# PARSIMONY

Consider the sequences

$S_1$: A A T G       L

$S_2$: A G T G       S

$S_3$: A G T C       S
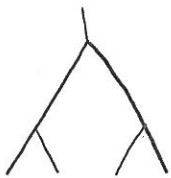
$S_4$: A A T C       L

$\underbrace{\qquad\qquad}_{\text{DNA}}$ $\underbrace{\qquad}_{\text{morphological}}$

Which of the 4 (out of $(2n-3)!! = 5!! = 15$) rooted trees, do we prefer? and why?
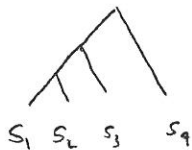


$S_1$ $S_2$ $S_3$ $S_4$

PS = 5



$S_1$ $S_4$ $S_2$ $S_3$

PS = ?



$S_1$ $S_3$ $S_4$ $S_2$

PS = 6



$S_1$ $S_2$ $S_3$ $S_4$

These numbers are called the PARSIMONY SCORES of $T$, $ps(T)$

Subtle point: Parsimony scores are really counts for unrooted trees.

For example, trees 1 and 4 are rooted versions of $S_4 \diagdown\!\!\diagup \begin{smallmatrix}S_1 & S_3\end{smallmatrix} S_4$ and $ps(T_1) = ps(T_4)$

The three other rooted versions of $\begin{smallmatrix}S_1\\S_2\end{smallmatrix}\!\diagup\!\diagdown$ would also have the same ps.

Rooting just provides a way (= algorithm) for computing ps.

These numbers are known as the PARSIMONY SCORES of T

Theorem: $ps(T^p) = ps(T)$

    i.e. the parsimony score of a rooted tree

    = " " " " its unrooted version

$\Rightarrow$ PARSIMONY SCORES compare unrooted trees only

The underlying principle: The best tree (or trees in the case of ties) is the one requiring the minimal amount of changes (substitutions/ state changes)

Defn: A column in aligned sequences is called a CHARACTER (usually denoted by $X$) and a character can take on any of

    s STATES

        Eg: DNA character $s=4$      morphological char $s=2$

Data:    $X_1 \; X_2 \; X_3 \cdots$

How should we compute parsimony scores?
      v
    (or have a computer)

Do alg development + Eg simultaneously

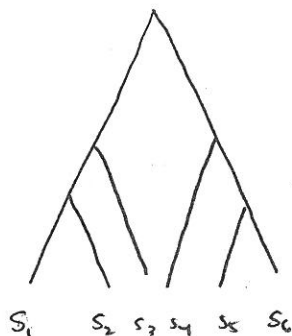Fix a binary tree $T$

## Alg.

Sum over all characters $X_i$, $i=1,\ldots,m$

find $ps_{X_i}(T)$ by

0. If unrooted, arbitrarily root tree
1. place $X_i$ at leaves of tree

2. Starting at the leaves,

label each ancestral node by

- the union of the two sets below if they are disjoint
  in which case the parsimony count is augmented by 1
- the intersection of the two state sets below if they share elts

  (and do **not** augment the parsimony count)

3. continue up tree until reach root    Return $ps_{X_i}(T)$.

```
        A A
        A A
        T A
        T G
        G G
        C C

S₁   S₂ S₃ S₄ S₅ S₆
```

## Comments:

1. This is called the **FITCH - HARTIGAN ALGORITHM**

and it **does** compute the $ps(T)$    Thm: FH count on $T = ps(T)$

2. Computer scientists would call this the **SMALL - PARSIMONY PROBLEM**

( fix **one** single tree and compute its $ps$ )

How complex is this?    If $T$ is an $X$-tree with $n = |X|$ leaves

Roughly,    $n$ internal nodes to visit
           each node might have any of $s$ states    and
           must loop over $m$ characters

$$\Rightarrow \approx nsm \equiv \Theta(nsm) \quad \text{"big-}\Theta\text{"}$$

$$\Rightarrow \text{time is} \quad \leq \quad (\text{constant})(nsm)$$

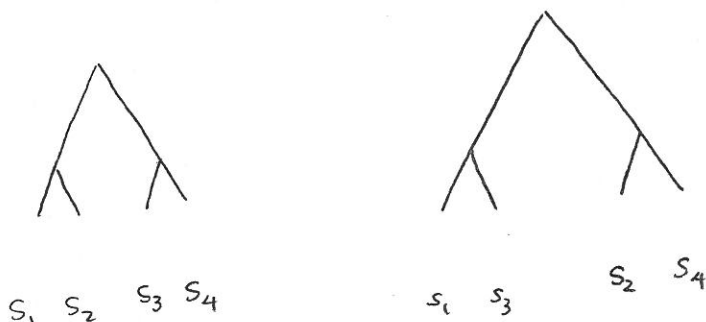3. The LARGE- PARSIMONY PROBLEM is the one of interest

