State-space exploration of Tajima Trees

Julia A. Palacios Joint with A. Véber, J. Wakeley and S. Ramachandran

Department of Statistics Department of Biomedical Data Science Stanford University

> CIRM July 2018

> > < 日 > < 同 > < 回 > < 回 > < □ > <

1/54

Motivation: Estimation of Effective Population Size

Effective Population Size Trajectory Ne(t)

 $N_e(t)$ is a measure of relative genetic diversity over time.

Motivation: Estimation of Effective Population Size

Effective Population Size Trajectory Ne(t)

 $N_e(t)$ is a measure of relative genetic diversity over time.

Why is it important?



Motivation: Estimation of Effective Population Size

Effective Population Size Trajectory Ne(t)

 $N_e(t)$ is a measure of relative genetic diversity over time.

Why is it important?

Viral Gene Sequences Reveal the Variable History of Hepatitis C Virus Infection among Countries

Tatsunori Nakano,¹ Ling Lu,^{2,*} Pengbo Liu,³ and Oliver G. Pybus⁴

Department of Internet Medicine, Ichinomiya Nishi Hospital, Ichinomiya, Akhir, Japan; 'Division of Digestive Disease, Department of Medicine, and 'Department of Pathology and Laboratory Medicine, Emory University School of Medicine, Atlanta, Georgia; 'Department of Zoology, University of Oxford, Xoridy, University Medicine, Emory University School of Medicine, Atlanta, Georgia; 'Department of Zoology,

Example 1: Hepatitis C Virus



Prevalence of HCV - WHO 1999

- Identified in 1989
- Spread by blood to blood contact
- ≈3% of infected population worldwide
- 8,000 10,000 deaths per year in the USA
- Egypt has the highest prevalence

- 62 samples in 1993 from the E1 gene (411bp)
- Parenteral antischistosomal therapy (PAT) was practiced from 1920s to 1980s
- In the 1970s started a transition from the intravenous to the oral administration of the PAT

Example 1: Hepatitis C Virus



- 62 samples in 1993 from the E1 gene (411bp)
- Parenteral antischistosomal therapy (PAT) was practiced from 1920s to 1980s
- In the 1970s started a transition from the intravenous to the oral administration of the PAT

Present-day DNA data inform us about the past



Present-day DNA data inform us about the past



• In humans, mutation rate is estimated to be $\approx 10^{-8}$ per base per generation.

Present-day DNA data inform us about the past



- In humans, mutation rate is estimated to be $\approx 10^{-8}$ per base per generation.
- Recent ancestry indicates small population sizes.

Goal: Estimation of Effective Population Size

Coalescent-based model



Goal: Estimation of Effective Population Size

Coalescent-based model



Ancestral process: coalescent process of genealogies

 Mutation process: Poisson process along the branches of the genealogy

Goal: Estimation of Effective Population Size

Coalescent-based model



• Ancestral process: coalescent process of genealogies.

 Mutation process: Poisson process along the branches of the genealogy.

 Population process: Effective population size trajectory over time.



 $P(N_e(t), \mathbf{G}, \mathbf{Q}, \tau \mid \mathbf{Y}) \propto \underline{P(\mathbf{Y} \mid \mathbf{G}, \mathbf{Q})} \underline{P(\mathbf{G} \mid N_e(t))} P(\mathbf{Q}) \underline{P(N_e(t) \mid \tau)} P(\tau)$

prior

 $\log \mathsf{GP}(0, \mathbf{C}(\tau))$



 $P(N_{e}(t), \mathbf{G}, \mathbf{Q}, \tau \mid \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} \mid \mathbf{G}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{G} \mid N_{e}(t))}_{\text{Coalescent prior}} P(\mathbf{Q}) \underbrace{P(N_{e}(t) \mid \tau)}_{\text{log } \mathbf{GP}(0, \mathbf{C}(\tau))} P(\tau)$

• The likelihood $P(\mathbf{Y} | \mathbf{G}, \mathbf{Q})$ is tractable.



