

ORIGINAL ARTICLE

High-dimensional, multiscale online changepoint detection

Yudong Chen¹ | Tengyao Wang^{2,3}  | Richard J. Samworth¹¹University of Cambridge, Cambridge, Cambridgeshire, UK²London School of Economics and Political Science, London, UK³University College London, London, UK**Correspondence**

Tengyao Wang, Department of Statistics, London School of Economics and Political Science, Columbia House, 69 Aldwych, London, WC2B 4RR, UK.
Email: t.wang59@lse.ac.uk

Funding information

EPSRC, Grant/Award Number: EP/T02772X/1, EP/P031447/1 and EP/N031938/1

Abstract

We introduce a new method for high-dimensional, online changepoint detection in settings where a p -variate Gaussian data stream may undergo a change in mean. The procedure works by performing likelihood ratio tests against simple alternatives of different scales in each coordinate, and then aggregating test statistics across scales and coordinates. The algorithm is online in the sense that both its storage requirements and worst-case computational complexity per new observation are independent of the number of previous observations; in practice, it may even be significantly faster than this. We prove that the patience, or average run length under the null, of our procedure is at least at the desired nominal level, and provide guarantees on its response delay under the alternative that depend on the sparsity of the vector of mean change. Simulations confirm the practical effectiveness of our proposal, which is implemented in the R package `ocd`, and we also demonstrate its utility on a seismology data set.

KEYWORDS

average run length, detection delay, high-dimensional changepoint detection, online algorithm, sequential method

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2022 The Authors. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* published by John Wiley & Sons Ltd on behalf of Royal Statistical Society

1 | INTRODUCTION

Modern technology has not only allowed the collection of data sets of unprecedented size, but has also facilitated the real-time monitoring of many types of evolving processes of interest. Wearable health devices, astronomical survey telescopes, self-driving cars and transport network load-tracking systems are just a few examples of new technologies that collect large quantities of streaming data, and that provide new challenges and opportunities for statisticians.

Very often, a key feature of interest in the monitoring of a data stream is a *change point*; that is, a moment in time at which the data-generating mechanism undergoes a change. Such times often represent events of interest, for example, a change in heart function, and moreover, the accurate identification of change points often facilitates the decomposition of a data stream into stationary segments.

Historically, it has tended to be univariate time series that have been monitored and studied, within the well-established field of statistical process control (e.g. Barnard, 1959; Duncan, 1952; Fearnhead & Liu, 2007; Oakland, 2007; Page, 1954; Tartakovsky et al., 2014). These days, however, it is frequently the case that many data processes are measured simultaneously. In the context of change point detection, this introduces the new challenge of borrowing strength across the different component series in an attempt to detect much smaller changes than would be possible through the observation of any individual series alone.

The field of change point detection and estimation also has a long history (e.g. Page, 1955), but has been undergoing a marked renaissance in recent years; entry points to the field include Csörgő and Horváth (1997) and Horváth and Rice (2014). However, the vast majority of this ever-growing literature has focused on the offline change point problem, where, after the entire data stream is observed, the statistician is asked to identify any change points retrospectively. For univariate, offline change point estimation, state-of-the-art methods include the pruned exact linear time method (PELT) (Killick et al., 2012), narrowest-over-threshold (NOT) (Baranowski et al., 2019), simultaneous multiscale change point estimator (SMUCE) (Frick et al., 2014) and ℓ_0 -penalisation (Wang et al., 2018), while work on multivariate and high-dimensional offline change points includes the double CUSUM method of Cho (2016), the `inspect` algorithm of Wang and Samworth (2018), as well as Enikeeva and Harchaoui (2019), Liu et al. (2021) and Padilla et al. (2019).

Despite this rich literature on offline change point problems, it is the online version of the problem that is arguably the more important for many applications: one would like to be able to detect a change as soon as possible after it has occurred. Of course, one option here is to apply an offline method after seeing every new observation (or batch of observations). However, this is unlikely to be a successful strategy: not only is there a difficult and highly dependent multiple testing issue to handle when using the method repeatedly on data sets of increasing size (see also Chu et al. (1996) for further discussion of this point), but moreover, the storage and running time costs may frequently be prohibitive.

In this work, we are interested in algorithms for detecting change points in high-dimensional data that are observed sequentially. In order to avoid the trap mentioned in the previous paragraph and ensure that any methods we consider can be applied to large data streams, we will focus our attention on *online algorithms*. By this, we mean that the computational complexity for processing a new observation, as well as the storage requirements, depend only on the number of bits needed to represent the new observation.¹ Importantly, they are not allowed to

¹For the purpose of this definition, we ignore the errors in rounding real numbers to machine precision. Thus, when we later work with observations having Gaussian (or other absolutely continuous) distributions, we do not distinguish between these distributions and quantised versions where the data have been rounded to machine precision.

depend on the number of previously observed data points. This turns out to be a very stringent requirement, in the sense that finding online algorithms with good statistical performance is typically extremely challenging. Online algorithms must necessarily store only compact summaries of the historical observations, so the class of all possible procedures is severely restricted.

To set the scene for our contributions, let X_1, X_2, \dots be a sequence of independent random vectors in \mathbb{R}^p . Assume that for some unknown, deterministic time $z \in \mathbb{N} \cup \{0\}$, the sequence is generated according to

$$X_1, \dots, X_z \sim \mathcal{N}_p(\mu_-, I_p) \quad \text{and} \quad X_{z+1}, X_{z+2}, \dots \sim \mathcal{N}_p(\mu_+, I_p), \quad (1)$$

for some $\mu_-, \mu_+ \in \mathbb{R}^p$. When $\mu_+ \neq \mu_-$, we say that there is a changepoint at time z . In many applications, such as in industrial quality control where the distribution of relevant properties of goods in a manufacturing process under regular conditions may be well understood, we may assume that the mean before the change is known (or at least can be estimated to high accuracy using historical data). However, the vector of change, $\theta := \mu_+ - \mu_-$, is typically unknown. Thus, for simplicity, we will work in the setting where $\mu_- = 0$ and $\mu_+ = \theta$. Let $\mathbb{P}_{z,\theta}$ denote the joint distribution of $(X_n)_{n=1}^\infty$ under (1) and $\mathbb{E}_{z,\theta}$ the expectation under this distribution. Note that when $\theta = 0$, the joint distribution of the data does not depend on z , and we therefore let $\mathbb{P}_0 = \mathbb{P}_{z,0}$ denote this joint distribution (with corresponding expectation \mathbb{E}_0). We will then say that the data is generated *under the null*. By contrast, if $\theta \neq 0$, we will say that the data is generated *under the alternative*, though we emphasise that in fact the alternative is composite, being indexed by $z \in \mathbb{N} \cup \{0\}$ and $\theta \in \mathbb{R}^p \setminus \{0\}$. In practice, in order for our procedure to have uniformly non-trivial power, it will be necessary to work with a subset of the alternative hypothesis parameter space that is well-separated from the null, in the sense that the ℓ_2 -norm of the vector of mean change, $\vartheta := \|\theta\|_2$, is at least a known lower bound $\beta > 0$.

A sequential changepoint procedure is an extended stopping time² N (with respect to the natural filtration) taking values in $\mathbb{N} \cup \{\infty\}$. Equivalently, we can think of it as a family of $\{0,1\}$ -valued estimators $(\hat{H}_n)_{n=1}^\infty$, where $\hat{H}_n = \hat{H}_n(X_1, \dots, X_n)$, and where the sequence is increasing in the sense that $\hat{H}_m(X_1, \dots, X_m) \leq \hat{H}_n(X_1, \dots, X_n)$ for $m \leq n$. Here, the correspondence arises from $\hat{H}_n = \mathbb{1}_{\{N \leq n\}}$ and $N = \inf\{n \in \mathbb{N} : \hat{H}_n = 1\}$, with the usual convention that $\inf \emptyset := \infty$.

We measure the performance of a sequential changepoint procedure via its responsiveness subject to a given upper bound on the false alarm rate, or equivalently, a lower bound on the average run length in the absence of change. Specifically, following the concepts introduced by Lorden (1971), we define the *patience* (this is sometimes referred to as the average run length under the null or average run length to false alarm in the literature) of a sequential changepoint procedure N to be $\mathbb{E}_0(N)$, and its *worst-case response delay* (likewise, this is sometimes referred to as the worst-worst-case average detection delay) to be

$$\overline{\mathbb{E}}_\theta^{\text{wc}}(N) := \sup_{z \in \mathbb{N} \cup \{0\}} \text{ess sup } \mathbb{E}_{z,\theta} \{(N - z) \vee 0 | X_1, \dots, X_z\}.$$

While controlling the worst-case response delay provides a very strong theoretical guarantee of the average detection delay of the procedure, even under the worst possible pre-change

²A random variable τ taking values in $\mathbb{N} \cup \{\infty\}$ is an *extended stopping time* with respect to the filtration $(\mathcal{F}_n)_{n \in \mathbb{N}}$ if $\{\tau = n\} \in \mathcal{F}_n$ for all $n \in \mathbb{N}$.

data sequence, obtaining a good bound for this quantity is often difficult. We therefore also consider the average-case response delay, or simply the *response delay* of a procedure N , defined as

$$\bar{\mathbb{E}}_{\theta}(N) := \sup_{z \in \mathbb{N} \cup \{0\}} \mathbb{E}_{z, \theta} \{(N - z) \vee 0\}.$$

We note that $\bar{\mathbb{E}}_{\theta}(N) \leq \bar{\mathbb{E}}_{\theta}^{\text{wc}}(N)$. A good sequential changepoint procedure should have small worst- and average-case response delays, uniformly over the relevant class of alternatives $\{\mathbb{P}_{z, \theta} : (z, \theta) \in (\mathbb{N} \cup \{0\}) \times \mathbb{R}^p, \|\theta\|_2 \geq \beta\}$, subject to its patience being at least some suitably large, pre-determined $\gamma > 0$. Finally, as mentioned above, we are interested in sequential changepoint procedures that are online, so that the computational complexity per additional observation should be a function of p only.

Our main contribution in this work is to propose, in Section 2, a new algorithm called `ocd` (short for online changepoint detection), for high-dimensional, online changepoint detection in the above setting. The procedure works by performing likelihood ratio tests against simple alternatives of different scales in each coordinate, and then aggregating test statistics across scales and coordinates for changepoint detection. The `ocd` algorithm has worst-case computational complexity $O(p^2 \log(ep))$ per new observation, so satisfies our requirement for being an online algorithm. In fact, as we explain in Section 2.1, the algorithmic complexity is often even better than this. Moreover, as we illustrate in Section 4, it has extremely effective empirical performance. In terms of theoretical guarantees, it turns out to be more convenient to analyse a slight variant of our initial algorithm, which we refer to as `ocd'`. This has the same order of computational complexity per new observation as `ocd`, but enables us to ensure that whenever we are yet to declare that a change has occurred, only post-change observations contribute to the running test statistics. In practice, the original `ocd` algorithm also appears to have this property for typical pre-change sequences, and we argue heuristically that there is a sense in which it is more efficient than `ocd'` by a factor of at most 2.

Our theoretical analysis in Section 3 initially considers separately versions of the `ocd'` algorithm best tuned towards settings where the vector θ of change is dense, and where it is sparse in an appropriate sense. We then present results for a combined, adaptive procedure that seeks the best of both worlds. In all cases, the appropriate version of `ocd'` has guaranteed patience, at least at the desired nominal level. In the (small-change) regime of primary interest, and when ϑ is of the same order as β , the response delay of `ocd'` is of order at most \sqrt{p}/ϑ^2 in the dense case, up to a polylogarithmic factor; this can be improved to order s/ϑ^2 , again up to a polylogarithmic factor, when the effective sparsity of θ is $s < \sqrt{p}$.

As alluded to above, there is a paucity of prior literature on multivariate, online changepoint problems, though exceptions include Tartakovsky et al. (2006), Mei (2010) and Zou et al. (2015). These works focus either on the case where both the pre- and post-change distributions are exactly known, or where, for each coordinate, both the sign and a lower bound on the magnitude of change, are known in advance. A number of methods have also been proposed that involve scanning a moving window of fixed size for changes (Chan, 2017; Soh & Chandrasekaran, 2017; Xie & Siegmund, 2013). Such methods can be effective when the signal-to-noise ratio is large enough that the change can be detected within the prescribed window, but may experience excessive response delay in other cases. Of course, the window size may be increased to compensate, but this correspondingly increases the computational

complexity and storage requirements, so allowing the window size to vary with the signal strength would fail to satisfy our definition of an online algorithm. Recently, Gösmann et al. (2020) presented a new monitoring procedure for changepoints in the mean structure of a high-dimensional time series. Their method is sequential but not online as we defined it here, because at each time point, their monitoring statistic is computed using the entire data history up to that point. We also mention that online changepoint detection has been studied in the econometrics literature, where the problem is often referred to as that of monitoring structural breaks (Chu et al., 1996; Leisch et al., 2000; Zeileis et al., 2005). These works have studied low-dimensional regression settings, and asymptotic theory has been provided on the probability of eventual declaration of change.

Numerical results illustrate the performance of our `ocd` algorithm in Section 4. Proofs of our main results are given in Section 5. All the auxiliary lemmas and their proofs are provided in the online supplementary material Chen et al. (2021).

1.1 | Notation

We write \mathbb{N}_0 for the set of all non-negative integers. For $d \in \mathbb{N}$, we write $[d] := \{1, \dots, d\}$. Given $a, b \in \mathbb{R}$, we denote $a \vee b := \max(a, b)$ and $a \wedge b := \min(a, b)$. For a set S , we use $\mathbb{1}_S$ and $|S|$ to denote its indicator function and cardinality, respectively. For a real-valued function f on a totally ordered set S , we write $\text{sargmax}_{x \in S} f(x) := \min \arg\max_{x \in S} f(x)$, the smallest maximiser of f in S . For a vector $v = (v^1, \dots, v^M)^\top \in \mathbb{R}^M$, we define $\|v\|_0 := \sum_{i=1}^M \mathbb{1}_{\{v^i \neq 0\}}$, $\|v\|_2 := \left\{ \sum_{i=1}^M (v^i)^2 \right\}^{1/2}$ and $\|v\|_\infty := \max_{i \in [M]} |v^i|$. In addition, we define $v^{-j} := (v^1, \dots, v^{j-1}, v^{j+1}, \dots, v^M)^\top \in \mathbb{R}^{M-1}$. For a matrix $A = (A^{i,j}) \in \mathbb{R}^{d_1 \times d_2}$ and $j \in [d_2]$, we write $A^{\cdot,j} := (A^{1,j}, \dots, A^{d_1,j})^\top \in \mathbb{R}^{d_1}$ and $A^{-j,j} := (A^{1,j}, \dots, A^{j-1,j}, A^{j+1,j}, \dots, A^{d_1,j})^\top \in \mathbb{R}^{d_1-1}$. We use $\Phi(\cdot)$ and $\phi(\cdot)$ to denote the distribution function and density function of the standard normal distribution, respectively. For two real-valued random variables U and V , we write $U \succeq_{\text{st}} V$ if $\mathbb{P}(U \leq x) \leq \mathbb{P}(V \leq x)$ for all $x \in \mathbb{R}$. We adopt the convention that an empty sum is 0.

