



Multiple Change-Point Estimation With a Total Variation Penalty

Z. Harchaoui & C. Lévy-Leduc

To cite this article: Z. Harchaoui & C. Lévy-Leduc (2010) Multiple Change-Point Estimation With a Total Variation Penalty, Journal of the American Statistical Association, 105:492, 1480-1493, DOI: [10.1198/jasa.2010.tm09181](https://doi.org/10.1198/jasa.2010.tm09181)

To link to this article: <http://dx.doi.org/10.1198/jasa.2010.tm09181>



Published online: 01 Jan 2012.



Submit your article to this journal [↗](#)



Article views: 548



Citing articles: 7 View citing articles [↗](#)

Multiple Change-Point Estimation With a Total Variation Penalty

Z. HARCHAOU and C. LÉVY-LEDUC

We propose a new approach for dealing with the estimation of the location of change-points in one-dimensional piecewise constant signals observed in white noise. Our approach consists in reframing this task in a variable selection context. We use a penalized least-square criterion with a ℓ_1 -type penalty for this purpose. We explain how to implement this method in practice by using the LARS/LASSO algorithm. We then prove that, in an appropriate asymptotic framework, this method provides consistent estimators of the change points with an almost optimal rate. We finally provide an improved practical version of this method by combining it with a reduced version of the dynamic programming algorithm and we successfully compare it with classical methods.

KEY WORDS: Change-point estimation; LARS; LASSO; ℓ_1 -type penalty; Sparsity.

1. INTRODUCTION

Retrospective multiple change-point estimation consists in partitioning a nonstationary series of observations into several contiguous stationary segments of variable durations; see Brodsky and Darkhovsky (1993, 2000). It is particularly appropriate for analyzing *a posteriori* time series in which the quantity driving the behavior of the time series jumps from one level to another different level at random instants called change points. Such a task, also known as temporal signal segmentation in signal processing, arises in many applications, ranging from EEG to speech processing and network intrusion detection (Basseville and Nikiforov 1993; Ruanaidh and Fitzgerald 1996; Lévy-Leduc and Roueff 2009).

As argued by both Carlstein, Müller, and Siegmund (1994) and Brodsky and Darkhovsky (2000), in most cases detecting changes of a time-evolving statistical quantity may be reduced to the detection of changes in the mean of a new sequence derived from the initial one. Thus, we are interested in the estimation of the change-point locations t_k^* in the following model:

$$Y_t = \mu_k^* + \varepsilon_t, \\ t_{k-1}^* \leq t \leq t_k^* - 1, k = 1, \dots, K^* + 1, t = 1, \dots, n, \quad (1)$$

with the convention $t_0^* = 1$ and $t_{K^*+1}^* = n + 1$ and where the $\{\varepsilon_t\}_{0 \leq t \leq n}$ are iid zero-mean random variables, having a sub-Gaussian distribution.

This problem has recently received much attention on the theoretical side, both in a nonasymptotic and in an asymptotic setting by Massart (2005) and Yao and Au (1989), Lavielle and Moulines (2000), Boysen et al. (2009), respectively. From

a practical point of view, the standard approach for estimating the change-point locations is based on least-square fitting, performed via a dynamic programming algorithm (DP), coupled with an informational criterion such as the Schwarz criterion (Yao and Au 1989) for choosing the unknown number of change points. Indeed, for a given number of change points K , the dynamic programming algorithm, proposed by Fisher (1958) and Bellman (1961), takes advantage of the intrinsic additive nature of the least-square objective to recursively compute the optimal change-points locations with a complexity of $O(Kn^2)$ in time. Then selecting the number of change points is usually performed thanks to a Schwarz-like penalty $\lambda_n K$, where λ_n is often calibrated on data (Lavielle and Moulines 2000; Lavielle 2005), or a penalty $K(a + b \log(n/K))$ as in Lebarbier (2005), Massart (2005), where a and b are data-driven as well. We should also mention that an abundant literature tackles both change-point estimation and model selection issues from a Bayesian point of view; see Ruanaidh and Fitzgerald (1996), Fearnhead (2006), and references therein; we shall not adopt such a point of view in this work.

While optimal from a maximum likelihood point of view in the case of Gaussian noise, the application of the standard least-square approach, called LS in the remainder, is seriously harmed by a quadratic time complexity in the total duration of the series of observations in its exact implementation. Yet approximate dynamic programming procedures were devised in other contexts, such as for Dynamic Time Warping or the Viterbi algorithm (Kolesnikov and Fraenti 2003; Gales and Young 2008). Moreover, as pointed in Hawkins (2001), a computationally efficient dynamic programming algorithm for change-point estimation may be devised when a prior assumption of order structure between the segments is satisfied and therefore consists in restricting the change-point locations search to a prespecified set. Yet, designing a computationally efficient dynamic programming algorithm for change-point estimation under general assumptions is still an open problem.

Therefore, an alternative formulation might be profitable from a computational point of view, while keeping comparable performance when compared to the least-square method. A natural way to lower the time complexity of a ℓ_0 -penalized

Zaid Harchaoui is an INRIA (Institut National de Recherche en Informatique et en Automatique) research scientist (Chargé de Recherche) in the LEAR project team at Laboratoire Jean Kuntzmann, 655 avenue de l'Europe, 38334 Saint Ismier Cedex, France (E-mail: zaid.harchaoui@inria.fr). Céline Lévy-Leduc is a CNRS (Centre National de la Recherche Scientifique) research scientist (Chargée de Recherche) at Laboratoire Traitement et Communication de l'Information, which is a joint lab with TELECOM ParisTech; in the Signal and Image Processing (TSI) department, 46 Rue Barrault, 75634 Paris Cedex 13, France (E-mail: celine.levy-leduc@telecom-paristech.fr). We would like to thank Éric Moulines for drawing our attention to this problem and for helpful discussions. We would also like to thank Olivier Cappé and Jean-Philippe Vert for fruitful discussions related to this work. A large portion of this work was performed while Zaid Harchaoui was at LTCL. This work was supported by the Agence Nationale de la Recherche under grant ANR-06-BLAN-0078 KERNISIG.

least-square problem is to relax the ℓ_0 -penalty to an ℓ_1 -penalty. This strategy has proved to be appropriate in other statistical problems such as sparse PCA, sparse LDA; see d'Aspremont, Bach, and El Ghaoui (2008) and Moghaddam, Weiss, and Avidan (2006). Hence, it boils down to estimating the change-point locations by solving

$$\underset{\mathbf{u} \in \mathbb{R}^n}{\text{Minimize}} \frac{1}{n} \sum_{i=1}^n (Y_i - u_i)^2 + \lambda_n \sum_{i=1}^{n-1} |u_{i+1} - u_i|, \quad (2)$$

and recovering the change-point locations from the jumps in the $\{\hat{u}_i\}_{i=1, \dots, n}$ minimizing the criterion in Equation (2). This alternative formulation yields a subquadratic time complexity in the length of the sequence of observations, and still remains asymptotically consistent in terms of change-point estimation. Note that Tibshirani and Wang (2008) introduced the “fused lasso,” which corresponds to a two-step procedure where the first step is a least-square change-point estimation with a total-variation penalty and the second is a thresholding one to discard small jumps from the zero mean, a method specifically designed for spatial smoothing and hot spot detection in CGH data.

This article is organized as follows. In Section 2, we describe how Equation (2) is related to the well-known Least Absolute Shrinkage eStimatOr (LASSO) in least-square regression of Tibshirani (1996), usually used for efficient variable selection. We show that it turns out to be also useful for change-point estimation as well when used with a particular design matrix. We take advantage of this relationship to devise a subquadratic change-point estimation algorithm, called LS-TV for Least-Square with Total Variation penalty. In Section 3, we give theoretical results concerning the estimation of the underlying piecewise constant function and the estimation of the change-point locations. More precisely, we provide rates of convergence for the underlying piecewise constant function and for the change-point instants and we show that we can attain almost optimal rates of convergence in both cases. In Section 4, we run numerical experiments to assess the empirical behavior of LS-TV, and propose an enhanced version LS-TV* with better empirical performance.

2. METHODOLOGY

In this section, we describe the least-square change-point estimation with a total variation penalty LS-TV. In Section 2.1, we show how to recast the multiple change-point estimation problem into a particular variable selection problem. Then in Section 2.2, we describe a LAR-based implementation of LS-TV, and derive its time complexity. The theoretical properties of LS-TV are given in Section 3.

2.1 From Change-Point Estimation to Variable Selection

The multiple change-point estimation problem may be relaxed into a LASSO-type problem using appropriate auxiliary variables.

Recall the multiple change-point model (Yao and Au 1989):

$$Y_t = u_t^* + \varepsilon_t, \quad t = 1, \dots, n, \quad (3)$$

where $u_t^* = \mu_k^*$ for $t_{k-1}^* \leq t \leq t_k^* - 1$, $k = 1, \dots, K^* + 1$. We shall always assume in the remainder of this section that the true number of change points K^* is known. The issue of dealing with

an unknown number of change points will be addressed later in Sections 3 and 4.

The least-square estimation method LS, which may also be viewed as the maximum-likelihood approach in the case of Gaussian white noise, solves the following problem:

$$\begin{aligned} &\underset{\mathbf{u} \in \mathbb{R}^n}{\text{Minimize}} \frac{1}{n} \sum_{i=1}^n (Y_i - u_i)^2 \\ &\text{subject to} \quad \sum_{i=1}^{n-1} \mathbf{1}\{u_{i+1} - u_i\} = K^*. \end{aligned} \quad (4)$$

We propose here to relax the above ℓ_0 constraint into an ℓ_1 constraint on the magnitude of the jumps as follows:

$$\begin{aligned} &\underset{\mathbf{u} \in \mathbb{R}^n}{\text{Minimize}} \frac{1}{n} \sum_{i=1}^n (Y_i - u_i)^2 \\ &\text{subject to} \quad \sum_{i=1}^{n-1} |u_{i+1} - u_i| \leq K^* J_{\max}^*, \end{aligned} \quad (5)$$

where $J_{\max}^* = \max_{1 \leq k \leq K^*} |u_{k+1}^* - u_k^*|$. This alternative setting was previously elusively mentioned several times, in, for example, Mammen and Van De Geer (1997) and Boysen et al. (2009). In order to further understand the behavior of the solution $\hat{\mathbf{u}} = (\hat{u}_1, \dots, \hat{u}_n)$ of this criterion, let us denote by \mathbf{X}_n the $n \times n$ lower triangular matrix with nonzero elements equal to one.

Then, by straightforward algebra, the problem in Equation (5) may be rewritten as

$$\begin{aligned} &\underset{\boldsymbol{\beta} \in \mathbb{R}^n}{\text{Minimize}} \frac{1}{n} \sum_{i=1}^n (Y_i - (\mathbf{X}_n \boldsymbol{\beta})_i)^2 \\ &\text{subject to} \quad \sum_{i=1}^n |\beta_i| \leq K^* J_{\max}^*. \end{aligned} \quad (6)$$

The underpinning insight is the sparsity-enforcing property of the ℓ_1 -constraint, which is expected to give a sparse vector $\hat{\boldsymbol{\beta}}^n$, whose nonzero components would match with change-points locations.

A major feature of Equation (6) is that it exactly corresponds to the well-known Least Absolute Shrinkage eStimatOr (LASSO) in least-square regression of Tibshirani (1996), used for efficient variable selection. However, as far as we know, neither thorough practical implementation nor theoretical grounding has been given so far to support such an approach for change-point estimation. Actually, the corresponding minimization can be solved by using the LAR/LASSO algorithm described in Efron, Hastie, and Tibshirani (2004) and Hesterberg et al. (2008).

2.2 Implementation With Least-Angle Regression

In this section, we detail the process of the Least-Angle Regression (LAR) algorithm of Efron, Hastie, and Tibshirani (2004). For the sake of generality, we shall describe here this algorithm when we look for K_{\max} change points, K_{\max} being a known upper bound on the true number of change points. When implemented with care, we get a time complexity in

$O(K_{\max}n \log(n))$ of the LAR/LASSO algorithm in the particular case of our model. This substantial reduction of the computational complexity has to be contrasted with the complexity $O(K_{\max}n^2)$ of DP. We use in this section standard notation given for instance in [Cormen et al. \(2001\)](#).

The process is described in Table 1, the different notations involved being explained in the following in the description of each step of the algorithm. It essentially involves four steps, each of them being solved in subquadratic time complexity with respect to the number of observations n . Suppose we have performed $k - 1$ iterations in the main loop of the algorithm, then the current set of estimated change points, that is, the *active set* in the variable selection framework, is $\hat{\mathcal{T}}_{n,k-1} = \{\hat{t}_1, \dots, \hat{t}_{k-1}\}$ and the current set of estimated segment levels is $\{\hat{u}_1(k-1), \dots, \hat{u}_n(k-1)\}$. We are now describing the computational requirements of the k th iteration of the algorithm.

First, we look for the next change point \hat{t}_k to add to $\hat{\mathcal{T}}_{n,k-1}$ yielding the largest discrepancy with the true signal. This requires, given $\{\hat{u}_1(k-1), \dots, \hat{u}_n(k-1)\}$, the computation of the n cumulative sums $\{\sum_{i=j}^n \hat{u}_i(k-1)\}_{j=1, \dots, n}$. These cumulative sums may actually be computed in $O(n)$ operations in time, using the simple recursion $\sum_{i=j}^n \hat{u}_i(k-1) = \sum_{i=j+1}^n \hat{u}_i(k-1) + \hat{u}_j(k-1)$. Besides, to be included in the current set of change-point estimates (“active set”), we need to locate the new change-point estimate with regard to the other change-point estimates, which is formally equivalent to sort the set of observations.

