

# CoViD-19: an automatic, semiparametric estimation method for the population infected in Italy

Livio Fenga

ISTAT, Rome, Italy

## ABSTRACT

To date, official data on the number of people infected with the SARS-CoV-2—responsible for the Covid-19—have been released by the Italian Government just on the basis of a non-representative sample of population which tested positive for the swab. However a reliable estimation of the number of infected, including asymptomatic people, turns out to be crucial in the preparation of operational schemes and to estimate the future number of people, who will require, to different extents, medical attentions. In order to overcome the current data shortcoming, this article proposes a bootstrap-driven, estimation procedure for the number of people infected with the SARS-CoV-2. This method is designed to be robust, automatic and suitable to generate estimations at regional level. Obtained results show that, while official data at March the 12th report 12.839 cases in Italy, people infected with the SARS-CoV-2 could be as high as 105.789.

**Subjects** Epidemiology, Global Health, Health Policy, Infectious Diseases, Statistics

**Keywords** Autoregressive metric, Covid-19, Maximum entropy bootstrap, Model uncertainty, Number of Italian people infected

## INTRODUCTION

Covid-19 epidemic has severely hit Italy, and its spread throughout Europe is expected soon. In such a scenario, the availability of reliable information related to its spread plays a significant role in many regards. In fact, many targeted measures, such as the coordination among emergency services or the implementation of operative actions (extensive or local lock-downs or even curfew) can only be efficiently taken when reliable estimates of the epidemic spread are available at the population level.

At the moment, official data on the infection in Italy are based on non-random, non-representative samples of the population: people are tested for Covid-19 on the condition that some symptoms related to the virus are present. These data can ensure a proper estimation of the number of both deaths and hospitalizations due to the virus and are crucial for the optimization of the available resources. Nonetheless, from a statistical point of view, the number of people tested positive for Covid-19 represents a simple count which is not suitable to provide a reliable assessment of the “true”, unknown, number of infected people (thereafter “positive cases” ). In addition to the strong bias components induced by this testing strategy, there is at least another major obstacle to the construction of a valid estimator: the small sample size available. These issues are

Submitted 16 April 2020  
Accepted 2 January 2021  
Published 4 March 2021

Corresponding author  
Livio Fenga, fenga@istat.it

Academic editor  
Sharif Aly

Additional Information and  
Declarations can be found on  
page 15

DOI 10.7717/peerj.10819

© Copyright  
2021 Fenga

Distributed under  
Creative Commons CC-BY 4.0

**OPEN ACCESS**

considered in the available literature: *Feinstein & Esdaile (1987)* point out how the statistical information in many cases can contain gross violations of epidemiological principles as well as of scientific standards for credible evidence. On the other hand, a substantial corpus of theory and methods are available to epidemiologists and/or the statisticians working on the field of epidemiology—see, for example, *Kahn, Kahn & Sempos (1989)* and, more recently, *Clayton & Hills (2013)* and *Lawson (2013)*. Therefore, a “reasonable” trade-off between goodness of the outcomes of a statistical analysis and the available data, in some cases, is the best we can hope for. In the case of the present article, the shortness of the time series of interest is simply something that, at an early stage of an epidemic, cannot be avoided. It is well known that the shortness of the time series of interest might lead to a strong bias in the asymptotic results and therefore to the construction of biased confidence intervals. However, the results obtained in this article can be considered reliable as the approach used has been specifically designed to mitigate these negative effects. To confirm that, the estimates provided by this method have been proved to be in line with those published by official entities and have been reported on a number of nationally distributed daily newspaper published in Italy.

Based on the number of the deaths and of the observed positive cases and improving on an estimation equation proposed by *Pueyo (2020)*, this article aims at estimating the “true” number of people infected by the Covid-19 in each of the 20 Italian regions. Presently, to the best of the author’s knowledge, Pueyo’s equation does not appear in the literature nevertheless its validity in the present context will be discussed later in “Data and Contagion Indicator”. In more details, the presented procedure is designed to reduce the impact of the biasing components on the parameter estimations, by employing a resampling scheme, called Maximum Entropy Bootstrap (MEBOOT) proposed by *Vinod & López-de Lacalle (2009)*. This bootstrap method is particularly suitable in this context: as it will be outlined in the sequel it is designed to work with a broad class of time series (including non stationary ones) and—by virtue of its inherent simplicity—is able to generate bona fide replications in the case of short time series. In fact, unlike other schemes, long time series are not required. For example, in the case of the sieve bootstrap method *Andrés, Pena & Romo (2002)*, a lengthy series is needed in order to estimate an high order autoregressive model from which the bootstrap replications are generated. In conjunction with MEBOOT, a distance measure—based on the theory of stochastic processes and proposed by *Piccolo (1990)*—has been used to find pairs of similar regions. As it will be explained later, this has been done to maintain the same methodology in those cases where one of the variable employed in the model—that is, the number of deaths—was missing.

## AN OVERVIEW OF THE PROPOSED METHOD

In small data sets it is essential to save degrees of freedom (DOF) which are inevitably lost in an amount correlated with the complexity of the statistical model entertained (see, for example, *Faes et al. (2009)* and *Barnard & Rubin (1999)*). With this in mind, the proposed method is of the type semiparametric and consists of two parts: a purely non-parametric and a parametric one. The non-parametric part refers to the maximum

entropy resampling method, which will be used to generate more robust estimations. On the other hand, a parametric approach has been chosen to select certain regions on the basis of a similarity function, as it will be explained at the end of the following “Data and Contagion Indicator”. While the former does not pose problems in terms of DOF, the latter clearly does. However, the sacrifice in terms of DOF is very limited as an autoregressive model of order 1 (employed in a suitable distance function, as below illustrated) has proved sufficient for the purpose. DOF—saving strategy is also the driving force behind the choice to not consider exogenous variables such as the regions geolocation or their population—for example, in a regression-like scheme—but to implicitly assume these (and other) variables are embedded in the dynamic of the time series considered.

## DATA AND CONTAGION INDICATOR

The paper makes use of official data, published by the Italian Authorities, related to the following two variables employed in the proposed method, that is, the number of

1. deaths from Covid-19 (denoted by the symbol  $M_t$ )
2. currently positive cases which have been recorded as a result of the administration of the test (denoted by the symbol  $C_t$ ).

