



Citation for published version:

Pagel, C & Yates, CA 2021, 'Tackling the pandemic with (biased) data', *Science*, vol. 374, no. 6566, pp. 403-404. <https://doi.org/10.1126/science.abi6602>

DOI:

[10.1126/science.abi6602](https://doi.org/10.1126/science.abi6602)

Publication date:

2021

Document Version

Peer reviewed version

[Link to publication](#)

This is the author's version of the work. It is posted here by permission of the AAAS for personal use, not for redistribution. The definitive version was published in *Science* on Vol. 374, Issue 6566 on 21/10/2021, DOI: [10.1126/science.abi6602](https://doi.org/10.1126/science.abi6602)

University of Bath

Alternative formats

If you require this document in an alternative format, please contact:
openaccess@bath.ac.uk

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

VIEWPOINT: COVID-19

Tackling the pandemic with (biased) data

Data are crucial for understanding and addressing the pandemic, but there are pitfalls

By Christina Pagel¹ and Christian A. Yates²

Accurate and near-real time data about the trajectory of the COVID-19 pandemic have been crucial in informing mitigation policies. Because choosing the right mitigation policies relies on an accurate assessment of the current state of the local epidemic, the potential ramifications of misinterpreting the data are serious. Each data source has inherent biases and pitfalls in interpretation. The more data sources that are interpreted in combination, the easier it is to detect genuine changes in the course of the epidemic. Recently, in many countries, this has involved disentangling the varying impact of rising, but heterogeneous, vaccination rates, relaxation of mitigations and the rise of new variants such as Delta.

The exact data collected, and their accuracy will vary by country. Typical data common to many countries are: numbers of tests, confirmed cases, hospital and intensive care unit (ICU) admissions/occupancy, deaths and vaccinations (1). Many countries additionally sequence a proportion of new positive tests to identify and track emerging new variants. Some countries also now collect and publish data on infections, hospitalisations and deaths by vaccination status (e.g. Israel, UK). Stratifying all available data by different demographic factors (e.g. age, location, measures of deprivation, ethnicity) is crucial for understanding patterns of spread, potential impact of policies and efficacy of vaccines (age, timing of breakthrough infections and prevalent variants).

We must also be aware of what data is not being collected. For instance, persistent symptoms of COVID-19 ("Long COVID") were recognised as a long-term adverse outcome by the autumn of 2020. However, no simple diagnostic test has been associated with the up to 200 different symptoms (2). Counting Long COVID relies on a clinical diagnosis, based on a history of having had COVID-19 and a failure to fully recover, with development of some characteristic symptoms, and with no obvious alternative cause (3). These features make it very difficult to measure routinely and so it rarely is. As a result, Long COVID is often neglected in epidemic decision-making. Failure to account for the disease load associated with Long

COVID may lead to unnecessary long-term societal health burden.

The feedback between different types of outcomes, different COVID variants, different mitigation policies (including vaccination) and individual risks (a combination of exposure and clinical risk) is complex and must be factored into both interpretation of data and the development of policy. Using all available data to quantify transmission is crucial to ensuring rapid and effective responses to early phases of renewed exponential growth and to evaluating how well mitigation measures are working. Relying too much on a single data source, or without disaggregating data, risks fundamentally misunderstanding the state of the epidemic.

The inherent biases and lags in data are particularly important to understand from the point of view of policy makers. Because of the natural timescales of COVID-19 disease progression (Figure 1), policy changes can take several weeks to show in the data while purely reactive policy making is likely to be ineffective. When cases are rising, increases in hospital admissions and deaths will follow. When a new variant is outcompeting existing strains, it is likely to become dominant without action to suppress. The precautionary principle suggests acting early and emphatically. Conversely, when releasing restrictions, it is vital that governments wait long enough to assess the impact before continuing with re-opening.

The most up to date indicator of the state of the epidemic is typically the number of confirmed cases, as ascertained through testing of both symptomatic individuals and those tested frequently regardless of symptoms. Symptom-based testing is likely to pick up more adults and fewer younger individuals (4). Other biases include test accessibility, reporting lags, and the ability to act lawfully upon receiving a positive result.

Substantial changes in the number of people seeking tests may further confound case figures (5). Case positivity rates may provide a more accurate reflection of the state of the epidemic (6) but are themselves dependent on the mix of symptomatic and asymptomatic people being tested.

