

近畿大学・農学部・生命情報学

分子系統学基礎

2009年5月12日(火)

奈良先端大・情報・蛋白質機能予測学講座

川端 猛

takawaba@is.naist.jp

<http://isw3.naist.jp/IS/Kawabata-lab/home-ja.html>

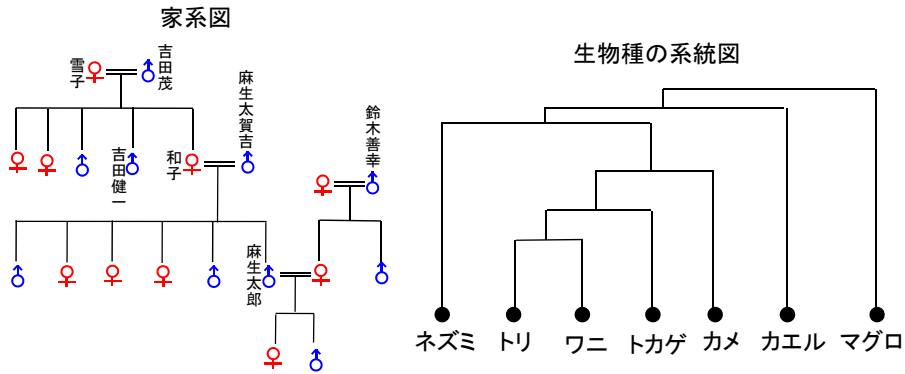
平成21年度「生命情報学&生命情報学実習」講義日程

2009.4.14

	講義	生命情報学	演習	生命情報学演習
4/7	川端1	配列決定とバイオインフォマティクス		
4/14	川端2	ペアワイズアライメントと配列相同性検索	川端	主要WEBデータベースの使用法(BLAST)
4/21	川端3	マルチプルアライメントとその応用	中村	ChemOfficeを用いた計算化学演習
4/28	川端4	蛋白質の物理化学的性質と配列解析		
5/12	川端5	分子系統学基礎	中村	系統樹作成演習(ClustalX)
5/19	川端6	蛋白質立体構造データの情報解析	川端	蛋白質立体構造データの可視化(RasMol)
5/26	川端7	>>試験<<		
6/2	金谷1	ポストゲノム解析入門(トランスクリプトーム解析)		
6/9	金谷2	ポストゲノム解析入門(インタラクトーム解析)	金谷1	発現プロフィール解析演習
6/16	金谷3	ポストゲノム解析(統合解析)	金谷2	インタラクトーム・代謝物解析演習
6/23	金谷4	メタボローム解析(その1)		
6/30	金谷5	メタボローム解析(その2)		
7/7	金谷6	メタボローム解析(その3)		
7/14	金谷7	>>試験<<		

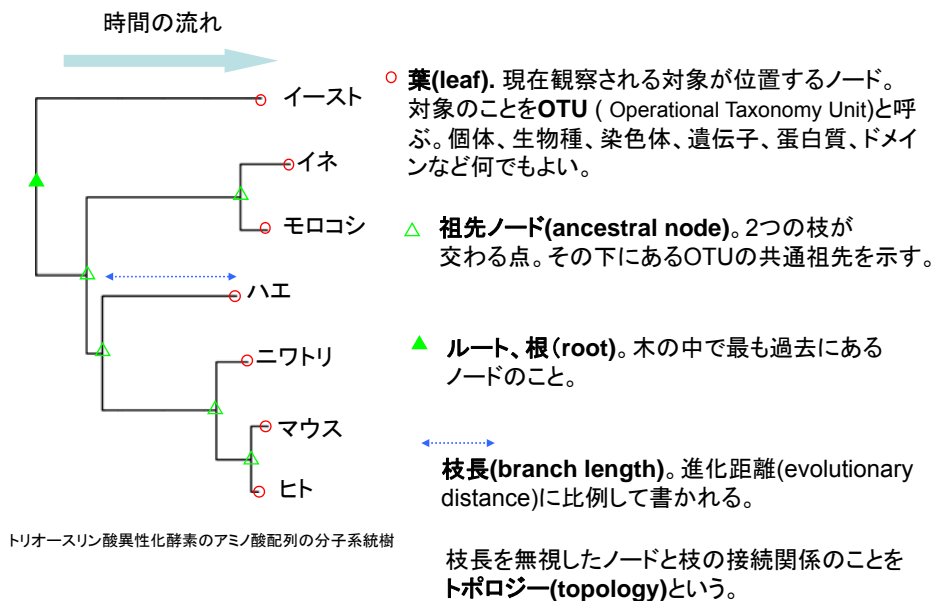
系統樹(phylogenetic tree)

対象物が生成される過程(歴史、進化史)を木構造で示したもの

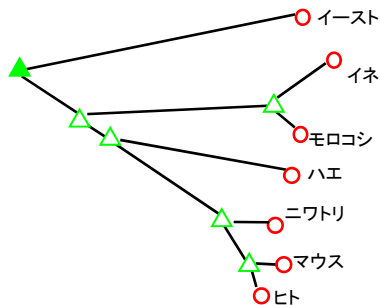


- ・何を対象にするかはいろいろ(個体、生物種、染色体、遺伝子、タンパク質)
- ・「系統樹を書く」→「過去(歴史)を推定する」
- ・「分類」(似ているものをまとめること)と「系統推定」の手続きは似ている
- ・様々な「分類法」が在り得るが、「系統樹」には唯一つの歴史的真相があるはず。

系統樹の用語



系統樹(二分岐樹)のデータ構造



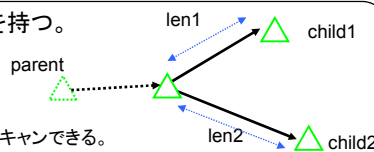
ノード(node)と枝(branch)からなるグラフ

- ・ノードには葉(leaf)ノードと祖先ノード(ancestor)ノードの2種がある。
- ・祖先ノード(ancestor)ノードから2つの子孫ノードへ枝が引かれる
- ・葉(leaf)ノードは、子孫ノードを持たない。
- ・ルートノードは、親ノードを持たない。

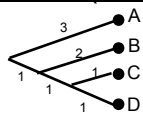
各ノードが、2つの子ノードへのポイントと、枝長を持つ。

```
struct NODE{
    struct NODE *child1,*child2;
    double len1, len2;};
```

ルートノードからスタートして再帰呼び出しすれば全ノードをスキャンできる。



・Newick(New Hampshire)フォーマット: 系統樹を括弧やカンマで記述

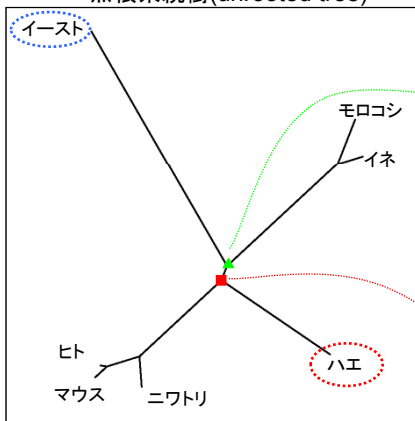


枝長なし (A,(B,(C,D)));

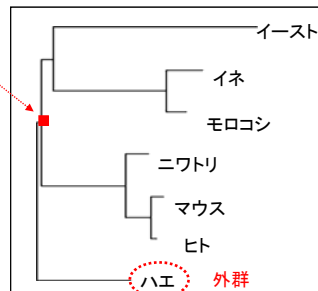
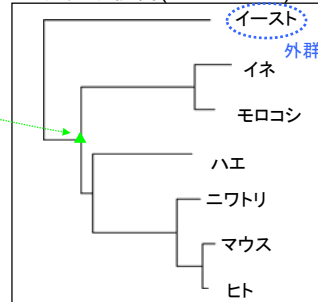
枝長つき (A:3,(B:2,(C:1,D:1):1):1);

無根と有根の系統樹

無根系統樹(unrooted tree)



