

Multilocus Evolution II : the successive mutations regime

(1)

Last time, we showed that ~~the~~ generalizing our serial dilution model to a genome w/ $L > 1$ sites leads to a system of coupled stochastic differential equations for each genotype frequency, $f(\vec{g})$, which has the general form:

$$\frac{df(\vec{g}, t)}{dt} = \underbrace{\left[X(\vec{g}) - \bar{X}(t) \right] f(\vec{g})}_{\text{natural selection}} + \underbrace{\sum_{\vec{g}'} \left[M_{\vec{g}' \rightarrow \vec{g}} f(\vec{g}') - M_{\vec{g} \rightarrow \vec{g}'} f(\vec{g}) \right]}_{\text{mutation (matrix of rates, } M)}$$

$$+ \underbrace{\rho \left[-f(\vec{g}) + \sum_{\vec{g}_F, \vec{g}_D} T_{\vec{g}_F, \vec{g}_D \rightarrow \vec{g}} f(\vec{g}_F) f(\vec{g}_D) \right]}_{\text{recombination (tensor of rates, } T_{\vec{g}_F, \vec{g}_D \rightarrow \vec{g}})}$$

$$+ \underbrace{\left[\sqrt{\frac{f(\vec{g})}{N_e}} \eta(\vec{g}, t) - f(\vec{g}) \sum_{\vec{g}'} \sqrt{\frac{f(\vec{g}')}{N_e}} \eta(\vec{g}', t) \right]}_{\text{genetic drift}}$$

2

In contrast to single-locus case ($L=1$), there is no closed form solution of this model — even for simple things like equilibrium distributions or fixation probabilities — for the general multi-locus case, even for L as small as $L=2$!

(similar to multi-particle Schrödinger eq. in Quantum Mechanics)

⇒ instead, the only way we can make progress in understanding the $L>1$ case is by considering behavior in different asymptotic limits. (this is why we spent so much time discussing them earlier in the course... they are our only hope now!)

⇒ the next several lectures will focus on a few different limits where the behavior is reasonably well understood, and when we expect them to apply. (e.g. to data)

⇒ ~~although~~ a lot is already known, but it is far from a complete picture, and a lot of theoretical work is still being done to understand consequences of this basic model.

Q: Given parameters (knobs), $L, N, X(\vec{g}), M, e, T,$
what are some limits where we might be able to understand behavior of these equations?

① Obvious answer: $L=1$. that's cheating! need $L \geq 2$
(but basic idea is a good one... when can we reduce ~~system~~ system to something that looks like $L=1$ → then can use what we already know...)

② In ~~physics~~ physics, already primed to take $N \rightarrow \infty$ limit since at least then we expect noise to go away...

⇒ is this a good approximation here? for $L=2$ maybe... (@ long times)

⇒ but for $L \gg 1 \Rightarrow 2^L$ is eventually $\gg N$, even for enormous N , e.g. $L \sim 1000$ (1 gene) - $2^L = 10^{300}$!

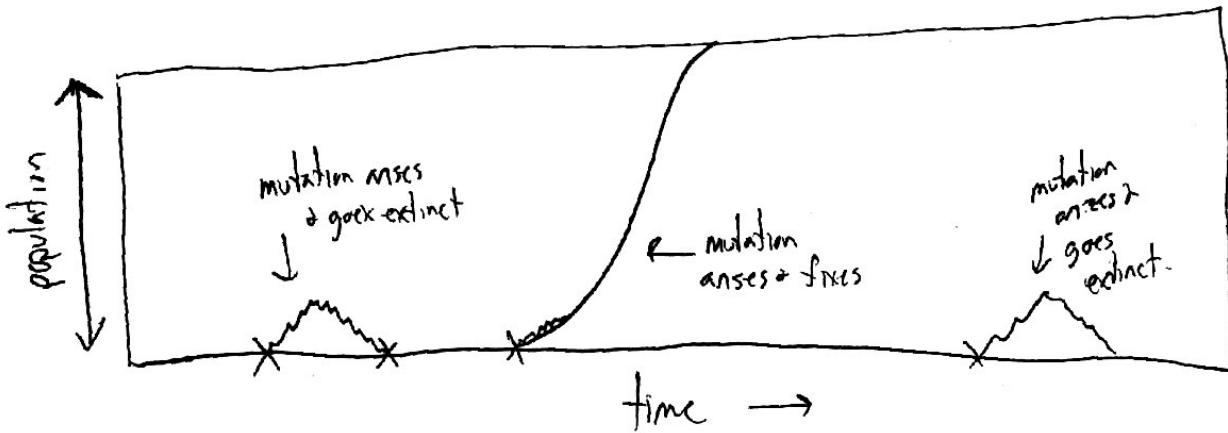
⇒ always will be some genotypes that will be unoccupied, and mutations into these genotypes will be strongly influenced by noise... so large N doesn't help here.

