

Announcements : ① PSET 3 DUE 3/9/21 (only 3 problems this time,
 ⇒ but #1 is long!)

② weekly comments

Last time : Multi-locus models of evolution ($L \gg 1$)

$$\frac{df(\vec{g})}{dt} = [X(\vec{g}) - \bar{X}(\vec{f})] f(\vec{g}) + \sum_{\vec{g}'} M(\vec{g} \rightarrow \vec{g}') f(\vec{g}') - M(\vec{g}' \rightarrow \vec{g}) f(\vec{g}) + r \sum_{\vec{g}_1, \vec{g}_2} T(\vec{g}_1, \vec{g}_2 \rightarrow \vec{g}) f(\vec{g}_1) f(\vec{g}_2) - r f(\vec{g}) + \sqrt{\frac{f(\vec{g})}{N}} \eta(\vec{g}) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N}} \eta(\vec{g}')$$

No closed-form sol'n in the general case...

⇒ Need asymptotic approximations (i.e., limits of N, L, r, M , etc...)

① Large N / "mean field" (genetic drift is negligible)

↳ ~~X~~ inconsistent when $L \gg 1$! (e.g. $\frac{\# \text{individuals}}{\# \text{genotypes}} \approx 1$)

Today :

② Successive mutations regime
 (i.e., mutation is small correction)

$$\frac{df(\vec{g})}{dt} = \sim (x - \bar{x}) + \cancel{\sim L \mu} + \sim r + \sim \frac{\sigma^2}{JN}$$

\Rightarrow i.e. new mutations fix or go extinct before next one occurs...

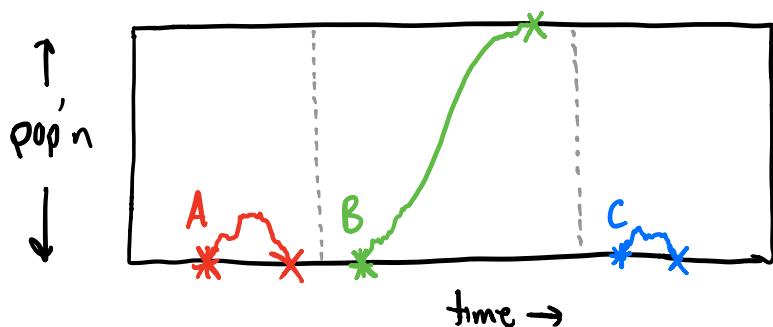


\Rightarrow @ any given time, only 2 genotypes present:

$$\begin{array}{ll} \text{"current wildtype"} & \vec{g}_0 = (1, 0, 1, 1, 0, 0, 0) \\ & \downarrow \text{mutation @ site } l \\ \text{"single mutant"} & \vec{g}_m = (1, 0, 1, 1, 0, 1, 0) \end{array}$$

\Rightarrow what can **recombination** do? Nothing! (on average)

\Rightarrow then each mutation looks like $L=1$ model w/



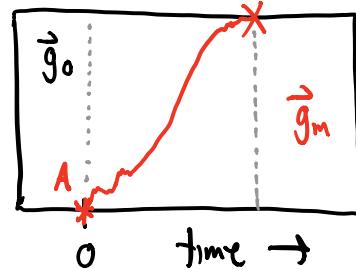
$$S_{\text{eff}} \equiv X(\vec{g}_m) - X(\vec{g}_0)$$

\Rightarrow in this case, know exactly what happens:

$$(i) \text{ w/ prob } P_{\text{fix}} = \frac{2s}{1-e^{-2Ns}},$$

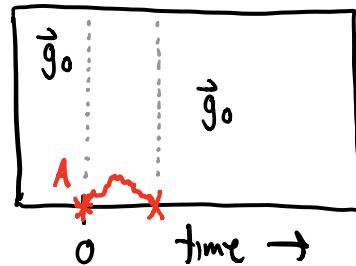
\Rightarrow mutation fixes ("sweeps")

$\Rightarrow \vec{g}_0 \rightarrow \vec{g}_m$; repeat!