Search over all $(2n-5)!!$ unrooted trees

    compute the parsimony score for each

    return those with best score

    UAF grad Ron Graham + Foulds proved this is NP-hard

4. Ways to speed this up



$S_1$   $S_2$   $S_3$ $S_4$        $S_1$   $S_3$     $S_2$   $S_4$

$S_1$ : A A A A A A
$S_2$ : A G C C A A
$S_3$ : A A T C C T
$S_4$ : A A G T C T

Remove the constant sites or, more generally, the non-informative sites

Defn: If $X$ is a character, with $S$ states, then a PARSIMONY – INFORMATIVE

CHARACTER in which 2 states appear at least twice

A A A       S S
A A C       L S
T A G       L L
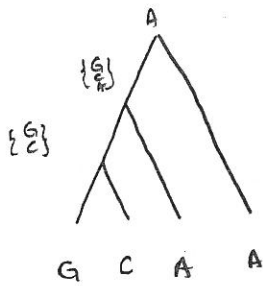T T T       L L

To speed up parsimony, preprocess to remove non-informative sites

recognize that all informative patterns, for instance,
X
X
Y
Y
contribute equally to parsimony count

i.e. particular states don't matter, only pattern of states.
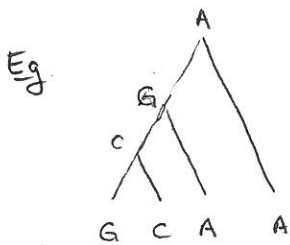
# 5. Reconstructing ancestral states
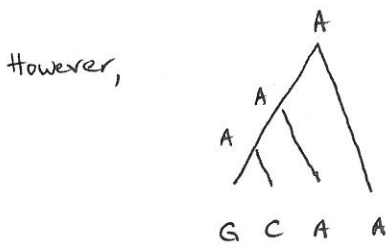
(Beware FH.)



② 

A EXTENSION $\tilde{X}$ of $X$ is a assignment of states to the internal nodes of $T$ consistent with $X$ at the leaves.

$\tilde{X}$ is minimal, if it computes $ps_X(T)$

Eg



$\tilde{X}$ is minimal ext. as is

of $X$

However,



This extension also gives rise to $ps = 2$, but does not arise from FH

Bad News: Not all labellings that give rise to $ps$ for $X_A$ on $T$ arise from FH algorithm

extensions

Good news: ② You can get some by a second pass down the tree.

Modify FH to get all ...

or use Sankoff algorithm

① The state set assigned to $\rho$ by FH is exactly the state-set for any minimal ext $\tilde{X}$ of character $X$

WEIGHTED PARSIMONY   or the   Sankoff algorithm

Underlying idea:

        not all characters are equal       → Character weighting

        not all state changes are equal     → Sankoff alg.

Character-weighting:        CODONS:    2 3 1   2 3 1

$$ps(T) = 2\, ps_{x_1}(T) + ps_{x_3}(T) + ps_{x_4}(T)$$

→ $x_1\ x_2\ x_3$

twice as important

Any scheme for character-weighting is legitimate if coefficients = weights > 0

In practice, software usually wants integer weights

Weighting State Changes.

A transition is a state change    $A \leftrightarrow G$, $C \leftrightarrow T$

transversion        between purines and pyrimidines

$$\{A, G\} \longleftrightarrow \{C, T\}$$

Empirically, transitions are observed much more frequently than transversions

    FH all state changes are weighted equally

Perhaps better:    introduce a WEIGHT MATRIX, a COST MATRIX, a STEPMATRIX

Eg.   $W =$

|  |  | to |  |  |  |
|---|---|---|---|---|---|
|  |  | A | G | C | T |
| from | A | 0 | 1 | 2 | 2 |
|  | G | 1 | 0 | 2 | 2 |
|  | C | 2 | 2 | 0 | 1 |
|  | T | 2 | 2 | 1 | 0 |

Comments: this weight matrix is symmetric, but this is not necessary

We will use $w_{ij} = $ cost of changing from state $i$ to state $j$.

