

Stochastic processes and inference in population genetics

①

explosion of sequencing data

human genome: $\sim 3 \times 10^9$ bases $\sim \$1000$

bacterial genome: $\sim 3 \times 10^6$ bases $\sim \$10$

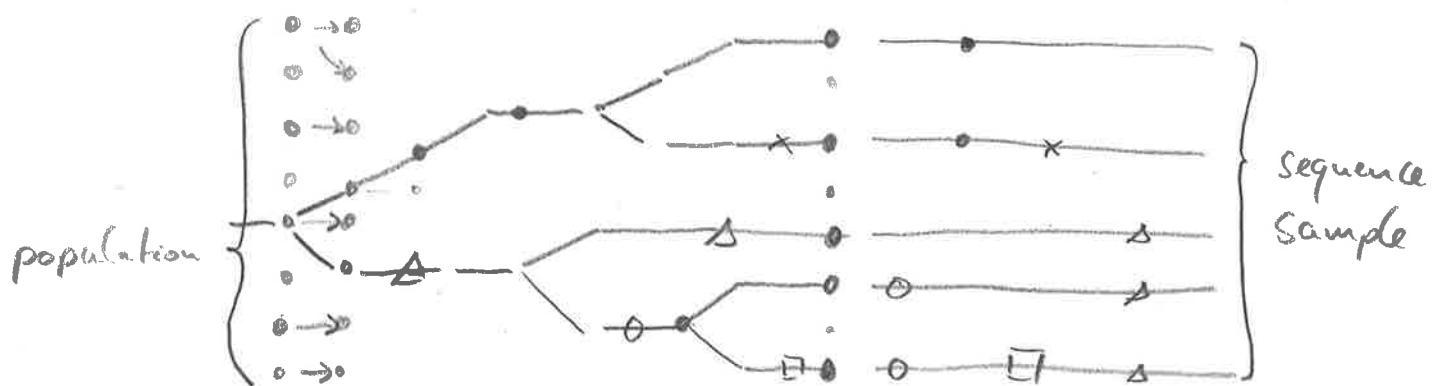
virus genome: $\sim 10^4$ bases

→ 1000s of human genomes are sequenced every day

→ what do we learn from all those data?

Population genetics

What determines and shapes genetic diversity?



- locally, sequences are related by a tree
- genetic diversity depends on the tree
- the tree depends on how one generation seeds the next
- several simple models are used as null models for inference

Coalescent model

(2)

- backwards in time
- explicit model of trees
- offspring distribution Ω_m ($\text{const pop} \Leftrightarrow \langle m \rangle = 1$)
- a randomly chosen pair has the same parent with prob:

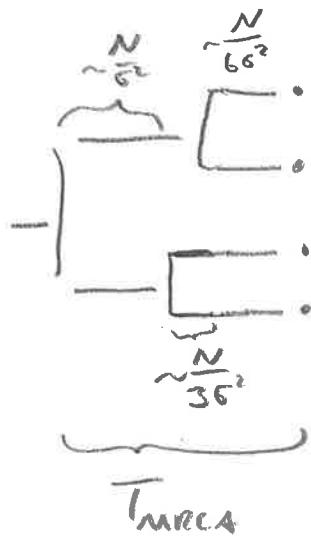
$$N \sum_m \frac{m(m-1)}{N^2} \Omega_m = \frac{\sigma^2}{N}$$

- sample of $k \rightarrow \frac{k(k-1)}{2}$ possible pairs

- \rightarrow merger with rate $\frac{\sigma^2 k(k-1)}{2N}$

$$\rightarrow k \rightarrow k-1$$

$$\rightarrow \dots \rightarrow 1$$



$$\cdot \langle \bar{T}_{\text{MRCAs}} \rangle = \frac{2N}{\sigma^2} \sum_{n=2}^N \frac{1}{k(k-1)}$$

$$= \frac{N}{\sigma^2} \left[1 + \frac{1}{2} + \frac{1}{3} + \dots \right]$$

$$= \frac{2N}{\sigma^2} \frac{n-1}{n} \quad n = \text{sample size}$$

- average pairwise difference: $2\mu \langle \bar{T}_2 \rangle = 2\mu N$

- # of mutations in sample: $\mu \langle \bar{T}_{\text{total}} \rangle = \frac{N}{\sigma^2} \sum_n k \bar{T}_n$

$$= \frac{N}{\sigma^2} \sum_k \frac{1}{k-1}$$

$$\approx \frac{N}{\sigma^2} \log n$$

(3)

