Some problems with genealogies

Simon C. Harris

University of Auckland New Zealand

Based on joint work with: Samuel Johnston (Dublin) & Matt Roberts (Bath)

> CIMAT, Mexico August 24th 2018

Our fundamental question

What does the ancestral family tree look like for a sample of individuals chosen from some population?

The question needs more precise formulation:

- a suitable model for the population's evolution (eg. fixed or randomly varying population size?)
- a sampling mechanism (eg. choose uniformly at random from entire population?)

A question of genealogy

Our fundamental question

What does the ancestral family tree look like for a sample of individuals chosen from some population?

The question needs more precise formulation:

- a suitable model for the population's evolution (eg. fixed or randomly varying population size?)
- a sampling mechanism (eg. choose uniformly at random from entire population?)



A question of genealogy

Our fundamental question

What does the ancestral family tree look like for a sample of individuals chosen from some population?

The question needs more precise formulation:

- a suitable model for the population's evolution (eg. fixed or randomly varying population size?)
- a sampling mechanism (eg. choose uniformly at random from entire population?)



S.C.Harris (University of Auckland)

Some problems with genealogies

A population model

Offspring distribution. *L* is a random variable taking values in $\mathbb{Z}^+ := \{0, 1, 2, ...\}$.

$$p_k := \mathbb{P}(L=k), \qquad m := \mathbb{E}(L) = \sum_{i=0}^{\infty} i p_i < \infty.$$



Continuous time Galton-Watson process

 N_t represents the number of individuals in a population that are alive at time $t \ge 0$.

- any individual alive branches at rate r > 0, independently of others;
- when it branches, the parent dies and is replaced by a random number of offspring given by an independent realisation of offspring distribution L;
- once born, offspring evolve independently as above, and so on.

Some more basic properties of GW processes

Super-critical case (m > 1): Kesten-Stigum theorem. Suppose $1 < m < \infty$

If $E(L \log L) < \infty$, martingale $Z_t := e^{-r(m-1)t}N_t \to Z_\infty$ a.s. & in L^1 with $\mathbb{E}(Z_\infty) = 1$. Further, $\{Z_\infty > 0\} = \{N_t \ge 1, \forall t \ge 0\}$ a.s., $\Rightarrow \mathbb{P}(\text{survival}) = \mathbb{P}(Z_\infty > 0) > 0$. When the population survives it grows exponentially: $N_t \sim Z_\infty e^{r(m-1)t}$ a.s.

Critical case (m = 1): Conditioning on surviving a long time.

- $\mathbb{P}(N_t > 0) \sim \frac{c}{t} \to 0$
- Critical Yaglom theorem. Let $\sigma^2 := \mathbb{E}(L^2) 1$. For each x > 0

$$\mathbb{P}\left(\frac{N_t}{t} > x \middle| N_t > 0\right) \to \exp\left(-\frac{2}{r\sigma^2}x\right) \quad \text{as } t \to \infty. \quad \text{Huge population} \approx \mathbf{t} \times \mathbf{RV}$$

Sub-critical case (m < 1): Conditioning on surviving a long time.

• $\mathbb{P}(N_t > 0) \sim c\mathbb{E}(N_t) = ce^{-r(1-m)t} \rightarrow 0$

• Sub-critical Yaglom theorem. For $\theta > 0$, for some RV Z taking values in \mathbb{Z}^+

 $\mathbb{E}(e^{-\theta N_t}|N_t>0) \to E(e^{-\theta Z}) \quad \text{as } t\to\infty. \qquad \text{Finite population}$

Uniform sampling from a Galton-Watson process

Our question

Consider a continuous time Galton-Watson process along with its genealogical tree.

- Fix a large time T > 0 (we will sometimes let $T \to \infty$)
- Condition on the event that at least k individuals are alive at time T.
- Choose k individuals uniformly at random (without replacement) from those alive at time T.

What does the ancestral tree drawn out by this sample of k individuals look like?



Existing literature. Critical case (k = 2): see Durrett (1978) and Athreya (2012). Near critical (k = 2), see O'Connell (1995). Also see related works by Aldous & Popovic (2005), Lambert & Stadler (2013), Lambert (2003, 2010 & recent...)

