

outbreakerR

disease outbreak reconstruction using genetic data

Thibaut Jombart

MRC Centre for Outbreak Analysis and Modelling – Imperial College London

GenEpi — LSHTM, London
10-04-2013

Investigating disease outbreaks using genetic data

Background

- here, “*outbreak*” = small, localised epidemic
- **small-scale** → **dense sampling** possible
- **pathogen genomes sequences increasingly available**



Investigating disease outbreaks using genetic data

Background

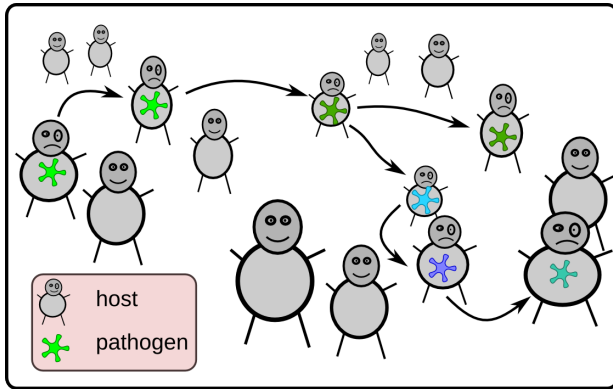
- here, “*outbreak*” = small, localised epidemic
- **small-scale** → **dense sampling** possible
- **pathogen genomes sequences increasingly available**



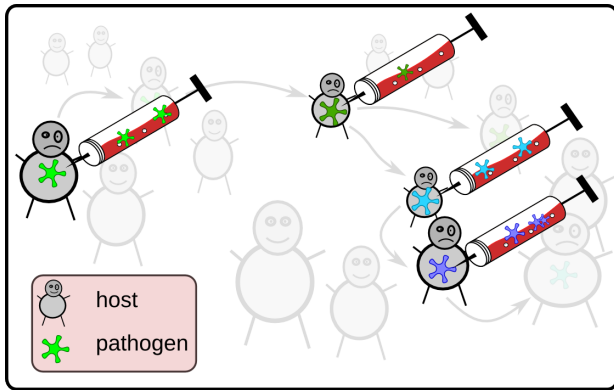
Objectives

- exploit genetic information to reconstruct outbreaks
- infer **transmission trees**, dates of infection, R , ...
- design a tool for **retrospective** or **near real-time** analysis

