

Research



Cite this article: Helekal D, Keeling M, Grad YH, Didelot X. 2023 Estimating the fitness cost and benefit of antimicrobial resistance from pathogen genomic data. *J. R. Soc. Interface* **20**: 20230074.

<https://doi.org/10.1098/rsif.2023.0074>

Received: 15 February 2023

Accepted: 22 May 2023

Subject Category:

Life Sciences—Mathematics interface

Subject Areas:

bioinformatics, evolution, systems biology

Keywords:

genomic epidemiology, phylodynamics, antimicrobial resistance, resistance fitness cost

Author for correspondence:

Xavier Didelot

e-mail: xavier.didelot@warwick.ac.uk

Electronic supplementary material is available online at <https://doi.org/10.6084/m9.figshare.c.6673522>.

Estimating the fitness cost and benefit of antimicrobial resistance from pathogen genomic data

David Helekal¹, Matt Keeling², Yonatan H. Grad⁴ and Xavier Didelot³

¹Centre for Doctoral Training in Mathematics for Real-World Systems,²Mathematics Institute and School of Life Sciences, and ³School of Life Sciences and Department of Statistics, University of Warwick, Coventry, UK

⁴Department of Immunology and Infectious Diseases, TH Chan School of Public Health, Harvard University, Boston, MA, USA

MK, 0000-0003-4639-4765; YHG, 0000-0001-5646-1314; XD, 0000-0003-1885-500X

Increasing levels of antibiotic resistance in many bacterial pathogen populations are a major threat to public health. Resistance to an antibiotic provides a fitness benefit when the bacteria are exposed to this antibiotic, but resistance also often comes at a cost to the resistant pathogen relative to susceptible counterparts. We lack a good understanding of these benefits and costs of resistance for many bacterial pathogens and antibiotics, but estimating them could lead to better use of antibiotics in a way that reduces or prevents the spread of resistance. Here, we propose a new model for the joint epidemiology of susceptible and resistant variants, which includes explicit parameters for the cost and benefit of resistance. We show how Bayesian inference can be performed under this model using phylogenetic data from susceptible and resistant lineages and that by combining data from both we are able to disentangle and estimate the resistance cost and benefit parameters separately. We applied our inferential methodology to several simulated datasets to demonstrate good scalability and accuracy. We analysed a dataset of *Neisseria gonorrhoeae* genomes collected between 2000 and 2013 in the USA. We found that two unrelated lineages resistant to fluoroquinolones shared similar epidemic dynamics and resistance parameters. Fluoroquinolones were abandoned for the treatment of gonorrhoea due to increasing levels of resistance, but our results suggest that they could be used to treat a minority of around 10% of cases without causing resistance to grow again.

1. Introduction

The levels of antimicrobial resistance of many pathogens have risen worryingly over the past few decades. In a report on the threat posed by antibiotic resistance published by the CDC (Centres for Disease Control and Protection), three microorganisms including *Neisseria gonorrhoeae* are classified as posing an urgent threat level, and twelve more represent a serious threat to public health [1]. A review on antimicrobial resistance estimated that resistance claims at least 700 000 lives per year worldwide and that the death toll could go up to 10 million per year by 2050 if current trends are allowed to continue [2], and a recent study estimated that there were almost 5 million deaths associated with resistance in 2019 [3]. Few new antimicrobials have been developed and deployed since the 1970s, whereas resistance to new drugs often emerges soon after initial introduction [4], so that several pathogens are dangerously close to becoming completely untreatable. Effectively tackling antimicrobial resistance requires greater understanding of epidemiological and evolutionary factors leading to emergence of resistance and the spread of resistance through pathogen populations. Achieving this goal requires development of mathematical models of antimicrobial resistance and robust statistical analysis

of epidemiological models with informative observations. This modelling approach to resistance was initiated in the late 1990s [5,6] and has led to the development of many models, appropriate for different organisms, mode of spread, study scale and context [7].

Resistance brings a clear fitness benefit to pathogens acquiring it in the presence of antimicrobials. The net value of this fitness benefit therefore increases with the frequency with which the specific antimicrobial is employed, either against the pathogen itself or more generally in the case of a pathogen that can be carried asymptotically. However, resistance also typically comes with a fitness cost to the pathogen [8]. The simplest demonstration of this effect is when discontinued use of an antimicrobial leads to reductions in resistance rates. The fitness costs and benefits of resistance remain poorly understood for many pathogens and antimicrobials [9]. A better quantification of resistance benefits and costs is required to provide a solid basis for evaluating the potential effectiveness of public health intervention measures proposed to exploit fitness costs in the hope of stopping or even reversing the spread of resistance [9]. For example, the numbers of gonorrhoea cases sensitive and resistant to cefixime in England over a decade were recently analysed to quantify the cost and benefit associated with resistance to this antibiotic [10]. These estimates were used to predict that cefixime could be reintroduced to treat a minority (approx. 25%) of gonorrhoea cases without causing an increase in cefixime resistance levels, which would reduce the risk of emergence of resistance to the currently used antibiotics. Moreover, the extent of the fitness cost of resistance can vary by genomic background [11], such that the effect of interventions that seek to capitalize on the fitness costs of resistance may be lineage dependent. Therefore, it is necessary to estimate fitness costs at the per lineage level. The aim of this study is to quantify the contribution that changes in prescription policy have on the population dynamics of particular resistant lineages. This is in contrast to studies that are interested in the overall ecology of resistance or the eventual fate of resistant phenotypes (e.g. [12]).

Pathogen genomic data have great potential to help us understand the evolutionary and epidemiological dynamics of infectious disease [13]. An important advantage of this phylogenetic approach is that analysis of genomic data is less sensitive to sampling biases, especially when using a coalescent framework which describes the ancestry process conditional on sampling [14]. A few studies have used this approach to shed light on the fitness cost associated with antimicrobial resistance. For example, a study showed the association between the growth rate of a methicillin-resistant *Staphylococcus aureus* lineage and consumption of beta-lactams [15]. Other studies quantified the relative transmission fitness of resistance mutations in HIV [16] and *Mycobacterium tuberculosis* [17]. Here, we take a different approach by modelling explicitly the phylogenetic trajectories of the sensitive and resistant lineages as a function of the fitness cost, which is constant, and the fitness benefit, which depends on the antimicrobial consumption. Our method therefore requires three inputs: the amount of antimicrobial being used over time, genomic data from a sensitive lineage and genomic data from a resistant lineage. From this, we disentangle the fitness cost and benefit of resistance, thereby providing the parameters needed to predict phylogenetic trajectories and inform recommendations on how to use antimicrobials without worsening the resistance threat.

Overall, the scenario we are interested in is that of overall resistance dynamics at a large population level. In such a scenario, the bulk of incidence is going to be caused by local transmission rather than imports. We do not intend for the methods presented in this paper to be applicable to small populations dominated by imports and complex, heterogeneous routes of transmission, such as nosocomial infections in a hospital setting. For such a scenario, a different approach using birth–death type models would be more appropriate [16,17].

2. Methods

2.1. Overall approach

Pathogen phylogenetic data contain information about past population size dynamics of the pathogen under study [13,18]. Under assumptions of the epidemic process being characterized well enough by a simple compartmental epidemic model, this information about population size dynamics can be translated into epidemic trajectories [19,20]. These epidemic trajectories can be described using an epidemic model which accounts for the effects of a fitness cost and benefit of resistance to a specific antimicrobial. As the use of this antimicrobial changes through time, so will the net fitness of the particular lineage in consideration. This will in turn lead to changes in the behaviour of the epidemic trajectory. However, not all changes in the behaviour of the epidemic trajectory will be due to changes in the fitness of the resistant phenotype. Confounding factors, such as depletion of susceptibles or changes in host behaviour, will also affect the epidemic trajectory. Under relatively mild assumptions detailed below changes in these confounding factors will affect other lineages equally. We can therefore use as ‘control’ some data from a susceptible lineage, ideally closely related and with the same resistance profile to other antimicrobials used in significant amounts as primary treatment. Differences between the trajectories of the sensitive and resistant lineages can then be ascribed specifically to resistance, allowing us to estimate the associated fitness cost and benefit parameters.

Let us consider a pathogen causing infections at the level of a large population that are or were treated with a certain antimicrobial compound. We assume that at some point in the past one or several lineages with resistance to this antimicrobial compound have arisen. Our aim is to quantify the fitness cost and benefit of the resistance to this antimicrobial for a given lineage as a function of use of the antimicrobial of interest through time. To this end, we need data that quantify the use over time of the given antimicrobial to treat infections caused by this pathogen, as well as a reasonable sample of sequenced case isolates from infections caused by the pathogen over time. Furthermore, we need information that characterizes the resistance profiles of the individual isolates, which can be either obtained by resistance screening *in vitro* or predicted from the sequences *in silico* [21]. A dated phylogeny of these samples is estimated, for example using BEAST [22], BEAST2 [23] or BactDating [24]. This phylogeny is then used as the starting point for analysis [25], to identify which samples belong to resistant and susceptible lineages and to select related lineages for further study that are wholly resistant or susceptible to the antimicrobial of interest, but otherwise similar in their resistance profiles. Note that for simplicity resistance is treated as a binary trait, with samples being either resistant or susceptible to antimicrobials, as is usually the case in resistance modelling studies [7].

