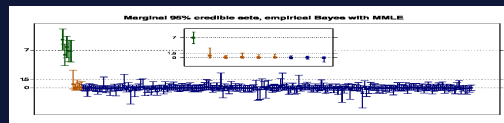


# *Bayesian Statistics in High Dimensions*

## *Lecture 2: Sparsity*

**Aad van der Vaart**

Universiteit Leiden, Netherlands



47th John H. Barrett Memorial Lectures, Knoxville, Tennessee, May 2017

# Contents

**Sparsity**

**Bayesian Sparsity**

**Frequentist Bayes**

**Model Selection Prior**

**Horseshoe Prior**

Sparsity

## Sparsity — sequence model

A **sparse model** has many parameters, but most of them are (nearly) zero.

# Sparsity — sequence model

A **sparse model** has many parameters, but most of them are (nearly) zero.

In this lecture, (Bayesian) theory for:

$$Y^n \sim N_n(\theta, I), \quad \text{for } \theta = (\theta_1, \dots, \theta_n) \in \mathbb{R}^n.$$

- $n$  independent observations  $Y_1^n, \dots, Y_n^n$ .
- $n$  unknowns  $\theta_1, \dots, \theta_n$ .
- $Y_i^n = \theta_i + \varepsilon_i$ , for standard normal noise  $\varepsilon_i$ .
- $n$  is large.
- many of  $\theta_1, \dots, \theta_n$  are (almost) zero.

# History

Given  $Y^n | \theta \sim N_n(\theta, I)$  the obvious estimator of  $\theta$  is  $Y^n$ .

*ML, UMVU, best-equivariant, minimax*

# History

Given  $Y^n | \theta \sim N_n(\theta, I)$  the obvious estimator of  $\theta$  is  $Y^n$ .  
*ML, UMVU, best-equivariant, minimax, but inadmissible.*

**Theorem** [Stein, 1956]

If  $n \geq 3$ , then there exists  $T$  such that  $\forall \theta \in \mathbb{R}^n$ ,

$$\mathbb{E}_\theta \|T(Y^n) - \theta\|^2 < \mathbb{E}_\theta \|Y^n - \theta\|^2.$$



# History

Given  $Y^n | \theta \sim N_n(\theta, I)$  the obvious estimator of  $\theta$  is  $Y^n$ .  
*ML, UMVU, best-equivariant, minimax, but inadmissible.*

**Theorem** [Stein, 1956]

If  $n \geq 3$ , then there exists  $T$  such that  $\forall \theta \in \mathbb{R}^n$ ,

$$\mathbb{E}_\theta \|T(Y^n) - \theta\|^2 < \mathbb{E}_\theta \|Y^n - \theta\|^2.$$



**Empirical Bayes method** [Robbins, 1960s]:

- Working hypothesis:  $\theta_1, \dots, \theta_n \stackrel{\text{iid}}{\sim} G$ .
- Estimate  $G$  using  $Y^n$ , pretending this is true.
- $T(Y^n) := \mathbb{E}_{\hat{G}(Y^n)}(\theta | Y^n)$ .





# History

Given  $Y^n | \theta \sim N_n(\theta, I)$  the obvious estimator of  $\theta$  is  $Y^n$ .  
*ML, UMVU, best-equivariant, minimax, but inadmissible.*

**Theorem** [Stein, 1956]

If  $n \geq 3$ , then there exists  $T$  such that  $\forall \theta \in \mathbb{R}^n$ ,

$$\mathbb{E}_\theta \|T(Y^n) - \theta\|^2 < \mathbb{E}_\theta \|Y^n - \theta\|^2.$$



**Empirical Bayes method** [Robbins, 1960s]:

- Working hypothesis:  $\theta_1, \dots, \theta_n \stackrel{\text{iid}}{\sim} G$ .
- Estimate  $G$  using  $Y^n$ , pretending this is true.
- $T(Y^n) := \mathbb{E}_{\hat{G}(Y^n)}(\theta | Y^n)$ .



For large  $n$  the gain can be substantial,  
("borrowing of strength"),  
and be targeted to special subsets of  $\mathbb{R}^n$ , e.g. sparse vectors.

# Sparsity — regression

$$Y^n | \theta \sim N_n(X_{n \times p} \theta, \sigma^2 I), \text{ for } \theta = (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$$

- $Y_i^n$ : measurement on individual  $i = 1, \dots, n$ .
- $X_{ij}$  score of individual  $i$  on feature  $j = 1, \dots, p$ .
- $\theta_j$  effect of feature  $j$ .
- sparse if only few features matter.

If  $p > n$ , then *sparsity is necessary* to recover  $\theta$ ,  
and  $X$  must be *sparse-invertible*, e.g.:

*Compatibility:*

$$\inf_{\theta: |S_\theta| \leq 5s_n} \frac{\|X\theta\|_2 \sqrt{|S_\theta|}}{\|X\| \|\theta\|_1} \gg 0.$$

*Mutual coherence:*

$$s_n \max_{i \neq j} |\text{cor}(X_{.i}, X_{.j})| \ll 1.$$

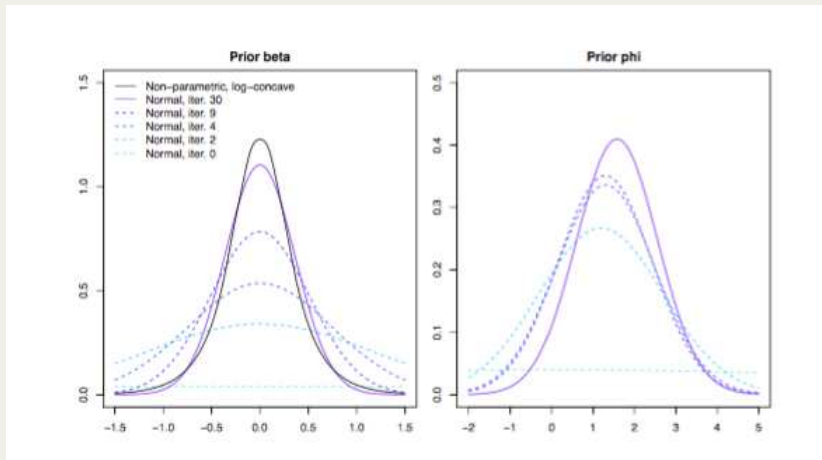
$s_n$  = true number of nonzero coordinates.

