

Probabilistic Numerical Computation: A Role for Statistical Science in Numerical Analysis?

Mark Girolami

Sir Kirby Laing Chair of Civil Engineering
&
Royal Academy of Engineering Research Chair

Department of Engineering
University of Cambridge

The Alan Turing Institute
Programme Director



Lloyd's Register
Foundation



MLSS2019

Collaborators in Probabilistic Numerical Computation



Chris J. Oates
Newcastle



Jon Cockayne
Turing



F-X Briol
UCL



Tim Sullivan
F.U. Berlin



Philipp Hennig
MPI Tuebingen



Mike Osborne
Oxford



Dino Sejdinovic
Oxford



Andrew Stuart
Caltech

Probabilistic Numerical Computation

Consider the following school boy and girl differential equation

$$\frac{du}{dt} = \theta u, \quad u(t=0) = 1.$$

This is the simplest example model used to describe Malthusian population growth e.g. bacterial growth and radioactive decay. Simplest representation of compound interest in finance.

Every school boy and girl knows the solution:

$$u(t; \theta) = \exp(\theta t)$$

Despite the function $u(t; \theta)$ being implicitly defined it is a fully deterministic object.

Given the initial value then $u(t; \theta)$ at any point in the future is fully determined.

Consider the following school boy and girl differential equation

$$\frac{du}{dt} = \theta u, \quad u(t=0) = 1.$$

This is the simplest example model used to describe Malthusian population growth e.g. bacterial growth and radioactive decay. Simplest representation of compound interest in finance.

Every school boy and girl knows the solution:

$$u(t; \theta) = \exp(\theta t)$$

Despite the function $u(t; \theta)$ being implicitly defined it is a fully deterministic object.

Given the initial value then $u(t; \theta)$ at any point in the future is fully determined.

Consider the following school boy and girl differential equation

$$\frac{du}{dt} = \theta u, \quad u(t=0) = 1.$$

This is the simplest example model used to describe Malthusian population growth e.g. bacterial growth and radioactive decay. Simplest representation of compound interest in finance.

Every school boy and girl knows the solution:

$$u(t; \theta) = \exp(\theta t)$$

Despite the function $u(t; \theta)$ being implicitly defined it is a fully deterministic object.

Given the initial value then $u(t; \theta)$ at any point in the future is fully determined.

Consider the following school boy and girl differential equation

$$\frac{du}{dt} = \theta u, \quad u(t=0) = 1.$$

This is the simplest example model used to describe Malthusian population growth e.g. bacterial growth and radioactive decay. Simplest representation of compound interest in finance.

Every school boy and girl knows the solution:

$$u(t; \theta) = \exp(\theta t)$$

Despite the function $u(t; \theta)$ being implicitly defined it is a fully deterministic object.

Given the initial value then $u(t; \theta)$ at any point in the future is fully determined.

Consider the following school boy and girl differential equation

$$\frac{du}{dt} = \theta u, \quad u(t=0) = 1.$$

This is the simplest example model used to describe Malthusian population growth e.g. bacterial growth and radioactive decay. Simplest representation of compound interest in finance.

Every school boy and girl knows the solution:

$$u(t; \theta) = \exp(\theta t)$$

Despite the function $u(t; \theta)$ being implicitly defined it is a fully deterministic object.

Given the initial value then $u(t; \theta)$ at any point in the future is fully determined.

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

$$\text{Uncertainty } \theta \sim \mathcal{N}(\mu, \sigma) \quad \Rightarrow \quad u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

$$\text{Uncertainty } \theta \sim \mathcal{N}(\mu, \sigma) \quad \Rightarrow \quad u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

But wait.....

The rate parameter θ may be an empirically derived parameter.

This immediately introduces uncertainty into our deterministic world.

Our uncertainty in θ can be described using the calculus of probability

This uncertainty in θ propagates and induces uncertainty in $u(t; \theta)$

Uncertainty $\theta \sim \mathcal{N}(\mu, \sigma)$ \Rightarrow $u(t; \theta) \sim \log\mathcal{N}(\mu t, \sigma t)$

With uncertainty our deterministic object becomes a probabilistic object

Uncertainty can also enter by being unable to solve the differential equation analytically

What if the differential equation cannot be solved analytically ?

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a Known Unknown

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a **Known Unknown**

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a **Known Unknown**

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a **Known Unknown**

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a **Known Unknown**

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a **Known Unknown**

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Must resort to numerical methods to access approximations to the solution

The implicit function is unknown - we have a **Known Unknown**

For a general differential equation $\dot{u} = f(u; \theta)$ then the Euler method gives

$$U_{n+1} = U_n + hf(U_n; \theta)$$

For our school boy example with $U_0 = 1$ then

$$U_{n+1} = U_n + h\theta U_n = (1 + h\theta)^n$$

$$\theta \sim \log\mathcal{N}(\mu, \sigma)$$

$$\mathbb{E}\{U_n\} = \sum_{k=0}^{n-1} \binom{n-1}{k} h^k \mathbb{E}\{\theta^k\} \quad \mathbb{E}\{U_n^2\} = \sum_{k=0}^{2(n-1)} \binom{2(n-1)}{k} h^k \mathbb{E}\{\theta^k\}$$

The deterministic numerical procedure contributes further to uncertainty

The numerical procedure is now an inference procedure

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo
The numerical procedure is now an inference procedure
Defines a measure from which approximate solutions can be drawn

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo
The numerical procedure is now an inference procedure
Defines a measure from which approximate solutions can be drawn

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo
The numerical procedure is now an inference procedure
Defines a measure from which approximate solutions can be drawn

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo
The numerical procedure is now an inference procedure
Defines a measure from which approximate solutions can be drawn

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo

The numerical procedure is now an inference procedure

Defines a measure from which approximate solutions can be drawn

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo

The numerical procedure is now an inference procedure

Defines a measure from which approximate solutions can be drawn

Now then in the *unlikely* situation where we have complete knowledge of the initial value and value that θ takes we only have the **Known Unknown** to deal with.

Now everything is fully deterministic in the computation of our approximation. The evolution of the error is fully determined.

$$e_{n+1} = e_n + h[u(t_n) - U_n] + R$$

Nothing stochastic or random about this.

However we cannot compute the deterministic error or its equation of evolution - it is unknown

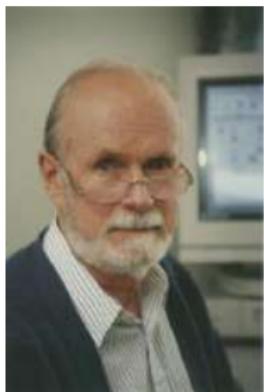
Subjectivist Probability - De Finetti, Ramsey, Jeffreys, Berger, Bernardo

The numerical procedure is now an inference procedure

Defines a measure from which approximate solutions can be drawn

Probabilistic Numerical Computation??

History of Probabilistic Numerical Methods

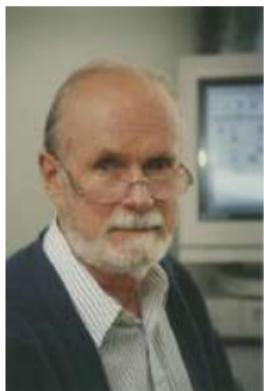


Tests of Probabilistic Models for Propagation of Roundoff Errors

T. E. HULL, University of Toronto; J. R. SWENSON, New York University (Ed: J. Traub) Communications of the ACM, 9(2):108 113, 1966.

In any prolonged computation it is generally assumed that the accumulated effect of roundoff errors is in some sense statistical. The purpose of this paper is to give precise descriptions of certain probabilistic models for roundoff error, and then to describe a series of experiments for testing the validity of these models. It is concluded that the models are in general very good. Discrepancies are both rare and mild. The test techniques can also be used to experiment with various types of special arithmetic.

History of Probabilistic Numerical Methods



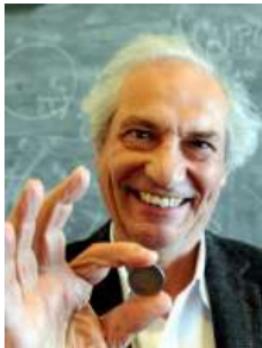
Tests of Probabilistic Models for Propagation of Roundoff Errors

T. E. HULL, University of Toronto; J. R. SWENSON, New York University (Ed: J. Traub) Communications of the ACM, 9(2):108 113, 1966.

In any prolonged computation it is generally assumed that the accumulated effect of roundoff errors is in some sense statistical. The purpose of this paper is to give precise descriptions of certain probabilistic models for roundoff error, and then to describe a series of experiments for testing the validity of these models. It is concluded that the models are in general very good. Discrepancies are both rare and mild. The test techniques can also be used to experiment with various types of special arithmetic.



Joseph Kadane
Kadane [1985]



Persi Diaconis
Diaconis [1988]



Tony O'Hagan
O'Hagan [1992]



John Skilling
Skilling [1991]

Question: “Is numerical computation a statistical inference problem?”

ROCKY MOUNTAIN
JOURNAL OF MATHEMATICS
Volume 2, Number 3, Summer 1972

GAUSSIAN MEASURE IN HILBERT SPACE AND APPLICATIONS IN NUMERICAL ANALYSIS

F. M. LARKIN

ABSTRACT. The numerical analyst is often called upon to estimate a function from a very limited knowledge of its properties (e.g. a finite number of ordinate values). This problem may be made well posed in a variety of ways, but an attractive approach is to regard the required function as a member of a linear space on which a probability measure is constructed, and then use established techniques of probability theory and statistics in order to infer properties of the function from the given information. This formulation agrees with established theory, for the problem of optimal linear approximation (using a Gaussian probability distribution), and also permits the estimation of nonlinear functionals, as well as extension to the case of "noisy" data.

History of Probabilistic Numerical Methods, F.M.Larkin



The numerical analyst is often called upon to estimate a function from a very limited knowledge of its properties (e.g. a finite number of ordinate values). This problem may be made well posed in a variety of ways, but an attractive approach is to regard the required function as a member of a linear space on which a probability measure is constructed, and then use established techniques of probability theory and statistics in order to infer properties of the function from the given information. This formulation agrees with established theory, for the problem of optimal linear approximation (using a Gaussian probability distribution), and also permits the estimation of nonlinear functionals, as well as extension to the case of “noisy” data.

History of Probabilistic Numerical Methods, F.M.Larkin



The numerical analyst is often called upon to estimate a function from a very limited knowledge of its properties (e.g. a finite number of ordinate values). This problem may be made well posed in a variety of ways, but an attractive approach is to regard the required function as a member of a linear space on which a probability measure is constructed, and then use established techniques of probability theory and statistics in order to infer properties of the function from the given information. This formulation agrees with established theory, for the problem of optimal linear approximation (using a Gaussian probability distribution), and also permits the estimation of nonlinear functionals, as well as extension to the case of “noisy” data.

What is Probabilistic Numerics?¹

Definition (Probabilistic Numerics)

Probabilistic Numerics **models the function uncertainty** and propagates a probabilistic description of this error through subsequent computations.

¹[Hennig, Osborne, Girolami., 2015]

What is Probabilistic Numerics?¹

Definition (Probabilistic Numerics)

Probabilistic Numerics **models the function uncertainty** and propagates a probabilistic description of this error through subsequent computations.

- Produces probability measures over all unknowns.
- Structure in residuals can be propagated through later computations.
- Analysis of variance to determine the computational sticking points.
- New perspective leads to design of new algorithms.
- Safeguards against unwarranted optimism for decision making

¹[Hennig, Osborne, Girolami., 2015]

Downloaded from <http://rspa.royalsocietypublishing.org/> on July 27, 2017

PROCEEDINGS A

rspa.royalsocietypublishing.org

Research



Cite this article: Hennig P, Osborne MA, Girolami M. 2015 Probabilistic numerics and uncertainty in computations. *Proc. R. Soc. A* **471**: 20150142.
<http://dx.doi.org/10.1098/rspa.2015.0142>.

Received: 2 March 2015

Accepted: 3 June 2015

Subject Areas:

statistics, computational mathematics,
artificial intelligence

Keywords:

numerical methods, probability, inference,
statistics

Probabilistic numerics and uncertainty in computations

Philipp Hennig¹, Michael A. Osborne²

and Mark Girolami³

¹Department of Empirical Inference, Max Planck Institute for Intelligent Systems, Tübingen, Germany

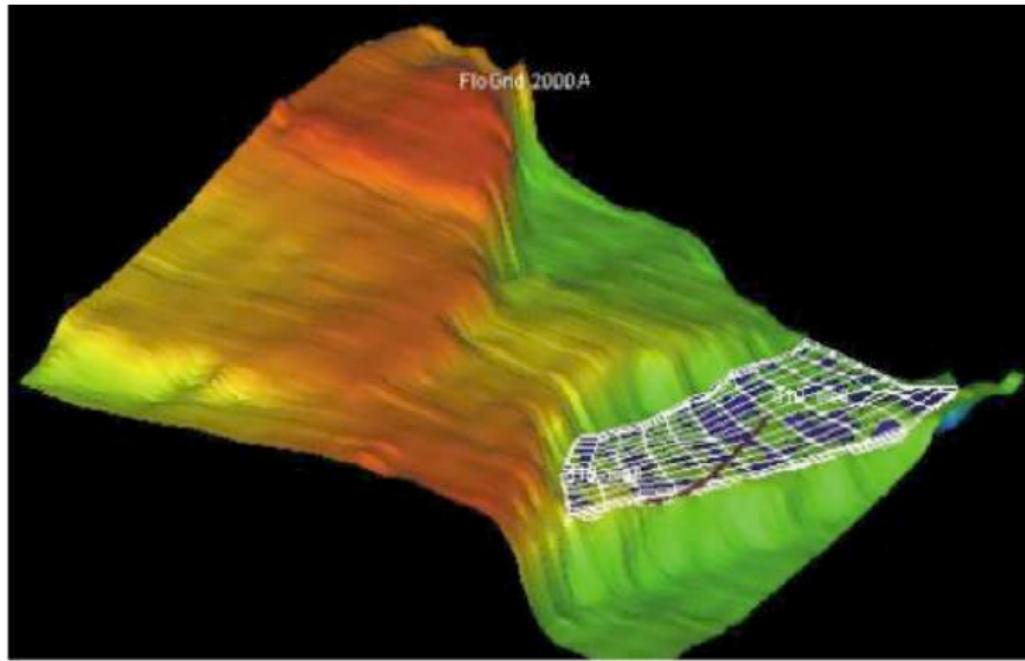
²Department of Engineering Science, University of Oxford, Oxford, UK

³Department of Statistics, University of Warwick, Warwick, UK

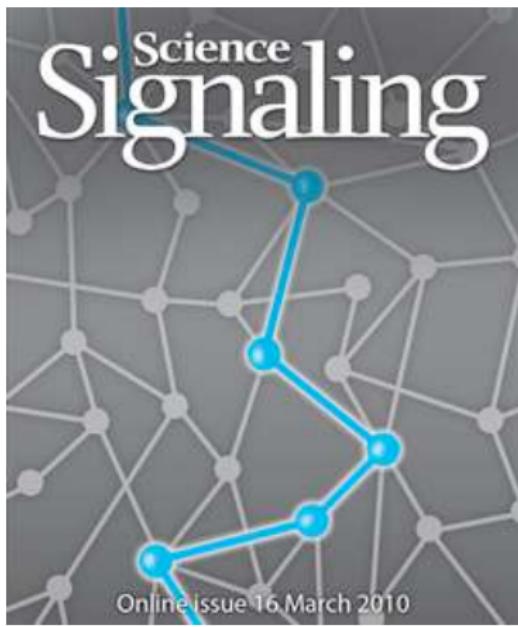
We deliver a call to arms for *probabilistic numerical methods*: algorithms for numerical tasks, including linear algebra, integration, optimization and solving differential equations, that return uncertainties in their calculations. Such uncertainties, arising from the loss of precision induced by numerical calculation with limited time or hardware, are important for much contemporary science and industry. Within applications such as climate science and astrophysics, the need to make decisions on the basis of computations with large and complex data have led to a renewed focus on the management of numerical uncertainty. We describe how several

Differential Equations

Motivation - Data Driven Engineering



Motivation - Data Informed Medical and Life Sciences



Motivation - Computational Social Science



Burglaries



Drugs



Traffic



Violence

Navier-Stokes Equations (2-d domain) [?]

