CS 466 Introduction to Bioinformatics Lecture 14

Mohammed El-Kebir

October 24, 2018



Course Announcements

HW 3 due Oct 29 by 11:59pm

Outline

- Recap additive distance
- Neighbor joining
- Character-based phylogeny (small)
- Application to cancer

Reading:

• Chapters 10.2, 10.5-10.8, 10.9 in Jones and Pevzner

Hierarchical Clustering

- 1. <u>Hierarchical Clustering</u> (*D*, *n*)
- 2. Form *n* clusters each with one element
- 3. Construct a graph **T** by assigning one vertex to each cluster
- 4. while there is more than one cluster
- 5. Find the two closest clusters C_1 and C_2
- 6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
- 7. Compute distance from *C* to all other clusters
- 8. Add a new vertex **C** to **T** and connect to vertices C_1 and C_2
- 9. Remove rows and columns of **D** corresponding to C_1 and C_2
- 10. Add a row and column to **D** corresponding to the new cluster **C**

11. return *T*

Definition of distance between clusters (or, linkage criterion) affects clustering! Organize elements into a tree s.t.:

- Leaves are elements
- Paths between leaves represent pairwise element distance
- Similar elements lie within same subtrees



Additive Distance Matrices

Matrix *D* is \longrightarrow ADDITIVE if there exists a tree *T* with $d_{ij}(T) = D_{ij}$







This is a constructive definition

A Small and a Large Problem

Small Additive Distance Phylogeny Problem: Given $n \times n$ distance matrix $D = [d_{i,j}]$ and unweighted tree T with n leaves, determine edge weights such that $d_T(i,j) = d_{i,j}$

Large Additive Distance Phylogeny Problem: Given $n \times n$ distance matrix $D = [d_{i,j}]$, find tree T with n leaves and edge weights such that $d_T(i,j) = d_{i,i}$

Both problems can be solved in polynomial time

Small Additive Distance Problem

- 1. Find neighboring leaves *i* and *j* with parent *k*
- 2. Remove the rows and columns of *i* and *j*
- 3. Add a new row and column corresponding to k, where the distance from k to any other leaf m is computed as



4. Repeat steps 1-3 until tree has only two vertices

Additive Phylogeny

AdditivePhylogeny(D) if D is a 2 x 2 matrix $T = \text{tree of a single edge of length } D_{1.2}$ return T if D is non-degenerate Compute trimming parameter δ Trim(D, δ) Find a triple *i*, *j*, *k* in *D* such that $D_{ii} + D_{ik} = D_{ik}$ $x = D_{ii}$ Remove j^{th} row and j^{th} column from D T = AdditivePhylogeny(D).Add a new vertex v to T at distance x from i to k Add j back to T by creating an edge (v,j) of length 0 for every leaf *I* in *T* if distance from *I* to *v* in the tree $\neq D_{Li}$ output "matrix is not additive" return Extend all "hanging" edges by length δ return T



Additive Distance Matrix

Four point condition of matrix $D = [d_{i,j}]$: Every four leaves (quartet) can be labeled as (i, j, k, l) such that $d_{i,j} + d_{k,l} \le d_{i,k} + d_{j,l} = d_{i,l} + d_{j,k}$

Theorem: Let D be an $n \times n$ matrix. The following statements are equivalent.

- 1. Matrix *D* is additive.
- 2. There exists a unique tree T (modulo isomorphism) s.t. $d_{i,j} = d_T(i,j)$ for all $(i,j) \in n^2$.
- 3. Four point condition holds for every quartet $(i, j, k, l) \in [n]^4$.

Outline

- Recap additive distance
- Neighbor joining
- Character-based phylogeny (small)
- Application to cancer

Reading:

• Chapters 10.2, 10.5-10.8, 10.9 in Jones and Pevzner

Distance Based Phylogeny Problem

Large Additive Distance Phylogeny Problem: Given $n \times n$ matrix $D = [d_{i,j}]$, find tree T with n leaves and edge weights such that $\max_{(i,j)\in[n]^2} |d_T(i,j) - d_{i,j}|$ is minimum.

