

**PHYLOGENETIC ESTIMATION OF CONTACT NETWORK PARAMETERS
WITH APPROXIMATE BAYESIAN COMPUTATION**

by

Rosemary Martha McCloskey

B.Sc., Simon Fraser University, 2014

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE
DEGREE OF

MASTER OF SCIENCE

in

The Faculty of Graduate and Postdoctoral Studies

(Bioinformatics Graduate Program)

THE UNIVERSITY OF BRITISH COLUMBIA
(Vancouver)

August 2016

Abstract

Models of the spread of disease in a population often make the simplifying assumption that the population is homogeneously mixed, or is divided into homogeneously mixed compartments. However, human populations have complex structures formed by social contacts, which can have a significant influence on the rate and pattern of epidemic spread. Contact ~~network-models~~ networks capture this structure by explicitly representing each contact that could possibly lead to a transmission. Contact network models parameterize the structure of these networks, but estimating their parameters from contact data requires extensive, often prohibitive, epidemiological investigation.

We developed a method based on approximate Bayesian computation (ABC) for estimating structural parameters of the contact network underlying an observed viral phylogeny. The method combines adaptive sequential Monte Carlo for ABC, Gillespie simulation for propagating epidemics through networks, and a previously developed kernel-based tree similarity score. Our method offers the potential to quantitatively investigate contact network structure from phylogenies derived from viral sequence data, complementing traditional epidemiological methods.

We applied our method to ~~fit~~ the Barabási-Albert network model. This model incorporates the preferential attachment mechanism observed in real world social and sexual networks, whereby individuals with more connections attract new contacts at an elevated rate (“the rich get richer”). ~~to-simulated transmission-trees-and-applied-it-to-viral-phylogenies-estimated-from-six-real-world-HIV-sequence-datasets.~~ Using simulated data, we found that the strength of preferential attachment and the number of infected nodes could often be accurately estimated. However, the mean degree of the network and the total number of nodes appeared to be weakly- or non-identifiable with ABC.

Finally, the Barabási-Albert model was fit to ~~six~~ eleven real world HIV datasets, and substantial heterogeneity in the parameter estimates was observed. ~~Point-estimates~~ Posterior means for the preferential attachment power were all sub-linear and ranged from ~~0.06-to-1.05~~ 0.27 to 0.9. Point estimates of the strength of preferential attachment were higher in injection drug user populations, potentially indicating that high-degree “hub” nodes may play a role in epidemics among this risk group. Our results underscore the importance of considering contact structures when ~~performing-phylogenetic-inference~~ investigating viral outbreaks.

Preface

The initial idea to use approximate Bayesian computation (ABC) to infer contact network model parameters was Dr. Poon's, based on his previous work using ABC to infer parameters of population genetic models. The tree kernel was originally developed by Dr. Poon, but the version used here was implemented by me to improve computational efficiency. The idea to apply sequential Monte Carlo was mine, but Dr. Alexandre Bouchard-Côté made me aware of the adaptive version used in this work. Dr. Sarah Otto suggested the experiments involving a network with a heterogeneous α parameter and peer-driven sampling. Dr. Richard Liang provided guidance in the development of the Gillespie simulation algorithm and statistical advice. The *netabc* program, and all supplementary analysis programs, were written by me.

A version of chapter 2 has been submitted for publication with the title “Reconstructing network parameters from viral phylogenies.” An oral presentation entitled “Phylogenetic inference of contact network parameters with kernel-ABC” was given based on chapter 2 to the 23rd HIV Dynamics and Evolution meeting on April 25, 2016, in Woods Hole, Massachusetts, USA (the presentation was delivered remotely). A poster based on chapter 2 entitled “Likelihood-free estimation of contact network parameters from viral phylogenies” is scheduled for presentation at the Intelligent Systems for Molecular Biology meeting on July 8, 2016, in Orlando, Florida, USA.

Use of the BC data is in accordance with an ethics application that was reviewed and approved by the UBC/Providence Health Care Research Ethics Board (H07-02559). Rosemary M. McCloskey completed the Tri-Council Policy Statement: Ethical Conduct for Research Involving Humans Course on Research Ethics (TCPS 2: CORE) tutorial on June 8, 2016.

[Source code for the *netabc* program is freely available at <https://github.com/rmcclosk/netabc> under the GPL-3 license. Scripts to run all computational experiments, as well as the source code for this thesis, are available at <https://github.com/rmcclosk/thesis>.](https://github.com/rmcclosk/netabc)

Table of Contents

Abstract	ii
Preface	iii
Table of Contents	v
List of Tables	vi
List of Figures	vii
List of Symbols	viii
List of Abbreviations	ix
Acknowledgements	x
1 Introduction	1
1.1 Objective	1
1.2 Phylogenetics and phylodynamics	3
1.2.1 Phylogenetic trees	3
1.2.2 Transmission trees	5
1.2.3 Phylodynamics: linking evolution and epidemiology	6
1.2.4 Tree shapes	7
1.3 Contact networks	9
1.3.1 Overview	9
1.3.2 Scale-free networks and preferential attachment	11
1.3.3 Relationship between network structure and transmission trees	13
1.4 Sequential Monte Carlo	14
1.4.1 Overview and notation	14
1.4.2 Sequential importance sampling	16
1.4.3 Sequential Monte Carlo	17
1.4.4 The sequential Monte Carlo sampler	19
1.5 Approximate Bayesian computation	21
1.5.1 Overview and motivation	21

1.5.2	Algorithms for ABC	24
1.6	Summary	26
2	Reconstructing contact network parameters from viral phylogenies	28
2.1	<i>Netabc</i> : a computer program for estimation of contact network parameters with ABC .	28
2.1.1	Simulation of transmission trees over contact networks	30
2.1.2	Phylogenetic kernel	32
2.1.3	Adaptive sequential Monte Carlo for Approximate Bayesian computation . . .	32
2.1.4	Justification for approach	33
2.2	Analysis of Barabási-Albert model with synthetic data	36
2.2.1	Motivation	36

List of Tables

List of Figures

1.1	Illustration of a rooted, ultrametric, time-scaled phylogeny.	4
1.2	Illustration of a contact network and transmission tree.	6
1.3	Examples of Barabási-Albert networks with preferential attachment power $\alpha = 0, 1$, and 2.	13
2.1	Graphical schematic of the ABC-SMC algorithm implemented in <i>netabc</i>	29
2.2	Illustration of an estimated transmission tree without labels and two possible underlying complete transmission trees with labels.	34

List of Symbols

- I number of infected nodes in a contact network at the time of transmission tree sampling.
- M number of simulated datasets per particle in ABC-SMC.
- N total number of nodes in a contact network.
- T a transmission tree.
- Θ parameter space for a given mathematical model.
- α_{ESS} per-iteration rate of expected sample size decay in adaptive ABC-SMC.
- α preferential attachment power parameter in Barabási-Albert networks.
- β transmission rate in susceptible-infected and susceptible-infected-removed epidemiological models.
- γ exponent of power-law degree distribution in scale-free networks.
- \mathcal{D} set of discordant edges in Gillespie simulation.
- \mathcal{I} set of infected nodes in Gillespie simulation.
- \mathcal{R} set of recovered nodes in Gillespie simulation.
- \mathcal{S} set of recovered nodes in Gillespie simulation.
- ν recovery rate in the susceptible-infected-removed epidemiological model.
- ρ distance function for approximate Bayesian computation.
- ε distance function for approximate Bayesian computation.
- m number of edges added per vertex when constructing a Barabási-Albert network.
- n number of particles used for sequential Monte Carlo.
- q proposal function for Metropolis-Hastings kernel.
- t_{max} user-defined cutoff time at which to stop Gillespie simulation.
- t time since initial infection of index case in an epidemic.

List of Abbreviations

ABC approximate Bayesian computation.

BA Barabási-Albert.

ER Erdős-Rényi.

ESS expected sample size.

GSL GNU scientific library.

HCV hepatitis C virus.

HIV human immunodeficiency virus.

HMM hidden Markov model.

IDU injection drug users.

IS importance sampling.

LTT lineages-through-time.

MCMC Markov chain Monte Carlo.

MH Metropolis-Hastings.

ML maximum likelihood.

MSM men who have sex with men.

nLTT normalized lineages-through-time.

PA preferential attachment.

POSIX Portable Operating System Interface.

SARS severe acute respiratory syndrome.

SI susceptible-infected.

SIR susceptible-infected-recovered.

SIS sequential importance sampling.

SMC sequential Monte Carlo.

SVM support vector machine.

TasP treatment as prevention.

WS Watts-Strogatz.

Acknowledgements

Chapter 1

Introduction

1.1 Objective

The spread of a disease is most often modelled by assuming either a homogeneously mixed population [1, 2], or a population divided into a small number of homogeneously mixed groups [3]. This assumption, also called *mass action* [4], or *panmixia*, implies that any two individuals in the same compartment are equally likely to come into contact making transmission possible at some predefined rate. Although this provides a reasonable approximation in many cases [5], the error introduced by assuming a panmictic population can be substantial when significant contact heterogeneity exists in the underlying population [6–8]. Contact network models provide an alternative to compartmental models which do not require the assumption of panmixia. In addition to more accurate predictions, the parameters of the networks themselves may be of interest from a public health perspective. For example, certain vaccination strategies may be more or less effective in curtailing an epidemic depending on the underlying network’s degree distribution [9, 10]. Phylodynamic methods, [which link viruses’ evolutionary and epidemiological dynamics](#), have been used to fit many different types of models to phylogenetic data [11, 12]. However, these models generally assume a panmictic population. The primary objective of this work is *to develop a method to fit contact network models, and thereby relax the assumption of homogeneous mixing, in a phylodynamic framework.*

[In this work, we take a Bayesian approach: our goal is to estimate the posterior distribution on model parameters given our data,](#)

$$\pi(\theta \mid T) = \frac{f(T \mid \theta)\pi(\theta)}{\int_{\Theta} f(T \mid \theta)\pi(\theta)d\theta},$$

[where \$f\(T \mid \theta\)\$ is the likelihood of the parameters given \$T\$, \$\pi\(\theta\)\$ is the prior on \$\theta\$, and \$\Theta\$ is the space of possible model parameters. The denominator on the right-hand side is the marginal probability of \$T\$ which acts as a normalizing constant on the posterior \(see ?? for a review of mathematical modeling and Bayesian inference, including definitions of these concepts\). As we shall show \(section 2.1.4\), estimating this distribution presents computational challenges beyond those usually encountered in Bayesian inference. Both the likelihood \$f\(T \mid \theta\)\$ and the normalizing constant seem to be intractable, which rules out the use of most common maximum likelihood and Bayesian methods.](#)

Calculating the likelihood of the parameters of a contact network model seems likely to be an intractable problem. We have not proven this is the case, but some intuition can be provided by examining the process involved in the likelihood calculation. Consider a contact network model with parameters θ and an estimated transmission tree T with n tips. In general, we do not know the labels of the internal nodes of T , only the labels of its tips. To fit this model using likelihood-based methods, we must calculate the likelihood of θ , that is, $\Pr(T|\theta)$. Let \mathcal{G} be the set of all possible contact networks, and \mathcal{N} be the set of all possible labellings of the internal nodes of T . We can write the likelihood as

$$\begin{aligned}\Pr(T|\theta) &= \sum_{v \in \mathcal{N}} \Pr(T, v|\theta) \\ &= \sum_{G \in \mathcal{G}} \sum_{v \in \mathcal{N}} \Pr(T, v|G, \theta) \Pr(G|\theta) \\ &= \sum_{G \in \mathcal{G}} \sum_{v \in \mathcal{N}} \Pr(T, v|G) \Pr(G|\theta),\end{aligned}\tag{1.1}$$

the last equality following from the fact that T and v depend only on G , not on θ . Although $\Pr(T, v|G)$ and $\Pr(G|\theta)$ may individually be straightforward to calculate, the number of possible directed graphs on N nodes is $2^{N(N-1)}$ [13], larger if the nodes and edges in the graph may have different labels or attributes. Hence, the number of terms in the sum is at least exponential in n , as there must be at least n nodes in the network. In addition, eq. (1.1) assumes that T is complete, meaning that all infected individuals were sampled. This is rarely the case in practice – most often, we only have access to a subset of the infected individuals. In this case, the likelihood calculation becomes even more complex, because we must also sum over all possible complete trees.

Depending on the network model studied, it is possible that eq. (1.1) could be simplified into a tractable expression. An alternative to likelihood-based methods, which could be applied to any network model, is provided by approximate Bayesian computation (ABC) [14–17]. All of the ingredients required to apply ABC to this problem are readily available. Simulating networks is straightforward under a variety of models. Epidemics on those networks, and the corresponding transmission trees, can also be easily simulated. As mentioned above, contact networks can profoundly affect transmission tree shape. Those shapes can be compared using a highly informative similarity measure called the “tree kernel” [18]; [similar kernel functions have been demonstrated to work well as distance functions in ABC \[19\]](#). ABC can be implemented with several algorithms, but sequential Monte Carlo (SMC) has advantages over others, [including improved accuracy in low-density regions and parallelizability](#) [20]. A recently-developed adaptive algorithm requiring minimal tuning on the part of the user makes SMC an even more attractive approach [21]. In summary, our method to infer contact network parameters will combine the following: stochastic simulation of epidemics on networks, the tree kernel, and adaptive ABC-SMC. [Our method will expand on the framework developed by \[22\], who combined ABC with the tree kernel to infer parameters of population genetic models from viral phylogenies using Markov chain Monte Carlo \(MCMC\).](#)

Empirical studies of sexual contact networks have found that these networks tend to be scale-free [23–26], meaning that their degree distributions follow a power law (although there has been some disagreement, see [6, 27]). Preferential attachment has been postulated as a mechanism by which scale-free networks could be generated [28]. The Barabási-Albert (BA) model [28] is one of the simplest

preferential attachment models, which makes it a natural choice to explore with our method. The second aim of this work is *to use simulations to investigate the parameters of the Barabási-Albert model, including whether they have a detectable impact on tree shape, and whether they can be accurately recovered using ABC.*

Due to its high global prevalence and fast mutation rate, human immunodeficiency virus (HIV) is one of the most commonly-studied viruses in a phylodynamic context. Consequently, a large volume of HIV sequence data is publicly available, more than for any other pathogen, and including sequences sampled from diverse geographic and demographic settings. At the time of this writing, there were 635,400 HIV sequences publicly available in GenBank, annotated with 172 distinct countries of origin. Since HIV is almost always spread through either sexual contact or sharing of injection drug supplies, the contact networks underlying HIV epidemics are driven by social dynamics and are therefore likely to be highly structured [26]. Moreover, since no cure yet exists, efforts to curtail the progression of an epidemic have relied on preventing further transmissions through measures such as treatment as prevention (TasP) and education leading to behaviour change. The effectiveness of this type of intervention can vary significantly based on the underlying structure of the network and the particular nodes to whom the intervention is targeted [29, 30]. Due to this combination of data availability and potential public health impact, HIV is an obvious context in which our method could be applied. Therefore, the third and final aim of this work is *to apply ABC to fit the Barabási-Albert model to existing HIV outbreaks.*

To summarize, this work has three objectives. First, we will develop a method which uses ABC to infer parameters of contact network models from observed transmission trees. Second, we will use simulations to characterize the parameters of the BA network model in terms of their effect on tree shape and how accurately they can be recovered with ABC. Finally, we will apply the method to fit the BA model to several real-world HIV datasets.

