Inferring the Past: Phylogenetic Trees (chapter 12)

- p The biological problem
- p Parsimony and distance methods
- p Models for mutations and estimation of distances
- p Maximum likelihood methods

Parsimony method

p The parsimony method finds the tree that explains the observed sequences with a minimal number of substitutions

p Method has two steps

- n Compute smallest number of substitutions for a given tree with a *parsimony algorithm*
- n Search for the tree with the minimal number of substitutions

Parsimony: an example

p Consider the following short sequences

- 1 ACTTT
- 2 ACATT
- 3 AACGT
- 4 AATGT
- 5 AATTT
- p There are 105 possible rooted trees for 5 sequences
- p Example: which of the following trees explains the sequences with least number of substitutions?



This tree explains the sequences with 4 substitutions



Computing parsimony

- Parsimony treats each site (position in a sequence) independently
- P Total parsimony cost is the sum of parsimony costs (=required substitutions) of each site
- P We can compute the minimal parsimony cost for a given tree by
 - First finding out possible assignments at each node, starting from leaves and proceeding towards the root
 - n Then, starting from the root, assign a letter at each node, proceeding towards leaves

Labelling tree nodes

- P An unrooted tree with n leaves contains 2n-1 nodes altogether
- p Assign the following labels to nodes in a rooted tree
 - n leaf nodes: 1, 2, ..., n
 - n internal nodes: n+1, n+2, ..., 2n-1
 - n root node: 2n-1
- P The label of a child node is always smaller than the label of the parent node



Parsimony algorithm: first phase

- P Find out possible assignments at every node for each site u independently. Denote site u in sequence i by s_{i,u.}
 - For i := 1, ..., n do
 - $F_i := \{s_{i,u}\}$ % possible assignments at node i
 - $L_i := 0$ % number of substitutions up to node i

For i := n+1, ..., 2n-1 do Let j and k be the children of node i If $F_j \cap F_k = \emptyset$ then $L_i := L_j + L_k + 1$, $F_i := F_j \cup F_k$ else $L_i := L_j + L_k$, $F_i := F_j \cap F_k$

Parsimony algorithm: first phase

Choose u = 3 (for example, in general we do this for all sites) $F_1 := \{T\}$ $L_1 := 0$ $F_2 := \{A\}$ $L_2 := 0$ $F_3 := \{C\}, L_3 := 0$ $F_4 := \{T\}, L_4 := 0$ $F_5 := \{T\}, L_5 := 0$ AACOT AATOT AATOT ACTT

Parsimony algorithm: first phase

 $F_6 := F_3 \cup F_4 = \{C, T\}$

 $L_6 := L_3 + L_4 + 1 = 1$

 $\begin{array}{l} {\sf F}_7 := {\sf F}_5 \cap {\sf F}_6 = \{{\sf T}\} \\ {\sf L}_7 := {\sf L}_5 + {\sf L}_6 = 1 \end{array}$

 $F_8 := F_1 \cup F_2 = \{A, T\}$ $L_8 := L_1 + L_2 + 1 = 1$

 $\begin{array}{l} \mathsf{F_9} := \, \mathsf{F_7} \cap \, \mathsf{F_8} = \{\mathsf{T}\} \\ \mathsf{L_9} := \, \mathsf{L_7} \, + \, \mathsf{L_8} = \, 2 \end{array}$

 Image: state of the state o

 \Rightarrow Parsimony cost for site 3 is 2

Parsimony algorithm: second phase

- **p** Backtrack from the root and assign $x \in F_i$ at each node
- **p** If we assigned y at parent of node i and y $\in F_i$, then assign y
- p Else assign $x \in F_i$ by random

Parsimony algorithm: second phase

At node 6, the algorithm assigns T because T was assigned to parent node 7 and T \in F₆. T is assigned to node 8 for the same reason.



The other nodes have only one possible letter to assign

Parsimony algorithm

First and second phase are repeated for each site in the sequences, summing the parsimony costs at each site



Properties of parsimony algorithm

- Parsimony algorithm requires that the sequences are of same length
 - First align the sequences against each other and, optionally, remove indels
 - n Then compute parsimony for the resulting sequences
 - n Indels (if present) considered as characters
- p Is the most parsimonious tree the correct tree?
 - n Not necessarily but it explains the sequences with least number of substitutions
 - n We can assume that the probability of having fewer mutations is higher than having many mutations

Finding the most parsimonious tree

- p Parsimony algorithm calculates the parsimony cost for a given tree...
- p ...but we still have the problem of finding the tree with the lowest cost
- p Exhaustive search (enumerating all trees) is in general impossible
- p More efficient methods exist, for example
 - n Probabilistic search
 - n Branch and bound

Branch and bound in parsimony

P We can exploit the fact that adding edges to a tree can only increase the parsimony cost {C, T}



Branch and bound in parsimony

Branch and bound is a general search strategy where

- p Each solution is potentially generated
- p Track is kept of the best solution found
- p If a partial solution cannot achieve better score, we abandon the current search path

In parsimony...

