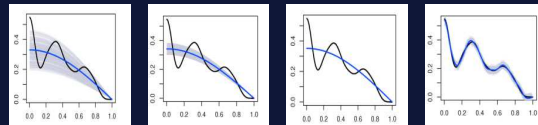


Nonparametric Bayesian Uncertainty Quantification

Lecture 1: Curve estimation

Aad van der Vaart

Universiteit Leiden, Netherlands



Hotelling Lectures, UNC, Chapel Hill, March 2017

Introduction

Recovery

Gaussian process priors

Dirichlet process mixtures

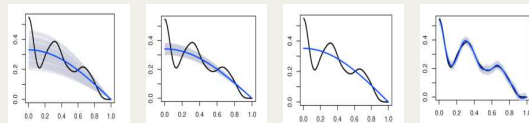
Uncertainty quantification

Priors of fixed regularity

Priors of flexible regularity

Nonparametric regression

Closing remarks



Introduction

The Bayesian paradigm



- A parameter θ is generated according to a **prior distribution** Π .
- Given θ the data X is generated according to a measure P_θ .

This gives a **joint distribution** of (X, θ) .

- Given observed data x the statistician computes the conditional distribution of θ given $X = x$, the **posterior distribution**:

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

The Bayesian paradigm



- A parameter θ is generated according to a **prior distribution** Π .
- Given θ the data X is generated according to a measure P_θ .

This gives a **joint distribution** of (X, θ) .

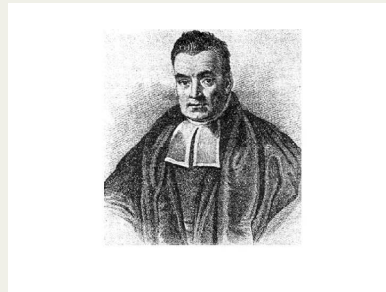
- Given observed data x the statistician computes the conditional distribution of θ given $X = x$, the **posterior distribution**:

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

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

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

Reverend Thomas

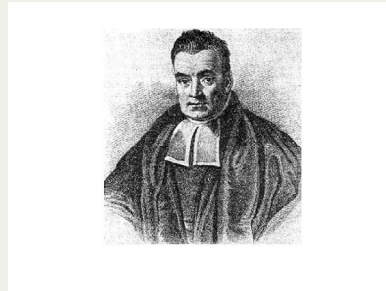


Thomas Bayes (1702–1761, 1763) followed this argument with Θ possessing the *uniform* distribution and X given $\Theta = \theta$ *binomial* (n, θ) .

Using his famous rule he computed that the posterior distribution is then *Beta* $(X + 1, n - X + 1)$.

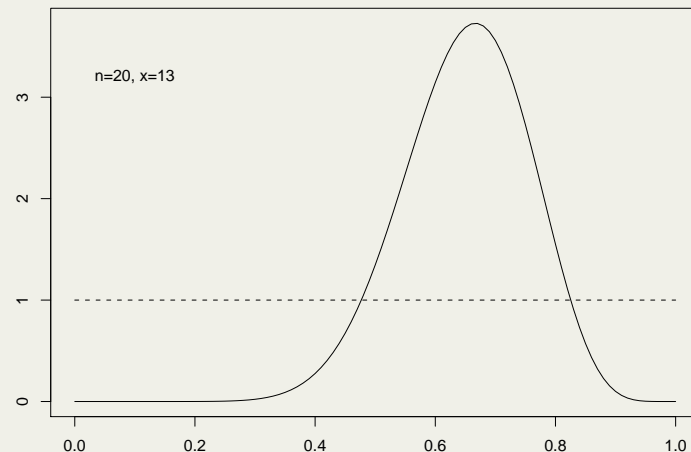
$$P(a \leq \Theta \leq b) = b - a, \quad 0 < a < b < 1,$$
$$P(X = x | \Theta = \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}, \quad x = 0, 1, \dots, n,$$
$$d\Pi(\theta | X) = \theta^X (1 - \theta)^{n-X} \cdot 1.$$

Reverend Thomas

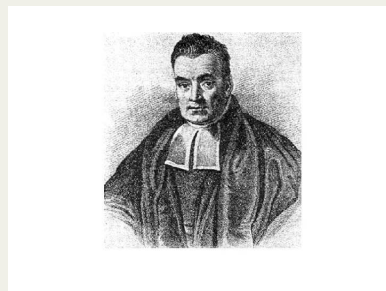


Thomas Bayes (1702–1761, 1763) followed this argument with Θ possessing the *uniform* distribution and X given $\Theta = \theta$ *binomial* (n, θ) .

Using his famous rule he computed that the posterior distribution is then *Beta* $(X + 1, n - X + 1)$.

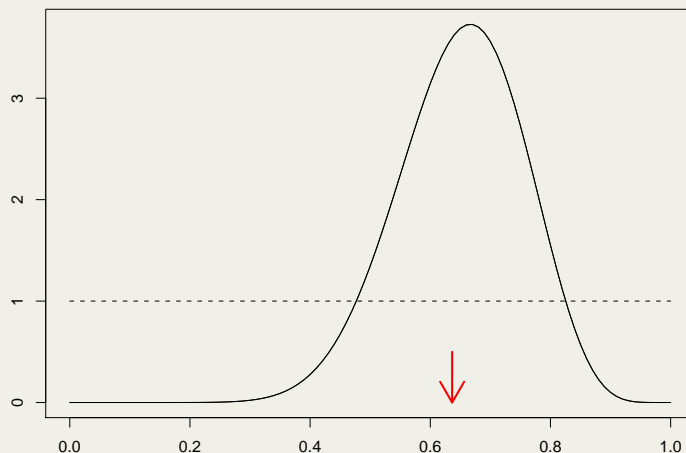


Reverend Thomas



Thomas Bayes (1702–1761, 1763) followed this argument with Θ possessing the *uniform* distribution and X given $\Theta = \theta$ *binomial* (n, θ) .

Using his famous rule he computed that the posterior distribution is then *Beta* $(X + 1, n - X + 1)$.

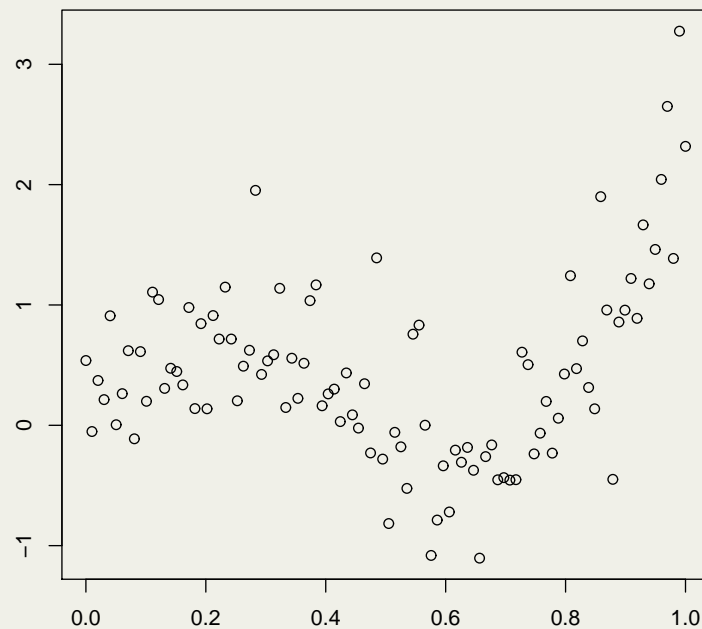


Nonparametric Bayes

If the parameter θ is a function, then the prior is a **probability distribution on an function space**. So is the posterior, given the data. Bayes's formula does not change:

$$d\Pi(\theta | X) \propto p_{\theta}(X) d\Pi(\theta).$$

Prior and posterior can be visualized by plotting functions that are simulated from these distributions.

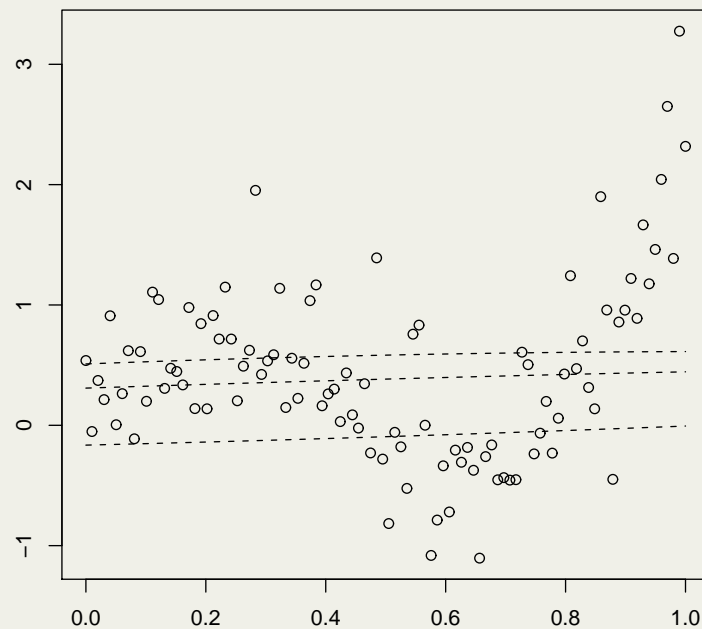


Nonparametric Bayes

If the parameter θ is a function, then the prior is a **probability distribution on an function space**. So is the posterior, given the data. Bayes's formula does not change:

$$d\Pi(\theta | X) \propto p_{\theta}(X) d\Pi(\theta).$$

Prior and posterior can be visualized by plotting functions that are simulated from these distributions.