Time-evolution of vorticity on a 2-d torus

angle of cross-section of ring (ρ)

angle of inner ring (θ)

PN for PDEs

A “widely used” linear PDE. Given g , κ , b find u

$$\begin{aligned}-\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D\end{aligned}$$

For general D , $u(\mathbf{x})$ this cannot be solved analytically.

The majority of PDE solvers produce an approximation like:

$$\hat{u}(\mathbf{x}) = \sum_{i=1}^N w_i \phi_i(\mathbf{x})$$

We want to quantify the error from finite N probabilistically.

PN for PDEs

A “widely used” linear PDE. Given g , κ , b find u

$$\begin{aligned}-\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D\end{aligned}$$

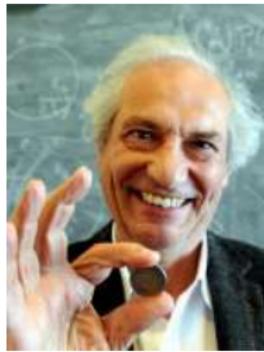
For general D , $u(\mathbf{x})$ this cannot be solved analytically.

The majority of PDE solvers produce an approximation like:

$$\hat{u}(\mathbf{x}) = \sum_{i=1}^N w_i \phi_i(\mathbf{x})$$

We want to quantify the error from finite N probabilistically.

History of Probabilistic Numerical Methods



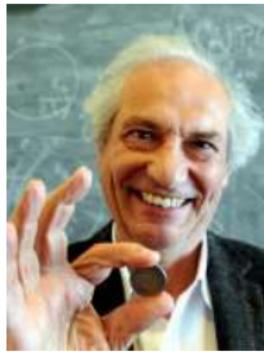
Bayesian Numerical Analysis

P. DIACONIS, Stanford University.

Statistical Decision Theory and Related Topics IV,
1, 163 175, 1988.

Seeing standard procedures emerge from the Bayesian approach may convince readers the argument isn't so crazy after all. The examples suggest the following program: Take standard numerical analysis procedures and see if they are Bayes (or admissible, or minimax). [...] The Bayesian approach yields more than the Bayes rule; it yields a posterior distribution. This can be used to give confidence sets as in Wahba (1983).

History of Probabilistic Numerical Methods



Bayesian Numerical Analysis

P. DIACONIS, Stanford University.

Statistical Decision Theory and Related Topics IV,
1, 163 175, 1988.

Seeing standard procedures emerge from the Bayesian approach may convince readers the argument isn't so crazy after all. The examples suggest the following program: Take standard numerical analysis procedures and see if they are Bayes (or admissible, or minimax). [...] The Bayesian approach yields more than the Bayes rule; it yields a posterior distribution. This can be used to give confidence sets as in Wahba (1983).

PN for PDEs

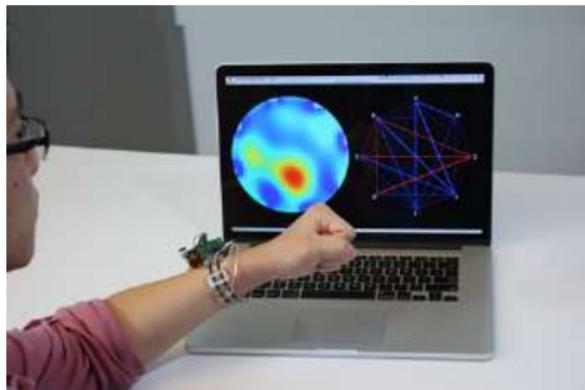
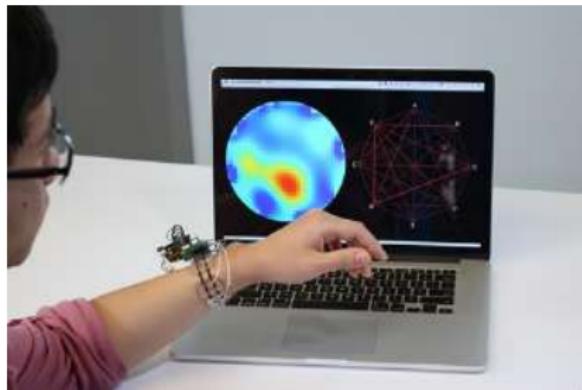
Inverse Problem: Given partial information of g, b, u find κ

$$\begin{aligned}-\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D\end{aligned}$$

PN for PDEs

Inverse Problem: Given partial information of g , b , u find κ

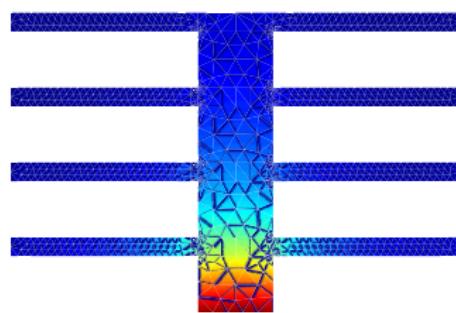
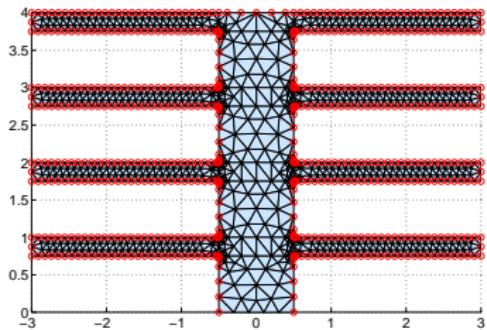
$$\begin{aligned}-\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D\end{aligned}$$



PN for PDEs

Inverse Problem: Given partial information of g, b, u find κ

$$\begin{aligned}-\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D\end{aligned}$$

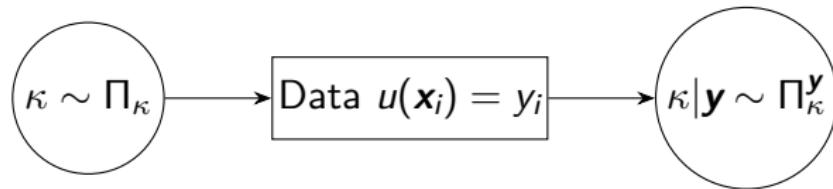


PN for PDEs

Inverse Problem: Given partial information of g, b, u find κ

$$\begin{aligned} -\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D \end{aligned}$$

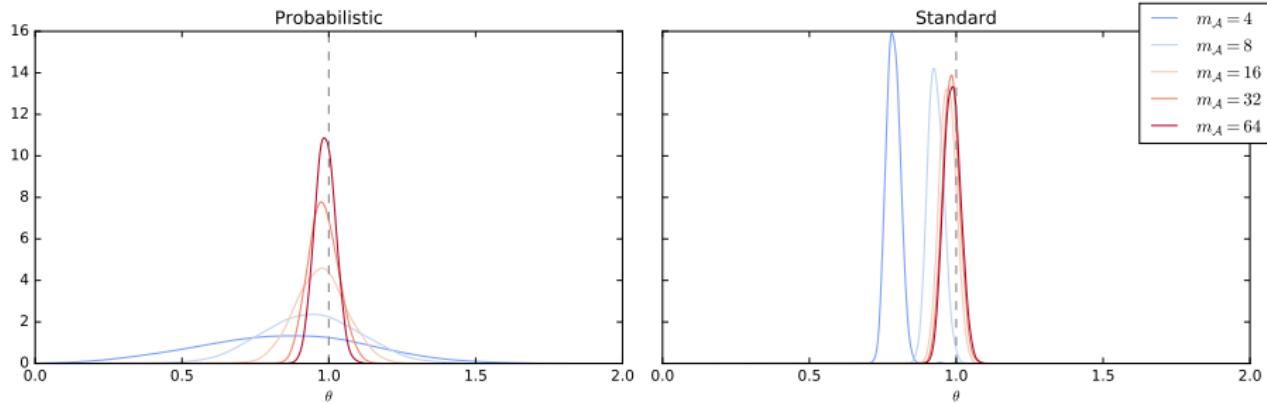
Bayesian Inverse Problem:



We want to account for an inaccurate forward solver in the inverse problem.

Why do this?

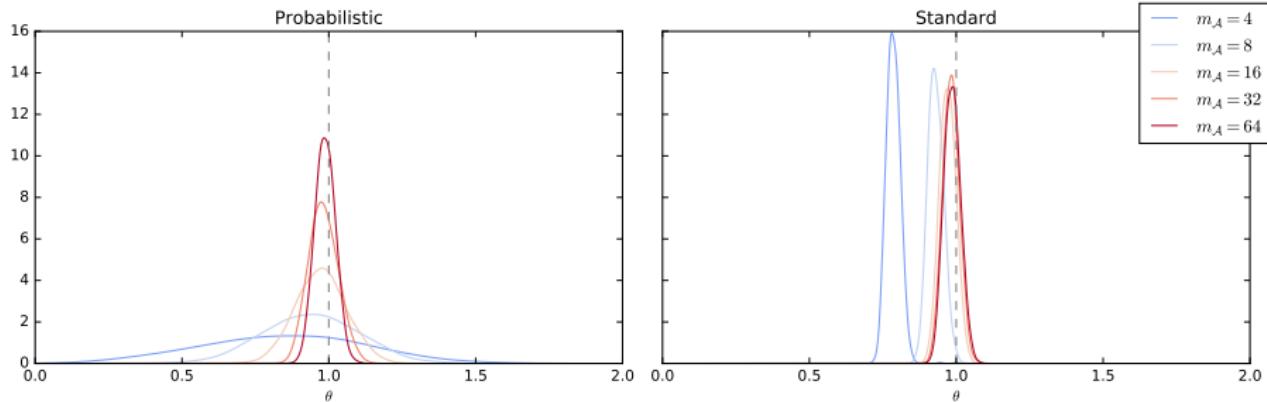
Using an inaccurate forward solver in an inverse problem can produce **biased** and **overconfident** posteriors.



Comparison of inverse problem posteriors produced using the Probabilistic Meshless Method (PMM) vs. symmetric collocation.

Why do this?

Using an inaccurate forward solver in an inverse problem can produce **biased** and **overconfident** posteriors.



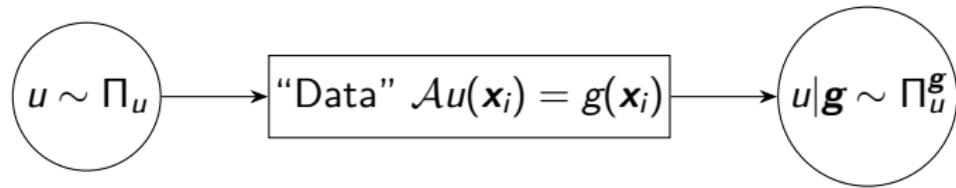
Comparison of inverse problem posteriors produced using the Probabilistic Meshless Method (PMM) vs. symmetric collocation.

Forward Problem

Abstract Formulation

$$\mathcal{A}u(\mathbf{x}) = g(\mathbf{x}) \quad \text{in } D$$

Forward inference procedure:



Posterior for the forward problem

Use a Gaussian Process prior $u \sim \Pi_u = \mathcal{GP}(0, k)$. Assuming linearity, the posterior Π_u^g is available in closed-form².

$$\Pi_u^g \sim \mathcal{GP}(m_1, \Sigma_1)$$

$$m_1(x) = \bar{\mathcal{A}}K(x, X) [\mathcal{A}\bar{\mathcal{A}}K(X, X)]^{-1} g$$

$$\Sigma_1(x, x') = k(x, x') - \bar{\mathcal{A}}K(x, X) [\mathcal{A}\bar{\mathcal{A}}K(X, X)]^{-1} \mathcal{A}K(X, x')$$

$\bar{\mathcal{A}}$ the adjoint of \mathcal{A}

Observation: The mean function is the same as in symmetric collocation!

