

Bayesian Hierarchical Model for Change Point Detection in Multivariate Sequences

Huaqing Jin and Guosheng Yin

Department of Statistics and Actuarial Science,
the University of Hong Kong

and

Binhang Yuan

Computer Science Department,
Rice University

and

Fei Jiang*

Department of Epidemiology and Biostatistics,
University of California, San Francisco

May 10, 2021

Abstract

Motivated by the wind turbine anomaly detection, we propose a Bayesian hierarchical model (BHM) for the mean-change detection in multivariate sequences. By combining the exchange random order distribution induced from the Poisson–Dirichlet process and nonlocal priors, BHM exhibits satisfactory performance for mean-shift detection with multivariate sequences under different error distributions. In particular, BHM yields the smallest detection error compared with other competitive methods considered in the paper. We utilize a local scan procedure to accelerate the computation, while the anomaly locations are determined by maximizing the posterior probability through dynamic programming. We establish consistency of the estimated number and locations of the change points and conduct extensive simulations to evaluate the BHM approach. Among the popular change point detection algorithms, BHM yields the best performance for most of the datasets in terms of the F1 score for the wind turbine anomaly detection. Supplementary materials for this article are available online.

Keywords: Change points, Multivariate data, Nonlocal prior, Non-maximum suppression, Poisson–Dirichlet process

*Corresponding author: Fei.Jiang@ucsf.edu

1 Introduction

The wind turbine blade is the core component of the wind power generation system, and hence the failure of the wind turbine blade directly leads to the dysfunction of the whole system. One leading cause of the wind turbine failure is the accumulated ice that covers the turbine as shown in the left panel of Figure 1. Typically, the Supervisory Control and Data Acquisition (SCADA) system is used to examine the ice level in real time by monitoring several signal sequences simultaneously, such as wind speed, environment temperature and accelerated speed. Once the ice on the wind turbine reaches a certain level that may lead to the failure of the system, some of the signals would exhibit mean shifts as shown in the right panel of Figure 1.

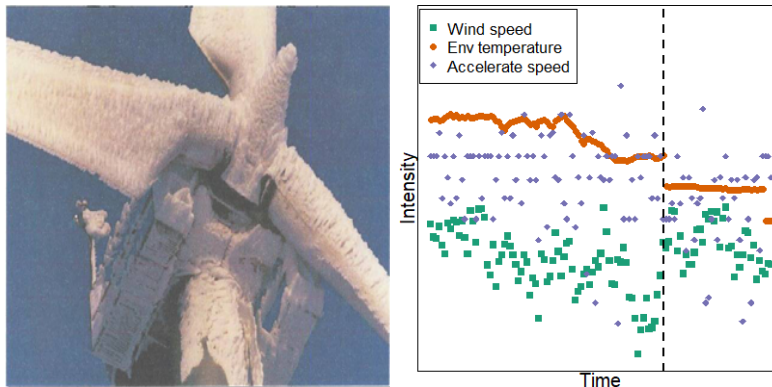


Figure 1: The wind turbine blade covered by ice (left) and the mean shifts of the wind speed, environment temperature and accelerated speed (right). The vertical dashed black line highlights the change point location.

To timely capture the anomalies on the wind turbines, we apply several existing mean detection algorithms to the signals generated by the SCADA system, including E-division (ECP) ([Matteson and James, 2014](#)), the dynamic programming based maximum likelihood

estimation (DPMLE) (Maboudou-Tchao and Hawkins, 2013), and the popular structure change detection algorithm called AutoPlait (Matsubara, Sakurai, and Faloutsos, 2014). The vertical red shadows in Figure 2 are the manually labeled change point locations based on the observed wind turbine blade failure times. Although these labels may not include all change points, they may serve as the basis for the comparison among different methods. A method would be considered the best if the detected change points constitute the smallest set that contains all the labels. As shown in Figure 2, none of the existing methods under consideration provides satisfactory results: ECP tends to select a large number of spurious change points, while DPMLE and AutoPlait miss almost all the change points.

This phenomenon motivates us to propose a novel algorithm based on a Bayesian hierarchical model (BHM) to detect mean shifts among multiple sequences. The BHM naturally borrows information across multiple sequences by using a prior induced from the Poisson–Dirichlet process, and hence it can identify true change points more effectively. In addition, BHM utilizes the nonlocal priors to control the false discoveries, which is shown to yield the smallest detection error in comparison with the existing methods under consideration. To reduce the computational burden, we introduce a initial local scan procedure in our algorithm, and then utilize the dynamic programming (Bellman and Roth, 1969; Du et al., 2016) to identify the change point locations by optimizing the posterior probability.

The contributions of our work are three-fold: (i) We develop a novel Bayesian method to estimate both the number and locations of the change points in an integrative manner. Our method is shown to outperform the competitive ones in both simulation studies and real application to the wind turbine data. (ii) We explore the advantages of using nonlocal

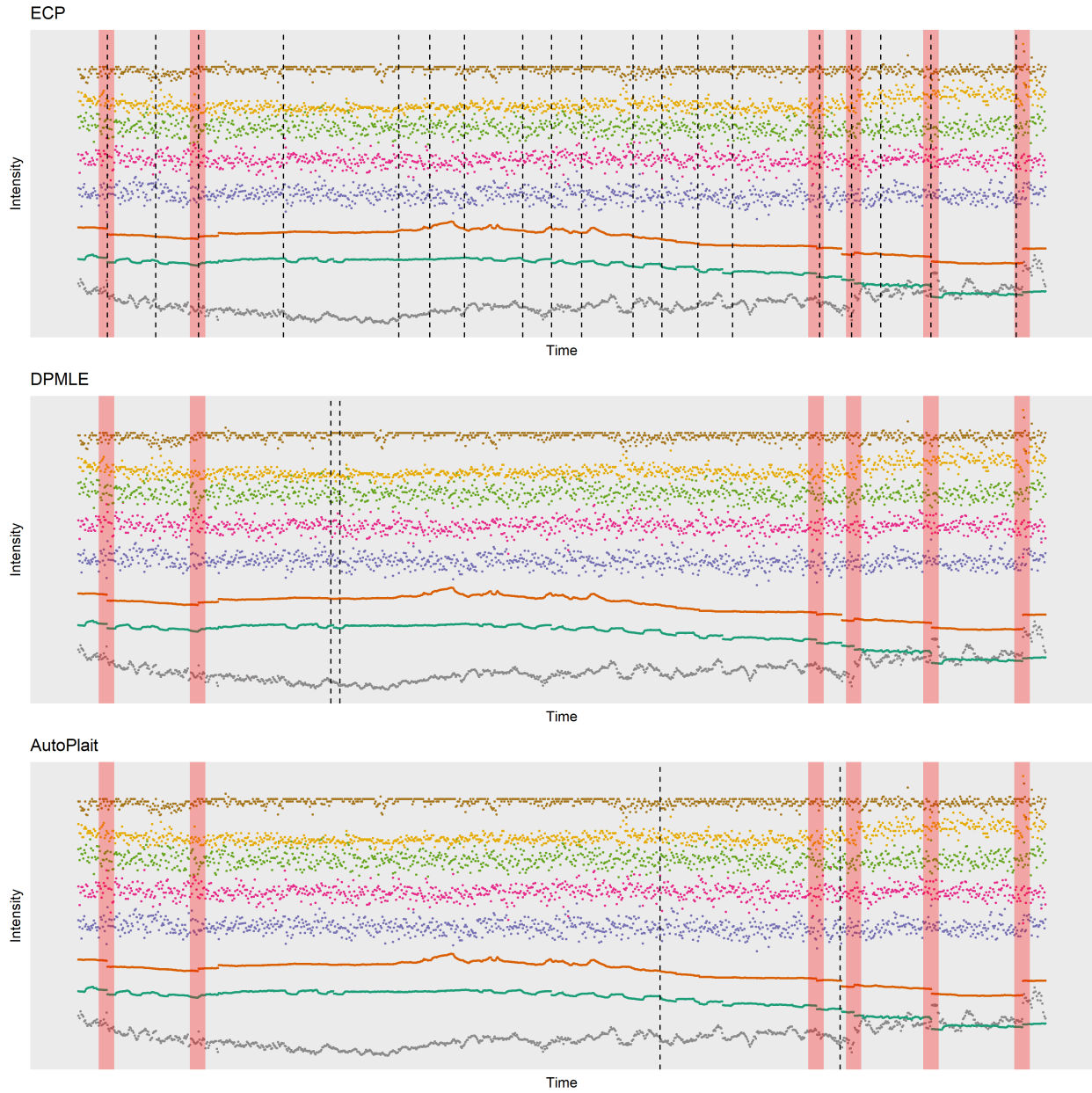


Figure 2: The detection results of three existing methods, namely ECP, DPMLE, AutoPlait, for the wind turbine dataset 1 with $n = 8$ sequences. The black dashed lines are the estimated state shift points and the red shadows are the m_I -neighbourhood of the ground truth.

priors in BHM for reducing the detection errors of the change points with multivariate sequences. (iii) We establish the consistency of our BHM method, in that it can identify the correct number and locations of change points asymptotically, providing the sample size is sufficiently large.

The rest of the paper is organized as follows. We discuss the related work in Section 2, and provide the BHM method and its theoretical properties in Section 3 and Section 4, respectively. In Section 5, we conduct extensive simulation studies to compare our proposal with some existing methods. Finally, we apply the proposed method to the wind turbine data in Section 6.

2 Related work

The wind turbine anomaly detection problem has been extensively studied (Tautz-Weinert and Watson, 2016) under different contexts. Other than the standard SCADA system, novel anomaly physical detectors (Muñoz et al., 2018) with higher signal-to-noise ratios have been developed aiming at amplifying the abnormal signals of the wind turbine failure. However, the practical use of these new physical detectors is limited due to the high cost (Yang et al., 2013). Moreover, many anomaly detection algorithms have been developed to capture the anomalies using the widely used SCADA system (Tautz-Weinert and Watson, 2016). Yang et al. (2013) proposed a trending method using bin averaging with output power, wind speed or generator speed, and a quantifying criterion was introduced based on a correlation model of historical and present data. Kim et al. (2011) applied an artificial neural network self-organising map approach to the SCADA data, which detected the

anomalies through clustering. Relying upon the correlation analysis and physics of the system, high-order polynomial models were developed for the anomaly detection (Wilkinson et al., 2014). Gray and Watson (2010) introduced a damage model based on a physical understanding of the particular failure mode of interest for damage calculation as well as failure probability estimation. However, all the aforementioned methods rely on additional knowledge about other wind turbine features, historical data or domain experience, which are not available in the current wind turbine dataset. This motivates us to develop a new change point detection algorithm to identify anomalies solely based on the signal patterns.

Change point detection algorithms have been widely used to detect anomalies (Muggeo and Adelfio, 2010). Under the frequentist framework, the change point detection relies on optimizing certain objective functions, such as a parametric or nonparametric log-likelihood function (Hawkins, 2001; Zou et al., 2014), quadratic loss (Rigail, 2015) and cumulative sums (Hinkley, 1971; Manogaran and Lopez, 2018). The Bayesian information criterion (BIC) (Yao, 1988) and its variants (Yao and Au, 1989; Zhang and Siegmund, 2007) are commonly used for determining the number of change points.

On the other side, the Bayesian algorithms identify the change point locations by maximizing the posterior distribution (Martínez and Mena, 2014) or the marginal likelihood (Du et al., 2016). The optimization routines are performed through the Markov chain Monte Carlo (MCMC) (Barry and Hartigan, 1993; Martínez and Mena, 2014) or dynamic programming (Du et al., 2016). More recently, Hopfield’s network has been advocated for use to identify change points (Fuentes-García et al., 2019). By considering the randomness in the model parameters and incorporating prior distributions, Bayesian methods auto-

matically add penalties on the number of change points ([Lavielle, 2005](#)). As a result, the number of change points can be determined seamlessly in conjunction with their locations ([Truong et al., 2018](#)).

Change point problems can also be considered under the context of the statistical process control (SPC) ([Qiu, 2013](#)). The conventional methods for change point detection under SPC are the Shewhart and cumulative sum charts as well as the exponential weighted moving average (EWMA) chart ([Hawkins et al., 2003](#)). [Tsiamyrtzis and Hawkins \(2005\)](#) introduced a dynamic model to handle the short-run process with the aim to detect the mean shift. Using a sequential estimation technique, [Zamba and Hawkins \(2006\)](#) considered the multivariate change point problem for SPC when the in-control parameters were unknown. Using an autoregression model, [Tsiamyrtzis and Hawkins \(2008\)](#) worked on autocorrelated processes with mean shift under a Bayesian EWMA method. The goal of these change point detection methods is to detect a single change in the sequence, while they can be adapted to identify multiple changes with an unknown number of change points.

To construct the posterior distribution or the marginal likelihood under the Bayesian paradigm, one crucial step is to specify the prior distributions. The commonly used priors for the mean differences include the local priors ([Bertolino et al., 2000](#)), such as normal prior distributions ([Du et al., 2016](#)), and nonlocal priors, such as the moment prior and inverse moment prior distributions ([Johnson and Rossell, 2010](#)). The nonlocal priors were first proposed in the Bayesian hypothesis testing framework to improve the speed of the accumulation of the evidence in favour of the true null model ([Johnson and Rossell, 2010](#)). [Jiang, Yin, and Dominici \(2018\)](#) applied the nonlocal prior to identify the change points

in a single sequence of data, which leads to a faster convergence rate compared with the algorithms based on the local priors.

Because the change point detection can be regarded as a special clustering problem, we can utilize the exchangeable random partition distribution (ERPD) (Pitman, 1995) to construct the prior distribution of the segments. The ERPD has been widely used for clustering problems (Lau and Green, 2007; Wade et al., 2018), which penalizes the model complexity (Pitman, 2002) and automatically selects the number of clusters (McCullagh and Yang, 2008). However, ERPD is not directly applicable to the change point detection, because it does not account for the order constraints in the change point problem. Martínez and Mena (2014) proposed to use a modified ERPD, namely the exchangeable random order distribution (EROD), as the prior distribution specifically for the change point detection, which inherits the symmetry and automatic penalization properties of ERPD.

In addition to the mean-change detection, many other applications of change points have been carried out. Matsubara, Sakurai, and Faloutsos (2014) proposed the Auto-Plait method to detect the changes in the periodic sequences for identifying the structure changes in the signals. Gharghabi et al. (2017) proposed a fast, low-cost online semantic segmentation for the structure change detection in cyclic data. Bouchard and Badler (2007) introduced a Laban movement based segmentation method for the motion data. Gong, Medioni, and Zhao (2014) utilized the kernelized temporal cut method to recognize action changes in the continuous monocular motion sequences.

3 Bayesian multivariate change point detection

3.1 Probability model

Suppose there are n sequences of signals measured over time, where correlations exist both between sequences and within sequences. Let Y_{ik} represent the strength of the i th signal at time k , $i = 1, \dots, n$; $k = 1, \dots, T$. Define $\mathbf{Y}_k = (Y_{1k}, \dots, Y_{nk})^\top$, and $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_T)$ is an $n \times T$ observation matrix and $\mathbf{Y}_{(a,b)} = (\mathbf{Y}_{a+1}, \dots, \mathbf{Y}_b)$ represents the matrix containing the $(a + 1)$ th to b th columns of \mathbf{Y} . Let $\mathcal{K} = \{\kappa_0, \dots, \kappa_{p+1}\}$ be a generic notation for a set of p change points, with $\kappa_0 = 0$ and $\kappa_{p+1} = T$, and let $\mathcal{K}_0 = \{\kappa_{00}, \dots, \kappa_{0,p_0+1}\}$ be the true change point set, with $\kappa_{00} = 0$ and $\kappa_{0,p_0+1} = T$. In addition, let $\{N_s, s = 0, \dots, p\} = \{(\kappa_{s+1} - \kappa_s), s = 0, \dots, p\}$ be a collection of numbers of observations between consecutive change points.