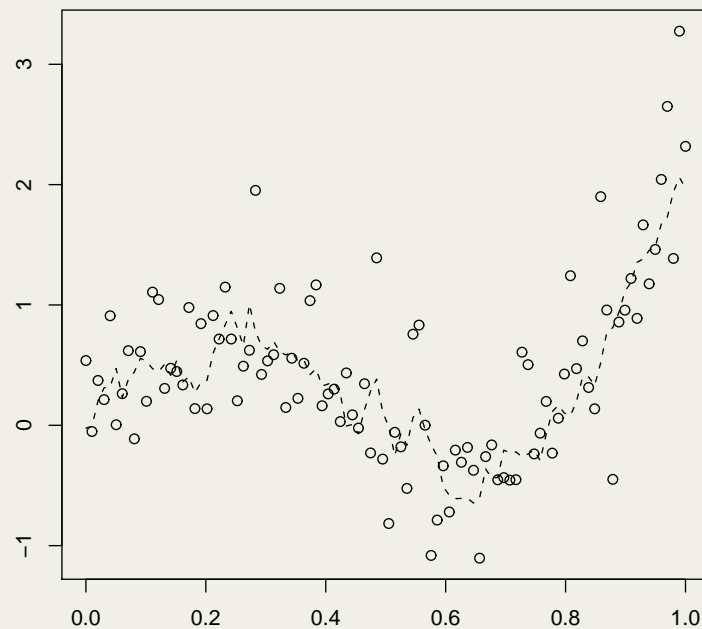


Nonparametric Bayes

If the parameter θ is a function, then the prior is a **probability distribution on an function space**. So is the posterior, given the data. Bayes's formula does not change:

$$d\Pi(\theta | X) \propto p_{\theta}(X) d\Pi(\theta).$$

Prior and posterior can be visualized by plotting functions that are simulated from these distributions.

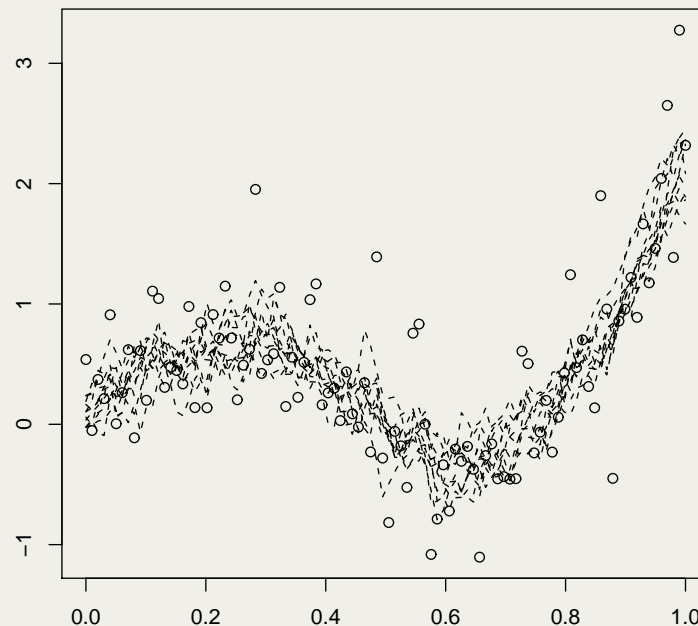


Nonparametric Bayes

If the parameter θ is a function, then the prior is a **probability distribution on an function space**. So is the posterior, given the data. Bayes's formula does not change:

$$d\Pi(\theta | X) \propto p_{\theta}(X) d\Pi(\theta).$$

Prior and posterior can be visualized by plotting functions that are simulated from these distributions.

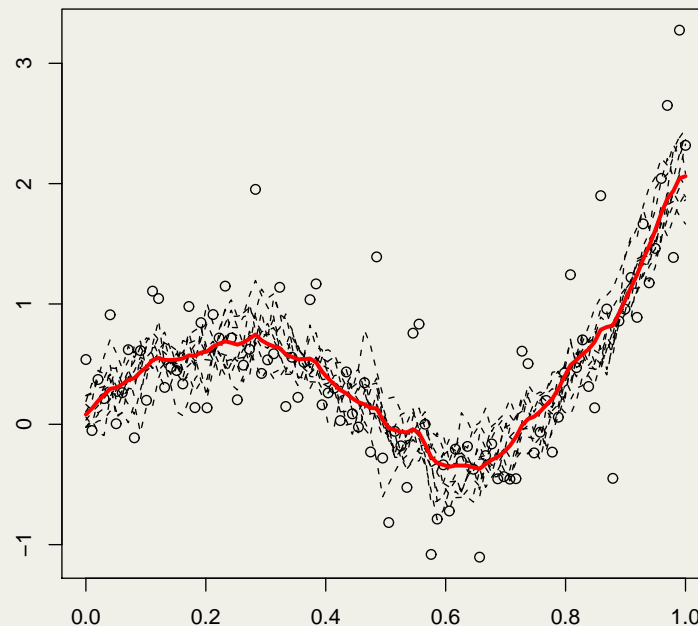


Nonparametric Bayes

If the parameter θ is a function, then the prior is a **probability distribution on an function space**. So is the posterior, given the data. Bayes's formula does not change:

$$d\Pi(\theta | X) \propto p_{\theta}(X) d\Pi(\theta).$$

Prior and posterior can be visualized by plotting functions that are simulated from these distributions.

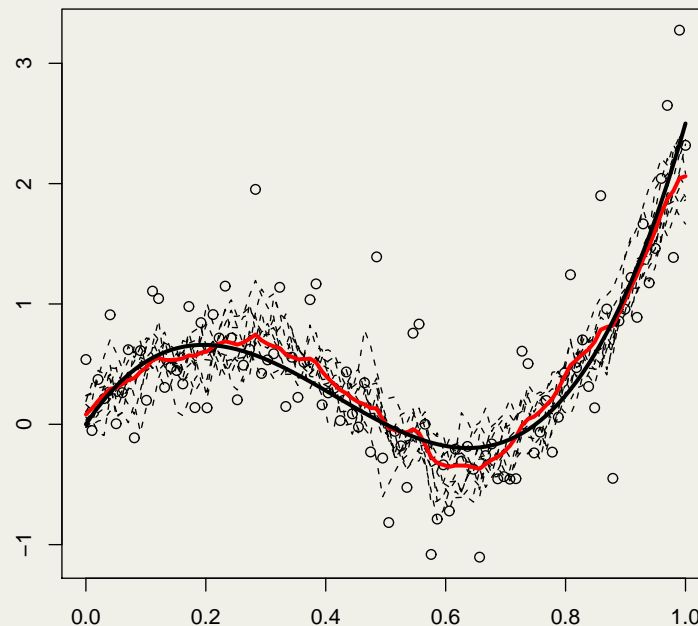


Nonparametric Bayes

If the parameter θ is a function, then the prior is a **probability distribution on an function space**. So is the posterior, given the data. Bayes's formula does not change:

$$d\Pi(\theta | X) \propto p_{\theta}(X) d\Pi(\theta).$$

Prior and posterior can be visualized by plotting functions that are simulated from these distributions.



Frequentist Bayesian

Assume that the data X is generated according to a **given parameter** θ_0 and consider the posterior $\Pi(\theta \in \cdot | X)$ as a random measure on the parameter set dependent on X .

RECOVERY

We like $\Pi(\theta \in \cdot | X)$ to put “most” of its mass near θ_0 for “most” X .

Frequentist Bayesian

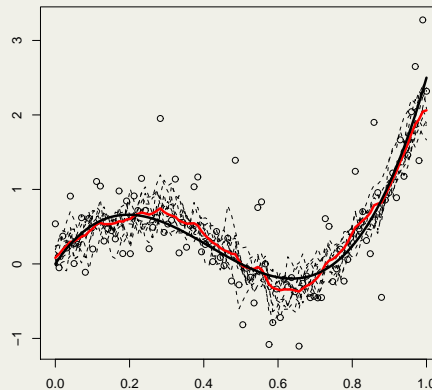
Assume that the data X is generated according to a **given parameter** θ_0 and consider the posterior $\Pi(\theta \in \cdot | X)$ as a random measure on the parameter set dependent on X .

RECOVERY

We like $\Pi(\theta \in \cdot | X)$ to put “most” of its mass near θ_0 for “most” X .

UNCERTAINTY QUANTIFICATION

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



Frequentist Bayesian

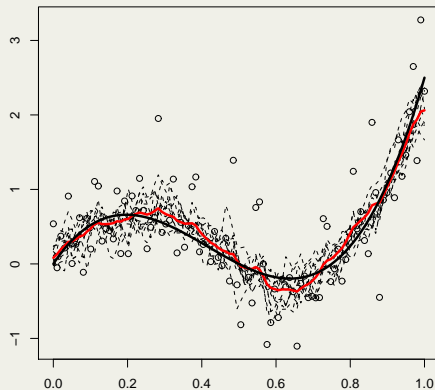
Assume that the data X is generated according to a **given parameter** θ_0 and consider the posterior $\Pi(\theta \in \cdot | X)$ as a random measure on the parameter set dependent on X .

RECOVERY

We like $\Pi(\theta \in \cdot | X)$ to put “most” of its mass near θ_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 the *coverage* of a set of large posterior mass to be large.

