

The Coalescent

Chris Cannings,

University of Sheffield,

c.cannings@shef.ac.uk

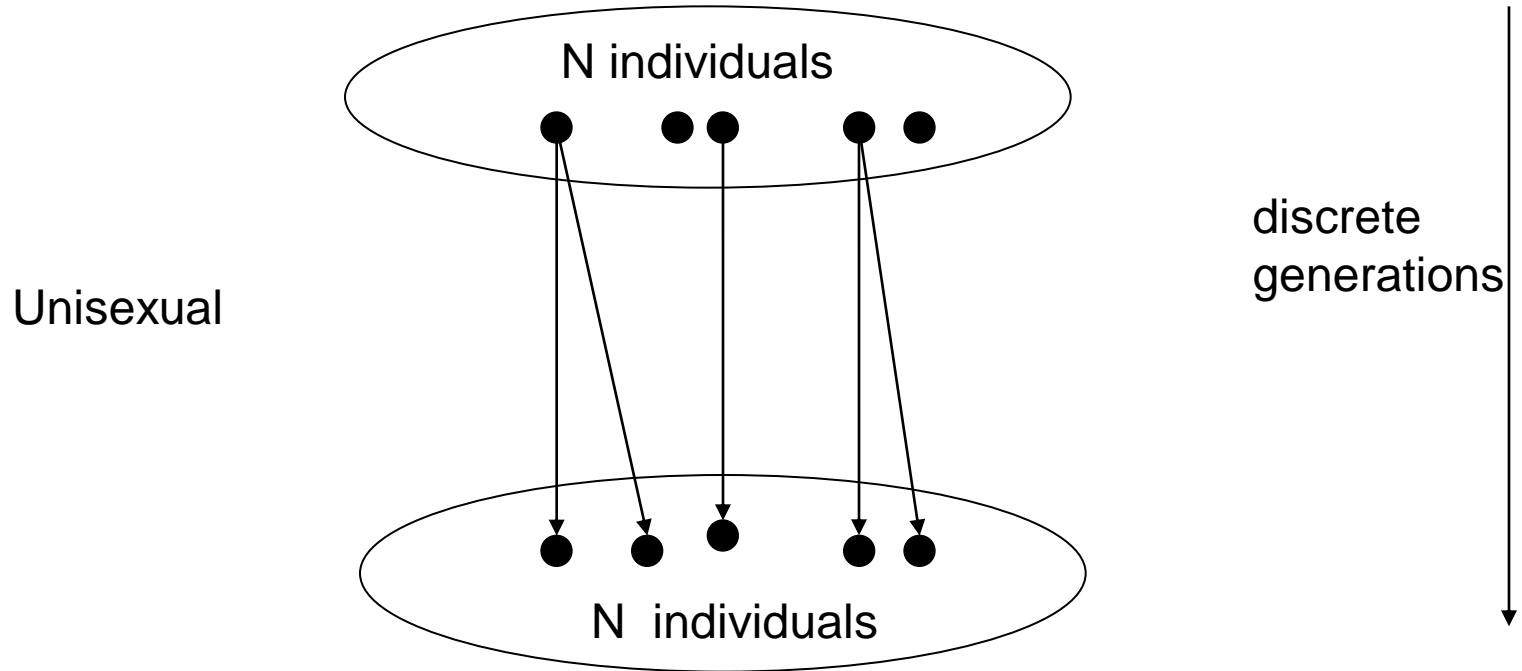
<http://www.amorph.group.shef.ac.uk/>

Singapore, 19.05.2011

The Coalescent

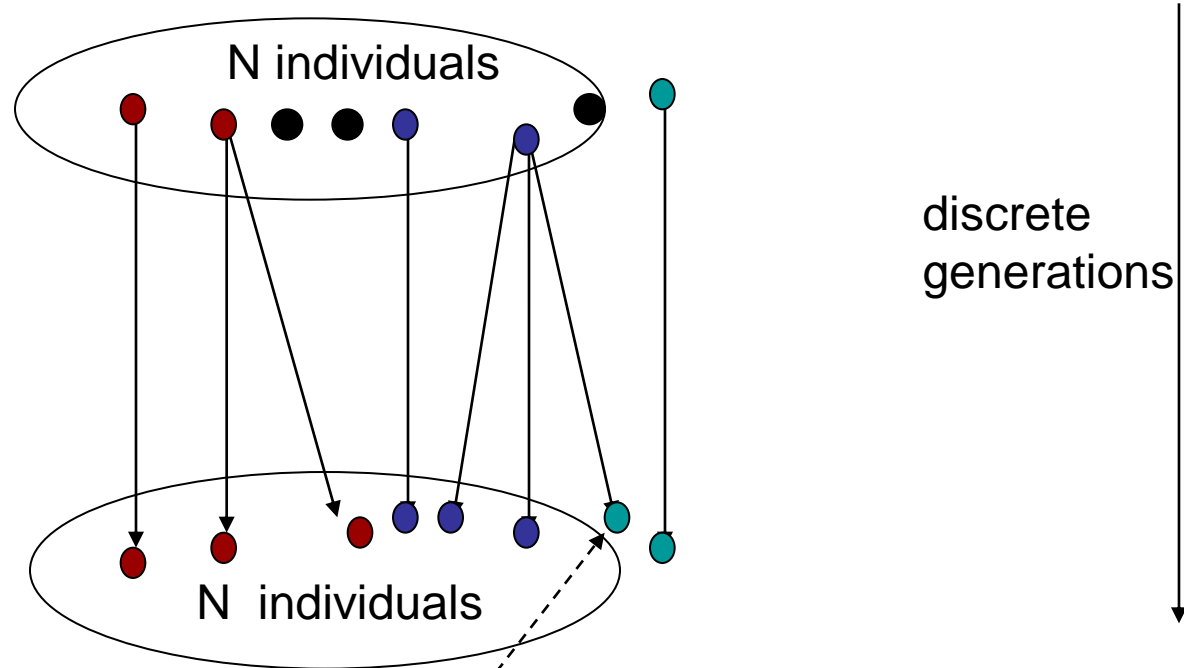
- (1) **Genetic Drift.** Stochastic model of evolution of population of neutral genes.
- (2) **The Coalescent.** Reversing time.
- (3) Tree topology. **Catalan** numbers.
- (4) **External Edges.** Distribution of lengths.
- (5) **Stepwise Mutation Model.** Distributions and moments. Bell Polynomials.

Genetic Drift



Genetic Drift

Individuals of varying types. Possibly mutation, Recombination, etc.



Mutation (offspring type different from parent type)

Wright-Fisher Model

- Each individual in the m 'th generation is the offspring of an individual in the $(m-1)$ 'th generation selected at random (prob= $1/N$) independently of all others.

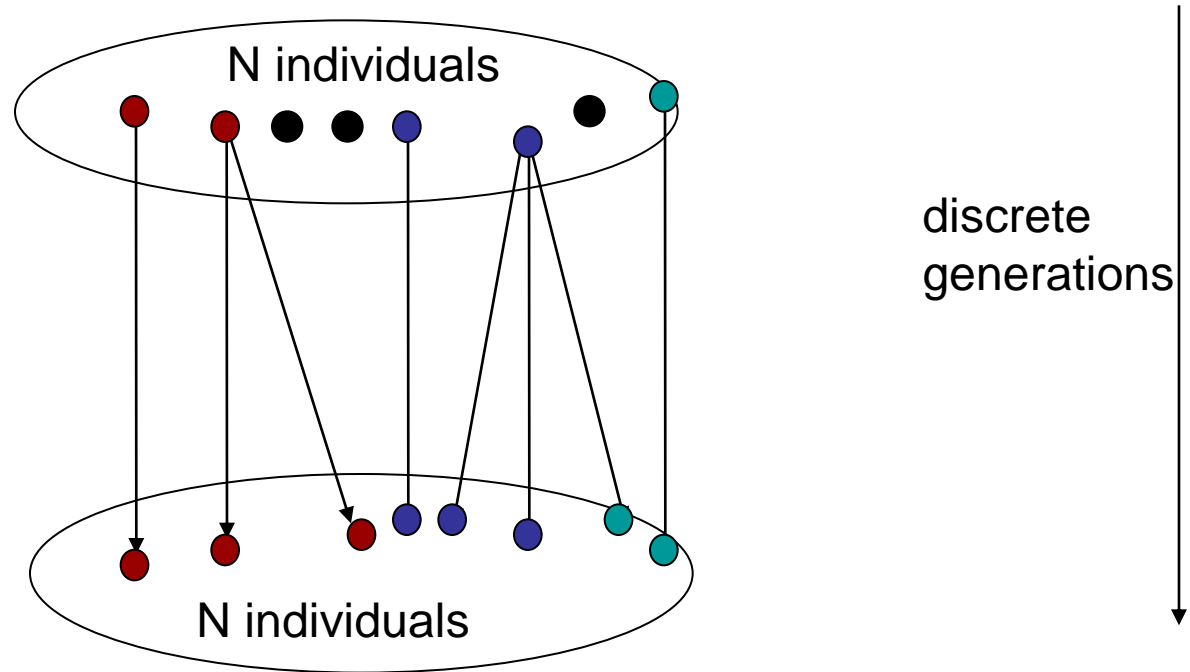
Moran Model

- One individual dies and one (possibly the same one) gives birth.

