

# Accurate Covid-19 prevalence measurement in the field

Sotiris Georganas,<sup>a</sup> Alina Velias,<sup>a</sup> Sotiris Vandalos<sup>b,c,\*</sup>

<sup>a</sup> City, University of London, London, UK

<sup>b</sup> King's College London, London, UK

<sup>c</sup> Harvard T.H. Chan School of Public Health, Boston, USA

\* Correspondence to: King's College London, Bush House, 30 Aldwych, London WC2B 4BG, United Kingdom. Email: [s.vandalos@kcl.ac.uk](mailto:s.vandalos@kcl.ac.uk)

## Abstract:

Accurate epidemic prevalence measurement is a necessary condition for informed policy decision-making. In the Covid-19 pandemic especially, wrong prevalence measurement can lead to tremendous waste, be that in life years or economic output. A number of countries offer random Covid-19 tests to estimate the prevalence of the virus in the population, and report daily positivity rates. However, since virus testing *has* to be voluntary, all tests done in the field, even if supposedly random, suffer from selection bias. This bias, unlike standard biases in polling, is not limited to having a representative sample, and thus cannot be corrected by the usual methods (quota sampling etc). The issue is that people who *feel* they have symptoms (or other reasons to suspect they are carrying the virus), are up to 38 times more likely to volunteer to get tested, and testing stations cannot readily correct this by oversampling (i.e. selecting people without symptoms to test). Using controlled, incentivized online experiments with over 500 subjects of all ages in a European country, we show that this difference in testing propensities leads to sizeable bias; “random” tests in the field inflate infection figures by up to five times. We suggest ways to correct the bias of the testing stations, but even better, a cleaner way to sample the population to avoid the bias altogether. Our methodology is relevant for covid-19, but also any other epidemic where carriers can have informative beliefs about their own carrier status.

*Keywords:* Covid-19; testing; bias; prevalence measurement; infection rates; experiment

## 1. Background

Covid-19 has already caused over 1.7 million confirmed deaths globally (Johns Hopkins 2020). Apart from deaths directly attributed to Covid-19, there are reports of excess mortality over and above officially reported Covid-19 deaths,<sup>1</sup> as well as indirect health effects and deaths. These include those suffering from other diseases not receiving medical attention or being undiagnosed<sup>2-3</sup> and suicides.<sup>4-5</sup> The pandemic has crippled economic activity, leading to increasing unemployment rates and shrinking national income worldwide<sup>6</sup> – which in turn can lead to further deterioration of health.<sup>7-8</sup>

Tackling the pandemic is of paramount importance for obvious health and financial reasons. Any suggested policy responses and their implementation (such as social distancing rules) will inevitably be inefficient if we are not aware of the real number of active cases, and in which areas and age groups these occur. Observing mortality rates or the number of hospitalisations and patients in ICU provides an estimate of how many people caught Covid-19 weeks earlier (although estimating the fatality rate is also challenging).<sup>9</sup> It is important to know the number actual cases at present, which, apart from designing policy responses, also provides an estimate of hospitalisations and mortality in the next weeks.

Community testing, often conducted in the high street and in neighbourhoods, is widely considered a useful tool to monitor incidence and trends (e.g. the ECDC<sup>10</sup> lists “[to] reliably monitor SARS-CoV-2 transmission rates and severity” among five objectives of testing). It also publishes weekly testing data and “positivity rates” by EU State.<sup>11</sup> However, we show that such testing is highly unlikely to provide accurate estimates of Covid-19 prevalence, and the main problem is not related the typical issues that arise in population sampling, such as age group structure. Testing has to be voluntary, and we expect that people are more likely self-select into testing if they have reasons to believe they have a good chance of having Covid-19 (such as, e.g. if they have symptoms or if they are exposed to a high-risk environment). This

*self-selection bias* is likely to increase with waiting times and any other cost associated with testing. Furthermore, the bias is expected to be time-varying, because it depends non-linearly on time varying parameters. For example, when cases rise steeply, people might be more likely to want to test out of fear – but at the same time this will also affect their behaviour and thus the likelihood of catching Covid-19.

The objective of this paper is to examine whether and to what extent bias occurs in Covid-19 testing; and to offer a debiasing solution to accurately estimate Covid-19 prevalence in the field. To measure the bias, we employ incentivised controlled experiments with more than 500 subjects in Greece, using standard experimental methods (such methods have been used, for example, to estimate the demand for HIV testing in an influential paper by R. Thornton.<sup>12</sup>

## **2. Data and Methods**

Data collection took place over a week, from 11 till 18 December 2020. The majority of the responses were collected online, via the *QualtricsTX* platform. To enable greater representativeness of the sample, 94 responses (16%) from elder people (median age = 63) were collected using phone interviews. Out of 608 participants starting the online study, 24 (4.7) dropped out mostly after the first few questions, resulting in the final sample of 578 observations.

Median age for the sample was 39 years (median for Greece 45.6), and the age distribution is shown in Figure 1.

