

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Journal of Econometrics

journal homepage: [www.elsevier.com/locate/jeconom](http://www.elsevier.com/locate/jeconom)

# Nonparametric comparison of epidemic time trends: The case of COVID-19

Marina Khismatullina<sup>a,\*</sup>, Michael Vogt<sup>b,1</sup>

<sup>a</sup> Bonn Graduate School of Economics, University of Bonn, 53113 Bonn, Germany

<sup>b</sup> Institute of Statistics, Department of Mathematics and Economics, Ulm University, 89081 Ulm, Germany



## ARTICLE INFO

### Article history:

Received 4 August 2020

Received in revised form 13 February 2021

Accepted 8 April 2021

Available online 24 May 2021

### JEL classification:

C12

C23

C54

### Keywords:

Simultaneous hypothesis testing

Multiscale test

Time trend

Panel data

COVID-19

## ABSTRACT

The COVID-19 pandemic is one of the most pressing issues at present. A question which is particularly important for governments and policy makers is the following: Does the virus spread in the same way in different countries? Or are there significant differences in the development of the epidemic? In this paper, we devise new inference methods that allow to detect differences in the development of the COVID-19 epidemic across countries in a statistically rigorous way. In our empirical study, we use the methods to compare the outbreak patterns of the epidemic in a number of European countries.

© 2021 Published by Elsevier B.V.

## 1. Introduction

There are many questions surrounding the current COVID-19 pandemic that are not well understood yet. A question which is particularly important for governments and policy makers is the following: How do the outbreak patterns of COVID-19 compare across countries? Are the time trends of daily new infections more or less the same across countries, or is the virus spreading differently in different regions of the world? Identifying differences between countries may help, for instance, to better understand which government policies have been more effective in containing the virus than others. The main aim of this paper is to develop new inference methods that allow to detect differences between time trends of COVID-19 infections in a statistically rigorous way.

Let  $X_{it}$  be the number of new infections on day  $t$  in country  $i$  and suppose we observe a sample of data  $\mathcal{X}_i = \{X_{it} : 1 \leq t \leq T\}$  for  $n$  different countries  $i$ . In order to make the data comparable across countries, we take the starting date  $t = 1$  to be the first Monday after reaching 100 confirmed cases in each country. Considering the dates after reaching a certain level of confirmed cases is a common practice of “normalizing” the data (see e.g. [Cohen and Kupferschmidt, 2020](#)). Starting on a Monday additionally aligns the data across countries by the day of the week. This allows us to take care of possible weekly cycles in the data which are produced by delays in reporting new cases over the weekend. A simple way

\* Corresponding author.

E-mail addresses: [marina.k@uni-bonn.de](mailto:marina.k@uni-bonn.de) (M. Khismatullina), [m.vogt@uni-ulm.de](mailto:m.vogt@uni-ulm.de) (M. Vogt).

<sup>1</sup> Financial support by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation), Germany – grant VO 2503/1-1, project number 430668955 – is gratefully acknowledged.

to model the count data  $X_{it}$  is to use a Poisson distribution. Specifically, we may assume that the random variables  $X_{it}$  are Poisson distributed with time-varying intensity parameter  $\lambda_i(t/T)$ , that is,  $X_{it} \sim P_{\lambda_i(t/T)}$ . Since  $\lambda_i(t/T) = \mathbb{E}[X_{it}] = \text{Var}(X_{it})$ , we can model the observations  $X_{it}$  by the nonparametric regression equation

$$X_{it} = \lambda_i\left(\frac{t}{T}\right) + u_{it} \tag{1.1}$$

for  $1 \leq t \leq T$ , where  $u_{it} = X_{it} - \mathbb{E}[X_{it}]$  with  $\mathbb{E}[u_{it}] = 0$  and  $\text{Var}(u_{it}) = \lambda_i(t/T)$ . As usual in nonparametric regression (see e.g. [Robinson, 1989](#)), we let the regression function  $\lambda_i$  in model (1.1) depend on rescaled time  $t/T$  rather than on real time  $t$ . Hence,  $\lambda_i : [0, 1] \rightarrow \mathbb{R}$  can be regarded as a function on the unit interval, which allows us to estimate it by standard techniques from nonparametric regression. Since  $\lambda_i$  is a function of rescaled time  $t/T$ , the variables  $X_{it}$  in model (1.1) depend on the time series length  $T$  in general, that is,  $X_{it} = X_{it,T}$ . To keep the notation simple, we however suppress this dependence throughout the paper. In Section 2, we introduce the model setting in detail which underlies our analysis. As we will see there, it is a generalized version of the Poisson model (1.1).

In model (1.1), the time trend of new COVID-19 infections in country  $i$  is described by the intensity function  $\lambda_i$  of the underlying Poisson distribution. Hence, the question whether the time trends are comparable across countries amounts to the question whether the intensity functions  $\lambda_i$  have the same shape across countries  $i$ . In this paper, we construct a multiscale test which allows to *identify* and *locate* the differences between the functions  $\lambda_i$ . More specifically, let  $\mathcal{F} = \{\mathcal{I}_k \subseteq [0, 1] : 1 \leq k \leq K\}$  be a family of (rescaled) time intervals  $\mathcal{I}_k$  and let  $H_0^{(ijk)}$  be the hypothesis that the functions  $\lambda_i$  and  $\lambda_j$  are the same on the interval  $\mathcal{I}_k$ , that is,

$$H_0^{(ijk)} : \lambda_i(w) = \lambda_j(w) \text{ for all } w \in \mathcal{I}_k.$$

We design a method to test the hypothesis  $H_0^{(ijk)}$  *simultaneously* for all pairs of countries  $i$  and  $j$  under consideration and for all intervals  $\mathcal{I}_k$  in the family  $\mathcal{F}$ . The main theoretical result of the paper shows that the method controls the familywise error rate, that is, the probability of wrongly rejecting at least one null hypothesis  $H_0^{(ijk)}$ . As we will see, this allows us to make simultaneous confidence statements of the following form for a given significance level  $\alpha \in (0, 1)$ :

*With probability at least  $1 - \alpha$ , the functions  $\lambda_i$  and  $\lambda_j$  differ on the interval  $\mathcal{I}_k$  for every  $(i, j, k)$  for which the test rejects  $H_0^{(ijk)}$ .*

Hence, the method allows us to make simultaneous confidence statements (a) about which time trend functions differ from each other and (b) about where, that is, in which time intervals  $\mathcal{I}_k$  they differ.

Even though our multiscale test is motivated by the current COVID-19 crisis, its applicability is by no means restricted to this specific event. It is a general method to compare nonparametric trends in epidemiological (count) data. It thus contributes to the literature on statistical tests for equality of nonparametric regression and trend curves. Examples of such tests can be found in [Härdle and Marron \(1990\)](#), [Hall and Hart \(1990\)](#), [King et al. \(1991\)](#), [Delgado \(1993\)](#), [Kulasekera \(1995\)](#), [Young and Bowman \(1995\)](#), [Munk and Dette \(1998\)](#), [Lavergne \(2001\)](#), [Neumeier and Dette \(2003\)](#) and [Pardo-Fernández et al. \(2007\)](#). More recent approaches were developed in [Degras et al. \(2012\)](#), [Zhang et al. \(2012\)](#), [Hidalgo and Lee \(2014\)](#) and [Chen and Wu \(2019\)](#). Compared to existing methods, our test has the following crucial advantage: it is much more informative. Most existing procedures allow to test *whether* the regression or trend curves under consideration are all the same or not. However, they do not allow to infer *which* curves are different and *where* (that is, in which parts of the support) they differ. Our multiscale approach, in contrast, conveys this information. Indeed, it even allows to make rigorous confidence statements about which curves  $\lambda_i$  are different and where they differ. To the best of our knowledge, there is no other method available in the literature which allows to make such simultaneous confidence statements. As far as we know, the only other multiscale test for comparing trend curves has been developed in [Park et al. \(2009\)](#). However, their analysis is mainly methodological and not backed up by a general theory. In particular, theory is only available for the special case  $n = 2$ . Moreover, the theoretical results are only valid under very severe restrictions on the family of time intervals  $\mathcal{F}$ .

The paper is structured as follows. As already mentioned above, Section 2 details the model setting which underlies our analysis. The multiscale test is developed step by step in Section 3. To keep the presentation as clear as possible, the technical details are deferred to the [Appendix](#) and the Supplementary Material. Section 4 contains the empirical part of the paper. There, we run some simulation experiments to demonstrate that the multiscale test has the formal properties predicted by the theory. Moreover, we use the test to compare the outbreak patterns of the COVID-19 epidemic in a number of European countries.

## 2. Model setting

As already discussed in the Introduction, the assumption that  $X_{it} \sim P_{\lambda_i(t/T)}$  leads to a nonparametric regression model of the form

$$X_{it} = \lambda_i\left(\frac{t}{T}\right) + u_{it} \quad \text{with} \quad u_{it} = \sqrt{\lambda_i\left(\frac{t}{T}\right)}\eta_{it}, \tag{2.1}$$

where  $\eta_{it}$  has zero mean and unit variance. In this model, both the mean and the variance are described by the same function  $\lambda_i$ . In empirical applications, however, the variance often tends to be much larger than the mean. To deal with this

issue, which has been known for a long time in the literature (Cox, 1983) and which is commonly called overdispersion, so-called quasi-Poisson models (McCullagh and Nelder, 1989; Efron, 1986) are frequently used. In our context, a quasi-Poisson model of  $X_{it}$  has the form

$$X_{it} = \lambda_i\left(\frac{t}{T}\right) + \varepsilon_{it} \quad \text{with} \quad \varepsilon_{it} = \sigma \sqrt{\lambda_i\left(\frac{t}{T}\right)} \eta_{it}, \quad (2.2)$$

where  $\sigma$  is a scaling factor that allows the variance to be a multiple of the mean function  $\lambda_i$ . In what follows, we assume that the observed data  $X_{it}$  are produced by model (2.2), where the noise residuals  $\eta_{it}$  have zero mean and unit variance but we do not impose any further distributional assumptions on them.

Poisson and quasi-Poisson models are often used in the literature on epidemic modelling. De Salazar et al. (2020), for example, assume that the observed COVID-19 case count in country  $i$  follows a Poisson distribution with parameter  $\lambda_i$  being a linear function of some covariate  $Z_i$ , that is,  $\lambda_i = \beta Z_i$ . Pellis et al. (2020) consider a quasi-Poisson model for the number of new COVID-19 cases. They in particular examine (a) a version of the model where the mean function is parametrically restricted to be exponentially growing with a constant growth rate and (b) a version where the mean function is modelled nonparametrically by splines. Tobías et al. (2020) analyse data on the accumulated number of cases using quasi-Poisson regression, where the mean function is modelled parametrically as a piecewise linear curve with known change points.

In order to derive our theoretical results, we impose the following regularity conditions on model (2.2):

- (C1) The functions  $\lambda_i$  are uniformly Lipschitz continuous, that is,  $|\lambda_i(u) - \lambda_i(v)| \leq L|u - v|$  for all  $u, v \in [0, 1]$ , where the constant  $L$  does not depend on  $i$ . Moreover, they are uniformly bounded away from zero and infinity, that is, there exist constants  $\lambda_{\min}$  and  $\lambda_{\max}$  with  $0 < \lambda_{\min} \leq \min_{w \in [0, 1]} \lambda_i(w) \leq \max_{w \in [0, 1]} \lambda_i(w) \leq \lambda_{\max} < \infty$  for all  $i$ .
- (C2) The random variables  $\eta_{it}$  are independent both across  $i$  and  $t$ . Moreover, for any  $i$  and  $t$ , it holds that  $\mathbb{E}[\eta_{it}] = 0$ ,  $\mathbb{E}[\eta_{it}^2] = 1$  and  $\mathbb{E}[|\eta_{it}|^\theta] \leq C_\theta < \infty$  for some  $\theta > 4$ .

We briefly comment on the above conditions.

- (C1) imposes some standard-type regularity conditions on the functions  $\lambda_i$ . In particular, the functions are assumed to be smooth, bounded from above and bounded away from zero. The latter restriction is required because the noise variance in model (2.2) equals 0 if  $\lambda_i$  is equal to 0. Since we normalize our test statistics by an estimate of the noise variance as detailed in Section 3, we need this variance and thus the functions  $\lambda_i$  to be bounded away from zero.
- (C2) assumes the noise terms  $\eta_{it}$  to fulfill some mild moment conditions and to be independent both across countries  $i$  and time  $t$ . We require the independence assumptions of (C2) in order to apply the Gaussian approximation results for hyperrectangles from Chernozhukov et al. (2017) in our proofs. We in particular need independence across  $t$ , but it would in principle be possible to allow for certain forms of dependence across  $i$  at the cost of a more complicated test procedure and more involved technical arguments.
- In the current COVID-19 crisis, independence across countries  $i$  seems to be a fairly reasonable assumption due to severe travel restrictions, the closure of borders, etc. Note that this assumption can in principle be tested, for example, with the help of the tests in Blum et al. (1961), Sinha and Wieand (1977) and Bakirov et al. (2006).
- Independence across time  $t$  is more debatable than independence across countries  $i$ , but it is by no means unreasonable in our model framework: The time series process  $\mathcal{X}_i = \{X_{it} : 1 \leq t \leq T\}$  produced by model (2.2) is nonstationary for each  $i$ . Specifically, both the mean  $\mathbb{E}[X_{it}] = \lambda_i(t/T)$  and the variance  $\text{Var}(X_{it}) = \sigma^2 \lambda_i(t/T)$  are time-varying. A well-known fact in the time series literature is that nonstationarities such as a time-varying mean may produce spurious sample autocorrelations (see e.g. Mikosch and Stărică, 2004; Fryzlewicz et al., 2008). Hence, the observed persistence of a time series (captured by the sample autocorrelation function) may be due to nonstationarities rather than real autocorrelations. This insight has led researchers to prefer simple nonstationary models over intricate stationary time series models in some application areas such as finance (see e.g. Mikosch and Stărică, 2000, 2004; Fryzlewicz et al., 2006; Hafner and Linton, 2010). In a similar vein, our model accounts for the persistence in the observed time series  $\mathcal{X}_i$  via nonstationarities rather than autocorrelations in the error terms.

