\top} \mathbf{b}}{\mathbf{s}_i^{\top} \mathbf{A} \mathbf{s}_i}$$

can be computed in $O(N^2)$. The total computational cost is $O(nN^2)$.



Recap: Conjugate Gradient Method

Let $\langle \cdot, \cdot \rangle_{\mathbf{A}}$ be the inner-product induced by $\| \cdot \|_{\mathbf{A}}$. (i.e. $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = \mathbf{z}^{\top} \mathbf{A} \tilde{\mathbf{z}}$.)

Call $\mathbf{z}, \tilde{\mathbf{z}} \in \mathbb{R}^N$ conjugate (w.r.t. \mathbf{A}) if $\langle \mathbf{z}, \tilde{\mathbf{z}} \rangle_{\mathbf{A}} = 0$.

Suppose that we have a conjugate basis $\{\mathbf{s}_1, \dots, \mathbf{s}_N\}$ for \mathbb{R}^N (i.e. $\langle \mathbf{s}_i, \mathbf{s}_j \rangle_{\mathbf{A}} = 0$ for all $i \neq j$).

Consider the natural sequence of approximations

$$\hat{\mathbf{x}}_n = \sum_{i=1}^n \alpha_i \mathbf{s}_i$$

where each

$$\alpha_i = \frac{\langle \mathbf{s}_i, \mathbf{x} \rangle_{\mathbf{A}}}{\langle \mathbf{s}_i, \mathbf{s}_i \rangle_{\mathbf{A}}} = \frac{\mathbf{s}_i^{\top} \mathbf{b}}{\mathbf{s}_i^{\top} \mathbf{A} \mathbf{s}_i}$$

can be computed in $O(N^2)$. The total computational cost is $O(nN^2)$.

So what is needed to proceed?

- Need a smart choice of $\{\mathbf{s}_1, \dots, \mathbf{s}_N\}$.
- For theory, need to bound $\|\mathbf{x} - \hat{\mathbf{x}}_n\| = \|\sum_{i=n+1}^N \alpha_i \mathbf{x}_i\|$ in your favourite $\| \cdot \|$.

Recap: Conjugate Gradient Method

Aim is to adaptively select s_n based on the computation up to iteration $n - 1$.

Gradient Descent: Notice that x is a minimum of

$$f(x) = \frac{1}{2}x^\top Ax - x^\top b.$$

This suggests to select $s_n = -\nabla f(\hat{x}_{n-1})$ which is equal to $r_{n-1} = b - A\hat{x}_{n-1}$. However, this does not ensure $\{s_1, \dots, s_n\}$ is a conjugate set.

Conjugate Gradient: A more delicate procedure selects

$$s_n = r_{n-1} - \sum_{i < n} \frac{s_i^\top A r_{n-1}}{s_i^\top A s_i} s_i$$

i.e. gradient descent plus Gram-Schmidt orthogonalisation w.r.t $\langle \cdot, \cdot \rangle_A$ to subtract off components in the directions $\{s_1, \dots, s_{n-1}\}$ already used.

Recap: Conjugate Gradient Method

Aim is to adaptively select s_n based on the computation up to iteration $n - 1$.

Gradient Descent: Notice that x is a minimum of

$$f(x) = \frac{1}{2}x^\top Ax - x^\top b.$$

This suggests to select $s_n = -\nabla f(\hat{x}_{n-1})$ which is equal to $r_{n-1} = b - A\hat{x}_{n-1}$. However, this does not ensure $\{s_1, \dots, s_n\}$ is a conjugate set.

Conjugate Gradient: A more delicate procedure selects

$$s_n = r_{n-1} - \sum_{i < n} \frac{s_i^\top A r_{n-1}}{s_i^\top A s_i} s_i$$

i.e. gradient descent plus Gram-Schmidt orthogonalisation w.r.t $\langle \cdot, \cdot \rangle_A$ to subtract off components in the directions $\{s_1, \dots, s_{n-1}\}$ already used.

Recap: Conjugate Gradient Method

Aim is to adaptively select s_n based on the computation up to iteration $n - 1$.

