

Multiscale Inference and Long-Run Variance Estimation in Nonparametric Regression with Time Series Errors

Marina Khismatullina¹
University of Bonn

Michael Vogt²
University of Bonn

July 23, 2019

In this paper, we develop new multiscale methods to test qualitative hypotheses about the function m in the nonparametric regression model $Y_{t,T} = m(t/T) + \varepsilon_t$ with time series errors ε_t . In time series applications, m represents a nonparametric time trend. Practitioners are often interested in whether the trend m has certain shape properties. For example, they would like to know whether m is constant or whether it is increasing/decreasing in certain time intervals. Our multiscale methods allow to test for such shape properties of the trend m . In order to perform the methods, we require an estimator of the long-run error variance $\sigma^2 = \sum_{\ell=-\infty}^{\infty} \text{Cov}(\varepsilon_0, \varepsilon_\ell)$. We propose a new difference-based estimator of σ^2 for the case that $\{\varepsilon_t\}$ belongs to the class of $\text{AR}(\infty)$ processes. In the technical part of the paper, we derive asymptotic theory for the proposed multiscale test and the estimator of the long-run error variance. The theory is complemented by a simulation study and an empirical application to climate data.

Key words: Multiscale statistics; long-run variance; 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 a wide range of situations, practitioners are particularly interested in certain shape properties of the trend. They raise questions such as the following: Does the observed time series have a trend at all? If so, is the trend increasing/decreasing in certain time intervals? Can one identify the intervals of increase/decrease? 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 learning about

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.

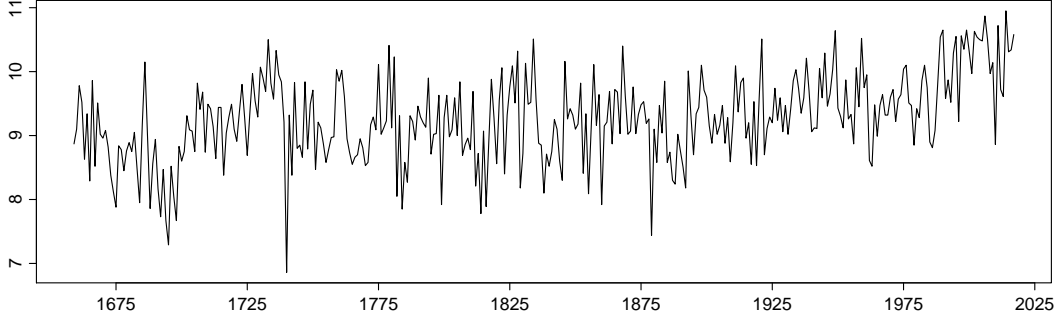


Figure 1: Yearly mean temperature in Central England from 1659 to 2017 measured in $^{\circ}\text{C}$.

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 this paper, we develop new methods to test for certain shape properties of a nonparametric time trend. We in particular construct a multiscale test which allows to identify local increases/decreases of the trend function. We develop our test in the context of the following model setting: We observe a time series $\{Y_{t,T} : 1 \leq t \leq T\}$ of the form

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

for $1 \leq t \leq T$, where $m : [0, 1] \rightarrow \mathbb{R}$ is an unknown nonparametric regression function and the error terms ε_t form a stationary time series process with $\mathbb{E}[\varepsilon_t] = 0$. In a time series context, the design points t/T represent the time points of observation and m is a nonparametric time trend. As usual in nonparametric regression, we let the function m depend on rescaled time t/T rather than on real time t . A detailed description of model (1.1) is provided in Section 2.

Our multiscale test is developed step by step in Section 3. Roughly speaking, the procedure can be outlined as follows: Let $H_0(u, h)$ be the hypothesis that m is constant in the time window $[u - h, u + h] \subseteq [0, 1]$, where u is the midpoint and $2h$ the size of the window. In a first step, we set up a test statistic $\widehat{\varphi}_T(u, h)$ for the hypothesis $H_0(u, h)$. In a second step, we aggregate the statistics $\widehat{\varphi}_T(u, h)$ for a large number of different time windows $[u - h, u + h]$. We thereby construct a multiscale statistic which allows to test the hypothesis $H_0(u, h)$ simultaneously for many time windows $[u - h, u + h]$. In the technical part of the paper, we derive the theoretical properties of the resulting multiscale test. To do so, we come up with a proof strategy which combines strong approximation results for dependent processes with anti-concentration bounds for Gaussian random vectors. This strategy is of interest in itself and may be applied to other multiscale test problems for dependent data. As shown by our theoretical analysis, our multiscale test is a rigorous level- α -test of the overall null hypothesis

H_0 that $H_0(u, h)$ is simultaneously fulfilled for all time windows $[u - h, u + h]$ under consideration. Moreover, for a given significance level $\alpha \in (0, 1)$, the test allows to make simultaneous confidence statements of the following form: We can claim, with statistical confidence $1 - \alpha$, that there is an increase/decrease in the trend m on all time windows $[u - h, u + h]$ for which the hypothesis $H_0(u, h)$ is rejected. Hence, the test allows to identify, with a pre-specified statistical confidence, time intervals where the trend m is increasing/decreasing.

For independent data, multiscale tests have been developed in a variety of different contexts in recent years. In the regression context, Chaudhuri and Marron (1999, 2000) introduced the so-called SiZer method which has been extended in various directions; see e.g. Hannig and Marron (2006) where a refined distribution theory for SiZer is derived. Hall and Heckman (2000) constructed a multiscale test on monotonicity of a regression function. Dümbgen and Spokoiny (2001) developed a multiscale approach which works with additively corrected supremum statistics and derived theoretical results in the context of a continuous Gaussian white noise model. Rank-based multiscale tests for nonparametric regression were proposed in Dümbgen (2002) and Rohde (2008). More recently, Proksch et al. (2018) have constructed multiscale tests for inverse regression models. In the context of density estimation, multiscale tests have been investigated in Dümbgen and Walther (2008), Rufibach and Walther (2010), Schmidt-Hieber et al. (2013) and Eckle et al. (2017) among others.

Whereas a large number of multiscale tests for independent data have been developed in recent years, multiscale tests for dependent data are much rarer. Most notably, there are some extensions of the SiZer approach to a time series context. Park et al. (2004) and Rondonotti et al. (2007) have introduced SiZer methods for dependent data which can be used to find local increases/decreases of a trend and which may thus be regarded as an alternative to our multiscale test. However, these SiZer methods are mainly designed for data exploration rather than for rigorous statistical inference. Our multiscale method, in contrast, is a rigorous level- α -test of the hypothesis H_0 which allows to make simultaneous confidence statements about the time intervals where the trend m is increasing/decreasing. Some theoretical results for dependent SiZer methods have been derived in Park et al. (2009), but only under a quite severe restriction: Only time windows $[u - h, u + h]$ with window sizes or scales h are taken into account that remain bounded away from zero as the sample size T grows. Scales h that converge to zero as T increases are excluded. This effectively means that only large time windows $[u - h, u + h]$ are taken into consideration. Our theory, in contrast, allows to simultaneously consider scales h of fixed size and scales h that converge to zero at various different rates. We are thus able to take into account time windows of many different sizes. In Section 3.4, we compare our approach to SiZer methods for dependent data in more detail.

Our multiscale approach is also related to Wavelet-based methods: Similar to the latter,

it takes into account different locations u and resolution levels or scales h simultaneously. However, while our multiscale approach is designed to test for local increases/decreases of a nonparametric trend, Wavelet methods are commonly used for other purposes. Among other things, they are employed for estimating/reconstructing nonparametric regression curves [see e.g. Donoho et al. (1995) or Von Sachs and MacGibbon (2000)] and for change point detection [see e.g. Cho and Fryzlewicz (2012)].

The test statistic of our multiscale method depends on the long-run error variance $\sigma^2 = \sum_{\ell=-\infty}^{\infty} \text{Cov}(\varepsilon_0, \varepsilon_\ell)$, which is usually unknown in practice. To carry out our multiscale test, we thus require an estimator of σ^2 . Indeed, such an estimator is required for virtually all inferential procedures in the context of model (1.1). Hence, the problem of estimating σ^2 in model (1.1) is of broader interest and has received a lot of attention in the literature; see Müller and Stadtmüller (1988), Herrmann et al. (1992) and Hall and Van Keilegom (2003) among many others. In Section 4, we introduce a new difference-based estimator of σ^2 for the case that $\{\varepsilon_t\}$ belongs to the class of $\text{AR}(\infty)$ processes. This estimator improves on existing methods in several respects.

The methodological and theoretical analysis of the paper is complemented by a simulation study in Section 5 and two empirical applications in Section 6. In the simulation study, we examine the finite sample properties of our multiscale test and compare it to the dependent SiZer methods introduced in Park et al. (2004) and Rondonotti et al. (2007). Moreover, we investigate the small sample performance of our estimator of σ^2 and compare it to the estimator of Hall and Van Keilegom (2003). In Section 6, we use our methods to analyse the temperature data from Figure 1 as well as a sample of global temperature data.

2 The model

We now describe the model setting in detail which was briefly outlined in the Introduction. We observe a time series $\{Y_{t,T} : 1 \leq t \leq T\}$ of length T which satisfies the nonparametric regression equation

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

for $1 \leq t \leq T$. Here, m is an unknown nonparametric 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.

3 The multiscale test

In this section, we introduce our multiscale method to test for local increases/decreases of the trend function m and analyse its theoretical properties. We assume throughout that m is continuously differentiable on $[0, 1]$. The test problem under consideration can be formulated as follows: Let $H_0(u, h)$ be the hypothesis that m is constant on the interval $[u - h, u + h]$. Since m is continuously differentiable, $H_0(u, h)$ can be reformulated as

$$H_0(u, h) : m'(w) = 0 \text{ for all } w \in [u - h, u + h],$$

where m' is the first derivative of m . We want to test the hypothesis $H_0(u, h)$ not only for a single interval $[u - h, u + h]$ but simultaneously for many different intervals. The overall null hypothesis is thus given by

$$H_0 : \text{The hypothesis } H_0(u, h) \text{ holds true for all } (u, h) \in \mathcal{G}_T,$$

where \mathcal{G}_T is some large set of points (u, h) . The details on the set \mathcal{G}_T are discussed at the end of Section 3.1 below. Note that \mathcal{G}_T in general depends on the sample size T , implying that the null hypothesis $H_0 = H_{0,T}$ depends on T as well. We thus consider a sequence of null hypotheses $\{H_{0,T} : T = 1, 2, \dots\}$ as T increases. For simplicity of notation, we however suppress the dependence of H_0 on T . In Sections 3.1 and 3.2, we step by step construct the multiscale test of the hypothesis H_0 . The theoretical properties of the test are analysed in Section 3.3.

3.1 Construction of the multiscale statistic

We first construct a test statistic for the hypothesis $H_0(u, h)$, where $[u - h, u + h]$ is a given interval. To do so, we consider the kernel average

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

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

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

where

$$\Lambda_{t,T}(u, h) = K\left(\frac{\frac{t}{T} - u}{h}\right) \left[S_{T,0}(u, h) \left(\frac{\frac{t}{T} - u}{h}\right) - S_{T,1}(u, h) \right],$$

$S_{T,\ell}(u, h) = (Th)^{-1} \sum_{t=1}^T K\left(\frac{\frac{t}{T} - u}{h}\right) \left(\frac{\frac{t}{T} - u}{h}\right)^\ell$ 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 $C > 0$.

The kernel average $\widehat{\psi}_T(u, h)$ is nothing else than a rescaled local linear estimator of the derivative $m'(u)$ with bandwidth h .³

A test statistic for the hypothesis $H_0(u, h)$ is given by the normalized kernel average $\widehat{\psi}_T(u, h)/\widehat{\sigma}$, where $\widehat{\sigma}^2$ is an estimator of the long-run variance $\sigma^2 = \sum_{\ell=-\infty}^{\infty} \text{Cov}(\varepsilon_0, \varepsilon_\ell)$ of the error process $\{\varepsilon_t\}$. The problem of estimating σ^2 is discussed in detail in Section 4. For the time being, we suppose that $\widehat{\sigma}^2$ is an estimator with reasonable theoretical properties. Specifically, we assume that $\widehat{\sigma}^2 = \sigma^2 + o_p(\rho_T)$ with $\rho_T = o(1/\log T)$. This is

³Alternatively to the local linear weights defined in (3.1), we could also work with the weights $w_{t,T}(u, h) = K'(h^{-1}[u - t/T])/\{\sum_{t=1}^T K'(h^{-1}[u - t/T])^2\}^{1/2}$, where the kernel function K is assumed to be differentiable and K' is its derivative. We however prefer to use local linear weights as these have superior theoretical properties at the boundary.

a fairly weak condition which is in particular satisfied by the estimator of σ^2 analysed in Section 4. The kernel weights $w_{t,T}(u, h)$ are chosen such that in the case of independent errors ε_t , $\text{Var}(\hat{\psi}_T(u, h)) = \sigma^2$ for any location u and bandwidth h , where the long-run error variance σ^2 simplifies to $\sigma^2 = \text{Var}(\varepsilon_t)$. In the more general case that the error terms satisfy the weak dependence conditions from Section 2, $\text{Var}(\hat{\psi}_T(u, h)) = \sigma^2 + o(1)$ for any u and h under consideration. Hence, for sufficiently large sample sizes T , the test statistic $\hat{\psi}_T(u, h)/\hat{\sigma}$ has approximately unit variance.

We now combine the test statistics $\hat{\psi}_T(u, h)/\hat{\sigma}$ for a wide range of different locations u and bandwidths or scales h . There are different ways to do so, leading to different types of multiscale statistics. Our multiscale statistic 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 given 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 calibrate the statistics $\hat{\psi}_T(u, h)/\hat{\sigma}$ that correspond to the bandwidth h by subtracting the additive correction term $\lambda(h)$. This approach was pioneered by Dümbgen and Spokoiny (2001) and has been used in numerous other studies since then; see e.g. Dümbgen (2002), Rohde (2008), Dümbgen and Walther (2008), Rufibach and Walther (2010), Schmidt-Hieber et al. (2013) and Eckle et al. (2017).

