

Adaptive Bayesian estimation and its self-informative limit in an indirect sequence space model

Xavier Loizeau, joint work with Jan Johannes

The Gaussian sequence space model

Consider an indirect Gaussian sequence space model consisting of:

- ▶ an unknown parameter of interest $(\theta_j^\circ)_{j \in \mathbb{N}} = \theta^\circ$,
- ▶ a decreasing multiplicative sequence $(\lambda_j)_{j \in \mathbb{N}} = \lambda$ converging to 0,
- ▶ observations $(Y_j)_{j \in \mathbb{N}} = Y$, contaminated by an additive independent centered Gaussian noise with variance n^{-1} ,

$$Y = \left(\theta_j^\circ \cdot \lambda_j + \sqrt{n^{-1}} \cdot \xi_j \right)_{j \in \mathbb{N}}, \quad (\xi_j)_{j \in \mathbb{N}} \sim_{iid} \mathcal{N}(0, 1).$$

The goal is to recover θ° and derive an upper bound.

The frequentist model selection

For any index j , an unbiased estimator of θ_j° is Y_j/λ_j . Hence, an intuitive class of estimators are the projection estimators:

$\hat{\theta}^m = \left(Y_j/\lambda_j \mathbb{1}_{\{j \leq m\}} \right)_{j \in \mathbb{N}}$ with m in \mathbb{N} . The model selection method offers a data driven way to select m in this context:

$$G_n := \max \left\{ 1 \leq j \leq n : n^{-1} \lambda_j^{-2} \leq \lambda_j^{-2} \right\},$$

$$\hat{m} := \arg \min_{m \in \llbracket 1, G_n \rrbracket} \left\{ 3m - \sum_{j=1}^m Y_j^2 \right\}, \quad \hat{\theta} := (\hat{\theta}_j^{\hat{m}})_{j \in \mathbb{N}}.$$

It is shown in Massart [2003], in the direct case, that this estimator is **consistent**, converges in probability and \mathbb{L}^2 -norm, noted $\|\cdot\|$, with **minimax optimal rate** over some Sobolev ellipsoid:

$$\Theta^\circ := \Theta^\circ(a, L^\circ) \left\{ \theta : \sum_{j=1}^\infty \frac{1}{a_j} \theta_j^2 < L^\circ \right\}.$$

Bayesian paradigm, iterated posterior distribution and self informative limit

We adopt a **Bayesian point of view**:

- ▶ the parameter θ is a random variable with prior \mathbb{P}_θ ,
- ▶ given θ , the likelihood of Y is $\mathbb{P}_Y^n | \theta = \mathcal{N}(\theta \lambda, n^{-1} \mathbb{I})$,
- ▶ we are interested in the posterior distribution $\mathbb{P}_{\theta^n | Y} \propto \mathbb{P}_Y^n | \theta \cdot \mathbb{P}_\theta$.

In the spirit of Bunke and Johannes [2005], we then generate a posterior family by introducing an **iteration parameter η** :

- ▶ for $\eta = 1$, the prior distribution is $\mathbb{P}_{\theta^1} = \mathbb{P}_\theta$, the likelihood $\mathbb{P}_{Y^1 | \theta^1} = \mathbb{P}_Y^n | \theta$ and the posterior distribution is $\mathbb{P}_{\theta^1 | Y^1} = \mathbb{P}_{\theta | Y}^n$,
- ▶ for $\eta = 2$, we take the posterior for $\eta = 1$ as prior, hence, the prior distribution is $\mathbb{P}_{\theta^2} = \mathbb{P}_{\theta^1 | Y^1}$, the likelihood is kept the same $\mathbb{P}_{Y^2 | \theta^2} = \mathbb{P}_Y^n | \theta$ and we compute the posterior distribution with the same observations Y , which we note $\mathbb{P}_{\theta^2 | Y^2}$,

- ▶ ...
- ▶ for any value of $\eta > 1$, the prior is $\mathbb{P}_{\theta^\eta} = \mathbb{P}_{\theta^{\eta-1} | Y^{\eta-1}}$ and we compute the posterior with the same likelihood $\mathbb{P}_{Y^\eta | \theta^\eta} = \mathbb{P}_Y^n | \theta$ and same observation Y which gives $\mathbb{P}_{\theta^\eta | Y^\eta}$.