2.2. Transmission model derivation

In order to estimate the fitness cost and benefit of antimicrobial resistance, a transmission model needs to be specified. We focus on estimating the fitness parameters of a particular lineage

harbouring a certain treatment resistant phenotype when previous infection does not confer immunity against reinfection. Under the simplifying assumptions that the host population is unstructured and that past infections do not confer any immunity, the multi-lineage susceptible–infected–susceptible (SIS) is a reasonable model [26,27]. This model is more commonly referred to as multi-strain SIS. Fluctuations in the carriage levels of different lineages can also be due to external factors, such as changes in host demography or behaviours. Left unaccounted, such fluctuations would bias estimates of the fitness cost and benefit of resistance to a given antimicrobial. Therefore, we modify the model with time-varying transmission rate $\beta(t)$ and population size $N(t)$. This leads to an n -lineage model described by a system of the following n -coupled ordinary differential equations (ODEs):

$$\left. \begin{aligned} \frac{dI_1(t)}{dt} &= \frac{\beta(t)S(t)I_1(t)}{N(t)} - \gamma_1(t)I_1(t), \\ \frac{dI_2(t)}{dt} &= \frac{\beta(t)S(t)I_2(t)}{N(t)} - \gamma_2(t)I_2(t), \\ &\vdots \\ \text{and } \frac{dI_n(t)}{dt} &= \frac{\beta(t)S(t)I_n(t)}{N(t)} - \gamma_n(t)I_n(t), \end{aligned} \right\} \quad (2.1)$$

where $I_j(t)$ denotes the number of people infected with the j th lineage at time t . $\beta(t)$ is the transmission rate that varies with time due for example to changes that are not specific to any lineage, for example host behaviour. $N(t)$ is the host population size which may also change with time due to demographic factors. $\gamma_j(t)$ is the recovery rate of the j th lineage at time t . These may or may not vary with time through their dependency on the antimicrobial usage which changes with time. Finally, $S(t)$ denotes the number of susceptible hosts:

$$S(t) = \left(N(t) - \sum_{j=1}^n I_j(t) \right). \quad (2.2)$$

Typically, this model could simply be reduced to a two lineage model, averaging over all lineages that are phenotypically similar in their resistance profiles. However, this is undesirable, as some of the lineages with the same resistance phenotype could differ in fitness due to different genomic background which would confound our estimates. Furthermore, this sort of model would not be readily tractable in a genomic framework, because phylogenetic data are generally going to be informative about the dynamics of a particular lineage only. Note that this also means that the analysis produced is valid for the lineages being studied, and cannot be extrapolated to the overall dynamics of resistance for a given pathogen.

We therefore need to focus on the resolution of individual lineages. We note that environmental effects such as fluctuations in host population size or behaviour affect all lineages equally, if the population is well mixed. We denote the combination of these effects as $b(t) = \beta(t)S(t)/N(t)$. Conditional on the knowledge trajectory of $b(t)$ the ODEs in equation (2.1) become uncoupled, and this allows us to reduce the system to uncoupled equations corresponding to the lineage we will be focusing on. As such, we will treat $b(t)$ as a random object that needs to be inferred. We further assume that for the susceptible lineages the average recovery rate denoted γ_s does not change over time, whereas for the resistant lineage it takes one of two values: $\gamma_T = q_T + \gamma_s$ if a given patient is treated with the antimicrobial of interest, or $\gamma_U = q_U + \gamma_s$ otherwise. If we also consider the known proportion of registered cases treated with the antimicrobial of interest $u(t)$, this fully determines the average recovery rate of the resistant lineages as

$$\gamma_r(t) = u(t)\gamma_T + (1 - u(t))\gamma_U. \quad (2.3)$$

We can now fully write down the equations of the model we will be using for the sensitive and resistant lineages, respectively:

$$\left. \begin{aligned} \frac{dI_s(t)}{dt} &= b(t)I_s(t) - \gamma_s I_s(t) \\ \text{and } \frac{dI_r(t)}{dt} &= b(t)I_r(t) - [u(t)\gamma_T + (1 - u(t))\gamma_U]I_r(t). \end{aligned} \right\} \quad (2.4)$$

In practice, we are interested in the difference in recovery rates between the susceptible and the resistant lineages when every case gets treated with the antimicrobial of interest, and when the antimicrobial of interest is not used at all. We denote these by

$$\text{and } \left. \begin{aligned} q_T &= \gamma_T - \gamma_s \\ q_U &= \gamma_U - \gamma_s. \end{aligned} \right\} \quad (2.5)$$

The interpretation is therefore that q_T captures the fitness benefit of resistance in the case $q_T < 0$ and q_U captures the fitness cost of resistance in the case $q_U > 0$.

This model can be applied to any number of resistant and sensitive lineages, simply by adding lineage-associated terms to the likelihood and adding required parameters. This is straightforward as the individual lineages are independent conditional on $b(t)$, but for simplicity the remainder of methods description focuses on the case of a single sensitive and a single resistant lineage, with the general case being a straightforward extension.

2.3. Link to phylogenies

Having defined the epidemiological model, we can now link it to the phylogenetic process. Based on [19,28], the instantaneous coalescent rates for a single pair of lineages can be derived as

$$\lambda_s(t) = \frac{2b(t)}{I_s(t)} \quad \text{and} \quad \lambda_r(t) = \frac{2b(t)}{I_r(t)} \quad (2.6)$$

in the susceptible and resistant populations, respectively. The likelihood of a dated phylogeny \mathbf{g} with n leaves at times $s_1 < \dots < s_n$ and $n - 1$ coalescent events at times $c_1 < \dots < c_{n-1}$ and $A(t)$ lineages at time t is therefore given by Griffiths & Tavaré [29]:

$$p(\mathbf{g}|\lambda(t)) = \exp\left(-\int_{-\infty}^{\infty} \mathbb{1}[A(t) \geq 2] \binom{A(t)}{2} \lambda(t) dt\right) \prod_{i=1}^{n-1} \lambda(c_i), \quad (2.7)$$

where $\lambda(t) = \lambda_s(t)$ and $\lambda(t) = \lambda_r(t)$ for the susceptible and resistant phylogenies, respectively. However, in most cases, and indeed in our case, the integral in equation (2.7) is not analytically intractable. Furthermore, the antibiotic use data are unlikely to span the entire phylogeny. Therefore, we define the approximate likelihood for the phylogeny truncated to $[t_{\min}, t_{\max}]$, which is the intersection interval spanned by the antibiotic use data and the phylogenies under study.

As such we resort to the standard way of approximating coalescent likelihoods [30], partitioning the interval $[t_{\min}, t_{\max}]$ into a fine mesh $t_{\min} = t_1 < t_2 < t_3 < \dots < t_N = t_{\max}$ such that $t_i - t_{i-1} < \Delta_t$ and that all sampling and coalescent times between t_{\min} and t_{\max} are included in the mesh:

$$p(\mathbf{g}|\lambda(t)) = \exp\left(-\sum_{i=2}^N (t_i - t_{i-1}) \binom{A(t_{i-1})}{2} \lambda(t_{i-1})\right) \prod_{i=1}^{n-1} \mathbb{1}[c_i \in [t_{\min}, t_{\max}]] \lambda(c_i). \quad (2.8)$$

We note that the approach of how we treat the relationship between the phylogenies and epidemic is effectively a structured coalescent with no migration and time varying $N_e(t)$ determined by the deterministic epidemic model. Approaches reminiscent of ours have been used to formally study the expected age of a mutation in both the presence or absence of selection [31]. However, in that case the populations correspond to different alleles, and the $N_e(t)$

Table 1. Summary of the parameters and priors used in the model.

parameter	symbol	prior
susceptible lineage recovery rate	γ_s	$\text{lognormal}(\log \gamma^*, \sigma)$
resistant lineage recovery rate if treated with focal antibiotic	γ_T	$\text{normal}(\gamma_s, 0.3\gamma^*)\mathbb{1}[x > 0]$
resistant lineage recovery rate if treated with other antibiotic	γ_U	$\text{normal}(\gamma_s, 0.3\gamma^*)\mathbb{1}[x > 0]$
initial prevalence of sensitive lineage	I_{s0}	$\text{lognormal}(6, 2)$
initial prevalence of resistant lineage	I_{r0}	$\text{lognormal}(6, 2)$
GP kernel marginal variance	α	$\text{gamma}(4, 4)$
GP kernel length scale	ρ	$\text{inverse} - \text{gamma}(4.63, 2.21)$
approximate GP functions	$f_{1:m}$	$\mathcal{N}(0, 1)$

curves follow the proportion of population with a given allele as determined by Wright–Fisher diffusion forwards in time. Migration between the demes corresponding to individual alleles can also further be added corresponding to recombination [32].

2.4. Bayesian inference

We first re-scale time from the interval $[t_{\min}, t_{\max}]$ to $[-1, 1]$. Denoting the scale factor $D = (t_{\max} - t_{\min})/2$ associated with this re-scaling, we account for this in the model by defining $\tilde{\gamma}_s = \gamma_s D$.

The model consists of independent first-order linear homogeneous ODEs for each lineage with time-varying coefficients. The solutions at time t subject to initial conditions $I_s(0) = I_{s0}$ and $I_r(0) = I_{r0}$ can be obtained in terms of the integral of the instantaneous rates up to time t :

$$\left. \begin{aligned} I_s(t) &= I_{s0} \exp\left\{\int_0^t b(\tau) - \gamma_s d\tau\right\} \\ \text{and } I_r(t) &= I_{r0} \exp\left\{\int_0^t b(\tau) - [u(\tau)\gamma_T + (1-u(\tau))\gamma_U] d\tau\right\}. \end{aligned} \right\} \quad (2.9)$$

As it stands, this model would not be well suited for performing inference under, primarily due to the difficulty in choosing a sensible prior on $b(t)$, and a very complicated dependency structure between the initial conditions and $b(t)$. As such we re-parameterize the model by directly modelling the logarithm of $I_s(t)$ as a Gaussian process:

$$C(t) = \log I_s(t) - \mu_s, \quad (2.10)$$

where $C(t)$ is an appropriately chosen zero mean Gaussian process, and μ_s is the susceptible intercept which relates to the susceptible initial condition I_{s0} as follows:

$$\mu_s = \log I_{s0} - C(0). \quad (2.11)$$

We use this formulation principally to loosen the coupling between the intercept parameter and the Gaussian process in order to speed up sampling. From this, we can compute $b(t)$ and $\log I_r(t)$ as

$$b(t) = \frac{d}{dt}C(t) + \gamma_s \quad (2.12)$$

and

$$\begin{aligned} \log I_r(t) &= C(t) + \mu_r + \int_0^t \gamma_s d\tau - \int_0^t u(\tau)\gamma_T d\tau - \int_0^t (1-u(\tau))\gamma_U d\tau \\ &= C(t) + \mu_r + \int_0^t \gamma_s - u(\tau)(\gamma_T - \gamma_U) - \gamma_U d\tau \\ &= C(t) + \mu_r + (\gamma_s - \gamma_U)t - (\gamma_T - \gamma_U) \int_0^t u(\tau) d\tau. \end{aligned} \quad (2.13)$$

