

Week 5 Phylogenomics

Friday, 18 February 2022 8:02 AM

African ancestry (American SW)

20x coverage LCT ; 16 4988235 SNP

70% G (ref) ; 30% A (alternate)

① Compute genotype likelihoods. ←

② Which is the most likely genotype at this locus?

TRUTH

GG

AA

AG

$$\begin{aligned} \text{GENOTYPE LIKELIHOOD} \\ {}^n C_k \epsilon^k (1-\epsilon)^{n-k} \\ = P(D|GG) \end{aligned}$$

$${}^n C_{n-k} \epsilon^{n-k} (1-\epsilon)^k = P(D|AA)$$

$${}^n C_k \left(\frac{1}{2^n} \right) = P(D|AG)$$

$n = \text{total \# of reads}$
(20)
 $k = 14$

$$\textcircled{1} P(D|GG) = {}^{20}C_6 \left(\frac{6}{20} \right)^6 \left(\frac{14}{20} \right)^{14} = 0.191639 \leftarrow$$

$$\textcircled{2} P(D|AA) = {}^{20}C_{14} \left(\frac{14}{20} \right)^{14} \left(\frac{6}{20} \right)^6 = 0.002181$$

$$\textcircled{3} P(D|AG) = {}^n C_k \left(\frac{1}{2^n} \right) = {}^{20}C_{14} \left(\frac{1}{2^{20}} \right) = 0.0369$$

$$\begin{aligned} P(GG|D) &\propto P(D|GG) \times P(GG) \\ &\propto 0.191639 \times 0.705 = 0.135 \leftarrow \end{aligned}$$

$$\begin{aligned} P(AA|D) &\propto P(D|AA) \times P(AA) \\ &\propto 0.002181 \times 0.049 = 0.000106 \end{aligned}$$

$$\begin{aligned} P(AG|D) &\propto P(D|AG) \times P(AG) \\ &\propto 0.0369 \times 0.246 = 0.0090774 \end{aligned}$$

2 0.0369 x 10^-4

II Mutations

PURINES - A, G
PYRIMIDINES - C, U, T

SUBSTITUTION (SNPs)

↳ transitions (PUR ↔ PUR, PYR ↔ PYR)
↳ transversions (PUR ↔ PYR)

transitions more likely

synonymous

vs

non-synonymous

synonymous more likely

II INDELS

Indels are less likely than subs.

MSA

sequence similarity

LOTS - PARSIMONY

LESS

some degree of sequence similarity

YES

DISTANCE - BASED
(NJ, UPGMA)

NO

LIKELIHOOD - BASED
(Bayesian)

PARSIMONY - MINIMUM

EVOLUTION

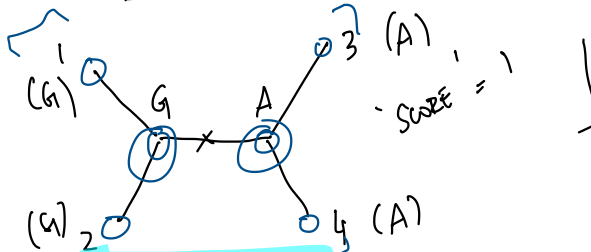
	1	2	3	4
INDIV 1	A	G	T	A
2	C	G	T	G
3	T	A	T	G
4	G	A	T	G