### 3. The multiscale test

Let  $S \subseteq \{(i, j) : 1 \leq i < j \leq n\}$  be the set of all pairs of countries  $(i, j)$  whose trend functions  $\lambda_i$  and  $\lambda_j$  we want to compare. Moreover, as already introduced above, let  $\mathcal{F} = \{\mathcal{I}_k : 1 \leq k \leq K\}$  be the family of (rescaled) time intervals under consideration. Finally, write  $\mathcal{M} := S \times \{1, \dots, K\}$  and let  $p := |\mathcal{M}|$  be the cardinality of  $\mathcal{M}$ . In this section, we devise a method to test the null hypothesis  $H_0^{(ijk)}$  simultaneously for all pairs of countries  $(i, j) \in S$  and all time intervals  $\mathcal{I}_k \in \mathcal{F}$ , that is, for all  $(i, j, k) \in \mathcal{M}$ . The value  $p = |\mathcal{M}|$  is the dimensionality of the simultaneous test problem we are dealing with. It amounts to the number of tests that we carry out simultaneously. As shown by our theoretical results in the Appendix,  $p$  may grow as a polynomial  $T^\gamma$  of the time series length  $T$ , where the exponent  $\gamma$  depends on the number of error moments  $\theta$  defined in (C2) and on the minimal length of the rescaled time intervals in the family  $\mathcal{F}$ . Precise conditions on the exponent  $\gamma$  are given in the statement of Theorem A.1. These conditions show that  $\gamma$  can be very large provided that the error terms have sufficiently many moments  $\theta$ . Consequently,  $p$  may be much larger than the time series length  $T$ , which means that the simultaneous test problem under consideration can be very high-dimensional.

### 3.1. Construction of the test statistics

A statistic to test the hypothesis  $H_0^{(ijk)}$  for a given triple  $(i, j, k)$  can be constructed as follows. To start with, we consider the expression

$$\hat{s}_{ijk,T} = \frac{1}{Th_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (X_{it} - X_{jt}),$$

where  $h_k$  is the length of the time interval  $\mathcal{I}_k$ ,  $\mathbf{1}(\cdot)$  denotes the indicator function and  $\mathbf{1}(t/T \in \mathcal{I}_k)$  can be regarded as a rectangular kernel weight. Inserting the model equation (2.2) into the definition of  $\hat{s}_{ijk,T}$  yields that  $\hat{s}_{ijk,T} = \Delta_{ijk,T} + R_{ijk,T}$ , where

$$\Delta_{ijk,T} = \frac{1}{Th_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) \left\{ \lambda_i\left(\frac{t}{T}\right) - \lambda_j\left(\frac{t}{T}\right) \right\}$$

is the average distance between the functions  $\lambda_i$  and  $\lambda_j$  on the interval  $\mathcal{I}_k$  and

$$R_{ijk,T} = \frac{1}{Th_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) \sigma \left\{ \sqrt{\lambda_i\left(\frac{t}{T}\right)} \eta_{it} - \sqrt{\lambda_j\left(\frac{t}{T}\right)} \eta_{jt} \right\}$$

is a remainder term that is asymptotically negligible in the following sense: A simple application of the law of large numbers for a fixed  $i$  gives that  $(Th_k)^{-1} \sum_{t=1}^T \mathbf{1}(t/T \in \mathcal{I}_k) \sqrt{\lambda_i(t/T)} \eta_{it} = o_p(1)$ , which in turn implies that  $R_{ijk,T} = o_p(1)$ . Hence, for any fixed triple  $(i, j, k)$ , we obtain that

$$\hat{s}_{ijk,T} = \Delta_{ijk,T} + o_p(1),$$

which means that the statistic  $\hat{s}_{ijk,T}$  estimates the average distance  $\Delta_{ijk,T}$  between the functions  $\lambda_i$  and  $\lambda_j$  on the interval  $\mathcal{I}_k$ .

We next have a closer look at the variance of the statistic  $\hat{s}_{ijk,T}$ . Under (C2), the variance of  $\hat{s}_{ijk,T}$  is given by

$$v_{ijk,T}^2 := \text{Var}(\hat{s}_{ijk,T}) = \frac{\sigma^2}{(Th_k)^2} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) \left\{ \lambda_i\left(\frac{t}{T}\right) + \lambda_j\left(\frac{t}{T}\right) \right\}$$

and can be estimated by

$$\hat{v}_{ijk,T}^2 = \frac{\hat{\sigma}^2}{(Th_k)^2} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (X_{it} + X_{jt}).$$

Here,  $\hat{\sigma}^2$  is an estimator of  $\sigma^2$  which is defined as  $\hat{\sigma}^2 = |C|^{-1} \sum_{i \in C} \hat{\sigma}_i^2$ , where  $C = \{\ell : \ell = i \text{ or } \ell = j \text{ for some } (i, j) \in S\}$  denotes the set of countries that are taken into account by our test and

$$\hat{\sigma}_i^2 = \frac{\sum_{t=2}^T (X_{it} - X_{it-1})^2}{2 \sum_{t=1}^T X_{it}}$$

for each country  $i$ . The idea behind the estimator  $\hat{\sigma}_i^2$  is as follows: We can write

$$X_{it} - X_{it-1} = \sigma \sqrt{\lambda_i\left(\frac{t}{T}\right)} (\eta_{it} - \eta_{it-1}) + r_{it} \tag{3.1}$$

with

$$r_{it} = \lambda_i\left(\frac{t}{T}\right) - \lambda_i\left(\frac{t-1}{T}\right) + \sigma \left\{ \sqrt{\lambda_i\left(\frac{t}{T}\right)} - \sqrt{\lambda_i\left(\frac{t-1}{T}\right)} \right\} \eta_{it-1}.$$

By the triangle inequality and since  $\lambda_i$  is Lipschitz continuous, we have that

$$\begin{aligned} |r_{it}| &\leq \left| \lambda_i\left(\frac{t}{T}\right) - \lambda_i\left(\frac{t-1}{T}\right) \right| + \sigma \left| \sqrt{\lambda_i\left(\frac{t}{T}\right)} - \sqrt{\lambda_i\left(\frac{t-1}{T}\right)} \right| |\eta_{it-1}| \\ &\leq \frac{L}{T} + \frac{\sigma L}{2\sqrt{\lambda_{\min} T}} |\eta_{it-1}| \leq \frac{C(1 + |\eta_{it-1}|)}{T}, \end{aligned} \tag{3.2}$$

where  $C = \max\{L, \sigma L/(2\sqrt{\lambda_{\min}})\}$  with  $L$  and  $\lambda_{\min}$  defined in (C1). From (3.1) and (3.2), we can infer that

$$\frac{1}{T} \sum_{t=2}^T (X_{it} - X_{it-1})^2 = 2\sigma^2 \left\{ \frac{1}{T} \sum_{t=2}^T \lambda_i\left(\frac{t}{T}\right) \right\} + o_p(1).$$

Moreover, since  $T^{-1} \sum_{t=1}^T X_{it} = T^{-1} \sum_{t=1}^T \lambda_i(t/T) + o_p(1)$ , we get that  $\hat{\sigma}_i^2 = \sigma^2 + o_p(1)$  for any fixed  $i$ . In Lemma S.1 of the Supplement, we further show that  $\hat{\sigma}^2 = \sigma^2 + o_p(1)$  under our regularity conditions. Hence,  $\hat{\sigma}^2$  is a consistent estimator of  $\sigma^2$ .

We now replace the statistic  $\hat{s}_{ijk,T}$  by a normalized version whose variance is approximately equal to 1. To achieve this, we simply divide  $\hat{s}_{ijk,T}$  by its estimated standard deviation  $\hat{v}_{ijk,T}$ . This results in the expression

$$\hat{\psi}_{ijk,T} := \frac{\hat{s}_{ijk,T}}{\hat{v}_{ijk,T}} = \frac{\sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k)(X_{it} - X_{jt})}{\hat{\sigma} \{ \sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k)(X_{it} + X_{jt}) \}^{1/2}}, \tag{3.3}$$

which serves as our test statistic of the hypothesis  $H_0^{(ijk)}$ . In addition to  $\hat{\psi}_{ijk,T}$ , we introduce the auxiliary statistic

$$\hat{\psi}_{ijk,T}^0 = \frac{\sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k) \sigma \bar{\lambda}_{ij}^{-1/2}(\frac{t}{T})(\eta_{it} - \eta_{jt})}{\hat{\sigma} \{ \sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k)(X_{it} + X_{jt}) \}^{1/2}} \tag{3.4}$$

with  $\bar{\lambda}_{ij}(u) = \{\lambda_i(u) + \lambda_j(u)\}/2$ , which by construction is identical to  $\hat{\psi}_{ijk,T}$  under  $H_0^{(ijk)}$ . This auxiliary statistic is needed to define the critical values of our multiscale test in what follows.

### 3.2. Construction of the test

Our multiscale test is carried out as follows: For a given significance level  $\alpha \in (0, 1)$  and each  $(i, j, k) \in \mathcal{M}$ , we reject  $H_0^{(ijk)}$  if

$$|\hat{\psi}_{ijk,T}| > c_{ijk,T}(\alpha),$$

where  $c_{ijk,T}(\alpha)$  is the critical value for the  $(i, j, k)$ th test problem. The critical values  $c_{ijk,T}(\alpha)$  are chosen such that the familywise error rate (FWER) is controlled at level  $\alpha$ , which is defined as the probability of wrongly rejecting  $H_0^{(ijk)}$  for at least one  $(i, j, k)$ . More formally speaking, for a given significance level  $\alpha \in (0, 1)$ , the FWER is

$$\begin{aligned} \text{FWER}(\alpha) &= \mathbb{P}\left(\exists(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| > c_{ijk,T}(\alpha)\right) \\ &= 1 - \mathbb{P}\left(\forall(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| \leq c_{ijk,T}(\alpha)\right), \end{aligned}$$

where  $\mathcal{M}_0 \subseteq \mathcal{M}$  is the set of triples  $(i, j, k)$  for which  $H_0^{(ijk)}$  holds true.

There are different ways to construct critical values  $c_{ijk,T}(\alpha)$  that ensure control of the FWER at level  $\alpha$ . In the traditional approach, the same critical value  $c_T(\alpha) = c_{ijk,T}(\alpha)$  is used for all  $(i, j, k)$ . In this case, controlling the FWER at the level  $\alpha$  requires to determine the critical value  $c_T(\alpha)$  such that

$$\begin{aligned} \text{FWER}(\alpha) &= 1 - \mathbb{P}\left(\forall(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| \leq c_T(\alpha)\right) \\ &= 1 - \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}_0} |\hat{\psi}_{ijk,T}| \leq c_T(\alpha)\right) \leq \alpha. \end{aligned} \tag{3.5}$$

This can be achieved by choosing  $c_T(\alpha)$  as the  $(1 - \alpha)$ -quantile of the statistic

$$\tilde{\psi}_T = \max_{(i,j,k) \in \mathcal{M}} |\hat{\psi}_{ijk,T}^0|,$$

where the auxiliary statistic  $\hat{\psi}_{ijk,T}^0$  was introduced in (3.4) and is equal to  $\hat{\psi}_{ijk,T}$  under the null  $H_0^{(ijk)}$  by construction.<sup>2</sup>

A more modern approach assigns different critical values  $c_{ijk,T}(\alpha)$  to the test problems  $(i, j, k)$ . In particular, the critical value for the hypothesis  $H_0^{(ijk)}$  is allowed to depend on the length  $h_k$  of the time interval  $\mathcal{I}_k$ , that is, on the scale of the test problem. A general approach to construct scale-dependent critical values was pioneered by [Dümbgen and Spokoiny \(2001\)](#) and has been used in many other studies since then; see e.g. [Rohde \(2008\)](#), [Dümbgen and Walther \(2008\)](#), [Rufibach and Walther \(2010\)](#), [Schmidt-Hieber et al. \(2013\)](#), [Eckle et al. \(2017\)](#) and [Dunker et al. \(2019\)](#). In our context, the approach of [Dümbgen and Spokoiny \(2001\)](#) leads to the critical values

$$c_{ijk,T}(\alpha) = c_T(\alpha, h_k) := b_k + q_T(\alpha)/a_k,$$

where  $a_k = \{\log(e/h_k)\}^{1/2} / \log \log(e/h_k)$  and  $b_k = \sqrt{2 \log(1/h_k)}$  are scale-dependent constants and the quantity  $q_T(\alpha)$  is determined by the following consideration: Since

$$\begin{aligned} \text{FWER}(\alpha) &= \mathbb{P}\left(\exists(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| > c_T(\alpha, h_k)\right) \\ &= 1 - \mathbb{P}\left(\forall(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| \leq c_T(\alpha, h_k)\right) \end{aligned}$$

<sup>2</sup> Note that both the statistic  $\tilde{\psi}_T$  and the quantile  $c_T(\alpha)$  depend on the dimensionality  $p$  of the test problem in general. To keep the notation simple, we however suppress this dependence throughout the paper. We use the same convention for all other quantities that are defined in the sequel.

$$\begin{aligned}
 &= 1 - \mathbb{P}\left(\forall(i, j, k) \in \mathcal{M}_0 : a_k(|\hat{\psi}_{ijk,T}| - b_k) \leq q_T(\alpha)\right) \\
 &= 1 - \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}_0} a_k(|\hat{\psi}_{ijk,T}| - b_k) \leq q_T(\alpha)\right), \tag{3.6}
 \end{aligned}$$

we need to choose the quantity  $q_T(\alpha)$  as the  $(1 - \alpha)$ -quantile of the statistic

$$\hat{\Psi}_T = \max_{(i,j,k) \in \mathcal{M}} a_k(|\hat{\psi}_{ijk,T}^0| - b_k)$$

in order to ensure control of the FWER at level  $\alpha$ . Comparing (3.6) with (3.5), the current approach can be seen to differ from the traditional one in the following respect: the maximum statistic  $\hat{\Psi}_T$  is replaced by the rescaled version  $\hat{\Psi}_T$  which re-weights the individual statistics  $\hat{\psi}_{ijk,T}^0$  by the scale-dependent constants  $a_k$  and  $b_k$ . As demonstrated above, this translates into scale-dependent critical values  $c_{ijk,T}(\alpha) = c_T(\alpha, h_k)$ .

