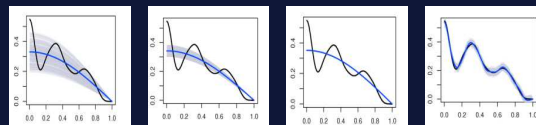


# *Nonparametric Bayes: review and challenges*

**Aad van der Vaart**

Universiteit Leiden, Netherlands



European Meeting of Statisticians  
Palermo, Italy, July 2019

**Introduction**

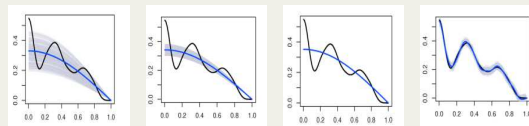
**Recovery**

**Uncertainty quantification**

**Uncertainty quantification for curve estimation**

**Uncertainty quantification for sparse high-dimensional parameters**

**Closing remarks**



# Introduction

# The Bayesian paradigm



- A parameter  $\theta$  is generated according to a **prior distribution**  $\Pi$
- Given  $\theta$  the data  $X$  is generated according to a measure  $P_\theta$

This gives a **joint distribution** of  $(X, \theta)$

- Given observed data  $X$  the statistician computes the conditional distribution of  $\theta$  given  $X$ , the **posterior distribution**:

$$\Pi(\theta \in B | X).$$

# The Bayesian paradigm



- A parameter  $\theta$  is generated according to a **prior distribution**  $\Pi$
- Given  $\theta$  the data  $X$  is generated according to a measure  $P_\theta$

This gives a **joint distribution** of  $(X, \theta)$

- Given observed data  $X$  the statistician computes the conditional distribution of  $\theta$  given  $X$ , the **posterior distribution**:

$$\Pi(\theta \in B | X).$$

If  $P_\theta$  is given by a density  $x \mapsto p_\theta(x)$ , then **Bayes's rule** gives

$$d\Pi(\theta | X) \propto p_\theta(X) d\Pi(\theta)$$

# Frequentist Bayes

Assume the data  $X$  is generated according to a **given parameter**  $\theta_0$   
Consider the posterior  $\Pi(\theta \in \cdot | X)$  as a given random measure

# Frequentist Bayes

Assume the data  $X$  is generated according to a **given parameter**  $\theta_0$   
Consider the posterior  $\Pi(\theta \in \cdot | X)$  as a given random measure

## Recovery

We like  $\Pi(\theta \in \cdot | X)$  to put “most” of its mass near  $\theta_0$  for “most”  $X$

# Frequentist Bayes

Assume the data  $X$  is generated according to a **given parameter**  $\theta_0$   
Consider the posterior  $\Pi(\theta \in \cdot | X)$  as a given random measure

## Recovery

We like  $\Pi(\theta \in \cdot | X)$  to put “most” of its mass near  $\theta_0$  for “most”  $X$

## Uncertainty quantification

We like the “spread” of  $\Pi(\theta \in \cdot | X)$  to indicate remaining uncertainty



# Frequentist Bayes

Assume the data  $X$  is generated according to a **given parameter**  $\theta_0$   
Consider the posterior  $\Pi(\theta \in \cdot | X)$  as a given random measure

## Recovery

We like  $\Pi(\theta \in \cdot | X)$  to put “most” of its mass near  $\theta_0$  for “most”  $X$

## Uncertainty quantification

We like the “spread” of  $\Pi(\theta \in \cdot | X)$  to indicate remaining uncertainty

### Asymptotic setting:

Data  $X^{(n)}$  where the information increases as  $n \rightarrow \infty$

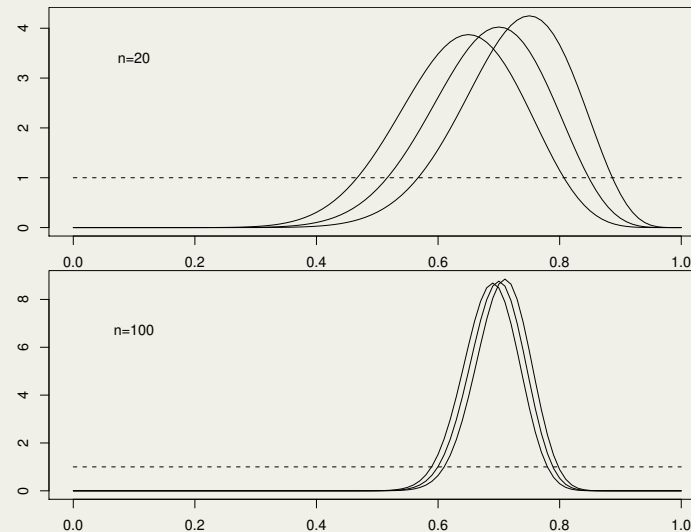
- We want  $\Pi_n(\cdot | X^{(n)}) \rightsquigarrow \delta_{\theta_0}$ , at a good rate
- We like a set of large posterior mass to *cover*

**Data:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_\theta$   
 $\mathbb{R}^d \ni \theta \mapsto p_\theta$  smooth and **identifiable**

**Thm** Under  $\theta_0$ , for **any prior** with positive density,

$$\left\| \Pi(\cdot | X_1, \dots, X_n) - N_d\left(\tilde{\theta}_n, \frac{1}{n} I_{\theta_0}^{-1}\right)(\cdot) \right\|_{TV} \rightarrow 0$$

Here  $\tilde{\theta}_n$  are estimators with  $\sqrt{n}(\tilde{\theta}_n - \theta_0) \rightsquigarrow N(0, I_{\theta_0}^{-1})$



**Data:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_\theta$   
 $\mathbb{R}^d \ni \theta \mapsto p_\theta$  smooth and **identifiable**

**Thm** Under  $\theta_0$ , for **any prior** with positive density,

$$\left\| \Pi(\cdot | X_1, \dots, X_n) - N_d(\tilde{\theta}_n, \frac{1}{n} I_{\theta_0}^{-1})(\cdot) \right\|_{TV} \rightarrow 0$$

Here  $\tilde{\theta}_n$  are estimators with  $\sqrt{n}(\tilde{\theta}_n - \theta_0) \rightsquigarrow N(0, I_{\theta_0}^{-1})$

## Recovery:

The posterior distribution concentrates most of its mass on balls of radius  $O(1/\sqrt{n})$  around  $\theta_0$

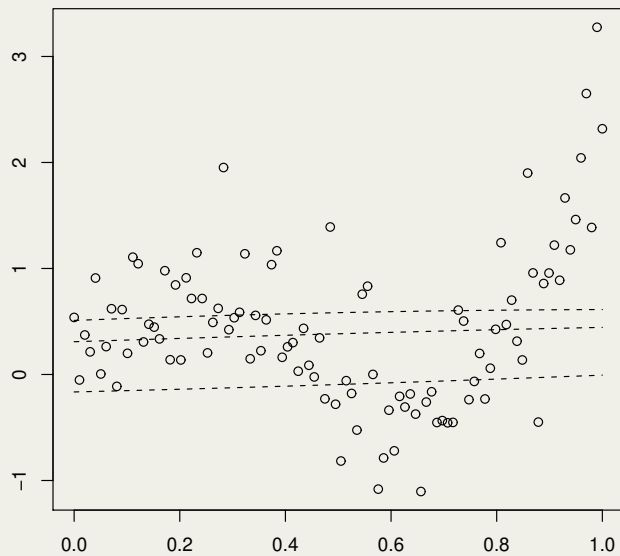
## Uncertainty quantification:

A central set of posterior probability 95 % is equivalent to the usual Wald confidence set  $\{\theta: n(\theta - \tilde{\theta}_n)^T I_{\tilde{\theta}_n} (\theta - \tilde{\theta}_n) \leq \chi_{d,1-\alpha}^2\}$

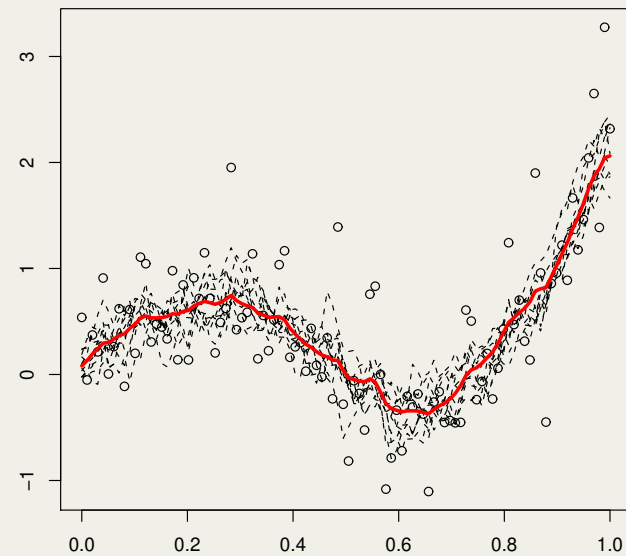
# Bayesian curve estimation

A prior and posterior of a **function** can be visualized by plotting functions that are simulated from the prior and posterior distributions

prior



posterior



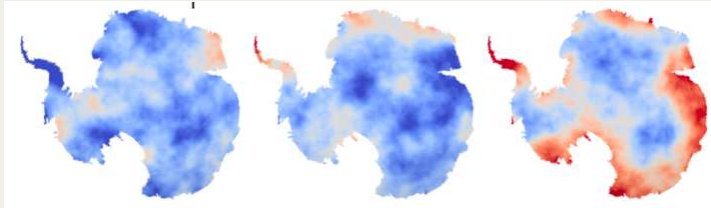
Many examples of priors: Dirichlet, Gaussian, random series,...

