

# ESTIMATOR SELECTION FOR REGRESSION FUNCTIONS IN EXPONENTIAL FAMILIES WITH APPLICATION TO CHANGEPOINT DETECTION

JUNTONG CHEN

ABSTRACT. We observe  $n$  independent pairs of random variables  $(W_i, Y_i)$  for which the conditional distribution of  $Y_i$  given  $W_i = w_i$  belongs to a one-parameter exponential family with parameter  $\gamma^*(w_i) \in \mathbb{R}$  and our aim is to estimate the regression function  $\gamma^*$ . Our estimation strategy is as follows. We start with an arbitrary collection of piecewise constant candidate estimators based on our observations and by means of the same observations, we select an estimator among the collection. Our approach is agnostic to the dependencies of the candidate estimators with respect to the data and can therefore be unknown. From this point of view, our procedure contrasts with other alternative selection methods based on data splitting, cross validation, hold-out etc. To illustrate its theoretical performance, we establish a non-asymptotic risk bound for the selected estimator. We then explain how to apply our procedure to the changepoint detection problem in exponential families. The practical performance of the proposed algorithm is illustrated by a comparative simulation study under different scenarios and on two real datasets from the copy numbers of DNA and British coal disasters records.

## 1. INTRODUCTION

We observe  $n$  pairs of independent (but not necessarily i.i.d.) random variables, i.e.  $X_i = (W_i, Y_i)$  for  $i = 1, \dots, n$ , with values in a measurable product space  $(\mathcal{W} \times \mathcal{Y}, \mathcal{W} \otimes \mathcal{Y})$ . For each  $i$ , we assume that the conditional distribution of  $Y_i$  given  $W_i = w_i$  exists which we denote as  $R_i^*(w_i)$  and  $R_i^*(w_i)$  belongs to a one-parameter exponential family with parameter  $\gamma^*(w_i) \in \mathbb{R}$ . The aim of the present paper is to estimate the  $n$  conditional distributions  $R_i^*(w_i)$  of  $Y_i$  given  $W_i = w_i$ , i.e. to estimate the unknown function  $\gamma^*$  on  $\mathcal{W}$ , on the basis of the observations  $\mathbf{X} = (X_1, \dots, X_n)$ .

---

*Date:* January 19, 2023.

*2020 Mathematics Subject Classification.* Primary 62G05, 62G35; Secondary 62P10.

*Key words and phrases.* Exponential families, estimator selection, model selection, changepoint detection.

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement N° 811017.

Under the above statistical setting, we are only aware of a small amount of related papers tackling this estimation problem and establishing risk bounds for the proposed estimators. Based on one model, [Baraud and Chen \(2020\)](#) proposed an estimation procedure, where their idea comes from the  $\rho$ -estimation ([Baraud et al. \(2017\)](#) and [Baraud and Birgé \(2018\)](#)). The risk bound for their estimator  $\hat{\gamma}$  can be written as, up to a constant, the sum of an approximation term and a complexity term of the model adopted for estimation. Such an approach performs well if one knows a suitable model for the function  $\gamma^*$  in advance which means a model can provide a good enough approximation of  $\gamma^*$  and is also not too complicated. But such a model can be difficult to design in some situations where only few prior information is known. A safer strategy then considered by [Chen \(2022\)](#) in which the problem was solved by a model selection procedure. However, one defect of this strategy is the expensive numerical cost especially when the number of the models becomes large. At this point, an interesting problem could be can we come up with a new estimation strategy overcoming the above two limitations in the literature? This is to say the desired strategy should be capable of comparing estimators from several different models with a reasonable numerical cost.

In particular, when  $W_i = (i - 1)/n$  (or  $i/n$  in some literature) are deterministic for all  $i \in \{1, \dots, n\}$  and  $\gamma^*$  is an unknown function on  $[0, 1)$  (or  $(0, 1]$  respectively), more work has been done in the statistical literature. We only observe  $\mathbf{Y} = (Y_1, \dots, Y_n) \in \mathcal{Y}^n$  in this case ordered by some covariate such as time or position along a chromosome. [Antoniadis and Sapatinas \(2001\)](#) considered one-parameter natural exponential families with quadratic variance functions (i.e. the variance of the distribution is a quadratic function of its mean) which cover Gaussian, Poisson, gamma, binomial, negative binomial and generalized hyperbolic secant distributions and proposed their estimator based on wavelet shrinkage estimation. Then [Antoniadis et al. \(2001\)](#) extended such a wavelet based methodology to the families with cubic variance functions. [Brown et al. \(2010\)](#) also focused on one-parameter natural exponential families with quadratic variance functions. When the exponential family is parametrized by its mean, they suggest to use a mean-match variance stabilizing transformation so that turn the original problem into a standard homoscedastic Gaussian regression problem.

In a more specific situation where  $W_i = (i - 1)/n$  are deterministic and  $\gamma^* : [0, 1) \rightarrow I \subset \mathbb{R}$  is a right-continuous step function with an unknown number  $N - 1$  of changepoints (i.e.  $N$  segments,  $N \geq 1$ ), the estimation problem we consider here immediately fits the setting of changepoint detection problem in exponential families. In this context, [Frick et al. \(2013\)](#) proposed a simultaneous multiscale changepoint estimator (SMUCE for short). More precisely, for each candidate estimator, [Frick et al. \(2013\)](#) designed a multiscale statistic to evaluate the maximum over the local likelihood ratio statistics on all discrete intervals such that the estimator is constant on these

intervals with some value. Then provided a threshold  $q$ , the quantity  $N$  is estimated by  $\widehat{N}(q)$  which is the number of segments of the estimators satisfying their threshold condition with the minimal segments. Finally, their estimator is the likelihood maximizer over a constrained set in which all the estimators satisfy the threshold condition with exact  $\widehat{N}(q)$  segments. [Cleyen and Lebarbier \(2014, 2017\)](#) considered partitions given by the pruned dynamic programming algorithm ([Rigaill \(2015\)](#)) and proposed a penalized log-likelihood estimator following the works of constructing the penalty function done by L. Birgé and P. Massart (see [Barron et al. \(1999\)](#) and [Birgé and Massart \(1997\)](#) for instance). They showed that their resulting estimator satisfies some oracle inequalities. One common feature of the above mentioned two methods is that they are both, more or less, based on the maximum likelihood estimation. When the data set contains a small amount of outliers, both of the two procedures infer extra changepoints to fit the outliers while identifying the true ones in the signal. For this point, we shall illustrate it in a more straightforward way in the simulation part of this paper. A natural question is can we find a procedure to enhance the stability of their estimators?

Besides the two procedures specially designed for the changepoint detection problem in exponential families, detecting changes in the characteristics of a sequence of observed random variables has a long history and experienced a renaissance in recent years boosted by a flourishing development in bioinformatics (e.g. [Olshen et al. \(2004\)](#), [Huang et al. \(2005\)](#), [Tibshirani and Wang \(2007\)](#), [Zhang and Siegmund \(2007\)](#) and [Muggeo and Adelfio \(2010\)](#)). It also has attracted attention from other fields including climatology (e.g. [Reeves et al. \(2007\)](#) and [Gallagher et al. \(2013\)](#)), financial econometrics (e.g. [Spokoiny \(2009\)](#)) and signal processing (e.g. [Blythe et al. \(2012\)](#) and [Hotz et al. \(2013\)](#)), among many others. Within the regime of univariate mean changepoint detection, theoretical analysis has been established recently by [Verzelen et al. \(2020\)](#) and [Wang et al. \(2020\)](#). A recently selective review of the related literature can be found in [Truong et al. \(2020\)](#). We only mention some representative procedures here. [Scott and Knott \(1974\)](#) proposed a binary segmentation (BS for short) method to detect the changes in means. A modified procedure circular binary segmentation (CBS for short) was provided by [Olshen et al. \(2004\)](#) then a faster algorithm was given in [Venkatraman and Olshen \(2007\)](#) which has achieved a big success in genome analysis. Later, to enhance the robustness to departures from standard model assumptions, another method (denoted as cumSeg in the sequel) had been tailor-made by [Muggeo and Adelfio \(2010\)](#) to detect changes in genomic sequences. In the direction of reducing the complexity of computation, the pruned exact linear time (PELT for short) method was proposed by [Killick et al. \(2012\)](#) where they also showed PELT leads to a substantially more accurate result than BS. Wild binary segmentation (WBS for short) is an approach proposed by [Fryzlewicz \(2014\)](#) based on a development of

BS and it becomes quite popular nowadays due to its nice performance and an easy implementation. Aimed at improving SMUCE (Frick et al. (2013)) especially under the situation with low signal-to-noise ratio or with many changepoints compared to the length of observations, Li et al. (2016) proposed an alternative multiscale segmentation method (denoted as FDR in the sequel) by controlling the false discovery rate of their whole segmentation procedure. In the direction of being robust in the presence of outliers, Fearnhead and Rigaiil (2019) proposed an algorithm (denoted as robseg in the sequel) based on the idea of adapting existing penalized cost methods to some loss functions which are less sensitive to the outliers. Two examples of the loss functions to which their procedure applies are the Huber loss and biweight loss. In practice, based on the same observations, different approaches mentioned above may give different estimators. As it was pointed out by the comparison study in Fearnhead and Rigaiil (2020), it is rather rare that one particular method uniformly outperforms another. Given so many experts' suggestions, a realistic and also interesting question is which one we should pick? Or in another word, can we let the data decide the preference of several (possibly random) estimators case by case so that finally we can always achieve a nearly optimal performance among the candidates taken into consideration?

These three problems mentioned above are the main motivations to propose this paper. In fact, we shall see that all of them can be solved simultaneously by an estimation strategy based on a data-driven estimator selection (denoted as ES in the sequel). More precisely, given the observations  $\mathbf{X} = (X_1, \dots, X_n)$ , we assume to have at disposal an arbitrary but at most countable collection of piecewise constant (possibly random) candidates for  $\gamma^*$  which we denote as  $\hat{\Gamma}(\mathbf{X})$ . The dependency of each candidate in  $\hat{\Gamma}$  on the observations  $\mathbf{X}$  can be unknown. We design an algorithm to compare these candidates in  $\hat{\Gamma}$  pair by pair based on the same observations  $\mathbf{X}$  and let the data choose the desired one. A non-asymptotic risk bound for the selected estimator is established, where we compare it with the infimum of the risks over the collection  $\hat{\Gamma}(\mathbf{X})$ .

The paper is organized as follows. We give a specific description of the statistical framework in Section 2. Our estimator selection procedure and the theoretical properties of the resulting estimator are presented in Section 3. In Section 4, we explain how to apply this procedure to the changepoint detection problem in exponential families. Section 5 is devoted to a comparative simulation study for illustrating the practical performance of the selected estimator. Its performance on two real datasets (DNA copy numbers and British coal disasters) is exhibited in Section 6. Finally, all the proofs in this paper are left to Section 7, and details of the test signals used in Section 5 are provided in the Appendix.

## 2. THE STATISTICAL SETTING

We observe  $n$  pairs of independent (but not necessary i.i.d.) random variables  $X_i = (W_i, Y_i)$ , for  $i \in \{1, \dots, n\}$  with values in a measurable product space  $(\mathcal{X}, \mathcal{X}) = (\mathcal{W} \times \mathcal{Y}, \mathcal{W} \otimes \mathcal{Y})$ . We denote  $\mathcal{R}$  the set of all probabilities on  $(\mathcal{Y}, \mathcal{Y})$  and equip it with the Borel  $\sigma$ -algebra  $\mathcal{R}$  associated to the Hellinger distance (which induces the same topology as the total variation distance one). Recall that the Hellinger distance between two probabilities  $P = p \cdot \mu$  and  $Q = q \cdot \mu$  dominated by a reference measure  $\mu$  on a measurable space  $(A, \mathcal{A})$  is given by the formula

$$(1) \quad h(P, Q) = \left[ \frac{1}{2} \int_A (\sqrt{p} - \sqrt{q})^2 d\mu \right]^{1/2},$$

which is independent with respect to the choice of the dominated reference measure  $\mu$ . For each  $i \in \{1, \dots, n\}$ , we assume the conditional distribution of  $Y_i$  given  $W_i = w_i$  exists and can be written as  $R_i^*(w_i)$  with  $R_i^*$  a measurable function from  $(\mathcal{W}, \mathcal{W})$  to  $\mathcal{R}$ . Let us remark that with the choice of  $\mathcal{R}$ , for any  $R \in \mathcal{R}$  and all  $i \in \{1, \dots, n\}$ , the mapping  $w \mapsto h^2(R, R_i^*(w))$  on  $(\mathcal{W}, \mathcal{W})$  is measurable. Before introducing our statistical setting, we also recall the following definition.

**Definition 1.** *Let  $I$  be a non-trivial interval of  $\mathbb{R}$  (i.e.  $I \neq \emptyset$ ). We call a family of probabilities  $\mathcal{Q}_0 = \{R_\gamma, \gamma \in I\}$  on the measured space  $(\mathcal{Y}, \mathcal{Y}, \mu)$  an exponential family under its general form, if  $\mathcal{Q}_0$  is a family of probabilities on  $(\mathcal{Y}, \mathcal{Y})$  admitting densities  $\bar{r}_\gamma$  with respect to  $\mu$  of the form, for all  $\gamma \in I$*

$$\bar{r}_\gamma(y) = e^{u(\gamma)T(y) - A(\gamma)} h(y), \quad \text{for all } y \in \mathcal{Y},$$

where  $T$  is a real-valued measurable function on  $(\mathcal{Y}, \mathcal{Y})$  which does not coincide with a constant  $\nu = h \cdot \mu$ -a.e.,  $u$  is a continuous, strictly monotone function on  $I$ ,  $h$  is a nonnegative function on  $\mathcal{Y}$  and

$$A(\gamma) = \log \left[ \int_{\mathcal{Y}} e^{u(\gamma)T(y)} h(y) d\mu(y) \right].$$

Definition 1 covers many interesting and useful distributions including Gaussian (with a known variance), Poisson, binomial (with a known number of trials, e.g. Bernoulli) and gamma (with a fixed shape parameter, e.g. exponential) distributions, among many others. In the sequel, for simplicity, we rewrite  $\mathcal{Q}_0$  as  $\{R_\gamma = r_\gamma \cdot \nu, \gamma \in I\}$ , with the notation

$$(2) \quad r_\gamma(y) = e^{u(\gamma)T(y) - A(\gamma)}, \quad \text{for all } y \in \mathcal{Y} \text{ and } \gamma \in I.$$

The statistical framework we consider here is more general than the one we have mentioned in the introduction. More precisely, based on the observations  $\mathbf{X} = (X_1, \dots, X_n)$ , we would like to estimate the  $n$  conditional distributions  $R_i^*(w_i)$  of  $Y_i$  given  $W_i = w_i$  for  $i \in \{1, \dots, n\}$ . To do so, we presume (even this may not be true) that there exists a piecewise constant

function  $\gamma^*$  on  $\mathcal{W}$  such that for all  $i \in \{1, \dots, n\}$ , the conditional distributions  $R_i^*(w_i)$  of  $Y_i$  given  $W_i = w_i$  are of the form  $R_{\gamma^*(w_i)}$  or at least not far away from it with respect to the Hellinger distance defined by (1). Such a statistical setting includes the following situations:

- (i) the ideal case, where the conditional distributions  $R_i^*(w_i) = R_{\gamma^*(w_i)}$  for all  $i \in \{1, \dots, n\}$ . We then refer to this  $\gamma^*$  as the *regression function* which is a natural generalization from Gaussian regression.
- (ii) the model is slightly misspecified. This includes the situation where for some  $i \in \{1, \dots, n\}$ , the conditional distributions  $R_i^*(w_i)$  are slightly different from the presumed ones  $R_{\gamma^*(w_i)}$  or the situation where the data set contains a small amount of outliers.

Let  $\mathcal{R}_{\mathcal{W}}$  be the collection of all measurable mappings from  $(\mathcal{W}, \mathcal{W})$  into  $(\mathcal{R}, \mathcal{R})$  and set  $\mathcal{R}_{\mathcal{W}} = \mathcal{R}_{\mathcal{W}}^n$ . We denote  $\mathbf{R}^*$  the  $n$ -tuple  $(R_1^*, \dots, R_n^*)$  and hence  $\mathbf{R}^* \in \mathcal{R}_{\mathcal{W}}$ . Our goal is to estimate this  $\mathbf{R}^*$  based on the observations  $\mathbf{X} = (X_1, \dots, X_n)$ . Our estimation strategy is as follows. We suppose that we have at disposal an arbitrary but at most countable collection of (possibly random) piecewise constant candidates of  $\gamma^*$  mapping  $\mathcal{W}$  into  $I$  written as  $\hat{\Gamma} = \{\hat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$ . We then design a procedure based on the same data  $\mathbf{X}$  to select among  $\hat{\Gamma}$  and denote the selected one as  $\hat{\gamma}_{\hat{\lambda}}(\mathbf{X})$  (or  $\hat{\gamma}_{\hat{\lambda}}$  for short). Once obtaining  $\hat{\gamma}_{\hat{\lambda}}$ , our estimator of  $\mathbf{R}^*$  is given by the mapping  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}} : \mathbf{w} = (w_1, \dots, w_n) \in \mathcal{W}^n \mapsto (R_{\hat{\gamma}_{\hat{\lambda}}(w_1)}, \dots, R_{\hat{\gamma}_{\hat{\lambda}}(w_n)})$  taking values in  $\mathcal{Q}_0^n$  with  $R_{\hat{\gamma}_{\hat{\lambda}}} \in \mathcal{R}_{\mathcal{W}}$ . With a slight abuse of language, sometimes in this paper we also call  $\hat{\gamma}_{\hat{\lambda}}$  an estimator of  $\gamma^*$  though we know that such a regression function  $\gamma^*$  does not necessarily exist. It is worth emphasizing that besides the independence, we assume nothing about the distributions of the covariates  $W_i$  which therefore can be unknown.

