

Estimating the COVID-19 Prevalence in Spain with Indirect Reporting via Open Surveys

Augusto Garcia-Agundez^{1,*}, Oluwasegun Ojo², Harold Hernandez³, Carlos Baquero⁴, Davide Frey⁵, Chryssis Georgiou⁶, Mathieu Goessens⁷, Rosa Lillo³, Raquel Menezes⁸, Nicolas Nicolaou⁹, Antonio Ortega¹⁰, Efstathios Stavrakis⁹ and Antonio Fernandez Anta²

¹Multimedia Communications Lab, etit, TU Darmstadt, Darmstadt, Germany

²IMDEA Networks Institute, Madrid, Spain

³Department of Statistics, University Carlos III de Madrid, Madrid, Spain

⁴Departamento de Informatica, University of Minho, Braga, Portugal

⁵Inria Centre de Recherche Rennes Bretagne Atlantique, Rennes, France

⁶Department of Computer Science, University of Cyprus, Nicosia, Cyprus

⁷IMT Atlantique, Nantes, France

⁸Departamento de Matematica, University of Minho, Braga, Portugal

⁹Algolysis Ltd, Limassol, Cyprus

¹⁰Department of Electrical and Computer Engineering, USC Viterbi School of Engineering, Los Angeles, CA, USA

Correspondence*:

Augusto Garcia-Agundez

augusto.garcia@kom.tu-darmstadt.de

2 ABSTRACT

3 During the initial phases of the COVID-19 pandemic, accurate tracking has proven unfeasible.
4 Initial estimation methods pointed towards case numbers that were much higher than officially
5 reported. In the CoronaSurveys project, we have been addressing this issue using open online
6 surveys with indirect reporting. We compare our estimates with the results of a serology study for
7 Spain, obtaining high correlations (R squared 0.89). In our view, these results strongly support
8 the idea of using open surveys with indirect reporting as a method to broadly sense the progress
9 of a pandemic.

10 **Keywords:** COVID-19, pandemic, serology, survey, indirect reporting, sensing

1 INTRODUCTION

11 During the initial phases of the COVID-19 pandemic, progress tracking via massive serology testing has
12 proven to be unfeasible. However, initial estimation methods suggested that the real numbers of COVID-19
13 cases were significantly higher than those officially reported (1). For instance, by April 30th, 2020, the
14 number of confirmed fatalities due to COVID-19 in the US was 66,028, and the number of confirmed cases
15 was 1,080,303. However, with that number of fatalities the number of cases must have been no less than
16 4,784,637, by simply using the Case-fatality Ratio (CFR) of 1.38% measured in Wuhan (2).

17 In the case of Spain, the discrepancy seems to be even higher. Preliminary studies point towards only one
18 in 53 cases being reported during the first days of the pandemic (3). Although recent availability of massive

19 testing has reduced this discrepancy, demographic statistics still indicate a degree of underreporting to this
20 day, which can be seen among others in mortality numbers: all-cause mortality statistics in Spain point to
21 two periods of significant excess of deaths in the country over the predicted values in 2020: March and
22 April (44, 599 deaths in excess) and August to December (26, 186 deaths in excess) (4). These numbers
23 contrast with the officially reported number of deaths due to COVID-19, which rests at 50, 837 (5). This
24 discrepancy is corroborated in publications from official government authorities, which indicate an ongoing
25 estimated underreporting of 20% to 40% (6).

26 A potential method to address this limitation is to use online surveys during the initial stages of pandemics.
27 Online surveys can be deployed quickly and are cost-effective, but show potential weaknesses in sampling,
28 confidentiality, and other ethical issues (7). In spite of these weaknesses, online surveys have already been
29 successfully implemented in scenarios such as influenza tracking (8).

30 In the CoronaSurveys project, (9) we aim to track the progress of the COVID-19 pandemic using online,
31 open, anonymous surveys with indirect reporting. Other recent articles have also suggested the use of
32 surveys to monitor this pandemic, both for Spain (10, 11) and globally (12). However, to our knowledge,
33 all surveys conducted in Spain have employed direct reporting only, asking participants about themselves.
34 CoronaSurveys implements the network scale-up method of indirect reporting instead, allowing us to
35 collect data on a wide fraction of the population with a small number of responses and in a very short
36 time-frame (13). In this article, we compare the accuracy of CoronaSurveys with a gold standard: serology
37 testing data collected by the Spanish government in the ENE-COVID study (14).

