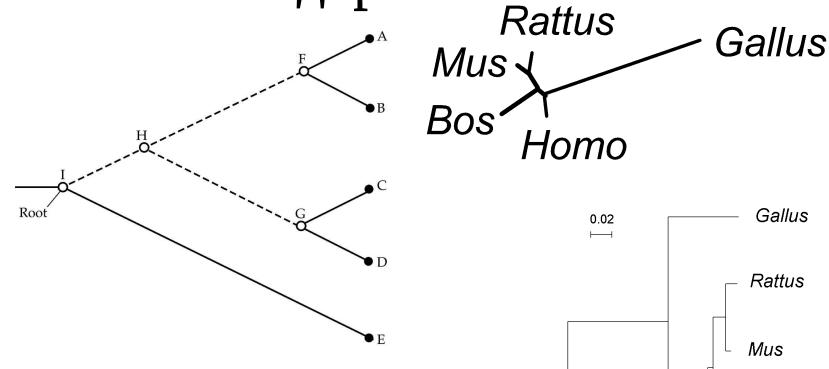
Построение филогенетических деревьев

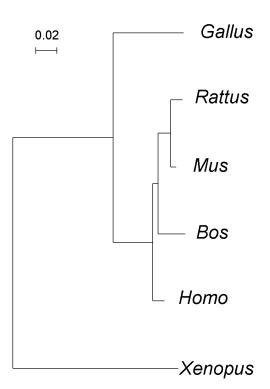
Особенности молекулярной эволюции

1. Скорость эволюции любого белка, выраженная через число аминокислотных замен на сайт в год, приблизительно постоянна и одинакова в разных филогенетических линиях, если только функция и третичная структура этого белка остаются в основном неизменными.

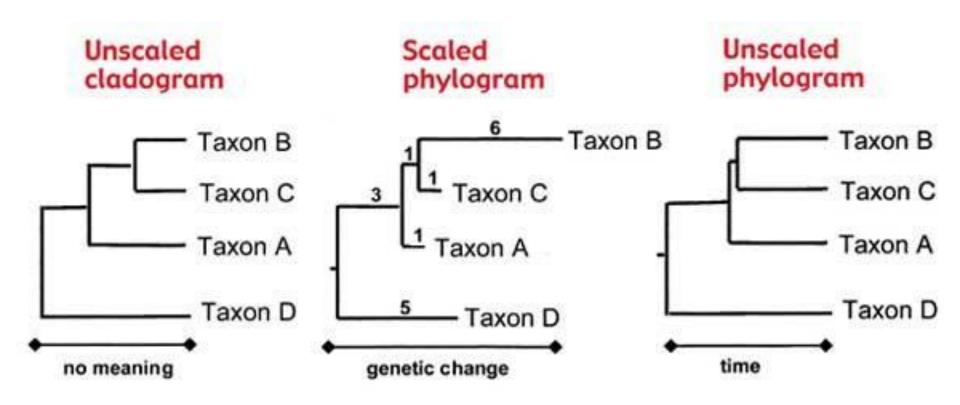
Что такое филогенетические деревья?



Дерево — это граф, в котором два соседних узла соединены только одним ребром.

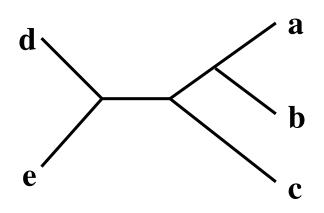


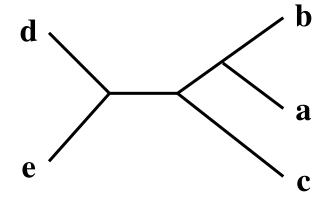
Кладограммы и филограммы

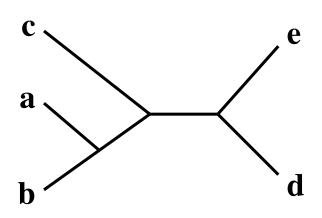


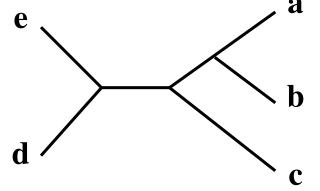
Кладограммы отражают только порядок ветвления, филограммы — ещё и длину ветвей

Сколько здесь разных кладограмм?