2 | AN ONLINE CHANGEPOINT PROCEDURE

2.1 | The `ocd` algorithm

In this section, we describe our online changepoint procedure, `ocd`, in more detail. As mentioned in the introduction, the procedure aggregates likelihood ratio test statistics against simple alternatives of different scales in different coordinates. For $i \in [n]$ and $j \in [p]$, we write X_i^j for the j th coordinate of X_i . If we want to test a null of $\mathcal{N}(0, 1)$ against a simple post-change alternative distribution of $\mathcal{N}(b, 1)$ for some $b \neq 0$ in coordinate $j \in [p]$, by Page (1954), the optimal online changepoint procedure is to declare that a change has occurred by time n when the test statistic

$$R_{n,b}^j := \max_{0 \leq h \leq n} \sum_{i=n-h+1}^n b(X_i^j - b/2) \quad (2)$$

exceeds a certain threshold. Note that $\sum_{i=n-h+1}^n b(X_i^j - b/2)$ can be viewed as the likelihood ratio test statistic between the null and this simple alternative using the tail sequence X_{n-h+1}, \dots, X_n . Thus, $R_{n,b}^j$ can be regarded as the most extreme of these likelihood ratio statistics, over all possible starting points for the tail sequence. Write

$$t_{n,b}^j := \operatorname{sargmax}_{0 \leq h \leq n} \sum_{i=n-h+1}^n b(X_i^j - b/2) \quad (3)$$

for the length of the tail sequence in which the associated likelihood ratio statistic (in the j th coordinate) is maximised. One way to aggregate across the p coordinates would be to use $\sum_{j=1}^p R_{n,b}^j$ as a test statistic. However, this approach is not ideal for two reasons. First, the exact distribution of the tail likelihood ratio statistic $R_{n,b}^j$ is hard to obtain, making it difficult to analyse the aggregated statistic under the null. More importantly, this aggregated statistic uses the same simple alternative $\mathcal{N}(b, 1)$ in all coordinates, and so even after varying the magnitude of b , it is only effective against a very limited set of alternative distributions in $\{\mathbb{P}_{z,\theta} : z \in \mathbb{N}, \|\theta\|_2 \geq \beta\}$, namely those for which the change is of very similar magnitude in all coordinates. In order to overcome these problems, our procedure uses the coordinate-wise statistics $(R_{n,b}^j : j \in [p])$, which we call ‘diagonal statistics’, to detect changes that have a large proportion of their signal concentrated in one coordinate. To detect denser changes, for each $j \in [p]$, we also compute tail partial sums of length $t_{n,b}^j$ in all other coordinates $j' \neq j$, given by

$$A_{n,b}^{j'j} := \sum_{i=n-t_{n,b}^j+1}^n X_i^{j'},$$

and aggregate them to form an ‘off-diagonal statistic’ anchored at coordinate j . Note that the number of summands in $A_{n,b}^{j'j}$ depends only on the observed data in the j th coordinate, and not on the data being aggregated in the j' th coordinate. These off-diagonal statistics are used to detect changes whose signal is not concentrated in a single coordinate. Intuitively, if a change has occurred and $\theta^j/b \geq 1$, then we can expect the tail length in coordinate j to be roughly of order $n-z$ for sufficiently large n , and this will ensure that the off-diagonal statistic anchored at coordinate j is close to the generalised likelihood ratio test statistic between the null and the composite alternative $\{\mathbb{P}_{z,\theta} : \|\theta\|_2 \neq 0\}$. If, in addition, a non-trivial proportion of the signal is contained in coordinates $[p] \setminus \{j\}$, then this statistic will be powerful for detecting the change.

The full description of the `ocd` procedure is given in Algorithm 1. Note that for notational simplicity, we have suppressed the time dependence of many variables as they are updated recursively in the algorithm. In the following, when necessary, we will make this dependence explicit by writing $A_{n,b}, t_{n,b}, Q_{n,b}, S_n^{\text{diag}}$ and S_n^{off} for the relevant quantities at the end of the n th iteration of the repeat loop.

By Lemma 1, $bA_{n,b}^{jj} - b^2 t_{n,b}^j/2$, as defined in the algorithm, is equal to the quantity $R_{n,b}^j$ defined in Equation (2) (we will also suppress its n dependence when it is clear from the context). Moreover, by Lemma 2, the two definitions of $t_{n,b}^j$ from Algorithm 1 and Equation (3) coincide. In the algorithm, we allow b to range over the (signed) dyadic grid $\mathcal{B} \cup \mathcal{B}_0$, since the maximal signal strength in individual coordinates, $\|\theta\|_\infty$, can range from ϑ/\sqrt{p} to ϑ . In this way, the algorithm automatically adapts to different signal strengths in each coordinate. Here, the inclusion of \mathcal{B}_0 and the extra logarithmic factors in the denominators of elements of $\mathcal{B} \cup \mathcal{B}_0$ appear due to technical reasons in the theoretical analysis of the algorithm.

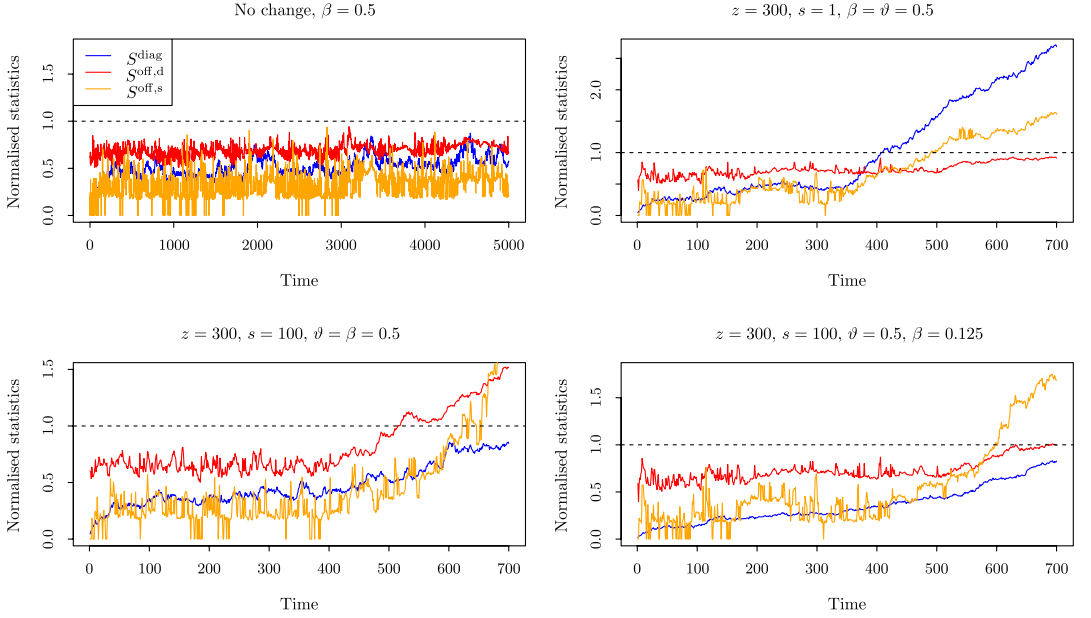


FIGURE 1 Behaviour of the three normalised statistics in `ocd` under the null and under the alternative with different signal strength, sparsity level and assumed lower bound. A change is declared as soon as one of these three normalised statistics exceeds 1. The data were generated in the top-left panel according to \mathbb{P}_0 , and, in the other panels, according to $\mathbb{P}_{z,\theta}$ with $p = 100$, $z = 300$ and $\theta = \vartheta U$, where U is uniformly distributed on the union of all s -sparse unit spheres in \mathbb{R}^p (see Section 4.2 for a more detailed description)

Algorithm 1: Pseudo-code of the `ocd` algorithm

Input: $X_1, X_2 \dots \in \mathbb{R}^p$ observed sequentially, $\beta > 0$, $a \geq 0$, $T^{\text{diag}} > 0$ and $T^{\text{off}} > 0$

Set: $\mathcal{B} = \left\{ \pm \frac{\beta}{\sqrt{2^\ell \log_2(2p)}} : \ell = 0, \dots, \lfloor \log_2 p \rfloor \right\}$, $\mathcal{B}_0 = \left\{ \pm \frac{\beta}{\sqrt{2^{\lfloor \log_2 p \rfloor + 1} \log_2(2p)}} \right\}$, $n = 0$,

$A_b = \mathbf{0} \in \mathbb{R}^{p \times p}$ and $t_b = 0 \in \mathbb{R}^p$ for all $b \in \mathcal{B} \cup \mathcal{B}_0$

repeat

$n \leftarrow n + 1$

 observe new data vector X_n

for $(j, b) \in [p] \times (\mathcal{B} \cup \mathcal{B}_0)$ **do**

$t_b^j \leftarrow t_b^j + 1$

$A_b^{j,j} \leftarrow A_b^{j,j} + X_n$

if $bA_b^{j,j} - b^2 t_b^j / 2 \leq 0$ **then**

$t_b^j \leftarrow 0$ and $A_b^{j,j} \leftarrow 0$

 compute $Q_b^j \leftarrow \sum_{j' \in [p]: j' \neq j} \frac{(A_b^{j',j})^2}{t_b^{j'} \vee 1} \mathbb{1}_{\{|A_b^{j',j}| \geq a\sqrt{t_b^j}\}}$

$S^{\text{diag}} \leftarrow \max_{(j,b) \in [p] \times (\mathcal{B} \cup \mathcal{B}_0)} (bA_b^{j,j} - b^2 t_b^j / 2)$

$S^{\text{off}} \leftarrow \max_{(j,b) \in [p] \times \mathcal{B}} Q_b^j$

until $S^{\text{diag}} \geq T^{\text{diag}}$ or $S^{\text{off}} \geq T^{\text{off}}$,

Output: $N = n$

Algorithm 1 uses S^{diag} and S^{off} to aggregate diagonal and off-diagonal statistics, respectively, as mentioned above, and declares that a change has occurred as soon as either of these quantities exceeds its own pre-determined threshold. As mentioned previously, S^{diag} tracks the maximum of R_b^j over all scales b and coordinates j . Before introducing S^{off} , we first discuss the off-diagonal statistics Q_b^j in Algorithm 1, which are ℓ_2 aggregations of normalised tail sums $A_b^{j',j} / \sqrt{t_b^j \vee 1}$, each

hard-thresholded at level a . The hard thresholding level can be chosen to detect dense or sparse signals θ ; in the sparse case a non-zero a facilitates an aggregation that aims to exclude coordinates with negligible change (thereby reducing the variance of the normalised tail sums). Finally, S^{off} is computed as the maximum of the Q_b^j over all anchoring coordinates $j \in [p]$ and scales $b \in \mathcal{B}$.

Although the off-diagonal statistics described in the previous paragraph are effective for detecting changes when the signal sparsity is known, it is desirable to the practitioner to have a combined procedure that adapts to the sparsity level. This may be computed straightforwardly by tracking S^{off} for $a = a^{\text{dense}}$ and $a = a^{\text{sparse}}$, as well as S^{diag} , and declaring a change when any of these three statistics exceeds a suitable threshold. Figure 1 illustrates the performance of this adaptive procedure, together with the time evolution of normalised versions of all three statistics tracked, in synthetic data sets both with and without a change. This adaptive procedure is analysed theoretically in Section 3 and empirically in Section 4.

The `ocd` procedure satisfies our definition of an online algorithm. Indeed, for each new observation X_n , `ocd` updates $t_{n,b} \in \mathbb{R}^p$ and $A_{n,b} \in \mathbb{R}^{p \times p}$ for $O(p^2 \log(ep))$ different values of b . It then computes S_n^{diag} and S_n^{off} via $A_{n,b}$. These steps require $O(p^2 \log(ep))$ operations. Moreover, the total storage used is $O(p^2 \log(ep))$ throughout the algorithm.

In fact, the computational complexity of `ocd` can often be reduced, because typically $\mathcal{T} := \{t_b^j : j \in [p], b \in \mathcal{B}\}$ has cardinality much less than $p|\mathcal{B}|$ (which is the worst case, when all elements are distinct). Correspondingly, at each time step, we need only store the $p \times |\mathcal{T}|$ matrix $(B^{k,t})_{k \in [p], t \in \mathcal{T}}$ given by $B^{k,t_b} := A_b^{k,j}$, resulting in an improved per-iteration computational complexity and storage for `ocd` of $O(p|\mathcal{T}|)$. For simplicity of exposition, we have not presented this computational speed-up in Algorithm 1, and it appears to be difficult to provide theoretical guarantees on $|\mathcal{T}|$. Nevertheless, we have implemented the algorithm in this form in the `R` package `ocd` (Chen et al., 2020), and have found it to provide substantial computational savings in practice.

2.2 | A slight variant of `ocd`

While the `ocd` algorithm performs very well numerically, it turns out to be easier theoretically to analyse a slight variant, which we call `ocd'`, and describe in Algorithm 2. Again, we have suppressed the time dependence n of many variables including $\tau_{n,b}$, $\tilde{\tau}_{n,b}$, $\Lambda_{n,b}$ and $\tilde{\Lambda}_{n,b}$ in the algorithm. The main difference between these two algorithms is that in `ocd'`, the off-diagonal statistics Q_b^j are computed using tail partial sums of length τ_b^j instead of t_b^j . These new tail partial sums are recorded in $\Lambda_b \in \mathbb{R}^{p \times p}$.

By Lemma 9, we always have

$$t_b^j/2 \leq \tau_b^j < 3t_b^j/4 \quad (4)$$

whenever $t_b^j \geq 2$. In this sense, the tail sample size used by `ocd'` is smaller than that of `ocd` by a factor of at most 2. The benefit of using a shorter tail in `ocd'` is that when n exceeds a known, deterministic threshold, we can be sure that whenever we have not declared that a change has occurred by time z , the tail partial sum consists exclusively of post-change observations. In practice, we observe that even in Algorithm 1, the tail lengths $t_{z,b}^j$ at the changepoint are generally very short for many coordinates, so the inclusion of a few pre-change observations in the tail partial sum calculation does not significantly affect the efficacy of the changepoint detection procedure. The practical performance of Algorithm 1 is statistically

Algorithm 2: Pseudo-code of the ocd' algorithm, a slight variant of ocd

Input: $X_1, X_2, \dots \in \mathbb{R}^p$ observed sequentially, $\beta > 0$, $a \geq 0$, $T^{\text{diag}} > 0$ and $T^{\text{off}} > 0$.

Set: $\mathcal{B} = \left\{ \pm \frac{\beta}{\sqrt{2^\ell \log_2(2p)}} : \ell = 0, \dots, \lfloor \log_2 p \rfloor \right\}$, $\mathcal{B}_0 = \left\{ \pm \frac{\beta}{\sqrt{2^{\lfloor \log_2 p \rfloor + 1} \log_2(2p)}} \right\}$, $n = 0$,

$A_b = \Lambda_b = \tilde{\Lambda}_b = \mathbf{0} \in \mathbb{R}^{p \times p}$ and $t_b = \tau_b = \tilde{\tau}_b = 0 \in \mathbb{R}^p$ for all $b \in \mathcal{B} \cup \mathcal{B}_0$

repeat

$n \leftarrow n + 1$

 observe new data vector X_n

for $(j, b) \in [p] \times (\mathcal{B} \cup \mathcal{B}_0)$ **do**

$t_b^j \leftarrow t_b^j + 1$ and $A_b^{j,j} \leftarrow A_b^{j,j} + X_n$

 set $\delta = 0$ if t_b^j is a power of 2 and $\delta = 1$ otherwise.

$\tau_b^j \leftarrow \tau_b^j \delta + \tilde{\tau}_b^j (1 - \delta) + 1$ and $\Lambda_b^{j,j} \leftarrow \Lambda_b^{j,j} \delta + \tilde{\Lambda}_b^{j,j} (1 - \delta) + X_n$

$\tilde{\tau}_b^j \leftarrow (\tilde{\tau}_b^j + 1) \delta$ and $\tilde{\Lambda}_b^{j,j} \leftarrow (\tilde{\Lambda}_b^{j,j} + X_n) \delta$.

if $bA_b^{j,j} - b^2 t_b^j / 2 \leq 0$ **then**

$t_b^j \leftarrow \tau_b^j \leftarrow \tilde{\tau}_b^j \leftarrow 0$

$A_b^{j,j} \leftarrow \Lambda_b^{j,j} \leftarrow \tilde{\Lambda}_b^{j,j} \leftarrow 0$

 compute $Q_b^j \leftarrow \sum_{j' \in [p]: j' \neq j} \frac{(\Lambda_b^{j',j})^2}{\tau_b^j \vee 1} \mathbb{1}_{\{|\Lambda_b^{j',j}| \geq a \sqrt{\tau_b^j}\}}$

$S^{\text{diag}} \leftarrow \max_{(j,b) \in [p] \times (\mathcal{B} \cup \mathcal{B}_0)} (bA_b^{j,j} - b^2 t_b^j / 2)$

$S^{\text{off}} \leftarrow \max_{(j,b) \in [p] \times \mathcal{B}} Q_b^j$

until $S^{\text{diag}} \geq T^{\text{diag}}$ or $S^{\text{off}} \geq T^{\text{off}}$,

Output: $N = n$

more efficient than Algorithm 2 in many settings by a factor of between 4/3 and 2, as suggested by Equation (4). By construction, τ_b^j and $\Lambda_b^{j,j}$ are computable online, through auxiliary variables $\tilde{\tau}_b^j$ and $\tilde{\Lambda}_b^{j,j}$. Indeed, Algorithm 2 is also an online algorithm, with overall computational complexity per observation and storage remaining at $O(p^2 \log(ep))$ in the worst case; similar computational improvements to those mentioned for ocd at the end of Section 2.1 are also possible here.