Let $\bar{Y}_{i,\kappa_s} = (\kappa_s - \kappa_{s-1})^{-1} \sum_{k=\kappa_{s-1}+1}^{\kappa_s} Y_{ik}$, with $\bar{Y}_{i,\kappa_0} = 0$. We assume

$$\frac{Y_{ik} - \bar{Y}_{i,\kappa_s} - \mu_{is}}{\omega_s} \Big| \mathcal{K}, \mu_{is}, \omega_s \sim \pi_0, \quad \text{for } k \in (\kappa_s, \kappa_{s+1}],$$

$$\mu_{is} \sim \pi_\mu(\mu_{is}),$$

$$\omega_s \sim \pi_\omega(\omega_s),$$

$$(N_0, \dots, N_p) \sim \pi_k(\mathcal{K}), \tag{1}$$

where π_0 is the likelihood function selected according to the data distribution. It is worth emphasizing that we model the difference $Y_{ik} - \bar{Y}_{i,\kappa_s}$ rather than Y_{ik} , and such a modeling strategy allows us to better control the convergence rates by selecting suitable priors on the mean differences. In model (1), we assume that all the sequences share the same variance for each segment, which can help to reduce the computational burden and numerical error

in the optimization procedure. As shown in Table A.1, using the same variance parameter does not undermine the performance of our algorithm. In practice, when the sequences have very distinct variances, we may standardize the variance for each sequence before the implementation of our detection procedure. Typically, we use a normal likelihood if the data do not contain outliers, and use a t distribution if the data are contaminated with a substantial amount of outliers. In practice, the existence of outliers can be determined by the generalized extreme studentized deviate test (Rosner, 1983), as discussed in Section A.3 of the Supplementary Materials.

3.2 Prior specifications

Change point detection is closely related to a special clustering algorithm in which each cluster only contains the neighborhood points. This motivates us to consider ERPD (Pitman, 1995; Gnedin and Pitman, 2006) induced from the Poisson–Dirichlet process as the prior distribution for (N_0, \dots, N_p) , which has been widely used in clustering problems (Broderick et al., 2013). By choosing $\sigma = 0$ or $\alpha = 0$, the Poisson–Dirichlet process reduces to the Dirichlet or the normalized stable processes, respectively (Martínez and Mena, 2014). The ERPD strikes a good balance between the generalization and complexity. As a special case of Gibbs-type priors, it places a tradeoff on the prior distribution between being informative and noninformative about the number of change points (Lijoi et al., 2007). However, this prior may classify non-neighborhood signals into the same group, and hence the resulting clusters would contradict with the property of the change points. To account

for this neighborhood constraint, the probability mass function of ERPD is modified as

$$\pi_k(\mathcal{K}) = \frac{T!}{(p+1)! \prod_{h=0}^p N_h!} \cdot \frac{\prod_{s=1}^p (\alpha + s\sigma)}{\prod_{j=1}^{T-1} (\alpha + j)} \prod_{h=0}^p \prod_{i=1}^{N_h-1} (i - \sigma), \quad \sigma \in [0, 1), \alpha > -\sigma, \quad (2)$$

which corresponds to the probability mass function of EROD (Pitman, 2002; Martínez and Mena, 2014). The first term in (2) accounts for the neighborhood constraint.

Under the prior distribution in (2), the marginal distribution of the number of change points p can be derived as

$$\Pr(p = l) = \frac{\prod_{i=1}^l (\alpha + i\sigma)}{\sigma^{l+1} \prod_{i=1}^{n-1} (\alpha + i)} \frac{1}{(l+1)!} \sum_{j=0}^{l+1} (-1)^j \binom{l+1}{j} \prod_{i=0}^{n-1} (-j\sigma + i).$$

We can select the values of (α, σ) via encoding our prior belief in the number of change points for the data. As discussed in Martínez and Mena (2014), the parameter σ plays a more important role for detecting the change points, and thus it deserves more attention in practice.

The prior for the mean difference, π_μ , is critical for reducing the detection error of the change points. We consider the local prior, nonlocal moment prior and inverse moment prior, respectively defined as follows:

$$\begin{aligned} \pi_{\mu,L}(\mu) &= \text{N}(0, \psi^2), \\ \pi_{\mu,M}(\mu) &= \frac{\mu^{2v}}{C_M} \frac{1}{\sqrt{2\pi}} \exp(-\mu^2/2), \\ \pi_{\mu,I}(\mu) &= \frac{r\phi^{q/2}}{\Gamma(q/2r)} |\mu|^{-(q+1)} \exp\{-(\mu^2/\phi)^{-r}\}, \end{aligned}$$

where $v \geq 1$, $\psi, r, q, \phi > 0$ and C_M is the normalizing constant. For the nuisance parameters ω_s , we take $\pi_\omega(\omega_s)$ to be a Gamma distribution.

3.3 Posterior distribution and change point detection

Based on model (1), the marginal likelihood given a change point set \mathcal{K} can be written as

$$\Pr(\mathbf{Y}|\mathcal{K}) = \prod_{s=0}^p \int \prod_{i=1}^n \int \prod_{k \in (\kappa_s, \kappa_{s+1}]} \frac{1}{\omega_s} \pi_0 \left(\frac{Y_{ik} - \bar{Y}_{i, \kappa_s} - \mu_{is}}{\omega_s} \right) \pi_\mu(\mu_{is}) d\mu_{is} \pi_\omega(\omega_s) d\omega_s. \quad (3)$$

Using the prior in (2), the posterior probability of \mathcal{K} is

$$\Pr(\mathcal{K}|\mathbf{Y}) \propto \Pr(\mathbf{Y}|\mathcal{K}) \pi_k(\mathcal{K}). \quad (4)$$

The standard procedure to optimize the posterior is through dynamic programming (Bellman and Roth, 1969; Du et al., 2016). However, because the dynamic programming evaluates the signals at every time point, the computation time grows in the order of $O(MT^2)$ where M is the upper bound of the number of change points (Rigaiil, 2015; Du et al., 2016). When T is large, the computational burden is prohibitive. To alleviate the computational burden, we propose a screening procedure to reduce the search space of the change points.

While we adopt a Bayesian mechanism for change point detection, it is not a fully Bayesian approach. The prior is used to induce penalties on the parameters so as to identify the best set of change points. More specifically, by using a prior on \mathcal{K} , the algorithm induces a penalty on the locations and number of change points. Hence, optimizing the posterior distribution of \mathcal{K} automatically provides the best estimates for the locations and number of change points. Furthermore, the prior on the mean difference induces a penalty on μ_{is} , which reduces the false positive rate of detection.

3.4 Screening candidate points

Through the screening step, we can reduce the search space of the change points to a subset of time points which is guaranteed to cover the true change points asymptotically. This leads to a substantial reduction in the computational time when there are much fewer candidate points than the total measurement times in the data.

Let $\bar{Y}_{ik} = m_I^{-1} \sum_{l=k-m_I+1}^k Y_{il}$, where m_I is a prespecified window size and the m_I -neighborhood is defined as the set $\{\mathbf{Y}_l : l \in (k - m_I, k + m_I)\}$. We construct a local scan statistic,

$$R_k = \prod_{i=1}^n \frac{\int \prod_{l=k+1}^{k+m_I} \exp\{-(Y_{il} - \bar{Y}_{ik} - \mu)^2\} \pi_\mu(\mu) d\mu}{\prod_{l=k+1}^{k+m_I} \exp\{-(Y_{il} - \bar{Y}_{ik})^2\}}. \quad (5)$$

A large value of R_k favors the k th signal vector to be the only change point in its m_I -neighborhood. Thus, $R_k \rightarrow \infty$ if k is a true change point, and $R_k \rightarrow 0$ if there is no change point in the m_I -neighborhood of the k th signal. Based on this property, we develop Algorithm 1 for selecting the candidate points.

Algorithm 1 Screening Candidate Points

- (i) For each $k \in [m_I, T - m_I]$, compute R_k .
 - (ii) If $R_k = \max\{R_j, j \in (k - m_I, k + m_I)\}$, then k is selected as a candidate point.
-

3.5 Change point detection with candidate points

Let $\mathcal{H}(m_I) = \{\tau_0, \dots, \tau_{N+1}\}$ be the set containing the candidate points where $\{\tau_i\}_{i=1}^N$ are obtained by Algorithm 1 and $\tau_0 = 0$ and $\tau_{N+1} = T$. Given \mathcal{K} , we can define the utility

function, $U(\mathcal{K}|\mathbf{Y}) = \prod_{s=0}^p u(\mathbf{Y}_{(\kappa_s, \kappa_{s+1}]}, s)$, where

$$u(\mathbf{Y}_{(\kappa_s, \kappa_{s+1}]}, s) = \int \prod_{i=1}^n \int \prod_{k \in (\kappa_s, \kappa_{s+1}]} \frac{1}{\omega_s} \pi_0 \left(\frac{Y_{ik} - \bar{Y}_{i\kappa_s} - \mu_{is}}{\omega_s} \right) \pi_\mu(\mu_{is}) d\mu_{is} \pi_\omega(\omega_s) d\omega_s \\ \times \frac{(\alpha + s\sigma)}{(s+1)} \frac{\prod_{i=1}^{N_s-1} (i - \sigma)}{N_s!}, \quad (6)$$

and $u(\mathbf{Y}_{(\kappa_s, \kappa_{s+1}]}, s)$ can be regarded as the utility function for segment $\mathbf{Y}_{(\kappa_s, \kappa_{s+1}]}$.

The estimator for the set of change points is given by

$$\hat{\mathcal{K}} = \operatorname{argmax}_{\mathcal{K} \subseteq \mathcal{H}(m_I)} U(\mathcal{K}|\mathbf{Y}).$$

To optimize $U(\mathcal{K}|\mathbf{Y})$ over $\mathcal{H}(m_I)$ via dynamic programming, we require that $u(\mathbf{Y}_{(\kappa_s, \kappa_{s+1}]}, s)$ only depends on $\{\kappa_s, \kappa_{s+1}, s\}$ given $\mathcal{H}(m_I)$, so we replace \bar{Y}_{i, κ_s} in (6) by $\tilde{Y}_{i, \kappa_s} = (\tau_l - \tau_{l-1})^{-1} \sum_{k \in (\tau_{l-1}, \tau_l]} Y_{ik}$ with $\tau_l = \kappa_s$, $\kappa_s \in \mathcal{K}$. Note that \tilde{Y}_{i, κ_s} is the sample mean of the segment $(\tau_{l-1}, \tau_l]$, while \bar{Y}_{i, κ_s} is the sample mean of the segment $(\kappa_{s-1}, \kappa_s]$. Both are consistent estimators of the mean of signals in the segment $(\kappa_{s-1}, \kappa_s]$. The dynamic programming procedure is presented as Algorithm A.1 in the Supplementary Materials, which has computational time $O(MN^2)$, in comparison with $O(MT^2)$ of the existing algorithms (Rigaill, 2015; Du et al., 2016).

4 Theoretical Properties

In this section, we present the theoretical properties of the BHM method. Lemma 1 below shows that $\mathcal{H}(m_I)$ covers \mathcal{K}_0 with probability one.

Lemma 1. *Suppose that regularity conditions (1)–(2) in Section C of the Supplementary Materials hold. Let η_{is}^2 be the variance of Y_{ik} in the s th segment based on the true change*

points. Assume $m_I^{1/2}\delta_I/\eta \rightarrow \infty$ where $\eta = \max_{\{i=1,\dots,n;s=0,\dots,p_0\}} \eta_{is}$, δ_I is defined in condition (2) and $m_I/(\log(T))^{1+\epsilon} \rightarrow c > 0$ for $\epsilon > 0$. If $\min_{\{i=0,\dots,p_0\}}(\kappa_{0,i+1} - \kappa_{0,i}) > m_I$, then for each $\kappa_{0j} \in \mathcal{K}_0$, there is a $\tau \in \mathcal{H}(m_I)$, such that $\Pr\{\kappa_{0j} \in (\tau - m_I, \tau + m_I)\} \rightarrow 1$ as $T \rightarrow \infty$.

The proof of Lemma 1 is given in Section C of the Supplementary Materials. The condition $m_I/(\log(T))^{1+\epsilon} \rightarrow c > 0$ regulates the selection of m_I , which grows no slower than $\log(T)$. Lemma 1 indicates that $\mathcal{H}(m_I)$ should cover \mathcal{K}_0 asymptotically, while the cardinality of $\mathcal{H}(m_I)$, $N + 2$, is far less than T . Therefore, when performing the dynamic programming on the smaller set $\mathcal{H}(m_I)$, the algorithm guarantees a positive probability to select the true change points and at the same time improves the computational efficiency.

Furthermore, Theorem 1 shows that the screening step would accelerate the computational speed without sacrificing the statistical consistency. Let $\mathcal{K}_0 = \{\kappa_{01}, \dots, \kappa_{0p_0}\}$ be the true set of change points with $\min_{\{i=0,\dots,p_0\}}(\kappa_{0,i+1} - \kappa_{0,i}) > m_I$, and let $\widehat{\mathcal{K}} = \{\hat{\kappa}_1, \dots, \hat{\kappa}_{\hat{p}}\}$ be the estimated set of change points.

Theorem 1. *Suppose that regularity conditions (1)–(4) in Section C of the Supplementary Materials hold. Assume $m_I^{1/2}\delta_I/\eta \rightarrow \infty$ and $m_I/(\log(T))^{1+\epsilon} \rightarrow c > 0$ for $\epsilon > 0$, and \mathcal{K}_0 is a subset of $\mathcal{H}(m_I)$. Then, as $T \rightarrow \infty$,*

$$\hat{p} \xrightarrow{\mathcal{P}} p_0 \quad \text{and} \quad \sup_{b \in \mathcal{K}_0} \inf_{a \in \widehat{\mathcal{K}}} |a - b| = O_p(1).$$

The proof of Theorem 1 is delineated in Section C of the Supplementary Materials. Theorem 1 establishes the consistency of the estimated change points when $\mathcal{H}(m_I)$ covers \mathcal{K}_0 . Combined with the result in Lemma 1 that each true change point falls in the m_I -

neighborhood of at least one candidate point, we have

$$\hat{p} \xrightarrow{\mathcal{P}} p_0 \quad \text{and} \quad \sup_{b \in \mathcal{K}_0} \inf_{a \in \hat{\mathcal{K}}} |a - b| = O_p(m_I).$$

Hence, Lemma 1 and Theorem 1 imply that as $T \rightarrow \infty$, the estimated change points are guaranteed to fall in the m_I -neighborhood of the true change points.

5 Simulation Studies

5.1 Simulation settings

We conduct simulation experiments to evaluate the properties of the BHM method in comparison with two existing methods. First, we introduce the two main assessment metrics: the over-segmentation error,

$$d(\hat{\mathcal{K}}|\mathcal{K}_0) = \sup_{b \in \mathcal{K}_0} \inf_{a \in \hat{\mathcal{K}}} |a - b|,$$

and the under-segmentation error,

$$d(\mathcal{K}_0|\hat{\mathcal{K}}) = \sup_{b \in \hat{\mathcal{K}}} \inf_{a \in \mathcal{K}_0} |a - b|.$$

The over- or under-segmentation errors would be larger if we select fewer or more change points than the truth, respectively. The maximum segmentation error is $\max\{d(\hat{\mathcal{K}}|\mathcal{K}_0), d(\mathcal{K}_0|\hat{\mathcal{K}})\}$, and the estimation error for p_0 is $|\hat{p} - p_0|$.