[The remainder of this background chapter is organized in four sections. The first section introduces phylogenies and transmission trees, which are the input data from which our method aims to make statistical inferences. This section also introduces phylodynamics, a family of methods that, like ours, aim to infer epidemiological parameters from evolutionary data. The second section focuses on contact networks and network models, whose parameters we are attempting to infer. The relationship between contact networks and transmission trees is also discussed. The third and fourth sections introduce SMC and ABC respectively, which are the two algorithmic components of the method we will implement. In particular, ABC refers to the general approach of using simulations to replace likelihood calculations in a Bayesian setting, while SMC is a particular algorithm which can be used to implement ABC.](#)

1.2 Phylogenetics and phylodynamics

1.2.1 Phylogenetic trees

In evolutionary biology, a *phylogeny*, or *phylogenetic tree*, is a graphical representation of the evolutionary relationships among a group of organisms or species (generally, *taxa*) [31]. The *tips* of a phylogeny, that is, the nodes without any descendants, correspond to *extant*, or observed, taxa. The *internal nodes*

correspond to their common ancestors, [usually extinct \(although occasionally the internal nodes may be observed as well, eg. \[32\]\)](#). The edges or *branches* of the phylogeny connect ancestors to their descendants. Phylogenies may have a *root*, which is a node with no descendants distinguished as the most recent common ancestor of all the extant taxa [33]. When such a root exists, the tree is referred to as being *rooted*; otherwise, it is *unrooted*. The structural arrangement of nodes and edges in the tree is referred to as its *topology* [34].

The branches of the tree may have associated lengths, representing either evolutionary distance or calendar time between ancestors and their descendants. The term “evolutionary distance” is used here imprecisely to mean any sort of quantitative measure of evolution, such as the number of differences between the DNA sequences of an ancestor and its descendant, or the difference in average body mass or height. A phylogeny with branch lengths in calendar time units is often referred to as *time-scaled*. In a time-scaled phylogeny, the internal nodes can be mapped onto a timeline by using the tips of the tree, which usually correspond to the present day, reference points [35]. The corresponding points on the timeline are called *branching times*, and the rate of their accumulation is referred to as the *branching rate*. Rooted trees whose tips are all the same distance from the root are called *ultrametric* trees [36]. These concepts are illustrated in fig. 1.1.

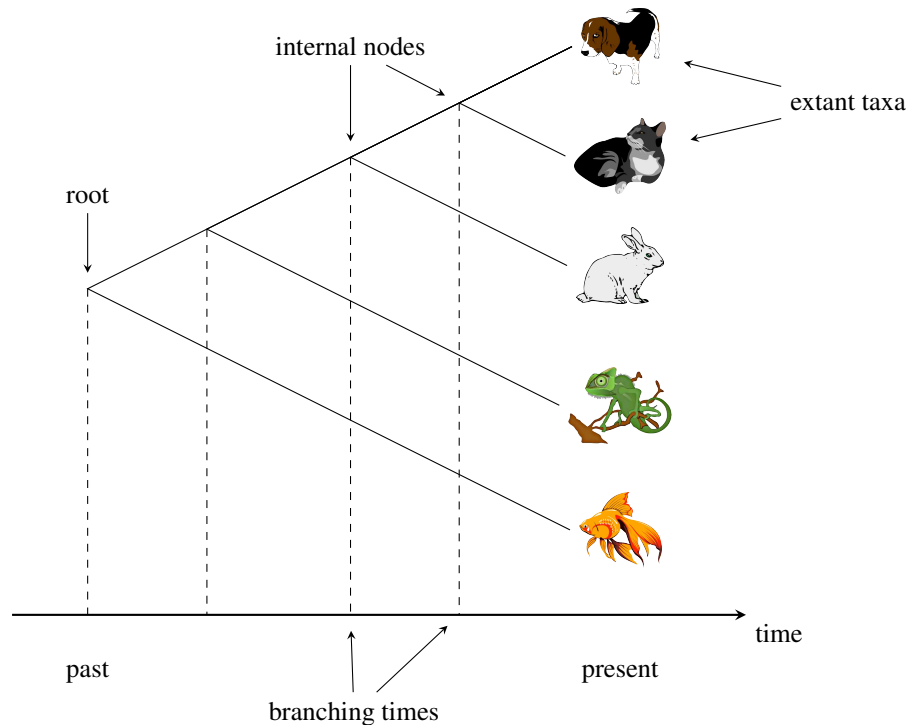


Figure 1.1: Illustration of a rooted, ultrametric, time-scaled phylogeny. The tips of the tree, which represent extant taxa, are placed at the present day on the time axis. Internal nodes, representing extinct common ancestors to the extant taxa, fall in the past. The topology of the tree indicates that cats and dogs are the most closely related pair of species, whereas fish is most distantly related to any other taxon in the tree.

1.2.2 Transmission trees

In epidemiology, a *transmission tree* is a graphical representation of an epidemic’s progress through a population [37]. Like phylogenies, transmission trees have tips, nodes, edges, and branch lengths. However, rather than recording an evolutionary process (speciation), they record an epidemiological process (transmission). The tips of a transmission tree represent the removal by sampling of infected hosts, while internal nodes correspond to transmissions from one host to another. Transmission trees generally have branch lengths in units of calendar time, with branching times indicating times of transmission. The root of a transmission tree corresponds to the initially infected patient who introduced the epidemic into the network, also known as the *index case*. The internal nodes may be labelled with the donor of the transmission pair, if this is known. The tips of the tree, rather than being fixed at the present day, are placed at the time at which the individual was removed from the epidemic, such as by death, recovery, isolation, behaviour change, or migration [38]. Consequently, the transmission tree may not be ultrametric, but may have tips located at varying distances from the root. Such trees are said to have *heterochronous* taxa [39], in contrast to the *isochronous* taxa found in most phylogenies of macro-organisms. A transmission tree is illustrated in fig. 1.2 (right). The object on the right of the figure is called a *contact network*, which depicts the entire susceptible population along with all possible routes of disease transmission. Contact networks, and their relationships to transmission trees, will be discussed further in section 1.3.

Each infected individual in an epidemic may appear at nodes of the transmission tree more than once. This is different from the transmission *network*, in which each infected individual appears exactly once, and edges are in one-to-one correspondence with transmissions [8, 40]. The distinction between the two objects is illustrated in fig. 1.2. However, since transmission networks generally have no cycles (unless re-infection occurs), they are trees in the graph theoretical sense, and hence are sometimes also referred to as transmission trees [*e.g.* 41]. In this work, we reserve the term “transmission tree” for the objects depicted on the right side of fig. 1.2, following *e.g.* [38]. The term “transmission network” is taken to mean the subgraph of the contact network along which transmissions occurred, following *e.g.* [8, 40].

Since transmission trees are essentially a detailed record of an epidemic’s progress, they contain substantial epidemiological information. As a basic example, the lineages-through-time (LTT) plot [35], which plots the number of lineages in a phylogeny against time, can be used to quantify the incidence of new infections over the course of an epidemic [42]. However, in all but the most well-studied of epidemics, transmission trees are not possible to assemble through traditional epidemiological methods [40]. The time and effort to conduct detailed interviews and contact tracing of a sufficient number of infected individuals is usually prohibitive, and may additionally be confounded by misreporting and other challenges [43]. However, it turns out that for viral epidemics, some of the epidemiological information contained in the transmission tree leaves a mark on the viral genetic material circulating in the population. A family of methods called *phylodynamics* [44] addresses the challenge of estimating epidemiological parameters from viral sequence data [12].

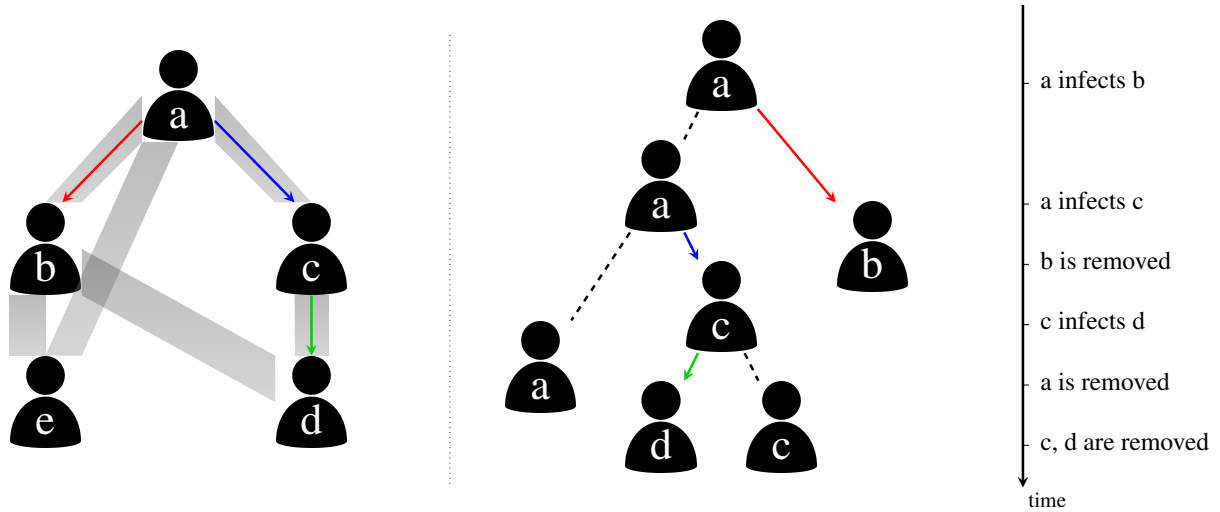


Figure 1.2: Illustration of epidemic spread over a contact network, and the corresponding transmission tree. (Left) A contact network with five hosts, labelled *a* through *e*. Thick shaded edges indicate symmetric contacts among the hosts. The transmission network is indicated by coloured arrows. The epidemic began with node *a*, who transmitted to nodes *b* and *c*. Node *c* further transmitted to node *d*. Node *e* was not infected. (Right) The transmission tree corresponding to this scenario, with a timeline of transmission and removal times.

1.2.3 Phylodynamics: linking evolution and epidemiology

The basis of phylodynamics is the fact that, for RNA viruses, epidemiological and evolutionary processes occur on similar time scales [39]. In fact, these two processes interact, such that it is possible to detect the influence of host epidemiology on the evolutionary history of the virus as recorded in an *inter-host viral phylogeny*. Phylodynamic methods aim to detect and quantify the signatures of epidemiological processes in these phylogenies [11, 12], which relate one representative viral genotype from each host in an infected population. These methods have been used to investigate parameters such as transmission rate, recovery rate, and basic reproductive number [11, 12]. The majority of phylodynamic studies attempt to infer the parameters of an epidemiological model for which the likelihood of an observed phylogeny can be calculated. Most often, this is some variation of the birth-death [45, 46] or coalescent [47, 48] models. These methods either assume the viral phylogeny is known, as we do in this work, or (more commonly) integrate over phylogenetic uncertainty in a Bayesian framework. Phylogenetic inference is a complex topic which we shall not discuss here; see *e.g.* [49] for a full review.

Due to the relationship between the aforementioned processes, there is a degree of correspondence between viral phylogenies and transmission trees [37, 41, 50, 51]. In particular, the transmission process is quite similar to *allopatric speciation* [52], where genetic divergence follows the geographic isolation of a sub-population of organisms. Thus, transmission, which is represented as branching in the transmission tree, causes branching in the viral phylogeny as well [53]. Similarly, the removal of an individual from the transmission tree causes the extinction of their viral lineage in the phylogeny. Consequently, the topology of the viral phylogeny is sometimes used as a proxy for the topology of the transmission tree [54]. Modern likelihood-based methods of phylogenetic reconstruction [*e.g.* 55, 56] produce

unrooted trees whose branch lengths measure genetic distance in units of expected substitutions per site. On the other hand, transmission trees are rooted, and have branches measuring calendar time [11]. Therefore, estimating a transmission tree from a viral phylogeny requires the phylogeny to be rooted and time-scaled. Methods for performing this process include root-to-tip regression [57–59], which we apply in this work, and least-square dating [60]. Alternatively, the tree may be rooted separately with an outgroup [61] before time-scaling.

A caveat of estimating transmission trees in this manner is that the correspondence between the topologies of the viral phylogeny and transmission tree is far from exact [37, 62]. Due to intra-host diversity, the viral strain which is transmitted may have split from another lineage within the donor long before the transmission event occurred. Hence, the branching point in the viral phylogeny may be much earlier than that in the transmission tree. Another possibility is that one host transmitted to two or more recipients in one order, but the transmitted lineages originated within the donor in a different order. In this case, the topology of the transmission tree and the viral phylogeny will be mismatched. In practice, this discordance has not proven an insurmountable problem: for example, Leitner et al. [63] and Paraskevis et al. [64] were able to accurately recover known transmission trees using viral phylogenies. The problem of accurately estimating transmission trees is an ongoing area of research [32, 54, 65–68]. For example, Hall, Woolhouse, and Rambaut [54] developed a Bayesian method to jointly estimate a transmission tree and viral phylogeny by combining models of agent-based transmission, within-host population dynamics, and sequence evolution.

1.2.4 Tree shapes

To perform phylodynamic inference, we must be able to extract quantitative information from viral phylogenies. What is informative about a phylogeny, beyond the demographic characteristics of the individuals it relates, is its *shape*. The shape of a phylogeny has two components: the topology, and the distribution of branch lengths [69]. Methods of quantifying tree shape fall into two categories: summary statistics, and pairwise measures. Summary statistics assign a numeric value to each individual tree, while pairwise measures quantify the similarity between pairs of trees.

