

1 Assessing the importance of demographic risk factors across  
2 two waves of SARS-CoV-2 using fine-scale case data

3 A.J. Wood<sup>1</sup>, A.R. Sanchez<sup>1</sup>, P.R. Bessell<sup>1</sup>, R. Wightman<sup>2</sup>, and R.R. Kao<sup>\*1,3</sup>

4 <sup>1</sup>Roslin Institute, University of Edinburgh

5 <sup>2</sup>Edinburgh Medical School, University of Edinburgh

6 <sup>3</sup>Royal (Dick) School of Veterinary Studies, University of Edinburgh

7 September 8, 2023

8 **Abstract**

9 For the long term control of an infectious disease such as COVID-19, it is crucial to identify  
10 the most likely individuals to become infected and the role that differences in demographic  
11 characteristics play in the observed patterns of infection. As high-volume surveillance winds  
12 down, testing data from earlier periods are invaluable for studying risk factors for infection  
13 in detail. Observed changes in time during these periods may then inform how stable the  
14 pattern will be in the long term.

15 To this end we analyse the distribution of cases of COVID-19 across Scotland in 2021,  
16 where the location (census areas of order 500–1,000 residents) and reporting date of cases are  
17 known. We consider over 450,000 individually recorded cases, in two infection waves triggered  
18 by different lineages: B.1.1.529 (“Omicron”) and B.1.617.2 (“Delta”). We use random forests,  
19 informed by measures of geography, demography, testing and vaccination. We show that the  
20 distributions are only adequately explained when considering multiple explanatory variables,  
21 implying that case heterogeneity arose from a combination of individual behaviour, immunity,  
22 and testing frequency.

23 Despite differences in virus lineage, time of year, and interventions in place, we find the  
24 risk factors remained broadly consistent between the two waves. Many of the observed smaller  
25 differences could be reasonably explained by changes in control measures.

\*Corresponding author: [rowland.kao@ed.ac.uk](mailto:rowland.kao@ed.ac.uk)  
**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

## 26 1 Introduction

27 A key challenge in the long term control of an infectious disease is to identify predictable patterns of  
28 incidence. The emergence and spread of the SARS-CoV-2 virus saw restrictions imposed globally  
29 on everyday life to control the spread of COVID-19 infection, and to protect individuals at highest  
30 risk of severe disease. While as of March 2023 few to no restrictions remain in place in Scotland,  
31 as in the rest of the UK, randomised testing [1] and hospital admissions [2] indicate continued  
32 widespread transmission. The winding down of community testing and other surveillance is making  
33 it more difficult to track the transmission patterns of COVID-19 in detail.

34 Typically, identifying risk factors for infection rely on disease surveillance studies. While these  
35 studies can be powerful and provide important insights [3, 4, 5, 6], they are often expensive,  
36 laborious and time consuming. “Big Data” in the health sciences offers an opportunity to gain  
37 some of the same insights using routinely collected data. The availability of COVID-19 case  
38 data at fine spatial scales with detailed metadata enables us to identify important health-related  
39 risks, with the data collected during the pandemic being made available to researchers in close to  
40 real-time.

41 In this work we aim to identify risk factors for COVID-19 cases in Scotland, and their change  
42 over time, to serve as an indicator for how the longer-term profile of infection may evolve. We  
43 fit the case distributions of two different waves of COVID-19, with a machine learning model  
44 informed by a range of explanatory variables relating to geography and demographics.

45 The first COVID-19 case in Scotland was identified on 1<sup>st</sup> March 2020 [7]. The Scottish  
46 Government imposed strict “lockdown” non-pharmaceutical interventions (NPIs) on 23<sup>rd</sup> March  
47 2020 [8]. While initially applied at the national level, following the initial lockdown period NPIs  
48 were adjusted by local authority (administrative areas with populations ranging between 22,540–  
49 635,130) through a “levels”-based system [9]. The seeding and rapid spread of the B.1.1.7 lineage  
50 (termed the “Alpha” variant) in December 2020 led to a tightening of NPIs and a second lock-  
51 down [10, 11]. A mass vaccination programme began in December 2020 [12, 13], prioritising the  
52 elderly and healthcare workers, with all adults eventually eligible.

53 We focus on case data gathered between May 2021 and January 2022, a period that saw the  
54 steady relaxation of nearly all NPIs [14]. This period had two major waves of infection: the first  
55 from May 2021 triggered by the B.1.617.2 lineage (“Delta”), and a second wave from November  
56 2021 by the B.1.1.529 B.A.1 lineage (“Omicron”). The deletion of two specific amino acids in the  
57 Omicron sub-variant distinguished it from most co-circulating variants including Delta, in PCR

58 tests that have an accompanying “S-gene” test result [15]. A high-capacity testing programme  
59 was in place throughout, with free-of-charge lateral flow testing strongly encouraged, and PCR  
60 testing mandated for those with symptoms, or a lateral flow positive.

61 Earlier work has exploited finely-grained case data to highlight risk factors for cases and se-  
62 vere outcomes including (but not limited to) sex [16, 17, 18], population density [19, 20, 21],  
63 deprivation [22, 23, 24, 25], occupation [26, 27, 28], and age [29, 30, 31]. Similar studies have  
64 incorporated movement data [32] to demonstrate the protective impact of NPIs that restrict mo-  
65 bility [21, 33, 34, 35, 36, 37]. Many of these studies focus on “first wave” of infection, during which  
66 strict NPIs were imposed and no population immunity had been established. This study focuses  
67 on a more advanced period moving away from NPIs, and the conditions for disease spread com-  
68 paratively less “exceptional”. This is especially the case for the Omicron wave. A unique feature  
69 of our model is the inclusion of lateral flow test taking *frequency*. The proportion of infections that  
70 end up reported is likely to depend on testing propensity, and we consider how that may lead to  
71 distortions in the case distribution.

72 Our main finding is that the risk factors for cases remained broadly consistent across both  
73 waves. Differences between the two waves either offer relatively small scale changes in demographic  
74 risk or are consistent with the impact of changes in approaches to control.

## 75 **2 Results**

76 The period November 15<sup>th</sup> 2021 – January 6<sup>th</sup> 2022 covers the first outbreak and peak of the  
77 B.1.1.529 lineage (BA.1 sublineage, hereafter referred to as the Omicron variant) (S-gene “dropout”  
78 test signature). Prior to this, the B.1.617.2 lineage (Delta variant) (S-gene positive test signature)  
79 was dominant. From 15<sup>th</sup> November 2021, S-gene dropout cases consistently rise, and all subse-  
80 quent “dropout” cases are assumed Omicron. Remaining S-gene positive cases are presumed to  
81 be Delta, consistent with nationwide sequence data [38].

### 82 **2.1 Time evolution and early patterns of spread**

83 We identified 385,558 cases between November 15<sup>th</sup> 2021 and January 6<sup>th</sup> 2022, of which 227,286  
84 were likely Omicron. From 1<sup>st</sup> May 2021 to 7<sup>th</sup> September 2021 we identified 269,838 cases, of  
85 which 229,073 were likely Delta. The remaining cases in these periods (those with no S-gene result,  
86 or a different result) are excluded. The start date for each of these periods is the first date from

87 which there are consistent rises in cases that are likely the new variant.

88 Omicron cases had a doubling time (the time taken for *newly reported daily* cases to double) of  
89 2.9 days over the first 28 days, compared to 6.2 days for Delta (Supplementary Material, Fig. S1).  
90 Over half of all DZs had reported an Omicron case in the wave within 29 days, whereas for Delta  
91 this took 39 days (Supplementary Material, Fig. S2).

92 The reproduction number  $R_t$  consistently rose for Omicron, peaking at above 2 for nearly all  
93 local authorities 28 days in to the outbreak, and only consistently falling below 1 after 50 days  
94 (Supplementary Material, Fig. S3). Reproduction numbers for Delta are less consistent between  
95 LAs; while the number generally remains above 1 for most LAs in the period, there is no coherent  
96 peak at the start of the wave.

97 In the intermediate period during which Omicron became dominant and Delta declined, the  
98 *age* distributions by variant differed. Taking the mid-points of the five-year age brackets, the  
99 mean ages of the Delta-type cases was 3.9 years lower than the Omicron-type cases (31.8 years  
100 compared to 35.7 years). A Student's  $t$  test shows this difference to be statistically significant  
101 ( $t = -52.2$ ,  $p < 0.001$ ). This was the case from relatively early on when Omicron accounted for at  
102 least 5% of cases (Supplementary Material, Fig. S4A). However, the median ages are equal (both  
103 32.5 years), as in the Omicron-type cases there is a trough in those aged 0–14, with fewer than  
104 50% of cases in this age group Omicron, but then a peak in the 20–29 age group (Supplementary  
105 Material, Fig. S4B).

## 106 2.2 Case distribution and model fit

107 Fig. 1, shows the distribution of COVID-19 cases for the Omicron and Delta waves broken down  
108 by age, sex, prior cases (serving as a proxy for prior immunity from infection), deprivation and  
109 health board. Omicron case rates were highest in younger adults, peaking at 90 cases/1,000 in ages  
110 20–24. There was only a small difference in rates between men and women. Case rates were much  
111 lower amongst those that had tested positive for COVID-19 previously. Fig. 2 shows case rates  
112 per DZ. Geographically, case rates fall with increasing rurality, most notably in Orkney, Shetland  
113 and the Western Isles (all island communities). The trend with respect to multiple deprivation  
114 decile is bimodal, with higher rates towards the highest and lowest deciles.

115 The fit case rates from our random forest regression models are overlaid onto Fig. 1. We achieve  
116 a good fit to these larger-scale trends. The model slightly under-fits the age ranges 15–24, where  
117 case rates were the highest overall. Variable importance outputs are presented in Supplementary

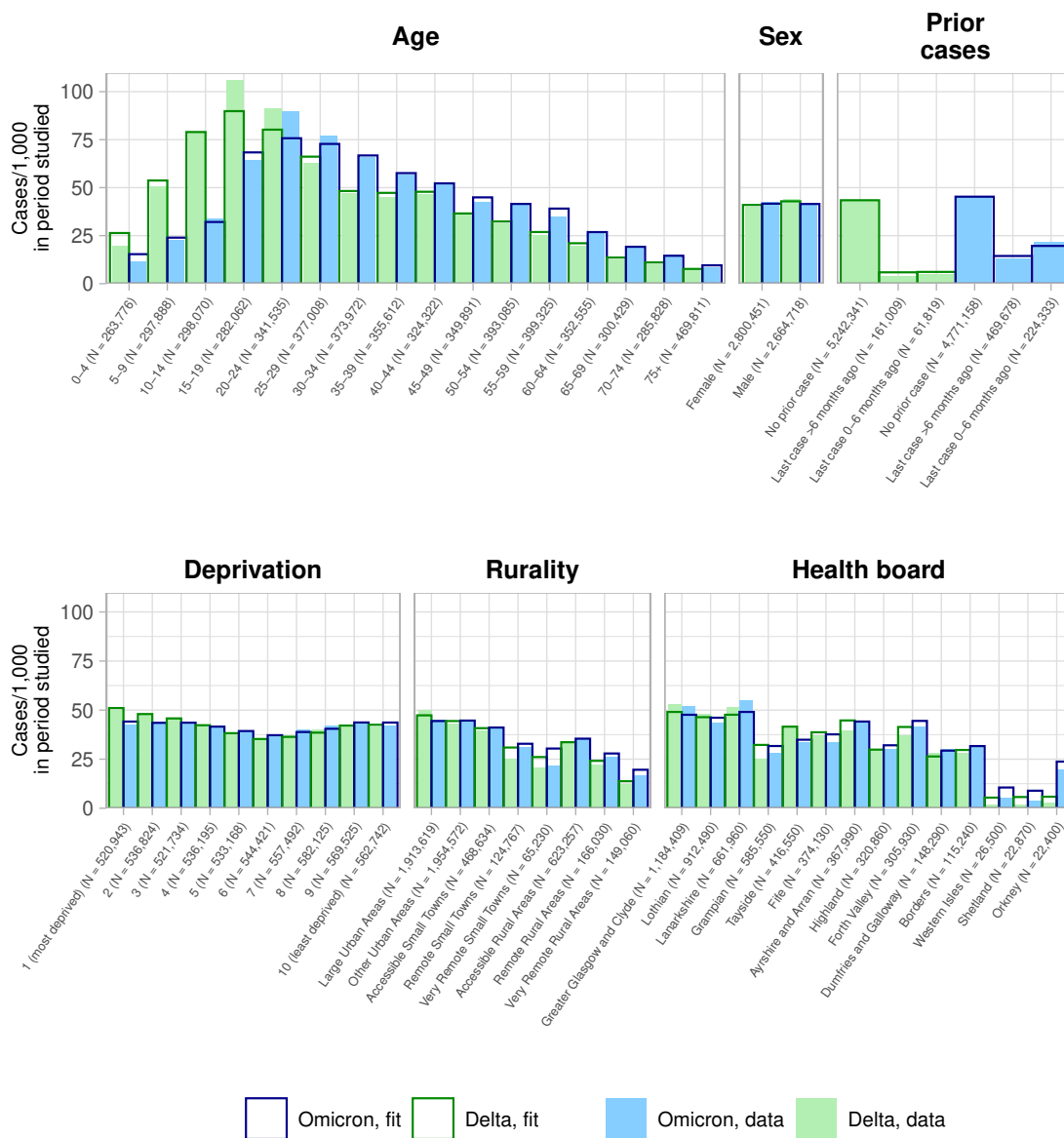


Figure 1: Summary of 227,286 Omicron COVID-19 cases in Scotland between November 15<sup>th</sup> 2021 and January 6<sup>th</sup> 2022 (blue, filled), and 229,073 Delta cases from 1<sup>st</sup> May 2021 to 7<sup>th</sup> September 2021 (green, filled). The full population ( $N = 5,465,169$ ) is broken down by *age range*, *prior case status* (whether a person had previously reported a COVID-19 case prior to that specific wave, and when), *deprivation* (of place of residence, per the SIMD decile, with 1 the most deprived), *rurality* (of place of residence, per the census Urban/Rural Classification) and *location* (at the level of Scottish health board). Cases are given per 1,000 people in that group (with subpopulation  $N$  recorded on the axis labels). The corresponding case rates as fit by our models are superimposed. Note that the subpopulations in the *prior case status* plot change across waves, due to being at different points in time.

