Coalescent Theory: An Introduction for Phylogenetics

Laura Salter Kubatko Departments of Statistics and Evolution, Ecology, and Organismal Biology The Ohio State University Ikubatko@stat.ohio-state.edu

May 11, 2010

(4回) (1日) (日)

Why study population genetics? Wright-Fisher Model The Coalescent

Relationship between population genetics and phylogenetics

Population genetics: Study of genetic variation within a population

- Population genetics: Study of genetic variation within a population
- Phylogenetics: Use genetic variation between taxa (species, populations) to infer evolutionary relationships

イロト イヨト イヨト イヨト

- Population genetics: Study of genetic variation within a population
- Phylogenetics: Use genetic variation between taxa (species, populations) to infer evolutionary relationships
- So far, we've assumed:
 - Each taxon is represented by a single sequence this is often called "exemplar sampling"
 - We have data for a single gene and wish to estimate the evolutionary history for that gene (the gene tree or gene phylogeny)

Given current technology, we could do much more:

- Sample many individuals within each taxon (species, population, etc.)
- Sequence many genes for all individuals

イロト イヨト イヨト イヨト

- Given current technology, we could do much more:
 - Sample many individuals within each taxon (species, population, etc.)
 - Sequence many genes for all individuals
- Need models at two levels:
 - Model what happens within each population (standard population genetics)
 - Apply within-population models to each population represented on a phylogeny (more recent work)

Wright-Fisher Model

Why study population genetics? Wright-Fisher Model The Coalescent

Assumptions:

- Population of 2N gene copies
- Discrete, non-overlapping generations of equal size
- Parents of next generation of 2N genes are picked randomly with replacement from preceding generation (genetic differences have no fitness consequences)
- Probability of a specific parent for a gene in the next generation is ¹/_{2N}

イロト イポト イヨト イヨト

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

æ

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

æ

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

-2

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

-2

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

æ

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

-

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

Why study population genetics? Wright-Fisher Model The Coalescent

Wright-Fisher Model



Stat 882: Statistical Phylogenetics - Lecture 6

Why study population genetics? Wright-Fisher Model The Coalescent

The Coalescent Model

- Discrete Time Coalescent
 - P(two genes have same parent in the previous generation) is $\frac{1}{2N}$
 - Number of generations since two genes first shared a common ancestor \sim Geometric($\frac{1}{2N}$)
 - Number of generations since at least two genes in a sample of k shared a common ancestor ~ Geometric(^{k(k-1)}/_{4N})

・ロン ・回 と ・ ヨ と ・ ヨ と

▶ Define G_{k,k} to be the probability that k genes have k distinct ancestors in the previous generation. Then

Define G_{k,k} to be the probability that k genes have k distinct ancestors in the previous generation. Then

$$G_{k,k} = \left(\frac{2N-1}{2N}\right) \left(\frac{2N-2}{2N}\right) \cdots \left(\frac{2N-(k-1)}{2N}\right)$$

Define G_{k,k} to be the probability that k genes have k distinct ancestors in the previous generation. Then

$$G_{k,k} = \left(\frac{2N-1}{2N}\right) \left(\frac{2N-2}{2N}\right) \cdots \left(\frac{2N-(k-1)}{2N}\right)$$
$$= \left(1 - \frac{1}{2N}\right) \left(1 - \frac{2}{2N}\right) \cdots \left(1 - \frac{k-1}{2N}\right)$$

▶ Define G_{k,k} to be the probability that k genes have k distinct ancestors in the previous generation. Then

$$G_{k,k} = \left(\frac{2N-1}{2N}\right) \left(\frac{2N-2}{2N}\right) \cdots \left(\frac{2N-(k-1)}{2N}\right)$$
$$= \left(1 - \frac{1}{2N}\right) \left(1 - \frac{2}{2N}\right) \cdots \left(1 - \frac{k-1}{2N}\right)$$
$$= 1 - \left(\frac{1+2+3+\dots+(k-1)}{2N}\right) + \mathcal{O}\left(\frac{1}{N^2}\right)$$