Table 1. Description of the adaptation of LAR/LASSO algorithm for solving the LS-TV problem

<i>LS-TV with LAR/LASSO</i>	
Initialization, $k = 0$.	
(a) Set $\hat{\mathcal{T}}_{n,0} = \emptyset$.	
(b) Set $\hat{u}_i(0) = 0$, for all $i = 1, \dots, n$.	
While $k < K_{\max}$.	
(a) <i>Change-point addition</i> :	
Find \hat{t}_k such that	
$\hat{t}_k = \underset{t \in \{1, \dots, n\} \setminus \hat{\mathcal{T}}_{n,k-1}}{\text{Arg max}} \left \sum_{i=t}^n Y_i - \sum_{i=t}^n \hat{u}_i(k-1) \right .$	
(b) <i>Descent direction computation</i> :	
Compute	
$\mathbf{w}_k = (\mathbf{X}_k^T \mathbf{X}_k)^{-1} \mathbf{1}_k.$	
(c) <i>Descent step search</i> :	
Search for $\hat{\gamma}$ such that	
$\hat{\gamma} = \min_{t \in \{1, \dots, n\} \setminus \hat{\mathcal{T}}_{n,k}} \left(\frac{\sum_{i=t}^n Y_i - \sum_{i=t}^n \hat{u}_i(k)}{1 - \sum_{i=t}^n w_{k,i}}, \frac{\sum_{i=t}^n Y_i + \sum_{i=t}^n \hat{u}_i(k)}{1 + \sum_{i=t}^n w_{k,i}} \right).$	
(d) <i>Zero-crossing check</i> :	
If	
$\hat{\gamma} > \hat{\gamma}^{\text{def}} \equiv \min_j (\alpha_j w_{k,j})^{-1} \left(\sum_{i=j}^n \hat{u}_i(k) \right),$	
then, decrease $\hat{\gamma}$ down to $\hat{\gamma} = \hat{\gamma}^{\text{def}}$, and remove \tilde{t} from $\hat{\mathcal{T}}_{n,k}$, where	
$\tilde{t} \stackrel{\text{def}}{=} \underset{j}{\text{Arg min}} (\alpha_j w_{k,j})^{-1} \left(\sum_{i=j}^n \hat{u}_i(k) \right).$	

Therefore, the “change-point addition” step in Table 1 has a $O(n + n \log(n))$ time complexity.

Second, we have to compute the descent direction, which involves the multiplication of the inverse of a $(k \times k)$ -matrix by a k -long vector. Indeed, \mathbf{X}_k is a matrix which consists of the columns of \mathbf{X} indexed by the elements of $\hat{\mathcal{T}}_{n,k}$ and $\mathbf{1}_k$ denotes a vector of dimension k with each component equal to one. Given the current set of change points $\hat{\mathcal{T}}_{n,k}$, the inverse may be computed in $O(k^2)$ operations, since the entries of the inverse matrix of size $(k \times k)$ are available in close form beforehand; see (A.2) in the [Appendix](#). Then, the multiplication of the $(k \times k)$ -inverse by $\mathbf{1}_k$ is computed in $O(k^2)$ operations. If $k < K_{\max}$, then the time complexity of “descent direction computation” step is upper bounded by $O(K_{\max}^2)$.

Third, we search for the descent step. For similar reasons as for the first step, the “descent step search” step may be performed in linear time $O(n)$ time complexity. Indeed, again, this step involves the computation of n cumulative sums, which may be computed recursively.

Fourth, we check the zero crossing of the coefficients to exactly track the regularization path of the LASSO. In this step, $\alpha_j = \text{sign}(\hat{u}_{j+1}(k) - \hat{u}_j(k))$. Again, all computations involved in this step hinge on cumulative sums as previously in the first step, and therefore may be performed in $O(n)$ time complexity. Note that the maximum number of iterations N needed in practice to decrease $\hat{\gamma}$ to a small enough value to satisfy $\hat{\gamma} = \hat{\gamma}^{\text{def}}$ is unknown in general, and no theoretically grounded upper bound on N was provided in the literature so far. In practice, we set $N < K_{\max}$ in our implementation, and we never encountered any numerical issue which demanded a different (larger) setting of N . Hence, the “zero-crossing” step has at most $O(K_{\max}n)$ time complexity.

Thus, the implementation of LS-TV based upon the LAR/LASSO algorithm runs in at most $O(K_{\max}^3 + K_{\max}n \log(n))$ in time.

3. THEORETICAL RESULTS

In this section, we give some theoretical results providing justification on the relevance of LS-TV for multiple change-point estimation. First, in Section 3.1, we prove that LS-TV is consistent in terms of estimation of the underlying signal. Second, in Section 3.2, we show that LS-TV is also consistent in terms of change-point estimation.

The main point of both Section 3.1 and Section 3.2 is the following. While the equivalence of LS-TV to a particular LASSO problem is fruitful from a computational point of view, it turns out to be less relevant for theoretical analysis. To get optimal results for LS-TV both in terms of means and change-points estimation, the original formulation (2) is more useful than the LASSO formulation.

3.1 Estimation of the Means

We consider here the multiple changes in the mean problem as described in (1). Our purpose is to estimate the unknown means $\mu_1^*, \dots, \mu_{K^*+1}^*$ together with the change points from observations Y_1, \dots, Y_n .

Let us first work with the LASSO formulation to establish the consistency in terms of means estimation. The model (1) can be rewritten as

$$\mathbf{Y}^n = \mathbf{X}_n \boldsymbol{\beta}^n + \boldsymbol{\varepsilon}^n, \quad (7)$$

where $\mathbf{Y}^n = (Y_1, \dots, Y_n)^T$ is the $n \times 1$ vector of observations, \mathbf{X}_n the $n \times n$ triangular matrix with nonzero elements equal to one and $\boldsymbol{\varepsilon}^n = (\varepsilon_1^n, \dots, \varepsilon_n^n)^T$ is a zero mean random vector such that the $\varepsilon_1^n, \dots, \varepsilon_n^n$ are iid random variables with finite variance equal to σ^2 . As for $\boldsymbol{\beta}^n$, it is a $n \times 1$ vector having all its components equal to zero except those corresponding to the change-points instants.

Let us denote by \mathcal{A} the set of nonzero components of $\boldsymbol{\beta}^n$ and by $\bar{\mathcal{A}}$ its complementary set defined as follows:

$$\mathcal{A} = \{k, \beta_k^n \neq 0\} \quad \text{and} \quad \bar{\mathcal{A}} = \{1, \dots, n\} \setminus \mathcal{A}. \quad (8)$$

With the reformulation (7), the evaluation of the means estimation rate amounts to finding the rate of convergence of $\|\mathbf{X}_n(\hat{\boldsymbol{\beta}}^n(\lambda_n) - \boldsymbol{\beta}^n)\|_n$ to zero, $\hat{\boldsymbol{\beta}}^n(\lambda_n)$ satisfying:

$$\begin{aligned} \hat{\boldsymbol{\beta}}^n(\lambda_n) &= (\hat{\beta}_1(\lambda_n), \dots, \hat{\beta}_n(\lambda_n))^T \\ &= \underset{\boldsymbol{\beta} \in \mathbb{R}^n}{\text{Arg min}} \{ \|\mathbf{Y}^n - \mathbf{X}_n \boldsymbol{\beta}\|_n^2 + \lambda_n \|\boldsymbol{\beta}\|_1 \}, \end{aligned} \quad (9)$$

where $\|\mathbf{u}\|_n$ and $\|\mathbf{u}\|_1$ are defined for a vector $\mathbf{u} = (u_1, \dots, u_n) \in \mathbb{R}^n$ by $\|\mathbf{u}\|_n = n^{-1} \sum_{j=1}^n u_j^2$ and $\|\mathbf{u}\|_1 = \sum_{j=1}^n |u_j|$, respectively. Hence, within this framework, we are able to prove the following result regarding the consistency in means estimation of LS-TV.

Proposition 1. Consider Y_1, \dots, Y_n a set of observations following the model described in (7). Assume that the $\varepsilon_1^n, \dots, \varepsilon_n^n$ are centered iid Gaussian random variables with variance $\sigma^2 > 0$. Assume also that there exists β_{\max} such that for all k in \mathcal{A} , $|\beta_k^n| \leq \beta_{\max}$, the set \mathcal{A} being defined in (8). Then, for all $n \geq 1$ and $C > 2\sqrt{2}$, we obtain that with a probability larger than $1 - n^{1-C^2/8}$, if $\lambda_n = C\sigma\sqrt{\log n/n}$,

$$\|\mathbf{X}_n(\hat{\boldsymbol{\beta}}^n(\lambda_n) - \boldsymbol{\beta}^n)\|_n \leq (2C\sigma\beta_{\max}K^*)^{1/2} \left(\frac{\log n}{n} \right)^{1/4}.$$

The proof, which follows similar lines as Bickel, Ritov, and Tsybakov (2009), is postponed to Section 7. Note that in Proposition 1, where no upper bound on the number of change points is assumed to be known, we do not attain the known (parametric) optimal rate which is of order $1/\sqrt{n}$ derived by Yao and Au (1989) where an upper bound for the number of change points is available. But, as we shall see in Proposition 2, the rate of Proposition 1 can be improved if the model and the criterion are rewritten in a different way and if an upper bound for the number of change points is available.

Indeed, let us now work in the standard formulation of LS-TV instead of its LASSO counterpart, and write model (1) as

$$Y_t = u_t^* + \varepsilon_t, \quad t = 1, \dots, n, \quad (10)$$

where $u_t^* = \mu_k^*$ for $t_{k-1}^* \leq t \leq t_k^* - 1$, $k = 1, \dots, K^* + 1$ and estimate the vector (u_1^*, \dots, u_n^*) by using a criterion based on a total variation penalty as in Mammen and Van De Geer (1997):

$$\begin{aligned} \hat{\mathbf{u}}(\lambda_n) &= (\hat{u}_1(\lambda_n), \dots, \hat{u}_n(\lambda_n)) \\ &= \underset{\mathbf{u} \in \mathbb{R}^n}{\text{Arg min}} \left\{ \|\mathbf{Y}^n - \mathbf{u}\|_n^2 + \lambda_n \sum_{i=1}^{n-1} |u_{i+1} - u_i| \right\}. \end{aligned} \quad (11)$$

The following proposition gives the rate of convergence of $\hat{\mathbf{u}}(\lambda_n)$ when an upper bound for the number of change points is known and equal to K_{\max} .

Proposition 2. Consider Y_1, \dots, Y_n a set of observations following the model described in (10) where the $\varepsilon_1, \dots, \varepsilon_n$ are zero-mean iid Gaussian random variables with a variance $\sigma^2 > 0$. Assume also that $\hat{\mathbf{u}}$ defined in (11) belongs to a set of dimension at most $K_{\max} - 1$. Then, for all $n \geq 1$, A in $(0, 1)$ and $B > 0$, if $\lambda_n = \sigma(A\sqrt{B}/2)(K_{\max} \log n)^{1/2}n^{-3/2} - \sigma(2K_{\max} + 1)^{1/2}n^{-3/2}$,

$$\begin{aligned} \mathbb{P}(\|\hat{\mathbf{u}} - \mathbf{u}^*\|_n \geq \sigma(BK_{\max} \log n/n)^{1/2}) \\ \leq K_{\max} n^{\{1-B(1-A)^2/8\}K_{\max}}. \end{aligned} \quad (12)$$

The proof of this proposition is postponed to Section 7. The rate of convergence that we obtain for the estimation of the means is almost optimal up to a logarithmic factor since the optimal rate derived by Yao and Au (1989) is $O(n^{-1/2})$.

Let us now study the consistency in terms of change-point estimation, which is more of interest in this article. Again, we shall see that the LASSO formulation is less relevant than the standard formulation for establishing the change-point estimation consistency.

3.2 Estimation of the Change-Point Locations

In this section, we aim at estimating the change-point locations from the observations (Y_1, \dots, Y_n) satisfying Model (7). The change-point estimates that we propose to study are obtained from the $\hat{\beta}_i(\lambda_n)$'s satisfying the criterion (9) as follows. Let us define the set of active variables by

$$\hat{\mathcal{A}}(\lambda_n) = \{i \in \{1, \dots, n\}, \hat{\beta}_i(\lambda_n) \neq 0\}. \quad (13)$$

Then, we define the change-point estimates by $\hat{t}_i(\lambda_n)$ satisfying

$$\begin{aligned} \hat{\mathcal{A}}(\lambda_n) &= \{\hat{t}_1(\lambda_n), \dots, \hat{t}_{|\hat{\mathcal{A}}(\lambda_n)|}(\lambda_n)\}, \\ &\text{where } \hat{t}_1(\lambda_n) < \dots < \hat{t}_{|\hat{\mathcal{A}}(\lambda_n)|}(\lambda_n), \end{aligned} \quad (14)$$

$|\hat{\mathcal{A}}(\lambda_n)|$ denoting the cardinal of the set $\hat{\mathcal{A}}(\lambda_n)$.

Discussion and Related Works. With such a reformulation of the change point in the mean problem, the change-point estimates can be seen as Lasso-type estimates in a sparse framework. But, many classical assumptions under which the asymptotic properties of the Lasso estimates have been studied are not satisfied.