The data set includes 18 daily data points collected at the regional level during the period of February 24th to March 12th 2020. The total number of Italian regions considered is 20. However, one special administrative area (Trentino Alto Adige) is divided in two subregions, that is, Trento and Bolzano. Therefore, the set containing all the Italian regions—called  $\Omega$ —has cardinality  $|\Omega| = 21$  (the cardinality function is denoted by the symbol  $|\cdot|$ ). Two different subsets are built from  $\Omega$  that is,  $\Omega^\bullet$ —containing the regions for which at least one death, out of the group of tested people, has been recorded and  $\Omega^\circ$  (no recorded deaths). Those two sets are now specified:

1.  $\{\Omega^\bullet\} \equiv \text{Piemonte, Lombardia, Veneto, Friuli, Liguria, Emilia, Toscana, Marche, Lazio, Abruzzo, ValleAosta, Bolzano, Campania, Puglia, Sicilia}$
2.  $\{\Omega^\circ\} \equiv \text{Trento, Umbria, Molise, Basilicata, Calabria, Sardegna,}$

where  $\Omega \equiv \Omega^\bullet \cup \Omega^\circ$ . In what follows, the two superscripts  $\bullet$  and  $\circ$  will be always used respectively with reference to the regions  $\{r_1, r_2, \dots, r_{15}\} \in \Omega^\bullet$  and  $\{s_1, s_2, \dots, s_6\} \in \Omega^\circ$ . The time span is denoted as  $\{1, 2, \dots, T\}$ . In the case of the regions included in the set  $\Omega^\bullet$ , following [Pueyo \(2020\)](#), the total number of positive is estimated as follows:

$$y_{j,T}^\bullet = w_T * 2^{\frac{\tau}{\delta}} \quad (1)$$

$$w_T = \frac{C_T}{M_T} \quad (2)$$

Here,  $w_T$  (Eq. 2) is the ratio between the current positive cases ( $C$ ) and the number of deaths ( $M$ ) whereas, in Eq. (1),  $\tau$  is the average doubling time for the Covid-19 (i.e., the

average span of time needed for the virus to double the cases) and  $\delta$  the average time needed for an infected person to die. These two constant terms have been kept fixed as estimated according to the data so far available as reported and justified in the above mentioned Puejo's web document. They are as follows:  $\tau = 17.3$  and  $\delta = 6.2$ .

By construction, Eqs. (1) and (2) are able to properly describe the spread of the virus at the population level, as they are based on the key parameters average doubling  $\tau$  and time to death ( $\delta$ ). To make this clear, suppose a situation where  $\tau = \delta$  (i.e., all the subjects, in average, die the following day after the disease has been contracted). In this case, Eq. (1) reduces to  $y_{j,T}^{\bullet} = 2 * w_T$ , that is we will have the total positives equal to twice the mortality rate. As for the constants chosen, they appear to be in line with the data released by the Italian public authority.

The case of the regions belonging to  $\Omega^{\circ}$  is more complicated. The related estimation procedure has been carried out as detailed below (the subscript  $t$  will be omitted for the sake of simplicity):

1. given the series  $s_j \in \Omega^{\circ}$ , a series  $c^{\pi} \in \Omega^{\bullet}$  minimizer of a suitable distance function—denoted by the Greek letter  $\pi(\cdot)$ —is found. In symbols:

$$c^{\pi} = \underset{(c \in \Omega^{\bullet})}{\operatorname{argmin}} \pi(s, c) \quad (3)$$

2. the estimated number of positives at the population level—already found for  $c^{\pi}$ , say  $I_{c^{\pi}}$ —becomes the weight for which the total cases recorded for  $s_j$  are multiplied. Therefore, the estimate of the variable of interest for this case becomes

$$y_{j,T}^{\circ} = \frac{I_{c^{\pi}} * C_{s_j}}{C_{r_j}} \quad (4)$$

The distance function adopted  $\pi(\cdot)$  (Eq. 3), called AR-distance, has been introduced by [Piccolo \(2007\)](#). Briefly, this metric can be applied if and only if the pair of series of interest are assumed to be realizations of two (possibly of different orders) Autoregressive Moving Average (ARMA) models (see, e.g., [Makridakis & Hibon \(1997\)](#)). Under this condition, each series can be expressed as an autoregressive model of infinite order, that is,  $AR(\infty)$ , whose (infinite) sequence of AR parameters is denoted by  $\{\alpha_j\}_j^{\infty} \equiv \alpha_1, \alpha_2, \dots$

Without loss of generality, the distance between the series  $s$  and  $c$ , that is,  $\pi(s, c)$  (Eq. 4), under  $(s, c) \sim ARMA(\alpha, \beta)$ , being  $\alpha$  and  $\beta$  respectively the autoregressive and moving average parameters, is expressed as

$$\pi(s, c) = \left( \sum_{j=1}^{\infty} \alpha_j (s) - \alpha_j (c) \right)^{1/2} \quad (5)$$

Equation 5 asymptotically converges under stationary condition of the autoregressive parameters, as proved in [Piccolo \(2010\)](#). In other words, considering for brevity only the autoregression in  $\alpha_j$ , the roots of the polynomial  $\Phi(z) := 1 - \sum_{j=1}^S \alpha_j z^{S-j}$  must lie

outside the unit circle, that is, each root  $z_j$  must satisfy  $|z_1| > 1$ . For other asymptotic properties the reader is referred to [Corduas & Piccolo \(2008\)](#). It is well known that, with small sample sizes, the asymptotic properties of the ARMA parameters tend to deteriorate and therefore the statistical model might not perform optimally. However, in the present context their use is justified at least for two reasons: firstly the ARMA models have been here employed only for the construction of a simple distance measure used to build a similarity ranking of the Italian regions. As a simple way to pick a suitable “donor” (see the explanation below), that ARMA models tend to not perform optimally in such conditions can be considered a crucial issues. The second reason refers to the fact that, epidemics are an emergency situations and the the typical case where only a few (all the more so likely to be noisy) data points are available. Finally, in order to reach stationarity and thus correctly assess the distance functions, all the models have been estimated on properly differentiated time series.

## THE RESAMPLING METHOD

The bootstrap scheme adopted proved to be adequate for the problem at hand. Given the pivotal role played in the proposed method, it will be briefly presented. In essence, the choice of the most appropriate resampling method is far from being an easy task, especially when the identical and independent distribution (*iid*) assumption (used in Efron’s initial bootstrap method) is violated. Under dependance structures embedded in the data, simple sampling with replacement has been proved—see, for example [Carlstein \(1986\)](#)—to yield suboptimal results. As a matter of fact, *iid*—based bootstrap schemes are not designed to capture, and therefore replicate, dependance structures. This is especially true under the actual conditions (small sample sizes) where the selection of the “right” resampling scheme becomes a particularly challenging task. Several ad hoc methods have been therefore proposed, many of which now freely and publicly available in the form of powerful routines working under software package such as Python<sup>®</sup> or R<sup>®</sup>. In more details, while in the classic bootstrap an ensemble  $\Gamma$  represents the population of reference the observed time series is drawn from, in *MEB* a large number of ensembles (subsets), say  $\{\gamma_1, \dots, \gamma_N\}$  becomes the elements belonging to  $\Gamma$ , each of them containing a large number of replicates  $\{x_1, \dots, x_J\}$ . Perhaps, the most important characteristic of the *MEB* algorithm is that its design guarantees the inference process to satisfy the ergodic theorem. Formally, recalling the symbol  $|\cdot|$  to denote the cardinality function (counting function) of a given ensemble of time series  $\{x_t \in \gamma_i; i = 1, \dots, N\}$ , the *MEB* procedure generates a set of disjoint subsets  $\Gamma_N \equiv \gamma_1 \cap \gamma_1 \dots \cap \gamma_N$  s.t.  $\mathbb{E}\Gamma_N \approx \mu(x_t)$ , being  $\mu(\cdot)$  the sample mean. Furthermore, basic shape and probabilistic structure (dependency) is guaranteed to be retained  $\forall x_{t,j}^* \in \gamma_i \subset \Gamma$ .