Stat 882: Statistical Phylogenetics - Lecture 6

▶ Define G_{k,k} to be the probability that k genes have k distinct ancestors in the previous generation. Then

$$G_{k,k} = \left(\frac{2N-1}{2N}\right) \left(\frac{2N-2}{2N}\right) \cdots \left(\frac{2N-(k-1)}{2N}\right)$$
$$= \left(1 - \frac{1}{2N}\right) \left(1 - \frac{2}{2N}\right) \cdots \left(1 - \frac{k-1}{2N}\right)$$
$$= 1 - \left(\frac{1+2+3+\dots+(k-1)}{2N}\right) + \mathcal{O}\left(\frac{1}{N^2}\right)$$
$$= 1 - \frac{k(k-1)}{4N} + \mathcal{O}\left(\frac{1}{N^2}\right)$$

 Therefore, the probability that at least two genes share a common ancestor in the previous generation is

$$1-G_{k,k}=\frac{k(k-1)}{4N}+\mathcal{O}\left(\frac{1}{N^2}\right)$$

 Therefore, the probability that at least two genes share a common ancestor in the previous generation is

$$1-G_{k,k}=\frac{k(k-1)}{4N}+\mathcal{O}\bigg(\frac{1}{N^2}\bigg)$$

Since this is the same in each generation, we have that the number of generations until at least two genes in a sample of k shared a common ancestor ~ Geometric(\frac{k(k-1)}{4N} \])

・ロト ・回ト ・ヨト ・ヨト

- Kingman (1982a, b, c) considered the case where N (population size) is very large relative to k (sample size).
- ► Then, we can ignore the terms that are O(1/N²) this amounts to assuming that three or more genes coalescing in the same generation happens relatively rarely in comparison to two genes coalescing in one generation.

イロト イポト イヨト イヨト

- Kingman (1982a, b, c) considered the case where N (population size) is very large relative to k (sample size).
- ► Then, we can ignore the terms that are O(1/N²) this amounts to assuming that three or more genes coalescing in the same generation happens relatively rarely in comparison to two genes coalescing in one generation.
- We have
 - $\blacktriangleright\,$ Time since two gene copies had a common ancestor $\sim\,$ exponential ($\mu=2N$)
 - ► Time to coalescence of k gene copies into k − 1 ~ exponential(µ = 4N/(k(k − 1)))

where time, T, is measured in number of generations.

(日) (同) (E) (E) (E)

- Kingman (1982a, b, c) considered the case where N (population size) is very large relative to k (sample size).
- ► Then, we can ignore the terms that are O(1/N²) this amounts to assuming that three or more genes coalescing in the same generation happens relatively rarely in comparison to two genes coalescing in one generation.
- We have
 - $\blacktriangleright\,$ Time since two gene copies had a common ancestor $\sim\,$ exponential ($\mu=2N$)
 - ► Time to coalescence of k gene copies into k 1 ~ exponential(µ = 4N/(k(k - 1)))

where time, T, is measured in number of generations.

This is generally a very good approximation, provided N is large enough.

・ロン ・回 と ・ ヨ と ・ ヨ と

- To generate a genealogy of k genes under Kingman's coalescent:
 - ▶ Draw an observation from an exponential distribution with mean µ = 4N/(k(k − 1)). This will be the time of the first coalescent event (looking from the present backwards in time).

- To generate a genealogy of k genes under Kingman's coalescent:
 - ▶ Draw an observation from an exponential distribution with mean $\mu = 4N/(k(k-1))$. This will be the time of the first coalescent event (looking from the present backwards in time).
 - Pick two lineages at random to coalescence.

- To generate a genealogy of k genes under Kingman's coalescent:
 - ▶ Draw an observation from an exponential distribution with mean $\mu = 4N/(k(k-1))$. This will be the time of the first coalescent event (looking from the present backwards in time).
 - Pick two lineages at random to coalescence.
 - Decrease *k* by 1.

- To generate a genealogy of k genes under Kingman's coalescent:
 - ▶ Draw an observation from an exponential distribution with mean µ = 4N/(k(k − 1)). This will be the time of the first coalescent event (looking from the present backwards in time).
 - Pick two lineages at random to coalescence.
 - Decrease k by 1.
 - If k = 1, stop. Otherwise, repeat these steps.