 $P(N_{e}(t), \mathbf{G}, \mathbf{Q}, \tau \mid \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} \mid \mathbf{G}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{G} \mid N_{e}(t))}_{\text{Coalescent prior}} P(\mathbf{Q}) \underbrace{P(N_{e}(t) \mid \tau)}_{\text{log } \mathbf{GP}(0, \mathbf{C}(\tau))} P(\tau)$

• The likelihood $P(\mathbf{Y} | \mathbf{G}, \mathbf{Q})$ is tractable.

The state space of genealogies \mathcal{G}

•
$$\mathcal{G} = \mathcal{T}_n \otimes \mathbb{R}^{+n-1}$$



 $P(N_{e}(t), \mathbf{G}, \mathbf{Q}, \tau \mid \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} \mid \mathbf{G}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{G} \mid N_{e}(t))}_{\text{Coalescent prior}} P(\mathbf{Q}) \underbrace{P(N_{e}(t) \mid \tau)}_{\text{log GP}(0, \mathbf{C}(\tau))} P(\tau)$

• The likelihood $P(\mathbf{Y} | \mathbf{G}, \mathbf{Q})$ is tractable.

The state space of genealogies \mathcal{G}

•
$$\mathcal{G} = \mathcal{T}_n \otimes \mathbb{R}^{+n-1}$$

•
$$|\mathcal{T}_n| = n!(n-1)!/2^{n-1}$$

•
$$|\mathcal{T}_{100}| \approx 10^{284}$$



 $P(N_{e}(t), \mathbf{G}, \mathbf{Q}, \tau \mid \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} \mid \mathbf{G}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{G} \mid N_{e}(t))}_{\text{Coalescent prior}} P(\mathbf{Q}) \underbrace{P(N_{e}(t) \mid \tau)}_{\text{log GP}(0, \mathbf{C}(\tau))} P(\tau)$

• The likelihood $P(\mathbf{Y} | \mathbf{G}, \mathbf{Q})$ is tractable.

The state space of genealogies \mathcal{G}

•
$$\mathcal{G} = \mathcal{T}_n \otimes \mathbb{R}^{+n-2}$$

•
$$|\mathcal{T}_n| = n!(n-1)!/2^{n-1}$$

•
$$|\mathcal{T}_{100}| \approx 10^{284}$$

- $\approx 10^{80}$ atoms in the universe
- $\approx 4.4 \times 10^{17}$ seconds since the Big Bang

Coalescent times alone are sufficient statistics for N(t)



Coalescent Density:

$$P(\mathbf{G}|N_e(t)) \propto \prod_{k=2}^n P[t_{k-1}|t_k, N_e(t)].$$

 $\mathbf{t} = (t_2, t_3, \dots, t_n)$ are sufficient statistics for inferring $N_e(t)$



What if we replace **G** with **t**, the vector of coalescent times? Posteriors:

 $P(N_e(t), \mathbf{G}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y}) \text{ vs } P(N_e(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y})$



What if we replace **G** with **t**, the vector of coalescent times? Posteriors:

 $P(N_e(t), \mathbf{G}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y}) \text{ vs } P(N_e(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y})$

Using only coalescent times t:

 $P(N_e(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} | \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} | \mathbf{t}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{t} | N_e(t))}_{\text{Coalescent prior}} P(\mathbf{Q}, N_e(t), \boldsymbol{\tau})$



What if we replace **G** with **t**, the vector of coalescent times? Posteriors:

 $P(N_e(t), \mathbf{G}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y}) \text{ vs } P(N_e(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y})$

Using only coalescent times t:

 $P(N_{e}(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} | \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} | \mathbf{t}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{t} | N_{e}(t))}_{\text{Coalescent prior}} P(\mathbf{Q}, N_{e}(t), \boldsymbol{\tau})$

• $P(\mathbf{Y} | \mathbf{t}, \mathbf{Q}) = \sum_{\mathcal{T}} P(\mathbf{Y}, \mathcal{T} | \mathbf{t}, \mathbf{Q})$ is not "practical" since we need to sum over all possible topologies.



