# Brittleness and Robustness of Bayesian Inference

**Tim Sullivan**

with Houman Owhadi and Clint Scovel (Caltech)

Mathematics Institute, University of Warwick

*Predictive Modelling Seminar*
**University of Warwick, U.K.**
15 January 2015

THE UNIVERSITY OF
WARWICK

http://www.tjsullivan.org.uk/pdf/2015-01-15_wcpm.pdf

## What Do We Mean by 'Bayesian Brittleness'?

- Bayesian procedures give posterior distributions for quantities of interest in the form of Bayes' rule

$$p(\text{parameters}|\text{data}) \propto L(\text{data}|\text{parameters})p(\text{parameters})$$

given the following data:

  - a prior probability distribution on parameters — later denoted $u \in \mathbb{U}$;
  - a likelihood function;
  - observations / data — later denoted $y \in \mathbb{Y}$.

- It is natural to ask about the robustness, stability, and accuracy of such procedures.

- This is a subtle topic, with both positive and negative results, especially for large/complex systems, with fine geometrical and topological considerations playing a key role.

# What Do We Mean by 'Bayesian Brittleness'?

$$p(\text{parameters}|\text{data}) \propto L(\text{data}|\text{parameters})p(\text{parameters})$$

- Frequentist questions: If the data are generated from some 'true' distribution, will the posterior eventually/asymptotically identify the 'true' value? Are Bayesian credible sets also frequentist confidence sets? What if the model class doesn't even contain the 'truth'?

- Numerical analysis questions: Is Bayesian inference a well-posed problem, in the sense that small perturbations of the prior, likelihood, or data (e.g. those arising from numerical discretization) lead to small changes in the posterior? Can effective estimates be given?

- For us, 'brittleness' simply means the strongest possible negative result: under arbitrarily small perturbations of the problem setup the posterior conclusions change as much as possible — i.e. extreme discontinuity. (More precise definition later on.)

# Overview

# Overview

## Bayesian Modelling Setup

- Parameter space $\mathbb{U}$, equipped with a prior $\pi \in \mathcal{P}(\mathbb{U})$.
- Observed data with values in $\mathbb{Y}$ are explained using a likelihood model, i.e. a function $L \colon \mathbb{U} \to \mathcal{P}(\mathbb{Y})$ with

$$L(E|u) = \mathbb{P}\big[\, y \in E \mid u \,\big].$$

- This defines a (non-product) joint measure $\mu$ on $\mathbb{U} \times \mathbb{Y}$ by

$$\mu(E) := \mathbb{E}_{u \sim \pi, y \sim L(\,\cdot\,|u)}\big[\mathbb{1}_E(u, y)\big] \equiv \int_{\mathbb{U}} \int_{\mathbb{Y}} \mathbb{1}_E(u, y)\, L(\mathrm{d}y|u)\, \pi(\mathrm{d}u).$$

- The Bayesian posterior on $\mathbb{U}$ is just $\mu$ conditioned on a $\mathbb{Y}$-fibre, and re-normalized to be a probability measure. Bayes' Rule gives this as

$$\pi(E|y) = \frac{\mathbb{E}_{u \sim \pi}[\mathbb{1}_E(u) L(y|u)]}{\mathbb{E}_{u \sim \pi}[L(y|u)]}.$$

## Bayesian Modelling Setup

- Parameter space $\mathbb{U}$, equipped with a prior $\pi \in \mathcal{P}(\mathbb{U})$.
- Observed data with values in $\mathbb{Y}$ are explained using a likelihood model, i.e. a function $L \colon \mathbb{U} \to \mathcal{P}(\mathbb{Y})$ with

$$L(E|u) = \mathbb{P}\big[\, y \in E \mid u \,\big].$$

- This defines a (non-product) joint measure $\mu$ on $\mathbb{U} \times \mathbb{Y}$ by

$$\mu(E) := \mathbb{E}_{u \sim \pi, y \sim L(\,\cdot\,|u)}\big[\mathbb{1}_E(u, y)\big] \equiv \int_{\mathbb{U}} \int_{\mathbb{Y}} \mathbb{1}_E(u, y)\, L(\mathrm{d}y|u)\, \pi(\mathrm{d}u).$$