Gradient Descent: Notice that x is a minimum of

$$f(x) = \frac{1}{2}x^\top Ax - x^\top b.$$

This suggests to select $s_n = -\nabla f(\hat{x}_{n-1})$ which is equal to $r_{n-1} = b - A\hat{x}_{n-1}$. However, this does not ensure $\{s_1, \dots, s_n\}$ is a conjugate set.

Conjugate Gradient: A more delicate procedure selects

$$s_n = r_{n-1} - \sum_{i < n} \frac{s_i^\top A r_{n-1}}{s_i^\top A s_i} s_i$$

i.e. gradient descent plus Gram-Schmidt orthogonalisation w.r.t $\langle \cdot, \cdot \rangle_A$ to subtract off components in the directions $\{s_1, \dots, s_{n-1}\}$ already used.

Recap: Conjugate Gradient Method

Aim is to adaptively select s_n based on the computation up to iteration $n - 1$.

Gradient Descent: Notice that x is a minimum of

$$f(x) = \frac{1}{2}x^\top Ax - x^\top b.$$

This suggests to select $s_n = -\nabla f(\hat{x}_{n-1})$ which is equal to $r_{n-1} = b - A\hat{x}_{n-1}$. However, this does not ensure $\{s_1, \dots, s_n\}$ is a conjugate set.

Conjugate Gradient: A more delicate procedure selects

$$s_n = r_{n-1} - \sum_{i < n} \frac{s_i^\top A r_{n-1}}{s_i^\top A s_i} s_i$$

i.e. gradient descent plus Gram-Schmidt orthogonalisation w.r.t $\langle \cdot, \cdot \rangle_A$ to subtract off components in the directions $\{s_1, \dots, s_{n-1}\}$ already used.

For either method, the computational cost of selecting s_n is $O(N^2)$, so the overall computational overhead added is $O(nN^2)$; the same order as random projections.

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n$$

Deploy full Bayesian inference for x :

- Prior $p(x)$
- Likelihood $\prod_{i=1}^n \delta(\mathbf{s}_i^\top \mathbf{A}x - \mathbf{s}_i^\top \mathbf{b})$
- Posterior $p(x|\mathcal{D})$

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{s}_i^\top \mathbf{A}, \mathbf{s}_i^\top \mathbf{b})\}_{i=1}^n$$

Deploy full Bayesian inference for \mathbf{x} :

- Prior $p(\mathbf{x})$
- Likelihood $\prod_{i=1}^n \delta(\mathbf{s}_i^\top \mathbf{A} \mathbf{x} - \mathbf{s}_i^\top \mathbf{b})$
- Posterior $p(\mathbf{x} | \mathcal{D})$

Calculations for the conjugate set-up:

•

•

•

•

•

Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$



Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$
- Prior:

•

•

•

Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$
- Prior: $\mathbf{x}|\lambda \sim \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.

Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$
- Prior: $\mathbf{x}|\lambda \sim \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:



Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$
- Prior: $\mathbf{x}|\lambda \sim \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\mathbf{x}|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{b}, \right. \\ \left. \frac{1}{n} (\mathbf{b}^\top \mathbf{S}^\top \mathbf{A}^\top \mathbf{S}^\top (\mathbf{I} + \mathbf{S} \mathbf{A} \mathbf{A}^\top \mathbf{S})^{-1} \mathbf{S} \mathbf{A} \mathbf{S} \mathbf{b}) (\mathbf{I} + \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A})^{-1}, n \right)$$



Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$
- Prior: $\mathbf{x}|\lambda \sim \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\mathbf{x}|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{b}, \right. \\ \left. \frac{1}{n} (\mathbf{b}^\top \mathbf{S}^\top \mathbf{A}^\top \mathbf{S}^\top (\mathbf{I} + \mathbf{S} \mathbf{A} \mathbf{A}^\top \mathbf{S})^{-1} \mathbf{S} \mathbf{A} \mathbf{S} \mathbf{b}) (\mathbf{I} + \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A})^{-1}, n \right)$$

- This is for general \mathbf{S} .

