# Linked selection from "classic selective sweeps"
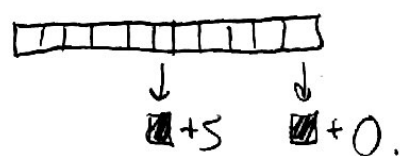
Last time, we saw that the appealing & commonly used QLE/ "independent sites" approx breaks down when $s \gtrsim r\Delta \ell$,

which is not so uncommon for strongly selected mutations $(s \gtrsim 10^{-5})$ & realistic recombination rates $(r \sim \mu \sim 10^{-8} - 10^{-10})$.
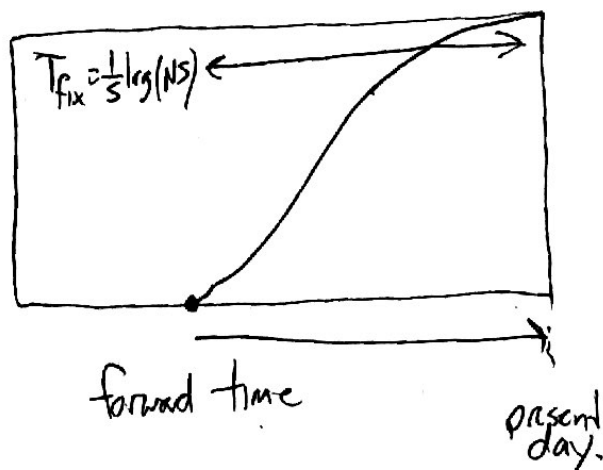
Today we'll discuss some methods for treating cases where $s \gg r\Delta\ell$.

$\Rightarrow$ as before, the simplest scenario to start with is where we have a single strongly beneficial mutation $(s \gg 1/N,)$ on a genome w/ other neutral sites:
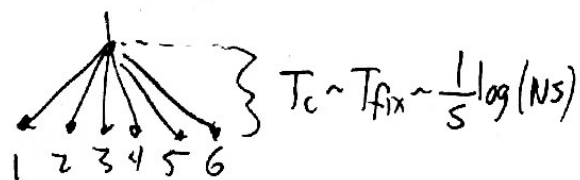


$1+s \qquad 1+0$.

$\Rightarrow$ in this case, we know exactly how the beneficial mutation behaves: it will establish & grow as



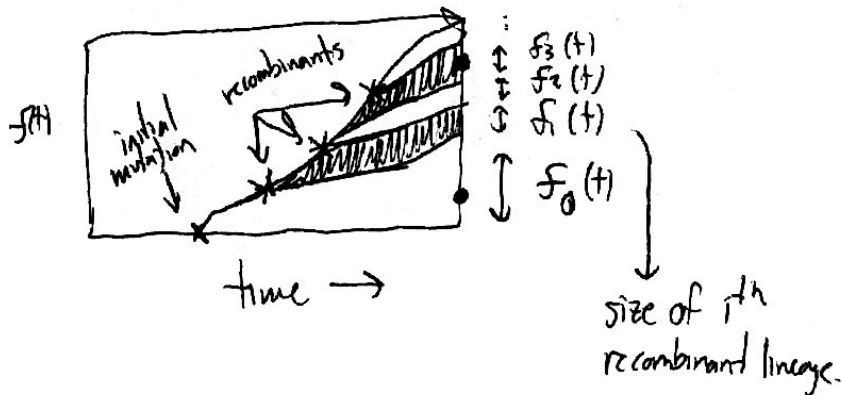$$f(t) = \frac{\frac{1}{Ns} e^{st}}{1 + \frac{1}{Ns} e^{st}}$$

$T_{fix} = \frac{1}{s} \log(Ns)$

forward time

present day.

we can then think about the genealogy of a sample
of 2 (or more individuals) sampled at the present day,
conditioned on this trajectory.

$\Longrightarrow$ in Lecture 15 $(p.6-8)$, we worked out genealogy
at selected locus: everyone descended from original mutation
event w/ star-like tree:

$$T_c \sim T_{fix} \sim \frac{1}{s} \log(Ns)$$

1 2 3 4 5 6

$\Longrightarrow$ for other ~~loci~~ loci $(r>0)$, we have to account for the
fact that some lineages may not have descended from the original
mutation event, but a <u>recombinant lineage</u> produced by recombination
between the selected genetic background + wildtype population.

i.e., <u>lineage structure</u> of
population might look like:



$f_3(t)$
$f_2(t)$
$f_1(t)$
$f_0(t)$

recombinants

initial
mutation

time $\rightarrow$

size of $i^{th}$
recombinant lineage.

