

1 Model design for non-parametric phylodynamic inference
2 and applications to pathogen surveillance

3 Xavier Didelot^{1,*}, Lily Geidelberg², The COVID-19 Genomics UK (COG-UK) consortium³,
4 and Erik M Volz²

5 ¹ School of Life Sciences and Department of Statistics, University of Warwick, United Kingdom

6

7 ² Department of Infectious Disease Epidemiology, School of Public Health, Imperial College London,
8 United Kingdom

9

10 ³ Full list of consortium names and affiliations are in the appendix. <https://www.cogconsortium.uk>

11 * Corresponding author. Tel: 0044 (0)2476 572827. Email: xavier.didelot@gmail.com

12 ABSTRACT

13 Inference of effective population size from genomic data can provide unique information about
14 demographic history, and when applied to pathogen genetic data can also provide insights into
15 epidemiological dynamics. The combination of non-parametric models for population dynamics with
16 molecular clock models which relate genetic data to time has enabled phylodynamic inference based
17 on large sets of time-stamped genetic sequence data. The methodology for non-parametric inference
18 of effective population size is well-developed in the Bayesian setting, but here we develop a frequentist
19 approach based on non-parametric latent process models of population size dynamics. We appeal to
20 statistical principles based on out-of-sample prediction accuracy in order to optimize parameters that
21 control shape and smoothness of the population size over time. We demonstrate the flexibility and
22 speed of this approach in a series of simulation experiments, and apply the methodology to reconstruct
23 the previously described waves in the seventh pandemic of cholera. We also estimate the impact of non-
24 pharmaceutical interventions for COVID-19 in England using thousands of SARS-CoV-2 sequences.
25 By incorporating a measure of the strength of these interventions over time within the phylodynamic
26 model, we estimate the impact of the first national lockdown in the UK on the epidemic reproduction
27 number.

28 INTRODUCTION

29 Past fluctuation in the size of a population are reflected in the genealogy of a sample of individuals
30 from that population. For example, under the coalescent model, two distinct lines of ancestry coalesce
31 (i.e. find a common ancestor) at a rate that is inversely proportional to the effective population size at
32 any given time (Kingman 1982; Griffiths and Tavaré 1994; Donnelly and Tavaré 1995). More coalescent
33 events are therefore likely when the population size is small compared to when the population size is
34 large. This causal effect of population size on genealogies can be reversed in an inferential framework
35 to recover past population size dynamics from a given pathogen genealogy. This approach to inference
36 of past demographic changes was first proposed 20 years ago (Pybus et al. 2000, 2001; Strimmer and
37 Pybus 2001) and has been fruitfully applied to many disease systems (Pybus and Rambaut 2009; Ho
38 and Shapiro 2011; Baele et al. 2016).

39 Population size analysis is often performed within the Bayesian BEAST framework (Suchard et al. 2018;
40 Bouckaert et al. 2019) which jointly infers a phylogeny and demographic history from genetic data. Here
41 we focus on an alternative approach in which the dated phylogeny is inferred first, for example using
42 *treedater* (Volz and Frost 2017), *TreeTime* (Sagulenko et al. 2018) or *BactDating* (Didelot et al. 2018),
43 and demography is investigated on the basis of the phylogeny. Although potentially less sensitive,
44 this approach has the advantage of scalability to very large sequence datasets. This post-processing
45 approach also allows more focus on models and assumptions involved in the demographic inference
46 itself as previously noted in studies following the same strategy (Lan et al. 2015; Karcher et al. 2017;
47 Volz and Didelot 2018; Volz et al. 2020). However, some of the methodology and results we describe
48 here should be applicable in a joint inferential setting as well.

49 The reconstruction of past population size dynamics is usually based on a non-parametric model, since
50 the choice of any parametric function for the past population size would cause restrictions and be
51 hard to justify in many real-life applications (Drummond et al. 2005; Ho and Shapiro 2011). However,
52 even if a non-parametric approach offers a lot more flexibility than a parametric one, it does not fully
53 circumvent the question of how to design the demographic model to use as the basis of inference. For
54 example, the *skygrid* model considers that the logarithm of the effective population size is piecewise
55 constant, with values following a Gaussian Markov chain, in which each value is normally distributed
56 around neighbouring values and standard deviation determined by a precision hyperparameter (Gill

57 et al. 2013). This model can be justified as the discretisation of a continuous *skyride* model in which the
58 logarithm of the population size is ruled by a Brownian motion (Minin et al. 2008). Alternatively, the
59 *skygrowth* model is a similar Gaussian Markov chain on the growth rate of the population size (Volz
60 and Didelot 2018). Both models can be conveniently extended to explore the association between
61 population size dynamics and covariate data (Gill et al. 2016; Volz and Didelot 2018).

62 The *skygrid*, *skygrowth* or other similar models can be assumed when performing the inference of
63 the demographic function, and the effect of this model choice has not been formally investigated.
64 Furthermore, these non-parametric models require several model design choices which are often
65 given little consideration in practice. This includes the number of pieces in the piecewise constant
66 demographic function, the location of boundaries between pieces, and the prior expectation for the
67 difference from one piece to another. All of these model design choices may have significant effect on
68 the inference results. Here we propose several statistical procedures to optimise these variables. In
69 particular, the parameter controlling the smoothness of the population size function is usually assumed
70 to have an arbitrary non-informative prior distribution in a Bayesian inferential setting (Minin et al.
71 2008; Gill et al. 2013), whereas we show here that it can be selected using a frequentist statistical
72 approach based on out-of-sample prediction accuracy. We tested the effect of these procedures on
73 simulated datasets, where the correct demographic function is known and can be used to assess the
74 relative accuracy of inference under various conditions. We applied our methodology to a previously
75 published dataset of *Vibrio cholerae*, the causative agent of cholera. We also analysed a state-of-the-art
76 real dataset and show how our methodology can be used to estimate the impact of non-pharmaceutical
77 interventions for SARS-CoV-2 in England.

78 MATERIALS AND METHODS

79 Demographic Models

80 Let the demographic function $N_e(t)$ denote the effective population size of a pathogen at time t . Let
81 us consider that $N_e(t)$ is piecewise linear with R pieces of equal lengths h over the timescale of interest.
82 Let γ_i denote the logarithm of the effective population size in the i -th piece. In the *skygrid* model
83 (Gill et al. 2013), the values of γ_i follow a Gaussian Markov chain, with the conditional distribution
84 of γ_{i+1} given γ_i equal to:

$$\gamma_{i+1} \sim \mathcal{N}(\gamma_i, h/\tau) \quad (1)$$

85 By contrast, the *skygrowth* model (Volz and Didelot 2018) is defined using the effective population size
86 growth rates ρ_i which are assumed constant in each interval and are equal to:

$$\rho_i = \frac{\exp(\gamma_{i+1}) - \exp(\gamma_i)}{h \exp(\gamma_i)} \quad (2)$$

87 These growth rate values form a Gaussian Markov chain, with:

$$\rho_{i+1} \sim \mathcal{N}(\rho_i, h/\tau) \quad (3)$$

88 We also define a new model which we call *skysigma* based on the values σ_i of the second order differences
89 of the logarithm of the effective population size:

$$\sigma_i = (\gamma_{i+1} - \gamma_i) - (\gamma_i - \gamma_{i-1}) = \gamma_{i+1} - 2\gamma_i + \gamma_{i-1} \quad (4)$$

90 Once again we consider a Gaussian Markov chain in which:

$$\sigma_{i+1} \sim \mathcal{N}(\sigma_i, h/\tau) \quad (5)$$

91 Dependency on known covariate time series can be easily incorporated into these models as previously

described (Gill et al. 2016; Volz and Didelot 2018). Let there be a $m \times p$ matrix $X_{1:m,1:p}$ of p covariate measurements for each of m time points. Ideally these time points would correspond to the $R + 1$ boundaries between pieces of the demographic function, but otherwise linear interpolation can be used to make it so. We model the effect of this covariate data as a modification of the expected change in the demographic variables defined above (γ_i, ρ_i or σ_i). For example, in the *skysigma* model (Equation 5), the kernel of the Markov chain becomes:

$$\sigma_{i+1} \sim \mathcal{N}(\sigma_i + (X_{i+1,1:p} - X_{i,1:p})\beta, h/\tau) \quad (6)$$

where $\beta_{1:p}$ is a vector of coefficients for a linear model of the covariate data on the expected value of the increments. Note in particular that if a term in the β vector is equal to zero, then this covariate measurement has no effect on the demographic function, so that to test the significance of covariate requires to test whether the corresponding value in the β vector is non-zero.

