

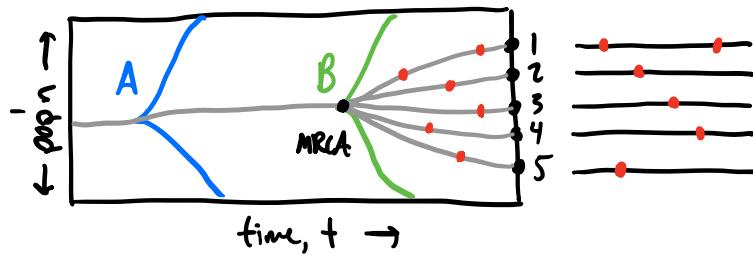
Chapter 13

The independent sites approximation

Last time: Extending coalescent approaches for:

① Selection

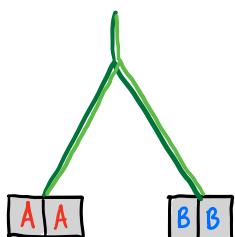
(successive mutations
regime)



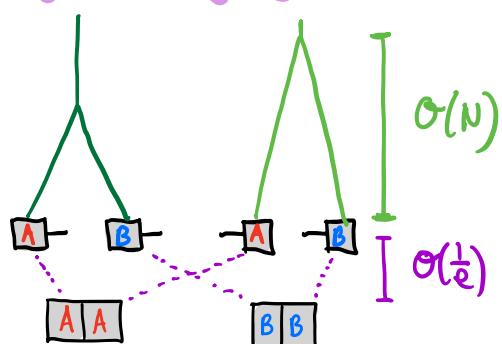
- OR -

② Recombination

$N\mu \ll 1$ (effectively
ascaval)



$N\mu \gg 1$ (effectively
independent)



Today: can we apply similar idea in forward-time picture?

(working directly w/ $f(\vec{g}, t)$)

\Rightarrow i.e., can we find a **high recombination limit** of our multi-locus model of evolution?

$$\frac{d\vec{g}(t)}{dt} = \sim(x - \bar{x}) + \sim L^* \mu + e + \sim \frac{z}{\sqrt{N}}$$

$e = \infty ?$

and therefore treat **sel'n** & **recomb'n**
 — and **mut'n** & **drift** —
 @ the same time?

Answer: Yes!

\Rightarrow to start, consider 2-locus model

w/o **selection** or **mutation** (i.e. genotypes already exist)

\Rightarrow 4 genotypes: $\vec{g} = (0,0), (1,0), (0,1), (1,1)$

\Rightarrow 4 genotype freqs: $f_{00}, f_{10}, f_{01}, f_{11}$

Multi-locus SDEs reduce to:

$$(i) \quad \frac{\partial f_{11}}{\partial t} = e \left[f_{10}f_{01} - f_{11}f_{00} \right] + \sqrt{\frac{f_{11}}{N}} \eta_{11} - f_{11} \sum_{j \neq 1} \sqrt{\frac{f_j}{N}} \eta_j$$

recombination

genetic drift

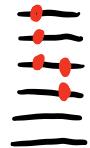
$$(ii) \quad \frac{\partial f_{10}}{\partial t} = e \left[f_{11}f_{00} - f_{10}f_{01} \right] + \sqrt{\frac{f_{10}}{N}} \eta_{10} - f_{10} \sum_{j \neq 1} \sqrt{\frac{f_j}{N}} \eta_j$$

$$(iii) \quad \frac{\partial f_{01}}{\partial t} = e \left[f_{11}f_{00} - f_{10}f_{01} \right] + \sqrt{\frac{f_{01}}{N}} \eta_{01} - f_{01} \sum_{j \neq 1} \sqrt{\frac{f_j}{N}} \eta_j$$

$$(iv) \quad \frac{\partial f_{00}}{\partial t} = e \left[f_{10}f_{01} - f_{11}f_{00} \right] + \sqrt{\frac{f_{00}}{N}} \eta_{00} - f_{00} \sum_{j \neq 1} \sqrt{\frac{f_j}{N}} \eta_j$$