To evaluate the performance of the selected estimator  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}}$ , we need to introduce a loss function. Since we focus on estimating the  $n$  conditional distributions, it is natural to consider a loss function based on the Hellinger distance. More precisely, we endow the space  $\mathcal{R}_{\mathcal{W}}$  with a pseudo Hellinger distance  $\mathbf{h}$  defined for any  $\mathbf{R} = (R_1, \dots, R_n)$  and  $\mathbf{R}' = (R'_1, \dots, R'_n)$  in  $\mathcal{R}_{\mathcal{W}}$  by

$$(3) \quad \mathbf{h}^2(\mathbf{R}, \mathbf{R}') = \mathbb{E} \left[ \sum_{i=1}^n h^2(R_i(W_i), R'_i(W_i)) \right] \\ = \sum_{i=1}^n \int_{\mathcal{W}} h^2(R_i(w), R'_i(w)) dP_{W_i}(w),$$

where  $h$  is the Hellinger distance introduced in (1). Whenever the regression function  $\gamma^*$  exists, we automatically deduce a performance of  $\hat{\gamma}_{\hat{\lambda}}$  with respect to  $\gamma^*$  by the distance  $d(\gamma^*, \hat{\gamma}_{\hat{\lambda}}) = \mathbf{h}(\mathbf{R}_{\gamma^*}, \mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}})$ . In particular, in the context of changepoint detection problem in exponential families where  $W_i$

are deterministic, the loss function  $\mathbf{h}$  is the sum of the Hellinger distances between each two probabilities  $R_i$  and  $R'_i$ . Such a loss function has been considered in several literature, for instance [Le Cam \(1986\)](#) and [Le Cam and Yang \(1990\)](#). From this point of view, unlike the typical methods detecting changes for some parameter of a distribution (for example detecting changes in means for Gaussian and Poisson distributions), our approach validates the changes along the sequence if there are abrupt variations with respect to the distribution level.

### 3. A STRATEGY BASED ON ESTIMATOR SELECTION

As already mentioned, given the observations  $\mathbf{X} = (X_1, \dots, X_n)$ , we assume that we have at disposal an arbitrary but at most countable (possibly random) candidates  $\hat{\Gamma} = \{\hat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$  for  $\gamma^*$ , where for each  $\lambda \in \Lambda$ ,  $\hat{\gamma}_\lambda$  is piecewise constant on  $\mathscr{W}$ . This  $\hat{\Gamma}$  may contain the estimators based on the minimization of some criterions, estimators based on Bayes procedures or just simple guesses by some experts. The dependency of these estimators with respect to the observations  $\mathbf{X}$  can be unknown. Our goal is to select some  $\hat{\gamma}_{\hat{\lambda}}(\mathbf{X})$  among the family  $\hat{\Gamma} = \{\hat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$  based on the same observations  $\mathbf{X}$  such that the risk of our estimator is as close as possible to the quantity  $\inf_{\lambda \in \Lambda} \mathbb{E} [\mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda})]$ .

**3.1. Estimator selection procedure.** Let  $\mathcal{M}$  be a finite or countable set of partitions on  $\mathscr{W}$ . We begin with a family of collections  $\{\Gamma_m, m \in \mathcal{M}\}$  indexed by the partition  $m$  on  $\mathscr{W}$ , where for each  $m \in \mathcal{M}$ ,  $\Gamma_m$  stands for an at most countable collection of piecewise constant functions on  $\mathscr{W}$  with values in  $I$  based on the partition  $m$ . Setting the notation  $\tilde{\Gamma} = \cup_{m \in \mathcal{M}} \Gamma_m$ , we assume the family of (possibly random) candidates  $\hat{\Gamma} = \{\hat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$  for  $\gamma^*$  (may not exist) with values in  $\tilde{\Gamma}$ . This is to say, for each  $\lambda \in \Lambda$ , there is a (possibly random) partition  $\hat{m}(\lambda) \in \mathcal{M}$  such that  $\hat{\gamma}_\lambda \in \Gamma_{\hat{m}(\lambda)}$ . For any  $\gamma \in \tilde{\Gamma}$ , we define

$$\mathcal{M}(\gamma) = \{m \in \mathcal{M}, \gamma \in \Gamma_m\},$$

therefore naturally we have  $\hat{m}(\lambda) \in \mathcal{M}(\hat{\gamma}_\lambda)$ .

Let  $\Delta(\cdot)$  be a map from  $\mathcal{M}$  to  $\mathbb{R}_+ = [0, +\infty)$ . For each  $m \in \mathcal{M}$ , we associate it with a nonnegative weight  $\Delta(m)$  and assume the following holds true.

**Assumption 1.** *There exists a positive number  $\Sigma$  such that*

$$(4) \quad \Sigma = \sum_{m \in \mathcal{M}} e^{-\Delta(m)} < +\infty.$$

We remark that when  $\Sigma = 1$ , the weights  $\Delta(m)$  define a prior distribution on the collection of partitions  $\mathcal{M}$ , which gives a Bayesian flavour to our selection procedure.

Given two partitions  $m_1, m_2 \in \mathcal{M}$ , we define a refined partition  $m_1 \vee m_2$  on  $\mathscr{W}$  generated by  $m_1, m_2$  as

$$m_1 \vee m_2 = \{K_1 \cap K_2 \mid K_1 \in m_1, K_2 \in m_2, K_1 \cap K_2 \neq \emptyset\}.$$

For any partition  $m$  on  $\mathscr{W}$ , we denote the number of its segments by  $|m|$ . To define our selection procedure, we also make the following assumption on the family  $\mathcal{M}$ .

**Assumption 2.** *There exists some constant  $\alpha \geq 1$  such that  $|m_1 \vee m_2| \leq \alpha(|m_1| + |m_2|)$ , for all  $m_1, m_2 \in \mathcal{M}$ .*

We give some examples of the family  $\mathcal{M}$  here satisfying Assumption 2. When  $\mathscr{W}$  is either  $\mathbb{R}$  or some subinterval of  $\mathbb{R}$ , for any finite or countable family  $\mathcal{M}$  of partitions on  $\mathscr{W}$ , it is easy to observe that Assumption 2 is satisfied with  $\alpha = 1$ . Another example can be the nested partitions, i.e. the family  $\mathcal{M}$  is ordered for the inclusion. In this situation,  $m_1 \vee m_2$  either equals to  $m_1$  or  $m_2$  so that Assumption 2 also holds true with  $\alpha = 1$ . Besides, when  $\mathscr{W} = [0, 1]^d$  with  $d \geq 2$ , a specific example satisfying Assumption 2 with  $\alpha = 2$  has been introduced in Example 3 of Baraud and Birgé (2009).

Our selection procedure is based on a pair-by-pair comparison of the candidates, where the selection mechanism is inspired by a series of work of  $\rho$ -estimation (Baraud et al. (2017) and Baraud and Birgé (2018)). However, unlike the above literature, we generalize the comparison device into the situation where the elements in  $\tilde{\Gamma}$  can be random.

Let us first introduce a monotone increasing function  $\psi$  from  $[0, +\infty]$  into  $[-1, 1]$  defined as

$$\psi(x) = \begin{cases} \frac{x-1}{x+1} & , \quad x \in [0, +\infty), \\ 1 & , \quad x = +\infty. \end{cases}$$

For any  $\gamma, \gamma' \in \tilde{\Gamma}$ , we define the  $\mathbf{T}$ -statistic as

$$\mathbf{T}(\mathbf{X}, \gamma, \gamma') = \sum_{i=1}^n \psi \left( \sqrt{\frac{r_{\gamma'}(W_i)(Y_i)}{r_{\gamma}(W_i)(Y_i)}} \right)$$

with the conventions  $0/0 = 1$  and  $a/0 = +\infty$  for all  $a > 0$ . Let  $D_n$  be a map from  $\mathcal{M}$  to  $\mathbb{R}_+$  defined as, for any  $m \in \mathcal{M}$ ,

$$D_n(m) = |m| \left[ 9.11 + \log_+ \left( \frac{n}{|m|} \right) \right],$$

where  $\log_+(x) = \max\{\log(x), 0\}$ . We define the penalty function from  $\tilde{\Gamma}$  to  $\mathbb{R}_+$  such that for all  $\gamma \in \tilde{\Gamma}$ ,

$$(5) \quad \mathbf{pen}(\gamma) \geq C_0 \left( 2\alpha + \frac{1}{2} \right) \inf_{m \in \mathcal{M}(\gamma)} [D_n(m) + \Delta(m)],$$



where  $C_0 > 0$  is a universal constant. For each  $\lambda \in \Lambda$ , we set

$$\mathbf{v}(\mathbf{X}, \hat{\gamma}_\lambda) = \sup_{\lambda' \in \Lambda} [\mathbf{T}(\mathbf{X}, \hat{\gamma}_\lambda, \hat{\gamma}_{\lambda'}) - \mathbf{pen}(\hat{\gamma}_{\lambda'})] + \mathbf{pen}(\hat{\gamma}_\lambda).$$

We select  $\hat{\gamma}_{\hat{\lambda}}$  as any measurable element of the random (and non-void) set

$$(6) \quad \mathcal{E}(\mathbf{X}) = \left\{ \hat{\gamma}_\lambda \in \hat{\Gamma} \text{ such that } \mathbf{v}(\mathbf{X}, \hat{\gamma}_\lambda) \leq \inf_{\lambda' \in \Lambda} \mathbf{v}(\mathbf{X}, \hat{\gamma}_{\lambda'}) + 1 \right\}.$$

The final selected estimator  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}}$  of  $\mathbf{R}^*$  is given by  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}} = (R_{\hat{\gamma}_{\hat{\lambda}}}, \dots, R_{\hat{\gamma}_{\hat{\lambda}}})$ .

We comment that the number 1 in (6) does not play any role, therefore can be substituted by any small number  $\delta > 0$ . We choose  $\delta = 1$  here just for enhancing the legibility of our results. Moreover, to improve the performance of the selected estimator  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}}$ , the choice of a  $\hat{\gamma}_{\hat{\lambda}}$  such that  $\mathbf{v}(\mathbf{X}, \hat{\gamma}_{\hat{\lambda}}) = \inf_{\lambda \in \Lambda} \mathbf{v}(\mathbf{X}, \hat{\gamma}_\lambda)$  should be preferred whenever available, which is the case when  $\hat{\Gamma}$  is a finite set.

**3.2. The performance of the selected estimator.** In this section, we establish non-asymptotic exponential inequalities of deviations between the selected estimator  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}}$  and the truth  $\mathbf{R}^*$ .

**Theorem 1.** *Under Assumption 1 and 2, whatever the conditional distributions  $\mathbf{R}^* = (R_1^*, \dots, R_n^*)$  of  $Y_i$  given  $W_i$  and the distributions of  $W_i$ , there exists a universal constant  $C_0 > 0$  such that the selected estimator  $\mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}}$  given by the procedure in Section 3.1 among a family of (possibly random) candidates  $\hat{\Gamma} = \{\hat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$  based on the observations  $\mathbf{X} = (X_1, \dots, X_n)$  satisfies for any  $\xi > 0$ , on a set of probability larger than  $1 - \Sigma^2 e^{-\xi}$*

$$(7) \quad \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_{\hat{\lambda}}}) \leq \inf_{\lambda \in \Lambda} [c_1 \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) + c_2 \mathbf{pen}(\hat{\gamma}_\lambda)] + c_3 (1.471 + \xi),$$

where  $c_1 = 91.4$ ,  $c_2 = 42.7$  and  $c_3 = 12666.9$ .

The proof of Theorem 1 is postponed to Section 7. We hereby give a short discussion of the numerical constant  $C_0$  appearing in the penalty function (5). In the proof of Theorem 1, we show that there does exist a numerical constant  $C_0 > 0$  such that for all the penalties satisfying (5), our procedure defined in Section 3.1 results in a selected estimator fulfilling the performance stated in Theorem 1. Unfortunately, this theoretical constant  $C_0$  turns out to be quite large and we do not have enough information about the smallest value of  $C_0$  which validates the non-asymptotic exponential inequalities in (7). In practice, when we implement our estimator selection procedure we regard this  $C_0$  as a tuning parameter instead of using the theoretical value. For this point, we will make it more clear in the simulation study, where it also turns out the value of  $C_0$  in theory seems to be too pessimistic.

To comment on the performance of the selected estimator further, we integrate (7) with respect to  $\xi$  and obtain the following risk bound.

**Corollary 1.** *Under Assumption 1 and 2, whatever the conditional distributions  $\mathbf{R}^* = (R_1^*, \dots, R_n^*)$  of  $Y_i$  given  $W_i$  and the distributions of  $W_i$ , there exists a universal constant  $C_0 > 0$  such that the selected estimator  $\mathbf{R}_{\hat{\gamma}_\lambda}$  given by the procedure in Section 3.1 among  $\hat{\Gamma} = \{\hat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$  satisfies*

$$\begin{aligned} \mathbb{E} \left[ \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) \right] &\leq \mathbb{E} \left[ \inf_{\lambda \in \Lambda} (c_1 \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) + c_2 \text{pen}(\hat{\gamma}_\lambda)) \right] + c_3 (\Sigma^2 + 1.471) \\ &\leq \inf_{\lambda \in \Lambda} \{ \mathbb{E} [c_1 \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) + c_2 \text{pen}(\hat{\gamma}_\lambda)] \} + c_3 (\Sigma^2 + 1.471). \end{aligned}$$

In particular, if the equality in (5) holds,

$$(8) \quad \mathbb{E} \left[ \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) \right] \leq C_{\alpha, \Sigma} \inf_{\lambda \in \Lambda} \{ \mathbb{E} [\mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda})] + \mathbb{E} [\Xi(\hat{\gamma}_\lambda)] \},$$

where for all  $\lambda \in \Lambda$ ,

$$\begin{aligned} \Xi(\hat{\gamma}_\lambda) &= \inf_{m \in \mathcal{M}(\hat{\gamma}_\lambda)} \left[ |m| \left( 9.11 + \log_+ \left( \frac{n}{|m|} \right) \right) + \Delta(m) \right] \\ &\leq |\hat{m}(\lambda)| \left[ 9.11 + \log_+ \left( \frac{n}{|\hat{m}(\lambda)|} \right) \right] + \Delta(\hat{m}(\lambda)) \end{aligned}$$

and

$$C_{\alpha, \Sigma} = \left[ c_2 C_0 \left( 2\alpha + \frac{1}{2} \right) + \frac{c_3 (\Sigma^2 + 1.471)}{9.11} \right] \vee c_1.$$

The result given in (8) compares the risk of the selected estimator  $\mathbf{R}_{\hat{\gamma}_\lambda}$  to those of  $\mathbf{R}_{\hat{\gamma}_\lambda}$  plus an additional nonnegative term  $\mathbb{E} [\Xi(\hat{\gamma}_\lambda)]$ . One nice feature of this approach implied by (8) lies in the fact that the risk bound does not depend on the cardinality of the set  $\hat{\Gamma}$ . This entails that if we enlarge the collection of our candidates by keeping  $\mathcal{M}$  unchanged (so that  $\Delta(m)$  will not change), the risk bound for the selected estimator only decreases over the larger collection of candidates. On the other hand, our procedure is based on  $\mathcal{O}(|\hat{\Gamma}|^2)$  times of pair-by-pair comparisons. Therefore, the payment for enlarging set  $\hat{\Gamma}$  is the computation time.

The risk bound (8) in Corollary 1 also accounts for the stability of our selection procedure under a slight misspecification framework. To illustrate, let us first consider the ideal situation where  $\mathbf{R}^* = \mathbf{R}_{\gamma^*} = (R_{\gamma^*}, \dots, R_{\gamma^*})$  with  $\gamma^*$  a piecewise constant function based on the partition  $m^*$  of  $\mathcal{W}$ . We denote  $\bar{\Gamma}_{m^*}$  the class of all piecewise constant functions with values in  $I \subset \mathbb{R}$  based on the partition  $m^*$  and assume for simplicity  $\hat{\Gamma} = \Gamma_{m^*}$ , where  $\Gamma_{m^*}$  stands for a dense (for the topology of the pointwise convergence) and countable subset of  $\bar{\Gamma}_{m^*}$ . Taking  $\Delta(m^*) = 0$ , we deduce from (8) that the estimator  $\mathbf{R}_{\hat{\gamma}}$  based on the selection among  $\Gamma_{m^*}$  satisfies for  $C > 0$  being a numerical constant,

$$(9) \quad \mathbb{E} \left[ \mathbf{h}^2(\mathbf{R}_{\gamma^*}, \mathbf{R}_{\hat{\gamma}}) \right] \leq C |m^*| \left[ 1 + \log_+ \left( \frac{n}{|m^*|} \right) \right],$$

which is, up to a logarithm term, the expected magnitude of  $|m^*|$  for the quantity  $\mathbf{h}^2(\mathbf{R}_{\gamma^*}, \mathbf{R}_{\hat{\gamma}})$ . If it is not the ideal case, an approximation error  $\mathbf{h}^2(\mathbf{R}^*, \mathcal{Q}_{m^*})$  with  $\mathcal{Q}_{m^*} = \{\mathbf{R}_{\gamma}, \gamma \in \bar{\Gamma}_{m^*}\}$ , will be added into the right hand side of (9) according to (8). However, as long as this bias term remains small, the performance of our selected estimator will not deteriorate too much as compared to the ideal situation.

**3.3. Connection to model selection.** The work done in this paper differs from the corresponding result (12) given by a model selection procedure in Chen (2022). In fact, one can regard Corollary 1 in this paper as a more general result of the one in Chen (2022). We illustrate this connection as follows.

We consider the particular application of our selection procedure in the context of model selection. For simplicity, let the equality holds in (5). We take  $\Lambda = \{1, \dots, |\tilde{\Gamma}|\}$  which is the index set of all the functions belonging to  $\tilde{\Gamma} = \cup_{m \in \mathcal{M}} \Gamma_m$  so that in this case,  $\hat{\Gamma} = \tilde{\Gamma} = \{\gamma_\lambda, \lambda \in \Lambda\}$  is a collection of deterministic candidates. Moreover, for each  $\lambda \in \Lambda$ , there exists a deterministic  $m(\lambda) \in \mathcal{M}$  such that  $\gamma_\lambda \in \Gamma_{m(\lambda)}$ . Let us denote  $\mathcal{Q}_m = \{\mathbf{R}_{\gamma}, \gamma \in \Gamma_m\}$ , for all  $m \in \mathcal{M}$ . We can immediately deduce from (8) that the estimator  $\mathbf{R}_{\hat{\gamma}}$  based on the selection among the family  $\{\Gamma_m, m \in \mathcal{M}\}$  satisfies

$$\mathbb{E}[\mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}})] \leq C_{\alpha, \Sigma} \inf_{m \in \mathcal{M}} [\mathbf{h}^2(\mathbf{R}^*, \mathcal{Q}_m) + D_n(m) + \Delta(m)],$$

which is, up to constants, the result (12) of Chen (2022) when one takes  $\Gamma_m$  in Chen (2022) as the collection of piecewise constant functions on  $\mathcal{W}$ . The difference is their model selection procedure, on the one hand, does not require Assumption 2 to be satisfied and can be applied to other types of models to approximate the potential  $\gamma^*$  besides piecewise constant ones. On the other hand, when the number of models becomes large, model selection strategy is more of theoretical interest due to its expensive numerical cost. Our estimator selection strategy, however, allows to deal with random partitions which can be obtained for example from dynamic programming algorithm (e.g. Rigail (2015)) or CART algorithm (e.g. Breiman et al. (1984)). Efficiently reducing the cardinality of  $\hat{\Gamma}$ , these algorithms together with our estimator selection procedure take the model selection strategy into practice. Moreover, the idea that selecting among random candidates set makes the selection between estimators given by different model selection strategies possible.

