

Multiscale Inference for Nonparametric Time Trends

Marina Khismatullina¹

University of Bonn

Michael Vogt²

University of Bonn

April 30, 2018

We develop multiscale methods to test qualitative hypotheses about nonparametric time trends. In many applications, practitioners are interested in whether the observed time series has a time trend at all, that is, whether the trend function is non-constant. Moreover, they would like to get further information about the shape of the trend function. Among other things, they would like to know in which time regions there is an upward/downward movement in the trend. When multiple time series are observed, another important question is whether the observed time series all have the same time trend. We design multiscale tests to formally approach these questions. We derive asymptotic theory for the proposed tests and investigate their finite sample performance by means of simulations. In addition, we illustrate the methods by two applications to temperature data.

Key words: Multiscale statistics; nonparametric regression; time series errors; shape constraints; strong approximations; anti-concentration bounds.

AMS 2010 subject classifications: 62E20; 62G10; 62G20; 62M10.

1 Introduction

The analysis of time trends is an important aspect of many time series applications. In this paper, we develop new methods to analyse nonparametric time trends. We consider two different model settings, depending on whether a single or multiple time series are observed. When the observations come from a single time series $\{Y_t : 1 \leq t \leq T\}$, we consider the model

$$Y_t = m\left(\frac{t}{T}\right) + \varepsilon_t \quad (1.1)$$

for $1 \leq t \leq T$, where $m : [0, 1] \rightarrow \mathbb{R}$ is an unknown nonparametric trend function and the error terms ε_t form a time series process with $\mathbb{E}[\varepsilon_t] = 0$ for all t . When several time series $\mathcal{Y}_i = \{Y_{it} : 1 \leq t \leq T\}$ are observed for $1 \leq i \leq n$, we similarly model each time series \mathcal{Y}_i by the equation

$$Y_{it} = m_i\left(\frac{t}{T}\right) + \alpha_i + \varepsilon_{it} \quad (1.2)$$

for $1 \leq t \leq T$, where m_i is a nonparametric time trend, α_i is a (random or deterministic) intercept and ε_{it} are time series errors with $\mathbb{E}[\varepsilon_{it}] = 0$ for all t . As usual in nonparametric

¹Address: Bonn Graduate School of Economics, University of Bonn, 53113 Bonn, Germany. Email: marina.k@uni-bonn.de.

²Corresponding author. Address: Department of Economics and Hausdorff Center for Mathematics, University of Bonn, 53113 Bonn, Germany. Email: michael.vogt@uni-bonn.de.



Figure 1: Yearly mean temperature in Central England for the time period from 1659 to 2017 measured in $^{\circ}\text{C}$.

regression, we let the trend functions in (1.1) and (1.2) depend on rescaled time t/T rather than on real time t . A detailed description of models (1.1) and (1.2) is provided in Section 2.

Let us first have a closer look at the situation where a single time series is observed. In this case, practitioners are interested in questions such as the following: Does the observed time series have a trend at all? If so, which are the time regions where there is a strong trend? Is the trend decreasing or increasing in these regions? As an example, consider the time series plotted in Figure 1 which shows the yearly mean temperature in Central England from 1659 to 2017. Climatologists are very much interested in analysing the trending behaviour of temperature time series like this; see e.g. Benner (1999) and Rahmstorf et al. (2017). Among other things, they would like to know whether there is an upward trend in the Central England mean temperature towards the end of the sample as visual inspection might suggest. In Section 4, we develop a statistical procedure to approach questions like this. Specifically, we construct a method to test the null hypothesis that there is no time trend in the data. Importantly, the proposed method does not only allow to test whether the null hypothesis of no trend is violated. It also allows to identify, with a pre-specified statistical confidence, time regions where there is some upward or downward movement in the trend. As regards the temperature time series in Figure 1, we can for example claim, with a statistical confidence of approximately 95%, that there is some upward movement of the trend in the time period between ?? and ??. This is one of the results obtained from a detailed analysis of the time series as conducted in Section 8.

Let us now turn to the situation where multiple time series of the form (1.2) are observed. An important question in many applications is whether the time trends m_i of the various time series are all the same. When some of the time trends are different, there may still be groups of time series with the same time trend. In this case, it is often of interest to estimate the unknown groups from the data. In addition, when two time trends m_i and m_j are not the same, it may also be relevant to know in which time regions they differ from each other. In Section 5, we construct statistical methods to approach these questions. In particular, we develop a test of the hypothesis that all

time trends in model (1.2) are the same, that is, $m_1 = m_2 = \dots = m_n$. Similar as above, our method does not only allow to test whether the null hypothesis is violated. It also allows to detect, with a given statistical confidence, which time trends are different and in which time regions they differ from each other. Based on our test method, we further construct an algorithm which clusters the observed time series into groups with the same time trend.

We develop our methods and the underlying theory step by step in Sections 3–5. In Section 3, we introduce our methods in the context of a simple baseline case. We in particular discuss the problem of testing the simple hypothesis $H_0 : m = 0$ in model (1.1). In Sections 4 and 5, we adapt the methods and the theory to the test problems we are actually interested in. To construct our methods, we build on ideas from statistical multiscale testing as developed in Chaudhuri and Marron (1999, 2000), Hall and Heckman (2000) and Dümbgen and Spokoiny (2001) among others. Considering the hypothesis $H_0 : m = 0$ in model (1.1), our test procedure can be outlined roughly as follows: In a first step, we set up statistics $\widehat{s}_T(u, h)$ to test the hypothesis $H_0(u, h)$ that $m = 0$ on the interval $[u - h, u + h]$. In a second step, we aggregate these statistics $\widehat{s}_T(u, h)$ for a wide range of different intervals $[u - h, u + h]$. We thus construct a multiscale statistic which allows to test the hypothesis $H_0(u, h)$ simultaneously for many intervals $[u - h, u + h]$. A simple approach to aggregate the statistics $\widehat{s}_T(u, h)$ is to take their supremum $\sup_{(u, h) \in \mathcal{G}_T} |\widehat{s}_T(u, h)|$, where \mathcal{G}_T denotes the set of points (u, h) which are taken into account. As shown in the seminal work of Dümbgen and Spokoiny (2001), this approach is suboptimal in some sense. Following their lead, we define our multiscale statistic by $\widehat{\Psi}_T = \sup_{(u, h) \in \mathcal{G}_T} \{|\widehat{s}_T(u, h)| - \lambda(h)\}$, where $\lambda(h)$ are additive correction terms. The idea behind this additively corrected supremum statistic is discussed in detail in Section 3.

In recent years, multiscale approaches have been developed for a variety of test problems. The general idea of all these approaches is to simultaneously consider a family of test statistics for a wide range of locations u and scales or bandwidths h . This idea has been put to work in different ways, thus resulting in different multiscale test approaches. In the regression context, Chaudhuri and Marron (1999, 2000) have developed the so-called SiZer method. Hall and Heckman (2000) have constructed a multiscale test on monotonicity of a regression function. As already discussed above, Dümbgen and Spokoiny (2001) have developed a multiscale approach which works with additively corrected supremum statistics. They derive theoretical results in the context of a continuous Gaussian white noise model. Their theory can be extended in a fairly straightforward way to a nonparametric regression model $Y_t = m(t/T) + \varepsilon_t$ with i.i.d. Subgaussian errors ε_t . However, it is far from trivial to extend their theory to our setting, that is, to a nonparametric regression model $Y_t = m(t/T) + \varepsilon_t$ with a general weakly dependent error process $\{\varepsilon_t\}$. To derive our theoretical results, we come up with a proof strategy which is quite different from that of Dümbgen and Spokoiny (2001).

This strategy is of interest in itself and may be applied to other multiscale test problems. It combines strong approximation results for dependent processes as developed in Berkes et al. (2014) with anti-concentration bounds for Gaussian vectors as derived in Chernozhukov et al. (2015). The details are described in Section 3.

Our multiscale tests complement already existing tools for analysing nonparametric time trends. The test method developed in Section 4 provides an alternative to dependent SiZer methods as introduced in Park et al. (2004) and Rondonotti et al. (2007). Whereas the focus of these papers is mainly methodological, we back up our multiscale test by a complete asymptotic theory which characterizes its size and power properties. The multiscale test from Section 5 is a valuable alternative to other procedures for testing equality of time trends. Park et al. (2009) extend the SiZer approach to the problem of comparing multiple time trends and provide some theory for the special case that the number of observed time series is $n = 2$. Besides, there are several non-multiscale approaches to test for equality of time trends. Vogelsang and Franses (2005) and Lyubchich and Gel (2016), for example, construct tests for comparing parametric time trends, whereas Degras et al. (2012) and Chen and Wu (2018) develop L^2 -type tests for comparing nonparametric trends. One of the advantages of our multiscale approach is that it is very informative: It does not only allow to test whether the null hypothesis is violated or not. It also gives information on where violations occur, in particular, which time trends are different and in which time regions they differ. It thus provides valuable additional information for practitioners.

We complement the theoretical analysis of the paper by a simulation study and two application examples in Sections 7 and 8. The simulation study investigates the finite sample properties of the test methods and the clustering algorithm from Sections 4 and 5. In the first application example, we examine the temperature time series from Figure 1 with the help of the methods developed in Section 4. In the second example, we analyse temperature time series measured at 34 different weather stations located in Great Britain. We in particular apply our procedure from Section 5 to test whether the different time series have the same trend.

2 The model

We now describe the two model settings in detail which were briefly outlined in the Introduction. The model for the test problems considered in Sections 3 and 4 is as follows: We observe a single time series $\{Y_t : 1 \leq t \leq T\}$ of length T which satisfies the model equation

$$Y_t = m\left(\frac{t}{T}\right) + \varepsilon_t \quad (2.1)$$

for $1 \leq t \leq T$. Here, m is an unknown nonparametric trend function defined on $[0, 1]$ and $\{\varepsilon_t : 1 \leq t \leq T\}$ is a zero-mean stationary error process. For simplicity, we restrict

attention to equidistant design points $x_t = t/T$. However, our methods and theory can also be carried over to non-equidistant designs. The stationary error process $\{\varepsilon_t\}$ is assumed to have the following properties:

(C1) The variables ε_t allow for the representation $\varepsilon_t = G(\dots, \eta_{t-1}, \eta_t, \eta_{t+1}, \dots)$, where η_t are i.i.d. random variables and $G : \mathbb{R}^{\mathbb{Z}} \rightarrow \mathbb{R}$ is a measurable function.

(C2) It holds that $\|\varepsilon_t\|_q < \infty$ for some $q > 4$, where $\|\varepsilon_t\|_q = (\mathbb{E}|\varepsilon_t|^q)^{1/q}$.

Following Wu (2005), we impose conditions on the dependence structure of the error process $\{\varepsilon_t\}$ in terms of the physical dependence measure $d_{t,q} = \|\varepsilon_t - \varepsilon'_t\|_q$, where $\varepsilon'_t = G(\dots, \eta_{-1}, \eta'_0, \eta_1, \dots, \eta_{t-1}, \eta_t, \eta_{t+1}, \dots)$ with $\{\eta'_t\}$ being an i.i.d. copy of $\{\eta_t\}$. In particular, we assume the following:

(C3) Define $\Theta_{t,q} = \sum_{|s| \geq t} d_{s,q}$ for $t \geq 0$. It holds that

$$\Theta_{t,q} = O(t^{-\tau_q}(\log t)^{-A}),$$

where $A > \frac{2}{3}(1/q + 1 + \tau_q)$ and $\tau_q = \{q^2 - 4 + (q - 2)\sqrt{q^2 + 20q + 4}\}/8q$.

The conditions (C1)–(C3) are fulfilled by a wide range of stationary processes $\{\varepsilon_t\}$. As a first example, consider linear processes of the form $\varepsilon_t = \sum_{i=0}^{\infty} c_i \eta_{t-i}$ with $\|\varepsilon_t\|_q < \infty$, where c_i are absolutely summable coefficients and η_t are i.i.d. innovations with $\mathbb{E}[\eta_t] = 0$ and $\|\eta_t\|_q < \infty$. Trivially, (C1) and (C2) are fulfilled in this case. Moreover, if $|c_i| = O(\rho^i)$ for some $\rho \in (0, 1)$, then (C3) is easily seen to be satisfied as well. As a special case, consider an ARMA process $\{\varepsilon_t\}$ of the form $\varepsilon_t + \sum_{i=1}^p a_i \varepsilon_{t-i} = \eta_t + \sum_{j=1}^r b_j \eta_{t-j}$ with $\|\varepsilon_t\|_q < \infty$, where a_1, \dots, a_p and b_1, \dots, b_r are real-valued parameters. As before, we let η_t be i.i.d. innovations with $\mathbb{E}[\eta_t] = 0$ and $\|\eta_t\|_q < \infty$. Moreover, as usual, we suppose that the complex polynomials $A(z) = 1 + \sum_{j=1}^p a_j z^j$ and $B(z) = 1 + \sum_{j=1}^r b_j z^j$ do not have any roots in common. If $A(z)$ does not have any roots inside the unit disc, then the ARMA process $\{\varepsilon_t\}$ is stationary and causal. Specifically, it has the representation $\varepsilon_t = \sum_{i=0}^{\infty} c_i \eta_{t-i}$ with $|c_i| = O(\rho^i)$ for some $\rho \in (0, 1)$, implying that (C1)–(C3) are fulfilled. The results in Wu and Shao (2004) show that condition (C3) (as well as the other two conditions) is not only fulfilled for linear time series processes but also for a variety of non-linear processes.

The model setting for the test problem analysed in Section 5 is closely related to the setting discussed above. The main difference is that we observe multiple rather than only one time series. In particular, we observe time series $\mathcal{Y}_i = \{Y_{it} : 1 \leq t \leq T\}$ of length T for $1 \leq i \leq n$. Each time series \mathcal{Y}_i satisfies the model equation

$$Y_{it} = m_i\left(\frac{t}{T}\right) + \alpha_i + \varepsilon_{it} \quad (2.2)$$

for $1 \leq t \leq T$, where m_i is an unknown nonparametric trend function defined on $[0, 1]$, α_i is a (deterministic or random) intercept term and $\mathcal{E}_i = \{\varepsilon_{it} : 1 \leq t \leq T\}$ is a

zero-mean stationary error process. For identification, we normalize the functions m_i such that $\int_0^1 m_i(u)du = 0$ for all $1 \leq i \leq n$. The term α_i can also be regarded as an additional error component. In the econometrics literature, it is commonly called a fixed effect error term. It can be interpreted as capturing unobserved characteristics of the time series \mathcal{Y}_i which remain constant over time. We allow the error terms α_i to be dependent across i in an arbitrary way. Hence, by including them in model equation (2.2), we allow the n time series \mathcal{Y}_i in our sample to be correlated with each other. Whereas the terms α_i may be correlated, the error processes \mathcal{E}_i are assumed to be independent across i . In addition, each process \mathcal{E}_i is supposed to satisfy the conditions (C1)–(C3). Finally note that throughout the paper, we restrict attention to the case where the number of time series n in model (2.2) is bounded. It is however possible to extend our theoretical results to the case where n slowly grows with the sample size T .

3 The multiscale method

In this section, we introduce our multiscale test method and the underlying theory for the simple hypothesis $H_0 : m = 0$ in model (2.1). As we will see in Sections 4 and 5, both the method and the theory for this simple case can be easily adapted to more interesting test problems, in particular to the test problems discussed in the Introduction.