For instance, the *irrepresentable condition* as defined in (Meinshausen and Yu 2009, p. 5) which ensures the *sign consistency* as defined in Zhao and Yu (2006) is not satisfied in the change point in the mean problem. More precisely, *sign consistency* ensures that $\mathbb{P}(\text{sign}(\hat{\boldsymbol{\beta}}^n(\lambda_n)) = \text{sign}(\boldsymbol{\beta}^n))$ tends to one as n tends to infinity and the *irrepresentable condition* is a condition on the covariance matrix \mathbf{C}^n defined by

$$\mathbf{C}^n = n^{-1} \mathbf{X}_n^T \mathbf{X}_n,$$

which requires that the following inequality holds element-wise:

$$|\mathbf{C}_{\bar{\mathcal{A}}\mathcal{A}}^n (\mathbf{C}_{\mathcal{A}\mathcal{A}}^n)^{-1} \text{sign}(\boldsymbol{\beta}_{\mathcal{A}}^n)| < 1, \quad (15)$$

where \mathbf{C}_{IJ}^n is a submatrix of \mathbf{C}^n obtained by keeping rows with index in the set I and columns with index in J . The vector $\boldsymbol{\beta}_{\mathcal{A}}^n$ is defined by $\boldsymbol{\beta}_{\mathcal{A}}^n = (\beta_k^n)_{k \in \mathcal{A}}$ and *sign* denotes a function mapping positive entries of a vector to 1, negative entries to -1 and

null entries to zero. In our case, there exists at least one component i_0 such that

$$(|\mathbf{C}_{\mathcal{A}\mathcal{A}}^n (\mathbf{C}_{\mathcal{A}\mathcal{A}}^n)^{-1} \text{sign}(\boldsymbol{\beta}_{\mathcal{A}}^n)|)_{i_0} = 1.$$

This can be proved by computing explicitly the matrices $\mathbf{C}_{\mathcal{A}\mathcal{A}}^n$ and $(\mathbf{C}_{\mathcal{A}\mathcal{A}}^n)^{-1}$, see the [Appendix](#) for further details. In terms of change-point estimation, it means, as already known; see, for example, [Yao and Au \(1989\)](#) or [Lavielle and Moulines \(2000\)](#), that we cannot have a perfect estimation of the change points.

Note that [Meinshausen and Yu \(2009\)](#) brought to light some less restrictive conditions than the irrerepresentable condition on the matrix \mathbf{C}^n under which the Lasso estimates can be proved to be consistent in the ℓ_2 -norm sense. The main assumption consists in assuming a m_n -incoherent design which means

$$\liminf_{n \rightarrow \infty} \phi_{\min}(m_n) > 0,$$

$$\text{where } \phi_{\min}(m) = \min_{\boldsymbol{\beta}: \|\boldsymbol{\beta}\|_{\ell_0} \leq m} \frac{\boldsymbol{\beta}^T \mathbf{C}^n \boldsymbol{\beta}}{\boldsymbol{\beta}^T \boldsymbol{\beta}}, \quad (16)$$

with $m_n = s_n \log n$, s_n being the sparsity of the model that is the number of nonzero coefficients. In other words, a design is called m_n -incoherent if the minimal eigenvalue of a collection of m_n variables is bounded from below by a positive constant. In our setting, if the distance between two consecutive indices of nonnull coefficients is equal to one, then for all $n \geq 1$

$$\phi_{\min}(m_n) \leq 1/n,$$

this making the condition (16) not satisfied in our case. A justification of this statement is given in the [Appendix](#).

These particularities of the change point in the mean model prevent us from using the techniques recently devised to study the asymptotic properties of the Lasso estimates in a general regression framework. However, the consistency of the $\hat{t}_i(\lambda_n)$ defined in (14) is established in Proposition 3.

Let us now detail the assumptions under which our theoretical results are established. Define

$$I_{\min}^* = \min_{1 \leq k \leq K^*} |t_{k+1}^* - t_k^*|, \quad J_{\min}^* = \min_{1 \leq k \leq K^*} |\mu_{k+1}^* - \mu_k^*|,$$

$$J_{\max}^* = \max_{1 \leq k \leq K^*} |\mu_{k+1}^* - \mu_k^*|,$$

which are respectively the minimum interval length, the minimum and maximum jump sizes. From now on, we shall work under the following assumptions:

(A1) The $\varepsilon_1, \dots, \varepsilon_n$ are iid zero-mean random variables with $\text{var}[\varepsilon_1] = \sigma^2$ satisfying: there exists a positive constant β such that for all $v \in \mathbb{R}$, $\mathbb{E}\{\exp(v\varepsilon_1)\} \leq \exp(\beta v^2)$.

(A2) The sequence $\{\delta_n\}_{n \geq 1}$ is a nonincreasing and positive sequence tending to zero as n tends to infinity and satisfying $n\delta_n(J_{\min}^*)^2/\log(n) \rightarrow \infty$.

(A3) The change points $t_1^*, \dots, t_{K^*}^*$ satisfy $I_{\min}^* \geq n\delta_n$, for all $n \geq 1$.

(A4) The sequence of regularization parameters $\{\lambda_n\}_{n \geq 1}$ is such that $(n\delta_n J_{\min}^*)^{-1} n\lambda_n \rightarrow 0$, as n tends to infinity.

We first state a lemma arising from the Karush–Kuhn–Tucker conditions of the optimization problem stated in (9) which will be useful in the proof of the consistency of our procedure.

Lemma 1. Consider Y_1, \dots, Y_n a set of observations following the model described in (10). Then, $(\hat{t}_1(\lambda_n), \dots, \hat{t}_n(\lambda_n))$ defined by (14) and $(\hat{u}_1(\lambda_n), \dots, \hat{u}_n(\lambda_n))$ defined by $\hat{u}_i(\lambda_n) = (\mathbf{X}_n \hat{\boldsymbol{\beta}}^n(\lambda_n))_i$, where \mathbf{X}_n is a $n \times n$ lower triangular matrix with nonzero elements equal to one and the $(\hat{\beta}_i(\lambda_n))_{1 \leq i \leq n}$ are obtained in (9), satisfy

$$\sum_{i=\hat{t}_\ell(\lambda_n)}^n Y_i - \sum_{i=\hat{t}_\ell(\lambda_n)}^n \hat{u}_i = \frac{n\lambda_n}{2} \hat{\alpha}_\ell \quad \text{for all } \ell = 1, \dots, |\hat{\mathcal{A}}(\lambda_n)| \quad (17)$$

and

$$\left| \sum_{i=j}^n Y_i - \sum_{i=j}^n \hat{u}_i \right| \leq \frac{n\lambda_n}{2} \quad \text{for all } j = 1, \dots, n, \quad (18)$$

using the convention $\hat{\alpha}_\ell = +1$, if $\hat{u}_{\hat{t}_\ell(\lambda_n)} > \hat{u}_{\hat{t}_\ell(\lambda_n)-1}$ and $\hat{\alpha}_\ell = -1$, otherwise. The vector $(\hat{u}_1(\lambda_n), \dots, \hat{u}_n(\lambda_n))$ has the following additional property:

$$\hat{u}_t(\lambda_n) = \hat{\mu}_k \quad \text{for } \hat{t}_{k-1}(\lambda_n) \leq t \leq \hat{t}_k(\lambda_n) - 1, \quad k = 1, \dots, |\hat{\mathcal{A}}(\lambda_n)| + 1, \quad (19)$$

where $|\hat{\mathcal{A}}(\lambda_n)|$ denotes the cardinal of the set $\hat{\mathcal{A}}(\lambda_n)$ defined in (14).

The proof of Lemma 1 is given in Section 7. Then, we state a lemma which allows us to control the supremum of the average of the noise and which will also be useful for proving the consistency of our estimation criterion.

Lemma 2. Let $(\varepsilon_i)_{1 \leq i \leq n}$ be a sequence of random variables satisfying Assumption (A1). If $\{v_n\}_{n \geq 1}$ and $\{x_n\}_{n \geq 1}$ are two positive sequences such that $v_n x_n^2 / \log(n) \rightarrow \infty$, then

$$\mathbb{P} \left(\max_{\substack{1 \leq r_n < s_n \leq n \\ |r_n - s_n| \geq v_n}} \left| (s_n - r_n)^{-1} \sum_{i=r_n}^{s_n-1} \varepsilon_i \right| \geq x_n \right) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

The proof of Lemma 2 is postponed to Section 7.

Proposition 3. Let Y_1, \dots, Y_n be a set of observations satisfying Model (1) then under Assumptions (A1)–(A4), the change-points estimators $\{\hat{t}_1(\lambda_n), \dots, \hat{t}_{|\hat{\mathcal{A}}(\lambda_n)|}(\lambda_n)\}_{n \geq 1}$ defined by (14), satisfy, if $|\hat{\mathcal{A}}(\lambda_n)| = K^*$ with probability tending to one:

$$\mathbb{P} \left(\max_{1 \leq k \leq K^*} |\hat{t}_k - t_k^*| \leq n\delta_n \right) \rightarrow 1, \quad \text{as } n \rightarrow \infty. \quad (20)$$

The proof of Proposition 3 is given in Section 7.

Under the assumptions of Proposition 3, the \hat{t}_k 's defined for all $k \in \{1, \dots, K^*\}$ by $\hat{t}_k = [n\hat{\tau}_k]$ are consistent estimators of the τ_k^* 's defined by $t_k^* = [n\tau_k^*]$, for all $k \in \{1, \dots, K^*\}$ with the rate δ_n .

Note that with $\delta_n = (\log n)^2/n$, $J_{\min}^* \geq (\log n)^{1/4}$, $\lambda_n = \sqrt{\log n/n}$ or $\lambda_n = \sqrt{\log n/n^{3/2}}$, the Assumptions (A2)–(A4) are satisfied leading thus to a rate of order $(\log n)^2/n$ for the estimation of the \hat{t}_k . With this choice of parameters, we obtain an almost optimal rate for the estimation of the τ_k^* (up to a logarithmic factor) since the optimal rate is of order $1/n$ according to [Yao and Au \(1989\)](#).

This result has also to be compared with the work by [Lavielle and Moulines \(2000\)](#). They also obtained a rate in $1/n$ using a least-square approach in the case where the (ε_t) are not

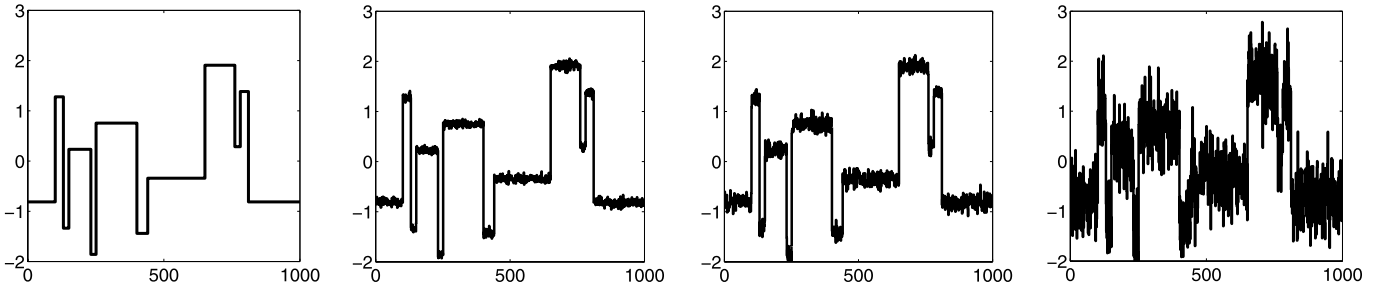


Figure 1. The Blocks dataset, subsampled to 1000 observations, and rescaled to mean zero and variance one, displayed without noise (on the far left) and with, respectively, low noise, medium noise, and high noise (from left to right).

necessarily independent random variables but with more restrictive assumptions than ours on J_{\min}^* and J_{\min}^* . Indeed, it is assumed in theorem 7 of [Lavielle and Moulines \(2000\)](#), that $\min_{1 \leq k \leq K^*} |\tau_{k+1}^* - \tau_k^*| = \Delta_\tau^*$ where Δ_τ^* is a positive constant and that J_{\min}^* is a positive constant.

In Proposition 3, the number of estimated change points is assumed to be equal to the true number of change points. Since this information is not in general available, we propose to evaluate the distance between the set $\widehat{T}_{n,K} = \{\hat{t}_1, \dots, \hat{t}_K\}$ of K estimated change points and the set of the true change points $T_n^* = \{t_1^*, \dots, t_{K^*}^*\}$ by using as in [Boysen et al. \(2009\)](#) the two quantities $\mathcal{E}(\widehat{T}_{n,K} \| T_n^*)$ and $\mathcal{E}(T_n^* \| \widehat{T}_{n,K})$, where $\mathcal{E}(\cdot \| \cdot)$ is defined for two sets A and B by

$$\mathcal{E}(A \| B) = \sup_{b \in B} \inf_{a \in A} |a - b|. \quad (21)$$

Note that we recover the Hausdorff distance between the sets A and B with

$$\Delta(A, B) = \sup\{\mathcal{E}(A \| B); \mathcal{E}(B \| A)\}.$$

Obviously, when $K = K^*$, Proposition 3 implies that, under the same assumptions, $\mathcal{E}(\widehat{T}_{n,K^*} \| T_n^*) \leq n\delta_n$ and $\mathcal{E}(T_n^* \| \widehat{T}_{n,K^*}) \leq n\delta_n$ with probability tending to one as n tends to infinity. In the case where $K > K^*$, we prove in Proposition 4 that $\mathcal{E}(\widehat{T}_{n,K} \| T_n^*) \leq n\delta_n$ with probability tending to one as n tends to infinity.

Proposition 4. Let Y_1, \dots, Y_n be a set of observations satisfying Model (1) then under Assumptions (A1), (A3), (A4) and if $n\delta_n J_{\min}^* / \log(n^3/\lambda_n^2) \rightarrow \infty$, the change-points estimators $\{\hat{t}_1(\lambda_n), \dots, \hat{t}_{|\hat{A}(\lambda_n)|}(\lambda_n)\}_{n \geq 1}$ defined by (14), satisfy, if $|\hat{A}(\lambda_n)| \geq K^*$ with probability tending to one:

$$\mathbb{P}(\mathcal{E}(\widehat{T}_{n,|\hat{A}(\lambda_n)|} \| T_n^*) \leq n\delta_n) \rightarrow 1, \quad \text{as } n \rightarrow \infty. \quad (22)$$

Note that with $\delta_n = (\log n)^2/n$, $J_{\min}^* \geq (\log n)^{1/4}$, $\lambda_n = \sqrt{\log n/n}$ or $\lambda_n = \sqrt{\log n/n^{3/2}}$, the Assumptions (A3), (A4), and $n\delta_n J_{\min}^* / \log(n^3/\lambda_n^2) \rightarrow \infty$ of Proposition 4 are fulfilled.

