

1 **Title:** Predicting social distancing index during COVID-19 outbreak through online  
2 search engines trends

3 P. C. Lins-Filho<sup>1</sup>, M. M. S. Araújo<sup>1</sup>, T. S. Macêdo<sup>1</sup>, A. K. A. Ferreira<sup>1</sup>, M. C. F. Melo<sup>1</sup>,  
4 E. L. M. S. Silva<sup>1</sup>, J. L. M. Freitas<sup>1</sup>, A. F. Caldas Jr<sup>1,2</sup>.

5 <sup>1</sup> Universidade Federal de Pernambuco, Recife, Pernambuco, Brazil

6 <sup>2</sup> Universidade de Pernambuco, Camaragibe, Pernambuco, Brazil

7 **Corresponding Author:**

8 PhD Arnaldo de França Caldas Jr

9 Address for correspondence: Estrada de Aldeia, Km 13, Prive Portal de Aldeia, Aldeia,  
10 Camaragibe, Pernambuco, Brasil. Tel: +55 (81) 999713652 E-mail:  
11 caldasjr@aldeia.com.br

12 **Running Title:** Predicting SDI with online search volume

## 13 **Summary**

14 Online-available information has been considered an accessory tool to estimate  
15 epidemiology and collect data on diseases and population behavior patterns. This study  
16 aimed to explore the potential use of Google and YouTube relative search volume to  
17 predict social distancing index in Brazil during COVID-19 outbreak and verify the  
18 correlation between social distancing measures with the course of the epidemic. Data  
19 concerning the social distancing index, epidemiological data on COVID-19 in Brazil and  
20 the search engines trends for “Coronavirus” were retrieved from online databases.  
21 Multiple linear regression was performed and resulted in a statistically significant model  
22 evidencing that Google and YouTube relative search volumes are predictors of the social  
23 distancing index. The Spearman correlation test revealed a weak correlation between  
24 social distancing measures and the course of the COVID-19 epidemic. Health authorities  
25 can apply these data to define the proper timing and location for practicing appropriate  
26 risk communication strategies.

27 **Keywords:** COVID-19; Social Isolation; Internet; Consumer Health Information.

## 28 **Introduction**

29           The World Health Organization have recently declared South America as the new  
30 coronavirus epicenter, mainly because of the situation in Brazil that registers the most  
31 cases and deaths in Latin America [1]. Given the COVID-19 pandemic, robust risk  
32 communication is urgently needed particularly in the most affected countries [2]. Internet  
33 query platforms, which allows to interact with internet-based data, have been considered  
34 a source of potentially useful and accessible resources, especially aimed to identify  
35 outbreaks and implement intervention strategies [3, 4]. Online-available information has  
36 been considered as a surrogate tool for estimating epidemiology and gathering data about  
37 patterns of disease and population behavior [5-7].

38           As the online queries on COVID-19 increases globally reflecting the interest of  
39 people to be aware about this emerging infectious disease, mining online data and search  
40 patterns on electronic resources might provide a better support to manage this worldwide  
41 health crisis [8]. Internet searches and social media data have been reported to correlate  
42 with traditional surveillance data and can even predict the outbreak of disease epidemics  
43 several days or weeks earlier [9]. Evidence points that Google Trends could potentially  
44 define the proper timing and location for practicing appropriate risk communication  
45 strategies for affected populations and be employed to predict outbreak trends of the novel  
46 coronavirus [2, 5, 10].

47           Previous investigations reported the use of Internet search engines as source of data  
48 for public health surveillance and diseases incidence prediction worldwide, as zika, in  
49 Brazil and Colombia [11]; influenza, in the United States [12]; malaria, dengue fever and  
50 chikungunya, in India [13]; and Middle East respiratory syndrome, in Korea [14].  
51 However, there are no reports on assessing the relative search volume (RSV) of search  
52 engines to predict social distancing behavior during infectious diseases outbreaks.

