

近畿大学・農学部・生命情報学

マルチプルアライメントと 分子系統学基礎

2008年5月20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>

授業予定			
日付	担当	講義	演習
4/8(火)	黒川	バイオインフォマティクス概論	
4/15(火)	黒川	配列解析1	IMCを使ったゲノム解析
4/22(火)	黒川	配列解析2	IMCを使った比較ゲノム解析
5/13(火)	川端	ペアワイズアライメントと配列相同性解析	
5/20(火)	川端	マルチプルアライメントと分子系統学基礎	配列相同性解析と系統樹作成演習
5/27(火)	川端	タンパク質配列の分類と機能推定	
6/3(火)	川端	タンパク質立体構造データの情報解析	タンパク質立体構造データの可視化演習
6/10(火)	川端	<試験>	
6/17(火)	金谷	ポストゲノム解析入門(トランスクリプトーム解析)	
6/24(火)	金谷	ポストゲノム解析入門(インタラクトーム解析)	発現プロファイル解析演習
7/1(火)	金谷	ポストゲノム解析入門(統合解析)	インタラクトーム解析演習・代謝物解析演習
7/8(火)	金谷	メタボローム解析(その1)	
7/15(火)	金谷	メタボローム解析(その2)	
7/22(火)	金谷	<試験>	

マルチプルアライメント

(multiple sequence alignment
多重配列整列)

マルチプルアライメント(多重配列整列)とは

3本以上の配列を進化的な対応関係に従って並べること

```
>lnshA
SRPTEATERCIESLIAVPQKYAGKDGHSVTLKTEFLSPMNTLEAAFTKNQKDPGVLDKRLDLSNDSGQLDFQKQFL
NLIIGLAVAESFVKAAPPQKRF
>1j55A
MTELETAMGMIIDVFSRYSSEGSTQTLTKGELKVLMEKELPGFLDAVKLLKDLNDANGDAQVDFSEFIVFAAITS
ACHKYFEKAL
>1ig5A
KSPPEELKGI FEKYAAKEGDPNQLSKEELKLLQLTEFFPSLLKGPSTLDELFEELDKNGDGEVSFEFQVLVKKISQ
>1qx2A
MKSPPEELKGAPEVFAAKEGDPNQLSKEELKLVQMQLGPGSLLKGMSTLDEMIREVDKNGDGEVSFEFLVMMKISQ
```

↓

```
CLUSTAL W (1.83) multiple sequence alignment

lnshA      SRPTEATERCIESLIAVPQKYAGKDGHSVTLKTEFLSPMNTLEAAFTKNQKDPGVLDKRLD
1j55A      --MTELETAMGMIIDVFSRYSSEGSTQTLTKGELKVLMEKELPGFLD-----AVDKLL
1ig5A      -----KSPPEELKGI FEKYAAKEGDPNQLSKEELKLLQLTEFFPSLLKGPSTLDELFE
1qx2A      -----MKSPPEELKGAPEVFAAKEGDPNQLSKEELKLVQMQLGPGSLLKGMSTLDEMI
          . : * . : : : * . : : * * : : . : . : . : : : : : : : : * : :

lnshA      KKLDLNSDGQLDFQKQFLNLIIGLAVACHESFVKAAPPQKRF
1j55A      KDLNDANGDAQVDFSEFIVFAAITSACHKYFEKAGL-----
1ig5A      EELDKNGDGEVSFEFQVLVKKISQ-----
1qx2A      EEVDKNGDGEVSFEFLVMMKISQ-----
          . : * * * . : : * * * : : : : :
```

マルチプルアライメントの目的

```
lnshA      SRPTEATERCIESLIAVPQKYAGKDGHSVTLKTEFLSPMNTLEAAFTKNQKDPGVLDKRLD
1j55A      --MTELETAMGMIIDVFSRYSSEGSTQTLTKGELKVLMEKELPGFLD-----AVDKLL
1ig5A      -----KSPPEELKGI FEKYAAKEGDPNQLSKEELKLLQLTEFFPSLLKGPSTLDELFE
1qx2A      -----MKSPPEELKGAPEVFAAKEGDPNQLSKEELKLVQMQLGPGSLLKGMSTLDEMI
          . : * . : : : * . : : * * : : . : . : . : : : : : : : : * : :

lnshA      KKLDLNSDGQLDFQKQFLNLIIGLAVACHESFVKAAPPQKRF
1j55A      KDLNDANGDAQVDFSEFIVFAAITSACHKYFEKAGL-----
1ig5A      EELDKNGDGEVSFEFQVLVKKISQ-----
1qx2A      EEVDKNGDGEVSFEFLVMMKISQ-----
          . : * * * . : : * * * : : : : :
```