3.1 Construction of the test statistic

To construct a multiscale test statistic for the hypothesis $H_0 : m = 0$ in model (2.1), we consider the kernel averages

$$\widehat{\psi}_T(u, h) = \sum_{t=1}^T w_{t,T}(u, h) Y_t,$$

where $w_{t,T}(u, h)$ is a kernel weight with $u \in [0, 1]$ and the bandwidth parameter h . In order to avoid boundary issues, we work with a local linear weighting scheme. We in particular set

$$w_{t,T}(u, h) = \frac{\Lambda_{t,T}(u, h)}{\{\sum_{t=1}^T \Lambda_{t,T}^2(u, h)\}^{1/2}}, \quad (3.1)$$

where

$$\Lambda_{t,T}(u, h) = K\left(\frac{\frac{t}{T} - u}{h}\right) \left[S_{T,2}(u, h) - S_{T,1}(u, h) \left(\frac{\frac{t}{T} - u}{h}\right) \right],$$

$S_{T,\ell}(u, h) = (Th)^{-1} \sum_{t=1}^T K\left(\frac{\frac{t}{T} - u}{h}\right) \left(\frac{\frac{t}{T} - u}{h}\right)^\ell$ for $\ell = 0, 1, 2$ and K is a kernel function with the following properties:

(C4) The kernel K is non-negative, symmetric about zero and integrates to one. Moreover, it has compact support $[-1, 1]$ and is Lipschitz continuous, that is, $|K(v) - K(w)| \leq C|v - w|$ for any $v, w \in \mathbb{R}$ and some constant $C > 0$.

Alternatively to the local linear weights defined in (3.1), we could also work with local constant weights which are defined analogously with $\Lambda_{t,T}(u, h) = K(\frac{t}{h} - \frac{u}{h})$. We however prefer to use local linear weights as these have superior theoretical properties at the boundary.

The kernel average $\hat{\psi}_T(u, h)$ is a local average of the observations Y_1, \dots, Y_T which gives positive weight only to data points Y_t with $t/T \in [u - h, u + h]$. Hence, only observations Y_t with t/T close to the location u are taken into account, the amount of localization being determined by the bandwidth h . With the weights defined in (3.1), the kernel average $\hat{\psi}_T(u, h)$ is nothing else than a rescaled local linear estimator of $m(u)$ with bandwidth h . The weights are chosen such that in the case of independent error terms ε_t , $\text{Var}(\hat{\psi}_T(u, h)) = \sigma^2$ for any location u and bandwidth h , where $\sigma^2 = \text{Var}(\varepsilon_t)$. In the more general case that the error terms satisfy the weak dependence conditions from Section 2, it holds that $\text{Var}(\hat{\psi}_T(u, h)) = \sigma^2 + o(1)$ for any location u and any bandwidth h with $h \rightarrow 0$ and $Th \rightarrow \infty$, where $\sigma^2 = \sum_{\ell=-\infty}^{\infty} \text{Cov}(\varepsilon_0, \varepsilon_\ell)$ is the long-run variance of the error terms. Hence, the statistics $\hat{\psi}_T(u, h)$ have approximately the same variance across u and h for sufficiently large sample sizes T . In what follows, we consider normalized versions $\hat{\psi}_T(u, h)/\hat{\sigma}$ of the kernel averages $\hat{\psi}_T(u, h)$, where $\hat{\sigma}^2$ is an estimator of the long-run error variance σ^2 . The problem of estimating σ^2 is discussed in detail in Section 6. There, we construct estimators $\hat{\sigma}^2$ with the property that $\hat{\sigma}^2 = \sigma^2 + O_p(1/\sqrt{T})$ under appropriate conditions. For the time being, we suppose that $\hat{\sigma}^2$ is an estimator with reasonable theoretical properties. We in particular assume that $\hat{\sigma}^2 = \sigma^2 + o_p(\rho_T)$ with $\rho_T = o(1/\log T)$. The convergence rate ρ_T is thus allowed to be much slower than $1/\sqrt{T}$.

Our multiscale statistic combines the kernel averages $\hat{\psi}_T(u, h)$ for a wide range of different locations u and bandwidths or scales h . Specifically, it is defined as

$$\hat{\Psi}_T = \max_{(u, h) \in \mathcal{G}_T} \left\{ \left| \frac{\hat{\psi}_T(u, h)}{\hat{\sigma}} \right| - \lambda(h) \right\}, \quad (3.2)$$

where $\lambda(h) = \sqrt{2 \log\{1/(2h)\}}$ and \mathcal{G}_T is the set of points (u, h) that are taken into consideration. The details on the set \mathcal{G}_T are discussed below. As can be seen, the statistic $\hat{\Psi}_T$ does not simply aggregate the individual statistics $\hat{\psi}_T(u, h)/\hat{\sigma}$ by taking the supremum over all points $(u, h) \in \mathcal{G}_T$ as in more traditional multiscale approaches. We rather follow the approach pioneered by Dümbgen and Spokoiny (2001) and subtract the additive correction term $\lambda(h)$ from the statistics $\hat{\psi}_T(u, h)/\hat{\sigma}$ that correspond to the bandwidth level h . To see the heuristic idea behind the additive correction $\lambda(h)$,

consider for a moment the uncorrected statistic

$$\widehat{\Psi}_{T,\text{uncorrected}} = \max_{(u,h) \in \mathcal{G}_T} \left| \frac{\widehat{\psi}_T(u,h)}{\widehat{\sigma}} \right|$$

and suppose that the null hypothesis $H_0 : m = 0$ holds true. For simplicity, assume that the errors ε_t are i.i.d. normally distributed and neglect the estimation error in $\widehat{\sigma}$, that is, set $\widehat{\sigma} = \sigma$. Moreover, suppose that the set \mathcal{G}_T only consists of the points $(u_k, h_\ell) = ((2k-1)h_\ell, h_\ell)$ with $k = 1, \dots, \lfloor 1/2h_\ell \rfloor$ and $\ell = 1, \dots, L$. In this case, we can write

$$\widehat{\Psi}_{T,\text{uncorrected}} = \max_{1 \leq \ell \leq L} \max_{1 \leq k \leq \lfloor 1/2h_\ell \rfloor} \left| \frac{\widehat{\psi}_T(u_k, h_\ell)}{\sigma} \right|.$$

Under our simplifying assumptions, the statistics $\widehat{\psi}_T(u_k, h_\ell)/\sigma$ with $k = 1, \dots, \lfloor 1/2h_\ell \rfloor$ are independent and standard normal for any given bandwidth h_ℓ . Since the maximum over $\lfloor 1/2h \rfloor$ independent standard normal random variables is $\lambda(h) + o_p(1)$ as $h \rightarrow 0$, we obtain that $\max_k \widehat{\psi}_T(u_k, h_\ell)/\sigma$ is approximately of size $\lambda(h_\ell)$ for small bandwidths h_ℓ . As $\lambda(h) \rightarrow \infty$ for $h \rightarrow 0$, this implies that $\max_k \widehat{\psi}_T(u_k, h_\ell)/\sigma$ tends to be much larger in size for small than for large bandwidths h_ℓ . As a result, the stochastic behaviour of the uncorrected statistic $\widehat{\Psi}_{T,\text{uncorrected}}$ tends to be dominated by the statistics $\widehat{\psi}_T(u_k, h_\ell)$ corresponding to small bandwidths h_ℓ . The additively corrected statistic $\widehat{\Psi}_T$, in contrast, puts the statistics $\widehat{\psi}_T(u_k, h_\ell)$ corresponding to different bandwidths h_ℓ on a more equal footing, thus counteracting the dominance of small bandwidth values.

The multiscale statistic $\widehat{\Psi}_T$ simultaneously takes into account all locations u and bandwidths h with $(u, h) \in \mathcal{G}_T$. Throughout the paper, we suppose that \mathcal{G}_T is some subset of $\mathcal{G}_T^{\text{full}} = \{(u, h) : u = t/T \text{ for some } 1 \leq t \leq T \text{ and } h \in [h_{\min}, h_{\max}]\}$, where h_{\min} and h_{\max} denote some minimal and maximal bandwidth value, respectively. For our theory to work, we require the following conditions to hold:

(C5) $|\mathcal{G}_T| = O(T^\theta)$ for some arbitrarily large but fixed constant $\theta > 0$, where $|\mathcal{G}_T|$ denotes the cardinality of \mathcal{G}_T .

(C6) $h_{\min} \gg T^{-(1-\frac{2}{q})} \log T$, that is, $h_{\min}/\{T^{-(1-\frac{2}{q})} \log T\} \rightarrow \infty$ with $q > 4$ defined in (C2) and $h_{\max} = o(1)$.

According to (C5), the number of points (u, h) in \mathcal{G}_T should not grow faster than T^θ for some arbitrarily large but fixed $\theta > 0$. This is a fairly weak restriction as it allows the set \mathcal{G}_T to be extremely large as compared to the sample size T . For example, we may work with the set

$$\begin{aligned} \mathcal{G}_T &= \{(u, h) : u = t/T \text{ for some } 1 \leq t \leq T \text{ and } h \in [h_{\min}, h_{\max}]\} \\ &\quad \text{with } h = t/T \text{ for some } 1 \leq t \leq T, \end{aligned}$$

which contains more than enough points (u, h) for most practical applications. Condition (C6) imposes some restrictions on the minimal and maximal bandwidths h_{\min} and h_{\max} . These conditions are fairly weak, allowing us to choose the bandwidth window $[h_{\min}, h_{\max}]$ extremely large. In particular, we can choose the minimal bandwidth h_{\min} to be of the order $T^{-1/2}$ for any $q > 4$, which means that we can let h_{\min} converge to 0 very quickly. Moreover, the maximal bandwidth h_{\max} is allowed to converge to 0 arbitrarily slowly, which implies that we can pick it very large.

3.2 The test procedure

In order to formulate a test for the hypothesis $H_0 : m = 0$, we still need to specify a critical value. To do so, we define the statistic

$$\Phi_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\phi_T(u, h)}{\sigma} \right| - \lambda(h) \right\}, \quad (3.3)$$

where $\phi_T(u, h) = \sum_{t=1}^T w_{t,T}(u, h) \sigma Z_t$ and Z_t are independent standard normal random variables. The statistic Φ_T can be regarded as a Gaussian version of the test statistic $\hat{\Psi}_T$ under the null hypothesis H_0 . Let $q_T(\alpha)$ be the $(1 - \alpha)$ -quantile of Φ_T . Importantly, the quantile $q_T(\alpha)$ can be computed by Monte Carlo simulations and can thus be regarded as known. Our multiscale test of the hypothesis $H_0 : m = 0$ is now defined as follows: For a given significance level $\alpha \in (0, 1)$, we reject H_0 if $\hat{\Psi}_T > q_T(\alpha)$.

3.3 Theoretical properties of the test

In order to examine the theoretical properties of our multiscale test, we introduce the statistic

$$\hat{\Phi}_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\hat{\phi}_T(u, h)}{\hat{\sigma}} \right| - \lambda(h) \right\} \quad (3.4)$$

with $\hat{\phi}_T(u, h) = \hat{\psi}_T(u, h) - \mathbb{E}[\hat{\psi}_T(u, h)] = \sum_{t=1}^T w_{t,T}(u, h) \varepsilon_t$. According to the following theorem, the (known) quantile $q_T(\alpha)$ of Φ_T defined in Section 3.2 can be used as a proxy for the $(1 - \alpha)$ -quantile of the statistic $\hat{\Phi}_T$.

Theorem 3.1. *Let (C1)–(C6) be fulfilled and assume that $\hat{\sigma}^2 = \sigma^2 + o_p(\rho_T)$ with $\rho_T = o(1/\log T)$. Then*

$$\mathbb{P}(\hat{\Phi}_T \leq q_T(\alpha)) = (1 - \alpha) + o(1).$$

A full proof of Theorem 3.1 is given in the Appendix. We here shortly outline the proof strategy, which splits up into two main steps. In the first, we replace the statistic $\hat{\Phi}_T$ for each $T \geq 1$ by a statistic $\tilde{\Phi}_T$ with the same distribution as $\hat{\Phi}_T$ and the property

that

$$|\tilde{\Phi}_T - \Phi_T| = o_p(\delta_T), \quad (3.5)$$

where $\delta_T = o(1)$ and the Gaussian statistic Φ_T is defined in Section 3.2. We thus replace the statistic $\hat{\Phi}_T$ by an identically distributed version which is close to a Gaussian statistic whose distribution is known. To do so, we make use of strong approximation theory for dependent processes as derived in Berkes et al. (2014). In the second step, we show that

$$\sup_{x \in \mathbb{R}} |\mathbb{P}(\tilde{\Phi}_T \leq x) - \mathbb{P}(\Phi_T \leq x)| = o(1), \quad (3.6)$$

which immediately implies the statement of Theorem 3.1. Importantly, the convergence result (3.5) is not sufficient for establishing (3.6). Put differently, the fact that $\tilde{\Phi}_T$ can be approximated by Φ_T in the sense that $\tilde{\Phi}_T - \Phi_T = o_p(\delta_T)$ does not imply that the distribution of $\tilde{\Phi}_T$ is close to that of Φ_T in the sense of (3.6). For (3.6) to hold, we additionally require the distribution of Φ_T to have some sort of continuity property. Specifically, we prove that

$$\sup_{x \in \mathbb{R}} \mathbb{P}(|\Phi_T - x| \leq \delta_T) = o(1), \quad (3.7)$$

which says that Φ_T does not concentrate too strongly in small regions of the form $[x - \delta_T, x + \delta_T]$. The main tool for verifying (3.7) are anti-concentration results for Gaussian random vectors as derived in Chernozhukov et al. (2015). The claim (3.6) can be proven by combining (3.5) and (3.7), which in turn yields Theorem 3.1.

With the help of Theorem 3.1, we can investigate the theoretical properties of our multiscale test. The first result is an immediate consequence of Theorem 3.1. It says that the test has the correct (asymptotic) size.

Proposition 3.2. *Let the conditions of Theorem 3.1 be satisfied. Under the null hypothesis $H_0 : m = 0$, it holds that*

$$\mathbb{P}(\hat{\Psi}_T \leq q_T(\alpha)) = (1 - \alpha) + o(1).$$

The second result characterizes the power of the multiscale test against local alternatives. To formulate it, we consider any sequence of functions $m = m_T$ with the following property: There exists $(u, h) \in \mathcal{G}_T$ with $[u - h, u + h] \subseteq [0, 1]$ such that

$$m_T(w) \geq c_T \sqrt{\frac{\log T}{Th}} \quad \text{for all } w \in [u - h, u + h], \quad (3.8)$$

where $\{c_T\}$ is any sequence of positive numbers with $c_T \rightarrow \infty$. Alternatively to (3.8), we may also assume that $-m_T(w) \geq c_T \sqrt{\log T / (Th)}$ for all $w \in [u - h, u + h]$. According to the following result, our test has asymptotic power 1 against local alternatives of the form (3.8).

Proposition 3.3. *Let the conditions of Theorem 3.1 be satisfied and consider any sequence of functions m_T with the property (3.8). Then*

$$\mathbb{P}(\widehat{\Psi}_T \leq q_T(\alpha)) = o(1).$$

The proof of Proposition 3.3 can be found in the Appendix. To formulate the next result, we define

$$\Pi_T = \{I_{u,h} = [u - h, u + h] : (u, h) \in \mathcal{A}_T\}$$

with

$$\mathcal{A}_T = \left\{ (u, h) \in \mathcal{G}_T : \left| \frac{\widehat{\psi}_T(u, h)}{\widehat{\sigma}} \right| - \lambda(h) > q_T(\alpha) \right\}.$$

Π_T is the collection of intervals $I_{u,h} = [u - h, u + h]$ for which the (corrected) test statistic $|\widehat{\psi}_T(u, h)/\widehat{\sigma}| - \lambda(h)$ lies above the critical value $q_T(\alpha)$. With this notation at hand, we consider the event

$$E_T = \left\{ \forall I_{u,h} \in \Pi_T : m(v) \neq 0 \text{ for some } v \in I_{u,h} = [u - h, u + h] \right\}.$$

This is the event that the null hypothesis is violated on all intervals $I_{u,h}$ for which the (corrected) test statistic $|\widehat{\psi}_T(u, h)/\widehat{\sigma}| - \lambda(h)$ is above the critical value $q_T(\alpha)$. We can make the following formal statement about the event E_T whose proof is given in the Appendix.