*MEB* resampling scheme has significant advantages over many of the available bootstrap methods: it does not require complicated tune up procedures (unavoidable, for example, in the case of resampling methods of the type Block Bootstrap) and it is effective under non-stationarity. *MEB* method relies on the entropy theory and the related

concept of (un)informativeness of a system. In particular, the Maximum Entropy of a given density  $\rho(x)$ , is chosen so that the expectation of the Shannon Information  $\mathcal{H} = \mathbb{E}(-\log \rho(x))$ , is maximized, that is,

$$\max_{(\rho)} \mathcal{H} = \mathbb{E}(-\log \rho(x))$$

Under mass and mean preserving constraints, this resampling scheme generates an ensemble of time series from a density function satisfying (4). Technically, MEB algorithm can be broken down, following *Koutris, Heracleous & Spanos (2008)*, in 8 steps. They are:

1. a sorting matrix of dimension  $T \times 2$ , say  $\mathcal{S}_1$ , accommodates in its first column the time series of interest  $x_t$  and an Index Set—that is,  $I_{ind} = \{2, 3, \dots, T\}$ —in the other one;
2.  $\mathcal{S}_1$  is sorted according to the numbers placed in the first column. As a result, the order statistics  $\mathbf{x}_{(t)}$  and the vector  $I_{ord}$  of sorted  $I_{ind}$  are generated and respectively placed in the first and second column;
3. compute “intermediate points”, averaging over successive order statistics, that is,  $c_t = \frac{x_{(t)} + x_{(t+1)}}{2}$ ,  $t = 1, \dots, T - 1$  and define intervals  $I_t$  constructed on  $c_t$  and  $r_t$ , using ad hoc weights obtained by solving the following set of equations:

$$(i) \ g(x) = \frac{1}{r_1} \exp\left(\frac{[x - c_1]}{r_1}\right); \quad x \in I_1; \ r_1 = \frac{3x_{(1)}}{4} + \frac{x_{(2)}}{4}$$

$$(ii) \ g(x) = \frac{1}{c_k - c_{k-1}}; \quad x \in (c_k; c_{k+1}),$$

$$r_k = \frac{x_{(k-1)}}{4} + \frac{x_{(k)}}{2} + \frac{x_{(k+1)}}{4}; \quad k = 1, \dots, T - 1;$$

$$(iii) \ g(x) = \frac{1}{r_T} \exp\left(\frac{[c_{T-1} - x]}{r_T}\right); \quad x \in I_T; \quad r_T = \frac{x_{T-1}}{4} + \frac{3x_T}{4}$$

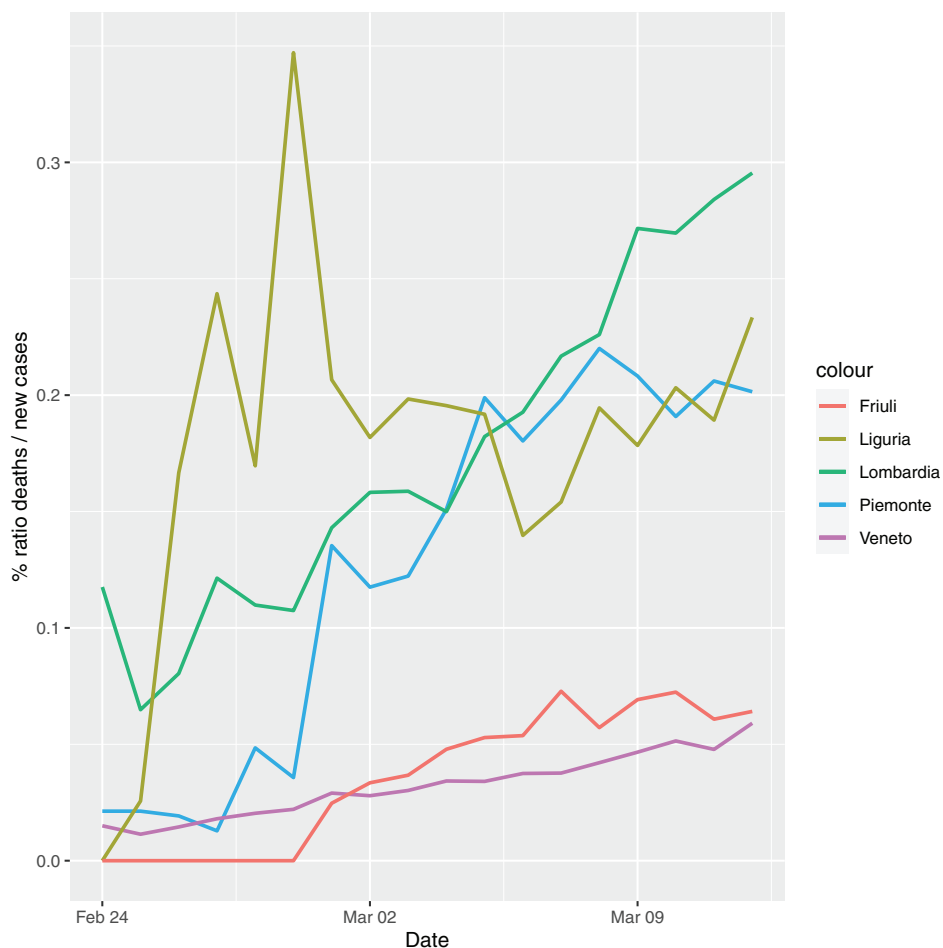
4. from a uniform distribution in  $[0,1]$ , generate  $T$  pseudorandom numbers and define the interval  $R_t = (t/T; t + 1/T)$  for  $t = 0, 1, \dots, T - 1$ , in which each  $p_j$  falls;
5. create a matching between  $R_t$  and  $I_t$  according to the following equations:

$$x_{j,t,me} = c_{T-1} - |\theta| \ln(1 - p_j) \quad \text{if } p_j \in R_0$$

$$x_{j,t,me} = c_1 - |\theta| |\ln(1 - p_j)| \quad \text{if } p_j \in R_{T-1}$$

so that a set of  $T$  values  $\{x_{j,t}\}$ , as the  $j$ th resample is obtained. Here  $\theta$  is the mean of the standard exponential distribution;

6. a new  $T \times 2$  sorting matrix  $\mathcal{S}_2$  is defined and the  $T$  members of the set  $\{x_{j,t}\}$  for the  $j$ th resample obtained in Step 5 is reordered in an increasing order of magnitude and placed in column 1. The sorted  $I_{ord}$  values (Step 2) are placed in column 2 of  $\mathcal{S}_2$ ;
7. matrix  $\mathcal{S}_2$  is sorted according to the second column so that the order  $\{1, 2, \dots, T\}$  is there restored. The jointly sorted elements of column 1 is denoted by  $\{x_{\mathcal{S}_j,t}\}$ , where  $\mathcal{S}$  recalls the sorting step;
8. Repeat Steps 1–7 a large number of times.