S.C.Harris (University of Auckland)

Some problems with genealogies

A small taste: birth-death process (k = 2)

Birth-death process: individuals branch into two at rate $\beta > 0$ and die at rate $\alpha > 0$ Offspring distribution $\mathbb{P}(L = 0) = \frac{\alpha}{\alpha + \beta}$, $\mathbb{P}(L = 2) = \frac{\beta}{\alpha + \beta}$ Branching rate $r := \alpha + \beta$ Mean $m := \mathbb{E}(L) = \frac{2\beta}{(\alpha + \beta)}$

Conditional on $\{N_T \ge 2\}$, let S_T be the **time of most recent common ancestor** of a pair of individuals $U_T, V_T \in N_T$ chosen *uniformly at random (without replacement)*.

Supercritical birth-death ($\beta > \alpha$) - related near start!

The law of S_T conditional on $N_T \ge 2$ converges to a non-trivial distribution as $T \to \infty$, with tail satisfying

$$\lim_{T\to\infty}\mathbb{P}(S_{\mathcal{T}}\geq s\,|\,N_{\mathcal{T}}\geq 2)\sim 2(\beta-\alpha)s\,e^{-(\beta-\alpha)s}$$



Sub-critical birth-death ($\beta < \alpha$) - related near end!

The law of $T - S_T$ conditional on $N_T \ge 2$ converges to a non-trivial distribution as $T \to \infty$, with tail satisfying

$$\lim_{t\to\infty} \mathbb{P}(T - S_T \ge s \mid N_T \ge 2) \sim \left(1 - \frac{2\beta}{3\alpha}\right) e^{-(\alpha - \beta)s}$$

S.C.Harris (University of Auckland)

A bigger taste: general critical GW process (k = 2)

General critical Galton-Watson process.

Assume offspring distribution L satisfies $m = \mathbb{E}(L) = 1$ and $\mathbb{E}(L^2) < \infty$.

On the event $\{N_T \ge 2\}$, let S_T be the **time of most recent common ancestor** of a pair of individuals $U_T, V_T \in N_T$ chosen *uniformly at random (without replacement)*.

Critical GW case (m = 1) - related anywhere!

The law of S_T/T conditional on $N_T \ge 2$ converges as $T \to \infty$ to a non-trivial distribution on [0, 1] satisfying

$$\lim_{T \to \infty} \mathbb{P}\Big(\frac{S_T}{T} \ge t \ \Big| \ N_T \ge 2\Big) = \frac{2(1-t)}{t^2} \Big(\log\Big(\frac{1}{1-t}\Big) - t\Big)$$



Harris-Johnston-Roberts (2017+) also give genealogy of k individuals explicitly for:

- super-critical birth-death processes at *fixed* times T
- general critical and *near*-critical GW processes as $T \to \infty$.

Further results: also see Johnston (2017+)



Genealogy of uniform sample of k individuals: We want to jointly characterise:

- times of the k-1 mergers of family lines, $(S_1^k(T), S_2^k(T), \ldots, S_{k-1}^k(T));$
- shape of the genealogical tree.

Critical GW scaling limit: genealogy of k individuals Harris-Johnston-Roberts ('17)

As $T \to \infty$, the scaled merger times $\left(\frac{S_1^k(T)}{T}, \ldots, \frac{S_{k-1}^k(T)}{T}\right)$ conditional on $N_T \ge k$, converge in distribution to $(S_1^k, \ldots, S_{k-1}^k)$ where, for any $0 \le s_1 \le \cdots \le s_{k-1} \le 1$,

$$\mathbb{P}(S_1^k \ge s_1, \dots, S_{k-1}^k \ge s_{k-1}) = k! \prod_{i=1}^{k-1} \frac{s_i - 1}{s_i} - k! \sum_{j=1}^{k-1} \frac{1 - s_j}{s_j^2} \left(\prod_{\substack{i=1 \ i \neq j}}^{k-1} \frac{1 - s_i}{s_j - s_i}\right) \log(1 - s_j)$$

As $T \to \infty$, the **shape of the tree** is asymptotically independent of the merger times where each pair of family lines is equally likely to be the one to coalesce at the next merger time.

IMPORTANT: as $\mathbb{E}L = 1$ and $\mathbb{E}L^2 < \infty$, only pairwise mergers observed in limit!



NB. This is *same* tree topology as Kingman coalescent...

Kingman coalescent. Every pair of family lines merges at rate 1. That is, if there are currently *i* separate family lines, the next merger occurs at rate i(i-1)/2 and is equally likely to be any two of the family lines that merge.