③ Successive mutations regime

④

what if mutation rates are so low that only 0, or 1 mutations are present in population @ any given time.

~~what if~~ \rightarrow i.e., mutation occurs and then either fixes or goes extinct before next mutation occurs.



then @ any given time, only 2 genotypes present in population:

~~current~~ "current wildtype", $\vec{g}_0 = (1, 0, 1, 1, 0, 0, 0)$
 single mutation away from WT eg. @ site l . $\vec{g}_m = (1, 0, 1, 1, 0, \underline{1}, 0)$ \rightarrow site l .

what can recombination do?
 $\vec{g}_0 = (1, 0, 1, 1, 0, 0, 0)$ $\xrightarrow{\text{no change!}}$ creates another \vec{g}_m
 $\vec{g}_m = \boxed{(1, 0, 1)} \mid \boxed{(0, 1, 0)}$ \otimes (or vice versa)

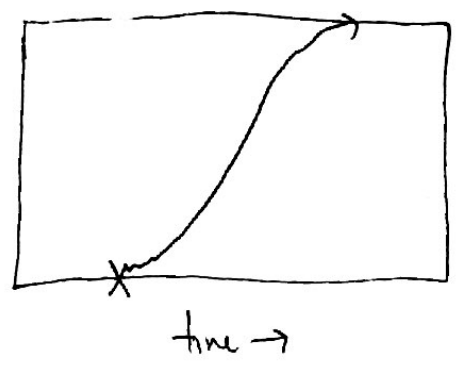
\Rightarrow recombination drops out.

then right after a mutation occurs, we're left w/
something that looks like a single locus model

$$w/ \quad s \equiv X(\vec{g}_m) - X(\vec{g}_0)$$

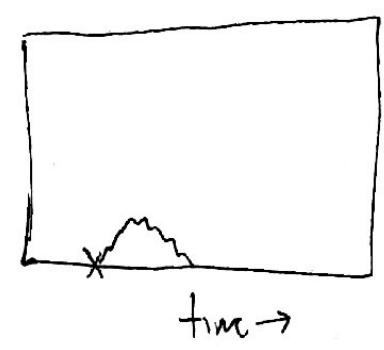
⇒ in this case, we know exactly what happens:

(i) w/ probability $p_{fix}(s) = \frac{2s}{1 - e^{-2Nes}}$, the mutation takes over ("sweeps")



and $\vec{g}_0 \rightarrow \vec{g}_m$; process repeats.

(ii) otherwise, mutation goes extinct
and \vec{g}_0 stays put.



to account for all the ~~ways~~ ways that this
mutation could have occurred, let's assume for simplicity
that $\mu_e = \nu_e$ (forward & back mutations occur @ same rate)

$$\Rightarrow \text{total mutation rate for all genotypes is } U \equiv \sum_{l=1}^L \mu_e$$

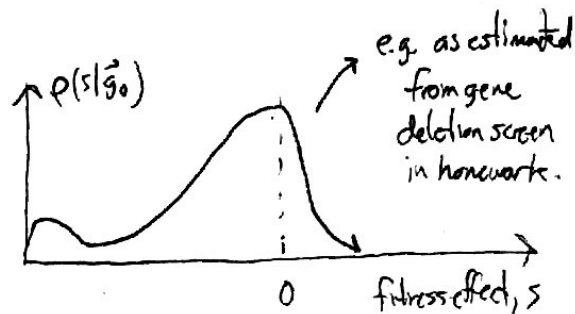
the fitness effect of each mutation is given by

6

$$S_e = X(\vec{g}_0 + \text{mutation @ site } e) - X(\vec{g}_0)$$

and entire collection $\{(\mu_e, S_e)\}_{e=1}^L$ is sometimes summarized as a distribution of fitness effects (DFE)

$$P(s | \vec{g}_0) \equiv \frac{1}{U} \sum_{e=1}^L \mu_e \delta(s - S_e)$$



↑
technically, depends on current w/ type

↳ probability of drawing a mutation w/ effect, s.

⇒ mutations w/ fitness effect s are produced as a Poisson process @ rate $NU P(s | \vec{g})$. If each one is successful w/ probability $P_{\text{fix}}(s)$, then ~~rate~~ occurrence of successful* mutations is Poisson process w/ total rate

$$R \equiv \int NU P(s | \vec{g}_0) P_{\text{fix}}(s) ds.$$

(*technically, "mutations destined to be successful")

the time until the next successful mutation is born is

$$T_{\text{est}} \sim \text{Exponential}\left(\frac{1}{R}\right)$$

⇒ when a successful mutation occurs, the probability it was site l is given by

~~Pe~~ ~~Pe~~
$$P_l = \frac{N \mu_l P_{\text{fix}}(s_l)}{R} = \frac{\mu_l P_{\text{fix}}(s_l)}{\sum_l \mu_l P_{\text{fix}}(s_l)}$$

⇒ alternatively, probability it has fitness effect $s \pm ds$ is

$$P_{\text{fix}}(s) \propto P_{\text{fix}}(s) e(s)$$

⇒ so we randomly select a site to fix, ~~take~~
take $\vec{g}_0 \rightarrow \vec{g}_M$, and repeat!

when is this a good approximation?

⇒ ~~need~~ need $T_{\text{est}} \sim \frac{1}{R} \gg T_{\text{fix}} \sim \begin{cases} N & \text{if } |s| \ll \frac{1}{N} \\ \frac{1}{s} \log(Ns) & \text{if } |s| \gg \frac{1}{N} \end{cases}$

thus, we see that T_{fix} is bounded from above by N

(8)

\Rightarrow if we have $\frac{1}{R} \gg N \geq T_{\text{fix}}$, or $\int NU N p_{\text{fix}}(s) p(s) ds \ll 1$

\Rightarrow then successive mutations picture will apply.

\Rightarrow known as "weak mutation limit" (better "rare mutation limit")

Since we can always find a U low enough that the WM condition applies. How low? Different behavior depending on s 's that dominate $p_{\text{fix}}(s)p(s)$ integral.

~~scribble~~

(i) if $p_{\text{fix}}(s)p(s)$ dominated by neutral mutations ($s \approx 0$)

then $\int NU N p_{\text{fix}}(s)p(s) ds \approx NU_0$ (where U_0 is neutral mutation rate, typically $\theta(U)$.)

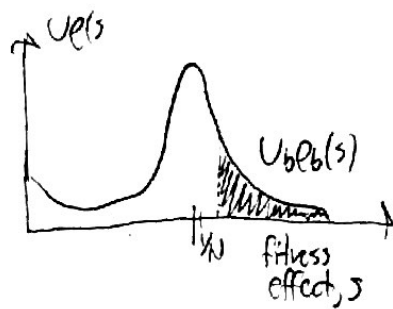
\Rightarrow weak mutation limit when $NU_0 \ll 1$.

(ii) if $p_{\text{fix}}(s)p(s)$ dominated by $s \gg \frac{1}{N}$, then

$\int NU N p_{\text{fix}}(s)p(s) ds = \int 2NU_b N s p_b(s) ds$

where $U_b \rho_b(s)$ is the strongly beneficial portion of the original OPE ($U_b \ll U$)

where $P_{\text{fix}}(s) \approx 2s$.



