

# Supplementary material for Segmentation of the mean of heteroscedastic data via cross-validation

Sylvain Arlot and Alain Celisse

September 20, 2010

## 1 Proofs of theoretical results

For completeness, we report here proofs (or sketch of proofs) of the theoretical results stated in the main paper, that are coming from previous works.

### 1.1 Proof of Lemma 1 of the main paper

Let us first state the following lemma.

**Lemma 1.**

$$\mathbb{E} [S_{\lambda,1}^2] = n_\lambda (\sigma_\lambda^r)^2 + n_\lambda^2 \beta_\lambda^2, \quad \text{and} \quad \mathbb{E} [S_{\lambda,2}] = n_\lambda \left( (\sigma_\lambda^r)^2 + (\sigma_\lambda^b)^2 + \beta_\lambda^2 \right),$$

with

$$\begin{aligned} S_{\lambda,1} &= \frac{\sum_{i=1}^n Y_i \mathbf{1}_\lambda(t_i)}{n_\lambda}, \quad \text{and} \quad S_{\lambda,2} = \frac{\sum_{i=1}^n Y_i^2 \mathbf{1}_\lambda(t_i)}{n_\lambda}, \\ (\sigma_\lambda^r)^2 &= \frac{\sum_{i=1}^n \sigma(t_i)^2 \mathbf{1}_\lambda(t_i)}{n_\lambda}, \quad \text{and} \quad (\sigma_\lambda^b)^2 = \frac{\sum_{i=1}^n (s(t_i) - \beta_\lambda)^2 \mathbf{1}_\lambda(t_i)}{n_\lambda}, \\ \beta_\lambda &= \frac{\sum_{i=1}^n s(t_i) \mathbf{1}_\lambda(t_i)}{n_\lambda}, \quad \text{and} \quad \hat{\beta}_\lambda = \frac{\sum_{i=1}^n Y_i \mathbf{1}_\lambda(t_i)}{n_\lambda}, \end{aligned}$$

where  $n_\lambda = \text{Card}(\{j \mid t_j \in \lambda\})$ .

With the notation of Lemma 1, simple calculations lead to

$$\begin{aligned} \mathbb{E} [P_n \gamma(\hat{s}_m)] &= \sum_{\lambda \in \Lambda_m} \left\{ \frac{1}{n} \sum_{i=1}^n \left( \sigma(t_i)^2 + \mathbb{E} [s(t_i) - \hat{\beta}_\lambda]^2 \right) \mathbf{1}_\lambda(t_i) + \frac{1}{n} \sum_{i=1}^n \mathbb{E} [2\varepsilon_i (s(t_i) - \hat{\beta}_\lambda)] \mathbf{1}_\lambda(t_i) \right\} \\ &= \sum_{\lambda \in \Lambda_m} \left\{ \frac{1}{n} \left( n_\lambda (\sigma_\lambda^r)^2 + \sum_{i=1}^n \mathbb{E} [s(t_i) - \hat{\beta}_\lambda]^2 \mathbf{1}_\lambda(t_i) \right) - \frac{2}{n} (\sigma_\lambda^r)^2 \right\}. \end{aligned}$$

Moreover, by use of Lemma 1, it comes

$$\mathbb{E} [s(t_i) - \hat{\beta}_\lambda]^2 = [s(t_i) - s_m(t_i)]^2 + \frac{1}{n_\lambda} (\sigma_\lambda^r)^2,$$

which enables to conclude

$$\mathbb{E} [P_n \gamma(\widehat{s}_m)] = \|s - s_m\|_n^2 - \frac{1}{n} \sum_{\lambda \in \Lambda_m} (\sigma_\lambda^r)^2 + \frac{1}{n} \sum_{i=1}^n \sigma(t_i)^2 .$$

The expectation of the loss of  $\widehat{s}_m$  results from similar calculations.

## 1.2 Proof of Theorem 1 of the main paper

Since the proof of Theorem 1 of the main paper is quite technical, it is not reproduced here. For the detailed calculations, we refer interested readers to the Ph.D. thesis of the second author [6, Section 3.5.3], which is freely available online.

## 1.3 Proof of Proposition 1 of the main paper

On the basis of Theorem 1 of the main paper and assuming that for every  $\lambda \in \Lambda_m$ ,  $n_\lambda \geq 3$  and  $p \leq n_\lambda + 1$ , it comes that the expectation of the Lpo risk is

$$\begin{aligned} \mathbb{E} \left[ \widehat{R}_{\text{Lpo}_p} \right] &= \sum_{\lambda \in \Lambda_m} \frac{1}{p} \left( (\sigma_\lambda^r)^2 + \frac{n_\lambda}{n_\lambda - 1} (\sigma_\lambda^b)^2 \right) \left[ n_\lambda - 1 - \frac{n_\lambda(n-p)}{n} + \frac{n}{n-p} V_\lambda(1)V_\lambda(-1) \right] . \\ &= \sum_{\lambda \in \Lambda_m} (\sigma_\lambda^r)^2 \left[ \frac{n_\lambda}{n} + \frac{1}{p} \left( \frac{n}{n-p} V_\lambda(1)V_\lambda(-1) - 1 \right) \right] \\ &\quad + \sum_{\lambda \in \Lambda_m} n_\lambda (\sigma_\lambda^b)^2 \frac{1}{p} \left[ 1 - \frac{n-p}{n} \frac{n_\lambda}{(n_\lambda - 1)} + \frac{n}{n-p} \frac{V_\lambda(1)V_\lambda(-1)}{n_\lambda - 1} \right] . \end{aligned}$$

Thanks to Lemma 2, it comes that

$$\frac{1}{n-p} \leq \frac{1}{p} \left( \frac{n}{n-p} V_\lambda(1)V_\lambda(-1) - 1 \right) \leq \frac{p}{n-p} + \frac{n}{n-p} \frac{1}{n_\lambda} + o \left( \frac{n}{n-p} \frac{1}{n_\lambda} \right) + o \left( \frac{p}{n-p} \right) .$$

With the assumptions of Lemma 2 and assuming  $p \rightarrow +\infty$  as  $n$  tends to infinity, one gets

$$\frac{1}{n-p} \leq \frac{1}{p} \left( \frac{n}{n-p} V_\lambda(1)V_\lambda(-1) - 1 \right) \leq \frac{p}{n-p} + o \left( \frac{p}{n-p} \right) .$$

Then,

$$\begin{aligned} \frac{\sum_{\lambda \in \Lambda_m} (\sigma_\lambda^r)^2}{n-p} &\leq \sum_{\lambda \in \Lambda_m} (\sigma_\lambda^r)^2 \left[ \frac{n_\lambda}{n} + \frac{1}{p} \left( \frac{n}{n-p} V_\lambda(1)V_\lambda(-1) - 1 \right) \right] - \frac{1}{n} \sum_{i=1}^n \sigma(t_i)^2 \\ &\leq \frac{p \left( \sum_{\lambda \in \Lambda_m} (\sigma_\lambda^r)^2 + o(1) \right)}{n-p} , \end{aligned}$$

under the additional assumption **(BN)** that  $\sigma(t_i)^2$  is uniformly bounded.

Similarly, another use of Lemma 2 implies

$$\begin{aligned} \frac{n}{(n-p)n_\lambda} &\leq \frac{n}{(n-p)n_\lambda} V_\lambda(1)V_\lambda(-1) \leq \frac{n}{(n-p)n_\lambda} + \frac{np}{(n-p)n_\lambda^2} + o \left( \frac{np}{(n-p)n_\lambda^2} \right) \\ &\quad + \frac{p(p-1)}{(n-p)n_\lambda} + o \left( \frac{p(p-1)}{(n-p)n_\lambda} \right) . \end{aligned}$$

After some calculations, it yields

$$n_\lambda \left(\sigma_\lambda^b\right)^2 \frac{1}{p} \left[ 1 - \frac{n-p}{n} \frac{n_\lambda}{(n_\lambda-1)} + \frac{n}{n-p} \frac{V_\lambda(1)V_\lambda(-1)}{n_\lambda-1} \right] = n_\lambda \left(\sigma_\lambda^b\right)^2 \left[ \frac{1}{n} + o\left(\frac{1}{n}\right) \right].$$

Moreover, if one assumes **(BV)** that there exists  $M > 0$  such that for every  $i$ ,  $(s(t_i) - s_m(t_i))^2 \leq M$ , then

$$\sum_{\lambda \in \Lambda_m} n_\lambda \left(\sigma_\lambda^b\right)^2 \frac{1}{p} \left[ 1 - \frac{n-p}{n} \frac{n_\lambda}{(n_\lambda-1)} + \frac{n}{n-p} \frac{V_\lambda(1)V_\lambda(-1)}{n_\lambda-1} \right] = \|s - s_m\|_n^2 + o\left(\frac{1}{n}\right),$$

which concludes the proof.