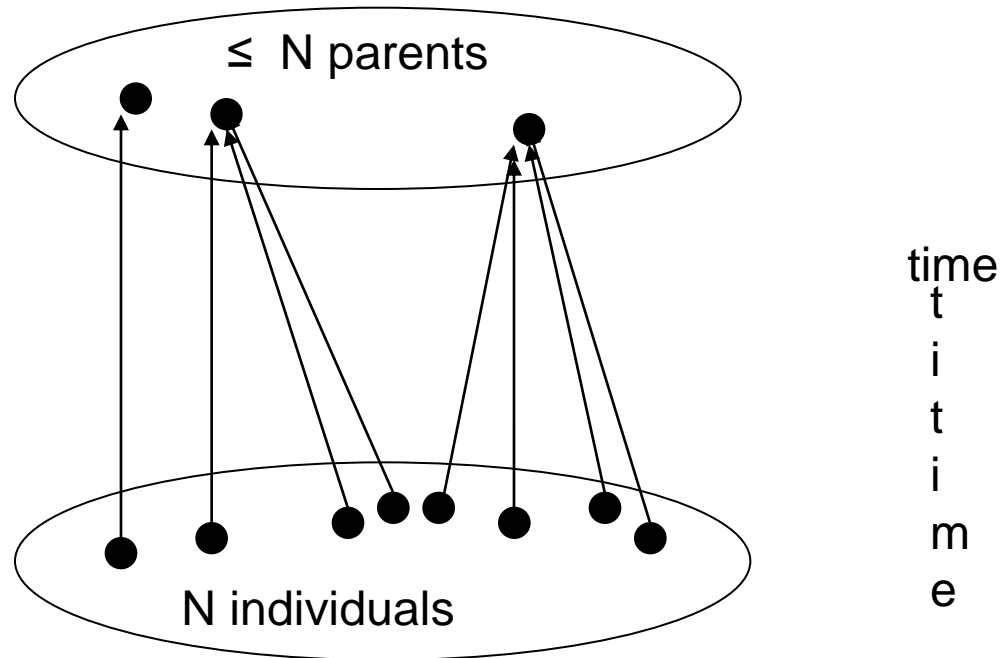
Cannings's Model

- Cannings(1974) the n individuals produce X_1, X_2, \dots, X_n offspring these being exchangeable r.v.'s. Markov Chain
- This model (which encompasses the Wright-Fisher & Moran model) gives rise to neat formulae for the eigenvalues and eigenvectors of the process.

Genetic Drift (multiple types)



The Coalescent

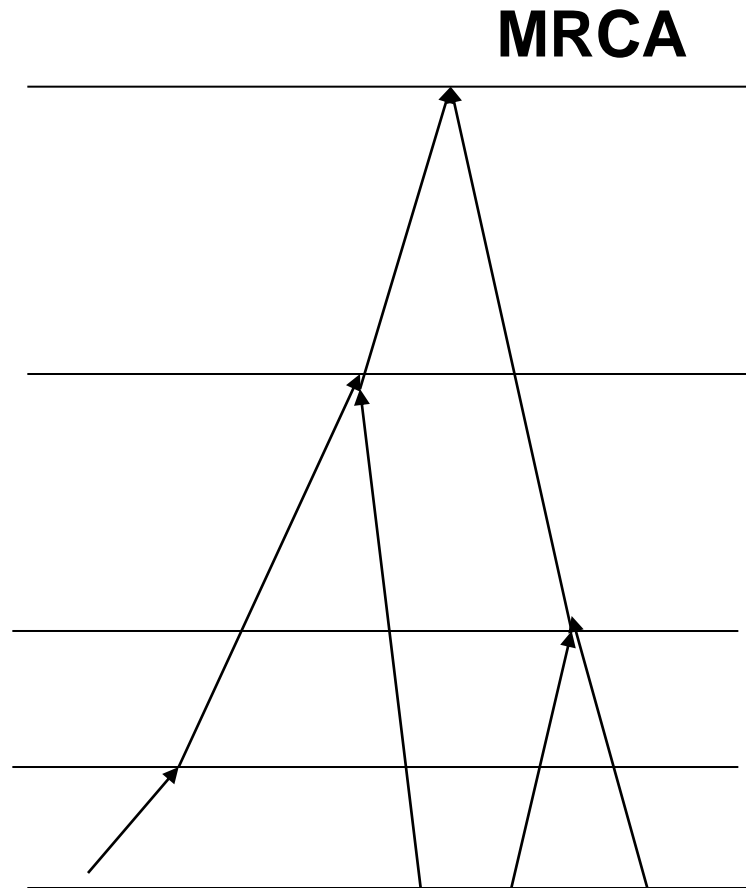


Kingman, 1982

Coalescent

- For a finite population undergoing random reproduction (i.e. not just one offspring for each adult) the current generation will all be descended from a single common ancestor (provided the population has been running for sufficient time) the Most Recent Common Ancestor (MRCA).
- Tracing back to that MRCA produces a tree, and in this tree lines coalesce as we run backwards in time.

Coalescent Tree



Coalescent

- The coalescent allows one to track the behaviour of a population without the necessity of keeping track of all N individuals through the generations.
- Process separated into two pieces
(1) produce the **coalescent tree** (time backwards)
(2) run types, with **mutation, recombination** (time forward).

Coalescent

- At each point in time (which is run backwards) keep track only of the individuals who are ancestors of the n individuals one is interested in (or sampled) at time 0.

The Coalescent

- Now suppose that for a population of N individuals we use the Wright-Fisher model. The assumption of this model is that each offspring “selects” its parent at random from the N in the previous generation.
- Consider 2 individuals in generation 0 (time will be measured backwards) and their lineages back through time.

The Coalescent

- These two lineages will at some time coalesce. When? The probability that they coalesce at time 1 is $1/N$, i.e.

$$\text{Prob}\{\text{Not coalesced by } t=1\}=(1-1/N).$$

- Each step back is independent so

$$\begin{aligned}\text{Prob}\{\text{Not coalesced by } t\}&=(1-1/N)^t \\ &=(1-1/N)^{N\tau}\end{aligned}$$

(where $\tau=t/N$)

$$\text{approx}=\text{e}^{-\tau}$$

The Coalescent

- Thus the
Prob{2 lineages not coalesced by τ } = $e^{-\tau}$
approx, which is a negative exponential
with rate 1, and the expected time to
coalescence is 1 unit of time ($=N$
generations).