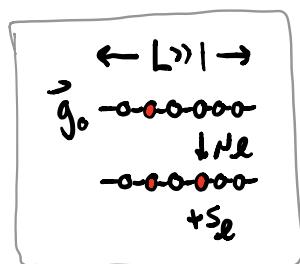


(ii) otherwise, mutation goes extinct

$\Rightarrow \vec{g}_0$ stays put.

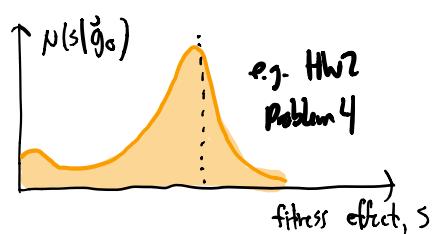


\Rightarrow New: When $L > 1$, multiple different mutations are possible...



\Rightarrow useful to define fitness effects

$$s_e \equiv X(\vec{g}_0 + \underset{\text{mut'n}}{\text{at site } e}) - X(\vec{g}_0)$$



along w/ distribution of fitness effects ("DFE")

$$N(s|\vec{g}_0) \equiv \sum_{e=1}^L N_e \delta(s - s_e)$$

I
 technically
 depends on \vec{g}_0 ↓
 prob. of drawing
 a mutation w/ effect s .

\Rightarrow mutations w/ fitness effect s are produced as

Poisson process w/ rate $N_N(s)$.

\Rightarrow if each successful w/ prob $P_{\text{fix}}(s)$,

then successful mutations also Poisson Process

$$\text{w/ total rate } R \equiv \int_0^{\infty} N_N(s) \cdot P_{\text{fix}}(s) ds = \sum_e N_{N_e} P_{\text{fix}}(s_e)$$

\Rightarrow ① time until next successful mutation is born is:

$$T_{\text{est}} \sim \text{Exponential}\left(\frac{1}{R}\right) \quad \left[\begin{array}{l} \text{similar to } N_N \rightarrow 0 \\ \text{case in lecture 10} \end{array} \right]$$

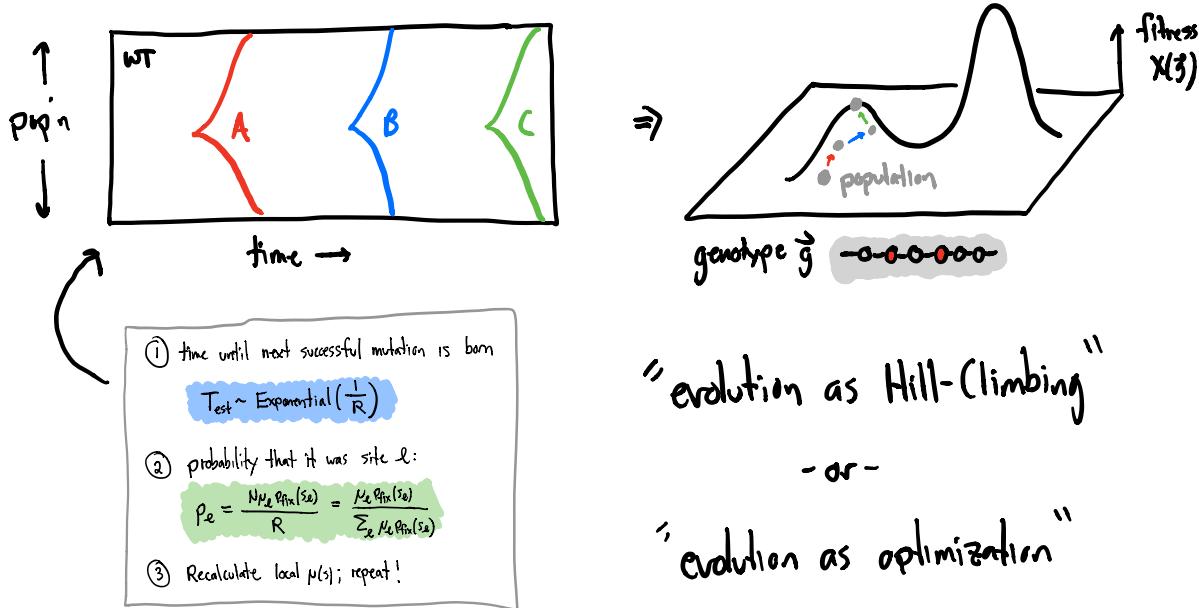
\Rightarrow ② probability that it was site e :

