

# DUST: A Duality-Based Pruning Method For Exact Multiple Change-Point Detection

Vincent Runge<sup>\*1</sup>, Charles Truong<sup>2</sup>, and Simon Querné<sup>3,4</sup>

<sup>1</sup>Université Paris-Saclay, CNRS, Univ Evry, Laboratoire de Mathématiques et Modélisation d’Evry, 91037, Evry-Courcouronnes, France.

<sup>2</sup>Centre Borelli, Université Paris-Saclay, CNRS, ENS Paris-Saclay, 4 avenue des Sciences, 91190, Gif-sur-Yvette, France

<sup>3</sup>Laboratoire de mathématiques de Versailles, Université Paris-Saclay, UVSQ, CNRS, 45 avenue des États-Unis, 78000, Versailles, France

<sup>4</sup>IFPEN, 1-4 Av. du Bois Préau, 92852, Rueil-Malmaison, France

## Abstract

We tackle the challenge of detecting multiple change points in large time series by optimising a penalised likelihood derived from exponential family models. While dynamic programming algorithms can solve this task exactly with at most quadratic time complexity, recent years have seen the development of pruning strategies to improve their time efficiency. However, the two existing approaches have notable limitations: PELT struggles with pruning efficiency in sparse-change scenarios, while FPOP’s structure is ill-suited for multi-parametric settings. To address these issues, we introduce the DUal Simple Test (DUST) framework, which prunes candidates by evaluating a dual function against a threshold. This approach is highly flexible and broadly applicable to parametric models of any dimension. Under mild assumptions, we establish strong duality for the underlying non-convex pruning problem. We demonstrate DUST’s effectiveness across various change point regimes and models. In particular, for one-parametric models, DUST matches the simplicity of PELT with the efficiency of FPOP. Its use is especially advantageous for non-Gaussian models. Its use is especially advantageous for non-Gaussian models. Finally, we apply DUST to mouse monitoring time series under a change-in-variance model, illustrating its ability to recover the optimal change point structure efficiently.

**Keywords:** Offline change-point detection, dynamic programming, pruning, duality theory, time efficiency

---

<sup>\*</sup>Corresponding author: vincent.runge@univ-evry.fr

# 1 Introduction

Single and multiple change-point detection are well-established unsupervised machine learning tasks within the field of time-series analysis, with foundational work dating back to the 1950s [25, 24]. Over the past decades, the topic has been the subject of extensive research, resulting in numerous monographs [36, 5, 11, 8] and comprehensive review articles [17, 1, 38]. Till recently, the longstanding focus of the scientific community has been the statistical modelling and calibration challenge associated with change-point detection. With the rise of big data, the demand for computationally efficient algorithms has become increasingly pressing. Time efficiency is particularly crucial in many application domains, including genomics [19, 22], econometrics [3, 13], climatology [28, 37], speech processing [14, 7], and network analysis [44, 4], to name just a few.

This work addresses the computational challenge of recovering multiple change points in a time series of fixed length. We consider change-point problems based on optimizing a penalised likelihood whose penalty is proportional to the number of changes. Although algorithms with quadratic time complexity have been available for some time [2, 16], the central objective is to approach quasi-linear execution time as closely as possible, while preserving the exact resolution of the underlying optimisation problem.

Detecting multiple change points requires carefully designed algorithmic strategies often balancing exactness with computational efficiency. The most used approximate algorithm in this category is the quasi-linear binary segmentation (BS) algorithm [34, 35, 40], which is often competitive with exact sub-quadratic complexity algorithms using dynamic programming (DP). We illustrate this with a simple example. With a time budget of 10 seconds, we estimate the maximum data length that an algorithm can segment when the signal consists of 10 segments of equal length. The observations follow  $y_t \sim \mathcal{N}(\mu_t, 1)$ , with piecewise constant means  $\mu_t \in \{0, 1\}$ . Binary Segmentation (BS), stopping at 10 segments, can process up to approximately  $n = 75 \times 10^6$  data points. In contrast, the classical exact Optimal Partitioning (OP) algorithm [16] with BIC penalty [42], which has quadratic time complexity, is limited to around  $n = 120 \times 10^3$ . However, improved exact algorithms such as FPOP can handle up to  $n = 30 \times 10^6$  data points within the same time budget, and the DUST algorithm achieves around  $n = 42 \times 10^6$  data points<sup>1</sup>.

DP methods have experienced a period of renaissance with the development of accelerating pruning strategies, making their execution time competitive with BS on simulations (as just illustrated) and on real data sets (see Figure 6 in [23]). Among them, inequality-based pruning (PELT) [20] and functional-based pruning (FPOP) [23] are the two extreme pruning strategies available on the “pruning scale”. In the latest developments, DP with functional pruning has made possible the inference of complex structured models that constrain the successive segment parameter values (throughout a graph of constraints) [15, 33]. One-parametric model with exponential decays [18] or data modelled with auto-correlation and random drift [31] can also be considered, among others.

---

<sup>1</sup>Simulations were conducted on a MacBook Pro equipped with an Apple M1 chip (8-core CPU: 4 performance cores and four efficiency cores), 16GB of unified memory, running macOS Sequoia 15.5. The code was implemented in R and executed using our dust Rcpp package (for DUST and OP) and the fpop packages provided in [23] (for FPOP and BS).

PELT pruning is efficient for time series with many changes and is simple to code, but does not prune well when there are only a few changes. FPOP pruning is well-suited for problems with a one-parametric functional description. However, extensions to multivariate problems are uneasy as the parameter space description made of intervals in a one-parametric cost function is no longer possible. In this case, the partition of the parameter space is made of potential non-convex and unconnected elements [32]. Approximations of the parameter space with sphere-like or rectangle-like sets have been tested and significantly improve the time complexity, but only for small dimensions in the independent Gaussian multivariate setting [26].

This work proposes a new pruning method called DUST that goes beyond these limitations. To do so, we recast the pruning task as a constrained optimisation problem. Within this framework, we explicitly derive the dual formulation for multivariate data drawn from the exponential family. We analyse its properties and, in particular, establish a strong duality result that underpins the effectiveness of our pruning rule. The proposed DUST rule consists of evaluating a dual function and comparing it to a threshold – yielding a test as simple as that used in PELT. For one-parameter cost functions (corresponding essentially to univariate data), we further derive a simple inequality-based rule by maximising the dual function explicitly. The PELT criterion corresponds to evaluating the dual at zero and is provably upper-bounded by the DUST value. Normalising the dual into a decision function simplifies its interpretation and opens the door to more sophisticated DUST rules applicable to multivariate time series. A comprehensive simulation study confirms the efficiency and versatility of DUST, in particular its time robustness to all change-point regimes and to misspecification. We apply DUST on mouse monitoring data with a force platform to quantify muscle fatigue.

This paper has the following structure. In Section 2 we present the functional problem for pruning in the context of multiple change-point detection. In Section 3 we describe DUST in the one-constraint case and illustrate DUST on a first simple example. Astonishingly, the one-parametric case leads to a simple closed formula for the positivity test associated with the dual. The general presentation of the duality method is given in Section 4. The simulation study of Section 5 explores the efficiency of DUST. We eventually describe an example on a real-world data set in Section 6. If not in the main body of the article, proofs are given in the appendix.

## 2 Change-point problem and its functional description

### 2.1 Model and optimisation problem

We consider a time series consisting of  $n$  data points, denoted by  $y_1, y_2, \dots, y_n$ , where each point is drawn independently from a distribution belonging to the exponential family. A segment, denoted by  $y_{ab}$ , refers to a contiguous sub-sequence of the data,  $y_{a+1}, \dots, y_b$ , where  $0 \leq a < b \leq n$ . In its canonical form with a minimal representation,

the likelihood of a segment  $y_{ab}$  can be expressed as:

$$\begin{aligned} f(y_{ab}; \boldsymbol{\theta}) &= \prod_{i=a+1}^b \left[ h(y_i) \exp \left( \sum_{j=1}^d \theta_j \cdot T_j(y_i) - A(\boldsymbol{\theta}) \right) \right], \\ &= \prod_{i=a+1}^b \left[ h(y_i) \right] \exp \left( \boldsymbol{\theta} \cdot \sum_{i=a+1}^b \mathbf{T}(y_i) - (b-a)A(\boldsymbol{\theta}) \right), \end{aligned}$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)^T$  is the natural parameter, belonging to a convex domain  $\Theta \subset \mathbb{R}^d$ . The function  $A$  is the log-partition function, which is strictly convex due to the minimal representation. The functions  $T_1, \dots, T_d$  are the sufficient statistics, aggregated in the vector  $\mathbf{T} = (T_1, \dots, T_d)^T$ , and  $h$  is the base measure (a normalising term). The dot sign denotes the scalar product.

To define a parametric cost function, we take the negative log-likelihood of a segment and omit the data-dependent term  $\prod h(y_i)$ , which is constant with respect to the segmentation structure (as seen in Equation (4)). For a segment  $y_{ab}$ , we define the sufficient statistic sum  $\mathbf{S}_{ab} = \sum_{i=a+1}^b \mathbf{T}(y_i) \in \mathbb{R}^d$ , and get the cost function:

$$c(y_{ab}; \boldsymbol{\theta}) = (b-a)A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{S}_{ab}. \quad (1)$$

An important property of this cost function is its additivity over disjoint segments. Specifically, for any  $a < b < t$ , we have:

$$c(y_{at}; \boldsymbol{\theta}) - c(y_{bt}; \boldsymbol{\theta}) = c(y_{ab}; \boldsymbol{\theta}). \quad (2)$$

We present a few well-known examples of exponential family models and their corresponding segment cost functions.

**Example 1.** Let  $S_{ab} = \sum_{i=a+1}^b y_i$  denote the sum of univariate data points over the segment  $y_{ab}$ . With a Poisson model, its one-parametric cost function is given by  $c(y_{ab}; \theta) = (b-a)\exp(\theta) - S_{ab}\theta$ ; for exponential model,  $c(y_{ab}; \theta) = (b-a)(-\log(-\theta)) - S_{ab}\theta$ ; for a binomial distribution,  $c(y_{ab}; \theta) = (b-a)\log(1 + \exp(\theta)) - S_{ab}\theta$ ; while for a Gaussian distribution with unitary variance,  $c(y_{ab}; \theta) = (b-a)\frac{\theta^2}{2} - S_{ab}\theta$ .

**Example 2.** For the Gaussian distribution with unknown mean and variance, the canonical form leads to a bi-parametric cost function:

$$c(y_{ab}; [\theta_1; \theta_2]) = (b-a) \left( -\frac{\theta_1^2}{4\theta_2} + \frac{1}{2} \log \left( -\frac{1}{2\theta_2} \right) \right) - \theta_1 \left( \sum_{i=a+1}^b y_i \right) - \theta_2 \left( \sum_{i=a+1}^b y_i^2 \right), \quad (3)$$

which is quadratic in parameter  $\theta_1$  only and  $(\theta_1, \theta_2) \in \mathbb{R} \times \mathbb{R}^-$ .

In this work, we address the problem of multiple change-point detection via penalised maximum likelihood, where a penalty proportional to the number of segments is introduced to control model complexity. Specifically, each segment is assigned a positive penalty value  $\beta$ , referred to as the unitary penalty. A larger  $\beta$  encourages sparser segmentations, i.e., fewer change points. The optimal penalised cost over the entire time series is defined as:

$$Q_n = \min_{\tau \in \mathcal{T}} \left( \sum_{k=0}^K \left[ \min_{\boldsymbol{\theta} \in \Theta} c(y_{\tau_k \tau_{k+1}}; \boldsymbol{\theta}) + \beta \right] \right), \quad (4)$$

where  $\tau = (\tau_0 = 0, \tau_1, \dots, \tau_K, \tau_{K+1} = n)$  is a change-point vector, and  $\mathcal{T}$  denotes the set of all admissible segmentations  $\mathcal{T} = \{\tau \in \mathbb{N}^{K+2}, K \in \mathbb{N}, 0 = \tau_0 < \tau_1 < \dots < \tau_K < \tau_{K+1} = n\}$ . Here,  $K$  – the number of change points – is not fixed in advance but is instead inferred in the optimisation, depending on the choice of penalty value  $\beta$ . The quantity  $Q_n$  is often called the global cost.

This framework, based on the exponential family with independent segments and a linear penalty, is widely used in the literature (e.g., [16, 20]). Non-linear penalty terms have also been investigated in various settings (see [43, 9, 39]). Although our current work focuses on a specific case, we believe that the proposed pruning strategy, which leverages dual functions, is highly versatile and has the potential to be extended to a broader class of models. These include, for example, online change-point detection [27], models with dependencies across segments [12, 31, 33], and non-linear penalisation schemes [29]. We leave the development of these extensions for future work.

## 2.2 Functional cost for pruning

The exact solution to the optimisation problem (4) can be obtained using the optimal partitioning algorithm [16]. In our case of the exponential family, the term  $\min_{\boldsymbol{\theta} \in \Theta} c(y_{st}; \boldsymbol{\theta})$  can be computed in constant time. As a result, the global cost  $Q_t$  can be evaluated in quadratic time using the following recursion over the last segment position using the initial value  $Q_0 = 0$ :

$$Q_t = \min_{0 \leq s < t} \left\{ Q_s + \min_{\boldsymbol{\theta} \in \Theta} c(y_{st}; \boldsymbol{\theta}) + \beta \right\}. \quad (5)$$

This dynamic programming approach has been significantly accelerated through pruning strategies, most notably PELT [20] and, more recently, FPOP [23]. Pruning aims to reduce the number of candidate indices  $s$  considered in the minimisation of (5), by applying conditions that ensure no optimal solution is discarded – thus preserving exactness. In the FPOP framework, the optimisation over the natural parameter  $\boldsymbol{\theta}$  in recursion (5) is postponed. That is, we search for the best last change-point location in truncated data  $y_{0:t}$  for  $t$  from 1 to  $n$  for each value of the parameter  $\boldsymbol{\theta}$ :

$$Q_t(\boldsymbol{\theta}) = \min_{0 \leq s < t} (Q_s + c(y_{st}; \boldsymbol{\theta}) + \beta). \quad (6)$$

Here, the global cost  $Q_t$  has been transformed into a functional cost by treating the parameter  $\boldsymbol{\theta}$  as a free variable. We recover its value by  $Q_t = \min_{\boldsymbol{\theta} \in \Theta} Q_t(\boldsymbol{\theta})$ . Defining inner cost function  $q_t^s$  as:

$$q_t^s(\boldsymbol{\theta}) = Q_s + c(y_{st}; \boldsymbol{\theta}) + \beta, \quad (7)$$

we can write  $Q_t(\boldsymbol{\theta}) = \min_{s \in \mathcal{T}_t} \{q_t^s(\boldsymbol{\theta})\}$  with  $\mathcal{T}_t \subset \{0, \dots, t-1\}$  being the set of non-pruned indices. An effective pruning rule is both fast to test and capable of eliminating a substantial number of indices from  $\{0, \dots, t-1\}$ . For example, the simple pruning criterion used in PELT yields a set  $\mathcal{T}_t$  of bounded size, even for large time series, under the assumption that the number of change points grows proportionally with the data length (see Theorem 3.2 in [20]). A more efficient pruning strategy, FPOP, relies on the following principle.

**Definition 1.** (*Functional pruning principle*) Considering an index  $s$  in  $\{0, \dots, t-1\}$ , if for all  $\boldsymbol{\theta}$  in  $\Theta$  there exists  $r$  (depending on  $\boldsymbol{\theta}$ ) such that  $q_t^r(\boldsymbol{\theta}) < q_t^s(\boldsymbol{\theta})$  then  $s$  can be removed from  $\mathcal{T}_{t'}$  for all  $t' > t$ .

In other words, it means that inner functions which are unseen in the minimisation at some time step  $t$  for all values in the parametric space  $\Theta$  will never be seen again later and can therefore be safely discarded. This insight follows directly from the additive property (2), which implies the equivalence  $q_t^r(\boldsymbol{\theta}) < q_t^s(\boldsymbol{\theta}) \iff q_{t'}^r(\boldsymbol{\theta}) < q_{t'}^s(\boldsymbol{\theta})$  for  $t' > t$ . Functional pruning is the most efficient pruning strategy.

**Proposition 1.** (*FPOP maximal pruning*) *The FPOP-like pruning presented in Definition 1 is the maximal possible pruning for the optimisation problem of type (4) with additive cost property (2). Say differently, at each time step  $t$ , the index set  $\mathcal{T}_t$  obtained by FPOP is minimal: a smaller  $\mathcal{T}_t$  would potentially lead to an under-optimal solution.*

The proof relies on the following argument. If we remove an index  $s_0$  in  $\mathcal{T}_t$  such that  $Q_t(\boldsymbol{\theta}) = q_t^{s_0}(\boldsymbol{\theta})$  in a neighbourhood of  $\boldsymbol{\theta}_0$ , we can prove that, by adding data points "centred on  $\boldsymbol{\theta}_0$ " at further iterations (with unitary cost  $A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \nabla A(\boldsymbol{\theta}_0)$ ), the last segment in (6) would start at index  $s_0 + 1$  at some later time  $t_0$ . However, since  $s_0$  was pruned, the solution of the DP algorithm is no longer exact.

It is important to note that, although pruning rules can drastically reduce execution time in practice, there exist examples of data for which no pruning rule could discard any index, leading back to the worst-case quadratic complexity of the optimal partitioning algorithm. We construct such examples using increasing time series where all inner cost functions attain the same minimal value at time  $n$ . We first illustrate this phenomenon in the simple univariate Gaussian case (Proof in Appendix C.1), and then extend the result to a broader class of continuous distributions within the natural exponential family in Appendix C.2.

**Proposition 2.** *We consider the Gaussian univariate model with fixed variance. Its inner cost functions are given by relations  $q_t^s(\theta) = Q_s + \frac{t-s}{2}\theta^2 - S_{st}\theta + \beta$  with  $S_{st} = \sum_{i=s+1}^t y_i$  and  $y_i \in \mathbb{R}$ . If we observe the following data points:*

$$y_t = \sqrt{\frac{\beta}{n}} \left( \sqrt{n-1} - \sqrt{t(n-t)} + \sqrt{(t-1)(n-t+1)} \right), \quad t = 1, \dots, n \quad (8)$$

*then no pruning happens, that is  $\mathcal{T}_n = \{0, \dots, n-1\}$ .*

### 2.3 Functional pruning rule as an optimisation problem

In the evolution of the dynamic programming algorithm, if an inner function  $q_t^s$  is still accessible, it can be discarded only by the last introduced function at time  $t+1$ , that is  $Q_t + c(y_{t(t+1)}; \boldsymbol{\theta}) + \beta$  (due to the additive property (2)). This naturally leads to the following pruning strategy: compare each inner function at time  $t$  with the constant threshold  $Q_t + \beta$ . We denote by  $\Theta_t^s \subset \Theta$  the set of parameter values for which the function  $q_t^s$  attains the minimum among all candidate functions at time  $t$ , defined explicitly as:

$$\Theta_t^s := \{\boldsymbol{\theta} \in \Theta \mid \forall r \in \mathcal{T}_t \setminus \{s\}, q_t^s(\boldsymbol{\theta}) \leq q_t^r(\boldsymbol{\theta})\}. \quad (9)$$

The criterion " $\Theta_t^s = \emptyset$ ?" forms the basis of the functional pruning rule introduced in FPOP [23] for univariate data, and has since been adopted in several one-dimensional extensions [18, 31]. For multivariate problems, a recent method proposes to approximate the set  $\Theta_t^s$  using simple geometric shapes [26]. However, this approach becomes

challenging to implement efficiently for non-Gaussian cost functions, resulting in relatively slow updates (see also [32]). Instead, with DUST, we propose to evaluate the minimal value of  $q_t^s$  over its visibility region  $\Theta_t^s$ , and compare it directly to the threshold  $Q_t + \beta$ . While this may appear to introduce additional computation, it avoids the need to characterise the potentially complex shape of  $\Theta_t^s$  explicitly. Thus, the task reduces to finding a single value – or a lower bound on it – which is often more tractable in practice.