Equivalently, find additive matrix D' closest to input matrix D

12

Neighbor Joining Algorithm (Saitou and Nei 1987)

- Constructs binary unrooted trees.
- Recall: leaves *a* and *b* are neighbors if they have a common parent
- Recall: closest leaves are not necessarily neighbors
- NJ: Find pair of leaves that are "close" to each other but "far" from other leaves

Two advantages: (1) reproduces correct tree for additive matrix, and (2) otherwise gives good approximation of correct tree



Distance Trees as Hierarchical Clustering



Distance Trees as Hierarchical Clustering



Distance Trees as Hierarchical Clustering

- 1. <u>Hierarchical Clustering</u> (*D*, *n*)
- 2. Form *n* clusters each with one element
- 3. Construct a graph **T** by assigning one vertex to each cluster
- 4. while there is more than one cluster
- 5. Find the two closest clusters C_1 and C_2
- 6. Merge C_1 and C_2 into new cluster C with $|C_1| + |C_2|$ elements
- 7. Compute distance from *C* to all other clusters
- 8. Add a new vertex **C** to **T** and connect to vertices C_1 and C_2
- 9. Remove rows and columns of **D** corresponding to C_1 and C_2
- 10. Add a row and column to **D** corresponding to the new cluster **C**

11. return *T*

Selection criterion: distance between clusters affects clustering!



Neighbor Joining: Selection Criterion

 $U_{1} = D(1,3) + D(1,2) + D(1,4)$



Let $C = \{1, \dots, n\}$ be <u>current clusters</u>/leaves.

Define $u = \sum_k D(i, k)$.

Intuitively, *u_i* measures separation of *i* from other leaves.

Goal: Minimize D(i, j) and maximize $u_i + u_j$. Separation

Solution: Find pair (i, j) that minimizes: $S_D(i, j) = (n - 2) D(i, j) - u_i - u_j$

Claim: Given additive matrix D. $S_D(x, y) = \min S_D(i, j)$ if and only if x and y are neighbors in tree T with $d_T = D$.

Neighboring Joining: Algorithm

Form *n* clusters $C_1, C_2, ..., C_n$, one for each leaf node.

Define tree T to be the set of leaf nodes, one per sequence.

Iteration: (D is $m \times m$) Pick *i*, *j* such that $S_D(i, j) = (m - 2) D(i, j) - u_i - u_j$ is minimal. Merge *i* and *j* into new node [*ij*] in *T*. Assign length ½ ($D(i, j) + 1/(m-2) (u_i - u_j)$) to edge (*i*, [*ij*]) Assign length ½ ($D(i, j) + 1/(m-2) (u_j - u_i)$) to edge (*j*, [*ij*])

Remove rows and columns from D corresponding to *i* and *j*. Add row and column to D for new vertex [*ij*]. Set D([*ij*], *m*) = $\frac{1}{2}$ [D(*i*, *m*) + D(*j*, *m*) - D(*i*,*j*)]

Termination:

When only one cluster

Neighboring Joining: Example

Initialization:

Form *n* clusters $C_1, C_2, ..., C_n$, one for each leaf node.

Define tree T to be the set of leaf nodes, one per sequence.

Iteration: (D is $m \times m$)

Pick *i*, *j* such that $S_D(i, j) = (m - 2) D(i, j) - u_i - u_j$ is minimal.

Merge *i* and *j* into new node [*ij*] in *T*. Assign length $\frac{1}{2} (D(i, j) + 1/(m-2) (u_i - u_j))$ to edge (*i*, [*ij*]) Assign length $\frac{1}{2} (D(i, j) + 1/(m-2) (u_j - u_i))$ to edge (*j*, [*ij*])

Remove rows and columns from D corresponding to *i* and *j*. Add row and column to D for new vertex [*ij*]. Set D([*ij*], *m*) = $\frac{1}{2}$ [D(*i*, *m*) + D(*j*, *m*) - D(*i*,*j*)]

Termination:

When only one cluster

	Α	В	С	D
Α	0	4	10	9
B	4	0	8	7
C	10	8	0	9
D	9	7	9	0



Advantages of Neighbor Joining

Max . [d; - d;]

Theorem: Let D be an $n \times n$ matrix. If matrix D is additive then neighbor joining produces the unique phylogenetic tree T (modulo isomorphism) such that $d_{i,j} = d_T(i,j)$ for all $(i,j) \in n^2$.