- The Bayesian posterior on $\mathbb{U}$ is just $\mu$ conditioned on a $\mathbb{Y}$-fibre, and re-normalized to be a probability measure. Bayes' Rule gives this as

$$\frac{\mathrm{d}\pi(\,\cdot\,|y)}{\mathrm{d}\pi} \propto L(y|\,\cdot\,).$$

# Bayesian Modelling Setup

### Traditional Setting

$\mathbb{U}$ is a finite set or $\mathbb{R}^d$ for small $d \in \mathbb{N}$.

### More Modern Applications

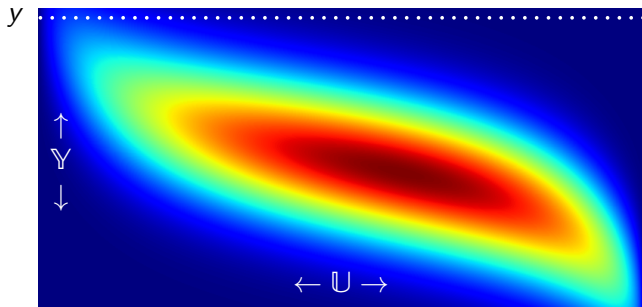A very high-dimensional or infinite-dimensional $\mathbb{U}$, e.g. an inverse problem for a PDE:

$$-\nabla \cdot (u\nabla p) = f,$$
$$\text{boundary conditions}(p) = 0.$$

in which we attempt to infer the permeability $u$ from e.g. some noisy point observations of the pressure/head $p$.

Prior measure $\pi$ on $\mathbb{U}$:



Joint measure $\mu$ on $\mathbb{U} \times \mathbb{Y}$:

$y$



$\uparrow$

$\mathbb{Y}$

$\downarrow$

$\leftarrow \mathbb{U} \rightarrow$

Posterior measure $\pi(\,\cdot\,|y) \propto \mu|_{\mathbb{U} \times \{y\}}$ on $\mathbb{U}$:

## Specification of Bayesian Models

- Parameter space $\mathbb{U}$, equipped with a prior $\pi \in \mathcal{P}(\mathbb{U})$.
- Observed data with values in $\mathbb{Y}$ are explained using a likelihood model, i.e. a function $L \colon \mathbb{U} \to \mathcal{P}(\mathbb{Y})$ with

$$L(E|u) = \mathbb{P}\big[\, y \in E \mid u \,\big].$$

### Definition (Frequentist well-specification)

If data are generated according to $\mu^{\dagger} \in \mathcal{P}(\mathbb{Y})$, then the Bayesian model is called well-specified if there is some $u^{\dagger} \in \mathbb{U}$ such that $\mu^{\dagger} = L(\,\cdot\,|u^{\dagger})$; otherwise, the model is called misspecified.

## Consistency of Bayesian Models

Suppose that the observed data consists of a sequence of independent $\mu^\dagger$-distributed samples $(y_1, y_2, \dots)$, and let

$$\pi^{(n)}(u) := \pi(u|y_1, \dots, y_n) \propto L(y_1, \dots, y_n|u)\pi(u)$$

be the posterior measure obtained by conditioning the prior $\pi$ with respect to the first $n$ observations using Bayes' rule.

### Definition (Frequentist consistency)

A well-specified Bayesian model with $\mu^\dagger = L(\,\cdot\,|u^\dagger)$ is called consistent (in an appropriate topology on $\mathcal{P}(\mathbb{U})$) if

$$\lim_{n\to\infty} \pi^{(n)} = \delta_{u^\dagger},$$

i.e. the posterior asymptotically gives full mass to the true parameter value.