#### 4. APPLICATION TO CHANGEPOINT DETECTION IN EXPONENTIAL FAMILIES

In this section, we consider the application of our estimator selection procedure to changepoint detection problem in exponential families. In this

context, people usually assume the exponential family  $\mathcal{Q}_0 = \{R_\gamma, \gamma \in I\}$  has been parametrized in its natural form which entails  $u$  is taken as the identity function in (2) and  $A(\gamma) = \log [\int_{\mathcal{Y}} \exp(\gamma T(y)) d\nu(y)]$ . We observe a sequence  $\mathbf{Y} = (Y_1, \dots, Y_n)$  with values in  $\mathcal{Y}^n$  and assume that there exists a vector  $\gamma^* = (\gamma_1^*, \dots, \gamma_n^*) \in I^n$  with  $N - 1$  changepoints,  $N \geq 1$  such that within each segment, the values of  $\gamma^*$  remain a constant and for each  $i \in \{1, \dots, n\}$ , the distribution of  $Y_i$  is given by  $R_{\gamma_i^*}$ . This corresponds to the situation in our setting when  $W_i = (i - 1)/n$  are deterministic, for all  $i \in \{1, \dots, n\}$  so that  $\mathcal{W} = [0, 1)$  and the function  $\gamma^* : [0, 1) \rightarrow I \subset \mathbb{R}$  is a right-continuous step function with  $N \geq 1$  segments. For consistency with the previous paragraphs, we take  $W_i = (i - 1)/n$  throughout this section and use the function notation  $\gamma^*$  rather than the vector  $\gamma^* \in I^n$  in the sequel.

For each  $1 \leq k \leq n$ , let  $\mathcal{M}_k$  stand for the collection of all possible partitions of the sequence  $1, \dots, n$  into  $k$  segments and denote  $\mathcal{M} = \cup_{1 \leq k \leq n} \mathcal{M}_k$ . In changepoint detection problem, for each  $m \in \mathcal{M}$ , we assign its weight as

$$(10) \quad \Delta(m) = \log \binom{n-1}{|m|-1} + |m|.$$

With (10), a basic computation leads to  $\Sigma = \sum_{m \in \mathcal{M}} \exp[-\Delta(m)] \leq 1/(e-1)$  which entails Assumption 1 is satisfied. Moreover, since  $\mathcal{W} = [0, 1) \subset \mathbb{R}$ , for any  $m_1, m_2 \in \mathcal{M}$ ,  $|m_1 \vee m_2| \leq |m_1| + |m_2| - 1$ , Assumption 2 also holds true with  $\alpha = 1$ .

Supposing that we have a finite but arbitrary collection of (possibly random) piecewise constant candidates  $\widehat{\Gamma} = \{\widehat{\gamma}_\lambda(\mathbf{X}), \lambda \in \Lambda\}$ , we associate each  $\widehat{\gamma}_\lambda(\mathbf{X})$  with the penalty

$$\mathbf{pen}(\widehat{\gamma}_\lambda) = \kappa \left\{ |\widehat{m}(\lambda)| \left[ 10.11 + \log \left( \frac{n}{|\widehat{m}(\lambda)|} \right) \right] + \log \binom{n-1}{|\widehat{m}(\lambda)|-1} \right\},$$

where  $\kappa = 2.5C_0$  is the parameter to be tuned later. Once the value of  $\kappa$  is given, our estimator selection procedure can be implemented by running Algorithm 1.

---

**Algorithm 1** Estimator selection

---

**Input:**

$\mathbf{X} = (X_1, \dots, X_n)$ : the observations.

**Output:**  $\widehat{\gamma}_{\widehat{\lambda}}$

- 1: Collect  $\widehat{\Gamma} = \{\widehat{\gamma}_\lambda, \lambda \in \Lambda\}$  based on  $\mathbf{X}$ .
  - 2: **for**  $\lambda \in \Lambda$  **do**
  - 3:    $\mathbf{v}(\mathbf{X}, \widehat{\gamma}_\lambda) \leftarrow \sup_{\lambda' \in \Lambda} [\mathbf{T}(\mathbf{X}, \widehat{\gamma}_\lambda, \widehat{\gamma}_{\lambda'}) - \mathbf{pen}(\widehat{\gamma}_{\lambda'})] + \mathbf{pen}(\widehat{\gamma}_\lambda)$ .
  - 4: **end for**
  - 5:  $\widehat{\lambda} \leftarrow \operatorname{argmin}_{\lambda \in \Lambda} \mathbf{v}(\mathbf{X}, \widehat{\gamma}_\lambda)$ .
  - 6: Return  $\widehat{\gamma}_{\widehat{\lambda}}$ .
-

**4.1. Calibrating the value of  $\kappa$ .** We take  $\kappa = 0.08$  uniformly over all the exponential families. The reason for this choice of  $\kappa$  is explained in this section.

The idea to calibrate the value of  $\kappa$  is rather simple. Roughly speaking, we first simulate data of size  $n$  and prepare a collection of candidates  $\widehat{\Gamma}$  which can be done by running the algorithm in R package `Segmentor3IsBack` (implementing the procedure proposed by [Cleynen and Lebarbier \(2014, 2017\)](#)). Then we take different values of  $\kappa$  to design our penalty function and obtain a sequence of the selected  $\widehat{\gamma}_{\kappa, \widehat{\lambda}}$  among  $\widehat{\Gamma}$  associated to various  $\kappa$ . For each value of  $\kappa$ , we repeat the experiment under each simulation setting 100 times and finally evaluate the risk  $\mathbb{E} \left[ \mathbf{h}^2 \left( \mathbf{R}^*, \mathbf{R}_{\widehat{\gamma}_{\kappa, \widehat{\lambda}}} \right) \right]$  of the selected estimator  $\mathbf{R}_{\widehat{\gamma}_{\kappa, \widehat{\lambda}}}$  by its empirical mean, i.e. we compute

$$\widehat{R}_n \left( \widehat{\gamma}_{\kappa, \widehat{\lambda}} \right) = \frac{1}{100} \sum_{l=1}^{100} \left[ \sum_{i=1}^n h^2 \left( R_i^*, R_{\widehat{\gamma}_{\kappa, \widehat{\lambda}}^l} \left( \frac{i-1}{n} \right) \right) \right],$$

where  $\widehat{\gamma}_{\kappa, \widehat{\lambda}}^l$  is the  $l$ -th realisation of the selected estimator associated to a fixed  $\kappa$ .

**4.1.1. Simulating data.** We carry out experiments for three models: Gaussian, Poisson and exponential changepoint detection.

Let  $\gamma^*$  be piecewise constant on  $[0, 1)$  with  $N$  segments and  $\mathbf{R}^* = \mathbf{R}_{\gamma^*}$ . For each model, we design the experiments under three settings where for all the settings  $n = 500$ , but  $N = 5$ ,  $N = 10$  and  $N = 20$  respectively. For all the three settings, the changepoints are uniformly located, i.e. every 100 data-points for the first setting, every 50 data-points for the second setting and every 25 data-points for the third setting.

— Under all the settings of Gaussian model, for  $1 \leq i \leq n$ , if  $Y_i$  locates at the  $j$ -th segment with  $1 \leq j \leq N$ ,  $Y_i$  follows a Gaussian distribution with mean  $(j + 1)/2$ , variance  $\sigma^2 = 1$ .

— Under all the settings of Poisson model, for  $1 \leq i \leq n$ , if  $Y_i$  locates at the  $j$ -th segment with  $1 \leq j \leq N$ ,  $Y_i$  follows a Poisson distribution with mean  $j$  which means  $\gamma^*$  takes value  $\log(j)$  on the  $j$ -th segment.

— Under all the settings of exponential model, for  $1 \leq i \leq n$ , if  $Y_i$  locates at the  $j$ -th segment with  $1 \leq j \leq N$ ,  $Y_i$  follows an exponential distribution with natural parameter  $0.01j$ .

Figure 1 exhibits one example of the simulated data (when  $N = 10$ ) and the true value of the regression function  $\gamma^*$  (or a suitable transformation of  $\gamma^*$ ) on each segment.

**4.1.2. Collecting candidates in  $\widehat{\Gamma}$ .** In the work of [Cleynen and Lebarbier \(2014, 2017\)](#), they solved this problem by a model selection procedure via

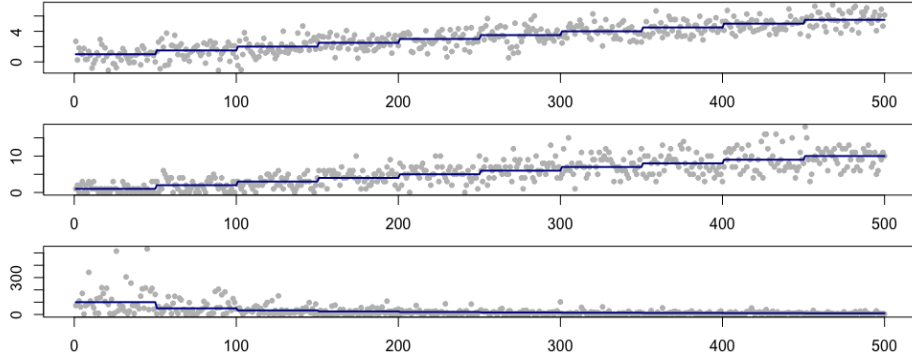


FIGURE 1. The 1st graph (top) corresponds to one profile of the simulated data (dots) and  $\gamma^*$  (solid line) for Gaussian model; The 2nd graph (middle) corresponds to one profile of the simulated data (dots) and  $\exp(\gamma^*)$  (solid line) for Poisson model; The 3rd graph (bottom) corresponds to one profile of the simulated data (dots) and  $1/\gamma^*$  (solid line) for exponential model.

some suitable penalty function based on the partitions given by the pruned dynamic programming algorithm (PDPA for short) proposed by [Rigaill \(2015\)](#). Given  $N_{\max}$  the maximum number of segments for consideration, for each integer  $\lambda$  with  $1 \leq \lambda \leq N_{\max}$ , PDPA searches the optimal partition with exact  $\lambda$  segments. We set  $N_{\max} = 30$  hence 30 different partitions of the sequence  $1, \dots, n$  are returned by PDPA. Provided a partition, the value of  $\gamma^*$  on each segment is given by maximum likelihood estimation as it was done in [Cleynen and Lebarbier \(2014, 2017\)](#). By doing so, we collect 30 candidates which we denote as  $\hat{\Gamma}_c = \{\hat{\gamma}_\lambda, 1 \leq \lambda \leq N_{\max}\}$ .

4.1.3. *Results.* Under each setting of all the three models, one experiment means we simulate  $n = 500$  observations with  $N$  segments based on the corresponding  $\gamma^*$  introduced in Section 4.1.1. We then select the estimator among the candidate ones  $\hat{\Gamma}_c$  by Algorithm 1 via the penalty functions associated to different values of  $\kappa$ . Finally we observe the quantity  $\hat{R}_n(\hat{\gamma}_{\kappa, \hat{\lambda}})$  and regard it as the criterion to calibrate a suitable value of  $\kappa$ . The results for all nine settings are shown in Figure 2, where the horizontal axis represents the value of  $\kappa$  and the vertical axis indicates the quantity  $\hat{R}_n(\hat{\gamma}_{\kappa, \hat{\lambda}})$ .

In Figure 2, the quantities  $\hat{R}_n(\hat{\gamma}_{\kappa, \hat{\lambda}})$  in all nine settings have a tendency to first decrease and then increase with respect to the increasing of  $\kappa$ , which is consistent to the theoretical results. When  $\kappa$  is too small, the penalty function is relatively small for the complexed models therefore the overfitting issue may happen. However, when  $\kappa$  is too large, the penalty function is excessively large for the complexed models which will cause an overpenalization. Moreover, the minimizers of  $\kappa$  for the quantities  $\hat{R}_n(\hat{\gamma}_{\kappa, \hat{\lambda}})$  in all nine settings are very close to each other and all concentrate within a short

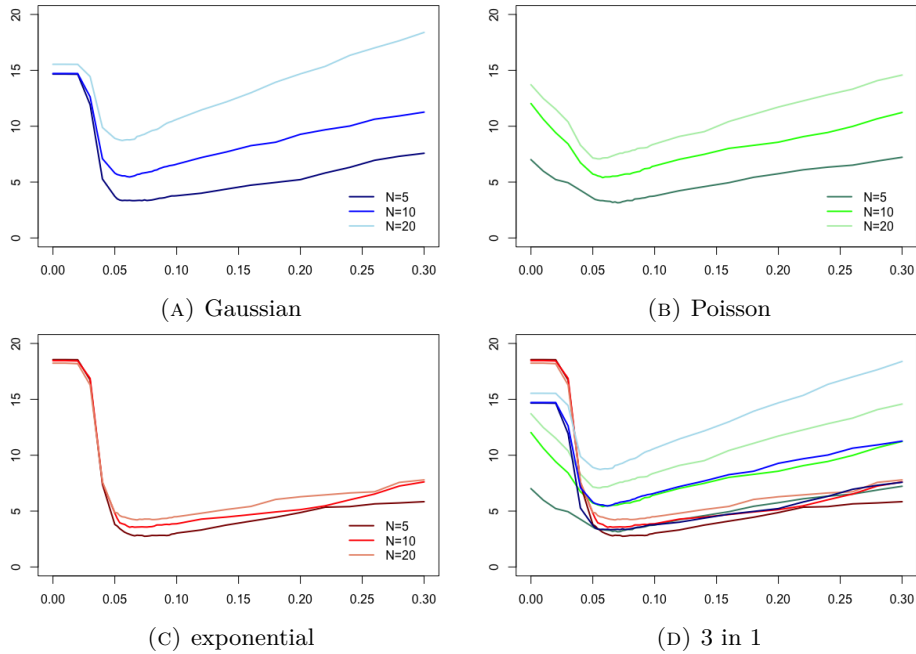


FIGURE 2.  $\widehat{R}_n(\widehat{\gamma}_{\kappa, \widehat{\lambda}})$  with respect to  $\kappa$  under nine settings.

interval  $[0.05, 0.1]$ . Considering the optimal performance of all the settings and also being safe with respect to overfitting, we choose  $\kappa$  as the largest minimizer of  $\widehat{R}_n(\widehat{\gamma}_{\kappa, \widehat{\lambda}})$  among nine settings which approximately equals to 0.08 and implement our procedure with  $\kappa = 0.08$  in later studies.

## 5. SIMULATION STUDY AND DISCUSSION

Throughout this section, we carry out a comparative simulation study with the state-of-art competitors available in R packages for changepoint detection problem in exponential families. Unless otherwise specified, the competitors are implemented under the default settings in their packages. For Gaussian model, some of our competitors use the estimated value of the standard deviation  $\sigma$ . To make the comparison as fair as possible, we also implement the median absolute deviation estimator for  $\sigma$  while running our procedure, which is the one adopted in [Killick et al. \(2012\)](#) and [Fearnhead and Rigail \(2019\)](#).

To evaluate the performance of each estimator, besides the empirical risk  $\widehat{R}_n(\cdot)$  obtained from replications, we also record  $\widehat{N} - N$  which computes the difference between the estimated number of segments and the truth for each replication.

**5.1. Accuracy.** In this section, we study the changepoint detection problem for Gaussian model where numerous literature can be found tackling this issue. We construct our candidates set  $\hat{\Gamma}$  as a collection of some cutting-edge estimators with implemented R packages and these ones are also regarded as the competitors of our estimator ES. More precisely, the competing packages we consider are: `PSCBS`, which implements the CBS procedure proposed in [Olshen et al. \(2004\)](#); `cumSeg`, which performs the method given by [Muggeo and Adelfio \(2010\)](#); `changepoint`, which implements the PELT approach provided by [Killick et al. \(2012\)](#); `StepR`, which implements the SMUCE given by [Frick et al. \(2013\)](#); `Segmentor3IsBack`, which implements CL proposed by [Cleynen and Lebarbier \(2014, 2017\)](#); `wbs`, which implements the wild binary segmentation methodology proposed in [Fryzlewicz \(2014\)](#); `FDRSeg`, which implements the approach given in [Li et al. \(2016\)](#); `robseg`, which implements the procedure proposed by [Fearnhead and Rigaiil \(2019\)](#). We would like to study the performance of our estimator ES based on the selection among these state-of-art ones.

We follow the test signals considered by [Fryzlewicz \(2014\)](#) and then by [Fearnhead and Rigaiil \(2019\)](#) which involves 5 different formats of signals with length from  $n = 140$  to 2048: (1) `blocks`, (2) `fms`, (3) `mix`, (4) `teeth10` and (5) `stairs10`. The specific settings of these signals including the sample sizes and noise standard deviations are given in Appendix B of [Fryzlewicz \(2014\)](#). Following the experiments done in [Fearnhead and Rigaiil \(2019\)](#), we also consider an additional signal setting by changing the standard deviation of (2) `fms` from 0.3 into 0.2, which is also one of the settings studied in [Frick et al. \(2013\)](#). An example of one profile of the simulated data and the underlying signals  $\gamma^*$  are plotted in [Figure 3](#). For each signal, the experiment has been replicated 1000 times. The results are shown in [Table 1](#). The performance of each estimator is stated as follows.

**CBS and cumSeg.** The CBS and cumSeg in general behave poorly compared with other procedures. The CBS only has satisfactory performance of detecting changes for `blocks` and `fms` ( $\sigma = 0.2$ ) but it turns out CBS always results in a relatively large empirical risk  $\hat{R}_n(\cdot)$ . Except acceptable performance for `fms` ( $\sigma = 0.2$ ) and `stairs10`, cumSeg always tends to underestimate the number of changes and also yields an estimator with quite large empirical risk.

**PELT.** The PELT has excellent performance for both of the `fms` signals and `stairs10`. For `blocks` signal, it is above the average but does not belong to the first class among all. As for `mix` and `teeth10`, it performs rather average.

**SMUCE.** The SMUCE has very excellent performance for `fms` ( $\sigma = 0.2$ ). However, it behaves poorly for all the other signals.



