## BAYESIAN PROBABILISTIC NUMERICAL METHODS

J. Cockayne<sup>1</sup>

M. Girolami<sup>1,2</sup>

C. J. Oates<sup>1,3</sup>

T. J. Sullivan<sup>4,5</sup>

Neyman Statistics Seminar UC Berkeley, CA, US 29 January 2020

<sup>1</sup>Alan Turing Institute, London, UK <sup>2</sup>University of Cambridge, UK <sup>3</sup>Newcastle University, UK <sup>4</sup>Freie Universität Berlin, DE <sup>5</sup>Zuse Institute Berlin, DE "Numerical analysts and statisticians are both in the business of estimating parameter values from incomplete information. The two disciplines have separately developed their own approaches to formalizing strangely similar problems and their own solution techniques; the author believes they have much to offer each other."

— F. M. Larkin (1979b)



## OVERVIEW: PNMs and BIPs

- There are many reasons to consider a probabilistic/statistical perspective on the analysis and design of numerical methods, and even to return probabilistic solutions to deterministic forward problems like quadrature / DE solution.
- In various forms, these ideas have a long history.

 $\rightarrow$  Oates and Sullivan (2019) Stat. Comp. arXiv:1901.04457

- What are probabilistic numerical methods (PNMs) and in what sense can they be Bayesian? → Cockayne et al. (2019) SIAM Rev. arXiv:1702.03673
- A Bayesian interpretation of forward problems is especially appealing for Bayesian inverse problems (BIPs), since then both the forward and inverse problem "speak the same language", without spurious posterior over-concentration.
- How does their use connect to established theory for BIPs?

 $\rightarrow$  Lie et al. (2018) SIAM/ASA JUQ arXiv:1712.05717

## MOTIVATING EXAMPLE: FITZHUGH-NAGUMO ODE INFERENCE

• Nonlinear FitzHugh–Nagumo oscillator  $u: [0, T] \rightarrow \mathbb{R}^2$ :

$$\frac{\mathrm{d}u}{\mathrm{d}t} = f(\mathbf{u}) \coloneqq \begin{bmatrix} u_1 - \frac{u_1^3}{3} + u_2 \\ -\frac{1}{\theta_3}(u_1 - \theta_1 + \theta_2 u_2) \end{bmatrix}$$

- Aim: recover  $\theta \in \mathbb{R}^3_{>0}$  from observations  $y_i = u(t_i^{obs}) + \eta_i$  taken at discrete times  $t_i^{obs} = 0, 1, \dots, 40$ , with  $\eta_i \sim \mathcal{N}(0, 10^{-3}l)$  i.i.d.
- Take ground truth u(0) = (-1, 1) and  $\theta = (0.2, 0.2, 3)$ ; generate data y from a reference trajectory using RK4 with time step  $\tau = \tau_{ref} = 10^{-3}$ .
- Infer  $\theta$  using PN-Euler solvers with local noise  $\xi$  of variance  $\propto \sigma \tau^3$  and hence strong error  $\mathbb{E}[\sup_{0 \le t \le \tau} ||u(t) u^{\mathsf{PN}}(t)||^2] \le C\tau^2$  (Conrad et al., 2016; Lie et al., 2019).
- Log-normal prior for  $\theta \rightsquigarrow$  marginal Bayesian posterior  $\mathbb{E}_{\xi} [\mathbb{P}[\theta|y, \sigma, \tau, \xi]]$  for various  $\tau \gg \tau_{\text{ref}}$  and  $\sigma \ge 0$ .

#### MOTIVATING EXAMPLE: FITZHUGH-NAGUMO ODE INFERENCE



**Figure 1:** The deterministic posteriors (i.e.  $\sigma = 0$ ) are over-confident at all values of the time step  $\tau = 0.1, 0.05, 0.02, 0.01, 0.005$ , often do not overlap, and are biased.

#### MOTIVATING EXAMPLE: FITZHUGH-NAGUMO ODE INFERENCE



**Figure 1:** In contrast, the PN-Euler posteriors (here with  $\sigma = 1/5$ ) for  $\tau = 0.1, 0.05, 0.02, 0.01, 0.005$  are less confident and overlap more; admittedly they are still biased.

#### OUTLINE

A little history

Numerics: An inference perspective

Optimal information operators

Disintegration

Coherent pipelines of PNMs, and Bayesian inverse problems

Applications

Closing remarks

# A LITTLE HISTORY

"Je suppose que l'on sache a priori que la fonction f(x) est développable, dans une certain domaine, suivant les puissances croissantes des x,

$$f(x) = A_0 + A_1 x + \dots$$

Nous ne savons rien sur les A, sauf que la probabilité pour que l'un d'eux,  $A_i$ , soit compris entre certaines limites, y et y + dy, est

$$\sqrt{\frac{h_i}{\pi}}e^{-h_iy^2}\,\mathrm{d}y.$$

Nous connaissons par n observations

$$f(a_1) = B_1, \qquad f(a_2) = B_2, \qquad \cdots \cdots \cdots \qquad f(a_n) = B_n.$$

Nous cherchons la valeur probable de f(x) pour une autre valeur de x."

#### **ROUND-OFF ERROR**

- What about probabilistic numerical methods for use on a computer?
- The limited nature of the earliest computers led authors to focus initially on the phenomenon of round-off error (Henrici, 1962; Hull and Swenson, 1966; von Neumann and Goldstine, 1947), whether of fixed-point or floating-point type, without any particular statistical *inferential* motivation; indeed, this aspect is still alive (Barlow and Bareiss, 1985; Chatelin and Brunet, 1990; Tienari, 1970).
- One early, utilitarian view is that probabilistic models in computation are just useful shortcuts:

"[Round-off errors] are strictly very complicated but uniquely defined number theoretical functions [of the inputs], yet our ignorance of their true nature is such that we best treat them as random variables."

- von Neumann and Goldstine (1947, p. 1027)

# AL'BERT VALENTINOVICH SUL'DIN

- One of the earliest attempts to statistically motivate a numerical algorithm was due to A. V.
   Sul'din (1924–1996), working at Kazan State University in the USSR.
- After first making contributions to the study of Lie algebras, towards the end of the 1950s
   Sul'din turned his attention to computational and applied mathematics, and in particular to probabilistic and statistical methodology.
- His work led to the establishment of the Faculty of Computational Mathematics and Cybernetics in Kazan, of which he was the founding Dean.



Al'bert Valentinovich Sul'din (1924–1996) © Kazan Federal University, reproduced with permission.

# FREDERICK MICHAEL ("MIKE") LARKIN

- On the other side of the Iron Curtain, between 1957 and 1969, Frederick Michael ("Mike") Larkin (1936–1982) worked for the UK Atomic Energy Authority in its laboratories at Harwell and Culham, as well as working for two years at Rolls Royce; from 1969, he was at Queen's University in Kingston, Ontario, Canada.
- Following a parallel path to that of Sul'din, over the next decade Larkin would further blend numerical analysis and statistical thinking (Kuelbs et al., 1972; Larkin, 1969, 1972, 1974, 1979a,b,c), arguably laying the foundations of modern PN as a kind of inference.



Frederick Michael Larkin (1936–1982) © (Larkin et al., 1967, reproduced with permission).

#### LARKIN

- Larkin worked on building some of the first graphical calculators, called GHOST (short for graphical output system), and the GHOUL (graphical output language) – perhaps a motivation for seeking a richer description of numerical error.
- The perspective developed by Larkin was fundamentally statistical and, in modern terminology, the probabilistic numerical methods he developed would be described as *Bayesian* — though Larkin used the term *relative likelihood* for the prior.
- Larkin's perspective on quadrature: consider the Wiener measure as a prior, the information  $(t_j, u(t_j))_{j=1}^J$  as (noiseless) data, and output the posterior marginal for  $\int_a^b u(t) dt$  what we would now recognise as a probabilistic numerical method: "[This] permits, at least in principle, the derivation of joint probability density functions for [both observed and unobserved] functionals on the space and also allows us to evaluate confidence limits on the estimate of a required functional (in terms of given values of other functionals)." Larkin (1972) <sup>10/53</sup>

- We wish to approximate the definite integral  $\int_a^b u(t) dt$  of  $u \in \mathcal{U} := C^0([a, b]; \mathbb{R})$  under a statistical assumption that  $(u(t) u(a))_{t \in [a,b]}$  follows a standard Brownian motion (Wiener measure,  $\mu_W$ ).
- We receive pointwise data about u in the form of the values of u at  $J \in \mathbb{N}$  nodes  $a = t_1 < t_2 < \cdots < t_J = b$ .
- In more statistical language, anticipating the terminology of Cockayne et al. (2019):
  - we have a latent quantity (integrand) u living in a space  $\mathcal{U}$ ,
  - our observed data or information concerning u is  $y := (t_j, u(t_j))_{j=1}^j$ , living in the space  $\mathcal{Y} := ([a, b] \times \mathbb{R})^j$ ,
  - and we care about the quantity of interest  $Q(u) \coloneqq \int_a^b u(t) dt$ , living in  $\mathcal{Q} \coloneqq \mathbb{R}$ .

## LARKIN V. SUL'DIN ON UNIVARIATE QUADRATURE II

• Sul'din (1959, 1960, 1963) showed by direct calculation that the quadrature rule B:  $\mathcal{Y} \to \mathbb{R}$  that minimises the mean squared error

$$\int_{\boldsymbol{\mathcal{U}}} \left| \int_{a}^{b} \boldsymbol{u}(t) \, \mathrm{d}t - \mathsf{B}\big( (t_j, \boldsymbol{u}(t_j))_{j=1}^{J} \big) \right|^2 \mu_{\mathsf{W}}(\mathrm{d}\boldsymbol{u})$$

is the classical trapezoidal rule

$$\mathsf{B}_{\mathsf{tr}}\big((t_j, z_j)_{j=1}^J\big) \coloneqq \frac{1}{2} \sum_{j=1}^{J-1} (z_{j+1} + z_j)(t_{j+1} - t_j) = z_1 \frac{t_2 - t_1}{2} + \sum_{j=2}^{J-1} z_j \frac{t_{j+1} - t_{j-1}}{2} + z_J \frac{t_j - t_{j-1}}{2},$$

i.e. the definite integral of the piecewise linear interpolant of the observed data.