The Coalescent

- Now consider k lineages (perhaps corresponding to a sample from the current generation)

Prob{none coalesce in the previous gen}=

Prob{they have k parents}=

$$(1-1/N)(1-2/N)\dots(1-(k-1)/N) = \prod_{i=0}^{k-1} (1-i/N)$$

$$\text{approx} = 1 - {}_k C_2 / N + O(1/N^2).$$

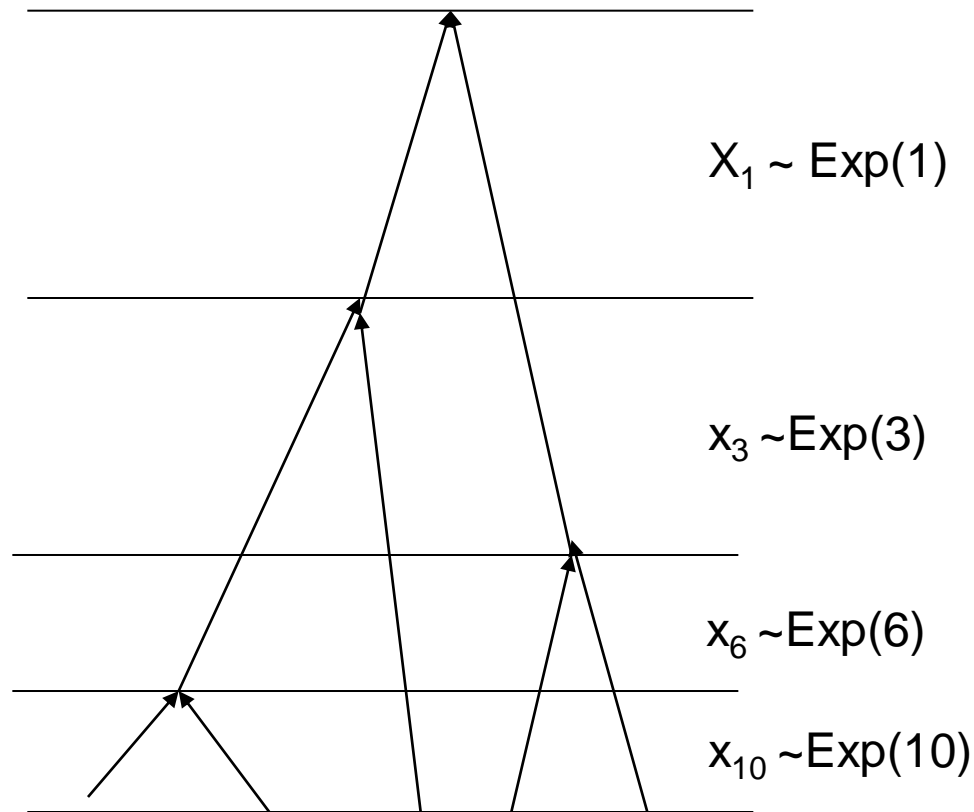
The Coalescent

- Now the probability that more than two coalesce in any generation is negligible.
- Let $T(k)$ =(scaled) time to first coalescence then, as earlier, $T(k)$ is neg. exp. with mean $2/(k(k-1))$, and these coalescences take place one at a time.

Coalescent (summary)

- At each point in time (which is run backwards) keep track only of the individuals who are ancestors of the n individuals at time 0.
- Replace **discrete** generations with **continuous** time, and assume that any pair of individuals at time t will be the offspring of one individual at $t - \delta t$ with same prob ($\lambda \delta t$) as any other pair, independently
- Scale time

The Coalescent (n=5)



Assume no multiple coalescents, only pairwise allowed.

The Coalescent

- The times we are interested in will be sums of negative exponentials (see last slide)
- There is a general expression for a sum of negative exponential random variables.

$\sum_i \eta_i$ for independent $\eta_i \sim \text{Exp}(\lambda_i)$

$\sum_i^l \eta_i$ has density function

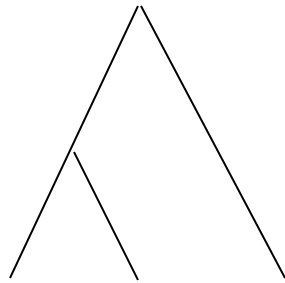
$$\sum_i^l \lambda_i \exp(-\lambda_i t) \left\{ \prod_{j \neq i}^l \lambda_j / (\lambda_j - \lambda_i) \right\}.$$

Tree topology

- We need to find the probabilities for the possible tree topologies (shapes)

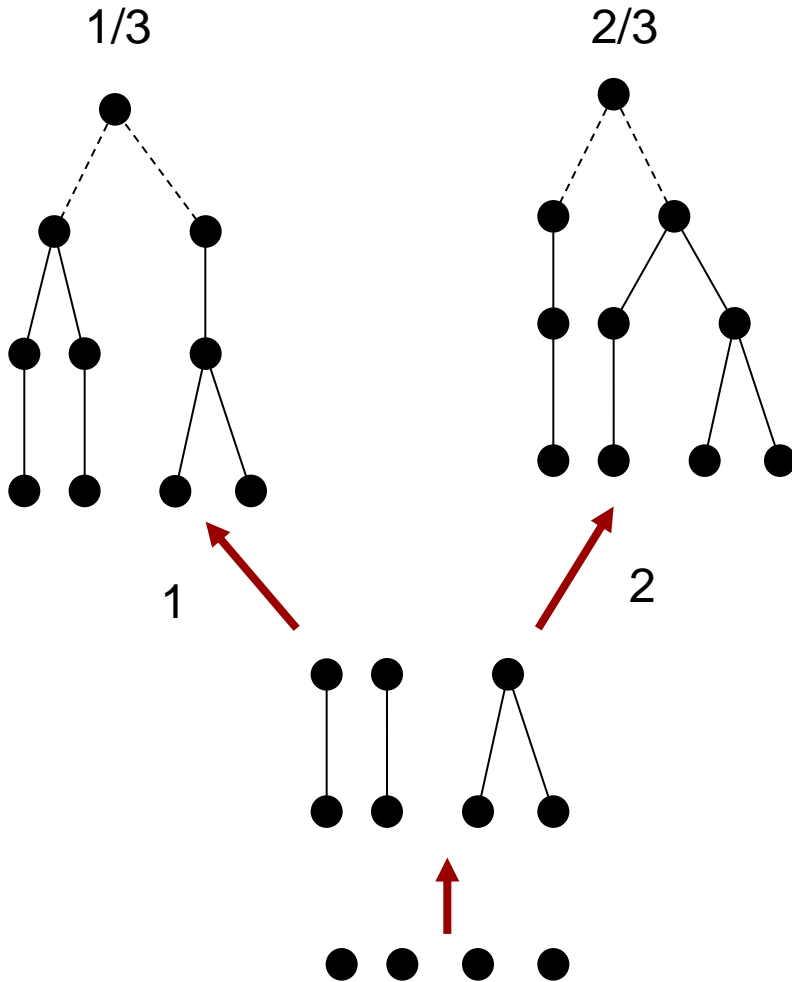
$n=3$

- $n=3$



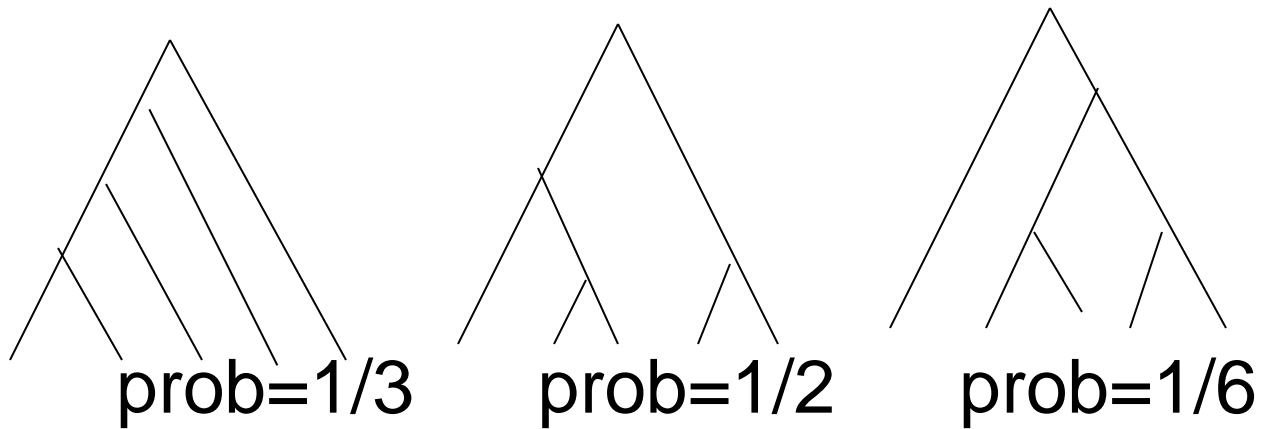
Only one possible shape.

n=4



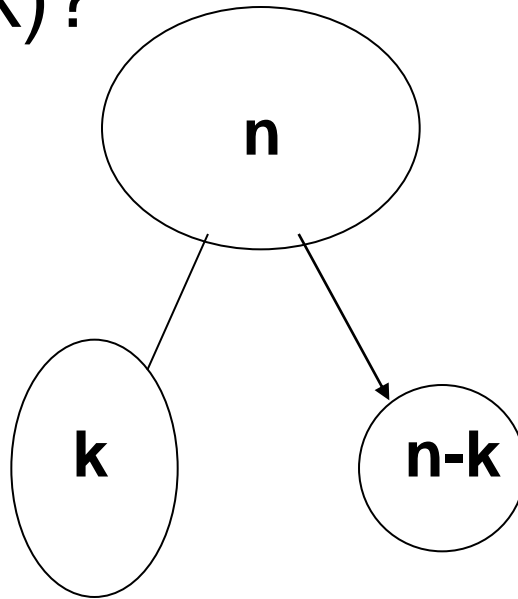
The Coalescent

$n=5$



Distribution of Tree Topology

- Suppose we have a population of size n ; what is probability that we have a split at top of $(k, n-k)$?



