

# 1 Epidemiological SIR model to study 'infodemics' about child vaccination

2 J. C. Onofre<sup>1</sup>, J. L. Acebal<sup>1</sup>

3 <sup>1</sup>*Federal Centre for Technological Education of Minas Gerais*  
4 *Av. Amazonas, 5253, 30421-169, Belo Horizonte, MG, Brazil*  
5 *julio.onofre@gmail.com, acebal@cefetmg.br*

## 6 Abstract.

7 Technological development has made the internet increasingly accessible in recent decades. However, despite  
8 the benefits, overexposure to the dissemination of information has become a social concern. The rapid dissemina-  
9 tion of information with fewer criteria plagued by rumours and fake news has compromised its accuracy, clarity  
10 and reliability. The process of rapid and massive dissemination of unreliable information has been called an info-  
11 demic because it behaves an epidemic. Google Trends is a tool developed to show the relative number of searches  
12 for a given term of interest available on the Google platform. In this work, we apply the classic SIR model of  
13 mathematical epidemiology to study the case of the time series of the frequency of searches on the controversial  
14 term 'childhood vaccination'. The aim of the approach is to assess how accurately the simple form of the epidemi-  
15 ological model can describe the infodemic process. The matter is treated as complex open system and the model  
16 parameters like infection and recovery rates from infection are supposed to vary but keeping stable along intervals  
17 of engagement and disengagement. We show that the model parameters can provide useful information about the  
18 social phenomenon in periods of social media engagement and disengagement in controversial news.

19 **Keywords:** SIR Model, infodemics, Google Trends, child vaccination

## 20 1 Introduction

21 Since December 30, 2019, when the Wuhan-China market was closed after an epidemiological alert, the  
22 SARS-Cov-2 virus, the etiological agent of Covid-19, has spread rapidly and has gone from isolated cases to a  
23 pandemic. On January 24, 2020, the first confirmed case in Europe occurred in France. And, on January 30, 2020,  
24 WHO declares pandemic status. On March 20, 2020, the Ministry of Health recognises community transmission  
25 in Brazil [1]. The pandemic has caused enormous losses to the world population, both social and economic. Since  
26 then, researchers from all over the world have made efforts to mitigate and possibly prevent the spread of the  
27 virus. The critical period between the emergence of epidemics, the development of vaccines and the (reasonable)  
28 distribution of vaccines among countries was the subject of great engagement on social media.

29 Amid discussions about possible treatments for Covid-19, the first vaccines emerged in less than a year.  
30 Quickly, the effectiveness of treatments and vaccines became the subject of intense debates marked by misinforma-  
31 tion and politicisation. In addition to a process of spreading fake news through social networks, news without  
32 criteria was given by media vehicles with sensationalism. In this environment, fake news quickly became the topic  
33 of the day which confirms observations that fake news spreads significantly faster and more widely than true news  
34 in many categories of information [2, 3]. In particular, if this category involves public health, the impact of uncer-  
35 tainties, rumours and fake news on information is a feedback factor in the process of spreading the news and the  
36 debate, as scientific sources do not always have the desired reach, are understood or even believed. Engagement  
37 and its social effects were enormously intensified when childhood vaccination emerged as a possibility.

38 According the WHO (World Health Organisation), the social phenomenon characterised by an excess of  
39 accurate or inaccurate information which obfuscates access to reliable sources and guidelines is called an infodemic  
40 [4]. Infodemia gets its name because of the similarity between the dissemination of information and the spread  
41 of a disease in an epidemic [5]. This intense availability of information via news distribution channels and social  
42 networks is strongly fuelled by polemics, debates and the legitimate search for reliable sources. Such ebullition  
43 results in great social engagement on the topic. As a consequence of this process, a considerable increase in  
44 searches for information on the internet is expected. Therefore, the increase in searches for themes through the  
45 Google platform can be considered an indicator of the infodemic process associated with these themes. A Recent

works have been studied infodemic processes with use of Sir model [6].

Google Trends is a feature present on the Google services platform for evaluating trends in searches carried out on a topic of interest. It is a new and open-access tool that allows users to interact with research data on the Internet, which can provide deep insights into population behaviour and health-related phenomena [7]. The tool shows the relative frequency over time as a ratio of the maximum frequency reached (fixed at 100 points) in the period of interest for a geographic region that can be selected. Hence, Google Trends serves as a time series recorder of relative search frequencies for a term in this region. Through this time series, it is possible to evidence pulses of engagement and disengagement triggered by social media facts and, following, dumped by information campaigns. Consequently, the tool provides a method for evaluating the infodemic phenomenon.

