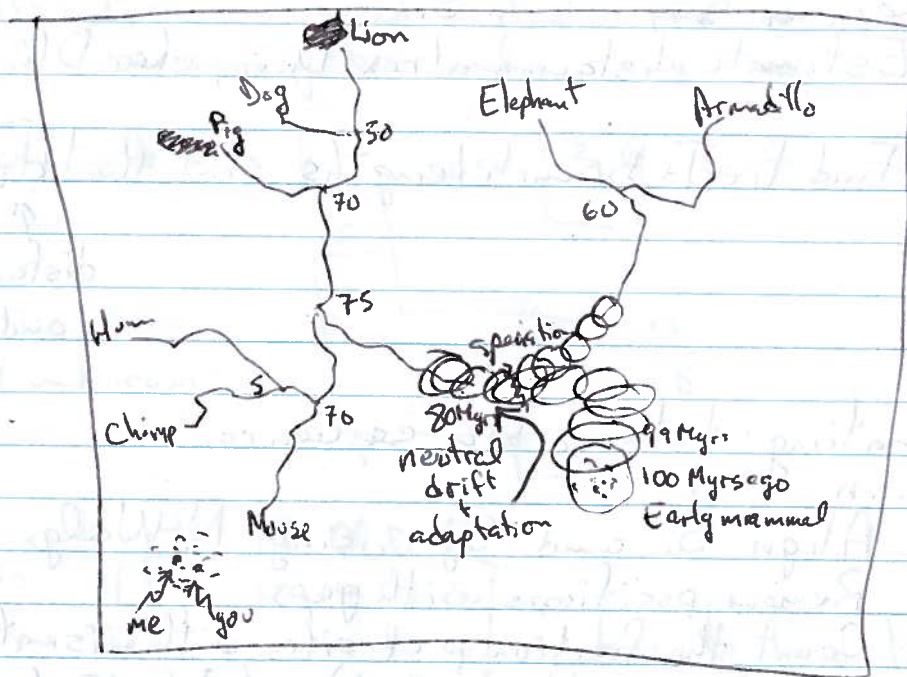


Sequence Evolution + Phylogenetics

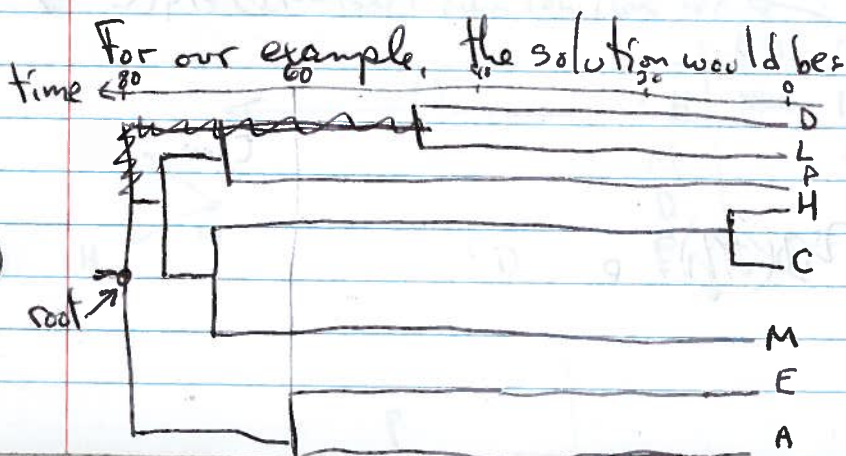


Phylogenetic Inference Problem (Biological version)

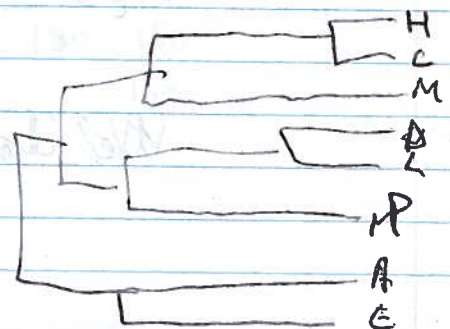
Given: DNA sequences from multiple species

Find: Evolutionary tree that describes their evolution

→ Topology, Dating of speciation events



Note: This is the same topology as



Distance-based Phylo. Inference methods

Idea: Given seq. S_1, \dots, S_n

① Estimate distance matrix $D_{n \times n}$, where $D(i, j) = \text{distance}$ btw S_i, S_j

② Find tree + Branch lengths such that $d_T(i, j) \approx D(i, j)$
distance btw nodes i and j in tree T

③ Estimating distance btw sequences

For $i, j = 1 \dots n$

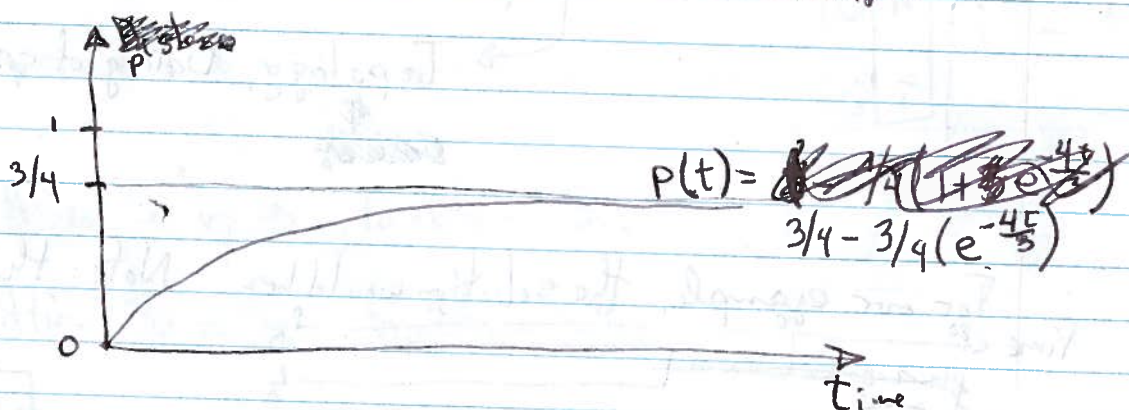
1. Align S_i and S_j using N-W algo
2. Remove positions with gaps
3. Count the fraction p of sites with mismatches
4. $D(i, j) = -\frac{3}{4} \left(\log \left(1 - \frac{4}{3} p \right) \right)$ (Jukes-Cantor model)

Example $S_i: A C - - T G A C T G$
 $S_j: A G T T T A - C A G$ $\Rightarrow p = 3/7$

$$D(i, j) = 0.63$$

Note:

If the two seqs.
are very long



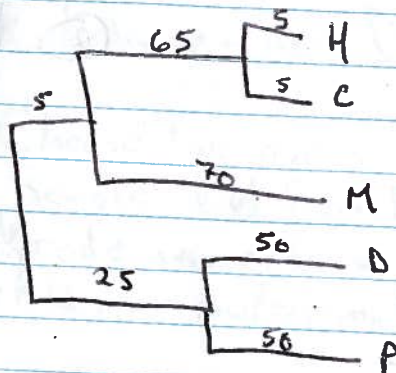
We choose $D(i, j) = p$

Under Kimura 2-parameter model with rates μ transition & ν transversion

Notes: More advanced distance measures consider separately transitions and transversions

Example:

True tree
(unknown)



$$D = \begin{matrix} & \begin{matrix} H & C & M & D & P \end{matrix} \\ \begin{matrix} H \\ C \\ M \\ D \\ P \end{matrix} & \begin{bmatrix} - & 10 & 140 & 150 & 150 \\ & - & 140 & 150 & 150 \\ & & - & 150 & 150 \\ & & & - & 100 \\ & & & & - \end{bmatrix} \end{matrix}$$

How to find tree?

Incorrect

See appendix for corrected version and example

UPGMA

- ① Choose two closest nodes according to UPGMA rule
- ② Merge two nodes into a new node w
- ③ Remove row/col u, v from matrix
- ④ Insert new row/col for w : $D(w, i) = D(u, i) - \frac{1}{2} D(u, v)$

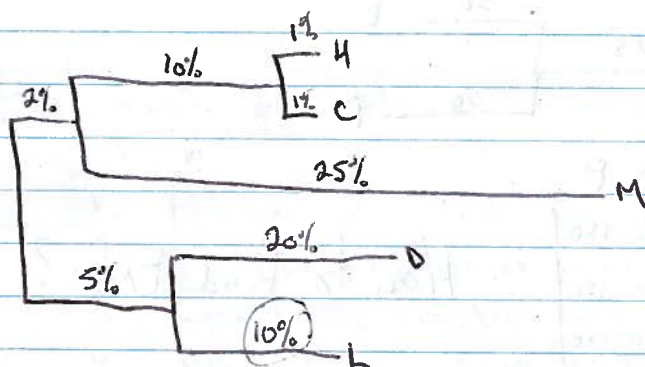
				CH	M	D	P
				CH	135	145	145
				M		150	150
				D			100
				P			

Running time: ?

Problem: Only works if

- ① Distances are estimated perfectly accurately
- ② Mutation rate is constant

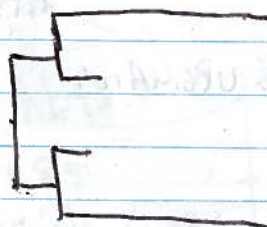
More realistic tree:



Mutation rate varies

⇒ Distances btw seq cannot be calculated in M₉₅ but instead in Expected substitution per site

→ This means that ~~40% of positions will have been mutated~~ the expected # of subst. per site is 0.1



⇒ UPGMA will produce the wrong tree