3 | THEORETICAL ANALYSIS

As mentioned in Section 2, the input a in Algorithms 1 and 2 allows users to detect changepoints of different sparsity levels. More precisely, for any $\theta \in \mathbb{R}^p$, we have by Lemma 8 that there exists a smallest $s(\theta) \in \{2^0, 2^1, \dots, 2^{\lfloor \log_2 p \rfloor}\}$ such that the set

$$S(\theta) := \left\{ j \in [p] : |\theta^j| \geq \frac{\|\theta\|_2}{\sqrt{s(\theta) \log_2(2p)}} \right\}$$

has cardinality at least $s(\theta)$. On the other hand, we also have $|S(\theta)| \leq s(\theta) \log_2(2p)$. We call $s(\theta)$ the *effective sparsity* of the vector θ and $S(\theta)$ its *effective support*. Intuitively, the sum of squares of coordinates in the effective support of θ has the same order of magnitude as $\|\theta\|_2^2$, up to logarithmic factors. Moreover, if θ is an s -sparse vector in the sense that $\|\theta\|_0 \leq s$, then $s(\theta) \leq s$, and the equality is attained when, for example, all non-zero coordinates have the same magnitude.

In this section, we initially analyse the theoretical performance of Algorithm 2 for two different choices of a in $S^{\text{off}} = S^{\text{off}}(a)$, namely $a = 0$ and $a = \sqrt{8 \log(p-1)}$. We then present our combined, adaptive procedure and its performance guarantees.

Define $N^{\text{diag}} := \inf\{n: S_n^{\text{diag}} \geq T^{\text{diag}}\}$ and $N^{\text{off}} = N^{\text{off}}(a) := \inf\{n: S_n^{\text{off}}(a) \geq T^{\text{off}}\}$. Then the stopping time for our changepoint detection procedure is simply $N = N(a) = N^{\text{diag}} \wedge N^{\text{off}}(a)$.

3.1 | Dense case

Here, we analyse the changepoint detection procedure $N = N(0)$, which, as we will see, is most suitable for detecting dense mean changes in the sense that $s(\theta) \geq \sqrt{p}$ (though we do not assume this in our theory). In this case, when $p \geq 2$ and conditionally on τ_b^j , the quantity Q_b^j follows a chi-squared distribution with $p - 1$ degrees of freedom under the null, provided that τ_b^j is positive. (When $p = 1$, we have that $Q_b^j = 0$ for all $j \in [p]$ and $b \in \mathcal{B}$, so $S^{\text{off}} = 0$ and the off-diagonal statistic never triggers the declaration of a change. Similarly, if $p \geq 2$ but $\tau_{n,b}^j = 0$, then we also have $Q_{n,b}^j = 0$.) Motivated by the chi-squared tail bound of Laurent and Massart (2000, Lemma 1), we choose a threshold of the form

$$T^{\text{off}} := p - 1 + \tilde{T}^{\text{off}} + \sqrt{2(p-1)\tilde{T}^{\text{off}}} =: \psi(\tilde{T}^{\text{off}}), \quad (5)$$

say, for some $\tilde{T}^{\text{off}} > 0$.

The following theorem provides control of the patience of ocd' .

Theorem 1 *Let X_1, X_2, \dots be generated according to \mathbb{P}_0 . For any $\gamma \geq 1$, let $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $a = 0$, $T^{\text{diag}} = \log\{16p\gamma \log_2(4p)\}$ and $T^{\text{off}} = \psi(\tilde{T}^{\text{off}})$ with $\tilde{T}^{\text{off}} = 2 \log\{16p\gamma \log_2(2p)\}$ be the inputs of Algorithm 2, with corresponding output N . Then $\mathbb{E}_0(N) \geq \gamma$.*

We note that either of the two statistics S^{diag} and S^{off} may trigger a false alarm under the null. The two threshold levels T^{diag} and T^{off} are chosen so that $\mathbb{E}_0(N^{\text{diag}})$ and $\mathbb{E}_0(N^{\text{off}})$ have comparable upper bounds. We also remark that although Theorem 1 as stated only controls the expected value of N under the null, careful examination of the proof reveals that we can also control $\mathbb{P}_0(N \leq m)$ for every $m \in \mathbb{N}$. More precisely, from Equations (15) and (16) in the proof, we can deduce that

$$\mathbb{P}_0(N \leq m) \leq \frac{m}{4\gamma}$$

for every $m \in \mathbb{N}$. The same bound holds for our other patience control results below, though we omit formal statements for brevity.

Our next result controls the response delay of ocd' in both worst-case and average senses.

Theorem 2 *Assume that X_1, X_2, \dots are generated according to $\mathbb{P}_{z,\theta}$ for some z and θ such that $\|\theta\|_2 = \vartheta \geq \beta > 0$ and that θ has an effective sparsity of $s := s(\theta)$. Then there exists a universal constant $C > 0$, such that the output N from Algorithm 2, with inputs $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $a = 0$, $T^{\text{diag}} = \log\{16p\gamma \log_2(4p)\}$ and $T^{\text{off}} = \psi(\tilde{T}^{\text{off}})$ with $\tilde{T}^{\text{off}} = 2 \log\{16p\gamma \log_2(2p)\}$, satisfies*

$$\bar{\mathbb{E}}_{\theta}^{\text{wc}}(N) \leq C \left\{ \frac{\sqrt{p} \log(ep\gamma)}{\vartheta^2} \vee \frac{s \log(ep\gamma) \log(ep)}{\beta^2} \vee 1 \right\}. \quad (6)$$

Furthermore, there exists $\beta_0(s) > 0$, depending only on s , such that for all $\beta \leq \beta_0(s)$, the output N satisfies

$$\bar{\mathbb{E}}_{\theta}(N) \leq C \left\{ \frac{\sqrt{p} \log(ep\gamma)}{\vartheta^2} \vee \frac{\sqrt{s} \log(ep/\beta) \log(ep)}{\beta^2} \vee 1 \right\}, \quad (7)$$

for $s \geq 2$, and

$$\bar{\mathbb{E}}_{\theta}(N) \leq C \left\{ \frac{\log(ep\gamma) \log(ep)}{\beta \vartheta} \vee 1 \right\}, \quad (8)$$

for $s = 1$.

We defer detailed discussion of our response delay bounds until after we have presented our adaptive procedure in Section 3.3.

3.2 | Sparse case

We now assume that $p \geq 2$, and analyse the performance of $N = N(\sqrt{8 \log(p-1)})$; in other words, we choose $a = \sqrt{8 \log(p-1)}$. This choice turns out to work particularly well when the vector of mean change is sparse in the sense that $s(\theta) \leq \sqrt{p}$, though again we do not assume this in our theory. The motivation for this choice of a comes from the fact that, for fixed b and j , we have $\Lambda_b^{j',j} | \tau_b^j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau_b^j)$ for $j' \in [p] \setminus \{j\}$ under the null. Since a is the threshold level for $|\Lambda_b^{j',j}| / \sqrt{\tau_b^j}$, it is therefore natural to choose a to be of the same order as the maximum absolute value of $p-1$ independent and identically distributed $\mathcal{N}(0, 1)$ random variables. The declaration threshold T^{off} is determined based on Lemma 10. Theorem 3 below shows that, in the sparse case, the patience of our procedure is also guaranteed to be at least at the nominal level $\gamma > 0$. In addition, as in the dense case, we can also control the response delay of ocd' according to Theorem 4.

Theorem 3 *Let X_1, X_2, \dots be generated according to \mathbb{P}_0 . For any $\gamma \geq 1$, let $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $a = \sqrt{8 \log(p-1)}$, $T^{\text{diag}} = \log\{16p\gamma \log_2(4p)\}$ and $T^{\text{off}} = 8 \log\{16p\gamma \log_2(2p)\}$ be the inputs of Algorithm 2, with corresponding output N . Then $\mathbb{E}_0(N) \geq \gamma$.*

Theorem 4 *Assume that X_1, X_2, \dots are generated according to $\mathbb{P}_{z,\theta}$ for some z and θ such that $\|\theta\|_2 = \vartheta \geq \beta > 0$ and that θ has an effective sparsity of $s := s(\theta)$. Then there exists a universal constant $C > 0$, such that the output N from Algorithm 2, with inputs $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $a = \sqrt{8 \log(p-1)}$, $T^{\text{diag}} = \log\{16p\gamma \log_2(4p)\}$ and $T^{\text{off}} = 8 \log\{16p\gamma \log_2(2p)\}$, satisfies*

$$\bar{\mathbb{E}}_{\theta}(N) \leq \bar{\mathbb{E}}_{\theta}^{\text{wc}}(N) \leq C \left\{ \frac{s \log(ep\gamma) \log(ep)}{\beta^2} \vee 1 \right\}. \quad (9)$$

Comparing Theorems 2 and 4, we see that the thresholding induced by the non-zero choice of $a = \sqrt{8 \log(p-1)}$ in Theorem 4 facilitates an improved dependence on the effective sparsity s in the bound on the response delay, whenever s is of smaller order than \sqrt{p} .

3.3 | Adaptive procedure

To adapt to different sparsity levels s , we can run ocd (or ocd') with two values of a simultaneously: we choose $a = a^{\text{dense}} = 0$ to form the off-diagonal dense statistic $S^{\text{off,d}} = S^{\text{off}}(a^{\text{dense}})$, and $a = a^{\text{sparse}} = \sqrt{8 \log(p-1)}$ to form the off-diagonal sparse statistic $S^{\text{off,s}} = S^{\text{off}}(a^{\text{sparse}})$. We recall that the diagonal statistic S^{diag} does not depend on the choice of a . For clarity, we redefine the three stopping times here: $N^{\text{diag}} := \inf\{n : S_n^{\text{diag}} \geq T^{\text{diag}}\}$, $N^{\text{off,d}} := \inf\{n : S_n^{\text{off,d}} \geq T^{\text{off,d}}\}$ and $N^{\text{off,s}} := \inf\{n : S_n^{\text{off,s}} \geq T^{\text{off,s}}\}$, for appropriately chosen thresholds T^{diag} , $T^{\text{off,d}}$ and $T^{\text{off,s}}$. The output of this adaptive procedure is thus $N = N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}}$.

The following results provide patience and response delay guarantees for this adaptive procedure.

Theorem 5 *Let X_1, X_2, \dots be generated according to \mathbb{P}_0 . For any $\gamma \geq 1$, let $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $T^{\text{diag}} = \log\{24p\gamma \log_2(4p)\}$, $T^{\text{off,d}} = \psi(\tilde{T}^{\text{off,d}})$ with $\tilde{T}^{\text{off,d}} = 2 \log\{24p\gamma \log_2(2p)\}$ and $T^{\text{off,s}} = 8 \log\{24p\gamma \log_2(2p)\}$ be the inputs of the adaptive version of Algorithm 2, with corresponding output N . Then $\mathbb{E}_0(N) \geq \gamma$.*

Theorem 6 *Assume that X_1, X_2, \dots are generated according to $\mathbb{P}_{z,\theta}$ for some z and θ such that $\|\theta\|_2 = \vartheta \geq \beta > 0$ and that θ has an effective sparsity of $s := s(\theta)$. Then there exists a universal constant $C > 0$, such that the output N from the adaptive version of Algorithm 2, with inputs $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $T^{\text{diag}} = \log\{24p\gamma \log_2(4p)\}$, $T^{\text{off,d}} = \psi(\tilde{T}^{\text{off,d}})$ with $\tilde{T}^{\text{off,d}} = 2 \log\{24p\gamma \log_2(2p)\}$ and $T^{\text{off,s}} = 8 \log\{24p\gamma \log_2(2p)\}$, satisfies*

$$\overline{\mathbb{E}}_{\theta}^{\text{wc}}(N) \leq C \left\{ \frac{s \log(ep\gamma) \log(ep)}{\beta^2} \vee 1 \right\}. \quad (10)$$

Furthermore, there exists $\beta_0(s) \in (0, 1/2]$, depending only on s , such that for all $\beta \leq \beta_0(s)$, the output N satisfies

$$\overline{\mathbb{E}}_{\theta}(N) \leq C \left\{ \left(\frac{\sqrt{p} \log(ep\gamma)}{\vartheta^2} \vee \frac{\sqrt{s} \log(ep\beta^{-1}) \log(ep)}{\beta^2} \right) \wedge \frac{s \log(ep\gamma) \log(ep)}{\beta^2} \right\}, \quad (11)$$

for $s \geq 2$, and

$$\overline{\mathbb{E}}_{\theta}(N) \leq \frac{C \log(ep\gamma) \log(ep)}{\beta^2}, \quad (12)$$

for $s = 1$.

Comparing these two results with the corresponding theorems in Sections 3.1 and 3.2, we see that by choosing slightly more conservative thresholds, the adaptive procedure retains the nominal patience control while (up to constant factors) achieving the best of both worlds in terms of its response delay guarantees under different sparsity regimes.

To better understand the worst-case and average-case response delay bounds in Theorem 6, it is helpful to assume that $\vartheta/C_1 \leq \beta \leq \vartheta \leq C_1$ and $\log(\gamma/\beta) \leq C_2 \log p$ for some $C_1, C_2 > 0$. Under these additional assumptions, the result of Theorem 6 takes the simpler form that for some $C > 0$, depending only on C_1 and C_2 , we have

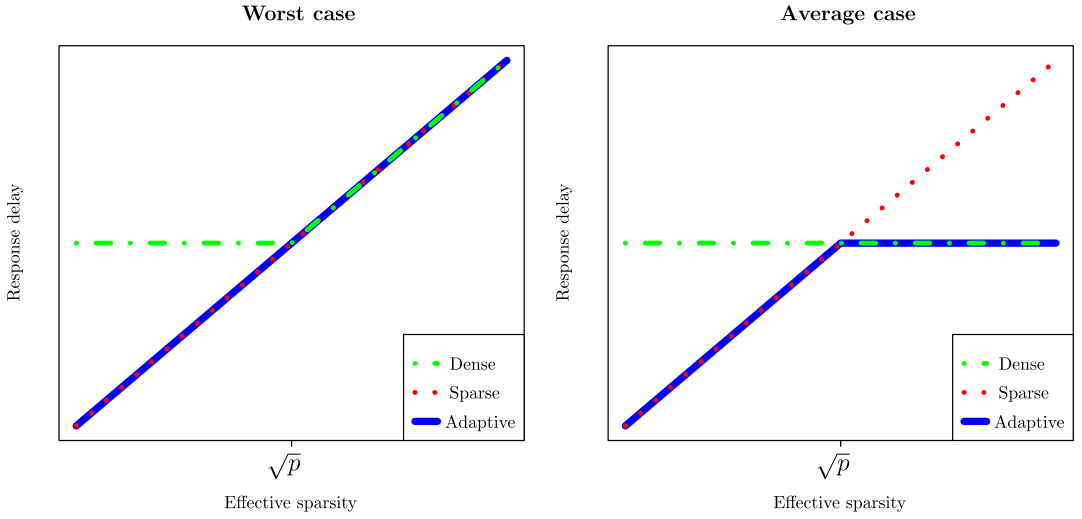


FIGURE 2 Illustration of the dependencies on sparsity of the worst-case and average-case response delays for the dense, sparse and adaptive versions of ocd' , as given by Theorems 2, 4 and 6

$$\bar{\mathbb{E}}_{\theta}^{\text{wc}}(N) \leq \frac{Cs \log^2(ep)}{\vartheta^2} \quad \text{and} \quad \bar{\mathbb{E}}_{\theta}(N) \leq \frac{C(s \wedge p^{1/2}) \log^2(ep)}{\vartheta^2}.$$

In particular, the average-case response delay upper bound exhibits a phase transition when the effective sparsity level s is of order \sqrt{p} , which is the boundary between the sparse and dense cases. Similar sparsity-related elbow effects have been observed in the minimax rate for high-dimensional Gaussian mean testing (Collier et al., 2017) and the corresponding offline changepoint detection problem (Liu et al., 2021). On the other hand, we note that quadratic dependence on ϑ in the denominator, and the logarithmic dependence on γ in the numerator, are known to be optimal in the case when $p = 1$ (Lorden, 1971, Theorem 3). The different dependencies on sparsity of the worst-case and average-case response delays for the dense, sparse and adaptive versions of ocd' are illustrated in Figure 2.