$\Longrightarrow$ 2 lineages coalesce during sweep if drawn from
same lineage $\left( P_c = \sum_{i=0}^{n} (f_i)^2 \right)$, $T_c \leq \frac{1}{s} \log(Ns)$

if the individuals aren't from the same sweep lineage, then their ancestors will remain as uncoalesced lineages before origin of sweep. Since no more selection, will coalesce neutrally from here on back $(T_c \sim N \gg T_{fix})$. thus, the

coalescence process looks like:

① if same $\overset{\text{sweep}}{\uparrow}$ lineage, coalesce "immediately" $(T_c \ll N)$

② otherwise, coalesce neutrally $(T_c \sim N)$

$\Rightarrow$ crucial step is to determine the sweep lineage sizes, $f_i(t)$
$\quad\Rightarrow$ we will do this now using a heuristic analysis.

① at short times $(t \ll T_{fix} \sim \frac{1}{S} \log(Ns))$, the selected mutation will still be at low frequency $(f_S(t) \ll 1)$

$\quad\Rightarrow$ most recombinations will occur between a selected haplotype and the wildtype ~~recombinant~~ haplotype, at rate $e f_{sel}(t)$, producing new recombinant lineages.

$\quad\Rightarrow$ each of these recombinant lineages will satisfy the SDE $(t \ll T_{fix})$

$$\frac{df_i}{dt} = s f_i(t) - r f_i(t) + \sqrt{\frac{f_i}{N}} \eta_i$$

we know exactly how these recombinant lineages behave:

$\Rightarrow$ w/ probability $P_{est} \sim (s-r) \sim s$, $\qquad \swarrow$ since $r \ll s$

the lineage will establish and grow · as $f_k(t) \sim \frac{1}{Ns} e^{(s-\rho)(t-\tau_k)}$

where $\tau_k$ is the establishment time of $k^{th}$ recombinant

$\Rightarrow$ since all recombinants grow at the same rate, their relative sizes, $f_k(t)/f_{k-1}(t) = e^{(s-\rho)(\tau_{k-1}-\tau_k)}$,

will only depend on the difference in establishment times and will be "frozen in" for remainder of the sweep, even when $f_{sel}(t)$ is no longer rare.

$\Rightarrow$ thus, understanding what happens during the early phase (when $f_{sel}(t)$ is rare) will be useful for understanding the lineage sizes @ end of the sweep: $f_k(\infty) = \dfrac{f_k(t)}{\sum\limits_{k'} f_{k'}(t)} = \dfrac{e^{-s\tau_k}}{\sum\limits_{k'} e^{-s\tau_{k'}}}$

$\Rightarrow$ crucial step is to calculate the establishment times, $\tau_k$, of the $k^{th}$ recombinant lineage. By convention, will take $\tau_0 = 0$ (establishment time of original mutant lineage)

from our discussion above, <u>successful</u> recombinants

are produced at rate $\quad N \rho f_i(t) \cdot s = N\rho \cdot \frac{1}{Ns} e^{st} \cdot s = \rho e^{st}$

$\underbrace{N \rho f_i(t)}$ → total # recombinants produced @ gen t.

$\cdot s$ → prob that recombinant survives genetic drift.

$\rho e^{st}$ → very similar to Luria delbrück problem from PS1

$\Rightarrow \quad$ avg # successful recombinants produced by time t $\quad = \int_0^t dt' \, \rho e^{st'}$

$\Rightarrow \quad$ heuristically, time to <u>first</u> successful recombinant $(\tau_1)$ occurs when $\quad \int_0^{\tau_*} \rho e^{st} \, dt \sim 1$

$\Rightarrow \quad 1 \sim \frac{\rho}{s}(e^{s\tau_1} - 1) \quad \Rightarrow \quad \tau_1 \approx \frac{1}{s} \log\left(\frac{s}{\rho}\right) \quad (s \gg \rho)$

$\gg \frac{1}{s} \quad \left(\begin{array}{l}\text{i.e., must } \cancel{\text{wait}} \text{ wait} \\ \text{quite a while for first} \\ \text{recombinant to occur}\end{array}\right)$

$\Rightarrow \quad$ similarly, time to $k^{th}$ successful recombinant $(\tau_k)$ occurs when

$\int_0^{\tau_k} \rho e^{st} \, dt \sim k \quad \longrightarrow \quad \tau_k = \frac{1}{s} \log\left(\frac{sk}{\rho}\right) \approx \tau_1 + \frac{1}{s} \log(k)$

Since recombinant lineages grow as $f_i(t) = \frac{1}{Ns} e^{(s-\rho)(t-\tau_k)}$

we have: $f_0(t) = \frac{1}{Ns} e^{(s-\rho)t}$, $f_k(t) = \frac{1}{Ns} e^{(s-r)t} \left(\frac{\rho}{sk}\right)^{1-\frac{\rho}{s}} = f_0(t) \left(\frac{\rho}{sk}\right)^{1-\frac{\rho}{s}}$

extra tiny bit will be important in a bit.