The locations of change points are $\mathcal{K}_0 = \{\lfloor T \times r_i \rfloor\}_{i=1}^{p_0}$, with $p_0 = 10$ and

$$\{r_i\}_{i=1}^{10} = \{0.025, 0.155, 0.220, 0.365, 0.395, 0.495, 0.630, 0.725, 0.865, 0.975\}.$$

We adopt $n = 2$ for parameter tuning and $n = 8$ for comparing the BHM method with other methods. We define

$$\mathbf{g}(i, k) = \left[\frac{\{1 + \text{sgn}(k - \kappa_{0j})\}(1 - \mathbb{I}_{\{j \in A_i\}})}{2}, j = 1, \dots, p_0 \right]^\top, \quad i = 1, \dots, n; k = 1, \dots, T,$$

where $A_i = \{(3l+i) \bmod p_0, l = 0, 1, 2\}$, $\text{sgn}(\cdot)$ is the sign function and $\mathbb{I}_{\{\cdot\}}$ is the indicator function.

The data are generated with different dimensions from the base model,

$$Y_{ik} = 2 + \mathbf{d}^\top \mathbf{g}(i, k) + \xi_{ik} \prod_{j=1}^{\mathbf{1}_{p_0}^\top \mathbf{g}(i, k)} v_j, \quad i = 1, \dots, n; k = 1, \dots, T, \quad (7)$$

where $\mathbf{d} = (2.5, -2.8, 2.4, 2.6, -3, -2.9, 3.1, -2.5, -2.7, 2.6)^\top$ are the mean differences between consecutive segments, $\mathbf{1}_{p_0}$ is a p_0 -dimensional vector of 1's, ξ_{ik} 's represent the errors and $[v_j]_{j=1}^{p_0}$ controls the homogeneity of the variances for different segments. If $[v_j]_{j=1}^{p_0} = \mathbf{1}_{p_0}$, the model yields a homogeneous variance across the segments.

By selecting distinct $[v_j]_{j=1}^{p_0}$ and different error distributions, we can obtain different data-generating models. The models used in the simulation studies are summarized as follows,

I : $[v_j]_{j=1}^{p_0} = \mathbf{1}_{p_0}$ and independent standard normal errors.

II : $[v_j]_{j=1}^{p_0} = (0.6, 2, 2/3, 0.6, 2, 2/3, 0.6, 2, 2/3, 0.6)$ and independent standard normal errors.

III : $[v_j]_{j=1}^{p_0} = \mathbf{1}_{p_0}$ and independent $t(5)$ errors with unit variance.

IV : $[v_j]_{j=1}^{p_0} = \mathbf{1}_{p_0}$ and independent skewed normal errors with slant parameter 4 and variance 1 (O'Hagan and Leonard, 1976).

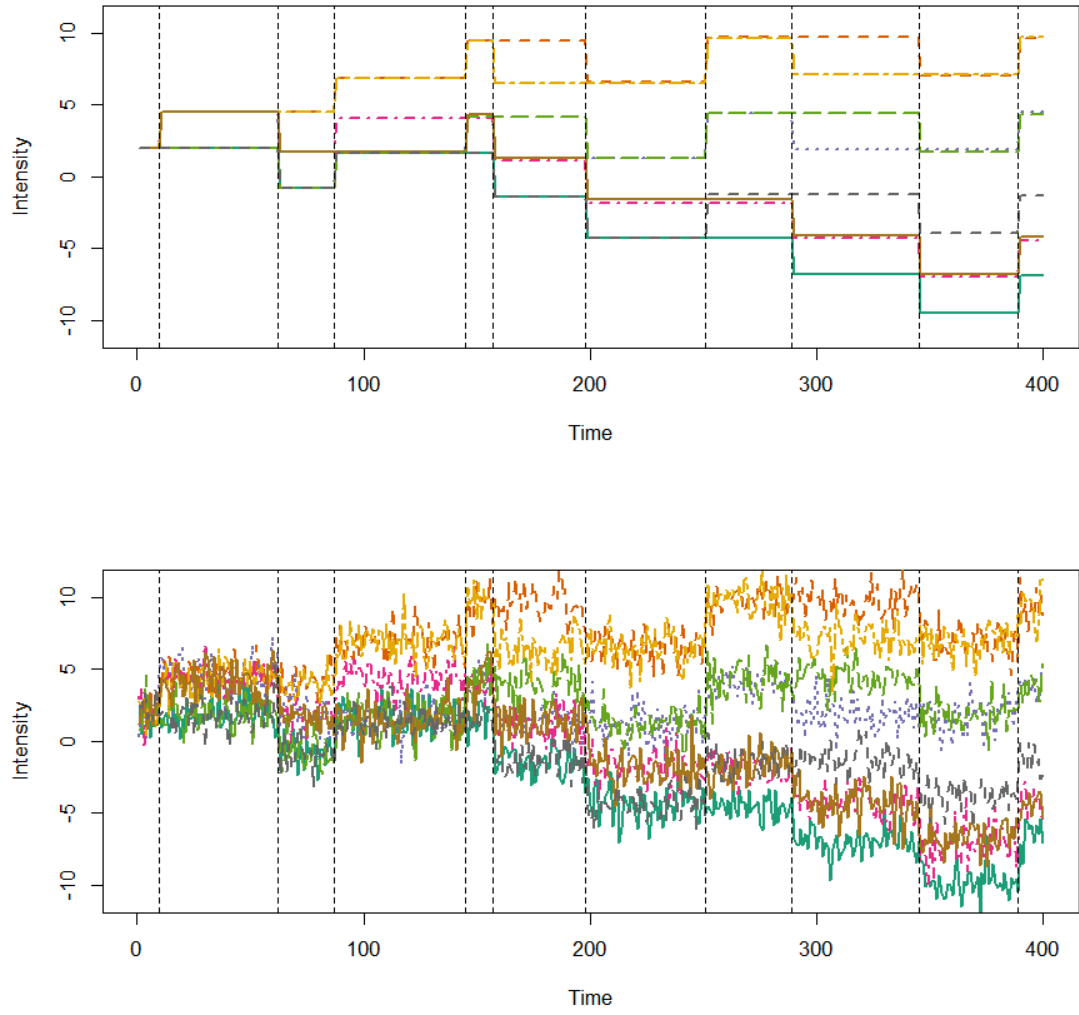


Figure 3: The mean of the simulated data (top) and randomly generated data under model I (bottom) with sample size $T = 400$ and $n = 8$. The vertical dashed lines indicate the true change points.

V : $[v_j]_{j=1}^{p_0} = \mathbf{1}_{p_0}$ and autocorrelated standard normal errors under a moving average (MA) model, $\xi_{ik} = a_k - 1.5a_{k-1}$ with $a_k \sim N(0, 4/13)$.

VI : $[v_j]_{j=1}^{p_0} = \mathbf{1}_{p_0}$ and standard normal errors with correlated sequences of correlation coefficient 0.3.

The top panel of Figure 3 presents the mean of the simulated data under the base model (7) with $T = 400$ and $n = 8$. At some change points, the mean shifts only happen in certain sequences while the rest keep unchanged. For illustration, we also display the simulated sequences under model I in the bottom panel of Figure 3.

We choose $\alpha = 1$, $\sigma = 0.3$ in the prior π_k , as they yield the smallest maximal segmentation errors under models I and II as shown in Figure A.1 of the Supplementary Materials. We choose π_ω to be Gamma(1, 1).

5.2 Tuning parameters

We use the simulation to find a suitable window size m_I for our method. Lemma 1 indicates that m_I should grow no slower than $\log(T)$. Hence, we select $m_I = \{\log(T)\}^{1.5} h$, where h is a constant to be determined numerically. Figure 4 exhibits the relationship between h and $|\hat{p} - p_0|$, and that between h and the maximum segmentation error under models I and II with $T = 400$ and $n = 2$, respectively. Clearly, $h = 0.55$ leads to the overall smallest errors for both models. In general, BHM works well for $h \in [0.5, 0.6]$.

We also compare the performances of using different priors for the mean difference under model I, including the normal prior ($\pi_{\mu,L}$) with $\psi = 2$, moment prior ($\pi_{\mu,M}$) with $v = 1$ and inverse moment prior ($\pi_{\mu,I}$) with $q = \phi = 2$ and $r = 0.6$. For fair comparisons, we

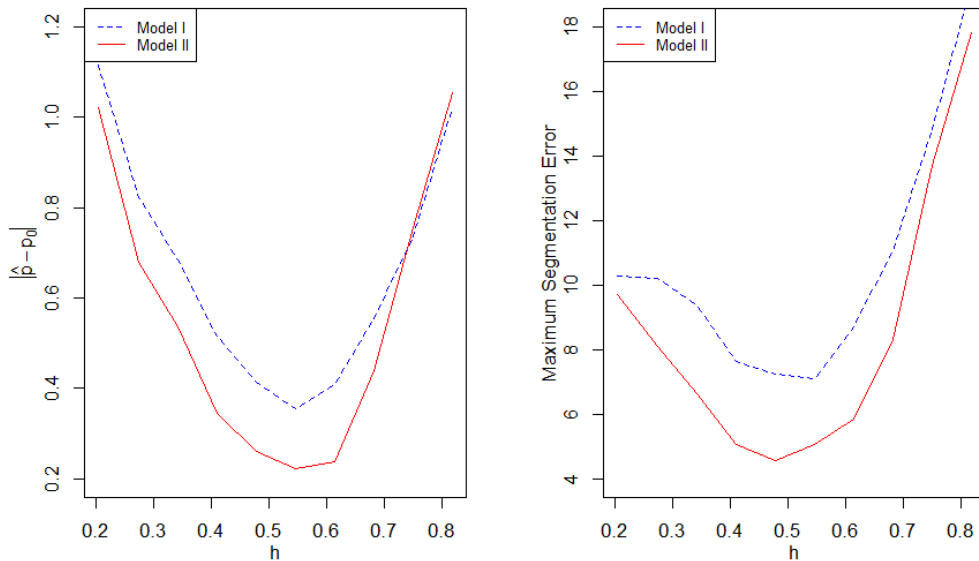


Figure 4: The absolute difference $|\hat{p} - p_0|$ (left) and the maximum segmentation error (right) versus h over 500 simulations with sample size $T = 400$, $n = 2$ under models I and II, respectively.

select the best tuning parameters for each prior, so as to achieve the lowest segmentation error under model I, as shown in Figure A.2.

Figure A.3 shows the simulation results under model I with $n = 2$, where both $|\hat{p} - p_0|$ and the maximum segmentation error decrease as the sample size increases. Further, the two nonlocal priors have similar performances and both outperform the local prior by yielding smaller errors, especially when the sample size is larger than 350.

Table 1: Comparison of BHM-FIX, BHM-MPP, ECP and DPMLE when $n = 8$ under the six data-generating models over 500 simulations. Standard deviations are given in parentheses.

Data-generating Model	Method	$\hat{p} - p_0$							Segmentation Error	
		≤ -3	-2	-1	0	1	2	≥ 3	$d(\hat{\mathcal{K}} \mathcal{K}_0)$	$d(\mathcal{K}_0 \hat{\mathcal{K}})$
I	BHM-FIX	0	0	0	500	0	0	0	0.14 (0.38)	0.14 (0.38)
	BHM-MPP	0	0	0	500	0	0	0	0.08 (0.29)	0.08 (0.29)
	ECP	0	0	0	471	26	3	0	0.18 (0.43)	1.10 (4.11)
	DPMLE	500	0	0	0	0	0	0	78.32 (20.14)	0.08 (0.32)
II	BHM-FIX	0	0	0	500	0	0	0	0.05 (0.23)	0.05 (0.23)
	BHM-MPP	0	0	0	500	0	0	0	0.02 (0.15)	0.02 (0.13)
	ECP	0	0	0	477	21	2	0	0.08 (0.29)	0.12 (0.35)
	DPMLE	499	1	0	0	0	0	0	52.14 (1.97)	3.85 (0.36)
III	BHM-FIX	0	0	0	500	0	0	0	0.17 (0.40)	0.17 (0.40)
	BHM-MPP	0	0	0	500	0	0	0	0.11 (0.32)	0.11 (0.32)
	ECP	0	0	0	469	28	3	0	0.20 (0.45)	1.20 (4.41)
	DPMLE	500	0	0	0	0	0	0	72.02 (21.76)	0.12 (0.37)
IV	BHM-FIX	0	0	0	500	0	0	0	0.10 (0.31)	0.11 (0.34)
	BHM-MPP	0	0	0	500	0	0	0	0.10 (0.31)	0.09 (0.30)
	ECP	0	0	0	474	22	4	0	0.15 (0.40)	0.20 (0.45)
	DPMLE	500	0	0	0	0	0	0	75.60 (21.87)	5.80 (0.51)
V	BHM-FIX	0	0	0	500	0	0	0	0.10 (0.32)	0.11 (0.32)
	BHM-MPP	0	0	0	500	0	0	0	0.10 (0.27)	0.08 (0.27)
	ECP	0	0	0	500	0	0	0	0.15 (0.38)	0.17 (0.38)
	DPMLE	500	0	0	0	0	0	0	73.56 (21.39)	0.07 (0.26)
VI	BHM-FIX	0	0	0	486	14	0	0	1.02 (1.65)	1.60 (3.86)
	BHM-MPP	0	0	0	488	12	0	0	0.67 (0.83)	1.04 (2.57)
	ECP	0	0	0	476	20	4	0	0.45 (0.64)	1.25 (3.86)
	DPMLE	500	0	0	0	0	0	0	77.70 (22.67)	0.10 (0.31)

5.3 Comparison with other methods

We compare BHM with the ECP and DPMLE mean-change detection methods under all the six data-generating models with sample size $T = 400$ and $n = 8$. For the ECP and DPMLE methods, we adopt the default parameters from the original papers, respectively. The BHM method utilises a normal likelihood for all models except for model III where a t -distribution likelihood is used. The moment prior $\pi_{\mu, M}$ is chosen for mean differences according to the results in Figure A.3. We adopt two methods to select parameters (α, σ, v, h) , as they are essential for BHM. (i) We set $(\alpha, \sigma, v, h) = (1, 0.3, 1, 0.55)$ as they deliver satisfactory performances in the simulation studies (denoted as BHM-FIX). (ii) We tune parameters (α, σ, v, h) to maximize the posterior probability $\Pr(\mathcal{K}|\mathbf{Y})$ (denoted as BHM-MPP). Table 1 shows the results for $n = 8$ under the six models. The BHM-MPP consistently outperforms BHM-FIX under all the six models, indicating that maximization of the posterior probability is a more effective method to select hyper-parameters for the BHM method in practice. Nevertheless, the performance of BHM-FIX is only slightly inferior to that of BHM-MPP under all six settings, suggesting that the selected parameters in Section 5.2 achieve high estimation accuracy with small $|\hat{p} - p_0|$ and segmentation errors. In practice, $(v, h) = (1, 0.55)$ can be set as the default parameters while the values of (α, σ) should be selected via the prior belief on the number of change points or the MPP method, as they rely on the number and locations of change points.

Overall, the proposed BHM-FIX and BHM-MPP are superior to other methods by yielding the smallest $|\hat{p} - p_0|$ and segmentation errors under all the six models. Although incorrect models are adopted under models IV, V and VI, the BHM-FIX and BHM-MPP

still perform better than the competitive methods due to the robustness property.

