

Inference of breakpoints in high-dimensional time series

Likai Chen

Department of Mathematics and Statistics, Washington University in St. Louis,

Weining Wang

Department of Economics and Related Studies, University of York

Wei Biao Wu

Department of Statistics, University of Chicago

October 21, 2020

Abstract

For multiple change-points detection of high-dimensional time series, we provide asymptotic theory concerning the consistency and the asymptotic distribution of the breakpoint statistics and estimated break sizes. The theory backs up a simple two-step procedure for detecting and estimating multiple change-points. The proposed two-step procedure involves the maximum of a MOSUM (moving sum) type statistics in the first step and a CUSUM (cumulative sum) refinement step on an aggregated time series in the second step. Thus, for a fixed time-point, we can capture both the biggest break across different coordinates and aggregating simultaneous breaks over multiple coordinates. Extending the existing high-dimensional Gaussian approximation theorem to dependent data with jumps, the theory allows us to characterize the size and power of our multiple change-point test asymptotically. Moreover, we can make inferences on the breakpoints estimates when the break sizes are small. Our theoretical setup incorporates both weak temporal and strong or weak cross-sectional dependence and is suitable for heavy-tailed innovations. A robust long-run covariance matrix estimation is proposed, which can be of independent interest. An application on detecting structural changes of the U.S. unemployment rate is considered to illustrate the usefulness of our method.

Keywords: multiple change points detection; temporal and cross-sectional dependence; Gaussian approximation; inference of break locations

1 Introduction

Statistical inference of structural breaks in mean is an important subject to study, and involves estimating the trend functions, detecting and locating abnormal changes and making inferences on the break estimators. Breaks may arise in various applications in different fields, such as in network analysis, biology, engineering, economics and finance, among others. Specific examples are anomaly of network traffic data caused by attacks (Lévy-Leduc and Roueff (2009)), recurrent DNA copy number variants in multiple samples (Zhang et al. (2010)), abrupt changes in household electrical power consumption (Harlé et al. (2016)) and minimum wage policy changes analysis (Chen et al. (2020)), etc. In those data scenarios, temporal and cross-sectional dependence for large-dimensional data might pose challenges to statistical analysis.

To formulate our problem, we assume that observation vectors $Y_1; Y_2; \dots; Y_n$ follow the model,

$$Y_t = (\mu_t) + \epsilon_t; \quad t = 1; 2; \dots; n; \quad (1)$$

where $(\epsilon_t)_t$ is a sequence of zero-mean p -dimensional stationary noise vectors and $(\mu) = (\mu_1; \mu_2; \dots; \mu_p)^\top : [0; 1] \rightarrow \mathbb{R}^p$ is a vector of unknown trend functions. In this way, the data generating process is trend stationary. We will model breaks occurring on the vector of trend functions (μ_t) . Notably, we assume that the trend function satisfies

$$(\mu) = f(\mu) + \sum_{i=1}^{K_0} \mathbf{1}_{\mu} \mu_i; \quad (2)$$

where K_0 is an unknown integer representing the number of breaks; $f(\cdot) = (f_1(\cdot); f_2(\cdot); \dots; f_p(\cdot))^\top : [0; 1] \rightarrow \mathbb{R}^p$ is a vector of smooth trend functions; μ_k s with $0 < \mu_1 < \mu_2 < \dots < \mu_{K_0} < 1$ are the time stamps of the change-points with $j\mu_i = \mu_j + b$, where b is the

bandwidth parameter; and $u_k \in \mathbb{R}^p$ are the jump vectors with size $\|u_k\|_1$ ($\|\cdot\|_1$ is the infinity norm) at point u_k . Note that the jump sizes are characterized in terms of the infinity norm; therefore, we do not require simultaneous jumps for all entities $1 \leq j \leq p$, and some coordinates of u_k can be zero. Namely, we will focus on the largest jump (i.e., $\|u_k\|_1$) happening in the cross-sectional dimension for any fixed time point k (cf. Theorem 2), and this is of particular interest when the jumps are sparse. In case many series jump at the same time, we further propose a refined method, which aggregates all the contemporaneous jumps (cf. Theorem 4). In most of the change-point settings, the smooth part of the trend functions is zero (i.e., $f = 0$). This means that the trend functions are piecewise constant for each coordinate. In contrast, our model is more flexible and realistic, since we allow the mean functions to vary smoothly instead of staying at the same level between break-points.

The goal of this paper is to provide theory for structural break inference. We first detect the existence of breaks. We then deliver theorems to test for the existence of breaks, identify their change-point u_k , calibrate sizes of the breaks, i.e. $\|u_k\|_1; 1 \leq k \leq K_0$, and construct confidence intervals for the estimated break points. Our theorem allows us to consider a multiple change-point test based on a threshold method on the maximum of generalized MOSUM statistics. We derive the asymptotic distribution of the test statistics including estimated breaks sizes, and the estimated breakpoint locations (cf. Theorem 3, 4 ii)). The results provide solid foundations for conducting statistical inferences for multiple change-point estimation in high dimensional time series. Moreover, we consider a further aggregation step targeting at simultaneous breaks, and also this step gives us finer consistency rates of the break location estimation.

Multiple change-point detection can be classified into two categories, i.e. model selection and testing. The traditional model selection method, for example BIC, has the drawback of computational inefficiency, which can be improved by some modified penaliza-

tion procedure, see for example, Killick et al. (2012), and LASSO (Least absolute shrinkage and selection operator) type penalization such as by Tibshirani and Wang (2007), Li et al. (2016) and Lee et al. (2016). Regarding multiple change-point detection via testing, a classical method utilises an exhaustive search, which examines all the possible breakpoints combination. An exhaustive search is very time consuming and some dynamic technique and improved versions are invented, see for instance Bai and Perron (1998, 2003) and Jackson et al. (2005). A very popular approach is the binary segmentation introduced in Scott and Knott (1974). However its power might suffer for certain alternatives. This drawback can be handled by the wild binary segmentation algorithm developed in Fryzlewicz (2014) and sparsified binary segmentation as in Cho and Fryzlewicz (2015). Moreover, Fryzlewicz (2018) recently introduces a bottom-up algorithm to overcome the disadvantage of the classical binary segmentation. Besides, Wu and Zhou (2019) propose multiscale abrupt change estimation under complex temporal dynamics.

Detection using the MOSUM (moving sum) statistics is another popular way for multiple change-point analysis; see, for example, Hušková and Slabý (2001) for i.i.d data; Wu and Zhao (2007) and Eichinger and Kirch (2018) for general temporal dependent data. Preuss et al. (2015) deal with multivariate time series for structural breaks in covariance. A MOSUM procedure has the advantage of computation simplicity and can avoid issues due to multiple testing in multiple break inference. A possible drawback is that MOSUM introduces a new bandwidth parameter. Such an issue can be dealt with through a multi-scale MOSUM, which uses multiple bandwidths; see, for instance Meier et al. (2019). Eichinger and Kirch (2018) provide a comprehensive theoretical analysis of multiple change-point detection using MOSUM analysis including the distribution theory of the estimated breakpoint. Our work can be viewed as a generalization of their work on the high-dimensional case as we adopt a MOSUM type of statistics in our first step.

Change-point detection for high-dimensional time series has recently drawn a lot of attention due to the increasing number of applications. In particular, we shall consider the case of $p \rightarrow \infty$: Even in the simplest setup of a mean-shift model, large p may pose challenge to change-point detection. It is common to consider aggregation, either over the original time series or certain transformed statistics of individual time series and to convert the problem to a one-dimensional analysis. For instance, targeting at sparse breaks, Cho and Fryzlewicz (2015) propose a sparse binary segmentation which concerns an l_1 -based aggregation with a hard threshold, and Wang and Samworth (2018) consider sparse singular value decomposition based on the CUSUM (cumulative sum) statistics. Moreover, there are a few other work looking at l_2 -based aggregation of statistic: Bai (2010) evaluates the performance of a least square estimation of a single breakpoint with distribution theory on the break location estimates without assuming cross sectional dependency; Zhang et al. (2010) extend the method in Olshen et al. (2004); Enikeeva and Harchaoui (2019) and Liu et al. (2019) regard the detection of change-points in a high-dimensional mean vector as a minimax testing problem. For a single break point in time and targeting at sparse break coordinates, Jirak (2015) studies a CUSUM type statistic for each coordinate and then takes maximum of them, and asymptotic theory is provided to facilitate the simultaneous inferences of the breakpoint estimation. Cho (2016) proposes a double-CUSUM algorithm, etc. For a single change-point in time, distribution theory is still available in a few works, see for example Bai (2010). However, Bai (2010) is only concerning cross-sectional independent data. When it comes to multiple change points detection, the majority of the aforementioned work focus on developing novel algorithms, and a complete distribution theory is not readily available due to the complexity of the problem. An exception is Jirak (2015). Compared to Jirak (2015), we are taking a different path in terms of an algorithm using the MOSUM and an aggregation step with refined rates of estimator achieved.

We thus provide a new angle to conduct inferences in multiple change-point analysis for high-dimensional time series.

It shall be noted that as there are already many novel algorithms developed, we do not claim a total superiority of ours. The algorithm proposed here is a generalization or modification of the existing methods, which facilitates us to obtain a complete theory and good theoretical rates. Nevertheless, our aggregation step is different and complement to existing algorithms. For example, one main difference with the aggregation step is that our project is based on the estimates in the first step. Cho and Fryzlewicz (2015) and Wang and Samworth (2018) use other approaches to find the best projection direction.

To summarize, we provide theory for a two-step multiple change-point procedure. We prove consistency results as well as distribution theorems for breakpoint location estimation, which is crucial for inference of breakpoints. The aggregation step can help us to achieve good rates of the breakpoint estimation. We deliver general theoretical results that allow heavy-tailed distribution and general spatial-temporal dependency assumption on the error term, and we do not require the mean function to be piece-wise constant (i.e. $f \neq 0$). The detection procedure is not computationally expensive, as we only need to evaluate the statistic once for each point t . Additionally, we consider the estimation of the long-run covariance matrices. This paper is structured as follows. Section 2 constructs a test and delivers its asymptotic performance for testing the existence of change-points. Section 3 introduces the two-step algorithm for inference on break estimation. The associated consistency and asymptotic distribution theorems are also covered in this section. Long-run covariance matrix estimation is derived in Section 4. Simulation results are in Section A in supplementary materials and an application on U.S. unemployment rate is given in Section 5. Detailed proofs are presented in Section B in the supplementary materials.

Notations: For a constant $k \geq \mathbb{N}$ and a vector $v = (v_1; \dots; v_d)^> \in \mathbb{R}^d$; we denote

$\|A\|_k = (\sum_{i=1}^d |v_{ij}|^k)^{1/k}$, $\|A\|_1 = \sum_j |v_{ij}|$ and $\|A\|_\infty = \max_i \sum_j |v_{ij}|$. For a matrix $A = (a_{ij})_{1 \leq i \leq m, 1 \leq j \leq n}$, we define the spectral norm $\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$ and the max norm $\|A\|_{\max} = \max_{i,j} |a_{ij}|$. For a function f ; we denote $\|f\|_1 = \int |f(x)| dx$. We set (a_n) and (b_n) to be positive number sequences. We write $a_n = O(b_n)$ or $a_n \leq C b_n$ (resp. $a_n \sim b_n$) if there exists a positive constant C such that $a_n \leq C b_n$ (resp. $a_n/b_n \rightarrow 1$) for all large n , and we denote $a_n = o(b_n)$ (resp. $a_n \ll b_n$), if $a_n/b_n \rightarrow 0$ (resp. $a_n/b_n \rightarrow 1$). For two sequences of random variables (X_n) and (Y_n) ; we write $X_n = o_p(Y_n)$; if $X_n/Y_n \rightarrow 0$ in probability.

2 Testing the existence of change-points

In this section, we provide a test for the existence of breaks. Considering our observations generated by the model in (1) and (2), we would like to test the null hypothesis,

$$H_0: \mu_1 = \mu_2 = \dots = \mu_{K_0} = 0;$$

which corresponds to the case of no breaks, against the alternative of the existence of at least one break i.e. $H_A: \exists k \geq 1; \mu_k \neq 0$. It shall be noted that we do not need to assume the number of breaks (K_0) to be bounded, but to rather restrict on the separation between breakpoints (c.f. Assumption 2.4).

In Subsection 2.1, we derive our test statistic. Its asymptotic property is given in Subsection 2.2. In Subsection 2.3, we derive the performance of the test based on Gaussian approximation, which provides the theoretical foundation for calculating the size and power of the test.

2.1 Test statistic

In this subsection, we introduce the test statistics and some intuition. Recall that our trend function $f(u)$ can be disentangled into two parts, namely a smooth transition part $\tilde{f}(u)$ and a jump part $\mathbf{1}_{u = u_j}$. We can define the jump vector at point u as a gap between the right-side function $f^{(r)}(u)$ and the left-side function $f^{(l)}(u)$, which is

$$J(u) = f^{(r)}(u) - f^{(l)}(u); \quad \text{where we define } f^{(r)}(u) = \lim_{t \neq u} f(t) \quad \text{and} \quad f^{(l)}(u) = \lim_{t \rightarrow u^-} f(t):$$

Due to the smoothness of the constitutes of $f(\cdot)$; the gap function $J(u) = 0$ when there is no jump, and $J(u) = \delta_k$ when $u = u_k$. A natural way to test the existence of change-points is to check whether the gap is zero (i.e. $J(u) = 0$). To this end, we need $\hat{f}^{(r)}(u)$ and $\hat{f}^{(l)}(u)$; which are estimates of $f^{(r)}(u)$ and $f^{(l)}(u)$: We propose to adopt the local linear estimation technique, see Fan and Gijbels (1996).

The local linear estimates of $\hat{f}^{(l)}(u)$ and $\hat{f}^{(r)}(u)$ at the point $u = i/n$ are of the following weighted form

$$\hat{f}_i^{(l)} := \hat{f}^{(l)}(i/n) = \sum_{t=i}^{i-1} w_{i-t} Y_t \quad \text{and} \quad \hat{f}_i^{(r)} := \hat{f}^{(r)}(i/n) = \sum_{t=i+1}^{i+bn} w_{t-i} Y_t; \quad (3)$$

with weights

$$w_i = w_{i,b} = w_b(0; i/n); \quad i = 1; \quad w_0 = 0; \quad (4)$$

The weight functions are defined as

$$w_b(u; v) = \frac{K((v-u)=b)[S_2(u) - (u-v)S_1(u)]}{S_2(u)S_0(u) - S_1^2(u)}; \quad S_l(u) = \sum_{i=1}^n (u - i/n)^l K((i/n - u)=b); \quad (5)$$

where $K(\cdot)$ is a kernel function and b is a bandwidth with $b \rightarrow 0$ and $bn \rightarrow 1$. It is worth noting that the estimator in (3) is equivalent to adopting a one-sided kernel function, i.e. $K(u)\mathbf{1}_{u \geq 0}$ to fix the boundary estimation issue for the kernel estimation method.

If there is no jump around the time point $u = i/n$, the gap estimate $\hat{J}(i/n) = \hat{\rho}_i^{(l)} - \hat{\rho}_i^{(r)}$ would be small for all coordinates. Otherwise if for some entity $1 \leq j \leq p$, the gap estimate $j\hat{J}_j(i/n)j$ is large, there might exist a jump around time i/n at coordinate j . Note that the test statistics is in fact of a MOSUM type, and we replace the uniform kernel for MOSUM by a local linear one to adapt for slowly varying trends $f(u)$ in (2).

To conduct the breakpoint detection with $p \rightarrow 1$, we consider the maximum of the gap statistics. Furthermore, we need to standardize our test statistics in order to get a regular limit distribution. To obtain the long-run variance matrix involved in the standardization, we need to specify the error process, as in model (1). We would like to make a general temporal and cross-sectional dependence assumption. This is a crucial issue, since for time series data, dependence is the rule rather than the exception. Specifically, we let

$$\varepsilon_t = \sum_{k=0}^{\rho} A_k \varepsilon_{t-k}; \quad (6)$$

where $\varepsilon_t \in \mathbb{R}^p$ are independent and identically distributed (i.i.d.) random vectors with zero mean and an identity covariance matrix. $A_k; k \geq 0$; are coefficient matrices in $\mathbb{R}^{p \times p}$ such that ε_t is a proper random vector, and $\rho = \tilde{\rho} + c_p \rho$; for some constant $c_p > 1$. If $A_i = 0$ for all $i \geq 1$, then the noise sequences are temporally independent; if $\rho = \tilde{\rho}$ and matrices A_i are diagonal, then the sequences become the model in Bai (2010), which is spatially independent. The VMA(1) process is very general and includes many important time series models such as a vector autoregressive moving averages (VARMA) model, i.e.

$$(1 - \sum_{j=1}^s \Theta_j B^j) X_i = X_i - \sum_{j=1}^s \Theta_j X_{i-j} = \sum_{k=1}^t \Xi_k \varepsilon_{i-k};$$

where Θ_j and Ξ_k are real matrices such that $\det(1 - \sum_{j=1}^s \Theta_j z^j)$ is not zero for all $|z| = 1$ and B is the backshift operator.

Correspondingly, we define the sum of the coefficient matrix to be $S = \sum_{k=0}^{\infty} A_k$. The long run covariance matrix for the error process is

$$\Sigma = SS^>: \tag{7}$$

We denote $\Sigma = (\sigma_{ij}); 1 \leq i, j \leq p$ and

$$\Lambda = \text{diag}(\sigma_{1,1}^{-1}, \sigma_{2,2}^{-1}, \dots, \sigma_{p,p}^{-1}): \tag{8}$$

Following the previous intuition of the effect of jumps on the gap statistics $\hat{J}(\cdot)$, we consider the test statistic

$$T_n = \max_{1 \leq i \leq n} |V_i|^{-1} \hat{J}_i^{(l)} - \hat{J}_i^{(r)}; \quad \text{where } V_i = \Lambda^{-1}(\hat{J}_i^{(l)} - \hat{J}_i^{(r)}): \tag{9}$$

We adopt a supreme type of statistics as it shares good property under certain alternatives, see for example Bai and Saranadasa (1996). However, we do not claim the strict superiority of our test statistics. When the majority of locations exhibit simultaneous jumps, an l_2 type statistics tends to have better power.

We exclude Y_i in V_i , because that the weights in front of Y_i would be the same for the right side and the left side estimator, and will be canceled when taking the difference. Note that we consider the normalized statistic as multiplying the jump estimates $\hat{J}(i=n) = \hat{J}_i^{(l)} - \hat{J}_i^{(r)}$ by Λ^{-1} since the long-run variances $\sigma_{j,j}$ for different coordinates $1 \leq j \leq p$ can be very different. We refer to T_n as an infeasible test statistic since Λ is unknown. The estimation of Λ is deferred to Section 4.

2.2 Properties of the test statistics

We shall show the asymptotic properties of our test statistics T_n in (9) in this subsection. First we analyze the mean of the normalized jump estimators, i.e. $E V_j$. Intuitively, we can decompose the level of our jump estimator $E V_j$ into two parts, one is the commonly encountered bias term for the nonparametric kernel estimators of the smooth trend functions, and the other is induced by jumps on the deterministic trend, which is denoted as d_j . Recall the definition of w_j in (4) for $i = 1; 2; \dots; bn$, and $w_j = 0$ for $i = 0$ and $i > bn$: We denote the location of breaks as $u_k = nu_k$ and Ω_j as a set of indices indicating the break locations within the bn neighborhood around time j , namely $\Omega_j = \{k | |j - u_k| \leq bn\}$. If $|u_i - u_j| = n(u_i - u_j) > n$; for any i, j ; then for large n ; the cardinality of Ω_j is at most one, i.e. $|\Omega_j| \leq 1$: Actually this condition can be relaxed to $\min_{i \in \Omega_j} |u_i - u_j| > bn$: For a time point j where $|\Omega_j| > 1$, we define the weighted break sizes to be,

$$d_j = \left(1 - \sum_{t=1}^{j-1} w_t\right) \Lambda^{-1} u_k; \quad k = \operatorname{argmin}_{j-2 \leq i \leq j} |j - u_i| \quad (10)$$

and for the rest of locations j ; let $d_j = 0$: We further stack d_j over all breakpoints that are of interest, which is denoted by $\underline{d} = (d_{bn+1}^>; d_{bn+2}^>; \dots; d_n^>)^>$. It should be noted that under the null, $\underline{d} = 0$.