What if we replace **G** with **t**, the vector of coalescent times? Posteriors:

 $P(N_e(t), \mathbf{G}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y}) \text{ vs } P(N_e(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} \mid \mathbf{Y})$

Using only coalescent times t:

 $P(N_{e}(t), \mathbf{t}, \mathbf{Q}, \boldsymbol{\tau} | \mathbf{Y}) \propto \underbrace{P(\mathbf{Y} | \mathbf{t}, \mathbf{Q})}_{\text{Likelihood}} \underbrace{P(\mathbf{t} | N_{e}(t))}_{\text{Coalescent prior}} P(\mathbf{Q}, N_{e}(t), \boldsymbol{\tau})$

- $P(\mathbf{Y} | \mathbf{t}, \mathbf{Q}) = \sum_{\mathcal{T}} P(\mathbf{Y}, \mathcal{T} | \mathbf{t}, \mathbf{Q})$ is not "practical" since we need to sum over all possible topologies.
- The times alone t are not sufficient for inferring N_e(t) for whole genomes [Palacios et al., 2015]

Different resolutions of the coalescent:

Finding the best resolution for the Kingman-Tajima coalescent: theory and applications

R. Sainudiin · T. Stadler · A. Véber

[J. Math. Biol, 2015]

^a considering the optimal resolution with respect to a given statistic can (i) lead to significant computational savings in terms of time complexity by directly sampling from a much smaller hidden space and (ii) help generate samples from the conditional hidden space (given the observed statistics) by controlling the sampling in such a way that only trees or shapes in the hidden space that are compatible with the observed statistics are drawn^{*}



Fig. 2 Example for a ranked labeled tree with leaf label set $\mathfrak{L} = \{1, 2, 3, 4, 5\}$, a labeled tree with $\mathfrak{L} = \{1, 2, 3, 4, 5\}$, a ranked tree shape and a tree shape (from left to right).

Kingman's genealogies vs Tajima's genealogies



Kingman's genealogies vs Tajima's genealogies



Ranked tree shapes

For n = 5



In parenthetical notation, the first tree would be represented by

$$4: (3: (1: (,),), 2: (,))$$
(1)

< ロ > < 同 > < 回 > < 回 >

27/54

Inference with Tajima's coalescent





+

With Kingman's coalescent:

- Goal: $P(N_e(t), \mathbf{G}, \mathbf{Q}, \tau \mid \mathbf{Y})$
- The state space of genealogies \mathcal{G} $\mathcal{G} = \mathcal{T}_n \otimes \mathbb{R}^{+n-1}$
- The likelihood *P*(**Y** | **G**, **Q**) is tractable.

Mutations ² • ATTTCCCCCCA • AATTCCCCCCCA • AATTCCCCCCCA • TTAAGGGGGGTT • TTAAGGGGGGTT

With Tajima's coalescent:

- Goal: $P(N_e(t), \mathbf{G}^T, \mathbf{Q}, \tau \mid \mathbf{Y})$
- The state space of Tajima's genealogies \mathcal{G}^T $\mathcal{G}^T = \mathcal{R}_n \otimes \mathbb{R}^{+n-1}$
- $P(\mathbf{Y} | \mathbf{G}^{T}, \mathbf{Q})$ directly with infinite sites mutations model

Inference with Tajima's coalescent:

- Felsenstein-Tajima conditional likelihood
- Bayesian model for inferring N(t)
- MCMC Algorithm for Posterior inference
 - Sampling of ranked tree shapes (F)
 - Sampling of coalescent times (t)
 - Sampling of *N*(*t*)
- Results
- Summary and future directions

Goal:

$$P(\mathbf{Y} \mid G^T = {\mathbf{F}, \mathbf{t}}, \mu)$$

Assumptions:

- We assume the infinite-sites mutation model
- We know the ancestral state at each polymorphic site
- Our data can be represented as sequences of 0s and 1s
- There is a one-to-one correspondence between data consistent with ISM and a perfect phylogeny (gene tree)
 P(Y | G^T = {F, t}, μ) = P(PP | G^T = {F, t}, μ)



Data as perfect phylogeny (gene tree)

A bigger example (Dan Gusfield, 1991):



| | | _ | _ | | _ | _ | _ | | _ | | | |
|-----------|-----------|---|---|---|---|---|---|---|---|---|---|--|
| Haplotype | Frequency | а | b | b | с | d | е | е | f | f | f | |
| 1 | 2 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2 | 2 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | |
| 3 | 2 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | |
| 4 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | |
| 5 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 6 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 7 | 2 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 8 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 9 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| | 16 | | | | | | | | | | | |

B Perfect Phylogeny (\mathcal{T})





Goal:

$$P(\mathbf{Y} \mid G^T = {\mathbf{F}, \mathbf{t}}, \mu)$$

The tricky part: once we remove the labels, we can have more than one assignment of mutations to branches



With Kingman's coalescent

 $\mathbf{P}\left(\mathbf{Y}\mid\mathbf{G},\boldsymbol{\mu}\right)=\mathbf{P}\left(\mathbf{0},\mathbf{0}\mid\mathbf{G},\boldsymbol{\mu}\right)$



33/54

The computational helper: We use the structure of the perfect phylogeny to generate a probabilistic DAG model (Bayes network, Bayes nets, belief networks, Bayesian graphical model)



We merge sister leaves with the same number of descendants into a single leaf in the DAG.

We equip the DAG with mutations ...



We merge sister leaves with the same number of descendants into a single leaf in the DAG.

We augment the DAG with allocation of mutations (A_i) along \mathbf{G}^T





 $\begin{array}{c} Z_{0} & z \\ Z_{1} & Z_{2} & Z_{3} \\ Z_{4} & Z_{5} & Z_{6} & Z_{7} \\ Z_{8} & Z_{9} & Z_{10} & Z_{11} \\ Z_{10} & Z_{10} & Z_{11} & Z_{2} \\ Z_{10} & Z_{10} & Z_{11} & Z_{2} \\ Z_{10} & Z_{10} & Z_{10} & Z_{11} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} & Z_{10} \\ Z_{10} & Z_{1$

B DAG corresponding to A

$$\begin{array}{l} z_0 = (d_0 = 16, a_0 = (b_5, b_4, b_{14})) \\ z_1 = (d_1 = 7, x_1 = 1, a_1 = ((b_{12}, b_9), b_{10})) \\ z_2 = (d_2 = 7, x_2 = 0, a_2 = ((b_8, b_{13}), b_{11})) \\ z_3 = (d_3 = 2, x_3 = 2) \\ z_4 = (d_4 = 2, x_4 = (1, 1)) \\ z_5 = (d_6 = 2, x_6 = (1, 2)) \\ z_6 = (d_6 = 2, x_6 = (1, 2)) \\ z_7 = (d_7 = 3, x_7 = 1, a_7 = (b_{15}, b_{25})) \\ z_8 = (d_8 = 2, x_8 = 2) \\ z_9 = (d_9 = 1, x_9 = 3) \\ z_{10} = (d_{10} = 2, x_{10} = 1) \\ z_{11} = (d_{11} = 1, x_{11} = 1) \end{array}$$