(9)

\Rightarrow in this case, weak mutation limit applies when

$$NU_b N\bar{s}_b \ll 1 \quad \Rightarrow \quad NU_b \ll \frac{1}{N\bar{s}_b} \quad (\bar{s}_b \equiv \int s \rho_b(s) ds)$$

(very small)

\Rightarrow this condition is maybe too stringent, since $T_{\text{fix}} \sim N$ only for neutral mutations. The mutations that dominate R

have $T_{\text{fix}} \sim \frac{1}{\bar{s}_b} \log(N\bar{s}_b) \ll N$. If we only care about

these, then the successive mutations picture applies when

$$\frac{1}{\bar{s}_b} \log(N\bar{s}_b) \ll \frac{1}{2NU_b\bar{s}_b} \quad \Rightarrow \quad NU_b \ll \frac{1}{\log(N\bar{s}_b)} \lesssim 1$$

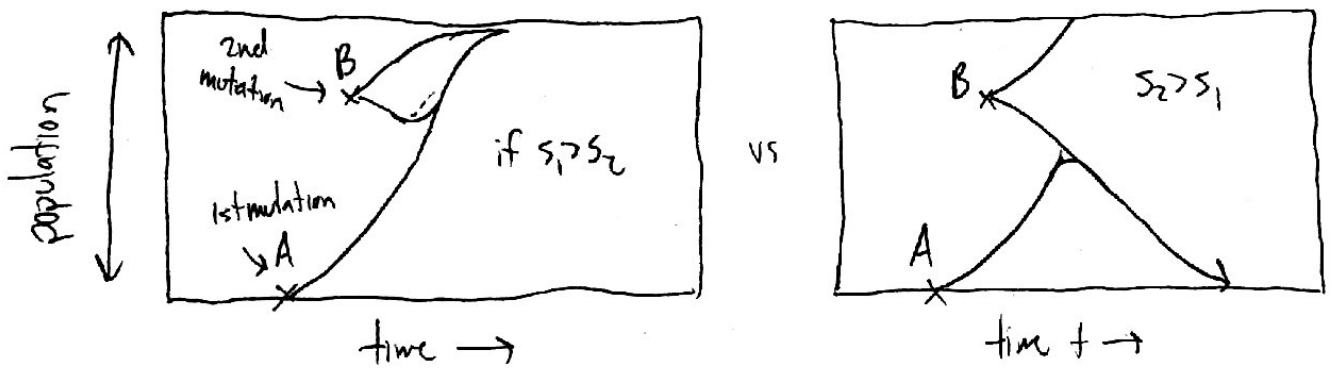
\Rightarrow this is known as the strong-selection weak-mutation (SSWM) regime, because only strongly beneficial mutations have a chance of fixing.

What does this condition look like for some real parameter values? (10)

$\Rightarrow U_b e_b(s)$ is hardest to constrain; from Problem 4 last week, you estimated $U_b \sim 5 \times 10^{-6}$, $s_b \sim 2e^{-02}$ just from L.O.F. mutations (probably a lower bound on U_b ...)

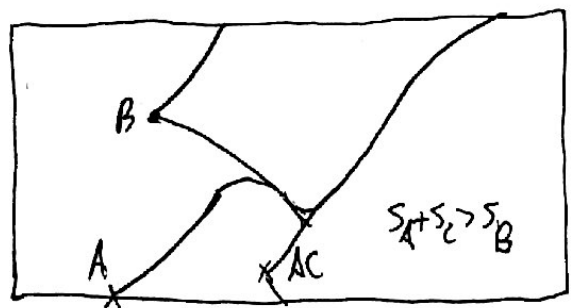
\Rightarrow if $N = 10^5$, then $2NU_b \log(Ns_b) \approx 3$, so already outside region of validity.

\Rightarrow what does behavior look like instead?



\Rightarrow either way, a successful mutation has been wasted. known as "clonal interference".

can get arbitrarily complicated when more mutations involved:

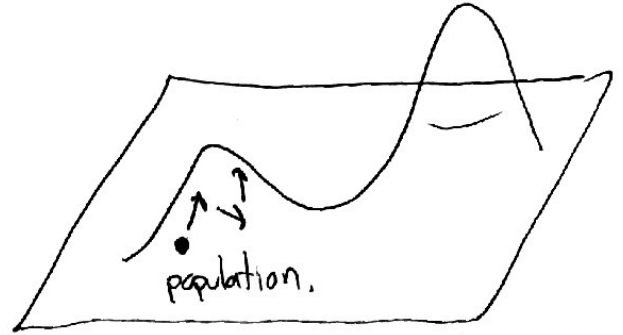


However, for small enough N , the SSWM limit will apply.

