

Local prevalence of transmissible SARS-CoV-2 infection : an integrative causal model for debiasing fine-scale targeted testing data

George Nicholson^{1,9,*}, Briec Lehmann^{1,9,*}, Tullia Padellini^{2,9}, Koen B Pouwels^{3,4},
5 Radka Jersakova^{5,9}, James Lomax^{5,9}, Ruairidh E King^{6,9}, Ann-Marie Mallon^{6,9},
Peter J Diggle^{7,9}, Sylvia Richardson^{8,9}, Marta Blangiardo^{2,9}, and Chris Holmes^{1,5,6,9}

¹University of Oxford, UK

²MRC Centre for Environment and Health, Dept of Epidemiology and
Biostatistics, Imperial College London

10 ³Health Economics Research Centre, Nuffield Department of Population Health,
University of Oxford, UK

⁴The National Institute for Health Research Health Protection Research Unit in
Healthcare Associated Infections and Antimicrobial Resistance at the University of
Oxford, University of Oxford, Oxford, UK

15 ⁵The Alan Turing Institute, London, UK

⁶MRC Harwell Institute, Harwell, UK

⁷CHICAS, Lancaster Medical School, Lancaster University, UK

⁸MRC Biostatistics Unit, University of Cambridge, UK

⁹Joint Biosecurity Centre, Turing and Royal Statistical Society Laboratory

20 *These authors contributed equally to this research.

Abstract

Targeted surveillance testing schemes for SARS-CoV-2 focus on certain subsets of the population, such as individuals experiencing one or more of a prescribed list of symptoms. These schemes have routinely been used to monitor the spread of SARS-CoV-2 in countries across the world. The number of positive tests in a given region can provide local insights into important epidemiological parameters, such as prevalence and effective reproduction number. Moreover, targeted testing data has been used to inform the deployment of localised non-pharmaceutical interventions. However, surveillance schemes typically suffer from ascertainment bias; the individuals who are tested are not necessarily representative of the wider population of interest. Here, we show that data from randomised testing schemes, such as the REACT study in the UK, can be used to debias fine-scale targeted testing data in order to provide accurate localised estimates of the number of infectious individuals. We develop a novel, integrative causal framework that explicitly models the process underlying the selection of individuals for targeted testing. The output from our model can readily be incorporated into longitudinal analyses to provide local estimates of the reproduction number. We apply our model to characterise the size of the infectious population in England between June 2020 and January 2021. Our local estimates of the effective reproduction number are predictive of future changes in positive case numbers. We also capture local increases in both prevalence and effective reproductive number in the South East from November 2020 to December 2020, reflecting the spread of the Kent variant. Our results illustrate the complementary roles of randomised and targeted testing schemes. Preparations for future epidemics should ensure the rapid deployment of both types of schemes to accurately monitor the spread of emerging and ongoing infectious diseases.

Introduction

The spread of the new severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) and the ensuing outbreaks of coronavirus disease 2019 (COVID-19) have placed a significant burden on public health in the United Kingdom (UK). As of 12 April 2021, the number of people who have died within 28 days of a positive SARS-CoV-2 test was 127,100 [1, 2]. In response to the ongoing epidemic, the UK government has implemented a number of non-pharmaceutical interventions (NPIs) to reduce transmission of SARS-CoV-2, ranging from localised measures, such as the closures of bars and restaurants, to full national lockdowns [3]. The localised measures have been employed through a regional tier system, with areas being placed under varying levels of restrictions according to data such as the number of positive polymerase chain reaction (PCR) tests returned there over a seven-day interval (or local *weekly* positive tests) [4]. Following a third national lockdown that began on the 6th January 2021, the UK is currently undergoing a staged relaxation of restrictions [5]. Accurate local measures of prevalence and incidence are needed to assess the need for any changes to this plan and importantly to measure the relative impact of the individual stages thereby providing crucial information for future waves and pandemics both in the UK and

more globally.

In the UK, there are two major, ongoing studies that undertake randomised testing to provide
60 an insight into the prevalence of SARS-CoV-2. Since April 2020, the Office for National Statistics
(ONS) COVID-19 Infection Survey (CIS) tests a random sample of people living in the community
with longitudinal follow-up. [6]. The survey is designed to be representative of the UK population,
with individuals aged 2 years and over in private households randomly selected from address
lists and previous ONS surveys, though it does not explicitly cover care homes, the sheltering
65 population, student halls or individuals currently being hospitalised. The REal-time Assessment of
Community Transmission (REACT) study is a second nationally representative prevalence survey
of SARS-CoV-2 based on repeated cross-sectional samples from a representative subpopulation
defined via (stratified) random sampling from England’s National Health Service patient register
[7]. Importantly, both surveys recruit participants regardless of symptom status and are thus able
70 to largely avoid issues arising from ascertainment bias when estimating prevalence. The ONS CIS
uses multilevel regression and poststratification to account for any residual ascertainment effects
due to non-reponse [6] while REACT uses survey weights for this purpose.

While randomised surveillance testing readily provides an accurate statistical estimate of preva-
lence of PCR positivity, precision can be low at finer spatiotemporal scales (e.g. the lower tier local
75 authority (LTLA) level), even in large studies such as the ONS CIS and REACT surveys. The
major goal there is to unlock the information in non-randomised testing under arbitrary, unknown
ascertainment bias. While we expect the methods to apply broadly, here we focus in on Pillar
1 and Pillar 2 PCR tests conducted in England between 31st May 2020 and 24th January 2021
(lateral flow device, LFD, tests are not included; further details in *Methods–Data*). As Pillar 1
80 tests refer to “*all swab tests performed in Public Health England (PHE) labs and National Health
Service (NHS) hospitals for those with a clinical need, and health and care workers*”, and Pillar 2
comprises “*swab testing for the wider population*”, Pillar 1+2 testing has more capacity than the
randomised programs, but the protocol incurs ascertainment bias as those at elevated risk of being
infected are tested, such as frontline workers, contacts traced to a COVID-19 case, or the sub
85 population presenting with COVID-19 symptoms, such as loss of taste and smell [8]. Hence, raw
prevalence estimates from Pillar 1+2 data (as a proportion of tested population) will tend to be
biased upwards and cannot directly be used to estimate the unknown infection rate in a region
(in contrast, as a proportion of the whole population the bias is downwards as not all infected
individuals in the area are captured). Also they tend not to capture asymptomatic infection, while
90 there is evidence that asymptomatic individuals can contribute to spread of the virus [9, 10].

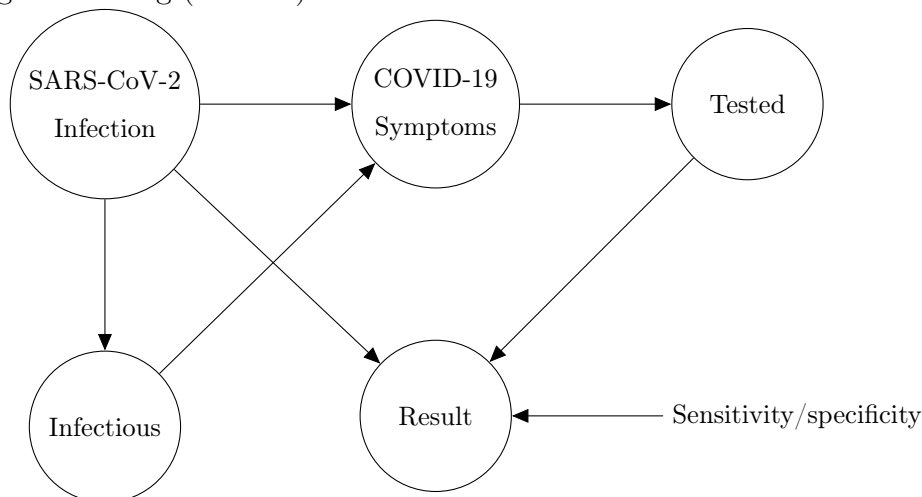
Combining data from multiple surveillance schemes can improve estimates for prevalence. For
example, Manzi et al. incorporate information from multiple, biased, commercial surveys to provide
more accurate and precise estimates of smoking prevalence in local authorities across the east of
England [11]. A number of geostatistical frameworks for infectious disease modelling based on
95 multiple diagnostic tests have been developed [12, 13, 14]. These accommodate different sources of
heterogeneity among the tests to deliver more reliable and precise inferences on disease prevalence.

To understand the ascertainment bias problem and a statistical approach to correction, it is helpful to consider a simplified causal model for Pillar 1+2 data. This is represented by a directed acyclic graph (DAG), shown in Figure 1(a), that charts the dependencies of an individual from infection status to test result. The circles indicate the binary (Yes/No) states of an individual. The DAG characterises the joint distribution of the major factors leading to the observed data. Throughout the paper we use the term *targeted* testing data to refer to data gathered under some ascertainment process distinct from (stratified) random sampling, with an exemplar being selection for testing of the subpopulation with COVID-19 symptoms, which comprises a sizeable proportion of Pillar 1+2 tests. The Pillar 1+2 DAG can be compared to that of a randomised surveillance study (shown in Figure 1(b)). The randomised nature of the test allocations in REACT renders Tested conditionally independent of Symptoms given Infected, yielding unbiased estimates of Infection rates. The DAG explicitly characterises statistically why we cannot use Pillar 1+2 data directly. The DAG also points to a potential solution if the statistical dependencies as indicated by the arrows in Figure 1(a) can be modelled, then Pillar 1+2 data *can* be used. In this paper, we describe an approach allowing characterisation and adjustment for the ascertainment bias inherent in Pillar 1+2 data.

In addition to prevalence, there are a number of epidemiological parameters that may be useful for informing localised NPIs. For example, one particular variable of interest is the (time-varying) effective reproductive number \mathcal{R}_t , defined roughly as the average number of infections caused by an infectious individual; when $\mathcal{R}_t > 1$, the epidemic will continue to spread. Estimation of these parameters relies on careful mathematical modelling supported by relevant data such as surveillance testing results or number of hospitalisations [15]. Incorporating multiple sources of data can produce more reliable parameter estimates [16], though doing so in a computationally efficient manner may be nontrivial [17]. For example, a stochastic epidemic model of the 2009 influenza outbreak in Finland that used data on hospitalisations, lab tests and vaccination data provided insights into the time-varying nature of \mathcal{R}_t but took over a month of compute time to run [18]. Given the time-sensitive nature of the current epidemic, one important modelling consideration is the timely inference of parameters; the work we present here has been developed with both accuracy and computational efficiency in mind.

The current pandemic has spurred the development of a number of models that also aim to incorporate multiple sources of data in order to estimate important epidemiological parameters, in particular \mathcal{R}_t [19]. Abbott et al. [20] generate daily estimates of \mathcal{R}_t at a national and PHE region level by incorporating case counts and death notifications, building on a model to estimate \mathcal{R}_t from incidence time series [21]. Birrell et al. [22] estimate daily, PHE regional \mathcal{R}_t using ONS CIS data, death notifications and serological data within an age-stratified transmission model. Colman et al. [23] develop methods to combine Pillar 1+2 alongside ONS CIS data to estimate the proportion of infections that result in a positive diagnosis, outputting estimates of the true incidence of infections over time. Methods have also been proposed to use serological data to adjust for the effects of biases in testing data in order to estimate the infection to fatality ratio [24, 25].

(a) Targeted testing (Pillar 2)



(b) Randomized surveillance (e.g. REACT, ONS)

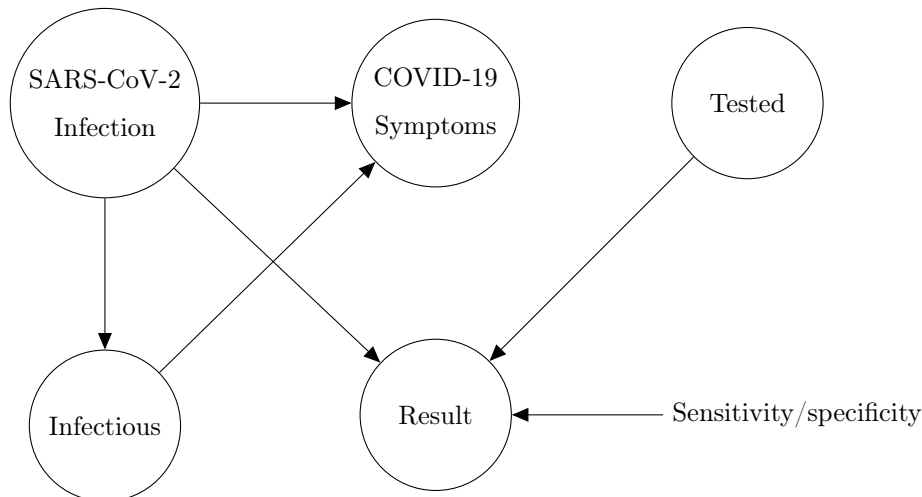


Figure 1: DAGs representing causal models underlying SARS-CoV-2 swab testing data for (a) Targeted test-and-trace data (Pillar 1+2); and (b) Randomised surveillance data (e.g. REACT). In (b), randomisation breaks the causal link between COVID-19 symptoms and swab testing. The nodes represent binary (yes/no) states for an individual in the relevant population.

The research referenced so far infers epidemiological parameters at spatially coarse scales, such as PHE region. To extend this body of work, we focus here on accurate estimation of prevalence (along with other epidemiological parameters) at a more local level, such as the LTLA.

There are two very useful websites providing LTLA-level up-to-date estimates and predictions of some epidemiological parameters (but not prevalence).¹ A team at Oxford has produced a local Covid map² as part of the Royal Society Data Evaluation and Learning for Viral Epidemics (DELVE) initiative.³ Their methods take as input Pillar 1+2 daily counts and commuter flow data, and output local estimates and predictions of \mathcal{R}_t and positive case numbers. An Imperial College team has produced a COVID-19 UK map and table.⁴ Their methodology takes as input: daily cases data, weekly deaths data, as well as daily infections from the ONS CIS and REACT data sets [26]. They output estimates and predictions of \mathcal{R}_t , positive case numbers, and change in new infections. Their results are based on the `epidemia` software [27], which is an extension of the Bayesian semi-mechanistic model introduced in [28], though detailed methods are not yet available. We have downloaded their \mathcal{R}_t estimates and find a high level of consistency in local \mathcal{R}_t estimates between our model and theirs. A group from Lancaster University has estimated daily case prevalence (proportion of infected population), incidence, and \mathcal{R}_t at the Local Authority District level in England (315 areas in total) by building an epidemic model incorporating measures of human mobility and Pillar 1 and 2 tests across England [29]. An important aspect of this approach is that it assumes each infection is eventually reflected in the Pillar 1 and 2 case reports and so does not account for possible ascertainment bias in targeted testing. Furthermore, the model requires substantial computational resources to obtain timely estimates.

Within this urgent and fast developing area of research, it is clearly important to define the aspects in which our method contributes novelty. Firstly, we have developed methods to infer local prevalence, I_t , accurately from targeted testing data. Here we work with weekly period prevalence, and explicitly target the number of infectious individuals via a correction to the estimated PCR-positive numbers. This is all novel and important in its own right – being able to estimate local prevalence accurately from targeted testing data adds an important facet to existing COVID-19 monitoring capabilities. Second, our method outputs bias-adjusted cross-sectional prevalence likelihoods $p(n_t \text{ of } N_t \mid I_t)$, where n_t and N_t are positive and total targeted test counts. This allows prevalence information from targeted data to be coherently embedded in a modular way into complex spatiotemporal epidemiological models, including those synthesising multiple data types. We exemplify this by implementing an Susceptible-Infectious-Recovered (SIR) model around our ascertainment model likelihood. Third, our local ascertainment model is based on targeted testing data alone with, uniquely to our knowledge, both the number of positive *and total* tests being modelled (n_t *and* N_t). This has two important benefits: spatiotemporal variation in testing uptake

¹Technical details of the methodology driving these websites is not yet available at the time of writing, but in both cases the peer-review process is underway.

²<https://localcovid.info/>

³<https://rs-delve.github.io/>

⁴<https://imperialcollegelondon.github.io/covid19local/#map>

and capacity is explicitly conditioned on (via N_t), and differential test specificity and sensitivity can be naturally incorporated into our causal ascertainment model.

Results

Correcting for ascertainment bias in targeted testing data

175 Figure 2(a-b) displays the percentage of positive Pillar 1+2 tests (as a proportion of those tested) against accurate prevalence estimates from the REACT study, showing a clear upward bias (each point corresponds to a single LTLA). Here we introduce a bias-correction method that aims to provide accurate estimates of prevalence at the local level as displayed in Figure 2(c-d), based on the posterior cross-sectional prevalence $p(I_t | n_t \text{ of } N_t)$.