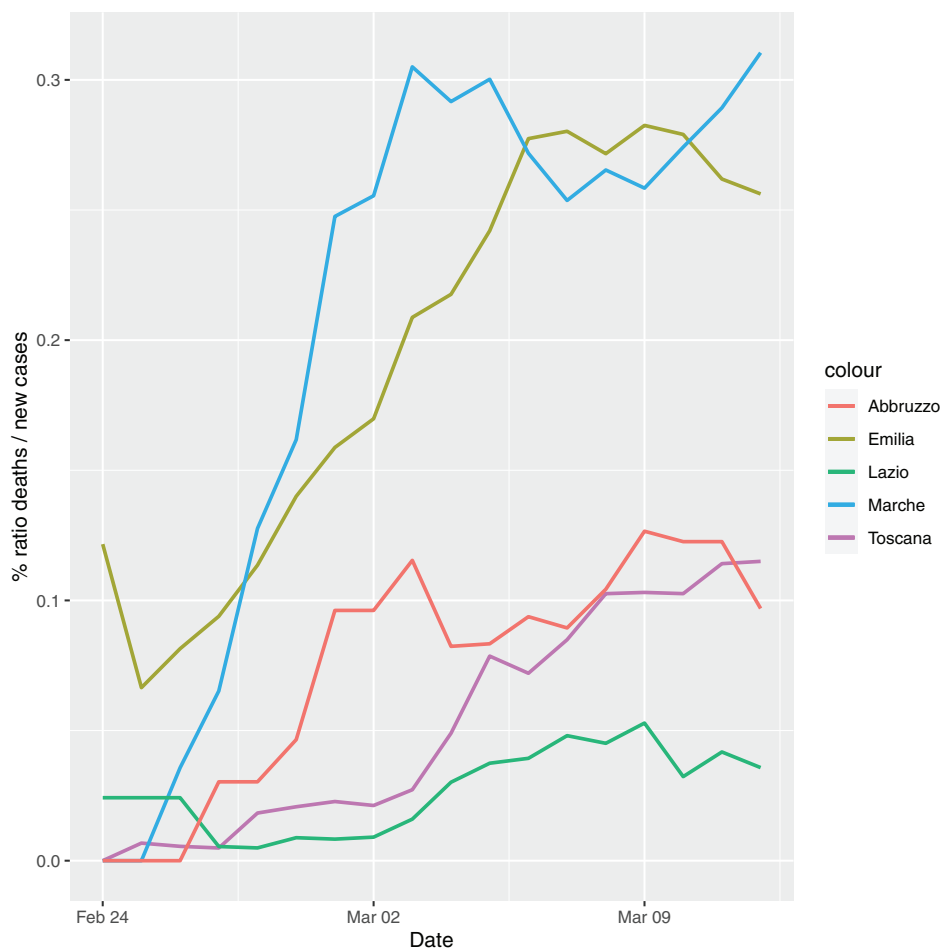


**Figure 1** Percentage ratio deaths/new cases for the following Italian regions: Piemonte, Lombardia, Veneto, Liguria and Friuli-Venezia-Giulia. [Full-size !\[\]\(fd7fe780e8fd8eece60268c87d0c3e04\_img.jpg\) DOI: 10.7717/peerj.10819/fig-1](https://doi.org/10.7717/peerj.10819/fig-1)

## THE APPLICATION OF THE MAXIMUM ENTROPY BOOTSTRAP

In what follows, the proposed procedure is presented in a step-by-step fashion.

1. For each time series  $y_t^\bullet$  and  $y_t^\circ$  the bootstrap procedure is applied so that  $B = 100$  “bona fide” replications are available as a result, that is,  $\tilde{y}_{t,b}^\bullet; b = 1, 2, \dots, B$  and  $\tilde{y}_{t,b}^\circ; b = 1, 2, \dots, B$ ;
2. for both the series, the row vector related to the last observation  $T$  is extracted, that is,  $\{v^\circ = \tilde{y}_{T,1}^\circ, \tilde{y}_{T,2}^\circ \dots \tilde{y}_{T,B}^\circ\}$  and  $\{v^\bullet = \tilde{y}_{T,1}^\bullet, \tilde{y}_{T,2}^\bullet \dots \tilde{y}_{T,B}^\bullet\}$ ;
3. the expected values, that is,  $\mathbb{E}(v^\bullet)$  and  $\mathbb{E}(v^\circ)$ , are then extracted along with the  $\approx 95\%$  confidence intervals ( $CI^\bullet$  and  $CI^\circ$ ), which are computed according to the  $t$ -percentile method. In essence, through this method, suitable quantiles of an ordered bootstrap sample of  $t$ -statistics are selected and, as a result, the critical values for the construction of an appropriate confidence interval become available. A thorough



**Figure 2** Percentage ratio deaths/new cases for the following Italian regions Emilia, Toscana, Marche, Lazio and Abruzzo. [Full-size !\[\]\(5fd6ef84f97f42d7f8b34275f1b65312\_img.jpg\) DOI: 10.7717/peerj.10819/fig-2](https://doi.org/10.7717/peerj.10819/fig-2)

explanation of the  $t$ -percentile method goes beyond the scope of this article, therefore the interested reader is referred to the excellent article by [Berkowitz & Kilian \(2000\)](#).