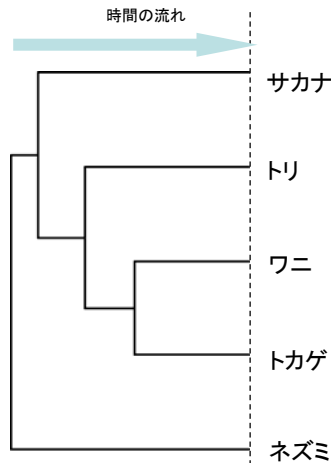
有根系統樹(rooted tree)



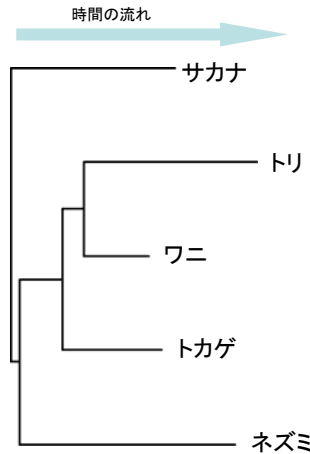
- ・近隣結合法等のアルゴリズムは、根を指定しない無根系統樹を生成する
- ・どの枝に根を置かかによって、様々な有根系統樹が生成可能。
- ・根は適当な外群(out group)の選択で決める。
外群: 他の全てのOTUと十分遠いと考えられるOTU

進化速度の同一を仮定する場合・しない場合

$$\text{進化速度} = [\text{進化距離}] / [\text{時間}]$$



進化速度が一定の場合
(UPGMA法で作成)
全てのOTU(葉ノード)が一列に揃う

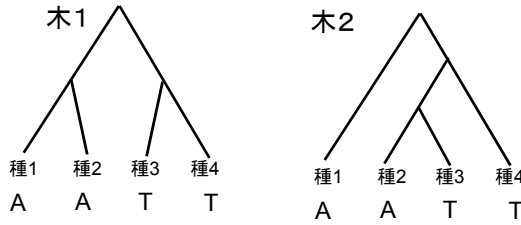


進化速度が一定でない場合
(近隣結合法で作成)
OTU(葉ノード)は一列に揃わない

分子配列からの系統樹の推定法

方法	解析方法	出力する木	計算速度	特徴
最節約法	サイト(特徴)単位	有根	遅い	アイデアは単純。分子データ以外の質的特徴にも適用可能
UPGMA法	距離行列	有根	速い	分子速度の一定性を仮定。重心間距離のクラスター解析と等価。
近隣結合法	距離行列	無根	速い	最小進化の法則を距離行列に適用。分子速度の一定性を仮定しない。
最尤法	サイト単位	有根	遅い	分子進化の確率モデルに従う。数学的な厳密さは高い。

最節約法(maximum parsimony)

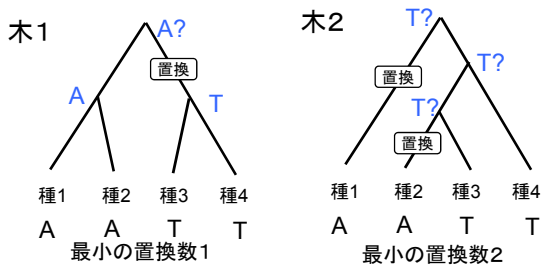


4つの生物種のある1つのサイトのDNA配列がわかったとする。

どちらの木が尤もらしいか？

(1) 総置換数が最小になるように、祖先形質を推定

(2) 総置換数が最小の木が尤もらしいとする



木1のほうが、置換数が少ない
→木1のほうが木2より尤もらしい

最節約の考え(最小進化の法則)

現在の生物の形質を表現する仮説(系統樹)の中で、進化による変化の回数がか最も少ない仮説が正しい。

最小進化の法則(minimum evolution principle)、オッカムの剃刀(Ockham's razor)

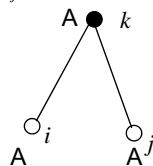
最節約法による最少置換数の推定アルゴリズム (traditional parsimony)

[初期化]
 $Cost=0, k=2n-1$ (ルートノード)

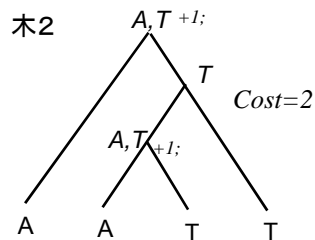
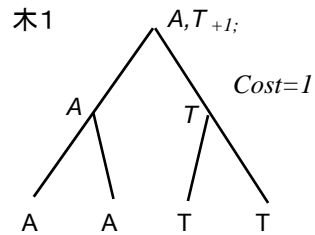
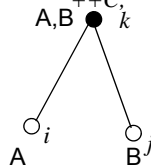
[再帰的実行]
 k が葉ノードなら、
 $R_k = x_k$
 k が葉ノードでないなら、 i, j を k の子ノードとすると、子ノードの R_i, R_j が計算されていないなら、
 R_i, R_j を計算(再帰呼び出し)。
計算されているなら、以下のように R_k を計算
 $R_i \cap R_j$ が空でないなら、 $R_k = R_i \cap R_j$
 $R_i \cap R_j$ が空なら、 $R_k = R_i \cup R_j, Cost$ に1加算

[終了処理]
 $Cost$ が最小コスト

$R_i \cap R_j$ が空でないなら、
 $R_k = R_i \cap R_j$



$R_i \cap R_j$ が空なら、
 $R_k = R_i \cup R_j, Cost$ に1加算

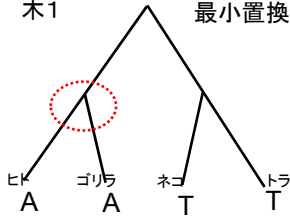


最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

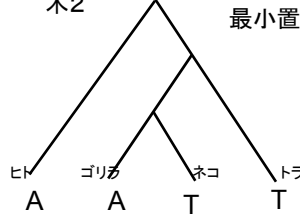
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、C=A and B
 A and Bが空なら、C=A or B Costに1加算

木1 最小置換数: _____

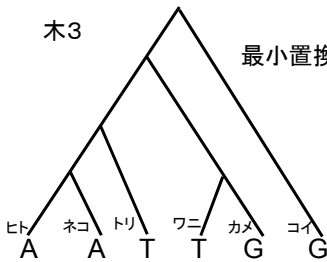


Cost = 0

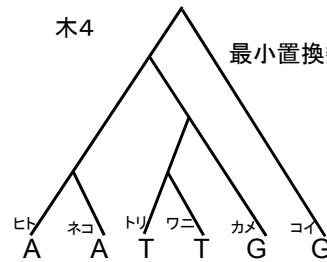
木2 最小置換数: _____



木3 最小置換数: _____



木4 最小置換数: _____

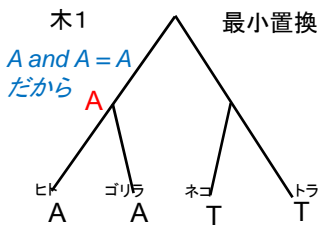


最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

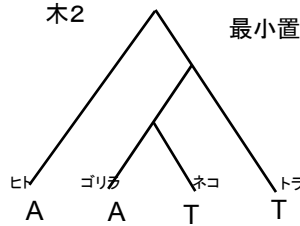
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、C=A and B
 A and Bが空なら、C=A or B Costに1加算