Proposition 3.4. *Under the conditions of Theorem 3.1, it holds that*

$$\mathbb{P}(E_T) \geq (1 - \alpha) + o(1).$$

According to Proposition 3.4, our test procedure allows to make uniform confidence statements of the following form: With (asymptotic) probability $\geq (1 - \alpha)$, the null hypothesis $H_0 : m = 0$ is violated on all intervals $I_{u,h} \in \Pi_T$. Hence, our multiscale test does not only allow to check whether the null hypothesis is violated. It also allows to identify regions where violations occur with a pre-specified level of confidence.

The statement of Proposition 3.4 suggests to graphically present the results of our multiscale test by plotting the intervals $I_{u,h} \in \Pi_T$, that is, by plotting the intervals where (with asymptotic confidence $\geq 1 - \alpha$) our test detects a violation of the null hypothesis. The drawback of this graphical presentation is that the number of intervals in Π_T is often quite large. To obtain a better graphical summary of the results, we replace Π_T by a subset Π_T^{\min} which is constructed as follows: As in Dümbgen (2002), we call an interval $I_{u,h} \in \Pi_T$ minimal if there is no other interval $I_{u',h'} \in \Pi_T$ with $I_{u',h'} \subset I_{u,h}$. Let Π_T^{\min} be the set of all minimal intervals in Π_T and define the event

$$E_T^{\min} = \left\{ \forall I_{u,h} \in \Pi_T^{\min} : m(v) \neq 0 \text{ for some } v \in I_{u,h} = [u - h, u + h] \right\}.$$

It is easily seen that $E_T = E_T^{\min}$. Hence, by Proposition 3.4, it holds that

$$\mathbb{P}(E_T^{\min}) \geq (1 - \alpha) + o(1).$$

This suggests to plot the minimal intervals in Π_T^{\min} rather than the whole collection of intervals Π_T as a graphical summary of the test results. We in particular use this way of presenting the test results in our application examples of Section 8.

4 Testing for the presence of a time trend

In what follows, we construct a multiscale test for the null hypothesis that the trend function m in model (2.1) is constant. To achieve this, we adapt the methodology developed in Section 3. Importantly, the resulting multiscale procedure does not only allow to test whether the null hypothesis is violated. As we will see, it also allows to identify, with a pre-specified statistical confidence, time regions where violations occur. Put differently, it allows to identify, with a given confidence, intervals $I_{u,h} = [u-h, u+h]$ where m is not constant over time. It thus provides information on where the time trend is increasing/decreasing, which is important knowledge in many applications.

4.1 Construction of the test statistic

Throughout the section, we suppose that the trend m is continuously differentiable. The null hypothesis that m is constant can be formulated as $H_0 : m' = 0$, where m' denotes the first derivative of m . To construct a test statistic for the hypothesis H_0 , we proceed analogously as in Section 3.1. To start with, we introduce the kernel averages

$$\widehat{\psi}'_T(u, h) = \sum_{t=1}^T w'_{t,T}(u, h) Y_t,$$

where the kernel weights $w'_{t,T}(u, h)$ are given by

$$w'_{t,T}(u, h) = \frac{\Lambda'_{t,T}(u, h)}{\{\sum_{t=1}^T \Lambda'_{t,T}(u, h)^2\}^{1/2}} \quad (4.1)$$

with

$$\Lambda'_{t,T}(u, h) = K\left(\frac{\frac{t}{T} - u}{h}\right) \left[S_{T,0}(u, h) \left(\frac{\frac{t}{T} - u}{h}\right) - S_{T,1}(u, h) \right].$$

Here, $S_{T,\ell}(u, h)$ is defined as in Section 3.1 and K is a kernel function which satisfies (C4). The kernel average $\widehat{\psi}'_T(u, h)$ is a rescaled version of the local linear estimator of the derivative $m'(u)$ with bandwidth h . Alternatively to the local linear weights defined in (4.1), we could employ the weights $w'_{t,T}(u, h) = K'(\frac{u - \frac{t}{T}}{h}) / \{\sum_{t=1}^T K'(\frac{u - \frac{t}{T}}{h})^2\}^{1/2}$, where the kernel function K is assumed to be differentiable and K' is its derivative. To

avoid boundary problems, we however work with the local linear weights from (4.1) throughout the paper. Our multiscale statistic is defined as

$$\widehat{\Psi}'_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\widehat{\psi}'_T(u,h)}{\widehat{\sigma}} \right| - \lambda(h) \right\},$$

where $\lambda(h) = \sqrt{2 \log\{1/(2h)\}}$ and the set \mathcal{G}_T has been introduced in Section 3.1. As can be seen, the statistic $\widehat{\Psi}'_T$ is very similar to that from Section 3. Only the kernel averages $\widehat{\psi}'_T(u,h)$ have a somewhat different form.

4.2 The test procedure

As in Section 3.2, we define a Gaussian version Φ'_T of the test statistic $\widehat{\Psi}'_T$ under the null hypothesis H_0 by

$$\Phi'_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\phi'_T(u,h)}{\sigma} \right| - \lambda(h) \right\},$$

where $\phi'_T(u,h) = \sum_{t=1}^T w'_{t,T}(u,h) \sigma Z_t$ and Z_t are independent standard normal random variables. Denoting the $(1 - \alpha)$ -quantile of Φ'_T by $q'_T(\alpha)$, our multiscale test of the hypothesis $H_0: m' = 0$ is defined as follows: For a given significance level $\alpha \in (0, 1)$, we reject H_0 if $\widehat{\Psi}'_T > q'_T(\alpha)$.

4.3 The theoretical properties of the test

The theoretical analysis parallels that of Section 3.3. We first investigate the theoretical properties of the auxiliary statistic

$$\widehat{\Phi}'_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\widehat{\phi}'_T(u,h)}{\widehat{\sigma}} \right| - \lambda(h) \right\},$$

where $\widehat{\phi}'_T(u,h) = \sum_{t=1}^T w'_{t,T}(u,h) \varepsilon_t$. The following result adapts Theorem 3.1 to our current test problem.

Theorem 4.1. *Let (C1)–(C6) be fulfilled and assume that $\widehat{\sigma}^2 = \sigma^2 + o_p(\rho_T)$ with $\rho_T = o(1/\log T)$. Then*

$$\mathbb{P}(\widehat{\Phi}'_T \leq q'_T(\alpha)) = (1 - \alpha) + o(1).$$

The proof of Theorem 4.1 is essentially the same as that of Theorem 3.1 and thus omitted. With the help of Theorem 4.1, we can derive the following theoretical properties of our multiscale test.

Proposition 4.2. *Let the conditions of Theorem 4.1 be satisfied.*

(a) *Under the null hypothesis H_0 , it holds that*

$$\mathbb{P}(\widehat{\Psi}'_T \leq q'_T(\alpha)) = (1 - \alpha) + o(1).$$

(b) *Consider any sequence of functions $m = m_T$ with the following property: There exists $(u, h) \in \mathcal{G}_T$ with $[u - h, u + h] \subseteq [0, 1]$ such that $m'_T(w) \geq c_T \sqrt{\log T / (Th^3)}$ for all $w \in [u - h, u + h]$ or $-m'_T(w) \geq c_T \sqrt{\log T / (Th^3)}$ for all $w \in [u - h, u + h]$, where $\{c_T\}$ is any sequence of positive numbers with $c_T \rightarrow \infty$. Then*

$$\mathbb{P}(\widehat{\Psi}'_T \leq q'_T(\alpha)) = o(1).$$

Part (a) of Proposition 4.2 is a simple consequence of Theorem 4.1. Part (b) can be proven by similar arguments as Proposition 3.3. The details are given in the Supplementary Material. Taken together, the two parts of Proposition 4.2 show that our multiscale test has the correct (asymptotic) size and that it is able to detect certain local alternatives with probability tending to 1. We next consider the events

$$\begin{aligned} E_T^+ &= \left\{ \forall I_{u,h} \in \Pi_T^+ : m'(v) > 0 \text{ for some } v \in I_{u,h} = [u - h, u + h] \right\} \\ E_T^- &= \left\{ \forall I_{u,h} \in \Pi_T^- : m'(v) < 0 \text{ for some } v \in I_{u,h} = [u - h, u + h] \right\}, \end{aligned}$$

where the sets Π_T^+ and Π_T^- are given by

$$\begin{aligned} \Pi_T^+ &= \{I_{u,h} = [u - h, u + h] : (u, h) \in \mathcal{A}_T^+ \text{ and } I_{u,h} \subseteq [0, 1]\} \\ \Pi_T^- &= \{I_{u,h} = [u - h, u + h] : (u, h) \in \mathcal{A}_T^- \text{ and } I_{u,h} \subseteq [0, 1]\} \end{aligned}$$

with

$$\begin{aligned} \mathcal{A}_T^+ &= \left\{ (u, h) \in \mathcal{G}_T : \frac{\widehat{\psi}'_T(u, h)}{\widehat{\sigma}} > q'_T(\alpha) + \lambda(h) \right\} \\ \mathcal{A}_T^- &= \left\{ (u, h) \in \mathcal{G}_T : -\frac{\widehat{\psi}'_T(u, h)}{\widehat{\sigma}} > q'_T(\alpha) + \lambda(h) \right\}. \end{aligned}$$

E_T^+ is the event that for each interval $I_{u,h} \in \Pi_T^+$, there is a subset $J_{u,h} \subseteq I_{u,h}$ with m being an increasing function on $J_{u,h}$. An analogous description applies to the event E_T^- . The following result shows that the events E_T^+ and E_T^- occur with asymptotic probability $\geq 1 - \alpha$.

Proposition 4.3. *Under the conditions of Theorem 4.1, it holds that*

$$\begin{aligned}\mathbb{P}(E_T^+) &\geq (1 - \alpha) + o(1) \\ \mathbb{P}(E_T^-) &\geq (1 - \alpha) + o(1).\end{aligned}$$

The proof of Proposition 4.3 parallels that of Proposition 3.4. The details are provided in the Supplementary Material. The statement of Proposition 4.3 can be summarized as follows: With asymptotic probability $\geq 1 - \alpha$, there is a subset $J_{u,h} \subseteq I_{u,h}$ for each interval $I_{u,h} \in \Pi_T^+$ such that m is an increasing function on $J_{u,h}$. Put differently, with asymptotic probability $\geq 1 - \alpha$, the trend m is increasing on some part of the interval $I_{u,h}$ for any $I_{u,h} \in \Pi_T^+$. An analogous statement holds for the intervals in the set Π_T^- . Our multiscale procedure thus allows to identify, with a pre-specified confidence, time regions where there is an increase/decrease in the time trend m .

We close the section with some additional remarks on Proposition 4.3: (i) The statement of Proposition 4.3 remains to hold true when we replace the sets Π_T^+ and Π_T^- by the corresponding sets of minimal intervals. (ii) In the sets Π_T^+ and Π_T^- , we only take into account intervals $I_{u,h} = [u - h, u + h]$ which are subsets of $[0, 1]$. We thus exclude points $(u, h) \in \mathcal{A}_T^+$ and $(u, h) \in \mathcal{A}_T^-$ which lie at the boundary, that is, for which $I_{u,h} \not\subseteq [0, 1]$. The reason is as follows: Let $(u, h) \in \mathcal{A}_T^+$ with $I_{u,h} \not\subseteq [0, 1]$. Our technical arguments allow us to say, with asymptotic confidence $\geq 1 - \alpha$, that $m'(v) \neq 0$ for some $v \in I_{u,h}$. However, we cannot say whether $m'(v) > 0$ or $m'(v) < 0$, that is, we cannot make confidence statements about the sign. Roughly speaking, the problem is that the local linear weights $w'_{t,T}(u, h)$ behave quite differently at boundary points (u, h) with $I_{u,h} \not\subseteq [0, 1]$. If we are only interested in whether there is some movement in the trend on an interval $I_{u,h}$ but we do not care whether it is an upward or downward movement, we may also consider the event $E_T^\pm = \{\forall I_{u,h} \in \Pi_T^\pm : m'(v) \neq 0 \text{ for some } v \in I_{u,h}\}$, where the set $\Pi_T^\pm = \{I_{u,h} : (u, h) \in \mathcal{A}_T^+ \cup \mathcal{A}_T^-\}$ contains all intervals $I_{u,h}$ with $(u, h) \in \mathcal{A}_T^+ \cup \mathcal{A}_T^-$, in particular those with $I_{u,h} \not\subseteq [0, 1]$. With the help of the technical arguments for Proposition 4.3, it follows that $\mathbb{P}(E_T^\pm) \geq (1 - \alpha) + o(1)$.

5 Testing for equality of time trends

In this section, we adapt the multiscale method developed in Section 3 to test the hypothesis that the trend functions in model (2.2) are all the same. More formally, we test the null hypothesis $H_0 : m_1 = m_2 = \dots = m_n$ in model (2.2). As we will see, the proposed multiscale method does not only allow to test whether the null hypothesis is violated. It also provides information on where violations occur. More specifically, it allows to identify, with a pre-specified confidence, (i) trend functions which are different from each other and (ii) time intervals where these trend functions differ.

5.1 Construction of the test statistic

To start with, we introduce some notation. The i -th time series in model (2.2) satisfies the equation $Y_{it} = m_i(\frac{t}{T}) + \alpha_i + \varepsilon_{it}$, where ε_{it} are zero-mean error terms and α_i are (random or deterministic) intercepts. Defining $Y_{it}^\circ = Y_{it} - \alpha_i$, this equation can be rewritten as $Y_{it}^\circ = m_i(\frac{t}{T}) + \varepsilon_{it}$, which is a standard nonparametric regression equation. The variables Y_{it}° are not observed, but they can be approximated by $\hat{Y}_{it} = Y_{it} - \hat{\alpha}_i$, where $\hat{\alpha}_i = T^{-1} \sum_{t=1}^T Y_{it}$ is an estimator of the intercept α_i . By construction, $\hat{\alpha}_i - \alpha_i = T^{-1} \sum_{t=1}^T \varepsilon_{it} + T^{-1} \sum_{t=1}^T m_i(\frac{t}{T}) = O_p(T^{-1/2}) + T^{-1} \sum_{t=1}^T m_i(\frac{t}{T})$. Hence, $\hat{\alpha}_i$ is a reasonable estimator of α_i if $T^{-1} \sum_{t=1}^T m_i(\frac{t}{T})$ converges to zero as $T \rightarrow \infty$. To ensure this, we suppose throughout the section that the functions m_i are Lipschitz continuous, that is, $|m_i(v) - m_i(w)| \leq L|v - w|$ for all $v, w \in [0, 1]$ and some constant $L < \infty$. Since $\int_0^1 m_i(u) du = 0$ by normalization, this implies that $T^{-1} \sum_{t=1}^T m_i(\frac{t}{T}) = O(T^{-1})$. We further let $\hat{\sigma}_i^2$ be an estimator of the long-run error variance $\sigma_i^2 = \sum_{\ell=-\infty}^{\infty} \text{Cov}(\varepsilon_{i0}, \varepsilon_{i\ell})$ which is computed from the constructed sample $\{\hat{Y}_{it} : 1 \leq t \leq T\}$. We thus regard $\hat{\sigma}_i^2 = \hat{\sigma}_i^2(\hat{Y}_{i1}, \dots, \hat{Y}_{iT})$ as a function of the variables \hat{Y}_{it} for $1 \leq t \leq T$. Throughout the section, we assume that $\hat{\sigma}_i^2 = \sigma_i^2 + o_p(\rho_T)$ with $\rho_T = o(1/\log T)$. Details on how to construct estimators of σ_i^2 are deferred to Section 6.