$$P_e = \frac{N_{N_e} P_{\text{fix}}(s_e)}{R} = \frac{N_e P_{\text{fix}}(s_e)}{\sum_e N_e P_{\text{fix}}(s_e)}$$

$\Rightarrow ③ \vec{g}_0 \rightarrow \vec{g}_m \Rightarrow$ recalculate $\mu(s|\vec{g}_0) \Rightarrow$ repeat!

When approx is valid: will check carefully below...

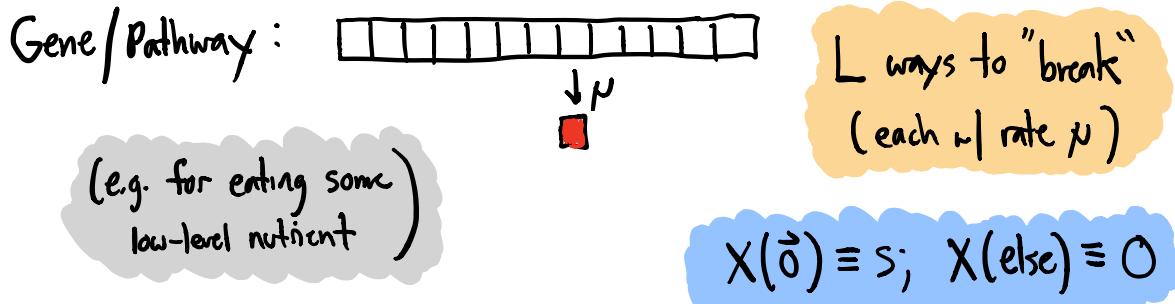
\Rightarrow simple algorithm for modeling evolution (not just pop gen)



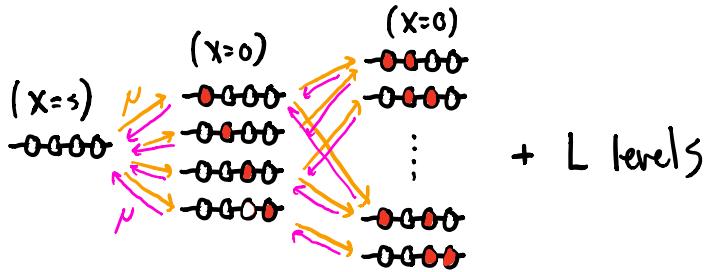
Note: even in these simplified settings,

\Rightarrow fundamental limits to optimization picture...

Example: maintaining a useful function



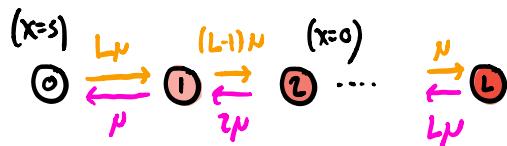
Genotype network:



Key simplification:

"equivalence class"

$$f_k \equiv \sum_{|\vec{g}|=k} f(\vec{g})$$



↳ can track: $P_k(t) \equiv \Pr [S_k(t)=1] = \langle f_k(t) \rangle$

How can $P_k(t)$ change?