## LARKIN V. SUL'DIN ON UNIVARIATE QUADRATURE III

- Thus, Sul'din described the trapezoidal rule  $B_{tr}$  as a frequentist point estimator obtained from minimising the mean square error, which "just happens" to produce an unbiased estimator with variance  $\frac{1}{12}\sum_{j=1}^{j-1}(t_{j+1}-t_j)^3$ .
- However, Larkin saw the normal distribution

$$\mathcal{N}\Big(\mathsf{B}_{\mathsf{tr}}\big((t_j, z_j)_{j=1}^J\big), \frac{1}{12}\sum_{j=1}^{J-1}(t_{j+1} - t_j)^3\Big)$$

on  $\mathbb{R}$  as the measure-valued output of a probabilistic quadrature rule, of which  $B_{tr}((t_j, z_j)_{j=1}^{J})$  is a convenient point summary. *En passant* he made fundamental contributions to the study of Gaussian measures (Kuelbs et al., 1972; Larkin, 1972).

• Neither Larkin nor Sul'din would have had access to the computing resources needed to realise their more general (nonlinear, non-Gaussian) vision.

## Optimal numerical methods are Bayes rules (1980–1990) i

- The average-case analysis (ACA) of numerical methods received interest and built on the work of Kolmogorov (1936) and Sard (1963).
- In ACA the performance of a numerical method is assessed in terms of its *average error* with respect to a probability measure over the problem set; a prime example is univariate quadrature with the average quadratic loss given earlier.
- A traditional (deterministic) NM can also be regarded as a decision rule and the probability measure used in ACA can be used to instantiate the Bayesian decision-theoretic framework (Berger, 1985). The average error is then recognised as the *expected loss*, also called the *risk*. ACA is mathematically equivalent to Bayesian decision theory restricted to the case of an experiment that produces a deterministic dataset (Kimeldorf and Wahba, 1970a,b; Parzen, 1970; Larkin, 1970).

ACA optimal methods are Bayes rules or Bayes acts in the decision-theoretic context. Kadane and Wasilkowski (1985) had the insight that ACA-optimal methods coincide with (non-randomised) Bayes rules when the measure used to define the MSE is the Bayesian prior. Recently it has become clear that ACA and Bayesian optimality differ in general (Cockayne et al., 2019; Oates et al., 2019b).

#### INFORMATION-BASED COMPLEXITY

- Information-based complexity (IBC) developed simultaneously with ACA, with the aim of relating the computational complexity and optimality properties of algorithms to the available information on the unknowns.
- Smale (1985) compared the accuracies (with respect to mean absolute error) for a given cost of the Riemann sum, trapezoidal, and Simpson quadrature rules; in the same paper, Smale also considered root-finding, optimisation via linear programming, and the solution of systems of linear equations.
- Diaconis (1988) repeated Sul'din's observation that the posterior mean for \$\int\_a^b u(t) dt\$ under the Wiener measure prior is the trapezoidal method, which is a ACA-optimal, and posed a further question: can other numerical methods for other tasks be similarly recovered as Bayes rules in a decision-theoretic framework? For linear cubature methods, a positive and constructive answer was recently provided by Karvonen et al. (2018), but the general question remains open. 16/53

- Research interest in PN was revived by contributions from on quadrature (Minka, 2000; O'Hagan, 1991; Rasmussen and Ghahramani, 2003), each to a greater or lesser extent a rediscovery of earlier work due to Larkin (1972). In each case the algorithmic output was considered to be a probability distribution over the quantity of interest.
- The 1990s saw an expansion in the PN agenda, first with early work on an area that would become Bayesian optimisation (Močkus, 1975, 1977, 1989).
- Skilling (1992) presented a (partially) Bayesian perspective on the numerical solution of ODE initial value problems, explicitly considering, e.g., the role of regularity assumptions on the vector field, prior and likelihood choice, and sampling strategies.

 Skilling himself considered his then-new explicit emphasis on a Bayesian statistical approach to be quite natural:

"This paper arose from long exposure to Laplace/Cox/Jaynes probabilistic reasoning, combined with the University of Cambridge's desire that the author teach some (traditional) numerical analysis. The rest is common sense. [...] Simply, Bayesian ideas are 'in the air'." — Skilling (1992)

- The machine learning community took up the ODE theme again ≈ 5 years ago (Schober et al., 2014), provoking further mathematical analysis (Conrad et al., 2016; Lie et al., 2019) and then an explosion of more general studies.
- Gaussian process techniques also work well for some PDEs (Cockayne et al., 2016, 2017; Raissi et al., 2018).

#### Conceptual evolution – A summary

- In the traditional setting of numerical analysis, c. 1950, all objects and operations are seen as being strictly deterministic. These deterministic objects are sometimes exceedingly complicated, to the extent that they may be treated as being stochastic.
- 2. Sard and Sul'din consider the questions of optimal performance of a numerical method in, respectively, the worst-case and the average-case context. Some of the average-case performance measures amount to variances of point estimators but are not *viewed* as such; probabilistic aspects are not a motivating factor.
- 3. Larkin's innovation, 1960s–1970s, is to formulate numerical tasks in terms of a joint distribution over latent quantities and quantities of interest; the quantity of interest is a stochastic object. Larkin summarises his posterior distributions using a point estimator accompanied by a credible interval.
- 4. The fully modern viewpoint, 2017+, is to explicitly think of the output as a probability measure to be realised, sampled, and possibly summarised. 19/53

# AN INFERENCE PERSPECTIVE ON NUMERICAL TASKS

An abstraction of a numerical task consists of three spaces and three functions:

- $\mathcal{U}$ , where an unknown/variable object u lives;
- $\mathcal{Q}$ , with a quantity of interest Q:  $\mathcal{U} \rightarrow \mathcal{Q}$ ;
- $\mathcal{Y}$ , where we observe information Y(u), via a function  $Y: \mathcal{U} \to \mathcal{Y}$ . dim  $\mathcal{Y} < \infty$

 $\dim \mathcal{U} = \infty$ 

An abstraction of a numerical task consists of three spaces and three functions:

- $\mathcal{U}$ , where an unknown/variable object u lives;
- $\mathcal{Q}$ , with a quantity of interest Q:  $\mathcal{U} \rightarrow \mathcal{Q}$ ;
- $\mathcal{Y}$ , where we observe information Y(u), via a function  $Y: \mathcal{U} \to \mathcal{Y}$ . dim  $\mathcal{Y} < \infty$

# Example (Quadrature)

$$\mathcal{U} = C^0([0,1];\mathbb{R}) \qquad \qquad \mathcal{Y} = ([0,1] \times \mathbb{R})^m \qquad \qquad \mathcal{Q} = \mathbb{R}$$
$$Y(\mathbf{u}) = (t_i, \mathbf{u}(t_i))_{i=1}^m \qquad \qquad Q(\mathbf{u}) = \int_0^1 \mathbf{u}(t) \, dt$$

 $\dim \mathcal{U} = \infty$ 

An abstraction of a numerical task consists of three spaces and three functions:

- $\mathcal{U}$ , where an unknown/variable object u lives;
- $\mathcal{Q}$ , with a quantity of interest Q:  $\mathcal{U} \rightarrow \mathcal{Q}$ ;
- $\mathcal{Y}$ , where we observe information Y(u), via a function  $Y: \mathcal{U} \to \mathcal{Y}$ . dim  $\mathcal{Y} < \infty$

 $\dim \mathcal{U} = \infty$ 

## Example (Quadrature)

$$\mathcal{U} = C^{0}([0,1];\mathbb{R}) \qquad \qquad \mathcal{Y} = ([0,1] \times \mathbb{R})^{m} \qquad \qquad \mathcal{Q} = \mathbb{R}$$
$$Y(\mathbf{u}) = (t_{i}, \mathbf{u}(t_{i}))_{i=1}^{m} \qquad \qquad Q(\mathbf{u}) = \int_{0}^{1} \mathbf{u}(t) \, \mathrm{d}t$$

• Conventional numerical methods are cleverly-designed functions B:  $\mathcal{Y} \rightarrow \mathcal{Q}$ : such a method **"believes"** that  $Q(u) \approx B(Y(u))$ . 20/53

#### AN ABSTRACT VIEW OF NUMERICAL METHODS II

## Example (Quadrature)

$$\mathcal{U} = C^{0}([0,1];\mathbb{R}) \qquad \qquad \mathcal{Y} = ([0,1] \times \mathbb{R})^{m} \qquad \qquad \mathcal{Q} = \mathbb{R}$$
$$Y(u) = (t_{i}, u(t_{i}))_{i=1}^{m} \qquad \qquad Q(u) = \int_{0}^{1} u(t) dt$$

- Some but not all methods  $B: \mathcal{Y} \to \mathcal{Q}$  try to invert Y, estimate u, then apply Q.
  - E.g. the trapezoidal rule does estimate *u*, via piecewise linear interpolation:

$$\mathsf{B}_{\mathsf{tr}}\big((t_j, z_j)_{j=1}^j\big) \coloneqq \sum_{j=1}^{J-1} \frac{z_{j+1} + z_j}{2} (t_{j+1} - t_j) = z_1 \frac{t_2 - t_1}{2} + \sum_{j=2}^{J-1} z_j \frac{t_{j+1} - t_{j-1}}{2} + z_J \frac{t_j - t_{j-1}}{2}.$$

E.g. vanilla Monte Carlo does not estimate <u>u</u>! (cf. O'Hagan, 1987)

$$B_{MC}((t_i, z_i)_{i=1}^n) \coloneqq \frac{1}{n} \sum_{i=1}^n z_i$$
 21/53

#### AN ABSTRACT VIEW OF NUMERICAL METHODS III

- Question: What makes for a "good" numerical method? (Larkin, 1970)
- Answer 1, Gauss:  $B \circ Y = Q$  on a "large" finite-dimensional subspace of  $\mathcal{U}$ .
- Answer 2, Sard (1949): residual  $B \circ Y Q$  is "small" on  $\mathcal{U}$ . In what sense?
  - The worst-case error:

$$e_{\mathsf{WC}} \coloneqq \sup_{u \in \mathcal{U}} \|\mathsf{B}(\mathsf{Y}(u)) - \mathsf{Q}(u)\|_{\mathcal{Q}}.$$

• The average-case error (Ritter, 2000) with respect to a probability measure  $\mu \in \mathcal{P}_{\mathcal{U}}$ :

$$e_{\mathsf{AC}} \coloneqq \int_{\mathcal{U}} \|\mathsf{B}(\mathsf{Y}(u)) - \mathsf{Q}(u)\|_{\mathcal{Q}} \, \mu(\mathrm{d} u).$$

To a Bayesian, seeing the additional structure of μ, there is only one way forward: if u ~ μ, then B(Y(u)) should be obtained by conditioning μ and then applying Q. But is this Bayesian solution always well-defined, and what are its error properties?





