

Spring 2024
BIOS 477/877
Bioinformatics and Molecular Evolution
Lecture 23

BIOS477/877 L.23 - 1

1

TODAY'S TOPICS

- Phylogenetic reconstruction
 - Maximum parsimony method: use examples
 - Maximum likelihood method
 - Tree searching

(Exhaustive, Branch-and-Bound, Heuristic)

- Branch support (bootstrap analysis, etc.)

BIOS477/877 L.23 - 2

2

Maximum Parsimony (MP): examples

BIOS477/877 L.23 - 3

3

Maximum Parsimony (MP): examples

Inference of single-cell phylogenies from lineage tracing data using *Cassiopeia*

Matthew G. Jones^{1,2,3,4,5}, Alex Rhoades^{6,7}, Jeffrey J. Quast^{8,9,10}, Michelle M. Chan¹¹, Jeffrey A. Hausman^{12,13}, Robert Wang¹⁴, Chering Tu¹⁵, Jonathan S. Weisman^{16,17} and Nir Yossef^{18,19}

Jones et al. (2020) *Genome Biol.* 21: 92

BIOS477/877 L.23 - 4

4

Phylogenetic methods

➤ Data types and tree-building methods

		[Data types]	
		Distances	Characters
[Tree-building methods]	Clustering	UPGMA Neighbor joining	
	Optimality criterion	Minimum evolution Fitch-Margoliash	Maximum parsimony Maximum likelihood (Bayesian inference)

Examine all possible topologies based on a certain criterion

Yang and Rannala (2012) *Molecular phylogenetics: principles and practice. Nature Reviews Genetics* 13: 303-314. BIOS477/877 L.23 - 5

5

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**

- Chooses the tree that makes the observed data the most probable evolutionary outcome
- Likelihood = Conditional probability of obtaining the observed sequences given a hypothesis (substitution model and tree)

$$L(\tau, \theta) = \text{Prob}(\text{Data} \mid \tau, \theta) = \text{Prob}(\text{Aligned sequences} \mid \text{tree, model of evolution})$$

τ : a tree (including its topology, branch lengths)
 θ : a set of parameters for a substitution model

➔ The likelihood calculated is not the probability of tree or evolutionary model! It is the probability of the data.

BIOS477/877 L.23 - 6

6

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**
 $L(\tau, \theta) = \text{Prob}(\text{Data} | \tau, \theta)$

Find the topology that gives the maximum $L(\tau, \theta)$, and simultaneously estimate all required parameters

- To compute the likelihood of a given tree,
 - τ : the topology and the maximum likelihood estimates for the tree's branch lengths (d_1, d_2, \dots) need to be found
 - θ : the best values for the parameters for the evolutionary model need to be found, too

$$\begin{cases} r_i = \frac{1}{4} + \frac{3}{4} e^{-4\alpha} \\ s_i = \frac{1}{4} - \frac{1}{4} e^{-4\alpha} \end{cases}$$

BIOS477/877 L23 - 7

7

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**

Data

1 ATATT
2 ATCGT
3 GCAGT
4 GCCGT

Tree1

$L(\text{AAGG} | \text{Tree}) = \text{Prob} \left(\begin{matrix} A & G \\ A & G \end{matrix} \right)$

All possible combinations

BIOS477/877 L23 - 8

8

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**

Data

1 ATATT
2 ATCGT
3 GCAGT
4 GCCGT

Tree1

$L(\text{AAGG} | \text{Tree}) = \text{Prob} \left(\begin{matrix} A & G \\ A & G \end{matrix} \right) + \text{Prob} \left(\begin{matrix} A & G \\ G & G \end{matrix} \right) + \dots$

4x4 combinations

BIOS477/877 L23 - 9

9

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**

Data

1 ATATT
2 ATCGT
3 GCAGT
4 GCCGT

Tree1

$L(\text{AAGG} | \text{Tree}) = \text{Prob} \left(\begin{matrix} A & G \\ A & G \end{matrix} \right) + \dots$

$P_{AA}(d_1) \times P_{AA}(d_2) \times P_{AA}(d_3) \times P_{AG}(d_4) \times P_{AG}(d_4)$
 $P_{ij}(d_i)$: the probability that two sequences separated by d_i would have character states i and j