[Insert Figure 1 here]

Subjects were invited to participate in a study on Covid-19 and related behaviours. Upon signing a consent form, the participant was first asked about general and Covid-19-related health. We then elicited hypothetical willingness to wait (WTW) to take a rapid test for Covid-19, *conditional* on (i) feeling healthy, (ii) having flu-like symptoms, (iii) having Covid-19 like symptoms. For all three hypothetical scenarios, the test was being offered by the national health authority (EODY) while the participant was walking down the street. This was done to reduce the (hypothetical) travel costs and reliability-related concerns.

After eliciting the hypothetical WTW, we asked the subjects several control questions, including exposure to Covid-19 risky environments (e.g. taking public transport or working face-to-face with many people) and socio-demographics. After completing the compulsory part of the study, the participants were offered an optional task for which they were randomly allocated to one of the two prize treatments. In treatment *Book*, the participant would enter a 1/30 chance lottery for a voucher for the local large-scale bookshop chain (“Public”), worth €80. In treatment *Test*, the participant would enter the same 1/30 chance lottery for a voucher for a home-administered Covid-19 test. For both prizes, the delivery was guaranteed within next 36 hours. All 578 participants completed the hypothetical elicitation and the control questions (left part of Figure 2).

[Insert Figure 2 here]

As was partly expected, a substantial part of the sample (n=174) did not continue to the optional task. A major part of it (n=78) was the elder people subsample. We are not very concerned that the inconvenience of the waiting task over the phone was the issue, since the participants came from the sample that previously participated in a study involving a real effort task over the phone (Georganas, Laliotis and Velias 2020, working paper). For n=38

participants, a software glitch in Qualtrics, in the first five hours of the study resulted in missing recording of the treatment allocation, so we had to drop their data despite completion of the optional task.

The participants then read the description of the optional task. They learned that it involved waiting in front of their screen for some time (target) that would be revealed in the next screen, and the lottery draw for the prize would take place right after the wait. They also learned that to ensure that they are waiting, a button would appear at random times and they would need to press it within 4 seconds to avoid being disqualified. Among the 303 participants who read the description of the optional task, 241 continued to the next screen which revealed the waiting target. At this stage, they were randomly allocated to one of the four waiting target conditions {300, 600, 900, 1200}. Upon learning the wait time, further 59 participants dropped out instantly (median target time 900 seconds). Among the 241 waiting, 69 dropped out before completing the full wait (median target time 900 seconds). In total, 172 participants completed the waiting target (median target 600 seconds).

Upon completing the waiting task, each participant was randomly allocated to one of the four *Cash* conditions, {€20, €35, €50, €65}. The participant was offered a choice to enter the lottery for: (a) the original prize (*Book, Test*), or (b) the displayed *Cash* amount. Out of the 172 participants, 112 chose to swap the original prize for the cash amount, whilst 60 chose to stay with the original prize (median cash value €35 for both). A total of 7 participants won the lottery.

### **3. Results**

Table 1 shows the ratio of willingness to test between people with symptoms and those without. The figure ranges between 1.5 and 38, depending on the age group and waiting times. People under 30 with symptoms are 1.532 times more likely to test when there is no waiting

time, compared to those without symptoms. This figure increases to 2.882 when there is a short wait of 5-15 minutes; 4.423 with a 15-30 minute wait; 15.5 with a 30-60 minute wait and 38 with a 1-2 hour wait. The ratio for 30-50 year-olds ranges between 1.517 for no wait and 16 for a 1-2 hour wait. For over 50-year-olds, the ratio ranges between 1.708 and 11.333. Overall, there is a bias even for no waiting time at all, which increases steeply for long waiting times in all age groups. Note that the bias also varies by other observable characteristics, for example, for waiting times of two hours and more, it is 84% higher for men than for women. Also, the propensity bias is 50% higher for obese people than for people in the healthy range, which indicates that people at risk not only have higher propensity to test (as is to be expected) but also react stronger to symptoms.

[Insert Table 1 here]

The *propensity to test bias* translates to a biased virus *prevalence estimate* ( $\beta$ ) which is also time varying. Crucially it depends on symptom prevalence, which, given the exponential spread of Covid-19, can change massively in a short period of time. This means that the estimate depends on symptom prevalence, but the bias itself also depends on it – so the bias is time variant.

Apart from waiting times, self-selecting into testing also depends on the cost associated with it (if applicable – costs can vary from time to monetary value, travel etc). We found that the bias is associated with willingness to pay for the test (Table A2 in Online Appendix A). Of those who won a test voucher, 83.8% swapped it for cash, as opposed to 48.9% of those who won the book voucher, indicating that the majority of subjects would not be willing to pay to receive a test. However, the scope of this article is to correct bias for free tests subject to

different waiting times, and further experiments are needed to reach concrete conclusions on willingness to pay.

We have launched an online calculator that provides estimates on the testing bias (available at <http://georgana.net/sotiris/task/atten/covid.php> ). The bias calculations that lead to the formula on which the calculator is based is provided in Online Appendix B. The estimates on the testing bias depend on (a) the percentage of tests yielding positive results; (b) the percentage of the general population that reports symptoms; (c) the relative likelihood of having Covid-19 for those with symptoms compared to those without symptoms; and (d) how more likely are people with symptoms to self-select into testing than those without symptoms. According to our methodology, it is possible to calculate these figures and thus estimate the bias. (a) is provided by the results of community testing; (b) is provided by surveying; (c) can be obtained by asking people a simple question before testing them for Covid-19; and (d) is provided by surveying.