Coalescent framework

Each of the models above defines a demographic function $N_e(t)$ from which the likelihood of the genealogy \mathcal{G} can be calculated as briefly described below. Let n denote the number of tips in \mathcal{G} , let $s_{1:n}$ denote the dates of the leaves and $c_{1:(n-1)}$ denote the dates of the internal nodes. Let $A(t)$ denote the number of extant lineages at time t in \mathcal{G} which is easily computed as the number of leaves dated after t minus the number of internal nodes dated after t :

$$A(t) = \sum_{i=1}^n \mathbb{1}[s_i > t] - \sum_{i=1}^{n-1} \mathbb{1}[c_i > t] \quad (7)$$

This quantity is important because in the coalescent model, each pair of lineages finds a common ancestor at rate $1/N_e(t)$, so that the total coalescent rate at time t is equal to:

$$\lambda(t) = \begin{cases} \frac{A(t)(A(t)-1)}{2N_e(t)}, & \text{if } A(t) \geq 2 \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

110 The full likelihood of the coalescent process is therefore computed as (Griffiths and Tavaré 1994;
 111 Donnelly and Tavaré 1995):

$$L(\mathcal{G}|N_e(t)) = \exp\left(-\int_{-\infty}^{\infty} \mathbb{1}[A(t) \geq 2] \frac{A(t)(A(t)-1)}{2N_e(t)} dt\right) \prod_{i=1}^{n-1} \frac{1}{N_e(c_i)} \quad (9)$$

112 This computation is straightforward for the models considered here where the demographic function
 113 $N_e(t)$ is piecewise constant.

114 Selection of the precision parameter

The demographic models described above (*skygrid*, *skygrowth* and *skysigma*) all rely on a precision parameter τ (also known as the 'smoothing' parameter). The value of τ controls how much consecutive values of the effective population size will vary when the data is uninformative. The selection of this parameter is therefore shaped by competing aims of optimising the fit to observed data and maximizing explanatory power and avoidance of overfitting. In frequentist statistics, a standard approach to selecting smoothing parameters is to minimize the out-of-sample prediction error. Here, we pursue a k -fold cross-validation strategy where genealogical data is partitioned into k sets, $k-1$ of which are used for fitting, and the last one is used for prediction. This procedure is equivalent to maximizing the following objective function:

$$f(\tau) = \prod_{j=1}^k L(\mathcal{G} \setminus X_j | \hat{N}_e(X_j, \tau)), \quad (10)$$

115 where $\hat{N}_e(X_j, \tau)$ is the maximum likelihood estimates of N_e on the partial data $X_j \subset \mathcal{G}$ and assuming
 116 the precision parameter is τ . In this case $X_{j=1:k}$ represents a subset of the sample times and internal
 117 node times of the genealogy \mathcal{G} .

118 This is a standard formulation of the cross-validation method, but the implementation depends on how
 119 genealogical data is partitioned. We use the strategy of discretizing the coalescent likelihood (Equation
 120 9) into intervals bordered by the time of nodes (tips s_i or internal nodes c_i of the tree) and/or the $R-1$
 121 times when the piecewise-constant N_e changes value. Given $R-1$ change points, n tips, and $n-1$
 122 internal nodes of \mathcal{G} , there are $R+2n-3$ intervals $(\iota_1, \dots, \iota_{R+2n-3})$. Each cross-validation training set

123 is formed by taking a staggered sequence of intervals and collecting the genealogical data contained in
124 each, so that $X_k = \{\iota_{j=1:R+2n-3} \mid \text{modulo}(j, k) \neq 0\}$.

125 Selection of the grid resolution

126 Before any of the non-parametric models described above can be fitted, the number R of pieces in the
127 piecewise demographic function needs to be specified. Setting R too low may lead to an oversimplified
128 output that does not capture all the information on past population changes suggested by the genealogy,
129 whereas setting R too high can lead to overfitting.

130 We therefore propose to use well established statistical methods to select the optimal value of R . First
131 the model is fitted for multiple proposed values of R , and then for each output we compute the Akaike
132 information criterion (AIC), which is equal to:

$$\text{AIC}_R = 2R - 2\log(L_R) \quad (11)$$

133 where L_R is the maximum value of the likelihood when using R pieces. The value of R giving the
134 smallest value of AIC_R is selected. We also implemented the Bayesian information criterion (BIC),
135 which is equal to:

$$\text{BIC}_R = R\log(n-1) - 2\log(L_R) \quad (12)$$

136 Simulation of testing data

137 In order to test the accuracy of our methodology, we implemented a new simulator of coalescent
138 genealogies given sampling dates and a past demographic function $N_e(t)$. When the demographic
139 function is constant, the simulation of coalescent genealogies is equivalent to simulating from a
140 homogeneous Poisson process, in which the waiting times from one event to the next are exponentially
141 distributed. To extend this to the situation where the demographic function is non-constant requires to
142 simulate from an equivalent non-homogeneous Poisson process. The approach we used to achieve this
143 is to consider a homogeneous Poisson process with a population size N_m which is lower than any value
144 of $N_e(t)$, i.e. $\forall t, N_e(t) \geq N_m$. We simulate this process using exponential waiting times, but filter an

145 event happening at time t according to the ratio $N_m/N_e(t)$. Specifically, we draw $u \sim \text{Unif}(0, 1)$ and
146 if $u < N_m/N_e(t)$ the event is accepted and otherwise rejected. The resulting filtered Poisson process
147 simulates from the non-homogeneous Poisson process as required (Ross 2014). The disadvantage of
148 this approach over other methods of simulations is that there may be many rejections if $N_e(t)$ takes
149 small values so that N_m needs to be small too. However, efficiency of simulation is not important for
150 our purpose here, and this method has the advantage to avoid the computation of integrals on the
151 $N_e(t)$ function which other methods would require.

152 **Implementation**

153 We implemented the simulation and inference methods described in this paper into a new R
154 package entitled *mlesky* which is available at <https://github.com/emvolz-phylogenetics/mlesky>. The
155 optimisation of the demographic function makes use of the quasi-Newton Broyden-Fletcher-Goldfarb-
156 Shanno (BFGS) method implemented in the `optim` command (Nash 2014). Confident intervals are
157 computed based on an approximation of the curvature of the likelihood surface around its maximum.
158 If multiple CPU cores are available, these resources are exploited within the procedure of selection
159 of the smoothing parameter where the computation can be split between the different cross values
160 in the cross-validation. Multicore processing is also applied in the procedure of selection of the grid
161 resolution where computation can be split between different values of the resolution parameter R .
162 All the code and data needed to reproduce our results on simulated and real datasets is available at
163 <https://github.com/mrc-ide/mlesky-experiments>.

164 **RESULTS**

165 **Application to simulated phylogeny with constant population size**

166 A dated phylogeny was simulated with 200 tips sampled at regular intervals between 2000 and 2020,
167 and a constant past population size function $N_e(t) = 20$ (Figure S1). To illustrate the importance of
168 the resolution R and precision τ parameters, we inferred the demographic function under the *skygrid*

169 model (cf Equation 1) for a grid of values with $R \in \{5, 20, 50\}$ and $\tau \in \{1, 10, 20\}$ (Figure 1). The
170 results look quite different depending on the parameters used, and in particular when R is large and τ
171 is small, fluctuation in the population size are incorrectly inferred. When applying the AIC procedure
172 to this dataset, the correct value of $R = 1$ was inferred for which the parameter τ becomes irrelevant.
173 In these conditions the effective population size was estimated to be 19.65 with confidence interval
174 ranging from 17.10 to 22.57 which includes the correct value of 20 used in the simulation. We repeated
175 the AIC procedure for 100 different phylogenies all which had been simulated under the same constant
176 population size conditions described above. For 65 of these phylogenies the AIC procedure selected
177 $R = 1$, with the third quartile falling on $R = 3$ and 94% of the simulations giving $R \leq 5$. We also
178 applied the BIC procedure for the same 100 phylogenies, and found that $R = 1$ was selected in all but
179 one instance for which $R = 2$ was inferred. However, the BIC is well known to be overly conservative
180 (Kuha 2004; Weakliem 1999) and so the rest of results make use of the AIC procedure.