3.4 | Relaxation of assumptions

The setting we consider for our theoretical results, with independent Gaussian observations having identity covariance matrix, is convenient for facilitating a relatively clean presentation and to clarify the main ideas behind the ocd procedure. Nevertheless, it is of interest to consider more general data-generating mechanisms, where these assumptions are relaxed. Focusing on the dense case for simplicity of exposition, the Gaussianity assumption ensures that our aggregated statistics have chi-squared distributions (under the null) or non-central chi-squared distributions (under the alternative), so we can apply existing sharp tail bounds. If, instead, our observations have sub-Gaussian distributions, then the corresponding statistics would have sub-Gamma distributions, in the terminology of Boucheron et al. (2013), so Bernstein's inequality could be applied to give alternative bounds in this setting. Another place where we make use of the Gaussianity assumption is in comparing the trajectories of our test statistics with a Brownian

motion with drift (see, for instance, the proof of Lemma 6). Since we can view these trajectories as discrete Gaussian random walks, we can establish direct inequalities in this comparison. If we were to relax the Gaussianity, then we would need to rely on Donsker's invariance principle, or preferably its finite-sample version given by the Hungarian embedding (Komlós et al., 1976).

In cases where the covariance matrix of the observations were unknown, it may be possible to estimate this using a training sample, known to come from the null hypothesis, and use this to pre-whiten the data. The form of the estimator to be used should be chosen to exploit any known dependence structure (e.g. banding, Toeplitz or tapering) between the different coordinates. Similar remarks apply when there is short-range serial (temporal) dependence between successive observations. In Section 4.4, we demonstrate one way of handling temporal dependence with real data, by studying the residuals of the fit of an autoregressive model.

4 | NUMERICAL STUDIES

In this section, we study the empirical performance of the `ocd` algorithm and compare it with other online changepoint detection methods. Recall that the (adaptive) `ocd` algorithm declares a change when any of the three statistics S^{diag} , $S^{\text{off,d}}$ and $S^{\text{off,s}}$ exceeds their respective thresholds T^{diag} , $T^{\text{off,d}}$ and $T^{\text{off,s}}$. If a priori knowledge about the signal sparsity is available, it may be slightly preferable to use $N^{\text{diag}} \wedge N^{\text{off,d}}$ in the dense case, and $N^{\text{diag}} \wedge N^{\text{off,s}}$ in the sparse case, but for simplicity of exposition, we will focus on the adaptive version of our `ocd` procedure throughout the remainder of this section. While the threshold choices given in Theorem 5 guarantee that the patience of (adaptive) `ocd` will be at least at the nominal level, in practice, they may be conservative. We therefore describe a scheme for practical choice of thresholds in Section 4.1. Recall that, in order to form $S^{\text{off,d}}$ and $S^{\text{off,s}}$, two different entrywise hard thresholds for $A_b^{j',j} / \sqrt{t_b^j} \vee 1$ need to be specified. For $S^{\text{off,d}}$, we choose $a = 0$ for both theoretical analysis and practical usage. For $S^{\text{off,s}}$, the theoretical choice is $a = \sqrt{8 \log(p-1)}$, but since this is also slightly conservative, the choice of $a = \sqrt{2 \log p}$ is used in our practical implementation of the algorithm, and our numerical simulations below.

4.1 | Practical choice of declaration thresholds

The purpose of this section is to introduce an alternative to using the theoretical thresholds T^{diag} , $T^{\text{off,d}}$ and $T^{\text{off,s}}$ provided by Theorem 5, namely to determine the thresholds through Monte Carlo simulation. The basic idea is that since the null distribution is known, we can simulate from it to determine the patience for any given choice of thresholds. A complicating issue is the fact that the choices of the three thresholds T^{diag} , $T^{\text{off,d}}$ and $T^{\text{off,s}}$ are related, so that we may be able to achieve the same patience by increasing T^{diag} and decreasing $T^{\text{off,d}}$, for example. To handle this, we first argue that the renewal nature of the processes involved means that, at least for moderately large thresholds, the times to exceedance for each of the three statistics S^{diag} , $S^{\text{off,d}}$ and $S^{\text{off,s}}$ are approximately exponentially distributed. Evidence to support this is provided by Figure 3, where we present QQ-plots of $N^{\text{diag}}/m(N^{\text{diag}})$, $N^{\text{off,d}}/m(N^{\text{off,d}})$ and $N^{\text{off,s}}/m(N^{\text{off,s}})$, where the $m(N)$ statistics are empirical medians of the corresponding N statistics (divided by $\log 2$) over 200 repetitions.

We can therefore set an individual Monte Carlo threshold for S^{diag} as follows (the other two statistics can be handled in identical fashion): for $r \in [B]$, simulate $X_1^{(r)}, \dots, X_p^{(r)} \stackrel{\text{iid}}{\sim} \mathcal{N}_p(0, I_p)$

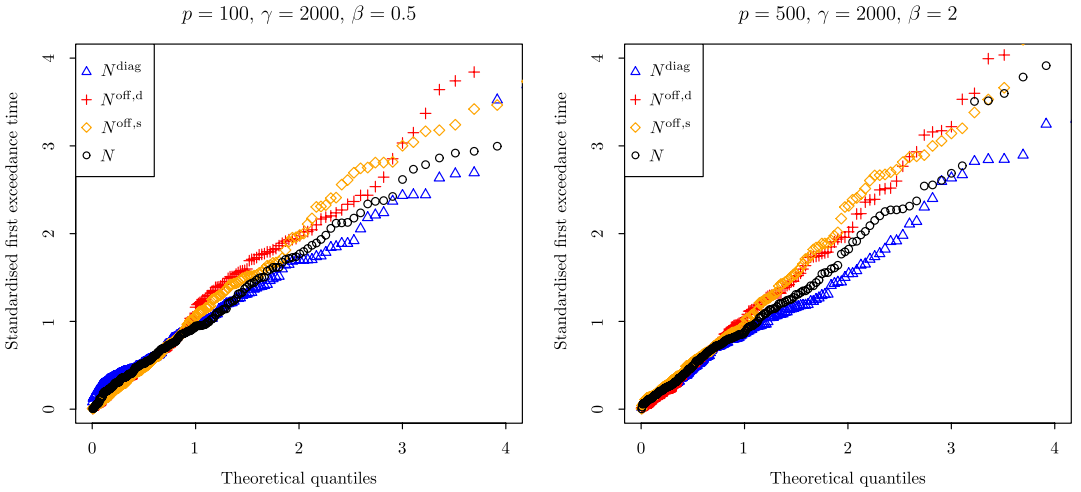


FIGURE 3 QQ-plots of standardised versions of N^{diag} , $N^{\text{off,d}}$ and $N^{\text{off,s}}$, as well as $N = N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}}$, against theoretical $\text{Exp}(1)$ quantiles

TABLE 1 Estimated run lengths under the null using the Monte Carlo thresholds described in Section 4.1 over 500 repetitions, with desired patience level $\gamma = 5000$. Algorithm is terminated after 20000 data points for each repetition. Each reported value is the average run length taken over the repetitions which have already declared prior to time 20000. For reference, $\mathbb{E}(X|X < 20000) \approx 4626.9$ when $X \sim \text{Exp}(1/5000)$

	$p = 100$	$p = 1000$
$\beta = 2$	4606.2	4480.8
$\beta = 1/2$	5291.5	4383.6

and for each $n \in [\gamma]$, compute the diagonal statistic $S_n^{\text{diag},(r)}$ on the r th sample. Now compute $V^{(r)} := \max_{1 \leq n \leq \gamma} S_n^{\text{diag},(r)}$, and take \tilde{T}^{diag} to be the $(1/e)$ th quantile of $\{V^{(r)} : r \in [B]\}$. The rationale for the final step here is that if $\mathbb{P}_0(V^{(1)} < \tilde{T}^{\text{diag}}) = 1/e$, then $\mathbb{P}_0(\tilde{N}^{\text{diag}} > \gamma) = 1/e$, where $\tilde{N}^{\text{diag}} := \min\{n : S_n^{\text{diag}} \geq \tilde{T}^{\text{diag}}\}$. Thus, under an exponential distribution for \tilde{N}^{diag} , we have that \tilde{N}^{diag} has individual patience γ .

Having determined appropriate thresholds \tilde{T}^{diag} , $\tilde{T}^{\text{off,d}}$ and $\tilde{T}^{\text{off,s}}$, we can then use similar ideas to set a suitable combined threshold T^{comb} . In particular, we also argue that $N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}}$ has an approximate exponential distribution; see Figure 3 for supporting evidence. We therefore proceed as follows: for $r \in [B]$, simulate $\tilde{X}_1^{(r)}, \dots, \tilde{X}_\gamma^{(r)} \stackrel{\text{iid}}{\sim} \mathcal{N}_p(0, I_p)$ and use this new data to compute $\tilde{S}_n^{\text{diag},(r)} := S_n^{\text{diag},(r)} / \tilde{T}^{\text{diag}}$, $\tilde{S}_n^{\text{off,d},(r)} := S_n^{\text{off,d},(r)} / \tilde{T}^{\text{off,d}}$ and $\tilde{S}_n^{\text{off,s},(r)} := S_n^{\text{off,s},(r)} / \tilde{T}^{\text{off,s}}$ for each $n \in [\gamma]$, and set $W^{(r)} := \max\{\tilde{S}_n^{\text{diag},(r)} \vee \tilde{S}_n^{\text{off,d},(r)} \vee \tilde{S}_n^{\text{off,s},(r)} : n \in [\gamma]\}$ on the r th sample. Now take T^{comb} to be the $(1/e)$ th quantile of $\{W^{(r)} : r \in [B]\}$. Similar to before, our reasoning here is that if $\mathbb{P}_0(W^{(1)} < T^{\text{comb}}) = 1/e$, then $N^{\text{diag}} := \min\{n : S_n^{\text{diag}} \geq \tilde{T}^{\text{diag}} T^{\text{comb}}\}$, $N^{\text{off,d}} := \min\{n : S_n^{\text{off,d}} \geq \tilde{T}^{\text{off,d}} T^{\text{comb}}\}$ and $N^{\text{off,s}} := \min\{n : S_n^{\text{off,s}} \geq \tilde{T}^{\text{off,s}} T^{\text{comb}}\}$ satisfy

$$\mathbb{P}_0(N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}} > \gamma) = 1/e.$$

TABLE 2 Estimated response delays over 200 repetitions for N^{diag} , $N^{\text{off,d}}$ and $N^{\text{off,s}}$ and the response delay of the combined declaration time N for ocd , with the percentages of repetitions on which each statistics triggers the declaration first (or equal first) shown in parentheses. The quickest response in each setting is given in bold. Other parameters: $p = 100$, $\gamma = 5000$, $z = 0$ and $\theta = \vartheta U$, where the distribution of U is described in Section 4.2

s	ϑ	$\beta = \vartheta$			
		N^{diag}	$N^{\text{off,d}}$	$N^{\text{off,s}}$	N
1	2	11.5 (83.5)	19.4 (1.5)	13.0 (35)	11.2
1	1	40.6 (79.5)	74.4 (1.5)	47.4 (19)	39.1
1	0.5	136.3 (82)	305.2 (1)	169.2 (17)	129.7
1	0.25	455.4 (83)	1124.5 (1)	635.0 (16)	433.6
10	2	20.1 (9.5)	19.2 (9.5)	14.7 (88)	14.3
10	1	69.7 (15.5)	72.6 (12)	52.4 (73.5)	50.4
10	0.5	240.4 (29.5)	308.0 (3)	207.7 (68)	197.1
10	0.25	723.3 (56.5)	1124.3 (6)	760.7 (37.5)	648.4
100	2	53.3 (0.5)	19.7 (92)	27.4 (10)	19.5
100	1	169.9 (2)	75.2 (85)	94.9 (14.5)	73.1
100	0.5	544.1 (9)	300.6 (75.5)	345.1 (15.5)	278.9
100	0.25	1493.6 (28.5)	1206.0 (51.5)	1420.2 (20)	1065.4

Thus, under an exponential distribution for $N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}}$, it again has the desired nominal patience. Our practical thresholds, therefore, are $T^{\text{diag}} = \tilde{T}^{\text{diag}} T^{\text{comb}}$, $T^{\text{off,d}} = \tilde{T}^{\text{off,d}} T^{\text{comb}}$ and $T^{\text{off,s}} = \tilde{T}^{\text{off,s}} T^{\text{comb}}$ for S^{diag} , $S^{\text{off,d}}$ and $S^{\text{off,s}}$, respectively. Table 1 confirms that, with these choices of Monte Carlo thresholds, the patience of the adaptive ocd algorithm remains at approximately the desired nominal level.

4.2 | Numerical performance of ocd

In this section, we study the empirical performance of ocd . As shown in Figure 1, under the alternative, all three statistics S^{diag} , $S^{\text{off,d}}$ and $S^{\text{off,s}}$ in ocd can be the first to trigger a declaration that a mean change has occurred. We thus examine different settings under which each of these three statistics can, respectively, be the quickest to react to a change. Our simulations were run for $p = 100$, $s \in \{1, \lfloor p^{1/2} \rfloor, p\}$, $z \in \{0, 1000\}$, $\gamma = 5000$, $\vartheta \in \{2, 1, 0.5, 0.25\}$, $\beta \in \{\vartheta, 4\vartheta, \vartheta/4\}$. In all cases, θ was generated as ϑU , where U is uniformly distributed on the union of all s -sparse unit spheres in \mathbb{R}^p . By this, we mean that we first generate a uniformly random subset S of $[p]$ of cardinality s , then set $U := Z/\|Z\|_2$, where $Z = (Z^1, \dots, Z^p)^\top$ has independent components satisfying $Z^j \sim \mathcal{N}(0, 1)\mathbb{1}_{\{j \in S\}}$. Instead of terminating the ocd procedure once one of the three statistics declares a change (as we would in practice), we run the procedure until all three statistics have exceeded their respective thresholds. Tables 2 and 3 summarise the performance of the three statistics for $z = 0$. Simulation results for $z = 1000$ were similar, and are therefore not included here.

We first discuss the case when β is correctly specified (Table 2). When the sparsity s is small or moderate and ϑ is small, the diagonal statistic S^{diag} is likely to be the first to declare a change. The response delay of S^{diag} increases with s , which means that the off-diagonal sparse statistic $S^{\text{off,s}}$ typically reacts quickest to a change when the s is moderate to large and ϑ is not too small.

TABLE 3 Estimated response delays over 200 repetitions for N^{diag} , $N^{\text{off,d}}$ and $N^{\text{off,s}}$ and the response delay of the combined declaration time N for ocd . Settings where β is both over- and under-specified are given. The quickest response in each setting is given in bold. Other parameters: $p = 100$, $\gamma = 5000$, $z = 0$ and $\theta = \vartheta U$, where the distribution of U is described in Section 4.2

s	ϑ	$\beta = 4\vartheta$				$\beta = \vartheta/4$			
		N^{diag}	$N^{\text{off,d}}$	$N^{\text{off,s}}$	N	N^{diag}	$N^{\text{off,d}}$	$N^{\text{off,s}}$	N
1	2	7.7	19.5	12.8	7.6	30.3	19.5	12.6	12.6
1	1	27.8	77.7	48.3	27.6	98.3	73.7	45.2	45.1
1	0.5	92.9	288.9	162.0	92.3	304.8	304.9	171.8	171.1
1	0.25	351.7	1148.7	657.2	342.8	746.7	1158.1	614.0	586.7
10	2	16.7	19.0	14.9	13.7	50.0	20.4	15.1	15.0
10	1	57.6	72.9	51.2	46.5	161.9	76.5	54.7	54.5
10	0.5	228.3	286.4	201.0	180.5	509.0	314.7	203.6	201.8
10	0.25	739.3	1175.1	787.9	645.1	1208.2	1189.6	725.1	715.9
100	2	59.2	18.9	25.3	18.7	110.8	21.2	27.2	20.5
100	1	213.9	73.0	92.4	71.0	347.4	76.8	95.5	74.2
100	0.5	696.5	307.0	385.0	284.8	1029.0	310.2	352.5	289.3
100	0.25	1811.5	1218.1	1327.4	967.1	2149.9	1091.9	1175.9	957.8

On the other hand, the stopping time $N^{\text{off,d}}$, which is driven by the off-diagonal dense statistic, is not significantly affected by s (in agreement with our average-case bound in Theorem 2), and is usually the dominant statistic when the signal is dense. A further observation is that the three individual response delays, as well as the combined response delay, are all approximately proportional to ϑ^{-2} , a phenomenon which is supported by Theorem 6.