#### 1.4 Magnitude of $V_\lambda(1)V_\lambda(-1)$

**Lemma 2.** *Let  $m \in \mathcal{M}_n$ . Assume that some  $K \in (0, 1)$  exists such that for every  $\lambda \in \Lambda_m$ ,  $n_\lambda \geq Kn$ , and  $n_\lambda > p$  with  $p^2/n \rightarrow 0$  as  $n$  tends to infinity. Then,*

$$1 \leq V_\lambda(1)V_\lambda(-1) \leq 1 + \frac{p}{n_\lambda} + o\left(\frac{p}{n_\lambda}\right) + \frac{p(p-1)}{n} + o\left(\frac{p(p-1)}{n}\right).$$

*Lemma 2.* The lower bound comes from Jensen's inequality, while the upper bound results from

$$\begin{aligned} V_\lambda(1)V_\lambda(-1) &\leq \frac{n_\lambda}{n_\lambda - p} \frac{n^p}{n(n-1)(n-(p-1))} \\ &\leq 1 + \left( \frac{p}{n_\lambda} + \frac{p(p-1)}{n} + o\left(\frac{p}{n_\lambda} + \frac{p(p-1)}{n}\right) \right). \end{aligned}$$

□

## 2 Random frameworks

In order to assess the generality of the results of Section 5.3 of the main paper, the procedures  $\llbracket \text{ERM}, \text{VF}_5 \rrbracket$ ,  $\llbracket \text{Loo}, \text{VF}_5 \rrbracket$ ,  $\llbracket \text{Lpo}_{20}, \text{VF}_5 \rrbracket$ ,  $\llbracket \text{ERM}, \text{BM} \rrbracket$ , BGH, ZS, and PML have been compared in three random settings. The following process has been repeated  $N = 10\,000$  times. First, piecewise constant functions  $s$  and  $\sigma$  are randomly chosen (see Section 2.2 for details). Then, given  $s$  and  $\sigma$ , a data sample  $(t_i, Y_i)_{1 \leq i \leq n}$  is generated as described in Section 3.3.1 of the main paper, and the same collection of models is used. Finally, each procedure  $\mathfrak{P}$  is applied to the sample  $(t_i, Y_i)_{1 \leq i \leq n}$ , and its loss  $\|s - \tilde{s}_{\mathfrak{P}}(P_n)\|_n^2$  is measured, as well as the loss of the oracle  $\inf_{m \in \mathcal{M}_n} \left\{ \|s - \hat{s}_m\|_n^2 \right\}$ .

To summarize the results, the quality of each procedure is measured by the ratio

$$C_{\text{or}}^{(R)}(\mathfrak{P}) = \frac{\mathbb{E}_{s, \sigma, \varepsilon_1, \dots, \varepsilon_n} \left[ \|s - \tilde{s}_{\mathfrak{P}}(P_n)\|_n^2 \right]}{\mathbb{E}_{s, \sigma, \varepsilon_1, \dots, \varepsilon_n} \left[ \inf_{m \in \mathcal{M}_n} \left\{ \|s - \hat{s}_m\|_n^2 \right\} \right]}.$$

The notation  $C_{\text{or}}^{(R)}(\mathfrak{P})$  differs from  $C_{\text{or}}(\mathfrak{P})$  to emphasize that each expectation includes the randomness of  $s$  and  $\sigma$ , in addition to the one of  $(\varepsilon_i)_{1 \leq i \leq n}$ .

Framework	A	B	C
[[ERM, BM]]	4.68 ± 0.03	4.63 ± 0.03	7.34 ± 0.05
[[ERM, VF <sub>5</sub> ]]	<b>4.67</b> ± 0.03	4.61 ± 0.03	6.49 ± 0.04
[[Loo, VF <sub>5</sub> ]]	<b>4.63</b> ± 0.03	<b>4.52</b> ± 0.03	6.04 ± 0.04
[[Lpo <sub>20</sub> , VF <sub>5</sub> ]]	5.40 ± 0.03	5.25 ± 0.03	6.57 ± 0.05
[[Lpo <sub>50</sub> , VF <sub>5</sub> ]]	6.23 ± 0.03	5.94 ± 0.04	7.56 ± 0.05
BGH	4.85 ± 0.03	4.78 ± 0.03	7.70 ± 0.05
ZS	5.34 ± 0.04	5.25 ± 0.04	7.28 ± 0.05
PML	5.03 ± 0.03	4.99 ± 0.03	<b>4.84</b> ± 0.04

Table 1: Performance  $C_{\text{or}}^{(R)}(\mathfrak{P})$  of several model selection procedures  $\mathfrak{P}$  in frameworks A, B, C with sample size  $n = 100$ . In each framework,  $N = 10\,000$  independent samples have been considered. Next to each value is indicated the corresponding empirical standard deviation divided by  $\sqrt{N}$ .

## 2.1 Results of the simulation experiments

The results of this experiment—which are reported in Tables 1—mostly confirm the results of Section 5.3 of the main paper (except that all the frameworks are heteroscedastic here), that is, whatever  $p$ ,  $[[\text{Lpo}_p, \text{VF}_5]]$  outperforms  $[[\text{ERM}, \text{VF}_5]]$ , which strongly outperforms  $[[\text{ERM}, \text{BM}]]$ .

Moreover, the difference between the performances of  $[[\text{Lpo}_p, \text{VF}_5]]$  and  $[[\text{ERM}, \text{VF}_5]]$  is the largest in setting C and the smallest in setting A. This fact confirms the interpretation given in Section 3 of the main paper for the failure of ERM for localizing a given number of change-points. Indeed, the main differences between frameworks A, B and C—which are precisely defined in Section 2.2—can be sketched as follows:

- A. the partitions on which  $s$  is built is often close to regular, and  $\sigma$  is chosen independently from  $s$ .
- B. the partitions on which  $s$  is built are often irregular, and  $\sigma$  is chosen independently from  $s$ .
- C. the partitions on which  $s$  is built are often irregular, and  $\sigma$  depends on  $s$ , so that the noise-level is smaller where  $s$  jumps more often.

In other words, frameworks A, B and C have been built so that for any  $D \in \mathcal{D}_n$ , the largest variations over  $\mathcal{M}_n(D)$  of  $V(m)$  (defined by Eq. (7) of the main paper) occur in framework C, and the smallest variations occur in framework A. As a consequence, variations of the performance of  $[[\text{ERM}, \text{VF}_5]]$  compared to  $[[\text{Lpo}_p, \text{VF}_5]]$  according to the framework certainly come from the local overfitting phenomenon presented in Section 3 of the main paper.

## 2.2 Detailed definition of the random frameworks

Let us now detail how piecewise constant functions  $s$  and  $\sigma$  have been generated in the frameworks A, B and C. In each framework,  $s$  and  $\sigma$  are of the form

$$s(x) = \sum_{j=0}^{K_s-1} \alpha_j \mathbb{1}_{[a_j; a_{j+1})} + \alpha_{K_s} \mathbb{1}_{[a_{K_s}; a_{K_s+1})} \quad \text{with } a_0 = 0 < a_1 < \dots < a_{K_s+1} = 1$$

$$\sigma(x) = \sum_{j=0}^{K_\sigma-1} \beta_j \mathbb{1}_{[b_j; b_{j+1})} + \beta_{K_\sigma} \mathbb{1}_{[b_{K_\sigma}; b_{K_\sigma+1})} \quad \text{with } b_0 = 0 < b_1 < \dots < b_{K_\sigma+1} = 1$$

for some positive integers  $K_s, K_\sigma$  and real numbers  $\alpha_0, \dots, \alpha_{K_s} \in \mathbb{R}$  and  $\beta_0, \dots, \beta_{K_\sigma} > 0$ .