Transition probability matrix for JC model (see Lec. 20 & "Derivation of JC equation")

$$P(t) = \begin{bmatrix} r_1 & r_2 & r_3 & r_4 \\ r_2 & r_1 & r_3 & r_4 \\ r_3 & r_3 & r_1 & r_4 \\ r_4 & r_4 & r_4 & r_1 \end{bmatrix} \begin{cases} r_1 = \frac{1}{4} + \frac{3}{4} e^{-4\alpha t} \\ r_2 = \frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \\ r_3 = \frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \\ r_4 = \frac{1}{4} - \frac{1}{4} e^{-4\alpha t} \end{cases}$$

$d_i = \alpha t_i$, and $P_{ij}(d_i) = \begin{cases} r_{ii} & \text{if } i=j \\ s_{ij} & \text{if } i \neq j \end{cases}$

BIOS477/877 L23 - 10

10

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**

Data

1 ATATT
2 ATCGT
3 GCAGT
4 GCCGT

Tree1

$L_{\text{Tree1}} = L(\text{AAGG} | \text{Tree1}) \times L(\text{TTCC} | \text{Tree1}) \times L(\text{ACAC} | \text{Tree1}) \times L(\text{TGGG} | \text{Tree1}) \times L(\text{TTTT} | \text{Tree1})$

$= \prod_i L_i$ where $L_i = L(\text{site}_i | \text{Tree1})$

$\ln(L_{\text{Tree1}}) = \ln[L(\text{AAGG} | \text{Tree1})] + \ln[L(\text{TTCC} | \text{Tree1})] + \ln[L(\text{ACAC} | \text{Tree1})] + \ln[L(\text{TGGG} | \text{Tree1})] + \ln[L(\text{TTTT} | \text{Tree1})]$

$= \sum_i \ln[L(\text{site}_i | \text{Tree1})]$

Likelihood of the data (multiple alignment) given Tree1

BIOS477/877 L23 - 11

11

Phylogenetic methods (Character-based)

➤ **Maximum Likelihood (ML)**
 $L(\tau, \theta) = \text{Prob}(\text{Data} | \tau, \theta)$

- To compute the likelihood of a given tree,
 - τ : the maximum likelihood estimates for the tree's branch lengths (d_1, d_2, \dots) need to be found
 - θ : the best values for the parameters for the evolutionary model need to be found, too

→ [JC model] α (single evolutionary rate)

→ [2-parameter model] transition/transversion ratio (TS/TV)

base composition, gamma shape parameter (α)

etc. etc.

→ Number of the parameters depends on the model

$$\begin{cases} r_1 = \frac{1}{4} + \frac{3}{4} e^{-4\alpha} \\ r_2 = \frac{1}{4} - \frac{1}{4} e^{-4\alpha} \\ r_3 = \frac{1}{4} - \frac{1}{4} e^{-4\alpha} \\ r_4 = \frac{1}{4} - \frac{1}{4} e^{-4\alpha} \end{cases}$$

Find the topology that gives the maximum $L(\tau, \theta)$, and simultaneously estimate all required parameters

BIOS477/877 L23 - 12

12

Phylogenetic methods

	Distances		Characters
Clustering	UPGMA		
	NJ		
Optimality criterion	ME	MP	
	EM	ML	

➤ **Clustering vs. search methods**

- **Clustering methods (UPGMA, Neighbor-joining)**
 - Do not search all possible topologies
 - Very fast
 - Produce only one tree
- **Search methods**
 - Use optimality criterion (minimum evolution, maximum parsimony, maximum likelihood)
 - Exhaustive search for all possible topologies is not possible for a large number of taxa
 - A heuristic search algorithm needs to be used

How can we search trees?

BIOS477/877 L23 - 13

13

Phylogenetic methods: tree searching

➤ **Number of possible tree topologies**

Number of OTUs	Number of rooted trees	Number of unrooted trees
2	1	1

BIOS477/877 L23 - 14

14

Phylogenetic methods: tree searching

➤ **Number of possible tree topologies**

Number of OTUs	Number of rooted trees	Number of unrooted trees
2	1	1
3	3	1

BIOS477/877 L23 - 15

15

Phylogenetic methods: tree searching

➤ **Number of possible tree topologies**

Number of OTUs	Number of rooted trees	Number of unrooted trees
2	1	1
3	3	1
4	15	3

5 rooted trees 5 rooted trees 5 rooted trees

BIOS477/877 L23 - 16

16

Phylogenetic methods: tree searching

➤ **Number of possible tree topologies**

Number of OTUs	Number of rooted trees	Number of unrooted trees
2	1	1
3	3	1
4	15	3
5	105	15
6	954	105
7	10,395	954
8	135,135	10,395
9	2,027,025	135,135
10	34,459,425	2,027,025

Impossible to examine all of the possible trees!