**Proposition 3** (Functional pruning with minimum value rule). *We define  $R_t^s := \min_{\boldsymbol{\theta} \in \Theta_t^s} q_t^s(\boldsymbol{\theta})$  with  $\Theta_t^s$  given in (9). By convention, if  $\Theta_t^s = \emptyset$ , then  $R_t^s = +\infty$ . The functional pruning condition is as follows. If there exists  $s$  in  $\mathcal{T}_t$  satisfying:*

$$R_t^s > Q_t + \beta, \quad (10)$$

*then  $s$  can never be the last change point of a segmentation of  $y_{0t'}$  for all  $t' > t$ : that is  $s \notin \mathcal{T}_{t'}$ .*

*Proof.* Let  $\boldsymbol{\theta} \notin \Theta_t^s$ . By definition of  $\Theta_t^s$  there exists  $r$  such that  $q_t^s(\boldsymbol{\theta}) > q_t^r(\boldsymbol{\theta})$  and therefore  $q_{t'}^s(\boldsymbol{\theta}) > q_{t'}^r(\boldsymbol{\theta})$ . Let  $\boldsymbol{\theta} \in \Theta_t^s$ . If the functional pruning condition (10) holds, we have  $q_{t'}^s(\boldsymbol{\theta}) = q_t^s(\boldsymbol{\theta}) + c(y_{tt'}, \boldsymbol{\theta}) > Q_t + \beta + c(y_{tt'}, \boldsymbol{\theta}) = q_{t'}^t(\boldsymbol{\theta})$  which means that index  $s$  is never seen in the minimisation problem (it is hidden by index  $t$ ) and can be discarded.  $\square$

The central challenge in pruning for multiple change-point detection is a problem of constrained optimisation. Its dual formulation forms the core of the DUST pruning method.

**Definition 2** (Pruning problem). *The pruning problem for multiple change-point detection at time  $t$  for testing change point  $s$  is the following problem of optimisation under constraints:*

$$\begin{cases} \min_{\boldsymbol{\theta} \in \Theta} q_t^s(\boldsymbol{\theta}), \\ \text{s.t. } q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta}) \leq 0 \quad \text{for all } r \in \mathcal{T}_t \setminus \{s\} \subset \{0, \dots, t-1\}. \end{cases} \quad (11)$$

*A result greater than threshold  $Q_t + \beta$  implies that the function  $q_t^s(\boldsymbol{\theta})$  can be removed and the index  $s$  pruned.*

We aim to compute either the exact solution or a lower bound, but with a very efficient algorithm. We begin by considering a simplified version of the pruning problem with a single constraint to derive a lower bound. This setting offers a clear framework to introduce and illustrate our dual approach. Furthermore, a very efficient pruning test can be designed for one-parametric cost functions.

## 2.4 Notations for dual and decision functions

For clarity, we compile below all the notations related to dual and decision functions used in the sequel. The function  $\mathcal{D}^*$  is given by  $\mathcal{D}^*(x) = x \cdot (\nabla A)^{-1}(x) - A((\nabla A)^{-1}(x))$  and is strictly convex. We introduce the following mean values

$$\bar{\mathbf{S}}_{rs} = \begin{cases} \frac{\mathbf{S}_{rs}}{s-r} & \text{if } r < s, \\ \frac{\mathbf{S}_{sr}}{r-s} & \text{if } s < r, \end{cases} \quad \text{and index indicator } \psi_{rs} = \begin{cases} 1 & \text{if } r < s, \\ -1 & \text{if } s < r. \end{cases}$$

We also use a mean value between global costs,  $\bar{Q}_{rs} = \frac{Q_s - Q_r}{s - r}$ , as well as a difference operator between mean values, which incorporates the order between the first two indices:

$$\Delta \bar{\mathbf{S}}_{rst} = \psi_{rs}(\bar{\mathbf{S}}_{st} - \bar{\mathbf{S}}_{rs}) \quad , \quad \Delta \bar{Q}_{rst} = \psi_{rs}(\bar{Q}_{st} - \bar{Q}_{rs}) .$$

The absence of bold notation (e.g.,  $S_{rs}, \bar{S}_{rs}, \Delta \bar{S}_{rst}, \dots$ ) indicates that the data is univariate. In this case, some standard notation will also be used for means and variance when  $r < s$ :

$$\bar{S}_{rs} = \bar{y}_{rs} = \frac{1}{s - r} \sum_{i=r+1}^s y_i, \quad \bar{S}_{rs}^2 = \bar{y}_{rs}^2 = \frac{1}{s - r} \sum_{i=r+1}^s y_i^2,$$

and  $V(y_{rs}) = \bar{y}_{rs}^2 - (\bar{y}_{rs})^2$ . Eventually, we define the following linear function:

$$\sigma_{\mathcal{R}}(\mathbf{x}) = \bar{\mathbf{S}}_{st} + \sum_{r \in \mathcal{R}} x_r \Delta \bar{\mathbf{S}}_{rst} \quad , \quad \phi_{\mathcal{R}}(\mathbf{x}) = \bar{Q}_{st} + \sum_{r \in \mathcal{R}} x_r \Delta \bar{Q}_{rst} .$$

The vector  $\mathbf{x} = (x_r)_{r \in \mathcal{R}}$  is identified with  $(x_r)_{r=1, \dots, |\mathcal{R}|}$  depending on the context (same with vector  $\mu$ ). For a vector  $\mathbf{v} \in \mathbb{R}^d$ , its Euclidean norm is denoted by  $\|\mathbf{v}\|$ .

### 3 DUST method with one constraint

We present the dual formulation for the single-constraint case without delving into the theoretical analysis of the dual, which is deferred to Section 4 in a generic framework. We provide a detailed description of the DUST change-point algorithm. As a concrete illustration, we explore the change-in-mean-and-variance problem in depth. This simplified presentation is intended to make our methodology more accessible and to promote its application in other settings.

#### 3.1 A one-dimensional dual function

Optimisation problem (11) with one constraint is as follows:

$$\begin{cases} \min_{\boldsymbol{\theta} \in \Theta} q_t^s(\boldsymbol{\theta}), \\ \text{s.t. } q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta}) \leq 0 . \end{cases} \quad (12)$$

where we consider that  $r < s$ . With  $r > s$ , the obtained (convex) problem (12) leads to a very inefficient pruning. The optimal value of (12) is defined as  $R_t^{rs}$ . We can easily write down the Lagrangian function as:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}, \mu) &= q_t^s(\boldsymbol{\theta}) + \mu(q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta})), \\ &= \left( (t - s) - \mu(s - r) \right) \left( A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \frac{\mathbf{S}_{st} - \mu \mathbf{S}_{rs}}{(t - s) - \mu(s - r)} \right) + \mu(Q_s - Q_r) + Q_s + \beta, \end{aligned}$$

or with a change of variable  $\mu \rightarrow \mu \frac{t-s}{s-r}$  and the parametric mean  $\mathbf{m}(\mu) = \frac{\bar{\mathbf{S}}_{st} - \mu \bar{\mathbf{S}}_{rs}}{1-\mu}$ :

$$\mathcal{L}(\boldsymbol{\theta}, \mu) = (t - s) \left( (1 - \mu)(A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{m}(\mu)) + \mu \frac{Q_s - Q_r}{s - r} \right) + Q_s + \beta .$$

The dual function presented in the following proposition serves as a lower bound for the target quantity  $R_t^s$  (see Proposition 3).

**Proposition 4.** *The dual function  $\mathcal{D} : [0, \mu_{\max}] \rightarrow \mathbb{R}$  to Problem (12) with inner cost functions (7) and cost (1) for pruning index  $s$  using index  $r$  with  $r < s$  is given by:*

$$\mathcal{D}(\mu) = (t - s) \left( - (1 - \mu) \mathcal{D}^*(\mathbf{m}(\mu)) + \mu \bar{Q}_{rs} \right) + Q_s + \beta, \quad (13)$$

where  $\mathbf{m}(\mu) = \frac{\bar{\mathbf{S}}_{st} - \mu \bar{\mathbf{S}}_{rs}}{1 - \mu}$ . The domain of the dual is a segment bounded by a value  $\mu_{\max}$  which is the largest value for which  $\mathcal{D}^*(\mathbf{m}(\mu))$  is finite or can be computed (e.g. when we have log terms). The value  $\mu_{\max}$  is always smaller than or equal to 1.

At time  $t$ , we prune index  $s$  using a candidate index  $r < s$  if  $\mathcal{D}(\mu_0) > Q_t + \beta$ , for some  $\mu_0 \in [0, \mu_{\max}]$ . The solution  $R_t^{rs}$  to Problem (12) is smaller than the solution  $R_t^s$  to Problem (11), and is lower bounded by the dual function: for all  $\mu \in [0, \mu_{\max}]$ , we have:  $\mathcal{D}(\mu) \leq R_t^{rs} \leq R_t^s$ . Therefore, using any value  $\mu_0 \in [0, \mu_{\max}]$ , pruning based on the condition  $\mathcal{D}(\mu_0) > Q_t + \beta$  is safe: it guarantees that  $R_t^s > Q_t + \beta$ , so no index is removed by mistake. Ideally, we would evaluate the dual function at its maximum or some value close to it.

### 3.2 DUST decision rule

Pruning relies on the existence of a value  $\mu_0$  such that  $\mathcal{D}(\mu_0) - (Q_t + \beta) > 0$ . Since dividing the left-hand side by any positive function does not change the decision, we normalise the dual function  $\mathcal{D}$  from Proposition 4 to allow its direct maximisation (with a closed formula) in some particular cases. The new function  $\mathbb{D}$ , called the decision function, has this simple property: pruning of index  $s$  occurs as soon as  $\mathbb{D}$  returns a positive value for some point in its domain.

**Proposition 5.** *The decision function  $\mathbb{D} : [0, x_{\max}] \rightarrow \mathbb{R}$  to Problem (12) with inner cost functions (7) and cost (1) for pruning index  $s$  using index  $r$  with  $r < s$  is given by:*

$$\mathbb{D}(x) = -\mathcal{D}^*(\bar{\mathbf{S}}_{st} + x \Delta \bar{\mathbf{S}}_{rst}) - (\bar{Q}_{st} + x \Delta \bar{Q}_{rst}). \quad (14)$$

The value  $x_{\max}$  is derived from  $\mu_{\max}$ .

*Proof.* We shift the dual  $\mathcal{D}$  of Equation (13) by the quantity  $Q_t + \beta$  to transform the dual pruning threshold test into a sign test. We then divide the obtained function by  $(t - s)(1 - \mu)$ , which is always a positive value. We eventually introduce the variable  $x = \frac{\mu}{1 - \mu}$ . This is summarised by the transformation:

$$\mathbb{D}(x) = \left( (t - s)(1 - \frac{x}{1 + x}) \right)^{-1} \left( \mathcal{D}\left(\frac{x}{1 + x}\right) - (Q_t + \beta) \right).$$

□

A typical setting in change-point detection is that of univariate data drawn from the exponential family [8, 33, 30]. For this setting, we can explicitly compute the maximum of the decision function.

**Theorem 6.** *The decision function  $\mathbb{D} : [0, x_{\max}] \rightarrow \mathbb{R}$  to Problem (12) given in Proposition 5 admits a closed-form maximum for the single-constraint case and one-parametric*

cost functions. If the argument of the maximum  $x^*$  of this function is not located on the frontier of the domain, then we prune index  $s$  using index  $r$  with  $r < s$  when:

$$q_t^s \left( -\frac{\Delta \bar{Q}_{rst}}{\Delta \bar{S}_{rst}} \right) > Q_t + \beta. \quad (15)$$

Note that certain special cases must be handled before testing inequality (15), such as  $x_{\max} = 0$ ,  $\bar{S}_{st} = \bar{S}_{rs}$  (i.e., the decision function is linear) and  $x^* \notin (0, x_{\max})$ . This theorem replaces the PELT pruning rule, which was given by  $q_t^s((\nabla A)^{-1}(\bar{S}_{st})) > Q_t + \beta$ . We will highlight the effectiveness of the DUST rule in simulation Section 5.

### 3.3 DUST algorithm

We have designed our new pruning method as straightforward and practical as the PELT rule. The DUST multiple change-point detection algorithm is described in Algorithm 1. Here, the value  $c(y_{st})$  is defined as the minimum of the segment cost function  $\boldsymbol{\theta} \mapsto c(y_{st}; \boldsymbol{\theta})$ . For the smallest index in  $\mathcal{T}_t$ , we use the standard PELT pruning rule by default. The index  $r$  is selected from a given probability distribution  $p_r$  over the indices in  $\mathcal{T}_t$  that are smaller than  $s$ , while the value  $\mu_0$  is sampled from a distribution  $p_\mu$  supported on the segment  $[0, \mu_{\max} \leq 1]$ . For one-parametric cost functions, we can consider a Dirac data-dependent  $p_\mu$  which matches relation (15) in addition to its limit cases.

Simulations (Section 5) explore several choices for the distributions  $p_r$  and  $p_\mu$ . In the single-constraint setting, an effective strategy is to select  $r$  as the closest index below  $s$ , and to choose  $\mu_0$  as the argument that maximises the dual function, using a simple iterative maximisation procedure or relation (15) for one-parametric cost functions. For simplicity, the algorithm is presented with the dual  $\mathcal{D}$  with bounded domain, but a similar algorithm can be written with decision function  $\mathbb{D}$ . From the output of Algorithm 1, we can easily recover the optimal segmentation for the initial problem (4) via a standard backtracking step described in Algorithm 2.

**Remark 1.** *DUST extends PELT by evaluating the dual function beyond zero. While PELT applies the test  $\mathcal{D}(0) > Q_t + \beta$ , derived from the unconstrained form of Problem (12), DUST samples  $\mu_0$  in  $[0, \mu_{\max}]$ , potentially yielding stronger pruning. This makes DUST especially effective in time series with few change points, where PELT tends to prune poorly. Despite its improved pruning, DUST maintains a complexity close to PELT, with minimal overhead from computing  $\mu_{\max}$  and selecting  $r$  and  $\mu_0$ .*

### 3.4 Example of the change-in-mean-and-variance problem

Detecting changes in both the mean and variance of a large Gaussian time series is challenging. To the best of our knowledge, no efficient implementation currently exists in this setting due to two main difficulties. First, the problem involves two parameters, which limit the applicability of methods like FPOP [23]. Second, a logarithmic term in the cost function complicates root-finding, making methods such as geomFPOP [26] difficult to apply. We address this challenge using the DUST method. We write down the likelihood of segment  $y_{st}$  with mean  $m$  and variance  $\sigma^2$ ,

$$\mathcal{L}(y_{st}; [m; \sigma^2]) = \prod_{i=s+1}^t \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{(m - y_i)^2}{2\sigma^2} \right) \right],$$

---

**Algorithm 1** DUST algorithm

---

**Input:** Time series  $y_{0n}$ , penalty value  $\beta > 0$ , rules  $p_r, p_\mu$

**Output:** Sequences of global costs  $Q_{0n}$  and last changes  $\hat{s}_{0n}$

```

1:  $Q_0 \leftarrow 0, \quad \mathcal{T}_1 \leftarrow \{0\}$                                 ▷ OP step
2: for  $t = 1, \dots, n$  do
3:    $Q_t \leftarrow \min_{s \in \mathcal{T}_t} \{Q_s + c(y_{st}) + \beta\}$ 
4:    $\hat{s}_t \leftarrow \arg \min_{s \in \mathcal{T}_t} \{Q_s + c(y_{st}) + \beta\}$ 
5:   for  $s \in \mathcal{T}_t$  do                                              ▷ DUST pruning step
6:     Draw  $r$  in  $\mathcal{T}_t$  such that  $r < s$  from distribution  $p_r$ 
7:     Compute  $\mu_{max} = \mu_{max}^{r,s,t}$  (data dependent)
8:     Draw  $\mu_0$  in  $[0, \mu_{max}]$  from distribution  $p_\mu$ 
9:     if  $D(\mu_0) > Q_t + \beta$  then                                         ▷ DUST test
10:     $\mathcal{T}_t \leftarrow \mathcal{T}_t \setminus \{s\}$ 
11:   end if
12:   end for
13:    $\mathcal{T}_{t+1} \leftarrow \mathcal{T}_t \cup \{t\}$ 
14: end for

```

---

**Algorithm 2** Backtracking the change-point locations

---

**Input:**  $\hat{s}_{0n}$

**Output:** Set of optimal change point indices  $\widehat{\mathcal{T}} = \{\tau_1, \tau_2, \dots\}$

```

1:  $\tau \leftarrow n, \widehat{\mathcal{T}} \leftarrow \emptyset$ 
2: while  $\tau > 0$  do
3:    $\widehat{\mathcal{T}} \leftarrow (\tau, \widehat{\mathcal{T}}), t \leftarrow \hat{s}_\tau$ 
4: end while

```

---

and its associated cost:

$$c(y_{st}; [m; \sigma^2]) = (t - s) \left( \frac{(m - \bar{y}_{st})^2 + V(y_{st})}{2\sigma^2} + \frac{1}{2} \log(\sigma^2) \right).$$

With the change of variable  $(\theta_1, \theta_2) = (\frac{m}{\sigma^2}; -\frac{1}{2\sigma^2}) \in \mathbb{R} \times \mathbb{R}^-$ , we obtain the canonical form of Example 2 with  $A(\theta_1, \theta_2) = -\frac{\theta_1^2}{4\theta_2} + \frac{1}{2} \log(-\frac{1}{2\theta_2})$ . We can now derive the exact maximum value of the associated decision function.