COVID-19 variants have been an important driver of local epidemics in 2021. The four main SARS-CoV-2 variants of concern to date have been B.1.1.7 (Alpha),

B.1.351 (Beta), P.1 (Gamma) and B.1.617.2 (Delta). Some have been more transmissible (Alpha), some have substantial resistance to previous infection or vaccines (Beta) and some have elements of both (Gamma and Delta) (7). At the time of writing, Delta's high transmissibility combined with some immune evasion has made it the world's dominant variant. Determining which variants pose a significant threat is difficult and takes time, particularly where many variants co-circulate. This is especially true for situations where a dominant variant is declining, and a new one growing. While the declining variant remains dominant, its decrease masks increases in the new variant, as case numbers remain unchanged or fall overall. Only when a new variant becomes dominant does its growth become apparent in the aggregated case data, by which time it is, by definition, too late to contain its spread. We have seen exactly this dynamic play out across the world with Delta over the second and third quarters of 2021.

With multiple variants circulating, there are, effectively, multiple epidemics occurring in parallel and must be tracked separately. This typically requires the availability of sequencing data, unfortunately rare in most countries. Sequencing takes time and so it is typically a few weeks out of date. These lags, and the uncertainty in sampling can lead to hesitancy in acting. The rapid path to dominance of the Delta variant in the UK highlights the need for action when a rapidly growing variant represents only a few percent (or less) of overall case load.

Hospital admissions or occupancy data do not have the biases associated with testing behaviours and provide unequivocal evidence of widespread transmission, its geography and demographics. However, hospital admissions lag infections more than reported cases, rendering these data less useful for proactive decision making. Hospital data are also biased towards older people who are more likely to suffer severe COVID-19, and now, unvaccinated populations. Intensive care occupancy data show a younger age profile since younger patients have a better chance of benefitting from the invasive treatment procedures (8).

Deaths are the most lagged indicator – typically occurring 3 or more weeks post infection and with an additional lag in regis-

¹University College London; London, UK. ²University of Bath; Bath, UK. Email: c.pagel@ucl.ac.uk, c.yates@bath.ac.uk

1 tration and reporting. Death data should
2 never be used to inform real-time policy de-
3 cisions. Instead, deaths are an unambiguous
4 eventual measure of the success of a coun-
5 try's epidemic strategy and implementation.

6 The age distribution of those who even-
7 tually die from COVID-19 is different again
8 from other metrics of the epidemic –
9 skewed furthest towards older age groups
10 (9). Those with clinical risk factors (immu-
11 nodeficiency, obesity, existing lung condi-
12 tions etc), high exposure (healthcare work-
13 ers, low-income workers) and the
14 unvaccinated are over-represented in
15 COVID death figures.

16 In countries with high vaccination rates
17 it is clear that vaccination has had a signifi-
18 cant impact - reducing COVID-19 cases,
19 hospitalisations and deaths. However, when
20 looking at the raw numbers in highly vac-
21 cinated populations it can be the case that
22 more fully vaccinated people are dying of
23 COVID-19 than unvaccinated. If these raw
24 statistics are misinterpreted, or worse de-
25 liberately misused, they can damage vaccine
26 confidence. In reality, more vaccinated peo-
27 ple may die than unvaccinated because such
28 a high proportion of people are vaccinated
29 (10). This does not mean vaccines are not
30 effective at preventing death. Looking at the
31 rates of death in vaccinated and unvaccinat-
32 ed individuals separately demonstrates that
33 vaccines provide significant protection
34 against severe disease and death. This ex-
35 ample illustrates how important it is to cu-
36 rate and manage the way in which data is
37 presented in the midst of an epidemic.

38 Each country has established its own
39 vaccination priority lists and dosing sched-
40 ules in order to best achieve its goals
41 (11,12). Each of these strategies will mani-
42 fest differently in the data. Additionally,
43 many countries are using multiple vaccines
44 in tandem and employing them differently
45 for different demographics. Some countries
46 are vaccinating adolescents and others are
47 not or not offering them the full approved
48 dose. Most vaccines require two doses,
49 spaced between 3 and 12 weeks apart, ex-
50 cept for the Johnson & Johnson single dose
51 vaccine. This matters, particularly as differ-
52 ent variants spread, because different vac-
53 cines have different effectiveness after 1
54 and 2 doses, different timelines to full effec-
55 tiveness. and different effectiveness against
56 variants (for instance, mRNA vaccine-
57 mediated immunity is less impacted by the
58 Beta variant than immunity from vaccines
59 based on adenoviruses (13)).

Data published on the vaccination deliv-
ery itself must thus go beyond the raw

numbers of people vaccinated. Vaccine up-
take must be reported by whether fully or
partially (1-dose in a 2-dose regimen) vac-
cinated and using the whole population as a
denominator. It is vital to disaggregate vac-
cine data by age, gender and ethnicity as
well as location so that it is possible, for ex-
ample, to understand the impact of depriva-
tion on vaccine coverage or vaccine hesitan-
cy in particular demographics. When
interpreting vaccination data in the context
of immunity provided it is important to re-
member there is a lag between delivery and
the build-up of immunity.

Data on re-infection and post-
vaccination (breakthrough) infection are al-
so important in order to determine the rela-
tive benefits of infection-mediated and vac-
cine-mediated immunity and the length of
protection offered.

Studies which show that those who were
immunized earlier were catching covid with
higher rates than those vaccinated later
may, at face value, suggest waning vaccine
protection (14). Such studies have already
been used as justification for vaccine boost-
er programmes. However, any study sug-
gesting waning immunity must be extreme-
ly careful to ensure the 'early' and 'late'
subgroups are properly controlled. Differ-
ences in prior exposure, affluence, educa-
tion-level, age and other demographic fac-
tors between early and late cohorts may be
enough to explain the disparities in covid in-
fection rates even in the absence of waning
immunity. Waning must also be reported
separately for different outcomes: for in-
stance there might be waning in terms of
preventing symptomatic infection but far
less or none in preventing death (15). In ad-
dition, there are clear ethical concerns sur-
rounding mass-booster programmes in rich
countries whilst many poorer countries
have been unable procure vaccines to pro-
tect the majority of their populations. The
evidence for boosters, if it is to be acted up-
on, must be unequivocal.

As we move into the vaccination era, re-
ported cases, hospitalisations and deaths
should also be disaggregated by vaccination
status (and by which vaccine), which will be
easier in countries where national linked
datasets exist. Whilst we already have ac-
cess to many sources of data, this finer-
grained information would help under-
standing of emerging issues including
breakthrough infection, reinfection, new
variants and waning immunity. Addition-
ally, incorporating Long COVID into routine
reporting and policy making is crucial. Con-
sistent diagnostic criteria and well-

controlled studies will be vital to this effort.
These elusive data will be of crucial im-
portance as we navigate our way out of the
epidemic.

REFERENCES AND NOTES

1. M Roser *et al.*, Our World Data (2021) <https://bit.ly/3kepLgw>
2. H.E. Davis *et al.*, *E. Clin. Med.* 2021; 38 101019
3. M. Sivan *et al.*, *BMJ* 2020; 371: m4938.
4. S.M. Moghadas *et al.*, *Proc. Natl. Acad. Sci.* 117 17513 (2020).
5. J. Wise *BMJ*. 370 m3678 (2020).
6. M. Hartman, COVID-19 Testing: Understanding the "Percent Positive" (2021) <https://bit.ly/3CeN8wl>.
7. C.E. Gómez *et al.*, *Vaccines* 9 243 (2021).
8. A.B. Docherty *et al.*, *BMJ*. 369 m1985 (2020).
9. ONS, Deaths registered weekly in England and Wales by age and sex: covid-19 (2021) <https://bit.ly/3Ci2obS>.
10. C Yates Significant proportions of people admitted to hospital, or dying from covid-19 in England are vaccinated—this doesn't mean the vaccines don't work (2021) <https://bit.ly/3kfql0H>
11. CDC COVID-19 Vaccine Information for Specif2c Groups (2021) <https://bit.ly/39aijwp>
12. JCVI Priority groups for coronavirus (COVID-19) vaccination: advice from the JCVI (2020) <https://bit.ly/2VlhfwC>
13. J.P. Moore. *JAMA* 325 1251 (2021).
14. Y. Goldberg *et al.* Waning immunity of the BNT162b2 vaccine: A nationwide study from Israel (2021) <https://bit.ly/3kgDaV9>
15. PHE, Duration of Protection of COVID-19 vaccines against clinical disease (2021) <https://bit.ly/3CCoVAq>

ACKNOWLEDGMENTS (OPTIONAL)

CP and CY are both members of Independent SAGE: <https://www.independentsage.org/>

DOI

PHOTO CREDIT GOES HERE

COVID-19 infection progression

An approximate timeline from infection with COVID-19 to death including the times at which these figures are expected to show up in the different data sources. We are illustrating death as one end point because these are the outcomes that are most relevant to the statistics, however it should be noted that the vast majority of people infected with COVID-19 will survive.