# Quasi Linkage Equilibrium (QLE)
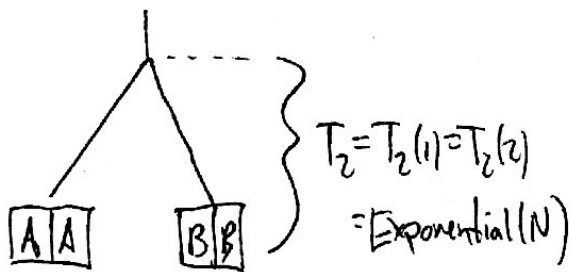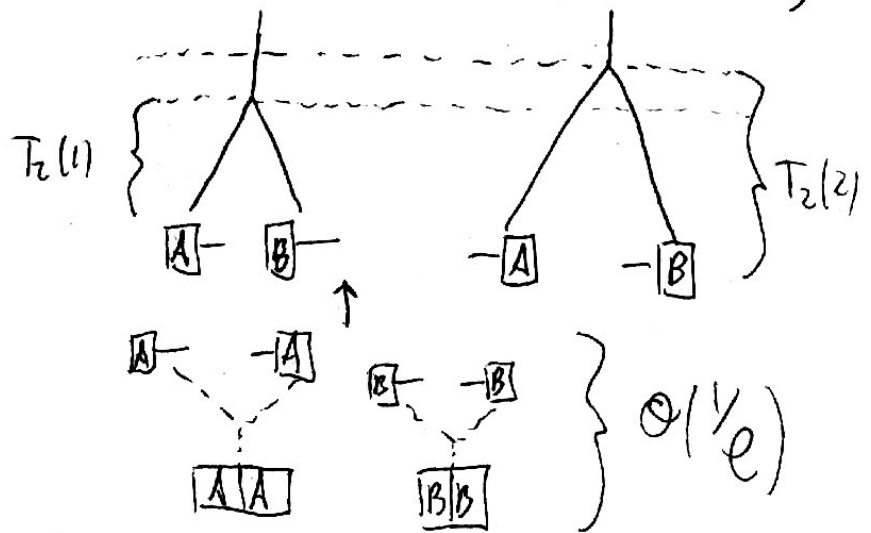
Last time, we talked about coalescent models for neutral recombining genomes, and saw that this led to 2 extreme limits:

$Ne \ll 1$

(effectively asexual)



$\left.\begin{array}{c}\end{array}\right\} T_2 = T_2(1) = T_2(2)$
$= \text{Exponential}(N)$

⇑

not enough time for recombination to occur.

$Ne \gg 1$ (effectively independent)

$T_2(1) \{$

$\} T_2(2)$



$\left.\begin{array}{c}\end{array}\right\} \Theta(1/e)$

⇑

recombination will occur very fast, coarse grain over recomb. timescale.

⟹ in between, $Ne \sim \Theta(1)$ very complicated (ARG)

@ end of day, want mutations in sample, not trees...



genome 1
⋮
genome n

⟹ also hard to get selection into this picture...

So today, want to talk about same concepts from a forward-time perspective, based on genotype frequencies, $f(\vec{g})$

$\Longrightarrow$ to start, let's consider a simple 2-locus model for genotypes $\vec{g} = (0,0), (1,0), (0,1), (1,1)$, w/o selection

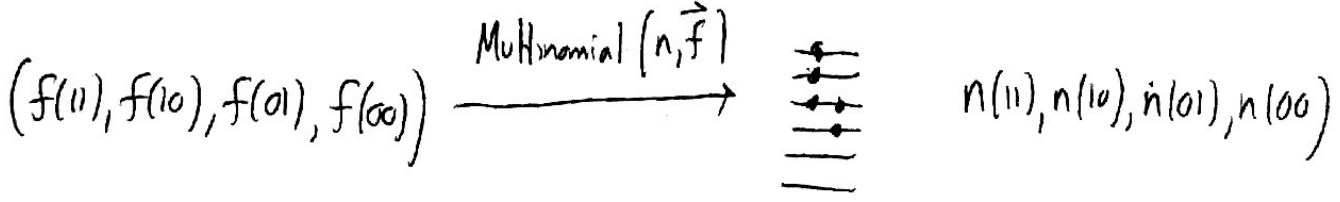then genotype freqs satisfy the system of SDEs:

$$\frac{df(11)}{dt} = \rho\left[f(10)f(01) - f(11)f(00)\right] + \sqrt{\frac{f(11)}{N}}\eta(11) - f(11)\sum_{\vec{g}'}\sqrt{\frac{f(\vec{g}')}{N}}\eta(\vec{g}')$$

$$\frac{df(10)}{dt} = \rho\left[f(11)f(00) - f(10)f(01)\right] + \sqrt{\frac{f(10)}{N}}\eta(10) - f(10)\sum_{\vec{g}'}\sqrt{\frac{f(\vec{g}')}{N}}\eta(\vec{g}')$$

$$\frac{df(01)}{dt} = \rho\left[f(11)f(00) - f(10)f(01)\right] + \sqrt{\frac{f(01)}{N}}\eta(01) - f(01)\sum_{\vec{g}'}\sqrt{\frac{f(\vec{g}')}{N}}\eta(\vec{g}')$$

$$\frac{df(00)}{dt} = \rho\left[f(10)f(01) - f(11)f(00)\right] + \sqrt{\frac{f(00)}{N}}\eta(00) - f(00)\sum_{\vec{g}'}\sqrt{\frac{f(\vec{g}')}{N}}\eta(\vec{g}')$$

and sample comes from:

$$\left(f(11), f(10), f(01), f(00)\right) \xrightarrow{\text{Multinomial}(n,\vec{f})} \quad\quad n(11), n(10), n(01), n(00)$$

this system really has only 3 independent eqs,
since $f(11) + f(10) + f(01) + f(00) = 1$

$\implies$ • e.g. can eliminate $f(00) = 1 - f(11) - f(10) - f(01)$
   & work w/ $f(11), f(10), f(01)$.

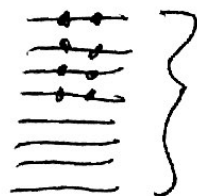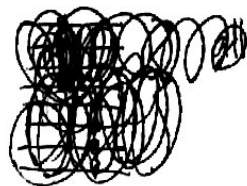However, $(f(11), f(10), f(01))$ is not the only basis we could work with.
   $\implies$ free to choose any other combination.

$\implies$ one combination that is often used:

$$\underbrace{f_1 \equiv f(11) + f(10)}_{\substack{\text{marginal allele freq} \\ \text{of mutations @} \\ \text{site 1}}} \quad , \quad \underbrace{f_2 \equiv f(11) + f(01)}_{\substack{\text{allele freq} \\ \text{@ site 2}}} \quad , \quad \underbrace{D \equiv f(11) - f_1 f_2 = f(11)f(00) - f(10)f(01)}_{\text{"Linkage disequilibrium" (LD)}}$$

$\implies$ LD is measure of how much double mutant frequency deviates
   from a model where mutations are totally independent.

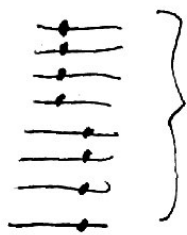e.g. one high-LD scenario
   ($D = $ large & positive)



mutations always appear
together on same genomes.

$$D = \tfrac{1}{2} - \tfrac{1}{2}\tfrac{1}{2} = \tfrac{1}{4}$$

e.g. another high-LD scenario

(LD = large & negative)

mutations always on separate chromosomes

$D = 0 - \frac{1}{2}\frac{1}{2} = -\frac{1}{4}$

$\implies$ Sometimes people measure as correlation coefficient:

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

e.g. in example 1, $r = +1$

example 2, $r = -1$

$\implies$ mutation @ site 1 gives you perfect info about site 2, & vice versa.