180 With reference to the causal DAG in figure 1, we define the essential bias parameter, δ , as

$$\delta := \log \left(\frac{\text{Odds}(\text{Tested} | \text{Infected})}{\text{Odds}(\text{Tested} | \text{Not Infected})} \right) \quad (1)$$

i.e. the log odds ratio of being tested in the infected versus non-infected populations. Larger values of δ generally correspond to higher levels of ascertainment bias, i.e. a higher chance of an infected individual being selected for testing, relative to a non-infected.

Our approach combines randomised surveillance data (REACT) and targeted surveillance data
185 (Pillars 1 and 2) to infer δ at the coarse geographical level (PHE region). We then integrate this information by specifying a temporally smooth empirical Bayes (EB) prior on $\delta_{1:T}$, applied to each constituent local region (LTLA) in the local prevalence analyses. Figure 3 shows the resulting EB priors on δ ; there is potentially more variation in δ across regions early on in the sampling period (pre-September 2020), though the prior credible intervals (CIs) are quite broad and often
190 overlapping. The data provide more information on δ from October 2020 onwards, and there is a consistent upward trend for all nine PHE regions.

Cross-sectional local prevalence from targeted testing data

De-biased likelihood for modular sharing of prevalence information

Equipped with a coarse-scale (PHE-region level) EB prior on bias δ we evaluate a fine-scale
195 (LTLA-level) δ -marginalized likelihood of the form $p(n_t \text{ of } N_t | I_t, \hat{\nu}_t)$, as described at (17) in Methods—*Cross-sectional inference on local prevalence*. This de-biased prevalence likelihood can be readily exported and modularly incorporated into more complex models, as we illustrate below in Results—*Longitudinal local prevalence and transmission*.

Cross-sectional prevalence posterior

200 The δ -marginalized likelihood can be inputted directly into cross-sectional Bayesian inference, outputting the prevalence posterior $p(I_t | n_t \text{ of } N_t, \hat{\nu}_t)$ for each time point at which such count data are available. Figure 4 plots these cross-sectional prevalence posteriors beneath the raw counts

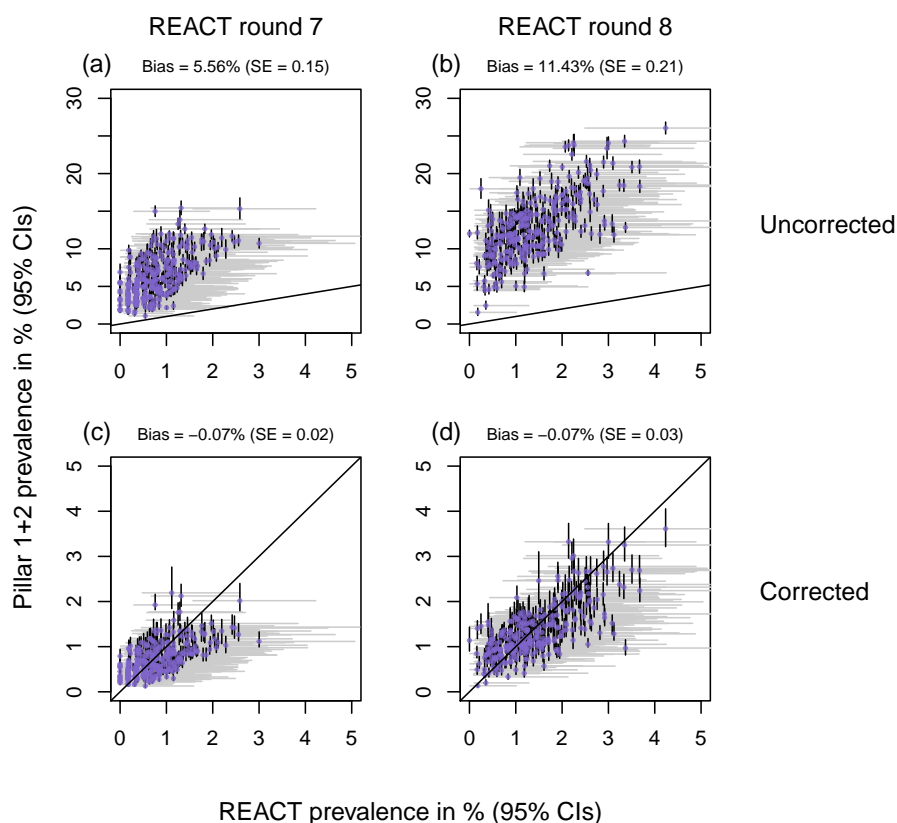


Figure 2: Uncorrected and corrected Pillar 1+2 PCR-positive prevalence estimates against (gold-standard) REACT estimates from randomised surveillance. Each point corresponds to an LTLA. Each scatter plot compares Pillar 1+2 prevalence estimates against unbiased estimates from the REACT study. Panels (a,c) show REACT round 7 data (13th Nov - 3rd Dec 2021), and (b,d) show round 8 (6th-22nd Jan 2021). Uncorrected results are shown in panels (a-b) and bias-corrected cross-sectional estimates in (c-d). Horizontal grey lines are 95% exact binomial confidence intervals from the REACT data. Vertical black lines in panels (a) and (b) are 95% exact binomial confidence intervals for from the raw, non-debiased Pillar 1+2 data. Vertical black lines in panels (c) and (d) are 95% posterior credible intervals from the debiased Pillar 1+2 data. Neither set of prevalence estimates has been corrected for false positives/negatives. Note that in panels (c) and (d), the CI widths are systematically tighter for the debiased Pillar 1+2 compared to the REACT data, pointing to the useful information content in debiased Pillar 1+2 data.

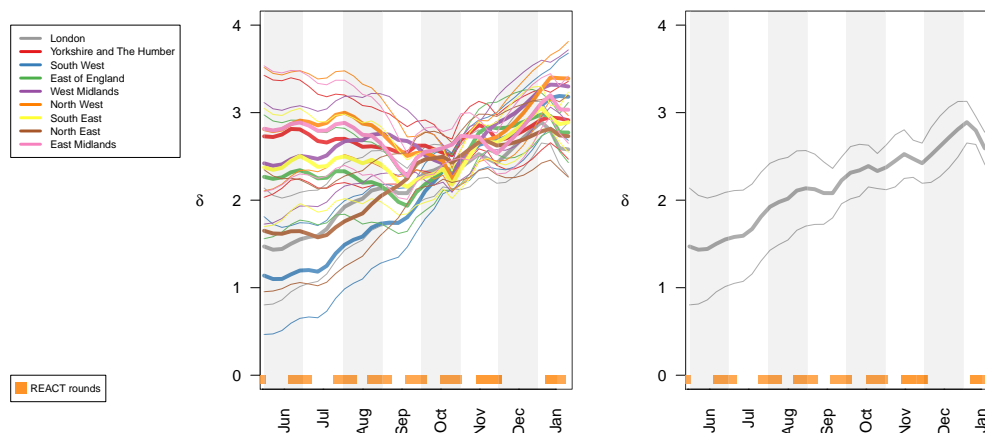


Figure 3: Smooth EB priors on bias parameters $\delta_{1:T}$. Left: Showing heterogeneous bias across the nine PHE regions. Right: London only. 95% CIs shown. Note that δ is the log odds ratio, so for example $\delta = 3$ implies that the odds of being Tested are $e^3 \approx 20$ times higher in Infected compared to Not Infected individuals

for a subset of LTLAs across the nine PHE regions. REACT sampling periods are plotted at the base of each panel, and local prevalence estimates from REACT round 7 (November 2020) and round 8 (January 2021) are also superimposed. The corrected cross-sectional prevalence estimates are consistent with the gold-standard REACT estimates, but are more precise, as expected from Bayesian principles of data synthesis. Inspecting the bias-corrected estimates within LTLA across time points, they tend to show relatively wider 95% CIs in time intervals between REACT sampling rounds (particularly in the December 2020 period between round 7 and 8) reflecting the dependence on the REACT data for good inference.

Longitudinal local prevalence and transmission

The cross-sectional de-biased likelihood can be introduced modularly into a wide variety of downstream epidemiological models. We illustrate this by using the likelihood as an input to a simple SIR epidemic model (see Methods—*Full Bayesian inference under a stochastic SIR epidemic model* and Figure 12). Figure 5(a) plots estimated prevalence against \mathcal{R}_t number at the most recent time point (the week of 2021-01-24), with each point corresponding to a single LTLA. The scatter plot provides a quick visual representation of regions where transmission rates and/or prevalence are relatively high – to illustrate, we label five LTLAs with high prevalence and/or \mathcal{R}_t estimates. The estimated longitudinal prevalence and \mathcal{R}_t for this subset of LTLAs (Figure 5(b-c)) can help further to characterise the longitudinal dynamics of prevalence and transmission in the time interval

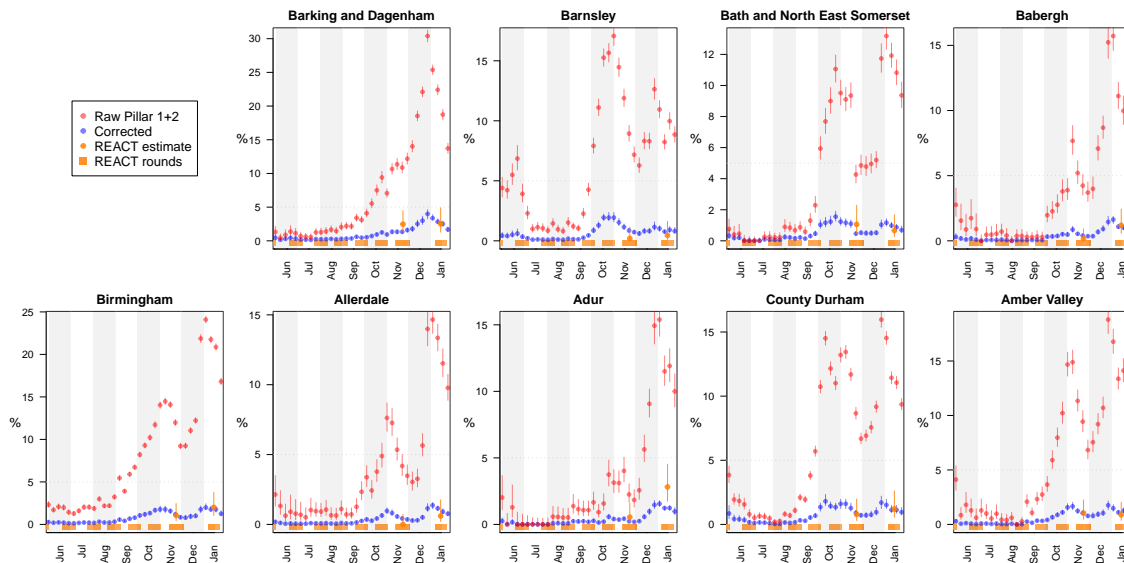


Figure 4: LTLA-level prevalence estimates: Raw Pillar 1+2 estimates, cross-sectionally corrected Pillar 1+2, and gold-standard REACT estimates (see legend). For each of the nine PHE regions, we present the constituent LTLA whose name is ranked top alphabetically.

leading up to 2021-01-24, in particular showing the estimated rate of change in prevalence and separately indicating whether \mathcal{R}_t is increasing or decreasing.

Figures 6 and 7 display spatiotemporal local prevalence and \mathcal{R}_t respectively, using a sequence of weekly maps, with each LTLA coloured according to its weekly prevalence estimate. Zoom-in boxes display the local fine-scale structure for expanded areas including London and the North West.

Relating local prevalence and transmission to spread of the UK variant

One striking feature of the maps in Figure 6 is the increasing prevalence in the London area throughout November to December 2020. This is consistent with the known arrival of the UK variant of concern (VoC) 202012/01 (lineage B.1.1.7), that emerged in the South East of England in November 2020 and which has been estimated to have a 43–90% higher reproduction number than preexisting variants [30].

We investigate this hypothesis similarly to [30], by characterising the relationship between estimated local \mathcal{R}_t and the frequency of VoC 202012/01, as approximated by the frequency of S gene target failure (SGTF) in the Taqpath sequencing assay used over this time period [31]. Figure 8 illustrates the spatial distributions of VoC 202012/01 against estimated prevalence and estimated \mathcal{R}_t from mid-November 2020 to mid-December 2020. The increase in frequency of the VoC was initially isolated to the South-East but then spread outwards, accompanied by a corresponding increase in both local estimated prevalence and \mathcal{R}_t . We observe a strong positive association between the local VoC frequency and estimated local \mathcal{R}_t , consistent with the increased transmissibility identified in [30].

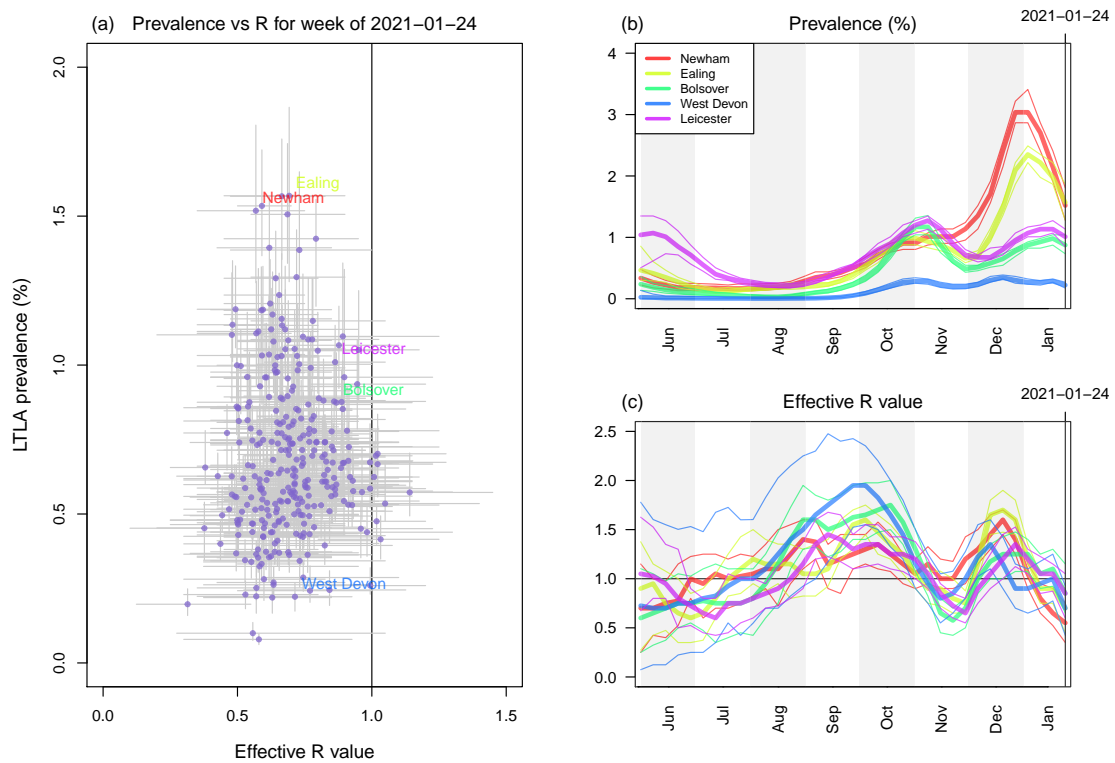


Figure 5: Outputs of longitudinal local prevalence model. (a) Scatterplot of prevalence against effective R number (each point corresponds to one LTLA). (b) Longitudinal posteriors for prevalence at a selection of LTLAs. (c) Longitudinal posteriors for \mathcal{R}_t at a selection of LTLAs.

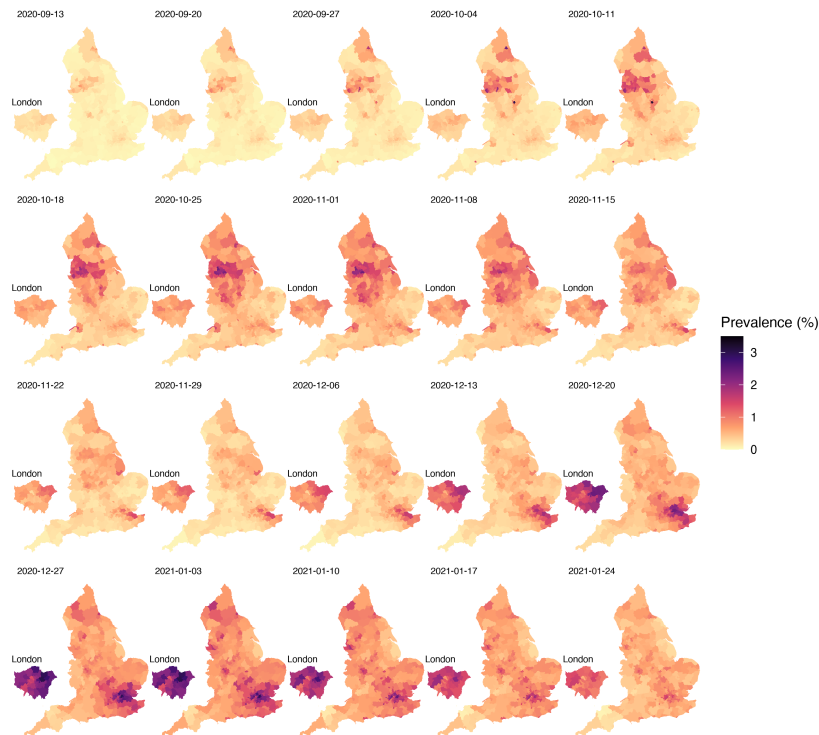


Figure 6: Longitudinal maps of estimated local prevalence from 13th September 2020 to 24th January 2021.