**Lemma 1.** *The one-dimensional decision function for the change-in-mean-and-variance problem with a single constraint from the cost function (3) is given by the function:*

$$\mathbb{D}(x) = \frac{1}{2} \left[ 1 + \log \left( \overline{S}_{st}^2 + x\Delta\overline{S}_{rst}^2 - (\overline{S}_{st} + x\Delta\overline{S}_{rst})^2 \right) \right] - \left( \overline{Q}_{st} + x\Delta\overline{Q}_{rst} \right),$$

with  $x \in (x_0 - \sqrt{x_1}, x_0 + \sqrt{x_1})$  and

$$x_0 = \frac{1}{2} \left( \frac{V(y_{st}) - V(y_{rs})}{(\Delta\overline{S}_{rst})^2} - 1 \right) \quad x_1 = x_0^2 + \frac{V(y_{st})}{(\Delta\overline{S}_{rst})^2}$$

and its maximum, with notation  $x_2 = \Delta\overline{Q}_{rst}$ , is evaluated in:

$$x^* = \max \left\{ 0, x_0 + (2x_2)^{-1} - \text{sign}(x_2)\sqrt{x_1 + (2x_2)^{-2}} \right\}$$

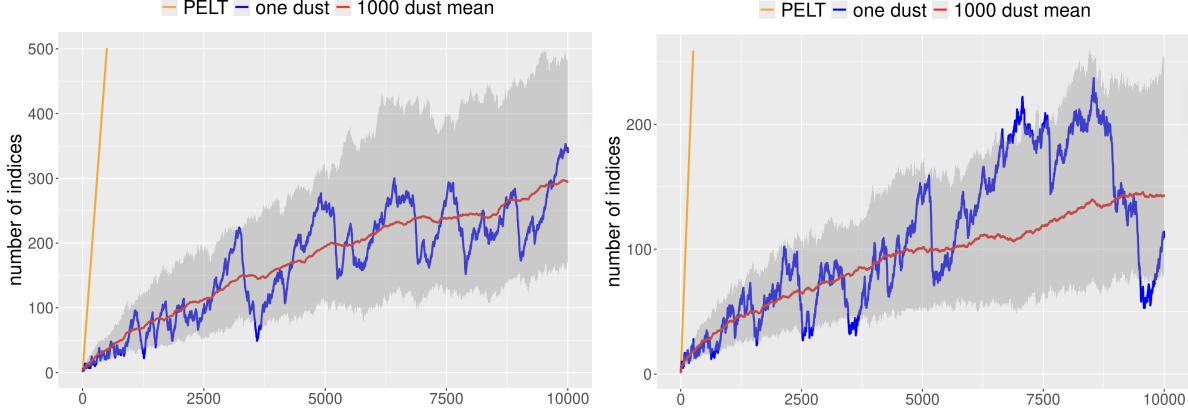


Figure 1: For data with no change, the number of indices saved by DUST over time is consistently much less than PELT (no pruning), for the one-constraint DUST (left) and the two-constraint case (right)

In Figure 1, we evaluate the efficiency of DUST for this example. For each index  $s$  considered for pruning, we choose the biggest non-pruned index smaller than  $s$  and evaluate the dual at its maximum position  $x^*$ . With  $n = 10^4$ ,  $\beta = 4 \log(n)$  and data with no change ( $\mathcal{N}(0, 1)$ ), only 2.95% of the indices remain non-pruned and the overall number of indices to be considered is reduced by a factor 28 compared to PELT (left panel).

A higher lower bound for  $R_t^s$  is obtained when we consider up to two well-chosen constraints:  $q_t^s - q_t^{r_1} \leq 0$  and  $q_t^s - q_t^{r_2} \leq 0$ . This leads to a bi-parametric decision function:

$$\begin{aligned} \mathbb{D}(x_1, x_2) = & \frac{1}{2} \left[ 1 + \log \left( \overline{S_{st}^2} + x_1 \Delta \overline{S_{r_1 st}^2} + x_2 \Delta \overline{S_{r_2 st}^2} \right) - (\overline{S}_{st} + x_1 \overline{S}_{r_1 st} + x_2 \overline{S}_{r_2 st})^2 \right] \\ & - \left( \overline{Q}_{st} + x_1 \Delta \overline{Q}_{r_1 st} + x_2 \Delta \overline{Q}_{r_2 st} \right). \end{aligned}$$

Such a general dual and its properties are studied in the next Section. In this mean-and-variance problem, explicit maximal values are available for the three values to be tested:  $\max_{x_1} \mathbb{D}(x_1, 0)$ ,  $\max_{x_2} \mathbb{D}(0, x_2)$  and  $\max_{x_1, x_2} \mathbb{D}(x_1, x_2)$  with positive  $x_1, x_2$ . In Figure 1 (right panel) we run the same simulation as in the left panel but with a bi-parametric dual, built from the two highest indices smaller than  $s$  ( $r_1 < r_2 < s$ ). Only 1.42% of the indices remain non-pruned and the overall number of indices to be considered is reduced by a factor 54 compared to PELT. In this example, 0 indice is pruned by PELT, 36% of indices are pruned by the  $r_2$ -dual, 23% of indices are pruned by the  $r_1$ -dual and 40% by the maximum of the  $(r_1, r_2)$ -dual (function  $\mathbb{D}(x_1, x_2)$  with  $x_1 > 0$  and  $x_2 > 0$ ). We recall that this DUST pruning procedure does not involve any iterative routine, but relies solely on three inequality tests (as PELT) evaluated at three points of the decision function, each computed from a closed-form expression. With  $10^6$  data point, the percent of indices drop to 0.5%.

## 4 Duality function in change-point problems

### 4.1 The shape of the dual

We consider the pruning problem (11):  $\min_{\boldsymbol{\theta} \in \Theta} q_t^s(\boldsymbol{\theta})$  under constraints  $q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta}) \leq 0$  for all indices  $r \neq s$  in  $\{0, \dots, t-1\}$ . As for the single-constraint case, the dual function can be explicitly written in closed form.

**Proposition 7.** *The dual function,  $\mathcal{D}_{st} : (\mathbb{R}^+)^{t-1} \rightarrow \mathbb{R}$  for analysing change index  $s$  at time  $t > s$  is given by expression:*

$$\mathcal{D}_{st}(\mu) = -(t-s) \left[ l(\mu) \mathcal{D}^*(\mathbf{m}(\mu)) + \sum_{r \neq s} (-1)^{[r < s]} \mu_r \bar{Q}_{rs} \right] + Q_s + \beta,$$

with  $\mu = (\mu_r)_{0 \leq r < t, r \neq s}$ ,  $[r < s] = 1$  if  $r < s$  and 0 otherwise. We also have:

$$\mathbf{m}(\mu) = \frac{\bar{\mathbf{S}}_{st} + \sum_r \mu_{r \neq s} (-1)^{[r < s]} \bar{\mathbf{S}}_{rs}}{1 + \sum_r \mu_{r \neq s} (-1)^{[r < s]}} \in \mathbb{R}^d, \quad l(\mu) = 1 + \sum_{r \neq s} \mu_r (-1)^{[r < s]}. \quad (16)$$

Notice that  $\mathbf{m}(\mu)$  is a vector if the statistics  $\mathbf{S}_{st}$  are vector-valued. The case  $\bar{\mathbf{S}}_{st} = \bar{\mathbf{S}}_{rs}$  for all  $r \neq s$  can be left apart. It corresponds to  $\mathbf{m}(\mu) = \bar{\mathbf{S}}_{st}$  and the dual  $\mathcal{D}_{st}$  becomes a linear function.

*Proof.* We differentiate in  $\boldsymbol{\theta}$  the Lagrangian function  $\mathcal{L}(\boldsymbol{\theta}, \bar{\mu}) = q_t^s(\boldsymbol{\theta}) + \sum_{r \neq s} \bar{\mu}_r (q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta}))$  with  $q_t^r(\boldsymbol{\theta}) = Q_r + (t-r)A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{S}_{rt} + \beta$  and inject its solution  $\boldsymbol{\theta}^*(\bar{\mu}) = (\nabla A)^{-1}(\mathbf{m}(\bar{\mu}))$  to get  $\mathcal{D}_{st}(\bar{\mu}) = \mathcal{L}(\boldsymbol{\theta}^*(\bar{\mu}), \bar{\mu})$ . The final result comes after the change of variable  $\bar{\mu}_r = \mu_r \frac{t-s}{|s-r|}$ .  $\square$

**Remark 2.** With constraints of type "r < s" only, the domain  $\Omega_\mu$  of  $\mathcal{D}_{st}$  is bounded and included in the simplex of dimension d.

We also define the decision function  $\mathbb{D}_{st}$  (as done in previous section) with the change of variable  $x_r = \mu_r / (1 + \sum_r \mu_r (-1)^{[r < s]})$  and notation  $\mathbf{x} = (x_r)_{0 \leq r < t, r \neq s}$ . It is given by a simpler relation:

$$\begin{aligned} \mathbb{D}_{st}(\mathbf{x}) &= -\mathcal{D}^* \left( \bar{\mathbf{S}}_{st} + \sum_{r \neq s} x_r \Delta \bar{\mathbf{S}}_{rst} \right) - \left( \bar{Q}_{st} + \sum_{r \neq s} x_r \Delta \bar{Q}_{rst} \right), \\ &= -\mathcal{D}^*(\sigma(\mathbf{x})) - \phi(\mathbf{x}), \end{aligned} \quad (17)$$

leading to the simple pruning rule  $\mathbb{D}_{st}(\mathbf{x}_0) > 0$ . The values of the decision function can be interpreted as evaluations of  $q_t^s$  along a linear path from its minimum, where  $\sigma(\mathbf{x}_0) = \boldsymbol{\theta}_0$  and the minimum of  $q_t^s$  is attained at  $\sigma(0) = \bar{\mathbf{S}}_{st}$ .

Appendix G presents the form of the function  $\mathcal{D}^*$  for many distributions. We give here the explicit form of the decision function for two examples in dimension d with q constraints (in index set  $\mathcal{R} \subset \{0, \dots, t-1\}$ ) with  $|\mathcal{R}| = q \leq d$ .

(i) Multivariate Gaussian:

$$\begin{aligned} \mathbb{D}_{st}^N(\mathbf{x}) &= -\frac{1}{2} \sum_{i=1}^d \left( (\bar{\mathbf{S}}_{st})_i + \sum_{r \in \mathcal{R}} x_r (\Delta \bar{\mathbf{S}}_{rst})_i \right)^2 - \left( \Delta Q_{st} + \sum_{r \in \mathcal{R}} x_r \Delta \bar{Q}_{rst} \right) \\ &= -\frac{1}{2} \sum_{i=1}^d \sigma_i(\mathbf{x})^2 - \phi(\mathbf{x}). \end{aligned}$$

(ii) Multivariate Bernoulli:

$$\mathbb{D}_{st}^B(\mathbf{x}) = - \sum_{i=1}^d \left[ \sigma_i(\mathbf{x}) \log(\sigma_i(\mathbf{x})) - (1 - \sigma_i(\mathbf{x})) \log(1 - \sigma_i(\mathbf{x})) \right] - \phi(\mathbf{x}).$$

In this case, its domain  $\Omega_{\mathbf{x}}$  is at the intersection of  $3d$  half-spaces in dimension  $q$  given by relations  $\sigma_i(\mathbf{x}) \geq 0$ ,  $1 - \sigma_i(\mathbf{x}) \geq 0$  and the positive orthant  $\mathbf{x} > 0$ . The shape of the domain depends on the underlying distribution used.

**Definition 3** (DUST pruning rule with multiple constraints). *We consider the dual function  $\mathcal{D}_{st} : \Omega_{\mu} \rightarrow \mathbb{R}$  and the decision function  $\mathbb{D}_{st} : \Omega_{\mathbf{x}} \rightarrow \mathbb{R}$ . At time  $t$ , we prune index  $s$  using candidate indices from a set  $\mathcal{R} \subset \{0, \dots, t-1\} \setminus \{s\}$  to construct the dual, if there exists  $\mu_0 \in \Omega_{\mu}$  such that  $\mathcal{D}_{st}(\mu_0) > Q_t + \beta$ , or equivalently, if there exists  $\mathbf{x}_0 \in \Omega_{\mathbf{x}}$  such that  $\mathbb{D}_{st}(\mathbf{x}_0) > 0$ .*

This definition leaves implicit three important choices that must be addressed: (i) the selection of candidate indices in  $\mathcal{R}$ , (ii) the method for identifying a suitable evaluation point ( $\mu_0$  or  $\mathbf{x}_0$ ), and (iii) a precise characterisation of the domain of definition ( $\Omega_{\mu}$  or  $\Omega_{\mathbf{x}}$ ). Since these domains always lie within the positive orthant, solving the equation  $\nabla \mathbb{D}_{st}(\mathbf{x}) = 0$  is never a valid option. For now, we address the problem using a quasi-Newton iterative algorithm or random sampling, based on a carefully chosen distribution over  $\Omega_{\mu}$ . This domain is defined through the mean parameter space  $\mathcal{M}$  [41]. We can show that the correct number  $q$  of constraints to consider is upper bounded by the dimension  $d$  of the parameter space  $\Theta \subset \mathbb{R}^d$ .

## 4.2 d-Strong duality

Geometrically, the solution to our optimisation problem (11) lies at the intersection of at most  $d+1$  inner functions in a parameter space of dimension  $d$ . This implies that no more than  $d$  constraints can be active at the solution point. Identifying the correct subset of  $q \leq d$  active constraints is computationally intractable. Nevertheless, focusing on at most  $d$  constraints among the available  $t-1$  is particularly appealing: in this case, the maximum of the dual function coincides with the solution of the original optimisation problem.

**Theorem 8.** *There is no duality gap:*

$$R_t^s := \min_{\boldsymbol{\theta} \in \Theta_t^s(\mathcal{R}) \subset \mathbb{R}^d} \left\{ q_t^s(\boldsymbol{\theta}) \right\} = \max_{\mu \in \Omega_{\mu} \subset \mathbb{R}^q} \left\{ \mathcal{D}_{st}(\mu) \right\},$$

where  $\Theta_t^s(\mathcal{R})$  is the feasible set defined by the  $q$  constraints with  $q \leq d$  such that  $q_t^s(\boldsymbol{\theta}) - q_t^{r_i}(\boldsymbol{\theta}) \leq 0$  with  $r_i \in \mathcal{R}$ ,  $i = 1, \dots, q = |\mathcal{R}|$ , and  $\mathcal{D}_{st}$  is built from these constraints.

We consider  $d$  arbitrary constraints. In case,  $q < d$ , we can still add  $d-q$  constraints of type  $q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta}) \leq 0$  that we know would be unused (no equality at the optimum point). With  $d$  constraints, we consider the geometric object  $\mathcal{O}$  in the Euclidean space  $\mathbb{R} \times \mathbb{R}^d$  whose elements  $(x_0, x_1, \dots, x_d)$  are given by a parametric representation:

$$\mathcal{O} = \left\{ (x_0, \dots, x_d), \exists \boldsymbol{\theta} \in \mathbb{R}^d, (x_0, \dots, x_d) = (q_t^s(\boldsymbol{\theta}), (q_t^s - q_t^{r_1})(\boldsymbol{\theta}), \dots, (q_t^s - q_t^{r_d})(\boldsymbol{\theta})) \right\}.$$

The proof provided in Appendix E.1 is based on the geometric interpretation of duality. Specifically, we show that the epigraph of the objective function  $\mathcal{O}$  is convex. This result relies primarily on the fact that the functions  $q_t^j$  are similar. Notably, the duality gap is zero – that is, strong duality holds – even in the presence of concave constraints (i.e., when  $r < s$ ).

We illustrate the failure of strong duality when  $q > d$  by showing that it does not hold with  $d + 1$  constraints. This is demonstrated using a simple example of a one-dimensional Gaussian cost with three data points,  $(y_1, y_2, y_3) = (2, -1, 0)$ , penalty  $\beta = 2$ , and initial cost  $Q_0 = -\beta$ . These data points yield the following functions:  $q_3^0(\theta) = \frac{3}{2}\theta^2 - \theta$ ,  $q_3^1(\theta) = \theta^2 + \theta$  and  $q_3^2(\theta) = \frac{1}{2}\theta^2 + \frac{3}{2}$ . In this example, the dual function (without normalising the Lagrangian parameters) for testing index 2 with indices 0 and 1 is given by:

$$\mathcal{D}(\mu_1, \mu_2) = \frac{1}{2} \frac{(\mu_1 - \mu_2)^2}{1 - 2\mu_1 - \mu_2} + \frac{3}{2}(1 + \mu_1 + \mu_2),$$

which attains a maximum value of 2.5. The solution to the corresponding optimisation problem is  $\frac{10+\sqrt{7}}{4}$ , as illustrated in Figure 2.

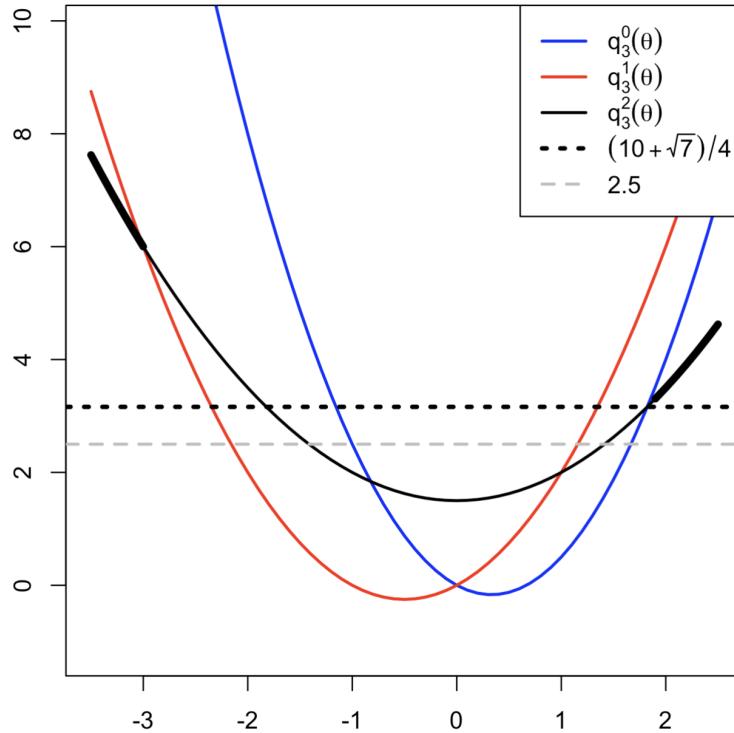


Figure 2: Illustration of the dual maximum value (2.5) and the true threshold ( $\frac{10+\sqrt{7}}{4}$ ), solution of the optimisation problem (11).

### 4.3 Gaussian case

The maximum of the dual can be explicitly expressed in closed form in the single-constraint case. When no changes are present in the time series and the penalty is

well-chosen, we use the relation  $Q_u = -\frac{u}{2} \|\bar{\mathbf{S}}_{0u}\|^2$  to simplify the expression further.

**Proposition 9.** *In d-variate Gaussian model the maximum value of the dual with one constraint ( $r < s < t$ ),  $\max_{\mu \in [0,1]} \left\{ \mathcal{D}^{gauss}(\mu) \right\}$ , is given by expression:*

$$\begin{cases} -\frac{t-s}{2} \|\bar{\mathbf{S}}_{st}\|^2 + Q_s + \beta & \text{if } \|\Delta \bar{\mathbf{S}}_{rst}\| \geq R_{rs} \quad (\text{case } \mu = 0), \\ -\frac{t-s}{2} \|\bar{\mathbf{S}}_{st}\|^2 + Q_s + \beta + \frac{t-s}{2} (\|\Delta \bar{\mathbf{S}}_{rst}\| - R_{rs})^2 & \text{if } \|\Delta \bar{\mathbf{S}}_{rst}\| < R_{rs}, \end{cases}$$

where:  $R_{rs} = \sqrt{\|\bar{\mathbf{S}}_{rs}\|^2 + 2\bar{Q}_{rs}}$ . Index  $s$  is pruned at time  $t$  by index  $r$  when this maximum is greater than  $Q_t + \beta$ . If no change is in the data, we get the value:

$$\frac{1}{2} \left( \sqrt{\|\Delta \bar{\mathbf{S}}_{rst}\|^2} - \sqrt{\frac{r}{s} \|\Delta \bar{\mathbf{S}}_{0rs}\|^2} \right)^2 - \frac{1}{2} \left( \sqrt{\frac{s}{t} \|\Delta \bar{\mathbf{S}}_{0st}\|^2} \right)^2.$$

The proof is straightforward by derivation of the dual. The proof for the no-change case uses variance expressions in Appendix A. For such case, we prune when

$$\sqrt{\frac{r}{s} \|\Delta \bar{\mathbf{S}}_{0rs}\|^2} - \|\Delta \bar{\mathbf{S}}_{rst}\| > \sqrt{\frac{s}{t} \|\Delta \bar{\mathbf{S}}_{0st}\|^2}.$$

With  $d$  (independent) constraints, it is still possible to get the maximum value of the dual with simple linear algebra, but only if the maximum is strictly inside the domain of definition of the dual. If the obtained maximum is outside, solving inside the positive orthant for  $\mu$  is a challenging problem.