Distribution of Tree Topology

- There will be

$$\binom{n}{2} \binom{n-1}{2} \binom{n-2}{2} \cdots \binom{3}{2} \binom{2}{2} \binom{1}{2} = \binom{n}{2} \binom{n-1}{2} \cdots \binom{3}{2} \binom{2}{2} \binom{1}{2} = \frac{n(n-1)(n-2)\cdots 3 \cdot 2 \cdot 1}{2 \cdot 2 \cdots 2 \cdot 2} = \frac{n!}{2^{n-1}}$$

orders for choosing the pairs, and

$$\binom{k}{2} \binom{k-1}{2} \cdots \binom{2}{2} \binom{1}{2} = \frac{k!}{2^{k-1}} \text{ within the } k \text{ and}$$

$$\binom{n-k}{2} \binom{n-k-1}{2} \cdots \binom{2}{2} \binom{1}{2} = \frac{(n-k)!}{2^{(n-k)-1}} \text{ within the } (n-k).$$

Distribution of Tree Topology

Now the $k-1$ joins in the k set and the $(n-k-1)$ in the $(n-k)$ set can be ordered in ${}_{k-1}C_{k-2}$ ways so the probability of $k/(n-k)$ is

$$\frac{{}_n C_k {}_{n-k} C_{k-2} k((k-1)!)^2/2^k (n-k)((n-k)!)^2/2^{n-k}}{n((n-1)!)^2/2^n} = 1/(n-1)$$

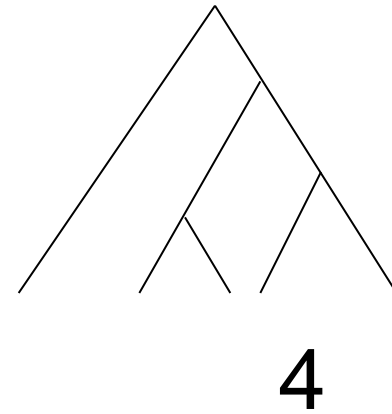
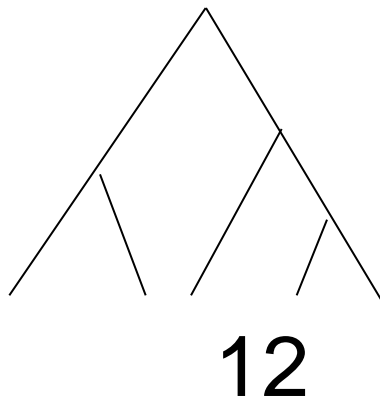
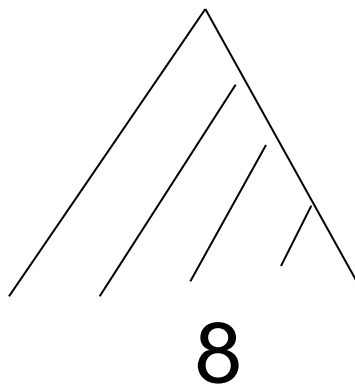
Distribution of Tree Topology

Top-down (for n leaves)

First sub-trees have r and $(n-r)$ leaves with probabilities $1/(n-1)$, then their subtrees split similarly, but you need to allow for the choice of next subtree to split.

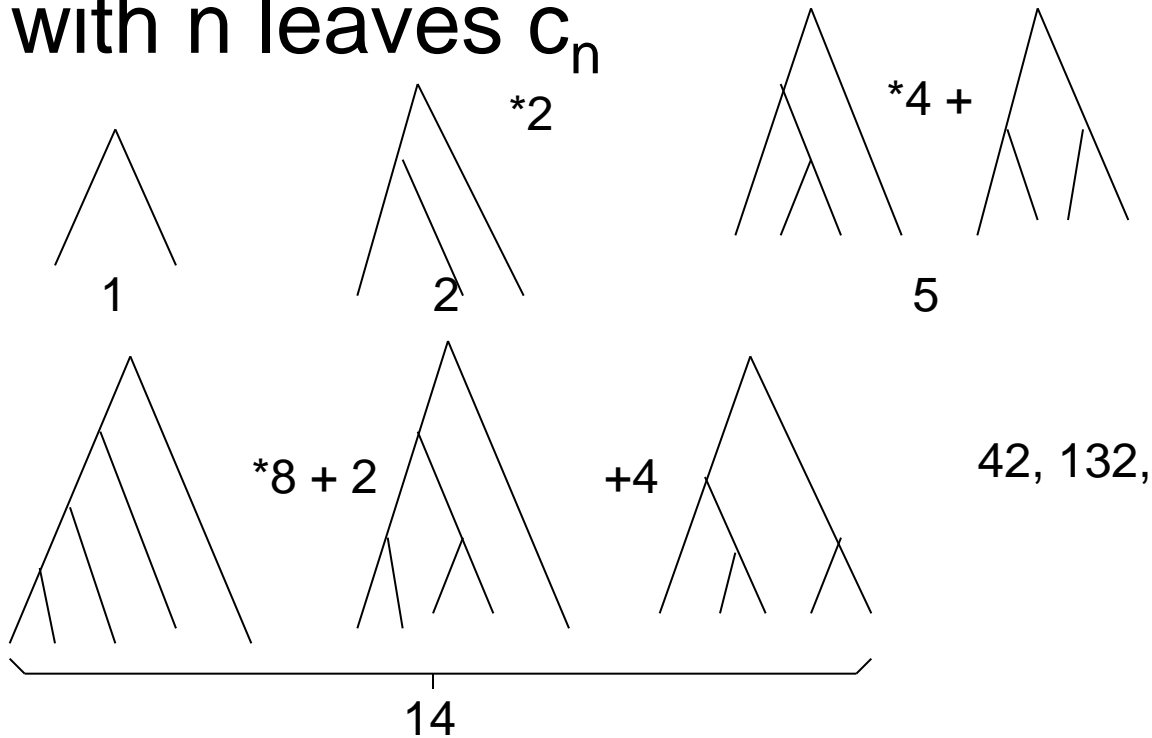
Distribution of Tree Topology

- $n=5$



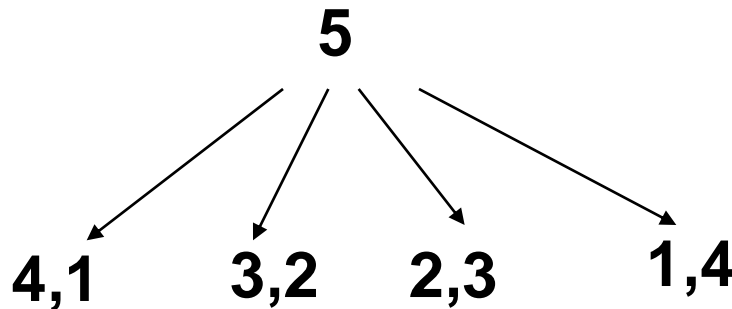
Catalan Numbers (a brief detour)

Number of distinct rooted tree topologies
with n leaves c_n



Catalan Numbers (a brief detour)