Why is $f_1, f_2, D$ a good basis? Rewrite SDEs:

$$\frac{\partial f_1}{\partial t} \equiv \frac{\partial f_{(11)}}{\partial t} + \frac{\partial f_{(10)}}{\partial t} = \rho\left[f_{(0)}f_{(01)} - f_{(00)}f_{(11)}\right] + \rho\left[f_{(00)}f_{(11)} - f_{(10)}f_{(01)}\right]^{\to 0}$$

$+$ noise (later)

$$\frac{\partial f_2}{\partial t} = 0 + \text{noise}, \qquad \frac{\partial D}{\partial t} = \frac{\partial f_{11}}{\partial t} - \frac{\partial f_1}{\partial t}(f_2)^{\to 0} - f_1 \frac{\partial f_2}{\partial t}^{\to 0}$$
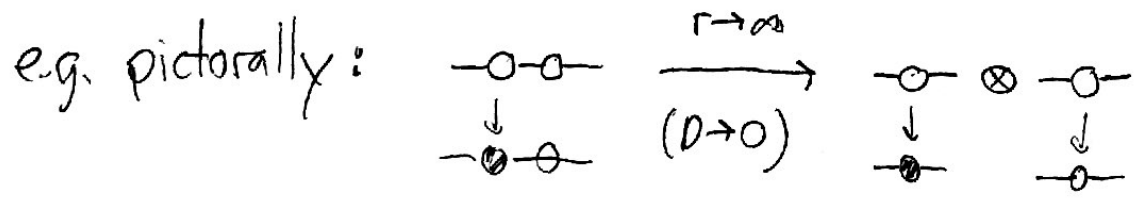
$$= -\rho D + \text{noise.}$$

In other words, ~~recombination~~ recombination cannot change allele frequencies, can only change linkage disequilibrium. change in linkage disequilibrium has very simple form (deterministically)

$$\Rightarrow \frac{dD}{dt} = -rD \Rightarrow D(t) = D(0)e^{-\rho t} \Rightarrow \text{decays to zero exponentially fast in } \rho$$

$\Rightarrow$ suggests that if $\rho \to \infty$ (compared to what?)

$D \to 0$ and maybe we can treat 2 locus system as direct product of single locus systems :

e.g. pictorally :

we call this limit "Linkage equilibrium", "free recombination", "independent sites", etc.

we can check this assumption using our now-familiar self consistency argument: assume that $D$ is small & calculate next order correction $\Rightarrow$ this correction is called

"Quasi-linkage equilibrium"

(QLE)

Easiest to see QLE if we focus on rare mutations, $f_1, f_2 \ll 1$. Then SDEs reduce to (w/ noise now)

$$\frac{df_1}{dt} = \sqrt{\frac{f^{(11)}}{N}} \eta^{(11)} + \sqrt{\frac{f^{(10)}}{N}} \eta^{(10)} = \sqrt{\frac{f_1 f_2 + D}{N}} \eta^{(11)} + \sqrt{\frac{f_1(1-f_2)-D}{N}} \eta^{(10)}$$

$$\cong \sqrt{\frac{f_1}{N}} \tilde{\eta}_1$$

$$\frac{df_2}{dt} = \sqrt{\frac{f^{(11)}}{N}} \eta^{(11)} + \sqrt{\frac{f^{(01)}}{N}} \eta^{(01)} = \sqrt{\frac{f_1 f_2 + D}{N}} \eta^{(11)} + \sqrt{\frac{f_2(1-f_1)f_2 (1-f_1)f_2}{N}} \eta^{(01)} \equiv \sqrt{\frac{f_2}{N}} \tilde{\eta}_2$$

$$\frac{df_{11}}{dt} = -\rho D + \sqrt{\frac{f^{(11)}}{N}} \eta^{(11)} = \quad \rho f_1 f_2 - \rho f_{11} + \sqrt{\frac{f_{11}}{N}} \eta_{11}$$

w/ $\langle \eta_{11} \tilde{\eta}_1 \rangle = \sqrt{\frac{f_{11}}{f_1}}$, $\langle \eta_{11} \tilde{\eta}_2 \rangle = \sqrt{\frac{f_{11}}{f_2}}$, $\langle \tilde{\eta}_1 \tilde{\eta}_2 \rangle = \frac{f_{11}}{\sqrt{f_1 f_2}}$

$\Rightarrow$ in QLE, we will assume that dynamics of $D$ ($f_{11}$) will relax much faster than dynamics of $f_1, f_2$

$\Rightarrow$ since neutral, know that $f_1, f_2$ change on timescale

$$T_{drift} \sim N f_1, N f_2$$

$\Rightarrow$ on timescales $\ll T_{drift}$, $f_1$ & $f_2$ are effectively const.