\Rightarrow Present day sample = Multinomial (\vec{n}, \vec{f})

$(n_{11}, n_{10}, n_{01}, n_{00})$



\Rightarrow Note: only 3 independent eqs (since $f_{11} + f_{10} + f_{01} + f_{00} = 1$)

\Rightarrow can eliminate $f_{00} = 1 - f_{11} - f_{10} - f_{01}$

& work w/ f_{11}, f_{10}, f_{01}

key idea: f_{11}, f_{10}, f_{01} is not only basis we can work with...

\Rightarrow one alternative that is often used:

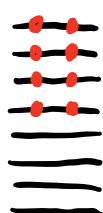
"allele freqs" $\begin{cases} f_1 \equiv f_{11} + f_{10} & \Rightarrow \text{total freq of mutants @ site 1} \\ f_2 \equiv f_{11} + f_{01} & \Rightarrow \quad " \quad \quad \quad \text{site 2} \end{cases}$

$$D \equiv f_{11} - f_1 f_2 \equiv f_{11} f_{00} - f_{10} f_{01} \Rightarrow \text{"linkage disequilibrium" (LD)}$$

\Rightarrow Why? LD is measure of how double mutant deviates from model where mut'n's are independent

e.g. one high-LD scenario:
(D large + positive)

$$\boxed{r^2 = 1}$$



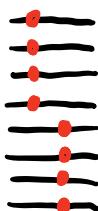
$$f_1 = \frac{1}{2}, f_2 = \frac{1}{2}$$

$$f_{11} = \frac{1}{2}$$

$$D = \frac{1}{2} - \frac{1}{4} = +\frac{1}{4}$$

e.g. another high LD scenario:
(D large & negative)

$$r^2 = \frac{D^2}{f_1(1-f_1)f_2(1-f_2)} = \frac{\frac{1}{16}}{\frac{1}{2}\frac{1}{2}\cdot\frac{1}{2}\cdot\frac{1}{2}} = 1$$

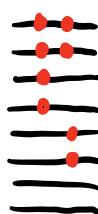


$$f_1 = \frac{1}{2}, f_2 = \frac{1}{2}$$

$$f_{11} = 0$$

$$\Rightarrow D = 0 - \frac{1}{4} = -\frac{1}{4}$$

e.g. a low LD scenario:
(D ≈ 0) (r² = 0)



$$f_1 = f_2 = \frac{1}{2}, f_{11} = \frac{1}{4}$$

$$D = \frac{1}{4} - \frac{1}{4} = 0$$

⇒ sometimes write as correlation coefficient:

$$r \equiv \frac{D}{\sqrt{f_1(1-f_1)f_2(1-f_2)}}$$

$$\Rightarrow r^2$$

(doesn't care about sign)

⇒ why is f_1, f_2, D a good basis?

⇒ let's rewrite our SDEs using def'ns:

$$f_1 \equiv f_{11} + f_{10}, \quad f_2 \equiv f_{11} + f_{01}, \quad D = f_{11} - f_1 f_2$$

$$\frac{\partial f_1}{\partial t} = \frac{\partial f_{11}}{\partial t} + \frac{\partial f_{10}}{\partial t} = \rho [f_{10}f_{01} - f_{11}f_{00}] + \rho [f_{11}f_{00} - f_{10}f_{01}] + \text{noise}$$

~~$f_{10}f_{01} - f_{11}f_{00}$~~ ~~$f_{11}f_{00} - f_{10}f_{01}$~~

$$= 0 + \text{noise}$$

$$\Rightarrow \frac{\partial f_2}{\partial t} = 0 + \text{noise}$$

$$\Rightarrow \frac{\partial D}{\partial t} = \frac{\partial}{\partial t} [f_{11} - f_1 f_2] = \frac{\partial f_{11}}{\partial t} - f_2 \frac{\partial f_1}{\partial t} - f_1 \frac{\partial f_2}{\partial t}$$