#### 4.4 Maximum of the non-constrained decision function

The decision function admits a closed-form expression for its critical point thanks to its simple structure. However, this point may fall outside the positive orthant. Enforcing this constraint analytically remains intractable at present, necessitating an iterative optimisation algorithm.

**Proposition 10.** *In d-variate model the maximum value of the decision function  $\max_{\mathbf{x} \in \mathbb{R}^q} \mathbb{D}(\mathbf{x})$  with  $q$  constraints is the solution of the following system of equations with  $\mathbf{x} \in \mathbb{R}^q$  and instrumental variable  $\mathbf{y} \in \mathbb{R}^d$ :*

$$\begin{cases} (\Delta \bar{\mathbf{S}}_{\bullet st}) \mathbf{x} = \nabla A(\mathbf{y}) - \bar{\mathbf{S}}_{st}, \\ (\Delta \bar{\mathbf{S}}_{\bullet st})^T \mathbf{y} = -\Delta \bar{Q}_{\bullet st}, \end{cases}$$

where  $(\Delta \bar{\mathbf{S}}_{\bullet st})$  is the matrix of size  $d \times |\mathcal{R}|$  gathering all vectors  $\Delta \bar{\mathbf{S}}_{rst}$  in columns and  $\Delta \bar{Q}_{\bullet st} \in \mathbb{R}^{|\mathcal{R}|}$ . If  $q = d$  and the matrix  $(\Delta \bar{\mathbf{S}}_{\bullet st})$  is invertible, we get the solution:

$$\mathbf{x} = (\Delta \bar{\mathbf{S}}_{\bullet st})^{-1} \left( \nabla A \left( -((\Delta \bar{\mathbf{S}}_{\bullet st})^T)^{-1} \Delta \bar{Q}_{\bullet st} \right) - \bar{\mathbf{S}}_{st} \right). \quad (18)$$

The explicit solution (18) is promising, but it fails to incorporate constraint inequalities, particularly the simple positivity constraint on  $\mathbf{x}$ . The following section will also consider pruning capacity for multivariate time series. However, unlike the univariate case, computational efficiency is not addressed here and is reserved for future work.

## 5 Simulation study

This section presents DUST’s performance on simulated data and compares its efficiency against state-of-the-art algorithms. (This first draft version does not include the study of multivariate signals.)

### 5.1 Framework

#### 5.1.1 Data

We simulate time series of length  $\exp(n)$ , where  $n$  is a sequence of regularly spaced values between  $\log 10^2$  and  $\log 10^8$ . We have 8 available models from the exponential family: Gauss (G), Poisson (P), Exponential (E), Geometric (G), Bernoulli (Be), Binomial (Bi), Negative Binomial (NB), and Variance (V). The simulated data consists of alternating two segments of length  $k$ . We often set  $k = n$ , ensuring that only the first segment is observed. This no-change regime presents the most challenging scenario for pruning and computational efficiency. The typical parameters used are outlined in Table 1.

Table 1: Parameter values for simulations

Model	Penalty scale factor	Parameter	Values
gauss	1	$\mu$ ( $\sigma = 1$ fixed)	{0, 1}
poisson	2/3	$\lambda$	{3, 4}
exponential	3/4	$\lambda$	{1, 0.5}
geometric	2/3	$p$	{0.5, 0.7}
bernoulli	2/3	$p$	{0.5, 0.7}
binomial	1/6	$p$	{0.5, 0.7}
negative binomial	1/10	$p$	{0.5, 0.7}
variance	1	$\sigma$ ( $\mu = 0$ fixed)	{1, 2}

When considering  $d$ -variate time series, we generate  $d$  time series with identical change-point locations using the previous model and concatenate the  $d$  copies. Further, we study three different configurations. (i) (with  $k = n$ ) Pruning and time capacity of DUST over time and against competitors. (ii) (with  $k = n$ ) Pruning and time robustness of DUST with respect to the penalty value at fixed data length (iii) Pruning and time robustness in the presence of changes (against competitors). For the multivariate setting, we also study these configurations with respect to the dimension.

We run each configuration 50 to 100 times at fixed data length depending on computational cost.

### 5.1.2 Baseline and parameters

We compare DUST to the PELT<sup>2</sup> [20] and FPOP<sup>3</sup> [23] algorithms for univariate data. In the multivariate setting, there is only one competitor and only for the Gaussian cost: GeomFPOP<sup>4</sup> [26]. Each algorithm is provided with the same simulated data to ensure proper head-to-head comparisons.

The penalty factor for each model is calibrated to achieve similar segmentation behaviour across the range of penalty for data with no change point. With the Gaussian model as a baseline, we apply a scaling factor to the penalty for each model. This scaling factor is determined by finding the smallest penalty value where change-point detection correctly produces no change point on data of length  $10^3$ . The ratio between this value and the corresponding value in the Gaussian model gives us the scaling factor. For example, suppose we execute DUST on Gaussian data of length  $n$  with penalty  $2 \log n$ . In that case, an equivalent Negative binomial simulation is done with penalty  $2a \log n$  with  $a = 1/10$  (see Table 1).

### 5.1.3 Metrics

We measure the execution time for parsing each simulated time series to compare the computational cost of each algorithm in each configuration. We further record the number of candidate indices during the execution of each DUST. The remaining number of indices at the end of the algorithm is a measure of the algorithm's pruning capacity.

### 5.1.4 DUST variants

DUST algorithm comes with variants depending on the dual evaluation algorithm and index selection method used.

The dual evaluation can be made:

1. at its maximum using the closed formula (for one-parametric cost functions only);
2. at a random point uniformly drawn in the dual domain;
3. at zero (equivalent to PELT test);
4. at its maximum using the Quasi-Newton algorithm.

Evaluation in the one-parameter-cost case is naturally performed using a closed-form formula. Other methods are tested for multi-parameter cases to balance time complexity with pruning efficiency. Index selection for building the dual comprises two parts: indices below the index  $s$ , and indices above. Indices can be chosen randomly (uniformly) or deterministically (the largest available below  $s$  and the smallest available after  $s$ ). Here, we consider only the deterministic case.

---

<sup>2</sup><https://cran.r-project.org/web/packages/changepoint/index.html>

<sup>3</sup><https://cran.r-project.org/web/packages/fpopw/index.html>

<sup>4</sup><https://github.com/computorg/pishchagina-change-point>

## 5.2 Univariate signals

### 5.2.1 Pruning capacity

Figure 3 displays the number of non-pruned indices over time for one time series of length  $10^4$  (left column) and length  $10^8$  (right column) and the median and confidence intervals (over 100 repetitions). We present the results for two cost models: Gauss, top row; Negbin, bottom row. DUST shows a persistent pruning efficiency that is robust to the number of data points, with low variations being recorded. With  $n = 10^8$ , the number of remaining indices is 50 with the exact evaluation method, and 500 with the random method.

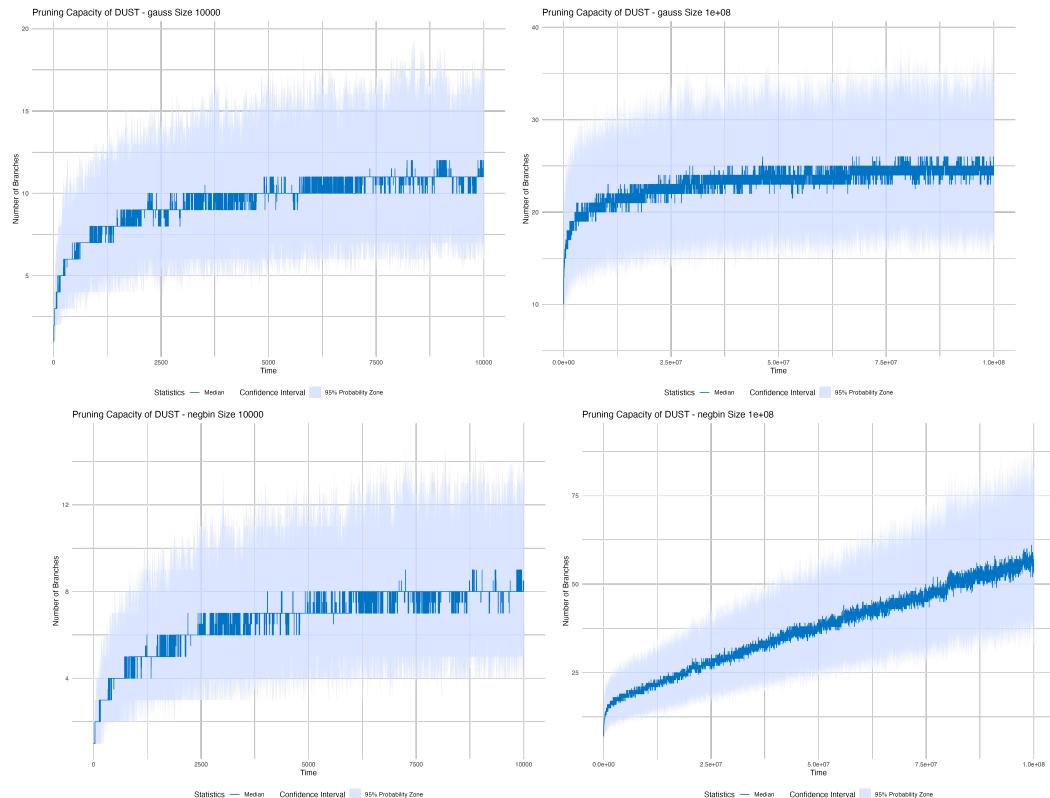


Figure 3: Median number of remaining candidate indices at each time for 100 DUST simulations for the Gaussian (top row) and negative binomial (bottom row) models, data lengths  $10^4$  (left column) and  $10^8$  (right column). The shaded region shows the interval between the 0.025 % and 0.975 % quantiles.

We plot the number of remaining candidate indices for different data lengths on a log-log scale in Figure 4. We run 100 simulations for each data length among 100 regularly spaced (in the logarithmic scale) lengths between  $10^2$  and  $10^6$ . The shaded region shows the interval between the 0.025 % and 0.975 % quantiles. The solid line shows the fitted linear regression with a 95% prediction confidence interval as dotted lines on either side. The slope value of either model suggests that the mean number of candidate indices remaining upon exit of the DUST algorithm is of order  $n^\alpha$ ,  $\alpha < 0.15$ , which indicates a time complexity of order  $\mathcal{O}(n^{1+\alpha})$ , with minor variations in the

coefficient depending on the model being considered.

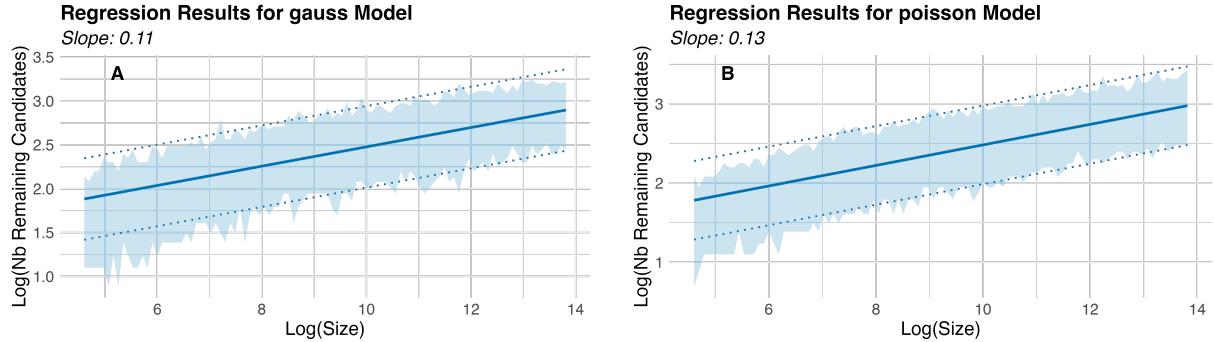


Figure 4: Log-log comparison of number of remaining indices over times as a function of data size for the Gaussian (panel A) and Poisson (panel B) models. We run 100 simulations on 100 data lengths between  $10^2$  and  $10^6$ , regularly spaced in the log-scale. The shaded region shows the interval between the 0.025 % and 0.975 % quantiles. The solid line shows the fitted linear regression with a 95% prediction confidence interval as dotted lines on either side. All simulations are performed on data without a change point.

### 5.2.2 Time competition

Figure 5 displays a log-log comparison of execution times between DUST and FPOP as a function of data size for the Gaussian (panel A) and Poisson (panel B) models. 100 simulations were performed on 100 data lengths between  $10^2$  and  $10^6$ , regularly spaced in the log-scale. The shaded regions show the interval between each algorithm's 0.025 % and 0.975 % quantiles. The solid lines show the fitted linear regression for each algorithm for sizes 3125 through  $10^6$ , with 95% prediction confidence intervals as dotted lines on either side. The threshold at 3125 was introduced as the relation between execution time and data length stabilizes by that point. Analysis of the interaction coefficient between data length and algorithm used shows that FPOP's slope is significantly greater than DUST's under either model, which implies that DUST's execution time scales better on larger data. Under the Gaussian model, DUST runs much slower on smaller data sizes, but DUST reliably beats FPOP beyond a break-even point at around 5,000, and is on average 1.19 times faster on data of length  $10^6$ . On non-Gaussian models such as the Poisson model, however, the DUST algorithm outperforms FPOP even at the smallest sizes, achieving speeds 5.88 times up to 8.15 times greater than FPOP, at sizes  $10^2$  and  $10^6$  respectively. Comparison of execution time is presented in Figure 6 for relatively small data length.

### 5.2.3 Pruning exploration

Figure 7 displays the median number of candidate indices upon exit of the DUST algorithm for 100 executions on fixed-length data with no change point for 100 different penalty factors, under the Gaussian (panel A) and Poisson (panel B) models. Data length is  $n = 10^7$  and penalty factors are of the form  $\beta = a \log n$ , with 100 different

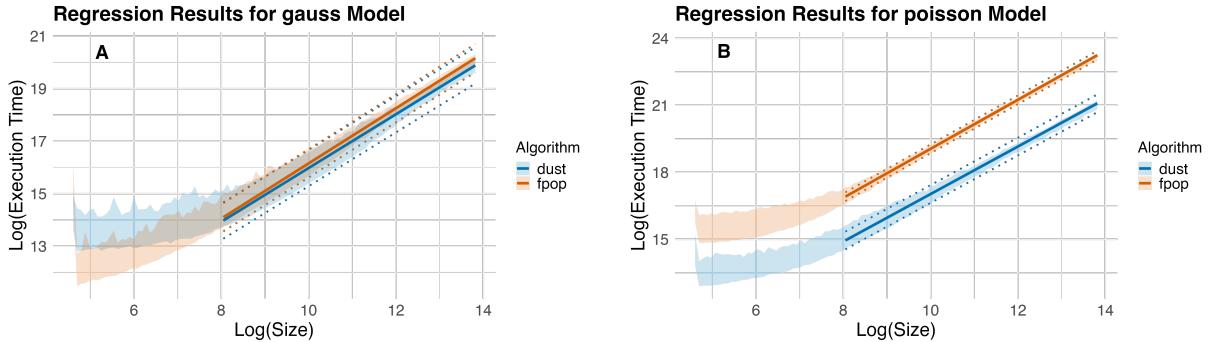


Figure 5: Log-log comparison of execution times between DUST and FPOP as a function of data size for the Gaussian (panel A) and Poisson (panel B) models. 100 simulations were performed on  $10^2$  data lengths between 100 and  $10^6$ , regularly spaced in the log-scale. The shaded regions show the interval between the 0.025 % and 0.975 % quantiles. The solid lines show the fitted linear regression for each algorithm for sizes 3125 through  $10^6$ , with 95% prediction confidence intervals as dotted lines on either side. All simulations are performed on data without a change point.

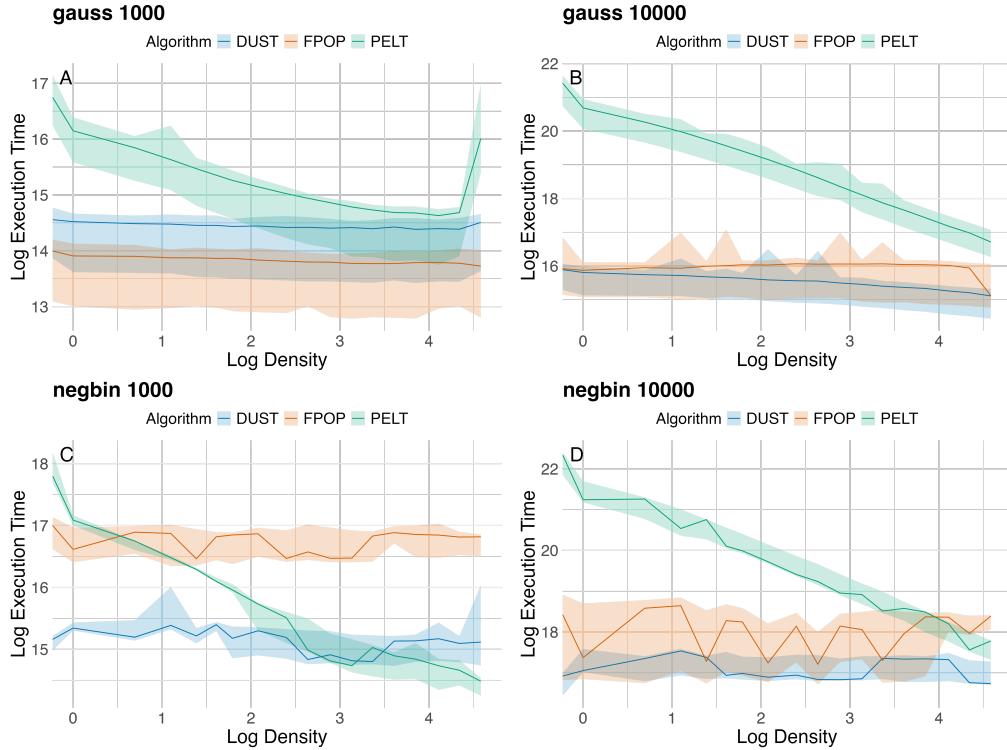


Figure 6: Log-log comparison of execution times for varying numbers of change points in input data of lengths  $10^3$  (panels A, C) and  $10^4$  (panels B, D) under the Gaussian (panels A, B) and the Negative binomial (panels C, D) models.

$a$  values regularly spaced between 0.001 and 20 in the logarithmic scale. The shaded region shows the interval between the 0.025 % and 0.975 % quantiles. Lower penalty

values produce a high number of change points. Values beyond 0.75 predominantly produce the correct number of change points. As for the number of remaining candidates, we observe a polynomial relation with the logarithmic values of  $a$  up to  $a = 1$ , and a negative linear relation beyond that point, with a maximum median number of candidates of 24 under the Gaussian model, and 28 under the Poisson model. This suggests that pruning is strong regardless of the penalty chosen, which guarantees a strongly reduced computational cost even on non-normalised data or with a poorly tuned penalty factor.