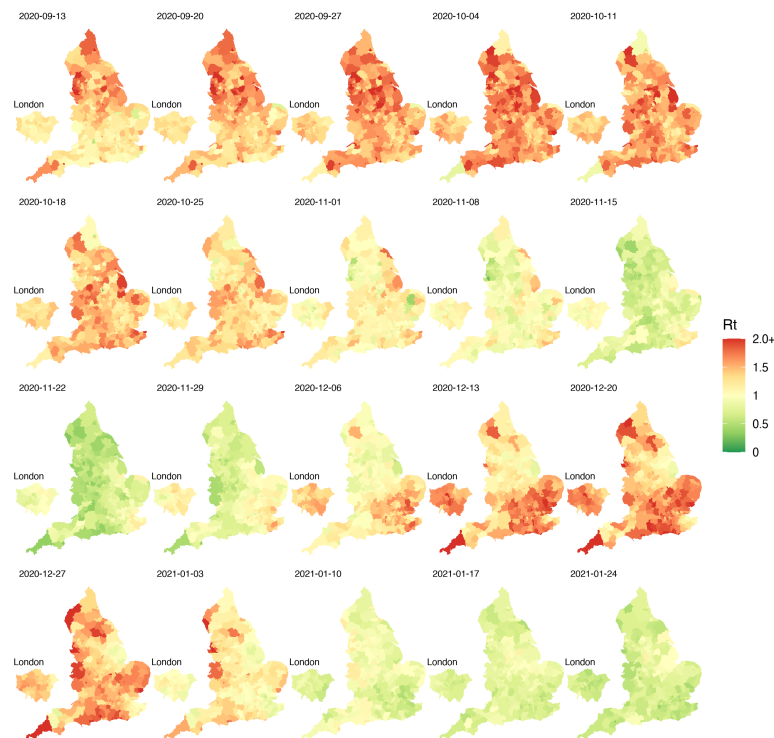


Figure 7: Longitudinal maps of estimated local R_t from 13th September 2020 to 24th January 2021.

Validation 1 (accuracy)

We qualitatively assess the performance of de-biased fine-scale (LTLA-level) prevalence estimates by measuring how well they predict LTLA-level REACT data. The validation is best described in terms of coarse-scale REACT training data and contemporaneous fine-scale REACT test data. The training data inputted are REACT PHE region-level and Pillar 1+2 LTLA-level positive (and number of) test counts for the week at the centre of the corresponding REACT round to be predicted. The test data are REACT LTLA-level positive (and number of) test counts aggregated across the relevant REACT sampling round. Figure 2(c-d) visually compares cross-sectional LTLA prevalence estimates from de-biased targeted data (i.e. based only on the training data) with accurate gold-standard estimates from REACT LTLA-level test data. The average estimated bias is reduced to low levels for comparisons with both REACT round 7 (-0.07%, SE = 0.02) and round 8 (-0.07%, SE = 0.03).

Validation 2 (prediction)

The effective reproduction number, \mathcal{R}_t , measures whether the number of infectious individuals is increasing, $\mathcal{R}_t > 1$, or decreasing, $\mathcal{R}_t < 1$, in the population at time point t . Figure 9 compares LTLA \mathcal{R}_t estimates with the future change in local case numbers. For validation purposes, here we are doing one-step-ahead at a time prediction and comparing predictions with out-of-training-sample observed statistics (fold change in raw case numbers from baseline). The results

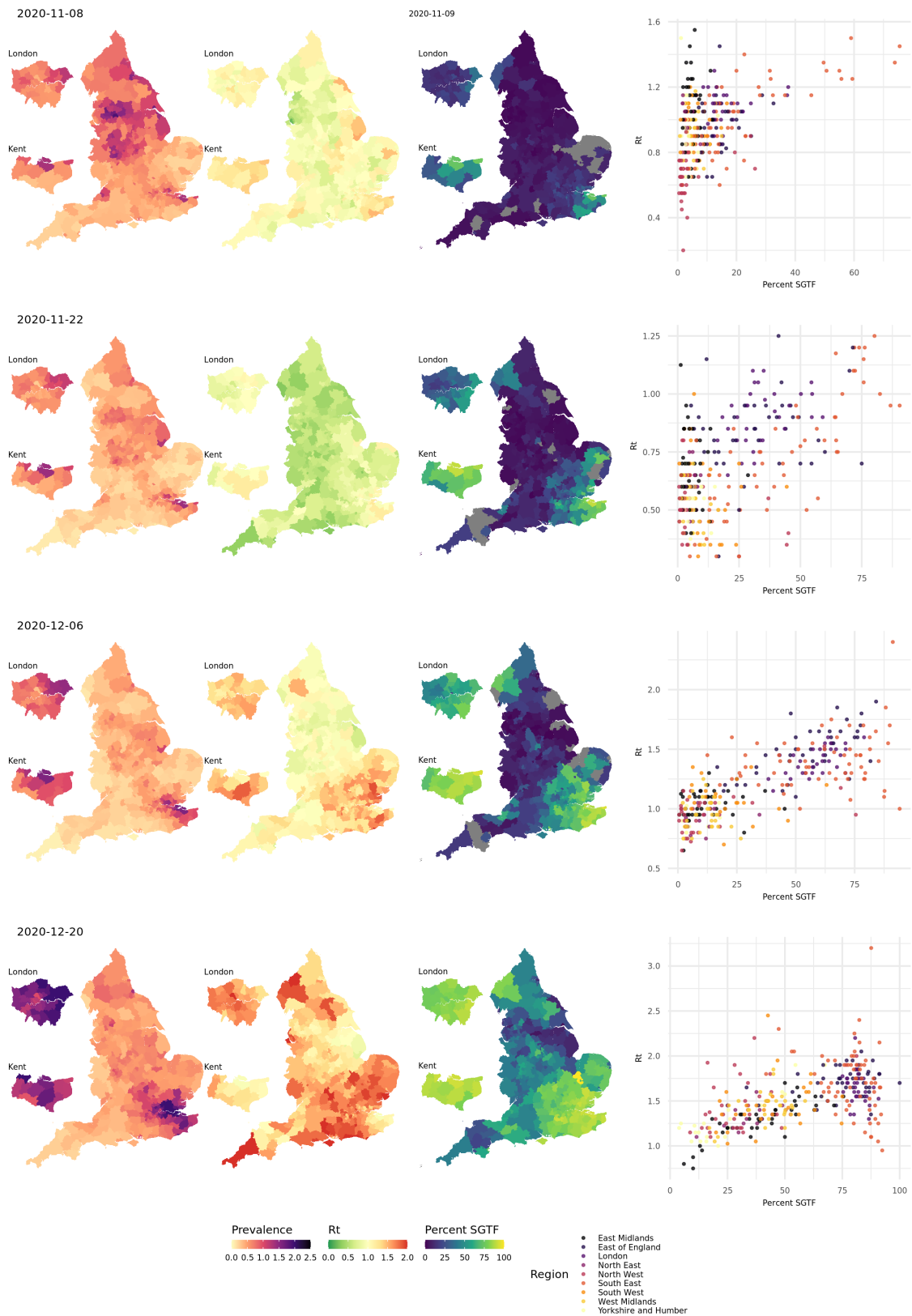


Figure 8: Maps of estimated local prevalence (left), estimated local R_t (middle), and frequency of S gene target failure (SGTF; right), and scatter plot of SGTF frequency against estimated R_t .

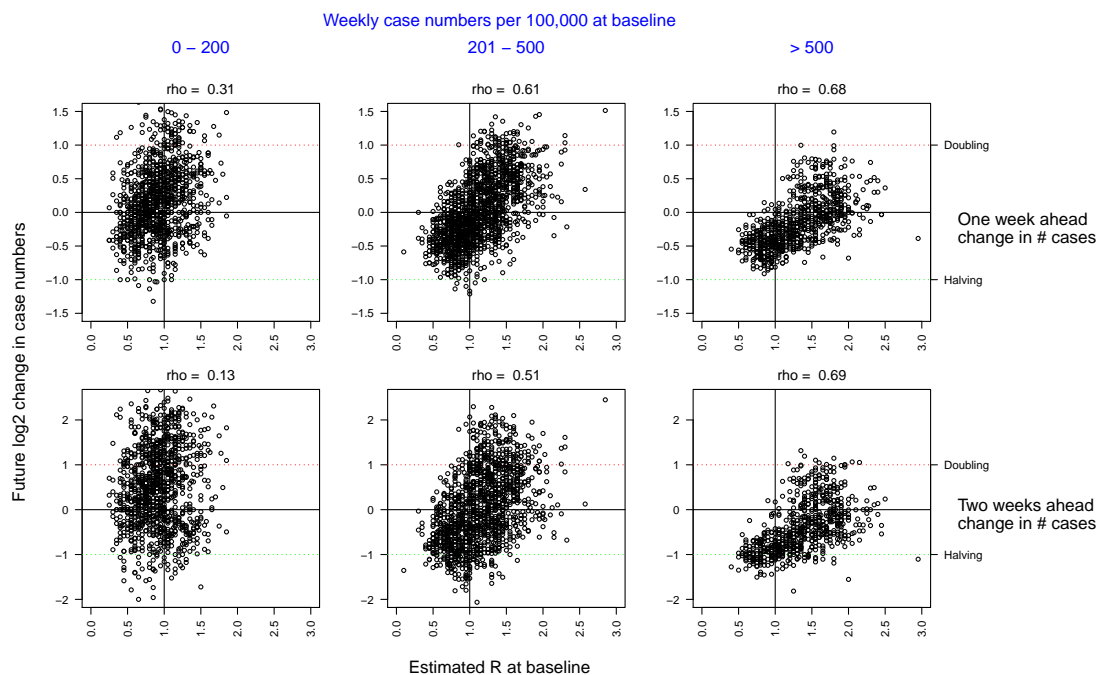


Figure 9: Predicting future change in case numbers from current estimated \mathcal{R}_t . Each point corresponds to an (LTLA, week) pair, predicting future case numbers in the LTLA using \mathcal{R}_t for that week. Future case numbers are represented by forward-in-time \log_2 fold change $\log_2(n_{t+k}/n_t)$. Case data underlying the plot are from the period 2020-10-18 - 2021-01-24. Note the number of points in each column differs based on how many LTLA-week pairs have baseline case numbers in the intervals in blue shown at the top of the plot.

260 are stratified according to baseline case numbers, and we examine predictions one week and two weeks ahead. Each point corresponds to an (LTLA, week) pair, and the results are for the period 2020-10-18 - 2021-01-24. Across each of the six scenarios presented, there is strong evidence of association between \mathcal{R}_t and future change in case numbers ($p < 2 \times 10^{-14}$). The strength of association between \mathcal{R}_t and one week ahead case numbers has Spearman's $\rho = 0.68$ for the high
 265 baseline case group (>500 cases per 100,000), decreasing to $\rho = 0.31$ in the low baseline group (≤ 200 cases per 100,000). The association remains strong when predicting caseloads two weeks ahead, with for example $\rho = 0.69$ (Spearman) for the high baseline case group.

Validation 3 (external)

270 We extracted estimates of effective reproduction number \mathcal{R}_t based on our de-biasing model likelihood implemented within a standard SIR model, illustrated in Figure 12. We compare the results to the local \mathcal{R}_t estimates outputted by at the Imperial College COVID-19 website.⁵ A cross-method comparison of longitudinal traces of \mathcal{R}_t for a subset of LTLAs is shown in Figure 10. Encouragingly for both approaches, the estimates generally display good concordance, with credi-

⁵<https://imperialcollegelondon.github.io/covid19local/#map>

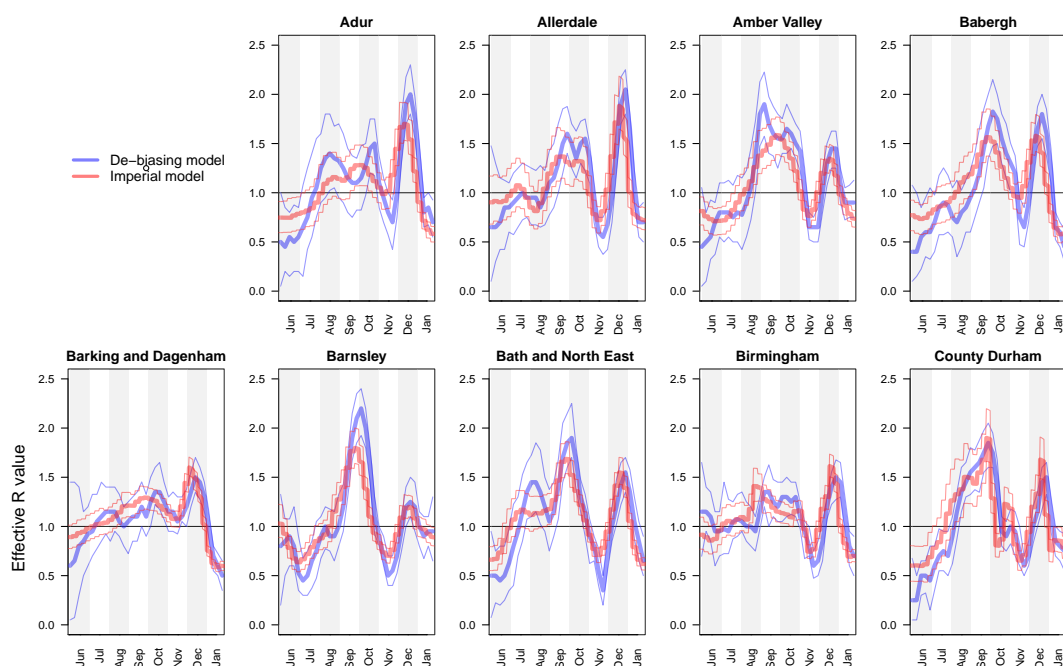


Figure 10: Comparison of \mathcal{R}_t estimates between de-biasing model and Imperial model [32]. For each of the nine PHE regions, we present the constituent LTLA whose name is ranked top alphabetically.

ble intervals overlapping appropriately, despite being based on different data and models.⁶

275 Methods

Observational models for surveillance data

The primary target of inference is prevalence, I out of M , being the unknown number of individuals infected at a particular time-point in the local population of known size M . Our method estimates two types of prevalence: 1) the number of individuals that would test PCR positive (\tilde{I}), and
 280 2) the number of individuals that are infectious (I); see Methods—*Focusing prevalence on the infectious subpopulation*. We clarify below the distinction between the PCR positive and infectious subpopulations, and how we target the latter.

Randomised surveillance data, u of U .

Suppose that out of a total U randomised surveillance (e.g. REACT, ONS CIS) tests, we observe
 285 u positive tests. The randomised testing (e.g. REACT, ONS CIS) likelihood is

$$\mathbb{P}(u \text{ of } U \mid \tilde{I}) = \text{HyperGeom}(u \mid M, \tilde{I}, U), \quad (2)$$

⁶The imperial model uses daily cases data, weekly deaths data, as well as daily infections from the ONS CIS and REACT data sets

and this allows direct, accurate statistical inference on \tilde{I} , the proportion of the population that would return a positive PCR test.

Focusing prevalence on the infectious subpopulation

PCR tests are sensitive, and can detect the presence of SARS-CoV-2 both days before and weeks after an individual is infectious. It is usually desirable for prevalence to represent the proportion of a population that is infectious. We can obtain a likelihood for the number of infectious individuals I as follows:

$$\mathbb{P}(u \text{ of } U \mid I) = \int_{\tilde{I}} \mathbb{P}(u \text{ of } U \mid \tilde{I}) \mathbb{P}(\tilde{I} \mid I) d\tilde{I} \quad (3)$$

where I and \tilde{I} are the number of infectious and PCR-positive individuals respectively.

