Cesky Krumlov

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

Bayesian coalescent inference of population size history

Alexei Drummond University of Auckland

Workshop on Population and Speciation Genomics Cesky Krumlov, 2016

1st February 2016

Population inference methods using BEAST2 we will cover

- 1. Extended Bayesian skyline plot (EBSP)
- 2. Structured coalescent (MultiTypeTree)
- 3. SNP-based species coalescent inference (SNAPP)

Not being covered

- Divergence time estimation
- relaxed molecular clocks
- multi-locus species coalescent inference (*BEAST)
- Epidemiological phylodynamic models
- Fossilized Birth-death model, sampled ancestors
- Programming new models in BEAST

Cesky Krumlov

Coalescent theory

Data: a small genetic sample from a large background population.

The coalescent

- is a model of the ancestral relationships of a sample of individuals taken from a larger population.
- describes a probability distribution on ancestral genealogies (trees) given a population history, N(t).
 - Therefore the coalescent can convert information from ancestral genealogies into information about population history and vice versa.
- a model of ancestral genealogies, not sequences, and its simplest form assumes neutral evolution.
- can be thought of as a prior on the tree, in a Bayesian setting.

Cesky Krumlov

Cesky Krumlov

Theoretical population genetics

Most of theoretical population genetics is based on the idealized Wright-Fisher model of population which assumes

- Constant population size N
- Discrete generations
- Complete mixing

For the purposes of this presentation the population will be assumed to be haploid, as is the case for many pathogens.



Kingman's n-coalescent

Consider tracing the ancestry of a sample of k individuals from the present, back into the past.

This process eventually *coalesces* to a single common ancestor (*concestor*) of the sample of individuals.

Kingman's n-coalescent describes the statistical properties of such an ancestry when k is small compared to the total population size N.



Cesky Krumlov

The coalescence of two ancestral lineages

- First, consider two random members from a population of fixed size N.
- ► By perfect mixing, the probability they share a *concestor* in the previous generation is 1/N.
- The probability the concestor is t generations back is

$$\Pr\{t\} = \frac{1}{N}(1 - \frac{1}{N})^{t-1}.$$

► It follows that g = t - 1, has a geometric distribution with a success rate of $\lambda = 1/N$, and so has mean N and variance of $N^3/(N-1)$.



The coalescence of k lineages

With k lineages the time to the first coalescence is derived in the same way, only now there are $\binom{k}{2}$ possible pairs that may coalesce, resulting in a success rate of $\lambda = \binom{k}{2}/N$ and mean time to first coalescence (t_k) of

$$\mathsf{E}[t_k] = \frac{\mathsf{N}}{\binom{k}{2}}$$

This implicitly assumes that N is much larger than $O(k^2)$, so that the probability of two coalescent events in the same generation is small.



The coalescent likelihood for a genealogy

For a genealogy with known coalescent times $\mathbf{t}=\{t_2,t_3,...,t_n\}$ we can write the likelihood:

$$f(\mathbf{t}|\mathbf{N}) = \frac{1}{\mathbf{N}^{n-1}} \prod_{k=2}^{n} \exp\left(-\frac{\binom{k}{2}\mathbf{t}_{k}}{\mathbf{N}}\right) \,.$$



Cesky Krumlov

The coalescent with serial samples

Many epidemiological agents, like RNA viruses, evolve very rapidly, so that the effect of sampling the population at different times becomes important. Ancient DNA also requires care handling of sampling times.



Constant size

Exponential growth

Cesky Krumlov

Bayesian integration of uncertainty in genealogies



How similar are these two trees? Both of them are plausible given the data. We can use Bayesian Markov-chain Monte Carlo to average the coalescent over all plausible trees. Cesky Krumlov

Hepatitis C in Egypt

Hepatitis C in Egypt

A case study of the coalescent approach to molecular epidemiology

Cesky Krumlov

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot



Hepatitis C (HCV)

- Identified in 1989
- ▶ 9.6kb single-stranded RNA genome
- Polyprotein cleaved by proteases
- Tissue culture system only recently developed



How important is Hepatitis C?

> 185 million people infected worldwide



- ► ~80% infections are chronic
- Liver cirrhosis and cancer risk
- ▶ 10,000 deaths per year in USA
- No protective immunity?

Cesky Krumlov

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

HCV Transmission

By percutaneous exposure to infected blood

- Blood transfusion / blood products
- Injecting and nasal drug use
- Sexual and vertical transmission
- Unsafe injections
- Unidentified routes



Coalescent population inference of HCV

Pybus et al (2003) Molecular Biology and Evolution

Egyptian HCV gene sequences n=61 E1 gene, 411bp

- All sequence contemporaneous
- Egypt has highest prevalence of HCV worldwide (10-20%)
- But low prevalence in neighbouring states
- Why is Egypt so seriously affected?
- Parenteral antischistosomal therapy (PAT)?