6 Wind turbine data

For illustration, we apply our BHM method to detect the changes in the wind turbine data available from a wind turbine anomaly detection contest. For the details of the contest, refer to the link (http://www.caict.ac.cn/kxyj/qwfb/bps/201804/t20180426_158519.htm). It includes a total of seven datasets which are manually labeled for the wind turbine failure times serving as the ground truth. Each dataset contains eight sequences ($n = 8$), corresponding to wind speed, cabin temperature, environment temperature, accelerated speed along horizontal and vertical directions and the ng5 temperatures from three pitches. In each dataset, there are 4 to 8 change points, i.e., $p_0 \in [4, 8]$ in the sequences. The lengths of the sequences are from 400 to 1000, i.e., $T \in [400, 1000]$. Under the moment prior, we choose (α, σ, v, h) to yield the largest posterior probability $\Pr(\mathcal{K}|\mathbf{Y})$ under the BHM-MPP. We first utilize an outlier detection method introduced in Section A.3 of the Supplementary Materials to evaluate the data distribution. As there are significant amount of outliers in the data and thus the normal likelihood may not fit the data well, we adopt the t distribution to construct the likelihood in the BHM method. For comparison, we also implement other mean-change point detection methods, including ECP and DPML, and the popular pattern detection method called AutoPlait (Matsubara, Sakurai, and Faloutsos, 2014).

We use three metrics to compare different methods, while considering an estimator within (or outside) the m_I -neighbourhood of a true change point as a true positive (or false positive) detection.

Table 2: The running times based on Intel core i7-7700K CPU and detection results of BHM-MPP, ECP, DPMLE and AutoPlait on seven wind turbine datasets, and the overall results are the average of those for the seven datasets.

Metrics	BHM-MPP	ECP	DPMLE	AutoPlait	BHM-MPP	ECP	DPMLE	AutoPlait
	Dataset 1				Dataset 2			
Time (s)	161.1	21.9	435.2	2.8	32.9	6.3	58.6	1.2
Precision	0.625	0.316	0	0	0.500	0.300	0	1.000
Recall	0.833	1.000	0	0	0.750	0.750	0	0.250
F1 score	0.715	0.480	0	0	0.600	0.429	0	0.400
	Dataset 3				Dataset 4			
Time (s)	31.2	5.2	58.6	2.2	28.2	7.2	58.7	1.8
Precision	0.429	0.300	0.444	1.000	0.333	0.214	0	0
Recall	0.750	0.750	1.000	0.250	0.750	0.750	0	0
F1 score	0.545	0.429	0.615	0.400	0.462	0.333	0	0
	Dataset 5				Dataset 6			
Time (s)	70.3	13.3	156.2	2.8	30.3	5.8	58.8	2.6
Precision	0.375	0.176	0	0	0.875	0.667	0.727	0
Recall	0.750	0.750	0	0	0.875	1.000	1.000	0
F1 score	0.500	0.286	0	0	0.875	0.800	0.842	0
	Dataset 7				Overall			
Time (s)	15.8	2.0	13.0	1.4	52.8	8.8	119.9	2.1
Precision	0.429	0.375	0.429	0	0.509	0.322	0.441	0.500
Recall	0.750	0.750	0.750	0	0.794	0.853	0.441	0.059
F1 score	0.545	0.500	0.545	0	0.620	0.468	0.441	0.105

1. Precision (P): proportion of the estimated change points that are true change points.

If the method yields no estimated change point, the precision is defined as 0.

2. Recall (R): proportion of true change points detected by an algorithm.

3. F1 score: $F1 = 2RP/(R + P)$. When $R = P = 0$, the F1 score is defined as 0.

Table 2 shows the results from the BHM-MPP, ECP, DPMLE, AutoPlait methods on

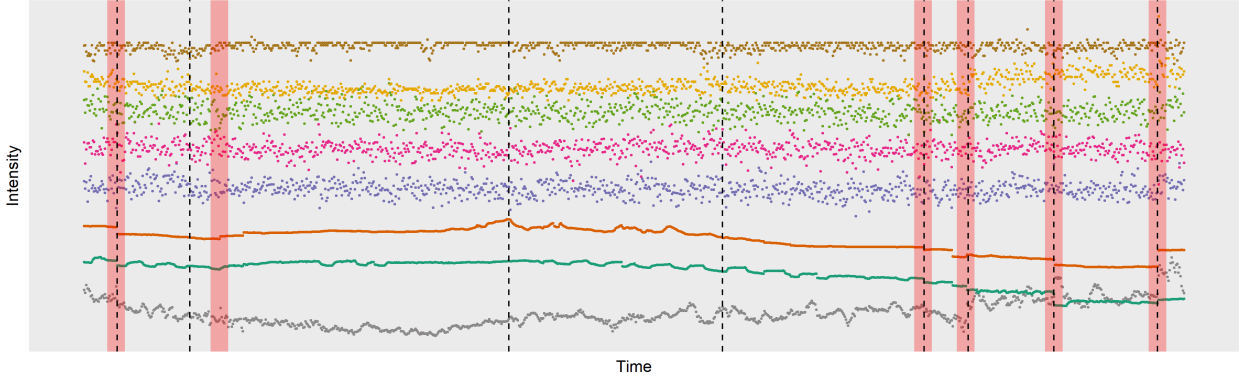


Figure 5: Detection results of the BHM-MPP method for the wind turbine dataset 1 with $n = 8$ sequences. The black dashed lines are the estimated state shift points and the red shadows are the m_I -neighbourhood of the ground truth.

the seven datasets. Among the four methods, BHM-MPP yields the best performance in five out of seven datasets (i.e., datasets 1, 2, 4, 5, 6) in terms of the F1 score. ECP yields a competitive recall but a much lower precision; DPMLE does not provide satisfactory precision or recall, especially for datasets 1, 2, 4, 5, where the algorithm misses all the change points; AutoPlait barely captures any true change points in all the datasets. Figure 5 shows the eight data sequences of the wind turbine dataset 1 and the detection results using the BHM-MPP method. Compared with the results of other methods in Figure 2, the BHM-MPP method clearly yields the best performance.

We also report the running time of the four methods for each dataset based on Intel core i7-7700k CPU in Table 2. While the AutoPlait leads to unsatisfactory performance, it consumes the shortest running time. Among the three mean-change detection methods, the nonparametric ECP is fastest. Due to the dynamic programming procedure, the computational burden of the BHM and DPMLE is relatively heavier, requiring longer com-

putational time compared to the other two methods. However, it is worth noting that all AutoPlait, ECP and DPMLE are implemented with C/C++ languages while our BHM method is implemented with the R language.

7 Conclusion

Motivated by the wind turbine data, we propose a BHM-based algorithm to detect mean changes for multivariate data sequences. Our method borrows the information across different sequences using the exchangeable random order prior. Further, BHM reduces the detection errors by applying the nonlocal priors to the mean difference. It also eases the computational burden by employing an initial screening stage for selecting the candidate points. We show the asymptotic consistency of the proposed method from both theoretical and numerical perspectives. As an illustration, we apply the BHM method to detect the anomalies in the wind turbine data where the BHM method shows robust outcomes and yields the best performance in most of the datasets in terms of the F1 score compared with other competitive methods considered in the paper.

Supplementary Materials

Supplementary Materials: PDF file containing additional simulation results, the dynamic programming algorithm as well as the proofs of Lemma 1 and Theorem 1.

RCode_and_data: Zip file containing the R code for implementing the BHM method as well as the real data used in the article in RData form.

References

- Barry, D. and Hartigan, J. A. (1993), “A Bayesian analysis for change point problems,” *Journal of the American Statistical Association*, 88, 309–319. [6](#)
- Bellman, R. and Roth, R. (1969), “Curve fitting by segmented straight lines,” *Journal of the American Statistical Association*, 64, 1079–1084. [3](#), [12](#)
- Bertolino, F., Racugno, W., and Moreno, E. (2000), “Bayesian model selection approach to analysis of variance under heteroscedasticity,” *Journal of the Royal Statistical Society: Series D (The Statistician)*, 49, 495–502. [7](#)
- Bouchard, D. and Badler, N. (2007), “Semantic segmentation of motion capture using laban movement analysis,” in *International Workshop on Intelligent Virtual Agents*, Springer. [8](#)
- Broderick, T., Jordan, M. I., and Pitman, J. (2013), “Cluster and Feature Modeling from Combinatorial Stochastic Processes,” *Statistical Science*, 28, 289–312. [10](#)
- Du, C., Kao, C. L. M., and Kou, S. C. (2016), “Stepwise Signal Extraction via Marginal Likelihood,” *Journal of the American Statistical Association*, 111, 314–330. [3](#), [6](#), [7](#), [12](#), [14](#)
- Fuentes-García, R., Mena, R. H., and Walker, S. G. (2019), “Modal posterior clustering motivated by Hopfield’s network,” *Computational Statistics & Data Analysis*, 137, 92–100. [6](#)

- Gharghabi, S., Ding, Y., Yeh, C. C. M., Kamgar, K., Ulanova, L., and Keogh, E. (2017), “Matrix profile VIII: Domain agnostic online semantic segmentation at superhuman performance levels,” *Proceedings - IEEE International Conference on Data Mining, ICDM*, 2017, 117–126. [8](#)
- Gnedin, A. and Pitman, J. (2006), “Exchangeable Gibbs partitions and Stirling triangles,” *Journal of Mathematical Sciences*, 138, 5674–5685. [10](#)
- Gong, D., Medioni, G., and Zhao, X. (2014), “Structured time series analysis for human action segmentation and recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36, 1414–1427. [8](#)
- Gray, C. S. and Watson, S. J. (2010), “Physics of failure approach to wind turbine condition based maintenance,” *Wind Energy*, 13, 395–405. [6](#)
- Hawkins, D. M. (2001), “Fitting multiple change-point models to data,” *Computational Statistics & Data Analysis*, 37, 323–341. [6](#)
- Hawkins, D. M., Qiu, P., and Kang, C. W. (2003), “The changepoint model for statistical process control,” *Journal of Quality Technology*, 35, 355–366. [7](#)
- Hinkley, D. V. (1971), “Inference about the change-point from cumulative sum tests,” *Biometrika*, 58, 509–523. [6](#)
- Jiang, F., Yin, G., and Dominici, F. (2018), “Bayesian model selection approach to boundary detection with non-local priors,” in *Advances in Neural Information Processing Systems*. [7](#)

- Johnson, V. E. and Rossell, D. (2010), “On the use of non-local prior densities in Bayesian hypothesis tests,” *Journal of the Royal Statistical Society: Series B*, 72, 143–170. [7](#)
- Kim, K., Parthasarathy, G., Uluycu, O., Foslien, W., Sheng, S., and Fleming, P. (2011), “Use of SCADA data for failure detection in wind turbines,” in *Energy Sustainability*, volume 54686. [5](#)
- Lau, J. W. and Green, P. J. (2007), “Bayesian model-based clustering procedures,” *Journal of Computational and Graphical Statistics*, 16, 526–558. [8](#)
- Lavielle, M. (2005), “Using penalized contrasts for the change-point problem,” *Signal Processing*, 85, 1501–1510. [7](#)
- Lijoi, A., Mena, R. H., and Prünster, I. (2007), “Controlling the reinforcement in Bayesian non-parametric mixture models,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69, 715–740. [10](#)
- Maboudou-Tchao, E. M. and Hawkins, D. M. (2013), “Detection of multiple change-points in multivariate data,” *Journal of Applied Statistics*, 40, 1979–1995. [3](#)
- Manogaran, G. and Lopez, D. (2018), “Spatial cumulative sum algorithm with big data analytics for climate change detection,” *Computers & Electrical Engineering*, 65, 207–221. [6](#)
- Martínez, A. F. and Mena, R. H. (2014), “On a Nonparametric Change Point Detection Model in Markovian Regimes,” *Bayesian Analysis*, 9, 823–858. [6](#), [8](#), [10](#), [11](#)

- Matsubara, Y., Sakurai, Y., and Faloutsos, C. (2014), “Autoplait: Automatic mining of co-evolving time sequences,” in *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*, ACM. [3](#), [8](#), [23](#)
- Matteson, D. S. and James, N. A. (2014), “A nonparametric approach for multiple change point analysis of multivariate data,” *Journal of the American Statistical Association*, [109](#), 334–345. [2](#)
- McCullagh, P. and Yang, J. (2008), “How many clusters?” *Bayesian Analysis*, [3](#), 1–19. [8](#)
- Muggeo, V. M. and Adelfio, G. (2010), “Efficient change point detection for genomic sequences of continuous measurements,” *Bioinformatics*, [27](#), 161–166. [6](#)
- Muñoz, C. Q. G., Jiménez, A. A., and Márquez, F. P. G. (2018), “Wavelet transforms and pattern recognition on ultrasonic guides waves for frozen surface state diagnosis,” *Renewable Energy*, [116](#), 42–54. [5](#)
- O’Hagan, A. and Leonard, T. (1976), “Bayes estimation subject to uncertainty about parameter constraints,” *Biometrika*, [63](#), 201–203. [17](#)
- Pitman, J. (1995), “Exchangeable and partially exchangeable random partitions,” *Probability Theory and Related Fields*, [102](#), 145–158. [8](#), [10](#)
- (2002), “Combinatorial stochastic processes,” Technical report, Dept. Statistics, UC Berkeley, 2002. [8](#), [11](#)
- Qiu, P. (2013), *Introduction to statistical process control*, CRC press. [7](#)

- Rigaill, G. (2015), “A pruned dynamic programming algorithm to recover the best segmentations with 1 to K_{\max} change-points.” *Journal de la Société Française de Statistique*, 156, 180–205. [6](#), [12](#), [14](#)
- Rosner, B. (1983), “Percentage points for a generalized ESD many-outlier procedure,” *Technometrics*, 25, 165–172. [10](#)
- Tautz-Weinert, J. and Watson, S. J. (2016), “Using SCADA data for wind turbine condition monitoring—a review,” *IET Renewable Power Generation*, 11, 382–394. [5](#)
- Truong, C., Oudre, L., and Vayatis, N. (2018), “A review of change point detection methods,” *arXiv preprint arXiv:1801.00718*. [7](#)
- Tsiamyrtzis, P. and Hawkins, D. M. (2005), “A Bayesian scheme to detect changes in the mean of a short-run process,” *Technometrics*, 47, 446–456. [7](#)
- (2008), “A Bayesian EWMA method to detect jumps at the start-up phase of a process,” *Quality and Reliability Engineering International*, 24, 721–735. [7](#)
- Wade, S., Ghahramani, Z., et al. (2018), “Bayesian cluster analysis: Point estimation and credible balls (with discussion),” *Bayesian Analysis*, 13, 559–626. [8](#)
- Wilkinson, M., Darnell, B., Van Delft, T., and Harman, K. (2014), “Comparison of methods for wind turbine condition monitoring with SCADA data,” *IET Renewable Power Generation*, 8, 390–397. [6](#)
- Yang, W., Court, R., and Jiang, J. (2013), “Wind turbine condition monitoring by the approach of SCADA data analysis,” *Renewable Energy*, 53, 365–376. [5](#)

- Yao, Y.-C. (1988), “Estimating the number of change-points via Schwarz’criterion,” *Statistics & Probability Letters*, 6, 181–189. 6
- Yao, Y.-C. and Au, S.-T. (1989), “Least-squares estimation of a step function,” *Sankhyā: The Indian Journal of Statistics, Series A*, 51, 370–381. 6
- Zamba, K. and Hawkins, D. M. (2006), “A multivariate change-point model for statistical process control,” *Technometrics*, 48, 539–549. 7
- Zhang, N. R. and Siegmund, D. O. (2007), “A modified Bayes information criterion with applications to the analysis of comparative genomic hybridization data,” *Biometrics*, 63, 22–32. 6
- Zou, C., Yin, G., Feng, L., and Wang, Z. (2014), “Nonparametric maximum likelihood approach to multiple change-point problems,” *The Annals of Statistics*, 42, 970–1002. 6

Supplementary Materials for “Bayesian Hierarchical Model for Change Point Detection in Multivariate Sequences”

May 10, 2021

A Additional numerical results

A.1 Variance parameter

We assess the impact of using the same ω_s across the sequences in our BHM model. We generate two-dimensional data from model I (homogeneous errors) and model II (heteroscedastic errors). Under each model, the data are generated with the same ($N_2(\mathbf{0}, \mathbf{I}_2)$) and distinct ($N_2(\mathbf{0}, \mathbf{\Sigma})$) variance parameters across the sequences, where \mathbf{I}_2 is an identity covariance matrix and $\mathbf{\Sigma} = \text{diag}(0.8, 1.2)$. We estimate the change point locations by assuming the two sequences share the same variance or use different variance parameters. The results in Table A.1 show that the two strategies yield similar performances across all scenarios. Hence, we suggest to use the same ω_s in practice which helps to reduce the computational burden and numerical errors as well as facilitating the information borrowing across the sequences.

Table A.1: Comparison results over 500 simulations when using the same or different variance parameters across the sequences with $n = 2$ under model I and model II, respectively. Standard deviations are given in parentheses.

Data-generating Model	Variance Parameter	$\hat{p} - p_0$							Segmentation Error	
		≤ -3	-2	-1	0	1	2	≥ 3	$d(\hat{\mathcal{K}} \mathcal{K}_0)$	$d(\mathcal{K}_0 \hat{\mathcal{K}})$
I, $(\xi_{1k}, \xi_{2k}) \sim N_2(\mathbf{0}, \mathbf{I}_2)$	Same	0	0	18	477	5	0	0	2.56 (2.83)	2.59 (3.86)
	Different	0	0	25	473	2	0	0	2.56 (2.83)	2.26 (2.92)
I, $(\xi_{1k}, \xi_{2k}) \sim N_2(\mathbf{0}, \mathbf{\Sigma})$	Same	0	0	23	474	3	0	0	2.62 (2.91)	2.43 (3.42)
	Different	0	0	25	472	3	0	0	2.62 (2.91)	2.36 (3.21)
II, $(\xi_{1k}, \xi_{2k}) \sim N_2(\mathbf{0}, \mathbf{I}_2)$	Same	0	0	3	341	135	21	0	1.23 (1.59)	6.32 (7.83)
	Different	0	0	3	345	137	14	0	1.23 (1.59)	6.27 (7.88)
II, $(\xi_{1k}, \xi_{2k}) \sim N_2(\mathbf{0}, \mathbf{\Sigma})$	Same	0	0	19	479	2	0	0	2.28 (2.74)	2.03 (2.82)
	Different	0	0	19	453	28	0	0	2.28 (2.74)	2.73 (4.07)

A.2 Parameter tuning in the priors of BHM

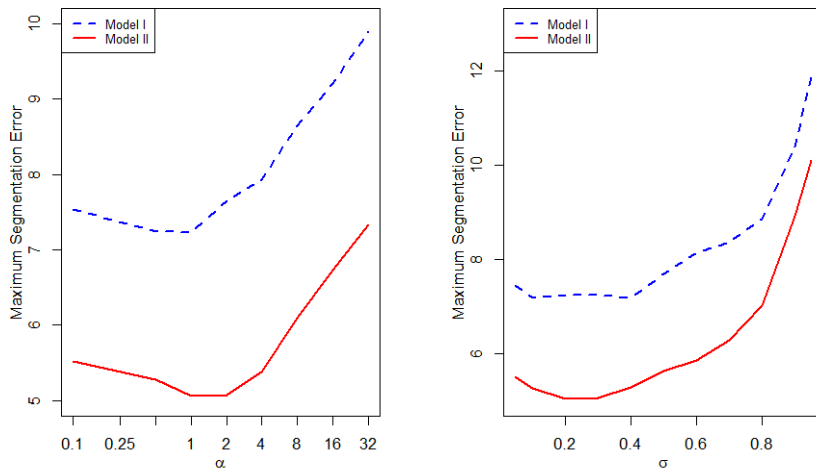


Figure A.1: The maximum segmentation error versus α (left) and σ (right) over 500 simulations with sample size $T = 400$, $n = 2$ under models I and II, respectively.

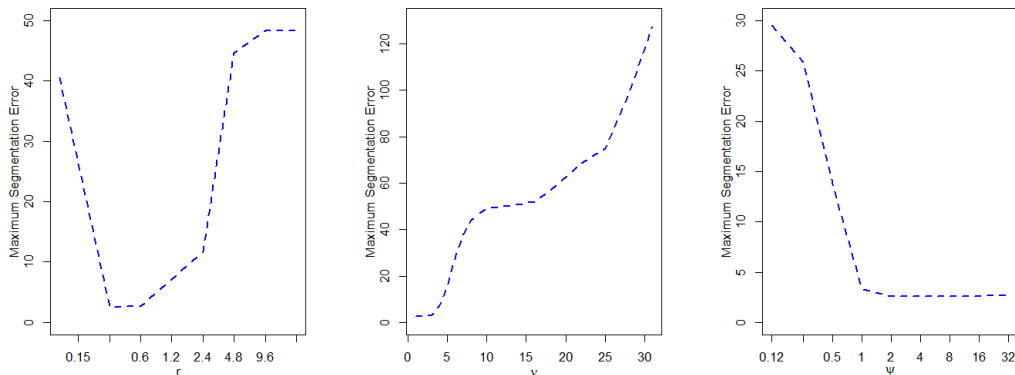


Figure A.2: The maximum segmentation errors versus the corresponding parameters for the inverse moment prior (left), moment prior (middle) and local prior (right) over 500 simulations with sample size $T = 400$, $n = 2$ under model I.

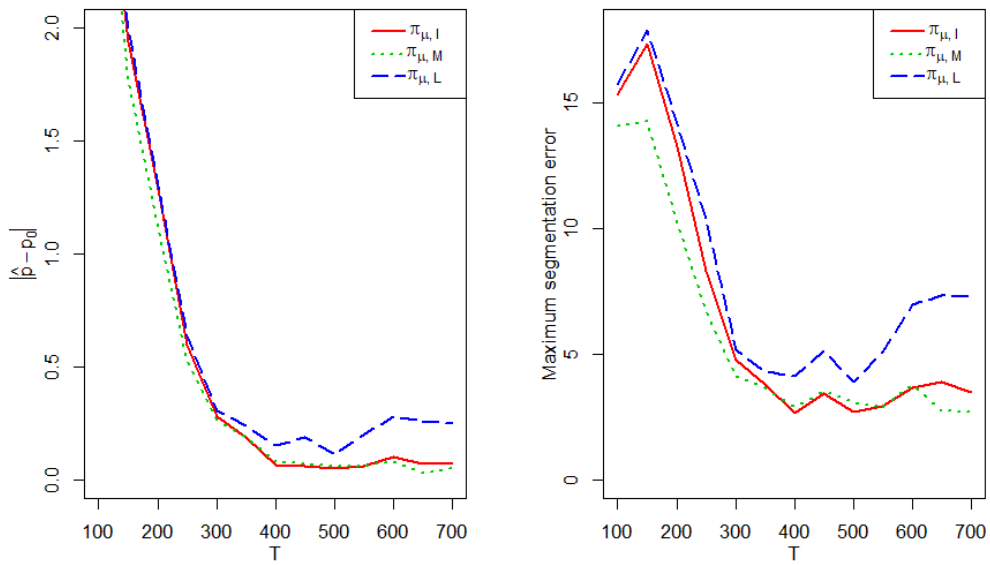


Figure A.3: The absolute difference $|\hat{p} - p_0|$ (left) and the maximum segmentation error (right) versus T over 500 simulations under three priors: the nonlocal inverse moment prior $\pi_{\mu,I}$ with $r = 0.6$, $\phi = q = 2$, nonlocal moment prior $\pi_{\mu,M}$ with $v = 1$ and local prior $\pi_{\mu,L}$ with $\psi = 2$ under model I with $n = 2$.

A.3 Determination of outliers

In the real data application, we determine the existence of the outliers as follows,

- (1) Select a candidate point set $\mathcal{H}(m_I)$ for the dataset with the screening algorithm in the BHM method.
- (2) Divide the dataset into segments based on the candidate point set and in this case, we can assume the data points in each segment has homogeneous distribution.
- (3) Conduct the generalized extreme studentized deviate (ESD) test (Rosner, 1983) for each segment.
- (4) If more than 12% of the segments contain outliers, we adopt the t likelihood; otherwise a normal likelihood is adopted.

In our experiments, this procedure works well as shown in Figure A.4.

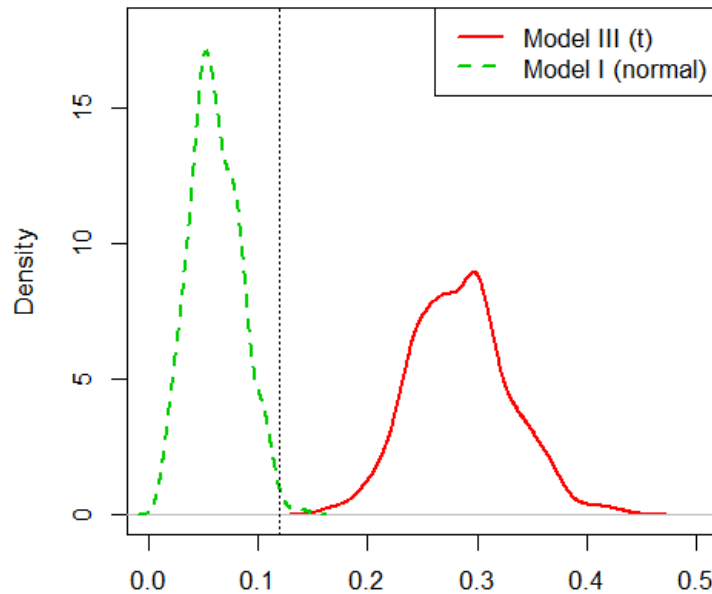


Figure A.4: The densities of proportions of segments containing outliers under model I (normal error) and model III (t error). The black dotted line indicates 12%.

A.4 Autocorrelation in the wind turbine data

While we have shown robustness of the BHM method for the moving-average error (refer to model V of Table ??), we also check the autocorrelation function (ACF) of the real data and simulate the datasets with ACF similar to the real data. The ACFs of dataset 1 in the wind turbine data are shown in Figure A.5. From the plots, it is clear that for the wind turbine data, the first three sequences show significant autocorrelation while the other five sequences are not significantly autocorrelated.

Thus, we conduct the simulation studies with a mixed error model where the first three sequences of the simulated dataset are with autoregressive (AR) errors while the other five sequences follows model I in Section ???. Specifically, we adopt the AR(2) model, i.e., $\xi_{ik} = 0.5\xi_{i,k-1} + 0.2\xi_{i,k-2} + a_k$ with $a_k \sim N(0, 3/5)$. The ACFs of the simulated dataset are shown in Figure A.6 which has similar patterns to Figure A.5. We use an independent normal likelihood in the BHM method and repeat the simulation for 500 times and the results are presented in Table A.2. Based on the results, the BHM-FIX and BHM-MPP methods still yield satisfactory results under the mixed error datasets. However, the nonparameteric ECP method deteriorates dramatically compared with the results in Table ???.

Table A.2: Comparison results over 500 simulations among BHM-FIX, BHM-MPP, ECP and DPMLE when $n = 8$ under the mixed error model. Standard deviations are given in parentheses.

Data-generating Model	Method	$\hat{p} - p_0$							Segmentation Error	
		≤ -3	-2	-1	0	1	2	≥ 3	$d(\hat{\mathcal{K}} \mathcal{K}_0)$	$d(\mathcal{K}_0 \hat{\mathcal{K}})$
Mixed errors	BHM-FIX	0	0	0	486	14	0	0	0.08 (0.31)	0.54 (2.80)
	BHM-MPP	0	0	0	484	16	0	0	0.03 (0.18)	0.52 (2.77)
	ECP	0	0	0	24	49	101	326	0.22 (0.53)	20.34 (6.59)
	DPMLE	500	0	0	0	0	0	0	59.79 (15.64)	0.12 (1.23)

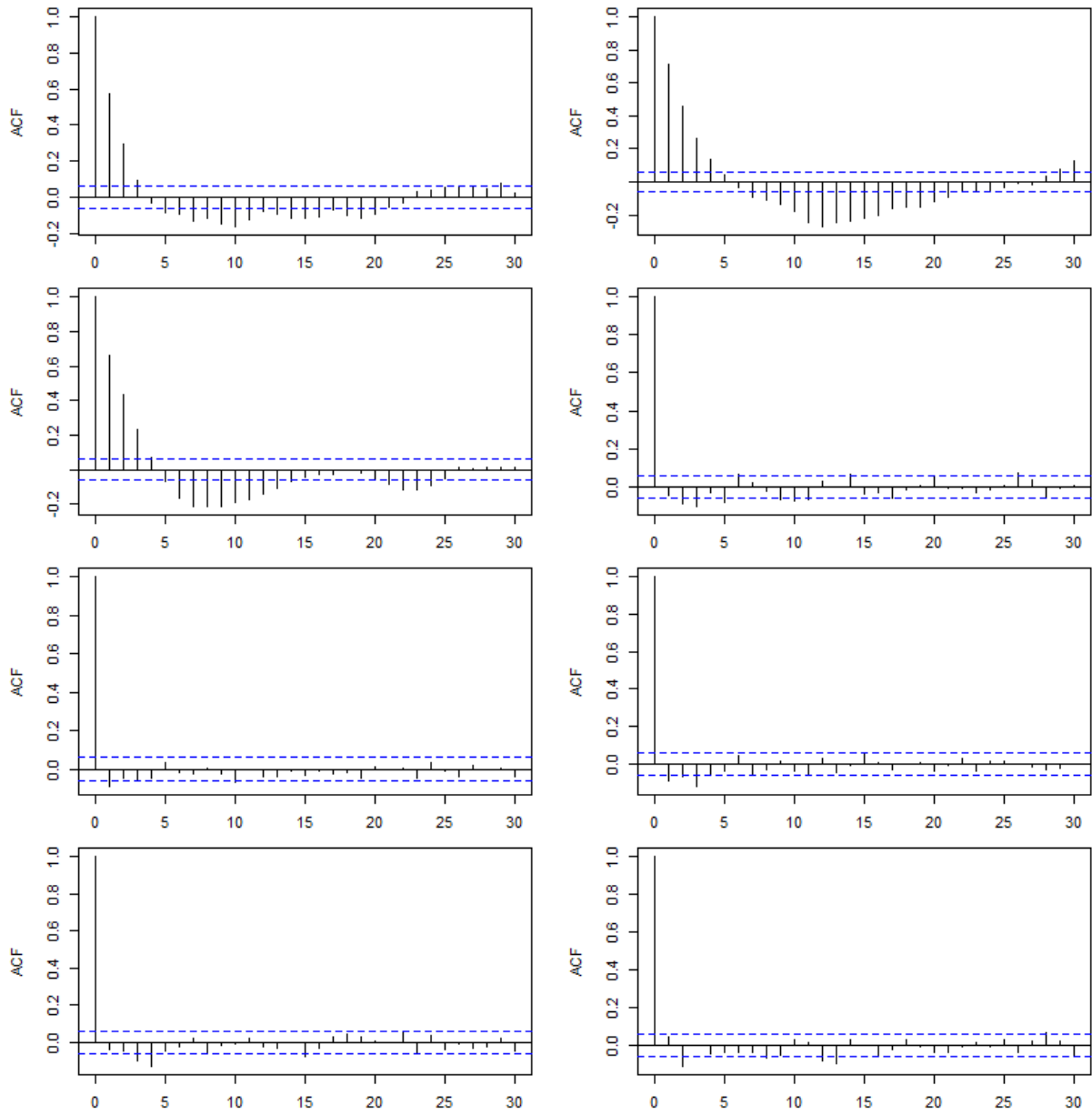


Figure A.5: The autocorrelation functions of dataset 1 in the wind turbine data.