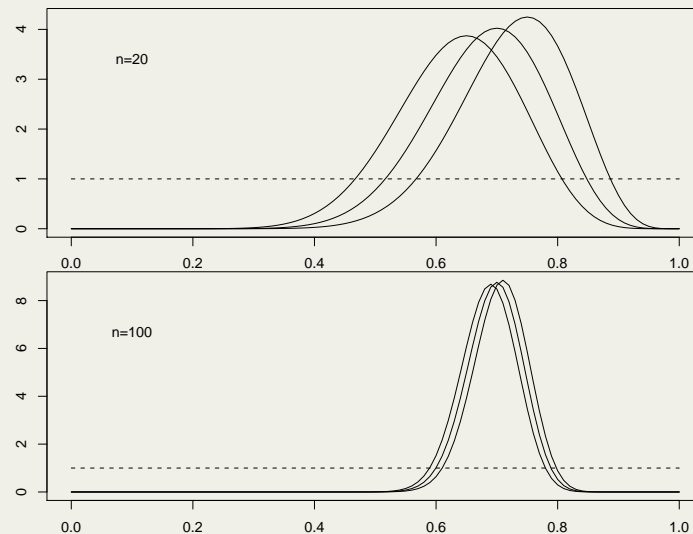
Parametric models Laplace, Bernstein, von Mises, Le Cam 1989

Suppose the data are a random sample X_1, \dots, X_n from a density $x \mapsto p_\theta(x)$ that is smoothly and **identifiably** parametrized by a vector $\theta \in \mathbb{R}^d$ (e.g. $\theta \mapsto \sqrt{p_\theta}$ continuously differentiable as map in $L_2(\mu)$).

Theorem. Under $P_{\theta_0}^n$, for *any prior* with positive density,

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

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



Parametric models Laplace, Bernstein, von Mises, Le Cam 1989

Suppose the data are a random sample X_1, \dots, X_n from a density $x \mapsto p_\theta(x)$ that is smoothly and **identifiably** parametrized by a vector $\theta \in \mathbb{R}^d$ (e.g. $\theta \mapsto \sqrt{p_\theta}$ continuously differentiable as map in $L_2(\mu)$).

Theorem. Under $P_{\theta_0}^n$, for *any prior* with positive density,

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

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 θ_0 .

UNCERTAINTY QUANTIFICATION:

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

These lectures

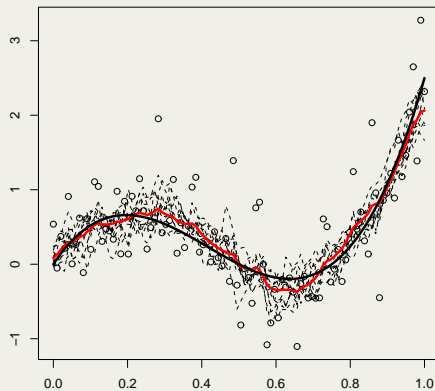
Recovery and uncertainty quantification for **nonparametric** models.

LECTURE 1: Curve fitting.

LECTURE 2: High dimensional inference and sparsity.

Point of view:

How does the posterior distribution for **natural priors** behave, in particular for priors that **adapt** to complexity in the data.



Consistency

- $X^{(n)}$ observation in sample space $(\mathcal{X}^{(n)}, \mathcal{X}^{(n)})$ with distribution $P_\theta^{(n)}$.
- θ belongs to metric space (Θ, d) .

Definition. *Posterior consistency at θ_0 means that for every $\epsilon > 0$,*

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

The main result on consistency is **Schwartz's theorem (1965)**.
This was adapted to nonparametric estimation in the 1990s.

Rate of contraction

- $X^{(n)}$ observation in sample space $(\mathfrak{X}^{(n)}, \mathcal{X}^{(n)})$ with distribution $P_\theta^{(n)}$.
- θ belongs to metric space (Θ, d) .

Definition. *The posterior contraction rate at θ_0 is $\epsilon_n \rightarrow 0$ such that, for every $M_n \rightarrow \infty$,*

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

Rate of contraction

- $X^{(n)}$ observation in sample space $(\mathcal{X}^{(n)}, \mathcal{X}^{(n)})$ with distribution $P_\theta^{(n)}$.
- θ belongs to metric space (Θ, d) .

Definition. The posterior contraction rate at θ_0 is $\epsilon_n \rightarrow 0$ such that, for every $M_n \rightarrow \infty$,

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

Benchmark rate for curve fitting: A function θ of d variables that has bounded derivatives of order β is estimable based on n observations at rate

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

Proposition. If the contraction rate at θ_0 is ϵ_n , then the center $\hat{\theta}_n$ of a (nearly) smallest ball of posterior mass $\geq 1/2$ satisfies $d(\hat{\theta}_n, \theta_0) = O_P(\epsilon_n)$.

Basic contraction theorem (Ghosal, Ghosh, vdV 2000)

- $p \sim \Pi$, prior on set of densities \mathcal{P} .
- $X_1, \dots, X_n | p \stackrel{\text{iid}}{\sim} p$.

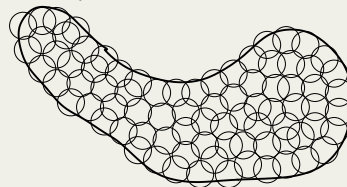
$$K(p_0; p) = P_0 \log \frac{p_0}{p}, \quad V(p_0; p) = P_0 \left(\log \frac{p_0}{p} \right)^2.$$

Theorem. Let d convex metric bounded above by Hellinger metric such that that there exist $\mathcal{P}_n \subset \mathcal{P}$ and $C > 0$ with

- (i) $\Pi_n(p: K(p_0; p) < \epsilon_n^2, V(p_0; p) < \epsilon_n^2) \geq e^{-Cn\epsilon_n^2}$, *(prior mass)*
- (ii) $\log N(\epsilon_n, \mathcal{P}_n, d) \leq n\epsilon_n^2$. *(complexity)*
- (iii) $\Pi_n(\mathcal{P}_n^c) \leq e^{-(C+4)n\epsilon_n^2}$.

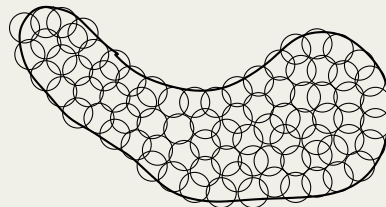
Then the posterior rate of contraction is $\epsilon_n \vee n^{-1/2}$.

The **covering number** $N(\epsilon, \mathcal{P}, d)$ is the minimal number of d -balls of radius ϵ needed to cover \mathcal{P} .



Interpretation

Let p_1, \dots, p_N in \mathcal{P} be a maximal set with $d(p_i, p_j) \geq \epsilon_n$.



Hence, under the complexity bound,

$$N \asymp N(\epsilon_n, \mathcal{P}, d) \geq e^{cn\epsilon_n^2}.$$

If prior mass were evenly distributed, then each ball of radius $\epsilon_n/2$ would have mass of order

$$1/N \leq e^{-cn\epsilon_n^2}.$$

This is the order of the prior mass bound.

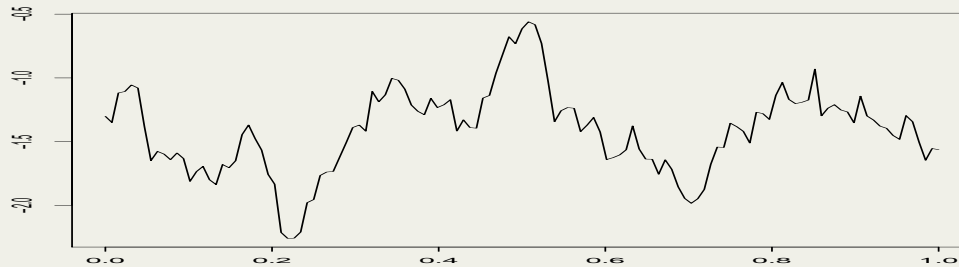
Suggestion:

The conditions can be satisfied for every $p_0 \in \mathcal{P}$ if the prior “*distributes its mass uniformly over \mathcal{P} , at discretization level ϵ_n* ”.

Gaussian process priors

Gaussian process prior

The law of a stochastic process $W = (W_t: t \in T)$ is a prior distribution on the space of functions $\theta: T \rightarrow \mathbb{R}$.



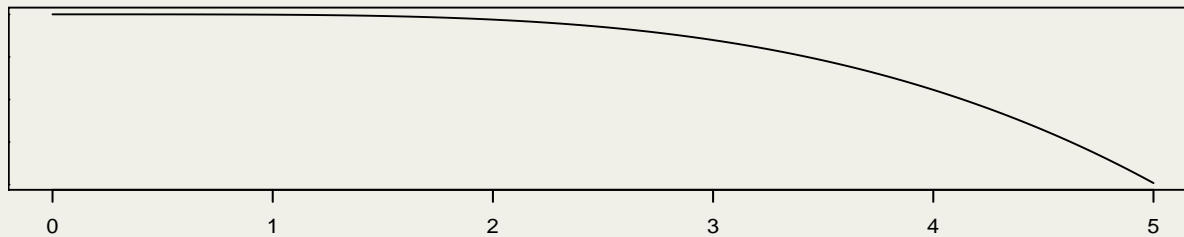
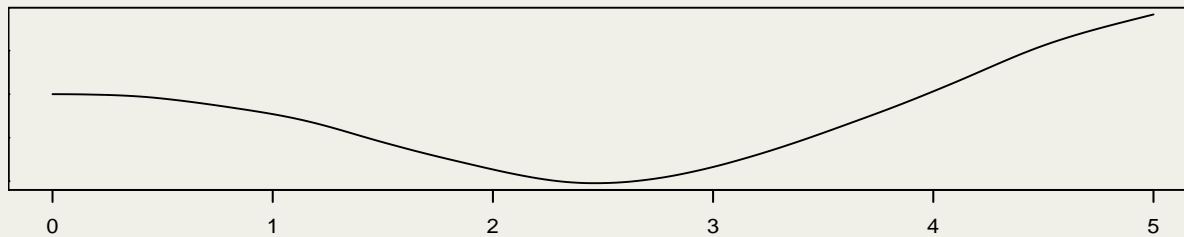
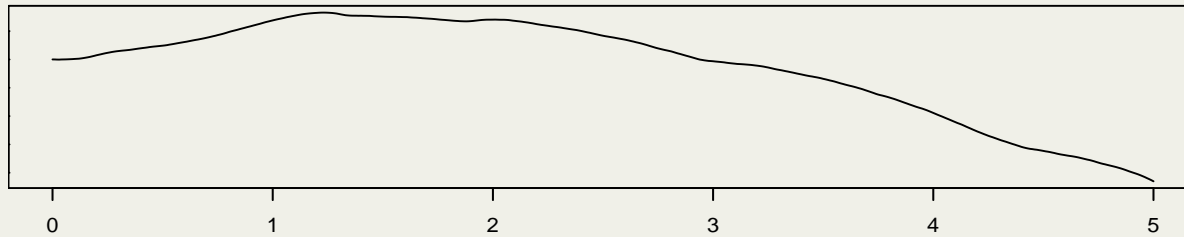
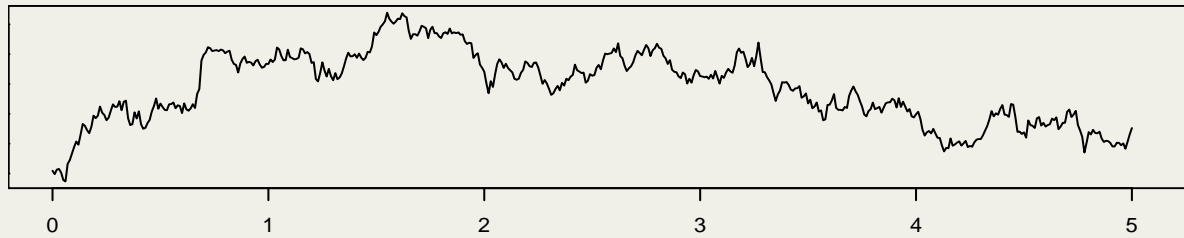
W is a **Gaussian process** if