We are now ready to introduce the multiscale statistic for testing the hypothesis $H_0 : m_1 = m_2 = \dots = m_n$. For any pair of time series i and j , we define the kernel averages

$$\hat{\psi}_{ij,T}(u, h) = \sum_{t=1}^T w_{t,T}(u, h)(\hat{Y}_{it} - \hat{Y}_{jt}),$$

where the kernel weights $w_{t,T}(u, h)$ are defined as in (3.1). The kernel average $\hat{\psi}_{ij,T}(u, h)$ can be regarded as measuring the distance between the two trend curves m_i and m_j on the interval $[u - h, u + h]$. Similar as in Section 3.1, we aggregate the kernel averages $\hat{\psi}_{ij,T}(u, h)$ for all $(u, h) \in \mathcal{G}_T$ by the multiscale statistic

$$\hat{\Psi}_{ij,T} = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\hat{\psi}_{ij,T}(u, h)}{(\hat{\sigma}_i^2 + \hat{\sigma}_j^2)^{1/2}} \right| - \lambda(h) \right\},$$

where $\lambda(h) = \sqrt{2 \log\{1/(2h)\}}$ and the set \mathcal{G}_T has been introduced in Section 3.1. The statistic $\hat{\Psi}_{ij,T}$ can be interpreted as a distance measure between the two curves m_i and m_j . We finally define the multiscale statistic for testing the null hypothesis $H_0 : m_1 = m_2 = \dots = m_n$ as

$$\hat{\Psi}_{n,T} = \max_{1 \leq i < j \leq n} \hat{\Psi}_{ij,T},$$

that is, we define it as the maximal distance $\hat{\Psi}_{ij,T}$ between any pair of curves m_i and m_j with $i \neq j$.

5.2 The test procedure

Let Z_{it} for $1 \leq t \leq T$ and $1 \leq i \leq n$ be independent standard normal random variables which are independent of the error terms ε_{it} . Denote the empirical average of the variables Z_{i1}, \dots, Z_{iT} by $\bar{Z}_{i,T} = T^{-1} \sum_{t=1}^T Z_{it}$. To simplify notation, we write $\bar{Z}_i = \bar{Z}_{i,T}$ in what follows. For each i and j , we introduce the Gaussian statistic

$$\Phi_{ij,T} = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\phi_{ij,T}(u,h)}{(\hat{\sigma}_i^2 + \hat{\sigma}_j^2)^{1/2}} \right| - \lambda(h) \right\},$$

where $\phi_{ij,T}(u,h) = \sum_{t=1}^T w_{t,T}(u,h) \{ \hat{\sigma}_i(Z_{it} - \bar{Z}_i) - \hat{\sigma}_j(Z_{jt} - \bar{Z}_j) \}$. Moreover, we define the statistic

$$\Phi_{n,T} = \max_{1 \leq i < j \leq n} \Phi_{ij,T}$$

and denote its $(1 - \alpha)$ -quantile by $q_{n,T}(\alpha)$. Our multiscale test of the hypothesis $H_0 : m_1 = m_2 = \dots = m_n$ is defined as follows: For a given significance level $\alpha \in (0, 1)$, we reject H_0 if $\hat{\Psi}_{n,T} > q_{n,T}(\alpha)$.

5.3 The theoretical properties of the test

To start with, we introduce the auxiliary statistic

$$\hat{\Phi}_{n,T} = \max_{1 \leq i < j \leq n} \hat{\Phi}_{ij,T},$$

where

$$\hat{\Phi}_{ij,T} = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\hat{\phi}_{ij,T}(u,h)}{\{\hat{\sigma}_i^2 + \hat{\sigma}_j^2\}^{1/2}} \right| - \lambda(h) \right\}$$

and $\hat{\phi}_{ij,T}(u,h) = \sum_{t=1}^T w_{t,T}(u,h) \{ (\varepsilon_{it} - \bar{\varepsilon}_i) - (\varepsilon_{jt} - \bar{\varepsilon}_j) \}$ with $\bar{\varepsilon}_i = \bar{\varepsilon}_{i,T} = T^{-1} \sum_{t=1}^T \varepsilon_{it}$. Our first theoretical result characterizes the asymptotic behaviour of the statistic $\hat{\Phi}_{n,T}$ and parallels Theorem 3.1 from Section 3.

Theorem 5.1. *Suppose that the error processes $\mathcal{E}_i = \{\varepsilon_{it} : 1 \leq t \leq T\}$ are independent across i and satisfy (C1)–(C3) for each i . Moreover, let (C4)–(C6) be fulfilled and assume that $\hat{\sigma}_i^2 = \sigma_i^2 + o_p(\rho_T)$ with $\rho_T = o(1/\log T)$ for each i . Then*

$$\mathbb{P}(\hat{\Phi}_{n,T} \leq q_{n,T}(\alpha)) = (1 - \alpha) + o(1).$$

Theorem 5.1 is the main stepping stone to derive the theoretical properties of our multiscale test. It can be proven by slightly modifying the arguments for Theorem 3.1. The details are provided in the Supplementary Material. The following proposition characterizes the behaviour of our multiscale test under the null hypothesis and under local alternatives.

Proposition 5.2. *Let the conditions of Theorem 5.1 be satisfied.*

(a) *Under the null hypothesis $H_0 : m_1 = m_2 = \dots = m_n$, it holds that*

$$\mathbb{P}(\widehat{\Psi}_{n,T} \leq q_{n,T}(\alpha)) = (1 - \alpha) + o(1).$$

(b) *Let $m_i = m_{i,T}$ be a Lipschitz continuous function with $\int_0^1 m_{i,T}(w)dw = 0$ for any i . In particular, suppose that $|m_{i,T}(v) - m_{i,T}(w)| \leq L|v - w|$ for all $v, w \in [0, 1]$ and some fixed constant $L < \infty$ which does not depend on T . Moreover, assume that for some pair of indices i and j , the functions $m_{i,T}$ and $m_{j,T}$ have the following property: There exists $(u, h) \in \mathcal{G}_T$ with $[u - h, u + h] \subseteq [0, 1]$ such that $m_{i,T}(w) - m_{j,T}(w) \geq c_T \sqrt{\log T / (Th)}$ for all $w \in [u - h, u + h]$ or $m_{j,T}(w) - m_{i,T}(w) \geq c_T \sqrt{\log T / (Th)}$ for all $w \in [u - h, u + h]$, where $\{c_T\}$ is any sequence of positive numbers with $c_T \rightarrow \infty$. Then*

$$\mathbb{P}(\widehat{\Psi}_{n,T} \leq q_{n,T}(\alpha)) = o(1).$$

Part (a) of Proposition 5.2 is a direct consequence of Theorem 5.1. The proof of part (b) is very similar to that of Proposition 3.3 and thus omitted.

5.4 Clustering of time trends

Consider a situation in which the null hypothesis $H_0 : m_1 = m_2 = \dots = m_n$ is violated. Even though some of the trend functions are different in this case, part of them may still be the same. Put differently, there may be groups of time series which have the same time trend. Formally speaking, we define a group structure as follows: There exist sets or groups of time series G_1, \dots, G_N with $N \leq n$ and $\{1, \dots, n\} = \dot{\bigcup}_{\ell=1}^N G_\ell$ such that for each $1 \leq \ell \leq N$,

$$m_i = g_\ell \quad \text{for all } i \in G_\ell,$$

where g_ℓ are group-specific trend functions. Hence, the time series which belong to the group G_ℓ all have the same time trend g_ℓ . Throughout the section, we suppose that the group-specific trend functions g_ℓ have the following properties: For each ℓ , $g_\ell = g_{\ell,T}$ is a Lipschitz continuous function with $\int_0^1 g_{\ell,T}(w)dw = 0$. In particular, it holds that $|g_{\ell,T}(v) - g_{\ell,T}(w)| \leq L|v - w|$ for all $v, w \in [0, 1]$ and some constant $L < \infty$ that does not depend on T . Moreover, for any $\ell \neq \ell'$, the trends $g_{\ell,T}$ and $g_{\ell',T}$ are assumed to differ in the following sense: There exists $(u, h) \in \mathcal{G}_T$ with $[u - h, u + h] \subseteq [0, 1]$ such that $g_{\ell,T}(w) - g_{\ell',T}(w) \geq c_T \sqrt{\log T / (Th)}$ for all $w \in [u - h, u + h]$ or $g_{\ell',T}(w) - g_{\ell,T}(w) \geq c_T \sqrt{\log T / (Th)}$ for all $w \in [u - h, u + h]$, where $0 < c_T \rightarrow \infty$.

In many applications, it is natural to suppose that there is a group structure in the data. In this case, a particular interest lies in estimating the unknown groups from the

data at hand. In what follows, we combine our multiscale methods with a clustering algorithm to achieve this. More specifically, we use the multiscale statistics $\hat{\Psi}_{ij,T}$ as distance measures which are fed into a hierarchical clustering algorithm. To describe the algorithm, we first need to introduce the notion of a dissimilarity measure: Let $S \subseteq \{1, \dots, n\}$ and $S' \subseteq \{1, \dots, n\}$ be two sets of time series from our sample. We define a dissimilarity measure between S and S' by setting

$$\hat{\Delta}(S, S') = \max_{\substack{i \in S, \\ j \in S'}} \hat{\Psi}_{ij,T}. \quad (5.1)$$

This is commonly called a complete linkage measure of dissimilarity. Alternatively, we may work with an average or a single linkage measure. We now combine the dissimilarity measure $\hat{\Delta}$ with a hierarchical agglomerative clustering (HAC) algorithm which proceeds as follows:

Step 0 (Initialization): Let $\hat{G}_i^{[0]} = \{i\}$ denote the i -th singleton cluster for $1 \leq i \leq n$ and define $\{\hat{G}_1^{[0]}, \dots, \hat{G}_n^{[0]}\}$ to be the initial partition of time series into clusters.

Step r (Iteration): Let $\hat{G}_1^{[r-1]}, \dots, \hat{G}_{n-(r-1)}^{[r-1]}$ be the $n - (r - 1)$ clusters from the previous step. Determine the pair of clusters $\hat{G}_\ell^{[r-1]}$ and $\hat{G}_{\ell'}^{[r-1]}$ for which

$$\hat{\Delta}(\hat{G}_\ell^{[r-1]}, \hat{G}_{\ell'}^{[r-1]}) = \min_{1 \leq k < k' \leq n-(r-1)} \hat{\Delta}(\hat{G}_k^{[r-1]}, \hat{G}_{k'}^{[r-1]})$$

and merge them into a new cluster.

Iterating this procedure for $r = 1, \dots, n - 1$ yields a tree of nested partitions $\{\hat{G}_1^{[r]}, \dots, \hat{G}_{n-r}^{[r]}\}$, which can be graphically represented by a dendrogram. Roughly speaking, the HAC algorithm merges the n singleton clusters $\hat{G}_i^{[0]} = \{i\}$ step by step until we end up with the cluster $\{1, \dots, n\}$. In each step of the algorithm, the closest two clusters are merged, where the distance between clusters is measured in terms of the dissimilarity $\hat{\Delta}$. We refer the reader to Section 14.3.12 in Hastie et al. (2009) for an overview of hierarchical clustering methods.

When the number of groups N is known, we estimate the group structure $\{G_1, \dots, G_N\}$ by the N -partition $\{\hat{G}_1^{[n-N]}, \dots, \hat{G}_N^{[n-N]}\}$ produced by the HAC algorithm. When N is unknown, we estimate it by the \hat{N} -partition $\{\hat{G}_1^{[n-\hat{N}]}, \dots, \hat{G}_{\hat{N}}^{[n-\hat{N}]}\}$, where \hat{N} is an estimator of N . The latter is defined as

$$\hat{N} = \min \left\{ r = 1, 2, \dots \mid \max_{1 \leq \ell \leq r} \hat{\Delta}(\hat{G}_\ell^{[n-r]}) \leq q_{n,T}(\alpha) \right\},$$

where we write $\hat{\Delta}(S) = \hat{\Delta}(S, S)$ for short and $q_{n,T}(\alpha)$ is the $(1 - \alpha)$ -quantile of $\Phi_{n,T}$ defined in Section 5.2.

The following proposition summarizes the theoretical properties of the estimators \widehat{N} and $\{\widehat{G}_1, \dots, \widehat{G}_{\widehat{N}}\}$, where we use the shorthand $\widehat{G}_\ell = \widehat{G}_\ell^{[n-\widehat{N}]}$ for $1 \leq \ell \leq \widehat{N}$.

Proposition 5.3. *Let the conditions of Theorem 5.1 be satisfied. Then*

$$\mathbb{P}\left(\{\widehat{G}_1, \dots, \widehat{G}_{\widehat{N}}\} = \{G_1, \dots, G_N\}\right) \geq (1 - \alpha) + o(1)$$

and

$$\mathbb{P}(\widehat{N} = N) \geq (1 - \alpha) + o(1).$$

This result allows us to make statistical confidence statements about the estimated clusters $\{\widehat{G}_1, \dots, \widehat{G}_{\widehat{N}}\}$ and their number \widehat{N} . In particular, we can claim with asymptotic confidence $\geq 1 - \alpha$ that the estimated group structure is identical to the true group structure. Note that it is possible to let the significance level α depend on the sample size T in Proposition 5.3. In particular, we can allow $\alpha = \alpha_T$ to converge slowly to zero as $T \rightarrow \infty$, in which case we obtain that $\mathbb{P}(\{\widehat{G}_1, \dots, \widehat{G}_{\widehat{N}}\} = \{G_1, \dots, G_N\}) \rightarrow 1$ and $\mathbb{P}(\widehat{N} = N) \rightarrow 1$. The proof of Proposition 5.3 can be found in the Supplementary Material.