$\overset{0+\text{noise}}{\nearrow}$ $\overset{0+\text{noise}}{\nearrow}$

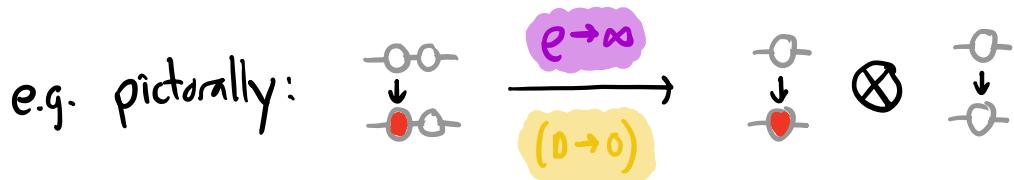
$$= -\rho D + \text{noise}$$

\Rightarrow i.e. recombination cannot change mutation frequencies

\Rightarrow can only change linkage disequilibrium!

\Rightarrow simple deterministic behavior: $D(t) = D(0) e^{-\rho t}$

\Rightarrow suggests that if $\rho \rightarrow \infty \Rightarrow D(t) \approx 0$,
 and maybe can treat 2-locus system as
direct product of independent single locus models:



Mathematically: $f_{11}(t) \approx f_1(t)f_2(t)$, $f_{10}(t) \approx f_1(t)[1 - f_2(t)]$, ...

$$\text{w/ } \frac{df_i}{dt} \approx \sqrt{\frac{f_i(1-f_i)}{N}} \eta_i(t) \quad (i=1,2)$$

\Rightarrow known as "linkage equilibrium", "free recombination",
 "independent sites", "unlinked", etc.

\Rightarrow can check validity using self-consistency argument:

\Rightarrow assume $D \equiv f_{11} - f_{10}f_{01}$ is small &
 calculate next order correction

\Rightarrow correction is known as Quasi-Linkage Equilibrium (QLE)

\Rightarrow easiest to see QLE for rare mutations ($f_1, f_2 \ll 1$)

\Rightarrow then SDEs reduce to:

$$\begin{aligned}\frac{df_1}{dt} &= \frac{df_{11}}{dt} + \frac{df_{10}}{dt} = \sqrt{\frac{f_{11}}{N}} \eta_{11}(t) + \sqrt{\frac{f_{10}}{N}} \eta_{10}(t) \\ &= \sqrt{\frac{f_1 f_2 + D}{N}} \eta_{11}(t) + \sqrt{\frac{f_1 - f_1 f_2 - D}{N}} \eta_{10}(t) \equiv \sqrt{\frac{f_1}{N}} \eta_1(t)\end{aligned}$$

$$\Rightarrow \text{i.e., define } \eta_1(t) = \sqrt{\frac{N}{f_1}} \left[\sqrt{\frac{f_1 f_2 + D}{N}} \eta_{11} + \sqrt{\frac{f_1 - f_1 f_2 - D}{N}} \eta_{10} \right]$$

(which satisfies $\langle \eta_1 \rangle = 0$, $\langle \eta_1^2 \rangle = 1$)

\Rightarrow Similarly,

$$\frac{df_2}{dt} = \sqrt{\frac{f_2}{N}} \eta_2(t) \quad \text{w/} \quad \eta_2 = \sqrt{\frac{N}{f_2}} \left[\sqrt{\frac{f_1 f_2 + D}{N}} \eta_{11} + \sqrt{\frac{f_2 - f_1 f_2 - D}{N}} \eta_{10} \right]$$

$$\Rightarrow \langle \eta_1 \eta_2 \rangle = \sqrt{\frac{N^2}{f_1 f_2} \cdot \left(\frac{f_1 f_2 + D}{N} \right)} = \sqrt{\frac{(f_1 f_2 + D)^2}{f_1 f_2}}$$

\Rightarrow Finally, trickiest one:

$$\frac{dD}{dt} = \frac{df_{11}}{dt} - \frac{d}{dt}(f_1 f_2) \approx -eD + \sqrt{\frac{f_{11}}{N}} \eta_{11} - \left\langle \left(\frac{df_1}{dt} \right)_{\text{drift}} \left| \frac{df_2}{dt} \right. \right\rangle_{\text{drift}}$$

$$= -eD - \frac{f_1 f_2 + D}{N} + \sqrt{\frac{f_{11}}{N}} \eta_{11} (+)$$

$$\omega | \quad \langle \eta_{11} \eta_1 \rangle = \sqrt{\frac{f_1 f_2 + D}{f_1}}$$

\Rightarrow in QLE ($f_{11} \approx f_1 f_2 + \text{small correction}$) $\Rightarrow D \ll f_1 f_2$

$$\Rightarrow \langle \eta_1 \eta_2 \rangle = \sqrt{\frac{(f_1 f_2 + D)^2}{f_1 f_2}} \approx \sqrt{f_1 f_2} \ll 1 \quad (\text{since } f_1, f_2 \ll 1)$$

$$\Rightarrow \langle \eta_1 \eta_{11} \rangle = \sqrt{\frac{f_1 f_2 + D}{f_1}} \approx \sqrt{f_2} \ll 1$$

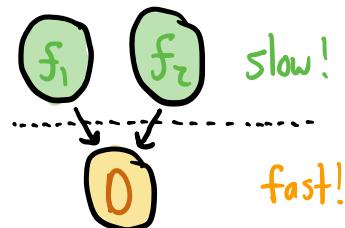
$\Rightarrow f_1 \circ f_2$ independent of each other ($+ D$)

$\Rightarrow f_1 + f_2$ change on drift timescale $T_{\text{drift}} \sim Nf_1, Nf_2$

\Rightarrow When $D \ll f_1 f_2$, LD equation reduces to:

$$\frac{\partial D}{\partial t} = -\frac{f_1 f_2}{N} - \rho D + \sqrt{\frac{f_1 f_2}{N}} \eta_{\parallel}(t)$$

\Rightarrow key idea: dynamics of D relax much faster than f_1, f_2 (since depends on ρ)



\Rightarrow Looks like Brownian particle in quadratic potential (Lecture 6)

$$\text{w/ } \bar{x}_{\text{eff}} = -\frac{f_1 f_2}{N\rho}; \quad r_{\text{eff}} = \ell; \quad D_{\text{eff}} = \frac{f_1 f_2}{N}$$

Solution: (i) $\langle D(t) \rangle = D(0) e^{-\rho t} - \frac{f_1 f_2}{N\rho} (1 - e^{-\rho t}) \xrightarrow{\rho \gg \frac{1}{N}} -\frac{f_1 f_2}{N\rho}$

(ii) $\text{Var}(D(t)) \xrightarrow{\rho \gg \frac{1}{N}} \frac{D_{\text{eff}}}{2r_{\text{eff}}} = \frac{\frac{f_1 f_2}{N}}{\frac{\ell}{2\rho}} = \frac{f_1 f_2}{2N\rho}$

(iii) $\text{Cov}(D(t+\tau), D(t)) = \text{Var}(D(t)) e^{-\rho \tau}$

\Rightarrow QLE ($f_{11} \approx f_1 f_2 + \text{small correction}$) self-consistent if:

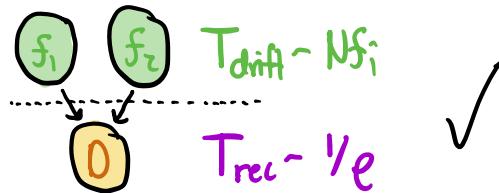
$$\Rightarrow |\langle D \rangle \pm \sqrt{\text{Var}(D)}| \ll f_1 f_2$$

$$\left| -\frac{f_1 f_2}{N e} \pm \sqrt{\frac{f_1 f_2}{2 N e}} \right| \ll f_1 f_2$$

$$\Rightarrow N e \gg \frac{1}{f_1 f_2} \gg 1$$

(like coalescent result, but now depends on f_1, f_2 !)

\Rightarrow separation of timescales
also self-consistent since



\Rightarrow we're done! showed that:

$$\begin{array}{ccc} \text{---} & \xrightarrow{\epsilon \gg 1/N f_1 f_2} & \text{---} \\ \text{---} \downarrow & & \downarrow \\ \text{---} & \otimes & \text{---} \\ \text{---} & & \text{---} \end{array}$$

$(D \ll f_1 f_2)$

$f_{11}(t) \approx f_1(t) f_2(t)$

i.e.

$$\frac{df(\vec{q})}{dt} = \sim \epsilon + \sim \frac{\vec{q}}{J_N} \quad \rightarrow \quad \frac{\partial f_1}{\partial t} = \sqrt{\frac{f_1}{N}} \eta_1(t) \quad \otimes \quad \frac{\partial f_2}{\partial t} = \sqrt{\frac{f_2}{N}} \eta_2(t)$$

\Rightarrow can use same argument for selection too!

\Rightarrow e.g. if $X(\vec{g}) \equiv S_1 g_1 + S_2 g_2$, can show:

$$(i) \frac{\partial f_1}{\partial t} = \frac{\partial S_{11}}{\partial t} + \frac{\partial S_{10}}{\partial t} \approx S_1 f_1 + S_2 f_{11} + \text{noise}$$

$$(ii) \frac{\partial f_2}{\partial t} = S_2 f_2 + S_1 f_{11} + \text{noise}$$

$$(iii) \frac{\partial D}{\partial t} = \frac{\partial S_{11}}{\partial t} - S_1 \frac{\partial f_2}{\partial t} - f_2 \frac{\partial f_1}{\partial t} \approx (S_1 + S_2 - \rho) D + \text{noise}$$

\Rightarrow if $\rho \gg S_1 + S_2 \Rightarrow D(+ \rightarrow 0)$

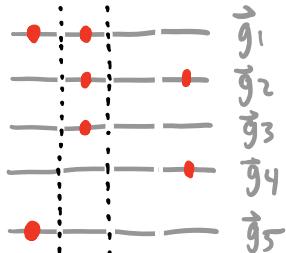
\Rightarrow More generally, if recombination is faster than all other timescales \Rightarrow sites evolve independently

\Rightarrow in practice, people often take this argument & run w/ it for entire genome (rarely check, since QLE is hard!)

Linkage equilibrium approx ("independent sites")



\Rightarrow Present-day sample:



mutations assigned
independently
based on $\{f_e(+)\}$

\Rightarrow Data completely summarized by $\{n_e\}$: "mutation counts"

$$\Pr[n_e = k] = \int \binom{n}{k} f_e^k (1-f_e)^{n-k} p_e(f_e) df_e$$

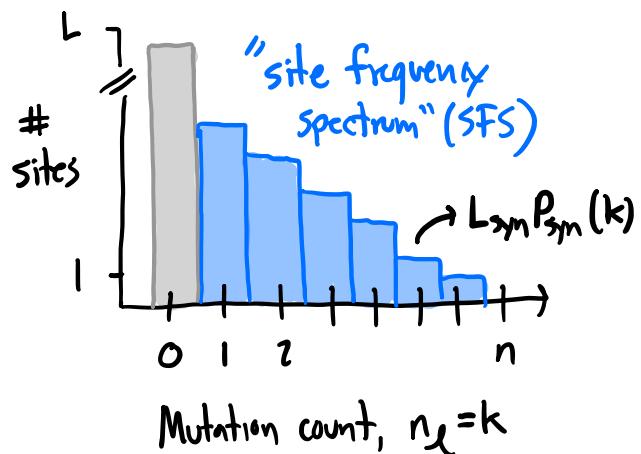
$$\frac{df_e}{dt} = S_e f_e (1-f_e) + \mu_e (1-2f_e) + \sqrt{\frac{f_e(1-f_e)}{N(+)}} \eta(+)$$