To see the heuristic idea behind the additive correction $\lambda(h)$, consider for a moment the uncorrected statistic

$$\hat{\Psi}_{T, \text{uncorrected}} = \max_{(u,h) \in \mathcal{G}_T} \left| \frac{\hat{\psi}_T(u, h)}{\hat{\sigma}} \right| \quad (3.3)$$

and suppose that the hypothesis $H_0(u, h)$ is true for all $(u, h) \in \mathcal{G}_T$. For simplicity, assume that the errors ε_t are i.i.d. normally distributed and neglect the estimation error in $\hat{\sigma}$, that is, set $\hat{\sigma} = \sigma$. 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

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

Under our simplifying assumptions, the statistics $\hat{\psi}_T(u_k, h_\ell)/\sigma$ 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 $\max_k \hat{\psi}_T(u_k, h_\ell)/\sigma$ is approximately of size $\lambda(h_\ell)$ for small bandwidths h_ℓ . As $\lambda(h) \rightarrow \infty$ for $h \rightarrow 0$, this implies that $\max_k \hat{\psi}_T(u_k, h_\ell)/\sigma$ 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 $h_{\max} < 1/2$.

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 compared to the sample size T . For example, we may work with the set

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

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. The lower bound on h_{\min} depends on the parameter q defined in (C2) which specifies the number of existing moments for the error terms ε_t . As one can see, we can choose h_{\min} to be of the order $T^{-1/2}$ for any $q > 4$. Hence, we can let h_{\min} converge to 0 very quickly even if only the first few moments of the error terms ε_t exist. If all moments exist (i.e. $q = \infty$), h_{\min} may converge to 0 almost as quickly as $T^{-1} \log T$. Furthermore, the maximal bandwidth h_{\max} is not even required to converge to 0, which implies that we can pick it very large.

Remark 3.1. *The above construction of the multiscale statistic can be easily adapted to hypotheses other than H_0 . To do so, one simply needs to replace the kernel weights $w_{t,T}(u, h)$ defined in (3.1) by appropriate versions which are suited to test the hypothesis of interest. For example, if one wants to test for local convexity/concavity of m , one may define the kernel weights $w_{t,T}(u, h)$ such that the kernel average $\widehat{\psi}_T(u, h)$ is a (rescaled) estimator of the second derivative of m at the location u with bandwidth h .*

3.2 The test procedure

In order to formulate a test for the null hypothesis H_0 , we still need to specify a critical value. To do so, we define the statistic

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

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 $\widehat{\Psi}_T$ 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 is now defined as follows: For a given significance level $\alpha \in (0,1)$, we reject the overall null hypothesis H_0 if $\widehat{\Psi}_T > q_T(\alpha)$. In particular, for any $(u,h) \in \mathcal{G}_T$, we reject $H_0(u,h)$ if the (corrected) test statistic $|\widehat{\psi}_T(u,h)/\widehat{\sigma}| - \lambda(h)$ lies above the critical value $q_T(\alpha)$, that is, if $|\widehat{\psi}_T(u,h)/\widehat{\sigma}| > q_T(\alpha) + \lambda(h)$.

3.3 The theoretical properties of the test

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

$$\widehat{\Phi}_T = \max_{(u,h) \in \mathcal{G}_T} \left\{ \left| \frac{\widehat{\phi}_T(u,h)}{\widehat{\sigma}} \right| - \lambda(h) \right\} \quad (3.5)$$

with $\widehat{\phi}_T(u,h) = \widehat{\psi}_T(u,h) - \mathbb{E}[\widehat{\psi}_T(u,h)] = \sum_{t=1}^T w_{t,T}(u,h) \varepsilon_t$. The following result is central to the theoretical analysis of our multiscale test. According to it, the (known) quantile $q_T(\alpha)$ of the Gaussian statistic Φ_T defined in Section 3.2 can be used as a proxy for the $(1-\alpha)$ -quantile of the multiscale statistic $\widehat{\Phi}_T$.

Theorem 3.1. *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*

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

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

$$|\widetilde{\Phi}_T - \Phi_T| = o_p(\delta_T), \quad (3.6)$$

where $\delta_T = o(1)$ and the Gaussian statistic Φ_T is defined in Section 3.2. We thus replace the statistic $\widehat{\Phi}_T$ 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

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

which immediately implies the statement of Theorem 3.1. Importantly, the convergence result (3.6) is not sufficient for establishing (3.7). 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 (3.7). For (3.7) to hold, we additionally require the distribution of Φ_T to have some sort of continuity property. Specifically, we prove that

$$\sup_{x \in \mathbb{R}} \mathbb{P}(|\Phi_T - x| \leq \delta_T) = o(1), \quad (3.8)$$

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 (3.8) are anti-concentration results for Gaussian random vectors as derived in Chernozhukov et al. (2015). The claim (3.7) can be proven by using (3.6) together with (3.8), which in turn yields Theorem 3.1.

The main idea of our proof strategy is to combine strong approximation theory with anti-concentration bounds for Gaussian random vectors to show that the quantiles of the multiscale statistic $\hat{\Phi}_T$ can be proxied by those of a Gaussian analogue. This strategy is quite general in nature and may be applied to other multiscale problems for dependent data. Strong approximation theory has also been used to investigate multiscale tests for independent data; see e.g. Schmidt-Hieber et al. (2013). However, it has not been combined with anti-concentration results to approximate the quantiles of the multiscale statistic. As an alternative to strong approximation theory, Eckle et al. (2017) and Proksch et al. (2018) have recently used Gaussian approximation results derived in Chernozhukov et al. (2014, 2017) to analyse multiscale tests for independent data. Even though it might be possible to adapt these techniques to the case of dependent data, this is not trivial at all as part of the technical arguments and the Gaussian approximation tools strongly rely on the assumption of independence.

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

Proposition 3.1. *Let the conditions of Theorem 3.1 be satisfied. Under the null hypothesis H_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

$$m'_T(w) \geq c_T \sqrt{\frac{\log T}{Th^3}} \quad \text{for all } w \in [u - h, u + h], \quad (3.9)$$

where $\{c_T\}$ is any sequence of positive numbers with $c_T \rightarrow \infty$. Alternatively to (3.9), we may also assume that $-m'_T(w) \geq c_T \sqrt{\log T/(Th^3)}$ for all $w \in [u - h, u + h]$.

Proposition 3.2. *Let the conditions of Theorem 3.1 be satisfied and consider any sequence of functions m_T with the property (3.9). Then*

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

According to Proposition 3.2, our test has asymptotic power 1 against local alternatives of the form (3.9). The proof can be found in the Supplementary Material.

The next result formally shows that we can make simultaneous confidence statements about the time intervals where the trend m is increasing/decreasing. To formulate it, we define

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

where

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

The object Π_T^\pm can be interpreted as follows: Our multiscale test rejects the null hypothesis $H_0(u, h)$ if $|\widehat{\psi}_T(u, h)/\widehat{\sigma}| > q_T(\alpha) + \lambda(h)$. Put differently, it rejects $H_0(u, h)$ for all $(u, h) \in \mathcal{A}_T^\pm$. Hence, Π_T^\pm is the collection of time intervals $I_{u,h} = [u - h, u + h]$ for which our test rejects $H_0(u, h)$. The objects Π_T^+ and Π_T^- can be interpreted analogously: If $\widehat{\psi}_T(u, h)/\widehat{\sigma} > q_T(\alpha) + \lambda(h)$, that is, if $(u, h) \in \mathcal{A}_T^+$, then our test rejects $H_0(u, h)$ and indicates an increase in the trend m on the interval $I_{u,h}$, taking into account the positive sign of the statistic $\widehat{\psi}_T(u, h)/\widehat{\sigma}$. Hence, Π_T^+ is the collection of time intervals $I_{u,h}$ for which our test indicates an increase in the trend m . Likewise, Π_T^- is the collection of intervals for which the test indicates a decrease. Note that Π_T^\pm (as well as Π_T^+ and Π_T^-) is a random collection of intervals: Whether our test rejects $H_0(u, h)$ for some (u, h) depends on the realization of the random vector $(Y_{1,T}, \dots, Y_{T,T})$. Hence, whether an

interval $I_{u,h}$ belongs to Π_T^\pm depends on this realization as well. Having defined the objects Π_T^\pm , Π_T^+ and Π_T^- , we now consider the events

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

E_T^\pm (E_T^+ , E_T^-) is the event that the function m is non-constant (increasing, decreasing) on all intervals $I_{u,h} \in \Pi_T^\pm$ (Π_T^+ , Π_T^-). More precisely, E_T^\pm (E_T^+ , E_T^-) is the event that for each interval $I_{u,h} \in \Pi_T^\pm$ (Π_T^+ , Π_T^-), there is a subset $J_{u,h} \subseteq I_{u,h}$ with m being a non-constant (increasing, decreasing) function on $J_{u,h}$. We can make the following formal statement about the events E_T^\pm , E_T^+ and E_T^- , whose proof is given in the Supplement.

Proposition 3.3. *Let the conditions of Theorem 3.1 be fulfilled. Then for $\ell \in \{\pm, +, -\}$, it holds that*

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

According to Proposition 3.3, we can make simultaneous confidence statements of the following form: With (asymptotic) probability $\geq (1 - \alpha)$, the trend function m is non-constant (increasing, decreasing) on each interval $I_{u,h} \in \Pi_T^\pm$ (Π_T^+ , Π_T^-). Hence, our multiscale procedure allows to identify, with a pre-specified confidence, time intervals where there is an increase/decrease in the trend m .

Remark 3.2. *Unlike Π_T^\pm , the sets Π_T^+ and Π_T^- only contain 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. Crudely speaking, the problem is that the local linear weights $w_{t,T}(u, h)$ behave quite differently at boundary points (u, h) with $I_{u,h} \not\subseteq [0, 1]$. As a consequence, we can include boundary points (u, h) in Π_T^\pm but not in Π_T^+ and Π_T^- .*

Remark 3.3. *The statement of Proposition 3.3 suggests to graphically present the results of our multiscale test by plotting the intervals $I_{u,h} \in \Pi_T^\ell$ for $\ell \in \{\pm, +, -\}$, 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^ℓ is often quite large. To obtain a better graphical summary of the results, we replace Π_T^ℓ by a subset $\Pi_T^{\ell, \min}$ which is constructed as follows: As in Dümbgen (2002), we call an interval $I_{u,h} \in \Pi_T^\ell$ minimal if there is no other interval $I_{u',h'} \in \Pi_T^\ell$ with $I_{u',h'} \subset I_{u,h}$. Let $\Pi_T^{\ell, \min}$ be the set of all minimal intervals in Π_T^ℓ for*

$\ell \in \{\pm, +, -\}$ and define the events

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

It is easily seen that $E_T^\ell = E_T^{\ell, \min}$ for $\ell \in \{\pm, +, -\}$. Hence, by Proposition 3.3, it holds that

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

for $\ell \in \{\pm, +, -\}$. This suggests to plot the minimal intervals in $\Pi_T^{\ell, \min}$ rather than the whole collection of intervals Π_T^ℓ as a graphical summary of the test results. We in particular use this way of presenting the test results in our application in Section 6.

Proposition 3.3 allows to make confidence statements for a fixed significance level $\alpha \in (0, 1)$. In some situations, one may be interested in letting $\alpha = \alpha_T \in (0, 1)$ tend to zero as $T \rightarrow \infty$. This situation is considered in the following corollary to Proposition 3.3, whose proof can be found in the Supplementary Material.

Corollary 3.1. *Let the conditions of Theorem 3.1 be fulfilled and let $\alpha = \alpha_T \in (0, 1)$ go to zero as $T \rightarrow \infty$. Then $\mathbb{P}(E_T^\ell) \rightarrow 1$ for $\ell \in \{\pm, +, -\}$.*

Corollary 3.1 can be interpreted as a consistency result: If we let the significance level $\alpha = \alpha_T$ go to zero, then the event E_T^\pm (E_T^+ , E_T^-) occurs with probability tending to 1, that is, the trend m is non-constant (increasing, decreasing) on each interval $I_{u,h} \in \Pi_T^\pm$ (Π_T^+ , Π_T^-) with probability tending to 1.

3.4 Comparison to SiZer methods

As already mentioned in the Introduction, some SiZer methods for dependent data have been introduced in Park et al. (2004) and Rondonotti et al. (2007), which we refer to as dependent SiZer for short. Informally speaking, both our approach and dependent SiZer are methods to test for local increases/decreases of a nonparametric trend function m . The formal problem is to test the hypothesis $H_0(u, h)$ simultaneously for all $(u, h) \in \mathcal{G}_T$, where in this section, we let $\mathcal{G}_T = U_T \times H_T$ with U_T being the set of locations and H_T the set of bandwidths or scales. In what follows, we compare our approach to dependent SiZer and point out the most important differences.

Dependent SiZer is based on the statistics $\widehat{s}_T(u, h) = \widehat{m}'(u, h) / \widehat{\text{sd}}(\widehat{m}'(u, h))$, where $\widehat{m}'(u, h)$ is a local linear kernel estimator of $m'(u)$ with bandwidth h and $\widehat{\text{sd}}(\widehat{m}'(u, h))$ is an estimator of its standard deviation. The statistic $\widehat{s}_T(u, h)$ parallels the statistic $\widehat{\psi}_T(u, h) / \widehat{\sigma}$ in our approach. In particular, both can be regarded as test statistics of the hypothesis $H_0(u, h)$. There are two versions of dependent SiZer:

- (a) The global version aggregates the individual statistics $\widehat{s}_T(u, h)$ into the overall statistic $\widehat{S}_T = \max_{h \in H_T} \widehat{S}_T(h)$, where $\widehat{S}_T(h) = \max_{u \in U_T} |\widehat{s}_T(u, h)|$. The statistic \widehat{S}_T is the counterpart to the multiscale statistic $\widehat{\Psi}_T$ in our approach.
- (b) The row-wise version considers each scale $h \in H_T$ separately. In particular, for each bandwidth $h \in H_T$, a test is carried out based on the statistic $\widehat{S}_T(h)$. A row-wise analogue of our approach would be obtained by carrying out a test for each scale $h \in H_T$ separately based on the statistic $\widehat{\Psi}_T(h) = \max_{u \in U_T} |\widehat{\psi}_T(u, h)/\widehat{\sigma}|$.⁴

In practice, SiZer is commonly implemented in its row-wise form. The main reason is that it has more power than the global version by construction. However, this gain of power comes at a cost: Row-wise SiZer carries out a test *separately* for each scale $h \in H_T$, thus ignoring the simultaneous test problem across scales h . Hence, it is not a rigorous level- α -test of the null H_0 . For this reason, we focus on global SiZer in the rest of this section.