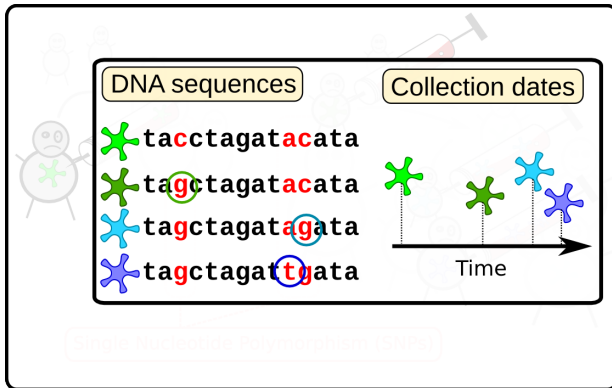
The data



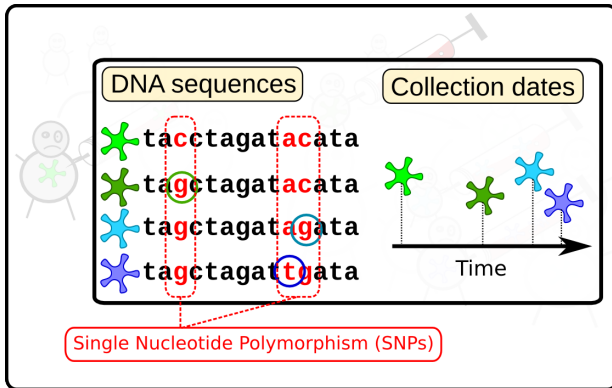
The data



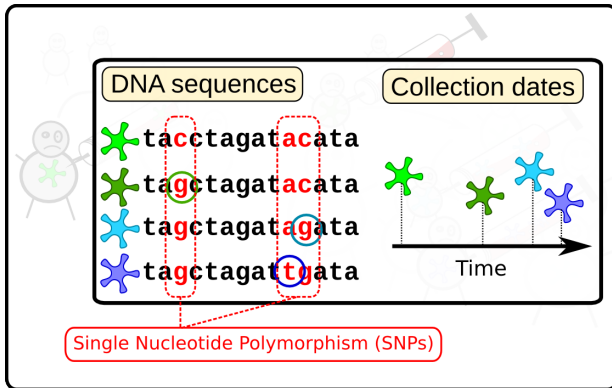
The data



The data



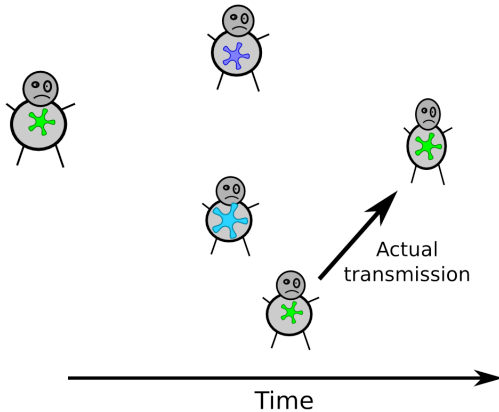
The data



Generic model integrating **genomic data** and **collection dates**

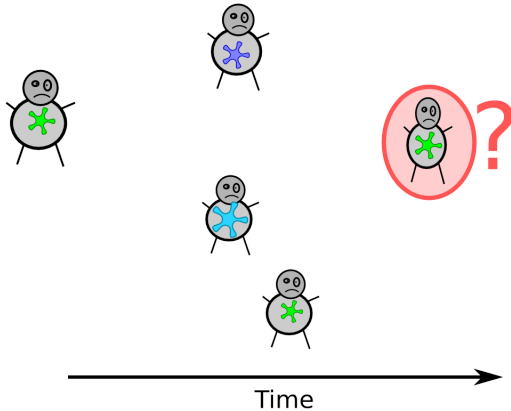
Rationale of the method

How to infer the infector of a given individual?



Rationale of the method

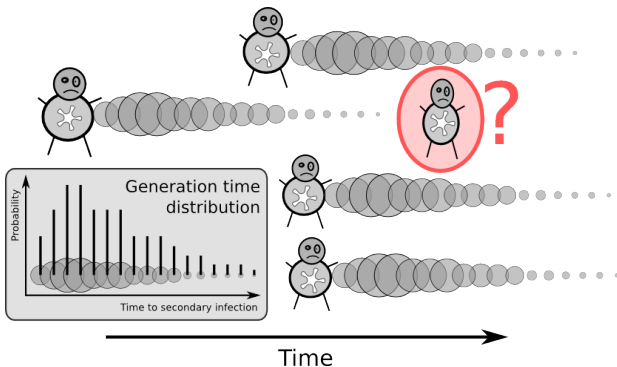
How to infer the infector of a given individual?



Rationale of the method

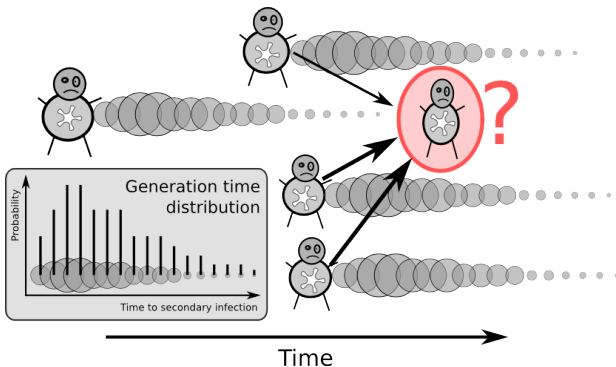
Approach based on generation time distribution

(Ferguson *et al.* 2001, Nature; Wallinga & Teunis 2004, Am. J. Epidemiol.)