Diffusion models

- a dual process - forward in time
- considers the frequency of mutations
- if a mutation has freq. p_t at time t , what is p_{t+1}
 - depends on offspring distribution
 - here, we assume the mutation has no effect

$$\langle p_{t+1} \rangle = \langle p_t \rangle \quad \langle Sp_{t+1} \rangle = \frac{5^2}{N} p_t (1-p_t)$$

→ essentially binomial resampling of the mutation.

$$\cdot \frac{\partial}{\partial t} P(p, l) = \frac{5^2}{2N} \frac{\partial}{\partial p^2} p(1-p) P(p, l)$$

$$\cdot \text{separation of variables: } P(p, l) = a(p) b(l)$$

$$\rightarrow \frac{1}{b(l)} \frac{\partial}{\partial t} b(l) = -\lambda = \frac{5^2}{2N a(p)} \frac{\partial}{\partial p^2} p(1-p) a(p)$$

→ hypergeometric equation

→ BC $a(0)$ & $a(1)$ finite → discrete spectrum

$$\lambda = \frac{5^2 h(h+1)}{N} \quad \text{Same as coalescent rate?}$$

$$P(p, t) = \sum_h c_h T_{h,1}(p) e^{-\frac{5^2 h(h+1)}{N} t}$$

↳ Gegenbauer polynomial.

Branching processes

- stochastic dynamics of rare mutations

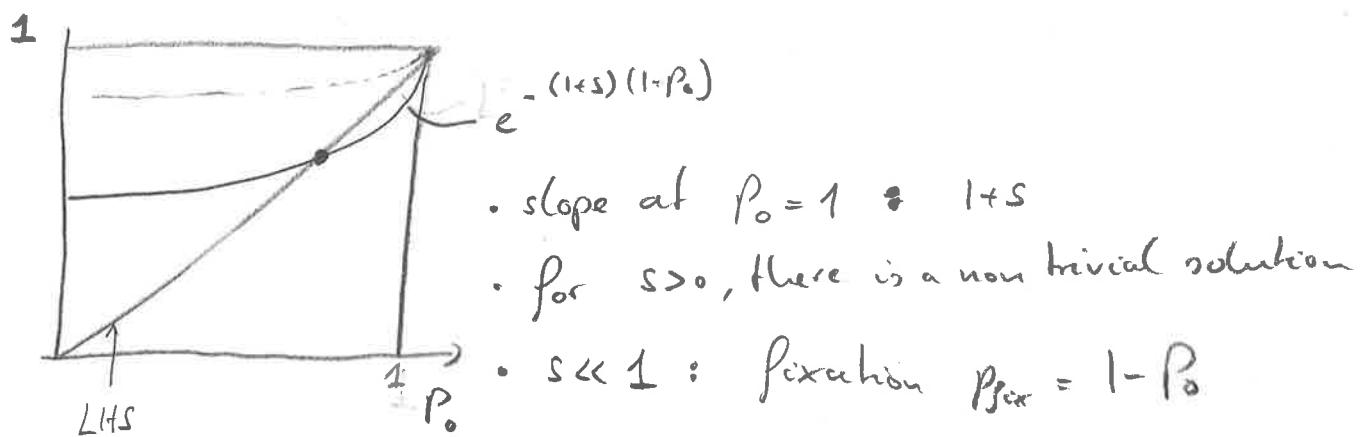
Q: how frequently is a beneficial mutation lost?

assume offspring distribution is Poisson with mean $1+s$

$$P_0 = \sum_{m=0}^{\infty} \frac{e^{-(1+s)}}{m!} (1+s)^m P_0^m = e^{-(1+s)(1-P_0)}$$

loss if all offspring are lost

graphic solution



$$P_{fix} \approx 2s$$

- more generally: what is the distribution of the copy number n at time t of a clone with birth rate b & death rate d ?

$$P(n, t+dt) = (1 - dt(b+d))_n P(n, t) + (n-1)b P(n-1, t) + (n+1)d P(n+1, t)$$

$$\Rightarrow \frac{\partial}{\partial t} P(n, t) = -(b+d)_n P(n, t) + b(n-1) P(n-1, t) + d(n+1) P(n+1, t) \quad | \times \sum_n z^n$$