In fact, to get the explicit distribution function we first found the following density. Critical GW scaling limit: density for genealogy of k individuals

As $T \to \infty$, the scaled merger times $\left(\frac{S_k^k(T)}{T}, \dots, \frac{S_{k-1}^k(T)}{T}\right)$ conditional on $N_T \ge k$, converge in distribution to $(S_1^k, \dots, S_{k-1}^k)$ where, for any $0 \le s_1 \le \dots \le s_{k-1} \le 1$,

$$\mathbb{P}(S_1^k \in \mathrm{d} s_1, \dots, S_{k-1}^k \in \mathrm{d} s_{k-1}) = \int_0^\infty \frac{k}{(1+\theta)^2} \left(\prod_{i=1}^{k-1} \frac{\theta}{(1+\theta(1-s_i))^2} \, \mathrm{d} s_i\right) \mathrm{d} \theta$$

Limiting coalescent times as a mixture of IID times:

First choose a mixture random variable M_k then, conditional on $M_k = \theta$, the k - 1 coalescent times are IID on [0, 1] with a density depending on θ .



A construction of limiting coalescent times:
In fact, (S₁^k,..., S_{k-1}^k) can be constructed by:
Let X₁, X₂,..., X_k be a sequence of IID RVs on (0,∞) with PDF (1 + x)⁻².

- Renormalise X_1, \ldots, X_k by the maximum $M_k := \max\{X_1, \ldots, X_k\}$
- Ignoring the maximum value 1, the remaining ordered k 1 renormalised RVs have the same distribution as (1 S₁^k,..., 1 S_{k-1}^k).

Limiting coalescent times as a mixture of IID times:

First choose a mixture random variable M_k then, conditional on $M_k = \theta$, the k - 1 coalescent times are IID on [0, 1] with a density depending on θ .



A construction of limiting coalescent times:
In fact, (S₁^k,..., S_{k-1}^k) can be constructed by:
Let X₁, X₂,..., X_k be a sequence of IID RVs on (0,∞) with PDF (1 + x)⁻².

- Renormalise X_1, \ldots, X_k by the maximum $M_k := \max\{X_1, \ldots, X_k\}$
- Ignoring the maximum value 1, the remaining ordered k 1 renormalised RVs have the same distribution as (1 S₁^k,..., 1 S_{k-1}^k).

Kingman coalescent often appears as limit when population is *constant*, but GW population size *varies* randomly...

A 'slow' Kingman, with mergers only at rate i - 1 when there are i individuals, would have coalescent times of k individuals given by k - 1 IID exponentials.

Genealogy of uniform sample from a critical GW is like a mixture of time changed 'slow' Kingman coalescents. Interpret M_k like biologists' *effective population size*...

uniform chorice of K GW process with & spines watescence KSpiner. branching 2.2.2 = 8 · GN process under P · <u>Extend</u> to GW tree with k spries under IP^(K): 22.23-6 (1) Choose & spines independently (1) Af birth events each spine g. IP(u, u, u, e)(4,) = 1.1.1 follows path uniformly from offspring IP (u e P (M) = 1 II II LV U,,, un are spinis drosen i=1 VKUL'A Probabilitio of spine durices follspring oncestors of Ui at node V

A useful change of measure

[Skespines alviele] IT II Lv (district at huit) spines VKS. dip(w) dip(w) ay(w) H_{2} $\left(\mathcal{N}_{\tau} (\mathcal{N}_{\tau} - i) \dots (\mathcal{N}_{\tau} - k + i) \right)$ spive + Ger mformation GW information Note: ZT K Note: T K prives dishut (1 TT L M) Note: T K sprives dishut (1 TT L M) N K >= TTLV = S P(ues at). Is u dohad II TI Lv u, uen T = N_T (N_T-()... (N_T-KF()) # choices of & ridudend from N_T inthout replacement