# Sparsity — RNA sequencing

- $Y_{i,j}$ : RNA expression count of tag  $j = 1, \dots, p$  in tissue  $i = 1, \dots, n$ ,
- $x_i$ : covariate(s) of tissue  $i$ , e.g. 0 or 1 for normal or cancer.
- sparse if only few tags (genes) matter.

$Y_{i,j} \sim$  (zero-inflated) *negative binomial*, with

$$\mathbb{E}Y_{i,j} = e^{\alpha_j + \beta_j x_i}, \quad \text{var } Y_{i,j} = \mathbb{E}Y_{i,j} (1 + \mathbb{E}Y_{i,j} e^{-\phi_j}).$$



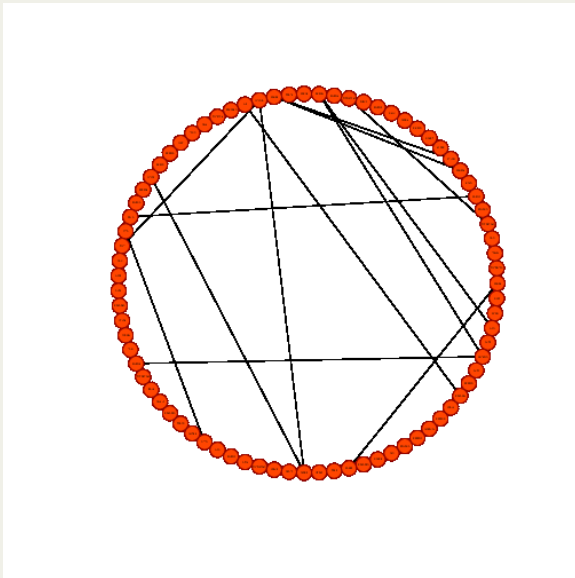
distribution of  $\beta_j$ 's and  $\phi_j$ 's estimated by Empirical Bayes

# Sparsity — Gaussian graphical model

Data  $n$  iid copies of  $Y^p | \theta \sim N_p(0, \theta^{-1})$

- $Y_j^p$  = value of individual on feature  $j$ .
- **precision matrix**  $\theta$  gives **partial correlations**:

$$\text{cor}(Y_i^p, Y_j^p | Y_k^p : k \neq i, j) = -\frac{\theta_{i,j}}{\sqrt{\theta_{i,i}\theta_{j,j}}}.$$



Apoptosis network

- nodes  $1, 2, \dots, p$
- edge  $(i, j)$  present iff  $\theta_{i,j} \neq 0$
- sparse if few edges

# Bayesian Sparsity

## Bayesian sparsity

A **sparse model** has many parameters, but most of them are (nearly) zero.

# Bayesian sparsity

A **sparse model** has many parameters, but most of them are (nearly) zero.



We express this in the prior, and  
apply **the standard (full or empirical) Bayesian machine**.

# Bayesian sparsity

A **sparse model** has many parameters, but most of them are (nearly) zero.



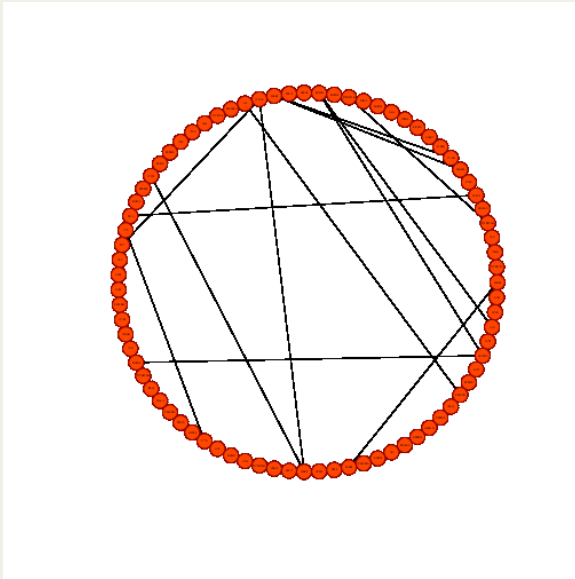
We express this in the prior, and  
apply **the standard (full or empirical) Bayesian machine**.

Prior  $\theta \sim \Pi$ , and data  $Y^n | \theta \sim p(\cdot | \theta)$ , give posterior:

$$d\Pi(\theta | Y^n) \propto p(Y^n | \theta) d\Pi(\theta).$$



# Bayesian sparsity — Gaussian graphical model



Apoptosis network

- nodes  $1, 2, \dots, p$
- edge  $(i, j)$  present iff  $\theta_{i,j} \neq 0$
- sparse if few edges

For given *incidence matrix*  $(P_{i,j})$  use different priors for

$$(\theta_{i,j}: P_{i,j} = 0) \quad \text{and} \quad (\theta_{i,j}: P_{i,j} = 1).$$

# Model selection prior

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

- (1) Choose  $s$  from prior on  $\{0, 1, 2, \dots, p\}$ .
- (2) Choose  $S \subset \{0, 1, \dots, p\}$  of size  $|S| = s$  at random.
- (3) Choose  $\theta_S = (\theta_i: i \in S) \sim g_S$  and set  $\theta_{S^c} = 0$ .

[Mitchell & Beachamp (88), George, George & McCulloch, Yuan, Berger, Johnstone & Silverman, Richardson et al., Johnson & Rossell,

Chao Gao, ...]

# Model selection prior

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

- (1) Choose  $s$  from prior on  $\{0, 1, 2, \dots, p\}$ .
- (2) Choose  $S \subset \{0, 1, \dots, p\}$  of size  $|S| = s$  at random.
- (3) Choose  $\theta_S = (\theta_i: i \in S) \sim g_S$  and set  $\theta_{S^c} = 0$ .

Example: **spike and slab**

- Choose  $\theta_1, \dots, \theta_p$  i.i.d. from  $\tau\delta_0 + (1 - \tau)G$ .
- Put a prior on  $\tau$ , e.g.  $\text{Beta}(1, p + 1)$ .

Then  $s \sim \text{binomial}$  and  $g_S = \otimes_{i \in S} g$ .

