

Evolving genealogies in neutral populations

Peter Pfaffelhuber

University of Munich,
Joint work with Anita Winter and Andreas Greven (Erlangen)

May 27th, 2007

Neutral models of constant size

Populations of constant size have been **modelled** by

- ▶ **Markov Chains** (Wright-Fisher-model, Moran model)
- ▶ **Diffusion approximations** (Fisher-Wright diffusion)

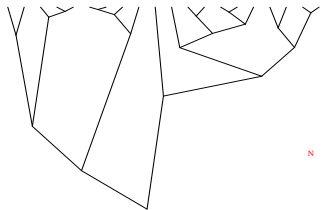
$$dX = \sqrt{X(1-X)}dW.$$

- ▶ **Measure-valued diffusions** (Fleming-Viot superprocess)
Process $(\mu_t)_{t \geq 0}$ with state space $\mathcal{P}(K)$ such that $(\mu_t(A))_{t \geq 0}$ follows a Wright-Fisher diffusion.

Kingman's coalescent

Genealogies relating individuals are known to be distributed according to **Kingman's coalescent**:

- ▶ start with n lines
- ▶ if there are k lines the coalescence rate is $\binom{k}{2}$.
- ▶ stop when reaching **one** line.

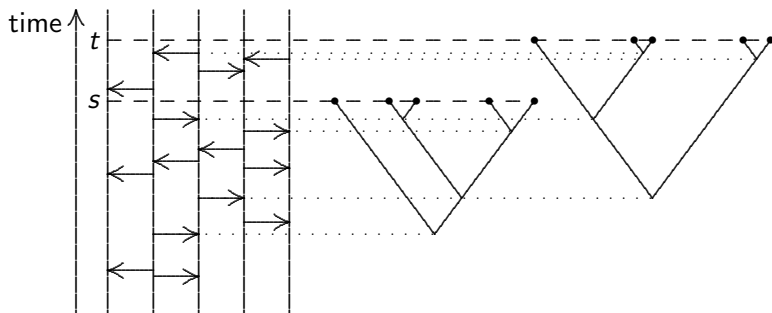


Example: Allelic distribution of a sample of size n at time t is uniquely given by the allelic distribution of ancestors at time 0

Evolving coalescents

Goal: construct a tree-valued stochastic process $\mathcal{U} = (\mathcal{U}_t)_{t \geq 0}$ which

- ▶ describes genealogical relationships **dynamically**
- ▶ makes **forward** and **coalescent** picture implicit



Formalizing genealogical trees

- ▶ **Leaves in Kingman's coalescent** form an ultrametric space.

State space of \mathcal{U} :

$\mathbb{U} := \{\text{isometry class of } (U, r, \mu) : (U, r) \text{ complete and separable } \mathbf{ultrametric} \text{ space, } \mu \in \mathcal{P}(U)\}.$

- ▶ $r(u, v)$ defines the genealogical distance of individuals u and v
- ▶ μ marks currently living individuals

The Gromov-weak topology on \mathbb{U}

- ▶ Π : Algebra generated by functions on \mathbb{U} of the form

$$\Phi(U, r, \mu) := \int \varphi(r(\underline{u}, \underline{u})) \mu^{\otimes n}(d\underline{u})$$

for $\underline{u} = (u_1, \dots, u_n)$, $\varphi \in \mathcal{C}_b(\mathbb{R}^{\binom{n}{2}})$

Π **separates points** in \mathbb{U} .

- ▶ **Gromov-weak topology:**

$$(U_n, r_n, \mu_n) \rightarrow (U, r, \mu) \iff \forall \Phi \in \Pi : \Phi(U_n, r_n, \mu_n) \rightarrow \Phi(U, r, \mu)$$

- ▶ **Proposition** (Gromov; Vershik; Greven, P, Winter)

- ▶ The space $(\mathbb{U}, \mathcal{O}_{\mathbb{U}})$ is **Polish**.

Martingale Problem

- ▶ Given: operator Ω on Π
 $\mathcal{U} = (\mathcal{U}_t)_{t \geq 0}$ solution of the (Ω, Π) -**martingale problem** if for all $\Phi \in \Pi$,

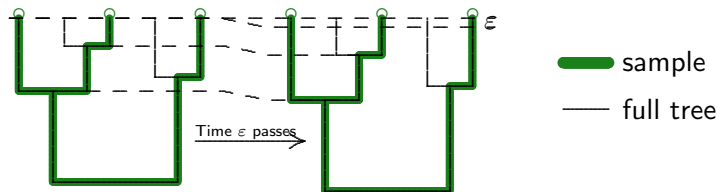
$$\left(\Phi(U_t) - \int_0^t ds \Omega \Phi(U_s) \right)_{t \geq 0}$$

is a martingale. It is **well-posed** if there is exactly one such process.

- ▶ **Program:**
 - ▶ Define pregenerator for finite populations.
 - ▶ Establish **existence** by generator convergence and tightness.
 - ▶ Establish **uniqueness** by duality.