\Rightarrow common to group "similar" sites together

e.g. all synonymous sites

$$\downarrow s_e \approx 0$$

$$P_{\text{syn}}(k) = \int \binom{n}{k} f^k (1-f)^{n-k} p(f|s=0) df$$



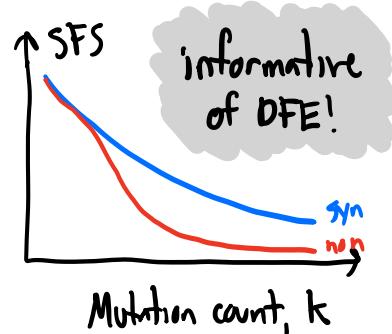
$$\Rightarrow \text{e.g. if } N(t) \approx N \Rightarrow P_{\text{syn}}(k) = \frac{2N\mu}{k}$$

\Rightarrow can do same thing for nonsynonymous mut'ns:

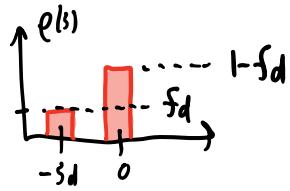
$$P_{\text{non}}(k) = \iint \binom{n}{k} f^k (1-f)^{n-k} p(f|s) \rho(s) df ds$$

$$\approx \int_0^\infty \frac{2N\mu e^{-s} ds}{k(1+2Ns/n)^k} \quad \text{if } N(t) \approx N + k \ll n$$

$$\Rightarrow \frac{P_{\text{non}}(k)}{P_{\text{syn}}(k)} = \int_0^\infty \frac{e^{-s} ds}{(1+2Ns/n)^k}$$



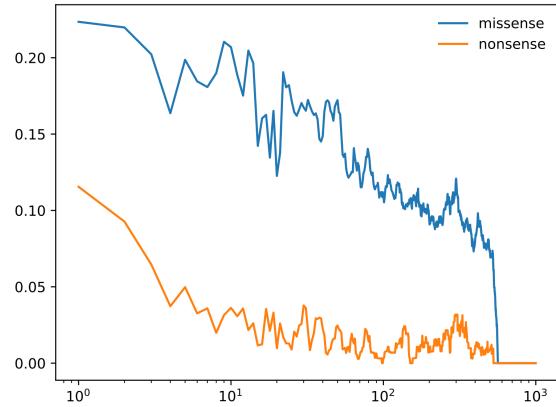
e.g. if $\rho(s) = (1-f_d)\delta(s) + f_d\delta(s+s_d)$



$$\Rightarrow \frac{P_{non}(k)}{P_{syn}(k)} = (1-f_d) + f_d e^{-k \log(1+2Ns/\eta)}$$

e.g. real data from
a human gut bacterium:
(*Bacteroides fragilis*)

$$\frac{P_n(k)}{P_s(k)}$$



Mutation count, k

\Rightarrow informative of "constraint" (strong negative sel'n)

\Rightarrow can coarse-grain over smaller subsets of sites
to look for constraint on smaller regions (e.g. genes)

\Rightarrow why? strongly constrained \approx important for organism

\Rightarrow when do we expect independent sites approx to work?

\Rightarrow need $\rho_{\text{eff}} = r \Delta l$ large for all pairs of SNVs

\downarrow
recombination
rate per site

distance between SNVs
 $(\approx 1/\pi = 1/2N_N)$

$$\Rightarrow N\rho_{\text{eff}} \gg 1 \Rightarrow \frac{\Gamma}{N} \gg 1$$

(in most organisms we've measured, $\frac{\Gamma}{N} \sim \mathcal{O}(1)!$)

$$\Rightarrow \rho_{\text{eff}} \gg s \Rightarrow Ns \ll \frac{\Gamma}{N} \sim \mathcal{O}(1)$$

\Rightarrow bad approximation for strong selection!

\Rightarrow need to turn to other approaches...