Algorithmic Approaches for Biological Data, Lecture #24

Katherine St. John

City University of New York American Museum of Natural History

4 May 2016

<ロト <四ト <注入 <注下 <注下 <

Outline



- Approximation Algorithms: UPGMA & Neighbor Joining
- Searching for Optimal Trees
- Early Project Presentations

3

• • • • • • • • • • • •



Idea:

 Compute Hamming distance between sequences.



- Compute Hamming distance between sequences.
- Base a tree only on the distances.



- Compute Hamming distance between sequences.
- Base a tree only on the distances.
- Common Algorithms:



- Compute Hamming distance between sequences.
- Base a tree only on the distances.
- Common Algorithms:
 - Unweighted Pair Group Method with Arithmetic Mean (UPGMA) (Sokal & Michener 1958)



- Compute Hamming distance between sequences.
- Base a tree only on the distances.
- Common Algorithms:
 - Unweighted Pair Group Method with Arithmetic Mean (UPGMA) (Sokal & Michener 1958)
 - Neighbor Joining (Saitou-Nei 1987)

Unweighted Pair Group Method with Arithmetic Mean (UPGMA) (Sokal & Michener 1958)

イロン イヨン イヨン イヨン

• Compute the distance between every pair of leaves.

イロト イポト イヨト イヨト

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.

(日) (同) (三) (三)

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

- 4 同 6 4 日 6 4 日 6

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

$$d(A \cup B, X) = \frac{\sum_{x \in X} \sum_{a \in A} d(x, a) + \sum_{x \in X} \sum_{b \in B} d(x, b)}{|X|(|A| + |B|)}$$

- 4 同 6 4 日 6 4 日 6

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

$$d(A \cup B, X) = \frac{\sum_{x \in X} \sum_{a \in A} d(x, a) + \sum_{x \in X} \sum_{b \in B} d(x, b)}{|X|(|A| + |B|)}$$

• Repeat until only one cluster.

A B b

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

$$d(A \cup B, X) = \frac{\sum_{x \in X} \sum_{a \in A} d(x, a) + \sum_{x \in X} \sum_{b \in B} d(x, b)}{|X|(|A| + |B|)}$$

• Repeat until only one cluster.

| | а | b | c | d | е |
|---|----|----|----|----|----|
| a | 0 | 17 | 21 | 31 | 23 |
| b | 17 | 0 | 30 | 34 | 21 |
| c | 21 | 30 | 0 | 28 | 39 |
| d | 31 | 34 | 28 | 0 | 43 |
| e | 23 | 21 | 39 | 43 | 0 |

A B b

Unweighted Pair Group Method with Arithmetic Mean (UPGMA) (Sokal & Michener 1958)

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

$$d(A \cup B, X) = \frac{\sum_{x \in X} \sum_{a \in A} d(x, a) + \sum_{x \in X} \sum_{b \in B} d(x, b)}{|X|(|A| + |B|)}$$

• Repeat until only one cluster.

| | | a | - 1 | b | | • | c | | e | |
|------|----------|-----|-----|----|----|---|-----|-----|----|---|
| a | | 0 | 1 | 7 | 2 | 1 | 3 | 1 | 2 | 3 |
| b | 1 | 17 | | D | 3 | 0 | 3 | 4 | 2 | 1 |
| c | 2 | 21 | 3 | 0 | (|) | 2 | в | 3 | 9 |
| d | 3 | 31 | 3 | 14 | 2 | 8 | C | | 43 | 3 |
| e | 2 | 23 | 2 | 1 | 3 | 9 | 4 | 3 | 0 | |
| (0) | • | (a, | b) | 0 | ; | | d | | e | |
| (a,i |) | U | - | 20 | .5 | 3 | 2.5 | | 2 | |
| c | | 25 | .5 | (|) | 2 | 28 | - 2 | 39 | |
| d | | 32 | .5 | 2 | 8 | | D | | 13 | |

22 39 43 0

< ∃ > <

Unweighted Pair Group Method with Arithmetic Mean (UPGMA) (Sokal & Michener 1958)

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

$$d(A \cup B, X) = \frac{\sum_{x \in X} \sum_{a \in A} d(x, a) + \sum_{x \in X} \sum_{b \in B} d(x, b)}{|X|(|A| + |B|)}$$

• Repeat until only one cluster.

| | а | b | c | d | е |
|---|----|----|----|----|----|
| а | 0 | 17 | 21 | 31 | 23 |
| b | 17 | 0 | 30 | 34 | 21 |
| c | 21 | 30 | 0 | 28 | 39 |
| d | 31 | 34 | 28 | 0 | 43 |
| e | 23 | 21 | 39 | 43 | 0 |

| | (a,b) | c | d | е |
|-------|-------|------|------|----|
| (a,b) | 0 | 25.5 | 32.5 | 22 |
| c | 25.5 | 0 | 28 | 39 |
| d | 32.5 | 28 | 0 | 43 |
| е | 22 | 39 | 43 | 0 |

| | ((a,b),e) c | | d |
|-----------|-------------|----|----|
| ((a,b),e) | 0 | 30 | 36 |
| c | 30 | 0 | 28 |
| d | 36 | 28 | 0 |

.

