

Spring 2024
BIOS 477/877
Bioinformatics and Molecular Evolution
Lecture 21

BIOS477/877 L21 - 1

1

TODAY'S TOPICS

- Distance estimation
 - Synonymous & nonsynonymous distances
 - Phylogenetic reconstruction
- Introduction (terminologies, rooting, etc.)
 - Distance methods (UPGMA, FM, ME)

BIOS477/877 L21 - 2

2

Universal Genetic Code

		T	C	A	G
T	T	TTT Phe	TCT Ser	TAT Tyr	TGT Cys
	C	TTC Phe	TCC Ser	TAC Tyr	TGC Cys
	A	TTA Stop	TCA Ser	TAA Stop	TGA Stop
	G	TTG Trp	TCG Ser	TAG Stop	TGG Trp
C	T	CTT Leu	CCT Pro	CAT His	CGT Arg
	C	CTC Leu	CCC Pro	CAC His	CGC Arg
	A	CTA Leu	CCA Pro	CAA Gln	CGA Arg
	G	CTG Leu	CCG Pro	CAG Gln	CGG Arg
A	T	ATT Ile	ACT Thr	AAT Asn	AGT Ser
	C	ATC Ile	ACC Thr	AAC Asn	AGC Ser
	A	ATA Ile	ACA Thr	AAA Lys	AGA Arg
	G	ATG Met	ACG Thr	AAG Lys	AGG Arg
G	T	GTT Val	GCT Ala	GAT Asp	GGT Gly
	C	GTC Val	GCC Ala	GAC Asp	GGC Gly
	A	GTA Val	GCA Ala	GAA Glu	GGA Gly
	G	GTG Val	GCG Ala	GAG Glu	GGG Gly

• **Synonymous (silent)** substitutions DOES NOT change amino acids
 • **Nonsynonymous (replacement)** substitutions DOES change amino acids

BIOS477/877 L21 - 3

3

Synonymous/nonsynonymous distance methods

- **Nei-Gojobori method** (Nei and Gojobori, 1986)
 - Number of synonymous differences: S_d
 - Number of nonsynonymous differences: N_d
 - Proportion of synonymous differences: p_S
 - Proportion of nonsynonymous differences: p_N

$\rightarrow p_S = S_d/S, p_N = N_d/N$
 S : Number of synonymous sites
 N : Number of nonsynonymous sites

- Jukes-Cantor correction for multiple-hits
 - $\rightarrow d_S = -3/4 \ln(1 - 4p_S/3)$
 - $\rightarrow d_N = -3/4 \ln(1 - 4p_N/3)$

K2P or Tajima-Nei (1-parameter+base freq.) correction is also used in modified versions
 BIOS477/877 L21 - 4

4

Synonymous/nonsynonymous distance methods

➤ How to count synonymous/nonsynonymous differences

(Val) (Phe) (Leu)

GTT	TTT	TTG
GTA	GTA	AGA

(Val) (Val) (Arg)

One synonymous substitution (blue arrow) Six possible pathways!

Two possible pathways (red arrows):

1/2 Nonsynonymous: TTT (Phe) ↔ GTA (Val)

1/2 Synonymous: TTT (Phe) ↔ TTA (Leu) ↔ GTA (Val)

1/2 + 1/2 Nonsynonymous

[Total: 1.5 Nonsynonymous + 0.5 Synonymous]

BIOS477/877 L21 - 5

5

Synonymous/nonsynonymous distance methods

➤ How to count synonymous/nonsynonymous sites

(Phe)

TTT

ATT	TAT	TTA
CTT	TCT	TTC ← synonymous
GTT	TGT	TTG

Synonymous sites (S): 0 + 0 + 1/3 = 1/3
 Nonsynonymous sites (N): 3/3 + 3/3 + 2/3 = 8/3

- Count the number of sites from each codon and sum up for each sequence. Take the average from two sequences.