- Number of mutations

We augment the DAG with allocation of mutations (A_i) along \mathbf{G}^T

 Z_4

 Z_8

 $Z_j =$





B DAG corresponding to A

$$\begin{array}{c} z_0 = (d_0 = 16, a_0 = (b_5, b_4, b_{14})) \\ z_1 = (d_1 = 7, x_1 = 1, a_1 = ((b_{12}, b_3), b_{10}) \\ z_2 = (d_2 = 7, x_2 = 0, a_2 = ((b_8, b_{13}), b_{11}) \\ z_3 = (d_2 = 2, x_3 = 2) \\ z_4 = (d_4 = 2, x_4 = (1, 1)) \\ z_5 = (d_5 = 3, x_5 = 0, a_5 = (b_{16}, b_{24})) \\ z_6 = (d_6 = 2, x_6 = (1, 2)) \\ z_7 = (d_7 = 3, x_7 = 1, a_7 = (b_{15}, b_{25})) \\ z_8 = (d_8 = 2, x_8 = 2) \\ z_9 = (d_9 = 1, x_9 = 3) \\ z_9 = (d_9 = 1, x_9 = 3) \\ z_{10} = (d_{11} = 1, x_{11} = 1) \\ (D_p, X_j) \quad j \in \mathcal{V}_I \end{array}$$

- Number of mutations

1

$$P(\mathbf{Y} \mid G^{T}, \mu) \propto \sum_{A_{0}} \sum_{A_{1}} \dots \sum_{A_{n_{l}}} P(\mathbf{D} \mid G^{T}, \mu) e^{-\mu \mathcal{L}}$$

$$= e^{-\mu \mathcal{L}} \sum_{A_{0}} \sum_{A_{1}} \dots \sum_{A_{n_{l}}} P(Z_{0}, \dots, Z_{n_{l}+n_{L}} \mid G^{T}, \mu)$$

$$= e^{-\mu \mathcal{L}} \sum_{A_{0}} \sum_{A_{1}} \dots \sum_{A_{n_{l}}} \prod_{i=1}^{n_{l}+n_{L}} P(Z_{i} \mid Z_{pa(i)}, G^{T}, \mu)$$

・ロト・雪ト・雪ト・雪・ 今々で

37/54

Bayesian model for inferring N(t)



• $\gamma = \log N(t) \sim GP(0, C(\tau)), \tau \sim Gamma(\alpha, \beta)$

Bayesian model for inferring N(t)



- $\gamma = \log N(t) \sim GP(0, C(\tau)), \tau \sim Gamma(\alpha, \beta)$
- Tajima coalescent prior

$$P[G^{T} = \{\mathbf{F}, \mathbf{t}\} \mid N_{e}(t)] = P(\mathbf{F}) \prod_{k=2}^{n} P[t_{k} \mid t_{k+1}, N_{e}(t)], \quad (2)$$

and

$$P(\mathbf{F}) = \frac{2^{n-c-1}}{(n-1)!},$$
(3)

$$P[t_{k-1}|t_k, N_e(t)] = \frac{C_{k-1}}{N_e(t_{k-1})} \exp\left[-\int_{t_k}^{t_{k-1}} \frac{C_{k-1}dt}{N_e(t)}\right], \quad (4)$$

where $C_k = \binom{k}{2}$ is the coalescent factor that depends on the number of lineages k = 2, ..., n.

39/54

Goal:

$$P[\boldsymbol{\gamma}, \boldsymbol{G}^{T}, \tau \mid \boldsymbol{Y}, \mu] \propto P(\boldsymbol{Y} \mid \boldsymbol{G}^{T}, \mu) P[\boldsymbol{G}^{T} \mid \boldsymbol{\gamma}] P[\boldsymbol{\gamma} \mid \tau] P(\tau)$$
 (5)

Metropolis-Hastings

- splitHMC [Lan et al., Bioinformatics 2015] to sample γ, τ
- HMC to sample t
- *q*(*F* | **Y**) proposal for ranked tree shapes exploiting the DAG representation of perfect phylogeny

Results

Simulation



Only 548 ranked tree shapes are compatible with the data

Simulation: Sampling coalescent times



Posterior of coalescent times

Simulation: Sampling Ne



Simulation

Time

Results

Comparison with BEAST [Suchard et al (2018)]



- Tajima's coalescent is a more **efficient** lower resolution model for inference of effective population sizes.
- A priori, the hidden state space of ranked tree shapes is much **smaller** than the **space** of labeled topologies.
- Current implementation (in R phylodyn) is limited to the infinite-sites mutation model.
- Tajima's inference can be extended for modeling recombination under the sequential Markovian coalescent.
- The Felsenstein-Tajima conditional likelihood calculation can be extended for tree shapes.