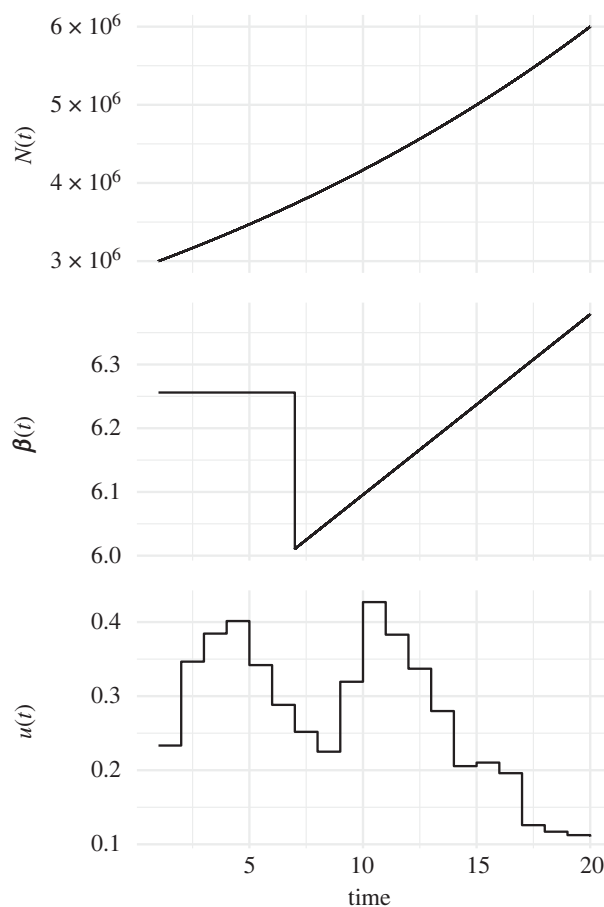


Figure 1. Host population size function $N(t)$, transmission rate over time $\beta(t)$ and antibiotic usage function $u(t)$ used in the simulated datasets.

Once again we follow the same reasoning for the resistant trajectory intercept μ_r , relating it to I_{r0} as

$$\mu_r = \log I_{r0} - C(0). \quad (2.14)$$

Note that $(d/dt)C(t)$ exists as long as the associated covariance kernel is sufficiently smooth such as in the case of the radial basis function (RBF) kernel [33] which we used. Evaluating a full-rank, Gaussian process with differentiable trajectories on the entirety of the mesh would be prohibitively expensive due to the $O(n^3)$ computational complexity, where n is the number of grid points. Such a high computational cost would make the model infeasible. Instead, we work with a low-rank representation of $C(t)$ based on the framework introduced in [34]. This leads to the representation of the low-rank projection of $C(t)$,

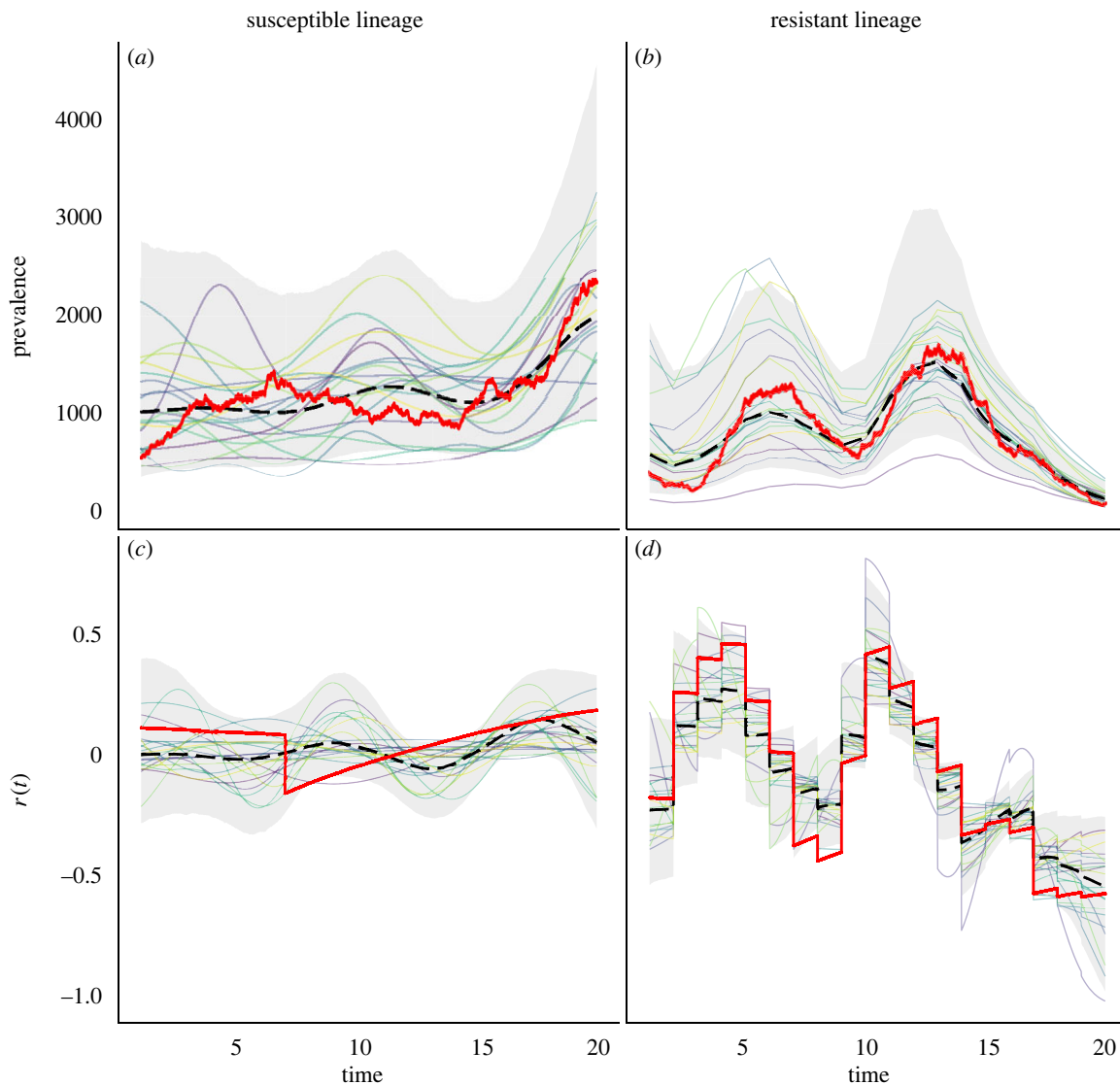


Figure 2. Posterior summary of dynamics for the sensitive (*a,c*) and resistant (*b,d*) lineages, showing prevalence (*a,b*) and reproduction number (*c,d*). Bold solid red lines indicate simulated values. Posterior median in bold dashed black line. Shaded bands indicate 95% posterior credible intervals. Solid light lines represent posterior draws.

denoted by $\hat{C}(t)$

$$\hat{C}(t) = \sum_{j=1}^m S_{\text{RBF}}\left(\sqrt{\frac{j\pi}{2L}}; \rho, \alpha\right) \sqrt{\frac{1}{L}} \sin\left(\frac{j\pi}{2L}(t+L)\right) f_j \quad (2.15)$$

and

$$\frac{d}{dt} \hat{C}(t) = \sum_{j=1}^m S_{\text{RBF}}\left(\sqrt{\frac{j\pi}{2L}}; \rho, \alpha\right) \sqrt{\frac{1}{L}} \frac{j\pi}{2L} \cos\left(\frac{j\pi}{2L}(t+L)\right) f_j. \quad (2.16)$$

This reduces the evaluation complexity of the Gaussian process prior from $O(n^3)$ to $O(nm)$. L and m are approximation parameters that need to be specified *a priori* (see [34] for details). In practice, we used the Hilbert space Gaussian process (HSGP) approximation with parameters $L=6.5$ and $m=60$. These approximation parameters are appropriate for the 99% interval of the length-scale prior used as per [34]. Here f_j are independent and identically distributed random variables following the standard Gaussian distribution, $S_{\text{RBF}}(\cdot; \cdot, \cdot)$ is the appropriate spectral density for the RBF kernel, ρ is the kernel length scale and α is the marginal standard deviation of the kernel [34].

Denote by $\theta = (\gamma_s, \gamma_U, \gamma_T, I_{s0}, I_{r0}, \hat{C}(t))$ the parameters of the pathogen dynamics model. We can now factorize the model

posterior $\pi(\theta, \alpha, \rho, f_{1:m} | \mathbf{g}_s, \mathbf{g}_r)$, suppressing dependency on t where appropriate:

$$\pi(\theta, \alpha, \rho, f_{1:m} | \mathbf{g}_s, \mathbf{g}_r) \propto \pi(\mathbf{g}_s | \lambda_s) \pi(\mathbf{g}_r | \lambda_r) \pi(\lambda_s | \theta) \pi(\lambda_r | \theta) \pi(\theta, \alpha, \rho, f_{1:m}). \quad (2.17)$$

The first two terms are computed using the coalescent likelihood in equation (2.7). The third term is given by combining equations (2.6), (2.10) and (2.12). The fourth term is obtained by combining equations (2.6), (2.12) and (2.13). Finally, the last term is given by

$$\pi(\theta, \alpha, \rho, f_{1:m}) = \pi(\hat{C}(t) | \alpha, \rho, f_{1:m}) \pi(\gamma_T | \gamma_s) \pi(\gamma_U | \gamma_s) \pi(\gamma_s) \pi(I_{s0}) \pi(I_{r0}) \pi(\alpha) \pi(\rho) \pi(f_{1:m}), \quad (2.18)$$

where the first term is given by the Gaussian process (equations (2.15) and (2.16)) and the remaining terms correspond to the prior distributions listed below.