118 Material Fig. S7, with node purity and accuracy loss.

119 Fig. 3 (top) shows model performance at DZ level, comparing observed cases to fit cases.

120 Beginning with Omicron cases, our full model explains 70% (fit: 71%, test: 62%) of local variation  
121 in the case distribution (R-squared for case numbers, aggregated at a DZ level), with a poorer fit  
122 for cohorts with very high case counts. A “reduced” random forest model informed by population  
123 and population density alone explained 59% (fit: 60%, test: 55%) of variation. A model informed  
124 by only population/deprivation rank explained 53% (fit: 53%, test: 51%), and one informed by  
125 only population/age explained 48% (fit: 48%, test: 51%). Fig. 3 shows further deviation of the  
126 data-fit slopes away from the diagonal for these “reduced” models.

127 Considering now earlier Delta cases from 1<sup>st</sup> May to 7<sup>th</sup> September 2021, the geographical  
128 distribution (Fig. 2) is visually similar, with a concentration of high case rates in the denser  
129 “central belt”. Cases skewed slightly younger (Fig. 1), with the highest rates within ages 15–19.  
130 The distribution with respect to deprivation decile remains bimodal, with higher rates in both the  
131 most and least deprived DZs. Model performance was similar, explaining 72% (fit: 73%, test:  
132 61%) of DZ-level variation.

133 Fig. 3 (bottom) shows for both the Delta and Omicron models, autocorrelation of residuals (as  
134 measured by the Moran’s I statistic, Section 4.5) within 1km is 0.35, falling to 0.15 at 5km, and  
135 0.05 at 50km. The reduced models exhibit much higher residual autocorrelation, with the density-  
136 only model performing best, but persisting over larger distances (see Supplementary Material,  
137 Fig. S6 for a map view of residuals).

### 138 **2.3 Accumulated local effects**

139 Fig. 4 shows the accumulated local effects (ALEs) of all explanatory variables in the model (see  
140 Section 4.4 for definition).

141 Population, age, sex, and prior case status have ALEs that follow the empirical distributions  
142 observed in Fig. 1; ALEs are strongly positive for ages between 15–40, and those that had never  
143 reported a case before.

144 Beyond these variables, Fig. 4 shows that features such as low population density, high vac-  
145 cination uptake, a low mean household size, and a low rate of negative LFD test reporting are  
146 protective. We note that for vaccination uptake, the protective value at zero is likely an artefact  
147 arising from cohorts with ages 0–9 that were not eligible.

148 The effects for many variables associated with social deprivation such as the ratio of working  
149 age people with no qualifications and the rate of income deprivation (see Supplementary Material,  
150 Section B.2 for full descriptions) are weaker. This is consistent with the small degree of deprivation-

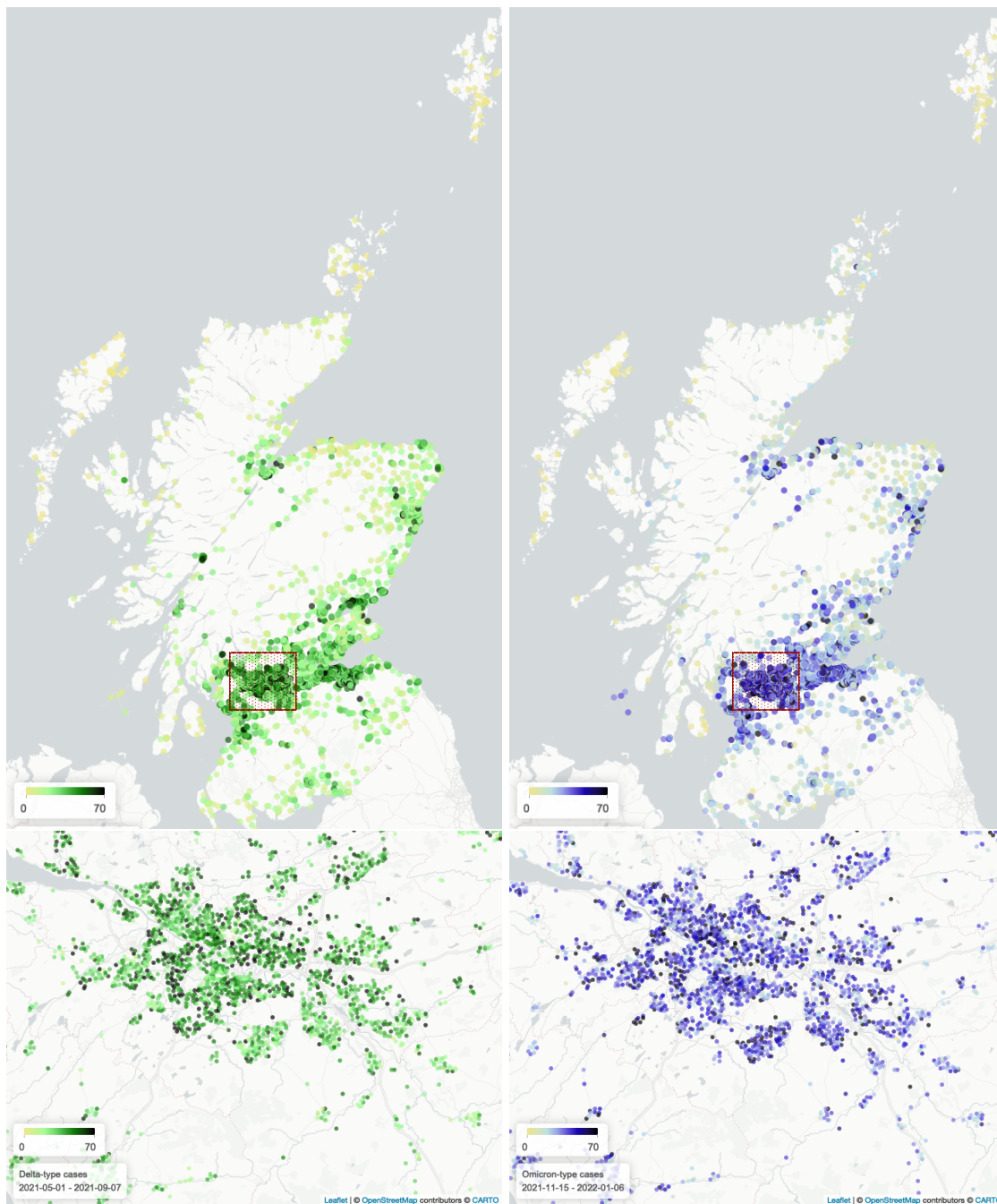


Figure 2: COVID-19 cases in Scotland (top) over the Omicron period (right, blue) as compared to Delta cases (left, green), with focus on the Greater Glasgow region (boxed, bottom). Each point indicates the population *centroid* of a DZ, with the colour representing the number of cases reported.

151 level variation seen in Fig. 1.

152 The directionality of the ALEs remain broadly consistent across both waves. Some risk factors  
153 were more pronounced in the Delta model, including in mean household size, population density

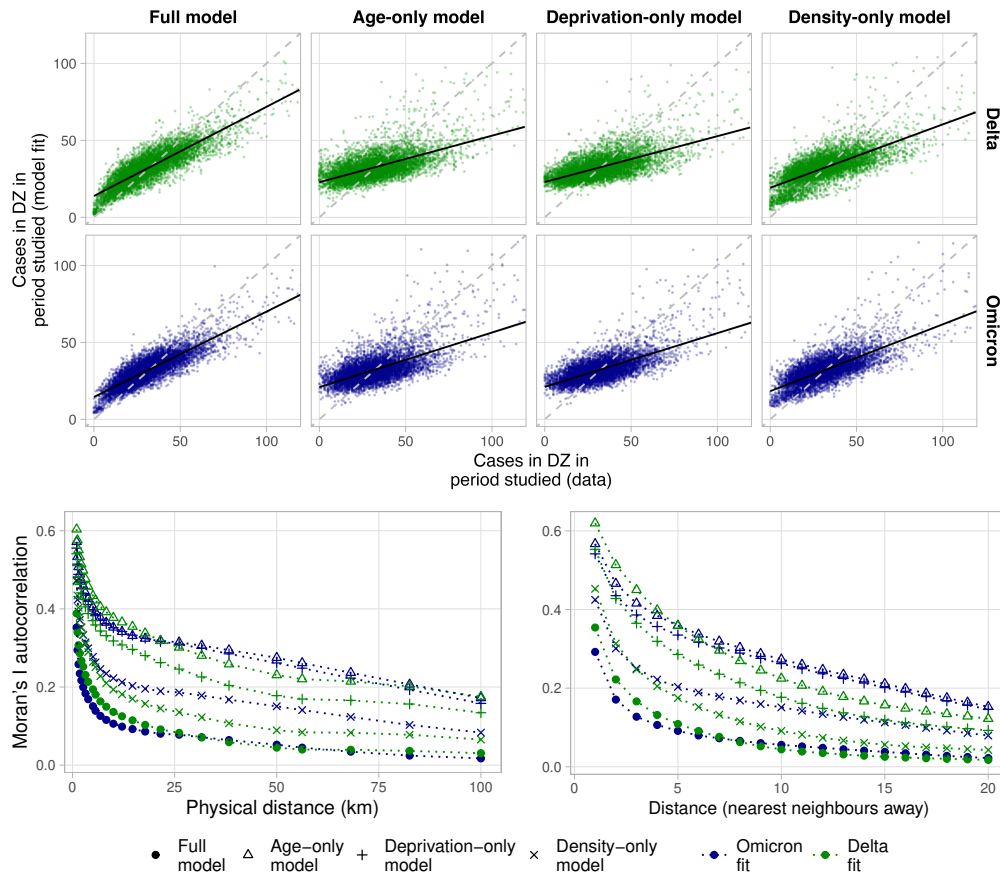


Figure 3: Top: performance of different models, comparing observed cases to fit cases at DZ level. Each point represents a DZ. Points deviating from the diagonal indicate DZs with less accurate fits. The full model is compared with performance of reduced models informed with only population, and one of either age, overall deprivation rank, or population density. Bottom: residual clustering as measured by the Moran's I statistic, at different physical (left) and network-based distances (right). Higher values represent higher autocorrelation between model residuals, when comparing DZs sitting within a given locus. DZs are defined as nearest neighbours of one another if they share a boundary.

154 and the proportion of individuals belonging to a black or minority ethnicity. Conversely, cohorts  
 155 with very high student populations were associated more strongly with high case rates in the  
 156 Omicron fit.