Our theory allows us to work with both the traditional choice  $c_{ijk,T}(\alpha) = c_T(\alpha)$  and the more modern, scale-dependent choice  $c_{ijk,T}(\alpha) = c_T(\alpha, h_k)$ . Since the latter choice produces a test approach with better theoretical properties in general (see Dümmbgen and Spokoiny, 2001), we restrict attention to the critical values  $c_T(\alpha, h_k)$  in the sequel. There is one complication we need to deal with: As the quantiles  $q_T(\alpha)$  are not known in practice, we cannot compute the critical values  $c_T(\alpha, h_k)$  exactly in practice but need to approximate them. This can be achieved as follows: Under appropriate regularity conditions, it can be shown that

$$\hat{\psi}_{ijk,T}^0 \approx \frac{1}{\sqrt{2T}h_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) \{\eta_{it} - \eta_{jt}\}.$$

A Gaussian version of the statistic displayed on the right-hand side above is given by

$$\phi_{ijk,T} = \frac{1}{\sqrt{2T}h_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) \{Z_{it} - Z_{jt}\},$$

where  $Z_{it}$  are independent standard normal random variables for  $1 \leq t \leq T$  and  $1 \leq i \leq n$ . Hence, the statistic

$$\Phi_T = \max_{(i,j,k) \in \mathcal{M}} a_k(|\phi_{ijk,T}| - b_k)$$

can be regarded as a Gaussian version of the statistic  $\hat{\Psi}_T$ . We approximate the unknown quantile  $q_T(\alpha)$  by the  $(1 - \alpha)$ -quantile  $q_{T,\text{Gauss}}(\alpha)$  of  $\Phi_T$ , which can be computed (approximately) by Monte Carlo simulations and can thus be treated as known.

To summarize, we propose the following procedure to simultaneously test the hypothesis  $H_0^{(ijk)}$  for all  $(i, j, k) \in \mathcal{M}$  at the significance level  $\alpha \in (0, 1)$ :

$$\text{For each } (i, j, k) \in \mathcal{M}, \text{ reject } H_0^{(ijk)} \text{ if } |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k), \tag{3.7}$$

where  $c_{T,\text{Gauss}}(\alpha, h_k) = b_k + q_{T,\text{Gauss}}(\alpha)/a_k$  with  $a_k = \{\log(e/h_k)\}^{1/2} / \log \log(e/h_k)$  and  $b_k = \sqrt{2 \log(1/h_k)}$ .

### 3.3. Formal properties of the test

In Theorem A.1 of the Appendix, we prove that under appropriate regularity conditions, the test defined in (3.7) (asymptotically) controls the familywise error rate  $\text{FWER}(\alpha)$  for each pre-specified significance level  $\alpha$ . As shown in Corollary A.1, this has the following implication:

$$\mathbb{P}\left(\forall(i, j, k) \in \mathcal{R} : (i, j, k) \notin \mathcal{M}_0\right) \geq 1 - \alpha + o(1), \tag{3.8}$$

where  $\mathcal{R} = \{(i, j, k) \in \mathcal{M} \text{ with } |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\}$  is the set of triples  $(i, j, k)$  for which our test rejects the null  $H_0^{(ijk)}$  and  $\mathcal{M}_0$  is the set of triples  $(i, j, k)$  for which  $H_0^{(ijk)}$  holds true. Verbally, (3.8) can be expressed as follows:

$$\text{With (asymptotic) probability at least } 1 - \alpha, \text{ the null hypothesis } H_0^{(ijk)} \text{ is violated for all } (i, j, k) \in \mathcal{M} \text{ for which the test rejects } H_0^{(ijk)}. \tag{3.9}$$

In other words:

$$\text{With (asymptotic) probability at least } 1 - \alpha, \text{ the functions } \lambda_i \text{ and } \lambda_j \text{ differ on the interval } \mathcal{I}_k \text{ for all } (i, j, k) \in \mathcal{M} \text{ for which the test rejects } H_0^{(ijk)}. \tag{3.10}$$

Hence, the test allows us to make simultaneous confidence statements (a) about which pairs of countries  $(i, j)$  have different trend functions and (b) about where, that is, in which time intervals  $\mathcal{I}_k$  the functions differ.

According to (3.8), our test does not produce any false positives with high probability. In addition, we would like the test not to produce any false negatives either. Put differently, we would like the test to have high power against deviations from the null. In Proposition A.1 in the Appendix, we derive the power properties of the test against a certain class of local



alternatives. To summarize, we show the following: Let  $\lambda_i = \lambda_{i,T}$  and  $\lambda_j = \lambda_{j,T}$  be functions whose difference  $\lambda_{i,T} - \lambda_{j,T}$  converges to zero as  $T \rightarrow \infty$ . Moreover, let  $\mathcal{M}_1$  be the set of triples  $(i, j, k)$  such that either

$$\lambda_{i,T}(w) - \lambda_{j,T}(w) \geq \kappa_T \sqrt{\log T / (Th_k)} \quad \text{for all } w \in \mathcal{I}_k \tag{3.11}$$

or

$$\lambda_{j,T}(w) - \lambda_{i,T}(w) \geq \kappa_T \sqrt{\log T / (Th_k)} \quad \text{for all } w \in \mathcal{I}_k, \tag{3.12}$$

where  $\{\kappa_T\}$  is any sequence of positive numbers which diverges at a faster rate than  $\{\sqrt{\log T} \sqrt{\log \log T} / \log \log \log T\}$ . According to Proposition A.1, it holds that

$$\mathbb{P}\left(\forall (i, j, k) \in \mathcal{M}_1 : |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\right) = 1 - o(1). \tag{3.13}$$

Hence, the test detects any local deviation from the null of the form (3.11) or (3.12) with probability tending to 1.

### 3.4. Implementation of the test in practice

For a given significance level  $\alpha \in (0, 1)$ , the test procedure defined in (3.7) is implemented as follows in practice:

- Step 1. Compute the quantile  $q_{T,\text{Gauss}}(\alpha)$  by Monte Carlo simulations. Specifically, draw a large number  $N$  (say  $N = 5000$ ) of samples of independent standard normal random variables  $\{Z_{it}^{(\ell)} : 1 \leq t \leq T, 1 \leq i \leq n\}$  for  $1 \leq \ell \leq N$ . Compute the value  $\Phi_T^{(\ell)}$  of the Gaussian statistic  $\Phi_T$  for each sample  $\ell$  and calculate the empirical  $(1 - \alpha)$ -quantile  $\hat{q}_{T,\text{Gauss}}(\alpha)$  from the values  $\{\Phi_T^{(\ell)} : 1 \leq \ell \leq N\}$ . Use  $\hat{q}_{T,\text{Gauss}}(\alpha)$  as an approximation of the quantile  $q_{T,\text{Gauss}}(\alpha)$ .
- Step 2. Compute the critical values  $c_{T,\text{Gauss}}(\alpha, h_k)$  for  $1 \leq k \leq K$  based on the approximation  $\hat{q}_{T,\text{Gauss}}(\alpha)$ .
- Step 3. Carry out the test for each  $(i, j, k) \in \mathcal{M}$  and store the test results in the variable  $r_{ijk,T} = \mathbf{1}(|\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k))$  for each  $(i, j, k) \in \mathcal{M}$ , that is, let  $r_{ijk,T} = 1$  if the hypothesis  $H_0^{(ijk)}$  is rejected and  $r_{ijk,T} = 0$  otherwise.

To graphically present the test results, we produce a plot for each pair of countries  $(i, j) \in \mathcal{S}$  that shows the intervals  $\mathcal{I}_k$  for which the test rejects the null  $H_0^{(ijk)}$ , that is, the intervals in the set  $\mathcal{F}_{\text{reject}}(i, j) = \{\mathcal{I}_k \in \mathcal{F} : r_{ijk,T} = 1\}$ . The plot is designed such that it graphically highlights the subset of intervals  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j) = \{\mathcal{I}_k \in \mathcal{F}_{\text{reject}}(i, j) : \text{there exists no } \mathcal{I}_{k'} \in \mathcal{F}_{\text{reject}}(i, j) \text{ with } \mathcal{I}_{k'} \subset \mathcal{I}_k\}$ . The elements of  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j)$  are called minimal intervals. By definition, there is no other interval  $\mathcal{I}_{k'}$  in  $\mathcal{F}_{\text{reject}}(i, j)$  which is a proper subset of a minimal interval  $\mathcal{I}_k$ . Hence, the minimal intervals can be regarded as those intervals in  $\mathcal{F}_{\text{reject}}(i, j)$  which are most informative about the precise location of the differences between the trends  $\lambda_i$  and  $\lambda_j$ . In Section 4, we use the graphical device just described to present the test results of our empirical application; see panels (d) in Figs. 3–6.

According to (3.8), we can make the following simultaneous confidence statement about the intervals in  $\mathcal{F}_{\text{reject}}(i, j)$  for  $(i, j) \in \mathcal{S}$ :

$$\text{With (asymptotic) probability at least } 1 - \alpha, \text{ it holds that for every pair of countries } (i, j) \in \mathcal{S}, \text{ the functions } \lambda_i \text{ and } \lambda_j \text{ differ on each interval in } \mathcal{F}_{\text{reject}}(i, j). \tag{3.14}$$

Hence, we can claim with statistical confidence at least  $1 - \alpha$  that the functions  $\lambda_i$  and  $\lambda_j$  differ on each time interval which is depicted in the plots of our graphical device. Since  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j) \subseteq \mathcal{F}_{\text{reject}}(i, j)$  for any  $(i, j) \in \mathcal{S}$ , the confidence statement (3.14) trivially remains to hold true when  $\mathcal{F}_{\text{reject}}(i, j)$  is replaced by  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j)$ .

The graphical device described above is of course not the only way to present the test results. Another object which is helpful in summarizing the test results for a given pair of countries  $(i, j)$  is the union of minimal intervals  $U_{ij} = \cup_{\mathcal{I}_k \in \mathcal{F}_{\text{reject}}^{\text{min}}(i, j)} \mathcal{I}_k$ . One can formally show that the union  $U_{ij}$  is closely related to the set  $U_{ij}^* = \{u \in [0, 1] : \lambda_i(u) \neq \lambda_j(u)\}$  of time points where the two functions  $\lambda_i$  and  $\lambda_j$  differ from each other. For a precise mathematical statement and the technical details, we refer to Lemma S.2 of the Supplementary Material.

## 4. Empirical application to COVID-19 data

We now use our test to analyse the outbreak patterns of the COVID-19 epidemic. We proceed in two steps. In Section 4.1, we assess the finite sample performance of our test by Monte-Carlo experiments. Specifically, we run a series of experiments which show that the test controls the FWER at level  $\alpha$  as predicted by the theory and that it has good power properties. In Section 4.2, we then apply the test to a sample of COVID-19 data from different European countries. Our multiscale test is implemented in the R package `multiscale`, available on GitHub at <https://github.com/marina-khi/multiscale>.

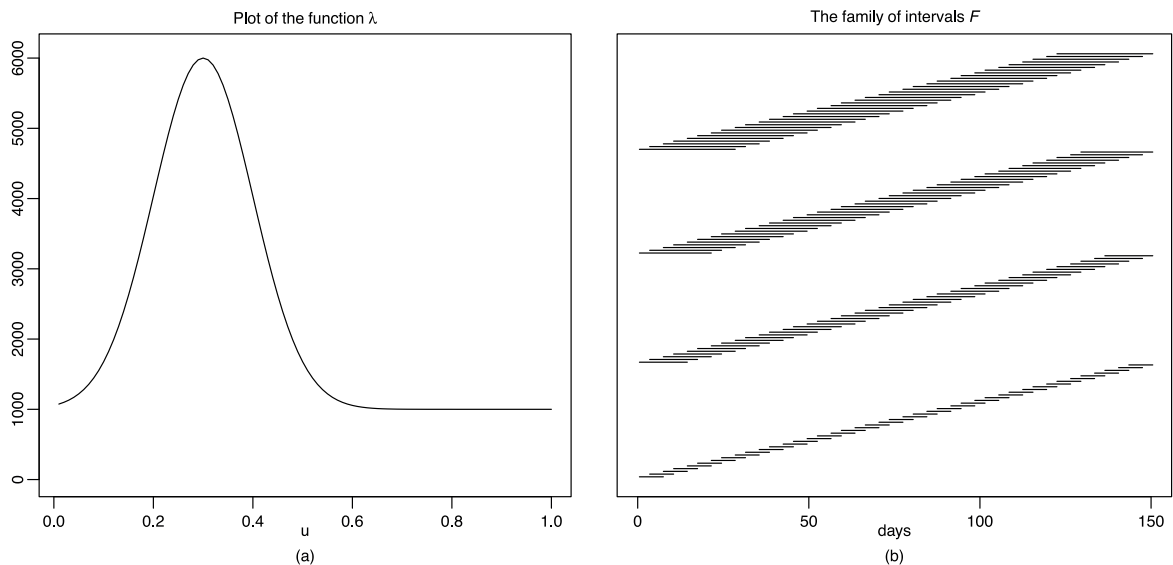


Fig. 1. (a) Plot of the function  $\lambda$ ; (b) plot of the family of intervals  $\mathcal{F}$ .

4.1. Simulation experiments

We simulate count data  $\mathcal{X} = \{X_{it} : 1 \leq i \leq n, 1 \leq t \leq T\}$  by drawing the observations  $X_{it}$  independently from a negative binomial distribution with mean  $\lambda_i(t/T)$  and variance  $\sigma^2 \lambda_i(t/T)$ . By definition,  $X_{it}$  has a negative binomial distribution with parameters  $q$  and  $r$  if  $\mathbb{P}(X_{it} = m) = \Gamma(m+r)/(\Gamma(r)m!)q^r(1-q)^m$  for each  $m \in \mathbb{N} \cup \{0\}$ . Since  $\mathbb{E}[X_{it}] = r(1-q)/q$  and  $\text{Var}(X_{it}) = r(1-q)/q^2$ , we can use the parametrization  $q = 1/\sigma^2$  and  $r = \lambda_i(t/T)/(\sigma^2 - 1)$  to obtain that  $\mathbb{E}[X_{it}] = \lambda_i(t/T)$  and  $\text{Var}(X_{it}) = \sigma^2 \lambda_i(t/T)$ . With this parametrization, the simulated data follow a nonparametric regression model of the form

$$X_{it} = \lambda_i\left(\frac{t}{T}\right) + \sigma \sqrt{\lambda_i\left(\frac{t}{T}\right)} \eta_{it},$$

where the noise variables  $\eta_{it}$  have zero mean and unit variance. The functions  $\lambda_i$  are specified below. The overdispersion parameter is set to  $\sigma = 15$ , which is similar to the estimate  $\hat{\sigma} = 14.82$  obtained in the empirical application of Section 4.2. Robustness checks with  $\sigma = 10$  and  $\sigma = 20$  are provided in the Supplement.

We consider different values for  $T$  and  $n$ , in particular,  $T \in \{100, 250, 500\}$  and  $n \in \{5, 10, 50\}$ . Note that in the application, we have  $T = 150$  and  $n = 5$ . We let  $\mathcal{S} = \{(i, j) : 1 \leq i < j \leq n\}$ , that is, we compare all pairs of countries  $(i, j)$  with  $i < j$ . Moreover, we choose  $\mathcal{F}$  to be a family of time intervals  $\mathcal{I}_k$  with length  $h_k \in \{7/T, 14/T, 21/T, 28/T\}$ . Hence, the intervals in  $\mathcal{F}$  have length either 7, 14, 21 or 28 days (i.e., 1, 2, 3 or 4 weeks). For each length  $h_k$ , we include all intervals that start at days  $t = 1 + 7(j - 1)$  and  $t = 4 + 7(j - 1)$  for  $j = 1, 2, \dots$ . A graphical presentation of the family  $\mathcal{F}$  for  $T = 150$  (as in the application) is given in Fig. 1b. All our simulation experiments are based on  $R = 5000$  simulation runs.

In the first part of the simulation study, we examine whether our test controls the FWER as predicted by the theory. To do so, we assume that the hypothesis  $H_0^{(ijk)}$  holds true for all  $(i, j, k)$  under consideration, which implies that  $\lambda_i = \lambda$  for all  $i$ . We consider the function

$$\lambda(u) = 5000 \exp\left(-\frac{(10u - 3)^2}{2}\right) + 1000, \tag{4.1}$$

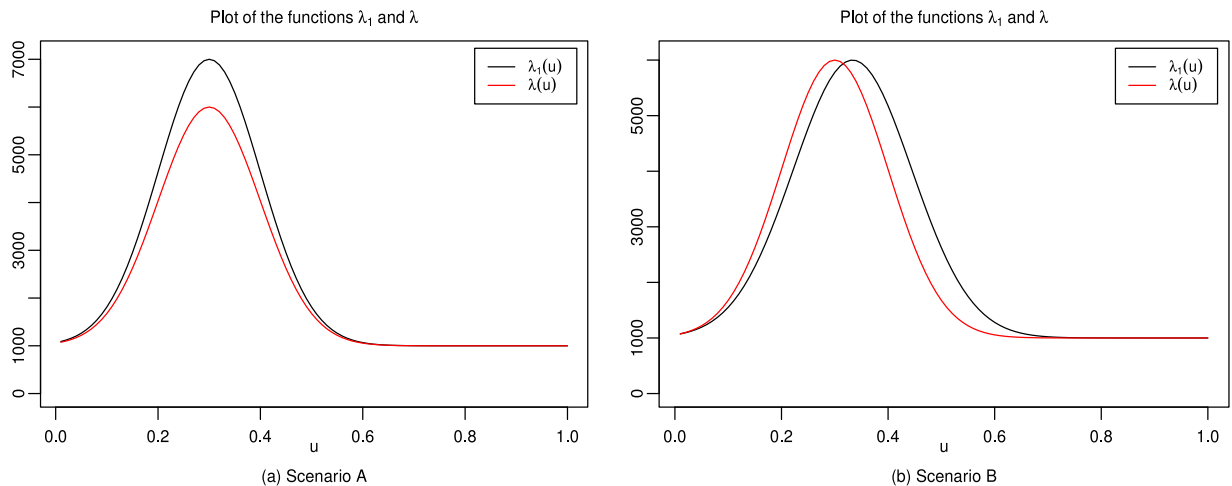
which is similar in shape to some of the estimated trend curves in the application of Section 4.2. A plot of the function  $\lambda$  is provided in Fig. 1a. To evaluate whether the test controls the FWER at level  $\alpha$ , we compare the empirical size of the test with the target  $\alpha$ . The empirical size is computed as the percentage of simulation runs in which the test falsely rejects at least one null hypothesis  $H_0^{(ijk)}$ .

The simulation results are reported in Table 1. As can be seen, the empirical size gives a reasonable approximation to the target  $\alpha$  in all scenarios under investigation, even though the size numbers have a slight downward bias. This bias gets larger as the number of time series  $n$  increases, which reflects the fact that the test problem becomes more difficult for larger  $n$ . Already for  $n = 5$ , the number  $p$  of hypotheses to be tested is quite high, in particular,  $p = 960, 2680, 5560$  for  $T = 100, 250, 500$ . This number increases to  $p = 117\,600, 328\,300, 681\,100$  when  $n = 50$ . Hence, the dimensionality and thus the complexity of the test problem increases considerably as  $n$  gets larger. On first sight, it may seem astonishing



**Table 1**  
Empirical size of the test for different values of  $n$  and  $T$ .

	$n = 5$			$n = 10$			$n = 50$		
	Significance level $\alpha$			Significance level $\alpha$			Significance level $\alpha$		
	0.01	0.05	0.1	0.01	0.05	0.1	0.01	0.05	0.1
$T = 100$	0.011	0.047	0.093	0.010	0.044	0.087	0.008	0.037	0.075
$T = 250$	0.009	0.047	0.091	0.009	0.046	0.087	0.008	0.035	0.069
$T = 500$	0.010	0.044	0.083	0.008	0.048	0.093	0.007	0.035	0.077



**Fig. 2.** Plot of the functions  $\lambda_1$  (black) and  $\lambda$  (red) in the simulation scenarios A and B. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

that the downward bias does not diminish notably as the time series length  $T$  increases. This, however, has a simple explanation: The interval lengths  $h_k$  remain the same (7, 14, 21 or 28 days) as  $T$  increases, which implies that the effective sample size for computing the test statistics  $\hat{\psi}_{ijk,T}$  does not change. To summarize, even though slightly conservative, the test controls the FWER quite accurately in the simulation setting at hand.

In the second part of the simulation study, we investigate the power properties of the test. To do so, we assume that  $\lambda_i = \lambda$  for all  $i > 1$  and that  $\lambda_1 \neq \lambda$ , where  $\lambda$  is defined in (4.1). Hence, only the first mean function  $\lambda_1$  is different from the others. This implies that the hypothesis  $H_0^{(ijk)}$  holds true for all  $(i, j, k)$  with  $i > 1$  and  $j > 1$ , while  $H_0^{(ijk)}$  does not hold true for any pair  $(i, j)$  with either  $i = 1$  or  $j = 1$ . We consider two different simulation scenarios. In Scenario A, the function  $\lambda_1$  has the form

$$\lambda_1(u) = 6000 \exp\left(-\frac{(10u - 3)^2}{2}\right) + 1000$$

and is plotted together with  $\lambda$  in Fig. 2a. As can be seen, the two functions  $\lambda_1$  and  $\lambda$  peak at the same point in time, but the peak of  $\lambda_1$  is higher than that of  $\lambda$ . In Scenario B, we let

$$\lambda_1(u) = 5000 \exp\left(-\frac{(9u - 3)^2}{2}\right) + 1000.$$

Fig. 2b shows that the peaks of  $\lambda_1$  and  $\lambda$  have the same height but are reached at different points in time. To evaluate the power properties of the test in Scenarios A and B, we compute the percentage of simulation runs where the test (i) correctly detects differences between  $\lambda_1$  and at least one of the other mean functions and (ii) does not spuriously detect differences between the other mean functions. Put differently, we calculate the percentage of simulation runs where (i) the set  $\mathcal{F}_{\text{reject}}(1, j)$  is non-empty at least for one  $j \in \{2, \dots, n\}$  and (ii) all other sets  $\mathcal{F}_{\text{reject}}(i, j)$  with  $2 \leq i < j \leq n$  are empty. We call this percentage number the (empirical) power of the test. We thus use the term “power” a bit differently than usual.

The results for Scenario A (see Fig. 2a) are presented in Table 2 and those for Scenario B (see Fig. 2b) in Table 3. As can be seen, the test has substantial power in all the considered simulation settings. It is more powerful in Scenario B than in Scenario A, which is most presumably due to the fact that the differences  $|\lambda_1(u) - \lambda(u)|$  are much larger in Scenario B. Moreover, it is less powerful for larger numbers of time series  $n$ , which reflects the fact that the test problem gets more high-dimensional and thus more difficult as  $n$  increases. As one would expect, the power numbers tend to become larger as the time series length  $T$  and the significance level  $\alpha$  increase. In Scenario B (mostly for  $T = 250$  and  $T = 500$ ), however,

**Table 2**Power of the test for different values of  $n$  and  $T$  in Scenario A.

	$n = 5$			$n = 10$			$n = 50$		
	Significance level $\alpha$			Significance level $\alpha$			Significance level $\alpha$		
	0.01	0.05	0.1	0.01	0.05	0.1	0.01	0.05	0.1
$T = 100$	0.335	0.518	0.597	0.306	0.474	0.545	0.212	0.352	0.418
$T = 250$	0.615	0.790	0.836	0.580	0.764	0.800	0.470	0.648	0.705
$T = 500$	0.736	0.905	0.917	0.738	0.884	0.890	0.636	0.799	0.830

**Table 3**Power of the test for different values of  $n$  and  $T$  in Scenario B.

	$n = 5$			$n = 10$			$n = 50$		
	Significance level $\alpha$			Significance level $\alpha$			Significance level $\alpha$		
	0.01	0.05	0.1	0.01	0.05	0.1	0.01	0.05	0.1
$T = 100$	0.824	0.910	0.903	0.812	0.893	0.890	0.738	0.847	0.857
$T = 250$	0.991	0.972	0.941	0.991	0.960	0.920	0.991	0.965	0.933
$T = 500$	0.997	0.973	0.949	0.995	0.961	0.923	0.996	0.969	0.932

the power numbers drop down a bit as  $\alpha$  gets larger. This reverse dependence can be explained by the way we calculate power: we exclude simulation runs where the test spuriously detects differences between the trends in countries  $i$  and  $j$  with  $i, j > 1$ . The number of spurious findings increases as we make the significance level  $\alpha$  larger, which presumably causes the slight drop in power.

#### 4.2. Analysis of COVID-19 data

The COVID-19 pandemic is one of the most pressing issues at present. The first outbreak occurred in Wuhan, China, in December 2019. On 30 January 2020, the World Health Organization (WHO) declared that the outbreak constitutes a Public Health Emergency of International Concern, and on 11 March 2020, the WHO characterized it as a pandemic. As of 2 February 2021, more than 102 million cases of COVID-19 infections have been reported worldwide, resulting in more than 2 million deaths.

There are many open questions surrounding the current COVID-19 pandemic. A question which is particularly relevant for governments and policy makers is whether the pandemic has developed similarly in different countries or whether there are notable differences. Identifying these differences may give some insight into which government policies have been more effective in containing the virus than others. In what follows, we use our multiscale test to compare the development of COVID-19 in several European countries. It is important to emphasize that our test allows to identify differences in the development of the epidemic across countries in a statistically rigorous way, but it does not tell what causes these differences. By distinguishing statistically significant differences from artefacts of the sampling noise, the test provides the basis for a further investigation into the causes. Such an investigation, however, presumably goes beyond a mere statistical analysis.

##### 4.2.1. Data

We analyse data from five European countries: Germany, Italy, Spain, France and the United Kingdom. For each country  $i$ , we observe a time series  $\mathcal{X}_i = \{X_{it} : 1 \leq t \leq T\}$ , where  $X_{it}$  is the number of newly confirmed COVID-19 cases in country  $i$  on day  $t$ . The data are freely available on the homepage of the European Center for Disease Prevention and Control (<https://www.ecdc.europa.eu>) and were downloaded on 2 February 2021.<sup>3</sup> As already mentioned in the Introduction, we take the first Monday after reaching 100 confirmed cases in each country as the starting date  $t = 1$ . Beginning the time series of each country on the day when that country reached 100 confirmed cases is a common way of “normalizing” the data (see e.g. [Cohen and Kupferschmidt, 2020](#)). Additionally aligning the data by Monday allows to take care of possible weekly cycles in the data which are produced by delays in reporting new cases over the weekend. The time series length  $T$  is taken to be equal to 150, which covers the first wave of the pandemic in all of the considered countries. The resulting dataset thus consists of  $n = 5$  time series, each with  $T = 150$  observations. Some of the time series contain negative values which we replaced by 0. Overall, this resulted in 4 replacements. Plots of the observed time series are presented in the upper panels (a) of [Figs. 3–6](#). As a robustness check, we have repeated the data analysis for the longer time span  $T = 200$ . The results are reported in Section S.4 of the Supplement.

To interpret the results produced by our multiscale test, we consider the Government Response Index (GRI) from the Oxford COVID-19 Government Response Tracker (OxCGRT) ([Hale et al., 2020b](#)). The GRI measures how severe the actions

<sup>3</sup> ECDC switched to a weekly reporting schedule for the COVID-19 situation on 17 December 2020. Hence, all daily updates have been discontinued from 14 December. The downloaded daily data set presents historical data until 14 December 2020.

are that are taken by a country's government to contain the virus. It is calculated based on several common government policies such as school closures and travel restrictions. The GRI ranges from 0 to 100, with 0 corresponding to no response from the government at all and 100 corresponding to full lockdown, closure of schools and workplaces, ban on travelling, etc. Detailed information on the collection of the data for government responses and the methodology for calculating the GRI is provided in Hale et al. (2020a). Plots of the GRI time series are given in panels (c) of Figs. 3–6.

#### 4.2.2. Test results

We assume that the data  $X_{it}$  of each country  $i$  in our sample follow the nonparametric trend model

$$X_{it} = \lambda_i\left(\frac{t}{T}\right) + \sigma \sqrt{\lambda_i\left(\frac{t}{T}\right)} \eta_{it},$$

which was introduced in Eq. (2.2). The overdispersion parameter  $\sigma$  is estimated by the procedure described in Section 3.1, which yields the estimate  $\hat{\sigma} = 14.82$ . 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 4.1. In particular, we let  $\mathcal{S} = \{(i, j) : 1 \leq i < j \leq 5\}$ , that is, we compare all pairs of countries  $(i, j)$  with  $i < j$ , and we choose  $\mathcal{F}$  to be the family of time intervals plotted in Fig. 1b. Hence, all intervals in  $\mathcal{F}$  have length either 7, 14, 21 or 28 days.

With the help of our multiscale method, we simultaneously test the null hypothesis  $H_0^{(ijk)}$  that  $\lambda_i = \lambda_j$  on the interval  $\mathcal{I}_k$  for each  $(i, j, k) \in \mathcal{M}$ . The results are presented in Figs. 3–6, each figure comparing a specific pair of countries  $(i, j)$  from our sample. For the sake of brevity, we only show the results for the pairwise comparisons of Germany with each of the four other countries. The remaining figures can be found in Section S.3 of the Supplementary Material. Each figure splits into four panels (a)–(d). Panel (a) shows the observed time series for the two countries  $i$  and  $j$  that are compared. Panel (b) presents smoothed versions of the time series from (a), that is, it shows nonparametric kernel estimates (specifically, Nadaraya–Watson estimates) of the two trend functions  $\lambda_i$  and  $\lambda_j$ , where the bandwidth is set to 7 days and a rectangular kernel is used. Panel (c) displays the Government Response Index (GRI) of the two countries. Finally, panel (d) presents the results produced by our test: it depicts in grey the set  $\mathcal{F}_{\text{reject}}(i, j)$  of all the intervals  $\mathcal{I}_k$  for which the test rejects the null  $H_0^{(ijk)}$ . The minimal intervals in the subset  $\mathcal{F}_{\text{reject}}^{\min}(i, j)$  are highlighted by a black frame. Note that according to (3.8), we can make the following simultaneous confidence statement about the intervals plotted in panels (d) of Figs. 3–6: we can claim, with confidence of about 95%, that there is a difference between the functions  $\lambda_i$  and  $\lambda_j$  on each of these intervals.

We now have a closer look at the results in Figs. 3–6. Fig. 3 presents the comparison of Germany with Italy. The two time series of daily new cases in panel (a) can be seen to be very similar until approximately day 40. Thereafter, the German time series appears to trend downwards more strongly than the Italian one. The smoothed data in panel (b) give a similar visual impression: the kernel estimates of the German and Italian trend curves  $\lambda_i$  and  $\lambda_j$  are very close to each other until approximately day 40 but then start to differ. It is however not clear whether the differences between the two curve estimates reflect differences between the underlying trend curves or whether these are mere artefacts of sampling noise. Our test allows to clarify this issue. Inspecting panel (d), we see that the test detects significant differences between the trend curves in the time period between day 36 and 91. However, it does not find any significant differences up to day 36. Taken together, our results provide evidence that the epidemic developed very similarly in Germany and Italy until a peak was reached around day 40. Thereafter, however, the German time series exhibits a significantly stronger downward trend than the Italian one.

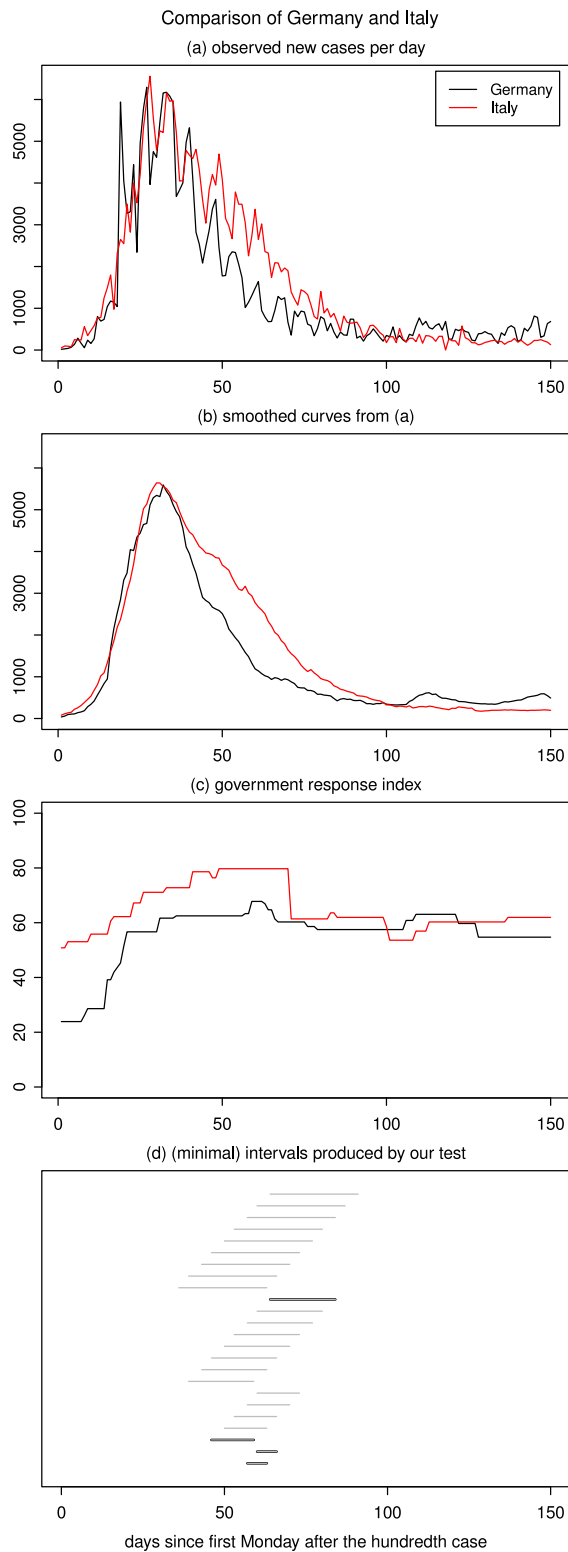
A quite different picture arises when Germany is compared with Spain and France. As can be seen in Figs. 4 and 5, the test detects significant differences between the German trend and the trends in Spain and France up to (approximately) day 50. This indicates that the time trends evolved differently during the outbreak of the crisis. However, the test does not find any differences in the time period between (approximately) days 50 and 120. The trends thus appear to decrease in more or less the same fashion after the first peak was reached. As can be seen in Fig. 4, the test detects additional differences between the German and Spanish trends after day 120. This reflects the fact that the number of daily new cases in Spain picked up again after day 120, foreshadowing the second wave, whereas the numbers in Germany were still quite stable.

Finally, the comparison of Germany with the UK in Fig. 6 reveals significant differences between the time trends in the period from (approximately) day 40 to 120. Similar to the comparison with Italy in Fig. 3, this indicates that the trend decayed in a different fashion in Germany than in the UK after a first peak was reached. However, we do not find any significant differences between the trends during the onset of the pandemic.

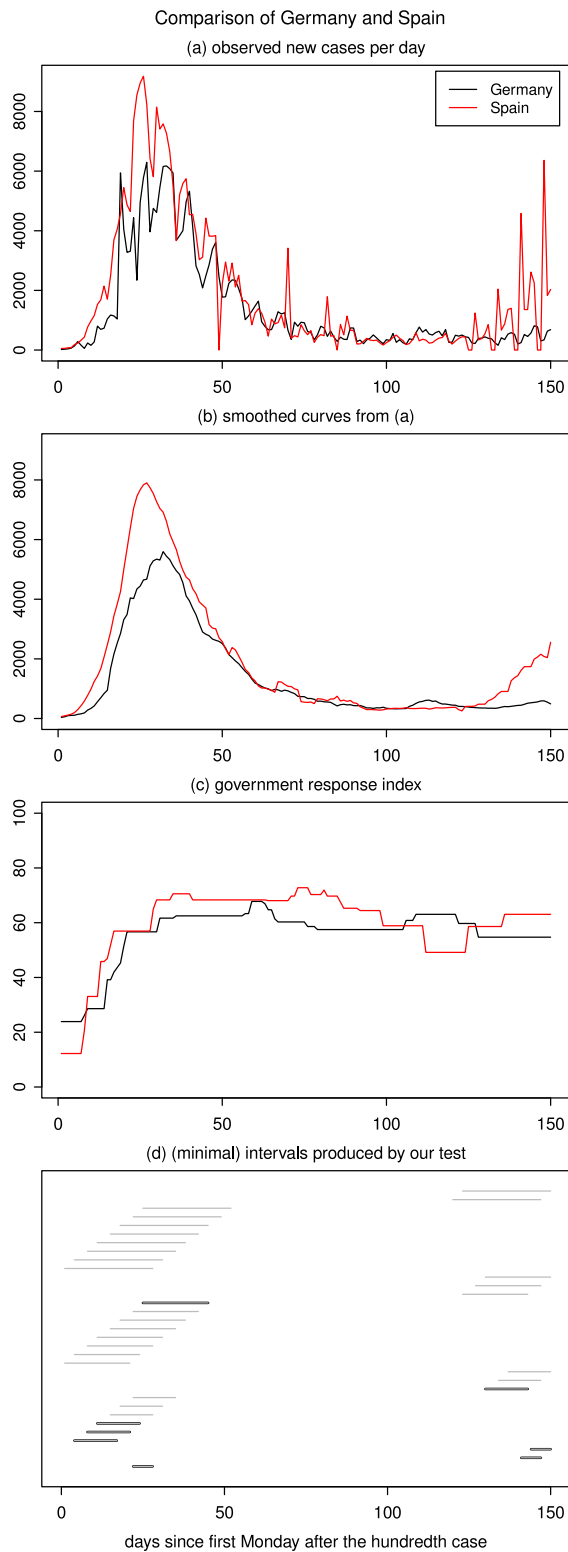
#### 4.2.3. Discussion

Having identified significant differences between the epidemic trends in the five countries under consideration, one may ask next what are the causes of these differences. As already mentioned at the beginning of this section, this question cannot be answered by our test. Rather, a further analysis which presumably goes beyond pure statistics is needed to shed some light on it. We here do not attempt to provide any answers. We merely discuss some observations which become apparent upon considering our test results in the light of the Government Response Index (GRI). For reasons of brevity, we focus on the comparison of Germany with Italy and Spain in Figs. 3 and 4.

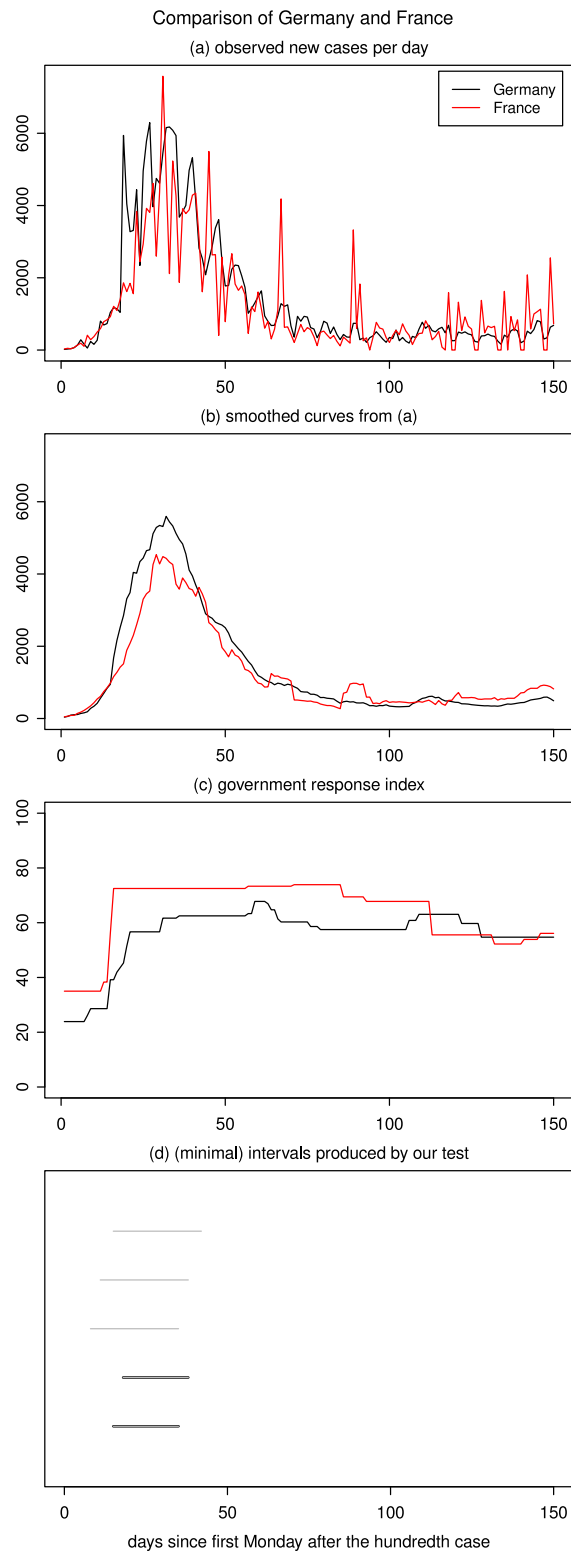
According to our test results in Fig. 4, there are significant differences between the trends in Germany and Spain during the onset of the epidemic up to about day 50, with Spain having more new cases of infections than Germany on most



**Fig. 3.** Test results for the comparison of Germany and Italy. Note: Panel (a) shows the two observed time series, panel (b) smoothed versions of the time series, and panel (c) the corresponding Government Response Index (GRI). Panel (d) depicts the set of intervals  $\mathcal{F}_{\text{reject}}(i, j)$  in grey and the subset of minimal intervals  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j)$  with a black frame.

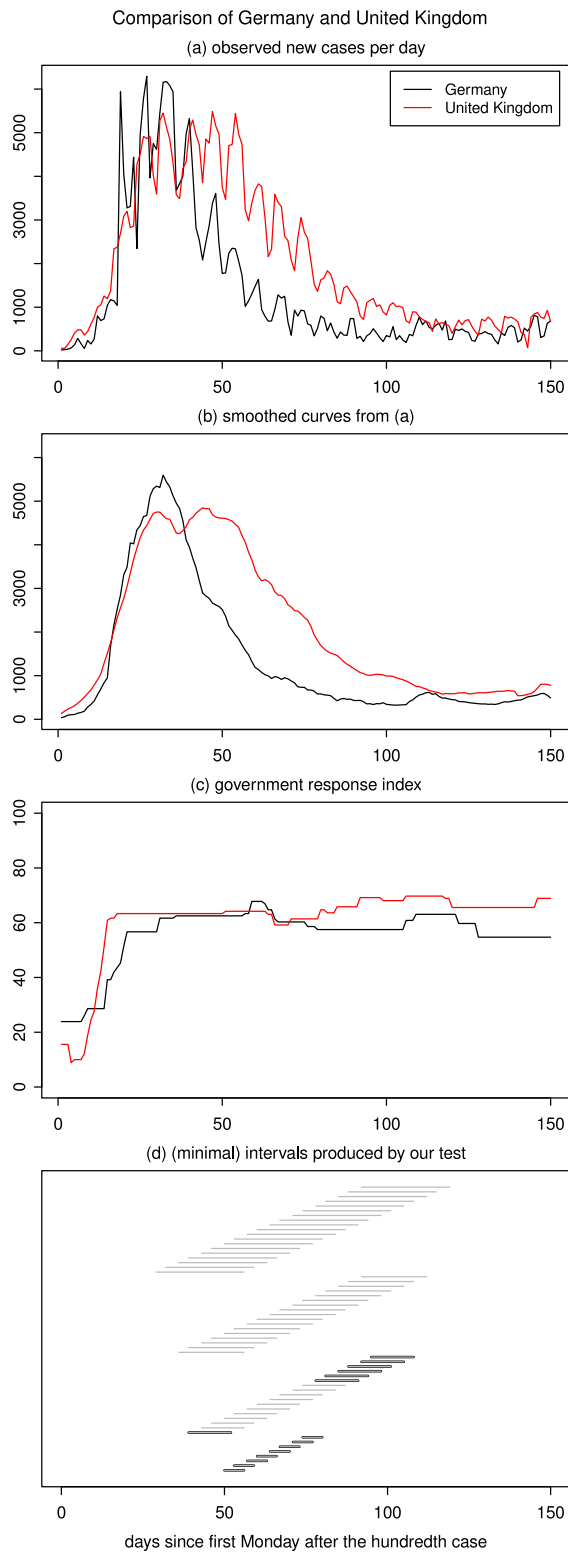


**Fig. 4.** Test results for the comparison of Germany and Spain. Note: Panel (a) shows the two observed time series, panel (b) smoothed versions of the time series, and panel (c) the corresponding Government Response Index (GRI). Panel (d) depicts the set of intervals  $\mathcal{F}_{\text{reject}}(i, j)$  in grey and the subset of minimal intervals  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j)$  with a black frame.



**Fig. 5.** Test results for the comparison of Germany and France. Note: Panel (a) shows the two observed time series, panel (b) smoothed versions of the time series, and panel (c) the corresponding Government Response Index (GRI). Panel (d) depicts the set of intervals  $\mathcal{F}_{\text{reject}}(i, j)$  in grey and the subset of minimal intervals  $\mathcal{F}_{\text{reject}}^{\text{min}}(i, j)$  with a black frame.





**Fig. 6.** Test results for the comparison of Germany and the UK. Note: Panel (a) shows the two observed time series, panel (b) smoothed versions of the time series, and panel (c) the corresponding Government Response Index (GRI). Panel (d) depicts the set of intervals  $\mathcal{F}_{\text{reject}}(i, j)$  in grey and the subset of minimal intervals  $\mathcal{F}_{\text{reject}}^{\min}(i, j)$  with a black frame.

days. After day 50, the trends become quite similar and start to decrease at approximately the same rate until around day 120. This may be due to the fact that Spain in general introduced more severe measures of lockdown than Germany (as can be seen upon inspecting the GRI in panel (c) of Fig. 4), which may have helped to battle the spread of infection. Furthermore, around days 110–120, the measures in Spain were less strict than in Germany, which could be a reason for the detected differences between the trends towards the end of the sample. However, a much more thorough analysis is of course needed to find out whether this is indeed the case or whether other factors were mainly responsible.

Turning to the comparison of Germany and Italy, we found that the German trend drops down significantly faster than the Italian one after approximately day 40. Interestingly, the GRI of Italy almost always lies above that of Germany. Hence, even though Italy has in general taken more severe and restrictive measures against the virus than Germany, it appears that the virus could be contained better in Germany (in the sense that the trend of daily new cases went down significantly faster in Germany than in Italy). This suggests that there are indeed important factors besides the level of government response to the pandemic which substantially influence the trend of new COVID-19 cases.

This brief discussion already indicates that it is extremely difficult to determine the exact causes of the differences in epidemic trends across countries. Since even similar countries such as those in our sample differ in a variety of aspects that are relevant for the spread of the virus, it is very challenging to pin down these causes. One issue that is often discussed in the context of cross-country comparisons are country-specific strategies to test for the coronavirus. The argument is that differences between epidemic trends may be spuriously produced by country-specific test procedures.

Even though we can of course not fully exclude this possibility, our test results are presumably not driven by different test regimes in the countries under consideration. To see this, we consider again the comparison of Germany and Italy: The test regimes in these two countries are arguably quite different. Germany is often cited as the country that employed early, widespread testing with more than 100 000 tests per week even in the beginning of the pandemic (Cohen and Kupferschmidt, 2020), while testing in Italy became widespread only in the later stages of the pandemic. Nevertheless, visual inspection of the raw and smoothed data in panels (a) and (b) of Fig. 3 suggest that the underlying time trends are very similar up to day 36. This is confirmed by our multiscale test which does not find any significant differences before that day. Hence, the different test regimes in Germany and Italy towards the beginning of the pandemic do not appear to have an overly strong effect and to produce spurious differences between the time trends. This suggests that the differences detected by our multiscale test indeed reflect differences in the way the virus spread in Germany and Italy rather than being mere artefacts of different test regimes.

### Acknowledgements

We thank the Editor, the Associate Editor and two referees for their constructive comments on an earlier version of the paper.

### Appendix

In what follows, we state and prove the main theoretical results on the multiscale test developed in Section 3. Throughout the Appendix, we let  $C$  be a generic positive constant that may take a different value on each occurrence. Unless stated differently,  $C$  depends neither on the time series length  $T$  nor on the dimension  $p$  of the test problem. We further use the symbols  $h_{\min} := \min_{1 \leq k \leq K} h_k$  and  $h_{\max} := \max_{1 \leq k \leq K} h_k$  to denote the smallest and largest interval length in the family  $\mathcal{F}$ , respectively.

The first result shows that the multiscale test asymptotically controls the FWER at level  $\alpha$ .

**Theorem A.1.** *Let (C1) and (C2) be satisfied. Moreover, assume that (i)  $h_{\max} = o(1/\{\log T\}^2)$ , (ii)  $h_{\min} \geq CT^{-b}$  for some  $b \in (0, 1)$ , and (iii)  $p = O(T^{(\theta/2)(1-b)-(1+\delta)})$  for some small  $\delta > 0$ . Then for any given  $\alpha \in (0, 1)$ ,*

$$\text{FWER}(\alpha) := \mathbb{P}\left(\exists(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\right) \leq \alpha + o(1),$$

where  $\mathcal{M}_0 \subseteq \mathcal{M}$  is the set of all  $(i, j, k) \in \mathcal{M}$  for which  $H_0^{(ijk)}$  holds true.

We briefly discuss the conditions (i)–(iii) on  $h_{\min}$ ,  $h_{\max}$  and  $p$ . Restriction (i) allows the maximal interval length  $h_{\max}$  to converge to zero very slowly, which means that  $h_{\max}$  can be picked very large in practice. According to restriction (ii), the minimal interval length  $h_{\min}$  can be chosen to go to zero as any polynomial  $T^{-b}$  with some  $b \in (0, 1)$ . Restriction (iii) allows the dimension  $p$  of the test problem to grow polynomially in  $T$ . Specifically,  $p$  may grow at most as the polynomial  $T^\gamma$  with  $\gamma = (\theta/2)(1-b)-(1+\delta)$ . As one can see, the exponent  $\gamma$  depends on the number of error moments  $\theta$  defined in (C2) and the parameter  $b$  that specifies the minimal interval length  $h_{\min}$ . In particular, for any given  $b \in (0, 1)$ , the exponent  $\gamma$  gets larger as  $\theta$  increases. Hence, the larger the number of error moments  $\theta$ , the faster  $p$  may grow in comparison to  $T$ . In the extreme case where all error moments exist, that is, where  $\theta$  can be made as large as desired,  $p$  may grow as any polynomial of  $T$ , no matter how we pick  $b \in (0, 1)$ . Thus, if the error terms have sufficiently many moments, the dimension  $p$  can be extremely large in comparison to  $T$  and the minimal interval length  $h_{\min}$  can be chosen very small.

The following corollary is an immediate consequence of Theorem A.1. It provides the theoretical justification needed to make simultaneous confidence statements of the form (3.9), (3.10) and (3.14).

**Corollary A.1.** Under the conditions of [Theorem A.1](#),

$$\mathbb{P}\left(\forall(i, j, k) \in \mathcal{R} : (i, j, k) \notin \mathcal{M}_0\right) \geq 1 - \alpha + o(1),$$

where  $\mathcal{R} = \{(i, j, k) \in \mathcal{M} \text{ with } |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\}$  is the set of triples  $(i, j, k)$  for which the test rejects the null  $H_0^{(ijk)}$ .

The next result specifies the power of the multiscale test against a certain class of local alternatives. To formulate it, we allow the functions  $\lambda_{i,T}$  and  $\lambda_{j,T}$  to depend on  $T$ , that is, we consider sequences of functions  $\{\lambda_{i,T}\}$  and  $\{\lambda_{j,T}\}$  rather than fixed functions  $\lambda_i$  and  $\lambda_j$ .

**Proposition A.1.** Let the conditions of [Theorem A.1](#) be satisfied and let  $\mathcal{M}_1$  be the set of triples  $(i, j, k) \in \mathcal{M}$  such that either

$$\lambda_{i,T}(w) - \lambda_{j,T}(w) \geq \kappa_T \sqrt{\log T / (Th_k)} \quad \text{for all } w \in \mathcal{I}_k \tag{A.1}$$

or

$$\lambda_{j,T}(w) - \lambda_{i,T}(w) \geq \kappa_T \sqrt{\log T / (Th_k)} \quad \text{for all } w \in \mathcal{I}_k, \tag{A.2}$$

where  $\{\kappa_T\}$  is any sequence of positive numbers for which  $\kappa_T / \ell_T \rightarrow \infty$  with  $\ell_T = \sqrt{\log T} \sqrt{\log \log T} / \log \log \log T$ . Then

$$\mathbb{P}\left(\forall(i, j, k) \in \mathcal{M}_1 : |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\right) = 1 - o(1)$$

for any given  $\alpha \in (0, 1)$ .

**Proof of Theorem A.1.** The proof proceeds in several steps.

*Step 1.* Let  $\hat{\Psi}_T = \max_{(i,j,k) \in \mathcal{M}} a_k(|\hat{\psi}_{ijk,T}^0| - b_k)$  with  $\hat{\psi}_{ijk,T}^0$  introduced in [\(3.4\)](#) and define  $\Psi_T = \max_{(i,j,k) \in \mathcal{M}} a_k(|\psi_{ijk,T}^0| - b_k)$  with

$$\psi_{ijk,T}^0 = \frac{1}{\sqrt{2Th_k}} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (\eta_{it} - \eta_{jt}).$$

To start with, we prove that

$$|\hat{\Psi}_T - \Psi_T| = o_p(r_T), \tag{A.3}$$

where  $\{r_T\}$  is any null sequence that converges more slowly to zero than  $\rho_T = \sqrt{\log T} \{\log p / \sqrt{Th_{\min}} + h_{\max} \sqrt{\log p}\}$ , that is,  $\rho_T / r_T \rightarrow 0$  as  $T \rightarrow \infty$ . Since the proof of [\(A.3\)](#) is rather technical and lengthy, the details are provided in the [Supplementary Material](#).

*Step 2.* We next prove that

$$\sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Psi_T \leq q) - \mathbb{P}(\Phi_T \leq q) \right| = o(1). \tag{A.4}$$

To do so, we rewrite the statistics  $\Psi_T$  and  $\Phi_T$  as follows: Define

$$V_t^{(ijk)} = V_{t,T}^{(ijk)} := \sqrt{\frac{T}{2Th_k}} \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (\eta_{it} - \eta_{jt})$$

for  $(i, j, k) \in \mathcal{M}$  and let  $\mathbf{V}_t = (V_t^{(ijk)} : (i, j, k) \in \mathcal{M})$  be the  $p$ -dimensional random vector with the entries  $V_t^{(ijk)}$ . With this notation, we get that  $\psi_{ijk,T}^0 = T^{-1/2} \sum_{t=1}^T V_t^{(ijk)}$  and thus

$$\begin{aligned} \Psi_T &= \max_{(i,j,k) \in \mathcal{M}} a_k(|\psi_{ijk,T}^0| - b_k) \\ &= \max_{(i,j,k) \in \mathcal{M}} a_k \left\{ \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T V_t^{(ijk)} \right| - b_k \right\}. \end{aligned}$$

Analogously, we define

$$W_t^{(ijk)} = W_{t,T}^{(ijk)} := \sqrt{\frac{T}{2Th_k}} \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (Z_{it} - Z_{jt})$$

with  $Z_{it}$  i.i.d. standard normal and let  $\mathbf{W}_t = (W_t^{(ijk)} : (i, j, k) \in \mathcal{M})$ . The vector  $\mathbf{W}_t$  is a Gaussian version of  $\mathbf{V}_t$  with the same mean and variance. In particular,  $\mathbb{E}[\mathbf{W}_t] = \mathbb{E}[\mathbf{V}_t] = 0$  and  $\mathbb{E}[\mathbf{W}_t \mathbf{W}_t^\top] = \mathbb{E}[\mathbf{V}_t \mathbf{V}_t^\top]$ . Similarly as before, we can write  $\phi_{ijk,T} = T^{-1/2} \sum_{t=1}^T W_t^{(ijk)}$  and

$$\begin{aligned} \Phi_T &= \max_{(i,j,k) \in \mathcal{M}} a_k (|\phi_{ijk,T}| - b_k) \\ &= \max_{(i,j,k) \in \mathcal{M}} a_k \left\{ \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T W_t^{(ijk)} \right| - b_k \right\}. \end{aligned}$$

For any  $q \in \mathbb{R}$ , it holds that

$$\begin{aligned} \mathbb{P}(\Psi_T \leq q) &= \mathbb{P}\left( \max_{(i,j,k) \in \mathcal{M}} a_k \left\{ \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T V_t^{(ijk)} \right| - b_k \right\} \leq q \right) \\ &= \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T V_t^{(ijk)} \right| \leq c_{ijk}(q) \text{ for all } (i, j, k) \in \mathcal{M} \right) \\ &= \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{V}_t \right| \leq \mathbf{c}(q) \right), \end{aligned}$$

where  $\mathbf{c}(q) = (c_{ijk}(q) : (i, j, k) \in \mathcal{M})$  is the  $\mathbb{R}^p$ -vector with the entries  $c_{ijk}(q) = q/a_k + b_k$ , we use the notation  $|v| = (|v_1|, \dots, |v_p|)^\top$  for vectors  $v \in \mathbb{R}^p$  and the inequality  $v \leq w$  is to be understood componentwise for  $v, w \in \mathbb{R}^p$ . Analogously, we have

$$\mathbb{P}(\Phi_T \leq q) = \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{W}_t \right| \leq \mathbf{c}(q) \right).$$

With this notation at hand, we can make use of Proposition 2.1 from Chernozhukov et al. (2017). In our context, this proposition can be stated as follows:

**Proposition A.2.** Assume that

- (a)  $T^{-1} \sum_{t=1}^T \mathbb{E}(V_t^{(ijk)})^2 \geq \delta > 0$  for all  $(i, j, k) \in \mathcal{M}$ .
- (b)  $T^{-1} \sum_{t=1}^T \mathbb{E}[|V_t^{(ijk)}|^{2+r}] \leq B_T^r$  for all  $(i, j, k) \in \mathcal{M}$  and  $r = 1, 2$ , where  $B_T \geq 1$  are constants that may tend to infinity as  $T \rightarrow \infty$ .
- (c)  $\mathbb{E}[\{\max_{(i,j,k) \in \mathcal{M}} |V_t^{(ijk)}|/B_T\}^\theta] \leq 2$  for all  $t$  and some  $\theta > 4$ .

Then

$$\sup_{\mathbf{c} \in \mathbb{R}^p} \left| \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{V}_t \right| \leq \mathbf{c} \right) - \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{W}_t \right| \leq \mathbf{c} \right) \right| \leq C \left\{ \left( \frac{B_T^2 \log^7(pT)}{T} \right)^{1/6} + \left( \frac{B_T^2 \log^3(pT)}{T^{1-2/\theta}} \right)^{1/3} \right\}, \tag{A.5}$$

where  $C$  depends only on  $\delta$  and  $\theta$ .

It is straightforward to verify that assumptions (a)–(c) are satisfied under the conditions of Theorem A.1 for sufficiently large  $T$ , where  $B_T$  can be chosen as  $B_T = Cp^{1/\theta} h_{\min}^{-1/2}$  with  $C$  sufficiently large. Moreover, it can be shown that the right-hand side of (A.5) is  $o(1)$  for this choice of  $B_T$ . Hence, Proposition A.2 yields that

$$\sup_{\mathbf{c} \in \mathbb{R}^p} \left| \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{V}_t \right| \leq \mathbf{c} \right) - \mathbb{P}\left( \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{W}_t \right| \leq \mathbf{c} \right) \right| = o(1),$$

which in turn implies (A.4).

Step 3. With the help of (A.3) and (A.4), we now show that

$$\sup_{q \in \mathbb{R}} \left| \mathbb{P}(\hat{\Psi}_T \leq q) - \mathbb{P}(\Phi_T \leq q) \right| = o(1). \tag{A.6}$$

To start with, the above supremum can be bounded by

$$\begin{aligned}
 & \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\hat{\Psi}_T \leq q) - \mathbb{P}(\Phi_T \leq q) \right| \\
 &= \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Psi_T \leq q + \{\Psi_T - \hat{\Psi}_T\}) - \mathbb{P}(\Phi_T \leq q) \right| \\
 &\leq \sup_{q \in \mathbb{R}} \max \left\{ \left| \mathbb{P}(\Psi_T \leq q + |\Psi_T - \hat{\Psi}_T|) - \mathbb{P}(\Phi_T \leq q) \right|, \right. \\
 &\quad \left. \left| \mathbb{P}(\Psi_T \leq q - |\Psi_T - \hat{\Psi}_T|) - \mathbb{P}(\Phi_T \leq q) \right| \right\} \\
 &\leq \sup_{q \in \mathbb{R}} \max \left\{ \left| \mathbb{P}(\Psi_T \leq q + r_T) - \mathbb{P}(\Phi_T \leq q) \right| + \mathbb{P}(|\Psi_T - \hat{\Psi}_T| > r_T), \right. \\
 &\quad \left. \left| \mathbb{P}(\Psi_T \leq q - r_T) - \mathbb{P}(\Phi_T \leq q) \right| + \mathbb{P}(|\Psi_T - \hat{\Psi}_T| > r_T) \right\} \\
 &\leq \max_{\ell=0,1} \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Psi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q) \right| + \mathbb{P}(|\Psi_T - \hat{\Psi}_T| > r_T) \\
 &= \max_{\ell=0,1} \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Psi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q) \right| + o(1), \tag{A.7}
 \end{aligned}$$

where the last line is by (A.3). Moreover, for  $\ell = 0, 1$ ,

$$\begin{aligned}
 & \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Psi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q) \right| \\
 &\leq \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Psi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q + (-1)^\ell r_T) \right| \\
 &\quad + \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Phi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q) \right| \\
 &= \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Phi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q) \right| + o(1), \tag{A.8}
 \end{aligned}$$

the last line following from (A.4). Finally, by Nazarov’s inequality (Nazarov, 2003, and Lemma A.1 in Chernozhukov et al., 2017), we have that for  $\ell = 0, 1$ ,

$$\begin{aligned}
 & \sup_{q \in \mathbb{R}} \left| \mathbb{P}(\Phi_T \leq q + (-1)^\ell r_T) - \mathbb{P}(\Phi_T \leq q) \right| \\
 &= \sup_{q \in \mathbb{R}} \left| \mathbb{P}\left(\left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{w}_t \right| \leq \mathbf{c}(q + (-1)^\ell r_T)\right) - \mathbb{P}\left(\left| \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{w}_t \right| \leq \mathbf{c}(q)\right) \right| \\
 &\leq C \frac{r_T \sqrt{\log(2p)}}{\min_{1 \leq k \leq K} a_k} \leq Cr_T \sqrt{\log \log T} \sqrt{\log(2p)}, \tag{A.9}
 \end{aligned}$$

where  $C$  is a constant that depends only on the parameter  $\delta$  defined in condition (a) of Proposition A.2 and we have used the fact that  $\min_k a_k \geq c/\sqrt{\log \log T}$  for some  $c > 0$ . Inserting (A.8) and (A.9) into Eq. (A.7) completes the proof of (A.6).

*Step 4.* By definition of the quantile  $q_{T,\text{Gauss}}(\alpha)$ , it holds that  $\mathbb{P}(\Phi_T \leq q_{T,\text{Gauss}}(\alpha)) \geq 1 - \alpha$ . As shown in the Supplementary Material, we even have that

$$\mathbb{P}(\Phi_T \leq q_{T,\text{Gauss}}(\alpha)) = 1 - \alpha \tag{A.10}$$

for any  $\alpha \in (0, 1)$ . From this and (A.6), it immediately follows that

$$\mathbb{P}(\hat{\Psi}_T \leq q_{T,\text{Gauss}}(\alpha)) = 1 - \alpha + o(1), \tag{A.11}$$

which in turn implies that

$$\begin{aligned} \text{FWER}(\alpha) &= \mathbb{P}\left(\exists(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\right) \\ &= \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}_0} a_k(|\hat{\psi}_{ijk,T}| - b_k) > q_{T,\text{Gauss}}(\alpha)\right) \\ &= \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}_0} a_k(|\hat{\psi}_{ijk,T}^0| - b_k) > q_{T,\text{Gauss}}(\alpha)\right) \\ &\leq \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}} a_k(|\hat{\psi}_{ijk,T}^0| - b_k) > q_{T,\text{Gauss}}(\alpha)\right) \\ &= \mathbb{P}(\hat{\Psi}_T > q_{T,\text{Gauss}}(\alpha)) = \alpha + o(1). \end{aligned}$$

This completes the proof of [Theorem A.1](#).  $\square$

**Proof of Corollary A.1.** With the help of [Theorem A.1](#), we obtain that

$$\begin{aligned} 1 - \alpha + o(1) &\leq 1 - \text{FWER}(\alpha) \\ &= \mathbb{P}\left(\forall(i, j, k) \in \mathcal{M}_0 : |\hat{\psi}_{ijk,T}| \leq c_{T,\text{Gauss}}(\alpha, h_k)\right) \\ &\leq \mathbb{P}\left(\forall(i, j, k) \in \mathcal{R} : (i, j, k) \notin \mathcal{M}_0\right), \end{aligned}$$

which is the statement of [Corollary A.1](#).  $\square$

**Proof of Proposition A.1.** To start with, note that

$$c \frac{1}{\sqrt{\log \log T}} \leq a_k \leq C \frac{\sqrt{\log T}}{\log \log T} \quad \text{and} \quad b_k \leq C\sqrt{\log T} \tag{A.12}$$

with appropriately chosen constants  $c$  and  $C$ . We decompose the statistics  $\hat{\psi}_{ijk,T}$  into two parts. In particular, we write  $\hat{\psi}_{ijk,T} = \hat{\psi}_{ijk,T}^A + \hat{\psi}_{ijk,T}^B$  with

$$\begin{aligned} \hat{\psi}_{ijk,T}^A &= \frac{\sigma \sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k) \left( \sqrt{\lambda_i(\frac{t}{T})} \eta_{it} - \sqrt{\lambda_j(\frac{t}{T})} \eta_{jt} \right)}{\hat{\sigma} \left\{ \sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k) (X_{it} + X_{jt}) \right\}^{1/2}} \\ \hat{\psi}_{ijk,T}^B &= \frac{\sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k) \left( \lambda_i(\frac{t}{T}) - \lambda_j(\frac{t}{T}) \right)}{\hat{\sigma} \left\{ \sum_{t=1}^T \mathbf{1}(\frac{t}{T} \in \mathcal{I}_k) (X_{it} + X_{jt}) \right\}^{1/2}}. \end{aligned}$$

As we will prove below, it holds that

$$\min_{(i,j,k) \in \mathcal{M}_1} |\hat{\psi}_{ijk,T}^B| \geq C\kappa_T \sqrt{\log T} \quad \text{with prob. approaching } 1 \tag{A.13}$$

for some sufficiently small constant  $C > 0$  and

$$\max_{(i,j,k) \in \mathcal{M}} |\hat{\psi}_{ijk,T}^A| = O_p(\sqrt{\log T}). \tag{A.14}$$

From [\(A.13\)](#) and [\(A.14\)](#), it follows that

$$\min_{(i,j,k) \in \mathcal{M}_1} a_k(|\hat{\psi}_{ijk,T}| - b_k) \geq \min_{(i,j,k) \in \mathcal{M}_1} a_k |\hat{\psi}_{ijk,T}^B| - \max_{(i,j,k) \in \mathcal{M}} a_k (|\hat{\psi}_{ijk,T}^A| + b_k) \geq C \frac{\kappa_T \sqrt{\log T}}{\sqrt{\log \log T}} \tag{A.15}$$

with probability tending to 1, where we have used the bounds on  $a_k$  and  $b_k$  from [\(A.12\)](#) and the assumption that  $\kappa_T/\ell_T \rightarrow \infty$  with  $\ell_T$  defined in [Proposition A.1](#). It further holds that

$$q_{T,\text{Gauss}}(\alpha) \leq C \frac{\log T}{\log \log \log T} \tag{A.16}$$



with a sufficiently large constant  $C$ , since

$$\begin{aligned} & \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}} a_k(|\phi_{ijk,T}| - b_k) \leq C \frac{\log T}{\log \log \log T}\right) \\ & \geq \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}} a_k \max_{(i,j,k) \in \mathcal{M}} (|\phi_{ijk,T}| - b_k) \leq C \frac{\log T}{\log \log \log T}\right) \\ & \geq \mathbb{P}\left(\frac{\sqrt{\log T}}{\log \log \log T} \max_{(i,j,k) \in \mathcal{M}} |\phi_{ijk,T}| \leq C \frac{\log T}{\log \log \log T}\right) \\ & = \mathbb{P}\left(\max_{(i,j,k) \in \mathcal{M}} |\phi_{ijk,T}| \leq C\sqrt{\log T}\right) \geq 1 - \alpha, \end{aligned}$$

where the last inequality is a consequence of the fact that the terms  $\phi_{ijk,T}$  are normally distributed random variables and  $|\mathcal{M}| = p \leq CT^\gamma$ . From (A.15), (A.16) and the assumption that  $\kappa_T/\ell_T \rightarrow \infty$ , we can finally conclude that

$$\begin{aligned} & \mathbb{P}\left(\forall (i, j, k) \in \mathcal{M}_1 : |\hat{\psi}_{ijk,T}| > c_{T,\text{Gauss}}(\alpha, h_k)\right) \\ & = \mathbb{P}\left(\min_{(i,j,k) \in \mathcal{M}_1} a_k(|\hat{\psi}_{ijk,T}| - b_k) > q_{T,\text{Gauss}}(\alpha)\right) = 1 - o(1), \end{aligned}$$

which is the statement of Proposition A.1.

It remains to prove (A.13) and (A.14). From (S.10) in the Supplementary Material together with some straightforward arguments, it follows that for any fixed  $\delta > 0$ ,

$$\min_{(i,j,k) \in \mathcal{M}} \left\{ \frac{1}{Th_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (X_{it} + X_{jt}) \right\} \geq (2 - \delta)\lambda_{\min} \tag{A.17}$$

$$\max_{(i,j,k) \in \mathcal{M}} \left\{ \frac{1}{Th_k} \sum_{t=1}^T \mathbf{1}\left(\frac{t}{T} \in \mathcal{I}_k\right) (X_{it} + X_{jt}) \right\} \leq (2 + \delta)\lambda_{\max} \tag{A.18}$$

with probability tending to 1. Since  $\hat{\sigma}^2 = \sigma^2 + O_p(\sqrt{\log p/T})$  by Lemma S.1, it further holds that with probability tending to 1,

$$(1 - \delta)\sigma \leq \hat{\sigma} \leq (1 + \delta)\sigma \tag{A.19}$$

for any fixed  $\delta > 0$ . Taking into account that for any  $(i, j, k) \in \mathcal{M}_1$ , either  $\lambda_{i,T}(w) - \lambda_{j,T}(w) \geq \kappa_T\sqrt{\log T/(Th_k)}$  or  $\lambda_{j,T}(w) - \lambda_{i,T}(w) \geq \kappa_T\sqrt{\log T/(Th_k)}$  for all  $w \in \mathcal{I}_k$ , we can use (A.18) and (A.19) to obtain that

$$\min_{(i,j,k) \in \mathcal{M}_1} |\hat{\psi}_{ijk,T}^B| \geq \frac{\kappa_T\sqrt{\log T}}{(1 + \delta)\sigma\sqrt{(2 + \delta)\lambda_{\max}}} = C\kappa_T\sqrt{\log T}$$

with probability tending to 1. Moreover, with the help of (A.17), (A.19) and analogous arguments as for the proof of (S.6) in the Supplementary Material, we can show that

$$\max_{(i,j,k) \in \mathcal{M}} |\hat{\psi}_{ijk,T}^A| = O_p(\sqrt{\log p}) = O_p(\sqrt{\log T}),$$

where the last equation is due to the fact that  $p = O(T^\gamma)$  for a fixed  $\gamma > 0$ .  $\square$

### Appendix B. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jeconom.2021.04.010>.

### References

Bakirov, N.K., Rizzo, M.L., Székely, G.J., 2006. A multivariate nonparametric test of independence. *J. Multivariate Anal.* 97 (8), 1742–1756.  
 Blum, J.R., Kiefer, J., Rosenblatt, M., 1961. Distribution free tests of independence based on the sample distribution function. *Ann. Math. Stat.* 32, 485–498.  
 Chen, L., Wu, W.B., 2019. Testing for trends in high-dimensional time series. *J. Amer. Statist. Assoc.* 114, 869–881.  
 Chernozhukov, V., Chetverikov, D., Kato, K., 2017. Central limit theorems and bootstrap in high dimensions. *Ann. Probab.* 45, 2309–2352.  
 Cohen, J., Kupferschmidt, K., 2020. Countries test tactics in ‘war’ against COVID-19. *Science* 367 (6484), 1287–1288.  
 Cox, D.R., 1983. Some remarks on overdispersion. *Biometrika* 70, 269–274.  
 De Salazar, P.M., Niehus, R., Taylor, A., Buckee, C., Lipsitch, M., 2020. Using predicted imports of 2019-nCoV cases to determine locations that may not be identifying all imported cases. medRxiv.  
 Degras, D., Xu, Z., Zhang, T., Wu, W.B., 2012. Testing for parallelism among trends in multiple time series. *IEEE Trans. Signal Process.* 60, 1087–1097.  
 Delgado, M.A., 1993. Testing the equality of nonparametric regression curves. *Statist. Probab. Lett.* 17, 199–204.  
 Dümbgen, L., Spokoiny, V.G., 2001. Multiscale testing of qualitative hypotheses. *Ann. Statist.* 29, 124–152.  
 Dümbgen, L., Walther, G., 2008. Multiscale inference about a density. *Ann. Statist.* 36, 1758–1785.

- Dunker, F., Eckle, K., Proksch, K., Schmidt-Hieber, J., 2019. Tests for qualitative features in the random coefficients model. *Electron. J. Stat.* 13, 2257–2306.
- Eckle, K., Bissantz, N., Dette, H., 2017. Multiscale inference for multivariate deconvolution. *Electron. J. Stat.* 11, 4179–4219.
- Efron, B., 1986. Double exponential families and their use in generalized linear regression. *J. Amer. Statist. Assoc.* 81, 709–721.
- Fryzlewicz, P., Sapatinas, T., Subba Rao, S., 2006. A Haar-Fisz technique for locally stationary volatility estimation. *Biometrika* 93, 687–704.
- Fryzlewicz, P., Sapatinas, T., Subba Rao, S., 2008. Normalized least-squares estimation in time-varying ARCH models. *Ann. Statist.* 36, 742–786.
- Hafner, C.M., Linton, O., 2010. Efficient estimation of a multivariate multiplicative volatility model. *J. Econometrics* 159, 55–73.
- Hale, T., Petherick, A., Phillips, T., Webster, S., 2020a. Variation in government responses to COVID-19. Blavatnik school of government working paper 31.
- Hale, T., Webster, S., Petherick, A., Phillips, T., Kira, B., 2020b. Oxford COVID-19 government response tracker. Blavatnik school of government. <http://www.bsg.ox.ac.uk/covidtracker>.
- Hall, P., Hart, J.D., 1990. Bootstrap test for difference between means in nonparametric regression. *J. Amer. Statist. Assoc.* 85, 1039–1049.
- Härdle, W., Marron, J.S., 1990. Semiparametric comparison of regression curves. *Ann. Statist.* 18, 63–89.
- Hidalgo, J., Lee, J., 2014. A CUSUM test for common trends in large heterogeneous panels. In: *Essays in Honor of Peter C. B. Phillips*. Emerald Group Publishing Limited, pp. 303–345.
- King, E.C., Hart, J.D., Wehrly, T.E., 1991. Testing the equality of regression curves using linear smoothers. *Statist. Probab. Lett.* 12, 239–247.
- Kulasekera, K.B., 1995. Comparison of regression curves using quasi-residuals. *J. Amer. Statist. Assoc.* 90, 1085–1093.
- Lavergne, P., 2001. An equality test across nonparametric regressions. *J. Econometrics* 103, 307–344.
- McCullagh, P., Nelder, J., 1989. *Generalized Linear Models*. Chapman and Hall.
- Mikosch, T., Stărică, C., 2000. Is it really long memory we see in financial returns?. In: Embrechts, P. (Ed.), *Extremes and Integrated Risk Management*. pp. 149–168.
- Mikosch, T., Stărică, C., 2004. Non-stationarities in financial time series, the long-range dependence, and IGARCH effects. *Rev. Econ. Stat.* 86, 378–390.
- Munk, A., Dette, H., 1998. Nonparametric comparison of several regression functions: exact and asymptotic theory. *Ann. Statist.* 26, 2339–2368.
- Nazarov, F., 2003. On the maximal perimeter of a convex set in  $\mathbb{R}^n$  with respect to a Gaussian measure. In: *Geometric Aspects of Functional Analysis*. In: *Lecture Notes in Mathematics*, vol. 1807, Springer, pp. 169–187.
- Neumeyer, N., Dette, H., 2003. Nonparametric comparison of regression curves: an empirical process approach. *Ann. Statist.* 31, 880–920.
- Pardo-Fernández, J.C., Van Keilegom, I., González-Manteiga, W., 2007. Testing for the equality of  $k$  regression curves. *Statist. Sinica* 17, 1115–1137.
- Park, C., Vaughan, A., Hannig, J., Kang, K.-H., 2009. SiZer analysis for the comparison of time series. *J. Statist. Plann. Inference* 139, 3974–3988.
- Pellis, L., Scarabel, F., Stage, H.B., Overton, C.E., Chappell, L.H., Lythgoe, K.A., Fearon, E., Bennett, E., Curran-Sebastian, J., Das, R., et al., 2020. Challenges in control of Covid-19: short doubling time and long delay to effect of interventions. *arXiv*.
- Robinson, P.M., 1989. Nonparametric estimation of time-varying parameters. In: Hackl, P. (Ed.), *Statistical Analysis and Forecasting of Economic Structural Change*. Springer, pp. 253–264.
- Rohde, A., 2008. Adaptive goodness-of-fit tests based on signed ranks. *Ann. Statist.* 36, 1346–1374.
- Rufibach, K., Walther, G., 2010. The block criterion for multiscale inference about a density, with applications to other multiscale problems. *J. Comput. Graph. Statist.* 19, 175–190.
- Schmidt-Hieber, J., Munk, A., Dümbgen, L., 2013. Multiscale methods for shape constraints in deconvolution: confidence statements for qualitative features. *Ann. Statist.* 41, 1299–1328.
- Sinha, B.K., Wieand, H.S., 1977. Multivariate nonparametric tests for independence. *J. Multivariate Anal.* 7, 572–583.
- Tobías, A., Valls, J., Satorra, P., Tebé, C., 2020. COVID19-tracker: a shiny app to produce comprehensive data visualization for SARS-CoV-2 epidemic in Spain. *medRxiv*.
- Young, S.G., Bowman, A.W., 1995. Nonparametric analysis of covariance. *Biometrics* 51, 920–931.
- Zhang, Y., Su, L., Phillips, P.C.B., 2012. Testing for common trends in semi-parametric panel data models with fixed effects. *Econom. J.* 15, 56–100.