Now, we shall investigate the empirical behavior of LS-TV on simulated data. In the remainder, we focus on the so-called Blocks dataset introduced in [Donoho and Johnstone \(1995\)](#) which contains $K^* = 11$ change points. One may indeed consider the Blocks dataset as a typically difficult dataset for multiple change-point estimation, since both segment levels and segment lengths are highly heterogeneous.

4. EXPERIMENTAL RESULTS

4.1 Specified Number of Change Points

The Blocks dataset introduced in the article ([Donoho and Johnstone 1995](#), page 1201, table 1) was subsampled down to 1000 points as depicted in Figure 1, and corrupted with Gaussian white noise at three different levels: *low-noise* when $\sigma = 0.05$, *medium-noise* when $\sigma = 0.10$, and *high-noise* when $\sigma = 0.50$.

To assess our large-sample consistency result which stated that $n^{-1}\mathcal{E}(\widehat{T}_{n,K^*} \| T^*) = o_P(1)$, as n tends to infinity, we ran Monte Carlo simulations to investigate the empirical performance of LS-TV in terms of $n^{-1}\mathcal{E}(\widehat{T}_{n,K^*} \| T^*) = n^{-1} \max_{k=1, \dots, K^*} |\hat{t}_k - t_k^*|$ in the three different noise settings. For each noise setting, we generated 100 replications of the Blocks dataset corrupted with Gaussian white noise. The results are displayed in Table 2. In all noise conditions, the large-sample change-point estimation consistency of LS-TV is confirmed. In high-noise conditions, even for medium-scale samples, that is for $n = 1000$, the change-point detection ability of LS-TV remains satisfactory. For large-scale samples, that is, for $n = 5000$, the performance continue improving both on average and standard deviation.

Since, in general, the number of change points is unknown, we shall investigate in the next section the impact of misspecifying the number of change points. For this purpose, we study the evolution of both $\{\mathcal{E}(\widehat{T}_{n,K} \| T^*), \mathcal{E}(T^* \| \widehat{T}_{n,K})\}$, as $K = 1, \dots, 3K^*$ in the three different noise settings.

4.2 Unspecified Number of Change Points

4.2.1 Performance of LS-TV. We consider here the performance of LS-TV on the Blocks dataset corrupted with three different levels of noise, when the true number of change points is unknown. For each noise level, we generated 100 replications of

Table 2. Performance in terms of $\mathcal{E}(\widehat{T}_{n,K^*} \| T^*)$ of LS-TV on the Blocks dataset corrupted with low noise ($\sigma = 0.05$), medium noise ($\sigma = 0.10$), and high noise ($\sigma = 0.50$). The values after \pm correspond to the standard deviations

	Low noise	Medium noise	High noise
$n = 1000$	0.0200 ± 0.0068	0.0200 ± 0.0098	0.0230 ± 0.0185
$n = 5000$	0.0127 ± 0.0059	0.0127 ± 0.0082	0.0127 ± 0.0169

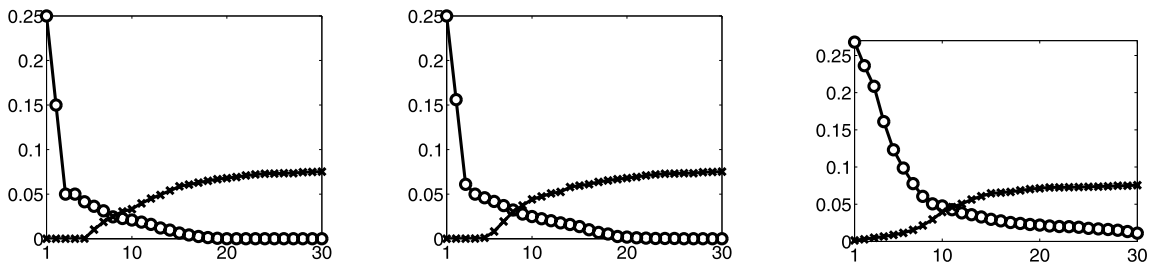


Figure 2. The evolution of the two types of error as $K = 1, \dots, 3K^*$, that is, $\{\mathcal{E}(\hat{T}_{n,K}^{\text{LS-TV}} \parallel T^*)\}_{K=1, \dots, 3K^*}$ ("o") and $\{\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS-TV}})\}_{K=1, \dots, 3K^*}$ ("x"), in different noise settings (low, medium, high noise from left to right).

corrupted versions of the Blocks dataset. For each noise replication, we measured both $\mathcal{E}(\hat{T}_{n,K} \parallel T^*)$ and $\mathcal{E}(T^* \parallel \hat{T}_{n,K})$ for all $K = 1, \dots, 3K^*$. We display in Figure 2 the results averaged over all replications for both errors. Also, note that the optimal trade-off between the two types of error is reached almost exactly at the true number of change points $K = K^*$.

4.2.2 Comparison With the Standard Least-Square (LS) Approach. Let us now compare the performance of LS-TV with the performance of the standard least-square estimation of multiple change points theoretically studied by Yao and Au (1989). The latter criterion provides a number of K change points for the model (1) by

$$\text{Minimize } \sum_{t_1 < \dots < t_K} \sum_{k=1}^K \sum_{i=t_{k-1}+1}^{t_k} (Y_i - \bar{\mu}_k)^2,$$

$$\text{where } \bar{\mu}_k \stackrel{\text{def}}{=} (t_k - t_{k-1})^{-1} \sum_{i=t_{k-1}+1}^{t_k} Y_i.$$

A computationally efficient way of solving this minimization is based on a dynamic programming algorithm (DP), originally proposed by Fisher (1958); Bellman (1961) and described in Kay (1993), chapter 12. While a naive approach would require a $O(2^n)$ time complexity, DP has a time complexity of $O(Kn^2)$ if we look for at most K change points within the signal. For a fair comparison, we used exactly the same settings for both methods LS-TV and LS.

From Table 3 displayed in Section 5, we can see that LS-TV reaches satisfactory performance, in terms of both types of errors, in all noise settings as well as LS. It is worthwhile to emphasize that, while LS has a $O(Kn^2)$ time complexity when implemented with the DP algorithm, our method LS-TV has at

most $O(Kn \log(n))$ time complexity. We can also remark that the variance of $\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS-TV}})$ is larger than the variance of $\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS}})$. It is then interesting to remedy this issue, without harming the subquadratic time complexity of LS-TV.

In the next section, we show how LS-TV may be further enhanced, both in mean and variance in both types of errors, when combined with an additional step based on a reduced-search dynamic programming algorithm.

5. AN ENHANCED VERSION OF LS-TV: LS-TV*

We now propose an enhanced version of LS-TV, called LS-TV*, which combines two steps. First, we run LS-TV with $K = K_{\max}$ larger than K^* , and get a set of change-point estimates $\hat{T}_{n,K_{\max}}$. Second, we run a reduced version of DP searching $L < K_{\max}$ change points over the set $\hat{T}_{n,K_{\max}}$, instead of $\{1, \dots, n\}$ as in the raw DP algorithm, which finally yields a new set of change-point estimates $\mathcal{S}_{n,L} \subseteq \hat{T}_{n,K_{\max}}$.

From Table 3, we observe that for $K = 30$ the error $\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS}})$ becomes larger than $\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS-TV}})$ in all noise settings. This suggests that one type of error made by LS-TV stabilizes in the over-segmentation regime, that is, when $K \gg K^*$, whereas the same type of error made by LS still increases. Therefore, one might think of running LS-TV to look for an *a priori* much larger set of change points than the true number of change points, that is, to look for $K_{\max} \gg K^*$ change points with $K_{\max} \ll n$, and then propose a way of selecting the best change-point estimates within the large set of change-point estimates obtained by LS-TV.

We suggest running a dynamic programming algorithm to perform this postselection. More precisely, we aim at minimiz-

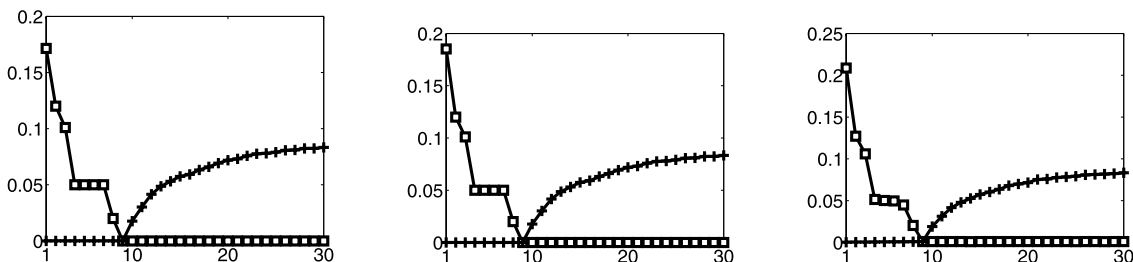


Figure 3. The evolution of the two types of error as $K = 1, \dots, 3K^*$, that is, $\{\mathcal{E}(\hat{T}_{n,K}^{\text{LS-TV}^*} \parallel T^*)\}_{K=1, \dots, 3K^*}$ displayed with squares and $\{\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS-TV}^*})\}_{K=1, \dots, 3K^*}$ ("x"), in different noise settings (low, medium, and high noise from left to right).

Table 3. Performance in terms of $\mathcal{E}(\hat{T}_{n,K} \parallel T^*)$ and $\mathcal{E}(T^* \parallel \hat{T}_{n,K})$ for different values of K of LS, LS-TV and LS-TV* on the Blocks dataset corrupted with low noise ($\sigma = 0.05$), medium noise ($\sigma = 0.10$), and high noise ($\sigma = 0.50$). For each method, the first and second lines correspond to the mean and standard deviation of $\mathcal{E}(\hat{T}_{n,K} \parallel T^*)$, respectively, and the third and fourth lines correspond to the mean and standard deviation of $\mathcal{E}(T^* \parallel \hat{T}_{n,K})$, respectively. K_{\max} was set to 30 in all experiments

	$K = 1$			$K = 11$			$K = 20$			$K = 30$		
	Low	Medium	High	Low	Medium	High	Low	Medium	High	Low	Medium	High
LS	0.169	0.169	0.169	0.000	0.000	0.001	0.000	0.000	0.001	0.000	0.000	0.001
	(0.014)	(0.033)	(0.041)	(0.000)	(0.000)	(0.001)	(0.000)	(0.000)	(0.001)	(0.000)	(0.000)	(0.001)
	0.000	0.000	0.000	0.023	0.023	0.027	0.072	0.072	0.072	0.086	0.086	0.086
	(0.000)	(0.000)	(0.001)	(0.025)	(0.025)	(0.025)	(0.018)	(0.018)	(0.017)	(0.014)	(0.014)	(0.014)
LS-TV	0.250	0.250	0.250	0.020	0.020	0.040	0.000	0.000	0.020	0.000	0.000	0.019
	(0.000)	(0.000)	(0.042)	(0.006)	(0.006)	(0.009)	(0.000)	(0.000)	(0.007)	(0.000)	(0.000)	(0.009)
	0.000	0.000	0.001	0.034	0.041	0.042	0.071	0.071	0.075	0.081	0.081	0.081
	(0.000)	(0.000)	(0.002)	(0.030)	(0.031)	(0.028)	(0.025)	(0.025)	(0.022)	(0.020)	(0.020)	(0.020)
LS-TV*	0.169	0.169	0.169	0.000	0.000	0.001	0.000	0.000	0.001	0.000	0.000	0.001
	(0.014)	(0.033)	(0.041)	(0.000)	(0.000)	(0.005)	(0.000)	(0.000)	(0.001)	(0.000)	(0.000)	(0.001)
	0.000	0.000	0.000	0.029	0.029	0.033	0.071	0.071	0.072	0.082	0.082	0.082
	(0.000)	(0.000)	(0.001)	(0.023)	(0.023)	(0.024)	(0.015)	(0.015)	(0.014)	(0.013)	(0.013)	(0.013)

ing, for each K in $\{1, \dots, K_{\max}\}$:

$$\begin{aligned} & \text{Minimize}_{t_1 < \dots < t_K} \sum_{k=1}^K \sum_{i=t_{k-1}+1}^{t_k} (Y_i - \bar{\mu}_k)^2, \\ & \text{s.t. } t_1, \dots, t_K \in \hat{T}_{n, K_{\max}} \\ & \text{where } \bar{\mu}_k \stackrel{\text{def}}{=} (t_k - t_{k-1})^{-1} \sum_{i=t_{k-1}+1}^{t_k} Y_i. \quad (23) \end{aligned}$$

The above algorithm, subsequently called rDP, outputs for each $K = 1, \dots, K_{\max}$ a new set of change-point estimates $\mathcal{S}_{n,K} \subsetneq \hat{T}_{n, K_{\max}}$. We call LS-TV* the method which combines LS-TV with a postselection based on rDP.

First, we investigate how LS-TV* improves on LS-TV in terms of error variance. The settings are the same as previously. We observe in Table 3 that the postselection step indeed consistently reduces the variance of both errors obtained by LS-TV.

