

Gene Genealogies Coalescence Theory



Nature Reviews | Genetics

Annabelle Haudry Glasgow, July 2009

Gen



What could tell a gene genealogy?

- How much diversity in the population?
- Has the demographic size of the population changed? How? When?
- Is that gene under selection?
- Is/was there gene flow between sub-populations?

Could I describe the evolutionary model under which the species has been evolving?

Wright-Fisher: a simple model

Finite and constant population size N diploid individuals, 2N genes Random mating Non overlapping generations No selection, structure... Only Drift



















Backward in Time



Coalescent Events



Mathematical Theory: Kingman

Wright-Fisher Model



Coalescence probability of 2 genes at a given generation

p(2, 2N) = 1/2N

Mathematical Theory: Kingman

Wright-Fisher Model



Coalescence probability of 2 genes at a given generation

p(2, 2N) = 1/2N

Coalescence probability of 2 genes *j* generations back, within a sample of *n* genes

$$p(T_2 = j) = \frac{1}{2N} \left(1 - \frac{1}{2N} \right)^{j-1} \frac{n(n-1)}{2}$$



Mathematical Theory: Kingman

Wright-Fisher Model

Coalescence probability of 2 genes at a given generation

p(2, 2N) = 1/2N

Coalescence probability of 2 genes *j* generations back, within a sample of *n* genes



Mathematical Theory: Kingman Wright-Fisher Model Coalescence probability of 2 genes at a given generation p(2, 2N) = 1/2NCoalescence probability of 2 genes *j* generations back, within a sample of *n* genes $p(T_2 = j) = \frac{1}{2N} \left(1 - \frac{1}{2N} \right)^{j-1} \frac{n(n-1)}{2}$ *j* generations $T_2 \sim$ geometric distribution For $n \ll N$, approximation by an exponential distribution *Tn* coalescence time between 2 genes: *Tn* ~ exponential distribution: E(Tn) = 4N/(n(n-1)) so $E(T_2) = 2N$

Highly Variable Process



All trees have same probability

Coalescent with Mutations



mutations added on each branch of the tree following a *Poisson* process *mi*: $E(mi) = \mu t$ μ : mutation rate/gene/generation *t*: length of branch

How to summarize a genealogy?



How to summarize a genealogy?



How to summarize a genealogy?



 $\pi = 0.37$ S = 5 $\theta_{S} = 0.31$

Estimations de 4Nµ

- **<u>nucleotide diversity</u>** π (Tajima 1983) π = average # of pairwise differences

$$E(\pi) = \mu t = \mu . 2.2N = 4N\mu$$

- <u>theta θ_{s} </u> (Watterson 1975) S = # segregating sites

$$\Theta_{S} = \frac{S}{\sum_{i=1}^{n-1} \frac{1}{i}}$$

And a set of genealogies?

Expected distributions of summary statistics in a population under a standard model, 10,000 simulations with $\theta = 5$



And a set of genealogies?

Expected distributions of summary statistics in a population under a standard model, 10,000 simulations with $\theta = 5$



Effect of Demography

Gene genealogy affected by N variations



- Formalization of the evolutionary history of populations
- Coalescence theory will describe expected diversity in a sample under different evolution scenarii

Effect of Demography



Affects polymorphism patterns ($\neq \theta_S$ and π)



Molecular signatures of variations in population size



$$D = \frac{\pi - \theta_{\rm S}}{\sqrt{Var(\pi - \theta_{\rm S})}}$$

 $D \sim \mathcal{N}(0,1)$ under standard coalescent



$$D = \frac{\pi - \theta_{\rm S}}{\sqrt{Var(\pi - \theta_{\rm S})}}$$

 $D \sim \mathcal{N}(0,1)$ under standard coalescent



Coalescent Simulations

Hudson (2002):

Program *ms* based on coalescence theory to generate simulated gene samples.

Simplest model = Wright-Fisher population

+ complex: - changes in population size

- population structure
- migration
- recombination
- partial selfing
- selection...

Define the best model's parameter





Example of Maize Domestication

Wright et al. (2005)







Gene history ≠ Species history

More details in:

Hein, J., Schierup, M.H. & Wiuf, C. *Gene Genealogies, Variation and Evolution: A Primer in Coalescent Theory*. (Oxford University Press, USA: 2005).