$\Rightarrow$ recombinant lineages are all smaller than
  $f_0$ by a factor of $\rho/s \ll 1$.

(hint that linkage to $f_0(t)$ is important for dynamics)

How many recombinant lineages do we have to keep track of?
(since $f_k \propto 1/k$, $\sum_{k=0}^{\infty} f_k$ diverges!)

$\Rightarrow$ once selected mutations become common, many recombination events will occur between 2 selected haplotypes. $\Rightarrow$ doesn't produce recombinant lineage (from perspective of coalescence during sweep)

$\Rightarrow$ rate of recombinants becomes $\rho f_{sel}(t)[1-f_{sel}(t)]$

$\Rightarrow$ since $f_{sel}(t)$ grows logistically, know that this rate rapidly drops to zero w/in $O(\frac{1}{s})$ generations of $f_{sel}(t) \approx \frac{1}{2}$

$\Rightarrow T_{max} \approx \frac{1}{s} \log(Ns) + O(\frac{1}{s}) \approx \frac{1}{s} \log(Ns)$

$\Rightarrow$ last recombinants occur when $\tau_{K_{max}} \approx \frac{1}{s}\log\left(\frac{sK_{max}}{\rho}\right) \approx \frac{1}{s}\log(Ns)$

$\Rightarrow K_{max} \approx N\rho$

Can now see that there are 2 regimes depending on value of $N\rho$:

① If $N\rho \ll 1 \Rightarrow K_{max} \ll 1 \Rightarrow$ typically no recombinants before mutation sweeps

$\Rightarrow$ "effectively asexual"

② If $N\rho \gg 1 \Rightarrow K_{max} \gg 1 \Rightarrow$ many recombinant lineages contribute @ end of sweep.

$\Rightarrow$ total size is $\displaystyle\sum_{k=0}^{K_{max}} f_k(t) = f_0(t)\left[1 + \sum_{k=1}^{K_{max}}\left(\frac{\rho}{sk}\right)^{1-\frac{\rho}{s}}\right]$

$\approx f_0(t)\left[1 + \int_1^{K_{max}}\left(\frac{\rho k}{s}\right)^{-1+\frac{\rho}{s}} dk\right]$

$\approx f_0(t) \; \exp\left[+\frac{\rho}{s}\log(Ns)\right]$

$\Rightarrow f_0(\infty) = e^{-\frac{\rho}{s}\log(Ns)}, \quad f_k(\infty) = \left(\frac{\rho}{sk}\right)^{1-\frac{\rho}{s}} e^{-\frac{\rho}{s}\log(Ns)}$

Probability that 2 individuals share the same sweep lineage: ⑧

$$P_c = \sum_{k=0}^{K_{max}} f_k(\infty)^2 = e^{-\frac{2\rho}{s}\log(Ns)}\left[1 + \sum_{k=1}^{K_{max}}\left(\frac{\rho}{sk}\right)^{2(1-\frac{\rho}{s})}\right] \simeq e^{-\frac{2\rho}{s}\log(Ns)}$$

$\longrightarrow O\left(\frac{\rho}{s}\right) \ll 1$

$\longrightarrow$ dominated by probability of descending from initial mutant lineage.

Similarly, probability that n individuals share the same lineage:

$$P_c(n) = \sum_{k=0}^{K_{max}} f_k(\infty)^n = e^{-\frac{n\rho}{s}\log(Ns)}\left[1 + \sum_{k=1}^{K_{max}}\left(\frac{\rho}{sk}\right)^{n(1-\frac{\rho}{s})}\right] \simeq e^{-\frac{n\rho}{s}\log(Ns)}$$

$\Longrightarrow$ can use this to calculate pairwise coalescence time :
    ( conditioned on sweep just fixing now )

$$\langle T_2 \rangle = \underbrace{\frac{2}{s}\log(Ns)\, P_c(2)}_{\substack{\text{coalesce during}\\\text{sweep}}} + \underbrace{2N(1-P_c(2))}_{\text{coalesce afterwards}} \approx \underbrace{2N\left(1-e^{-\frac{2\rho}{s}\log(Ns)}\right)}_{\substack{\text{since } \frac{2}{s}\log(Ns) \ll N,\, Ne \gg 1}}$$

if sweep hasn't just fixed, also need
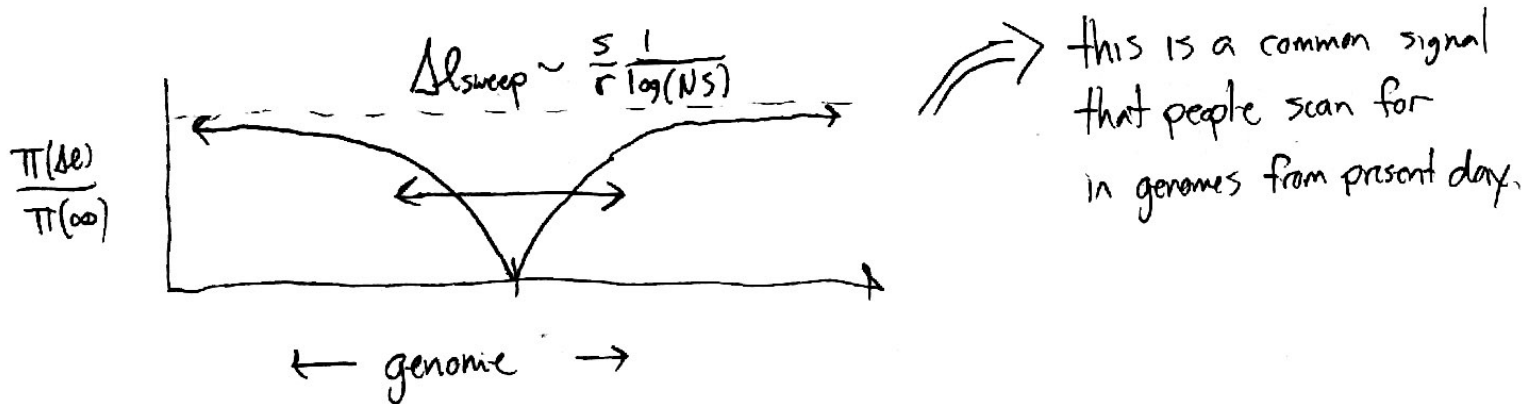to integrate over establishment time $T_{est} \sim Exp\left(\frac{1}{2Ns}\right)$

2 regimes as usual:

① if $T_{est} \gg N$ ($N \nu N s \ll 1$)

$\Rightarrow$ neutral coalescence before sweep, $T_2 \sim N$

② if $T_{est} \ll N$ ($N \nu N s \gg 1$)

$\Rightarrow$ little chance of coalescence until sweep happens

$\Rightarrow \langle T_2 \rangle$ same as before.

$\Rightarrow$ leads to reduction in pairwise heterozygosity ($\pi$) near sweep.

e.g for neutral site distance $\Delta \ell$ from sweep, $\rho = r \Delta \ell$



$$\frac{\pi(\Delta \ell)}{\pi(\infty)} = \left( 1 - e^{-\frac{2r\Delta \ell}{s} \log(Ns)} \right)$$

$\underbrace{\pi(\infty)}_{neutral}$

$\Delta \ell_{sweep} \sim \frac{s}{r} \frac{1}{\log(Ns)}$

$\frac{\pi(\Delta \ell)}{\pi(\infty)}$

$\leftarrow$ genomic $\rightarrow$

$\rightarrow$ this is a common signal that people scan for in genomes from present day.
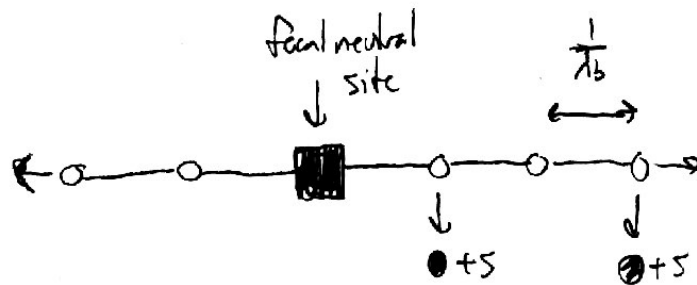
"Recurrent sweeps": so far we have focused on single selected site. can extend this picture to lots of selected sites, ~~provided~~ provided that they act like SSWM sweeps like this one.