2.5. Choice of prior and parameterization

The model is parameterized with the priors summarized in table 1. The data are not expected to be very informative about the value of γ_s . As such, we impose a fairly informative prior on this parameter, centred around a guess γ^* which must be known and supplied *a priori*. σ then governs how informative the prior is. We typically

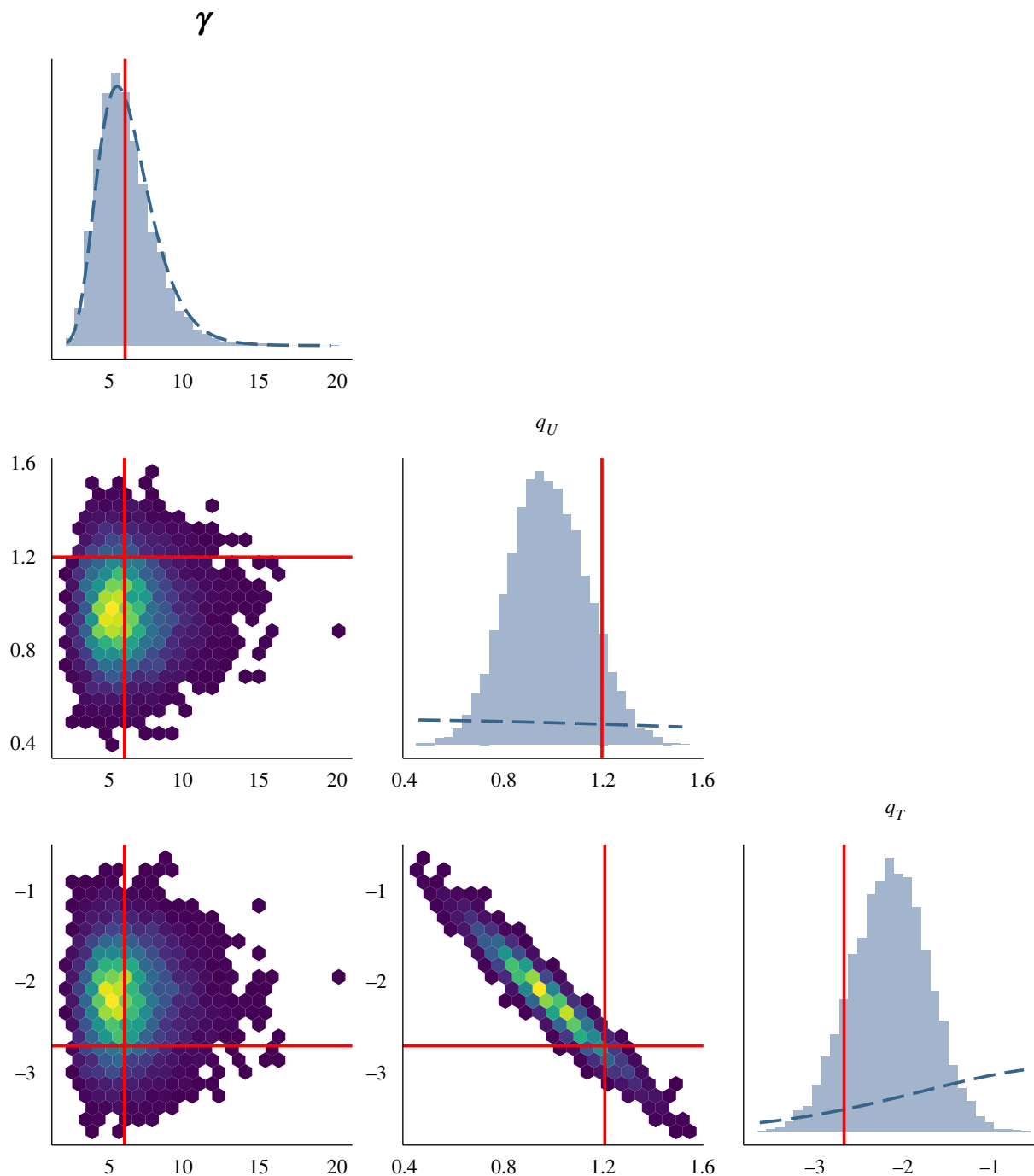


Figure 3. Marginal and joint posterior distributions for the recovery rate of the sensitive lineage (γ_s), fitness cost (q_U) and fitness benefit (q_T) of resistance. Bold red solid lines indicate simulation values. Bold blue dashed lines indicate prior density values.

use a value of $\sigma = 0.3$, which includes relative fluctuations of over 50% in its 95% interval. The higher the value of σ , the more complicated the geometry and subsequently sampling of the posterior becomes. γ_T and γ_U represent the recovery rates for the resistant lineage when the resistant lineage is treated with the focal antibiotic of interest, or another antibiotic, respectively. A normal distribution centred at γ_s and truncated to positive values only is a natural choice. We choose its standard deviation to be $0.3\gamma_s^*$ as this puts greater than 99% of the weight within $2\gamma_s^*$ thus making implausibly large fluctuations unlikely. Such large fluctuations are hardly of interest here since they would lead to a very rapid selective sweep or extinction. The recovery rates γ_T and γ_U are related to the absolute changes in recovery and therefore fitness parameters using equation (2.5). $\gamma_U > \gamma_s$ corresponds to faster recovery when the resistant lineage is treated with an antimicrobial it is sensitive to and therefore a cost of resistance. $\gamma_T < \gamma_s$ corresponds to slower recovery when the resistant lineage is treated with the antimicrobial of interest and therefore a

benefit of resistance. If instead a large proportion of posterior probability mass has $\gamma_U < \gamma_s$ or $\gamma_T > \gamma_s$, we conclude that the result is consistent with either the cost or the benefit of resistance not being significantly present. The prior on ρ was chosen so that approximately 1% of mass lies on values of $\rho < 0.2$ and approximately 1% of mass lies on $\rho > 2$. The lower bound was chosen to avoid overfitting, and the upper bound to suppress length scales that exceed the range of data and thus cannot be informed about by the data.

In practice, due to our choice of a sampling approach we need to parameterize γ_U and γ_T on an unconstrained space, and ideally also weaken the dependency on γ_s . To do so, we introduce parameters \tilde{q}_U and \tilde{q}_T , and define γ_U and γ_T to be a deterministic transformation of these:

$$\left. \begin{aligned} \gamma_U &= \log(1 + \exp\{\tilde{q}_U + \log \gamma_s\}) \\ \gamma_T &= \log(1 + \exp\{\tilde{q}_T + \log \gamma_s\}) \end{aligned} \right\} \quad (2.19)$$

and

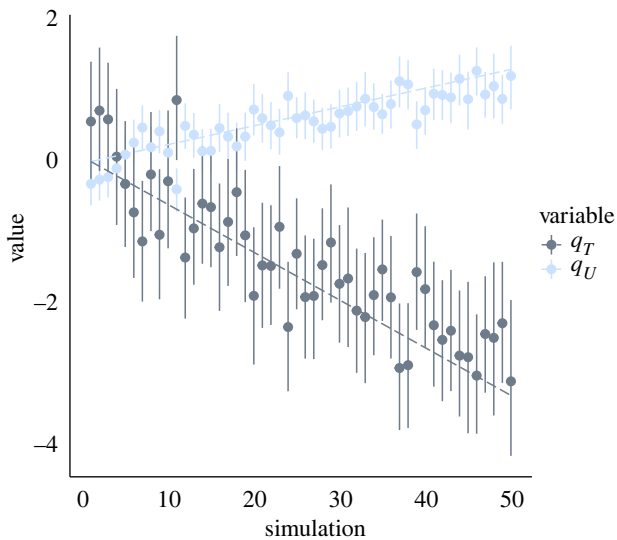


Figure 4. Inferred parameters versus correct values. A total of 50 simulated datasets were generated, with decreasing values of q_T and increasing values of q_U as shown by the dotted lines in grey and blue, respectively. For each simulated dataset, we applied our inference method. The grey and blue dots show the mean inferred values of q_T and q_U , respectively, with vertical bars representing the 95% credible intervals for both parameters.

The Jacobian adjustment to the likelihood associated with this transformation is proportional to

$$|\det J_q| \propto (1 + \exp\{-\tilde{q}_T - \log \gamma_s\})^{-1} (1 + \exp\{-\tilde{q}_U - \log \gamma_s\})^{-1}. \quad (2.20)$$

2.6. Computational implementation

The posterior in equation (2.17) is a high-dimensional distribution and we expect many parameters to have a high degree of interdependency. In order to sample from this distribution, we use dynamic Hamiltonian Monte Carlo (HMC), a HMC sampler available in Stan [35]. HMC is a Markov chain Monte Carlo approach that due to possessing energy conserving properties is able to take large steps between individual states while maintaining high acceptance rates. This makes it efficient at sampling from moderately high dimensional posterior distributions with differentiable likelihoods, while requiring a much lower number of iterations. We implemented the model and inference method in an R package which is available at <https://github.com/dhelekal/ResistPhy/>. All results shown used four chains with 2000 iterations for warmup and 2000 iterations for sampling. For all model parameters and all analysis, the bulk effective sample size (bulk-ESS) was always greater than 500, and all \hat{R} statistics were lower than 1.05 [36], values that indicate no issues with mixing. We also checked that there were no divergent transitions at least during the sampling phase.

2.7. Use of simulated and real datasets

For all simulations, we use a stochastic, discrete state-space version of the multi-lineage SIS in equation (2.1). The system is simulated using tau-leaping [37]. More specifically, we consider a scenario with three lineages simulated over the course of 19 years. Two lineages are set to be susceptible and thus unaffected by antibiotic usage fluctuations and one is set to be resistant. The first lineage aims to represent the unobserved bulk of the population and thus is set to start at much higher prevalence. Conditional on the trajectories of the two lineages, we sample phylogenies under Kingman's coalescent with varying effective population size $N_e(t)$ following equation (2.6) conditional on

the trajectories [28]. The parameters for the simulation were selected as to consistently provide a reasonable range of plausible behaviours so that resistant lineages would reach prevalence with orders of magnitude between 10^2 and 10^4 .

