

Evolutionary Trees.

- Unrooted / rooted evolutionary trees.
- Binary weighted trees where every internal vertex has degree equal to 3 and every edge has an assigned +ve wt.

★ Distance-based tree reconstruction:

Problem statement:

Reconstruct an evolutionary tree from a distance matrix.

I/P: $\{D_{i,j}\}_{n \times n}$ Assuming $D_{i,j}$ satisfies metric conditions
 O/P: A weighted unrooted tree T with n leaves fitting D v.e. $d_{i,j}(T) = D_{i,j}$,
 $1 \leq i < j \leq n$ if $D_{i,j}$ is additive (w.e. \exists a solution).

- T will have $2n-3$ edges. Hence $2n-2$ vertices, out of which n leaves & $n-2$ internal vertices.

- 3-vertex solution:



$$\therefore d_{i,c} = \frac{D_{i,j} + D_{i,k} - D_{j,k}}{2}$$

$$d_{j,c} = \frac{D_{j,i} + D_{j,k} - D_{i,k}}{2}$$

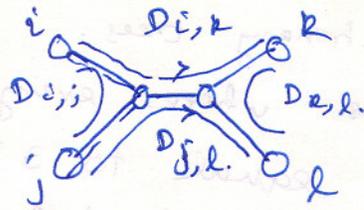
$$d_{k,c} = \frac{D_{k,i} + D_{k,j} - D_{i,j}}{2}$$

- For n leaves: no. of variables $(2n-3)$ (no. of edges) & no. of eqns. available

$\binom{n}{2}$ [overconstrained]

→ \exists may not be any solution.
 → \exists if \exists a solution, $(D_{i,j})_{n \times n}$ is called an additive matrix.

Four point condition for checking additive matrix



Compute 3 sums:

$$D_{i,j} + D_{k,l}$$

$$D_{i,k} + D_{j,l}$$

$$D_{i,l} + D_{j,k}$$

Two of these would be same & one will be smaller than the other 2.

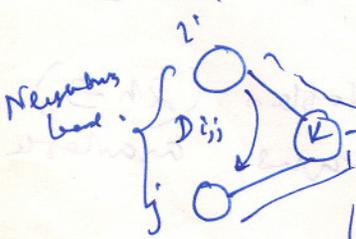
Th: An $n \times n$ matrix D is additive iff the four point conditions hold for every 4 distinct elements $1 \leq i, j, k, l \leq n$.

Least-Square Distance based Phylogeny problem:

If the distance matrix is not additive, one might want instead to find a tree that approximates D using the sum of squared errors $\sum \sum (d_{i,j}(T) - D_{i,j})^2$ as a measure of the quality of approximation. This leads to the (NP-hard) least squares distance based phylogeny problem.

Algorithms for Additive Phylogenetic Tree Reconstruction:

Collapsing neighboring leaves:

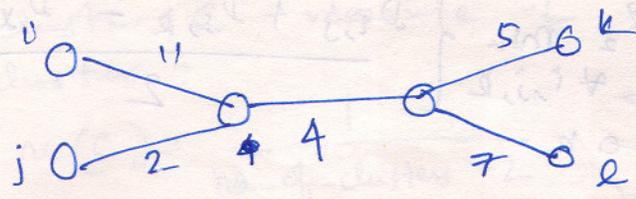


$$D_{k,m} = \frac{D_{i,m} + D_{j,m} - D_{i,j}}{2}$$

$T_n \Rightarrow$ Tree of n -vertices.

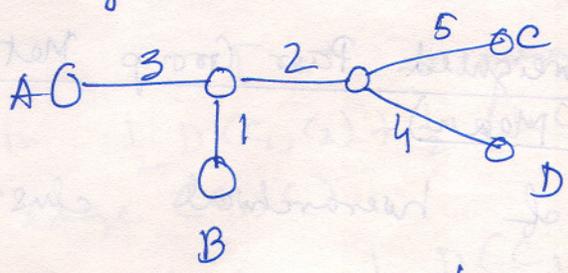
\therefore iteratively collapse neighboring vertices & update distances reduce tree T_n to T_{n-1} & continue till arrived at T_1 .