290 The conditional distribution $\mathbb{P}(\tilde{I} \mid I)$ can be specified on the basis of external knowledge of the average length of time spent PCR-positive vs infectious. Our approach to estimating this quantity imports information on the timing of COVID-19 transmission [33] and the interval of PCR positivity in SARS-CoV-2 infected individuals [34]. More precisely, we specify the infectious time interval for an average infected individual in the population to span the interval 1 to 11 days
 295 post infection (the empirical range of generation time from Figure 1A of [33]). We then calculate the posterior probability of a positive PCR occurring 1 to 11 days post-infection (Figure 1A of [34]). We incorporate the effects of changing incidence in the calculations; this is important because, for example, if incidence is rising steeply, the majority of people who would test PCR positive in the population are those that are relatively recently infected. Full details can be found
 300 in Supplementary Information–*PCR positive to infectious mapping – method details*

Targeted surveillance data, n of N .

In contrast to the randomised surveillance likelihood at (2), the targeted likelihood can be expressed in terms of the observation of n of N positive targeted (e.g. Pillar 1+2) tests as follows:

$$\begin{aligned} \mathbb{P}(n \text{ of } N \mid I, \delta, \nu) &= \text{Binomial}(n \mid I, \mathbb{P}(\text{Tested} \mid \text{Infected})) \\ &\times \text{Binomial}(N - n \mid M - I, \mathbb{P}(\text{Tested} \mid \text{Not Infected})) \end{aligned} \quad (4)$$

where $\mathbb{P}(\text{Tested} \mid \text{Infected})$ and $\mathbb{P}(\text{Tested} \mid \text{Not Infected})$ are the probabilities of an infected (re-
 305 spectively non-infected) individual being tested on date t .

Bias parameters, δ and ν .

We introduce the following parameters:

$$\delta := \log \left(\frac{\text{Odds}(\text{Tested} \mid \text{Infected})}{\text{Odds}(\text{Tested} \mid \text{Not Infected})} \right) \quad (5)$$

$$\nu := \log \text{Odds}(\text{Tested} \mid \text{Not Infected}) , \quad (6)$$

leading to the targeted swab testing likelihood being represented as

$$\begin{aligned} \mathbb{P}(n \text{ of } N \mid I, \delta, \nu) &= \text{Binomial}(n \mid I, \text{logit}^{-1}(\delta + \nu)) \\ &\times \text{Binomial}(N - n \mid M - I, \text{logit}^{-1}\nu). \end{aligned} \quad (7)$$

The unknown parameter requiring special care to infer is δ , i.e. the log odds ratio of being tested in the infected versus the non-infected subpopulations. The other parameter, ν , is directly estimable from the targeted data: $\hat{\nu} := \text{logit}[(N - n)/M]$ is a precise estimator with little bias when prevalence is low.

Test sensitivity and specificity.

The likelihood at (7) assumes a perfect antigen test. If the test procedure has false-positive rate α , and false-negative rate β , the targeted likelihood is instead

$$\mathbb{P}(n \text{ of } N \mid I, \delta, \nu) = \sum_{z=0}^{\min\{I, N\}} \mathbb{P}(z \text{ of } N \mid I, \delta, \nu) \mathbb{P}(n \mid z \text{ of } N), \quad (8)$$

where z denotes the unknown number of truly infected individuals that were tested. The first term in the sum at (8) is obtained by substituting z in (7), while the second term is

$$\mathbb{P}(n \mid z \text{ of } N) = \sum_{n_\beta = \max\{0, z - n\}}^{\min\{z, N - n\}} \text{Binomial}(n_\beta \mid z, \beta) \text{Binomial}(n_\beta + n - z \mid N - z, \alpha), \quad (9)$$

with n_β denoting the number of false-negative test results. An analogous adjustment can be made to the randomised surveillance likelihood at (2).

Cross-sectional inference on local prevalence

We leverage spatially coarse-scale randomised surveillance data to specify an EB prior on bias parameters $p(\delta)$ at coarse-scale (PHE region), and thereby infer prevalence accurately from targeted data at fine scale (LTLA j within PHE region J_j). We explicitly use the superscripts LTLA (j) in PHE region (J_j) in step 4 below where notation from both coarse and fine scale appear together. All quantities in steps 1-3 are implicitly superscripted (J_j) but these are suppressed for notational clarity. For computational efficiency we handle prevalence in a reduced-dimension space of bins as described in Supplementary Information (SI) section *Interval-based prevalence inference – set-up and assumptions*. The method in detail is as follows:

1. **Infer prevalence from unbiased testing data.** At a coarse geographic level (PHE region J_j), estimate prevalence from randomised surveillance data u_t of U_t . Represent the posterior at time t in mass function

$$\hat{p}_t(I_t) := \mathbb{P}(I_t \mid u_t \text{ of } U_t) \quad (10)$$

where $\hat{p}_t : \{0, \dots, M\} \rightarrow [0, 1]$ need only be available at a subset $t \in \mathcal{T} \subseteq \{1, \dots, T\}$ of time points.

2. **Learn δ_t from accurate prevalence.** At a coarse geographic level, for each $t \in \mathcal{T}$, we estimate bias parameter δ_t by coupling biased data n_t of N_t with accurate prevalence information \hat{p}_t . With ν_t fixed at $\hat{\nu}_t := \text{logit}[(N_t - n_t)/M]$

$$p(\delta_t | n_t \text{ of } N_t, \hat{p}_t, \hat{\nu}_t) = \sum_{I_t} p(\delta_t | n_t \text{ of } N_t, I_t, \hat{\nu}_t) \hat{p}_t(I_t) \quad (11)$$

$$\approx N(\delta_t | \hat{\mu}_t, \hat{\sigma}_t^2) \quad (12)$$

where a moment-matched Gaussian approximation is performed at (12). The posterior density in the sum at (11), $p(\delta_t | n_t \text{ of } N_t, I_t, \hat{\nu}_t)$ is conjugate under a $\text{Beta}(a, b)$ prior on $\text{logit}^{-1}(\nu_t + \delta_t) \equiv \mathbb{P}(\text{Tested} | \text{Infected})$, and so can be evaluated as

$$\mathbb{P}(\delta_t \leq \text{logit}(x) - \hat{\nu}_t | n_t \text{ of } N_t, I_t, \hat{\nu}_t) = \text{BetaCDF}(x | n_t + a, I_t - n_t + b). \quad (13)$$

3. **Specify smooth EB prior on $\delta_{1:T}$.** A smooth prior on $\delta_{1:T}$ is specified as follows

$$p(\boldsymbol{\delta}) \propto N(\boldsymbol{\delta} | \mathbf{0}, \boldsymbol{\Sigma}_\delta) \prod_{t \in \mathcal{T}} N(\delta_t | \hat{\mu}_t, \hat{\sigma}_t^2) \prod_{t \notin \mathcal{T}} N(\delta_t | 0, \sigma_{\text{flat}}^2) \quad (14)$$

where $N(\boldsymbol{\delta} | \mathbf{0}, \boldsymbol{\Sigma}_\delta)$ imparts a user-specified degree of longitudinal smoothness, thereby sharing information on δ across time points. Ignorance of δ_t , in the absence of random surveillance data, is encapsulated in a Gaussian with large variance σ_{flat}^2 . A standard choice for $N(\boldsymbol{\delta} | \mathbf{0}, \boldsymbol{\Sigma}_\delta)$ corresponds to a stationary autoregressive, AR(1), process of the form

$$\delta_t = c + \psi \delta_{t-1} + \varepsilon_t \quad (15)$$

with a diffuse Gaussian prior $c \sim N(0, \sigma_{\text{flat}}^2)$ and with smoothing tuned by $0 < \psi < 1$ and white noise variance σ_ε^2 . The normalised form of the prior at (14) is

$$p(\boldsymbol{\delta}) = N\left(\boldsymbol{\delta} \mid (\boldsymbol{\Sigma}_\delta^{-1} + \mathbf{D}^{-1})^{-1} \mathbf{D}^{-1} \hat{\boldsymbol{\mu}}, (\boldsymbol{\Sigma}_\delta^{-1} + \mathbf{D}^{-1})^{-1}\right) \quad (16)$$

with $(\hat{\boldsymbol{\mu}}, \text{diagonal matrix } \mathbf{D}_{T \times T})$ having elements $(\hat{\mu}_t, \hat{\sigma}_t^2)$ for $t \in \mathcal{T}$ and $(0, \sigma_{\text{flat}}^2)$ for $t \notin \mathcal{T}$.

4. **Infer cross-sectional local prevalence from biased testing data.** At a fine-scale geographic level (LTLA j in PHE region J_j), having observed $n_t^{(j)}$ of $N_t^{(j)}$ positive test results (a subset of the $n_t^{(J_j)}$ of $N_t^{(J_j)}$ observed at the coarse-scale level above), calculate the posterior for $I_t^{(j)}$ separately at each time point t :

$$p(I_t^{(j)} | n_t^{(j)} \text{ of } N_t^{(j)}) \propto p(I_t^{(j)}) p(n_t^{(j)} \text{ of } N_t^{(j)} | I_t^{(j)}, \hat{\nu}_t^{(j)}) \quad (17)$$

$$= p(I_t^{(j)}) \int_{\delta_t^{(j)}} p(n_t^{(j)} \text{ of } N_t^{(j)} | I_t^{(j)}, \hat{\nu}_t^{(j)}, \delta_t^{(j)}) p(\delta_t^{(j)}) d\delta_t^{(j)} \quad (18)$$

where $\hat{\nu}_t^{(j)} := \text{logit}[(N_t^{(j)} - n_t^{(j)})/M_t^{(j)}]$, the likelihood in the integral at (18) is available at (7), and the prior $p(\delta_t^{(j)})$ is time-point t 's marginal Gaussian from (16).

Debiasing lateral flow device (LFD) tests with PCR surveillance (or vice versa)

The methods can be adapted straightforwardly to the situation in which the randomised surveillance study uses a different assay to the targeted testing. For a concrete example we could use

REACT PCR prevalence posterior $\hat{p}_t(\tilde{I}_t)$ from (10) to debias Pillar 1+2 LFD test data n_t of N_t . Equation (11) can be adjusted to estimate ascertainment bias δ pertaining to LFD data as follows:

$$p(\delta_t | n_t \text{ of } N_t, \hat{p}_t, \hat{\nu}_t) = \sum_{\bar{I}_t} \left\{ p(\delta_t | n_t \text{ of } N_t, \bar{I}_t, \hat{\nu}_t) \sum_{\tilde{I}_t} \mathbb{P}(\bar{I}_t | \tilde{I}_t) \hat{p}_t(\tilde{I}_t) \right\}, \quad (19)$$

where \bar{I}_t and \tilde{I}_t are the unobserved LFD- and PCR-positive prevalence respectively, and the conditional distribution $\mathbb{P}(\bar{I}_t | \tilde{I}_t)$ can be estimated on the basis of external knowledge of the average length of time spent PCR-positive vs LFD-positive, analogously to as described in Methods—*Focusing prevalence on the infectious subpopulation*. The remaining computations, from (12) onwards, are unchanged, with the outputted fine-scale marginal likelihood $p(n_t^{(j)} \text{ of } N_t^{(j)} | I_t^{(j)}, \hat{\nu}_t^{(j)})$ at (17) to be interpreted as targeting the local LFD-positive prevalence $\bar{I}_t^{(j)}$.

365 Full Bayesian inference under a stochastic SIR epidemic model

The cross-sectional analysis described in *Cross-sectional inference on local prevalence* generates the δ -marginalised likelihood, $p(n_t^{(j)} \text{ of } N_t^{(j)} | I_t^{(j)}, \hat{\nu}_t)$ at (17), at each time point for which targeted data are available. These likelihoods can be used as input for longitudinal models to obtain better prevalence estimates and to infer epidemiological parameters such as \mathcal{R}_t .

370 We illustrate this via a Bayesian implementation of a stochastic epidemic model whereby individuals become immune through population vaccination and/or exposure to COVID-19 (Figure 11). We incorporate known population vaccination counts into a standard discrete time Markov chain (DTMC) SIR model ([35], Chapter 3). Details of the transition probability calculations are given in SI section *SIR model details*, and assumptions in Supplementary Information—*SIR model*
375 – *discussion, assumptions and caveats*.

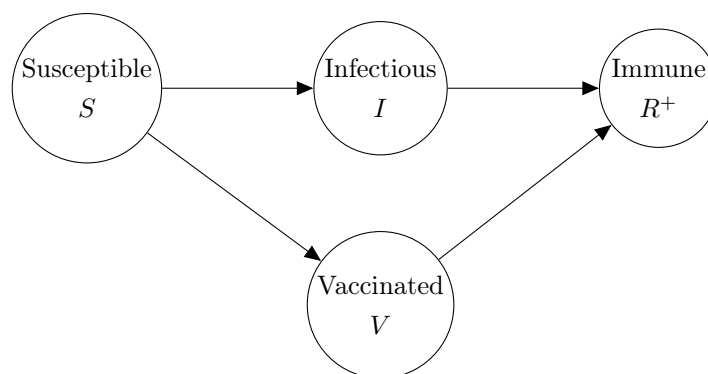


Figure 11: SIR/V epidemic model compartmental diagram.

Priors on \mathcal{R}, I, R^+

We place priors on I, R^+ measured as a proportion of the population; this proportion then gets mapped to prevalence intervals on subpopulation counts as described in *Interval-based prevalence*

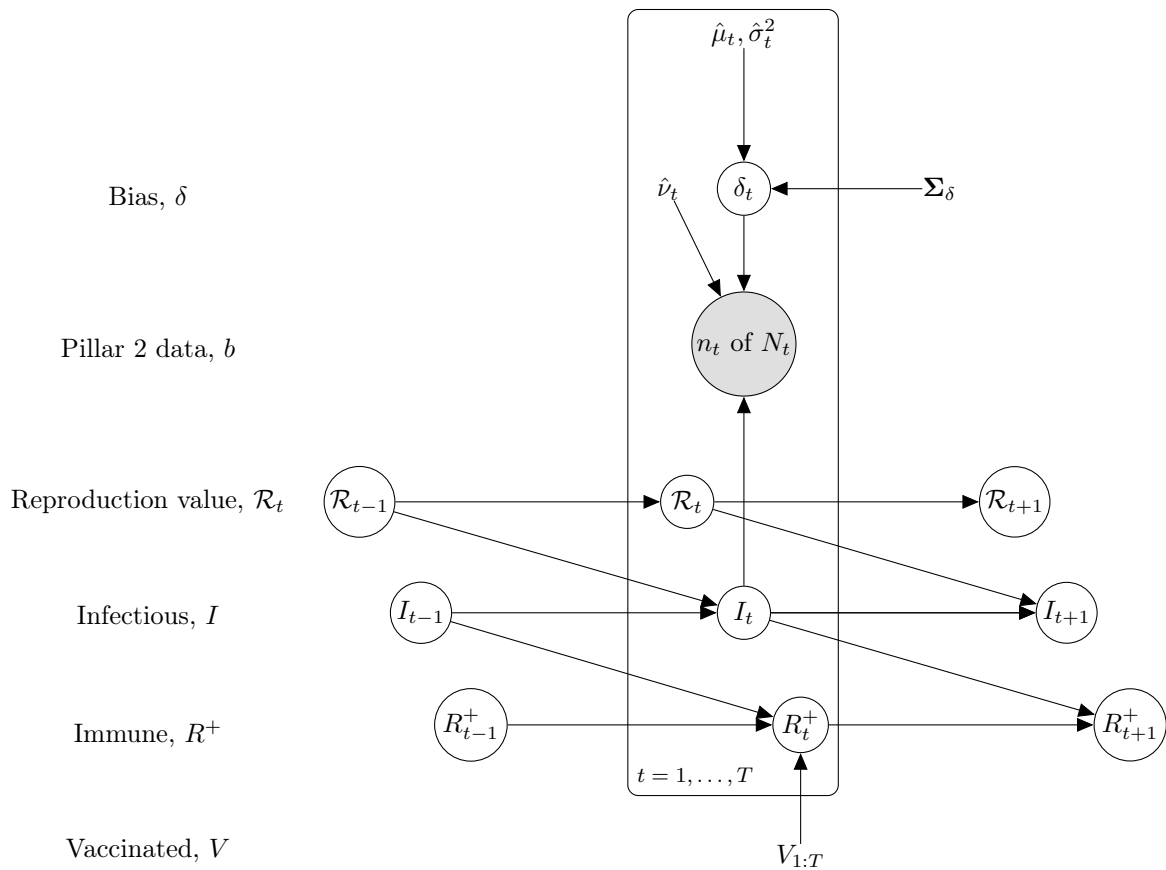


Figure 12: Longitudinal model DAG for SIR epidemic model at local level (e.g. LTLA). Directed paths characterise conditional probability distributions, in contrast to the paths showing transitions between model compartments in Figure 11. Inference is for a region, e.g. an LTLA, based only on targeted test data collected in this region, n_t of N_t . A prior on δ_t parameterized $(\hat{\mu}_t, \hat{\sigma}_t^2)$ brings information on the Pillar 2 ascertainment bias learned from randomized surveillance testing data available for the PHE region in which the LTLA lies. The $T \times T$ covariance matrix Σ_δ imparts temporal smoothness on $\delta_{1:T}$. Effective reproduction numbers are denoted $\mathcal{R}_{1:T}$, number of infectious individuals by $I_{1:T}$, and the number of immune individuals by $R_{1:T}^+$.

inference – set-up and assumptions. Specifically, we use truncated, discretized Gaussian distributions on the proportion of the population immune and infectious. For example, on number of infectious individuals I_t at each timepoint t , we specify the prior (suitably normalized over its support)

$$\mathbb{P}(I_t = j) \propto \int_{(j-1)/M}^{j/M} \mathcal{N}(x | \mu_I, \sigma_I^2) dx \quad \text{for } j/M \in [p_{\min}, \dots, p_{\max}] , \quad (20)$$

with an example weakly informative hyperparameter setting being $\mu_I = 0.5\%$, $\sigma_I = 1\%$, $p_{\min} = 0\%$, $p_{\max} = 4\%$. To ensure meaningful inference on $R_{1:T}^+$, we place an informative prior that reflects the state of knowledge of the immune population size; we do this using an informative truncated Gaussian prior on R_1^+ , and non-informative priors on $R_{2:T}^+$. We place a noninformative uniform prior on each \mathcal{R}_t , e.g. a $\text{Uniform}(0.5, 2.5)$.