Our multiscale methods do not only allow us to compute estimators of the unknown groups G_1, \dots, G_N . They also provide information on the locations where two group-specific trend functions g_ℓ and $g_{\ell'}$ differ from each other. To turn this claim into a mathematically precise statement, we need to introduce some notation. First of all, note that the indexing of the estimators $\widehat{G}_1, \dots, \widehat{G}_{\widehat{N}}$ is completely arbitrary. We could, for example, change the indexing according to the rule $\ell \mapsto \widehat{N} - \ell + 1$. In what follows, we suppose that the estimated groups are indexed such that $P(\widehat{G}_\ell = G_\ell \text{ for all } \ell) \geq (1 - \alpha) + o(1)$. Theorem 5.3 implies that this is possible without loss of generality. Keeping this convention in mind, we define the sets

$$\mathcal{A}_{n,T}(\ell, \ell') = \left\{ (u, h) \in \mathcal{G}_T : \left| \frac{\widehat{\psi}_{ij,T}(u, h)}{(\widehat{\sigma}_i^2 + \widehat{\sigma}_j^2)^{1/2}} \right| > q_{n,T}(\alpha) + \lambda(h) \text{ for some } i \in \widehat{G}_\ell, j \in \widehat{G}_{\ell'} \right\}$$

and

$$\Pi_{n,T}(\ell, \ell') = \{I_{u,h} = [u - h, u + h] : (u, h) \in \mathcal{A}_{n,T}(\ell, \ell')\}$$

for $1 \leq \ell < \ell' \leq \widehat{N}$. An interval $I_{u,h}$ is contained in $\Pi_{n,T}(\ell, \ell')$ if our multiscale test indicates a significant difference between the trends m_i and m_j on the interval $I_{u,h}$ for some $i \in \widehat{G}_\ell$ and $j \in \widehat{G}_{\ell'}$. Put differently, $I_{u,h} \in \Pi_{n,T}(\ell, \ell')$ if the test suggests a significant difference between the trends of the ℓ -th and the ℓ' -th group on the interval $I_{u,h}$. We further let

$$E_{n,T}(\ell, \ell') = \left\{ \forall I_{u,h} \in \Pi_{n,T}(\ell, \ell') : g_\ell(v) \neq g_{\ell'}(v) \text{ for some } v \in I_{u,h} = [u - h, u + h] \right\}$$

be the event that the group-specific time trends g_ℓ and $g_{\ell'}$ differ on all intervals $I_{u,h} \in \Pi_{n,T}(\ell, \ell')$.

$\Pi_{n,T}(\ell, \ell')$. With this notation at hand, we can make the following formal statement whose proof is given in the Supplementary Material.

Proposition 5.4. *Under the conditions of Theorem 5.1, the event*

$$E_{n,T} = \left\{ \bigcap_{1 \leq \ell < \ell' \leq \hat{N}} E_{n,T}(\ell, \ell') \right\} \cap \left\{ \hat{N} = N \text{ and } \hat{G}_\ell = G_\ell \text{ for all } \ell \right\}$$

asymptotically occurs with probability $\geq 1 - \alpha$, that is,

$$\mathbb{P}(E_{n,T}) \geq (1 - \alpha) + o(1).$$

The statement of Proposition 5.4 remains to hold true when the sets of intervals $\Pi_{n,T}(\ell, \ell')$ are replaced by the corresponding sets of minimal intervals. According to Proposition 5.4, the sets $\Pi_{n,T}(\ell, \ell')$ allow us to locate, with a pre-specified confidence, time regions where the group-specific trend functions g_ℓ and $g_{\ell'}$ differ from each other. In particular, we can claim with asymptotic confidence $\geq 1 - \alpha$ that the trend functions g_ℓ and $g_{\ell'}$ differ on all intervals $I_{u,h} \in \Pi_{n,T}(\ell, \ell')$.

6 Estimation of the long-run error variance

We now discuss how to estimate the long-run error variance $\sigma^2 = \sum_{\ell=-\infty}^{\infty} \gamma(\ell)$ with $\gamma(\ell) = \text{Cov}(\varepsilon_0, \varepsilon_\ell)$ in model (2.1). The same methods can be applied in the context of model (2.2). A number of different methods have been established in the literature to estimate the long-run error variance σ^2 in the trend model (2.1) under various assumptions on the error terms. In what follows, we give a brief overview of estimation methods which are suitable for our purposes. We in particular focus attention on difference-based methods as these have the following advantage: They do not involve a nonparametric estimator of the function m and thus do not require to specify a smoothing parameter for the estimation of m .

In principle, it is possible to construct an estimator of σ^2 under the general conditions on the error process laid out in Section 2 (or at least under somewhat stronger versions of these conditions). However, as is well-known, it is quite involved to estimate the long-run variance of a time series process under general conditions, the resulting estimators often tending to be quite imprecise. From a practical point of view, one might thus prefer to impose some time series model on the error terms and to estimate σ^2 under the restrictions of this model. Of course, this will create some bias due to misspecification. However, as long as the model gives a reasonable approximation to the true error process, this bias may very well be less severe than the error stemming from the instable behaviour of a general estimator of σ^2 . In what follows, we consider an autoregressive (AR) model for the error terms since this error model is widely used in practice and is also appropriate for our applications in Section 8.

6.1 Independent error terms

Before we discuss the case of autoregressive error terms, we introduce the idea of difference-based methods for estimating σ^2 in the simple case of i.i.d. errors ε_t . In this case, σ^2 is identical to the variance of the random variables ε_t , that is, $\sigma^2 = \text{Var}(\varepsilon_t)$. Let $D_\ell Y_t = Y_t - Y_{t-\ell}$ denote the difference between Y_t and $Y_{t-\ell}$ and suppose that m is sufficiently smooth. In particular, assume that m is Lipschitz continuous on $[0, 1]$, that is, $|m(u) - m(v)| \leq C|u - v|$ for all $u, v \in [0, 1]$ and some constant $C < \infty$. Under these conditions, it holds that $|m(\frac{t}{T}) - m(\frac{t-\ell}{T})| \leq C\ell/T$, which implies that $D_\ell Y_t = D_\ell \varepsilon_t + O(\ell/T)$ uniformly over t . Hence, the observed differences $D_\ell Y_t$ approximate the unobserved differences of the error terms $D_\ell \varepsilon_t$. This together with the fact that $\mathbb{E}[\{D_\ell \varepsilon_t\}^2]/2 = \sigma^2$ suggests to estimate σ^2 by $\hat{\sigma}^2 = (T - \ell)^{-1} \sum_{t=\ell+1}^T \{D_\ell Y_t\}^2/2$, where most commonly $\ell = 1$. As can be easily verified, the estimator $\hat{\sigma}^2$ has the property that $\hat{\sigma}^2 = \sigma^2 + O_p(T^{-1/2})$.

6.2 Autoregressive error terms

The differencing approach presented above can be extended to more complicated error structures. For the case of k -dependent error terms, estimators for σ^2 have been proposed by Müller and Stadtmüller (1988), Herrmann et al. (1992) and Tecuapetla-Gómez and Munk (2017) among others. We here focus attention on the case of autoregressive error terms. Specifically, we suppose that $\{\varepsilon_t\}$ is an $\text{AR}(p)$ process of the form

$$\varepsilon_t = \sum_{j=1}^p a_j \varepsilon_{t-j} + \eta_t,$$

where a_1, \dots, a_p are unknown parameters and η_t are i.i.d. innovations with $\mathbb{E}[\eta_t] = 0$ and $\mathbb{E}[\eta_t^2] = \sigma_\eta^2$. Throughout the discussion, we assume that $\{\varepsilon_t\}$ is a stationary and causal $\text{AR}(p)$ process of known order p with finite fourth moment $\mathbb{E}[\varepsilon_t^4] < \infty$. A difference-based method to estimate the long-run variance σ^2 of the $\text{AR}(p)$ error process $\{\varepsilon_t\}$ in model (2.1) has been developed in Hall and Van Keilegom (2003). Their estimator $\hat{\sigma}^2$ is constructed in the following three steps:

Step 1. We first set up an estimator of the autocovariance $\gamma(\ell) = \text{Cov}(\varepsilon_t, \varepsilon_{t+\ell})$ for a given lag ℓ . As in the case of independent errors, it holds that $D_\ell Y_t = D_\ell \varepsilon_t + O(\ell/T)$ uniformly over t provided that m is Lipschitz. This together with the fact that $\mathbb{E}[\{D_\ell \varepsilon_t\}^2]/2 = \gamma(0) - \gamma(\ell)$ motivates to estimate $\gamma(0)$ by

$$\hat{\gamma}(0) = \frac{1}{L_2 - L_1 + 1} \sum_{r=L_1}^{L_2} \frac{1}{2(T-r)} \sum_{t=r+1}^T \{D_r Y_t\}^2,$$

where $L_1 \leq L_2$ are tuning parameters which are discussed in more detail below. More-

over, an estimator of $\gamma(\ell)$ for $1 \leq \ell \leq p$ is given by

$$\hat{\gamma}(\ell) = \hat{\gamma}(0) - \frac{1}{2(T-\ell)} \sum_{t=\ell+1}^T \{D_\ell Y_t\}^2.$$

As $\gamma(\ell) = \gamma(-\ell)$, we finally set $\hat{\gamma}(-\ell) = \hat{\gamma}(\ell)$ for $1 \leq \ell \leq p$.

Step 2. We next estimate the AR coefficients $(a_1, \dots, a_p)^\top$ by the Yule-Walker estimators $(\hat{a}_1, \dots, \hat{a}_p)^\top = \hat{\Gamma}^{-1}(\hat{\gamma}(1), \dots, \hat{\gamma}(p))^\top$, where the matrix $\hat{\Gamma}$ is given by $\hat{\Gamma} = \{\hat{\gamma}(|k-\ell|)\}_{1 \leq k, \ell \leq p}$.

Step 3. Let $\hat{d}_0 = 1$ and define the parameters $\hat{d}_1, \hat{d}_2, \dots$ by the equation $1 + \sum_{\ell=1}^{\infty} \hat{d}_\ell z^\ell = (1 - \sum_{j=1}^p \hat{a}_j z^j)^{-1}$. In the AR(1) case $\varepsilon_t = a\varepsilon_{t-1} + \eta_t$, for instance, it holds that $\sum_{\ell=0}^{\infty} \hat{a}^\ell z^\ell = (1 - \hat{a}z)^{-1}$ and thus $\hat{d}_\ell = \hat{a}^\ell$ for $\ell \geq 1$. The variance $\sigma_\eta^2 = \mathbb{E}[\eta_t^2]$ of the innovations can be estimated by $\hat{\sigma}_\eta^2 = \hat{\gamma}(0)/(\sum_{\ell=0}^{\infty} \hat{d}_\ell^2)$. With this notation at hand, we define

$$\hat{\sigma}^2 = \hat{\sigma}_\eta^2 \left(1 - \sum_{j=1}^p \hat{a}_j\right)^{-2}$$

to be our estimator of the long-run error variance σ^2 .

The estimator $\hat{\sigma}^2$ depends on the two tuning parameters L_1 and L_2 which are required to compute $\hat{\gamma}(0)$. To better understand the role of these tuning parameters, let us have a closer look at the estimator $\hat{\gamma}(0)$. As $\mathbb{E}[\{D_\ell Y_t\}^2]/2 = \mathbb{E}[\{D_\ell \varepsilon_t\}^2]/2 + O(\{\ell/T\}^2) = \gamma(0) - \gamma(\ell) + O(\{\ell/T\}^2)$, it can be easily shown that

$$\mathbb{E}[\hat{\gamma}(0)] = \gamma(0) - \frac{1}{L_2 - L_1 + 1} \sum_{r=L_1}^{L_2} \gamma(r) + O\left(\left\{\frac{L_2}{T}\right\}^2\right).$$

The two bias terms $\sum_{r=L_1}^{L_2} \gamma(r)/(L_2 - L_1 + 1)$ and $O(\{L_2/T\}^2)$ can be asymptotically neglected if we choose the tuning parameters L_1 and L_2 appropriately. Since $\{\varepsilon_t\}$ is an AR(p) process, the autocovariances $\gamma(r)$ decay exponentially fast to zero as $r \rightarrow \infty$. Hence, the bias term $\sum_{r=L_1}^{L_2} \gamma(r)/(L_2 - L_1 + 1)$ is asymptotically negligible if L_1 grows sufficiently fast with the sample size T . Due to the exponential decay of the autocovariances, it in particular suffices to assume that $L_1/\log T \rightarrow \infty$. For the second bias term $O(\{L_2/T\}^2)$ to be asymptotically negligible, we need to assume that L_2 grows more slowly than the sample size T . In practice, L_1 should be chosen so large that the autocovariances $\gamma(\ell)$ with $\ell \geq L_1$ can be expected to be close to zero, ensuring that the bias term $\sum_{r=L_1}^{L_2} \gamma(r)/(L_2 - L_1 + 1)$ is sufficiently small. The choice of L_2 can be expected to be less important in practice than that of L_1 as long as we do not pick L_2 too close to the sample size T . As pointed out in Hall and Van Keilegom (2003), it can be shown that $\hat{\sigma}^2 = \sigma^2 + O_p(T^{-1/2})$ provided that $L_1/\log T \rightarrow \infty$ and $L_2 = O(T^{1/2})$.

7 Simulations

To assess the finite sample performance of the methods from Sections 4 and 5, we conduct a number of simulations. We first investigate the test procedure from Section 4. The simulation design is set up to mimic the situation in the application example of Section 8.1: We generate data from the model $Y_t = m(\frac{t}{T}) + \varepsilon_t$ for different time series lengths T . The errors ε_t are drawn from the AR(1) process $\varepsilon_t = a\varepsilon_{t-1} + \eta_t$, where η_t are independent and normally distributed with mean 0 and variance σ_η^2 . We set $a = 0.267$ and $\sigma_\eta^2 = 0.35$, thus matching the estimated values obtained in the application of Section 8.1. To simulate data under the null $H_0 : m' = 0$, we let m be a constant function. In particular, we set $m = 0$ without loss of generality. To generate data under the alternative, we consider the trend functions $m(u) = \beta(u - 0.6)1(0.6 \leq u \leq 1)$ with $\beta = 1.25, 1.875, 2.5$. These functions are broken lines with a kink at $u = 0.6$ and different slopes β . The slope parameter β corresponds to a trend with the value $m(1) = 0.4\beta$ at the right endpoint $u = 1$. We thus consider broken lines with the values $m(1) = 0.5, 0.75, 1.0$. Inspecting the middle panel of Figure 2, the broken line with the slope $\beta = 2.5$ can be seen to resemble the local linear trend estimates in the real-data example of Section 8.1 (where we neglect the nonlinearities of the local linear fits at the beginning of the observation period), whereas the trend functions with smaller values of slope $\beta = 1.25, 1.875$ are closer to the null making it harder for our test to detect the difference.

To implement our test methods, we choose K to be an Epanechnikov kernel and define the set \mathcal{G}_T of location-scale points (u, h) as

$$\mathcal{G}_T = \left\{ (u, h) : u = 5k/T \text{ for some } 1 \leq k \leq T/5 \text{ and } h = (3 + 5\ell)/T \text{ for some } 0 \leq \ell \leq T/20 \right\}. \quad (7.1)$$

For the bandwidth value $h = (3 + 5\ell)/T$, the local linear weights $w'_{t,T}(u, h)$ give a non-zero weight to exactly $5 + 5\ell$ observations. Hence, the bandwidths h considered in \mathcal{G}_T correspond to effective sample sizes of 5, 10, 15, \dots up to approximately $T/2$ data points. Moreover, we take into account all rescaled time points $u \in [0, 1]$ on an equidistant grid with step length $5/T$. The long-run error variance σ^2 is estimated by the procedure from Section 6.2, setting the tuning parameters L_1 and L_2 to $\lfloor \sqrt{T} \rfloor$ and $\lfloor 2\sqrt{T} \rfloor$, respectively. To compute the critical values of the test, we simulate 1000 values of the statistic Φ'_T defined in Section 4.2 and compute their empirical $(1 - \alpha)$ quantile $q'_T(\alpha)$.

Tables 1 and 2 report the simulation results for the sample sizes $T = 250, 350, 500, 1000$ and the confidence levels $\alpha = 0.01, 0.05, 0.10$. The sample size $T = 350$ is approximately equal to the time series length 359 in the real-data example of Section 8.1. To produce our simulation results, we generate $S = 1000$ samples for each time series length T and carry out the multiscale test for each simulated sample. The entries of Tables 1

Table 1: Size of the multiscale test from Section 4 for different sample sizes T and nominal sizes α .

T	nominal size α		
	0.01	0.05	0.1
250	0.004	0.022	0.082
350	0.007	0.031	0.067
500	0.011	0.056	0.086
1000	0.011	0.060	0.098

Table 2: Power of the multiscale test from Section 4 for different sample sizes T and nominal sizes α . Each panel corresponds to a different slope parameter β .

(a) $\beta = 1.25$				(b) $\beta = 1.875$				(c) $\beta = 2.5$			
T	nominal size α			T	nominal size α			T	nominal size α		
	0.01	0.05	0.1		0.01	0.05	0.1		0.01	0.05	0.1
250	0.107	0.223	0.358	250	0.365	0.582	0.709	250	0.699	0.876	0.928
350	0.216	0.374	0.500	350	0.644	0.779	0.845	350	0.943	0.973	0.992
500	0.280	0.554	0.678	500	0.784	0.942	0.976	500	0.984	1.000	1.000
1000	0.756	0.910	0.935	1000	0.997	1.000	1.000	1000	1.000	1.000	1.000

and 2 are computed as the number of simulations in which the test rejects divided by the total number of simulations. As can be seen from Table 1, the actual size of the test is fairly close to the nominal target α even for small values of T . Hence, the test has approximately the correct size. Inspecting Table 2, one can further see that the test has reasonable power properties. For the smallest value $\beta = 1.25$, the deviation from the null is quite small, making it hard for the test to detect the alternative. As a consequence, the power is only moderate for $T = 250$ and $T = 350$. When we move further away from the null by increasing the slope parameter β , the power of the test quickly increases. It can also be seen to rapidly get larger as the sample size grows. For the slope $\beta = 2.5$ and the sample size $T = 350$, which are the values that resemble the real-life data the most, the power of the test is above 93% for all significance levels α considered and thus comes quite close to 1.

We next turn to the test methods from Section 5. The simulation design extends the setup from above. We generate data from the model $Y_{it} = m_i(\frac{t}{T}) + \varepsilon_{it}$, where the number of time series is set to $n = 15$ and we consider different time series lengths T . For each i , the errors ε_{it} are drawn from the AR(1) model $\varepsilon_{it} = a\varepsilon_{i,t-1} + \eta_{it}$, where as before $a = 0.267$ and the innovations η_{it} are i.i.d. normally distributed with mean 0 and variance 0.35. To generate data under the null $H_0 : m_1 = \dots = m_n$, we let $m_i = 0$ for all i without loss of generality. To produce data under the alternative, we define $m_1(u) = \beta(u - 0.5)$ with $\beta = 0.75, 1, 1.25$ and set $m_i = 0$ for all $i \neq 1$. Hence, all trend functions are the same except for m_1 which is an increasing linear function with

Table 4: Size of the multiscale test from Section 5 for $n = 15$ time series, different sample sizes T and nominal sizes α .

T	nominal size α		
	0.01	0.05	0.1
250	0.096	0.096	0.096
300	0.096	0.096	0.096
500	0.096	0.096	0.096
1000	0.096	0.096	0.096

Table 5: Power of the multiscale test from Section 5 for $n = 15$ time series, different sample sizes T and nominal sizes α . Each panel corresponds to a different slope parameter β .

(a) $\beta = 0.75$				(b) $\beta = 1.00$				(c) $\beta = 1.25$			
T	nominal size α			T	nominal size α			T	nominal size α		
	0.01	0.05	0.1		0.01	0.05	0.1		0.01	0.05	0.1
250	0.354	0.557	0.687	250	0.758	0.895	0.946	250	0.961	0.990	0.997
300	0.349	0.659	0.772	300	0.790	0.957	0.978	300	0.991	1.000	1.000
500	0.859	0.946	0.964	500	0.997	0.999	0.999	500	1.000	1.000	1.000
1000	0.997	1.000	1.000	1000	1.000	1.000	1.000	1000	1.000	1.000	1.000

$m_1(0.5) = 0$. We here use a linear function rather than a broken line in order to satisfy the normalization constraint $\int_0^1 m_1(u)du = 0$.

The test is implemented analogously as above. As before, we work with an Epanechnikov kernel, we define the grid \mathcal{G}_T as in (7.1) and we set the two tuning parameters L_1 and L_2 to $\lfloor \sqrt{T} \rfloor$ and $\lfloor 2\sqrt{T} \rfloor$ respectively. In order to compute the critical values of the test, we simulate 1000 values of the statistic $\Phi_{n,T}$ defined in Section 5.2 and compute their empirical $(1 - \alpha)$ quantile $q_{n,T}(\alpha)$. Note that the statistic $\Phi_{n,T}$ depends on the estimators $\hat{\sigma}_i^2$ of the long-run error variances σ_i^2 . This implies that for each simulated sample, we have to recompute the quantile $q_{n,T}(\alpha)$ and thus the critical value of the test. This is of course computationally extremely expensive. In order to circumvent this issue, we make the additional assumption that the long-run error variance is known to be the same across i , that is, $\sigma_i^2 = \sigma^2$ for all i . Under this assumption, we can estimate σ^2 by $\hat{\sigma}^2 = (\sum_{i=1}^n \hat{\sigma}_i^2)/n$, and the Gaussian statistic $\Phi_{n,T}$ simplifies to $\Phi_{n,T} = \max_{1 \leq i < j \leq n} \Phi_{ij,T}$ with $\Phi_{ij,T} = \max_{(u,h) \in \mathcal{G}_T} \{|\phi_{ij,T}(u,h)| - \lambda(h)\}$ and $\phi_{ij,T}(u,h) = \sum_{t=1}^T w_{t,T}(u,h) \{(Z_{it} - \bar{Z}_i) - (Z_{jt} - \bar{Z}_j)\}$. This statistic does not depend on the estimators $\hat{\sigma}_i^2$ anymore. We thus do not need to recompute the critical values for each simulated sample, which decreases the running time significantly.

The simulation results are reported in Tables 4 and 5. The entries of the tables are computed in the same way as those in Tables 1 and 2. Inspecting Table 4, the actual size of the test can be seen to approximate the nominal target α quite well even for small values of T . Moreover, as can be seen from Table 5, the test also has reasonable

Table 7: Clustering results for different sample sizes T and nominal sizes α .

(a) Empirical probabilities that $\hat{N} = N$					(b) Empirical probabilities that $\{\hat{G}_1, \dots, \hat{G}_{\hat{N}}\} = \{G_1, G_2, G_3\}$				
nominal size α					nominal size α				
T	0.01	0.05	0.1		T	0.01	0.05	0.1	
250	0.711	0.911	0.944		250	0.581	0.747	0.776	
300	0.815	0.963	0.961		300	0.741	0.886	0.889	
500	0.990	0.978	0.969		500	0.984	0.974	0.966	
1000	0.998	0.987	0.972		1000	0.998	0.987	0.972	

power against the alternatives considered. For the smallest slope $\beta = 0.75$ and the smaller sample sizes $T = 250, 350$, the power is only moderate, reflecting the fact that the alternative is not very far away from the null. However, as we increase the slope β and the sample size T , the power quickly increases. For the largest slope $\beta = 1.25$ and $T = 350$, we already reach a power of 1.00.

We finally investigate the finite sample performance of the clustering algorithm from Section 5.4. To do so, we partition the $n = 15$ time series into $N = 3$ groups, each containing 5 time series. Specifically, we set $G_1 = \{1, \dots, 5\}$, $G_2 = \{6, \dots, 10\}$ and $G_3 = \{11, \dots, 15\}$. Moreover, we define the group-specific trend functions g_1 , g_2 and g_3 by $g_1(u) = 0$, $g_2(u) = 1 \cdot (u - 0.5)$ and $g_3(u) = (-1) \cdot (u - 0.5)$. In order to compute our estimators of the groups G_1 , G_2 , G_3 and their number $N = 3$, we use the same implementation as before followed by the clustering procedure from Section 5.4. The estimation results are reported in Table 7. The entries in Table 7a are computed as the number of simulations for which $\hat{N} = N$ divided by the total number of simulations $S = 1000$. They thus specify the empirical probabilities with which the estimated number of groups \hat{N} is equal to the true number $N = 3$. Analogously, the entries of Table 7b give the empirical probabilities with which the estimated group structure $\{\hat{G}_1, \dots, \hat{G}_{\hat{N}}\}$ equals the true one $\{G_1, G_2, G_3\}$.

The simulation results nicely illustrate the theoretical properties of our clustering algorithm. According to Proposition 5.3, the probability that $\hat{N} = N$ and $\{\hat{G}_1, \dots, \hat{G}_{\hat{N}}\} = \{G_1, G_2, G_3\}$ should be at least $(1 - \alpha)$ asymptotically. For the sample sizes $T = 500$ and $T = 1000$, the empirical probabilities reported in Table 7 can indeed be seen to exceed the value $(1 - \alpha)$ as predicted by Proposition 5.3. For the smaller sample sizes $T = 250$ and $T = 350$, in contrast, the empirical probabilities are mostly below $(1 - \alpha)$. This reflects the asymptotic nature of Proposition 5.3 and is not very surprising. It simply mirrors the fact that for small sample sizes, the effective noise level is very high. Even though below the target of $(1 - \alpha)$, the empirical probabilities for $T = 250$ and $T = 350$ are still quite substantial. Hence, even for these small sample sizes, our estimates \hat{N} and $\{\hat{G}_1, \dots, \hat{G}_{\hat{N}}\}$ are equal to their true counterparts in a large number of simulations.

8 Applications

In what follows, we illustrate the multiscale methods from Sections 4 and 5 by two real-data examples. In the first example, we apply the test method from Section 4 to a long time series of temperature data from Central England. In the second, we analyse a sample of temperature time series from 34 different weather stations in Great Britain with the help of the methods from Section 5.

8.1 Analysis of Central England temperature data

The analysis of time trends in long temperature records is an important task in climatology. Information on the shape of the trend is needed in order to better understand long-term climate variability. The Central England temperature record is the longest instrumental temperature time series in the world. It is a valuable asset for analysing climate variability over the last few hundred years. The data is publicly available on the webpage of the UK Met Office. A detailed description of the data can be found in Parker et al. (1992). For our analysis, we use the dataset of yearly mean temperatures which consists of $T = 359$ observations covering the years from 1659 to 2017. We assume that the data follow the nonparametric trend model

$$Y_t = m\left(\frac{t}{T}\right) + \varepsilon_t,$$

where m is the unknown time trend of interest. The error process $\{\varepsilon_t\}$ is supposed to have the AR(1) structure $\varepsilon_t = a\varepsilon_{t-1} + \eta_t$, where η_t are i.i.d. innovations with mean 0 and variance σ_η^2 . As pointed out in Mudelsee (2010) among others, this is the most widely used error model for discrete climate time series. We estimate the unknown parameters a and σ_η^2 by the procedure from Section 6.2 which yields the estimates $\hat{a} \approx 0.267$ and $\hat{\sigma}_\eta^2 \approx 0.35$.

With the help of our multiscale method from Section 4, we test the null hypothesis $H_0 : m' = 0$, that is, the hypothesis that m is constant. To do so, we set the significance level to $\alpha = 0.05$ and implement the test in exactly the same way as in the simulations of Section 7. The results are presented in Figure 2. The upper panel shows the raw temperature time series, whereas the middle panel depicts local linear kernel estimates of the trend m for different bandwidths h . As one can see, the shape of the estimated time trend strongly differs with the chosen bandwidth. When the bandwidth is small, there are many local increases and decreases in the estimated trend. When the bandwidth is large, most of these local variations get smoothed out. Hence, by themselves, the nonparametric fits do not give much information on whether the trend m is increasing or decreasing in certain time regions.

Our multiscale test provides this kind of information, which is summarized in the lower panel of Figure 2. The plot depicts the minimal intervals contained in the set Π_T^+ which

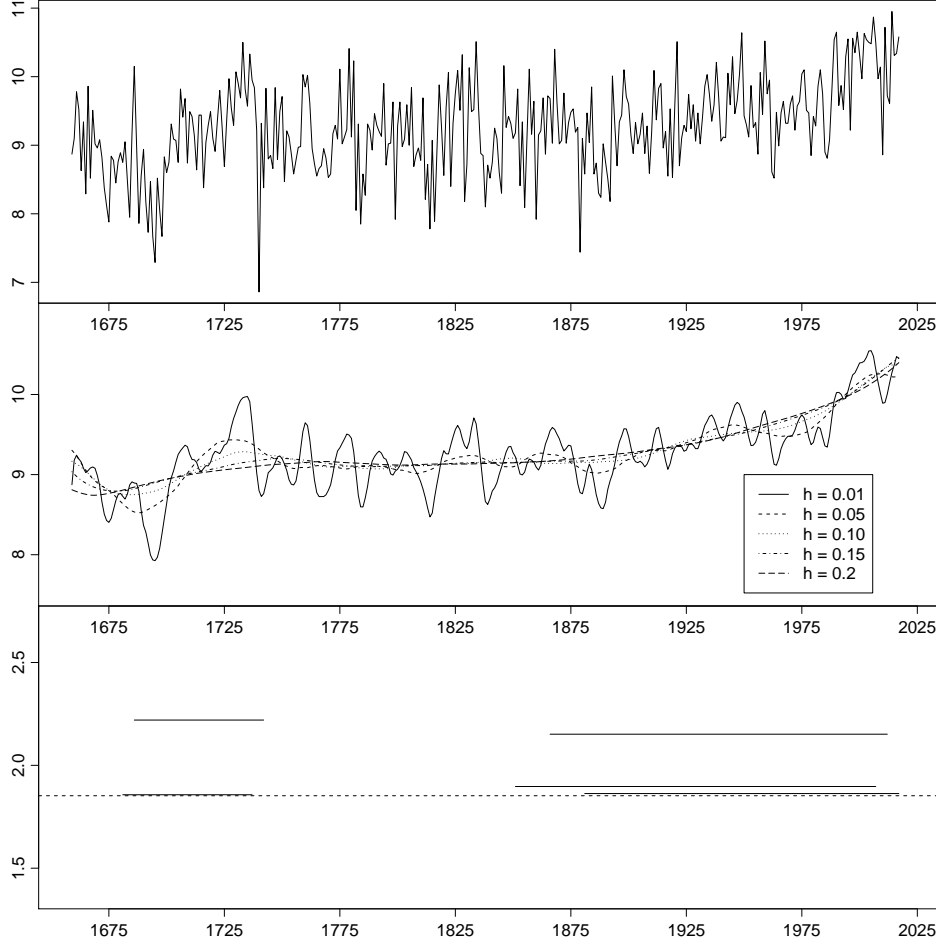


Figure 2: Summary of the application results from Section 8.1. The upper panel shows the Central England mean temperature time series. The middle panel depicts local linear kernel estimates of the time trend for a number of different bandwidths h . The lower panel presents the minimal intervals in the set Π_T^+ produced by the multiscale test. These are $[1681, 1737]$, $[1686, 1742]$, $[1851, 2007]$, $[1866, 2012]$ and $[1881, 2017]$.

is defined in Section 4.3. The set of intervals Π_T^- is empty in the present case. The height at which a minimal interval $I_{u,h} = [u-h, u+h] \in \Pi_t^+$ is plotted indicates the value of the corresponding (additively corrected) test statistic $\hat{\psi}_T'(u, h)/\hat{\sigma} - \lambda(h)$. The dashed line specifies the critical value $q'_T(\alpha)$, where $\alpha = 0.05$ as already mentioned above. According to Proposition 4.3, we can make the following simultaneous confidence statement about the collection of minimal intervals in Π_T^+ . We can claim, with confidence of about 95%, that the trend function m has some increase on each minimal interval. More specifically, we can claim with this confidence that there has been some upward movement in the trend both in the period from around 1680 to 1740 and in the period from around 1880 onwards. Hence, our test in particular provides evidence that there has been some warming trend in the period over the last 150 years or so. On the other hand, as the set Π_T^- is empty, there is no evidence of any downward movement of the trend.