Theorem: Let D be an $n \times n$ matrix. If there exists an additive matrix D' such that $|D - D'|_{\infty} \leq 0.5$ then neighbor joining applied to D reconstructs the unique tree T (modulo isomorphism) such that $d'_{i,j} = d_T(i,j)$ for all $(i,j) \in n^2$.

Is this normalizer.

Atteson 1991

Neighbor Joining in Practice



Neighbor Joining tree relating copy number profiles from single cells in a tumor.

Outline

- Recap additive distance
- Neighbor joining
- Character-based phylogeny (small)
- Application to cancer

Reading:

• Chapters 10.2, 10.5-10.8, 10.9 in Jones and Pevzner

Character-Based Tree Reconstruction

- Characters may be morphological features
 - Shape of beak {generalist, insect catching, ...}
 - Number of legs {2,3,4, ..}
 - Hibernation {yes, no}
- Character may be nucleotides/amino acids
 - {A, T, C, G}
 - 20 amino acids
- Values of a character are called states
 - We assume discrete states



Character-Based Phylogeny Reconstruction



Question: What is optimal?

Want: Optimization criterion



Character-Based Phylogeny Reconstruction



Question: What is optimal?

Want: Optimization criterion

Question: How to optimize this criterion?

Want: Algorithm



Character-Based Phylogeny Reconstruction: Input

Characters / states	State 1	State 2
Mouth M	Smile S	Frown F
Eyebrows E	Normal 从	Pointed P



Character-Based Phylogeny Reconstruction: Criterion



Question: Which tree is better?

Character-Based Phylogeny Reconstruction: Criterion



(a) Parsimony Score=3

(b) Parsimony Score=2

Parsimony: minimize number of changes on edges of tree

Why Parsimony?

- Ockham's razor: "simplest" explanation for data
- Assumes that observed character differences resulted from the fewest possible mutations
- Seeks tree with the lowest **parsimony score**, i.e. the sum of all (costs of) mutations in the tree.



Again, a Small and a Large Problem

Small Maximum Parsimony Phylogeny Problem: Given $m \times n$ matrix $A = [a_{i,j}]$ and tree T with m leaves, find assignment of character states to each internal vertex of Twith minimum parsimony score.

Large Additive Distance Phylogeny Problem: Given $m \times n$ matrix $A = [a_{i,j}]$, find a tree T with m leaves labeled according to A and an assignment of character states to each internal vertex of T with minimum parsimony score.

Question: Are both problems easy (i.e. in P)?

Again, a Small and a Large Problem

Small Maximum Parsimony Phylogeny Problem: Given $m \times n$ matrix $A = [a_{i,j}]$ and tree T with m leaves, find assignment of character states to each internal vertex of Twith minimum parsimony score.

Large Additive Distance Phylogeny Problem:

Given $m \times n$ matrix $A = [a_{i,j}]$, find a tree T with m leaves labeled according to A and an assignment of character states to each internal vertex of T with minimum parsimony score.

Question: Are both problems easy (i.e. in P)?

Small Maximum Parsimony Phylogeny Problem ACCC ACC ACCA ACCG ACČA ATCC ATCG ACCG ATCG ATCC More Less Parsimonious Parsimonious Score: 6 Score: 5

Question: There are n = 4 characters in the m = 2 taxa (leaves). Can we solve each character separately?