²Larkin 1972, [Cockayne et al., 2016, Owhadi, 2014]

Posterior for the forward problem

Use a Gaussian Process prior $u \sim \Pi_u = \mathcal{GP}(0, k)$. Assuming linearity, the posterior Π_u^g is available in closed-form².

$$\Pi_u^g \sim \mathcal{GP}(m_1, \Sigma_1)$$

$$m_1(\mathbf{x}) = \bar{\mathcal{A}}K(\mathbf{x}, X) [\mathcal{A}\bar{\mathcal{A}}K(X, X)]^{-1} \mathbf{g}$$

$$\Sigma_1(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') - \bar{\mathcal{A}}K(\mathbf{x}, X) [\mathcal{A}\bar{\mathcal{A}}K(X, X)]^{-1} \mathcal{A}K(X, \mathbf{x}')$$

$\bar{\mathcal{A}}$ the adjoint of \mathcal{A}

Observation: The mean function is the same as in symmetric collocation!

²Larkin 1972, [Cockayne et al., 2016, Owhadi, 2014]

Theoretical Results

Theorem (Forward Contraction)

For a ball $B_\epsilon(u_0)$ of radius ϵ centered on the true solution u_0 of the PDE, we have

$$1 - \Pi_u^g[B_\epsilon(u_0)] = \mathcal{O}\left(\frac{h^{2\beta-2\rho-d}}{\epsilon}\right)$$

- h the fill distance
- β the smoothness of the prior
- $\rho < \beta - d/2$ the order of the PDE
- d the input dimension

Toy Example

$$\begin{aligned}-\nabla^2 u(x) &= g(x) & x \in (0, 1) \\ u(x) &= 0 & x = 0, 1\end{aligned}$$

To associate with the notation from before...

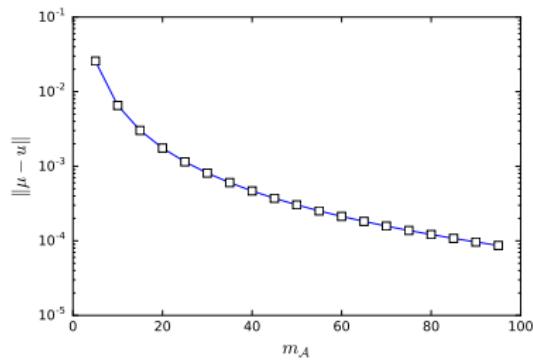
$$\Pi_u \sim \mathcal{GP}(0, k(x, y))$$

$$\mathcal{A} = -\frac{d^2}{dx^2} \quad \bar{\mathcal{A}} = -\frac{d^2}{dy^2}$$

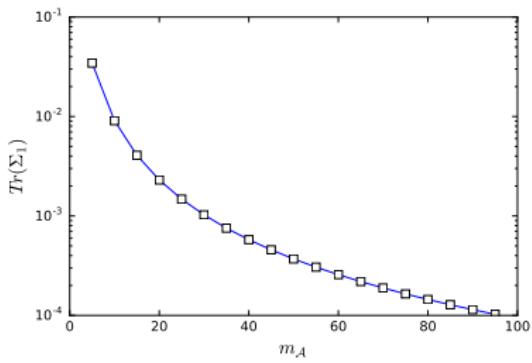
Forward problem: posterior samples

$$g(x) = \sin(2\pi x)$$

Forward problem: convergence



(a) Mean error from truth



(b) Trace of posterior covariance

Figure: Convergence

Inverse Problem

Recap

$$\begin{aligned}-\nabla \cdot (\kappa(\mathbf{x}) \nabla u(\mathbf{x})) &= g(\mathbf{x}) \quad \text{in } D \\ u(\mathbf{x}) &= b(\mathbf{x}) \quad \text{on } \partial D\end{aligned}$$

Now we need to incorporate the forward posterior measure Π_u^g into the posterior measure for the inverse problem, κ

Incorporation of Forward Measure

Assuming the data in the inverse problem is:

$$\begin{aligned}y_i &= u(\mathbf{x}_i) + \xi_i \quad i = 1, \dots, n \\ \boldsymbol{\xi} &\sim N(\mathbf{0}, \boldsymbol{\Gamma})\end{aligned}$$

implies the **standard** likelihood:

$$p(\mathbf{y}|\kappa, \mathbf{u}) \sim N(\mathbf{y}; \mathbf{u}, \boldsymbol{\Gamma})$$

But we don't know \mathbf{u}

Marginalise the forward posterior Π_u^g to obtain a “PN” likelihood:

$$\begin{aligned}p_{\text{PN}}(\mathbf{y}|\kappa) &\propto \int p(\mathbf{y}|\kappa, \mathbf{u}) d\Pi_u^g \\ &\sim N(\mathbf{y}; \mathbf{m}_1, \boldsymbol{\Gamma} + \boldsymbol{\Sigma}_1)\end{aligned}$$

Incorporation of Forward Measure

Assuming the data in the inverse problem is:

$$\begin{aligned}y_i &= u(\mathbf{x}_i) + \xi_i \quad i = 1, \dots, n \\ \boldsymbol{\xi} &\sim N(\mathbf{0}, \boldsymbol{\Gamma})\end{aligned}$$

implies the **standard** likelihood:

$$p(\mathbf{y}|\kappa, \mathbf{u}) \sim N(\mathbf{y}; \mathbf{u}, \boldsymbol{\Gamma})$$

But we don't know \mathbf{u}

Marginalise the forward posterior Π_u^g to obtain a “**PN**” likelihood:

$$\begin{aligned}p_{\text{PN}}(\mathbf{y}|\kappa) &\propto \int p(\mathbf{y}|\kappa, \mathbf{u}) d\Pi_u^g \\ &\sim N(\mathbf{y}; \mathbf{m}_1, \boldsymbol{\Gamma} + \boldsymbol{\Sigma}_1)\end{aligned}$$

Inverse Contraction

Denote by Π_κ^y the posterior for κ from likelihood p , and by $\Pi_{\kappa,PN}^y$ the posterior for κ from likelihood p_{PN} .

Theorem (Inverse Contraction)

Assume $\Pi_\kappa^y \rightarrow \delta(\kappa_0)$ as $n \rightarrow \infty$.

Then $\Pi_{\kappa,PN}^y \rightarrow \delta(\kappa_0)$ provided

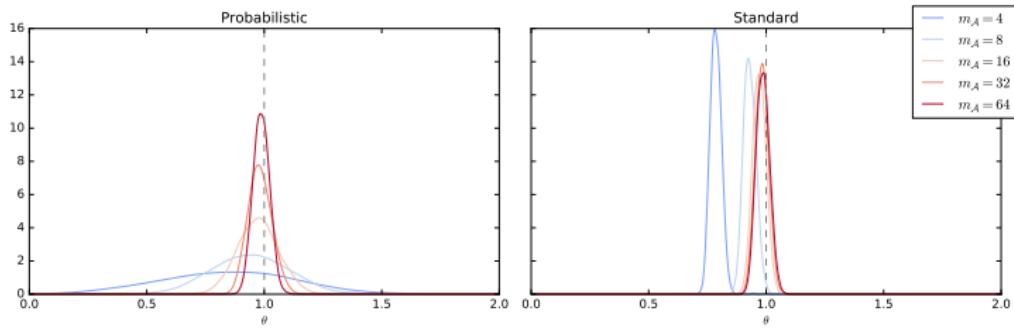
$$h = o(n^{-1/(\beta-\rho-d/2)})$$

Back to the Toy Example

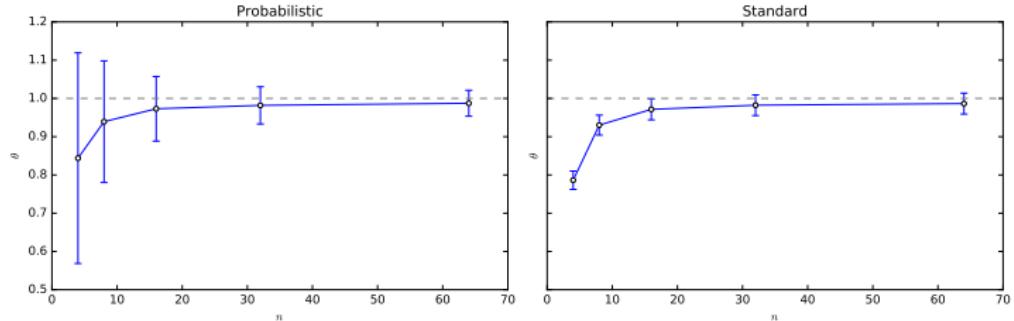
$$\begin{aligned}-\nabla \cdot (\kappa \nabla u(x)) &= \sin(2\pi x) & x \in (0, 1) \\ u(x) &= 0 & x = 0, 1\end{aligned}$$

Infer $\kappa \in \mathbb{R}^+$; data generated for $\kappa = 1$ at $x = 0.25, 0.75$.
Corrupted with independent Gaussian noise $\xi \sim N(0, 0.01^2)$

Posteriors for κ



(a) Posterior Distributions for different numbers of design points.



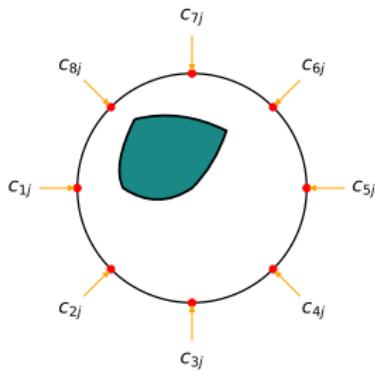
(b) Convergence of posterior distributions with number of design points.

Electrical Impedance Tomography

A medical imaging technique. Goal: reconstruct **interior conductivity field** of a patient, to detect tumors.

Electrical Impedance Tomography

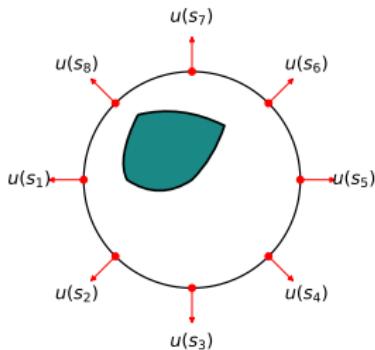
A medical imaging technique. Goal: reconstruct **interior conductivity field** of a patient, to detect tumors.



Many patterns of current $c_{ij}, j = 1, \dots, N_c$ injected through **boundary electrodes** $t_i^{\text{obs}}, i = 1, \dots, N_s$

Electrical Impedance Tomography

A medical imaging technique. Goal: reconstruct **interior conductivity field** of a patient, to detect tumors.



Resulting voltage measured: $y_i = x(t_i^{\text{obs}}) - x(t_{\text{ref}}) + \epsilon_i$

Electrical Impedance Tomography

A medical imaging technique. Goal: reconstruct **interior conductivity field** of a patient, to detect tumors.

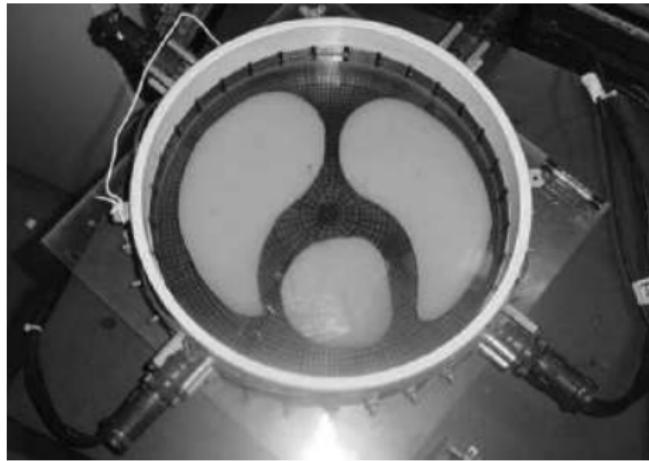
Governing equations are essentially Darcy's law:

$$-\nabla \cdot (\theta(t) \nabla x(t)) = 0 \quad t \in D$$

$$\theta(t_i^{\text{obs}}) \frac{\partial x}{\partial n}(t_i^{\text{obs}}) = c_{ij} \quad i = 1, \dots, N_S$$

Experimental Set-Up

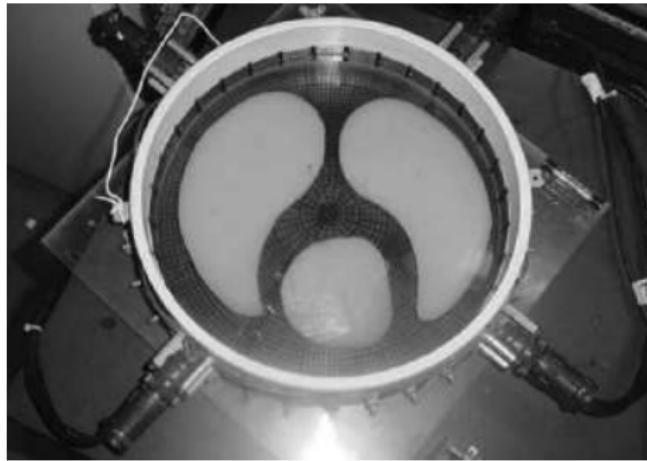
Experiments due to Isaacson 2004.



- Tank filled with saline.
- Three targets:
 - “Heart shaped”: higher conductivity.
 - “Lung shaped”: lower conductivity.
- 32 equally spaced electrodes.
- Simultaneously stimulated for 31 different stimulation patterns.

Experimental Set-Up

Experiments due to Isaacson 2004.



- Tank filled with saline.
- Three targets:
 - “Heart shaped”: higher conductivity.
 - “Lung shaped”: lower conductivity.
- 32 equally spaced electrodes.
- Simultaneously stimulated for 31 different stimulation patterns.

A Hard Problem...

- High dimensional (992) observations.
- Observations are only of the boundary - weak information.
- Target $\theta(\cdot)$ is infinite-dimensional.
- The “ideal” likelihood $\mathcal{L}(\theta; \mathbf{y})$ requires exact solution of the PDE.

Posteriors obtained using the PN likelihood

$$\begin{aligned}\mathcal{L}_n(\theta; \mathbf{y}) &\propto \int p(\mathbf{y}|\theta, \mathbf{x}) dP_{\mathbf{x}|a} \\ \implies \mathbf{y}|\theta &\sim N(\mathbf{m}_1, \Gamma + \Sigma_1).\end{aligned}$$

Focus on varying the number n of points in $T = \{t_i\}_{i=1}^n$ that are used.

Computation facilitated with Markov chain Monte Carlo, based on the preconditioned Crank-Nicholson proposal.

A Hard Problem...