Even though related, our methods and theory are markedly different from those of the SiZer approach:

- (i) Theory for SiZer is derived under the assumption that $H_T \subseteq H$ for all T , where H is a compact subset of $(0, \infty)$. As already pointed out in Chaudhuri and Marron (2000) on p.420, this is a quite severe restriction: Only bandwidths h are taken into account that remain bounded away from zero as the sample size T increases. Bandwidths h that converge to zero are excluded. Our theory, in contrast, allows to simultaneously consider bandwidths h of fixed size and bandwidths h that converge to zero at different rates. To achieve this, we come up with a proof strategy which is very different from that in the SiZer literature: As proven in Chaudhuri and Marron (2000) for the i.i.d. case and in Park et al. (2009) for the dependent data case, \widehat{S}_T weakly converges to some limit process S under the overall null hypothesis H_0 . This is the central technical result on which the theoretical properties of SiZer are based. In contrast to this, our proof strategy (which combines strong approximation theory with anti-concentration bounds as outlined in Section 3.3) does not even require the statistic $\widehat{\Psi}_T$ to have a weak limit and is thus not restricted by the limitations of classic weak convergence theory.
- (ii) There are different ways to combine the test statistics $\widehat{S}_T(h) = \max_{u \in U_T} |\widehat{s}_T(u, h)|$ for different scales $h \in H_T$. One way is to take their maximum, which leads to the SiZer statistic $\widehat{S}_T = \max_{h \in H_T} \widehat{S}_T(h)$. We could proceed analogously and consider the statistic $\widehat{\Psi}_{T, \text{uncorrected}} = \max_{h \in H_T} \widehat{\Psi}_T(h) = \max_{(u, h) \in U_T \times H_T} |\widehat{\psi}_T(u, h)/\widehat{\sigma}|$. However, as argued in Dümbgen and Spokoiny (2001) and as discussed in Section 3.1, this aggregation scheme is not optimal when the set H_T contains scales h of many

⁴Note that we can drop the correction term $\lambda(h)$ in this case as it is a fixed constant if only a single bandwidth h is taken into account.

different rates. Following the lead of Dümbgen and Spokoiny (2001), we consider the test statistic $\widehat{\Psi}_T = \max_{(u,h) \in U_T \times H_T} \{|\widehat{\psi}_T(u,h)/\widehat{\sigma}| - \lambda(h)\}$ with the additive correction terms $\lambda(h)$. Hence, even though related, our multiscale test statistic $\widehat{\Psi}_T$ differs from the SiZer statistic \widehat{S}_T in important ways.

- (iii) The main complication in carrying out both our multiscale test and SiZer is to determine the critical values, that is, the quantiles of the test statistics $\widehat{\Psi}_T$ and \widehat{S}_T under H_0 . In order to approximate the quantiles, we proceed quite differently than in the SiZer literature. The quantiles of the SiZer statistic \widehat{S}_T can be approximated by those of the weak limit process S . Usually, however, the quantiles of S cannot be determined analytically but have to be approximated themselves (e.g. by the bootstrap procedures of Chaudhuri and Marron (1999, 2000)). Alternatively, the quantiles of \widehat{S}_T can be approximated by procedures based on extreme value theory (as proposed in Hannig and Marron (2006) and Park et al. (2009)). In our approach, the quantiles of $\widehat{\Psi}_T$ under H_0 are approximated by those of a suitably constructed Gaussian analogue of $\widehat{\Psi}_T$. It is far from obvious that this Gaussian approximation is valid when the data are dependent. To see this, deep strong approximation theory for dependent data (as derived in Berkes et al. (2014)) is needed. It is important to note that our Gaussian approximation procedure is not the same as the bootstrap procedures proposed in Chaudhuri and Marron (1999, 2000). Both procedures can of course be regarded as resampling methods. However, the resampling is done in a quite different way in our case.

4 Estimation of the long-run error variance

In this section, we discuss how to estimate the long-run variance $\sigma^2 = \sum_{\ell=-\infty}^{\infty} \text{Cov}(\varepsilon_0, \varepsilon_\ell)$ of the error terms in model (2.1). There are two broad classes of estimators: residual- and difference-based estimators. In residual-based approaches, σ^2 is estimated from the residuals $\widehat{\varepsilon}_t = Y_{t,T} - \widehat{m}_h(t/T)$, where \widehat{m}_h is a nonparametric estimator of m with the bandwidth or smoothing parameter h . Difference-based methods proceed by estimating σ^2 from the ℓ -th differences $Y_{t,T} - Y_{t-\ell,T}$ of the observed time series $\{Y_{t,T}\}$ for certain orders ℓ . In what follows, we focus attention on difference-based methods as these do not involve a nonparametric estimator of the function m and thus do not require to specify a bandwidth h for the estimation of m .

So far, we have assumed that $\{\varepsilon_t\}$ is a general stationary error process which fulfills the weak dependence conditions (C3). Estimating the long-run error variance σ^2 in model (2.1) under general weak dependence conditions is a notoriously difficult problem. Estimators of σ^2 often tend to be quite imprecise. To circumvent this issue in practice, it may be beneficial to impose a time series model on the error process $\{\varepsilon_t\}$. Estimating σ^2 under the restrictions of such a model may of course create some misspecification

bias. However, as long as the model gives a reasonable approximation to the true error process, the produced estimates of σ^2 can be expected to be fairly reliable even though they are a bit biased.

Estimators of the long-run error variance σ^2 in model (2.1) have been developed for different kinds of error models. A number of authors have analysed the case of $\text{MA}(m)$ or, more generally, m -dependent error terms. Difference-based estimators of σ^2 for this case were proposed in Müller and Stadtmüller (1988), Herrmann et al. (1992) and Tecuapetla-Gómez and Munk (2017) among others. Presumably the most widely used error model in practice is an $\text{AR}(p)$ process. Residual-based methods to estimate σ^2 in model (2.1) with $\text{AR}(p)$ errors can be found for example in Truong (1991), Shao and Yang (2011) and Qiu et al. (2013). A difference-based method was proposed in Hall and Van Keilegom (2003).

We consider the class of $\text{AR}(\infty)$ processes as an error model, which is a quite large and important subclass of linear time series processes. Formally speaking, we let $\{\varepsilon_t\}$ be a process of the form

$$\varepsilon_t = \sum_{j=1}^{\infty} a_j \varepsilon_{t-j} + \eta_t, \quad (4.1)$$

where a_1, a_2, a_3, \dots are unknown coefficients and η_t are i.i.d. with $\mathbb{E}[\eta_t] = 0$ and $\mathbb{E}[\eta_t^2] = \nu^2$. We assume that $A(z) := 1 - \sum_{j=1}^{\infty} a_j z^j \neq 0$ for all complex numbers $|z| \leq 1 + \delta$ with some small $\delta > 0$, which has the following implications: (i) $\{\varepsilon_t\}$ is stationary and causal. (ii) The coefficients a_j decay to zero exponentially fast, that is, $|a_j| \leq C\xi^j$ with some $C > 0$ and $\xi \in (0, 1)$. (iii) $\{\varepsilon_t\}$ has an $\text{MA}(\infty)$ representation of the form $\varepsilon_t = \sum_{k=0}^{\infty} c_k \eta_{t-k}$. The coefficients c_k can be computed iteratively from the equations

$$c_k - \sum_{j=1}^k a_j c_{k-j} = b_k \quad (4.2)$$

for $k = 0, 1, 2, \dots$, where $b_0 = 1$ and $b_k = 0$ for $k > 0$. Moreover, they decay to zero exponentially fast, that is, $|c_k| \leq C\xi^k$ with some $C > 0$ and $\xi \in (0, 1)$. Notably, the error model (4.1) nests $\text{AR}(p^*)$ processes of any finite order p^* as a special case: If $a_{p^*} \neq 0$ and $a_j = 0$ for all $j > p^*$, then $\{\varepsilon_t\}$ is an AR process of order p^* . In the sequel, we let $p^* \in \mathbb{N} \cup \{\infty\}$ denote the true AR order of $\{\varepsilon_t\}$ which may be finite or infinite. We can thus rewrite (4.1) as

$$\varepsilon_t = \sum_{j=1}^{p^*} a_j \varepsilon_{t-j} + \eta_t, \quad (4.3)$$

where the AR order p^* is treated as unknown.

We now construct a difference-based estimator of σ^2 for the case that $\{\varepsilon_t\}$ is an $\text{AR}(p^*)$ process of the form (4.3). To do so, we will fit $\text{AR}(p)$ type models to $\{\varepsilon_t\}$, where we distinguish between the following two cases:

(A) We do not know the precise AR order p^* but we know an upper bound p on it. In this case, p is a fixed natural number with $p \geq p^*$.

(B) We neither know p^* nor an upper bound on it. In this case, we let $p = p_T \rightarrow \infty$ as $T \rightarrow \infty$, where formal conditions on the growth of $p = p_T$ are specified later on.

To simplify notation, we let $\Delta_\ell Z_t = Z_t - Z_{t-\ell}$ denote the ℓ -th differences of a general time series $\{Z_t\}$. Our estimation method relies on the following simple observation: If $\{\varepsilon_t\}$ is an $\text{AR}(p^*)$ process of the form (4.3), then the time series $\{\Delta_q \varepsilon_t\}$ of the differences $\Delta_q \varepsilon_t = \varepsilon_t - \varepsilon_{t-q}$ is an $\text{ARMA}(p^*, q)$ process of the form

$$\Delta_q \varepsilon_t - \sum_{j=1}^{p^*} a_j \Delta_q \varepsilon_{t-j} = \eta_t - \eta_{t-q}. \quad (4.4)$$

As m is Lipschitz, the differences $\Delta_q \varepsilon_t$ of the unobserved error process are close to the differences $\Delta_q Y_{t,T}$ of the observed time series in the sense that

$$\Delta_q Y_{t,T} = [\varepsilon_t - \varepsilon_{t-q}] + \left[m\left(\frac{t}{T}\right) - m\left(\frac{t-q}{T}\right) \right] = \Delta_q \varepsilon_t + O\left(\frac{q}{T}\right). \quad (4.5)$$

Taken together, (4.4) and (4.5) imply that the differenced time series $\{\Delta_q Y_{t,T}\}$ is approximately an $\text{ARMA}(p^*, q)$ process of the form (4.4). It is precisely this point which is exploited by our estimation method.

We first describe our procedure to estimate the AR parameters a_j . For any $q \geq 1$, the $\text{ARMA}(p^*, q)$ process $\{\Delta_q \varepsilon_t\}$ satisfies the Yule-Walker equations

$$\gamma_q(\ell) - \sum_{j=1}^{p^*} a_j \gamma_q(\ell - j) = \begin{cases} -\nu^2 c_{q-\ell} & \text{for } 1 \leq \ell < q+1 \\ 0 & \text{for } \ell \geq q+1, \end{cases} \quad (4.6)$$

where $\gamma_q(\ell) = \text{Cov}(\Delta_q \varepsilon_t, \Delta_q \varepsilon_{t-\ell})$ and c_k are the coefficients from the $\text{MA}(\infty)$ expansion of $\{\varepsilon_t\}$. Combining the equations (4.6) for $\ell = 1, \dots, p$, we get that

$$\mathbf{\Gamma}_q \mathbf{a} = \boldsymbol{\gamma}_q + \nu^2 \mathbf{c}_q - \boldsymbol{\rho}_q, \quad (4.7)$$

where $\mathbf{a} = (a_1, \dots, a_p)^\top$, $\boldsymbol{\gamma}_q = (\gamma_q(1), \dots, \gamma_q(p))^\top$ and $\mathbf{\Gamma}_q$ denotes the $p \times p$ covariance matrix $\mathbf{\Gamma}_q = (\gamma_q(i-j) : 1 \leq i, j \leq p)$. Moreover, $\mathbf{c}_q = (c_{q-1}, \dots, c_{q-p})^\top$ and $\boldsymbol{\rho}_q = (\rho_q(1), \dots, \rho_q(p))^\top$ with $\rho_q(\ell) = \sum_{j=p+1}^{p^*} a_j \gamma_q(\ell-j)$. Since the AR coefficients a_j as well as the MA coefficients c_k decay exponentially fast to zero, $\boldsymbol{\rho}_q \approx \mathbf{0}$ and $\mathbf{c}_q \approx \mathbf{0}$ for large values of q , implying that $\mathbf{\Gamma}_q \mathbf{a} \approx \boldsymbol{\gamma}_q$. This suggests to estimate \mathbf{a} by

$$\tilde{\mathbf{a}}_q = \hat{\mathbf{\Gamma}}_q^{-1} \hat{\boldsymbol{\gamma}}_q, \quad (4.8)$$

where $\hat{\mathbf{\Gamma}}_q$ and $\hat{\boldsymbol{\gamma}}_q$ are defined analogously as $\mathbf{\Gamma}_q$ and $\boldsymbol{\gamma}_q$ with $\gamma_q(\ell)$ replaced by the sample

autocovariances $\hat{\gamma}_q(\ell) = (T - q)^{-1} \sum_{t=q+\ell+1}^T \Delta_q Y_{t,T} \Delta_q Y_{t-\ell,T}$ and $q = q_T$ goes to infinity as $T \rightarrow \infty$. For our theory to work, we require that $q/p \rightarrow \infty$, that is, q needs to grow faster than p . Formal conditions on the growth of q are given later on.

The estimator $\tilde{\mathbf{a}}_q$ depends on the tuning parameter q , that is, on the order of the differences $\Delta_q Y_{t,T}$. An appropriate choice of q needs to take care of the following two points: (i) q should be chosen large enough to ensure that the vector $\mathbf{c}_q = (c_{q-1}, \dots, c_{q-p})^\top$ is close to zero. As we have already seen, the constants c_k decay to zero exponentially fast and can be computed from the recursive equations (4.2) for given parameters a_1, a_2, a_3, \dots . In the special case of an AR(1) process, for example, one can readily calculate that $c_k \leq 0.0035$ for any $k \geq 20$ and any $|a_1| \leq 0.75$. Hence, if we have an AR(1) model for the errors ε_t and the error process is not too persistent, choosing $q \geq 20$ should make sure that \mathbf{c}_q is close to zero. Generally speaking, the recursive equations (4.2) can be used to get some idea for which values of q the vector \mathbf{c}_q can be expected to be approximately zero. (ii) q should not be chosen too large in order to ensure that the trend m is appropriately eliminated by taking q -th differences. As long as the trend m is not very strong, the two requirements (i) and (ii) can be fulfilled without much difficulty. For example, by choosing $q = 20$ in the AR(1) case just discussed, we do not only take care of (i) but also make sure that moderate trends m are differenced out appropriately.