2 METHODS

38 The survey deployed in the CoronaSurveys project can be answered via browser or mobile app. After the
39 participant indicates the region (Spanish autonomous community) for which information will be provided,
40 two additional questions are presented:

- 41 1. *How many people do you know in your area for which you know their health condition?* The answer to
42 this question by participant i is the *Reach* r_i .
- 43 2. *How many of those were diagnosed with or have symptoms of COVID-19?* The answer to this question
44 by participant i is the *Cumulative Number of Cases* c_i .

45 In the CoronaSurveys project we have focused on simplicity and brevity to maximize interest and retain
46 users that would consistently provide data every few days. For that reason the total number of questions
47 in the survey has been kept small at all times. Our approach yielded good initial results with about 200
48 responses per week. The survey has been promoted via social networks, direct contacts, and, more recently,
49 with paid advertising.

50 To ensure total anonymity, the surveys are hosted on a private instance of LimeSurvey (15). Data is
51 aggregated daily, and in this process the responses are shuffled so no single entry can be back-traced to its
52 user. All the data is published in a public Github repository. The study design was reviewed and approved
53 by the ethics committee of the IMDEA Networks Institute. The survey includes an informed consent.

54 Once the data is collected, we remove outlier responses. A response is considered an outlier if (1) r_i is
55 outside 1.5 times the interquartile range above the upper quartile (which for the data in this paper means
56 $r_i > 175$) or if (2) c_i/r_i is greater than $1/3$ (to exclude participants with an exceptionally high contact
57 with cases). Although participants may choose to provide information for the whole country, in this paper
58 we only consider responses in which participants provide information for their specific region. Hence,

59 the data is aggregated by region for all participants, to obtain the estimator of COVID-19 prevalence
60 $(\sum_i c_i)/(\sum_i r_i)$ (13).

3 RESULTS

61 To assess the accuracy of this method in estimating the cumulative number of cases of COVID-19, we
62 compare our cross-sectional survey estimates with the results of the serology study of Pollán et al. (14) for
63 Spain. We exclude Ceuta and Melilla due to lack of data on our part. Conducted between April 27 and
64 May 11, 2020, the serology study provides data for $n = 61,075$ participants ($0.1787\% \pm 0.0984\%$ of the
65 regional population, and 0.1299% of the national population). We consider as positive cases those that
66 tested positive to the point-of-care or immunoassay IgG tests (Supplementary Table 6 in Pollán et al. (14),
67 column *Either test positive*).

68 For our estimates, we consider the (up to) 100 most recent survey responses per region on April 20. The
69 date is chosen because the mean period between illness onset and a 95% confidence of IgG antibodies
70 presence is 14 days (16). This results in $n = 999$ responses (59 ± 35 per region) across Spanish regions,
71 with a cumulative reach of $\sum_i r_i = 67,199$ ($0.1827\% \pm 0.0701\%$ of the regional population, and 0.1434%
72 of the national population). On average, participants provide information for $r_i = 74.6219 \pm 38.0291$
73 members in their social circle, which is coherent with Dunbar's acquaintance group and related studies that
74 take social networks into consideration (17). Within this dataset, our outlier removal methods excluded
75 $19.8883\% \pm 9.2692\%$ of responses, including spurious contributions as the original average reach per
76 participant before filtering was greater than $5 \cdot 10^{15}$.

77 The Bland-Altman plot in Figure 1B shows a high correlation between the CoronaSurveys estimates and
78 the gold standard. A direct comparison of crude percentages, depicted in Figure 1B, also yields excellent
79 results ($R^2 = 0.8994$). Table 1 presents a detailed comparison of the estimates per region obtained in the
80 different studies.