One of the most widely used tree summary statistics is Sackin’s index [70], which measures the imbalance or asymmetry in a rooted tree. For the i th tip of the tree, we define N_i to be the number of branches between that tip and the root. The unnormalized Sackin’s index is defined as the sum of all N_i . It is called unnormalized because it does not account for the number of tips in the tree. Among two trees having the same number of tips, the least-balanced tree will have the highest Sackin’s index. However, among two equally balanced trees, the larger tree will have a higher Sackin’s index. This makes it challenging to compare balances among trees of different sizes. To correct this, Kirkpatrick and Slatkin [71] derive the expected value of Sackin’s index under the Yule model [72]. Dividing by this expected value normalizes Sackin’s index, so that it can be used to compare trees of different sizes. An example of a pairwise measure is the normalized lineages-through-time (nLTT) [73], which compares the LTT [35] plots of two trees. Specifically, the two LTT plots are normalized so that they begin at (0,0) and end at (1,1), and the absolute difference between the two plots is integrated between 0 and 1. In the context of

infectious diseases, the LTT is related to the prevalence [42], so large values may indicate that the trees being compared were produced by different epidemic trajectories [73].

Poon et al. [18] developed an alternative pairwise measure which applies the concept of a *kernel function* to phylogenies. Kernel functions, originally developed for support vector machines (SVMs) [74], compare objects in a space \mathcal{X} by mapping them into a feature space \mathcal{F} of high or infinite dimension via a function ϕ . The similarity between the objects is defined as

$$K(x, x') = \langle \phi(x), \phi(x') \rangle,$$

that is, the inner product of the objects' representations in the feature space. Computing $\phi(x)$ may be computationally prohibitive due to the dimension of \mathcal{F} . The utility of a kernel function K is that it is constructed in such a way that it can compute the inner product without explicitly computing $\phi(x)$. The kernel function developed in [18] will henceforth be referred to as the *tree kernel*. This kernel maps trees into the space of all possible possible *subset trees*, which are subtrees that do not necessarily extend all the way to the tips. The subset-tree kernel was originally developed for comparing parse trees in natural language processing [75] and did not incorporate branch length information. The version developed by Poon et al. [18] includes a radial basis function to compare the differences in branch lengths, thus incorporating both the trees' topologies and their branch lengths in a single similarity score.

The kernel score of a pair of trees, denoted $K(T_1, T_2)$, is defined as a sum over all pairs of nodes (n_1, n_2) , where n_1 is a node in T_1 and n_2 is a node in T_2 . Following Poon et al. [18], let $N(T)$ denote the set of all nodes in T , $\text{nc}(n)$ be the number of children of node n , c_n^j be the j th child of node n , and l_n be the vector of branch lengths connecting node n to its $\text{nc}(n)$ children. Furthermore, let $\text{nl}(n)$ be the number of children of n which are leaves (we always have $\text{nl}(n) \leq \text{nc}(n)$). The production rule of n is the pair $(\text{nc}(n), \text{nl}(n))$. That is, if two nodes have the same number of children and among these, the same number of leaves, then they have the same production rule. Let $k_G(x, y)$ be a Gaussian radial basis function of the vectors x and y ,

$$k_G(x, y) = \exp\left(-\frac{1}{2\sigma} \|x - y\|_2^2\right),$$

where $\|\cdot\|_2$ is the Euclidean norm and σ is a variance parameter. The tree kernel is defined as

$$K(T_1, T_2) = \sum_{n_1 \in N(T_1)} \sum_{n_2 \in N(T_2)} \Delta(n_1, n_2), \quad (1.2)$$

where

$$\Delta(n_1, n_2) = \begin{cases} \lambda & n_1 \text{ and } n_2 \text{ are leaves} \\ \lambda k_G(l_{n_1}, l_{n_2}) \prod_{j=1}^{\text{nc}(n_1)} (1 + \Delta(c_{n_1}^j, c_{n_2}^j)) & n_1 \text{ and } n_2 \text{ have the same production rule} \\ 0 & \text{otherwise.} \end{cases}$$

The parameter λ in the above expression, is called the *decay factor* [76], and takes a value between 0

and 1. Without this parameter, terms in the sum 1.2 corresponding to large subset trees with the same topology would be similarly large and tend to dominate the kernel score. λ penalizes Δ more strongly as the number of recursive calls increases, which downweights the largest matching substructures and allows smaller matches to contribute more to the kernel score. In this work, we refer to the parameters λ and σ as *meta-parameters*, to avoid confusing them with model parameters we are trying to estimate. When evaluating the tree kernel, it is helpful to reorder the children of each internal node such that the larger of the two subtrees is on the right-hand side. If the two subtrees have equal sizes, then the child with the longer branch length can be put on the right-hand side. This operation is referred to as *ladderizing*. Since the ordering of children is arbitrary in phylogenies, this operation ensures that a maximal number of matching subset trees are counted by the tree kernel without making meaningful changes to the trees.

The tree kernel was later shown to be highly effective in differentiating trees simulated under a compartmental model with two risk groups of varying contact rates [22]. In that paper, Poon used the tree kernel as the distance function in approximate Bayesian computation (ABC) (see section 1.5), to fit epidemiological models to observed trees.

1.3 Contact networks

1.3.1 Overview

Epidemics spread through populations of hosts through *contacts* between those hosts. The definition of contact depends on the mode of transmission of the pathogen in question. For an airborne pathogen like influenza, a contact may be simple physical proximity, while for human immunodeficiency virus (HIV), contact could be via unprotected sexual relations or blood-to-blood contact (such as through needle sharing). A *contact network* is a graphical representation of a host population and the contacts among its members [8, 77, 78]. The *nodes* in the network represent hosts, and *edges* or *links* represent contacts between them. A contact network is shown in fig. 1.2 (left). Contact networks are a particular type of *social network* [79, 80], which is a network in which edges may represent any kind of social or economic relationship. Social networks are frequently used in the social sciences to study phenomena where relationships between people or entities are important [for a review see 81].

Edges in a contact networks may be *directed*, representing one-way transmission risk, or *undirected*, representing symmetric transmission risk. For example, a network for an airborne epidemic would use undirected edges, because the same physical proximity is required for a host to infect or to become infected. However, an infection which may be spread through blood-to-blood contact through transfusions would use directed edges, since the recipient has no chance of transmitting to the donor. Directed edges are also useful when the transmission risk is not equal between the hosts, such as with HIV transmission among men who have sex with men (MSM), where the receptive partner carries a higher risk of infection than the insertive partner [82]. In this case, a contact could be represented by two directed edges, one in each direction between the two hosts, with the edges annotated by what kind of risk they imply [81]. An undirected contact network is equivalent to a directed network where each contact is represented by two

symmetric directed edges. The *degree* of a node in the network is how many contacts it has. In directed networks, we may make the distinction between *in-degree* and *out-degree*, which count respectively the number incoming and outgoing edges. The *degree distribution* of a network denotes the probability that a node has any given number of links. The set of edges attached to a node are referred to as its *incident* edges.

Epidemiological models most often assume some form of contact homogeneity. The simplest models, such as the susceptible-infected-recovered (SIR) model [5], assume a completely homogeneously mixed population, where every pair of contacts is equally likely. More sophisticated models partition the population into groups with different contact rates between and among each group [83]. However, these models still assume that every possible contact between a member of group i and a member of group j is equally likely. This assumption is clearly unrealistic for the majority of human communities and can lead to errors in predicted epidemic trajectories when there is substantial heterogeneity present [6, 84, 85]. Contact networks provide a way to relax this assumption by representing individuals and their contacts explicitly. It is important to note that, although panmixia is an unrealistic modelling assumption, it has not proven a substantial hurdle to epidemic modelling in practice [5]. Using this assumption, researchers have been able to derive estimates of the transmission rate and the basic reproductive number of various outbreaks, which have agreed with values obtained by on-the-ground data collection [86]. Therefore, if one is interested only in these population-level variables, the additional complexity of contact network models may not be warranted. Rather, these models are most useful when we are interested in properties of the network itself, such as centrality, structural balance, and transitivity [81].

From a public health perspective, knowledge of contact networks has the potential to be extremely useful. On a population level, network structure can dramatically affect the speed and pattern of epidemic spread [*e.g.* 7, 87]. For example, epidemics are expected to spread more rapidly in networks having the “small world” property, where the average path length between two nodes in the network is relatively low [88]. Some sexually transmitted infections would not be expected to survive in a homogeneously mixed population, but their long-term persistence can be explained by contact heterogeneity [5, 89]. Hence, the contact network can provide an idea of what to expect as an epidemic unfolds. In terms of actionable information, the efficacy of different vaccination strategies may depend on the topology of the network [8–10, 90]. On a local level, contact networks can be informative about the groups or individuals who are at highest risk of acquiring or transmitting infection who would therefore benefit most from public health interventions [29, 30].

Contact networks are a challenging type of data to collect, requiring extensive epidemiological investigation in the form of contact tracing [8, 40, 43, 78]. Therefore, it has been necessary to explore less resource-intensive alternatives which still contain information about population structure. For instance, it is possible to obtain limited information about the contact network by individual interviews without contact tracing. Variables which can be estimated in this fashion are referred to as *node-level* measures [81]. One of the most well-studied of these is the degree distribution mentioned above, which can theoretically be estimated by simply asking each person how many contacts they had in some interval of time. However, the degree distributions often observed in real-world sexual networks are heavy-

tailed [23–25], so dense or respondent-driven sampling [91] would be needed to capture the high-degree nodes characterizing the tail of the distribution.

An alternative approach has been the analysis of other types of network, which can be directly estimated with phylogenetic methods from viral sequence data. Some work focuses on the *phylogenetic network*, in which two nodes are connected if the genetic distance between their viral sequences is below some threshold. Primarily, this work has focused on the detection of *phylogenetic clusters*, which are groups of individuals whose viral sequences are significantly more similar to each other’s than to the general population’s. The phylogenetic network is informative about “hotspots” of transmission and can be used to identify demographic groups to whom targeted interventions are likely to have the greatest effect [92]. However, this network may show little to no agreement with contact data obtained through epidemiological methods [93–95] and therefore may be a poor proxy for the contact network. Other studies [96] have investigated the *transmission network*, which is the subgraph of the contact network consisting of infected nodes and the edges that led to their infections [40] (fig. 1.2, left). It is possible to estimate the transmission network phylogenetically, although the methods required for doing so are more sophisticated than for estimating the phylogenetic network [96]. These studies again mostly focus on clustering and also on degree distributions.

Other statistical methods have been developed to infer contact network parameters strictly from the timeline of an epidemic, using neither genetic data nor reported contacts. Britton and O’Neill [97] developed a Bayesian method to infer the p parameter of an Erdős-Rényi (ER) network, along with the transmission and removal rate parameters of the susceptible-infected (SI) model, using observed infection and optionally removal times. However, it was designed for only a small number of observations, and was unable to estimate p independently from the transmission rate. Groendyke, Welch, and Hunter [98] significantly updated and extended the methodology of Britton and O’Neill and applied it to a measles outbreak affecting 188 individuals. They were able to obtain a much more informative estimate of p , although this data set included both symptom onset and recovery times for all individuals and was unusual in that the entire contact network was presumed to be infected. Volz [87] developed differential equations describing the dynamics of the SIR model on a wide variety of random networks defined by their degree distributions. Although the topic of estimation was not addressed in the original paper, Volz’s method could in principle be used to fit such models to observed epidemic trajectories, similar to what is done with the ordinary SIR model. Volz and Meyers [84] later extended the method to dynamic contact networks and applied it to a sexual network relating 99 individuals investigated during a syphilis outbreak.

1.3.2 Scale-free networks and preferential attachment

A *scale-free* network is one whose degree distribution follows a power law, meaning that the number of nodes in the network with degree k is proportional to $k^{-\gamma}$ for some constant γ [28]. Scale-free networks are characterized by a large number of nodes of low degree, with relatively few “hub” nodes of very high degree. Epidemiological surveys have indicated that human sexual networks tend to be scale-free [23–26]. Interestingly, many other types of network, including computer networks [89], biological

metabolic networks [99], and academic co-author networks [100], also have the scale-free property.

Several properties of scale-free networks are relevant in epidemiology. The high-degree hub nodes are known as *superspreaders* [101], which have been postulated to contribute in varying degree to the spread of diseases such as HIV [38] and severe acute respiratory syndrome (SARS) [102]. Scale-free networks have no epidemic threshold [89], meaning that diseases with arbitrarily low transmissibility [can persist have a chance, however small, of persisting](#) at low levels indefinitely. This is in contrast with homogeneously mixed populations, in which transmissibility below the epidemic threshold would result in exponential decay in the number of infected individuals and eventual extinction of the pathogen [5].

One mechanism which has been shown to lead to scale-free networks is *preferential attachment* [28, 103]. The simplest preferential attachment model is known as the Barabási-Albert (BA) model after its inventors [28]. Under this model, networks are formed by starting with a small number m_0 of nodes. New nodes are added one at a time until there are a total of N in the network. Each time a new node is added, $m \geq 1$ edges are added from it to other nodes in the graph. In the original formulation, the partners of the new node are chosen with probability linearly proportional to their degree plus one.

There has been some contention over the idea that contact networks are scale-free. Handcock and Jones [27] fit several stochastic models of partner formation to empirical degree distributions derived from population surveys of sexual behaviour. They found that a negative binomial distribution, rather than a power law, was the best fit to five out of six datasets, although the difference in goodness of fit was extremely small in four out of these five. Bansal, Grenfell, and Meyers [6] found that an exponential distribution, rather than a power law, was the best fit to degree distributions of six social and sexual networks. [Dombrowski et al. \[104\] contend that sexual networks are shaped more by homophily \(“like attract like”\) than by preferential attachment, but find that injection drug users \(IDU\) network do demonstrate a scale-free structure.](#)

In the paper describing the BA model, Barabási and Albert suggest an extension where the probability of choosing a partner of degree d is proportional to $d^\alpha + 1$ for some constant α . [When \$\alpha \neq 1\$, the degree distribution no longer follows a power law \[105\]. For \$\alpha < 1\$, the distribution is a stretched exponential, meaning that the number of nodes of degree \$k\$ is proportional to \$\exp\(-k^\beta\)\$ for some constant \$\beta\$. For \$\alpha > 1\$, the distribution takes on a characteristic called *gelation*, where a one or a few high-degree hub nodes are connected to nearly every other node in the graph. We do not believe these departures from the power law affect the applicability of the model to real world networks. In fact, de Blasio, Svensson, and Liljeros \[106\] were able to estimate the preferential attachment power from partner count data collected from the same individuals for consecutive time intervals, and found a value less than one in all cases. It is also worth noting that, in addition to the BA model, other investigations of the interaction between contact networks and transmission trees have studied the Erdős-Rényi and Watts-Strogatz models \[107\], whose degree distributions do not generally follow a power law under any parameter settings.](#)

When $m = 1$, the network takes on the distinctive shape of a tree, that is, it does not contain any cycles. Cycles are present in the network for all other m values. Examples of BA networks with three different values of the preferential attachment power α are shown in fig. 1.3.