Example



$$C_n = \sum_{i=1}^{n-1} C_i C_{n-i}$$



Catalan Numbers

- Generating function

$$C(x) = \sum_n c_n x^n$$



Catalan Numbers

- Generating function

$$C(x) = \sum_n c_n x^n$$

now clearly for $n > 1$

$$c_n = \sum_{i=1}^{n-1} c_i c_{n-i}$$



Catalan Numbers

- Generating function

$$C(x) = \sum_n c_n x^n$$

now clearly for $n > 1$

$$c_n = \sum_{i=1}^{n-1} c_i c_{n-i}$$

so

$$[C(x)]^2 = C(x) - x$$



Catalan Numbers

- Generating function

$$C(x) = \sum_n c_n x^n$$

now clearly for $n > 1$

$$\text{so } c_n = \sum_{i=1}^{n-1} c_i c_{n-i}$$

$$[C(x)]^2 = C(x) - x$$

and so

$$C(x) = (1 - (1 - 4x)^{1/2})$$

Catalan Numbers

- Generating function

$$C(x) = \sum_n c_n x^n$$

now clearly for $n > 1$

$$c_n = \sum_{i=1}^{n-1} c_i c_{n-i}$$

$$[C(x)]^2 = C(x) - x$$

$$C(x) = \frac{1 - (1 - 4x)^{1/2}}{2x}$$

and so

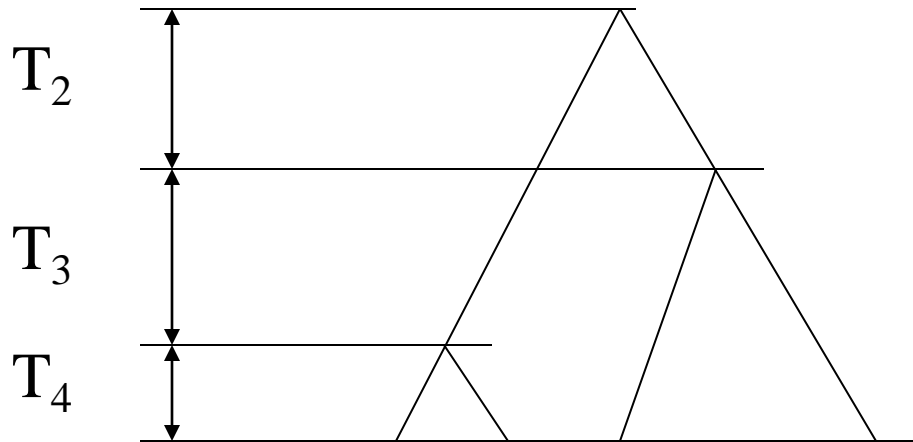
$$c_n = \frac{(2n)!}{n!(n+1)!}$$

The Coalescent. MRCA

- Thus if we start with k we have $T(k), T(k-1), \dots, T(3), T(2)$ as the times back to the MRCA (Most Recent Common Ancestor).
- Take $TT(k) = \sum_{i=2}^k T(i)$, so the expected time $E(TT(k)) = \sum_{i=2}^k 2/(i(i-1))$
 $= \sum_{i=2}^k 2 (1/(i-1) - 1/i)$
 $= 2(1-1/k)$ approx = 2 for large k .

The Coalescent.MRCA

- Thus $E(TT(2)) = 1$ is more than $\frac{1}{2}$ of the expected time for the tree for any k .
- Typical tree $n=4$



$$E(T_2) = 1$$

$$E(T_3) = 1/3$$

$$E(T_4) = 1/6$$



The Coalescent

- Variances

$$V(T_i) = (2/(i(i-1)))^2, \quad V(T_2) = 1,$$

$$TT(k) = \sum_{i=2}^k T(i), \text{ so}$$

$$\begin{aligned} V(TT(k)) &= \sum_{i=2}^k (2/(i(i-1)))^2 \\ &= \sum_{i=2}^k 2^2 (1/(i-1) - 1/i)^2 \\ &= 4 \sum_{i=2}^k (1/(i-1)^2 + 1/i^2 - \\ &\quad 2/(i(i-1))) \end{aligned}$$

The Coalescent

- $$= 4 \sum_{i=2}^k \left(\frac{1}{(i-1)^2} + \frac{1}{i^2} - \frac{2}{i(i-1)} \right)$$
$$= 8 \sum_{i=1}^{k-1} \frac{1}{i^2} + \frac{4}{k^2} - 4 - 8\left(1 - \frac{1}{k}\right)$$
from earlier
$$= 8 \sum_{i=1}^{k-1} \frac{1}{i^2} - \frac{(3k+1)(k-1)}{k^2}$$

The Coalescent. $V(TT(k))$

- $V(TT(k)) \leq 8\pi^2/6 - 12 \text{ approx} = 1.16$
so that $V(T_2)$ accounts for most of the variance.
- The following shows six realisations for a sample of size five.

The Coalescent.

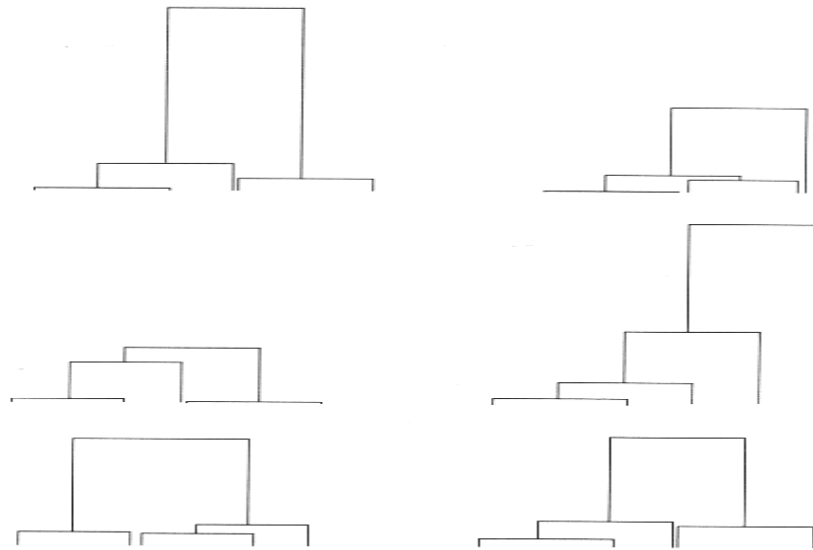


Fig. 2.3. Six realizations, drawn on the same scale, of coalescent trees for a sample of $n = 5$. (In each tree the labels 1,2,3,4,5 should be assigned at random to the leaves.)

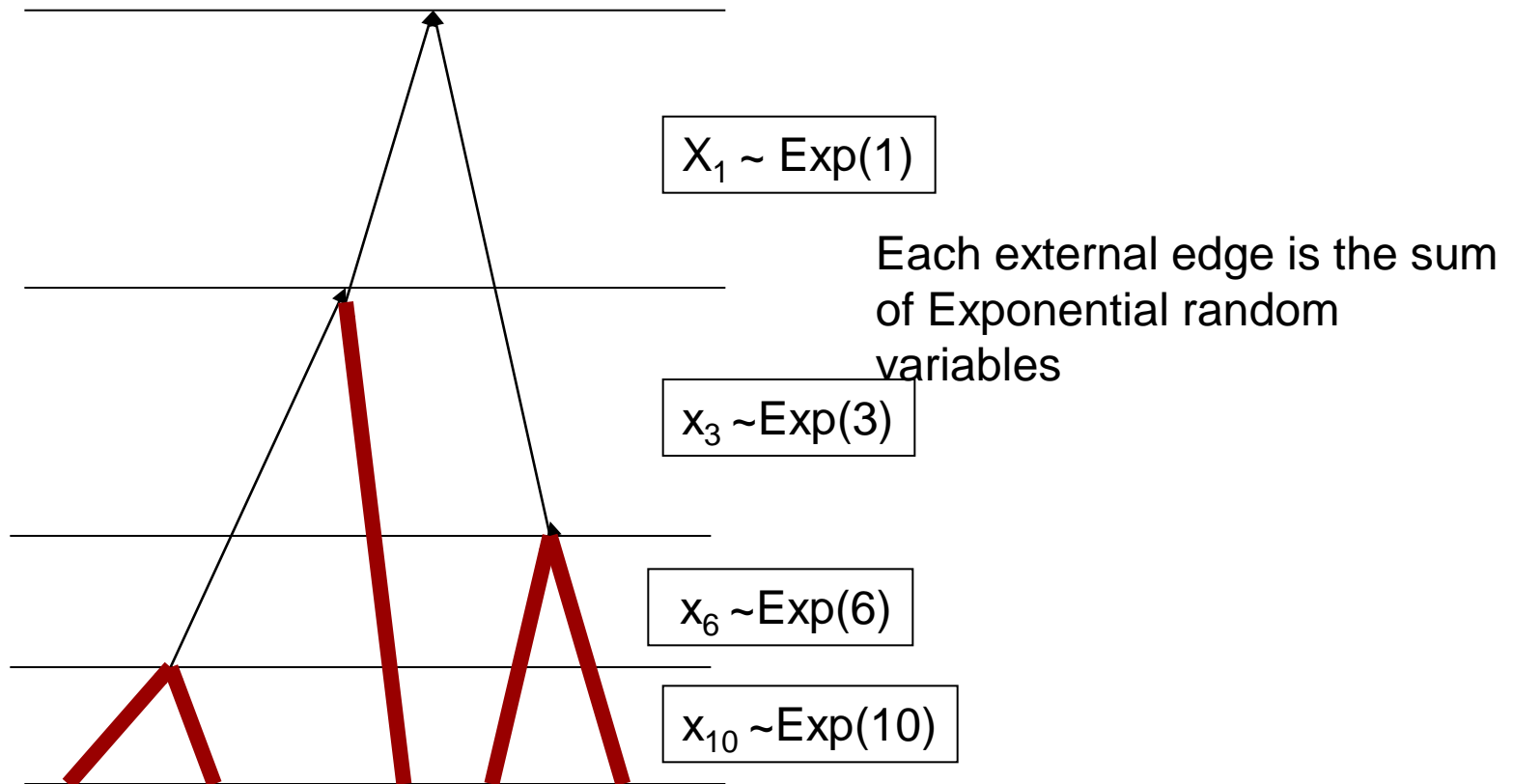
© S. Tavaré and P. Donnelly. DRAFT of September 15, 1999
Duplication for commercial purposes prohibited