In particular, the lower (upper) CIs will be the lower (upper) bounds of our estimator while the quantities  $\mathbb{E}(v^\bullet)$   $\mathbb{E}(v^\circ)$  are estimated through the mean operator, that is,

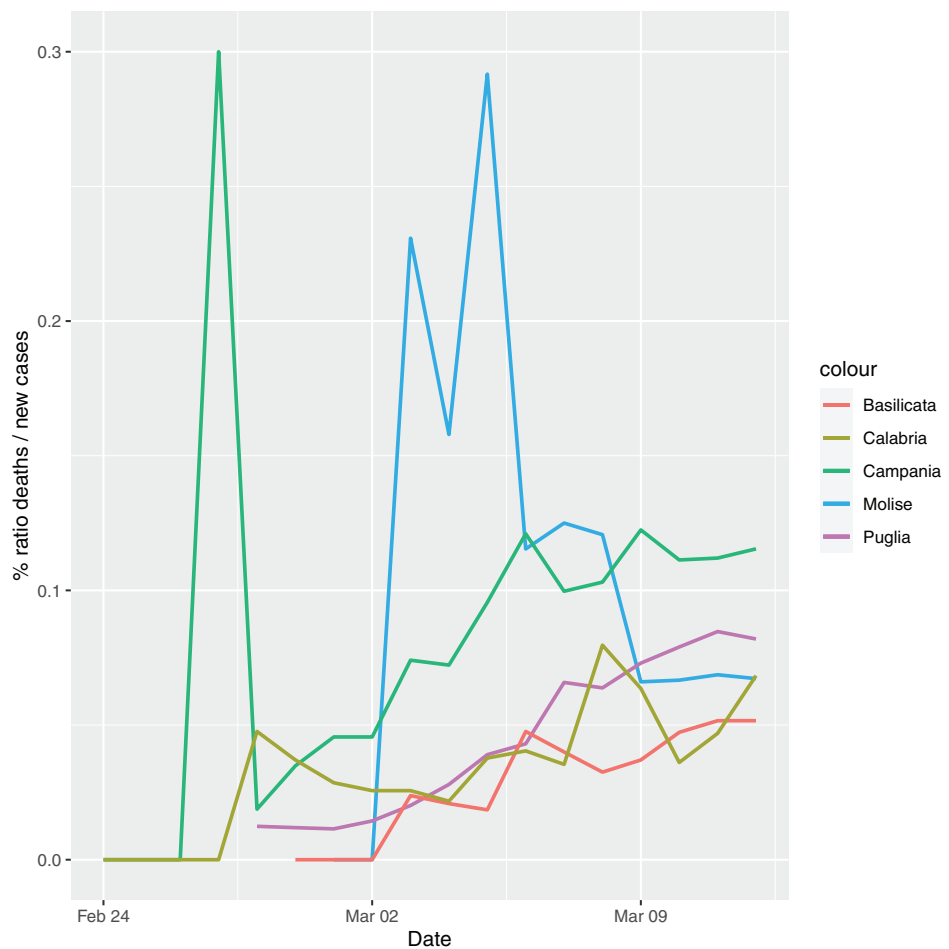
$$\mu^\circ = \sum_{j=1}^6 v_j^\circ \quad (6)$$

and

$$\mu^\bullet = \sum_{j=1}^6 v_j^\bullet \quad (7)$$

At this point, it is worth emphasizing that the procedure not only, as just seen, requires very little in terms of input data (only the time series of the positives and the



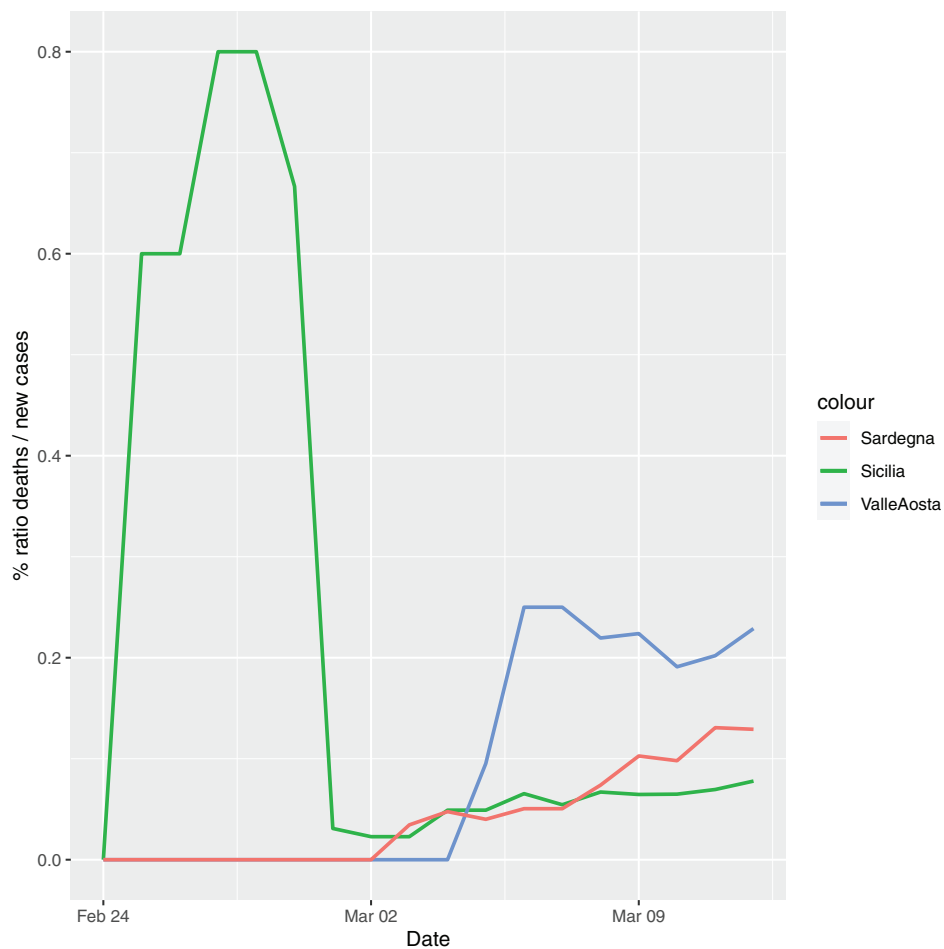


**Figure 3** Percentage ratio deaths/new cases for the following Italian regions: Molise, Campania, Puglia, Basilicata and Calabria. [Full-size !\[\]\(fcc3264021d438d9732560e78099f674\_img.jpg\) DOI: 10.7717/peerj.10819/fig-3](https://doi.org/10.7717/peerj.10819/fig-3)

deaths are required) but also can be performed in an automatic fashion. In fact, once the data become available, one has just to properly assign the time series to the subsets  $\Omega^\circ$  and  $\Omega^\bullet$  and the code will process the new data in an automatic way. The procedure is also very fast, as the computing time needed for the generation of the bootstrap samples requires—for the sample size in question—less than 2 min. Both code and data-set employed in this article have been uploaded as [Supplemental Files](#). However, the data can also be downloaded free of charge at the following web address: <https://github.com/pcm-dpc/COVID-19/tree/master/dati-regioni> (the file name is dpc-covid19-ita-regioni-20200323.csv).