- ファミリー内の機能的な重要部位の検出
- ファミリーを特徴付けるモチーフの発見
- プロフィール法による遠縁のホモログ発見
- 分子系統解析の第一ステップとして不可欠
- 進化的追跡法(evolutionary trace method)

多重整列のスコア

(1) SP (sum-of-pairs) スコア

複数の文字列間のスコアを
ペアワイズのアミノ酸置換スコア $s(a,b)$ の和で表す

$$S(m_i) = \sum_{k < l} s(m_i^k, m_i^l)$$

m_i^k : k 番目の配列の i 番目の文字

```
RCIAVF
TAMDVF
KSPGIF
```

$S(m_i) = s(R,T) + s(T,K) + s(R,K)$

理論的にはおかしい: $S(A,B) + S(B,C) + S(A,C) = \log \frac{P(A,B)P(B,C)P(A,C)}{P(A)^2 P(B)^2 P(C)^2} \neq \log \frac{P(A,B,C)}{P(A)P(B)P(C)}$

BLOSUM62

A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	B	Z	X	*	
A	4	-1	-2	-2	0	-1	-1	0	-2	-1	-1	-1	-2	-1	1	0	-3	-2	0	-2	-1	0	-4	
R	-1	5	0	-2	-3	1	0	-2	0	-3	-2	-1	-3	-2	-1	-3	-2	-3	-1	0	-1	-4		
N	-2	0	6	1	-3	0	0	1	-3	-3	0	-2	-3	-2	1	0	-4	-2	-3	3	0	-1	-4	
D	-2	-2	1	6	-3	0	2	-1	-1	-3	-4	-1	-3	-3	-1	0	-1	-4	-3	-3	4	1	-1	-4
C	0	-3	-3	-3	9	-3	-4	-3	-3	-1	-1	-3	-1	-2	-3	-1	-1	-2	-2	-1	-3	-3	-2	-4
Q	-1	1	0	0	-3	5	2	-2	0	-3	-2	1	0	-3	-1	0	-1	-2	-1	-2	0	3	-1	-4
E	-1	0	0	2	-4	2	5	-2	0	-3	-3	1	-2	-3	-1	0	-1	-3	-2	-2	1	4	-1	-4
G	0	-2	0	-1	-3	-2	-2	6	-2	-4	-4	-2	-3	-3	-2	0	-2	-2	-3	-3	-1	-2	-1	-4
H	-2	0	1	-1	-3	0	0	-2	8	-3	-3	-1	-2	-1	-2	-1	-2	-2	2	3	0	0	-1	-4
I	-1	-3	-3	-1	-3	-3	-4	-3	4	2	3	1	0	-3	-2	-1	-3	-1	3	-3	-3	-1	-4	
L	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4	-2	2	0	-3	-2	-1	-2	-1	1	-4	-3	-1	-4
K	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5	-1	-3	-1	0	-1	-3	-2	-2	0	1	-1	-4
M	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5	0	-2	-1	-1	-1	-1	1	-3	-1	-1	-4
F	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6	-4	-2	-2	1	3	-1	-3	-3	-1	-4
P	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7	-1	-1	-4	-3	-2	-2	-1	-2	-4
S	0	-1	1	0	-1	0	0	-1	-2	-2	0	-1	-2	-1	4	1	-3	-2	-2	0	0	0	0	-4
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-2	-1	1	5	-2	-2	0	-1	-1	0	0	-4
W	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-3	-1	1	-4	-3	-2	11	2	-3	-4	-3	-2	-4	-4
Y	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-2	-1	3	-3	-2	-2	2	7	-1	-3	-2	-1	-4	-4
V	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4	-3	-2	-1	-4
B	-2	-1	3	4	-3	0	1	-1	0	-3	-4	0	-3	-2	0	-1	-4	-3	-3	4	1	-1	-4	
Z	-1	0	0	-1	-3	3	4	-2	0	-3	-3	1	-1	-3	-1	0	-1	-3	-2	-2	1	4	-1	-4
X	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	-1	-1	-1	-1	-2	0	0	-2	-1	-1	-1	-1	-1	-4
*	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	1