Go Probabilistic!

$$\mu \in \mathcal{P}_{\mathcal{U}}$$
$$(\mathsf{Y}_{\sharp}\mu)(E) := \mu(\mathsf{Y}^{-1}(E)$$



## Example (Quadrature)

$$\mathcal{U} = C^{0}([0,1];\mathbb{R}) \qquad \qquad \mathcal{Y} = ([0,1] \times \mathbb{R})^{m} \qquad \qquad \mathcal{Q} = \mathbb{R}$$
$$Y(\boldsymbol{u}) = (t_{i}, \boldsymbol{u}(t_{i}))_{i=1}^{m} \qquad \qquad Q(\boldsymbol{u}) = \int_{0}^{1} \boldsymbol{u}(t) \, \mathrm{d}t$$

A classical numerical method B uses only the spaces and data to produce a point estimate of Q(u). We could engage in average-case analysis of B against  $\mu$ .



Go Probabilistic!

$$\begin{split} \mu \in \mathcal{P}_{\mathcal{U}} \\ (\mathsf{Y}_{\sharp} \mu)(E) \coloneqq \mu(\mathsf{Y}^{-1}(E)) \end{split}$$



## Example (Quadrature)

$$\mathcal{U} = C^{0}([0,1];\mathbb{R}) \qquad \qquad \mathcal{Y} = ([0,1] \times \mathbb{R})^{m} \qquad \qquad \mathcal{Q} = \mathbb{R}$$
$$Y(\boldsymbol{u}) = (t_{i}, \boldsymbol{u}(t_{i}))_{i=1}^{m} \qquad \qquad Q(\boldsymbol{u}) = \int_{0}^{1} \boldsymbol{u}(t) \, \mathrm{d}t$$

A classical numerical method B uses only the spaces and data to produce a point estimate of Q(u). A probabilistic numerical method converts an additional belief  $\mu \in \mathcal{P}_{\mathcal{U}}$  about u into a belief  $\beta(\mu, Y(u)) \in \mathcal{P}_{\mathcal{Q}}$  about Q(u).



Go Probabilistic!

 $\mu \in \mathcal{P}_{\mathcal{U}}$ 

$$(\mathsf{Y}_{\sharp}\mu)(E) \coloneqq \mu(\mathsf{Y}^{-1}(E))$$



## Definition (Bayesian PNM)

A PNM  $\beta(\mu, \cdot): \mathcal{Y} \to \mathcal{P}_{\mathcal{Q}}$  with prior  $\mu \in \mathcal{P}_{\mathcal{U}}$  is **Bayesian** for a Qol Q:  $\mathcal{U} \to \mathcal{Q}$  and information operator Y:  $\mathcal{U} \to \mathcal{Y}$  if the bottom-left  $\mathcal{Y} - \mathcal{P}_{\mathcal{U}} - \mathcal{P}_{\mathcal{Q}}$  triangle commutes, i.e. the output of  $\beta$  is the push-forward of the conditional distribution  $\mu^{\mathcal{Y}}$  through Q:

$$\beta(\mu, y) = Q_{\sharp}\mu^{y}, \quad \text{ for } Y_{\sharp}\mu\text{-almost all } y \in \mathcal{Y}.$$

# Definition (Bayesian PNM)

A PNM  $\beta$  with prior  $\mu \in \mathcal{P}_{\mathcal{U}}$  is **Bayesian** for a quantity of interest Q and information Y if its output is exactly the image of the conditional distribution<sup>\*</sup>  $\mu^{\gamma} = \mu | [Y = y]$  under Q:

 $\beta(\mu, y) = Q_{\sharp}\mu^{y}, \quad \text{for } Y_{\sharp}\mu\text{-almost all } y \in \mathcal{Y}.$ 

(\* conditioning in the sense of a disintegration - more later)

# Definition (Bayesian PNM)

A PNM  $\beta$  with prior  $\mu \in \mathcal{P}_{\mathcal{U}}$  is **Bayesian** for a quantity of interest Q and information Y if its output is exactly the image of the conditional distribution<sup>\*</sup>  $\mu^{y} = \mu | [Y = y]$  under Q:

 $\beta(\mu, y) = Q_{\sharp}\mu^{y}, \quad \text{for } Y_{\sharp}\mu\text{-almost all } y \in \mathcal{Y}.$ 

(\* conditioning in the sense of a disintegration - more later)

#### Example

- Under the Gaussian Brownian motion prior on U = C<sup>0</sup>([0, 1]; ℝ), the posterior mean
   / MAP estimator for the definite integral is the trapezoidal rule, i.e. integration using linear interpolation (Sul'din, 1959, 1960).
- Integrated Brownian motion prior  $\leftrightarrow$  integration using cubic spline interpolation<sub>24/55</sub>

# A ROGUE'S GALLERY OF BAYESIAN AND NON-BAYESIAN PNMS (2017)

Method	<b>QoI</b> $Q(x)$	Information $A(x)$	Non-Bayesian PNMs	Bayesian PNMs <sup>1</sup>
Integrator	$\int x(t)\nu(\mathrm{d}t)$	${x(t_i)}_{i=1}^n$	Approximate Bayesian Quadrature Methods [Os-	Bayesian Quadrature [Diaconis, 1988, O'Hagan,
			borne et al., 2012b,a, Gunter et al., 2014]	1991, Ghahramani and Rasmussen, 2002, Briol
	$\int f(t) r(dt)$	$f_{t} : \mathbb{I}^{n} \to \mathbb{S}^{t} \to \mathbb{C}^{n}$	Kong et al [2003] Tan [2004] Kong et al [2007]	et al., 2016]
	$\int x_1(t) x_2(\mathrm{d}t)$	$ \{(t_i, x_1(t_i))\}_{i=1}^n \text{ s.t. } t_i \sim x_2 $	rong et al. [2003], Tan [2004], Hong et al. [2007]	Oates et al. [2016]
Optimiser	$\arg \min x(t)$	$\{x(t_i)\}_{i=1}^n$		Bayesian Optimisation [Mockus, 1989] <sup>6</sup>
-		$\{ abla x(t_i)\}_{i=1}^n$		Hennig and Kiefel [2013]
		$\{(x(t_i), \nabla x(t_i)\}_{i=1}^n$		Probabilistic Line Search [Mahsereci and Hennig,
		(m), 13 m		2015]
		$\{\mathbb{I}[t_{\min} < t_i]\}_{i=1}^n$		Probabilistic Bisection Algorithm Horstein,
		$\{I[t_{n}) < t_{n}\} + \operatorname{error} \{n\}$	Waeher et al [2013]	1903
	-11	$\left[1\left[t_{\min} < t_{i}\right] + \operatorname{error}_{i=1}\right]$		
Linear Solver	x ~0	$\{xt_i\}_{i=1}$		Probabilistic Linear Solvers Hennig, 2015, Bartels
				and Hennig, 2016
ODE Solver	x	$\{\nabla x(t_i)\}_{i=1}^n$	Filtering Methods for IVPs [Schober et al., 2014,	Skilling [1992]
			Chkrebtii et al., 2016, Kersting and Hennig, 2016,	
			Teymur et al., 2016, Schober et al., 2016] <sup>4</sup> Finite	
			Difference Methods [John and Wu, 2017] <sup>7</sup>	
		$\nabla x$ + rounding error	Hull and Swenson [1966], Mosbach and Turner	
			$[2009]^2$	
	$x(t_{ m end})$	$\{\nabla x(t_i)\}_{i=1}^n$	Stochastic Euler [Krebs, 2016]	
PDE Solver	x	${Dx(t_i)}_{i=1}^n$	Chkrebtii et al. [2016]	Probabilistic Meshless Methods [Owhadi,
				2015a,b, Cockayne et al., 2016, Raissi et al., 2016
		Dx + discretisation error	Conrad et al. [2016] <sup>3</sup>	

Optimal information operators: the Worst, the Average, and the Bayesian
Suppose we have a loss function  $L: \mathcal{Q} \times \mathcal{Q} \to \mathbb{R}$ , e.g.  $L(q,q') \coloneqq ||q-q'||_{\mathcal{Q}}^2$ .

 $\blacksquare$  The worst-case loss for a classical numerical method B:  $\mathcal{Y} \to \mathcal{Q}$  is

$$e_{WC}(Y,B) \coloneqq \sup_{u \in \mathcal{U}} L(B(Y(u)), Q(u)).$$

• The average-case loss under a probability measure  $\mu \in \mathcal{P}_{\mathcal{U}}$  is

$$\begin{split} e_{\mathsf{AC}}(\mathsf{Y},\mathsf{B}) &\coloneqq \int_{\mathcal{U}} L\big(\mathsf{B}(\mathsf{Y}(\boldsymbol{u})),\mathsf{Q}(\boldsymbol{u})\big)\,\mu(\mathrm{d}\boldsymbol{u}),\\ e_{\mathsf{AC}}(\mathsf{Y},\boldsymbol{\beta}) &\coloneqq \int_{\mathcal{U}} \left[\int_{\mathcal{Q}} L\big(q,\mathsf{Q}(\boldsymbol{u})\big)\,\boldsymbol{\beta}(\mu,\mathsf{Y}(\boldsymbol{u}))(\mathrm{d}q)\right] \mu(\mathrm{d}\boldsymbol{u}). \end{split}$$

For a BPNM β, we must have β(μ, y) = Q<sub>↓</sub>μ<sup>y</sup> once Y(u) = y is given; optimality of Y means minimising the Bayesian loss

$$e_{\mathsf{BPN}}(\mathsf{Y}) \coloneqq \int_{\mathcal{U}} \left[ \int_{\mathcal{Q}} L(q, \mathsf{Q}(\boldsymbol{u})) \left( \mathsf{Q}_{\sharp} \boldsymbol{\mu}^{\mathsf{Y}(\boldsymbol{u})} \right) (\mathrm{d}q) \right] \boldsymbol{\mu}(\mathrm{d}\boldsymbol{u}).$$
<sup>26/53</sup>

Kadane and Wasilkowski (1985) show that *e*<sub>AC</sub>-minimisers are *deterministic* decision rules B, and the minimiser Y is "optimal information" for this task. But what if we restrict to *Bayesian* β?

## Theorem (AC = BPN for quadratic loss; Cockayne et al., 2019)

For a quadratic loss  $L(q, q') := ||q - q'||_{\mathcal{Q}}^2$  on a Hilbert space  $\mathcal{Q}$ , optimal information operators Y for BPNM and AC coincide (though the minimal values may differ).