We denote the smooth part of the local linear estimate as

$$\hat{f}_i^{(l)} = \sum_{t=i-bn}^{i-1} w_{t-i} f(t/n) \quad \text{and} \quad \hat{f}_i^{(r)} = \sum_{t=i+1}^{i+bn} w_{t-i} f(t/n):$$

By Fan and Gijbels (1996), under some smoothness conditions, the bias part of the estimated smooth functions would be of the order b^2 , which goes to zero by assumption, i.e.

$$\max_{bn+1 \leq i \leq n-bn} |j \Lambda^{-1} (\hat{f}_i^{(l)} - \hat{f}_i^{(r)})| = O(b^2): \quad (11)$$

Given the definition of our model $Y_i = (i=T) + i; d_i$ can be expressed as

$$d_i = E\{\Lambda^{-1}((\hat{f}_i^{(r)} - \hat{f}_i^{(l)}) - (\hat{f}_i^{(r)} - \hat{f}_i^{(l)}))\} = E\{V_i - \Lambda^{-1}(\hat{f}_i^{(r)} - \hat{f}_i^{(l)})\}: \quad (12)$$

Combining (11) and (12), $E V_i$ can be approximated by the part induced by jumps κ s, as

$$jE V_i - d_i j_{\gamma} = j\Lambda^{-1}(\hat{f}_i^{(r)} - \hat{f}_i^{(l)})j_{\gamma} = O(b^2): \quad (13)$$

Let us now consider the $V_i - E V_i$ part. We observe that the centered statistics can be expressed as a weighted sum of the error term, namely

$$V_i - E V_i = \sum_{l=i}^{i-1} w_{l,i} \Lambda^{-1} \epsilon_l - \sum_{l=i+1}^{i+bn} w_{l,i} \Lambda^{-1} \epsilon_l: \quad (14)$$

To approximate its distribution, we introduce a scaling matrix for variance of the limit distribution. Recall $S = \sum_{k=0} A_k$ and define a block matrix $G = (G_{i;l})_{bn+1 \times i \times n \times bn; 1 \times n \times 2} \in \mathbb{R}^{(n-2bn)\rho \times n\rho}$ with components as $\rho \times \rho$ dimension matrices,

$$G_{i;l} = \begin{cases} w_{l,i} \Lambda^{-1} S; & \text{if } i - bn \leq l \leq i - 1; \\ w_{l,i} \Lambda^{-1} S; & \text{if } i + 1 \leq l \leq i + bn; \end{cases} \quad (15)$$

and elsewhere zero. Let \underline{Z} be a Gaussian vector in $\mathbb{R}^{n\rho}$ with zero mean and identity covariance matrix. We set $G_{i;\cdot}$ to be $(G_{i,1}; G_{i,2}; \dots; G_{i,n})$: It can be shown that $G_{i;\cdot} \underline{Z}$ has a similar covariance structure as $V_i - E V_i$. We shall use the distribution of $jG_{i;\cdot} \underline{Z} j_{\gamma}$ to approximate the distribution of $jV_i - E V_i j_{\gamma}$: Combining this approximation with the bias term in (13), we shall expect that for each time point i , our normalized break test statistics can be approximated by the maximum of a Gaussian vector centered at d_i , i.e.,

$$P(jV_i j_{\gamma} \leq u) \approx P(jd_i + G_{i;\cdot} \underline{Z} j_{\gamma} \leq u):$$

We now let the statistics go over all the time points, and recall $T_n = \max_{bn+1 \leq i \leq n} jV_{ij_1}$. Then we shall expect

$$P(T_n \leq u) = P(j\underline{d} + G \underline{z}_{j_1} \leq u); \quad (16)$$

and equivalently

$$P(T_n \leq u) = P(j\underline{d} + Z_{j_1} \leq u); \quad (17)$$

where $Z = (Z_{bn+1}^>, Z_{bn+2}^>, \dots, Z_n^>)^>$ and $(Z_i)_{bn+1 \leq i \leq n}$ is a sequence of centered Gaussian vectors in \mathbb{R}^p with covariance matrices $\text{Cov}(Z_i; Z_j) = Q_{ij}$ of the following form:

$$Q_{ij} = \mathcal{S}_{ij} \Lambda^{-1} \Sigma \Lambda^{-1} \quad \text{and} \quad \mathcal{S}_{ij} = \sum_{l=1}^n w_{ji} w_{lj} \text{sign}(i-l) \text{sign}(j-l); \quad (18)$$

To see the equivalence between (16) and (17), let

$$Q = (Q_{ij})_{bn+1 \leq i, j \leq n} = G G^>:$$

Then Z is a Gaussian vector with zero mean and covariance matrix Q : Note that

$$Z_i \stackrel{d}{=} G_{i, \underline{z}} \quad \text{and} \quad Z \stackrel{d}{=} G \underline{z}; \quad (19)$$

This transformation from $G \underline{z}$ to Z is to show that the involved Gaussian process only depends on the long-run covariance matrix Σ and the weight functions. We note that Z are not element-wise independent, but with dependency governed by G . The above argument will be rigorously formulated in Theorem 1 in the next subsection.

2.3 Gaussian approximation

In this subsection, we provide the formal theory supporting our test. We first present necessary assumptions. The following is to guarantee the smoothness of the trend functions $j(u)$ when no break occurs.

ASSUMPTION 2.1. Function $f_j \in C^2[0;1]$ with $\max_{1 \leq j \leq p} \int_0^1 |f_j''(x)| dx \leq c_f$; $\max_{1 \leq j \leq p} \int_0^1 |f_j'(x)| dx \leq c_f$ for some constant $c_f > 0$:

Additionally, to ensure the property of our kernel estimation, we need conditions on the kernel function.

ASSUMPTION 2.2. The kernel $K(\cdot)$ is symmetric with support $[-1;1]$, assume $\int_{-1}^1 K(x) dx < 1$ and $\int_{-1}^1 x K(x) dx = 0$. Also assume $K(x)$ has first-order derivative with $\int_{-1}^1 |K'(x)| dx < 1$ on $(-1;1)$. Let $b \neq 0$ and $bn \rightarrow 1$: Denote $\mu_i = \int_0^1 x^i K(x) dx$. Assume $\mu_2 \neq 0$.

We also set conditions on the regularity of the long-run covariance matrix and the dependency strength of the noise sequence.

ASSUMPTION 2.3. (Lower bound for the long run variance) $\lambda_{jj} \geq c$; $1 \leq j \leq p$ for some finite constant $c > 0$:

We need enough separation between adjacent breakpoints.

ASSUMPTION 2.4. (Separation) Assume $\min_{1 \leq i < j \leq K_0} |j_i - j_j| \geq bn$:

It is worth noting that Assumption 2.4 implies that the number of breaks K_0 shall not exceed the order $1/b$.

ASSUMPTION 2.5. (Dependence strength) $\max_{1 \leq j \leq p} \sum_k |A_{k,j}| \leq c_s(i-1)$; where $c_s > 0$ is some constant and $A_{k,j}$ is the j th row of A_k :

Assumption 2.5 is a very general spatial and temporal dependence condition and embraces many interesting processes. It requires an algebraic decay rate of the temporal dependence. However, the cross-sectional dependence does not need to be weak; and in fact, it can be strong such that it has a factor structure. We provide an example as follows.

EXAMPLE 1. Assume that $\varepsilon_{tj} \in \mathbb{R}^p$ are i.i.d random vectors with zero mean and covariance matrix I_p . Let

$$\varepsilon_t = F_t + Z_t; \text{ with } Z_t = \sum_{k=0}^{\infty} \Lambda_k \varepsilon_{t-k} \text{ and } F_t = \sum_{k=0}^{\infty} v f_k^> \varepsilon_{t-k}; \quad (20)$$

where $\Lambda_k = \text{diag}(f_{k,1}, \dots, f_{k,p})$, $v = (v_1, \dots, v_p)^>$ and $f_k = (f_{k,1}, \dots, f_{k,p})^>$. Here F_t is the factor term and Z_{tj} are independent for different j . Then the long-run variances for Z_{tj} and F_{tj} are $\Sigma_{Z,j} = (\sum_{k=0}^{\infty} f_{k,j})^2$ and $\Sigma_{F,j} = v_j \sum_{k=0}^{\infty} f_{k,j}^2 v_j^2$, respectively. If for some constant $c > 0$:

$$\sum_{k=0}^{\infty} f_{k,j}^2 = \Sigma_{Z,j}^{-1} < cI \quad \text{and} \quad \sum_{k=0}^{\infty} v_j f_{k,j}^2 v_j = \Sigma_{F,j}^{-1} < cI; \quad (21)$$

then Assumption 2.5 holds with $\Sigma = \Sigma_{Z,j}^{-1}$. To see this, we note $j A_{k,j} j_2 = (f_{k,j}^2 + v_j f_{k,j}^2 v_j^2)^{1/2}$; and $j_2 j = \Sigma_{Z,j}^{-1} + \Sigma_{F,j}^{-1}$. Hence,

$$\sum_{k=0}^{\infty} j A_{k,j} j_2 \leq \sum_{k=0}^{\infty} (f_{k,j}^2 + v_j f_{k,j}^2 v_j^2)^{1/2} < cI \quad \left(\Sigma_{Z,j}^{-1} + \Sigma_{F,j}^{-1} \right)^{1/2} < \sqrt{2} cI \quad \Sigma_{Z,j}^{-1}.$$

ASSUMPTION 2.6. (Finite moment) The innovations ε_{ij} are i.i.d. with $\mathbb{E}|\varepsilon_{ij}|^q < \infty$ for some $q > 4$:

ASSUMPTION 2.7. (Sub-exponential) The innovations ε_{ij} are i.i.d. with $\mathbb{E}e^{a_0|\varepsilon_{ij}|} < \infty$; for some $a_0 > 0$:

Assumptions 2.6 and 2.7 put tail assumptions on the distribution of the noise sequences. Given the above-mentioned conditions, we provide the main Gaussian approximation theorem, which is essential for the asymptotic distribution of our test statistics T_n . Our theorem extends the Gaussian approximation theory in Chernozhukov et al. (2013a, 2017), which

build on the Stein's method and the anti-concentration bounds. Markedly, our theory is developed for modeling dependent data. To this aim, one important technical non-triviality lies in handling the spatial-temporal dependency of the trend stationary high-dimensional processes. We derive the corresponding concentration inequalities based on m -dependence approximation of the underlying processes. Compared to the existing results on Gaussian approximation for time series, for example Zhang et al. (2017), our setting works for non-centered Gaussian approximation that accommodates our interest for time series with breaks.

Theorem 1 (Gaussian approximation). *Under Assumptions 2.1-2.5,*

(i) *if Assumption 2.6 holds, and $np(bn)^{-q-2}(\log(np))^{3q-2} = o(1)$; then for*

$$\Delta = (bn)^{-1-6} \log^{7-6}(pn) + (bn)^{-(3+3)} \log^{2-3}(np) + b^{5-2} n^{1-2} \log^{1-2}(np);$$

we have

$$\sup_{u \in \mathbb{R}} |P(T_n \leq u) - P(j\mathcal{d} + Zj_1 \leq u)| \leq \Delta + ((np)^{2-q} (bn))^{1-3} \log(pn);$$

(ii) *if Assumption 2.7 holds, and $(bn)^{-1}(\log(np))^{\max\{7, 2(1+q)\}} = g = o(1)$; we have*

$$\sup_{u \in \mathbb{R}} |P(T_n \leq u) - P(j\mathcal{d} + Zj_1 \leq u)| \leq \Delta;$$

where the constants in Δ are independent of $n; p; b$.

If in addition

$$b^5 n \log(np) = o(1); \tag{22}$$

then under both cases, we have

$$\sup_{u \in \mathbb{R}} |P(T_n \leq u) - P(j\mathcal{d} + Zj_1 \leq u)| \rightarrow 0; \tag{23}$$

REMARK 1. (Allowed dimension) One key theoretical insight is that we explicitly show the trade-off between the tail assumption of the innovations and the allowed dimension of the time series ρ relative to the sample size n in the above theorem. In particular, when we have exponential tail assumption on the distribution of the innovations, we allow an ultra high dimension setup indicating ρ to be at an exponential rate with respect to n . And when we have only finite moment assumptions, we can allow ρ to be at a polynomial order with respect to n . Specifically, for Theorem 1 case (i), we allow ρ to be of some polynomial order of n , and its order depends on the value of q : For some $\alpha_1 > 0$ and $0 < \alpha_2 < 1=2$; assume $\rho = n^{\alpha_1}$ and $b = n^{-\alpha_2}$: If $\alpha_1 < (1 - \alpha_2)q=2 - 1$ and $\alpha_2 > 1=5$; then (23) holds. It is easy to see that the bigger the moment q is, the larger the allowance of the dimension ρ : The moment condition 2.6 depends on q which characterizes the heavy tailedness of the noise, larger q means thinner tails. For case (ii), we can allow ρ to be exponential in n , i.e. the ultra high dimensional scenario. For instance, for some $\alpha_1 > 0$ and $1=5 < \alpha_2 < 1$; we can set $\rho = e^{n^{\alpha_1}}$ and $b = n^{-\alpha_2}$: If $\alpha_1 < 5 - \alpha_2 - 1$ and $\alpha_1 \max\{7; 2(1 + \alpha_2)\} = g < 1 - \alpha_2$; then (23) holds.

It is not hard to understand the size and power implication of Theorem 1 to our test. Under the null hypothesis, we have $\underline{d} = 0$, then for any prefixed significant level $\alpha \in (0; 1)$; we have the critical value of our test as q i.e. the quantile of the Gaussian limit distribution,

$$q = \inf_{r > 0} \{r : P(|Z| > r) = \alpha\} \quad (24)$$

As from the Gaussian approximation result in (23), we have the approximated sizes of the test statistics,

$$\left| P(T_n > q) - P(|Z| > q) \right| \rightarrow 0:$$

We shall reject the null hypothesis at the significant level α ; if the test statistics exceed the critical value i.e. $T_n > q$:

To evaluate our testing power, consider the alternative that if not all $\kappa = 0$; then \underline{d} is non-zero. We have the following corollary for the power, which is a straightforward consequence of Theorem 1.

Corollary 1. *(Power) Under conditions in Theorem 1 (i) or (ii). The testing power satisfies*

$$= P(\underline{d} + Zj_{\gamma} \geq q) + o(1):$$

Thus, we can see that the power of our test would depend on the vector \underline{d} . The size of it is determined by the true jump sizes i.e. κ s. Since the covariance matrix for Z is $Q = (Q_{i,j})$; where $Q_{i,j} = \mathcal{S}_{i,j} \Lambda^{-1} \Sigma \Lambda^{-1}$ with $\mathcal{S}_{i,j}$ defined in (18). It can be calculated that $\mathcal{S}_{i,i} = (bn)^{-1}$, therefore $jZj_{\gamma} = O_p((bn)^{-1/2} \log(np))$; which tends to zero by Assumption 2.6 and 2.7. Thus if $\underline{d}j_{\gamma} = (bn)^{-1/2} \log(np)$; $\rightarrow 1$ by Corollary 1.

3 Estimation and inference of breaks

In this section, we show how to estimate the number of change-points, the time stamps, the spatial coordinates and the sizes of the structural breaks. We summarize the key steps of the adopted two step procedure for the multiple change-point detection. The main reason for a two-step estimation is to achieve an optimal rate of consistency for our break estimation. The first step can be regarded as an extension of the MOSUM I_{γ} aggregation. Namely, in our first step, we conduct a “rough” estimation through a MOSUM type statistic as in Equation (9), and we can draw a conclusion on the existence of a break. In case it exists, we obtain a “rough” estimate of the change-points locations. In the second step, we refine our jump estimates based on a one-dimensional aggregated time series. The aggregation can be viewed as a projection using information on the jump estimators from the first step.

To be more specific, within each time region around the k th breakpoint, we can aggregate data by a weighted sum of different coordinates whose weights are determined by the first step jump size estimators (\hat{v}_k) . Instead of looking at the biggest break at one time point, the aggregated change-point statistics carry more information regarding significant jumps across contemporaneous locations, and would thus provide better precision. In the following, we introduce the first “rough” estimation step and its properties in Subsection 3.1. We further improve the first step in Subsection 3.1 through an aggregated statistics, which is proposed and analyzed in Subsection 3.2.

3.1 The “rough” estimation step

We define the sizes of the breakpoints at time k as

$$j\Lambda^{-1} \kappa_j^1 :$$

Here, we normalize κ_j by the long-run standard deviations for the same reason as V_i in (9). Intuitively, the noise fluctuation levels for different locations can be very different, and at one location, a break can be significant due to purely high noise level without normalization. We define the minimum size of breaks over time as

$$= \min_{1 \leq k \leq K_0} j\Lambda^{-1} \kappa_j^1 : \tag{25}$$

In the following, we outline the steps of our testing, detecting and estimation procedure.

Step 1. For significance level α , we test the existence of jumps based on the critical value q in (24). If we find no significant breaks, then we cannot reject the null H_0 . In case our test statistic exceeds the critical value, we reject H_0 and acknowledge the existence of breaks, then we proceed to step 2.

Step 2. To detect the change-points, we collect all the time stamps with the jump statistics $jV j_1$ exceeding a threshold value w , namely, $A_1 = \{bn + 1, \dots, n - bn : jV j_1 > w\}$; where V is defined in (9). Let $\hat{\tau}_1$ be the time point in A_1 that maximizes the test statistics $jV j_1$. We further eliminate a $2bn$ neighborhood of time points around $\hat{\tau}_1$ from A_1 to create A_2 . Then we find the next point in A_2 that maximize $jV j_1$, and repeat the same operation until the set A_k is empty. Namely, for $k \geq 1$, we let the k th estimated break point be denoted as $\hat{\tau}_k = \operatorname{argmax}_{t \in A_k} jV j_1$ and $A_{k+1} = A_k \setminus \{j : |j - \hat{\tau}_k| \leq 2bn\}$. We denote the maximum number of breakpoints as \hat{K}_0 , with $\hat{K}_0 = \max_{k \geq 1} \{k : A_k \neq \emptyset\}$. It is worth noting that we have chosen $2bn$ to exclude both bn neighborhood of $\hat{\tau}_k$ and $\hat{\tau}_{k+1}$.

Step 3. Given the detected breakpoints in Step 2, we calculate the break sizes over time. We denote the window size to be $M = bn$;

$$\hat{\tau}_k = \hat{\tau}_k^{(l)} \wedge_M \hat{\tau}_k^{(r)} \quad \text{and} \quad \hat{\tau}_k = \min_{1 \leq k \leq \hat{K}_0} j\Lambda^{-1} \hat{\tau}_k j_1 : \quad (26)$$

It is worth noting that in this algorithm, we only need to calculate the gap statistics $jV j_1$ once for each point. Hence, it is not time consuming regardless of the true number of breakpoints. In Step 1, we test the existence of the breaks. In Step 2, we use the estimated $jV j_1$ for all the points from $bn + 1$ to $n - bn$ and select the points that are beyond the threshold w . Intuitively, the points in A_1 would contain the break indices, as well as points in their neighborhood where estimates are contaminated by the breaks. Therefore in Step 2, we find the local maximums and discard points around them. In Step 3, we estimate the sizes of the change-points and calculate their minimum values.

In the following, we shall provide consistency results of estimates of the break numbers, locations and break sizes in Theorem 2; and derive asymptotic distribution of break sizes in Theorem 3.

We need to first impose the minimum jump size condition on the break size as

ASSUMPTION 3.1. Assume the break size satisfies $\max \{ \sqrt{\log(pn)/(bn)}; b \}$:

It can be seen that the break size requirement is related to the dimensionality of the time series, the number of observations available and the bandwidth parameter. The larger the sample n , the smaller the requirement for δ due to the better approximation of the trends. In the following theorem, we show that we would asymptotically obtain the right number of breaks. Moreover, we can bound the errors of the estimated break locations and the break sizes. The threshold δ^y shall be set as a high quantile of its limited Gaussian distribution to ensure the consistent estimation of the breaks.

Theorem 2. We assume conditions in Theorem 1 (i) or (ii), Assumption 3.1, and (22) hold. If $\min f \delta^y; \delta^y g \leq 2c_w^0 \sqrt{\log(pn)/(bn)}$; where c_w^0 is the constant defined as the limit $(bn \sum_{i=0}^n w_i^2)^{1/2} \rightarrow c_w^0$, then

(i) $P(\hat{K}_0 = K_0) \rightarrow 1$:

(ii) under Theorem 1 (i), $j_{\hat{k}} - k_j = O_p(f(np)^{2-q} = \delta^2 g)$; and under Theorem 1 (ii), $j_{\hat{k}} - k_j = O_p(f \log^2(np) = \delta^2 g)$; uniformly over k ; where $k = \operatorname{argmin}_j j_{\hat{k}} - j$.

(iii) $j_{\Lambda^{-1}(\hat{k} - k)} j_1 = O_p((bn)^{-1/2} \log(np)^{1/2} + b)$; uniformly over k ; which indicates $j_{\hat{k}} - j = O_p((bn)^{-1/2} \log(np)^{1/2} + b)$:

Result (i) indicates that the number of breaks can be consistently estimated, (ii) suggests that the estimated break dates u_k can be consistently determined in view of $u_k = k/n$ and (iii) shows that the break sizes can be consistently recovered. The convergence rate of the break sizes depends on the bandwidth b , sample size n and the dimension of the time series p . It is worth noting that the bias is of order b in (iii), as the difference is taken with a gap of $2M$ as in equation (26). It shall be noted that the consistency rate of \hat{k} depends on the break size δ , which depends only on the maximum break size for any fixed time.