BIOS477/877 L21 - 6

6

Synonymous/nonsynonymous distance methods

➤ Nei-Gojobori method (Nei and Gojobori, 1986)

- Number of synonymous differences: S_d
- Number of nonsynonymous differences: N_d
- Proportion of synonymous differences: p_S
- Proportion of nonsynonymous differences: p_N

$$\rightarrow p_S = S_d/S, p_N = N_d/N$$

S : Number of synonymous sites

N : Number of nonsynonymous sites

- Jukes-Cantor correction for multiple-hits

$$\rightarrow d_S = -3/4 \ln(1 - 4p_S/3)$$

$$\rightarrow d_N = -3/4 \ln(1 - 4p_N/3)$$

K2P or Tajima-Nei (1-parameter+base freq.) correction is also used in modified versions

BIOS477/877 L21 - 7

7

Nucleotide substitution patterns

Table 2. Comparisons among the methods for estimating synonymous and nonsynonymous substitution numbers per site

Method	Gene: length (species compared with <i>D. melanogaster</i>)			
	GC:60% Adh: 816 bp (<i>D. teissieri</i>)	GC:80% Adh: 762 bp (<i>D. ps. bogotana</i>)	GC:4% Col: 1497 bp (<i>D. yakuba</i>)	
	Synonymous	Nonsynonymous	Synonymous	Nonsynonymous
NG	0.402 ± 0.060	0.009 ± 0.004	0.604 ± 0.080	0.053 ± 0.010
LWL	0.394 ± 0.058	0.009 ± 0.004	0.599 ± 0.080	0.054 ± 0.010
PBL	0.328 ± 0.052	0.009 ± 0.004	0.561 ± 0.078	0.054 ± 0.010
			0.380 ± 0.041	0.007 ± 0.003
			0.364 ± 0.040	0.007 ± 0.002
			0.401 ± 0.051	0.007 ± 0.003

Moriyama and Powell (1997) J Mol Evol 45:378-391

NG: Nei-Gojobori method (Nei & Gojobori 1986); based on JC model

LWL: Li-Wu-Luo method (Li et al. 1985); based on K2P model

PBL or Li93: Pamilo-Bianchi-Li method (Pamilo and Bianchi 1993; Li 1993)

Kumar method (available in MEGA; modification to PBL)

NG method underestimates the number of synonymous sites: S

LWL method overestimates the number of synonymous sub.: S_d

PBL method corrected problems found in both NG and LWL methods

BIOS477/877 L21 - 8

8

Available distance method programs

- MEGA X <http://www.megasoftware.net/>
→ Includes synonymous & nonsynonymous distances
- PAML <http://abacus.gene.ucl.ac.uk/software/paml.html>
→ Includes Yang and Nielsen (2000) method [yn00]
- SNAP <https://www.hiv.lanl.gov/content/sequence/SNAP/SNAP.html>
→ Synonymous & nonsynonymous (Nei-Gojobori) distance only
- Ape (R package for Analysis of Phylogenetics and Evolution)
→ Includes many distance methods <https://emmanuelparadis.github.io/index.html>
<https://cran.r-project.org/web/packages/ape/index.html>
- Phylip3.698 <http://evolution.genetics.washington.edu/phylip.html>
→ JC, K2P, F84 (HKY85), LogDet, gamma distances
→ Dayhoff's PAM, JTT, PMB (Probability Matrix from Blocks), Kimura's PAM approximation, gamma distances
→ On the Web: <http://phylemon.bioinfo.cipf.es> (ver. 3.68)
→ In EMBOSS: <http://emboss.sourceforge.net/cgi-bin/emboss/> (found in Phylogeny sections)
→ See "How to use Phylip" on Canvas

! ClustalW2 (ClustalX2) → K2P for DNA, hybrid between Kimura and PAM for protein!

$p \leq 0.75$ Use Kimura's PAM distance approximation method
 $0.75 < p \leq 0.93$ Use a conversion table with 0.01 interval (.75, .751, ...)
 $0.93 < p$ $k = 10.0$ [arbitrary-constant]

BIOS477/877 L21 - 9

9

Introduction to phylogeny

➤ Phylogeny (phylogenetic tree)

→ a graphic representation of evolutionary relationships among genes or organisms

- True phylogeny cannot be known

We cannot actually observe the long-term evolution!

- Phylogenetic relationships can be only **inferred**
- Phylogenetic relationships are **reconstructed** based on the information available (*e.g.*, sequences)

→ represents a **hypothesis of evolutionary relationships** among gene or protein sequences: **gene tree**

→ Organismal relationships are inferred based on **phylogenetic analysis: species tree**

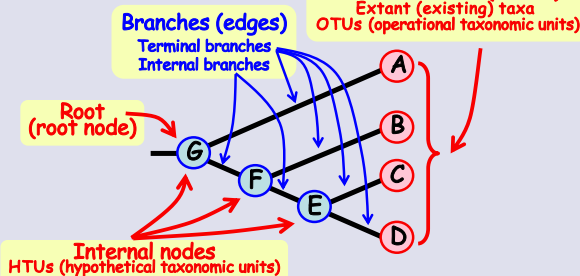
Note: Gene trees do not always represent species trees!