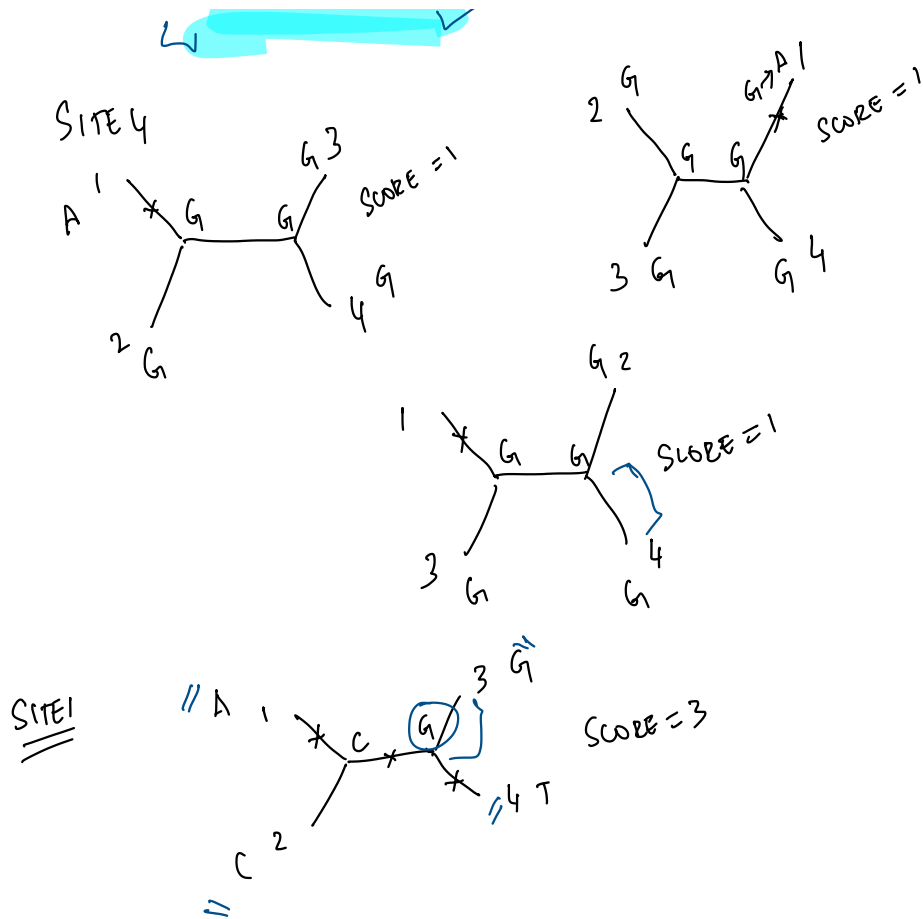
SITE 3 - invariant

SITES 1, 2, 4

- SEGREGATING

SITE 2



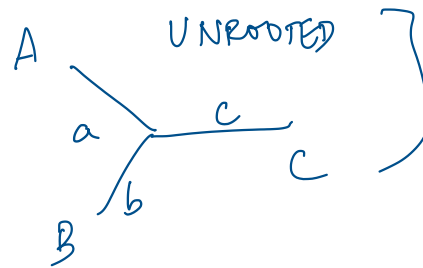


	1	2	3	4
1	0	2	3	3
2		0	2	2
3			0	1
4				0

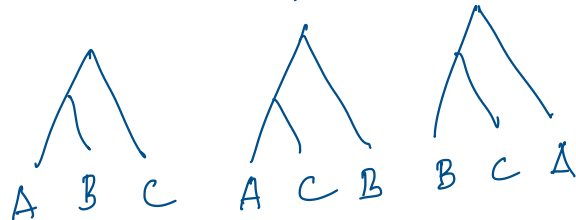
FITCH & MARGOLISAN

- ① MOLECULAR CLOCK → mutations accumulate at a constant rate along a branch
- ② INFINITE SITES MODEL → a mutation occurs just once at a site

	A	B	C
A	-	22	39
B		-	41
C			-



ROOTED



NEIGHBOR-JOINING

① Always start tree with "STAR" TOPOLOGY with 3 branches

②

$$a + b = 22 \quad \text{--- ①}$$

$$a + c = 39 \quad \text{--- ②}$$

$$b + c = 41 \quad \text{--- ③}$$

Subtract ② from ①

$$\cancel{a} + b - \cancel{a} - c = 22 - 39$$

$$b - c = -17 \quad \text{--- ④}$$

add ③ and ④

$$\cancel{b} + \cancel{c} + b - \cancel{c} = 41 - 17$$

$$2b = 24$$

$$b = 12$$

\Rightarrow put this in (3)