→ Problem of finding neighbouring vertices, which is non-trivial. e.g.

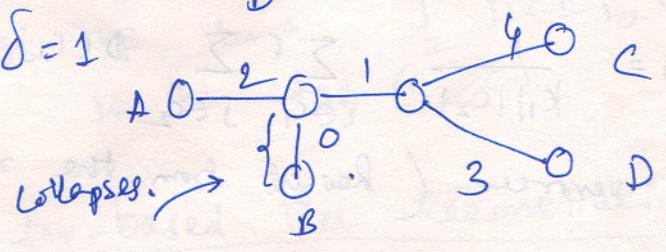


Threshold
 $D_{j,k} < D_{i,j}$
 j is not neighbour of k but of i.

- Collapsing isolated vertices (from hanging edge) by iterative pruning of all the ^{leaf} edges (to leaves) by the same amount (δ) e.g. 1 at any stage.



$\delta = 1$



→ A triple of distinct elements $1 \leq i, j, k \leq n$ is called degenerate, if $D_{ij} + D_{jk} = D_{ik}$.

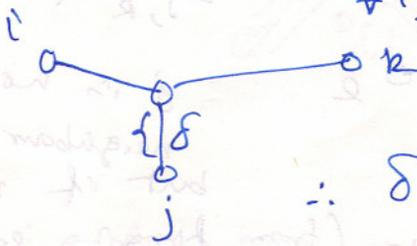
We call the entire matrix D is ~~is~~ degenerate if it has at least one degenerate triple. Convert $[D]_{n \times n}$ to $[D]_{(n-1) \times (n-1)}$ by removing j^{th} column & i^{th} row by making size of $(n-1)$ vertices.

→ Iteratively collapse every vertex to get intermediate nodes & connections (dist. for the nodes).

→ δ : trimming parameters.

Min. δ ?

$\delta = \min_{i,j,k} \left\{ \frac{D_{i,j} + D_{j,k} - D_{i,k}}{2} \right\}$



[J, R].

$\therefore \delta = \frac{D_{i,j} + D_{j,k} - D_{i,k}}{2}$

★ Evolutionary trees and hierarchical clustering:

• UPGMA (Unweighted Pair Group Method with Arithmetic Mean)

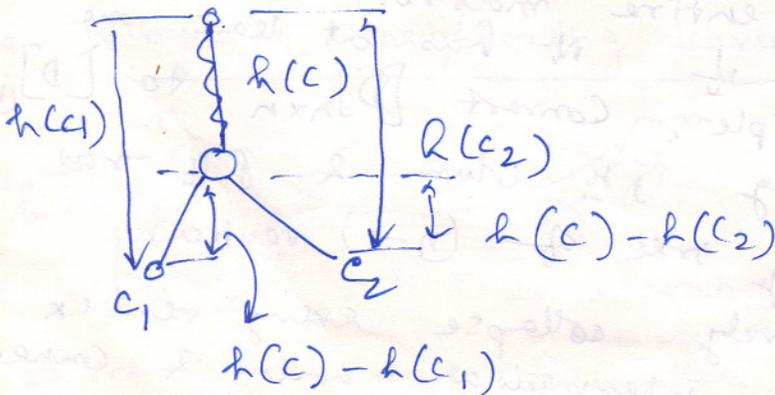
→ A variant of hierarchical clustering

→ $D(c_1, c_2) = \frac{1}{k_1 + k_2} \sum_{i \in c_1} \sum_{j \in c_2} D(i, j)$

→ "dating" vertices (height from the root).

→ $h(c) = D(c_1, c_2)$.

→ $\begin{cases} h(c) - h(c_1) & \text{to the edge } (c_1, c) \\ h(c) - h(c_2) & \text{to the edge } (c_2, c). \end{cases}$



The ultrametric tree i.e. the distances from root to leaves are same for all leaves.