木1 最小置換数: _____

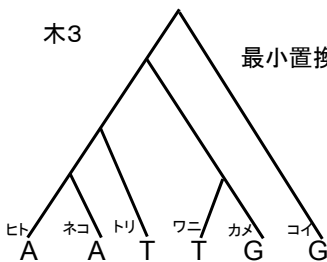


Cost = 0

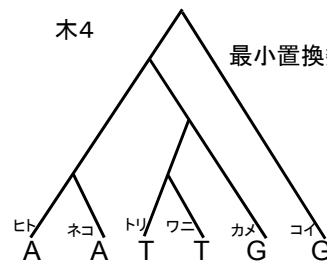
木2 最小置換数: _____



木3 最小置換数: _____



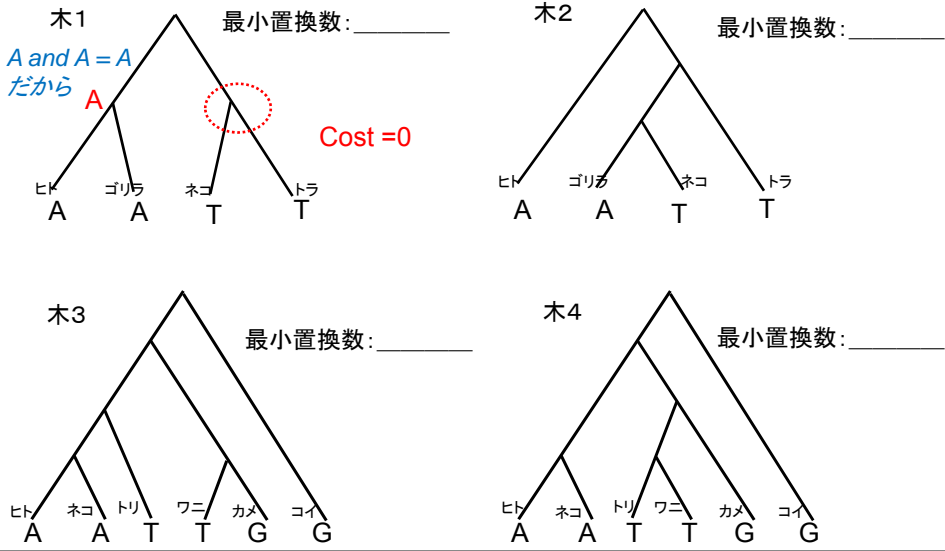
木4 最小置換数: _____



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

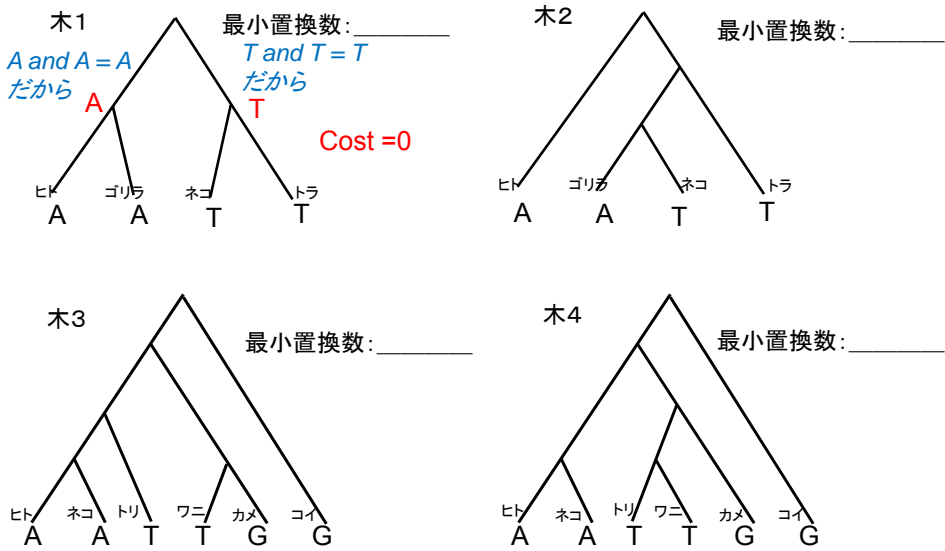
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A \text{ and } B$
 A and Bが空なら、 $C=A \text{ or } B$ Costに1加算



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

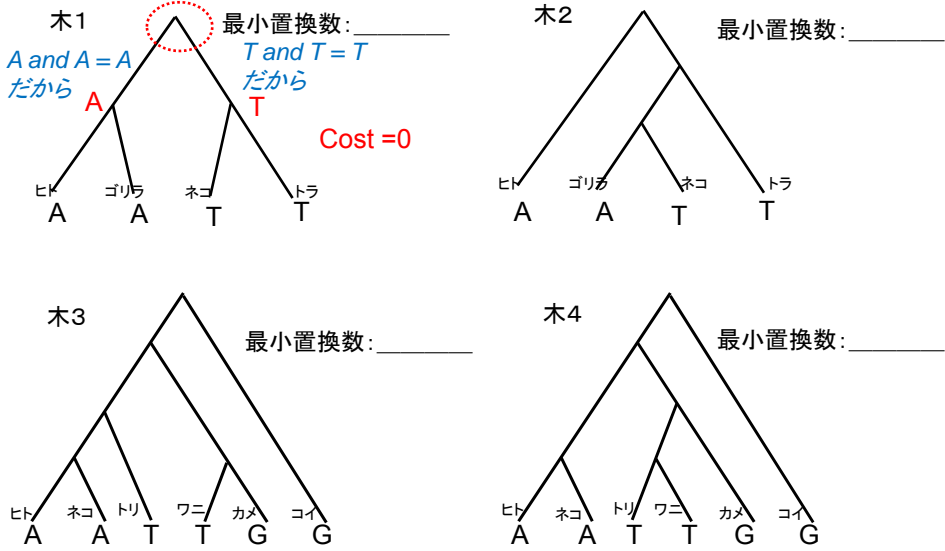
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A \text{ and } B$
 A and Bが空なら、 $C=A \text{ or } B$ Costに1加算



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

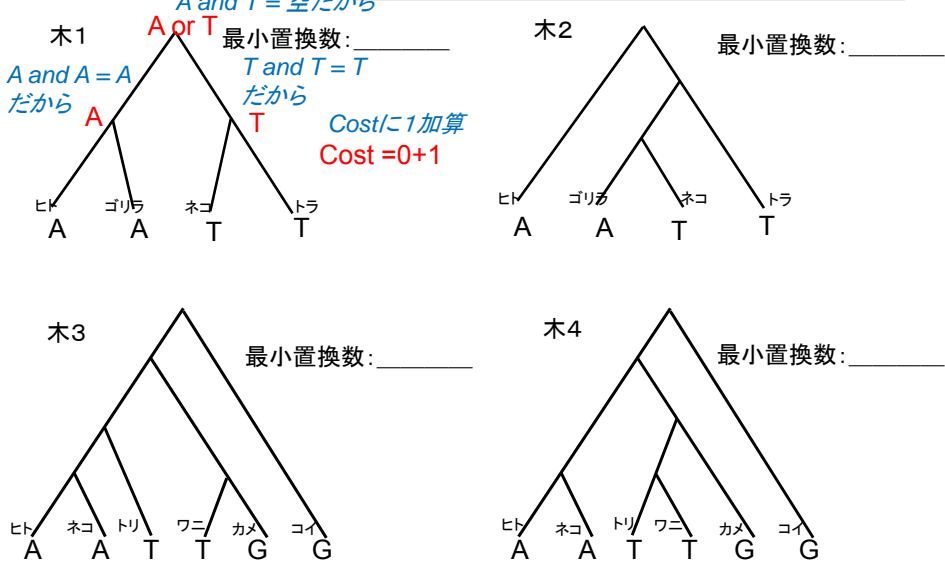
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A \text{ and } B$
 A and Bが空なら、 $C=A \text{ or } B$ Costに1加算



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

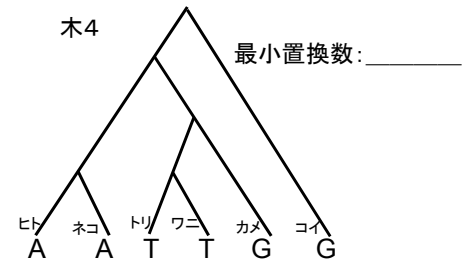
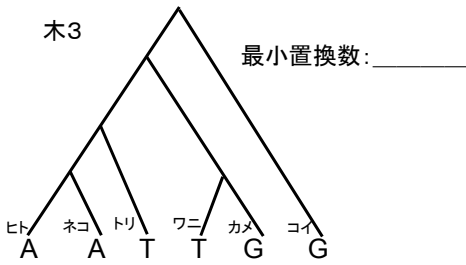
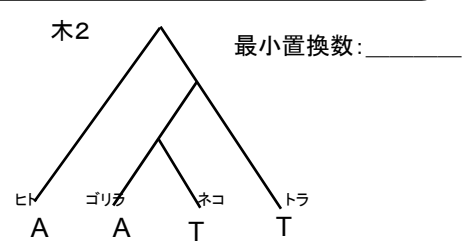
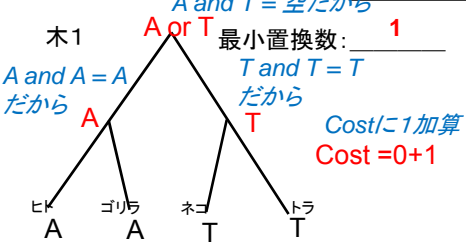
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A \text{ and } B$
 A and Bが空なら、 $C=A \text{ or } B$ Costに1加算



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

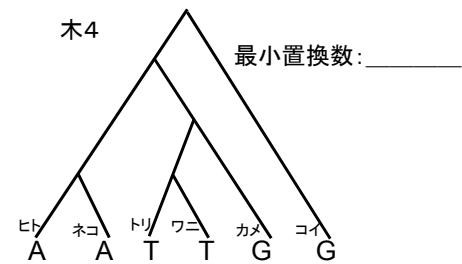
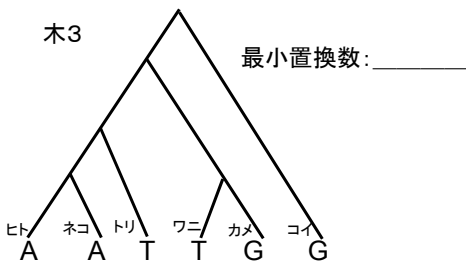
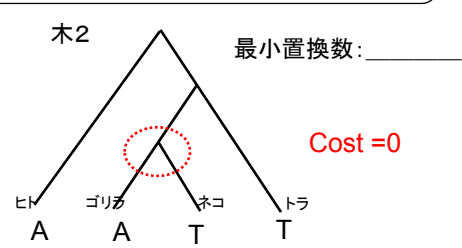
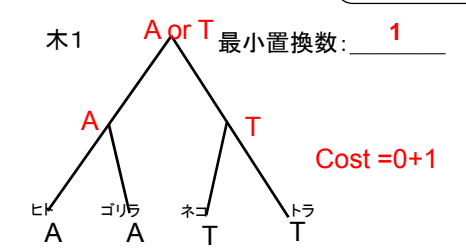
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A$ and B
 A and Bが空なら、 $C=A$ or B Costに1加算



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A$ and B
 A and Bが空なら、 $C=A$ or B Costに1加算

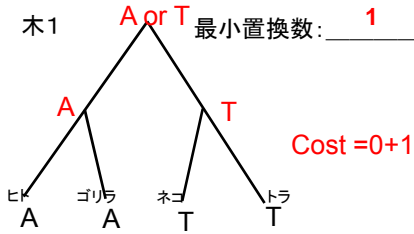


最節約法による最小置換数

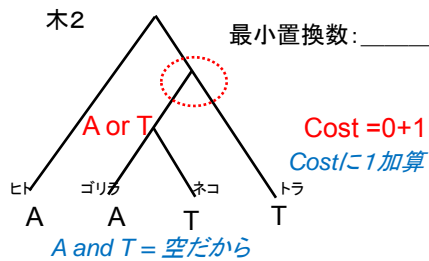
最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A$ and B
 A and Bが空なら、 $C=A$ or B Costに1加算

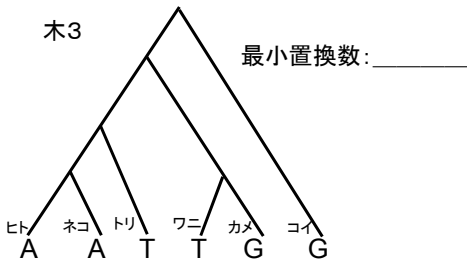
木1 最小置換数: 1



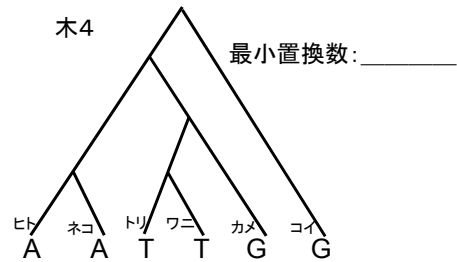
木2 最小置換数: _____



木3 最小置換数: _____



木4 最小置換数: _____

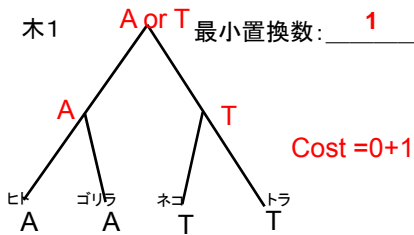


最節約法による最小置換数

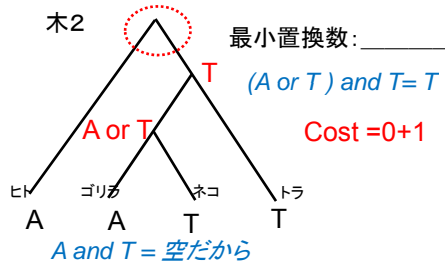
最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A$ and B
 A and Bが空なら、 $C=A$ or B Costに1加算

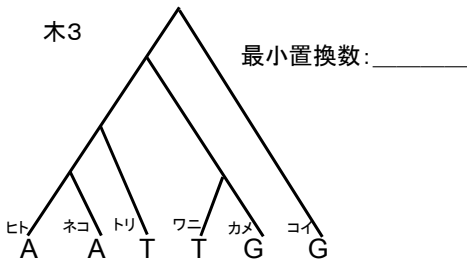
木1 最小置換数: 1



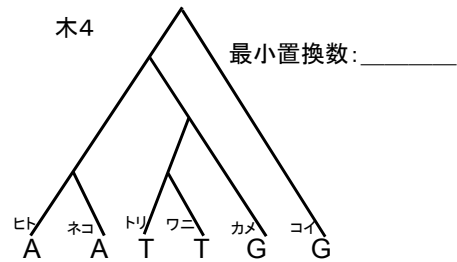
木2 最小置換数: _____



木3 最小置換数: _____



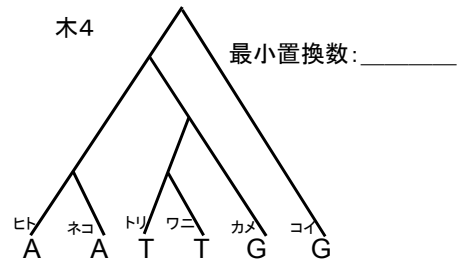
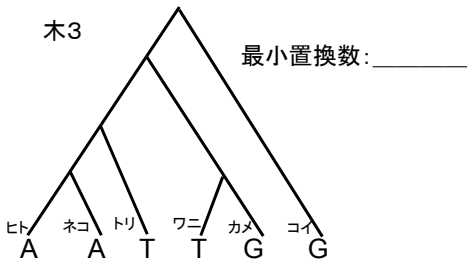
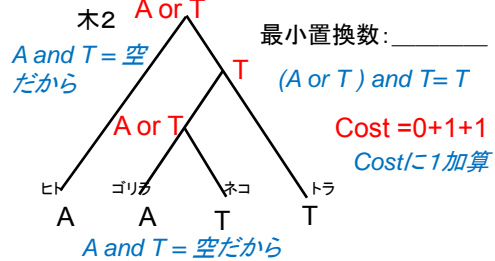
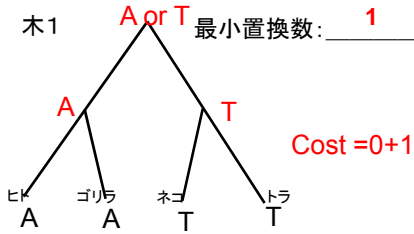
木4 最小置換数: _____



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

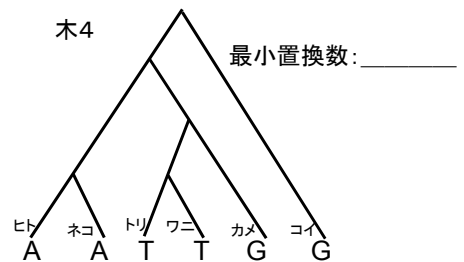
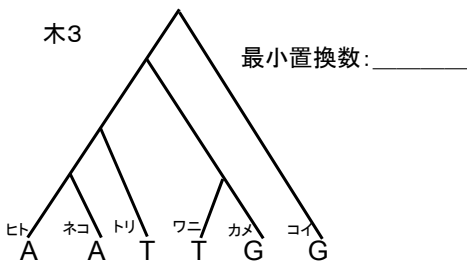
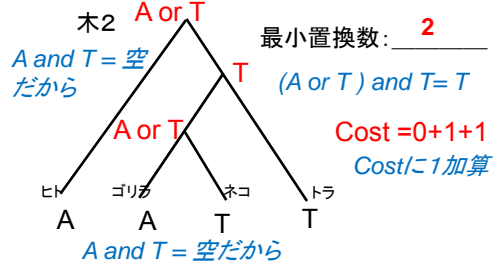
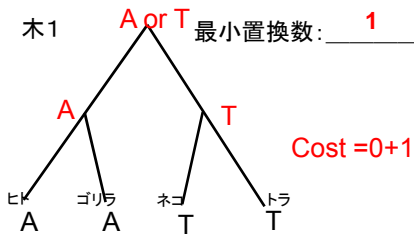
子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A \text{ and } B$
 A and Bが空なら、 $C=A \text{ or } B$ Costに1加算



最節約法による最小置換数

最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A \text{ and } B$
 A and Bが空なら、 $C=A \text{ or } B$ Costに1加算

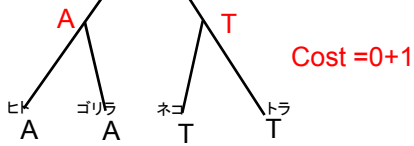


最節約法による最小置換数

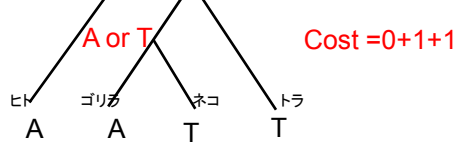
最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A$ and B
 A and Bが空なら、 $C=A$ or B Costに1加算

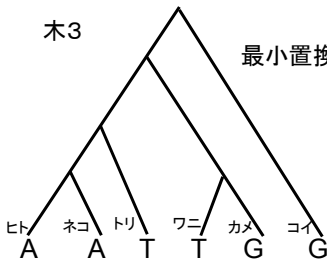
木1 **A or T** 最小置換数: 1



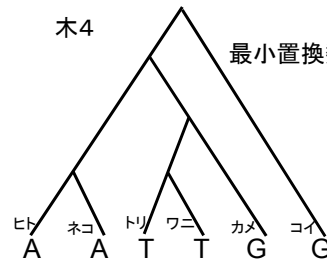
木2 **A or T** 最小置換数: 2



木3 最小置換数: _____



木4 最小置換数: _____

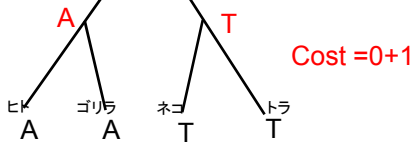


最節約法による最小置換数

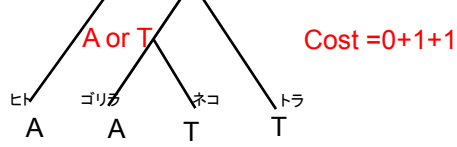
最節約法を用いて以下の系統樹の祖先形質を推定し、最小置換数を求めなさい。

子ノードがA,Bなら、親ノードCは
 A and Bが空でないなら、 $C=A$ and B
 A and Bが空なら、 $C=A$ or B Costに1加算

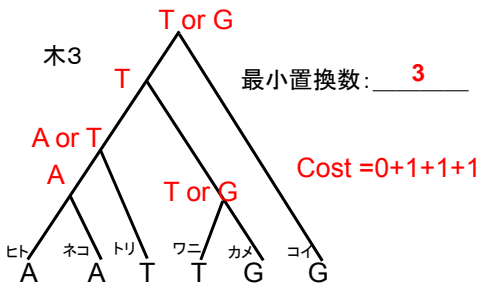
木1 **A or T** 最小置換数: 1



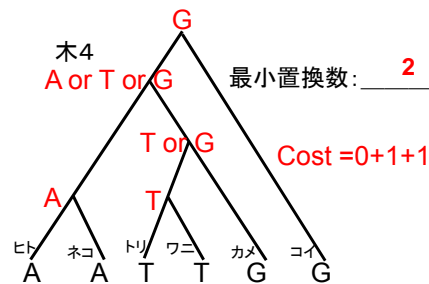
木2 **A or T** 最小置換数: 2



木3 最小置換数: 3



木4 最小置換数: 2



Traditional Parsimonyの使用上の注意

- Traditional Parsimonyはコストは正しく計算される。しかし、祖先形質は可能な組み合わせの一部しか計算されない。

→ コストだけを知りたい場合、あるいは祖先形質の一部の解だけを(手計算で)知りたいときに有効

→ より本格的な計算にはWeighted Parsimonyを用いて(計算機で)計算すべき

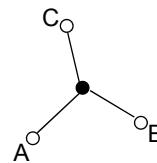
参考文献: Durbin R., Eddy S., Krogh A., Mitchson G. "Biological Sequence analysis", Cambridge University Press, 1998. Chapter 7

可能な木のトポロジーの数

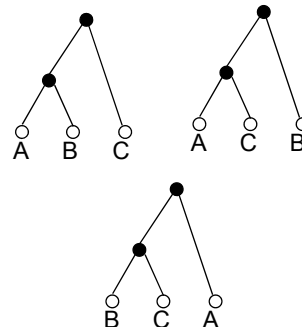
$$\prod_{k=3}^N (2k-5) \quad \prod_{k=3}^N (2k-3)$$

OTU数 N	無根系統樹	有根系統樹
3	1	3
4	3	15
5	15	105
6	105	945
7	945	10395
8	10395	135135
9	135135	2027025
10	2027025	34459425

N=3の場合の無根系統樹のトポロジー



N=3の場合の有根系統樹のトポロジー



最節約法の特徴

- 分子データに限らず、様々な形質に対して適用可能
骨、化石など生物の形態から系統樹を推定する唯一の方法
- 祖先形質の推定が可能
- 「最節約 / 最小進化」という考え方は、全ての系統推定の基本
- 配列・特徴の数が増えた場合、膨大な計算時間が必要となる
祖先形質の推定が必要。トポロジー探索は全回探索が基本。配列数が10を超える場合、分岐限定法あるいはヒューリスティック検索の適用が必須。
- 各特徴が独立・無相関であることが前提
- 多重置換等、複雑な進化のモデルを扱えない



	塩基配列	羽毛	二足歩行	心臓	体温
種1	A G G G	ない	不可能	1心房1心室	変温
種2	A G A A	ない	不可能	2心房1心室	変温
種3	T G A A	ない	不可能	2心房2心室	変温
種4	T A G A	ある	可能	2心房2心室	恒温

距離行列法

なんらかの方法でOTU間の距離(進化距離)を定義し、距離行列を作成。
その距離をできるだけ満たすような木を計算する方法

アライメント

配列 1 AAAAA
配列 2 AAAAT
配列 3 TAATA
配列 4 TAATT

距離行列 d_{ij}
(不一致サイト数)

	1	2	3	4
1	0	1	2	3
2	1	0	2	2
3	2	2	0	1
4	3	2	1	0

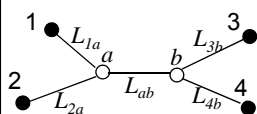
とか

距離行列 d_{ij}
(p距離)

	1	2	3	4
1	0.0	0.2	0.4	0.6
2	0.2	0.0	0.4	0.4
3	0.4	0.4	0.0	0.2
4	0.6	0.4	0.2	0.0

※距離行列の大きさは配列の本数だけに依存、
それぞれの配列の長さには依存しない。

$$p\text{距離} = \frac{[\text{不一致のサイト数}]}{[\text{比較したサイト数}]}$$



$$d_{12} \doteq L_{1a} + L_{2a} \quad d_{34} \doteq L_{3b} + L_{4b}$$

$$d_{13} \doteq L_{1a} + L_{ab} + L_{3b} \quad d_{14} \doteq L_{1a} + L_{ab} + L_{4b}$$

$$d_{23} \doteq L_{2a} + L_{ab} + L_{3b} \quad d_{24} \doteq L_{2a} + L_{ab} + L_{4b}$$

木の枝長の和
が距離行列の
値になるように木の
トポロジーと枝長を推定

配列データからの進化距離の推定

進化距離: 1サイトあたりに受けた置換の回数

分子時計:

DNAやアミノ酸配列の違いが生じる速度(進化速度)は近似的に一定であること。

分子進化の中立説(木村資生、1968)

DNAやアミノ酸配列が進化の過程で受ける変異のほとんどは、自然選択の上からは、よくも悪くもない“中立的”なものであるという仮説。

p -距離 : 最も単純な進化距離の推定法

$$p\text{-距離} = n_d / n$$

n : 比較したサイトの数
 n_d : 配列が異なっていたサイトの数

GAALSTLLS
GGVVSTLVA

$$p\text{-距離} = 4 / 10 = 0.4$$

多重置換の影響を考慮した距離

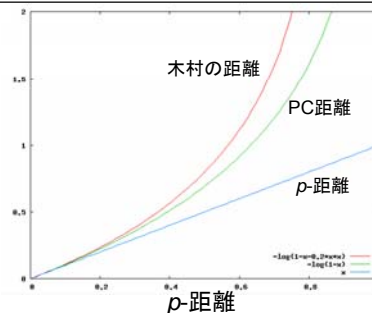
p -距離

0: AAAAAAAAAA	0.0
1: AKAAAAAAAA	0.1
2: PKAAAAAAAA	0.2
3: PKAAMAAAAA	0.3
4: PKAAMAIAAA	0.4
5: PKAAMAIARA	0.5
6: PKAAMADARA	0.5
7: PKAAMADARR	0.6
8: PKAAMADATR	0.6
9: PKAAMADRTR	0.7
10: PKAANADRTR	0.7
11: PKAANADWTR	0.7
12: PKVANADWTR	0.8
13: PKVAADWTR	0.7
14: NKVAAADWTR	0.7

多重置換 : 進化時間が長いときに、同じサイトに複数回の置換が起こること。

$$\text{PC距離 (Poisson Correction)} = -\log(1-p)$$

$$\text{木村の距離} = -\log(1 - p - 0.2p^2)$$



UPGMA法

Unweighted Pair-Group Method with Arithmetic mean

[初期化]

全ての配列間の距離 d_{ij} を計算。それぞれの配列 i が一つのクラスタ C_i を構成するとする。

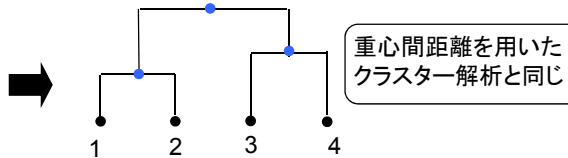
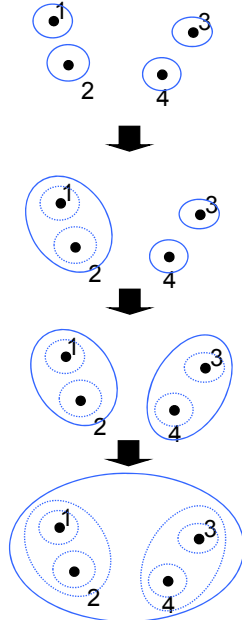
[反復]

(1) 全てのクラスタのペアの中で距離 d_{ij} が最小のペア C_i と C_j を選び、融合して新しいクラスタ $C_k = C_i \cup C_j$ を作る。このとき、 C_i と C_j を子にもつ親ノードを枝長の高さが $d_{ij}/2$ になるように作る

(2) 距離行列を更新する。クラスタ間の距離は、属する配列間の平均距離で定義する。

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

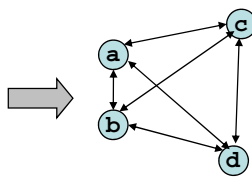
クラスタ数が1つになるまで反復する。



UPGMA法による系統樹の計算例(1)

不一致文字数を距離とする

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



距離行列

	a	b	c	d
a	0			
b	X	0		
c	X	X	0	
d	X	X	X	0

最小距離のペアを選んで融合

距離行列

	0			
	X	0		
	X	X	0	

クラスタと配列の距離は、配列間平均の距離とする

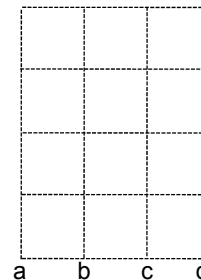
最小距離のペアを選んで融合

距離行列

	0		
	X	0	

クラスタとクラスタの距離は、クラスタのメンバーの配列間の平均の距離とする

系統樹



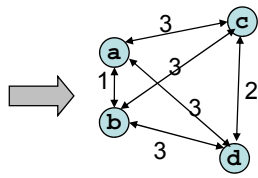
距離の半分が枝長

UPGMA法による系統樹の計算例(2)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離のペアを選んで融合

距離行列

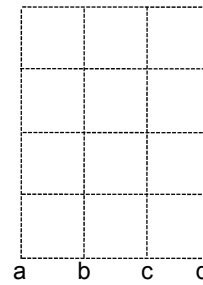
	a	b	c	d
a	0			
b	X	0		
c	X	X	0	
d				0

最小距離のペアを選んで融合

距離行列

	a	b	c	d
a	0			
b	X	0		
c			0	
d				0

系統樹



クラスタと配列の距離は、配列間平均の距離とする

クラスタとクラスタの距離は、クラスタのメンバーの配列間の平均の距離とする

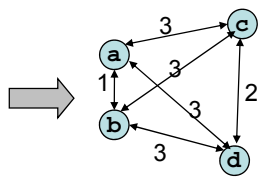
距離の半分が枝長

UPGMA法による系統樹の計算例(3)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離のペアを選んで融合

距離行列

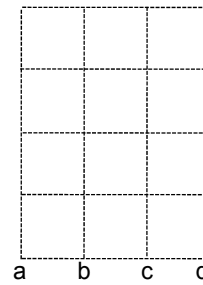
	a	b	c	d
a	0			
b	X	0		
c	X	X	0	
d				0

最小距離のペアを選んで融合

距離行列

	a	b	c	d
a	0			
b	X	0		
c			0	
d				0

系統樹



クラスタと配列の距離は、配列間平均の距離とする

クラスタとクラスタの距離は、クラスタのメンバーの配列間の平均の距離とする

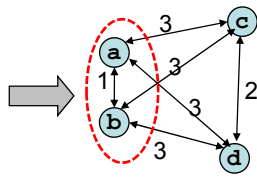
距離の半分が枝長

UPGMA法による系統樹の計算例(4)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離
のペアを
選んで融合

距離行列

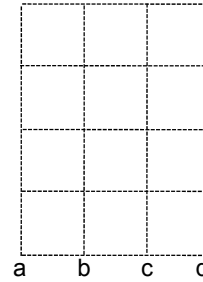
最小距離
のペアを
選んで融合

距離行列

	a,b	c	d
a,b	0		
c	X	0	
d	X	X	0

	a,b	c
a,b	0	
c	X	0

系統樹



クラスタと配列の距離は、
配列間平均の距離とする

クラスタとクラスタの
距離は、クラスタの
メンバーの配列間の
平均の距離とする

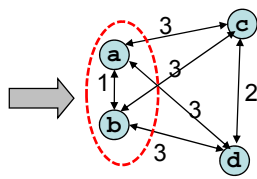
距離の半分が枝長

UPGMA法による系統樹の計算例(5)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離
のペアを
選んで融合