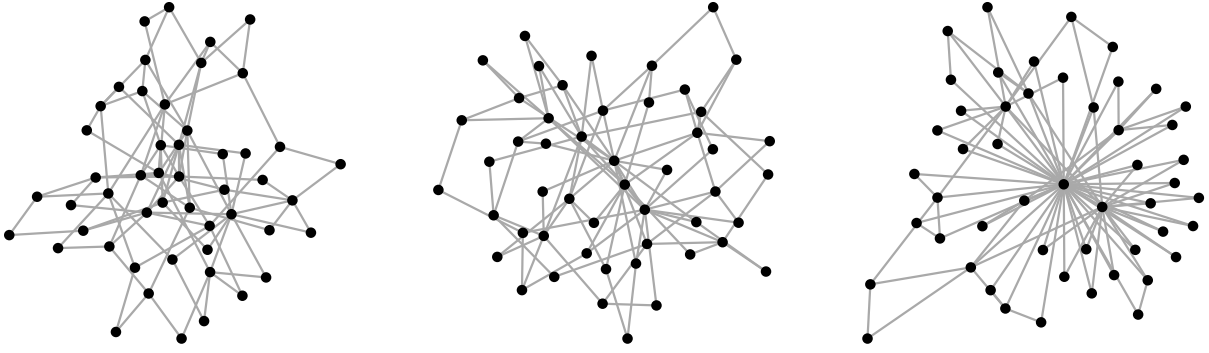


Figure 1.3: Examples of Barabási-Albert networks with preferential attachment power $\alpha = 0$ (left), 1 (centre), and 2 (right). All networks have $N = 50$ nodes and were constructed with $m = 2$ edges per vertex. When $\alpha = 0$, attachments are formed at random and most nodes have low degree. When $\alpha = 1$, preferential attachment is linear and several higher-degree nodes are observable. When $\alpha = 2$, preferential attachment is quadratic and nearly every vertex is attached to a small number of hub nodes.

1.3.3 Relationship between network structure and transmission trees

The contact network underlying an epidemic constrains the shape of the transmission network, which in turn determines the topology of the transmission tree relating the infected hosts (fig. 1.2). The index case who introduces the epidemic into the network becomes the root of the tree. Each time a transmission occurs, the lineage corresponding to the donor host in the tree splits into two, representing the recipient lineage and the continuation of the donor lineage. Figure 1.2 illustrates this correspondence. It must be emphasized that, although the order and timing of transmissions determines the tree topology uniquely, the converse does not hold. That is, for any given topology, there are in general many transmission networks which would lead to that topology. In other words, it is impossible to distinguish who transmitted to whom from a transmission tree alone [108].

A number of studies have made progress in quantifying the relationship between contact networks and transmission trees. O’Dea and Wilke [109] simulated epidemics over networks with four types of degree distribution. They then estimated the Bayesian skyride [110] population size trajectory in two ways: from the phylogeny, using MCMC; and from the incidence and prevalence trajectories, using the method developed by Volz et al. [53]. The concordance between the two skyrides, as well as the relationship between the skyride and prevalence curve, was qualitatively different for each degree distribution. Leventhal et al. [107] investigated the relationship between transmission tree imbalance and several epidemic parameters under four contact network models and found that these relationships varied considerably depending on which model was being considered. The authors also investigated a real-world HIV phylogeny and found a level of imbalance inconsistent with a randomly mixing population. Welch [111] simulated transmission trees over networks with varying degrees of community structure. They found that transmission trees simulated under networks with low clustering could not generally be distinguished from those simulated under highly clustered networks and concluded that contact network clusters do not affect transmission tree shape. However, more recently, Villandre et al. [112] investigated the correspondence between contact network clusters and transmission tree clusters and did find a

moderate correspondence between the two in some cases. Goodreau [113] combined a dynamic contact network model with a model of within-host viral evolution to simulate viral phylogenies over eight types of contact network. Estimates of prevalence and effective population size were calculated for each simulated phylogeny under three models of epidemic growth. The author found that estimates for networks with a small high-risk subgroup and networks involving commercial sex workers were substantially different than estimates for random networks or networks with segregated equal-risk groups.

1.4 Sequential Monte Carlo

1.4.1 Overview and notation

Recall that the primary objective of our work is to develop a statistical inference method for estimating contact network parameters from transmission trees. For a network model with parameters θ and an input transmission tree T , our goal is to estimate statistics, such as means and credible intervals, of the posterior distribution

$$\pi(\theta | T) = \frac{f(T | \theta)\pi(\theta)}{\int_{\Theta} f(T | \theta)\pi(\theta) d\theta}. \quad (1.3)$$

As alluded to in section 1.1, both the likelihood $f(T | \theta)$ and the normalizing constant are likely computationally intractable (this will be discussed further in section 2.1.4). Hence, rather than computing the posterior distribution analytically, we will approximate it using a *Monte Carlo* approach. The fundamental idea behind Monte Carlo methods is succinctly expressed by Liu, Chen, and Logvinenko [114]:

Monte Carlo’s view of the world is that any probability distribution π , regardless of its complexity, can always be *represented* by a discrete sample from it. By “represented”, we mean that any computation of expectations using π can be replaced to an acceptable degree of accuracy by using the empirical distribution resulting from the discrete sample.

In other words, if we are able to sample enough points from a distribution of interest, we will be able to make reasonably accurate statements about the distribution itself. For example, the expected value of the distribution can be estimated by the sample’s population mean. The reason Monte Carlo methods will be useful in this work is that algorithms exist for obtaining samples from distributions that are analytically intractable and from which direct sampling is not possible (for a review see [115]). Sequential Monte Carlo (SMC) [116–118] is one such algorithm.

SMC considers a population of points or “particles”, here denoted $\{x^{(k)}\}$ and indexed by an integer k . The particles are associated with weights, $\{w(x^{(k)})\}$. After the algorithm has been run, the weighted particles are a Monte Carlo representation for the target distribution. For example, the expected value is approximated by

$$\frac{1}{n} \sum_{k=1}^n x^{(k)} w(x^{(k)}),$$

where n is the number of particles. Initially, the particles do not represent the target distribution but rather a more tractable distribution from which direct sampling is straightforward. The word “sequential” is used to describe the iterative process of perturbation, resampling, and reweighting applied to the particles in such a way that they converge, collectively, to a representation of the target. SMC is also known as the *particle filter*.

In this work, the distribution of interest is the posterior distribution 1.3. The particles are particular values of the parameters θ of the contact network model being studied. If we were taking a typical Bayesian Monte Carlo approach to this problem, the particles would end up weighted by their posterior probability and distributed in such a way that the weighted population was a reasonable representation of $\pi(\theta | T)$. In our case, due to the intractable likelihood, we will need to consider an approximation to the posterior (see section 1.5). However, for now, nothing is lost by assuming that our target distribution is the posterior itself.

Here we describe an algorithm called the SMC sampler, developed by Del Moral, Doucet, and Jasra [119], which forms the basis of the adaptive ABC-SMC algorithm we apply toward the main objective of this work. We begin by describing sequential importance sampling (SIS), which is a precursor to SMC that samples from a sequence of distributions defined on spaces of increasing dimension. We then describe SMC itself, which extends SIS with a resampling step to fight particle degeneracy. Finally, we outline the SMC sampler, which allows SMC to be applied to sequences of distributions all defined on the same space. This terminology will become clear as the methods are described.

~~SMC is the name for a family of statistical inference methods that rely on approximating probability distributions of interest with large collections of particles, here denoted $\{x^{(k)}\}$ [116–118]. These collections or populations are constructed to form a Monte Carlo approximation to some distribution of interest π , meaning that the empirical distribution of the particles converges in distribution to π as the population size gets large [114]. The word *sequential* is used because the particle populations are modified in an iterative fashion over time, for example, to incorporate new evidence.~~

To fully describe SMC, we will introduce some notation and terminology. The definitions of these terms will become clearer as they are used. For a sequence x_1, \dots, x_d , we will write \mathbf{x}_i to mean the partial sequence x_1, \dots, x_i . The subscript $^{(k)}$ will be used to indicate the k th particle in a population. To ease the notational burden we will omit the superscripts and subscripts on the weight functions w . We define a *Markov kernel* as the continuous analogue of the transition matrix in a finite-state Markov model. For some spaces X and Y , $K : X \times Y \rightarrow [0, 1]$ such that

$$\int_Y K(x, y) dy = 1 \quad (1.4)$$

for all $x \in X$. This is an “operational” definition of Markov kernel which will be suitable for our purposes. A more rigorous definition can be found in e.g. [120]. Note that Markov kernels have nothing to do with the kernel functions defined in section 1.2.4, other than sharing a name (the word “kernel” is ubiquitous in mathematics). Also, the variable π is used here for a generic target distribution to describe the algorithm, and does not refer to the prior or posterior distributions we will eventually want to work with.

1.4.2 Sequential importance sampling

Sequential importance sampling (SIS) [121] [is a particle-based method](#) whose aim is to sample from a distribution π on an high-dimensional space, say $\pi(\mathbf{x}) = \pi(x_1, \dots, x_d)$. The basis of SIS is importance sampling (IS), which is a method of estimating summary statistics of distributions which are known only up to a normalizing constant, and therefore cannot be sampled from directly. That is, if π is such a distribution and f is any real-valued function, IS is concerned with estimating

$$\pi(f) = \int f(x)\pi(x)dx = \int f(x)\frac{\gamma(x)}{Z}dx,$$

where the integral is over the space on which π is defined, $\gamma(x)$ is known pointwise, and $Z = \int \gamma(x)dx$ is the unknown normalizing constant. Suppose we have at hand another distribution η , called the *importance distribution*, from which we are able to sample. Define the *importance weight* as the ratio $w(x) = \gamma(x)/\eta(x)$. We can write the expectation of interest as

$$\int f(x)\pi(x)dx = \frac{\int f(x)\gamma(x)dx}{\int w(x)\eta(x)dx} = \frac{1}{Z} \int w(x)\eta(x)f(x)dx. \quad (1.5)$$

[Since \$\eta\$ can be sampled from exactly, and \$\gamma\$ and \$f\$ can both be evaluated pointwise, the integral \$\int w\(x\)\eta\(x\)f\(x\)dx\$ can be approximated by a Monte Carlo estimate. We simply take a sample from \$\eta\$, multiply each point in the sample by \$f\$ and \$\gamma\$ evaluated at that point, and sum the results. Moreover, the normalizing constant \$Z\$ can be expressed in terms of the importance weight and distribution, \$Z = \int w\(x\)\eta\(x\)dx\$. Therefore, we have all the ingredients we need to obtain an estimate of \$\pi\(f\)\$ using eq. \(1.5\).](#) Although this is a simple and elegant approach, the drawback is that the variance of the estimate is proportional to the variance of the importance weights [118], which may be quite large if η and γ are very different. Therefore, the practical use of IS on its own is limited, since it depends on finding an importance distribution similar to π , which we usually know very little about *a priori*.

The objective of SIS is to build up an importance distribution η for π sequentially. By the general product rule, $\pi(\mathbf{x})$ can be decomposed as

$$\pi(\mathbf{x}) = \pi(x_1)\pi(x_2 | x_1) \cdots \pi(x_{d-1} | \mathbf{x}_{d-2})\pi(x_d | \mathbf{x}_{d-1}).$$

This decomposition is natural in many contexts, particularly for on-line estimation. For example, in a stateful model like an hidden Markov model (HMM), x_i may represent the state at time i , with $\pi(\mathbf{x})$ being the posterior distribution over possible paths. The importance distribution η for π will be constructed using a similar decomposition,

$$\eta(\mathbf{x}) = \eta(x_1)\eta(x_2 | x_1) \cdots \eta(x_{d-1} | \mathbf{x}_{d-2})\eta(x_d | \mathbf{x}_{d-1}).$$

The importance weights for η can be written recursively as

$$\begin{aligned}
w(\mathbf{x}_i) &= \frac{\pi(\mathbf{x}_i)}{\eta(\mathbf{x}_i)} && \text{definition of importance weight} \\
&= \frac{\pi(x_i | \mathbf{x}_{i-1})\pi(\mathbf{x}_{i-1})}{\eta(x_i | \mathbf{x}_{i-1})\eta(\mathbf{x}_{i-1})} && \text{definition of conditional probability} \\
&= \frac{\pi(x_i | \mathbf{x}_{i-1})}{\eta(x_i | \mathbf{x}_{i-1})} \cdot w(\mathbf{x}_{i-1}) && \text{definition of importance weight.}
\end{aligned} \tag{1.6}$$

Thus, we can choose $\eta(x_i | \mathbf{x}_{i-1})$ such that the variance of the importance weights is as small as possible at every step, eventually arriving at a full importance distribution. This choice is made on a problem-specific basis, taking any available information about $\pi(x_i | \mathbf{x}_{i-1})$ into account (see *e.g.* [118, 122] for many examples). One potential choice for $\eta(x_i | \mathbf{x}_{i-1})$ is simply $\pi(x_i | \mathbf{x}_{i-1})$, if it is possible to compute. In a Bayesian setting, the prior distribution may be used. The exact form of $\eta(x_i | \mathbf{x}_{i-1})$ which minimizes the variance of the weights is called the *optimal kernel* [123], the name deriving from the fact that $k(x_i, \mathbf{x}_{i-1}) = \eta(x_i | \mathbf{x}_{i-1})$ is a Markov kernel. In some applications, it is possible to approximate the optimal kernel or even compute it explicitly.