*Remark 1.* The frameworks A, B and C depend on the sample size  $n$ , through the distribution of  $K_s$ ,  $K_\sigma$ , and of the size of the intervals  $[a_j; a_{j+1})$  and  $[b_j; b_{j+1})$ . This ensures that the signal-to-noise ratio remains rather small, so that the quadratic risk remains an adequate performance measure for change-point detection.

When the signal-to-noise ratio is larger (that is, when all jumps of  $s$  are much larger than the noise-level, and the number of jumps of  $s$  is small compared to the sample size), the change-point detection problem is of different nature. In particular, the number of change-points would be better estimated with procedures targeting identification (such as BIC, or even larger penalties) than efficiency (such as VFCV).

### 2.2.1 Framework A

In framework A,  $s$  and  $\sigma$  are generated as follows:

- $K_s$ , the number of jumps of  $s$ , has uniform distribution over  $\{3, \dots, \lfloor \sqrt{n} \rfloor\}$ .
- For  $0 \leq j \leq K_s$ ,

$$a_{j+1} - a_j = \Delta_{\min}^s + \frac{(1 - (K_s + 1)\Delta_{\min}^s)U_j}{\sum_{k=0}^{K_s} U_k}$$

with  $\Delta_{\min}^s = \min\{5/n, 1/(K_s + 1)\}$  and  $U_0, \dots, U_{K_s}$  are i.i.d. with uniform distribution over  $[0; 1]$ .

- $\alpha_0 = V_0$  and for  $1 \leq j \leq K_s$ ,  $\alpha_j = \alpha_{j-1} + V_j$  where  $V_0, \dots, V_{K_s}$  are i.i.d. with uniform distribution over  $[-1; -0.1] \cup [0.1; 1]$ .
- $K_\sigma$ , the number of jumps of  $\sigma$ , has uniform distribution in  $\{5, \dots, \lfloor \sqrt{n} \rfloor\}$ .
- For  $0 \leq j \leq K_\sigma$ ,

$$b_{j+1} - b_j = \Delta_{\min}^\sigma + \frac{(1 - (K_\sigma + 1)\Delta_{\min}^\sigma)U'_j}{\sum_{k=0}^{K_\sigma} U'_k}$$

with  $\Delta_{\min}^\sigma = \min\{5/n, 1/(K_\sigma + 1)\}$  and  $U'_0, \dots, U'_{K_\sigma}$  are i.i.d. with uniform distribution over  $[0; 1]$ .

- $\beta_0, \dots, \beta_{K_\sigma}$  are i.i.d. with uniform distribution over  $[0.05; 0.5]$ .

Two examples of a function  $s$  and a sample  $(t_i, Y_i)$  generated in framework A are plotted on Figure 1.

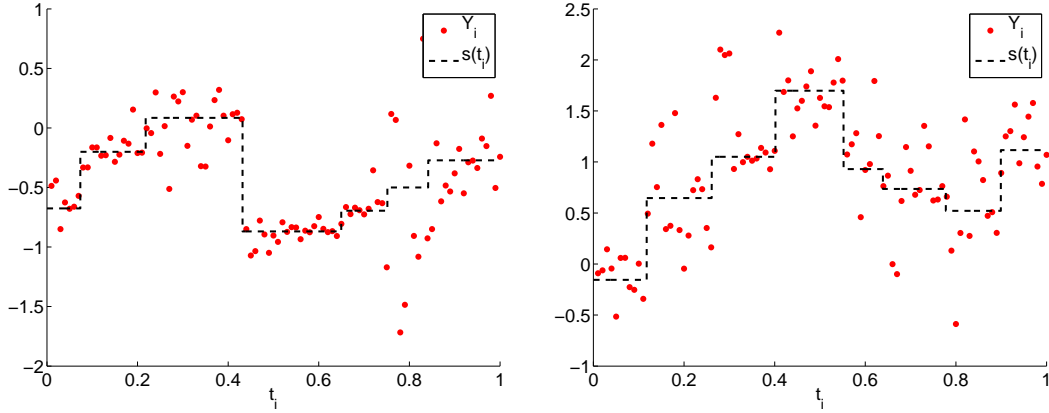


Figure 1: Random framework A: two examples of a sample  $(t_i, Y_i)_{1 \leq i \leq 100}$  and the corresponding regression function  $s$ .

### 2.2.2 Framework B

The only difference with framework A is that  $U_0, \dots, U_{K_s}$  are i.i.d. with the same distribution as  $Z = |10Z_1 + Z_2|$  where  $Z_1$  has Bernoulli distribution with parameter  $1/2$  and  $Z_2$  has a standard Gaussian distribution. Two examples of a function  $s$  and a sample  $(t_i, Y_i)$  generated in framework B are plotted on Figure 2.

### 2.2.3 Framework C

The main difference between frameworks C and B is that  $[0; 1]$  is split into two regions:  $a_{K_{s,1}+1} = 1/2$  and  $K_s = K_{s,1} + K_{s,2} + 1$  for some positive integers  $K_{s,1}, K_{s,2}$ , and the bounds of the distribution of  $\beta_j$  are larger when  $b_j \geq 1/2$  and smaller when  $b_j < 1/2$ . Two examples of a function  $s$  and a sample  $(t_i, Y_i)$  generated in framework C are plotted on Figure 3. More precisely,  $s$  and  $\sigma$  are generated as follows:

- $K_{s,1}$  has uniform distribution over  $\{2, \dots, K_{\max,1}\}$  with  $K_{\max,1} = \lfloor \sqrt{n} \rfloor - 1 - \lfloor (\lfloor \sqrt{n} \rfloor - 1) / 3 \rfloor$ .
- $K_{s,2}$  has uniform distribution over  $\{0, \dots, K_{\max,2}\}$  with  $K_{\max,2} = \lfloor (\lfloor \sqrt{n} \rfloor - 1) / 3 \rfloor$ .
- Let  $U_0, \dots, U_{K_s}$  be i.i.d. random variables with the same distribution as  $Z = |10Z_1 + Z_2|$  where  $Z_1$  has Bernoulli distribution with parameter  $1/2$  and  $Z_2$  has a standard Gaussian distribution.
- For  $0 \leq j \leq K_{s,1}$ ,

$$a_{j+1} - a_j = \Delta_{\min}^{s,1} + \frac{(1 - (K_{s,1} + 1)\Delta_{\min}^{s,1})U_j}{\sum_{k=0}^{K_{s,1}} U_k}$$

with  $\Delta_{\min}^{s,1} = \min \{5/n, 1/(K_{s,1} + 1)\}$ .

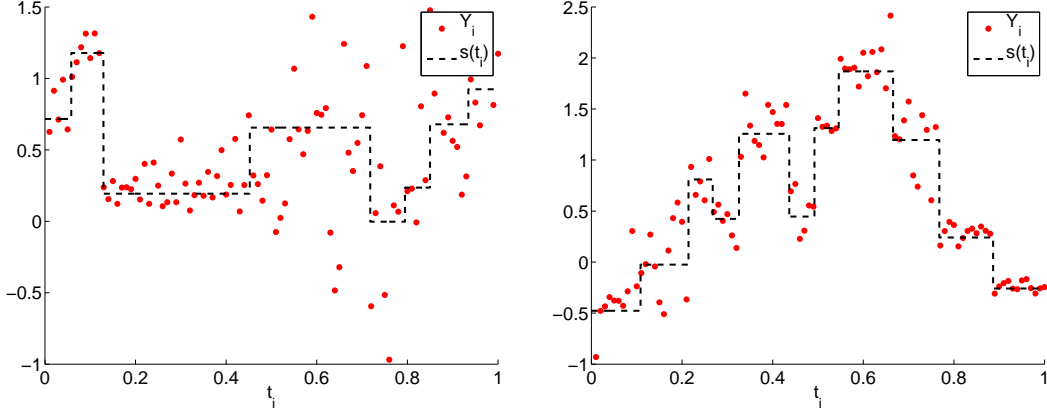


Figure 2: Random framework B: two examples of a sample  $(t_i, Y_i)_{1 \leq i \leq 100}$  and the corresponding regression function  $s$ .