距離行列

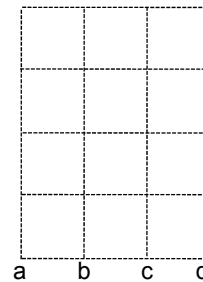
最小距離
のペアを
選んで融合

距離行列

	a,b	c	d
a,b	0	3	3
c	X	0	2
d	X	X	0

	a,b	c
a,b	0	
c	X	0

系統樹



$$(3+3)/2=3 \quad (3+3)/2=3$$

クラスタと配列の距離は、
配列間平均の距離とする

クラスタとクラスタの
距離は、クラスタの
メンバーの配列間の
平均の距離とする

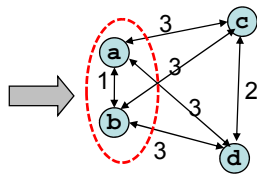
距離の半分が枝長

UPGMA法による系統樹の計算例(6)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離
のペアを
選んで融合

距離行列

	a,b	c	d
a,b	0	3	3
c	X	0	2
d	X	X	0

$$(3+3)/2=3 \quad (3+3)/2=3$$

クラスタと配列の距離は、
配列間平均の距離とする

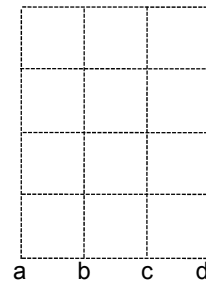
最小距離
のペアを
選んで融合

距離行列

	a,b	c,d
a,b	0	
c,d	X	0

クラスタとクラスタの
距離は、クラスタの
メンバーの配列間の
平均の距離とする

系統樹



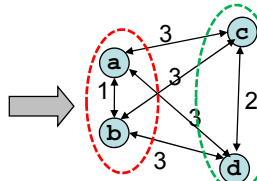
距離の半分が枝長

UPGMA法による系統樹の計算例(7)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離
のペアを
選んで融合

距離行列

	a,b	c	d
a,b	0	3	3
c	X	0	2
d	X	X	0

$$(3+3)/2=3 \quad (3+3)/2=3$$

クラスタと配列の距離は、
配列間平均の距離とする

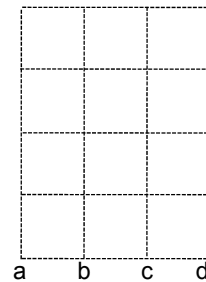
最小距離
のペアを
選んで融合

距離行列

	a,b	c,d
a,b	0	
c,d	X	0

クラスタとクラスタの
距離は、クラスタの
メンバーの配列間の
平均の距離とする

系統樹



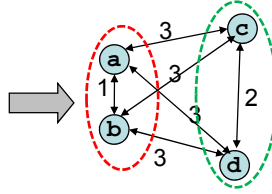
距離の半分が枝長

UPGMA法による系統樹の計算例(8)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離のペアを選んで融合

距離行列

最小距離のペアを選んで融合

距離行列

	a,b	c	d
a,b	0	3	3
c	X	0	2
d	X	X	0

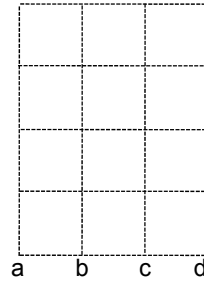
	a,b	c,d
a,b	0	3
c,d	X	0

$(3+3)/2=3$ $(3+3)/2=3$

$(3+3+3+3)/4=3$
クラスタとクラスタの距離は、クラスタのメンバーの配列間の平均の距離とする

クラスタと配列の距離は、配列間平均の距離とする

系統樹



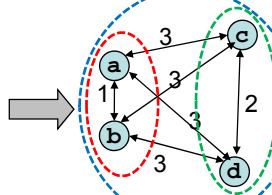
距離の半分が枝長

UPGMA法による系統樹の計算例(9)

不一致文字数を距離とする

距離行列

配列a GACT
配列b GTCT
配列c CCAT
配列d CGTT



	a	b	c	d
a	0	1	3	3
b	X	0	3	3
c	X	X	0	2
d	X	X	X	0

最小距離のペアを選んで融合

距離行列

最小距離のペアを選んで融合

距離行列

	a,b	c	d
a,b	0	3	3
c	X	0	2
d	X	X	0

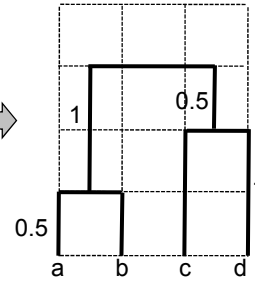
	a,b	c,d
a,b	0	3
c,d	X	0

$(3+3)/2=3$ $(3+3)/2=3$

$(3+3+3+3)/4=3$
クラスタとクラスタの距離は、クラスタのメンバーの配列間の平均の距離とする

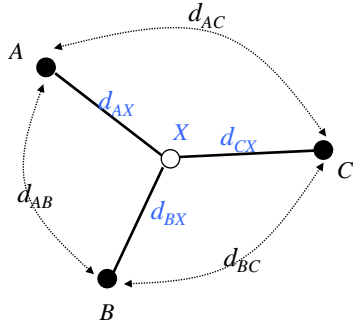
クラスタと配列の距離は、配列間平均の距離とする

系統樹



距離の半分が枝長

Fitch-Margoliashの式



もとの距離行列 d_{ij} を再現することを3つのOTUについて考える。

OTUが3つA,B,Cの場合、その間の3つの距離 d_{AB} , d_{BC} , d_{AC} を満たすように、祖先ノードXを作成して、木を作成する。

連立1次方程式
$$\begin{cases} d_{AX} + d_{BX} = d_{AB} \\ d_{BX} + d_{CX} = d_{BC} \\ d_{AX} + d_{CX} = d_{AC} \end{cases}$$
 を解くと、

OTUが3つの場合、この式で、距離行列を完全に満たす枝長を求めることができる。

$$d_{AX} = (d_{AB} + d_{AC} - d_{BC})/2$$

$$d_{BX} = (d_{AB} + d_{BC} - d_{AC})/2$$

$$d_{CX} = (d_{AC} + d_{BC} - d_{AB})/2$$

近隣結合法 (Neighbor-Joining法、NJ法)

Saito, N., Nei, N. Mol. Biol. Evol. 4, 406-425, 1987.

[初期化]

L (相互結合したノード集合) をOTUの集合とする。

[反復]

(1) $d_{ij} - r_i - r_j$ が最小となる i, j を L から選択。

$$r_i = \frac{1}{|L| - 2} \sum_{m \in L} d_{im} \quad \text{他のノードへの平均距離のような値}$$

子ノード i, j を持つ親ノード k を作成し、 L に加える。

また、 L からノード i, j を除く。

(2) 距離行列を更新する。

新ノード k の距離行列は、Fitch-Margoliashの式から、

$$d_{mk} = (d_{im} + d_{jm} - d_{ij}) / 2$$

$$d_{ik} = (d_{ij} + d_{im} - d_{jm}) / 2$$

$$d_{jk} = (d_{ij} + d_{jm} - d_{im}) / 2$$

で定義。ただし、木の枝長となる d_{ik}, d_{jk} については、 L に属する全ての m についての平均の枝長を用いる。

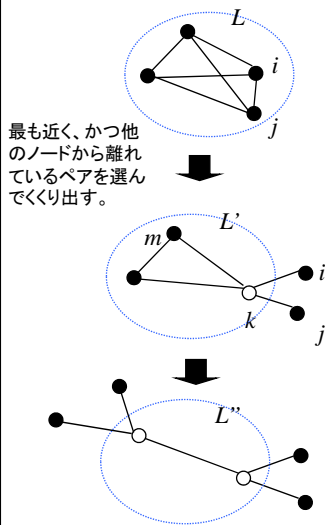
$$d_{ik} = \langle (d_{ij} + d_{im} - d_{jm}) / 2 \rangle_m = (d_{ij} + r_i - r_j) / 2$$

$$d_{jk} = \langle (d_{ij} + d_{jm} - d_{im}) / 2 \rangle_m = (d_{ij} + r_j - r_i) / 2$$

[終了処理]

L が2つのノードを含むだけになったら終了

残ったノードのどちらかを木のルートノード(3分岐)とする。

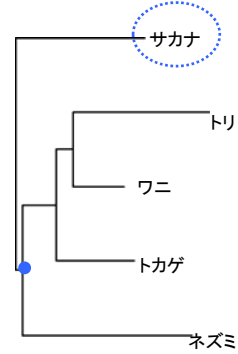
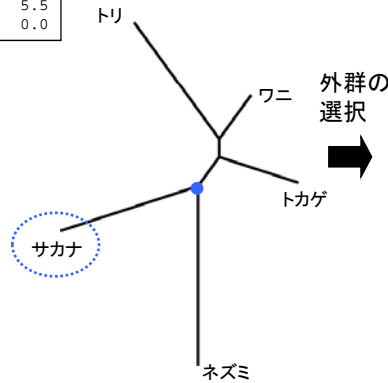
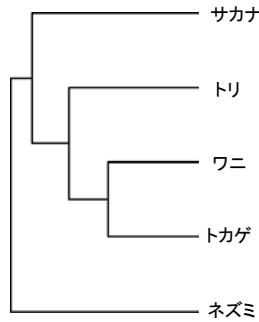


最も近く、かつ他のノードから離れているペアを選んでくり出す。

UPGMA法とNJ法の樹形の違い

距離行列

sakana	0.0	9.0	7.3	7.0	9.5
nezumi	9.0	0.0	8.3	8.0	10.5
tokage	7.3	8.3	0.0	4.3	6.8
wani	7.0	8.0	4.3	0.0	5.5
tori	9.5	10.5	6.8	5.5	0.0

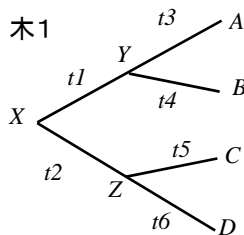


・無根系統樹から有根系統樹への変換: OTUの中から適切な外群(out group)を選べばよい。

外群の選択基準: (1)他の全てのOTUと相同、(2)他のどのOTUとも十分遠縁

最尤法(maximum likelihood)

分子進化に関する確率モデルを立て、葉ノードの形質を最もよく説明する(最も尤度が高い)系統樹を推定する。



$P_{ab}(t)$: 時間 t の間に a から b に変異する確率

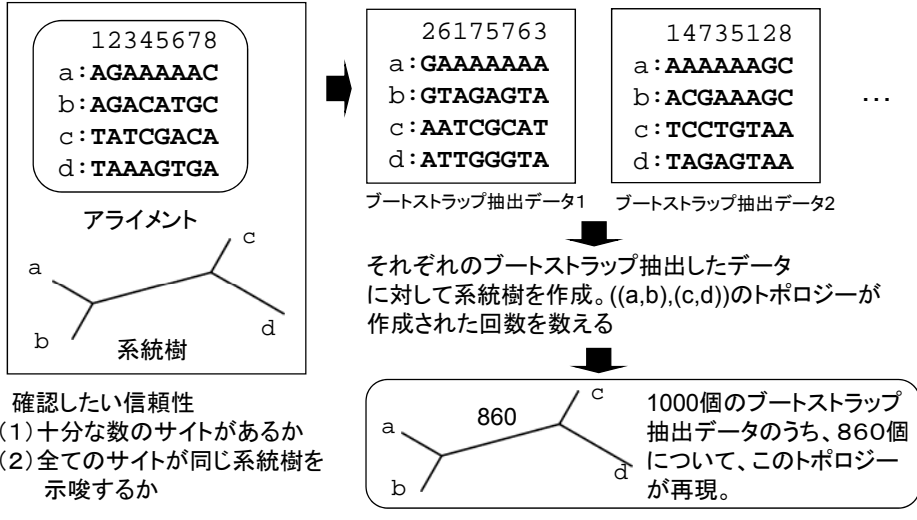
木1が起こる確率 L は以下で表される。

$$L = P(G) \cdot P_{XY}(t1) \cdot P_{YA}(t3) \cdot P_{YB}(t4) \cdot P_{XZ}(t2) \cdot P_{ZC}(t5) \cdot P_{ZD}(t6)$$

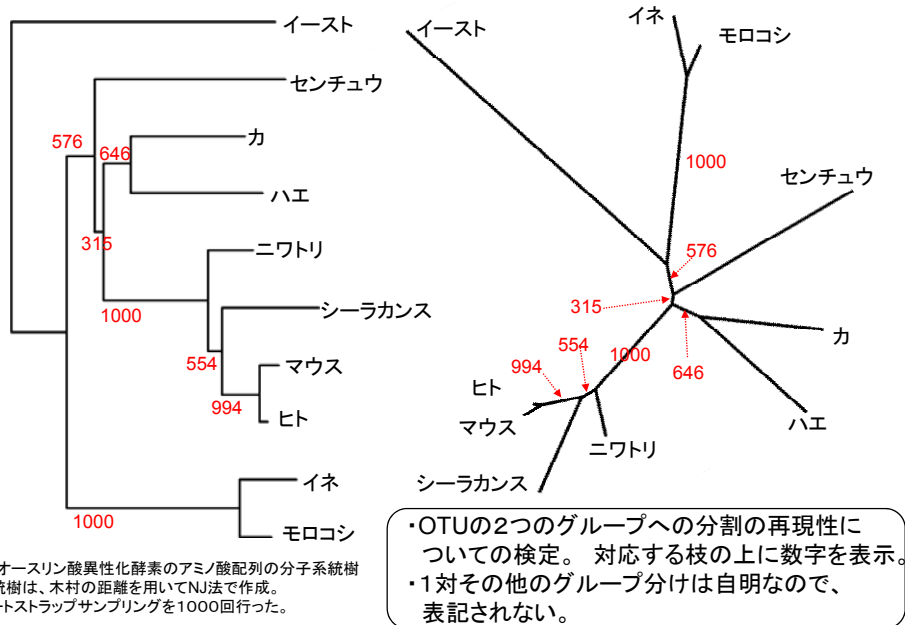
- ・あるトポロジーについて L を最大化するように枝長($t1, t2, \dots$)と祖先形質(X, Y, \dots)を計算
- ・尤度 L が最も高いトポロジーを探索する
- ・最節約法と同程度の長い計算時間を必要

系統樹のトポロジーの信頼性の検定

ブートストラップ(bootstrap)抽出を行い多数の擬似データを作成
ランダムにサイトを元の数だけ選ぶ。同じサイトを複数回選んでもかまわない。



ブートストラップ値付きの系統樹の例



分子系統樹作成のためのソフトウェア

- ClustalW/ClustalX

マルチプルアライメントのソフトだが、NJ法による系統樹作成の機能が付属。ブートストラップ計算にも対応。

- Phylip <http://evolution.genetics.washington.edu/phylip.html>

様々な系統樹作成のためのプログラムのセット。最節約法、NJ法、最尤法など多くのアルゴリズムに対応。UNIX, DOS, Macに対応。

- MEGA <http://www.megasoftware.net>

様々な系統樹作成のためのプログラムのセット。最節約法、NJ法、など多くのアルゴリズムに対応。Windows/DOS/Macに対応。

- PAUP <http://paup.csit.fsu.edu>

最節約法を中心とした系統樹作成ソフト。分子以外の形態データにも対応。有料。

分子系統樹表示のためのソフトウェア

- NJplot <http://pbil.univ-lyon1.fr/software/njplot.html>

簡素な有根系統樹の描画ソフト。

- TreeView/TreeViewX

<http://taxonomy.zoology.gla.ac.uk/rod/treeview.html>

<http://darwin.zoology.gla.ac.uk/~rpage/treeviewx/index.html>

多機能な系統樹の描画ソフト

参考文献

- 長谷川政美、岸野洋久「分子系統学」岩波書店(1996)
- 根井正利、S.クマー「分子進化と分子系統学」(2006)培風館
- 斎藤成也「ゲノム進化学入門」(2007) 共立出版
- Durbin R., Eddy S., Krogh A., Mitchson, G. "Biological Sequence analysis", Cambridge University Press, 1998. Chapter 7, 8.
- R. Durbin 他著、阿久津達也他訳「バイオインフォマティクス - 確率モデルによる遺伝子解析」医学出版、2001年、9800円