# Nowcasting Italian Covid-19 epidemic
# User guide

Alaimo Di Loro P.          Mingione M.          StatGroup19

May 21, 2020

## CoviData19 - An app for all

The CoviData19 Shinyapp is the result of the common project of a group of statisticians who share the same commitment to the social role of statistics, but are aware of the pitfalls that can stem from its miscommunication. In this regard, the final aim of the app is to provide the general public with an accessible platform for consulting and monitoring information about the Italian COVID-19 epidemic in an interactive and transparent way. All data and graphics are automatically updated at every user access with the most recent version available in the Italian Protezione Civile (IPC) Github repository (https://github.com/pcm-dpc/COVID-19).

The web-app shows both descriptive and model-based analysis, allowing the user to customize the report as he deasires. In particular, our COVID-19 Shinyapp is composed of 4 panels: (i) *"Overview"*, which provides a general description of the Italian epidemic; (ii) *"Short-term forecast"*, which allows the modeling and short-term forecast of several indicators separately, at national and regional level; (iii) *"ICU Nowcasting"*, which is specifically built to provide robust and trustworty 1-day ahead intensive care unit (ICU) hospitalizations forecast; (iv) *"Info and Credits"*, which includes some useful information about the research team and the app itself.

## Introductory notes and warnings

### Vocabulary and definitions

First of all, it is important to explain how the indicators are built and what is their meaning. A first distinction must be made between **cumulative** and **daily** indicators. The formers are always positives and increasing over time, as they are reported as the cumulative sum of the daily values. Therefore, also their variations are always positives. However, for the reasons that will be explained in the next paragraph, some of these monotone indicators may decrease at regional level. The cumulative indicators provided

by IPC are: cumulative positives, discharged recovered, deceased, swabs and tested cases[1]. Daily indicators are always positives too, but they could increase or decrease over time. As a result, their variations could be both positive or negative. Secondly, some of the indicators can be obtained by summing other basic indicators together. In particular:

- **Cumulative positives** = current positives + deceased + discharged recovered

- **Current positives** = hospitalized with symptoms + intensive care + home isolation

- **New positives** = cumulative positives of today - cumulative positives of yesterday

- **Current hospitalized** = hospitalized with symptoms + intensive care

Further important quantities that can be derived and should be monitored during a pandemic are the mortality rate, defined as the ratio between the total number of deaths and the population size; and the fatality rate, defined as the ratio between the total number of deaths and the total number of positives. Last, but not least, the user must also be aware that the data as they are reported, refer to the number of people that have been found positive to SARS-Cov-2. This does not imply that the same number of people developed COVID-19 disease.

### Data issues

The user must be aware that the data present several issues which severely affect the data quality. The general problem of the database is that information are gathered and reported at a regional level and each regional healthcare organization has different transmission and data collection systems. Moreover, the national health data system was not meant to be updated daily. This implies that there is a severe measurement error issue, also including time-dependent propensity to diagnostic testing administration, delays in reporting, saturation of diagnostic and treatment facilities, transfer of patients without notification, and plain data entry errors. Least reliable indicators include positive status (which could have up to 30 days delay in reporting), recovery (which for hospitalized patients corresponded to discharge without ascertainment of negative status in certain regions, and which was severely undercounted in other regions where access to diagnostic test was given least priority to home isolated patients). In our opinion, the most reliable indicator is the count of ICU occupancy. Other issues are structural and mostly due to the suddenness and drama of the epidemic event that caught the country healthcare system mostly unprepared. For all these reasons, discordant data transmissions and lack of coordination are fairly common events; hence, when using the app, everyone must take into account the above and the following warnings:

---

[1]The difference between swabs and tested cases corresponds to control swabs, which are done on the same individual to confirm the healing or for other necessities.

1. 1. Swabs and positive cases are not time aligned. For example, in countries like Singapore (`https://www.moh.gov.sg/covid-19`) daily data include information on total swabs tested, total unique persons swabbed as well as total swabs per 1,000,000 total population and total unique persons swabbed per 1,000,000 total population. In Italy only the total number of daily swabs is available and no linkage between swabs and tested individuals is kept in the data repository. Hence it is not possible to use the information on swabs in any modeling effort.

2. Deaths are aligned part to the date of death and part to the date of communication to the national health institute. Hence some care must be taken when daily data are considered, some variations maybe due to delayed data transmissions and not to a real peak in deaths.