### 157 3 Discussion

158 Scotland's programme of free community testing was an invaluable tool for tracking the spread of  
 159 COVID-19 infection up to early 2022. With the ending of detailed surveillance since, it is more  
 160 difficult to monitor the precise patterns of infection amongst the population and how that will



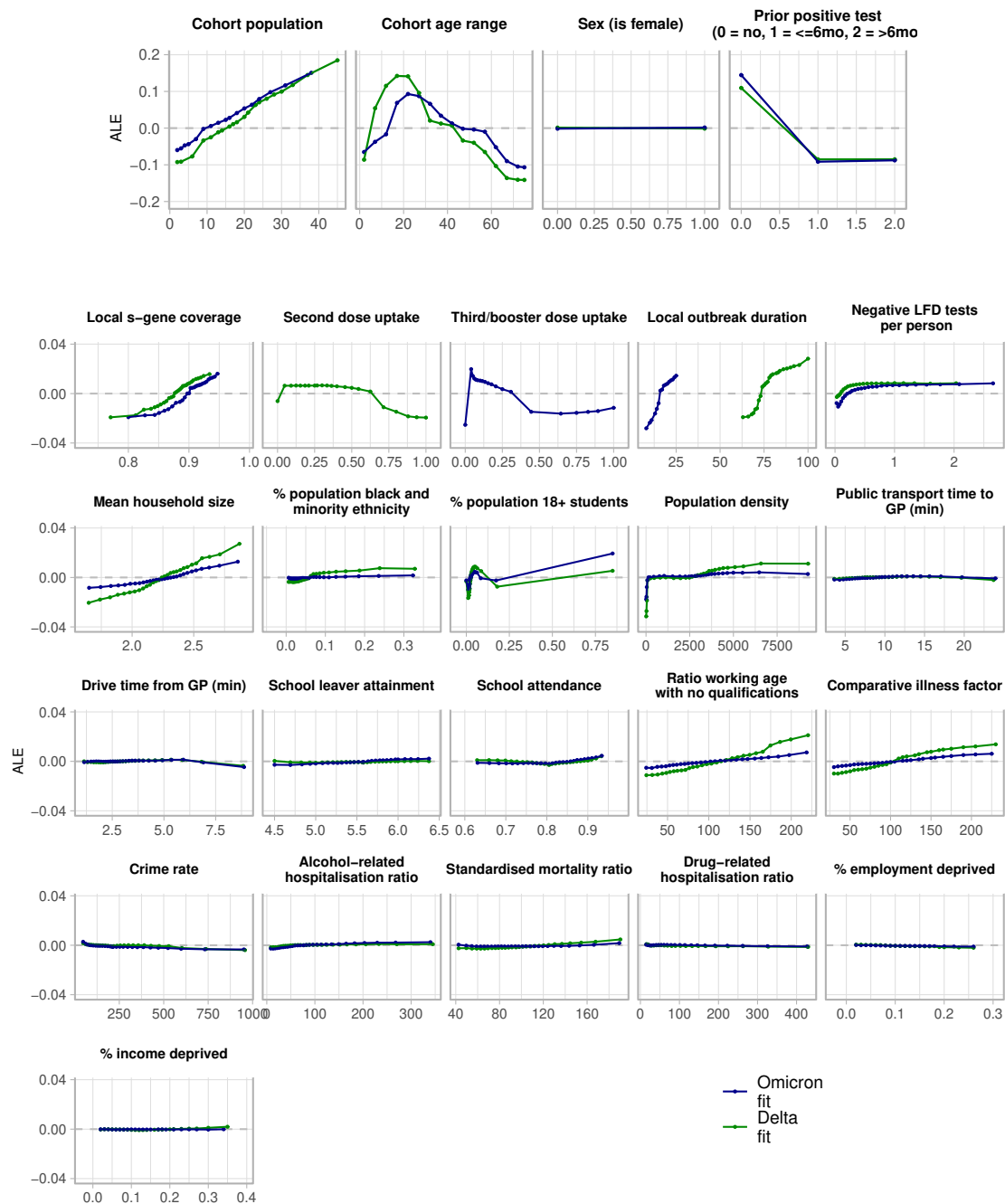


Figure 4: Accumulated local effects across all explanatory variables. For each variable, the  $x$ -axis represents the range of values of that variable in the data, and the  $y$ -axis (note scale differences for *population*, *age*, *sex* and *prior case status*) is the ALE for that variable value. The overall magnitude of the ALE represents the relative size of the effect.

161 evolve over time, especially with respect to different variants.

162 The aim of this study was to compare the pattern of cases across two waves of COVID-19  
163 in Scotland in 2021, during which non pharmaceutical interventions (NPIs) were being relaxed  
164 but testing remained mandatory and a mass vaccination rollout was in progress. We analysed  
165 the distribution of cases during the B.1.617.2 “Delta” wave from May 2021, and the B.1.1.529  
166 “Omicron” wave from November 2021. We have shown that case heterogeneity was associated  
167 with broad factors such as age structure and residual immunity from earlier cases, but also with  
168 factors relating to testing, vaccination, geography and demographics. Despite differences in the  
169 severity of interventions in place, time of year, vaccination uptake and virus phenotype, these risk  
170 factors remain broadly consistent across both waves.

171 Our models accurately capture the case distributions (Fig. 1). However, not all variation is  
172 explained, and residual autocorrelation persists at <5km scales (Fig. 3). A reason for this may be  
173 that our model is not informed by mobility, thus explicit links between communities are not known  
174 to the model. We also do not include meteorological data (such as in e.g. [33]). This could have  
175 explained further variation as our waves occur in different seasons, where the characteristic routes  
176 of transmission may have differed. Last, the fit cases are also time-aggregated, and therefore do  
177 not account for changes in risk factors *during* each wave.

178 The inclusion of the *local outbreak duration* for each DZ (the time the first case was detected  
179 in the DZs wider intermediate zone, typically containing 4-6 DZs) accounts in part for local inter-  
180 actions between neighbouring communities, in the absence of explicit mobility data. A weakness  
181 of this is that the local outbreak duration correlates with the total number of cases, given the  
182 relatively short periods studied. We suspect this is less influential in the Omicron model where  
183 geographical spread was more rapid. The regression models applied here may be better suited  
184 to scenarios where an infectious disease is already well established in the population. For future  
185 analyses on cases at the very beginning of an outbreak with fewer cases, this approach may be  
186 adapted to instead fit case rates per day, from when the first case was identified locally.

## 187 Risk factors

188 We presented the accumulated local effects (Fig. 4), revealing broad indicators for higher or lower  
189 case rates, and how they changed between waves. It is difficult to fully disentangle whether a  
190 difference was caused by a change in control measures, or a change in virus strain. Nonetheless,  
191 our analyses provide some important insights.

192 To begin, high mean household size emerges as a risk factor, consistent with the high secondary  
193 attack rates for SARS-CoV-2 [39, 40], and increased risk of inter-household transmission relative  
194 to contacts outside of the home [41]. That this, and high population density are both stronger risk  
195 factors for Delta may reflect the stronger NPIs at this time increasing the proportion of within-DZ  
196 or within-household transmissions.

197 High vaccine uptake (amongst those eligible) is also protective, more so with Delta, consistent  
198 with higher rates of immune breakthrough with the Omicron variant as compared to Delta [42,  
199 43, 44]. We do not know the specific vaccination status of those in the test data, however, and  
200 linked data may show a stronger protective effect.

201 For Delta, a high proportion of individuals of black and minority ethnicity is a stronger risk  
202 factor. In the UK, this is also a risk factor for severe COVID-19 outcomes [45, 46, 47] but  
203 without detailed, linked data, it is difficult to firmly establish drivers for a *heightened* risk during  
204 the Delta wave. Differences may emerge from known variations in vaccination uptake [48] and  
205 occupation [49] (thus ability to work from home or effectively physical distance), and the relative  
206 impacts of those factors changing across the two waves.

207 Finally, living in a deprived community was suggested from early on [50] and has since also  
208 emerged as a risk factor for *severe* COVID-19 disease [51, 52, 53, 54, 55, 56]. However, the  
209 corresponding ALEs for the variables associated with deprivation are small. Deprivation effects  
210 may be captured by proxy with other variables that correlate with deprivation such as age [57]  
211 and vaccine uptake [58, 59].

## 212 **Testing frequency**

213 The low case rate variation with deprivation (Fig. 1) contrasts with observed inequalities over  
214 severe outcomes [60, 22, 23, 24, 25], suggesting that those living in more deprived communities  
215 experience a higher inherent case-hospitalisation rate. We suspect that a lower proportion of case  
216 *ascertainment*, however, may also be a factor.

217 An important and unique variable in our model is the rate at which *negative* LFD tests were  
218 reported throughout the period. We found high rates of *negative* test reporting to be a *risk* factor.  
219 This suggests a variation in case ascertainment across different demographics, which may in turn  
220 lead to skews in the observed case distribution [61, 62, 22].

221 Further work (Supplementary Material Table S7, Fig. S8) shows that up to February 2023,  
222 the rate of LFD testing and positivity varied substantially across deprivation (quintile 1: 3.6

223 tests/person, 4.61% positive; quintile 5: 6.7 tests/person, 3.57% positive) as well as sex (M: 3.7  
224 tests/person, 4.82% positive; F: 7.0 tests/person 3.30% positive). If demographic differences in  
225 testing behaviour correspond to differences in case ascertainment, the profile of all infections may  
226 then be biased from reported cases, and testing rates may be obscuring the true patterns of  
227 infection over sex and deprivation.

228 In addition, the magnitude of the risk factor (as seen in the ALE, Fig. 4) plateaus beyond a  
229 certain rate ( $>\sim 1$  test/person in each period). This hints at a deeper relationship between true  
230 incidence, the frequency of testing (and whom amongst the population is taking those tests), and  
231 the proportion of infections that are ascertained.

232 Our model is unique in including negative test reporting, and has revealed strong differences  
233 between different demographics that may bias the profile of cases. Beyond the work presented  
234 here, further analysis of reported cases need to be considered with these strong skews in testing  
235 behaviour in mind.

## 236 Conclusion

237 The COVID-19 data studied here are remarkable in terms of volume and resolution, and has  
238 allowed us to assess a national-level epidemic at extremely fine scale. However, regardless of  
239 resolution, cases only partially represent the full underlying pattern of infection. Variations in  
240 testing frequency and known trends in severe outcomes suggest that the distribution of infections  
241 may have been very different to that of reported cases. By incorporating trends on cases, testing  
242 behaviour, and severe outcomes more closely linked to infection (hospitalisation, ICU admission  
243 and mortality), it may be possible to build a much more comprehensive retrospective picture of  
244 how infections were distributed amongst the population.

245 Importantly, while our access to such finely-grained data was exceptional, it can be expected  
246 that such data are likely to become more common in the future, and may become available in  
247 real time. As such, our demonstration of the utility of such data points the way to an impor-  
248 tant approach to improving data analysis supporting control policy response to infectious disease  
249 emergencies in the future.

## 250 4 Data and methods

### 251 4.1 Preparation of case data

252 We use COVID-19 testing data from Public Health Scotland’s *electronic Data Research and Inno-*  
253 *vation Service* (eDRIS) system, dated from July 14<sup>th</sup> 2022. The data include individual tests by  
254 type (polymerase chain reaction (PCR) or rapid lateral flow device (LFD)), test result (positive,  
255 negative, void, inconclusive), test date, S-gene test result if known (positive, dropout, inconclu-  
256 sive), age, sex, and residing data zone (*DZ*, a census area typically comprising 500–1,000 individ-  
257 uals). De-identified IDs link repeat tests by the same individual. We reduce the raw test data to  
258 cases by removing duplicate tests by the same individual within 60 days (taking the date of the first  
259 PCR positive as the case date, or the first LFD in the absence of any PCR). These metadata — in  
260 particular the *DZ*, specifying location to within an area as small as 0.1km<sup>2</sup> in densely populated  
261 areas — therefore identify cases at a fine spatio-temporal scale. Data on vaccine administrations  
262 are also provided by eDRIS.

263 This analysis considers the BA.1 sub-variant of the Omicron lineage only. The sub-variant  
264 BA.2/B.1.1.529.2 later replaced BA.1, becoming dominant in Scotland from around 25<sup>th</sup> February  
265 2022. This variant, like Delta, has an S-gene positive test signature. However by the end of the  
266 period studied the BA.2 variant was only being identified in fewer than 1% of fully sequenced  
267 cases in the UK [63], and here we assume all remaining S-gene positive cases to be Delta.

268 Prior to January 6<sup>th</sup> 2022 in Scotland, positive LFD tests (typically taken at home) required  
269 PCR confirmation. Approximately 90% of cases in this period have a definitive S-gene result. A  
270 policy change then dropped this PCR requirement [64], after which cases with S-gene results fell  
271 to about 50% by February 2022 (per eDRIS data).

272 For Omicron cases, we gather from the data S-gene dropout cases between 15<sup>th</sup> November 2021  
273 and 6<sup>th</sup> January 2022, and for the Delta outbreak, S-gene positive cases between 1<sup>st</sup> May and 7<sup>th</sup>  
274 September 2021 (choosing this end date to have a similar number of cases in each set). We exclude  
275 cases that have a different, or no S-gene result.

276 Using the linked historical tests, we label cases based on whether the individual had either:  
277 never tested positive before; had tested positive in the last six months prior to the start of that  
278 wave, or; last tested positive over six months prior to the start of that wave. We denote this the  
279 *prior case status*, as a proxy for infection-based immunity.

280 Finally to prepare the cases data to be fit, we group individuals that have the same age range,

281 sex, residing datazone, and *prior case status*, terming these subsets of individuals *cohorts*. As  
282 an illustrative example, a cohort may be a population of 38 males aged between 50–54 residing  
283 in a given datazone “X”, that have never tested positive for COVID-19 before, among whom 9  
284 Omicron COVID-19 cases were identified. This is the highest practical resolution we can achieve  
285 using the eDRIS case data, and our model (Section 4.3) fits case counts at this resolution.

