

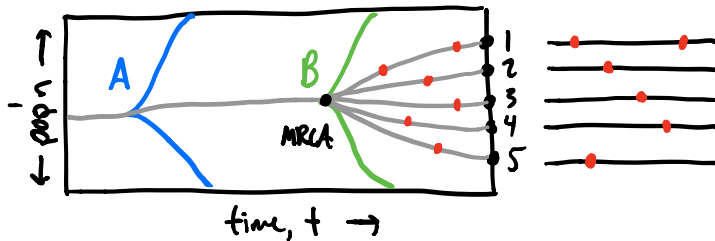
Announcements: Final Project Change! (details in email or Slack)

⇒ Problem set 4 DUE today or Thurs ↷

Last time: Extending coalescent approaches for:

① Selection

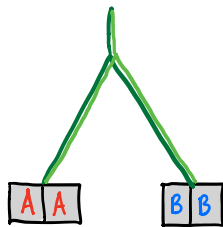
(Successive mutations regime)



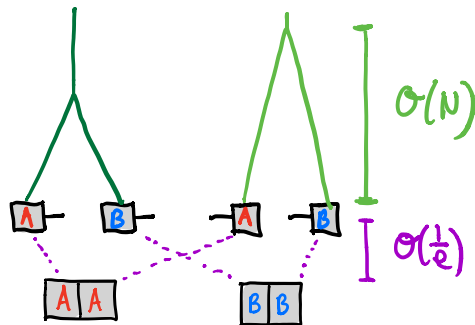
- OR -

② Recombination

$N_e \ll 1$ (effectively asexual)



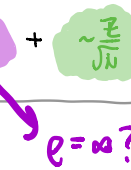
$N_e \gg 1$ (effectively independent)



Today: can we apply similar idea in forward-time picture?
(working directly w/ $f(\vec{g}, t)$)

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

$$\frac{d\vec{g}}{dt} = \underbrace{\sim(x-\bar{x})}_{\text{blue}} + \underbrace{\sim L\mu}_{\text{orange}} + \underbrace{e}_{\text{purple}} + \underbrace{\sim \frac{\sigma^2}{\sqrt{N}}}_{\text{green}}$$



 $e = \infty?$

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

Answer: Yes!

⇒ to start, consider 2-locus model

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

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

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

Multi-locus SDEs reduce to:


$$(i) \quad \frac{df_{11}}{dt} = e \left[f_{10}f_{01} - f_{11}f_{00} \right] + \sqrt{\frac{f_{11}}{N}} \eta_{11} - f_{11} \sum_{\vec{j}} \sqrt{\frac{f_{\vec{j}}}{N}} \eta_{\vec{j}} \quad \text{genetic drift}$$

recombination

$$(ii) \quad \frac{df_{10}}{dt} = e \left[f_{11}f_{00} - f_{10}f_{01} \right] + \sqrt{\frac{f_{10}}{N}} \eta_{10} - f_{10} \sum_{\vec{j}} \sqrt{\frac{f_{\vec{j}}}{N}} \eta_{\vec{j}}$$

$$(iii) \quad \frac{df_{01}}{dt} = e \left[f_{11}f_{00} - f_{10}f_{01} \right] + \sqrt{\frac{f_{01}}{N}} \eta_{01} - f_{01} \sum_{\vec{j}} \sqrt{\frac{f_{\vec{j}}}{N}} \eta_{\vec{j}}$$

$$(iv) \quad \frac{df_{00}}{dt} = e \left[f_{10}f_{01} - f_{11}f_{00} \right] + \sqrt{\frac{f_{00}}{N}} \eta_{00} - f_{00} \sum_{\vec{j}} \sqrt{\frac{f_{\vec{j}}}{N}} \eta_{\vec{j}}$$

\Rightarrow Present day sample = Multinomial (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...