181 **Application to simulated phylogeny with varying population size**

182 Next we simulated a dated phylogeny with the same number and dates of the tips as previously,
183 but using a demographic function $N_e(t)$ that was sinusoidal with minimum 2 and maximum 22, with
184 period 6.28 years. Figure S2 shows both the demographic function used and the resulting simulated
185 phylogeny. We attempted to reconstruct the demographic function based on the phylogeny under the
186 three models *skygrid*, *skygrowth* and *skysigma* described in Equations 1, 3 and 5, respectively. For
187 each model the precision parameter τ was optimised using our new cross-validation procedure and the
188 number of pieces was set to be $R = 20$ for ease of comparison. The results obtained in these conditions
189 were very similar under the three models (Figure 2). This suggests that when the precision parameter is
190 optimised using the cross-validation method, the choice between these three models becomes relatively
191 unimportant. The same conclusions when reached when comparing the results of inference based on
192 the three models to other simulated phylogenies. The choice of using one model rather than another is
193 therefore mostly guided by the presence of covariate data and whether these are expected to correlate
194 with the effective population size directly or some other function of it such as the population growth
195 rates (Gill et al. 2016; Volz and Didelot 2018).

196 One situation in which all models are expected to perform poorly is when then there are sudden changes

197 to the demographic function. To exemplify this, we simulated another dated phylogeny with the same
198 and dates of the tips as before, but using a bottleneck function for $N_e(t)$ which was equal to 10 at all
199 times except between 2005 and 2010 when it was equal to 1 (Figure 3A). The phylogeny simulated
200 using this bottleneck function is shown in Figure 3B. We reconstructed the demographic function using
201 the *skygrid* model. The lowest value of the AIC was obtained for $R = 14$, and the precision parameter
202 was optimised using the cross-validation procedure to $\tau = 0.87$. The inferred demographic function is
203 shown in Figure 3C, where the bottleneck between 2005 and 2010 has been accurately detected.

204 Application to simulated phylogeny with covariate data

205 Finally, we used simulations to test our procedure for the analysis of association between demography
206 and covariate data. An example is shown in Figure S3 where the covariate data follows a simple
207 quadratic function in order to create a boom and bust dynamic (Figure S3A). The growth rate of
208 the population however does not follow exactly this function, and is subjected to monthly Gaussian
209 noise with standard deviation 0.4 in this case (Figure S3B). From this growth rate we compute the
210 effective population size function over time (Figure S3C) and simulate a phylogenetic tree as previously,
211 with 200 tips sampled at regular intervals between 2000 and 2020 (Figure S3D). We then analysed
212 this simulated phylogeny alongside the covariate data, and found in this case a strong association
213 with coefficient $\beta = 0.77$. We repeated this procedure 100 times with increasing values of the noise
214 standard deviation and the results are summarised in Figure S4. As expected, we found that as the
215 noise increases, the coefficient of association β between growth rate and the covariate decreases, and
216 eventually the association becomes non-significant with an estimated coefficient of association close to
217 zero.

218 Application to *Vibrio cholerae* dataset

219 We applied our methodology to a previously described collection of 260 genomes from the seventh
220 pandemic of *Vibrio cholerae* (Didelot et al. 2015). A genealogy was estimated in this previous study
221 using an early version of BactDating (Didelot et al. 2018), and it is reproduced in Figure 4A. We
222 applied the AIC procedure to determine that the demographic function would be modelled using

223 $R = 16$ pieces. The precision parameter was optimised to a value of $\tau = 1.84$ using the cross-validation
224 procedure. The whole analysis took less than 20 seconds on a standard laptop computer. The inferred
225 demographic function is shown in Figure 4B. A first peak was detected in the 1960s, followed by a
226 second peak in the 1970s and finally a third peak in the 1990s. This demographic function follows
227 closely on the previously described three “waves” of cholera spreading globally from the Bay of Bengal
228 (Mutreja et al. 2011; Didelot et al. 2015; Weill et al. 2017). However, these three waves had previously
229 been described based on phylogeographic reconstructions of the spread of the pandemic around the
230 world. The fact that we found a similar wave pattern in our analysis which did not include any
231 information about the geographical origin of the genomes provides further support for the validity of
232 this phylodynamic reconstruction.

233 **Estimating the impact of non-pharmaceutical interventions for COVID-19** 234 **in England**

235 We applied our methodology to the SARS-CoV-2 epidemic in England using data from the first
236 epidemic wave spanning the spring of 2020. By incorporating data on timing of public health measures
237 such as lockdowns, we estimated the association on non-pharmaceutical interventions (NPIs) with viral
238 transmission. The COVID-19 Genomics UK Consortium (COG-UK) was established on 23rd March
239 2020 and has coordinated a large-scale sequencing and bioinformatics effort to assist with COVID-19
240 surveillance and response (COG-UK Consortium 2020). The proportion of cases with a virus genome
241 has varied over time and increased rapidly in April 2020 following the establishment of large-scale
242 national sequencing laboratories. In order to facilitate molecular clock dating, we carried out a stratified
243 random sample of genomes between 1st January and 30th April 2020 ensuring good representation of
244 sequences across a wide range of calendar time. Sequences were ordered by sample date, binned by
245 day, and randomly selected from each bin. Duplicate sequences were removed. Repeating this process
246 ten times resulted in ten distinct sequence sub-samples with a mean of 4,217 sequences each.

247 As part of the COG-UK bioinformatics pipeline, a maximum likelihood tree is estimated at regular
248 intervals on the MRC-CLIMB infrastructure (Nicholls et al. 2021). We pruned these trees to retain
249 samples in each of our sequence sub-samples. Each of these sub-trees was then converted into time-
250 scaled phylogenies using *treedater* v0.5.1 (Volz and Frost 2017) by randomly resolving polytomies in the

251 tree and sampling a molecular clock rate of evolution from a normal distribution with mean 5.91×10^{-4}
252 substitutions per site per year and standard deviation 1.92×10^{-5} , based on previous analysis of SARS-
253 CoV-2 in the UK (Volz et al. 2021). In all, 100 time trees were estimated representing uncertainty in
254 phylogenetic dating and sampling variation. The *skysigma* model was fitted to the trees by maximizing
255 the combined (average) likelihood. Maximum likelihood estimates were also computed for each tree
256 and 95% quantiles were used to quantify the uncertainty in the parameter estimates. The results are
257 shown in Figure 5A for the estimation of the effective population size function and in Figure 5B for the
258 estimation of the basic reproduction number over time. The latter is calculated as $R(t) = \rho(t)\Psi + 1$
259 where $\rho(t)$ is the growth rate of the effective population size $N_e(t)$ estimated through time and Ψ is
260 the mean of the serial interval (Wallinga and Lipsitch 2007; Volz and Didelot 2018). The value $\Psi = 6.5$
261 days was used based on previous studies of infector-infectees pairs (Chan et al. 2020; Bi et al. 2020;
262 Wu et al. 2020).

263 The estimated peak of the epidemic occurred on 1st April 2020, eight days after the imposition of the
264 first national lockdown, illustrated by the red boxes in Figure 5. The rise and fall in $N_e(t)$ precedes a
265 similar dynamic in the number of confirmed cases by several weeks (Figure 5A), which is as expected
266 since the case ascertainment rate was initially very low and improved dramatically in April. On the
267 other hand, the estimated $N_e(t)$ is approximately consistent with the number of genomic sequences
268 available over time (Figure 5C). The estimated $R(t)$ decreased gradually in the three weeks preceding
269 the start of the national lockdown (Figure 5B). This may be due to changing behaviour prior to the
270 national lockdown and a changing proportion of cases due to travel-related importation. Travel-linked
271 cases declined while internal transmission increased throughout March and April (du Plessis et al.
272 2021).

273 To test for association between growth rates and NPIs, we also fitted the model to both the genealogical
274 data and the OxCGRT health containment index (Hale et al. 2020), a time series representing the
275 intensity of the public health response. A higher value of this index indicates more stringent NPIs.
276 The model was fitted under the assumption that the differential of the logarithm of $N_e(t)$ follows
277 differential of the OxCGRT index, which approximately corresponds to the hypothesis that the basic
278 reproduction number $R(t)$ follows the daily change in the index (Volz and Didelot 2018). The median
279 estimated epidemic trajectories are very similar when including this covariate (Figure 5A), and we
280 observe improved precision in the estimate of the reproduction number (Figure 5B).

281 Figure 6 shows the coefficient β which represents the estimated strength of association between the
282 reproduction number and the daily change in the OxCGRT index. A negative value of β indicates
283 a negative association between changes in NPIs and the reproduction number. We investigated how
284 the estimate of β depends on a lag (back-shift) between the value of the index and the demographic
285 function. The largest effects (most negative values) were found when shifting the index back 8 days,
286 which means relating the values of the index to the reproduction number 8 days later. In contrast,
287 when changes in NPIs are compared to growth rates that precede them by several days (negative
288 delay), the coefficient β is not significant.

289 DISCUSSION

290 Non-parametric phylodynamic inference of population size dynamics is usually carried out in a Bayesian
291 framework (Drummond et al. 2005; Minin et al. 2008; Gill et al. 2013). Here we presented methods
292 for performing such inference in a frequentist setting with a particular view towards model selection
293 and avoiding over-fitting. Optimal smoothing can be obtained in a natural way using standard cross-
294 validation methods, and the optimal resolution of the discretised demographic function is achieved
295 using the well-established AIC criterion. This approach can be advantageous when prior distributions
296 are difficult to design or results are sensitive to arbitrarily chosen priors. Methods based on likelihood
297 maximization are also fast and scalable to datasets much larger than is conventionally studied with
298 Bayesian methods, and the selection of smoothing parameters does not require arbitrarily chosen
299 hyperparameters. Conventional AIC metrics also alleviate the difficulty of model selection. In most of
300 our simulations, we find relatively little difference in our estimates when parameterizing the model in
301 terms of $\log(N_e(t))$ (Equation 1), the growth rate of $N_e(t)$ (Equation 3) or the second order variation
302 of $\log(N_e(t))$ (Equation 5), as long as the precision parameter τ for each model is optimized as we
303 proposed.

304 Our methodology assumed that a dated phylogeny has been previously reconstructed from the genetic
305 data. It is therefore well suited for the post-processing analysis of the outputs from *treedater* (Volz
306 and Frost 2017) or *TreeTime* (Sagulenko et al. 2018). A key assumption of our method, as with
307 its Bayesian counterparts, is that all samples in the phylogeny come from a single population ruled

308 by a unique demographic function. To ensure that this is indeed the case, complementary methods
309 are emerging that can test for the presence or asymmetry or hidden population structure in dated
310 phylogenies (Dearlove and Frost 2015; Volz et al. 2020). Conversely, if multiple phylogenies follow the
311 same demographic dynamic, they can be analysed jointly to provide a more precise reconstruction
312 of the demographic function and epidemiological parameters (Xu et al. 2019), and our software
313 implementation is able to perform such a joint analysis when appropriate.

314 Past variations in the effective population size of a pathogen population can reveal key insights into
315 past epidemiological dynamics and help make predictions about the future. It is important to note
316 that the effective population size is not generally equal to or even proportional to the number of
317 infections over time (Volz et al. 2009; Dearlove and Wilson 2013). On the other hand, the growth rate
318 of the effective population size can be used to estimate the basic reproduction number over time $R(t)$
319 (Wallinga and Lipsitch 2007; Volz et al. 2013; Volz and Didelot 2018) as we used in our application
320 to COVID-19 in England. Having good estimates of this quantity is especially important for assessing
321 the effect of infectious disease control measures (Fraser 2007), and phylodynamic approaches provide
322 a useful complementary approach to more traditional methods of estimation based on case report data
323 (Cori et al. 2013).

324 **Acknowledgements**

325 XD acknowledges funding from the National Institute for Health Research (NIHR) Health Protection
326 Research Unit in Genomics and Enabling Data. LG acknowledges funding from the MRC Doctoral
327 Training Partnership. EMV was supported by NIH R01-AI135970. EMV acknowledges the MRC
328 Centre for Global Infectious Disease Analysis (MR/R015600/1). All authors thank the COVID-19 UK
329 Consortium for providing data and bioinformatics resources. We thank all partners and contributors
330 to the COG-UK consortium who are listed at <https://www.cogconsortium.uk/about/>. COG-UK is
331 supported by funding from the Medical Research Council (MRC) part of UK Research & Innovation
332 (UKRI), the National Institute of Health Research (NIHR) and Genome Research Limited, operating
333 as the Wellcome Sanger Institute.

334 References

- 335 Baele G, Suchard MA, Rambaut A, Lemey P. 2016. Emerging concepts of data integration in pathogen
336 phylogenetics. *Syst. Biol.* 00:1–24.
- 337 Bi Q, Wu Y, Mei S, Ye C, Zou X, Zhang Z, Liu X, Wei L, Truelove SA, Zhang T, et al. (11 co-authors).
338 2020. Epidemiology and transmission of covid-19 in 391 cases and 1286 of their close contacts in
339 shenzhen, china: a retrospective cohort study. *The Lancet Infectious Diseases.* 20:911–919.
- 340 Bouckaert R, Vaughan TG, Fourment M, Gavryushkina A, Heled J, Denise K, Maio ND, Matschiner
341 M, Ogilvie H, Plessis L, et al. (11 co-authors). 2019. BEAST 2.5 : An Advanced Software Platform
342 for Bayesian Evolutionary Analysis. *PLoS Comput. Biol.* 15:e1006650.
- 343 Chan YWD, Flasche S, Lam TLT, Leung MHJ, Wong ML, Lam HY, Chuang SK. 2020. Transmission
344 dynamics, serial interval and epidemiology of covid-19 diseases in hong kong under different control
345 measures. *Wellcome Open Research.* 5:91.
- 346 COG-UK Consortium. 2020. An integrated national scale SARS-CoV-2 genomic surveillance network.
347 *The Lancet Microbe.* 1:e99.
- 348 Cori A, Ferguson NM, Fraser C, Cauchemez S. 2013. A new framework and software to estimate
349 time-varying reproduction numbers during epidemics. *Am. J. Epidemiol.* 178:1505–12.
- 350 Dearlove B, Wilson D. 2013. Coalescent inference for infectious disease: meta-analysis of hepatitis C.
351 *Philos. Trans. R. Soc. B.* 368:20120314.
- 352 Dearlove BL, Frost SDW. 2015. Measuring Asymmetry in Time-Stamped Phylogenies. *PLoS Comput.*
353 *Biol.* 11:e1004312.
- 354 Didelot X, Croucher NJ, Bentley SD, Harris SR, Wilson DJ. 2018. Bayesian inference of ancestral
355 dates on bacterial phylogenetic trees. *Nucleic Acids Res.* 46:e134.
- 356 Didelot X, Pang B, Zhou Z, McCann A, Ni P, Li D, Achtman M, Kan B. 2015. The Role of China in
357 the Global Spread of the Current Cholera Pandemic. *PLoS Genet.* 11:e1005072.
- 358 Donnelly P, Tavaré S. 1995. Coalescents and genealogical structure under neutrality. *Annu. Rev.*
359 *Genet.* 29:401–21.

- 360 Drummond AJ, Rambaut A, Shapiro B, Pybus OG. 2005. Bayesian coalescent inference of past
361 population dynamics from molecular sequences. *Mol. Biol. Evol.* 22:1185–92.
- 362 du Plessis L, McCrone JT, Zarebski AE, Hill V, Ruis C, Gutierrez B, Raghwani J, Ashworth J,
363 Colquhoun R, Connor TR, et al. (11 co-authors). 2021. Establishment and lineage dynamics of the
364 SARS-CoV-2 epidemic in the UK. *Science*. 371:708–712.
- 365 Fraser C. 2007. Estimating individual and household reproduction numbers in an emerging epidemic.
366 *PLoS One*. 2:e758.
- 367 Gill MS, Lemey P, Bennett SN, Biek R, Suchard MA. 2016. Understanding Past Population Dynamics
368 : Bayesian Coalescent-Based Modeling with Covariates. *Syst. Biol.* 65:1041–1056.
- 369 Gill MS, Lemey P, Faria NR, Rambaut A, Shapiro B, Suchard MA. 2013. Improving bayesian
370 population dynamics inference: A coalescent-based model for multiple loci. *Mol. Biol. Evol.* 30:713–
371 724.
- 372 Griffiths R, Tavaré S. 1994. Sampling theory for neutral alleles in a varying environment. *Philos.*
373 *Trans. R. Soc. B.* 344:403–410.
- 374 Hale T, Webster S, Petherick A, Phillips T, Kira B. 2020. Oxford covid-19 government response tracker
375 (OxCGRT). last updated. 8:30.
- 376 Ho SYW, Shapiro B. 2011. Skyline-plot methods for estimating demographic history from nucleotide
377 sequences. *Mol. Ecol. Resour.* 11:423–434.
- 378 Karcher MD, Palacios JA, Lan S, Minin VN. 2017. *phylodyn*: an R package for phylodynamic
379 simulation and inference. *Mol. Ecol. Resour.* 17:96–100.
- 380 Kingman J. 1982. The coalescent. *Stoch. Process. their Appl.* 13:235–248.
- 381 Kuha J. 2004. AIC and BIC: Comparisons of assumptions and performance. *Sociol. Methods Res.*
382 33:188–229.
- 383 Lan S, Palacios JA, Karcher M, Minin VN, Shahbaba B. 2015. An efficient Bayesian inference
384 framework for coalescent-based nonparametric phylodynamics. *Bioinformatics*. 31:3282–3289.
- 385 Minin VN, Bloomquist EW, Suchard MA. 2008. Smooth skyride through a rough skyline: Bayesian
386 coalescent-based inference of population dynamics. *Mol. Biol. Evol.* 25:1459–1471.