# Horseshoe prior

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

- (1) Generate  $\tau \sim \text{Cauchy}^+(0, \sigma)$  (?)
- (2) Generate  $\sqrt{\psi_1}, \dots, \sqrt{\psi_p}$  iid from  $\text{Cauchy}^+(0, \tau)$ .
- (3) Generate independent  $\theta_i \sim N(0, \psi_i)$ .

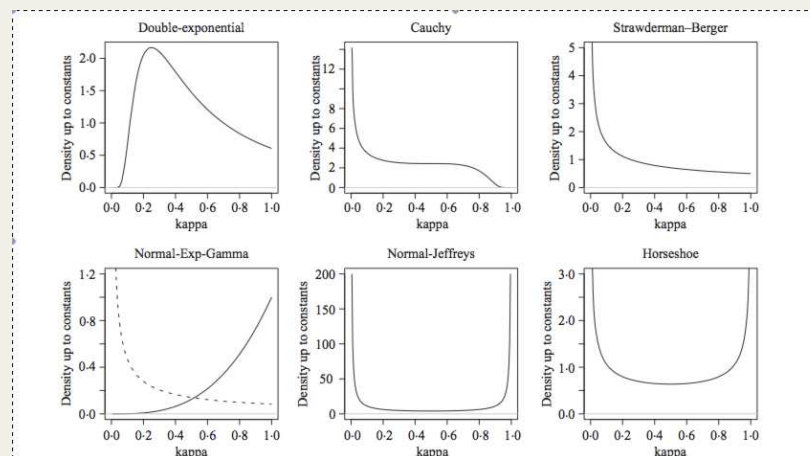
# Horseshoe prior

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

- (1) Generate  $\tau \sim \text{Cauchy}^+(0, \sigma)$  (?)
- (2) Generate  $\sqrt{\psi_1}, \dots, \sqrt{\psi_p}$  iid from  $\text{Cauchy}^+(0, \tau)$ .
- (3) Generate independent  $\theta_i \sim N(0, \psi_i)$ .

## Motivation

if  $\theta \sim N(0, \psi)$  and  $Y | \theta \sim N(\theta, 1)$ ,  
then  $\theta | Y, \psi \sim N((1 - \kappa)Y, 1 - \kappa)$  for  $\kappa = 1/(1 + \psi)$ .  
*This suggests a prior for  $\kappa$  that concentrates near 0 or 1.*



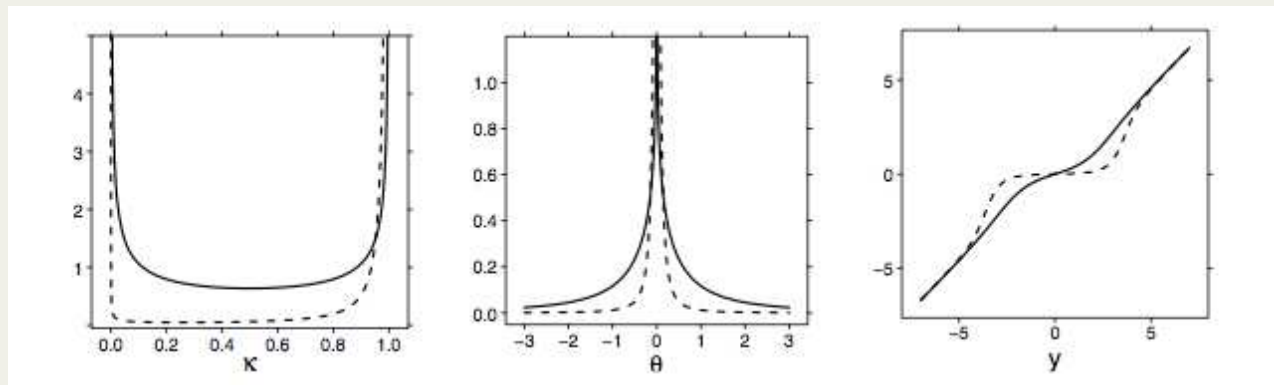
# Horseshoe prior

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

- (1) Generate  $\tau \sim \text{Cauchy}^+(0, \sigma)$  (?)
- (2) Generate  $\sqrt{\psi_1}, \dots, \sqrt{\psi_p}$  iid from  $\text{Cauchy}^+(0, \tau)$ .
- (3) Generate independent  $\theta_i \sim N(0, \psi_i)$ .

## Motivation

if  $\theta \sim N(0, \psi)$  and  $Y | \theta \sim N(\theta, 1)$ ,  
then  $\theta | Y, \psi \sim N((1 - \kappa)Y, 1 - \kappa)$  for  $\kappa = 1/(1 + \psi)$ .  
*This suggests a prior for  $\kappa$  that concentrates near 0 or 1.*

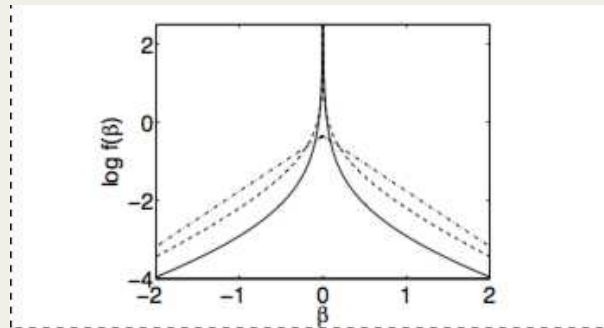


prior shrinkage factor

prior of  $\theta_i$

posterior mean of  $\theta_i$  as function of  $Y_i$

# Other sparsity priors



- *Bayesian LASSO*:  $\theta_1, \dots, \theta_p$  iid from a mixture of Laplace ( $\lambda$ ) distributions over  $\lambda \sim \sqrt{\Gamma(a, b)}$ .
- *Bayesian bridge*: Same but with Laplace replaced with a density  $\propto e^{-|\lambda y|^\alpha}$ .
- *Normal-Gamma*:  $\theta_1, \dots, \theta_p$  iid from a Gamma scale mixture of Gaussians. *Correlated multivariate normal-Gamma*:  $\theta = C\phi$  for a  $p \times k$ -matrix  $C$  and  $\phi$  with independent normal-Gamma ( $a_i, 1/2$ ) coordinates.
- *Horseshoe*.
- *Horseshoe+*.
- *Normal spike*.
- *Scalar multiple of Dirichlet*.
- *Nonparametric Dirichlet*.
- ...