53 Predictions might support in health resource management and planning for  
54 prevention purposes [5]. As COVID-19 treatment protocol is still uncertain it is especially  
55 important to prevent the virus dissemination in society [15]. Currently, the virus spread  
56 prevention approaches focus on hand hygiene, social distancing and quarantine [16].  
57 Social distancing is designed to reduce interactions between people in a broader  
58 community, in which individuals may be infectious but have not yet been identified hence  
59 not yet isolated. This measure is particularly useful in settings where community  
60 transmission is believed to have occurred, but where the linkages between cases is  
61 unclear, and where restrictions placed only on persons known to have been exposed is  
62 considered insufficient to prevent further transmission [15, 17, 18]. Collective infection  
63 control measures can reduce the disease incidence, though at the price of a prolongation  
64 of the epidemic period [19]. Therefore, it is important to raise information on these  
65 measures.

66 A recent study that assessed the impact of online information on the individual-  
67 level intention to voluntarily self-isolate during the pandemic concluded that in order to  
68 enhance individuals' motivation to adopt preventive measures such as social distancing,  
69 actions should focus on raising consciousness on the severity of the situation, in addition,  
70 information overload had a significant impact on individuals' threat and coping  
71 perceptions, and through them on self-isolation intention [20]. Thus, the aim of the  
72 present investigation was to predict social distancing index (SDI) through Google and  
73 YouTube search trends and investigate the correlation between the SDI with  
74 epidemiological data on COVID-19 outbreak in Brazil.

## 75 **Methods**

76 In Brazil, the current market share of Google among the search engines is over 97%  
77 [21]. Google Trends (<https://trends.google.com/trends/>) data is a randomly collected

78 sample of Google search queries, each piece of data is categorized and tagged with a  
79 topic. Each data point is divided by the total searches in a specific location over a time  
80 period to compare relative popularity. Google Trends portrays search frequency output  
81 as a normalized data series and the resulting numbers are scaled on a range of 0 to 100  
82 based on a topic proportion to all searches on all topics. Scores represents search interest  
83 relative to the highest point on the graph for that time period and geographic location. A  
84 value of 100 is the peak popularity of a term. A value of 50 indicates that the term is half  
85 as popular as it was at its peak of popularity [22]. This same methodology was applied to  
86 YouTube search trends. To the present investigation the Brazilian Portuguese  
87 correspondent topic for “coronavirus”, which held most popularity, was used.

88 The Social Distancing Index was created to help combat the spread of COVID-19,  
89 since its launch it has been improved with the sole objective of providing an increasingly  
90 accurate data for public authorities and research institutes. In order to achieve the index,  
91 highly accurate geolocation data was treated with a distance algorithm. Polygons from all  
92 regions of the Brazilian Institute of Geography and Statistics were adopted in order to  
93 ensure a more accurate categorization and more reliable data [23]. Data is available on  
94 Inloco website, displayed as map and chart.

95 Epidemiological data concerning COVID-19 outbreak in Brazil were collected  
96 from the Brazilian government Health Ministry database, available online [24]. Statistical  
97 data on daily new cases, cumulative number of cases, cumulative number of deaths and  
98 recovered cases were retrieved.

99 All databases were assessed for data collection on 23 May 2020, and the  
100 information corresponds to the period from 23 February to 20 May 2020.

101 Data were submitted to statistical analysis, all tests were applied considering an  
102 error of 5% and the confidence interval of 95%, and the analyzes were carried out using  
103 SPSS software version 23.0 (SPSS Inc. Chicago, IL, USA). Although the hypothesis of  
104 normal distribution of data was not confirmed by the Kolmogorov-Smirnov test, the  
105 statistical analysis was performed by the application of nonparametric tests. The strength  
106 of the association between distinct measures was tested with Spearman rank correlation.  
107 Multiple linear regression was performed to verify whether Google and YouTube relative  
108 search volumes are predictors of the social distancing index in Brazil.

## 109 **Results**

110 The multiple linear regression analysis resulted in a statistically significant model  
111 [F (2,85) = 32,045; p<0.001; R<sup>2</sup>=0.430]. Therefore, Google RSV ( $\beta = 1.226$ ; t = 7.887;  
112 p<0.001) and YouTube RSV ( $\beta = -0.930$ ; t = -5.980; p<0.001) are predictors of the social  
113 distancing index in Brazil. The equation that describes this relationship is (SDI) = 34.347  
114 + 0.422 (Google RSV) + (-0.359) (YouTube RSV).