- For  $K_{s,1} + 1 \leq j \leq K_s$ ,

$$a_{j+1} - a_j = \Delta_{\min}^{s,2} + \frac{(1 - (K_{s,2} + 1)\Delta_{\min}^{s,2})U_j}{\sum_{k=K_{s,1}+1}^{K_s} U_k}$$

with  $\Delta_{\min}^{s,2} = \min \{5/n, 1/(K_{s,2} + 1)\}$ .

- $\alpha_0 = V_0$  and for  $1 \leq j \leq K_s$ ,  $\alpha_j = \alpha_{j-1} + V_j$  where  $V_0, \dots, V_{K_s}$  are i.i.d. with uniform distribution over  $[-1; -0.1] \cup [0.1; 1]$ .
- $K_\sigma, (b_{j+1} - b_j)_{0 \leq j \leq K_\sigma}$  are distributed as in frameworks A and B.
- $\beta_0, \dots, \beta_{K_\sigma}$  are independent.  
When  $b_j < 1/2$ ,  $\beta_j$  has uniform distribution over  $[0.025; 0.2]$ .  
When  $b_j \geq 1/2$ ,  $\beta_j$  has uniform distribution over  $[0.1; 0.8]$ .

### 3 Restrictions on the model collection with cross-validation

The use of CV induces some restrictions on the collection of models we consider.

First, when used for choosing the best segmentation of each dimension (Section 3 of the main paper), CV estimators require that each interval of the considered segmentation contains at least two observations: one belonging to the training sample, and the other to the validation sample. This implies a first constraint of the largest dimension:  $D \leq n/2$ . Note that there is only one possible segmentation with two points in each interval (with  $n$  even).

Second, when the estimation of the number  $D - 1$  of change-points is addressed in Section 4 of the main paper,  $\text{crit}_{\text{VF}_V}(D)$  is computed by considering the models  $\tilde{S}_D((t_i)_{i \notin B_k})$  for  $k = 1, \dots, V$ , which are non-empty only when  $D \leq (n - \max_k \{\text{Card}(B_k)\})/2 \approx$

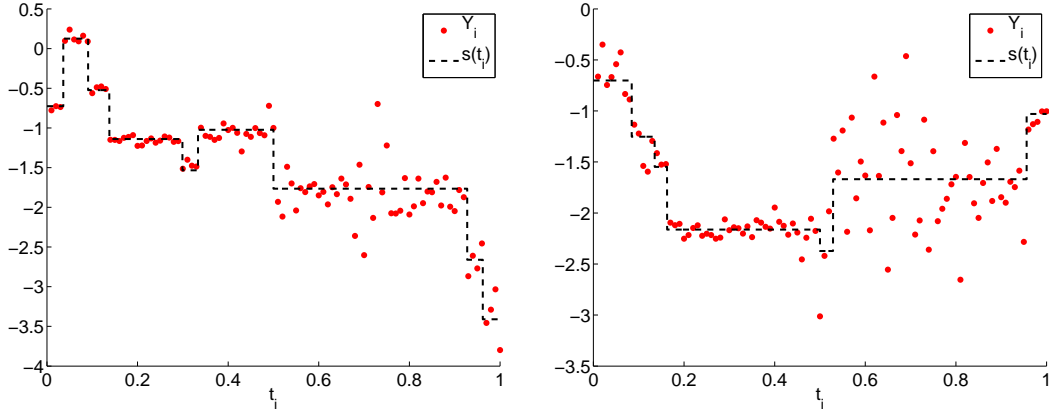


Figure 3: Random framework C: two examples of a sample  $(t_i, Y_i)_{1 \leq i \leq 100}$  and the corresponding regression function  $s$ .

$n(V-1)/(2V)$ . Indeed, at each round of the  $V$ -fold process,  $\max_k \{\text{Card}(B_k)\}$  points are removed from the whole sample. Since, moreover, Lpo is used at each round to choose the best segmentations, the largest possible dimension is  $n(V-1)/(2V)$ .

From a practical point of view, the blocks  $(B_k)_{1 \leq k \leq V}$  are always chosen such that  $\forall i, k$ ,  $\{t_i, t_{i+1}\} \cap B_k^c \neq \emptyset$ , so that  $\tilde{S}_D((t_i)_{i \notin B_k})$  is not too different from  $\tilde{S}_D((t_i)_{1 \leq i \leq n})$ . Despite this careful choice of  $(B_k)_{1 \leq k \leq V}$ ,  $\tilde{S}_D((t_i)_{i \notin B_k})$  is still much smaller than  $\tilde{S}_D((t_i)_{1 \leq i \leq n})$  when  $D$  is close to  $n(V-1)/(2V)$ . Since the curves on Figure 4 of the main paper reflect the maximum estimation error over all the segmentations in  $\tilde{S}_D$  (up to the bias), the smaller  $\tilde{S}_D((t_i)_{i \notin B_k})$ , the lower the value of the curve. This explains why  $\text{crit}_{\text{VFV}}(D)$  decreases for the largest values of  $D$  on Figure 4 of the main paper.

In practical applications, since this artefact could lead  $\text{crit}_{\text{VFV}}(D)$  to underestimate  $\|s - \hat{s}_{\tilde{m}(D)}\|_n^2$  for values of  $D$  close to  $n(V-1)/(2V)$ , we suggest to discard values of  $D > 9n(V-1)/(20V)$ . This restriction of the collection of models has been made, for all procedures, in the main paper. Since we took  $V = 5$  and  $n = 100$  in all our experiments, this leads to the choice

$$\mathcal{D}_n = \left\{ 1, \dots, \frac{9n(V-1)}{20V} \right\} = \{1, \dots, 36\}$$

instead of  $\mathcal{D}_n = \{1, \dots, 40\}$ .

## 4 Calibration of Birgé and Massart's penalization

Birgé and Massart's penalization makes use of the penalty

$$\text{pen}_{\text{BM}}(D) := \frac{\hat{C}D}{n} \left( 5 + 2 \log \left( \frac{n}{D} \right) \right) .$$

In a previous version of this work [6, Chapter 7],  $\hat{C}$  was defined as suggested in [7, 8], that is,  $\hat{C} = 2\hat{K}_{\text{max.jump}}$  with the notation below. This yielded poor performances, which



seemed related to the definition of  $\widehat{C}$ . Therefore, alternative definitions for  $\widehat{C}$  have been investigated, leading to the choice  $\widehat{C} = \widehat{\sigma}^2$  defined by (1) below. The present section intends to motivate this choice.

Two main approaches have been considered in the literature for defining  $\widehat{C}$  in the penalty  $\text{pen}_{\text{BM}}$ :

- Use  $\widehat{C} = \widehat{\sigma}^2$  any estimate of the (average) noise-level, for instance,

$$\widehat{\sigma}^2 := \frac{1}{n} \sum_{i=1}^{n/2} (Y_{2i} - Y_{2i-1})^2, \quad (1)$$

assuming  $n$  is even and  $t_1 < \dots < t_n$ .

- Use Birgé and Massart's *slope heuristics*, that is, compute the sequence

$$\widehat{D}(K) := \arg \min_{D \in \mathcal{D}_n} \left\{ P_n \gamma(\widehat{s}_{\widehat{m}_{\text{ERM}}(D)}) + \frac{KD}{n} \left( 5 + 2 \log \left( \frac{n}{D} \right) \right) \right\},$$

find the (unique)  $K = \widehat{K}_{\text{jump}}$  at which  $\widehat{D}(K)$  jumps from large to small values, and define  $\widehat{C} = 2\widehat{K}_{\text{jump}}$ .