$\Rightarrow$ equation for $f_{11}$ looks like Branching process w/ mutation w/ effective params: $\mu_e = \rho f_1 f_2$, $S_e = -\rho$

$\Rightarrow$ solution from Lecture 8, p.4:

Gamma Dist'n w/ shape: ~~~~~ $\alpha = 2N\rho f_1 f_2$

and ~~~~~ $f_{11, max} = \dfrac{1 - e^{-\rho t}}{2N\rho} \rightarrow \dfrac{1}{2N\rho}$

$\Rightarrow$ $\langle f_{11} \rangle = f_1 f_2$, $Var(f_{11}) = \dfrac{f_1 f_2}{2N\rho}$

$\Rightarrow$ $\langle D \rangle = 0$ $\quad Var(D) = \dfrac{f_1 f_2}{2N\rho}$

we also know that relaxation timescale is $\sim 1/\rho$

$\Rightarrow$ Now can check our assumptions: $\frac{1}{\rho} \ll T_{drift} \sim Nf_1, Nf_2$

$\Rightarrow$ QLE holds if $N\rho f_1, N\rho f_2 \gg 1$.

$\quad\quad\quad\quad\quad\quad\quad\quad$ $\llcorner\rightarrow$ similar to coalescent picture, except now explicit $f$ dependence.

$\Rightarrow$ if this is true, then $\langle \hat{\eta}_2 \hat{\eta}_1 \rangle = \sqrt{\frac{f_{11}^2}{f_1 f_2}} = \sqrt{\frac{1}{(N\rho f_1)(N\rho f_2)}} \ll 1$

and $\frac{df_1}{dt}, \frac{df_2}{dt}$ equations really decouple! $\checkmark$.

so we showed that $\quad$ —o—o— $\longrightarrow$ —o—⊗—o—
$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ 1 $\;$ 2 $\quad\quad\quad\quad\quad\quad$ 1 $\quad\quad$ 2

but only ~~correlates~~ for sufficiently high freqs where can coarse-grain over recombination time.

$\Rightarrow$ can do same argument for selection, e.g. $X(\vec{g}) = S_1 g_1 + S_2 g_2$

then can show:

$$\frac{\partial f_1}{\partial t} = S_1 f_1 + S_2 D + \text{noise}$$

$$\frac{\partial f_2}{\partial t} = S_2 f_2 + S_1 D + \text{noise}$$

$$\frac{\partial D}{\partial t} = \frac{\partial f_{11}}{\partial t} - f_2 \frac{\partial f_1}{\partial t} - f_1 \frac{\partial f_2}{\partial t} = \text{〰〰〰〰〰}$$

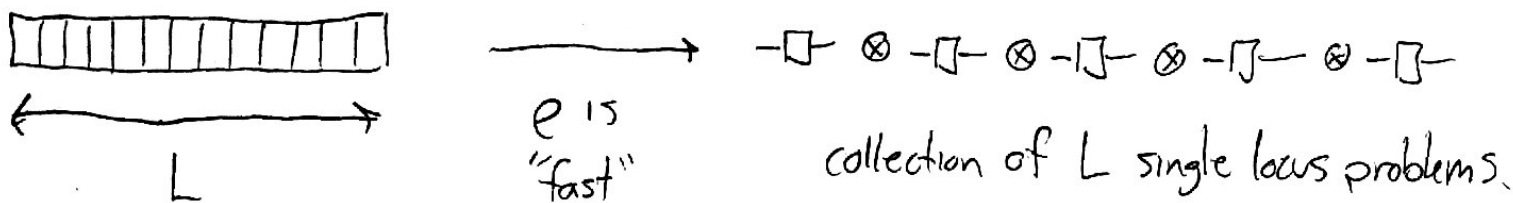$$= (S_1 + S_2 - r) D + \text{noise}.$$

$\Rightarrow$ so if $r > S_1 + S_2 \Rightarrow D(t) \to 0$ w/ time.

$\Rightarrow$ General conclusion: if $r$ is faster than
all other timescales in system
then loci evolve independently!