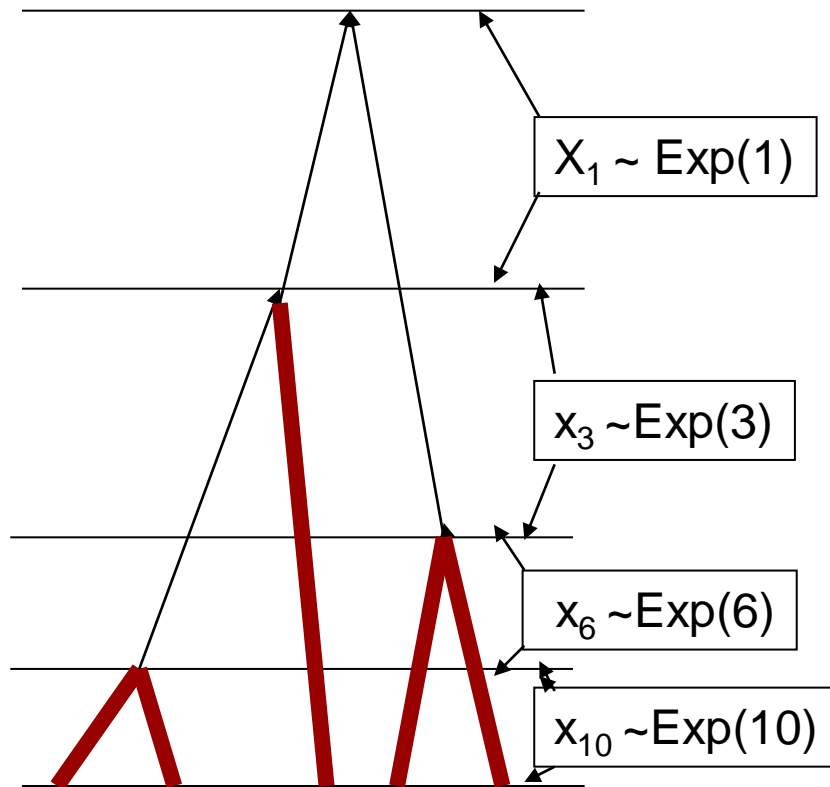
The Coalescent.

- Adding recombination to the forward process on the tree is somewhat difficult, and requires MCMC (Monte Carlo Markov Chain) a numerical estimation method.

External Edges



External Edges



Here a random edge is x_{10} , $x_{10}+x_6$, or $x_{10}+x_6+x_3$ with probabilities $2/5$, $2/5$ and $1/5$

NB. Sums of consecutive x_i 's from n to k

$\sum_i \eta_i$ for independent $\eta_i \sim \text{Exp}(\lambda_i)$

$\sum_i^l \eta_i$ has density function

$$\sum_i^l \lambda_i \exp(-\lambda_i t) \left\{ \prod_{j \neq i}^l \lambda_j / (\lambda_j - \lambda_i) \right\}.$$

$\sum_i n_i$ for λ 's = 1, 3, 6, 10

$$(3 \cdot 6 \cdot 10) / (2 \cdot 5 \cdot 9) = 2$$

$$(1 \cdot 3 \cdot 6) / (-9 \cdot -7 \cdot -4) = -1/14$$

$$(1 \cdot 6 \cdot 10) / (-2 \cdot 3 \cdot 7) = -10/7$$

$$(1 \cdot 3 \cdot 10) / (-5 \cdot -3 \cdot 4) = 1/2$$

$$2E(1) - (10/7)E(3) + (1/2)E(6) - (1/14)E(10),$$

where $E(i) = \lambda_i \exp(-\lambda_i t)$.

NB. Coefficients alternate in sign.

Prob of trees

- We can look at probability of the different tree topologies.
- Easier to find distribution of edge probabilities.

Probability(edge is $\Sigma n \rightarrow k$)

- Pick one of the n individuals then its external edge is $\text{Exp}({}_n C_2)$ with $\text{prob} = 2/n = (n-1) / {}_n C_2$.
- Edge is $\text{Exp}({}_n C_2) + \text{Exp}({}_{(n-1)} C_2)$ with $\text{prob} = (1 - 2/n) * 2/(n-1) = (n-2) / {}_n C_2$.
- Edge is $\Sigma_{i=n,k} \text{Exp}({}_i C_2)$ with $\text{prob} (k-1) / {}_n C_2$

n=5

Table 1. n=5

	$\lambda's$	p_k	$a(10exp(-10t))$	$a(6exp(-6t))$	$a(3exp(-3t))$	$exp(-t)$
1	10	4/10	1			
2	10, 6	3/10	-6/4	10/4		
3	10, 6, 3	2/10	9/14	-35/14	40/14	
4	10, 60, 3, 1	1/10	-1/14	7/14	-20/14	28/14
	<i>totals</i>	1	5/70	21/70	30/70	14/70

Distributions $n=2(1)6$

Table 2

n	e^{-t}	$3e^{-3t}$	$6e^{-6t}$	$10e^{-10t}$	$15e^{-15t}$
2	1				
3	1/2	1/2			
4	3/10	5/10	2/10		
5	14/70	30/70	21/70	5/70	
6	6/42	15/42	14/42	6/42	1/42

NB. All the coefficients are positive.

Result on sequence of exp

- Suppose that we have a set of λ_i 's and the coefficient of $\exp(\lambda_k)$ is a_k in the distribution of the random variable which is the sum of the exponential random variables. Now add an extra term with μ then the coefficient of $\exp(\lambda_k)$ is simply $a_k \times \mu / (\lambda_k - \mu)$ irrespective of the details of the sequence of λ_i 's.

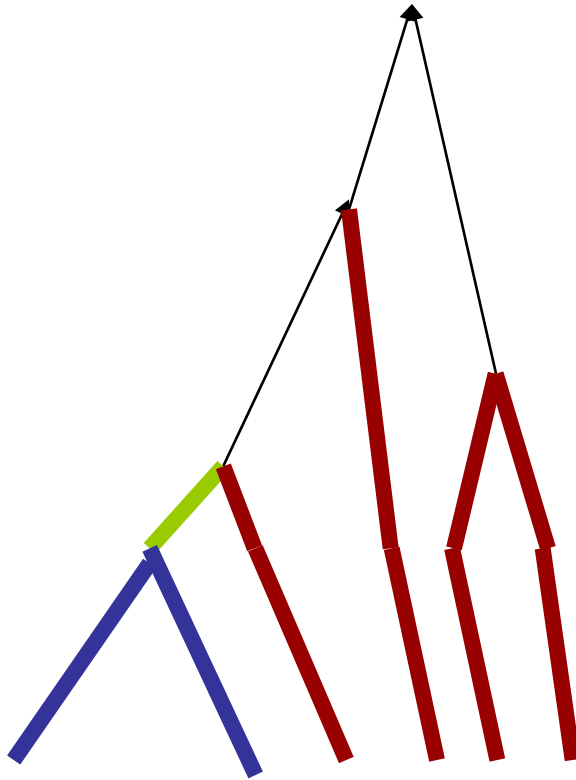
Adding an Extra $\exp(\mu)$

$$\sum_i^l \lambda_i \exp(-\lambda_i t) \left\{ \prod_{j \neq i}^l \lambda_j / (\lambda_j - \lambda_i) \right\}.$$