# Bernstein–von Mises Theorem

The classic positive result regarding posterior consistency is the Bernstein–von Mises theorem or Bayesian CLT, historically first envisioned by **Laplace** (1810) and first rigorously proved by **Le Cam** (1953):

## Theorem (Bernstein–von Mises)

*If $\mathbb{U}$ and $\mathbb{Y}$ are finite-dimensional, then, subject to regularity assumptions on $L$ and $\pi$, any well-specified Bayesian model is consistent provided $u^\dagger \in \mathrm{supp}(\pi)$. Furthermore, $\pi^{(n)}$ is asymptotically normal about $\widehat{u}_n^{\mathsf{MLE}} \to u^\dagger$, with precision proportional to the Fisher information $\mathcal{I}(u^\dagger)$:*

$$\mathbb{P}_{y_i \sim \mu^\dagger}\left[\left\|\pi^{(n)} - \mathcal{N}\left(\widehat{u}_n^{\mathsf{MLE}}, \frac{\mathcal{I}(u^\dagger)^{-1}}{n}\right)\right\|_{\mathsf{TV}} > \varepsilon\right] \xrightarrow[n \to \infty]{} 0,$$

*where* $\qquad \mathcal{I}(u^\dagger)_{ij} = \mathbb{E}_{y \sim L(\,\cdot\,|u^\dagger)}\left[\left.\frac{\partial \log L(y|u)}{\partial u_i}\frac{\partial \log L(y|u)}{\partial u_j}\right|_{u=u^\dagger}\right].$

## Bernstein–von Mises Theorem

- Informally, the BvM theorem says that a well-specified model is capable of learning any 'truth' in the support of the prior.
- If we obey Cromwell's Rule

    *"I beseech you, in the bowels of Christ, think it possible that you may be mistaken."*

    by choosing a globally supported prior $\pi$, then everything should turn out OK — and the limiting posterior should be independent of $\pi$.
- Unfortunately, the BvM theorem is not always true if dim $\mathbb{U} = \infty$, even for globally supported priors — but nor is it always false.
- Applications of Bayesian methods in function spaces are increasingly popular, so it is important to understand the precise circumstances in which we do or do not have the BvM property.

# Some Positive and Negative Consistency Results

## Positive

- **Barron, Schervish & Wasserman** (1999): K–L and Hellinger
- **Castillo & Rousseau** and **Nickl & Castillo** (2013): Gaussian seq. space model, <u>modified</u> $\ell^2$ balls
- **Szabó, Van der Vaart, Van Zanten** (2014)
- **Stuart** & al. (2010+): Gaussian / Besov measures
- Dirichlet processes

## Negative

- **Freedman** (1963, 1965): prior supported on $\mathcal{P}(\mathbb{N}_0)$ sees i.i.d. $y_i \sim \text{Geom}(\frac{1}{4})$, but posterior $\rightarrow \text{Geom}(\frac{3}{4})$
- **Diaconis & Freedman** (1998): such 'bad' priors are of small measure, but are topologically generic
- **Johnstone** (2010) and **Leahu** (2011): further Freedman-type results
- $\rightarrow$ **Owhadi, Scovel & S.**

Main moral: the geometry and topology play a critical role in consistency.

# Consistency of Misspecified Bayesian Models

- By definition, if the model is mis-specified, then we cannot hope for posterior consistency in the sense that $\pi^{(n)} \to \delta_{u^\dagger}$ where $L(\,\cdot\,|u^\dagger) = \mu^\dagger$, because no such $u^\dagger \in \mathbb{U}$ exists.

- However, we can still hope that $\pi^{(n)} \to \delta_{\widehat{u}}$ for some 'meaningful' $\widehat{u} \in \mathbb{U}$, and that we get consistent estimates for the values of suitable quantities of interest, e.g. the posterior asymptotically puts all mass on $\widehat{u} \in \mathbb{U}$ such that $L(\,\cdot\,|\widehat{u})$ matches the mean and variance of $\mu^\dagger$, if not the exact distribution.

- For example, **Berk** (1966, 1970), **Kleijn & Van der Vaart** (2006), **Shalizi** (2009) have results of the type:

## Theorem (Minimum relative entropy)

*Under suitable regularity assumptions, the posterior concentrates on*

$$\widehat{u} \in \arg\min\Big\{ D_{\mathsf{KL}}\big(\mu^\dagger \big\| L(\,\cdot\,|u)\big) \Big| u \in \mathsf{supp}(\pi_0) \Big\}.$$

# Overview

## Setup for Brittleness

- For simplicity, $\mathbb{U}$ and $\mathbb{Y}$ will be complete and separable metric spaces — see arXiv:1304.6772 for weaker but more verbose assumptions.