BIOS477/877 L23 - 17

17

Tree-searching methods (unrooted)

➤ **Exhaustive search**

- Possible only for a few taxa (11 or fewer)
- Number of possible unrooted tree:

$$B(t) = \prod_{i=3}^t (2i - 5) \quad t : \text{number of taxa}$$

e.g., $B(7) = 1 \times 3 \times 5 \times 7 \times 9 = 945$
 $B(10) = B(7) \times 11 \times 13 \times 15 > 2 \times 10^6$
 $B(20) > 2 \times 10^{20}$
- An algorithm is required to guarantee generation of all possible trees

BIOS477/877 L23 - 18

18

Tree-searching methods (unrooted)

➤ **Exhaustive search**

If there are 2 taxa → Only one topology possible

There is only one place to add the third taxon

If there are 3 taxa There are three places to add the fourth taxon
Three topologies are possible

BIOS477/877 L.23 - 19

19

Tree-searching methods (unrooted)

➤ **Exhaustive search**

If there are 4 taxa There are five places to add the fifth taxon
Five topologies are possible

BIOS477/877 L.23 - 20

20

Tree-searching methods (unrooted)

➤ **Exhaustive search for 5-taxon trees**

3 ways to add another branch 5 ways to add another branch

5 ways to add another branch 5 ways to add another branch

5 ways to add another branch 5 ways to add another branch

5 ways to add another branch 5 ways to add another branch

Each of the possible topologies (3 x 5 = 15 trees for 5-taxon tree) will be evaluated to identify optimal trees (ME, MP, or ML trees)

BIOS477/877 L.23 - 21

21

Tree-searching methods (unrooted)

➤ **Exhaustive search**

- searches all possible trees
- guarantees to find the optimal trees
- Impractical for many taxa

BIOS477/877 L.23 - 22

22

Tree-searching methods (unrooted)

➤ **Branch-and-bound method**

- Alternative exact method
- Useful for 25 or so taxa
- Implicitly evaluates all possible trees
- But cuts off search paths that do not lead to optimal trees
 - Reduces the number of trees to evaluate
 - Will find the optimal tree(s)

BIOS477/877 L.23 - 23

23

Tree-searching methods (unrooted)

➤ **Branch-and-bound method**

[To search 6-taxon MP tree]

(7 6-taxon trees)

241 becomes the upper bound ($L_u=241$) for searching 6-taxon trees (MP example: searching the shortest tree length, L)

BIOS477/877 L.23 - 24

24

Tree-searching methods (unrooted)

➤ **Branch and bound search**

- searches only the paths under the upper bound (L_u)
- guarantees to find the optimal trees

BIOS477/877 L.23 - 31

31

Tree-searching methods (unrooted)

➤ **Heuristic (approximation) method**

- Stepwise addition: greedy algorithm

BIOS477/877 L.23 - 32

32

Tree-searching methods (unrooted)

➤ **Heuristic (approximation) method**

- Stepwise addition: greedy algorithm

BIOS477/877 L.23 - 33

33

Tree-searching methods (unrooted)

➤ **Heuristic (approximation) method**

- Stepwise addition: greedy algorithm
 - Rarely finds the global optima
 - Can be improved (refinement)
- Branch-swapping
 - Cutting off one or more subtrees and reassembling them to generate locally different trees
 - Repeat this process many times to find better topologies
 - Still can be entrapped in local optima
- Random addition + branch-swapping (many times)
 - Multiple optimal islands can be identified

BIOS477/877 L.23 - 34

34

Tree-searching landscape

BIOS477/877 L.23 - 35

35

Tree-searching landscape

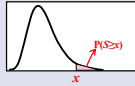
BIOS477/877 L.23 - 36

36

Reliability of inferred trees

➤ Bootstrap analysis

- Evaluates the reliability of each cluster
- A sampling technique to estimate the statistical errors when underlying sampling distribution is unknown



→ Approximates the underlying distribution by resampling from the original dataset

- First applied by Felsenstein (1985) for phylogenetic analysis

BIOS477/877 L.23 - 37

37

Reliability of inferred trees

➤ Bootstrap analysis

→ Resampling is done on sequence data

	1	2	3	4	5	6	7	8	Site #
S1	G	C	A	G	T	A	C	T	
S2	G	T	A	G	T	A	C	T	
S3	A	C	A	A	T	A	C	C	
S4	A	C	A	A	C	A	C	T	
S5	G	C	G	G	C	A	T	T	