When the trend m is very pronounced, in contrast, even moderate values of q may be too large to eliminate the trend appropriately. As a result, the estimator $\tilde{\mathbf{a}}_q$ will have a strong bias. In order to reduce this bias, we refine our estimation procedure as follows: By solving the recursive equations (4.2) with \mathbf{a} replaced by $\tilde{\mathbf{a}}_q$, we can compute estimators \tilde{c}_k of the coefficients c_k and thus estimators $\tilde{\mathbf{c}}_r$ of the vectors \mathbf{c}_r for any $r \geq 1$. Moreover, the innovation variance ν^2 can be estimated by $\tilde{\nu}^2 = (2T)^{-1} \sum_{t=p+2}^T \tilde{r}_{t,T}^2$, where $\tilde{r}_{t,T} = \Delta_1 Y_{t,T} - \sum_{j=1}^p \tilde{a}_j \Delta_1 Y_{t-j,T}$ and \tilde{a}_j is the j -th entry of the vector $\tilde{\mathbf{a}}_q$. Plugging the expressions $\hat{\Gamma}_r$, $\hat{\gamma}_r$, $\tilde{\mathbf{c}}_r$ and $\tilde{\nu}^2$ into (4.7), we can estimate \mathbf{a} by

$$\hat{\mathbf{a}}_r = \hat{\Gamma}_r^{-1} (\hat{\gamma}_r + \tilde{\nu}^2 \tilde{\mathbf{c}}_r), \quad (4.9)$$

where r is a much smaller differencing order than q . Specifically, in case (A), we can choose r to be any fixed number $r \geq 1$. Unlike q , the parameter r thus remains bounded as T increases. In case (B), our theory allows to choose any number r with $r \geq (1 + \delta)p$ for some small $\delta > 0$. Since $q/p \rightarrow \infty$, it holds that $q/r \rightarrow \infty$ as well, which means that r is of smaller order than q . Hence, in both cases (A) and (B), the estimator $\hat{\mathbf{a}}_r$ is based on a differencing order r that is much smaller than q ; only the pilot estimator $\tilde{\mathbf{a}}_q$ relies on differences of the larger order q . As a consequence, $\hat{\mathbf{a}}_r$ should eliminate the trend m more appropriately and should thus be less biased than the pilot estimator $\tilde{\mathbf{a}}_q$. In order to make the method more robust against estimation errors in $\tilde{\mathbf{c}}_r$, we finally

average the estimators $\hat{\mathbf{a}}_r$ for a few values of r . In particular, we define

$$\hat{\mathbf{a}} = \frac{1}{\bar{r} - \underline{r} + 1} \sum_{r=\underline{r}}^{\bar{r}} \hat{\mathbf{a}}_r, \quad (4.10)$$

where \underline{r} and \bar{r} are chosen as follows: In case (A), we let \underline{r} and \bar{r} be small natural numbers. In case (B), we set $\underline{r} = (1 - \delta)p$ for some small $\delta > 0$ and choose \bar{r} such that $\bar{r} - \underline{r}$ remains bounded. For ease of notation, we suppress the dependence of $\hat{\mathbf{a}}$ on the parameters \underline{r} and \bar{r} . Once $\hat{\mathbf{a}} = (\hat{a}_1, \dots, \hat{a}_p)^\top$ is computed, the long-run variance σ^2 can be estimated by

$$\hat{\sigma}^2 = \frac{\hat{\nu}^2}{(1 - \sum_{j=1}^p \hat{a}_j)^2}, \quad (4.11)$$

where $\hat{\nu}^2 = (2T)^{-1} \sum_{t=p+2}^T \hat{r}_{t,T}^2$ with $\hat{r}_{t,T} = \Delta_1 Y_{t,T} - \sum_{j=1}^p \hat{a}_j \Delta_1 Y_{t-j,T}$ is an estimator of the innovation variance ν^2 and we make use of the fact that $\sigma^2 = \nu^2 / (1 - \sum_{j=1}^{p^*} a_j)^2$ for the $\text{AR}(p^*)$ process $\{\varepsilon_t\}$.

We briefly compare the estimator $\hat{\mathbf{a}}$ to competing methods. Presumably closest to our approach is that of Hall and Van Keilegom (2003) which is designed for $\text{AR}(p^*)$ processes of known finite order p^* . For comparing the two methods, we thus assume p^* to be known and set $p = p^*$. The two main advantages of our method are as follows:

- (a) Our estimator produces accurate estimation results even when the AR process $\{\varepsilon_t\}$ is quite persistent, that is, even when the AR polynomial $A(z) = 1 - \sum_{j=1}^{p^*} a_j z^j$ has a root close to the unit circle. The estimator of Hall and Van Keilegom (2003), in contrast, may have very high variance and may thus produce unreliable results when the AR polynomial $A(z)$ is close to having a unit root. This difference in behaviour can be explained as follows: Our pilot estimator $\tilde{\mathbf{a}}_q = (\tilde{a}_1, \dots, \tilde{a}_{p^*})^\top$ has the property that the estimated AR polynomial $\tilde{A}(z) = 1 - \sum_{j=1}^{p^*} \tilde{a}_j z^j$ has no root inside the unit disc, that is, $\tilde{A}(z) \neq 0$ for all complex numbers z with $|z| \leq 1$.⁵ Hence, the fitted AR model with the coefficients $\tilde{\mathbf{a}}_q$ is ensured to be stationary and causal. Even though this may seem to be a minor technical detail, it has a huge effect on the performance of the estimator $\tilde{\mathbf{a}}_q$: It keeps the estimator stable even when the AR process is very persistent and the AR polynomial $A(z)$ has almost a unit root. This in turn results in a reliable behaviour of the estimator $\hat{\mathbf{a}}$ in the case of high persistence. The estimator of Hall and Van Keilegom (2003), in contrast, may produce non-causal results when the AR polynomial $A(z)$ is close to having a unit root. As a consequence, it may have unnecessarily high variance in the case of high persistence. We illustrate this difference between the estimators by the simulation exercises in Section 5.2. A striking example is Figure 6, which

⁵More precisely, $\tilde{A}(z) \neq 0$ for all z with $|z| \leq 1$, whenever the covariance matrix $(\hat{\gamma}_q(i-j) : 1 \leq i, j \leq p^* + 1)$ is non-singular. Moreover, $(\hat{\gamma}_q(i-j) : 1 \leq i, j \leq p^* + 1)$ is non-singular whenever $\hat{\gamma}_q(0) > 0$, which is the generic case.

presents the simulation results for the case of an AR(1) process $\varepsilon_t = a_1\varepsilon_{t-1} + \eta_t$ with $a_1 = -0.95$ and clearly shows the much better performance of our method.

- (b) Both our pilot estimator $\tilde{\mathbf{a}}_q$ and the estimator of Hall and Van Keilegom (2003) tend to have a substantial bias when the trend m is pronounced. Our estimator $\hat{\mathbf{a}}$ reduces this bias considerably as demonstrated in the simulations of Section 5.2. Unlike the estimator of Hall and Van Keilegom (2003), it thus produces accurate results even in the presence of a very strong trend.

We close this section by deriving some basic asymptotic properties of the estimators $\tilde{\mathbf{a}}_q$, $\hat{\mathbf{a}}$ and $\hat{\sigma}^2$. To formulate the following result, we use the shorthand $v_T \ll w_T$ which means that $v_T/w_T \rightarrow 0$ as $T \rightarrow \infty$.

Proposition 4.1. *Let m be Lipschitz continuous and suppose that $\{\varepsilon_t\}$ is an $AR(p^*)$ process of the form (4.3) with the following properties: $A(z) \neq 0$ for all $|z| \leq 1 + \delta$ with some small $\delta > 0$ and the innovations η_t have a finite fourth moment. Assume that p , q , \underline{r} and \bar{r} satisfy the following conditions: In case (A), p , \underline{r} and \bar{r} are fixed natural numbers and $\log T \ll q \ll \sqrt{T}$. In case (B), $C \log T \leq p \ll \min\{T^{1/5}, q\}$ for some sufficiently large C , $q \ll \sqrt{T}$, $\underline{r} = (1 + \delta)p$ for some small $\delta > 0$ and $\bar{r} - \underline{r}$ remains bounded. Under these conditions, $\tilde{\mathbf{a}}_q - \mathbf{a} = O_p(\sqrt{p/T})$ as well as $\hat{\mathbf{a}} - \mathbf{a} = O_p(\sqrt{p^3/T})$ and $\hat{\sigma}^2 - \sigma^2 = O_p(\sqrt{p^4/T})$.*

The proof is provided in the Supplementary Material. As one can see, the convergence rate of the second-step estimator $\hat{\mathbf{a}}$ is somewhat slower than that of the pilot estimator $\tilde{\mathbf{a}}_q$. Hence, from an asymptotic perspective, there is no gain from using the second-step estimator. Nevertheless, in finite samples, the estimator $\hat{\mathbf{a}}$ vastly outperforms $\tilde{\mathbf{a}}_q$ since it considerably reduces the bias of the latter.

5 Simulations

5.1 Small sample properties of the multiscale test

In what follows, we investigate the performance of our multiscale test and compare it to the dependent SiZer methods from Park et al. (2004), Rondonotti et al. (2007) and Park et al. (2009). We consider the following versions of our multiscale test and SiZer:

\mathcal{T}_{MS} : our multiscale test with the statistic $\hat{\Psi}_T = \max_{h \in H_T} \{\hat{\Psi}_T(h) - \lambda(h)\}$, where $\hat{\Psi}_T(h) = \max_{u \in U_T} |\hat{\psi}_T(u, h)/\hat{\sigma}|$. Here and in what follows, we write $\mathcal{G}_T = U_T \times H_T$, where U_T is the set of locations and H_T the set of bandwidths.

\mathcal{T}_{UC} : the uncorrected version of our multiscale test with the test statistic $\hat{\Psi}_{T, \text{uncorrected}} = \max_{h \in H_T} \hat{\Psi}_T(h)$, which was already introduced in (3.3). The uncorrected test is carried out in exactly the same way as \mathcal{T}_{MS} . The only difference is that the correction terms $\lambda(h)$ are removed.

\mathcal{T}_{RW} : a row-wise (or scale-wise) version of our multiscale test as briefly mentioned in Section 3.4. This version carries out a test scale-wise, that is, separately for each scale $h \in H_T$ based on the statistic $\widehat{\Psi}_T(h)$. Note: (i) For each $h \in H_T$, the test based on $\widehat{\Psi}_T(h)$ can be performed in the same way as the multiscale test \mathcal{T}_{MS} , since it is a degenerate version of the latter with the set of scales H_T replaced by the singleton $\{h\}$. (ii) It does not matter whether we correct the statistic $\widehat{\Psi}_T(h)$ by subtracting $\lambda(h)$ or not, since $\lambda(h)$ acts as a fixed constant when only one bandwidth h is taken into account.

$\mathcal{T}_{\text{SiZer}}$: the row-wise version of dependent SiZer from Park et al. (2004), Rondonotti et al. (2007) and Park et al. (2009). We do not consider a global version of dependent SiZer as such a version was not introduced in the aforementioned papers.

The simulation setup is as follows: We generate data from the model $Y_{t,T} = m(t/T) + \varepsilon_t$ for different trends m , error processes $\{\varepsilon_t\}$ and sample sizes T . The error terms are supposed to have the AR(1) structure $\varepsilon_t = a_1 \varepsilon_{t-1} + \eta_t$, where $a_1 \in \{-0.9, -0.5, -0.25, 0.25, 0.5, 0.9\}$, η_t are i.i.d. standard normal and the AR order $p^* = 1$ is treated as known. To simulate data under the null, 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 different non-constant trend functions which are specified below. For each model specification, we simulate $S = 1000$ data samples and carry out the tests \mathcal{T}_{MS} , \mathcal{T}_{UC} , \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ for each simulated sample.

To implement our multiscale test \mathcal{T}_{MS} , we choose K to be an Epanechnikov kernel and let $\mathcal{G}_T = U_T \times H_T$ with

$$U_T = \left\{ u \in [0, 1] : u = \frac{5t}{T} \text{ for some } t \in \mathbb{N} \right\}$$

$$H_T = \left\{ h \in \left[\frac{\log T}{T}, \frac{1}{4} \right] : h = \frac{5\ell}{T} \text{ for some } \ell \in \mathbb{N} \right\}.$$

We thus take into account all locations u on an equidistant grid U_T with step length $5/T$ and all bandwidths $h = 5/T, 10/T, 15/T, \dots$ with $\log T/T \leq h \leq 1/4$. Note that the lower bound $\log T/T$ is motivated by (C6) which requires that $\log T/T \ll h_{\min}$ (given that all moments of ε_t exist). As a robustness check, we have re-run the simulations for a number of other grids. As the results are very similar, we do however not report them here. To estimate the long-run error variance σ^2 , we apply the procedure from Section 4 with $\underline{r} = 1$, $\bar{r} = 10$ and the following choices of q : For $a_1 \in \{-0.5, -0.25, 0.25, 0.5\}$, we set $q = 25$. As already discussed in Section 4, this should be an appropriate choice for AR(1) errors that are not too strongly correlated, in particular, for $a_1 \in \{-0.5, -0.25, 0.25, 0.5\}$. When the errors are very strongly correlated, larger values of q are required to produce precise estimates of σ^2 . In the case of AR(1) errors with $a_1 \in \{-0.9, 0.9\}$, we thus set $q = 50$. The dependence of our long-run variance estimator on the tuning parameters q , \underline{r} and \bar{r} is explored more systematically

Table 1: Size of \mathcal{T}_{MS} for the AR parameters $a_1 \in \{-0.5, -0.25, 0.25, 0.5\}$.

	$a_1 = -0.5$			$a_1 = -0.25$			$a_1 = 0.25$			$a_1 = 0.5$		
	nominal size α			nominal size α			nominal size α			nominal size α		
	0.01	0.05	0.1	0.01	0.05	0.1	0.01	0.05	0.1	0.01	0.05	0.1
$T = 250$	0.013	0.040	0.086	0.016	0.054	0.106	0.009	0.045	0.094	0.014	0.058	0.106
$T = 500$	0.013	0.044	0.102	0.008	0.041	0.089	0.013	0.057	0.107	0.014	0.056	0.101
$T = 1000$	0.011	0.052	0.090	0.007	0.057	0.114	0.011	0.049	0.106	0.007	0.050	0.098

Table 2: Size of \mathcal{T}_{MS} for the AR parameters $a_1 \in \{-0.9, 0.9\}$.

	$a_1 = -0.9$					$a_1 = 0.9$				
	sample size T					sample size T				
	250	500	1000	2000	3000	250	500	1000	2000	3000
$\alpha = 0.01$	0.040	0.032	0.017	0.009	0.012	0.003	0.016	0.015	0.021	0.017
$\alpha = 0.05$	0.137	0.093	0.067	0.061	0.047	0.017	0.038	0.055	0.059	0.057
$\alpha = 0.1$	0.218	0.160	0.124	0.108	0.098	0.040	0.054	0.095	0.096	0.106

in Section 5.2. To compute the critical values of the multiscale test \mathcal{T}_{MS} , we simulate 5000 values of the statistic Φ_T defined in Section 3.2 and compute their empirical $(1 - \alpha)$ quantile $q_T(\alpha)$. The uncorrected and row-wise versions \mathcal{T}_{UC} and \mathcal{T}_{RW} of our multiscale test are implemented analogously. The SiZer test is implemented as described in Park et al. (2009). The details are summarized in Section S.3 of the Supplementary Material.

5.1.1 Size simulations

The first part of our simulation study investigates the size properties of the four tests \mathcal{T}_{MS} , \mathcal{T}_{UC} , \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ under the null that the trend m is constant. To start with, we focus on the multiscale test \mathcal{T}_{MS} . Table 1 reports the actual size of \mathcal{T}_{MS} for the AR parameters $a_1 \in \{-0.5, -0.25, 0.25, 0.5\}$, which is computed as the number of simulations in which \mathcal{T}_{MS} rejects the null divided by the total number of simulations. As can be seen, the actual size of the multiscale test \mathcal{T}_{MS} is fairly close to the nominal target α for all the considered AR parameters and sample sizes. Hence, the test has approximately the correct size.

In Table 1, we have explored the size of \mathcal{T}_{MS} when the errors are moderately autocorrelated. The case of strongly autocorrelated errors is investigated in Table 2, where we consider AR(1) errors with $a_1 \in \{-0.9, 0.9\}$. We first discuss the results for the positive AR parameter $a_1 = 0.9$. As can be seen, the size numbers are substantially downward biased for small sample sizes, in particular, for $T = 250$ and $T = 500$. As the sample size increases, this downward bias diminishes and the size numbers stabilize around their target α . In particular, for $T \geq 1000$, the size numbers give a decent approximation to α . An analogous picture arises for the negative AR parameter $a_1 = -0.9$. The size numbers, however, are upward rather than downward biased for small sample sizes

Table 3: Global size comparisons for the significance level $\alpha = 0.05$.

	$a_1 = -0.5$				$a_1 = 0.5$			
	\mathcal{T}_{MS}	\mathcal{T}_{UC}	\mathcal{T}_{RW}	$\mathcal{T}_{\text{SiZer}}$	\mathcal{T}_{MS}	\mathcal{T}_{UC}	\mathcal{T}_{RW}	$\mathcal{T}_{\text{SiZer}}$
$T = 250$	0.069	0.065	0.230	0.333	0.049	0.048	0.143	0.289
$T = 500$	0.054	0.065	0.288	0.448	0.042	0.026	0.187	0.397
$T = 1000$	0.046	0.051	0.318	0.522	0.052	0.049	0.276	0.509

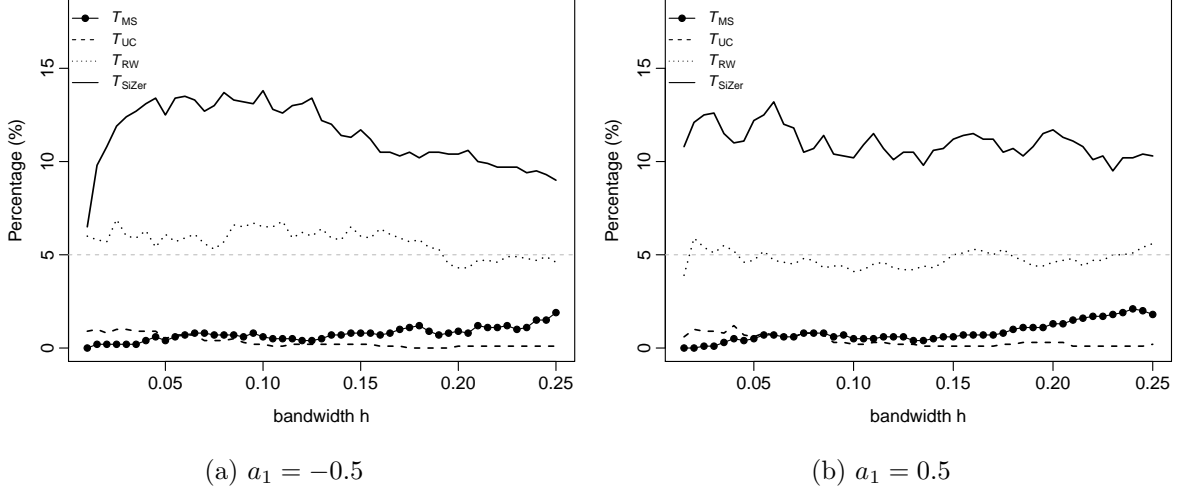


Figure 2: Row-wise size comparisons for the significance level $\alpha = 5\%$ and the sample size $T = 1000$. Subfigure (a) corresponds to the case with $a_1 = -0.5$, subfigure (b) to the case with $a_1 = 0.5$. Each curve in the two subfigures shows the row-wise size (given in percentage % on the y -axis) as a function of the bandwidth h (specified on the x -axis) for one of the four tests \mathcal{T}_{MS} , \mathcal{T}_{UC} , \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$.

T and the size distortions appear to vanish a bit more slowly as T increases. To summarize, in the case of strongly autocorrelated errors, our multiscale test has good size properties only for sufficiently large sample sizes. This is not very surprising: Statistical inference in the presence of strongly autocorrelated data is a very difficult problem in general and satisfying results can only be expected for fairly large sample sizes.

We next compare our multiscale test \mathcal{T}_{MS} with \mathcal{T}_{UC} , \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ in terms of size. There is an important difference between \mathcal{T}_{MS} and \mathcal{T}_{UC} on the one hand and \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ on the other. \mathcal{T}_{MS} and its uncorrected version \mathcal{T}_{UC} are *global* test procedures: They test $H_0(u, h)$ simultaneously for all locations $u \in U_T$ and scales $h \in H_T$. Hence, they control the size simultaneously over both locations u and scales h . The methods \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$, in contrast, are *row-wise* (or *scale-wise*) in nature: They test the hypothesis $H_0(u, h)$ simultaneously for all $u \in U_T$ but separately for each scale $h \in H_T$. Hence, they control the size for each scale $h \in H_T$ separately.

We conduct some simulation exercises to illustrate this important distinction. To keep the simulation study to a reasonable length, we restrict attention to the significance level $\alpha = 0.05$ and the AR parameters $a_1 \in \{-0.5, 0.5\}$. To simplify the implementation of

$\mathcal{T}_{\text{SiZer}}$, we assume that the autocovariance function of the error process and thus the long-run error variance σ^2 is known. To keep the comparison fair, we treat σ^2 as known also when implementing \mathcal{T}_{MS} , \mathcal{T}_{UC} and \mathcal{T}_{RW} . Moreover, we use exactly the same location-scale grid for all four methods. To achieve this, we start off with the grid $\mathcal{G}_T = U_T \times H_T$ with U_T and H_T defined above. We then follow Rondonotti et al. (2007) and Park et al. (2009) and restrict attention to those points $(u, h) \in \mathcal{G}_T$ for which the effective sample size $\text{ESS}^*(u, h)$ for correlated data is not smaller than 5. This yields the grid $\mathcal{G}_T^* = \{(u, h) \in \mathcal{G}_T : \text{ESS}^*(u, h) \geq 5\}$. A definition of $\text{ESS}^*(u, h)$ is given in Section S.3 of the Supplement.

For our simulation exercises, we distinguish between global and row-wise (or scale-wise) size: Global size is defined as the percentage of simulations in which the test under consideration rejects $H_0(u, h)$ for some $(u, h) \in \mathcal{G}_T^*$. Hence, it is identical to the size as computed in Tables 1 and 2. Row-wise size for scale $h^* \in H_T$, in contrast, is the percentage of simulations in which the test rejects $H_0(u, h^*)$ for some $(u, h^*) \in \mathcal{G}_T^*$. Table 3 reports the global size of the four tests. As can be seen, the size numbers of our multiscale test \mathcal{T}_{MS} and its uncorrected version \mathcal{T}_{UC} are reasonably close to the target $\alpha = 0.05$. The global size numbers of the row-wise methods \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$, in contrast, are much larger than the target $\alpha = 0.05$. Since the number of scales h in the grid \mathcal{G}_T^* increases with T , they even move away from α as the sample size T increases. To summarize, as expected, the global tests \mathcal{T}_{MS} and \mathcal{T}_{UC} hold the size reasonably well, whereas the row-wise methods \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ are much too liberal.

Figure 2 reports the row-wise size of the four tests by so-called parallel coordinate plots [Inselberg (1985)] for the sample size $T = 1000$. Each curve in the figure specifies the row-wise size of one of the tests for the scales h under consideration. As can be seen, the row-wise version \mathcal{T}_{RW} of our multiscale test holds the size quite accurately across scales. The row-wise size of $\mathcal{T}_{\text{SiZer}}$ also gives an acceptable approximation to the target $\alpha = 5\%$, even though the size numbers are upward biased quite a bit. The global tests \mathcal{T}_{MS} and \mathcal{T}_{UC} , in contrast, have a row-wise size much smaller than the target $\alpha = 5\%$, which reflects the fact that they control global rather than row-wise size.

5.1.2 Power comparisons

In the second part of our simulation study, we compare the tests \mathcal{T}_{MS} , \mathcal{T}_{UC} , \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ in terms of power. As above, we use the location-scale grid \mathcal{G}_T^* and treat the autocovariance function of the error terms as known when implementing the tests. Moreover, we restrict attention to the significance level $\alpha = 0.05$ and the AR parameters $a_1 \in \{-0.5, 0.5\}$. Our simulation exercises investigate the ability of the four tests to detect local increases in the trend m . (The same could of course be done for decreases.) The tests indicate a local increase in m according to the following decision rules: For each $(u, h) \in \mathcal{G}_T^*$,

Table 4: Global power and global spurious power comparisons for $\alpha = 0.05$.

		$a_1 = -0.5$				$a_1 = 0.5$			
		\mathcal{T}_{MS}	\mathcal{T}_{UC}	\mathcal{T}_{RW}	$\mathcal{T}_{\text{SiZer}}$	\mathcal{T}_{MS}	\mathcal{T}_{UC}	\mathcal{T}_{RW}	$\mathcal{T}_{\text{SiZer}}$
$T = 250$	Power	0.094	0.070	0.215	0.294	0.100	0.081	0.197	0.307
	Spurious power	0.021	0.032	0.109	0.166	0.012	0.017	0.054	0.131
$T = 500$	Power	0.186	0.138	0.418	0.571	0.194	0.168	0.431	0.602
	Spurious power	0.020	0.024	0.137	0.212	0.016	0.016	0.082	0.192
$T = 1000$	Power	0.504	0.360	0.760	0.872	0.550	0.427	0.795	0.895
	Spurious power	0.023	0.024	0.158	0.283	0.020	0.019	0.123	0.252

$$\begin{aligned}
 \mathcal{T}_{\text{MS}} \text{ indicates an increase on } [u - h, u + h] &\iff \hat{\psi}_T(u, h)/\hat{\sigma} > q_T(\alpha) + \lambda(h) \\
 \mathcal{T}_{\text{UC}} \text{ indicates an increase on } [u - h, u + h] &\iff \hat{\psi}_T(u, h)/\hat{\sigma} > q_T^{\text{UC}}(\alpha) \\
 \mathcal{T}_{\text{RW}} \text{ indicates an increase on } [u - h, u + h] &\iff \hat{\psi}_T(u, h)/\hat{\sigma} > q_T^{\text{RW}}(\alpha, h) \\
 \mathcal{T}_{\text{SiZer}} \text{ indicates an increase on } [u - h, u + h] &\iff \hat{s}_T(u, h) > q_T^{\text{SiZer}}(\alpha, h),
 \end{aligned}$$

where $q_T^{\text{UC}}(\alpha)$, $q_T^{\text{RW}}(\alpha, h)$, $q_T^{\text{SiZer}}(\alpha, h)$ are the critical values of \mathcal{T}_{UC} , \mathcal{T}_{RW} , $\mathcal{T}_{\text{SiZer}}$, respectively. Note that the critical values of \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ depend on the scale h as these are row-wise procedures.

To be able to make systematic power comparisons, we consider a very simple trend function m . More complicated signals m are analysed in Section S.3 of the Supplementary Material. The trend function we are considering here is defined as $m(u) = c \cdot 1(u \in [0.45, 0.55]) \cdot (1 - \{\frac{u-0.5}{0.05}\}^2)^2$, where $c = 0.85$ in the AR case with $a_1 = -0.5$ and $c = 2.65$ in the case with $a_1 = 0.5$. The function m is increasing on $I^+ = (0.45, 0.5)$, decreasing on $I^- = (0.5, 0.55)$ and constant elsewhere. The two upper panels of Figure 3 give a graphical illustration of m , where the grey line in the background is the time series path of a representative simulated data sample. As can be seen, m is a small bump around $u = 0.5$, where c determines the height of the bump. The constant c is chosen such that the bump is difficult but not impossible to detect for the four tests. We distinguish between the following types of power for the tests \mathcal{T}_j with $j \in \{\text{MS}, \text{UC}, \text{RW}, \text{SiZer}\}$, where we restrict attention to increases in m :

- (i) global power: the percentage of simulation runs in which the test \mathcal{T}_j indicates an increase on some interval $I_{u,h} = [u - h, u + h]$ where m is indeed increasing, that is, on some $I_{u,h}$ with $I_{u,h} \cap I^+ \neq \emptyset$.
- (ii) spurious global power: the percentage of simulation runs in which the test \mathcal{T}_j indicates an increase on some interval $I_{u,h} = [u - h, u + h]$ where m is not increasing, that is, on some $I_{u,h}$ with $I_{u,h} \cap I^+ = \emptyset$.
- (iii) row-wise power on scale h^* : the percentage of simulation runs in which the test \mathcal{T}_j indicates an increase on some interval $I_{u,h^*} = [u - h^*, u + h^*]$ where m is indeed increasing, that is, on some I_{u,h^*} with $I_{u,h^*} \cap I^+ \neq \emptyset$.

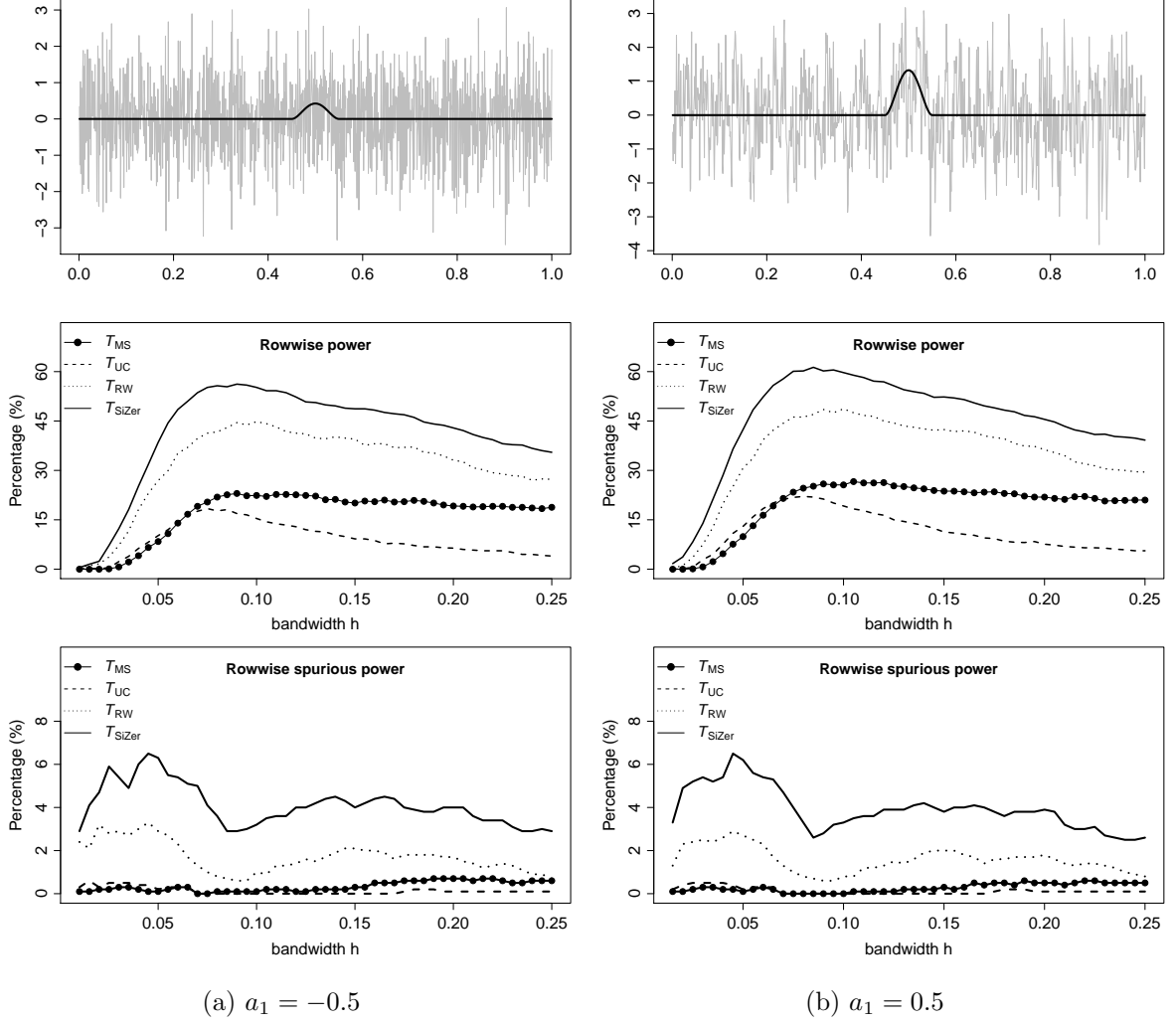


Figure 3: Row-wise power and row-wise spurious power comparisons for $\alpha = 5\%$ and $T = 1000$. Subfigure (a) corresponds to the case with $a_1 = -0.5$, subfigure (b) to the case with $a_1 = 0.5$. The upper panel of each subfigure shows the bump function m with a representative data sample in the background. The parallel coordinate plot in the middle panel reports row-wise power. In particular, each curve shows the row-wise power (given in percentage % on the y -axis) as a function of the bandwidth h (specified on the x -axis) for one of the four tests \mathcal{T}_{MS} , \mathcal{T}_{UC} , \mathcal{T}_{RW} and \mathcal{T}_{SiZer} . The parallel coordinate plot in the lower panel reports row-wise spurious power in an analogous fashion.

- (iv) spurious row-wise power on scale h^* : the percentage of simulation runs in which the test \mathcal{T}_j indicates an increase on some interval $I_{u,h^*} = [u - h^*, u + h^*]$ where m is not increasing, that is, on some I_{u,h^*} with $I_{u,h^*} \cap I^+ = \emptyset$.

Table 4 reports the global power and global spurious power of the four tests. As can be seen, our multiscale test \mathcal{T}_{MS} has higher power than the uncorrected version \mathcal{T}_{UC} . This confirms the theoretical optimality theory in Dümbgen and Spokoiny (2001) [see also Dümbgen and Walther (2008) and Rufibach and Walther (2010)] according to which the aggregation scheme of \mathcal{T}_{MS} with its additive correction term should yield better power properties than the simpler scheme of \mathcal{T}_{UC} . As expected, the row-wise methods \mathcal{T}_{RW}

and $\mathcal{T}_{\text{SiZer}}$ have substantially more power than the global tests. Indeed, $\mathcal{T}_{\text{SiZer}}$ is even a bit more powerful than \mathcal{T}_{RW} , which is presumably due to the fact that it is somewhat too liberal in terms of row-wise size as observed in Figure 2. The higher power of the row-wise procedures comes at some cost: Their spurious global power is much higher than that of the global tests. For the sample size $T = 1000$ and the AR parameter $a_1 = -0.5$, for example, $\mathcal{T}_{\text{SiZer}}$ spuriously finds an increase in the trend m in more than 28% of the simulations, \mathcal{T}_{RW} in more than 15%. The multiscale test \mathcal{T}_{MS} (as well as its uncorrected version \mathcal{T}_{UC}), in contrast, controls the probability of finding a spurious increase. In particular, as implied by Proposition 3.3, its spurious global power is below $100 \cdot \alpha\% = 5\%$.

Figure 3 gives a more detailed picture of the power properties of the four tests for the sample size $T = 1000$. The parallel coordinate plots of the figure show how power and spurious power are distributed across scales h . Let us first have a look at the row-wise methods. As can be seen, $\mathcal{T}_{\text{SiZer}}$ is more powerful than \mathcal{T}_{RW} on all scales under consideration. As already mentioned when discussing the global power results, this is presumably due to the fact that $\mathcal{T}_{\text{SiZer}}$ is a bit too liberal in terms of row-wise size. Comparing the power curves of the two global methods gives an interesting insight: Our multiscale test \mathcal{T}_{MS} has substantially more power than the uncorrected version \mathcal{T}_{UC} on medium and large scales. On small scales, in contrast, it is slightly less powerful than \mathcal{T}_{UC} . This again illustrates the theoretical optimality theory in Dümbgen and Spokoiny (2001) which suggests that, asymptotically, the multiscale test \mathcal{T}_{MS} should be as powerful as \mathcal{T}_{UC} on small scales but more powerful on large scales. This is essentially what we see in the two middle panels of Figure 3. Of course, \mathcal{T}_{MS} does not have exactly as much power as \mathcal{T}_{UC} on fine scales. However, the loss of power on fine scales is very small compared to the gain of power on larger scales (which is also reflected by the fact that \mathcal{T}_{MS} has more global power than \mathcal{T}_{UC}).

The main findings of our simulation exercises can be summarized as follows: If one is interested in an exploratory data tool for finding local increases/decreases of a trend, the row-wise methods \mathcal{T}_{RW} and $\mathcal{T}_{\text{SiZer}}$ both do a good job. However, if one wants to make rigorous statistical inference simultaneously across locations and scales, one needs to go for a global method. Our simulation exercises have demonstrated that our multiscale test \mathcal{T}_{MS} is a global method which enjoys good size and power properties. In particular, as predicted by the theory, it is a more effective test than the uncorrected version \mathcal{T}_{UC} .

5.2 Small sample properties of the long-run variance estimator

In the final part of our simulation study, we analyse the estimators of the AR parameters and of the long-run error variance from Section 4 and compare them to the estimators of Hall and Van Keilegom (2003). We simulate data from the model $Y_{t,T} = m(t/T) + \varepsilon_t$, where $\{\varepsilon_t\}$ is an AR(1) process of the form $\varepsilon_t = a_1 \varepsilon_{t-1} + \eta_t$. We consider the

AR parameters $a_1 \in \{-0.95, -0.75, -0.5, -0.25, 0.25, 0.5, 0.75, 0.95\}$ and let η_t be i.i.d. standard normal innovation terms. Throughout the simulation study, the AR order $p^* = 1$ is treated as known. We report our findings for the sample size $T = 500$, the results for other sample sizes being very similar. For simplicity, m is chosen to be a linear function of the form $m(u) = \beta u$ with the slope parameter β . For each value of a_1 , we consider two different slopes β , one corresponding to a moderate and one to a pronounced trend m . In particular, we let $\beta = s_\beta \sqrt{\text{Var}(\varepsilon_t)}$ with $s_\beta \in \{1, 10\}$. When $s_\beta = 1$, the slope β is equal to the standard deviation $\sqrt{\text{Var}(\varepsilon_t)}$ of the error process, which yields a moderate trend m . When $s_\beta = 10$, in contrast, the slope β is 10 times as large as $\sqrt{\text{Var}(\varepsilon_t)}$, which results in a quite pronounced trend m .

For each model specification, we generate $S = 1000$ data samples and compute the following quantities for each simulated sample:

- (i) the pilot estimator \tilde{a}_q from (4.8) with the tuning parameter q , the estimator \hat{a} from (4.10) with the tuning parameters (\underline{r}, \bar{r}) and the long-run variance estimator $\hat{\sigma}^2$ from (4.11).
- (ii) the estimators of a_1 and σ^2 from Hall and Van Keilegom (2003), which are denoted by \hat{a}_{HVK} and $\hat{\sigma}_{\text{HVK}}^2$. The estimator \hat{a}_{HVK} is computed as described in Section 2.2 of Hall and Van Keilegom (2003) and $\hat{\sigma}_{\text{HVK}}^2$ as defined at the bottom of p.447 in Section 2.3. The estimator \hat{a}_{HVK} (as well as $\hat{\sigma}_{\text{HVK}}^2$) depends on two tuning parameters which we denote by m_1 and m_2 as in Hall and Van Keilegom (2003).
- (iii) oracle estimators \hat{a}_{oracle} and $\hat{\sigma}_{\text{oracle}}^2$ of a_1 and σ^2 , which are constructed under the assumption that the error process $\{\varepsilon_t\}$ is observed. For each simulation run, we compute \hat{a}_{oracle} as the maximum likelihood estimator of a_1 from the time series of simulated error terms $\varepsilon_1, \dots, \varepsilon_T$. We then calculate the residuals $r_t = \varepsilon_t - \hat{a}_{\text{oracle}} \varepsilon_{t-1}$ and estimate the innovation variance $\nu^2 = \mathbb{E}[\eta_t^2]$ by $\hat{\nu}_{\text{oracle}}^2 = (T - 1)^{-1} \sum_{t=2}^T r_t^2$. Finally, we set $\hat{\sigma}_{\text{oracle}}^2 = \hat{\nu}_{\text{oracle}}^2 / (1 - \hat{a}_{\text{oracle}})^2$.

Throughout the section, we set $q = 25$, $(\underline{r}, \bar{r}) = (1, 10)$ and $(m_1, m_2) = (20, 30)$. We in particular choose q to be in the middle of m_1 and m_2 to make the tuning parameters of the estimators \tilde{a}_q and \hat{a}_{HVK} more or less comparable. In order to assess how sensitive our estimators are to the choice of q and (\underline{r}, \bar{r}) , we carry out a number of robustness checks, considering a range of different values for q and (\underline{r}, \bar{r}) . In addition, we vary the tuning parameters m_1 and m_2 of the estimators from Hall and Van Keilegom (2003) to make sure that the results of our comparison study are not driven by the particular choice of any of the involved tuning parameters. The results of our robustness checks are reported in Section S.3 of the Supplement. They show that the results of our comparison study are robust to different choices of the parameters q , (\underline{r}, \bar{r}) and (m_1, m_2) .

For each estimator \hat{a} , \hat{a}_{HVK} , \hat{a}_{oracle} and $\hat{\sigma}^2$, $\hat{\sigma}_{\text{HVK}}^2$, $\hat{\sigma}_{\text{oracle}}^2$ and for each model specification, the simulation output consists in a vector of length $S = 1000$ which contains the 1000

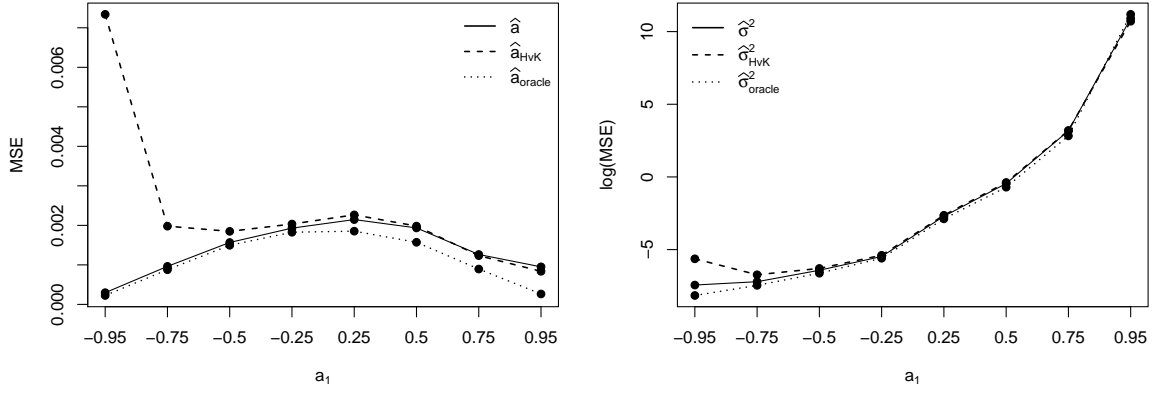


Figure 4: MSE values for the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and $\hat{\sigma}^2$, $\hat{\sigma}_{HvK}^2$, $\hat{\sigma}_{oracle}^2$ in the simulation scenarios with a moderate trend ($s_\beta = 1$).

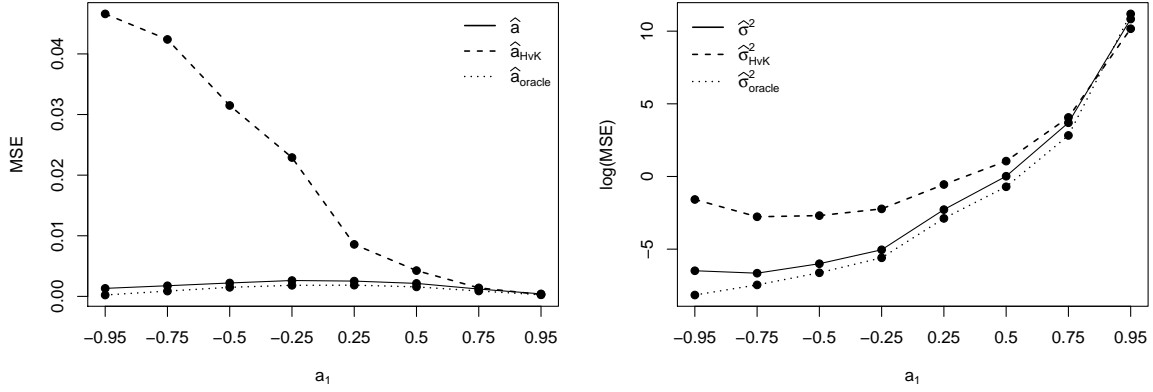


Figure 5: MSE values for the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and $\hat{\sigma}^2$, $\hat{\sigma}_{HvK}^2$, $\hat{\sigma}_{oracle}^2$ in the simulation scenarios with a pronounced trend ($s_\beta = 10$).

simulated values of the respective estimator. Figures 4 and 5 report the mean squared error (MSE) of these 1000 simulated values for each estimator. On the x -axis of each plot, the various values of the AR parameter a_1 are listed which are considered. The solid line in each plot gives the MSE values of our estimators. The dashed and dotted lines specify the MSE values of the HvK and the oracle estimators, respectively. Note that for the long-run variance estimators, the plots report the logarithm of the MSE rather than the MSE itself since the MSE values are too different across simulation scenarios to obtain a reasonable graphical presentation. In addition to the MSE values presented in Figures 4 and 5, we depict histograms of the 1000 simulated values produced by the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and $\hat{\sigma}^2$, $\hat{\sigma}_{HvK}^2$, $\hat{\sigma}_{oracle}^2$ for two specific simulation scenarios in Figures 6 and 7. The main findings can be summarized as follows:

- (a) In the simulation scenarios with a moderate trend ($s_\beta = 1$), the estimators \hat{a}_{HvK} and $\hat{\sigma}_{HvK}^2$ of Hall and Van Keilegom (2003) exhibit a similar performance as our estimators \hat{a} and $\hat{\sigma}^2$ as long as the AR parameter a_1 is not too close to -1 . For strongly negative values of a_1 (in particular for $a_1 = -0.75$ and $a_1 = -0.95$),

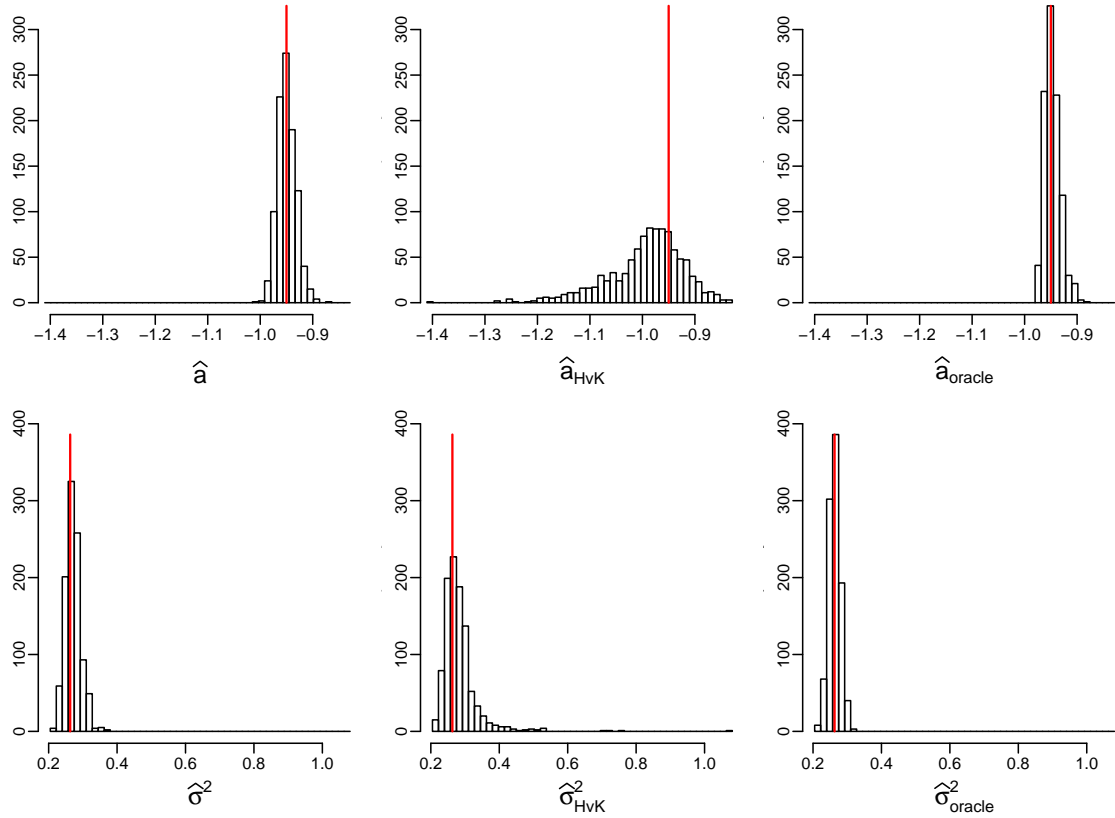


Figure 6: Histograms of the simulated values produced by the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and $\hat{\sigma}^2$, $\hat{\sigma}_{HvK}^2$, $\hat{\sigma}_{oracle}^2$ in the scenario with $a_1 = -0.95$ and $s_\beta = 1$. The vertical red lines indicate the true values of a_1 and σ^2 .

the estimators perform much worse than ours. This can be clearly seen from the much larger MSE values of the estimators \hat{a}_{HvK} and $\hat{\sigma}_{HvK}^2$ for $a_1 = -0.75$ and $a_1 = -0.95$ in Figure 4. Figure 6 gives some further insights into what is happening here. It shows the histograms of the simulated values produced by the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and the corresponding long-run variance estimators in the scenario with $a_1 = -0.95$ and $s_\beta = 1$. As can be seen, the estimator \hat{a}_{HvK} does not obey the causality restriction $|a_1| < 1$ but frequently takes values substantially smaller than -1 . This results in a very large spread of the histogram and thus in a disastrous performance of the estimator.⁶ A similar point applies to the histogram of the long-run variance estimator $\hat{\sigma}_{HvK}^2$. Our estimators \hat{a} and $\hat{\sigma}^2$, in contrast, exhibit a stable behaviour in this case.

Interestingly, the estimator \hat{a}_{HvK} (as well as the corresponding long-run variance estimator $\hat{\sigma}_{HvK}^2$) performs much worse than ours for large negative values but not for large positive values of a_1 . This can be explained as follows: In the special case of an AR(1) process, the estimator \hat{a}_{HvK} may produce estimates smaller than -1

⁶One could of course set \hat{a}_{HvK} to $-(1 - \delta)$ for some small $\delta > 0$ whenever it takes a value ≤ -1 . This modified estimator, however, is still far from performing in a satisfactory way when a_1 is close to -1 .

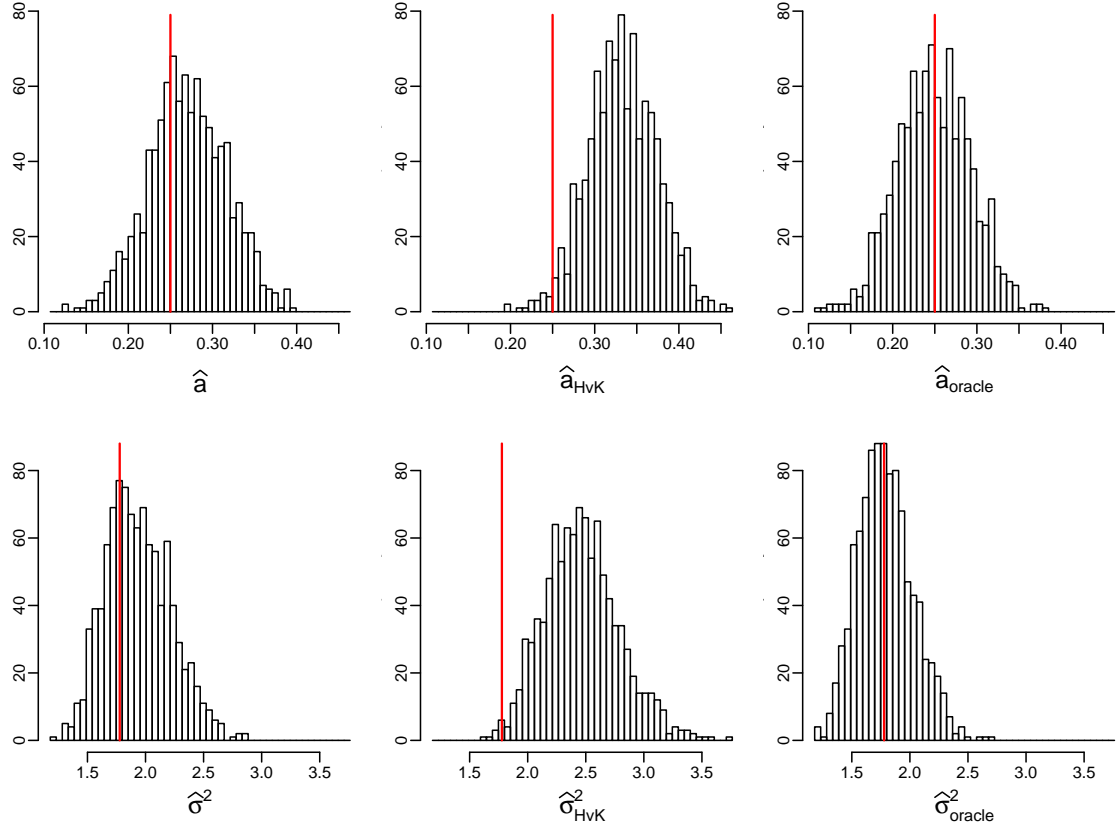


Figure 7: Histograms of the simulated values produced by the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and $\hat{\sigma}^2$, $\hat{\sigma}_{HvK}^2$, $\hat{\sigma}_{oracle}^2$ in the scenario with $a_1 = 0.25$ and $s_\beta = 10$. The vertical red lines indicate the true values of a_1 and σ^2 .

but it cannot become larger than 1. This can be easily seen upon inspecting the definition of the estimator. Hence, for large positive values of a_1 , the estimator \hat{a}_{HvK} performs well as it satisfies the causality restriction that the estimated AR parameter should be smaller than 1.

- (b) In the simulation scenarios with a pronounced trend ($s_\beta = 10$), the estimators of Hall and Van Keilegom (2003) are clearly outperformed by ours for most of the AR parameters a_1 under consideration. In particular, their MSE values reported in Figure 5 are much larger than the values produced by our estimators for most parameter values a_1 . The reason is the following: The HvK estimators have a strong bias since the pronounced trend with $s_\beta = 10$ is not eliminated appropriately by the underlying differencing methods. This point is illustrated by Figure 7 which shows histograms of the simulated values for the estimators \hat{a} , \hat{a}_{HvK} , \hat{a}_{oracle} and the corresponding long-run variance estimators in the scenario with $a_1 = 0.25$ and $s_\beta = 10$. As can be seen, the histogram produced by our estimator \hat{a} is approximately centred around the true value $a_1 = 0.25$, whereas that of \hat{a}_{HvK} is strongly biased upwards. A similar picture arises for the long-run variance estimators $\hat{\sigma}^2$ and $\hat{\sigma}_{HvK}^2$. Whereas the methods of Hall and Van Keilegom (2003) perform much worse than

ours for negative and moderately positive values of a_1 , the performance (in terms of MSE) is fairly similar for large values of a_1 . This can be explained as follows: When the trend m is not eliminated appropriately by taking differences, this creates spurious persistence in the data. Hence, the estimator \hat{a}_{HvK} tends to overestimate the AR parameter a_1 , that is, \hat{a}_{HvK} tends to be larger in absolute value than a_1 . Very loosely speaking, when the parameter a_1 is close to 1, say $a_1 = 0.95$, there is not much room for overestimation since \hat{a}_{HvK} cannot become larger than 1. Consequently, the effect of not eliminating the trend appropriately has a much smaller impact on \hat{a}_{HvK} for large positive values of a_1 .

6 Application

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. In what follows, we use our multiscale test \mathcal{T}_{MS} to analyse two long-term temperature records. Throughout the section, we set the significance level to $\alpha = 0.05$ and implement the multiscale test in exactly the same way as in the simulation study of Section 5.

6.1 Analysis of the Central England temperature record

The Central England temperature record is the longest instrumental temperature time series in the world. The data are 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 $Y_{t,T}$ covering the years from 1659 to 2017. A plot of the time series is given in panel (a) of Figure 8. We assume that the temperature data $Y_{t,T}$ follow the nonparametric trend model $Y_{t,T} = m(t/T) + \varepsilon_t$, where m is the unknown time trend of interest. The error process $\{\varepsilon_t\}$ is supposed to have the $\text{AR}(p^*)$ structure $\varepsilon_t = \sum_{j=1}^{p^*} a_j \varepsilon_{t-j} + \eta_t$, where η_t are i.i.d. innovations with mean 0 and variance ν^2 . As pointed out in Mudelsee (2010) among others, this is the most widely used error model for discrete climate time series. We select the AR order p^* by the Bayesian information criterion (BIC), which yields $p^* = 2$.⁷ We then estimate the $\text{AR}(2)$ parameters $\mathbf{a} = (a_1, a_2)$ and the long-run error variance σ^2 by the procedures from Section 4 with $q = 25$ and $(\underline{r}, \bar{r}) = (1, 10)$. This gives the estimators $\hat{a}_1 = 0.164$, $\hat{a}_2 = 0.175$ and $\hat{\sigma}^2 = 0.737$.

⁷More precisely, we proceed as follows: We estimate the AR parameters and the corresponding variance of the innovation terms for different AR orders by the methods from Section 4 and then choose p^* as the minimizer of the Bayesian information criterion (BIC). As a robustness check, we have repeated this procedure for a wide range of the tuning parameters q and (\underline{r}, \bar{r}) , which produces the value $p^* = 2$ throughout. Moreover, we have considered other information criteria such as FPE, AIC and AICC, which gives the AR order $p^* = 2$ for almost all values of q and (\underline{r}, \bar{r}) .

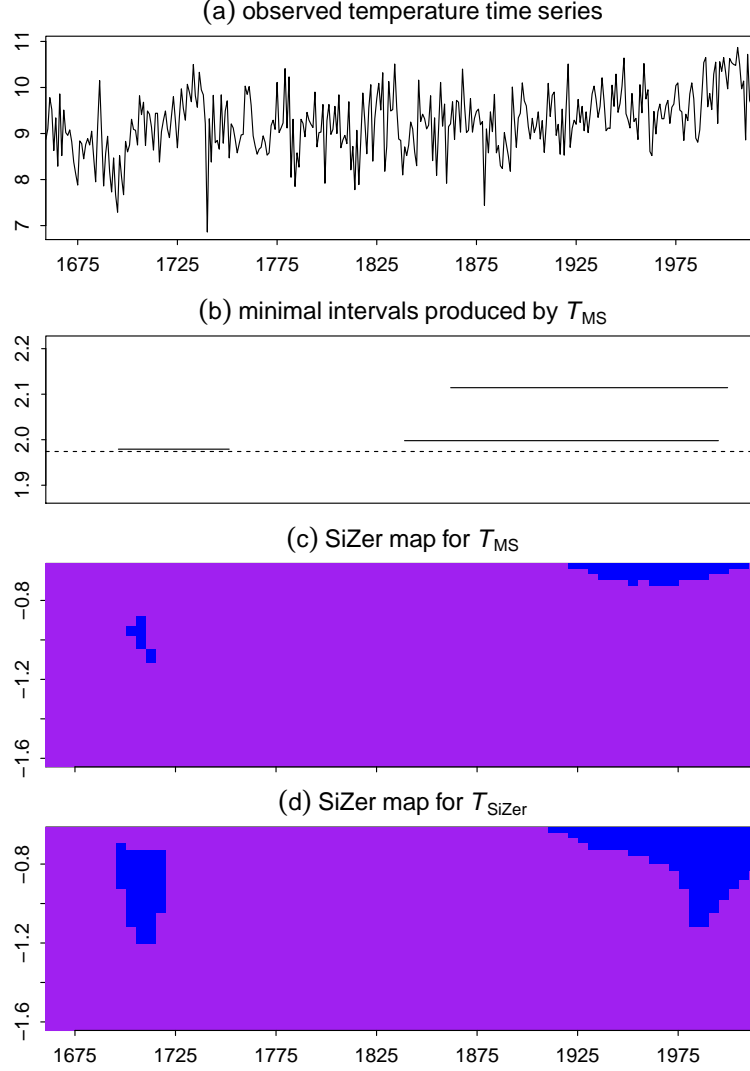


Figure 8: Summary of the results for the Central England temperature record. Panel (a) shows the observed temperature time series. Panel (b) depicts the minimal intervals in the set Π_T^+ produced by our multiscale test. These are [1684, 1744], [1839, 2009] and [1864, 2014]. Panels (c) and (d) present the SiZer maps produced by our multiscale test and SiZer.

With the help of our multiscale method, we now test the null hypothesis H_0 that m is constant on all intervals $[u - h, u + h]$ with $(u, h) \in \mathcal{G}_T^*$, where the grid \mathcal{G}_T^* is defined in the same way as in Section 5. The results are presented in Figure 8. Panel (b) depicts the minimal intervals in the set Π_T^+ which is produced by our multiscale test \mathcal{T}_{MS} . 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 $\widehat{\psi}_T(u, h)/\widehat{\sigma} - \lambda(h)$. The dashed line specifies the critical value $q_T(\alpha)$, where $\alpha = 0.05$ as already mentioned above. According to Proposition 3.3, we can make the following simultaneous confidence statement about the collection of minimal intervals plotted in panel (b). We can claim, with confidence of about 95%, that the trend 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 about 1870 onwards. Hence, our test in particular provides evidence that there has been some warming trend in the period over approximately the last 150 years. On the other hand, as the set Π_T^- is empty, there is no evidence of any downward movement of the trend.

Panel (c) presents the SiZer map produced by our multiscale test \mathcal{T}_{MS} . For comparison, the SiZer map of the dependent SiZer test $\mathcal{T}_{\text{SiZer}}$ is shown in panel (d). To produce panel (d), we have implemented SiZer as described in Section S.3 of the Supplement, where the autocovariance function of the errors $\{\varepsilon_t\}$ is estimated with the help of our procedures from Section 4 under the assumption that $\{\varepsilon_t\}$ is an AR(2) process. The SiZer maps of panels (c) and (d) are to be read as follows: Each pixel of the map corresponds to a location-scale point (u, h) , or put differently, to a time interval $[u - h, u + h]$. The pixel (u, h) is coloured blue if the test indicates an increase in the trend m on the interval $[u - h, u + h]$, red if the test indicates a decrease and purple if the test does not reject the null hypothesis that m is constant on $[u - h, u + h]$. As can be seen, the two SiZer maps in panels (c) and (d) have a similar structure. Both our multiscale test and SiZer indicate increases in the trend m during a short time period around 1700 and towards the end of the sample. However, in contrast to SiZer, our method allows to make formal confidence statements about the regions of blue pixels in the SiZer map. In particular, as the set of blue pixels in panel (c) exactly corresponds to the collection of intervals Π_T^+ , we can claim, with confidence of about 95%, that the trend m has an increase on each time interval represented by a blue pixel in panel (c).

6.2 Analysis of global temperature data

We next analyse a data set which consists of annual global temperature anomalies from 1850 onwards. The data are publicly available on the webpage <https://cdiac.ess-dive.lbl.gov/trends/temp/jonescru/jones.html> and are plotted in panel (a) of Figure 9. As before, we assume that the data come from the model $Y_{t,T} = m(t/T) + \varepsilon_t$, where m is the trend and $\{\varepsilon_t\}$ the noise process. We apply our multiscale methods to test the null hypothesis H_0 that m is constant on all time intervals $[u - h, u + h]$ with $(u, h) \in \mathcal{G}_T$, where the grid \mathcal{G}_T is defined as in Section 5. We compare our results with those obtained by Wu et al. (2001) who developed a method for testing the hypothesis that m is constant on $[0, 1]$ against the alternative that m is an arbitrary monotonic function. For comparability reasons, we use exactly the same data as in Wu et al. (2001), in particular, the yearly temperature anomalies from 1856 to 1998. Moreover, we use their estimate of the long-run error variance σ^2 which amounts to 0.01558. As we do not have an estimate available from Wu et al. (2001) for the autocovariance function of the error process, we do not consider dependent SiZer in the application example at hand.

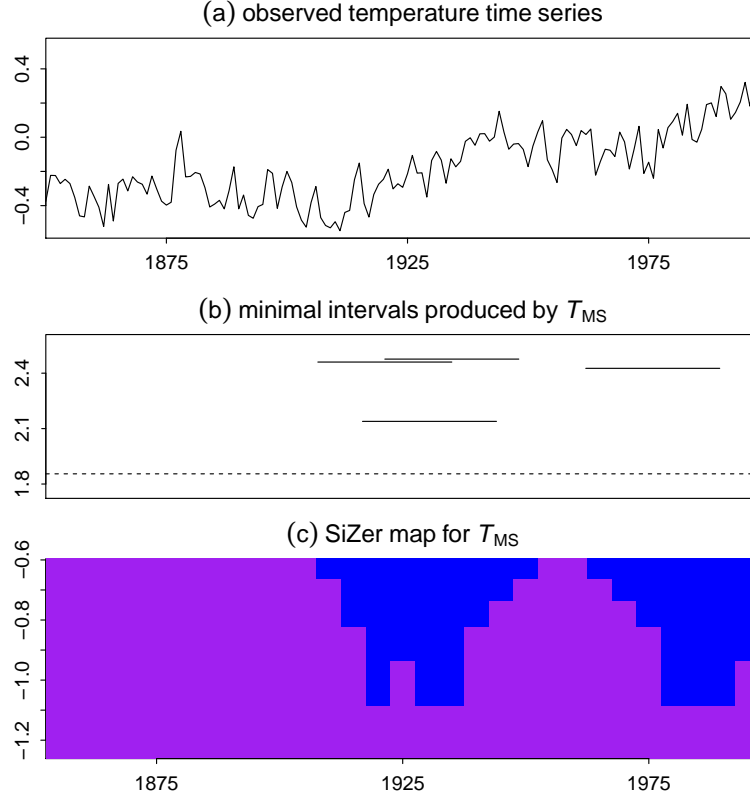


Figure 9: Summary of the results for the global temperature anomalies. Panel (a) shows the observed temperature time series. Panel (b) depicts the minimal intervals in the set Π_T^+ produced by the multiscale test. These are $[1905, 1935]$, $[1915, 1945]$, $[1920, 1950]$ and $[1965, 1995]$. Panel (c) presents the SiZer map of our test.

The results produced by our multiscale test are reported in Figure 9. Panel (b) shows the minimal intervals in Π_T^+ and panel (c) the SiZer map of the test. As can be clearly seen from both panels (b) and (c), the test indicates an increase in the trend m during the first half of the 20th century followed by another increase during the second half. These findings are in line with those in Wu et al. (2001) who reject the null hypothesis that m is constant. In contrast to the test of Wu et al. (2001), however, our multiscale method does not only allow to test whether the null is violated. It also allows to make formal confidence statements about where violations occur, that is, about where the trend m is increasing. In particular, we can claim, with confidence of about 95%, that the trend has an increase on each interval plotted in panel (b) of Figure 9.

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-Tusnádý 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.
- CHERNOZHUKOV, V., CHETVERIKOV, D. and KATO, K. (2014). Gaussian approximation of suprema of empirical processes. *Annals of Statistics*, **42** 1564–1597.
- 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.
- CHERNOZHUKOV, V., CHETVERIKOV, D. and KATO, K. (2017). Central limit theorems and bootstrap in high dimensions. *Annals of Probability*, **45** 2309–2352.
- CHO, H. and FRYZLEWICZ, P. (2012). Multiscale and multilevel technique for consistent segmentation of nonstationary time series. *Statistica Sinica*, **22** 207–229.
- DONOHU, D., JOHNSTONE, I., KERKYACHARIAN, G. and PICARD, D. (1995). Wavelet shrinkage: Asymptopia? *Journal of the Royal Statistical Society: Series B*, **57** 301–369.
- 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.
- DÜMBGEN, L. and WALTHER, G. (2008). Multiscale inference about a density. *Annals of Statistics*, **36** 1758–1785.
- ECKLE, K., BISSANTZ, N. and DETTE, H. (2017). Multiscale inference for multivariate deconvolution. *Electronic Journal of Statistics*, **11** 4179–4219.
- 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.
- HANNIG, J. and MARRON, J. S. (2006). Advanced distribution theory for SiZer. *Journal of the American Statistical Association*, **101** 484–499.
- HERRMANN, E., GASSER, T. and KNEIP, A. (1992). Choice of bandwidth for kernel regression when residuals are correlated. *Biometrika*, **79** 783–795.
- INSELBERG, A. (1985). The plane with parallel coordinates. *The Visual Computer*, **1** 69–91.
- 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., , HANNIG, J. and KANG, K.-H. (2009). Improved SiZer for time series. *Statistica Sinica*, **19** 1511–1530.
- 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.
- 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.
- PROKSCH, K., WERNER, F. and MUNK, A. (2018). Multiscale scanning in inverse problems. *Forthcoming in Annals of Statistics*.
- QIU, D., SHAO, Q. and YANG, L. (2013). Efficient inference for autoregressive coefficients in the presence of trends. *Journal of Multivariate Analysis*, **114** 40–53.
- RAHMSTORF, S., FOSTER, G. and CAHILL, N. (2017). Global temperature evolution: recent trends and some pitfalls. *Environmental Research Letters*, **12**.
- ROHDE, A. (2008). Adaptive goodness-of-fit tests based on signed ranks. *Annals of Statistics*, **36** 1346–1374.
- 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.
- RUFIBACH, K. and WALTHER, G. (2010). The block criterion for multiscale inference about a density, with applications to other multiscale problems. *Journal of Computational and Graphical Statistics*, **19** 175–190.
- SCHMIDT-HIEBER, J., MUNK, A. and DÜMBGEN, L. (2013). Multiscale methods for shape constraints in deconvolution: confidence statements for qualitative features. *Annals of Statistics*, **41** 1299–1328.
- SHAO, Q. and YANG, L. J. (2011). Autoregressive coefficient estimation in nonparametric analysis. *Journal of Time Series Analysis*, **32** 587–597.
- 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.
- TRUONG, Y. K. (1991). Nonparametric curve estimation with time series errors. *Journal of Statistical Planning and Inference*, **28** 167–183.
- VON SACHS, R. and MACGIBBON, B. (2000). Non-parametric curve estimation by Wavelet thresholding with locally stationary errors. *Scandinavian Journal of Statistics*, **27** 475–499.
- 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.
- WU, W. B., WOODROOFE, M. and MENTZ, G. (2001). Isotonic regression: another look at the changepoint problem. *Biometrika*, **88** 793–804.