- p Start from a tree with 1 sequence
- p Add a sequence to the tree and calculate parsimony cost
- p If the tree is complete, check if found the best tree so far
- p If tree is not complete and cost exceeds best tree cost, do not continue adding edges to this tree

Branch and bound example

Complete tree: compute parsimony cost



Distance methods

- p The parsimony method works on sequence (character string) data
- p We can also build phylogenetic trees in a more general setting
- p Distance methods work on a set of pairwise distances d_{ii} for the data
- p Distances can be obtained from phenotypes as well as from genotypes (sequences)

Distances in a phylogenetic tree

- p Distance matrix D = (d_{ij}) gives pairwise distances for *leaves* of the phylogenetic tree
- P In addition, the phylogenetic tree will now specify distances between leaves and internal nodes
 - n Denote these with d_{ij} as well



Distance d_{ij} states how far apart species i and j are evolutionary (e.g., number of mismatches in aligned sequences)

Distances in evolutionary context

p Distances d_{ij} in evolutionary context satisfy the following conditions

n Symmetry: $d_{ij} = d_{ji}$ for each i, j

n Distinguishability: $d_{ij} \neq 0$ if and only if $i \neq j$

n Triangle inequality: $d_{ij} \le d_{ik} + d_{kj}$ for each i, j, k

- p Distances satisfying these conditions are called metric
- In addition, evolutionary mechanisms may impose additional constraints on the distances
 additive and *ultrametric* distances

Additive trees

P A tree is called *additive*, if the distance between any pair of leaves (i, j) is the sum of the distances between the leaves and a node k on the shortest path from i to j in the tree

$$d_{ij} = d_{ik} + d_{jk}$$

p "Follow the path from the leaf i to the leaf j to find the exact distance d_{ij} between the leaves."

Additive trees: example



Ultrametric trees

P A rooted additive tree is called an *ultrametric tree*, if the distances between any two leaves i and j, and their common ancestor k are equal

 $d_{ik} = d_{jk}$

- Edge length d_{ij} corresponds to the time elapsed since divergence of i and j from the common parent
- P In other words, edge lengths are measured by a molecular clock with a constant rate

Identifying ultrametric data

- p We can identify distances to be ultrametric by the three-point condition:
 - D corresponds to an ultrametric tree if and only if for any three species i, j and k, the distances satisfy $d_{ij} \leq max(d_{ik}, d_{kj})$

p If we find out that the data is ultrametric, we can utilise a simple algorithm to find the corresponding tree









UPGMA algorithm

- p UPGMA (unweighted pair group method using arithmetic averages) constructs a phylogenetic tree via clustering
- p The algorithm works by at the same time
 - n Merging two clusters
 - n Creating a new node on the tree
- p The tree is built from leaves towards the root
- p UPGMA produces a ultrametric tree

Cluster distances

p Let distance d_{ij} between clusters C_i and C_j be $d_{ij} = \frac{1}{|C_i||C_j|} \sum_{p \in C_i, q \in C_j} d_{pq}$

that is, the average distance between points (species) in the cluster.

UPGMA algorithm

- p Initialisation
 - n Assign each point i to its own cluster C_i
 - n Define one leaf for each sequence, and place it at height zero
- p Iteration
 - n Find clusters i and j for which d_{ii} is minimal
 - **n** Define new cluster k by $C_k = C_i \cup C_j$, and define d_{kl} for all I
 - n Define a node k with children i and j. Place k at height $d_{ij}/2$
 - n Remove clusters i and j
- **p** Termination:
 - n When only two clusters i and j remain, place root at height $d_{ii}/2$











UPGMA implementation

- p In naive implementation, each iteration takes O(n²) time with n sequences => algorithm takes O(n³) time
- p The algorithm can be implemented to take only O(n²) time (see Gronau & Moran, 2006, for a survey)

Problem solved?