Multiple alignment

BIOS477/877 L.23 - 38

38

Reliability of inferred trees

➤ Bootstrap analysis

→ Resampling is done on sequence data

	1	2	3	4	5	6	7	8	
8 independent samples	G	C	A	G	T	A	C	T	
	G	T	A	G	T	A	C	T	
	A	C	A	A	T	A	C	C	
	A	C	A	A	C	A	C	T	
	G	C	G	G	C	A	T	T	

Configuration for each site is not changed

Each column of the multiple alignment is treated as an independent sample

BIOS477/877 L.23 - 39

39

Reliability of inferred trees

➤ Bootstrap analysis

1	2	3	4	5	6	7	8
G	C	A	G	T	A	C	T
G	T	A	G	T	A	C	T
A	C	A	A	T	A	C	C
A	C	A	A	C	A	C	T
G	C	G	G	C	A	T	T

1
S1 G
S2 G
S3 A
S4 A
S5 G

Columns are randomly sampled

BIOS477/877 L.23 - 40

40

Reliability of inferred trees

➤ Bootstrap analysis

1	2	3	4	5	6	7	8
G	C	A	G	T	A	C	T
G	T	A	G	T	A	C	T
A	C	A	A	T	A	C	C
A	C	A	A	C	A	C	T
G	C	G	G	C	A	T	T

14
S1 GG
S2 GG
S3 AA
S4 AA
S5 GG

Columns are randomly sampled

BIOS477/877 L.23 - 41

41

Reliability of inferred trees

➤ Bootstrap analysis

1	2	3	4	5	6	7	8
G	C	A	G	T	A	C	T
G	T	A	G	T	A	C	T
A	C	A	A	T	A	C	C
A	C	A	A	C	A	C	T
G	C	G	G	C	A	T	T

146
S1 GGA
S2 GGA
S3 AAA
S4 AAA
S5 GGA

Columns are randomly sampled

BIOS477/877 L.23 - 42

42

Reliability of inferred trees

➤ **Bootstrap analysis**

1	2	3	4	5	6	7	8
G	C	A	G	T	A	C	T
G	T	A	A	T	A	C	T
A	C	A	A	T	A	C	T
A	C	A	A	C	A	C	T
G	C	G	G	C	A	T	T

1461
S1 GGAG
S2 GGAG
S3 AAAA
S4 AAAA
S5 GGAG

Multiple samplings are allowed from the same column

BIOS477/877 L.23 - 43

43

Reliability of inferred trees

➤ **Bootstrap analysis**

1	2	3	4	5	6	7	8
G	C	A	G	T	A	C	T
G	T	A	A	T	A	C	T
A	C	A	A	T	A	C	T
A	C	A	A	C	A	C	T
G	C	G	G	C	A	T	T

8 samples (sites)

A bootstrap replicate with the same length (the same number of sites) as the original alignment

14614853
S1 GGAGGTTA
S2 GGAGGTTA
S3 AAAAATA
S4 AAAAATCA
S5 GGAGGTCG

8 samples (sites)

BIOS477/877 L.23 - 44

44

Reliability of inferred trees

➤ **Bootstrap analysis**

1	2	3	4	5	6	7	8
G	C	A	G	T	A	C	T
G	T	A	A	T	A	C	T
A	C	A	A	T	A	C	T
A	C	A	A	C	A	C	T
G	C	G	G	C	A	T	T

14614853
S1 GGAGGTTA
S2 GGAGGTTA
S3 AAAAATA
S4 AAAAATCA
S5 GGAGGTCG

74761232
S1 CGCAGCAC
S2 CGCAGTAT
S3 CACAACAC
S4 CACAACAC
S5 TGTAGCCG

85851124
S1 TTTTGGCG
S2 TTTTGGTG
S3 CTCTAACA
S4 TCTCAACA
S5 TCTCGCGG

Many resamplings are done (~ 1000 times)

BIOS477/877 L.23 - 45

45

Reliability of inferred trees

➤ **Bootstrap analysis**

14614853
S1 GGAGGTTA
S2 GGAGGTTA
S3 AAAAATA
S4 AAAAATCA
S5 GGAGGTCG

74761232
S1 CGCAGCAC
S2 CGCAGTAT
S3 CACAACAC
S4 CACAACAC
S5 TGTAGCCG

85851124
S1 TTTTGGCG
S2 TTTTGGTG
S3 CTCTAACA
S4 TCTCAACA
S5 TCTCGCGG

... (~1000 alignments)