81 Figure 2A presents how the number of responses per region affects the resulting value of R^2 . This
82 analysis indicates that 50 responses per region can already offer a reasonable estimation of cases. Including
83 more responses may further increase accuracy, but the numbers remain reasonably stable. Naturally, it is
84 important that responses are well distributed across all regions. Figure 2B depicts the effect of the day
85 limit on R^2 if we consider a date of \pm one week. Theoretically, a bell curve centered on the 20th should be
86 expected, as estimating too early would imply too few cases are reported, and estimating too late would
87 include more cases. We indeed observe an impact on accuracy, and the left half of the bell curve is more
88 visible. The change in accuracy is mostly due to new daily responses collected on April 16th. The lack of
89 the right half of the bell curve is due to the low number of new daily responses after April 16th, which
90 implies that the daily estimates are computed with sets of responses with large intersections. Interestingly,
91 a similarly high number of responses was collected on April 14th, with nearly no impact on accuracy.

4 DISCUSSION

92 The linear regression equation in Figure 1A points to CoronaSurveys very consistently underestimating
93 the number of cases by a factor of approximately 46%, possibly due to asymptomatic cases. This ratio
94 is consistent with the estimates of the Covid19Impact study of Oliver et al. (11), which used more than
95 140,000 direct survey responses collected on March 28th-30th. It is also consistent with the reported data
96 on asymptomatic cases reported by Pollán et al. (14), which found that around a third of the seropositive
97 participants were asymptomatic (see Table 1).

98 Concerning the impact of the number of responses as depicted in Figure 2, we observe how once the
99 minimum number is reached, further responses will not significantly increase accuracy unless these come
100 from underreported regions. As depicted in Figure 3, additional responses from regions where many are
101 already available will barely have an impact on the global result. As the great majority of contributions for
102 April 14th were for Madrid, where we already had many responses available, the 77 new daily responses on
103 April 14th barely had any impact, while the contributions on April 16th significantly increase the accuracy
104 of our estimation.

105 Our study presents a number of limitations. Firstly, as presented in Table 1, our number of responses
106 in some regions was limited (e.g., 9 responses in La Rioja or 16 in Navarra and Cantabria). Our own
107 analysis suggests this is not enough to offer reliable data for these three regions. Additionally, our criteria
108 to eliminate outliers is heuristic, and may change in the future as we collect more data.

109 Nevertheless, despite these limitations, the estimates obtained in CoronaSurveys show high correlation
110 with serology tests. Moreover, since the underestimation of our method over all regions is homogeneous,
111 and consistent with the one third fraction of asymptomatic reported by Pollán et al. (14), these estimates
112 can be “corrected” to provide an accurate cumulative number of cases for each region. We will further
113 evaluate the robustness of our model as Pollán et al. publish the results of their three additional serology
114 studies.

115 In summary, we believe these results strongly support using open surveys with indirect reporting as a
116 method to broadly sense the progress of a pandemic.

CONFLICT OF INTEREST STATEMENT

117 The authors declare that the research was conducted in the absence of any commercial or financial
118 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

119 The analysis presented in this article was conducted by Augusto Garcia-Agundez and Antonio Fernandez
120 Anta with support and feedback from all remaining co-authors. The data acquisition and processing
121 techniques were developed by all co-authors.

FUNDING

122 At the time of writing this article, CoronaSurveys has received no public funding. Social networks surveys
123 have been partially funded via donations through our website. CoronaSurveys received an award from the
124 UMD/CMU COVID-19 Symptom Data Challenge.

ACKNOWLEDGMENTS

125 We would like to thank all CoronaSurveys researchers and collaborators for their contribution to this
126 project: <https://coronasurveys.org/team/>.

DATA AVAILABILITY STATEMENT

127 The datasets generated and analyzed for this study can be found in the CoronaSurveys Github Repository
128 at <https://github.com/GCGImdea/coronasurveys>.

REFERENCES

129 1 .Maxmen A. How much is coronavirus spreading under the radar. *Nature* **10** (2020).