*Recovery well understood*

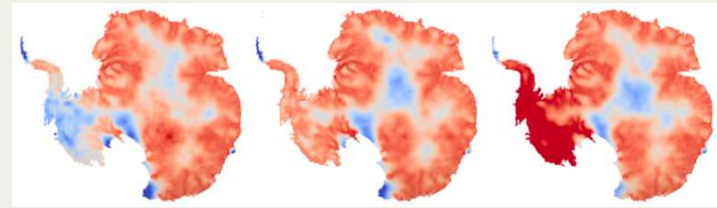
# Bayesian inverse problems — data assimilation

A prior and posterior of a [surface](#) can be visualized by plotting surfaces that are simulated from the prior and posterior distributions.

prior



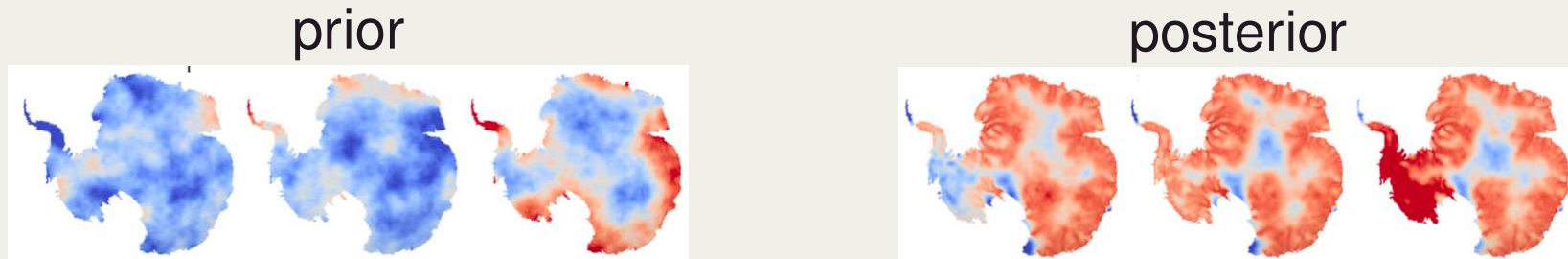
posterior



From Stadler et al., 2017

# Bayesian inverse problems — data assimilation

A prior and posterior of a **surface** can be visualized by plotting surfaces that are simulated from the prior and posterior distributions.



From Stadler et al., 2017

$$\text{Data} := B_\theta u + \text{noise}, \quad \text{for } D_\theta u = 0, \quad u \in G_\theta$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left[ L(\theta, \text{Data}) + \text{penalty}(\theta) \right]$$

$$\begin{aligned} -\nabla \cdot [2\eta(u, n) \dot{\epsilon}_u - Ip] &= \rho g && \text{in } \Omega \\ \nabla \cdot u &= 0 && \text{in } \Omega \\ \sigma_u n &= 0 && \text{on } \Gamma_t \\ u \cdot n = 0, T\sigma_u n + \exp(\beta)Tu &= 0 && \text{on } \Gamma_b \end{aligned}$$

## Thm

Dashti et al 2013

If  $\text{penalty}(\theta) = \|\theta\|_{\mathbb{H}}^2$  for RKHS norm,  $\hat{\theta}$  is **posterior mode** for Gaussian prior

Connects to applied analysis

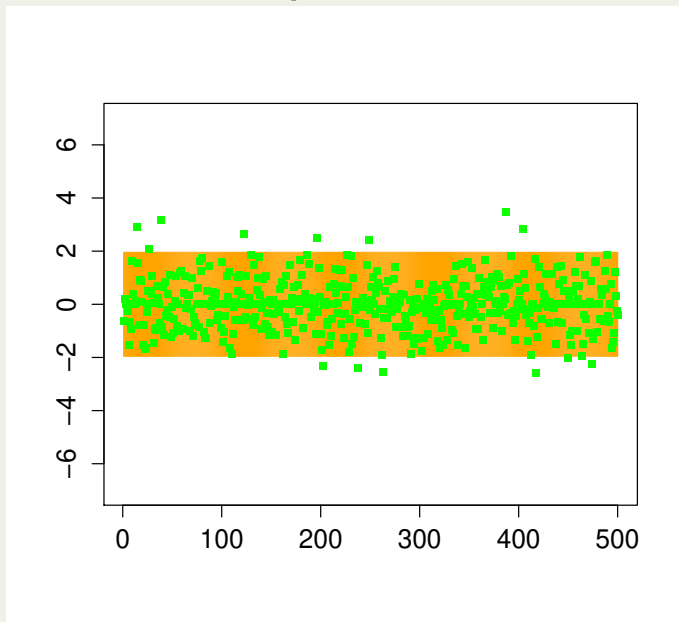
*Bayesian understanding starting to develop*

Stuart, Agapiou, Nickl,...

# High-dimensional Bayes

A **high-dimensional parameter** vector (or matrix) may be visualized through a plot of marginal distributions versus an index

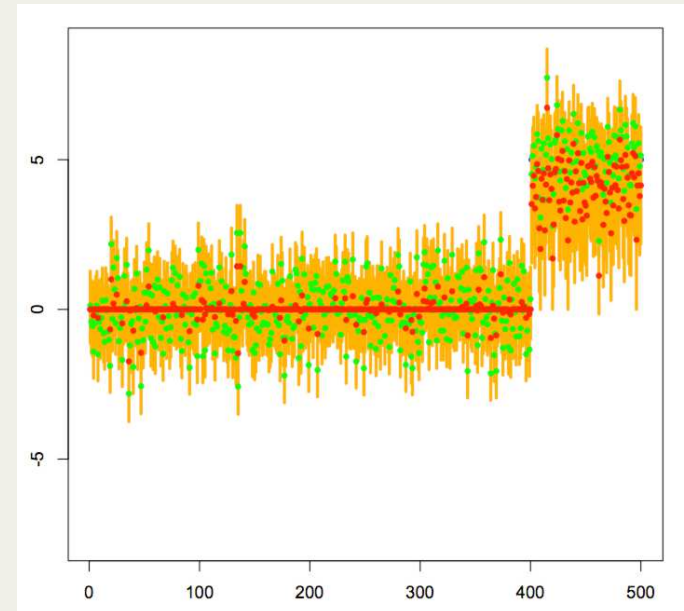
prior



Parameters  $\theta_1, \dots, \theta_{500}$  versus index  $1, \dots, 500$   
Orange: marginal prior predictive intervals

Green dots: prior draws

posterior



Parameters  $\theta_1, \dots, \theta_{500}$  versus index  $1, \dots, 500$   
Red dots: marginal posterior medians  
Orange: marginal credible intervals

Green dots: data points

Connects to Empirical Bayes and Large Scale Inference

*Recent progress*

Recovery

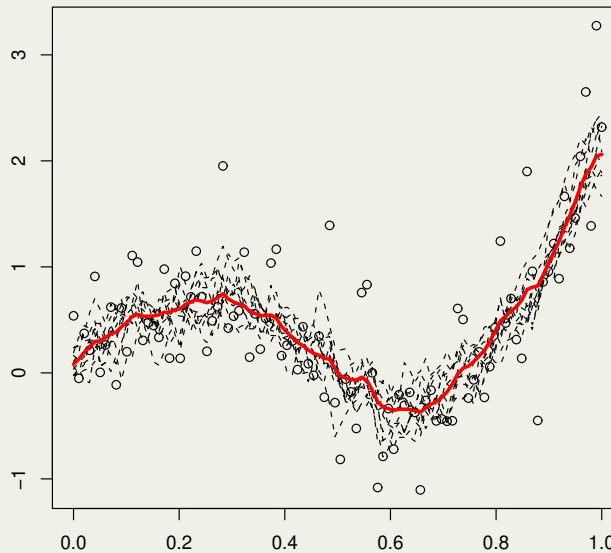


# Rate of contraction

$$\text{Data: } X^{(n)} \sim P_{\theta}^{(n)} \quad \theta \in (\Theta, d)$$

**Def** contraction rate at  $\theta_0$  is  $\epsilon_n$  if, for large  $M$ ,

$$E_{\theta_0} \Pi_n(\theta: d(\theta, \theta_0) > M\epsilon_n | X^{(n)}) \rightarrow 0, \quad n \rightarrow \infty$$