Indeed, the weight matrix is judgment on the cost of changes
↑ user-defined

$W =$ transversions cost $2x$ transitions

Alternative: $\chi$    3-state character

|       | $S_1$    | $S_2$    | $S_3$ |
|-------|----------|----------|-------|
| $S_1$ | 0        | 1        | 2     |
| $S_2$ | $\infty$ | 0        | 1     |
| $S_3$ | $\infty$ | $\infty$ | 0     |

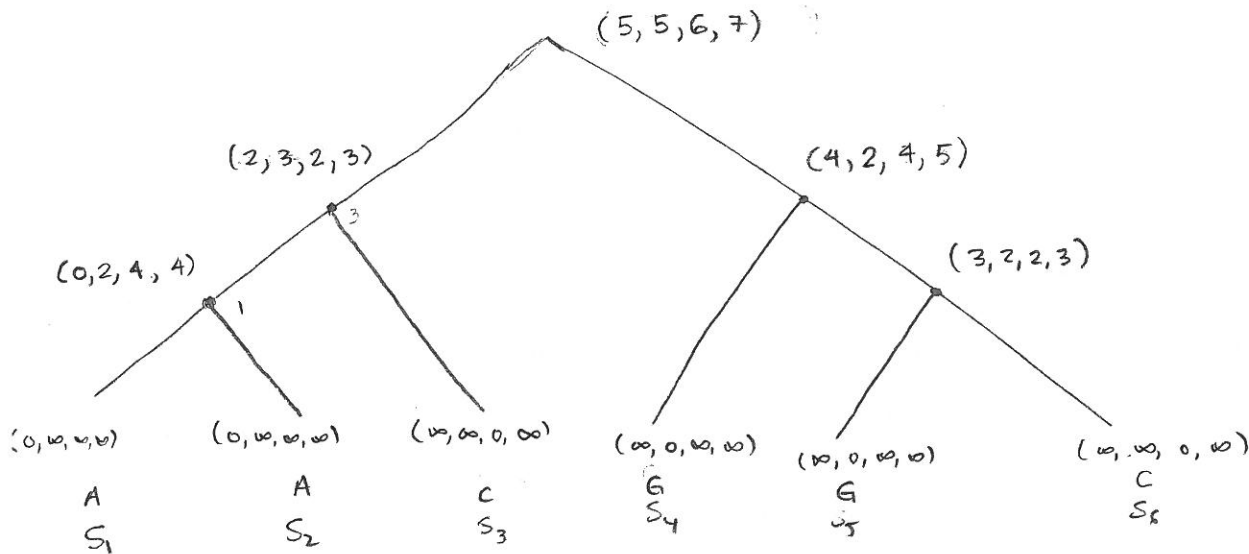$S_1 \longrightarrow S_2 \longrightarrow S_3$

See for eg
PAUP* manual on web
pp 16-20

The Sankoff algorithm for weighted-parsimony by example:   from
with transversions twice the cost of transition.   $W =$   to:

                    (5, 5, 6, 7)

        (2, 3, 2, 3)                    (4, 2, 4, 5)

(0, 2, 4, 4)          3                    (3, 2, 2, 3)

            1

(0,∞,∞,∞)  (0,∞,∞,∞)  (∞,∞,0,∞)   (∞,0,∞,∞)  (∞,0,∞,∞)   (∞,∞,0,∞)

   A          A          C            G          G            C
   $S_1$      $S_2$      $S_3$        $S_4$      $S_5$        $S_6$

Main idea: starting at the leaves, and moving up toward root $\rho$

Place a $\sigma = 4$ element vector with each entry containing the minimal
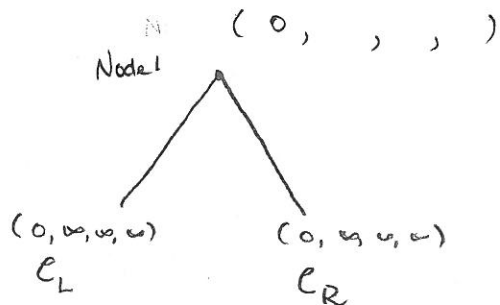cost of labelling that node with state $S_i$

For this particular example, we order the vectors
                              purines first, pyrimidines 2nd

$$\begin{array}{|c|c|c|c|} \hline \phantom{x} & \phantom{x} & \phantom{x} & \phantom{x} \\ \hline \end{array}$$
  A    G    C    T

At each tip, form its count vector by placing a $0$ in the $i$th slot if the tip is in state $i$; an $\infty$ elsewhere.

Move "up" through the tree. focusing on an internal node and its immediate children $e_L + e_R$

N $(0, \quad , \quad , \quad )$

Node1

$(0, \infty, \infty, \infty)$
$e_L$

$(0, \infty, \infty, \infty)$
$e_R$

Use the two children and weights $W_{ij}$ to compute the minimal count vector.