This iteration procedure corresponds to giving more and more weight to the observations and make the prior knowledge vanish.

Within this framework we define the family of estimators:

$$\hat{\theta}^{(\eta)} := \mathbb{E}_{\theta^\eta | Y^\eta}[\theta],$$

and call **self-informative limit** the limit of the estimate with $\eta \rightarrow \infty$.

We are interested in the behavior of the family $(\mathbb{P}_{\theta^\eta | Y^\eta})_{\eta \in \mathbb{N}^*}$ as n and/or η tend to infinite.

In particular, the question of oracle and minimax concentration (resp. convergence) is answered for any element of the family of posterior distributions (resp. posterior means), including when η tends to infinite.

Hierarchical prior

- ▶ Consider a **random hyper-parameter M** , with values in a subset of \mathbb{N} , acting like a threshold:

$$\forall j > m, \quad \mathbb{P}_{\theta_j | M=m} = \delta_0,$$

$$\forall j \leq m, \quad \mathbb{P}_{\theta_j | M=m} = \mathcal{N}(0, 1).$$

- ▶ if we denote \mathbb{P}_M the distribution of M (to be specified later), then

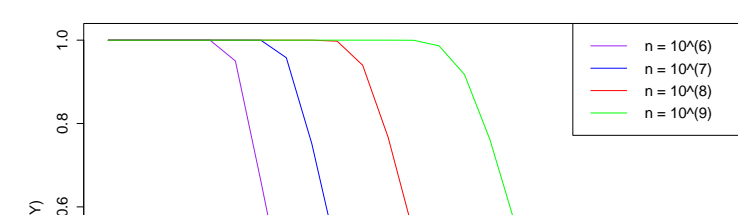
$$\mathbb{P}_{\theta | Y}^n = \sum_{m \in \mathbb{N}} \mathbb{P}_{\theta | M=m, Y}^n \cdot \mathbb{P}_{M=m | Y}^n.$$

- ▶ Hence, given M , the posterior is

$$\forall j > m, \quad \theta_j | M = m, Y \sim \delta_0,$$

$$\forall j \leq m, \quad \theta_j | M = m, Y \sim \mathcal{N}\left(\frac{Y_j \cdot n \cdot \lambda_j}{1 + n \cdot \lambda_j^2}, \frac{1}{1 + n \cdot \lambda_j^2}\right).$$

Remark: the family of hierarchical priors with deterministic threshold M is called family of sieve priors.



Existing results

In Johannes et al. [2016], under a **pragmatic Bayesian** point of view; that is, the existence of a true parameter θ° is accepted; it is shown that, by choosing \mathbb{P}_M suitably:

- ▶ the estimator $\hat{\theta}^{(1)}$ **converges with**,
- ▶ **oracle optimal rate** for the quadratic risk which means,

$$\forall \theta^\circ \in \Theta^\circ, \exists C^\circ \in [1, \infty[: \forall n \in \mathbb{N}, \exists \Phi_n^\circ \in \mathbb{R} :$$

$$\inf_{m \in \mathbb{N}} \mathbb{E}_{\theta^\circ}^n \left[\left\| \hat{\theta}^m - \theta^\circ \right\|^2 \right] \geq \Phi_n^\circ,$$

$$\mathbb{E}_{\theta^\circ}^n \left[\left\| \hat{\theta}^{(1)} - \theta^\circ \right\|^2 \right] \leq C^\circ \Phi_n^\circ;$$