The recursive definition 1.6 suggests an algorithm for obtaining a sample from [η and using it to obtain an approximate sample from π by IS](#) (algorithm 1). We begin with n particles, which have been sampled from the importance distribution $\eta(x_1)$ for $\pi(x_1)$. ~~The particles are updated and reweighted d times, corresponding to the d elements of the decomposition of π . At the i th step, each particle is extended to include x_i drawn according to the chosen $\eta(x_i | \mathbf{x}_{i-1})$, and the importance weights are recalculated and normalized.~~ [Each particle is extended by drawing \$x_2\$ from \$\eta\(x_2 | x_1\)\$, and the importance weights are updated by eq. \(1.6\). This procedure is repeated \$d\$ times until the particles have dimension equal to that of \$\pi\$.](#)

Algorithm 1 Sequential importance sampling.

```

for  $k = 1$  to  $n$  do
  Sample  $x_1^{(k)}$  from  $\eta(x_1)$                                 ▷ Initialize the  $k$ th particle
   $w^{(k)} \leftarrow \pi(x_1^{(k)}) / \eta(x_1^{(k)})$                 ▷ Initialize importance weight
end for
for  $i = 2$  to  $d$  do
  for  $k = 1$  to  $n$  do
    Sample  $x_i^{(k)}$  from  $\eta(x_i | \mathbf{x}_{i-1}^{(k)})$                 ▷ Extend the  $k$ th particle
     $w^{(k)} \leftarrow [\pi(x_i^{(k)} | \mathbf{x}_{i-1}^{(k)}) / \eta(x_i^{(k)} | \mathbf{x}_{i-1}^{(k)})] \cdot w^{(k)}$     ▷ Update importance weight
  end for
  Normalize the weights so that  $\sum w = 1$ 
end for

```

1.4.3 Sequential Monte Carlo

The importance distribution η constructed with SIS is merely an approximation to π , and may be a fairly poor one in practice depending on the application. Try as we might to keep the variances of the

weights low, the cumulative errors at each sequential step tend to push many of the weights to very low values [116]. This results in a poor approximation to π , since only a few particles retain high importance weights after all d sequential steps, a problem known as *particle degeneracy*. To mitigate this problem, Gordon, Salmond, and Smith [121] introduced technique they called the *bootstrap filter*, which involves a resampling of the population of particles after each sequential step in accordance with their importance weights. A similar idea, termed *particle rejuvenation*, was proposed by Liu and Chen [124]. These approaches cause particles with high importance weights to be replicated in the population, while particles with low weights may be removed. After each resampling step, the importance weights for all particles are set equal.

The resampling step was formally integrated with SIS by Doucet, Godsill, and Andrieu [116] to form the first SMC algorithm (algorithm 2). Rather than resample at every step as the bootstrap filter proposed, the authors use a criterion based on the expected sample size (ESS) the particle population to determine when resampling is necessary. The ESS of the population of particles is defined as

$$\text{ESS}(w) = \frac{n}{1 + \text{Var}(w)},$$

where n is the number of particles in the population. Resampling is triggered when the ESS drops below the threshold (conventionally $n/2$ [118]). ~~This results in the removal of low-weight particles from the population, and also equalizes all the weights.~~ Various resampling strategies beyond basic sampling with replacement have been proposed [125], but we will not discuss those here.

Algorithm 2 Sequential Monte Carlo [116].

```

for  $k = 1$  to  $n$  do
    Sample  $x_1^{(k)}$  from  $\eta(x_1)$                                 ▷ Initialize the  $k$ th particle
     $w^{(k)} \leftarrow \pi(x_1^{(k)}) / \eta(x_1^{(k)})$                     ▷ Initialize importance weight
end for
for  $i = 2$  to  $d$  do
    for  $k = 1$  to  $n$  do
        Sample  $x_i^{(k)}$  from  $\eta(x_i | \mathbf{x}_{1:i-1}^{(k)})$                 ▷ Extend the  $k$ th particle
         $w^{(k)} \leftarrow [\pi(x_i^{(k)} | \mathbf{x}_{1:i-1}^{(k)}) / \eta(x_i^{(k)} | \mathbf{x}_{1:i-1}^{(k)})] \cdot w^{(k)}$     ▷ Update importance weight
    end for
    if  $\text{ESS}(w) < T$  then                                    ▷  $T$  is a user-defined threshold
        Resample the particles according to  $w$ 
        for  $k = 1$  to  $n$  do
             $w^{(k)} \leftarrow 1/n$ 
        end for
    end if
end for

```

1.4.4 The sequential Monte Carlo sampler

The SIS and SMC algorithms described above aim to sample from a high-dimensional distribution $\pi(\mathbf{x})$, by sequentially sampling from d distributions of lower but increasing dimension. Del Moral, Doucet, and Jasra [119] developed the *SMC sampler* with an alternative objective: to sample sequentially from d distributions π_1, \dots, π_d , all of the *same* dimension and defined on the same space. The π_i are assumed to form a related sequence, such as posterior distributions attained by sequentially considering new evidence. As with SIS, we assume that $\pi_i(x) = \gamma_i(x)/Z_i$, where each γ_i is known pointwise and the normalizing constants Z_i are unknown.

Both algorithms involve progression through a sequence of related distributions. For SIS and SMC, these distributions are lower-dimensional marginals of the target distribution, while for the SMC sampler, they are of the same dimension and constitute a smooth progression from an initial to a final distribution. In both cases, the neighbouring distributions in the sequence are related to each other in some way, and we can take advantage of that relationship to create a sequence of importance distributions alongside the sequence of targets. In SIS, the neighbouring marginals $\pi(\mathbf{x}_i)$ and $\pi(\mathbf{x}_{i+1})$ were related by the conditional density $\pi(x_i | \mathbf{x}_{i-1})$, which we used to inform the importance distribution. In SMC, the relationship between subsequent distributions is less explicit, but it is assumed that they are related closely enough that an importance distribution for π_i can be easily transformed into one for π_{i+1} . In particular, the sequence of importance distributions η_i is constructed as

$$\eta_i(x') = \int \eta_{i-1}(x) K_i(x, x') dx, \quad (1.7)$$

where K_i is a Markov kernel and the integral is over the space on which the π_i are defined. The choice of K_i should be based on the perceived relationship between π_{i-1} and π_i . Del Moral, Doucet, and Jasra [119] propose the use of a MCMC kernel with equilibrium distribution π_i . That is,

$$K_i(x, x') = \max \left(1, \frac{q(x', x) \pi_i(x)}{q(x, x') \pi_i(x')} \right),$$

where $q(x, x')$ is a proposal function such as a Gaussian distribution centred at x [from which \$x'\$ is drawn](#) (see ??).

Although this method of building up η appears straightforward, the drawback is that the importance distribution itself becomes intractable. In particular, evaluating $\eta_i(x)$ involves a i -dimensional integral of the type in eq. (1.7). As it is necessary to evaluate $\eta(x)$ pointwise to perform IS, this construction appears to have defeated the purpose of providing an importance distribution for each π_i . Del Moral, Doucet, and Jasra [119] overcome this problem with two “artificial” objects. First, they propose the existence of *backward* Markov kernels $L_{i-1}(x_i, x_{i-1})$. For now, these kernels are arbitrary; they will later be precisely defined on a problem-specific basis. Second, the authors define an alternative sequence of target distributions

$$\tilde{\pi}_i(\mathbf{x}_i) = \pi_i(x_i) \prod_{k=1}^{i-1} L_k(x_{k+1}, x_k)$$

of increasing dimension. That is, $\tilde{\pi}_i$ has dimension equal to the original dimension of π_i raised to the power of i . This brings us back to the SIS setting described above (section 1.4.2), namely of building up an importance distribution sequentially through lower-dimensional distributions. The dimension of the importance distributions η is similarly augmented with the forward kernels,

$$\eta_i(\mathbf{x}_i) = \eta_1(x_1) \prod_{k=1}^{i-1} K_k(x_k, x_{k+1}) \quad (1.8)$$

Thanks to the backwards kernels, we can write $\tilde{\pi}_i$ in terms of $\tilde{\pi}_{i-1}$ as follows.

$$\frac{\tilde{\pi}_i(\mathbf{x}_i)}{\tilde{\pi}_{i-1}(\mathbf{x}_{i-1})} = \frac{\pi_i(x_i) \prod_{k=1}^{i-1} L(x_{k+1}, x_k)}{\pi_{i-1}(x_{i-1}) \prod_{k=1}^{i-2} L(x_{k+1}, x_k)} = \frac{\pi_i(x_i) L(x_i, x_{i-1})}{\pi_{i-1}(x_{i-1})},$$

and hence

$$\tilde{\pi}_i = \frac{\pi_i(x_i) L(x_i, x_{i-1})}{\pi_{i-1}(x_{i-1})} \cdot \tilde{\pi}_{i-1}.$$

$$\begin{aligned} \tilde{\pi}_i(\mathbf{x}_i) &= \pi_i(x_i) \prod_{k=1}^{i-1} L_k(x_{k+1}, x_k) && \text{definition of } \tilde{\pi}_i \\ &= \pi_i(x_i) L_{i-1}(x_{i-1}, x_i) \prod_{k=1}^{i-2} L_k(x_{k+1}, x_k) && \text{pull out } k = i-1 \text{ term from product} \\ &= \frac{\pi_i(x_i) L_{i-1}(x_{i-1}, x_i)}{\pi_{i-1}(x_{i-1})} \cdot \pi_{i-1}(x_{i-1}) \prod_{k=1}^{i-2} L_k(x_{k+1}, x_k) && \text{multiply and divide by } \pi_{i-1}(x_{i-1}) \\ &= \frac{\pi_i(x_i) L_{i-1}(x_{i-1}, x_i)}{\pi_{i-1}(x_{i-1})} \cdot \tilde{\pi}_{i-1}(\mathbf{x}_{i-1}) && \text{definition of } \tilde{\pi}_{i-1}. \end{aligned}$$

The importance distributions can also be expressed recursively using the forward kernels, which follows directly from eq. (1.8),

$$\eta_i(\mathbf{x}_i) = \eta_{i-1}(\mathbf{x}_{i-1}) K_i(x_{i-1}, x_i).$$

Therefore, the importance weights for these new targets are defined recursively as

$$\begin{aligned} w(\mathbf{x}_i) &= \frac{\tilde{\pi}_i(\mathbf{x}_i)}{\eta_i(\mathbf{x}_i)} && \text{definition of importance weight} \\ &= \frac{\tilde{\pi}_i(\mathbf{x}_i)}{\eta_{i-1}(\mathbf{x}_{i-1}) K_i(x_{i-1}, x_i)} && \text{recursive definition of } \eta_i \\ &= \frac{\tilde{\pi}_{i-1}(\mathbf{x}_{i-1}) \pi_i(x_i) L_{i-1}(x_i, x_{i-1})}{\eta_{i-1}(\mathbf{x}_{i-1}) \pi_{i-1}(x_{i-1}) K_i(x_{i-1}, x_i)} && \text{recursive definition of } \pi_i \\ &= w(\mathbf{x}_{i-1}) \cdot \frac{\pi_i(x_i) L_{i-1}(x_i, x_{i-1})}{\pi_{i-1}(x_{i-1}) K_i(x_{i-1}, x_i)} && \text{definition of importance weight} \\ &\propto w(\mathbf{x}_{i-1}) \cdot \frac{\gamma_i(x_i) L_{i-1}(x_i, x_{i-1})}{\gamma_{i-1}(x_{i-1}) K_i(x_{i-1}, x_i)} && \text{remove normalizing constant } Z_i/Z_{i-1} \end{aligned} \quad (1.9)$$

The final key piece of information is to notice that, because the L_i are Markov kernels, π_i is simply the marginal in \mathbf{x}_{i-1} of $\tilde{\pi}$. Therefore, a sample from $\tilde{\pi}_i$ automatically gets us a sample from π_i , by considering only the i th component of \mathbf{x}_i . In fact, since the weight update eq. (1.9) depends only on the i th and $i - 1$ st components of each particle, we do not even need to keep track of the complete particles if we are only interested in the final distribution. ~~These are all the ingredients we need to apply SIS.~~ The sequences of kernels L and K should be chosen based on the problem at hand to minimize the variance in the importance weights as well as possible. For a fixed choice of K , the backward kernels which minimize this variance are called the *optimal* backward kernels. The full SMC sampler algorithm is presented as algorithm 3. A resampling step is applied whenever the ESS of the population drops too low, as discussed in the previous section.

Algorithm 3 Sequential Monte Carlo sampler of Del Moral, Doucet, and Jasra [119].

```

for  $k = 1$  to  $n$  do
  Sample  $x_1^{(k)}$  from  $\eta_1(x_1)$                                 ▷ Initialize the  $k$ th particle
   $w^{(k)} \leftarrow \gamma_1(x_1^{(k)})/\eta_1(x_1^{(k)})$                 ▷ Initialize the importance weights
  Normalize the weights so that  $\sum w = 1$ 
end for
for  $i = 2$  to  $d$  do
  for  $k = 1$  to  $n$  do
    Sample  $x_i^{(k)}$  from  $K(x_{i-1}^{(k)}, x_i)$                     ▷ Extend the  $k$ th particle
     $w^{(k)} \leftarrow w^{(k)} \cdot \frac{\gamma_i(x_i)L_{i-1}(x_i, x_{i-1})}{\gamma_{i-1}(x_{i-1})K_i(x_{i-1}, x_i)}$     ▷ Update the importance weights
  end for
  Normalize the weights so that  $\sum w = 1$ 
  if  $\text{ESS}(w) < T$  then                                ▷  $T$  is a user-defined threshold
    Resample the particles according to  $w$ 
    for  $k = 1$  to  $n$  do
       $w^{(k)} \leftarrow 1/n$ 
    end for
  end if
end for

```

1.5 Approximate Bayesian computation

1.5.1 Overview and motivation

Sequential Monte Carlo, and the SMC sampler, were developed for sampling from distributions which can be evaluated up to a normalizing constant. We claim, and shall argue more thoroughly below (section 2.1.4) that the posterior distribution

$$\pi(\theta \mid T) \propto f(T \mid \theta)\pi(\theta)$$

for a contact network model with parameters θ and an input transmission tree T does not fall in this

category. Therefore, SMC, and other Bayesian and maximum likelihood (ML) techniques for fitting mathematical models (see ??), cannot be directly applied to our problem. In particular, MCMC and the SMC sampler are designed for distributions π which can be evaluated up to a normalizing constant Z , that is, $\pi(x) = \gamma(x)/Z$. Both algorithms calculate a ratio of the form $\pi(x)/\pi(x') \propto \gamma(x)/\gamma(x')$ for a current value x and proposed updated value x' - for MCMC, this is part of the Metropolis-Hastings ratio, while for the SMC sampler, a similar ratio is required to calculate the importance weights. In the context of Bayesian inference, this ratio contains a likelihood ratio, which must be calculated by computing the individual likelihoods and dividing them. If the likelihood is intractable, this is clearly not a viable approach.

Approximate Bayesian computation (ABC) [14–16] was developed to estimate posterior distributions with intractable likelihoods, which have arisen frequently in the domain of population genetics [17, 126]. ABC navigates around the intractable likelihood by replacing the posterior as the target of inference by an *approximate* posterior. This distribution is constructed in such a way that the ratios required for MCMC and the SMC sampler can be computed, conveniently allowing us to apply those algorithms with minimal changes. In the next section, we shall demonstrate how this is done, but first we give the definition of the approximate posterior.

~~Most mathematical models are amenable to fitting via one or both of the approaches, ML or Bayesian inference, discussed above. However, there are some, particularly in the domain of population genetics [17, 126], for which calculation of either the likelihood or the product of the likelihood and the prior may be infeasible. For example, one or both of these quantities may be expressible only as an intractable integral. ABC is designed for such cases, where standard likelihood-based techniques for model fitting cannot be applied.~~

~~Ordinarily, Bayesian inference targets the posterior distribution $\pi(\theta | y)$. That is, in the Bayesian framework, By targeting the posterior distribution, Bayesian inference makes the assertion that~~ model parameters with higher posterior density are “better” in the sense that they offer a more credible explanation for the observed data. The approximate posterior targeted by ABC uses an alternative metric for parameter credibility: the similarity of simulated datasets to the observed data. If datasets simulated under the model closely resemble the real data, it follows that the model is a reasonable approximation to the real-world process generating the observed data. More formally, ~~let y be the observed data to which we are trying to fit a model with parameters θ . In the case of this work, the data is a transmission tree T , but we shall stick with the generic variable y for now.~~ Suppose we have a distance measure ρ defined on the space of all possible data our model could generate. ABC aims to sample from the joint posterior distribution of model parameters and simulated datasets z which are within some small distance ε of the observed data y ,

$$\pi_\varepsilon(\theta, z | y) = \frac{\pi(\theta)f(z | \theta)\mathbb{I}_{A_{\varepsilon,y}}(z)}{\int_{A_{\varepsilon,y} \times \Theta} \pi(\theta)f(z | \theta)d\theta}. \quad (1.10)$$

Here, $A_{\varepsilon,y}$ is an ε -ball around y with respect to ρ , Θ is the space of all possible model parameters, and \mathbb{I} is the indicator function [127]. The distribution $\pi_\varepsilon(\theta, z | y)$ will be referred to as the *ABC target*

distribution. The term $f(z | \theta)$ appears to be the bothersome likelihood again, but this will turn out not to be a problem because we are simulating z ourselves.

To return to the context of this thesis, the observed data y is an estimated transmission tree for a viral epidemic under investigation. The model in question is a contact network model with parameters θ . The simulated dataset z is a transmission tree, obtained by first generating a contact network under the model, and then simulating the spread of an epidemic over that network. A transmission tree can be constructed by keeping track of who infected whom during the simulated epidemic (further details will be given in section 2.1.1). The distance function ρ must compare two transmission trees - the observed tree, which takes the place of y , and the simulated tree z . We will define this distance function using the tree kernel discussed in section 1.2.4. In words, the approximate posterior we consider here is a distribution which assigns a joint probability density to model parameters and simulated transmission trees under those parameters. The probability density is proportional to the product of the prior on the parameters, and the likelihood of the simulated transmission tree under those parameters, but only if the simulated transmission tree is sufficiently close to the true tree. Otherwise, the probability density is zero.

~~As we shall see in the next section, this distribution can be sampled from exactly. The word “approximate” derives from assumption that, for a suitably chosen distance ρ and a small enough ϵ , the marginal in z of this distribution approximates the posterior of interest [127]. That is, In fact, it is not the ABC target distribution itself, but rather its marginal in z , which approximates the posterior distribution. In other words, we claim that~~

$$\int \pi_{\epsilon}(\theta, z | y) dz \approx \pi(\theta | y).$$

The intuition for why this approximation might be reasonable comes from considering the case when $\epsilon = 0$ and ρ has the property that $\rho(z, y)$ if and only if $x = y$. In that case, ~~the fact that, when $\epsilon = 0$,~~ the ϵ -ball around y should contain only y itself, hence the integral on the left is exactly equal to the posterior. Thus, by taking ϵ small, we should attain something close to the posterior if ρ captures the similarity between datasets reasonably well. However, the accuracy of the ABC approximation depends heavily on the choice of distance function [128, 129].

Distance functions and summary statistics in ABC

In many applications (eg. [16, 130]) , ρ is defined as $\rho(S(\cdot), S(\cdot))$ where S is a function which maps data points into a vector of summary statistics. In the context of ABC, a summary statistic S is called sufficient if

$$\pi(\theta | y) = \pi(\theta | S(y)).$$

That is, sufficiency implies that the data can be replaced with the summary statistic without losing any information about the posterior distribution [131]. For most problems, it is not possible to find sufficient summary statistics [131]. A number of sophisticated methods have been developed for selecting and weighting summary statistics based on various optimality criteria [128, 129, and references therein]. We do not apply these methods in this work, instead focusing on a distance function which is not based on

[summary statistics.](#)

~~Summary statistics can be useful if the data are high-dimensional or of a complex type, but~~ [Although summary statistics are often presented as a fundamental part of ABC](#), they are not strictly necessary. For instance, if the data are numeric and of low dimension, the distance function may simply be the Euclidean distance [132]. [Park et al. \[19\] proposed the use of a kernel function \(as defined in section 1.2.4\) in place of a distance function.](#) The authors referred to their approach as “double-kernel ABC” due to [the use of a second \(unrelated\) kernel function to compute the weights of the particles.](#) The work by [Poon \[22\]](#), upon which ours is based, employed a similar approach, replacing the likelihood ratio in Bayesian MCMC with a ratio of kernel scores.

1.5.2 Algorithms for ABC

Algorithms for performing ABC can be grouped into three categories: rejection, MCMC, and SMC [127]. To simplify the notation, we shall restrict the descriptions of these algorithms to the case of one simulated dataset per parameter particle (the meaning of this will become clear shortly). The extension to multiple datasets per particle is straightforward and will be given at the end of the section. We use the variable x to refer to the pair (θ, z) , so that the ABC target distribution can be written $\pi_\epsilon(x | y)$. [This makes our notation consistent with section 1.4.](#)

Rejection ABC is the simplest method, and also the one which was first proposed [14, 15]. The algorithm, outlined in algorithm 4, repeats the following steps until a desired number of samples from the target distribution are obtained. Parameter values θ are sampled according to the prior distribution $\pi(\theta)$. Then, a simulated dataset z is generated from the model with the sampled parameter values. By definition, the probability density of obtaining the particular dataset z is $f(z | \theta)$. Finally, the parameters are sampled if the distance of z from the observed data y is less than ϵ , that is, with probability $\mathbb{I}_{A_{\epsilon,y}}(z)$. Putting this all together, the parameters θ are sampled with probability proportional to

$$\pi(\theta)f(z | \theta)\mathbb{I}_{A_{\epsilon,y}}(z),$$

which is exactly the numerator of the ABC target distribution. Thus, θ represents an unbiased sample from the approximate posterior.

Algorithm 4 Rejection ABC.

```
loop
  Draw  $\theta$  according to  $\pi(\theta)$ 
  Simulate a dataset  $z$  from the model with parameters  $\theta$ 
  if  $\rho(y, z) < \epsilon$  then
    Sample  $\theta$ 
  end if
end loop
```

Rejection ABC is easy to understand and implement, but it is not generally computationally feasible. If the posterior is very different from the prior, a very large number of samples may need to be taken

in order to find a simulated dataset which is close to z . The inefficiency is compounded by the curse of dimensionality - the measure of the ε -ball around y decreases exponentially with the number of dimensions. ABC-MCMC (algorithm 5) was designed to overcome these hurdles [133]. The approach is similar to ordinary Bayesian MCMC (??), except that a distance cutoff replaces the likelihood ratio. That is, the transition probability between states x and x' is defined as

$$\min \left(1, \frac{\pi(\theta')q(\theta', \theta)}{\pi(\theta)q(\theta, \theta')} \cdot \mathbb{I}_{A_{\varepsilon,y}}(z') \right).$$

Algorithm 5 ABC-MCMC.

Draw θ according to $\pi(\theta)$

loop

Propose θ' according to $q(\theta, \theta')$

Simulate a dataset z' according to the model with parameters θ

Accept $\theta \leftarrow \theta'$ with probability $\min \left(1, \frac{\pi(\theta')q(\theta', \theta)}{\pi(\theta)q(\theta, \theta')} \cdot \mathbb{I}_{A_{\varepsilon,y}}(z') \right)$

end loop

Some of the same computational inefficiencies arise with ABC-MCMC as with rejection. For example, in regions of low posterior density, the probability to simulate a dataset proximal to the observed data is low. Various strategies have been developed to mitigate this, including reducing the tolerance level ε as the chain progresses [134].

The most recently developed class of algorithm for ABC is ABC-SMC [132, 135]. As with ABC-MCMC, the algorithm is a straightforward modification of an existing Bayesian inference method, in this case the SMC sampler (section 1.4.4). The sequence of target distributions is defined as $\pi_i(x) = \pi_{\varepsilon_i}(x | y)$ for a decreasing sequence of tolerances ε_i . The intention is for the algorithm to progress smoothly through a sequence of target distributions which ends at the ABC approximation to the posterior. [The initial value \$\varepsilon_1\$ is set to \$\infty\$, which makes the first distribution in the sequence](#)

$$\pi_1(\theta, z) = \frac{\pi(\theta)f(z | \theta)}{\int_{\mathbb{R} \times \Theta} \pi(\theta)f(z | \theta) d\theta dz}.$$

[This initial distribution does not depend on the observed data \$y\$. In the terminology of the SMC sampler \(algorithm 2\), the numerator is the first of the \$\gamma\$'s, that is, \$\gamma_1 = \pi\(\theta\)f\(z | \theta\)\$. Sampling in proportion to \$\gamma_1\$ is straightforward and was already demonstrated for rejection ABC above. Because the sampling is exact, the initial importance weights are all set equal to 1 and normalized to \$1/n\$ where \$n\$ is the number of particles.](#)

As discussed in section 1.4.4, the choices of the kernels K and L is problem-specific, and so appropriate kernels must be chosen for ABC. Several options have been proposed [21, 132, 135]. [With an appropriately chosen kernel, the weight update \(eq. \(1.9\)\) will simplify into a computable expression. Sisson, Fan, and Tanaka \[132\] and Beaumont et al. \[135\] both suggest random walk kernels for \$K\$, where each particle is perturbed according to a Gaussian distribution. The backwards kernels \$L\$ are chosen to](#)

approximately minimize the variance in importance weights. In this thesis, we use an MCMC kernel, as proposed by Del Moral, Doucet, and Jasra [21]. The associated backwards kernels and weight updates are discussed in the next chapter (section 2.1.3). With generic kernels K and L , the ABC-SMC algorithm is almost identical to the SMC sampler (algorithm 3), with γ_i replaced with π_i . Therefore we will not repeat the full algorithm here.

All the algorithms discussed in this section can be straightforwardly extended to sample from the joint distribution

$$\pi_\varepsilon(\theta, z_1, \dots, z_M \mid y),$$

which is equivalent to associating M simulated datasets to each parameter particle instead of just one. The simulated dataset z is replaced by $z = z_1, \dots, z_M$, and the indicator function for the ε -ball around y is replaced by

$$\sum_{k=1}^M \mathbb{I}_{A_{\varepsilon, y}}(z_k).$$

For ABC-MCMC and ABC-SMC, the proposal distribution $q(\theta, \theta')f(z \mid \theta')$ is replaced by

$$q_i(\theta, \theta') \prod_{k=1}^M f(z_k \mid \theta').$$

1.6 Summary

In section 1.1, we outlined three research objectives that will be addressed in this thesis. First, we aim to develop a method for fitting contact network models to estimated transmission trees. Although transmission trees can be estimated for any epidemic by thorough contact tracing, viral diseases are the most useful context for this method due to the possibility of inferring the trees from sequence data. Transmission and sequence evolution occur on similar time scales for RNA viruses, resulting in viral phylogenies whose shapes are heavily constrained by the transmission process. The study of this interaction between evolution and epidemiology is called *phylodynamics*; phylodynamic methods make it possible to estimate transmission trees from viral sequence data. These estimated trees form the input data for our method.

The desired output of our method is a posterior distribution of the parameters of a contact network model. Rather than assuming a homogeneously mixed population, as most epidemiological models do, network models take the more realistic view that human populations are structured. That is, contacts which allow for transmission may occur in a nonrandom way, rather than between every pair of individuals in the population. A network model parameterizes this structure. In particular, our second research objective is to characterize our ability to fit the Barabási-Albert (BA) model, which incorporates preferential attachment to generate networks with realistic degree distributions.

To fit these models, our method will apply approximate Bayesian computation (ABC), a simulation-based approach. Simulating a transmission tree according to a network model is straightforward: a network can be generated according to the model, and the spread of an epidemic can be simulated over the

network and recorded in a transmission tree. ABC uses the concordance between these simulated transmission trees and the “true” estimated tree as an indicator of parameter credibility. The closer the simulated transmission trees appear to the true tree, the more weight is assigned to the associated network parameters.

ABC can be implemented by at least three classes of algorithm, but the one we choose to apply in this work is sequential Monte Carlo (SMC). SMC uses a population of parameter “particles” to approximate a distribution of interest, in this case the approximate posterior distribution targeted by ABC. After running the ABC-SMC algorithm, statistics on the model parameters can be approximated by the weighted population of particles. For example, a weighted average would give an approximate expected value for each parameter.

Thus, our method integrates four research topics: phylogenetics, contact networks, sequential Monte Carlo, and approximate Bayesian computation. The first two topics together form the problem domain. Phylogenetic data is the input to our method, while estimates of the parameters of contact network models are the desired output. The latter two topics define the algorithm and statistical framework that our inference method will use.

Chapter 2

Reconstructing contact network parameters from viral phylogenies

In this chapter, we will address the three research aims of this thesis introduced in section 1.1. First, in section 2.1, we describe *netabc*, a computer program that implements an approximate Bayesian computation-based algorithm to fit contact network models to phylogenetic data. We also provide a justification for the use of ABC for this problem by arguing that the likelihood functions required to fit these models by more conventional means are likely to be computationally intractable. Second, in section 2.2, we perform a simulation study to investigate the Barabási-Albert network model, which uses a preferential attachment mechanism to generate networks with the power law degree distributions observed in real world social and sexual networks. We progress through two exploratory analyses testing the identifiability of the model's parameters, and conclude by testing *netabc*'s ability to recover the parameters from simulated transmission trees. Third, in ??, we apply *netabc* to fit the BA model to six real world HIV datasets, with the understanding of the model parameters' identifiability gained through the simulation experiments. We conclude the chapter with a unified discussion of the three research aims, including interpretation of the results of both the simulated and real data experiments, as well as an examination of the limitations of our approach and opportunities for future investigation.

2.1 *Netabc*: a computer program for estimation of contact network parameters with ABC

Netabc is a computer program to perform statistical inference of contact network parameters from an estimated transmission tree using ABC. As discussed in section 1.1, the principal statistical algorithm used by *netabc* is adaptive ABC-SMC [21]. In addition, there are two supplementary components which are specific to the domain of phylogenetics and contact networks: Gillespie simulation [136], to simulate transmission trees on contact networks; and the tree kernel [18], which is used as the distance function in ABC to compare transmission trees [22] (see section 1.5). We give a high-level overview of the program here, before describing these components in detail. *Netabc* takes as input an estimated transmission tree,

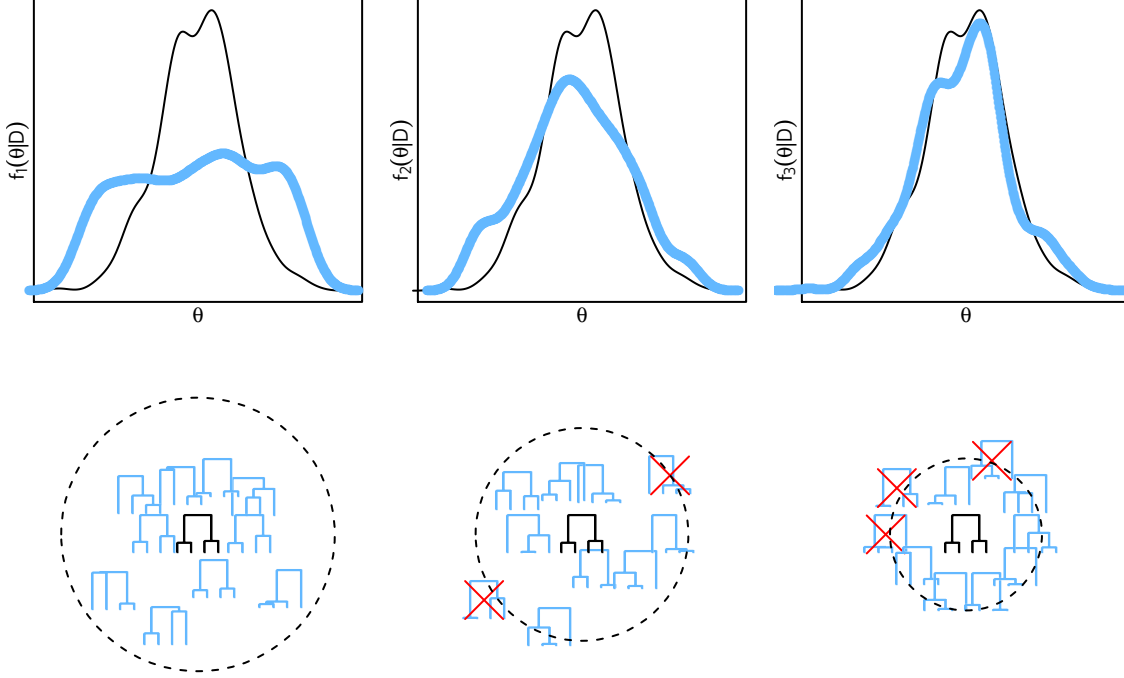


Figure 2.1: Graphical schematic of the ABC-SMC algorithm implemented in *netabc*. Particles are initially drawn from their prior distributions, making the initial population a Monte Carlo approximation to the prior. At each iteration, particles are perturbed, and a distance threshold around the true tree contracts. Particles are rejected, and eventually resampled, when all their associated simulated trees lie outside the threshold. As the algorithm progresses, the population smoothly approaches a Monte Carlo approximation of the ABC target distribution, which is assumed to resemble the posterior.

which can be derived from a viral phylogeny by rooting and time-scaling as described in section 1.2.3 or estimated by other methods [32, 54, 65–68]. We variously refer to this estimated transmission tree as the observed tree, input tree, or true tree.

As described in section 1.4, *netabc* keeps track of a population of particles $x^{(k)}$ indexed by an integer k , each of which contains particular parameter values $\theta^{(k)}$ for the [contact network](#) model we are trying to fit [to the input tree](#). A small number M of contact networks $z^{(k,i)}$, $1 \leq i \leq M$, are generated under the model for each particle in accordance with that particle’s parameters. An epidemic is simulated over each of these networks using Gillespie simulation, and by keeping track of its progress, a transmission tree is obtained. Thus, each particle becomes associated with several simulated transmission trees. These trees are compared to the input tree using the tree kernel. Particles are weighted according to the similarity of their associated simulated trees with the true tree, with more similar trees receiving higher weights. The particles are iteratively perturbed to explore the parameter space, and particles with simulated trees too distant from the true tree are periodically dropped and resampled. Once a convergence criterion is attained, the final set of particles is used as a Monte Carlo approximation to the target distribution of ABC, which is assumed to resemble the posterior distribution on model parameters (see section 1.5). A graphical schematic of this algorithm is shown in fig. 2.1.

Netabc is written in the C programming language. The *igraph* library [137] is used to generate and

store contact networks and phylogenies. Judy arrays [138] are used for hash tables and dynamic programming matrices. The GNU scientific library (GSL) [139] is used to generate random draws from probability distributions, and to ~~perform the bisection step~~ [solve for the next \$\epsilon\$ by bisection](#) in the adaptive ABC-SMC algorithm. Parallelization is implemented with Portable Operating System Interface (POSIX) threads [140]. In addition to the *netabc* binary to perform ABC, we provide three additional stand-alone utilities: *trekernel*, to calculate the tree kernel; *nettree*, to simulate a transmission tree over a contact network; and *treestat*, to compute various summary statistics of phylogenies. The programs are freely available at <https://github.com/rmcclosk/netabc>.

To check that our implementation of Gillespie simulation was correct, we reproduced Figure 1A of Leventhal et al. [107] (our ??), which plots the imbalance of transmission trees simulated over four network models at various levels of pathogen transmissibility. Our implementation of adaptive ABC-SMC was tested by applying it to the same mixture of Gaussians used by Del Moral, Doucet, and Jasra [21] to demonstrate their method (originally used by Sisson, Fan, and Tanaka [132]). We were able to obtain a close approximation to the function (see ??), and attained the stopping condition used by the authors in a comparable number of steps. To check that the algorithm would converge to a bimodal distribution, we also applied it to a mixture of two Gaussians with means ± 4 and variances 1. The algorithm was able to recover both peaks (??).

2.1.1 Simulation of transmission trees over contact networks

The simulation of epidemics, and the corresponding transmission trees, over contact networks is performed in *netabc* using the Gillespie simulation algorithm [136]. This method has been independently implemented and applied by several authors [*e.g.* 95, 98, 107, 109, 112]. Groendyke, Welch, and Hunter [98] published their implementation as an *R* package, but since the SMC algorithm is quite computationally intensive, we chose to implement our own version in *C* [as part of *netabc*](#).

Let $G = (V, E)$ be a directed contact network. We assume the individual nodes and edges of G follow the dynamics of the SIR model [2]. Each directed edge $e = (u, v)$ in the network is associated with a transmission rate β_e , which indicates that, once u becomes infected, the waiting time until u infects v is distributed as $\text{Exponential}(\beta_e)$. Note that v may become infected before this time has elapsed, if v has other incoming edges. v also has a removal rate ν_v , so that the waiting time until removal of v from the population is $\text{Exponential}(\nu_v)$. Removal may correspond to death or recovery with immunity, or a combination of both, but in our implementation recovered nodes never re-enter the susceptible population. We define a *discordant edge* as an edge (u, v) where u is infected and v has never been infected. [In the epidemiology literature, the symbol \$\gamma\$ is usually used in place of \$\nu\$; we use \$\nu\$ here to distinguish the recovery rate from the power law exponent of scale free networks \(see section 1.3.2\).](#)

To describe the algorithm, we introduce some notation and variables. Let $\text{inc}(v)$ be the set of incoming edges to v , and $\text{out}(v)$ be the set of outgoing edges from v . Let \mathcal{I} be the set of infected nodes in the network, \mathcal{R} be the set of removed nodes, \mathcal{S} be the set of susceptible nodes, and \mathcal{D} be the set of discordant edges in the network. Let β be the total transmission rate over all discordant edges, and ν be

the total removal rate of all infected nodes,

$$\beta = \sum_{e \in \mathcal{D}} \beta_e, \quad \nu = \sum_{v \in \mathcal{I}} \nu_v.$$

The variables \mathcal{S} , \mathcal{I} , \mathcal{R} , \mathcal{D} , β , and ν are all updated as the simulation progresses. When a node v becomes infected, it is deleted from \mathcal{S} and added to \mathcal{I} . Any formerly discordant edges in $\text{inc}(v)$ are deleted from \mathcal{D} , and edges in $\text{out}(v)$ to nodes in \mathcal{S} are added to \mathcal{D} . If v is later removed, it is deleted from \mathcal{I} and added to \mathcal{R} , and any discordant edges in $\text{out}(v)$ are deleted from \mathcal{D} . At the time of either infection or removal, the variables β and ν are updated to reflect the changes in the network. ~~The updates to \mathcal{S} , \mathcal{I} , \mathcal{R} , \mathcal{D} , β , and ν are straightforward and are not written explicitly in the algorithm.~~

The Gillespie simulation algorithm is given as section 2.1.1. The transmission tree T is simulated along with the epidemic. We keep a map called “*tip*”, which maps infected nodes in \mathcal{I} to the tips of T . The simulation continues until either there are no discordant edges left in the network, or we reach a user-defined cutoff of time (t_{\max}) or number of infections (I). We use the notation $\text{Uniform}(0, 1)$ to indicate a number drawn from a uniform distribution on $(0, 1)$, and $\text{Exponential}(\lambda)$ to indicate a number drawn from an exponential distribution with rate λ . The combined number of internal nodes and tips in T is denoted $|T|$. The updates to \mathcal{S} , \mathcal{I} , \mathcal{R} , \mathcal{D} , β , and ν described in the previous paragraph are not written explicitly in section 2.1.1, as they are quite straightforward and would only obfuscate the pseudocode.

Algorithm 6 Simulation of an epidemic and transmission tree over a contact network

```

infect a node  $v$  at random, updating  $\mathcal{S}$ ,  $\mathcal{I}$ ,  $\mathcal{D}$ ,  $\beta$ , and  $\gamma$ 
 $T \leftarrow$  a single node with label 1
 $\text{tip}[v] \leftarrow 1$ 
 $t \leftarrow 0$ 
while  $\mathcal{D} \neq \emptyset$  and  $|\mathcal{I}| + |\mathcal{R}| < I$  and  $t < t_{\max}$  do
     $s \leftarrow \min(t_{\max} - t, \text{Exponential}(\beta + \nu))$ 
    for  $v \in \text{tip}$  do
        extend the branch length of  $\text{tip}[v]$  by  $s$ 
    end for
     $t \leftarrow t + s$ 
    if  $t < t_{\max}$  then
        if  $\text{Uniform}(0, \beta + \nu) < \beta$  then
            choose an edge  $e = (u, v)$  from  $\mathcal{D}$  with probability  $\beta_e / \beta$  and infect  $v$ 
             $\text{tip}[v] \leftarrow |T| + 1$  ▷ add new tips to tree and tip array
             $\text{tip}[u] \leftarrow |T| + 2$  ▷ corresponding to  $u$  and  $v$ 
            add tips with labels  $(|T| + 1)$  and  $(|T| + 2)$  to  $T$ 
            connect the new nodes to  $\text{tip}[v]$  in  $T$ , with branch lengths 0
        else
            choose a node  $v$  from  $\mathcal{I}$  with probability  $\nu_v / \nu$  and remove  $v$ 
            delete  $v$  from  $\text{tip}$ 
        end if
        update  $\mathcal{S}$ ,  $\mathcal{I}$ ,  $\mathcal{R}$ ,  $\mathcal{D}$ ,  $\beta$ , and  $\nu$ 
    end if
end while

```

2.1.2 Phylogenetic kernel

The tree kernel developed by Poon et al. [18] provides a comprehensive similarity score between two phylogenetic trees, via the dot-product of the two trees' feature vectors in the space of all possible subset trees with branch lengths (see section 1.2.4). Because the branch lengths are continuous, there are infinitely many possible subset trees; hence, the feature space is infinite-dimensional. The kernel was implemented using the fast algorithm developed by Moschitti [76]. First, the production rule of each node, which is the total number of children and the number of leaf children, is recorded. The nodes of both trees are ordered by production rule, and a list of pairs of nodes sharing the same production rule is created. These are the nodes for which the value of the tree kernel must be computed - all other pairs have a value of zero. The pairs to be compared are then re-ordered so that the child nodes are always evaluated before their parents. Due to its recursive definition, ordering the pairs in this way allows the tree kernel to be computed by dynamic programming. The complexity of this implementation is $O(|T_1||T_2|)$ for the two trees T_1 and T_2 being compared.

The tree kernel cannot be used directly as a distance measure for ABC, since it is maximized, not minimized, when the two trees being compared are the same. Therefore, we defined the distance between two trees as

$$\rho(T_1, T_2) = 1 - \frac{K(T_1, T_2)}{\sqrt{K(T_1, T_1)K(T_2, T_2)}},$$

which is a number between 0 and 1 that is minimized when $T_1 = T_2$. This is similar to the normalization used by Poon et al. [18] and Collins and Duffy [75].