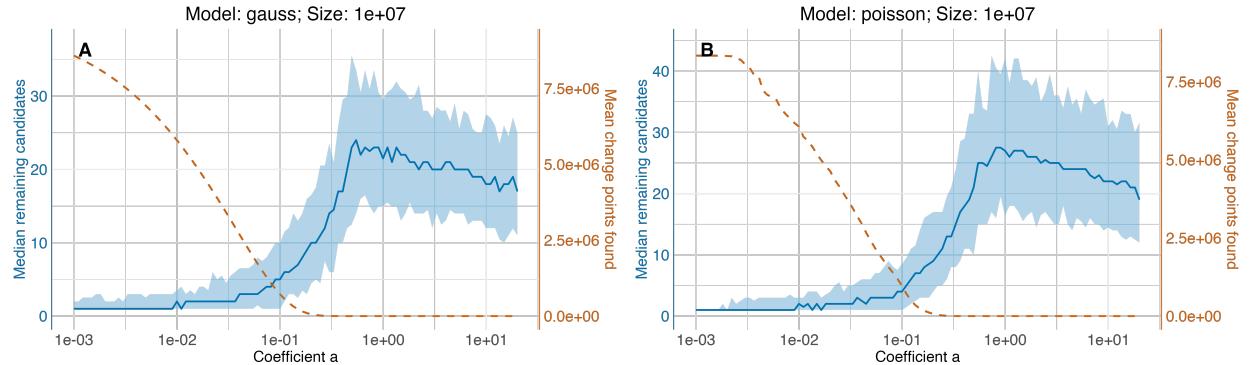


Figure 7: Median number of candidate indices upon exit of the DUST algorithm for 100 executions on fixed-length data with no change point for 100 different penalty factors, under the Gaussian (panel A) and Poisson (panel B) models. Data length is  $n = 10^7$  and Penalty factors are of the form  $\beta = a_i \log n$ , with 100 different  $a$  values regularly spaced between 0.001 and 20 in the logarithmic scale. The shaded region shows the interval between the 0.025 % and 0.975 % quantiles.

### 5.3 Multivariate signals

Coming soon!

### 5.4 The DUST package

DUST package is available on GitHub<sup>5</sup>.

## 6 Application to mouse monitoring

This section explores the use of DUST on a real-world application taken from [21].

**Context** The neuromuscular junction (NMJ) is the synapse responsible for the chemical transmission of electrical impulses from motor neurons to muscle cells. Studying the NMJ is crucial because it explains how nerves communicate with muscles to move.

---

<sup>5</sup><https://github.com/vrunge/dust.git>

Since this synapse is a disease-prone synapse [10], studying the NMJ also helps understand and treat disorders like myasthenia gravis and congenital myasthenic syndrome.

In [21], the authors find a correlation between a partial inhibition of the NMJ and the level of muscle fatigue in mice. They measure muscle fatigue by quantifying the number and duration of active periods over several days.

**Data** A typical approach to estimate resting and active states is to use a force platform, a scale measuring the ground reaction force generated by the mouse in its cage. To compare mice, [21] uses several features, including the total duration of activity, the mean duration of activity, and the number of active periods. The procedure to detect whether a mouse is active or inactive consists of several steps, including filtering and simple signal transformations. Changes are detected on a simplified signal because the signals have too many samples.

Here we consider two mice with two different genetic modifications (Mouse ColQ and Mouse A7), monitored over 2 nights and 1 day. The cage is placed on four force sensors, and we record the evolution over time of the sum of the ground reaction forces at 10 Hz. If the sum of forces is constant (low variance), the mouse is not active (resting state). If the sum varies, the mouse is active, regardless of the activity.

**Proposed approach and results** In this section, we have a more straightforward approach: we use DUST on the raw data. Each segment is then classified as active or inactive by thresholding the variance. Since DUST is fast, we can process 12 hours ( $\sim 400k$  samples) of time series in seconds without much preprocessing. The statistical model is piecewise Gaussian with a fixed mean equal to 0 and piecewise constant variance.

From a qualitative standpoint, Figure 8 shows parts of the time series and the segmentation into active/resting stages. Activity is defined by a high variance.

From Table 2, which summarizes with simple features the level of activity of each mouse, we draw two conclusions. First, mice are more active during the night by a margin. They have longer stretches of activity and spend more time in an active state. We recover here a well-known fact in mouse behaviour analysis. Second, Mouse A7 is more active than Mouse ColQ during the night, but equally active during the day. This observation is in line with the results of [21], albeit with a smaller number of mice. However, our approach is simpler than the one in [21] because we can find variance changes on long time series in seconds.

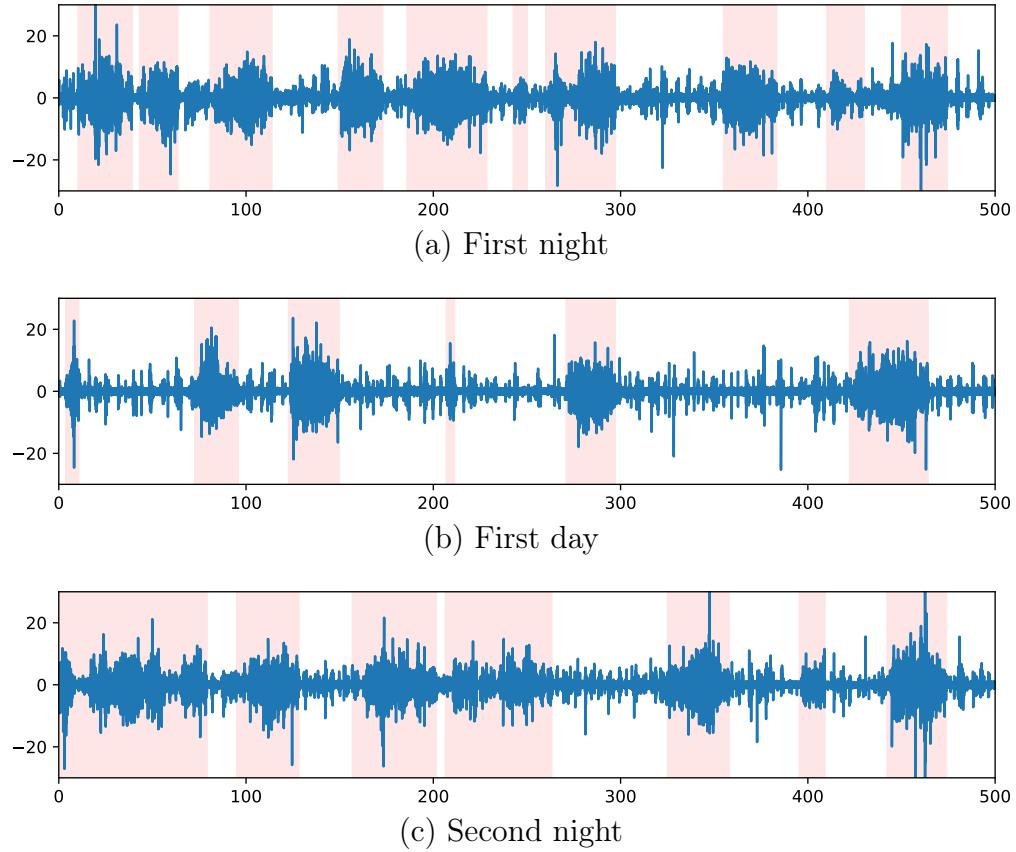


Figure 8: For Mouse ColQ, the activity measured by a force platform. The x-axis is the time (in minutes), and the y-axis is the ground reaction forces generated by a mouse in its cage (in arbitrary units). Red areas indicate activity periods found by our change-point approach; others are rest periods.

Table 2: Summary of activity segments for each period.

	# of active segments	Avg. active duration (min)	Total active time (h)
Mouse ColQ			
First night	18	21.22	6.37
First day	21	11.03	3.86
Second night	16	24.84	6.62
Mouse A7			
First night	20	25.82	8.61
First day	21	9.03	3.16
Second night	15	35.73	8.93

## 7 Conclusion

A large class of change-point problems can be solved exactly using dynamic programming algorithms with functional pruning rules. Here, pruning refers to identifying the functions  $q_t^i(\cdot)$  that are dominated by the minimum operator  $\min_j \{q_t^j(\cdot)\}$  in problems with a functional structure that preserves pointwise ordering. Specifically, if there exist indices  $i, j$  and a parameter  $\boldsymbol{\theta}$  such that  $q_t^i(\boldsymbol{\theta}) < q_t^j(\boldsymbol{\theta})$ , this inequality remains valid for all  $u > t$ .

In the DUST framework, pruning is reformulated as a (generally non-convex) optimisation problem, along with its relaxations obtained by removing constraints, which we analyse via their dual representations. Pruning occurs when the computed value exceeds a specified threshold. We prove there is no duality gap when the number of constraints does not exceed the dimension of the parameter space of the cost function derived from the exponential family.

We focused on the one-parameter case, essentially corresponding to univariate data, for which an inequality-based pruning test as simple as PELT can be designed. Unlike PELT, however, this test remains highly efficient across all change-point regimes, whether or not there are many or few changes. This new approach can be viewed as a natural extension of PELT, where PELT corresponds to the dual evaluation at zero. In this simple case, DUST evaluates a decision function with a closed formula, equivalent to the maximum evaluation of the dual.

FPOP-like algorithms can efficiently tackle one-parametric problems. DUST has been proven to be very efficient in our simulation study while saving more functions. This is due to the iterative root searching algorithm used by FPOP for most of the non-Gaussian cost in updating the functional cost (saving this piecewise function by a succession of intervals).

DUST extensions to higher dimensions are straightforward and still efficient in small dimensions with a unique constraint. An optimal approach considering multiple constraints still needs to be designed using new theoretical tools or leading an extensive simulation study. Current results are promising; this is left for future work.

Using a DUST pruning rule is a natural step that could be introduced in many change-point detection algorithms to efficiently prune indices in all change-point regimes. Even if the explicit maximum of the dual or decision function cannot be explicitly derived, evaluating at one point only (e.g., randomly) is a simple rule, fast to test, with high potential for pruning. The exact evaluation in the one-parameter case provides a very efficient inequality-based test; this test should be further studied with the hope of proving time complexity bounds for very generic underlying data processes.

## A Useful relations between means

**Proposition 11.** *For all  $i \in \{1, \dots, t-1\}$  we have*

$$\begin{aligned} tV(y_{0t}) - (t-i)V(y_{it}) - iV(y_{0i}) \\ = \frac{(t-i)i}{t} (\bar{y}_{it} - \bar{y}_{0i})^2 \end{aligned} \tag{19}$$

$$= \frac{ti}{t-i} (\bar{y}_{0t} - \bar{y}_{0i})^2 \tag{20}$$

$$= \frac{t(t-i)}{i} (\bar{y}_{0t} - \bar{y}_{it})^2 \tag{21}$$

$$= i(\bar{y}_{0t} - \bar{y}_{0i})^2 + (t-i)(\bar{y}_{0t} - \bar{y}_{it})^2 \tag{22}$$

*Proof.* For any  $(a, b) \in (\mathbb{N}^*)^2$  with  $a < b$  we easily get the following expressions:

$$\bar{y}_{0b} - \bar{y}_{0a} = \frac{b-a}{b} (\bar{y}_{ab} - \bar{y}_{0a}) , \tag{23}$$

$$\bar{y}_{0b} - \bar{y}_{ab} = \frac{a}{b} (\bar{y}_{0a} - \bar{y}_{ab}) . \tag{24}$$

With  $a = i$  and  $b = t$ , we derive from (19) the relation (20) using (23) and also derive relation (21) using (24). Expression (22) is the result of  $\frac{t-i}{t}(20) + \frac{i}{t}(21)$ . Therefore, we only need to prove relation (19). Expanding the squares in variance expressions we get:

$$\begin{aligned} S &= tV(y_{0t}) - (t-i)V(y_{it}) - iV(y_{0i}) \\ &= -t(\bar{y}_{0t})^2 + (t-i)(\bar{y}_{it})^2 + i(\bar{y}_{0i})^2 \\ &= i((\bar{y}_{0i})^2 - (\bar{y}_{0t})^2) + (t-i)((\bar{y}_{it})^2 - (\bar{y}_{0t})^2) \\ &= i(\bar{y}_{0i} - \bar{y}_{0t})(\bar{y}_{0i} + \bar{y}_{0t}) + (t-i)(\bar{y}_{it} - \bar{y}_{0t})(\bar{y}_{it} + \bar{y}_{0t}) . \end{aligned}$$

Using again relations (23) and (24), we have

$$S = \frac{(t-i)i}{t} (\bar{y}_{0i} - \bar{y}_{it})(\bar{y}_{0i} + \bar{y}_{0t}) + (\bar{y}_{it} - \bar{y}_{0i})(\bar{y}_{it} + \bar{y}_{0t}) ,$$

and we get the expression (19) by simplifying this last expression.  $\square$

## B Proof of Proposition 1

If we would remove index  $s_0$  from  $\mathcal{T}_t$  while still having a point  $\boldsymbol{\theta}_0 \in \Theta$  such that:

$$\begin{cases} Q_t(\boldsymbol{\theta}_0) = q_t^{s_0}(\boldsymbol{\theta}_0) < q_t^s(\boldsymbol{\theta}_0) & \text{for } s \neq s_0 , \\ Q_t + \beta > q_t^{s_0}(\boldsymbol{\theta}_0) , \end{cases} \tag{25}$$

we show that removing  $s_0$  could lead to an under-optimal solution for (6) and then for (4) at some data time  $t_0 > t$ . We will choose data points after time  $t$  to create such a solution, i.e., with  $\arg \min_{\boldsymbol{\theta}} Q_{t_0}(\boldsymbol{\theta})$  attained by a value on function  $q_{t_0}^{s_0}$ . To that end, for all further iteration  $t' > t$  we choose data points  $y_{t'}$  with unitary cost:

$$c(y_{t'}; \boldsymbol{\theta}) = A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{T}(y_{t'}) = A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \nabla A(\boldsymbol{\theta}_0) = c(\boldsymbol{\theta}).$$

*Case 1: non-optimal indices  $s \neq s_0 \in \{0, \dots, t-1\}$ .* By continuity of function  $Q_t(\cdot)$  there exists a ball centred on  $\boldsymbol{\theta}_0$  with small radius  $\epsilon$ ,  $B(\boldsymbol{\theta}_0, \epsilon)$ , such that  $Q_t(\boldsymbol{\theta}) = q_t^{s_0}(\boldsymbol{\theta}) < q_t^s(\boldsymbol{\theta})$  for points  $\boldsymbol{\theta}$  in this ball. We consider  $\boldsymbol{\theta}$  outside  $B(\boldsymbol{\theta}_0, \epsilon)$ . By strong convexity of  $A$ , there exists a constant  $a > 0$  such that:

$$c(\boldsymbol{\theta}) > c(\boldsymbol{\theta}_0) + a\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|^2 \geq c(\boldsymbol{\theta}_0) + a\epsilon^2,$$

for all  $\boldsymbol{\theta}$  outside  $B(\boldsymbol{\theta}_0, \epsilon)$ . We also write  $\Delta = q_t^{s_0}(\boldsymbol{\theta}_0) - Q_t \geq 0$ . We consider  $T = \lceil \frac{\Delta}{a\epsilon^2} \rceil$ , then:

$$Tc(\boldsymbol{\theta}) > Tc(\boldsymbol{\theta}_0) + aT\epsilon^2 \geq Tc(\boldsymbol{\theta}_0) + \Delta.$$

Thus,

$$Q_t + Tc(\boldsymbol{\theta}) > Q_t + Tc(\boldsymbol{\theta}_0) + \Delta = Tc(\boldsymbol{\theta}_0) + q_t^{s_0}(\boldsymbol{\theta}_0),$$

for all  $\boldsymbol{\theta}$  outside  $B(\boldsymbol{\theta}_0, \epsilon)$ . Thus, the minimum of  $Q_{t_0}(\cdot) = Q_{t+T}(\cdot)$  can only be in the ball. All points outside the ball corresponding to indices  $s \neq s_0$  in  $\{0, \dots, t-1\}$  are not optimal, which proves case 1.

*Case 2: non-optimal indices  $s \in \{t, \dots, t+T\}$ .* We first consider  $s = t$ . The second assumption in Equation (25) ( $Q_t + \beta > q_t^{s_0}(\boldsymbol{\theta}_0)$ ) ensures that values visible for function  $q_t^{s_0}$  are not all too large and pruned by function  $q_{t+1}^t$ . We have:

$$\begin{aligned} q_{t_0}^t(\boldsymbol{\theta}) &= Q_t + Tc(\boldsymbol{\theta}) + \beta > q_t^{s_0}(\boldsymbol{\theta}_0) + Tc(\boldsymbol{\theta}), \\ &\geq q_t^{s_0}(\boldsymbol{\theta}_0) + Tc(\boldsymbol{\theta}_0), \\ &= q_{t+T}^{s_0}(\boldsymbol{\theta}_0) = q_{t_0}^{s_0}(\boldsymbol{\theta}_0). \end{aligned}$$

All points on the function  $q_{t_0}^t$  are higher than  $q_{t_0}^{s_0}(\boldsymbol{\theta}_0)$ , consequently index  $t$  can not be optimal. For  $s \in \{t+1, \dots, t+T\}$ , we have for index  $\bar{s} < t$ :

$$\min_{\boldsymbol{\theta}} c(y_{\bar{s}t}; \boldsymbol{\theta}) + (s-t)c(\boldsymbol{\theta}_0) \leq \min_{\boldsymbol{\theta}} c(y_{\bar{s}s}; \boldsymbol{\theta}),$$

or,

$$\min_{\boldsymbol{\theta}} c(y_{\bar{s}t}; \boldsymbol{\theta}) + Tc(\boldsymbol{\theta}_0) \leq \min_{\boldsymbol{\theta}} c(y_{\bar{s}s}; \boldsymbol{\theta}) + (t_0-s)c(\boldsymbol{\theta}_0).$$

Using for  $\bar{s}$  the best index for last change point in  $Q_s$  (written  $\bar{s}_0$ ) we have  $Q_{\bar{s}_0} + \min_{\boldsymbol{\theta}} c(y_{\bar{s}_0s}; \boldsymbol{\theta}) + \beta = Q_s$  and we get:

$$Q_{\bar{s}_0} + \min_{\boldsymbol{\theta}} c(y_{\bar{s}_0t}; \boldsymbol{\theta}) + \beta + Tc(\boldsymbol{\theta}_0) \leq Q_s + (t_0-s)c(\boldsymbol{\theta}_0).$$

By definition of  $Q_t$  (see (6)) we have  $Q_t \leq Q_{\bar{s}_0} + \min_{\boldsymbol{\theta}} c(y_{\bar{s}_0t}; \boldsymbol{\theta}) + \beta$  and eventually:

$$\begin{aligned} q_{t_0}^t(\boldsymbol{\theta}_0) &= Q_t + Tc(\boldsymbol{\theta}_0) + \beta \leq Q_s + (t_0-s)c(\boldsymbol{\theta}) + \beta, \\ &\leq Q_s + (t_0-s)c(\boldsymbol{\theta}) + \beta = q_{t_0}^s(\boldsymbol{\theta}). \end{aligned}$$

The global minimum of the function  $Q_{t_0}(\cdot)$  cannot be attained in index  $s$ , as we proved that it cannot either in index  $t$ . By exclusion of all indices, except  $s_0$ , we have shown that the global minimum is attained at a value of the function  $q_{t_0}^{s_0}$ ; consequently, removing  $s_0$  from  $\mathcal{T}_t$  would lead to an under-optimal solution.