$(W_{t_1}, \dots, W_{t_k})$ is multivariate Gaussian, for every t_1, \dots, t_k .

Mean and covariance function:

$$t \mapsto \mathbb{E}W_t, \quad \text{and} \quad (s, t) \mapsto \text{cov}(W_s, W_t), \quad s, t \in T.$$

Brownian motion and its primitives



0, 1, 2 and 3 times integrated Brownian motion

View Gaussian process W as map into Banach space $(\mathbb{B}, \|\cdot\|)$.

Theorem. *If statistical distances on the model combine appropriately with the norm $\|\cdot\|$ of \mathbb{B} , then the posterior rate is ε_n if*

$$P(\|W - w_0\| < \varepsilon_n) \geq e^{-n\varepsilon_n^2}.$$

View Gaussian process W as map into Banach space $(\mathbb{B}, \|\cdot\|)$.

Theorem. *If statistical distances on the model combine appropriately with the norm $\|\cdot\|$ of \mathbb{B} , then the posterior rate is ε_n if*

$$P(\|W - w_0\| < \varepsilon_n) \geq e^{-n\varepsilon_n^2}.$$

Proof.

- The stated condition is prior mass.
- Complexity can be shown automatic due to concentration of Gaussian processes.

View Gaussian process W as map into Banach space $(\mathbb{B}, \|\cdot\|)$.

Theorem. *If statistical distances on the model combine appropriately with the norm $\|\cdot\|$ of \mathbb{B} , then the posterior rate is ε_n if*

$$P(\|W - w_0\| < \varepsilon_n) \geq e^{-n\varepsilon_n^2}.$$

An equivalent condition is, for $(\mathbb{H}, \|\cdot\|_{\mathbb{H}})$ the **RKHS**,

$$\phi_0(\varepsilon_n) \leq n\varepsilon_n^2 \quad \text{AND} \quad \inf_{h \in \mathbb{H}: \|h - w_0\| < \varepsilon_n} \|h\|_{\mathbb{H}}^2 \leq n\varepsilon_n^2,$$

where $\phi_0(\varepsilon) = -\log \Pi(\|W\| < \varepsilon)$ is the **small ball exponent**.

- *Both inequalities give lower bound on ε_n .*
- *The first depends on W and not on w_0 .*

Settings

Density estimation

X_1, \dots, X_n iid in $[0, 1]$,

$$p_\theta(x) = \frac{e^{\theta(x)}}{\int_0^1 e^{\theta(t)} dt}.$$

Classification

$(X_1, Y_1), \dots, (X_n, Y_n)$ iid in $[0, 1] \times \{0, 1\}$

$$P_\theta(Y = 1 | X = x) = \frac{1}{1 + e^{-\theta(x)}}.$$

Regression

Y_1, \dots, Y_n independent $N(\theta(x_i), \sigma^2)$, for fixed design points x_1, \dots, x_n .

Ergodic diffusions

$(X_t: t \in [0, n])$, ergodic, recurrent:

$$dX_t = \theta(X_t) dt + \sigma(X_t) dB_t.$$

- Distance on parameter: **Hellinger on p_θ** .
- Norm on W : **uniform**.
- Distance on parameter: **$L_2(G)$ on P_θ** . (G marginal of X_i .)
- Norm on W : **$L_2(G)$** .
- Distance on parameter: **empirical L_2 -distance on θ** .
- Norm on W : **empirical L_2 -distance**.
- Distance on parameter: **random Hellinger h_n** ($\approx \|\cdot / \sigma\|_{\mu_0, 2}$).
- Norm on W : **$L_2(\mu_0)$** . (μ_0 stationary measure.)

Brownian Motion prior

Theorem. *If $\theta_0 \in C^\beta[0, 1]$, then the rate for Brownian motion is $n^{-\beta/2}$ if $\beta \leq 1/2$ and $n^{-1/4}$ for every $\beta \geq 1/2$.*

The rate is $n^{-\beta/(2\beta+1)}$ iff $\beta = 1/2$.

Brownian Motion prior

Theorem. If $\theta_0 \in C^\beta[0, 1]$, then the rate for Brownian motion is $n^{-\beta/2}$ if $\beta \leq 1/2$ and $n^{-1/4}$ for every $\beta \geq 1/2$.

The rate is $n^{-\beta/(2\beta+1)}$ iff $\beta = 1/2$.



The *small ball probability* of Brownian motion is

$$\mathbb{P}(\|W\|_\infty < \varepsilon) \sim e^{-(1/\varepsilon)^2}, \quad \varepsilon \downarrow 0.$$

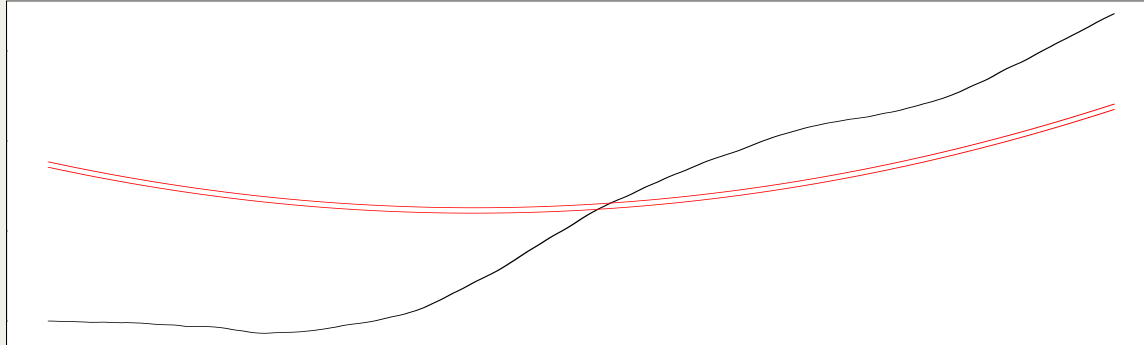
This causes a $n^{-1/4}$ -rate even for very smooth truths.

Integrated Brownian Motion prior

Theorem.

If $\theta_0 \in C^\beta[0, 1]$, then the rate for $(\alpha - 1/2)$ -times integrated Brownian is $n^{-(\alpha \wedge \beta)/(2\alpha + d)}$.

The rate is $n^{-\beta/(2\beta+1)}$ iff $\beta = \alpha$.



The small ball probability of integrated Brownian motion is much bigger

Integrated Brownian motion prior — adaptation by random scaling

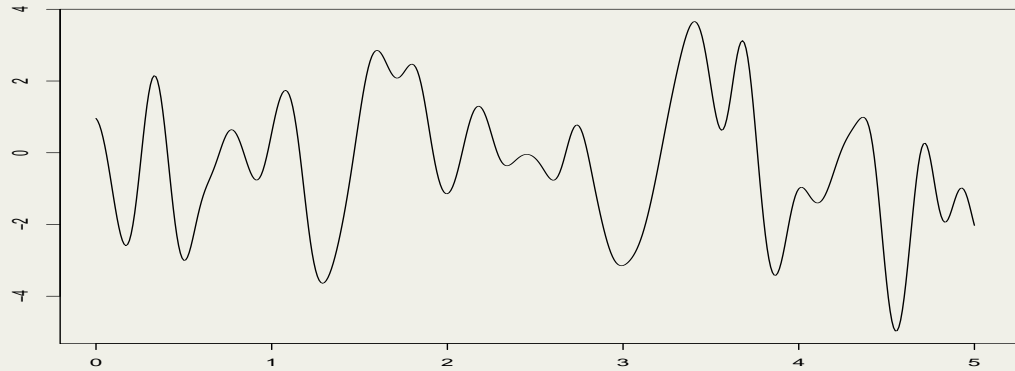
- $1/c \sim \Gamma(a, b)$.
- $(G_t: t > 0)$ k -times integrated Brownian motion “released at zero”,
- $W_t \sim \sqrt{c} G_t$.