$\Rightarrow$ again, can analyze w/ self consistency argument:

Focus on neutral site surrounded by selected sites @ density $\lambda_b$:

$\quad\quad\quad$ ↳ fraction of sites that are strongly beneficial.

$\Rightarrow$ per generation probability of generating a sweep that leads to pairwise coalescence:

$$P_c = \int_0^\infty e^{-\frac{2r\Delta\ell}{s}\log(Ns)} \underbrace{2N\mu\lambda_b \cdot s \, d\Delta\ell} \approx \frac{N\mu\lambda_b s^2}{r\log(Ns)}$$

$\quad\quad\quad\quad$ ↳ dominated by probability of having a really close sweep w/ $P_c(z) \approx \mathcal{O}(1)$.

$$\Delta\ell \lesssim \ell^* = \frac{s}{r}\frac{1}{\log(Ns)}$$

$\Rightarrow$ again, if $P_c \gg \frac{1}{N}$ $\Rightarrow$ $\langle T_2 \rangle \approx \frac{r\log(Ns)}{N\mu\lambda_b s^2} \approx \frac{1}{N U_{eff} s}$

$\quad\quad$ "linkage block, $\ell^* = \frac{s}{r}\frac{1}{\log(Ns)}$" $\quad\longleftarrow\quad$ like asexual case w/ $U_{b,eff} = \mu\lambda_b \ell^*$

Note, however, that coalescent process looks very different from neutral coalescence w/ $N_e \approx \frac{1}{N U_{eff} s}$ (or any $N_e(t)$)
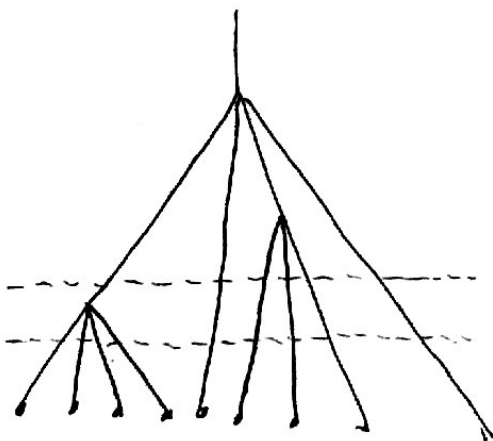
$\hookrightarrow$ common misconception in pop-gen. world @ moment

$\Rightarrow$ per generation probability of generating a sweep that leads to coalescence of all ~~lineages~~ $n$ lineages:

$$P_c(n) = \int_0^\infty e^{-\frac{nr\Delta\ell}{s}\log(Ns)} \, 2N\mu\lambda_b \, s \, d\Delta\ell \approx \frac{P_c(2)}{n} \rightarrow \text{decays very slowly in } n.$$

$$\Big[ \text{compare to } P_c(n) = P_c(2)^{\wedge} \text{ for neutral coalescent} \Big]$$

$\Rightarrow$ this means that conditioned on 2 lineages coalescing, very likely to have several coalescing @ once:



$N_e(t)$  $\rightarrow t$

$\updownarrow$

$\Big\{$ even for transient bottleneck, rare for all $k$ lineages to coalesce into 1, rather than $\binom{k}{2}$ disparate pairs.

Finally, when is this picture of recurrent non-Interfering sweeps a good approximation? $\Rightarrow$ need to check self-consistency.

① 2 sweeps cannot occur w/in $\ell^*$ of each other w/in a single fixation time, $\frac{1}{s}\log(Ns)$: (i.e, SSWM w/in $\ell^*$)

$$\Rightarrow N\nu\lambda_b\ell^* \cdot s \cdot \frac{1}{s}\log(Ns) \ll 1 \quad\Rightarrow\quad ❶ \frac{\nu}{r}\cdot\lambda_b\cdot Ns \ll 1$$

② when 2 sweeps do occur in same sweep time, should have $s \ll r\hat{\ell}$. (unliked even before establishment time.)

$$\Rightarrow N\nu\lambda_b\hat{\ell}\cdot s \cdot \frac{1}{s}\log(Ns) \sim 1 \quad\Rightarrow\quad \hat{\ell} = \frac{1}{N\nu\lambda_b\log(Ns)}$$

$$\Rightarrow s \ll \frac{r}{N\nu\lambda_b\log(Ns)} \quad\Rightarrow\quad \boxed{\frac{\nu}{r}\cdot\lambda_b\cdot Ns\log(Ns) \ll 1}$$

$\hookrightarrow$ slightly more stringent.

$\Rightarrow$ when this condition is not met, multiple beneficial mutations will interfere w/ each other. $\Rightarrow$ generic for large enough $N$!