The first model in mathematical epidemiology of the spread of smallpox was proposed by Daniel Bernoulli in 1760. In 1906, W. H. Hamer suggested that the spread of infection should depend on the number of susceptible and infected individuals. Finally, in 1927, Kermack and McKendrick proposed the SIR model, a model involving differential equations, to describe the spread of infectious diseases among populations of susceptible, infected, and recovered individuals [8, 9]. Since then, the SIR model has been successfully used to predict the effect of epidemics, as well as to provide a tool to study control actions such as vaccination and social distancing.

In the present work, a case study of an infodemic associated with engagement occurred on the term 'child vaccination' <sup>1</sup> is carried out. The time series of searches for the term on Google obtained via Google Trends is partitioned into subintervals associated with pulses or infodemic waves. In addition, each subinterval is also divided into two periods, the period of engagement associated with increasing values of social engagement, followed by the period of disengagement, characterised by a decrease in the social engagement. The rise of engagement is associated with polemics, discussions with the dissemination of non-reliable or non sufficiently verified information. The disengagement period is associated with the progressive action of information media, organisations and health agencies clarifying the matter. The system is considered open since the total number  $N$  of individuals can vary, and complex (non autonomous), since the SIR parameters of the *i.e.* model, the contagion and recovery rates can also vary over time. However, they must be constant during periods of engagement or disengagement. Under such assumptions, we evaluate the description of the system with the SIR model in the study case of the infodemics associated to the term 'childhood vaccination' in Brazil. In section 1, we present basic information about the SIR model. Data and methods are discussed in the Section 3. Results and discussion are provided in Section 4. In Section 5, we draw the conclusions.

## 2 The SIR model

The SIR model is a mathematical description for the evolution of epidemics from a class called compartmental or population models. The model and its variations have been useful for decades for its simplicity and accuracy in describing the spread of diseases in populations and for assisting in inferences about epidemics and methods of control.

In the SIR model, there are three populations  $S(t)$ ,  $I(t)$  and  $R(t)$ , respectively, the population of susceptible individuals, the population infected by the disease and the population recovered from the disease. When the system is considered closed, the sum of the three populations equals a constant,  $S(t) + I(t) + R(t) = N$ , the total population number. The evolution of the populations of the SIR-model in time is governed by dynamic rules that can be written in continuous time ( $t \in \mathbb{R}$ ), by using differential equations, or in discrete time, ( $t \in \mathbb{Z}$ ), using difference equations.

Figure 1 shows a discrete-time SIR model. The transitions of individuals between populations in a time interval correspond to fractions of these populations. Thus, at a time  $t$ , the nonlinear interaction of infection indicated by the dashed line occurs between individuals from populations  $S(t)$  and  $I(t)$ . The values of these populations make up the rate  $\frac{\beta S(t)I(t)}{N}$  of susceptible individuals that become infected in the time interval and change to the infected population. The  $\beta$  parameter is proportional to the rate of social contact times the fraction of contacts that actually produce contagion. Likewise, the population  $I(t)$  gives up a fraction of its individuals at a rate  $\gamma I(t)$  to the population of recoveries. The  $\gamma$  parameter is the recovery rate, which means it is the inverse of the mean recovery time from the disease.

<sup>1</sup>Search performed for the term in portuguese: 'vacinação infantil'.