## Theorem (AC $\neq$ BPN in general; Oates et al., 2019b)

If *U* can be partitioned into three sets of positive probability, then there exists a choice of Q and L so that optimal information operators Y for BPNM and AC differ.



DISINTEGRATION: EXACT AND NUMERICAL

#### DEFINING THE POSTERIOR

• The posterior  $\mu^{y}$  is subtle to define precisely, since *heuristically* it is given by

 $\mu^{\mathbf{y}}(\mathbf{d}\mathbf{u}) \propto \mathbb{I}[\mathbf{Y}(\mathbf{u}) = \mathbf{y}] \, \mu(\mathbf{d}\mathbf{u})$ 

- We have a 0-1 likelihood, and moreover the likelihood is zero  $\mu$ -a.e.!
  - The posterior  $\mu^{y}$  is singular w.r.t.  $\mu$ , supported on the null event  $[Y = y] \subseteq \mathcal{U}$ .
  - Why? Numerical analysts usually think of function evaluations as noiseless, in contrast to the noisy observations that are typical in statistics.
  - E.g. what is the prior probability that a Brownian path interpolates given data?
- We cannot even express Bayes' formula in the form favoured by Stuart (2010),

$$\frac{\mathrm{d}\mu^{y}}{\mathrm{d}\mu}(\boldsymbol{u}) = \frac{\mathbb{I}[Y(\boldsymbol{u}) = y]}{Z(y)}.$$

because Z(y) = 0.

• One way to consistently condition on events of measure zero is to define the conditioning operation in terms of **disintegration**.

## DISINTEGRATION I

# Definition (Disintegration)

A disintegration of  $\mu \in \mathcal{P}_{\mathcal{U}}$  w.r.t.  $\forall : \mathcal{U} \to \mathcal{Y}$  is a map  $\mathcal{Y} \to \mathcal{P}_{\mathcal{U}}$ ,  $y \mapsto \mu^{y}$ , such that

• (support)  $\mu^{y}(\{u \in \mathcal{U} \mid Y(u) = y\}) = 1$  for  $Y_{\#}\mu$ -almost all  $y \in \mathcal{Y}$ ;

and, for each measurable  $f \colon \mathcal{U} \to [0,\infty)$ ,

$$(f = \mathbb{I}_E, E \subseteq \mathcal{U} \text{ will do})$$

- (measurability)  $y \mapsto \int_{\mathcal{U}} f(u) \mu^{y}(du)$  is measurable
- (conditioning/reconstruction/law of total probability)

$$\int_{\mathcal{U}} f(\mathbf{u}) \, \mu(\mathrm{d}\mathbf{u}) = \int_{\mathcal{Y}} \left[ \int_{\mathcal{U}} f(\mathbf{u}) \, \mu^{\mathbf{y}}(\mathrm{d}\mathbf{u}) \right] (\mathsf{Y}_{\#}\mu)(\mathrm{d}\mathbf{y}).$$

(Closely related concept: a regular conditional probability is basically the same thing, but in a different coordinate system.) 29/53

## DISINTEGRATION II

## Theorem (Disintegration theorem (Chang and Pollard, 1997, Thm. 1))

Let  $\mathcal{U}$  be a metric space and let  $\mu \in \mathcal{P}_{\mathcal{U}}$  be inner regular. If the Borel  $\sigma$ -algebra on  $\mathcal{U}$  is countably generated and contains all singletons  $\{y\}$  for  $y \in \mathcal{Y}$ , then there is an essentially unique disintegration  $\{\mu^{\mathcal{Y}}\}_{\mathcal{Y} \in \mathcal{Y}}$  of  $\mu$  with respect to Y. (If  $\{\nu^{\mathcal{Y}}\}_{\mathcal{Y} \in \mathcal{Y}}$  is another such disintegration, then  $\{y \in \mathcal{Y} \mid \mu^{\mathcal{Y}} \neq \nu^{\mathcal{Y}}\}$  is an Y<sub>#</sub> $\mu$ -null set.)

- The familiar conditional densities for a probability density on ℝ<sup>n</sup> conditioned on a "nice" subset such as a lower-dimensional submanifold M ⊂ ℝ<sup>n</sup> are disintegrations.
- In particular, the familiar Woodbury formula for the conditioning of Gaussian measures subject to linear constraints is a disintegration (Owhadi and Scovel, 2015).
- But, in general, disintegrations cannot be computed exactly we have to work approximately.

### NUMERICAL DISINTEGRATION I

- The exact disintegration " $\mu^{y}(du) \propto \mathbb{I}[Y(u) = y] \mu(du)$ " can be accessed numerically via relaxation, with approximation guarantees provided  $y \mapsto \mu^{y}$  is "nice", e.g.  $Y_{\sharp}\mu \in \mathcal{P}_{\mathcal{V}}$  has a smooth Lebesgue density.
- Consider relaxed posterior  $\mu_{\delta}^{\mathcal{Y}}(\mathrm{d} u) \propto \phi(\|Y(u) y\|_{\mathcal{Y}}/\delta) \mu(\mathrm{d} u)$  with  $0 < \delta \ll 1$ .
  - Essentially any  $\phi \colon [0,\infty) \to [0,1]$  tending continuously to 1 at 0 and decaying quickly enough to 0 at  $\infty$  will do.
  - E.g.  $\phi(r) \coloneqq \mathbb{I}[r < 1] \text{ or } \phi(r) \coloneqq \exp(-r^2).$
- The integral probability metric on  $\mathcal{P}_{\mathcal{U}}$  associated to a normed space  $\mathcal{F}$  of test functions  $f: \mathcal{U} \to \mathbb{R}$  is

$$d_{\mathcal{F}}(\mu,\nu) \coloneqq \sup \big\{ |\mu(f) - \nu(f)| \big| ||f||_{\mathcal{F}} \le 1 \big\}.$$

- $\mathcal{F} =$  bounded continuous functions with uniform norm  $\leftrightarrow$  total variation.
- $\mathcal{F} = \text{bounded Lipschitz continuous functions with Lipschitz norm} \leftrightarrow \text{Wasserstein}$ .
- $\mathcal{F} = \mathsf{RKHS}$  of functions  $\leftrightarrow$  maximum mean discrepancy.

  $\mu_{\delta}^{\mathcal{Y}}(\mathrm{d} u) \propto \phi(\|\mathsf{Y}(u) - y\|_{\mathcal{Y}}/\delta) \,\mu(\mathrm{d} u)$ 

## Theorem (Cockayne et al., 2019, Theorem 4.4)

For any normed space  $\mathcal{F}$ , if  $y \mapsto \mu^{y}$  is  $\gamma$ -Hölder from  $(\mathcal{Y}, \|\cdot\|_{\mathcal{Y}})$  into  $(\mathcal{P}_{\mathcal{U}}, d_{\mathcal{F}})$ , then so too is the approximation  $\mu_{\delta}^{y} \approx \mu^{y}$  as a function of  $\delta$ . That is,

$$\begin{aligned} & d_{\mathcal{F}}(\mu^{y}, \mu^{y'}) \leq C \cdot \|y - y'\|^{\gamma} & \text{for } y, y' \in \mathcal{Y} \\ \implies & d_{\mathcal{F}}(\mu^{y}, \mu^{y}_{\delta}) \leq C \cdot C_{\phi} \cdot \delta^{\gamma} & \text{for } Y_{\sharp}\mu\text{-almost all } y \in \mathcal{Y}. \end{aligned}$$

Open question: when does the hypothesis, a quantitative version of the Tjur property (Tjur, 1980), actually hold? (Fixed y and free y' is easy; both y and y' free is hard.)  $_{32/53}$ 

A simple boundary value problem with multiple solutions:



**Figure 2:** The two solutions of Painlevé's first transcendental and their spectra in the orthonormal Chebyshev polynomial basis over [0, 10].

A simple boundary value problem with multiple solutions:





**Figure 2:** The two solutions of Painlevé's first transcendental and their spectra in the orthonormal Chebyshev polynomial basis over [0, 10].

# Example: Painlevé's first transcendental ii

- Try a Gaussian prior on *u*; we obtain qualitatively similar results for a heavy-tailed Cauchy prior (Sullivan, 2017).
- Parallel tempered pCN-MCMC with 100  $\delta$ -values log-spaced from  $\delta = 10$  to  $\delta = 10^{-4}$  and  $10^8$  iterations recovers both solutions in approximately the same proportions as the posterior densities at the two exact solutions.



# Example: Painlevé's first transcendental ii

- Try a Gaussian prior on *u*; we obtain qualitatively similar results for a heavy-tailed Cauchy prior (Sullivan, 2017).
- Parallel tempered pCN-MCMC with 100  $\delta$ -values log-spaced from  $\delta = 10$  to  $\delta = 10^{-4}$  and  $10^8$  iterations recovers both solutions in approximately the same proportions as the posterior densities at the two exact solutions.



# Example: Painlevé's first transcendental ii

- Try a Gaussian prior on *u*; we obtain qualitatively similar results for a heavy-tailed Cauchy prior (Sullivan, 2017).
- Parallel tempered pCN-MCMC with 100  $\delta$ -values log-spaced from  $\delta = 10$  to  $\delta = 10^{-4}$  and  $10^8$  iterations recovers both solutions in approximately the same proportions as the posterior densities at the two exact solutions.



What ways do we have around the cost of MCMC for BPNM?

- Of course, we have explicit conditioning in the (unimodal!) linear Gaussian case, e.g. the probabilistic meshless solver of Cockayne et al. (2016, 2017) for elliptic PDE. For even mildly smooth solution objects, the "screening effect" enables near-linear computational complexity (Schäfer et al., 2017).
- SMC seems to reliably but transiently detect the existence of multiple solutions but then suffers extinction problems.
- High-order quadrature (QMC) (e.g. Dick et al., 2014) and Laplace approximations for highly-concentrated posteriors (Schillings et al., 2019).
- Kernel and conditional mean embedding (Berlinet and Thomas-Agnan, 2004; Smola et al., 2007; ...; Klebanov et al., 2019) of probability distributions into RKHSs, with applications to nonlinear statistics (Nava-Yazdani et al., 2020).

COHERENT PIPELINES OF PNMS, AND BAYESIAN INVERSE PROBLEMS

#### **COMPUTATIONAL PIPELINES**

- Numerical methods usually form part of pipelines.
- Prime example: a PDE solve is a forward model in an inverse problem.
- Motivation for PNMs in the context of Bayesian inverse problems:

Make the forward and inverse problem speak the same statistical language!