Выбор последовательностей

- Последовательности должны быть гомологичны! Программа выровняет любые последовательности => нужно проверить с помощью Blast
- Затем нужно выровнять последовательности, и по получившемуся выравниванию, определить, какие последовательности

DIZHIAHIIMI D AHAHIXA

Choosing appropriate molecular markers

Amplification, sequencing, assembly

Alignment

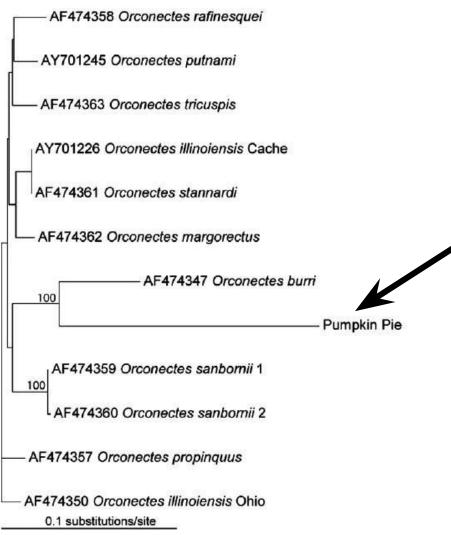
Evolutionary model

Phylogenetic analysis

Tree construction

Evaluation of phylogenetic tree

«Эффект тыквенного пирога»



Рецепт тыквенного пирога на филогенетическом дереве креветок.

Выбор

Поспеловательностей

```
DS FGDLS SASAIMGNAKVKAHGKKI
     DLSSASAIMGNPKVKAH
         TPDAVMGNPKVKA
  FGDLSSPDAVMGNPKVKAHGKKV
DK F G N L S S A L A I M G N P R I R A H G K K V
  F G N L S S A Q A I M G N P R I K A H G K K V
  FGNLSSASAIMGNPKVKAHGKK\
  F G N L S S A S A I M G N P K V K A H G K K V
  FGNLS S P S A I L GN P K V K A H G K K V
       LS S A S A I M G N P R V K A H G K K V
     DLHP----GSAQLRAHGS
       LHH----GSQQLRAHGFK
        S P - - - - - G S S Q V R A H G Q K
PHF-DLSH----GSAQVKGHGKKV
     DVSH----GSAQVKGHGKKV
```

Особенности молекулярной эволюции

- 2. Функционально менее важные молекулы или их части эволюционируют (накапливая эволюционные замены) быстрее, чем более важные
- 3. Мутационные замены, приводящие к меньшим нарушениям структуры и функции молекулы (консервативные замены), в ходе эволюции происходят чаще тех, которые вызывают существенное нарушение структуры и функции этой

Различия между деревом генов и деревом видов

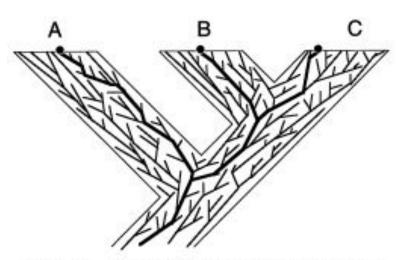


FIGURE 1. A gene tree contained within a species tree leading to three extant species: A, B, and C. Bold branches of gene tree show relationships among the sampled copies of the gene (●). Sampled copies from sister species B and C are sister copies.

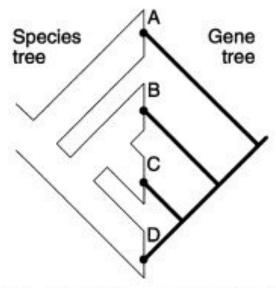
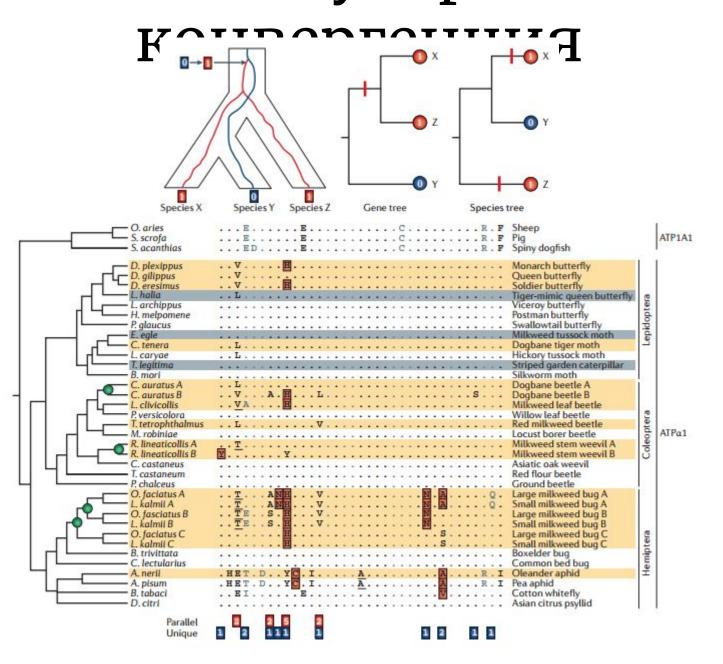


FIGURE 2. Discord between gene and species trees. At left is the species tree of four species, A, B, C, and D, and at right is the tree of a gene sampled one copy per species. Species B and C are sister species, but their gene copies are not sister copies.

Проблема: ортологи и паралоги

молекулярная



Филогенетические маркёры Свойства:

- •Гены, которые представлены одной копией в геноме лучше, чем те, у которых множество копий.
- •Длина гена не должна варьировать у разных организмов
- •Скорость изменения гена должна соответствовать скорости эволюции таксонов заданного уровня
- •Должны легко подбираться специфические праймеры

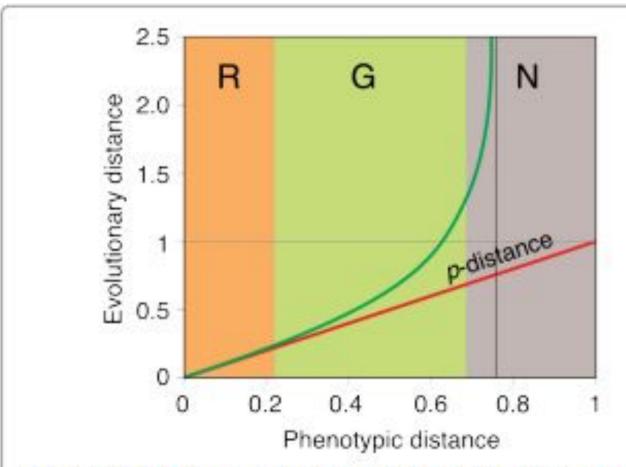


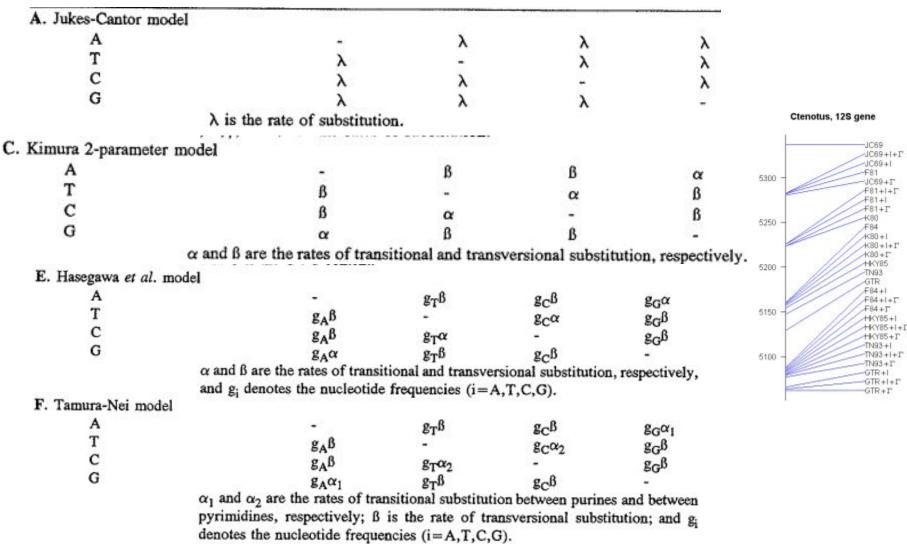
Figure 2: Regions of phenotypic distance corresponding to different estimates of evolutionary distance. R: "Parsimony zone", region where evolutionary distance is accurately approximated by phenotypic distance d. G: "Probabilistic zone", region where the p-distance under-estimates evolutionary distance. N: "Mutational saturation zone", region where the evolutionary distance cannot be estimated because of loss of phylogenetic information.

Филогенетические маркёры

- Рибосомальные гены
- Митохондриальные гены (COI/II, 12s RNA, cyt b)
- Хлоропластные гены
- Гены домашнего хозяйства и некоторые другие ядерные

Gene	Description	Reference
EF-1a	Elongation factor-1α, Role in protein synthesis.	[52]
rpoA gene	Encoding the alpha subunit of RNA polymerase	[53]
atpB	Encode the beta subunit of ATP synthase	[54]
dnaA	involved in DNA synthesis initiation	[55]
ftsZ	Role in cell division	[56]
gapA	Codes for glyceraldehyde phosphate dehydrogenase	[57]
groEL	Encodes bacterial heat shock protein.	[58]
gltA	Encoding citrate synthase	[59]
ITS	Piece of non-functional RNA situated between structural	[60]
	ribosomal RNAs precursor transcript.	8(8)
lux Gene	encode proteins involved in luminescence	[61]
PEPCK	Codes for phosphoenolpyruvate carboxykinase	[62]
pyrH genes	Codes for uridine monophosphate (UMP) kinases	[63]
recA	Role in recombination	[64]
U2 snRNA	Component of the spliceosome	[65]
Wsp gene	Encodes a major cell surface coat protein	[66]
Nuclear H3	Codes for protein which is associated with DNA	[67]
trnH-psbA	Non-coding intergenic spacer region located in plastid genome	[68]
rpoB, rpoC1	Coding region located in plastid genome	[69]

Выбор модели замен



Результаты вычисления эволюционных дистанций будут отличаться в зависимости от выбранной

Выбор модели замен

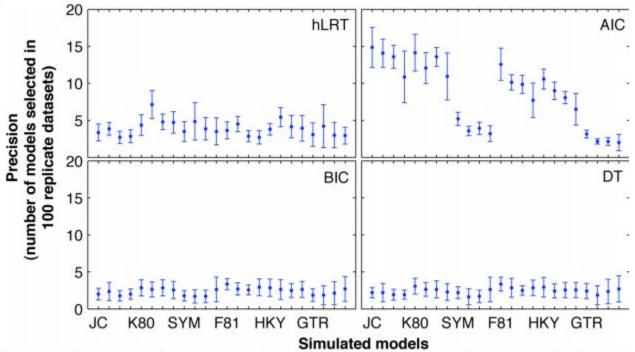


Figure 3 Precision of the four criteria corresponding to 24 simulated models. Categories along the x-axis represent the 24 simulated models. For the sake of clarity, only seven models are labelled, and each one is followed by three similar ones (e.g., JC is followed by JC + I, $JC + \Gamma$, and $JC + I + \Gamma$). The y-axis represents the means and standard deviations of precision values for each simulated model across the 14 simulations, which are different statistical results from those in Additional file 2. The markers denote the means, while lengths of error bars denote the standard deviation values.

- AIC Akaike's Information Criterion. Быстрее
- BIC Bayesian information criteria. He «любит» более сложные модели
- DT decision theory
- LRT тест соотношения вероятностей. «Любит» более сложные модели.

Types of data used in phylogenetic inference:

Characters

Species	A	ATGGCTATTCTTATAGTACG
Species	B	ATCGCTAGTCTTATATTACA
Species	C	TTCACTAGACCTGTGGTCCA
Species	D	TTGACCAGACCTGTGGTCCG
Species	E	TTGACCAGTTCTCTAGTTCG