BIOS477/877 L21 - 10

10

Introduction to phylogeny

➤ Tree terminology

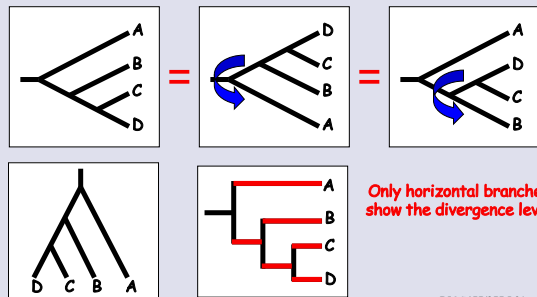


BIOS477/877 L21 - 11

11

Introduction to phylogeny

➤ Many ways of drawing trees



BIOS477/877 L21 - 12

12

Introduction to phylogeny

➤ Many ways of drawing trees

Only vertical branches show the divergence level

BIOS477/877 L21 - 13

13

Introduction to phylogeny

➤ Three different types of trees

Cladogram
Relative recency of common ancestry (or branching order)
No quantitative information

Additive tree (phenogram)
Branch lengths show the amount of evolutionary changes

Ultrametric tree
shows evolutionary time

In ultrametric trees, end nodes are all equidistant from the root of the tree
→ possible only assuming molecular clock (constant evolutionary rate)

BIOS477/877 L21 - 14

14

Introduction to phylogeny

➤ Three different types of trees

Cladogram
Branch length has no information

Additive tree (phenogram)
Ultrametric tree
Branch length shows the amount of divergence

BIOS477/877 L21 - 15

15

Introduction to phylogeny

➤ Resolution of trees

Star tree
No resolution

Partially resolved

Fully resolved (bifurcating tree)

BIOS477/877 L21 - 16

16

Introduction to phylogeny

➤ Nested parentheses format: **Newick format**

Rooted **Unrooted**

Rooted
((A,B), (C, (D, E)));
((A:2, B:1.5):2, C:3, (D:1, E:1):1);

Unrooted
((A,B), C, (D, E));
((A:2, B:1.5):2, C:3, (D:1, E:1):1);

Branch lengths

BIOS477/877 L21 - 17

17

Introduction to phylogeny

➤ Nested parentheses format: **Newick format**

Rooted **Unrooted**

Rooted
((A,B), (C, (D, E)));
2 clusters divided by the root

Unrooted
((A,B), C, (D, E));
3 clusters and no root

BIOS477/877 L21 - 18

18

Introduction to phylogeny

Rooted tree

Unrooted tree

BIOS477/877 L21 - 19

19

Introduction to phylogeny

Rooted tree

Unrooted tree

BIOS477/877 L21 - 20

20

Introduction to phylogeny

Rooted tree

Unrooted tree

BIOS477/877 L21 - 21

21

Introduction to phylogeny

Rooted tree

Unrooted tree

BIOS477/877 L21 - 22

22

Introduction to phylogeny

Rooted tree

Unrooted tree

BIOS477/877 L21 - 23

23

Introduction to phylogeny

Topology is the same!

Unrooted tree

BIOS477/877 L21 - 24

24

Introduction to phylogeny

➤ **Outgroup**

- used to "root" unrooted tree
- Biological information is required to choose appropriate outgroup

Unrooted tree

Root for the ingroup

BIOS477/877 L21 - 25

25

Introduction to phylogeny

➤ **Outgroup**

Which species can be used as the outgroup?

BIOS477/877 L21 - 26

26

Introduction to phylogeny

➤ **Outgroup**

- used to "root" unrooted tree
- Biological information is required to choose appropriate outgroup

Most appropriate outgroup?

BIOS477/877 L21 - 27

27

Introduction to phylogeny

➤ **Outgroup**

- used to "root" unrooted tree
- Biological information is required to choose appropriate outgroup

Which tree has the most appropriate outgroup?