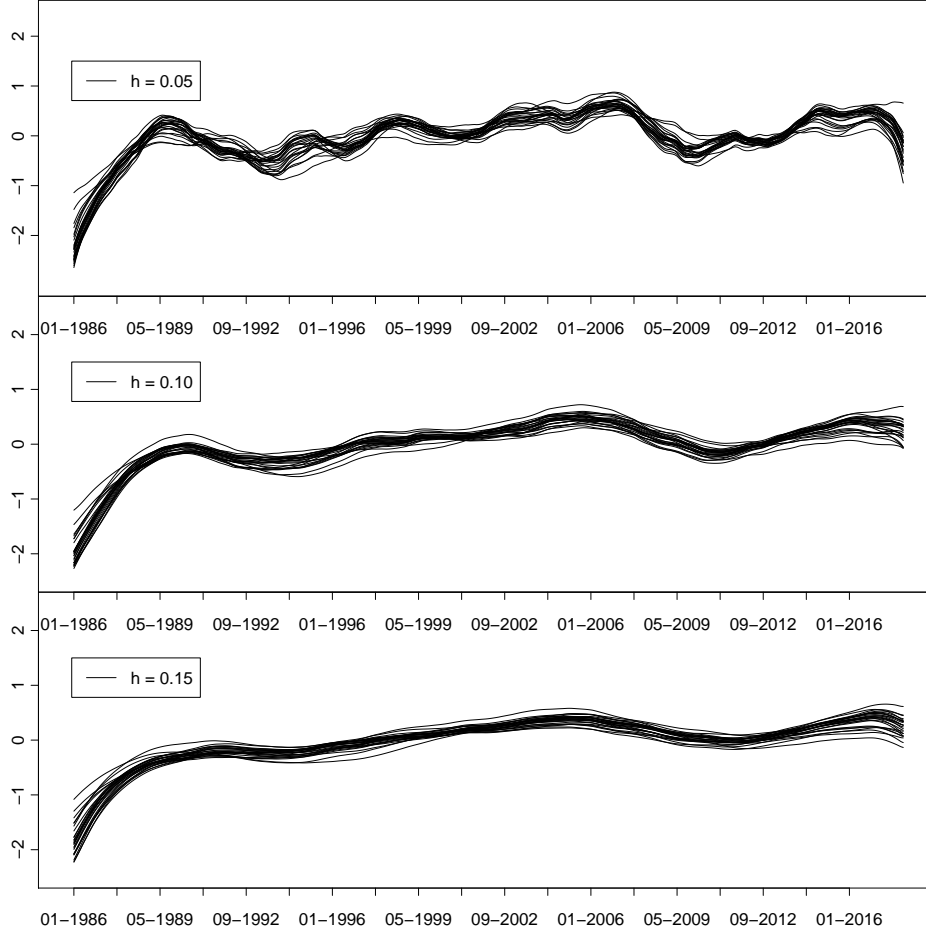


Figure 3: Local linear kernel estimates of the $n = 25$ time trends from the application of Section 8.2. Each panel shows the estimates for a different bandwidth h .

8.2 Analysis of UK weather station data

To illustrate our test method from Section 5, we examine a dataset of monthly mean temperatures from 34 different UK weather stations. The data are publicly available on the webpage of the UK Met Office. We use a subset of 25 stations for which data are available over the time span from 1986 to 2017. We thus observe a time series $\mathcal{Y}_i = \{Y_{it} : 1 \leq t \leq T\}$ of length $T = 396$ for each station $i \in \{1, \dots, 25\}$. The time series \mathcal{Y}_i is assumed to follow the model

$$Y_{it} = \alpha_i(t) + m_i\left(\frac{t}{T}\right) + \varepsilon_{it}, \quad (8.1)$$

where m_i is an unknown nonparametric time trend and $\alpha_i(t)$ is a month-specific intercept which captures the seasonality pattern in the data. We suppose that $\alpha_i(t) = \alpha(t + 12\ell)$ for any integer ℓ , that is, we have a different intercept $\alpha_i(k)$ for each month $k = 1, \dots, 12$. The test method and the underlying theory from Section 5 can be easily adapted to model (8.1), which is a slight extension of model (2.2). The details are

provided below. As in Section 8.1, the error process $\mathcal{E}_i = \{\varepsilon_{it} : 1 \leq t \leq T\}$ is assumed to have the AR(1) structure $\varepsilon_{it} = a_i \varepsilon_{i,t-1} + \eta_{it}$ for each i , where η_{it} are i.i.d. innovations with mean zero.

We aim to test whether the time trend m_i is the same at each of the 25 weather stations. In other words, we want to test the null hypothesis $H_0 : m_1 = \dots = m_n$ with $n = 25$ in model (8.1). To do so, we apply the multiscale test from Section 5 with two minor modifications: (i) We define $\hat{Y}_{it} = Y_{it} - \hat{\alpha}_i(t)$, where $\hat{\alpha}_i(t)$ is an estimator of $\alpha_i(t)$. In particular, we set $\hat{\alpha}_i(k) = T_k^{-1} \sum_{t=1}^T 1_k(t) Y_{it}$, where $1_k(t) = 1(t = k + \lfloor (t-1)/12 \rfloor \cdot 12)$ and $T_k = \sum_{t=1}^T 1_k(t)$. (ii) We define the Gaussian statistic $\Phi_{n,T}$ as in Section 5.2 with $\phi_{ij,T}(u, h) = \sum_{t=1}^T w_{t,T}(u, h) \{\hat{\sigma}_i(Z_{it} - \bar{Z}_i(t)) - \hat{\sigma}_j(Z_{jt} - \bar{Z}_j(t))\}$, where $\bar{Z}_i(t) = \sum_{k=1}^{12} 1_k(t) \{T_k^{-1} \sum_{s=1}^T 1_k(s) Z_{is}\}$. Apart from these two modifications, the multiscale test is constructed exactly as described in Section 5. We implement the test in the same way as in the simulations of Section 7.

We are now ready to apply the test procedure to the data. Figure 3 depicts the local linear estimates of the trend functions m_i for the $n = 25$ different stations. Each panel corresponds to a different bandwidth h . As can be seen, for a given bandwidth h , the fits look very similar to each other. Visual inspection thus suggests that there are no strong differences between the time trends m_i . Our test confirms this impression. It does not reject the null hypothesis at the most common levels $\alpha = 0.01, 0.05, 0.1$. Hence, the test does not provide any evidence for a violation of the null.

Appendix

In what follows, we prove the theoretical results from Section 3. The proofs of the results from Sections 4 and 5 are deferred to the Supplementary Material. Throughout the Appendix, we use the following notation: The symbol C denotes a universal real constant which may take a different value on each occurrence. For $a, b \in \mathbb{R}$, we write $a_+ = \max\{0, a\}$ and $a \vee b = \max\{a, b\}$. For any set A , the symbol $|A|$ denotes the cardinality of A . The notation $X \stackrel{\mathcal{D}}{=} Y$ means that the two random variables X and Y have the same distribution. Finally, $f_0(\cdot)$ and $F_0(\cdot)$ denote the density and distribution function of the standard normal distribution, respectively.

Auxiliary results using strong approximation theory

The main purpose of this section is to prove that there is a version of the multiscale statistic $\widehat{\Phi}_T$ defined in (3.4) which is close to a Gaussian statistic whose distribution is known. More specifically, we prove the following result.

Proposition A.1. *Under the conditions of Theorem 3.1, there exist statistics $\widetilde{\Phi}_T$ for $T = 1, 2, \dots$ with the following two properties: (i) $\widetilde{\Phi}_T$ has the same distribution as $\widehat{\Phi}_T$ for any T , and (ii)*

$$|\widetilde{\Phi}_T - \Phi_T| = o_p\left(\frac{T^{1/q}}{\sqrt{Th_{\min}}} + \rho_T \sqrt{\log T}\right),$$

where Φ_T is a Gaussian statistic as defined in (3.3).

Proof of Proposition A.1. For the proof, we draw on strong approximation theory for stationary processes $\{\varepsilon_t\}$ that fulfill the conditions (C1)–(C3). By Theorem 2.1 and Corollary 2.1 in Berkes et al. (2014), the following strong approximation result holds true: On a richer probability space, there exist a standard Brownian motion \mathbb{B} and a sequence $\{\widetilde{\varepsilon}_t : t \in \mathbb{N}\}$ such that $[\widetilde{\varepsilon}_1, \dots, \widetilde{\varepsilon}_T] \stackrel{\mathcal{D}}{=} [\varepsilon_1, \dots, \varepsilon_T]$ for each T and

$$\max_{1 \leq t \leq T} \left| \sum_{s=1}^t \widetilde{\varepsilon}_s - \sigma \mathbb{B}(t) \right| = o(T^{1/q}) \quad \text{a.s.}, \quad (\text{A.1})$$

where $\sigma^2 = \sum_{k \in \mathbb{Z}} \text{Cov}(\varepsilon_0, \varepsilon_k)$ denotes the long-run error variance. To apply this result, we define

$$\widetilde{\Phi}_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\widetilde{\phi}_T(u,h)}{\widetilde{\sigma}} \right| - \lambda(h) \right\},$$

where $\widetilde{\phi}_T(u,h) = \sum_{t=1}^T w_{t,T}(u,h) \widetilde{\varepsilon}_t$ and $\widetilde{\sigma}^2$ is the same estimator as $\widehat{\sigma}^2$ with $Y_t = m(t/T) + \varepsilon_t$ replaced by $\widetilde{Y}_t = m(t/T) + \widetilde{\varepsilon}_t$ for $1 \leq t \leq T$. In addition, we let

$$\Phi_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\phi_T(u,h)}{\sigma} \right| - \lambda(h) \right\}$$

$$\Phi_T^\diamond = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\phi_T(u,h)}{\tilde{\sigma}} \right| - \lambda(h) \right\}$$

with $\phi_T(u,h) = \sum_{t=1}^T w_{t,T}(u,h) \sigma Z_t$ and $Z_t = \mathbb{B}(t) - \mathbb{B}(t-1)$. With this notation, we can write

$$|\tilde{\Phi}_T - \Phi_T| \leq |\tilde{\Phi}_T - \Phi_T^\diamond| + |\Phi_T^\diamond - \Phi_T| = |\tilde{\Phi}_T - \Phi_T^\diamond| + o_p(\rho_T \sqrt{\log T}), \quad (\text{A.2})$$

where the last equality follows by taking into account that $\phi_T(u,h) \sim N(0, \sigma^2)$ for all $(u,h) \in \mathcal{G}_T$, $|\mathcal{G}_T| = O(T^\theta)$ for some large but fixed constant θ and $\tilde{\sigma}^2 = \sigma^2 + o_p(\rho_T)$. Straightforward calculations yield that

$$|\tilde{\Phi}_T - \Phi_T^\diamond| \leq \tilde{\sigma}^{-1} \max_{(u,h) \in \mathcal{G}_T} |\tilde{\phi}_T(u,h) - \phi_T(u,h)|.$$

Using summation by parts, we further obtain that

$$\begin{aligned} |\tilde{\phi}_T(u,h) - \phi_T(u,h)| &\leq W_T(u,h) \max_{1 \leq t \leq T} \left| \sum_{s=1}^t \tilde{\varepsilon}_s - \sigma \sum_{s=1}^t \{\mathbb{B}(s) - \mathbb{B}(s-1)\} \right| \\ &= W_T(u,h) \max_{1 \leq t \leq T} \left| \sum_{s=1}^t \tilde{\varepsilon}_s - \sigma \mathbb{B}(t) \right|, \end{aligned}$$

where

$$W_T(u,h) = \sum_{t=1}^{T-1} |w_{t+1,T}(u,h) - w_{t,T}(u,h)| + |w_{T,T}(u,h)|.$$

Standard arguments show that $\max_{(u,h) \in \mathcal{G}_T} W_T(u,h) = O(1/\sqrt{Th_{\min}})$. Applying the strong approximation result (A.1), we can thus infer that

$$\begin{aligned} |\tilde{\Phi}_T - \Phi_T^\diamond| &\leq \tilde{\sigma}^{-1} \max_{(u,h) \in \mathcal{G}_T} |\tilde{\phi}_T(u,h) - \phi_T(u,h)| \\ &\leq \tilde{\sigma}^{-1} \max_{(u,h) \in \mathcal{G}_T} W_T(u,h) \max_{1 \leq t \leq T} \left| \sum_{s=1}^t \tilde{\varepsilon}_s - \sigma \mathbb{B}(t) \right| \\ &= o_p\left(\frac{T^{1/q}}{\sqrt{Th_{\min}}}\right). \end{aligned} \quad (\text{A.3})$$

Plugging (A.3) into (A.2) completes the proof. \square

Auxiliary results using anti-concentration bounds

In this section, we establish some properties of the Gaussian statistic Φ_T defined in (3.3). We in particular show that Φ_T does not concentrate too strongly in small regions of the form $[x - \delta_T, x + \delta_T]$ with δ_T converging to zero.

Proposition A.2. *Under the conditions of Theorem 3.1, it holds that*

$$\sup_{x \in \mathbb{R}} \mathbb{P}(|\Phi_T - x| \leq \delta_T) = o(1),$$

where $\delta_T = T^{1/q}/\sqrt{Th_{\min}} + \rho_T\sqrt{\log T}$.

Proof of Proposition A.2. The main technical tool for proving Proposition A.2 are anti-concentration bounds for Gaussian random vectors. The following proposition slightly generalizes anti-concentration results derived in Chernozhukov et al. (2015), in particular Theorem 3 therein.

Proposition A.3. *Let $(X_1, \dots, X_p)^\top$ be a Gaussian random vector in \mathbb{R}^p with $\mathbb{E}[X_j] = \mu_j$ and $\text{Var}(X_j) = \sigma_j^2 > 0$ for $1 \leq j \leq p$. Define $\bar{\mu} = \max_{1 \leq j \leq p} |\mu_j|$ together with $\underline{\sigma} = \min_{1 \leq j \leq p} \sigma_j$ and $\bar{\sigma} = \max_{1 \leq j \leq p} \sigma_j$. Moreover, set $a_p = \mathbb{E}[\max_{1 \leq j \leq p} (X_j - \mu_j)/\sigma_j]$ and $b_p = \mathbb{E}[\max_{1 \leq j \leq p} (X_j - \mu_j)]$. For every $\delta > 0$, it holds that*

$$\sup_{x \in \mathbb{R}} \mathbb{P}\left(\left|\max_{1 \leq j \leq p} X_j - x\right| \leq \delta\right) \leq C\delta\{\bar{\mu} + a_p + b_p + \sqrt{1 \vee \log(\underline{\sigma}/\delta)}\},$$

where $C > 0$ depends only on $\underline{\sigma}$ and $\bar{\sigma}$.