Therefore having several large breaks simultaneously would not improve the break size estimation. With respect to the condition $\min_{j \in \mathcal{J}} \frac{1}{\rho} \frac{1}{\log n} \frac{1}{nb} \geq 2c_w^d \sqrt{\log(\rho n) = (bn)}$; relative to the summary in the Table 1 in Cho (2016), our break size is comparable up to the weakest condition $(nb)^{1-2} \geq 1$ up to a logarithmic factor. Moreover, when $\rho = 1$, our requirement of breaksize is similar to the rate as in Theorem 3.2 in Wu and Zhou (2019), namely $\frac{\rho}{\log n} = \frac{\rho}{nb}$.

Given the consistency of the breakpoints, we can obtain a distribution theory that facilitates us in making inferences on the break sizes. Let \tilde{Z} be a Gaussian vector in \mathbb{R}^p with zero mean and covariance matrix

$$\tilde{Q} := Q_{bn+1;bn+1} = 2 \sum_{t=1}^{bn} w_t^2 \Lambda^{-1} \Sigma \Lambda^{-1}; \quad (27)$$

Theorem 3. (*Break size inference*) Assume conditions in Theorem 2 and $b^3 n \log(np) = o(1)$. We have

$$\sup_{u \in \mathbb{R}} \mathbb{P}(\Lambda^{-1}(\hat{\alpha}_k - \alpha_k) \leq u) = \mathbb{P}(j \tilde{Z} \leq u) \geq 0; \text{ where } k = \operatorname{argmin}_j \hat{\alpha}_k - \alpha_j;$$

This theorem indicates that the maximum of the difference between the estimated jump size $\hat{\alpha}_k$ and the true jump size α_k can be approximated by the maximum of a Gaussian random vector with the same asymptotic variance-covariance structure. Based on Theorem 2 (ii) and Theorem 3, we can construct joint confidence interval for $\alpha_{k:j}$. We set

$$q = \mathbb{P}(j \tilde{Z} \leq q) \text{ and } q = (q_{1,1}^{1=2}, q_{2,2}^{1=2}, \dots, q_{p,p}^{1=2})^>; \quad (28)$$

some $\delta > 0$: Then as $n \rightarrow \infty$ with probability close to 1; we have

$$q_{j:j}^{1=2} + \hat{\alpha}_{k:j} - \alpha_{k:j} \leq q_{j:j}^{1=2} + \hat{\alpha}_{k:j} - \delta_j; \quad (29)$$

Theorem 3 can be extended to hold uniformly over k by stacking the statistics over all k s. In addition, we see that Theorem 2 and Theorem 3 are closely connected in the sense that we can reach the same threshold by stacking $\hat{\alpha}_k$ over all k s.

3.2 The refined aggregation step

The estimation in the first step is only driven by $\max_j |j_{k,j}|$, i.e. the maximum size of jumps at a time point t_k . Therefore it is only sensitive to the biggest jump across all the time series at the same time. The l_1 type test potentially have more power for certain alternatives than the l_2 type statistics, see for example Bai and Saranadasa (1996). However, if the majority or all of the entities exhibit simultaneous jumps, the supremum statistic tends to have lower power than the l_2 statistic.

In case there are multiple simultaneous time series jumps, it would be beneficial to modify our procedure to aggregate all of the series with a jump. This enlightens us to propose a two-stage method: first, we follow the steps in the previous subsections to detect the “rough” timing of the jumps and the estimated jump sizes; second, for each bn neighborhood of a change-point estimate \hat{t}_k obtained from step one, we update the change-point estimates according to a newly aggregated time series. The time series is calculated with a weighted sum of simultaneous observations corresponding to significant jump locations and the weights are based on the jump size estimates in the first step. The aggregation returns a one-dimensional time series with richer information on the cross-sectional jumps.

We denote S_k to be the set of series that jump at location t_k ; that is

$$S_k = \{j \mid \rho_j |j_{k,j}| \neq 0\}; \quad (30)$$

where $j_{k,j}$ is the j th coordinate of \hat{t}_k : Detailed steps of the aggregation are formulated as follows:

Stage 1. Apply Steps 1-3 in Subsection 3.1 to obtain \hat{t}_k and $\hat{\rho}_k$; $k = 1; 2; \dots; \hat{K}_0$: For some $w^\nu > 0$; let the estimation of S_k be

$$\hat{S}_k = \{j \mid \rho_j |(\Lambda^{-1} \hat{t}_k)_j| \geq w^\nu\}; \quad (31)$$

In practice, w' can be chosen to be large enough to ensure that we can detect all the jumps with probability 1 as in Theorem 2.

Stage 2. For $j \in \hat{S}_k$, $j \geq 2bn$; we let

$$X_t = \sum_{j \in \hat{S}_k} (\Lambda^{-1} \hat{\Lambda}_k)_j (\Lambda^{-1} Y_t)_j \quad (32)$$

Note that after the modification, for all the jump locations, the new time series X_t would only contain positive sized jumps i.e. $\sum_{j \in \hat{S}_k} (\Lambda^{-1} \hat{\Lambda}_k)_j^2$. This step can be understood as a projection of the high-dimensional observations $\Lambda^{-1} Y_t$ according to the direction of $\Lambda^{-1} \hat{\Lambda}_k$ ($j \in \hat{S}_k$). This is similar to the idea of Wang and Samworth (2018).

Based on the aggregated time series X_t , the refined change-point locations can be detected through a CUSUM type of test statistics, for $k = 1; 2; \dots; \hat{K}_0$:

$$\tilde{\tau}_k = \operatorname{argmax}_{j \in \hat{S}_k} \left(\sum_{s=\hat{\tau}_k}^{\hat{\tau}_k+2bn} X_s \frac{t - \hat{\tau}_k + 2bn}{4bn + 1} - \sum_{s=\hat{\tau}_k}^{t-1} X_s \right) \sqrt{\frac{4bn + 1}{(t - (\hat{\tau}_k - 2bn) + 1)(\hat{\tau}_k + 2bn - t)}} \quad (33)$$

After we update the break points estimation, we can construct confidence intervals for the updated breakpoints estimates $\tilde{\tau}_k$. We denote the long-run correlation matrix to be $(\tilde{\Sigma}_{ij})_{i,j} = \Lambda^{-1} \Sigma \Lambda^{-1}$; where Σ is the long run covariance matrix for t . We let $\tilde{\Sigma}_k = (\tilde{\Sigma}_{ij})_{i,j \in \hat{S}_k}$ be the sub covariance matrix corresponding to coordinates in \hat{S}_k at time $\hat{\tau}_k$ and let the standardized significant break sizes $\tilde{a}_k = (\Lambda^{-1} \hat{\Lambda}_k)_{i \in \hat{S}_k}$: We define two objects involved in the limit distributions of the breaks, i.e.,

$$a_k = \tilde{a}_k^2 \quad \text{and} \quad \mathfrak{L}_k^2 = \tilde{\Sigma}_k^{-1} \tilde{a}_k \quad (34)$$

Then \mathfrak{L}_k^2 is the long-run variance for the sequence $\sum_{j \in \hat{S}_k} (\Lambda^{-1} \hat{\Lambda}_k)_j (\Lambda^{-1} Y_t)_j$. For the aggregated jump estimation, we alternatively define the minimum jump size across different

locations and time points as

$$j_k = \min_{1 \leq j \leq K_0} \min_{2S_k} j(\Lambda^{-1} \Sigma \Lambda^{-1})_k j;$$

Then j_k and its functions similarly as \tilde{a}_k to capture the identifiable jump size of the time series. We shall put the same assumption on j_k as on \tilde{a}_k . It is worth noting that j_k is the minimum jump size to ensure the consistency of our break estimation.

ASSUMPTION 3.2. Let $j_k \geq \max \{ \sqrt{\log(pn)/(bn)}, b \}$:

In the following corollary, we show that we can consistently recover the locations of the series with a jump for each change-point. It can be directly derived from Theorem 2 (iii).

Corollary 2. *We assume conditions in Theorem 1 (i) or (ii) hold, and Assumption 3.2. If $j_k \geq \sqrt{\log(pn)/(bn)} + b$; then we have*

$$P(\hat{S}_k = S_k; 1 \leq k \leq K_0) \rightarrow 1;$$

In addition, we provide a theorem that allows us to make inference on the estimated break-dates \tilde{a}_k .

Theorem 4. (Aggregated break estimation) *Assume conditions in Corollary 2, and that for some constants $c_1, c_2 > 0$;*

$$c_1 \leq \min_k (\Lambda^{-1} \Sigma \Lambda^{-1})_k \leq c_2; \quad (35)$$

Recall the definition of \tilde{a}_k and \hat{a}_k in (34). Then we have for any fixed $1 \leq k \leq K_0$;

(i) $\hat{a}_k - \tilde{a}_k = O_p(\hat{a}_k^2 - \tilde{a}_k^2)$;

(ii) In addition, if Assumption 2.5 holds with $\beta > 1$, and $\hat{a}_k^2 - \tilde{a}_k^2 \leq bn$; then we have

$$\hat{a}_k - \tilde{a}_k \stackrel{P}{\rightarrow} (\hat{a}_k - \tilde{a}_k)^2 \arg \max_r (\hat{a}_k^2 - \tilde{a}_k^2 + W(r));$$

where $W(r)$ is a two-sided Brownian motion. That is $W(r) = W_1(r)$; if $r > 0$, and $W(r) = W_2(-r)$; if $r < 0$. W_1, W_2 are two independent Brownian motions.

REMARK 2. We shall note that the consistency rate of $\tilde{\kappa}$ is improved compared to the results for $\hat{\kappa}$ in Theorem 2 ii). a_k which is an l_2 aggregation of simultaneous significant break sizes, plays a role in the rate of convergence of $\tilde{\kappa}$. For instance, if we assume that there are s breaks which are of size $\delta > 0$ in the cross-sectional dimension, then $a_k = s^2$. If moreover there is no cross-sectional correlation, i.e., $\tilde{\Sigma}_k = I$, then we may expect $\tilde{\kappa}$ to be consistent so long that $1/(s^2) \neq 0$, while $\hat{\kappa}$ can be not consistent. Thus the rate of $\tilde{\kappa}$ will be better than $\hat{\kappa}$. Moreover, the long-run variance also plays a critical role in the rate of convergence. For example, when the variance part of the limit distribution satisfies $\delta_k^2 = \sum_{k,j} \tilde{\Sigma}_{kj} a_k$, if $\sum_{k,j} \tilde{\Sigma}_{kj} a_k = o(1)$ then by Theorem 4 (i), we have $\tilde{\kappa} \neq \kappa$ in probability. This corresponds to the insight of Bai (2010) and Hansen (2000). We can also see that when the breaks are truly sparse in the cross sectional dimension or the break size for each time series is very small, the l_2 aggregation cannot improve the performance compared to the previous step. Also when there are strong cross-sectional dependence l_2 aggregation will not improve the break estimation performance. Moreover, we also need the aggregated breaksize to shrink to zero ($1/\delta_k^2 = a_k^2$) to obtain the limit distribution. If $\tilde{\Sigma}_k$ is a d -banded matrix, $\sum_{k,j} \tilde{\Sigma}_{kj} = (\sum_{k,j} \tilde{\Sigma}_{kj})^{1/2} = d$: We can derive that $\sum_{k,j} \tilde{\Sigma}_{kj} = O_p(d a_k)$:

To illustrate the insight of Remark 2, we compare the performance of a simple model with $N(0;0,1)$ and one breakpoint placed at $\tau_0 = 50$. Figure 3.2 shows the histogram of $\hat{\kappa}_{\tau_0}$ and $\tilde{\kappa}_{\tau_0}$ respectively. The jump-size for all breaks are the same, with the value $1.6\sqrt{\log(np)}(\tau_0)^{-1/3}$. As the dimension p grows, we see the significant improvement of the performance of $\tilde{\kappa}$ relative to that of $\hat{\kappa}$.

From Theorem 4, with estimates of a_k and δ_k , we can construct a $100(1 - \alpha)\%$ confidence interval for $\tilde{\kappa}$:

$$(\tilde{\kappa} - b\hat{q}^{\alpha/2} \leq C \leq \tilde{\kappa} + b\hat{q}^{\alpha/2} + 1); \quad (36)$$

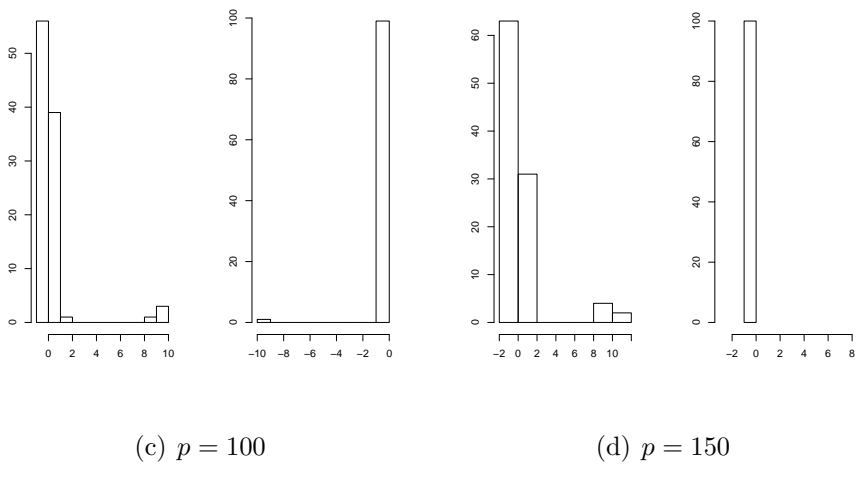
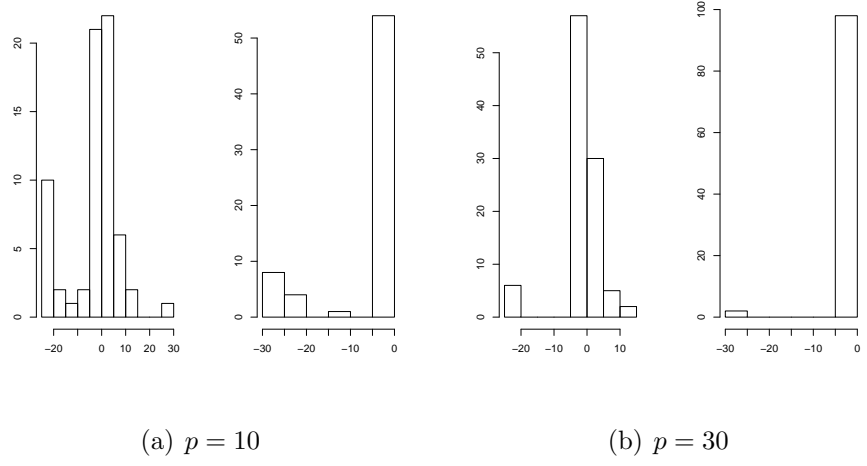


Figure 1: Histogram of $\hat{\alpha}_0$ (left) and $\tilde{\alpha}_0$ (right) for $n = 100, p = 10; 30; 100; 150, \mathcal{K}_0 = 1$. The number of breaks in the cross-sectional dimension are $s = 1; 5; 20; 30$ respectively, and there are 100 simulation samples. (a) describes the case with $p = 10; s = 1$; (b) describes the case with $p = 30; s = 5$; (c) describes the case with $p = 100; s = 20$; and, (d) describes the case with $p = 150; s = 30$.

where $q_{1-\alpha}^0$ (q_{α}^0) is $1-\alpha$ (α)th quantile of the limit distribution of the break point $\tilde{\tau}_k$, i.e. $\arg\max_r f^{-1}(\alpha) + \lambda_k W(r)g$ and $\hat{q}_{1-\alpha}^0$ (\hat{q}_{α}^0) are estimates of the quantiles. $\lfloor \cdot \rfloor$ denotes the floor function. $q_{1-\alpha}^0$ (q_{α}^0) can be calculated following Stryhn (1996). Alternatively, we can also simulate the critical values.

4 Long-run covariance matrix

In the previous sections, we assume that Σ is known. However, this is unrealistic in practice, as we mostly do not know the long-run covariance matrix. Thus, an estimation of the long-run covariance matrix is needed in Gaussian approximation. A simpler version of this problem was considered by Politis et al. (1999) and Lahiri (2003), who allow for a constant mean of the random vector. More generally, Chen and Wu (2019) consider the high-dimensional situation with smooth trends. However, this does not fit directly to our interest due to the possible existence of the breakpoints. We then propose a robust covariance matrix estimation motivating from the M-estimation method in Catoni (2012). It is worth noting that due to the jumps, our method shall be different from the classical covariance matrix estimation. Our long-run variance-covariance matrix estimation is complementary to the recent article on high-dimensional robust matrix method under independence settings in Fan et al. (2017).

First of all, to account for temporal dependency, we group our observations into blocks of the same size m , for some $m \geq \mathbb{N}$. We denote the number of blocks $N_1 = \lfloor n/m \rfloor = mc$; and the observation indices within a block k is set to be $A_k = \{t \in \mathbb{N} : km+1 \leq t \leq (k+1)m\}$, and we let

$$\bar{y}_k = \frac{1}{|A_k|} \sum_{t \in A_k} Y_t$$

be the average observations within the block A_k : Without jumps, a natural estimate of the

long-run covariance matrix is

$$\sum_{k=1}^{N_1} (m=2) (k \quad k-1) (k \quad k-1)^{\top} = N_1;$$

Note that we take the difference $y_k - y_{k-1}$ to cancel out the trends, as the trend function $f(t)$ is smooth, and the aggregated difference between two consecutive blocks can be shown to be of order m/n , which vanishes when $m/n \rightarrow 0$. However, this estimator can be greatly contaminated by the jumps, as jumps are not smooth and cannot be canceled out by taking difference. Thus a robust covariance matrix estimation is needed. We borrow the framework of Catoni (2012), who considers a new robust M -estimation method. We extend the method for estimating our long run covariance matrix.

We denote $y_k = (y_{k,1}, y_{k,2}, \dots, y_{k,p})^{\top}$ and let

$$\hat{h}_{ij;k} = m (y_{k,i} - y_{k-1,i}) (y_{k,j} - y_{k-1,j}) = 2; \quad k = 1, 2, \dots, N_1; \quad (37)$$

For some $\rho_{ij} > 0$; we denote the M -estimation zero function of our variance-covariance matrix to be

$$h_{ij}(u) = \sum_{k=1}^{N_1} \rho_{ij} (\hat{h}_{ij;k} - u) = N_1; \quad (38)$$

where $\rho(x) = \rho(x^2)$ and

$$\rho(x) = \begin{cases} \log(2); & |x| \leq 1; \\ \log(1 - x + x^2); & 0 < |x| < 1; \\ \log(1 + x + x^2); & 1 < |x| < 0; \\ \log(2); & |x| \geq 1; \end{cases} \quad (39)$$

REMARK 3. Function $j(\cdot)$ is bounded by $\log(2)$ and is Lipschitz continuous with the Lipschitz constant bounded by 1. Also note that the function has envelopes of nice form,

$$\log(1 - x + x^2) \leq j(x) \leq \log(1 + x + x^2). \quad (40)$$

We set the estimates of the components of the long-run covariance matrix $\hat{\Sigma}_{i,j}$ to be the solution to $h_{i,j}(u) = 0$ (if more than one root, pick one of them). We can collect all the estimates of the variance and covariances and organize them into the variance covariance matrix,

$$\hat{\Sigma} = (\hat{\Sigma}_{i,j})_{1 \leq i,j \leq p} \quad \text{and} \quad \hat{\Lambda} = \text{diag}(\hat{\Sigma}_{1,1}^{1=2}, \hat{\Sigma}_{2,2}^{1=2}, \dots, \hat{\Sigma}_{p,p}^{1=2}). \quad (41)$$

We denote $\hat{\Sigma}_{i,i} = 2 \sum_{N_1=4}^{k=3N_1=4} \hat{\Sigma}_{i,i;k=N_1}$ and let the $\hat{\Sigma}_{i,j}$ in (38) be $\hat{\Sigma}_{i,i}^{1=2} \hat{\Sigma}_{j,j}^{1=2} (m=n)^{1=2}$.