115 In Brazil for the time span analyzed the mean SDI score was approximately 43%,  
116 the maximum of social distancing observed during this period was 62.2%. In mean scores  
117 over 3312 new cases of COVID-19 were confirmed daily. In the moment of data  
118 collection, the reported total number of cumulative deaths, confirmed cases and recovered  
119 cases were 18859, 291579 and 116683, respectively. The mean values of search engines  
120 RSV are shown in table 1.

121 Correlation between SDI and the other studied variables was found to be varying  
122 from weak to moderate, statistically significant correlation was found with all measures  
123 tested except for cumulative recovered cases, as shown in table 2.

## 124 **Discussion**

125 Evidence suggests that collective isolation measures have been highly effective in  
126 controlling the spread of the COVID-19 [16, 25]. However, maintaining isolation for  
127 many months may have even worse consequences than an epidemic wave that runs an  
128 acute course, the isolation measures should be thoughtfully planned and executed based  
129 on current stage of pandemic [26]. As observed in table 1 the mean score for SDI was  
130 43.126 with a discrete standard deviation when compared with the standard deviation of  
131 the daily new cases. A weak correlation was found between the isolation measure with  
132 epidemiological data from COVID-19, this may represent that in Brazil social isolation  
133 measures are poorly associated with the course of the disease in the country. In addition,  
134 these findings may be correlated with failure to control the increase in cases that were  
135 added daily, on average, by 3312.26 new cases, during the period covered by this  
136 investigation. The absence of concise social distancing policies and, furthermore, the  
137 political instability at the center of the Brazilian government poses a deadly distraction in  
138 the middle of a public health emergency [1].

139 The weak correlation between social distancing and the course of the disease can  
140 also be observed in the time series pattern seen in figure 1 (section A), where while the  
141 number of cases per day and cumulative deaths show an ascending pattern, the social  
142 distancing index remains with slight changes, with the exception of a slight increase at  
143 the end of March, when cases start to show an ascending profile. At the end of the time  
144 series, when the daily new cases are at peak, SDI was under de mean observed in the  
145 evaluated period.

146 A correlation was found between search engines RSV and SDI (table 2). To further  
147 investigate this correlation a multiple linear regression was performed, data extracted  
148 from this analysis showed that Google RSV and YouTube RSV are predictors to social  
149 distancing. In figure 1 (section B) it is possible to observe de matching behavior of this

150 measures along the time series, with a similar pattern of peaks and decrease starting at the  
151 end of April.

152 The positive correlation found with Google RSV may be associated with the access  
153 to information, since raising awareness on the severity of the situation and the importance  
154 on following the advice from health organizations is a key point on achieving self-  
155 isolation intention [20].

156 The negative association observed between YouTube RSV and SDI through de  
157 multiple linear regression may correlate with the low quality of YouTube content on  
158 COVID-19 reported in previous studies [27-29], since misinformation can hinder the  
159 communication of health professionals and organizations with general public and even  
160 reduce compliance with treatments or medical advices [30, 31].

161 The findings of the present study support the evidence that online-available  
162 information can potentially assist conventional epidemiologic tools for estimating data  
163 about patterns of disease and population behavior [5, 6]. Relative search volume of  
164 Google and YouTube could define the proper timing and location for practicing  
165 appropriate risk communication strategies. Health authorities might apply these data to  
166 measure the effect of the transmission of information on the population and to obtain  
167 feedback from research statistics.

## 168 **Acknowledgments**

169 P.C., T.M., M.A. and E.L. were supported by a PhD scholarship from Coordenação de  
170 Aperfeiçoamento de Pessoal de Nível Superior (CAPES).

## 171 **Declaration of interest**

172 None



173 **Funding**

174 No funding was received for this work.

175 **References**

176 (1) **Anon.** COVID-19 in Brazil: "So what?"[Editorial]. *The Lancet* 2020; **395**: 1461.

177 (2) **Husnayain A, Fuad A, Chia-Yu ES.** Applications of Google Search Trends for Risk  
178 Communication in Infectious Disease Management: A Case Study of the COVID-19  
179 Outbreak in Taiwan. *International journal of infectious diseases : IJID : official*  
180 *publication of the International Society for Infectious Diseases* 2020; **95**: 221-223.

181 (3) **Barros JM, Duggan J, Rebholz-Schuhmann D.** The Application of Internet-Based  
182 Sources for Public Health Surveillance (Infoveillance): Systematic Review, *Journal of*  
183 *Medical Internet Research* 2020; **22**: e13680.

184 (4) **Bhattacharya S.** Predicting emerging and re-emerging disease outbreaks through  
185 internet search trends: An analysis from India. *AIMS Public Health* 2019; **6**: 1-3.

186 (5) **Ayyoubzadeh SM, et al.** Predicting COVID-19 Incidence Through Analysis of  
187 Google Trends Data in Iran: Data Mining and Deep Learning Pilot Study. *JMIR public*  
188 *health and surveillance* 2020; **6**: e18828.

189 (6) **Walker A, Hopkins C, Surda P.** The Use of Google Trends to Investigate the Loss  
190 of Smell Related Searches During COVID-19 Outbreak. *International forum of allergy*  
191 *& rhinology* 2020.

192 (7) **Cervellin G, Comelli I, Lippi G.** Is Google Trends a reliable tool for digital  
193 epidemiology? Insights from different clinical settings. *Journal of Epidemiology and*  
194 *Global Health* 2017; **7**: 185-189.

195 (8) **Effenberger M, et al.** Association of the COVID-19 pandemic with Internet Search  
196 Volumes: A Google Trends™ Analysis. *International Journal of Infectious Diseases*  
197 2020; **95**:192-197.

198 (9) **Ortiz-Martínez Y, et al.** Can Google® trends predict COVID-19 incidence and help  
199 preparedness? The situation in Colombia. *Travel Medicine and Infectious Disease* 2020.

- 200 (10) **Higgins TS, et al.** Correlations of Online Search Engine Trends With Coronavirus  
201 Disease (COVID-19) Incidence: Infodemiology Study. *JMIR Public Health and*  
202 *surveillance* 2020; **6**: e19702.
- 203 (11) **Morsy S, et al.** Prediction of Zika-confirmed Cases in Brazil and Colombia Using  
204 Google Trends. *Epidemiology and infection* 2018; **146**: 1625-1627.
- 205 (12) **Pei S, et al.** Forecasting the Spatial Transmission of Influenza in the United States.  
206 *Proceedings of the National Academy of Sciences of the United States of America* 2018;  
207 **115**: 2752-2757.
- 208 (13) **Verma M, et al.** Google Search Trends Predicting Disease Outbreaks: An Analysis  
209 from India. *Healthcare Informatics Research* 2018; **24**: 300-308.
- 210 (14) **Shin S-Y, et al.** High correlation of Middle East respiratory syndrome spread with  
211 Google search and Twitter trends in Korea. *Scientific Reports* 2016; **6**: 1-7.
- 212 (15) **Guner R, Hasanoglu I, Aktas F.** COVID-19: Prevention and control measures in  
213 community. *Turkish Journal of Medical Sciences* 2020; **50**: 571-577.
- 214 (16) **Taghrir MH, Akbarialiabad H, Marzaleh MA.** Efficacy of Mass Quarantine as  
215 Leverage of Health System Governance During COVID-19 Outbreak: A Mini Policy  
216 Review. *Archives of Iranian medicine* 2020; **23**: 265-267.
- 217 (17) **Wilder-Smith A, Chiew CJ, Lee VJ.** Can we contain the COVID-19 outbreak with  
218 the same measures as for SARS?. *The Lancet Infectious Diseases* 2020; **20**: e102-e107.
- 219 (18) **Wilder-Smith A, Freedman DO.** Isolation, quarantine, social distancing and  
220 community containment: pivotal role for old-style public health measures in the novel  
221 coronavirus (2019-nCoV) outbreak. *Journal of Travel Medicine* 2020; **27**: taaa020.
- 222 (19) **Raoult D, et al.** Coronavirus infections: Epidemiological, clinical and  
223 immunological features and hypotheses. *Cell Stress* 2020; **4**: 66-74.
- 224 (20) **Faroog A, Laato S, Islam AKMN.** Impact of Online Information on Self-Isolation  
225 Intention During the COVID-19 Pandemic: Cross-Sectional Study. *Journal of medical*  
226 *Internet research* 2020; **22**: e19128.

- 227 (21) **Statcounter Search Engine Market Share Brazil See**  
228 (<https://gs.statcounter.com/search-engine-market-share/all/brazil>). Accessed 20 April  
229 2020.
- 230 (22) **Mavragani A, et al.** Assessing the Methods, Tools, and Statistical Approaches in  
231 Google Trends Research: Systematic Review. *J Med Internet Res* 2018; **20**: e270.
- 232 (23) **Inloco Mapa brasileiro da COVID-19**  
233 (<https://mapabrasileirodacovid.inloco.com.br/pt/>). Accessed 23 May 2020.
- 234 (24) **Brazil Painel Coronavírus** (<https://covid.saude.gov.br/>). Accessed 23 May 2020.
- 235 (25) **Sjördin H, et al.** Only Strict Quarantine Measures Can Curb the Coronavirus Disease  
236 (COVID-19) Outbreak in Italy, 2020. *Eurosurveillance* 2020; **25**: 2000280.
- 237 (26) **Ioannidis, JPA.** Coronavirus disease 2019: The harms of exaggerated information  
238 and non-evidence-based measures. *European Journal of Clinical Investigation* 2020,  
239 e13223.
- 240 (27) **Basch CH, et al.** Preventive Behaviors Conveyed on YouTube to Mitigate  
241 Transmission of COVID-19: Cross-Sectional Study. *JMIR public health and surveillance*  
242 2020; **6**: e18807.
- 243 (28) **Basch CE, et al.** The Role of YouTube and the Entertainment Industry in Saving  
244 Lives by Educating and Mobilizing the Public to Adopt Behaviors for Community  
245 Mitigation of COVID-19: Successive Sampling Design Study. *JMIR public health and*  
246 *surveillance* 2020; **6**: e19145.
- 247 (29) **Li HOY, et al.** YouTube as a source of information on COVID-19: a pandemic of  
248 misinformation? *BMJ Global Health* 2020; **5**: e002604.
- 249 (30) **Lu X, et al.** Relationship Between Internet Health Information and Patient  
250 Compliance Based on Trust: Empirical Study. *Journal of medical Internet research* 2018;  
251 **20**: e253.
- 252 (31) **Lu X, Zhang R.** Impact of Physician-Patient Communication in Online Health  
253 Communities on Patient Compliance: Cross-Sectional Questionnaire Study. *Journal of*  
254 *medical Internet research* 2019; **21**: e12891.

255

256 **Table 1.** Frequency measures of SDI, search engines RSV and COVID-19 daily new  
 257 cases for the evaluated period

	Mean <sup>(SD)</sup>	Median	Minimum	Maximum
Social distancing index	43.126 <sup>(8.130)</sup>	43.30	26.60	62.20
Google RSV	37.470 <sup>(23.649)</sup>	34.50	2	100
YouTube RSV	19.520 <sup>(21.035)</sup>	11.50	2	100
Daily new cases	3312.26 <sup>(4564.24)</sup>	1180	0	19951

258 \*RSV – Relative Search Volume

259

260 **Table 2.** Correlation between SDI measures, COVID-19 epidemiological data and search  
 261 engines RSV

	Social distancing index	p-value
Google RSV	$\rho$ 0.567	<0.001
YouTube RSV	$\rho$ 0.367	<0.001
Daily new cases	$\rho$ 0.370	<0.001
Cumulative cases	$\rho$ 0.391	<0.001
Cumulative deaths	$\rho$ 0.396	<0.001
Recovered cases	$\rho$ 0.650	0.547

262 \*RSV – Relative Search Volume

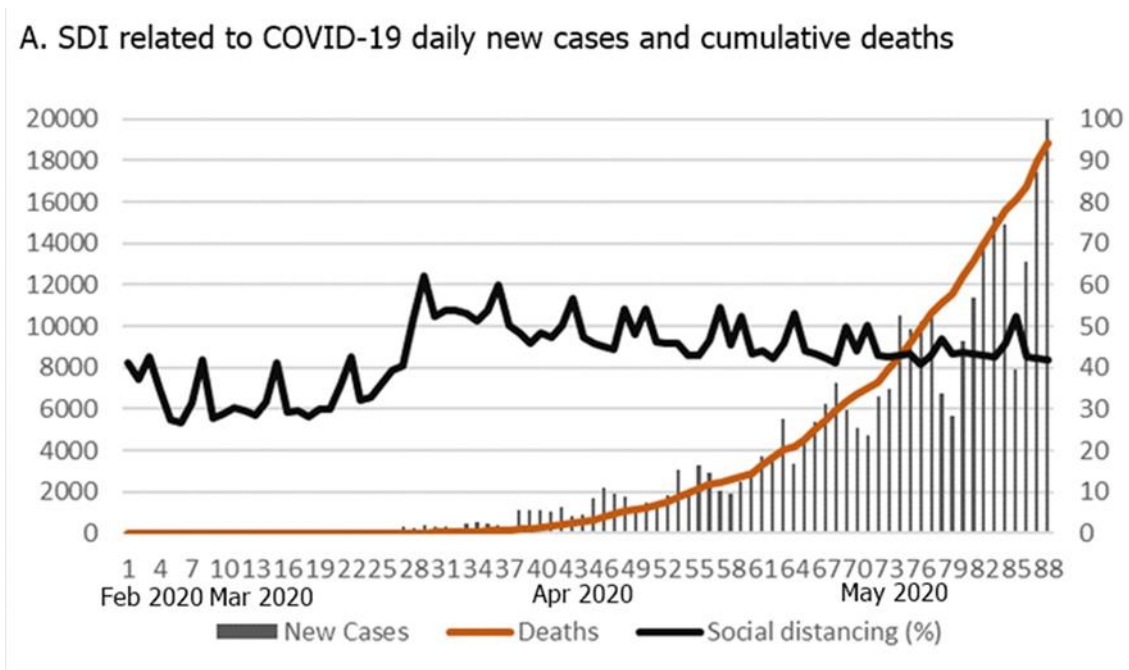
263 Mann-Whitney U test

264

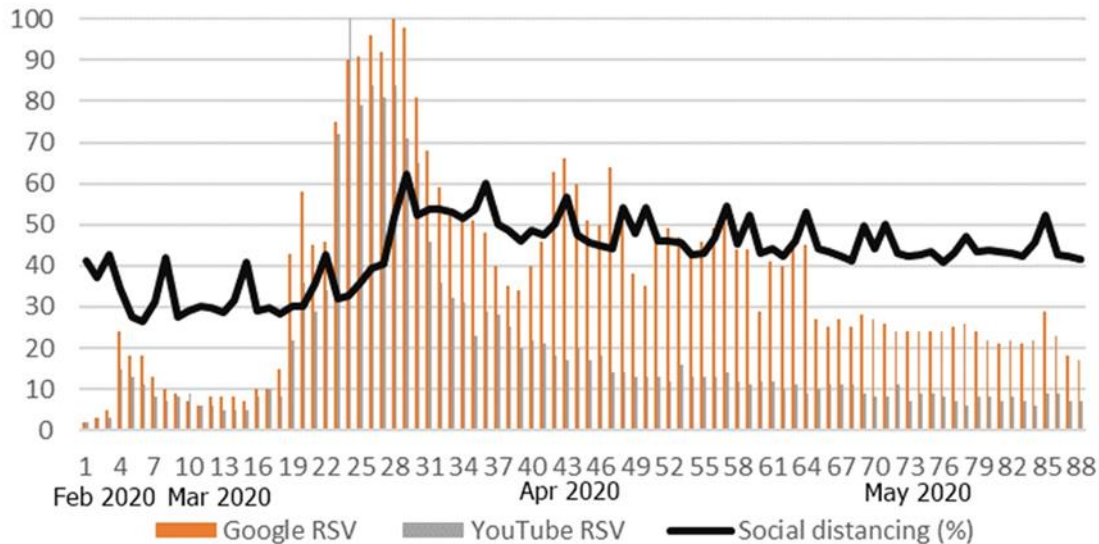
265

266

267 **Figure 1.** Time series of social distancing index related to Google and YouTube RSVs  
268 and COVID-19 new cases and cumulative deaths in Brazil.



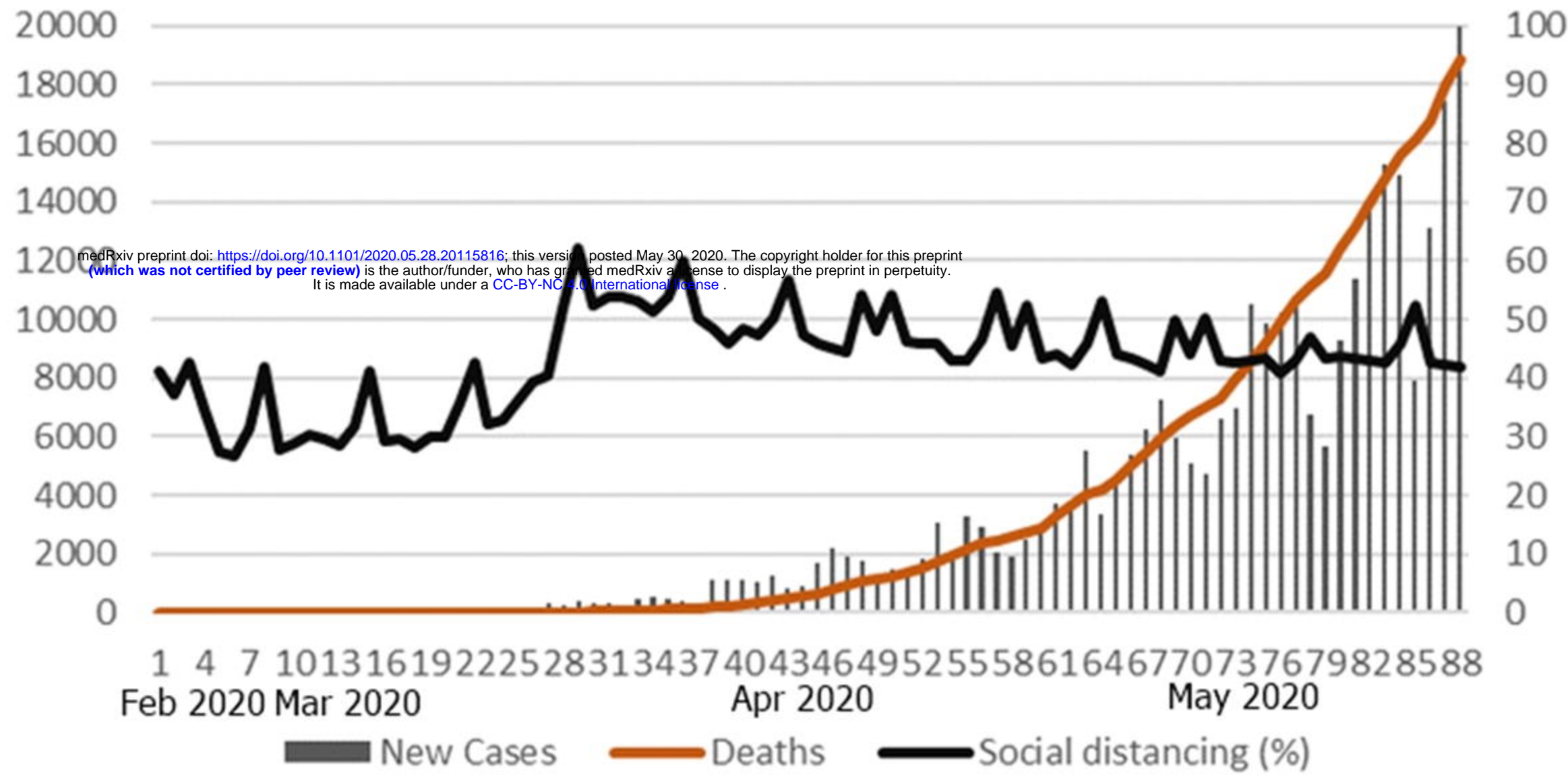
**B. SDI related to search engines RSVs**



269



### A. SDI related to COVID-19 daily new cases and cumulative deaths



### B. SDI related to search engines RSVs