- P We now have a simple algorithm which finds a ultrametric tree
 - n If the data is ultrametric, then there is exactly one ultrametric tree corresponding to the data (we skip the proof)
 - n The tree found is then the "correct" solution to the phylogeny problem, if the assumptions hold
- p Unfortunately, the data is not ultrametric in practice
 - n Measurement errors distort distances
 - n Basic assumption of a molecular clock does not hold usually very well

Incorrect reconstruction of nonultrametric data by UPGMA



Tree which corresponds to non-ultrametric distances Incorrect ultrametric reconstruction by UPGMA algorithm

Checking for additivity

p How can we check if our data is additive?
p Let i, j, k and I be four *distinct* species
p Compute 3 sums: d_{ij} + d_{kl}, d_{ik} + d_{jl}, d_{il} + d_{jk}

Four-point condition



p The sums are represented by the three figures

- n Left and middle sum cover all edges, right sum does not
- *Four-point condition*: i, j, k and I satisfy the fourpoint condition if two of the sums d_{ij} + d_{kl}, d_{ik} + d_{jl}, d_{il} + d_{jk} are the same, and the third one is smaller than these two

Checking for additivity

p An n x n matrix D is additive if and only if the four point condition holds for every 4 distinct elements 1 ≤ i, j, k, l ≤ n

Finding an additive phylogenetic tree

- Additive trees can be found with, for example, the neighbor joining method (Saitou & Nei, 1987)
- P The neighbor joining method produces unrooted trees, which have to be rooted by other means
 - n A common way to root the tree is to use an outgroup
 - n Outgroup is a species that is known to be more distantly related to every other species than they are to each other
 - n Root node candidate: position where the outgroup would join the phylogenetic tree
- P However, in real-world data, even additivity usually does not hold very well

- p Neighbor joining works in a similar fashion to UPGMA
 - n Find clusters C_1 and C_2 that minimise a function $f(C_1, C_2)$
 - n Join the two clusters C₁ and C₂ into a new cluster C
 - n Add a node to the tree corresponding to C
 - n Assign distances to the new branches
- p Differences in
 - n The choice of function $f(C_1, C_2)$
 - n How to assign the distances

p Recall that the distance d_{ij} for clusters C_i and C_j was

$$d_{ij} = \frac{1}{|C_i||C_j|} \sum_{p \in C_i, q \in C_j} d_{pq}$$

p Let u(C_i) be the separation of cluster C_i from other clusters defined by

$$u(C_i) = \frac{1}{n-2} \sum_{C_j} d_{ij}$$

where n is the number of clusters.

- p Instead of trying to choose the clusters C_i and C_j closest to each other, neighbor joining at the same time
 - n Minimises the distance between clusters C_{i} and C_{j} and
 - ${\bf n}$ Maximises the separation of both ${\bf C}_{i}$ and ${\bf C}_{j}$ from other clusters

- p Initialisation as in UPGMA
- p Iteration
 - n Find clusters i and j for which $d_{ij} u(C_i) u(C_j)$ is minimal
 - **n** Define new cluster k by $C_k = C_i \cup C_j$, and define d_{kl} for all I
 - n Define a node k with edges to i and j. Remove clusters i and j
 - n Assign length $\frac{1}{2} d_{ij} + \frac{1}{2} (u(C_i) u(C_j))$ to the edge i -> k
 - n Assign length $\frac{1}{2} d_{ij} + \frac{1}{2} (u(C_j) u(C_i))$ to the edge j -> k
- **p** Termination:
 - n When only one cluster remains

Neighbor joining algorithm: example

	a	b	C (<u> </u>	i	u(i)
а	0	6	7 !	5	а	(6+7+5)/2 = 9
b		0 1	1 9	9	b	(6+11+9)/2 = 13
С			0	5	С	(7+11+6)/2 = 12
d				С	d	(5+9+6)/2 = 10



Neighbor joining algorithm: example

	a	b	С	d	_	i	u(i)
а	0	6	7	5		a	(6+7+5)/2 = 9
b		0	11	9		b	(6+11+9)/2 = 13
С			0	6		С	(7+11+6)/2 = 12
d				0		d	(5+9+6)/2 = 10



This is the first step only...

Inferring the Past: Phylogenetic Trees (chapter 12)

- p The biological problem
- p Parsimony and distance methods
- p Models for mutations and estimation of distances
- p Maximum likelihood methods
 - n These parts of the book is skipped on this course (see slides of 2007 course for material on these topics)
 - n No questions in exams on these topics!

Problems with tree-building

p Assumptions

- n Sites evolve independently of one other
- n (Sites evolve according to the same stochastic model; not really covered this year)
- n The tree is rooted
- n The sequences are aligned
- n Vertical inheritance

Additional material on phylogenetic trees

- p Durbin, Eddy, Krogh, Mitchison: Biological sequence analysis
- p Jones, Pevzner: An introduction to bioinformatics algorithms
- p Gusfield: Algorithms on strings, trees, and sequences
- p Course on phylogenetic analyses in Spring 2009