A total of 1102 genomes were collected between 2000 and 2013 by the CDC Gonococcal Isolate Surveillance Project (GISP) [38]. A maximum-likelihood phylogeny was computed using PhyML [39], which was corrected for recombination using ClonalFrameML [40] and dated using BactDating [24]. This dated phylogeny is the same as previously used in an analysis of hidden population structure [41]. The distribution of primary antimicrobial drugs used to treat gonorrhoea among participants of the GISP between 1988 and 2019 was obtained from the GISP reports available at <https://www.cdc.gov/std/statistics/archive.htm>. Note that usages of ciprofloxacin and ofloxacin were combined into a single fluoroquinolone category. All the data and code used in the simulated and real dataset analyses are available at <https://github.com/dhelekal/ResistPhy/tree/main/run>.

3. Results

3.1. Detailed analysis of a single simulated dataset

To validate the performance of this model, we first resort to simulation from a three-lineages stochastic SIS with population size $N(t)$, transmission rate $\beta(t)$ and antimicrobial usage function $u(t)$ varying over the past 20 years, as illustrated in figure 1. The first two lineages are susceptible and thus unaffected by fluctuations in antimicrobial usage, whereas the third lineage is resistant and therefore affected. The first lineage represents the bulk of the susceptible lineages and is thus left unobserved. The remaining two lineages represent the observed lineages, susceptible and resistant, respectively. The per-day recovery rate of the sensitive lineage was set to $\gamma_s = 1/60$, the fitness cost of resistance to $q_U = 1.25$ and the fitness benefit of resistance to $q_T = -2.7$. From each of these two observed lineages, a dated phylogeny with 200 leaves was simulated. The sampling dates were randomly assigned to one of the first 6 years, with the relative probability of a particular year being chosen proportional to the total prevalence in that year. We performed inference on this simulated dataset; the traces are shown in electronic supplementary material, figure S1, and the posterior distribution of the kernel parameters in electronic supplementary material, figure S2. The prevalence and reproduction number $R(t)$ of both the susceptible and resistant lineages are shown in figure 2. As expected, the inferred values followed the correct values used in the simulation. The inferred values of the susceptible lineage recovery rate γ_s and the cost and benefit of resistance q_U and q_T were also found to be close to their correct values, as shown in figure 3. The posterior distribution of γ_s was almost identical to the prior, which was centred on the correct value $1/60$, reflecting the fact that the data are uninformative about this parameter and stressing the importance of using an informative prior. There was a strong negative correlation between the inferred values of q_U and q_T , as expected since these two parameters play opposite roles in the overall fitness of the resistant lineage relative to the sensitive lineage. Nevertheless, we detected both the cost and the benefit associated with resistance, since the ranges of inferred values for q_U and q_T were, respectively, above and below one, contrary to their lognormal priors with mean one (figure 3). Finally, we computed the posterior predictive distribution [42] for the

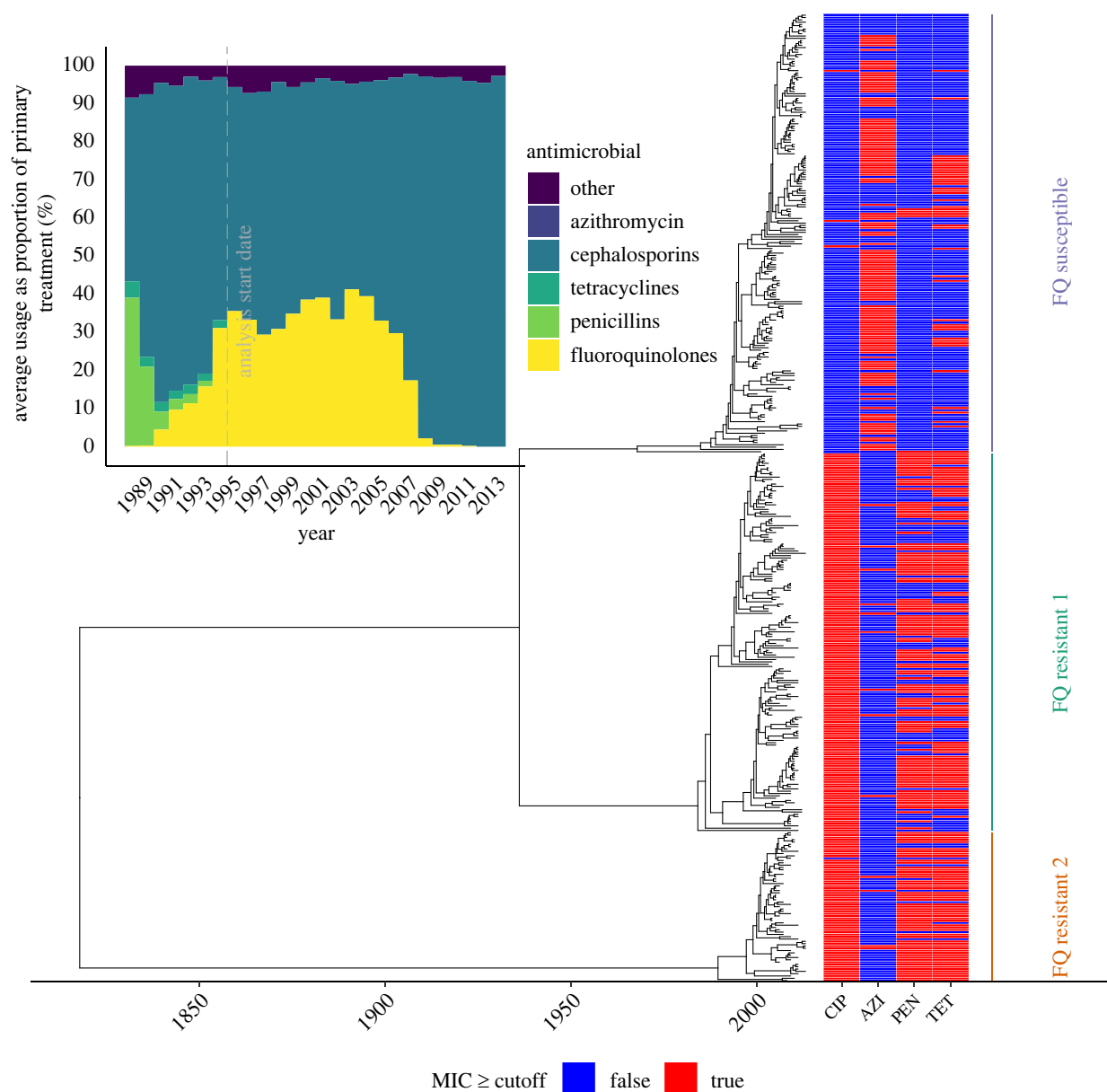


Figure 5. Antibiotic usage data and phylogeny used for the application to fluoroquinolone resistant *N. gonorrhoeae*.

number of ancestral lineages through time $A(t)$ and compared this with the input phylogenetic data (electronic supplementary material, figure S3). The data and posterior predictive trajectories were similar, indicating a good fit of the model to the data as indeed would be expected here since the same model was used for simulation and inference.

3.2. Benchmark using multiple simulated datasets

We repeated the same application of our inference method to data simulated in the same conditions as described above and illustrated in figure 1, except the values of the fitness cost and benefit of resistance were varied. A total of 50 simulated datasets were generated and analysed, with the fitness cost q_U increasing linearly from 1 to 1.2, and the fitness benefit q_T decreasing linearly from 1 to 0.5. The prevalences of the susceptible and resistant lineages in these simulations are shown in electronic supplementary material, figure S4. The results of inference are illustrated in figure 4 and show that in almost all cases, the posterior 95% credible intervals covered the correct values of the fitness cost and benefit of resistance used in the simulations.

3.3. Application to fluoroquinolone resistant *N. gonorrhoeae* in USA

We demonstrate the use of our model and inferential framework by estimating the cost and benefit of fluoroquinolone resistance in *N. gonorrhoeae*. Based on the 1102 genomes collected between 2000 and 2013 by the CDC GISP [38], a recombination-corrected tree was constructed using Clonal-FrameML [40] and dated using BactDating [24]. As there are two major fluoroquinolone resistant lineages present in this phylogeny [38], we decided to do a comparative study. The two fluoroquinolone resistant lineages and one fluoroquinolone susceptible lineage were selected based on similar resistance profiles against other relevant antibiotics. By inspecting the antibiotic usage data and the resistance profiles for the three lineages (figure 5), we can see that the resistance profiles match for antimicrobials that were in use as primary treatment at significant levels after 1995. As such, this is the year we set as the analysis start date ($t_{\min} = 1995$) and the end date is the date when the last genomes were collected ($t_{\max} = 2013$). Note that a subclade within the susceptible lineage that displayed a de novo gain of resistance