BIOS477/877 L21 - 28

28

Phylogenetic methods

➤ **Data types and tree-building methods**

[Data types]

	Distances	Characters
Clustering	UPGMA	
	Neighbor joining	
Optimality criterion	Minimum evolution	Maximum parsimony
	Fitch-Margoliash	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 L21 - 30

30

Phylogenetic methods (Distance)

➤ **UPGMA: unweighted pair-group method with arithmetic mean**

- reconstructs **ultrametric trees**
 - all terminal nodes are equidistant from the root of the tree
 - equivalent to assuming a **molecular clock**
 - assumes all sequences evolve at the same rate
- reconstructs a **rooted tree**
- extremely sensitive to unequal rates in different lineages → could result in a wrong topology

BIOS477/877 L21 - 31

31

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

Example: a distance matrix for 5 sequences

The smallest distance pair is grouped together

	B	C	D	E
A	.53	.99	1.02	.82
B		.80	.93	.73
C			.65	.81
D				.94

$C \beta$

BIOS477/877 L21 - 32

32

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

	B	C	D	E
A	.53	.99	1.02	.82
B		.80	.93	.73
C			.65	.81
D				.94

	C	D	E
A/B			
C			
D			

Distance matrix is recalculated with A and B as one group

$C \beta$

BIOS477/877 L21 - 33

33

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

	B	C	D	E
A	.53	.99	1.02	.82
B		.80	.93	.73
C			.65	.81
D				.94

	C	D	E
A/B	.895	.975	.775
C		.65	.81
D			.94

Distance matrix is recalculated with A and B as one group

$d(A/B, C) = [d(A, C) + d(B, C)]/2$
 $d(A/B, D) = [d(A, D) + d(B, D)]/2$
 $d(A/B, E) = [d(A, E) + d(B, E)]/2$

$C \beta$

BIOS477/877 L21 - 34

34

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

	B	C	D	E
A	.53	.99	1.02	.82
B		.80	.93	.73
C			.65	.81
D				.94

	C	D	E
A/B	.895	.975	.775
C		.65	.81
D			.94

$C \beta$ ← The smallest distance pair

$d(A/B, C) = [d(A, C) + d(B, C)]/2$
 $d(A/B, D) = [d(A, D) + d(B, D)]/2$
 $d(A/B, E) = [d(A, E) + d(B, E)]/2$

$C \beta$

BIOS477/877 L21 - 35

35

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

	B	C	D	E
A	.53	.99	1.02	.82
B		.80	.93	.73
C			.65	.81
D				.94

	C	D	E
A/B	.895	.975	.775
C		.65	.81
D			.94

	C/D	E
A/B	.935	.775
C/D		.875

If ij is the new cluster:
 $d(i, j, k) = \frac{n_i}{n_i + n_j} d(i, k) + \frac{n_j}{n_i + n_j} d(j, k)$
 (n_i and n_j are the numbers of taxa in the groups)

Or from the original matrix, calculated also as

 $d(C/D, A/B) = [d(C, A/B) + d(D, A/B)]/2$
 $d(C/D, A/B) = [d(C, A) + d(C, B) + d(D, A) + d(D, B)]/4$
 $d(C/D, E) = [d(C, E) + d(D, E)]/2$

2-ways to calculate

$C \beta$

BIOS477/877 L21 - 36

36

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

	B	C	D	E
A	.53	.99	1.02	.82
B		.80	.93	.73
C			.65	.81
D				.94

	C	D	E
A/B	.895	.975	.775
C		.65	.81
D			.94

	C/D	E
A/B	.935	.775
C/D		.875

← The smallest distance pair

Or from the original matrix, calculated also as

 $d(C/D, A/B) = [d(C, A/B) + d(D, A/B)]/2$
 $d(C/D, A/B) = [d(C, A) + d(C, B) + d(D, A) + d(D, B)]/4$
 $d(C/D, E) = [d(C, E) + d(D, E)]/2$

2-ways to calculate

$C \beta$

BIOS477/877 L21 - 37

37

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

Or from the original matrix, calculated also as $d(A/B/E, C/D) = [d(A,C)+d(A,D)+d(B,C)+d(B,D)+d(E,C)+d(E,D)]/3(2)$

$d(i,j,k) = \frac{n_i}{n_i+n_j}d(i,k) + \frac{n_j}{n_i+n_j}d(j,k)$

$d(A/B/E, C/D) = (2/3)d(A/B,C/D) + (1/3)d(E,C/D)$

2-ways to calculate

BIOS477/877 L21 - 38

38

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

Each branch length is a half of the distance

BIOS477/877 L21 - 39

39

Phylogenetic methods (Distance)

➤ **UPGMA:** unweighted pair-group method with arithmetic mean

- reconstructs **ultrametric trees**
- all terminal nodes are equidistant from the root of the tree
- equivalent to assuming a **molecular clock**
- assumes all sequences evolve at the same rate
- reconstructs a **rooted tree**
- extremely sensitive to unequal rates in different lineages → could result in a wrong topology