Theorem. *The prior $W = (\sqrt{c} G_t: 0 \leq t \leq 1)$ gives contraction rate $n^{-\beta/(2\beta+1)}$ for $\theta_0 \in C^\beta[0, 1]$, for any $\beta \in (0, k + 1]$.*

Bayes solves the bandwidth problem.

Square exponential prior

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



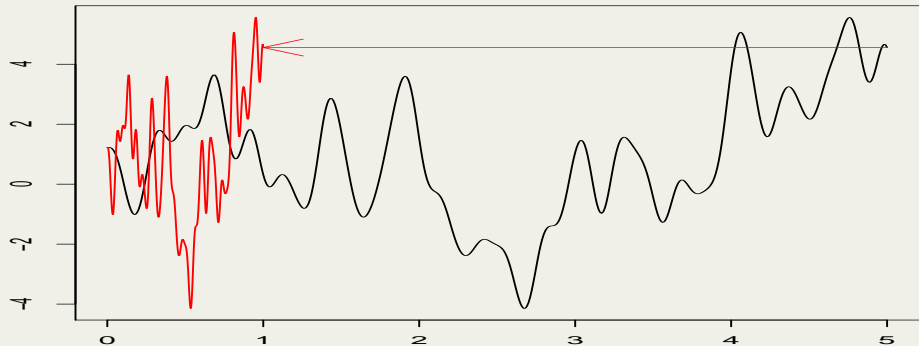
Theorem. *The prior G gives a rate $(\log n)^\gamma / \sqrt{n}$ if θ_0 is analytic, but may give a rate $(\log n)^{-\gamma'}$ if θ_0 is only ordinary smooth.*

Square exponential prior — adaptation by random time scaling

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

Theorem.

- if $\theta_0 \in C^\beta[0, 1]^d$, then the rate of contraction is nearly $n^{-\beta/(2\beta+d)}$.
- if θ_0 is supersmooth, then the rate is nearly $n^{-1/2}$.



Gaussian processes: summary



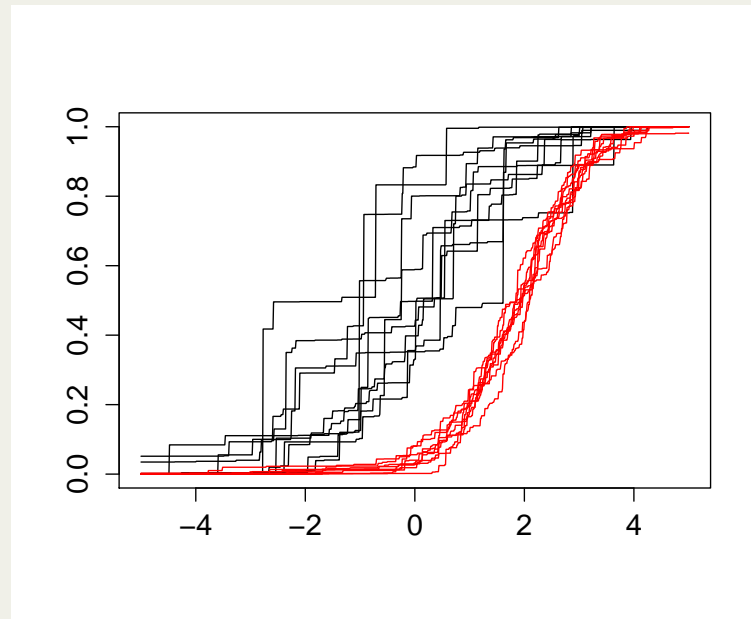
- Recovery is best if prior ‘matches’ truth.
- Mismatch slows down, but does not prevent, recovery.
- Mismatch can be prevented by using hyperparameters.

Dirichlet process mixtures

Dirichlet process [Ferguson 1973]

Definition. A **Dirichlet process** is a random measure P on $(\mathcal{X}, \mathcal{X})$ such that for every partition A_1, \dots, A_k of \mathcal{X} ,

$$(P(A_1), \dots, P(A_k)) \sim \text{Dir}(k; \alpha(A_1), \dots, \alpha(A_k)).$$

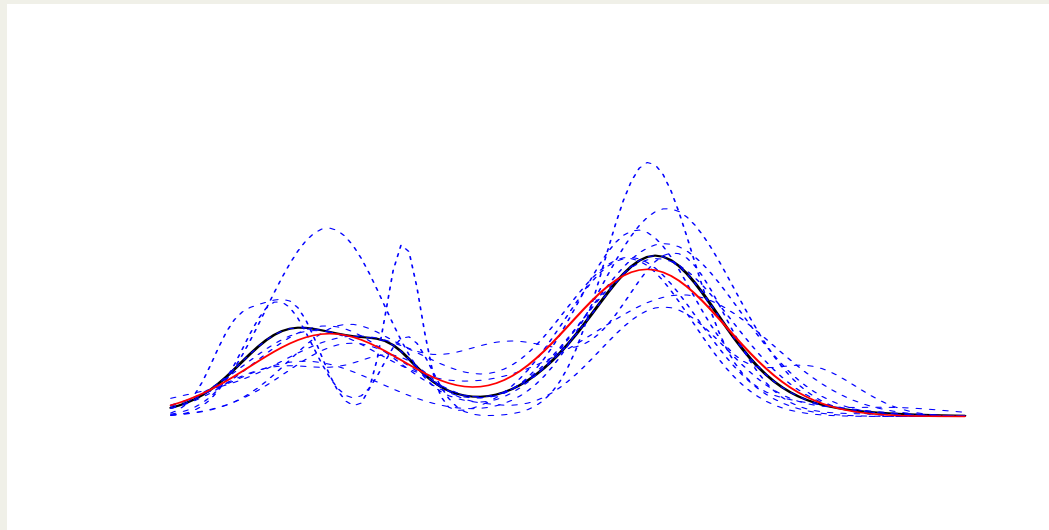


Draws from Dirichlet prior (black) and posterior based on random sample from P (red).

Dirichlet normal mixtures [Ghosal, vdV, Rousseau, Kruijer, Tokdar, Shen, 2001–2013]

- $F \sim$ Dirichlet process (α) , independent of $1/c \sim \Gamma(a, b)$.
- **Data:** $X_1, \dots, X_n | F, c \stackrel{\text{iid}}{\sim} p_{F,c}$, for

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



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

- $F \sim$ Dirichlet process (α) , independent of $1/c \sim \Gamma(a, b)$.
- **Data:** $X_1, \dots, X_n | F, c \stackrel{\text{iid}}{\sim} p_{F,c}$, for

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

Theorem. *Hellinger rate of contraction for $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_0$ is*

- *nearly $n^{-1/2}$ if $p_0 = p_{F_0, c_0}$, some F_0, c_0 .*
- *nearly $n^{-\beta/(2\beta+1)}$ if p_0 has β derivatives and exponentially small tails.*

- $F \sim$ Dirichlet process (α) , independent of $1/c \sim \Gamma(a, b)$.
- **Data:** $X_1, \dots, X_n | F, c \stackrel{\text{iid}}{\sim} p_{F,c}$, for

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

Theorem. *Hellinger rate of contraction for $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_0$ is*

- *nearly $n^{-1/2}$ if $p_0 = p_{F_0, c_0}$, some F_0, c_0 .*
- *nearly $n^{-\beta/(2\beta+1)}$ if p_0 has β derivatives and exponentially small tails.*

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

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

Credible sets



- A parameter Θ is generated according to a **prior distribution** Π .
- Given $\Theta = \theta$ the data X is generated according to a measure P_θ .

This gives a **joint distribution** of (X, Θ) .

- Given observed data x the statistician computes the conditional distribution of Θ given $X = x$, the **posterior distribution**:

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

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

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

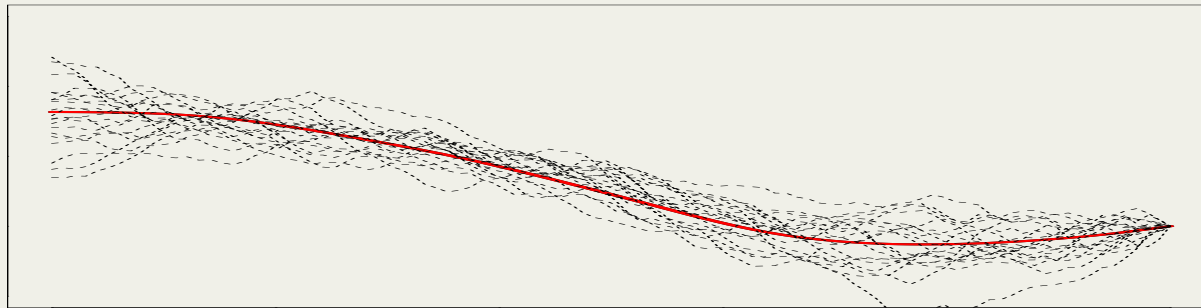
Nonparametric credible sets

Nonparametric credible sets are sets in function space.

They can take many forms:

- Plots of realizations from the posterior distribution.
- Credible bands.
- Credible balls.

They are routinely produced from MCMC output.



20 realizations from the posterior.

Do credible sets correctly quantify *remaining uncertainty*?

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

Does

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

imply

$$P_{\theta_0}(\theta_0 \in C_n(X)) = 0.95?$$

Do credible sets correctly quantify *remaining uncertainty*?