Then when extra $\exp(\mu)$ term in λ_i changes to

$$\lambda_i \exp(-\lambda_i t) \prod_{j \neq i} \lambda_j / (\lambda_j - \lambda_i) \times \mu / (\mu - \lambda_i)$$

Adding a layer



Coefficients of $10\exp(-10t)$

$n=5$		$n=6$
$a(10\exp(-10t))$		$a(10\exp(-10t))$
1		3
$-6/4$	$*3 =$	$-27/6$
$9/14$		$81/42$
$-1/14$		$-9/42$
$5/70$		$6/42$

$n=5 \rightarrow 6$ (i.e. add 15)

Table 1. $n=5$

	$\lambda's$	p_k	$a(10exp(-10t))$	$a(6exp(-6t))$	$a(3exp(-3t))$	$exp(-t)$
1	10	4/10	1			
2	10, 6	3/10	-6/4	10/4		
3	10, 6, 3	2/10	9/14	-35/14	40/14	
4	10, 6, 3, 1	1/10	-1/14	7/14	-20/14	28/14
	<i>totals</i>	1	5/70	21/70	30/70	14/70

n=6

Table 1b. n=6

k	λ_s	p_k	$a(15\exp(-15t))$	$a(10\exp(-10t))$	$a(6\exp(-6t))$	$a(3\exp(-3t))$	$\exp(-t)$
1	15	5/15	1				
2	15, 10	4/15	-2	3			
3	15, 10, 6	3/15	8/6	-27/6	25/6		
4	15, 10, 6, 3	2/15	-14/42	81/42	-175/42	150/42	
5	15, 10, 6, 3, 1	1/15	1/42	-9/42	35/42	-75/42	90/42
	<i>totals</i>	1	1/42	6/42	14/42	15/42	6/42

From sum row=1



Prob for Exponential with $\lambda = k C_2$ for tree of depth n , $f(n, k)$

$$\begin{aligned}
 f(n, k) &= \sum_{j=1}^{k-1} \frac{j}{n C_2} \frac{\prod_{u=j+1, u \neq k}^n C_2}{\prod_{u=j+1, u \neq k}^n \{u C_2 - k C_2\}} \\
 &= \prod_{u=k+1}^n \frac{u C_2}{\{u C_2 - k C_2\}} \frac{k C_2}{n C_2} + \sum_{j=1}^{k-1} \frac{j}{k C_2} \frac{\prod_{u=j+1, u \neq k}^{k-1} C_2}{\prod_{u=j+1, u \neq k}^{k-1} \{u C_2 - k C_2\}} \\
 &= \prod_{u=k+1}^n \frac{u C_2}{\{u C_2 - k C_2\}} \frac{k C_2}{n C_2} + f(k, k)
 \end{aligned}$$

Prob for $\exp({}_k C_2)$ for n

$$\begin{aligned}
 f(n, k) &= \sum_{j=1}^{k-1} \frac{j}{n C_2} \frac{\prod_{u=j+1, u \neq k}^n C_2}{\prod_{u=j+1, u \neq k}^n \{u C_2 - k C_2\}} \\
 &= \prod_{u=k+1}^n \frac{u C_2}{\{u C_2 - k C_2\}} \frac{k C_2}{n C_2} + \sum_{j=1}^{k-1} \frac{j}{k C_2} \frac{\prod_{u=j+1, u \neq k}^{k-1} C_2}{\prod_{u=j+1, u \neq k}^{k-1} \{u C_2 - k C_2\}} \\
 &= \prod_{u=k+1}^n \frac{u C_2}{\{u C_2 - k C_2\}} \frac{k C_2}{n C_2} + f(k, k)
 \end{aligned}$$

Prob for $\exp({}_k C_2)$ for n

$$\begin{aligned}
 f(n, k) &= \sum_{j=1}^{k-1} \frac{j}{n C_2} \frac{\prod_{u=j+1, u \neq k}^n C_2}{\prod_{u=j+1, u \neq k}^n \{u C_2 - k C_2\}} \\
 &= \prod_{u=k+1}^n \frac{u C_2}{\{u C_2 - k C_2\}} \frac{k C_2}{n C_2} * \sum_{j=1}^{k-1} \frac{j}{k C_2} \frac{\prod_{u=j+1, u \neq k}^{k-1} C_2}{\prod_{u=j+1, u \neq k}^{k-1} \{u C_2 - k C_2\}} \\
 &= \prod_{u=k+1}^n \frac{u C_2}{\{u C_2 - k C_2\}} \frac{k C_2}{n C_2} * f(k, k)
 \end{aligned}$$

>0



$f(k, k)$ = coefficient of $\exp(-_k C_2)$ for k

$$f(k, k) = \sum_{j=1}^{k-1} \frac{j}{_k C_2} \frac{\prod_{u=j+1, u \neq k}^{k-1} u C_2}{\prod_{u=j+1, u \neq k}^{k-1} \{u C_2 - k C_2\}}$$

$$= \frac{1}{_k C_2} \frac{\sum_{j=1}^{k-1} j \prod_{u=j+1, u \neq k}^{k-1} u C_2 \prod_{u=2}^j \{u C_2 - k C_2\}}{\prod_{u=2, u \neq k}^{k-1} \{u C_2 - k C_2\}}$$

Denominator of $f(k,k)$

$$\begin{aligned} \prod_{u=2}^{k-1} \{ {}_k C_2 - {}_u C_2 \} &= \prod_{u=1}^{k-2} \{ {}_k C_2 - {}_{k-u} C_2 \} \\ &= \prod_1^{k-2} u(2k - u - 1)/2 = \frac{(k-2)!(2k-2)!}{k!2^{k-2}} \\ &= \frac{[2(k-1)]!}{k!(k-1)!} \frac{(k-1)!(k-2)!}{2^{k-2}} \end{aligned}$$

Denominator of $f(k,k)$

$$= \frac{[2(k-1)]! (k-1)!(k-2)!}{k!(k-1)! 2^{k-2}} = c_n * \Delta_k$$

Where c_n is Catalan number and Δ_k products of triangular numbers

Numerator of $f(k, k)$

$$\sum_{j=1}^{k-1} j \Pi_{u-j+1, u-j+1}^{k-1} C_2 \Pi_{u-2}^j \{u C_2 - k C_2\}$$
$$= \Delta_k$$

$$f(k, k)$$

- Thus $f(k, k) = 1/c_n$.
- $f(n, k) > 0$ and we have **convex** combinations of exponentials.