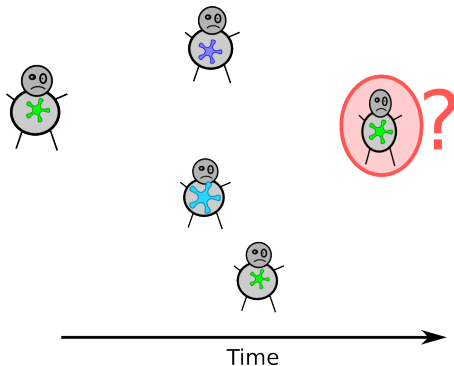
Rationale of the method

Approach based on generation time distribution
(Wallinga & Teunis 2004, Am. J. Epidemiol.)



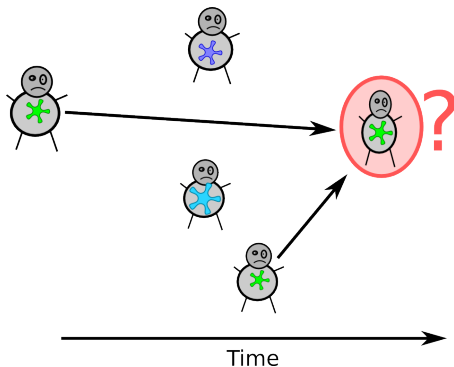
Rationale of the method

Approach based on genetic data (*SeqTrack*)
(Jombart *et al.* 2010, *Heredity*.)



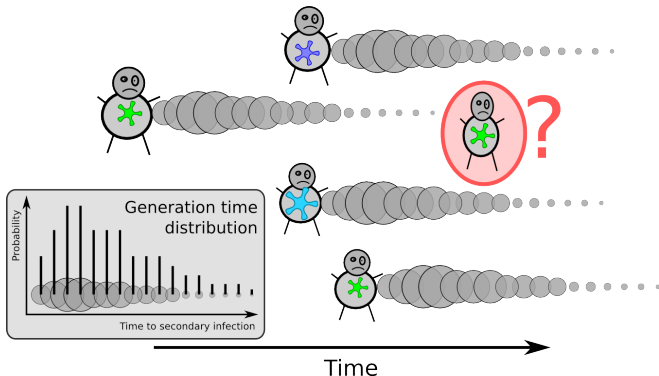
Rationale of the method

Approach based on genetic data (*SeqTrack*)
(Jombart *et al.* 2010, *Heredity*.)



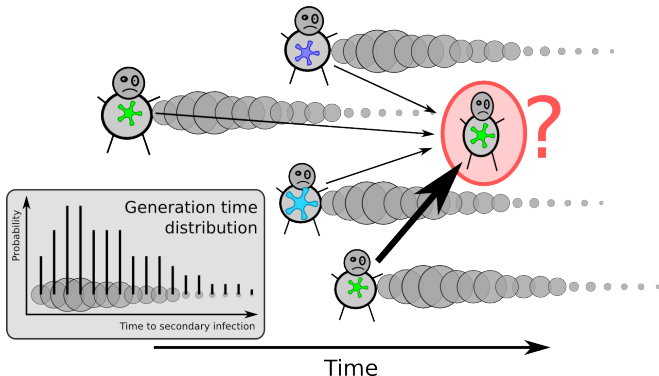
Rationale of the method

outbreaker: use generation time distribution and genetic data:



Rationale of the method

outbreaker: use generation time distribution and genetic data:



Outline of the model

Likelihood

- branching process:

$$p(\text{transmission tree}) = \prod_{\text{all branches}} p(\text{branches})$$

- $p(\text{branch}) =$
 $p(\text{infection/collection dates}) \times p(\text{genetic differences})$

Outline of the model

Likelihood

- branching process:


$$p(\text{transmission tree}) = \prod_{\text{all branches}} p(\text{branches})$$

- $p(\text{branch}) =$
 $p(\text{infection/collection dates}) \times p(\text{genetic differences})$