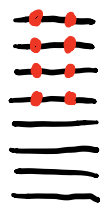
⇒ one alternative that is often used:

"allele freqs" $\left\{ \begin{array}{l} f_1 \equiv f_{11} + f_{10} \Rightarrow \text{total freq of mutants @ site 1} \\ f_2 \equiv f_{11} + f_{01} \Rightarrow \text{" " " site 2} \end{array} \right.$

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

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

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



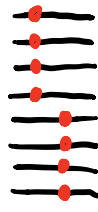
$$f_1 = \frac{1}{2}, f_2 = \frac{1}{2}$$

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

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

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

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



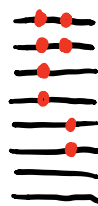
$$f_1 = \frac{1}{2}, f_2 = \frac{1}{2}$$

$$f_{11} = 0$$

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

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

e.g. a low LD scenario:
($D \approx 0$) ($r^2 = 0$)



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

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

\Rightarrow 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)

\Rightarrow why is f_1, f_2, D a good basis?

\Rightarrow let's rewrite our SDEs using def'ns:

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

$$\frac{df_1}{dt} \equiv \frac{df_{11}}{dt} + \frac{df_{10}}{dt} = \cancel{e[f_{10}f_{01} - f_{11}f_{00}]} + \cancel{e[f_{11}f_{00} - f_{10}f_{01}]} + \text{noise}$$

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

$$\Rightarrow \frac{df_2}{dt} = 0 + \text{noise}$$

$$\Rightarrow \frac{dD}{dt} \equiv \frac{d}{dt} [f_{11} - f_1 f_2] = \frac{df_{11}}{dt} - f_2 \frac{df_1}{dt} - f_1 \frac{df_2}{dt}$$

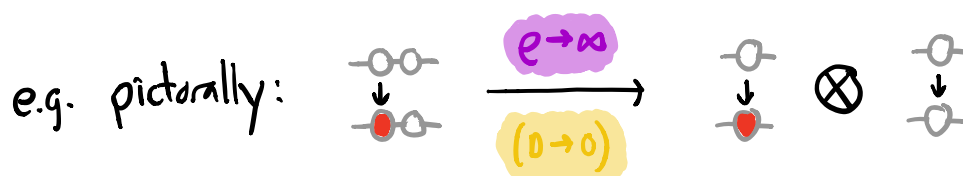
$$= -eD + \text{noise}$$

\Rightarrow i.e. recombination cannot change mutation frequencies

\Rightarrow can only change linkage disequilibrium!

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

\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:

$$f_{11}(t) \approx f_1(t)f_2(t), \quad f_{10}(t) \approx f_1(t)[1 - f_2(t)], \dots$$

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

\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

⇒ correction is known as Quasi-Linkage Equilibrium (QLE)

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

⇒ then SDEs reduce to:

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

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

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

⇒ Similarly,

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

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

\Rightarrow Finally, trickiest one:

$$\begin{aligned} \frac{dD}{dt} &\equiv \frac{df_{11}}{dt} - \frac{d}{dt}(f_1 f_2) \approx -eD + \sqrt{\frac{f_{11}}{N}} \eta_{11} - \left\langle \left(\frac{df_1}{dt} \right)_{\text{diff}} \left(\frac{df_2}{dt} \right)_{\text{diff}} \right\rangle \\ &= -eD - \frac{f_1 f_2 + D}{N} + \sqrt{\frac{f_{11}}{N}} \eta_{11}(t) \end{aligned}$$

$$\text{w/ } \langle \eta_{11} \eta_1 \rangle = \sqrt{\frac{f_1 f_2 + D}{f_1}}$$

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

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

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

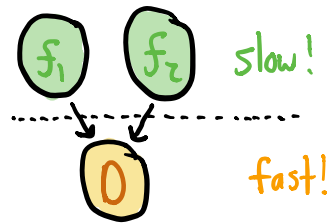
$\Rightarrow f_1 \neq f_2$ independent of each other ($\neq D$)

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

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

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

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



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

$$w/ \bar{X}_{\text{eff}} = -\frac{f_1 f_2}{Ne} ; r_{\text{eff}} = e ; D_{\text{eff}} = \frac{f_1 f_2}{N}$$

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

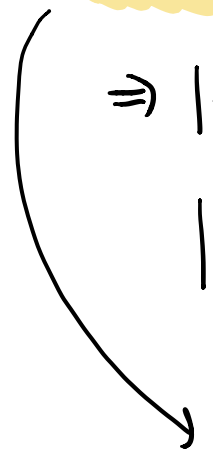
(ii) $\text{Var}(D(t)) \xrightarrow{t \gg \tau} \frac{D_{\text{eff}}}{2r_{\text{eff}}} = \frac{\frac{f_1 f_2}{N}}{2e} = \frac{f_1 f_2}{2Ne}$

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

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

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

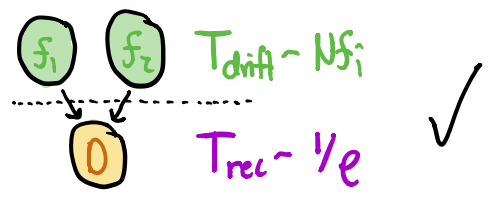
$$\left| -\frac{f_1 f_2}{N_e} \pm \sqrt{\frac{f_1 f_2}{2N_e}} \right| \ll f_1 f_2$$



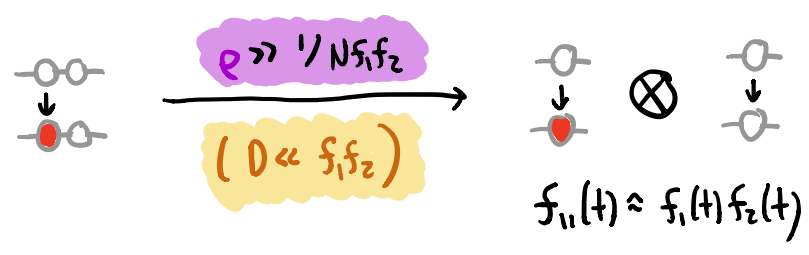
$$\Rightarrow N_e \gg \frac{1}{f_1 f_2} \gg 1$$

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

⇒ separation of timescales also self-consistent since



⇒ we're done! showed that:



i.e.

$$\frac{df_i(t)}{dt} = \sim e + \sim \frac{e}{\sqrt{N}} \rightarrow \frac{df_1}{dt} = \sqrt{\frac{f_1}{N}} \eta_1(t) \otimes \frac{df_2}{dt} = \sqrt{\frac{f_2}{N}} \eta_2(t)$$

⇒ can use same argument for **selection** too!

⇒ e.g. if $X(\vec{g}) \equiv s_1 g_1 + s_2 g_2$, can show:

$$(i) \frac{df_1}{dt} = \frac{ds_{11}}{dt} + \frac{ds_{10}}{dt} \approx s_1 f_1 + s_2 f_{11} + \text{noise}$$

$$(ii) \frac{df_2}{dt} = s_2 f_2 + s_1 f_{11} + \text{noise}$$

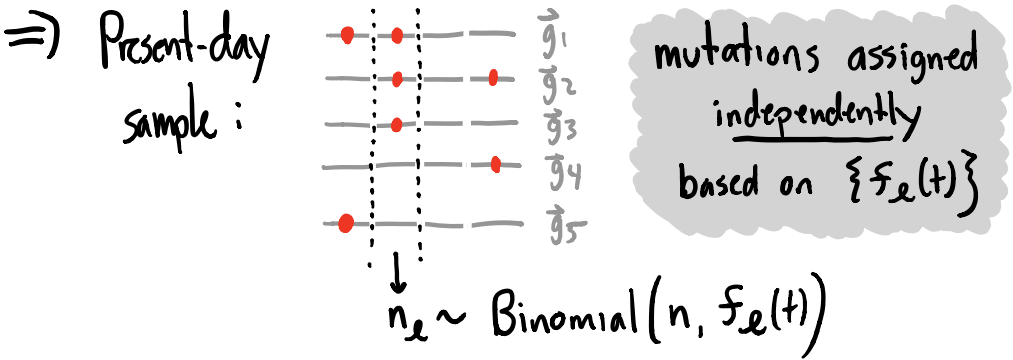
$$(iii) \frac{dD}{dt} = \frac{ds_{11}}{dt} - s_1 \frac{df_2}{dt} - f_2 \frac{df_1}{dt} \approx (s_1 + s_2 - e) D + \text{noise}$$