A simple example is the following: Assume community testing led to 10% positive results, and 5% of the population reported symptoms. Without waiting time, if those with symptoms are 5 times more likely to have the virus than those without symptoms, then the results of community testing exaggerate by 27.71%, and the true prevalence in the population is 7.83% (instead of the reported 10%). At a 30-60 minute waiting time, the bias increases to 106.95%, meaning that the true prevalence in the population is 4.83%.

To further illustrate our results, Figure 1 depicts our best estimate of the virus prevalence bias, i.e. the ratio between reported prevalence and actual, depending on symptoms prevalence and waiting time, for the three age groups.

Based on these estimates, we can simulate how different demographic structures would affect the prevalence bias. In the following graph we depict the results from 3 million draws from the plausible parameter space (we assume symptoms prevalence of 5%, and allow the

testing bias parameter to vary uniformly within the 95% confidence interval gained from the experiments in Greece) applied to three countries, with different demographic structures: Nigeria (with one of the youngest populations globally), Italy (heavily ageing population) and the USA (between the two extremes).

The simulation shows that demography matters: a young country like Nigeria could have a substantially higher prevalence bias than Italy. However, it is also clear that the waiting times are more important. Lowering waiting times would result in a low bias for all countries.

#### **4. Discussion**

Using a survey-based experiment, we found that the probability of taking a Covid-19 test for those who have symptoms (or believe they are more likely to have caught the virus) is many times higher than those who do not. Our results show that people who *feel* they have symptoms (or other reasons to suspect they are carrying the virus), are up to 38 times more likely to volunteer to get tested. In our sample, this testing propensity bias ranged from 1.5 times (for people under 30 years with no waiting time) to 38 times (for people under 30 and a 2-hour waiting time). The bias becomes larger with longer waiting times, and any cost associated with taking the test. Testing stations cannot readily correct this by oversampling (i.e. selecting people without symptoms to test).

As demonstrated in our results, demographics influence the testing propensity bias, which means that different areas (or countries) will have different biases depending on the age composition. Furthermore, there have been reports of very long waiting times in community testing, which further exacerbate the bias. It is important to note that the bias is time variant, and is likely to depend on the actual virus prevalence.

Our findings suggest that results from community testing sites are heavily biased, and the bias goes beyond the usual issues of age groups etc. Rather, it relates to self-selection into



testing for those who are more likely to have Covid-19. This makes the aggregate results of community testing unreliable, when it comes to drawing conclusions on the prevalence of Covid-19 in the population.

We recognise the importance of giving people the opportunity to test, as this identifies positive cases, thus allowing them to self-isolate and stop spreading the disease. If the goal of testing sites is to allow random people to have a quick and free test, then this possibly meets its goal. Note, however, that random testing is not the economically or epidemiologically most efficient solution: subsidising tests specifically for populations with a high risk of getting infected and infecting others would probably save more lives at lower costs (say, tests for young people working in service industries and living with their parents). Such questions remain open for future research. However, we have shown that “random” voluntary testing is not really random. As such, it does not provide accurate information on disease prevalence, which is important to design and implement urgent policy responses to the pandemic, in terms of type, intensity and geographic area. Since voluntary testing is always biased, aggregate results on prevalence should be corrected. Debiasing can be performed using our methodology, as long as there are good estimates for some parameters, namely (a) the percentage of tests yielding positive results; (b) the percentage of the general population that reports symptoms; (c) the relative likelihood of having Covid-19 for those with symptoms compared to those without symptoms; and (d) how more likely are people with symptoms to self-select into testing than those without symptoms. If the probability of having covid given that one has symptoms is known (e.g. from asking a simple question at testing sites) then all that is needed is regularly polling a representative sample to get symptoms prevalence, which is much simpler and cost effective than random street testing.

Our methodology is not limited to correcting the results of community testing. The confirmed cases reported daily is also biased, as some people might not test because of costs,

or the inconvenience of going to a testing site, or even due to being afraid of losing income. According to our results, even at no monetary cost and no wait, 3.98% of people with symptoms would not get tested – which increases to 9.51% even for the slightest waiting time, rising even further when tests have a non-negligible cost to the citizen. Using polling results from a representative sample can correct this error.

Mass testing, extending to a very large part of the population is extremely useful as it can provide more accurate figures, and also identifies positive cases. It has been used, among others, in Liverpool, Slovakia and South Korea.<sup>13-15</sup> However, mass testing is very expensive, and might not be possible, especially at frequent intervals, due to capacity or other technical reasons. In those cases, debiasing the estimates is of paramount importance for health and the economy. Underestimating disease prevalence can trigger inadequate measures and further spread of disease, while overestimating can be detrimental to economic activity. We thus urge policy makers to redesign “random” testing as a matter of priority in the effort to tackle the pandemic.

As a final note, our methodology can be applied to the prevalence measurement of any epidemic, when carriers have informative private information about their health status. Fighting disease is hard, even without the added complication of not knowing the location and magnitude of the fight. Our work offers tools to be able to measure prevalence in real time. Further work is needed though, to estimate specific selection-bias parameters for every disease, as they are necessarily related to the health burden and life expectancy reduction caused by the specific pathogen.