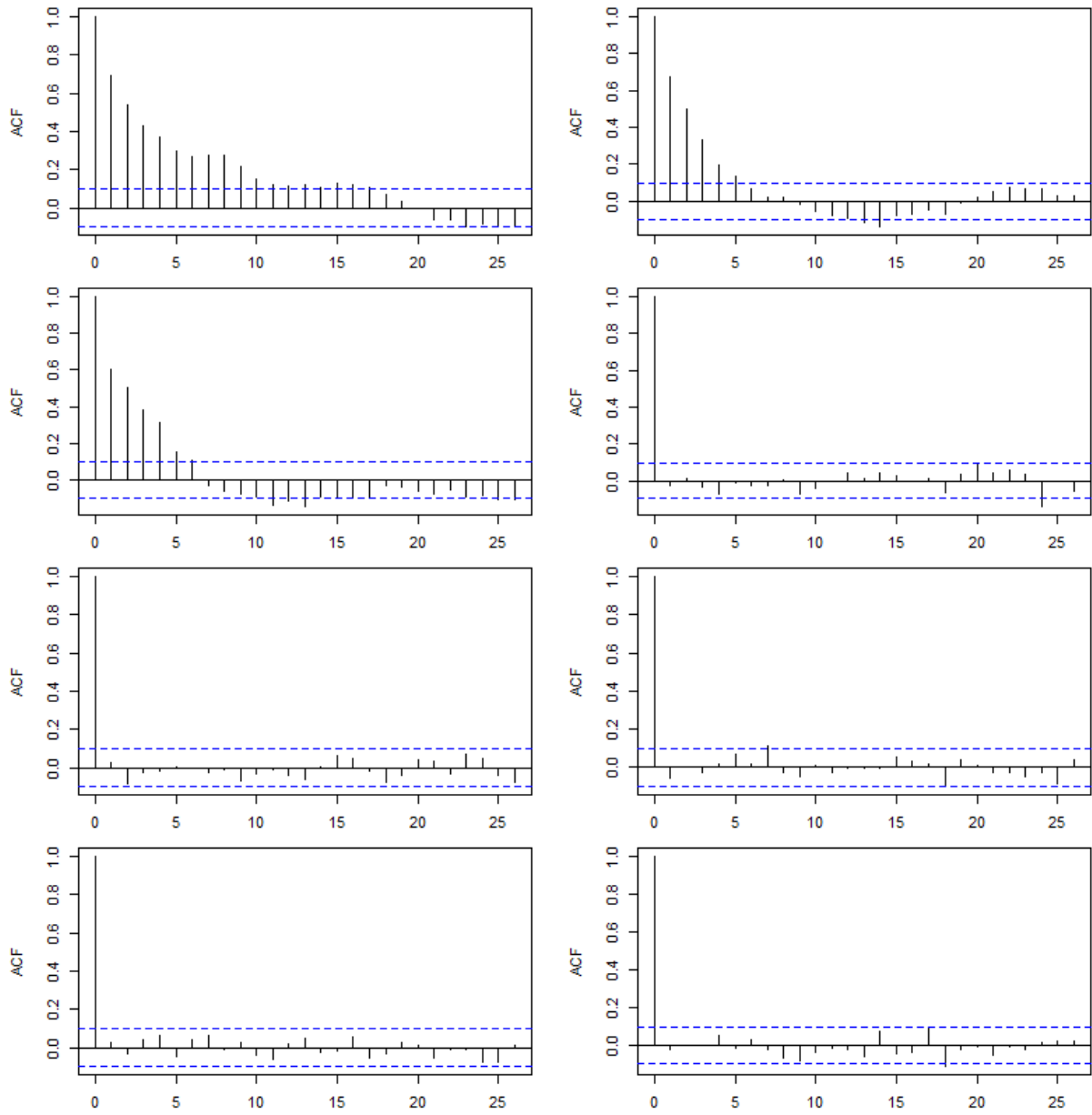


Figure A.6: The autocorrelation functions of the simulated data.

B Dynamic programming

For $l = 0, 1, \dots, N$, we define

$$H(j|l) \equiv \max_{\mathcal{K} \subseteq \{\tau_0, \tau_1, \dots, \tau_j, \tau_{j+1}\}, |\mathcal{K}|=l} U(\mathcal{K} | \mathbf{Y}_{(\tau_0, \tau_{j+1})}).$$

The dynamic programming algorithm is given as Algorithm A.1.

Algorithm A.1 Dynamic programming

Input:

The upper bound of the number of change points M , dataset \mathbf{Y} , candidate point set $\mathcal{H}(m_I)$.

- 1: Let A be an empty $M \times N$ matrix.
- 2: **for** $i = 0, \dots, N$ **do**
- 3: $H(i|0) \leftarrow u(\mathbf{Y}_{(\tau_0, \tau_{i+1})}, s = 0)$
- 4: **end for**
- 5: **for** $l = 1, \dots, M$ **do**
- 6: **for** $i = l, \dots, N$ **do**
- 7: $A_{l,i} \leftarrow \operatorname{argmax}_{l-1 \leq k \leq i-1} \{H(k|l-1)u(\mathbf{Y}_{(\tau_{k+1}, \tau_{i+1})}, s = l)\}$
- 8: $H(i|l) \leftarrow \max_{l-1 \leq k \leq i-1} \{H(k|l-1)u(\mathbf{Y}_{(\tau_{k+1}, \tau_{i+1})}, s = l)\}$
- 9: **end for**
- 10: **end for**
- 11: $\hat{p} \leftarrow \operatorname{argmax}_{l=0,1,\dots,M} H(N|l)$
- 12: **if** $\hat{p} = 0$ **then**
- 13: **return** \emptyset
- 14: **else**
- 15: $s \leftarrow \hat{p}$
- 16: $t \leftarrow N$
- 17: $E \leftarrow \emptyset$
- 18: **while** $s \neq 0$ **do**
- 19: $E \leftarrow E \cup \{A_{s,t} + 1\}$
- 20: $s \leftarrow s - 1$
- 21: $t \leftarrow A_{s,t}$
- 22: **end while**
- 23: **return** $\hat{\mathcal{K}} = \{\tau_i, i \in E\}$
- 24: **end if**

Output:

Estimated change point set $\hat{\mathcal{K}}$.

C Proofs

We denote \mathcal{K}_0 as the true change point set with p_0 change points, $\widehat{\mathcal{K}}$ as the estimated change point set with \widehat{p} estimated change points. The m_I -neighbourhood of a time point \mathbf{Y}_k is defined as $\{\mathbf{Y}_l : l \in (k - m_I, k + m_I)\}$. Given an interval $\mathbf{Y}_{(a,b]}$, we denote $p_{(a,b]}(\boldsymbol{\theta}) = \prod_{k \in (a,b]} f(\mathbf{Y}_k | \boldsymbol{\theta})$ as the likelihood function, the corresponding log-likelihood function is $l_{(a,b]}(\boldsymbol{\theta}) = \log p_{(a,b]}(\boldsymbol{\theta})$. We also let $\widehat{\boldsymbol{\theta}}_{(a,b]}$ be the maximum likelihood estimator (MLE) based on $l_{(a,b]}(\boldsymbol{\theta})$, and $\boldsymbol{\theta}_{(a,b]}$ be the true parameters on $(a, b]$. We denote $\widehat{\sigma}_{(a,b]}^2 = \{-E(\frac{\partial^2 l_{(a,b]}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top})|_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_{(a,b]}}\}^{-1}$ and let $J(\boldsymbol{\theta}_0) = -E\frac{\partial^2 \log f(\mathbf{Y}_1 | \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$ be the Fisher information for one observation. We also define d as the dimension of $\boldsymbol{\theta}$. Finally, denote

$$C(\mathbf{Y}_{(\kappa_s, \kappa_{s+1}]}) = \int p_{(\kappa_s, \kappa_{s+1}]}(\boldsymbol{\theta}_s) \pi(\boldsymbol{\theta}_s) d\boldsymbol{\theta}_s.$$

We list the regularity conditions as follows.

- (1) The prior for mean difference $\pi_\mu(\mu)$ is continuous with bounded first and second derivatives.
- (2) For a segment between two true change points κ_{0j} and $\kappa_{0,j+1}$ ($j = 1, \dots, p_0$) with parameters $(\mu_{1,j}, \dots, \mu_{n,j})$, there exists $\delta_I > 0$ such that for any $i \in \{1, \dots, n\}$, $|\mu_{i,j}|$ is either greater than δ_I or equal to 0. Further, there is $i \in \{1, \dots, n\}$ such that $|\mu_{i,j}| > \delta_I$.
- (3) The generic prior $\pi(\boldsymbol{\theta})$ is continuous and positive at all $\boldsymbol{\theta}_i$ ($i = 0, \dots, p_0$), where $\boldsymbol{\theta}_i$ is the true parameters for interval $(\kappa_{0i}, \kappa_{0,i+1}]$.
- (4) The regularity conditions (A1)–(A5) and (B1)–(B4).

Regularity conditions (A1)–(A5) and (B1)–(B4) are listed as follows. All the conditions are multivariate extensions from (Du et al., 2016).

- (A1) Θ is a closed set, and $\Theta \subseteq \mathcal{R}^d$.
- (A2) The set of points $\{\mathbf{x} : f(\mathbf{x} | \boldsymbol{\theta}) > 0\}$ is independent of $\boldsymbol{\theta}$. We denote this set by \mathcal{X} .
- (A3) If $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2$ are two distinct points in Θ , then the Lebesgue measure of $\mu\{\mathbf{x} : f(\mathbf{x} | \boldsymbol{\theta}_1) \neq f(\mathbf{x} | \boldsymbol{\theta}_2)\} > 0$.
- (A4) Let $\mathbf{x} \in \mathcal{X}$, $\boldsymbol{\theta}' \in \Theta$. Then for all $\boldsymbol{\theta}$ such that $\|\boldsymbol{\theta} - \boldsymbol{\theta}'\| < \delta$, where $\|\cdot\|$ is the L_2 norm, with δ sufficiently small,

$$|\log f(\mathbf{x} | \boldsymbol{\theta}) - \log f(\mathbf{x} | \boldsymbol{\theta}')| < H_\delta(\mathbf{x}, \boldsymbol{\theta}'),$$

where

$$\lim_{\delta \rightarrow 0} H_\delta(\mathbf{x}, \boldsymbol{\theta}') = 0,$$

and for any $\boldsymbol{\theta}_0 \in \Theta$,

$$\lim_{\delta \rightarrow 0} \int_{\mathcal{X}} H_\delta(\mathbf{x}, \boldsymbol{\theta}') f(\mathbf{x} | \boldsymbol{\theta}_0) d\mu = 0.$$

(A5) If Θ is not bounded, then for any $\boldsymbol{\theta}_0 \in \Theta$, and sufficiently large Δ ,

$$\log f(\mathbf{x}|\boldsymbol{\theta}) - \log f(\mathbf{x}|\boldsymbol{\theta}_0) < K_\Delta(\mathbf{x}, \boldsymbol{\theta}_0),$$

whenever $\|\boldsymbol{\theta}\| > \Delta$, where

$$\lim_{\Delta \rightarrow \infty} \int_{\mathcal{X}} K_\Delta(\mathbf{x}, \boldsymbol{\theta}_0) f(\mathbf{x}|\boldsymbol{\theta}_0) d\mu < 0.$$

(B1) $\log f(\mathbf{x}|\boldsymbol{\theta})$ is twice differentiable with respect to $\boldsymbol{\theta}$ in some neighborhood of $\boldsymbol{\theta}_0$.

(B2) Let

$$J(\boldsymbol{\theta}_0) = \left[\int_{\mathcal{X}} f_0 \frac{\partial \log f_0}{\partial \theta_{0i}} \frac{\partial \log f_0}{\partial \theta_{0j}} d\mu \right]_{i,j=1}^d,$$

where f_0 denotes $f(\mathbf{x}|\boldsymbol{\theta}_0)$, θ_{0i} is the i th element of $\boldsymbol{\theta}_0$. Then $J(\boldsymbol{\theta}_0)$ is positive definite.

(B3) For any $1 \leq i, j \leq d$,

$$\int_{\mathcal{X}} \frac{\partial f_0}{\partial \theta_{0i}} d\mu = \int_{\mathcal{X}} \frac{\partial^2 f_0}{\partial \theta_{0i} \partial \theta_{0j}} d\mu = 0.$$

(B4) For $\delta > 0$, if $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| < \delta$, where δ is small enough, then

$$\left\| \frac{\partial^2 \log f(\mathbf{x}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} - \frac{\partial^2 \log f(\mathbf{x}|\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right\| < M_\delta(\mathbf{x}, \boldsymbol{\theta}_0),$$

where $\lim_{\delta \rightarrow 0} \int M_\delta(\mathbf{x}, \boldsymbol{\theta}_0) f(\mathbf{x}|\boldsymbol{\theta}_0) d\mu = 0$.

Proof of Lemma ??:

We define j as an m_I -flat point if there is no change point in $(j - m_I, j + m_I)$. Let \mathcal{F} be the set of all m_I -flat points. So $|\mathcal{F}| = T - p_0(2m_I - 1)$, where $|\mathcal{F}|$ denotes the cardinality of set \mathcal{F} . To prove Lemma ??, it is sufficient to show

$$\Pr \left(\min_{k \in \mathcal{K}_0} R_k > \max_{l \in \mathcal{F}} R_l \right) \rightarrow 1,$$

as $T \rightarrow \infty$. Note that

$$\Pr \left(\min_{k \in \mathcal{K}_0} R_k > \max_{l \in \mathcal{F}} R_l \right) \geq \Pr \left(\min_{k \in \mathcal{K}_0} R_k > b_T > \max_{l \in \mathcal{F}} R_l \right),$$

where b_T is a positive sequence with respect to T . It follows that

$$\begin{aligned} & \Pr \left(\min_{k \in \mathcal{K}_0} R_k > b_T > \max_{l \in \mathcal{F}} R_l \right) \\ &= \Pr \{ (\cap_{k \in \mathcal{K}_0} \{R_k > b_T\}) \cap (\cap_{l \in \mathcal{F}} \{R_l < b_T\}) \} \\ &= 1 - \Pr \{ (\cup_{k \in \mathcal{K}_0} \{R_k \leq b_T\}) \cup (\cup_{l \in \mathcal{F}} \{R_l \geq b_T\}) \} \\ &\geq 1 - \{ \Pr(\cup_{k \in \mathcal{K}_0} \{R_k \leq b_T\}) + \Pr(\cup_{l \in \mathcal{F}} \{R_l \geq b_T\}) \} \\ &\geq 1 - \left\{ \sum_{k \in \mathcal{K}_0} \Pr(\{R_k \leq b_T\}) + \sum_{l \in \mathcal{F}} \Pr(\{R_l \geq b_T\}) \right\}. \end{aligned}$$

We define

$$R_{ij} = \frac{\int \prod_{l=j+1}^{j+m_I} \exp\{-(Y_{il} - \bar{Y}_{ij} - \mu)^2\} \pi(\mu) d\mu}{\prod_{l=j+1}^{j+m_I} \exp\{-(Y_{il} - \bar{Y}_{ij})^2\}}, \text{ for } i = 1, \dots, n,$$

where $\bar{Y}_{ij} = m_I^{-1} \sum_{l=j-m_I+1}^j Y_{il}$. Clearly, $R_j = \prod_{i=1}^n R_{ij}$.

For any change point $k \in \mathcal{K}_0$, assume n_x sequences have mean shifts at this change point. By regularity condition (2), we know $n_x \geq 1$ and the absolute change of mean is greater than δ_I .