BIOS477/877 L21 - 40

40

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 L21 - 41

41

Phylogenetic methods (Distance)

➤ Fitch-Margoliash method (weighted least-square) (Fitch and Margoliash, 1967)

- Initial tree: constructed by clustering 2 OTUs with shortest distances → similar to UPGMA
- No constant rate assumption (additive trees)
- Reconstructs **unrooted trees**
- Alternative trees are tested to identify the best tree based on the smallest percent standard deviation (PSD):

$$PSD = \sqrt{\frac{2 \sum_{ij} \{(D_{ij} - E_{ij}) / D_{ij}\}^2}{n(n-1)}} \times 100$$

n : number of taxa in the tree
 D_{ij} : observed distances between i and j
 E_{ij} : estimated distances between i and j (calculated by branch lengths)

BIOS477/877 L21 - 42

42

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: 3 taxa (Fitch and Margoliash, 1967)

d_{AB}, d_{AC}, d_{BC} : distances between sequences A, B, and C
 a, b, c : branch lengths

$$\begin{cases} d_{AB} = a + b & a = (d_{AB} + d_{AC} - d_{BC}) / 2 \\ d_{AC} = a + c & b = (d_{AB} + d_{BC} - d_{AC}) / 2 \\ d_{BC} = b + c & c = (d_{AC} + d_{BC} - d_{AB}) / 2 \end{cases}$$

BIOS477/877 L21 - 43

43

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

	B	C	D
A	d_{AB}	d_{AC}	d_{AD}
B		d_{BC}	d_{BD}
C			d_{CD}

- Choose two taxa (e.g., with the smallest distance) → A and B

BIOS477/877 L21 - 44

44

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

- Choose two taxa (e.g., with the smallest distance)
- Remaining taxa are combined into a single composite taxon → X

BIOS477/877 L21 - 45

45

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

Recalculate the distance matrix (as shown in UPGMA)

	B	X
A	d_{AB}	d_{AX}
B		d_{BX}

BIOS477/877 L21 - 46

46

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

$$d_{AX} = (d_{AC} + d_{AD}) / 2$$

$$d_{BX} = (d_{BC} + d_{BD}) / 2$$

→ a and b can be calculated as before

	B	X
A	d_{AB}	d_{AX}
B		d_{BX}

BIOS477/877 L21 - 47

47

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

$$\begin{cases} d_{AB} = a + b \\ d_{AX} = a + x \\ d_{BX} = b + x \end{cases} \Rightarrow \begin{cases} a = (d_{AB} + d_{AX} - d_{BX}) / 2 \\ b = (d_{AB} + d_{BX} - d_{AX}) / 2 \end{cases}$$

	B	X
A	d_{AB}	d_{AX}
B		d_{BX}

BIOS477/877 L21 - 48

48

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

- The two taxa A and B are combined into a single composite taxon → Y
- Recalculate the distances between Y and other taxa (as in UPGMA method)

	C	D
Y	d_{YC}	d_{YD}
C		d_{CD}

BIOS477/877 L21 - 49

49

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

$d_{YC} = (d_{AC} + d_{BC}) / 2$
 $d_{YD} = (d_{AD} + d_{BD}) / 2$

	C	D
Y	d_{YC}	d_{YD}
C		d_{CD}

→ c and d can be calculated as before

BIOS477/877 L21 - 50

50

Phylogenetic methods (Distance)

➤ Estimation of branch lengths: more than 3 taxa
(Fitch and Margoliash, 1967)

- If no other unresolved taxon remains,
→ z (internal branch length) can be calculated:
e.g., $z = d_{AC} - a - c$
- If there are still more taxa,
→ choose two (*e.g.*, closest) taxa and repeat the above procedure

BIOS477/877 L21 - 51

51

Phylogenetic methods (Distance)

➤ Fitch-Margoliash method (weighted least-square)
(Fitch and Margoliash, 1967)

- Initial tree: constructed by clustering 2 OTUs with shortest distances
→ similar to UPGMA
- No constant rate assumption (additive trees)
- Reconstructs **unrooted trees**
- Alternative trees are tested to identify the best tree based on the smallest percent standard deviation (PSD):

$$PSD = \sqrt{\frac{2 \sum_{ij} \{(D_{ij} - E_{ij}) / D_{ij}\}^2}{n(n-1)}} \times 100$$

n: number of taxa in the tree
D_{ij}: observed distances between *i* and *j*
E_{ij}: estimated distances between *i* and *j* (calculated by branch lengths)

BIOS477/877 L21 - 52

52