Table 3 presents corresponding results when β is both over- and under-specified. We note that both $N^{\text{off,d}}$ and $N^{\text{off,s}}$ are almost unaffected by either type of misspecification. For N^{diag} , a mild over-misspecification of β helps it to react faster, while an under-misspecification causes it to have increased response delay. However, since we can also observe that N^{diag} rarely declares first by a large margin, the performance of ocd is highly robust to misspecification of β , especially when s is not too small.

4.3 | Comparison with other methods

We now compare our adaptive ocd algorithm with other online changepoint detection algorithms proposed in the literature, namely those of Mei (2010), Xie and Siegmund (2013) and Chan (2017). Since we were unable to find publicly available implementations of any of these algorithms, we briefly describe below their methodology and the small adaptations that we made in order to allow them to be used in our settings.

Mei (2010) assumes knowledge of θ , and, on observing each new data point, aggregates likelihood ratio tests in each coordinate of the null $\mathcal{N}(0, 1)$ against an alternative of $\mathcal{N}(\theta^j, 1)$ in the j th coordinate. More precisely, in the notation of Equation (2), the algorithm declares a change when either $\sum_{j \in [p]} R_{n, \theta^j}^j$ or $\max_{j \in [p]} R_{n, \theta^j}^j$ exceeds given thresholds. In our setting where we do not know θ and only assume that $\|\theta\|_2 \geq \beta$, we replace $\sum_{j \in [p]} R_{n, \theta^j}^j$ and $\max_{j \in [p]} R_{n, \theta^j}^j$ with

$$\max \left\{ \sum_{j=1}^p R_{n,\beta/\sqrt{p}}^j, \sum_{j=1}^p R_{n,-\beta/\sqrt{p}}^j \right\} \quad \text{and} \quad \max \left\{ \max_{j \in [p]} R_{n,\beta/\sqrt{p}}^j, \max_{j \in [p]} R_{n,-\beta/\sqrt{p}}^j \right\},$$

respectively.

The algorithms of Xie and Siegmund (2013) and Chan (2017) have a similar flavour. The idea is to test the null $\mathcal{N}_p(0, I_p)$ distribution against an alternative where the j th coordinate has a $(1 - p_0)\mathcal{N}(0, 1) + p_0\mathcal{N}(\mu^j, 1)$ mixture distribution, for some known $p_0 \in [0, 1]$ and unknown $\mu^j \in \mathbb{R}$. After specifying a window size w , both algorithms search for the strongest evidence against the null from the past $r \in [w]$ observations. Specifically, writing $Z_{n,r}^j := r^{-1/2} \sum_{i=n-r+1}^n X_i^j$ for $n \in \mathbb{N}$, $r \in [n]$ and $j \in [p]$, the test statistics are of the form

$$S_{\text{XS,C}}^+(p_0, \lambda, \kappa, w) := \max_{r \in [w \wedge n]} \sum_{j=1}^p \log \left(1 - p_0 + \lambda p_0 e^{(Z_{n,r}^j \vee 0)^2 / \kappa} \right),$$

where Xie and Siegmund (2013) take $(\lambda, \kappa, w) = (1, 2, 200)$ and Chan (2017) takes $(\lambda, \kappa, w) = (2\sqrt{2} - 2, 4, 200)$. Since such a test statistic is only effective when $\sum_{j \in [p]} (\mu^j \vee 0)^2$ is large, we considered statistics of the form $S_{\text{XS,C}}^+(p_0, \lambda, \kappa, w) \vee S_{\text{XS,C}}^-(p_0, \lambda, \kappa, w)$, where $S_{\text{XS,C}}^-(p_0, \lambda, \kappa, w)$ replaces the exponent $Z_{n,r}^j \vee 0$ with $Z_{n,r}^j \wedge 0$. An adaptive choice of p_0 is not provided by the authors, but the choices $p_0 \in \{1/\sqrt{p}, 0.1, 1\}$ have been considered; we found the choice $p_0 = 1/\sqrt{p}$ to be the most competitive overall, so for simplicity of exposition, present only that choice in our results.

For each of the Mei (2010), Xie and Siegmund (2013) and Chan (2017) algorithms, we determined appropriate thresholds using Monte Carlo simulation, as suggested by the authors, and in a similar fashion to the way in which we set the `ocd` thresholds as described in Section 4.1. This guarantees that the algorithms have approximately the nominal patience, and so allows us to compare the methods by means of the response delay.

Table 4 displays the response delays for the `ocd` algorithm, as well as the alternative methods described above, for $p \in \{100, 2000\}$, $s \in \{5, \lfloor \sqrt{p} \rfloor, p\}$ and $\vartheta \in \{2, 1, 0.5, 0.25\}$. In fact, we also ran simulations for $p = 1000$, $s \in \{1, p/2\}$ and $\vartheta = 0.125$, but the results are qualitatively similar and are therefore omitted. Overall, the results reveal that `ocd` performs very well in comparison with existing methods, across a wide range of scenarios; in several cases it is by far the most responsive procedure, and in none of the settings considered is it outperformed by much. The Xie and Siegmund (2013) and Chan (2017) algorithms perform similarly to each other, and in most settings are both more competitive than the Mei (2010) method described above. We note that the performance of the Xie and Siegmund (2013) and Chan (2017) algorithms is relatively better when the signal-to-noise ratio ϑ is high; in these scenarios, the default window size $w = 200$ is large enough that sufficient evidence against the null can typically be accumulated within the moving window. For lower signal-to-noise ratios, this ceases to be the case, and from time $z+w$ onwards, the test statistic has the same marginal distribution (with no positive drift). This explains the relative deterioration in performance for those algorithms in the harder settings considered. As mentioned in the introduction, if the change in mean were known to be small, then the window size could be increased to compensate, but at additional computational expense; a further advantage of `ocd`, then, is that the computational time only depends on p (and not on β or other problem parameters).

TABLE 4 Estimated response delay for oed , as well as the algorithms of Mei (2010) (Mei), Xie and Siegmund (2013) (XS) and Chan (2017) (Chan) over 200 repetitions, with $z = 0$, $\gamma = 5000$ and θ generated as described in Section 4.2. The smallest response delay is given in bold

p	s	θ	oed	Mei	XS	Chan
100	5	2	13.7	36.3	13.1	11.9
100	5	1	46.9	125.9	47.3	42.0
100	5	0.5	174.8	383.1	194.3	163.7
100	5	0.25	583.5	970.4	2147	1888.8
100	10	2	14.9	44.1	15.2	14.5
100	10	1	53.8	150.1	52.9	51.5
100	10	0.5	194.4	458.2	255.8	245.6
100	10	0.25	629.7	1171.3	2730.7	2484.9
100	100	2	19.4	72.7	23.6	27.5
100	100	1	74.4	268.3	89.6	102.1
100	100	0.5	287.9	834.9	526.8	756.0
100	100	0.25	1005.8	1912.9	3598.3	3406.6
2000	5	2	19.0	130.5	20.8	15.6
2000	5	1	67.3	316.7	79.5	59.5
2000	5	0.5	247.3	680.2	607.7	285.0
2000	5	0.25	851.3	1384.8	4459.2	3856.9
2000	44	2	37.5	247.7	40.2	37.7
2000	44	1	136.0	596.1	149.1	145.0
2000	44	0.5	479.1	1270.8	2945.5	2751.4
2000	44	0.25	1584.2	2428.8	4457.8	5049.7
2000	2000	2	97.1	949.9	103.2	136.7
2000	2000	1	360.7	2126.5	1020.0	2074.7
2000	2000	0.5	1296.0	3428.1	4669.3	4672.7
2000	2000	0.25	3436.7	4140.4	5063.7	5233.5

4.4 | Real data example

We consider a seismic signal detection problem, using a data set from the High Resolution Seismic Network, operated by the Berkeley Seismological Laboratory. Ground motion sensor measurements were recorded using geophones at a frequency of 250 Hz in three mutually perpendicular directions, at 13 stations near Parkfield, California for a total of 740 s from 2 AM on 23 December 2004. This data set was also studied by Xie et al. (2019), and was obtained from <http://service.ncedc.org/fdsnws/dataselect/1/>. To begin, we removed the linear trend in each coordinate and applied a 2–16 Hz bandpass filter to the data using the GISMO toolbox³; these are standard pre-processing steps in the seismology literature (e.g. Caudron et al., 2018; Xie et al., 2019). In order to reduce the effects of temporal dependence, we computed a root mean square

³Available at: <http://geoscience-community-codes.github.io/GISMO/>

amplitude envelope, downsampled to 16 Hz, and then extracted the residuals from the fit of an autoregressive model of order 1. The processed data are available as a built-in data set in the `ocd` R package. The first 4 min of the series were used to estimate the baseline mean and variance for each sensor, and we plot the standardised data from 2:04 AM onwards in Figure 4. When applying our `ocd` algorithm to this data, we specified the patience level to be $\gamma = 1.35 \times 10^6$, corresponding to a patience of one day, and $\beta = 150$. The `ocd` algorithm declared a change at 02:10:03.84, and was triggered by $S^{\text{off},d}$. According to the Northern California Earthquake Catalog⁴, an earthquake of magnitude 1.47 Md hit near Atascadero, California (50 km away from Parkfield) at 02:09:54.01, so the delay was 9.8 s. It is known⁵ that P waves, which are the primary preliminary wave and arrive first after an earthquake, travel at up to 6 km/s in the Earth's crust, which is consistent with this delay.

5 | PROOFS OF MAIN RESULTS

5.1 | Proofs from Section 3.1

Proof of Theorem 1 Define $m := \lfloor 2\gamma \rfloor$. It suffices to prove that (a) $\mathbb{P}_0(N^{\text{off}} \leq m) \leq 1/4$ and (b) $\mathbb{P}_0(N^{\text{diag}} \leq m) \leq 1/4$, since then we have

$$\begin{aligned} \mathbb{E}_0(N) &= \mathbb{E}_0(N^{\text{off}} \wedge N^{\text{diag}}) \geq 2\gamma \mathbb{P}_0(N^{\text{off}} \wedge N^{\text{diag}} > 2\gamma) \\ &\geq 2\gamma \{1 - \mathbb{P}_0(N^{\text{off}} \leq m) - \mathbb{P}_0(N^{\text{diag}} \leq m)\} \geq \gamma. \end{aligned}$$

We prove the two claims below.

(a) By Equation (5) and a union bound, we have

$$\mathbb{P}_0(N^{\text{off}} \leq m) \leq \sum_{\substack{n \in [m], j \in [p] \\ b \in \mathcal{B}}} \mathbb{P}_0(Q_{n,b}^j \geq T^{\text{off}}) = \sum_{\substack{n \in [m], j \in [p] \\ b \in \mathcal{B}}} \mathbb{E}_0 \left[\mathbb{P}_0(Q_{n,b}^j \geq T^{\text{off}} \mid \tau_{n,b}^j) \right]. \quad (13)$$

Recall that under the null, $\Lambda_b^{k,j} \mid \tau_b^j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau_b^j)$ for all $b \in \mathcal{B}, j \in [p]$ and $k \in [p] \setminus \{j\}$, which implies that $Q_b^j \mid \tau_b^j \sim \chi_{p-1}^2 \mathbb{1}_{\{\tau_b^j > 0\}}$. Thus, we have by Laurent and Massart (2000, Lemma 1) that for all $n \in [m]$, $j \in [p]$ and $b \in \mathcal{B}$,

$$\mathbb{P}_0(Q_{n,b}^j \geq T^{\text{off}} \mid \tau_{n,b}^j) \leq e^{-\tilde{T}^{\text{off}}/2}. \quad (14)$$

Combining Equations (13) and (14), we have

$$\mathbb{P}_0(N^{\text{off}} \leq m) \leq |\mathcal{B}| m p e^{-\tilde{T}^{\text{off}}/2} \leq 1/4. \quad (15)$$

(b) For $j \in [p]$ and $b \in \mathcal{B} \cup \mathcal{B}_0$, denote $N_b^j := \inf\{n : R_{n,b}^j \geq T^{\text{diag}}\}$, where $R_{n,b}^j$ is defined by

⁴Available at: <http://www.ncedc.org/ncedc/catalog-search.html>.

⁵One source for this information is <https://www.usgs.gov/natural-hazards/earthquake-hazards/science/seismographs-keeping-track-earthquakes>.

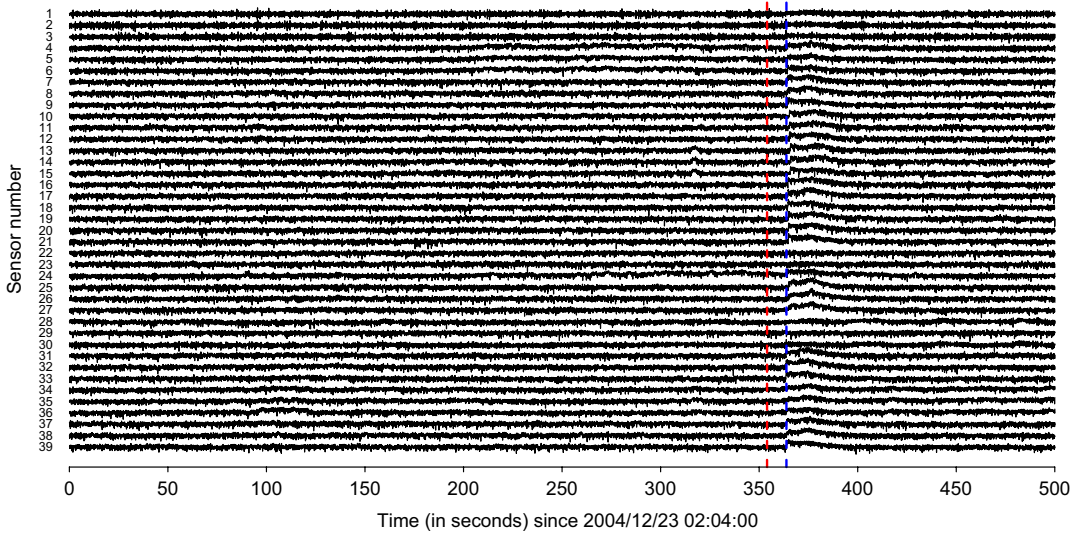


FIGURE 4 Standardised, pre-processed earthquake data from 39 sensors. The time of the 1.47 Md earthquake is given by the vertical red dashed line, while time of ocd declaration of change is given as a blue dashed line

Equation (2). By Lemma 1, we have that $R_{n,b}^j = \{R_{n-1,b}^j + b(X_n^j - b/2)\} \vee 0$, and that this process is always non-negative. Then $N^{\text{diag}} = \min \left\{ N_b^j : j \in [p], b \in \mathcal{B} \cup \mathcal{B}_0 \right\}$.