(11)

⇒ in this case, we have a well-defined limit for thinking about evolution on a long genome:

⇒ the population is represented by a single point in genotype space, and it undergoes a biased random walk to higher fitness values, as dictated by the DFE

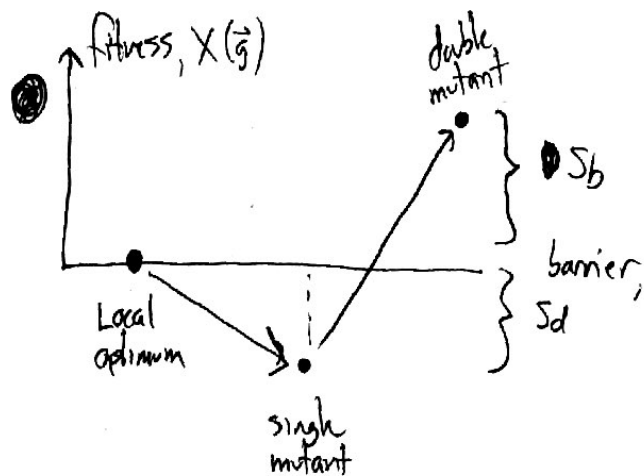


⇒ this is the classic picture of evolution as hill-climbing that many of us have in our heads. Here we see it emerges only in a very simple limit when we ignore the complexities of what's going on w/in a population.

⇒ in this hill-climbing picture, all of the interesting behavior boils down to epistasis: how long will the population wander before it gets trapped in a local fitness peak. how many peaks are there? ("ruggedness")

of course, when the population does reach a local optimum, the SSWM limit will not apply in its original sense, because $U_b e_b(s) = 0!$

To escape, population must cross a "fitness valley":



⇒ How fast does this happen?

⇒ when $Ns_d \gg 1$, very low rate of fixing deleterious intermediate
($N_{fix} \sim Ns_d e^{-Ns_d}$)

⇒ however, each deleterious mutation has some chance of producing successful mutation before it dies off.

⇒ can model this w/ 2-locus branching process
(since all interesting behavior will happen when ~~some~~ $f_1, f_2 \ll 1$)

$$\frac{df_1}{dt} = -s_d f_1 - \mu_{1 \rightarrow 2} f_1 + \sqrt{\frac{f_1}{N}} \eta_1(t)$$

$$f_1(0) = \frac{1}{N}, f_2(0) = 0$$

$$\frac{df_2}{dt} = s_b f_2 + \mu_{1 \rightarrow 2} f_1 + \sqrt{\frac{f_2}{N}} \eta_2(t)$$

↳ what happens @ long times?

this is similar to our single locus B.P. model, except that now f_2 is coupled to random input μf_1 .

\Rightarrow as before, can solve by turning to generating function, this time for 2 variables: $H(z_1, z_2, t) = \langle e^{-z_1 f_1(t) - z_2 f_2(t)} \rangle$

\Rightarrow @ time $t=0$: $H(z_1, z_2, 0) = e^{-z_1 \frac{1}{N}} \approx 1 - \frac{1}{N} z_1$

\Rightarrow @ long times, $H(z_1, z_2, t) \rightarrow e^{-0} \cdot P_{ext} + e^{-\infty} \cdot (1 - P_{ext})$
 \downarrow beneficial mutant takes off.

\Rightarrow using same SDE manipulations as before, can show that generating function satisfies PDE: (exercise for reader)

$$\frac{\partial H}{\partial t} = \left[-s_1 z_1 + N_2 z_2 - \frac{z_1^2}{2N} \right] \frac{\partial H}{\partial z_1} + \left[s_2 z_2 - \frac{z_2^2}{2N} \right] \frac{\partial H}{\partial z_2}$$

\Rightarrow can solve completely w/ method of characteristics (exercise) but can also look for z_1^*, z_2^* s.t. $\frac{\partial H}{\partial t} = 0$.

\Rightarrow from above: $z_1^* = N p_{fix}$ \leftarrow probability of generating successful double mutant before going extinct.

system of quadratic equations:

$$s_b z_2^* - \frac{z_2^{*2}}{2N} = 0$$

$$-s_d z_1^* + s_b z_2^* - \frac{z_1^{*2}}{2N} = 0$$

note:
 \Rightarrow same trick works
 for arbitrary # of types
 coupled together!