- High dimensional (992) observations.
- Observations are only of the boundary - weak information.
- Target $\theta(\cdot)$ is infinite-dimensional.
- The “ideal” likelihood $\mathcal{L}(\theta; \mathbf{y})$ requires exact solution of the PDE.

Posteriors obtained using the PN likelihood

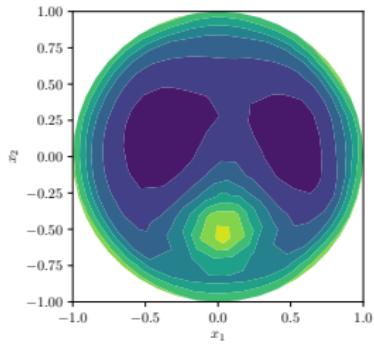
$$\begin{aligned}\mathcal{L}_n(\theta; \mathbf{y}) &\propto \int p(\mathbf{y}|\theta, x) dP_{x|a} \\ \implies \mathbf{y}|\theta &\sim N(\mathbf{m}_1, \Gamma + \Sigma_1).\end{aligned}$$

Focus on varying the number n of points in $T = \{t_i\}_{i=1}^n$ that are used.

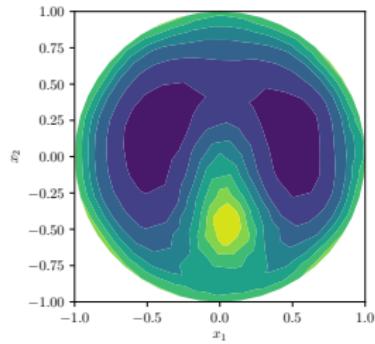
Computation facilitated with Markov chain Monte Carlo, based on the preconditioned Crank-Nicholson proposal.

Recovered Fields

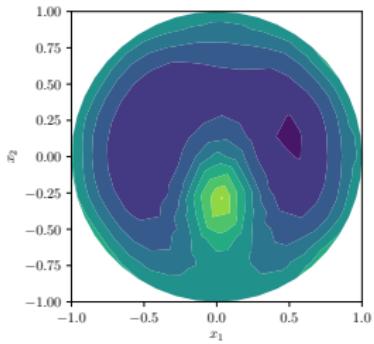
Posterior means $m(t) = \mathbb{E}_{\mathbf{y}}[\theta(t)]$:



(a) $n = 96$



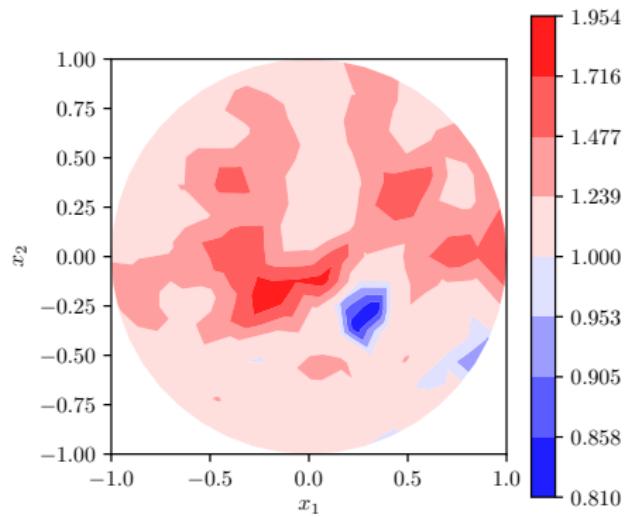
(b) $n = 127$



(c) $n = 165$

Variance Analysis

Ratio of (pointwise) posterior variance $v(t) = \mathbb{V}_y[\theta(t)]$ computed from the PN posterior based on \mathcal{L}_n and the “standard” posterior based on $\hat{\mathcal{L}}_N$ with $n = N = 96$:

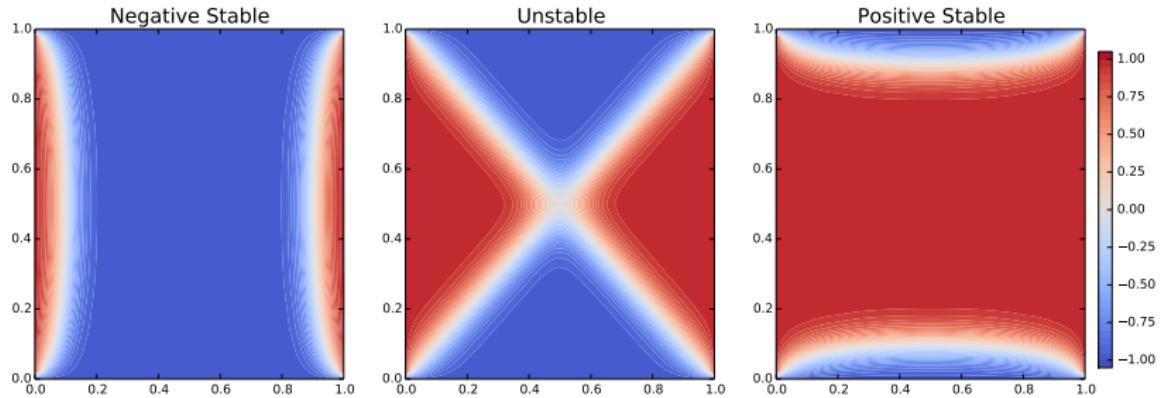


Allen–Cahn

A prototypical nonlinear model.

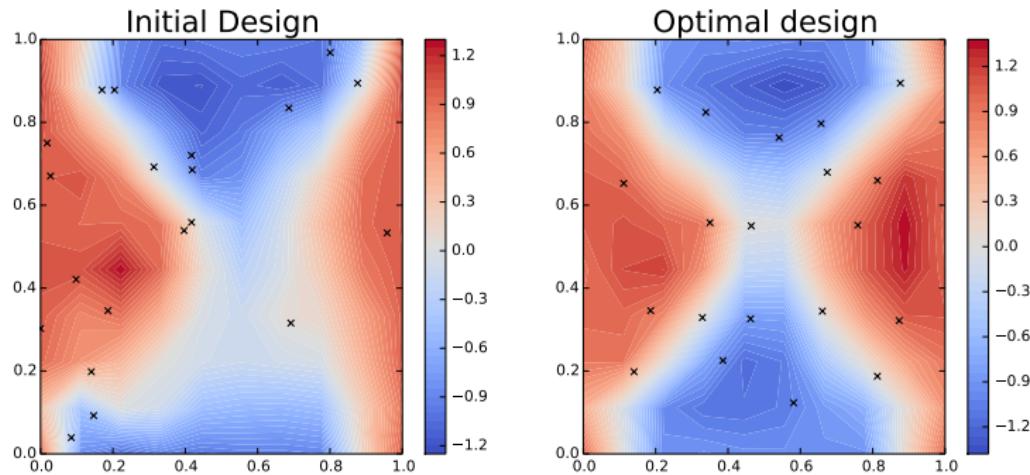
$$\begin{aligned}
 -\theta \nabla^2 u(\mathbf{x}) + \theta^{-1} (u(\mathbf{x})^3 - u(\mathbf{x})) &= 0 & \mathbf{x} \in (0, 1)^2 \\
 u(\mathbf{x}) &= 1 & x_1 \in \{0, 1\}; 0 < x_2 < 1 \\
 u(\mathbf{x}) &= -1 & x_2 \in \{0, 1\}; 0 < x_1 < 1
 \end{aligned}$$

Goal: infer θ

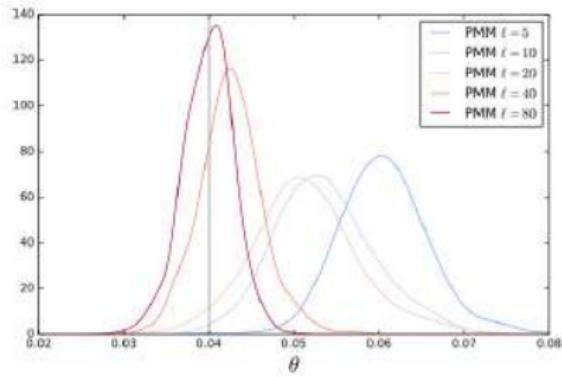


Allen–Cahn: Forward Solutions

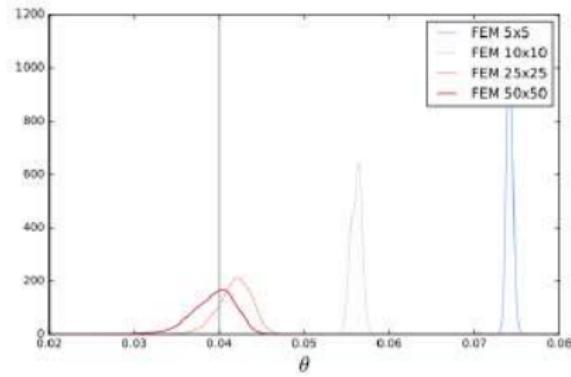
Nonlinear PDE - non-GP posterior sampling schemes required, see [Cockayne et al., 2016].



Allen–Cahn: Inverse Problem



(a) Probabilistic.



(b) Standard

Parabolic Systems

Parabolic systems are of the form

$$\frac{\partial u}{\partial t} = \mathcal{A}u + g$$

Solution in the spirit of Probabilistic ODE solvers³; discretise the time operator and model the error **probabilistically**:

$$u(\mathbf{x}, t_{n+1}) = u(\mathbf{x}, t_n) + [\mathcal{A}u(\mathbf{x}, t_{n+1}) + g(\mathbf{x}, t_{n+1})] \Delta t + \xi_n$$

³[????]

A Statistical Model of Burglary [?]

$A(\mathbf{x}, t)$ the attractiveness of the domain to criminals. $\rho(\mathbf{x}, t)$ the concentration of criminals. These evolve according to

$$\frac{\partial A}{\partial t} = \frac{\eta D}{z} \nabla^2 A - \omega(A - A^0) + \epsilon D \rho A$$

$$\frac{\partial \rho}{\partial t} = \frac{D}{z} \nabla \cdot \left[\nabla \rho - \frac{2\rho}{A} \nabla A \right] - \rho A + \gamma$$

with initial conditions $A(\mathbf{x}, 0)$, $\rho(\mathbf{x}, 0)$ known.

This is a nonlinear system thanks to terms in ρA ; must sample from the posterior distribution.

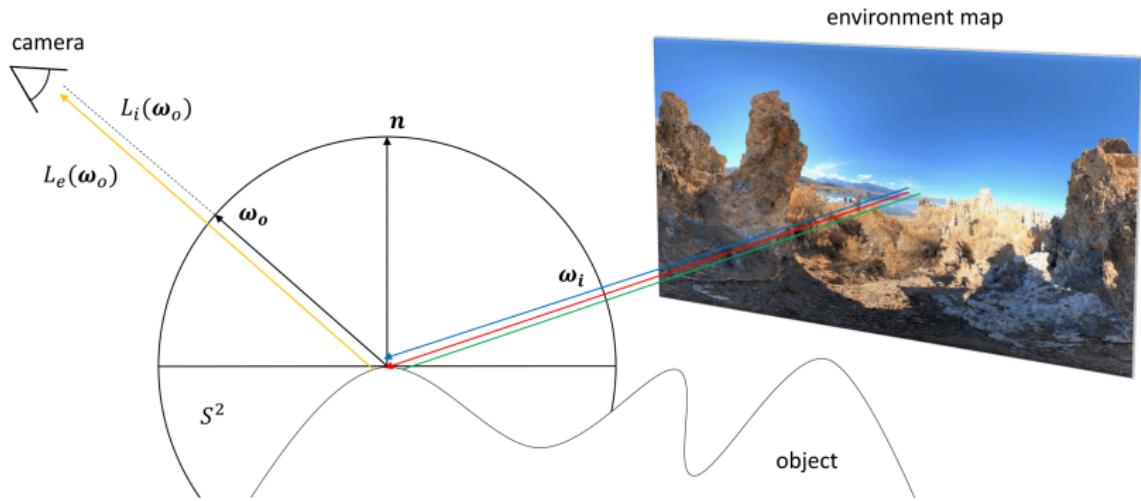


Example: A Statistical Model of Burglary (1-d domain)

Time-evolution of the attractiveness field $A(x)$ and criminal density $\rho(x)$ on a 1-d interval (street)

Integration

Illustrative Application - Integral over Manifold



Integrals Over Manifolds

$$L_o(\omega_o) = L_e(\omega_o) + \int_{\mathbb{S}^2} L_i(\omega_i) \rho(\omega_i, \omega_o) [\omega_i \cdot \mathbf{n}]_+ d\pi(\omega_i)$$

- $L_o(\omega_o)$ = outgoing radiance
- $L_e(\omega_o)$ = amount of light emitted by the object itself
- $L_i(\omega_i)$ = amount of light reaching object from direction ω_i
- ρ = bidirectional reflectance distribution function
- π = uniform distribution on \mathbb{S}^2

To be computed

- for each pixel, and
- for each RGB channel.

Integrals Over Manifolds

$$L_o(\omega_o) = L_e(\omega_o) + \int_{\mathbb{S}^2} L_i(\omega_i) \rho(\omega_i, \omega_o) [\omega_i \cdot \mathbf{n}]_+ d\pi(\omega_i)$$

- $L_o(\omega_o)$ = outgoing radiance
- $L_e(\omega_o)$ = amount of light emitted by the object itself
- $L_i(\omega_i)$ = amount of light reaching object from direction ω_i
- ρ = bidirectional reflectance distribution function
- π = uniform distribution on \mathbb{S}^2

To be computed

- for each pixel, and
- for each RGB channel.

The Problem

Let f be continuous and square-integrable, Π be a probability measure and $\mathcal{X} \subseteq \mathbb{R}^d$. We want to compute (numerically):

$$\Pi[f] = \int_{\mathcal{X}} f d\Pi \approx \sum_{i=1}^n w_i f(\mathbf{x}_i) = \hat{\Pi}[f] \quad (1)$$

High numerical uncertainty when f is expensive or n is small!