References

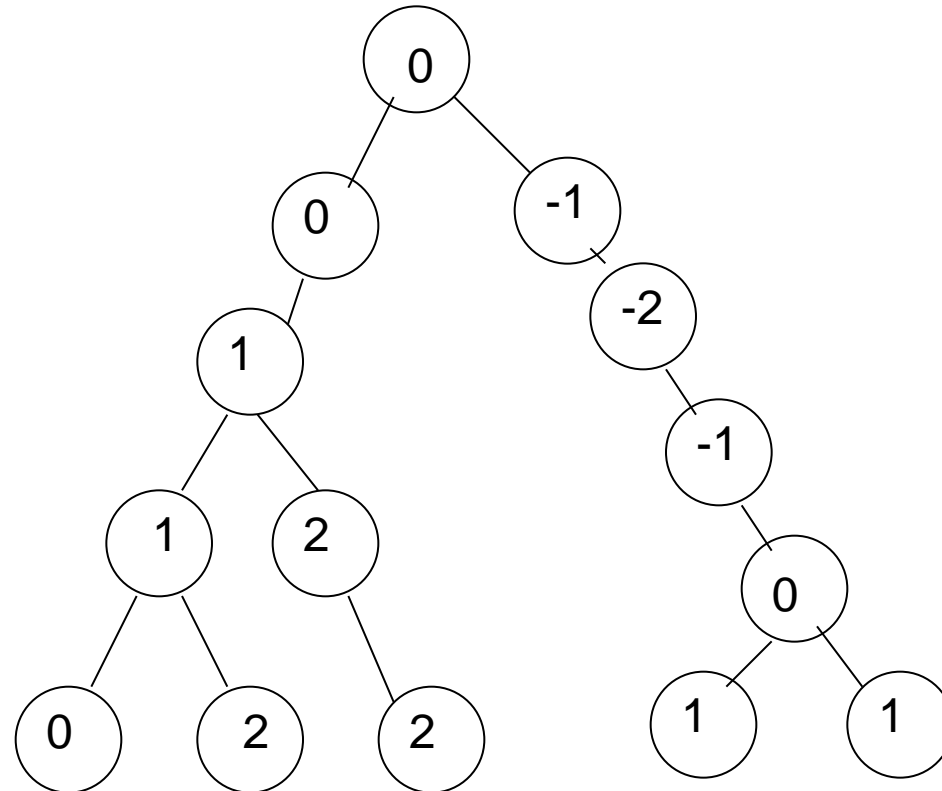
References

- Caliebe A, Neininger R, Krawczak M & Rösler U (2007) On the length distribution of external branches in coalescence trees: genetic diversity within species. *Theoret Pop Biol*, 72, 245-252.
- Cannings C (1974) The latent roots of certain Markov chains arising in genetics: A new approach, I haploid models. *Adv Appl Prob* 6, 260-290.
- Fu Y-X & Li W-H (1993) Statistical tests of neutrality of mutations. *Genetics* 133, 693-709.
- Kingman JFC (1982) The coalescent. *Stoch Proc and Appl*, 13, 235-248.
- Nordborg M (2007) Coalescent theory. *Handbook of Statistical Genetics*, Wiley, Chichester. pgs 843-877. Editors DJ Balding, M Bishop, C Cannings.

Stepwise Mutation

- Suppose that we have variants of some gene which can be mapped to the integers (they may be the number of copies of some sequence [acct] [acct].... [acct]). Suppose further that the copying from parent to offspring is subject to error so that each offspring produced by a parent with copy-number x will be $(x-1)$, x or $(x+1)$ with probabilities $\mu/2$, $(1-\mu)$, $\mu/2$.

Example



Stepwise Mutation

- How can we keep track of the distribution of values in a population of n undergoing reproduction as per the Wright-Fisher model?
- Suppose we start at time 0 with some value 0 (arbitrary). Now there will be n reproductions along the line from the founder to an individual in generation n . Thus the value will be simply $Z = \sum_{i=1, n} X_i$ where X_i is the step at i . Thus $E(Z) = 0$, $\text{Var}(Z) = n\mu$ so that $\text{Var} \rightarrow \text{inf}$ as $n \rightarrow \text{inf}$.

Stepwise Mutation

- Now consider the values normalised by some specified individual i.e. $Z_i - Z_n$. Now we need to know how many mutational events there have been between two individuals in generation n . Since there are N individuals in generation n then the chance that they share a parent in the previous generation is $1/N$, and $(1 - 1/N)$ that they have distinct parents.

Stepwise Mutation

- Thus the probability that they have a MRCA (most recent common ancestor) k generations ago is $(1-1/N)^{(k-1)} * 1/N$ for $k=1,2,\dots,(n-1)$, and $(1-1/N)^{n-1}$ for $k=n$.
- The probability generating function for the random variable in question , where

$$x = (N - 1) / N$$

$$g(s) = \frac{s(1 - (xs)^{n-1})}{N(1 - xs)} + s^n x^{n-1}$$

Stepwise Mutation

- Given k generations until join there will be $2(k-1)$ mutation “opportunities” so the pgf for this variable is $g(s^2)$.
- The pgf for the mutation step is

$$f(s) = (1 - \mu) + s\mu + s^{-1}\mu$$

so we then examine $g(f(s))$ for the separation of the two individuals.

Moments of K with pgf $h(s)$

$$E(K_{[r]}) = [d^r h(s) / ds^r]_{s=1}$$

$$d^r g(f(s)) / ds^r$$

- Here

$$h(s)=g(f(s))$$

$$d^r g(f(s)) / ds^r$$

We write $g_i = \frac{d^i g}{ds^i}$ and similarly for f_i and h_i .

Now $h = g(f)$ so $h_1 = g_1 f_1$,

$$h_2 = g_2 (f_1)^2 + g_1 f_2,$$

$$h_3 = g_3 (f_1)^3 + 3g_2 f_1 f_2 + g_1 f_3,$$

$$h_4 = g_4 (f_1)^4 + 6g_3 (f_1)^2 f_2 + 3g_2 (f_2)^2 + 4g_2 f_1 f_3 + g_1 f_4$$

$$d^r g(f(s)) / ds^r$$

- Now each time we differentiate a g_i we get $g_{i+1}f_1$ and each time we differentiate an f_i we get an f_{i+1} . So each term is of the form

$$g^{x_1} f_1^{x_2} \dots f_t^{x_t} \text{ where } \sum_{i=1}^t x_i = r$$

$$d^r g(f(s)) / ds^r$$

We write $g_i = \frac{d^i g}{ds^i}$ and similarly for f_i and h_i .

$$\text{Now } h = g(f) \text{ so } h_1 = g_1 f_1, \quad \mathbf{1}$$

$$h_2 = g_2 (f_1)^2 + g_1 f_2, \quad \mathbf{1,1 \text{ or } 2}$$

$$h_3 = g_3 (f_1)^3 + 3g_2 f_1 f_2 + g_1 f_3, \quad \mathbf{1,1,1 \text{ or } 2,1 \text{ or } 3}$$

$$h_4 = g_4 (f_1)^4 + 6g_3 (f_1)^2 f_2 + 3g_2 (f_2)^2 + 4g_2 f_1 f_3 + g_1 f_4$$

$$\mathbf{1,1,1,1 \text{ or } 2,1,1 \text{ or } 2,2 \text{ or } 3,1 \text{ or } 4}$$

$$d^r g(f(s)) / ds^r$$

- We can see that the process of differentiating n times just corresponds to placing n balls in boxes. For example with $n=5$ the term $2,2,1$ corresponds to $f_1(f_2)^2$ and then we add the

$$g_3 \quad \text{to give} \quad g_3 f_1(f_2)^2$$

- The number of terms correspond to labelled balls in boxes; Bell numbers/polynomials.

Bell numbers

- Suppose we want the coefficient for $g_3 f_1 (f_2)^2$ then we need to know how many ways we can put 5 labelled balls into 3 unlabelled boxes with 1 in one box and 2 in each of the others; just 5 ways to choose the singleton and then ${}_4C_2/2$ to choose the two dividing by two since we do not distinguish the boxes, i.e. 15.

Stepwise Mutation

- Now since $f(s) = (1 - \mu) + s\mu + s^{-1}\mu$

$$f(1) = 1, f_1(s) = \mu - \mu / s^2, f_1(1) = 0$$

Stepwise Mutation

- Thus the terms with f_1 will be eliminated from the expressions for the moments of K .
- The equivalent in terms of the balls in boxes is to not allow any box to have a single ball.

$$d^r g(f(s)) / ds^r$$

We write $g_i = \frac{d^i g}{ds^i}$ and similarly for f_i and h_i .

Now $h = g(f)$ so $h_1 = g_1 f_1$

$$h_2 = g_2 (f_1)^2 + g_1 f_2$$

$$h_3 = g_3 (f_1)^3 + 3g_2 f_1 f_2 + g_1 f_3$$

$$h_4 = g_4 (f_1)^4 + 6g_3 (f_1)^2 f_2 + 3g_2 (f_2)^2 + 4g_2 f_1 f_3 + g_1 f_4$$

All $f_i=0$ for i odd

$$d^r g(f(s)) / ds^r$$

We can only include even f 's. Thus we get

$$h_2 = g_1 f_2$$

$$h_4 = 3g_2 (f_2)^2 + g_1 f_4$$

$$h_6 = 15g_3 (f_2)^3 + 15g_2 f_2 f_4 + g_1 f_6$$

$$h_8 = 105g_4 (f_2)^4 + 210g_3 (f_2)^2 f_4 + 35g_2 (f_4)^2 + 28g_2 f_2 f_6 + g_1 f_8$$

Stepwise Mutation

- From the above theory we can easily obtain the moments of the random variable in question.
- We know the mean=0 and can prove the variance $=2\mu N(1-(1-1/N)^n)$.

Higher Order Variables

- We might look at the time back from some number greater than two to their common ancestor.

Thoughts

- By reversing time we switch attention from a complex process to a relatively simple one on a tree.
- Different questions seem to throw up distinct tree topologies in an interesting way.