•

Calculations for the conjugate set-up:

- Let: $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_n]^\top$
- Prior: $\mathbf{x}|\lambda \sim \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\mathbf{x}|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{b}, \frac{1}{n} (\mathbf{b}^\top \mathbf{S}^\top \mathbf{A}^\top \mathbf{S}^\top (\mathbf{I} + \mathbf{S} \mathbf{A} \mathbf{A}^\top \mathbf{S})^{-1} \mathbf{S} \mathbf{A} \mathbf{S} \mathbf{b}) (\mathbf{I} + \mathbf{A}^\top \mathbf{S}^\top \mathbf{S} \mathbf{A})^{-1}, n \right)$$

- This is for general \mathbf{S} .
- For the conjugate gradient method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$ we have the orthogonality equation $\mathbf{S} \mathbf{A} \mathbf{A}^\top \mathbf{S}^\top = \mathbf{I}$ and the above can be further simplified.

[no video for this one!]

- Approximate linear solvers used extensively in engineering applications.
- Also relevant in statistics, e.g. simulation of spatial random fields.
- It turns out that the posterior mean in our construction coincides with the classical conjugate gradient (CG) method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$.
- Thus the classical error bounds for CG are inherited.
- The full posterior can be computed in $O(nN^2)$.

- **Approximate linear solvers used extensively in engineering applications.**
- Also relevant in statistics, e.g. simulation of spatial random fields.
- It turns out that the posterior mean in our construction coincides with the classical conjugate gradient (CG) method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$.
- Thus the classical error bounds for CG are inherited.
- The full posterior can be computed in $O(nN^2)$.

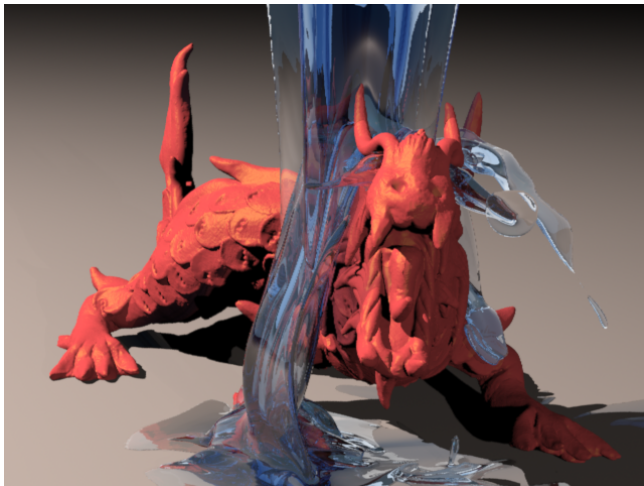
- Approximate linear solvers used extensively in engineering applications.
- Also relevant in statistics, e.g. simulation of spatial random fields.
- It turns out that the posterior mean in our construction coincides with the classical conjugate gradient (CG) method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$.
- Thus the classical error bounds for CG are inherited.
- The full posterior can be computed in $O(nN^2)$.

- Approximate linear solvers used extensively in engineering applications.
- Also relevant in statistics, e.g. simulation of spatial random fields.
- It turns out that the posterior mean in our construction coincides with the classical conjugate gradient (CG) method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$.
- Thus the classical error bounds for CG are inherited.
- The full posterior can be computed in $O(nN^2)$.

- Approximate linear solvers used extensively in engineering applications.
- Also relevant in statistics, e.g. simulation of spatial random fields.
- It turns out that the posterior mean in our construction coincides with the classical conjugate gradient (CG) method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$.
- Thus the classical error bounds for CG are inherited.
- The full posterior can be computed in $O(nN^2)$.

- Approximate linear solvers used extensively in engineering applications.
- Also relevant in statistics, e.g. simulation of spatial random fields.
- It turns out that the posterior mean in our construction coincides with the classical conjugate gradient (CG) method applied to the pre-conditioned system $\mathbf{A}^\top \mathbf{A} \mathbf{x} = \mathbf{A}^\top \mathbf{b}$.
- Thus the classical error bounds for CG are inherited.
- The full posterior can be computed in $O(nN^2)$.

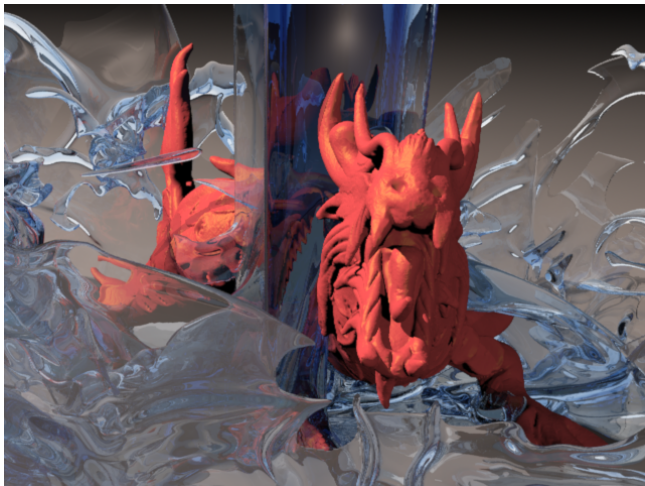
From McAdams *et al.*, SIGGRAPH 2010:



From McAdams *et al.*, SIGGRAPH 2010:



From McAdams *et al.*, SIGGRAPH 2010:



Solution of Differential Equations

E.g.:

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

- Well-defined:
 - $\mathcal{X} \subset \mathbb{R}^d$ be $C^{1,1}$
 - $f \in L^p(\mathcal{X})$, $p > n/2$
 - Cor. 9.18 in Gilbarg and Trudinger ensures $\exists!$ solution $u \in W_{\text{loc}}^{2,p}(\mathcal{X}) \cap C^0(\mathcal{X} \cup \partial\mathcal{X})$
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs $x_1, \dots, x_n \in \mathcal{X} \cup \partial\mathcal{X}$ which you can select.
 - Aim to minimise $\int_{\mathcal{X}} \|\hat{u}(x) - u(x)\|_2^2 dx$.

E.g.:

$$\begin{aligned}\Delta u(x) &= f(x), & x \in \mathcal{X} \\ u(x) &= 0, & x \in \partial\mathcal{X}\end{aligned}$$

- Well-defined:

- $\mathcal{X} \subset \mathbb{R}^d$ be $C^{1,1}$
- $f \in L^p(\mathcal{X})$, $p > n/2$
- Cor. 9.18 in Gilbarg and Trudinger ensures $\exists!$ solution $u \in W_{\text{loc}}^{2,p}(\mathcal{X}) \cap C^0(\mathcal{X} \cup \partial\mathcal{X})$

- Well-posed:

- Allowed n evaluations of $f(\cdot)$ at inputs $x_1, \dots, x_n \in \mathcal{X} \cup \partial\mathcal{X}$ which you can select.
- Aim to minimise $\int_{\mathcal{X}} \|\hat{u}(x) - u(x)\|_2^2 dx$.

E.g.:

$$\begin{aligned}\Delta u(x) &= f(x), & x \in \mathcal{X} \\ u(x) &= 0, & x \in \partial\mathcal{X}\end{aligned}$$

- Well-defined:
 - $\mathcal{X} \subset \mathbb{R}^d$ be $C^{1,1}$
 - $f \in L^p(\mathcal{X})$, $p > n/2$
 - Cor. 9.18 in Gilbarg and Trudinger ensures $\exists!$ solution $u \in W_{\text{loc}}^{2,p}(\mathcal{X}) \cap C^0(\mathcal{X} \cup \partial\mathcal{X})$
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs $x_1, \dots, x_n \in \mathcal{X} \cup \partial\mathcal{X}$ which you can select.
 - Aim to minimise $\int_{\mathcal{X}} \|\hat{u}(x) - u(x)\|_2^2 dx$.

E.g.:

$$\begin{aligned}\Delta u(x) &= f(x), & x \in \mathcal{X} \\ u(x) &= 0, & x \in \partial\mathcal{X}\end{aligned}$$

- Well-defined:
 - $\mathcal{X} \subset \mathbb{R}^d$ be $C^{1,1}$
 - $f \in L^p(\mathcal{X})$, $p > n/2$
 - Cor. 9.18 in Gilbarg and Trudinger ensures $\exists!$ solution $u \in W_{\text{loc}}^{2,p}(\mathcal{X}) \cap C^0(\mathcal{X} \cup \partial\mathcal{X})$
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs $x_1, \dots, x_n \in \mathcal{X} \cup \partial\mathcal{X}$ which you can select.
 - Aim to minimise $\int_{\mathcal{X}} \|\hat{u}(x) - u(x)\|_2^2 dx$.

E.g.:

$$\begin{aligned}\Delta u(x) &= f(x), & x \in \mathcal{X} \\ u(x) &= 0, & x \in \partial\mathcal{X}\end{aligned}$$

- Well-defined:
 - $\mathcal{X} \subset \mathbb{R}^d$ be $C^{1,1}$
 - $f \in L^p(\mathcal{X})$, $p > n/2$
 - Cor. 9.18 in Gilbarg and Trudinger ensures $\exists!$ solution $u \in W_{\text{loc}}^{2,p}(\mathcal{X}) \cap C^0(\mathcal{X} \cup \partial\mathcal{X})$
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs $x_1, \dots, x_n \in \mathcal{X} \cup \partial\mathcal{X}$ which you can select.
 - Aim to minimise $\int_{\mathcal{X}} \|\hat{u}(x) - u(x)\|_2^2 dx$.

E.g.:

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

- Well-defined:
 - $\mathcal{X} \subset \mathbb{R}^d$ be $C^{1,1}$
 - $f \in L^p(\mathcal{X})$, $p > n/2$
 - Cor. 9.18 in Gilbarg and Trudinger ensures $\exists!$ solution $u \in W_{\text{loc}}^{2,p}(\mathcal{X}) \cap C^0(\mathcal{X} \cup \partial\mathcal{X})$
- Well-posed:
 - Allowed n evaluations of $f(\cdot)$ at inputs $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$ which you can select.
 - Aim to minimise $\int_{\mathcal{X}} \|\hat{u}(\mathbf{x}) - u(\mathbf{x})\|_2^2 d\mathbf{x}$.

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(x_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(x_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(x_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(x_i, f(x_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(x_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(\mathbf{x}_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(\mathbf{x}_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(\mathbf{x}_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

$$\begin{aligned}\Delta u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} \in \mathcal{X} \\ u(\mathbf{x}) &= 0, & \mathbf{x} \in \partial\mathcal{X}\end{aligned}$$

Two distinct requirements:

- A method to select the function evaluation locations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{X} \cup \partial\mathcal{X}$.
 - Corners of a mesh on $\mathcal{X} \cup \partial\mathcal{X}$?
 - An adaptive method?
- An estimator $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n \mapsto \hat{u}(\cdot)$.
 - Linear interpolation of the $f(\mathbf{x}_i)$ and then solution of the PDE?
 - Something better?
- Key idea (again): Estimator uncertainty quantification!

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$$

Bayesian linear regression onto a basis $\{\phi_i\}_{i=1}^m$:

$$f(\mathbf{x}) = \beta_1\phi_1(\mathbf{x}) + \cdots + \beta_m\phi_m(\mathbf{x})$$

with $n \leq m \in \mathbb{N} \cup \{\infty\}$.

- Prior $p(\beta_1, \dots, \beta_m)$
- Likelihood $\prod_{i=1}^n \delta(f(\mathbf{x}_i) - \beta_1\phi_1(\mathbf{x}_i) - \cdots - \beta_m\phi_m(\mathbf{x}_i))$
- Posterior $p(\beta_1, \dots, \beta_m | \mathcal{D})$
- Posterior marginal $p(u(\cdot) | \mathcal{D})$

Start with the data that have been collected:

$$\mathcal{D} = \{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$$

Bayesian linear regression onto a basis $\{\phi_i\}_{i=1}^m$:

$$f(\mathbf{x}) = \beta_1\phi_1(\mathbf{x}) + \cdots + \beta_m\phi_m(\mathbf{x})$$

with $n \leq m \in \mathbb{N} \cup \{\infty\}$.