多重配列のスコア(続き)

(2) 配列への重み付きのSum-of-pair関数 (ClustalW)

$$S(m_i) = \sum_{k < l} w_k \cdot w_l \cdot s(m_i^k, m_i^l)$$



(3) エントロピー関数の最小化

各サイトのアミノ酸の頻度 $p_i(a)$ を推定し、そのエントロピーの和を求める

$$S(m_i) = - \sum p_i(a) \log p_i(a)$$

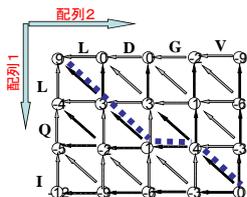
サイト	$P_i(a)$	$S(m_i)$
1	$P_i(L)=1.0$	0.00
2	$P_i(G)=0.5, P_i(A)=0.5$	0.69
3	$P_i(V)=0.25, P_i(I)=0.25, P_i(A)=0.5$	1.04

(4) 対アライメントライブラリの重複による部位特異的スコア (T-COFFEE)

どうやって並べるか?

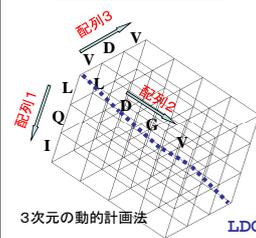
多次元DPIによる多重配列の厳密解

2本の配列のアライメント



2次元の動的計画法
LDGV
LQ-I
メモリ・計算時間 $O(L^2)$

3本の配列のアライメント



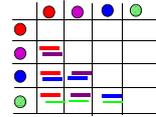
3次元の動的計画法
LDGV
LQ-I
VD-V
メモリ・計算時間 $O(L^3)$

N本の配列のアライメントのメモリ・計算時間は $O(L^N)$ → 非現実的
長さ100の2本のアライメントが1秒でできても、10本に増やすと 100^8 秒かかる。

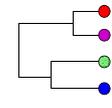
プログレッシブ・アライメント (progressive alignment, 累進法)

Feng and Doolittle (1987)

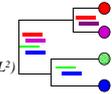
(1) 全ての配列ペアのペアワイスアライメントを計算する



(2) ペアワイスアライメントによる距離行列を計算し、樹形図を計算する。



(3) 樹形図の葉から、ペアワイスアライメントを組み上げていく



ステップ1に最も計算時間がかかる。全体の計算量はほぼ $O(NL^2)$

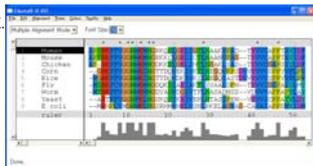
ClustalW / ClustalX

UNIX/Windows/Mac版: <http://ftp.ebi.ac.uk/pub/software/clustalw2>
WEBサーバ: <http://www.ebi.ac.uk/Tools/clustalw2>

- 現在、最も一般的な多重配列のプログラム
- アルゴリズムは累進法。ペアワイスアライメントはグローバルアライメントを用い、ガイド木はNJ法で作成。スコアは配列の重みを導入したSum-of-pairs。置換スコア行列の選択、ギャップペナルティ等に様々な経験的な工夫が見られる。

- GUI版はClustalW, GUI版はClustalX. UNIX, Windows, MACでも動作する。

- NJ法による系統樹計算機能付き。



Thompson, J.D., Higgins, D.G., Gibson T.J. "CLUSTALW: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice". Nucleic Acids Research, 1994, 22, 4673-4680.

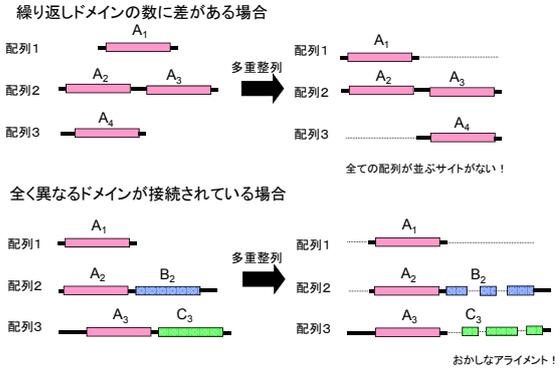
主要なマルチプルアライメントのプログラム

	WEBサイト	アルゴリズム	特徴
ClustalW・ClustalX	http://www.ebi.ac.uk/Tools/clustalw2	累進法。重み付きSPスコアを使用。置換スコア行列の選択、ギャップペナルティ等に様々な工夫	もっとも広く使われている標準的なプログラム
T-COFFEE	http://www.ebi.ac.uk/t-coffee/	ペアワイスアライメントをローカル、グローバル、進展を用いて多数生成。それらの集合から、位置特異的スコアを作成し、累進法を実行する。	計算時間がかかるが精度は高い。配列の本数が100本以下の場合に向いている。
MAFFT	http://align.bmr.kyushu-u.ac.jp/mafft/online/server/	高速フーリエ変換(FFT)を用いて、高速にペアワイスアライメントを実装、それを利用して、累進法、あるいは反復改善法を実行する。	計算時間は高速なので、配列の本数が100~500本程度でも、計算可能。

マルチプルアライメントを行う上での注意点

- 対象とする配列群が相同であることの確認
 - 他と全く似ていない配列が混入していると意味のない比較になる
- 対象とする配列群のほぼ全長どうしが対応することの確認
 - ClustalW等主要な多重整列プログラムはグローバルアライメントなので、全長どうしに対応することがアルゴリズムの前提
 - マルチドメイン構造、繰り返し構造になっていないかを確認
 - そもそも、配列長が著しく異なる場合は、ほぼ間違いない問題が生じる
 - 配列の一部しか、対応しないなら、その部分だけ切り出して入力する
- 計算されたマルチプルアライメントの結果の吟味
 - 既知の機能部位がきちんと保存されているか
 - 長すぎるギャップはないか(マルチドメインの可能性)
 - 保存部位が、非保存の配列はないか(ホモログでない可能性)
 - 立体構造が既知のものが含まれているなら、立体構造アライメントも参照

マルチドメインのときのアライメントの問題点



マルチプルアライメントから何を読み取るか?

```

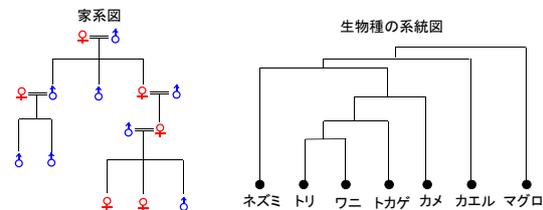
5p21- MTEYKLVVVGAGGVGKSALTIQLIQNHFVDEYDPTIEDSY
1ctqA MTEYKLVVVGAGGVGKSALTIQLIQNHFVDEYDPTIEDSY
1c1yA MREYKLVVVGSGGVGKSALTVQFVQGI FVEKYDPTIEDSY
1kao- MREYKVVVLGSGGVGKSALTVQFVTGTFIEKYDPTIEDFY
1huqA --QFKLVLLGESAVGKSSLVLRFRVKGQFHEYQESTIGAAF
1g16A ----KILLIGDSGVGKSCLLVRFVE-- --DKFNPI--DFK
1ek0A VTSIKLVLLGEAAVGGKSSIVLRFVSNDFEAENKEPTIGAAF
3rabA ---FKLLIQNSVVGKTSFLFRYADDSFTPAFVSTVGIDF
1mh1- ----KCVVVG DGAVGKTCLLISYTTNAFPGEYIPTVFDNY
2ngrA MQTIKCVVVG DGAVGKTCLLISYTTNKFPSEYVPTVFDNY
1tx4B ----KLVI VGDGACGKTCLLIVNSKDQE---YVPTVFENY
1i2mA --QFKLVLVGDGGTGKTFVVRHLKKYVATEVHPLVFHTN
1d5cA --KYKLVFLGEQAVGKTSI-ITRFYDTFDNNYQSTIGDFL
    
```

サイトごとに保存の度合いに差がある。 [AG]-x(4)-G-K-[ST]_
 サイトごとにアミノ酸の出現傾向に差がある

分子系統学基礎

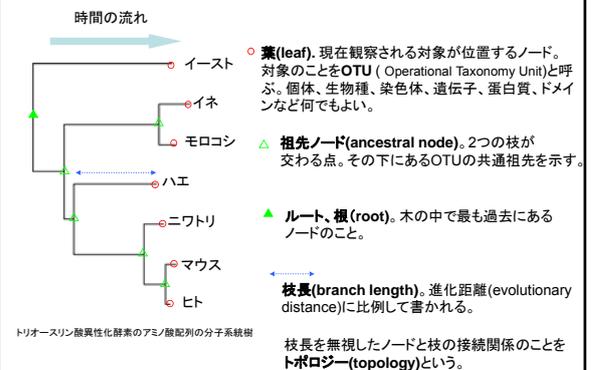
系統樹(phylogenetic tree)

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

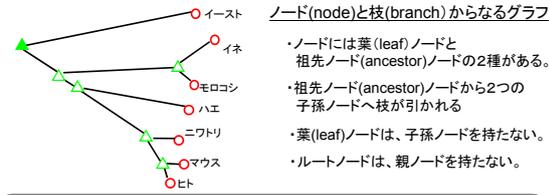


- 何を対象にするかはいろいろ(個体、生物種、染色体、遺伝子)
- 「系統樹を書く」→「過去(歴史)を推定する」
- 「分類」(似ているものをまとめること)と「系統推定」の手続きは似ている
- 様々な「分類法」が在り得るが、「系統樹」には唯一の歴史的真相があるはず。

系統樹の用語



系統樹(二分岐樹)のデータ構造

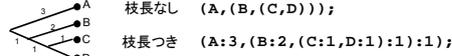


```

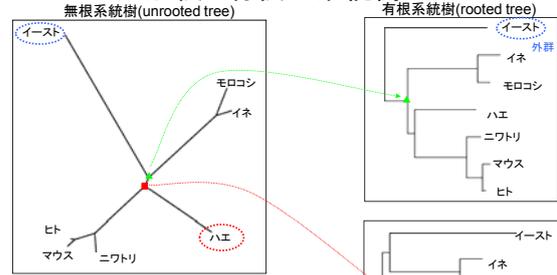
各ノードが、2つの子ノードへのポインタと、枝長を持つ。
struct NODE{
  struct NODE *child1,*child2;
  double len1, len2;
};
    
```

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

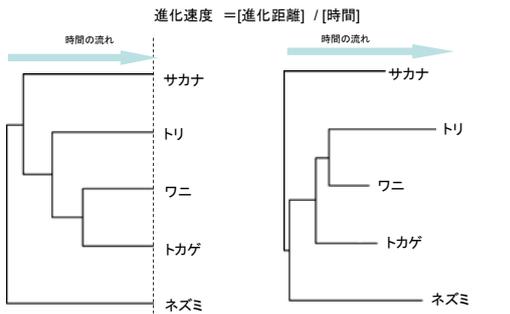


無根と有根の系統樹



- ・NJ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
- ・どの枝に根を置くかによって、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



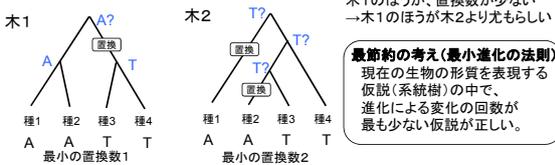
分子配列からの系統樹の推定法

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)



- (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



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_k, 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

最節約法のアルゴリズムのキーポイント

子ノードの R_i, R_j が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

「 \cap 」、「 \cup 」、「空である」:などは集合の用語

$A \cap B$: 積集合。共通部分。2つの集合A, Bの共通要素

例 $(a, b, c) \cap (b, c, d) = (b, c)$, $(a, b, c) \cap (a) = (a)$, $(a) \cap (b) = \text{空}$

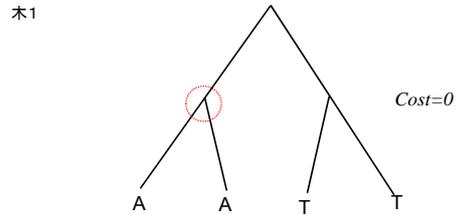
$A \cup B$: 和集合。合併集合。2つの集合A, Bのどちらかに属する要素

例 $(a, b, c) \cup (b, c, d) = (a, b, c, d)$, $(a, b, c) \cup (a) = (a, b, c, d)$, $(a) \cup (b) = (a, b)$

Aが空である: 集合Aに属する要素が一つもないこと。

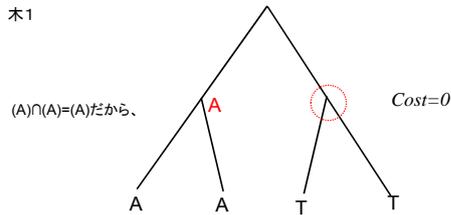
置換数の推定の例:木1(1)

子ノードの R_i, R_j が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



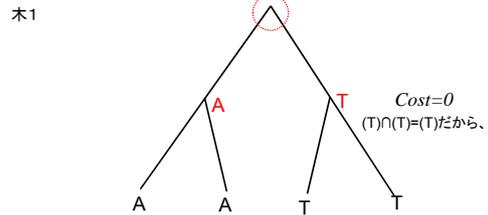
置換数の推定の例:木1(2)

子ノードの R_i, R_j が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



置換数の推定の例:木1(3)

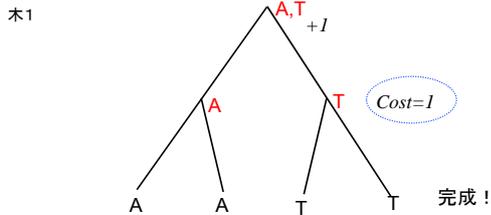
子ノードの R_i, R_j が計算されているなら、以下のように R_k を計算
 $R_i \cap R_j$ が空でないなら、 $R_k = R_i \cap R_j$
 $R_i \cap R_j$ が空なら、 $R_k = R_i \cup R_j$, Costに1加算



置換数の推定の例:木1(4)

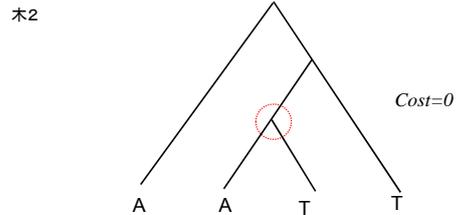
子ノードの R_i, R_j が計算されているなら、以下のように R_k を計算
 $R_i \cap R_j$ が空でないなら、 $R_k = R_i \cap R_j$
 $R_i \cap R_j$ が空なら、 $R_k = R_i \cup R_j$, Costに1加算

$(A) \cap (T) = \text{空}$ だから、 $(A) \cup (T) = (A, T)$ を祖先形質とする。コストを1増やす



置換数の推定の例:木2(1)

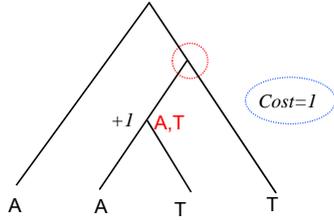
子ノードの R_i, R_j が計算されているなら、以下のように R_k を計算
 $R_i \cap R_j$ が空でないなら、 $R_k = R_i \cap R_j$
 $R_i \cap R_j$ が空なら、 $R_k = R_i \cup R_j$, Costに1加算



置換数の推定の例:木2(2)

子ノードの R_i, R_j が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

木2

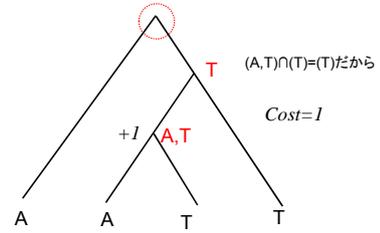


(A)∩(T)=空だから、(A)∪(T)=(A,T)を祖先形質とする。コストを1増やす

置換数の推定の例:木2(3)

子ノードの R_i, R_j が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

木2

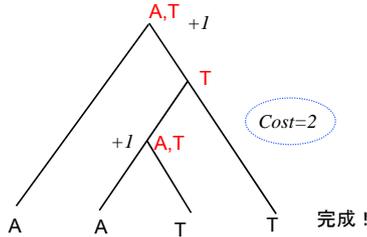


置換数の推定の例:木2(4)

子ノードの R_i, R_j が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

(A)∩(T)=空だから、(A)∪(T)=(A,T)を祖先形質とする。コストを1増やす。

木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

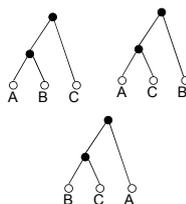
可能な木のトポロジーの数

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

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



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

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	AAAA T
配列 3	T AA T A
配列 4	T AA T T

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

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

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

$d_{12} \doteq L_{1a} + L_{2a}$ $d_{34} \doteq L_{3b} + L_{4b}$
 $d_{13} \doteq L_{1a} + L_{ab} + L_{3b}$ $d_{14} \doteq L_{1a} + L_{ab} + L_{4b}$
 $d_{23} \doteq L_{2a} + L_{ab} + L_{3b}$ $d_{24} \doteq L_{2a} + L_{ab} + L_{4b}$

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 p-距離 = 4 / 10 = 0.4
GGVVSTLVA

多重置換の影響を考慮した距離

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。

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

木村の距離 = $-\log(1-p-0.2p^2)$

0	:AAAAA	0.0
1	:AKAAAA	0.1
2	:PKAAAA	0.2
3	:PKAAMAAA	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	:NKVAADWTR	0.7

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。それぞれの配列が一つのクラスター C_i を構成するとする。

[反復]

(1) 全てのクラスターのペアの中で距離 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ー $C_k = C_i \cup C_j$ 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。クラスター間の距離は、属する配列間の平均距離で定義する。

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

クラスター数が1つになるまで反復する。

重心間距離を用いたクラスター解析と同じ

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

	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

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

クラスターと配列の距離は、配列間平均の距離とする

クラスターと配列の距離は、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(2)

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

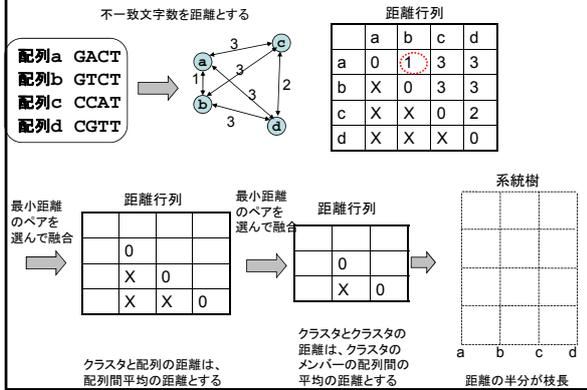
クラスターと配列の距離は、配列間平均の距離とする

クラスターと配列の距離は、配列間平均の距離とする

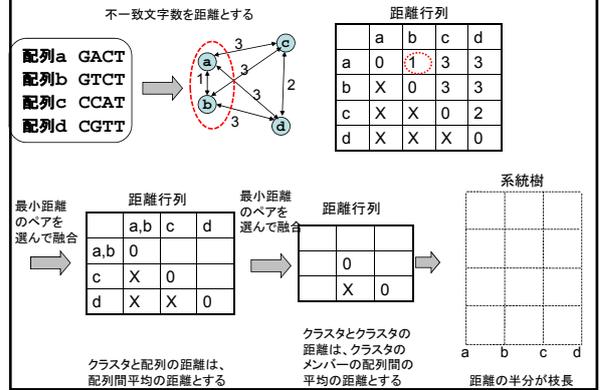
系統樹

距離の半分が枝長

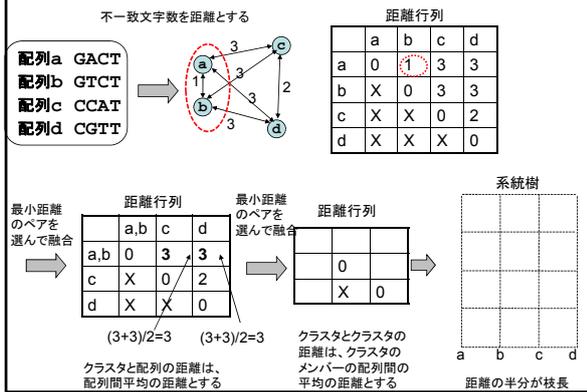
UPGMA法による系統樹の計算例(3)



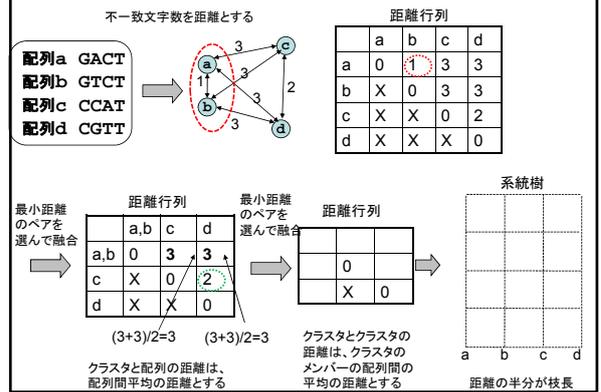
UPGMA法による系統樹の計算例(4)



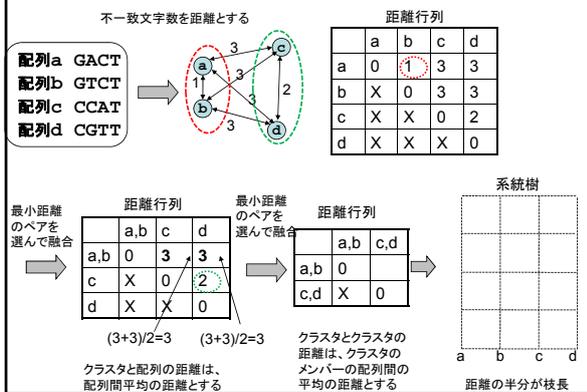
UPGMA法による系統樹の計算例(5)



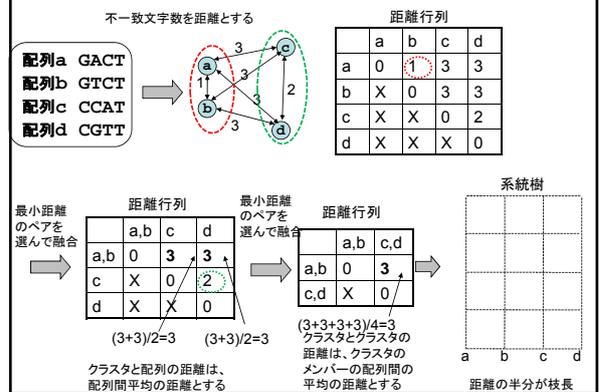
UPGMA法による系統樹の計算例(6)



UPGMA法による系統樹の計算例(7)



UPGMA法による系統樹の計算例(8)



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

	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

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

クラスタと配列の距離は、配列間平均の距離とする

距離行列

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

$(3+3+3+3)/4=3$
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

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、

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

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

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

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。

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。
 $r_i = \frac{1}{|L|} \sum_{m \in L} d_{im}$ 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、
 $d_{mk} = (d_{im} + d_{jm}) / 2$
 $d_{jk} = (d_{ij} + d_{im} - d_{jm}) / 2$
 $d_{ik} = (d_{ij} + d_{jm} - d_{im}) / 2$
で定義。ただし、木の枝長となる d_{ij}, 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。
 $d_{ij} = \langle (d_{ij} + d_{im} - d_{jm}) / 2 \rangle_m = (d_{ij} + r_i - r_j) / 2$
 $d_{jk} = \langle (d_{ij} + d_{jm} - d_{im}) / 2 \rangle_m = (d_{ij} + r_j - r_i) / 2$

[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。

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
	torii	9.5	10.5	6.8	5.5	0.0

UPGMA法

NJ法(無根)

NJ法(有根)

外群の選択

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。

木1

$P_{ab}(t)$: 時間 t の間に a から b に変異する確率

木1が起る確率 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。

アライメント

```

12345678
a: AGAAAAAC
b: AGACATGC
c: TATCGACA
d: TAAAGTGA

```

ブートストラップ抽出データ1

```

26175763
a: GAAAAAAA
b: GTAGAGTA
c: AATCGCAT
d: ATTGGGTA

```

ブートストラップ抽出データ2

```

14735128
a: AAAAAAGC
b: ACGAAAGC
c: TCCTGTAA
d: TAGAGTAA

```

それぞれのブートストラップ抽出したデータに対して系統樹を作成。(a,b),(c,d)のトポロジーが作成された回数を数える

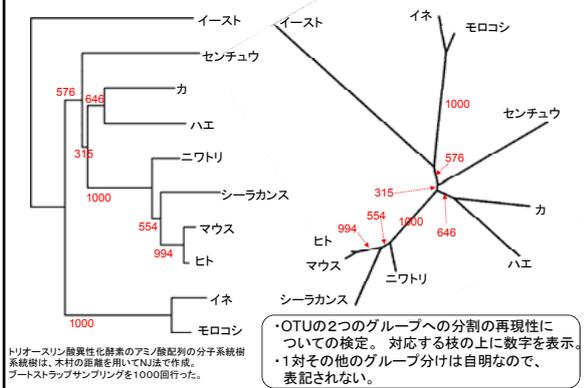
系統樹

確認したい信頼性

(1) 十分な数のサイトがあるか
(2) 全てのサイトが同じ系統樹を示唆するか

1000個のブートストラップ抽出データのうち、860個について、このトポロジーが再現。

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
マルチプリ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に对应。
- 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に对应。
- 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

- 金久實 著「ポストゲノム情報への招待」(2001) 共立出版
- Arthur M.Lesk(岡崎康司、坊農秀雄 監訳)「バイオインフォマティクス基礎講義 一歩進んだ発想をみがぐために」(2003)、メディカル・サイエンス・インターナショナル
- 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