・ロン ・回と ・ヨン・

Why study population genetics? Wright-Fisher Model The Coalescent

Example Genealogies Under Kingman's Coalescent



Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies

- Two measures of the size of a genealogy are commonly defined:
 - ► T_{MRCA} = the time of the most recent common ancestor of all lineages sampled
 - T_{total} = the total time represented by the geneaology
- Of interest are the mean, variance, and probability distribution of these.

(日) (同) (E) (E) (E)

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies

Define T_i to be the time in the history of the sample during which there were exactly i ancestral lineages.

• Note that
$$T_{MRCA} = \sum_{i=1}^{k} T_i$$
 and $T_{total} = \sum_{i=1}^{k} iT_i$



イロト イヨト イヨト イヨト

-2

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies - T_{MRCA}

▶ Note that
$$T_{MRCA} = \sum_{i=1}^{k} T_i$$
 and $T_i \sim Exp(\mu = \frac{4N}{i(i-1)})$
▶ Therefore, the mean is

Therefore, the mean is

$$E(T_{MRCA}) = \sum_{i=2}^{k} E(T_i) = \sum_{i=2}^{k} \frac{4N}{i(i-1)}$$
$$= 4N \sum_{i=2}^{k} \left(\frac{1}{i-1} - \frac{1}{i}\right)$$
$$= 4N \left(1 - \frac{1}{k}\right)$$

If time is measured in units of 2N generations (coalescent units), then the mean is 2(1 − ¹/_k)

(日) (同) (E) (E) (E)
Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies - T_{MRCA}

Mean time to coalescence of all lineages is

$$4N\left(1-\frac{1}{k}\right)$$

Notes:

- When k is large, it takes $\approx 4N$ generations to reach the MRCA
- When k = 2, it takes $\approx 2N$ generations to reach the MRCA
- For a large sample, much of the total time represented in the genealogy will be spent waiting for the last coalescence to occur.

・ロト ・回ト ・ヨト ・ヨト

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies - T_{MRCA}

- We can also show that $Var(T_{MRCA}) = (4N^2) \sum_{i=2}^{k} \frac{1}{i^2(i-1)^2}$
- ▶ We can show that as the sample size $k \to \infty$, $Var(T_{MRCA})$ converges to $4\pi^2/3 12 \approx 1.16$ (in coalescent units).
- Since T_{MRCA} is the sum of k − 1 independent exponential random variables T_i, we have the following distribution for T_{MRCA}:

$$f_{T_{MRCA}}(t) = \sum_{i=2}^{k} {\binom{i}{2}} e^{-\binom{i}{2}t} \prod_{j=2, j \neq i}^{k} \frac{\binom{j}{2}}{\binom{j}{2} - \binom{i}{2}}$$

イロト イポト イヨト イヨト

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies - T_{total}

$$E(T_{total}) = \sum_{i=2}^{k} i E(T_i) = \sum_{i=2}^{k} i \frac{4N}{i(i-1)}$$
$$= 4N \sum_{i=1}^{k-1} \frac{1}{i}$$

If time is measured in units of 2N generations (coalescent units), then the mean is 2∑_{i=1}^{k-1} 1/i

(日) (同) (E) (E) (E)

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies - T_{total}

• We can also show that
$$Var(T_{total}) = (2N^2) \left[4 \sum_{i=1}^{k-1} \frac{1}{i^2} \right]$$

- ▶ Note that as the sample size $k \to \infty$, $Var(T_{total})$ converges to $2\pi^2/3 \approx 6.58$ (in coalescent units).
- Since *T_{total}* is the sum of *k* − 1 independent exponential random variables *iT_i*, we have the following distribution for *T_{total}*:

$$f_{\mathcal{T}_{total}}(t) = \sum_{i=2}^{k} \frac{i-1}{2} e^{-\frac{i-1}{2}t} \prod_{j=2, j \neq i}^{k} \frac{j-1}{j-i}$$

イロト イポト イヨト イヨト

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies - T_{MRCA} and T_{total}



Stat 882: Statistical Phylogenetics - Lecture 6

▲ロン ▲圖> ▲屋> ▲屋>

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies

We need one more quantity to be able to link our population genetics model to our phylogenetic model – the probability that a specified number of coalescent events have occurred in a fixed amount of time, t.