$\Rightarrow$ in practice, people take this argument & run w/ it
for entire genome : (though people rarely check it like
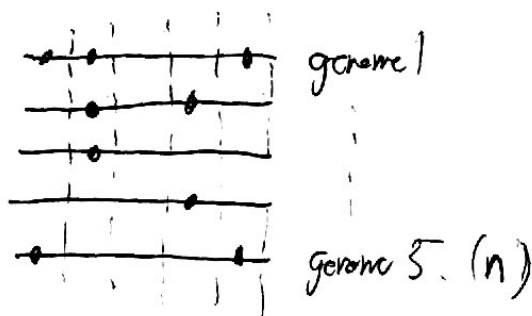we just did because QLE is hard...)

# Linkage equilibrium approx ("independent sites")

L

$\ell$ is "fast"

collection of L single locus problems.

$\Rightarrow$ when it works, one of the most powerful approximations in pop gen. Since it lets us use single locus results to look @ real data

$\hookrightarrow$ crazy if you think about it.

now when we draw a sample of individuals, we can assign mutations independently given current allele frequencies, $\{f_\ell(t)\}$

genome 1

genome 5. (n)

$\Rightarrow$ by definition, $D \hat{=} 0$ & no information in haplotype structure.

$\Rightarrow$ instead, data can be completely summarized by $n_\ell = $ # of individuals w/ mutation @ site $\ell$.
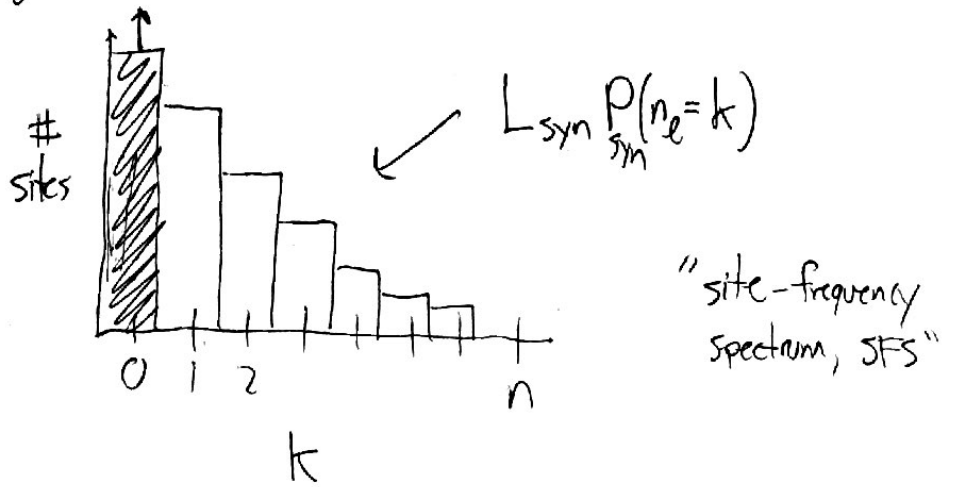
$\Rightarrow$

$\Rightarrow$ $Pr(n_\ell = k) = \int \binom{n}{k} f_\ell^k (1-f_\ell)^{n-k} \rho_\ell(f_\ell) df_\ell$

@ this point, people often gap "similar" sites together.

eg. all synonymous sites
  (putatively neutral, )
  
  so $S_e = 0$



$L_{syn} \, \underset{syn}{P}(n_e = k)$

"site-frequency spectrum, SFS"

$\Rightarrow$ since there are lots of synonymous sites, & each one is an independent draw from $P_{syn}(n_e = k)$, then across genome, we get a self-averaged version of $P_{syn}(n_e = k)$, even from just 1 population! this means we can estimate demography, since

$$P_e(f_e) \Longleftrightarrow \frac{df_e}{dt} = \mu_e + \sqrt{\frac{f_e(1-f_e)}{N(t)}} \, \eta(t)$$

( mapping often not possible in closed form, but can do numerically )

$\Rightarrow$ if $N(t) = N \Rightarrow Pr(n_e = k) = \frac{2N\mu}{k} \rightarrow \frac{2N\mu}{f}$ when n large.

$\searrow$ $2N\mu$ when n=2.

$\Rightarrow$ similarly, $\Pi_{syn}$ self averages to $2N\mu$, even w/ just 2 samples!

can do same thing for <u>non-synonymous mutations</u>

$$Pr[n_e = k] = \iint \binom{n}{k} f_e^k (1-f_e)^{n-k} p(f_e|s_e) \rho(s_e) df_e ds_e$$

↓ random freq @ locus e  → selection coeff @ locus e

$\Rightarrow$ e.g. if $\rho(s) =$ ~~⦿⦿⦿⦿⦿⦿⦿⦿⦿⦿~~ $(1-\lambda_d)\delta(s) + \lambda_d \delta(s+s_d)$

↓ neutral mutations   ↓ strongly deleterious mut. (e.g. LOF mutant)

$\Rightarrow p(f) = \dfrac{2N\mu(1-\lambda_d)}{f} + \dfrac{2N\mu \lambda_d e^{-Nsf}}{f} = \begin{cases} \dfrac{2N\mu}{f} & \text{for } f \ll \frac{1}{NS} \\ & \text{(like synonymous)} \\ \dfrac{2N\mu(1-\lambda_d)}{f} & \text{for } f \gg \frac{1}{NS} \end{cases}$

$\Rightarrow$ similarly for $\pi$:

$$\pi_{non} = (1-\lambda_d)2N\mu + \lambda_d \frac{2N\mu}{2NS}$$

$\Rightarrow$ typically we don't know $2N\mu$, but can estimate it from $\pi_{syn}$

$$\Rightarrow \pi_{non}/\pi_{syn} = (1-\lambda_d) + \frac{\lambda_d}{NS} \qquad (NS \gg 1)$$

$\dfrac{\Pi_n}{\Pi_s}$ provides a measure of "constraint", how much negative selection is going on in population w/in $T_{MRCA}$.

e.g. in bacteria (E. coli in 2 different people's guts) often find $\Pi_{non}/\Pi_{syn} \simeq 0.1$.

⟹ this is pretty crazy... suggests that $\lambda_d \gtrsim 0.9$

i.e., 90% of all amino acid changes in bacteria ~~bbb~~
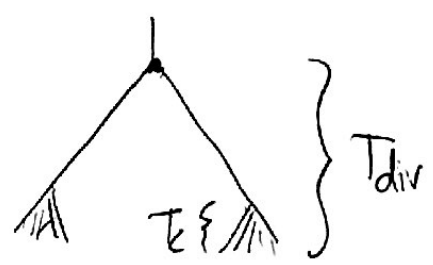
~~[scribble]~~ are sufficiently strongly deleterious

~~that they are~~ ~~so~~ strongly selected against.

⟹ in practice, people often coarse grain sites & look for constraint on even smaller portions of genome.

⟹ reason is that strongly constrained $\approx$ important for organism
(interesting biology)

⟹ can also do same thing w/ substitutions between 2 species ($dN/dS$)

⟹ more time for mutations to occur, better signal


$\left.\vphantom{\begin{array}{c}\\ \\ \\ \end{array}}\right\} T_{div}$

when do we expect QLE to work?

$\Rightarrow$ • think about collection of 2 locus mini problems.

$$\rho_{eff} = r \cdot \Delta \ell \rightarrow \text{distance between sites}$$
$$\hookrightarrow \text{recombination rate per site.}$$

if guestimate $N$ from $\pi_{syn}$ :    $N = \dfrac{\pi_{syn}}{2\nu}$

$\Rightarrow$ $D_{max} \sim \dfrac{1}{N\rho} \sim \dfrac{\nu}{r(\pi \Delta \ell)} \approx \dfrac{\nu}{r}$   if $\Delta \ell \approx \dfrac{1}{\pi}$

(neighboring SNPs)

in most organisms we've measured, $\dfrac{\nu}{r} \sim \mathcal{O}(1)$

(weird right?) $\rightarrow$ so neutral muts right on boundary of ok.

$\Rightarrow$ selected mutations needs: $\dfrac{s}{r} \ll 1$

$\Rightarrow$ if $r \sim \nu \sim 10^{-8} - 10^{-10}$ $\Rightarrow$ need $s \ll 10^{-8} - 10^{-10}$

$\Rightarrow$ bad approximation most of time. need some
other method to predict evolution in this case.