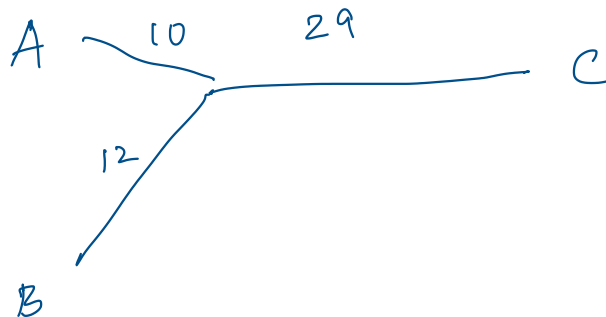
$$12 + c = 41$$

$$\Rightarrow c = 29$$

put $b = 12$ in (1)

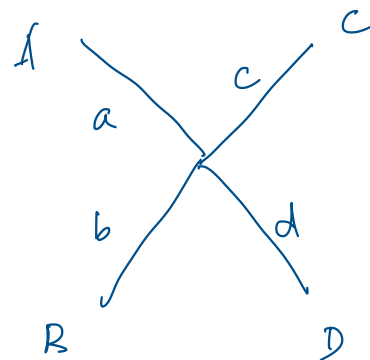
$$\Rightarrow a + 12 = 22$$

$$\Rightarrow a = 10$$



eg.

	A	B	C	D
A		2	3	3 //
B			2	2
C				1
D				



$$C \begin{matrix} 0.5 \\ c \end{matrix} \quad x = 2 \quad (1, 2)$$

' NJ ' STEP

A		2	$\frac{AC+AD}{2} = \frac{2}{2}$
B			$\frac{BC+BD}{2} = \frac{2+2}{2} = 2$
(C,D)			