Theorem 5. (*Long-run variance precision*) *We assume that Assumption 2.5 holds with 1.5 and let*

$$\& = j \Lambda^{-1} (\hat{\Sigma} - \Sigma) \Lambda^{-1} j_{\max}.$$

Then for K_0 finite, we have $\& = O_P(n^{-1/4} \log(np))$ under either one of the following two conditions:

- (i) *Assuming conditions in Theorem 1 (i), $p = cn^v$ with $v < q=8$ and some $c > 0$; we take $m = \min\{n^{1-8v/(q-4)}, n^{1=2}\}g$;*
- (ii) *Assuming conditions in Theorem 1 (ii), we take $m = n^{1=2}$;*

By the above theorem, for the diagonal values, we have $\max_{1 \leq i \leq p} \hat{\Sigma}_{i,i} - \Sigma_{i,i} = o_P(1)$. Let \hat{Q} be the same as Q in (18), with Σ therein replaced by $\hat{\Sigma}$ in (41). We denote \hat{Z} as

the Gaussian vector with covariance matrix \hat{Q} ; then by Theorem 5 and Lemma 3, $j\hat{Z} + dj_1$ converges to $jZ + dj_1$ in distribution. Thus, all previous results are still valid with $\hat{\Sigma}$ as well.

5 Application

As an application, we analyze the monthly the unemployment rate data in 20 U.S. states (namely, Alabama, Arizona, California, Colorado, Florida, Georgia, Illinois, Indiana, Kentucky, Michigan, Mississippi, New Jersey, New York, North Carolina, Ohio, Pennsylvania, Texas, Virginia, Washington and Wisconsin). The data time span is from January 1976 to September 2018 ($n = 513$), and the data source is Bureau of Labor Statistics from Department of Labor in the United States (<https://www.bls.gov/>). Figure 2 displays the 20 time series of unemployment rate. Although from a long time span and on an overall level, we do not see obvious abrupt structural changes, it would be still of great interest to consider detected changes induced by some well-known exogenous shocks, such as the subprime mortgage crisis in 2007-2008. It is understood that there will be likely a smooth cyclical trend associated with the unemployment time series, as they mostly rise during a recession and fall during periods of economics prosperity, following the business cycle. Further studies on whether the shock induced by recessions creates a significant structural change in the unemployment rate should be performed. We select b according to a cross validation method, $m = 10$, and $\hat{\rho}_{ij} = \frac{-1=2}{i;i} \frac{-1=2}{j;j} (m=bn)^{1=2}$ which varies over different i, j . We have used the estimated 0.999 quantile of the maximum of the Gaussian random variables (as in equation (19) with correlation matrix replaced by its estimator), which in our case is estimated as 2:10. We refer to the guidance of the selection of tuning parameters as in Remark 4 in the Supplementary materials.

Figure 3 shows the estimated robust long-run correlation matrix using the method in Section 4. One sees some significant values in the correlations between residuals in different states. We can see that the correlations across different locations are not negligible, however our method is robust against the underlying spatial-temporal dependency.

Figure 4 plots the estimated breakpoints and the confidence intervals around them. We see that the estimated breaks \tilde{k} using the CUSUM statistics in Equation (33) pick up the breaks earlier than the estimates obtained from the non-aggregated method i.e. \hat{k} . We can see that our method can identify important dates such as the financial crisis period starting in Jan, 2009. Moreover, \tilde{k} tends to detect earlier dates of structure changes than the observed averaged peaks in the time series. Other time-points with significant jumps detected are January 1977, October 1981, January 1991 and October 2001. There are a few recession periods documented by the national bureau of Economics Research, namely November 1973 to March 1975, July 1981 to November 1982, July 1990 to March 1991, July 1981 to November 1982, July 1990 to March 1991 and March 2001 to November 2001. All the break-dates of the unemployment structure happen during or slightly before the recession periods, featuring a close relationship between the structure change of unemployment rate and the economic cycles. This implies that economic recessions indeed bring significant structural changes in unemployment rates across all the states.

Figure 2: Plot of Unemployment rate of 20 U.S. states

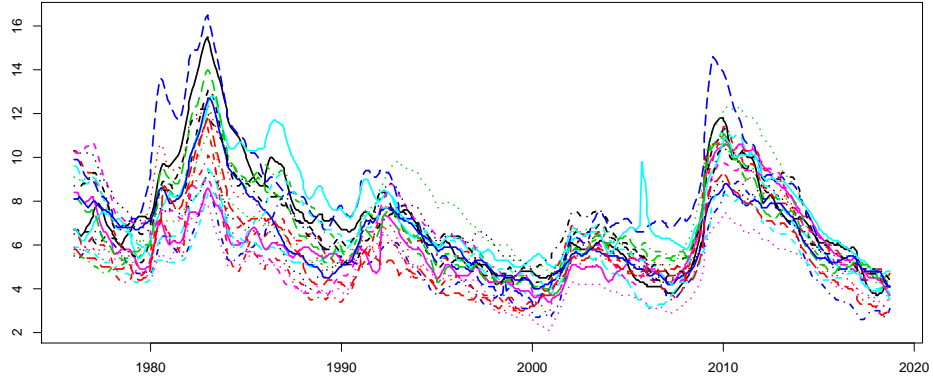


Figure 3: Plot of estimation of the robust long-run correlation matrix; $m = 10$.

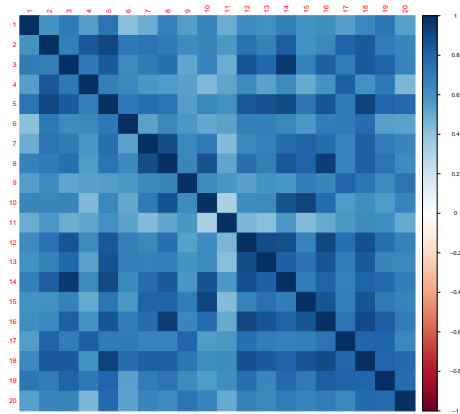
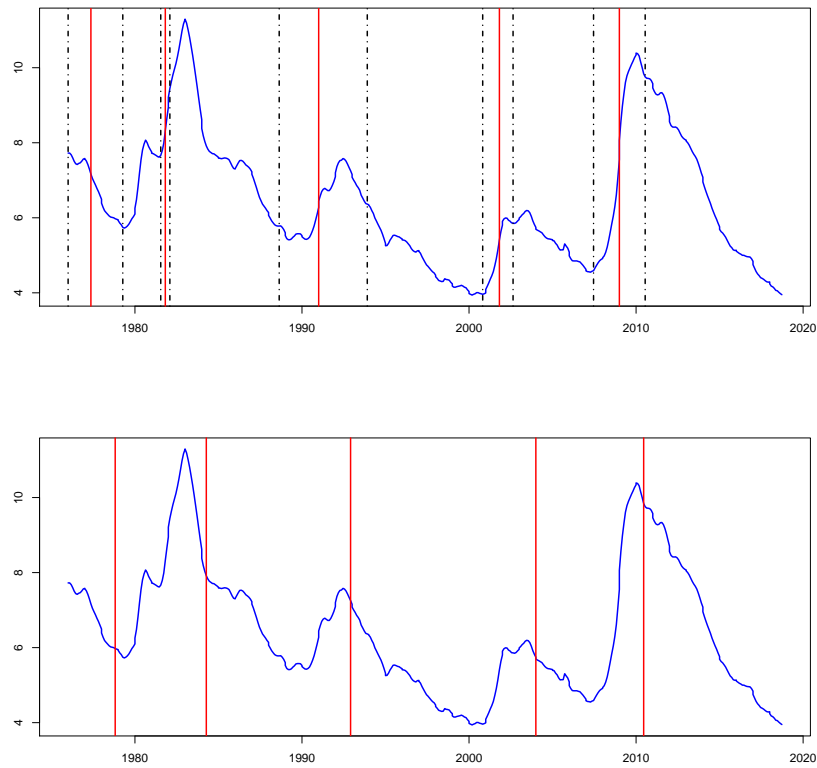


Figure 4: Plot of estimated breakpoints $\tilde{\kappa}(\hat{\kappa})$ (red lines) and their confidence intervals (dotted black lines). $\tilde{\kappa}$ (upper panel), $\hat{\kappa}$ (lower panel). The blue time series line represents the average unemployment rate over states under consideration.



References

- Bai, J. (2010). Common breaks in means and variances for panel data. *J. Econometrics* 157(1), 78–92.
- Bai, J. and P. Perron (1998). Estimating and testing linear models with multiple structural changes. *Econometrica* 66(1), 47–78.
- Bai, J. and P. Perron (2003). Computation and analysis of multiple structural change models. *Journal of applied econometrics* 18(1), 1–22.
- Bai, Z. and H. Saranadasa (1996). Effect of high dimension: by an example of a two sample problem. *Statistica Sinica*, 311–329.
- Berkes, I., W. Liu, and W. B. Wu (2014). Komlós-Major-Tusnády approximation under dependence. *Ann. Probab.* 42(2), 794–817.
- Burkholder, D. L. (1988). Sharp inequalities for martingales and stochastic integrals. *Asterisque* (157-158), 75–94. Colloque Paul Lévy sur les Processus Stochastiques (Palaiseau, 1987).
- Catoni, O. (2012). Challenging the empirical mean and empirical variance: a deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.* 48(4), 1148–1185.
- Chen, L., W. Wang, and W. B. Wu (2020). Dynamic semiparametric factor model with structural breaks. *Journal of Business & Economic Statistics*, 1–15.
- Chen, L. and W. B. Wu (2019). Testing for trends in high-dimensional time series. *Journal of the American Statistical Association* 114(526), 869–881.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2013a). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann. Statist.* 41(6), 2786–2819.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2013b). Testing many moment inequalities. *arXiv preprint arXiv:1312.7614*.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2015). Comparison and anti-concentration bounds for maxima of Gaussian random vectors. *Probab. Theory Related Fields* 162(1-2), 47–70.

- Chernozhukov, V., D. Chetverikov, and K. Kato (2017). Central limit theorems and bootstrap in high dimensions. *Ann. Probab.* 45(4), 2309–2352.
- Cho, H. (2016). Change-point detection in panel data via double CUSUM statistic. *Electron. J. Stat.* 10(2), 2000–2038.
- Cho, H. and P. Fryzlewicz (2015). Multiple-change-point detection for high dimensional time series via sparsified binary segmentation. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* 77(2), 475–507.
- Eichinger, B. and C. Kirch (2018). A MOSUM procedure for the estimation of multiple random change points. *Bernoulli* 24(1), 526–564.
- El Machkouri, M., D. Volný, and W. B. Wu (2013). A central limit theorem for stationary random fields. *Stochastic Process. Appl.* 123(1), 1–14.
- Enikeeva, F. and Z. Harchaoui (2019). High-dimensional change-point detection under sparse alternatives. *Ann. Statist.* 47(4), 2051–2079.
- Fan, J. and I. Gijbels (1996). *Local polynomial modelling and its applications*, Volume 66 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London.
- Fan, J., Q. Li, and Y. Wang (2017). Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* 79(1), 247–265.
- Fryzlewicz, P. (2014). Wild binary segmentation for multiple change-point detection. *Ann. Statist.* 42(6), 2243–2281.
- Fryzlewicz, P. (2018). Tail-greedy bottom-up data decompositions and fast multiple change-point detection. *Ann. Statist.* 46(6B), 3390–3421.
- Hansen, B. E. (2000). Sample splitting and threshold estimation. *Econometrica* 68(3), 575–603.
- Harlé, F., F. Chatelain, C. Gouy-Pailler, and S. Achard (2016). Bayesian model for multiple change-points detection in multivariate time series. *IEEE Trans. Signal Process.* 64(16), 4351–4362.
- Hušková, M. and A. Slabý (2001). Permutation tests for multiple changes. *Kybernetika* 37(5), 605–622.

- Imbens, G. W. and T. Lemieux (2008). Regression discontinuity designs: A guide to practice. *Journal of econometrics* 142(2), 615–635.
- Jackson, B., J. D. Scargle, D. Barnes, S. Arabhi, A. Alt, P. Gioumouisis, E. Gwin, P. Sangtrakulcharoen, L. Tan, and T. T. Tsai (2005). An algorithm for optimal partitioning of data on an interval. *IEEE Signal Processing Letters* 12(2), 105–108.
- Jirak, M. (2015). Uniform change point tests in high dimension. *Ann. Statist.* 43(6), 2451–2483.
- Killick, R., P. Fearnhead, and I. A. Eckley (2012). Optimal detection of changepoints with a linear computational cost. *J. Amer. Statist. Assoc.* 107(500), 1590–1598.
- Lahiri, S. N. (2003). *Resampling methods for dependent data*. Springer Series in Statistics. Springer-Verlag, New York.
- Lee, S., M. H. Seo, and Y. Shin (2016). The lasso for high dimensional regression with a possible change point. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 78(1), 193–210.
- Lévy-Leduc, C. and F. Roueff (2009). Detection and localization of change-points in high-dimensional network traffic data. *Ann. Appl. Stat.* 3(2), 637–662.
- Li, D., J. Qian, and L. Su (2016). Panel data models with interactive fixed effects and multiple structural breaks. *Journal of the American Statistical Association* 111(516), 1804–1819.
- Liu, H., C. Gao, and R. J. Samworth (2019). Minimax rates in sparse, high-dimensional changepoint detection. *arXiv preprint arXiv:1907.10012*.
- Meier, A., C. Kirch, and H. Cho (2019). mosum: A package for moving sums in change-point analysis. *Preprint*.
- Nazarov, F. (2003). On the maximal perimeter of a convex set in \mathbb{R}^n with respect to a Gaussian measure. In *Geometric aspects of functional analysis*, Volume 1807 of *Lecture Notes in Math.*, pp. 169–187. Springer, Berlin.
- Olshen, A. B., E. Venkatraman, R. Lucito, and M. Wigler (2004). Circular binary segmentation for the analysis of array-based dna copy number data. *Biostatistics* 5(4), 557–572.

- Politis, D. N., J. P. Romano, and M. Wolf (1999). *Subsampling*. Springer Series in Statistics. Springer-Verlag, New York.
- Preuss, P., R. Puchstein, and H. Dette (2015). Detection of multiple structural breaks in multivariate time series. *J. Amer. Statist. Assoc.* 110(510), 654–668.
- Rio, E. (2009). Moment inequalities for sums of dependent random variables under projective conditions. *J. Theoret. Probab.* 22(1), 146–163.
- Scott, A. J. and M. Knott (1974). A cluster analysis method for grouping means in the analysis of variance. *Biometrics*, 507–512.
- Stryhn, H. (1996). The location of the maximum of asymmetric two-sided brownian motion with triangular drift. *Statistics & Probability Letters* 29(3), 279 – 284.
- Tibshirani, R. and P. Wang (2007). Spatial smoothing and hot spot detection for cgh data using the fused lasso. *Biostatistics* 9(1), 18–29.
- Wang, T. and R. J. Samworth (2018). High dimensional change point estimation via sparse projection. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* 80(1), 57–83.
- Wu, W. and Z. Zhou (2019). Mace: Multiscale abrupt change estimation under complex temporal dynamics. *arXiv preprint arXiv:1909.06307*.
- Wu, W.-B. and Y. N. Wu (2016). Performance bounds for parameter estimates of high-dimensional linear models with correlated errors. *Electron. J. Stat.* 10(1), 352–379.
- Wu, W. B. and Z. Zhao (2007). Inference of trends in time series. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 69(3), 391–410.
- Zhang, D., W. B. Wu, et al. (2017). Gaussian approximation for high dimensional time series. *The Annals of Statistics* 45(5), 1895–1919.
- Zhang, N. R., D. O. Siegmund, H. Ji, and J. Z. Li (2010). Detecting simultaneous changepoints in multiple sequences. *Biometrika* 97(3), 631–645. With supplementary data available online.

SUPPLEMENTARY MATERIAL

A Simulation

In this section, we conduct a simulation study to evaluate the accuracy of our method. We define $u_{it} = \varepsilon_{it}(t=T)$. The discrete version of the model can be written as:

$$y_{it} = u_{it} + \sum_{j=1}^{K_0} \beta_{ijt} \mathbf{1}_{t=j} + \varepsilon_{it} \quad (42)$$

$i = 1; \dots; p, t = 1; \dots; n.$

We use cross validation to select the bandwidth and the block parameter. The detailed testing procedure is summarized as follows in line with the descriptions in Section 2.

Step 1 (Long-run covariance estimation.) We estimate the long-run covariance matrix $\hat{\Sigma} = (\hat{\Sigma}_{i,j})$ and its diagonal matrix $\hat{\Lambda}$: We first calculate $\hat{\Sigma}_{i,j;k}$ in (37) and we let $\hat{\Sigma}_{i,j}$ be the solution of $h_{i,j}(u) = 0$ as in (38).

Step 2 (Q matrix relates to critical values.) We construct the block matrix $\hat{Q} = (\hat{Q}_{i,j})$; where $\hat{Q}_{i,j}$ is $Q_{i,j}$ in (18), with Σ and Λ therein replaced by $\hat{\Sigma}$ and $\hat{\Lambda}$ respectively.

Step 3 (Calculating critical values.) We generate i.i.d. Gaussian vectors $\hat{Z}^{(l)}$; $i = 1; 2; \dots; N$; with the covariance matrix \hat{Q} ; and we obtain \hat{q} which is the empirical $(1 - \alpha)$ quantile of the $j \hat{Z}^{(l)} j_1$ over several samples and it can be viewed as an estimate of q in (24).

Step 4 (Testing the existence of jump.) We construct \hat{T}_n as T_n in (9) with Λ replaced by $\hat{\Lambda}$: We reject the null hypothesis that there is no jump at level α if \hat{T}_n is larger than \hat{q} :

Step 5 (**Detecting significant break-points.**) Supposing that we reject the null in Step 4, we will continue with the following steps. To detect the significant jumps, we construct $j\hat{V}_{j_1}$ for $t = bn + 1; bn + 2; \dots; n - bn$; where \hat{V}_t is the same as V_t in (9) with Λ therein replaced by $\hat{\Lambda}$: Let $A_1 = \inf_{j_1} j\hat{V}_{j_1} > w^\nu g$: w^ν can be set as \hat{q} with \hat{q} to be small (e.g. $\hat{q} = 0.0001$).

Step 6 (**Stamping multiple breaks**) In the case of multiple significant breaks in Step 5, we sequentially locate the multiple change-points following steps in Section 3.1. To be more specific, for $k = 1$, we let $\hat{k} = \arg\max_{2A_k} j\hat{V}_{j_1}$ and $A_{k+1} = A_k \wedge \inf_{j_1} j\hat{V}_{j_1} > 2bn g$: Then the estimate of the number of breaks is $\hat{K}_0 = \max_{k=1, \dots, K_0} \hat{k}$:

Step 7 (**Estimating the sizes of breaks**) We construct \hat{k} as in Step 3 in Subsection 3.1. We set the estimates of the sizes of the jumps as $\hat{\alpha}_k = j\hat{\Lambda}^{-1} \hat{k} j_1$ and their minimum as $\hat{\alpha} = \min_{k=1, \dots, \hat{K}_0} \hat{\alpha}_k$:

Step 8 (**Constructing confidence intervals for the sizes**) We construct \tilde{q} as in (28).

Let

$\hat{\alpha} = (\hat{\alpha}_{1,1}^{1=2}, \hat{\alpha}_{2,2}^{1=2}, \dots, \hat{\alpha}_{p,p}^{1=2})$: Then the confidence interval for vector $\hat{\alpha}_k$ at level $2 - \alpha$ is $(\hat{\alpha}_k - \tilde{q} \hat{\alpha}; \hat{\alpha}_k + \tilde{q} \hat{\alpha})$.

Step 9 (**Aggregated jump location estimation and confidence interval construction**) Construct aggregated jump location estimates $\tilde{\alpha}_k$ as in (33). The confidence interval for $\hat{\alpha}_k$ is $(\tilde{\alpha}_k - \chi; \tilde{\alpha}_k + \chi)$; where χ is the $1 - \alpha/2$ quantile of the distribution $\arg\max_r (2^{-1} \hat{\alpha}_k^2 j r j + \hat{\alpha}_k W(r))$ and $\hat{\alpha}_k$ (resp. $\hat{\alpha}_k$) is α_k (resp. α_k) with Λ ; Σ and $\hat{\alpha}_k$ replaced by their estimations.

REMARK 4. We notice that there are a few tuning parameters involved in the procedure. The rigorous study of the selection of tuning parameters deserves further research. We

make some suggestions as follows. For b we can use either a plug-in approach or a cross validation method following section 5.2 of Imbens and Lemieux (2008). We select the maximum bandwidth over all locations for each time. $\hat{\rho}_{ij}$ is chosen according to the guidance specified below equation (41). m can be set up as $(bn)^{1/2}$ as in Theorem 5 (ii) initially. And later can be updated by a grid search method which minimizing the out of sample prediction error. We suggest to eliminate the breakpoints and its bn neighborhood when we calculate the prediction error.