[Park & Casella 08, Polson & Scott, Griffin & Brown 10, 12, Carvalho & Polson & Scott, 10, George & Rockova 13, Bhattacharya et al.

# LASSO is not Bayesian

$$\hat{\theta}_{\text{LASSO}} = \underset{\theta}{\operatorname{argmin}} \left[ \|Y^n - X\theta\|^2 + \lambda_n \sum_{i=1}^p |\theta_i| \right].$$

**posterior mode** for prior  $\theta_i \stackrel{\text{iid}}{\sim} \text{Laplace}(\lambda_n)$ ,  
works great,  
but the full **posterior distribution** is useless.



# LASSO is not Bayesian

$$\hat{\theta}_{\text{LASSO}} = \underset{\theta}{\operatorname{argmin}} \left[ \|Y^n - X\theta\|^2 + \lambda_n \sum_{i=1}^p |\theta_i| \right].$$

posterior mode for prior  $\theta_i \stackrel{\text{iid}}{\sim} \text{Laplace}(\lambda_n)$ ,  
works great,  
but the full posterior distribution is useless.

**Theorem** If  $\sqrt{n}/\lambda_n \rightarrow \infty$  then

$$\mathbb{E}_0 \Pi_n(\|\theta\|_2 \lesssim \sqrt{n}/\lambda_n | Y^n) \rightarrow 0.$$

$\lambda_n = \sqrt{2 \log n}$  gives almost no “Bayesian shrinkage”.

**Trouble:**  $\lambda_n$  must be large to shrink  $\theta_i$  to 0, but small to model nonzero  $\theta_i$ .

# Frequentist Bayes

# Frequentist Bayes

Assume data  $Y^n$  follows a **given parameter**  $\theta_0$ .

Consider posterior  $\Pi(\theta \in \cdot | Y^n)$  as *random measure* on parameter set.

We like  $\Pi(\theta \in \cdot | Y^n)$ :

- to put “most” of its mass near  $\theta_0$  for “most”  $Y^n$ .
- to have a spread that expresses “remaining uncertainty”.
- to select the model defined by the nonzero parameters of  $\theta_0$ .

We evaluate this by probabilities or expectations, given  $\theta_0$ .

# Benchmarks for recovery — sequence model

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

$$\|\theta\|_0 = \#\{1 \leq i \leq n: \theta_i \neq 0\},$$

$$\|\theta\|_2^2 = \sum_{i=1}^n |\theta_i|^2.$$

*Frequentist benchmark*: **minimax rate** relative to  $\|\cdot\|_2$  over:

- **black bodies**  $\{\theta: \|\theta\|_0 \leq s_n\}$ :

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

# Benchmarks for recovery — sequence model

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

$$\|\theta\|_0 = \#\{1 \leq i \leq n: \theta_i \neq 0\},$$

$$\|\theta\|_q^q = \sum_{i=1}^n |\theta_i|^q, \quad 0 < q \leq 2.$$

*Frequentist benchmarks*: **minimax rate** relative to  $\|\cdot\|_q$  over:

- **black bodies**  $\{\theta: \|\theta\|_0 \leq s_n\}$ :

$$s_n^{1/q} \sqrt{\log(n/s_n)}.$$

# Benchmarks for recovery — sequence model

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

$$\|\theta\|_0 = \#\{1 \leq i \leq n: \theta_i \neq 0\},$$

$$\|\theta\|_q^q = \sum_{i=1}^n |\theta_i|^q, \quad 0 < q \leq 2.$$

*Frequentist benchmarks:* **minimax rate** relative to  $\|\cdot\|_q$  over:

- **black bodies**  $\{\theta: \|\theta\|_0 \leq s_n\}$ :

$$s_n^{1/q} \sqrt{\log(n/s_n)}.$$

- **weak  $\ell_r$ -balls**  $m_r[s_n] := \{\theta: \max_i i |\theta_{[i]}|^r \leq n(s_n/n)^r\}$ :

$$n^{1/q} (s_n/n)^{r/q} \sqrt{\log(n/s_n)}^{1-r/q}.$$

# Model Selection Prior

## Model selection prior

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

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

- (1) Choose  $s$  from prior  $\pi_n$  on  $\{0, 1, 2, \dots, n\}$ .
- (2) Choose  $S \subset \{0, 1, \dots, n\}$  of size  $|S| = s$  at random.
- (3) Choose  $\theta_S = (\theta_i: i \in S)$  from density  $g_S$  on  $\mathbb{R}^S$  and set  $\theta_{S^c} = 0$ .



# Model selection prior

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

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

- (1) Choose  $s$  from prior  $\pi_n$  on  $\{0, 1, 2, \dots, n\}$ .
- (2) Choose  $S \subset \{0, 1, \dots, n\}$  of size  $|S| = s$  at random.
- (3) Choose  $\theta_S = (\theta_i: i \in S)$  from density  $g_S$  on  $\mathbb{R}^S$  and set  $\theta_{S^c} = 0$ .