$$a + b = 2 \quad \text{--- (1)}$$

$$a + y = 3 \quad \text{--- (2)}$$

$$b + y = 2 \quad \text{--- (3)}$$

$$\text{(1) - (2):}$$

$$a + b - a - y = 2 - 3$$

$$b - y = -1$$

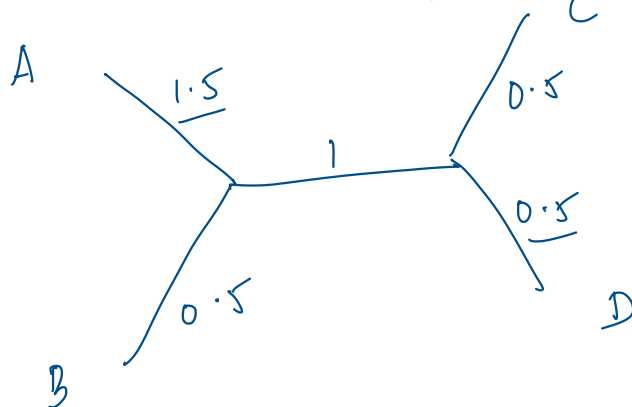
$$b + y = 2$$

$$\hline 2b = 1$$

$$\Rightarrow b = 0.5$$

$$a = 1.5$$

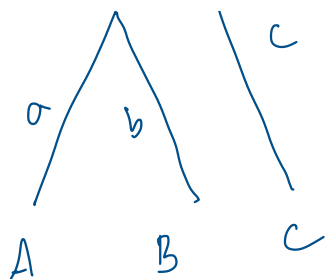
$$y = 1.5$$



f

$$a = b$$

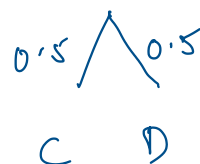
UPGMA



$$a + x = c$$

ULTRAMETRIC TREES

	A	B	C	D
A		2	3	3
B			2	2
C				1
D				

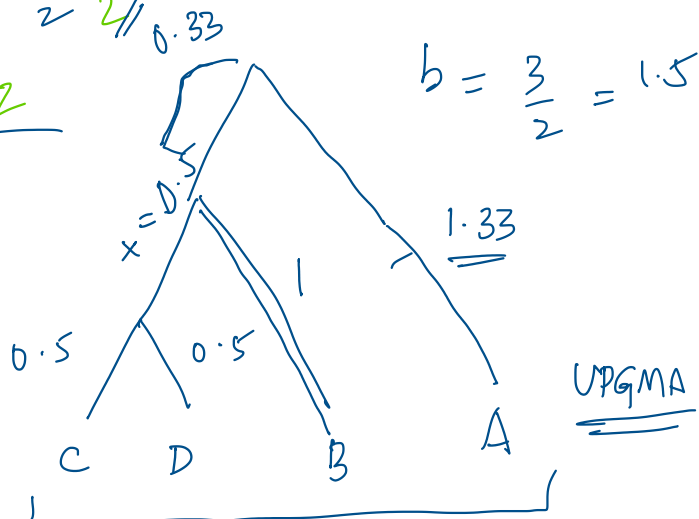


$$c + d = 1$$

$$c = d$$

$$\Rightarrow c = d = 0.5$$

	(C,D)	A	B
(C,D)		$\frac{AC+AD}{2} = 3$	$\frac{BC+BD}{2} = 2$
A			2
B			



$$b = \frac{3}{2} = 1.5$$

UPGMA

$$x + 0.5 = b$$

	A	(B,C,D)
(B,C,D)	$\frac{AB+AC+AD}{3}$	-

$$\frac{AB+AC+AD}{3}$$

$$= 2 + 3 + 3$$

A		

$$= \frac{8}{3} = \underline{\underline{2.66}}$$

Observed data = D

$H_1 = \text{TOPOLOGY 1}$, $H_2 = \text{TOPOLOGY 2}$

$$\textcircled{1} - P(H_1 | D) = \frac{P(D | H_1) \times P(H_1)}{P(D)} \quad \text{Bayes theorem}$$

$$\textcircled{2} - P(H_2 | D) = \frac{P(D | H_2) \times P(H_2)}{P(D)}$$

$$\textcircled{1} \div \textcircled{2}$$

$$\Rightarrow \left[\frac{P(H_1 | D)}{P(H_2 | D)} \right] = \left[\frac{P(D | H_1) \times P(H_1)}{P(D | H_2) \times P(H_2)} \right]$$

$$\text{Bias} = P(H) = p \quad \Rightarrow \quad P(T) = 1 - p$$

D = HHTTHTHTTT

$$L = P(D | p) = \left[p \times p \times (1-p) \times (1-p) \times (1-p)^3 \right]$$

$$x p \times (1-p) \times p \times p \times \dots$$

$$L = p^5 (1-p)^6 \Rightarrow$$

$$5p^4 (1-p)^6 - 6p^5 (1-p)^4 = 0$$

$$p^4 (1-p)^5 [5(1-p) - 6p] = 0$$

$$\Rightarrow \hat{p} = 5/11$$

SOLUTION 1

$$\frac{dL}{dp}$$

natural log on both sides

SOLUTION 2

$$\ln L = 5 \ln p + 6 \ln (1-p)$$

$$\frac{d(\ln L)}{dp} = \frac{5}{p} - \frac{6}{1-p} = 0$$

$$\Rightarrow \hat{p} = 5/11$$

ASSUME :

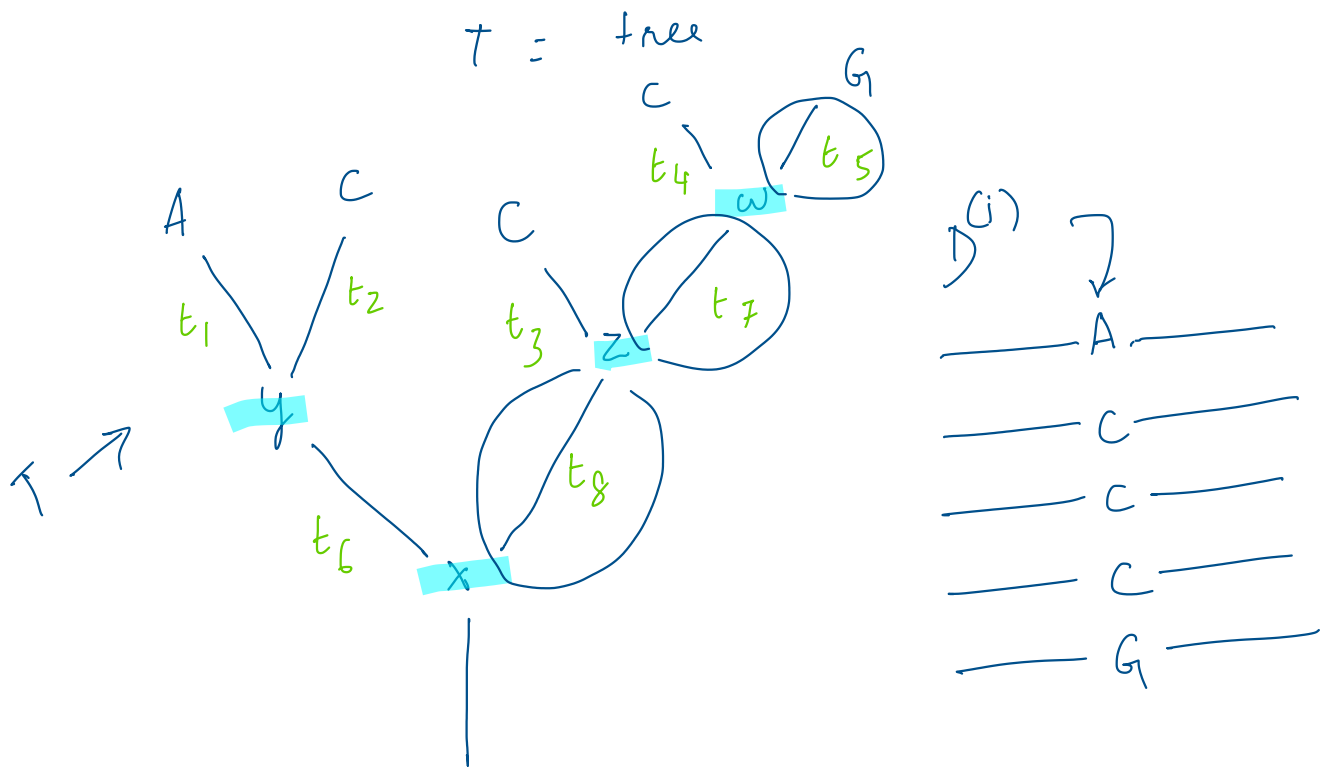
(1) evolution on different sites
are independent of each other

② evolution on different lineages
are independent of each other.]

T = tree (topology) ; D = genetic data

$$\Rightarrow L = P(D|T) = \prod_{i=1}^m P(D^{(i)}|T) \leftarrow$$

$D^{(i)}$ = data at site i



$$P(D^{(i)}|T) = \sum_x \sum_y \sum_z \sum_w P(A, C, C, C, G, x, y, z, w | T)$$

$$x, y, z, w = \{A, C, G, T\}$$

$$= \underbrace{p(x)} \times \overbrace{p(y|x, t_6)} \times A(A|y, t_1) \\ \times p(C|y, t_2) \times p(z|x, t_8) \\ \times p(C|z, t_3) \times p(w|z, t_7) \\ \times p(C|w, t_4) \times p(G|w, t_5)$$