## C Worst-case time series

### C.1 Proof of Proposition 2

The minimum of the function  $q_t^s$  is defined as  $m_t^s$ . A sufficient condition for no pruning is to meet the condition  $m_n^0 = m_n^1 = \dots = m_n^{n-1}$  obtained at distinct values. This is the case for a strictly increasing time series. We assume that the global minimum is always attained for index 0 and that the time series is increasing; thus, using the variance notation in Appendix A, we solve:

$$m_n^0 = nV(y_{0n}) = (n-t)V(y_{tn}) + tV(y_{0t}) + \beta, \quad \text{for all } t = 1, \dots, n-1.$$

By relations in Appendix A, this leads to equations

$$\frac{nt}{n-t} (\bar{y}_{0n} - \bar{y}_{0t})^2 = \beta, \quad \text{for all } t = 1, \dots, n-1.$$

With  $t = 1$ , we obtain  $\bar{y}_{0n} = \sqrt{\beta \frac{n-1}{n}}$  having chosen first data value  $y_1 = 0$ , and therefore, with other indices, we get:

$$y_1 + \dots + y_t = t \sqrt{\beta \frac{n-1}{n}} - \sqrt{\beta \frac{t(n-t)}{n}}.$$

We then find the proposed expression easily. We need to verify that the time series is strictly increasing and that the global minimum is associated with the first data point at each time step. Indeed,

$$y_{t+1} - y_t = \sqrt{\frac{\beta}{n}} \left( -\sqrt{(t+1)(n-t-1)} + 2\sqrt{t(n-t)} - \sqrt{(t-1)(n-t+1)} \right)$$

is positive as  $t \mapsto -\sqrt{t(n-t)}$  is strictly convex on interval  $[0, n]$ . The last change point is always associated with the first data point, as we have an increasing sequence of data with equality of the minima at the ending instant  $n$  ( $m_n^0 = m_n^1 = \dots = m_n^{n-1}$ ).

### C.2 No pruning in the exponential family

We say that a time series increases if it is increasing coordinate by coordinate. Under some mild assumptions, we can build a no-pruning example of increasing time series in a more general setting.

**Proposition 12.** *We consider costs of type "aA( $\theta$ ) + b ·  $\theta$  + c" derived from the natural exponential family with continuous density. We consider the following equations in variable  $x$ :*

$$E(t) : \quad g_t(x) = \frac{\beta}{n} + \mathcal{D}^*(Y), \quad t = 1, \dots, n-1, \tag{26}$$

with

$$\begin{cases} \mathcal{D}^*(x) = x \cdot (\nabla A)^{-1}(x) - A((\nabla A)^{-1}(x)), & x \in \Omega, \\ g_t(x) = \frac{t}{n} \mathcal{D}^*(x) + \frac{n-t}{n} \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right), & x \in \Omega_t \subset \Omega. \end{cases}$$

If  $Y \in \Omega_0 = \cap_{t=1}^{n-1} \Omega_t$  and  $\beta$  is chosen such that:

$$\sup_{x \in \Omega_0, x \leq Y} \{g_1(x)\} > \frac{\beta}{n} + \mathcal{D}^*(Y),$$

there exists an increasing sequence of data  $y_1, \dots, y_n$  with zero pruning solving (4) which verifies:

$$\begin{cases} Y = \bar{y}_{0n} = \frac{1}{n} \sum_{t=1}^n y_t, \\ x = s_t \text{ solution of } E(t), \\ y_t = ts_t - (t-1)s_{t-1} \quad t = 1, \dots, n, \quad \text{with } s_0 = 0. \end{cases}$$

*Proof.* The proof is an extension of the evidence of Proposition 2 and is based on a similar strategy. We look for data configurations for which the minimum of functions  $q_t^s$ , denoted  $m_t^s$ , are all the same:  $m_n^0 = m_n^1 = \dots = m_n^{n-1}$ . There is no possible pruning if we find an increasing time series realising such conditions. For density from the natural exponential family, we get the relations:

$$m_n^t = Q_t + (n-t)(A(\rho_t) - \bar{y}_{tn} \cdot \rho_t) + \beta = Q_0 + n(A(\rho_0) - \bar{y}_{0n} \cdot \rho_0) + \beta = m_n^0,$$

with  $\rho_i = (\nabla A)^{-1}(\bar{y}_{in})$ . We have also  $Q_0 = 0$  and  $Q_t = t(A((\nabla A)^{-1}(\bar{y}_{0t})) - \bar{y}_{0t} \cdot (\nabla A)^{-1}(\bar{y}_{0t})) + \beta$  as data is not split into segments (we have only one large segment and the minimum for  $Q_t(\cdot)$  is attained on  $q_t^0$ ). This leads to:

$$\frac{t}{n} \mathcal{D}^*(\bar{y}_{0t}) + \frac{n-t}{n} \mathcal{D}^*(\bar{y}_{tn}) = \frac{\beta}{n} + \mathcal{D}^*(\bar{y}_{0n}),$$

with  $\mathcal{D}^*(x) = x \cdot (\nabla A)^{-1}(x) - A((\nabla A)^{-1}(x))$ . We get Equations (26) denoted  $E(t)$  with  $x = \bar{y}_{0t}$ ,  $Y = \bar{y}_{0n}$  and relation  $\frac{nY - tx}{n-t} = \bar{y}_{tn}$ . It remains to find conditions for the existence of a solution. We define  $g_t(x) = \frac{t}{n} \mathcal{D}^*(x) + \frac{n-t}{n} \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right)$  and study this function. We have:

$$\begin{cases} g_t(Y) = \mathcal{D}^*(Y), \\ \nabla g_t(Y) = 0, \\ g_t \text{ strictly convex}, \\ g_{t+1} > g_t \text{ except in point } Y. \end{cases}$$

Evaluation of  $g_t$  in  $Y$  is obvious, for the gradient we get:

$$\nabla g_t(x) = \frac{t}{n} \left( \nabla \mathcal{D}^*(x) - \nabla \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right) \right),$$

and then  $\nabla g_t(Y) = 0$ . Moreover, function  $g_t$  is strictly convex:  $\nabla^2 g_t(x) = \frac{t}{n} \nabla^2 \mathcal{D}^*(x) + \frac{t^2}{n(n-t)} \nabla^2 \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right)$ ; its Hessian is definite positive since the Hessian of  $\mathcal{D}^*$  is. We prove  $g_{t+1} > g_t$  considering  $t$  as a continuous variable. We define  $h_x(t) = g_t(x)$  and we have:

$$\nabla h_x(t) = \frac{1}{n} \mathcal{D}^*(x) - \frac{1}{n} \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right) + \frac{Y - x}{n-t} \cdot \nabla \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right).$$

However, by a characterisation of the strict convexity of  $\mathcal{D}^*$  we have for  $x \neq Y$ :

$$\mathcal{D}^*(x) > \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right) + n\frac{Y-x}{n-t} \cdot \nabla \mathcal{D}^*\left(\frac{nY - tx}{n-t}\right),$$

as  $n\frac{Y-x}{n-t} = x - \frac{nY-tx}{n-t}$ . This leads obviously to  $\nabla h_x(t) > 0$  and consequently to  $h_x(t+1) - h_x(t) > 0$  or with the  $g$  notation:  $g_{t+1}(x) > g_t(x)$  for all  $t$  and all  $x \neq Y$ .

Now, with  $\Omega_t$  the convex domain of  $g_t$ , if  $\Omega_0 = \cap_{t=1}^{n-1} \Omega_t$  is not empty, we can choose  $Y$  in  $\Omega_0$ . We also choose  $\beta$  such that  $\sup_{x \in \Omega_0, x \leq Y} \{g_1(x)\} > \frac{\beta}{n} + \mathcal{D}^*(Y)$ . By continuity of  $g_t$  we can find  $x = s_1$  solving  $E(1)$ . Using the fact that  $g_2 > g_1$  and  $\min_{x \in \Omega_0} g_2(x) = \mathcal{D}^*(Y)$  we can now find  $x = s_2 > s_1$  solving  $E(2)$ . Step by step, we build the sequence  $(s_t)$  corresponding to an increasing sequence of means  $\bar{y}_{0t}$ . From this point, it is easy to extract the increasing time series  $(y_t)$ , which concludes the proof.  $\square$

Notice that the multiple independent Poisson model leads to the expression (26) written as:

$$\frac{t}{n}x \log(x) - Y \log(Y) + \left(Y - \frac{t}{n}x\right) \log\left(\frac{Y - \frac{t}{n}x}{1 - \frac{t}{n}}\right) = \frac{\beta}{n},$$

where  $x \in \Omega_0 = \left(0, \frac{n}{n-1}Y\right) = \left(0, \frac{n}{n-1}Y_1\right) \times \cdots \times \left(0, \frac{n}{n-1}Y_d\right) \subset \mathbb{R}^d$ . The set  $\Omega_0$  is not empty and the limit of  $g_1$  in zero leads to the condition for beta:  $\beta \leq n \log\left(\frac{n}{n-1}\right) \left(\sum_{i=1}^d Y_i\right)$ . This bound is much lower than usual values chosen for penalty, bounded by a  $\log(n)$  term [9]. As the Poisson model needs integer data points, the solution we get has to be approximated by  $\tilde{y}_t = \lceil y_t \rceil$  and the resulting data and functional cost  $Q_n(\cdot)$  have some of their indices pruned as illustrated in Figure 9. If  $Y$  is chosen very large, the approximation improves and the obtained time series converges to the Gaussian data with no pruning.

In Table 4 we provide the expression for  $A$ ,  $(\nabla A)^{-1}$ , and its domain  $\Omega$  for a few distributions of the natural exponential family. For instance, the exponential model leads to the expression (26) written as:

$$\log(Y) - \frac{t}{n} \log(x) - \left(1 - \frac{t}{n}\right) \log\left(\frac{Y - \frac{t}{n}x}{1 - \frac{t}{n}}\right) = \frac{\beta}{n},$$

where  $x \in \Omega_0 = \left(0, \frac{n}{n-1}Y\right)$ . For a multiple independent exponential model, the dual  $\mathcal{D}^*$  can be decomposed as the sum of dual functions, one for each dimension:  $\mathcal{D}^*(x) = \sum_{i=1}^d \mathcal{D}_i^*(x_i)$ , with  $x = (x_1, \dots, x_d)^T$ , justifying the following rectangular form for  $\Omega_0 = \left(0, \frac{n}{n-1}Y_1\right) \times \cdots \times \left(0, \frac{n}{n-1}Y_d\right) \subset \mathbb{R}^d$ . The set  $\Omega_0$  is not empty, and the limit of  $g_1$  is infinite; thus, there is no restriction on the  $\beta$  value. A more general result for the exponential family seems possible. Still, it requires solving Equation (26) for  $x = \frac{1}{t} \sum_{i=1}^t \mathbf{T}(y_i)$  and inverting  $\mathbf{T}$  to explicit the time series  $(y_t)$ . This step depends on the chosen distribution and requires additional effort, which is left for further work.

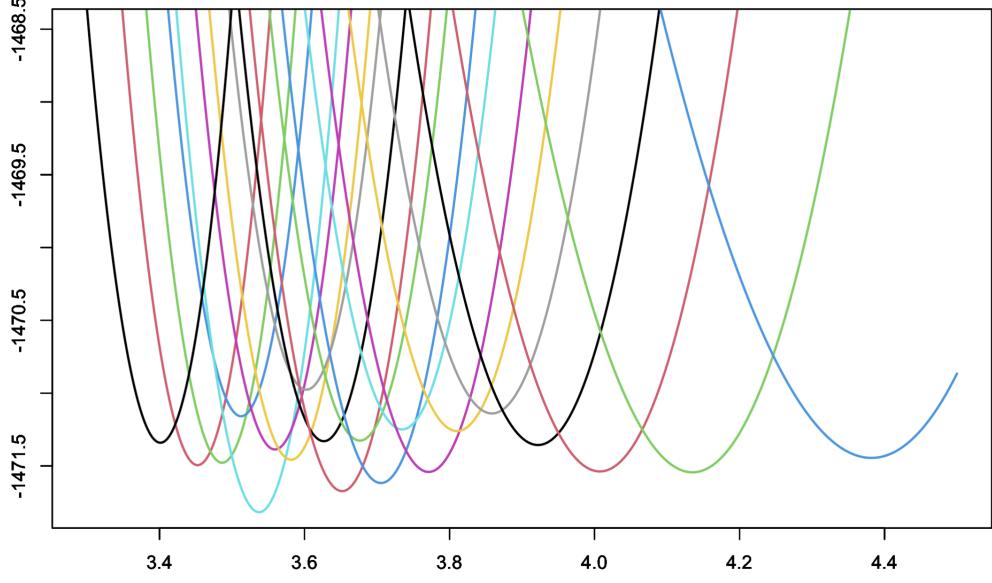


Figure 9: Functional cost  $Q_{20}(\cdot)$  obtained from 20 data points. We identify easily the 20 different inner functions from which the functional cost is built. We choose  $Y = 30$  and  $\beta = 0.995Y \log(n/(n-1))$ . Due to the integer approximation for  $\tilde{y}_t$ , only 14 inner functions are visible in the functional cost.

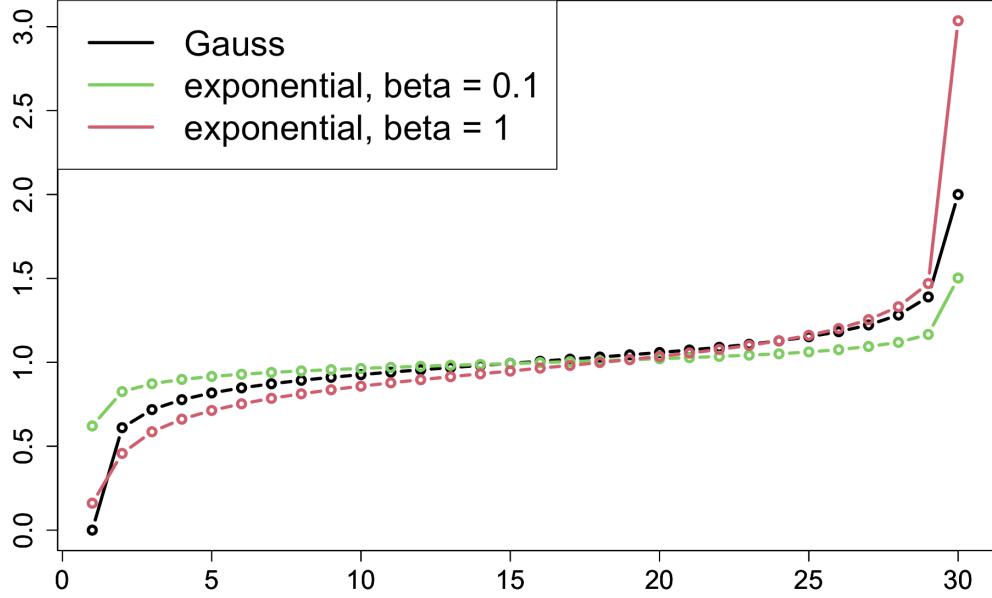


Figure 10: Examples with 30 data points in worst case complexity (no pruning) with Gaussian model (middle curve in black) and two examples with exponential model with two different penalty values (beta). We chose  $Y = 1$ .

## D Proofs of Section 3

### D.1 Proof of Proposition 4

The primal Lagrangian function is given for a non-negative multiplier  $\mu$  by the relation:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}, \mu) &= q_t^s(\boldsymbol{\theta}) + \mu(q_t^s(\boldsymbol{\theta}) - q_t^r(\boldsymbol{\theta})), \quad 31 \\ &= (t-s)A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{S}_{st} + Q_s + \beta - \mu((s-r)A(\boldsymbol{\theta}) - \boldsymbol{\theta} \cdot \mathbf{S}_{rs} + Q_r - Q_s), \\ &\quad \left( \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right) = \mathbf{S}_{st} - \mu \mathbf{S}_{rs} = 0, \quad \left( \frac{\partial \mathcal{L}}{\partial \mu} \right) = Q_s - \mu Q_r = 0, \end{aligned}$$

We get the critical point solving  $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mu) = 0$ :

$$\boldsymbol{\theta}^*(\mu) = (\nabla A)^{-1} \left( \frac{\mathbf{S}_{st} - \mu \mathbf{S}_{rs}}{(t-s) - \mu(s-r)} \right) = (\nabla A)^{-1} \left( \mathbf{m}(\mu) \right),$$

which leads to the proposed expression:

$$\mathcal{L}(\boldsymbol{\theta}^*(\mu), \mu) = \mathcal{D}(\mu) = - \left( (t-s) - \mu(s-r) \right) \mathcal{D}^*(\mathbf{m}(\mu)) + \beta + (1+\mu)Q_s - \mu Q_r,$$

where  $\mathcal{D}^*(x) = x \cdot (\nabla A)^{-1}(x) - A((\nabla A)^{-1}(x))$ . We have the constraint:

$$\nabla A(\boldsymbol{\theta}) = \frac{\mathbf{S}_{st} - \mu \mathbf{S}_{rs}}{(t-s) - \mu(s-r)} \in \mathcal{M}^o,$$

where  $\mathcal{M}^o$  is an open convex set ([6, 41]) including the value  $\mathbf{S}_{st}$  and thus  $\mu = 0$  is in the definition domain and  $\mu = \frac{t-s}{s-r}$  is the maximum feasible value (but it can be a smaller value, depending on the shape of  $\mathcal{M}^o$ ). Theoretical results on this dual formulation can be found in [6, 41].

## D.2 Dual for the change-in-mean-and-variance problem

We obtain the canonical form of the exponential family, for  $(\theta_1, \theta_2) \in \mathbb{R} \times \mathbb{R}^-$ :

$$c(y_{st}; [\theta_1; \theta_2]) = (t-s)A(\theta_1, \theta_2) - \theta_1 \left( \sum_{i=s+1}^t y_i \right) - \theta_2 \left( \sum_{i=s+1}^t y_i^2 \right). \quad (27)$$

where

$$A(\theta_1, \theta_2) = -\frac{\theta_1^2}{4\theta_2} + \frac{1}{2} \log \left( -\frac{1}{2\theta_2} \right).$$

Computing the dual of the one-constraint optimisation problem. The gradient is:

$$\nabla A(\theta) = \left( -\frac{\theta_1}{2\theta_2}, \frac{\theta_1^2}{4\theta_2^2} - \frac{1}{2\theta_2} \right),$$

and

$$\nabla A(\theta) = (u, v) \iff (\theta_1, \theta_2) = \left( -\frac{u}{u^2-v}, \frac{1}{2(u^2-v)} \right).$$

The minimum cost is then

$$c(y_{st}) = \frac{t-s}{2} \left( 1 + \log(v - u^2) \right) \quad \text{with} \quad u = \bar{y}_{st} \quad v = \bar{y}_{st}^2. \quad (28)$$

We obtain the dual, considering the case  $r < s$  (the case producing an efficient pruning) with constraint  $q_t^s - q_t^r \leq 0$ :

$$\begin{aligned} \mathcal{D}(\mu) &= \frac{1}{2} \left( (t-s) - \mu(s-r) \right) \left[ 1 + \log \left( \frac{S_{st}^2 - \mu S_{rs}^2}{(t-s) - \mu(s-r)} - \left( \frac{S_{st} - \mu S_{rs}}{(t-s) - \mu(s-r)} \right)^2 \right) \right] \\ &\quad + \beta + (1+\mu)Q_s - \mu Q_r. \end{aligned}$$

The dual is transformed into its decision function:

$$\mathbb{D}(x) = \frac{1}{2} \left[ 1 + \log \left( \overline{S_{st}^2} + x\Delta\overline{S_{rst}^2} - (\overline{S}_{st} + x\Delta\overline{S}_{rst})^2 \right) \right] - \left( \overline{Q}_{st} + x\Delta\overline{Q}_{rst} \right),$$

its domain is given by the non-negative  $x$  satisfying the condition  $\overline{S_{st}^2} + x\Delta\overline{S_{rst}^2} - (\overline{S}_{st} + x\Delta\overline{S}_{rst})^2 > 0$ .

Let's consider the particular limit cases. If  $V_\alpha = 0$ , then the data is constant:  $y_{s+1}, \dots, y_t = c$  and we solve directly  $q_t^s = q_t^r$  to find a solution. If  $V_\beta = 0$  we can use the formula for  $\mu_{max}$  to obtain

$$\mu_{max}(V_\beta = 0) = \frac{t-s}{s-r} \frac{V_\alpha}{V_\alpha + (\bar{y}_\alpha - \bar{y}_\beta)^2}.$$

In all cases, the result for  $\mu_{max}$  is always smaller than  $\frac{t-s}{s-r}$ ; this is why the solution  $(t-s) - \mu(s-r) = 0$  has been discarded.

## E Proofs of Section 4

### E.1 Proof of Theorem 8

For simplicity, we rename the indices as follows:  $s \rightarrow 0$ ,  $r_i \rightarrow i$  for  $i = 1, \dots, d$ , skip the  $t$  variable, and move the upper indices to a lower position. Therefore we solve:  $\min_{\theta \in \Theta} q_0(\theta)$  under constraints  $q_0(\theta) - q_i(\theta) \leq 0$  for all indices  $i \in \{1, \dots, d\}$ . We also use the vector function notation  $q = (q_1, \dots, q_d)$ .

The proof for strong duality relies on the following three lemmas.

**Lemma 2.** *We consider a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  for which the evaluation at point  $c \in \mathbb{R}^p$  can only be done through an auxiliary bijection  $g : \mathbb{R}^p \rightarrow \mathbb{R}^p$  such that there exists  $\theta_c \in \mathbb{R}^p$  with  $g(\theta_c) = c$ . Function  $f$  is convex if for any  $a, b$  in the convex domain  $\Omega$  of  $f$  there exists a parametric curve  $g_{ab}$  joining in a straight line the points  $a$  to  $b$  such that  $g_{ab}(\theta_{(1-\alpha)a+\alpha b}) = (1-\alpha)a + \alpha b$  and with notation  $h(\alpha) = \theta_{(1-\alpha)a+\alpha b}$  ( $h : (0, 1) \rightarrow \mathbb{R}^p$ ) we have:*

$$\nabla(fg_{ab}h(0)) \cdot (h'(0)) \leq fg_{ab}h(1) - fg_{ab}h(0).$$

It's a generalisation of a well-known result. Function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  is convex if for all  $a, b$  in its domain:  $\nabla f(a) \cdot (b-a) \leq f(b) - f(a)$ . When  $g_{ab}$  is the identity, we return to this result.

*Proof.* (of Lemma 2) The convexity of  $f$  is the same as the convexity of  $fg_h$  by definition of  $fg_h$ . Using the definition of convexity for  $fg_h$ , we have for all  $\alpha \in (0, 1)$ :

$$fgh(\alpha) \leq (1-\alpha)fgh(0) + \alpha fgh(1) \quad \text{or} \quad fgh(\alpha) - fgh(0) \leq \alpha(fgh(1) - fgh(0)).$$

Dividing by  $\alpha$  and taking the limit  $\alpha \rightarrow 0$ :  $\nabla(fgh(0)) \cdot (h'(0)) \leq fgh(1) - fgh(0)$ .  $\square$

**Lemma 3.** *Consider the  $d$  constraints. Let  $\theta_0$  be in  $\mathbb{R}^d$  such that  $(q_0 - q)(\theta_0) = c$  with  $c \in \mathbb{R}^p$ , then  $\theta_0 \in \{x^+u + w, x^-u + w\}$  with  $x^+, x^- \in \mathbb{R}$  and  $u, w \in \mathbb{R}^p$ . They are the two intersection points between a straight line and a level curve  $(q_0 - q_1) = c_1$ .*

*Proof.* (of Lemma 3) We have for  $i = 1, \dots, p$ :

$$(q_0 - q_i)(\boldsymbol{\theta}_0) = (r_i - s)A(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0 \cdot \bar{\mathbf{S}}_{r_i s} + Q_s - Q_{r_i} = c_i,$$

or

$$A(\boldsymbol{\theta}_0) - \boldsymbol{\theta}_0 \cdot \bar{\mathbf{S}}_{r_i s} = \frac{c_i}{r_i - s} + \frac{Q_{r_i} - Q_s}{r_i - s}.$$

Combining equations with indices  $i$  and  $i + 1$  for  $i = 1, \dots, p - 1$  to remove  $A(\boldsymbol{\theta}_0)$  we get:

$$\boldsymbol{\theta}_0 \cdot (\bar{\mathbf{S}}_{r_{i+1}s} - \bar{\mathbf{S}}_{r_i s}) = \frac{c_i}{r_i - s} - \frac{c_{i+1}}{r_{i+1} - s} + \frac{Q_{r_i} - Q_s}{r_i - s} - \frac{Q_{r_{i+1}} - Q_s}{r_{i+1} - s}.$$

With the assumption that the points  $\bar{\mathbf{S}}_{r_i s}$  are in general position (no redundant equation or unsolvable system), we get  $\boldsymbol{\theta}_0(x) = xu + w$  for potential solutions where  $x \in \mathbb{R}$  is a parameter and  $u, w \in \mathbb{R}^p$  are fixed values. As we considered that a solution exists and as  $A$  is convex, this straight line intersects in two points the level curve  $(q_0 - q_1) = c_1$  (counted with their multiplicity).  $\square$

**Remark 3.** Using relation  $(\bar{\mathbf{S}}_{r_{i+1}s} - \bar{\mathbf{S}}_{r_i s}) \cdot u = 0$ , we easily get that for all indices  $i$  and a fixed intersection value  $\boldsymbol{\theta}_0$  we have expressions  $(\nabla A(\boldsymbol{\theta}_0) - \bar{\mathbf{S}}_{r_i s}) \cdot u$  is a constant.

**Lemma 4.** The orthogonal projection of  $\mathcal{O}$  on its last  $d$  variables  $(\mathcal{O} \cdot (\emptyset, \mathbb{R}^d))$  is a convex object.

*Proof.* (of Lemma 4) Let  $\boldsymbol{\theta}_a$  and  $\boldsymbol{\theta}_b$  be in  $\mathbb{R}^d$  such that  $(q_0 - q)(\boldsymbol{\theta}_a) = a$  and  $(q_0 - q)(\boldsymbol{\theta}_b) = b$  with  $a, b \in \mathbb{R}^p$ . We need to prove that for all  $\alpha \in (0, 1)$ , there exists  $\boldsymbol{\theta}_\alpha$  such that  $(q_0 - q)(\boldsymbol{\theta}_\alpha) = (1 - \alpha)a + \alpha b$ . Using Lemma 3 there are only two points for each  $\alpha$  (counting multiplicity) in  $\mathcal{O} \cdot (\emptyset, \mathbb{R}^d)$  solving the equations:  $x^-(\alpha)u + \alpha v + w$  and  $x^+(\alpha)u + \alpha v + w$ . They intersect in  $\alpha = 0$  and  $\alpha = 1$  the level curve  $(q_0 - q_1) = a$  and  $(q_0 - q_1) = b$ , respectively. We must show that the straight line also intersects  $(q_0 - q_1) = (1 - \alpha)a + \alpha b$ . We increase by one the size of the problem, considering  $\alpha$  to be one of the variables. the plane  $xu + \alpha v + w$  intersect  $f(\boldsymbol{\theta}, \alpha) = (q_0 - q_1)(\boldsymbol{\theta}) - \alpha(b - a) - a$  in  $\alpha = 0$  and  $\alpha = 1$ .  $f(\boldsymbol{\theta}, \alpha) \leq 0$  or  $f(\boldsymbol{\theta}, \alpha) \geq 0$  is convex and bounded. Thus, the intersection also occurs for all  $\alpha$  in  $(0, 1)$ .  $\square$

**Proof of the Theorem:** We need to prove the convexity of the epigraph, that is, for the lowest values of  $q_0$  for points on the projection  $\mathcal{O} \cdot (\emptyset, \mathbb{R}^d)$ . This is not equivalent to the convexity of  $q_0$  as the evaluation is done here on the initial parameter's functions (the constraint values) instead of  $\boldsymbol{\theta}$ . Let  $\boldsymbol{\theta}_a$  and  $\boldsymbol{\theta}_b$  be in  $\mathbb{R}^d$  and  $a, b$  be in  $\mathcal{O} \cdot (\emptyset, \mathbb{R}^d)$  such that  $(q_0 - q)(\boldsymbol{\theta}_a) = a$  and  $(q_0 - q)(\boldsymbol{\theta}_b) = b$ .

We have:

$$h^- : \mathbb{R} \rightarrow \mathbb{R}^p, \alpha \mapsto x^-(\alpha)u + \alpha v + w, \quad h^+ : \mathbb{R} \rightarrow \mathbb{R}^p, \alpha \mapsto x^+(\alpha)u + \alpha v + w.$$

such that  $\boldsymbol{\theta}_a \in \{h^-(0), h^+(0)\}$ ,  $\boldsymbol{\theta}_b \in \{h^-(1), h^+(1)\}$  and  $g(h^-(\alpha)) = g(h^+(\alpha)) = (1 - \alpha)a + \alpha b$ . Function  $h^-$  corresponds to the solution for  $q_0$  returning the smallest value:

$$q_0(g_{ab}(h^-(\alpha))) \leq q_0(g_{ab}(h^+(\alpha))).$$

$h^-$  and  $h^+$  can be discontinuous, however  $h^-$  is differentiable in 0 (no jump in minimum value at 0). Applying Lemma 2 with  $h = h^-$  (and  $x = x^-$ ) we need to prove:

$$\nabla(fg_{ab}h(0)) \cdot (h'(0)) \leq fg_{ab}h(1) - fg_{ab}h(0).$$

That is:

$$\left( \nabla A(x(0)u + w) - \bar{\mathbf{S}}_{st} \right) \cdot \left( \frac{dx}{d\alpha}(0)u + v \right) \leq \frac{q_0(\boldsymbol{\theta}_b) - q_0(\boldsymbol{\theta}_a)}{t-s} \quad (29)$$

Notice that we decide to choose  $u$  such that  $(\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot u \geq 0$  (as  $u$  can be replaced by  $-u$  if we have the wrong sign). We can also derive two interesting equalities. We consider the intersection between the line and the level curve (see Lemma 3), which also intersects the curves  $q_0 = q_0(\boldsymbol{\theta}_a)$  and  $q_0 = q_0(\boldsymbol{\theta}_b)$ . For  $\alpha = 0$  and  $\alpha = 1$  we get four relations:

$$\begin{cases} A(\boldsymbol{\theta}_a) - \boldsymbol{\theta}_a \cdot \bar{\mathbf{S}}_{st} + \frac{Q_s + \beta}{t-s} = \frac{q_0(\boldsymbol{\theta}_a)}{t-s}, \\ A(\boldsymbol{\theta}_b) - \boldsymbol{\theta}_b \cdot \bar{\mathbf{S}}_{st} + \frac{Q_s + \beta}{t-s} = \frac{q_0(\boldsymbol{\theta}_b)}{t-s}, \\ A(\boldsymbol{\theta}_a) - \boldsymbol{\theta}_a \cdot \bar{\mathbf{S}}_{s_1s} - \frac{Q_{s_1} - Q_s}{s_1 - s} = \frac{(q_0 - q_1)(\boldsymbol{\theta}_a)}{s_1 - s}, \\ A(\boldsymbol{\theta}_b) - \boldsymbol{\theta}_b \cdot \bar{\mathbf{S}}_{s_1s} - \frac{Q_{s_1} - Q_s}{s_1 - s} = \frac{(q_0 - q_1)(\boldsymbol{\theta}_b)}{s_1 - s}, \end{cases}$$

leading to relation:

$$(\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot (\boldsymbol{\theta}_b - \boldsymbol{\theta}_a) = \frac{q_0(\boldsymbol{\theta}_b) - q_0(\boldsymbol{\theta}_a)}{s_1 - s} - \frac{(q_0 - q_1)(\boldsymbol{\theta}_b) - (q_0 - q_1)(\boldsymbol{\theta}_a)}{s_1 - s}. \quad (30)$$

The second relation is the dynamic equation for points  $x^\pm(\alpha)u + \alpha v + w$ . We use relations:

$$(q_0 - q_1)(x(0)u + w) = (q_0 - q_1)(\boldsymbol{\theta}_a)$$

and

$$(q_0 - q_1)(x(\epsilon)u + \epsilon v + w) = (1 - \epsilon)(q_0 - q_1)(\boldsymbol{\theta}_a) + \epsilon(q_0 - q_1)(\boldsymbol{\theta}_b)$$

We derive, taking the difference with  $\epsilon \rightarrow 0$ :

$$\left( \nabla A(x(0)u + w) - \bar{\mathbf{S}}_{s_1s} \right) \cdot \left( \frac{dx}{d\alpha}(0)u + v \right) = \frac{(q_0 - q_1)(\boldsymbol{\theta}_b) - (q_0 - q_1)(\boldsymbol{\theta}_a)}{s_1 - s}. \quad (31)$$

Using (30) and (31) we can now reformulate (29) as:

$$(\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot \left( \frac{dx}{d\alpha}(0)u + v \right) \leq (\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot (\boldsymbol{\theta}_b - \boldsymbol{\theta}_a),$$

or

$$\frac{dx}{d\alpha}(0)(\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot u \leq (x(1) - x(0))(\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot u.$$

As we set  $(\bar{\mathbf{S}}_{s_1s} - \bar{\mathbf{S}}_{st}) \cdot u \geq 0$ , this gives us the relation:

$$\nabla x(0) \cdot (1 - 0) \leq x(1) - x(0),$$

to be proven. This is true if  $\alpha \mapsto x(\alpha)$  is convex. Generalising Equation (31) for all  $\alpha \in (0, 1)$  we get:

$$\left( \nabla A(x(\alpha)u + \alpha v + w) - \bar{\mathbf{S}}_{s_1s} \right) \cdot \left( \frac{dx}{d\alpha}(\alpha)u + v \right) = \frac{(q_0 - q_1)(\boldsymbol{\theta}_b) - (q_0 - q_1)(\boldsymbol{\theta}_a)}{s_1 - s},$$

and differentiating in  $\alpha$ :

$$h'(\alpha) \cdot \left( \nabla A^2(h(\alpha)) \right) \cdot h'(\alpha) = -\frac{d^2x}{d\alpha^2}(\alpha) \left( \nabla A(h(\alpha)) - \bar{\mathbf{S}}_{s_1 s} \right) \cdot u.$$

However,  $\left( \nabla A(h(\alpha)) - \bar{\mathbf{S}}_{s_1 s} \right) \cdot u < 0$ , knowing the choice of  $u$  we made and Remark 3. We exclude the nullity case as it is possible only for a point on the boundary of  $\mathcal{O} \cdot (\emptyset, \mathbb{R}^d)$ . We can restrict points  $a$  and  $b$  to be in the interior and still prove the same result. Using the convexity of  $A$ , we obtain  $\frac{d^2x}{d\alpha^2}(\alpha) > 0$  and we find that  $\alpha \mapsto x(\alpha)$  is convex.

## E.2 Proof of Proposition 10

The gradient of  $\mathcal{D}^*$  is equal to  $(\nabla A)^{-1}$  which leads to relation:

$$\nabla \mathcal{D}^* \left( \bar{\mathbf{S}}_{st} + \sum_{r \neq s} x_r \Delta \bar{\mathbf{S}}_{rst} \right) = (\Delta \bar{\mathbf{S}}_{\bullet st})^T \left( (\nabla A)^{-1} \left( \bar{\mathbf{S}}_{st} + \sum_{r \neq s} x_r \Delta \bar{\mathbf{S}}_{rst} \right) \right).$$

We introduce notation  $\mathbf{y} = (\nabla A)^{-1} \left( \bar{\mathbf{S}}_{st} + \sum_{r \neq s} x_r \Delta \bar{\mathbf{S}}_{rst} \right)$  which leads to the first system of equations:  $(\Delta \bar{\mathbf{S}}_{\bullet st}) \mathbf{x} = \nabla A(\mathbf{y}) - \bar{\mathbf{S}}_{st}$ . With such notation, the critical point (the maximum, function  $\mathbb{D}$  being concave), is the solution of  $\nabla \mathbb{D}(\mathbf{x}) = 0$  and gives:

$$-(\Delta \bar{\mathbf{S}}_{\bullet st})^T \mathbf{y} - \Delta \bar{Q}_{\bullet st} = 0.$$

## F Quadratic cost function

### F.1 The General case

We consider the following functions of the type:

$$q_t^s(\theta_1, \theta_2) = A_{st}\theta_1^2 + 2B_{st}\theta_1\theta_2 + C_{st}\theta_2^2 + 2D_{st}\theta_1 + 2E_{st}\theta_2 + F_{st}.$$

A natural restriction is to study the case when all the cost functions have a unique finite minimum value and a unique argument for the minimum. Thus, we can define the cost of a segment and associate with it a non-ambiguous parameter value. For any function  $q_t^k$  it is equivalent to condition  $A_{kt}C_{kt} - B_{kt}^2 > 0$  with  $A_{kt} > 0$  (such functions are then strictly convex). With the Lagrangian with one constraint given by inequality  $q_t^s - q_t^r \leq 0$  (always considering that  $r < s$ , unless otherwise specified) we obtain the coefficients:

$$A(\mu) = A_{st} + \mu(A_{st} - A_{rt}), B(\mu) = B_{st} + \mu(B_{st} - B_{rt}) \dots$$

After computation, we have the following dual function given by the expression

$$\mathcal{D}(\mu) = \frac{2B(\mu)D(\mu)E(\mu) - A(\mu)E^2(\mu) - C(\mu)D^2(\mu)}{A(\mu)C(\mu) - B^2(\mu)} + F(\mu). \quad (32)$$

Suppose there is no possible reduction of the rational function defining the dual. When the underlying process generating the time-series and therefore the functions  $q_t^k$  is continuous, the possibility of a reduction is certainly an event of measure zero. Thus, we are looking for the first positive value  $\mu_{max}$  such that  $A(\mu_{max})C(\mu_{max}) - B^2(\mu_{max}) = 0$ . This leads to the following result.