# Rate of contraction

**Data:**  $X^{(n)} \sim P_{\theta}^{(n)}$       $\theta \in (\Theta, d)$

**Def**     *contraction rate* at  $\theta_0$  is  $\epsilon_n$  if, for large  $M$ ,

$$\mathbb{E}_{\theta_0} \Pi_n(\theta: d(\theta, \theta_0) > M\epsilon_n | X^{(n)}) \rightarrow 0, \quad n \rightarrow \infty$$

**Benchmark rate for (inverse) curve fitting:**

A function  $\theta$  of  $d$  variables with bounded derivatives of order  $\beta$  is estimable based on  $n$  observations at rate

$$n^{-\beta/(2\beta+d+2p)}.$$

**Benchmark rate for sparse estimation:**

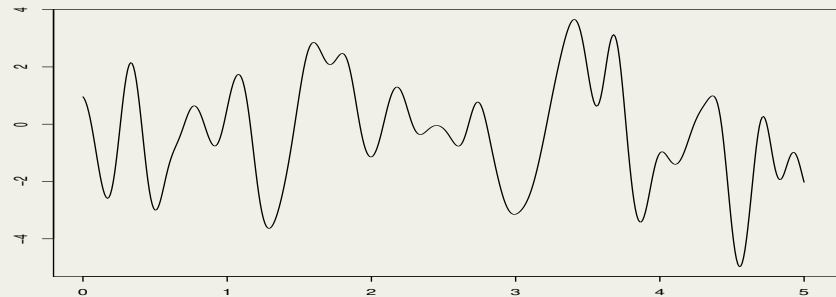
A vector  $\theta$  in  $\mathbb{R}^n$  of which  $s \ll n$  coordinates are nonzero is estimable based on 1 observation per parameter at rate

$$\sqrt{s \log(n/s)}.$$

**Data:** Sample of size  $n$  in regression model or from density

**Prior** on regression function or log density  $\theta$ : centered **Gaussian process** with

$$\text{cov}(\theta_s, \theta_t) = e^{-\|s-t\|^2}, \quad s, t \in \mathbb{R}^d$$



**Thm** Rate of contraction is  $(\log n)^\gamma / \sqrt{n}$  if  $\theta_0$  is analytic, but is  $(1/\log n)^k$  if  $\theta_0$  is only ordinary smooth

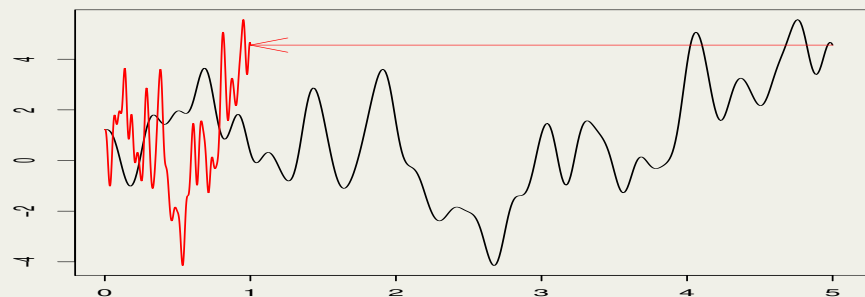
$$\mathbb{P}(\|\theta\|_\infty < \varepsilon) \gtrsim e^{-C(\log \varepsilon^{-1})^{1+d/2}}$$

# Curve fitting with square exponential prior — adaptation

**Data:** Sample of size  $n$  in regression model or from density

**Prior** on regression function or log density  $\theta$ :

- $c^d \sim \Gamma(a, b)$
- $(G_t: t > 0)$  square exponential process
- $\theta_t \sim G_{ct}$



**Thm** Rate of contraction is:

- if  $\theta_0 \in C^\beta[0, 1]^d$ , then nearly  $n^{-\beta/(2\beta+d)}$
- if  $\theta_0$  is analytic, then nearly  $n^{-1/2}$

**Data:**  $X^{(n)} = K\theta + n^{-1/2}\dot{W}$ , for white noise  $\dot{W}$

- $K$  compact operator with eigen basis  $(e_i)$  and eigenvalues  $\kappa_i \asymp i^{-p}$
- **Prior:**  $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , with  $\theta_i | \alpha \stackrel{\text{ind}}{\sim} N(0, i^{-2\alpha-1})$

**Data:**  $X^{(n)} = K\theta + n^{-1/2}\dot{W}$ , for white noise  $\dot{W}$

- $K$  compact operator with eigen basis  $(e_i)$  and eigenvalues  $\kappa_i \asymp i^{-p}$
- **Prior:**  $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , with  $\theta_i | \alpha \stackrel{\text{ind}}{\sim} N(0, i^{-2\alpha-1})$

**Thm** If  $\sum_{i=1}^{\infty} i^{2\beta} \theta_{i,0}^2 < \infty$ , then rate:

- $n^{-\beta/(2\alpha+2p+1)}$ , if  $\beta \leq \alpha$
- $n^{-\alpha/(2\alpha+2p+1)}$ , if  $\beta \geq \alpha$

**Data:**  $X^{(n)} = K\theta + n^{-1/2}\dot{W}$ , for white noise  $\dot{W}$

- $K$  compact operator with eigen basis  $(e_i)$  and eigenvalues  $\kappa_i \asymp i^{-p}$
- **Prior:**  $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , with  $\theta_i | \alpha \stackrel{\text{ind}}{\sim} N(0, i^{-2\alpha-1})$

**Thm** If  $\sum_{i=1}^{\infty} i^{2\beta} \theta_{i,0}^2 < \infty$ , then rate:

- $n^{-\beta/(2\alpha+2p+1)}$ , if  $\beta \leq \alpha$
- $n^{-\alpha/(2\alpha+2p+1)}$ , if  $\beta \geq \alpha$

- **Prior on  $\alpha$**

**Thm** If  $\sum_{i=1}^{\infty} i^{2\beta} \theta_{0,i}^2 < \infty$  and eigenvalues  $\kappa_i \asymp i^{-p}$ , then rate:

- $n^{-\beta/(2\beta+2p+1)}$ , any  $\beta > 0$

# Gaussian process priors



When using a Gaussian process as prior for a function:

Recovery is best if prior 'matches' truth  
Mismatch slows down, but does not prevent, recovery  
Mismatch can be prevented by using hyperparameters

(Generalizes to non-Gaussian priors

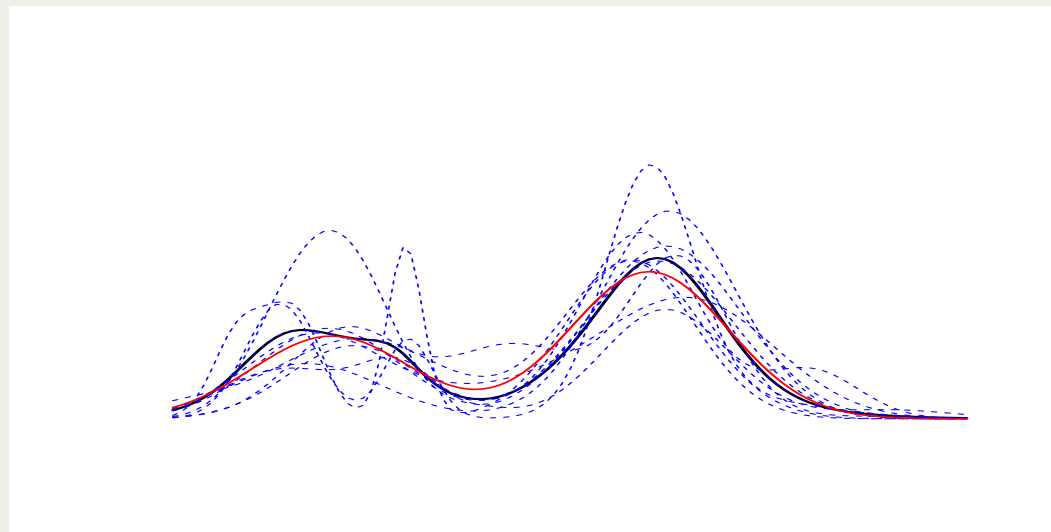
Ray, Yan, Agapiou,...)