Unweighted Pair Group Method with Arithmetic Mean (UPGMA) (Sokal & Michener 1958)

- Compute the distance between every pair of leaves.
- Make a cluster for each leaf.
- Join the two clusters with lowest distance.
- Average the weights to get the new distances to other clusters.

$$d(A \cup B, X) = \frac{\sum_{x \in X} \sum_{a \in A} d(x, a) + \sum_{x \in X} \sum_{b \in B} d(x, b)}{|X|(|A| + |B|)}$$

• Repeat until only one cluster.

| | a | b | c | d | е |
|---|----|----|----|----|----|
| a | 0 | 17 | 21 | 31 | 23 |
| b | 17 | 0 | 30 | 34 | 21 |
| с | 21 | 30 | 0 | 28 | 39 |
| d | 31 | 34 | 28 | 0 | 43 |
| e | 23 | 21 | 39 | 43 | 0 |

| | (a,b) | c | d | е |
|-------|-------|------|------|----|
| (a,b) | 0 | 25.5 | 32.5 | 22 |
| c | 25.5 | 0 | 28 | 39 |
| d | 32.5 | 28 | 0 | 43 |
| е | 22 | 39 | 43 | 0 |

| | ((a,b),e) | c | d |
|-----------|-----------|----|----|
| ((a,b),e) | 0 | 30 | 36 |
| c | 30 | 0 | 28 |
| d | 36 | 28 | 0 |

| | ((a,b),e) | (c,d) |
|-----------|-----------|-------|
| ((a,b),e) | 0 | 33 |
| (c,d) | 33 | 0 |

wiki

Neighbor Joining (Saitou & Nei 1987)

• A similar idea, but a different matrix is computed:

$$Q(i,j) = (n-2)d(i,j) - \sum_{k=1}^{n} d(i,k) - \sum_{k=1}^{n} d(j,k)$$

< A

- ₹ 🗦 🕨





Neighbor Joining (Saitou & Nei 1987)

 A similar idea, but a different matrix is computed:

$$Q(i,j) = (n-2)d(i,j) - \sum_{k=1}^{n} d(i,k) - \sum_{k=1}^{n} d(j,k)$$

• If A and B are joined to new node u, then the distance from k to u

$$d(u,k) = \frac{1}{2}[d(A,k) + d(B,k) - d(A,B)]$$

< ∃ >



Neighbor Joining (Saitou & Nei 1987)

 A similar idea, but a different matrix is computed:

$$Q(i,j) = (n-2)d(i,j) - \sum_{k=1}^{n} d(i,k) - \sum_{k=1}^{n} d(j,k)$$

• If A and B are joined to new node u, then the distance from k to u

$$d(u,k) = \frac{1}{2}[d(A,k) + d(B,k) - d(A,B)]$$

• This relaxes the assumption of UPGMA that all lineages evolve at the same rate.

- **4 A b**

-∢ ∃ ▶



Neighbor Joining (Saitou & Nei 1987)

 A similar idea, but a different matrix is computed:

$$Q(i,j) = (n-2)d(i,j) - \sum_{k=1}^{n} d(i,k) - \sum_{k=1}^{n} d(j,k)$$

• If A and B are joined to new node u, then the distance from k to u

$$d(u,k) = \frac{1}{2}[d(A,k) + d(B,k) - d(A,B)]$$

• This relaxes the assumption of UPGMA that all lineages evolve at the same rate.

- **4 A b**

- ∢ ∃ ▶

In Pairs: Distance Methods

| | a | b | С | d | е |
|---|---|----|----|----|---|
| a | 0 | 5 | 9 | 9 | 8 |
| b | 5 | 0 | 10 | 10 | 9 |
| С | 9 | 10 | 0 | 8 | 7 |
| d | 9 | 10 | 8 | 0 | 3 |
| е | 8 | 9 | 7 | 3 | 0 |

wiki

- If you run the UPGMA algorithm on the matrix, what is the resulting tree?
- What is an upper bound on the running time of UPGMA?
- What is an upper bound on the running time of Neighbor Joining?

.

treeSearch(sequences):



treeSearch(sequences):
Input: A set of n sequences of length k



treeSearch(sequences):

Input: A set of n sequences of length k **Output:** The best scoring tree we could find after 1000 steps through the NNI search space.





treeSearch(sequences):

Input: A set of n sequences of length k **Output:** The best scoring tree we could find after 1000 steps through the NNI search space.

