

COVID-19 mild cases determination from correlating COVID-line calls to reported cases

Ezequiel Alvarez^{(a)†} and Franco Marsico^{(b)★}

^(a) *International Center for Advanced Studies (ICAS) and CONICET, UNSAM, Campus Miguelete, 25 de Mayo y Francia, CP1650, San Martín, Buenos Aires, Argentina*

^(b) *Ministerio de Salud de la Provincia de Buenos Aires, Av. 51 1120, CP1900, La Plata, Buenos Aires, Argentina*

What is already known on this topic: There are always more COVID-19 cases than those reported.

What this study adds: A simple algorithm to inexpensively estimate the total number of symptomatic COVID-19 cases through a correlation in COVID-19 line phone calls data.

Abstract

Background: One of the most challenging keys to understand COVID-19 evolution is to have a measure on those mild cases which are never tested because their few symptoms are soft and/or fade away soon. The problem is not only that they are difficult to identify and test, but also that it is believed that they may constitute the bulk of the cases and could be crucial in the pandemic equation.

Methods: We present a novel algorithm to extract the number of these mild cases by correlating a COVID-line calls to reported cases in given districts. The key assumption is to realize that, being a highly contagious disease, the number of calls by mild cases should be proportional to the number of reported cases. Whereas a background of calls not related to infected people should be proportional to the district population.

Results: We find that for Buenos Aires Province, in addition to the background, there are in signal 6.6 ± 0.4 calls per each reported COVID-19 case. Using this we estimate in Buenos Aires Province 20 ± 2 COVID-19 symptomatic cases for each one reported.

Conclusions: A very simple algorithm that models the COVID-line calls as sum of signal plus background allows to estimate the crucial number of the rate of symptomatic to reported COVID-19 cases in a given district. The result from this method is an early and inexpensive estimate and should be contrasted to other methods such as serology and/or massive testing.

E-mail: † sequi@unsam.edu.ar, ★ fmarsico@mincyt.gob.ar,

COVID-19 Pandemic is impacting on World’s health and economy with an unprecedented strength [1–3]. Among the major challenges in mitigating the pandemic effect is assessing the real number of infected people at any time [3–5]. This key information is not only useful for determining health policies, but also for estimating the level of immunity in society which provides a reference framework to decide the re-opening of economic and other activities. Although the natural method for obtaining this information including mild symptom cases¹ would be to test persons upon the minimal symptom, COVID-19 is a disease with a large fraction of mild cases and therefore its cost and logistics is usually beyond the affordable. In this work we address a novel method to estimate the total number of symptomatic infected people, including mild cases, by correlating COVID-line phone calls to lab-confirmed reported cases. The main idea is that, since this is a highly contagious disease, then calls coming from infected people are proportional to the number of lab-confirmed people in that area and time, whereas other calls correspond to a background proportional to the population in the area. By measuring number of calls in different scenarios, we can fit the proportionality coefficient and estimate the total number of infected people, even though a fraction of these will not reach the threshold to be derived to a laboratory diagnostic. The idea of distinguishing a signal in a large dataset of queries is present in many schemes such as Google Flu [6] and others [7]. However, the present method is not only considerably simpler and straightforward to be implemented in any country or region, but also is specially designed for a contagious disease with the particular feature of the mild cases such as the COVID-19. The algorithm is based on a few reasonable hypotheses, is useful in any sub-testing scenario –as is the case in most of the countries–, and the presented general framework can be used in tackling other diseases and/or catastrophes, beyond the COVID-19 Pandemic. The method, being statistical, does not allow to identify the mild cases, but to estimate their number. In the following paragraphs we present the details of the algorithm, and we apply it to a real case scenario in Buenos Aires Province (PBA for its acronym in Spanish).

Along this work we model the number of calls to a COVID-line in a given district and during a given period of time by using a simple assumption of signal and background. We define signal to those calls due to real COVID-19 cases, and background to those calls due to similar symptoms and/or other causes but that do not correspond to real infections. It is key to clarify that the real COVID-19 cases that constitute the signal calls do not necessarily correspond to people whose symptoms will drive them to have a laboratory confirmation of their condition. The people who has symptoms to place the COVID-call, but that their symptoms and evolution does not reach the threshold to have a laboratory test confirmation is what we call infected with mild symptoms. As a matter of fact, we assume the compelling hypothesis that *the number of signal calls is proportional to the number of lab-confirmed cases*. On the other hand, we assume that the number of background calls is solely proportional to the district population. This assumption is reasonable as far as the studied populations have similar social behavior and there are not major changes in the social conditions –as for

¹Observe that our precise definition of mild case (see below) may not agree with others found in the literature.

instance temperature and weather— during the analyzed period of time.