**Proposition 13.** *The maximal value  $\mu_{max}$  of the dual function given by Equation (32), if no possible reduction of the fraction, is the smallest positive root of:*

$$A(\mu)C(\mu) - B^2(\mu) = (\omega_1^2 - 2\Delta + \omega_2^2)\mu^2 - 2(\Delta - \omega_1^2)\mu + \omega_1^2 = 0,$$

with

$$\omega_1^2 = A_{st}C_{st} - B_{st}^2, \quad 2\Delta = A_{st}C_{rt} + A_{rt}C_{st} - 2B_{st}B_{rt} \text{ and } \omega_2^2 = A_{rt}C_{rt} - B_{rt}^2.$$

We have  $\omega_1^2 > 0$  and  $\omega_2^2 > 0$  and we consider that  $\omega_1^2 \neq \omega_2^2$ . The discriminant is also positive:  $4(\Delta^2 - \omega_1^2\omega_2^2) > 0$ . The maximal value is then given by:

$$\mu_{max} = \begin{cases} \frac{\Delta - \omega_1^2 - \sqrt{\Delta^2 - \omega_1^2\omega_2^2}}{\omega_1^2 - 2\Delta + \omega_2^2}, & \text{if } 2\Delta > \omega_1^2 + \omega_2^2, \\ & \text{or } 2\Delta < \omega_1^2 + \omega_2^2, \omega_1^2 < \Delta, \\ \frac{\omega_1^2}{\omega_2^2 - \omega_1^2}, & \text{if } 2\Delta = \omega_1^2 + \omega_2^2, \omega_1^2 < \Delta, \\ +\infty, & \text{if } 2\Delta \leq \omega_1^2 + \omega_2^2, \omega_1^2 > \Delta. \end{cases} \quad (33)$$

*Proof.* First, we propose that the discriminant is positive. We can write:

$$\Delta^2 - \omega_1^2\omega_2^2 = A_{st}C_{st} \left( A_{st}C_{st}(d_C - d_A)^2 + 4B_{st}^2d_Ad_C \right),$$

with

$$d_A = \frac{B_{st}}{B_{rt}} - \frac{A_{st}}{A_{rt}} \quad \text{and} \quad d_C = \frac{B_{st}}{B_{rt}} - \frac{C_{st}}{C_{rt}}.$$

As we have  $A_{st}C_{st} > B_{st}^2$  and  $(d_C - d_A)^2 > -4d_Ad_C$ , multiplying these two inequalities give us a positive determinant. Second, we determine the smallest positive root by studying the sign of the product of the roots, which is equal to  $\omega_1^2(\omega_1^2 - 2\Delta + \omega_2^2)^{-1}$ . With  $2\Delta > \omega_1^2 + \omega_2^2$ , the product is negative and the biggest root is the only one positive. If on the contrary  $2\Delta < \omega_1^2 + \omega_2^2$ , the roots have the same sign and the sign of  $\Delta - \omega_1^2$  determines whether both roots are positive (in that case the smallest root is the solution) or both negative (in that case  $\mu_{max}$  is infinite). The degenerated case with no quadratic term is obvious.  $\square$

It is also possible to consider the quadratic form in dimension  $p$  more than 2. However, determining the boundary for  $\mu$  is challenging to determine explicitly: this is the first mu value such that  $\det(A - \mu B) = 0$  with  $A$  and  $B$  two  $p \times p$  matrices.

## F.2 Changes in simple regression

A direct application of the previous result is the problem of detecting a change point in simple regression. In this setting, we gather two-dimensional  $(x_t, y_t)$  points over time. At a change location, the linear bound linking  $x_t$  and  $y_t$  changes. To be specific, this corresponds to a model with data points  $(x_t, y_t)$  generated by a succession of simple linear models (choosing  $x_t$  and then getting the response  $y_t$  through the regression)

$$y_t = a_i x_t + b_i + \epsilon_t, \quad t = \tau_i + 1, \dots, \tau_{i+1}, \quad i = 0, \dots, K,$$

with  $t \mapsto a_t$  and  $t \mapsto b_t$  piecewise constant time series and  $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$  identically and independently distributed. The vector  $(\tau_1, \dots, \tau_K)$  of strictly increasing integers subdivides the time-series into  $K+1$  consecutive segments with natural notations  $\tau_0 = 0$  and  $\tau_{K+1} = n$ . Using the maximum likelihood approach, we see that this leads to considering the cost function  $q_t^s(\theta_1, \theta_2) = Q_s + c(y_{st}; [\theta_1; \theta_2]) + \beta$  with

$$\begin{aligned} c(y_{st}; [\theta_1; \theta_2]) &= \sum_{j=s+1}^t (y_j - (\theta_1 x_j + \theta_2))^2, \\ &= (t-s) \left( \overline{x_{st}^2} \theta_1^2 + 2\overline{x_{st}} \theta_1 \theta_2 + \theta_2^2 - 2\overline{(xy)_{st}} \theta_1 - 2\overline{y_{st}} \theta_2 + \overline{y_{st}^2} \right). \end{aligned}$$

Identifying the coefficients term by term, we introduce the notations:

$$q_t^s(\theta_1, \theta_2) = A_{st} \theta_1^2 + 2B_{st} \theta_1 \theta_2 + C_{st} \theta_2^2 + 2D_{st} \theta_1 + 2E_{st} \theta_2 + F_{st}.$$

Writing as previously  $A(\mu) = A_{st} + \mu(A_{st} - A_{rt}), \dots$  we obtain the dual function in expression (32).

## G Examples of dual change point functions

We consider at time  $t$  the function  $q_t^s$  constrained by the function  $q_t^r$ . The shape of the dual function  $\mathcal{D}$  for some distributions of the exponential family in the one-constraint case can be easily expressed as:

$$\mathcal{D}(\mu) = -(t-s) \left( l(\mu) \mathcal{D}^*(\mathbf{m}(\mu)) + (-1)^{[r < s]} \mu \frac{Q_s - Q_r}{s-r} \right) + Q_s + \beta, \quad (34)$$

where

$$\mathbf{m}(\mu) = \frac{\overline{\mathbf{S}}_{st} + \mu(-1)^{[r < s]} \overline{\mathbf{S}}_{rs}}{1 + \mu(-1)^{[r < s]}}, \quad l(\mu) = 1 + \mu(-1)^{[r < s]}.$$

Equation (34) is the re-scaled version of the dual, replacing  $\frac{|s-r|}{t-s} \mu$  by  $\mu$ . Only (convex) functions  $\mathcal{D}^*$  as well as their associated domain in  $\mu$  parameter for  $\mathcal{D}$  are distribution-dependent; we give their form in Table 3. Function  $\mathcal{D}^*$  can be expressed through the log-partition function  $A$  as  $\mathcal{D}^*(x) = x(\nabla A)^{-1}(x) - A((\nabla A)^{-1}(x))$ . The right bound of the support for  $\mathcal{D}$  is the largest value  $\mu_{max}$  such that the dual is defined on  $(0, \mu_{max})$ . For simplicity, we use notation  $\mathbf{m}(\mu) = \frac{\sigma_1 + (-1)^{[r < s]} \sigma_2 \mu}{1 + (-1)^{[r < s]} \mu}$  in Table 3. Notice that we focus here on the case  $r < s$ . When  $r > s$  there is no limit for parameter  $\mu$ :  $\mu_{max} = +\infty$ .

It is interesting to notice that  $l(\mu_{max}) \mathcal{D}^*(\mathbf{m}(\mu_{max}))$  is equal to zero except for the following cases: Gauss, exponential (if  $\mu_{max} = \frac{\sigma_1}{\sigma_2}$ ), Poisson (if  $\mu_{max} = 1$ ). Details for computing  $\mathcal{D}^*$  are given in Table 4. More details on the properties of  $\mathcal{D}^*$  can be found in Chapter 3 of [41].

Table 3: Distribution, their dual function and maximum dual parameter in case  $r < s$

Distribution	Function $x \mapsto \mathcal{D}^*(x)$	Maximal value $\mu_{max}$
Gauss	$\frac{1}{2}x^2$	1
Exponential	$-\log x - 1$	$\min(1, \frac{\sigma_1}{\sigma_2})$
Poisson	$x(\log x - 1)$	$\min(1, \frac{\sigma_1}{\sigma_2})$
Geometric	$(x-1)\log(x-1) - x \log x$	$\min(1, \frac{\sigma_1-1}{\sigma_2-1})$
Bernoulli/Binomial	$x \log x + (1-x) \log(1-x)$	$\min(\frac{\sigma_1}{\sigma_2}, \frac{1-\sigma_1}{1-\sigma_2})$
Negative Binomial	$x \log x - (1+x) \log(1+x)$	$\min(1, \frac{\sigma_1}{\sigma_2})$

In the Gaussian case, data is standardised: divided by the estimated standard deviation. In binomial and negative binomial cases, data is divided by the estimated number of trials and successes (respectively), and the logarithmic values can be computed.

Table 4: Functions  $A$  and  $(\nabla A)^{-1}$  used to compute  $\mathcal{D}^*(x) = x(\nabla A)^{-1}(x) - A((\nabla A)^{-1}(x))$  for some standard distributions in exponential family

Distribution	$y \mapsto A(y)$	$x \mapsto (\nabla A)^{-1}(x)$	$x \in \Omega$
Gauss	$\frac{1}{2}y^2$	$x$	$\Omega = \mathbb{R}$
Exponential	$-\log(-y)$	$-\frac{1}{x}$	$\Omega = (0, +\infty)$
Poisson	$\exp(y)$	$\log(x)$	$\Omega = (0, +\infty)$
Geometric	$-\log(e^{-y} - 1)$	$\log\left(\frac{x-1}{x}\right)$	$\Omega = (1, +\infty)$
Bernoulli/Binomial	$\log(1 + e^y)$	$\log\left(\frac{x}{1-x}\right)$	$\Omega = (0, 1)$
Negative Binomial	$-\log(1 - e^y)$	$\log\left(\frac{x}{1+x}\right)$	$\Omega = (0, +\infty)$

**Remark 4.** From the table 3 we can easily write down the dual functions in a multivariate setting using relations (16). In that case, we sum over all dimensions the dual obtained for each univariate time-series (except for constant  $Q_s + \beta$ ).

For the sake of completeness, we provide details about the limiting cases, defined by the limit values for  $\bar{\mathbf{S}}_{st}$  or  $\bar{\mathbf{S}}_{rs}$ , with  $\bar{\mathbf{S}}_{st} \neq \bar{\mathbf{S}}_{rs}$ , on the boundary of  $\Omega$ . These limit values are denoted  $\partial\Omega$  (equal to 0 or 1 with the given distributions in Table 4). If  $\bar{\mathbf{S}}_{st} = \partial\Omega$ , then  $\mu_{max} = 0$  for all distributions (except for Gauss). If  $\bar{\mathbf{S}}_{rs} = \partial\Omega$ , then  $\mu_{max} = 1$  except in the case Bernoulli/Binomial where  $\mu_{max} = \bar{\mathbf{S}}_{st}$  if  $\bar{\mathbf{S}}_{rs} = 1$  and  $\mu_{max} = 1 - \bar{\mathbf{S}}_{st}$  if  $\bar{\mathbf{S}}_{rs} = 0$ .

**Remark 5.** If  $\bar{\mathbf{S}}_{st} = \bar{\mathbf{S}}_{rs} = \mathbf{S} \in \Omega \cup \partial\Omega$ , we get  $\mathcal{D}^*(\mathbf{m}(\mu)) = \mathcal{D}^*(\mathbf{S})$ . Therefore, the dual  $\mathcal{D}$  is a linear function with maximum in  $\mu = 0$  or  $\mu = 1$  (the maximal value  $\mu_{max}$  in that case, see Table 3), depending on the sign of the slope.

## H Additional simulations

### Competing interests

No competing interest is declared.

### References

- [1] Samaneh Aminikhangahi and Diane J Cook. A survey of methods for time series change point detection. *Knowledge and information systems*, 51(2):339–367, 2017.
- [2] Ivan E Auger and Charles E Lawrence. Algorithms for the optimal identification of segment neighborhoods. *Bulletin of mathematical biology*, 51(1):39–54, 1989.
- [3] Jushan Bai and Pierre Perron. Computation and analysis of multiple structural change models. *Journal of applied econometrics*, 18(1):1–22, 2003.
- [4] Sayantan Banerjee and Kousik Guhathakurta. Change-point analysis in financial networks. *Stat*, 9(1):e269, 2020.
- [5] Emily Brodsky and Boris S Darkhovsky. *Nonparametric methods in change point problems*, volume 243. Springer Science & Business Media, 2013.
- [6] Lawrence D Brown. Fundamentals of statistical exponential families: with applications in statistical decision theory. In *Fundamentals of statistical exponential families: with applications in statistical decision theory*. Ims, 1986.
- [7] Kuo-Ching Chang, Chui-Liang Chiang, and Chung-Bow Lee. A fast two-stage dynamic programming algorithm for change-points model with application in speech signal. In *2010 Fourth international conference on genetic and evolutionary computing*, pages 366–370. IEEE, 2010.
- [8] Jie Chen, Arjun K Gupta, and AK Gupta. *Parametric statistical change point analysis*, volume 192. Springer, 2000.
- [9] Alice Cleynen and Émilie Lebarbier. Model selection for the segmentation of multiparameter exponential family distributions. *Electronic Journal of Statistics*, 11(1):800 – 842, 2017.
- [10] P. M. Rodríguez Cruz, J. Cossins, D. Beeson, and A. Vincent. The neuromuscular junction in health and disease: Molecular mechanisms governing synaptic formation and homeostasis. *Frontiers in Molecular Neuroscience*, 13, 12 2020.
- [11] Miklós Csörgő and Lajos Horváth. Limit theorems in change-point analysis. (*No Title*), 1997.
- [12] Paul Fearnhead, Robert Maidstone, and Adam Letchford. Detecting changes in slope with an  $l_0$  penalty. *Journal of Computational and Graphical Statistics*, 28(2):265–275, 2019.

- [13] William Gu, Jaesik Choi, Ming Gu, Horst Simon, and Kesheng Wu. Fast change point detection for electricity market analysis. In *2013 IEEE International Conference on Big Data*, pages 50–57. IEEE, 2013.
- [14] Zaid Harchaoui, Félicien Vallet, Alexandre Lung-Yut-Fong, and Olivier Cappé. A regularized kernel-based approach to unsupervised audio segmentation. In *2009 IEEE international conference on acoustics, speech and signal processing*, pages 1665–1668. IEEE, 2009.
- [15] Toby Dylan Hocking, Guillem Rigaill, Paul Fearnhead, and Guillaume Bourque. Constrained dynamic programming and supervised penalty learning algorithms for peak detection in genomic data. *Journal of Machine Learning Research*, 21(87):1–40, 2020.
- [16] Brad Jackson, Jeffrey D Scargle, David Barnes, Sundararajan Arabhi, Alina Alt, Peter Gioumousis, Elyus Gwin, Paungkaew Sangtrakulcharoen, Linda Tan, and Tun Tao Tsai. An algorithm for optimal partitioning of data on an interval. *IEEE Signal Processing Letters*, 12(2):105–108, 2005.
- [17] Venkata Jandhyala, Stergios Fotopoulos, Ian MacNeill, and Pengyu Liu. Inference for single and multiple change-points in time series. *Journal of Time Series Analysis*, 34(4):423–446, 2013.
- [18] Sean W Jewell, Toby Dylan Hocking, Paul Fearnhead, and Daniela M Witten. Fast nonconvex deconvolution of calcium imaging data. *Biostatistics*, 21(4):709–726, 2020.
- [19] Shengji Jia and Lei Shi. Efficient change-points detection for genomic sequences via cumulative segmented regression. *Bioinformatics*, 38(2):311–317, 2022.
- [20] Rebecca Killick, Paul Fearnhead, and Idris A Eckley. Optimal detection of change-points with a linear computational cost. *Journal of the American Statistical Association*, 107(500):1590–1598, 2012.
- [21] E. Krejci, P. Konstantin, O. Lenina, V. Bernard, T. Germain, C. Truong, L. Nureullin, G. Sibgatullina, D. Samigullin, and K. Ohno. An  $\alpha 7$  nicotinic and gabab receptor-mediated pathway controls acetylcholine release in the tripartite neuromuscular junction. *The Journal of Physiology*, 2024.
- [22] Arnaud Liehrmann, Etienne Delannoy, Alexandra Launay-Avon, Elodie Gilbault, Olivier Loudet, Benoît Castanet, and Guillem Rigaill. Diffsegr: an rna-seq data driven method for differential expression analysis using changepoint detection. *NAR Genomics and Bioinformatics*, 5(4):lqad098, 2023.
- [23] Robert Maidstone, Toby Hocking, Guillem Rigaill, and Paul Fearnhead. On optimal multiple changepoint algorithms for large data. *Statistics and computing*, 27:519–533, 2017.
- [24] ES Page. A test for a change in a parameter occurring at an unknown point. *Biometrika*, 42(3/4):523–527, 1955.

- [25] Ewan S Page. Continuous inspection schemes. *Biometrika*, 41(1/2):100–115, 1954.
- [26] Liudmila Pishchagina, Guillem Rigaill, and Vincent Runge. Geometric-based pruning rules for change point detection in multiple independent time series. *Computo*, 2024.
- [27] Liudmila Pishchagina, Gaetano Romano, Paul Fearnhead, Vincent Runge, and Guillem Rigaill. Online multivariate changepoint detection: Leveraging links with computational geometry. *arXiv preprint arXiv:2311.01174*, 2023.
- [28] Jaxk Reeves, Jien Chen, Xiaolan L Wang, Robert Lund, and Qi Qi Lu. A review and comparison of changepoint detection techniques for climate data. *Journal of applied meteorology and climatology*, 46(6):900–915, 2007.
- [29] Guillem Rigaill. A pruned dynamic programming algorithm to recover the best segmentations with 1 to  $k_{\max}$  change-points. *Journal de la société française de statistique*, 156(4):180–205, 2015.
- [30] Gaetano Romano, Idris A Eckley, Paul Fearnhead, and Guillem Rigaill. Fast online changepoint detection via functional pruning cusum statistics. *Journal of Machine Learning Research*, 24(81):1–36, 2023.
- [31] Gaetano Romano, Guillem Rigaill, Vincent Runge, and Paul Fearnhead. Detecting abrupt changes in the presence of local fluctuations and autocorrelated noise. *Journal of the American Statistical Association*, 117(540):2147–2162, 2022.
- [32] Vincent Runge. Is a finite intersection of balls covered by a finite union of balls in euclidean spaces? *Journal of Optimization Theory and Applications*, 187(2):431–447, 2020.
- [33] Vincent Runge, Toby Dylan Hocking, Gaetano Romano, Fatemeh Afghah, Paul Fearnhead, and Guillem Rigaill. gfpop: An r package for univariate graph-constrained change-point detection. *Journal of Statistical Software*, 106(6):1–39, 2023.
- [34] Andrew Jhon Scott and Martin Knott. A cluster analysis method for grouping means in the analysis of variance. *Biometrics*, pages 507–512, 1974.
- [35] Ashish Sen and Muni S Srivastava. On tests for detecting change in mean. *The Annals of statistics*, pages 98–108, 1975.
- [36] IMS Lecture Notes-Monograph Series. Change-point problems ims lecture notes-monograph series (volume 23, 1994) a nonparametric test for homogeneity: Applications to parameter estimation by k. ghoudi and d. mcdonald. *Change-point Problems*, 23:149, 1994.
- [37] Xueheng Shi, Claudio Beaulieu, Rebecca Killick, and Robert Lund. Changepoint detection: An analysis of the central england temperature series. *Journal of Climate*, 35(19):6329–6342, 2022.
- [38] Charles Truong, Laurent Oudre, and Nicolas Vayatis. Selective review of offline change point detection methods. *Signal Processing*, 167:107299, 2020.

- [39] Nicolas Verzelen, Magalie Fromont, Matthieu Lerasle, and Patricia Reynaud-Bouret. Optimal change-point detection and localization., 2020.
- [40] L. Yu. Vostrikova. Detecting "disorder" in multidimensional random processes. *Sov. Math., Dokl.*, 24:55–59, 1981.
- [41] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [42] Yi-Ching Yao. Estimating the number of change-points via schwarz'criterion. *Statistics & Probability Letters*, 6(3):181–189, 1988.
- [43] Nancy Zhang and David Siegmund. A modified bayes information criterion with applications to the analysis of comparative genomic hybridization data. *Biometrics*, 63:22–32, 04 2007.
- [44] Doudou Zhou and Hao Chen. Asymptotic distribution-free change-point detection for modern data based on a new ranking scheme. *arXiv preprint arXiv:2206.03038*, 2022.