Now, fix some $j \in [p]$ and $b \in \mathcal{B} \cup \mathcal{B}_0$. Define $U_0 := 0$ and $U_h := \inf \left\{ n > U_{h-1} : R_{n,b}^j \notin (0, T^{\text{diag}}) \right\}$ for $h \in \mathbb{N}$, and let $H := \inf \left\{ h : R_{U_h,b}^j \geq T^{\text{diag}} \right\}$. Then

$$N_b^j = U_H \geq H.$$

To study the distribution of H , consider the one-sided sequential probability ratio test of $H_{0,Z} : Z_1, Z_2, \dots \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ against $H_{1,Z} : Z_1, Z_2, \dots \stackrel{\text{iid}}{\sim} \mathcal{N}(b, 1)$ with log-boundaries T^{diag} and $-\infty$. The associated stopping time for this test is

$$N_{\text{os}} := \inf \left\{ n \in \mathbb{N} : b \sum_{t=1}^n (Z_t - b/2) \geq T^{\text{diag}} \right\}.$$

Since $(R_{n,b}^j)_n$ is a Markov process that renews itself every time it hits 0, H follows a geometric distribution with success probability

$$\mathbb{P}_0 \left(R_{U_1,b}^j \geq T^{\text{diag}} \right) \leq \mathbb{P}_{H_{0,Z}}(N_{\text{os}} < \infty) \leq e^{-T^{\text{diag}}},$$

where the last inequality follows from Lemma 3. Consequently,

$$\mathbb{P}_0 \left(N_b^j \leq m \right) \leq \mathbb{P}_0(H \leq m) \leq 1 - \left(1 - e^{-T^{\text{diag}}} \right)^m.$$

As the above inequality holds for all $j \in [p]$ and $b \in \mathcal{B} \cup \mathcal{B}_0$, we have that

$$\begin{aligned} \mathbb{P}_0(N^{\text{diag}} > m) &= \mathbb{P}_0 \left(\bigcap_{j \in [p], b \in \mathcal{B} \cup \mathcal{B}_0} \{N_b^j > m\} \right) = \prod_{j \in [p]} \left\{ 1 - \mathbb{P}_0 \left(\bigcup_{b \in \mathcal{B} \cup \mathcal{B}_0} \{N_b^j \leq m\} \right) \right\} \\ &\geq \left[1 - |\mathcal{B} \cup \mathcal{B}_0| \left\{ 1 - \left(1 - e^{-T^{\text{diag}}} \right)^m \right\} \right]^p \geq 1 - mp|\mathcal{B} \cup \mathcal{B}_0|e^{-T^{\text{diag}}} \geq 3/4, \end{aligned} \quad (16)$$

as desired, where in the penultimate inequality, we twice used the fact that $(1-x)^\alpha \geq 1-\alpha x$ for all $\alpha \geq 1$ and $x \in [0,1]$.

The proof of Theorem 2 is quite involved. We first define some relevant quantities, and then state and prove some preliminary results. For $\theta \in \mathbb{R}^p$ with effective sparsity $s(\theta)$, there is at most one coordinate in θ of magnitude larger than $\vartheta/\sqrt{2}$, so there exists $b_* \in \left\{ \beta/\sqrt{s(\theta)\log_2(2p)}, -\beta/\sqrt{s(\theta)\log_2(2p)} \right\} \subseteq \mathcal{B}$ such that

$$\mathcal{J} := \left\{ j \in [p] : \theta^j/b_* \geq 1 \text{ and } |\theta^j| \leq \vartheta/\sqrt{2} \right\} \quad (17)$$

has cardinality at least $s(\theta)/2$ (note that the condition $\theta^j/b_* \geq 1$ above ensures that $\{\theta^j : j \in \mathcal{J}\}$ all have the same sign as b_*). Both b_* and \mathcal{J} can be chosen as functions of θ . Now, given any sequence $X_1, X_2, \dots \in \mathbb{R}^p$ and $\theta \in \mathbb{R}^p$, define for any $\alpha \in (0,1]$ the function

$$q(\alpha) = q(\alpha; X_1, \dots, X_z, \theta) := \inf \left\{ y \in \mathbb{R} : \left| \{j \in \mathcal{J} : t_{z,b_*}^j \leq y\} \right| \geq \alpha |\mathcal{J}| \right\}, \quad (18)$$

where t_{z,b_*}^j is obtained by running Algorithm 2 up to time z with $a = 0$ and $T^{\text{diag}} = T^{\text{off}} = \infty$. In other words, $q(\alpha)$ is the empirical α -quantile of the tail lengths $(t_{z,b_*}^j : j \in \mathcal{J})$ when we run the algorithm without declaring any change up to time z . Recall the definition of the function ψ in Equation (5).

Proposition 7 Assume that X_1, X_2, \dots are generated according to $\mathbb{P}_{z,\theta}$ for some z and θ such that $\|\theta\|_2 = \vartheta \geq \beta > 0$ and that θ has an effective sparsity of $s := s(\theta) \geq 2$. Then the output N from Algorithm 2, with input $(X_t)_{t \in \mathbb{N}}$, $\beta > 0$, $a = 0$, $T^{\text{diag}} \geq 1$ and $T^{\text{off}} = \psi(\tilde{T}^{\text{off}})$ for $\tilde{T}^{\text{off}} \geq \log(ep)$, satisfies

$$\mathbb{E}_{z,\theta} \left\{ (N - z) \vee 0 \mid X_1, \dots, X_z \right\} \leq \frac{396\tilde{T}^{\text{off}} + 65\sqrt{p\tilde{T}^{\text{off}}}}{\vartheta^2} + \frac{24\log_2(2p)}{\alpha\beta^2} + 3q(\alpha) + 2, \quad (19)$$

for any $\alpha \in (0,1]$.

Proof Since the bound in Equation (19) is positive, we may, throughout the proof and for arbitrary $z \in \mathbb{N}$, restrict attention to realisations $X_1 = x_1, \dots, X_z = x_z$ for which we have not declared a change by time z . In other words, we have $N > z$. This restriction also ensures that $q(\alpha)$ defined in Equation (18) is now indeed the empirical α -quantile of the tail lengths $(t_{z,b_*}^j : j \in \mathcal{J})$ at the changepoint. Denote $\mathcal{J}_\alpha := \left\{ j \in \mathcal{J} : t_{z,b_*}^j \leq q(\alpha) \right\}$. Then we have $|\mathcal{J}_\alpha| \geq \alpha |\mathcal{J}| \geq \alpha s/2$.

We now fix some

$$r \geq \left\lceil \frac{12 \left(\tilde{T}^{\text{off}} + \sqrt{2(p-1)\tilde{T}^{\text{off}}} \right)}{\vartheta^2} \vee 3q(\alpha) \right\rceil + 2 =: r_0. \quad (20)$$

Note that $r_0 > 3q(\alpha) \geq 3t_{z,b_*}^j$ for all $j \in \mathcal{J}_\alpha$. For $j \in \mathcal{J}_\alpha$, we define the event

$$\Omega_r^j := \left\{ t_{z+[r],b_*}^j > 2 \lfloor r \rfloor / 3 \right\}.$$

Under $\mathbb{P}_{z,\theta}$, conditional on $X_1 = x_1, \dots, X_z = x_z$, we know that $X_{z+1}, X_{z+2}, \dots \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\theta, I_p)$. Hence, by using Lemma 2 and applying Lemma 6(b) to $t_{z+[r],b_*}^j \wedge \lfloor r \rfloor$ for $j \in \mathcal{J}_\alpha$, we obtain

$$\mathbb{P}_{z,\theta} \left(\bigcap_{j \in \mathcal{J}_\alpha} (\Omega_r^j)^c \middle| X_1 = x_1, \dots, X_z = x_z \right) \leq \exp \left\{ -|\mathcal{J}_\alpha| b_*^2 \lfloor r \rfloor / 12 \right\} \leq \exp \left\{ -\alpha s b_*^2 \lfloor r \rfloor / 24 \right\}. \quad (21)$$

We now work on the event Ω_r^j for some $j \in \mathcal{J}_\alpha$. We note that Equation (20) guarantees that $r \geq 2$, and thus $t_{z+[r],b_*}^j \geq \lfloor r \rfloor / 3 \geq 2$. Then, by Lemma 9 and the fact that $r_0 > 3t_{z,b_*}^j$, we have that

$$\frac{\lfloor r \rfloor}{3} < \frac{t_{z+[r],b_*}^j}{2} \leq \tau_{z+[r],b_*}^j \leq \frac{3t_{z+[r],b_*}^j}{4} \leq \frac{3(t_{z,b_*}^j + r)}{4} < r.$$

Hence we conclude that on the event Ω_r^j ,

$$2/3 \leq \lfloor r \rfloor / 3 < \tau_{z+[r],b_*}^j \leq \lfloor r \rfloor. \quad (22)$$

Recall that $\Lambda_{z+[r],b_*}^{j,j} \in \mathbb{R}^p$ records the tail CUSUM statistics with tail length $\tau_{z+[r],b_*}^j$. We observe by Equation (22) that on Ω_r^j , only post-change observations are included in $\Lambda_{z+[r],b_*}^{j,j}$. Hence we have that on the event Ω_r^j ,

$$\Lambda_{z+[r],b_*}^{k,j} \left| \left\{ \tau_{z+[r],b_*}^j, X_1 = x_1, \dots, X_z = x_z \right\} \right. \stackrel{\text{ind}}{\sim} \mathcal{N} \left(\theta^k \tau_{z+[r],b_*}^j, \tau_{z+[r],b_*}^j \right) \quad (23)$$

for $k \in [p] \setminus \{j\}$. Therefore, on the event Ω_r^j and conditional on $\tau_{z+[r],b_*}^j, X_1 = x_1, \dots, X_z = x_z$, the random variable $\frac{\|\Lambda_b^{-jj}\|_2^2}{\tau_{z+[r],b_*}^j \vee 1} = \frac{\|\Lambda_b^{-jj}\|_2^2}{\tau_{z+[r],b_*}^j}$ follows a non-central chi-squared distribution with $p-1$ degrees of freedom and noncentrality parameter $\|\theta^{-j}\|_2^2 \tau_{z+[r],b_*}^j$. Since $j \in \mathcal{J}$ and $s \geq 2$, we observe, by Equations (17) and (22) that $\|\theta^{-j}\|_2^2 \tau_{z+[r],b_*}^j \geq \vartheta^2 \lfloor r \rfloor / 6$ on Ω_r^j . Write

$$E_r^j := \left\{ \frac{\|\Lambda_{z+[r],b_*}^{-jj}\|_2^2}{\tau_{z+[r],b_*}^j \vee 1} < T^{\text{off}} \right\}.$$

Then by Birgé (2001, Lemma 8.1), we have

$$\mathbb{P}_{z,\theta} \left(E_r^j \cap \Omega_r^j \mid \tau_{z+\lfloor r \rfloor, b_*}^j, X_1 = x_1, \dots, X_z = x_z \right) \leq \exp \left\{ - \frac{\left(\vartheta^2 \lfloor r \rfloor / 6 - \tilde{T}^{\text{off}} - \sqrt{2(p-1)\tilde{T}^{\text{off}}} \right)^2}{4(p-1 + \vartheta^2 \lfloor r \rfloor / 3)} \right\}. \quad (24)$$

Combining Equations (21) and (24), we deduce that

$$\begin{aligned} \mathbb{P}_{z,\theta} (N > z + r \mid X_1 = x_1, \dots, X_z = x_z) &\leq \mathbb{P}_{z,\theta} (N > z + \lfloor r \rfloor \mid X_1 = x_1, \dots, X_z = x_z) \\ &\leq \mathbb{P}_{z,\theta} \left(\bigcap_{J \in J_\alpha} (\Omega_r^j)^c \mid X_1 = x_1, \dots, X_z = x_z \right) + \sum_{J \in J_\alpha} \mathbb{P}_{z,\theta} \left(E_r^j \cap \Omega_r^j \mid X_1 = x_1, \dots, X_z = x_z \right) \\ &\leq \exp \left\{ - \frac{\alpha s b_*^2 (r-1)}{24} \right\} + p \exp \left\{ - \frac{\left(\vartheta^2 (r-1) / 6 - \tilde{T}^{\text{off}} - \sqrt{2(p-1)\tilde{T}^{\text{off}}} \right)^2}{4(p-1 + \vartheta^2 (r-1) / 3)} \right\} \\ &\leq \exp \left\{ - \frac{\alpha s b_*^2 (r-1)}{24} \right\} + p \exp \left\{ - \frac{\vartheta^4 (r-1)^2}{576 (p-1 + \vartheta^2 (r-1) / 3)} \right\}, \end{aligned}$$

where the last inequality uses Equation (20). Therefore, we have

$$\begin{aligned} \mathbb{E}_{z,\theta} \{ (N-z) \vee 0 \mid X_1 = x_1, \dots, X_z = x_z \} &= \int_0^\infty \mathbb{P}_{z,\theta} (N > z + u \mid X_1 = x_1, \dots, X_z = x_z) du \\ &\leq r_0 + \int_{r_0-1}^\infty \left[\exp \left\{ - \frac{\alpha s b_*^2 u}{24} \right\} + p \exp \left\{ - \frac{\vartheta^4 u^2}{576 (p-1 + \vartheta^2 u / 3)} \right\} \right] \wedge 1 du \\ &\leq r_0 + \frac{24}{\alpha s b_*^2} + \int_0^\infty \left(p e^{-\vartheta^2 u / 384} \right) \wedge 1 du + \int_0^\infty \left(p e^{-\frac{\vartheta^4 u^2}{1152(p-1)}} \right) \wedge 1 du \\ &\leq r_0 + \frac{24}{\alpha s b_*^2} + \frac{384 \log(ep)}{\vartheta^2} + \frac{24 \sqrt{2(p-1) \log p}}{\vartheta^2} + \frac{12 \sqrt{2\pi(p-1)}}{\vartheta^2} \\ &\leq r_0 + \frac{24}{\alpha s b_*^2} + \frac{384 \log(ep)}{\vartheta^2} + \frac{48 \sqrt{(p-1) \log(ep)}}{\vartheta^2}, \end{aligned}$$

where the penultimate inequality follows from the fact that $1 - \Phi(x) \leq \frac{1}{2} e^{-x^2/2}$ for $x \geq 0$. The desired bound Equation (19) follows by substituting in the expressions for r_0 and b_* .

The following two propositions control the residual tail length quantile term $q(\alpha)$ in Equation (19) in the worst-case and average-case scenarios, respectively.

Proposition 8 *Let $X_1, X_2, \dots, z, \theta, s, a, p$ and N be defined as in Proposition 7. On the event $\{N > z\}$, we have*

$$q(1; X_1, \dots, X_z, \theta) \leq \frac{8T^{\text{diag}} s \log_2(2p)}{\beta^2}.$$

Proof We will show the stronger result that on the event $\{N > z\}$, we have

$$t_{z,b}^j < \frac{8T^{\text{diag}}}{b^2}$$

for all $b \in \mathcal{B}$ and $j \in [p]$. The desired result then follows immediately by taking $b = b_*$ and restricting to the subset $\mathcal{J} \subseteq [p]$.

Fix $b \in \mathcal{B}$ and $j \in [p]$. Recall from Equation (2) and Lemma 1 the definition of $R_{n,b}^j$ and the recursive relation $R_{n,b}^j = \left\{ R_{n-1,b}^j + b(X_n^j - b/2) \right\} \vee 0$. By the update procedure for $t_{n,b}^j$ in Algorithm 2 and Lemma 2, we have

$$R_{n,b}^j \begin{cases} = 0 & \text{when } n = z - t_{z,b}^j, \\ > 0 & \text{when } z - t_{z,b}^j < n \leq z. \end{cases} \quad (25)$$

We claim that

$$R_{n,b/2}^j \geq \frac{R_{n,b}^j}{2} + \frac{b^2(n - z + t_{z,b}^j)}{8}, \quad (26)$$

for all $n \in \{z - t_{z,b}^j, \dots, z\}$. To see this, the claim is true when $n = z - t_{z,b}^j$ since the right hand side of Equation (26) is 0 by Equation (25). Now, assume Equation (26) is true for some $n = m - 1$. Then,

$$\begin{aligned} R_{m,b/2}^j &\geq R_{m-1,b/2}^j + \frac{b}{2} \left(X_m^j - \frac{b}{4} \right) \geq \frac{R_{m-1,b}^j}{2} + \frac{b^2(m-1-z+t_{z,b}^j)}{8} + \frac{b}{2} \left(X_m^j - \frac{b}{4} \right) \\ &= \frac{R_{m,b}^j}{2} + \frac{b^2(m-z+t_{z,b}^j)}{8}. \end{aligned}$$

This proves the claim by induction. In particular, on the event $\{N > z\}$, we have $T^{\text{diag}} > R_{z,b/2}^j > b^2 t_{z,b}^j / 8$ as desired.