Choose randomly a tree to be bestSoFar.



treeSearch(sequences):

Input: A set of n sequences of length k **Output:** The best scoring tree we could find after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- Score the tree, bestSoFar.



treeSearch(sequences):

Input: A set of n sequences of length k **Output:** The best scoring tree we could find after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- 2 Score the tree, bestSoFar.
- Sor 1000 steps:



treeSearch(sequences):
Input: A set of n sequences of length k

Output: The best scoring tree we could find after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- Score the tree, bestSoFar.
- Sor 1000 steps:
- if any of the NNI neighbors of bestSoFar have better score,



treeSearch(sequences): Input: A set of n sequences of length k

Output: The best scoring tree we could find after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- Score the tree, bestSoFar.
- Sor 1000 steps:
- if any of the NNI neighbors of bestSoFar have better score,
 - choose it to be the current tree



treeSearch(sequences):

Input: A set of n sequences of length k **Output:** The best scoring tree we could find after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- 2 Score the tree, bestSoFar.
- I For 1000 steps:
- if any of the NNI neighbors of bestSoFar have better score,
 - choose it to be the current tree
- 6 else:
 - break (stuck in local minima)



treeSearch(sequences):

Input: A set of n sequences of length k **Output:** The best scoring tree we could find after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- 2 Score the tree, bestSoFar.
- I For 1000 steps:
- if any of the NNI neighbors of bestSoFar have better score,
 - ochoose it to be the current tree
- 6 else:
- break (stuck in local minima)
- Print number of steps, tree, and score.

In Pairs: Translating "To Do" into Pseudocode

Translate into pseudocode:



treeSearch(sequences):
Input: A set of n sequences of length k
Output: The best scoring tree we could find
after 1000 steps through the NNI search space.

- Choose randomly a tree to be bestSoFar.
- 2 Score the tree, bestSoFar.
- For 1000 steps:
- if any of the NNI neighbors of bestSoFar have better score,
 - ochoose it to be the current tree
- 6 else:
 - break (stuck in local minima)
- Print number of steps, tree, and score.

Put all leaves into a dictionary.



Image: A matrix

- ∢ ∃ ▶



- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.



- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.
- While len(parentless) > 1:



- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.
- While len(parentless) > 1:
 - Make a new internal node and choose two from parentless to be its children.



- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.
- While len(parentless) > 1:
 - Make a new internal node and choose two from parentless to be its children.
- Semove children from parentless.



- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.
- While len(parentless) > 1:
 - Make a new internal node and choose two from parentless to be its children.
- 6 Remove children from parentless.
- 6 Add new node to parentless.



- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.
- While len(parentless) > 1:
 - Make a new internal node and choose two from parentless to be its children.
- 6 Remove children from parentless.
- 6 Add new node to parentless.
- Rename remaining node in parentless to be root.

In Pairs: Randomly Building A Tree

Translate into pseudocode:



def randomBuild(sequenceList):

- Put all leaves into a dictionary.
- Also create a "To Do" list, parentless, of nodes without parents.
- While len(parentless) > 1:
 - Make a new internal node and choose two from parentless to be its children.
- Semove children from parentless.
- 6 Add new node to parentless.
- Rename remaining node in parentless to be root.



Input: A tree, t, and an internal node, node. **Output:** Two trees that result from switching the nodes' sibling with its children.

Let sib be the sibling of node.





Input: A tree, t, and an internal node, node. **Output:** Two trees that result from switching the nodes' sibling with its children.

- Let sib be the sibling of node.
- 2 Let kid1, kid2 be the children of node.





Input: A tree, t, and an internal node, node. **Output:** Two trees that result from switching the nodes' sibling with its children.

- Let sib be the sibling of node.
- 2 Let kid1, kid2 be the children of node.
- Make two (deep) copies of t: tree1 and tree2





Input: A tree, t, and an internal node, node. **Output:** Two trees that result from switching the nodes' sibling with its children.

- Let sib be the sibling of node.
- 2 Let kid1, kid2 be the children of node.



Switch sib and kid1 in tree1.





Input: A tree, t, and an internal node, node. **Output:** Two trees that result from switching the nodes' sibling with its children.

- Let sib be the sibling of node.
- 2 Let kid1, kid2 be the children of node.



- Switch sib and kid1 in tree1.
- Switch sib and kid2 in tree2.





Input: A tree, t, and an internal node, node. **Output:** Two trees that result from switching the nodes' sibling with its children.

- Let sib be the sibling of node.
- 2 Let kid1, kid2 be the children of node.



- Switch sib and kid1 in tree1.
- Switch sib and kid2 in tree2.
- Seturn tree1 and tree2.



In Pairs: Making NNI Move

Translate into pseudocode:



def nniNeighbors(t, node):
Input: A tree, t, and an internal node, node.
Output: Two trees that result from switching the nodes'
sibling with its children.

Let sib be the sibling of node.



- Make two (deep) copies of t: tree1 and tree2
- Switch sib and kid1 in tree1.
- Switch sib and kid2 in tree2.
- Return tree1 and tree2.



Early Project Presentaitons



K. St. John (CUNY & AMNH)

Algorithms #24

4 May 2016 13 / 14

- E - N

Image: A math a math



• Trees can be estimated used distance based methods.

Image: A matrix

Recap



- Trees can be estimated used distance based methods.
- More on searching treespace in lab today.

Recap



- Trees can be estimated used distance based methods.
- More on searching treespace in lab today.
- Email lab reports to kstjohn@amnh.org.

Recap



- Trees can be estimated used distance based methods.
- More on searching treespace in lab today.
- Email lab reports to kstjohn@amnh.org.
- Challenges available at rosalind.info.