- ▶ **minimax optimal rate** for the maximal risk over Θ° , that is to say, $\exists C^* \in [1, \infty[: \forall n \in \mathbb{N}, \exists \Phi_n^* \in \mathbb{R} :$

$$\inf_{\tilde{\theta}} \sup_{\theta^\circ \in \Theta^\circ} \mathbb{E}_{\tilde{\theta}^\circ}^n \left[\left\| \tilde{\theta} - \theta^\circ \right\|^2 \right] \geq \Phi_n^*,$$

$$\sup_{\theta^\circ \in \Theta^\circ} \mathbb{E}_{\theta^\circ}^n \left[\left\| \hat{\theta}^{(1)} - \theta^\circ \right\|^2 \right] \leq C^* \Phi_n^*,$$

where $\inf_{\tilde{\theta}}$ is taken over all possible estimators of θ° ;

- ▶ the posterior distribution **concentrates with**,
- ▶ **oracle optimal rate** for the quadratic loss which means,

$$\forall \theta^\circ \in \Theta^\circ, \exists K^\circ \in [1, \infty[:$$

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{P}_{\theta^1 | Y^1} \left(\left\| \theta - \theta^\circ \right\|^2 \leq K^\circ \Phi_n^\circ \right) \right] = 1;$$

- ▶ **minimax optimal rate** Θ° , that is to say, for any unbounded sequence $K_n \in \mathbb{R}^{\mathbb{N}} :$

$$\lim_{n \rightarrow \infty} \sup_{\theta^\circ \in \Theta^\circ} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{P}_{\theta^1 | Y^1} \left(\left\| \theta - \theta^\circ \right\|^2 \leq K_n \Phi_n^* \right) \right] = 1.$$

Iterated posterior distributions

Note that in the framework of our hierarchical prior, we have:

$$\mathbb{P}_{\theta^\eta | Y^\eta}^n = \sum_{m \in \mathbb{N}} \mathbb{P}_{\theta^\eta | M^\eta = m, Y^\eta}^n \cdot \mathbb{P}_{M^\eta = m | Y^\eta}^n,$$

$$\hat{\theta}^{(\eta)} = \left(\mathbb{E}_{\theta^\eta | M^\eta \geq j, Y^\eta}^n [\theta_j] \cdot \mathbb{P}_{M^\eta | Y^\eta}^n (M^\eta \geq j) \right)_{j \in \mathbb{N}}.$$

Hence, we first compute $\theta_j^\eta | M^\eta, Y^\eta$:

$$\forall j \in \mathbb{N}, \quad \theta_j^\eta | M^\eta \geq j, Y^\eta \sim \mathcal{N}\left(\frac{\eta \cdot Y_j \cdot n \cdot \lambda_j}{1 + \eta \cdot n \cdot \lambda_j^2}, \frac{1}{1 + \eta \cdot n \cdot \lambda_j^2}\right),$$

$$\theta_j^\eta | M^\eta < j, Y^\eta \sim \delta_0;$$

and then fix the distribution of M^1 : $\forall m \in \llbracket 1, G_n \rrbracket$,

$$\mathbb{P}_{M^1}(M = m) \propto \exp\left(-3 \cdot \eta \cdot \frac{m}{2}\right) \cdot \prod_{j=1}^m \left(1 + n \cdot \eta \cdot \lambda_j^2\right)^2.$$

Which gives the family of posterior distributions:

$$\mathbb{P}_{M^\eta | Y^\eta}^n(m) \propto \exp\left[-\frac{\eta}{2} \left(3m - \sum_{j=1}^m \frac{\eta \left(Y_j \cdot n \cdot \lambda_j^2\right)^2}{1 + \eta \cdot n \cdot \lambda_j^2}\right)\right].$$

Self informative limit and model selection

Consider the limit of the family of posteriors as η tends to infinite:

$$\lim_{\eta \rightarrow \infty} \mathbb{P}_{\theta^\eta | M^\eta = m, Y^\eta}^n = \delta_{\hat{\theta}^m},$$

