- §2.6: The Coalescent

- Let $\theta = 4Nu$
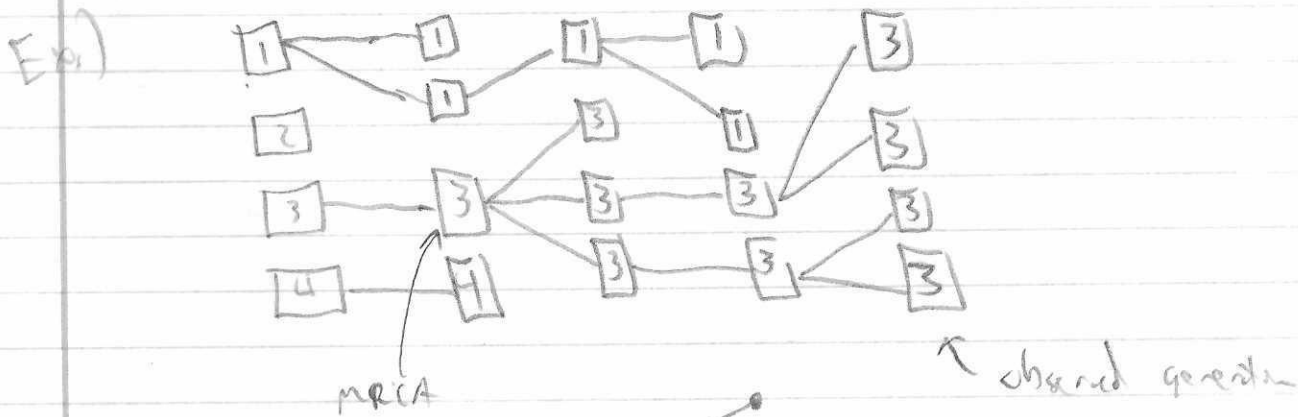
Recall, when drift = mutation, $H = \dfrac{4Nu}{1 + 4Nu}$

As $\theta \uparrow$, we get more heterozygosity (approaching 1)
(large pop or large mutation rate)

As $\theta \downarrow$, we get less heterozygosity
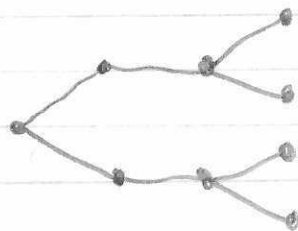(small pop or small mutation rate)

- Before, we estimated $H$ from observations to get estimate of $\theta$

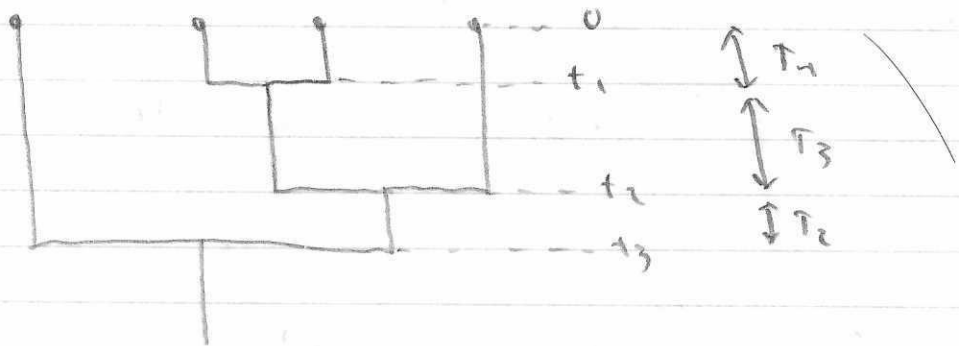- Here, we will estimate $\theta$ by coalescent theory.

- Coalescent: Lineage of alleles traced back in trees to most recent common ancestor

Ex.)



MRCA

← observed generation

Coalescent

- For typical $N$, we coalescent will look like this



$t_i$ is in units of _generations_ under wright-Fisher model.