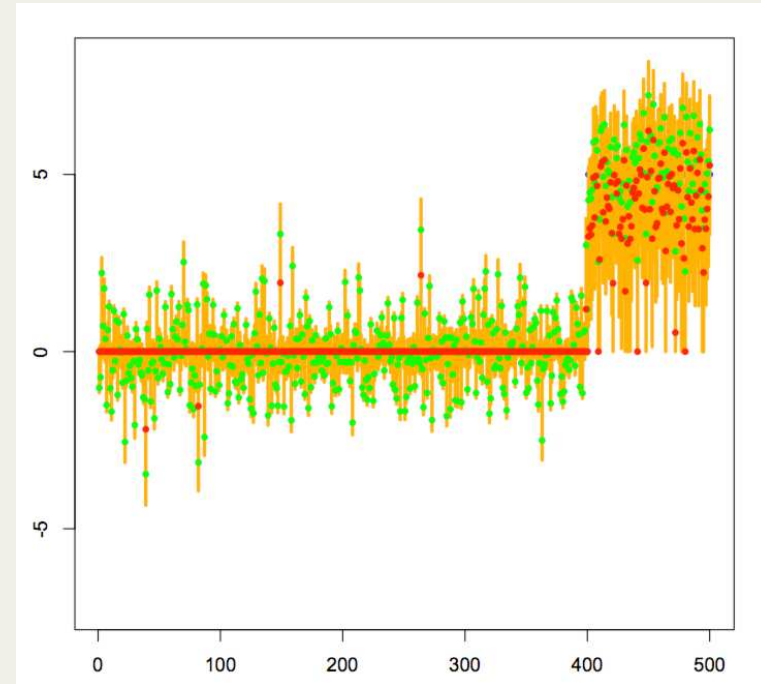
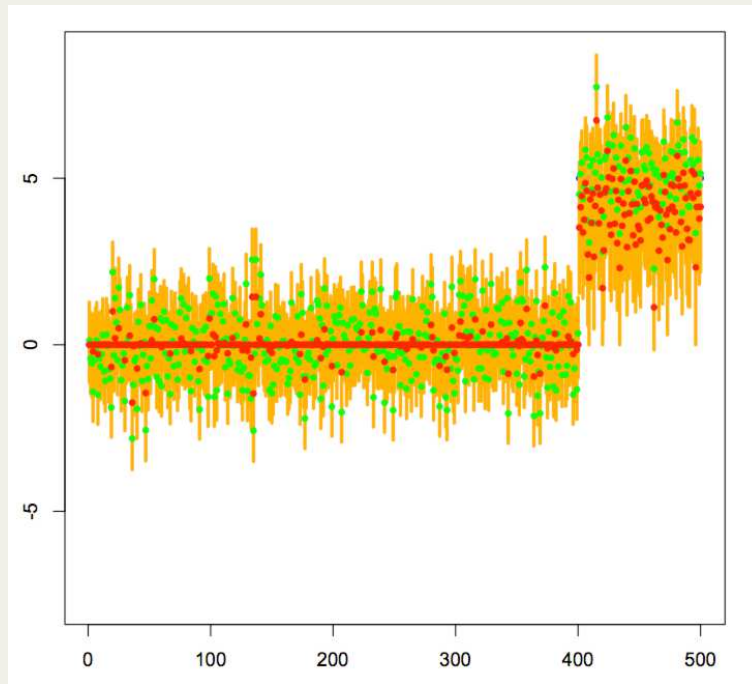
## Assume

- $\pi_n(s) \leq c \pi_n(s-1)$ , for some  $c < 1$  and every  $s$ .
- $g_S = \otimes_{i \in S} e^h$ , for uniformly Lipschitz  $h: \mathbb{R} \rightarrow \mathbb{R}$ .
- $s_n := \|\theta_0\|_0 \rightarrow \infty, n \rightarrow \infty, s_n/n \rightarrow 0$ .

## Examples:

- *complexity prior*:  $\pi_n(s) \propto e^{-as \log(bn/s)}$ .
- *spike and slab*:  $\theta_i \stackrel{\text{iid}}{\sim} \tau \delta_0 + (1 - \tau) \text{Lap}$  with  $\tau \sim B(1, n + 1)$ .

# Numbers



**Single data with  $\theta_0 = (0, \dots, 0, 5, \dots, 5)$  and  $n = 500$  and  $\|\theta_0\|_0 = 100$ .**

Red dots: marginal posterior medians

Orange: marginal credible intervals

Green dots: data points.

$g$  standard Laplace density.

$$\pi_n(k) \propto \binom{2n-k}{n}^{0.1} \text{ (left) and } \pi_n(k) \propto \binom{2n-k}{n} \text{ (right).}$$

# Dimensionality of posterior distribution

**Theorem** [*black body*]

There exists  $M$  such that

$$\sup_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: \|\theta\|_0 \geq Ms_n | Y^n) \rightarrow 0.$$

Outside the space in which  $\theta_0$  lives,  
the posterior is concentrated  
in low-dimensional subspaces along the coordinate axes.

# Recovery

## Theorem [black body]

For every  $0 < q \leq 2$  and large  $M$ ,

$$\sup_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: \|\theta - \theta_0\|_q > Mr_n s_n^{1/q-1/2} | Y^n) \rightarrow 0,$$

for  $r_n^2 = s_n \log(n/s_n) \vee \log(1/\pi_n(s_n))$ .

If  $\pi_n(s_n) \geq e^{-as_n \log(n/s_n)}$  minimax rate is attained.

# Selection

$$S_\theta := \{1 \leq i \leq n: \theta_i \neq 0\}.$$

**Theorem** *[No supersets]*

$$\sup_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: S_\theta \supset S_{\theta_0}, S_\theta \neq S_{\theta_0} | Y^n) \rightarrow 0.$$

# Selection

$$S_\theta := \{1 \leq i \leq n: \theta_i \neq 0\}.$$

**Theorem** *[No supersets]*

$$\sup_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: S_\theta \supset S_{\theta_0}, S_\theta \neq S_{\theta_0} | Y^n) \rightarrow 0.$$

**Theorem** *[Finds big signals]*

$$\inf_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: S_\theta \supset \{i: |\theta_{0,i}| \gtrsim \sqrt{\log n}\} | Y^n) \rightarrow 1.$$

# Selection

$$S_\theta := \{1 \leq i \leq n: \theta_i \neq 0\}.$$

**Theorem** *[No supersets]*

$$\sup_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: S_\theta \supset S_{\theta_0}, S_\theta \neq S_{\theta_0} | Y^n) \rightarrow 0.$$

**Theorem** *[Finds big signals]*

$$\inf_{\|\theta_0\|_0 \leq s_n} \mathbb{E}_{\theta_0} \Pi_n(\theta: S_\theta \supset \{i: |\theta_{0,i}| \gtrsim \sqrt{\log n}\} | Y^n) \rightarrow 1.$$

**Corollary:** if *all* nonzero  $|\theta_{0,i}|$  are suitably big,  
then posterior probability of true model  $S_{\theta_0}$  tends to 1.

# Bernstein-von Mises theorem

## Theorem

For spike-and-Laplace( $\lambda_n$ )-slab prior with  $\lambda_n \sqrt{\log n} / s_n \rightarrow 0$ , there are random weights  $\hat{w}_S$ ,

$$\mathbb{E}_{\theta_0} \left\| \Pi_n(\cdot | Y^n) - \sum_S \hat{w}_S N_{|S|}(Y_S^n, I) \otimes \delta_{S^c} \right\| \rightarrow 0.$$

## Theorem

Given consistent model selection, mixture can be replaced by  $N_{|S_0|}(Y_{S_0}, I) \otimes \delta_{S_0^c}$ .