Eg: Node 1:  A

Cost of changes
min cost at $\rightarrow$
$e_L$
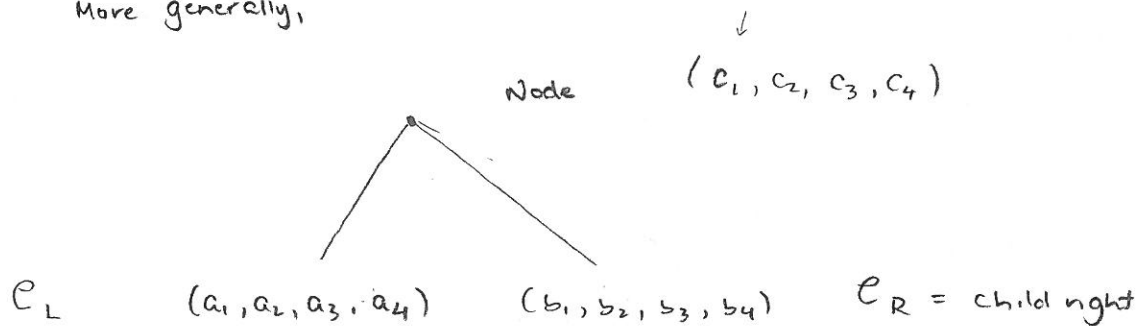
0  1  2  2
$(0, \infty, \infty, \infty)$

$(0, \infty, \infty, \infty)$

0  +  0

$\text{minimize} \left\{ (\text{cost of change } a \rightarrow j) + (\text{cost of } j \text{ at } e_L) \right\} +$

More generally,

Node $(C_1, C_2, C_3, C_4)$

$e_L$  $(a_1, a_2, a_3, a_4)$  $(b_1, b_2, b_3, b_4)$  $e_R = \text{child right}$

$\min \left\{ \text{Cost of } A \rightarrow j + \quad \right\} + \min \left\{ \text{Cost of } A \rightarrow k + \quad \right\}$

$\text{Cost of } j$
$\text{at } e_L$

$\text{Cost of } e_R(k)$

$= \min \left\{ w_{1j} + a_j \right\} + \min \left\{ w_{1k} + b_k \right\}$

Sankoff algorithm to compute the parsimony score for an s-state character

$X$ on $T$ with weight matrix $W$

    i.e. compute $PS_X(T)$ using weights $W = (w_{ij})$

    0. Arbitrarily root $T$ to get $T^{\rho}$

    1. At each tip of the tree, place an s-element vector
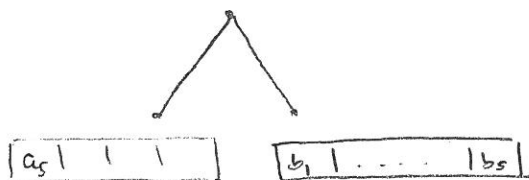
If tip is in state $i$, place a $0$ in the $i$-th position

Place an $\infty$ elsewhere.

    2. Move up the tree through all internal nodes.

At each internal node, place an s-element vector with entries
computed as follows:

$$\boxed{\; c_1 \mid c_2 \mid \quad \mid \quad \mid c_s \;}$$



$$\text{Let} \quad c_i = \min_{j \in S} \{ w_{ij} + a_j \} + \min_{j \in S} \{ w_{ij} + b_j \}$$

states
$1, 2, \ldots, s$

$$\boxed{\; a_s \mid \quad \mid \quad \mid \;} \qquad \boxed{\; b_1 \mid \ldots \ldots \mid b_s \;}$$

Cost of state
change $i \to j$

$+$ cost node in state $j$
    left child

    3. When you reach the root $\rho$,

Consider its vector $\boxed{\; c_1 \mid c_2 \mid \quad \mid \quad \mid c_s \;}$
   Count

Then $PS_X(T) = $ minimum entry.

Theorem: The Sankoff algorithm does compute $PS_x(T)$ with weights $W$

Comments: This is an example of a "dynamic programming algorithm"

One advantage: All minimal exts are computable by examining the partial count vectors.

• If the weight matrix $W$ is symmetric, then ps evaluates unrooted trees
is **not** symmetric, "    "    is a measure on **rooted** trees

Reflections:

1. There are many other types of parsimony possible     See Chap 7
Felsenstein

(Good project possibly?)

2.           Strengths              vs.        Weaknesses

heuristics are needed
since we can't explore
tree space     $(2n-5)!!$

Change is rare is
reasonable.

What if texon have evolved
over a "long" period of
time?    Is it reasineble
to avoid consideration of
multiple charger!
$$L \to M \to S$$
parallel evolution?

What is the meaning of
the parsimony scores?

For morphological data,
this seems a good choice

What is the meaning of
$W_{ij}$?