イロン イヨン イヨン イヨン

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies

- We need one more quantity to be able to link our population genetics model to our phylogenetic model – the probability that a specified number of coalescent events have occurred in a fixed amount of time, t.
- The probability that u lineages coalesce into v lineages in time t is given by (Tavare, 1984; Watterson, 1984; Takahata and Nei, 1985; Rosenberg, 2002)

$$P_{uv}(t) = \sum_{j=v}^{u} e^{-j(j-1)t/2} \frac{(2j-1)(-1)^{j-v}}{v!(j-v)!(v+j-1)} \prod_{y=0}^{j-1} \frac{(v+y)(u-y)}{u+y}$$

イロト イポト イヨト イヨト

Why study population genetics? Wright-Fisher Model The Coalescent

Properties of Genealogies

- When u and v are small, these are easy to compute. For example,
 - $P_{21}(t)$ = probability that 2 lineages coalescence to 1 lineage in time t
 - = probability of 1 coalescent event in time t when k=2

$$= P(T \leq t),$$
 where $T \sim Exp(\mu = \frac{4N}{2(2-1)})$

$$= \int_0^t \frac{1}{2N} e^{-\frac{x}{2N}} dx = 1 - e^{-\frac{t}{2N}}$$

Similarly,

$$P_{22}(t) = \text{prob. of no coalescence in time t when k=2}$$

= $P(T > t)$
= $\int_{t}^{\infty} \frac{1}{2N} e^{-\frac{x}{2N}} dx = e^{-\frac{t}{2N}}$

Stat 882: Statistical Phylogenetics - Lecture 6

イロン イヨン イヨン イヨン

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

The Coalescent Model Along a Species Tree

- So far, we've considered the coalescent process within a single population.
- A phylogenetic tree consists of many populations followed throughout evolutionary time:



- A 同 ト - A 三 ト - A 三 ト

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

The Coalescent Model Along a Species Tree

- Goal is to apply coalescent model across the phylogeny. The basic assumption is that events that occur in one population are independent of what happens in other populations within the phylogeny.
- More specifically, given the number of lineages entering and leaving a population, coalescent events within populations are independent of one another.
- It is also important to recall an assumption we "inherit" from our population genetics model: all pairs of lineages are equally likely to coalesce within a population.

・ロト ・回ト ・ヨト ・ヨト

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

The Coalescent Model Along a Species Tree

- When talking about gene tree distributions, there are two cases of interest:
 - The gene tree topology distribution
 - The joint distribution of topologies and branch lengths
- Start with the simple case of 3 species with 1 lineage sampled in each and look at the gene tree topology distribution

Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

Example of gene tree probability computation (for simplicity, let's use coalescent units for our time scale):

(a)
$$Prob = 1 - e^{-t}$$
; (b), (c), (d) $Prob = \frac{1}{3}e^{-t}$



イロト イポト イヨト イヨト

Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

- Thus, we have the following probabilities:
 - Gene tree (A,(B,C)): prob = $1 e^{-t} + \frac{1}{3}e^{-t} = 1 \frac{2}{3}e^{-t}$ Gene tree (B,(A,C)): prob = $\frac{1}{3}e^{-t}$

 - Gene tree (C,(A,B)): prob = $\frac{1}{2}e^{-t}$
- Note: There are two ways to get the first gene tree. We call these histories.
- The probability associated with a gene tree topology will be the sum over all histories that have that topology.

Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

What are these probabilities like as a function of t, the length of time between speciation events?



Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Example: A Slightly Larger Case

Consider 4 taxa – the human-chimp-gorilla problem



Stat 882: Statistical Phylogenetics - Lecture 6

◆□ > ◆□ > ◆臣 > ◆臣 > ○

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Coalescent Histories for the 4-taxon Example

There are 5 possibilities for this example:



Stat 882: Statistical Phylogenetics - Lecture 6

Computing the Topology Distribution by Enumerating Histories

► In the general case, we have the following:

The probability of gene tree g given species tree \mathcal{S} is given by

$$P\{G = g|S\} = \sum_{histories} P\{G = g, history|S\}$$

イロン イヨン イヨン イヨン

Computing the Topology Distribution by Enumerating Histories

▶ In the general case, we have the following:

The probability of gene tree g given species tree \mathcal{S} is given by

$$P\{G = g | S\} = \sum_{\text{histories}} P\{G = g, \text{history} | S\}$$
$$= \sum_{\text{histories}} \prod_{b} w_{b} P_{u(b), v(b)}(t_{b})$$

Degnan and Salter, Evolution, 2005

・ロン ・回 と ・ ヨ と ・ ヨ と

Computing the Topology Distribution by Enumerating Histories

• The probability of gene tree g given species tree S is given by

$$P\{G = g | S\} = \sum_{\text{histories}} P\{G = g, \text{history} | S\}$$
$$= \sum_{\text{histories}} \prod_{b} w_{b} P_{u(b),v(b)}(t_{b})$$

Number of terms only known in special cases (Rosenberg, 2007)

Computing the Topology Distribution by Enumerating Histories

• The probability of gene tree g given species tree S is given by

$$P\{G = g | S\} = \sum_{\text{histories}} P\{G = g, \text{history} | S\}$$
$$= \sum_{\text{histories}} \prod_{b} w_{b} P_{u(b), v(b)}(t_{b})$$

Multiply probabilities associated with history over internal branches (once the number of lineages entering and leaving a branch is known – which is what is given by the histories – coalescence happens independently along branches)

Computing the Topology Distribution by Enumerating Histories

• The probability of gene tree g given species tree S is given by

$$P\{G = g | S\} = \sum_{\text{histories}} P\{G = g, \text{history} | S\}$$
$$= \sum_{\text{histories}} \prod_{b} w_{b} P_{u(b),v(b)}(t_{b})$$

Probability of getting sequence of coalescent events that is consistent with g

Computing the Topology Distribution by Enumerating Histories

• The probability of gene tree g given species tree S is given by

$$P\{G = g | S\} = \sum_{\text{histories}} P\{G = g, \text{history} | S\}$$
$$= \sum_{\text{histories}} \prod_{b} w_{b} P_{u(b),v(b)}(t_{b})$$

Probability that u lineages coalescent into v in time t_b

・ロン ・回と ・ヨン・

Computing the Topology Distribution by Enumerating Histories

• The probability of gene tree g given species tree S is given by

$$P\{G = g | S\} = \sum_{\text{histories}} P\{G = g, \text{history} | S\}$$
$$= \sum_{\text{histories}} \prod_{b} w_{b} P_{u(b), v(b)}(t_{b})$$

Length of branch b

・ロト ・回ト ・ヨト ・ヨト

Computing the Topology Distribution by Enumerating Histories

TABLE 3. The number of valid coalescent histories when the gene tree and species tree have the same topology. The number of histories is also the number of terms in the outer sum in equation (12).

	Number of histories		
Taxa	Asymmetric trees	Symmetric trees	Number of topologies
4	5	4	15
5	14	10	105
6	42	25	945
7	132	65	10,395
8	429	169	135,135
9	1430	481	2,027,025
10	4862	1369	34,459,425
12	58,786	11,236	13,749,310,575
16	9,694,845	1,020,100	6.190×10^{15}
20	1,767,263,190	100,360,324	8.201×10^{21}

Degnan and Salter, Evolution, 2005

Applications of the Topology Distribution - Example 1

- Motivation: Paper by Ebersberger et al. 2007. Mol. Biol. Evol. 24:2266-2276
- Examined 23,210 distinct alignments for 5 primate taxa: Human, Chimp, Gorilla, Orangutan, Rhesus
- Looked at distribution of gene trees among these taxa observed strongly supported incongruence only among the Human-Chimp-Gorilla clade.

Applications of the Topology Distribution - Example 1



・ロン ・回と ・ヨン・

Applications of the Topology Distribution - Example 1



Observed proportions of each gene tree among ML phylogenies

・ロン ・回 と ・ ヨ と ・ ヨ と

Applications of the Topology Distribution - Example 1



Observed proportions of each gene tree among ML phylogenies

Predicted proportions using parameters from Rannala & Yang, 2003.