## References

- [1] Vandoros, S., 2020. Excess Mortality during the Covid-19 pandemic: Early evidence from England and Wales. *Social Science & Medicine*, p.113101.
- [2] Maringe, C., Spicer, J., Morris, M., Purushotham, A., Nolte, E., Sullivan, R., Rachet, B. and Aggarwal, A., 2020. The impact of the COVID-19 pandemic on cancer deaths due to delays in diagnosis in England, UK: a national, population-based, modelling study. *The Lancet Oncology*, 21(8), pp.1023-1034.
- [3] Solomon, M.D., McNulty, E.J., Rana, J.S., Leong, T.K., Lee, C., Sung, S.H., Ambrosy, A.P., Sidney, S. and Go, A.S., 2020. The Covid-19 Pandemic and the Incidence of Acute Myocardial Infarction. *New England Journal of Medicine*.
- [4] Ueda, M., Nordström, R. and Matsubayashi, T., 2020. Suicide and mental health during the COVID-19 pandemic in Japan. *medRxiv*.
- [5] Tanaka, T. and Okamoto, S., 2020. Suicide during the COVID-19 pandemic in Japan. *medRxiv*.
- [6] Chudik A, Mohaddes K, Perasan MH, Raissi M, Rebucci A. 2020. Economic consequences of Covid-19: A counterfactual multi-country analysis. VoxEU. Available at: <https://voxeu.org/article/economic-consequences-covid-19-multi-country-analysis>
- [7] McInerney, M. and Mellor, J.M., 2012. Recessions and seniors' health, health behaviors, and healthcare use: Analysis of the Medicare Current Beneficiary Survey. *Journal of Health Economics*, 31(5), pp.744-751.
- [8] Riumallo-Herl, C., Basu, S., Stuckler, D., Courtin, E. and Avendano, M., 2014. Job loss, wealth and depression during the Great Recession in the USA and Europe. *International Journal of Epidemiology*, 43(5), pp.1508-1517.
- [9] Atkeson, A., 2020. *How deadly is covid-19? understanding the difficulties with estimation of its fatality rate* (No. w26965). National Bureau of Economic Research.
- [10] ECDC, 2020. COVID-19 testing strategies and objectives Available at: [https://www.ecdc.europa.eu/sites/default/files/documents/TestingStrategy\\_Objective-Sept-2020.pdf](https://www.ecdc.europa.eu/sites/default/files/documents/TestingStrategy_Objective-Sept-2020.pdf)
- [11] ECDC, 2021. Data on testing for Covid-19 by week and country. Available at: <https://www.ecdc.europa.eu/en/publications-data/covid-19-testing>
- [12] Thornton, R., 2008. The Demand for, and Impact of, Learning HIV Status. *American Economic Review*. 98(5). 1829-63
- [13] Pavelka, M., van Zandvoort, K., Abbott, S., Sherratt, K., Majdan, M., Jarcuska, P., Krajci, M., Flasche, S., Funk, S. and CMMID COVID-19 working group, 2020. The effectiveness of

population-wide, rapid antigen test based screening in reducing SARS-CoV-2 infection prevalence in Slovakia. *medRxiv*.

[14] BBC, 2020. Covid: Mass testing in Liverpool sees ‘remarkable decline’ in cases.

Available at:

<https://www.bbc.com/news/uk-england-merseyside-55044488> Accessed 19 December 2020.

[15] Bloomberg, 2020. Seoul’s full cafes, apple store lines and show mass testing success.

Available at: <https://www.bloomberg.com/news/articles/2020-04-18/seoul-s-full-cafes-apple-store-lines-show-mass-testing-success>

Accessed 19 December 2020.

**Conflict of interest:** None

**Funding:** None

**Ethics approval:** Ethics approval was given by the Economics Research Ethics Committee of City, University of London. Approval date: 9/12/2020. Code: ETH2021-0749.

**Competing interests:** Authors declare no competing interests.

**Data and materials availability:** The data used in this study were collected via surveys.

Table 1. Bias ('survival' ratio of people with covid-19 symptoms to people with no symptoms) by waiting time for rapid test,  $N=510$

Age group	No wait/ Immediate	5-15 min	15-30 min	30-60min	1-2h	over 2h	$N$
Under 30	1.532	2.882	4.423	15.5	38	38 +	181
30-50	1.517	3.102	5.8	8.857	16	16 +	199
50+	1.708	2.706	3.571	11	11.333	11.333 +	130
Total	1.564	2.918	4.567	11.2	17.143	17.143 +	510

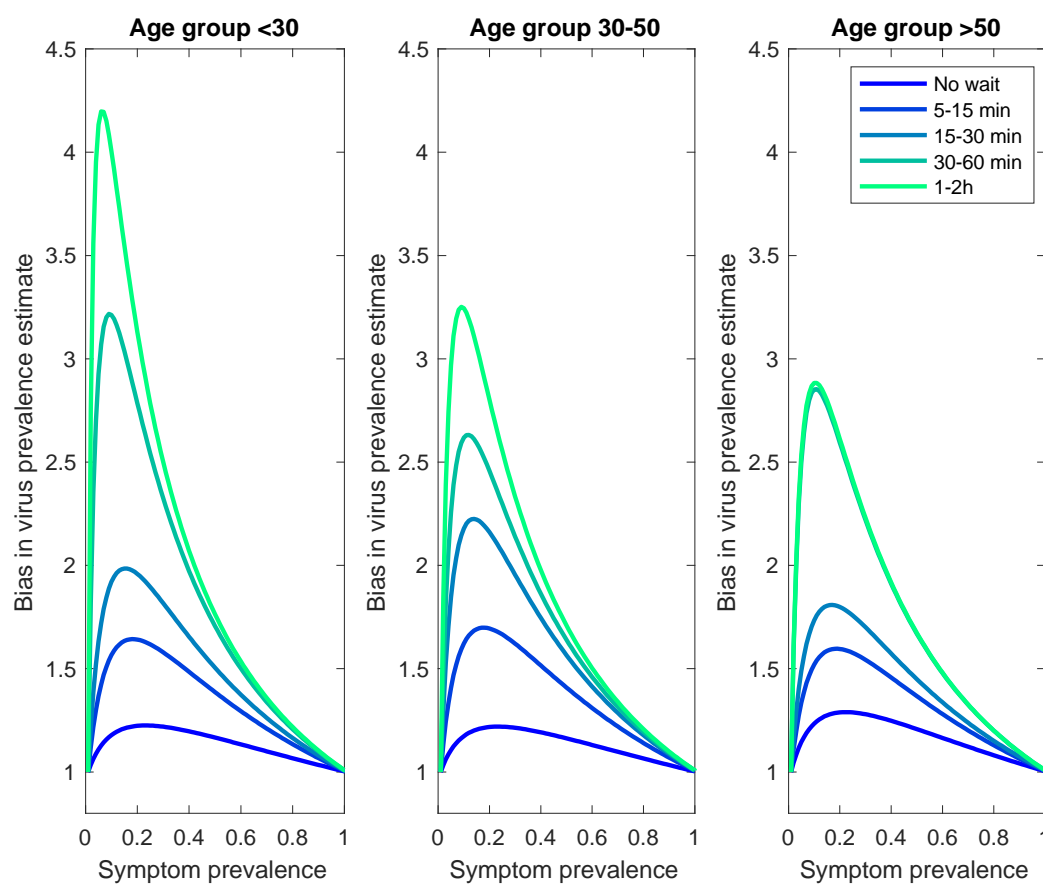


Figure 1. Best estimate of the virus prevalence bias: The ratio between reported prevalence and actual, depending on symptoms prevalence and waiting time, for the three age groups.

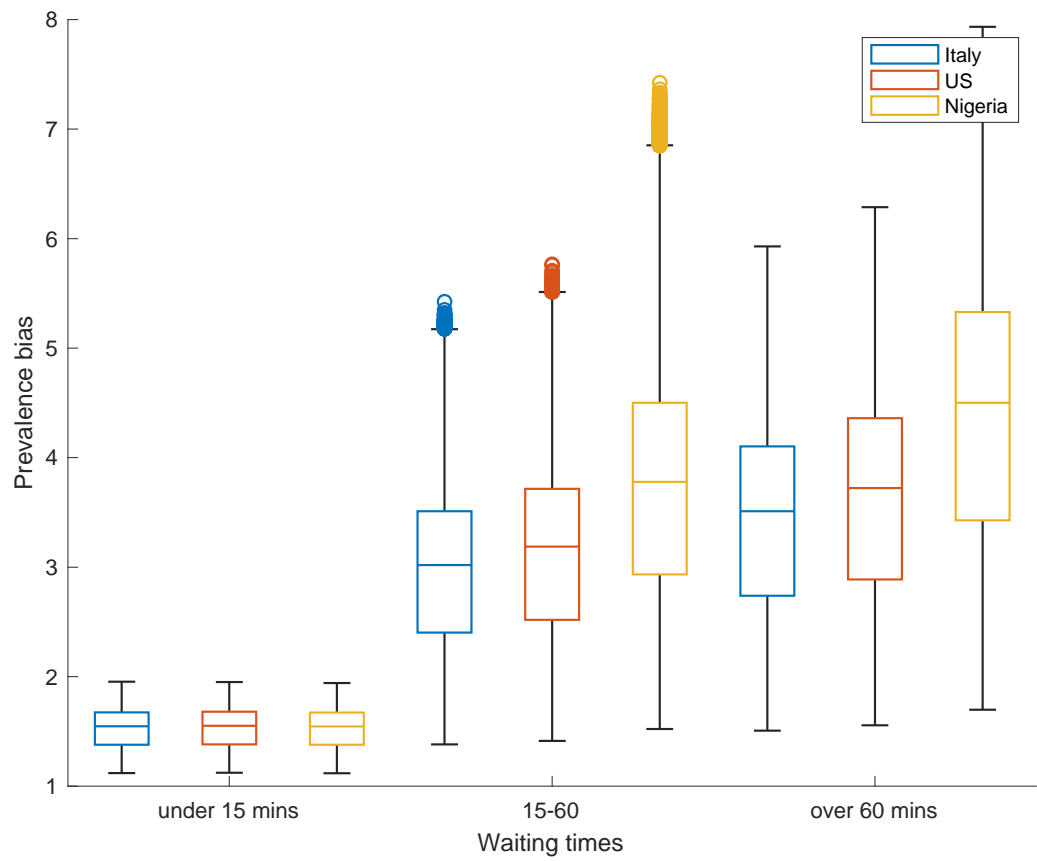


Figure 2. Simulation: How different demographic structures would affect the prevalence bias