Tree Growth

When no resampling occurs the tree **grows**



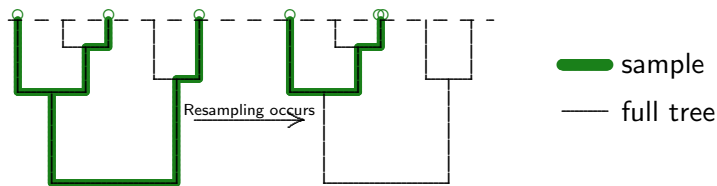
Distances in the sample grow

$$\Omega_{\text{grow}}^N \Phi(U, r, \mu) = \sum_{i,j} \int \frac{\partial \varphi}{\partial r_{ij}} (r(\underline{u}_i, \underline{u}_j)) \mu^{\otimes n}(d\underline{u}) + \mathcal{O}\left(\frac{1}{N}\right)$$

- ▶ $\frac{\partial \varphi}{\partial r_{ij}}$: measures the **change of φ** when $r(u_i, u_j)$ grows
- ▶ $\mathcal{O}\left(\frac{1}{N}\right)$: due to double sampling of the same individual

Resampling

Resampling does not change the tree



By resampling μ changes, so **different samples are picked**.

$$\Omega_{\text{res}}^N \Phi(U, r, \mu) = \frac{1}{2} \sum_{k,l} \int (\varphi(\theta_{kl} r(\underline{u}, \underline{u})) - \varphi(r(\underline{u}, \underline{u}))) \mu^{\otimes n}(d\underline{u}) + \mathcal{O}\left(\frac{1}{N}\right)$$

$$\text{where } \theta_{kl} r(u_i, u_j) = \begin{cases} r(u_k, u_j), & i = l \\ r(u_i, u_k), & j = l \\ r(u_i, u_j), & \text{else} \end{cases}$$

Martingale Problem

- ▶ **Pregenerator** for infinite system:

$$\begin{aligned}\Omega\Phi(U, r, \mu) &= \sum_{i,j} \int \frac{\partial\varphi}{\partial r_{i,j}}(r(\underline{u}, \underline{u})) \mu^{\otimes n}(d\underline{u}) \\ &\quad + \frac{1}{2} \sum_{k,l} \int (\varphi(\theta_{kl}r(\underline{u}, \underline{u})) - \varphi(r(\underline{u}, \underline{u}))) \mu^{\otimes n}(d\underline{u})\end{aligned}$$

- ▶ Existence: **Approximation** by finite systems
- ▶ Uniqueness: **Duality** to Kingman's coalescent

The tree-valued Fleming-Viot dynamics

Theorem

- ▶ Solution $\mathcal{U} = (\mathcal{U}_t)$ of (Ω, Π) -mp exists as limit of tree-valued Moran models and is unique
- ▶ Almost surely,
 - ▶ \mathcal{U} has **continuous** sample paths
 - ▶ \mathcal{U}_t **compact** for all $t > 0$
 - ▶ **Quadratic variation** of $\Phi(\mathcal{U})$:

$$d\langle \Phi(\mathcal{U}) \rangle_t = n^2 \langle \mu_t, (\rho - \langle \mu_t, \rho \rangle)^2 \rangle dt,$$

where

$$\rho(u_1) := \int \mu_t^{\otimes (n-1)} (d(u_2, \dots, u_n) \phi((r(u_i, u_j))_{1 \leq i < j \leq n}))$$

The tree-length process

- ▶ Given (U, r, μ) and u_1, u_2, \dots ,

$$L_n(\underline{u}) := \text{length of tree spanned by } u_1, \dots, u_n$$

- ▶ **Tree length distribution** of subsequentially sampled points:

$$\Lambda(U, r, \mu) := (L_2, L_3, \dots)_* \mu^{\otimes \mathbb{N}} \in \mathcal{P}(\mathbb{R}_+^{\mathbb{N}})$$

- ▶ **Theorem:** The process $(\Lambda(\mathcal{U}_t))_{t \geq 0}$ is **Markov**

Application: number of mutations on the tree

- ▶ Real data: **number of mutations** $S_n(t)$ of a sample at time t can be observed
- ▶ Length of sample tree $L_n(\underline{u}) \Rightarrow S_n \sim \mathbf{Pois}(\frac{\theta}{2}L_n(\underline{u}))$
- ▶ So,

$$\mathbb{E}[e^{-\lambda S_n(t)}] = g^n(t; \theta(1 - e^{-\lambda})).$$

for

$$g^n(t; \theta) := \mathbb{E}\left[\int \mu_t^{\otimes n}(d\underline{u}) \exp(-\theta L_n(\underline{u}))\right]$$

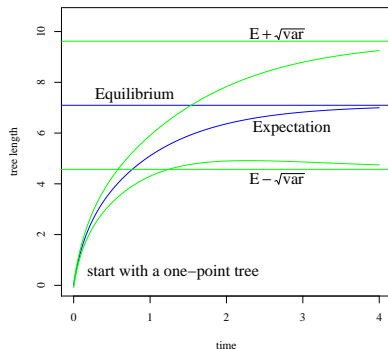
- ▶ (g^2, g^3, \dots) is solution of

$$\frac{d}{dt}g^n = -n\theta g^n + \binom{n}{2}(g^{n-1} - g^n)$$

- (g^2, g^3, \dots) given by

$$g^n(t; \theta) = \Gamma(n) \sum_{k=2}^n \frac{\binom{n}{k} (-1)^k (\theta + 2k - 1)}{\Gamma(\theta + n + k)} \cdot \left\{ e^{-k(\theta + (k-1))t} \sum_{m=2}^k \frac{\binom{k}{m} (-1)^m \Gamma(\theta + k + m - 1)}{\Gamma(m)} g^m(0; \theta) + (1 - e^{-k(\theta + (k-1))t}) (k-1)(\theta + k) \Gamma(\theta + k - 1) \right\}.$$

Application: tree lengths



Example: consider a sample of $n = 20$ which starts in $(\{\bullet\}, \delta_\bullet)$. Moments of tree lengths can be calculated.

Outlook

Further questions

- ▶ Is it possible to include **mutation, selection, recombination**?
- ▶ What does the mp for tree-valued **branching** processes look like?
- ▶ Which tree-evolutions can be defined via **sample evolution**?