- Prior $p(\beta_1, \dots, \beta_m)$
- Likelihood $\prod_{i=1}^n \delta(f(\mathbf{x}_i) - \beta_1\phi_1(\mathbf{x}_i) - \cdots - \beta_m\phi_m(\mathbf{x}_i))$
- Posterior $p(\beta_1, \dots, \beta_n | \mathcal{D})$
- Posterior marginal $p(u(\cdot) | \mathcal{D})$

Calculations for the conjugate set-up:

•

•

•

•

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$



Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$

- Prior:

-

-

Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta}|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.

Calculations for the conjugate set-up:

- Let: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)$
- Prior: $\boldsymbol{\beta}|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

5

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta|\lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\beta|\mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) (\mathbf{I} + \Phi^\top \Phi)^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\Phi]_{ij} = \phi_j(\mathbf{x}_i)$.

5

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta | \lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\beta | \mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) (\mathbf{I} + \Phi^\top \Phi)^{-1}, n \right)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$ and $[\Phi]_{ij} = \phi_j(\mathbf{x}_i)$.

- Posterior marginal:

Calculations for the conjugate set-up:

- Let: $\beta = (\beta_1, \dots, \beta_m)$
- Prior: $\beta | \lambda \sim \mathbf{N}(\mathbf{0}, \lambda \mathbf{I})$ and $\lambda \sim p(\lambda) \propto \lambda^{-1}$.
- Posterior:

$$\beta | \mathcal{D} \sim \text{MVT} \left((\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) (\mathbf{I} + \Phi^\top \Phi)^{-1}, n \right)$$

where $\mathbf{f} = (f(x_1), \dots, f(x_n))$ and $[\Phi]_{ij} = \phi_j(x_i)$.

- Posterior marginal:

$$u(x) | \mathcal{D} \sim \text{Student-T} \left(\mathbf{U}(x)^\top (\mathbf{I} + \Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f}, \frac{1}{n} (\mathbf{f}^\top \Phi^\top (\mathbf{I} + \Phi \Phi^\top)^{-1} \Phi \mathbf{f}) \mathbf{U}(x)^\top (\mathbf{I} + \Phi^\top \Phi)^{-1} \mathbf{U}(x), n \right)$$

where $[\mathbf{U}(x)]_i = u_i(x)$ and u_i solves $\Delta u = \phi_i$ on \mathcal{X} and $u = 0$ on $\partial \mathcal{X}$. [N.B. Don't need to explicitly compute the ϕ_i if you have the Green's function of the PDE.]

Solve the ODE $\frac{du}{dx} = f(x)$, $u(0) = u_0$ on $x \in [0, T]$:

- Posterior mean coincides with a classical “collocation” method.
- Generalises to GPs with the kernel trick.
- Theoretical results (for a method based on GPs) in Cockayne *et al.*, 2016:
 - The posterior mean converges in $\|\cdot\|_\infty$ at a rate

$$O(h^{\alpha-\rho-\frac{d}{2}}).$$

- The posterior mass for a ball of radius ϵ centred on the true solution $u(\cdot)$ scales as

$$1 - O\left(\frac{h^{2\alpha-2\rho-d}}{\epsilon}\right).$$

- Posterior mean coincides with a classical “collocation” method.
- Generalises to GPs with the kernel trick.
- Theoretical results (for a method based on GPs) in Cockayne *et al.*, 2016:
 - The posterior mean converges in $\|\cdot\|_\infty$ at a rate

$$O(h^{\alpha-\rho-\frac{d}{2}}).$$

- The posterior mass for a ball of radius ϵ centred on the true solution $u(\cdot)$ scales as