\Rightarrow in this case, can easily reduce to 1 quadratic:

$$z_1^{*2} + 2Ns_d z_1^* - 2Ns_b N = 0$$

$$\Rightarrow P_{fix} \approx \begin{cases} \frac{N}{s_d} \cdot 2s_b & \text{if } \frac{s_b N}{s_d^2} \ll 1 \\ \sqrt{2Ns_b} & \text{if } \frac{s_b N}{s_d^2} \gg 1 \end{cases}$$

\Rightarrow how can we understand this result?

\Rightarrow for each double mutant produced, probability of establishment is $\sim s_b$.

~~probability of producing double mutant~~ ~~before going extinct is~~

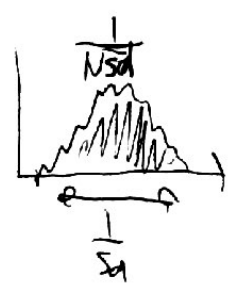
$$P_{fix} \approx \left(1 - e^{-\int_0^\infty N s_1(t) \cdot s_b \cdot dt} \right) \leftarrow \begin{matrix} \text{i.e., depends on} \\ \text{integral of random} \\ \text{function.} \end{matrix}$$

$f_1(t)$

we saw before that $f_i(t)$ can survive for as long as $\sim \frac{1}{s_d}$ generations and grow to size $\sim \frac{1}{N s_d}$ w/ probability s_d .

\Rightarrow when this happens, $\int f_i(t) dt \sim \frac{1}{N s_d} \cdot \frac{1}{s_d}$

\uparrow \uparrow
 height width
 (f_i, \max) (dt)



$\Rightarrow \int_0^\infty N \mu f_i s_b dt \approx \frac{N s_b}{s_d^2} \Rightarrow$ if $\ll 1$, then can Taylor expand exponential, and obtain:

$$P_{fix} \approx \left\langle 1 - e^{-\int_0^\infty N \mu f_i(t) dt s_b} \right\rangle \approx \left\langle \int_0^\infty N \mu f_i(t) s_b dt \right\rangle$$

$$= \int_0^\infty N \mu s_b \langle f_i(t) \rangle dt \approx \frac{N}{s_d} s_b \checkmark$$

$\hookrightarrow \frac{1}{N} e^{-s_d t}$

\Rightarrow on the other hand, if $\frac{N s_b}{s_d^2} \gg 1$, then by time $f_i(t)$ reaches $\sim \frac{1}{N s_d}$, it is already very likely to have produced several successful double mutants. Instead, probability is dominated

by lineages that just barely make it to $\int N \mu f_i(t) s_i dt \sim 1$
 & have a high probability of producing successful double mutants.

\Rightarrow recall that for $f_i \ll \frac{1}{Ns_d}$, lineage is essentially neutral:

w/ prob $\frac{1}{T}$ it will survive for $\sim T$ generations & reach ~~size~~

$$\text{size} \sim \frac{T}{N}. \Rightarrow \int_0^{\infty} f_i(t) dt \sim \frac{T^2}{N} \text{ w/ prob } \frac{1}{T}.$$

$$\Rightarrow \int_0^{\infty} N \mu f_i(t) s_b dt \approx N \mu \frac{T^2}{N} s_b \sim 1 \text{ when } T \sim \frac{1}{\sqrt{\mu s_b}}.$$

$$\Rightarrow \text{Pest} \sim (1 - e^{-1}) \frac{1}{T} \approx \sqrt{\mu s_b} \checkmark.$$

\Rightarrow total rate of valley crossing (summing over many downhill mutations)

$$\text{is then } N \mu_{0 \rightarrow 1} P_{\text{fix}} \sim \begin{cases} \frac{N \mu_{0 \rightarrow 1} \mu_{1 \rightarrow 2} s_b}{s_d} & \text{if } \frac{N \mu_{1 \rightarrow 2} s_b}{s_d^2} \ll 1 \\ N \mu_{0 \rightarrow 1} \sqrt{\mu_{1 \rightarrow 2} s_b} & \text{else} \rightarrow \text{happens when } s_b \text{ is large!} \end{cases}$$