**Corollary:** Given consistent model selection, credible sets for individual parameters are asymptotic confidence sets.



# Numbers: mean square errors

$p_n$	25			50			100		
	3	4	5	3	4	5	3	4	5
PM1	111	96	94	176	165	154	267	302	307
PM2	<b>106</b>	92	82	169	165	152	269	280	274
EBM	<b>103</b>	96	93	166	177	174	271	312	319
PMed1	129	<b>83</b>	73	205	149	<b>130</b>	<b>255</b>	279	283
PMed2	125	86	<b>68</b>	187	148	<b>129</b>	273	<b>254</b>	<b>245</b>
EBMed	110	<b>81</b>	72	<b>162</b>	148	142	<b>255</b>	294	300
HT	175	142	<b>70</b>	339	284	135	676	564	252
HTO	136	92	84	206	159	139	306	261	245

**Average  $\|\hat{\theta} - \theta\|^2$  over 100 data experiments.**

$n = 500; \theta_0 = (0, \dots, 0, A, \dots, A).$

*PM1, PM2*: posterior means for priors  $\pi_n(k) \propto e^{-k \log(3n/k)/10}, \binom{2n-k}{n}^{0.1}$ .

*PMed1, PMed2*: marginal posterior medians for the same priors

*EBM, EBMed*: empirical Bayes mean, median for Laplace prior (Johnstone et al.)

*HT, HTO*: thresholding at  $\sqrt{2 \log n}, \sqrt{2 \log(n/\|\theta_0\|_0)}$ .

*Short Summary: Bayesian method is neither better nor worse.*

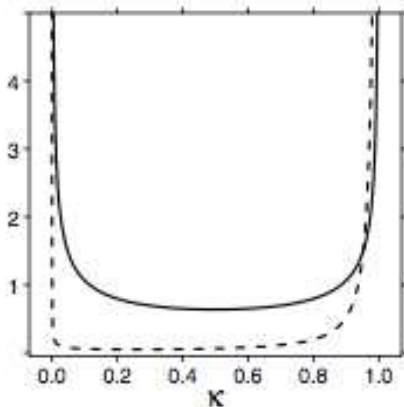
# Horseshoe Prior

# Horseshoe prior

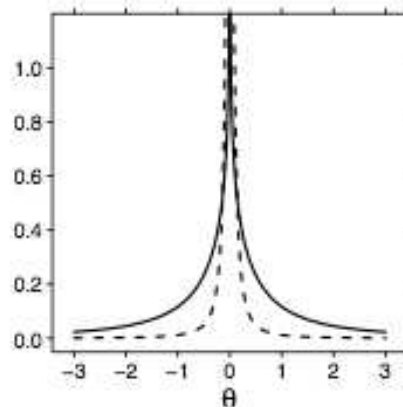
$$Y^n \sim N_n(\theta, I), \text{ for } \theta = (\theta_1, \dots, \theta_n) \in \mathbb{R}^n.$$

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

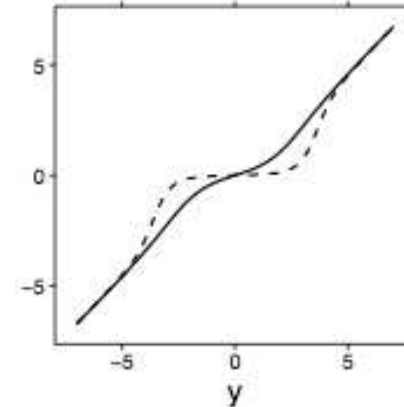
- (1) Choose “sparsity level”  $\hat{\tau}$ .
- (2) Generate  $\sqrt{\psi_1}, \dots, \sqrt{\psi_n}$  iid from  $\text{Cauchy}^+(0, \hat{\tau})$ .
- (3) Generate independent  $\theta_i \sim N(0, \psi_i)$ .



prior shrinkage factor



prior of  $\theta_i$



posterior mean of  $\theta_i$  as function of  $Y_i$

# Estimating $\tau$

Ad-hoc:

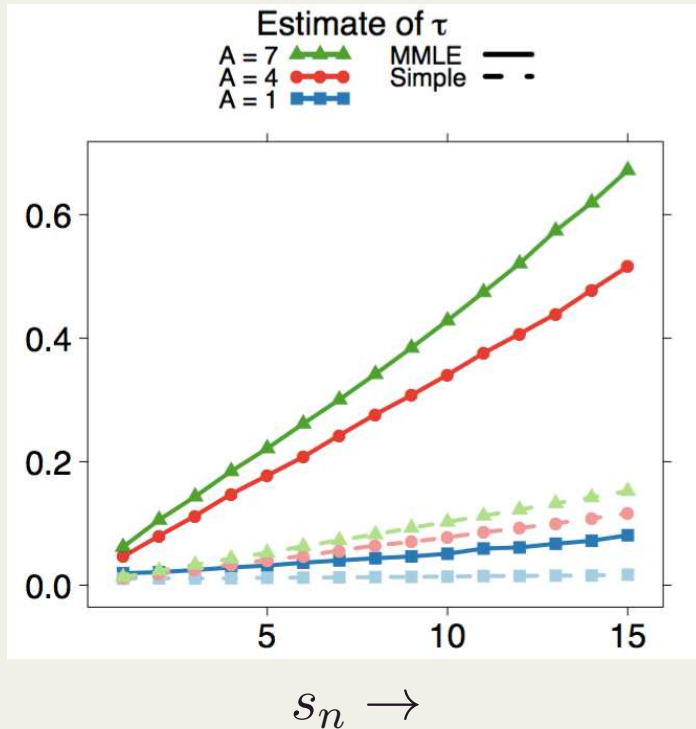
$$\hat{\tau}_n = \frac{\#\{|Y_i^n| \geq \sqrt{2 \log n}\}}{1.1n}.$$

Empirical Bayes: For  $g_\tau$  the prior of  $\theta_i$ ,

$$\hat{\tau}_n = \operatorname{argmax}_{\tau \in [1/n, 1]} \prod_{i=1}^n \int \phi(y_i - \theta) g_\tau(\theta) d\theta.$$

Full Bayes:  $\tau$  set by a “hyper prior” (supported on  $[1/n, 1]$ ).