$$1 - O\left(\frac{h^{2\alpha-2\rho-d}}{\epsilon}\right).$$

- Posterior mean coincides with a classical “collocation” method.
- Generalises to GPs with the kernel trick.
- Theoretical results (for a method based on GPs) in Cockayne *et al.*, 2016:
 - The posterior mean converges in $\|\cdot\|_\infty$ at a rate

$$O(h^{\alpha-\rho-\frac{d}{2}}).$$

- The posterior mass for a ball of radius ϵ centred on the true solution $u(\cdot)$ scales as

$$1 - O\left(\frac{h^{2\alpha-2\rho-d}}{\epsilon}\right).$$

- Posterior mean coincides with a classical “collocation” method.
- Generalises to GPs with the kernel trick.
- Theoretical results (for a method based on GPs) in Cockayne *et al.*, 2016:

- The posterior mean converges in $\|\cdot\|_\infty$ at a rate

$$O(h^{\alpha-\rho-\frac{d}{2}}).$$

- The posterior mass for a ball of radius ϵ centred on the true solution $u(\cdot)$ scales as

$$1 - O\left(\frac{h^{2\alpha-2\rho-d}}{\epsilon}\right).$$

- Posterior mean coincides with a classical “collocation” method.
- Generalises to GPs with the kernel trick.
- Theoretical results (for a method based on GPs) in Cockayne *et al.*, 2016:
 - The posterior mean converges in $\|\cdot\|_\infty$ at a rate

$$O(h^{\alpha-\rho-\frac{d}{2}}).$$

- The posterior mass for a ball of radius ϵ centred on the true solution $u(\cdot)$ scales as

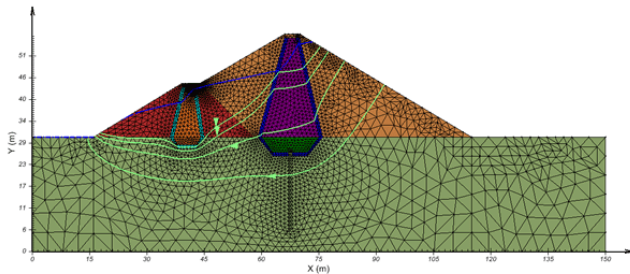
$$1 - O\left(\frac{h^{2\alpha-2\rho-d}}{\epsilon}\right).$$

- Posterior mean coincides with a classical “collocation” method.
- Generalises to GPs with the kernel trick.
- Theoretical results (for a method based on GPs) in Cockayne *et al.*, 2016:
 - The posterior mean converges in $\|\cdot\|_\infty$ at a rate

$$O(h^{\alpha-\rho-\frac{d}{2}}).$$

- The posterior mass for a ball of radius ϵ centred on the true solution $u(\cdot)$ scales as

$$1 - O\left(\frac{h^{2\alpha-2\rho-d}}{\epsilon}\right).$$



<http://www.svflux.com/subdomains/svflux.com/index.shtml>

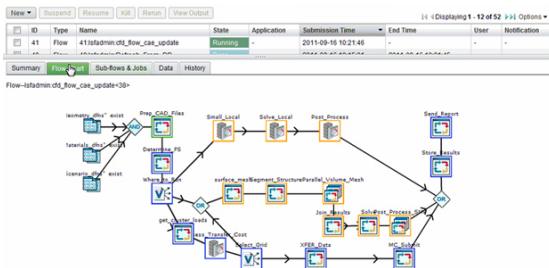
Summary

- General theory?
- Beyond linear and Gaussian assumptions?
- Experimental design?

Lots of work to do, but initial results in:

Cockayne J, Oates CJ, Sullivan T, Girolami M
Bayesian Probabilistic Numerical Methods
arXiv:1702.03673 (2017)

- Propagation of uncertainty through a computational workflow?
- Compatibility of multiple probabilistic numerical methods?



[Fig: IBM High Performance Computation]

The sophistication and scale of modern computer models creates an urgent need to better understand the propagation and accumulation of numerical error within arbitrary - often large - pipelines of computation, so that “numerical risk” to end-users can be controlled.