**Data:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p$

- $F \sim$  Dirichlet process
- $1/c \sim \Gamma(a, b)$ , independent of  $F$

$$p_{F,c}(x) = \int \frac{1}{c} \phi\left(\frac{x-z}{c}\right) dF(z)$$



Posterior mean (solid black) and 10 draws of the posterior distribution  
for a sample of size 50 from a mixture of two normals (red)

[Plot by DPpackage, Jara et al., 2011]

**Data:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p$

- $F \sim$  Dirichlet process
- $1/c \sim \Gamma(a, b)$ , independent of  $F$

$$p_{F,c}(x) = \int \frac{1}{c} \phi\left(\frac{x-z}{c}\right) dF(z)$$

- Thm** Hellinger rate of contraction for  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_0$  is, any  $\beta > 0$ ,
- nearly  $n^{-1/2}$  if  $p_0 = p_{F_0, c_0}$ , some  $F_0, c_0$
  - nearly  $n^{-\beta/(2\beta+d)}$  if  $p_0 \in C^\beta(\mathbb{R}^d)$  with exponentially small tails

**Data:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p$

- $F \sim$  Dirichlet process
- $1/c \sim \Gamma(a, b)$ , independent of  $F$

$$p_{F,c}(x) = \int \frac{1}{c} \phi\left(\frac{x-z}{c}\right) dF(z)$$

- Thm** Hellinger rate of contraction for  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_0$  is, any  $\beta > 0$ ,
- nearly  $n^{-1/2}$  if  $p_0 = p_{F_0, c_0}$ , some  $F_0, c_0$
  - nearly  $n^{-\beta/(2\beta+d)}$  if  $p_0 \in C^\beta(\mathbb{R}^d)$  with exponentially small tails

Adaptation to any smoothness with a **Gaussian** kernel!  
(Kernel density estimation needs *higher order* kernels)

$$\frac{1}{nc} \sum_{i=1}^n \phi\left(\frac{x - X_i}{c}\right) = p_{F_n, c}(x)$$

**Data:**  $Y^n \sim N_n(\theta, I)$ , for  $\theta = (\theta_1, \dots, \theta_n) \in \mathbb{R}^n$

- $\tau \sim B(1, n + 1)$
- $\theta_i \stackrel{\text{iid}}{\sim} (1 - \tau)\delta_0 + \tau G$ , e.g.  $G = \text{Laplace}$

**Thm** For  $s_n \rightarrow \infty$  with  $s_n \ll n$ ,

$$\sup_{\#(\theta_{0,i} \neq 0) \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: \|\theta - \theta_0\|_2^2 \gtrsim s_n \log(n/s_n) | Y^n) \rightarrow 0.$$

**Shrinkage** controlled by sparsity parameter  $\tau$   
(interpretation:  $\tau \approx (s/n) \sqrt{\log n/s}$ )

# Uncertainty quantification

# Credible sets



**Def** A **credible set** is a data-dependent set  $C(X)$  with

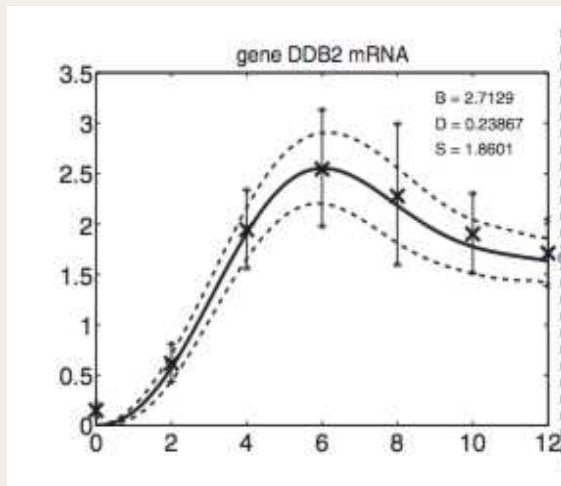
$$\Pi(\theta \in C(X) | X) = 0.95.$$

# Credible sets

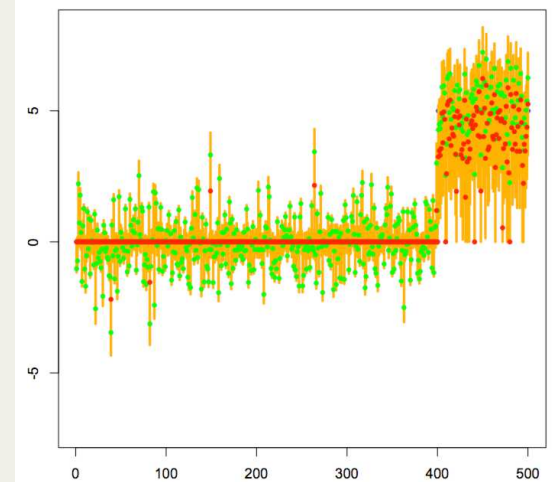


**Def** A **credible set** is a data-dependent set  $C(X)$  with

$$\Pi(\theta \in C(X) | X) = 0.95.$$



Estimated abundance of a transcription factor as function of time:  
posterior mean curve and 95% credible bands  
(Gao et al. *Bioinformatics*, 2008)



Red dots: marginal posterior medians  
Orange: marginal credible intervals

Green dots: data points

# Do credible sets correctly quantify *remaining uncertainty*?

Is a **credible set** a **confidence set**?

credible set

$$\Pi(\theta \in C(X) | X) = 0.95$$

confidence set

$$P_{\theta_0}(\theta_0 \in C(X)) = 0.95 \quad \forall \theta_0$$



# Do credible sets correctly quantify *remaining uncertainty*?

Is a **credible set** a **confidence set**?

credible set

$$\Pi(\theta \in C(X) | X) = 0.95$$

confidence set

$$P_{\theta_0}(\theta_0 \in C(X)) = 0.95 \quad \forall \theta_0$$

**Meta Thm**

Cox 1993, Freedman 2000, Leahu 2012

Only if some version of the Bernstein-von Mises theorem holds

# Do credible sets correctly quantify *remaining uncertainty*?

Is a **credible set** a **confidence set**?

credible set

$$\Pi(\theta \in C(X) | X) = 0.95$$

confidence set

$$P_{\theta_0}(\theta_0 \in C(X)) = 0.95 \quad \forall \theta_0$$

**Meta Thm**

Cox 1993, Freedman 2000, Leahu 2012

Only if some version of the Bernstein-von Mises theorem holds

Identify  $\theta$  with a set  $\Psi$  of smooth functionals  $\psi(\theta)$

A joint Bernstein-von Mises theorem

$$d\left(\Pi_n(\theta: \sqrt{n}(\psi(\theta) - \hat{\psi}_n)_{\psi \in \Psi} \in \cdot | X^{(n)}), P(Z(\psi)_{\psi \in \Psi} \in \cdot)\right) \rightarrow 0$$

may be used to get valid credible sets

Castillo, Nickl, Ray .. make this operational using weak norms

# Do credible sets correctly quantify *remaining uncertainty*?

Is a **credible set** a **confidence set**?

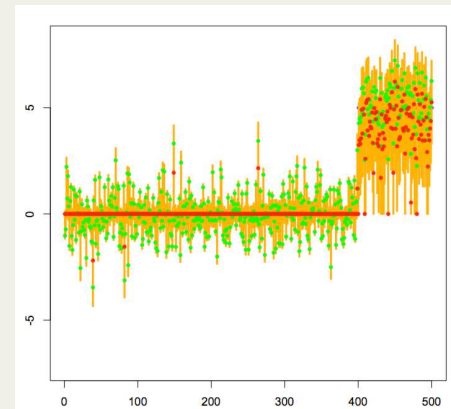
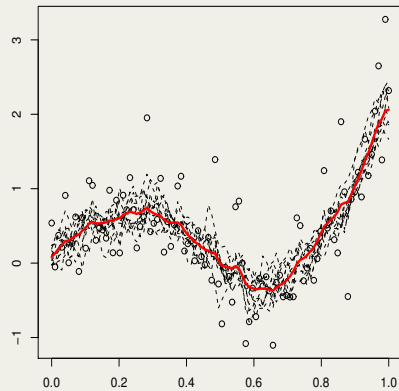
credible set

$$\Pi(\theta \in C(X) | X) = 0.95$$

confidence set

$$P_{\theta_0}(\theta_0 \in C(X)) = 0.95 \quad \forall \theta_0$$

Does the **spread in the posterior** give the correct order of the discrepancy between  $\theta_0$  and the posterior mean?