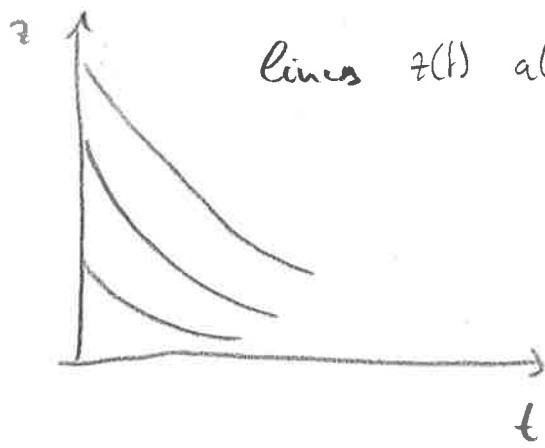
$$\frac{\partial}{\partial t} \sum_{n=0}^{\infty} z^n P(n, t) = \frac{\partial}{\partial t} \phi(z, t) = -(b+d)z \sum_{n=0}^{\infty} \partial_z z^n P_n + b z^2 \sum_{n=0}^{\infty} \partial_z z^{n-1} P_{n-1} + d \sum_{n=0}^{\infty} \partial_z z^{n+1} P_{n+1}$$

generating function

$$= [-(b+d)z + bz^2 + dz] \partial_z \phi$$

- PDE that can be solved by characteristics

⑤



lines $z(t)$ along which $\phi(z,t)$ is constant

$$\begin{aligned}\frac{d}{dt} \phi(z(t), t) &= \frac{dz}{dt} \frac{\partial \phi}{\partial z} + \frac{\partial \phi}{\partial t} \\ &= \frac{\partial \phi}{\partial z} \left[\frac{dz}{dt} - (b+d)z + b^2 + d \right] \\ &= 0\end{aligned}$$

the characteristic equation simplifies when substituting

$$x = 1-z$$

$$\begin{aligned}\frac{dx}{dt} &= -(b+d)(-x) + b(1-x)^2 + d \\ &= -(b-d)x + bx^2 \\ &= -sx + bx^2 \\ \Rightarrow x(t) &= \frac{e^{-st}}{C - b \int_0^t e^{-st'} dt'} = \frac{e^{-st}}{C + \frac{b}{s} [e^{-st} - 1]}\end{aligned}$$

$$\text{We want: } \phi(z, t)$$

we have: $x(t)$ along which $\phi(1-x(t), t)$ is constant

$$\phi(z, 0) = z$$

we need: $C(z, t)$ such that $x(t)$ goes through $(1-x, t)$

$$z = 1 - \frac{e^{-st}}{C(z, t) + \frac{b}{s}(e^{-st} - 1)}$$

$$\Rightarrow \phi(z, t) = 1 - \frac{1}{C(z, t)} = 1 - \frac{(1-z)e^{st}}{1 + (1-z)\frac{b}{s}[e^{st} - 1]}$$

Sanity checks

(6)

$$\phi(1,t) = \sum_{n=0}^{\infty} P(n,t) = 1 \quad \checkmark$$

$$\langle n \rangle = \sum_{n=0}^{\infty} n P(n,t) = \partial_z \phi(z,t) \Big|_{z=1} = e^{st} \quad \checkmark$$

Full distribution

$\phi(z,t)$ can be expressed as an geometric series

$$\frac{a+bz}{c+dz} = \frac{a}{c} + \frac{ad-bc}{dc} \sum_{n=1}^{\infty} \left(\frac{d}{c}\right)^n z^n$$

$$\phi(z,t) = \frac{(b-s)(e^{st}-1)}{be^{st} - (b+s)} + \frac{se^{st}}{(be^{st} - b+s)b(e^{st}-1)} \sum_{n=1}^{\infty} \left(\frac{b(e^{st}-1)}{be^{st} - b+s}\right)^n z^n$$

Probability of loss

$$P(0,t) = \phi(0,t) = 1 - \frac{e^{st}}{1 + \frac{b}{s}(e^{st}-1)} \rightarrow \begin{cases} 1 - \frac{s}{b} & \text{if } s > 0 \\ 1 - \frac{1}{e} & \text{if } s = 0 \\ 1 - \frac{e^{st}}{1 - \frac{b}{s}} & \text{if } s < 0 \end{cases}$$