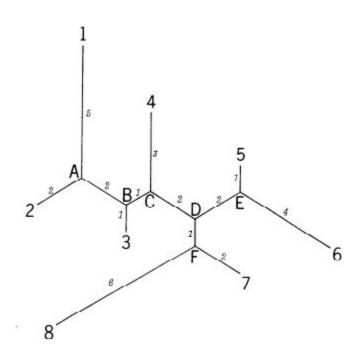
Distances

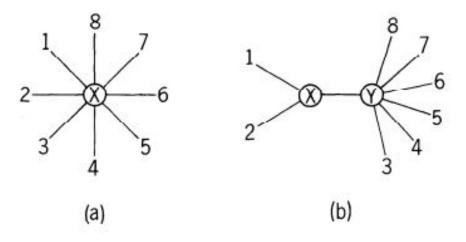
	A	В	С	D	E		
Species A		0.20	0.50	0.45	0.40]	
Species B	0.23		0.40	0.55	0.50		Example 1:
Species C	0.87	0.59		0.15	0.40	-	Uncorrected "p" distance
Species D	0.73	1.12	0.17		0.25	1	(=observed percent
Species E	0.59	0.89	0.61	0.31			sequence difference)

Example 2: Kimura 2-parameter distance (estimate of the true number of substitutions between taxa)

Методы реконструкции

Оилогении								
Дистанционные	Максимальной	— Максимальной Ворожние						
	экономии	вероятности						
Используют только	Используют	Используют все						
попарные	только	данные						
дистанции	символьные							
	данные							
Минимизация	Минимизация	Максимизация						
дистанции между	общей длины	вероятности						
ближайшими	дерева	заданного дерева с						
соседями	(минимизация	учётом заданных						
	числа мутаций	параметров						
Очень быстрые	Медленные	Очень медленные						
Ищут локальный	Неверны при	Сильно зависят от						
оптимум вместо	быстрой скорости	правильности						
глобального	эволюции	выбранной модели						
Хороши для	Лучший выбор	Хороши для очень						
чернового или	для подходящей	маленьких наборов						
препварительного	BPIQUEM(<30	панних и ппа опенки						





Начинаем с пары ветвей, которые меньше всего отличаются между собой

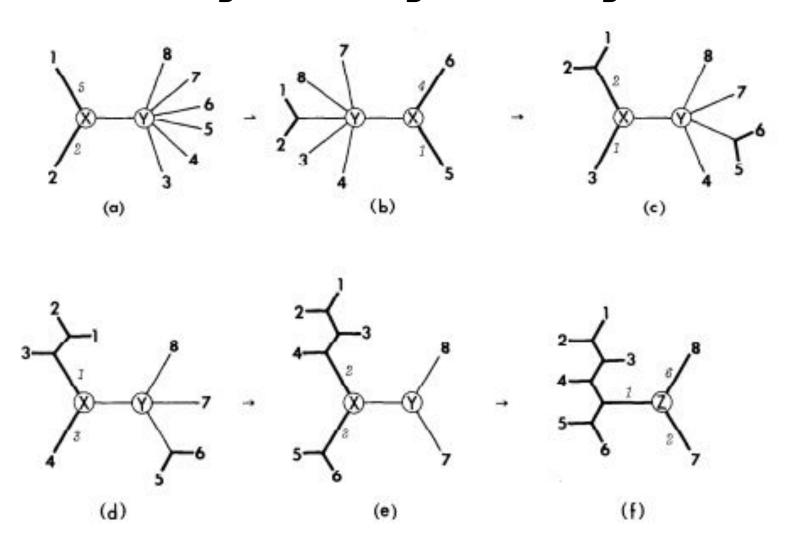
$$S_O = \sum_{i=1}^{N} L_{iX} = \frac{1}{N-1} \sum_{i < j} D_{ij},$$

$$L_{XY} = \frac{1}{2(N-2)} \left[\sum_{k=3}^{N} (D_{1k} + D_{2k}) - (N-2)(L_{1X} + L_{2X}) - 2 \sum_{i=3}^{N} L_{iY} \right].$$