## 286 4.2 Time series analysis

### 287 Time-dependent reproduction number

The time-dependent reproduction number  $R_i$  is the average number of forward infections caused by a person infected on day  $t_i$ . Define  $n_j$  as the number of new infections on day  $t_j$ . These new infections came from individuals infected on days on, or prior to  $t_j$ . Define  $A_{ij}$  as the number of new infections on day  $t_j$  *specifically* from those infected on day  $t_i \leq t_j$ :

$$A_{ij} = \frac{(n_i - \delta_{ij})P(t_j - t_i)}{\sum_{i' \leq j} (n_{i'} - \delta_{i'j})P(t_j - t_{i'})} n_j .$$

$P(\Delta t)$  is the probability of an individual passing on the infection,  $\Delta t$  days after being infected. The presence of the Kronecker delta  $\delta_{ij}$  excludes the possibility of infected individuals infecting themselves. The reproduction number  $R_i$  is then the average total of infections generated over all subsequent days [65]:

$$R_i = \frac{1}{n_i} \sum_{j \geq i} A_{ij} = \frac{1}{n_i} \sum_{j \geq i} \frac{n_j (n_i - \delta_{ij}) P(t_j - t_i)}{\sum_{i' \leq j} (n_{i'} - \delta_{i'j}) P(t_j - t_{i'})} .$$

We take  $P(\Delta t)$  to be

$$P(\Delta t) \sim e^{-\lambda \Delta t}$$

288 with  $\lambda^{-1}$  the mean infectious period. Individuals are equally infectious throughout the entire  
289 infection. In our calculations we estimate  $1/\lambda = 6.26$  days, using the posterior mean duration  
290 of infectiousness obtained from the *SCoVMod* compartmental model (for more detail see Refer-  
291 ence [56]).

292 As we estimate the infection reproduction number using the cases data, we implicitly assume  
293 that case ascertainment does not change over time, and does not account for the delay between  
294 infection, and registering a case.

295 In this work the reproductive number is measured at local authority level, the level at which

296 the Scottish Government monitored and adjusted NPIs.

### 297 **Case doubling time**

At the start of each wave we assume exponential growth of cases:

$$\text{new cases} \propto e^{rt}$$

where the gradient of a linear regression on  $\log(\text{new cases})$  against  $t$  returns the growth rate  $r$ .

The evolution of new cases can also be rewritten in terms of a doubling time  $t_D$ :

$$\text{new cases} \propto 2^{t/t_D}$$

298 where  $t_D = \frac{\log 2}{r}$ .

## 299 **4.3 Model**

300 Our statistical model is designed to explain variation in COVID-19 case numbers as prepared in  
301 Section 4.1, and identify risk factors amongst a broad range of variables, using random forest  
302 regression. We fit models to the distribution of Delta and Omicron cases respectively, allowing for  
303 comparison of risk factors across the two waves.

### 304 **Explanatory variables**

305 We include demographic factors (population, age, sex, ethnicity, student population), COVID-19  
306 related factors (testing volume, prior case status, vaccination uptake), geography (local population  
307 density and transport time to public services to serve as proxies for connectivity and geographic  
308 remoteness), as well as deprivation. Data on deprivation are taken from the *Scottish Indices of*  
309 *Multiple Deprivation* (SIMD) [66]. The SIMD ranks DZs in Scotland by “multiple” deprivation,  
310 incorporating measures relating to local health, housing, geographic access, employment, income,  
311 crime, and education. In our model we use the raw measures of deprivation as explanatory  
312 variables. To account for local spread of infection between neighbourhoods that are geographically  
313 close to one another, we include an *local outbreak duration* parameter, which specifies the date at  
314 which the *first* case of the variant was identified at the *intermediate zone* (IZ, an administrative  
315 area containing of order 4–6 DZs).

316 A comprehensive description of all individual variables used in given in Supplementary Mate-  
317 rial, Section [B.2](#).

### 318 **Random forest model**

319 We use random forest regression [\[67\]](#) on the distribution of COVID-19 cases, as it allows us to  
320 fit the distribution without specifying any prior analytical relation between the outcome variable  
321 (cases) and any of the explanatory variables, which may themselves be correlated. We fit the  
322 time-aggregated case distribution in *R* (version 4.1.0) [\[68\]](#), using the *randomForest* package [\[69\]](#)  
323 (version 4.6-14).

324 We fit the outcome variable  $\sqrt{\text{cases} + 1}$  at cohort level (with a *cohort* defined in Section [4.1](#)).  
325 The fit number of cases at other scales (such as DZ level) is then an aggregation of cases from  
326 their constituent cohorts.

327 We extract two metrics for variable importance from the *randomForest* function output: the  
328 node purity (a measure of how effective variables are at partitioning cohorts with differing numbers  
329 of cases in the tree), and the loss of model accuracy on effective removal of that variable from the  
330 model.

331 Model hyperparameters were chosen manually so as to maximise the variance explained by  
332 a subset of the data not used to fit the model. Full hyperparameter specification is included in  
333 Supplementary Material, Section [B.1](#). The model specifications for fitting the Omicron and Delta  
334 waves are identical with one exception: for the Omicron model, third/booster dose uptake is used,  
335 whereas for Delta, second dose uptake is used (third/booster doses were only administered later;  
336 see Supplementary Material, Section [B.3](#) for further details).

337 In addition to the full model, we fit for each of Omicron and Delta three “reduced” models,  
338 under equivalent hyperparameters to the full model and the same cohort structure, but informed  
339 only by population, and one of: age; the relative deprivation of the residing DZ, as defined by  
340 the overall SIMD deprivation *rank* [\[70\]](#), and; population density. These outputs illustrate how  
341 effective these variables are at alone at explaining case variation, relative to our full model.

## 342 **4.4 Accumulated local effects**

343 To identify risk factors amongst the explanatory variables used to inform the model, we calculate  
344 the *accumulated local effects* (ALEs) of each variable. The ALEs describe how the model fit value  
345 changes, in response to changing one variable value in isolation, averaged over many different



346 entries in the data [71]. In this context, ALEs indicate whether a variable value is associated with  
347 fewer or more cases in general over the data. If the ALE is greater than zero, the fit cases generally  
348 increases given that variable value.

## 349 4.5 Moran's I autocorrelation statistic

To probe geographical variation in cases *not* explained by the model, we measure the Moran's I autocorrelation [72, 73] on the residuals (the difference between the data and fit value), relating to their physical location. We compare local DZ-aggregated residuals over physical distances (from 1–100km), as well as network distance (number of nearest neighbours apart). For a set of  $N$  residuals  $y_i$ , the Moran's I is a measure of autocorrelation:

$$I = \frac{N}{\sum_{i=1}^N \sum_{j=1}^N w_{ij}} \frac{\sum_{i=1}^N \sum_{j=1}^N w_{ij} (y_i - \bar{y})(y_j - \bar{y})}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

350 with  $\bar{y}$  the mean of all residuals, and  $w_{i,j}$  is an associated *weight* of the pair of observations  
351  $(i, j)$ , with  $w_{i,i} = 0$ . To measure the autocorrelation between residuals within a separation  $d$   
352 (either a physical or network-based distance) of one another, we set  $w_{i,j} = 1$  if  $\text{dist}(i, j) \leq d$ ,  
353 and 0 otherwise. Fully correlated residuals would have  $I = 1$ , whereas  $I = 0$  would indicate no  
354 correlation.

355 This measure characterises how effective our models are at explaining geographical variation,  
356 and with different distances  $d$  shows over what length scales residual autocorrelation persists.

## 357 5 Acknowledgements

358 We thank Public Health Scotland's *electronic Data Research and Innovation Service* (eDRIS) for  
359 the provision of COVID-19 testing, vaccination and severe outcomes data. We also thank the  
360 reviewers for their feedback and suggestions, which has led to improvement of the article.

## 361 6 Author contributions

362 R.R.K. conceived the project. A.J.W. wrote the model code and performed the case distribution  
363 analysis. A.R.S. and P.R.B. performed the temporal analysis. R.W. performed analysis of lateral  
364 flow testing data. A.J.W. and R.R.K. wrote the manuscript.

## 365 7 Competing interests

366 The authors declare no competing interests.

## 367 8 Code availability

368 Analysis code is available at

369 <https://git.ecdf.ed.ac.uk/awood310/scotland-covid-case-distribution-random-forest-model>.

## 370 9 Data availability

371 The COVID-19 testing, vaccination and hospitalisation data utilised in this work are not publicly  
372 available. They are provided to the authors for academic research by Public Health Scotland's  
373 *electronic Data Research and Innovation Service*, under a data sharing agreement (*Spatial and*  
374 *Network Analysis of SARS-CoV-2 Sequences to Inform COVID-19 Control in Scotland*, and can  
375 be contacted via [phs.edris@phs.scot](mailto:phs.edris@phs.scot). The authors received no special privileges with respect to  
376 data access as compared to other researchers.

377 Data relating to deprivation are derived from the 2020 Scottish Index of Multiple Deprivation,  
378 and are publicly available (<https://simd.scot>). Fine-scale population statistics are drawn from  
379 publicly-available population estimates as of mid-2020 [74]. The population of students and those  
380 belonging to a minority ethnicity are drawn from public Scottish census data (tables KS501SC,  
381 LC2101SC respectively).

## 382 10 Funding statement

383 This work has been funded by the ESRC grant *Real-time monitoring and predictive modelling of*  
384 *the impact of human behaviour and vaccine characteristics on COVID-19 vaccination in Scotland*  
385 (ES/W001489/1). This work has also been funded the BBSRC Institute Strategic Programme  
386 grant to the Roslin Institute (BB/J004235/1).

## 387 References

388 [1] Office for National Statistics. Coronavirus (COVID-19) Infection Sur-  
389 vey: Scotland Dataset;. Available from <https://www.ons.gov.uk/>

- 390 [peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/](#)  
391 [datasets/covid19infectionsurveyscotland](#) (last accessed 15/08/2023).
- 392 [2] Office for National Statistics. Coronavirus (COVID-19) latest insights: Hospitals;. Available  
393 from [https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/](https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/articles/coronaviruscovid19latestinsights/hospitals)  
394 [conditionsanddiseases/articles/coronaviruscovid19latestinsights/hospitals](#)  
395 (last accessed 16/08/2023).
- 396 [3] Simpson CR, Robertson C, Vasileiou E, McMenamin J, Gunson R, Ritchie LD, et al. Early  
397 pandemic evaluation and enhanced surveillance of COVID-19 (EAVE II): protocol for an  
398 observational study using linked Scottish national data. *BMJ open*. 2020;10(6):e039097.
- 399 [4] Sheikh A, Kerr S, Woolhouse M, McMenamin J, Robertson C. Severity of Omicron variant of  
400 concern and vaccine effectiveness against symptomatic disease: national cohort with nested  
401 test negative design study in Scotland. 2021;.
- 402 [5] Canas LS, Sudre CH, Pujol JC, Polidori L, Murray B, Molteni E, et al. Early detection of  
403 COVID-19 in the UK using self-reported symptoms: a large-scale, prospective, epidemiolog-  
404 ical surveillance study. *The Lancet Digital Health*. 2021;3(9):e587–e598.
- 405 [6] Antonelli M, Penfold RS, Merino J, Sudre CH, Molteni E, Berry S, et al. Risk factors and  
406 disease profile of post-vaccination SARS-CoV-2 infection in UK users of the COVID Symptom  
407 Study app: a prospective, community-based, nested, case-control study. *The Lancet Infectious*  
408 *Diseases*. 2022;22(1):43–55.
- 409 [7] The Scottish Government. Coronavirus (COVID-19) confirmed in Scotland;. Available from  
410 <https://www.gov.scot/news/coronavirus-covid-19/> (last accessed 16/08/2023).
- 411 [8] The Scottish Government. Effective ‘lockdown’ to be introduced;. Available from  
412 <https://www.gov.scot/news/effective-lockdown-to-be-introduced/> (last accessed  
413 16/08/2023).
- 414 [9] The Scottish Government. Coronavirus (COVID-19): protection levels — re-  
415 views and evidence;. Available from [https://www.gov.scot/collections/](https://www.gov.scot/collections/coronavirus-covid-19-protection-levels-reviews-and-evidence/)  
416 [coronavirus-covid-19-protection-levels-reviews-and-evidence/](#) (last accessed  
417 16/08/2023).