The Problem

Let f be continuous and square-integrable, Π be a probability measure and $\mathcal{X} \subseteq \mathbb{R}^d$. We want to compute (numerically):

$$\Pi[f] = \int_{\mathcal{X}} f d\Pi \approx \sum_{i=1}^n w_i f(\mathbf{x}_i) = \hat{\Pi}[f] \quad (1)$$

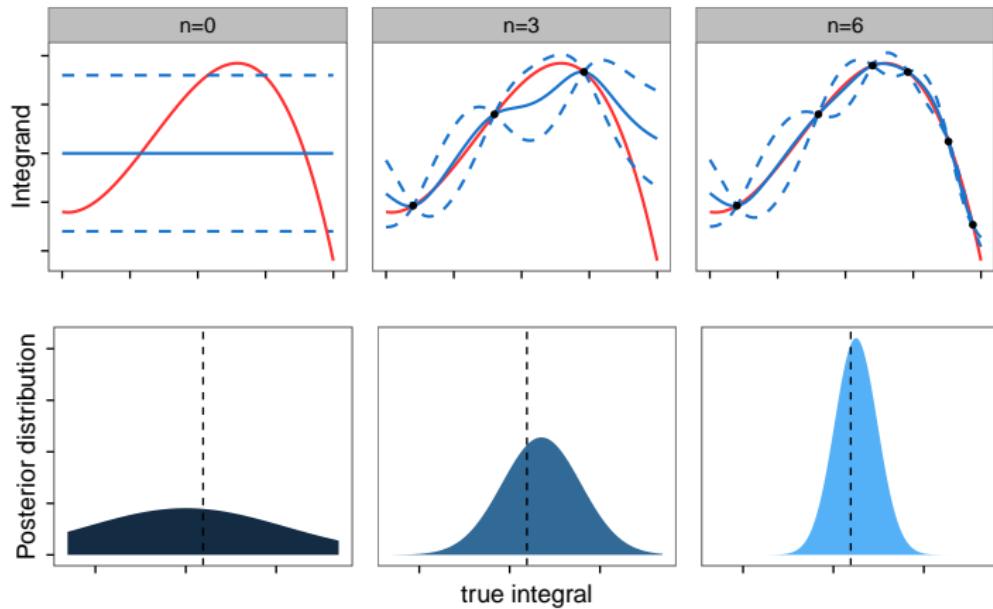
High numerical uncertainty when f is expensive or n is small!

Probabilistic Numerics Solution: Bayesian Quadrature⁴ (BQ) makes use of *prior information* about f to guide our choice of $\{\mathbf{x}_i, w_i\}_{i=1}^n$ (through a choice of function space/RKHS).

Measure on Integral push-forward of measure on function.

⁴[O'Hagan, 1991, Rasmussen and Ghahramani, 2002, Briol et al., 2015a,b]

Sketch of Bayesian Quadrature



$$\mathbb{E}_n[\Pi[f]] = \hat{\Pi}_{\text{BQ}} = \Pi[k(\cdot, \mathbf{X})] \mathbf{K}^{-1} \mathbf{f}$$

$$\mathbb{V}_n[\Pi[f]] = \Pi \Pi[k(\cdot, \cdot)] - \Pi[k(\cdot, \mathbf{X})] \mathbf{K}^{-1} \Pi[k(\mathbf{X}, \cdot)].$$

Theory for Bayesian Quadrature

We consider Sobolev spaces, which are RKHS \mathcal{H}^α of varying levels of smoothness α , which consist of functions in L_2 with associated inner product:

$$\langle f, g \rangle_{H^\alpha} := \sum_{m=0}^{\alpha} \left\langle \frac{d^m f}{dx^m}, \frac{d^m g}{dx^m} \right\rangle_{L_2}$$

and finite norm $\|f\|_{H^\alpha(\Pi)} := \langle f, f \rangle_{H^\alpha}^{1/2}$.

We study the performance of the method in terms of worst-case error:

$$e(\hat{\Pi}; \Pi, \mathcal{H}) = \sup_{f: \|f\|_{\mathcal{H}} \leq 1} |\Pi[f] - \hat{\Pi}[f]|.$$

Theory for Bayesian Quadrature

Theorem (BQ in Sobolev spaces [Briol et al., 2015b])

Let $\mathcal{X} = [0, 1]^d$, Π be $\text{Unif}(\mathcal{X})$ and Π_{BQ} be a BQ rule whose states $\{\mathbf{x}_i\}_{i=1}^n$ i.i.d. $\sim \pi$. Then, whenever $\alpha > d/2$, we have:

$$e(\hat{\Pi}_{BQ}; \Pi, \mathcal{H}) = \mathcal{O}_P(n^{-\alpha/d+\epsilon})$$

where $\epsilon > 0$ can be arbitrarily small. Furthermore, let $I_D = [\Pi[f] - D, \Pi[f] + D]$. Then:

$$\mathbb{P}_n[I_D^c] = o_P(\exp(-Cn^{2\alpha/d-\epsilon}))$$

Integrals Over Manifolds

Idea: Construct a RKHS of functions $x : \mathbb{S}^2 \rightarrow \mathbb{R}$.

One such kernel, that leads to a Sobolev space of smoothness $\frac{3}{2}$ on \mathbb{S}^2 :

$$k(t, t') = \frac{8}{3} - \|t - t'\|_2 \text{ for all } t, t' \in \mathbb{S}^2.$$

Integrals Over Manifolds

Idea: Construct a RKHS of functions $x : \mathbb{S}^2 \rightarrow \mathbb{R}$.

One such kernel, that leads to a Sobolev space of smoothness $\frac{3}{2}$ on \mathbb{S}^2 :

$$k(t, t') = \frac{8}{3} - \|t - t'\|_2 \text{ for all } t, t' \in \mathbb{S}^2.$$

Integrals Over Manifolds

Idea: Construct a RKHS of functions $x : \mathbb{S}^2 \rightarrow \mathbb{R}$.

One such kernel, that leads to a Sobolev space of smoothness $\frac{3}{2}$ on \mathbb{S}^2 :

$$k(t, t') = \frac{8}{3} - \|t - t'\|_2 \text{ for all } t, t' \in \mathbb{S}^2.$$

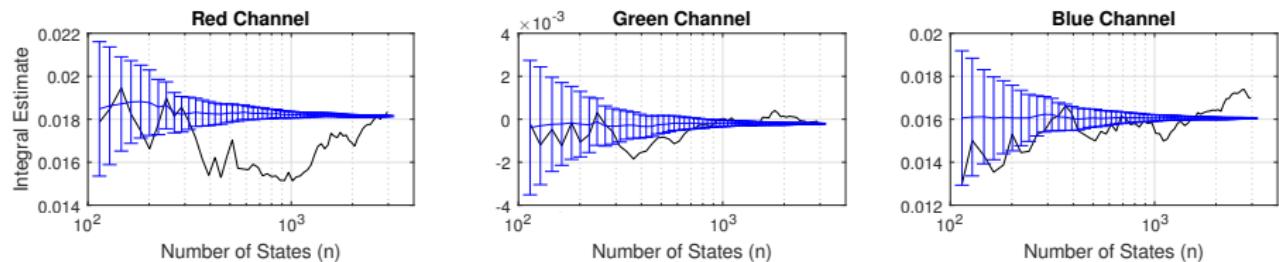
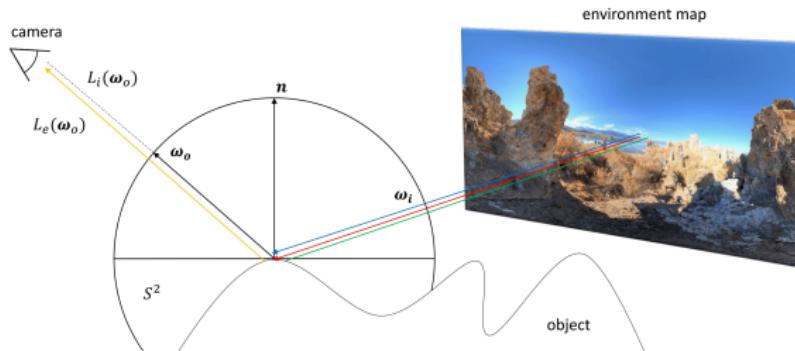
For a certain *spherical t-design* $\{t_i\}_{i=1}^n$, a convergence rate of $\epsilon_{\text{WCE}}(M) = O(n^{-\frac{3}{4}})$ is achieved by the method $M = (A, b)$ where b is the Bayesian Quadrature posterior mean - and this is worst-case optimal:



Integrals Over Manifolds

Full uncertainty quantification for integrals on manifolds:

Prob Integration in Comp Graphics [Briol et al., 2015b]



We provide rates of $\mathcal{O}_P(n^{-\frac{3}{4}})$ which is optimal for $\mathcal{H}^{\frac{3}{2}}(\mathbb{S}^2)!$

Integration with Intractable Densities

Intractable Densities and Stein's identity

What if $\pi(x)$ is only known up to a constant?

$$\pi(x) = \frac{\pi_c(x)}{c} \propto \pi_c(x)$$

In those cases $\Pi[k(\cdot, x)]$ is not available in closed form!

We can build an RKHS via kernel which takes into account information about π , but does not require us to know c^5 .

Let $\phi(x)$ be twice differentiable, we can use the Stein transformation

$$\mathcal{L}\phi(x) := \frac{\nabla[\phi(x)\pi(x)]}{\pi(x)}.$$

Obtain an RKHS taking account of smoothness of both integrand and density of distribution - **Control Functionals**

⁵Oates et al. [2017], Oates and Girolami [2016]



Journal of the Royal Statistical Society
Statistical Methodology
Series B

J. R. Statist. Soc. B (2017)
79, Part 3, pp. 695–718

Control functionals for Monte Carlo integration

Chris J. Oates,

University of Technology Sydney, Australia

Mark Girolami

University of Warwick, Coventry, and Alan Turing Institute, London, UK

and Nicolas Chopin

Centre de Recherche en Economie et Statistique and Ecole Nationale de la Statistique et de l'Administration Economique, Paris, France

[Received October 2014. Final revision February 2016]

Summary. A non-parametric extension of control variates is presented. These leverage gradient information on the sampling density to achieve substantial variance reduction. It is not required that the sampling density be normalized. The novel contribution of this work is based on two important insights: a trade-off between random sampling and deterministic approximation and a new gradient-based function space derived from Stein's identity. Unlike classical control variates, our estimators improve rates of convergence, often requiring orders of magnitude fewer simulations to achieve a fixed level of precision. Theoretical and empirical results are presented, the latter focusing on integration problems arising in hierarchical models and models based on non-linear ordinary differential equations.

Keywords: Control variates; Non-parametrics; Reproducing kernel; Stein's identity; Variance reduction

Theory for Control Functionals⁶

Theorem (Consistency of Control Functionals)

Suppose $\{x_i\}_{i=1}^n$ arise from a Markov chain that targets a density $\pi(x)$.

- Assume \mathcal{X} is bounded.
- Assume $\pi(x)$ is bounded away from 0 on \mathcal{X} .
- Assume $\pi \in C^{2a+1}(\mathcal{X})$ & $k \in C^{2b+2}(\mathcal{X} \times \mathcal{X})$.
- Assume k satisfies “certain boundary conditions”.
- Assume the Markov chain is uniformly ergodic.

Then, for $f \in \mathcal{H}_k$, there exists $h > 0$ such that

$$1_{h_n < h} (\Pi[f] - \hat{\Pi}[f])^2 = \mathcal{O}_P(n^{-1 - \frac{2(a \wedge b)}{d} + \epsilon}),$$

where $\epsilon > 0$ hides logarithmic factors.

⁶[Oates et al., 2016b]

Example: Computation of Marginal Likelihood

Consider computing the marginal likelihood for a non-linear ODE model

$$\frac{d^2x}{dt^2} - \theta(1 - x^2) \frac{dx}{dt} + x = 0$$

where $\theta \in \mathbb{R}$ is an unknown parameter indicating the non-linearity and the strength of damping.

Observations \mathbf{y} are made once every time unit, up to 10 units, and Gaussian measurement noise of standard deviation $\sigma = 0.1$ was added. A log-normal prior was placed on θ such that $\log(\theta) \sim N(0, 0.25)$.

Goal: Compute $p(\mathbf{y})$.

Example: Computation of Marginal Likelihood

Thermodynamic integration is based on the identity

$$\log p(\mathbf{y}) = \int_0^1 \mathbb{E}_{\theta|\mathbf{y},t} [\log p(\mathbf{y}|\theta)] dt.$$

where the “power posterior” for parameters θ given data \mathbf{y} is defined as $p(\theta|\mathbf{y}, t) \propto p(\mathbf{y}|\theta)^t p(\theta)$.

In TI, this integral is evaluated numerically over a discrete temperature ladder $0 = t_0 < t_1 < \dots < t_m = 1$. e.g.

$$\widehat{\log p(\mathbf{y})} := \sum_{i=0}^{m-1} \frac{(t_{i+1} - t_i)}{2} \{ \widehat{\mathbb{E}_{\theta|\mathbf{y},t_i}} [\log p(\mathbf{y}|\theta)] + \widehat{\mathbb{E}_{\theta|\mathbf{y},t_{i+1}}} [\log p(\mathbf{y}|\theta)] \}.$$

i.e. lots of integrals!

Example: Computation of Marginal Likelihood

Thermodynamic integration is based on the identity

$$\log p(\mathbf{y}) = \int_0^1 \mathbb{E}_{\theta|\mathbf{y},t} [\log p(\mathbf{y}|\theta)] dt.$$

where the “power posterior” for parameters θ given data \mathbf{y} is defined as $p(\theta|\mathbf{y}, t) \propto p(\mathbf{y}|\theta)^t p(\theta)$.

In TI, this integral is evaluated numerically over a discrete temperature ladder $0 = t_0 < t_1 < \dots < t_m = 1$. e.g.

$$\widehat{\log p(\mathbf{y})} := \sum_{i=0}^{m-1} \frac{(t_{i+1} - t_i)}{2} \{ \widehat{\mathbb{E}_{\theta|\mathbf{y},t_i}} [\log p(\mathbf{y}|\theta)] + \widehat{\mathbb{E}_{\theta|\mathbf{y},t_{i+1}}} [\log p(\mathbf{y}|\theta)] \}.$$

i.e. lots of integrals!

Example: Computation of Marginal Likelihood

Thermodynamic integration is based on the identity

$$\log p(\mathbf{y}) = \int_0^1 \mathbb{E}_{\theta|\mathbf{y},t} [\log p(\mathbf{y}|\theta)] dt.$$

where the “power posterior” for parameters θ given data \mathbf{y} is defined as $p(\theta|\mathbf{y}, t) \propto p(\mathbf{y}|\theta)^t p(\theta)$.