The first approach follows from theoretical and experimental results [4, 8] which show that  $\widehat{C}$  should be close to  $\sigma^2$  when the noise-level is constant; (1) is a classical estimator of the variance used for instance by Baraud [3] for model selection in a different setting.

The optimality (in terms of oracle inequalities) of the second approach has been proved for regression with *homoscedastic Gaussian noise* and possibly exponential collections of models [5], as well as in a heteroscedastic framework with polynomial collections of models [2]. In the context of change-point detection with *homoscedastic data*, Lavielle [7] and Lebarbier [8] showed that  $\widehat{C} = 2\widehat{K}_{\text{max,jump}}$  can even perform better than  $\widehat{C} = \sigma^2$  when  $\widehat{K}_{\text{max,jump}}$  corresponds to the highest jump of  $\widehat{D}(K)$ .

Alternatively, it was proposed in [2] to define  $\widehat{C} = 2\widehat{K}_{\text{thresh}}$ , where

$$\widehat{K}_{\text{thresh}} := \min \left\{ K \text{ s.t. } \widehat{D}(K) \leq D_{\text{thresh}} := \left\lfloor \frac{n}{\ln(n)} \right\rfloor \right\}. \quad (2)$$

These three definitions of  $\widehat{C}$  have been compared with  $\widehat{C} = \sigma_{\text{true}}^2 := n^{-1} \sum_{i=1}^n \sigma(t_i)^2$  in the settings of the paper. The results are reported in Table 2. The main conclusions are the following.

- $2\widehat{K}_{\text{thresh}}$  almost always beats  $2\widehat{K}_{\text{max,jump}}$ , even in homoscedastic settings. This confirms some simulation results reported in [2].
- $\sigma_{\text{true}}^2$  almost always beats slope heuristics-based definitions of  $\widehat{C}$ , but not always, as previously noticed by Lebarbier [8]. Differences of performance can be huge (in particular when  $\sigma = \sigma_s$ ), but not always in favour of  $\sigma_{\text{true}}^2$  (for instance, when  $s = s_3$ ).
- $\widehat{\sigma}^2$  yields significantly better performance than  $\sigma_{\text{true}}^2$  in most settings (but not all), with huge margins in some heteroscedastic settings.

$s.$	$\sigma.$	$2\widehat{K}_{\max,\text{jump}}$	$2\widehat{K}_{\text{thresh.}}$	$\widehat{\sigma}^2$	$\sigma_{\text{true}}^2$
1	c	17.12 $\pm$ 0.27	5.87 $\pm$ 0.03	<b>1.70</b> $\pm$ 0.02	1.97 $\pm$ 0.02
	pc,1	134.35 $\pm$ 1.50	11.00 $\pm$ 0.03	<b>1.24</b> $\pm$ 0.02	7.24 $\pm$ 0.04
	pc,2	35.60 $\pm$ 0.33	11.20 $\pm$ 0.03	<b>3.10</b> $\pm$ 0.03	7.63 $\pm$ 0.04
	pc,3	23.60 $\pm$ 0.20	11.28 $\pm$ 0.03	<b>3.81</b> $\pm$ 0.04	7.67 $\pm$ 0.04
	s	46.36 $\pm$ 0.48	10.14 $\pm$ 0.03	2.08 $\pm$ 0.02	<b>1.63</b> $\pm$ 0.02
2	c	6.34 $\pm$ 0.04	6.97 $\pm$ 0.03	3.58 $\pm$ 0.02	<b>3.54</b> $\pm$ 0.02
	pc,1	15.01 $\pm$ 0.08	17.63 $\pm$ 0.05	<b>6.87</b> $\pm$ 0.06	11.74 $\pm$ 0.06
	pc,2	15.25 $\pm$ 0.07	17.94 $\pm$ 0.05	<b>9.25</b> $\pm$ 0.06	13.00 $\pm$ 0.06
	pc,3	15.13 $\pm$ 0.07	17.59 $\pm$ 0.05	<b>8.79</b> $\pm$ 0.06	12.66 $\pm$ 0.07
	s	8.80 $\pm$ 0.04	10.05 $\pm$ 0.03	<b>4.76</b> $\pm$ 0.03	11.00 $\pm$ 0.02
3	c	5.17 $\pm$ 0.03	4.68 $\pm$ 0.01	4.67 $\pm$ 0.01	<b>4.19</b> $\pm$ 0.01
	pc,1	7.37 $\pm$ 0.11	5.98 $\pm$ 0.02	<b>4.55</b> $\pm$ 0.02	5.09 $\pm$ 0.02
	pc,2	7.15 $\pm$ 0.03	6.83 $\pm$ 0.02	<b>5.90</b> $\pm$ 0.02	5.95 $\pm$ 0.02
	pc,3	7.16 $\pm$ 0.02	7.19 $\pm$ 0.02	<b>6.24</b> $\pm$ 0.02	6.31 $\pm$ 0.02
	s	8.81 $\pm$ 0.08	7.38 $\pm$ 0.02	<b>5.64</b> $\pm$ 0.02	15.13 $\pm$ 0.04
4	c	17.73 $\pm$ 0.28	5.92 $\pm$ 0.03	2.05 $\pm$ 0.02	<b>1.97</b> $\pm$ 0.02
	pc,2	35.50 $\pm$ 0.34	11.40 $\pm$ 0.03	<b>5.54</b> $\pm$ 0.04	7.77 $\pm$ 0.04
	A	7.13 $\pm$ 0.04	7.34 $\pm$ 0.03	<b>4.68</b> $\pm$ 0.03	4.84 $\pm$ 0.03
	B	6.84 $\pm$ 0.04	7.16 $\pm$ 0.03	<b>4.63</b> $\pm$ 0.03	4.74 $\pm$ 0.03
	C	10.43 $\pm$ 0.06	12.25 $\pm$ 0.06	<b>7.34</b> $\pm$ 0.05	8.92 $\pm$ 0.05

Table 2: Performance  $C_{\text{or}}(\text{BM})$  with four different definitions of  $\widehat{C}$  (see text), in some of the simulation settings considered in the paper. See the text for the definition of  $s_4$ . In each setting,  $N = 10\,000$  independent samples have been generated. Next to each value is indicated the corresponding empirical standard deviation divided by  $\sqrt{N}$ , measuring the uncertainty of the estimated performance.

The latter result actually comes from an artefact, which can be explained by the bias of  $\widehat{\sigma^2}$  as an estimator of  $\sigma_{\text{true}}^2$ . Indeed,

$$\mathbb{E} \left[ \widehat{\sigma^2} \right] = \frac{1}{n} \sum_{i=1}^n \sigma(t_i)^2 + \frac{1}{n} \sum_{i=1}^{n/2} (s(t_{2i}) - s(t_{2i-1}))^2 \geq \frac{1}{n} \sum_{i=1}^n \sigma(t_i)^2 = \sigma_{\text{true}}^2 . \quad (3)$$

The difference between these expectations is not negligible in all the settings of the paper. For instance, when  $n = 100$ ,  $t_i = i/n$  and  $s = s_1$ ,  $n^{-1} \sum_{i=1}^n (s(t_{2i}) - s(t_{2i-1}))^2 = 0.04$  whereas  $\sigma_{\text{true}}^2$  varies between 0.015 (when  $\sigma = \sigma_{pc,1}$ ) and 0.093 (when  $\sigma = \sigma_{pc,3}$ ). To stress the influence of the bias on the performance, we introduced  $s_4$ , which is based on  $s_1$  with its four jumps shifted by  $1/200$  to the right (that is, with jumps at 0.205, 0.405, 0.605 and 0.805). When  $s = s_4$ ,  $\widehat{\sigma^2}$  is an unbiased estimator of  $\sigma_{\text{true}}^2$ , which significantly deteriorates the performances of BM (see Table 2), while the performances of other procedures do not change.