E (by areful choice (upportant feature of R (") QT (UES (JT) ~ PT (UES (JT) × ISu dutinet, dire TIT Lo at the T dire to used So, $Q_{\tau}^{(u)}(y \in \mathbb{F}[Y_{\tau}) = I[u_{1}, \dots u_{u} \in N_{\tau} \text{ hdistinct}]$ total through and C-W into (no spire identity) No (Ny-1) ... (Ny-k+1) ~ - choosing K from Ny inthat replacement : under QT, at time T the sphes are uniformy chosen inthout replacement from those alive at fine T functional along - functional along priver (out) Also find, tree of particles Un, -, Uk (Enly) $\mathbb{Q}_{\tau}^{(k)}\left(f\left(\frac{s}{r}\right)\left(\frac{s}{r}\right)=\frac{1}{N_{\tau}(N_{\tau}-1)\dots(N_{\tau}-k+1)}\sum_{\substack{\mathsf{N}_{\tau},\dots,\mathsf{U}_{u}\in\mathsf{N}_{\tau}\\\mathsf{V}_{i}\,\mathrm{dishind}}}\sum_{\substack{\mathsf{N}_{\tau},\dots,\mathsf{U}_{u}\in\mathsf{N}_{\tau}\\\mathsf{V}_{i}\,\mathrm{dishind}}}\sum_{\substack{\mathsf{N}_{\tau},\dots,\mathsf{V}_{u}\in\mathsf{N}_{\tau}\\\mathsf{V}_{i}\,\mathrm{dishind}}}$

Coalescent trie for Uniform samples in critical GU het Ti, Tz,..., The be the conferent finer of k particles chosen uniformly at random without replacement. Then: IP (T, edt, j...; Tk, EdSk, N, 2 k and choose k alwe at time T uniformly at random inthrout replacement) $= IP\left(\underbrace{I}_{N_{T}(N_{T}-1)...}(N_{T}-k+1)}\sum_{\substack{u_{1},...,u_{k}\in N_{T}}} I\left\{ \overline{\zeta_{k}^{u}} \in d_{x}, \right\} | N_{T} \ge k \right)$ $= Q_{\tau}^{(k)} \left(\frac{\prod \{Z_{i}^{2} \in dt_{i}\} \cdots \{Z_{k}^{2} \in dt_{k}\}}{N_{\tau}(N_{\tau}-i) \cdots (N_{\tau}-k+i)} \right) \frac{IE(N_{\tau}(N_{\tau}-i) \cdots (N_{\tau}-k+i))}{IP(N_{\tau} \ge k)}$ Com describe sprie behander completely under (DT) (Easy using Yaylom)

Application of Yaglom's theorem
Critical GW Offerning dust "L, IEL=1, o²:=War(L)<0,
Branching rate r
(a)
$$P(N_{\tau}>0) \sim \frac{\lambda}{T}$$
 as $T \rightarrow \infty$
(b) Conditional on $N_{\tau}>0$, $\frac{N_{\tau}}{T} \stackrel{2aw}{\rightarrow} Z \sim Exp(\lambda)$ as $T \rightarrow \infty$
• $E(e^{-PN_{\tau}}|N_{\tau}>0) \rightarrow IE(e^{-PZ}) = \frac{\lambda}{\lambda+Q}$ as $T \rightarrow \infty$
• $E(N_{\tau}(N_{\tau})...(N_{\tau}-k+1)|N_{\tau}>0) \rightarrow IE(Z^{k}) = \frac{k!}{\lambda^{k}}$
• $P(N_{\tau} \rightarrow \infty|N_{\tau}>0) = 1$, so $P(N_{\tau} \geq k) \sim P(N_{\tau}>0) \sim \frac{\lambda}{T}$
Hence, $IE(N_{\tau}(N_{\tau}r)...(N_{\tau}-k+1)|N_{\tau}>k) \sim T^{k} \frac{k!}{k^{k}}$

Understanding QT for large times $\cdot \mathcal{Q}_{\tau}^{(\mu)}\left(e^{-\mathcal{O}N_{\tau}}\right) = \frac{E\left(N_{\tau}\left(N_{\tau^{-1}}\right)...\left(N_{\tau^{-k+1}}\right)e^{-\mathcal{O}N_{\tau}}\right)}{E\left(N_{\tau}\left(N_{\tau^{-1}}\right)...\left(N_{\tau^{-k+1}}\right)e^{-\mathcal{O}N_{\tau}}\right)}$ Total population size under QT $IE(N_{\tau}(N_{\tau}-i)-(N_{\tau}-k_{\tau}-i))$ $\rightarrow I \underbrace{E(2^{k} - 0^{2})}_{I \underbrace{E(2^{k})}} = \begin{pmatrix} \lambda \\ \lambda + 0 \end{pmatrix}^{k+1} since \underbrace{N_{T}}_{T} \rightarrow Z \sim Exp(\lambda)$ under IP.. No under Q tends to Z~ Gamma(k+1, 2) under Q_0(k) equivalently, a sum of kt (, ridependent Exp(?)'s. (n particular, mass coming off single spine branch is $\Gamma(2, \lambda)$ -> We need joint distribution of (scaled) spine split times & (scaled)population size Z under Qo