2.1.3 Adaptive sequential Monte Carlo for Approximate Bayesian computation

We implemented the adaptive SMC algorithm for ABC developed by Del Moral, Doucet, and Jasra [21]. This algorithm is similar to the reference ABC-SMC algorithm described in section 1.5.2, except that the sequence of tolerances ε_i is automatically determined rather than specified in advance. The tolerances are chosen such that the ESS of the particle population, which indicates the quality of the Monte Carlo approximation (see section 1.4.2), decays at a controlled rate. A sudden precipitous drop in ESS would indicate that only a small number of particles had non-zero importance weights, which would result in a very poor Monte Carlo approximation to the target distribution. This situation is referred to as the collapse of the approximation or particle degeneracy (see section 1.4.3) and is mitigated by the adaptive approach. A single parameter ~~α (not to be confused with the BA model parameter)~~ controls the decay rate. In the original paper of Del Moral, Doucet, and Jasra [21], the parameter is called α , but to avoid confusion with the BA parameter of the same name we will refer to it here as α_{ESS} . The tolerance ε_i is chosen to satisfy

$$\text{ESS}(w_i) = \alpha_{\text{ESS}} \text{ESS}(w_{i-1}),$$

where, w_i is the vector of weights at the i th step. Note that, since w_i depends on ε_i , this equation solves for the updated weights and the updated tolerance simultaneously. As pointed out by Del Moral, Doucet, and Jasra [21], the equation has no analytic solution, but can be solved numerically by bisection. The forward kernels

K_i are taken to be MCMC kernels with stationary distributions π_{ϵ_i} and proposal distributions

$$q_i(\theta^{(k)}, \theta^{(k)'}) \prod_{j=1}^M f(z^{(j,k)' | \theta^{(k)'}}),$$

where $\theta^{(k)}$ is the vector of model parameters [associated with particle \$x^{\(k\)}\$](#) and $z^{(j,k)'}$, $1 \leq j \leq M$, are M datasets simulated according to $\theta^{(k)'}$. In our implementation, the q_i [are constructed component-wise for \$\theta\$ out of](#) Gaussian proposals for continuous parameters and Poisson proposals for discrete parameters. For the Poisson proposals, the number of [discrete](#) steps to move the particle is drawn from a Poisson distribution, and the direction in which to move the particle is chosen uniformly at random. The variance of each proposal distribution was set equal to twice the empirical variance of the particles, following [21, 135]. The backwards kernels are

$$L_{i-1}(x', x) = \frac{\pi_i(x) K_i(x, x')}{\pi_i(x')}.$$

Here we have written x' for x_i and x for x_{i-1} to emphasize that x_{i-1} is the current value of the particle and x_i is the proposed value. When substituted into eq. (1.9), the forward kernels $K_i(x, x')$ and densities $\pi_i(x') = \pi_{\epsilon_i}(x')$ cancel out, and we are left with the following weight update.

$$\begin{aligned} w_i(x) &\propto w_{i-1}(x) \frac{\pi_i(x') L_{i-1}(x', x)}{\pi_{i-1}(x) K_i(x, x')} && \text{importance weight update from SMC} \\ &= w_{i-1}(x) \frac{\pi_i(x') \pi_i(x) K_i(x, x')}{\pi_{i-1}(x) K_i(x, x') \pi_i(x')} && \text{substitute } L_{i-1} \\ &= w_{i-1}(x) \frac{\pi_i(x)}{\pi_{i-1}(x)} && \text{cancel } K_i(x, x') \text{ and } \pi_i(x') \\ &= w_{i-1}(x) \frac{\pi(\theta) \prod_{j=1}^M f(z^{(j)' | \theta}) \sum_{j=i}^M \mathbb{I}_{A_{\epsilon_i, y}}(z^{(j)})}{\pi(\theta) \prod_{j=1}^M f(z^{(j)' | \theta}) \sum_{j=i}^M \mathbb{I}_{A_{\epsilon_{i-1}, y}}(z^{(j)})} && \text{definition of } \pi_i(x) \\ &= w_{i-1}(x) \frac{\sum_{j=i}^M \mathbb{I}_{A_{\epsilon_i, y}}(z^{(j)})}{\sum_{j=i}^M \mathbb{I}_{A_{\epsilon_{i-1}, y}}(z^{(j)})} && \text{cancel prior and likelihood.} \end{aligned}$$

In other words, when the distance threshold ϵ_{i-1} is contracted to ϵ_i , the particles' weights are multiplied by the proportion of simulated datasets that are still inside the new threshold. The user may specify a final tolerance ϵ , or a final acceptance rate of the MCMC kernel, and the algorithm will be stopped when either of these termination conditions is reached. The latter condition stops the algorithm when the particles are not moving around very much, implying little change in the estimated target.

2.1.4 Justification for approach

[We present here a non-rigorous justification for the use of ABC for the problem at hand, as opposed to more frequently-used approaches for fitting mathematical models \(see ??\). Consider a contact network model with parameters \$\theta\$, and an estimated transmission tree \$T\$. Taking a Bayesian approach, our aim](#)

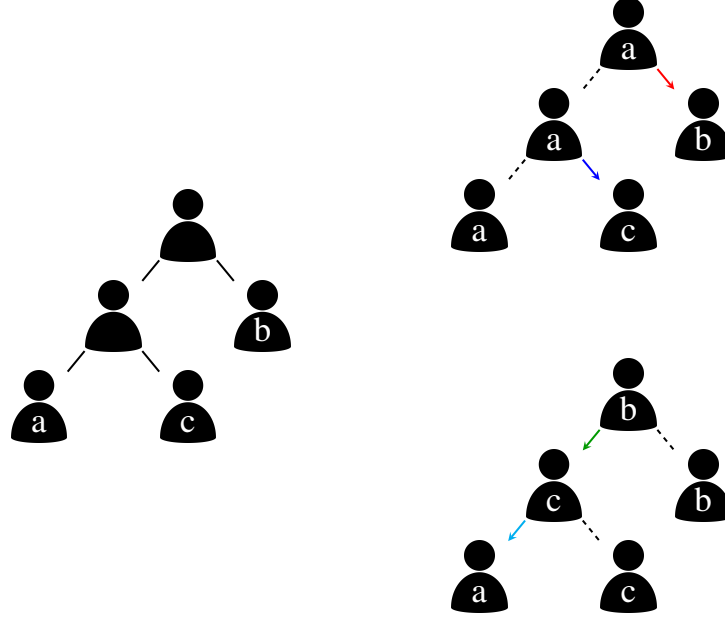


Figure 2.2: Illustration of an estimated transmission tree without labels (left) and two possible underlying complete transmission trees with labels (right). In the top right scenario, the epidemic began with node a who transmitted first to b and then to c . In the bottom right scenario, b was the index case; b infected c , who went on to infect a . A transmission tree estimated from a viral phylogeny would have the same topology and tip labels in both cases.

is to obtain a sample from the posterior distribution on the model's parameters given our data,

$$\pi(\theta \mid T) = \frac{f(T \mid \theta)\pi(\theta)}{\int_{\Theta} f(T \mid \theta)\pi(\theta)d\theta}.$$

For all but the simplest models, the normalizing constant in the denominator is an intractable integral. What we shall argue here is that, in contrast to most commonly studied mathematical models, the likelihood $f(T \mid \theta)$ is also likely to be intractable in our case.

As discussed in section 1.2.2, the internal nodes of transmission trees represent transmission events, and are labelled with the donor in the associated transmission pair. However, when we estimate a transmission tree from viral sequence data, we generally only know the labels of the tips of the tree, not the labels of the internal nodes. In viral phylogenies, the transmissions are at least partially preserved through the evolutionary relationships among the viruses, but the directionality of those transmissions is unknown. Thus, a single estimated transmission tree can correspond to many possible pathways of the epidemic through the network. Figure 2.2 illustrates this concept for a simple transmission tree with three tips. When calculating a likelihood given a transmission tree, we must sum over all possible labellings of the internal nodes. Let \mathcal{L} be the set of such labellings. Then

$$f(T \mid \theta) = \sum_{l \in \mathcal{L}} f(T, l \mid \theta). \quad (2.1)$$

A contact network model assigns a probability density to each possible contact network. Transmission trees are realized over particular contact networks, not over the model itself. Therefore, we must also sum over all contact networks which could be generated by the model. Let \mathcal{G} be the set of all possible contact networks. Summing eq. (2.1) over \mathcal{G} gives

$$f(T | \theta) = \sum_{G \in \mathcal{G}} \sum_{l \in \mathcal{L}} f(T, l | G, \theta) f(G | \theta). \quad (2.2)$$

This can be simplified somewhat by noticing that, given a specific contact network, the labelled transmission tree depends only on that network and not on the model that generated it. That is, $f(T, v | G, \theta) = f(T, l | G)$, and

$$f(T | \theta) = \sum_{G \in \mathcal{G}} f(G | \theta) \sum_{l \in \mathcal{L}} f(T, v | G) \quad (2.3)$$

Under the assumption that both transmission and removal are Poisson processes, calculating $f(T, l | G)$ can be accomplished by a straightforward modification of the Gillespie simulation algorithm (section 2.1.1). Rather than choosing transmission or removal events according to their probabilities, the events would be deterministically chosen based on the transmission tree and the probabilities of each event would be multiplied together. Assuming efficient data structures for storing lists of nodes and edges, the complexity of this calculation would be $O(|T|)$. The number of possible labellings of internal nodes of T is easily seen to be $2^{(|T|-1)/2}$ by noticing that each of the $(|T|-1)/2$ internal nodes must be labelled with the same label as either its right child or its left child. Although exponential calculations of this nature can often be simplified on trees using dynamic programming (e.g.[141]), it cannot be straightforwardly applied in this case because the subtrees' probabilities depend on the existing epidemic progress (their parents and siblings). Hence, calculating the inner sum over labels may take time $O(2^{(|T|-1)/2})$.

The outer sum, over all contact networks, is also difficult to evaluate in general. There are $2^{N(N-1)}$ directed graphs on N nodes [13]. There must be at least as many nodes in the contact network as the number of tips in the tree, which is $(|T|+1)/2$. Of course, it is very likely that there are more nodes in the network than observed tips because some individuals are never infected and/or some infected individuals are never sampled. The complexity of calculating $f(G | \theta)$ is obviously dependent on the particular model being investigated. For the BA model, we might have to sum over all possible orders in which the nodes could be added, and all assignments of edges to the nodes which generated them. However, even in the case that calculating $f(G | \theta)$ can be done in constant time, the sum (2.3) still has at least $O(2^{|T|^2})$ terms.

We have shown that both the normalizing constant $\int_{\Theta} f(T | \theta) \pi(\theta) d\theta$ and the likelihood $f(T | \theta)$ are likely computationally prohibitive to calculate. If this is the case, the problem of fitting contact network models to phylogenies seems to be of the *doubly-intractable* type [142], which would imply that these models are not amenable to neither ML nor Bayesian inference techniques. Although both methods are able to cope with an intractable normalizing constant (for example, by local search for ML

or Bayesian MCMC), neither can avoid the intractable likelihood calculations. This justifies the use of ABC, which is a likelihood-free method.

We have not proven here that eq. (2.3) is impossible to calculate in polynomial time - it could be possible to algebraically simplify the sum into a tractable expression. Furthermore, under certain models, a large proportion of \mathcal{G} may have zero probability, which would enable the simplification of the outer sum on a model-specific basis. It should also be noted that extensions of Bayesian MCMC have been developed for doubly-intractable problems [142, 143], which might be adaptable to the problem at hand. These have not been as widely used as ABC, nor are they as easily parallelizable as SMC.

2.2 Analysis of Barabási-Albert model with synthetic data

2.2.1 Motivation

We developed *netabc* with the objective of extracting useful, quantitative information about network structures from viral phylogenies.

There are three aspects to “usefulness”. The first is model specification and the biological or epidemiological interpretation of the parameters. We want the model to be somewhat realistic, so that . . . , but not too complicated that it becomes difficult to interpret. At least some of the parameters should be of interest from a theoretical or practical perspective. As this is a phylodynamic method, intended to be used with viral sequence data, we would also like to choose parameters which may be difficult to estimate with more standard methods; otherwise, our method provides no advantage. The Barabási-Albert (BA) model (section 1.3.2) satisfies these criteria, albeit some better than others. The purported realism of the model stems from the fat-tailed degree distributions it produces, which are similar to those observed in real world sexual networks [23–26, 144]. Moreover, the “rich get richer” phenomenon, where popular individuals attract new connections at an elevated rate, is intuitively reasonable for both sexual [106] and IDU [104] networks. However, the model is very simple, assuming that all nodes form the same number of links when added to the network and share the same preference for popular individuals.

In this thesis, we consider four parameters related to the BA model, denoted N , m , α , and I (see section 1.3.2). The first three of these are parameters of the model itself, while I is related to the simulation of transmission trees over the network. However, we will refer to all four as BA parameters. N denotes the total number of nodes in the network, or equivalently, susceptible individuals in the population. m is the number of new undirected edges added for each new vertex, or equivalently one-half of the average degree. α is the power of preferential attachment - new nodes are attached to existing nodes of degree d with probability proportional to $d^\alpha + 1$. Finally, I is the number of infected individuals at the time when sampling occurs.

From a public health standpoint, all four parameters are of some interest. The prevalence I can be used to estimate the resources required to combat an ongoing epidemic, while total susceptible population size N provides a similar metric for preventative measures. The average degree, controlled by m , may have some value in predicting how fast an epidemic will spread. When there are more connections in the network, the time intervals between transmission events decrease, as does the number

of transmissions required for the epidemic to travel from one part of the network to another. High average degrees may be reduced by interventions causing behaviour change, such as using condoms or caughing into one's sleeve; these types of interventions will be less helpful when the average number of contacts per person is already low. The preferential attachment power α is perhaps the least obviously useful, but it can offer some insight into the dynamics of the network and how it would react to the removal of nodes. Dombrowski et al. [104] found evidence of preferential attachment in IDU networks, and suggested that as a consequence of this characteristic, the removal of random nodes (such as through a police crackdown) might inadvertently make it easier for epidemics to spread. When individuals with only one or two connections lose them, they might tend to seek out well-known (that is, high-degree) members of the community, thus increasing those individuals' connectivity even further.

All four BA parameters can be estimated using more conventional approaches, but the difficulty of doing so varies between the parameters and based on other characteristics of the epidemic under consideration. All network parameters can theoretically be estimated by explicitly constructing the contact network, although this is highly resource intensive and is hampered by misreporting and other challenges [43]. In any event, if we are in a position to know the true contact network then the point of estimating network parameters is moot. Depending on the context, N may be estimable with a survey, but if the population in question is "hidden" (such as IDU in regions where drug use is illegal) estimation of numbers may be more challenging. Respondent-driven sampling [91] was developed to investigate such populations. The prevalence I may be estimated from reports of hospital admissions or engagement in care, but this may be very difficult for diseases such as hepatitis C virus (HCV) which may remain asymptomatic for decades. Routine population-level testing is another possibility, but is resource intensive and the barriers associated with hidden populations still apply. The average degree of the network ($2m$) can be ascertained by simply asking a sample of people how many contacts they had, but they might be unwilling to disclose this information or, in the case of airborne transmissions, they might not know. The estimation of α is even more complex, as this parameter is most strongly reflected in the connectivity of very high degree nodes, who are rare in the population. Even if full degree distributions are available, there are models other than preferential attachment which can produce scale-free networks [*e.g.* 145]. de Blasio, Svensson, and Liljeros [106] were able to estimate α by maximum likelihood using partner count data from several sequential time intervals, but they admit such detailed data are not usually available. Moreover, their dataset was constructed via a random survey, which would likely miss the few high-degree nodes characterizing a power law degree distribution. In short, each of these parameters may be estimated without phylodynamics, but there are sufficient difficulties that we believe an alternative method based on sequence data is warranted.