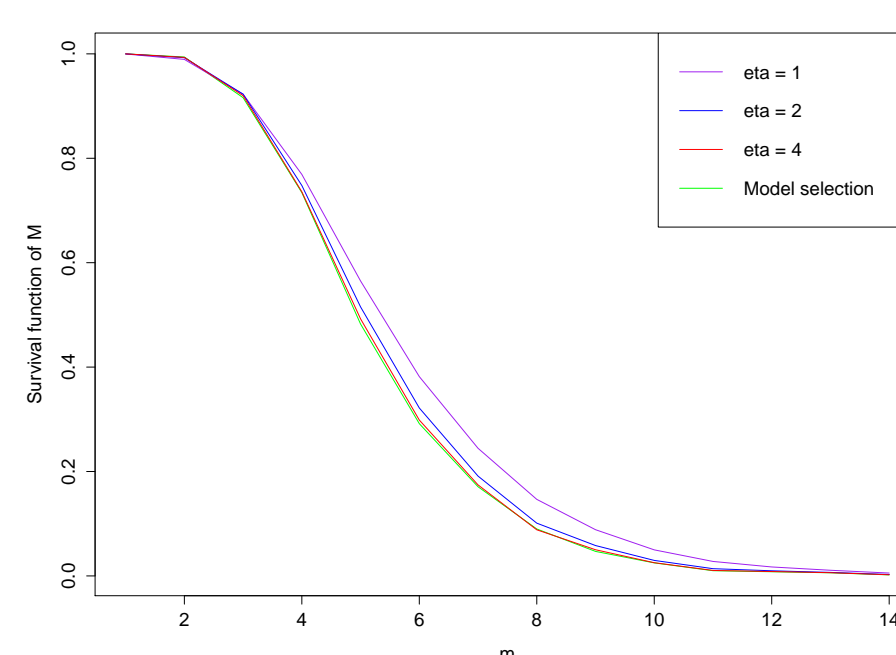
where $\hat{\theta}^m$ is the projection estimator on the first m dimensions. The distribution of M tends to a point mass:

$$\lim_{\eta \rightarrow \infty} \mathbb{P}_{M^\eta | Y^\eta}^n = \delta_{\hat{m}},$$

where \hat{m} is the choice given by the frequentist model selection presented earlier.

The **self-informative limit** is equal to the model selection estimator, $\hat{\theta}$, presented above.

Figure: Survival function of M for different values of η



Notations

Define the following quantities:

$$\mathfrak{b}_m := \sum_{j=m+1}^\infty (\theta^\circ)^2, \quad \Lambda_j := \lambda_j^{-2}, \quad m \cdot \bar{\Lambda}_m := \sum_{j=1}^m \Lambda_j,$$

$$m_n^\circ := \arg \min_{m \in \llbracket 1, G_n \rrbracket} \left[\mathfrak{b}_m \vee n^{-1} m \bar{\Lambda}_m \right], \quad \Phi_n^\circ := \left[\mathfrak{b}_{m_n^\circ} \vee n^{-1} m_n^\circ \bar{\Lambda}_{m_n^\circ} \right],$$

Set of assumptions

Define the following assumptions:

$$(\mathbb{H}_\lambda) : \exists a \in \mathbb{R}_+, c \geq 1 : \quad \forall j \in \mathbb{N}, \quad \left(\frac{1}{c} j^{-a} \leq \lambda_j \leq c j^{-a} \right)$$

$$(\mathbb{H}_1) : 0 < \inf_{n \in \mathbb{N}} \left\{ \frac{\left[\mathfrak{b}_{m_n^\circ} \wedge n^{-1} m_n^\circ \bar{\Lambda}_{m_n^\circ} \right]}{\left[\mathfrak{b}_{m_n^\circ} \vee n^{-1} m_n^\circ \bar{\Lambda}_{m_n^\circ} \right]} \right\} \leq 1$$

$$(\mathbb{H}_2) : 0 < \inf_{n \in \mathbb{N}} \left\{ \frac{\left[\mathfrak{a}_{m_n^*} \wedge n^{-1} m_n^* \bar{\Lambda}_{m_n^*} \right]}{\left[\mathfrak{a}_{m_n^*} \vee n^{-1} m_n^* \bar{\Lambda}_{m_n^*} \right]} \right\} \leq 1$$