We first report a few results with a known variance-covariance matrix. We put it under rather simple settings for checking the performance of the algorithm with respect to different $\rho = 20; 50; 100, \text{ and } 150$, $n = 100$. Therefore, Step 2 of the above-mentioned algorithm is omitted. We also include the cases with strong cross-sectional dependence with factor structure and no cross sectional dependence. In particular, we consider different kinds of data generating processes. We choose a) $f_i(u) = \rho^2 + u^2$ and b) $f_i(u) = \sin(2u + i\rho)$. Let $u_{it} = f_i(t=n)$. u_{it} is taken to follow 1) an i.i.d. standard normal distribution, 2) a VAR(1) model, with a randomly simulated coefficient matrix (maximum eigenvalue smaller than 1) and 3) a factor structure together with a VAR(1) noise. In case 3), the factor loading and factors are generated with i.i.d. $N(0;1)$. We set $K_0 = 10$ breaks for all cases, and we increase the number of breaks in the cross-sectional dimension as the dimension increases, i.e. $s = 1; 5; 20; 30$. We set the break size to be $3 = \log(np)$.

For the unknown covariance, we report results for the cases $n = 500; 1000$ and $\rho = 20; 30$. The break-locations are selected to start at time-point 100 and are distanced by 100, and the break sizes are set to be either i) $\rho_{jit} = 0.05$ for $i = 5; 10$ or ii) $\rho_{jit} = (\rho / \bar{t} = \log(\rho n))$. Figure 5 shows the simulated data with the model corresponding to the case a), i), ii). We evaluate our simulation performance over 1000 samples. We report the averaged difference between the estimated number of breaks and the true break points (AD) $(\hat{K}_0 - K_0)$ as in Table

Table 1: AD averaged over 1000 samples in different simulation scenarios, and their standard deviations in bracket.

	$\rho = 20; n = 100$		$\rho = 50; n = 100$	
	1)	3)	1)	3)
a)	0.152 (0.033)	0.187 (0.046)	0.121 (0.024)	0.125 (0.029)
b)	0.155 (0.041)	0.193 (0.054)	0.119 (0.025)	0.126 (0.028)
	$\rho = 100; n = 100$		$\rho = 150; n = 100$	
a)	0.098 (0.015)	0.117 (0.021)	0.072 (0.013)	0.078 (0.014)
b)	0.093 (0.017)	0.124 (0.025)	0.084 (0.016)	0.090 (0.019)

1. The averaged distances between the breaks $\sum_{k=1}^{K_0} \tilde{j}_k$ (AM) are shown in Table
2. We notice that as the dimension and the the number of breaks in the cross-sectional dimension grow, the estimation performance improves.

The estimation accuracy with covariance estimation is included in Tables 4 and 5. And the averaged coverage probabilities of the confidence interval for the breaks (AC) are in Table 6 at the confidence level of 90%.

When the tail of error distribution is light, our long run covariance estimator can still be consistent. Table 3 shows the AD and AM/n in cases: $\rho = 100$ and $n = 80$; $\rho = 150$ and $n = 100$. We set $s = 5$, $K_0 = 10$. We find that the method is slightly less accurate with the increasing dimension. To improve the estimation accuracy, a further extension of our estimation to regularized high-dimensional long-run variance estimation can be considered. As the sample sizes increase, the estimation precision is improved. We can see that our method is robust against different data simulation scenarios, and we can achieve good level

Table 2: AM/n averaged over 1000 samples in different simulation scenarios, and their standard deviations in bracket.

	$\rho = 20; n = 100$		$\rho = 50; n = 100$	
	1)	3)	1)	3)
a)	0.046 (0.016)	0.057 (0.025)	0.033 (0.012)	0.041 (0.015)
b)	0.048 (0.019)	0.060 (0.028)	0.036 (0.014)	0.040 (0.018)
	$\rho = 100; n = 100$		$\rho = 150; n = 100$	
a)	0.015 (0.007)	0.023 (0.009)	0.011 (0.003)	0.017 (0.008)
b)	0.019 (0.008)	0.028 (0.011)	0.012 (0.003)	0.020 (0.007)

of accuracy with our method. In particular, the spatial and temporal dependency in the error term would not affect our estimation.

Figure 6 shows the plot of the estimated robust long-run covariance matrix (right) against the true one (left). On an overall level, we see that the true correlation matrix has been precisely recovered, as the patterns of these two plots look the same. We also report the distance between our robustly estimated variance-covariance matrix and the true one in Table 7. The estimation precision of the long-run variance-covariance matrix is maintained across different data-generating processes.

Table 3: AD, AM/n, and the maximum norm of covariance matrix estimation error averaged over 1000 samples in different simulation scenarios, and their standard deviations in bracket.

	$p = 100; n = 80$		$p = 150; n = 100$	
	1)	2)	1)	2)
	AD			
a)	0.185 (0.041)	0.190 (0.044)	0.197 (0.058)	0.186 (0.062)
b)	0.189 (0.042)	0.194 (0.057)	0.192 (0.061)	0.187 (0.063)
	AM/n			
a)	0.057 (0.018)	0.082 (0.019)	0.093 (0.024)	0.099 (0.018)
b)	0.064 (0.021)	0.055 (0.026)	0.067 (0.023)	0.079 (0.027)
	Maximum Norm			
a)	0.080 (0.016)	0.075 (0.019)	0.076 (0.021)	0.083 (0.024)
b)	0.078 (0.034)	0.084 (0.041)	0.081 (0.037)	0.082 (0.022)

Table 4: AD averaged over 1000 samples in different simulation scenarios, and their standard deviations in bracket.

		$p = 20; n = 500$		$p = 30; n = 1000$	
		1)	2)	1)	2)
a)	i)	0.035 (0.010)	0.046 (0.014)	0.029 (0.005)	0.030 (0.007)
	ii)	0.028 (0.011)	0.043 (0.015)	0.024 (0.004)	0.027 (0.006)
b)	i)	0.038 (0.012)	0.039 (0.012)	0.023 (0.004)	0.029 (0.003)
	ii)	0.023 (0.008)	0.032 (0.011)	0.026 (0.003)	0.028 (0.002)

Table 5: AM/n averaged over 1000 samples in different simulation scenarios, and their standard deviations are in brackets.

		$p = 20; n = 500$		$p = 30; n = 1000$	
		1)	2)	1)	2)
a)	i)	0.044 (0.018)	0.056 (0.015)	0.033 (0.011)	0.038 (0.009)
	ii)	0.027 (0.014)	0.033 (0.013)	0.021 (0.008)	0.026 (0.007)
b)	i)	0.045 (0.017)	0.057 (0.018)	0.028 (0.004)	0.035 (0.006)
	ii)	0.039 (0.012)	0.037 (0.014)	0.016 (0.003)	0.023 (0.004)

Table 6: AC in different simulation scenarios over all the estimated break-points and samples.

		$\rho = 20; n = 500$		$\rho = 30; n = 1000$	
		1)	2)	1)	2)
a)	i)	0.692	0.676	0.833	0.824
	ii)	0.719	0.708	0.856	0.833
b)	i)	0.685	0.673	0.847	0.810
	ii)	0.741	0.715	0.889	0.872

Table 7: Averaged difference between the variance-covariance and the true matrix. (L_1 norm divided by $\rho(\rho - 1)=2$).

		$\rho = 20; T = 500$		$\rho = 30; T = 1000$	
		1)	2)	1)	2)
a)	i)	0.005	0.008	0.003	0.007
	ii)	0.004	0.006	0.003	0.005
b)	i)	0.006	0.009	0.003	0.005
	ii)	0.004	0.007	0.002	0.004

B Proof

B.1 Some useful Lemmas

LEMMA 1 (Basic properties of the weights). *We assume Assumption 2.2. Then by Fan and Gijbels (1996), the weights of the local linear estimator take the following form*

$$w_i = \frac{2}{2} \frac{1}{0} \frac{i-(bn)}{2} \frac{K(i-(bn))}{bn} + O((bn)^{-2});$$

We have the following results which holds uniformly over i . There exist strictly positive constants c_w, c_w^0, c_w^{00} only depending on kernel $K(\cdot)$; such that

$$\begin{aligned} & bn \max_{0 \leq i \leq n} |w_i| \leq c_w; \max_{j \leq m} |w_j| \leq c_w \frac{m}{(bn)^2}; \\ & (bn \sum_{i=0}^{bn} w_i^2)^{1/2} \leq c_w^0; \text{ and } \frac{bn}{k} \sum_{i=1}^k w_i \leq c_w^{00} \cdot k \leq bn; \end{aligned} \quad (43)$$

Proof. We only show the last one, since the rest are similar and easier. Note

$$bn \sum_{i=1}^k w_i = k = F(k/(bn)) + O((bn)^{-1}); \quad \text{where } F(t) = \frac{2 \int_0^t K(x) dx - \int_0^t x K(x) dx}{(2-1)t}.$$

Define the numerator function as $g(t) = 2 \int_0^t K(x) dx - \int_0^t x K(x) dx$: We can see that $g(0) = 0, g(1) > 0$; and the derivative function $g'(t) = (2-1)t K(t)$, which is strictly larger than 0 before $t = 1/2$ and less than 0 afterwards. Therefore we have $F(x) > 0$ on $(0, 1]$: In addition, we note $F(0+) = 2K(0) = (2-1) > 0$ and $F(1) = 1$: Thus $\inf_{t \in (0, 1]} F(t) > 0$ in view of $F(t)$ is a continuous function. \square

LEMMA 2 (Burkholder (1988), Rio (2009)). *Let $q > 1, q^0 = \min\{q, 2\}$. Let $M_T = \sum_{t=1}^T \epsilon_t$ where $\epsilon_t \in L^q$ are martingale differences. Then*

$$k M_T k_q^{q^0} \leq K_q^{q^0} \sum_{t=1}^T k_t k_q^{q^0}; \quad \text{where } K_q = \max\{(q-1)^{-1}, \sqrt{q-1}\};$$

B.2 Asymptotic results for Gaussian vector

LEMMA 3 (Comparison). Let $X = (X_1; X_2; \dots; X_v)^>$ and $Y = (Y_1; Y_2; \dots; Y_v)^>$ be two centered Gaussian vectors in \mathbb{R}^v and let $d = (d_1; d_2; \dots; d_v)^> \in \mathbb{R}^v$. We denote $\Delta = \max_{1 \leq i, j \leq v} \left| \frac{X_{ij}}{v} - \frac{Y_{ij}}{v} \right|$; where we define $\frac{X_{ij}}{v} = E(X_i X_j)$ (resp. $\frac{Y_{ij}}{v} = E(Y_i Y_j)$). Assume that Y_i s have the same variance $\sigma^2 > 0$. Then we have

$$\sup_{x \in \mathbb{R}^v} \left| P(jX + dj_1 \leq x) - P(jY + dj_1 \leq x) \right| \leq \Delta^{1+3} \log(v)^{2+3}; \quad (44)$$

where the constant involved in (44) only depends on σ :

Proof. It suffices to show for any $d \in \mathbb{R}^v$:

$$\sup_x \left| P\left(\max_{1 \leq i \leq v} (X_i + d_i) \leq x\right) - P\left(\max_{1 \leq i \leq v} (Y_i + d_i) \leq x\right) \right| \leq \Delta^{1+3} \log(v)^{2+3};$$

To this end, we define

$$F(z) = \frac{1}{v} \log \left(\sum_{j=1}^v \exp(-z_j + d_j) \right);$$

Replace the $F(\cdot)$ in the proof of Theorem 2 in Chernozhukov et al. (2015) by $F(z)$. Then by the argument in equation (10) in Chernozhukov et al. (2015), we have

$$P\left(\max_{1 \leq i \leq v} (X_i + d_i) \leq x\right) - P\left(\max_{1 \leq i \leq v} (Y_i + d_i) \leq x + \frac{1}{v} \log(v)\right) + c(\sigma^2 + \frac{1}{v})\Delta;$$

where c is some absolute constant. Then by Lemma 4, we have

$$P\left(\max_{1 \leq i \leq v} (X_i + d_i) \leq x\right) - P\left(\max_{1 \leq i \leq v} (Y_i + d_i) \leq x\right) \leq \left(\frac{1}{v} + \frac{1}{v} \log(v)\right) \sqrt{\log(v)} + \left(\sigma^2 + \frac{1}{v}\right)\Delta;$$

where the constant in (45) only depending on σ : Take $\frac{1}{v} = \frac{1}{v} \log(v)$ and $\frac{1}{v} = \log(v)^{1+6} \Delta^{1+3}$. Same argument can be applied in the other direction, and the desired result follows. \square

LEMMA 4 (Nazarov (2003)). Let $X = (X_1; X_2; \dots; X_v)^T$ be a centered Gaussian vector in \mathbb{R}^v : Assume $E(X_i^2) = b$ for some $b > 0$ and all $1 \leq i \leq v$: Then for any $\epsilon > 0$ and $d \in \mathbb{R}^v$:

$$\sup_{x \in \mathbb{R}^v} P\left(\left| \sum_{i=1}^v X_i d_i - x \right| \geq \epsilon\right) \leq c \sqrt{\log(v)}; \quad (45)$$

where c is some constant depending only on b .

B.3 Proof of Theorem 1

The proof of Theorem 1 is quite involved. We shall first provide some intuitive ideas of the proof strategy. We define

$$I := \max_{1 \leq i \leq n} \left| \sum_{t=i}^{i-1} W_t \Lambda^{-1} \mathbf{1}_t - \sum_{t=i+1}^{i+bn} W_t \Lambda^{-1} \mathbf{1}_t + d_i \right|; \quad (46)$$

By (13) and (14) we have

$$jT_n - I \leq \max_{1 \leq i \leq n} jE V_i - d_i = O(b^2); \quad (47)$$

Thus we only need to work on I : For some $m > 0$; let a truncated version of the error term be defined as

$$I_{t,m} = \sum_{k=0}^{m-1} A_k \mathbf{1}_{t-k};$$

Consider the m -dependent approximation $I_{\cdot,m}$ of I ; where $I_{\cdot,m}$ is I with $\mathbf{1}_t$ replaced by $I_{t,m}$: Then we have $I \leq I_{\cdot,m}$ for large m : Let $I_{Z,m}$ be $I_{\cdot,m}$ with $\mathbf{1}_t$ therein replaced by Z_t ; where $(Z_t; t \geq Z)$ are i.i.d. Gaussian vectors with zero mean and identity covariance matrix in \mathbb{R}^p : Since $I_{\cdot,m}$ can be rewritten into the format of the maximum of summation of independent vectors, by the Gaussian approximation theorem in Chernozhukov et al. (2017), the distributions of $I_{\cdot,m}$ and $I_{Z,m}$ are close. We complete the proof by showing that the distributions of $I_{Z,m}$ and $jZ + \underline{d}_1$ are close, and the continuity of the maximum of a non-centered Gaussian distribution.

Proof. We now proceed with the formal argument. We shall first focus on case (i). (We note that we use the same order of m in both case i) and ii).). Let $m = (bn)^{1-(1+\epsilon)}$; for any $\epsilon > 0$;

$$P((bn)^{1-2}T_n \leq u) = P((bn)^{1-2}jT_n \leq I_{:,mj} | u) + P((bn)^{1-2}I_{:,m} \leq u+)$$

and

$$P((bn)^{1-2}jZ + \underline{dj}_1 \leq u) = P((bn)^{1-2}jZ + \underline{dj}_1 \leq u+) - P(u < (bn)^{1-2}jZ + \underline{dj}_1 \leq u+):$$

Hence

$$\begin{aligned} & \sup_{u \in \mathbb{R}} \left[P((bn)^{1-2}T_n \leq u) - P((bn)^{1-2}jZ + \underline{dj}_1 \leq u) \right] \\ & P\left((bn)^{1-2}jT_n \leq I_{:,mj} \mid u\right) + \sup_{u \in \mathbb{R}} \left| P(I_{:,m} \leq u) - P(jZ + \underline{dj}_1 \leq u) \right| \\ & + \sup_{u \in \mathbb{R}} P\left(\left| (bn)^{1-2}jZ + \underline{dj}_1 \leq u \right| \right) =: I_1 + I_2 + I_3. \end{aligned}$$

For the I_1 part, $jT_n \leq I_{:,mj} \iff jT_n \leq jI + jI' \leq I_{:,mj}$: Recall (47), then $jT_n \leq jI \iff c_0 b^2$ for some constant $c_0 > 0$: We define $\epsilon = 2c_1(np)^{1-q}m^{-1}$; where the constant c_1 is the one to be defined in Lemma 5. Then by Lemma 5, we have

$$P((bn)^{1-2}jI \leq I_{:,mj} \mid u) \leq (bn)^{-q-2}.$$

Hence for $\epsilon = \epsilon + c_0(bn)^{1-2}b^2$; $I_1 = O((bn)^{-q-2})$:

For the I_2 part, note that

$$I_2 = \sup_{u \in \mathbb{R}} jP(I_{:,m} \leq u) - P(I_{Z,m} \leq u)j + \sup_{u \in \mathbb{R}} jP(I_{Z,m} \leq u) - P(jZ + \underline{dj}_1 \leq u)j =: I_{21} + I_{22}:$$

By Lemma 7 (1) for part $I_{21} = o(1)$ and Lemma 8 for I_{22} , we have

$$I_2 = O\left((bn)^{-1-6} \log^{7-6}(pn) + ((np)^{2-q} = (bn))^{1-3} \log(pn) + (m = (bn) + m) \right)^{1-3} \log(np)^{2-3}$$

For the I_3 part, the diagonal entities in bnQ take the same value i.e. $\sigma^2 = 2bn \sum_{i=1}^{bn} w_i^2$; which by (43), converges to $2c_w^2 > 0$, where c_w^d is a finite constant. By Lemma 4

$$I_3 \leq \log(np)^{1-2}.$$

The desired result follows by combining the I_1 - I_3 parts and a similar argument for the other side of the inequality.

For case (ii), we have the same decomposition I_1 - I_3 with $m = (bn)^{1-(1+\epsilon)}$. For the I_1 part, we define $\sigma^2 = c_1 \log(np)^{1-2} m^{-2} + c_0 (bn)^{1-2} b^2$; for some constant $c_1 > 0$. Then by Lemma 6, $I_1 = O((np)^{-q})$. For I_2 part, by Lemma 7 (2) and Lemma 8, we have

$$I_2 = O\left((bn)^{-1-6} \log^{7-6}(pn) + (m=(bn) + m^{-\epsilon})^{1-3} \log(np)^{2-3}\right)$$

For I_3 ; same argument can be applied. Combining the rates of I_1 - I_3 , we obtain the desired result. \square

Lemma 5 and 6 give us concentration inequalities for the m -dependent approximation of I .

LEMMA 5 (m -dependent approximation for polynomial case). *Assume conditions in Theorem 1 (i). For some $0 < m \leq n-2$ and $u > 0$; we have*

$$P\left(\left| \sum_{i=1}^m I_{i,m} - u \right| \geq c_1 (np)^{1-q} m^{-q}\right) \leq c_2 \left(e^{-c_3 u^2 m^2} + n p m^{-q} (bn)^{-q-2} u^{-q} \right);$$

where $c_1; c_2; c_3$ are some positive constants only depending on $q; C_p; C_w; C_S; \epsilon$.