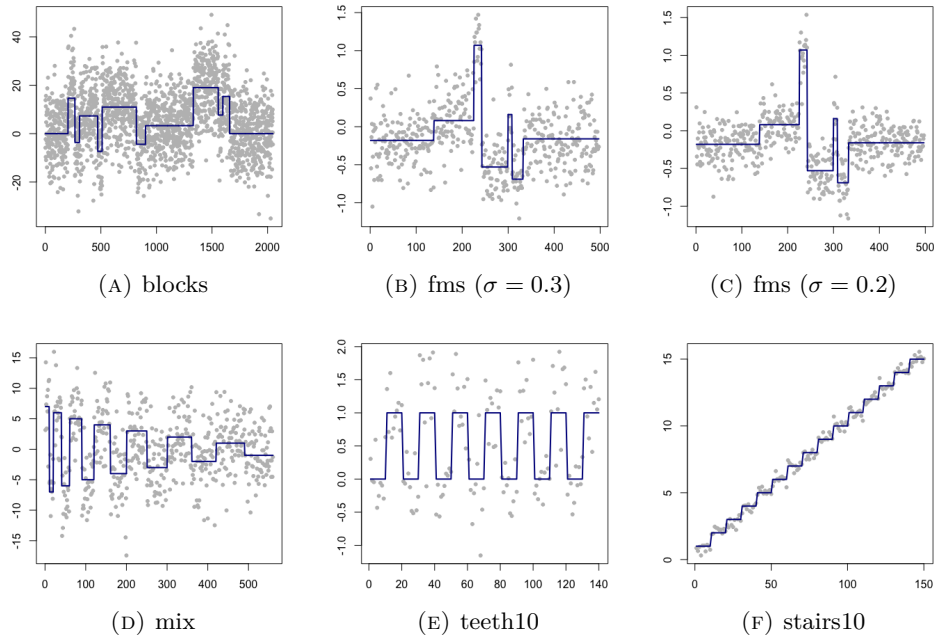


FIGURE 3. The six signals (solid line) and simulated data (dots).

CL. The CL has nice performance for `teeth10`. For `blocks` and `mix`, its performance is satisfactory though not belonging to the first class. For both of the `fms` signals, it shows rather average performance. The CL does not behave well for the `stairs10` signal where it tends to overestimate the number of changes compared to other methods. Let us remark here that the performance of CL in our simulation study is better than the corresponding context in [Fryżlewicz \(2014\)](#). This is because when implementing the package `Segmentor3IsBack`, users need to set the maximum number of segments  $N_{\max}$ . We set  $N_{\max} = 20$  for all the six signals considering the maximal number of changepoints (i.e.  $N - 1$ ) among six signals is 14 and they set  $N_{\max} = 15$  which resulted in a systematical underestimation of the number of changepoints for CL in their study.

WBS sSIC. We implement the package `wbs` combining the WBS method with the sSIC stopping criterion which, as it has been shown in [Fryżlewicz \(2014\)](#), is the overall winner compared to combining the WBS method with other thresholding stopping rules. The WBS sSIC has excellent performance for both of the `fms` signals and `teeth10`. However, it performs rather average for `blocks` and `mix`. As for `stairs10`, the performance of WBS sSIC is a little poor as a consequence of overestimating the number of changepoints. Such a result has also been confirmed by the study of WBS sSIC in [Fryżlewicz \(2014\)](#).

Method	Signal	$\widehat{N} - N$					$\widehat{R}_n(\cdot)$	Contribution
		$\leq -2$	-1	0	1	$\geq 2$		
ES	blocks	0.005	0.278	<b>0.656</b>	0.055	0.006	5.61 ± 0.12	-
CBS	blocks	0.006	0.090	0.575	0.184	0.145	7.57 ± 0.14	0.000
cumSeg	blocks	0.653	0.335	0.011	0.001	0.000	15.71 ± 0.40	0.000
PELT	blocks	0.014	0.389	0.574	0.020	0.003	5.69 ± 0.11	0.035
SMUCE	blocks	0.940	0.060	0.000	0.000	0.000	16.02 ± 0.37	0.010
CL	blocks	0.010	0.356	0.595	0.035	0.004	5.67 ± 0.12	0.533
WBS sSIC	blocks	0.021	0.412	0.532	0.032	0.003	6.11 ± 0.13	0.013
FDR( $\alpha = 0.05$ )	blocks	0.008	0.447	0.478	0.059	0.008	6.15 ± 0.13	0.332
robseg(Huber)	blocks	0.004	0.234	<b>0.674</b>	0.072	0.016	5.84 ± 0.12	0.063
robseg(biweight)	blocks	0.020	0.404	0.558	0.017	0.001	5.88 ± 0.12	0.014
ES	fms(0.3)	0.008	0.002	<b>0.915</b>	0.069	0.006	2.16 ± 0.07	-
CBS	fms(0.3)	0.007	0.012	0.796	0.139	0.046	5.10 ± 0.09	0.000
cumSeg	fms(0.3)	0.706	0.041	0.224	0.028	0.001	7.07 ± 0.44	0.000
PELT	fms(0.3)	0.007	0.003	<b>0.922</b>	0.061	0.007	2.15 ± 0.08	0.054
SMUCE	fms(0.3)	0.074	0.537	0.388	0.001	0.000	5.15 ± 0.18	0.293
CL	fms(0.3)	0.002	0.001	0.837	0.119	0.041	2.28 ± 0.08	0.199
WBS sSIC	fms(0.3)	0.007	0.003	<b>0.933</b>	0.048	0.009	2.26 ± 0.08	0.008
FDR( $\alpha = 0.05$ )	fms(0.3)	0.001	0.027	<b>0.879</b>	0.076	0.017	2.28 ± 0.09	0.409
robseg(Huber)	fms(0.3)	0.001	0.001	0.825	0.130	0.043	2.37 ± 0.08	0.007
robseg(biweight)	fms(0.3)	0.013	0.005	<b>0.928</b>	0.049	0.005	2.23 ± 0.08	0.030
ES	fms(0.2)	0.000	0.000	<b>0.923</b>	0.071	0.006	1.61 ± 0.06	-
CBS	fms(0.2)	0.000	0.000	0.871	0.086	0.043	5.79 ± 0.07	0.000
cumSeg	fms(0.2)	0.094	0.009	0.812	0.083	0.002	5.19 ± 0.22	0.002
PELT	fms(0.2)	0.000	0.000	<b>0.929</b>	0.060	0.011	1.59 ± 0.06	0.022
SMUCE	fms(0.2)	0.000	0.001	<b>0.994</b>	0.005	0.000	1.49 ± 0.06	0.734
CL	fms(0.2)	0.000	0.000	0.840	0.128	0.032	1.74 ± 0.07	0.102
WBS sSIC	fms(0.2)	0.000	0.000	<b>0.945</b>	0.050	0.005	1.65 ± 0.06	0.003
FDR( $\alpha = 0.05$ )	fms(0.2)	0.000	0.000	0.871	0.103	0.026	1.66 ± 0.06	0.115
robseg(Huber)	fms(0.2)	0.000	0.000	0.830	0.135	0.035	1.83 ± 0.07	0.008
robseg(biweight)	fms(0.2)	0.000	0.000	<b>0.937</b>	0.058	0.005	1.63 ± 0.06	0.014
ES	mix	0.264	0.243	<b>0.434</b>	0.056	0.003	5.91 ± 0.12	-
CBS	mix	0.313	0.201	0.324	0.109	0.053	11.18 ± 0.17	0.000
cumSeg	mix	0.999	0.001	0.000	0.000	0.000	32.61 ± 0.92	0.000
PELT	mix	0.375	0.270	0.321	0.032	0.002	6.11 ± 0.12	0.070
SMUCE	mix	0.922	0.076	0.002	0.000	0.000	12.59 ± 0.42	0.042
CL	mix	0.305	0.244	0.390	0.053	0.008	6.04 ± 0.12	0.585
WBS sSIC	mix	0.342	0.269	0.351	0.032	0.006	5.99 ± 0.12	0.029
FDR( $\alpha = 0.05$ )	mix	0.411	0.358	0.181	0.038	0.012	6.71 ± 0.13	0.190
robseg(Huber)	mix	0.209	0.240	<b>0.444</b>	0.088	0.019	6.10 ± 0.12	0.051
robseg(biweight)	mix	0.403	0.264	0.305	0.026	0.002	6.30 ± 0.12	0.033
ES	teeth10	0.215	0.025	<b>0.721</b>	0.037	0.002	5.69 ± 0.24	-
CBS	teeth10	0.999	0.000	0.001	0.000	0.000	24.69 ± 0.07	0.000
cumSeg	teeth10	1.000	0.000	0.000	0.000	0.000	24.85 ± 0.01	0.005
PELT	teeth10	0.274	0.029	0.657	0.037	0.003	6.03 ± 0.24	0.090
SMUCE	teeth10	0.984	0.013	0.003	0.000	0.000	20.11 ± 0.22	0.003
CL	teeth10	0.029	0.013	<b>0.679</b>	0.204	0.075	4.71 ± 0.13	0.321
WBS sSIC	teeth10	0.067	0.021	<b>0.752</b>	0.120	0.040	5.30 ± 0.26	0.010
FDR( $\alpha = 0.05$ )	teeth10	0.309	0.135	0.508	0.040	0.008	7.68 ± 0.32	0.356
robseg(Huber)	teeth10	0.105	0.026	<b>0.748</b>	0.102	0.019	4.94 ± 0.15	0.016
robseg(biweight)	teeth10	0.318	0.028	0.635	0.019	0.000	6.31 ± 0.25	0.199
ES	stairs10	0.00	0.004	<b>0.949</b>	0.044	0.003	3.33 ± 0.09	-
CBS	stairs10	0.012	0.172	0.789	0.027	0.000	13.81 ± 0.16	0.000
cumSeg	stairs10	0.024	0.090	0.819	0.067	0.000	8.61 ± 0.24	0.000
PELT	stairs10	0.000	0.004	<b>0.955</b>	0.039	0.002	3.32 ± 0.09	0.017
SMUCE	stairs10	0.801	0.137	0.062	0.000	0.000	22.26 ± 0.58	0.050
CL	stairs10	0.000	0.001	0.768	0.184	0.047	3.50 ± 0.09	0.178
WBS sSIC	stairs10	0.000	0.001	0.608	0.301	0.090	3.91 ± 0.10	0.004
FDR( $\alpha = 0.05$ )	stairs10	0.002	0.028	<b>0.896</b>	0.053	0.021	3.57 ± 0.12	0.703
robseg(Huber)	stairs10	0.000	0.000	0.867	0.110	0.023	3.45 ± 0.09	0.006
robseg(biweight)	stairs10	0.000	0.005	<b>0.964</b>	0.031	0.000	3.36 ± 0.09	0.042

TABLE 1. Frequencies of  $\widehat{N} - N$  and  $\widehat{R}_n(\cdot)$  of ES and its competitors for Gaussian model over 1000 simulated sample paths. Contribution denotes the frequency of each competitor being selected as ES. Bold: highest empirical frequency of  $\widehat{N} - N = 0$  and those with frequencies within 10% off the highest. The uncertainty is obtained by computing  $2\widehat{\sigma}/\sqrt{n_r}$ , where  $\widehat{\sigma}^2$  is the empirical variance and  $n_r$  is the number of replications.

FDR. The FDR with  $\alpha = 0.05$  performs well for `fms` ( $\sigma = 0.3$ ) and `satirs10` signals. For `fms` ( $\sigma = 0.2$ ), it has an average performance. But it behaves below the average under other test signals.

robseg. We consider Huber loss and biweight loss when implementing the package `robseg` which are the recommended ones (especially the biweight loss) according to [Fearnhead and Rigaiil \(2019\)](#). We adopted the suggested values in the Section 5.2 of their paper to set the parameters in their algorithms. The `robseg` (Huber) performs excellently for `blocks`, `mix` and `teeth10`. It behaves rather average for both of the `fms` signals and `stairs10`. The `robseg` (biweight) performs excellently for both of the `fms` signals and `stairs10`. As for `blocks`, `mix` and `teeth10` signals, it performs rather average.

ES. As we can observe from the column named ‘‘Contribution’’ in Table 1, under different test signals, our estimator selection procedure tends to allocate different preference to the candidates in  $\hat{\Gamma}$  based on their practical performance. For example, when SMUCE shows obvious outperformance for the signal `fms` ( $\sigma = 0.2$ ), we select it with a frequency 0.734 as our ES estimator. However, when SMUCE performs poorly under other signals we automatically reduce the frequency to select it as ES but prefer some more competitive ones. As a final result, our ES estimator shows a very competitive performance under all the test signals. The interesting point is that this cannot be achieved by any single candidate in  $\hat{\Gamma}$  since as we have seen above, each of them only outperforms others for some of the test signals but not all.

**5.2. Stability when outliers present.** As we have mentioned in the theoretical analysis part, our estimator selection procedure possesses the stability when there is a slight departure from the presumption  $\mathbf{R}^* = \mathbf{R}_{\gamma^*}$  with  $\gamma^*$  being piecewise constant on  $\mathscr{W}$ . One of the application scenario for this property is when there is a small proportion of outliers presenting in the observations which has attracted more and more attention recently in the changepoint detection. In this section, we test the practical performance of ES as well as its competitors when outliers present. We take the signal `fms` ( $\sigma = 0.2$ ) as an example since most of the existing methods behave rather well under this signal. Based on this signal, we add outliers by randomly choosing five points among the sequence of length  $n = 497$  and modifying the values of them into 3. The results of all the estimators are shown in Table 2.

We can observe from Table 2 that in such a scenario PELT, SMUCE, CL, WBS sSIC, FDR and `robseg` (Huber) are all not robust with respect to the outliers and they all overestimate the number of changepoints due to fitting the outliers. The CBS and `cumSeg` still systematically underestimate the number of changepoints. It is not that surprising `robseg` (biweight) proposed

Method	Signal	Outlier	$\widehat{N} - N$					$\widehat{R}_n(\cdot)$	Contribution
			$\leq -2$	$-1$	$0$	$1$	$\geq 2$		
ES	fms(0.2)	Yes	0.000	0.000	<b>0.956</b>	0.043	0.001	$1.64 \pm 0.06$	-
CBS	fms(0.2)	Yes	0.660	0.282	0.038	0.016	0.004	$34.55 \pm 0.79$	0.000
cumSeg	fms(0.2)	Yes	0.801	0.056	0.083	0.021	0.039	$16.96 \pm 0.51$	0.000
PELT	fms(0.2)	Yes	0.000	0.000	0.000	0.000	1.000	$7.27 \pm 0.07$	0.000
SMUCE	fms(0.2)	Yes	0.000	0.000	0.000	0.000	1.000	$8.02 \pm 0.11$	0.000
CL	fms(0.2)	Yes	0.000	0.000	0.000	0.000	1.000	$7.29 \pm 0.07$	0.000
WBS sSIC	fms(0.2)	Yes	0.000	0.000	0.000	0.000	1.000	$7.33 \pm 0.07$	0.000
FDR( $\alpha = 0.05$ )	fms(0.2)	Yes	0.000	0.000	0.000	0.000	1.000	$7.44 \pm 0.07$	0.000
robseg(Huber)	fms(0.2)	Yes	0.000	0.000	0.000	0.000	1.000	$7.51 \pm 0.08$	0.000
robseg(biweight)	fms(0.2)	Yes	0.000	0.000	<b>0.956</b>	0.043	0.001	$1.64 \pm 0.06$	1.000

TABLE 2. Frequencies of  $\widehat{N} - N$  and  $\widehat{R}_n(\cdot)$  of ES and its competitors for fms ( $\sigma = 0.2$ ) signal with 5 outliers over 1000 simulated sample paths. Contribution denotes the frequency of each competitor being selected as ES. Bold: highest empirical frequency of  $\widehat{N} - N = 0$ . The uncertainty is obtained by computing  $2\widehat{\sigma}/\sqrt{n_r}$ , where  $\widehat{\sigma}^2$  is the empirical variance and  $n_r$  is the number of replications.

in [Fearnhead and Rigail \(2019\)](#) is quite robust in this scenario since it was designed to handle such an issue. It shows a very high frequency 0.956 to recover the correct number of changepoints. Moreover, from the quantity of empirical risk  $\widehat{R}_n(\cdot)$ , it turns out robseg (biweight) outperforms all the other candidates significantly which also indicates an excellent performance of localising the changepoints as well as estimating the value of  $\gamma^*$  on each segment. Our selection procedure automatically gives the preference to robseg (biweight) in this case with frequency 1.000 which confirms the stability of our selection rule practically.

**5.3. From Gaussian to Poisson and exponential models.** As we have mentioned in the introduction, there are not too many works in the statistical literature addressing changepoint detection for Poisson and exponential models and establishing a theoretical guarantee for the proposed estimator. The CL method proposed by [Cleynen and Lebarbier \(2014, 2017\)](#) performs a model selection procedure based on the partitions given by [Rigail \(2015\)](#) and they have proved the resulting estimator satisfies some oracle inequality. The R package `Segmentor3IsBack` implements their procedure and tackle both of Poisson and exponential models. Another approach is given by [Frick et al. \(2013\)](#) with the R package `StepR` where algorithm is only available for Poisson segmentation.

A selection merely based on these two estimators is boring. Recall that in [Section 3.2](#), one feature of our selection procedure is enlarging the (possibly random) collection  $\widehat{\Gamma}$  but keeping  $\mathcal{M}$  unchanged, the risk bound for the selected estimator only decreases (or at least keeps unchanged) over the larger collection. Therefore, for Poisson and exponential models, besides CL and SMUCE (if available), we would like to recruit some reasonable estimators into our candidates set  $\widehat{\Gamma}$ . Although these estimators do not exist

in the literature and no quantitative or qualitative analysis for them, once they are selected as ES by our selection procedure, the theoretical guarantee we built in Section 3.2 indicates that, up to a constant, they perform better than the state-of-art ones (CL and SMUCE).