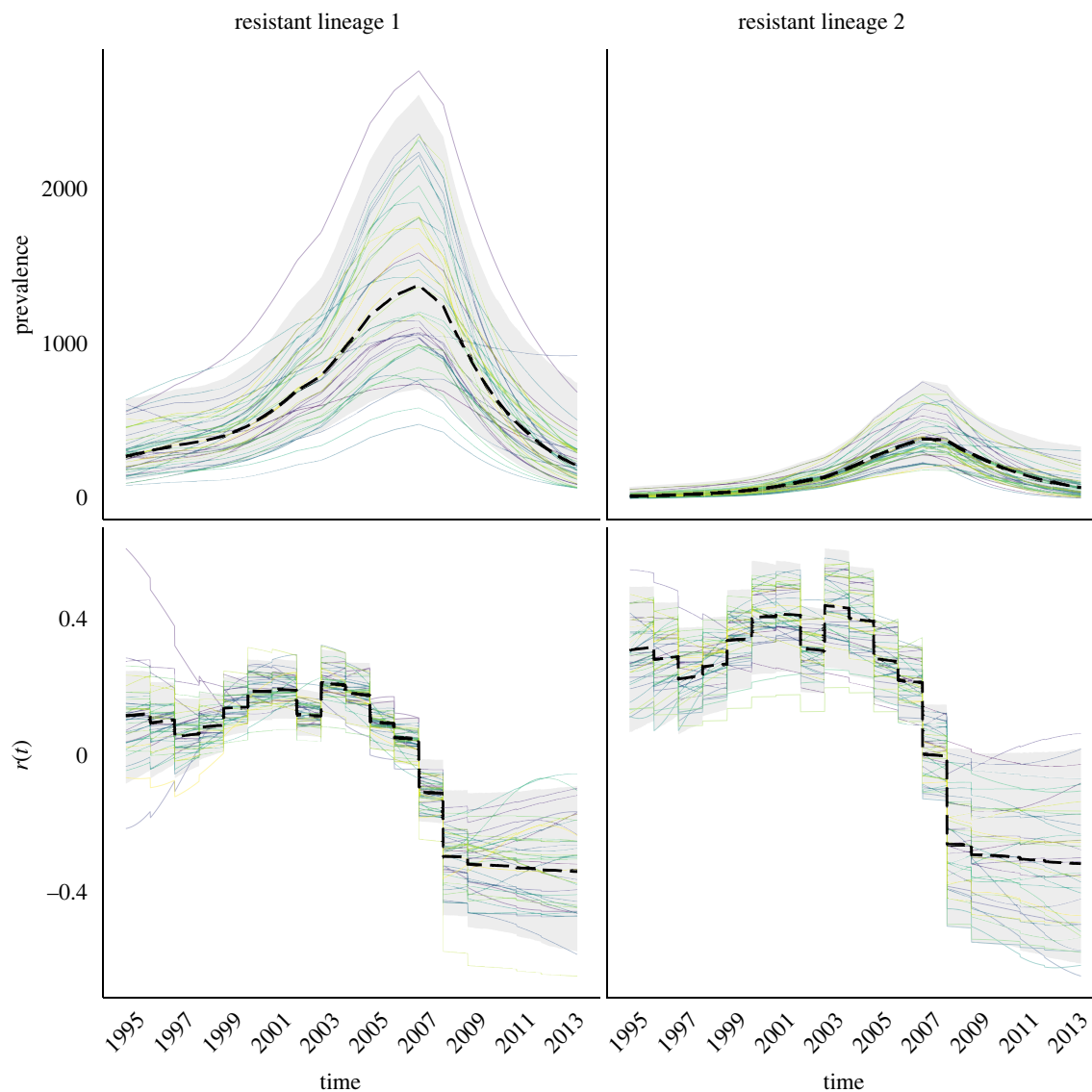


Figure 6. Posterior epidemic dynamics for both fluoroquinolone resistant lineages of *N. gonorrhoeae*.

to cefixime has been removed. The prior mean for the per-day recovery rate for the susceptible lineage was set to $\gamma^* = 1/90$ based on previous gonorrhoea modelling studies [10,43,44].

We performed inference for this dataset; the traces are shown in electronic supplementary material, figure S5, and the posterior distribution of kernel parameters in electronic supplementary material, figure S6. Figure 6 depicts the summary of posterior latent transmission dynamics for the two resistant lineages, whereas electronic supplementary material, figure S7, shows the same for the susceptible lineage. The two resistant lineages have similar dynamics, with a peak in prevalence around 2007, which corresponds to the moment when fluoroquinolone use dropped (figure 5). Figure 7 depicts the marginal and joint posterior distributions for the resistance parameters q_U and q_T for both resistant lineages. This is consistent with there being both a cost and benefit to fluoroquinolone resistance for both lineages, since both q_T and q_U are, respectively, localized below 1 and above 1, with high posterior probability. It is noteworthy that while both of these lineages come from distinct genetic background, their resistance profile is qualitatively very similar, indicating both of these lineages faced similar selective pressures and neither seems to have successfully adapted to overcome the fitness cost associated with fluoroquinolone

resistance. We used a posterior predictive approach to ensure that the model can explain the data appropriately [42]. Posterior predictive trajectories for the function of ancestral lineages through time $A(t)$ were simulated and found to be very similar to the ones implied by the phylogenetic data (electronic supplementary material, figure S8).

Under the assumption of perfect competition between lineages, if we want to ensure that a resistant lineage cannot establish, and its proportion decays sufficiently fast, we fix a decay factor $c > 0$ and aim to ensure that the growth rate of the resistant lineage is c units lower than that of the sensitive lineage, that is $r_s(t) - r_r(t) > c$. Note that $r(t)$ is the growth rate through time, not $R(t)$, the time varying reproduction number. We choose to work with growth rates as these are less sensitive to susceptible recovery rate misspecification. Given that the lineages have the same transmission rate function $b(t)$, this condition is equivalent to $\gamma_s(t) - \gamma_r(t) > c$, and using the definition of $\gamma_r(t)$ from equation (2.3), this is equivalent to $u(t)q_T + (1 - u(t))q_U > c$. We use this to estimate posterior probabilities. The differences in growth rates between the susceptible lineage and each of both resistant lineages exceed c as shown in figure 8. In order to be 95% certain that the resistant lineages remain at a lower fitness than the susceptible lineage, fluoroquinolone should not be prescribed to more than

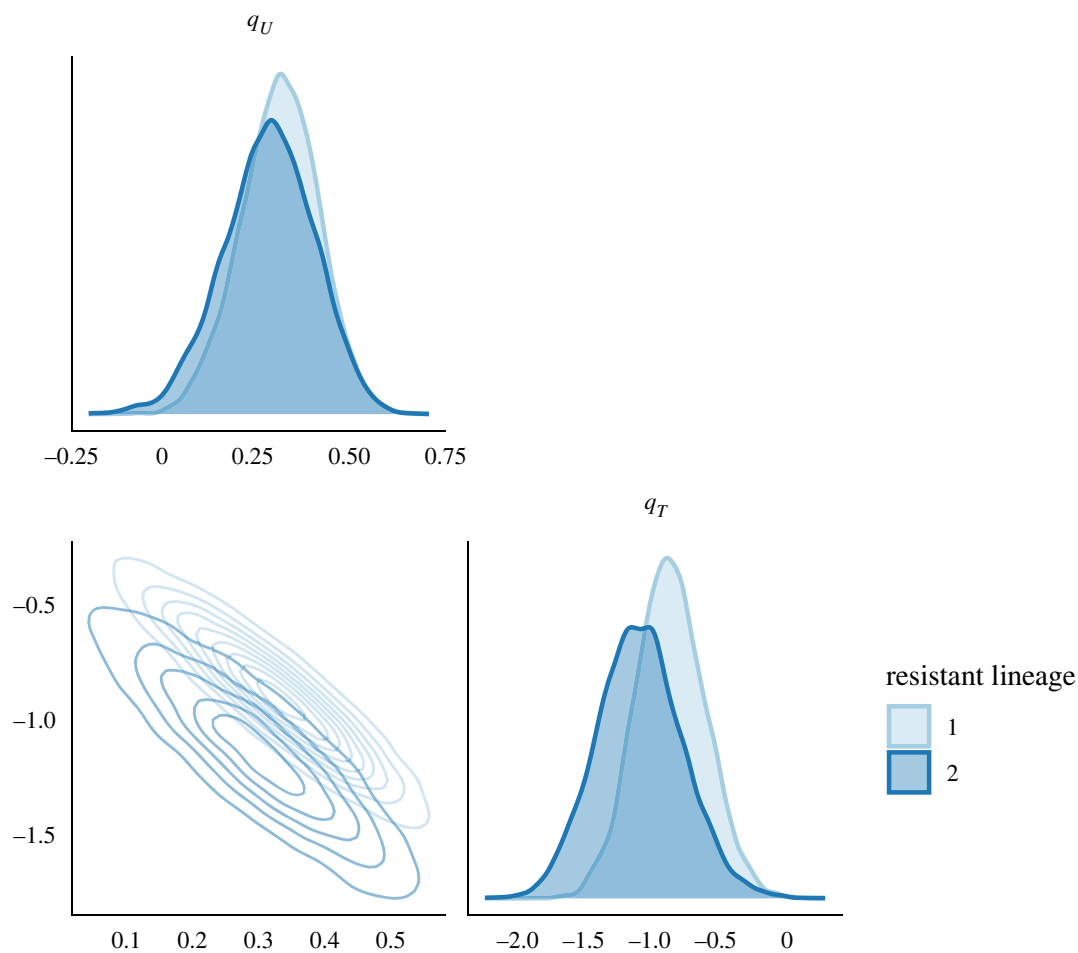


Figure 7. Marginal and joint posterior distribution for the cost (q_U) and benefit (q_T) of both fluoroquinolone resistant lineages of *N. gonorrhoeae*.