$$\underline{k=0}: \quad d_t P_0 = \underbrace{N \mu p_{fix}(s) P_1}_{\text{incoming}} - \underbrace{N L N \cdot p_{fix}(-s) \cdot P_0}_{\text{outgoing}} \xrightarrow{t \rightarrow 0} 0$$

$$\Rightarrow \frac{P_0}{P_1} = \frac{N \mu p_{fix}(s)}{N L N \cdot p_{fix}(-s)} = \frac{1}{L} e^{2Ns}$$

$$\underline{k=1}: \quad d_t P_1 = N L N p_{fix}(-s) P_0 - N \mu p_{fix} P_1 \quad (\text{from 0 class}) \\ + N \cdot 2 \mu \cdot \left(\frac{1}{N}\right) P_2 - N(L-1) N \left(\frac{1}{N}\right) P_1 \xrightarrow{t \rightarrow \infty} 0$$

$$\Rightarrow P_2 = \frac{(L-1)}{2} P_1$$

$$\underline{k=2}: \quad \Rightarrow P_3 = \frac{L-2}{3} \cdot P_2 = \frac{(L-1)(L-2)}{3 \cdot 2} \cdot P_1$$

$$\Rightarrow P_K = \frac{1}{L} \frac{L!}{k!(L-k)!} P_1$$

$$\Rightarrow 1 - P_0 = \sum_{k=1}^L P_k = \frac{1}{L} (2^L - 1) P_1$$

$$\Rightarrow f_1 = (1-f_0) \frac{L}{(2^L - 1)}$$

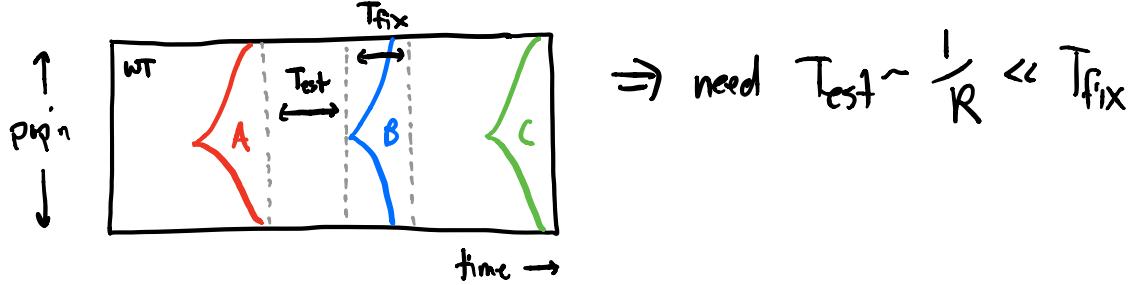
$$\Rightarrow \frac{P_0}{1-P_0} = \exp \left[\underbrace{2NS}_{\text{"drift barrier"}} - \underbrace{\log(2^L - 1)}_{\text{"entropy of generated space"}} \right] = \frac{\text{prob having function}}{\text{prob broken}}$$

$$\geq \frac{0.9}{0.1} = \exp(2) \leq \exp \left(2NS - \overbrace{\log(2^L - 1)}^{\alpha L} \right)$$

$\hookrightarrow S \geq \frac{1}{N}$ "drift barrier"

 e.g. compare to deterministic case: $\bar{f}_0 = 1 - \frac{L_N}{S}$

When is successive mutations regime a good approx?



$$\Rightarrow \text{need } T_{\text{fix}} - \frac{1}{R} \ll T_{\text{fix}}$$

E.g. Neutral mutations ($\mu(s) = U_n \delta(s)$)

$$P_{\text{fix}}(0) = \frac{1}{N} \Rightarrow R = NU_n \left(\frac{1}{N}\right) = U_n ; \quad T_{\text{fix}} \sim O(N)$$

$$\Rightarrow \text{need } \frac{1}{U_n} \gg N \Rightarrow \boxed{NU_n \ll 1} \quad \begin{matrix} \text{"weak mutation"} \\ \text{"weak selection"} \end{matrix}$$

E.g. Strongly beneficial mutations ($\mu(s) = U_b \delta(s-s_b) ; \quad Ns_b \gg 1$)