## EMPIRICAL EVIDENCES

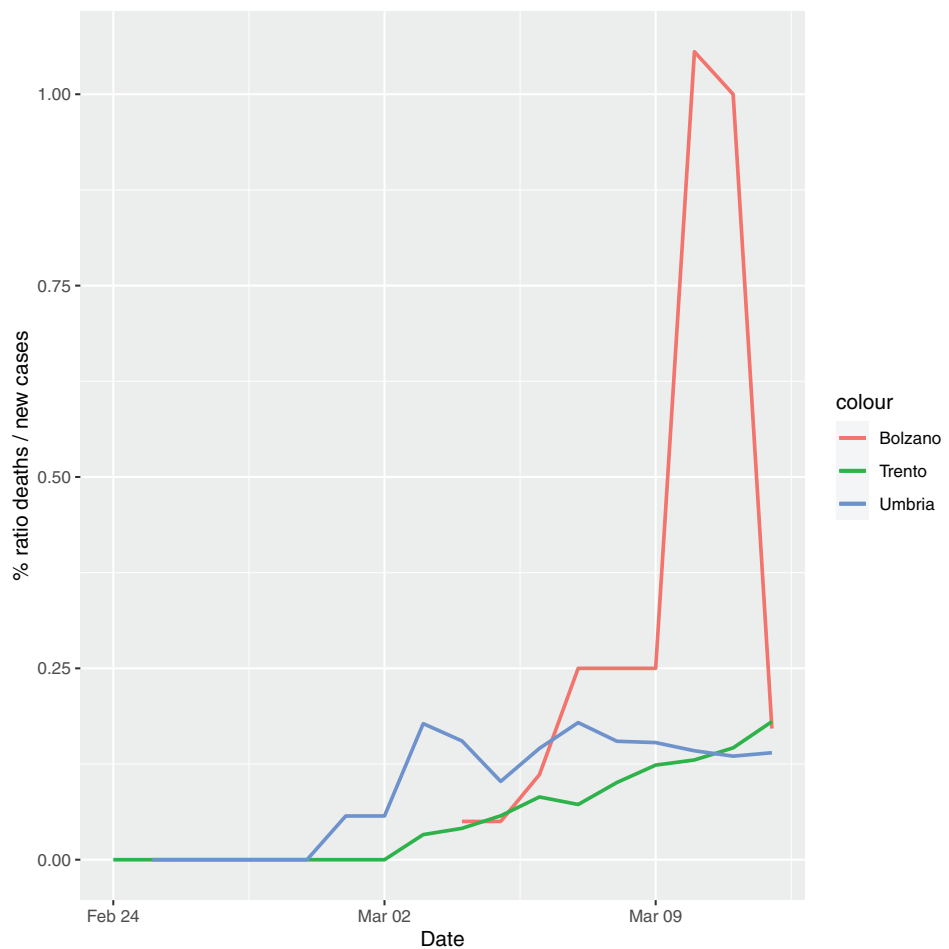
In order to give the reader the opportunity to gain a better insight on the different epidemic dynamical behaviors, in [Figure 1–5](#) the time series of the variable  $C$  (as defined in [Eq. 2](#)) is reported for each region. Note that the sudden variations noticeable in [Fig. 5](#) (Bolzano), [Fig. 4](#) (Valle D’Aosta) and [Fig. 3](#) (Molise and Campania) are due to the little number



**Figure 4** Percentage ratio deaths/new cases for the following Italian regions: Sicilia, Valle d'Aosta, Sardegna. [Full-size !\[\]\(1663bb69f307a960345edb0e712f8c02\_img.jpg\) DOI: 10.7717/peerj.10819/fig-4](https://doi.org/10.7717/peerj.10819/fig-4)

of tests administrated (i.e., the denominator of the variable  $w_T$  (2)) for these cases. In emergency situations the data are usually noisy, incomplete and might show large spikes, as in the case of Fig. 5.

That said, the main result of the article is summarized in Table 1, where three estimates of the number of positives are reported by region. The regions belonging to the set  $\Omega^\circ$  (no deaths) are in *Italics* whereas all the others, belonging to the set  $\Omega^\bullet$ , are in a standard format. In the columns “Mean” and “Lower (Upper) Bounds”, the bootstrap estimates computed according to Eqs. (6) and (7) and the Lower (Upper) Bounds the lower (upper) bootstrap CIs are respectively reported. The column denominating “Official Cases” accounts for the number of positives cases released by the Italian Authorities, whereas the column “Morbidity” expresses the percentage ratio between  $\mu^\bullet$  (6) or  $\mu^\circ$  (7) and the actual population of each region, as recorded by the Italian National Institute of Statistics. The latter source of data can be freely accessed at the web address [http://dati.istat.it/Index.aspx?DataSetCode=DCIS\\_POPRES1](http://dati.istat.it/Index.aspx?DataSetCode=DCIS_POPRES1).



**Figure 5** Percentage ratio deaths/new cases for the following Italian regions: Bolzano, Trento, Umbria. [Full-size !\[\]\(ba1b80118482ccef74a5d718ca4d7242\_img.jpg\) DOI: 10.7717/peerj.10819/fig-5](https://doi.org/10.7717/peerj.10819/fig-5)

By examining the data for the whole Country, it is clear how the data collected by the Italian Authorities on the positive cases cannot be indicative of the situation at the population level, which appear to be greater by a factor of 8. Such a consideration, straightforward from a statistical point of view, might be worth outlining as many sources of information (e.g., newspaper, TV) mainly focus on the simple count of the positive cases so that the general public might miss the magnitude of this disease. As expected, the top three regions in terms of number of infected persons are Lombardia, Emilia Romagna and Veneto, where the estimated infected population is respectively (bootstrap mean) around 45,020, 12,299 and 9,343.

On the other hand, the risk of contagion is relatively low in some regions—mostly located in the Southern part of Italy—and in the island of Sardinia.

Regarding the regions included in the subset  $\Omega^\circ$ , the application of the Piccolo distance ( $\pi$ ) has generated the associations reported in [Table 2](#).

**Table 1** Estimation of the number of people infected from Covid-19 by Italian regions. Lower and Upper Bounds are computed through the Bootstrap  $t$ -percentile method whereas the mean values is computed as in (6) and (7). The regions belonging to the set  $\Omega^\circ$  are in Italics.