Total population under Qa Time A construction of population under On $N_{\pm} \rightarrow Z \sim \Gamma(k_{\pm},\lambda)$ (c.f. Aldons' anst. Unitom Potal population size: Sample of Kingman coales) iduduets $Z = Z_{o} + \sum (1 - T_{i}) Z_{i}$ NB. Only splits into Zyroups in winit. k spines Coalescent fines: ti, tin, the IID $(1-\hat{\tau}_{2})\hat{z}_{2}$ ~Unit[0,1] U2 [, J'm Subpopulations: (-こ)え 2., 2, ..., 2, IID (い- てょ) える U3 61 $\sim \Gamma(2, \lambda)$ (1-24)24 U4 NB. (1-2) 2~ Exp(X) 2, 24 2 , O

Calculating the Roo expectation Then, as T-200, $\mathbb{Q}_{\tau}^{(\mathsf{W})} \left(\frac{\mathsf{T}^{\mathsf{k}}}{\mathsf{N}_{\tau}(\mathsf{N}_{\tau}-\mathsf{I})\dots(\mathsf{N}_{\tau}-\mathsf{k}+\mathsf{I})}; \frac{\mathsf{T}_{\tau}^{2}}{\mathsf{T}} edt_{i}; \ldots; \frac{\mathsf{T}_{k}^{2}}{\mathsf{T}} edt_{k,\tau} \right)$ $\mathsf{Trele:}$ $= Q_{\infty}^{(k)} \left(\frac{1}{2^{k}}; \hat{T}_{i} \in dt_{i}; ...; \hat{T}_{k-i} \in dt_{k-i} \right)$ $= Q_{\infty}^{(k)} \left(\hat{T} \in dt \right) \int_{0}^{\infty} \frac{O^{k-i}}{(k-i)!} Q_{\infty}^{(k)} \left(e^{-O2} \left(\overline{T} = t \right) d\theta \right)$ $= Q_{\infty}^{(k)} \left(\hat{T} \in dt \right) \int_{0}^{\infty} \frac{O^{k-i}}{(k-i)!} Q_{\infty}^{(k)} \left(e^{-O2} \left(\overline{T} = t \right) d\theta \right)$ $dt_{1}.dt_{2}...dt_{k-1}$ $Z = Z_{0} + \sum_{i=1}^{k-1} (i-t_{i})Z_{i}$ as $T_{i},...,T_{k-1}$ IID $U_{0,i}F[0,1]$ $Z_{0,...,Z_{k-1}}$ IID $P(2,\lambda)$ $= dt_{1}.dt_{2}...dt_{k-1}.\frac{1}{(k-1)!}\int O^{k-1}\left(\frac{\lambda}{\lambda+0}\right)^{2}\frac{k-1}{(k-1)!}\left(\frac{\lambda}{\lambda+0(1-k-1)}\right)^{2}dO$

Coalescent times for uniform sample in critical GW Now order the IP (Tiedtij...; Triedtur | uniform choice of Kindividuals at The The No ≥ R coalescent finies so $T_1 \leq T_2 - \cdot \cdot \leq T_{k-1}$ $\rightarrow IE(Z^{k}) \mathbb{R}_{\infty}^{(n)} \left(\frac{1}{Z^{k}}; \hat{\mathcal{C}}_{i}edt_{i}; ...; \hat{\mathcal{C}}_{u_{n}}edt_{u_{n}} \right)$ = $(k-1)! \frac{k!}{\lambda k} dt_{1}.dt_{2}...dt_{k-1}! \int_{0}^{\infty} O^{k-1} \left(\frac{\lambda}{\lambda + 0}\right)^{2} \frac{k-1}{(k-1)!} \left(\frac{\lambda}{\lambda + 0}\right)^{2} d\theta$ = k! $\int (\lambda + 0)^2 \prod_{i=1}^{k-1} \frac{\partial \lambda}{(\lambda + 0(i-t_i))^2} d\theta \cdot dt_{i} \dots dt_{k-1}$ (0=24) Can Compute $= |k| \int_{0}^{\infty} \frac{1}{(1+q)^{2}} \prod_{i=1}^{\infty} \frac{1}{(1+q)(1-t_{i})^{2}} dt dt_{i} dt_{i} = 0$ Explicitly!