- 130 2 .Verity R, Okell LC, Dorigatti I, Winskill P, Whittaker C, Imai N, et al. Estimates of the severity of
131 coronavirus disease 2019: a model-based analysis. *The Lancet infectious diseases* **20** (2020) 669–677.
- 132 3 .Krantz SG, Rao ASS. Level of underreporting including underdiagnosis before the first peak of
133 COVID-19 in various countries: Preliminary retrospective results based on wavelets and deterministic
134 modeling. *Infection Control & Hospital Epidemiology* (2020) 1–3.
- 135 4 .[Dataset] Centro Nacional de Epidemiología, Instituto de Salud
136 Carlos III. Informe MoMo. situación a 30 de diciembre de 2020.
137 <https://www.isciii.es/QueHacemos/Servicios/VigilanciaSaludPublicaRENAVE/Enfermedades>
138 [Transmisibles/MoMo/Paginas/Informes-MoMo-2020.aspx](https://www.isciii.es/QueHacemos/Servicios/VigilanciaSaludPublicaRENAVE/Enfermedades) (2020).
- 139 5 .[Dataset] Ministerio de Sanidad Gobierno de España. Actualización nº 282. enfermedad por el
140 coronavirus (COVID-19). [https://www.msbs.gob.es/profesionales/saludPublica/ccayes/alertasActual/nCov/documentos/](https://www.msbs.gob.es/profesionales/saludPublica/ccayes/alertasActual/nCov/documentos/Actualizacion_282_COVID19.pdf)
141 [Actualizacion_282_COVID19.pdf](https://www.msbs.gob.es/profesionales/saludPublica/ccayes/alertasActual/nCov/documentos/Actualizacion_282_COVID19.pdf) (2020).
- 142 6 .Moros MJS, Monge S, Rodríguez BS, San Miguel LG, Soria FS. COVID-19 in Spain: view from the
143 eye of the storm. *The Lancet Public Health* (2020).
- 144 7 .Nayak M, Narayan K. Strengths and weakness of online surveys. *IOSR Journal of Humanities and*
145 *Social Science* **24** (2019) 31–38.
- 146 8 .Dalton C, Durrheim D, Fejsa J, Francis L, Carlson S, d’Espaignet ET, et al. Flutracking: a weekly
147 australian community online survey of influenza-like illness in 2006, 2007 and 2008. *Communicable*
148 *diseases intelligence quarterly report* **33** (2009) 316–322.
- 149 9 .Ojo O, García-Agundez A, Girault B, Hernández H, Cabana E, García-García A, et al. Coronasurveys:
150 Using surveys with indirect reporting to estimate the incidence and evolution of epidemics.
151 *KDD Workshop Humanitarian Mapping, San Diego, California USA, August 24, 2020. ArXiv*
152 *preprint:2005.12783* (2020).
- 153 10 .Linares M, Garitano I, Santos L, Ramos JM. Estimando el número de casos de COVID-19 a tiempo real
154 utilizando un formulario web a través de las redes sociales: Proyecto COVID19-TRENDS. *Semergen*
155 (2020).
- 156 11 .Oliver N, Barber X, Roomp K, Roomp K. Assessing the impact of the COVID-19 pandemic in Spain:
157 Large-scale, online, self-reported population survey. *Journal of medical Internet research* **22** (2020)
158 e21319.
- 159 12 .[Dataset] Facebook Data for Good. COVID-19 symptom survey –
160 request for data access. [https://dataforgood.fb.com/docs/](https://dataforgood.fb.com/docs/covid-19-symptom-survey-request-for-data-access/)
161 [covid-19-symptom-survey-request-for-data-access/](https://dataforgood.fb.com/docs/covid-19-symptom-survey-request-for-data-access/) (2020). Accessed:
162 2021-01-24.
- 163 13 .Bernard HR, Hallett T, Iovita A, Johnsen EC, Lyerla R, McCarty C, et al. Counting hard-to-count
164 populations: the network scale-up method for public health. *Sex. Transm. Infect.* **86** (2010) ii11–ii15.
- 165 14 .Pollán M, Pérez-Gómez B, Pastor-Barriuso R, Oteo J, Hernán MA, Pérez-Olmeda M, et al. Prevalence
166 of SARS-CoV-2 in Spain (ENE-COVID): a nationwide, population-based seroepidemiological study.
167 *Lancet* **396** (2020) 535–544.
- 168 15 .LimeSurvey Project Team / Carsten Schmitz. *LimeSurvey: An Open Source survey tool*. LimeSurvey
169 Project, Hamburg, Germany (2012).
- 170 16 .Pallett SJ, Rayment M, Patel A, Fitzgerald-Smith SA, Denny SJ, Charani E, et al. Point-of-care
171 serological assays for delayed SARS-CoV-2 case identification among health-care workers in the UK:
172 a prospective multicentre cohort study. *Lancet Respir. Med.* **8** (2020) 885–894.
- 173 17 .Gonçalves B, Perra N, Vespignani A. Modeling users’ activity on twitter networks: Validation of
174 dunbar’s number. *PloS one* **6** (2011) e22656.