- We use mutations in the lineage to
  ① Estimate coalescent (called phylogenetics)
  ② Estimate $\theta$
       (polymorphic)
- Let $S_n = \#$ segregating sites in a population
- Goal: use $S_n$ to estimate $\theta$.
- Let $T_c = $ total tree in coalescent

[↑] In above figure, $T_c = 4t_1 + 3(t_2 - t_1) + 2(t_3 - t_2)$

Let $T_i = $ tree with $i$ alleles in pop

$$T_c = 4T_4 + 3T_3 + 2T_2$$

- Expected $\#$ mutations in whole coalescent is $\underbrace{T_c u}$

   mutations
   generation

- We will show that (for a sample to 4 alleles)

$$E[T_c] = 4N\left(1 + \tfrac{1}{2} + \tfrac{1}{3}\right) = \frac{44N}{6} \qquad \text{so } uT_c = \theta \frac{11}{6}$$

$$\theta = 4Nu$$

$$= u E[T_c]$$

- Thus $\quad E[S_+] = \theta \, 11/6$

$\uparrow$ easy way to estimate $\theta$

- Goal: Get $E[T_c]$ for any coalescent

- First interval $T_n$, second $T_{n-1}, \ldots, T_2$

- Consider $n$ alleles (pop size is $N$)

$Pr(\text{allele 1 and 2 have different parents}) = 1 - \frac{1}{2N} = \frac{2N-1}{2N}$

$Pr(\text{allele 3 diff} \mid 1 \text{ and } 2 \text{ diff}) = 1 - \frac{2}{2N} = \frac{2N-2}{2N} \quad Pr(\text{same parent})$

etc...

$Pr(\text{all diff}) = \left(1 - \frac{1}{2N}\right)\left(1 - \frac{2}{2N}\right) \cdots \left(1 - \frac{n-1}{2N}\right) \quad$ can't be one of the $n-1$ already chosen

$\approx 1 - \frac{1}{2N} - \frac{2}{2N} - \cdots - \frac{n-1}{2N} + \mathcal{O}\left(\frac{1}{N^2}\right)$

$Pr(\text{coalescent event}) = 1 - Pr(\text{all diff})$

$\approx \frac{1}{2N} + \frac{2}{2N} + \cdots + \frac{n-1}{2N} = \frac{1}{2N}\left(1 + 2 + \cdots + n-1\right)$

$= \frac{1}{2N} \frac{n(n-1)}{2}$

$= \frac{n(n-1)}{4N}$

- Each generation back, $Pr(\text{coalesce}) = \dfrac{n(n-1)}{4N}$

So time to coalescence $\sim \text{Geometric}\left(\dfrac{n(n-1)}{4N}\right)$

So $E(T_n) = \dfrac{4N}{n(n-1)}$     (Property of geometric distribution)

So if $T_c = \sum\limits_{i=2}^{n} i\, T_i$

$$E(T_c) = \sum_{i=2}^{n} i\, E[T_i]$$

$$= \sum_{i=2}^{n} i\, \frac{4N}{i(i-1)}$$

$$= 4N \sum_{i=2}^{n} \frac{1}{i-1}$$

$$= 4N \sum_{i=1}^{n-1} \frac{1}{i}$$

$$E[S_n] = u\, E[T_c] = 4Nu \sum_{i=1}^{n-1} \frac{1}{i} = \theta \sum_{i=1}^{n-1} \frac{1}{i}$$

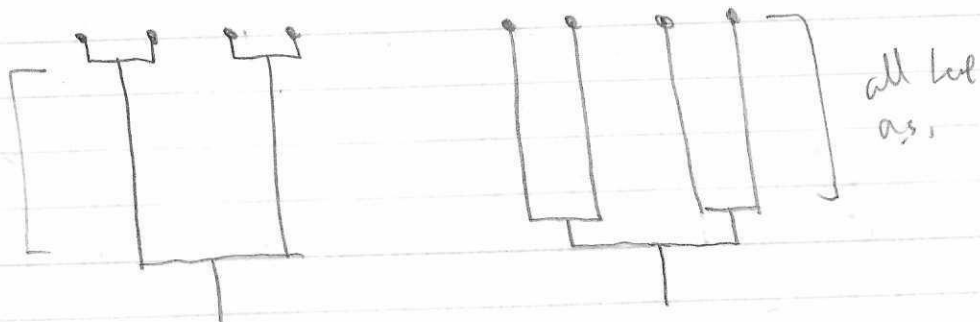$$\Rightarrow \hat{\theta} = \frac{S_n}{\sum\limits_{i=1}^{n-1} \frac{1}{i}}$$

---

- Tajima's D: Compare observed heterozygosity to that given by $\hat{\theta}$
  - tests for neutral model
  - If observed heterozygosity = that from $\hat{\theta}$, then no evidence against neutral model

$$\left(\frac{1}{4}\right)^2 + \left(\frac{3}{4}\right)^2 = \frac{1}{16} + \frac{9}{16} = \frac{5}{8}$$

Example: ①       ②

all here
as.

all here
as.

Put 5 random mutations

Case ① : Mutation is either in left or right
$p = \frac{1}{2}$ for all mutations     $H = 2\left(\frac{1}{2}\right)\left(1-\frac{1}{2}\right) = \frac{1}{2}$
↳ allele freq

Case ② : mutation is in one branch
$p = \frac{1}{4}$      heterozygosity $= 2\left(\frac{1}{4}\right)\left(1-\frac{1}{4}\right) = \frac{3}{8}$

$H = \Pr(\text{differ by state} \mid \text{drawn w/ replacement})$

• We will show that $E\left[\sum\limits_{i=1}^{\infty} 2p_i(1-p_i)\right] = \theta$
    ↳ later                ↳ under neutral loci

let $\pi = \sum\limits_{i=1}^{\infty} 2p_i(1-p_i)$
     ↳ unknown

$\hat{\pi} = \frac{n}{n-1} \sum\limits_{i=1}^{S_n} 2\hat{p}_i(1-\hat{p}_i)$
     ↳ sum over segregating sites.

Tajima $D$: $D_T = \dfrac{\hat{\pi} - \hat{\theta}}{C}$

$C$ chosen s.t. $\dfrac{\hat{\pi} - \hat{\theta}}{C} \sim$ Normal

$\hat{\pi}$ is estimate of $\theta$ from heterozygosity
$\hat{\theta}$ is estimate from coalescent

- $D_T > 0 \Rightarrow$ more heterozygosity than expected
  $D_T < 0 \Rightarrow$ less heterozygosity than expected

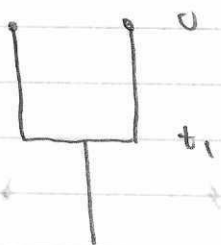- A proof that $E\left[\sum\limits_{i=1}^{n} 2p_i(1-p_i)\right] = \theta$:

$\pi$ is expected # of nucleotide differences between randomly selected pair of alleles.

$z_i :=$ pair of alleles differ at nucleotide $i$

$E[z_i | p_i] = 2p_i(1-p_i)$

$\pi = \sum\limits_{i=1}^{n} z_i = $ # differences b/t pair of alleles

$E[\pi | p] = E\left[\sum\limits_{i=1}^{n} z_i | p_i\right] = \sum\limits_{i=1}^{s} 2p_i(1-p_i)$



$E[t_i] = \dfrac{4N}{2(2-1)} = 2N$

# differences $= 2 t_i u$

$\Rightarrow E[\text{# differences}] = 4Nu = \theta$ //

- A coalescent derivation of

$$H = \frac{4Nu}{1 + 4Nu}$$

- **Previous derivation**: found $\Delta H$ in terms of mutation and drift, set $= 0$

① $\phi' = (1-u)^2 \left( \frac{1}{N} + \left(1 - \frac{1}{N}\right) \phi \right) = Pr(\text{sue state} \mid \text{diff origin})$

    Next        No mutation   Sue   diff   prev
    gen                                     gen

                            Sue gen - No mutation

② $H' = 1 - \phi'$

③ use $u^2 \approx 0$, $\frac{u}{N} \approx 0$

④ find $\Delta H = H' - H$

⑤ set $\Delta H = 0$, solve for $H$

Two alleles differ by state $\iff$ mutation after common ance



$Pr(\text{Coalesce in 1 generation}) = \frac{1}{2N}$

$Pr(\text{Mutation in 1 generation}) = 1 - (1-u)^2$

Note: $1 - (1-u)^2 = 1 - (1 - 2u + u^2)$

$$= 2u - u^2$$

$$\approx 2u \qquad \text{since } u^2 \text{ is small}$$

Heterozygous if mutation occurs first

Pr( mutation | mutation or coalesce)

$$= \frac{Pr[\text{mutation and (mutation or coalesce)}]}{Pr(\text{mutation or coalesce})}$$

$$= \frac{Pr(\text{mutation})}{Pr(\text{mutation or coalesce})}$$

$$= \frac{2u}{2u + \frac{1}{2N}}$$

$$= \frac{4Nu}{1 + 4Nu} \quad //$$