- We can compose PNMs in series, e.g.  $\beta_2(\beta_1(\mu, y_1), y_2)$  is formally  $\beta(\mu, (y_1, y_2))$ ... although figuring out what the spaces  $\mathcal{U}_i$ ,  $\mathcal{Y}_i$  and operators  $Y_i$  etc. are is a headache!
- A graphical approach is both more intuitive and useful for analysis.

## PIPELINE EXAMPLE: SPLIT INTEGRATION



- Integrate a function over [0,1] in two steps using nodes  $0 \le t_0 < \cdots < t_{m-1} < \frac{1}{2}$ ,  $t_m = \frac{1}{2}$ , and  $\frac{1}{2} < t_{m+1} < \cdots < t_{2m} \le 1$ .
- For example, the two nodal sets are very large, and so two are handled by two different processors with non-shared memory.
- A third processor handles the (easy!) task of aggregating the two estimates of the two integrals  $\int_0^{1/2} u(t) dt$  and  $\int_{1/2}^1 u(t) dt$  into an estimate of  $\int_0^1 u(t) dt$ .

## COHERENCE I

- We compose PNMs in a graphical way by allowing input information nodes (□) to feed into method nodes (■), which in turn output new information.
- N.B. one should at first think of having *deterministic* data at the left-most □ nodes, then *random variables* as outputs, *realisations* of which get fed into the next ■.



### COHERENCE I

- We compose PNMs in a graphical way by allowing input information nodes (□) to feed into method nodes (■), which in turn output new information.
- N.B. one should at first think of having *deterministic* data at the left-most □ nodes, then *random variables* as outputs, *realisations* of which get fed into the next ■.



- We define the corresponding dependency graph by replacing each → → → by → →, and number the □ vertices in an increasing fashion.
- The independence properties of the random variables at each node are crucial. 38/53



## Definition

A prior  $\mu$  and dependency graph are **coherent** if, when the "leaf" input nodes are  $Y_{\sharp}\mu$ -distributed and the remaining nodes are  $\beta(\mu, \text{parents})$ -distributed, each node is conditionally independent of all older non-parent nodes given its direct parents:

 $Y_k \perp \downarrow Y_{\{1,\dots,k-1\}\setminus \text{parents}(k)} \mid Y_{\text{parents}(k)}$ .

(Generalises the Markov condition for directed acyclic graphs of Lauritzen (1991).)

#### COHERENCE II



### Definition

A prior  $\mu$  and dependency graph are **coherent** if, when the "leaf" input nodes are  $Y_{\sharp}\mu$ -distributed and the remaining nodes are  $\beta(\mu, \text{parents})$ -distributed, each node is conditionally independent of all older non-parent nodes given its direct parents:

 $Y_k \perp \downarrow Y_{\{1,\ldots,k-1\}\setminus \text{parents}(k)} \mid Y_{\text{parents}(k)}$ .

(Generalises the Markov condition for directed acyclic graphs of Lauritzen (1991).)

#### COHERENCE II



### Definition

A prior  $\mu$  and dependency graph are **coherent** if, when the "leaf" input nodes are  $Y_{\sharp}\mu$ -distributed and the remaining nodes are  $\beta(\mu, \text{parents})$ -distributed, each node is conditionally independent of all older non-parent nodes given its direct parents:

 $Y_k \perp \downarrow Y_{\{1,\ldots,k-1\}\setminus \text{parents}(k)} \mid Y_{\text{parents}(k)}$ .

(Generalises the Markov condition for directed acyclic graphs of Lauritzen (1991).)

## **COHERENCY THEOREM**

## Theorem (Cockayne et al., 2019, Theorem 5.9)

If a pipeline of PNMs is coherent and the component PNMs are all Bayesian, then the pipeline is Bayesian as a whole, i.e. is equivalent to the Bayesian pipeline

data at leaves 
$$\beta(\mu, \cdot)$$
 final output

## **COHERENCY THEOREM**

## Theorem (Cockayne et al., 2019, Theorem 5.9)

If a pipeline of PNMs is coherent and the component PNMs are all Bayesian, then the pipeline is Bayesian as a whole, i.e. is equivalent to the Bayesian pipeline



- Redundant structure in the pipeline (recycled information) will break coherence, and hence Bayesianity of the pipeline.
- In principle, coherence and hence being Bayesian depend upon the prior.
- This should not be surprising as a loose analogy, one doesn't expect the trapezoidal rule to be a good way to integrate very smooth functions.
- Of course, non-Bayesian setups can be good in other ways (Jacob et al., 2017; Lie et al., 2018).

## Split integration: Coherence



- Integrate a function over [0,1] in two steps using nodes  $0 \le t_0 < \cdots < t_{m-1} < \frac{1}{2}$ ,  $t_m = \frac{1}{2}$ , and  $\frac{1}{2} < t_{m+1} < \cdots < t_{2m} \le 1$ .
- Is  $(\int_{1/2}^1 u(t) dt)$  independent of  $(u(t_0), \dots, u(t_{m-1}))$  given  $(u(t_m), \dots, u(t_{2m}))$ ?

## Split integration: Coherence



- Integrate a function over [0, 1] in two steps using nodes  $0 \le t_0 < \cdots < t_{m-1} < \frac{1}{2}$ ,  $t_m = \frac{1}{2}$ , and  $\frac{1}{2} < t_{m+1} < \cdots < t_{2m} \le 1$ .
- Is  $(\int_{1/2}^1 u(t) dt)$  independent of  $(u(t_0), \dots, u(t_{m-1}))$  given  $(u(t_m), \dots, u(t_{2m}))$ ?
- For a Brownian motion prior on *u*, yes. For an integrated BM prior *u*, no.

## Split integration: Coherence



- Integrate a function over [0, 1] in two steps using nodes  $0 \le t_0 < \cdots < t_{m-1} < \frac{1}{2}$ ,  $t_m = \frac{1}{2}$ , and  $\frac{1}{2} < t_{m+1} < \cdots < t_{2m} \le 1$ .
- Is  $(\int_{1/2}^1 u(t) dt)$  independent of  $(u(t_0), \dots, u(t_{m-1}))$  given  $(u(t_m), \dots, u(t_{2m}))$ ?
- For a Brownian motion prior on *u*, yes. For an integrated BM prior *u*, no.
- So how do we elicit an appropriate prior that respects the problem's structure, an in particular incorporates a-priori knowledge from numerical analysis?

# SHORT PIPELINES: (RANDOMISED) BAYESIAN INVERSE PROBLEMS I

Bayesian inverse problems are good examples of (short) pipelines of PNMs:

$$d \longrightarrow \beta_1(\mu, \cdot) \longrightarrow L(\cdot | d) \longrightarrow \beta_2(\mu, \cdot) \longrightarrow \theta$$

- A BIP is essentially a two-stage computational pipeline in which
  - β<sub>1</sub> converts data *d* into the likelihood function for parameters θ, and hence incorporates any forward model such as an O/PDE solver
  - $\beta_2$  converts the prior on  $\theta$  and the likelihood into a joint distribution for  $(\theta, d)$ , then conditions upon the actual observation it returns something in  $\mathcal{P}_{\Theta}$ .
- Conventionally, β<sub>1</sub> is a function from D into R<sup>Θ</sup>; a bona fide PNM would return a non-trivial probability distribution in P<sub>R<sup>Θ</sup></sub>, i.e. a randomised likelihood.

#### Lemma

Under the mild assumption that any randomisation in the forward model is independent of the prior on  $\theta$ , a BIP pipeline is always coherent.

- Even if some of the method nodes are non-Bayesian, we can assess how close the overall pipeline is to the Bayesian "ideal". In fact, some non-Bayesianity of component methods can confer robustness on the pipeline as a whole (Jacob et al., 2017; Owhadi et al., 2015).
- Lie et al. (2018) analyse, in terms of L<sup>p</sup> and Hellinger convergence, how the stochastic variability in the forward model / likelihood propagates to the (randomised or marginal) Bayesian posterior on θ.
- Alternative approach: assess sufficiency of forward solver accuracy for BIP purposes using Bayes factors (Capistrán et al., 2016; Christen et al., 2017).

APPLICATIONS

# Example: Hydrocyclones (Oates et al., 2019a)

- Hydrocyclones are used in industry as an alternative to centrifuges or filtration systems to separate fluids of different densities or particulate matter from a fluid.
- Monitoring is an essential control component, but usually cannot be achieved visually: Gutiérrez et al. (2000) propose electrical impedance tomography as an alternative.
- EIT is an indirect imaging technique in which the conductivity field in the interior — which correlates with many material properties of interest — is inferred from current and voltage boundary conditions.
- In its Bayesian formulation, this is a well-posed inverse problem (Dunlop and Stuart, 2016a,b) closely related to Calderón's problem (Uhlmann, 2009).



## COMPLETE ELECTRODE MODEL (CHENG ET AL., 1989; SOMERSALO ET AL., 1992)

The interior conductivity field  $\sigma$  and electrical potential field v and the applied boundary currents  $I_i$ , measured voltages  $V_i$ , and known contact impedances  $\zeta_i$  are related by

$$-\nabla \cdot \boldsymbol{\sigma}(\mathbf{x}) \nabla \mathbf{v}(\mathbf{x}) = 0 \qquad \mathbf{x} \in D; \qquad \int_{E_i} \boldsymbol{\sigma}(\mathbf{x}) \frac{\partial \mathbf{v}(\mathbf{x})}{\partial \hat{n}} \, \mathrm{d}\mathbf{u} = \mathbf{I}_i \qquad \mathbf{x} \in E_i, i = 1, \dots, m;$$
$$\mathbf{v}(\mathbf{x}) + \zeta_i \boldsymbol{\sigma}(\mathbf{x}) \frac{\partial \mathbf{v}(\mathbf{x})}{\partial \hat{n}} = \mathbf{V}_i \qquad \mathbf{x} \in E_i; \qquad \boldsymbol{\sigma}(\mathbf{x}) \frac{\partial \mathbf{v}(\mathbf{x})}{\partial \hat{n}} = 0 \qquad \mathbf{x} \in \partial D \setminus \bigcup_{i=1}^m E_i.$$

Furthermore, we consider a vector of such models, with multiple current stimulation patterns, at multiple points in time, for a time-dependent field  $\sigma(t, x)$ .