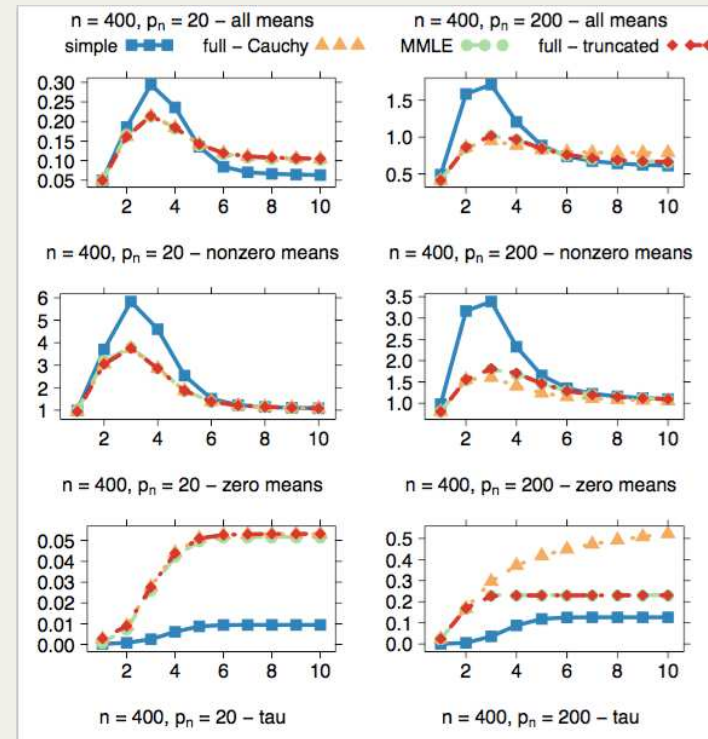
estimating  $\tau$



$n = 100$ ,  $s_n$  coordinates from  $N(0, 1/4)$ ,

$n - s_n$  coordinates from  $N(A, 1)$ .

MSE of posterior mean  
as function of nonzero parameter



" $p_n = s_n$ "

Short summary:  
Empirical Bayes and Full Bayes outperform ad-hoc estimator.

# Recovery

Horseshoe prior gives similar recovery as model selection prior.

# Recovery

Horseshoe prior gives similar recovery as model selection prior.

$\tau$  can be interpreted as  $(s_n/n) \sqrt{\log(n/s_n)}$ .

# Credible intervals

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 intervals

## 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 intervals

## Credible interval:

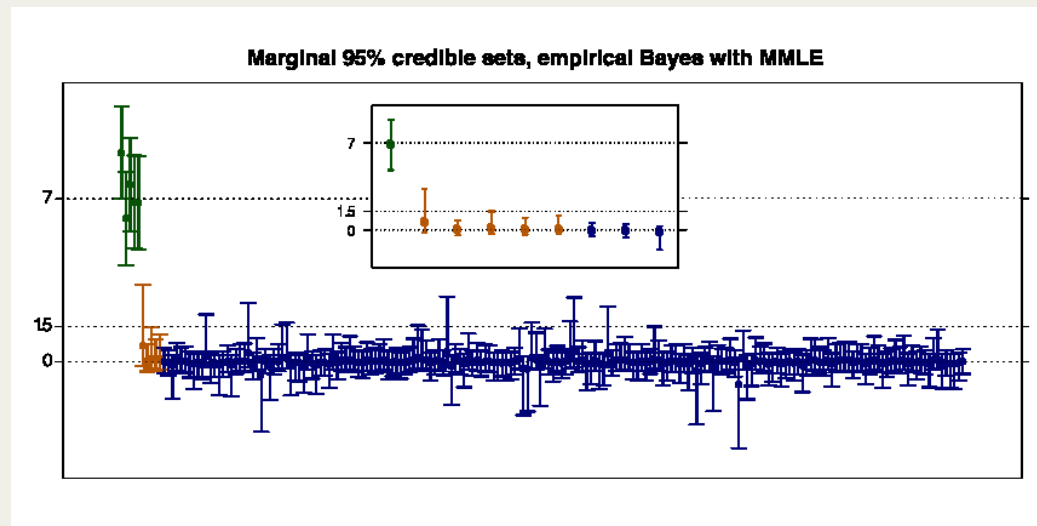
$$\hat{C}_{ni}(L) = \left\{ \theta_i : |\theta_i - \hat{\theta}_i| \leq L \hat{r}_i \right\}$$

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

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

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

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



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 intervals

## Credible interval:

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

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

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

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

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

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

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

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

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

Few false discoveries; most easy discoveries made.  
Intermediate discoveries not made.

# Simultaneous credible balls — impossibility of adaptation

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

# Simultaneous credible balls — impossibility of adaptation

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

**Theorem** [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  $\theta_0$ .

# Simultaneous credible balls — impossibility of adaptation

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

**Theorem** [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  $\theta_0$ .

**Theorem** [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  $\{\theta_0: \|\theta_0\|_0 \leq s_{2,n}\}$ .

# Simultaneous credible balls — impossibility of adaptation

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

**Theorem** [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  $\theta_0$ .

**Theorem** [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 sets *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.

**DEFINITION** [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.

**DEFINITION** [self-similarity]

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

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

**DEFINITION** [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 implies self-similarity.  
(Self-similarity allows to tighten up the sets **S**, **M**, **L**.)