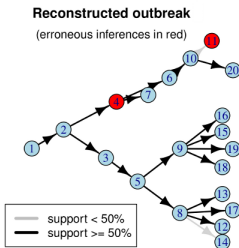
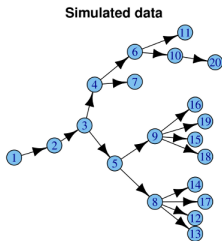
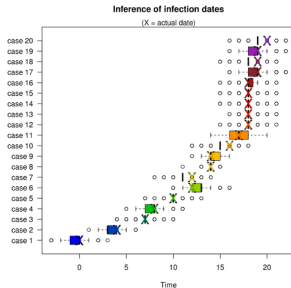
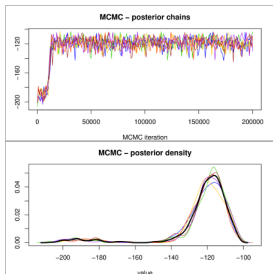
Implementation

- Bayesian framework
- augmented data for ancestries and unobserved cases
- MCMC (Metropolis-Hasting) for sampling from posterior

The package outbreakeR

- C implementation embedded within  package
- multi-platform: linux, MacOS X, Windows, Solaris, . . .
- supports parallelization
- post-processing of MCMC, simulations, graphics
- tested on wide range of simulations

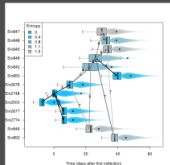
outbreakerR : some examples



Looking ahead: applications and developments

Applications

SARS 2003 Singapore outbreak



S. pneumoniae carriage Maeda camp, Thailand



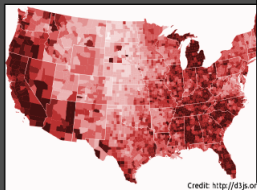
Model improvement

Spatial information, contact networks

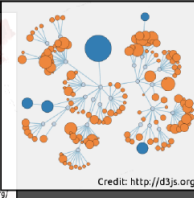


New visualization tools

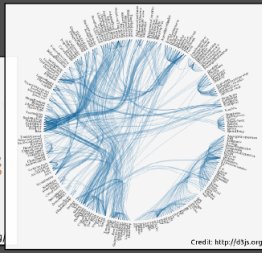
Web-based interactive graphics



Credit: <http://d3js.org/>



Credit: <http://d3js.org/>



Credit: <http://d3js.org/>

Looking ahead: the bigger picture

Applications

SARS 2003 Singapore outbreak



S. pneumoniae carriage
Maela camp, Thailand



Model improvement

Spatial information, contact networks



EpiEstim-package (EpiEstim)

The EpiEstim package

Version: 1.1.0

Description


Quantifying transmissibility during various infectious diseases is a challenge for public health researchers. The EpiEstim package provides a set of tools to estimate the effective reproduction number, R_{eff} , from case data. It is a full-scale model-based approach to estimate R_{eff} from case data. It is a full-scale model-based approach to estimate R_{eff} from case data. It is a full-scale model-based approach to estimate R_{eff} from case data.

About outbreakR

outbreakR is a collection of R packages that provide a framework for the analysis of infectious disease outbreaks. It is a full-scale model-based approach to estimate R_{eff} from case data. It is a full-scale model-based approach to estimate R_{eff} from case data. It is a full-scale model-based approach to estimate R_{eff} from case data.

New

Web-based analysis of disease outbreaks




outbreakerR

R Hackout:
a hackathon for the analysis of disease outbreaks in R

MRC Centre for Infectious Diseases and Modelling | Imperial College London | sourceforge

epibase

epibase is a collection of R packages that provide a framework for the analysis of infectious disease outbreaks. It is a full-scale model-based approach to estimate R_{eff} from case data. It is a full-scale model-based approach to estimate R_{eff} from case data. It is a full-scale model-based approach to estimate R_{eff} from case data.




...

...