3. ICU are relatively reliable although no official coordination at a national level exists on this point. Usually hospitals search for available beds in intensive care units by phone calling hospital offices directly. Then sometimes people have no access to ICU because of lack of information not because they do not need it.
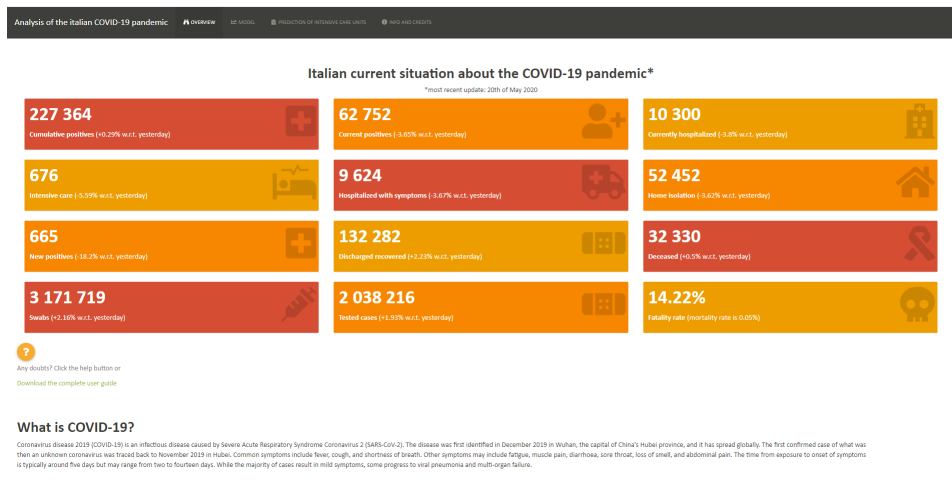
# 1 Overview



Figure 1: Overview - photograph of the Italian epidemic.

In the *overview* page, several variables are jointly recorded and visualized in order to provide an as accurate as possible picture of the Italian epidemic situation (both at national and regional level) based on data from the most recent update[2]. The page is organized in

---

[2]Note that every time a user starts the app, we check for updates and then, synchronise the data to the most recent version available on the Github repository of IPC.

six sections.

**Section 1** The first section, positioned on top of the home page (see Figure 1), presents an introductory overview in which the user is provided with all the basic information (reported values and changes with respect to the previous day) about the aggregate Italian trend of the epidemic. Definitions, meaning and computing formulas of all the indicators are readily available by clicking the help button just below it; next to it, there is the link to download this user guide. Then, a brief explanation about the COVID-19 and the SARS-Cov2 is provided.

If the interest of the user goes further, the following panels investigate in more details the temporal, compositional and spatial aspects of the available indicators.
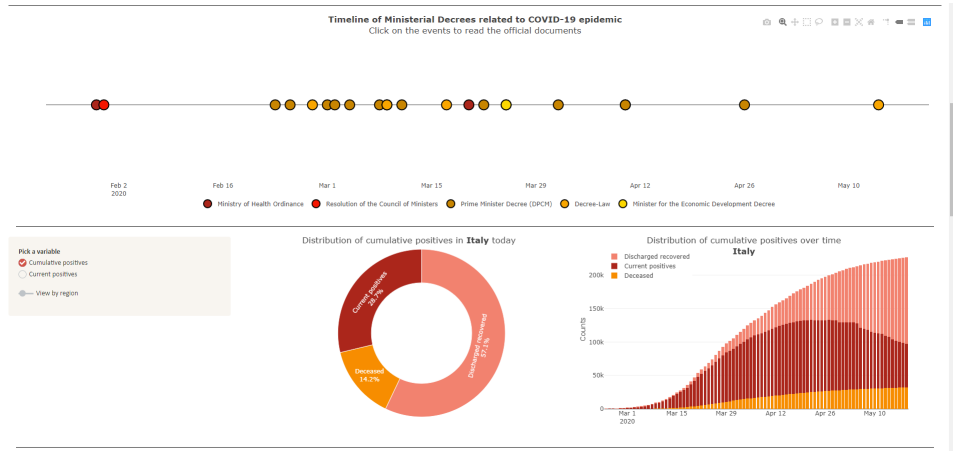


Figure 2: Overview - decree timeline and composition of derived indicators.

**Section 2** The second section (see the top part of Figure 2) shows a timeline of all the measures adopted by the government in order to keep under control the spread of the epidemic. This timeline is of the utmost importance in order to understand and appropriately interpret the temporal evolution of the indicators.