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

Does

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

imply

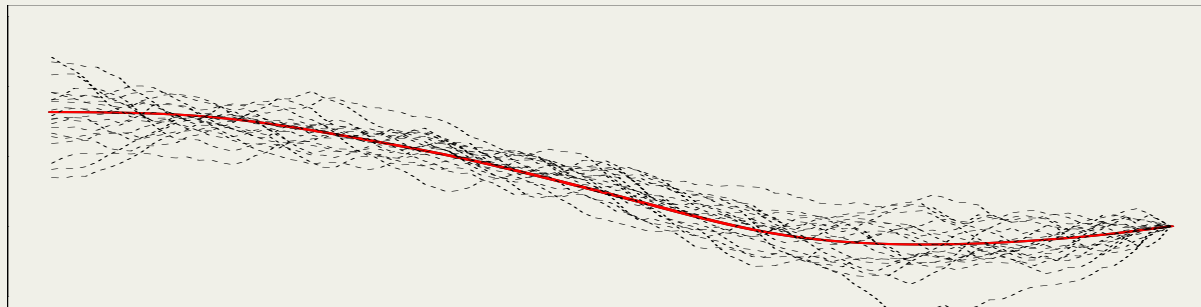
$$P_{\theta_0}(\theta_0 \in C_n(X)) = 0.95?$$

Rarely!

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

Do credible sets correctly quantify *remaining uncertainty*?

Does the **spread in the posterior** give the correct order of the discrepancy between θ_0 and the posterior mean?



20 realizations from the posterior.

Is this picture interesting?

Wahba, 1975

J. R. Statist. Soc. B (1983),
45, No. 1, pp. 133–150

Bayesian “Confidence Intervals” for the Cross-validated Smoothing Spline

By GRACE WAHBA

University of Wisconsin, USA

[Received August 1981. Revised August 1982]

SUMMARY

We consider the model $Y(t_i) = g(t_i) + \epsilon_i$, $i = 1, 2, \dots, n$, where $g(t)$, $t \in [0, 1]$ is a smooth function and the $\{\epsilon_i\}$ are independent $N(0, \sigma^2)$ errors with σ^2 unknown. The cross-validated smoothing spline can be used to estimate g non-parametrically from observations on $Y(t_i)$, $i = 1, 2, \dots, n$, and the purpose of this paper is to study confidence intervals for this estimate. Properties of smoothing splines as Bayes estimates are used to derive confidence intervals based on the posterior covariance function of the estimate. A small Monte Carlo study with the cubic smoothing spline is carried out to suggest by example to what extent the resulting 95 per cent confidence intervals can be expected to cover about 95 per cent of the true (but in practice unknown) values of $g(t_i)$, $i = 1, 2, \dots, n$. The method was also applied to one example of a two-dimensional thin plate smoothing spline. An asymptotic theoretical argument is presented to explain why the method can be expected to work on fixed smooth functions (like those tried), which are “smoother” than the sample functions from the prior distributions on which the confidence interval theory is based.

Keywords: SPLINE SMOOTHING; CROSS-VALIDATION; CONFIDENCE INTERVALS

1. INTRODUCTION

Consider the model

$$Y(t_i) = g(t_i) + \epsilon_i, \quad i = 1, 2, \dots, n, \quad t_i \in [0, 1], \quad (1.1)$$

where $\epsilon = (\epsilon_1, \dots, \epsilon_n)' \sim N(0, \sigma^2 I_{n \times n})$, σ^2 is unknown and $g(\cdot)$ is a fixed but unknown function with $m-1$ continuous derivatives and $\int_0^1 (g^{(m)}(t))^2 dt < \infty$. The smoothing spline estimate of g given $Y(t_i) = y_i$, $i = 1, 2, \dots, n$, which we will call $g_{n,\lambda}$, is the minimizer of

$$n^{-1} \sum_{i=1}^n (g(t_i) - y_i)^2 + \lambda \int_0^1 (g^{(m)}(t))^2 dt$$

Works great!

Cox, 1993

The Annals of Statistics
1993, Vol. 21, No. 2, 903–923

AN ANALYSIS OF BAYESIAN INFERENCE FOR NONPARAMETRIC REGRESSION¹

By DENNIS D. COX

Rice University

The observation model $y_i = \beta(i/n) + \epsilon_i$, $1 \leq i \leq n$, is considered, where the ϵ_i 's are i.i.d. with mean zero and variance σ^2 and β is an unknown smooth function. A Gaussian prior distribution is specified by assuming β is the solution of a high order stochastic differential equation. The estimation error $\delta = \beta - \hat{\beta}$ is analyzed, where $\hat{\beta}$ is the posterior expectation of β . Asymptotic posterior and sampling distributional approximations are given for $\|\delta\|^2$ when $\|\cdot\|$ is one of a family of norms natural to the problem. It is shown that the frequentist coverage probability of a variety of $(1 - \alpha)$ posterior probability regions tends to be larger than $1 - \alpha$, but will be infinitely often less than any $\epsilon > 0$ as $n \rightarrow \infty$ with prior probability 1. A related continuous time signal estimation problem is also studied.

1. Introduction. In this article we consider Bayesian inference for a class of nonparametric regression models. Suppose we observe

$$(1.1) \quad Y_{ni} = \beta(t_{ni}) + \epsilon_i, \quad 1 \leq i \leq n,$$

where $t_{ni} = i/n$, $\beta: [0, 1] \rightarrow \mathbb{R}$ is an unknown smooth function, and $\epsilon_1, \epsilon_2, \dots$ are i.i.d. random errors with mean 0 and known variance $\sigma^2 < \infty$. The ϵ_i are modeled as $N(0, \sigma^2)$. A Gaussian prior for β will now be specified. Let $m \geq 2$ and for some constants a_0, \dots, a_m with $a_m \neq 0$ let

$$L = \sum_{i=0}^m a_i D^i$$

Fails miserably!

Priors of fixed regularity

Coverage requires undersmoothing

In *nonparametric statistics*:

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

Coverage requires undersmoothing

In *nonparametric statistics*:

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

In *nonparametric Bayesian statistics*:

this occurs if the **prior produces too smooth functions**.

Coverage requires undersmoothing

In *nonparametric statistics*:

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

In *nonparametric Bayesian statistics*:

this occurs if the **prior produces too smooth functions**.

EXAMPLE

Truth:
$$\theta_0(x) = \sum_{i=1}^{\infty} \theta_{0,i} e_i(x), \quad \theta_{0,i} \asymp i^{-1-2\beta}.$$

Prior:
$$x \mapsto \sum_{i=1}^{\infty} \theta_i e_i(x), \quad \theta_i \stackrel{\text{ind}}{\sim} N(0, i^{-1-2\alpha}).$$

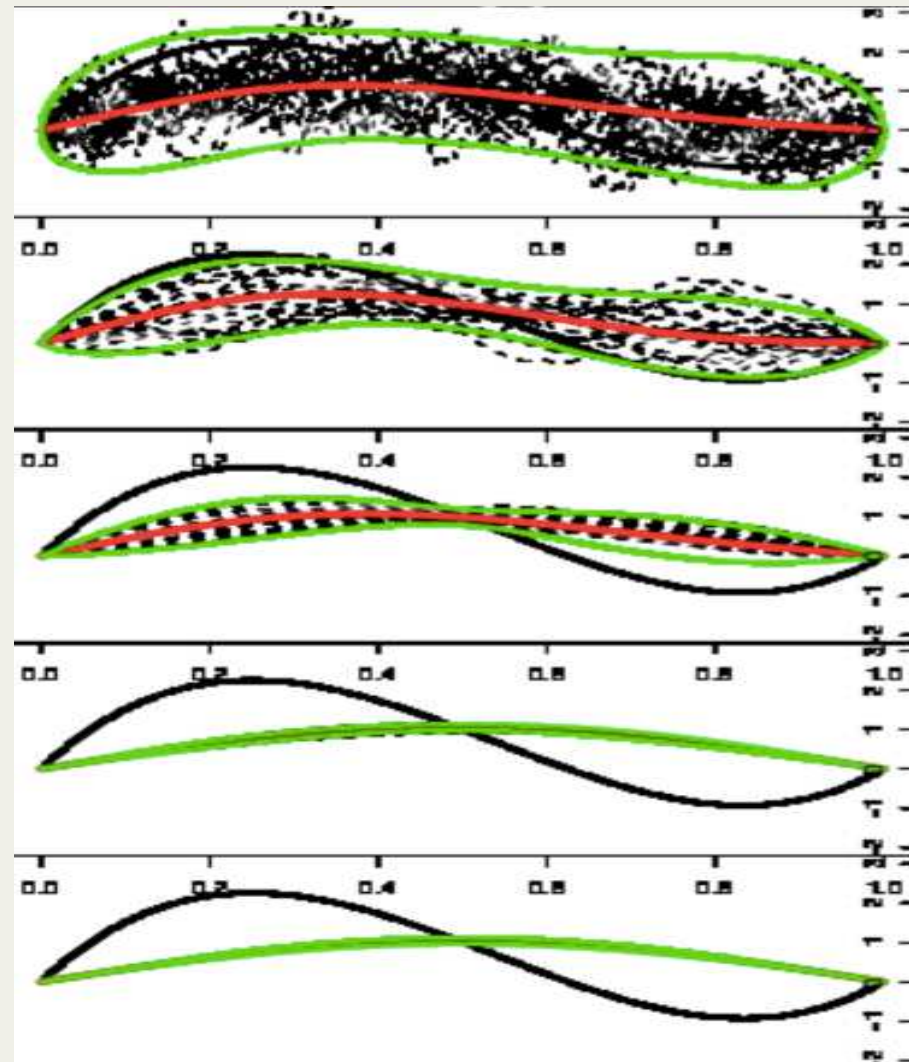
Interpretation:

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

$\alpha > \beta$: prior oversmooths.

$\alpha < \beta$: prior undersmooths.

Example: heat equation ($n=10\ 000$)



True θ_0 (black), posterior mean (red), 20 realizations from the posterior (dashed black), and posterior credible bands (green).
Left: $n = 10^4$; right: $n = 10^8$. Top to bottom: prior of increasing smoothness.

Priors of flexible regularity

Bayesian adaptation

Family of priors Π_c of varying smoothness; posteriors $\Pi_{n,c}(\cdot | Y_n)$.

Empirical Bayes:

- \hat{c}_n some “estimator”.
- Plug-in posterior $\Pi_{n,\hat{c}_n}(\cdot | Y_n)$.

Bayesian adaptation

Family of priors Π_c of varying smoothness; posteriors $\Pi_{n,c}(\cdot | Y_n)$.

Empirical Bayes:

- \hat{c}_n some “estimator”.
- Plug-in posterior $\Pi_{n,\hat{c}_n}(\cdot | Y_n)$.

Hierarchical Bayes:

- Full Bayes, with prior π on c .
- Posterior $\int \Pi_{n,c}(\cdot | Y_n) \pi_n(c | Y_n) dc$.

Both methods (in particular Hierarchical Bayes) are known to give **adaptive reconstructions** in some generality:

if the true function is smoother, then the reconstruction is better.

Bayesian adaptation

Family of priors Π_c of varying smoothness; posteriors $\Pi_{n,c}(\cdot | Y_n)$.

Empirical Bayes:

- \hat{c}_n some “estimator”.
- Plug-in posterior $\Pi_{n,\hat{c}_n}(\cdot | Y_n)$.

Hierarchical Bayes:

- Full Bayes, with prior π on c .
- Posterior $\int \Pi_{n,c}(\cdot | Y_n) \pi_n(c | Y_n) dc$.

Both methods (in particular Hierarchical Bayes) are known to give **adaptive reconstructions** in some generality:

if the true function is smoother, then the reconstruction is better.

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

Definition. $C_n(X^{(n)})$ is an *honest confidence set* over a model Θ if

$$P_{\theta_0}(C_n(X^{(n)}) \ni \theta_0) \geq 0.95, \quad \text{for all } \theta_0 \in \Theta.$$

Definition. $C_n(X^{(n)})$ is an *honest confidence set* over a model Θ if

$$P_{\theta_0}(C_n(X^{(n)}) \ni \theta_0) \geq 0.95, \quad \text{for all } \theta_0 \in \Theta.$$

Theorem. For any $\Theta_1 \subset \Theta$ the diameter of honest $C_n(X^{(n)})$ cannot be smaller, uniformly over Θ_1 , than:

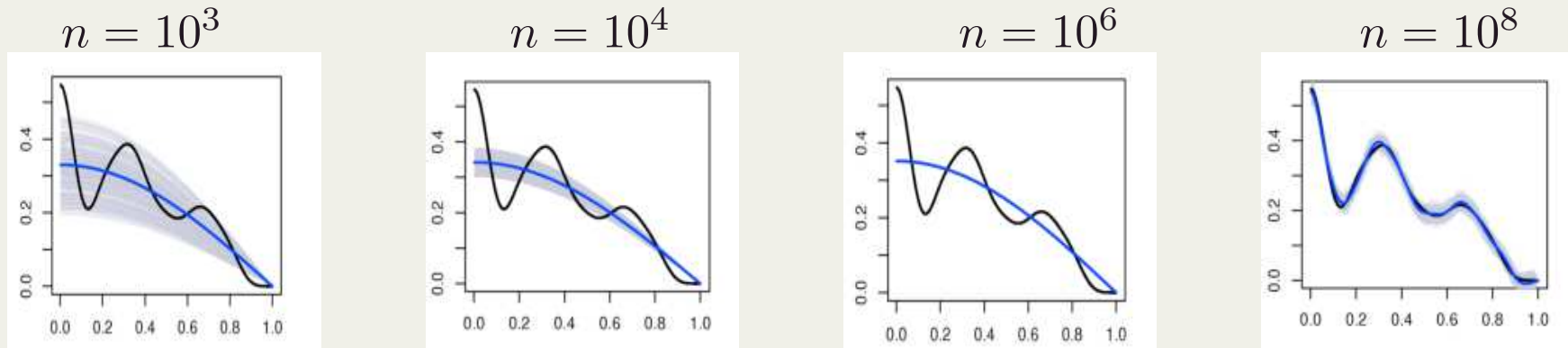
(a) ε_n such that, for any T_n ,

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

(b) rate ε_n of minimax testing of $H_0: \theta \in \Theta'_1$ versus $H_1: \theta \in \Theta, d(\theta, \Theta'_1) > \varepsilon_n$, for any given $\Theta'_1 \subset \Theta_1$.

- (a) typically gives minimax rate of estimation for model Θ_1 .
(b) is determined by biggest model Θ rather than Θ_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.

[Szabo, vdV, van Zanten, 2016.]

[*Not “asymptotical”*: for still bigger n it can become good and bad again!]

Credible balls — counter example — reconstructing a derivative

Theorem. For $n_1 \geq 2$ and $n_j \geq n_{j-1}^4$ for every j , and $\beta > 0$, 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 95%-credible ball \hat{C}_n centered at posterior mean and radius blown up by $L_n \ll n^\delta$ satisfies

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

Credible balls — counter example — reconstructing a derivative

Theorem. For $n_1 \geq 2$ and $n_j \geq n_{j-1}^4$ for every j , and $\beta > 0$, 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 95%-credible ball \hat{C}_n centered at posterior mean and radius blown up by $L_n \ll n^\delta$ satisfies

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

- Data allows inference on $\theta_1, \dots, \theta_N$ for an *effective dimension* $N = N_n$.
- Trouble if $\theta_1, \dots, \theta_N$ does not resemble $\theta_1, \theta_2, \dots$
- Example θ 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:

- The size of an honest confidence set is determined by the smallest possible regularity level.
- (Bayesian constructions can be misleading.)

Estimation versus uncertainty quantification

Adaptive estimation:

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

Uncertainty quantification:

- The size of an honest confidence set is determined by the smallest possible regularity level.
- (Bayesian constructions can be misleading.)

SOLUTION 1: *be honest*; only make conditional confidence statements.

Estimation versus uncertainty quantification

Adaptive estimation:

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

Uncertainty quantification:

- The size of an honest confidence set is determined by the smallest possible regularity level.
- (Bayesian constructions can be misleading.)

SOLUTION 1: *be honest*; only make conditional confidence statements.

SOLUTION 2: determine which θ cause the trouble; argue that these are implausible.

Nonparametric regression

Nonparametric regression

- $\theta: \mathcal{X} \rightarrow \mathbb{R}$; *design points* $x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .

Nonparametric regression

- $\theta: \mathcal{X} \rightarrow \mathbb{R}$; *design points* $x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .
- **Posterior:** $\vec{\theta}_n | Y_n \sim N_n(\hat{\theta}_{n,c}, I - \Sigma_{n,c}^{-1})$.

$$\hat{\theta}_{n,c} = (I - \Sigma_{n,c}^{-1})Y_n,$$

$$\Sigma_{n,c} = I + c \text{Cov}(\vec{W}_n).$$

Examples of processes W :

- Brownian motion
- discrete Laplacian $(n^2 L)^{-\alpha} \vec{W}_n \sim N_n(0, I)$, for $Lf(i) = \sum_{j:j \sim i} [f(j) - f(i)]$. [Kirichenko & van Zanten, 2015.]
- Brownian sheet
- eigenfunctions as Brownian sheet but “Sobolev eigenvalues”.

Empirical Bayes and hierarchical Bayes

- $\theta: \mathcal{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .
- **Posterior:** $\vec{\theta}_n | Y_n \sim N_n(\hat{\theta}_{n,c}, I - \Sigma_{n,c}^{-1})$.

RISK-BASED EMPIRICAL BAYES [Wahba, 1975]: plug in:

$$\hat{c}_n = \underset{c}{\operatorname{argmin}} \left[\underbrace{\operatorname{tr}((I - \Sigma_{n,c}^{-1})^2) - \operatorname{tr}(\Sigma_{n,c}^{-2}) + \vec{Y}_n^T \Sigma_{n,c}^{-2} \vec{Y}_n}_{\text{unbiased estimate of } \mathbb{E}_\theta \|\hat{\theta}_{n,c} - \vec{\theta}_n\|^2} \right].$$

Empirical Bayes and hierarchical Bayes

- $\theta: \mathcal{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .
- **Posterior:** $\vec{\theta}_n | Y_n \sim N_n(\hat{\theta}_{n,c}, I - \Sigma_{n,c}^{-1})$.

RISK-BASED EMPIRICAL BAYES [Wahba, 1975]: plug in:

$$\hat{c}_n = \operatorname{argmin}_c \underbrace{\left[\operatorname{tr}((I - \Sigma_{n,c}^{-1})^2) - \operatorname{tr}(\Sigma_{n,c}^{-2}) + \vec{Y}_n^T \Sigma_{n,c}^{-2} \vec{Y}_n \right]}_{\text{unbiased estimate of } \mathbb{E}_\theta \|\hat{\theta}_{n,c} - \vec{\theta}_n\|^2}.$$

- **Marginal distribution:** $Y_n | c \sim N_n(0, \Sigma_{n,c})$, $\Sigma_{n,c} = I + c \operatorname{Cov}(\vec{W}_n)$.

LIKELIHOOD-BASED EMPIRICAL BAYES: plug in MLE:

$$\hat{c}_n = \operatorname{argmin}_c \left[\log \det \Sigma_{n,c} + \vec{Y}_n^T \Sigma_{n,c}^{-1} \vec{Y}_n \right].$$

Empirical Bayes and hierarchical Bayes

- $\theta: \mathcal{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .
- **Posterior:** $\vec{\theta}_n | Y_n \sim N_n(\hat{\theta}_{n,c}, I - \Sigma_{n,c}^{-1})$.

RISK-BASED EMPIRICAL BAYES [Wahba, 1975]: plug in:

$$\hat{c}_n = \operatorname{argmin}_c \underbrace{\left[\operatorname{tr}((I - \Sigma_{n,c}^{-1})^2) - \operatorname{tr}(\Sigma_{n,c}^{-2}) + \vec{Y}_n^T \Sigma_{n,c}^{-2} \vec{Y}_n \right]}_{\text{unbiased estimate of } \mathbb{E}_\theta \|\hat{\theta}_{n,c} - \vec{\theta}_n\|^2}.$$

- **Marginal distribution:** $Y_n | c \sim N_n(0, \Sigma_{n,c})$, $\Sigma_{n,c} = I + c \operatorname{Cov}(\vec{W}_n)$.

LIKELIHOOD-BASED EMPIRICAL BAYES: plug in MLE:

$$\hat{c}_n = \operatorname{argmin}_c \left[\log \det \Sigma_{n,c} + \vec{Y}_n^T \Sigma_{n,c}^{-1} \vec{Y}_n \right].$$

HIERARCHICAL BAYES:

- **Prior:** $c^{-1} \sim \Gamma(a, b)$.

Polished tail sequences

Definition. $\theta \in \ell^2$ 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 .

Polished tail sequences

Definition. $\theta \in \ell^2$ 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...:

Polished tail sequences

Definition. $\theta \in \ell^2$ 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.

Polished tail sequences

Definition. $\theta \in \ell^2$ 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.

Polished tail sequences

Definition. $\theta \in \ell^2$ 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.

Polished tail arrays

- $\theta_{1,n}, \dots, \theta_{n,n}$ coefficients of $\vec{\theta}_n = (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$ in eigenbasis of $\text{Cov}(\vec{W}_n)$.
- $\lambda_{1,n}, \dots, \lambda_{n,n}$ eigenvalues of $\text{Cov}(\vec{W}_n)$.

Definition. θ satisfies the **discrete polished tail condition** if

$$\sum_{i:0.001 \leq c\lambda_{i,n} \leq 1} \theta_{i,n}^2 \geq 0.001 \sum_{i:c\lambda_{i,n} \leq 1} \theta_{i,n}^2, \quad \forall c > 0.$$

In the special case that $\lambda_{i,n} \asymp K_n/i^k$ this is close to polished tail.

Credible balls are honest over polished tail functions

- $\theta: \mathcal{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .

$$\hat{\theta}_{n,c}(x) = \mathbb{E}[\theta(x) | Y_n, c], \quad s_n^2(c) = \mathbb{E}[\|\vec{\theta}_n - \hat{\theta}_{n,c}\|^2 | Y_n, c].$$

CREDIBLE BALL:

$$\hat{C}_{n,M} = \{\theta: \|\vec{\theta}_n - \hat{\theta}_{n,\hat{c}_n}\| < M s_n(\hat{c}_n)\},$$

and similar for hierarchical Bayes.

Theorem. *For not too small M , uniformly in discrete polished tail functions θ ,*

$$P_\theta(\theta \in \hat{C}_{n,M}) \rightarrow 1.$$

Credible intervals are honest over polished tail functions

- $\theta: \mathfrak{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathfrak{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .

$$\hat{\theta}_{n,c}(x) = \mathbb{E}[\theta(x) | Y_n, c], \quad s_n^2(c, x) = \mathbb{E}[|\theta(x) - \hat{\theta}_{n,c}(x)|^2 | \vec{Y}_n, c].$$

CREDIBLE INTERVALS:

$$\hat{C}_{n,M}(x) = \{\theta: |\theta(x) - \hat{\theta}_{n,\hat{c}_n}(x)| < M s_n(\hat{c}_n, x)\},$$

and similar for hierarchical Bayes.

Theorem. *If $x_{j,n}$ “uniformly spread relative to the prior”, then for not too small M and all $\gamma < 1$, uniformly in discrete polished tail functions θ*

$$P_\theta \left(\frac{1}{n} \sum_{i=1}^n 1\{\theta \in \hat{C}_{n,M}(x_{i,n})\} \geq \gamma \right) \rightarrow 1.$$

Credible bands are honest for Bayesians

- $\theta: \mathcal{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathcal{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .

$$\hat{\theta}_{n,c}(x) = \mathbb{E}[\theta(x) | Y_n, c], \quad \tilde{s}_n^2(c) = \mathbb{E}[\sup_x |\theta(x) - \hat{\theta}_{n,c}(x)|^2 | \vec{Y}_n, c].$$

CREDIBLE BAND:

$$\tilde{C}_{n,M} = \bigcap_x \{\theta: |\theta(x) - \hat{\theta}_{n,\hat{c}_n}(x)| < M \tilde{s}_n(\hat{c}_n)\},$$

and similar for hierarchical Bayes.

Theorem. For almost any realization θ from a Gaussian process prior and not too small M ,

$$P_\theta(\theta \in \tilde{C}_{n,M}) \rightarrow 1.$$

Self-similarity [after Giné+Nickl, Hoffmann+Nickl, Bull, 2010-12]

Definition. A parameter $\theta \in \ell_2$ is **self-similar** of order β if

$$\sup_i i^{2\beta+1} \theta_i^2 \leq M,$$

$$\sum_{i=N}^{1000N} \theta_i^2 \geq 0.001 M N^{-2\beta}, \quad \forall N.$$

Interpretation:

θ has some energy at *any frequency* N relative to the total energy above N .

Credible bands can be honest for non-Bayesians

- $\theta: \mathfrak{X} \rightarrow \mathbb{R}; x_{1,n}, \dots, x_{n,n} \in \mathfrak{X}$.
- **Data:** $Y_n | \theta \sim N_n(\vec{\theta}_n, I)$, for $\vec{\theta}_n := (\theta(x_{1,n}), \dots, \theta(x_{n,n}))^T$.
- **Prior:** $\theta \sim \sqrt{c} W$, for Gaussian process W .

$$\hat{\theta}_{n,c}(x) = \mathbb{E}[\theta(x) | Y_n, c], \quad \tilde{s}_n^2(c, x) = \mathbb{E}[\sup_x |\theta(x) - \hat{\theta}_{n,c}(x)|^2 | \vec{Y}_n, c].$$

CREDIBLE BAND:

$$\tilde{C}_{n,M} = \bigcap_x \{\theta: |\theta(x) - \hat{\theta}_{n,\hat{c}_n}(x)| < M \tilde{s}_n(\hat{c}_n)\},$$

and similar for hierarchical Bayes.

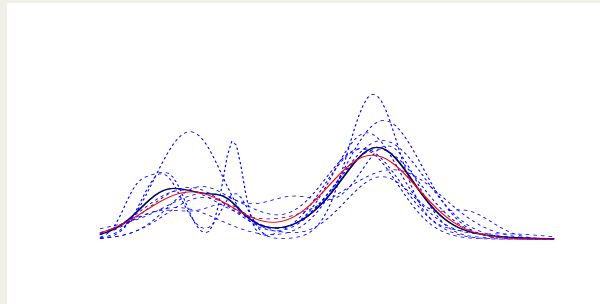
Theorem. *If θ is self-similar and Hölder of the same order, then for not too small M ,*

$$P_\theta(\theta \in \tilde{C}_{n,M}) \rightarrow 1.$$

Work in progress

Story appears to be generic, but conditions for good behaviour depend on prior and model.

There is further work [e.g. by Szabó et al.], but much is unknown.



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

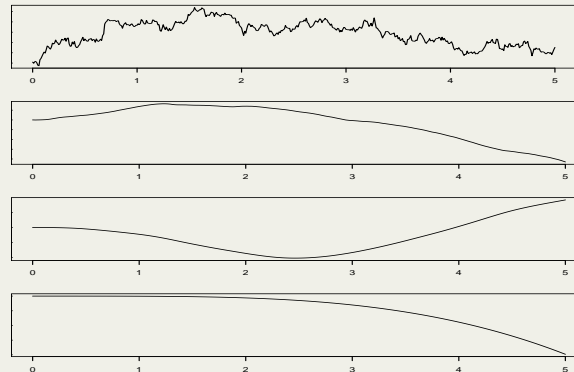
Summary

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.



Bayesians are perhaps more easily misled as they trust their priors. In nonparametrics they should not, as **the fine details of a prior are not obvious.**



Co-authors



Subhashis Ghoshal



Harry van Zanten



Ismael Castillo



Johannes Schmidt-Hieber



Stéphanie van der Pas



Botond Szabo



Willem Kruijer



Judith Rousseau



Bartek Knapik



Bas Kleijn



Suzanne Sniekers



Fengnan Gao

Example: heat equation

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

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

We **observe** a noisy version of the final curve: for Z white noise:

$$Y_n = K\theta + n^{-1/2}Z.$$

very ill-posed inverse problem: $Y_{n,i} | \theta_i \sim N(\kappa_i \theta_i, n^{-1})$ for

$$\kappa_i = e^{-i^2 \pi^2} \quad e_i = \sqrt{2} \sin(i\pi x),$$

$$(i = 1, 2, \dots).$$

Credible balls — counter example — reconstructing a derivative

The **Volterra operator** $K: L_2[0, 1] \rightarrow L_2[0, 1]$ is given by

$$K\theta(x) = \int_0^x \theta(s) ds.$$

We **observe** $(Y_n(x): x \in [0, 1])$, for W Brownian motion,

$$dY_n(x) = K\theta(x) dx + \frac{1}{\sqrt{n}} dW(x), \quad x \in [0, 1].$$

mildly ill-posed inverse problem: $Y_{n,i} | \theta_i \sim N(\kappa_i \theta_i, n^{-1})$ for

$$\kappa_i = \frac{1}{(i - 1/2)\pi} \quad e_i(x) = \sqrt{2} \cos((i - 1/2)\pi x),$$

$$(i = 0, 1, 2, \dots).$$