Without loss of generality, assume the first n_x sequences have mean changes. By Lemma 1 of [Jiang, Yin, and Dominici \(2018\)](#), we have

$$\lim_{T \rightarrow \infty} \Pr\{R_{ik} > \exp(Dm_I \delta_I)\} = 1, \quad (1)$$

when there is a mean shift in sequence i at change point k , where $D > 0$ is a constant. Then we set

$$b_T = \exp(D\delta_I m_I / 2).$$

For any $l \in \mathcal{F}$, by Lemmas 2, 3, 4 of [Jiang, Yin, and Dominici \(2018\)](#), there exist $c, C > 0$ such that

$$ca_T \leq R_{il} \leq Ca_T, \quad (2)$$

so

$$R_l = \prod_{i=1}^n R_{il} = O_p(a_T^n),$$

where $a_T = m_I^{-1/2}$, $m_I^{-v-1/2}$ and $\exp(-m_I^{s/(s+1)})$ correspond to the local prior, moment prior and inverse moment prior. Consequently, we have

$$\begin{aligned} \Pr(R_l \geq b_T) &= O\{a_T^n \exp(-D\delta_I m_I / 2)\}, \\ \sum_{l: l_i \in \mathcal{F}} \Pr(R_l \geq b_T) &= O\{T a_T^n \exp(-D\delta_I m_I / 2)\} = o(1), \end{aligned} \quad (3)$$

since $m_I / (\log T)^{1+\epsilon} \rightarrow c > 0$.

Next, for $k \in \mathcal{K}_0$, by (1), we know for $i = 1, \dots, n_x$, we have

$$\lim_{T \rightarrow \infty} \Pr\{R_{ik} > \exp(Dm_I \delta_I)\} = 1.$$

As a result,

$$\lim_{T \rightarrow \infty} \Pr\left\{\prod_{i=1}^{n_x} R_{ik} > \exp(n_x D m_I \delta_I)\right\} = 1. \quad (4)$$

Consequently, we obtain

$$\begin{aligned}
& \Pr(R_k \leq b_T) \\
&= \Pr\left(\prod_{i=1}^{n_x} R_{ik} \prod_{j=n_x+1}^n R_{jk} \leq b_T\right) \\
&= \Pr\left\{\prod_{i=1}^{n_x} R_{ik} \prod_{j=n_x+1}^n R_{jk} \leq b_T, \prod_{i=1}^{n_x} R_{ik} > \exp(n_x D \delta_I m_I)\right\} \\
&\quad + \Pr\left\{\prod_{i=1}^{n_x} R_{ik} \prod_{j=n_x+1}^n R_{jk} \leq b_T, \prod_{i=1}^{n_x} R_{ik} \leq \exp(n_x D \delta_I m_I)\right\} \\
&\leq \Pr\left\{\exp(n_x D \delta_I m_I) \prod_{j=n_x+1}^n R_{jk} \leq b_T\right\} + \Pr\left\{\prod_{i=1}^{n_x} R_{ik} \leq \exp(n_x D \delta_I m_I)\right\}.
\end{aligned}$$

Combining with (4),

$$\lim_{T \rightarrow \infty} \Pr(R_k \leq b_T) \leq \lim_{T \rightarrow \infty} \Pr\left\{\exp(n_x D m_I \delta_I) \prod_{j=n_x+1}^n R_{jk} \leq b_T\right\}.$$

For $j = n_x + 1, \dots, n$, by (2), $\exists c_1, C_1 > 0$ such that

$$c_1 a_T^{n-n_x} \leq \prod_{j=n_x+1}^n R_{jk} \leq C_1 a_T^{n-n_x},$$

and

$$C_1^{-1} a_T^{n_x-n} \leq \left(\prod_{j=n_x+1}^n R_{jk}\right)^{-1} \leq c_1^{-1} a_T^{n_x-n}.$$

This implies

$$\begin{aligned}
& \Pr\left\{\exp(n_x D m_I \delta_I) \prod_{j=n_x+1}^n R_{jk} \leq b_T\right\} \\
&= \Pr\left\{\left(\prod_{j=n_x+1}^n R_{jk}\right)^{-1} \geq \exp(n_x D m_I \delta_I) b_T^{-1}\right\} \\
&\leq \frac{c_1^{-1} a_T^{n_x-n}}{\exp(n_x D m_I \delta_I) b_T^{-1}} \\
&= c_1^{-1} a_T^{n_x-n} \exp(-n_x D m_I \delta_I) b_T,
\end{aligned}$$

where the second to the last inequality holds by the Markov inequality. Therefore

$$\Pr(R_k \leq b_T) = O\{a_T^{n_x-n} \exp(-n_x D m_I \delta_I) b_T\}.$$

Thus, we obtain

$$\sum_{k \in \mathcal{K}_0} \Pr(R_k \leq b_T) = O[p_0 a_T^{n_x - n} \exp\{-(n_x - 1/2) D m_I \delta_I\}] = o(1), \quad (5)$$

since $m_I / (\log T)^{1+\epsilon} \rightarrow c > 0$. Then using (3) and (5), we achieve:

$$\begin{aligned} & \Pr\left(\min_{k \in \mathcal{K}_0} R_k > \max_{l \in \mathcal{F}} R_l\right) \\ & \geq \Pr\left(\min_{k \in \mathcal{K}_0} R_k > b_T > \max_{l \in \mathcal{F}} R_l\right) \\ & = \Pr\left\{(\cap_{k \in \mathcal{K}_0} \{R_k > b_T\}) \cap (\cap_{l \in \mathcal{F}} \{R_l < b_T\})\right\} \\ & \geq 1 - \left\{\sum_{k \in \mathcal{K}_0} \Pr(\{R_k \leq b_T\}) + \sum_{l \in \mathcal{F}} \Pr(\{R_l \geq b_T\})\right\} \\ & = 1 - o(1). \end{aligned}$$

Finally, we know

$$\Pr\left(\min_{k \in \mathcal{K}_0} R_k > \max_{l \in \mathcal{F}} R_l\right) \rightarrow 1,$$

as $T \rightarrow \infty$. □

Lemma 3. *Under conditions (A1)–(A5) and (B1)–(B4), if there is no change point in the interval $(a, b]$, and the true value of parameter within this segment is $\boldsymbol{\theta}_{(a,b]}$, then as $(b - a) \rightarrow \infty$,*

1. *Let $N_0(\delta) = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}_{(a,b]}\| < \delta\}$ be a neighborhood of $\boldsymbol{\theta}_{(a,b]}$ contained in Θ , the parameter space, there exists a positive number $k_{\boldsymbol{\theta}_{(a,b]}}(\delta)$, depending on $\boldsymbol{\theta}_{(a,b]}$ and δ , such that*

$$\lim_{(b-a) \rightarrow \infty} \Pr\left\{\sup_{\boldsymbol{\theta} \notin N_0(\delta)} \frac{l_{(a,b]}(\boldsymbol{\theta}) - l_{(a,b]}(\boldsymbol{\theta}_{(a,b]})}{b - a} < -k_{\boldsymbol{\theta}_{(a,b]}}(\delta)\right\} = 1;$$

2. $l_{(a,b]}(\boldsymbol{\theta}_{(a,b]}) - l_{(a,b]}(\widehat{\boldsymbol{\theta}}_{(a,b]}) = O_p(1)$.

Proof:

The proof of Lemma 3 is a direct multi-dimensional extension from Theorem 1 of Walker (1969). □

The following result is Theorem 3.1 of Fraser and McDunnough (1984), and the regularity conditions (A1)–(A5) and (B1)–(B4) imply the three assumptions in Fraser and McDunnough (1984).

Lemma 4. *Suppose conditions (A1)–(A5) and (B1)–(B4) hold. For i.i.d samples $\{\mathbf{Y}_1, \dots, \mathbf{Y}_T\}$ from $f(\mathbf{Y}|\boldsymbol{\theta}_0)$, let $\hat{\sigma}^2 = \{-E(\frac{\partial^2 \log p(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}\}^{-1}$ and $p(\boldsymbol{\theta}) = \prod_{k:t_k \in (0,1]} f(\mathbf{Y}_k|\boldsymbol{\theta})$, where $\hat{\boldsymbol{\theta}}$ is*

the MLE of $\boldsymbol{\theta}_0$. If $w(\boldsymbol{\theta}) \geq 0$ for all $\boldsymbol{\theta} \in \Theta$ and satisfies $\int w(\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} < \infty$ and is continuous and nonzero at the true $\boldsymbol{\theta}_0$, then

$$\frac{\det(\hat{\sigma})w(\hat{\boldsymbol{\theta}})p(\hat{\boldsymbol{\theta}})}{\int w(\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}} \xrightarrow{a.s.} (2\pi)^{-d/2}.$$

Lemma 5. Assume conditions in Theorem ?? hold. Suppose that there are r change points in (a, b) , say $\{\kappa_1, \dots, \kappa_r\}$, with $\kappa_1 < \dots < \kappa_r$. Further assume $(\kappa_{i+1} - \kappa_i) \rightarrow \infty, i = 0, \dots, r$ (let $\kappa_0 = a, \kappa_{r+1} = b$) as $(b - a) \rightarrow \infty$. Then let $\underline{\kappa} = \min_{i=0, \dots, r} (\kappa_{i+1} - \kappa_i), \exists c_2 > 0$ such that

$$\frac{C(\mathbf{Y}_{(a,b]})}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} = O_p \left\{ (b - a)^{rd/2} \exp(-\underline{\kappa}c_2) \right\}.$$

Proof:

The r change points separate the sequences into $r + 1$ segments. We first assume all the $r + 1$ segments have different parameters denoted as $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{r+1}$. Then we can find a δ and define $N_i(\delta) = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}_i\| < \delta\}, i = 1, \dots, r + 1$ such that $N_i(\delta) \cap N_j(\delta) = \emptyset$ for $i \neq j$. We write

$$C(\mathbf{Y}_{(a,b]}) = \sum_{i=0}^{r+1} I_i$$

where

$$\begin{aligned} I_i &= \int_{N_i(\delta)} p_{(a,b]}(\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}, \quad \text{for } i = 1, \dots, r + 1, \\ I_0 &= \int_{\Theta - \cup_{i=1}^{r+1} N_i(\delta)} p_{(a,b]}(\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}. \end{aligned}$$

By Lemma 4, we have

$$\begin{aligned} C(\mathbf{Y}_{(\kappa_{i-1}, \kappa_i]}) &= p_{(\kappa_{i-1}, \kappa_i]}(\hat{\boldsymbol{\theta}}_{(\kappa_{i-1}, \kappa_i]})\pi(\hat{\boldsymbol{\theta}}_{(\kappa_{i-1}, \kappa_i]}) \det(\hat{\sigma}_{(\kappa_{i-1}, \kappa_i]}) O_p(1) \\ &\neq p_{(\kappa_{i-1}, \kappa_i]}(\hat{\boldsymbol{\theta}}_{(\kappa_{i-1}, \kappa_i]})\pi(\hat{\boldsymbol{\theta}}_{(\kappa_{i-1}, \kappa_i]}) \det(\hat{\sigma}_{(\kappa_{i-1}, \kappa_i]}) o_p(1). \end{aligned} \quad (6)$$

where $i = 1, \dots, r + 1$. Note that by definition of $\hat{\sigma}_{(\kappa_{i-1}, \kappa_i]}$,

$$\begin{aligned} \det(\hat{\sigma}_{(\kappa_{i-1}, \kappa_i]}) &= O_p\{(\kappa_i - \kappa_{i-1})^{-d/2}\}, \\ \det(\hat{\sigma}_{(\kappa_{i-1}, \kappa_i]}) &\neq o_p\{(\kappa_i - \kappa_{i-1})^{-d/2}\}. \end{aligned} \quad (7)$$

By (6), for $j \neq 0$, we obtain

$$\begin{aligned}
& \frac{I_j}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} \\
&= \frac{\int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} \\
&= \frac{O_p(1) \int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]}) \prod_{i \neq j} p_{(\kappa_{i-1},\kappa_i]}(\widehat{\boldsymbol{\theta}}_{(\kappa_{i-1},\kappa_i]}) \pi(\widehat{\boldsymbol{\theta}}_{(\kappa_{i-1},\kappa_i]}) \det(\widehat{\sigma}_{(\kappa_{i-1},\kappa_i]})} \\
&= \frac{O_p(1) \int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]}) \prod_{i \neq j} p_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}_i) \pi(\widehat{\boldsymbol{\theta}}_{(\kappa_{i-1},\kappa_i]}) \det(\widehat{\sigma}_{(\kappa_{i-1},\kappa_i]})} \\
&= \frac{\int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]}) \prod_{i \neq j} p_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}_i) \pi(\widehat{\boldsymbol{\theta}}_{(\kappa_{i-1},\kappa_i]}) O_p\{(\kappa_i - \kappa_{i-1})^{-d/2}\}} \\
&= \frac{\int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]}) \prod_{i \neq j} p_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}_i) \pi(\boldsymbol{\theta}_i) O_p\{(\kappa_i - \kappa_{i-1})^{-d/2}\}}, \tag{8}
\end{aligned}$$

where the third equality in the above equation is due to the second result of Lemma 3, the fourth is by (7), and the last one is by the continuous mapping theorem.

Using the first result of Lemma 3, $\exists k(\delta) > 0$ such that

$$\begin{aligned}
& \frac{\int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]}) \prod_{i \neq j} p_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}_i) \pi(\boldsymbol{\theta}_i)} \\
&= \frac{1}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]})} \int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_{j-1},\kappa_j]}(\boldsymbol{\theta}) \prod_{i \neq j} \exp\{l_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}) - l_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}_i)\} d\boldsymbol{\theta} \\
&< \frac{1}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]})} \prod_{i \neq j} \exp\{-(\kappa_i - \kappa_{i-1})k(\delta)\} \int_{N_j(\delta)} p_{(\kappa_{j-1},\kappa_j]}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \\
&\leq \frac{1}{C(\mathbf{Y}_{(\kappa_{j-1},\kappa_j]})} \prod_{i \neq j} \exp\{-(\kappa_i - \kappa_{i-1})k(\delta)\} \int_{\Theta} p_{(\kappa_{j-1},\kappa_j]}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \\
&= \prod_{i \neq j} \exp\{-(\kappa_i - \kappa_{i-1})k(\delta)\} \tag{9}
\end{aligned}$$

with probability tending to unit as $(b - a) \rightarrow \infty$. Combining (8) and (9), we achieve

$$\begin{aligned}
& \frac{I_j}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} \\
&= O_p \left[\prod_{i \neq j} (\kappa_i - \kappa_{i-1})^{d/2} \exp\{-(\kappa_i - \kappa_{i-1})k(\delta)\} \right] \\
&= O_p \left\{ (b - a)^{rd/2} \exp(-\underline{\kappa}k(\delta)) \right\}.
\end{aligned}$$

For I_0 , we apply the same argument, but note that the region $\Theta - \cup_{i=1}^{r+1} N_i(\delta)$ does not

contain the neighborhood of any $\boldsymbol{\theta}_i$, so

$$\frac{I_0}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})}$$

would have a faster convergence rate compared with

$$\frac{I_j}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})}.$$

Thus we achieve

$$\frac{C(\mathbf{Y}_{(a,b]})}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} = \frac{\sum_{i=0}^{r+1} I_i}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} = O_p \left\{ (b-a)^{rd/2} \exp(-\underline{\kappa}k(\delta)) \right\}.$$

If some segments share the same parameters, without loss of generality, we assume only $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_3$ then $N_1(\delta) = N_3(\delta)$. The argument is analogous when more than two segments share the same parameters.