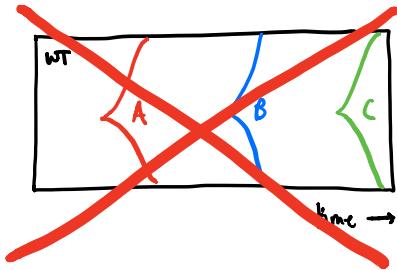
$$\Rightarrow P_{\text{fix}}(s) \approx 2s \Rightarrow R = 2NU_b s ; \quad T_{\text{fix}} = \frac{2}{3} \log(Ns)$$

$$\Rightarrow \text{need } \frac{1}{2NU_b s} \gg \frac{2}{3} \log(Ns) \Rightarrow \boxed{NU_b \ll \frac{1}{\log(Ns_b)}}^{< 1}$$

\Rightarrow What does this look like \Rightarrow e.g. HWZ problem 4
for some "real" parameter values?
 $U_b = 5 \times 10^{-6}, \quad s_b = 0.02$
just for L.O.F. muts.

$$\Rightarrow \text{e.g. if } N \sim 10^5 \Rightarrow 2Nu_b \lg(Ns) = 3$$

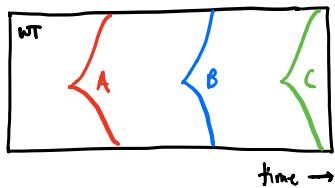
\Rightarrow successive mutations picture
does not apply!



\Rightarrow what do things look like instead?

Neutral theory & the Coalescent

Successive mutations:



$$\frac{dS(\vec{g})}{dt} = \sim(x-\bar{x}) + \cancel{\sim L\cdot\mu} \rightarrow \epsilon \\ + \sim\epsilon + \sim\frac{\pi}{JN}$$

⇒ ~1 variant present @ high freqs ⇒ solved by reducing to $L=1$ model

* But genomes in data separated by multiple mut'n's
(e.g. humans, 2 individuals differ by ~1 mut / 1000 bp)

⇒ need to understand what's going on in these cases...

$$\frac{dS(\vec{g})}{dt} = \cancel{\sim(x-\bar{x})} + \sim L\cdot\mu \rightarrow 0 \\ + \sim\epsilon + \sim\frac{\pi}{JN} \rightarrow 0$$

⇒ one other limit that's well understood:
neutral evolution in nonrecombining genome

when $X(\vec{g})=0 \wedge \rho=0$, left with: $(N_e = \nu_e)$

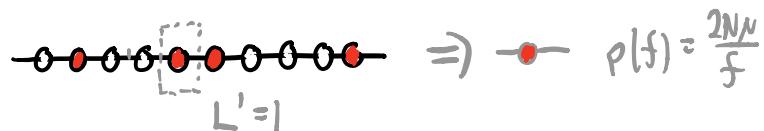
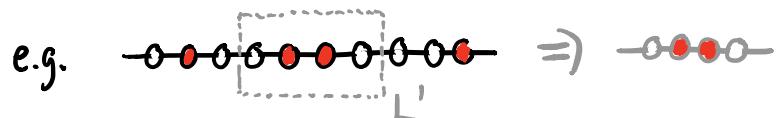
$$\frac{\partial f(\vec{g})}{\partial t} = \sum_e \sum_{|\vec{g}' - \vec{g}|=1} \mu_e f(\vec{g}') \left[g_e(1-g_e) + (1-g_e)g_e' \right] - \sum_e N_e f(\vec{g})$$

incoming mutations outgoing mutations

$$+ \sqrt{\frac{f(\vec{g})}{N}} \eta(\vec{g}) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N}} \eta(\vec{g}')$$