In TI, this integral is evaluated numerically over a discrete temperature ladder $0 = t_0 < t_1 < \dots < t_m = 1$. e.g.

$$\widehat{\log p(\mathbf{y})} := \sum_{i=0}^{m-1} \frac{(t_{i+1} - t_i)}{2} \{ \widehat{\mathbb{E}_{\theta|\mathbf{y},t_i}} [\log p(\mathbf{y}|\theta)] + \widehat{\mathbb{E}_{\theta|\mathbf{y},t_{i+1}}} [\log p(\mathbf{y}|\theta)] \}.$$

i.e. lots of integrals!

Example: Computation of Marginal Likelihood

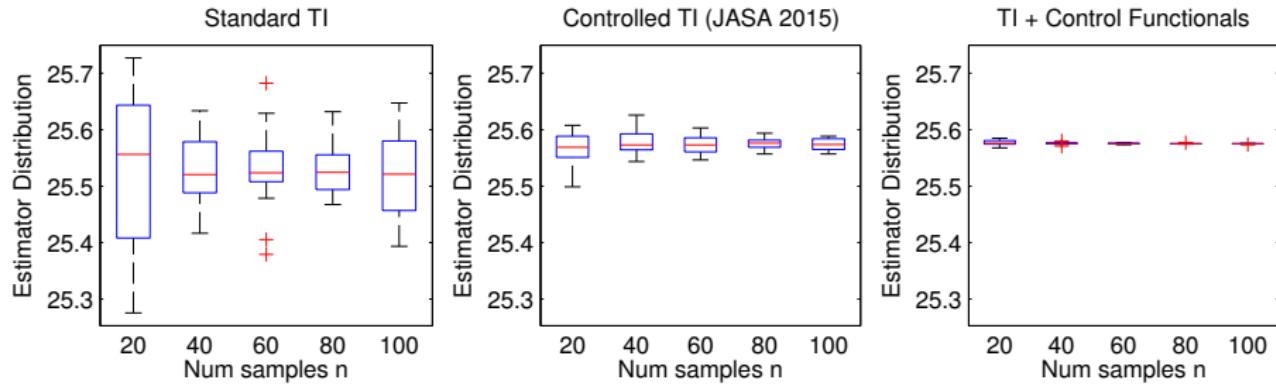


Figure: Computation of marginal likelihood for non-linear ordinary differential equations using thermodynamic integration (TI); van der Pol oscillator example. [Here we show the distribution of 100 independent realisations of each estimator for $\log p(\mathbf{y})$. “Standard TI” is based on arithmetic means. “Controlled TI” is based on ZV control variates.]

Stein's Method

Stein Characterisation

A distribution P is characterised by the pair $(\mathcal{A}, \mathcal{H})$, consisting of a Stein Operator \mathcal{A} and a Stein Class \mathcal{H} , if it holds that

$$X \sim P \quad \text{iff} \quad \mathbb{E}[\mathcal{A}h(X)] = 0 \quad \forall h \in \mathcal{H}.$$

Example presented in Stein, 1972:

- $\mathcal{X} = \mathbb{R}$
- $P = N(\mu, \sigma^2)$ with density function $p(x)$
- $\mathcal{A} : h \mapsto \nabla(hp)/p$
- $\mathcal{H} = \{h : \mathbb{R} \rightarrow \mathbb{R} \text{ s.t. } hp \in W^{1,1} \text{ and } \lim_{x \searrow -\infty} h(x)p(x) = \lim_{x \nearrow +\infty} h(x)p(x)\}.$

Stein's Method

Stein Characterisation

A distribution P is characterised by the pair $(\mathcal{A}, \mathcal{H})$, consisting of a Stein Operator \mathcal{A} and a Stein Class \mathcal{H} , if it holds that

$$X \sim P \quad \text{iff} \quad \mathbb{E}[\mathcal{A}h(X)] = 0 \quad \forall h \in \mathcal{H}.$$

Example presented in Stein, 1972:

- $\mathcal{X} = \mathbb{R}$
- $P = N(\mu, \sigma^2)$ with density function $p(x)$
- $\mathcal{A} : h \mapsto \nabla(hp)/p$
- $\mathcal{H} = \{h : \mathbb{R} \rightarrow \mathbb{R} \text{ s.t. } hp \in W^{1,1} \text{ and } \lim_{x \searrow -\infty} h(x)p(x) = \lim_{x \nearrow +\infty} h(x)p(x)\}.$

Stein's Method

Stein Characterisation

A distribution P is characterised by the pair $(\mathcal{A}, \mathcal{H})$, consisting of a Stein Operator \mathcal{A} and a Stein Class \mathcal{H} , if it holds that

$$X \sim P \quad \text{iff} \quad \mathbb{E}[\mathcal{A}h(X)] = 0 \quad \forall h \in \mathcal{H}.$$

Example presented in Stein, 1972:

- $\mathcal{X} = \mathbb{R}$
- $P = N(\mu, \sigma^2)$ with density function $p(x)$
- $\mathcal{A} : h \mapsto \nabla(hp)/p$ (can be computed based on $p(x) \propto \dots$)
- $\mathcal{H} = \{h : \mathbb{R} \rightarrow \mathbb{R} \text{ s.t. } hp \in W^{1,1} \text{ and } \lim_{x \searrow -\infty} h(x)p(x) = \lim_{x \nearrow +\infty} h(x)p(x)\}.$

Posterior Integration and Stein's Method

This suggests the following, due to

Stein's Method for Posterior Integration

A numerical integration method based on approximation of the integrand by an element

$$\hat{f} \in \mathcal{F} := \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$$

where \hat{f} is yet to be specified.

This could work because $\mathbb{E}_{X \sim P}[c + \mathcal{A}h(X)] = c$ is explicit.

However, it is not yet clear if this will work well:

- ① Is $\mathcal{F} = \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$ rich enough to approximate a generic integrand f ?
- ② Can a suitable element $\hat{f} \in \mathcal{F}$ be found based on finite information \mathcal{D} ?
- ③ Is the method convergent, and if so at what speed?

Posterior Integration and Stein's Method

This suggests the following, due to

Stein's Method for Posterior Integration

A numerical integration method based on approximation of the integrand by an element

$$\hat{f} \in \mathcal{F} := \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$$

where \hat{f} is yet to be specified.

This could work because $\mathbb{E}_{X \sim P}[c + \mathcal{A}h(X)] = c$ is explicit.

However, it is not yet clear if this will work well:

- ① Is $\mathcal{F} = \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$ rich enough to approximate a generic integrand f ?
- ② Can a suitable element $\hat{f} \in \mathcal{F}$ be found based on finite information \mathcal{D} ?
- ③ Is the method convergent, and if so at what speed?

Posterior Integration and Stein's Method

This suggests the following, due to

Stein's Method for Posterior Integration

A numerical integration method based on approximation of the integrand by an element

$$\hat{f} \in \mathcal{F} := \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$$

where \hat{f} is yet to be specified.

This could work because $\mathbb{E}_{X \sim P}[c + \mathcal{A}h(X)] = c$ is explicit.

However, it is not yet clear if this will work well:

- ① Is $\mathcal{F} = \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$ rich enough to approximate a generic integrand f ?
- ② Can a suitable element $\hat{f} \in \mathcal{F}$ be found based on finite information \mathcal{D} ?
- ③ Is the method convergent, and if so at what speed?

Posterior Integration and Stein's Method

This suggests the following, due to

Stein's Method for Posterior Integration

A numerical integration method based on approximation of the integrand by an element

$$\hat{f} \in \mathcal{F} := \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$$

where \hat{f} is yet to be specified.

This could work because $\mathbb{E}_{X \sim P}[c + \mathcal{A}h(X)] = c$ is explicit.

However, it is not yet clear if this will work well:

- ① Is $\mathcal{F} = \{\text{constant functions}\} + \mathcal{A}\mathcal{H}$ rich enough to approximate a generic integrand f ?
- ② Can a suitable element $\hat{f} \in \mathcal{F}$ be found based on finite information \mathcal{D} ?
- ③ Is the method convergent, and if so at what speed?

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- \mathcal{X} is a compact Riemannian manifold, with natural volume form dV .
- $\mathcal{H} \equiv \mathcal{H}(k)$ is a Sobolev space $H^s(\mathcal{X})$, some $s > 2 + \frac{1}{2}\dim(\mathcal{X})$.
- $\mathcal{A}h = \frac{1}{p}\nabla \cdot (p\nabla h)$, a divergence operator on the manifold [?].

Note that $\mathcal{A}h$ can be computed without the normalisation constant, and

$$\begin{aligned} \mathbb{E}_{X \sim P}[\mathcal{A}h(X)] &= \int_{\mathcal{X}} \nabla \cdot (p\nabla h) dV \\ &= 0 \end{aligned} \quad (\text{divergence theorem on the manifold}).$$

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- \mathcal{X} is a compact Riemannian manifold, with natural volume form dV .
- $\mathcal{H} \equiv \mathcal{H}(k)$ is a Sobolev space $H^s(\mathcal{X})$, some $s > 2 + \frac{1}{2}\dim(\mathcal{X})$.
- $\mathcal{A}h = \frac{1}{p}\nabla \cdot (p\nabla h)$, a divergence operator on the manifold [?].

Note that $\mathcal{A}h$ can be computed without the normalisation constant, and

$$\begin{aligned} \mathbb{E}_{X \sim P}[\mathcal{A}h(X)] &= \int_{\mathcal{X}} \nabla \cdot (p\nabla h) dV \\ &= 0 \end{aligned} \quad (\text{divergence theorem on the manifold}).$$

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- \mathcal{X} is a compact Riemannian manifold, with natural volume form dV .
- $\mathcal{H} \equiv \mathcal{H}(k)$ is a Sobolev space $H^s(\mathcal{X})$, some $s > 2 + \frac{1}{2}\dim(\mathcal{X})$.
- $\mathcal{A}h = \frac{1}{p}\nabla \cdot (p\nabla h)$, a divergence operator on the manifold [?].

Note that $\mathcal{A}h$ can be computed without the normalisation constant, and

$$\begin{aligned} \mathbb{E}_{X \sim P}[\mathcal{A}h(X)] &= \int_{\mathcal{X}} \nabla \cdot (p\nabla h) dV \\ &= 0 \end{aligned} \quad (\text{divergence theorem on the manifold}).$$

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- \mathcal{X} is a compact Riemannian manifold, with natural volume form dV .
- $\mathcal{H} \equiv \mathcal{H}(k)$ is a Sobolev space $H^s(\mathcal{X})$, some $s > 2 + \frac{1}{2}\dim(\mathcal{X})$.
- $\mathcal{A}h = \frac{1}{p}\nabla \cdot (p\nabla h)$, a divergence operator on the manifold [?].

Note that $\mathcal{A}h$ can be computed without the normalisation constant, and

$$\begin{aligned} \mathbb{E}_{X \sim P}[\mathcal{A}h(X)] &= \int_{\mathcal{X}} \nabla \cdot (p\nabla h) dV \\ &= 0 \end{aligned} \quad (\text{divergence theorem on the manifold}).$$

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- \mathcal{X} is a compact Riemannian manifold, with natural volume form dV .
- $\mathcal{H} \equiv \mathcal{H}(k)$ is a Sobolev space $H^s(\mathcal{X})$, some $s > 2 + \frac{1}{2}\dim(\mathcal{X})$.
- $\mathcal{A}h = \frac{1}{p}\nabla \cdot (p\nabla h)$, a divergence operator on the manifold [?].

Note that $\mathcal{A}h$ can be computed without the normalisation constant, and

$$\begin{aligned} \mathbb{E}_{X \sim P}[\mathcal{A}h(X)] &= \int_{\mathcal{X}} \nabla \cdot (p\nabla h) dV \\ &= 0 \end{aligned} \quad (\text{divergence theorem on the manifold}).$$

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- $R_1(c) = c^2$
- $R_2(h) = \inf\{\|h'\|_{\mathcal{H}}^2 : h' \in \mathcal{H} \text{ with } \mathcal{A}(h' - h) = 0\}$
- From the representer theorem, \hat{f} and its integral are explicit:

$$\lim_{\lambda_1 \rightarrow 0} \int \hat{f} dP = \frac{\mathbf{1}^\top \mathbf{K}_P^{-1} \mathbf{f}}{\mathbf{1}^\top \mathbf{K}_P^{-1} \mathbf{1}}$$

where $\mathbf{1} = [1, \dots, 1]^\top$, $[\mathbf{K}_P]_{i,j} = \mathcal{A}\bar{\mathcal{A}}k(x_i, x_j)$ and $\mathbf{f} = [f_1, \dots, f_n]^\top$.

- Thus we incur a $O(n^3)$ computational cost in running the method.

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- $R_1(c) = c^2$
- $R_2(h) = \inf\{\|h'\|_{\mathcal{H}}^2 : h' \in \mathcal{H} \text{ with } \mathcal{A}(h' - h) = 0\}$
- From the representer theorem, \hat{f} and its integral are explicit:

$$\lim_{\lambda_1 \rightarrow 0} \int \hat{f} dP = \frac{\mathbf{1}^\top \mathbf{K}_P^{-1} \mathbf{f}}{\mathbf{1}^\top \mathbf{K}_P^{-1} \mathbf{1}}$$

where $\mathbf{1} = [1, \dots, 1]^\top$, $[\mathbf{K}_P]_{i,j} = \mathcal{A} \bar{\mathcal{A}} k(x_i, x_j)$ and $\mathbf{f} = [f_1, \dots, f_n]^\top$.

- Thus we incur a $O(n^3)$ computational cost in running the method.

Posterior Integration and Stein's Method

Consider regularised least-squares:

$$\hat{f} \in \arg \inf_{\substack{c \in \mathbb{R} \\ h \in \mathcal{H}}} \sum_{i=1}^n [f_i - c - \mathcal{A}h(x_i)]^2 + \lambda_1 R_1(c) + \lambda_2 R_2(h)$$

for some regularisers $R_1(c)$ and $R_2(h)$.