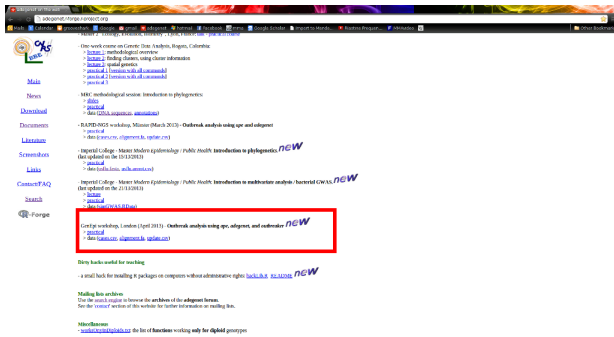
...

Acknowledgements

- **Organisers:** Anton Camacho & Marc Baguelin
- **Imperial College London:** Anne Cori, Simon Cauchemez, Xavier Didelot, Christophe Fraser, Neil Ferguson, David Aanensen
- **Wellcome Trust Sanger Institute:** Claire Chewapreecha, Sephen Bentley
- **Funding:** MIDAS
- **Thanks for your attention.**

Getting your hands dirty

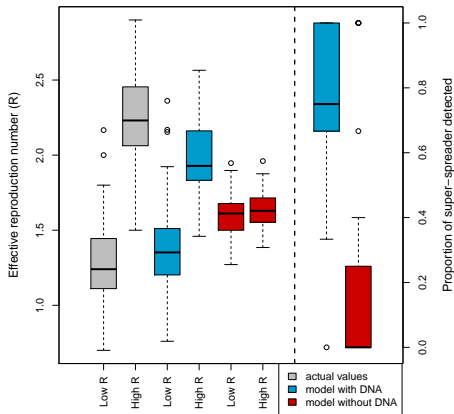
Google “adegenet” → “adegenet on the web” → “Documents”
(<http://adegenet.r-forge.r-project.org/>)



Genetic data reveal heterogeneous infectivity

Simulations with structured infectivity

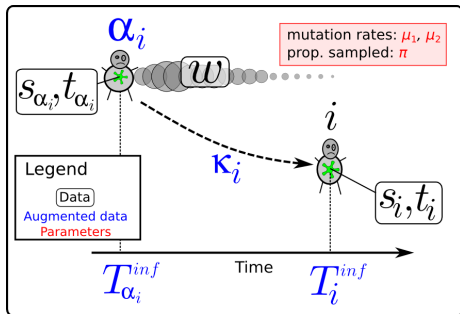
- left: 2 groups of hosts, low/high infectivity ($R_0 = 1.5/3$)
- right: rare super-spreaders ($R_0 = 1.5/20$)



outbreakerR : model notations

Data

- i : index of cases ($i = 1, \dots, n$)
- s_i : genetic sequence of i
- t_i : collection date for s_i
- w : generation time distribution



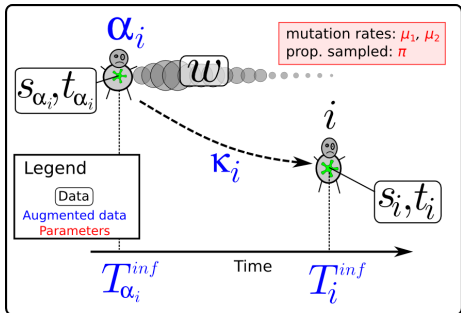
outbreakerR : model notations

Data

- i : index of cases ($i = 1, \dots, n$)
- s_i : genetic sequence of i
- t_i : collection date for s_i
- w : generation time distribution

Augmented data

- T_i^{inf} : date of infection of i
- α_i : infector of i
- κ_i : number of generations between i and α_i



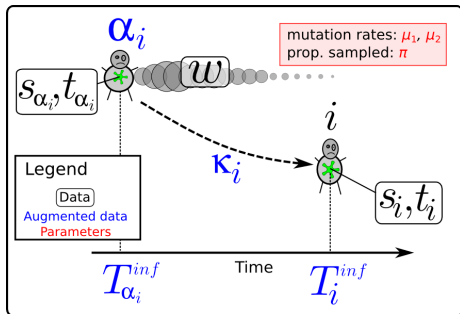
outbreakerR : model notations

Data

- i : index of cases ($i = 1, \dots, n$)
- s_i : genetic sequence of i
- t_i : collection date for s_i
- w : generation time distribution