The proof of Proposition A.3 is provided in the Supplementary Material. To apply Proposition A.3 to our setting at hand, we introduce the following notation: We write $x = (u, h)$ along with $\mathcal{G}_T = \{x : x \in \mathcal{G}_T\} = \{x_1, \dots, x_p\}$, where $p := |\mathcal{G}_T| \leq O(T^\theta)$ for some large but fixed $\theta > 0$ by our assumptions. Moreover, for $j = 1, \dots, p$, we set

$$\begin{aligned} X_{2j-1} &= \frac{\phi_T(x_{j1}, x_{j2})}{\sigma} - \lambda(x_{j2}) \\ X_{2j} &= -\frac{\phi_T(x_{j1}, x_{j2})}{\sigma} - \lambda(x_{j2}) \end{aligned}$$

with $x_j = (x_{j1}, x_{j2})$. This notation allows us to write

$$\Phi_T = \max_{1 \leq j \leq 2p} X_j,$$

where $(X_1, \dots, X_{2p})^\top$ is a Gaussian random vector with the following properties: (i) $\mu_j := \mathbb{E}[X_j] = -\lambda(x_{j2})$ and thus $\bar{\mu} = \max_{1 \leq j \leq 2p} |\mu_j| \leq C\sqrt{\log T}$, and (ii) $\sigma_j^2 := \text{Var}(X_j) = 1$ for all j . Since $\sigma_j = 1$ for all j , it holds that $a_{2p} = b_{2p}$. Moreover, as the variables $(X_j - \mu_j)/\sigma_j$ are standard normal, we have that $a_{2p} = b_{2p} \leq \sqrt{2\log(2p)} \leq C\sqrt{\log T}$. With this notation at hand, we can apply Proposition A.3 to obtain that

$$\sup_{x \in \mathbb{R}} \mathbb{P}(|\Phi_T - x| \leq \delta_T) \leq C\delta_T \left[\sqrt{\log T} + \sqrt{\log(1/\delta_T)} \right] = o(1)$$

with $\delta_T = T^{1/q}/\sqrt{Th_{\min}} + \rho_T\sqrt{\log T}$, which is the statement of Proposition A.2. \square

Proof of Theorem 3.1

To prove Theorem 3.1, we make use of the two auxiliary results derived above. By Proposition A.1, there exist statistics $\tilde{\Phi}_T$ for $T = 1, 2, \dots$ which are distributed as $\hat{\Phi}_T$ for any $T \geq 1$ and which have the property that

$$|\tilde{\Phi}_T - \Phi_T| = o_p\left(\frac{T^{1/q}}{\sqrt{Th_{\min}}} + \rho_T \sqrt{\log T}\right), \quad (\text{A.4})$$

where Φ_T is a Gaussian statistic as defined in (3.3). The approximation result (A.4) allows us to replace the multiscale statistic $\hat{\Phi}_T$ by an identically distributed version $\tilde{\Phi}_T$ which is close to the Gaussian statistic Φ_T . In the next step, we show that

$$\sup_{x \in \mathbb{R}} |\mathbb{P}(\tilde{\Phi}_T \leq x) - \mathbb{P}(\Phi_T \leq x)| = o(1), \quad (\text{A.5})$$

which immediately implies the statement of Theorem 3.1. For the proof of (A.5), we use the following simple lemma:

Lemma A.4. *Let V_T and W_T be real-valued random variables for $T = 1, 2, \dots$ such that $V_T - W_T = o_p(\delta_T)$ with some $\delta_T = o(1)$. If*

$$\sup_{x \in \mathbb{R}} \mathbb{P}(|V_T - x| \leq \delta_T) = o(1), \quad (\text{A.6})$$

then

$$\sup_{x \in \mathbb{R}} |\mathbb{P}(V_T \leq x) - \mathbb{P}(W_T \leq x)| = o(1). \quad (\text{A.7})$$

The statement of Lemma A.4 can be summarized as follows: If W_T can be approximated by V_T in the sense that $V_T - W_T = o_p(\delta_T)$ and if V_T does not concentrate too strongly in small regions of the form $[x - \delta_T, x + \delta_T]$ as assumed in (A.6), then the distribution of W_T can be approximated by that of V_T in the sense of (A.7).

Proof of Lemma A.4. It holds that

$$\begin{aligned} & |\mathbb{P}(V_T \leq x) - \mathbb{P}(W_T \leq x)| \\ &= |\mathbb{E}[1(V_T \leq x) - 1(W_T \leq x)]| \\ &\leq |\mathbb{E}[\{1(V_T \leq x) - 1(W_T \leq x)\}1(|V_T - W_T| \leq \delta_T)]| + |\mathbb{E}[1(|V_T - W_T| > \delta_T)]| \\ &\leq \mathbb{E}[1(|V_T - x| \leq \delta_T, |V_T - W_T| \leq \delta_T)] + o(1) \\ &\leq \mathbb{P}(|V_T - x| \leq \delta_T) + o(1). \end{aligned} \quad \square$$

We now apply this lemma with $V_T = \Phi_T$, $W_T = \tilde{\Phi}_T$ and $\delta_T = T^{1/q}/\sqrt{Th_{\min}} + \rho_T \sqrt{\log T}$: From (A.4), we already know that $\tilde{\Phi}_T - \Phi_T = o_p(\delta_T)$. Moreover, by Proposition A.2, it holds that

$$\sup_{x \in \mathbb{R}} \mathbb{P}(|\Phi_T - x| \leq \delta_T) = o(1). \quad (\text{A.8})$$

Hence, the conditions of Lemma A.4 are satisfied. Applying the lemma, we obtain (A.5), which completes the proof of Theorem 3.1.

Proof of Proposition 3.3

Write $\widehat{\psi}_T(u, h) = \widehat{\psi}_T^A(u, h) + \widehat{\psi}_T^B(u, h)$ with $\widehat{\psi}_T^A(u, h) = \sum_{t=1}^T w_{t,T}(u, h)\varepsilon_t$ and $\widehat{\psi}_T^B(u, h) = \sum_{t=1}^T w_{t,T}(u, h)m_T(\frac{t}{T})$. By assumption, there exists $(u_0, h_0) \in \mathcal{G}_T$ with $[u_0 - h_0, u_0 + h_0] \subseteq [0, 1]$ such that $m_T(w) \geq c_T \sqrt{\log T / (Th_0)}$ for all $w \in [u_0 - h_0, u_0 + h_0]$. Since the kernel K is symmetric and $u_0 = t/T$ for some t , it holds that $S_{T,1}(u_0, h_0) = 0$ and thus

$$w_{t,T}(u_0, h_0) = K\left(\frac{\frac{t}{T} - u_0}{h_0}\right) / \left\{ \sum_{t=1}^T K^2\left(\frac{\frac{t}{T} - u_0}{h_0}\right) \right\}^{1/2} \geq 0.$$

Together with the assumption that $m_T(w) \geq c_T \sqrt{\log T / (Th_0)}$ for all $w \in [u_0 - h_0, u_0 + h_0]$, this implies that

$$\widehat{\psi}_T^B(u_0, h_0) \geq c_T \sqrt{\frac{\log T}{Th_0}} \sum_{t=1}^T w_{t,T}(u_0, h_0). \quad (\text{A.9})$$

Standard calculations exploiting the Lipschitz continuity of the kernel K show that for any $(u, h) \in \mathcal{G}_T$ and any given natural number ℓ ,

$$\left| \frac{1}{Th} \sum_{t=1}^T K\left(\frac{\frac{t}{T} - u}{h}\right) \left(\frac{\frac{t}{T} - u}{h}\right)^\ell - \int_0^1 \frac{1}{h} K\left(\frac{w - u}{h}\right) \left(\frac{w - u}{h}\right)^\ell dw \right| \leq \frac{C}{Th}, \quad (\text{A.10})$$

where the constant C does not depend on u , h and T . With the help of (A.10), we obtain that for any $(u, h) \in \mathcal{G}_T$ with $[u - h, u + h] \subseteq [0, 1]$,

$$\left| \sum_{t=1}^T w_{t,T}(u, h) - \frac{\sqrt{Th}}{\kappa} \right| \leq \frac{C}{\sqrt{Th}}, \quad (\text{A.11})$$

where $\kappa = (\int K^2(\varphi)d\varphi)^{1/2}$ and the constant C does once again not depend on u , h and T . From (A.11), it follows that $\sum_{t=1}^T w_{t,T}(u, h) \geq \sqrt{Th}/(2\kappa)$ for sufficiently large T and any $(u, h) \in \mathcal{G}_T$ with $[u - h, u + h] \subseteq [0, 1]$. This together with (A.9) allows us to infer that

$$\widehat{\psi}_T^B(u_0, h_0) \geq \frac{c_T \sqrt{\log T}}{2\kappa} \quad (\text{A.12})$$

for sufficiently large T . Moreover, arguments very similar to those for the proof of Proposition A.1 yield that

$$\max_{(u,h) \in \mathcal{G}_T} |\widehat{\psi}_T^A(u, h)| = O_p(\sqrt{\log T}). \quad (\text{A.13})$$

With the help of (A.12), (A.13) and the fact that $\lambda(h) \leq \lambda(h_{\min}) \leq C\sqrt{\log T}$, we finally arrive at

$$\begin{aligned}\widehat{\Psi}_T &\geq \max_{(u,h) \in \mathcal{G}_T} \frac{|\widehat{\psi}_T^B(u,h)|}{\widehat{\sigma}} - \max_{(u,h) \in \mathcal{G}_T} \left\{ \frac{|\widehat{\psi}_T^A(u,h)|}{\widehat{\sigma}} + \lambda(h) \right\} \\ &= \max_{(u,h) \in \mathcal{G}_T} \frac{|\widehat{\psi}_T^B(u,h)|}{\widehat{\sigma}} + O_p(\sqrt{\log T}) \\ &\geq \frac{c_T \sqrt{\log T}}{2\kappa \widehat{\sigma}} + O_p(\sqrt{\log T}).\end{aligned}\tag{A.14}$$

Since $q_T(\alpha) = O(\sqrt{\log T})$ for any fixed $\alpha \in (0, 1)$, (A.14) immediately implies that $\mathbb{P}(\widehat{\Psi}_T \leq q_T(\alpha)) = o(1)$.

Proof of Proposition 3.4

The statement of Proposition 3.4 is a consequence of the following observation: For all $(u, h) \in \mathcal{G}_T$ with

$$\left| \frac{\widehat{\psi}_T(u, h) - \mathbb{E}\widehat{\psi}_T(u, h)}{\widehat{\sigma}} \right| - \lambda(h) \leq q_T(\alpha) \quad \text{and} \quad \left| \frac{\widehat{\psi}_T(u, h)}{\widehat{\sigma}} \right| - \lambda(h) > q_T(\alpha),$$

it holds that $\mathbb{E}[\widehat{\psi}_T(u, h)] \neq 0$, which in turn implies that $m(v) \neq 0$ for some $v \in I_{u,h}$. From this observation, we can infer the following: On the event

$$\{\widehat{\Phi}_T \leq q_T(\alpha)\} = \left\{ \max_{(u,h) \in \mathcal{G}_T} \left(\left| \frac{\widehat{\psi}_T(u, h) - \mathbb{E}\widehat{\psi}_T(u, h)}{\widehat{\sigma}} \right| - \lambda(h) \right) \leq q_T(\alpha) \right\},$$

it holds that for all $(u, h) \in \mathcal{A}_T$, $m(v) \neq 0$ for some $v \in I_{u,h}$. Hence, we obtain that

$$\{\widehat{\Phi}_T \leq q_T(\alpha)\} \subseteq E_T.$$

As a result, we arrive at

$$\mathbb{P}(E_T) \geq \mathbb{P}(\widehat{\Phi}_T \leq q_T(\alpha)) = (1 - \alpha) + o(1),$$

where the last equality holds by Theorem 3.1.

References

- BENNER, T. C. (1999). Central england temperatures: long-term variability and teleconnections. *International Journal of Climatology*, **19** 391–403.
- BERKES, I., LIU, W. and WU, W. B. (2014). Komlós-Major-Tusnády approximation under dependence. *Annals of Probability*, **42** 794–817.
- CHAUDHURI, P. and MARRON, J. S. (1999). SiZer for the exploration of structures in curves. *Journal of the American Statistical Association*, **94** 807–823.
- CHAUDHURI, P. and MARRON, J. S. (2000). Scale space view of curve estimation. *Annals of Statistics*, **28** 408–428.
- CHEN, L. and WU, W. B. (2018). Testing for trends in high-dimensional time series. *Forthcoming in Journal of the American Statistical Association*.
- CHERNOZHUKOV, V., CHETVERIKOV, D. and KATO, K. (2015). Comparison and anti-concentration bounds for maxima of Gaussian random vectors. *Probability Theory and Related Fields*, **162** 47–70.
- DEGRAS, D., XU, Z., ZHANG, T. and WU, W. B. (2012). Testing for parallelism among trends in multiple time series. *IEEE Transactions on Signal Processing*, **60** 1087–1097.
- DÜMBGEN, L. (2002). Application of local rank tests to nonparametric regression. *Journal of Nonparametric Statistics*, **14** 511–537.
- DÜMBGEN, L. and SPOKOINY, V. G. (2001). Multiscale testing of qualitative hypotheses. *Annals of Statistics*, **29** 124–152.
- HALL, P. and HECKMAN, N. E. (2000). Testing for monotonicity of a regression mean by calibrating for linear functions. *Annals of Statistics*, **28** 20–39.
- HALL, P. and VAN KEILEGOM, I. (2003). Using difference-based methods for inference in nonparametric regression with time series errors. *Journal of the Royal Statistical Society: Series B*, **65** 443–456.
- HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. (2009). *The Elements of Statistical Learning*. New York, Springer.
- HERRMANN, E., GASSER, T. and KNEIP, A. (1992). Choice of bandwidth for kernel regression when residuals are correlated. *Biometrika*, **79** 783–795.
- LYUBCHICH, V. and GEL, Y. R. (2016). A local factor nonparametric test for trend synchronism in multiple time series. *Journal of Multivariate Analysis*, **150** 91–104.
- MUDELSEE, M. (2010). *Climate time series analysis: classical statistical and bootstrap methods*. New York, Springer.
- MÜLLER, H.-G. and STADTMÜLLER, U. (1988). Detecting dependencies in smooth regression models. *Biometrika*, **75** 639–650.
- PARK, C., MARRON, J. S. and RONDONOTTI, V. (2004). Dependent SiZer: goodness-of-fit

- tests for time series models. *Journal of Applied Statistics*, **31** 999–1017.
- PARK, C., VAUGHAN, A., HANNIG, J. and KANG, K.-H. (2009). SiZer analysis for the comparison of time series. *Journal of Statistical Planning and Inference*, **139** 3974–3988.
- PARKER, D. E., LEGG, T. P. and FOLLAND, C. K. (1992). A new daily central england temperature series, 1772-1991. *International Journal of Climatology*, **12** 317–342.
- RAHMSTORF, S., FOSTER, G. and CAHILL, N. (2017). Global temperature evolution: recent trends and some pitfalls. *Environmental Research Letters*, **12**.
- RONDONOTTI, V., MARRON, J. S. and PARK, C. (2007). SiZer for time series: a new approach to the analysis of trends. *Electronic Journal of Statistics*, **1** 268–289.
- TECUAPETLA-GÓMEZ, I. and MUNK, A. (2017). Autocovariance estimation in regression with a discontinuous signal and m -dependent errors: a difference-based approach. *Scandinavian Journal of Statistics*, **44** 346–368.
- VOGELSANG, T. J. and FRANSES, P. H. (2005). Testing for common deterministic trend slopes. *Journal of Econometrics*, **126** 1–24.
- WU, W. B. (2005). Nonlinear system theory: another look at dependence. *Proc. Natn. Acad. Sci. USA*, **102** 14150–14154.
- WU, W. B. and SHAO, X. (2004). Limit theorems for iterated random functions. *Journal of Applied Probability* 425–436.