## Deterministic bandwidth: coverage requires undersmoothing

In *nonparametric statistics*:

**oversmoothing** gives **big bias** and **small variance** and hence **no coverage**

## Deterministic bandwidth: coverage requires undersmoothing

In *nonparametric statistics*:

oversmoothing gives **big bias** and **small variance** and hence **no coverage**

In *nonparametric Bayesian statistics*:

oversmoothing occurs if the **prior produces too smooth functions**

## \*\* Example: heat equation

**Data:**  $X^{(n)} = K\theta + n^{-1/2}\dot{W}$ , for white noise  $\dot{W}$

For given **initial heat curve**  $\theta: [0, 1] \rightarrow \mathbb{R}$  let  $K\theta = u(\cdot, 1)$  be the **final curve**:

$$\frac{\partial}{\partial t}u(x, t) = \frac{\partial^2}{\partial x^2}u(x, t), \quad u(\cdot, 0) = \theta, \quad u(0, t) = u(1, t) = 0$$

**Observe** noisy version  $(X_x^{(n)}: 0 \leq x \leq 1)$  of final curve

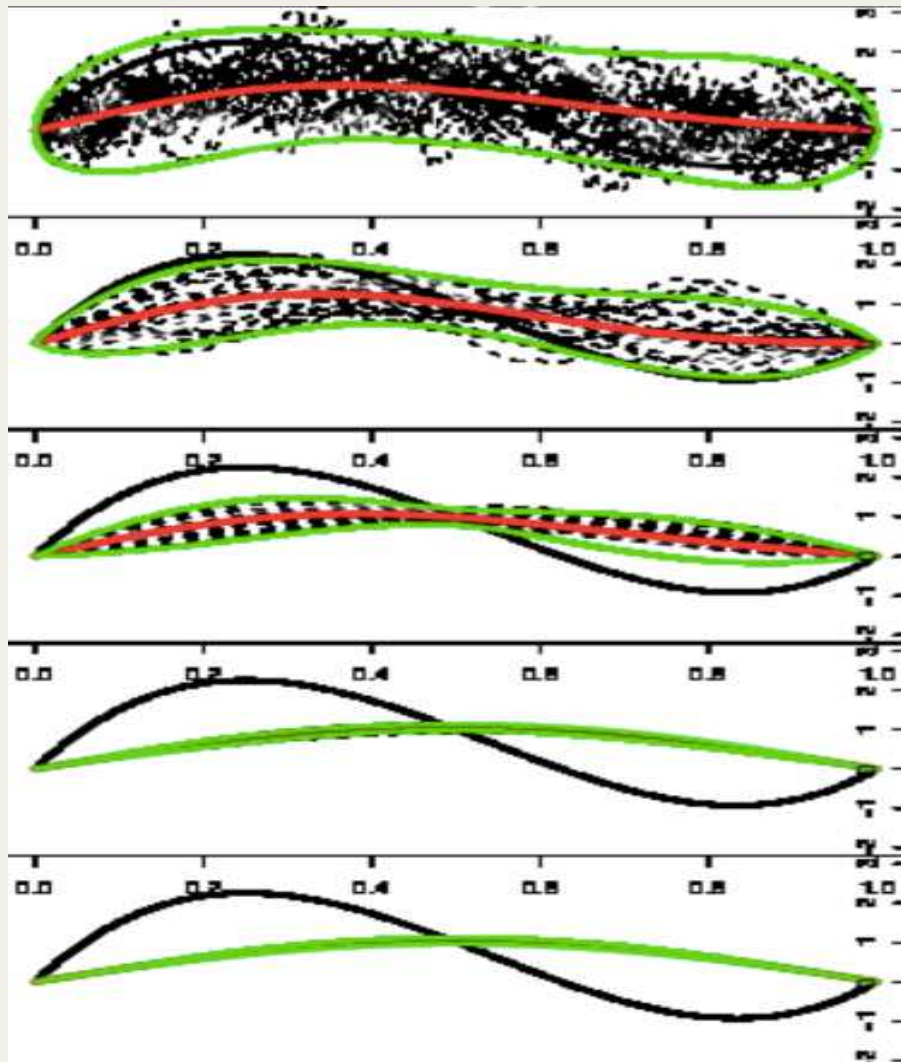
- $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , for  $e_i$  eigenbasis of  $K$
- **Truth:**  $\theta_{0,i} \asymp i^{-1-2\beta}$
- **Prior:**  $\theta_i \stackrel{\text{ind}}{\sim} N(0, i^{-1-2\alpha})$

### **Interpretation:**

$\alpha = \beta$ : prior and truth match

$\alpha > \beta$ : prior oversmooths

$\alpha < \beta$ : prior undersmooths



True  $\theta_0$  (black), posterior mean (red)

20 realizations from the posterior (dashed black)

posterior credible bands (green)

- $\theta = \sum_{i=1}^{\infty} \theta_i e_i$
- **Truth:**  
 $\theta_{0,i} \asymp i^{-1-2\beta}$
- **Prior:**  
 $\theta_i \stackrel{\text{ind}}{\sim} N(0, i^{-1-2\alpha})$

Top to bottom:  
increasing  $\alpha$

# Bayesian adaptation

Family of priors  $\Pi_\alpha$  of varying smoothness  $\alpha$ ; posteriors  $\Pi_\alpha(\cdot | X)$

## Examples

- $t \mapsto \sum_{i=1}^{\infty} \theta_i e_i(t)$ , for  $\theta_i \stackrel{\text{ind}}{\sim} N(0, i^{-1-2\alpha})$
- $t \mapsto G_{\alpha t}$ , for Gaussian process  $G$
- $t \mapsto \int \alpha^{-1} \phi(\alpha^{-1}(t - z)) dF(z)$ , with  $F \sim$  Dirichlet process
- $i \mapsto \theta_i \sim \alpha \delta_0 + (1 - \alpha) \text{Laplace}$



# Bayesian adaptation

Family of priors  $\Pi_\alpha$  of varying smoothness  $\alpha$ ; posteriors  $\Pi_\alpha(\cdot | X)$

## Examples

- $t \mapsto \sum_{i=1}^{\infty} \theta_i e_i(t)$ , for  $\theta_i \stackrel{\text{ind}}{\sim} N(0, i^{-1-2\alpha})$
- $t \mapsto G_{\alpha t}$ , for Gaussian process  $G$
- $t \mapsto \int \alpha^{-1} \phi(\alpha^{-1}(t - z)) dF(z)$ , with  $F \sim$  Dirichlet process
- $i \mapsto \theta_i \sim \alpha \delta_0 + (1 - \alpha) \text{Laplace}$

## Hierarchical Bayes:

- Prior on  $\alpha$
- Ordinary posterior

## Empirical Bayes:

- $\hat{\alpha} = \operatorname{argmax}_\alpha \int p(X | \theta) d\Pi_\alpha(\theta)$
- Plug-in posterior  $\Pi_{\hat{\alpha}}(\cdot | X)$

Both methods give adaptive reconstructions:  
for smoother true function better reconstruction

# Bayesian adaptation

Family of priors  $\Pi_\alpha$  of varying smoothness  $\alpha$ ; posteriors  $\Pi_\alpha(\cdot | X)$

## Examples

- $t \mapsto \sum_{i=1}^{\infty} \theta_i e_i(t)$ , for  $\theta_i \stackrel{\text{ind}}{\sim} N(0, i^{-1-2\alpha})$
- $t \mapsto G_{\alpha t}$ , for Gaussian process  $G$
- $t \mapsto \int \alpha^{-1} \phi(\alpha^{-1}(t - z)) dF(z)$ , with  $F \sim$  Dirichlet process
- $i \mapsto \theta_i \sim \alpha \delta_0 + (1 - \alpha) \text{Laplace}$

## Hierarchical Bayes:

- Prior on  $\alpha$
- Ordinary posterior

## Empirical Bayes:

- $\hat{\alpha} = \operatorname{argmax}_\alpha \int p(X | \theta) d\Pi_\alpha(\theta)$
- Plug-in posterior  $\Pi_{\hat{\alpha}}(\cdot | X)$

Both methods give **adaptive reconstructions**:  
for smoother true function better reconstruction

This implies that they **cannot** give **honest confidence sets**

**Def**  $C_n(X^{(n)})$  is a **(honest) confidence set** over a model  $\Theta$  if

$$P_{\theta_0}(C_n(X^{(n)}) \ni \theta_0) \geq 0.95, \quad \forall \theta_0 \in \Theta$$