The reason for this change of performance is that overpenalization improves the results of BM in most of the considered heteroscedastic settings, as shown by the right panel of Figure 4 of the main paper. Indeed,  $\text{pen}_{\text{BM}}$  is a poor penalty when data are heteroscedastic, underpenalizing dimensions close to the oracle but overpenalizing the largest dimensions. Then, in a setting like  $(s_2, \sigma_{pc,3})$  multiplying  $\text{pen}_{\text{BM}}$  by a factor  $C_{\text{over}} > 1$  helps decreasing the selected dimension; the same cause has opposite consequences in other settings, such as  $(s_1, \sigma_s)$  or  $(s_3, \sigma_c)$ . Nevertheless, even choosing  $\widehat{C}$  using both  $P_n$  and  $s$ ,  $(\text{crit}_{\text{BM}}(D))_{D>0}$  remains a poor estimate of  $\left( \|s - \widehat{s}_{\widehat{m}_{\text{ERM}}(D)}\|_n^2 \right)_{D>0}$  in most heteroscedastic settings (even up to an additive constant). Choosing the overpenalization factor from data is a difficult problem, especially without knowing *a priori* whether the signal is homoscedastic or heteroscedastic. This question deserves a specific extensive simulation experiment. To be completely fair with CV methods, such an experiment should also compare BM with overpenalization to  $V$ -fold penalization [1] with overpenalization, for choosing the number of change-points.

To conclude,  $\text{pen}_{\text{BM}}$  with  $\widehat{C} = \widehat{\sigma^2}$  is not a reliable change-point detection procedure, and the apparently good performances observed in Table 2 could be misleading.

Results of Table 2 for  $\widehat{C} = \sigma_{\text{true}}^2$  indicate how far the performances of BM could be improved without overpenalization. According to Tables 5 and 1, BM with  $\widehat{C} = \sigma_{\text{true}}^2$  only has significantly better performances than  $\llbracket \text{ERM}, \text{VF}_5 \rrbracket$  or  $\llbracket \text{Loo}, \text{VF}_5 \rrbracket$  in the three homoscedastic settings and in setting  $(s_1, \sigma_s)$ .

## 5 Additional results from the simulation study

In the next pages are presented extended versions of the Tables of the main paper, and two additional Figures.

## References

- [1] Sylvain Arlot.  $V$ -fold cross-validation improved:  $V$ -fold penalization, February 2008. arXiv:0802.0566v2.

- [2] Sylvain Arlot and Pascal Massart. Data-driven calibration of penalties for least-squares regression. *J. Mach. Learn. Res.*, 10:245–279 (electronic), 2009.
- [3] Yannick Baraud. Model selection for regression on a random design. *ESAIM Probab. Statist.*, 6:127–146 (electronic), 2002.
- [4] Lucien Birgé and Pascal Massart. Gaussian model selection. *J. European Math. Soc.*, 3(3):203–268, 2001.
- [5] Lucien Birgé and Pascal Massart. Minimal penalties for Gaussian model selection. *Probab. Theory Related Fields*, 138(1-2):33–73, 2007.
- [6] Alain Celisse. *Model selection via cross-validation in density estimation, regression and change-points detection*. PhD thesis, University Paris-Sud 11, December 2008. <http://tel.archives-ouvertes.fr/tel-00346320/>.
- [7] Marc Lavielle. Using penalized contrasts for the change-point problem. *Signal Proces.*, 85:1501–1510, 2005.
- [8] Émilie Lebarbier. Detecting multiple change-points in the mean of a Gaussian process by model selection. *Signal Proc.*, 85:717–736, 2005.

$s.$	$\sigma.$	ERM	Loo	Lpo <sub>20</sub>	Lpo <sub>50</sub>	PML
1	c	<b>1.56</b> $\pm$ 0.01	<b>1.56</b> $\pm$ 0.01	<b>1.56</b> $\pm$ 0.01	<b>1.55</b> $\pm$ 0.01	<b>1.57</b> $\pm$ 0.01
	pc,1	<b>1.12</b> $\pm$ 0.01	<b>1.14</b> $\pm$ 0.02	<b>1.14</b> $\pm$ 0.02	<b>1.13</b> $\pm$ 0.01	<b>1.14</b> $\pm$ 0.01
	pc,2	<b>1.82</b> $\pm$ 0.02	<b>1.81</b> $\pm$ 0.02	<b>1.82</b> $\pm$ 0.02	<b>1.80</b> $\pm$ 0.02	1.85 $\pm$ 0.02
	pc,3	<b>2.07</b> $\pm$ 0.02	<b>2.07</b> $\pm$ 0.02	<b>2.07</b> $\pm$ 0.02	<b>2.07</b> $\pm$ 0.02	2.19 $\pm$ 0.02
	s	<b>1.51</b> $\pm$ 0.02	<b>1.51</b> $\pm$ 0.02	<b>1.51</b> $\pm$ 0.02	<b>1.50</b> $\pm$ 0.02	1.56 $\pm$ 0.02
2	c	<b>2.87</b> $\pm$ 0.01	<b>2.89</b> $\pm$ 0.01	2.90 $\pm$ 0.01	2.96 $\pm$ 0.01	3.02 $\pm$ 0.01
	pc,1	1.33 $\pm$ 0.02	1.15 $\pm$ 0.02	1.14 $\pm$ 0.01	1.11 $\pm$ 0.01	<b>1.04</b> $\pm$ 0.01
	pc,2	2.91 $\pm$ 0.02	2.21 $\pm$ 0.02	2.15 $\pm$ 0.02	2.02 $\pm$ 0.02	<b>1.28</b> $\pm$ 0.01
	pc,3	3.14 $\pm$ 0.03	2.52 $\pm$ 0.02	2.47 $\pm$ 0.02	2.36 $\pm$ 0.02	<b>1.44</b> $\pm$ 0.01
	s	<b>2.98</b> $\pm$ 0.01	<b>2.98</b> $\pm$ 0.01	<b>3.00</b> $\pm$ 0.01	3.08 $\pm$ 0.01	3.17 $\pm$ 0.01
3	c	<b>3.18</b> $\pm$ 0.01	3.25 $\pm$ 0.01	3.29 $\pm$ 0.01	3.44 $\pm$ 0.01	3.71 $\pm$ 0.01
	pc,1	3.04 $\pm$ 0.02	2.70 $\pm$ 0.02	2.71 $\pm$ 0.02	2.79 $\pm$ 0.02	<b>2.29</b> $\pm$ 0.02
	pc,2	4.14 $\pm$ 0.02	3.69 $\pm$ 0.02	3.72 $\pm$ 0.02	3.86 $\pm$ 0.02	<b>3.12</b> $\pm$ 0.01
	pc,3	4.44 $\pm$ 0.02	3.98 $\pm$ 0.02	4.00 $\pm$ 0.02	4.14 $\pm$ 0.02	<b>3.24</b> $\pm$ 0.01
	s	3.92 $\pm$ 0.01	<b>3.72</b> $\pm$ 0.01	3.75 $\pm$ 0.01	3.88 $\pm$ 0.01	3.81 $\pm$ 0.01
A	<b>3.28</b> $\pm$ 0.02	<b>3.25</b> $\pm$ 0.02	<b>3.28</b> $\pm$ 0.02	3.39 $\pm$ 0.02	3.42 $\pm$ 0.02	
B	3.21 $\pm$ 0.02	<b>3.15</b> $\pm$ 0.02	<b>3.18</b> $\pm$ 0.02	3.28 $\pm$ 0.02	3.38 $\pm$ 0.02	
C	4.27 $\pm$ 0.03	3.79 $\pm$ 0.03	3.79 $\pm$ 0.03	3.78 $\pm$ 0.03	<b>3.14</b> $\pm$ 0.02	

Table 3: Average performance  $C_{\text{or}}(\mathfrak{P}, \text{Id})$  for change-point detection procedures  $\mathfrak{P}$  among ERM, Loo and Lpo <sub>$p$</sub>  with  $p = 20$  and  $p = 50$ . Several regression functions  $s$  and noise-level functions  $\sigma$  have been considered, each time with  $N = 10\,000$  independent samples. Next to each value is indicated the corresponding empirical standard deviation divided by  $\sqrt{N}$ , measuring the uncertainty of the estimated performance.

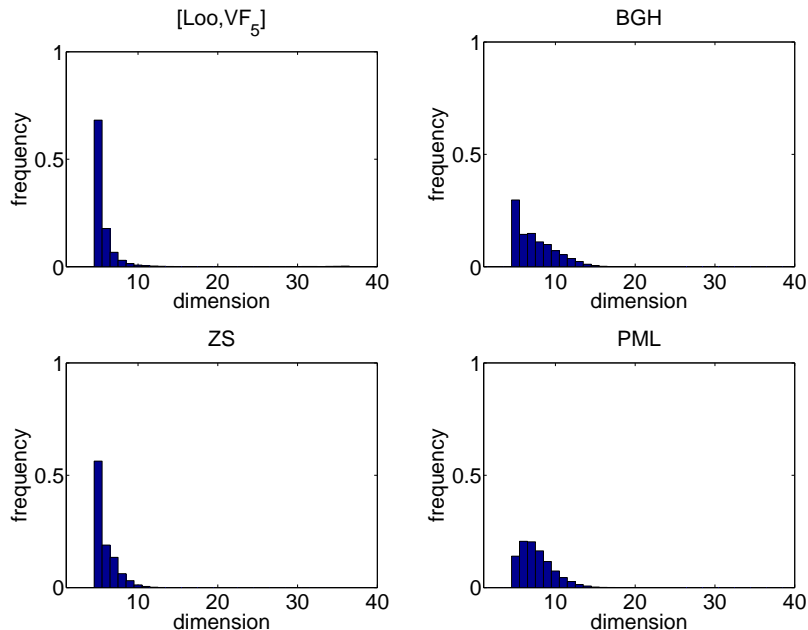


Figure 4: Empirical distribution of the selected dimension  $\hat{D}$  over  $N = 10\,000$  independent samples, for 4 procedures, in setting  $(s_1, \sigma_{pc,3})$ . The true model has dimension  $D = 5$ .

$s.$	$\sigma.$	Id	VF <sub>5</sub>	BM	BGH
1	c	1.56 ± 0.01	2.38 ± 0.02	<b>1.70</b> ± 0.02	1.85 ± 0.02
	pc,1	1.12 ± 0.01	2.67 ± 0.03	<b>1.24</b> ± 0.02	5.83 ± 0.05
	pc,2	1.82 ± 0.02	3.20 ± 0.03	<b>3.10</b> ± 0.03	6.38 ± 0.05
	pc,3	2.07 ± 0.02	<b>3.42</b> ± 0.03	3.81 ± 0.04	6.51 ± 0.04
	s	1.51 ± 0.02	2.73 ± 0.03	<b>2.08</b> ± 0.02	3.83 ± 0.03
2	c	2.87 ± 0.01	3.99 ± 0.02	3.58 ± 0.02	<b>3.52</b> ± 0.02
	pc,1	1.33 ± 0.02	<b>3.64</b> ± 0.05	6.87 ± 0.06	9.50 ± 0.07
	pc,2	2.91 ± 0.02	<b>5.62</b> ± 0.05	9.25 ± 0.06	10.13 ± 0.07
	pc,3	3.14 ± 0.03	<b>5.94</b> ± 0.06	8.79 ± 0.06	9.77 ± 0.07
	s	2.98 ± 0.01	<b>4.34</b> ± 0.03	4.76 ± 0.03	4.88 ± 0.03
3	c	3.18 ± 0.01	<b>4.31</b> ± 0.02	4.67 ± 0.01	4.47 ± 0.01
	pc,1	3.04 ± 0.02	4.65 ± 0.02	<b>4.55</b> ± 0.02	4.92 ± 0.02
	pc,2	4.14 ± 0.02	<b>5.82</b> ± 0.02	5.90 ± 0.02	5.93 ± 0.02
	pc,3	4.44 ± 0.02	<b>6.13</b> ± 0.02	6.24 ± 0.02	6.31 ± 0.02
	s	3.92 ± 0.01	<b>5.61</b> ± 0.02	<b>5.64</b> ± 0.02	<b>5.63</b> ± 0.02
A		3.28 ± 0.02	<b>4.67</b> ± 0.03	<b>4.68</b> ± 0.03	4.85 ± 0.03
B		3.21 ± 0.02	<b>4.61</b> ± 0.03	<b>4.63</b> ± 0.03	4.78 ± 0.03
C		4.27 ± 0.03	<b>6.49</b> ± 0.04	7.34 ± 0.05	7.70 ± 0.05

Table 4: Performance  $C_{\text{or}}(\llbracket \text{ERM}, \mathfrak{P} \rrbracket)$  for  $\mathfrak{P} = \text{Id}$  (that is, choosing the dimension  $D^* := \arg \min_{D \in \mathcal{D}_n} \left\{ \|s - \widehat{s}_{\widehat{m}_{\text{ERM}}(D)}\|_n^2 \right\}$ ),  $\mathfrak{P} = \text{VF}_V$  with  $V = 5$ ,  $\mathfrak{P} = \text{BM}$  or  $\mathfrak{P} = \text{BGH}$ . Several regression functions  $s$  and noise-level functions  $\sigma$  have been considered, each time with  $N = 10000$  independent samples. Next to each value is indicated the corresponding empirical standard deviation divided by  $\sqrt{N}$ , measuring the uncertainty of the estimated performance.

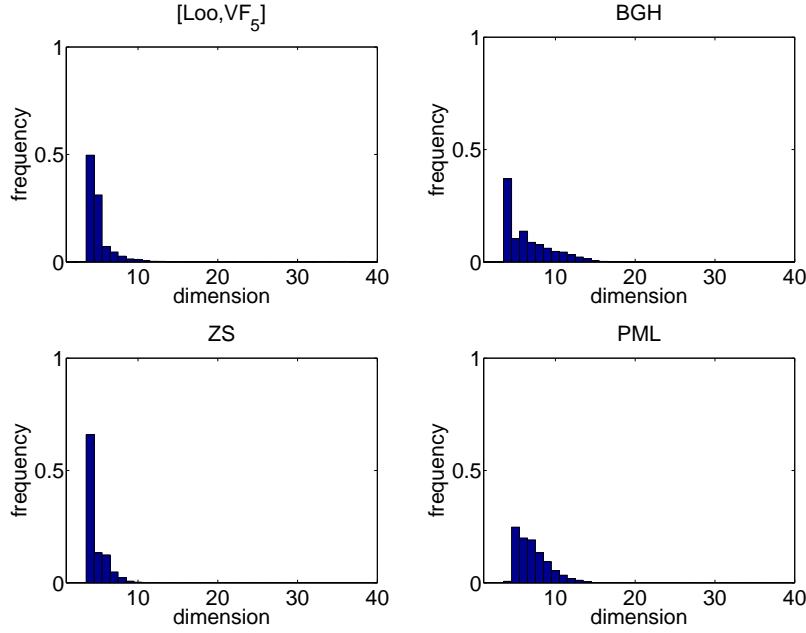


Figure 5: Same as Figure 4 in setting  $(s_2, \sigma_{pc,2})$ . The true model has dimension  $D = 5$ .

$(s, \sigma)$	$(s_1, \sigma_c)$	$(s_1, \sigma_{pc,1})$	$(s_1, \sigma_{pc,2})$	$(s_1, \sigma_{pc,3})$	$(s_1, \sigma_s)$
[[Loo, VF <sub>5</sub> ]]	2.40 ± 0.02	2.64 ± 0.03	3.17 ± 0.03	<b>3.40</b> ± 0.03	2.59 ± 0.03
[[Lpo <sub>20</sub> , VF <sub>5</sub> ]]	2.58 ± 0.02	3.78 ± 0.04	3.29 ± 0.03	3.55 ± 0.04	2.86 ± 0.03
[[Lpo <sub>50</sub> , VF <sub>5</sub> ]]	2.69 ± 0.03	3.97 ± 0.04	3.60 ± 0.04	3.77 ± 0.04	3.15 ± 0.04
[[ERM, VF <sub>5</sub> ]]	2.38 ± 0.02	2.67 ± 0.03	3.20 ± 0.03	<b>3.42</b> ± 0.03	2.73 ± 0.03
[[ERM, BM]]	<b>1.70</b> ± 0.02	<b>1.24</b> ± 0.02	<b>3.10</b> ± 0.03	3.81 ± 0.04	<b>2.08</b> ± 0.02
[[ERM, BM <sub>thr</sub> ]]	5.87 ± 0.03	11.00 ± 0.03	11.20 ± 0.03	11.28 ± 0.03	10.14 ± 0.03
BGH	1.85 ± 0.02	5.83 ± 0.05	6.38 ± 0.05	6.51 ± 0.04	3.83 ± 0.03
ZS	<b>1.71</b> ± 0.02	2.20 ± 0.03	3.92 ± 0.04	4.41 ± 0.04	2.46 ± 0.03
PML	2.79 ± 0.02	2.48 ± 0.03	3.42 ± 0.03	3.97 ± 0.04	2.72 ± 0.03

---

$(s, \sigma)$	$(s_2, \sigma_c)$	$(s_2, \sigma_{pc,1})$	$(s_2, \sigma_{pc,2})$	$(s_2, \sigma_{pc,3})$	$(s_2, \sigma_s)$
[[Loo, VF <sub>5</sub> ]]	4.02 ± 0.02	3.14 ± 0.04	4.95 ± 0.05	5.24 ± 0.05	4.32 ± 0.03
[[Lpo <sub>20</sub> , VF <sub>5</sub> ]]	4.10 ± 0.02	3.34 ± 0.05	5.00 ± 0.05	5.26 ± 0.05	4.49 ± 0.03
[[Lpo <sub>50</sub> , VF <sub>5</sub> ]]	4.44 ± 0.03	4.41 ± 0.06	5.20 ± 0.05	5.54 ± 0.06	4.86 ± 0.03
[[ERM, VF <sub>5</sub> ]]	3.99 ± 0.02	3.64 ± 0.05	5.62 ± 0.05	5.94 ± 0.06	4.34 ± 0.03
[[ERM, BM]]	3.58 ± 0.02	6.87 ± 0.06	9.25 ± 0.06	8.79 ± 0.06	4.76 ± 0.03
[[ERM, BM <sub>thr</sub> ]]	6.97 ± 0.03	17.63 ± 0.05	17.94 ± 0.05	17.59 ± 0.05	10.05 ± 0.03
BGH	<b>3.52</b> ± 0.02	9.50 ± 0.07	10.13 ± 0.07	9.77 ± 0.07	4.88 ± 0.03
ZS	3.62 ± 0.02	3.36 ± 0.05	6.50 ± 0.05	6.92 ± 0.06	<b>4.16</b> ± 0.02
PML	4.34 ± 0.02	<b>2.35</b> ± 0.03	<b>2.73</b> ± 0.03	<b>2.84</b> ± 0.03	4.32 ± 0.02

---

$(s, \sigma)$	$(s_3, \sigma_c)$	$(s_3, \sigma_{pc,1})$	$(s_3, \sigma_{pc,2})$	$(s_3, \sigma_{pc,3})$	$(s_3, \sigma_s)$
[[Loo, VF <sub>5</sub> ]]	4.42 ± 0.02	4.22 ± 0.02	5.24 ± 0.02	5.59 ± 0.02	<b>5.35</b> ± 0.02
[[Lpo <sub>20</sub> , VF <sub>5</sub> ]]	4.61 ± 0.02	4.62 ± 0.02	5.47 ± 0.02	5.78 ± 0.02	5.55 ± 0.02
[[Lpo <sub>50</sub> , VF <sub>5</sub> ]]	5.06 ± 0.02	4.84 ± 0.02	5.86 ± 0.02	6.25 ± 0.02	6.03 ± 0.02
[[ERM, VF <sub>5</sub> ]]	<b>4.31</b> ± 0.02	4.65 ± 0.02	5.82 ± 0.02	6.13 ± 0.02	5.61 ± 0.02
[[ERM, BM]]	4.67 ± 0.01	4.55 ± 0.02	5.90 ± 0.02	6.24 ± 0.02	5.64 ± 0.02
[[ERM, BM <sub>thr</sub> ]]	4.68 ± 0.01	5.98 ± 0.02	6.83 ± 0.02	7.19 ± 0.02	7.38 ± 0.02
BGH	4.47 ± 0.01	4.92 ± 0.02	5.93 ± 0.02	6.31 ± 0.02	5.63 ± 0.02
ZS	5.46 ± 0.02	4.44 ± 0.02	6.63 ± 0.02	6.61 ± 0.02	6.31 ± 0.02
PML	5.05 ± 0.02	<b>3.70</b> ± 0.03	<b>4.67</b> ± 0.03	<b>4.99</b> ± 0.03	5.52 ± 0.02

Table 5: Performance  $C_{\text{or}}(\mathfrak{P})$  for several change-point detection procedures  $\mathfrak{P}$ . ‘BM<sub>thr</sub>’ refers to BM with  $\widehat{C} = 2\widehat{K}_{\text{thresh}}$ . (see Section 4), whereas ‘BM’ correspond to  $\widehat{C} = 2\widehat{\sigma}^2$  as everywhere in the main paper. Several regression functions  $s$  and noise-level functions  $\sigma$  have been considered, each time with  $N = 10000$  independent samples. Next to each value is indicated the corresponding empirical standard deviation.

$(s, \sigma)$	$(s_4, \sigma_c)$	$(s_4, \sigma_{pc,2})$
[[Loo, VF <sub>5</sub> ]]	2.36 ± 0.02	<b>3.17</b> ± 0.03
[[Lpo <sub>20</sub> , VF <sub>5</sub> ]]	2.53 ± 0.02	3.32 ± 0.03
[[Lpo <sub>50</sub> , VF <sub>5</sub> ]]	2.64 ± 0.03	3.62 ± 0.04
[[ERM, VF <sub>5</sub> ]]	2.38 ± 0.02	3.24 ± 0.03
[[ERM, BM]]	2.05 ± 0.02	5.54 ± 0.04
[[ERM, BM <sub>thr</sub> ]]	5.92 ± 0.03	11.40 ± 0.03
BGH	1.85 ± 0.02	6.54 ± 0.05
ZS	<b>1.70</b> ± 0.02	4.00 ± 0.04
PML	2.77 ± 0.02	3.55 ± 0.04

Table 6: Same as Table 5 with  $s = s_4$ .

$(s, \sigma)$	$(s_2, \sigma_{pc,2})$	$(s_2, \sigma_{pc,3})$	$(s_3, \sigma_{pc,2})$	$(s_3, \sigma_{pc,3})$
[[Loo, VF <sub>5</sub> ]]	<b>4.47</b> ± 0.05	<b>4.69</b> ± 0.06	<b>4.80</b> ± 0.03	<b>5.11</b> ± 0.03
[[Lpo <sub>20</sub> , VF <sub>5</sub> ]]	4.62 ± 0.06	4.88 ± 0.06	5.03 ± 0.03	5.29 ± 0.03
[[Lpo <sub>50</sub> , VF <sub>5</sub> ]]	5.22 ± 0.07	5.54 ± 0.07	5.45 ± 0.03	5.83 ± 0.03
[[ERM, VF <sub>5</sub> ]]	5.98 ± 0.07	6.31 ± 0.07	5.82 ± 0.04	6.22 ± 0.04
[[ERM, BM]]	10.81 ± 0.09	10.31 ± 0.09	6.09 ± 0.04	6.45 ± 0.04
[[ERM, BM <sub>thr</sub> ]]	17.83 ± 0.09	17.26 ± 0.08	6.60 ± 0.03	6.99 ± 0.03
BGH	11.67 ± 0.09	11.15 ± 0.09	5.94 ± 0.03	6.42 ± 0.04
ZS	9.34 ± 0.09	9.29 ± 0.09	6.60 ± 0.04	6.83 ± 0.04
PML	5.04 ± 0.06	5.00 ± 0.06	5.17 ± 0.03	5.40 ± 0.03

Table 7: Same as Table 5 with  $\varepsilon_i$  having an exponential distribution.