- Fix a prior $\pi_0 \in \mathcal{P}(\mathbb{U})$ and likehood model $L_0$, and the induced joint measure (Bayesian model) $\mu_0$; we will consider other models $\mu_\alpha$ 'near' to $\mu_0$.

- Given $\pi_0$ and any quantity of interest $q \colon \mathbb{U} \to \mathbb{R}$,

$$\pi_0 \text{-}\underset{u \in \mathbb{U}}{\text{ess inf}}\, q(u) := \sup\big\{\, t \in \mathbb{R} \,\big|\, q(u) \geq t \ \pi_0\text{-a.s.} \,\big\},$$

$$\pi_0 \text{-}\underset{u \in \mathbb{U}}{\text{ess sup}}\, q(u) := \inf\big\{\, t \in \mathbb{R} \,\big|\, q(u) \leq t \ \pi_0\text{-a.s.} \,\big\}.$$

- To get around difficulties of data actually having measure zero, and with one eye on the fact that real-world data is always discretized to some precision level $0 < \delta < \infty$, we assume that our observation is actually that the 'exact' data lies in a metric ball $\mathbb{B}_\delta(y) \subseteq \mathbb{Y}$.

- Slight modification: $y$ could actually be $(y_1, \ldots, y_n) \in \mathbb{Y}^n$.

## Setup for Brittleness

The brittleness theorem covers three notions of closeness between models:

- total variation distance: for $\alpha > 0$ (small), $\|\mu_0 - \mu_\alpha\|_{\mathsf{TV}} < \alpha$; or
- Prohorov distance: for $\alpha > 0$ (small), $d_\Pi(\mu_0, \mu_\alpha) < \alpha$ (for separable $\mathbb{U}$, this metrizes the weak convergence topology on $\mathcal{P}(\mathbb{U})$); or
- common moments: for $\alpha \in \mathbb{N}$ (large), for prescribed measurable functions $\phi_1, \ldots, \phi_\alpha \colon \mathbb{U} \times \mathbb{Y} \to \mathbb{R}$,

$$\mathbb{E}_{\mu_0}[\phi_i] = \mathbb{E}_{\mu_\alpha}[\phi_i] \quad \text{for } i = 1, \ldots, \alpha,$$

or, for $\varepsilon_i > 0$,

$$\left| \mathbb{E}_{\mu_0}[\phi_i] - \mathbb{E}_{\mu_\alpha}[\phi_i] \right| \leq \varepsilon_i \quad \text{for } i = 1, \ldots, \alpha.$$

# Brittleness Theorem

## Theorem (Brittleness)

*Suppose that the original model $(\pi_0, L_0)$ permits observed data to be arbitrarily unlikely in the sense that*

$$\lim_{\delta \to 0} \sup_{\substack{y \in \mathbb{Y} \\ u \in \mathsf{supp}(\pi_0) \subseteq \mathbb{U}}} L_0\big(\mathbb{B}_\delta(y) \,\big|\, u\big) = 0, \tag{AU}$$

*and let $q \colon \mathbb{U} \to \mathbb{R}$ be any measurable function. Then, for all*

$$v \in \left[ \pi_0 \operatorname*{-ess\,inf}_{u \in \mathbb{U}} q(u), \pi_0 \operatorname*{-ess\,sup}_{u \in \mathbb{U}} q(u) \right],$$
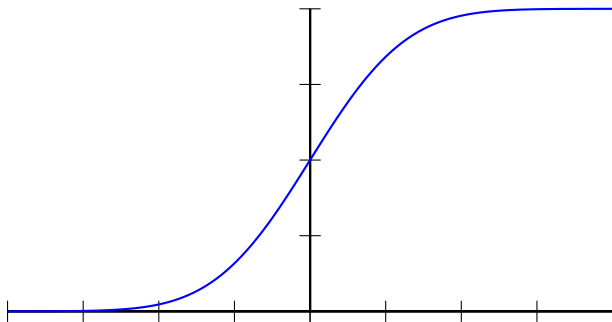
*and all $\alpha > 0$, there exists $\delta_*(\alpha) > 0$ and a model $\mu_\alpha$ 'α-close' to $\mu_0$ such that the posterior value $\mathbb{E}_{\pi_\alpha}\big[ q \,\big|\, \mathbb{B}_\delta(y) \big]$ for q given data of precision $0 < \delta < \delta_*(\alpha)$ is the chosen value v.*

# Brittleness Theorem — Proof Sketch

## Idea of Proof

- Optimize over the set $\mathcal{A}_\alpha$ of models that are $\alpha$-close to the original model. (*Cf.* construction of Bayesian least favourable priors and frequentist minimax estimators.)

- This involves understanding extreme points of $\mathcal{A}_\alpha$ and the optimization of affine functionals over such sets — Choquet theory and results of **von Weizsäcker & Winkler** — and previous work by S. and collaborators on Optimal UQ (*SIAM Rev.*, 2013).

- The three notions of closeness considered (moments, Prohorov, TV), plus the (AU) condition, together permit models $\mu_\alpha \in \mathcal{A}_\alpha$ to 'choose which data to trust' when forming the posterior.

- In our proof as written, the perturbations used to produce the 'bad' models use point masses; a slight variation would produce the same result using absolutely continuous perturbations.
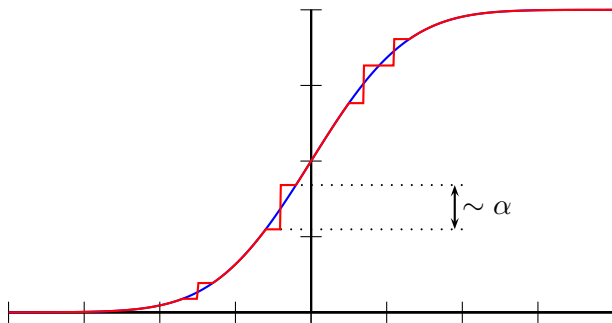
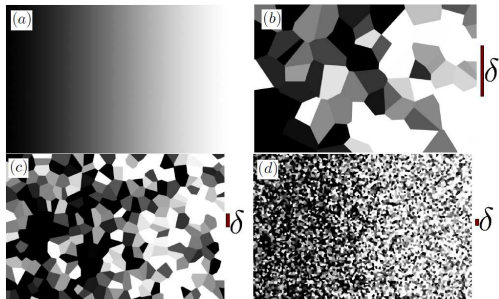Schematically, the perturbation from $\mu_0$ to $\mu_\alpha$ looks like

Schematically, the perturbation from $\mu_0$ to $\mu_\alpha$ looks like

## Brittleness Theorem — Interpretation

- Misspecification has profound consequences for Bayesian robustness on 'large' spaces — in fact, Bayesian inferences become extremely brittle as a function of measurement resolution $\delta$.
- If the model is misspecified, and there are possible observed data that are arbitrarily unlikely under the model, then under fine enough measurement resolution the posterior predictions of nearby priors differ as much as possible *regardless of the number of samples observed*.

Figure. As measurement resolution $\delta \to 0$, the smooth dependence of $\mathbb{E}_{\pi_0}[q]$ on the prior $\pi_0$ (top-left) shatters into a patchwork of diametrically opposed posterior values $\mathbb{E}_{\pi^{(n)}}[q] \equiv \mathbb{E}_{\pi_0}[q|\mathbb{B}_\delta(y)]$.

## Brittleness Rates — A Moment-Based Example

- Estimate the mean of a random variable $X$, taking values in $[0, 1]$, given a single observation $y$ of $X$
- Set $\mathcal{A}$ of admissible priors for the law of $X$: anything that gives uniform measure to the mean, uniform measure to the second moment given the mean, uniform measure to the third moment given the second, ... up to $k^{\text{th}}$ moment. (Note that $\dim \mathcal{A} = \infty$ but $\text{codim } \mathcal{A} = k$.)
- So, in particular, for any prior $\pi \in \mathcal{A}$, $\mathbb{E}_\pi[\mathbb{E}[X]] = \frac{1}{2}$.
- Can find priors $\pi_1, \pi_2 \in \mathcal{A}$ with

$$\mathbb{E}_{\pi_1}[\mathbb{E}[X]|y] \leq 4e \left[\frac{2k\delta}{e}\right]^{\frac{1}{2k+1}} \approx 0,$$

$$\mathbb{E}_{\pi_2}[\mathbb{E}[X]|y] \geq 1 - 4e \left[\frac{2k\delta}{e}\right]^{\frac{1}{2k+1}} \approx 1.$$

# Ways to Restore Robustness and Consistency
Or: What Would Break This Argument?

- Restrict to finite-precision data, i.e. keep $\delta$ bounded away from zero. Physically quite reasonable. The universe may be granular enough that $\delta^{1/(2k+1)} \gg 0$ for all 'practical' $\delta > 0$.

- Ask the robustness question before seeing the data, not after. This leads to a very large minimax problem, the computation of data-schema-specific optimal statistical estimators.

- Ask the robustness question about the limiting posterior, not each $\pi^{(n)}$ individually. The brittleness theorem and "$\lim_{n\to\infty}$" might not commute.

## Closing Remarks on Brittleness

- In contrast to the classical robustness and consistency results for Bayesian inference for discrete or finite-dimensional systems, the situation for infinite-dimensional spaces is *complicated*.
- Bayesian inference is extremely brittle in *some* topologies, and so cannot be consistent, and high-precision data only worsens things.
- Consistency *can* hold for complex systems, with *careful* choices of prior, geometry and topology — but, since the situation is so sensitive, all assumptions must be considered carefully.
- And, once a 'mathematical' prior is agreed upon, just as with classical numerical analysis of algorithms for ODEs and PDEs, the onus is on the algorithm designer to ensure that the 'numerical' prior is close to the 'mathematical' one in a 'good' topology.

# Overview

## Closing Remarks

- In contrast to the classical robustness and consistency results for Bayesian inference for discrete or finite-dimensional systems, the situation for infinite-dimensional spaces is *complicated*.

- Bayesian inference is extremely brittle in *some* topologies, and so cannot be consistent, and high-precision data only worsens things.

- Consistency *can* hold for complex systems, with *careful* choices of prior, geometry and topology — but, since the situation is so sensitive, all assumptions must be considered carefully.

- And, once a 'mathematical' prior is agreed upon, just as with classical numerical analysis of algorithms for ODEs and PDEs, the onus is on the algorithm designer to ensure that 'numerical' prior is close to the 'mathematical' one in a 'good' topology.

## Some Questions

- What happens if we do the physically reasonable thing of restricting to finite-precision data, i.e. keeping $\delta$ bounded away from zero? — Need quantitative versions of these theorems! The one-dimensional $k$-moments example suggests that the rate is not too bad, but what is the general picture?
- What happens if we ask the robustness question *before* seeing the data, not after?
- What happens if we ask the robustness question about the limiting posterior, not each $\pi^{(n)}$ individually? The brittleness theorem and "$\lim_{n\to\infty}$" might not commute.
- How does this relate to the phenomenon of Bayesian dilation?

# Thank You

- H. Owhadi & C. Scovel, arXiv:1304.7046
- H. Owhadi, C. Scovel & T. J. Sullivan, arXiv:1304.6772
- H. Owhadi, C. Scovel & T. J. Sullivan, arXiv:1308.6306
- H. Owhadi, C. Scovel, T. J. Sullivan, M. McKerns & M. Ortiz, *SIAM Rev.* **55**(2):271–345, 2013. arXiv:1009.0679
- T. J. Sullivan, M. McKerns, D. Meyer, F. Theil, H. Owhadi & M. Ortiz, *Math. Model. Numer. Anal.* **47**(6):1657–1689, 2013. arXiv:1202.1928