Second, we check whether LS-TV* improves, on average, the performance of LS-TV. As Figures 2, 4, and Table 3 show, not only LS-TV* yields much lower error rates than LS-TV in both types of errors, but LS-TV* also obtains similar error rates when compared to LS. Since the overall time complexity of LS-TV* is $O(K_{\max}^3 + K_{\max}n \log(n))$, and the overall time complexity of LS is $O(K_{\max}n^2)$, then, as long as $K^* < K_{\max} \ll n$, LS-TV* obtains the same performance results as LS at a much

lower computational cost. In order to give an idea to the reader of the actual computation times of LS-TV* and LS, we give in Table 4 the computation times of both methods when they are applied to the Blocks dataset for several values of n and K_{\max} .

Note that Figure 4 gives an appealing intuitive understanding of the statistical behavior of multiple change-point estimation methods. While the first type of error $\mathcal{E}(\hat{T} \parallel T^*)$ may be interpreted as the maximum error in the change-point location from estimated change points to true change points, the second type of error $\mathcal{E}(T^* \parallel \hat{T})$ may be interpreted as the maximum error in the change-point location from true change points to estimated change points. As the number of estimated change points increases, the first type of error $\mathcal{E}(\hat{T} \parallel T^*)$ decreases while the second type of error $\mathcal{E}(T^* \parallel \hat{T})$ increases. Finally, $\mathcal{E}(\hat{T} \parallel T^*)$ quantifies the over-segmentation error while $\mathcal{E}(T^* \parallel \hat{T})$ quantifies the under-segmentation error.

As presented here LS-TV* does not include a model selection part. A thorough practical version of LS-TV* should incorporate a data-driven way of choosing the optimal number of change points \hat{K} , and hence the optimal set of change-point estimates $\mathcal{S}_{n, \hat{K}}$. For interested readers, we proposed in Harchaoui and Lévy-Leduc (2008) an efficient practical approach to address this issue.

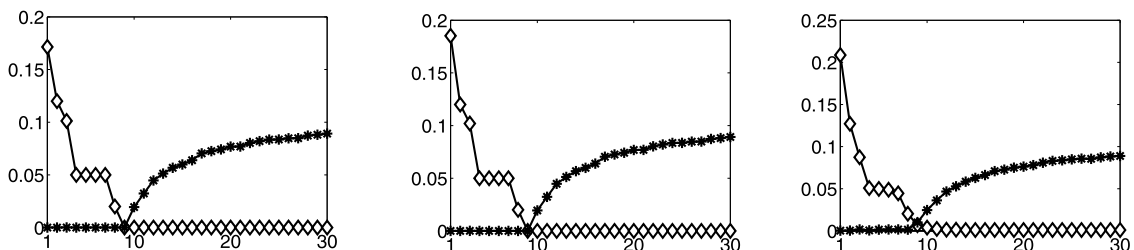


Figure 4. The evolution of the two types of error as $K = 1, \dots, 3K^*$, that is, $\{\mathcal{E}(\hat{T}_{n,K}^{\text{LS-TV}^*} \parallel T^*)\}_{K=1, \dots, 3K^*}$ (“ \diamond ”) and $\{\mathcal{E}(T^* \parallel \hat{T}_{n,K}^{\text{LS-TV}^*})\}_{K=1, \dots, 3K^*}$ (“ $*$ ”), in different noise settings (low, medium, and high noise from left to right).

Table 4. Computation times in seconds of LS and LS-TV* for several values of n and K_{\max}

(n, K_{\max})	(100, 5)	(500, 15)	(1000, 30)
LS	0.021 s	0.466 s	2.464 s
LS-TV*	0.005 s	0.119 s	0.689 s

6. CONCLUSION AND PROSPECTS

The standard least-square estimation approach LS suffers from an overwhelming time complexity for performing change-point estimation in long time series of observations. We showed, both theoretically and practically, that an alternative solution to the multiple change-point estimation problem, solved by a least-square fitting with a total variation penalty LS-TV, allowed us to get a lower time complexity while keeping competitive performance in terms of change-point estimation, even in high-noise settings.

We see several future research directions for this work. In the last section of the article, we proposed an enhanced version of LS-TV called LS-TV*, with better empirical performance and similar time complexity. We would like to provide thorough theoretical support to this method, which would involve a statistical analysis of the two steps LS-TV and reduced DP (rDP). Besides, since a lot of real datasets include a nonnegligible proportion of outliers, we would like to derive a robust version of both LS-TV and LS-TV*, and establish the corresponding theoretical results.

7. PROOFS

Proof of Proposition 1. By definition of $\hat{\beta}^n(\lambda_n)$ given by (9), we have

$$\|\mathbf{Y}^n - \mathbf{X}_n \hat{\beta}^n(\lambda_n)\|_n^2 + \lambda_n \|\hat{\beta}^n(\lambda_n)\|_1 \leq \|\mathbf{Y}^n - \mathbf{X}_n \beta^n\|_n^2 + \lambda_n \|\beta^n\|_1.$$

Using (7), we get

$$\begin{aligned} \|\mathbf{X}_n(\beta^n - \hat{\beta}^n(\lambda_n))\|_n^2 + \frac{2}{n}(\beta^n - \hat{\beta}^n(\lambda_n))^T \mathbf{X}_n^T \boldsymbol{\varepsilon}^n \\ + \lambda_n \sum_{k=1}^n |\hat{\beta}_k(\lambda_n)| \leq \lambda_n \sum_{k=1}^n |\beta_k^n|. \end{aligned}$$

Thus,

$$\begin{aligned} \|\mathbf{X}_n(\beta^n - \hat{\beta}^n(\lambda_n))\|_n^2 \\ \leq \frac{2}{n}(\hat{\beta}^n(\lambda_n) - \beta^n)^T \mathbf{X}_n^T \boldsymbol{\varepsilon}^n \\ + \lambda_n \sum_{j \in \mathcal{A}} (|\beta_j^n| - |\hat{\beta}_j(\lambda_n)|) - \lambda_n \sum_{j \in \bar{\mathcal{A}}} |\hat{\beta}_j(\lambda_n)|. \end{aligned}$$

Observe that

$$\frac{2}{n}(\hat{\beta}^n(\lambda_n) - \beta^n)^T \mathbf{X}_n^T \boldsymbol{\varepsilon}^n = 2 \sum_{j=1}^n (\hat{\beta}_j(\lambda_n) - \beta_j^n) \left(\frac{1}{n} \sum_{i=j}^n \varepsilon_i^n \right).$$

Let us define the event $E = \bigcap_{j=1}^n \{n^{-1} |\sum_{i=j}^n \varepsilon_i^n| \leq \lambda_n/2\}$. Then, given that the $\varepsilon_1^n, \dots, \varepsilon_n^n$ are iid zero-mean Gaussian random

variables with variance σ^2 , we obtain that

$$\begin{aligned} \mathbb{P}(\bar{E}) &\leq \sum_{j=1}^n \mathbb{P}\left(n^{-1} \left| \sum_{i=j}^n \varepsilon_i^n \right| > \lambda_n/2\right) \\ &\leq \sum_{j=1}^n \exp\left(-\frac{n^2 \lambda_n^2}{8\sigma^2(n-j+1)}\right). \end{aligned}$$

Thus, if $\lambda_n = C\sigma\sqrt{\log n/n}$,

$$\mathbb{P}(\bar{E}) \leq n^{1-C^2/8}.$$

With a probability larger than $1 - n^{1-C^2/8}$, we get

$$\begin{aligned} \|\mathbf{X}_n(\beta^n - \hat{\beta}^n(\lambda_n))\|_n^2 \\ \leq \lambda_n \sum_{j=1}^n |\hat{\beta}_j(\lambda_n) - \beta_j^n| \\ + \lambda_n \sum_{j \in \mathcal{A}} (|\beta_j^n| - |\hat{\beta}_j(\lambda_n)|) - \lambda_n \sum_{j \in \bar{\mathcal{A}}} |\hat{\beta}_j(\lambda_n)|, \end{aligned}$$

where \mathcal{A} and $\bar{\mathcal{A}}$ are defined in (8). Given that

$$\sum_{j=1}^n |\hat{\beta}_j(\lambda_n) - \beta_j^n| = \sum_{j \in \mathcal{A}} |\hat{\beta}_j(\lambda_n) - \beta_j^n| + \sum_{j \in \bar{\mathcal{A}}} |\hat{\beta}_j(\lambda_n)|,$$

we obtain that, with a probability larger than $1 - n^{1-C^2/8}$,

$$\begin{aligned} \|\mathbf{X}_n(\beta^n - \hat{\beta}^n(\lambda_n))\|_n^2 &\leq 2\lambda_n \sum_{j \in \mathcal{A}} |\beta_j^n| = 2C\sigma\sqrt{\frac{\log n}{n}} \sum_{j \in \mathcal{A}} |\beta_j^n| \\ &\leq 2C\sigma\beta_{\max} K^* \sqrt{\frac{\log n}{n}}. \end{aligned}$$

Proof of Proposition 2. For notational simplicity, we shall remove the dependence of $\hat{\mathbf{u}}$ in λ_n . By definition of $\hat{\mathbf{u}}$ as a minimizer of the criterion (11), we get:

$$\begin{aligned} \|\mathbf{Y}^n - \hat{\mathbf{u}}\|_n^2 + \lambda_n \sum_{i=1}^{n-1} |\hat{u}_{i+1} - \hat{u}_i| \\ \leq \|\mathbf{Y}^n - \mathbf{u}^*\|_n^2 + \lambda_n \sum_{i=1}^{n-1} |u_{i+1}^* - u_i^*|. \end{aligned}$$

Using Model (10), the previous inequality can be rewritten as follows:

$$\begin{aligned} \|\hat{\mathbf{u}} - \mathbf{u}^*\|_n^2 &\leq \lambda_n \left(\sum_{i=1}^{n-1} |u_{i+1}^* - u_i^*| - \sum_{i=1}^{n-1} |\hat{u}_{i+1} - \hat{u}_i| \right) \\ &\quad + \frac{2}{n} \sum_{i=1}^n \varepsilon_i (\hat{u}_i - u_i^*). \end{aligned}$$

Using the Cauchy Schwarz inequality, we obtain

$$\|\hat{\mathbf{u}} - \mathbf{u}^*\|_n^2 \leq 2n\lambda_n \|\hat{\mathbf{u}} - \mathbf{u}^*\|_n + \frac{2}{n} \sum_{i=1}^n \varepsilon_i (\hat{u}_i - u_i^*).$$

Thus, defining $G(\cdot)$ for \mathbf{v} in \mathbb{R}^n by $G(\mathbf{v}) = (\sum_{i=1}^n \varepsilon_i(v_i - u_i^*)) / (\sigma \sqrt{n} \|\mathbf{v} - \mathbf{u}^*\|_n)$, we get

$$\|\hat{\mathbf{u}} - \mathbf{u}^*\|_n^2 \leq 2n\lambda_n \|\hat{\mathbf{u}} - \mathbf{u}^*\|_n + \frac{2\sigma}{\sqrt{n}} \|\hat{\mathbf{u}} - \mathbf{u}^*\|_n G(\hat{\mathbf{u}}).$$

Let $\{S_K\}_{1 \leq K \leq K_{\max}}$ be the collection of linear spaces to which $\hat{\mathbf{u}}$ may belong, S_K denoting a space of dimension K . Then, given that the number of sets of dimension K is bounded by n^K , we obtain

$$\begin{aligned} \mathbb{P}(\|\hat{\mathbf{u}} - \mathbf{u}^*\|_n \geq \alpha_n) &\leq \mathbb{P}(n\lambda_n + \sigma n^{-1/2} G(\hat{\mathbf{u}}) \geq \alpha_n/2) \\ &\leq \sum_{K=1}^{K_{\max}} n^K \mathbb{P}\left(\sup_{\mathbf{v} \in S_K} G(\mathbf{v}) \geq n^{1/2} \sigma^{-1} \alpha_n/2 - n^{3/2} \sigma^{-1} \lambda_n\right). \end{aligned} \quad (24)$$

Using that, $\text{var}(G(\mathbf{v})) = 1$, for all \mathbf{v} in \mathbb{R}^n , we obtain by using an inequality due to Cirel'son, Ibragimov, and Sudakov in the same way as in the proof of theorem 1 in Birgé and Massart (2001), that for all $\beta > 0$,

$$\mathbb{P}\left(\sup_{\mathbf{v} \in S_K} G(\mathbf{v}) \geq \mathbb{E}\left[\sup_{\mathbf{v} \in S_K} G(\mathbf{v})\right] + \beta\right) \leq \exp(-\beta^2/2). \quad (25)$$

Let us now find an upper bound for $\mathbb{E}[\sup_{\mathbf{v} \in S_K} G(\mathbf{v})]$. Denoting by W the D -dimensional space to which $\mathbf{v} - \mathbf{u}^*$ belongs and some orthogonal basis ψ_1, \dots, ψ_D of W , we obtain

$$\begin{aligned} \sup_{\mathbf{v} \in S_K} G(\mathbf{v}) &\leq \sup_{\mathbf{w} \in W} \frac{\sum_{i=1}^n \varepsilon_i w_i}{\sigma \sqrt{n} \|\mathbf{w}\|_n} = \sup_{\alpha \in \mathbb{R}^D} \frac{\sum_{i=1}^n \varepsilon_i (\sum_{j=1}^D \alpha_j \psi_{j,i})}{\sigma \sqrt{n} \|\sum_{j=1}^D \alpha_j \psi_j\|_n} \\ &= \sup_{\alpha \in \mathbb{R}^D} \frac{\sum_{i=1}^n \varepsilon_i (\sum_{j=1}^D \alpha_j \psi_{j,i})}{\sigma \sqrt{n} (\sum_{j=1}^D \alpha_j^2)^{1/2}}. \end{aligned}$$

Using the Cauchy Schwarz inequality, we derive

$$\begin{aligned} \sup_{\mathbf{v} \in S_K} G(\mathbf{v}) &\leq \sup_{\alpha \in \mathbb{R}^D} \frac{\sum_{j=1}^D \alpha_j (\sum_{i=1}^n \varepsilon_i \psi_{j,i})}{\sigma \sqrt{n} (\sum_{j=1}^D \alpha_j^2)^{1/2}} \\ &\leq (\sigma^2 n)^{-1/2} \left\{ \sum_{j=1}^D \left(\sum_{i=1}^n \varepsilon_i \psi_{j,i} \right)^2 \right\}^{1/2}. \end{aligned}$$

By the concavity of the square-root function and by using that $D \leq K_{\max} + K^* + 1 \leq 2K_{\max} + 1$, we get

$$\mathbb{E}\left[\sup_{\mathbf{v} \in S_K} G(\mathbf{v})\right] \leq (2K_{\max} + 1)^{1/2}. \quad (26)$$

Using (24), (25), and (26) with $\beta = n^{1/2} \sigma^{-1} \alpha_n/2 - n^{3/2} \sigma^{-1} \lambda_n - (2K_{\max} + 1)^{1/2}$, we get

$$\begin{aligned} \mathbb{P}(\|\hat{\mathbf{u}} - \mathbf{u}^*\|_n \geq \alpha_n) &\leq K_{\max} \exp\left\{K_{\max} \log n \right. \\ &\quad \left. - \frac{1}{2} \left(\frac{n^{1/2} \alpha_n}{2\sigma} - n^{3/2} \sigma^{-1} \lambda_n - (2K_{\max} + 1)^{1/2} \right)^2 \right\}, \end{aligned}$$

which is valid only if $\beta = n^{1/2} \sigma^{-1} \alpha_n/2 - n^{3/2} \sigma^{-1} \lambda_n - (2K_{\max} + 1)^{1/2}$ is positive. Thus, writing for a constant A in $(0, 1)$,

$$n^{3/2} \sigma^{-1} \lambda_n + (2K_{\max} + 1)^{1/2} = A n^{1/2} \sigma^{-1} \alpha_n/2,$$

gives

$$\mathbb{P}(\|\hat{\mathbf{u}} - \mathbf{u}^*\|_n \geq \alpha_n) \leq K_{\max} \exp\left\{K_{\max} \log n - \frac{(1-A)^2 n \alpha_n^2}{8 \sigma^2}\right\}.$$

Thus, if $\alpha_n = (B \sigma^2 K_{\max} \log n / n)^{1/2}$, we obtain the expected result.

Proof of Lemma 1. A necessary and sufficient condition for a vector $\hat{\beta}$ in \mathbb{R}^n to minimize Φ defined by $\Phi(\beta) = \sum_{i=1}^n (Y_i - (X_n \beta)_i)^2 + n \lambda_n \sum_{i=1}^n |\beta_i|$, is that the zero vector in \mathbb{R}^n belongs to the subdifferential of Φ at point $\hat{\beta}$, that is,

$$(\mathbf{X}_n^T (\mathbf{Y}^n - \mathbf{X}_n \hat{\beta}))_j = \frac{n \lambda_n}{2} \text{sign}(\hat{\beta}_j), \quad \text{if } \hat{\beta}_j \neq 0,$$

$$|(\mathbf{X}_n^T (\mathbf{Y}^n - \mathbf{X}_n \hat{\beta}))_j| \leq \frac{n \lambda_n}{2}, \quad \text{if } \hat{\beta}_j = 0.$$

Using that $(\mathbf{X}_n^T \mathbf{Y}^n)_j = \sum_{k=j}^n Y_k$ and that $(\mathbf{X}_n^T \hat{\mathbf{u}})_j = \sum_{k=j}^n \hat{u}_k$, since \mathbf{X}_n is a $n \times n$ lower triangular matrix having all its nonzero elements equal to one, we obtain the expected result.

In the remainder, for any sequence of random variables, say, Z_1, \dots, Z_n , we shall use the following notation:

$$Z(r; s) \stackrel{\text{def}}{=} \sum_{i=r}^s Z_i \quad \text{for any } 1 \leq r < s \leq n. \quad (27)$$

Proof of Lemma 2. Using the notation introduced in (27), we obtain

$$\begin{aligned} \mathbb{P}\left(\max_{\substack{1 \leq r_n < s_n \leq n \\ |r_n - s_n| \geq v_n}} \left| \frac{\varepsilon(r_n; s_n - 1)}{s_n - r_n} \right| \geq x_n\right) &\leq \sum_{\substack{1 \leq r_n < s_n \leq n \\ |r_n - s_n| \geq v_n}} \mathbb{P}\left(\left| \frac{\varepsilon(r_n; s_n - 1)}{s_n - r_n} \right| \geq x_n\right). \end{aligned}$$

Using Assumption (A1), we get that for all $\eta > 0$,

$$\begin{aligned} \mathbb{P}\left(\frac{\varepsilon(r_n; s_n - 1)}{s_n - r_n} \geq x_n\right) &\leq \exp[-\eta(s_n - r_n)x_n] [\mathbb{E}\{\exp(\eta \varepsilon)\}]^{(s_n - r_n)} \\ &\leq \exp[-\eta(s_n - r_n)x_n + \beta \eta^2 (s_n - r_n)]. \end{aligned}$$

Since the sharpest bound holds for $\eta = x_n/2\beta$, we get

$$\mathbb{P}\left(\frac{\varepsilon(r_n; s_n - 1)}{s_n - r_n} \geq x_n\right) \leq \exp[-x_n^2 (s_n - r_n)/4\beta].$$

Since the same bound is valid when ε_i is replaced by $-\varepsilon_i$, we get that

$$\mathbb{P}\left(\left| \frac{\varepsilon(r_n; s_n - 1)}{s_n - r_n} \right| \geq x_n\right) \leq 2 \exp[-x_n^2 (s_n - r_n)/4\beta].$$

Hence, we obtain that

$$\mathbb{P}\left(\max_{\substack{1 \leq r_n < s_n \leq n \\ |r_n - s_n| \geq v_n}} \left| \frac{\varepsilon(r_n; s_n - 1)}{s_n - r_n} \right| \geq x_n\right) \leq 2n^2 \exp[-v_n x_n^2/4\beta],$$

which completes the proof.

Proof of Proposition 3. In this proof, we shall use the notation introduced in (27). Since $\mathbb{P}(\max_{1 \leq k \leq K^*} |\hat{t}_k - t_k^*| > n\delta_n) \leq \sum_{k=1}^{K^*} \mathbb{P}(|\hat{t}_k - t_k^*| > n\delta_n)$, it suffices to prove that for all $k =$

$1, \dots, K^*$, $\mathbb{P}(A_{n,k}) \rightarrow 0$, where $A_{n,k} = \{|\hat{t}_k - t_k^*| \geq n\delta_n\}$. Defining the set C_n by

$$C_n = \left\{ \max_{0 \leq k \leq K^*} |\hat{t}_k - t_k^*| < I_{\min}^*/2 \right\}, \quad (28)$$

it is enough to prove that $\mathbb{P}(A_{n,k} \cap C_n) \rightarrow 0$, and that $\mathbb{P}(A_{n,k} \cap \bar{C}_n) \rightarrow 0$. Let us first prove the first statement. Note that (28) implies that

$$t_{k-1}^* < \hat{t}_k < t_{k+1}^* \quad \text{for all } k \in \{1, \dots, K^*\}.$$

Let us first consider the case where $\hat{t}_k \leq t_k^*$. Applying (18) in Lemma 1 with $j = t_k^*$ and (17) in Lemma 1 with $\ell = k$ gives, respectively,

$$\left| \sum_{i=t_k^*}^n Y_i - \sum_{i=t_k^*}^n \hat{u}_i \right| \leq n\lambda_n/2 \quad \text{and} \\ \sum_{i=\hat{t}_k}^n Y_i - \sum_{i=\hat{t}_k}^n \hat{u}_i = n\hat{\alpha}_k \lambda_n/2.$$

This yields, using (19) in Lemma 1, that the event $C_{n,k}$ defined as follows, occurs with probability one:

$$C_{n,k} = \left\{ |(\hat{t}_k - t_k^*)(\mu_{k+1}^* - \mu_k^*) + (\hat{t}_k - t_k^*)(\hat{\mu}_{k+1} - \mu_{k+1}^*) + \varepsilon(\hat{t}_k; t_k^* - 1)| \leq n\lambda_n \right\}.$$

Using that $\mathbb{P}(A_{n,k} \cap C_n) = \mathbb{P}(A_{n,k} \cap C_{n,k} \cap C_n)$, we get

$$\begin{aligned} \mathbb{P}(A_{n,k} \cap C_n) &\leq \mathbb{P}(n\lambda_n/n\delta_n \geq |\mu_{k+1}^* - \mu_k^*|/3) \\ &\quad + \mathbb{P}(\{|\hat{\mu}_{k+1} - \mu_{k+1}^*| \geq |\mu_{k+1}^* - \mu_k^*|/3\} \cap C_n) \\ &\quad + \mathbb{P}\left(\left\{ \left| \frac{\varepsilon(\hat{t}_k; t_k^* - 1)}{t_k^* - \hat{t}_k} \right| \geq |\mu_{k+1}^* - \mu_k^*|/3 \right\} \cap A_{n,k}\right) \\ &\stackrel{\text{def}}{=} \mathbb{P}(A_{n,k,1}) + \mathbb{P}(A_{n,k,2}) + \mathbb{P}(A_{n,k,3}). \end{aligned}$$

By Assumption (A4), $n\lambda_n/(n\delta_n J_{\min}^*) < 1/3$, for n large enough, leading to $\mathbb{P}(A_{n,k,1}) \rightarrow 0$. By Lemma 2 with $x_n = J_{\min}^*/3$, $v_n = n\delta_n$ and Assumption (A2), $\mathbb{P}(A_{n,k,3}) \rightarrow 0$. Let us now address $\mathbb{P}(A_{n,k,2})$. Using (18) in Lemma 1 with $j = (t_k^* + t_{k+1}^*)/2$ and with $j = t_k^*$, and using the triangle inequality, we get

$$\left| \sum_{i=t_k^*}^{(t_k^* + t_{k+1}^*)/2 - 1} Y_i - \sum_{i=t_k^*}^{(t_k^* + t_{k+1}^*)/2 - 1} \hat{u}_i \right| \leq n\lambda_n.$$

Since we are in the event C_n and $\hat{t}_k \leq t_k^*$, $\hat{u}_i \equiv \hat{\mu}_{k+1}$ within the interval $[t_k^*, (t_k^* + t_{k+1}^*)/2 - 1]$, which gives $|(t_{k+1}^* - t_k^*)(\mu_{k+1}^* - \hat{\mu}_{k+1})/2 + \varepsilon(t_k^*; (t_k^* + t_{k+1}^*)/2 - 1)| \leq n\lambda_n$. This implies that

$$(t_{k+1}^* - t_k^*)|\mu_{k+1}^* - \hat{\mu}_{k+1}|/2 \leq n\lambda_n + |\varepsilon(t_k^*; (t_k^* + t_{k+1}^*)/2 - 1)|.$$

Therefore, we may upper bound $\mathbb{P}(A_{n,k,2})$ as follows:

$$\begin{aligned} \mathbb{P}(\{|\hat{\mu}_{k+1} - \mu_{k+1}^*| \geq |\mu_{k+1}^* - \mu_k^*|/3\} \cap C_n) &\leq \mathbb{P}(n\lambda_n \geq (t_{k+1}^* - t_k^*)|\mu_{k+1}^* - \mu_k^*|/12) \\ &\quad + \mathbb{P}\left(\left| \frac{\varepsilon(t_k^*; (t_k^* + t_{k+1}^*)/2 - 1)}{t_{k+1}^* - t_k^*} \right| \geq |\mu_{k+1}^* - \mu_k^*|/6\right), \end{aligned}$$

which is arbitrarily small if $n\lambda_n < I_{\min}^* \cdot J_{\min}^*/12$ for n large enough, and, by Lemma 2, if $I_{\min}^* (J_{\min}^*)^2 / \log(n) \rightarrow \infty$, as n tends to infinity. The last two conditions hold thanks to Assumptions (A2), (A3), and (A4). Since the proof in the case $\hat{t}_k \geq t_k^*$ follows from similar reasoning, we have proved that $\mathbb{P}(A_{n,k} \cap C_n) \rightarrow 0$, as n tends to infinity.