## **EIT FORWARD PROBLEM**

- Sampling from the posterior(s) requires repeatedly solving the forward PDE.
- We use the probabilistic meshless method (PMM) of Cockayne et al. (2016, 2017):
  - a Gaussian process extension of symmetric collocation;
  - a Bayesian PNM for a Gaussian prior and linear elliptic PDEs of this type.
- PMM allows us to:
  - account for uncertainty arising from the PDE having no explicit solution;
  - use coarser discretisations of the PDE to solve the problem faster while still providing meaningful UQ for the inverse problem, cf. Capistrán et al. (2016); Christen et al. (2017).



**Figure 3:** PMM imposes the PDE at  $n_A$  interior nodes and BCs at  $n_B$  boundary nodes.

46/53

• For the inverse problem we use a Karhunen–Loève series prior:

$$\log \sigma(t, x; \omega) = \sum_{k=1}^{\infty} k^{-\alpha} \psi_k(t; \omega) \phi_k(x),$$

with the  $\psi_k$  being a-priori independent Brownian motions in t.

- Like Dunlop and Stuart (2016a), we assume additive Gaussian observational noise with variance γ<sup>2</sup> > 0, independently on each E<sub>i</sub>.
- We adopt a filtering formulation, inferring  $\sigma(t_i, \cdot; \cdot)$  sequentially.
- Within each data assimilation step, the Bayesian update is performed by SMC with
   *P* ∈ ℕ weighted particles and a pCN transition kernel (which uses point evaluations
   of *σ* directly and avoids truncation of the KL expansion).
- Real-world data obtained at 49 regular time intervals: rapid injection between frames 10 and 11, followed by diffusion and rotation of the liquids.

### EIT STATIC RECOVERY I

**Figure 4:** A small number  $n_A + n_B = 71$  of collocation points was used to discretise the PDE. but the uncertainty due to discretisation was not modelled. The reference posterior distribution over the coefficients  $\psi_{k}$  is plotted (grey) and compared to the approximation to the posterior obtained when the PDE is discretised and the discretisation error is not modelled (blue, 'Non-PN'). The approximate posterior is highly biased.



#### EIT STATIC RECOVERY II



**Figure 5:** Posterior means and standard-deviations for the recovered conductivity field at t = 14. The first column shows the reference solution, obtained using symmetric collocation with a large number of collocation points. The remaining columns show the recovered field when PMM is used with  $n_A + n_B$  collocation points. 49/53
**Figure 6:** Posterior distribution over the coefficients  $\psi_k$  at the final time. A small number  $n_A + n_B = 71$  of collocation points was used to discretise the PDE. The reference posterior distribution over the coefficients  $\psi_k$  is plotted (grey) and compared to the approximation to the posterior obtained when discretisation of the PDE is not modelled (blue, 'Non-PN') and modelled (orange, 'PN').



### **EIT COMMENTS**

- Typically PDE discretisation error in BIPs is ignored, or its contribution is bounded through detailed numerical analysis. Theoretical bounds are difficult in the temporal setting due to propagation and accumulation of errors
- As a modelling choice, the PN approach eases these difficulties. As with the Painlevé example, this is a statistically correct implementation of the assumptions, but it is (at present) costly.
- Furthermore, Markov temporal evolution of the conductivity field was assumed; this is likely incorrect, since time derivatives of this field will vary continuously. Even a-priori knowledge about the spin direction is neglected at present.
- Again, we see a need for priors that are 'physically reasonable' and statistically/computationally appropriate.

P

CLOSING REMARKS

Aside from obvious improvements to computational cost and applications...

- The analysis of PNM pipelines uses directed acyclic graphs, i.e. excludes adaptivity. There are some recent advances in this direction for quadrature (Jagadeeswaran and Hickernell, 2019); new project starting soon for PDEs.
- Related question: automatisation of prior specification? Numerical analysis may be an ideal playground for empirical Bayesian methods.
- Questions of geometric measure theory continuity of disintegrations?
- Connections to category-theoretic interpretations of probability and probabilistic functional programming languages (Giry, 1982; ...; Fritz, 2019; Parzygnat, 2020)?
- If classical NMs correspond to maximum likelihood estimators, then point estimators for BPNMs are MAP estimators — connections to modern theory non-parametric MAP estimators (Dashti et al., 2013; Lie and Sullivan, 2018), and modern optimisation methods for finding them (DNNs?)?

#### **CLOSING REMARKS**

- $\,$  Numerical methods can be characterised in a Bayesian fashion, distinct from ACA.  $\checkmark$
- BPNMs can be composed into pipelines, e.g. for inverse problems.
- Bayes' rule as disintegration  $\rightarrow$  (expensive!) numerical implementation.
  - Lots of room to improve computational cost and bias.
  - Cost-accuracy tradeoff when leaving the "Bayesian gold standard".
- How to choose/design an appropriate (numerically-analytically right) prior?
- Foundations: Cockayne et al. (2019) arXiv:1702.03673
  Industrial applications: Cockayne et al. (2019, §6.3) & Oates et al. (2019a) arXiv:1707.06107
  History: Oates and Sullivan (2019) arXiv:1901.04457
  Optimality: Oates et al. (2019b) arXiv:1901.04326
  BIPs: Lie et al. (2018) arXiv:1712.05717

🖌 | 🗙

P

#### **CLOSING REMARKS**

- $\,$  Numerical methods can be characterised in a Bayesian fashion, distinct from ACA.  $\checkmark$
- BPNMs can be composed into pipelines, e.g. for inverse problems.
- Bayes' rule as disintegration  $\rightarrow$  (expensive!) numerical implementation.
  - Lots of room to improve computational cost and bias.
  - Cost-accuracy tradeoff when leaving the "Bayesian gold standard".
- How to choose/design an appropriate (numerically-analytically right) prior?
- Foundations:
- Industrial applications:
- History:
- Optimality:
- BIPs:

Cockayne et al. (2019) arXiv:1702.03673 Cockayne et al. (2019, §6.3) & Oates et al. (2019a) arXiv:1707.06107 Oates and Sullivan (2019) arXiv:1901.04457 Oates et al. (2019b) arXiv:1901.04326 Lie et al. (2018) arXiv:1712.05717

🖌 | 🗙

P

#### **REFERENCES** I

- J. L. Barlow and E. H. Bareiss. Probabilistic error analysis of Gaussian elimination in floating point and logarithmic arithmetic. *Computing*, 34(4):349–364, 1985. doi:10.1007/BF02251834.
- J. O. Berger. *Statistical Decision Theory and Bayesian Analysis*. Springer Series in Statistics. Springer-Verlag, New York, second edition, 1985. doi:10.1007/978-1-4757-4286-2.
- A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer Academic Publishers, 2004.
- M. A. Capistrán, J. A. Christen, and S. Donnet. Bayesian analysis of ODEs: solver optimal accuracy and Bayes factors. *SIAM/ASA J. Uncertain. Quantif.*, 4(1):829–849, 2016. doi:10.1137/140976777.
- J. T. Chang and D. Pollard. Conditioning as disintegration. *Statist. Neerlandica*, 51(3):287–317, 1997. doi:10.1111/1467-9574.00056.
- F. Chatelin and M.-C. Brunet. A probabilistic round-off error propagation model. Application to the eigenvalue problem. In *Reliable numerical computation*, Oxford Sci. Publ., pages 139–160. Oxford Univ. Press, New York, 1990.
- K.-S. Cheng, D. Isaacson, J. C. Newell, and D. G. Gisser. Electrode models for electric current computed tomography. *IEEE Trans. Biomed. Eng.*, 36(9), 1989. doi:10.1109/10.35300.

#### **REFERENCES II**

- J. A. Christen, M. A. Capistrán, M. L. Daza-Torres, H. Flores-Argüedas, and J. C. Montesinos-López. Posterior distribution existence and error control in Banach spaces in the Bayesian approach to UQ in inverse problems, 2017. arXiv:1712.03299.
- J. Cockayne, C. J. Oates, T. J. Sullivan, and M. Girolami. Probabilistic meshless methods for partial differential equations and Bayesian inverse problems, 2016. arXiv:1605.07811.
- J. Cockayne, C. J. Oates, T. J. Sullivan, and M. Girolami. Probabilistic numerical methods for PDE-constrained Bayesian inverse problems. In G. Verdoolaege, editor, *Proceedings of the 36<sup>th</sup> International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, volume 1853 of *AIP Conference Proceedings*, pages 060001–1–060001–8, 2017. doi:10.1063/1.4985359.
- J. Cockayne, C. J. Oates, T. J. Sullivan, and M. Girolami. Bayesian probabilistic numerical methods. SIAM Rev., 61(4):756–789, 2019. doi:10.1137/17M1139357.
- P. R. Conrad, M. Girolami, S. Särkkä, A. M. Stuart, and K. C. Zygalakis. Statistical analysis of differential equations: introducing probability measures on numerical solutions. *Stat. Comput.*, 27(4), 2016. doi:10.1007/s11222-016-9671-0.
- M. Dashti, K. J. H. Law, A. M. Stuart, and J. Voss. MAP estimators and their consistency in Bayesian nonparametric inverse problems. *Inv. Probl.*, 29(9):095017, 27, 2013. doi:10.1088/0266-5611/29/9/095017.

#### **REFERENCES III**

- P. Diaconis. Bayesian numerical analysis. In *Statistical Decision Theory and Related Topics, IV, Vol. 1 (West Lafayette, Ind., 1986)*, pages 163–175. Springer, New York, 1988.
- J. Dick, F. Y. Kuo, Q. T. Le Gia, D. Nuyens, and C. Schwab. Higher order QMC Petrov-Galerkin discretization for affine parametric operator equations with random field inputs. *SIAM J. Numer. Anal.*, 52(6):2676–2702, 2014. doi:10.1137/130943984.
- M. M. Dunlop and A. M. Stuart. The Bayesian formulation of EIT: analysis and algorithms. *Inv. Probl. Imaging*, 10(4): 1007–1036, 2016a. doi:10.3934/ipi.2016030.
- M. M. Dunlop and A. M. Stuart. MAP estimators for piecewise continuous inversion. *Inv. Probl.*, 32(10):105003, 50, 2016b. doi:10.1088/0266-5611/32/10/105003.
- T. Fritz. A synthetic approach to Markov kernels, conditional independence and theorems on sufficient statistics, 2019. arXiv:1908.07021.
- M. Giry. A categorical approach to probability theory. In *Categorical aspects of topology and analysis (Ottawa, Ont., 1980),* volume 915 of *Lecture Notes in Math.,* pages 68–85. Springer, Berlin-New York, 1982.
- J. Gutiérrez, T. Dyakowski, M. Beck, and R. Williams. Using electrical impedance tomography for controlling hydrocyclone underflow discharge. *Powder Tech.*, 108(2):180–184, 2000. doi:10.1016/S0032-5910(99)00218-1.
- P. Henrici. Discrete Variable Methods in Ordinary Differential Equations. John Wiley & Sons, Inc., New York-London, 1962.