- 387 Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, Croucher NJ, Choi SY, Harris
388 SR, Lebens M, et al. (21 co-authors). 2011. Evidence for several waves of global transmission in the
389 seventh cholera pandemic. *Nature*. 477:462–465.
- 390 Nash JC. 2014. On best practice optimization methods in r. *Journal of Statistical Software*. 60:1–14.
- 391 Nicholls SM, Poplawski R, Bull MJ, Underwood A, Chapman M, Abu-Dahab K, Taylor B, Colquhoun
392 RM, Rowe WP, Jackson B, et al. (11 co-authors). 2021. CLIMB-COVID: continuous integration
393 supporting decentralised sequencing for SARS-CoV-2 genomic surveillance.
- 394 Pybus OG, Charleston MA, Gupta S, Rambaut A, Holmes EC, Harvey PH. 2001. The Epidemic
395 Behavior of the Hepatitis C Virus. *Science*. 292:2323–2325.
- 396 Pybus OG, Rambaut A. 2009. Evolutionary analysis of the dynamics of viral infectious disease. *Nat.*
397 *Rev. Genet.* 10:540–50.
- 398 Pybus OG, Rambaut A, Harvey PH. 2000. An integrated framework for the inference of viral population
399 history from reconstructed genealogies. *Genetics*. 155:1429–1437.
- 400 Ross SM. 2014. Introduction to probability models. Academic press.
- 401 Sagulenko P, Puller V, Neher RA. 2018. TreeTime: Maximum likelihood phylodynamic analysis. *Virus*
402 *Evol.* 4:vex042.
- 403 Strimmer K, Pybus OG. 2001. Exploring the Demographic History of DNA Sequences Using the
404 Generalized Skyline Plot. *Mol. Biol. Evol.* 18:2298–2305.
- 405 Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. 2018. Bayesian phylogenetic
406 and phylodynamic data integration using BEAST 1.10. *Virus Evol.* 4:vey016.
- 407 Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, Hinsley WR, Laydon DJ, Dabrera
408 G, O’Toole Á, et al. (11 co-authors). 2021. Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7
409 in England. *Nature*. 593:266–269.
- 410 Volz EM, Didelot X. 2018. Modeling the Growth and Decline of Pathogen Effective Population Size
411 Provides Insight into Epidemic Dynamics and Drivers of Antimicrobial Resistance. *Syst. Biol.*
412 67:719–728.
- 413 Volz EM, Frost SDW. 2017. Scalable relaxed clock phylogenetic dating. *Virus Evol.* 3:vex025.

- 414 Volz EM, Koelle K, Bedford T. 2013. Viral Phylodynamics. *PLoS Comput. Biol.* 9:e1002947.
- 415 Volz EM, Kosakovsky Pond SL, Ward MJ, Leigh Brown AJ, Frost SDW. 2009. Phylodynamics of
416 infectious disease epidemics. *Genetics.* 183:1421–30.
- 417 Volz EM, Wiuf C, Grad YH, Frost SDW, Dennis AM, Didelot X. 2020. Identification of hidden
418 population structure in time-scaled phylogenies. *Syst. Biol.* 69:884–896.
- 419 Wallinga J, Lipsitch M. 2007. How generation intervals shape the relationship between growth rates
420 and reproductive numbers. *Proc. Biol. Sci.* 274:599–604.
- 421 Weakliem DL. 1999. A critique of the Bayesian information criterion for model selection.
- 422 Weill Fx, Domman D, Njamkepo E, Tarr C, Rauzier J, Fawal N, Keddy KH, Salje H, Moore S,
423 Mukhopadhyay AK, et al. (15 co-authors). 2017. Genomic history of the seventh pandemic of
424 cholera in Africa. *Science.* 789:785–789.
- 425 Wu JT, Leung K, Bushman M, Kishore N, Niehus R, de Salazar PM, Cowling BJ, Lipsitch M, Leung
426 GM. 2020. Estimating clinical severity of covid-19 from the transmission dynamics in wuhan, china.
427 *Nature medicine.* 26:506–510.
- 428 Xu Y, Cancino-Munoz I, Torres-Puente M, Villamayor LM, Borrás R, Borrás-Máñez M, Bosque M,
429 Camarena JJ, Colomer-Roig E, Colomina J, et al. (32 co-authors). 2019. High-resolution mapping
430 of tuberculosis transmission: Whole genome sequencing and phylogenetic modelling of a cohort from
431 Valencia Region, Spain. *PLoS Med.* 16:1–20.

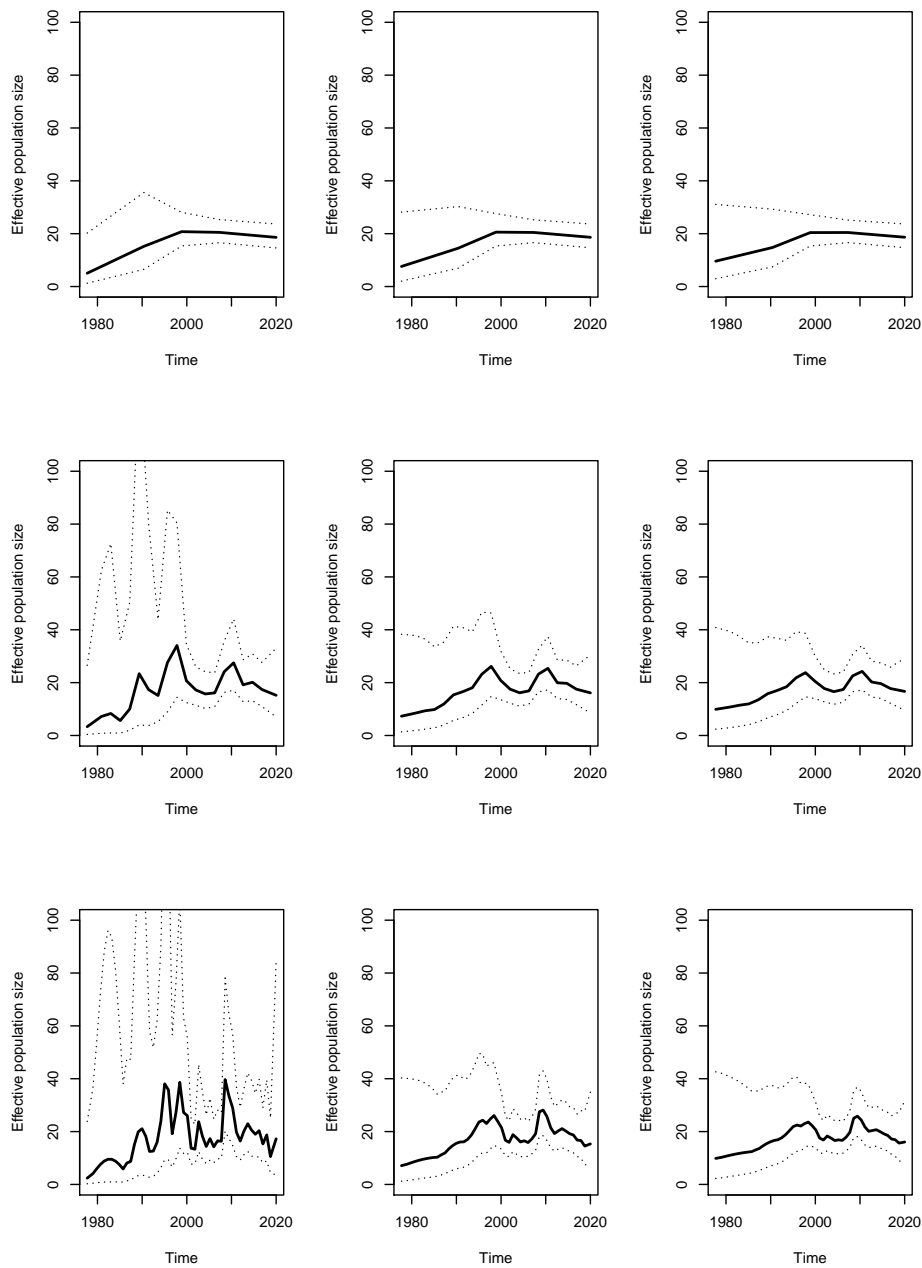


Figure 1: Result on simulated phylogeny shown in Figure S1 using the skyline model, from top to bottom $R = 5, 20, 50$ and from left to right $\tau = 1, 10, 20$.