イロト イポト イヨト イヨト

Applications of the Topology Distribution - Example 2

- In the previous example, one topology is clearly preferred
- Must the distribution always look this way?
- Examine entire distribution when the number of taxa is small

・ロン ・回 と ・ ヨ と ・ ヨ と

Applications of the Topology Distribution - Example 2

- Consider 4 taxa: A, B, C, and D
- Species tree:



 Look at probabilities of all 15 tree topologies for values of x, y, and z

イロト イヨト イヨト イヨト

Applications of the Topology Distribution - Example 2



・ロン ・回 と ・ ヨ と ・ ヨ と

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of the Topology Distribution - Example 2



・ロン ・回 と ・ ヨ と ・ ヨ と

Applications of the Topology Distribution - Example 2



・ロン ・回 と ・ ヨ と ・ ヨ と

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of the Topology Distribution - Example 2



Degnan and Rosenberg, *PLoS Genetics*, 2006 Rosenberg and Tao, *Systematic Biology*, 2008 The existence of anomalous gene trees has implications for the inference of species trees

イロト イヨト イヨト イヨト

Applications of the Topology Distribution - Example 3

- What about mutation? How does this affect data analysis?
- The coalescent gives a model for determining gene tree probabilities for each gene.
- ► View DNA sequence data as the result of a two-stage process:
 - Coalescent process generates a gene tree topology.
 - Given this gene tree topology, DNA sequences evolve along the tree.

・ロン ・回 と ・ ヨ と ・ ヨ と

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of the Topology Distribution - Example 3

Given this model, how should inference be carried out?
Applications of the Topology Distribution - Example 3

- Given this model, how should inference be carried out?
- Hypothesis: As more data (genes) are added, the process of estimating species trees from concatenated data can be statistically inconsistent
- May fail to converge to any single tree topology if there are many equally likely trees.
- May converge to the wrong tree when a gene tree that is topologically incongruent with the species tree has the highest probability.

・ロン ・回 と ・ ヨ と ・ ヨ と

Applications of the Topology Distribution - Example 3



Stat 882: Statistical Phylogenetics - Lecture 6

ヨン

Population Genetics Models Coalescent Theory for Phylogenetic Inference Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of the Topology Distribution - Example 3 Simulation Study 1



イロン イヨン イヨン イヨン

Population Genetics Models Coalescent Theory for Phylogenetic Inference Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of the Topology Distribution - Example 3 Simulation Study 2



Stat 882: Statistical Phylogenetics - Lecture 6

イロン イヨン イヨン イヨン

Applications of the Topology Distribution - Example 3

- Performance of the Concatenation Approach:
 - Can be statistically inconsistent when branch lengths in the species phylogeny are sufficiently small
 - May perform poorly even when branch lengths are only moderately short

・ロト ・回ト ・ヨト ・ヨト

Applications of the Topology Distribution - Example 3

- Performance of the Concatenation Approach:
 - Can be statistically inconsistent when branch lengths in the species phylogeny are sufficiently small
 - May perform poorly even when branch lengths are only moderately short
- What should we do? Need to design inference methods that incorporate the coalescent process. Dennis's lecture next week

・ロン ・回 と ・ ヨ と ・ ヨ と

Joint Density of Gene Tree Topology and Branch Lengths - An Example

Rannala and Yang 2003



$$2exp\{-2\frac{3(2)}{2}t_{ABCE}\}exp\{-2(t_{ABCE} - t_{ABCDE})\}$$

$$2exp\{-2\frac{3(2)}{2}t_{ABC}\}exp\{-2(\tau_1 - t_{ABC})\}$$

$$2exp\{-2\frac{3(2)}{2}t_{AB}\}exp\{-2(\tau_2 - t_{AB})\}$$

 $exp\{-2\tau_3\}$

Stat 882: Statistical Phylogenetics - Lecture 6

イロト イヨト イヨト イヨト

Gene Tree Branch Length Distributions

We now have the following distributions

We can thus, in theory, get the distributions of gene tree branches by simply manipulating these quantities:

$$f(\mathbf{t}|G,S) = rac{f(g,\mathbf{t}|S)}{p(g|S)}$$

- Integrating out branches which aren't of interest gives joint or marginal distributions
- Can even examine correlations between branch lengths

イロト イポト イヨト イヨト

Gene Tree Branch Length Distributions

- Complication: Region of integration will change for each history within a given gene tree
- Branch length densities are then a mixture over histories
- For the case of four taxa, James Degnan and I have worked out all joint and marginal distributions
- Simulate data and compare theoretical distribution to observed distribution
- Correlations are also well-approximated by simulation

・ロン ・回 と ・ ヨ と ・ ヨ と

Applications of Branch Length Distributions - Example 1

- Simulate 1,000,000 gene trees from species tree ((A:1.0,B:1.0):1.0,(C:1.5,D:1.5):0.5);
- Of these, 449,599 had the same topology as the species tree
- Compare observed distribution of branch length connecting (A,B) to root node to true distribution



 Good fit between observed and true distributions

イロト イポト イヨト イヨト

- Estimation of speciation times using information in gene trees is often desirable
- Under the coalescent model (with no gene flow following speciation), it must the case that gene divergence times pre-date speciation times

・ロト ・回ト ・ヨト ・ヨト

- Estimation of speciation times using information in gene trees is often desirable
- Under the coalescent model (with no gene flow following speciation), it must the case that gene divergence times pre-date speciation times



イロト イポト イヨト イヨト

- Estimation of speciation times using information in gene trees is often desirable
- Under the coalescent model (with no gene flow following speciation), it must the case that gene divergence times pre-date speciation times



- What is the distribution of this difference?
- How does it depend on species tree shape (e.g., symmetry) and species tree branch lengths?

イロト イヨト イヨト イヨト

When considering the distribution of the MRCA of a sample of lineages is of interest, we can simplify computation of the distribution

・ロト ・回ト ・ヨト ・ヨト

- When considering the distribution of the MRCA of a sample of lineages is of interest, we can simplify computation of the distribution
- Let T be the distribution of the difference between the speciation time and the time of the MRCA of all lineages. Note that

$$f_{T|\mathcal{S}}(t) = \sum_{n=2}^{k} Pr(L=n|\mathcal{S})P_{n1}(t)$$

where L is the random number of lineages available to coalesce above the root of the species tree

・ロト ・回ト ・ヨト ・ヨト

- When considering the distribution of the MRCA of a sample of lineages is of interest, we can simplify computation of the distribution
- Let T be the distribution of the difference between the speciation time and the time of the MRCA of all lineages. Note that

$$f_{T|\mathcal{S}}(t) = \sum_{n=2}^{k} Pr(L=n|\mathcal{S})P_{n1}(t)$$

where L is the random number of lineages available to coalesce above the root of the species tree

► Pr(L = n|S) can be computed recursively in a peeling-type algorithm

Efromovich and Kubatko, SAGMB, 2008

◆□▶ ◆□▶ ◆目▶ ◆目▶ ◆□ ● ● ●

Population Genetics Models Coalescent Theory for Phylogenetic Inference Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of Branch Length Distributions - Example 2

A particular example



イロト イヨト イヨト イヨト

-2

Applications of Branch Length Distributions - Example 2

A particular example

When τ₂ = τ₃ = 1.0, we have the following:





イロト イヨト イヨト イヨト

-2

- What can we conclude from this?
 - Shorter branches lead to more potential for incomplete lineage sorting, which results in longer times to the MRCA
 - This effect will be most pronounced for branches that are close to the root of the tree

イロン イ部ン イヨン イヨン 三連

- What can we conclude from this?
 - Shorter branches lead to more potential for incomplete lineage sorting, which results in longer times to the MRCA
 - This effect will be most pronounced for branches that are close to the root of the tree
- What about symmetry of the tree?

イロト イポト イヨト イヨト

Population Genetics Models Coalescent Theory for Phylogenetic Inference Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

Applications of Branch Length Distributions - Example 2

Consider a more symmetric tree:



- Two internal branches adjacent to the root node
- Suggests more possibility of incomplete lineage sorting longer time to MRCA

イロト イヨト イヨト イヨト

-2

Coalescent Theory

Gene Tree Topology Distributions Applications of the Gene Tree Topology Distribution Gene Tree Branch Length Distributions

- We now have the main ideas of the coalescent model and how to apply it to a phylogeny.
 - But there are many things we haven't discussed: migration, recombination, etc.
- Next week: How can we take this model and use it to infer a species-level phylogeny?
- Thursday's lab: Using the program COAL to compute gene tree topology probabilities.

・ロン ・回 と ・ ヨ と ・ ヨ と