FIGURE CAPTIONS

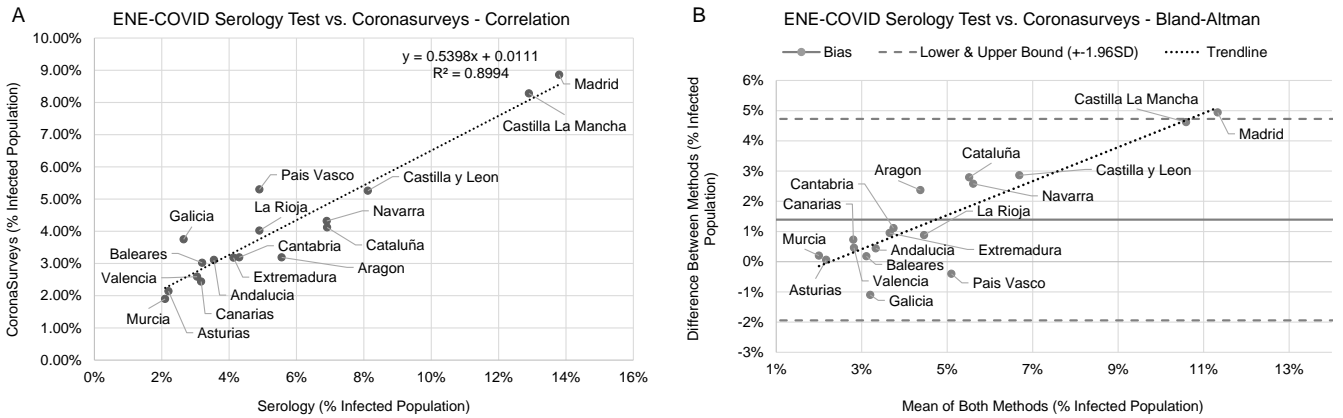


Figure 1. Comparison between the serology test and CoronaSurveys, Bland-Altman (A) and direct correlation (B)

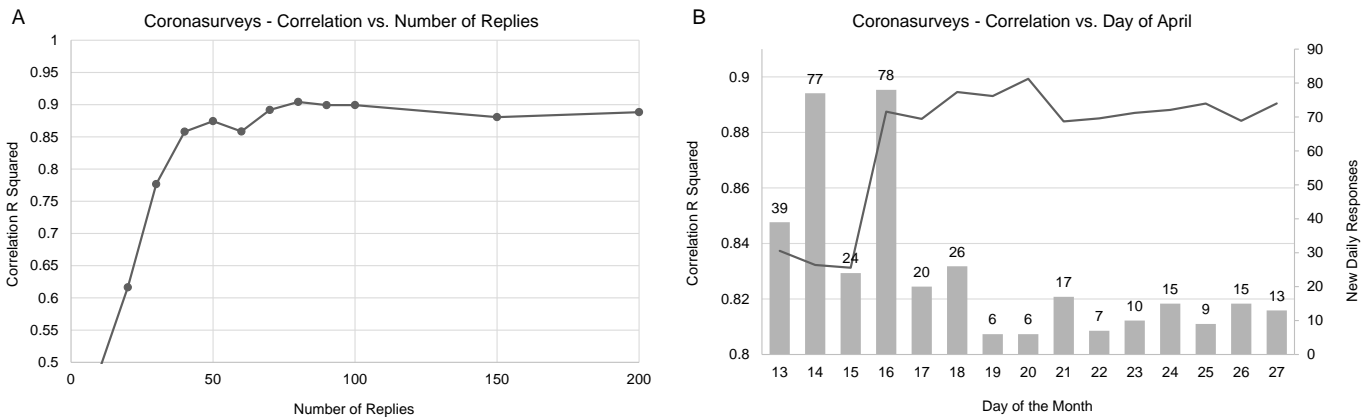


Figure 2. Convergence of correlation with number of responses (A) and day of the month (B). The line represents the resulting R squared correlation, the dots in the line represent sampling points. The bars represent the number of new daily responses