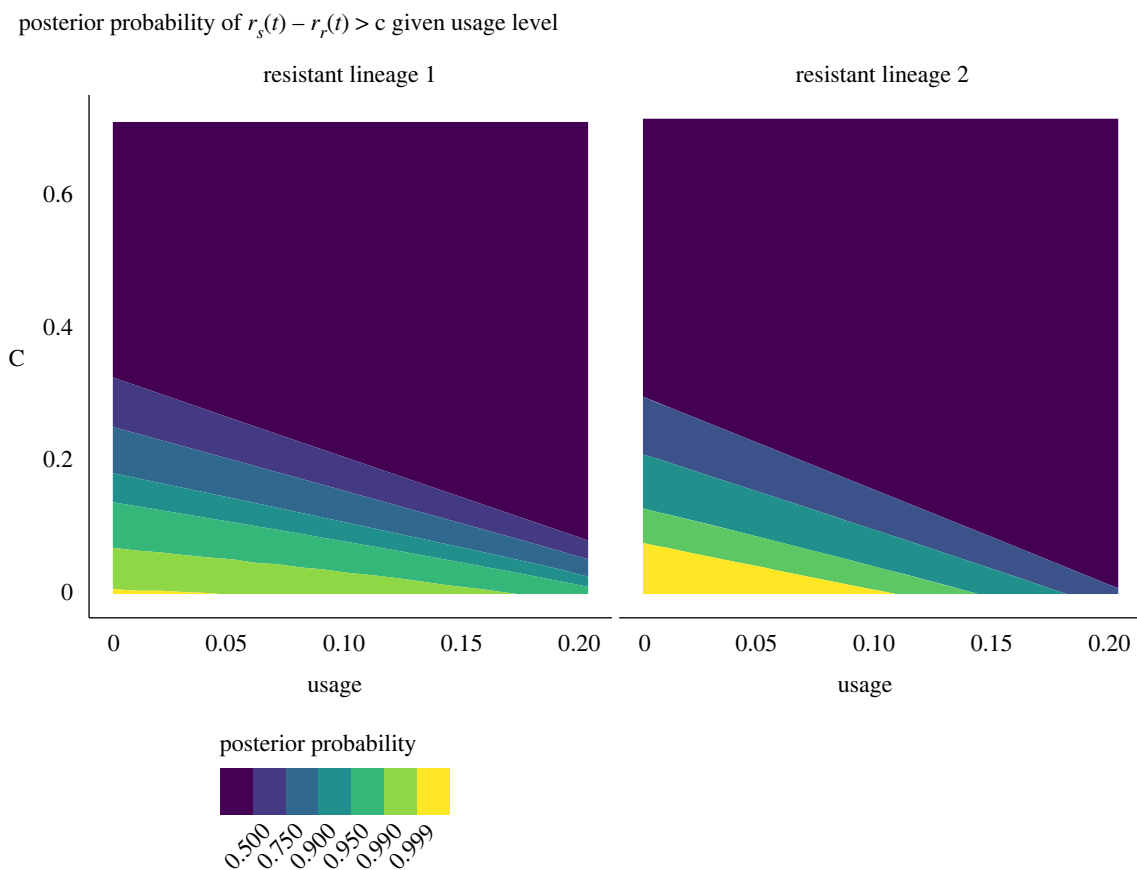


Figure 8. Posterior probabilities of $R_r(t)/R_s(t) < C$ given usage $u(t)$ in the x-axis and value of C in the y-axis, for both fluoroquinolone resistant lineages of *N. gonorrhoeae*.

approximately 20% and 15% of infected individuals, for resistant lineages 1 and 2, respectively.

4. Discussion

A bacterial pathogen lineage that is resistant to a given antibiotic incurs both a fitness cost and a fitness benefit compared to similar susceptible lineages [8]. When the antibiotic is used extensively, the benefit is likely to be greater than the cost. In that case, a resistant lineage has a selective advantage over susceptible lineages, and therefore grows at a faster rate. Conversely, if the antibiotic is used rarely or not at all, the benefit is likely to become smaller than the cost, which will lead to the resistant lineage decreasing in frequency. Estimating these parameters is therefore of primary importance to determine how antibiotics should be prescribed without causing an increase in resistance [9]. Here, we have shown how genome sequencing data coupled with data on antibiotic prescriptions can be used for this purpose, following on previous work that demonstrated the link between epidemic dynamics and phylogenetics [13,19,20,28]. By comparing the phylogenetic trajectories of susceptible and resistant lineages, and relating them with a known function of antibiotic use, we show that it is possible to estimate separately the parameters corresponding to the fitness cost and benefit of resistance. In particular, we reanalysed a large published collection of *N. gonorrhoeae* genomes [38]. We were able to infer these parameters for two lineages of *N. gonorrhoeae* resistant to fluoroquinolones, and found similar estimates of cost and benefit in both (figure 7). We were able to use this knowledge to make recommendations on antibiotic stewardship of fluoroquinolones (figure 8).

Dated phylogenies for both susceptible and resistant lineages are needed as input into our method. Several software tools can be used to produce this either from a sequence alignment, for example BEAST [22] and BEAST2 [23], or from an undated phylogeny, for example treedater [45] and BactDating [24]. Building such a dated phylogeny requires either the population to be measurably evolving over the sampling period [46,47], or a previous estimate of the molecular clock rate [48]. Another input required by our method is the antibiotic usage function over a relevant time-frame and geographical location. This may not always be available in all historical contexts, but efforts are increasingly being made to capture these data [49]. Finally, our method requires an informative prior of the recovery rate for the susceptible lineage (table 1), since this is typically not identifiable from the data, as in many similar compartmental epidemic models [50]. This prior needs to be chosen carefully depending on the infectious disease under study and based on the existing scientific literature.

Our inferential methodology is based on a well-defined and relatively simple epidemic model (equation (2.4)) which means making a number of assumptions the validity of which was considered before performing our analysis. Our model assumes multiple-lineage pathogen dynamics driven by person-to-person transmission in a well-mixed host population in the absence of any significant population structure, so that there is perfect competition between lineages. It also assumes that individuals become infectious as soon as they are infected, that their infectiousness remains constant until they recover, after which they become susceptible again

without any immunity being gained. This list of relatively strong assumptions may seem to preclude application to any real infectious disease, but they are necessary to obtain a model under which inference can be performed. Furthermore, violation of some of these assumptions does not necessarily invalidate the results of inference. For example, if infection causes immunity, this will effectively reduce the number $S(t)$ of susceptible individuals (equation (2.2)), but this number is not assumed to be constant in our model. In fact both the size $N(t)$ of the host population and the number $S(t)$ of susceptible individuals are integrated out as part of our parameterization in terms of the function $b(t)$ (cf. equation (2.4)), so the inference is robust as long as the immunity conferred applies to all lineages under study. Likewise, the assumption of an unstructured population may seem problematic, including in our application to *N. gonorrhoeae* throughout the USA, but for anything other than small local outbreaks the genomes available for analysis are sparsely sampled from the whole infected population [51]. In these conditions, any effect of the host population structure on phylodynamics is likely to be insignificant as long as an effective rather than actual number of infections is considered [52,53].

The compatibility of our model with the phylogenetic data under analysis can be tested using posterior predictive distribution checks (electronic supplementary material, figures S3 and S8). If these tests fail, or if the model assumptions are thought to be inappropriate, a solution may be to resort to other methods that postprocess a dated phylogeny [25] but make fewer assumptions, at the cost of not inferring directly the parameters of resistance. Alternative approaches includes non-parametric methods that detect differences in the branching patterns in different lineages [41,54] as well as methods parameterized in terms of the pathogen population size growth rather than underlying epidemiological drivers [15,55]. However, our model-based approach is both general and flexible, so that we expect it to be applicable in many settings using our software implementation which is available at <https://github.com/dhelekal/ResistPhy/>. We believe that this methodology, applied to the increasingly large genomic databases on many bacterial pathogens, will help quantify the exact link between antibiotic usage and resistance and therefore provide a much-needed evidence basis for the design of future antibiotic prescription strategies [9,56,57].

Data accessibility. All data and code used in this study are available from the GitHub digital repository: <https://github.com/dhelekal/ResistPhy/>.

The data are provided in electronic supplementary material [58].

Authors' contributions. D.H.: conceptualization, investigation, methodology, software, writing—original draft, writing—review and editing; M.K.: conceptualization, funding acquisition, investigation, supervision, writing—review and editing; Y.H.G.: conceptualization, investigation, writing—review and editing; X.D.: conceptualization, funding acquisition, investigation, methodology, supervision, writing—original draft, writing—review and editing.

All authors gave final approval for publication and agreed to be held accountable for the work performed therein.

Conflict of interest declaration. We declare we have no competing interests.

Funding. We acknowledge funding from the National Institute for Health Research (NIHR) Health Protection Research Unit in Genomics and Enabling Data. This work was supported by the UK Engineering and Physical Sciences Research Council (EPSRC) grant no. EP/S022244/1 for the EPSRC Centre for Doctoral Training in Mathematics for Real-World Systems II.