... (~1000 trees)

A phylogeny is reconstructed from each pseudoreplicate

BIOS477/877 L.23 - 46

46

Reliability of inferred trees

➤ **Bootstrap analysis**

... (~1000 trees)

(S1, S2), (S3, S4, S5)

(S1, S2), (S3, S4, S5)

(S1, S2), (S3, S4, S5)

Count the number of each cluster

(S1, S2), (S3, S4, S5): 3/3 = 100%

BIOS477/877 L.23 - 47

47

Reliability of inferred trees

➤ **Bootstrap analysis**

... (~1000 trees)

(S1, S2), (S3, S4, S5)

(S1, S2), (S3, S4, S5)

(S1, S2), (S3, S4, S5)

(S1, S2, S3), (S4, S5)

(S3, S4), (S1, S2, S5)

(S1, S2, S3), (S4, S5)

Count the number of each cluster

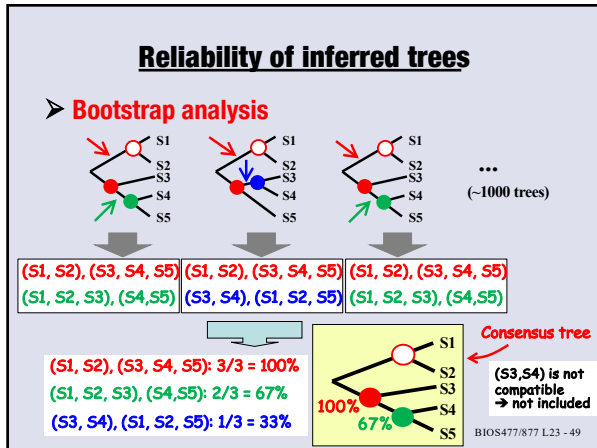
(S1, S2), (S3, S4, S5): 3/3 = 100%

(S1, S2, S3), (S4, S5): 2/3 = 67%

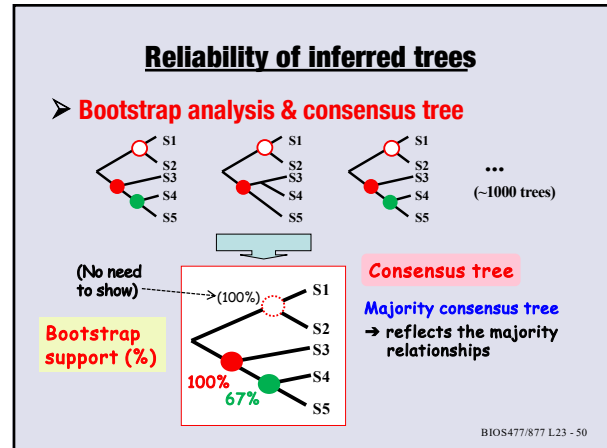
(S3, S4), (S1, S2, S5): 1/3 = 33%

BIOS477/877 L.23 - 48

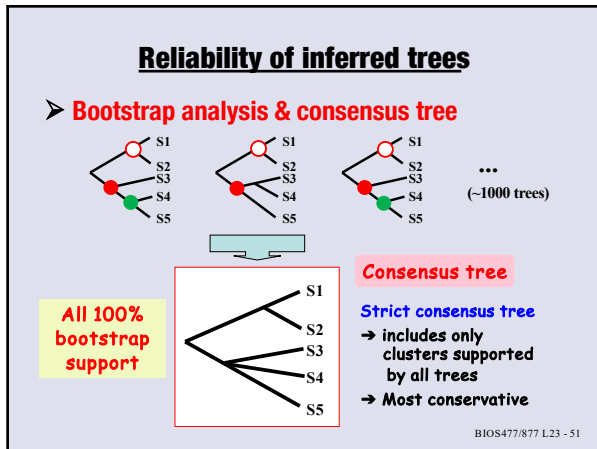
48



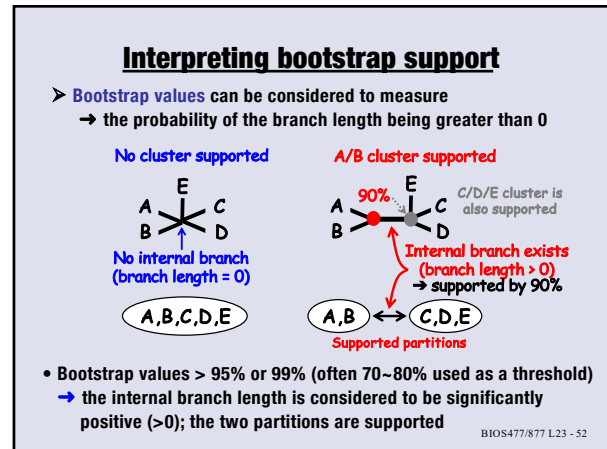
49



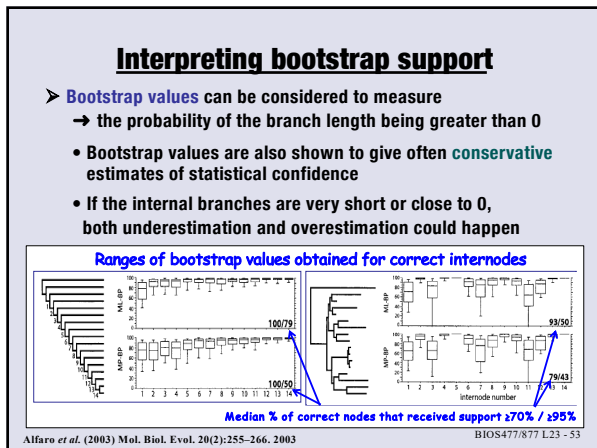
50



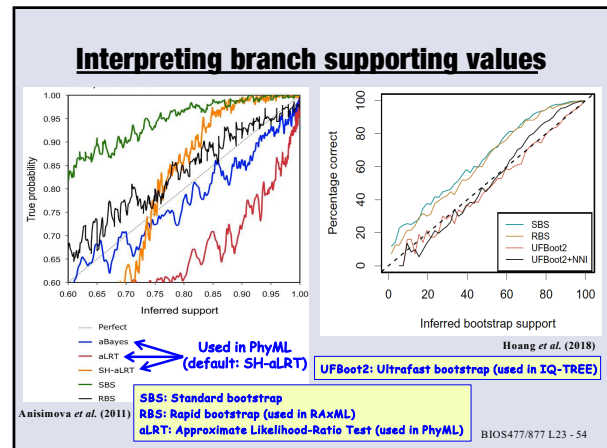
51



52



53



54

Renewing Felsenstein's phylogenetic bootstrap in the era of big data

F. Lemoine^{1,2}, J.-B. Domelevo Entfellner^{3,4}, E. Wilkinson⁵, D. Correia¹, M. Dávila Felipe^{1,6}, T. De Oliveira^{5,6} & O. Gascuel^{1,2*}

452 | NATURE | VOL 556 | 26 APRIL 2018

Transfer distance (or R distance):