MCMC sampling implementation

We perform inference under the model represented in the DAG at Figure 11. The likelihood is marginalised with respect to δ , and we use Markov chain Monte Carlo (MCMC) to draw samples from the posterior

$$p(\mathbf{I}, \mathbf{R}^+, \mathcal{R} | \mathbf{n}, \mathbf{N}) .$$

We sample \mathcal{R} and $(\mathbf{I}, \mathbf{R}^+)$ using separate Gibbs updates. For sampling $(\mathbf{I}, \mathbf{R}^+)$ we represent the joint full conditional as

$$p(\mathbf{I}, \mathbf{R}^+ | \mathcal{R}, \mathbf{n}, \mathbf{N}) = p(\mathbf{I} | \mathcal{R}, \mathbf{n}, \mathbf{N})p(\mathbf{R}^+ | \mathbf{I}) , \quad (21)$$

sampling \mathbf{I}^{new} from $p(\mathbf{I} | \mathcal{R}, \mathbf{n}, \mathbf{N})$, and then $\mathbf{R}^{+\text{new}}$ from $p(\mathbf{R}^+ | \mathbf{I}^{\text{new}})$.

Sampling from $p(\mathbf{I} | \mathcal{R}, \mathbf{n}, \mathbf{N})$

The sampling distribution on prevalence can be expressed:

$$\begin{aligned} p(\mathbf{I} | \mathcal{R}, \mathbf{n}, \mathbf{N}) &\propto p(\mathbf{n}, \mathbf{N} | \mathbf{I}, \mathcal{R})p(\mathbf{I} | \mathcal{R}) \\ &= p(n_1, N_1 | I_1)p(I_1) \prod_{t=2}^T p(n_t \text{ of } N_t | I_t)p(I_t | I_{t-1}, \mathcal{R}_{t-1}), \end{aligned} \quad (22)$$

which is an HMM with emission probabilities taken from the δ -marginalised likelihood at (18), and transition probabilities taken from (37).

Sampling from $p(\mathbf{R}^+ | \mathbf{I})$

We can express the full conditional for $\Delta R_{1:T}^+$ as

$$\mathbb{P}(R_{1:T}^+ | I_{1:T}) \propto \mathbb{P}(R_1^+ | V_1) \prod_{t=2}^T \mathbb{P}(R_t^+ | R_{t-1}^+, I_{t-1}, \Delta V_t)$$

and sample the $\Delta R_{1:T}^+$ sequentially, with $\mathbb{P}(R_t^+ | R_{t-1}^+, I_{t-1}, \Delta V_t)$ available at (39).

Sampling from $p(\mathcal{R} | \mathbf{I})$

The prior joint distribution of $\mathcal{R}_{1:T}$ is modelled using a random walk:

$$\mathcal{R}_t \sim \text{Normal}(\mathcal{R}_{t-1}, \sigma_{\mathcal{R}}^2), \quad (23)$$

where $\sigma_{\mathcal{R}}^2$ is a user-specified smoothness parameter.

395 The update involves sampling from

$$p(\mathcal{R} | \mathbf{I}) = p(\mathcal{R}_1) \prod_{t=2}^{T-1} p(\mathcal{R}_t | \mathcal{R}_{t-1}) \prod_{t=2}^T p(I_t | I_{t-1}, \mathcal{R}_{t-1}). \quad (24)$$

We discretize the space of \mathcal{R}_t into an evenly spaced grid and sample from the HMM defined at (24) [36]. The transition probabilities are given by (23) (suitably normalised over the discrete \mathcal{R}_t space) and the emission probabilities given by (37).

Data

400 With the exception of the variant of concern 202012/01 analysis, all data underlying the results presented here are publicly available. Randomised surveillance data comes from the REACT study [7].⁷ From REACT, we create weekly test counts at the spatially coarse-scale level (PHE region) and, for validation purposes but not model fitting, use round-aggregated counts at the fine-scale level (LTLA), for round 7 (13th Nov - 3rd Dec 2021) and round 8 (6th-22nd Jan 2021). The
405 combined weekly Pillar 1 and Pillar 2 data are publicly available for download.⁸

Analysis scripts

The R scripts used to generate the results in this manuscript are available at <https://github.com/alan-turing-institute/jbc-turing-rss-testdebiasing>.

Discussion

410 We have introduced and applied an integrative causal model allowing accurate inference from community testing data. The flexible probabilistic framework allows simultaneous and coherent incorporation of a number of important features, including:

- Adjustment for ascertainment bias caused by preferential testing based on symptom status, or on other confounders.
- 415 • Allowing for heterogeneous testing capacity by modelling the total number of tests conducted locally.
- Incorporation of multiple different SARS-CoV-2 testing assays, such a LFD and PCR, including adjustment for particular sensitivity/specificity.

⁷<https://github.com/mrc-ide/reactidd/tree/master/inst/extdata>

⁸<https://www.gov.uk/government/publications/nhs-test-and-trace-england-statistics-14-january-to-20-january-2021>. Note that lateral flow test results are not included in the these weekly summaries.

- Inference on the number of infectious individuals, when PCR tests pick up individuals at non-infectious stages.
- The model outputs week-specific debiased prevalence with uncertainty (via a marginal likelihood), that can be incorporated modularly into more complex models.
- An SIR epidemic model implementation allowing estimation of \mathcal{R}_t and adjustment for vaccination in the immune population

Because of the extensive Pillar 1+2 testing effort, there is a large amount of information contained in these targeted data at LTLA level, even for a single weekly timepoint, as shown by the narrow width of CIs in Figure 4. However, due to the strength of this information, the targeted data also have the potential to introduce bias into more complex models if they are incorporated as observed nodes without an appropriate ascertainment correction. Equipped with a well specified prior on δ , however, precise and accurate estimation at even finer scales may be feasible. Given the high volume of data, despite focusing on high resolution geographical units such as LTLAs, estimates do not seem to be affected by the classical issue of small area estimation, that is to say sample sizes too small to carry on inference without borrowing strength from neighbouring units. If interest is focused on ultra fine-scale geography and/or when prevalence is much lower, then spatial borrowing would be beneficial, with such additional smoothing subject to a variance-bias trade-off.

In addition to estimating at finer spatial scales, it may be desirable to estimate local prevalence from targeted data stratified by factors such as age and ethnicity. This would involve modelling the relevant confounders in Figure 1, rather than marginalising them out. One approach to doing this within the existing framework would be to perform a stratified analysis (e.g. performing the whole analysis using weekly PHE region and Pillar 1+2 Data stratified by age bands). A more sophisticated approach could model δ_t semi-parametrically with some spatiotemporal smoothness assumptions on the effect of age and other confounders. Of course, any approach would require appropriate metadata on any confounders that affect randomised surveillance and/or targeted sampling.

For cross-sectional prevalence estimation, a key dependency is the availability of a regular, up-to-date stream of randomised surveillance data at some level of spatial resolution. Here we deployed REACT data at the coarse PHE-region scale. The UK has led the way internationally in having regular national surveillance randomised surveys like REACT and ONS. This modelling work shows the importance of having both targeted testing and also a rolling randomised surveillance survey to be able to better track the epidemic. The methods are transferable beyond the UK wherever randomized testing data are being gathered. This could be built in an integrated way from the start as preparedness for pandemics, in particular for diseases where asymptomatic transmission plays an important role.

455 Funding

BL was supported by the UK Engineering and Physical Sciences Research Council through the Bayes4Health programme [Grant number EP/R018561/1] and gratefully acknowledges funding from Jesus College, Oxford. KBP is supported by the National Institute for Health Research Health Protection Research Unit (NIHR HPRU) in Healthcare Associated Infections and Antimicrobial Resistance at the University of Oxford in partnership with Public Health England (PHE) (NIHR200915). SR is supported by MRC programme grant MC_UU_00002/10; The Alan Turing Institute grant: TU/B/000092; EPSRC Bayes4Health programme grant: EP/R018561/1. MB acknowledges partial support from the MRC Centre for Environment and Health, which is currently funded by the Medical Research Council (MR/S019669/1). Infrastructure support for the Department of Epidemiology and Biostatistics is also provided by the NIHR Imperial BRC. The computational aspects of this research were supported by the Wellcome Trust Core Award Grant Number 203141/Z/16/Z and the NIHR Oxford BRC. The views expressed are those of the authors and not necessarily those of the National Health Service, NIHR, Department of Health, or PHE.

References

- 470 [1] COVID-19 EpiCell. PHE data series on deaths in people with COVID-19: technical summary - 12 august update (2020).
- [2] The official UK Government website for data and insights on Coronavirus (COVID-19). <https://coronavirus.data.gov.uk> (2021). [Online; accessed 15-February-2021].
- [3] Scientific Advisory Group for Emergencies. Summary of effectiveness and harms of npi (2020). URL <https://www.gov.uk/government/publications/summary-of-the-effectiveness-and-harms-of-different-non-pharmaceutical-interventions-16-september-2020>
- 475 [4] Prime Minister's Office, D. S. Prime Minister announces new local COVID Alert Levels. GOV. <https://www.gov.uk/government/news/prime-minister-announces-new-local-covid-alert-levels> (2020).
- 480 [5] UK Cabinet Office. COVID-19 Response - Spring 2021 (Summary). URL <https://www.gov.uk/government/publications/covid-19-response-spring-2021/covid-19-response-spring-2021-summary>.
- [6] Pouwels, K. B. et al. Community prevalence of SARS-CoV-2 in England from April to November, 2020: results from the ONS Coronavirus Infection Survey. *The Lancet Public Health* **6**, e30–e38 (2021).
- 485 [7] Riley, S. et al. Community prevalence of sars-cov-2 virus in england during may 2020: React study. *medRxiv* (2020).

- [8] Department of Health and Social Care (UK), COVID-19 testing data: methodology note. URL <https://www.gov.uk/government/publications/coronavirus-covid-19-testing-data-methodology/covid-19-testing-data-methodology-note>.
490
- [9] Byambasuren, O. *et al.* Estimating the extent of asymptomatic COVID-19 and its potential for community transmission: Systematic review and meta-analysis. Official Journal of the Association of Medical Microbiology and Infectious Disease Canada **5**, 223–234 (2020). URL <https://jammi.utpjournals.press/doi/10.3138/jammi-2020-0030>. Publisher: University of Toronto Press.
495
- [10] Subramanian, R., He, Q. & Pascual, M. Quantifying asymptomatic infection and transmission of COVID-19 in New York City using observed cases, serology, and testing capacity. Proceedings of the National Academy of Sciences **118** (2021). URL <https://www.pnas.org/content/118/9/e2019716118>. Publisher: National Academy of Sciences Section: Biological Sciences.
500
- [11] Manzi, G., Spiegelhalter, D. J., Turner, R. M., Flowers, J. & Thompson, S. G. Modelling bias in combining small area prevalence estimates from multiple surveys. Journal of the Royal Statistical Society: Series A (Statistics in Society) **174**, 31–50 (2011).
- [12] Giorgi, E., Sesay, S. S., Terlouw, D. J. & Diggle, P. J. Combining data from multiple spatially referenced prevalence surveys using generalized linear geostatistical models. Journal of the Royal Statistical Society. Series A (Statistics in Society) 445–464 (2015).
505
- [13] Amoah, B., Diggle, P. J. & Giorgi, E. A geostatistical framework for combining spatially referenced disease prevalence data from multiple diagnostics. Biometrics **76**, 158–170 (2020).
- [14] Crainiceanu, C. M., Diggle, P. J. & Rowlingson, B. Bivariate binomial spatial modeling of loa loa prevalence in tropical africa. Journal of the American Statistical Association **103**, 21–37 (2008).
510
- [15] Heesterbeek, H. *et al.* Modeling infectious disease dynamics in the complex landscape of global health. Science **347** (2015).
- [16] Birrell, P. J., De Angelis, D. & Presanis, A. M. Evidence synthesis for stochastic epidemic models. Statistical science: a review journal of the Institute of Mathematical Statistics **33**, 34 (2018).
515
- [17] De Angelis, D., Presanis, A. M., Birrell, P. J., Tomba, G. S. & House, T. Four key challenges in infectious disease modelling using data from multiple sources. Epidemics **10**, 83–87 (2015).
- [18] Shubin, M., Lebedev, A., Lyytikäinen, O. & Auranen, K. Revealing the true incidence of pandemic a (h1n1) pdm09 influenza in finland during the first two seasons—an analysis based on a dynamic transmission model. PLoS computational biology **12**, e1004803 (2016).
520

- [19] Funk, S. *et al.* Short-term forecasts to inform the response to the Covid-19 epidemic in the UK. *medRxiv* 2020.11.11.20220962 (2020). URL <http://medrxiv.org/content/early/2020/12/04/2020.11.11.20220962.abstract>.
525
- [20] Abbott, S. *et al.* Estimating the time-varying reproduction number of sars-cov-2 using national and subnational case counts. *Wellcome Open Research* **5**, 112 (2020).
- [21] Cori, A., Ferguson, N. M., Fraser, C. & Cauchemez, S. A new framework and software to estimate time-varying reproduction numbers during epidemics. *American journal of epidemiology* **178**, 1505–1512 (2013).
530
- [22] Birrell, P. J. *et al.* Real-time nowcasting and forecasting of covid-19 dynamics in england: the first wave? *medRxiv* (2020).
- [23] Colman, E., Enright, J., Puspitarani, G. A. & Kao, R. R. Estimating the proportion of SARS-CoV-2 infections reported through diagnostic testing. *medRxiv* 2021.02.09.21251411 (2021). URL <https://www.medrxiv.org/content/10.1101/2021.02.09.21251411v1>. Publisher: Cold Spring Harbor Laboratory Press.
535
- [24] Campbell, H. *et al.* Bayesian adjustment for preferential testing in estimating the covid-19 infection fatality rate: Theory and methods. *arXiv preprint arXiv:2005.08459* (2020).
- [25] Brazeau, N. *et al.* Report 34: Covid-19 infection fatality ratio: estimates from seroprevalence (2020). URL <https://doi.org/10.25561/83545>.
540
- [26] Mishra, S. *et al.* A COVID-19 Model for Local Authorities of the United Kingdom. *medRxiv* 2020.11.24.20236661 (2020). URL <https://www.medrxiv.org/content/10.1101/2020.11.24.20236661v1>.
- [27] Scott, J. A. *et al.* *epidemia: Modeling of Epidemics using Hierarchical Bayesian Models* (2020). URL <https://imperialcollegelondon.github.io/epidemia/>.
545
- [28] Flaxman, S. *et al.* Estimating the effects of non-pharmaceutical interventions on covid-19 in europe. *Nature* **584**, 257–261 (2020).
- [29] Jewell, C., Read, J., Roberts, G., Rowlington, B. & Suter, C. Bayesian stochastic model-based forecasting fro spatial Covid-19 risk in England. Technical Concept Note (2020). URL https://github.com/chrism0dwb/covid19uk/blob/master/doc/lancs_space_model_concept.pdf.
550
- [30] Davies, N. G. *et al.* Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England. *Science* eabg3055 (2021). URL <http://science.sciencemag.org/content/early/2021/03/03/science.abg3055.abstract>.
- [31] Investigation of novel SARS-COV-2 variant: Variant of Concern 202012/01. Tech. Rep., Public Health England (2020). URL www.gov.uk/government/publications/investigation-of-novel-sars-cov-2-variant-variant-of-concern-20201201.
555

- [32] Mishra, S. *et al.* A covid-19 model for local authorities of the united kingdom. *medRxiv* (2020).
- 560 [33] Ferretti, L. *et al.* The timing of COVID-19 transmission. *medRxiv* 2020.09.04.20188516 (2020).
URL <http://medrxiv.org/content/early/2020/09/16/2020.09.04.20188516.abstract>.
- [34] Hellewell, J. *et al.* Estimating the effectiveness of routine asymptomatic PCR testing at different frequencies for the detection of SARS-CoV-2 infections. *medRxiv* 2020.11.24.20229948 (2020). URL <https://www.medrxiv.org/content/10.1101/2020.11.24.20229948v1>. Publisher: Cold Spring Harbor Laboratory Press.
- 565 [35] Brauer, Fred, van den Driessche, Pauline & Wu, J. *Mathematical Epidemiology*. Mathematical Biosciences Subseries (Springer-Verlag Berlin Heidelberg, 2008).
- [36] Scott, S. L. Bayesian methods for hidden Markov models: Recursive computing in the 21st century. *Journal of the American Statistical Association* **97**, 337–351 (2002). URL <https://www.tandfonline.com/doi/abs/10.1198/016214502753479464>.
- 570 [37] Office of National Statistics. Covid-19 infection survey: methods and further information. <https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/methodologies/covid19infectionsurveypilotmethodsandfurtherinformation>.
- 575 [38] Pouwels, K. B. *et al.* Community prevalence of sars-cov-2 in england from april to november, 2020: results from the ons coronavirus infection survey. *The Lancet Public Health* **6**, e30–e38 (2021). URL <https://www.sciencedirect.com/science/article/pii/S2468266720302826>.
- [39] Department of Health and Social Care. Lateral flow device specificity in phase 4 (post-marketing) surveillance. <https://www.gov.uk/government/publications/lateral-flow-device-specificity-in-phase-4-post-marketing-surveillance>.
- 580 [40] Office of National Statistics. Coronavirus (covid-19) infection survey: antibody data for the uk, january 2021. <https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/articles/coronaviruscovid19infectionsinthecommunityinengland/antibodydatafortheukjanuary2021>.
- 585 [41] Overton, C. E. *et al.* Using statistics and mathematical modelling to understand infectious disease outbreaks: COVID-19 as an example. *Infectious Disease Modelling* **5**, 409–441 (2020). URL <https://www.sciencedirect.com/science/article/pii/S2468042720300245>.
- 590 [42] Keeling, M. J. *et al.* Predictions of COVID-19 dynamics in the UK: Short-term forecasting and analysis of potential exit strategies. *PLOS Computational Biology* **17**, e1008619

- (2021). URL <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1008619>. Publisher: Public Library of Science.
- 595 [43] Keeling, M. J. *et al.* Fitting to the UK COVID-19 outbreak, short-term forecasts and estimating the reproductive number. *medRxiv* 2020.08.04.20163782 (2020). URL <http://medrxiv.org/content/early/2020/09/29/2020.08.04.20163782.abstract>.
- [44] Brown, G. D., Porter, A. T., Oleson, J. J. & Hinman, J. A. Approximate Bayesian computation for spatial SEIR(S) epidemic models. *Spatial and Spatio-Temporal Epidemiology* **24**, 27–37 (2018).
- 600 [45] Schader, M. & Schmid, F. Two Rules of Thumb for the Approximation of the Binomial Distribution by the Normal Distribution. *The American Statistician* (2012). URL <https://www.tandfonline.com/doi/abs/10.1080/00031305.1989.10475601>. Publisher: Taylor & Francis Group.

Supplementary Information

605 Model parameters

A full list of model parameters, along with either their prior distribution or the value at which they were fixed, can be found in Table S1.

Table 1: Model parameters with specified prior distributions or fixed values

Parameter	Prior / Fixed value
Ascertainment bias, $\delta_{1:T}$	Empirical Bayes prior (see Eq. (16)): <ul style="list-style-type: none"> - AR(1) coefficient, $\psi = 0.99$ - Standard deviation, $\sigma_\epsilon = 1$ - Intercept, $c \sim \mathcal{N}(0, \sigma_{flat}^2)$ with $\sigma_{flat} = 10$
PCR false-positive rate, α	Fixed, $\alpha = 0.001$, taken from [37, 38]
PCR false-negative rate, β ,	Fixed, $\beta = 0.05$, taken from [39]
Expected time to recovery, $1/\gamma$	$T_{recovery} \sim \text{Exponential}(\gamma)$, with $\gamma = 1$ week
Effective reproduction number, \mathcal{R}_t	Random walk: $\mathcal{R}_t \sim \mathcal{N}(\mathcal{R}_{t-1}, \sigma_{\mathcal{R}}^2)$, with $\sigma_{\mathcal{R}}^2 = 0.2$
Proportion immune at $t = 0$, R_0^+/M	Truncated Gaussian (see Eq. (20), and reference [40]) <ul style="list-style-type: none"> - Mean, $\mu_R = 0.06$ - Standard deviation, $\sigma_R = 0.01$ - Minimum proportion, $p_{min} = 0$ - Maximum proportion, $p_{max} = 0.1$
Proportion infectious at each t , I_t/M	Truncated Gaussian (see Eq. (20)): <ul style="list-style-type: none"> - Mean, $\mu_I = 0.005$ - Standard deviation, $\sigma_I = 0.01$ - Minimum proportion, $p_{min} = 0$ - Maximum proportion, $p_{max} = 0.04$

Discussion of methodological assumptions and caveats

Interval-based prevalence inference – set-up and assumptions

610 The full prevalence state space comprises all potential numbers of infectious individuals in the population, i.e. $I \in \{0, \dots, M\}$. For computational tractability we define $B \ll M$ bins:⁹

$$\mathcal{B}_b := \{I : e_{b-1} \leq I < e_b\} \quad b = 1, \dots, B \quad (25)$$

having midpoints:

$$\check{I}_b := \left\lfloor \frac{e_{b-1} + e_b - 1}{2} \right\rfloor, \quad b = 1, \dots, B, \quad (26)$$

⁹Bins are equally sized on log scale, with interval edges are defined recursively as $e_0 = 0$, $e_b = \lceil e_{b-1}(1 + \varepsilon_B) \rceil$, and ε_B is a fixed constant giving B intervals.

and make three assumptions to allow computationally efficient inference on the B -dimensional space of bins, denoting these assumptions Interval-1:3 as follows:

615 Interval-1 The testing data likelihood, conditional on prevalence bin, is evaluated at the bin midpoint:

$$\mathbb{P}(n \text{ of } N, u \text{ of } U \mid I \in \mathcal{B}_b) := \mathbb{P}(n \text{ of } N, u \text{ of } U \mid I = \check{I}_b). \quad (27)$$

Interval-2 Prevalence I is uniformly distributed within each bin:

$$\mathbb{P}(I = k \mid I \in \mathcal{B}_b) := \begin{cases} \frac{1}{e_b - e_{b-1}} & k \in \mathcal{B}_b \\ 0 & \text{otherwise.} \end{cases} \quad (28)$$

Interval-3 The distribution of new infections, conditional on prevalence bin, is evaluated at the bin midpoint (with the same assumption applying to new recoveries):

$$\mathbb{P}(\# \text{ new infections} \mid I \in \mathcal{B}_b) := \mathbb{P}(\# \text{ new infections} \mid I = \check{I}_b) \quad (29)$$

$$\mathbb{P}(\# \text{ new recoveries} \mid I \in \mathcal{B}_b) := \mathbb{P}(\# \text{ new recoveries} \mid I = \check{I}_b). \quad (30)$$

620 Ascertainment bias model – assumptions and caveats

Debias-1 Spatial homogeneity of δ across LTLAs within a PHE region. The fact that we see relatively low variation in δ at each time point across PHE regions in Figure 3, particularly after October 2020, is consistent with a finer-scale spatial homogeneity assumption being reasonable.

625 Debias-2 We handle prevalence in a reduced-dimension space of bins as described in SI section *Interval-based prevalence inference – set-up and assumptions*

Debias-3 (In)stability of ascertainment mechanism. It is clear from Figure 3 that the ascertainment effects captured by δ can change rapidly and without obvious cause over time. Contemporaneous randomised surveillance data, such as REACT or ONS CIS, allow
630 estimation of δ . However, when predicting prevalence forward in time beyond availability of randomised surveillance data, we are making the implicit assumption that the ascertainment bias remains stable forwards in time, and such results should therefore be interpreted with caution.

PCR+ to infectious mapping – assumptions and caveats

635 For full details please see Supplementary Information—*PCR positive to infectious mapping – method details*.

Infectious-1 Pillar 1+2 positive test counts, across a four-week period, are used as an approximation to the true *relative* incidence over that time interval at coarse-scale level (e.g. PHE region).

640 Infectious-2 The probability (with credible intervals) of testing PCR positive when swabbed d days post infection is taken from Figure 1A of Hellewell et al. [34].

Infectious-3 The infectious interval for an average individual is defined to span days 1 to 11 post infection, based on Figure 1A of Ferretti et al. [33].

SIR model – discussion, assumptions and caveats

645 The illustrative epidemic model we implement here has one of the simplest SIR compartmental structures available, as summarised in Supplementary Information–*SIR model – discussion, assumptions and caveats* and particularly Assumption SIR-2. Other teams have developed more realistic and sophisticated compartmental models of transmission, reflecting for example that individuals are not immediately infectious after being infected [22, 41, 42, 43]. Importantly, these
650 are able to relate epidemiological disease dynamics to outcomes far downstream, such as hospitalisation and deaths. The fact that a large number and variety of models has been developed can be viewed as a strength, as demonstrated by efficacy of ensembles of multi-model forecasts to inform policy on future resource needs and population impacts [19]. One attractive feature of such model ensembles is that their forecasts may be relatively robust to changes in spatiotemporal and
655 compartmental dynamics over the course of an epidemic. Notably, the de-biased prevalence likelihood outputted in Results–*Cross-sectional local prevalence from targeted testing data* is agnostic to the downstream epidemic model, and so there might be benefits to incorporating it into such multi-compartment epidemic models.

SIR-1 The population is homogeneous within an LTLA, with each individual equally likely to be
660 infected

SIR-2 We assume individuals become instantly infectious and recover at a fixed rate $\gamma = 7$ days, i.e. with no spatiotemporal variation, and with recovery time distributed exponentially with mean $1/\gamma$.

SIR-3 Any projections forward in time are made under the implicit assumption that there is no
665 change in NPIs, such as tiering or lockdown status, affecting the LTLA.

SIR-4 We do not include age, ethnicity or deprivation indices in our model, and so epidemiological parameter estimates are to be interpreted as an average across these strata (with unknown weights).

SIR-5 We do not explicitly model transmission between regions or the demographic effects of births,
670 deaths and migration – the SIR model is fitted to each LTLA separately. While it would be possible to account for transmission between LTLAs [44], this dramatically increases the number of parameters to be estimated and consequently the computational burden of the model. Given that the study period here is almost all in lockdown, the effect of transmission between LTLAs is relatively small. In non-lockdown periods, epidemic models allowing for
675 inter-region transmission could be beneficial.

SIR-6 The number of new infections in the stochastic SIR model is modelled as a Poisson approximation, approximating the ‘true’ Binomial conditional distribution.

Gaussian approximation for δ

We approximate the cross-sectional component of the EB prior for δ using a moment-matched
680 Gaussian approximation (see (12)). Figure 13 illustrates the suitability of this approximation for PHE regions London and the North West across nine weeks.

SIR model details

We implement a DTMC SIR epidemic model based on the standard model as described in ([35], Chapter 3). As we choose Δt to be a day/week, we allow multiple infections and recoveries in
685 a time interval width Δt ; this requires derivation of Markov transition probabilities between all states (rather than just neighbouring ones), which we do below having established some notation.

Notation

Parameters are subscripted by timepoint index t (indexing week for the analyses presented, with Δt set to one week):

690 I_t : number of infectious individuals

R_t^+ : number of immune individuals (with infection- and/or vaccination-acquired immunity)

V_t : total number of vaccinated individuals in region (i.e. with vaccine-acquired immunity)

S_t : number of susceptible individuals ($S_t \equiv M - R_t^+ - I_t$)

ΔQ_t : number of *new* infections in interval $(t - \Delta t, t]$

695 ΔR_t : number of *new* recoveries in interval $(t - \Delta t, t]$

ΔV_t : number of vaccinations administered in interval $(t - \Delta t, t]$

$\Delta \tilde{V}_t$: number of vaccinations administered to susceptible individuals in interval $(t - \Delta t, t]$

β_t : transmission rate, i.e. the number of effective contacts in interval $(t - \Delta t, t]$

γ : recovery rate, with expected time to recovery $\mathbb{E}[T] = 1/\gamma$

700 γ_t : probability of recovery in interval $(t - \Delta t, t]$, i.e. $\gamma_t := \mathbb{P}(T \leq \Delta t)$ where $T \sim \text{Exp}(\gamma)$

\mathcal{R}_t^0 : basic reproduction number, $\mathcal{R}_t^0 \equiv \beta_t/\gamma_t$

\mathcal{R}_t : effective reproduction number, $\mathcal{R}_t \equiv \mathcal{R}_t^0 S_t/M$

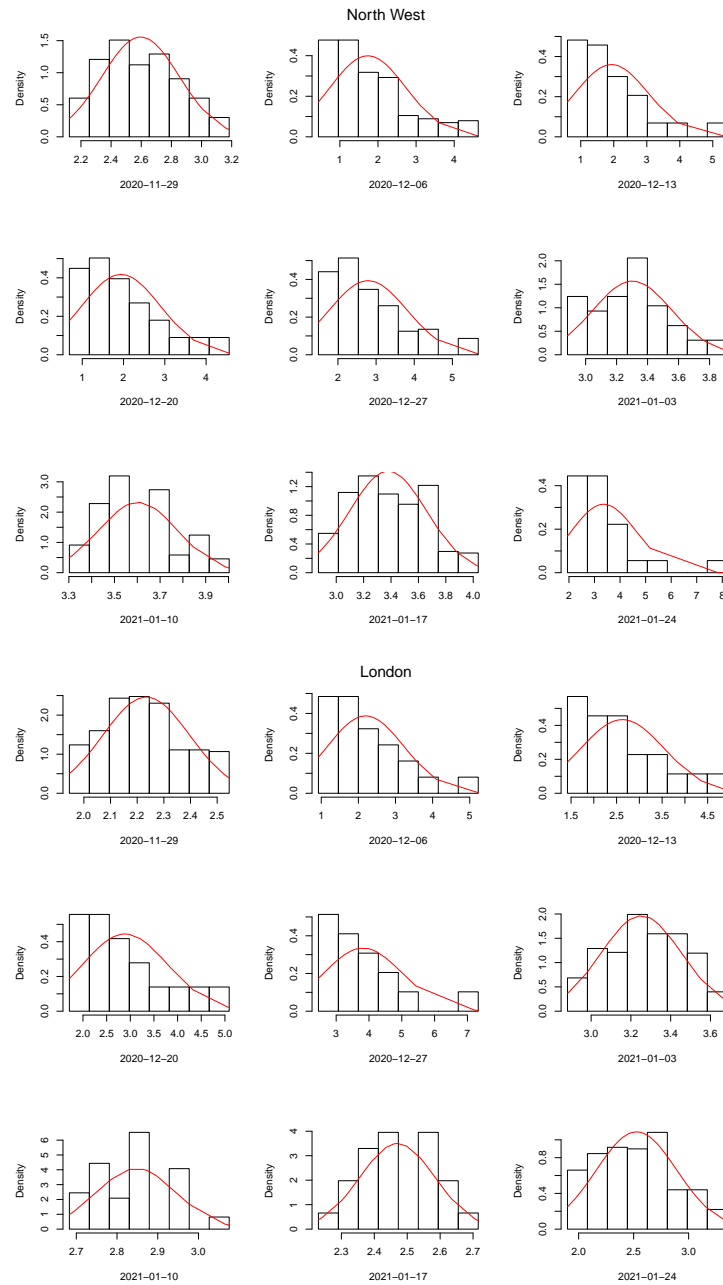


Figure 13: Comparison of moment-matched Gaussian EB prior (12) (red lines) with raw estimates (histograms) on δ for PHE regions North West (top) and London (bottom) from 29th November 2020 to 24th January 2021.

Distribution of the number of new infections ΔQ_t

Under the standard DTMC SIR model, the number of *new* infections, denoted here ΔQ_t , occurring in the time interval Δt up to time t has conditional distribution¹⁰

$$\mathbb{P}(\Delta Q_t | S_{t-1}, \beta_{t-1}, I_{t-1}) = \text{Binomial} \left(\Delta Q_t | S_{t-1}, \frac{\beta_{t-1} I_{t-1}}{M} \right). \quad (31)$$

The probability in (31) can be parameterised by the effective reproduction number, \mathcal{R}_t :

$$\mathcal{R}_t := \frac{\beta_t S_t}{\gamma_t M} \quad (32)$$

$$\mathbb{P}(\Delta Q_t | S_{t-1}, \mathcal{R}_{t-1}, I_{t-1}) \equiv \text{Binomial} \left(\Delta Q_t | S_{t-1}, \frac{\gamma_t \mathcal{R}_{t-1} I_{t-1}}{S_{t-1}} \right). \quad (33)$$

We approximate (33) with a Poisson distribution as follows [45]:¹¹

$$\mathbb{P}(\Delta Q_t | \mathcal{R}_{t-1}, I_{t-1}) := \text{Poisson}(\Delta Q_t | \gamma_t \mathcal{R}_{t-1} I_{t-1}). \quad (34)$$

Distribution of the number of new recoveries ΔR_t

705 The number of *new* recoveries, denoted ΔR_t , occurring in the time interval Δt up to time t is distributed

$$\mathbb{P}(\Delta R_t | I_{t-1}) = \text{Binomial}(\Delta R_t | I_{t-1}, \gamma_t). \quad (35)$$

Transition probabilities for the number of infectious individuals I_t

The change in the number of infectious individuals at time t , ΔI_t can then be expressed as

$$\Delta I_t = \Delta Q_t - \Delta R_t$$

¹⁰Based on each of $S_{t-1} \equiv M - R_{t-1}^+ - I_{t-1}$ susceptibles at time $t-1$ being infected independently with probability

$$\begin{aligned} & \mathbb{P}(\text{Susceptible infected} | \beta_{t-1} \text{ effective contacts in } (t - \Delta t, t]) \\ &= 1 - \mathbb{P}(\text{Susceptible is not infected} | \beta_{t-1} \text{ effective contacts}) \\ &= 1 - \mathbb{P}(\text{A random effective contact is with a noninfectious individual})^{\beta_{t-1}} \\ &= 1 - \left(1 - \frac{I_{t-1}}{M} \right)^{\beta_{t-1}} \\ &= \frac{\beta_{t-1} I_{t-1}}{M} + O \left(\left[\frac{I_{t-1}}{M} \right]^2 \right). \end{aligned}$$

¹¹According to Rule 2 in [45], the Poisson approximation is reasonable when both of these inequalities hold:

$$\begin{aligned} \gamma_t \mathcal{R}_{t-1} I_{t-1} &> 5 \\ \frac{\gamma_t \mathcal{R}_{t-1} I_{t-1}}{S_{t-1}} &< \frac{1}{2}. \end{aligned}$$

Of the two, the first is the least likely to obtain, but is still reasonable under most circumstances. For a simple example, if we set $\gamma_t = 1$ and $\mathcal{R}_{t-1} = 1$, the number of infectious individuals $I_{t-1} > 5$ is sufficient for the approximation to be reasonable.

this and so the conditional distribution for ΔI_t follows from (34) and (35):

$$\mathbb{P}(\Delta I_t | I_{t-1}, \mathcal{R}_{t-1}) = \sum_{\Delta R_t=0}^{I_{t-1}} \left\{ \text{Binomial}(\Delta R_t | I_{t-1}, \gamma_t) \times \text{Poisson}(\Delta I_t + \Delta R_t | \gamma_t \mathcal{R}_{t-1} I_{t-1}) \right\}. \quad (36)$$

Interval-to-interval transition probabilities are evaluated as

$$\begin{aligned} \mathbb{P}(I_t \in \mathcal{B}_{b'} | I_{t-1} \in \mathcal{B}_b, \mathcal{R}_{t-1}) &= \sum_{k \in \mathcal{B}_b} \mathbb{P}(I_{t-1} = k | I_{t-1} \in \mathcal{B}_b) \times \mathbb{P}(k + \Delta I_t \in \mathcal{B}_{b'} | I_{t-1} = k, \mathcal{R}_{t-1}) \\ &= \sum_{k \in \mathcal{B}_b} \frac{1}{e_b - e_{b-1}} \times \mathbb{P}(k + \Delta I_t \in \mathcal{B}_{b'} | I_{t-1} = \check{I}_b, \mathcal{R}_{t-1}) \end{aligned} \quad (37)$$

710 where the first term in the sum at (37) follows from Assumption 2 at (28), and the second term is conditional on prevalence at bin midpoint ($I_{t-1} = \check{I}_b$) based on Assumption 3 at (29)-(30), and can be evaluated using (36).

Transition probabilities for the number of immune individuals R_t^+

Denote by ΔV_t the number of vaccinations administered in interval $(t - \Delta t, t]$. Only a subgroup of those individuals vaccinated at time t may have been susceptible at time $t - \Delta t$; we denote the number in the subgroup by $\Delta \tilde{V}_t$ ($\leq \Delta V_t$), and evaluate its conditional distribution as follows:

$$\begin{aligned} \Delta \tilde{V}_t &:= \# \text{ susceptibles newly vaccinated in } (t - \Delta t, t] \\ \mathbb{P}(\Delta \tilde{V}_t | \Delta V_t, R_{t-1}^+, I_{t-1}) &= \text{HyperGeom}(\Delta \tilde{V}_t | M - V_t, M - R_{t-1}^+ - I_{t-1}, \Delta V_t), \end{aligned} \quad (38)$$

where V_t is the current number of vaccinated individuals in the population (with $\Delta V_t \equiv V_t - V_{t-1}$). The total number of immune, i.e. vaccinated and/or recovered, individuals at time t (denoted R_t^+) can then be represented by the recurrence

$$R_t^+ = R_{t-1}^+ + \Delta R_t + \Delta \tilde{V}_t.$$

This leads to the Markov conditional distribution for R_t^+ via convolution of (35) with (38)

$$\begin{aligned} \mathbb{P}(R_t^+ | R_{t-1}^+, I_{t-1}, \Delta V_t) &= \sum_{\Delta R_t=0}^{I_{t-1}} \left\{ \text{Binomial}(\Delta R_t | I_{t-1}, \gamma_t) \right. \\ &\quad \left. \times \text{HyperGeom}(R_t^+ - R_{t-1}^+ - \Delta R_t | M - V_{t-1}, M - R_{t-1}^+ - I_{t-1}, \Delta V_t) \right\}. \end{aligned} \quad (39)$$

715 The above treatment of immunity assumes individuals are made permanently immune immediately through either vaccination or infection. It would be straightforward to relax the above formulation to allow for more sophisticated treatment of immunity, for example specifying (a) a delay in vaccine effects, (b) incomplete vaccine efficacy (e.g. in the case of novel variants), or (c) decaying immunity over time.

720 Inference on the basic reproduction number

The basic reproduction number at time t , \mathcal{R}_t^0 is related to the effective reproduction number \mathcal{R}_t by the following equation,

$$\mathcal{R}_t^0 = \frac{S_t}{M} \mathcal{R}_t, \quad (40)$$

where M is the total number of individuals and S_t is the number of susceptible individuals at time t . Recall that $S_t \equiv M - R_t^+ - I_t$ where R_t^+ is the number of immune individuals and I_t is the number of infectious individuals, both of which are estimated by our DTMC SIR model. We can plug in these estimates into (40) to estimate \mathcal{R}_t^0 for a given LTLA. Figure 15 plots \mathcal{R}_t^0 and \mathcal{R}_t for a selection of LTLAs.

PCR positive to infectious mapping – method details

Recall we require $\mathbb{P}(\tilde{I} | I)$ in (3), which is the probability distribution on the number of PCR positive individuals \tilde{I} given the number of infectious individuals I . This can be expressed via Bayes' theorem as

$$\mathbb{P}(\tilde{I} | I) \propto \mathbb{P}(I | \tilde{I})\mathbb{P}(\tilde{I}) \quad (41)$$

where the likelihood is binomial:

$$\mathbb{P}(I | \tilde{I}) = \text{Binomial}(I | \tilde{I}, \mathbb{P}(\text{Infectious} | \text{PCR positive})) . \quad (42)$$

To target the $\mathbb{P}(\text{Infectious} | \text{PCR positive})$ success probability in (42), we introduce the following notation:

$$\text{Infected}_t \equiv \text{Individual becomes infected in week } t \quad (43)$$

$$\text{Infectious}_t \equiv \text{Individual is infectious in week } t \quad (44)$$

$$\text{PCR}_{+t} \equiv \text{Individual is PCR positive from swab taken in week } t \quad (45)$$

and proceed as follows:¹²

$$\mathbb{P}(\text{Infectious}_t | \text{PCR}_{+t}) \quad (46)$$

$$= \frac{\mathbb{P}(\text{Infectious}_t \wedge \text{PCR}_{+t})}{\mathbb{P}(\text{PCR}_{+t})} \quad (47)$$

$$= \frac{\sum_{k=0}^3 \mathbb{P}(\text{Infectious}_t \wedge \text{PCR}_{+t} | \text{Infected}_{t-k})\mathbb{P}(\text{Infected}_{t-k})}{\sum_{k=0}^3 \mathbb{P}(\text{PCR}_{+t} | \text{Infected}_{t-k})\mathbb{P}(\text{Infected}_{t-k})} \quad (48)$$

$$= \frac{\sum_{k=0}^3 \mathbb{P}(\text{Infectious}_t | \text{Infected}_{t-k})\mathbb{P}(\text{PCR}_{+t} | \text{Infected}_{t-k})\mathbb{P}(\text{Infected}_{t-k})}{\sum_{k=0}^3 \mathbb{P}(\text{PCR}_{+t} | \text{Infected}_{t-k})\mathbb{P}(\text{Infected}_{t-k})} , \quad (49)$$

where, at (49), we assumed conditional independence between Infectious_t and PCR_{+t} conditional on Infected_{t-k} . Also, at (48), we assumed that testing PCR positive implies that an individual was infected at most four weeks prior to being swabbed, which is consistent with Figure 1A of [34] (data input 2 below). We import three distinct data inputs to estimate the various terms in (49).

¹²We use \wedge to denote logical AND.

Data input 1 – Infectious interval

Figure 1A of Ferretti et al. [33] shows the estimated probability density function of the serial interval for SARS-CoV-2 transmission – **we denote this density function** $f_{\text{Fer}}(d)$. Noting the support of this density to be approximately $[1, 11]$, we specify that an average individual is infectious between days 1 to 11. Formally we define, independently for each individual in the population,

$$\begin{aligned} \mathbb{P}(\text{Infectious on } d\text{th day post-infection}) &:= \\ &\mathbb{I} \{ \mathbb{E}_X [\mathbb{P}(\text{individual } X \text{ Infectious on } d\text{th day post-infection})] > 0 \} \\ &\approx \begin{cases} 1 & \text{if } f_{\text{Fer}}(d) > 0, \text{ i.e. if } 1 \leq d \leq 11 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

where X denotes an individual selected uniformly at random from the population. We can use this to estimate the $\mathbb{P}(\text{Infectious}_t \mid \text{Infected}_{t-k})$ term appearing in the numerator of (49) as follows

$$\mathbb{P}(\text{Infectious}_t \mid \text{Infected}_{t-k}) \approx \begin{cases} 6/7 & k = 0 \\ 5/7 & k = 1 \\ 0 & k > 1. \end{cases} \quad (50)$$

Data input 2 – PCR positive interval

Figure 1A of Hellewell et al. [34] plots posterior probabilities (with credible intervals) of testing PCR positive when swabbed d days post infection. We denote this data input

$$\mathbb{P}_{\text{Hel}}(\text{PCR+} \mid \text{swabbed day } d \text{ after becoming infected}) \quad (51)$$

and use it to estimate the term $\mathbb{P}(\text{PCR+}_t \mid \text{Infected}_{t-k})$ appearing twice in (49), evaluating the following estimator for each $k = 0, \dots, 3$:

$$\mathbb{P}(\text{PCR+}_t \mid \text{Infected}_{t-k}) \approx \frac{1}{7} \sum_{d=7k}^{7(k+1)-1} \mathbb{P}_{\text{Hel}}(\text{PCR+} \mid \text{swabbed day } d \text{ after becoming infected}) \quad (52)$$

735 Hellewell et al. [34] helpfully provide reproducible scripts¹³ and we use these to extract the posterior distribution on $\mathbb{P}_{\text{Hel}}(\text{PCR+} \mid \text{swabbed day } d \text{ after becoming infected})$ from their Figure 1A, whose uncertainty we propagate to estimator (52) and onwards to (49), yielding a distribution on $\mathbb{P}(\text{Infectious}_t \mid \text{PCR+}_t)$ which we take forward approximated by a moment-matched Beta distribution (at each week t) to be used as an EB conjugate prior on the success probability in (42).

740 Data input 3 – Pillar 1+2 incidence

For the purposes of adjusting the PCR positive map to changing incidence, we use the raw regional weekly positive test counts $n_{0:T}$, where we denote weeks by $t = 0, \dots, T$. We use this data input

¹³<https://github.com/cmmid/pcr-profile>

to estimate the term $\mathbb{P}(\text{Infected}_{t-k})$ appearing twice in (49), evaluating the following estimator for each $k = 0, \dots, 3$:¹⁴

$$\mathbb{P}(\text{Infected}_{t-k}) = \mathbb{P}(\text{Infected}_{t-k} \wedge [\vee_{k'=0}^3 \text{Infected}_{t-k'}]) \quad (53)$$

$$= \mathbb{P}(\text{Infected}_{t-k} \mid \vee_{k'=0}^3 \text{Infected}_{t-k'}) \mathbb{P}(\vee_{k'=0}^3 \text{Infected}_{t-k'}) \quad (54)$$

$$\approx \frac{n_{t-k}}{\sum_{k'=0}^3 n_{t-k'}} \mathbb{P}(\vee_{k'=0}^3 \text{Infected}_{t-k'}) \quad (55)$$

which can be directly substituted for $\mathbb{P}(\text{Infected}_{t-k})$ in top and bottom of (49) with the second term on the right of (55) cancelling between numerator and denominator, and therefore not requiring evaluation. We note that we are using raw counts to model relative incidence over a relatively short period (four weeks), which is making the assumption that the bias is relatively stable over this timeframe (see Assumption Infectious-1 in *SI-PCR+ to infectious mapping – assumptions and caveats*).

Sensitivity analyses

Prior hyperparameters for δ

The EB prior for δ depends on two hyperparameters: σ_ϵ controls the variance of the white noise associated with each individual time point, while ψ controls the degree of autocorrelation from one time point to the next. Figures 16 and 17 show the estimates for prevalence and \mathcal{R}_t respectively of infectious individuals using different values of these two hyperparameters. Note that in the main text, we present results using $\sigma_\epsilon = 1$ and $\psi = 0.99$.

Sensitivity and specificity of PCR tests

PCR tests are not perfect and are subject to both false positives and false negatives. In our analysis, we account for imperfect testing via the false positive rate, α , and the false negative rate, β (see (8)). Figures 18 and 19 show the estimates for prevalence and \mathcal{R}_t respectively of infectious individuals using different values of these two hyperparameters. Note that in the main text, we present results using $\alpha = 0.001$ and $\beta = 0.05$.

¹⁴We use \vee to denote logical OR.

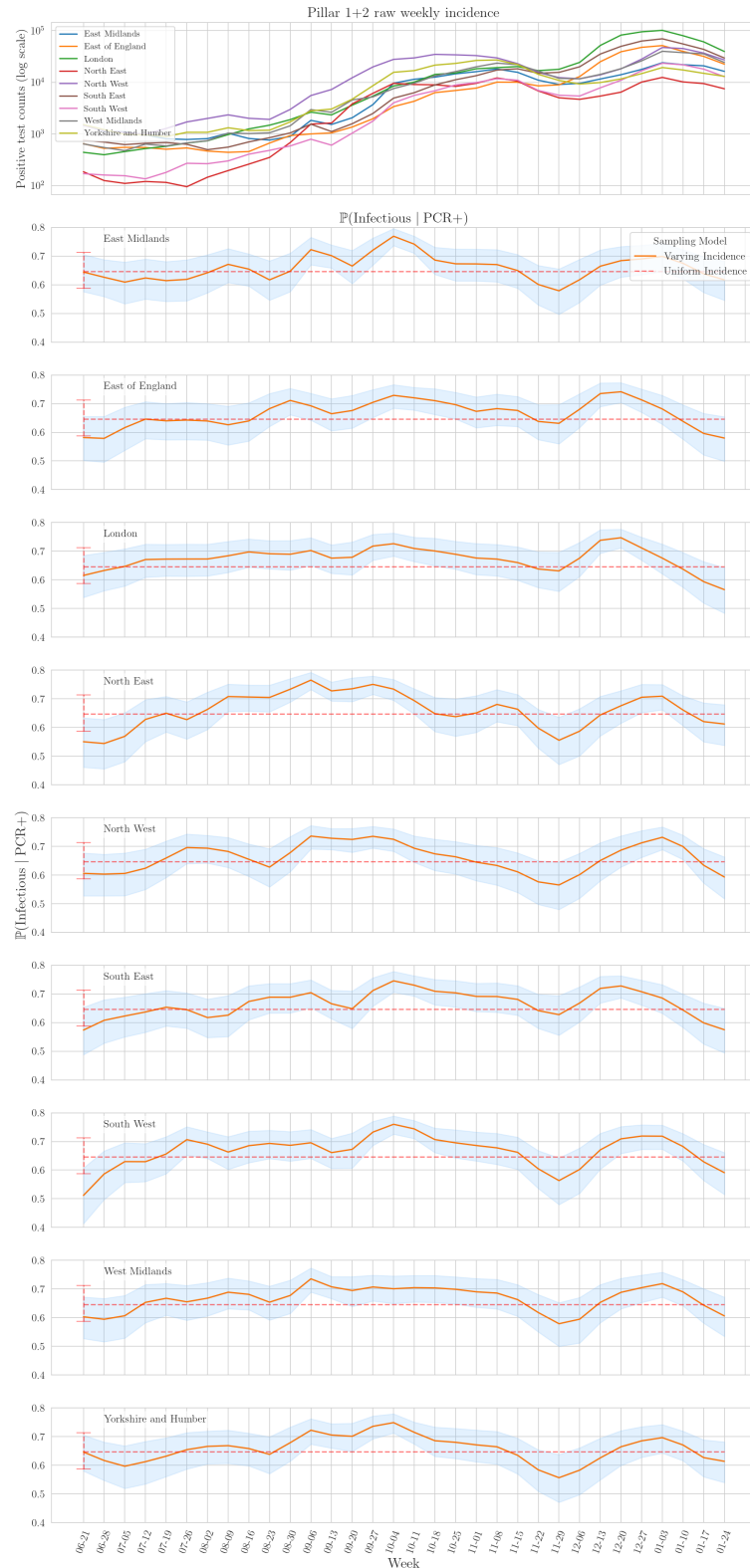


Figure 14: EB prior on $\mathbb{P}(\text{Infectious} \mid \text{PCR positive})$ by week and PHE region. The top panel shows raw weekly Pillar 1+2 incidence for the nine PHE regions; this is to provide intuition for the Varying incidence model in the panels below. The bottom nine panels display the prior we place on $\mathbb{P}(\text{Infectious}_t \mid \text{PCR}+_t)$, which is specific to week and region for the Varying incidence model, but is constant across weeks/regions for the Uniform incidence model (see legend in panel 2). Error bars (at left of panel for Uniform incidence; around curve for Varying incidence) represent 95% credible intervals.

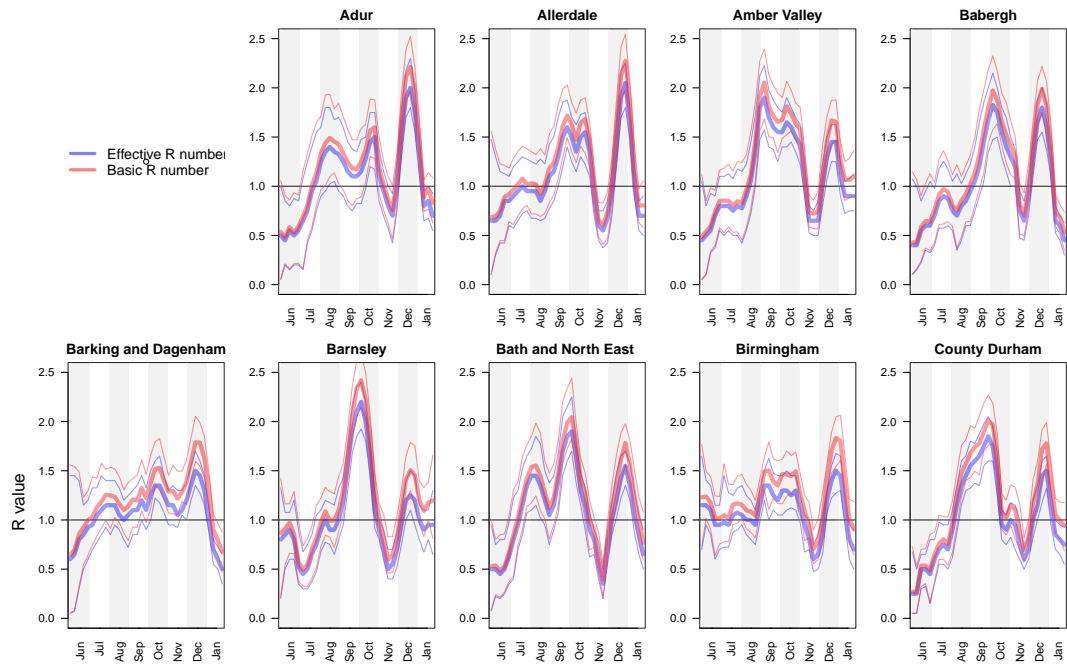


Figure 15: Comparison of \mathcal{R}_t^0 and \mathcal{R}_t estimates.

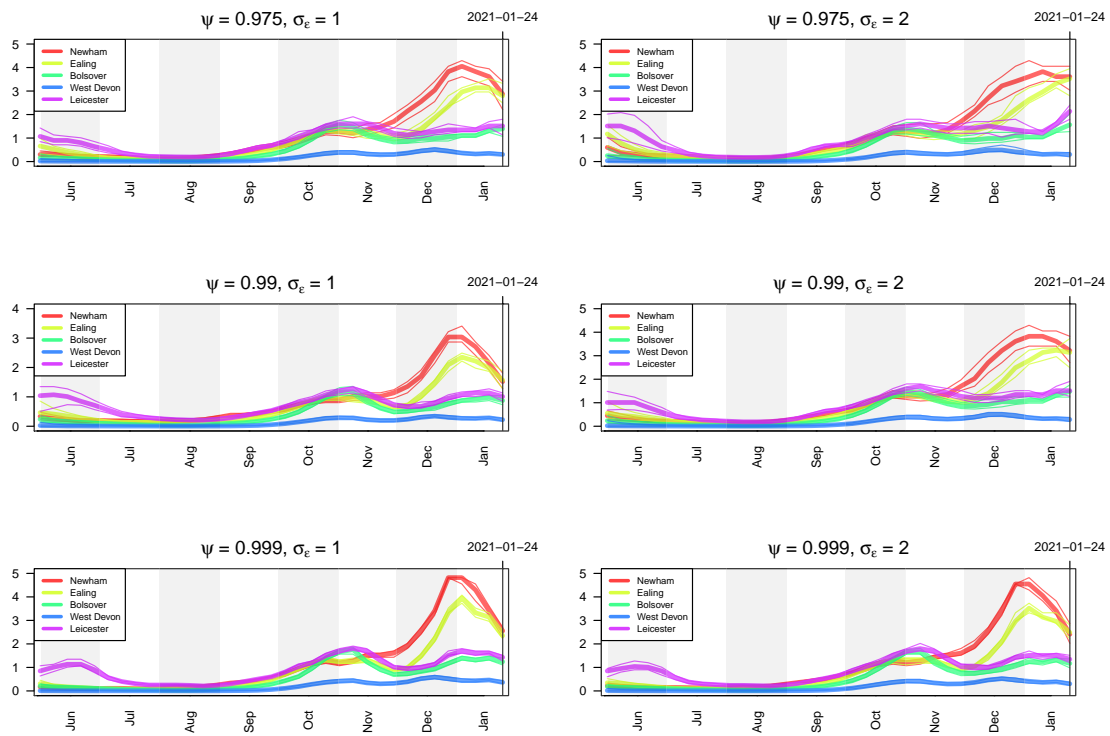


Figure 16: Estimates of prevalence of infectious individuals for five LTLAs using different values of the hyperparameters σ_ϵ and ψ controlling the smoothness of the bias parameter δ .

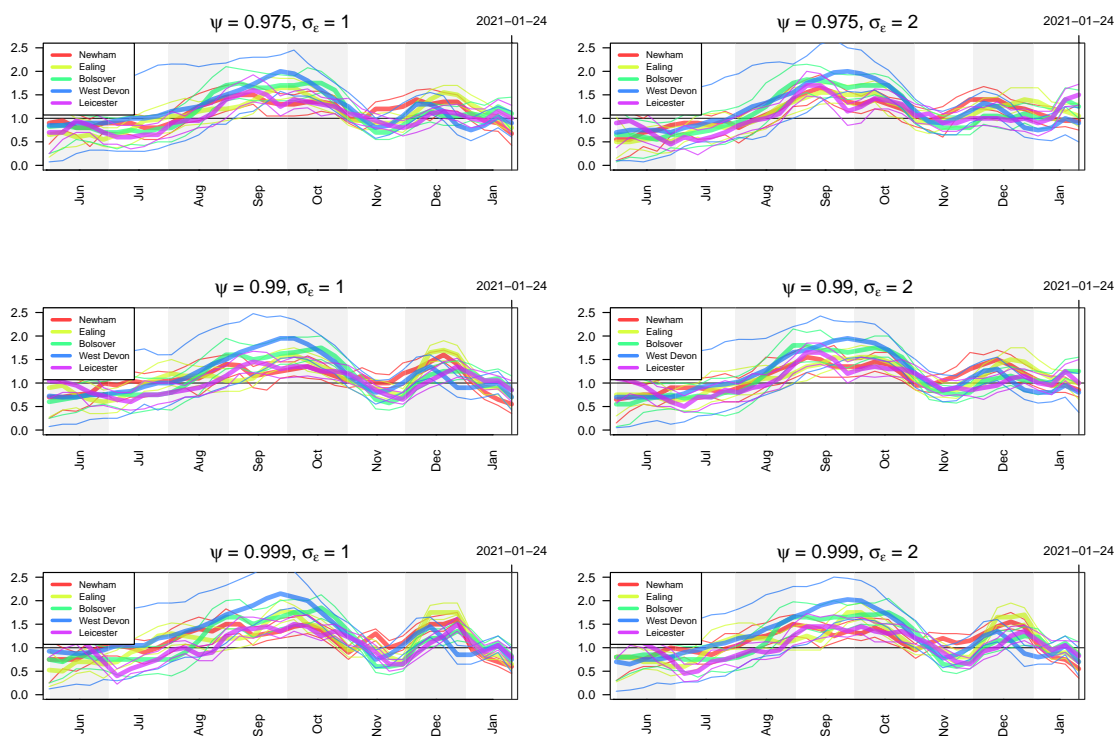


Figure 17: Estimates of prevalence of infectious individuals five LTLAs using different values of the false positive rate α and false negative rate β controlling the smoothness of the bias parameter δ .

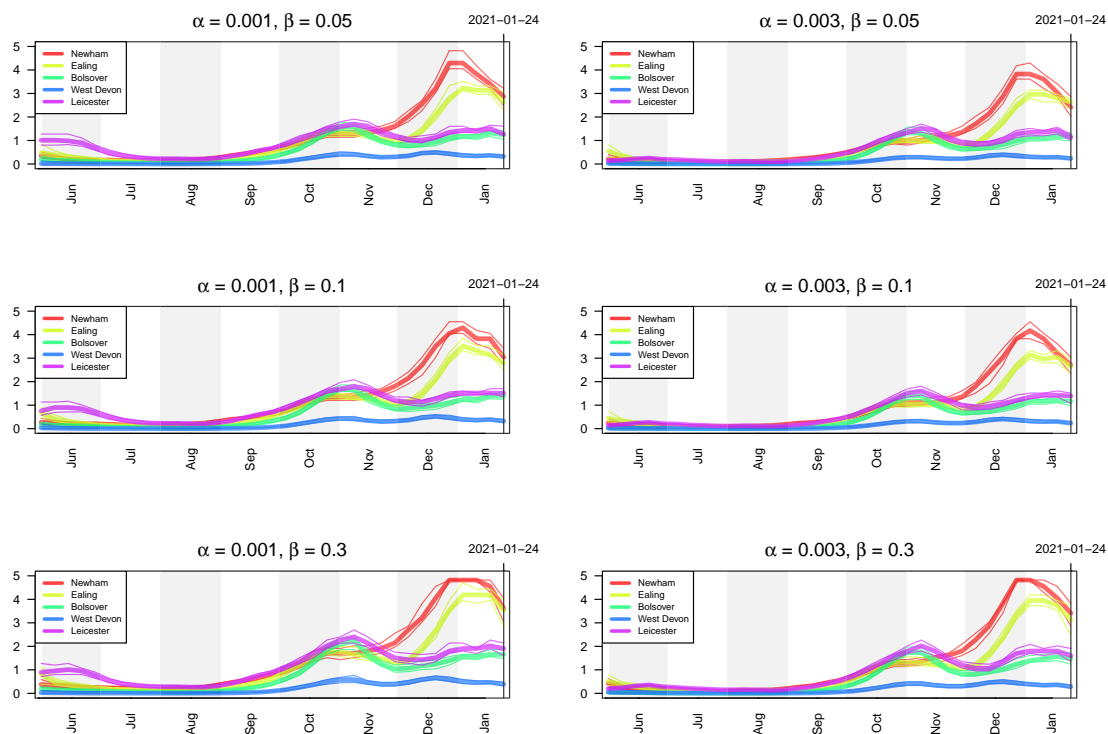


Figure 18: Estimates of \mathcal{R}_t of infectious individuals five LTLAs using different values of the false positive rate α and false negative rate β controlling the smoothness of the bias parameter δ .

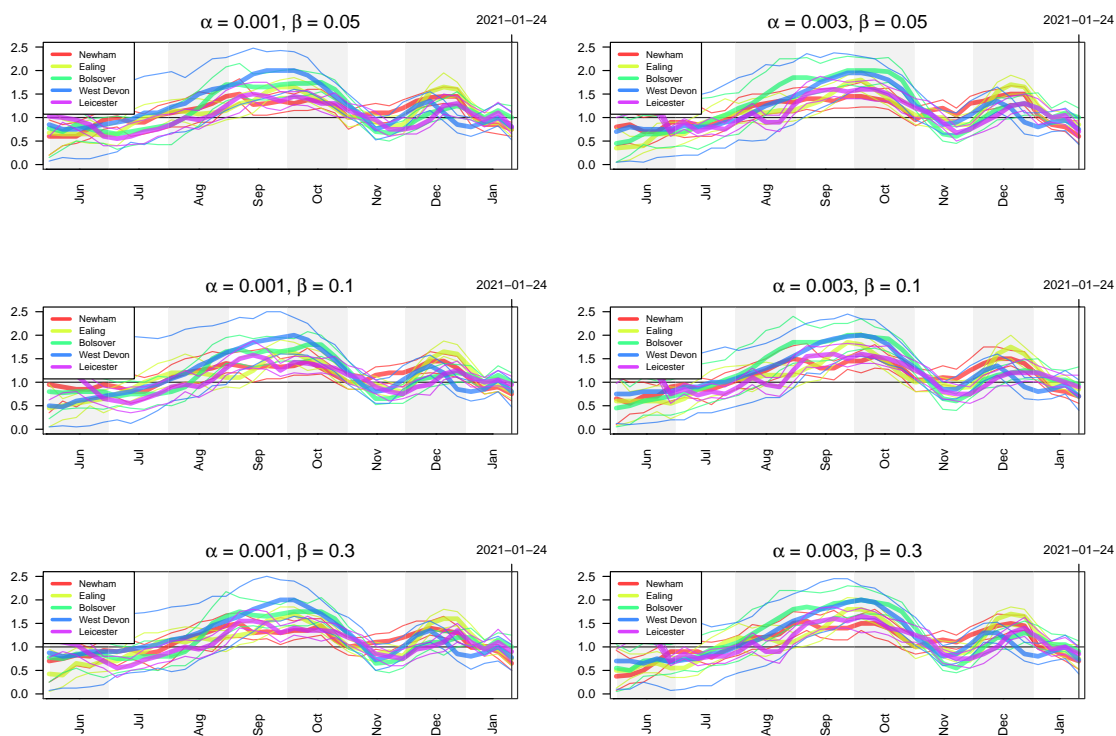


Figure 19: Estimates of prevalence of infectious individuals for five LTLAs using different values of the hyperparameters σ_ϵ and ψ controlling the smoothness of the bias parameter δ .