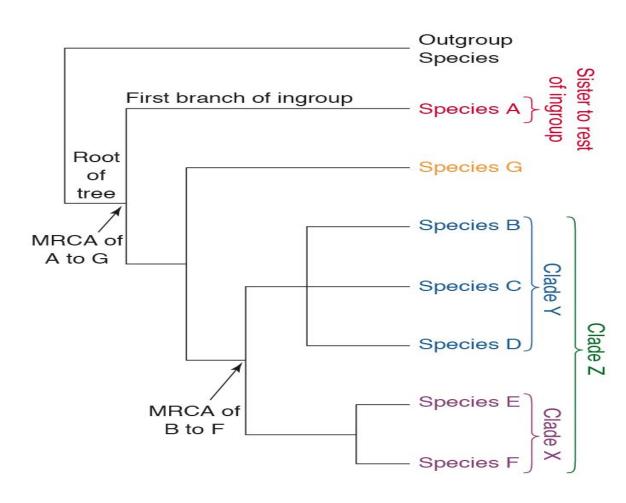
$$L_{1X} + L_{2X} = D_{12}$$

$$\sum_{i=3}^{N} L_{iY} = \frac{1}{N-3} \sum_{3 \le i < j} D_{ij}.$$

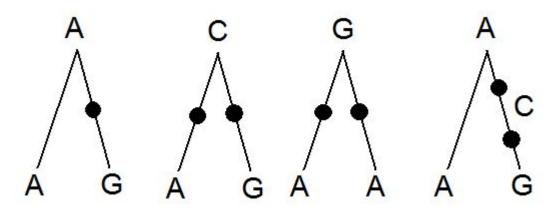
		A.	Cycle 1: Neig	ghbors = [1, 2]		And An
9 60 -				OTU			guest
OTU	1	2	3	4	5	6	ournals.drg/ by/gueston March 25, 2015
2	36.67	200000000					rch 2
3	38.33	38.33					3
4	39.00	39.00	38.67				201
5	40.33	40.33	40.00	39.67			C)
6	40.33	40.33	40.00	39.67	37.00		
7	40.17	40.17	39.83	39.50	38.83	38.83	58
8	40.17	40.17	39.83	39.50	38.83	38.83	37.67
		В.	Cycle 2: Neig	ghbors = [5, 6]		
				OTU			
OTU	1-2	3	4		5	6	7
3	31.50						
4	32.30	32.30					
5	33.90	33.90	33.70	0			
6	33.90	33.90	33.70	0	31.30		
7	33.70	33.70	33.50	0	33.10	33.10	
8	33.70	33.70	33.50	0	33.10	33.10	31.90



Зачем нужна аутгруппа



Молекулярнофилогенетические методы используют информацию о последовательностя х внешней группы (контроля), дистанция от которой для всех остальных последовательносте й заведомо выше, чем от других. Таким образом дерево «укореняется», а также внутри лерева убирается



- •Не учитываются обратные и параллельные замены
- => Мы считаем не настоящую дистанцию (расстояние), а редакционное расстояние.
- •Вычислительно более быстрые.
- •В большинстве случаев оценивают только топологию дерева, не воспроизводя исходную последовательность.
- •Если у нас будет бесконечная последовательность, то мы с вероятностью 100% получим истинное

Методы максимальной экономии

- •Минимизация числа замен символов
- •Всегда реконструируют предковые

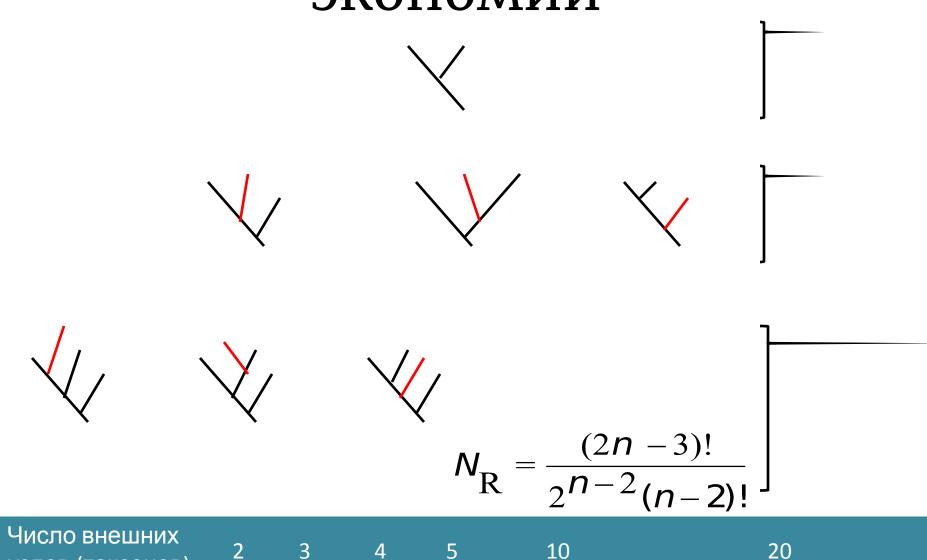
последовательности

•Лучше работает на небольших наборах последовательностей во многих случаях на больших объёмах данных работает хуж