Note that under (\mathbb{H}_λ) , there exist a constant L such that,

$$\forall m \in \mathbb{N}, \quad \Lambda_m \leq L \bar{\Lambda}_m.$$

Concentration results for the threshold parameter M

For any η in $\overline{\mathbb{N}}$, we have the following results:

- Under assumptions (\mathbb{H}_1) and (\mathbb{H}_λ) , define

$$G_n^- := \min \left\{ m \in \llbracket 1, m_n^\circ \rrbracket : \mathfrak{b}_m \leq 9L \Phi_n^\circ \right\},$$

$$G_n^+ := \max \left\{ m \in \llbracket m_n^\circ, G_n \rrbracket : (m - m_n^\circ) n^{-1} \leq 3 \Lambda_{m_n^\circ}^{-1} \Phi_n^\circ \right\},$$

and we then have the following concentration for M ,

$$\mathbb{P}_{M^\eta | Y^\eta}^n [M > G_n^+] \leq \exp \left[-\frac{5m_n^\circ}{9L} + \log(G_n) \right],$$

$$\mathbb{P}_{M^\eta | Y^\eta}^n [M < G_n^-] \leq \exp \left[-\frac{7m_n^\circ}{9} + \log(G_n) \right],$$

this means that M^η tends to select an oracle optimal threshold;

- whereas under (\mathbb{H}_2) and (\mathbb{H}_λ) , we define

$$G_n^{*-} := \min \left\{ m \in \llbracket 1, m_n^* \rrbracket : \mathfrak{b}_m \leq 9(1 \vee L^\circ) L \Phi_n^* \right\},$$

$$G_n^{*+} := \max \left\{ m \in \llbracket m_n^*, G_n \rrbracket : (m - m_n^*) n^{-1} \leq 3 \Lambda_{m_n^*}^{-1} (1 \vee L^\circ) \Phi_n^* \right\},$$

and the following concentration stands,

$$\mathbb{P}_{M^\eta | Y^\eta}^n [M > G_n^{*+}] \leq \exp \left[-\frac{5(1 \vee L^\circ) m_n^*}{9L} + \log(G_n) \right],$$

$$\mathbb{P}_{M^\eta | Y^\eta}^n [M < G_n^{*-}] \leq \exp \left[-\frac{7(1 \vee L^\circ) m_n^*}{9} + \log(G_n) \right],$$

which means that M^η tends to select a minimax optimal threshold.

Concentration results for θ

For any η in \mathbb{N} , we have the following results:

- under assumptions (\mathbb{H}_1) and (\mathbb{H}_λ) , for all θ° in Θ° , there exist $K^\circ \geq 1$ and $C^\circ > 1$ such that we have

$$\lim_{n \rightarrow \infty} \inf_{\mathbb{Q}_\theta} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{Q}_{\theta | Y}^n \left(\left\| \theta - \theta^\circ \right\|^2 \geq \Phi_n^\circ \right) \right] = 1,$$

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{P}_{\theta^\eta, M^\eta | Y^\eta}^n \left((K^\circ)^{-1} \Phi_n^\circ \leq \left\| \theta - \theta^\circ \right\|^2 \leq K^\circ \Phi_n^\circ \right) \right] = 1,$$

$$\mathbb{E}_{\theta^\circ}^n \left[\left\| \hat{\theta}^{(\eta)} - \theta^\circ \right\|^2 \right] \leq C^\circ \Phi_n^\circ,$$
- whereas under (\mathbb{H}_2) and (\mathbb{H}_λ) , for a finite constant $C^* \geq 1$ and any unbounded sequence K_n , we have