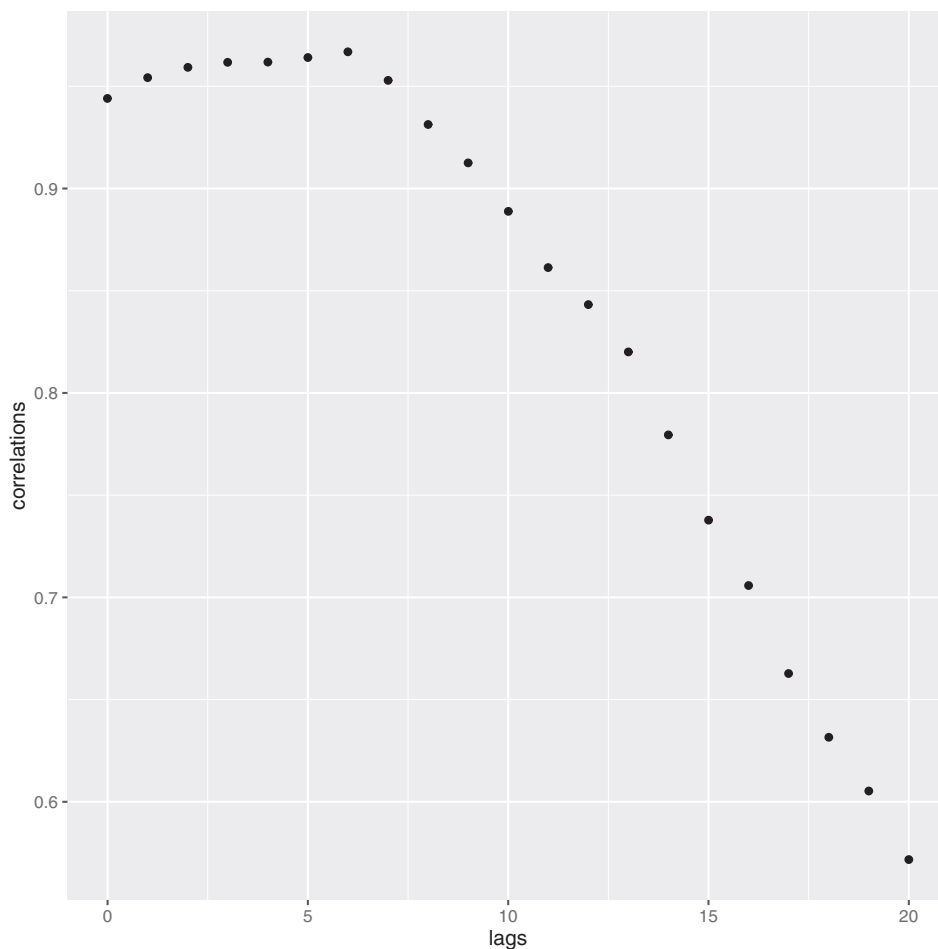
|                       | Lower bound | Mean   | Upper bound | Official cases | Population | Morbidity (%) |
|-----------------------|-------------|--------|-------------|----------------|------------|---------------|
| Abruzzo               | 526         | 600    | 807         | 78             | 1,311,580  | 0.06          |
| <i>Basilicata</i>     | 48          | 54     | 70          | 8              | 562,869    | 0.01          |
| Bolzano               | 697         | 730    | 795         | 103            | 531,178    | 0.15          |
| <i>Calabria</i>       | 182         | 238    | 493         | 32             | 1,947,131  | 0.03          |
| Campania              | 988         | 1,292  | 2,676       | 174            | 5,801,692  | 0.05          |
| Emilia Romagna        | 10,980      | 12,299 | 14,897      | 1,758          | 4,459,477  | 0.33          |
| Friuli Venezia Giulia | 983         | 1,201  | 2,514       | 148            | 1,215,220  | 0.21          |
| Lazio                 | 1,485       | 1,680  | 2,089       | 172            | 5,879,082  | 0.04          |
| Liguria               | 1,346       | 1,608  | 1,995       | 243            | 1,550,640  | 0.13          |
| Lombardia             | 37,744      | 45,020 | 49,723      | 6,896          | 10,060,574 | 0.49          |
| Marche                | 3,151       | 3,891  | 4,593       | 570            | 1,525,271  | 0.30          |
| <i>Molise</i>         | 119         | 134    | 167         | 16             | 305,617    | 0.05          |
| Piemonte              | 3,216       | 3,703  | 4,217       | 554            | 4,356,406  | 0.10          |
| Puglia                | 490         | 670    | 1,292       | 98             | 4,029,053  | 0.03          |
| <i>Sardegna</i>       | 244         | 278    | 375         | 39             | 1,639,591  | 0.02          |
| Sicilia               | 776         | 865    | 1,098       | 111            | 4,999,891  | 0.02          |
| Toscana               | 2,352       | 2,755  | 3,965       | 352            | 3,729,641  | 0.11          |
| <i>Trento</i>         | 670         | 764    | 1,028       | 102            | 541,098    | 0.19          |
| <i>Umbria</i>         | 432         | 481    | 611         | 62             | 882,015    | 0.07          |
| Valle Aosta           | 139         | 183    | 356         | 26             | 125,666    | 0.28          |
| Veneto                | 8,382       | 9,343  | 12,028      | 1,297          | 4,905,854  | 0.25          |
| Totale Italia         | 74,950      | 87,789 | 105,789     | 12,839         | 60,359,546 | 0.18          |

**Table 2** Association found between the regions belonging to  $\Omega^\circ$  and those in  $\Omega^\bullet$  according to the minimum distance  $\pi$ .

| $\Omega^\circ$ | $\Omega^\bullet$ | $\pi$   |
|----------------|------------------|---------|
| Basilicata     | Veneto           | 0.0389  |
| Calabria       | Campania         | 0.6211  |
| Molise         | Lazio            | 0.4212  |
| Sardegna       | Abruzzo          | 0.0157  |
| Trento         | Abruzzo          | 0.00186 |
| Umbria         | Sicilia          | 0.01398 |

## Model validation

The validation of the proposed approach is very simple and exploits the official Covid 19 mortality rate ( $K = \frac{\text{DEATH}}{\text{INFECTED}}$ ) issued by the WHO, which can be considered a well recognized and authoritative source. In essence, this constant—called  $K$ —has been used to make an estimate of the number of infected people (see Formula 8). Recalling that, in Italy, each and every person whose death was considered suspicious has been tested for



**Figure 6** Contraction of the infection—and time to death: delay structure.

Full-size  DOI: [10.7717/peerj.10819/fig-6](https://doi.org/10.7717/peerj.10819/fig-6)

Covid, it can be assumed the data related to these deaths to represent a population in itself (in other words, no inference procedures needed). The mortality rate, at the time of the writing of the paper, is  $K = 3.4\%$ . By applying the simple formula

$$P = \frac{\text{DEATH}}{K} \quad (8)$$

where DEATH refers to the number of deceased people, it is possible to have a rough estimate of the total positives ( $P$ ) at a population level. However, this is not the whole story. In fact, it is well known that the virus is not capable to kill a person instantly but it takes several days to do so. Therefore, Formula 8 is now rewritten to account for this temporal lag, that is,

$$P_t = \frac{\text{DEATH}_{t+h}}{K} \quad (9)$$

where  $h$  is the delay time, which can be easily estimated by considering the empirical correlation function at different lags. In Fig. 6 such a structure is reported until lag  $h = 20$ . As it can be noticed, the highest correlation is at the lag  $h = 6$ .

Recalling that in Italy the number of Covid-19 related deaths, at the date of March 12th 2020, reached the number of 2,978, by applying 9 and using  $h = 6$ , we have:  $\frac{2,978}{0.35} = 85,085.71$ . This number is very consistent with the estimate given in the article, which is 87.789.

Even considering higher lags, that is,  $h = 7, 8$ , Eq. (9) yields respectively the following number of deaths: 97,286 and 115,200. Both these results are still within the upper confidence interval given in the article ( $\approx 105,789$ ). Shorter lags can always be considered, even though the scientific community seems to exclude them.