#### **REFERENCES IV**

- T. E. Hull and J. R. Swenson. Tests of probabilistic models for the propagation of roundoff errors. *Comm. ACM*, 9:108–113, 1966. doi:10.1145/365170.365212.
- P. E. Jacob, L. M. Murray, C. C. Holmes, and C. P. Robert. Better together? Statistical learning in models made of modules, 2017. arXiv:1708.08719.
- R. Jagadeeswaran and F. J. Hickernell. Fast automatic Bayesian cubature using lattice sampling. *Stat. Comp.*, 29:1215–1229, 2019. doi:0.1007/s11222-019-09895-9.
- J. B. Kadane and G. W. Wasilkowski. Average case ε-complexity in computer science. A Bayesian view. In *Bayesian Statistics*, 2 (Valencia, 1983), pages 361–374. North-Holland, Amsterdam, 1985.
- T. Karvonen, C. J. Oates, and S. Särkkä. A Bayes–Sard cubature method. In 32nd Conference on Neural Information Processing Systems (NeurIPS 2018), 2018. papers.nips.cc/paper/7829-a-bayes-sard-cubature-method. Kazan Federal University. kpfu.ru/portal/docs/F 261937733/suldin2.jpg. Accessed December 2018.
- G. S. Kimeldorf and G. Wahba. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *Ann. Math. Stat.*, 41:495–502, 1970a. doi:10.1214/aoms/1177697089.
- G. S. Kimeldorf and G. Wahba. Spline functions and stochastic processes. *Sankhyā Ser. A*, 32:173–180, 1970b. www.jstor.org/stable/25049652.

I. Klebanov, I. Schuster, and T. J. Sullivan. A rigorous theory of conditional mean embeddings, 2019. arXiv:1912.00671.

#### **REFERENCES V**

- A. N. Kolmogorov. Über die beste Annäherung von Funktionen einer gegebenen Funktionenklasse. Ann. of Math. (2), 37(1): 107–110, 1936. doi:10.2307/1968691.
- J. Kuelbs, F. M. Larkin, and J. A. Williamson. Weak probability distributions on reproducing kernel Hilbert spaces. *Rocky Mountain J. Math.*, 2(3):369–378, 1972. doi:10.1216/RMJ-1972-2-3-369.
- F. M. Larkin. Estimation of a non-negative function. BIT Num. Math., 9(1):30-52, 1969. doi:10.1007/BF01933537.
- F. M. Larkin. Optimal approximation in Hilbert spaces with reproducing kernel functions. *Math. Comp.*, 24:911–921, 1970. doi:10.2307/2004625.
- F. M. Larkin. Gaussian measure in Hilbert space and applications in numerical analysis. *Rocky Mountain J. Math.*, 2(3): 379–421, 1972. doi:10.1216/RMJ-1972-2-3-379.
- F. M. Larkin. Probabilistic error estimates in spline interpolation and quadrature. In *Information Processing 74 (Proc. IFIP Congress, Stockholm, 1974)*, pages 605–609, Amsterdam, 1974. North-Holland.
- F. M. Larkin. A modification of the secant rule derived from a maximum likelihood principle. *BIT*, 19(2):214–222, 1979a. doi:10.1007/BF01930851.
- F. M. Larkin. Probabilistic estimation of poles or zeros of functions. J. Approx. Theory, 27(4):355–371, 1979b. doi:10.1016/0021-9045(79)90124-2.

#### **REFERENCES VI**

- F. M. Larkin. Bayesian estimation of zeros of analytic functions. Technical report, Queen's University of Kingston. Department of Computing and Information Science., 1979c.
- F. M. Larkin, C. E. Brown, K. W. Morton, and P. Bond. Worth a thousand words, 1967. www.amara.org/en/videos/7De21CeNlz8b/info/worth-a-thousand-words-1967/.
- S. Lauritzen. Graphical Models. Oxford University Press, 1991.
- H. C. Lie and T. J. Sullivan. Equivalence of weak and strong modes of measures on topological vector spaces. *Inv. Probl.*, 34 (11):115013, 22, 2018. doi:10.1088/1361-6420/aadef2.
- H. C. Lie, T. J. Sullivan, and A. L. Teckentrup. Random forward models and log-likelihoods in Bayesian inverse problems. *SIAM/ASA J. Uncertain. Quantif.*, 6(4):1600–1629, 2018. doi:10.1137/18M1166523.
- H. C. Lie, A. M. Stuart, and T. J. Sullivan. Strong convergence rates of probabilistic integrators for ordinary differential equations. *Stat. Comp.*, 29(6):1265–1283, 2019. doi:10.1007/s11222-019-09898-6.
- T. Minka. Deriving quadrature rules from Gaussian processes, 2000. www.microsoft.com/en-us/research/publication/deriving-quadrature-rules-gaussian-processes/.
- J. Močkus. On Bayesian methods for seeking the extremum. In *Optimization Techniques IFIP Technical Conference* Novosibirsk, July 1–7, 1974. Optimization Techniques 1974, volume 27 of Lecture Notes in Computer Science, pages 400–404. Springer, Berlin, Heidelberg, 1975. doi:10.1007/3-540-07165-2\_55.

#### **REFERENCES VII**

J. Močkus. On Bayesian methods for seeking the extremum and their application. In *Information Processing 77 (Proc. IFIP Congr., Toronto, Ont., 1977)*, pages 195–200. IFIP Congr. Ser., Vol. 7. North-Holland, Amsterdam, 1977.

- J. Močkus. Bayesian approach to global optimization, volume 37 of Mathematics and its Applications (Soviet Series). Kluwer Academic Publishers Group, Dordrecht, 1989. doi:10.1007/978-94-009-0909-0.
- E. Nava-Yazdani, H.-C. Hege, T. J. Sullivan, and C. von Tycowicz. Geodesic analysis in Kendall's shape space with epidemiological applications, 2020. To appear. arXiv:1906.11950.
- A. P. Norden, Y. I. Zabotin, L. D. Èskin, S. V. Grigor'ev, and E. A. Begovatov. Al'bert Valentinovich Sul'din (on the occasion of his fiftieth birthday). *Izv. Vysš. Učebn. Zaved. Mat.*, 12:3–5, 1978.
- E. Novak. Deterministic and Stochastic Error Bounds in Numerical Analysis, volume 1349 of Lecture Notes in Mathematics. Springer-Verlag, Berlin, 1988. doi:10.1007/BFb0079792.
- C. J. Oates and T. J. Sullivan. A modern retrospective on probabilistic numerics. *Stat. Comp.*, 29(6):1335–1351, 2019. doi:10.1007/s11222-019-09902-z.
- C. J. Oates, J. Cockayne, R. G. Aykroyd, and M. Girolami. Bayesian probabilistic numerical methods in time-dependent state estimation for industrial hydrocyclone equipment. *J. Amer. Stat. Assoc.*, pages 1–27, 2019a. doi:10.1080/01621459.2019.1574583.

#### **REFERENCES VIII**

- C. J. Oates, J. Cockayne, D. Prangle, T. J. Sullivan, and M. Girolami. Optimality criteria for probabilistic numerical methods. In F. J. Hickernell and P. Kritzer, editors, *Multivariate Algorithms and Information-Based Complexity*. Berlin/Boston: De Gruyter, 2019b. To appear. arXiv:1901.04326.
- A. O'Hagan. Monte Carlo is fundamentally unsound. Statistician, 36(2/3):247-249, 1987. doi:10.2307/2348519.
- A. O'Hagan. Bayes-Hermite quadrature. J. Stat. Plann. Inference, 29(3):245-260, 1991. doi:10.1016/0378-3758(91)90002-V.
- H. Owhadi and C. Scovel. Conditioning Gaussian measure on Hilbert space, 2015. arXiv:1506.04208.
- H. Owhadi, C. Scovel, and T. J. Sullivan. Brittleness of Bayesian inference under finite information in a continuous world. *Electron. J. Stat.*, 9(1):1–79, 2015. doi:10.1214/15-EJS989.
- E. Parzen. Statistical inference on time series by RKHS methods. Technical report, Stanford University of California, Department of Statistics, 1970.
- A. J. Parzygnat. Inverses, disintegrations, and Bayesian inversion in quantum Markov categories, 2020. arXiv:2001.08375. H. Poincaré. *Calcul des Probabilités*. Gauthier-Villars, second edition, 1912.
- M. Raissi, P. Perdikaris, and G. E. Karniadakis. Numerical Gaussian processes for time-dependent and nonlinear partial differential equations. *SIAM J. Sci. Comput.*, 40(1):A172–A198, 2018. doi:10.1137/17M1120762.
- C. E. Rasmussen and Z. Ghahramani. Bayesian Monte Carlo. In Advances in Neural Information Processing Systems 16, pages 505–512, 2003. papers.nips.cc/paper/2150-bayesian-monte-carlo.

#### **REFERENCES IX**

- K. Ritter. Average-Case Analysis of Numerical Problems, volume 1733 of Lecture Notes in Mathematics. Springer-Verlag, Berlin, 2000. doi:10.1007/BFb0103934.
- A. Sard. Best approximate integration formulas; best approximation formulas. *Amer. J. Math.*, 71:80–91, 1949. doi:10.2307/2372095.
- A. Sard. *Linear Approximation*. Number 9 in Mathematical Surveys. American Mathematical Society, Providence, RI, 1963. doi:10.1090/surv/009.
- F. Schäfer, T. J. Sullivan, and H. Owhadi. Compression, inversion, and approximate PCA of dense kernel matrices at near-linear computational complexity, 2017. arXiv:1706.02205.
- C. Schillings, B. Sprungk, and P. Wacker. On the convergence of the Laplace approximation and noise-level-robustness of Laplace-based Monte Carlo methods for Bayesian inverse problems, 2019. arXiv:1901.03958.
- M. Schober, D. K. Duvenaud, and P. Hennig. Probabilistic ODE solvers with Runge–Kutta means. In Advances in Neural Information Processing Systems 27, 2014.

papers.nips.cc/paper/5451-probabilistic-ode-solvers-with-runge-kutta-means.