One natural idea is to borrow the estimators for Gaussian model which is the case intensively studied. Inspired by Brown et al. (2010) where they implemented a mean-matching variance stabilizing transformation (MM-VST for short) to turn the problem of regression in exponential families into a standard homoscedastic Gaussian regression problem, we can perform a similar technique to the observations  $\mathbf{Y}$ . For more details of MM-VST, we refer Section 2 of Brown et al. (2010). Let us remark that while implementing MM-VST, we need to choose the value of  $m$  which corresponds to the number of data-points binned for transformation. Although it turns out that for regression problem, this  $m$  needs to be suitably chosen (see Section 4 of Brown et al. (2010)), we do not want this pre-process step presumes any information of the segmentation as we are in the context of changepoint detection. Therefore, we simply take  $m = 1$  in their transformation procedure and implement the formula  $Y'_i = 2\sqrt{Y_i + 1/4}$  for Poisson model and  $Y'_i = \log(2Y_i)$  for exponential model to derive new sequences of observations  $\mathbf{Y}' = (Y'_1, \dots, Y'_n)$ . We then apply the algorithms introduced in the last section to  $\mathbf{Y}'$  to get the locations of changepoints. Based on these locations, we associate  $\rho$ -estimators proposed in Baraud and Chen (2020) to the estimated values of  $\gamma^*$  on each segment to improve the performance. As it was shown in Baraud and Chen (2020), under some suitable conditions and when the model is exact,  $\rho$ -estimator recovers the accurate result given by MLE. Moreover, it possesses more robustness compared to MLE when there is a model misspecification and/or data contamination. To conclude, the candidates set for Poisson model is given by

$$(11) \quad \widehat{\Gamma} = \left\{ \text{SMUCE}, \text{CL}, \text{CBS}^t + \rho, \text{cumSeg}^t + \rho, \text{PELT}^t + \rho, \text{WBS sSIC}^t + \rho, \right. \\ \left. \text{FDR}^t(\alpha = 0.05) + \rho, \text{robseg}(\text{Huber})^t + \rho, \text{robseg}(\text{biweight})^t + \rho \right\},$$

where the character “t” indicates the procedure is implemented on the transformed data. For exponential model,  $\widehat{\Gamma}$  is constructed the same as (11) except we change SMUCE to  $\text{SMUCE}^t + \rho$  since it is no longer available.

To investigate the performance of ES and the candidates in  $\widehat{\Gamma}$ , we mimic the test signals `fms` and `mix` for Poisson model and `teeth10` and `stairs10` for exponential model. We also study the scenario when outliers present in the observations for the mimic signals `fms` (Poisson) and `teeth10` (exponential). We shall describe the specific settings of these signals as well as how we add outliers in Appendix. Figure 4 exhibits the four underlying signals together with one profile of the simulated data for each signal.

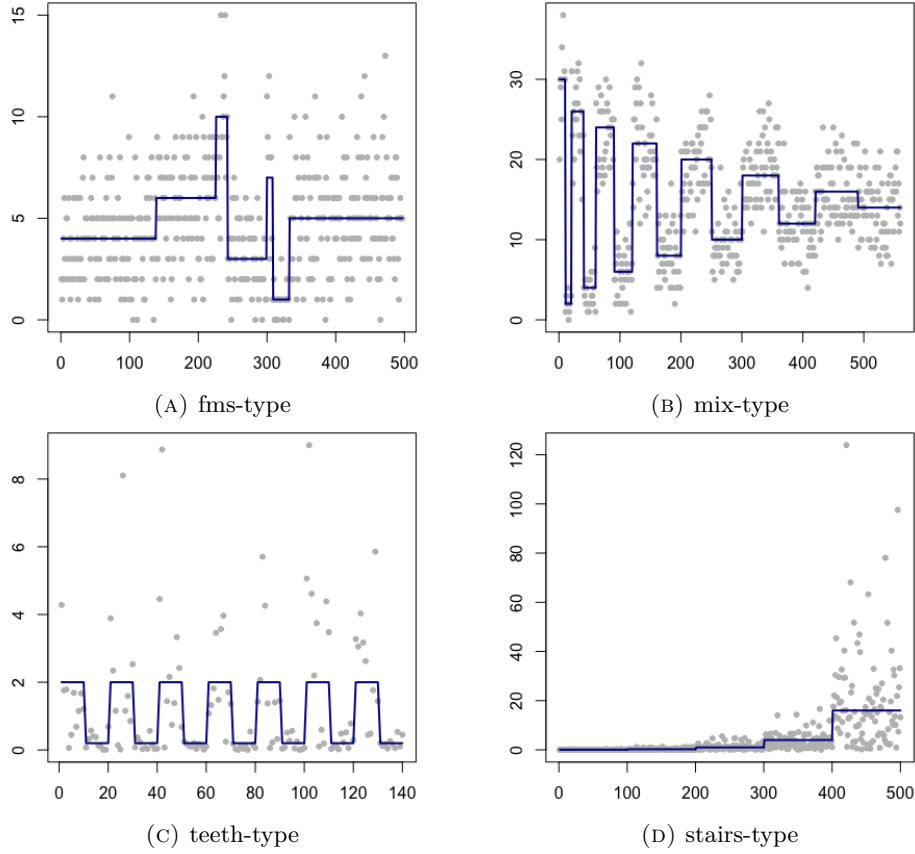


FIGURE 4. (A) and (B): the test signals of the form  $\exp(\gamma^*)$  (solid line) and simulated data (dots) for Poisson model. (C) and (D): the test signals of the form of  $1/\gamma^*$  (solid line) and simulated data (dots) for exponential model.

The results of Poisson model are shown in Table 3. Let us first comment the two existing estimators in the literature, namely SMUCE and CL. In both of the scenarios with or without outliers, the performance of SMUCE is quite poor at least under these two test signals. When no outlier presents in the observations, SMUCE has a tendency to underestimate the number of changepoints for both of the two signals `fms-type` and `mix-type`. When there are outliers, SMUCE is sensitive to them therefore overestimates the number of changepoints. The CL performs much better than SMUCE in the scenario where no outlier presents in the observations but it is also not robust with respect to the outliers. When no outlier presents, our estimator ES slightly improves the performance of CL on detecting changes under both of the two signals. When there are outliers presenting in the observations, ES obviously outperforms CL as a consequence of enjoying the excellent performance given by `robseg (biweight)`<sup>†</sup>. Interestingly, we find that when

Method	Signal	Outlier	$\widehat{N} - N$				$\widehat{R}_n(\cdot)$	Contribution	
			$\leq -2$	$-1$	<b>0</b>	$1$			$\geq 2$
ES	fms-type	No	0.002	0.050	<b>0.878</b>	0.062	0.008	$2.51 \pm 0.09$	-
SMUCE	fms-type	No	0.288	0.528	0.184	0.000	0.000	$6.21 \pm 0.19$	0.184
CL	fms-type	No	0.000	0.046	<b>0.854</b>	0.082	0.018	$2.54 \pm 0.10$	0.725
CBS <sup>t</sup> + $\rho$	fms-type	No	0.030	0.254	0.546	0.131	0.039	$5.34 \pm 0.11$	0.000
cumSeg <sup>t</sup> + $\rho$	fms-type	No	0.424	0.374	0.193	0.009	0.000	$7.26 \pm 0.20$	0.001
PELT <sup>t</sup> + $\rho$	fms-type	No	0.003	0.054	<b>0.867</b>	0.062	0.014	$2.56 \pm 0.10$	0.015
WBS sSIC <sup>t</sup> + $\rho$	fms-type	No	0.010	0.132	0.781	0.051	0.026	$3.08 \pm 0.12$	0.013
FDR <sup>t</sup> ( $\alpha = 0.05$ ) + $\rho$	fms-type	No	0.288	0.528	0.184	0.000	0.000	$5.97 \pm 0.18$	0.000
robseg(Huber) <sup>t</sup> + $\rho$	fms-type	No	0.001	0.035	<b>0.800</b>	0.130	0.034	$2.68 \pm 0.10$	0.032
robseg(biweight) <sup>t</sup> + $\rho$	fms-type	No	0.005	0.073	<b>0.867</b>	0.048	0.007	$2.63 \pm 0.10$	0.030
ES	fms-type	Yes	0.001	0.092	<b>0.825</b>	0.070	0.012	$3.78 \pm 0.11$	-
SMUCE	fms-type	Yes	0.000	0.000	0.000	0.000	1.000	$12.76 \pm 0.21$	0.000
CL	fms-type	Yes	0.000	0.000	0.000	0.000	1.000	$8.58 \pm 0.11$	0.000
CBS <sup>t</sup> + $\rho$	fms-type	Yes	0.521	0.354	0.086	0.035	0.004	$11.98 \pm 0.36$	0.000
cumSeg <sup>t</sup> + $\rho$	fms-type	Yes	0.795	0.164	0.038	0.003	0.000	$12.01 \pm 0.28$	0.009
PELT <sup>t</sup> + $\rho$	fms-type	Yes	0.000	0.000	0.000	0.000	1.000	$8.45 \pm 0.10$	0.000
WBS sSIC <sup>t</sup> + $\rho$	fms-type	Yes	0.000	0.000	0.000	0.000	1.000	$8.82 \pm 0.12$	0.000
FDR <sup>t</sup> ( $\alpha = 0.05$ ) + $\rho$	fms-type	Yes	0.000	0.000	0.000	0.000	1.000	$11.35 \pm 0.18$	0.001
robseg(Huber) <sup>t</sup> + $\rho$	fms-type	Yes	0.000	0.008	0.048	0.062	0.882	$6.13 \pm 0.13$	0.053
robseg(biweight) <sup>t</sup> + $\rho$	fms-type	Yes	0.000	0.092	<b>0.839</b>	0.066	0.003	$3.74 \pm 0.11$	0.937
ES	mix-type	No	0.005	0.371	<b>0.523</b>	0.091	0.010	$3.98 \pm 0.09$	-
SMUCE	mix-type	No	0.128	0.828	0.044	0.000	0.000	$4.67 \pm 0.13$	0.339
CL	mix-type	No	0.014	0.439	0.466	0.071	0.010	$3.99 \pm 0.09$	0.481
CBS <sup>t</sup> + $\rho$	mix-type	No	0.034	0.448	0.358	0.122	0.038	$13.39 \pm 0.11$	0.000
cumSeg <sup>t</sup> + $\rho$	mix-type	No	0.990	0.010	0.000	0.000	0.000	$31.18 \pm 0.45$	0.000
PELT <sup>t</sup> + $\rho$	mix-type	No	0.010	0.443	0.466	0.071	0.010	$4.03 \pm 0.09$	0.027
WBS sSIC <sup>t</sup> + $\rho$	mix-type	No	0.018	0.509	0.402	0.056	0.015	$4.03 \pm 0.09$	0.013
FDR <sup>t</sup> ( $\alpha = 0.05$ ) + $\rho$	mix-type	No	0.128	0.828	0.044	0.000	0.000	$4.67 \pm 0.12$	0.000
robseg(Huber) <sup>t</sup> + $\rho$	mix-type	No	0.003	0.293	<b>0.530</b>	0.149	0.025	$4.15 \pm 0.09$	0.099
robseg(biweight) <sup>t</sup> + $\rho$	mix-type	No	0.014	0.486	0.458	0.040	0.002	$4.06 \pm 0.09$	0.041

TABLE 3. Frequencies of  $\widehat{N} - N$  and  $\widehat{R}_n(\cdot)$  of ES and its competitors for Poisson model over 1000 simulated sample paths. Contribution denotes the frequency of each competitor being selected as ES. Bold: highest empirical frequency of  $\widehat{N} - N = 0$  and those with frequencies within 10% off the highest. The uncertainty is obtained by computing  $2\widehat{\sigma}/\sqrt{n_r}$ , where  $\widehat{\sigma}^2$  is the empirical variance and  $n_r$  is the number of replications.

there is no outlier presenting in the observations, the combinations PELT<sup>t</sup> +  $\rho$  and robseg<sup>t</sup> +  $\rho$  are perhaps nice choices at least under these two signals.

The results for exponential model are shown in Table 4. Under the teeth-type signal without an outlier, the ES obviously outperforms any single candidate by selecting mainly from CL and robseg (biweight)<sup>t</sup> +  $\rho$ . When there are outliers, robseg (biweight)<sup>t</sup> +  $\rho$  is the best one among all and we observe that ES improves the frequency of selecting robseg (biweight)<sup>t</sup> +  $\rho$  as the final estimator so that finally ES achieves a competitive performance compared to robseg (biweight)<sup>t</sup> +  $\rho$  and significantly outperforms the existing estimator CL. For stairs-type signal, CL performs quite nice but ES still slightly improves it by enjoying the contribution from other candidates in  $\widehat{\Gamma}$ .

Method	Signal	Outlier	$\widehat{N} - N$				$\widehat{R}_n(\cdot)$	Contribution	
			$\leq -2$	$-1$	$0$	$1$			$\geq 2$
ES	teeth-type	No	0.327	0.077	<b>0.468</b>	0.106	0.022	$7.69 \pm 0.25$	-
CL	teeth-type	No	0.381	0.055	0.411	0.116	0.037	$9.27 \pm 0.38$	0.766
SMUCE <sup>t</sup> + $\rho$	teeth-type	No	0.998	0.002	0.000	0.000	0.000	$20.29 \pm 0.14$	0.000
CBS <sup>t</sup> + $\rho$	teeth-type	No	1.000	0.000	0.000	0.000	0.000	$22.52 \pm 0.06$	0.000
cumSeg <sup>t</sup> + $\rho$	teeth-type	No	1.000	0.000	0.000	0.000	0.000	$22.56 \pm 0.05$	0.000
PELT <sup>t</sup> + $\rho$	teeth-type	No	0.134	0.068	0.241	0.191	0.366	$9.14 \pm 0.19$	0.015
WBS sSIC <sup>t</sup> + $\rho$	teeth-type	No	0.829	0.022	0.058	0.036	0.055	$18.62 \pm 0.37$	0.007
FDR <sup>t</sup> ( $\alpha = 0.05$ ) + $\rho$	teeth-type	No	0.998	0.002	0.000	0.000	0.000	$20.30 \pm 0.14$	0.000
robseg(Huber) <sup>t</sup> + $\rho$	teeth-type	No	0.076	0.096	0.263	0.227	0.338	$8.42 \pm 0.18$	0.023
robseg(biweight) <sup>t</sup> + $\rho$	teeth-type	No	0.435	0.122	0.348	0.082	0.013	$8.86 \pm 0.21$	0.189
ES	teeth-type	Yes	0.383	0.082	<b>0.303</b>	0.151	0.081	$9.38 \pm 0.25$	-
CL	teeth-type	Yes	0.500	0.048	0.169	0.128	0.155	$12.42 \pm 0.42$	0.534
SMUCE <sup>t</sup> + $\rho$	teeth-type	Yes	1.000	0.000	0.000	0.000	0.000	$22.02 \pm 0.14$	0.000
CBS <sup>t</sup> + $\rho$	teeth-type	Yes	1.000	0.000	0.000	0.000	0.000	$24.43 \pm 0.07$	0.001
cumSeg <sup>t</sup> + $\rho$	teeth-type	Yes	1.000	0.000	0.000	0.000	0.000	$24.49 \pm 0.06$	0.000
PELT <sup>t</sup> + $\rho$	teeth-type	Yes	0.090	0.069	0.131	0.162	0.548	$10.48 \pm 0.18$	0.023
WBS sSIC <sup>t</sup> + $\rho$	teeth-type	Yes	0.908	0.014	0.017	0.024	0.037	$21.82 \pm 0.32$	0.008
FDR <sup>t</sup> ( $\alpha = 0.05$ ) + $\rho$	teeth-type	Yes	1.000	0.000	0.000	0.000	0.000	$22.02 \pm 0.14$	0.001
robseg(Huber) <sup>t</sup> + $\rho$	teeth-type	Yes	0.090	0.074	0.202	0.222	0.412	$9.37 \pm 0.18$	0.105
robseg(biweight) <sup>t</sup> + $\rho$	teeth-type	Yes	0.456	0.106	<b>0.316</b>	0.105	0.017	$9.48 \pm 0.20$	0.328
ES	stairs-type	No	0.000	0.000	<b>0.923</b>	0.067	0.010	$2.09 \pm 0.08$	-
CL	stairs-type	No	0.000	0.000	<b>0.907</b>	0.075	0.018	$2.10 \pm 0.08$	0.977
SMUCE <sup>t</sup> + $\rho$	stairs-type	No	0.000	0.008	0.489	0.225	0.278	$4.28 \pm 0.19$	0.003
CBS <sup>t</sup> + $\rho$	stairs-type	No	0.006	0.134	0.594	0.193	0.073	$6.27 \pm 0.31$	0.000
cumSeg <sup>t</sup> + $\rho$	stairs-type	No	0.002	0.120	0.682	0.192	0.004	$7.54 \pm 0.32$	0.000
PELT <sup>t</sup> + $\rho$	stairs-type	No	0.000	0.000	0.032	0.041	0.927	$6.56 \pm 0.18$	0.000
WBS sSIC <sup>t</sup> + $\rho$	stairs-type	No	0.000	0.003	0.456	0.094	0.447	$4.63 \pm 0.17$	0.001
FDR <sup>t</sup> ( $\alpha = 0.05$ ) + $\rho$	stairs-type	No	0.000	0.008	0.489	0.225	0.278	$4.27 \pm 0.19$	0.002
robseg(Huber) <sup>t</sup> + $\rho$	stairs-type	No	0.000	0.000	0.207	0.144	0.649	$4.84 \pm 0.16$	0.001
robseg(biweight) <sup>t</sup> + $\rho$	stairs-type	No	0.000	0.000	0.699	0.183	0.118	$3.21 \pm 0.12$	0.016

TABLE 4. Frequencies of  $\widehat{N} - N$  and  $\widehat{R}_n(\cdot)$  of ES and its competitors for exponential model over 1000 simulated sample paths. Contribution denotes the frequency of each competitor being selected as ES. Bold: highest empirical frequency of  $\widehat{N} - N = 0$  and those with frequencies within 10% off the highest. The uncertainty is obtained by computing  $2\widehat{\sigma}/\sqrt{n_r}$ , where  $\widehat{\sigma}^2$  is the empirical variance and  $n_r$  is the number of replications.

## 6. REAL DATA EXAMPLES

In this section, we apply our estimator selection procedure to two real datasets and investigate its performance. The first one is the observations of DNA copy numbers from biological research where Gaussian model is considered to detect changes. The second one is the British coal disasters dataset to which Poisson model is applied.

**6.1. Detecting changes in DNA copy numbers.** In normal human cells, it is well known that the number of DNA copies is two. As it has been revealed by many works in biological research (see [Albertson and Pinkel \(2003\)](#) and [Redon et al. \(2006\)](#) for example), the pathogenesis of some diseases including various cancers and mental retardation is often associated



to chromosomal aberrations such as deletions, duplications and/or amplifications which finally result in the copy number of DNA from such regions differs from the normal number two. Including microarray and sequencing experiments, biologists have developed various techniques to measure DNA copy numbers of the selected genes on some genome and they record their experimental results as a sequence of observations  $\mathbf{Y} = (Y_1, \dots, Y_n)$ . The statistical interest lies in finding abrupt changes in the means of the observations. To address this issue, we consider Gaussian model with an estimated variance.