In what follows (in part for simplicity):

- $R_1(c) = c^2$
- $R_2(h) = \inf\{\|h'\|_{\mathcal{H}}^2 : h' \in \mathcal{H} \text{ with } \mathcal{A}(h' - h) = 0\}$
- From the representer theorem, \hat{f} and its integral are explicit:

$$\lim_{\lambda_1 \rightarrow 0} \int \hat{f} dP = \frac{\mathbf{1}^\top \mathbf{K}_P^{-1} \mathbf{f}}{\mathbf{1}^\top \mathbf{K}_P^{-1} \mathbf{1}}$$

where $\mathbf{1} = [1, \dots, 1]^\top$, $[\mathbf{K}_P]_{i,j} = \mathcal{A} \bar{\mathcal{A}} k(x_i, x_j)$ and $\mathbf{f} = [f_1, \dots, f_n]^\top$.

- Thus we incur a $O(n^3)$ computational cost in running the method.

Convergence Rate

Convergence Rate (Barp, Oates, Porcu, Girolami, in preparation.)

Under certain regularity assumptions, that include $p \in C^{s+1}(\mathcal{X})$, we have that

$$\|f - \hat{f}\|_{L_2(P)} \lesssim \rho(\{x_i\}_{i=1}^n)^s \|f\|_{W^{s,2}(\mathcal{X})}$$

where

$$\rho(\{x_i\}_{i=1}^n) := \sup_{x \in \mathcal{X}} \min_{i=1,\dots,n} d_{\mathcal{X}}(x, x_i)$$

is the covering distance on the manifold and $d_{\mathcal{X}}$ is the geodesic distance on the manifold.

For $x_i \sim P$ independent, we have $\rho(\{x_i\}_{i=1}^n) \lesssim n^{-\frac{1}{d}} \log(n)^{\frac{1}{d}}$ where $d = \dim(\mathcal{X})$. Thus we recover (up to log-terms) the optimal rate

$$\left| \int_{\mathcal{X}} f dP - \int_{\mathcal{X}} \hat{f} dP \right| \lesssim n^{-\frac{s}{d}} \log(n)^{\frac{s}{d}} \|f\|_{W^{s,2}(\mathcal{X})}.$$

Convergence Rate

Convergence Rate (Barp, Oates, Porcu, Girolami, in preparation.)

Under certain regularity assumptions, that include $p \in C^{s+1}(\mathcal{X})$, we have that

$$\|f - \hat{f}\|_{L_2(P)} \lesssim \rho(\{x_i\}_{i=1}^n)^s \|f\|_{W^{s,2}(\mathcal{X})}$$

where

$$\rho(\{x_i\}_{i=1}^n) := \sup_{x \in \mathcal{X}} \min_{i=1,\dots,n} d_{\mathcal{X}}(x, x_i)$$

is the covering distance on the manifold and $d_{\mathcal{X}}$ is the geodesic distance on the manifold.

For $x_i \sim P$ independent, we have $\rho(\{x_i\}_{i=1}^n) \lesssim n^{-\frac{1}{d}} \log(n)^{\frac{1}{d}}$ where $d = \dim(\mathcal{X})$. Thus we recover (up to log-terms) the optimal rate

$$\left| \int_{\mathcal{X}} f dP - \int_{\mathcal{X}} \hat{f} dP \right| \lesssim n^{-\frac{s}{d}} \log(n)^{\frac{s}{d}} \|f\|_{W^{s,2}(\mathcal{X})}.$$

Convergence Rate

Convergence Rate (Barp, Oates, Porcu, Girolami, in preparation.)

Under certain regularity assumptions, that include $p \in C^{s+1}(\mathcal{X})$, we have that

$$\|f - \hat{f}\|_{L_2(P)} \lesssim \rho(\{x_i\}_{i=1}^n)^s \|f\|_{W^{s,2}(\mathcal{X})}$$

where

$$\rho(\{x_i\}_{i=1}^n) := \sup_{x \in \mathcal{X}} \min_{i=1,\dots,n} d_{\mathcal{X}}(x, x_i)$$

is the covering distance on the manifold and $d_{\mathcal{X}}$ is the geodesic distance on the manifold.

For $x_i \sim P$ independent, we have $\rho(\{x_i\}_{i=1}^n) \lesssim n^{-\frac{1}{d}} \log(n)^{\frac{1}{d}}$ where $d = \dim(\mathcal{X})$. Thus we recover (up to log-terms) the optimal rate

$$\left| \int_{\mathcal{X}} f dP - \int_{\mathcal{X}} \hat{f} dP \right| \lesssim n^{-\frac{s}{d}} \log(n)^{\frac{s}{d}} \|f\|_{W^{s,2}(\mathcal{X})}.$$

Stein Integration

When is Stein integration more accurate than ergodic averages from MCMC?

- **Ergodic averages from MCMC:**

- the computational cost of obtaining n samples is $c = O(n)$
- the error of the ergodic average is $O_P(c^{-\frac{1}{2}})$

- **Stein integration:**

- the computational cost of working with n data points is $c = O(n^3)$
- the error of the estimator is $O((c^{\frac{1}{3}})^{-\frac{s}{d}})$

So Stein is asymptotically more accurate than MCMC iff

$$-\frac{1}{3} \times \frac{s}{d} < -\frac{1}{2}$$

i.e. iff $s > \frac{3}{2}d$.

i.e. “when the smoothness of the integrand is $> \frac{3}{2} \times$ the dimension of the manifold”.

Stein Integration

When is Stein integration more accurate than ergodic averages from MCMC?

- **Ergodic averages from MCMC:**

- the computational cost of obtaining n samples is $c = O(n)$
- the error of the ergodic average is $O_P(c^{-\frac{1}{2}})$

- **Stein integration:**

- the computational cost of working with n data points is $c = O(n^3)$
- the error of the estimator is $O((c^{\frac{1}{3}})^{-\frac{s}{d}})$

So Stein is asymptotically more accurate than MCMC iff

$$-\frac{1}{3} \times \frac{s}{d} < -\frac{1}{2}$$

i.e. iff $s > \frac{3}{2}d$.

i.e. “when the smoothness of the integrand is $> \frac{3}{2} \times$ the dimension of the manifold”.

Stein Integration

When is Stein integration more accurate than ergodic averages from MCMC?

- **Ergodic averages from MCMC:**

- the computational cost of obtaining n samples is $c = O(n)$
- the error of the ergodic average is $O_P(c^{-\frac{1}{2}})$

- **Stein integration:**

- the computational cost of working with n data points is $c = O(n^3)$
- the error of the estimator is $O((c^{\frac{1}{3}})^{-\frac{s}{d}})$

So Stein is asymptotically more accurate than MCMC iff

$$-\frac{1}{3} \times \frac{s}{d} < -\frac{1}{2}$$

i.e. iff $s > \frac{3}{2}d$.

i.e. “when the smoothness of the integrand is $> \frac{3}{2} \times$ the dimension of the manifold”.

Stein Integration

When is Stein integration more accurate than ergodic averages from MCMC?

- **Ergodic averages from MCMC:**

- the computational cost of obtaining n samples is $c = O(n)$
- the error of the ergodic average is $O_P(c^{-\frac{1}{2}})$

- **Stein integration:**

- the computational cost of working with n data points is $c = O(n^3)$
- the error of the estimator is $O((c^{\frac{1}{3}})^{-\frac{s}{d}})$

So Stein is asymptotically more accurate than MCMC iff

$$-\frac{1}{3} \times \frac{s}{d} < -\frac{1}{2}$$

i.e. iff $s > \frac{3}{2}d$.

i.e. “when the smoothness of the integrand is $> \frac{3}{2} \times$ the dimension of the manifold”.

Example: $\mathcal{X} = \mathbb{S}^2$

Recall that for a Riemannian manifold (\mathcal{X}, G) , in local coordinates (q_1, \dots, q_d) ,

$$\nabla h = \sum_{i,j=1}^d [G^{-1}]_{i,j} \frac{\partial h}{\partial q_j} \partial_{q_i}$$

is the gradient on the Riemannian manifold.

Thus we can derive the Stein operator on $\mathcal{X} = \mathbb{S}^2$:

$$G = \begin{pmatrix} \sin^2 q_2 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\mathcal{A}h = \frac{\cos q_2}{\sin q_2} \frac{\partial h}{\partial q_2} + \frac{1}{\sin^2 q_2} \left\{ \frac{1}{p} \frac{\partial p}{\partial q_1} \frac{\partial h}{\partial q_1} + \frac{\partial^2 h}{\partial q_1^2} \right\} + \left\{ \frac{1}{p} \frac{\partial p}{\partial q_2} \frac{\partial h}{\partial q_2} + \frac{\partial^2 h}{\partial q_2^2} \right\}$$

A reproducing kernel for the Sobolev space $H^\alpha(\mathbb{S}^2)$:

$$k(x, y) = C^{(1)} {}_3F_2 \left[\begin{matrix} \frac{3}{2} - \alpha, 1 - \alpha, \frac{3}{2} - \alpha \\ 2 - \alpha, 2 - 2\alpha \end{matrix}; \frac{1 - x \cdot y}{2} \right] + C^{(2)} \|x - y\|^{2\alpha - 2}$$

for $x, y \in \mathbb{S}^2$, $\alpha \in \mathbb{N} + \frac{1}{2}$.

Example: $\mathcal{X} = \mathbb{S}^2$

Recall that for a Riemannian manifold (\mathcal{X}, G) , in local coordinates (q_1, \dots, q_d) ,

$$\nabla h = \sum_{i,j=1}^d [G^{-1}]_{i,j} \frac{\partial h}{\partial q_j} \partial_{q_i}$$

is the gradient on the Riemannian manifold.

Thus we can derive the Stein operator on $\mathcal{X} = \mathbb{S}^2$:

$$G = \begin{pmatrix} \sin^2 q_2 & 0 \\ 0 & 1 \end{pmatrix}$$

$$Ah = \frac{\cos q_2}{\sin q_2} \frac{\partial h}{\partial q_2} + \frac{1}{\sin^2 q_2} \left\{ \frac{1}{p} \frac{\partial p}{\partial q_1} \frac{\partial h}{\partial q_1} + \frac{\partial^2 h}{\partial q_1^2} \right\} + \left\{ \frac{1}{p} \frac{\partial p}{\partial q_2} \frac{\partial h}{\partial q_2} + \frac{\partial^2 h}{\partial q_2^2} \right\}$$

A reproducing kernel for the Sobolev space $H^\alpha(\mathbb{S}^2)$:

$$k(x, y) = C^{(1)} {}_3F_2 \left[\begin{matrix} \frac{3}{2} - \alpha, 1 - \alpha, \frac{3}{2} - \alpha \\ 2 - \alpha, 2 - 2\alpha \end{matrix}; \frac{1 - x \cdot y}{2} \right] + C^{(2)} \|x - y\|^{2\alpha - 2}$$

for $x, y \in \mathbb{S}^2$, $\alpha \in \mathbb{N} + \frac{1}{2}$.

Example: $\mathcal{X} = \mathbb{S}^2$

Recall that for a Riemannian manifold (\mathcal{X}, G) , in local coordinates (q_1, \dots, q_d) ,

$$\nabla h = \sum_{i,j=1}^d [G^{-1}]_{i,j} \frac{\partial h}{\partial q_j} \partial_{q_i}$$

is the gradient on the Riemannian manifold.

Thus we can derive the Stein operator on $\mathcal{X} = \mathbb{S}^2$:

$$G = \begin{pmatrix} \sin^2 q_2 & 0 \\ 0 & 1 \end{pmatrix}$$

$$Ah = \frac{\cos q_2}{\sin q_2} \frac{\partial h}{\partial q_2} + \frac{1}{\sin^2 q_2} \left\{ \frac{1}{p} \frac{\partial p}{\partial q_1} \frac{\partial h}{\partial q_1} + \frac{\partial^2 h}{\partial q_1^2} \right\} + \left\{ \frac{1}{p} \frac{\partial p}{\partial q_2} \frac{\partial h}{\partial q_2} + \frac{\partial^2 h}{\partial q_2^2} \right\}$$

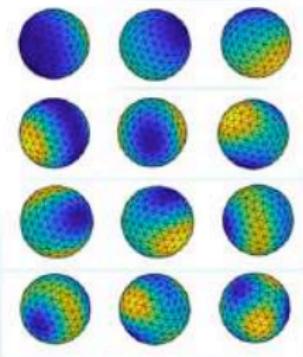
A reproducing kernel for the Sobolev space $H^\alpha(\mathbb{S}^2)$:

$$k(\mathbf{x}, \mathbf{y}) = C^{(1)} {}_3F_2 \left[\begin{matrix} \frac{3}{2} - \alpha, 1 - \alpha, \frac{3}{2} - \alpha \\ 2 - \alpha, 2 - 2\alpha \end{matrix}; \frac{1 - \mathbf{x} \cdot \mathbf{y}}{2} \right] + C^{(2)} \|\mathbf{x} - \mathbf{y}\|^{2\alpha - 2}$$

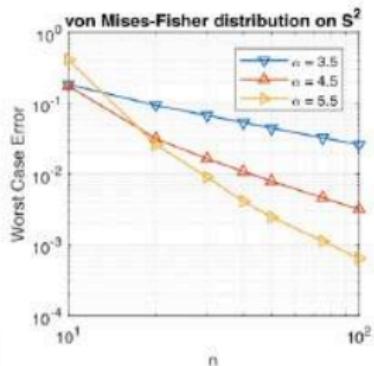
for $\mathbf{x}, \mathbf{y} \in \mathbb{S}^2$, $\alpha \in \mathbb{N} + \frac{1}{2}$.

Example: $\mathcal{X} = \mathbb{S}^2$

Consider the von Mises-Fisher distribution P whose density, with respect to dV , is $p(x) \propto \exp(c^\top x)$:



The first 12 eigenfunctions
of the kernel for \mathcal{AH} .



The worst case error

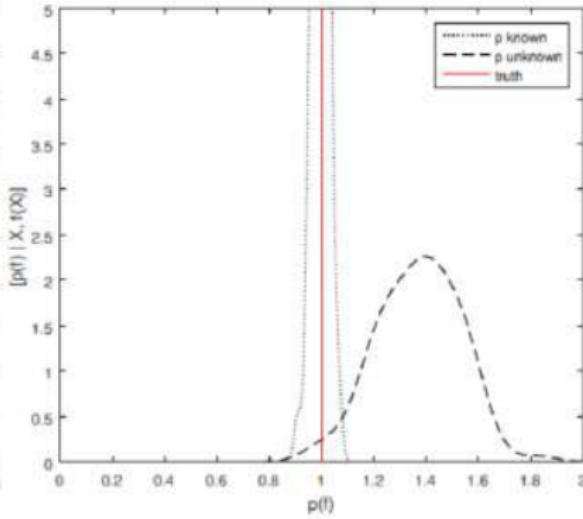
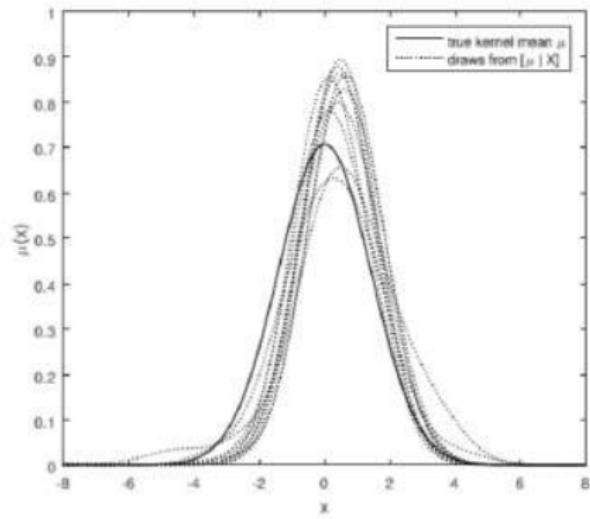
$$\sup_{\|f\|_{\mathcal{AH}} \leq 1} \left| \int_{\mathcal{X}} f dP - \int_{\mathcal{X}} \hat{f} dP \right|$$

plotted here for various smoothness s ($= \alpha$), controlling the differentiability of the kernel, and various n , the number of evaluations of the integrand. The $\{x_i\}$ were quasi-uniform over \mathbb{S}^2 .

Intractable Densities and the Cone of Probability Measures

Ongoing work: BQ for densities $\pi(x)$ only available via samples **Doubly Known Unknowns**, optimal approximating projection in convex cone [Oates et al., 2016a].

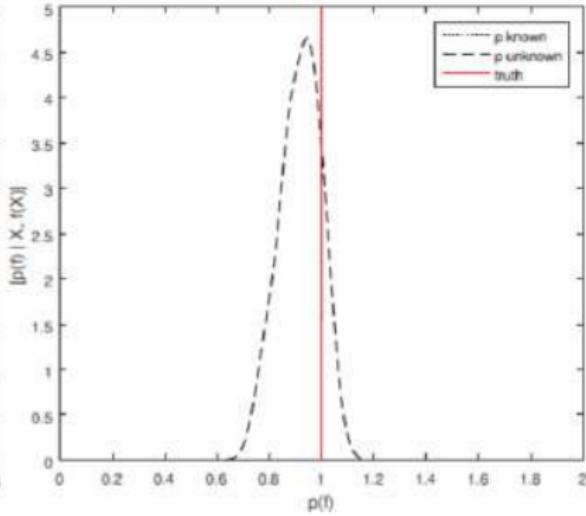
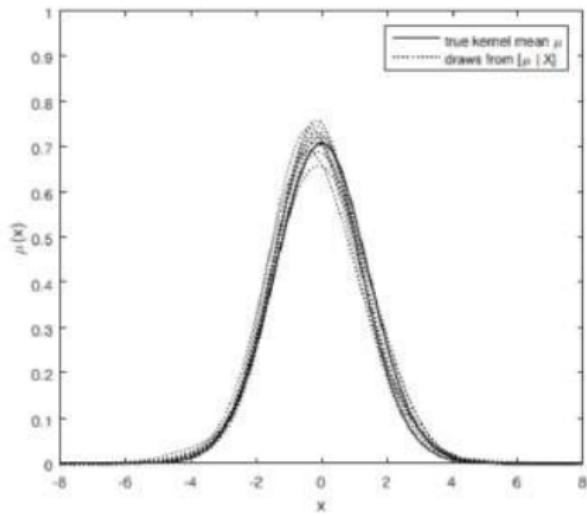
$n = 10$:



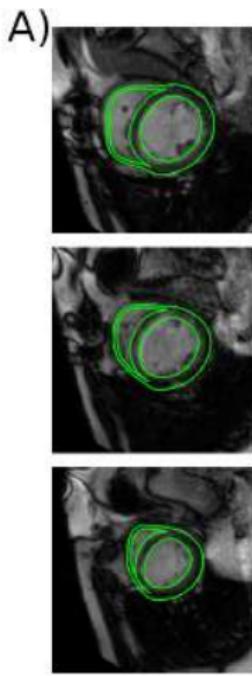
Intractable Densities and the Cone of Probability Measures

Ongoing work: BQ for densities $\pi(x)$ only available via samples, **Doubly Known Unknowns**, optimal approximating projection in convex cone [Oates et al., 2016a].

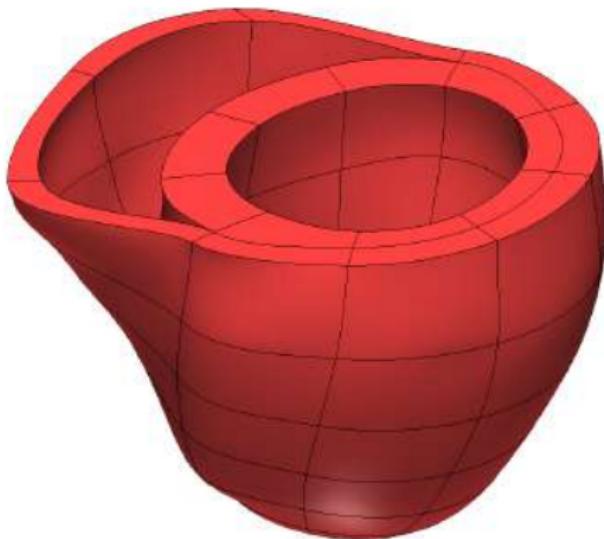
$n = 100$:



Motivation: Assessment of Cardiac Models

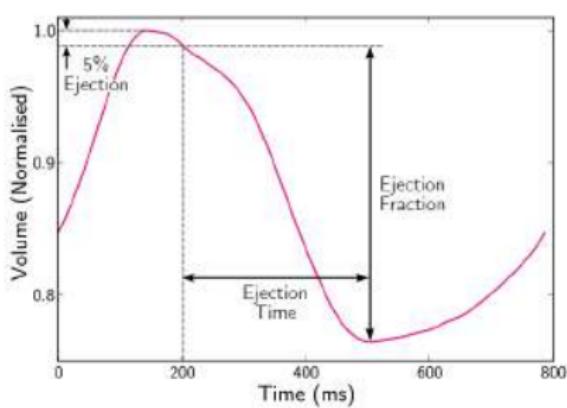
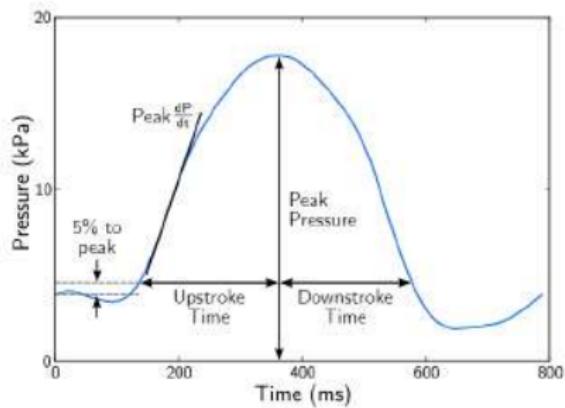


B)



Motivation: Assessment of Cardiac Models

C)



Motivation: Assessment of Cardiac Models

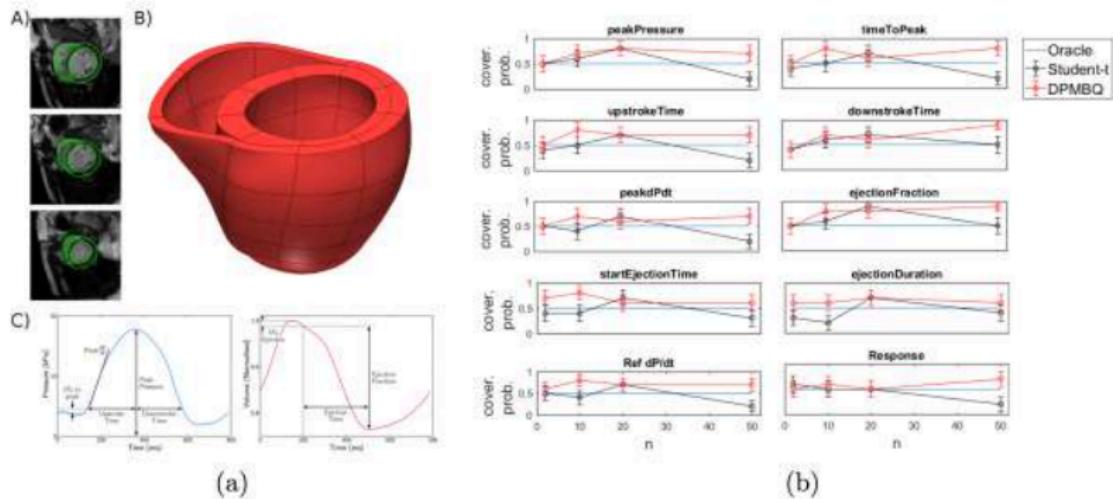


Figure 2: Cardiac model results: (a) Computational cardiac model. A) Segmentation of the cardiac MRI. B) Computational model of the left and right ventricles. C) Schematic image showing the features of pressure (left) and volume transient (right). (b) Comparison of coverage frequencies, for each of 10 numerical integration tasks defined by functionals g_j of the cardiac model output.

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - *come and join us !*

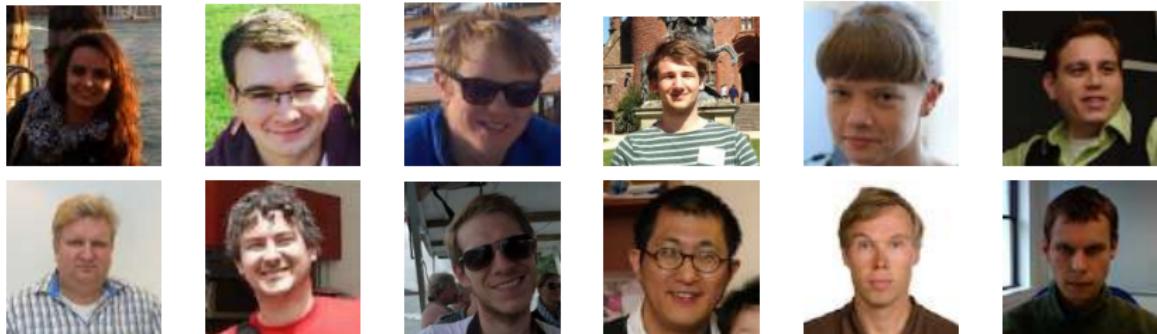
Conclusion

- A role for statistical science in numerical computation?
- A way to formally account for and quantify uncertainty in pipeline of computation
- Contemporary Sciences and Engineering reliant on increasingly sophisticated mathematical objects
- Numerical computation increasingly resorted to in methods and applications
- Quantifying, accounting for uncertainty fundamental to support reasoning and subsequent decision making under uncertainty
- Understanding the impacts of numerical uncertainty is essential for any application related to decision making and risk assessment
- An exciting research area emerging at the intersection of mathematics, statistics and computing science - **come and join us !**

Research Group at Imperial - PAPER SUBMITTED



Engineering and Physical Sciences
Research Council



Research Group at Imperial - PAPER ACCEPTED



Engineering and Physical Sciences
Research Council



- F.-X. Briol, C. J. Oates, M. Girolami, and M. A. Osborne. Frank-Wolfe Bayesian Quadrature: Probabilistic Integration with Theoretical Guarantees. In *Advances In Neural Information Processing Systems 28*, pages 1162–1170, 2015a.
- F.-X. Briol, C. J. Oates, M. Girolami, M. A. Osborne, and D. Sejdinovic. Probabilistic integration: A role for statisticians in numerical analysis? *arXiv preprint arXiv:1512.00933*, 2015b.
- J. Cockayne, C. J. Oates, T. Sullivan, and M. Girolami. Probabilistic Meshless Methods for Partial Differential Equations and Bayesian Inverse Problems. *arXiv preprint arXiv:1605.07811*, 2016.
- P. Diaconis. Bayesian numerical analysis. *Statistical decision theory and related topics IV*, 1:163–175, 1988.
- J. B. Kadane. Parallel and Sequential Computation: A Statistician's view. *Journal of Complexity*, 1:256–263, 1985.
- C. J. Oates and M. Girolami. Control Functionals for Quasi-Monte Carlo Integration. In *Proceedings of the 19th Conference on Artificial Intelligence and Statistics*, 2016.
- C. J. Oates, F.-X. Briol, and M. Girolami. Probabilistic Integration and Intractable Distributions. *arXiv preprint arXiv:1606.06841*, 2016a.
- C. J. Oates, J. Cockayne, F.-X. Briol, and M. Girolami. Convergence Rates for a Class of Estimators Based on Stein's Identity. *arXiv preprint arXiv:1603.03220*, 2016b.
- C. J. Oates, M. Girolami, and N. Chopin. Control functionals for Monte Carlo integration. *Journal of the Royal Statistical Society B: Statistical Methodology*, 2017.
- A. O'Hagan. Bayes–Hermite quadrature. *Journal of Statistical Planning and Inference*, 29:245–260, 1991.
- A. O'Hagan. Some Bayesian numerical analysis. *Bayesian Statistics*, 4:345–363, 1992.
- H. Owhadi. Bayesian numerical homogenization. *arXiv preprint arXiv:1406.6668*, 2014.
- C. Rasmussen and Z. Ghahramani. Bayesian Monte Carlo. In *Advances in Neural Information Processing Systems*, pages 489–496, 2002.
- J. Skilling. Bayesian solution of ordinary differential equations. *Maximum Entropy and Bayesian Methods*, 50:23–37, 1991.