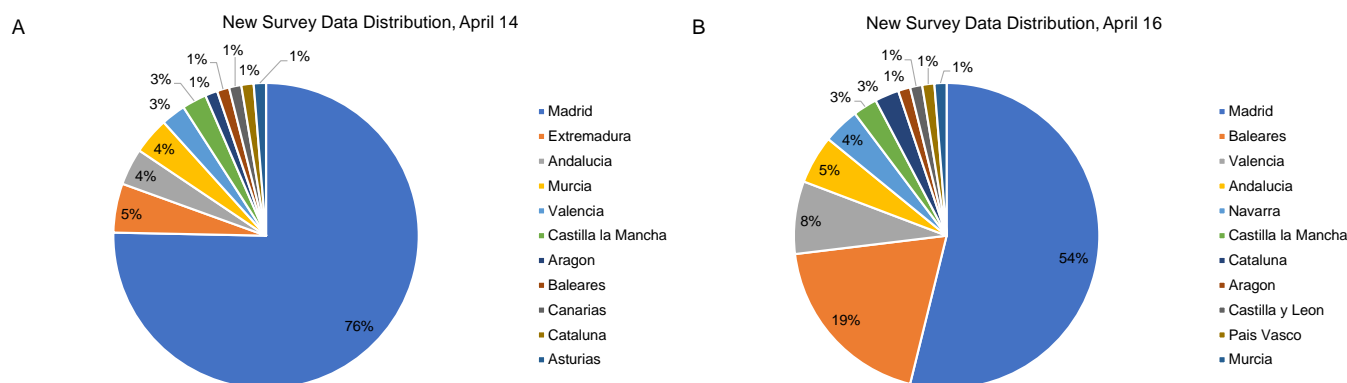


Figure 3. Distribution of new survey responses on April 14 (A) and April 16 (B)

Region	ENE-COVID (% Infected)	CoronaSurveys			Covid19Impact	
		(% Infected)	Responses	Reach	(% Infected)	Responses
Andalucia	3.55	3.11(±0.41)	100	6, 721	2.2(±0.3)	5, 691
Aragon	5.56	3.19(±0.41)	44	3, 045	2.0(±0.3)	1, 463
Asturias	2.20	2.14(±0.52)	42	2, 987	1.5(±0.3)	655
Cantabria	4.30	3.19(±0.96)	16	1, 285	2.8(±0.3)	497
Castilla y Leon	8.12	5.26(±0.58)	86	5, 763	3.7(±0.4)	1, 994
Castilla La Mancha	12.90	8.28(±0.68)	100	6, 399	8.0(±0.3)	3, 469
Canarias	3.17	2.44(±0.74)	26	1, 678	1.4(±0.2)	1, 052
Catalonia	6.91	4.12(±0.49)	100	6, 310	2.8(±0.3)	5, 088
Extremadura	4.13	3.18(±0.74)	32	2, 168	2.3(±0.4)	656
Galicia	2.65	3.75(±0.49)	85	5, 781	1.3(±0.3)	2, 257
Baleares	3.20	3.02(±0.76)	33	1, 955	1.9(±0.3)	1, 222
Murcia	2.10	1.90(±0.50)	45	2, 835	1.5(±0.3)	3, 566
Madrid	13.8	8.86(±0.67)	100	6, 850	6.1(±0.4)	10, 365
Navarra	6.90	4.32(±1.16)	16	1, 180	3.6(±0.4)	580
Basque Country	4.90	5.30(±0.65)	65	4, 511	1.9(±0.4)	1, 007
La Rioja	4.90	4.02(±1.72)	9	498	1.8(±0.4)	220
Valencia	3.05	2.59(±0.37)	100	7, 233	1.6(±0.3)	102, 021

Table 1. Percentage (and 95% confidence interval) of infected population per region according to the ENE-COVID serology study (14), CoronaSurveys and Covid19Impact (11) (symptom-only model).