For $j \neq 1$ or 3 , the argument is identical to the above. When $j = 1$ (and there is no I_3 , since $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_3$), following similar discussions for (8) and (9),

$$\begin{aligned} & \frac{I_1}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} \\ &= \frac{\int_{N_j(\delta)} \pi(\boldsymbol{\theta}) p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) \cdots p_{(\kappa_r,\kappa_{r+1}]}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_0,\kappa_1]}) C(\mathbf{Y}_{(\kappa_2,\kappa_3]}) \prod_{i \neq 1,3} p_{(\kappa_{i-1},\kappa_i]}(\boldsymbol{\theta}_i) \pi(\boldsymbol{\theta}_i) O_p \left\{ (\kappa_i - \kappa_{i-1})^{-d/2} \right\}} \\ &\leq \frac{\int_{\Theta} p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) p_{(\kappa_2,\kappa_3]}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_0,\kappa_1]}) C(\mathbf{Y}_{(\kappa_2,\kappa_3]})} \prod_{i \neq 1,3} \exp \left\{ -(\kappa_i - \kappa_{i-1})k(\delta) \right\} O_p \left\{ (\kappa_i - \kappa_{i-1})^{d/2} \right\}. \end{aligned} \quad (10)$$

Using similar discussion for (15),

$$\frac{\int_{\Theta} p_{(\kappa_0,\kappa_1]}(\boldsymbol{\theta}) p_{(\kappa_2,\kappa_3]}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}}{C(\mathbf{Y}_{(\kappa_0,\kappa_1]}) C(\mathbf{Y}_{(\kappa_2,\kappa_3]})} = O_p \left\{ \frac{(\kappa_3 - \kappa_2)(\kappa_1 - \kappa_0)}{(\kappa_3 - \kappa_2) + (\kappa_1 - \kappa_0)} \right\}^{d/2}. \quad (11)$$

Combining (10) and (11), we achieve

$$\begin{aligned} & \frac{I_1}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} \\ &= O_p \left\{ \frac{\prod_{i=1}^{r+1} (\kappa_i - \kappa_{i-1})}{(\kappa_3 - \kappa_2) + (\kappa_1 - \kappa_0)} \right\}^{d/2} O_p \left[\prod_{i \neq 1,3} \exp \left\{ -(\kappa_i - \kappa_{i-1})k(\delta) \right\} \right] \\ &= O_p \left\{ (b-a)^{rd/2} \exp(-\underline{\kappa}k(\delta)) \right\}. \end{aligned}$$

Then we obtain

$$\frac{C(\mathbf{Y}_{(a,b]})}{C(\mathbf{Y}_{(a,\kappa_1]}) \cdots C(\mathbf{Y}_{(\kappa_r,b]})} = O_p \left\{ (b-a)^{rd/2} \exp(-\underline{\kappa}k(\delta)) \right\}.$$

□

Lemma 6. Assume conditions in Theorem ?? hold, and \mathcal{K}_0 is a subset of $\mathcal{H}(m_I)$. Let $\widehat{\mathcal{K}}$ be the estimated change point set determined by our algorithm. Suppose that there exists a true change point $\kappa_{0j} \notin \widehat{\mathcal{K}}$. Let $\hat{\kappa}_i$ and $\hat{\kappa}_{i+1}$ be the estimated change point which sandwich κ_{0j} , and $\hat{\kappa}_i < \kappa_{0,j-l} < \dots < \kappa_{0j} < \dots < \kappa_{0,j+r} < \hat{\kappa}_{i+1}$, where $l, r \geq 0$. Considering a new estimated change point set

$$\widetilde{\mathcal{K}} = \{\hat{\kappa}_1, \dots, \hat{\kappa}_i, \kappa_{0,j-l}, \dots, \kappa_{0,j+r}, \hat{\kappa}_{i+1}, \dots, \hat{\kappa}_{\hat{p}}\},$$

then

$$\frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} = o_p(1).$$

Proof: Let $T_0 = \hat{\kappa}_{i+1} - \hat{\kappa}_i$, $t_1 = \kappa_{0,j-l} - \hat{\kappa}_i, \dots, t_{l+r+2} = \hat{\kappa}_{i+1} - \kappa_{0,j+r}$. By the Stirling formula, we have

$$\begin{aligned} & \frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} \\ &= \frac{\prod_{j=1}^{T_0-1} (j - \sigma) \prod_{h=1}^{l+r+2} t_h!}{T_0! \prod_{h=1}^{l+r+2} \prod_{j=1}^{t_h-1} (j - \sigma)} \prod_{s=\hat{p}+1}^{\hat{p}+l+r+1} \left(\frac{s+1}{\alpha + s\sigma} \right) \frac{C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}{C(\mathbf{Y}_{(\hat{\kappa}_i, \kappa_{0,j-l}]}) \cdots C(\mathbf{Y}_{(\kappa_{0,j+r}, \hat{\kappa}_{i+1}]})} \\ &= \frac{C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}{C(\mathbf{Y}_{(\hat{\kappa}_i, \kappa_{0,j-l}]}) \cdots C(\mathbf{Y}_{(\kappa_{0,j+r}, \hat{\kappa}_{i+1}]})} O \left(\frac{\prod_{j=1}^{l+r+2} t_j}{T_0} \right)^{1+\sigma}. \end{aligned}$$

By Lemma 5, let $\underline{\kappa} = \min_{i=1, \dots, l+r+2} t_i$, there exists $c_2 > 0$ such that

$$\frac{C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}{C(\mathbf{Y}_{(\hat{\kappa}_i, \kappa_{0,j-l}]}) \cdots C(\mathbf{Y}_{(\kappa_{0,j+r}, \hat{\kappa}_{i+1}]})} = O_p \left\{ T_0^{(l+r+1)d/2} \exp(-\underline{\kappa}c_2) \right\}.$$

Thus we obtain

$$\begin{aligned} & \frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} \\ &= O \left(\frac{\prod_{j=1}^{l+r+2} t_j}{T_0} \right)^{1+\sigma} \frac{C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}{C(\mathbf{Y}_{(\hat{\kappa}_i, \kappa_{0,j-l}]}) \cdots C(\mathbf{Y}_{(\kappa_{0,j+r}, \hat{\kappa}_{i+1}]})} \\ &= O \left(\frac{\prod_{j=1}^{l+r+2} t_j}{T_0} \right)^{1+\sigma} O_p \left\{ T_0^{(l+r+1)d/2} \exp(-\underline{\kappa}c_2) \right\} \\ &= O_p \left\{ T_0^{(l+r+1)(d/2+1+\sigma)} \exp(-\underline{\kappa}c_2) \right\}. \end{aligned} \tag{12}$$

By definition of $\mathcal{H}(m_I)$,

$$\underline{\kappa} \geq m_I \geq c \{\log(T)\}^{1+\epsilon}$$

for some $c > 0$ and $\epsilon > 0$ when T is large enough. Clearly, $T \geq T_0$. Thus as $T \rightarrow \infty$,

$$\begin{aligned}
& T_0^{(l+r+1)(d/2+1+\sigma)} \exp(-\underline{\kappa}c_2) \\
& \leq T^{(l+r+1)(d/2+1+\sigma)} \exp[-\{\log(T)\}^{1+\epsilon}cc_2] \\
& = \exp[\log(T)c_0 - \{\log(T)\}^{1+\epsilon}cc_2] \\
& = \exp(\log(T)[c_0 - \{\log(T)\}^\epsilon cc_2]) \longrightarrow 0
\end{aligned} \tag{13}$$

where $c_0 = (l+r+1)(d/2+1+\sigma)$. With (12) and (13), we achieve

$$\frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} = o_p(1).$$

□

Lemma 7. *Assume the conditions in Theorem ?? hold, and \mathcal{K}_0 is a subset of $\mathcal{H}(m_I)$. Let $\widehat{\mathcal{K}}$ be the estimated change point set determined by our algorithm. Suppose that there exists an estimated change point $\hat{\kappa}_i$, such that no true change point is within its m_I -neighbourhood, i.e., $\kappa_{0j} \notin (\hat{\kappa}_i - m_I, \hat{\kappa}_i + m_I)$ for all j . Considering a newly estimated change point set*

$$\widetilde{\mathcal{K}} = \{\hat{\kappa}_1, \dots, \hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}, \dots, \hat{\kappa}_{\hat{p}}\},$$

then

$$\frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} = o_p(1).$$

Proof: By the Stirling formula, we have

$$\begin{aligned}
& \frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} \\
& = \frac{\alpha + (\hat{p} + 1)\sigma}{\hat{p} + 2} \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})} \\
& \quad \times \frac{\prod_{j=1}^{\hat{\kappa}_{i+1}-\hat{\kappa}_i-1} (j - \sigma) / (\hat{\kappa}_{i+1} - \hat{\kappa}_i)!}{\prod_{j=1}^{\hat{\kappa}_{i+1}-\hat{\kappa}_{i-1}-1} (j - \sigma) / (\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})!} \prod_{j=1}^{\hat{\kappa}_i-\hat{\kappa}_{i-1}-1} (j - \sigma) / (\hat{\kappa}_i - \hat{\kappa}_{i-1})! \\
& = \frac{\alpha + (\hat{p} + 1)\sigma}{\hat{p} + 2} \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})} \\
& \quad \times \frac{\Gamma(\hat{\kappa}_i - \hat{\kappa}_{i-1} - \sigma)\Gamma(\hat{\kappa}_{i+1} - \hat{\kappa}_i - \sigma)(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})!}{\Gamma(1 - \sigma)\Gamma(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1} - \sigma)(\hat{\kappa}_i - \hat{\kappa}_{i-1})!(\hat{\kappa}_{i+1} - \hat{\kappa}_i)!} \\
& = O\left\{ \frac{(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{1+\sigma} \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})}. \tag{14}
\end{aligned}$$

By Lemma 6, every true change point κ_{0j} is in $\widetilde{\mathcal{K}}$. Thus, there is no true change point between $\hat{\kappa}_{i-1}$ and $\hat{\kappa}_{i+1}$. Then using Lemma 4, we obtain

$$\begin{aligned}
C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)}) & = p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\widehat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})\pi(\widehat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})\det(\widehat{\boldsymbol{\sigma}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})O_p(1), \\
C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})}) & = p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\widehat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})\pi(\widehat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})\det(\widehat{\boldsymbol{\sigma}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})O_p(1), \\
C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})}) & = p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\widehat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})\pi(\widehat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})\det(\widehat{\boldsymbol{\sigma}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})O_p(1).
\end{aligned}$$

It follows that

$$\begin{aligned}
& \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})} \\
&= \frac{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}) \det(\hat{\sigma}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]})}{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]})}) \times \frac{p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}) \det(\hat{\sigma}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})})}{\pi(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]})})} \\
& \times \frac{\pi(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]})\pi(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}))}{\det(\hat{\sigma}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} O_p(1).
\end{aligned}$$

As $\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]} \xrightarrow{P} \boldsymbol{\theta}_0$, $\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]} \xrightarrow{P} \boldsymbol{\theta}_0$ and $\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]} \xrightarrow{P} \boldsymbol{\theta}_0$ where $\boldsymbol{\theta}_0$ is the true parameter, and $\pi(\boldsymbol{\theta})$ is continuous, by the continuous mapping theorem, we know

$$\begin{aligned}
& \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})} \\
&= \frac{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}) \det(\hat{\sigma}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]})}{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} \times \frac{p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}) \det(\hat{\sigma}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})})}{\det(\hat{\sigma}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} O_p(1)
\end{aligned}$$

Further,

$$\begin{aligned}
& \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})} \\
&= \frac{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}) \det(\hat{\sigma}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]})}{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} \times \frac{p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}) \det(\hat{\sigma}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})})}{\det(\hat{\sigma}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} O_p(1) \\
&= \frac{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}) p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}))}{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} O_p \left\{ \frac{\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1}}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{d/2}. \quad (15)
\end{aligned}$$

Note that by the second result of Lemma 3,

$$\begin{aligned}
& \frac{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}) p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})})}{p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}))} \\
&= \frac{p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]})}) p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}) p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i]}(\boldsymbol{\theta}_0)^{-1}}{p_{(\hat{\kappa}_i, \hat{\kappa}_{i+1}]}(\boldsymbol{\theta}_0) p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\hat{\boldsymbol{\theta}}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]})}) p_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1}]}(\boldsymbol{\theta}_0)^{-1}} \\
&= O_p(1). \quad (16)
\end{aligned}$$

Thus, combining (14), (15) and (16), we obtain

$$\begin{aligned}
\frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} &= O \left\{ \frac{(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{1+\sigma} \frac{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_i)})C(\mathbf{Y}_{(\hat{\kappa}_i, \hat{\kappa}_{i+1})})}{C(\mathbf{Y}_{(\hat{\kappa}_{i-1}, \hat{\kappa}_{i+1})})} \\
&= O \left\{ \frac{(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{1+\sigma} O_p \left\{ \frac{\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1}}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{d/2} \\
&= O_p \left\{ \frac{\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1}}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{d/2+1+\sigma}.
\end{aligned}$$

Note that $(\hat{\kappa}_{i+1} - \hat{\kappa}_i) + (\hat{\kappa}_i - \hat{\kappa}_{i-1}) = (\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})$, thus either $(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})/(\hat{\kappa}_{i+1} - \hat{\kappa}_i)$ or $(\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1})/(\hat{\kappa}_i - \hat{\kappa}_{i-1})$ will go to a constant $c > 0$ as $T \rightarrow \infty$. Therefore,

$$\frac{\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1}}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \rightarrow 0,$$

since both $(\hat{\kappa}_i - \hat{\kappa}_{i-1})$ and $(\hat{\kappa}_{i+1} - \hat{\kappa}_i)$ go to infinity as $T \rightarrow \infty$. It follows that

$$\frac{\Pr(\widehat{\mathcal{K}}|\mathbf{Y})}{\Pr(\widetilde{\mathcal{K}}|\mathbf{Y})} = O_p \left\{ \frac{\hat{\kappa}_{i+1} - \hat{\kappa}_{i-1}}{(\hat{\kappa}_{i+1} - \hat{\kappa}_i)(\hat{\kappa}_i - \hat{\kappa}_{i-1})} \right\}^{d/2+1+\sigma} = o_p(1).$$

□

Proof of Theorem ??:

By Lemma 6, we know all the true change points will fall into $\widehat{\mathcal{K}}$ with probability one as $T \rightarrow \infty$. Lemme 7 implies that all the estimated change points out of m_I -neighbourhood of true change points can be removed in probability as $T \rightarrow \infty$. By the definition of set $\mathcal{H}(m_I)$, for any two points τ_i and τ_j with $\tau_i < \tau_j$ in $\mathcal{H}(m_I)$, $(\tau_j - \tau_i) > m_I$. We obtain $\widehat{\mathcal{K}}$ by optimizing over $\mathcal{H}(m_I)$. Thus for any true change point, there is one and only one point in $\widehat{\mathcal{K}}$ within its m_I -neighbourhood, i.e., the true change point itself. Therefore, with $T \rightarrow \infty$,

$$\hat{p} \xrightarrow{\mathcal{P}} p_0 \quad \text{and} \quad \sup_{b \in \widehat{\mathcal{K}_0}} \inf_{a \in \widehat{\mathcal{K}}} |a - b| = O_p(1).$$

□

References

- Du, C., Kao, C. L. M., and Kou, S. C. (2016), “Stepwise Signal Extraction via Marginal Likelihood,” *Journal of the American Statistical Association*, 111, 314–330. 9
- Fraser, D. A. S. and McDunnough, P. (1984), “Further Remarks on Asymptotic Normality of Likelihood and Conditional Analyses,” *The Canadian Journal of Statistics*, 12, 183–190. 13
- Jiang, F., Yin, G., and Dominici, F. (2018), “Bayesian model selection approach to boundary detection with non-local priors,” in *Advances in Neural Information Processing Systems*. 11
- Rosner, B. (1983), “Percentage points for a generalized ESD many-outlier procedure,” *Technometrics*, 25, 165–172. 4
- Walker, A. M. (1969), “On the Asymptotic Behaviour of Posterior Distributions,” *Journal of the Royal Statistical Society: Series B*, 31, 80–88. 13