- 418 [10] The Scottish Government. New guidance issued for the festive period;. Available from  
419 <https://www.gov.scot/news/new-guidance-issued-for-the-festive-period/> (last ac-  
420 cessed 16/08/2023).
- 421 [11] The Scottish Government. Scotland in Lockdown;. Available from [https://www.gov.scot/  
422 news/scotland-in-lockdown/](https://www.gov.scot/news/scotland-in-lockdown/) (last accessed 16/08/2023).
- 423 [12] The Scottish Government. First COVID-19 vaccinations in Scot-  
424 land take place;. Available from [https://www.gov.scot/news/  
425 first-covid-19-vaccinations-in-scotland-take-place/](https://www.gov.scot/news/first-covid-19-vaccinations-in-scotland-take-place/) (last accessed 16/08/2023).
- 426 [13] The Scottish Government. Coronavirus (COVID-19): vaccine deploy-  
427 ment plan 2021;. Available from [https://www.gov.scot/publications/  
428 coronavirus-covid-19-vaccine-deployment-plan-2021/](https://www.gov.scot/publications/coronavirus-covid-19-vaccine-deployment-plan-2021/) (last accessed 15/08/2023).
- 429 [14] Hale T, Angrist N, Kira B, Petherick A, Phillips T, Webster S. Variation in government  
430 responses to COVID-19. 2020;.
- 431 [15] McMillen T, Jani K, Robilotti EV, Kamboj M, Babady NE. The spike gene target failure  
432 (SGTF) genomic signature is highly accurate for the identification of Alpha and Omicron  
433 SARS-CoV-2 variants. *Scientific reports*. 2022;12(1):18968.
- 434 [16] Gebhard C, Regitz-Zagrosek V, Neuhauser HK, Morgan R, Klein SL. Impact of sex and  
435 gender on COVID-19 outcomes in Europe. *Biology of sex differences*. 2020;11:1–13.
- 436 [17] Galbadage T, Peterson BM, Awada J, Buck AS, Ramirez DA, Wilson J, et al. Systematic  
437 review and meta-analysis of sex-specific COVID-19 clinical outcomes. *Frontiers in medicine*.  
438 2020;7:348.
- 439 [18] Peckham H, de Gruijter NM, Raine C, Radziszewska A, Ciurtin C, Wedderburn LR, et al.  
440 Male sex identified by global COVID-19 meta-analysis as a risk factor for death and ICU  
441 admission. *Nature communications*. 2020;11(1):6317.
- 442 [19] Sartorius B, Lawson A, Pullan R. Modelling and predicting the spatio-temporal spread of  
443 COVID-19, associated deaths and impact of key risk factors in England. *Scientific reports*.  
444 2021;11(1):5378.
- 445 [20] Diao Y, Koder S, Anzai D, Gomez-Tames J, Rashed EA, Hirata A. Influence of population  
446 density, temperature, and absolute humidity on spread and decay durations of COVID-19:

- 447 A comparative study of scenarios in China, England, Germany, and Japan. *One Health*.  
448 2021;12:100203.
- 449 [21] Smith TP, Flaxman S, Gallinat AS, Kinosian SP, Stemkovski M, Unwin HJT, et al.  
450 Temperature and population density influence SARS-CoV-2 transmission in the absence  
451 of nonpharmaceutical interventions. *Proceedings of the National Academy of Sciences*.  
452 2021;118(25):e2019284118.
- 453 [22] Green MA, García-Fiñana M, Barr B, Burnside G, Cheyne CP, Hughes D, et al. Evaluating  
454 social and spatial inequalities of large scale rapid lateral flow SARS-CoV-2 antigen testing  
455 in COVID-19 management: An observational study of Liverpool, UK (November 2020 to  
456 January 2021). *The Lancet Regional Health-Europe*. 2021;6:100107.
- 457 [23] Meurisse M, Lajot A, Devleeschauwer B, Van Cauteren D, Van Oyen H, Van den Borre L,  
458 et al. The association between area deprivation and COVID-19 incidence: a municipality-level  
459 spatio-temporal study in Belgium, 2020–2021. *Archives of Public Health*. 2022;80(1):1–10.
- 460 [24] KC M, Oral E, Straif-Bourgeois S, Rung AL, Peters ES. The effect of area deprivation on  
461 COVID-19 risk in Louisiana. *PLoS One*. 2020;15(12):e0243028.
- 462 [25] Badr HS, Du H, Marshall M, Dong E, Squire MM, Gardner LM. Association between mobility  
463 patterns and COVID-19 transmission in the USA: a mathematical modelling study. *The*  
464 *Lancet Infectious Diseases*. 2020;20(11):1247–1254.
- 465 [26] Reuter M, Rigó M, Formazin M, Liebers F, Latza U, Castell S, et al. Occupation and SARS-  
466 CoV-2 infection risk among 108 960 workers during the first pandemic wave in Germany.  
467 *Scandinavian Journal of Work, Environment & Health*. 2022;48(6):446.
- 468 [27] Rhodes S, Wilkinson J, Pearce N, Mueller W, Cherrie M, Stocking K, et al. Occupational  
469 differences in SARS-CoV-2 infection: analysis of the UK ONS COVID-19 infection survey. *J*  
470 *Epidemiol Community Health*. 2022;76(10):841–846.
- 471 [28] Zhang M. Estimation of differential occupational risk of COVID-19 by comparing risk factors  
472 with case data by occupational group. *American journal of industrial medicine*. 2021;64(1):39–  
473 47.
- 474 [29] Chadeau-Hyam M, Bodinier B, Elliott J, Whitaker MD, Tzoulaki I, Vermeulen R, et al. Risk  
475 factors for positive and negative COVID-19 tests: a cautious and in-depth analysis of UK  
476 biobank data. *International journal of epidemiology*. 2020;49(5):1454–1467.

- 477 [30] Lau MS, Grenfell B, Thomas M, Bryan M, Nelson K, Lopman B. Characterizing superspread-  
478 ing events and age-specific infectiousness of SARS-CoV-2 transmission in Georgia, USA. Pro-  
479 ceedings of the National Academy of Sciences. 2020;117(36):22430–22435.
- 480 [31] Working group for the surveillance, control of COVID-19 in Spain, group for the surveillance  
481 W, control of COVID-19 in Spain, Redondo-Bravo L, Sierra Moros MJ, et al. The first wave  
482 of the COVID-19 pandemic in Spain: characterisation of cases and risk factors for severe  
483 outcomes, as at 27 April 2020. *Eurosurveillance*. 2020;25(50):2001431.
- 484 [32] Hu T, Wang S, She B, Zhang M, Huang X, Cui Y, et al. Human mobility data in the COVID-  
485 19 pandemic: characteristics, applications, and challenges. *International Journal of Digital*  
486 *Earth*. 2021;14(9):1126–1147.
- 487 [33] Ledebur K, Kaleta M, Chen J, Lindner SD, Matzhold C, Weidle F, et al. Meteorological  
488 factors and non-pharmaceutical interventions explain local differences in the spread of SARS-  
489 CoV-2 in Austria. *PLoS computational biology*. 2022;18(4):e1009973.
- 490 [34] Jia JS, Lu X, Yuan Y, Xu G, Jia J, Christakis NA. Population flow drives spatio-temporal  
491 distribution of COVID-19 in China. *Nature*. 2020;582(7812):389–394.
- 492 [35] Wang H, Ghosh A, Ding J, Sarkar R, Gao J. Heterogeneous interventions reduce the spread  
493 of COVID-19 in simulations on real mobility data. *Scientific reports*. 2021;11(1):7809.
- 494 [36] Hou X, Gao S, Li Q, Kang Y, Chen N, Chen K, et al. Intracounty modeling of COVID-19  
495 infection with human mobility: Assessing spatial heterogeneity with business traffic, age, and  
496 race. *Proceedings of the National Academy of Sciences*. 2021;118(24):e2020524118.
- 497 [37] Asem N, Ramadan A, Hassany M, Ghazy RM, Abdallah M, Ibrahim M, et al. Pattern and  
498 determinants of COVID-19 infection and mortality across countries: An ecological study.  
499 *Heliyon*. 2021;7(7).
- 500 [38] Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. Nextstrain: real-  
501 time tracking of pathogen evolution. *Bioinformatics*. 2018;34(23):4121–4123.
- 502 [39] Jalali N, Brustad HK, Frigessi A, MacDonald EA, Meijerink H, Feruglio SL, et al. Increased  
503 household transmission and immune escape of the SARS-CoV-2 Omicron variant compared to  
504 the Delta variant: evidence from Norwegian contact tracing and vaccination data. *medRxiv*.  
505 2022;.

- 506 [40] Fonager J, Bennedbæk M, Bager P, Wohlfahrt J, Ellegaard KM, Ingham AC, et al. Molecular epidemiology of the SARS-CoV-2 variant Omicron BA. 2 sub-lineage in Denmark, 29  
507 November 2021 to 2 January 2022. *Eurosurveillance*. 2022;27(10):2200181.
- 509 [41] Dupraz J, Butty A, Duperrex O, Estoppey S, Faivre V, Thabard J, et al. Prevalence of  
510 SARS-CoV-2 in household members and other close contacts of COVID-19 cases: a serologic  
511 study in canton of Vaud, Switzerland. In: *Open forum infectious diseases*. vol. 8. Oxford  
512 University Press US; 2021. p. ofab149.
- 513 [42] Andrews N, Stowe J, Kirsebom F, Toffa S, Rickeard T, Gallagher E, et al. Covid-19 vaccine  
514 effectiveness against the omicron (B. 1.1. 529) variant. *New England Journal of Medicine*.  
515 2022;.
- 516 [43] Cele S, Jackson L, Khoury DS, Khan K, Moyo-Gwete T, Tegally H, et al. Omicron extensively  
517 but incompletely escapes Pfizer BNT162b2 neutralization. *Nature*. 2022;602(7898):654–656.
- 518 [44] Vasileiou E, Simpson CR, Shi T, Kerr S, Agrawal U, Akbari A, et al. Interim findings  
519 from first-dose mass COVID-19 vaccination roll-out and COVID-19 hospital admissions in  
520 Scotland: a national prospective cohort study. *The Lancet*. 2021;397(10285):1646–1657.
- 521 [45] Office for National Statistics. Updating ethnic contrasts in deaths in-  
522 volving the coronavirus (COVID-19), England: 8 December 2020  
523 to 1 December 2021;. Available from: [https://www.ons.gov.uk/  
524 peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/articles/  
525 updatingethniccontrastsindeathsinvolvingthecoronaviruscovid19englandandwales/  
526 8december2020to1december2021](https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/articles/updatingethniccontrastsindeathsinvolvingthecoronaviruscovid19englandandwales/8december2020to1december2021) (last accessed 15/08/2023).
- 527 [46] Platt L, Warwick R. Are some ethnic groups more vulnerable to COVID-19 than others.  
528 *Institute for fiscal studies*. 2020;1(05):2020.
- 529 [47] Lo CH, Nguyen LH, Drew DA, Warner ET, Joshi AD, Graham MS, et al. Race, ethnicity,  
530 community-level socioeconomic factors, and risk of COVID-19 in the United States and the  
531 United Kingdom. *EClinicalMedicine*. 2021;38.
- 532 [48] Office for National Statistics. Coronavirus and vaccination rates in people aged 18  
533 years and over by socio-demographic characteristic and occupation, England: 8  
534 December 2020 to 31 December 2021;. Available from: <https://www.ons.gov.uk>.

- 535 [uk/peoplepopulationandcommunity/healthandsocialcare/healthinequalities/](#)  
536 [bulletins/coronavirusandvaccinationratesinpeopleaged18yearsandoverbysociodemographiccharacterist](#)  
537 [8december2020to31december2021](#) (last accessed 15/08/2023).
- 538 [49] National Records of Scotland. Census 2011: Release 3I - Detailed characteristics on Labour  
539 Market and Education in Scotland;. Available from: [https://www.nrscotland.gov.uk/](https://www.nrscotland.gov.uk/news/2014/census-2011-release-3i)  
540 [news/2014/census-2011-release-3i](#) (last accessed 15/08/2023).
- 541 [50] Khalatbari-Soltani S, Cumming RC, Delpierre C, Kelly-Irving M. Importance of collecting  
542 data on socioeconomic determinants from the early stage of the COVID-19 outbreak onwards.  
543 *J Epidemiol Community Health*. 2020;74(8):620–623.
- 544 [51] Lone NI, McPeake J, Stewart NI, Blayney MC, Seem RC, Donaldson L, et al. Influence of  
545 socioeconomic deprivation on interventions and outcomes for patients admitted with COVID-  
546 19 to critical care units in Scotland: a national cohort study. *The Lancet Regional Health-*  
547 *Europe*. 2021;1:100005.
- 548 [52] Blundell R, Costa Dias M, Joyce R, Xu X. COVID-19 and Inequalities. *Fiscal studies*.  
549 2020;41(2):291–319.
- 550 [53] Bambra C, Riordan R, Ford J, Matthews F. The COVID-19 pandemic and health inequalities.  
551 *J Epidemiol Community Health*. 2020;74(11):964–968.
- 552 [54] Baena-Díez JM, Barroso M, Cordeiro-Coelho SI, Díaz JL, Grau M. Impact of COVID-19 out-  
553 break by income: hitting hardest the most deprived. *Journal of Public Health*. 2020;42(4):698–  
554 703.
- 555 [55] McGurnaghan SJ, Weir A, Bishop J, Kennedy S, Blackburn LA, McAllister DA, et al. Risks  
556 of and risk factors for COVID-19 disease in people with diabetes: a cohort study of the total  
557 population of Scotland. *The lancet Diabetes & endocrinology*. 2021;9(2):82–93.
- 558 [56] Banks CJ, Colman E, Doherty T, Tearne O, Arnold M, Atkins KE, et al. SCoVMod—a spa-  
559 tially explicit mobility and deprivation adjusted model of first wave COVID-19 transmission  
560 dynamics. *Wellcome Open Research*. 2022;7(161):161.
- 561 [57] National Records of Scotland. Mid-2021 Small Area Population Estimates, Scot-  
562 land (Report);. Available from [https://www.nrscotland.gov.uk/files//statistics/](https://www.nrscotland.gov.uk/files//statistics/population-estimates/sape-2021/sape-21-report.pdf)  
563 [population-estimates/sape-2021/sape-21-report.pdf](#) (last accessed 15/08/2023).



- 564 [58] Office for National Statistics. Coronavirus (COVID-19) Infection Sur-  
565 vey technical article: Analysis of characteristics associated with  
566 vaccination uptake;. Available from [https://www.ons.gov.uk/  
567 peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/  
568 articles/coronaviruscovid19infectionsurveytechnicalarticleanalysisofcharacteristicsassociatedwith  
569 2021-11-15](https://www.ons.gov.uk/peoplepopulationandcommunity/healthandsocialcare/conditionsanddiseases/articles/coronaviruscovid19infectionsurveytechnicalarticleanalysisofcharacteristicsassociatedwith2021-11-15) (last accessed 15/08/2023).
- 570 [59] Wood AJ, MacKintosh AM, Stead M, Kao RR. Predicting future spatial patterns in COVID-  
571 19 booster vaccine uptake. medRxiv. 2022;p. 2022–08.
- 572 [60] Wood AJ, Kao RR. Empirical distributions of time intervals between COVID-19 cases and  
573 more severe outcomes in Scotland. PloS one. 2023;18(8):e0287397.
- 574 [61] Colman E, Puspitarani GA, Enright J, Kao RR. Ascertainment rate of SARS-CoV-2 in-  
575 fections from healthcare and community testing in the UK. Journal of Theoretical Biology.  
576 2022;p. 111333. Available from: [https://www.sciencedirect.com/science/article/pii/  
577 S0022519322003241](https://www.sciencedirect.com/science/article/pii/S0022519322003241).
- 578 [62] Nightingale ES, Abbott S, Russell TW. The local burden of disease during the first wave of  
579 the COVID-19 epidemic in England: estimation using different data sources from changing  
580 surveillance practices. BMC public health. 2022;22(1):1–14.
- 581 [63] The UK Health Security Agency. SARS-CoV-2 variants of concern and vari-  
582 ants under investigation in England: technical briefing 35;. Available from  
583 [https://assets.publishing.service.gov.uk/government/uploads/system/uploads/  
584 attachment\\_data/file/1050999/Technical-Briefing-35-28January2022.pdf](https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/1050999/Technical-Briefing-35-28January2022.pdf) (last  
585 accessed 15/08/2023).
- 586 [64] The Scottish Government. Self-Isolation and testing changes;. Available from [https://www.  
587 gov.scot/news/self-isolation-and-testing-changes/](https://www.gov.scot/news/self-isolation-and-testing-changes/) (last accessed 15/08/2023).
- 588 [65] Wallinga J, Teunis P. Different epidemic curves for severe acute respiratory syndrome reveal  
589 similar impacts of control measures. American Journal of epidemiology. 2004;160(6):509–516.
- 590 [66] The Scottish Government. SIMD 2020 Technical Notes;. Available from [https://www.gov.  
591 scot/publications/simd-2020-technical-notes/](https://www.gov.scot/publications/simd-2020-technical-notes/) (last accessed 15/08/2023).
- 592 [67] Breiman L. Random forests. Machine learning. 2001;45(1):5–32.

- 593 [68] R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria;  
594 2022. Available from: <https://www.R-project.org/>.
- 595 [69] Liaw A, Wiener M, et al. Classification and regression by randomForest. R news. 2002;2(3):18–  
596 22.
- 597 [70] The Scottish Government. Scottish Index of Multiple Depriva-  
598 tion 2020;. Available from [https://www.gov.scot/publications/  
599 scottish-index-of-multiple-deprivation-2020v2-indicator-data/](https://www.gov.scot/publications/scottish-index-of-multiple-deprivation-2020v2-indicator-data/) (last accessed  
600 15/08/2023).
- 601 [71] Apley D, Apley MD. Package ‘ALEPlot’. 2018;.
- 602 [72] Moran PA. Notes on continuous stochastic phenomena. Biometrika. 1950;37(1/2):17–23.
- 603 [73] Gittleman JL, Kot M. Adaptation: statistics and a null model for estimating phylogenetic  
604 effects. Systematic Zoology. 1990;39(3):227–241.
- 605 [74] National Records of Scotland. Mid-2020 Small Area Population Estimates for 2011 Data  
606 Zones;. Available from [https://www.nrscotland.gov.uk/statistics-and-data/  
607 statistics/statistics-by-theme/population/population-estimates/  
608 small-area-population-estimates-2011-data-zone-based/mid-2020/](https://www.nrscotland.gov.uk/statistics-and-data/statistics/statistics-by-theme/population/population-estimates/small-area-population-estimates-2011-data-zone-based/mid-2020/) (last accessed  
609 15/08/2023).
- 610 [75] The Scottish Government. Major milestone in vaccination programme;. Available  
611 from <https://www.gov.scot/news/major-milestone-in-vaccination-programme/> (last  
612 accessed 15/08/2023).
- 613 [76] The Cabinet Secretary for Health, Care S. Scotland’s autumn/winter vacci-  
614 nation strategy 2021;. Available from [https://www.gov.scot/publications/  
615 scotlands-autumn-winter-vaccination-strategy-2021/](https://www.gov.scot/publications/scotlands-autumn-winter-vaccination-strategy-2021/) (last accessed 15/08/2023).

616 **Supplementary material**

617 *Assessing the importance of demographic risk factors across two waves*  
618 *of SARS-CoV-2 using fine-scale case data*

619

620 A.J. Wood, A.R. Sanchez, P.R. Bessell, R. Wightman, R.R. Kao

621 **A** Supplementary plots for time evolution of cases

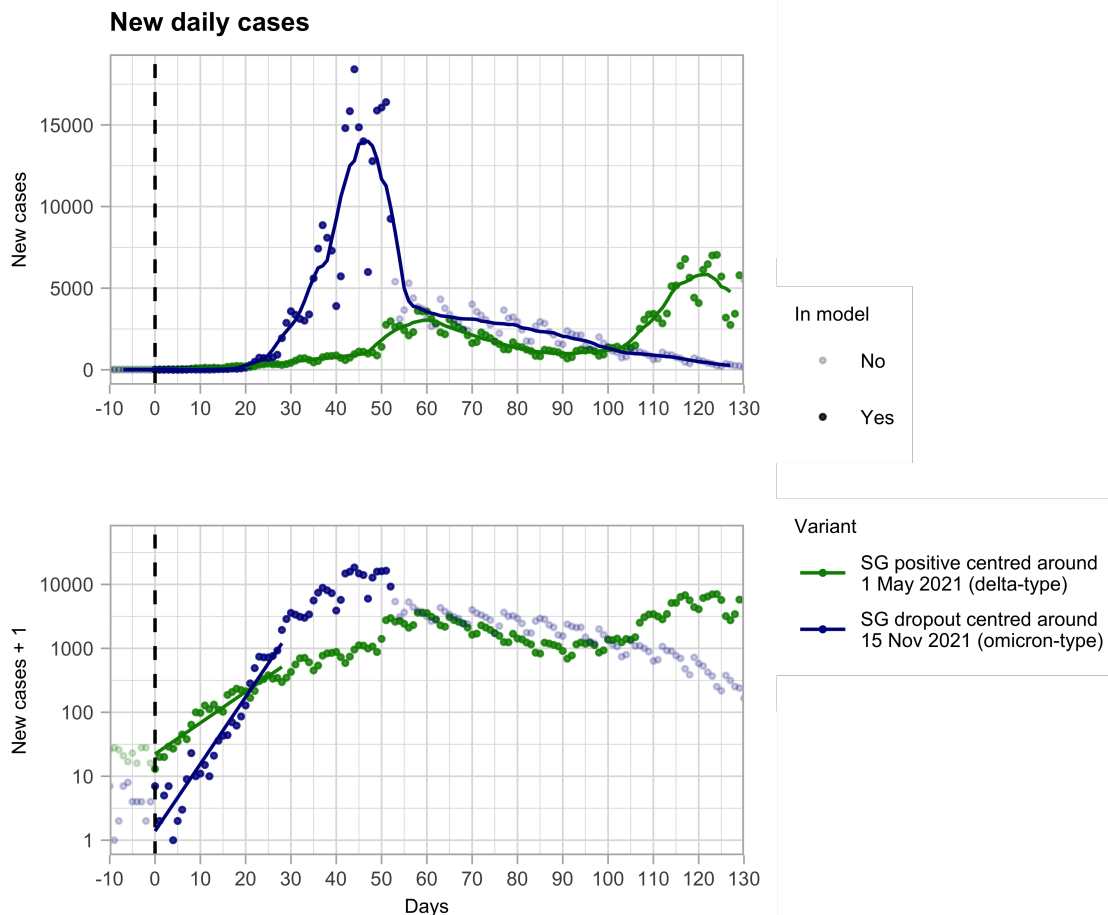


Figure S1: Timeseries of the initial outbreaks of the Delta and Omicron variants in terms of newly reported cases. The gradient of the linear regression (straight line) of the early trajectory of  $\log(\text{new cases} + 1)$  is inversely proportional to the case doubling time.

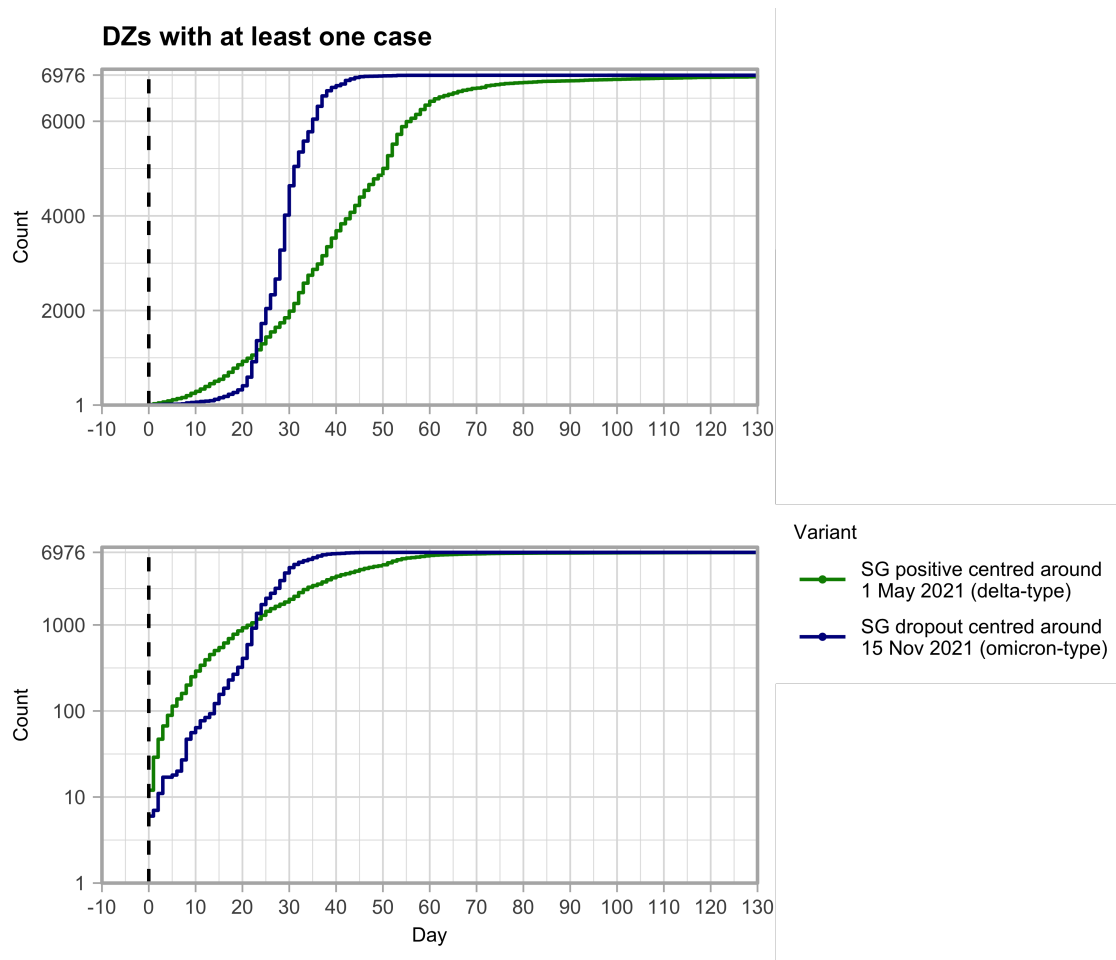


Figure S2: Timeseries of the initial outbreaks of the Delta and Omicron variants in terms of the cumulative number of DZs to have reported at least one case associated with the variant of interest.