Let S be a set of n sequences, each of length n, over a fixed alphabet Σ . Let T be a tree leaf-labelled by the set Sand with internal nodes labelled by sequences of length n over Σ . The length (or parsimony score) of T with this labelling is the sum, over all the edges, of the Hamming distances between the labels at the endpoints of the edge. (The Hamming distance between two strings of equal length is just the number of positions in which the two strings differ.) Thus the length of a tree is also the total number of point mutations along the edges of the tree. The Maximum Parsimony (MP) problem seeks the tree T leaf-labelled by S with the minimum length. While MP is NP-hard [4], constructing the optimal labeling of the internal nodes of a fixed tree T can be done in polynomial time [3].

TATE TOMBI MIGHT CATAMAN TO LICE

ЭКОНОМИИ



Число внешних узлов (таксонов)	2	3	4	5	10	20
Число возможных	1	3	15	105	34459425	8200794532637891559375
BOOMO/KIIBIX		<u> </u>	10	103	37733723	0200754552057051555575

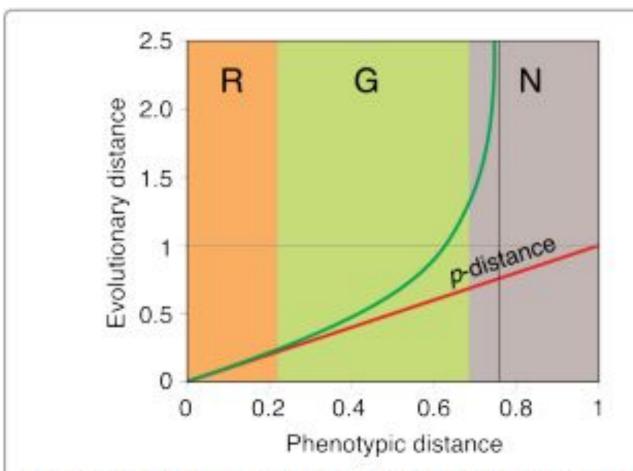
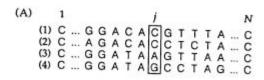


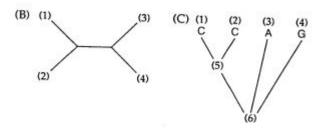
Figure 2: Regions of phenotypic distance corresponding to different estimates of evolutionary distance. R: "Parsimony zone", region where evolutionary distance is accurately approximated by phenotypic distance d. G: "Probabilistic zone", region where the p-distance under-estimates evolutionary distance. N: "Mutational saturation zone", region where the evolutionary distance cannot be estimated because of loss of phylogenetic information.

Методы максимальной вероятности

- •Так же, как и в случае с методами максимальной экономии, генерирует все возможные топологии деревьев
- •Предположение особой модели эволюции
- •В отличие от метода максимальной экономии может предполагать разную скорость эволюции и скорость замен в разных ветвях дерева
- •Поиск дерева с максимальной вероятностью существования, соответствующего данным
- •Чем больше последовательность, тем вероятнее найти истинное дерево
- •Самые медленные

Методы максимальной вероятности





(D)
$$L_{(j)} = \operatorname{Prob} \begin{pmatrix} C & C & A & G \\ A & A & A \end{pmatrix} + \operatorname{Prob} \begin{pmatrix} C & C & A & G \\ C & A & G \end{pmatrix}$$

$$+ \dots + \operatorname{Prob} \begin{pmatrix} C & C & A & G \\ G & A & G \end{pmatrix}$$

$$+ \dots + \operatorname{Prob} \begin{pmatrix} C & C & A & G \\ G & A & G \end{pmatrix}$$

(E)
$$L = L_{(1)} \cdot L_{(2)} \cdot \ldots \cdot L_{(N)} = \prod_{j=1}^{N} L_{(j)}$$

(F)
$$\ln L = \ln L_{(1)} + \ln L_{(2)} + ... + \ln L_{(N)} = \sum_{j=1}^{N} \ln L_{(j)}$$

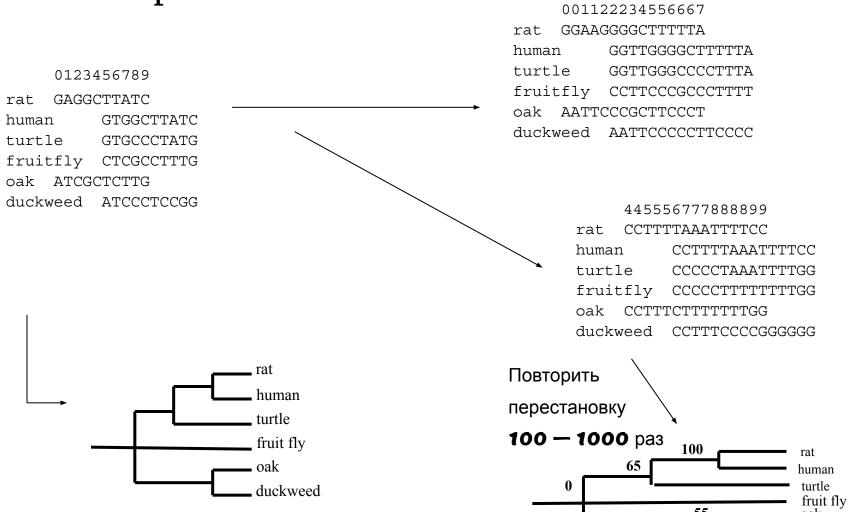
- В позиции ј для каждого внутреннего узла допустимы все четыре нуклеотида, значит всего 4*4=16 возможных деревьев.
- Каждое из деревьев это произведение вероятности возникновения какого-либо основания в корне дерева и вероятность его замены на тот, который в следующем узле. Т.е. частота нуклеотида умноженная на вероятность его мутации, если грубо.

A = 0.25 ог средняя частота A в последовательности зависит от модели) f A->C трансверсия = 10^{-6} and A->G транзиции = $2x10^{-6}$ f Вероятность T1 = 0.25 x 2x10-6 x 10-6 = 5x10-13

Оценка поддержки дерева

Bootstrap

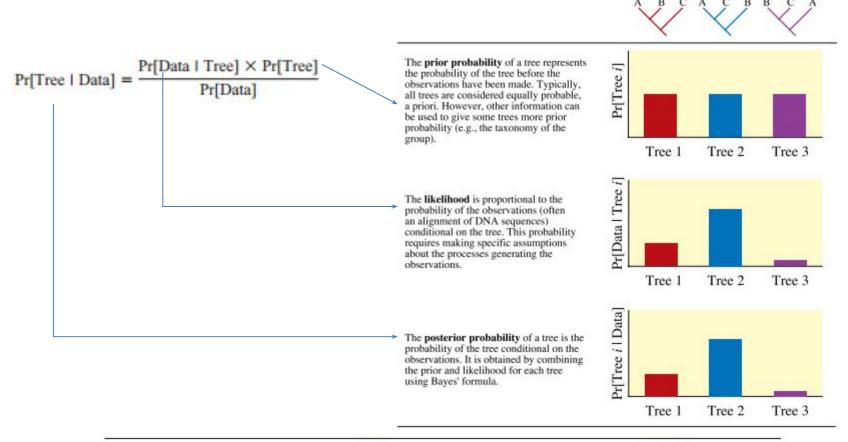
Inferred tree



duckweed

Оценка поддержки дерева

Bayes inference



Support category (%)		BA	YES		MLBOOT			
	Correct model		Incorrect model		Correct model		Incorrect model	
	N	% wrong	N	% wrong	N	% wrong	N	% wrong
70.1-80.0	418	17.3	262	54.6	612	10.8	634	15.3
80.1-90.0	542	10.2	337	37.4	713	4.5	622	8.5
90.1-95.0	409	6.2	292	40.4	399	2.0	301	4.7
95.1-100.0	1,216	1.7	2,622	15.6	505	0.2	293	2.0