The minimum number of elements to be transferred (or removed) to transform on partition into the other

(S1, S2), (S3, S4, S5) $\xleftrightarrow{\text{Distance}=0}$ (S1, S2), (S3, S4, S5)
 (S1, S2, S3), (S4, S5) $\xleftrightarrow{\text{Distance}=1}$ (S1, S2), (S3, S4, S5)

Transfer Bootstrap Expectation (TBE) for branch b :

$$TBE(b) = 1 - \frac{\phi(b, T^*)}{p-1}$$
 Average from all BS trees
 p : number of taxa from the smaller of the two clusters

BS: count only identical partitions
TBE: include most similar partitions (no need to be identical)

$R(b, b^*)$: Distance between branch b in tree T and branch b^* in BS tree T^*
 - Compare the two partitions defined by branch b in tree T and all branches in BS tree T^*
 - Find the b^* with the minimum $R(b, b^*) \rightarrow \phi(b, T^*)$
 $R(b, b^*)=1, R(b, b^*)=2$

BIOS477/877 L.23 - 55

55

Transfer bootstrap for big phylogenies

Renewing Felsenstein's phylogenetic bootstrap in the era of big data

A fast and memory-efficient implementation of the transfer bootstrap
 Sarah Lutteropp^{1*}, Alexey M. Kozlov^{2*} and Alexandros Stamatakis^{1,3}

Standard bootstrap proportion (SBP) > 70% **Transfer bootstrap expectation (TBE) > 70%**

- Deep branches in large phylogenies are often not supported by SBP
- TBE supports are higher without inducing falsely supported branches

BIOS477/877 L.23 - 56

56