Within the above hypothesis, we model the number of call received in district j in a given period of time as

$$n_C^{(j)} = \theta_P N_P^{(j)} + \theta_I N_I^{(j)}. \quad (1)$$

Where $N_P^{(j)}$ is the population and $N_I^{(j)}$ is the number of lab-confirmed infections at district j and at the given period of time. $n_C^{(j)}$ is the number of calls for district j predicted by the fit and that should be as similar as possible to the real number of calls $N_C^{(j)}$. The number $N_I^{(j)}$ corresponds to those cases whose record was opened in the studied period of time, regardless if the laboratory result was confirmed at some other time. Observe that proportionality coefficients $\boldsymbol{\theta} = (\theta_P, \theta_I)$ are independent of index j and should be fitted from the data.

There is important information to be extracted once the number of calls for each COVID-19 lab-confirmed case, θ_I , is fitted. Let f_c be the fraction of people with symptoms that contacts the Health Care System through the COVID-line, and let κ be the average of times each one of these persons calls the COVID-line. With these two variables we can now estimate the total number of lab-confirmed plus mild cases. In fact, we can now assert that the number of calls from different persons for each lab-confirmed case is θ_I/κ . Moreover, we can also estimate that for each mild calling the COVID-line, there exist other $1/f_c$ mild cases which are not calling. Henceforth, we obtain

$$\text{Total number of lab-confirmed plus mild} = \frac{\theta_I/\kappa}{f_c} N_I, \quad (2)$$

where N_I is the number of lab-confirmed cases in the studied district and period of time. The values for κ can be estimated by the telephone records or surveying on the COVID-line reported cases, whereas f_c is usually within the information available from the Health Care Administration. In any case, the outcome in Eq. 2 has many sources of intractable systematic uncertainties, and therefore its value should be understood within the corresponding caution.

Finally it is worth discussing a few details concerning the above ideas. First, observe that mild is not the same as asymptomatic. A mild case is defined as having enough symptoms to place a call to the Health Care System COVID-line, but less than the threshold required to be tested. Second, notice that, within this framework, the line dividing lab-confirmed and mild cases depends on the local Health Care System policies, since the division comes from the definition of the threshold needed to be tested. Therefore, the value of the factor accompanying N_I in Eq. 2 depends on the local Health Care System policies for each region. At last, observe that the division between mild and asymptomatic cases, although independent of policies, is rather a smooth division since depends person by person on their perception of the symptoms. The number that comes out from Eq. 2 does not consider cases which are purely asymptomatic.

As an example with a real application of the above proposal we study the Buenos Aires Province in Argentina (PBA) during the period of June 2020, in which there were not major

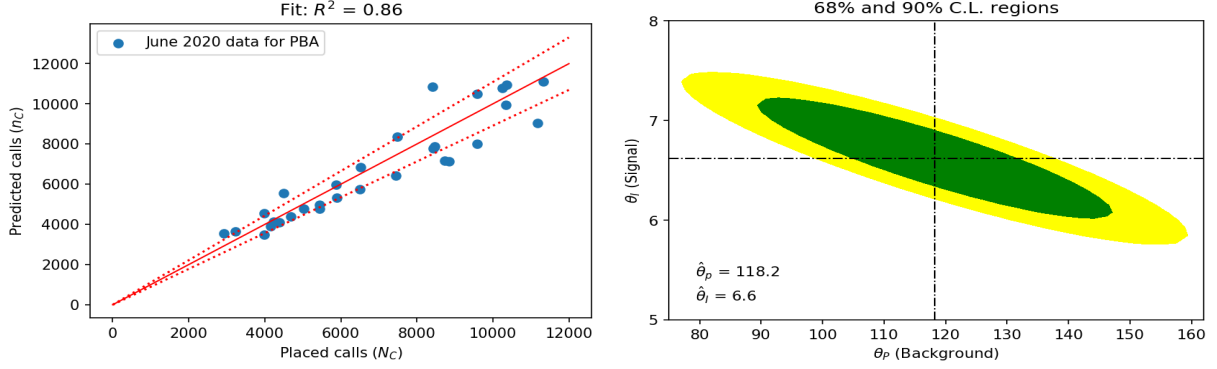


Figure 1: Left: Placed calls versus predicted calls using the result for the fit from Eq. 1 and the population and number of cases for PBA during June 2020. Each data point corresponds to a day in all PBA. Solid red line indicates the identity and the dashed red line indicates the average of the fraction of the deviation from identity for all data points, which is 11%. Right: Best fit point according to the maximum likelihood and its 68% and 90% C.L. regions (see text).