Recurrence
$$\leq$$
 set of states
Given the T with $|L(T)|$ betwees
 $G: L(T) \rightarrow \geq$
 $V \in S(v) = \{x, y\}$
 $S(v)$ is the set of children of v
 $v \in Y$
 $V = S(v) = \{x, y\}$
 $S(v)$ is the set of children of v
 $v \in Y$
 $V = S(v)$
 $V = S$

Solving the Recurrence



$$\mu(v,s) = \min \left\{ \begin{array}{l} 0, & \text{if } v \in 2(T) \text{ and } s \neq \sigma(v), \\ 0, & \text{if } v \in 2(T) \text{ and } s = \sigma(v), \\ \end{array} \right.$$

$$E \left(\begin{array}{c} 0, & \text{if } v \in 2(T) \text{ and } s = \sigma(v), \\ \end{array} \right) \left(\begin{array}{c} 1 + 2^{2H} \\ 1 + 2^{2H} \\ \end{array} \right) \left(\begin{array}{c} 1 + 2^{2H} \\ 1 & \text{if } s = f, \\ 1 & \text{if } s \neq f, \end{array} \right) \left(\begin{array}{c} 1 + 2^{2H} \\ \end{array} \right) \left(\begin{array}{c$$

 $|L(\bar{1})| = M$ Example Bendocode Filling out M Fill($T, r(T), \delta, \xi$) 1Sach +race (T, v, M) $Fill(T, v, \sigma, \Sigma) \qquad O(m|\Sigma|^2)$ $i \in v = v(\tau)$ $\mathcal{B}(r(\tau)) = \arg\min\left\{\mathcal{M}(r(\tau), s)\right\}$ ses SEZ $_{-6(u)}$ closeter a bethe provent of v and let s be the state G(v) = avg min $Sc(s,t) + \mu(v,t)g^{of u}$ $t \in \Sigma$ if VEL(T) than For se E iF S = G(v) then $\mu(v,s) = 0$ $\mu(v,s) = 0$ For w E S(V) Bachtraie (T, w, M) alse For $w \in \delta(v)$ // children FII(T, w, 6, 2) $\mathcal{M}(u, 5) = 0$ $\mu(v,s) \neq = \min \left\{ \sum_{t \in S} \sum_{t \in S} c(s,t) + \mu(w,t)^2 \right\}$ For w E S(U)

Sankoff Algorithm (Sankoff 1975)

Small Maximum Parsimony Phylogeny Problem: Given $m \times n$ matrix $A = [a_{i,j}]$ and tree T with m leaves, find assignment of character states to each internal vertex of Twith minimum parsimony score.



Outline

- Recap additive distance
- Neighbor joining
- Character-based phylogeny (small)
- Application to cancer

Reading:

• Chapters 10.2, 10.5-10.8, 10.9 in Jones and Pevzner

Tumorigenesis: (i) Cell Division, (ii) Mutation & (iii) Migration



Tumorigenesis: (i) Cell Division, (ii) Mutation & (iii) Migration



Tumorigenesis: (i) Cell Division, (ii) Mutation & (iii) Migration



Goal: Given phylogenetic tree *T*, find *parsimonious* vertex labeling *e* with fewest migrations

Slatkin, M. and Maddison, W. P. (1989). A cladistic measure of gene flow inferred from the phylogenies of alleles. *Genetics*, 123(3), 603–613.

Minimum Migration Analysis in Ovarian Cancer

McPherson et al. (2016). Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nature Genetics*.

• Instance of the maximum parsimony small phylogeny problem [Fitch, 1971; Sankoff, 1975]



Minimum Migration Analysis in Ovarian Cancer

McPherson et al. (2016). Divergent modes of clonal spread and intraperitoneal mixing in high-grade serous ovarian cancer. *Nature Genetics*.