• Neighbor joining Alg. (Saitou (1987) & Nei)

→ For a cluster c measure of separation from other clusters:

$$u(c) = \frac{1}{\text{no. of clusters} - 2} \sum_{\forall c'} D(c, c')$$

→ For merging two clusters c_1 & c_2

minimize $\rightarrow D(c_1, c_2) - u(c_1) - u(c_2)$

→ Use same hierarchical clustering.

$$\rightarrow D(c^*, c) = \frac{D(c_1, c) + D(c_2, c)}{2}$$

[c^* → merged cluster for c_1 & c_2]

→ Assign length $\frac{1}{2} D(c_1, c_2) + \frac{1}{2} (u(c_1) - u(c_2))$ to the edge (c_1, c)

→ Assign the length $\frac{1}{2} D(c_1, c_2) + \frac{1}{2} (u(c_2) - u(c_1))$ to the edge (c_2, c) .

★ Character-based Tree Reconstruction:

→ Character ↔ characteristic represented by m length DNA seq. (or Protein seq.)

→ For n species ~~with~~ each characterized by m ~~DNA~~ length DNA seq.

→ Each leaf node represent a species.

→ Internal node ~~should differ~~ characterized by hypothetical ' m ' length str.

→ Minimize $\sum_{\text{edge} \rightarrow (v_1, v_2) \in \text{Tree}} d_{HT}(v_1, v_2)$ ↔ Parsimony Score

③ Small Parsimony Problem:

I/P: Tree T with each leaf labeled by an m -character string.
O/P: Labelling of internal vertices of the tree T minimizing the parsimony score.

④ Weighted small parsimony problem:

Instead of $d_H(\vec{x}_1, \vec{x}_2)$ use $\sum_{i=1}^m \delta(x_{1i}, x_{2i})$
as a distance measure between two m -length strings, where, x_i 's \in Alphabet set e.g. $\{A, C, T, G\}$

I/P: Tree T with each leaf node labeled by elements of k -letter alphabet (of length m strings) and a scoring matrix (δ_{ij}) .

O/P: Labelling of internal vertices of the tree minimizing the parsimony score.

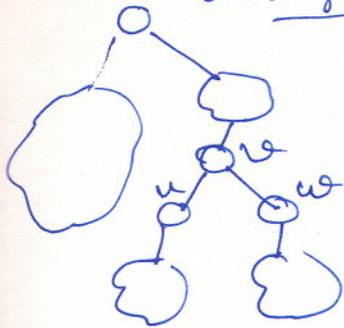
⑤ Large Parsimony Problem \rightarrow NP-Complete Problem

Tree is not given. Only, m -character strings for n -species are given & problem is to find 'Tree' with 'Parsimonious Labelling of vertices'.

I/P: An $n \times m$ matrix M describing n -spec.
O/P: A tree T with nodes labeled by n -rows of M & a labelling of internal vertices so that Min. score is obtained.

• Solutions of Small parsimony Problems

→ wt. Small parsimony tree (Sankoff (1975))
 • Single character alphabet strings:



Let $S_t(w) = \min$ parsimony score of the subtree under 'w', given label at v_t .

$$\therefore S_t(w) = \min_i \{ S_i(u) + \delta_{i,t} \} + \min_j \{ S_j(w) + \delta_{j,t} \}$$

- At leaves, $S_t(w) = 0$, if label is t else ∞ .
- Solve single alphabet assignment at every position of the string independently.
- Reconstruct the vertices by backtracking.
- At every node, computation is $O(k)$.
- Total runtime: $O(nk)$.
- as No. of vertices: $O(n)$.

→ unweighted small parsimony prob.
 (Fitch (1971)).

$$1. S_w = \begin{cases} S_u \cap S_v, & \text{if } S_u \cap S_v \neq \emptyset \\ S_u \cup S_v, & \text{otherwise.} \end{cases}$$

[From bottom to top]

2. From top to bottom:

- Root(v) → any label x from S_v .
- Label(w) → any label x from S_w if x is assigned to its parents else any label from S_v .

→ $O(nk)$.