**Section 3** This section, positioned just below the timeline (see Figure 2), presents the stacked barplot (over days) of the indicators' composition in terms of the counts contributing to their definition on the right-hand side. The user can decide to visualize either the *cumulative cases* or the *current positives*, both at national or regional level. On the left-hand side, a *Donut plot* shows the relative distribution (%) of the same indicator at the current date (corresponds to the last bar of the barplot).
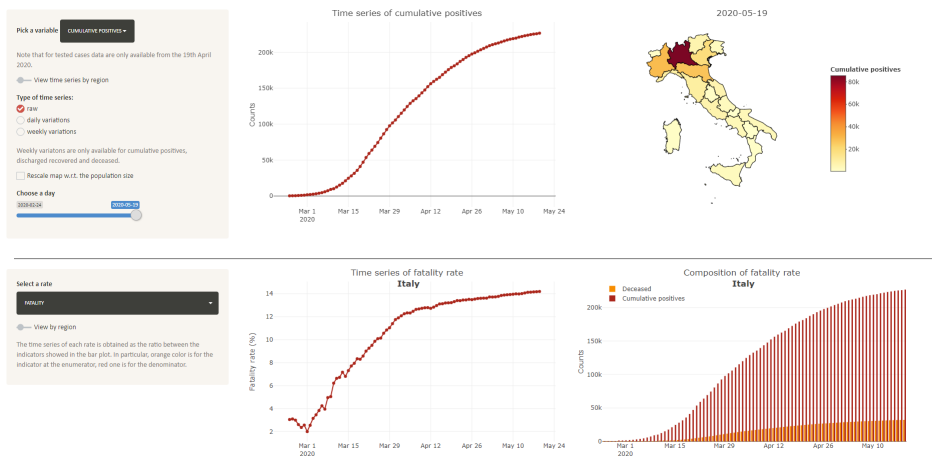
Figure 3: Overview - spatial and temporal comparison.

**Section 4** In this section (see top part of Figure 3), on the left-hand side, it is reported a time-series plot. The user can decide to visualize the temporal evolution of any of the available indicators by selecting the desired one from a drop-down menu. Since the user may want to compare the current situations of different regions, just next to it, a *Map* displays the regional distribution of the selected indicator at the current date. While the current date is the standard choice, the user can easily decide to visualize any other day through the control panel on the left. As the epidemic is known to spread out at different paces across regions, the user may want to focus on one or more specific area and compare the regional trends. With this purpose, it is possible to visualize the time series of the selected indicator at the regional level simply by checking a box on the control panel on the left and selecting the desired region. The possibility to include more than one region is available, so that trends of different regions can be easily compared on the same plot[3]. If the time series is visualized at the regional level, then the map on the right zooms on the selected region/regions and shows the distribution at the province level (only for cumulative positives, which is the only available indicator at the province level). The same control panel allows also to switch the visualization to the daily or weekly variations[4] (absolute or relative) just by checking a box, all others functionality being unchanged. As for the map, there is the possibility to rescale the distribution of the selected indicator with respect to the resident population (number of residents at 31st December 2020) so to get a relative measure the actual impact of the epidemic on the region.

**Section 5** This panel (see Figure 3 below) includes functionalities very similar to the previous one. One the left, there is a time series plot in which the user can decide to

---

[3]Comparison is allowed for up to 5 regions for visualization clarity

[4]only available for cumulative positives, discharged recovered and deceased.

visualize the temporal evolution of 4 ratios between indicators. In particular, the user can pick: (i) the ratio between new positives and daily swabs, called *positivity rate*, and its alternative formulation with the daily tested cases[5] at the denominator; (ii) the ratio between deceased and cumulative positives, called *fatality rate*; (iii) the ratio between discharged recovered and cumulative positives, called *healing rate*; (iv) the ratio between ICUs and current hospitalized, called *severity rate*. Ratios are relative measures and the user may want to get information also on its two components, separately. Hence, on the right-hand side, we included the barplot of the raw indicators by day. The same information are available at the regional level, too.

Figure 4: Overview - table with the raw data.

**Section 6** In the last section (see Figure 4), at the bottom of the page, the user can visualize the table of the regional raw data of a specific day. The first column identify the region, while the following six columns contain: cumulative positives, new cases, total deaths, new deaths, current positives and ICU. The new cases and new deaths columns are coloured according to the trend with respect to the previous day: a pale tonality for a decreasing trend and a vivid tonality for an increasing trend (respectively on a scale from yellow to red and from grey to black). The next four columns contain relative measured of the same indicators, where total cases, total deaths and current positives are reported for every thousand inhabitants while ICU are divided by the region ICU capacity. The user can decide to order the table according to any of the columns and/or to swap positions between columns if he wants to facilitate the comparison between two or more indicators.