Coalescent theory Hepatitis C in Egypt

Cesky Krumlov

Bayesian skyline plot

Demographic model for Hepatitis C in Egypt

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

- The coalescent can be extended to model any integrable function of varying population size.
- The model we used was a const-exp-const model.
- A Bayesian MCMC method was developed to sample the gene genealogy, the substitution model and demographic function simultaneously.



Estimated population history of HCV in Egypt

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot



Uncertainty in parameter estimates

Mutational parameters Demographic parameters 0.0020 5 Posterior probability density 0.0016 per year ji e 0.0012 substitutions per 0.0008 0.0004 cond codon position 0 0.0000 0 0.1 0.2 03 0.7 0.8 0.0 500000 1500000 2500000 3500000 4500000 5500000 Exponential growth rate, r MCMC state

Growth rate of the growth phase Grey box is the prior Rates at different codon positions, All significantly different Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

Cesky Krumlov

Full Bayesian Estimation



- Marginalized over uncertainty in genealogy and mutational processes
- Yellow band represents time over which PAT was employed in Egypt

Bayesian skyline plot

Virus Phylodynamics



Cesky Krumlov

"Skyline" coalescent model



The generalized skyline plot - simulated data

$$\label{eq:constant_population_size} \begin{split} \text{Constant population size,} & N(t) = N_0 \end{split}$$

Exponential growth, $N(t) = N_0 e^{-rt}$



Cesky Krumlov

The generalized skyline plot - HIV-1 group M

The tree used here was estimated in Yusim *et al* (2001) *Phil. Trans. Roy. Soc. Lond. B* **356**:855-866. The black curve is a parametric coalescent estimate obtained from the same data under the expansion model, $N(t) = (N_0 - N_A)e^{-rt} + N_A$



Cesky Krumlov

The Bayesian skyline plot

Drummond et al (2005) Molecular Biology and Evolution

Cesky Krumlov

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

The Bayesian skyline plot estimates a demographic function that has a certain fixed number of steps (in this example 15) and then integrates over all possible positions of the break points, and population sizes within each epoch.



Dengue-4 Bayesian skyline plot (15 epochs)

Validating the Bayesian skyline plot



Coalescent theory Hepatitis C in Egypt Bavesian skyline plot

Comparison of BSP to parametric coalescent

Hepatitis C in Egypt

100000 10000 Population size (Nt) 1000 100 10 0.0 50.0 100.0 150.0 200.0 Years (before 1993)

Cesky Krumlov

Modeling complex demographic history

Cesky Krumlov

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

Dengue 4 in Puerto Rico



- $\blacktriangleright N(t) = N_0 \exp(-rt)$
 - ► log marginal likelihood = -10566.421
- N(t) =scaled-translated case data
 - ▶ log marginal likelihood = -10478.572

Comparing BSP to incidence data

Dengue 4 in Puerto Rico



Extended BSP: Stochastic Variable Selection

Heled and Drummond (2008)

Population θş t₀ t₁ t₂ t₃ t₄ t₅ Time

Cesky Krumlov

Comparison of EBSP with BSP on Egypt HCV

Cesky Krumlov



Detecting evolutionary bottlenecks using EBSP

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

480 contemporaneous samples from a single locus



Detecting evolutionary bottlenecks using EBSP

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

4000 Population (red HPD, blue target) 3000 2000 1000 400 600 800 1000 1200 Time

16 contemporaneous samples from each of 32 loci



Detecting evolutionary bottlenecks using EBSP

Coalescent theory Hepatitis C in Egypt Bayesian skyline plot

480 samples sampled through time from a single locus



The population genetic dynamics of Influenza A

Rambaut et al (2008) Nature 453:615-620



Figure 1 | Population dynamics of genetic diversity in influenza A virus.

Bayesian skyline plots of the HA and NA segments for the A/H3N2 and A/H1N1 subtypes in New York state (top) and New Zealand (bottom). The horizontal shaded blocks represent the winter seasons. The *y*-axes represent a measure of relative genetic diversity (see Methods for details). The shorter timescale of New Zealand skyline plot is due to the shorter sampling period.

Cesky Krumlov

Conclusions

- Coalescent theory provides a mathematical framework for accurately modelling small genetic samples and the stochastic outcomes of genetic drift due to random mating.
- Coalescent theory provides backward probabilities for population parameters.
- In most implementations the coalescent does not directly model small populations or stochastically fluctuating populations.
- Coalescent theory can be extended to:
 - non-parametric fluctuating populations
 - structured populations
 - gene conversion

Cesky Krumlov

BEAST book

theory / practice / programming



ETHzürich

Taming the BEAST



Bayesian evolutionary analysis by sampling trees

A summer school on Bayesian phylogenetic and phylodynamic analyses in BEAST2 with invited talks, lectures and tutorials by leading and renowned experts in the field.

Alexei Drummond University of Auckland Tracy Heath Iowa State University Oliver Pybus University of Oxford Tanja Stadler ETH Zürich Tim Vaughan University of Auckland

26 June - 1 July 2016 Engelberg, Switzerland

© E_vo

For further information and to apply: http://www.bsse.ethz.ch/cevo/taming-the-beast.html



D-BSSE Department of Biosystems Science and Engineering