Augmented data

- T_i^{inf} : date of infection of i
- α_i : infector of i
- κ_i : number of generations between i and α_i



Parameters

- μ_1, μ_2 : rates of transitions and transversions ($\mu_2 = \gamma \times \mu_1$)
- π : proportion of the outbreak sampled

outbreakerR : model definition

- Posterior proportional to joint distribution:

$$\begin{aligned} & p(\{s_i, t_i, T_i^{inf}\}_{(i=1,\dots,n)}, \alpha, \kappa, w, \mu_1, \gamma, \pi) \\ = & \underbrace{p(\{s_i, t_i, T_i^{inf}, \alpha_i, \kappa_i\}_{(i=1,\dots,n)} | w, \mu_1, \gamma, \pi)}_{likelihood} \times \underbrace{p(w, \mu_1, \gamma, \pi)}_{prior} \end{aligned}$$

outbreakerR : model definition

- Posterior proportional to joint distribution:

$$\begin{aligned} & p(\{s_i, t_i, T_i^{inf}\}_{(i=1,\dots,n)}, \alpha, \kappa, w, \mu_1, \gamma, \pi) \\ = & \underbrace{p(\{s_i, t_i, T_i^{inf}, \alpha_i, \kappa_i\}_{(i=1,\dots,n)} | w, \mu_1, \gamma, \pi)}_{\text{likelihood}} \times \underbrace{p(w, \mu_1, \gamma, \pi)}_{\text{prior}} \end{aligned}$$

- Likelihood of case i decomposed as:

$$\underbrace{p(s_i | \alpha_i, s_{\alpha_i}, \kappa_i, \mu_1, \gamma)}_{\text{genetic}} \times \underbrace{p(t_i | T_i^{inf}, w) p(T_i^{inf} | \alpha_i, T_{\alpha_i}^{inf}, \kappa_i, w) p(\kappa_i | \pi)}_{\text{epidemiological}}$$

outbreakerR : model definition

- Posterior proportional to joint distribution:

$$\begin{aligned} & p(\{s_i, t_i, T_i^{inf}\}_{(i=1,\dots,n)}, \alpha, \kappa, w, \mu_1, \gamma, \pi) \\ = & \underbrace{p(\{s_i, t_i, T_i^{inf}, \alpha_i, \kappa_i\}_{(i=1,\dots,n)} | w, \mu_1, \gamma, \pi)}_{\text{likelihood}} \times \underbrace{p(w, \mu_1, \gamma, \pi)}_{\text{prior}} \end{aligned}$$

- Likelihood of case i decomposed as:

$$\underbrace{p(s_i | \alpha_i, s_{\alpha_i}, \kappa_i, \mu_1, \gamma)}_{\text{genetic}} \times \underbrace{p(t_i | T_i^{inf}, w) p(T_i^{inf} | \alpha_i, T_{\alpha_i}^{inf}, \kappa_i, w) p(\kappa_i | \pi)}_{\text{epidemiological}}$$

- Sampling from posterior distribution using MCMC (Metropolis-Hasting)

Model - detail of likelihoods

- Genetic likelihood:

$$\underbrace{\mathcal{B}(d(s_i, s_{\alpha_i}) | l(s_i, s_{\alpha_i}) \kappa_i, \mu_1)}_{\text{transitions}} \times \underbrace{\mathcal{B}(g(s_i, s_{\alpha_i}) | l(s_i, s_{\alpha_i}) \kappa_i, \gamma \mu_1)}_{\text{transversions}}$$

Model - detail of likelihoods

- Genetic likelihood:

$$\underbrace{\mathcal{B}(d(s_i, s_{\alpha_i}) | l(s_i, s_{\alpha_i}) \kappa_i, \mu_1)}_{\text{transitions}} \times \underbrace{\mathcal{B}(g(s_i, s_{\alpha_i}) | l(s_i, s_{\alpha_i}) \kappa_i, \gamma \mu_1)}_{\text{transversions}}$$