In R package `jointseg` (Pierre-Jean et al. (2015)), they provide two real datasets GSE11976 and GSE29172 to resample from, where the truth of changepoints is already known. However, since we do not have the information about the true value of  $\gamma^*$  on each segment, it is impossible to compute the pseudo Hellinger distance between each estimator and the truth. Note that for both GSE11976 and GSE29172 datasets, we need to choose the tumour fraction when resampling from them. We consider the tumour fraction levels 0.79 and 1 for the dataset GSE11976 and the levels 0.7 and 1 for GSE29172 which turns out to be the situations where the size of each jump at the changepoint is relatively large as indicated in Figure 9 of Fearnhead and Rigaiil (2019). Therefore, we can roughly evaluate the performance of each estimator by its frequency of correctly estimating the number of changepoints. Although our selection procedure can be applied in the scenario where small amount of outliers present in the observations, as we have seen in Section 5 some candidates in  $\hat{\Gamma}$  are sensitive to the outliers. To avoid the phenomenon that an estimator systematically underestimates the number of changepoints but due to the sensitivity to outliers it accidentally gives a correct number of segments, we run a smooth procedure on the data before applying all the estimation procedures by implementing the function `smooth.CNA` from the famous R package `DNAcopy`. Moreover, since we have seen in the simulation study that the performance of CBS and `cumSeg` is quite poor, we remove these two estimators from our candidates set  $\hat{\Gamma}$  for simplicity. For each dataset and each level of tumour fraction, we simulate 1000 profiles of length  $n = 1000$  with 5 changepoints where the length of each segment is at least 20. The results are shown in Table 5. As one can observe, among the state-of-art ones, `robseg` (biweight) is the best for correctly estimating the number of changepoints on this dataset. By running a data-driven procedure to select among the candidates set  $\hat{\Gamma}$ , our selected estimator ES shows a competitive performance in this situation as compared to the best one `robseg` (biweight).

**6.2. British coal disasters dataset.** To investigate the performance of ES for Poisson model in practice, we apply our procedure to British coal disasters dataset. This dataset is quite well-known in the context of Poisson segmentation see Green (1995), Yang and Kuo (2001), Fearnhead (2006)

Method	Dataset	Fraction	$\widehat{N} - N$							Contribution
			$\leq -3$	-2	-1	0	1	2	$\geq 3$	
ES	GSE11976	0.79	0.003	0.028	0.044	<b>0.771</b>	0.108	0.031	0.015	-
PELT	GSE11976	0.79	0.000	0.002	0.008	0.198	0.096	0.196	0.500	0.060
SMUCE	GSE11976	0.79	0.004	0.021	0.124	0.391	0.203	0.139	0.118	0.147
CL	GSE11976	0.79	0.011	0.066	0.053	0.550	0.117	0.118	0.085	0.393
WBS sSIC	GSE11976	0.79	0.005	0.031	0.066	0.508	0.066	0.174	0.150	0.100
FDR( $\alpha = 0.05$ )	GSE11976	0.79	0.000	0.005	0.011	0.096	0.056	0.126	0.706	0.020
robseg(Huber)	GSE11976	0.79	0.001	0.012	0.022	0.569	0.193	0.110	0.093	0.121
robseg(biweight)	GSE11976	0.79	0.002	0.046	0.045	<b>0.778</b>	0.102	0.019	0.008	0.159
ES	GSE11976	1.00	0.000	0.003	0.007	<b>0.790</b>	0.100	0.046	0.054	-
PELT	GSE11976	1.00	0.000	0.000	0.000	0.243	0.067	0.195	0.495	0.046
SMUCE	GSE11976	1.00	0.000	0.002	0.035	0.395	0.178	0.177	0.213	0.208
CL	GSE11976	1.00	0.001	0.011	0.011	0.604	0.098	0.170	0.105	0.357
WBS sSIC	GSE11976	1.00	0.000	0.003	0.008	0.536	0.060	0.225	0.168	0.075
FDR( $\alpha = 0.05$ )	GSE11976	1.00	0.000	0.000	0.004	0.138	0.059	0.126	0.673	0.018
robseg(Huber)	GSE11976	1.00	0.000	0.002	0.004	0.559	0.163	0.126	0.146	0.155
robseg(biweight)	GSE11976	1.00	0.000	0.010	0.006	<b>0.794</b>	0.101	0.043	0.046	0.141
ES	GSE29172	0.70	0.014	0.136	0.133	<b>0.596</b>	0.088	0.028	0.005	-
PELT	GSE29172	0.70	0.003	0.027	0.054	0.210	0.139	0.181	0.386	0.089
SMUCE	GSE29172	0.70	0.016	0.112	0.307	0.247	0.176	0.087	0.055	0.099
CL	GSE29172	0.70	0.035	0.159	0.155	0.305	0.129	0.126	0.091	0.302
WBS sSIC	GSE29172	0.70	0.022	0.105	0.155	0.290	0.113	0.173	0.142	0.046
FDR( $\alpha = 0.05$ )	GSE29172	0.70	0.003	0.024	0.075	0.133	0.112	0.133	0.520	0.032
robseg(Huber)	GSE29172	0.70	0.007	0.068	0.087	0.533	0.163	0.092	0.050	0.224
robseg(biweight)	GSE29172	0.70	0.018	0.168	0.153	<b>0.597</b>	0.052	0.012	0.000	0.208
ES	GSE29172	1.00	0.000	0.005	0.003	<b>0.828</b>	0.093	0.051	0.020	-
PELT	GSE29172	1.00	0.000	0.001	0.001	0.233	0.070	0.251	0.444	0.046
SMUCE	GSE29172	1.00	0.000	0.004	0.044	0.416	0.193	0.199	0.144	0.185
CL	GSE29172	1.00	0.001	0.009	0.006	0.684	0.077	0.163	0.060	0.427
WBS sSIC	GSE29172	1.00	0.000	0.006	0.009	0.576	0.051	0.230	0.128	0.070
FDR( $\alpha = 0.05$ )	GSE29172	1.00	0.000	0.001	0.002	0.119	0.063	0.133	0.682	0.018
robseg(Huber)	GSE29172	1.00	0.000	0.001	0.001	0.594	0.145	0.158	0.101	0.120
robseg(biweight)	GSE29172	1.00	0.000	0.007	0.006	<b>0.833</b>	0.098	0.043	0.013	0.134

TABLE 5. Frequencies of  $\widehat{N} - N$  of ES and its competitors for DNA copy numbers data. Contribution denotes the frequency of each competitor being selected as ES. Bold: highest empirical frequency of  $\widehat{N} - N = 0$  and those with frequencies within 10% off the highest.

and [Lloyd et al. \(2015\)](#) for example. We choose this dataset mainly because of two reasons. First, the changepoints have been studied by many different methods which makes it easier to understand our result. Besides, the sequence has a general tendency to decrease with the progress over time which can be correlated to implementing safety regulation in the history. Though pretty rough, we have some evidence to evaluate the changepoint detection procedures on this dataset.

The data at hand include the number of each year coal disasters in UK during the period from March 15th, 1851 to March 22nd, 1962 with length  $n = 112$ . In this situation, to detect changes along the sequence, Poisson model is considered together with the candidates set (11) described in Section 5.3. We conclude the results of different estimators as follows. Concerning to the changepoints, there are in total three suggestions:

- (1) 1 changepoint at the year 1891:  $\text{cumSeg}^t + \rho$ ,  $\text{PELT}^t + \rho$ ,  $\text{WBS sSIC}^t + \rho$ ,  $\text{FDR}^t(\alpha = 0.05) + \rho$  and  $\text{robseg}(\text{biweight})^t + \rho$ ;
- (2) 2 changepoints at the year 1891 and 1947: SMUCE and CL;
- (3) 3 changepoints at the year 1891, 1929 and 1942:  $\text{robseg}(\text{Huber})^t + \rho$ .

Our selection procedure finally choose SMUCE as ES, i.e. we support the suggestion with two changepoints at the year 1891 and 1947. The dataset as well as the result of ES (SMUCE) is plotted in Figure 5.

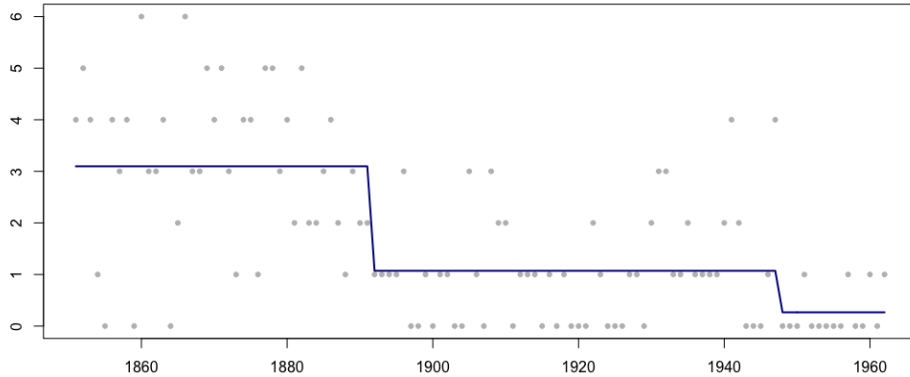


FIGURE 5. Coal mining disasters data (dots) and ES estimator (solid line).

Now we comment our result by comparing it with the existing ones in the literature. In [Green \(1995\)](#), they used the coal mining disasters data recorded per day and proposed a reversible jump MCMC approach to detect changepoints and estimate the intensity function. According to the Figure 2 in the same paper, the model with two changepoints has the highest posterior probability. Moreover, according to their Figure 3, in the two changepoints scenario, the posterior mode is approximately 14,000 days for the first changepoint and 35,000 days for the second one. This is very close to our result since counting from March 15th, 1851, 14,000 days is between the year 1889 and 1890 and 35,000 days is the time between the year 1946 and 1947. Later, a Bayesian binary segmentation procedure was proposed by [Yang and Kuo \(2001\)](#) to locate changepoints for Poisson process. Based on two different tests they adopted, their procedure obtained two different sets of changepoints (one changepoint for applying Bayes factor criterion and two for applying BIC approximation criterion) where the locations of changepoints for these two models are quite similar to the results (1) and (2) mentioned in the last paragraph. On the other hand, as it was pointed out in [Lloyd et al. \(2015\)](#), UK parliament passed several acts to improve the safety of mine works including the Coal Mines Regulation Acts of 1872 and 1887 and a further one in 1954 with mines and quarries acts. In general, it is reasonable to have a non-increasing expectation of the number of disasters after the year releasing these regulations. As it is shown in Figure 5, the

model with two changepoints meets the releasing regulation years 1887 and 1954. Considering the best fit with the time of released regulations and the results given in the literatures, we believe the two changepoints model for this dataset is the most reasonable one to the truth.

## 7. PROOFS

We first introduce some notations for later use. Recall that  $(\mathcal{X}, \mathcal{X}) = (\mathcal{W} \times \mathcal{Y}, \mathcal{W} \otimes \mathcal{Y})$ . We denote  $\mathcal{P}$  the set of all product probabilities on  $(\mathcal{X}^n, \mathcal{X}^{\otimes n})$ . For all  $i \in \{1, \dots, n\}$ , we denote the true distribution of  $X_i = (W_i, Y_i)$  by  $P_i^*$  and denote the true joint distribution of  $\mathbf{X} = (X_1, \dots, X_n)$  by  $\mathbf{P}^* = \otimes_{i=1}^n P_i^* \in \mathcal{P}$ . We denote  $\mathbf{P}_\gamma = \otimes_{i=1}^n P_{i,\gamma}$  the joint distribution of independent random variables  $(W_1, Y_1), \dots, (W_n, Y_n)$  for which the conditional distribution of  $Y_i$  given  $W_i = w_i$  is given by  $R_{\gamma(w_i)} \in \mathcal{Q}_0$  for all  $i \in \{1, \dots, n\}$ . Under such a notation setting, we have  $P_i^* = R_i^* \cdot P_{W_i}$ ,  $P_{i,\gamma} = R_\gamma \cdot P_{W_i}$  as well as the following equality

$$(12) \quad h^2(P_i^*, P_{i,\gamma}) = \int_{\mathcal{W}} h^2(R_i^*(w), R_{\gamma(w)}) dP_{W_i}(w).$$

We define the pseudo Hellinger distance  $\mathbf{h}$  between two probabilities  $\mathbf{P} = \otimes_{i=1}^n P_i$  and  $\mathbf{P}' = \otimes_{i=1}^n P'_i$  by

$$(13) \quad \mathbf{h}^2(\mathbf{P}, \mathbf{P}') = \sum_{i=1}^n h^2(P_i, P'_i).$$

As an immediate consequence of (12) and (13), for any  $\gamma \in \tilde{\Gamma}$ ,

$$(14) \quad \begin{aligned} \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_\gamma) &= \sum_{i=1}^n \int_{\mathcal{W}} h^2(R_i^*(w), R_{\gamma(w)}) dP_{W_i}(w) \\ &= \sum_{i=1}^n h^2(P_i^*, P_{i,\gamma}) = \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_\gamma). \end{aligned}$$

For each  $m \in \mathcal{M}$ , we define the set of probabilities  $\mathcal{P}_m = \{\mathbf{P}_\gamma, \gamma \in \Gamma_m\}$  and  $\tilde{\mathcal{P}} = \{\mathbf{P}_\gamma, \gamma \in \tilde{\Gamma}\}$  with  $\tilde{\Gamma} = \cup_{m \in \mathcal{M}} \Gamma_m$ . For any  $y > 0$ ,  $\mathbf{P}^* \in \mathcal{P}$ ,  $\mathcal{P}_{m_1}, \mathcal{P}_{m_2}$  with  $m_1, m_2 \in \mathcal{M}$ , we define the set

$$\begin{aligned} &\mathcal{B}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*, y) \\ &= \{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \mid \mathbf{P}_{\gamma_1} \in \mathcal{P}_{m_1}, \mathbf{P}_{\gamma_2} \in \mathcal{P}_{m_2}, \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2}) < y^2\} \end{aligned}$$

and for any  $\gamma_1, \gamma_2 \in \tilde{\Gamma}$ , we set

$$\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2) = \mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2) - \mathbb{E}[\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2)].$$

We then introduce below Proposition 45 of [Baraud et al. \(2017\)](#) which is an extensional version of Talagrand's Theorem on the supremum of empirical processes proved in [Massart \(2007\)](#).

**Proposition 1.** *Let  $T$  be some finite or countable set,  $U_1, \dots, U_n$  be independent centered random vectors with values in  $\mathbb{R}^T$  and let*

$$Z = \sup_{t \in T} \left| \sum_{i=1}^n U_{i,t} \right|.$$

*If for some positive numbers  $b$  and  $v$ ,*

$$\max_{i=1, \dots, n} |U_{i,t}| \leq b \quad \text{and} \quad \sum_{i=1}^n \mathbb{E}[U_{i,t}^2] \leq v^2 \quad \text{for all } t \in T,$$

*then, for all positive numbers  $c$  and  $x$ ,*

$$\mathbb{P}[Z \leq (1+c)\mathbb{E}(Z) + (8b)^{-1}cv^2 + 2(1+8c^{-1})bx] \geq 1 - e^{-x}.$$

**7.1. Elementary results and proofs.** Before proving our main theorem, we first present two preliminary results and their proofs in this section.

**Lemma 1.** *Let  $m_1, m_2 \in \mathcal{M}$  be any two partitions on  $\mathcal{W}$ . The class of functions*

$$\mathcal{F}(m_1, m_2) = \left\{ \frac{r_{\gamma_2}}{r_{\gamma_1}} : (w, y) \mapsto \frac{r_{\gamma_2(w)}(y)}{r_{\gamma_1(w)}(y)}, \gamma_1 \in \mathbf{\Gamma}_{m_1}, \gamma_2 \in \mathbf{\Gamma}_{m_2} \right\}$$

*on  $\mathcal{X} = \mathcal{W} \times \mathcal{Y}$  is a VC-subgraph class with dimension not larger than  $2|m_1 \vee m_2| + 1$ .*

*Proof.* For any  $\gamma_1 \in \mathbf{\Gamma}_{m_1}$  and  $\gamma_2 \in \mathbf{\Gamma}_{m_2}$ , we define function  $g_{\gamma_1, \gamma_2}$  on  $\mathcal{W} \times \mathcal{Y}$  as

$$g_{\gamma_1, \gamma_2}(w, y) = T(y) [u(\gamma_2(w)) - u(\gamma_1(w))] - [A(\gamma_2(w)) - A(\gamma_1(w))]$$

and define  $\mathcal{G}(m_1, m_2)$  the class of functions as

$$\mathcal{G}(m_1, m_2) = \{g_{\gamma_1, \gamma_2} \mid \gamma_1 \in \mathbf{\Gamma}_{m_1}, \gamma_2 \in \mathbf{\Gamma}_{m_2}\}.$$

With the fact that  $\mathcal{F}(m_1, m_2) = \{e^g, g \in \mathcal{G}(m_1, m_2)\}$  and the exponential function is monotone on  $\mathbb{R}$ , by [Baraud et al. \(2017\)](#)[Proposition 42-(ii)], it is enough to prove the conclusion holds for the class  $\mathcal{G}(m_1, m_2)$ .

Let  $K = |m_1 \vee m_2|$  be the number of segments given by the refined partition  $m_1 \vee m_2$  and  $\mathcal{I}_1, \dots, \mathcal{I}_K$  the resulted segments on  $\mathcal{W}$ . For any  $\gamma_1 \in \mathbf{\Gamma}_{m_1}$ , we can rewrite it as

$$\gamma_1(w) = \sum_{k=1}^K a_k \mathbb{1}_{\mathcal{I}_k}(w), \quad \text{where } (a_1, \dots, a_K) \in I^K$$

and any  $\gamma_2 \in \mathbf{\Gamma}_{m_2}$ ,

$$\gamma_2(w) = \sum_{k=1}^K b_k \mathbb{1}_{\mathcal{I}_k}(w), \quad \text{where } (b_1, \dots, b_K) \in I^K.$$

As an immediate consequence, for any  $g_{\gamma_1, \gamma_2} \in \mathcal{G}(m_1, m_2)$ , it can be rewritten as

$$g_{\gamma_1, \gamma_2}(w, y) = \sum_{k=1}^K [u(b_k) - u(a_k)] \mathbb{1}_{\mathcal{I}_k}(w) T(y) - \sum_{k=1}^K [A(b_k) - A(a_k)] \mathbb{1}_{\mathcal{I}_k}(w).$$

Therefore,  $\mathcal{G}(m_1, m_2)$  is contained in a  $2K$ -dimensional vector space spanned by  $\{\mathbb{1}_{\mathcal{I}_k}(w), T(y) \mathbb{1}_{\mathcal{I}_k}(w), k = 1, \dots, K\}$ . By Lemma 2.6.15 of [van der Vaart and Wellner \(1996\)](#), we conclude  $\mathcal{G}(m_1, m_2)$  is VC-subgraph on  $\mathcal{X} = \mathcal{W} \times \mathcal{Y}$  with dimension not larger than  $2K + 1$ .  $\square$

**Proposition 2.** *Let  $m_1, m_2 \in \mathcal{M}$  be any two partitions on  $\mathcal{W}$ . Under Assumption 2, for any  $\mathbf{P}^* \in \mathcal{P}$ ,  $\eta \geq 1$  and any  $y > 0$  satisfying*

$$y^2 \geq \eta [D_n(m_1) + D_n(m_2)],$$

we have

$$\begin{aligned} \mathbb{E} \left[ \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{P}^{m_1} \times \mathcal{P}^{m_2}(\mathbf{P}^*, y)} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \right] \\ \leq \left[ 9.77 \sqrt{\frac{2\alpha + 1/2}{\eta}} + \frac{90(2\alpha + 1/2)}{\eta} \right] y^2. \end{aligned}$$

*Proof.* We set  $\boldsymbol{\tau} = \otimes_{i=1}^n \tau_i$  with  $\tau_i = P_{W_i} \otimes \nu$ . For any  $\boldsymbol{\gamma} \in \tilde{\boldsymbol{\Gamma}}$ , we denote  $\mathbf{r}_\boldsymbol{\gamma}$  a density on  $\mathcal{X}^n = (\mathcal{W} \times \mathcal{Y})^n$  as

$$\mathbf{r}_\boldsymbol{\gamma}(x_1, \dots, x_n) = r_\gamma(x_1) \cdots r_\gamma(x_n), \quad \text{for all } (x_1, \dots, x_n) \in \mathcal{X}^n$$

so that for any  $\boldsymbol{\gamma} \in \tilde{\boldsymbol{\Gamma}}$ , we have  $\mathbf{P}_\boldsymbol{\gamma} = \mathbf{r}_\boldsymbol{\gamma} \cdot \boldsymbol{\tau}$ . For any  $y > 0$ , we define  $\mathcal{F}_y(m_1, m_2)$  the class of functions on  $\mathcal{X}$  as

$$\left\{ \psi \left( \sqrt{\frac{r_{\gamma_2}}{r_{\gamma_1}}} \right) \mid \gamma_1 \in \boldsymbol{\Gamma}_{m_1}, \gamma_2 \in \boldsymbol{\Gamma}_{m_2}, \mathbf{h}^2(\mathbf{P}^*, \mathbf{r}_{\gamma_1} \cdot \boldsymbol{\tau}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{r}_{\gamma_2} \cdot \boldsymbol{\tau}) < y^2 \right\}.$$

Since  $\mathcal{F}_y(m_1, m_2)$  is a subset of the collection

$$\left\{ \psi \left( \sqrt{\frac{r_{\gamma_2}}{r_{\gamma_1}}} \right) \mid \gamma_1 \in \boldsymbol{\Gamma}_{m_1}, \gamma_2 \in \boldsymbol{\Gamma}_{m_2} \right\}$$

and the function  $\psi$  is monotone, it follows from Lemma 1 and Proposition 42-(ii) of [Baraud et al. \(2017\)](#) that  $\mathcal{F}_y(m_1, m_2)$  is VC-subgraph on  $\mathcal{X}$  with dimension not larger than  $\bar{V} = 2|m_1 \vee m_2| + 1$ . Besides, by Proposition 3 of [Baraud and Birgé \(2018\)](#), our choice of the function  $\psi$  satisfies their Assumption 2 and more precisely (11) in their paper with  $a_2^2 = 3\sqrt{2}$  so that for any  $y > 0$ ,

$$(15) \quad \sup_{f \in \mathcal{F}_y(m_1, m_2)} n^{-1} \sum_{i=1}^n \mathbb{E} [f^2(X_i)] \leq \frac{a_2^2 y^2}{n}.$$

Moreover, since the function  $\psi$  takes values in  $[-1, 1]$ , we derive from (15) that

$$\sup_{f \in \mathcal{F}_y(m_1, m_2)} n^{-1} \sum_{i=1}^n \mathbb{E} [f^2(X_i)] \leq \left( \frac{a_2^2 y^2}{n} \right) \wedge 1 \leq 1.$$

To bound the expectation of the supremum of an empirical process over a VC-subgraph class, we apply Theorem 2 of [Baraud and Chen \(2020\)](#) to  $\mathcal{F}_y(m_1, m_2)$  and obtain

$$\begin{aligned} & \mathbb{E} \left[ \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1}} \times \mathcal{P}_{m_2}(\mathbf{P}^*, y)} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \right] \\ &= \mathbb{E} \left[ \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1}} \times \mathcal{P}_{m_2}(\mathbf{P}^*, y)} |\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2) - \mathbb{E}[\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2)]| \right] \\ &= \mathbb{E} \left[ \sup_{f \in \mathcal{F}_y(m_1, m_2)} \left| \sum_{i=1}^n (f(X_i) - \mathbb{E}[f(X_i)]) \right| \right] \\ (16) \quad & \leq 9.77y \sqrt{\bar{V} L_n(y)} + 90\bar{V} L_n(y), \end{aligned}$$

where  $L_n(y) = 9.11 + \log_+ [n / (3\sqrt{2}y^2)]$ . Under Assumption 2, there exists a constant  $\alpha \geq 1$  such that

$$(17) \quad \bar{V} = 2|m_1 \vee m_2| + 1 \leq 2\alpha(|m_1| + |m_2|) + 1 \leq \left(2\alpha + \frac{1}{2}\right) (|m_1| + |m_2|).$$

Therefore, combining (16) and (17), we obtain

$$\begin{aligned} & \mathbb{E} \left[ \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1}} \times \mathcal{P}_{m_2}(\mathbf{P}^*, y)} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \right] \\ (18) \quad & \leq 9.77y \sqrt{\left(2\alpha + \frac{1}{2}\right) (|m_1| + |m_2|) L_n(y)} + 90 \left(2\alpha + \frac{1}{2}\right) (|m_1| + |m_2|) L_n(y). \end{aligned}$$

Recall that  $D_n(m) = |m| [9.11 + \log_+(n/|m|)]$ . For any  $\eta \geq 1$ , provided  $y^2 \geq \eta [D_n(m_1) + D_n(m_2)]$ , on the one hand, we have

$$\begin{aligned} & y^2 \geq \eta |m_1| \left( 9.11 + \log_+ \left( \frac{n}{|m_1| + |m_2|} \right) \right) \\ & \quad + \eta |m_2| \left( 9.11 + \log_+ \left( \frac{n}{|m_1| + |m_2|} \right) \right) \\ (19) \quad & = \eta (|m_1| + |m_2|) \left[ 9.11 + \log_+ \left( \frac{n}{|m_1| + |m_2|} \right) \right]. \end{aligned}$$

On the other hand, (19) also implies  $y^2 \geq |m_1| + |m_2|$ . Therefore,

$$\begin{aligned} L_n(y) &= 9.11 + \log_+ \left( \frac{n}{3\sqrt{2}y^2} \right) \leq 9.11 + \log_+ \left[ \frac{n}{3\sqrt{2}(|m_1| + |m_2|)} \right] \\ (20) \quad &\leq 9.11 + \log_+ \left( \frac{n}{|m_1| + |m_2|} \right). \end{aligned}$$

Plugging (19) and (20) into (18), we complete the proof.  $\square$

**7.2. Proof of Theorem 1.** The proof of Theorem 1 is inspired by the proof of Theorem A.1 in Baraud and Birgé (2018). Before we start to prove Theorem 1, we first show the following result.

**Proposition 3.** *Let numbers  $a, \eta \geq 1$  and  $\delta, \vartheta > 1$  such that*

$$(21) \quad 2 \exp(-\vartheta) + \sum_{j=1}^{+\infty} \exp(-\vartheta \delta^j) \leq 1.$$

*Under Assumption 1 and 2, for any  $\xi > 0$  and for all  $m_1, m_2 \in \mathcal{M}$  simultaneously, with probability at least  $1 - \Sigma^2 e^{-\xi}$ ,*

$$\begin{aligned} &\sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{D}_{m_1} \times \mathcal{D}_{m_2}} [|\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| - k_1 [\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2})]] \\ &\leq k_0 a \{ \eta [D_n(m_1) + D_n(m_2)] \vee (\Delta(m_1) + \Delta(m_2) + \vartheta + \xi) \}, \end{aligned}$$

where

$$\begin{aligned} k_0 &= 16 \sqrt{\frac{9.77 \sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} + \frac{3\sqrt{2}}{16}}{2a}} + \frac{4}{a} \\ &\quad + \left( 9.77 \sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} \right), \\ k_1 &= 16 \sqrt{\frac{\delta \left( 9.77 \sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} + \frac{3\sqrt{2}}{16} \right)}{2a}} + \frac{4}{a} \\ &\quad + \left( 9.77 \sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} \right) \delta. \end{aligned}$$

*Proof.* Let  $\xi > 0$ ,  $\delta, \vartheta > 1$ ,  $a, \eta \geq 1$  and  $m_1, m_2 \in \mathcal{M}$  be fixed. For each  $j \in \mathbb{N}$ , we set

$$\begin{aligned} x_0(m_1, m_2) &= \eta (D_n(m_1) + D_n(m_2)) \vee (\Delta(m_1) + \Delta(m_2) + \vartheta + \xi), \\ x_j(m_1, m_2) &= \delta^j x_0(m_1, m_2), \quad y_j^2(m_1, m_2) = a x_j(m_1, m_2). \end{aligned}$$



For each  $j \in \mathbb{N}$ , we define the set

$$\begin{aligned} & \mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*) \\ &= \{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{P}_{m_1} \times \mathcal{P}_{m_2} \mid y_j^2 \leq \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2}) < y_{j+1}^2\} \end{aligned}$$

and set

$$\mathcal{Z}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X}) = \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)|.$$

For simplifying the notations, let us drop the dependency of  $x_j$  and  $y_j$  with respect to  $m_1, m_2$  for a while. Since  $\mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*) \subset \mathcal{B}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*, y_{j+1})$  and  $y_{j+1}^2 > y_0^2 = ax_0 \geq \eta [D_n(m_1) + D_n(m_2)]$ , under Assumption 2, applying Proposition 2 yields,

$$\begin{aligned} \mathbb{E} \left[ \mathcal{Z}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X}) \right] &= \mathbb{E} \left[ \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \right] \\ &\leq \mathbb{E} \left[ \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*, y_{j+1})} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \right] \\ (22) \quad &\leq \left( 9.77 \sqrt{\frac{2\alpha + 1/2}{\eta}} + \frac{90(2\alpha + 1/2)}{\eta} \right) y_{j+1}^2. \end{aligned}$$

For  $i \in \{1, \dots, n\}$ , we set

$$(23) \quad U_{i, (r_{\gamma_1}, r_{\gamma_2})} = \psi \left( \sqrt{\frac{r_{\gamma_2}(W_i)(Y_i)}{r_{\gamma_1}(W_i)(Y_i)}} \right) - \mathbb{E} \left[ \psi \left( \sqrt{\frac{r_{\gamma_2}(W_i)(Y_i)}{r_{\gamma_1}(W_i)(Y_i)}} \right) \right].$$

With the fact that  $\psi$  takes values in  $[-1, 1]$ , it is easy to observe that

$$\max_{i=1, \dots, n} |U_{i, (r_{\gamma_1}, r_{\gamma_2})}| \leq 2.$$

Moreover,  $\psi$  satisfies the Assumption 2 more precisely (11) in [Baraud and Birgé \(2018\)](#) with  $a_2^2 = 3\sqrt{2}$ , we derive for each  $j \in \mathbb{N}$ , all  $\gamma_1 \in \Gamma_{m_1}$ ,  $\gamma_2 \in \Gamma_{m_2}$  such that  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)$

$$\sum_{i=1}^n \mathbb{E} [U_{i, (r_{\gamma_1}, r_{\gamma_2})}^2] \leq \sum_{i=1}^n \mathbb{E} \left[ \psi^2 \left( \sqrt{\frac{r_{\gamma_2}(W_i)(Y_i)}{r_{\gamma_1}(W_i)(Y_i)}} \right) \right] \leq 3\sqrt{2}y_{j+1}^2.$$

Then, for each  $j \in \mathbb{N}$ , we can apply Proposition 1 with  $b = 2$ ,  $v^2 = 3\sqrt{2}y_{j+1}^2$  and  $T = \mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)$  and obtain that for all  $c > 0$  and for all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in$

$\mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)$  with probability at least  $1 - e^{-x_j}$ ,

$$\begin{aligned}
|\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| &\leq \mathcal{Z}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X}) \\
&\leq (1+c)\mathbb{E}\left[\mathcal{Z}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X})\right] + \frac{3\sqrt{2}y_{j+1}^2c}{16} + 4\left(1 + \frac{8}{c}\right)x_j \\
&\leq (1+c)\left(9.77\sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta}\right)y_{j+1}^2 \\
&\quad + \frac{3\sqrt{2}y_{j+1}^2c}{16} + 4\left(1 + \frac{8}{c}\right)x_j \\
&\leq (1+c)\left(9.77\sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta}\right)y_{j+1}^2 \\
&\quad + \frac{3\sqrt{2}y_{j+1}^2c}{16} + \frac{4}{a}\left(1 + \frac{8}{c}\right)y_j^2 \\
&\leq (1+c)\left(9.77\sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta}\right)\delta y_j^2 \\
&\quad + \left[\frac{3\sqrt{2}c\delta}{16} + \frac{4}{a}\left(1 + \frac{8}{c}\right)\right]y_j^2.
\end{aligned}$$

Taking

$$c = \sqrt{\frac{32}{\left(9.77\sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} + \frac{3\sqrt{2}}{16}\right)\delta a}}$$

to minimize the bracketed term yields for all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)$ , with probability at least  $1 - e^{-x_j}$

$$|\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \leq k_1 y_j^2.$$

By the definition of  $\mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)$ , we get for all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2})$  belonging to  $\mathcal{B}_j^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*)$ , with probability at least  $1 - e^{-x_j}$ ,

$$|\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \leq k_1 y_j^2 \leq k_1 [\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2})].$$

We define

$$\mathcal{Z}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X}) = \sup_{(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*, y_0)} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)|.$$

With an analogous argument by applying Proposition 1 to  $\mathcal{Z}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X})$  with  $x = x_0$ , we can obtain for all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*, y_0)$  and all

$c > 0$ , with probability at least  $1 - e^{-x_0}$ ,

$$\begin{aligned} |\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| &\leq \mathcal{Z}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X}) \\ &\leq (1+c) \left( 9.77 \sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} \right) y_0^2 \\ &\quad + \left[ \frac{3\sqrt{2}c}{16} + \frac{4}{a} \left( 1 + \frac{8}{c} \right) \right] y_0^2. \end{aligned}$$

To minimize the bracketed term, we take

$$c = \sqrt{\frac{32}{\left( 9.77 \sqrt{\frac{2\alpha+1/2}{\eta}} + \frac{90(2\alpha+1/2)}{\eta} + \frac{3\sqrt{2}}{16} \right) a}}$$

and therefore for all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{B}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{P}^*, y_0)$  with probability at least  $1 - e^{-x_0}$ ,

$$|\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \leq \mathcal{Z}^{\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}}(\mathbf{X}) \leq ak_0x_0.$$

Combining all the bounds together, we derive for all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \mathcal{P}_{m_1} \times \mathcal{P}_{m_2}$  simultaneously with probability at least  $1 - \varepsilon(m_1, m_2)$ ,

$$|\mathbf{Z}(\mathbf{X}, \gamma_1, \gamma_2)| \leq k_1 [\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2})] + ak_0x_0(m_1, m_2),$$

where

$$\varepsilon(m_1, m_2) = 2 \exp[-x_0(m_1, m_2)] + \sum_{j \geq 1} \exp[-x_j(m_1, m_2)].$$

By the definition of  $x_j(m_1, m_2)$ , we notice that for all  $j \in \mathbb{N}$ ,  $x_j(m_1, m_2) \geq \Delta(m_1) + \Delta(m_2) + \vartheta \delta^j + \xi$ . Hence, provided (21), we have

$$\begin{aligned} \varepsilon(m_1, m_2) &\leq \exp[-\xi - \Delta(m_1) - \Delta(m_2)] \left( 2 \exp(-\vartheta) + \sum_{j \geq 1} \exp(-\vartheta \delta^j) \right) \\ &\leq \exp[-\xi - \Delta(m_1) - \Delta(m_2)]. \end{aligned}$$

Finally we can extend this result to all  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2}) \in \widetilde{\mathcal{P}} \times \widetilde{\mathcal{P}}$  by summing these bounds over  $(m_1, m_2) \in \mathcal{M} \times \mathcal{M}$  and using (4).  $\square$

*Proof of Theorem 1.* We apply Proposition 3 with  $\delta = 1.175$ ,  $\vartheta = 1.47$  and as for the values of  $\eta$  and  $a$ , we shall choose them later such that  $k_1 = 3\beta/8$ , with some  $0 < \beta < 1$ . On a set  $\Omega_\xi$  the probability of which is at least

$1 - \Sigma^2 e^{-\xi}$ , for all  $\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2} \in \widetilde{\mathcal{P}}$  and all  $\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}$  containing  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2})$

$$\begin{aligned} \mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2) &\leq \mathbb{E}[\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2)] + \frac{3\beta}{8} [\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2})] \\ &\quad + k_0 a [\eta(D_n(m_1) + D_n(m_2)) \vee (\Delta(m_1) + \Delta(m_2) + \vartheta + \xi)] \\ &\leq \mathbb{E}[\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2)] + \frac{3\beta}{8} [\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2})] \\ &\quad + k_0 a [\eta D_n(m_1) + \eta D_n(m_2) + \Delta(m_1) + \Delta(m_2) + \vartheta + \xi]. \end{aligned}$$