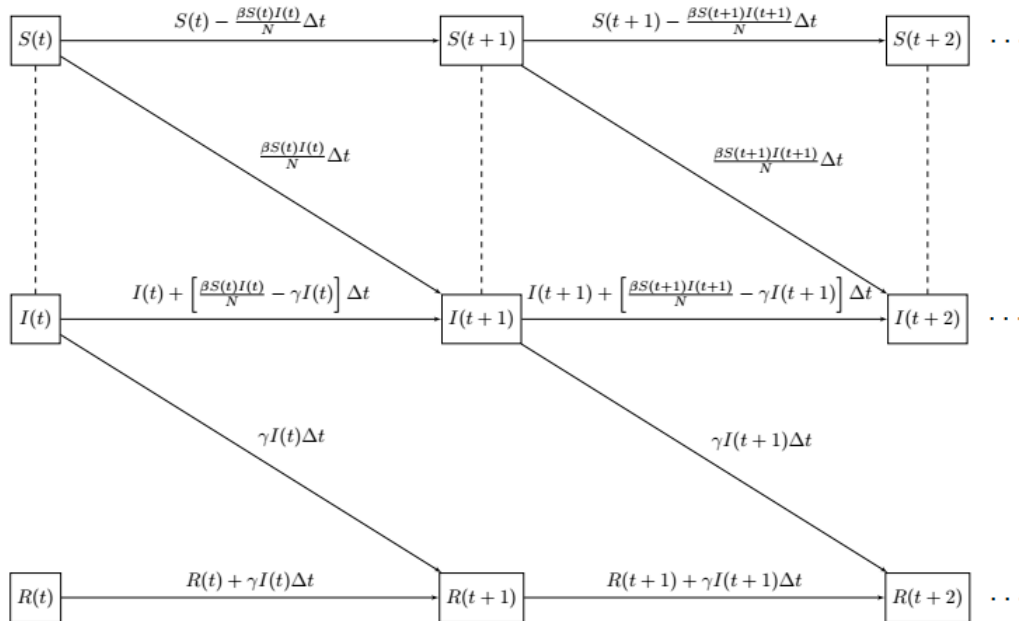


Figure 1. Schematic diagram of the SIR model showing the flow of individuals between populations, according to the parameters of the model and also of their populations.

94 The dynamics of the discrete-time SIR model can be represented through a system of difference equations:

$$\begin{cases} S(t+1) = S(t) - \left[ \frac{\beta S(t)I(t)}{N} \right] \Delta t, \\ I(t+1) = I(t) + \left[ \frac{\beta S(t)I(t)}{N} - \gamma I(t) \right] \Delta t, \\ R(t+1) = R(t) + [\gamma I(t)] \Delta t. \end{cases} \quad (1)$$

95 The time step  $\Delta t$  in the equations is, in principle, arbitrary and can be smaller than the time step of the data  
96 to be fitted.

In addition to the model parameters, to determine the intensity of the infectious process it is commonly used the *Basic Reproduction Number*,  $R_0 = \frac{I(1)}{I(0)}$ , which represents the average number of newly infected in a subsequent period from a single individual infected in an earlier period at the onset of the epidemic process[10]. Thus, the populations at this start are  $S(0) \approx N$ ,  $I(0) \approx 1$  and  $R(0) = 0$ . In this case, it follows from the equation for  $I(t)$  in the (1) and the from the definition of  $R_0$  that

$$R_0 = \frac{I(1)}{I(0)} = 1 + \left[ \frac{\beta S(0)}{N} - \gamma I(0) \right] \Delta t \approx 1 + \left[ \frac{\beta}{\gamma} - 1 \right] \gamma \Delta t, \quad (2)$$

in which we apply the condition of the beginning of the pandemic. When the time scale is chosen to have  $\gamma \Delta t = 1$ , we have,

$$R_0 = \frac{\beta}{\gamma}. \quad (3)$$

97 Hence, this number is expressed in terms of the model parameters. Although there are more general definitions  
98 for the reproduction number as  $R_t = \frac{I(t+1)}{I(t)}$  which expresses a proportion that varies with the evolution of the  
99 process, we will restrict ourselves to  $R_0$ , and we will understand that the process that occurs at any time  $t$ , this  
100 value is compatible with an epidemic that started in the past with  $R_0$  given by the equation (3).

### 101 3 Data and Method

102 To evaluate the infodemic process associated with the theme ‘childhood vaccination’, a study case was carried  
103 out with data collection from the Google Trends service with a geographic filter defined for Brazil, in a period

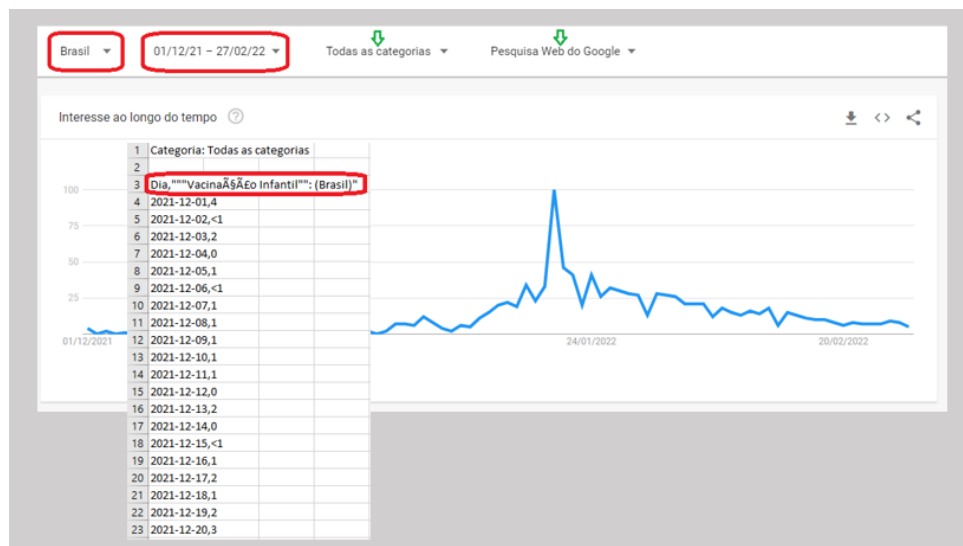


Figure 2. Google Trends Data Sheet and Google Trends Graph - Term "Vacinação Infantil"

104 between the dates 12/01/2021 to 27/ 02/2022. The time-series chart of the frequency of daily surveys on this  
 105 subject and the spreadsheet with the corresponding data were obtained (Fig. 2). The initial analysis of the data  
 106 revealed a certain degree of complexity without, at first, an identification of a typical epidemic process, but perhaps  
 107 the superposition of several wave pulses, composing a series of events (Fig.3).

108 The relative frequency of searches in Google was associated with the variable number of infected  $I(t)$  of an  
 109 epidemic in the SIR model (Section 2). The observation of the data raised the fundamental hypothesis that the  
 110 greater or lesser engagement in the process occurred in periods. In this way, the number  $N$  of people involved in  
 111 the process would vary from a period, according to social media events like the strength of engagement, campaigns  
 112 of information etc. The rates of the SIR model  $\beta$  and  $\gamma$  which are, respectively, the fraction of susceptible to engage  
 113 in the subject and the rate of disengagement also can vary at each period. This effect of changing parameters over  
 114 time is expected in real epidemics as control actions such as social distancing, sanitary measures and vaccines are  
 115 of common use whenever possible. In this scenario, the time series was partitioned into intervals whose logical  
 116 sequences are related to processes of engagement (increase) and disengagement (decrease) in each period (see  
 117 Fig.3).

118 The partition of data took place into 8 subintervals, each of them was divided into two subsequent periods  
 119 (see Table 1), corresponding to the engagement period, followed by the disengagement period. The engagement  
 120 period would correspond to the ascending data sequence that is supposedly the result of a media event that drives  
 121 an engagement process (infection rising). The disengagement period of these sub-intervals would correspond to  
 122 the disengagement process that would result from actions to clarify public opinion (epidemic control actions).

123 For each increasing or decreasing period of each subinterval, a curve of the infected population of the SIR  
 124 model is fitted to the data. Thus, for each case, a set of model parameters  $N$ ,  $\beta$ , and  $\gamma$  is obtained via optimisation  
 125 of the SIR model implemented in Python language. The method used to solve the ODE system was through  
 126 the *odeint* function of the *scipy.integrate* sub-package and the optimisation method used was the *Trust Region*  
 127 *Reflective* minimising the *least-square* objective function. After determining the parameters, the basal reproduction  
 128 number  $R_0$  of each case was obtained in order to characterise the intensity with which the infodemic process  
 129 developed in each analysed period.

Table 1. Description of the sub-intervals of the study period partition into engagement and disengagement pulses in terms of days of each sub-interval, the engagement and disengagement stage and their respective peak.

Period	Days	Peak	Period of
$p_1$	22° ao 24°	12/24/21 (24° day)	Engagement
$p_2$	24° ao 26°	12/24/21 (24° day)	Disengagement
$p_3$	26° ao 28°	12/28/21 (28° day)	Engagement
$p_4$	28° ao 31°	12/28/21 (28° day)	Disengagement
$p_5$	31° ao 37°	01/06/22 (37° day)	Engagement
$p_6$	37° ao 40°	01/06/22 (37° day)	Disengagement
$p_7$	40° ao 51°	01/20/22 (51° day)	Engagement
$p_8$	51° ao 89°	01/20/22 (51° day)	Disengagement

130 **4 Results**

131 The model parameters were obtained via optimisation in each sub-interval, distinguishing the periods of en-  
 132 gagement and disengagement, according to the initial assumption of having periods, respectively, of misinforma-  
 133 tion, rumours and fake news and another with educational actions of information through the media and agencies.  
 134 In each of these periods of engagement and disengagement, the parameters of the SIR model were optimised  
 135 via Trust Region Reflective optimisation method and had values of  $R_0$  and Root Mean Square Error (NRMSE)  
 136 normalised to mean of data in the period.