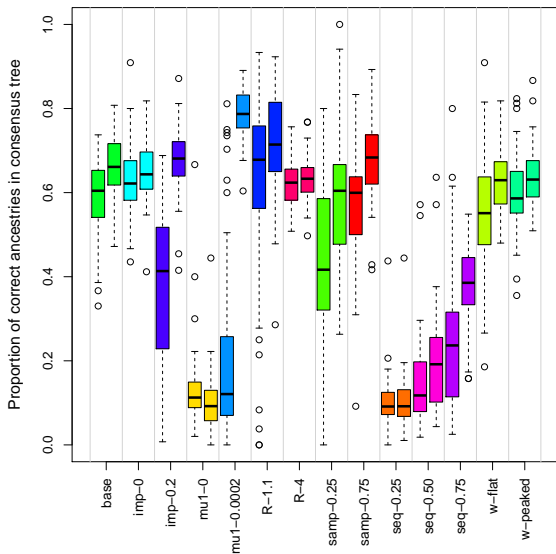
- Epidemiological likelihood:

$$p(t_i | T_i^{inf}, w) \times p(T_i^{inf} | \alpha_i, T_{\alpha_i}^{inf}, \kappa_i, w) \times p(\kappa_i | \pi) \\ w(t_i - T_i^{inf}) \times w^{(\kappa_i)}(T_i^{inf} - T_{\alpha_i}^{inf}) \times f_{\mathcal{NB}}(1 | \kappa_i - 1, \pi)$$

with:

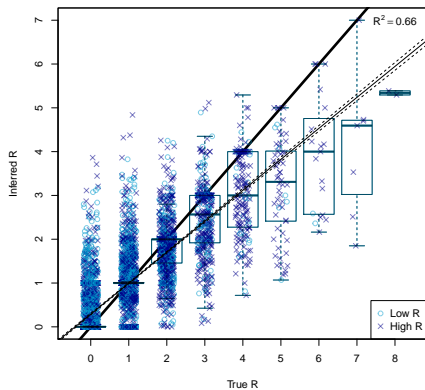
- $w(t_i - T_i^{inf})$: probability of sampling date (assumed prop. to infectiousness)
- $w^{(\kappa_i)}(T_i^{inf} - T_{\alpha_i}^{inf})$: probability of infection date ($w^{(\kappa_i)}$ denotes κ_i convolutions of w)
- $f_{\mathcal{NB}}(1 | \kappa_i - 1, \pi)$: probability of $\kappa_i - 1$ unsampled cases ($f_{\mathcal{NB}}$: density of the Negative Binomial distribution)

Simulation results: inference of correct ancestries

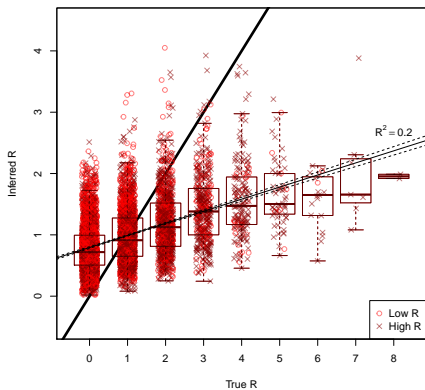


Inferring heterogeneous R : 2-groups simulations

With genetic data:

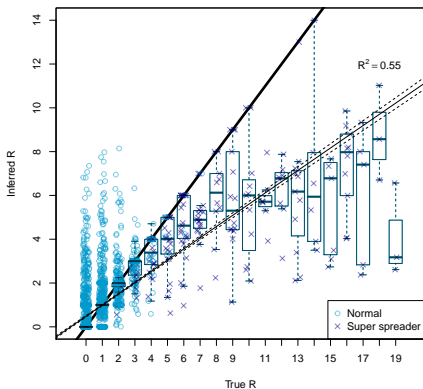


Without genetic data:



Inferring heterogeneous R : super-spreader simulations

With genetic data:



Without genetic data:

