

Kakusan3を用いた分子進化モデルの選択

モデル選択って何?

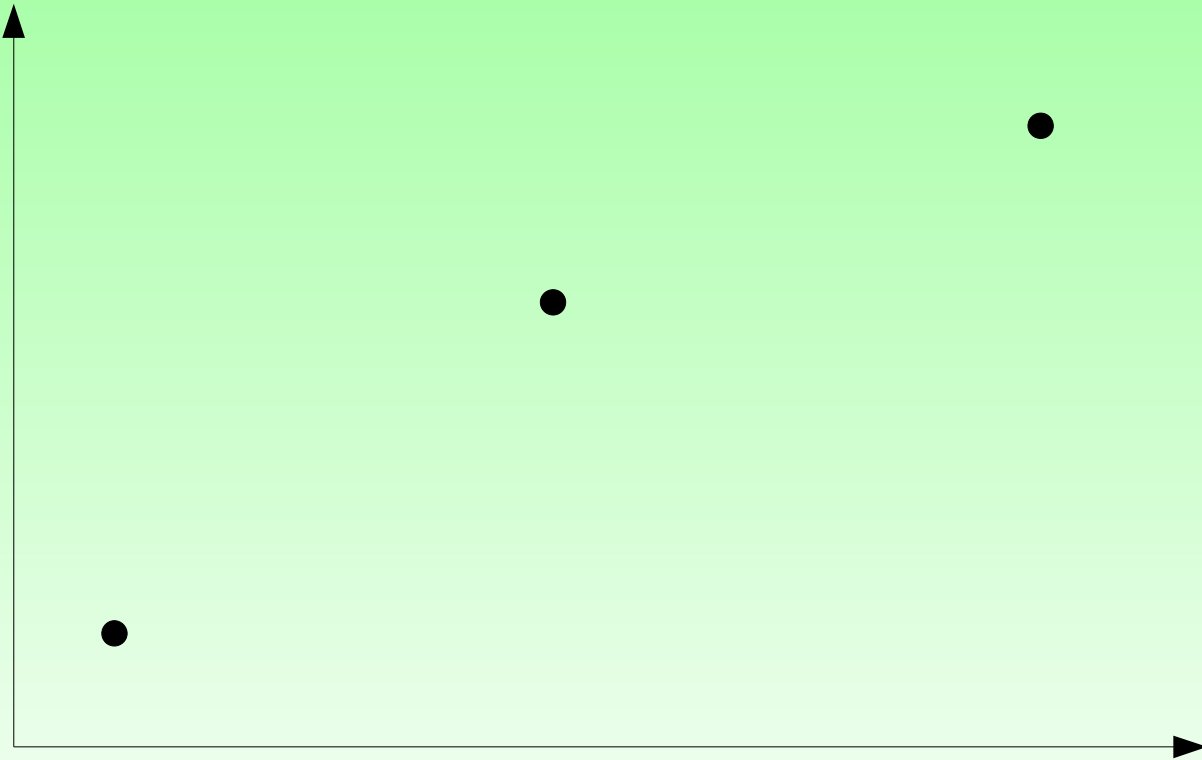
進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

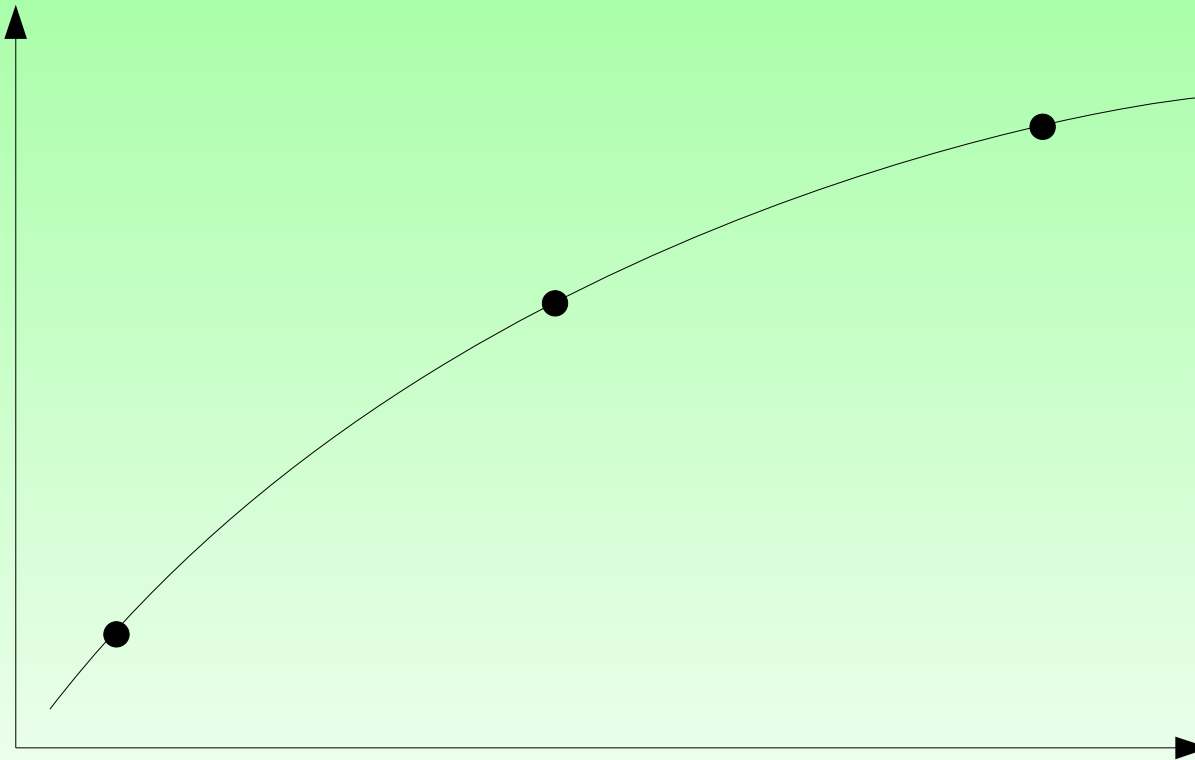
2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

モデル選択って何?

- 3点のデータに対する回帰分析を考える

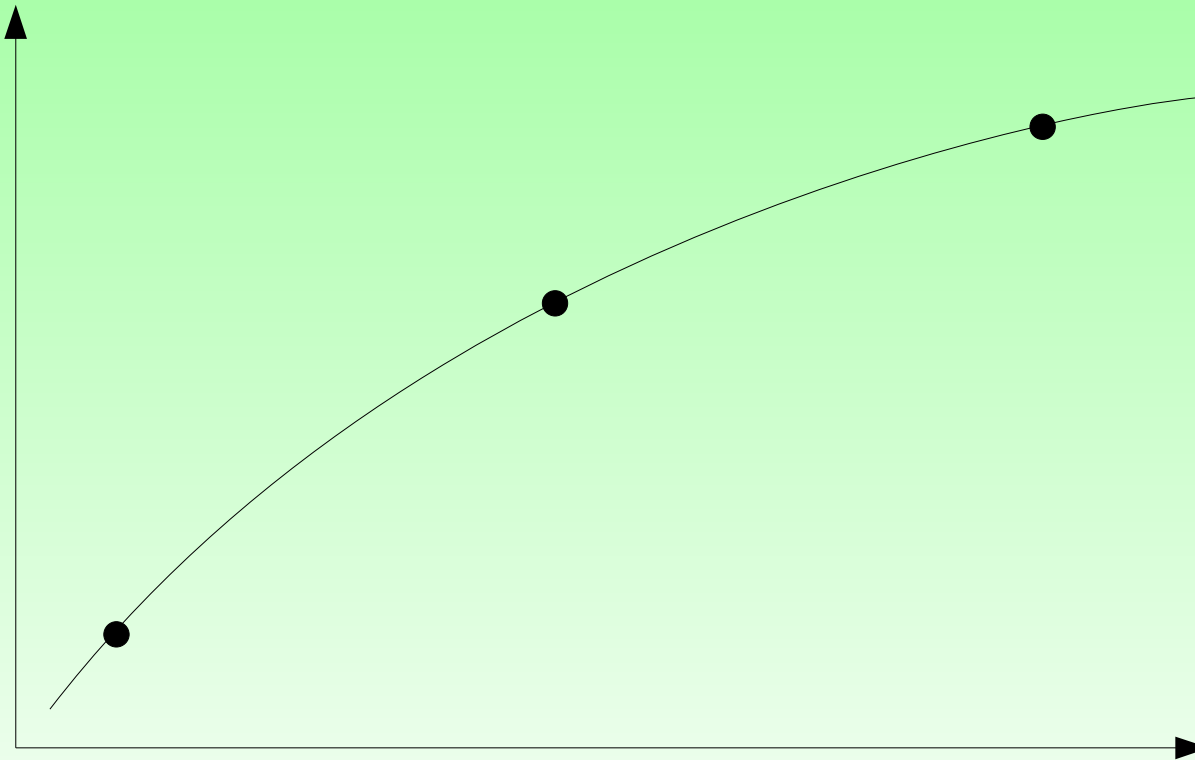


モデル選択って何?



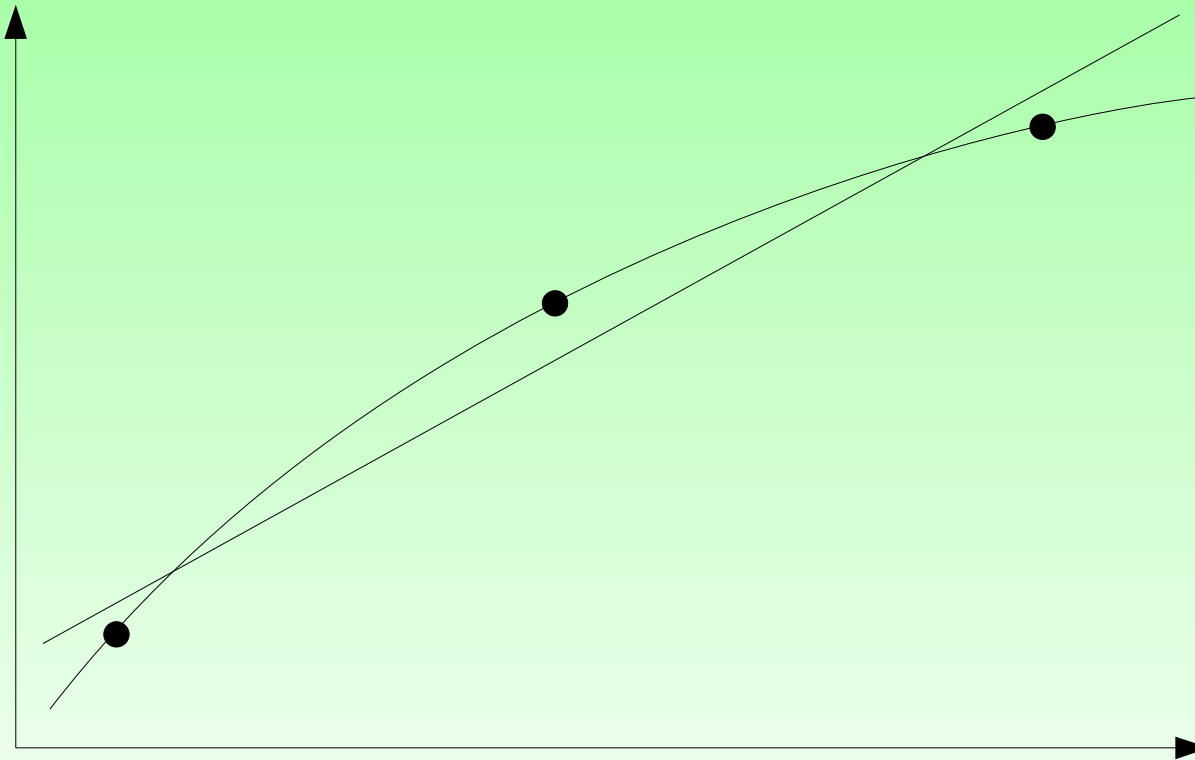
- 3点のデータに対する回帰分析を考える
- 2次曲線を当てはめれば全点を通せる = 尤度最大

モデル選択って何?



- 3点のデータに対する回帰分析を考える
- 2次曲線を当てはめれば全点を通せる = 尤度最大
- ただしデータから推定するパラメータは3つ

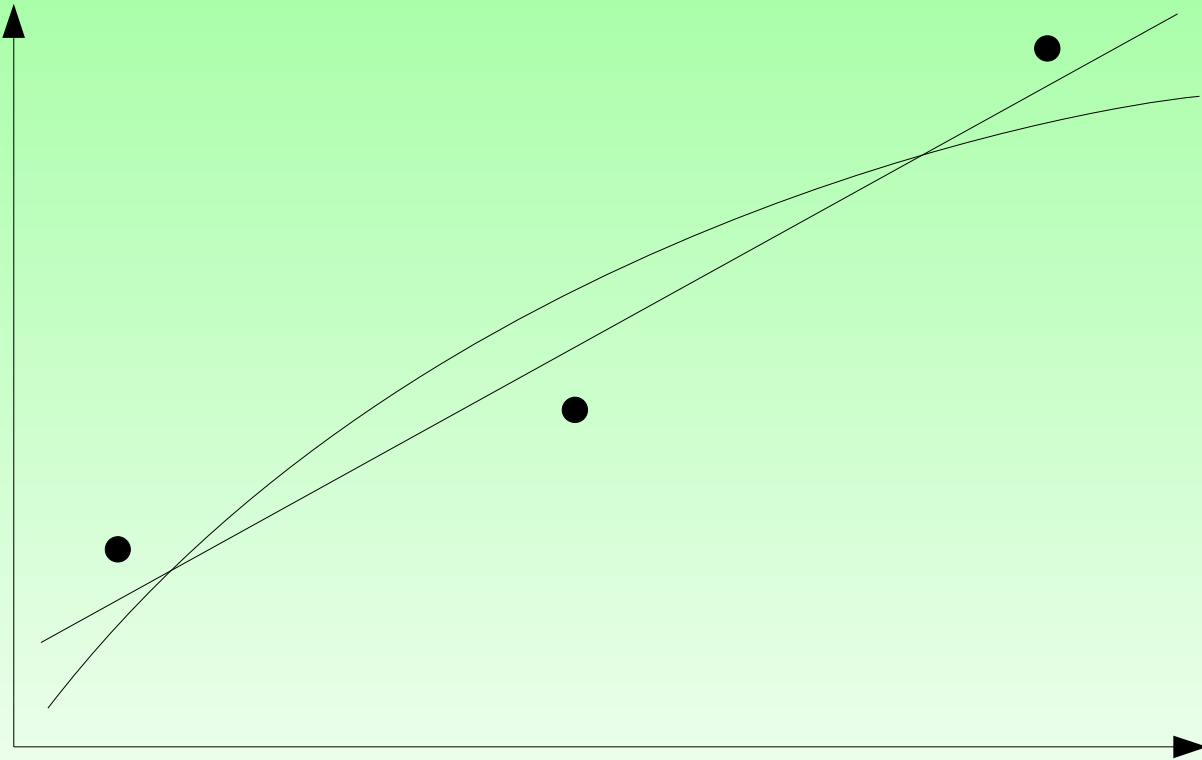
モデル選択って何?



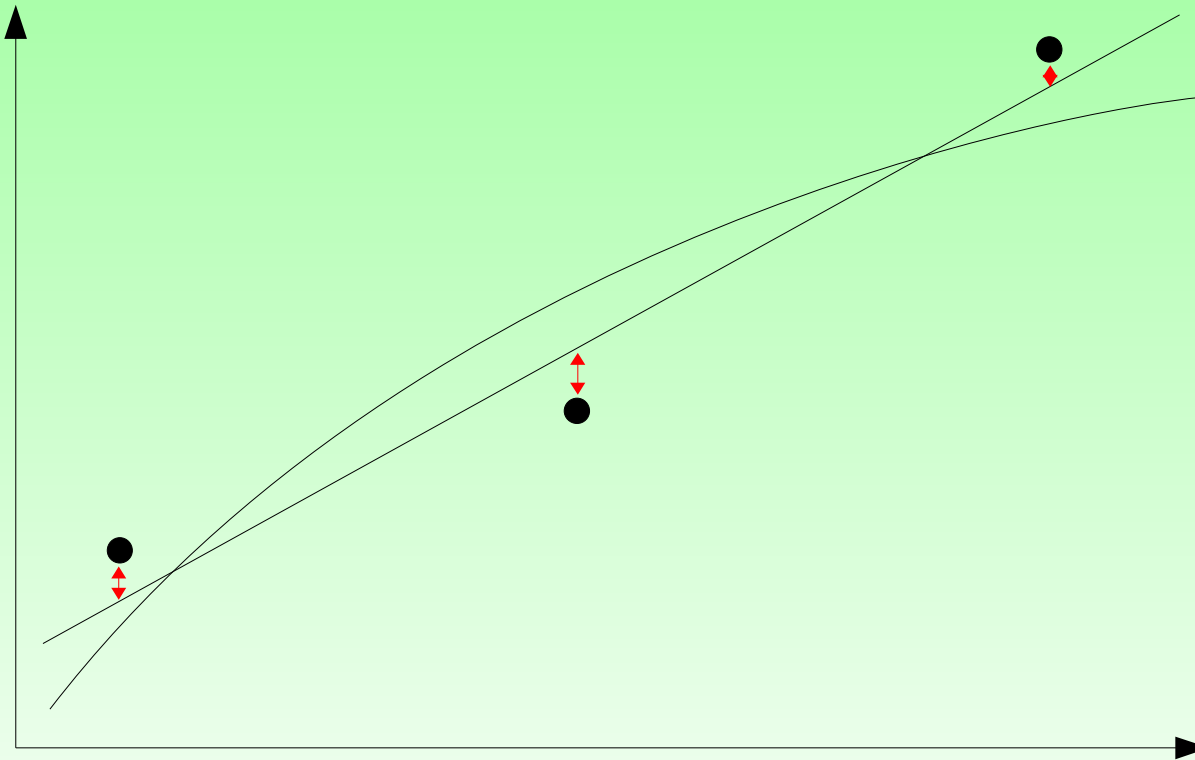
- 3点のデータに対する回帰分析を考える
- 2次曲線を当てはめれば全点を通せる=尤度最大
- ただしデータから推定するパラメータは3つ
- 1次直線を当てはめれば尤度は低下するがパラメータ数は減少する

モデル選択って何?

- データを取り直した場合を考える

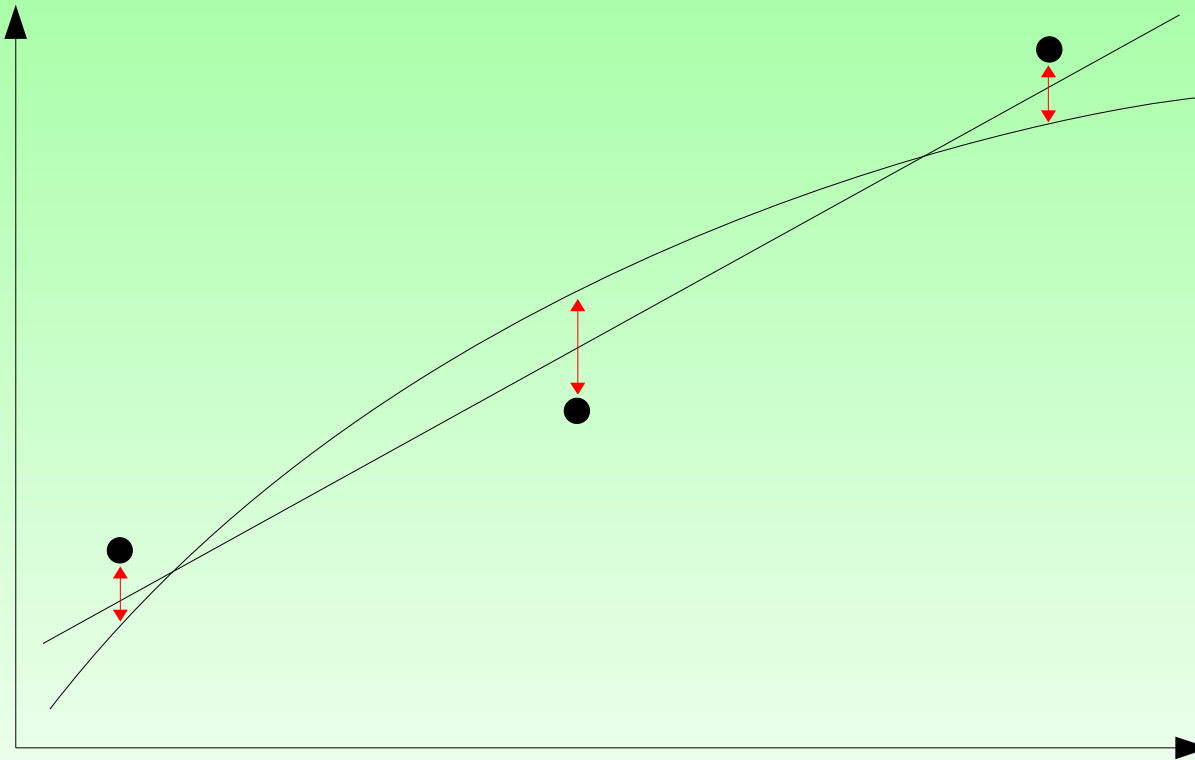


モデル選択って何?



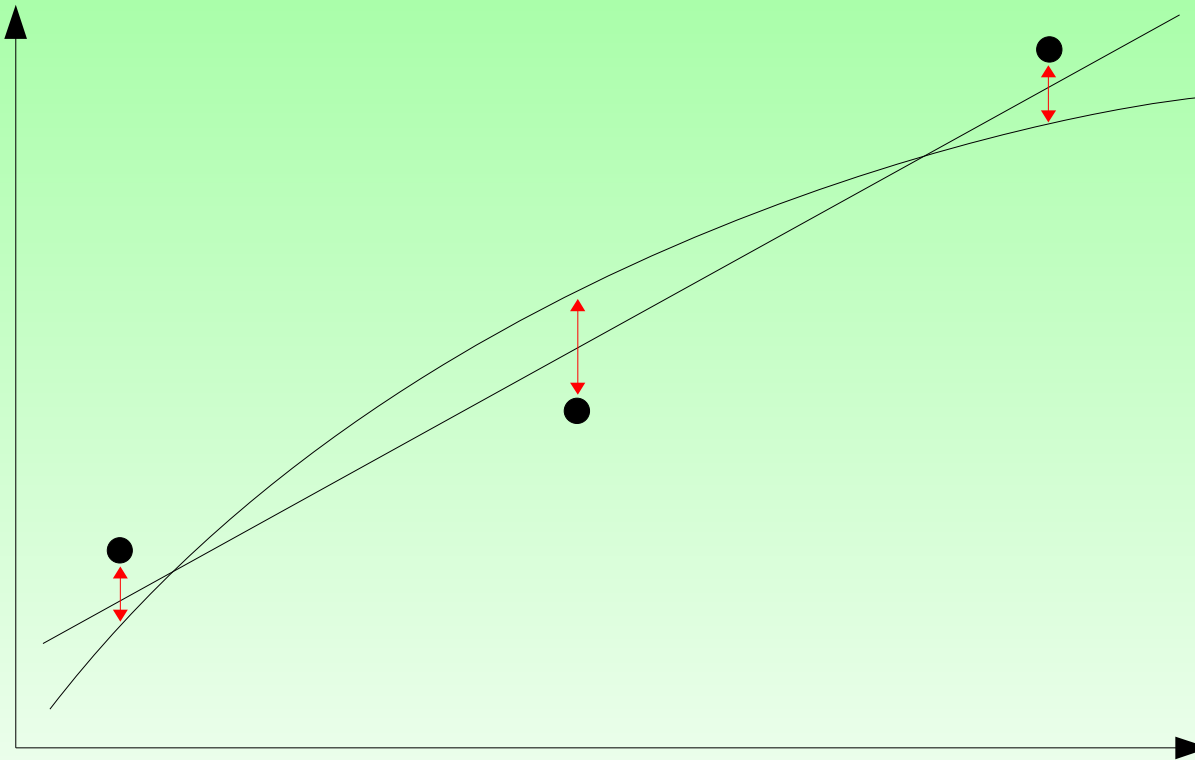
- データを取り直した場合を考える
- 1次直線は当てはまりが大きく変化しない

モデル選択って何?



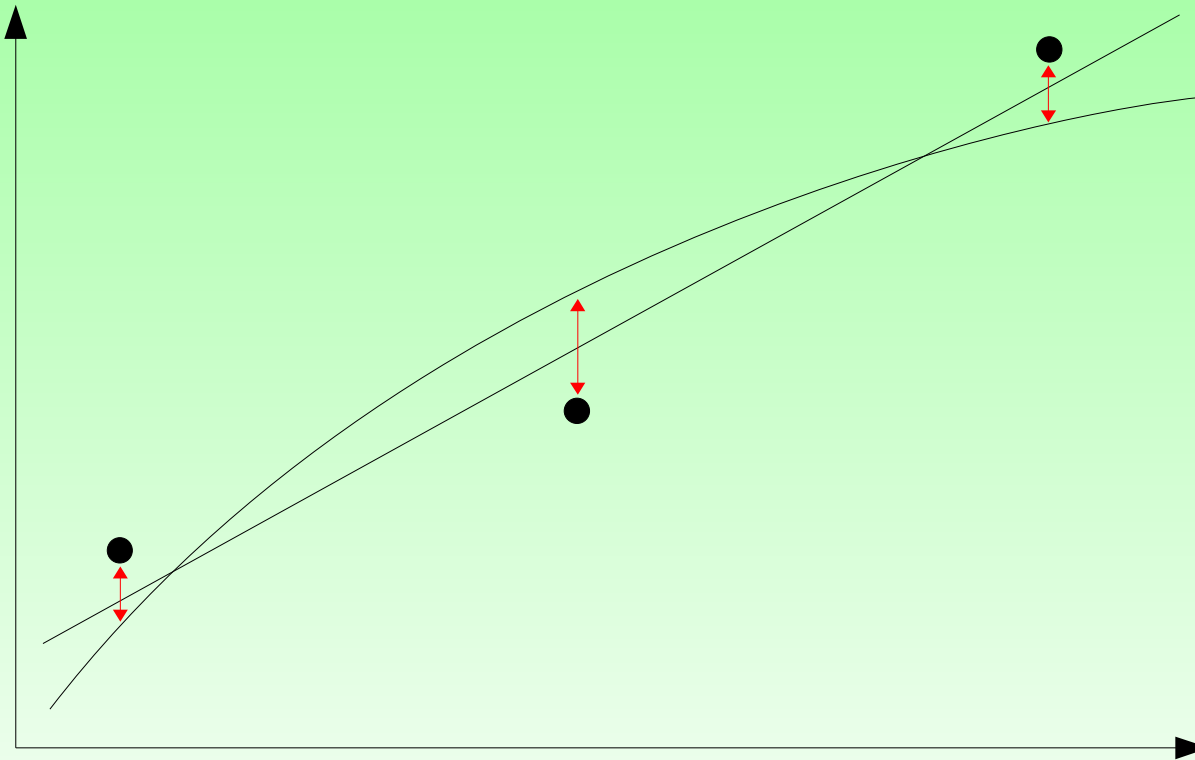
- データを取り直した場合を考える
- 1次直線は当てはまりが大きく変化しない
- 2次曲線の当てはまりは大きく低下する

モデル選択って何?



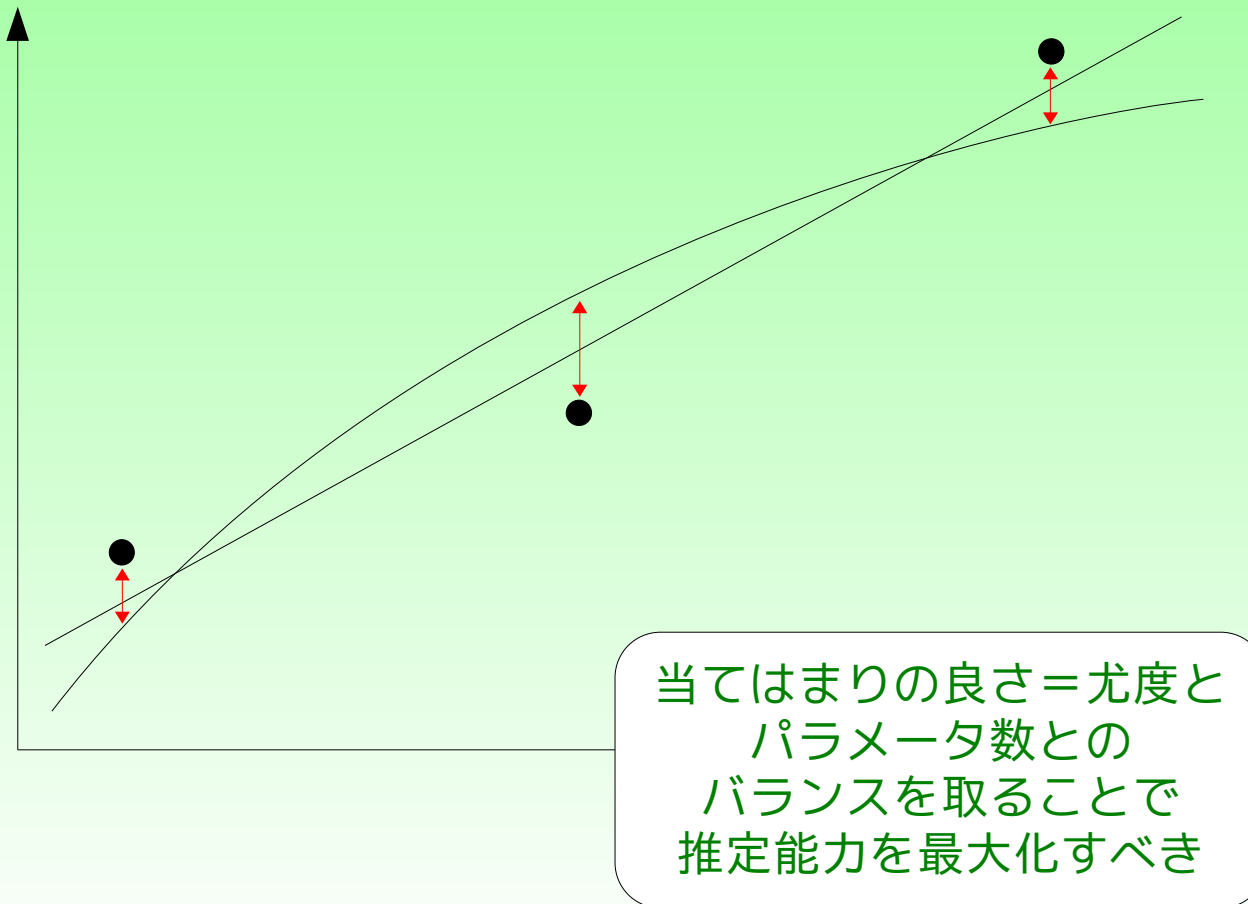
- データを取り直した場合を考える
- 1次直線は当てはまりが大きく変化しない
- 2次曲線の当てはまりは大きく低下する
- 1次直線の方が「母集団のパラメータを推定する能力」が高い

モデル選択って何?



- データを取り直した場合を考える
- 1次直線は当てはまりが大きく変化しない
- 2次曲線の当てはまりは大きく低下する
- 1次直線の方が「母集団のパラメータを推定する能力」が高い
- 2次曲線はデータの「ノイズ」にまでフィットすることで「推定能力」が低下している

モデル選択って何?



- データを取り直した場合を考える
- 1次直線は当てはまりが大きく変化しない
- 2次曲線の当てはまりは大きく低下する
- 1次直線の方が「母集団のパラメータを推定する能力」が高い
- 2次曲線はデータの「ノイズ」にまでフィットすることで「推定能力」が低下している

赤池情報量規準

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

赤池情報量規準

$$AIC = -2 \ln L + 2K$$

(L は最大化尤度・ K はパラメータ数)

$$AIC = -2 \ln L + 2K$$

(L は最大化尤度・ K はパラメータ数)

AIC最小となるモデル(回帰式)が
ベストな = 推定能力が最大化されるモデル

AICcとかBICなどもしばしば使われます

分子系統推定で用いられるデータ

site	12345M
OTU1	TGTTT	...	TTTTTC
OTU2	AGTAC	...	TTTTTC
OTU3	AGTAT	...	TTGTC
⋮	⋮		⋮
OTUN	AGTAT	...	ATTTC

分子系統推定で用いられるデータ

site	12345M
OTU1	TGTTT	..	TTTTTC
OTU2	AGTAC	..	TTTTTC
OTU3	AGTAT	..	TTGTC
⋮	⋮		⋮
OTUN	AGTAT	..	ATTTC

OTU名

分子系統推定で用いられるデータ

site	12345M
OTU1	TGTTT	..	TTTTTC
OTU2	AGTAC	..	TTTTTC
OTU3	AGTAT	..	TTGTC
⋮	⋮		⋮
OTUN	AGTAT	...	ATTTC

OTU名

アライメント済みデータ配列

分子系統推定で用いられるデータ

site	1	2	3	4	5	M
OTU1									TTC
OTU2									TTC
OTU3									GTC
⋮									⋮
⋮									⋮
⋮									⋮
OTUN			AGTAT			...			ATTTC

このデータマトリックス内の「変異」を利用して最適な系統樹を選択する
しかし、変異の起きやすさはタイプやサイトによって異なる

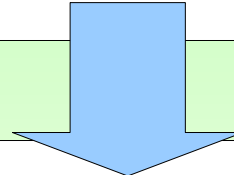
OTU名

アライメント済みデータ配列

分子系統推定で用いられるデータ

site	1	2	3	4	5	M
OTU1	T	C	T	T	T	T	T	T	TTC
OTU2	T	T	T	T	T	T	T	T	TTC
OTU3	T	T	T	T	T	T	T	T	GTC
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

このデータマトリックス内の「変異」を利用して最適な系統樹を選択する
しかし、変異の起きやすさはタイプやサイトによって異なる



分子進化モデルによるモデル化

OTU名

アライメント済みデータ配列

分子進化モデルがモデル化しているもの

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

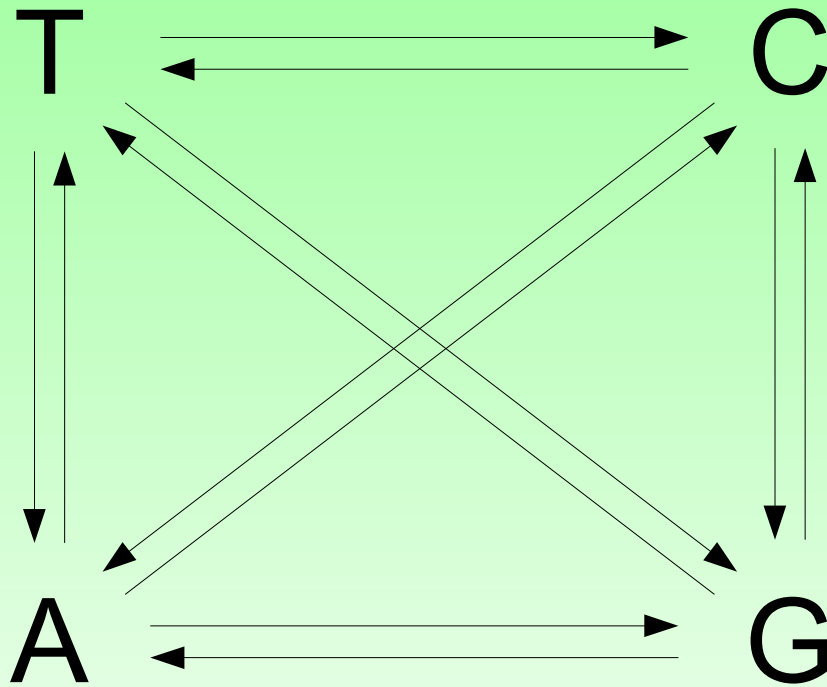
分子進化モデルがモデル化しているもの

- サイト内形質状態間の置換速度不均質性
 - どのタイプの変異が起きやすいかを表す

分子進化モデルがモデル化しているもの

- サイト内形質状態間の置換速度不均質性
 - どのタイプの変異が起きやすいかを表す
- サイト間の置換速度不均質性
 - どこで変異が起きやすいかを表す

DNAの塩基置換パターン



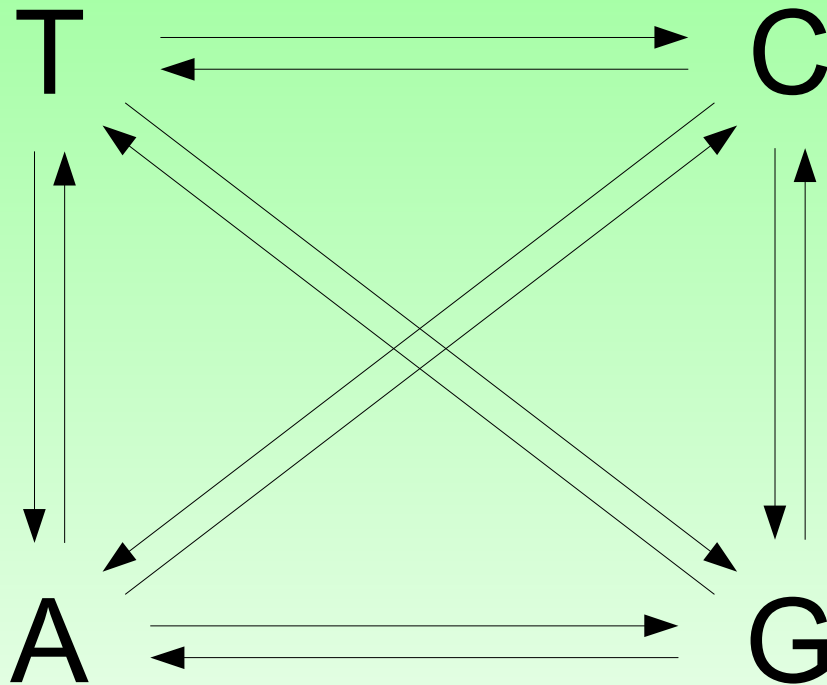
進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

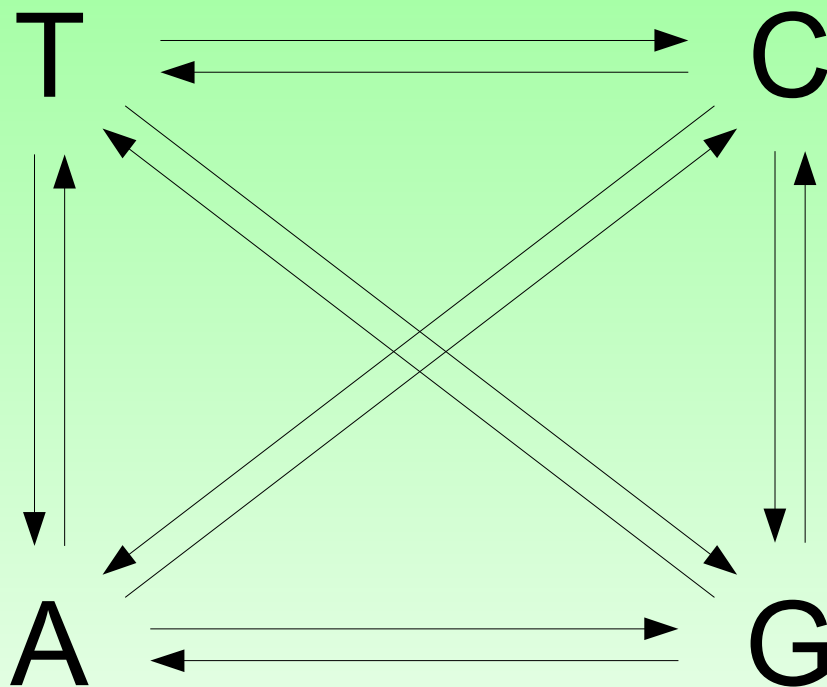
2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

DNAの塩基置換パターン

- 塩基の置換パターンは12通り



DNAの塩基置換パターン



- 塩基の置換パターンは12通り
- それぞれの起きる確率の違い = 形質状態間の推移速度の不均質性をモデル化する

塩基置換確率行列

- r_{XY} は塩基Xから塩基Yへの置換確率

From\To	A	C	G	T
A	-	r_{AC}	r_{AG}	r_{AT}
C	r_{CA}	-	r_{CG}	r_{CT}
G	r_{GA}	r_{GC}	-	r_{GT}
T	r_{TA}	r_{TC}	r_{TG}	-

塩基置換確率行列

- r_{XY} は塩基Xから塩基Yへの置換確率
- $r_{XY}=r_{YX}$ なモデルを時間反転可能という

From\To	A	C	G	T
A	-	r_{AC}	r_{AG}	r_{AT}
C	r_{CA}	-	r_{CG}	r_{CT}
G	r_{GA}	r_{GC}	-	r_{GT}
T	r_{TA}	r_{TC}	r_{TG}	-

サイト間の置換速度不均質性

site	12345M
OTU1	TGTTT	..	TTTTTC
OTU2	AGTAC	..	TTTTTC
OTU3	AGTAT	..	TTGTC
⋮	⋮		⋮
OTUN	AGTAT	..	ATTTC

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

サイト間の置換速度不均質性

site	1	2	3	4	5	M			
OTU1	T	G	T	T	T	T	T	T	T	C
OTU2	A	G	T	A	C	T	T	T	T	C
OTU3	A	G	T	A	T	T	T	G	T	C
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
OTUN	A	G	T	A	T	A	T	T	T	C

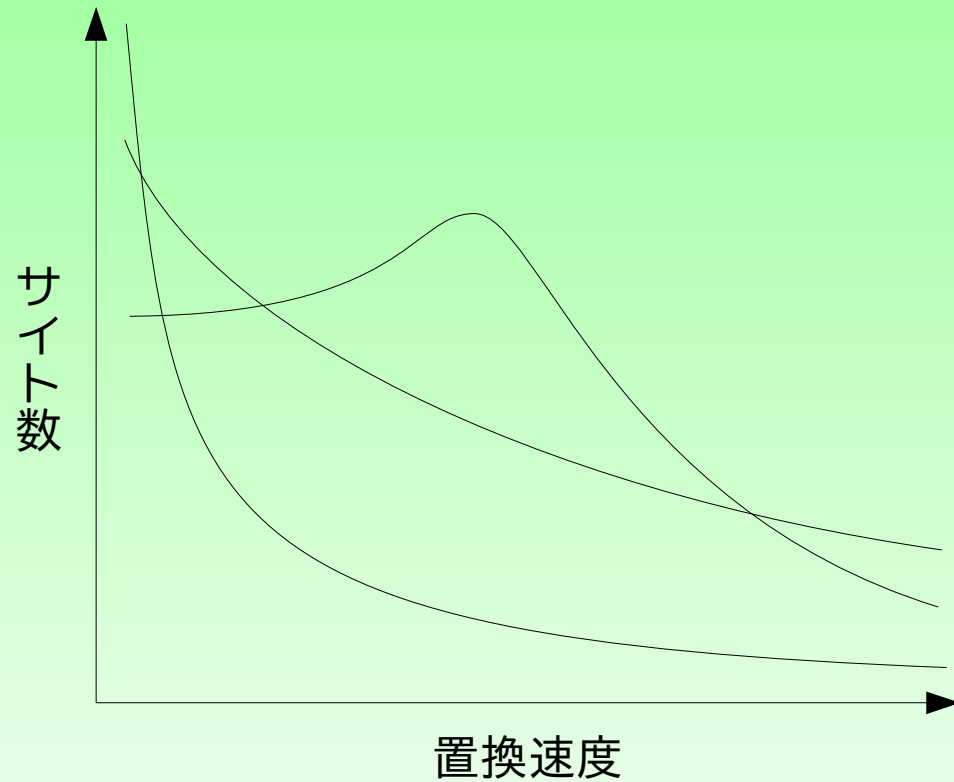
- 変異の多い=速いサイトと少ない=遅いサイトがある

サイト間の置換速度不均質性

site	1	2	3	4	5	M			
OTU1	T	G	T	T	T	T	T	T	T	C
OTU2	A	G	T	A	C	T	T	T	T	C
OTU3	A	G	T	A	T	T	T	G	T	C
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
OTUN	A	G	T	A	T	A	T	T	T	C

- 変異の多い=速いサイトと少ない=遅いサイトがある
- サイト間の置換速度の不均質性をモデル化する

サイト間の置換速度不均質性

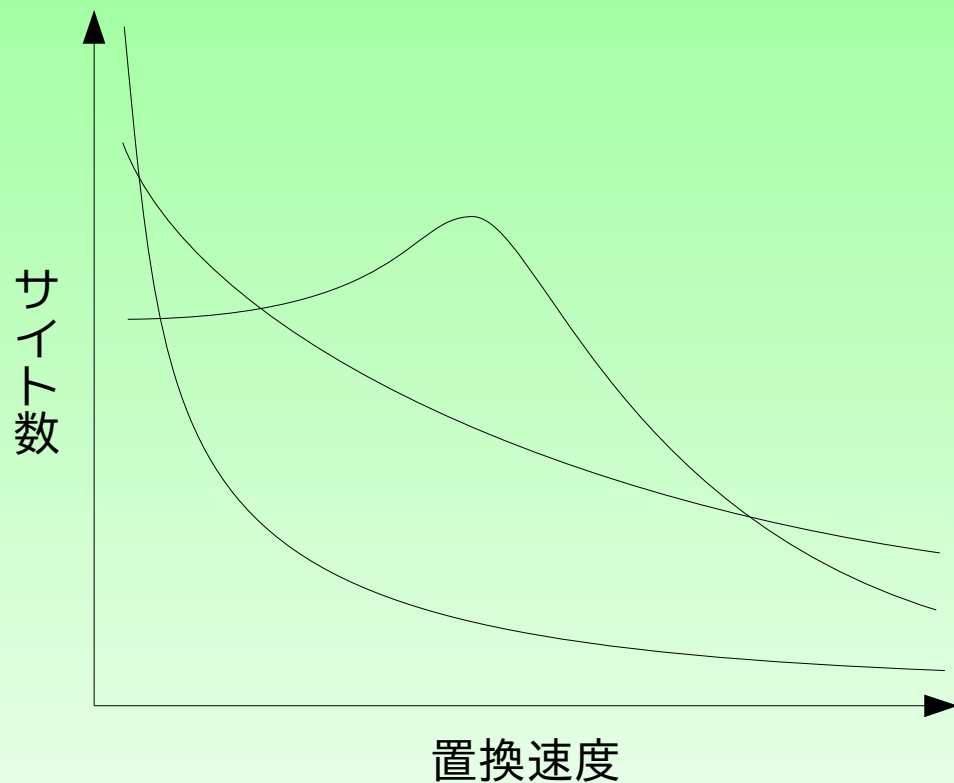


進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

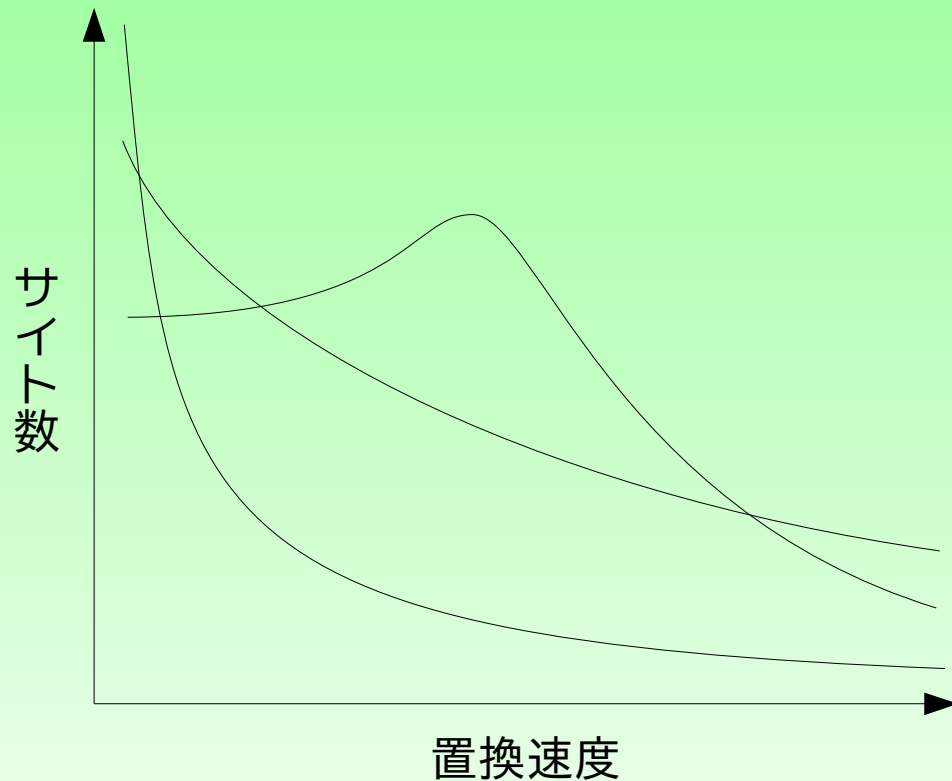
2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

サイト間の置換速度不均質性



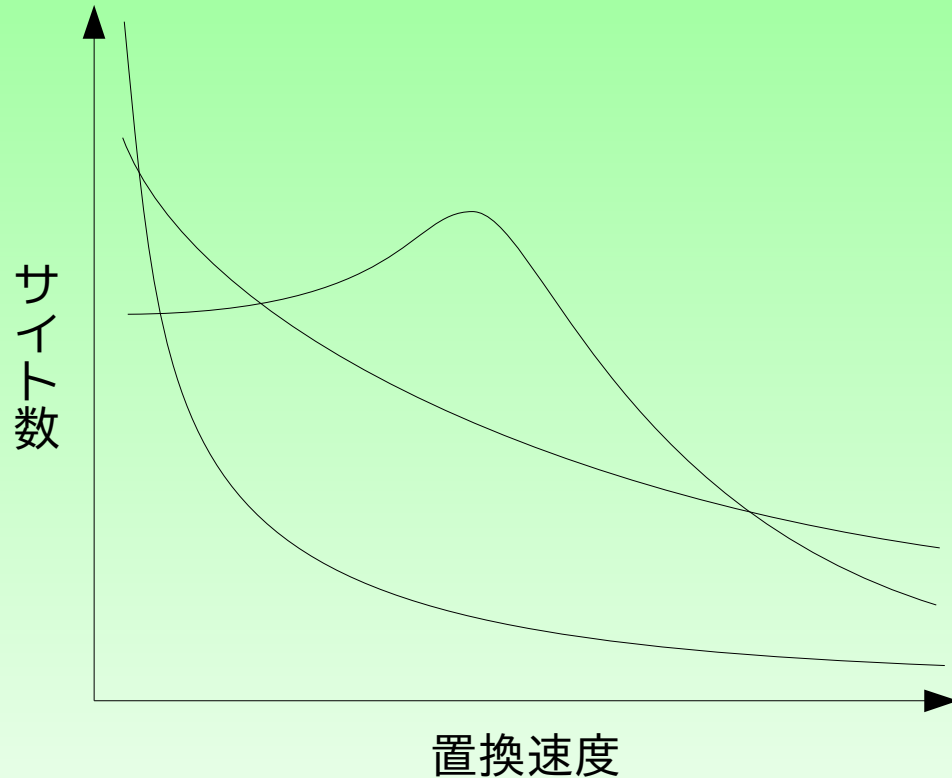
- 置換速度は連続データ
- 分布は非対称

サイト間の置換速度不均質性



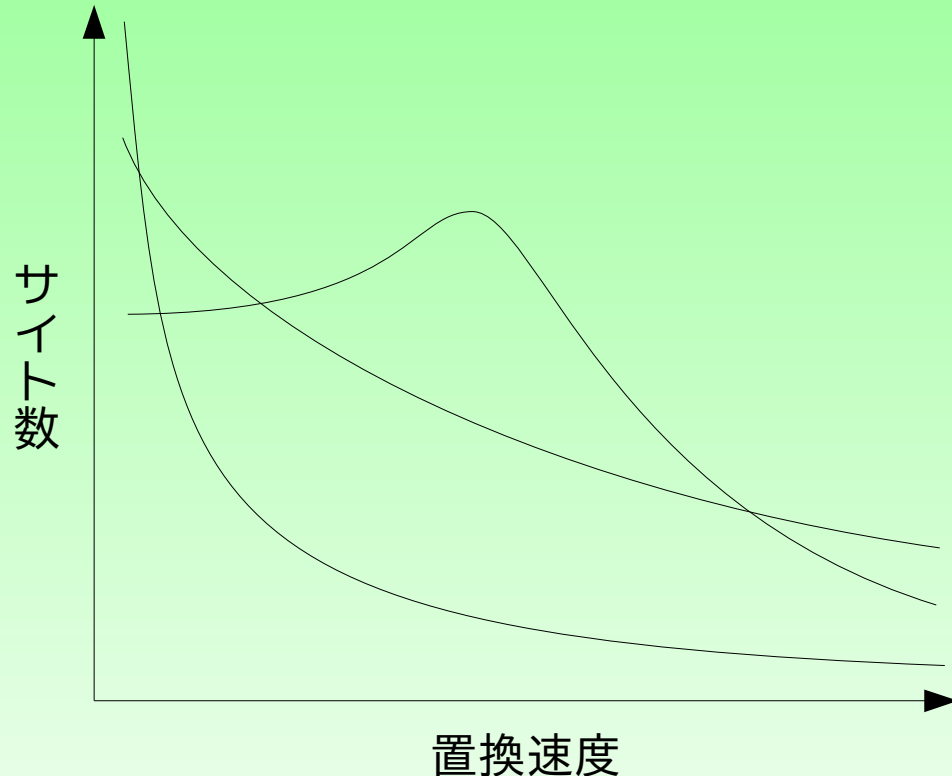
- 置換速度は連続データ
- 分布は非対称
- ガンマ分布を当てはめられる
 - パラメータ数は1

サイト間の置換速度不均質性



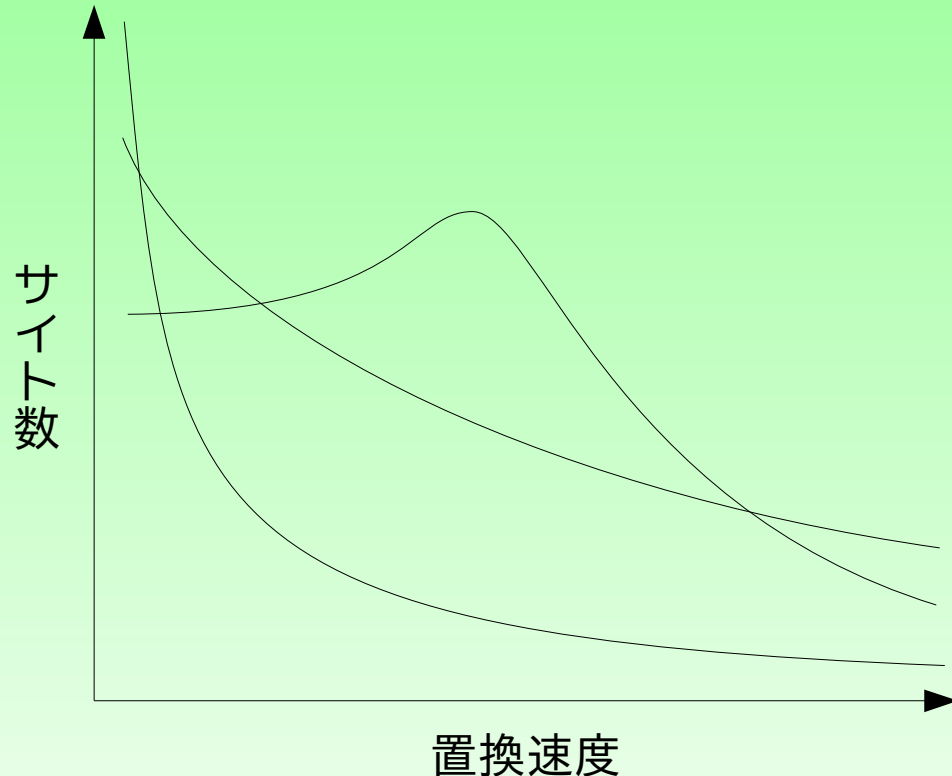
- 置換速度は連続データ
- 分布は非対称
- ガンマ分布を当てはめられる
 - パラメータ数は1
- 変異があるサイトと無いサイトのカテゴリ分けも用いられる
 - パラメータ数は1

サイト間の置換速度不均質性



- 置換速度は連続データ
- 分布は非対称
- ガンマ分布を当てはめられる
 - パラメータ数は1
- 変異があるサイトと無いサイトのカテゴリ分けも用いられる
 - パラメータ数は1
- 変異のあるサイトと無いサイトに分けた上で、変異のあるサイトにガンマ分布を当てはめる
 - パラメータ数は2

サイト間の置換速度不均質性



- 置換速度は連続データ
- 分布は非対称
- ガンマ分布を当てはめられる
 - パラメータ数は1
- 変異があるサイトと無いサイトのカテゴリ分けも用いられる
 - パラメータ数は1
- 変異のあるサイトと無いサイトに分けた上で、変異のあるサイトにガンマ分布を当てはめる
 - パラメータ数は2
- 3つのコドン位置ごとに異なる速度を当てはめる
 - パラメータ数は2

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

多領域データの場合

- 遺伝子A

OTU1	TGTTTTCTTT	...	TTTTTC
OTU2	AGTATTCTTC	...	TTTTTC
OTU3	AGTATTCTTT	...	TTTTTC
OTU4	AATATTTTTT	...	TTTTTC
OTU5	AGTATTTTTT	...	TTTTTC
OTU6	AGTATTCTCT	...	TTTCC
OTU7	AGTATTCTTT	...	TTTTTC
OTU8	AGTATTTTTT	...	TTTTTC
OTU9	AGTATTCTTT	...	TTTTTC

- 遺伝子B

OTU1	TCCTCTACTA	...	AGCTA
OTU2	TCTCCTACTA	...	AGCTA
OTU3	TCTACTACTA	...	AGCTA
OTU4	TCTTCCACTA	...	AGTTA
OTU5	TCTGCCGCTA	...	AGTTA
OTU6	TCTTTCACTA	...	AGTCA
OTU7	TCTTTTACTA	...	AGTCA
OTU8	TTTCCCCTG	...	AGCCA
OTU9	ATTTCCACTG	...	AGCCA

多領域データの場合

- 遺伝子A

OTU1	TGTTTTCTTT	...	TTTTTC
OTU2	AGTATTCTTC	...	TTTTTC
OTU3	AGTATTCTTT	...	TTTTTC
OTU4	AATATTTTTT	...	TTTTTC
OTU5	AGTATTTTTT	...	TTTTTC
OTU6	AGTATTCTCT	...	TTTCC
OTU7	AGTATTCTTT	...	TTTTTC
OTU8	AGTATTTTTT	...	TTTTTC
OTU9	AGTATTCTTT	...	TTTTTC

- 塩基置換確率行列はGTR
- サイト間の速度不均質性はガンマ分布でモデル化

- 遺伝子B

OTU1	TCCTCTACTA	...	AGCTA
OTU2	TCTCCTACTA	...	AGCTA
OTU3	TCTACTACTA	...	AGCTA
OTU4	TCTTCCACTA	...	AGTTA
OTU5	TCTGCCGCTA	...	AGTTA
OTU6	TCTTTCACTA	...	AGTCA
OTU7	TCTTTTACTA	...	AGTCA
OTU8	TTTCCCGCTG	...	AGCCA
OTU9	ATTTCCACTG	...	AGCCA

- 塩基置換確率行列はK2P
- サイト間の速度不均質性は可変/不変サイト比でモデル化

多領域データの場合

- 遺伝子A

OTU1 TGTTTTCTTT ... TTTTC
OTU2 AGTATTCTTC ... TTTTC
OTU3 AGTATTCTTT ... TTTTC
OTU4 AATATTTTTT ... TTTTC
OTU5 AGTATTTTTT ... TTTTC
OTU6 AGTATTCTCT ... TTTCC
OTU7 AGTATTCTTT ... TTTTC
OTU8
OTU9

- 遺伝子B

OTU1 TCCTCTACTA ... AGCTA
OTU2 TCTCCTACTA ... AGCTA
OTU3 TCTACTACTA ... AGCTA
OTU4 TCTTCCACTA ... AGTTA
OTU5 TCTGCCGCTA ... AGTTA
OTU6 TCTTTCACTA ... AGTCA
OTU7 TCTTTTACTA ... AGTCA
OTU8
OTU9

複数領域のそれぞれに異なる分子進化モデルを当てはめる

- 塩基置換確率行列はGTR
- サイト間の速度不均質性はガンマ分布でモデル化

- 塩基置換確率行列はK2P
- サイト間の速度不均質性は可変/不変サイト比でモデル化

多領域データの場合

- 遺伝子A

OTU1 TGTTTTCTTT ... TTTTC
OTU2 AGTATTCTTC ... TTTTC
OTU3 AGTATTCTTT ... TTTTC
OTU4 AATATTTTTT ... TTTTC
OTU5 AGTATTTTTT ... TTTTC
OTU6 AGTATTCTCT ... TTTCC
OTU7 AGTATTCTTT ... TTTTC
OTU8
OTU9

- 遺伝子B

OTU1 TCCTCTACTA ... AGCTA
OTU2 TCTCCTACTA ... AGCTA
OTU3 TCTACTACTA ... AGCTA
OTU4 TCTTCCACTA ... AGTTA
OTU5 TCTGCCGCTA ... AGTTA
OTU6 TCTTTCACTA ... AGTCA
OTU7 TCTTTTACTA ... AGTCA
OTU8
OTU9

複数領域のそれぞれに異なる分子進化モデルを当てはめる

- 塩基置換確率行列はGTR

- サイト間の速度不均質性はガンマ分布でモデル化

- 塩基置換確率行列はK2P

- サイト間の速度不均質性は可変/不変サイト比でモデル化

比例モデルと分離モデルという二通りの当てはめ方がある

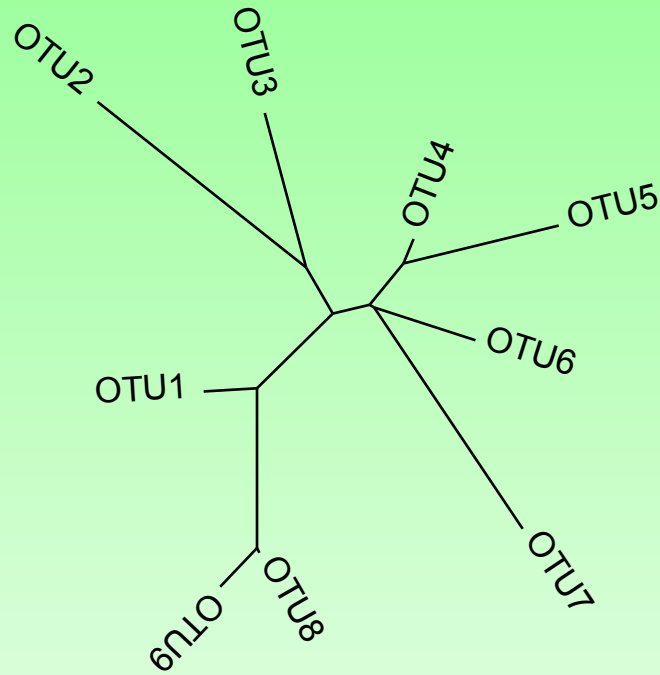
比例モデル

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

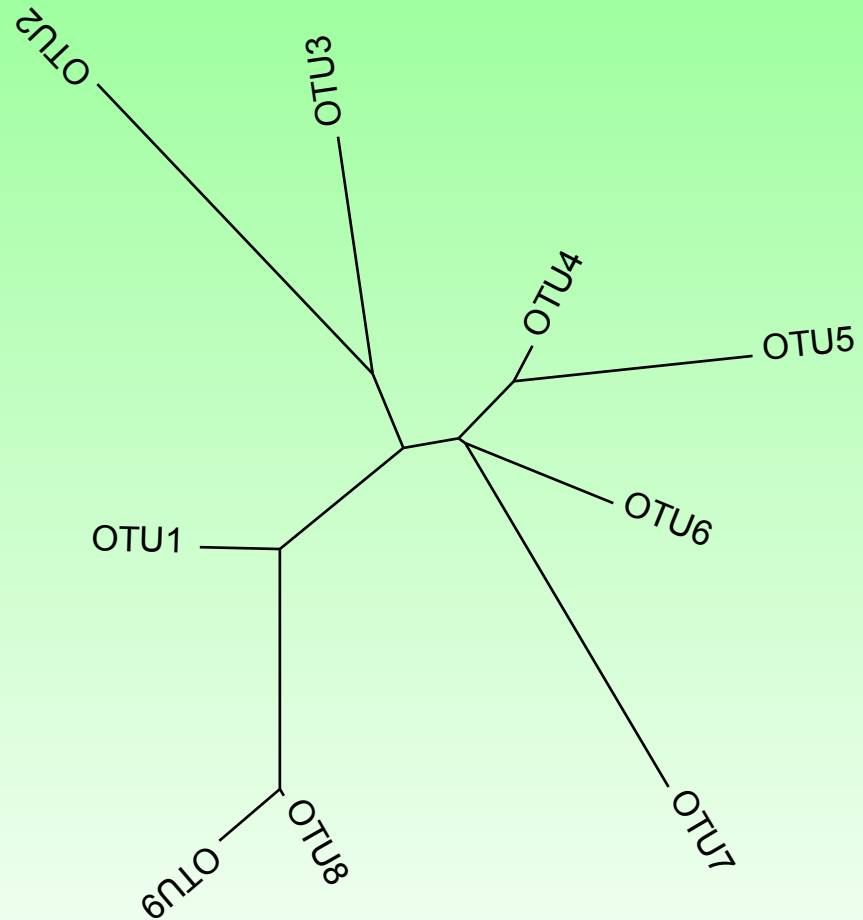
田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

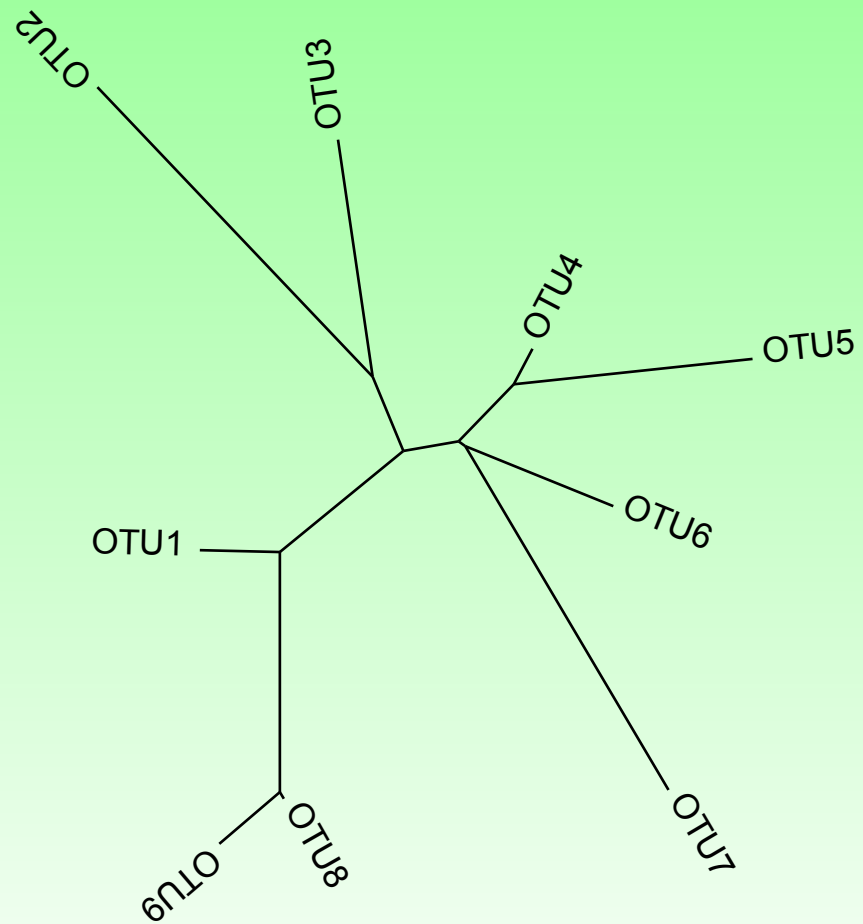
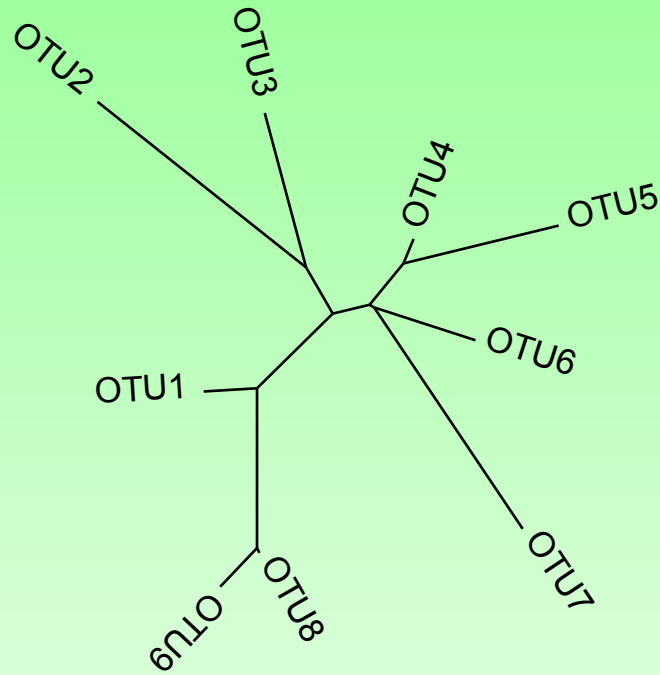
比例モデル



- 比例モデルでは系統樹は相似形



比例モデル



- 比例モデルでは系統樹は相似形
- パラメータ数は
 - 枝長 : $OTU数 \times 2 - 3$
 - 速度比(枝長比) : 領域数 - 1

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

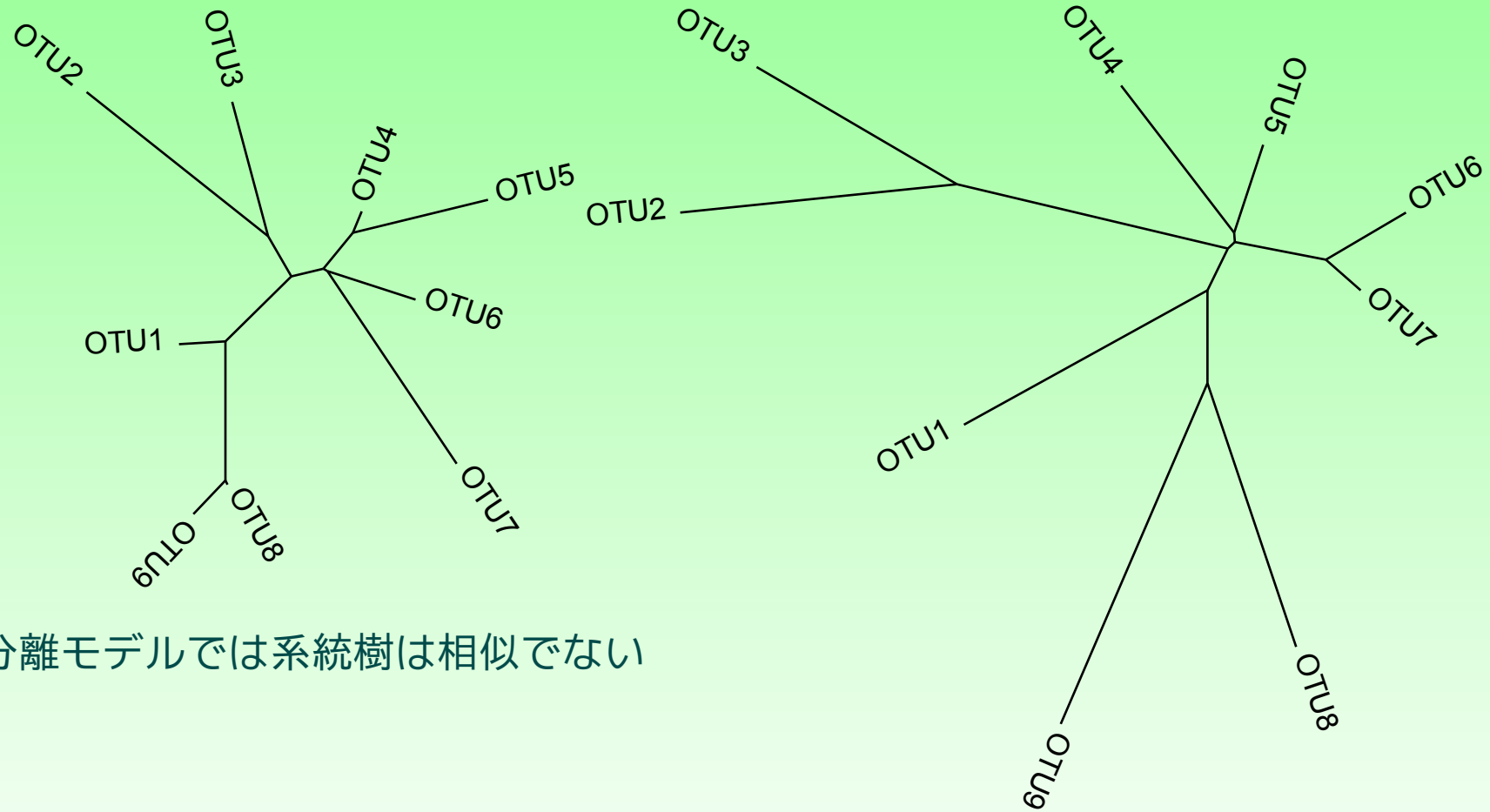
分離モデル

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

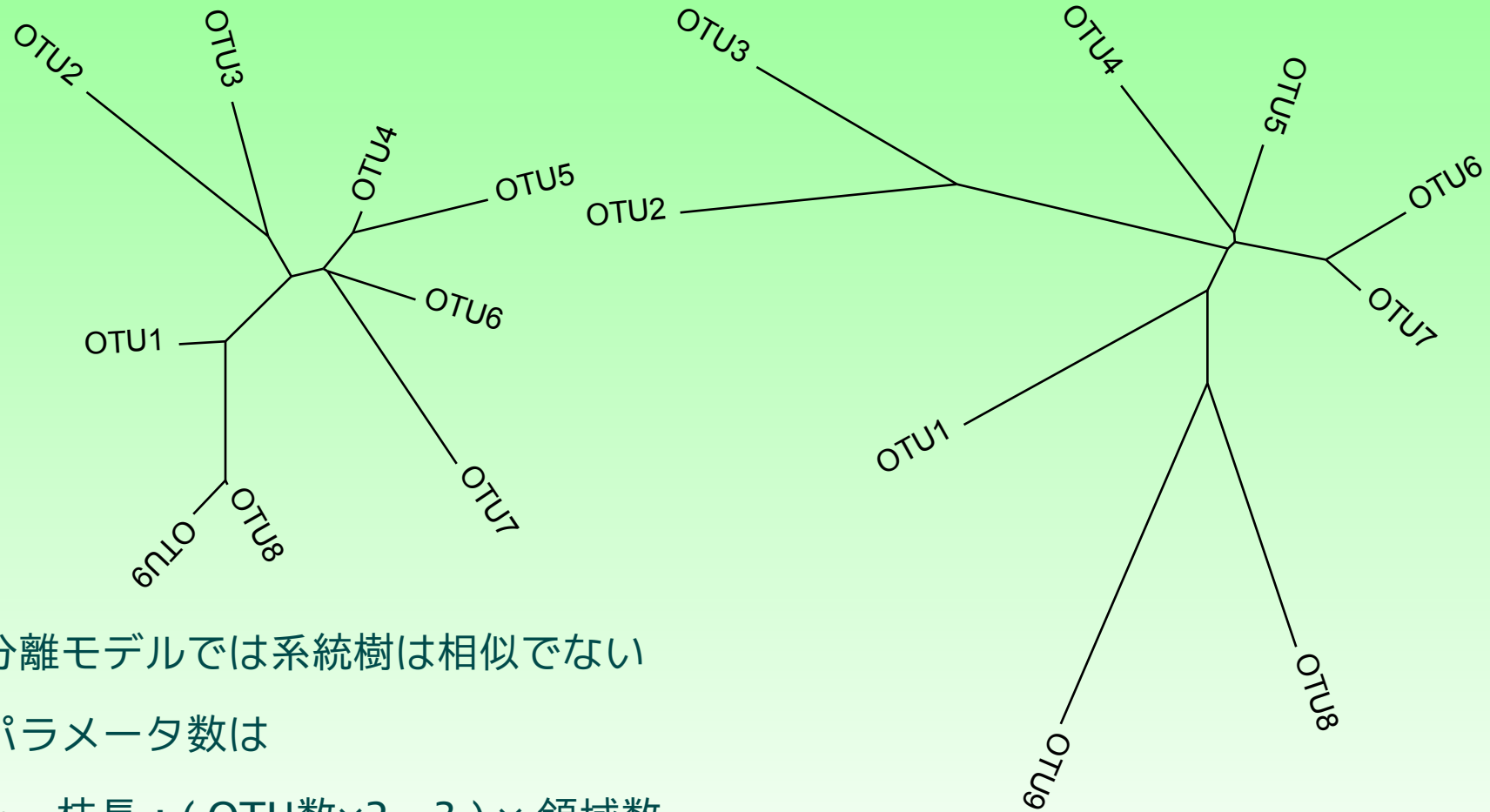
2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

分離モデル



- 分離モデルでは系統樹は相似でない

分離モデル



- 分離モデルでは系統樹は相似でない
- パラメータ数は
 - 枝長 : $(\text{OTU数} \times 2 - 3) \times \text{領域数}$

多領域データにおけるモデルとパラメータ数

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

多領域データにおけるモデルとパラメータ数

- 分子進化モデル
 - 置換確率行列
 - DNAでは領域ごとに0~8
 - サイト間置換速度不均質性
 - 領域ごとに0以上

多領域データにおけるモデルとパラメータ数

- 分子進化モデル
 - 置換確率行列
 - DNAでは領域ごとに0~8
 - サイト間置換速度不均質性
 - 領域ごとに0以上
- 系統樹
 - 単一モデルの場合
 - 枝長：OTU数 $\times 2 - 3$
 - 比例モデルの場合
 - 枝長：OTU数 $\times 2 - 3$
 - 枝長比：領域数 $- 1$
 - 分離モデルの場合
 - 枝長：(OTU数 $\times 2 - 3$)
× 領域数

多領域データにおけるモデルとパラメータ数

• 分子進化モデル

- 置換確率行列
 - DNAでは領域ごとに0~8
- サイト間置換速度不均質性
 - 領域ごとに0以上

各領域ごとに最適な分子進化モデルを選択し、そのモデルを適用して最適な樹形を選択する

• 系統樹

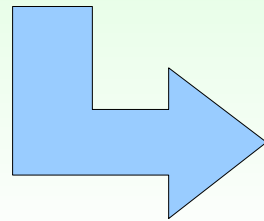
- 単一モデルの場合
 - 枝長：OTU数 $\times 2 - 3$
- 比例モデルの場合
 - 枝長：OTU数 $\times 2 - 3$
 - 枝長比：領域数 $- 1$
- 分離モデルの場合
 - 枝長：(OTU数 $\times 2 - 3$)
× 領域数

多領域データにおけるモデルとパラメータ数

- 分子進化モデル

- 置換確率行列
 - DNAでは領域ごとに0~8
- サイト間置換速度不均質性
 - 領域ごとに0以上

各領域ごとに最適な分子進化モデルを選択し、そのモデルを適用して最適な樹形を選択する



Kakusan3

- 系統樹

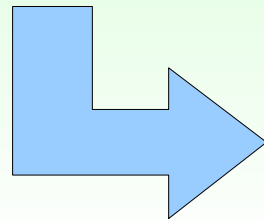
- 単一モデルの場合
 - 枝長 : $OTU数 \times 2 - 3$
- 比例モデルの場合
 - 枝長 : $OTU数 \times 2 - 3$
 - 枝長比 : 領域数 - 1
- 分離モデルの場合
 - 枝長 : $(OTU数 \times 2 - 3) \times 領域数$

多領域データにおけるモデルとパラメータ数

- 分子進化モデル

- 置換確率行列
 - DNAでは領域ごとに0~8
- サイト間置換速度不均質性
 - 領域ごとに0以上

各領域ごとに最適な分子進化モデルを選択し、そのモデルを適用して最適な樹形を選択する



Kakusan3

- 系統樹

- 単一モデルの場合
 - 枝長：OTU数 $\times 2 - 3$
- 比例モデルの場合
 - 枝長：OTU数 $\times 2 - 3$
 - 枝長比：領域数 $- 1$
- 分離モデルの場合
 - 枝長：(OTU数 $\times 2 - 3$)
 \times 領域数

単一・比例・分離モデル間の比較はしないことに注意

Kakusan3の概要

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

Kakusan3の概要

- PC・Macからクラスタ・スパコンでも動作可能
 - Windows・MacOS X・その他UNIX対応

Kakusan3の概要

- PC・Macからクラスタ・スパコンでも動作可能
 - Windows・MacOS X・その他UNIX対応
- 多くの形式のデータを読み込み可能
 - FASTA・GenBank・NEXUS・PHYML形式に対応

Kakusan3の概要

- PC・Macからクラスタ・スパコンでも動作可能
 - Windows・MacOS X・その他UNIX対応
- 多くの形式のデータを読み込み可能
 - FASTA・GenBank・NEXUS・PHYLIP形式に対応
- 解析ソフト用設定ファイルの出力
 - PAUP*4・MrBayes3・Treefinderに対応

Kakusan3の概要

- PC・Macからクラスタ・スパコンでも動作可能
 - Windows・MacOS X・その他UNIX対応
- 多くの形式のデータを読み込み可能
 - FASTA・GenBank・NEXUS・PHYML形式に対応
- 解析ソフト用設定ファイルの出力
 - PAUP*4・MrBayes3・Treefinderに対応
- 遺伝子座・コドン位置ごとにモデル選択
 - 各モデルの尤度最大化を並列処理し高速化

Kakusan3が中でやっていること

塩基配列データの入力

データ形式の変換
(ReadSeq / PAUP*)

塩基配列データの分割・結合

OTU間での塩基組成均一性を
 χ^2 独立性の検定で確認

近隣結合法にによる系統樹の生成
(PHYLIP / PAUP*)

各領域ごとに候補モデルを
当てはめて尤度最大化
(BASEML / Treefinder / PAUP*)

サンプルサイズの算出

AIC · AICc · BICの算出

結果の出力

Kakusan3

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

入力データファイルの用意

- A.phy

```
9 15
Outgroup TGTTTTCTTTTTTTC
Tax1     AGTATTCTTCTTTTC
Tax2     AGTATTCTTTTTTTC
Tax3     AATATTTTTTTTTTTC
Tax4     AGTATTTTTTTTTTTC
Tax5     AGTATTCTCTTTTCC
Tax6     AGTATTCTTTTTTTC
Tax7     AGTATTTTTTTTTTTC
Tax8     AGTATTCTTTTTTTC
```

- B_P.phy

```
9 15
Outgroup TCCTCTACTAAGCTA
Tax1     TCTCCTACTAAGCTA
Tax2     TCTACTACTAAGCTA
Tax3     TCTTCCACTAAGTTA
Tax4     TCTGCCGCTAAGTTA
Tax5     TCTTTCACTAAGTCA
Tax6     TCTTTTACTAAGTCA
Tax7     TTTCCCGCTGAGCCA
Tax8     ATTTCCACTGAGCCA
```

入力データファイルの用意

- A.phy

```
9 15
Outgroup TGTTTTCTTTTTTTC
Tax1      AGTATTCTTCTTTTC
Tax2      AGTATTCTTTTTTTC
Tax3      AATATTTTTTTTTTTC
Tax4      AGTATTTTTTTTTTTC
Tax5      AGTATTCTCTTTTCC
Tax6      AGTATTCTTTTTTTC
Tax7      AGTATTTTTTTTTTTC
Tax8      AGTATTCTTTTTTTC
```

- B_P.phy

```
9 15
Outgroup TCCTCTACTAAGCTA
Tax1      TCTCCTACTAAGCTA
Tax2      TCTACTACTAAGCTA
Tax3      TCTTCCACTAAGTTA
Tax4      TCTGCCGCTAAGTTA
Tax5      TCTTTCACTAAGTCA
Tax6      TCTTTTACTAAGTCA
Tax7      TTTCCCGCTGAGCCA
Tax8      ATTTCCACTGAGCCA
```

- 領域 (遺伝子座) ごとに異なるファイルに分ける

入力データファイルの用意

- A.phy

```
9 15
Outgroup TGTTTTCTTTTTTTC
Tax1     AGTATTCTTCTTTTC
Tax2     AGTATTCTTTTTTTC
Tax3     AATATTTTTTTTTTTC
Tax4     AGTATTTTTTTTTTTC
Tax5     AGTATTCTCTTTTCC
Tax6     AGTATTCTTTTTTTC
Tax7     AGTATTTTTTTTTTTC
Tax8     AGTATTCTTTTTTTC
```

- B_P.phy

```
9 15
Outgroup TCCTCTACTAAGCTA
Tax1     TCTCCTACTAAGCTA
Tax2     TCTACTACTAAGCTA
Tax3     TCTTCCACTAAGTTA
Tax4     TCTGCCGCTAAGTTA
Tax5     TCTTTCACTAAGTCA
Tax6     TCTTTTACTAAGTCA
Tax7     TTTCCCGCTGAGCCA
Tax8     ATTTCCACTGAGCCA
```

- 領域 (遺伝子座) ごとに異なるファイルに分ける
- OTU名はファイル間で統一する

入力データファイルの用意

- A.phy

```
9 15
Outgroup TGTTTTCTTTTTTTC
Tax1     AGTATTCTTCTTTTC
Tax2     AGTATTCTTTTTTTC
Tax3     AATATTTTTTTTTTTC
Tax4     AGTATTTTTTTTTTTC
Tax5     AGTATTCTCTTTTCC
Tax6     AGTATTCTTTTTTTC
Tax7     AGTATTTTTTTTTTTC
Tax8     AGTATTCTTTTTTTC
```

- B_P.phy

```
9 15
Outgroup TCCTCTACTAAGCTA
Tax1     TCTCCTACTAAGCTA
Tax2     TCTACTACTAAGCTA
Tax3     TCTTCCACTAAGTTA
Tax4     TCTGCCGCTAAGTTA
Tax5     TCTTTCACTAAGTCA
Tax6     TCTTTTACTAAGTCA
Tax7     TTTCCCGCTGAGCCA
Tax8     ATTTCCACTGAGCCA
```

- 領域 (遺伝子座) ごとに異なるファイルに分ける
- OTU名はファイル間で統一する
- タンパクコード領域のファイルはファイル名の末尾を「_P」にする

入力データファイルの用意

- A.phy

```
9 15
Outgroup TGTTTTCTTTTTTTC
Tax1     AGTATTCTTCTTTTC
Tax2     AGTATTCTTTTTTTC
Tax3     AATATTTTTTTTTTTC
Tax4     AGTATTTTTTTTTTTC
Tax5     AGTATTCTCTTTTCC
Tax6     AGTATTCTTTTTTTC
Tax7     AGTATTTTTTTTTTTC
Tax8     AGTATTCTTTTTTTC
```

- B_P.phy

```
9 15
Outgroup TCCTCTACTAAGCTA
Tax1     TCTCCTACTAAGCTA
Tax2     TCTACTACTAAGCTA
Tax3     TCTTCCACTAAGTTA
Tax4     TCTGCCGCTAAGTTA
Tax5     TCTTTCACTAAGTCA
Tax6     TCTTTTACTAAGTCA
Tax7     TTTCCCGCTGAGCCA
Tax8     ATTTCCACTGAGCCA
```

- 領域 (遺伝子座) ごとに異なるファイルに分ける
- OTU名はファイル間で統一する
- タンパクコード領域のファイルはファイル名の末尾を「_P」にする
- intronとexonは混ぜるな危険

Kakusan3の起動とデータの入力 (Windowsのみ)

名前 更新日時 種類 サイズ

12S.nex 16S.nex APP_i.nex

0. 遺伝子座ごとにばらばらのファイルを用意する
(タンパクコード領域ファイル名の末尾は「_P」にする)

フルパスをクリップボードにコピー
フルパスをダブルコーテーションで囲ってクリップボードにコピー
フルパスを¥を重ねてクリップボードにコピー

ESET NOD32 Antivirus で検査
詳細設定オプション

秀丸エディタで開く

書庫作成(A)

送る(N)

切り取り(T)
コピー(C)

ショートカットの作成(S)
削除(D)
名前の変更(M)
プロパティ(R)

Aminosan
B2
デスクトップ (ショートカットを作成)
Kakusan3
リムーバブルディスク (I:)
DVD RW ドライブ (Q:)
DVD RW ドライブ (R:)

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

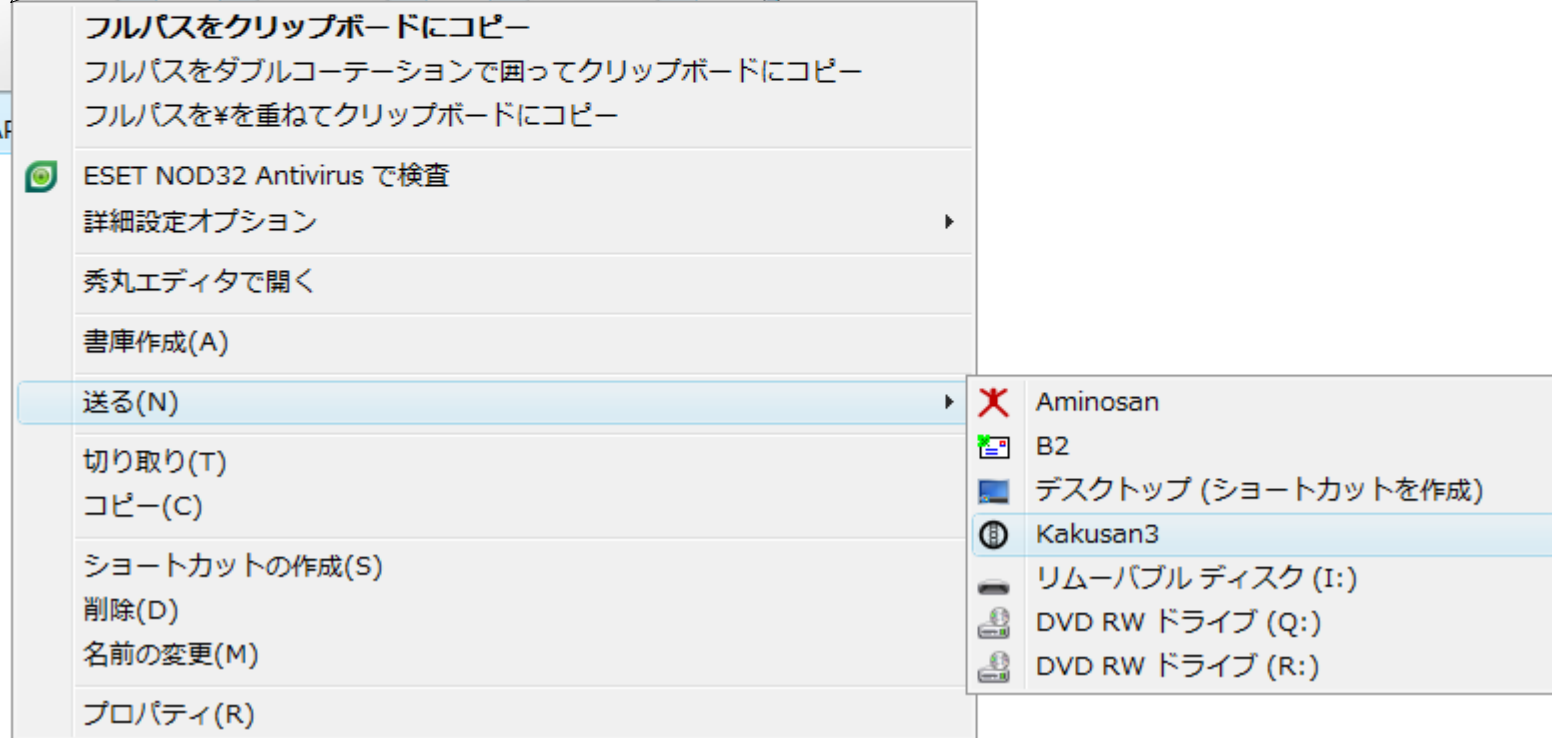
田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

Kakusan3の起動とデータの入力 (Windowsのみ)

1. データファイルを選択して右クリック

0. 遺伝子座ごとにばらばらのファイルを用意する
(タンパクコード領域ファイル名の末尾は「_P」にする)



進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

Kakusan3の起動とデータの入力 (Windowsのみ)

1. データファイルを選択して右クリック

0. 遺伝子座ごとにばらばらのファイルを用意する
(タンパクコード領域ファイル名の末尾は「_P」にする)

フルパスをクリップボードにコピー
フルパスをダブルコーテーションで囲ってクリップボードにコピー
フルパスを¥を重ねてクリップボードにコピー

2. 「送る」を選択

ESET NOD32 Antivirus
詳細設定オプション
秀丸エディタで開く
書庫作成(A)

送る(N)

切り取り(T)

コピー(C)

ショートカットの作成(S)

削除(D)

名前の変更(M)

プロパティ(R)

Aminosan
B2
デスクトップ (ショートカットを作成)
Kakusan3
リムーバブルディスク (I:)
DVD RW ドライブ (Q:)
DVD RW ドライブ (R:)

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

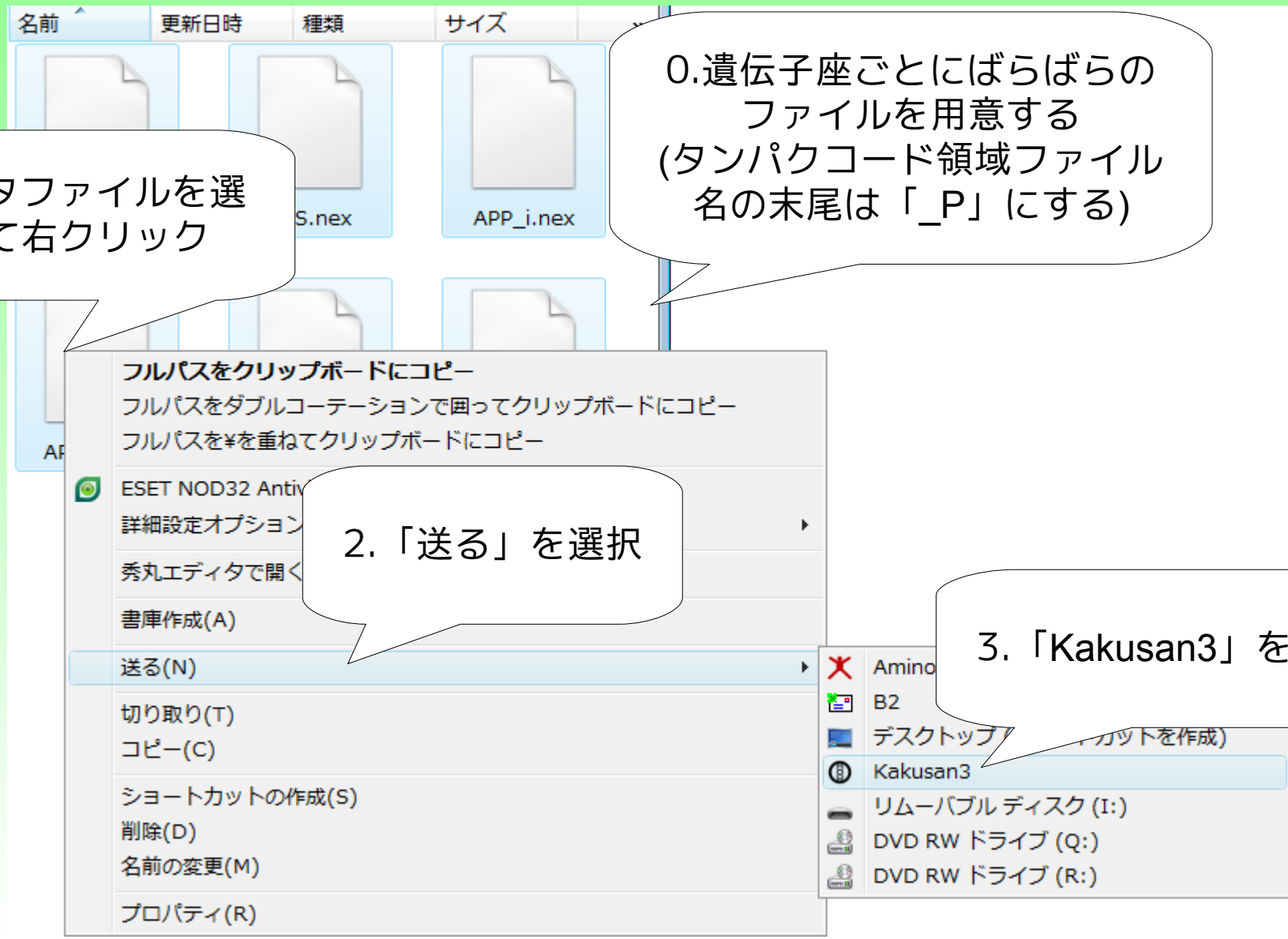
Kakusan3の起動とデータの入力 (Windowsのみ)

1. データファイルを選択して右クリック

0. 遺伝子座ごとにばらばらのファイルを用意する
(タンパクコード領域ファイル名の末尾は「_P」にする)

2. 「送る」を選択

3. 「Kakusan3」を選択

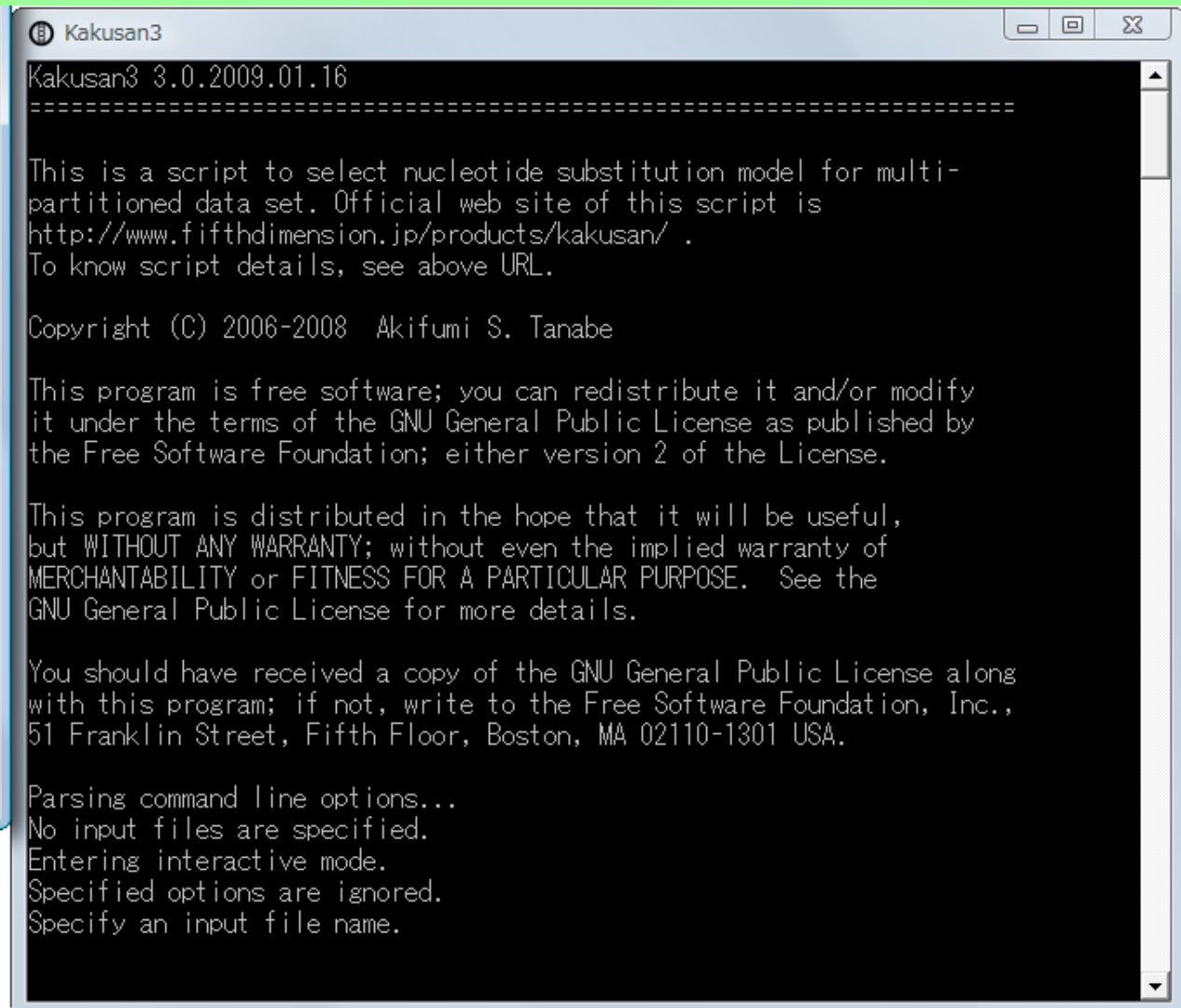
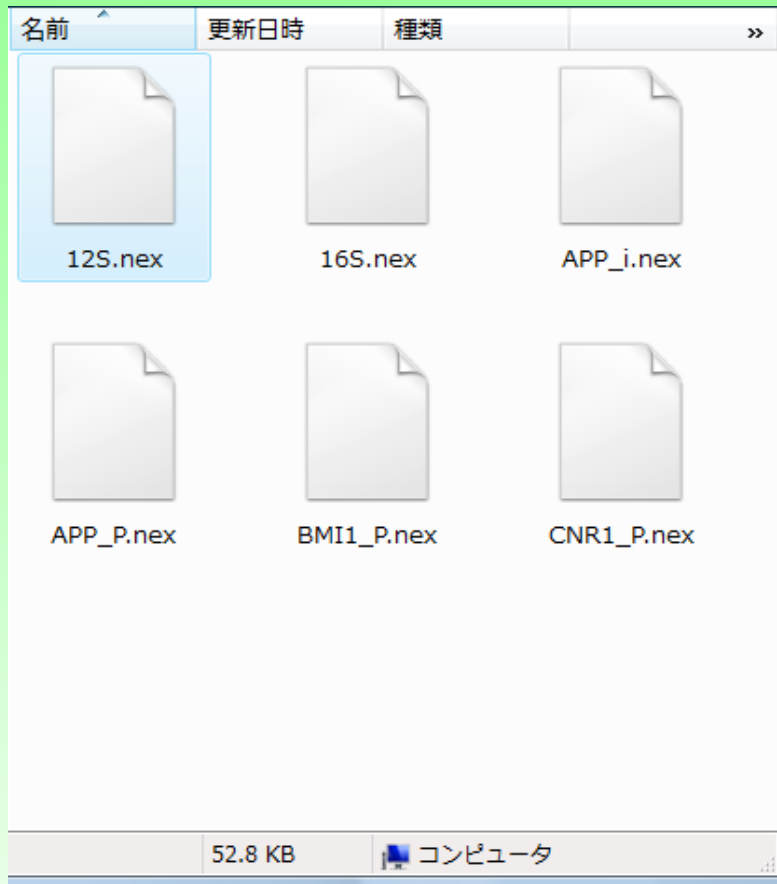


進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

Kakusan3の起動とデータの入力 (Windows・MacOS X共通) (XP以前)

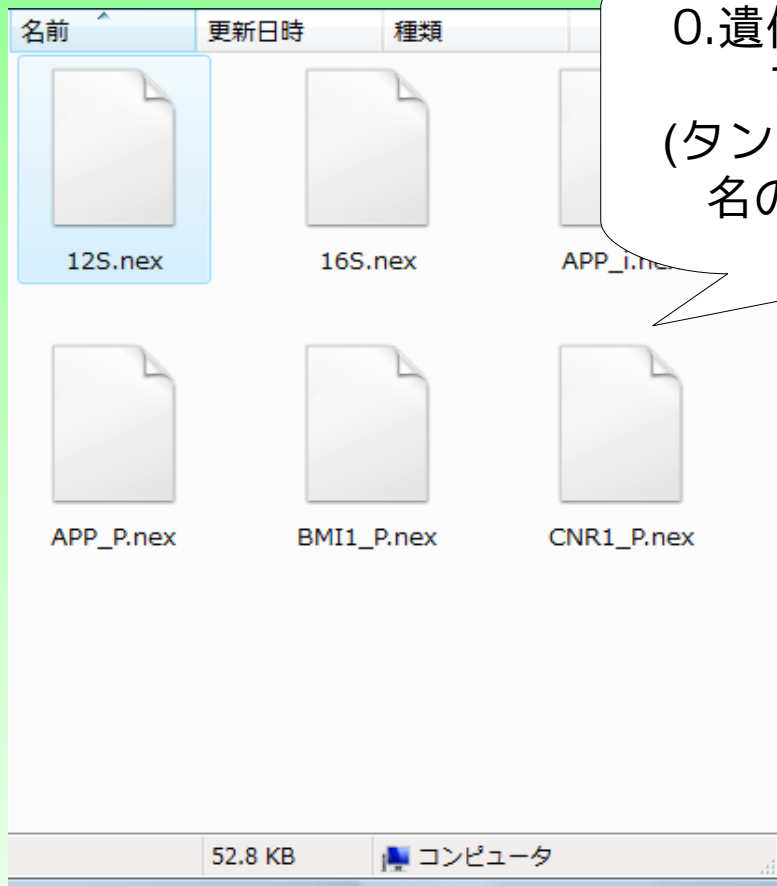


進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

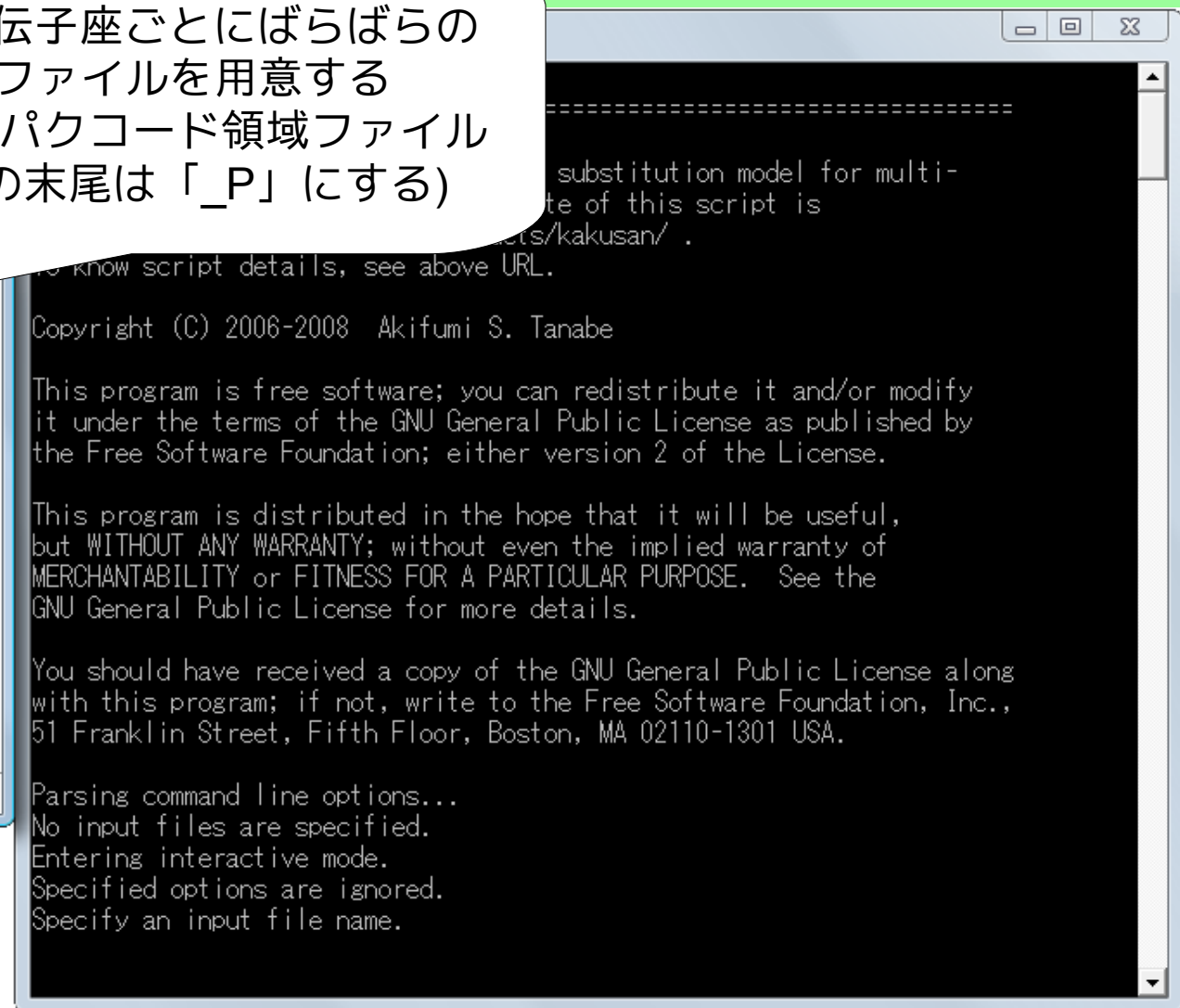
田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

Kakusan3の起動とデータの入力 (Windows・MacOS X共通) (XP以前)



0. 遺伝子座ごとにばらばらの
ファイルを用意する
(タンパクコード領域ファイル
名の末尾は「_P」にする)



進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

Kakusan3の起動とデータの入力 (Windows・MacOS X共通) (XP以前)

0. 遺伝子座ごとにばらばらの
ファイルを用意する
(タンパクコード領域ファイル
名の末尾は「_P」にする)

1. ファイルを1つずつド
ラッグアンドドロップ

```
=====
substitution model for multi-
te of this script is
.../kakusan/ .
to know script details, see above URL.

Copyright (C) 2006-2008 Akifumi S. Tanabe

This program is free software; you can redistribute it and/or modify
it under the terms of the GNU General Public License as published by
the Free Software Foundation; either version 2 of the License.

This program is distributed in the hope that it will be useful,
but WITHOUT ANY WARRANTY; without even the implied warranty of
MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
GNU General Public License for more details.

You should have received a copy of the GNU General Public License along
with this program; if not, write to the Free Software Foundation, Inc.,
51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA.

Parsing command line options...
No input files are specified.
Entering interactive mode.
Specified options are ignored.
Specify an input file name.
```

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

設定ファイル出力対象の指定

```
C:\Program Files\Kakusan3\kakusan3.exe
You should have received a copy of the GNU General Public License along
with this program; if not, write to the Free Software Foundation, Inc.,
51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA.

Parsing command line options...
Log, result and configuration files are output to "D:\Users\shimotsuki\Documents
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS
Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.
Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.
Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.
```


設定ファイル出力対象の指定

```
C:\Program Files\Kakusan3\kakusan3.exe
You should have received a copy of the GNU General Public License along
with this program; if not, write to the Free Software Foundation, Inc.,
51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA.

Parsing command line options...
Log, result and configuration files are output to "D:\Users\shimotsuki\Documents
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS

Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.
```

- PAUP*4・MrBayes3・Treefinderのそれぞれに選択されたモデルを当てはめて系統推定を行うための設定ファイルを出力できる

設定ファイル出力対象の指定

```
C:\Program Files\Kakusan3\kakusan3.exe
You should have received a copy of the GNU General Public License along
with this program; if not, write to the Free Software Foundation, Inc.,
51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA.

Parsing command line options...
Log, result and configuration files are output to "D:\Users\shimotsuki\Documents
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS
Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.
```

- PAUP*4・MrBayes3・Treefinderのそれぞれに選択されたモデルを当てはめて系統推定を行うための設定ファイルを出力できる
- 対象となるソフトの選択によって後の解析オプションが自動的に設定されます

設定ファイル出力対象の指定

```
C:\Program Files\Kakusan3\kakusan3.exe
You should have received a copy of the GNU General Public License along
with this program; if not, write to the Free Software Foundation, Inc.,
51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA.

Parsing command line options...
Log, result and configuration files are output to "D:\Users\shimotsuki\Documents
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS
Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.
```

- PAUP*4・MrBayes3・Treefinderのそれぞれに選択されたモデルを当てはめて系統推定を行うための設定ファイルを出力できる
- 対象となるソフトの選択によって後の解析オプションが自動的に設定されます
- そのため、1度の解析ではどれか1つのソフト用出力のみを有効にすることが望ましい

尤度最大化に用いるプログラムの設定

```
C:\Program Files\Kakusan3\kakusan3.exe
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS

Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.

ANALYSIS OPTIONS

Which do you want to use the program for likelihood calculation? (baseml/tf/paup
)
(default: baseml)
```

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

尤度最大化に用いるプログラムの設定

```
C:\Program Files\Kakusan3\kakusan3.exe
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS

Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.

ANALYSIS OPTIONS

Which do you want to use the program for likelihood calculation? (baseml/tf/paup
)
(default: baseml)
```

- PAUP*4用設定ファイルを書き出すなら paupを、

尤度最大化に用いるプログラムの設定

```
C:\Program Files\Kakusan3\kakusan3.exe
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS

Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.

ANALYSIS OPTIONS

Which do you want to use the program for likelihood calculation? (baseml/tf/paup
)
(default: baseml)
```

- PAUP*4用設定ファイルを書き出すなら paupを、
- MrBayes3用設定ファイルを書き出すなら basemlを、

尤度最大化に用いるプログラムの設定

```
C:\Program Files\Kakusan3\kakusan3.exe
¥2009¥0317自由集会¥test¥12S.nex.kakusan".

OUTPUT OPTIONS

Do you want the model configuration files for PAUP*? (y/n)
(default: y)
n
OK. The model configuration files for PAUP* will NOT be output.

Do you want the model configuration files for MrBayes? (y/n)
(default: y)
n
OK. The model configuration files for MrBayes will NOT be output.

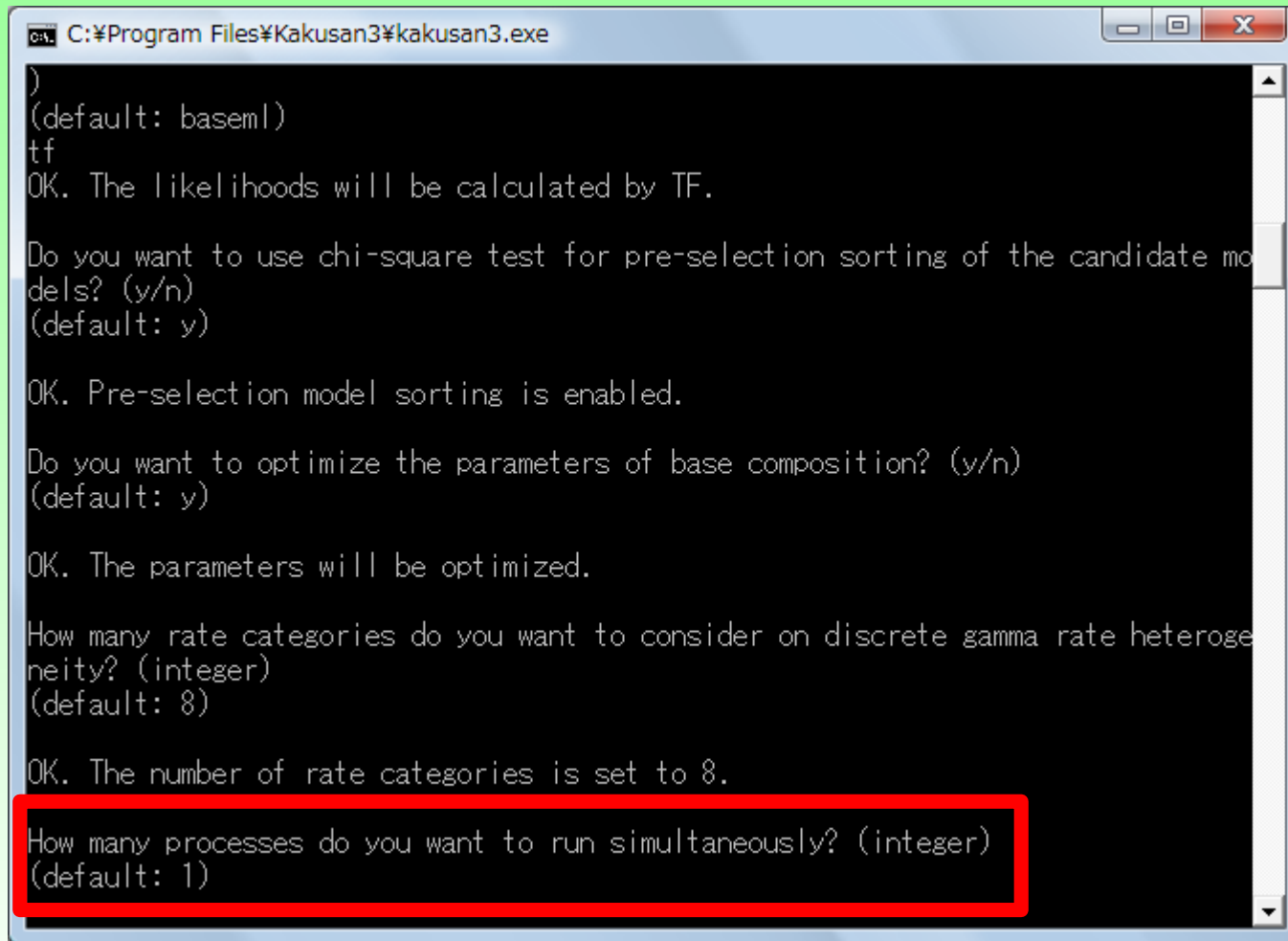
Do you want the model configuration files for Treefinder? (y/n)
(default: y)
y
OK. The model configuration files for Treefinder will be output.

ANALYSIS OPTIONS

Which do you want to use the program for likelihood calculation? (baseml/tf/paup
)
(default: baseml)
```

- PAUP*4用設定ファイルを書き出すなら paupを、
- MrBayes3用設定ファイルを書き出すなら basemlを、
- Treefinder用設定ファイルを書き出すなら paupかtfを指定することを推奨します

並列処理に関する設定



```
C:\Program Files\Kakusan3\kakusan3.exe
)
(default: baseml)
tf
OK. The likelihoods will be calculated by TF.

Do you want to use chi-square test for pre-selection sorting of the candidate models? (y/n)
(default: y)

OK. Pre-selection model sorting is enabled.

Do you want to optimize the parameters of base composition? (y/n)
(default: y)

OK. The parameters will be optimized.

How many rate categories do you want to consider on discrete gamma rate heterogeneity? (integer)
(default: 8)

OK. The number of rate categories is set to 8.

How many processes do you want to run simultaneously? (integer)
(default: 1)
```


並列処理に関する設定

```
C:\Program Files\Kakusan3\kakusan3.exe
)
(default: baseml)
tf
OK. The likelihoods will be calculated by TF.

Do you want to use chi-square test for pre-selection sorting of the candidate models? (y/n)
(default: y)

OK. Pre-selection model sorting is enabled.

Do you want to optimize the parameters of base composition? (y/n)
(default: y)

OK. The parameters will be optimized.

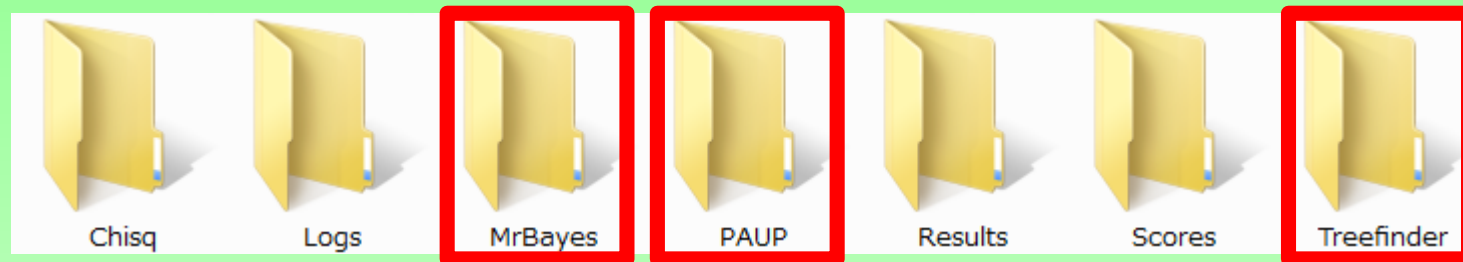
How many rate categories do you want to consider on discrete gamma rate heterogeneity? (integer)
(default: 8)

OK. The number of rate categories is set to 8.

How many processes do you want to run simultaneously? (integer)
(default: 1)
```

- 起動する並列プロセスの最大値を設定
- 搭載しているCPUの数以下にして下さい

出力フォルダの内容



- Chisq - 塩基組成のOTU間均一性の検定結果の出力フォルダ
- Logs - 各種実行ログファイルの保存フォルダ
- MrBayes - MrBayes3用設定ファイル出力フォルダ
- PAUP - PAUP*4用設定ファイル出力フォルダ
- Results - モデル選択結果のテキストファイル出力フォルダ
- Scores - 各モデルの尤度が記されたファイルが出力されるフォルダ
- Treefinder - Treefinder用設定ファイル出力フォルダ

どの設定ファイルを使うべきか - PAUP*4

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

どの設定ファイルを使うべきか - PAUP*4

- モデル選択規準 (AIC · AICc · BIC) とサンプルサイズ (1~6)
 - AICc4 (サイト数をサンプルサイズとするAICc)

どの設定ファイルを使うべきか - PAUP*4

- モデル選択規準 (AIC・AICc・BIC) とサンプルサイズ (1~6)
 - AICc4 (サイトをサンプルサイズとするAICc)
- 領域の区分 (区分無し・遺伝子座・コドン位置) と枝長の比例・分離
 - そもそもPAUP*4は対応していないので検討不要

どの設定ファイルを使うべきか - MrBayes3

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

どの設定ファイルを使うべきか - MrBayes3

- モデル選択規準 (AIC · AICc · BIC) とサンプルサイズ (1~6)
 - BIC4 (サイト数をサンプルサイズとするBIC)

どの設定ファイルを使うべきか - MrBayes3

- モデル選択規準 (AIC · AICc · BIC) とサンプルサイズ (1~6)
 - BIC4 (サイト数をサンプルサイズとするBIC)
- 領域の区分 (区分無し · 遺伝子座 · コドン位置) と枝長の比例 · 分離
 - `_proportional_codonshared` (コード領域を含む複数領域データの場合)
 - 遺伝子座ごとに異なる塩基置換モデル · 枝長は比例 · コドン位置区分無し
 - `_proportional` (コード領域を含まない複数領域データの場合)
 - 遺伝子座ごとに異なる塩基置換モデル · 枝長は比例
 - `_codonshared` (コード領域データの場合)
 - コドン位置区分無し
 - `_single` (非コード領域データの場合)
 - コドン位置区分無し

どの設定ファイルを使うべきか - Treefinder

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

どの設定ファイルを使うべきか - Treefinder

- モデル選択規準 (AIC · AICc · BIC) とサンプルサイズ (1~6)
 - AICc4 (サイト数をサンプルサイズとするAICc)

どの設定ファイルを使うべきか - Treefinder

- モデル選択規準 (AIC・AICc・BIC) とサンプルサイズ (1~6)
 - AICc4 (サイト数をサンプルサイズとするAICc)
- 領域の区分 (区分無し・遺伝子座・コドン位置) と枝長の比例・分離
 - `_proportional_codonproportional` (コード領域を含む複数領域データの場合)
 - 遺伝子座・コドン位置ごとに異なる塩基置換モデル・枝長は比例
 - `_proportional` (コード領域を含まない複数領域データの場合)
 - 遺伝子座ごとに異なる塩基置換モデル・枝長は比例
 - `_codonproportional` (コード領域データの場合)
 - コドン位置ごとに異なる塩基置換モデル・枝長は比例
 - `_single` (非コード領域データの場合)
 - コドン位置区分無し

まとめ

進化モデル選択とLikelihood Ratchet、系統樹から進化速度の変化を検出する

田辺 晶史 (東北大・院・生命科学)

2009年生態学会自由集会：進化生態学のための分子系統樹の推定と応用

まとめ

- より正確な系統推定のためには、モデル選択で選ばれたモデルを当てはめる必要がある

まとめ

- より正確な系統推定のためには、モデル選択で選ばれたモデルを当てはめる必要がある
- 複数領域・コード領域データでは、領域ごと・コドン位置ごとに異なる塩基置換モデルを当てはめることでモデル選択規準が改善することが多い

まとめ

- より正確な系統推定のためには、モデル選択で選ばれたモデルを当てはめる必要がある
- 複数領域・コード領域データでは、領域ごと・コドン位置ごとに異なる塩基置換モデルを当てはめることでモデル選択規準が改善することが多い
- Kakusan3により領域ごとに最適なモデルを選択し、それを当てはめた解析が可能

まとめ

- より正確な系統推定のためには、モデル選択で選ばれたモデルを当てはめる必要がある
- 複数領域・コード領域データでは、領域ごと・コドン位置ごとに異なる塩基置換モデルを当てはめることでモデル選択規準が改善することが多い
- Kakusan3により領域ごとに最適なモデルを選択し、それを当てはめた解析が可能
- Kakusan3は単一・比例・分離モデル間の比較は行わないが、これらを比較したい場合はそれぞれのモデルを当てはめた系統推定結果を比較すればよい