Since the last inequality is true for all the  $\mathcal{P}_{m_1} \times \mathcal{P}_{m_2}$  containing  $(\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2})$ , provided  $C_0(2\alpha + 1/2) \geq k_0 a \eta$ , we derive from (5) that with a probability at least  $1 - \Sigma^2 e^{-\xi}$ ,

$$\begin{aligned} \mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2) &\leq \mathbb{E}[\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2)] + \frac{3\beta}{8} [\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) + \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2})] \\ (24) \quad &\quad + \mathbf{pen}(\gamma_1) + \mathbf{pen}(\gamma_2) + k_0 a(\vartheta + \xi). \end{aligned}$$

According to Proposition 3 of Baraud and Birgé (2018), the function  $\psi$  satisfies Assumption 2 (more precisely (10)) in the same paper with  $a_0 = 4$  and  $a_1 = 3/8$ . As a consequence, for all  $\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2} \in \widetilde{\mathcal{P}}$  and  $\mathbf{P}^* \in \mathcal{P}$ ,

$$(25) \quad \mathbb{E}[\mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2)] \leq 4\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) - \frac{3}{8}\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2}).$$

Combining (24) and (25), we derive that for all  $\mathbf{P}_{\gamma_1}, \mathbf{P}_{\gamma_2} \in \widetilde{\mathcal{P}}$  and  $\mathbf{P}^* \in \mathcal{P}$ , with a probability at least  $1 - \Sigma^2 e^{-\xi}$ ,

$$\begin{aligned} \mathbf{T}(\mathbf{X}, \gamma_1, \gamma_2) &\leq (4 + \frac{3\beta}{8})\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_1}) - \frac{3(1-\beta)}{8}\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\gamma_2}) \\ (26) \quad &\quad + \mathbf{pen}(\gamma_1) + \mathbf{pen}(\gamma_2) + k_0 a(\vartheta + \xi). \end{aligned}$$

This entails that, for any (random) elements  $\mathbf{P}_{\widehat{\gamma}_\lambda}, \mathbf{P}_{\widehat{\gamma}_\lambda} \in \widetilde{\mathcal{P}}$ , on a set  $\Omega_\xi$  with probability at least  $1 - \Sigma^2 e^{-\xi}$

$$\begin{aligned} \mathbf{T}(\mathbf{X}, \widehat{\gamma}_\lambda, \widehat{\gamma}_\lambda) &\leq (4 + \frac{3\beta}{8})\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\widehat{\gamma}_\lambda}) - \frac{3(1-\beta)}{8}\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\widehat{\gamma}_\lambda}) \\ (27) \quad &\quad + \mathbf{pen}(\widehat{\gamma}_\lambda) + \mathbf{pen}(\widehat{\gamma}_\lambda) + k_0 a(\vartheta + \xi) \end{aligned}$$

and

$$\begin{aligned} \mathbf{v}(\mathbf{X}, \widehat{\gamma}_\lambda) &= \sup_{\lambda' \in \Lambda} [\mathbf{T}(\mathbf{X}, \widehat{\gamma}_\lambda, \widehat{\gamma}_{\lambda'}) - \mathbf{pen}(\widehat{\gamma}_{\lambda'})] + \mathbf{pen}(\widehat{\gamma}_\lambda) \\ (28) \quad &\leq (4 + \frac{3\beta}{8})\mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\widehat{\gamma}_\lambda}) - \frac{3(1-\beta)}{8} \inf_{\lambda' \in \Lambda} \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\widehat{\gamma}_{\lambda'}}) \\ &\quad + 2\mathbf{pen}(\widehat{\gamma}_\lambda) + k_0 a(\vartheta + \xi). \end{aligned}$$

By the construction of  $\psi$ ,  $\mathbf{T}(\mathbf{X}, \widehat{\gamma}_\lambda, \widehat{\gamma}_\lambda) = -\mathbf{T}(\mathbf{X}, \widehat{\gamma}_\lambda, \widehat{\gamma}_\lambda)$ . Combining (27), (28) and (6) leads to for any  $\lambda \in \Lambda$ , on a set  $\Omega_\xi$  with probability at least

$$\begin{aligned}
& 1 - \Sigma^2 e^{-\xi} \\
& \frac{3(1-\beta)}{8} \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) \leq (4 + \frac{3\beta}{8}) \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) - \mathbf{T}(\mathbf{X}, \hat{\gamma}_\lambda, \hat{\gamma}_\lambda) \\
& \quad + \mathbf{pen}(\hat{\gamma}_\lambda) + \mathbf{pen}(\hat{\gamma}_\lambda) + k_0 a(\vartheta + \xi) \\
& \leq (4 + \frac{3\beta}{8}) \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) + [\mathbf{T}(\mathbf{X}, \hat{\gamma}_\lambda, \hat{\gamma}_\lambda) - \mathbf{pen}(\hat{\gamma}_\lambda)] \\
& \quad + \mathbf{pen}(\hat{\gamma}_\lambda) + 2 \mathbf{pen}(\hat{\gamma}_\lambda) + k_0 a(\vartheta + \xi) \\
& \leq (4 + \frac{3\beta}{8}) \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) + \mathbf{v}(\mathbf{X}, \hat{\gamma}_\lambda) \\
& \quad + 2 \mathbf{pen}(\hat{\gamma}_\lambda) + k_0 a(\vartheta + \xi) \\
& \leq (4 + \frac{3\beta}{8}) \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) + \mathbf{v}(\mathbf{X}, \hat{\gamma}_\lambda) + 1 \\
(29) \quad & \quad + 2 \mathbf{pen}(\hat{\gamma}_\lambda) + k_0 a(\vartheta + \xi).
\end{aligned}$$

Plugging (28) into (29) yields, for any  $\lambda \in \Lambda$ , on a set  $\Omega_\xi$  with probability at least  $1 - \Sigma^2 e^{-\xi}$ ,

$$\frac{3(1-\beta)}{8} \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) \leq (8 + \frac{3\beta}{4}) \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) + 4 \mathbf{pen}(\hat{\gamma}_\lambda) + 2k_0 a(\vartheta + \xi) + 1.$$

Therefore, for any  $\lambda \in \Lambda$  on a set  $\Omega_\xi$  with probability at least  $1 - \Sigma^2 e^{-\xi}$ ,

$$(30) \quad \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) \leq \frac{64 + 6\beta}{3(1-\beta)} \mathbf{h}^2(\mathbf{P}^*, \mathbf{P}_{\hat{\gamma}_\lambda}) + \frac{32}{3(1-\beta)} \mathbf{pen}(\hat{\gamma}_\lambda) + \frac{16k_0 a(\vartheta + \xi) + 8}{3(1-\beta)}.$$

By the equality (14), we rewrite (30) as the following

$$(31) \quad \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) \leq \frac{64 + 6\beta}{3(1-\beta)} \mathbf{h}^2(\mathbf{R}^*, \mathbf{R}_{\hat{\gamma}_\lambda}) + \frac{32}{3(1-\beta)} \mathbf{pen}(\hat{\gamma}_\lambda) + \frac{16k_0 a(\vartheta + \xi) + 8}{3(1-\beta)}.$$

Taking  $\beta = 0.75$ ,  $\eta \approx 9947.13(2\alpha + 1/2)$ , we can compute the value of  $a \approx 2365.57$  such that  $k_1 = 3\beta/8$  and  $k_0 \approx 0.251$ . Therefore, provided  $C_0 \geq 5.9 \times 10^6$ , plugging the values of  $\beta$ ,  $k_0$ ,  $a$  and  $\vartheta$  into (31), we finally conclude.  $\square$

#### APPENDIX. SIGNALS FOR TESTING POISSON AND EXPONENTIAL MODELS

fms-type (Poisson):  $n = 497$ , changepoints are located at the positions

$$l_0 = \left( \frac{139}{497}, \frac{226}{497}, \frac{243}{497}, \frac{300}{497}, \frac{309}{497}, \frac{333}{497} \right).$$

The Poisson mean on each segment is 4, 6, 10, 3, 7, 1, 5 respectively, i.e.  $\gamma^*$  takes the value  $\log 4$ ,  $\log 6$ ,  $\log 10$ ,  $\log 3$ ,  $\log 7$ ,  $\log 1$ ,  $\log 5$  on each segment. For this signal, we also test the scenario when outliers present in the observations by randomly modifying five points in the observations into 30.

mix-type (Poisson):  $n = 560$  and  $\gamma^*$  is a piecewise constant function on  $[0, 1)$  with 13 changepoints at a sequence of locations

$$l_0 = \left( \frac{11}{560}, \frac{21}{560}, \frac{41}{560}, \frac{61}{560}, \frac{91}{560}, \frac{121}{560}, \frac{161}{560}, \frac{201}{560}, \frac{251}{560}, \frac{301}{560}, \frac{361}{560}, \frac{421}{560}, \frac{491}{560} \right)$$

and on each segment the Poisson mean  $e^{\gamma^*}$  is given by the value 30, 2, 26, 4, 24, 6, 22, 8, 20, 10, 18, 12, 16, 14 respectively.

teeth-type (exponential):  $n = 140$  and  $\gamma^*$  is a piecewise constant function on  $[0, 1)$  with 13 changepoints at a sequence of locations

$$l_0 = \left( \frac{11}{140}, \frac{21}{140}, \frac{31}{140}, \frac{41}{140}, \frac{51}{140}, \frac{61}{140}, \frac{71}{140}, \frac{81}{140}, \frac{91}{140}, \frac{101}{140}, \frac{111}{140}, \frac{121}{140}, \frac{131}{140} \right)$$

and on each segment the value of  $\gamma^*$  is given by 0.5, 5, 0.5, 5, 0.5, 5, 0.5, 5, 0.5, 5, 0.5, 5, 0.5, 5 respectively. For this signal, we also test the scenario when outliers present in the observations by randomly modifying two points in the observations into 20.

stairs-type (exponential):  $n = 500$  and  $\gamma^*$  is a piecewise constant function on  $[0, 1)$  with 4 changepoints at a sequence of locations

$$l_0 = \left( \frac{101}{500}, \frac{201}{500}, \frac{301}{500}, \frac{401}{500} \right)$$

and on each segment the value of  $\gamma^*$  is given by  $2^4$ ,  $2^2$ , 1,  $2^{-2}$ ,  $2^{-4}$  respectively.

#### ACKNOWLEDGEMENTS

The author is grateful to her supervisor Prof. Yannick Baraud for helpful discussions and constructive suggestions.

#### REFERENCES

- Albertson, D. G. and Pinkel, D. (2003). Genomic microarrays in human genetic disease and cancer. *Hum. Mol. Genet.*, **12**, 145–152.
- Antoniadis, A. and Sapatinas, T. (2001). Wavelet shrinkage for natural exponential families with quadratic variance functions. *Biometrika*, **88**, 805–820.
- Antoniadis, A., Besbeas, P. and Sapatinas, T. (2001). Wavelet shrinkage for natural exponential families with cubic variance functions. *Sankhyā: Indian J. Stat., Ser. A*, **63**, 309–327.
- Baraud, Y. and Birgé, L. (2009). Estimating the intensity of a random measure by histogram type estimators. *Probab. Theory Related Fields*, **143**, 239–284.
- Baraud, Y. and Birgé, L. (2018). Rho-estimators revisited: general theory and applications. *Ann. Statist.*, **46**, 3767–3804.

- Baraud, Y., Birgé, L., and Sart, M. (2017). A new method for estimation and model selection:  $\rho$ -estimation. *Invent. Math.*, **207**, 425–517.
- Baraud, Y. and Chen, J. (2020). Robust estimation of a regression function in exponential families. *arXiv preprint*, arXiv:2011.01657.
- Barron, A., Birgé, L. and Massart, P. (1999). Risk bounds for model selection via penalization. *Probab. Theory Related Fields*, **113**, 301–413.
- Birgé, L. and Massart, P. (1997). From model selection to adaptive estimation. In *Festschrift for Lucien Le Cam*, 55–87. Springer, New York.
- Blythe, D. A. J., von Bunau, P., Meinecke, F. C. and Muller, K.-R. (2012). Feature extraction for change-point detection using stationary subspace analysis. *IEEE Trans. Neural Netw. Learn. Syst.*, **23**, 631–643.
- Breiman, L., Friedman, J., Stone, C. J. and Olshen, R. A. (1984). *Classification and Regression Trees*. Taylor & Francis, New York.
- Brown, L. D., Cai, T. T., and Zhou, H. H. (2010). Nonparametric regression in exponential families. *Ann. Statist.*, **38**, 2005–2046.
- Chen, J. (2022). Estimating a regression function in exponential families by model selection. *arXiv preprint*, arXiv:2203.06656.
- Cleynen, A. and Lebarbier, E. (2014). Segmentation of the Poisson and negative binomial rate models: a penalized estimator. *ESAIM Probab. Stat.*, **18**, 750–769.
- Cleynen, A. and Lebarbier, E. (2017). Model selection for the segmentation of multiparameter exponential family distributions. *Electron. J. Stat.*, **11**, 800–842.
- Fearnhead, P. (2006). Exact and efficient Bayesian inference for multiple changepoint. *Statist. Comput.*, **16**, 203–213.
- Fearnhead, P. and Rigai, G. (2019). Changepoint detection in the presence of outliers. *J. Amer. Statist. Assoc.*, **114**, 169–183.
- Fearnhead, P. and Rigai, G. (2020). Relating and comparing methods for detecting changes in mean. *Stat*, e291.
- Frick, K., Munk, A. and Sieling, H. (2013). Multiscale change point inference. *J. Roy. Statist. Soc., Ser. B*, **76**, 495–580.
- Fryźlewicz, P. (2014). Wild binary segmentation for multiple change-point detection. *Ann. Statist.*, **42**, 2243–2281.
- Gallagher, C., Lund, R. and Robbins, M. (2013). Changepoint detection in climate time series with long-term trends. *J. Clim.*, **26**, 4994–5006.
- Green, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, **82**, 711–732.
- Hotz, T., Schütte, O. M., Sieling, H., Polupanow, T., Diederichsen, U., Steinem, C. and Munk, A. (2013). Idealizing ion channel recordings by a jump segmentation multiresolution filter. *IEEE Trans. NanoBioscience*, **12**, 376–386.
- Huang, T., Wu, B., Lizardi, P. and Zhao, H. (2005). Detection of DNA copy number alterations using penalized least squares regression. *Bioinformatics*, **21**, 3811–3817.

- Killick, R., Fearnhead, P. and Eckley, I. A. (2012). Optimal detection of changepoints with a linear computational cost. *J. Amer. Statist. Assoc.*, **107**, 1590–1598.
- Kolaczyk, E. D. and Nowak, R. D. (2005). Multiscale generalised linear models for nonparametric function estimation. *Biometrika*, **92**, 119–133.
- Le Cam, L. (1986). *Asymptotic Methods in Statistical Decision Theory*. Springer Series in Statistics. Springer, New York.
- Le Cam, L. and Yang, G. L. (1990). *Asymptotics in Statistics : Some Basic Concepts*. Springer Series in Statistics. Springer, New York.
- Li, H., Munk, A. and Sieling, H. (2016) FDR-control in multiscale change-point segmentation. *Electron. J. Stat.*, **10**, 918–959.
- Lloyd, C., Gunter, T., Osborne, M. A. and Roberts, S. J. (2015) Variational inference for Gaussian process modulated Poisson processes. In *International Conference on Machine Learning*, **37**, 1814–1822.
- Massart, P. (2007). *Concentration Inequalities and Model Selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003.
- Muggeo, V. M. R. and Adelfio, G. (2010). Efficient change point detection for genomic sequences of continuous measurements. *Bioinformatics*, **27**, 161–166.
- Olshen, A. B., Venkatraman, E., Lucito, R. and Wigler, M. (2004). Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics*, **5**, 557–572.
- Pierre-Jean, M., Rigaiil, G. and Neuviail, P. (2015). Performance evaluation of DNA copy number segmentation methods. *Brief. Bioinformatics*, **16**, 600–615.
- Redon, R., Ishikawa, S., Fitch, K., Feuk, L., Perry, G., Andrews, T., Fiegler, H., Shapero, M., Carson, A., Chen, W., Cho, E., Dallaire, S., Freeman, J., Gonzalez, J., Gratacòs, M., Huang, J., Kalaitzopoulos, D., Komura, D., Macdonald, J. and Hurles, M. (2006). Global variation in copy number in the human genome. *Nature*, **444**, 444–454.
- Reeves, J., Chen, J., Wang, X. L., Lund, R. and Lu, Q. Q. (2007). A review and comparison of changepoint detection techniques for climate data. *J. Appl. Meteorol. and Climatol.*, **46**, 900–915.
- Rigaiil, G. (2015). A pruned dynamic programming algorithm to recover the best segmentations with 1 to  $K_{max}$  change-points. *arXiv preprint*, arXiv:1004.0887.
- Scott, A. J. and Knott, M. (1974). A cluster analysis method for grouping means in the analysis of variance. *Biometrics*, **30**, 507–512.
- Spokoiny, V. (2009). Multiscale local change point detection with applications to value-at-risk. *Ann. Statist.*, **37**, 1405–1436.
- Tibshirani, R. and Wang, P. (2007). Spatial smoothing and hot spot detection for CGH data using the fused Lasso. *Biostatistics*, **9**, 18–29.



- Truong, C., Oudre, L. and Vayatis, N. (2020). Selective review of offline change point detection methods. *Signal Process.*, **167**, 107299.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes. With Applications to Statistics*. Springer Series in Statistics. Springer-Verlag, New York.
- Venkatraman, E. S. and Olshen, A. B. (2007). A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics*, **23**, 657–663.
- Verzelen, N., Fromont, M., Lerasle, M. and Reynaud-Bouret, P. (2020). Optimal change-point detection and localization. *arXiv preprint*, arXiv:2010.11470.
- Wang, D., Yu, Y. and Rinaldo, A. (2020). Univariate mean change point detection: penalization, CUSUM and optimality. *Electron. J. Stat.*, **14**, 1917–1961.
- Yang, T. Y. and Kuo, L. (2001). Bayesian binary segmentation procedure for a Poisson process with multiple changepoints. *J. Comput. Graph. Statist.*, **10**, 772–785.
- Zhang, N. R. and Siegmund, D. O. (2007). A modified Bayes information criterion with applications to the analysis of comparative genomic hybridization data. *Biometrics*, **63**, 22–32.

DEPARTMENT OF MATHEMATICS,  
UNIVERSITY OF LUXEMBOURG  
MAISON DU NOMBRE  
6 AVENUE DE LA FONTE  
L-4364 ESCH-SUR-ALZETTE  
GRAND DUCHY OF LUXEMBOURG  
*Email address:* juntong.chen@uni.lu