We now prove that $\mathbb{P}(A_{n,k} \cap \bar{C}_n) \rightarrow 0$. Recall that by definition of C_n given in (28), $\bar{C}_n = \{\max_{k \in \{1, \dots, K^*\}} |\hat{t}_k - t_k^*| \geq I_{\min}^*/2\}$. We now split $\mathbb{P}(A_{n,k} \cap \bar{C}_n)$ into three terms:

$$\begin{aligned} \mathbb{P}(A_{n,k} \cap \bar{C}_n) &= \mathbb{P}(A_{n,k} \cap D_n^{(l)}) + \mathbb{P}(A_{n,k} \cap D_n^{(m)}) \\ &\quad + \mathbb{P}(A_{n,k} \cap D_n^{(r)}), \end{aligned}$$

where

$$\begin{aligned} D_n^{(l)} &\stackrel{\text{def}}{=} \{\text{there exists } p \in \{1, \dots, K^*\}, \hat{t}_p \leq t_{p-1}^*\} \cap \bar{C}_n, \\ D_n^{(m)} &\stackrel{\text{def}}{=} \{\text{for all } k \in \{1, \dots, K^*\}, t_{k-1}^* < \hat{t}_k < t_{k+1}^*\} \cap \bar{C}_n, \\ D_n^{(r)} &\stackrel{\text{def}}{=} \{\text{there exists } p \in \{1, \dots, K^*\}, \hat{t}_p \geq t_{p+1}^*\} \cap \bar{C}_n. \end{aligned}$$

Let us first focus on $\mathbb{P}(A_{n,k} \cap D_n^{(m)})$ and consider the case where $\hat{t}_k \leq t_k^*$, since the other case can be addressed in a similar way. Note that

$$\begin{aligned} \mathbb{P}(A_{n,k} \cap D_n^{(m)}) &\leq \mathbb{P}(A_{n,k} \cap B_{k+1,k} \cap D_n^{(m)}) \\ &\quad + \sum_{l=k+1}^{K^*} \mathbb{P}(C_{l,l} \cap B_{l+1,l} \cap D_n^{(m)}), \quad (29) \end{aligned}$$

where $B_{p,q} = \{(\hat{t}_p - t_q^*) \geq I_{\min}^*/2\}$ with the convention $B_{K^*+1,K^*} = \{(n - t_{K^*}^*) \geq I_{\min}^*/2\}$ and $C_{p,q} = \{(t_p^* - \hat{t}_q) \geq I_{\min}^*/2\}$. Let us now prove that the first term in the right-hand side of (29) tends to zero as n tends to infinity, the arguments for addressing the other terms being similar. Using (18) and (17) in Lemma 1 with $j = t_k^*$ and $\ell = k$, on the one hand and (18) in Lemma 1 with $j = t_k^*$ and (17) in Lemma 1 with $\ell = k+1$ on the other hand, we obtain, respectively:

$$|\hat{t}_k - t_k^*||\hat{\mu}_{k+1} - \mu_k^*| \leq n\lambda_n + |\varepsilon(\hat{t}_k; t_k^* - 1)| \quad \text{and} \quad (30)$$

$$|\hat{t}_{k+1} - t_k^*||\hat{\mu}_{k+1} - \mu_{k+1}^*| \leq n\lambda_n + |\varepsilon(t_k^*; \hat{t}_{k+1} - 1)|.$$

Defining E_n by

$$\begin{aligned} E_n &= \{|\mu_{k+1}^* - \mu_k^*| \leq n\lambda_n/(n\delta_n) + 2n\lambda_n/I_{\min}^* \\ &\quad + (t_k^* - \hat{t}_k)^{-1}|\varepsilon(\hat{t}_k; t_k^* - 1)| \\ &\quad + (\hat{t}_{k+1} - t_k^*)^{-1}|\varepsilon(t_k^*; \hat{t}_{k+1} - 1)|\}, \end{aligned}$$

we obtain

$$\begin{aligned} \mathbb{P}(A_{n,k} \cap B_{k+1,k} \cap D_n^{(m)}) &\leq \mathbb{P}(E_n \cap \{(t_k^* - \hat{t}_k) \geq n\delta_n\} \cap \{(\hat{t}_{k+1} - t_k^*) \geq I_{\min}^*/2\}) \\ &\leq \mathbb{P}(n\lambda_n/(n\delta_n) \geq |\mu_{k+1}^* - \mu_k^*|/4) \\ &\quad + \mathbb{P}(2n\lambda_n/I_{\min}^* \geq |\mu_{k+1}^* - \mu_k^*|/4) \\ &\quad + \mathbb{P}(\{(t_k^* - \hat{t}_k)^{-1}|\varepsilon(\hat{t}_k; t_k^* - 1)| \geq |\mu_{k+1}^* - \mu_k^*|/4\} \\ &\quad \cap \{(t_k^* - \hat{t}_k) \geq n\delta_n\}) \\ &\quad + \mathbb{P}(\{(\hat{t}_{k+1} - t_k^*)^{-1}|\varepsilon(t_k^*; \hat{t}_{k+1} - 1)| \geq |\mu_{k+1}^* - \mu_k^*|/4\} \\ &\quad \cap \{(\hat{t}_{k+1} - t_k^*) \geq I_{\min}^*/2\}). \end{aligned}$$

By Assumptions (A2), (A3), and (A4), $\mathbb{P}(A_{n,k} \cap B_{k+1,k} \cap D_n^{(m)}) \rightarrow 0$, as n tends to infinity, which concludes that $\mathbb{P}(A_{n,k} \cap D_n^{(m)}) \rightarrow 0$.

Let us now focus on $\mathbb{P}(A_{n,k} \cap D_n^{(\ell)})$. The latter probability can be upper bounded by

$$\begin{aligned} \mathbb{P}(D_n^{(\ell)}) &\leq \sum_{k=1}^{K^*} 2^{k-1} \mathbb{P}(\max\{1 \leq l \leq K^*, \hat{t}_l \leq t_{l-1}^*\} = k) \\ &\leq 2^{K^*-1} \sum_{k=1}^{K^*-1} \sum_{m \geq k} \mathbb{P}(\{t_m^* - \hat{t}_m > I_{\min}^*/2\} \\ &\quad \cap \{\hat{t}_{m+1} - t_m^* > I_{\min}^*/2\}) \\ &\quad + 2^{K^*-1} \mathbb{P}(\{t_{K^*}^* - \hat{t}_{K^*} > I_{\min}^*/2\}). \end{aligned} \quad (31)$$

Consider one term of the sum in the right-hand side of (31). Using (30) with $k = m$, we get

$$\begin{aligned} &\mathbb{P}(\{t_m^* - \hat{t}_m > I_{\min}^*/2\} \cap \{\hat{t}_{m+1} - t_m^* > I_{\min}^*/2\}) \\ &\leq \mathbb{P}(4n\lambda_n/I_{\min}^* \geq |\mu_{m+1}^* - \mu_m^*|/3) \\ &\quad + \mathbb{P}(\{(t_m^* - \hat{t}_m)^{-1} |\varepsilon(\hat{t}_m; t_m^* - 1)| \geq |\mu_{m+1}^* - \mu_m^*|/3\} \\ &\quad \cap \{(t_m^* - \hat{t}_m) \geq I_{\min}^*/2\}) \\ &\quad + \mathbb{P}(\{(\hat{t}_{m+1} - t_m^*)^{-1} |\varepsilon(t_m^*; \hat{t}_{m+1} - 1)| \geq |\mu_{m+1}^* - \mu_m^*|/3\} \\ &\quad \cap \{\hat{t}_{m+1} - t_m^* \geq I_{\min}^*/2\}). \end{aligned}$$

By Assumptions (A2), (A3), and (A4), $\mathbb{P}(\{t_m^* - \hat{t}_m > I_{\min}^*/2\} \cap \{\hat{t}_{m+1} - t_m^* > I_{\min}^*/2\}) \rightarrow 0$, as n tends to infinity. Let us now consider the last term in the right-hand side of (31). Using (30) with $k = K^*$ leads to

$$\begin{aligned} &\mathbb{P}(\{t_{K^*}^* - \hat{t}_{K^*} > I_{\min}^*/2\}) \\ &\leq \mathbb{P}(3n\lambda_n/I_{\min}^* \geq |\mu_{K^*+1}^* - \mu_{K^*}^*|/3) \\ &\quad + \mathbb{P}(\{(t_{K^*}^* - \hat{t}_{K^*})^{-1} |\varepsilon(\hat{t}_{K^*}; t_{K^*}^* - 1)| \geq |\mu_{K^*+1}^* - \mu_{K^*}^*|/3\} \\ &\quad \cap \{(t_{K^*}^* - \hat{t}_{K^*}) \geq I_{\min}^*/2\}) \\ &\quad + \mathbb{P}(\{(n - t_{K^*}^* + 1)^{-1} |\varepsilon(t_{K^*}^*; n)| \geq |\mu_{K^*+1}^* - \mu_{K^*}^*|/3\}). \end{aligned}$$

By Assumptions (A2), (A3), and (A4), $\mathbb{P}(\{t_{K^*}^* - \hat{t}_{K^*} > I_{\min}^*/2\}) \rightarrow 0$, as n tends to infinity, which gives $\mathbb{P}(D_n^{(\ell)}) \rightarrow 0$. In a similar way, we can prove that $\mathbb{P}(D_n^{(r)}) \rightarrow 0$, as n tends to infinity which gives that $\mathbb{P}(A_{n,k} \cap \bar{C}_n) \rightarrow 0$ and concludes the proof.

Proof of Proposition 4. In this proof, we shall use the notation introduced in (27). By lemma 2 of Meinshausen and Yu (2009), we get that with probability tending to one

$$|\hat{A}(\lambda_n)| \leq C \frac{n}{\lambda_n^2}, \quad (32)$$

where C is a positive constant equal to $\sigma^2 + K^{*2} J_{\max}^2$. In order to prove that

$$\begin{aligned} &\mathbb{P}(\{\mathcal{E}(\hat{T}_{n,|\hat{A}(\lambda_n)|} \parallel \mathcal{T}_n^*) \geq n\delta_n\} \cap \{|\hat{A}(\lambda_n)| \geq K^*\}) \\ &\rightarrow 0, \quad \text{as } n \rightarrow \infty, \end{aligned}$$

it is enough to prove that

$$\begin{aligned} &\mathbb{P}(\{\mathcal{E}(\hat{T}_{n,|\hat{A}(\lambda_n)|} \parallel \mathcal{T}_n^*) \geq n\delta_n\} \cap \{K^* \leq |\hat{A}(\lambda_n)| \leq Cn/\lambda_n^2\}) \\ &\rightarrow 0, \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Note that

$$\begin{aligned} &\mathbb{P}(\{\mathcal{E}(\hat{T}_{n,|\hat{A}(\lambda_n)|} \parallel \mathcal{T}_n^*)\} \cap \{K^* \leq |\hat{A}(\lambda_n)| \leq Cn/\lambda_n^2\}) \\ &\leq \mathbb{P}(\mathcal{E}(\hat{T}_{n,K^*} \parallel \mathcal{T}_n^*) \geq n\delta_n) \\ &\quad + \sum_{K > K^*}^{Cn/\lambda_n^2} \mathbb{P}(\mathcal{E}(\hat{T}_{n,K} \parallel \mathcal{T}_n^*) \geq n\delta_n). \end{aligned} \quad (33)$$

The first term of the right-hand side of (33) tends to zero as $n \rightarrow \infty$ since it is upper bounded by $\mathbb{P}(\max_{1 \leq k \leq K^*} |\hat{t}_k - t_k^*| > n\delta_n)$ which tends to zero by Proposition 3. Let us now focus on the second term on the right-hand side of (33). Note that

$$\begin{aligned} &\sum_{K > K^*}^{Cn/\lambda_n^2} \mathbb{P}(\mathcal{E}(\hat{T}_{n,K} \parallel \mathcal{T}_n^*) \geq n\delta_n) \\ &\leq \sum_{K > K^*} \sum_{k=1}^{K^*} \mathbb{P}(\forall 1 \leq l \leq K, |\hat{t}_l - t_k^*| \geq n\delta_n) \\ &\stackrel{\text{def}}{=} \sum_{K > K^*} \sum_{k=1}^{K^*} \mathbb{P}(E_{n,k,1}) + \mathbb{P}(E_{n,k,2}) + \mathbb{P}(E_{n,k,3}), \end{aligned}$$

where

$$\begin{aligned} E_{n,k,1} &= \{\forall 1 \leq l \leq K, |\hat{t}_l - t_k^*| \geq n\delta_n \text{ and } \hat{t}_l < t_k^*, \\ E_{n,k,2} &= \{\forall 1 \leq l \leq K, |\hat{t}_l - t_k^*| \geq n\delta_n \text{ and } \hat{t}_l > t_k^*, \\ E_{n,k,3} &= \{\exists 1 \leq l \leq K-1, |\hat{t}_l - t_k^*| \geq n\delta_n, \\ &\quad |\hat{t}_{l+1} - t_k^*| \geq n\delta_n, \text{ and } \hat{t}_l < t_k^* < \hat{t}_{l+1}\}. \end{aligned}$$

Let us first upper bound $\mathbb{P}(E_{n,k,1})$. Remark that

$$\mathbb{P}(E_{n,k,1}) = \mathbb{P}(E_{n,k,1} \cap \{\hat{t}_K > t_{k-1}^*\}) + \mathbb{P}(E_{n,k,1} \cap \{\hat{t}_K \leq t_{k-1}^*\}).$$

Applying (18) in Lemma 1 with $j = t_k^*$ and (17) in Lemma 1 with $\ell = K$ in the case where $\hat{t}_K > t_{k-1}^*$ gives with probability one

$$\begin{aligned} &|(t_k^* - \hat{t}_K)(\mu_k^* - \mu_{k+1}^*) + (\mu_{k+1}^* - \hat{\mu}_{K+1})| \\ &\quad + \varepsilon(\hat{t}_K; t_k^* - 1) \leq n\lambda_n. \end{aligned}$$

Thus,

$$\begin{aligned} &\mathbb{P}(E_{n,k,1} \cap \{\hat{t}_K > t_{k-1}^*\}) \\ &\leq \mathbb{P}(n\lambda_n/(n\delta_n) \geq |\mu_k^* - \mu_{k+1}^*|/3) \\ &\quad + \mathbb{P}(|\mu_{k+1}^* - \hat{\mu}_{K+1}| \geq |\mu_k^* - \mu_{k+1}^*|/3) \\ &\quad + \mathbb{P}(\{(t_k^* - \hat{t}_K)^{-1} \varepsilon(\hat{t}_K; t_k^* - 1)| \geq |\mu_k^* - \mu_{k+1}^*|/3\} \\ &\quad \cap \{|\hat{t}_K - \hat{t}_K| \geq n\delta_n\}) \\ &\stackrel{\text{def}}{=} \mathbb{P}(E_{n,k,1}^{(1)}) + \mathbb{P}(E_{n,k,1}^{(2)}) + \mathbb{P}(E_{n,k,1}^{(3)}). \end{aligned}$$