⇒ if $e \gg s_1 + s_2 \Rightarrow D(t) \rightarrow 0$

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

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

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

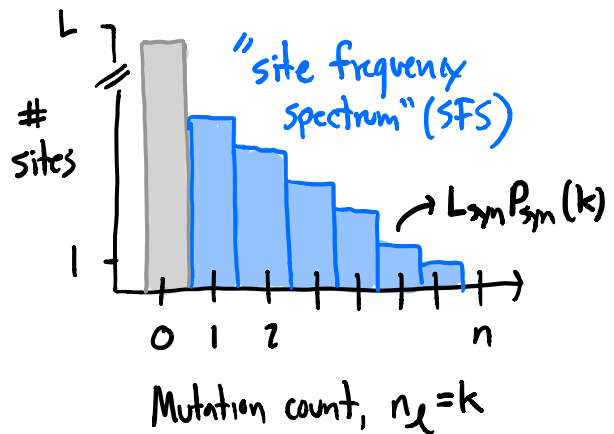
$$\frac{df_e}{dt} = s_e f_e (1-f_e) + \mu_e (1-2f_e) + \frac{\sqrt{f_e(1-f_e)}}{N(t)} \eta(t)$$

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

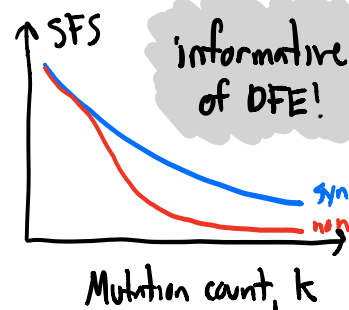
\Rightarrow can do same thing for nonsynonymous mut'ns:

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

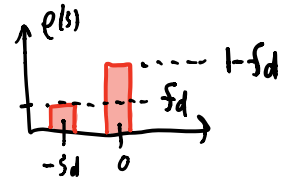
$$\approx \int_0^{\infty} \frac{2N\mu e(-s) ds}{k(1 + 2Ns/n)^k}$$

if $N(t) \approx N + k \ll n$

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



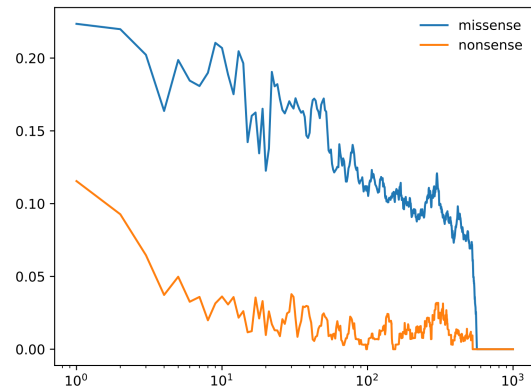
e.g. if $\rho(s) = (1-f_d)\delta(s) + f_d\delta(s+s_d)$



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

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

⇒ when do we expect **independent sites approx** to work?

⇒ need **$\rho_{\text{eff}} \approx r \Delta \ell$** large for all pairs of SNVs
recombination rate per site → distance between SNVs
($\approx 1/\pi \approx 1/2N\mu$)

⇒ $N\rho_{\text{eff}} \gg 1 \Rightarrow \frac{r}{\mu} \gg 1$

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

⇒ $\rho_{\text{eff}} \gg s \Rightarrow Ns \ll \frac{r}{\mu} \sim \mathcal{O}(1)$

⇒ bad approximation for strong selection!

⇒ need to turn to other approaches...