Additionally, to validate the number of deaths due to Covid-19, the number of deaths occurred in the first quarter of 2020 with the average number of the deaths recorded in the first quarters of the years between 2015 and 2019 have been compared. It turns out that the total number of deaths ascribable to the Covid-19 is roughly equal to the difference between these two quantities.

## CONCLUSIONS

It is widespread opinion in the scientific community that current official data on the diffusion of SARS-CoV-2, responsible of the correlated disease, COVID-19, among population, are likely to suffer from a strong downward bias. In this scenario, the aim of this article is twofold: on one hand, it generates realistic figures on the effective number of people infected with SARS-CoV-2 at a national and regional level; on the other hand, it provides a methodology representing a viable alternative to those interested to apply inference procedures on the diffusion of epidemics.

This article proposes a methodology—illustrated in “Data and Contagion Indicator”—based on simple counts, that is, the number of deaths and the number of people tested positive to the virus for Italy, to:

1. provide an estimation at the national and regional level of the number of infected people and the related confidence intervals;
2. extend Eqs. (1) and (2) to those regions exhibiting no deaths as a consequence of the contraction of the Covid-19.

The entire procedure has been written in the programming language R<sup>®</sup> and uses official data as published by the Italian National Institute of Health. The whole code is available as a [Supplemental Files](#) in the Journal’s repository.

The results obtained show that, while official data at March 12th report, for Italy, a total of 12,839 cases, the people infected with the Covid-19 could be as high as 105,789. This result, along with the estimated average doubling time for the Covid-19 ( $\approx 6.2$  days), confirms that this pandemic is to be regarded as much more dangerous than currently foreseen.

As it is well known, there are many critical questions about Covid-19 that remain unanswered. One of them is related to the mortality rate  $K$ —used in Formula 8—whose

accuracy should be carefully evaluated as more and more data become available. In fact, as observed by the WHO itself,  $K$  can be biased due to the denominator of the formula being not equal to the number of infected (clinical + asymptomatic).

## ACKNOWLEDGEMENTS

The author is deeply grateful to Dr. Luigi Di Landro for the generous help in the proofreading process.

## ADDITIONAL INFORMATION AND DECLARATIONS

### Funding

The author received no funding for this work.

### Competing Interests

The author declares that they have no competing interests.

### Author Contributions

- Livio Fenga conceived and designed the experiments, performed the experiments, analyzed the data, prepared figures and/or tables, authored or reviewed drafts of the paper, and approved the final draft.

### Data Availability

The following information was supplied regarding data availability:

Data and code are freely available at GitHub:

<https://github.com/pcm-dpc/COVID-19/tree/master/dati-regioni>.

### Supplemental Information

Supplemental information for this article can be found online at <http://dx.doi.org/10.7717/peerj.10819#supplemental-information>.

## REFERENCES

- Andre's MA, Pena D, Romo J. 2002. Forecasting time series with sieve bootstrap. *Journal of Statistical Planning and Inference* **100**(1):1–11 DOI [10.1016/S0378-3758\(01\)00092-1](https://doi.org/10.1016/S0378-3758(01)00092-1).
- Barnard J, Rubin DB. 1999. Miscellanea: small-sample degrees of freedom with multiple imputation. *Biometrika* **86**(4):948–955 DOI [10.1093/biomet/86.4.948](https://doi.org/10.1093/biomet/86.4.948).
- Berkowitz J, Kilian L. 2000. Recent developments in bootstrapping time series. *Econometric Reviews* **19**(1):1–48 DOI [10.1080/07474930008800457](https://doi.org/10.1080/07474930008800457).
- Carlstein E. 1986. The use of subseries values for estimating the variance of a general statistic from a stationary sequence. *Annals of Statistics* **14**(3):1171–1179 DOI [10.1214/aos/1176350057](https://doi.org/10.1214/aos/1176350057).
- Clayton D, Hills M. 2013. *Statistical models in epidemiology*. Oxford: Oxford University Press.
- Corduas M, Piccolo D. 2008. Time series clustering and classification by the autoregressive metric. *Computational Statistics & Data Analysis* **52**(4):1860–1872 DOI [10.1016/j.csda.2007.06.001](https://doi.org/10.1016/j.csda.2007.06.001).
- Faes C, Molenberghs G, Aerts M, Verbeke G, Kenward MG. 2009. The effective sample size and an alternative small-sample degrees-of-freedom method. *American Statistician* **63**(4):389–399 DOI [10.1198/tast.2009.08196](https://doi.org/10.1198/tast.2009.08196).

- Feinstein AR, Esdaile JM. 1987.** Incidence, prevalence, and evidence: scientific problems in epidemiologic statistics for the occurrence of cancer. *American Journal of Medicine* **82(1)**:113–123 DOI [10.1016/0002-9343\(87\)90386-X](https://doi.org/10.1016/0002-9343(87)90386-X).
- Kahn HA, Kahn HA, Sempos CT. 1989.** *Statistical methods in epidemiology*. Vol. 12. Oxford: Oxford University Press.
- Koutris A, Heracleous MS, Spanos A. 2008.** Testing for nonstationarity using maximum entropy resampling: a misspecification testing perspective. *Econometric Reviews* **27(4–6)**:363–384 DOI [10.1080/07474930801959776](https://doi.org/10.1080/07474930801959776).
- Lawson AB. 2013.** *Statistical methods in spatial epidemiology*. Hoboken: John Wiley & Sons.
- Makridakis S, Hibon M. 1997.** Arma models and the box-jenkins methodology. *Journal of Forecasting* **16(3)**:147–163 DOI [10.1002/\(SICI\)1099-131X\(199705\)16:3<147::AID-FOR652>3.0.CO;2-X](https://doi.org/10.1002/(SICI)1099-131X(199705)16:3<147::AID-FOR652>3.0.CO;2-X).
- Piccolo D. 1990.** A distance measure for classifying ARIMA models. *Journal of Time Series Analysis* **11(2)**:153–164 DOI [10.1111/j.1467-9892.1990.tb00048.x](https://doi.org/10.1111/j.1467-9892.1990.tb00048.x).
- Piccolo D. 2007.** Statistical issues on the AR metric in time series analysis. In: *Proceedings of the SIS, 2007 Intermediate Conference Risk and Prediction*. 221–232.
- Piccolo D. 2010.** The autoregressive metric for comparing time series models. *Statistica* **70(4)**:459–480.
- Pueyo T. 2020.** Coronavirus: why you must act now. Available at <https://medium.com/@tomaspueyo/coronavirus373-act-today-or-people-will-die-f4d3d9cd99ca>.
- Vinod HD, López-de Lacalle J. 2009.** Maximum entropy bootstrap for time series: the meboot R package. *Journal of Statistical Software* **29(5)**:1–19 DOI [10.18637/jss.v029.i05](https://doi.org/10.18637/jss.v029.i05).