$$\lim_{n \rightarrow \infty} \inf_{\mathbb{Q}_\theta} \sup_{\theta^\circ \in \Theta^\circ} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{Q}_{\theta | Y}^n \left(\left\| \theta - \theta^\circ \right\|^2 \geq \Phi_n^* \right) \right] = 1,$$

$$\lim_{n \rightarrow \infty} \sup_{\theta^\circ \in \Theta^\circ} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{P}_{\theta^\eta, M^\eta | Y^\eta}^n \left(\left\| \theta - \theta^\circ \right\|^2 \leq K_n \Phi_n^* \right) \right] = 1,$$

$$\sup_{\theta^\circ \in \Theta^\circ} \mathbb{E}_{\theta^\circ}^n \left[\left\| \hat{\theta}^{(\eta)} - \theta^\circ \right\|^2 \right] \leq C^* \Phi_n^*,$$

where $\inf_{\mathbb{Q}_\theta}$ is taken over all possible sieve priors; **establishing oracle optimal concentration and convergence of the posterior and Bayes estimate, respectively**;

- whereas under (\mathbb{H}_2) and (\mathbb{H}_λ) , for a finite constant $C^* \geq 1$ and any unbounded sequence K_n , we have

Note that in the case of $\eta \rightarrow \infty$, those results are still true and that the concentration corresponds to the convergence in probability as

$$\lim_{\eta \rightarrow \infty} \mathbb{E}_{\theta^\circ}^n \left[\mathbb{P}_{\theta^\eta, M^\eta | Y^\eta}^n \left(\left\| \theta - \theta^\circ \right\|^2 \leq K_n \Phi_n \right) \right] = \mathbb{P}_{\theta^\circ}^n \left[\left\| \hat{\theta} - \theta^\circ \right\|^2 \leq K_n \Phi_n \right].$$

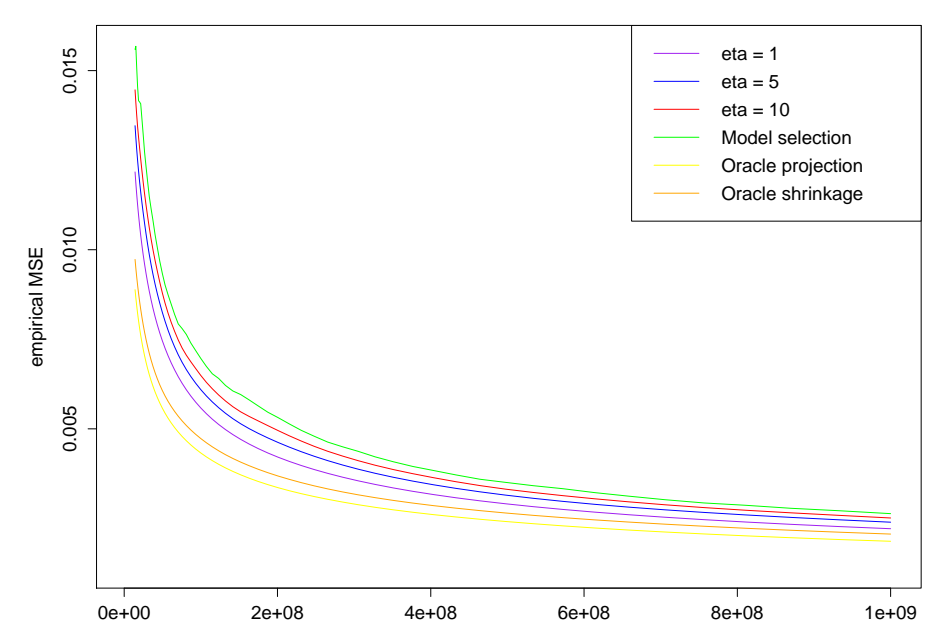


Figure: Estimated mean of the quadratic error of the Bayes estimate for θ° polynomial.

Bibliography

Olaf Bunke and Jan Johannes. Selfinformative limits of bayes estimates and generalized maximum likelihood. *Statistics*, 39(6):483–502, July 2005.