• Instance of the maximum parsimony small phylogeny problem [Fitch, 1971; Sankoff, 1975]



Minimum Migration History is Not Unique

• Enumerate all minimum-migration vertex labelings in the backtrace step



Comigrations: Simultaneous Migrations of Multiple Clones

- Multiple tumor cells migrate simultaneously through the blood stream [Cheung et al., 2016]
- Second objective: number γ of comigrations is the number of multi-edges in migration graph G^+





АрС	Appendix
LFTB	Left Fallopian Tube
LOv	Left Ovary
RFTA	Right Fallopian Tube
ROv	Right Ovary
SBwl	Small Bowel
Om	Omentum

+ Not necessarily true in the case of directed cycles

Comigrations: Simultaneous Migrations of Multiple Clones

- Multiple tumor cells migrate simultaneously through the blood stream [Cheung et al., 2016]
- Second objective: number γ of comigrations is the number of multi-edges in migration graph G^+



Constrained Multi-objective Optimization Problem

Parsimonious Migration History (PMH): Given a phylogenetic tree T and a set $\mathcal{P} \subseteq \{S, M, R\}$ of allowed migration patterns, find vertex labeling ℓ with minimum migration number $\mu^*(T)$ and smallest comigration number $\hat{\gamma}(T)$.



El-Kebir, M., Satas, G., & Raphael, B. J. (2018). Inferring parsimonious migration histories for metastatic cancers. *Nature Genetics*, 50(5), 718–726.

Results [El-Kebir, WABI 2018]

Parsimonious Migration History (PMH): Given a phylogenetic tree T and a set $\mathcal{P} \subseteq \{S, M, R\}$ of allowed migration patterns, find vertex labeling ℓ with minimum migration number $\mu^*(T)$ and smallest comigration number $\hat{\gamma}(T)$.



PMH is NP-hard when $\mathcal{P} = \{S\}$

3-SAT: Given $\varphi = \bigwedge_{i=1}^{k} (y_{i,1} \lor y_{i,2} \lor y_{i,3})$ with variables $\{x_1, \dots, x_n\}$ and k clauses, find $\varphi : [n] \rightarrow \{0,1\}$ satisfying φ



PMH is NP-hard when $\mathcal{P} = \{S\}$

3-SAT: Given $\varphi = \bigwedge_{i=1}^{k} (y_{i,1} \lor y_{i,2} \lor y_{i,3})$ with variables $\{x_1, \dots, x_n\}$ and k clauses, find $\varphi : [n] \rightarrow \{0,1\}$ satisfying φ

Three ideas:

- 1. Ensure that $(x, \neg x) \in E(G)$ or $(\neg x, x) \in E(G)$
- 2. Ensure that $\ell^*(r(T)) = \bot$
- 3. Ensure that ϕ is satisfiable if and only if ℓ^* encodes a satisfying truth assignment





PMH is NP-hard when $\mathcal{P} = \{S\}$

3-SAT: Given $\varphi = \bigwedge_{i=1}^{k} (y_{i,1} \lor y_{i,2} \lor y_{i,3})$ with variables $\{x_1, \dots, x_n\}$ and k clauses, find $\varphi : [n] \rightarrow \{0,1\}$ satisfying φ

Three ideas:

- 1. Ensure that $(x, \neg x) \in E(G)$ or $(\neg x, x) \in E(G)$
- 2. Ensure that $\ell^*(r(T)) = \bot$
- 3. Ensure that ϕ is satisfiable if and only if ℓ^* encodes a satisfying truth assignment





Lemma: Let B > 10k + 1 and A > 2Bn + 27k. Then, φ is satisfiable if and only if $\mu^*(T) = (B + 1)n + 25k$



Lemma: Let B > 10k + 1 and A > 2Bn + 27k. Then, φ is satisfiable if and only if $\mu^*(T) = (B + 1)n + 25k$

PMH is FPT in number m of locations when $\mathcal{P} = \{S\}$



Lemma: If (1) holds then ℓ^* is a minimum migration labeling consistent with \widehat{G} .



Lemma: If (1) holds then ℓ^* is a minimum migration labeling consistent with \widehat{G} .

Simulations



Available on: https://github.com/elkebir-group/PMH-S

Outline

- Recap additive distance
- Neighbor joining
- Character-based phylogeny (small)
- Application to cancer

Reading:

• Chapters 10.2, 10.5-10.8, 10.9 in Jones and Pevzner