**Def**  $C_n(X^{(n)})$  is a **(honest) confidence set** over a model  $\Theta$  if

$$P_{\theta_0}(C_n(X^{(n)}) \ni \theta_0) \geq 0.95, \quad \forall \theta_0 \in \Theta$$

**Thm** The diameter of  $C_n(X^{(n)})$  cannot be smaller, **uniformly in**  $\theta \in \Theta_1 \subset \Theta$ , than:

(a)  $\varepsilon_n$  such that, for any  $T_n$ ,

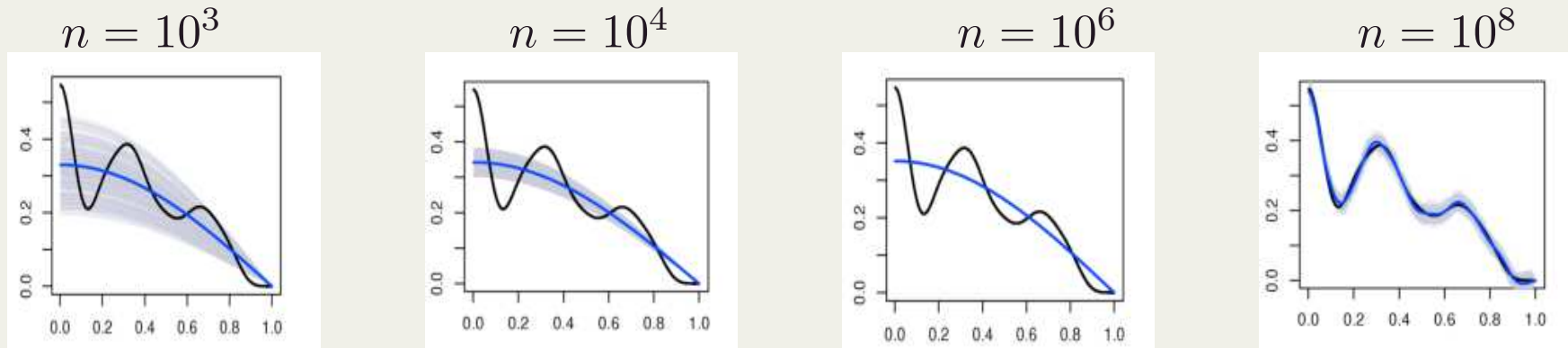
$$\liminf_{n \rightarrow \infty} \sup_{\theta \in \Theta_1} P_{\theta}(d(T_n, \theta) \geq \varepsilon_n) > 0.501$$

(b) rate  $\varepsilon_n$  of minimax testing, for any given  $\Theta'_1 \subset \Theta_1$  of

$H_0: \theta \in \Theta'_1$  versus  $H_1: \theta \in \Theta, d(\theta, \Theta'_1) > \varepsilon_n$

- (a) typically gives minimax rate of estimation for model  $\Theta_1$   
(b) is determined by biggest model  $\Theta$  rather than  $\Theta_1$

# \* Credible balls — counter example — reconstructing a derivative



Gaussian prior in white noise model of smoothness determined by empirical Bayes

Black: true curve. Blue: posterior mean. Grey: draws from posterior

The pictures show an *inconvenient truth*  
For some (most?) truths the results are good

## \*\* Credible balls — counter example — reconstructing a derivative

**Data:**  $X^{(n)} = \int_0^\cdot \theta(t) dt + n^{-1/2} \dot{W}$ , for white noise  $\dot{W}$

- **Prior:**  $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , with  $\theta_i | \alpha \stackrel{\text{ind}}{\sim} N(0, i^{-2\alpha-1})$
- **Prior** on  $\alpha$  or empirical Bayes  $\hat{\alpha}$

**Thm** For  $n_j \geq n_{j-1}^4$  for every  $j$ , define  $\theta = (\theta_1, \theta_2, \dots)$  by

$$\theta_i^2 = \begin{cases} n_j^{-\frac{1+2\beta}{1+2\beta+2p}}, & \text{if } n_j^{\frac{1}{1+2\beta+2p}} \leq i < 2n_j^{\frac{1}{1+2\beta+2p}}, \quad j = 1, 2, \dots, \\ 0, & \text{otherwise} \end{cases}$$

Then  $\sum_j j^{2\beta} \theta_j^2 \leq 1$ , but the central 95%-credible ball  $\hat{C}_n$ , blown up by  $L_n \ll n^\delta$ , satisfies

$$\liminf P_\theta(\theta \in \hat{C}_n) = 0$$

- Data allows inference only on  $\theta_1, \dots, \theta_{N_n}$
- Trouble if  $\theta_1, \dots, \theta_{N_n}$  does not resemble  $\theta_1, \theta_2, \dots$
- Example  $\theta$  has repeated runs of 0s of increasing lengths

# Estimation versus uncertainty quantification

## Adaptive estimation:

- Estimators can be simultaneously optimal for multiple regularities
- (Bayesian procedures are natural)

## Uncertainty quantification:

- Size of honest confidence set is determined by smallest considered regularity
- (Data-driven constructions can be misleading)

# Estimation versus uncertainty quantification

## Adaptive estimation:

- Estimators can be simultaneously optimal for multiple regularities
- (Bayesian procedures are natural)

## Uncertainty quantification:

- Size of honest confidence set is determined by smallest considered regularity
- (Data-driven constructions can be misleading)

SOLUTION 1: *be honest*

make conditional confidence statements



# Estimation versus uncertainty quantification

## Adaptive estimation:

- Estimators can be simultaneously optimal for multiple regularities
- (Bayesian procedures are natural)

## Uncertainty quantification:

- Size of honest confidence set is determined by smallest considered regularity
- (Data-driven constructions can be misleading)

SOLUTION 1: *be honest*

make conditional confidence statements

SOLUTION 2: *believe your prior*



# Estimation versus uncertainty quantification

## Adaptive estimation:

- Estimators can be simultaneously optimal for multiple regularities
- (Bayesian procedures are natural)

## Uncertainty quantification:

- Size of honest confidence set is determined by smallest considered regularity
- (Data-driven constructions can be misleading)

SOLUTION 1: *be honest*  
make conditional confidence statements

SOLUTION 2: believe your prior



SOLUTION 3: determine which  $\theta$  cause the trouble  
argue that these are implausible

# Uncertainty quantification for curve estimation

**Def**  $\theta \in \ell^2$  satisfies the *polished tail condition* if

$$\sum_{i=N}^{1000N} \theta_i^2 \geq 0.001 \sum_{i=N}^{\infty} \theta_i^2, \quad \forall \text{ large } N$$

**Interpretation:**

every block of frequencies  $(N, 1000N)$   
contains a fraction of the total energy above frequency  $N$

**Def**  $\theta \in \ell^2$  satisfies the *polished tail condition* if

$$\sum_{i=N}^{1000N} \theta_i^2 \geq 0.001 \sum_{i=N}^{\infty} \theta_i^2, \quad \forall \text{ large } N$$

“Everything” is polished tail...:

**Def**  $\theta \in \ell^2$  satisfies the *polished tail condition* if

$$\sum_{i=N}^{1000N} \theta_i^2 \geq 0.001 \sum_{i=N}^{\infty} \theta_i^2, \quad \forall \text{ large } N$$

“Everything” is polished tail...:

- For the *topologist* Giné+Nickl, 2010  
Non polished tail sequences are meagre in a natural topology

**Def**  $\theta \in \ell^2$  satisfies the *polished tail condition* if

$$\sum_{i=N}^{1000N} \theta_i^2 \geq 0.001 \sum_{i=N}^{\infty} \theta_i^2, \quad \forall \text{ large } N$$

“Everything” is polished tail...:

- For the *topologist* Giné+Nickl, 2010  
Non polished tail sequences are meagre in a natural topology
- For the *minimax expert*:  
Intersecting the usual models with polished tail sequences decreases the minimax risk by at most a logarithmic factor

**Def**  $\theta \in \ell^2$  satisfies the *polished tail condition* if

$$\sum_{i=N}^{1000N} \theta_i^2 \geq 0.001 \sum_{i=N}^{\infty} \theta_i^2, \quad \forall \text{ large } N$$

“Everything” is polished tail...:

- For the *topologist* Giné+Nickl, 2010  
Non polished tail sequences are meagre in a natural topology
- For the *minimax expert*:  
Intersecting the usual models with polished tail sequences decreases the minimax risk by at most a logarithmic factor
- For the *Bayesian*:  
Almost every parameter generated from a prior  $\theta_i \stackrel{\text{ind}}{\sim} N(0, ci^{-\alpha-1/2})$  is polished tail



# Credible balls in linear Gaussian inverse problems

**Data:**  $X^{(n)} = K\theta + n^{-1/2}\dot{W}$ , for white noise  $\dot{W}$

- $K$  compact operator with eigenvalues  $\kappa_i \asymp i^{-p}$  and eigen basis  $(e_i)$
- **Prior:**  $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , with  $\theta_i | \alpha \stackrel{\text{ind}}{\sim} N(0, i^{-2\alpha-1})$
- **Prior** on  $\alpha$

# Credible balls in linear Gaussian inverse problems

**Data:**  $X^{(n)} = K\theta + n^{-1/2}\dot{W}$ , for white noise  $\dot{W}$

- $K$  compact operator with eigenvalues  $\kappa_i \asymp i^{-p}$  and eigen basis  $(e_i)$
- **Prior:**  $\theta = \sum_{i=1}^{\infty} \theta_i e_i$ , with  $\theta_i | \alpha \stackrel{\text{ind}}{\sim} N(0, i^{-2\alpha-1})$
- **Prior** on  $\alpha$

Credible ball:

$$\hat{C}_n(M) := \{\theta: \|\theta - \hat{\theta}_n\| < Mr\}$$

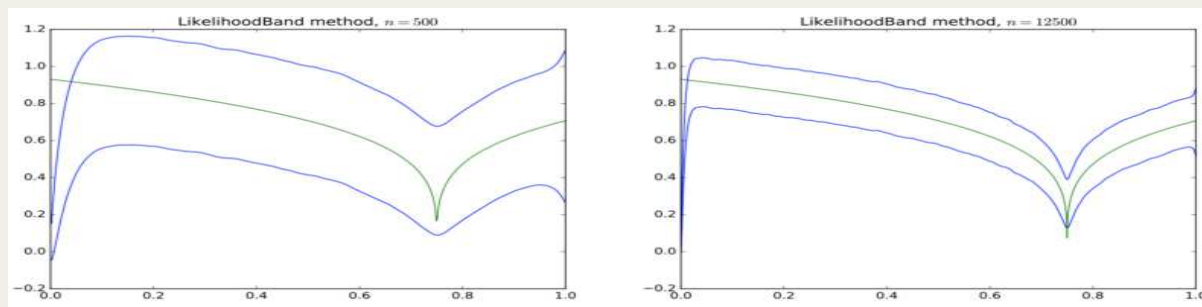
$$\begin{aligned} \hat{\theta}_n &= \mathbb{E}(\theta | X^{(n)}) \\ \Pi(\theta: \|\theta - \hat{\theta}_n\| < r | X^{(n)}) &= 0.95 \end{aligned}$$

**Thm** For not too small  $M$ , uniformly in polished tail functions  $\theta$ ,

$$P_{\theta}(\theta \in \hat{C}_n(M)) \rightarrow 1$$

Similar results for empirical Bayes

# Credible bands and other models



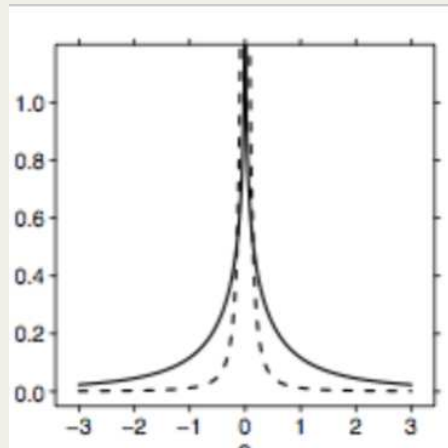
- Rousseau & Szabó, 2017-20: empirical Bayes and credible balls for general models
- Yoo 2017: 'Bayesian Lepski method' for adaptive credible bands; spline and wavelet priors
- Sniekers & vdV 2017, 19: bands for scaled Gaussian prior under 'self-similarity', 'good bias' and 'discrete polished tail'.
- Belitser & Nurushev 2015-19: general projection estimators; 'excessive-bias restriction'
- Ray 2017: intersect credible set from weak and strong norms
- Hadji& Szabó, 2019: supersmooth priors
- ...

# Uncertainty quantification for sparse high-dimensional parameters

**Data:**  $Y^n \sim N_n(\theta, I)$ , for  $\theta = (\theta_1, \dots, \theta_n) \in \mathbb{R}^n$

Constructive definition of prior  $\Pi$  for  $\theta \in \mathbb{R}^p$ :

- (1) Choose “sparsity level”  $\tau$  from prior or by empirical Bayes
- (2) Generate  $\sqrt{\psi_1}, \dots, \sqrt{\psi_n}$  iid from  $\text{Cauchy}^+(0, \tau)$
- (3) Generate independent  $\theta_i \sim N(0, \psi_i)$



prior density of  $\theta_i$

This prior gives optimal recovery of sparse vectors  $\theta$

## Credible interval:

$$\hat{C}_{ni}(L) = \left\{ \theta_i : |\theta_i - \hat{\theta}_i| \leq L \hat{r}_i \right\}$$

$$\begin{aligned} \hat{\theta} &= \mathbf{E}(\theta | Y^n) \\ \Pi(\theta_i : |\theta_i - \hat{\theta}_i| \leq \hat{r}_i | Y^n) &= 0.95 \end{aligned}$$

## Credible interval:

$$\hat{C}_{ni}(L) = \left\{ \theta_i : |\theta_i - \hat{\theta}_i| \leq L \hat{r}_i \right\}$$

$$\begin{aligned} \hat{\theta} &= \mathbf{E}(\theta | Y^n) \\ \Pi(\theta_i : |\theta_i - \hat{\theta}_i| \leq \hat{r}_i | Y^n) &= 0.95 \end{aligned}$$

$$\mathbf{S}_a := \{1 \leq i \leq n : |\theta_{0,i}| \leq 1/n\}$$

$$\mathbf{M}_a := \{1 \leq i \leq n : (s_n/n) \sqrt{\log(n/s_n)} \ll |\theta_{0,i}| \leq 0.99 \sqrt{2 \log(n/s_n)}\}$$

$$\mathbf{L}_a := \{1 \leq i \leq n : 1.001 \sqrt{2 \log n} \leq |\theta_{0,i}|\}$$

## Credible interval:

$$\hat{C}_{ni}(L) = \left\{ \theta_i : |\theta_i - \hat{\theta}_i| \leq L \hat{r}_i \right\}$$

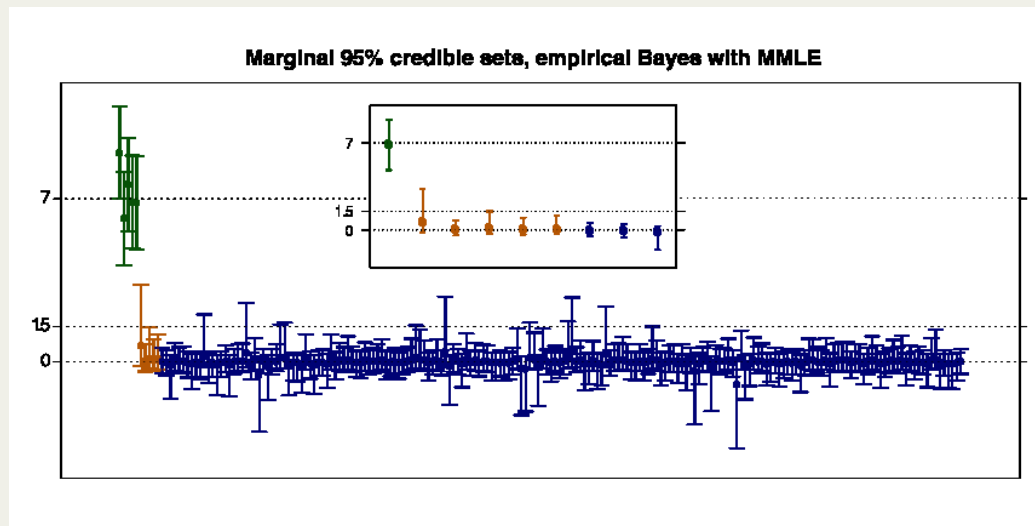
$$\hat{\theta} = \mathbb{E}(\theta | Y^n)$$

$$\Pi(\theta_i : |\theta_i - \hat{\theta}_i| \leq \hat{r}_i | Y^n) = 0.95$$

$$\mathbf{S}_a := \{1 \leq i \leq n : |\theta_{0,i}| \leq 1/n\}$$

$$\mathbf{M}_a := \{1 \leq i \leq n : (s_n/n) \sqrt{\log(n/s_n)} \ll |\theta_{0,i}| \leq 0.99 \sqrt{2 \log(n/s_n)}\}$$

$$\mathbf{L}_a := \{1 \leq i \leq n : 1.001 \sqrt{2 \log n} \leq |\theta_{0,i}|\}$$



marginal credible intervals for a single  $Y^n$  with  $n = 200$  and  $s_n = 10$

$\theta_1 = \dots = \theta_5 = 7, \theta_6 = \dots = \theta_{10} = 1.5$ . Insert: credible sets 5 to 13



## Credible interval:

$$\hat{C}_{ni}(L) = \left\{ \theta_i : |\theta_i - \hat{\theta}_i| \leq L \hat{r}_i \right\}$$

$$\begin{aligned} \hat{\theta} &= \mathbb{E}(\theta | Y^n) \\ \Pi(\theta_i : |\theta_i - \hat{\theta}_i| \leq \hat{r}_i | Y^n) &= 0.95 \end{aligned}$$

$$\mathbf{S}_a := \{1 \leq i \leq n : |\theta_{0,i}| \leq 1/n\}$$

$$\mathbf{M}_a := \{1 \leq i \leq n : (s_n/n) \sqrt{\log(n/s_n)} \ll |\theta_{0,i}| \leq 0.99 \sqrt{2 \log(n/s_n)}\}$$

$$\mathbf{L}_a := \{1 \leq i \leq n : 1.001 \sqrt{2 \log n} \leq |\theta_{0,i}|\}$$

**Thm** For any  $\gamma > 0$  and  $\|\theta_0\|_0 \leq s_n$ ,

$$P_{\theta_0} \left( \frac{1}{\#\mathbf{S}_a} \#\{i \in \mathbf{S}_a : \theta_{0,i} \in \hat{C}_{ni}(L_{S,\gamma})\} \geq 1 - \gamma \right) \rightarrow 1,$$

$$P_{\theta_0}(\theta_{0,i} \notin \hat{C}_{ni}(L)) \rightarrow 1, \quad \text{for any } i \in \mathbf{M}_a \quad \text{and any } L$$

$$P_{\theta_0} \left( \frac{1}{\#\mathbf{L}_a} \#\{i \in \mathbf{L}_a : \theta_{0,i} \in \hat{C}_{ni}(L_{L,\gamma})\} \geq 1 - \gamma \right) \rightarrow 1$$

Few false discoveries. Most easy discoveries made  
Intermediate discoveries not made

**Data:**  $Y^n \sim N_n(\theta, I)$ , for  $\theta = (\theta_1, \dots, \theta_n) \in \mathbb{R}^n$

- $\theta_i \stackrel{\text{iid}}{\sim} (1 - \tau)\delta_0 + \tau G$ , with  $G = \text{Laplace or Cauchy}$
- $\hat{\tau}$  determined by marginal empirical Bayes

$$\ell_\tau(x) = \Pi_\tau(\theta_i = 0 \mid X_i = x),$$

$$q_\tau(x) = \Pi_\tau(\theta_i = 0 \mid |X_i| \geq |x|)$$

Tests: Reject  $H_0: \theta_{0,i} = 0$  if  $\ell_{\hat{\tau}}(X_i) \leq t$  or  $q_{\hat{\tau}}(X_i) \leq t$

**Thm** For  $s_n \rightarrow \infty$  with  $s_n \ll n^\nu$ ,

$$\sup_{\#(\theta_{0,i} \neq 0) \leq s_n} \mathbb{E}_{\theta_0} \frac{\#(i: \theta_{0,i} = 0, \text{ rejected})}{\#(i: \text{ rejected}) \vee 1} \lesssim t \log \frac{1}{t}$$

## \* Credible balls

### A confidence ball

$$C_n(Y^n) = \{\theta \in \mathbb{R}^n: \|\theta - \hat{\theta}\| \leq \hat{r}\}$$

cannot have both:

- radius  $\hat{r}$  of order the adaptive benchmark  $\sqrt{s \log(n/s)}$  for sparsity,
- uniform coverage over multiple sparsity levels  $s$

### Meta Thm

A credible ball will cover “self-similar” parameters

## \*\* Simultaneous credible balls — impossibility of adaptation

General principle:  
size of **honest confidence set** is determined by **biggest model**

**Thm** [Li, 1987]

If  $P_{\theta_0}(C_n(Y^n) \ni \theta_0) \geq 0.95$ , all  $\theta_0 \in \mathbb{R}^n$ , then  $\text{diam}(C_n(Y^n)) \gtrsim n^{-1/4}$ , some  $\theta_0$

**Thm** [Nickl, van de Geer, 2013]

If  $s_{1,n} \ll s_{2,n}$  and  
 $\text{diam}(C_n(Y^n))$  is of optimal size, uniformly in  $\|\theta_0\|_0 \leq s_{i,n}$  for  $i = 1, 2$ ,  
then  $C_n(Y^n)$  cannot have uniform coverage over  $\{\theta_0: \|\theta_0\|_0 \leq s_{2,n}\}$ .

Since the Bayesian procedure adapts to sparsity,  
its credible balls *cannot* be honest confidence sets

[Optimal size is  $((s_{i,n}/n) \log(n/s_{i,n}))^{1/2}$ ]

## \*\* Simultaneous credible balls — impossibility of adaptation — restricting the parameter

Coverage only when  $\theta_0$  does not cause too much shrinkage

**Def** [self-similarity]

For  $s = \|\theta_0\|_0$  at least  $0.001s$  coordinates of  $\theta_0$  satisfy

$$|\theta_{0,i}| \geq 1.001 \sqrt{2 \log(n/s)}.$$

## \*\* Simultaneous credible balls — impossibility of adaptation — restricting the parameter

Coverage only when  $\theta_0$  does not cause too much shrinkage

**Def** [self-similarity]

For  $s = \|\theta_0\|_0$  at least  $0.001s$  coordinates of  $\theta_0$  satisfy

$$|\theta_{0,i}| \geq 1.001 \sqrt{2 \log(n/s)}.$$

**Def** [excessive-bias restriction, Belitser & Nurushev, 2015]

$\|\theta\|_0 \leq s$  and  $\exists \tilde{s}$  with  $\tilde{s} \asymp \#\{i: |\theta_{0,i}| \geq 1.001 \sqrt{2 \log(n/\tilde{s})}\}$  and

$$\sum_{i: |\theta_{0,i}| \leq 1.001 \sqrt{2 \log(n/\tilde{s})}} \theta_{0,i}^2 \lesssim \tilde{s} \log(n/\tilde{s})$$

Excessive-bias restriction weaker than self-similarity  
(Self-similarity allows to tighten up the sets  $S, M, L$ )

## \*\* Simultaneous credible balls

Credible ball:

$$\hat{C}_n(L) = \{\theta: \|\theta - \hat{\theta}\| \leq L\hat{r}\}$$

$$\begin{aligned} \hat{\theta} &= \mathbb{E}(\theta | Y^n) \\ \Pi(\theta: \|\theta - \hat{\theta}\| \leq \hat{r} | Y^n) &= 0.95 \end{aligned}$$

**Thm** If  $s_n/n \rightarrow 0$ , for sufficiently large  $L$ ,

$$\liminf_{n \rightarrow \infty} \inf_{\theta_0 \in \text{EBR}[s_n]} P_{\theta_0}(\theta_0 \in \hat{C}_n(L)) \geq 1 - \alpha$$

EBR $[s]$ : vectors  $\theta_0$  that satisfy excessive bias restriction

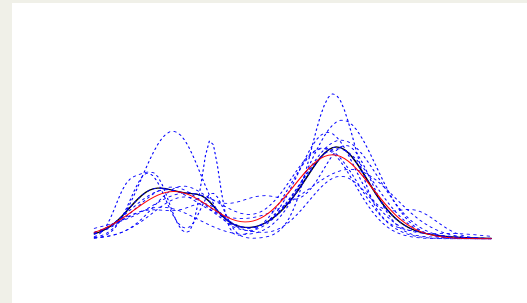
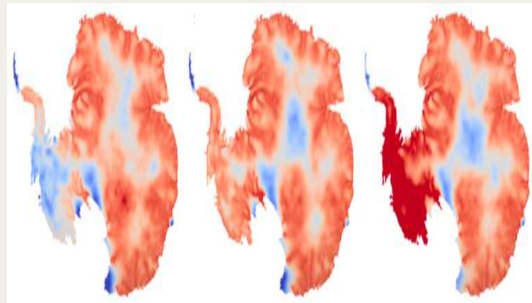
Closing remarks



## Closing remarks

In nonparametric statistics uncertainty quantification is problematic for both Bayesian and non-Bayesian methods

*It necessarily extrapolates into features of the world that cannot be seen in the data*



Adaptive methods seem reasonable, even though their confidence sets are **dishonest**