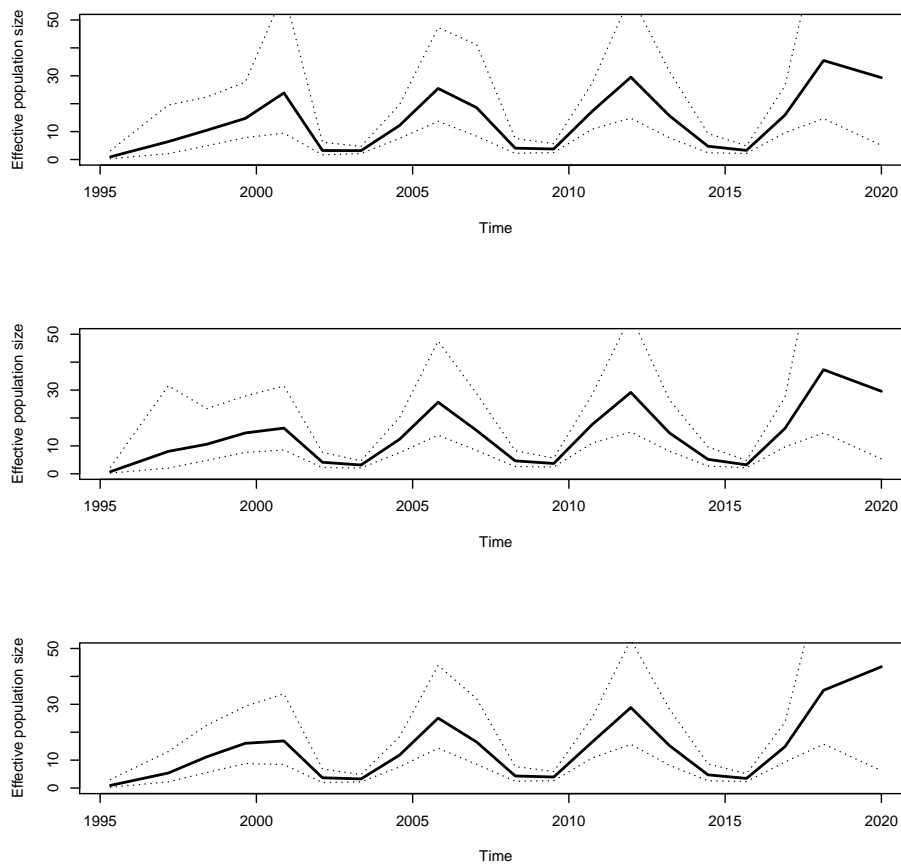


Figure 2: Result of applying the three different models (from top to bottom, *skygrid*, *skygrowth* and *skysigma*) to the phylogeny shown in Figure S2 which was simulated using a sinusoidal demographic function.

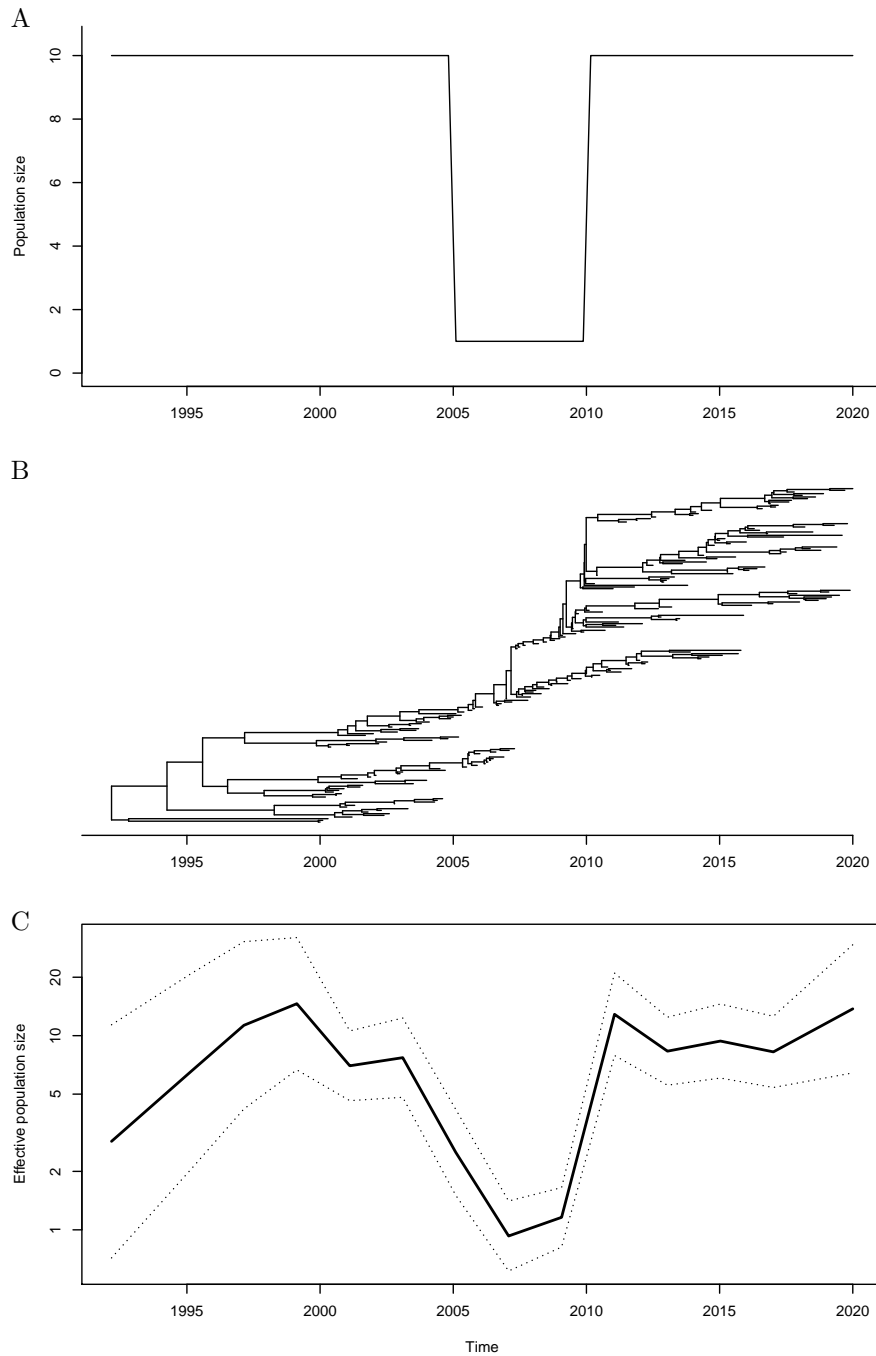
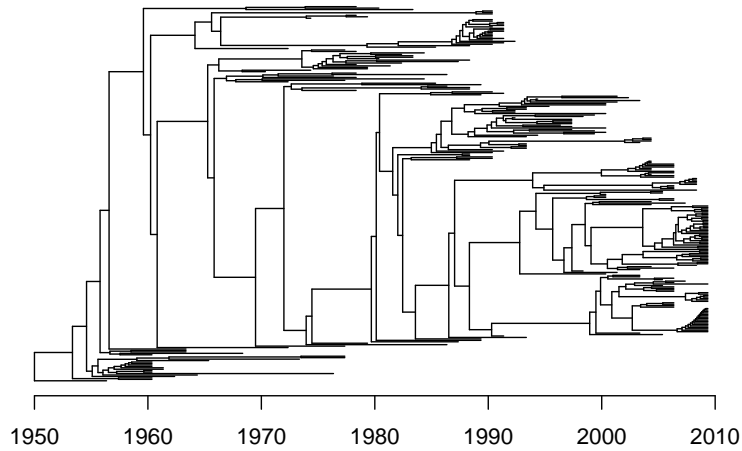


Figure 3: Demographic function (A), phylogeny (B) and inferred demographic function (C) for a simulated dataset under a bottleneck model.

A



B

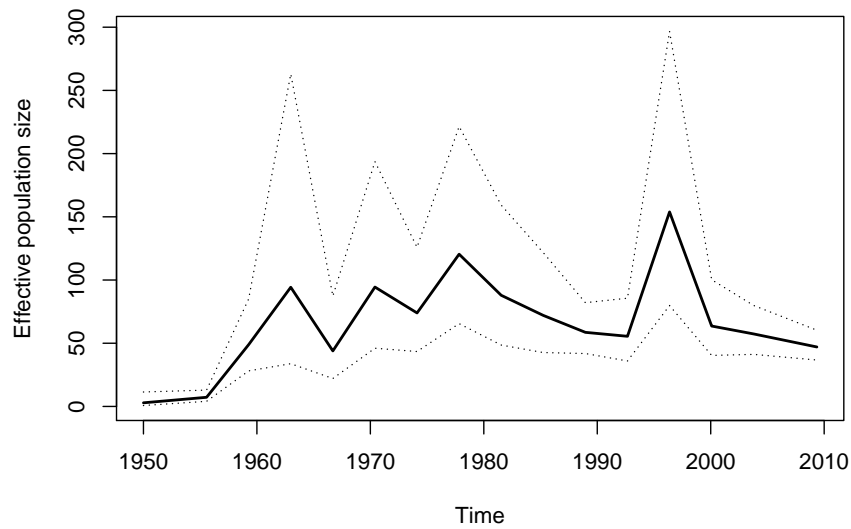


Figure 4: Analysis of the seventh pandemic of *Vibrio cholerae*. (A). Dated phylogeny used as the starting point of past population size inference. (B). Demographic function reconstructed based on the phylogeny above.

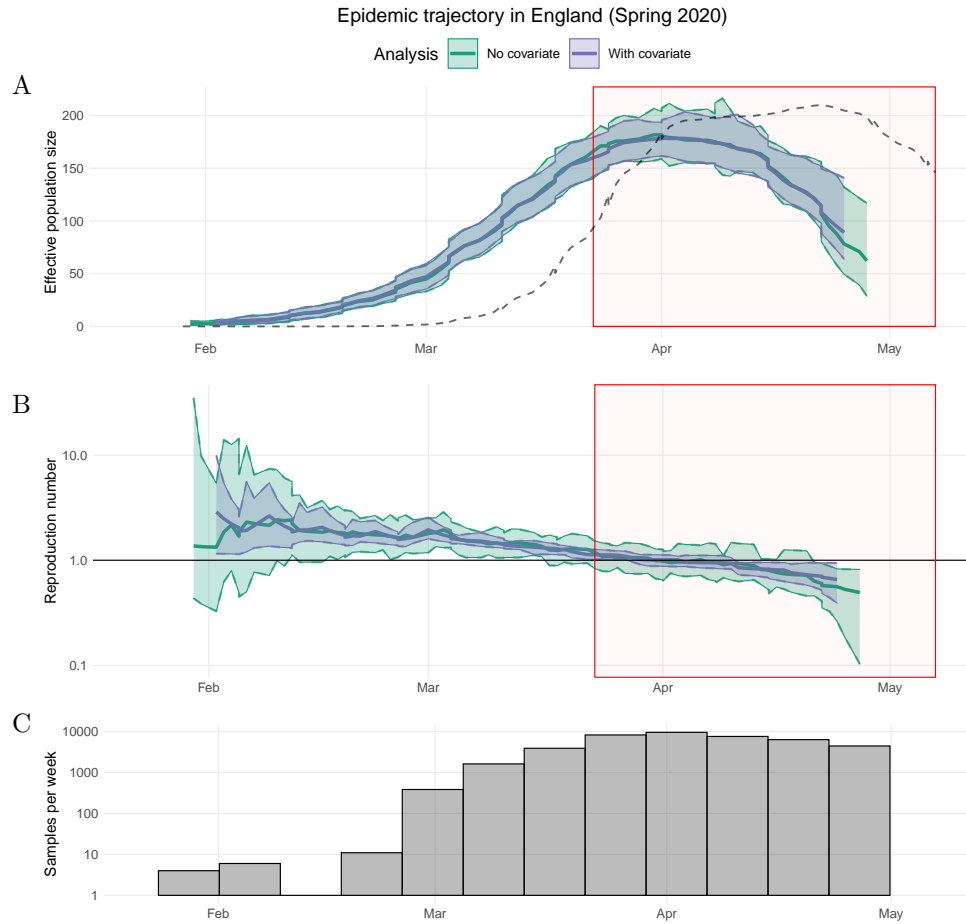


Figure 5: The epidemiological trajectory of SARS-CoV-2 in England during spring 2020. Thick solid lines and shaded areas represent the median and 95% quantiles of $N_e(t)$ with (purple) and without (green) the OxCGRT health containment index as a covariate of $N_e(t)$ growth rates. The model is fitted with no back-shift in the covariate. Red shaded area represents period of first national lockdown in England. Black dotted line represents daily confirmed cases (smoothed and rescaled). (A) Effective population size $N_e(t)$ through time. (B) Reproduction number $R(t)$ through time. (C) Frequencies of sample dates for tips in each sample week in the SARS-CoV-2 phylogenies.

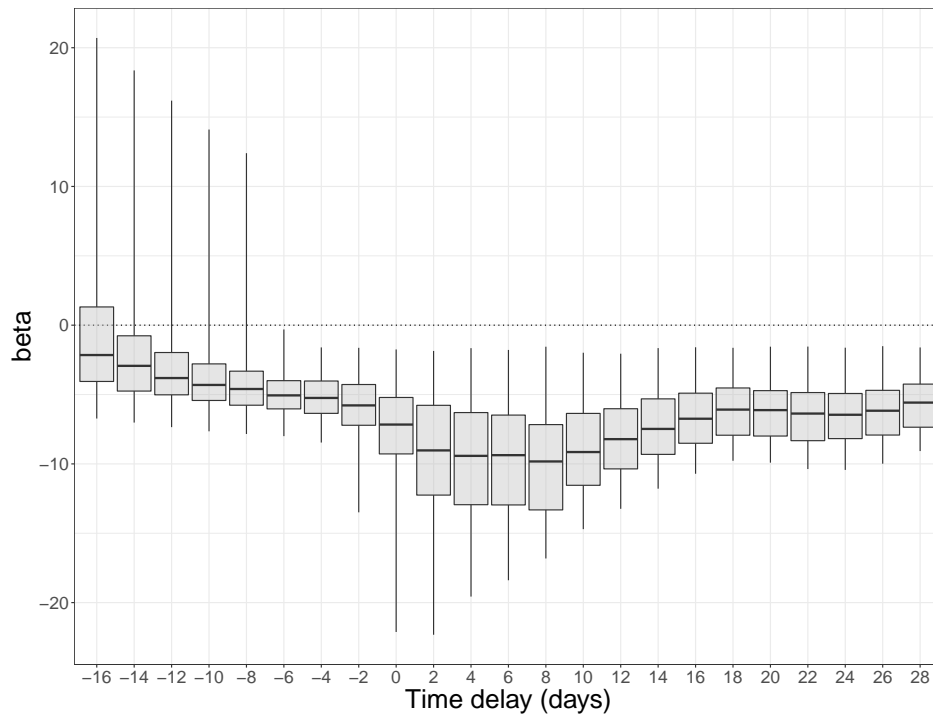


Figure 6: Distribution of association coefficients when testing for univariate association between daily changes in the OxCGRT health containment index and daily changes in the reproduction number of SARS-CoV-2 in England. Boxes represent the median and interquartile range; whiskers show 95% quantiles. A positive delay of 10 represents testing an association between the OxCGRT index at time t and the reproduction number at time $t + 10$ days.

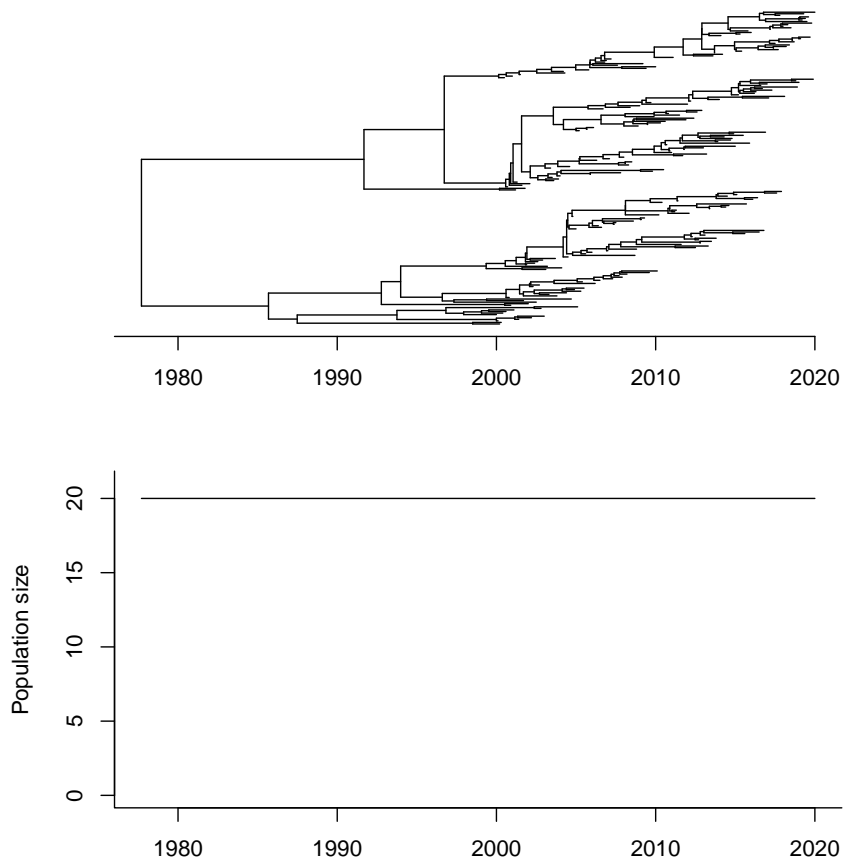


Figure S1: Simulated phylogeny using a constant demographic function.