137 The values obtained by optimisation via the Trust Region Reflective method were described in the 2 table. The  
 138 curves for the infected population fitted the data, at least, with a stronger correlation (Fig. 4). Finally, Table 2 shows  
 139 the values of  $R_0$  at each stage of engagement and disengagement of each sub-interval. In each period, the curve  
 140 behaviour is defined as a result of two tendencies of engagement and disengagement expresses mathematically  
 141 by the values of, respectively, the rate of infection  $\beta$  and recovering from infection  $\gamma$ , as well as the number  $N$   
 142 of individuals denotes a proportion of the number of individuals that are able to be interested in the idea. This  
 143 fact explains the high values of  $N$  in some periods. Although official actions on media tend to promote reliable  
 144 information and disengagement, the action of education may sometimes reach a greater number of individuals. In  
 145 some intervals, the high values of  $R_0$  are explained by the small number of recovery rates, indicating that in those  
 146 periods, the trend is markedly pointed to rapid engagement, with few or almost non-existent educational actions.

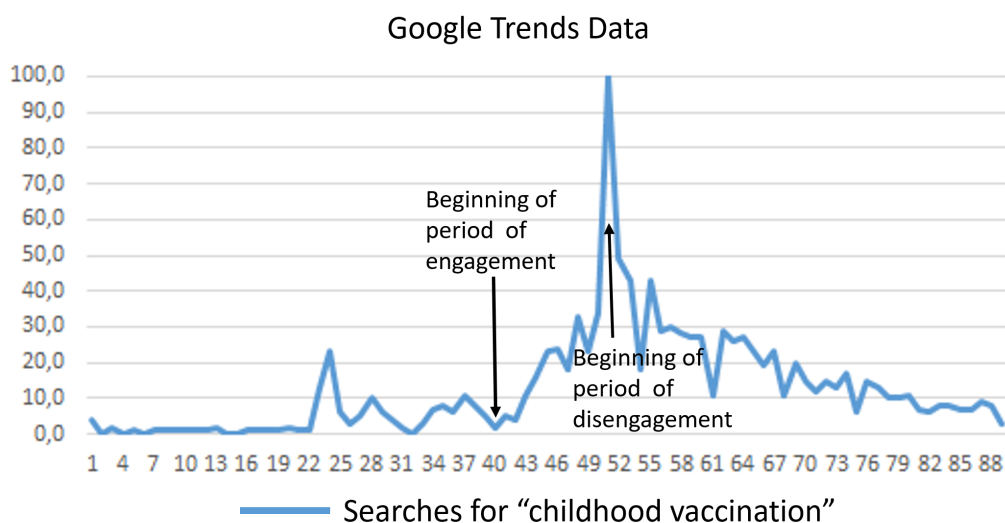


Figure 3. Google Trends search on the term "Childhood Vaccination".

Table 2. Parameter values  $N$ ,  $\beta$ ,  $\gamma$  optimised by calculating the approximate value of  $R_0$  in each sub-period and process steps.

	$p_1$		$p_2$		$p_3$		$p_4$	
$N$	43	4463	29	19	22	11	202	36705
$\beta$	2.311	$9.787 \cdot 10^{-4}$	$7.987 \cdot 10^{-1}$	0.8819	0.6429	0.6537	$3.399 \cdot 10^{-1}$	$1.148 \cdot 10^{-5}$
$\gamma$	0.4521	$8.2945 \cdot 10^{-1}$	$3.7566 \cdot 10^{-26}$	0.7359	0.145	0.659	$4.5478 \cdot 10^{-24}$	$7.6987 \cdot 10^{-2}$
$R_0$	5.113	0.0012	$2.126 \cdot 10^{25}$	1.198	4.435	0.992	$7.474 \cdot 10^{22}$	0.0001
NRMSE	0.4715	0.2314	0.1	0.1599	0.2614	0.2276	0.5702	0.4748
$\rho_S$	1.0	1.0	1.0	1.0	0.7	1.0	0.9317	0.9028

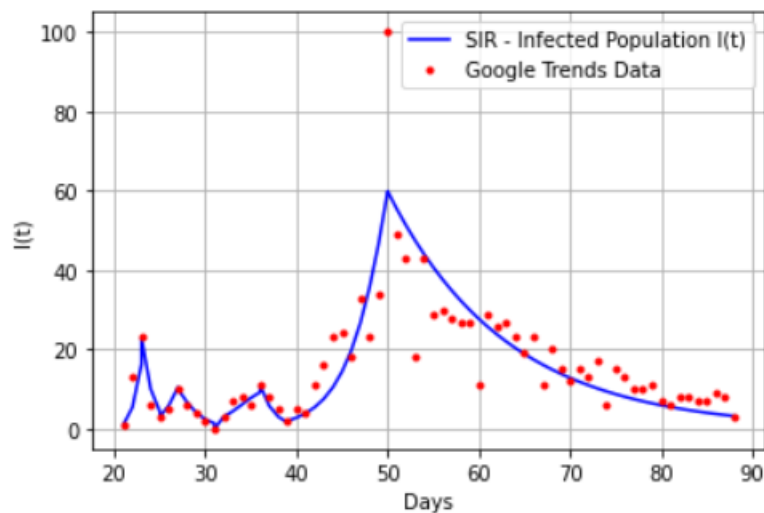


Figure 4. Frequency data of relative searches performed on Google in the analyzed period and curve adjusted by optimizing the parameters of the SIR model.

## 147 5 Conclusions

148 The SIR model, in its simplest form, was able to satisfactorily describe the infodemic process in this case  
 149 study. This fact can serve to corroborate the epidemic character of the dissemination of information.

150 The adopted assumptions with the partition of the complex process observation time into sub-intervals, and  
 151 into periods of engagement and disengagement proved to be effective. On the one hand, it was based on the  
 152 realistic varying characteristics of the epidemics process over time due to vaccination, control social distancing.  
 153 On the other hand, fitting the model to the data revealed that the hypothesis was adequate.

154 Other refinements can be made in the methodology of this work to make the method an instrument of social  
 155 analysis. One possible refinement is the automatic identification of sub-intervals with computational intelligence.  
 156 Also, other optimisation methods can be used with more complex data, such as Genetic Algorithms.

157 **Authorship statement.** This section is mandatory and should be positioned immediately before the References  
 158 section. The text should be exactly as follows: The authors hereby confirm that they are the sole liable persons  
 159 responsible for the authorship of this work, and that all material that has been herein included as part of the present  
 160 paper is either the property (and authorship) of the authors, or has the permission of the owners to be included here.

## 161 References

162 [1] N. H. THEY. Uma breve linha do tempo. *UFRGS Litoral*, 2020.

- 163 [2] Z. Zhao, J. Zhao, Y. Sano, O. Levy, H. Takayasu, M. Takayasu, D. Li, J. Wu, and S. Havlin. Fake news  
164 propagates differently from real news even at early stages of spreading. *EPJ data science*, vol. 9, n. 1, pp. 7, 2020.
- 165 [3] S. Vosoughi, D. Roy, and S. Aral. The spread of true and false news online. *science*, vol. 359, n. 6380, pp.  
166 1146–1151, 2018.
- 167 [4] W. H. O. (WHO) and others. Understanding the infodemic and misinformation in the fight against covid-19.  
168 *Retrieved March*, vol. 10, pp. 2021, 2020.
- 169 [5] J. Zarocostas. How to fight an infodemic. *The lancet*, vol. 395, n. 10225, pp. 676, 2020.
- 170 [6] M. Cinelli, W. Quattrocioni, A. Galeazzi, C. M. Valensise, E. Brugnoli, A. L. Schmidt, P. Zola, F. Zollo, and  
171 A. Scala. The covid-19 social media infodemic. *Scientific reports*, vol. 10, n. 1, pp. 1–10, 2020.
- 172 [7] S. V. Nuti, B. Wayda, I. Ranasinghe, S. Wang, R. P. Dreyer, S. I. Chen, and K. Murugiah. The use of google  
173 trends in health care research: a systematic review. *PloS one*, vol. 9, n. 10, pp. e109583, 2014.
- 174 [8] F. Brauer. Mathematical epidemiology: Past, present, and future. *Infectious Disease Modelling*, vol. 2, n. 2,  
175 pp. 113–127, 2017.
- 176 [9] C. Colombo and M. Diamanti. The smallpox vaccine: The dispute between bernoulli and d’alembert and the  
177 calculus of probabilities. *Lettera Matematica*, vol. 2, n. 4, pp. 185–192, 2015.
- 178 [10] S. Contreras, H. A. Villavicencio, D. Medina-Ortiz, C. P. Saavedra, and Á. Olivera-Nappa. Real-time esti-  
179 mation of  $r_t$  for supporting public-health policies against covid-19. *Frontiers in public health*, vol. 8, pp. 556689,  
180 2020.