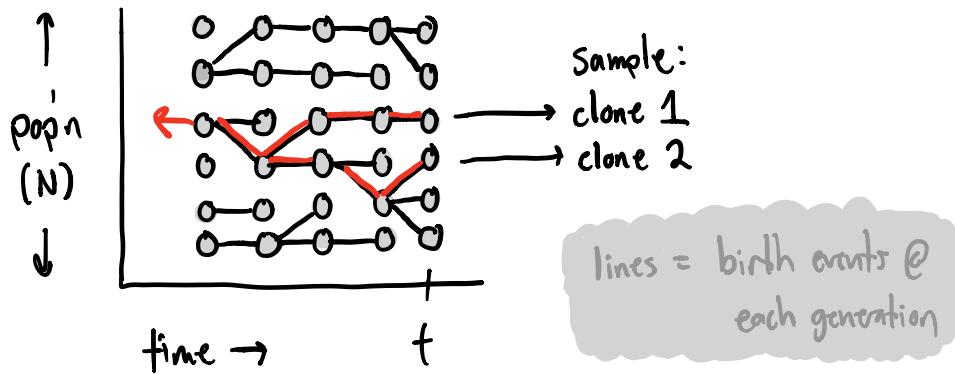
genetic drift

Key insight: sites don't actually influence each other (because neutral)

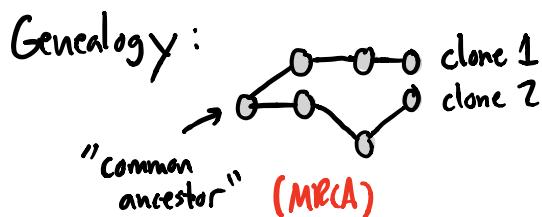


\Rightarrow 2nd key insight: can take $L'=0$ —

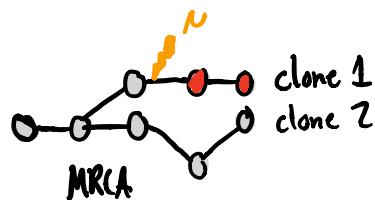
E.g. simulation of neutral pop'n in Wright-Fisher model:



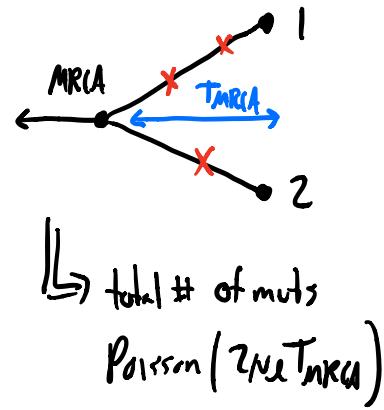
⇒ key insight: lines also = **genealogical relationships**
backward in time!



↓ differences between sampled individuals
 = mutations on genealogy



\Rightarrow Mutations occur as Poisson Process
w/ rate μ_e on each branch

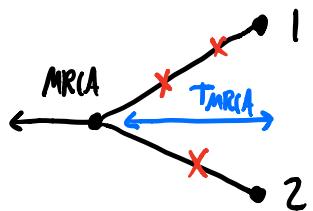


\Rightarrow 2 extreme limits:

(1) $\mu_e T_{\text{MRCAs}} \ll 1 \Rightarrow$ 0 or 1 mutations on whole tree

$$\Rightarrow \Pr[\text{genetic diff} @ \text{site } e] =$$

(2) $\mu_e T_{\text{MRCAs}} \gg 1 \Rightarrow$ lots of forward & backward mutations along each branch.

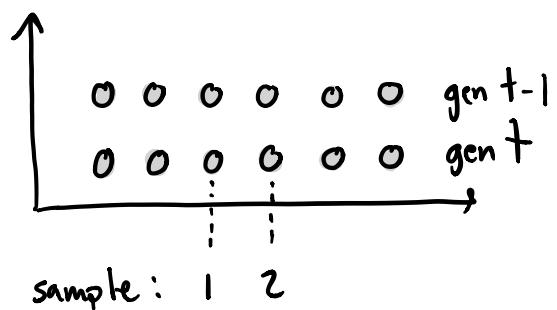


$$\Rightarrow \Pr[\text{genetic diff} @ \text{site } L] =$$

\Rightarrow Key Question: what sets T_{MRCAs} ?

Note: T_{MRCA} is random quantity (genealogy will differ from simulation to simulation)

\Rightarrow key insight: start from present & work backward in time:



Process repeats itself w/ next gen:

