

# CISC 636 Computational Biology & Bioinformatics (Fall, 2016)

## Phylogenetic Trees (III)

### Probabilistic methods

Bayes rule revisited.

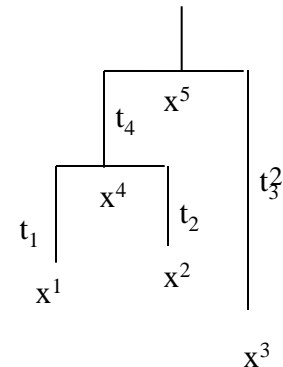
$$P(\text{model}|\text{data}) = P(\text{data}|\text{model}) P(\text{model})/P(\text{data})$$

- model includes
  - our evolution theory,
  - a specific phylogenetic tree (topology and edge lengths), and
  - assignment of sequences to the tree leaves.
- Data: a set of sequences that are used to infer phylogeny.

Let  $P(x|y, t)$  be the probability in that sequence  $x$  is evolved from an ancestral sequence  $y$  over an edge of length  $t$ .

$$P(x^1, x^2, x^3, x^4, x^5|T, t.) \\ = P(x^1|x^4, t_1) P(x^2|x^4, t_2) P(x^3|x^5, t_3) P(x^4|x^5, t_4) P(x^5)$$

A shorthand notation  $P(x'|T, t.)$  is used where  $x'$  stands for a set of sequences, and  $t.$  for edge lengths of the tree  $T$ .



In general, if we know

- $P(x|y, t)$  for any  $x$ ,  $y$ , and  $t$
- A tree  $T$ , and assignment of sequences to tree nodes,

Then we can compute the likelihood for observing the sequences as they are assigned onto the tree leaves.

Q: Given  $n$ , the number of leaves, there are  $(2n-3)!!$  different trees (plus many different ways to assign length to tree edges), which tree can best interpret the data?

A: The tree that gives the maximum likelihood (ML).

In practice, to implement the ML method, two issues we need to address

1. A model of evolution, which gives the conditional probabilities  $P(x|y, t)$
2. Method to find the maximum likelihood. For this, any optimization method may be utilized, such as
  - Descent gradient
  - Simulated annealing
  - Genetic algorithm

## Models of evolution

Independence assumption: mutations occur independently at different positions.

Therefore,

$$P(x|y, t) = \prod_u P(x_u|y_u, t),$$

where  $P(x_u|y_u, t)$  is the probability that residue  $x_u$  in sequence  $x$  mutates to residue  $y_u$  in sequence  $y$  over time  $t$ .

multiplicative assumption:

$$P(b|a, t + \Delta t) = \sum_c P(c|a, t) \cdot P(b|c, \Delta t).$$

For DNA sequences, probabilities for all possible mutations among four nucleotides during a given time period  $t$  form a 4 by 4 matrix

$$S(t) = \begin{pmatrix} P(A|A, t) & P(A|C, t) & P(A|G, t) & P(A|T, t) \\ P(C|A, t) & P(C|C, t) & P(C|G, t) & P(C|T, t) \\ P(G|A, t) & P(G|C, t) & P(G|G, t) & P(G|T, t) \\ P(T|A, t) & P(T|C, t) & P(T|G, t) & P(T|T, t) \end{pmatrix}$$

And, we have multiplicative property for these matrices

$$S(t) \cdot S(\Delta t) = S(t + \Delta t).$$

For each  $P(b|a, t)$  in the substitution matrix  $S(t)$ , it is reasonable to assume that no mutation can occur at zero time interval:

- $P(a|a, 0) = 1$
- $P(b|a, 0) = 0$ .

That is,

$$S(0) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Further, let's assume that  $P(b|a, \Delta t)$  over a infinitesimal interval  $\Delta t$  is proportional to  $\Delta t$  by a constant  $r$ , called mutation rate:

$$P(b|a, \Delta t) = r \Delta t.$$

**Jukes-Cantor model** for DNA sequences assumes that all nucleotide mutations have the same rate. That is,

$$S(\Delta t) = \begin{pmatrix} 1-3r & r & r & r \\ r & 1-3r & r & r \\ r & r & 1-3r & r \\ r & r & r & 1-3r \end{pmatrix} \quad \Delta t \equiv I + R \Delta t$$

Therefore,

$$S(t + \Delta t) = S(t) \cdot S(\Delta t) = S(t) \cdot [I + R\Delta t]$$

$$[S(t + \Delta t) - S(t)] / \Delta t = S(t) \cdot R$$

$$S'(t) = S(t) \cdot R \quad \text{when } \Delta t \rightarrow 0.$$

Solving the differential equations, we have

- $P(a|a, t) = 1/4 (1 + 3e^{-4rt})$  for  $a \in \{A, C, G, T\}$
- $P(b|a, t) = 1/4 (1 - e^{-4rt})$  for  $a \neq b, a \text{ and } b \in \{A, C, G, T\}$

In this model, when  $t = \infty$ ,  $P(b|a, \infty) = 1/4$ . That is, the nucleotide equilibrium frequencies are all equal.

**Kimura model** for DNA sequences assumes different rates for transitions ( $A \leftrightarrow G$  and  $C \leftrightarrow T$ ) and transversions ( $A \leftrightarrow T$ ,  $G \leftrightarrow T$ ,  $A \leftrightarrow C$ , and  $C \leftrightarrow G$ ).

That is,

$$S(\Delta t) = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{pmatrix} 1-2r-s & r & s & r \\ r & 1-2r-s & r & s \\ s & r & 1-2r-s & r \\ r & s & r & 1-2r-s \end{pmatrix} \end{matrix} \Delta t \equiv I + R \Delta t$$

Similar models are proposed for mutations among amino acids.

If we were able to quantify the “time” as how many number mutations have occurred, the substitute matrices in those models would correspond to PAM matrices at respective times.

# Maximum Likelihood

- The case of two sequences aligned with no gaps

$$P(x^1, x^2 | T, t_1, t_2) = \prod_{u=1 \text{ to } N} P(x_u^1, x_u^2 | T, t_1, t_2)$$

- Let  $x^1 = \text{CCGGCCGCGCG}$   
 $x^2 = \text{CGGGCCGCCCG}$

$$\begin{aligned} P(C, C | T, t_1, t_2) &= P(C|A, t_2) P(C|A, t_1) P(A) \\ &\quad + P(C|C, t_2) P(C|C, t_1) P(C) \\ &\quad + P(C|G, t_2) P(C|G, t_1) P(G) \\ &\quad + P(C|T, t_2) P(C|T, t_1) P(T) \\ &= \frac{1}{4} [3 \times \frac{1}{4} (1 - e^{-4rt_1}) \times \frac{1}{4} (1 - e^{-4rt_2}) + \\ &\quad \frac{1}{4} (1 + 3e^{-4rt_1}) \times \frac{1}{4} (1 + 3e^{-4rt_2})] \\ &= \frac{1}{16} (1 + 3e^{-4r(t_1 + t_2)}). \end{aligned}$$

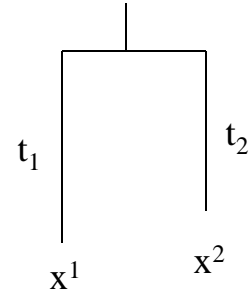
$$P(G, G | T, t_1, t_2) = \frac{1}{16} (1 + 3e^{-4r(t_1 + t_2)}).$$

$$P(C, G | T, t_1, t_2) = \frac{1}{16} (1 - e^{-4r(t_1 + t_2)}).$$

Therefore, for an alignment that has  $n_1$  identical sites and  $n_2$  mutational sites, we have

$$P(x^1, x^2 | T, t_1, t_2) = \frac{1}{16^{n_1 + n_2}} \times (1 + 3e^{-4r(t_1 + t_2)})^{n_1} \times (1 - e^{-4r(t_1 + t_2)})^{n_2},$$

which is a function of edge lengths in tree  $T$ .





In general, the probability can be computed by working up the tree from the leaves using post-order traversal. This is done by Felsenstein's algorithm (1981).

Once the probability is available, optimizing the assignments of edge lengths  $t$  in the tree amounts to

$$\left. \frac{\partial P}{\partial t} \right|_{t_m} = 0$$

Where  $t_m$  is the tree length assignment that maximize the likelihood.

## How to optimize tree topology?

- Discrete structure, therefore cannot take derivatives.

## Basic strategy for searching the tree space

- A tree generation algorithm that can generate trees
- Assess the likelihood
  - Accept
  - Reject

## A genetic algorithm implementation [Matsuda 1998]

# Genetic algorithm

## Input

- P, the population,
- r: the fraction of population to be replaced,
- f, a fitness,
- ft, the fitness\_threshold,
- m: the rate for mutation.

Initialize population (randomly)

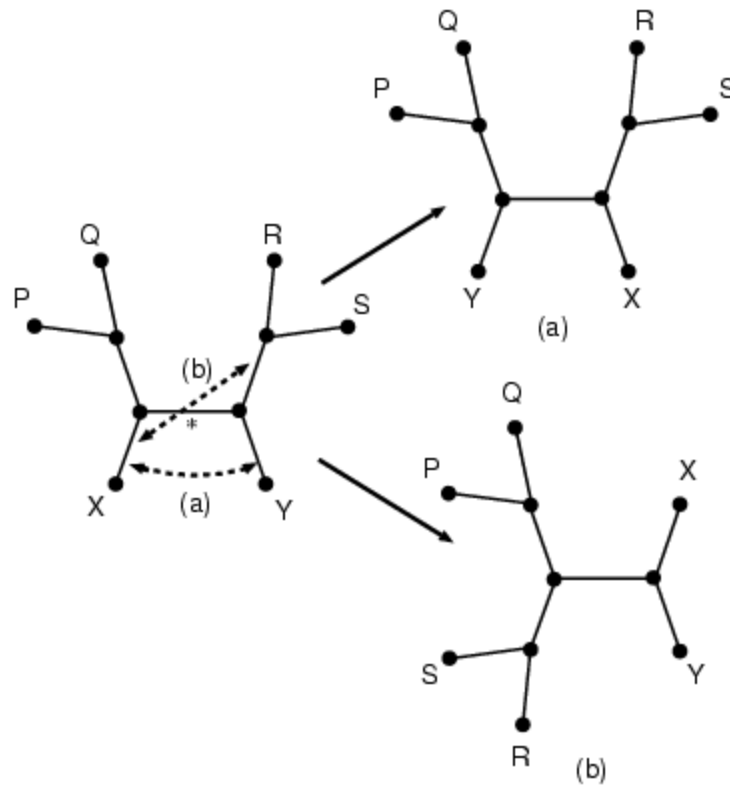
Evaluate: for each h in P, compute Fitness(h)

While [ $\text{Max}_h f(h)$ ] < ft

do

1. Select
2. Crossover
3. Mutate
4. Update P with the new generation Ps
5. Evaluate: f(h) for all h  $\in$  P

Return the h in P that has the best fitness



## Branch exchange in a phylogenetic tree

# Key components for implementing genetic algorithms

- Representing hypotheses (which are the trees here)

- New Hampshire format

(A,(F,(C,(B,D))))

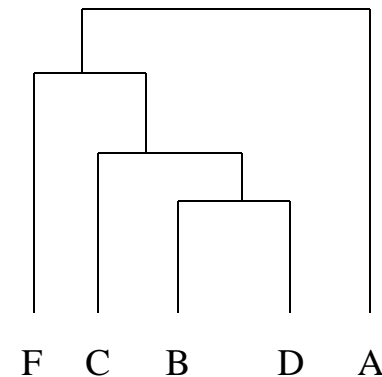
- Genetic operators

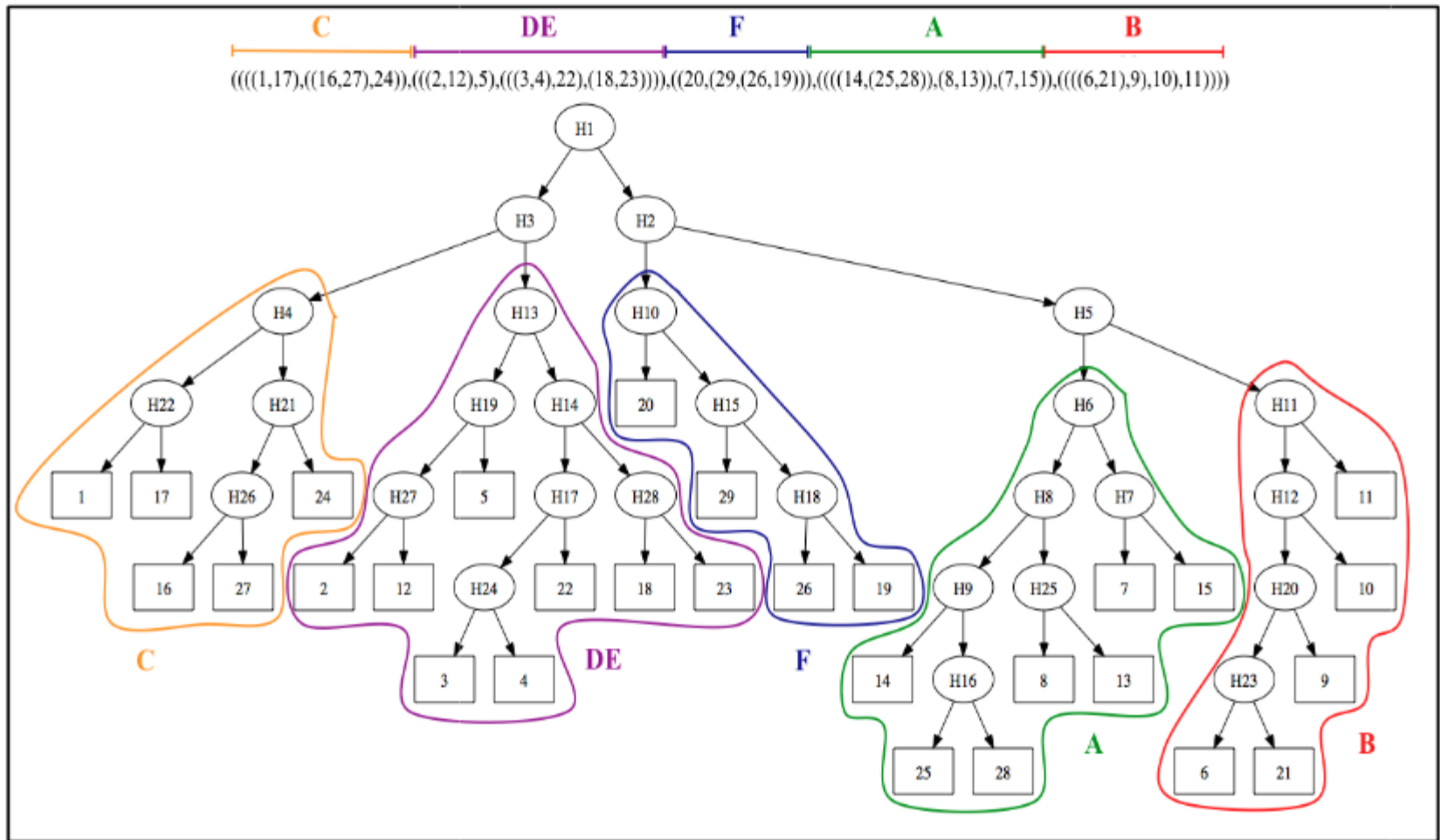
- Crossover:

- Single
    - Two point
    - uniform

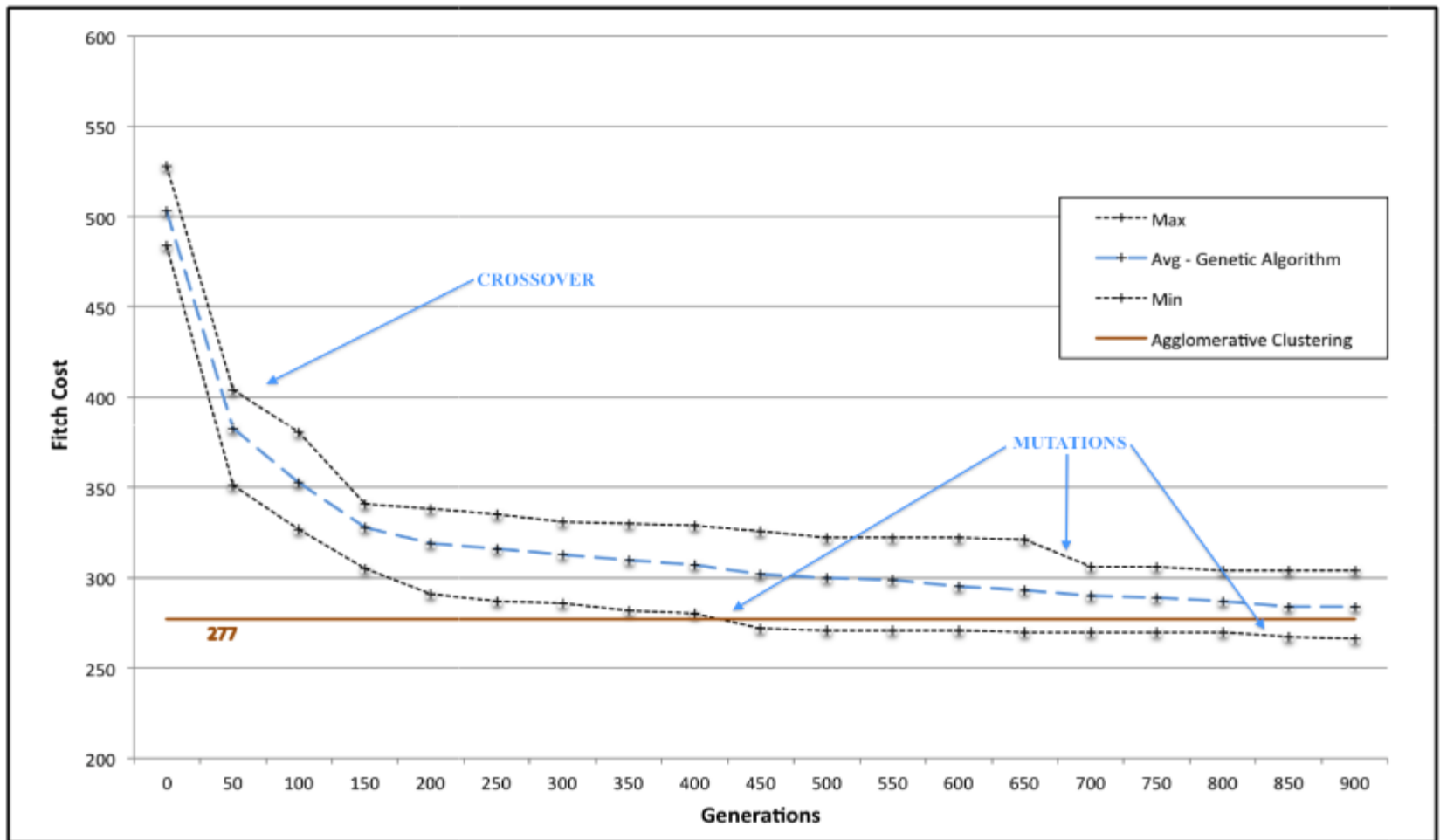
- Mutation: point

- Fitness function: we use the likelihood computed for each tree using





Blanchette, O;Keefe & Benuskova, 2012



Blanchette, O;Keefe & Benuskova, 2012

# More advanced topics in phylogenetic analysis

- Different heuristics for sampling the tree space
  - Monte Carlo
  - ...
- More realistic evolutionary models
  - With gaps
  - Non-uniform: different rates at different sites
  - ...
- Using different data sets and reconciliation
  - Sequences
  - Gene positions -> genome rearrangement [Nadeau & Taylor 1984, PNAS 81:814-818, Pavzner, Sankoff, ...]
  - ...



# More advanced topics in phylogenetic analysis

- Different heuristics for sampling the tree space
  - Monte Carlo
  - ...
- More realistic evolutionary models
  - With gaps
  - Non-uniform: different rates at different sites
  - ...
- Using different data sets and reconciliation
  - Sequences
  - Gene positions [Nadeau & Taylor 1984, PNAS 81:814-818, Pavzner, Sankoff, ...]
  - ...