Finally, for those who wish to analyse the raw data by themselves, the whole table (including values for all the available days) can be freely downloaded in a convenient .csv format by pushing the button on the bottom-left corner of the page.

---

[5]only available starting from the 19th of April 2020.

# 2 Short-term forecast

A simple description of the data at hand is sufficient for most users, but some other may be interested in "what is going to happen tomorrow?". Forecasting through appropriate statistical modelling is not available to everyone at a cost that is affordable. In the *Short-term forecast* section, we provide short-term (i.e. up to 15 days) forecasts of the main indicators, toghether with 99% confidence bounds to account for the estimation uncertainty. All the indicators currently provided by the *Protezione Civile Italiana* are counts, hence we considered as modeling options the Poisson $Poi(\theta_t)$ or the Negative binomial $NB(n, \theta_t)$ distribution.

The process mean $\theta_t$ is assumed to follow a pre-specified parametric form $\lambda_\gamma(t)$ over time and observations are considered independent fixed the mean function parameters $\gamma$.

In the case of cumulative counts (i.e. cumulative positives, deceased, ...), given the monotone increasing behaviour of the time series, $\lambda_\gamma(t)$ is assumed to follow the *Generalized Logistic Function* (Richard's Curve).

In the case of not cumulative counts (i.e. new positives, current hospitalized, ...), where the behavior is first increasing and then decreasing, the peak $T$ is determined in advance through kernel smoothing. Afterwards, $\lambda_\gamma(t)$ is assumed to follow a function with an increasing logistic behaviour before the observed peak and a decreasing logistic behaviour after the observed peak. Gradient and hessian have been computed analytically and coded for all the *distribution-mean function* combinations and parameters are estimated by maximizing the log-likelihood through gradient based optimization routines (i.e. gradient descent). Robust estimates of their asymptotic variance-covariance matrix are obtained using the Huber Sandwich estimator.

Prediction intervals, unavailable analytically, are obtained through a double bootstrap procedure that accounts for both the uncertainty of parameter estimation and the randomness of the observations. Model estimation and the forecasts are provided in an interactive way, with the estimation routine repeated each time the user clicks on the *Fit the model* button.



Figure 5: Short-term forecast - input.

As shown in Figure 5, the user can also decide to vary the fitting interval (exclude up to the last 15 days) and evaluate the out-of-sample predictive performance on the days excluded from the fit. While we consider this functionality both useful (always uses up-to-date data and allows for model evaluation) and enjoyable for the user experience, this does come at a cost. Indeed, sometimes the optimization routine may fail to converge to the global optimum. Generally, this is not an issue in terms of the point estimation and forecasts, but can cause extremely large (and meaningless) confidence bounds. In such cases, the user is warned through a message on the control panel on the left, can decide to not plot the confidence bounds, and is kindly encouraged to re-try and fit the same model again. On the same control panel on the left, after the model fit (see Figure 6), the user can choose the *prediction horizon term* from 1 to up to 15 days. Alternatively, he can let the whole forecast appear day after day with an animation by clicking on the *play* button on the same panel.

Schematically, the app window reports (see Figure 6 on the right):

- **Top-left**. Plot with observed and predicted values for the selected variable

- **Top-right**. Summary for the results, which includes:

  – Day of the estimated peak for the daily variations

  – Estimated upper asymptote (only for cumulative indicators)

  – Goodness of fit

  – Coverage

- **Bottom**. Table containing observed and predicted values of the past 2 days and the forecasts of the next 15 days
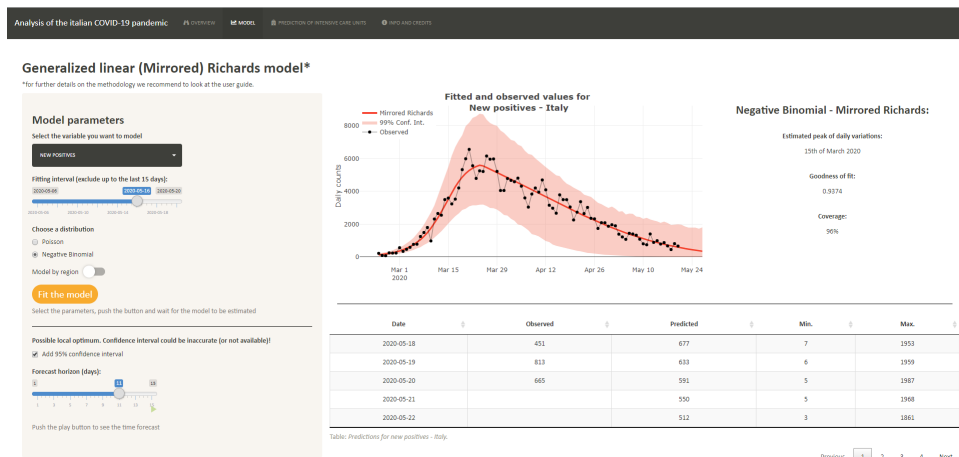


Figure 6: Short-term forecast - output.

# 3 ICU Nowcasting

Since the beginning of the epidemic, the COVID-19 has placed a huge stress on the National Health System facilities of the most hardly hit regions. As a matter of fact, a percentage between 5% and 20% of the Italian severe cases required an ICU bed, pushing the critical care capacity to (and over) the limit for regions like Lombardia and Piemonte. The overcrowding of hospital facilities and the consequent risk of breakdown of the National Health System is the greatest challenge this pandemic has put Italy through, hence monitoring the available ICU capacity is key in order to act timely and prevent this to happen.

The *ICU Nowcasting* section provides forecasts of the number of ICU beds that are likely to be required in each region on the next day. These predictions are obtained through an ensemble method, working at regional level, based on the optimal weighting of two models trained on data from the previous two weeks.

The first model is a *Poisson Generalized Linear Mixed Model* (GLMM) with region-specific random intercepts and region-specific random linear and quadratic trends. The number of residents in the region as on the 31st December 2019 is included as an offset.

The second model is a non-stationary INAR(1) model with Poisson innovations, fitted separately on data from each region, including a polynomial trend up to the cubic degree (a BIC-based selection procedure is used to drop uninformative terms region-wise).

Predictions and 99% CIs resulting from the two models are then put together by taking a weighted average, where the optimal weights are determined by a leave-last-out validation.
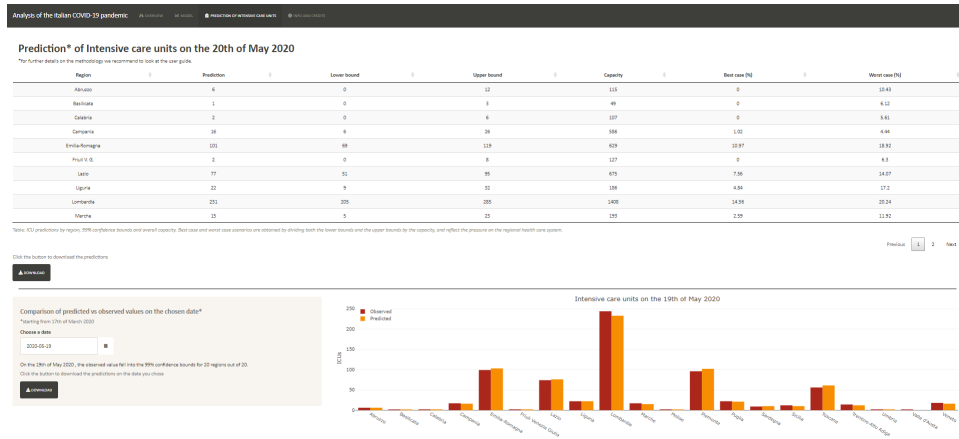


Figure 7: Prediction of ICUs by region.

The page is organized in two panels, as it is shown in Figure 7.

**Panel 1** In the first first panel, the user can observe a (downloadable) table containing the forecasts and the relative 99% confidence interval of the ICU beds to be dedicated to

COVID-19 patients on the next day, by region. The overall capacity of ICUs per region and the worst case/best case ratios are also reported.

**Panel 2**   The second panel provides the user with tools to evaluate the actual model performances on previous days. Indeed, the user can pick a date from the control panel on the left and generate a grouped bar-plot that compares predicted and observed values for the selected date. On the same control panel, just below the selected date, the user can read for how many of the regions the observed value was included in the forecast interval for the selected date. The user, if interested, can easily download the data generating the barplot by clicking the *download* button on the bottom-left corner of the page.

# 4   Info and Credits

In the last section, the user can find a brief description of StatGroup-19 and the assistants who collaborated with the research team for the web-app development. Links to the StatGroup-19's blog, facebook page and to the italian version of the ShinyApp are promptly available.

All data sources are referenced in the *Sources* paragraph, where an hyperlink connects directly to the IPC Github repository.

In the end, we wish to thank the *RStudio* software developers' team for providing such a flexible and powerful platform as *Shiny*, key tool in the production of this web-app. Moved by the aim of transparency and collaboration, all the codes used to generate the web-app are freely available in the Github repository `https://github.com/minmar94/StatGroup19`.