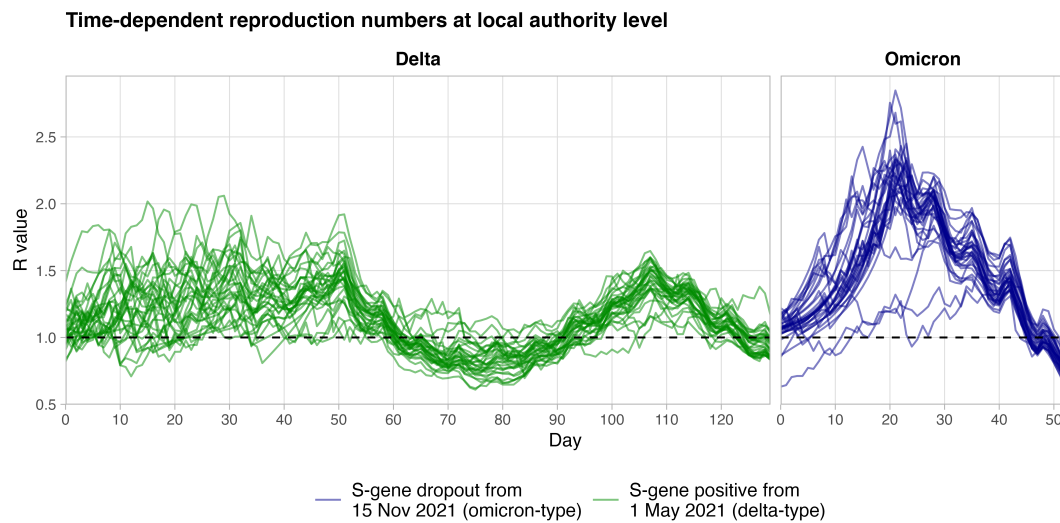


Figure S3: Time-dependent reproduction numbers for the Delta (left) and Omicron waves (right), over each of the 32 individual local authorities.

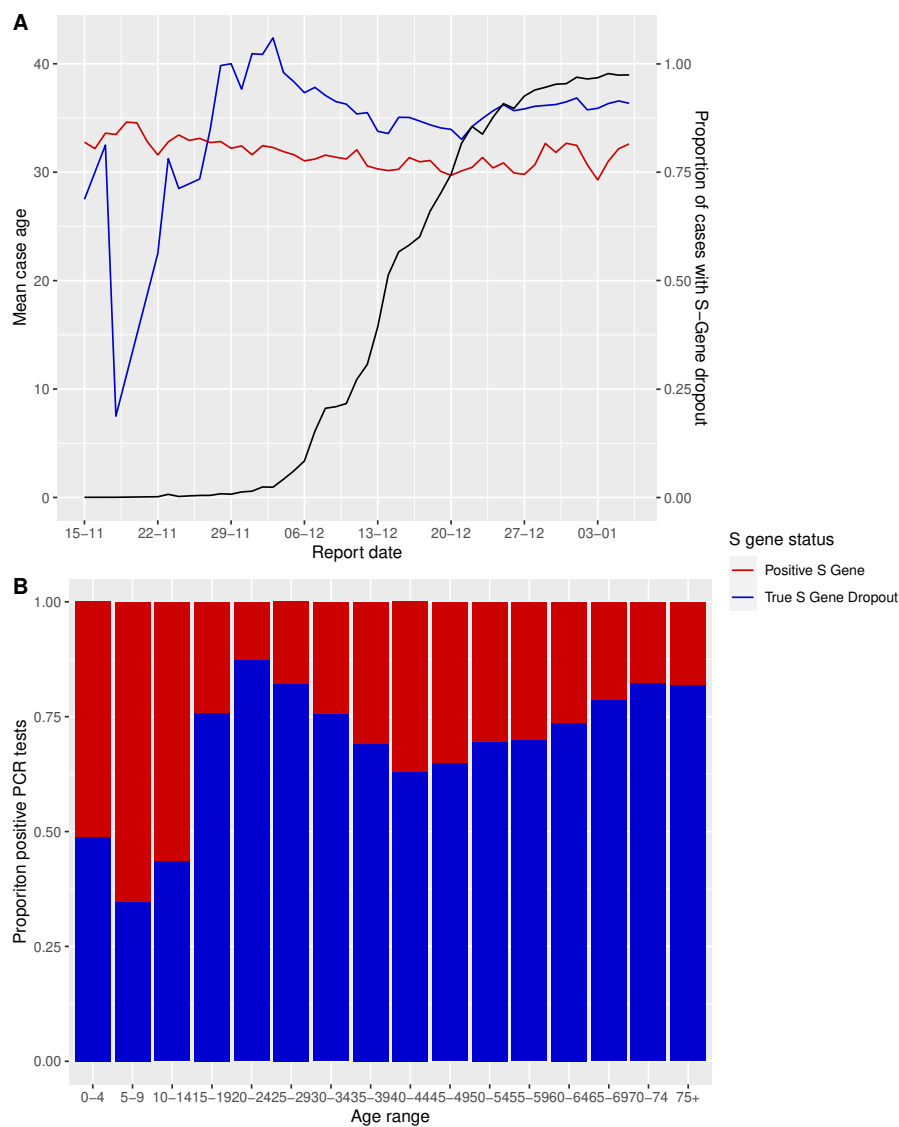


Figure S4: PCR positive cases over the period 15<sup>th</sup> November 2021 to 6<sup>th</sup> January 2022 that were S-gene dropout or true S-gene positive. (A) Daily mean case age for the two definite PCR S-gene outcomes (blue and red lines) against the proportion of the daily cases that were true S-gene dropout (presumed Omicron-type). (B) the proportion of the cases over the period by 5-year age bracket.

## 622 B Additional methodology details

### 623 B.1 Model hyperparameters

624 The random forest regression model is fit in *R* version 4.1.0 [68], using the *randomForest* pack-  
625 age [69] (version 4.6-14), and ALEs analysed using the *ALEPlot* package [71] (version 1.1).

626 From 6,976 DZs, 2 sexes, 16 age ranges, and 3 prior case states, there were a total of 669,696  
627 cohorts (of which a fraction will have population zero and are excluded). Cohorts from 90% of  
628 DZs were used for the fit, with 10% reserved to test model performance against data it explicitly  
629 did not fit. The fit was made to  $\sqrt{\text{cases} + 1}$ . The RF comprised 500 trees, with cohorts sampled  
630 for building each tree weighted by population. 5 variables were tested at each split, and each tree  
631 had a maximum of 30,000 terminal nodes, with a minimum node size of 300.

### 632 B.2 Explanatory variables used in random forest regression model

633 The models described in Section 4.3 are informed with the following data, first at cohort resolution:

- 634 • *Age range* (five-year windows: [0 – 4], [5 – 9], ..., [70 – 74], [75+]), using the numeric  
635 intermediate values 2, 7, ..., 72, and 75 for the 75+ category;
- 636 • *Sex*;
- 637 • *Prior case status*: the time of the last reported case, broken into three categories: never  
638 tested positive before, last tested positive in the 6 months prior to the first day of the  
639 outbreak, last tested positive over 6 months prior;
- 640 • *Cohort population* (derived using historical testing data for those testing positive before, and  
641 estimated populations as of mid-2020 collated by the National Records of Scotland [74], for  
642 the remainder that had not tested positive before).

643 At age/sex/DZ resolution, we then include:

- 644 • COVID-19 *vaccination uptake* (eDRIS) (see also Supplementary material B.3);
- 645 • *Ethnicity* (% population belonging to a minority ethnicity), as per the most recent Scottish  
646 census data (2011);
- 647 • The per-population, time-aggregated number of *negative LFD tests* reported in that period.

648 Finally included are the following at DZ resolution or broader:



- 649 ● Measures of DZ-level deprivation (obtained from Scottish census data, and the 2020 *Scottish*  
650 *Index of Multiple Deprivation* [70]);
- 651 ● *Local outbreak duration*: the difference between the final date of the period studied, and the  
652 date the variant was first detected in that cohort's corresponding *intermediate zone (IZ)*.  
653 An IZ typically contains 4–6 DZs, and 3,000–5,000 individuals, with this granularity chosen  
654 to give a reasonable proxy for when the variant was seeded locally;
- 655 ● *Student population* (% population being a full-time student aged 18 or over), also per 2011  
656 census data;
- 657 ● *Population density*, at IZ-level;
- 658 ● *S-gene coverage* (the proportion of cases with an accompanying S-gene result, required to  
659 associate a likely variant) at IZ level. S-gene coverage was 90% overall across mainland  
660 Scotland (per eDRIS data), but significantly lower in the LAs of Orkney Islands, Shetland  
661 Islands and Na h-Eileanan Siar (74%, 20% and 23% respectively).

662 The measures of DZ-level deprivation included are [66]:

- 663 ● *Drive time from GP*: Average drive time to a GP surgery in minutes;
- 664 ● *Public transport time to GP*: Public transport travel time to a GP surgery in minutes;
- 665 ● *% Income deprived*: Proportion of individuals in receipt of income support payments, such  
666 as Job Seekers Allowance;
- 667 ● *% Employment deprived*: Proportion of working age population claiming employment-related  
668 payments, such as Incapacity Benefit;
- 669 ● *Standardised mortality ratio*: Age/sex-standardised mortality rate as compared to the overall  
670 population;
- 671 ● *Comparative illness factor*: Proportion of individuals claiming from a variety of illness and  
672 disability-related payments as compared to the overall population;
- 673 ● *Drug-related hospitalisation ratio*: Rate of hospitalisations relating to drug use, as compared  
674 to the overall population;
- 675 ● *Alcohol-related hospitalisation ratio*: Rate of hospitalisations relating to alcohol use, as com-  
676 pared to the overall population;

- 677 • *Crime rate*: Rate of recorded crimes per population;
- 678 • *Attendance*: Percentage of pupils with school attendance of over 90%;
- 679 • *Attainment*: Measure for average attainment of school leavers from 2015–2018;
- 680 • *Ratio working age with no qualifications*: Proportion of working age people with no qualifi-  
681 cations, as compared to the overall population.

682 We do not use data on *PCR* negative tests. In the Omicron wave PCR positivity peaked  
683 at 30% (per eDRIS data), with testing capacity being reached (resulting in a policy change on  
684 5<sup>th</sup> January 2022 removing the need for a confirmatory PCR after an LFD positive [64]). Thus  
685 with this “ceiling” capacity being reached, we exclude negative PCR tests as a poorer proxy for  
686 propensity to test as compared to LFD negatives, and being too closely related to overall cases  
687 (requiring an S-gene sequenced positive PCR test).

### 688 **B.3 Vaccination uptake as an explanatory variable**

689 Scotland’s COVID-19 vaccination programme began on December 8<sup>th</sup> 2020, with initial priority  
690 given to healthcare workers, the elderly and those otherwise especially vulnerable to COVID-19,  
691 then generally by decreasing age [13]. All first doses had been offered and administered to willing  
692 adults by 18<sup>th</sup> July 2021 [75], with rates of first dose administration declining thereafter. By 15<sup>th</sup>  
693 November 2021, then, the first dose date may have differed between two individuals by up to 11  
694 months. This likely led to substantial variation in protection offered by the first dose at the time  
695 of the Omicron wave, given both evidence of efficacy waning over timescales of six months, and  
696 high rates of breakthrough for Omicron against vaccines originally designed against earlier “wild-  
697 type” SARS-CoV-2 lineages, particularly for non-mRNA vaccines [42, 43, 44]. This, combined  
698 with high uncertainty in the cohort-level population denominator used to determine uptake, leads  
699 us to exclude first and second dose uptake (being highly correlated with first dose uptake) as an  
700 explanatory variable for Omicron cases. We do, however, include third/booster dose uptake, as  
701 the proportion receiving a first dose to have *returned* for a third/booster dose by 15<sup>th</sup> November  
702 2021 (and zero if nobody in the cohort had yet received a first dose). This definition eliminates  
703 uncertainty in the underlying population. Prior to the detection of Omicron, those aged 50+ or  
704 otherwise vulnerable to COVID-19 were due to be offered a third or booster dose, twelve weeks  
705 after their second [76]. The booster programme began on September 20<sup>th</sup> 2021, and a snapshot  
706 on 15<sup>th</sup> November 2021 shows substantial variation between different cohorts, particularly by age.

707 With these doses being delivered more recently, as well as evidence of this dose proving more  
708 protective against Omicron [42, 4], we include this definition of third/booster dose uptake as a  
709 reasonable proxy for vaccine-induced protection against Omicron at the time.

710 The initial Delta wave occurred while the bulk of first and second doses were still being admin-  
711 istered, thus we include second dose uptake on 1<sup>st</sup> May 2021 as an explanatory variable, as the  
712 proportion of individuals that had returned for a second dose, having received a first (and zero, if  
713 nobody in the cohort had yet received their first dose).