changes in weather nor in threshold and methodology for laboratory testing. During this period of time the number of lab-confirmed cases was in the order of 1000 new cases per day. The PBA Health Care System runs a COVID-line for symptoms whose local phone number is 148. The access to this local number consists in an automatic menu that derives into an operator for those cases passing the automatic menu. We take as N_C the number of calls that choose from the automatic menu the COVID symptoms option. We count calls even if they hang up before their call is taken by an operator. Given the structure of this COVID-line, the call counting at this level does not differentiate the district from which the call was placed, therefore we can only take the whole PBA as one district for this specific analysis. To have many measurements with different N_I –which is key for the fitting– we take as the period of time of each data point for Eq. 1 each full day during June. Yielding a total of 30 data points.

In order to fit the parameters θ from the data we can use the method of maximum likelihood. We should maximize the likelihood $L(\theta)$ or equivalently minimize the $\chi^2(\theta)$ defined through

$$\chi^2(\theta) = -2 \ln L(\theta) = \sum_j \frac{\left(N_C^{(j)} - (\theta_P N_P^{(j)} + \theta_I N_I^{(j)}) \right)^2}{\sigma^{(j)2}}. \quad (3)$$

Observe that also the method of least squares could have been used and results would be similar. Statistically, the variance σ^2 should correspond to the Poissonian variance of the number of calls, however there are sources of systematic errors which dominate over the statistical one. In particular, the correspondence between the number of new records opened

a given day, and the number of calls that day, yields a systematic uncertainty since both may be shifted one or two days due to intractable causes. More systematic uncertainties may come from other uncontrollable behaviors. We find that assigning a 11% systematic uncertainty to the number of calls is self-consistent with the outcome of the fit. We obtain that the best fit for Eq. 3 is through

$$\begin{aligned}\hat{\theta}_P &= 118.2 \text{ calls per 1M people} \\ \hat{\theta}_I &= 6.6 \text{ calls per confirmed case,}\end{aligned}\tag{4}$$

and at this point we obtain $\chi^2(\hat{\theta}) = \chi_{min}^2 = 36.0$. We plot in Fig. 1a the result of this fit by comparing the real number of placed calls against the predicted number of calls coming from inserting the fitted values $\hat{\theta}$ into Eq. 1. The fit yields a coefficient of determination $R^2 = 0.86$, indicating that the fit is very good, but that there are still extra sources of uncertainties, as it can be seen in the plot. To find the uncertainty in $\hat{\theta}$ we use that the contour in parameter space defined by $\chi^2(\theta) \leq \chi_{min}^2 + 2.3$ has a 68% probability of covering the true value [8]. We plot the resulting 68% and 90% confidence level contour regions in Fig. 1b. If we disregard the value of θ_P , we find $\theta_I = 6.6 \pm 0.4$ calls per confirmed case.

In the above analyzed data for PBA we know from the records that from the total lab-confirmed cases, 22% correspond to cases that entered into the Health Care System through the COVID-line. We use this to estimate that the ratio of calls from infected people corresponds to $f_c = 0.22$. On the other hand we have determined through a simple survey that each confirmed person calling the COVID-line makes on average 1.5 calls ($\kappa = 1.5 \pm 0.1$). Using this into Eq. 2 we obtain

$$\text{Total number of lab-confirmed plus mild @ PBA} = (20 \pm 2) N_I.\tag{5}$$

Where N_I is the total reported cases in PBA. Observe that in deriving this result there may have been introduced additional systematic sources through the variables f_c and κ which have not been taken into account and, in contrast to θ_I , are not controlled by the goodness of fit. These systematic may consist, for instance, in infected who may have called the COVID-line, but finally entered the system through another path; or in a different ratio of mild than severe cases calling the COVID-line. The outcome of this algorithm should be understood as an early and inexpensive estimate of the rate of symptomatic to reported COVID-19 cases. This result should be complemented with serology and/or massive testing results.

The factor 20 ± 2 in Eq. 5 is compatible with results in other parts of the World on the ratio of total infected cases found through serology to reported cases. For instance, Germany has estimated 10 times more infected than those reported by lab-confirmation [9], Spain 15 [10] and London 45 [11]. Observe, however, that these last numbers may include the asymptomatic cases, which are beyond the scope of this work and whose role in the COVID-19 Pandemic is still controversial [12–14].

Summarizing, we have designed a novel method to estimate the number of mild cases in the COVID-19 Pandemic. This method, variations and/or digital adaptations based on the idea in Eq. 1, could be useful as alternatives or complements to massive testings while

considerably less expensive. We use the correlation between the COVID-line phone calls and the number of reported cases and population in each district to estimate how many calls are made for each reported case. Using in addition the fraction of cases entering the system through the COVID-line and how many time calls are repeated by the same person, we provide an estimation for the total number of reported plus mild cases. We apply the technique to Buenos Aires Province in Argentina and find numbers compatible with other countries in which serology antibody tests have been performed.

We thank R. Vaena, E. Roulet, L. da Rold, D. de Florian, M. Szewc, F. Lamagna, S. Crespo, N. Kreplak and L. H. Molinari for useful conversations on the results of this work. A.E. thanks Development Bank of Latin America CAF and CONICET for partially funding this research.

References

- [1] Kristian G. Andersen, Andrew Rambaut, W. Ian Lipkin, Edward C. Holmes and Robert F. Garry, ‘The proximal origin of SARS-CoV-2’. *Nat Med* 26, 450–452 (2020). doi.org/10.1038/s41591-020-0820-9.
- [2] Chih-Cheng Lai, Tzu-Ping Shih, Wen-Chien Ko, Hung-Jen Tang and Po-Ren Hsueh, ‘Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) and coronavirus disease-2019 (COVID-19): The epidemic and the challenges’. *Int. J. of Antimicrobial Agents*, 55, 3, March 2020, 105924. doi.org/10.1016/j.ijantimicag.2020.105924.
- [3] Anthony S. Fauci, M.D., H. Clifford Lane, M.D., and Robert R. Redfield, M.D, ‘Covid-19 — Navigating the Uncharted’. *N Engl J Med* 2020; 382:1268-1269, doi.org/10.1056/NEJMe2002387.
- [4] Daniel P. Oran, AM, Eric J. Topol, MD, ‘Prevalence of Asymptomatic SARS-CoV-2 Infection’. *Annals of Internal Medicine*, doi.org/10.7326/M20-3012.
- [5] Sang Woo Park, Daniel M. Cornforth, Jonathan Dushoff, Joshua S. Weitz, ‘The time scale of asymptomatic transmission affects estimates of epidemic potential in the COVID-19 outbreak’. *Epidemics*, 31, June 2020, 100392, doi.org/10.1016/j.epidem.2020.100392
- [6] Ginsberg, J., Mohebbi, M., Patel, R. et al. ‘Detecting influenza epidemics using search engine query data’. *Nature* 457, 1012–1014 (2009). doi.org/10.1038/nature07634.
- [7] S.J. Yan, A.A. Chughtai and C.R. Macintyre, ‘Utility and potential of rapid epidemic intelligence from internet-based sources’. *International Journal of Infectious Diseases*, 63, 2017, pp 77-87, doi.org/10.1016/j.ijid.2017.07.020.

- [8] M. Tanabashi *et al.* [Particle Data Group], ‘Review of Particle Physics,’ Phys. Rev. D **98** (2018) no.3, 030001 doi.org/10.1103/PhysRevD.98.030001
- [9] Deutsche Welle, May 4th 2020, ([article link](#)).
- [10] El Pais, Spain, April 8th 2020, ([article link](#)).
- [11] London, May 21st 2020, Health Secretary communication ([article link](#)).
- [12] World Health Organization report, June 11th 2020 ([link to bulletin](#)).
- [13] Kenji Mizumoto, Katsushi Kagaya, Alexander Zarebski, Gerardo Chowell. ‘Estimating the asymptomatic proportion of coronavirus disease 2019 (COVID-19) cases on board the Diamond Princess cruise ship’, Yokohama, Japan, 2020. Euro Surveill. 2020;25(10):pii=2000180. doi.org/10.2807/1560-7917.
- [14] Duaa W. Al-Sadeq, Gheyath K. Nasrallah, ‘The incidence of the novel coronavirus SARS-CoV-2 among asymptomatic patients: a systematic review’. Int J Infect Dis. 2020 Jul 2, doi.org/10.1016/j.ijid.2020.06.098