Proof. We note that $I = \sum_{i=1}^m I_{i,m}$ can be bounded by

$$|I - \sum_{i=1}^m I_{i,m}| \leq \max_{i=1, \dots, m} \left(\left| \sum_{t=i}^{i-1} w_t \Lambda^{-1}(t, m) \right| + \left| \sum_{t=i+1}^{i+bn} w_t \Lambda^{-1}(t, m) \right| \right) =: I_1 + I_2;$$

We let $E_{i;l} = \sum_{t=(i-1)bn-(l+m)}^{i-1} w_{i-t} \Lambda^{-1} A_t$, then I_1 can be rewritten into

$$I_1 = \max_{\substack{bn+1 \leq i \leq n \\ 1 \leq j_1 \leq p}} \left| \sum_{\substack{l=1 \\ 1 \leq j_2 \leq p}}^{i-m-1} E_{i;l;j_1;j_2} \right|; \quad (48)$$

where $E_{i;l;j_1;j_2}$ is the $(j_1; j_2)$ th entity of matrix $E_{i;l}$ and $\lambda_{l;j_2}$ is the j_2 th entity of λ_l . Since $\lambda_{l;j_2}$ s are independent for different $(l; j_2)$, by Lemma A.2 in Chernozhukov et al. (2013b), for $u > 0$:

$$P(\overline{\rho_{bn} I_1} \geq 2 \overline{\rho_{bn} E I_1} + u) \leq e^{-u^2/(3 \cdot 2)} + K_q u^{-q} H_q; \quad (49)$$

where K_q is some constant only depending on q ;

$$H_q = \max_{\substack{bn+1 \leq i \leq n \\ 1 \leq j_1 \leq p}} \sum_{\substack{l=1 \\ 1 \leq j_2 \leq p}}^{i-m-1} E(E_{i;l;j_1;j_2} \lambda_{l;j_2})^2;$$

and

$$H_q = (bn)^{q-2} \sum_{\substack{l=1 \\ 1 \leq j_2 \leq p}}^{i-m-1} E \left(\max_{\substack{(bn+1) \leq i \leq n \\ 1 \leq j_1 \leq p}}^{(l+m+1)} E_{i;l;j_1;j_2} \lambda_{l;j_2}^q \right);$$

Then we start to analyze the rates of the objects involved in (49). We define $E_{i;l;j_1}$ to be the j_1 th row of $E_{i;l}$. For the $\lambda_{l;j_2}$ part, by Assumption 2.5 and (43),

$$E_{i;l;j_1} \lambda_{l;j_2} = \sum_{t=l+m}^{i-1} w_{i-t} \lambda_{j_1;j_1}^{1=2} A_{t;l;j_1;j_2} \leq c_w c_s m \leq (bn); \quad (50)$$

and

$$\sum_{\substack{l=1 \\ 1 \leq j_2 \leq p}}^{i-m-1} E_{i;l;j_1;j_2} \lambda_{l;j_2} = \sum_{t=i}^{i-1} \sum_{\substack{bn \leq l \leq t \\ t \leq m}} w_{i-t} \lambda_{j_1;j_1}^{1=2} A_{t;l;j_1;j_2} \leq c_w c_s m; \quad (51)$$

Combining the above arguments and recall that $E \lambda_{i;j}^2 = 1$; we have

$$H_q \leq \max_{\substack{bn+1 \leq i \leq n \\ 1 \leq j_1 \leq p}} \left(\sum_{\substack{l=1 \\ 1 \leq j_2 \leq p}}^{i-m-1} E_{i;l;j_1;j_2} \max_{\substack{1 \leq i \leq n \\ 1 \leq j_1 \leq p}} E_{i;l;j_1;j_2} \right) (c_w c_s)^2 m^2; \quad (52)$$

For the H_q part, by Assumption 2.5 and (43), $\max_{i,j_1,j_2} j E_{i;l;j_1;j_2} j \leq c_W c_S ((1-l) - m) = (bn)$: Recall that $\tilde{\rho} = c_p \rho$: Then we have

$$\begin{aligned}
H_q &= (bn)^{q-2} \sum_{\substack{l=1 \\ n \\ j_2 \\ p}}^{bn-m-1} [c_W c_S ((1-l) - m) = (bn)]^q \frac{q}{q} \\
&= (c_W c_S)^q \frac{q}{q} (bn)^{q-2} \tilde{\rho} \left(\sum_{m \leq l \leq n-bn-m} m^q + \sum_{l < m} (1-l)^q \right) = c_0 (bn)^{q-2} n p m^q;
\end{aligned} \tag{53}$$

where $c_0 = 3c_p (c_W c_S)^q \frac{q}{q}$:

For EI_1 part, note that $EI_1 \leq kI_1 k_q$: By Lemma 2, we have

$$EI_1 = \left(\sum_{i,j_1} E(j \sum_{l,j_2} E_{i;l;j_1;j_2} l j_2 j^q) \right)^{1=q} = \left(\sum_{i,j_1} ((q-1) \sum_l j E_{i;l;j_1;j_2} \frac{j^2}{q})^{q=2} \right)^{1=q};$$

Thus by (50) and (51) we have

$$EI_1 \leq (bn)^{1=2} m^q (np)^{1=q}; \tag{54}$$

where the constant in (54) only depends on c_W, c_S, q, q : Our conclusions follows by applying (52), (53) and (54) into (49) and a similar argument for I_2 . \square

LEMMA 6 (m-dependent approximation for exponential case). *We assume conditions in Theorem 1 (ii). We have*

$$P((bn)^{1=2} j l \leq l; m j \leq u) \begin{cases} 2n p e^{-a_1 m^2 - u^2}; & \text{if } u < a_2 (bn)^{1=2} m; \\ 2n p e^{-a_3 m (bn)^{1=2} u}; & \text{if } u \geq a_2 (bn)^{1=2} m; \end{cases}$$

where a_1, a_2, a_3 are some positive constants only depending on a_0, c_W, c_S, e :

Proof. Recall the definition of I_1 and I_2 in the proof of Lemma 5. Let $e = c_w c_s m = (bn)$ and $c = a_0 = e$: Then by (50), $E(e^{cE_{i;l;j_1;j_2} l;j_2}) < 1$; for any $0 < c < c$; and we have

$$E(e^{cI_1}) = \sum_{\substack{bn+1 \\ 1}} \sum_{\substack{j_1 \\ 1}} \sum_{\substack{n \\ p}} \sum_{\substack{bn \\ 1}} E \left(\exp \left\{ c \sum_{\substack{l \\ 1}} \sum_{\substack{n \\ j_2}} \sum_{\substack{bn \\ p}} \sum_{\substack{m \\ 1}} E_{i;l;j_1;j_2} l;j_2 \right\} + \exp \left\{ c \sum_{\substack{l \\ 1}} \sum_{\substack{n \\ j_2}} \sum_{\substack{bn \\ p}} \sum_{\substack{m \\ 1}} E_{i;l;j_1;j_2} l;j_2 \right\} \right) \\ =: I_{11} + I_{12}:$$

Since $E_{i;j} = 0$; for $E_{i;l;j_1;j_2} \neq 0$; we have

$$E(e^{cE_{i;l;j_1;j_2} l;j_2}) = 1 + \frac{E(e^{cE_{i;l;j_1;j_2} l;j_2} - 1) c E_{i;l;j_1;j_2} l;j_2}{c^2 E_{i;l;j_1;j_2}^2} c^2 E_{i;l;j_1;j_2}^2 \\ = 1 + \frac{E(e^{cE_{i;l;j_1;j_2} l;j_2} - 1) c E_{i;l;j_1;j_2} l;j_2}{(c E_{i;l;j_1;j_2} l;j_2)^2} c^2 E_{i;l;j_1;j_2}^2 \\ = 1 + \frac{e}{a_0^2} c^2 E_{i;l;j_1;j_2}^2;$$

where the first inequality is because that for any $x > 0$, the function $g_x(t) = (e^{tx} - 1) = t^2$ increases on $t \geq (0; 1)$; and $e^t - t = e^{tj} - jtj$: We define $c^d = e = a_0^2$; and the rate of I_{11} is derived as follows,

$$I_{11} = \sum_{i;j_1} \prod_{l;j_2} (1 + c^d c^2 E_{i;l;j_1;j_2}^2) \sum_{\substack{bn+1 \\ 1}} \sum_{\substack{j_1 \\ 1}} \sum_{\substack{n \\ p}} \sum_{\substack{bn \\ 1}} \exp \left\{ c^d c^2 \sum_{\substack{l \\ 1}} \sum_{\substack{n \\ j_2}} \sum_{\substack{bn \\ p}} \sum_{\substack{m \\ 1}} E_{i;l;j_1;j_2}^2 \right\}; \\ np \exp \{ c^2 c_1 m^2 = (bn) \};$$

where $c_1 = c^d c_w^2 c_s^2$; the second inequality is due to $1 + x \leq e^x$ for any $x \geq 0$; and the last inequality is by (52). Same bound can be derived for I_{12} : We note that

$$P(I_1 \leq u) = e^{-cu} E(e^{cI_1}) = e^{-cu} (I_{11} + I_{12});$$

We define $c = bnm^2 u = (2c_1)$: Hence if $c < c$; then $P(I_1 \leq u) \geq 2npe^{-u^2 m^2 bn = (4c_1)}$; if $c > c$; then $P(I_1 \leq u) \leq 2npe^{-bn a_0 = (c_w c_s) bnm u}$: The proof for I_2 is similar and therefore omitted. \square

LEMMA 7. Let $m \geq 1$ and $m = (bn) \geq 0$:

(1) Assume conditions in Theorem 1 (i), we have

$$\sup_{u \in \mathbb{R}} |P(I_{:,m}; u) - P(I_{Z;m}; u)| \leq C \cdot (bn)^{1-6} \log^{7-6}(pn) + ((np)^{2-q} = (bn))^{1-3} \log(pn);$$

where the constant in C only depends on c_w, c_w^d, c_s and q .

(2) Assume conditions in Theorem 1 (ii), we have

$$\sup_{u \in \mathbb{R}} |P(I_{:,m}; u) - P(I_{Z;m}; u)| \leq C \cdot (bn)^{1-6} \log(pn)^{7-6};$$

where the constant in C only depends on c_w, c_w^d, c_s, a_0 and e .

Proof. First we consider the case of (1). Denote

$$D_{i;l} = \sum_{t=(i-bn)_-}^{(i-1) \wedge (m+l-1)} w_{i-t} \Lambda^{-1} A_{t+l}; \quad \text{and} \quad D_{i;l} = \sum_{t=(i+1)_-}^{(i+bn) \wedge (m+l-1)} w_{t-i} \Lambda^{-1} A_{t+l}; \quad (55)$$

Then $I_{:,m}$ can be rewritten into

$$I_{:,m} = \max_{bn+1 \leq i \leq n-bn} \left| \sum_{i=m+1}^{bn+l} D_{i;l} \quad \sum_{i=m+2}^{l+i+bn} D_{i;l} + d_i \right|_1;$$

Let $N_0 = (n - 2bn)p$ and $N_1 = (n + m - 1)p$: Let $G = (G_{i;l})_{i;l: bn+1 \leq i \leq n-bn; 2 \leq m+l \leq n}$; be a block matrix in $\mathbb{R}^{N_0 \times N_1}$ with

$$G_{i;l} = \begin{cases} D_{i;l} & \text{if } i = m+1 \leq bn+l \leq i = m+1; \\ D_{i;l} \quad D_{i;l} & \text{if } i = m+2 \leq l \leq i-1; \\ D_{i;l} & \text{if } i = l \leq i+bn; \end{cases} \quad (56)$$

and elsewhere zero. We define d_{ij_1} to be the j_1 th entity of d_i , $N_2 = bnN_1$ and $G_{i;l;j_1;j_2}$ be the $(j_1;j_2)$ th entity of $G_{i;l}$. Then

$$N_2^{1=2} I_{:,m} = \max_{\substack{bn+1 \\ 1}} \max_{\substack{j_1 \\ j_1}} \max_{\substack{i \\ j_1}} \max_{\substack{n \\ p}} \left| \sum_{\substack{2 \\ 1}} \sum_{\substack{m \\ j_2}} \sum_{\substack{l \\ p}} \sum_{\substack{i \\ n}} g_{i;l;j_1;j_2} + N_2^{1=2} d_{ij_1} \right|; \text{ where } g_{i;l;j_1;j_2} = N_2^{1=2} G_{i;l;j_1;j_2} \quad (57)$$

For any $r \geq q$; we denote

$$M_r := \max_{\substack{bn+1 \\ 1}} \max_{\substack{j_1 \\ j_1}} \max_{\substack{i \\ j_1}} \max_{\substack{n \\ p}} |j_{j_1;r}| \quad \text{where} \quad |j_{j_1;r}| := \sum_{\substack{2 \\ 1}} \sum_{\substack{m \\ j_2}} \sum_{\substack{l \\ p}} \sum_{\substack{i \\ n}} \mathbb{E} |j g_{i;l;j_1;j_2}|^r = N_1 = \sum_{l=2}^n \sum_{m=1}^l |j G_{i;l;j_1;j_2}|^r = N_2^{r=2} = N_1:$$

By Assumption 2.5 and (43), for any $r \geq 2$;

$$|j D_{i;l;j_1;j_r}| \leq |j D_{i;l;j_1;j_2}| \leq c_w c_s = (bn); \text{ and similarly } |j D_{i;l;j_1;j_r}| \leq c_w c_s = (bn); \quad (58)$$

Then by (56), $\max_{i;l;j_1} |j G_{i;l;j_1;j_r}| \leq 2c_w c_s = (bn)$: Since $G_{i;l}$ is zero for $l < i - m + 1 - bn$ or $l > i + bn$;

$$M_r \leq (4c_w c_s)^r (N_1 = bn)^{1=2 - 1=r}; \quad (59)$$

Especially, for $r = 2$; $M_r \leq c_1$ where $c_1 = 4c_w c_s = 2$: By (55), for $i - bn - l - i - m$; and $S_m = \sum_{k=0}^{m-1} A_k$;

$$G_{i;l} = D_{i;l} = \sum_{t=l}^{m+l-1} w_{i-t} \Lambda^{-1} A_{t-l} = w_{i-l} \Lambda^{-1} S + w_{i-l} \Lambda^{-1} (S_m - S) + \sum_{t=l}^{l+m-1} (w_{i-t} - w_{i-l}) \Lambda^{-1} A_{t-l};$$

Therefore by Assumption 2.5 and (43), we have

$$\begin{aligned} |j G_{i;l;j_1;j_2} - w_{i-l}| &= \sum_{k=m}^{1=2} |j A_{k;j_1;j_2}| + \sum_{t=l}^{l+m-1} |j w_{i-t} - w_{i-l}| |j A_{t-l;j_1;j_2}| \\ c_w c_s m &= (bn) + c_w c_s m = (bn)^2; \end{aligned}$$

Note $m=(bn) = o(1)$, thus by (43),

$$\min_{\substack{1 \leq i \leq bn \\ 1 \leq j_1 \leq p}} \min_{1 \leq j_2 \leq bn} \left(\sum_{l=i}^m j G_{i;l;j_1;j_2}^2 N_2=N_1 \right)^{1=2} c_w^d \quad o(1) \quad c_1; \quad (60)$$

some constant $c_1 > 0$: Since $\max_{i;l;j_1;j_2} j G_{i;l;j_1;j_2} \leq 2c_w c_s = (bn)$;

$$\max_{l;j_2} E(\max_{i;j_1} j g_{i;l;j_1;j_2}^q) = \max_{i;l;j_1;j_2} j G_{i;l;j_1;j_2} N_2^{1=2} j^q \quad (2c_w c_s)(N_1=(bn))^{q=2}:$$

Note

$$B_n := \max \left\{ M_3^3; M_4^2; \left(\max_{l;j_2} E(\max_{i;j_1} j g_{i;l;j_1;j_2}^q) \right)^{1=q} \right\} \cdot (N_1=(bn))^{1=2};$$

where the constant in \cdot only depends on $c_w; c_s; q$: By Proposition 2.1 in Chernozhukov et al. (2017) we have

$$\begin{aligned} & \sup_{u \in \mathbb{R}} |P(N_2^{1=2} I_{z;m}(u) - P(N_2^{1=2} I_{z;m}(u)))| \\ & \leq (B_n^2 \log^7(pn) = N_1)^{1=6} + (B_n^2 \log^3(pn) = N_1^{1-2=q})^{1=3} \\ & \leq (bn)^{1=6} \log^{7=6}(pn) + ((np)^{2=q} = (bn))^{1=3} \log(pn): \end{aligned}$$

For part (2), let $M = \log_2(e) - 1$; and $B_n^0 = (2c_w c_s M = a_0)(N_1=(bn))^{1=2}$: Since $\max_{i;l;j_1;j_2} j G_{i;l;j_1;j_2} \leq 2c_w c_s = (bn)$; we have

$$\max_{i;l;j_1;j_2} E(e^{g_{i;l;j_1;j_2} = B_n^0}) \leq 2:$$

Note

$$B_n := \max \{ M_3^3; M_4^2; B_n^0 \} \cdot (N_1=(bn))^{1=2};$$

Apply the same argument as for part (1) with this new B_n ; then Proposition 2.1 in Chernozhukov et al. (2017) leads to

$$\sup_{u \in \mathbb{R}} |P(N_2^{1=2} I_{z;m}(u) - P(N_2^{1=2} I_{z;m}(u)))| \leq (B_n^2 \log^7(pn) = N_1)^{1=6} \cdot (bn)^{1=6} \log^{7=6}(pn):$$

□

LEMMA 8. Assume conditions in Theorem 1 (i) or (ii), for $m \leq 1$; $m=(bn) \leq 0$;

$$\sup_{u \in \mathbb{R}} j\mathbb{P}(I_{z;m} \leq u) - \mathbb{P}(jZ + dj_1 \leq u)j. \quad (m=(bn) + m) \leq 3 \log(np)^{2=3};$$

where the constant in \cdot only depends on $c_w; c_w^0$ and c_s :

Proof. We recall that $D_{i;l}; D_{i;l}$ in (55), $G = (G_{i;l})$ in (56) and G in (15). It is not hard to see that the covariance matrix for $I_{z;m}$ is $GG^>$ and the covariance matrix for Z is $Q = G G^>$:

We let

$$H^0 = (G_{i;l})_{\substack{2 \ m \ l \ 0 \\ bn+1 \ i \ n \ bn}}; \quad \text{and} \quad H^1 = (G_{i;l})_{\substack{1 \ l \ n \\ bn+1 \ i \ n \ bn}}.$$

Then $G = (H^0; H^1)$ and

$$\begin{aligned} jGG^> &= G G^> j_{\max} \quad jH^0 H^0^> j_{\max} + 2j(H^1 \ G)G^> j_{\max} + j(H^1 \ G)(H^1 \ G)^> j_{\max} \\ &=: I_1 + I_2 + I_3: \end{aligned}$$

By (58), $\max_{i;l;j} jG_{i;l;j} j_2 \leq 2c_w c_s = (bn)$: Therefore

$$(bn)I_1 \leq (bn) \max_{i_1; i_2; j_1; j_2} \sum_{l=2}^0 jG_{i_1;l;j_1} j_2 jG_{i_2;l;j_2} j_2 \leq (2c_w c_s)^2 m = (bn):$$

Denote $\Delta_{i;l} = G_{i;l} - G_{i;l}$: For $i = m+1 \leq bn - l < i \leq bn$; $\Delta_{i;l} = D_{i;l}$; and thus $j\Delta_{i;l;j} j_2 \leq c_w c_s = (bn)$: For $i \leq bn - l - i = m+1$; we have

$$\Delta_{i;l} = D_{i;l} - w_{i-l} \Lambda^{-1} S = \sum_{t=l}^{m+l-1} (w_{i-t} - w_{i-l}) \Lambda^{-1} A_{t-l} - w_{i-l} \Lambda^{-1} \sum_{t=m}^l A_t: \quad (61)$$

Hence $j\Delta_{i;l;j} j_2 \leq c_w c_s m = (bn)^2 + c_w c_s m = (bn)$: For $i = m+1 - l - i = 1$; $\Delta_{i;l} = D_{i;l} - D_{i;l} - w_{i-l} \Lambda^{-1} S$: Then $j\Delta_{i;l;j} j_2 \leq 3c_w c_s = (bn)$: Similarly we can bound $j\Delta_{i;l;j} j_2$ for $i = l - i + bn$:

For the rest l ; $\Delta_{i;l} = 0$: We note that $jG_{i;l;j} j \leq c_w c_s = (bn)$: Consequently,

$$(bn)I_2 \leq (bn) \max_{i_1; i_2; j_1; j_2} \sum_{l=1}^n j\Delta_{i_1;l;j_1} j_2 jG_{i_2;l;j_2} j_2 \leq m = (bn) + m;$$

where the constant in (61) only depends on c_w, c_s . Similarly we have $(bn)I_3 \leq m=(bn) + m$. Combining I_1 - I_3 ,

$$(bn)jGG^> G G^>_{j_{\max}} \cdot m=(bn) + m$$

By (43), for any j we have $bnQ_{j,j} = 2bn \sum_{i=1}^n w_i^2 \leq 2c_w^2$. Then the desired result follows from Lemma 3. \square

B.4 Proof of Theorem 2

Proof of (i). Note $1 - \Phi(x) = (2\pi)^{-1/2} \int_x^\infty e^{-x^2/2} dx$; where $\Phi(\cdot)$ is the cumulative distribution function of a standard normal distribution. Recall $G_{i,l}$ in (15). Let $G_i = (G_{i,1}; G_{i,2}; \dots; G_{i,n})$ and \underline{z} be a Gaussian vector in \mathbb{R}^{np} with zero mean and identity covariance matrix. Let $G_{i,j}$ be the j th row of G_i ; then

$$\mathbb{P}((bn)^{1-2} jG_{i,j} \cdot \underline{z} \geq u) = \sum_{i=bn}^n \sum_{j=1}^p \mathbb{P}((bn)^{1-2} jG_{i,j} \cdot \underline{z} \geq u) \leq np(2\pi)^{-1/2} \int_{u/(bn)^{1-2}}^\infty e^{-u^2/(2(bn)^{1-2})} du; \quad (62)$$

where $u = (bn)^{1-2} jG_{i,j} \cdot \underline{z} = (2bn \sum_{t=1}^{bn} w_t^2)^{1/2} \cdot \underline{z}$; which by (43) converges to $2^{1-2} c_w^2 > 0$. Thus

$$\mathbb{P}(jG_{i,j} \cdot \underline{z} \geq 2c_w^2 \log(np)^{1-2} (bn)^{1-2}) \leq 0; \quad (63)$$

Let $S := \{i \in [n] : j_i \cdot \underline{z} \geq 2c_w^2 \log(np)^{1-2} (bn)^{1-2}\}$. For any $i \in S$; $d_i = 0$. Hence by Theorem 1,

$$\sup_{u \in \mathbb{R}} \left| \mathbb{P}(\max_{i \in S} jV_{i,j} \geq u) - \mathbb{P}(\max_{i \in S} jG_{i,j} \cdot \underline{z} \geq u) \right| \leq 0; \quad (64)$$

Since $\max_{i \in S} jG_{i,j} \cdot \underline{z} \geq jG_{i,j} \cdot \underline{z}$; by (63) and (64) we have $\mathbb{P}(\max_{i \in S} jV_{i,j} \geq u) \leq 0$. Thus we obtain

$$\lim_{n \rightarrow \infty} \mathbb{P}(S \subseteq A_1; \mathbb{P}(S \subseteq A_1) \geq 1 - \epsilon) = 1; \quad (65)$$

Recall that $d_k = \Lambda^{-1} k$: Since $jd_k + G_{k; \underline{z}j_1} = jd_{kj_1} + jG_{k; \underline{z}j_1}$; we have

$$\begin{aligned} P\left(\min_{1 \leq k \leq K_0} jd_k + G_{k; \underline{z}j_1} \leq !^y\right) &= P\left(\max_{1 \leq k \leq K_0} jG_{k; \underline{z}j_1} \leq \min_{1 \leq k \leq K_0} jd_{kj_1} \leq !^y\right) \\ &= P(jG_{\underline{z}j_1} \leq !^y); \end{aligned}$$

Therefore $P(\min_{1 \leq k \leq K_0} jd_k + G_{k; \underline{z}j_1} \leq !^y) \rightarrow 0$: Subsequently the break statistics will be bigger than the threshold at the points of break with probability approach 1,

$$P\left(\min_{1 \leq k \leq K_0} jV_{kj_1} \leq !^y\right) \rightarrow 0;$$

in view of

$$\sup_{u \in \mathbb{R}} \left| P(jV_{kj_1} \leq u) - P(jd_k + G_{k; \underline{z}j_1} \leq u) \right| \rightarrow 0;$$

Therefore we have

$$P(k \geq A_1; 1 \leq k \leq K_0) \rightarrow 1. \tag{66}$$

Let $B(k; r) = \{t: jt \leq j \leq rg\}$: By (65) and (66), we have

$$\lim_n P\left(\bigcap_{1 \leq k_1 < k_2 \leq K_0} g_{A_1} \in [1 \leq k \leq K_0] B(k; bn)\right) = 1;$$

Since for $k_1 \neq k_2; j_{k_1} \leq k_2 j \leq bn$; for any $k_1 \neq k_2$ and $t \in B(k_1; bn)$; for all large n ; $B(t; 2bn) \setminus B(k_2; 2bn) = \emptyset$: Thus we complete the proof. \square

Proof of (ii). First consider the case (1). Let $\hat{Y}_i^{(l)}$ (resp. $U_i^{(l)}$) be $\hat{Y}_i^{(l)}$ with Y_i therein replaced by $\hat{Y}_i^{(l)}$ (resp. $U_i^{(l)}$). Similarly we can define $\hat{Y}_i^{(r)}$ and $U_i^{(r)}$: Let $\Delta_i^{(l)} = \hat{Y}_i^{(l)} - U_i^{(l)}$ and $\Delta U_i = U_i^{(l)} - U_i^{(r)}$: Let $\Delta \hat{f}_i$ be $\Delta_i^{(l)}$ with \hat{Y}_i replaced by \hat{f}_i :

For any $1 \leq k \leq K_0$, and any t such that $jt \leq j \leq bn$; we have $\Delta_t = (1 - \sum_{i=1}^{jt} w_i) \Delta_k + \Delta f_t$: Hence

$$\begin{aligned} V_t &= \Lambda^{-1} \Delta_t + \Lambda^{-1} \Delta U_t \\ &= (1 - \sum_{i=1}^{jt} w_i) \Lambda^{-1} \Delta_k + \Lambda^{-1} \Delta f_t + \Lambda^{-1} \Delta U_k + (\Lambda^{-1} \Delta U_t - \Lambda^{-1} \Delta U_k): \end{aligned} \quad (67)$$

Note $\hat{k} = \operatorname{argmax}_{1 \leq k \leq bn} V_{tj}$: The proceeding proof contains three steps.

Step 1. Let $j_k = \operatorname{argmax}_j V_{kj}$; where V_{kj} is the j th entity of V_k : This step shows

$$\liminf_n \min_{1 \leq k \leq K_0} j(\Lambda^{-1})_{j_k} = 1:$$

We shall show by contradiction. By (47), $j \Delta f_{tj} = O(b^2)$: If there exists $1 \leq k \leq K_0$, such that $j(\Lambda^{-1})_{j_k} < c$; for some $c < 1$; then by (67), $j V_{kj} \leq c + O(b^2) + j \Lambda^{-1} \Delta U_{kj}$: Let \tilde{U}_t be U_t with U_{ij} replaced by Z_{ij} where Z_{ij} are i.i.d standard normal random variables. Then $\max_t j \Lambda^{-1} \Delta \tilde{U}_{tj} = O_p((bn)^{-1/2} \log(np)^{1/2})$: Then by Gaussian approximation Theorem 1,

$$\max_{1 \leq k \leq K_0} \max_{jt \leq j \leq bn} j \Lambda^{-1} \Delta U_{tj} = O_p((bn)^{-1/2} \log(np)^{1/2}):$$

Since $(bn)^{-1/2} \log(np)^{1/2} < c$; we have $j V_{kj} \leq c (1 + o_p(1))$: On the other hand, by (67), $j V_{kj} \geq c + O(b^2) - j \Lambda^{-1} \Delta U_{kj} = c (1 + o_p(1))$: These imply $P(V_{\hat{k}} < V_k) \rightarrow 1$; which is a contradiction.

Step 2. This step shows

$$\max_{1 \leq k \leq K_0} \max_{jt \leq j \leq bn} j \Lambda^{-1} \Delta U_k - \Lambda^{-1} \Delta U_{tj} = o_p(1) = O_p((np)^{-1/2} \log(np)^{1/2}):$$

Let $t = k$; the other direction can be similarly dealt with. Note that

$$\begin{aligned} \Delta U_t - \Delta U_k &= \sum_{i=t}^{k-1} W_{t+i} + \sum_{i=k}^{t-1} (W_{t+i} - W_{k+1+i}) - \sum_{i=t+1}^{k-1} (W_{t+i} + W_{k+i}) \\ &\quad - \sum_{i=k+1}^{t+bn} (W_{t+i} - W_{k+i}) + \sum_{i=t+bn+1}^{k+bn} W_{k+i} + W_{k+t} - W_{k+t-k} =: \sum_{k=1}^7 r_k. \end{aligned}$$

For r_1 we have

$$j\Lambda^{-1}r_1j_1 = \max_{1 \leq j_1 \leq p} \left| \sum_{l=k-1}^{k-1} E_{l;j_1;j_2} \right|; \text{ where } E_l = \sum_{i=(k-1)-l}^{k-1} W_{t+i} \Lambda^{-1} A_{i,j_1};$$

and $E_{l;j_1;j_2}$ is $(j_1;j_2)$ th entity of matrix E_l . Then

$$\max_{1 \leq k \leq j_1} \max_{K_0 \leq j_2 \leq j_1} j\Lambda^{-1}r_1j_1 \leq j_1 t^{-1} = O_p(f(np)^{1-q} = (bn)g);$$

uniformly over $t; k$: A similar argument leads to the same bound for r_3 and r_5 : For r_2 ; we can rewrite

$$j\Lambda^{-1}r_2j_1 = \max_{1 \leq j_2 \leq p} \left| \sum_{l=t-1} E_{l;j_1;j_2} \right|; \text{ where } E_l = \sum_{i=(k-1)-l}^{t-1} (W_{t+i} - W_{k+i}) \Lambda^{-1} A_{i,j_1};$$

Then similarly we have $\max_{1 \leq k \leq j_1} \max_{K_0 \leq j_2 \leq j_1} j\Lambda^{-1}r_2j_1 \leq j_1 t^{-2} = O_p(f(np)^{1-q} = (bn)^2g)$: We obtain the desired result by summing up all the above bounds.

Step 3. Without loss of generality, assume $j_k > 0$: Then by the argument in step 1, with probability tending to 1, $V_{k;j_k} > 0$: By Assumption 2.1, we have $j\Delta f_t - \Delta f j_1 = O(jt^{-1} j=n)$; uniformly over t : With probability tending to 1, by (67),

$$\begin{aligned} jV_{k;j_1} - jV_{k;j_1} &= V_{k;j_k} - V_{k;j_k} \\ &= \sum_{i=1}^{j_k} W_i (\Lambda^{-1})_{j_k} = O(j_k^{-1} j=n) = j\Lambda^{-1}\Delta U_k - \Lambda^{-1}\Delta U_{k;j_1}; \end{aligned}$$

By (43), we have $\sum_{i=1}^{jt} w_i = c_w^0 j t = j = (bn)$: For finite moment case, by Step 1 and Step 2 we further derive

$$jV_{k^j} j_1 = jV_{\wedge_k} j_1 = c_w^0 j_{\wedge_k} = O(j_{\wedge_k} = n) = O_P(j_{\wedge_k}^{1-2} (np)^{1-q} = (bn));$$

uniformly over k : Since $jV_{k^j} j_1 < jV_{\wedge_k} j_1$; we have

$$\max_{1 \leq k \leq K_0} j_{\wedge_k} = O_P f(np)^{2-q} = 2g;$$

Similar argument can be applied for sub exponential case. □

Proof of (iii). Recall the definition of $\binom{(\cdot)}{t}$; $\binom{(\cdot)}{t}$; $U_t^{(\cdot)}$ and $U_t^{(r)}$ in the proof of (ii) and $M = bn$: Since $M = \log(np) = 2$;

$$j \binom{(\cdot)}{\wedge_k M} ((\wedge_k M) = n) j_1 = j f_{\wedge_k M}^{(\cdot)} f((\wedge_k M) = n) j_1 = O(b^2):$$

Similarly $j \binom{(\cdot)}{\wedge_{k+M}} ((\wedge_{k+M}) = n) j_1 = O(b^2)$: Since $\max_{1 \leq j \leq p} j f_j^{(\cdot)}$ is bounded,

$$j \binom{(\cdot)}{(\wedge_k + M) = n} ((\wedge_k M) = n) j_1 = j f_{(\wedge_k + M) = n} f((\wedge_k M) = n) j_1 = O(M = n):$$

Hence

$$\begin{aligned} j \Lambda^{-1}(\wedge_k M) j_1 &= \left| \Lambda^{-1} \binom{(\cdot)}{\wedge_{k+M} M} \binom{(\cdot)}{\wedge_k M} + \Lambda^{-1} U_{\wedge_{k+M}}^{(r)} - \Lambda^{-1} U_{\wedge_k}^{(\cdot)} \right|_1 \\ &= O(b + M = n) + j \Lambda^{-1} U_{\wedge_k M}^{(\cdot)} - \Lambda^{-1} U_{\wedge_{k+M}}^{(r)} j_1 : \end{aligned} \quad (68)$$

By Gaussian approximation and (63) we have

$$P(j \Lambda^{-1} U_{\wedge_k M}^{(\cdot)} - \Lambda^{-1} U_{\wedge_{k+M}}^{(r)} j_1 \geq 2c_w^0 \log(np)^{1-2} = (bn)^{1-2}) \rightarrow 0:$$

Inserting the above equation into (68) and we obtain the desired result. □

B.5 Proof of Theorem 3

Proof. We recall $M = bn$: Similar to (68), we have $j\Lambda^{-1}(\hat{\kappa}_k) \Lambda^{-1}(U_{\hat{\kappa}_k+M}^{(r)} - U_{\hat{\kappa}_k}^{(l)})j_1$
 $cM=n$: Therefore

$$\begin{aligned} & \sup_{u \in \mathbb{R}} \left| \mathbb{P}((bn)^{1/2} j\Lambda^{-1}(\hat{\kappa}_k) j_1 - u) - \mathbb{P}((bn)^{1/2} j\tilde{Z}j_1 - u) \right| \\ & \sup_{u \in \mathbb{R}} \mathbb{P}(j(bn)^{1/2} j\tilde{Z}j_1 - u) - c(bn)^{1/2} M=n \\ & + \sup_{u \in \mathbb{R}} \left| \mathbb{P}(j\Lambda^{-1}(U_{\hat{\kappa}_k}^{(l)} - U_{\hat{\kappa}_k+M}^{(r)})j_1 - u) - \mathbb{P}(j\tilde{Z}j_1 - u) \right| = I_1 + I_2: \end{aligned}$$

We note that $(bn)^{1/2} \tilde{Z}_j$ are i.i.d with variance $2(bn) \sum_{t=1}^{bn} w_t^2$; which by (43) converges to $2c_w^2 > 0$. Therefore by Lemma 4,

$$I_1 = O(bn)^{1/2} (M=n) \log(np)^{1/2} g = o(1):$$

Let $\tilde{G} = (\tilde{G}_1; \tilde{G}_2; \dots; \tilde{G}_n)$; where $\tilde{G}_l = w_{\hat{\kappa}_k - M - l} \Lambda^{-1} S$; if $\hat{\kappa}_k - M - bn - l < \hat{\kappa}_k - M - 1$; and $\tilde{G}_l = w_l (\hat{\kappa}_k + M) \Lambda^{-1} S$; if $\hat{\kappa}_k + M + 1 - l < \hat{\kappa}_k + M + bn$ and elsewhere zero. Let \underline{z} be Gaussian vector in \mathbb{R}^{np} with zero mean and identity covariance matrix. Then $\tilde{G}\underline{z} \stackrel{d}{=} \tilde{Z}$: By the same argument as in Theorem 1 with $d_i = 0$ we have

$$\sup_{u \in \mathbb{R}} \left| \mathbb{P}(j\Lambda^{-1}(U_{\hat{\kappa}_k}^{(l)} - U_{\hat{\kappa}_k+M}^{(r)})j_1 - u) - \mathbb{P}(j\tilde{G}\underline{z}j_1 - u) \right| = o(1):$$

Thus $I_2 = o(1)$ and we complete the proof. \square

B.6 Proof of Theorem 4

Proof of (i). We shall condition on the event where $\hat{S}_k = S_k$ and $j\hat{\kappa}_k - \kappa_j - bn$: By Theorem 2 and Corollary 2, the event would take place with probability tending to 1.

Denote $"_t = \sum_{j \in 2S_k} (\Lambda^{-1} \hat{\kappa}_k)_j (\Lambda^{-1} t)_j$; and $\hat{\alpha}_k = \sum_{j \in 2S_k} (\Lambda^{-1} \hat{\kappa}_k)_j (\Lambda^{-1} \hat{\kappa}_k)_j$: Then we have

$$X_t = \hat{\alpha}_k \mathbf{1}_{t = \hat{\kappa}_k} + \sum_{j \in 2S_k} f_j(t=n) (\Lambda^{-1} \hat{\kappa}_k)_j + "_t:$$

Let $I(t) = \sqrt{(4bn+1)(t-t_0)(4bn+1-t-t_0)}$, $t_0 = \hat{k} - 2bn$, $t_1 = \hat{k} + 2bn$ and

$$D_t = \left(\sum_{s=t_0}^{t_1} X_s \frac{t-t_0}{4bn+1} - \sum_{s=t_0}^{t-1} X_s \right) I(t):$$

For any $r > 0$; we have

$$\begin{aligned} D_{k+r} - D_k &= \left(\sum_{s=t_0}^{t_1} X_s \frac{k-t_0}{4bn+1} - \sum_{s=t_0}^{k-1} X_s \right) (I(k+r) - I(k)) \\ &\quad + \left(\sum_{s=t_0}^{t_1} X_s \frac{r}{4bn+1} - \sum_{s=k}^{k+r-1} X_s \right) I(k+r) = I_1 + I_2: \end{aligned}$$

Denote $I_i(f)$ (resp. $I_i(\cdot)$, $I_i(a)$) to be I_i with X_t therein replaced by $f(t=n)$ (resp. \cdot , $\hat{a}_k \mathbf{1}_{t \leq k}$), $i = 1, 2$.

Firstly, consider the f part. Note $|I(k+r) - I(k)| \leq (bn)^{3/2} r$. Thus by the continuity of f_j ; for I_1 part, we have

$$\begin{aligned} \max_j \left| \sum_{s=t_0}^{t_1} f_j(s=n) \frac{k-t_0}{4bn+1} - \sum_{s=t_0}^{k-1} f_j(s=n) \right| |I(k+r) - I(k)| &= O((bn)^{3/2} r) \\ &= O(br(bn)^{1/2}): \end{aligned}$$

Similarly we can handle the f part in I_2 and therefore

$$|I_1(f) + I_2(f)| \leq O(br(bn)^{1/2} \tilde{j}_k^1):$$

Secondly, let us consider the drift part, for $r \leq bn$;

$$\begin{aligned} &I_1(a) + I_2(a) \\ &= (t_1 - k + 1) \hat{a}_k \frac{k-t_0}{4bn+1} (I(k+r) - I(k)) + \left((t_1 - k + 1) \hat{a}_k \frac{r}{4bn+1} - r \hat{a}_k \right) I(k+r) \\ &= (bn)^{1/2} \hat{a}_k r = 2(1 + o(1)): \end{aligned}$$

Thirdly, let us focus on the " part. By Theorem 2 (iii), $j(\Lambda^{-1} \hat{\Lambda}_k)_{j_2 S_k} j_2 = \tilde{j}_k j_2 (1 + o_P(1))$ and thus $\hat{\Delta}_k = \tilde{j}_k j_2 (1 + o_P(1))$: Then together with (35), we obtain that the long run variance for " k is $\hat{\Delta}_k^2 (1 + o_P(1))$: Hence by Theorem 1 in El Machkouri et al. (2013),

$$\left| \sum_{s=\hat{\Lambda}_k}^{\hat{\Lambda}_k + 2bn} "s \right| = O_P((bn)^{1=2} \hat{\Delta}_k \log(bn)^{1=2}) \quad \text{and} \quad \left| \sum_{s=k}^{k+r} "s \right| = O_P(r^{1=2} \hat{\Delta}_k):$$

Therefore

$$I_1(") = O_P((bn)^{1=2} \hat{\Delta}_k \log(bn)^{1=2} j_l(k+r) \quad l(k)j) = O_P((bn)^{-1} r \hat{\Delta}_k \log(bn)^{1=2}):$$

For I_2 part, we have

$$\begin{aligned} I_2(") &= \left(\sum_{s=t_0}^{t_1} "s \frac{r}{4bn+1} \quad \sum_{s=k}^{k+r} "s \right) l(k+r) \\ &= O_P\left((bn)^{-1} r \hat{\Delta}_k \log(bn)^{1=2} + r^{1=2} \hat{\Delta}_k (bn)^{-1=2} \right): \end{aligned}$$

Combining all the previous parts we have

$$\begin{aligned} &D_{k+r} \quad D_k \\ = &(bn)^{-1=2} a_k r (1 + o_P(1)) = 2 + O(br(bn)^{-1=2} \tilde{j}_k j_1) + O_P\left((bn)^{-1} r \hat{\Delta}_k \log(bn)^{1=2} + r^{1=2} \hat{\Delta}_k (bn)^{-1=2} \right) \end{aligned} \tag{69}$$

Note $\tilde{j}_k j_1 = j_{S_k} j_1^{1=2} \tilde{j}_k j_2$; thus

$$\tilde{j}_k j_1 b = j_{S_k} j_1^{1=2} \tilde{j}_k j_2 b = j_{S_k} j_1^{1=2} \tilde{j}_k j_2^y = \tilde{j}_k j_2^2 = a_k:$$

Therefore in (69), we have $O(br(bn)^{-1=2} \tilde{j}_k j_1) = o((bn)^{-1=2} a_k r)$:

Since $\tilde{\Sigma}_k$ is a covariance matrix with diagonal entities 1; $j_{S_k} j_2 = j_{S_k} j_2$: Note $\hat{\Delta}_k^2 = j_{S_k} j_2 a_k$: thus

$$\hat{\Delta}_k = j_{S_k} j_1^{1=2} a_k^{1=2}:$$

Then we have

$$\&_k(bn)^{1=2} \log(bn)^{1=2} \quad jS_{kj}^{1=2} a_k^{1=2} (bn)^{1=2} \log(bn)^{1=2} \quad jS_{kj}^{1=2} a_k^{1=2} \quad a_k;$$

where the last inequality is because $a_k = \tilde{j}_{kj}^2 \quad j^2 jS_{kj}$. Therefore $O_p((bn)^{-1} r \&_k \log(bn)^{1=2}) = o_p((bn)^{-1=2} a_k r)$ in (69). Inserting the above equations into (69) leads to

$$(bn)^{1=2} (D_{k+r} - D_k) = a_k r (1=2 + o_p(1)) + O_p(\&_k r^{1=2});$$

Since D_{-k} is the maximum, $D_{k+r} - D_k > 0$: Therefore $r = O_p(\&_k^2 = a_k^2)$: By a similar argument for the $r < 0$ part, the desired result follows. \square

Proof of (ii). Let F_t be the \mathbb{R} -field generated by $f_{s;j}; s = t; 1 \dots j = p$: Denote the projection operator $P_t = E(jF_t) - E(jF_{t-1})$: Let $\tilde{-}_{k;j} = (\Lambda^{-1} k)_j$; if $j \geq S_k$, $\tilde{-}_{k;j} = 0$ if $j \not\geq S_k$. Let (\tilde{t}) be an i.i.d copy of (t) . Then

$$c_t := k P_0 \tilde{t} k_4 \quad k^{-> \Lambda^{-1} A_t} (0 \quad \tilde{t}) k_4 \cdot j^{-> \Lambda^{-1} A_t} j_2 \quad 4;$$

where the last inequality is by Lemma 2 and that $\tilde{t}; 1 \dots j = p$ are i.i.d. By Assumption 2.5,

$$\sum_{s=m} c_s \cdot \sum_{s=m} j^{-> \Lambda^{-1} A_s} j_2 \quad \sum_{j=1}^p \sum_{s=m} \tilde{-}_{k;j} \quad j_j^{1=2} j A_{s;j} j_2 \cdot j^{-> \Lambda^{-1} A_s} j_1 m = \tilde{j}_{kj}^2 m \quad ;$$

Thus by Corollary 2.1 in Berkes et al. (2014), strong invariance principle holds for $\sum_{s=t} \tilde{t}_s$: Thus similar to (69), we have

$$(bn)^{1=2} (D_{k+r} - D_k) = \tilde{2}^{-1} a_k r (1 + o_p(1)) + I_2(\tilde{t}) \quad P \quad \tilde{2}^{-1} a_k r + \&_k B(r);$$

The $r < 0$ part can be similarly dealt with. \square

B.7 Proof of Theorem 5

Proof of (i). The main idea follows the proof of Proposition 2.4 in Catoni (2012), however due to the dependence and the break points, our result is much more involved. Let

$$S = \{k \mid A_k \text{ or } A_{k-1} \text{ contains break points}\}.$$

Then by assumption $|S| \geq 2K_0$: We look at estimators without the break point first.

$$\bar{h}_{ij}(u) = \sum_{k \in S} \hat{h}_{ij;k}(u) = N_2^{-1} \sum_{k \in S} \hat{h}_{ij;k}(u) \quad \text{where } N_2 = N_1 - |S|.$$

Let

$$\tilde{\mu}_{ij} = \sum_{k \in S} E \hat{h}_{ij;k}(u) \quad \text{and} \quad \tilde{v}_{ij}^2 = \sum_{k \in S} E \hat{h}_{ij;k}^2(u) - \tilde{\mu}_{ij}^2.$$

Define functions

$$\begin{aligned} B_{ij}^+(u; \chi) &= \tilde{\mu}_{ij} + u + \tilde{v}_{ij} [(\tilde{\mu}_{ij} - u)^2 + \tilde{v}_{ij}^2]^{-1/2} + \chi; \\ B_{ij}^-(u; \chi) &= \tilde{\mu}_{ij} - u - \tilde{v}_{ij} [(\tilde{\mu}_{ij} - u)^2 + \tilde{v}_{ij}^2]^{-1/2} - \chi; \end{aligned}$$

The proof contains four steps.

Step 1. This step shows that function $E \bar{h}_{ij}(u)$ for any i, j satisfies, the expected loss functions have upper and lower envelope functions,

$$B_{ij}^-(u; 0) \leq E \bar{h}_{ij}(u) \leq B_{ij}^+(u; 0).$$

By (40), $(x) = x + x^2/2$ and thus

$$E \bar{h}_{ij}(u) = \sum_{k \in S} (E \hat{h}_{ij;k}(u) + \tilde{v}_{ij} E (\hat{h}_{ij;k}(u) - \tilde{\mu}_{ij})^2/2) = N_2^{-1} \sum_{k \in S} (E \hat{h}_{ij;k}(u) + \tilde{v}_{ij} E (\hat{h}_{ij;k}(u) - \tilde{\mu}_{ij})^2/2) = B_{ij}^+(u; 0).$$

Similarly we can bound the other side.

Step 2. This step shows for any $\varepsilon > 0$, the estimated influence function $\bar{h}_{ij}(u)$ is highly concentrated around its mean, for $C_0 > 0$ and $\varepsilon \leq (N_2 \log(N_2))^{1/2}$;

$$\sum_{i,j=1}^p \mathbb{P} \left(\sup_{|u| \leq C_0} |\bar{h}_{ij}(u) - \mathbb{E} \bar{h}_{ij}(u)| \geq \varepsilon \right) \leq p^2 (N_2 \log(n))^{q-4} \varepsilon^{q-2} + e^{-\varepsilon^2/(cN_2)}; \quad (70)$$

where c and the constant in (70) are independent of n, p :

First introduce some notation. For any random variable X ; denote $E_0 X = X - \mathbb{E} X$; the centering operator. Let $F_k = (Y_t; t \geq [s, k] A_s)$ and $F_{k;fsg}; s \leq k$ be F_k with $Y_t; t \geq A_s$ therein replaced by Y_t^0 ; where Y_t^0 are i.i.d copy of Y_t ; For any random variable $X = h(F_k)$; let $X_{fig} = h(F_{k;fig})$; Denote $\Delta_k = X_k - X_{k-1}$; We now show that the temporal dependence measure decays with polynomial rate. Let $\Delta_{ij;k}(u) = \Delta_{ij;k}(\hat{Y}_{ij;k}(u))$; Since $\|Y_t^0\|_1 \leq 1$, we have for any $s \geq 2$ and any u ;

$$\begin{aligned} & \mathbb{E} \sup_u |\Delta_{ij;k}(u) - \Delta_{ij;k;fsg}(u)|_{q=2} \leq \mathbb{E} \|\Delta_{ij;k} - \Delta_{ij;k;fsg}\|_{q=2} \\ & \leq 2^{-1} m \left(\mathbb{E} \|\Delta_{k;i}(\Delta_{k;j} - \Delta_{k;j;fsg})\|_{q=2} + \mathbb{E} \|\Delta_{k;i}(\Delta_{k;j;fsg} - \Delta_{k;j;fsg})\|_{q=2} \right) \\ & =: 2^{-1} m(I_1 + I_2); \end{aligned}$$

Let $U_{k;i}$ (resp. $f_{k;i}$) be $Y_{k;i}$ with Y_t replaced by Y_t^0 (resp. $f(t=n)$). Then $U_{k;i} = U_{k;i} + f_{k;i}$ when there is no break. Let $\Delta U_{k;i} = U_{k;i} - U_{k-1;i}$ and $\Delta f_{k;i} = f_{k;i} - f_{k-1;i}$; Then we have

$$I_1 = \left\| \mathbb{E} \|\Delta f_{k;i}(\Delta U_{k;j} - \Delta U_{k;j;fsg})\|_{q=2} \right\| + \left\| \mathbb{E} \|\Delta U_{k;i}(\Delta U_{k;j} - \Delta U_{k;j;fsg})\|_{q=2} \right\| =: I_{11} + I_{12}; \quad (71)$$

Since $\max_j \|Y_t^0\|_1 \leq 1$;

$$\max_{1 \leq j \leq p} \|\Delta f_{k;j}\|_1 \leq m/n; \quad (72)$$

Let $E_{k;l;i} = \sum_{t=(k+1)_-l}^{(k+1)m} A_{t-l;i}$; where $A_{t-l;i}$ is the i th row of matrix A_{t-l} . Then

$$U_{k;i} = \sum_{l=(k+1)m} E_{k;l;i}$$

and

$$\Delta U_{k;i} \Delta U_{k;j;fk_{sg}} = \begin{cases} \sum_{l_2 A_{k-s}} (E_{k;l_1;i} E_{k-l_1;l_2;j}) \binom{l_1}{i} \binom{l_2}{j} = m; & s=1; \\ \sum_{l_2 A_{k-s}} E_{k;l_1;i} \binom{l_1}{i} \binom{l_2}{j} = m; & s=0; \end{cases} \quad (73)$$

Since l_i are i.i.d, by Lemma 2, (72) and (73),

$$k \Delta f_{k;i} (\Delta U_{k;j} \Delta U_{k;j;fk_{sg}})_{q=2} \leq 2f c_q \left(\sum_{l_2 A_{k-s}} (j E_{k;l_1;j_2} + E_{k-l_1;l_2;j_2})^2 \right)^{1=2} \quad q=2=n; \quad (74)$$

By Assumption 2.5, we have for any $s > 1$;

$$\sum_{l_2 A_{k-s}} j E_{k;l_1;j_2} \cdot m(m(s-1))^{1=2} \binom{l_1}{i} \binom{l_2}{j} \quad \text{and} \quad \sum_{l_1 (k+1)m} j E_{k;l_1;j_2} \cdot m \binom{l_1}{i} \binom{l_2}{j} \quad (75)$$

where the constant in \cdot only depending on $i; c_s$. Hence by (74) and (75),

$$I_{11} \leq m^{1=2} n^{-1} (m(s-1))^{1=2} \binom{l_1}{i} \binom{l_2}{j}; \quad (76)$$

where the constant in \cdot only depends on $i; c_s; q; q; f$. By Lemma 2 and (75)

$$\begin{aligned} k E_0 U_{k;i} (U_{k;j} U_{k;j;fk_{sg}})_{q=2} &= \left\| E_0 \left(\sum_{l_1 (k+1)m} E_{k;l_1;i} \binom{l_1}{i} \sum_{l_2 A_{k-s}} E_{k;l_2;j} \binom{l_2}{j} \right) \right\|_{q=2} m^2 \\ &\leq m^2 \left(\sum_{l_1 (k+1)m} \sum_{l_2 A_{k-s}} j E_{k;l_1;i} \binom{l_1}{i} j E_{k;l_2;j} \binom{l_2}{j} \right)^{1=2} \\ &\leq m^{-1} (m(s-1))^{1=2} \binom{l_1}{i} \binom{l_2}{j}; \end{aligned} \quad (77)$$

where the constant in \cdot only depends on $q; q; c_s$. Thus $I_{12} \leq m^{-1} (m(s-1))^{1=2} \binom{l_1}{i} \binom{l_2}{j}$.

By combining the bounds for I_{11} and I_{12} and a similar argument for I_2 ; we have

$$s := \max_k \left\| \sup_u j_{i;j;k}(u) \quad j_{i;j;k;fk_{sg}}(u) \right\|_{q=2} \leq ((ms) \mathbf{1}_{s>1} + \mathbf{1}_s)^{1=2} \binom{l_1}{i} \binom{l_2}{j}; \quad (78)$$

where the constant in \cdot only depends on $q; C_S; C; \cdot; q; f$:

Let $\cdot := \frac{1=2}{i;i} \frac{1=2}{j;j} \chi=(2N_2)$ and A_n be the net for $f_u : j_u \quad i;j \quad C_0g$: Denote $f(u) = \bar{h}_{i;j}(u) \quad E\bar{h}_{i;j}(u)$: Then by $j \quad \theta j_1 \quad 1$:

$$\sup_{j_v \quad i;j} \min_{C_0} \min_{u \in A_n} jf(u) \quad f(v)j \quad \cdot$$

Therefore $jA_nj = 2C_0 = O(n)$ and

$$P\left(\sup_{j_u \quad i;j \quad C_0} j\bar{h}_{i;j}(u) \quad E\bar{h}_{i;j}(u)j \quad \chi\left(\frac{1=2}{i;i} \frac{1=2}{j;j}\right)=N_2\right) \quad P\left(\max_{u \in A_n} j\bar{h}_{i;j}(u) \quad E\bar{h}_{i;j}(u)j \quad \chi\left(\frac{1=2}{i;i} \frac{1=2}{j;j}\right)=(2N_2)\right):$$

Desired result follows from Lemma 5.8 in Zhang et al. (2017).

Step 3. This step shows for the estimator

$$\max_{1 \leq i,j \leq p} j\tilde{f}_{i;j} \quad i;j = O(m^{-(+1)} \frac{1=2}{i;i} \frac{1=2}{j;j} + m^3=r^2); \quad \text{and} \quad v_{i;j}^2 = O\left(\frac{1=2}{i;i} \frac{1=2}{j;j}\right); \quad (79)$$

where the convergence is uniform for $1 \leq i,j \leq p$:

Let $\hat{f}_{i;j,k}$ be $\tilde{f}_{i;j,k}$ with Y_t replaced by \cdot_t and let $\hat{f}_{i;j} = E\hat{f}_{i;j,1}$: Then by (72),

$$j\tilde{f}_{i;j} \quad i;j \quad m \sum_{k \in S} j\Delta f_{k;i;j} j\Delta f_{k;j}j=(2N_2) = O(m^3=r^2): \quad (80)$$

Note the convergence in above $O(\cdot)$ and all the followings are uniform for i,j : Let $\hat{f}_{i;j,k} = E(\cdot_{0;i} \cdot_{k;j})$: Then for any $L < m$:

$$jmE(U_{1;i}U_{1;j}) \quad i;j = \left| m^{-1} \sum_{m < k < m} \binom{m-jk}{i;j;k} \hat{f}_{i;j,k} \quad i;j \right| = O\left(\sum_{j;k \leq L} j \hat{f}_{i;j;k} + Lm^{-1} \sum_{k \in Z} j \hat{f}_{i;j;k}\right):$$

By Assumption 2.5, $\sum_{j;k \leq L} j \hat{f}_{i;j;k} \quad \sum_{t \in Z; j;k \leq L} jA_{t;i} j_2 jA_{t+k;j} j_2 = O(L \frac{1=2}{i;i} \frac{1=2}{j;j})$: Take $L = m^{1-(+1)}$; then $jmE(U_{1;i}U_{1;j}) \quad i;j = O(m^{-(+1)} \frac{1=2}{i;i} \frac{1=2}{j;j})$: And similarly $jmE(U_{1;i}U_{2;j})j = O(m^{-(+1)} \frac{1=2}{i;i} \frac{1=2}{j;j})$: Hence

$$i;j = m(E(U_{1;i}U_{1;j}) + E(U_{2;i}U_{2;j}) \quad E(U_{1;i}U_{2;j}) \quad E(U_{2;i}U_{1;j}))=2 = i;j + O(m^{-(+1)} \frac{1=2}{i;i} \frac{1=2}{j;j}):$$

Together with (80) we obtain the first part in (79).

Since $\hat{v}_{ij;k} = m(\Delta f_{k;i} + \Delta U_{k;i})(\Delta f_{k;j} + \Delta U_{k;j})$; we have $E \hat{v}_{ij;k} = m\Delta f_{k;i}\Delta f_{k;j} + 2\Delta U_{k;i}\Delta U_{k;j}$.
By (72) and (80),

$$v_{ij}^2 = \sum_{k \in S} E \hat{v}_{ij;k}^2 = N_2 \quad \tilde{v}_{ij}^2 = \sum_{k \in S} \text{Var}(\hat{v}_{ij;k}) = N_2 + O(m^3 n^{-2} \frac{1}{i} \frac{1}{j}):$$

Note by (73) and (75) we have

$$m^2 \text{Var}(\Delta f_{k;i} \Delta U_{k;j}) = O(m^3 n^{-2} \frac{1}{j}); \quad \text{and} \quad m^2 \text{Var}(\Delta U_{k;i} \Delta U_{k;j}) = O(\frac{1}{i} \frac{1}{j}):$$

Thus $\text{Var}(\hat{v}_{ij;k}) = O(\frac{1}{i} \frac{1}{j})$ and the second part in (79) holds.

Step 4. Since $jS_j \leq 2K_0$; for any i, j ; and $j \geq \log(2)$;

$$jN_1 h_{ij}(u) = N_2 \quad \bar{h}_{ij}(u) \leq 2\log(2)K_0 = (\frac{1}{i} N_2): \quad (81)$$

Combining (81), Step 1 and step 2 with $x = N_2^{1-2} \log^{1-2}(np)$, then with probability tending 1, for all $1 \leq i, j \leq p$; and $ju \leq C_0$;

$$B_{ij}(u; \Delta) = N_1 N_2^{-1} h_{ij}(u) - B_{ij}^+(u; \Delta); \quad (82)$$

where

$$\Delta = h \frac{1}{i} \frac{1}{j} + 2\log(2)K_0 = (\frac{1}{i} N_2) \quad \text{and} \quad h = x N_2; \quad (83)$$

Note if

$$\frac{2}{ij} v_{ij}^2 + 2 \frac{1}{ij} \Delta \leq 1; \quad (84)$$

then $B_{ij}^+(u; \Delta)$ exists real roots. Denote the smaller one as u^+ , which satisfies $u^+ \leq \frac{1}{ij} + \frac{1}{ij} v_{ij}^2 + 2\Delta$: Take $\tilde{v}_{ij} = \frac{1}{ij} \frac{1}{i} \frac{1}{j}$: By Step 3 and Assumption 2.3, if (84), then

$$\frac{1}{i} \frac{1}{j} \frac{1}{i} \frac{1}{j} (u^+ \leq \tilde{v}_{ij}) = O\left\{ m^{-(+1)} + m^3 n^2 + \frac{1}{ij} + h + mK_0 = (\frac{1}{ij} n) \right\}; \quad (85)$$

Similar bound can be obtained for u as well. When (82) holds, $u = \hat{u}_{ij} + u^+$. Take $\hat{u}_{ij} = (m=n)^{1=2}$; then with probability tending to 1,

$$\hat{u}_{ij} = (m=n)^{1=2} + N_2^{1=2} \log(np);$$

where the convergence is uniform for all $1 \leq i, j \leq p$: Since there exists some constant $c_1, c_2 > 0$; such that $c_1 \leq \hat{u}_{ij} \leq c_2$; and any i, j with probability tending to 1. Thus the desired result follows. \square

Proof of (ii). Same argument as for (i), except that we need to replace Step 2 by Step 2' with $\chi = N_2 = \log(np)^{2.5}$. Then we obtain the desired result.

Step 2'. This step shows

$$\sum_{i,j=1}^p \mathbb{P} \left(\sup_{|u| \leq c_0} |\bar{h}_{ij}(u) - \mathbb{E} \bar{h}_{ij}(u)| \geq \chi (m=n)^{1=2} + N_2^{1=2} \right) \leq p^2 n e^{-c \chi N_2^{1=2}}; \quad (86)$$

where c and the constant in (86) are independent of n, p, i, j :

The proof follows similar argument as in Step 2 and Theorem 3 in Wu and Wu (2016). \square

Figure 5: Visualization of one sample of simulated data with jump in case a),2),ii). x_1 represent t and x_2 represents i .

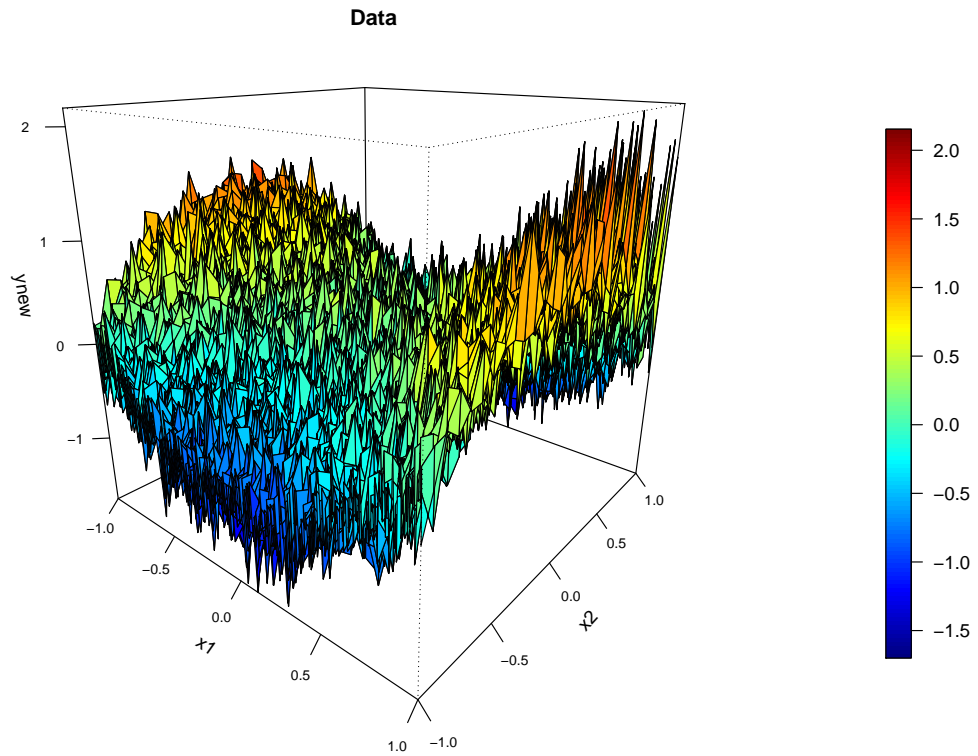


Figure 6: Visualization of the real (left) and estimated correlation matrix (using the robust estimation method).