Proposition 9 Let $X_1, X_2, \dots, X_z, \theta, s, a, p$ and N be defined as in Proposition 7. There exists a universal constant C and $\beta_0(s) > 0$, depending only on s , such that for all $\beta < \beta_0(s)$, we have

$$\mathbb{E}_{z,\theta} \left\{ q(s^{-1/2}, X_1, \dots, X_z, \theta) \right\} \leq \frac{Cs^{1/2} \log(16s^2 \beta^{-2} \log_2(2p)) \log_2(2p)}{\beta^2}.$$

Proof Recall the definition of b_* in Equation (17). We may assume without loss of generality that $b_* = \beta / \sqrt{s \log_2(2p)}$ (the case $b_* = -\beta / \sqrt{s \log_2(2p)}$ can be proved in essentially the same way). We first prove the result for sufficiently large $s > s_0$. Recall that $t_{z,b_*}^j = \arg\max_{0 \leq r \leq z} \sum_{i=z-r+1}^z (X_i^j - b_*/2)$. Define $Z_i := X_{z-i+1}$ for $i \in [z]$ and let $Z_{z+1}, Z_{z+2}, \dots \stackrel{\text{iid}}{\sim} \mathcal{N}_p(0, I_p)$ be independent from Z_1, \dots, Z_z . For each $j \in [p]$, let

$$S_r^j := \sum_{i=1}^r \left(Z_i^j - b_*/2 \right) \quad \text{and} \quad \tilde{S}_r^j := \sum_{i=1}^r Z_i^j$$

for $r \in \mathbb{N}$ and define $S_0^j := \tilde{S}_0^j := 0$. Writing $\xi_0^j := \operatorname{argmax}_{0 \leq r \leq \Delta b_*^{-2}} S_r^j$, $\xi^j := \operatorname{argmax}_{r \in \mathbb{N}_0} S_r^j$, and $\tilde{\xi}_0^j := \operatorname{argmax}_{0 \leq r \leq \Delta b_*^{-2}} \tilde{S}_r^j$, where $\Delta := 8 \log(2s)$, we note that like t_{z, b_*}^j , these three maxima are also uniquely attained almost surely (see the proof of Lemma 6). By construction, we have for each $j \in [p]$ that

$$t_{z, b_*}^j = \operatorname{argmax}_{0 \leq r \leq z} \sum_{i=z-r+1}^z (X_i^j - b_*/2) = \operatorname{argmax}_{0 \leq r \leq z} S_r^j \leq \operatorname{argmax}_{r \in \mathbb{N}_0} S_r^j = \xi^j.$$

Writing $q_\xi(\alpha) := \inf \{y : |\{j \in \mathcal{J} : \xi^j \leq y\}| \geq \alpha |\mathcal{J}|\}$ as the empirical α -quantile of $(\xi^j : j \in \mathcal{J})$, it follows that $q(\alpha) \leq q_\xi(\alpha)$ and so it suffices to control $\mathbb{E}\{q_\xi(s^{-1/2})\}$ instead of $\mathbb{E}\{q(s^{-1/2})\}$. To this end, we observe that $\{16\Delta s^{-1/2} b_*^2 < \xi^j \leq \Delta b_*^{-2}\} \subseteq \{16\Delta s^{-1/2} b_*^2 < \xi_0^j \leq \Delta b_*^{-2}\}$ and $\tilde{\xi}_0^j \geq \xi_0^j$, and thus

$$\begin{aligned} \mathbb{P}(\xi^j \leq 16\Delta s^{-1/2} b_*^{-2}) &\geq \mathbb{P}(\xi_0^j \leq 16\Delta s^{-1/2} b_*^{-2}) - \mathbb{P}(\xi^j > \Delta b_*^{-2}) \\ &\geq \mathbb{P}(\tilde{\xi}_0^j \leq 16\Delta s^{-1/2} b_*^{-2}) - \mathbb{P}(\xi^j > \Delta b_*^{-2}). \end{aligned} \quad (27)$$

For the first term on the right hand side of Equation (27), by Donsker's invariance principle and the continuity of the argmax map (see, e.g. van der Vaart & Wellner, 1996, Lemma 3.2.1 and Theorem 3.2.2), we have in the limit $\beta \searrow 0$ that $\Delta b_*^{-2} \rightarrow \infty$ and so

$$\frac{\tilde{\xi}_0^j}{\Delta b_*^{-2}} \xrightarrow[t \in [0,1]]{d} \operatorname{argmax}_t B_t,$$

where $(B_t)_{t \geq 0}$ denotes a standard Brownian motion. In particular, we can find $\beta_0(s) > 0$ depending only on s such that for $\beta \leq \beta_0(s)$ and $s > 256$, we have

$$\mathbb{P}(\tilde{\xi}_0^j \leq 16\Delta s^{-1/2} b_*^{-2}) \geq \frac{1}{2} \mathbb{P}\left(\operatorname{argmax}_{t \in [0,1]} B_t \leq 16s^{-1/2}\right) = \frac{1}{\pi} \arcsin(4s^{-1/4}) \geq \frac{4s^{-1/4}}{\pi}, \quad (28)$$

where in the second step we used the arcsine law for Brownian motion (see, e.g. Mörters & Peres, 2010, Theorem 5.26), and in the final step we used the fact that $4s^{-1/4} < 1$.

For the second term on the right-hand side of Equation (27), since $\Delta = 8 \log(2s)$, for sufficiently large $s \geq s_0$ and sufficiently small $\beta \leq \beta_0(s)$, we have by Lemma 6(d) that

$$\mathbb{P}(\xi^j > \Delta b_*^{-2}) \leq 2e^{-\Delta/8} = s^{-1}. \quad (29)$$

Substituting Equations (28) and (29) into Equation (27), we have, for all $j \in \mathcal{J}$, that

$$\mathbb{P}(\xi^j \leq 16\Delta s^{-1/2} b_*^{-2}) \geq s^{-1/4}.$$

As a result, $\left| \{j \in \mathcal{J} : \xi^j \leq 16\Delta s^{-1/2} b_*^{-2}\} \right|$ is stochastically larger than $\text{Bin}(|\mathcal{J}|, s^{-1/4})$. Thus, for $s \geq s_0$, we have,

$$\mathbb{P}_{z,\theta} \{q_\xi(s^{-1/2}) > 16\Delta s^{-1/2} b_*^{-2}\} \leq \mathbb{P} \{ \text{Bin}(|\mathcal{J}|, s^{-1/4}) \leq s^{-1/2} |\mathcal{J}| \} \leq e^{-s^{1/2}/2},$$

where we have used Hoeffding's inequality and the fact that $|\mathcal{J}| \geq s/2$ in the last step. On the other hand, for sufficiently large $s \geq s_0$ and sufficiently small $\beta \leq \beta_0(s)$, we have,

$$\begin{aligned} \mathbb{E}_{z,\theta} \left\{ q_\xi(s^{-1/2}) \mid q_\xi(s^{-1/2}) > 16\Delta s^{-1/2} b_*^{-2} \right\} &\leq \mathbb{E}_{z,\theta} \left\{ q_\xi(s^{-1/2}) \mid q_\xi(s^{-1/2}) \geq \Delta b_*^{-2} \right\} \\ &\leq \mathbb{E}_{z,\theta} \left\{ q_\xi(1) \mid q_\xi(|\mathcal{J}|^{-1}) \geq \Delta b_*^{-2} \right\} = \mathbb{E}_{z,\theta} \left\{ \max_{j \in \mathcal{J}} \xi^j \mid \min_{j \in \mathcal{J}} \xi^j \geq \Delta b_*^{-2} \right\} \leq \frac{61(\Delta + 4 \log(2/b_*))}{b_*^2}, \end{aligned}$$

where we have used Lemma 7(b) in the second inequality and Lemma 6(d) (with $\Delta/4$ taking the role of k and $b_*/2$ taking the role of b there) in the final inequality. As a result,

$$\begin{aligned} \mathbb{E}_{z,\theta} \{q(s^{-1/2})\} &\leq \mathbb{E}_{z,\theta} \{q_\xi(s^{-1/2})\} \leq 16\Delta s^{-1/2} b_*^{-2} + 61e^{-s^{1/2}/2} (\Delta + 4 \log(2/b_*)) b_*^{-2} \\ &\leq \frac{Cs^{1/2} \log(16s^2 \beta^{-2} \log_2(2p)) \log_2(2p)}{\beta^2}, \end{aligned}$$

where we have used in the final step the fact that $e^{-s^{1/2}/2} \leq s^{-1/2}/100$ for sufficiently large s . This proves the desired result for $s \geq s_0$.

Finally, for $s \leq 256$, we have by Lemma 6(c) that, for $\beta < \sqrt{s}/2$,

$$\begin{aligned} \mathbb{E}_{z,\theta} \{q(s^{-1/2})\} &\leq \mathbb{E}_{z,\theta} \left\{ \max_{j \in \mathcal{J}} \xi^j \right\} \leq \frac{32s \log(s^{3/2} \beta^{-1} \log_2^{1/2}(2p)) \log_2(2p)}{\beta^2} \\ &\leq \frac{Cs^{1/2} \log(16s^2 \beta^{-2} \log_2(2p)) \log_2(2p)}{\beta^2}, \end{aligned}$$

and the desired bound then follows.

We are now in a position to prove Theorem 2.

Proof of Theorem 2 The proof proceeds with different arguments for the case $s \geq 2$ and the case $s = 1$.

Case 1: $s \geq 2$. Combining Propositions 7 (applied with $\alpha = 1$) and 8, we have

$$\bar{\mathbb{E}}_\theta^{\text{wc}}(N) \leq \frac{396\tilde{T}^{\text{off}} + 65\sqrt{p\tilde{T}^{\text{off}}}}{g^2} + \frac{24\log_2(2p)}{\beta^2} + \frac{24T^{\text{diag}}s\log_2(2p)}{\beta^2} + 2.$$

The desired bound Equation (6) then follows by substituting in the expression for \tilde{T}^{off} . On the other hand, combining Propositions 7 (applied with $\alpha = s^{-1/2}$) and 9, we have

$$\bar{\mathbb{E}}_\theta(N) \leq \frac{396\tilde{T}^{\text{off}} + 65\sqrt{p\tilde{T}^{\text{off}}}}{g^2} + \frac{24\sqrt{s}\log_2(2p)}{\beta^2} + \frac{3Cs^{1/2} \log(16s^2 \beta^{-2} \log_2(2p)) \log_2(2p)}{\beta^2} + 2,$$

which proves Equation (7).

Case 2: $s = 1$. There exists $j_* \in [p]$ such that $|\theta^{j_*}| \geq \vartheta/\sqrt{\log_2(2p)}$, and recall from Equation (17) that $b_* := \text{sgn}(\theta^{j_*})\beta/\sqrt{\log_2(2p)} \in \mathcal{B}$. Note that $S_{n,1}^{\text{diag}} = \max_{(j,b) \in [p] \times (\mathcal{B} \cup \mathcal{B}_0)} R_{n,b}^j \geq R_{n,b_*}^{j_*}$. We define $\bar{R}_n := \sum_{i=z+1}^{z+n} b_*(X_i^{j_*} - b_*/2)$ for $n \in \mathbb{N}_0$. Since $R_{z,b_*}^{j_*} \geq 0 = \bar{R}_0$ and $R_n - R_{n-1} = b_*(X_{z+n}^{j_*} - b_*/2) \leq R_{z+n,b_*}^{j_*} - R_{z+n-1,b_*}^{j_*}$, it follows by induction that $R_{z+n,b_*}^{j_*} \geq \bar{R}_n$ for all $n \in \mathbb{N}_0$. Then, for $n \geq \lceil 4T^{\text{diag}}/(b_*\theta^{j_*}) \rceil =: n_0$, we have

$$\begin{aligned} \mathbb{P}_{z,\theta}(N > z+n | X_1=x_1, \dots, X_z=x_z) &\leq \mathbb{P}_{z,\theta}(R_{z+n,b_*}^{j_*} \leq T^{\text{diag}} | X_1=x_1, \dots, X_z=x_z) \\ &\leq \mathbb{P}_{z,\theta}(\bar{R}_n \leq T^{\text{diag}}) = \Phi\left(-\frac{b_*(\theta^{j_*} - b_*/2) - T^{\text{diag}}}{n^{1/2}b_*}\right) \\ &\leq \frac{1}{2} \exp\left\{-\frac{(b_*n\theta^{j_*}/2 - T^{\text{diag}})^2}{2nb_*^2}\right\} \leq \frac{1}{2} e^{-n(\theta^{j_*})^2/32}. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E}_{z,\theta}\{(N-z) \vee 0 | X_1=x_1, \dots, X_z=x_z\} &= \sum_{n=0}^{\infty} \mathbb{P}_{z,\theta}(N > z+n | X_1=x_1, \dots, X_z=x_z) \\ &\leq n_0 + \frac{1}{2} \sum_{n=n_0}^{\infty} e^{-n(\theta^{j_*})^2/32} \leq n_0 + \frac{1}{2} \int_0^{\infty} e^{-u(\theta^{j_*})^2/32} du \leq 1 + \frac{4T^{\text{diag}}}{b_*\theta^{j_*}} + \frac{16}{(\theta^{j_*})^2}. \end{aligned} \quad (30)$$

After substituting in the expressions for b_* , θ^{j_*} and T^{diag} , we see that

$$\bar{\mathbb{E}}_{\theta}(N) \leq \bar{\mathbb{E}}_{\theta}^{\text{wc}}(N) \leq 1 + \frac{4 \log(16p\gamma \log_2(4p)) \log_2(2p)}{\beta\vartheta} + \frac{16 \log_2(2p)}{\vartheta^2},$$

which proves both Equations (6) and (8).

5.2 | Proofs from Sections 3.2 and 3.3

Proof of Theorem 3 It suffices to only prove $\mathbb{P}_0(N^{\text{off}} \leq m) \leq 1/4$, since the remaining proof is identical to that of Theorem 1.

Since $\Lambda_b^{k,j} | \tau_b^j \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau_b^j)$ for all $b \in \mathcal{B}, j \in [p]$ and $k \in [p] \setminus \{j\}$ under the null, by the fact that $T^{\text{off}} \geq 12$ and Lemma 10, we have

$$\mathbb{P}_0(Q_{n,b}^j \geq T^{\text{off}} | \tau_{n,b}^j) \leq \mathbb{P}_0(Q_{n,b}^j \geq 6 + T^{\text{off}}/2 | \tau_{n,b}^j) \leq \exp(-T^{\text{off}}/8).$$

Hence, it follows that

$$\mathbb{P}_0(N^{\text{off}} \leq m) \leq |\mathcal{B}| m p e^{-T^{\text{off}}/8} \leq 1/4, \quad (31)$$

as desired.