By Assumption (A4), $n\lambda_n/(n\delta_n J_{\min}^*) < 1/3$, for n large enough, leading to $CnK^*/\lambda_n^2 \mathbb{P}(E_{n,k,1}^{(1)}) \rightarrow 0$. By Lemma 2

with $x_n = J_{\min}^*/3$, $v_n = n\delta_n$ and using that $n\delta_n J_{\min}^*/\log(n^3/\lambda_n^2) \rightarrow \infty$, $CnK^*/\lambda_n^2 \mathbb{P}(E_{n,k,1}^{(3)}) \rightarrow 0$. Let us now address $\mathbb{P}(E_{n,k,1}^{(2)})$. Using (18) in Lemma 1 with $j = t_k^*$ and with $j = t_{k+1}^*$, we get

$$(t_{k+1}^* - t_k^*)|\mu_{k+1}^* - \hat{\mu}_{K+1}| \leq n\lambda_n + |\varepsilon(t_k^*; t_{k+1}^* - 1)|.$$

Therefore, we may upper bound $\mathbb{P}(E_{n,k,1}^{(2)})$ as follows:

$$\begin{aligned} \mathbb{P}(|\mu_{k+1}^* - \hat{\mu}_{K+1}| \geq |\mu_k^* - \mu_{k+1}^*|/3) \\ \leq \mathbb{P}(n\lambda_n \geq (t_{k+1}^* - t_k^*)|\mu_k^* - \mu_{k+1}^*|/6) \\ + \mathbb{P}(|(t_{k+1}^* - t_k^*)^{-1}\varepsilon(t_k^*; t_{k+1}^* - 1)| \geq |\mu_k^* - \mu_{k+1}^*|/6). \end{aligned}$$

By using Assumptions (A2), (A3), and $n\delta_n J_{\min}^*/\log(n^3/\lambda_n^2) \rightarrow \infty$, we conclude as previously that $CnK^*/\lambda_n^2 \mathbb{P}(E_{n,k,1}^{(2)}) \rightarrow 0$. The same arguments can be used for addressing $\mathbb{P}(E_{n,k,1} \cap \{\hat{t}_K \leq t_{k-1}^*\})$. We can address in the same way the term $\mathbb{P}(E_{n,k,2})$.

Let us now focus on $\mathbb{P}(E_{n,k,3})$. Note that $\mathbb{P}(E_{n,k,3})$ can be split into four terms as follows:

$$\mathbb{P}(E_{n,k,3}) = \mathbb{P}(E_{n,k,3}^{(1)}) + \mathbb{P}(E_{n,k,3}^{(2)}) + \mathbb{P}(E_{n,k,3}^{(3)}) + \mathbb{P}(E_{n,k,3}^{(4)}),$$

where

$$\begin{aligned} E_{n,k,3}^{(1)} &= E_{n,k,3} \cap \{t_{k-1}^* < \hat{t}_l < \hat{t}_{l+1} < t_{k+1}^*\}, \\ E_{n,k,3}^{(2)} &= E_{n,k,3} \cap \{t_{k-1}^* < \hat{t}_l < t_{k+1}^*, \hat{t}_{l+1} \geq t_{k+1}^*\}, \\ E_{n,k,3}^{(3)} &= E_{n,k,3} \cap \{\hat{t}_l \leq t_{k-1}^*, t_{k-1}^* < \hat{t}_{l+1} < t_{k+1}^*\}, \\ E_{n,k,3}^{(4)} &= E_{n,k,3} \cap \{\hat{t}_l \leq t_{k-1}^*, t_{k+1}^* \leq \hat{t}_{l+1}\}. \end{aligned}$$

As for addressing $\mathbb{P}(E_{n,k,1} \cap \{\hat{t}_K > t_{k-1}^*\})$, we have to use twice Lemma 1. For $\mathbb{P}(E_{n,k,3}^{(1)})$, we first use (18) and (17) in Lemma 1 with $j = t_k^*$ and $\ell = l$, respectively. Second, we use (18) and (17) in Lemma 1 with $j = t_k^*$ and $\ell = l+1$, respectively. For $\mathbb{P}(E_{n,k,3}^{(2)})$, we first use Lemma 1 with $j = t_k^*$ and $\ell = l$. Second, we use Lemma 1 with $j = t_k^*$ and $j = t_{k+1}^*$. For $\mathbb{P}(E_{n,k,3}^{(3)})$, we first use Lemma 1 with $j = t_{k-1}^*$ and $j = t_k^*$. Second, we use Lemma 1 with $j = t_k^*$ and $\ell = l+1$. Finally, for $\mathbb{P}(E_{n,k,3}^{(4)})$, we first use Lemma 1 with $j = t_{k-1}^*$ and $j = t_k^*$. Second, we use Lemma 1 with $j = t_k^*$ and $j = t_{k+1}^*$.

APPENDIX

Discussion About Condition (15)

Let us compute the different matrices arising in (15). The matrix C_{AA}^n is a $K^* \times K^*$ matrix defined by

$$nC_{AA}^n = \begin{pmatrix} n-t_1^*+1 & n-t_2^*+1 & n-t_3^*+1 & \cdots & n-t_{K^*}^*+1 \\ n-t_2^*+1 & n-t_2^*+1 & n-t_3^*+1 & \cdots & n-t_{K^*}^*+1 \\ n-t_3^*+1 & n-t_3^*+1 & n-t_3^*+1 & \cdots & n-t_{K^*}^*+1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n-t_{K^*}^*+1 & n-t_{K^*}^*+1 & n-t_{K^*}^*+1 & \cdots & n-t_{K^*}^*+1 \end{pmatrix}. \quad (\text{A.1})$$

As for $(C_{AA}^n)^{-1}$, it is a $K^* \times K^*$ symmetric tridiagonal matrix satisfying

$$n^{-1}(C_{AA}^n)^{-1} = \begin{pmatrix} d_{2,1} & -d_{2,1} & 0 & 0 & \cdots \\ -d_{2,1} & d_{2,1} + d_{3,2} & -d_{3,2} & 0 & \cdots \\ 0 & -d_{3,2} & d_{3,2} + d_{4,3} & -d_{4,3} & 0 \\ 0 & 0 & \ddots & \ddots & \ddots \\ 0 & 0 & \cdots & 0 & 0 \\ \cdots & 0 & & & \\ \cdots & 0 & & & \\ \cdots & 0 & & & \\ & \vdots & & & \\ -d_{K^*,K^*-1} & d_{K^*,K^*-1} + d_{K^*,K^*-1} & & & \end{pmatrix}, \quad (\text{A.2})$$

where $d_{k,l} = (t_k^* - t_l^*)^{-1}$, for $1 \leq k, l \leq K^*$ and $d_{K^*,K^*-1} = (n - t_{K^*}^* + 1)^{-1}$.

Since $a_{1,1} = 1$ where $A = (a_{i,j})_{1 \leq i \leq n-K^*, 1 \leq j \leq K^*} = C_{AA}^n \times (C_{AA}^n)^{-1}$ and $a_{1,j} = 0$, for all $2 \leq j \leq K^*$, the irrerepresentable condition (15) is clearly not satisfied.

Discussion About Condition (16)

Let $\mathcal{M} = \{t_1, \dots, t_m\}$ be a set of indices of cardinal m . Using (A.2), one can see that, as soon as \mathcal{M} is such that $t_j - t_i = 1$ for all i and j such that $j - i = 1$, $n^{-1}(C_{AA}^n)^{-1}$ is a tridiagonal matrix with diagonal terms equal to 2 except the first one which is equal to 1 and extra diagonal terms equal to -1 . Such a matrix is symmetric and positive since all the determinants of its submatrices are equal to 1. Thus, the maximal eigenvalue of $(C_{AA}^n)^{-1}$ is larger than n implying that the minimal eigenvalue of C_{AA}^n is smaller than $1/n$. Hence, Condition (16) is not fulfilled.

[Received March 2009. Revised April 2010.]

REFERENCES

- Basseville, M., and Nikiforov, N. (1993), *The Detection of Abrupt Changes. Information and System Sciences Series*, Englewood Cliffs, NJ: Prentice-Hall. [1480]
- Bellman, R. (1961), "On the Approximation of Curves by Line Segments Using Dynamic Programming," *Communications of the ACM*, 4 (6), 284–294. [1480,1486]
- Bickel, P., Ritov, Y., and Tsybakov, A. (2009), "Simultaneous Analysis of Lasso and Dantzig Selector," *The Annals of Statistics*, 37 (4), 1705–1732. [1483]
- Birgé, L., and Massart, P. (2001), "Gaussian Model Selection," *Journal of the European Mathematical Society*, 3, 203–268. [1489]
- Boysen, L., Kempe, A., Munk, A., Liebscher, V., and Wittich, O. (2009), "Consistencies and Rates of Convergence of Jump Penalized Least Squares Estimators," *The Annals of Statistics*, 37 (1), 157–183. [1480,1481,1485]
- Brodsky, B., and Darkhovsky, B. (1993), *Nonparametric Methods in Change-Point Problems*, Dordrecht, The Netherlands: Kluwer Academic Publishers. [1480]
- (2000), *Non-Parametric Statistical Diagnosis: Problems and Methods*, Kluwer Academic Publishers. [1480]
- Carlstein, E., Müller, H., and Siegmund, D. (eds.) (1994), *Change-Point Problems. IMS Monograph*, Vol. 23, Hayward, CA: IMS. [1480]
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2001), *Introduction to Algorithms*, Cambridge, MA: MIT Press. [1482]
- d'Aspremont, A., Bach, F., and El Ghaoui, L. (2008), "Optimal Solutions for Sparse Principal Component Analysis," *Journal of Machine Learning Research*, 9, 1269–1294. [1481]
- Donoho, D., and I. Johnstone (1995), "Adapting to Unknown Smoothness via Wavelet Shrinkage," *Journal of the American Statistical Association*, 90 (432), 1200–1224. [1485]
- Efron, B., Hastie, T., and Tibshirani, R. (2004), "Least Angle Regression," *The Annals of Statistics*, 32, 407–499. [1481]

- Fearnhead, P. (2006), "Exact and Efficient Bayesian Inference for Multiple Change-point Problems," *Statistics and Computing*, 16, 203–213. [1480]
- Fisher, W. D. (1958), "On Grouping for Maximum Homogeneity," *Journal of the American Statistical Society*, 53, 789–798. [1480,1486]
- Gales, M., and Young, S. (2008), "The Application of Hidden Markov Models in Speech Recognition," *Foundations and Trends in Information Retrieval*, 1 (3), 195–304. [1480]
- Harchaoui, Z., and Lévy-Leduc, C. (2008), "Catching Change-Points With Lasso," in *Advances in Neural Information Processing Systems*, Vol. 20, Cambridge, MA: MIT Press. [1487]
- Hawkins, D. M. (2001), "Fitting Multiple Change-Point Models to Data," *Computational Statistics and Data Analysis*, 37, 323–341. [1480]
- Hesterberg, T. C., Choi, N. H., Meier, L., and Fraley, C. (2008), "Least Angle and ℓ_1 Penalized Regression: A Review," *Statistics Surveys*, 2, 61–93. [1481]
- Kay, S. M. (1993), *Fundamentals of Statistical Signal Processing: Detection Theory*, Prentice-Hall. [1486]
- Kolesnikov, A., and Fraenti, P. (2003), "Reduced-Search Dynamic Programming for Approximation of Polygonal Curves," *Pattern Recognition Letters*, 24 (14), 2243–2254. [1480]
- Lavielle, M. (2005), "Using Penalized Contrasts for the Change-Points Problems," *Signal Processing*, 85 (8), 1501–1510. [1480]
- Lavielle, M., and Moulines, E. (2000), "Least-Squares Estimation of an Unknown Number of Shifts in a Time Series," *Journal of Time Series Analysis*, 21 (1), 33–59. [1480,1484,1485]
- Lebarbier, E. (2005), "Detecting Multiple Change-Points in the Mean of a Gaussian Process by Model Selection," *Signal Processing*, 85 (4), 717–736. [1480]
- Lévy-Leduc, C., and Roueff, F. (2009), "Detection and Localization of Change-Points in High-Dimensional Network Traffic Data," *The Annals of Applied Statistics*, 3 (2), 637–662. [1480]
- Mammen, E., and Van De Geer, S. (1997), "Locally Adaptive Regression Splines," *The Annals of Statistics*, 25 (1), 387–413. [1481,1483]
- Massart, P. (2005), "A Non Asymptotic Theory for Model Selection," in *4th European Congress of Mathematics*, European Mathematical Society, pp. 309–323. [1480]
- Meinshausen, N., and Yu, B. (2009), "Lasso-Type Recovery of Sparse Representations for High-Dimensional Data," *The Annals of Statistics*, 37 (1), 246–270. [1483,1484,1491]
- Moghaddam, B., Weiss, Y., and Avidan, S. (2006), "Generalized Spectral Bounds for Sparse LDA," in *Proceedings of the 23rd International Conference on Machine Learning (ICML 2006)*, New York: ACM, pp. 641–648. [1481]
- Ruanaidh, J., and Fitzgerald, W. (1996), *Numerical Bayesian Methods Applied to Signal Processing. Statistics and Computing*, New York: Springer. [1480]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society, Ser. B*, 58 (1), 267–288. [1481]
- Tibshirani, R., and Wang, P. (2008), "Spatial Smoothing and Hot Spot Detection for CGH Data Using the Fused Lasso," *Biostatistics*, 9 (1), 18–29. [1481]
- Yao, Y., and Au, S. T. (1989), "Least-Squares Estimation of a Step Function," *Sankhyā: The Indian Journal of Statistics, Ser. A*, 51 (3), 370–381. [1480,1481,1483,1484,1486]
- Zhao, P., and B. Yu (2006), "On Model Selection Consistency of Lasso," *Journal of Machine Learning Research*, 7, 2541–2563. [1483]