# Simultaneous credible balls

Credible ball:

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

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

# Simultaneous credible balls

Credible ball:

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

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

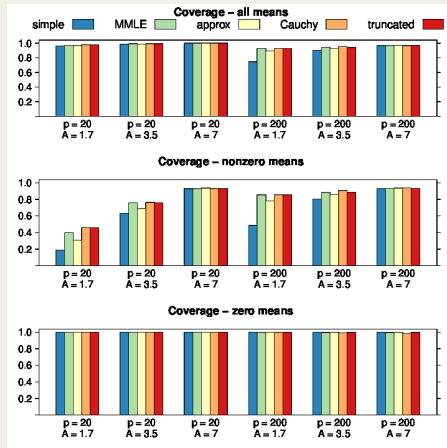
**Theorem**

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

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

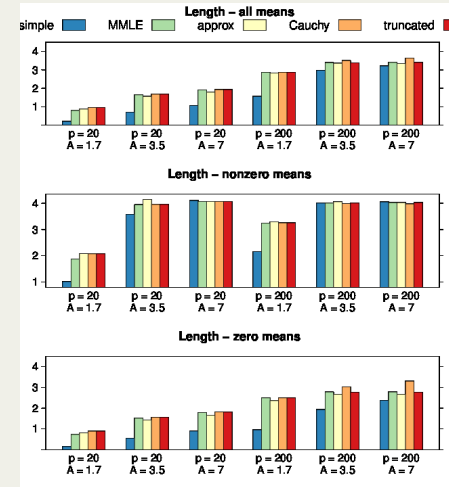
EBR[ $s$ ]: vectors  $\theta_0$  that satisfy excessive bias restriction.

## coverage



$n = 400$ .  $s_n$  ("=  $p$ ") nonzero means from  $\mathcal{N}(A, 1)$ .

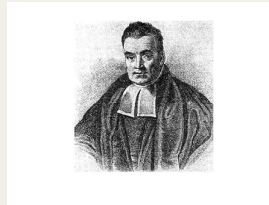
## average interval length



$n = 400$ .  $s_n$  ("=  $p$ ") nonzero means from  $\mathcal{N}(A, 1)$ .

Short summary:  
empirical and full Bayes work well

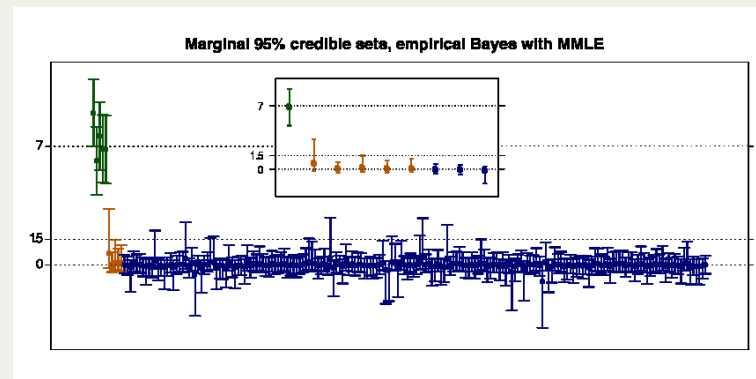
# Conclusions



Bayesian sparse estimation gives excellent recovery.

For valid simultaneous credible sets need a fraction of nonzero parameters above the “universal threshold”.

The danger of failing uncertainty quantification is *not* finding nonzero coordinates.  
Discoveries are real.



# Co-authors



Subhashis Ghoshal



Ismael Castillo



Stéphanie van der Pas



Willem Kruijer



Bartek Knapik



Suzanne Sniekers



Gwenael Leday



Mark van de Wiel



Harry van Zanten



Johannes Schmidt-Hieber



Botond Szabo



Magnus Münch



Bas Kleijn



Fengnan Gao



Gino Kpogbezan



Wessel van Wieringen

# Compatibility and coherence

$$\|X\| := \max_j \|X_{\cdot,j}\|.$$

*Compatibility number*  $\phi(S)$  for  $S \subset \{1, \dots, p\}$  is: 
$$\inf_{\|\theta_{S^c}\|_1 \leq \|\theta_S\|_1} \frac{\|X\theta\|_2 \sqrt{|S|}}{\|X\| \|\theta_S\|_1}.$$

*Compatibility in  $s_n$ -sparse vectors* means: 
$$\inf_{\theta: \|\theta\|_0 \leq s_n} \frac{\|X\theta\|_2 \sqrt{|S_\theta|}}{\|X\| \|\theta\|_1} \gg 0.$$

*Strong compatibility in  $s_n$ -sparse vectors* means: 
$$\inf_{\theta: \|\theta\|_0 \leq s_n} \frac{\|X\theta\|_2}{\|X\| \|\theta\|_2} \gg 0.$$

*Mutual coherence* means: 
$$s_n \max_{i \neq j} |\text{cor}(X_{\cdot,i}, X_{\cdot,j})| \ll 1.$$

## Compatibility and coherence — examples

Mutual coherence  $\Rightarrow$  Strong compatibility  $\Rightarrow$  Compatibility.

*Mutual coherence is easy to understand and gives best recovery results, but is very restrictive.*



## Compatibility and coherence — examples

Mutual coherence  $\Rightarrow$  Strong compatibility  $\Rightarrow$  Compatibility.

*Mutual coherence is easy to understand and gives best recovery results, but is very restrictive.*

Sequence model:

Completely compatible, with zero mutual coherence number.

# Compatibility and coherence — examples

Mutual coherence  $\Rightarrow$  Strong compatibility  $\Rightarrow$  Compatibility.

*Mutual coherence is easy to understand and gives best recovery results, but is very restrictive.*

Sequence model:

Completely compatible, with zero mutual coherence number.

Response model:

If  $X_{i,j}$  are i.i.d. random variables, then coherence if  $s_n \lesssim \sqrt{n / \log p}$ .

- if  $\log p = o(n)$  and  $X_{i,j}$  are bounded.
- if  $\log p = o(n^{\alpha/(4+\alpha)})$  and  $E^{tX_{i,n}^\alpha} < \infty$ .

# Compatibility and coherence — examples

Mutual coherence  $\Rightarrow$  Strong compatibility  $\Rightarrow$  Compatibility.

*Mutual coherence is easy to understand and gives best recovery results, but is very restrictive.*

Sequence model:

Completely compatible, with zero mutual coherence number.

Response model:

If  $X_{i,j}$  are i.i.d. random variables, then coherence if  $s_n \lesssim \sqrt{n/\log p}$ .

- if  $\log p = o(n)$  and  $X_{i,j}$  are bounded.
- if  $\log p = o(n^{\alpha/(4+\alpha)})$  and  $E^{tX_{i,n}^\alpha} < \infty$ .

$C = X^T X/n$ : Compatibility, but no coherence if

- $C_{i,j} = \rho^{|i-j|}$ , for  $0 < \rho < 1$ , and  $p = n$ .
- $C$  is block diagonal with fixed block sizes.