Proof of Theorem 4 We note that the case $s = 1$ in the proof of Theorem 2 does not rely on the off-diagonal statistics. Hence Equation (30) is still valid here with $a = \sqrt{8 \log(p-1)}$ and the last expression in Equation (30) again proves the desired bound (Equation 9). For the case $s \geq 2$, we follow exactly the proof of Proposition 7 until Equation (23), with the only exception that we now fix, instead of Equation (20),

$$r \geq \left\{ \frac{24T^{\text{off}} \log_2(2p)}{\vartheta^2} \vee \frac{96s \log_2(2p) \log p}{\vartheta^2} \vee 3q(\alpha) \right\} + 2 =: \tilde{r}_0. \quad (32)$$

By the definition of the effective sparsity of θ , for a fixed $j \in \mathcal{J}_\alpha$,

$$\mathcal{L}^j := \left\{ j' \in [p] : |\theta^{j'}| \geq \frac{\vartheta}{\sqrt{s \log_2(2p)}} \text{ and } j' \neq j \right\}$$

has cardinality at least $s - 1$. On the event Ω_r^j , we have, by Equation (22), that for all $k \in \mathcal{L}^j$,

$$|\theta^k| \sqrt{\tau_{z+[r],b_*}^j} \geq \sqrt{\frac{\vartheta^2 \lfloor r \rfloor}{3s \log_2(2p)}} =: \tilde{a}_r.$$

We then observe, by Equation (32), that

$$\tilde{a}_r \geq \sqrt{32 \log p} > 2a. \quad (33)$$

Now, from Equation (23) we have on the event Ω_r^j that, for all $k \in \mathcal{L}^j$,

$$\mathbb{P}_{z,\theta} \left(\Omega_r^j \cap \left\{ |\Lambda_{z+[r],b_*}^{kj}| < \frac{1}{2} \tilde{a}_r \sqrt{\tau_{z+[r],b_*}^j} \right\} \mid \tau_{z+[r],b_*}^j, X_1 = x_1, \dots, X_z = x_z \right) \leq \frac{1}{2} e^{-\tilde{a}_r^2/8} =: q_r.$$

We denote

$$U^j := \left| \left\{ k \in \mathcal{L}^j : \left\{ |\Lambda_{z+[r],b_*}^{kj}| < \frac{1}{2} \tilde{a}_r \sqrt{\tau_{z+[r],b_*}^j} \right\} \right\} \right|.$$

Then, by the Chernoff–Hoeffding binomial tail bound (Hoeffding, 1963, Equation (2.1)), we have

$$\begin{aligned} \mathbb{P}_{z,\theta} \left(\Omega_r^j \cap \{U^j \geq |\mathcal{L}^j|/2\} \mid \tau_{z+[r],b_*}^j, X_1 = x_1, \dots, X_z = x_z \right) &\leq \exp \left\{ -\frac{|\mathcal{L}^j|}{2} \log \left(\frac{1}{4q_r(1-q_r)} \right) \right\} \\ &\leq \exp \left\{ |\mathcal{L}^j| \left(\frac{\log 2}{2} - \frac{\tilde{a}_r^2}{16} \right) \right\} \leq \exp \left\{ -\frac{3|\mathcal{L}^j| \tilde{a}_r^2}{64} \right\} \leq \exp \left\{ -\frac{\vartheta^2 \lfloor r \rfloor}{128 \log_2(2p)} \right\}, \end{aligned} \quad (34)$$

where the penultimate inequality follows from Equation (33). Now, on the event $\Omega_r^j \cap \{U^j < |\mathcal{L}^j|/2\}$, we have

$$\begin{aligned} \sum_{j' \in [p] : j' \neq j} \frac{\left(\Lambda_{z+[r],b_*}^{j'j} \right)^2}{\tau_{z+[r],b_*}^j \vee 1} \mathbb{1}_{\left\{ |\Lambda_{z+[r],b_*}^{j'j}| \geq a \sqrt{\tau_{z+[r],b_*}^j} \right\}} &\geq \sum_{j' \in [p] : j' \neq j} \frac{\left(\Lambda_{z+[r],b_*}^{j'j} \right)^2}{\tau_{z+[r],b_*}^j \vee 1} \mathbb{1}_{\left\{ |\Lambda_{z+[r],b_*}^{j'j}| \geq \frac{\tilde{a}_r}{2} \sqrt{\tau_{z+[r],b_*}^j} \right\}} \\ &\geq \frac{\tilde{a}_r^2}{4} \left\{ |\mathcal{L}^j| - \left(\left\lceil \frac{|\mathcal{L}^j|}{2} \right\rceil - 1 \right) \right\} = \frac{\tilde{a}_r^2}{4} \left\lceil \frac{|\mathcal{L}^j| + 1}{2} \right\rceil \\ &\geq \frac{\vartheta^2 \lfloor r \rfloor}{24 \log_2(2p)} \geq T^{\text{off}}, \end{aligned} \quad (35)$$

where the penultimate inequality uses the fact that $|\mathcal{L}^j| \geq s - 1$ and the last inequality follows from Equation (32). We now denote

$$\tilde{E}_r^j := \left\{ \sum_{j' \in [p]: j' \neq j} \frac{\left(\Lambda_{z+\lfloor r \rfloor, b_*}^{j'j} \right)^2}{\tau_{z+\lfloor r \rfloor, b_*}^j \vee 1} \mathbb{1} \left\{ |\Lambda_{z+\lfloor r \rfloor, b_*}^{j'j}| \geq a \sqrt{\tau_{z+\lfloor r \rfloor, b_*}^j} \right\} < T^{\text{off}} \right\}.$$

Combining Equations (21), (34) and (35), we deduce that

$$\begin{aligned} \mathbb{P}_{z,\theta} (N > z + r \mid X_1 = x_1, \dots, X_z = x_z) &\leq \mathbb{P}_{z,\theta} (N > z + \lfloor r \rfloor \mid X_1 = x_1, \dots, X_z = x_z) \\ &\leq \mathbb{P}_{z,\theta} \left(\bigcap_{j \in J_\alpha} (\Omega_r^j)^c \mid X_1 = x_1, \dots, X_z = x_z \right) + \sum_{j \in J_\alpha} \mathbb{P}_{z,\theta} \left(\tilde{E}_r^j \cap \Omega_r^j \mid X_1 = x_1, \dots, X_z = x_z \right) \\ &\leq \mathbb{P}_{z,\theta} \left(\bigcap_{j \in J_\alpha} (\Omega_r^j)^c \mid X_1 = x_1, \dots, X_z = x_z \right) + \\ &\quad \sum_{j \in J_\alpha} \mathbb{P}_{z,\theta} \left(\Omega_r^j \cap \{U^j \geq |\mathcal{L}^j|/2\} \mid X_1 = x_1, \dots, X_z = x_z \right) \\ &\leq \exp \left\{ -\frac{\alpha s b_*^2 (r-1)}{24} \right\} + p \exp \left\{ -\frac{\vartheta^2 (r-1)}{128 \log_2(2p)} \right\}. \end{aligned}$$

Therefore we have

$$\begin{aligned} \mathbb{E}_{z,\theta} \{ (N - z) \vee 0 \mid X_1 = x_1, \dots, X_z = x_z \} &= \int_0^\infty \mathbb{P}_{z,\theta} (N > z + u \mid X_1 = x_1, \dots, X_z = x_z) du \\ &\leq \tilde{r}_0 + \int_{\tilde{r}_0-1}^\infty \left[\exp \left\{ -\frac{\alpha s b_*^2 u}{24} \right\} + p \exp \left\{ -\frac{\vartheta^2 u}{128 \log_2(2p)} \right\} \right] \wedge 1 du \\ &\leq \tilde{r}_0 + \frac{24}{\alpha s b_*^2} + \int_0^\infty \left(p e^{-\frac{\vartheta^2 u}{128 \log_2(2p)}} \right) \wedge 1 du \leq \tilde{r}_0 + \frac{24}{\alpha s b_*^2} + \frac{128 \log_2(2p) \log(ep)}{\vartheta^2} \\ &\leq \frac{24 T^{\text{off}} \log_2(2p) + 96 s \log_2(2p) \log p}{\vartheta^2} + 3q(\alpha) + \frac{24 \log_2(2p)}{\alpha \beta^2} + \frac{128 \log_2(2p) \log(ep)}{\vartheta^2} + 2. \end{aligned}$$

Combining this with Proposition 8 (applied with $\alpha = 1$), we have, by substituting in the expression for T^{off} , that

$$\bar{\mathbb{E}}_\theta(N) \leq \bar{\mathbb{E}}_\theta^{\text{wc}}(N) \leq C \left\{ \frac{s \log(ep\gamma) \log(ep)}{\beta^2} \vee 1 \right\},$$

for some universal constant $C > 0$, as desired.

Proof of Theorem 5 Let $T^{\text{off,d}} = \psi(\tilde{T}^{\text{off,d}})$. Then, similar to Equations (15), (16) and (31), we have

$$\begin{aligned} \mathbb{P}_0(N^{\text{diag}} \leq m) &\leq mp |B \cup B_0| e^{-T^{\text{diag}}} \leq 1/6, \\ \mathbb{P}_0(N^{\text{off,d}} \leq m) &\leq mp |B| e^{-\tilde{T}^{\text{off,d}}/2} \leq 1/6, \\ \mathbb{P}_0(N^{\text{off,s}} \leq m) &\leq mp |B| e^{-T^{\text{off,s}}/8} \leq 1/6. \end{aligned}$$

Hence,

$$\begin{aligned}\mathbb{E}_0(N) &= \mathbb{E}_0(N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}}) \geq 2\gamma \mathbb{P}_0(N^{\text{diag}} \wedge N^{\text{off,d}} \wedge N^{\text{off,s}} > 2\gamma) \\ &\geq 2\gamma \{1 - \mathbb{P}_0(N^{\text{diag}} \leq m) - \mathbb{P}_0(N^{\text{off,d}} \leq m) - \mathbb{P}_0(N^{\text{off,s}} \leq m)\} \geq \gamma,\end{aligned}$$

as desired.

Proof of Theorem 6 We observe that

$$\begin{aligned}\bar{\mathbb{E}}_{\theta}^{\text{wc}}(N) &= \bar{\mathbb{E}}_{\theta}^{\text{wc}}[(N^{\text{diag}} \wedge N^{\text{off,d}}) \wedge (N^{\text{diag}} \wedge N^{\text{off,s}})] \\ &\leq \bar{\mathbb{E}}_{\theta}^{\text{wc}}[(N^{\text{diag}} \wedge N^{\text{off,d}})] \wedge \bar{\mathbb{E}}_{\theta}^{\text{wc}}[(N^{\text{diag}} \wedge N^{\text{off,s}})],\end{aligned}$$

and similarly for $\bar{\mathbb{E}}_{\theta}(N)$. The desired bounds Equations (10), (11) and (12) are therefore direct consequences of Theorems 2 and 4 (note that the different constants in the thresholds only affect the value of the universal constant).

ACKNOWLEDGEMENTS

The research of TW was supported by EPSRC grant EP/T02772X/1 and that of RJS was supported by EPSRC grants EP/P031447/1 and EP/N031938/1. Data for this study come from the High Resolution Seismic Network (HRSN), doi: 10.7932/HRSN, operated by the UC Berkeley Seismological Laboratory, which is archived at the Northern California Earthquake Data Center (NCEDC), doi: 10.7932/NCEDC. We are grateful to the anonymous reviewers for constructive feedback, which helped to improve the paper.

ORCID

Tengyao Wang  <https://orcid.org/0000-0003-2072-6645>

REFERENCES

- Baranowski, R., Chen, Y. & Fryzlewicz, P. (2019) Narrowest-Over-Threshold detection of multiple change points and change-point-like features. *Journal of the Royal Statistical Society Series B*, 81, 649–672.
- Barnard, G. A. (1959) Control charts and stochastic processes. *Journal of the Royal Statistical Society Series B*, 21, 239–271.
- Birgé, L. (2001) An alternative point of view on Lepski's method. In *State of the art in probability and statistics (Leiden, 1999)*. Beachwood, OH: IMS, pp. 113–133.
- Boucheron, S., Lugosi, G. & Massart, P. (2013) *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford: Oxford University Press.
- Caudron, C., White, R. S., Green, R. G., Woods, J., Ágústssdóttir, T., Donaldson, C., et al. (2018) Seismic amplitude ratio analysis of the 2014–15 Bárðarbunga–Holuhraun dike propagation and eruption. *Journal of Geophysical Research: Solid Earth*, 123, 264–276.
- Chan, H. P. (2017) Optimal sequential detection in multi-stream data. *Annals of Statistics*, 45, 2736–2763.
- Chen, Y., Wang, T. & Samworth, R. J. (2020) ocd: high-dimensional, multiscale online changepoint detection. Available from <https://cran.r-project.org/web/packages/ocd/index.html>.
- Chen, Y., Wang, T. & Samworth, R. J. (2021) Online supplementary material to ‘High-dimensional, multiscale online changepoint detection’. Submitted.
- Cho, H. (2016) Change-point detection in panel data via double CUSUM statistic. *Electronic Journal of Statistics*, 10, 2000–2038.
- Chu, C.-S. J., Stinchcombe, M. & White, H. (1996) Monitoring structural change. *Econometrica*, 64, 1045–1065.

- Collier, O., Comminges, L. & Tsybakov, A. B. (2017) Minimax estimation of linear and quadratic functionals on sparsity classes. *Annals of Statistics*, 45, 923–958.
- Csörgő, M. & Horváth, L. (1997) *Limit Theorems in Change-Point Analysis*. New York: John Wiley and Sons.
- Duncan, A. J. (1952) *Quality Control and Industrial Statistics*. Chicago: Richard D. Irwin Professional Publishing Inc.
- Enikeeva, F. & Harchaoui, Z. (2019) High-dimensional change-point detection under sparse alternatives. *Annals of Statistics*, 47, 2051–2079.
- Fearnhead, P. & Liu, Z. (2007) On-line inference for multiple changepoint problems. *Journal of the Royal Statistical Society Series B*, 69, 589–605.
- Frick, K., Munk, A. & Sieling, H. (2014) Multiscale change point inference. *Journal of the Royal Statistical Society Series B*, 76, 495–580.
- Gösmann, J., Stoehr, C., Heiny, J. & Dette, H. (2020) Sequential change point detection in high dimensional time series. Available from <https://arxiv.org/abs/2006.00636>.
- Hoeffding, W. (1963) Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58, 13–30.
- Horváth, L. & Rice, G. (2014) Extensions of some classical methods in change point analysis. *TEST*, 23, 219–255.
- Killick, R., Fearnhead, P. & Eckley, I.A. (2012) Optimal detection of changepoints with a linear computational cost. *Journal of the American Statistical Association*, 107, 1590–1598.
- Komlós, J. Major, P. & Tusnády, G. (1976) An approximation of partial sums of independent RVs, and the sample DF. II. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 34, 33–58.
- Laurent, B. & Massart, P. (2000) Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, 28, 1302–1338.
- Leisch, F., Hornik, K. & Kuan, C.-M. (2000) Monitoring structural changes with the generalized fluctuation test. *Econometric Theory*, 16, 835–854.
- Liu, H., Gao, C. & Samworth, R. J. (2021) Minimax rates in sparse, high-dimensional changepoint detection. *Annals of Statistics*, 49, 1081–1112.
- Lorden, G. (1971) Procedures for reacting to a change in distribution. *Annals of Mathematical Statistics*, 42, 1897–1908.
- Mei, Y. (2010) Efficient scalable schemes for monitoring a large number of data streams. *Biometrika*, 97, 419–433.
- Mörters, P. & Peres, Y. (2010) *Brownian motion*. Cambridge: Cambridge University Press.
- Oakland, J.S. (2007) *Statistical Process Control* 6th edn. London: Routledge.
- Padilla, O.H.M., Yu, Y., Wang, D. & Rinaldo, A. (2019) Optimal nonparametric multivariate change point detection and localization. Available from <https://arxiv.org/abs/1910.13289>.
- Page, E.S. (1954) Continuous inspection schemes. *Biometrika*, 41, 100–115.
- Page, E.S. (1955) A test for a change in a parameter occurring at an unknown point. *Biometrika*, 42, 523–527.
- Soh, Y.S. & Chandrasekaran, V. (2017) High-dimensional change-point estimation: Combining filtering with convex optimization. *Applied and Computational Harmonic Analysis*, 43, 122–147.
- Tartakovsky, A.G., Rozovskii, B.L., Blažek, R.B. & Kim, H. (2006) Detection of intrusions in information systems by sequential change-point methods. *Statistical Methodology*, 3, 252–293.
- Tartakovsky, A., Nikiforov, I. & Basseville, M. (2014) *Sequential Analysis: Hypothesis Testing and Changepoint Detection*. London: Chapman and Hall.
- van der Vaart, A.W. & Wellner, J.A. (1996) *Weak Convergence and Empirical Processes*. New York: Springer.
- Wang, T. & Samworth, R. J. (2018) High dimensional change point estimation via sparse projection. *Journal of the Royal Statistical Society Series B*, 80, 57–83.
- Wang, D., Yu, Y. & Rinaldo, A. (2018) Univariate mean change point detection: penalization, CUSUM and optimality. Available from <https://arxiv.org/abs/1810.09498v4>.
- Xie, Y. & Siegmund, D. (2013) Sequential multi-sensor change-point detection. *Annals of Statistics*, 41, 670–692.
- Xie, L., Xie, Y. & Moustakides, G.V. (2019) Asynchronous multi-sensor change-point detection for seismic tremors. *IEEE International Symposium on Information Theory (ISIT)*, 787–791.
- Zeileis, A., Leisch, F., Kleiber, C. & Hornik, K. (2005) Monitoring structural change in dynamic econometric models. *Journal of Applied Econometrics*, 20, 99–121.

Zou, C., Wang, Z., Zi, X. & Jiang, W. (2015) An efficient online monitoring method for high-dimensional data streams. *Technometrics*, 57, 374–387.

SUPPORTING INFORMATION

Additional Supporting Information may be found in the online version of the article at the publisher's website.

How to cite this article: Chen Y, Wang T, Samworth RJ. High-dimensional, multiscale online changepoint detection. *J R Stat Soc Series B*. 2022;00:1–33. <https://doi.org/10.1111/rssb.12447>