J. Skilling. Bayesian solution of ordinary differential equations. In C. R. Smith, G. J. Erickson, and P. O. Neudorfer, editors, Maximum Entropy and Bayesian Methods, volume 50 of Fundamental Theories of Physics, pages 23–37. Springer, 1992. doi:10.1007/978-94-017-2219-3.

#### **REFERENCES** X

- S. Smale. On the efficiency of algorithms of analysis. *Bull. Amer. Math. Soc.* (N.S.), 13(2):87–121, 1985. doi:10.1090/S0273-0979-1985-15391-1.
- A. J. Smola, A. Gretton, L. Song, and B. Schölkopf. A Hilbert space embedding for distributions. In *Proceedings of the 18th International Conference on Algorithmic Learning Theory*, pages 13–31. Springer-Verlag, 2007.
- E. Somersalo, M. Cheney, and D. Isaacson. Existence and uniqueness for electrode models for electric current computed tomography. *SIAM J. Appl. Math.*, 52(4):1023–1040, 1992. doi:10.1137/0152060.
- A. M. Stuart. Inverse problems: a Bayesian perspective. Acta Numer., 19:451–559, 2010. doi:10.1017/S0962492910000061.
- A. V. Sul'din. Wiener measure and its applications to approximation methods. I. *Izv. Vysš. Učebn. Zaved. Mat.*, 6(13):145–158, 1959.
- A. V. Sul'din. Wiener measure and its applications to approximation methods. II. *Izv. Vysš. Učebn. Zaved. Mat.*, 5(18):165–179, 1960.
- A. V. Sul'din. On the distribution of the functional  $\int_0^1 x^2(t) dt$  where x(t) represents a certain Gaussian process. In *Kazan State Univ. Sci. Survey Conf. 1962 (Russian)*, pages 80–82. Izdat. Kazan. Univ., Kazan, 1963.
- T. J. Sullivan. Well-posed Bayesian inverse problems and heavy-tailed stable quasi-Banach space priors. *Inverse Probl. Imaging*, 11(5):857–874, 2017. doi:10.3934/ipi.2017040.

#### **REFERENCES XI**

- M. Tienari. A statistical model of roundoff error for varying length floating-point arithmetic. *Nordisk Tidskr. Informationsbehandling (BIT)*, 10:355–365, 1970. doi:10.1007/BF01934204.
- T. Tjur. *Probability Based on Radon Measures*. John Wiley & Sons, Ltd., Chichester, 1980. Wiley Series in Probability and Mathematical Statistics.
- J. F. Traub and H. Woźniakowsi. A General Theory of Optimal Algorithms. ACM Monograph Series. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London, 1980.
- J. F. Traub, G. W. Wasilkowski, and H. Woźniakowski. *Information, Uncertainty, Complexity*. Addison-Wesley Publishing Company, Advanced Book Program, Reading, MA, 1983.
- G. Uhlmann. Electrical impedance tomography and Calderón's problem. *Inv. Probl.*, 25(12):123011, 39, 2009. doi:10.1088/0266-5611/25/12/123011.
- J. von Neumann and H. H. Goldstine. Numerical inverting of matrices of high order. *Bull. Amer. Math. Soc.*, 53:1021–1099, 1947. doi:10.1090/S0002-9904-1947-08909-6.
- Y. I. Zabotin, N. K. Zamov, L. A. Aksent'ev, and T. N. Zemtseva. Al'bert Valentinovich Sul'din (obituary). *Izv. Vysš. Učebn. Zaved. Mat.*, 2(84), 1996.

## Theorem (AC $\neq$ BPN in general; Oates et al. (2019b))

If *U* can be partitioned into three sets of positive probability, then there exists a choice of QoI and loss so that optimal information for BPNM and AC differ.

# Example (AC $\neq$ BPN in general; Oates et al. (2019b))

Decide whether or not a card drawn fairly at random is  $\blacklozenge$ , incurring unit loss if you guess wrongly; can choose to be told whether the card is red (Y<sub>1</sub>) or is non- $\clubsuit$  (Y<sub>2</sub>).

$\mathcal{U} = \{ \clubsuit, \blacklozenge, \blacktriangledown, \bigstar \}$	$\mu = {\sf Unif}_{{m {\cal U}}}$	$\mathcal{Q} = \{0,1\} \subset \mathbb{R}$
$\mathcal{Y}_1 = \{0, 1\}$	$Y_1(\boldsymbol{u}) = \mathbb{I}[\boldsymbol{u} \in \{ \blacklozenge, \blacktriangledown \}]$	$Q(u) = \mathbb{I}[u = \blacklozenge]$
$\mathcal{Y}_2 = \{0, 1\}$	$Y_2(u) = \mathbb{I}[u \in \{\diamondsuit, \blacktriangledown, \clubsuit]]$	$L(q,q') = \mathbb{I}[q \neq q']$

Which information operator,  $Y_1$  or  $Y_2$ , is better? (Note that  $e_{WC}(Y_i, B) = 1$  for all deterministic *b*!)

$$U = \textcircled{\ } (Y_1, B) = \frac{1}{4} ( L(B(\blacksquare), 0) + L(B(\blacksquare), 1) + L(B(\blacksquare), 0) + L(B(\blacksquare), 0) )$$







<i>u</i> =	*	•		•		•	
$e_{AC}(Y_1,B) = \frac{1}{4} ($	$L(B(\blacksquare),0)$ +	$L(B(\blacksquare),1)$	+	$L(B(\blacksquare),0)$	+	$L(B(\blacksquare),0)$	)
$e_{AC}(Y_1,0) = \frac{1}{4} \big($	0 +	1	+	0	+	0	$) = \frac{1}{4}$
$e_{AC}(Y_1, id) = \frac{1}{4}($	0 +	0	+	1	+	0	$) = \frac{1}{4}$
$e_{AC}(Y_2,B) = \frac{1}{4}($	$L(B(\clubsuit),0)$ +	$L(B(\neg \clubsuit), 1)$	+	$L(B(\neg \clubsuit), 0)$	+	$L(B(\neg \clubsuit), 0)$	)
$e_{AC}(Y_2,0) = \frac{1}{4} ($	0 +	1	+	0	+	0	$) = \frac{1}{4}$

U =	*	•		•			
$e_{AC}(Y_1,B) = \frac{1}{4} ($	$L(B(\blacksquare),0)$ +	$L(B(\blacksquare),1)$	+	$L(B(\blacksquare),0)$	+	$L(B(\blacksquare),0)$	)
$e_{AC}(Y_1,0) = \frac{1}{4} ($	0 +	1	+	0	+	0	$)=\frac{1}{4}$
$e_{AC}(Y_1, id) = \frac{1}{4}($	0 +	0	+	1	+	0	$) = \frac{1}{4}$
$e_{AC}(Y_2,B) = \frac{1}{4}($	$L(B(\clubsuit),0)$ +	$L(B(\neg \clubsuit), 1)$	+	$L(B(\neg \clubsuit), 0)$	+	$L(B(\neg \clubsuit), 0)$	)
$e_{AC}(Y_2,0) = \frac{1}{4} ($	0 +	1	+	0	+	0	$)=\frac{1}{4}$
$e_{BPN}(Y_1) = \frac{1}{4} ($	$\mathbb{E}_{Q_{\sharp}\mu} \mathbf{I}(\cdot, 0) +$	$\mathbb{E}_{Q_{\sharp}\mu} L(\cdot, 1)$	+	$\mathbb{E}_{Q_{\sharp}\mu} L(\cdot, 0)$	+	$\mathbb{E}_{Q_{\sharp}\mu} L(\cdot, 0)$	)
$=\frac{1}{4}($	$(\frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 0) +$	$(\tfrac{1}{2} \cdot 0 + \tfrac{1}{2} \cdot 1)$	+	$(\frac{1}{2}\cdot 1 + \frac{1}{2}\cdot 0)$	+	$(\frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 0)$	$) = \frac{1}{4}$

#### ◀ Back

U =	*	•		•		<b>★</b>	
$e_{AC}(Y_1,B) = \frac{1}{4} \big($	$L(B(\blacksquare),0)$ +	$L(B(\blacksquare),1)$	+	$L(B(\blacksquare),0)$	+	$L(B(\blacksquare),0)$	)
$e_{AC}(Y_1,0) = \frac{1}{4} \big($	0 +	1	+	0	+	0	$) = \frac{1}{4}$
$e_{AC}(Y_1, id) = \frac{1}{4}($	0 +	0	+	1	+	0	$) = \frac{1}{4}$
$e_{AC}(Y_2,B) = \frac{1}{4}($	$L(B(\clubsuit),0)$ +	$L(B(\neg \clubsuit), 1)$	+	$L(B(\neg \clubsuit), 0)$	+	$L(B(\neg \clubsuit), 0)$	)
$e_{AC}(Y_2,0) = \frac{1}{4} ($	0 +	1	+	0	+	0	$) = \frac{1}{4}$
$e_{BPN}(Y_1) = \frac{1}{4} ($	$\mathbb{E}_{\mathbb{Q}_{\sharp}\mu} L(\cdot, 0) +$	$\mathbb{E}_{Q_{\sharp}\mu} L(\cdot, 1)$	+	$\mathbb{E}_{Q_{\sharp}\mu} L(\cdot, 0)$	+	$\mathbb{E}_{Q_{\sharp}\mu} L(\cdot, 0)$	)
$=\frac{1}{4}($	$(\frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 0) +$	$(\tfrac{1}{2} \cdot 0 + \tfrac{1}{2} \cdot 1)$	+	$\left(\frac{1}{2}\cdot 1 + \frac{1}{2}\cdot 0\right)$	+	$\left(\frac{1}{2}\cdot 0 + \frac{1}{2}\cdot 0\right)$	$)=rac{1}{4}$
$e_{BPN}(Y_2) = \frac{1}{4} ($	$\mathbb{E}_{\mathbb{Q}_{\sharp}\mu^{\bigstar}}L(\cdot,0) +$	$\mathbb{E}_{\mathbb{Q}_{\sharp}\mu^{\neg}} L(\cdot, 1)$	+	$\mathbb{E}_{\mathbf{Q}_{\sharp}\mu^{\neg} \bigstar} L(\cdot, 0)$	+	$\mathbb{E}_{Q_{\sharp}\mu^{\neg} \bigstar} L(\cdot,0)$	)
$=\frac{1}{4}($	$(1 \cdot 0) +$	$(\frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 1)$	) + (	$(\frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 0)$	) + (	$(\frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 0)$	$\left( \right) = \frac{1}{3}$

■ Back