- The cardinality of the hidden space of ranked tree shapes depends on your observed data.
- For our simulation example (n = 10 samples)
 - $|\mathcal{T}_{10}| = 2.5 \times 10^9$ labeled topologies

- The cardinality of the hidden space of ranked tree shapes depends on your observed data.
- For our simulation example (n = 10 samples)
 - $|\mathcal{T}_{10}| = 2.5 \times 10^9$ labeled topologies

•
$$|\mathcal{R}_{10}| = 7936$$

- The cardinality of the hidden space of ranked tree shapes depends on your observed data.
- For our simulation example (n = 10 samples)
 - $|\mathcal{T}_{10}| = 2.5 \times 10^9$ labeled topologies
 - $|\mathcal{R}_{10}| = 7936$
 - Only $|\mathcal{R}_{10} | \mathbf{Y}| = 548$ compatible with the data

- The cardinality of the hidden space of ranked tree shapes depends on your observed data.
- For our simulation example (n = 10 samples)
 - $|\mathcal{T}_{10}| = 2.5 \times 10^9$ labeled topologies
 - $|\mathcal{R}_{10}| = 7936$
 - Only $|\mathcal{R}_{10} | \mathbf{Y}| = 548$ compatible with the data

- The cardinality of the hidden space of ranked tree shapes depends on your observed data.
- For our simulation example (n = 10 samples)
 - $|\mathcal{T}_{10}| = 2.5 \times 10^9$ labeled topologies
 - $|\mathcal{R}_{10}| = 7936$
 - Only $|\mathcal{R}_{10} | \mathbf{Y}| = 548$ compatible with the data
- Let $\beta_x = (x + 1, n x + 1)$ be a binary perfect phylogeny with two leaves with x + 1 leaves on one side and n - x + 1on the other side.

- The cardinality of the hidden space of ranked tree shapes depends on your observed data.
- For our simulation example (n = 10 samples)
 - $|\mathcal{T}_{10}| = 2.5 \times 10^9$ labeled topologies
 - $|\mathcal{R}_{10}| = 7936$
 - Only $|\mathcal{R}_{10} | \mathbf{Y}| = 548$ compatible with the data
- Let $\beta_x = (x + 1, n x + 1)$ be a binary perfect phylogeny with two leaves with x + 1 leaves on one side and n - x + 1on the other side.
- The set of unlabeled histories (ranked tree shapes) compatible with β_x, assuming x < n/2 is

$$|R_{n+2} | \beta_x| = e_x \cdot e_{n-x} \cdot \binom{n}{x}, \qquad (6)$$

where e_i is the number of unlabeled histories with *i* internal nodes. The integer e_i is the *i*th Euler number:

$$\sum_{i=0}^{\infty} \frac{e_i z^i}{i!} = \frac{1}{\cos(z)} + \tan(z) \tag{7}$$

Theorem

The maximum number of ranked tree shapes of n leaves compatible with the perfect phylogeny occurs when the perfect phylogeny is a multifurcating phylogeny of degree n or n - 1.

Theorem

The maximum number of ranked tree shapes of n leaves compatible with the perfect phylogeny occurs when the perfect phylogeny is a multifurcating phylogeny of degree n or n - 1.

Theorem

The maximum number of ranked tree shapes of *n* leaves compatible with a binary perfect phylogeny occurs when there is a single bifurcation event dividing the samples in two groups of sizes (n - 2) and (2).

Acknowledgments

- Amandine Veber (Ecole Polytechnique)
- John Wakeley (Harvard University)
- Sohini Ramachandran (Brown University)
- Zhangyuan Wang (Stanford University)
- Noah Rosenberg (Stanford University)
- Filippo Disanto (Stanford University)
- Anand Bhaskar (Facebook)