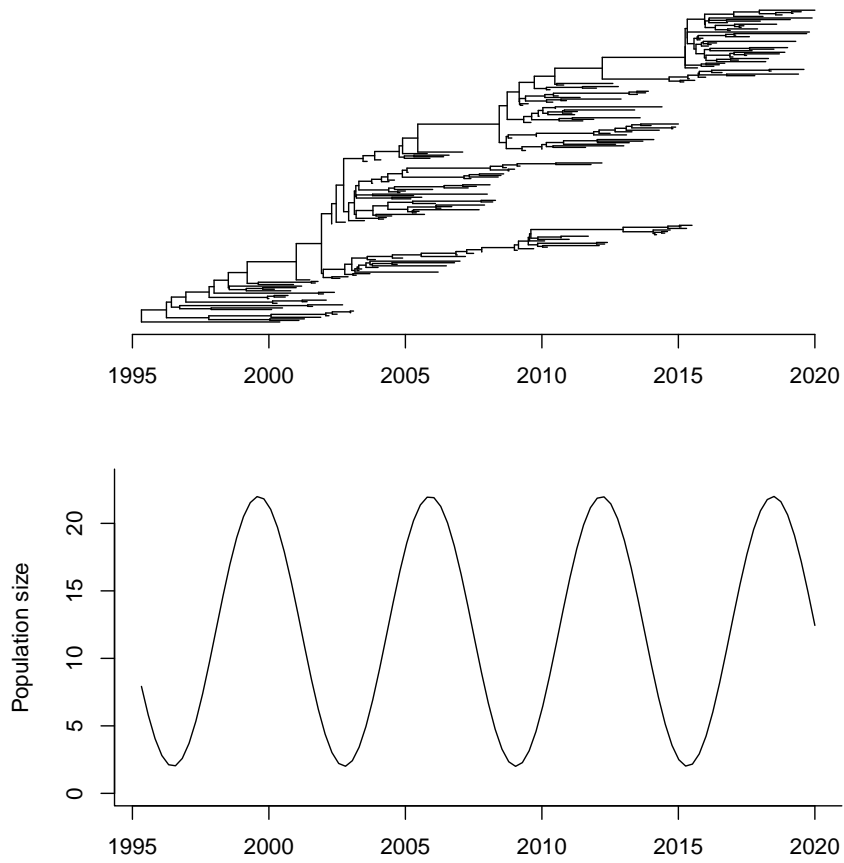
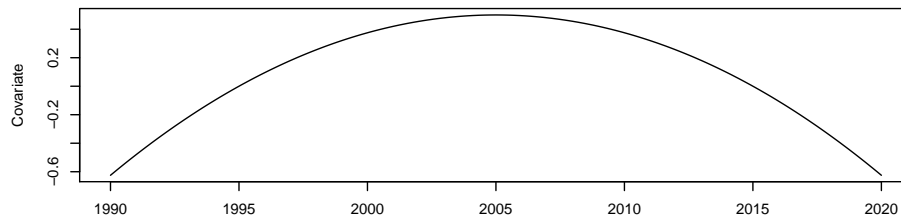
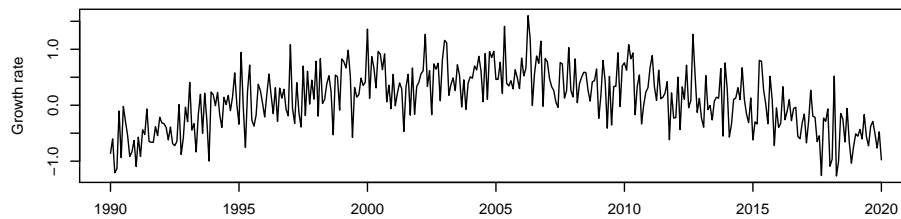


Figure S2: Simulated phylogeny using a sinusoidal demographic function.

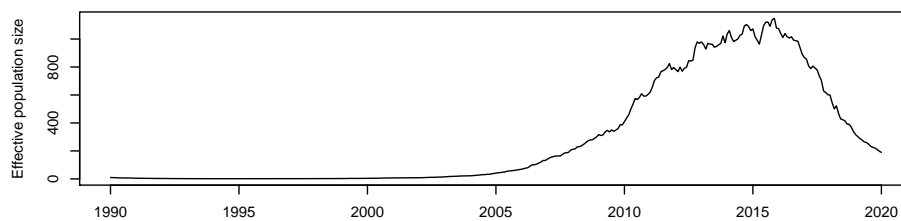
A



B



C



D

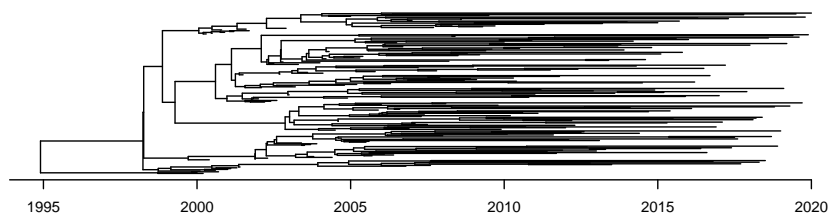


Figure S3: Example of simulation with covariate data driving the growth rate. (A) Covariate data following a quadratic function. (B) Growth rate equal to the covariate data plus some Gaussian noise. (C) Effective population size. (D) Dated phylogeny.

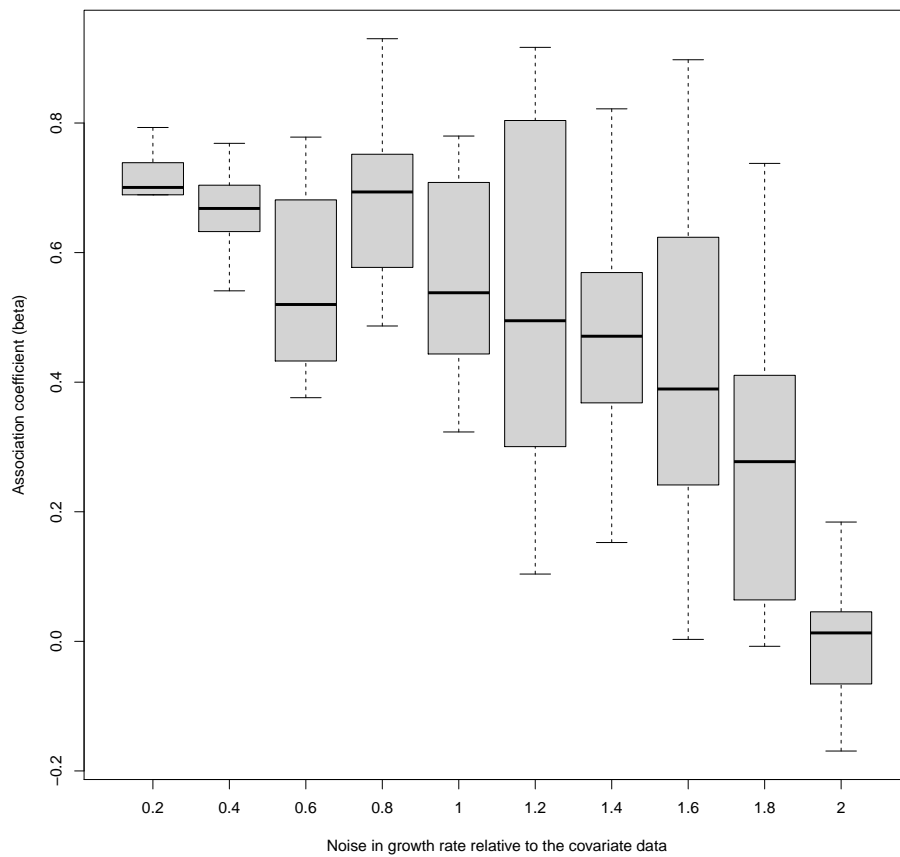


Figure S4: Results of the covariate analysis. For each value of the Gaussian noise (x-axis) ten simulations were performed and the inferred values of the association coefficient β are shown (y-axis) as boxplots.