References

1. CDC. 2013 *Antibiotic resistance threats in the United States, 2013*. CS239559-B. Atlanta, GA: CDC. See <https://www.cdc.gov/drugresistance/pdf/ar-threats-2013-508.pdf>.
2. O'Neill J. 2016 *Tackling drug-resistant infections globally: final report and recommendations*. London, UK: Wellcome Trust/HM Government.
3. Murray CJ *et al.* 2022 Global burden of bacterial antimicrobial resistance in 2019: a systematic analysis. *Lancet* **399**, 629–655. (doi:10.1016/S0140-6736(21)02724-0)
4. Clatworthy AE, Pierson E, Hung DT. 2007 Targeting virulence: a new paradigm for antimicrobial therapy. *Nat. Chem. Biol.* **3**, 541–548. (doi:10.1038/nchembio.2007.24)
5. Bonhoeffer S, Lipsitch M, Levin BR. 1997 Evaluating treatment protocols to prevent antibiotic resistance. *Proc. Natl Acad. Sci. USA* **94**, 12 106–12 111. (doi:10.1073/pnas.94.22.12106)
6. Austin DJ, Kristinsson KG, Anderson RM. 1999 The relationship between the volume of antimicrobial consumption in human communities and the frequency of resistance. *Proc. Natl Acad. Sci. USA* **96**, 1152–1156. (doi:10.1073/pnas.96.3.1152)
7. Spicknall IH, Foxman B, Marrs CF, Eisenberg JNS. 2013 A modeling framework for the evolution and spread of antibiotic resistance: literature review and model categorization. *Am. J. Epidemiol.* **178**, 508–520. (doi:10.1093/aje/kwt017)
8. Andersson DI, Levin BR. 1999 The biological cost of antibiotic resistance. *Curr. Opin. Microbiol.* **2**, 489–493. (doi:10.1016/S1369-5274(99)00005-3)
9. Andersson DI, Hughes D. 2010 Antibiotic resistance and its cost: is it possible to reverse resistance? *Nat. Rev. Microbiol.* **8**, 260–271. (doi:10.1038/nrmicro2319)
10. Whittles LK, White PJ, Didelot X. 2017 Estimating the fitness benefit and cost of cefixime resistance in *Neisseria gonorrhoeae* to inform prescription policy: a modelling study. *PLoS Med.* **14**, e1002416. (doi:10.1371/journal.pmed.1002416)
11. Rubin DH, Ma KC, Westervelt KA, Hullahalli K, Waldor MK, Grad YH. 2022 Variation in supplemental carbon dioxide requirements defines lineage-specific antibiotic resistance acquisition in *Neisseria gonorrhoeae*. *bioRxiv*. (doi:10.1101/2022.02.24.481660)
12. Lehtinen S, Blanquart F, Lipsitch M, Fraser C, Maela Pneumococcal Collaboration. 2019 On the evolutionary ecology of multidrug resistance in bacteria. *PLoS Pathog.* **15**, e1007763. (doi:10.1371/journal.ppat.1007763)
13. Pybus OG, Rambaut A. 2009 Evolutionary analysis of the dynamics of viral infectious disease. *Nat. Rev. Genet.* **10**, 540–550. (doi:10.1038/nrg2583)
14. Volz EM, Frost SDW. 2014 Sampling through time and phylodynamic inference with coalescent and birth? Death models. *J. R. Soc. Interface* **11**, 20140945. (doi:10.1098/rsif.2014.0945)
15. Volz EM, Didelot X. 2018 Modeling the growth and decline of pathogen effective population size provides insight into epidemic dynamics and drivers of antimicrobial resistance. *Syst. Biol.* **67**, 719–728. (doi:10.1093/sysbio/syy007)
16. Kühnert D *et al.* 2018 Quantifying the fitness cost of HIV-1 drug resistance mutations through phylodynamics. *PLoS Pathog.* **14**, e1006895. (doi:10.1371/journal.ppat.1006895)
17. Pecerska J, Kühnert D, Meehan CJ, Coscollá M, de Jong BC, Gagneux S, Stadler T. 2021 Quantifying transmission fitness costs of multi-drug resistant tuberculosis. *Epidemics* **36**, 100471. (doi:10.1016/j.epidem.2021.100471)
18. Ho SYW, Shapiro B. 2011 Skyline-plot methods for estimating demographic history from nucleotide sequences. *Mol. Ecol. Resour.* **11**, 423–434. (doi:10.1111/j.1755-0998.2011.02988.x)
19. Volz EM, Kosakovsky Pond SL, Ward MJ, Leigh Brown AJ, Frost SDW. 2009 Phylodynamics of infectious disease epidemics. *Genetics* **183**, 1421–30. (doi:10.1534/genetics.109.106021)
20. Dearlove B, Wilson D. 2013 Coalescent inference for infectious disease: meta-analysis of hepatitis C. *Phil. Trans. R. Soc. B* **368**, 20120314. (doi:10.1098/rstb.2012.0314)
21. Didelot X, Bowden R, Wilson DJ, Peto TEA, Crook DW. 2012 Transforming clinical microbiology with bacterial genome sequencing. *Nat. Rev. Genet.* **13**, 601–612. (doi:10.1038/nrg3226)
22. Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. 2018 Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol.* **4**, vey016. (doi:10.1093/ve/vey016)
23. Bouckaert R *et al.* 2019 BEAST 2.5: an advanced software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.* **15**, e1006650. (doi:10.1371/journal.pcbi.1006650)
24. Didelot X, Croucher NJ, Bentley SD, Harris SR, Wilson DJ. 2018 Bayesian inference of ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Res.* **46**, e134. (doi:10.1093/nar/gky783)
25. Didelot X, Parkhill J. 2022 A scalable analytical approach from bacterial genomes to epidemiology. *Phil. Trans. R. Soc. B* **377**, 20210246. (doi:10.1098/rstb.2021.0246)
26. Allen LJ, Kirupaharan N, Wilson SM. 2004 SIS epidemic models with multiple pathogen strains. *J. Differ. Equ. Appl.* **10**, 53–75. (doi:10.1080/10236190310001603680)
27. Keeling MJ, Rohani P. 2008 *Modeling infectious diseases in humans and animals*. Princeton, NJ: Princeton University Press.
28. Volz EM. 2012 Complex population dynamics and the coalescent under neutrality. *Genetics* **190**, 187–201. (doi:10.1534/genetics.111.134627)
29. Griffiths R, Tavaré S. 1994 Sampling theory for neutral alleles in a varying environment. *Phil. Trans. R. Soc. B* **344**, 403–410. (doi:10.1098/rstb.1994.0079)
30. Gill MS, Lemey P, Faria NR, Rambaut A, Shapiro B, Suchard MA. 2013 Improving Bayesian population dynamics inference: a coalescent-based model for multiple loci. *Mol. Biol. Evol.* **30**, 713–724. (doi:10.1093/molbev/mss265)
31. Griffiths RC. 2003 The frequency spectrum of a mutation, and its age, in a general diffusion model. *Theor. Popul. Biol.* **64**, 241–251. (doi:10.1016/S0040-5809(03)00075-3)
32. Etheridge A, Pfaffelhuber P, Wakolbinger A. 2006 An approximate sampling formula under genetic hitchhiking. *Ann. Appl. Probab.* **16**, 685–729. (doi:10.1214/105051606000000114)
33. Rasmussen CE. 2004 *Gaussian processes in machine learning*, pp. 63–71. Berlin, Germany: Springer. (doi:10.1007/978-3-540-28650-9_4)
34. Riutort-Mayol G, Bürkner P-C, Andersen MR, Solin A, Vehtari A. 2022 Practical Hilbert space approximate Bayesian Gaussian processes for probabilistic programming. (<https://arxiv.org/abs/2004.11408>)
35. Carpenter B *et al.* 2017 Stan: a probabilistic programming language. *J. Stat. Softw.* **76**, 1. (doi:10.18637/jss.v076.i01)
36. Vehtari A, Gelman A, Simpson D, Carpenter B, Burkner PC. 2021 Rank-normalization, folding, and localization: an improved R hat for assessing convergence of MCMC. *Bayesian Anal.* **16**, 667–718. (doi:10.1214/20-BA1221)
37. Gillespie DT. 2001 Approximate accelerated stochastic simulation of chemically reacting systems. *J. Chem. Phys.* **115**, 1716–1733. (doi:10.1063/1.1378322)
38. Grad YH *et al.* 2016 Genomic epidemiology of gonococcal resistance to extended spectrum Cephalosporins, Macrolides, and Fluoroquinolones in the US, 2000–2013. *J. Infect. Dis.* **214**, 1579–1587. (doi:10.1093/infdis/jiw420)
39. Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. 2010 New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.* **59**, 307–321. (doi:10.1093/sysbio/syq010)
40. Didelot X, Wilson DJ. 2015 ClonalFrameML: efficient inference of recombination in whole bacterial genomes. *PLoS Comput. Biol.* **11**, e1004041. (doi:10.1371/journal.pcbi.1004041)
41. Volz EM, Wiuf C, Grad YH, Frost SDW, Dennis AM, Didelot X. 2020 Identification of hidden population structure in time-scaled phylogenies. *Syst. Biol.* **69**, 884–896. (doi:10.1093/sysbio/syaa009)
42. Gelman A, Meng X, Stern H. 1996 Posterior predictive assessment of model fitness via realized discrepancies. *Stat. Sinica* **6**, 733–807.
43. Fingerhuth SM, Bonhoeffer S, Low N, Althaus CL. 2016 Antibiotic-resistant *Neisseria gonorrhoeae* spread faster with more treatment, not more sexual partners. *PLoS Pathog.* **12**, e1005611. (doi:10.1371/journal.ppat.1005611)

44. Whittles LK, White PJ, Didelot X. 2019 A dynamic power-law sexual network model of gonorrhoea outbreaks. *PLoS Comput. Biol.* **15**, e1006748. (doi:10.1371/journal.pcbi.1006748)
45. Volz EM, Frost SDW. 2017 Scalable relaxed clock phylogenetic dating. *Virus Evol.* **3**, vex025. (doi:10.1093/ve/vex025)
46. Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG. 2003 Measurably evolving populations. *Trends Ecol. Evol.* **18**, 481–488. (doi:10.1016/S0169-5347(03)00216-7)
47. Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. 2015 Measurably evolving pathogens in the genomic era. *Trends Ecol. Evol.* **30**, 306–313. (doi:10.1016/j.tree.2015.03.009)
48. Duchêne S, Holt KE, Weill F-X, Le Hello S, Hawkey J, Edwards DJ, Fourment M, Holmes EC. 2016 Genome-scale rates of evolutionary change in bacteria. *Microbial Genom.* **2**, e000094. (doi:10.1101/069492)
49. Curtis HJ, Goldacre B. 2018 Openprescribing: normalised data and software tool to research trends in English NHS primary care prescribing 1998–2016. *BMJ open* **8**, e019921. (doi:10.1136/bmjopen-2017-019921)
50. Roosa K, Chowell G. 2019 Assessing parameter identifiability in compartmental dynamic models using a computational approach: application to infectious disease transmission models. *Theor. Biol. Med. Modell.* **16**, 1. (doi:10.1186/s12976-018-0097-6)
51. Klinkenberg D, Colijn C, Didelot X. 2019 Methods for outbreaks using genomic data. In *Handbook of infectious disease data analysis* (eds L Held, N Hens, P O'Neill, J Wallinga), pp. 245–263. Boca Raton, FL: CRC Press.
52. Nordborg M. 1997 Structured coalescent processes on different time scales. *Genetics* **146**, 1501–1514. (doi:10.1093/genetics/146.4.1501)
53. Frost SDW, Volz EM. 2010 Viral phylodynamics and the search for an 'effective number of infections'. *Phil. Trans. R. Soc. B* **365**, 1879–1890. (doi:10.1098/rstb.2010.0060)
54. Dearlove BL, Xiang F, Frost SDW. 2017 Biased phylodynamic inferences from analysing clusters of viral sequences. *Virus Evol.* **3**, 1–10. (doi:10.1093/ve/vex020)
55. Helekal D, Ledda A, Volz E, Wyllie D, Didelot X. 2021 Bayesian inference of clonal expansions in a dated phylogeny. *Syst. Biol.* **71**, syab095. (doi:10.1093/sysbio/syab095)
56. Michael CA, Dominey-Howes D, Labbate M. 2014 The antimicrobial resistance crisis: causes, consequences, and management. *Front. Public Health* **2**, 145. (doi:10.3389/fpubh.2014.00145)
57. Holmes AH, Moore LS, Sundsfjord A, Steinbakk M, Regmi S, Karkey A, Guerin PJ, Piddock LJ. 2016 Understanding the mechanisms and drivers of antimicrobial resistance. *Lancet* **387**, 176–187. (doi:10.1016/S0140-6736(15)00473-0)
58. Helekal D, Keeling M, Grad YH, Didelot X. 2023 Estimating the fitness cost and benefit of antimicrobial resistance from pathogen genomic data. Figshare. (doi:10.6084/m9.figshare.c.6673522)