714 **C** Map views of population distribution, model residuals

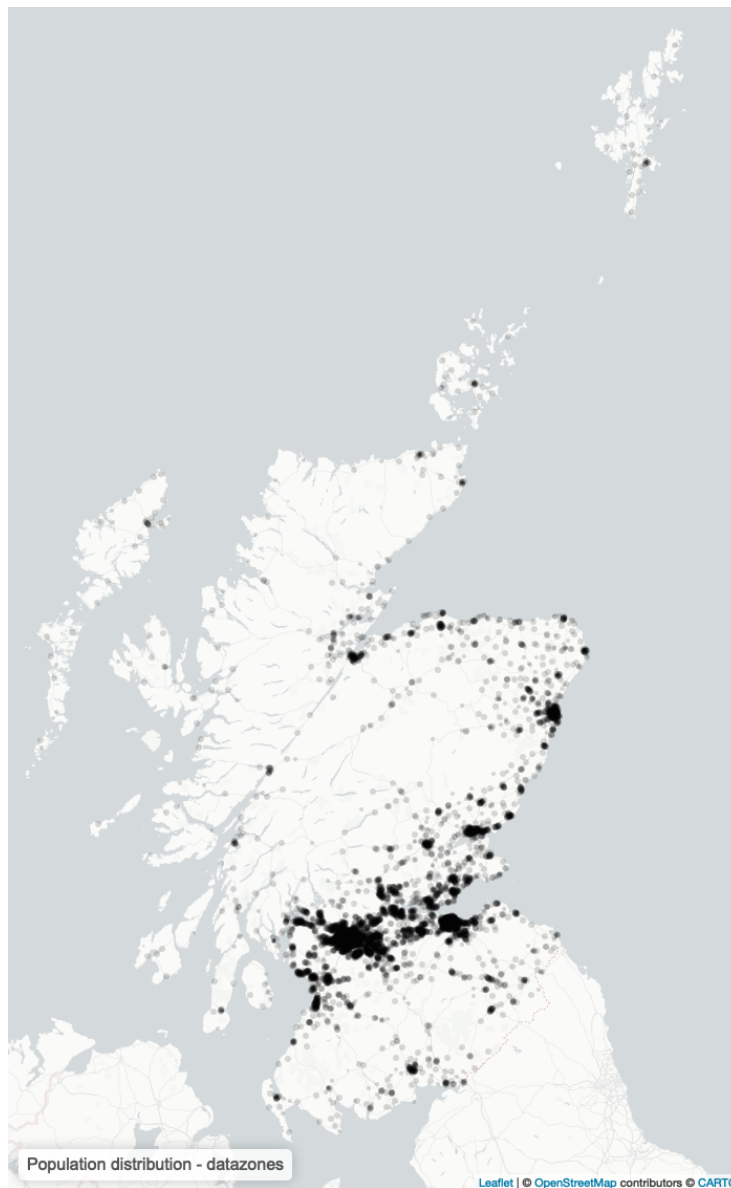


Figure S5: Distribution of population in Scotland. Each point indicates the population-weighted centroid of a datazone (DZ) of which there are 6,976 in total, with each representing a population of approximately 500-1,000 individuals.

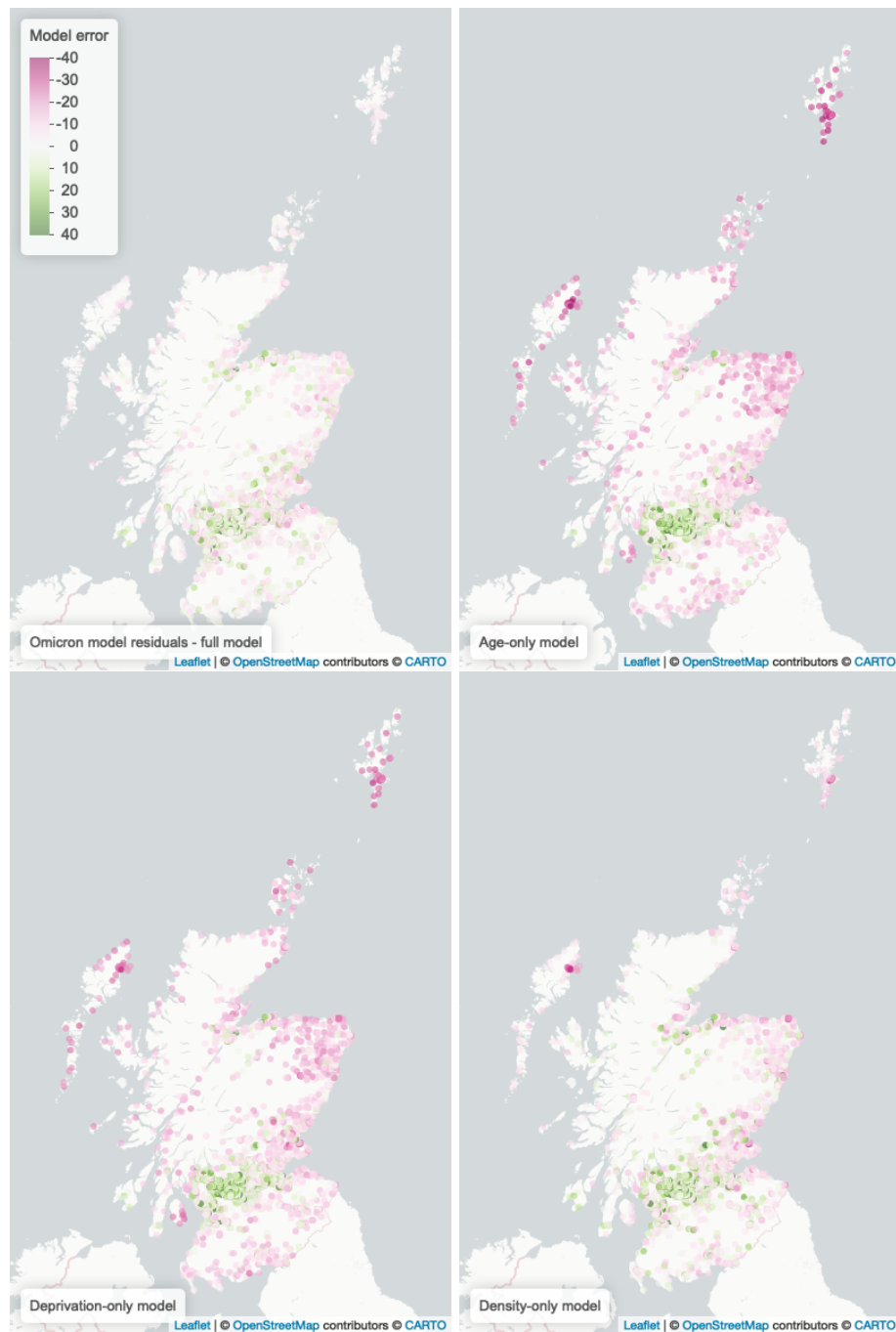


Figure S6: Comparing residuals for the distribution of Omicron cases, between the full model (top left), and the reduced models informed by population and one of age, deprivation and population density respectively. The colour scale indicates the DZ-level model error (model estimate - data), where purple points indicate DZs where the model overestimated the number of cases, and green points indicate where the model underestimated cases.

## 715 D Random forest variable importance

716 Fig. S7 shows variable importance measures extracted from the *RandomForest* function. Age,  
 717 population and prior case status have much higher node purity (Fig. S7, top) than the other  
 718 variables, indicating that splits in individual trees using values of these variables in particular are  
 719 characteristically more “effective” at separating cohorts with differing numbers of cases. Fig. S7,  
 720 bottom, then shows random permutation of each of the variables results in appreciable increase in  
 721 fit error, confirming that this larger collection of variables are important to explain finer patterns  
 722 in the data.

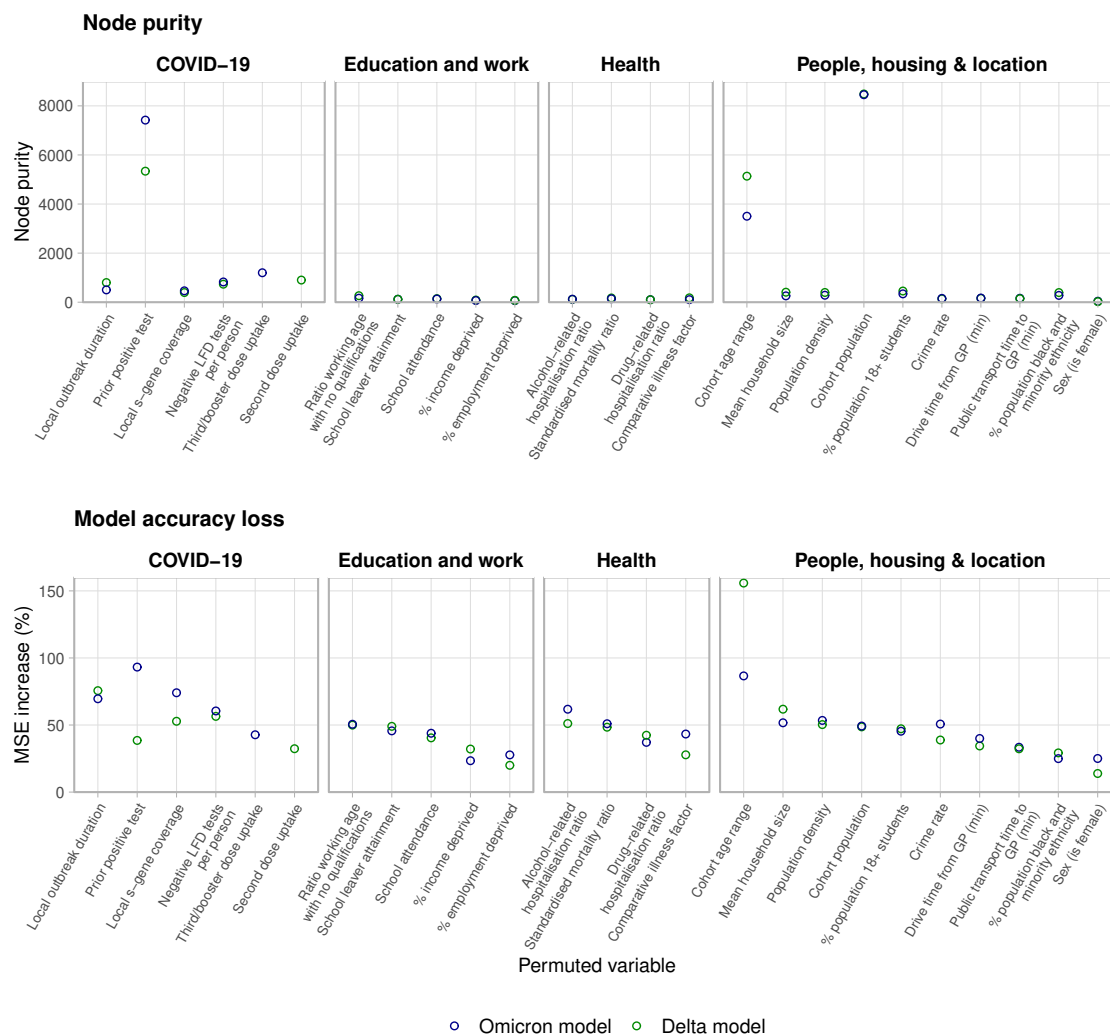


Figure S7: Feature importance outputs from the random forest regression models. Top: Node purity. Bottom: Explanatory variable mean squared error (MSE) increase on random permutation: for variable  $i$ , the increase in MSE on data not trained in each tree, if the entries of  $i$  were instead randomly permuted.

723 **E Frequency of lateral flow testing by sex, deprivation quin-**  
 724 **tile**

		Population	LFD tests	Tests per person	Positive LFD tests	Positivity
Total		5,466,000	29,508,794	5.40	1,123,210	3.81%
Sex	F	2,800,788	19,639,047	7.01	647,529	3.30%
	M	2,665,212	9,869,747	3.70	475,681	4.82%
Deprivation quintile (1: most deprived)	1	1,057,767	3,827,970	3.62	176,600	4.61%
	2	1,057,929	4,992,257	4.72	201,148	4.03%
	3	1,077,589	6,023,840	5.59	220,258	3.66%
	4	1,140,448	7,134,794	6.26	256,101	3.59%
	5	1,132,267	7,529,933	6.65	269,103	3.57%

Table S7: Summary statistics of lateral flow device (LFD) tests reported in Scotland from July 2020 to February 2023, broken down by sex, and deprivation quintile of the residing datazone of individuals as ranked by the 2020 Scottish Index of Multiple Deprivation, where the most deprived datazones are in quintile 1. The test positivity is the proportion of all tests of any result that were reported as positive.

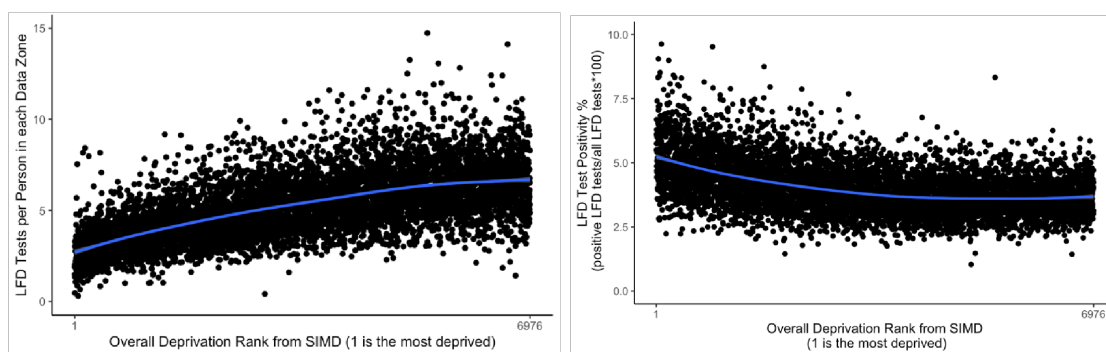


Figure S8: Lateral flow testing from July 2020 to February 2023 by datazone, ranked by deprivation per the 2020 Scottish Index of Multiple Deprivation, where the rank 1 is the datazone ranked as most deprived. Left: the number of LFD tests reported per person in each datazone. Right: the LFD test positivity, defined as the proportion of all reported LFD tests to have been positive.