# APPENDIX A

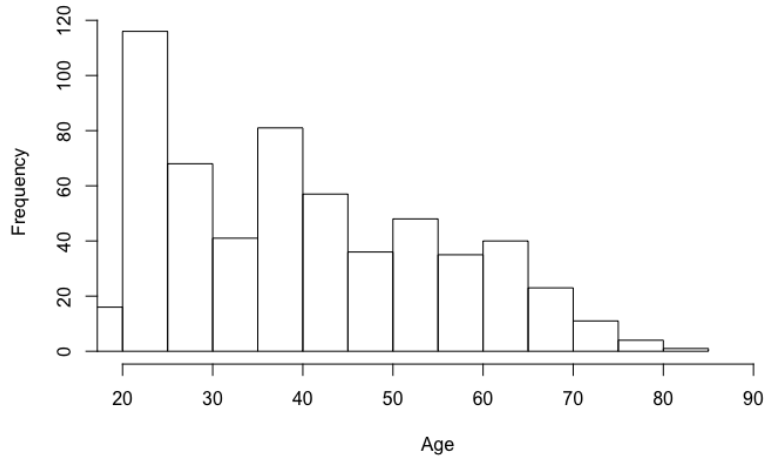


Figure A1. Age distribution of the experiment date, n=578.

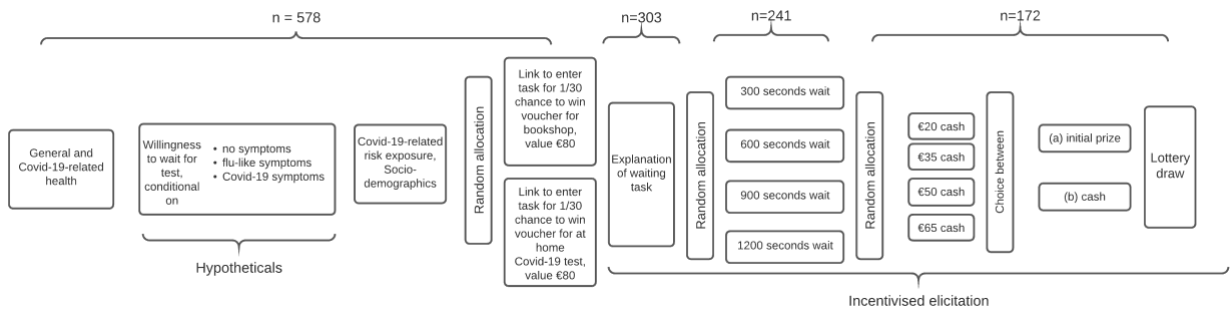


Figure A2. Flow of the experiment.

Table A1. Summary Statistics

Hypothetical willingness to test (N=578)			
<b><i>By symptoms</i></b>		<b><i>By waiting time</i></b>	
<b>No Symptoms</b>		<b>No Symptoms</b>	
Mean (SD)	2.96 (1.48)	Mean (SD)	2.39 (2.04)
Median [Min, Max]	3.00 [1.00, 5.00]	Median [Min, Max]	2.00 [0, 8.00]
<b>Flu Symptoms</b>		<b>Flu Symptoms</b>	
Mean (SD)	2.00 (1.20)	Mean (SD)	3.81 (2.26)
Median [Min, Max]	2.00 [1.00, 5.00]	Median [Min, Max]	4.00 [0, 8.00]
<b>Covid Symptoms</b>		<b>Covid Symptoms</b>	
Mean (SD)	1.46 (0.951)	Mean (SD)	5.19 (2.35)
Median [Min, Max]	1.00 [1.00, 5.00]	Median [Min, Max]	5.00 [0, 8.00]
1: certainly yes; 2: probably yes; 3: maybe; 4: probably no; 5: certainly no.		0: would not wait at all; 1 would only take it if available immediately; 3: 5-15 minutes; 4: 15-30 minutes; 5: 30-45 minutes; 6: up to an hour; 7: 1-2 hours; 8 over 2 hours	



## APPENDIX B

### Bias calculations for “Accurate Covid-19 prevalence measurement in the field“

The aim is to infer the percentage of sick people in the population from the “random” testing in the field figures, as released by several Health Agencies worldwide. The problem is that testing is voluntary, which leads to selection bias. How large is this bias?

To start, some people believe they have symptoms, some don't: call them S(ympntomatic) and H(ealthy). Note that the discussion below has to do with what people believe, not what they actually have. Also, we distinguish between people believing they have symptoms and those who do not, but the analysis readily extends to people having strong beliefs that they might be carrying the virus and those who do not.

Let the frequency of people who believe they have symptoms be  $p_s$ , or just  $p$ , with  $1-p$  being the frequency of people who do not think they have symptoms.

Of each group, some percentage turns out having the virus. Let  $v_s$  be the virus prevalence for those who believe they have symptoms,  $v_h$  for those who do not.

Of each group, some percentage are willing to take the test (for a given waiting time to take the test). Assume this only depends on symptoms, but not on actually having the virus (this assumption is mostly innocuous, unless there is a very large number of people in hospital). Let then  $t_s$  be the percentage of people who believe they have symptoms who actually take the test, and  $t_h$  for those who do not.

True prevalence is then

$$\tau = p_s v_s + (1-p_s) v_h \quad (1)$$

Given parameters, what number shows up positive in the sample (assuming that the test itself is perfect)

$$\pi = p_s t_s v_s + (1-p_s) t_h v_h \quad (2)$$

Divide by the total sampling rate

$$m = p_s t_s + (1-p_s) t_h \quad (3)$$

to get the sample prevalence (or virus frequency in the sample population)  $\phi$

Note that if  $t_s=t_h=t$ , then  $\pi = t (p_s v_s + (1-p_s) v_h)$  and  $\phi = t (p_s v_s + (1-p_s) v_h)/t = p_s v_s + (1-p_s) v_h = \tau$  which makes sense; if testing propensities are equal, there is no bias.

If on the other hand the testing propensities  $t$  are not the same, then the sample is selected leading to bias. Before we calculate the bias, express the propensities to test and be virus positive, for the people who believe they have symptoms, as a multiple of the propensities of those who do not:  $v_s = a v_h$ ,  $t_s = b t_h$ . Then, using these equations, rewrite (1), (2) and (3).

$$\begin{aligned} \tau &= p_s v_s + (1-p_s) v_h = a p_s v_h + (1-p_s) v_h = v_h (ap+1-p) \\ \pi &= p_s t_s v_s + (1-p_s) t_h v_h = ab p t_h v_h + (1-p) t_h v_h = t_h v_h (abp + 1-p) \\ m &= p_s t_s + (1-p_s) t_h = b p_s t_h + (1-p_s) t_h = t_h (bp+1-p) \end{aligned}$$

Simplify notation by writing  $p$  for  $p_s$  and calculate  $\phi = \pi/m = t_h v_h (abp + 1-p) / t_h (bp+1-p) = v_h (abp + 1-p) / (bp+1-p)$

Now, to find the size of the bias, divide  $\phi/\tau$   
 $v_h (abp + 1-p) / (bp+1-p) / v_h (ap+1-p)$

$$\Leftrightarrow \text{The bias in estimates } \beta = (abp + 1-p) / (bp + 1-p) / (ap + 1-p)$$

#### Examples

Suppose  $a=1$

$$(bp+1-p) / (bp+1-p) / (ap + 1-p) = 1/(p+1-p)=1$$

So, both  $a$  and  $b$  are necessary for the bias to exist, which makes sense.

Suppose  $a=b>1$

(conceptually it is not unlikely that the two propensities be of similar magnitude, since the higher the risk when I have symptoms, the more likely it should be that I seek testing)

The bias is then  $(a^2 p + 1-p) / (ap + 1-p)^2$

Let  $a=b=3$

$9 p+1-p / (3p + 1 - p)^2$

In this case, the  $p$  leading to the worst bias is around 0.3,  $\beta$  becomes 1.35.

Now suppose  $a=b=10$

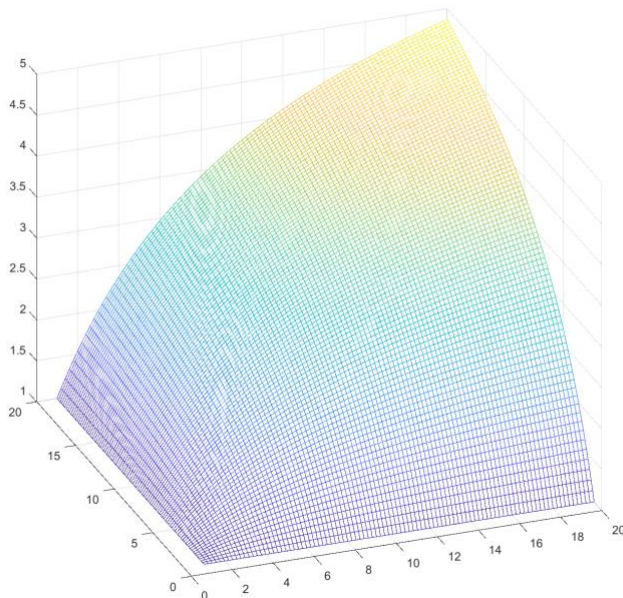
$100p+1-p / (10p+1p)^2$

$\beta=3$  at  $p$  0.1, at  $p$  0.05 it still is 2.7.

Suppose  $p=0.1$

Then  $\beta=(0.1ab + 0.9) / (0.1b+0.9) / (0.1a+0.9)$

This function is plotted in the next graph, with  $a$  in the  $x$  axis, and  $b$  in the  $y$  axis.



Meaning for  $a=b=20$ , street testing is overestimating the virus prevalence about 5 times.

## Getting $p_s$ from $\phi$

While we suggest to get  $p_s$  through random (unbiased) polling of people about their perceived symptoms, it is also possible to calculate it using  $a$ ,  $b$ ,  $v_h$  and  $\phi$  as follows.

Start with the definition of  $\phi=v_h (abp + 1-p) / (bp+1-p)$

$$\Leftrightarrow \phi(bp+1-p)=v(abp + 1-p)$$

$$\Leftrightarrow \phi bp+\phi-\phi p= vabp+v-pv$$

$$\Leftrightarrow \phi bp-\phi p+pv-vabp =v-\phi$$

$$\Leftrightarrow p(\phi b-\phi+v-vab)=(v-\phi)$$

$$\Leftrightarrow p=(v-\phi)/(\phi b-\phi+v-vab)$$

Note:

The denominator is negative

$$\phi b-\phi+v-vab<0$$

$$\Rightarrow \phi(b-1)>v(1-ab) \text{ (since } v(1-ab) \text{ is negative)}$$

$$\Rightarrow \varphi > v(1-ab)/(b-1)$$

Which is true since  $b-1$  is positive.

For  $v < \varphi$  the numerator is also negative, meaning  $p$  is positive.

If  $\varphi = v_h$  then symptoms prevalence is 0, all people in the sample have no symptoms, and  $v_h$  show positive in the test.

If  $\varphi = v_s = av_h$  then  $(v-av)/(avb-av+v-vab)=1$ , meaning everyone has symptoms,  $p=1$ .

Obviously  $\varphi$  cannot be above  $v_s$  (sample prevalence is highest if you only have people with symptoms in the sample, in which case not more than  $v_s$  can be positive)!

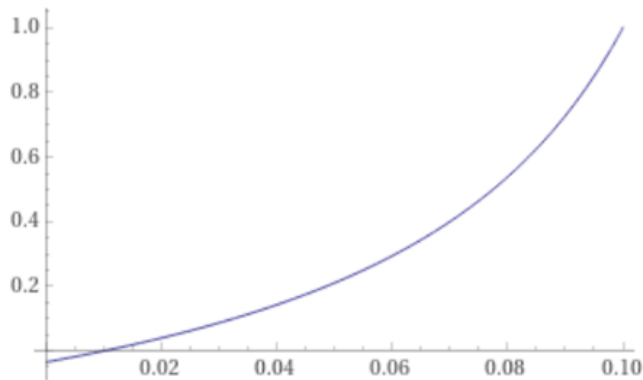
So we can debias the health agencies' numbers without knowing  $p_s$

Again, it is easier not to do street testing, but to use  $v_h$  and  $v_s$  and poll about  $p$ .

### Examples

Suppose  $a=10$ ,  $b=3$  and  $v_h=0.01$

$$\Leftrightarrow p = (0.01 - \varphi) / (3\varphi - \varphi + 0.01 - 0.3) = (0.01 - \varphi) / (2\varphi - 0.29)$$



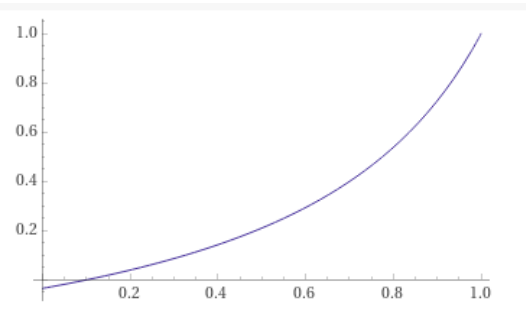
So, for example  $\varphi=0.1$  yields  $p = -0.09 / -0.09 = 1$ .

This makes sense. Everyone had symptoms, and  $v_s=0.1$  means that 10% had the virus, which is the proportion you will find in any sample. The interpretation is that with such a low true virus prevalence, the only way to get a relatively high  $\varphi$  is if there are *only* symptomatic people.

Suppose  $a=10$ ,  $b=3$  and  $v_h=0.1$

$$P = (0.1 - \varphi) / (3\varphi - \varphi + 0.1 - 3) = (0.1 - \varphi) / (2\varphi - 2.9)$$

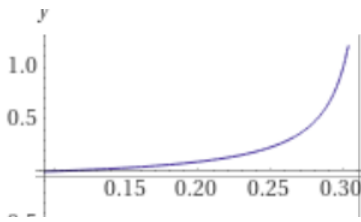
So, in this case,  $p$  is about half  $\varphi$  for many values.



Now, let  $a=3$ ,  $b=10$  and  $v_h=0.1$

The effect of  $a$  and  $b$  is not symmetric.

$$p = (0.1 - \varphi) / (10\varphi - \varphi - 2.9)$$



Suppose that some national agency is asking about (perceived) symptoms before testing. It is then easier to find the symptom prevalence in the general population  $p_s$

Symptom prevalence in the test would be  
 $\chi = \frac{pt_s}{(pt_s + (1-p)t_h)} = \frac{bpt_h}{(bpt_h + (1-p)t_h)}$   
 $\Rightarrow \chi = \frac{bp}{(bp + 1-p)} \Rightarrow \chi = \frac{bp}{bp + 1-p}$

So true symptom prevalence is

$$p = \frac{\chi}{(b + \chi - b\chi)}$$

For a relatively low  $b=3$

plot	$p = \frac{x}{-bx + b + x}$ where $b = 3$	$x = 0$ to $1$
------	---	----------------

Plot:

