If you don't clear your misunderstanding in time they become the reason for distance forever!!!



PHYLOGENETIC ANALYSIS

Amira A. AL-Hosary PhD of infectious diseases Department of Animal Medicine (Infectious Diseases) Faculty of Veterinary Medicine Assiut University-Egypt

Phylogenetic Basics:

- Biological sequence analysis is founded on solid evolutionary principles.
- Similarities and divergence among related biological sequences revealed by sequence alignment often have to be rationalized and visualized in the context of phylogenetic trees.
- Thus, molecular phylogenetic is a fundamental aspect of bioinformatics.

MOLECULAR EVOLUTION AND MOLECULAR PHYLOGENETICS

"What is evolution?"

- In the biological context, evolution can be defined as the development of a biological form from other pre existing forms or its origin to the current existing form through natural selections and modifications.
- The driving force behind evolution is natural selection in which "unfit" forms are eliminated through changes of environmental conditions or sexual selection so that only the fittest are selected.
- The underlying mechanism of evolution is genetic mutations that occur spontaneously.

Phylogenetic?

- It is the study of the evolutionary history of living organisms using tree like diagrams to represent pedigrees of these organisms.
 The tree branching patterns representing the evolutionary divergence are referred to
 - as phylogeny.

Methods of phylogenetic analysis:

- I. It is often studied using fossil records, which contain morphological information about <u>ancestors</u> of current species and the timeline of divergence.
- II. However, fossil records have many limitations; they may be available only for certain species.

Ancestors

For microorganisms,

Fossils are essentially nonexistent, which makes it impossible to study phylogeny with this approach.

- Molecular data: Like DNA and/or protein sequences can also provide very useful evolutionary perspectives of existing organisms because the genetic materials accumulate mutations over time causing phenotypic changes.
- Genes are the medium for recording the accumulated mutations, they can serve as molecular fossils.
- Molecular fossils

Through comparative analysis of the molecular fossils from a number of related organisms allow us to follow the evolutionary history of the genes and even the organisms.

The advantage of using molecular data:

- 1- Molecular data are more numerous and easier to obtain.
- 2- There is no sampling bias involved, which helps to mend the gaps in real fossil records.
- 3- More clear-cut and robust phylogenetic trees can be constructed with the molecular data.
- Therefore, they have become favorite and sometimes the only information available to reconstruct evolutionary history.
- The advent of the genomic era with tremendous amounts of molecular sequence data has led to the rapid development of molecular phylogenetic.

Molecular phylogenetic:

Defined as:

The study of evolutionary relationships of genes and other biological macromolecules by analyzing mutations at various positions in their sequences and developing hypotheses about the evolutionary relatedness of the biomolecules based on the sequence similarity

Major Assumptions

To use molecular data to reconstruct evolutionary history requires making a number of reasonable assumptions.

- 1. The molecular sequences used in phylogenetic construction are homologous, meaning that they share a common origin and subsequently diverged through time.
- 2. Phylogenetic divergence is assumed to be bifurcating, meaning that a parent branch splits into two daughter branches at any given point.
- 3. The variability among sequences is sufficiently informative for constructing unambiguous 23-4-2015 phylogenetic trees Dr. Amira A. AL-Hosary

TERMINOLOGY

- phylogenetic tree
- The lines in the tree are called branches.
- At the tips of the branches are presentday species or sequences known as taxa or Terminal point (the singular form is taxon).
- The connecting point where two adjacent branches join is called a node, which represents an inferred ancestor of extant taxa.
- The bifurcating point at the very bottom of the tree is the root node = (Common Ancestor) of all members of

23the tree.

×internal node branch A typical bifurcating phylogenetic tree showing root, internal nodes, terminal nodes and branches.

taxa (or terminal nodes)

D

В

Types of phylogenetic trees:





23-4-2015 Unrooted



Conversion of the un rooted tree to rooted:

In practice, however, it is often desirable to define the root of a tree.

There are two ways to define the root of a tree:

One is to use an <u>out group</u>, which is a sequence that is homologous to the sequences under consideration, but separated from those sequences at an early evolutionary time.

Out groups are generally determined from independent sources of information. *For example*, a bird sequence can be used as a root for the phylogenetic analysis of mammals based on multiple lines of evidence that indicate that birds branched off prior to all mammalian taxa in the in group.



Midpoint rooting approach

In the absence of a good out group, a tree can be rooted using *Midpoint rooting approach*, in which the midpoint of the two most divergent groups judged by overall branch lengths is assigned as the root.

This type of rooting assumes that divergence from root to tips for both branches is equal and follows the "molecular clock" hypothesis.

Molecular clock is an assumption by which molecular sequences evolve at constant rates so that the amount of accumulated mutations is proportional to evolutionary time.

Based on this hypothesis, branch lengths on a tree can be used to estimate divergence time.



GENE PHYLOGENY VS SPECIES PHYLOGENY

Gene phylogeny: (phylogeny inferred from a gene or protein sequence) only describes the evolution of that particular gene or encoded protein. This sequence may evolve more or less rapidly than other genes in the genome or may have a different evolutionary history from the rest of the genome.

The species evolution is the combined result of evolution by multiple genes in a genome.

In a species tree, the branching point at an internal node represents the speciation event whereas, in a gene tree, the internal node indicates a gene duplication event.

Thus, to obtain a species phylogeny, phylogenetic trees from a variety of gene families need to be constructed to give an overall <u>assessment of the species evolution</u>

FORMS OF TREE REPRESENTATION:

The topology of branches in a tree defines the relationships between the taxa. The trees can be drawn in different ways, such as a cladogram or a phylogram.

In a phylogram, the branch lengths represent the amount of evolutionary divergence. Such trees are said to be scaled. The scaled trees have the advantage of showing both the evolutionary relationships and information about the relative time of the divergence branches.

In a cladogram, however, the external taxa line up neatly in a row or column. Their branch lengths are not proportional to the number of evolutionary changes and thus have no phylogenetic meaning. In such an un scaled tree, only the topology of the tree matters, which shows the relative ordering of the taxa.

Phylogenetic trees drawn as cladograms (top) and phylograms (bottom). The branch lengths are un scaled in the cladograms and scaled in the phylograms. The trees can be drawn as angled form (left) or squared form (right).

PROCEDURE of Molecular phylogenetic tree construction :

It is divided into five steps:

(1) Choosing molecular markers.

(2) Performing multiple sequence alignment.

(3) Choosing a model of evolution.

(4) Determining a tree building method.

(5) Assessing tree reliability.

Choice of Molecular Markers

- 1. For constructing molecular phylogenetic trees, one can use either nucleotide or protein sequence data.
- 2. The choice of molecular markers is an important matter because it can make a major difference in obtaining a correct tree.
- 3. The decision to use nucleotide or protein sequences depends on the properties of the sequences and the purposes of the study.

For example, Studying very closely related organisms, nucleotide sequences, which evolve more rapidly than proteins.

Alignment

The second step in phylogenetic analysis: This is probably the most critical step in the procedure because it establishes positional correspondence in evolution.

Only the correct alignment produces correct phylogenetic inference because aligned positions are assumed to be genealogically related.

Incorrect alignment leads to systematic errors in the final tree or even a completely wrong

Choosing Substitution Models The statistical models used to correct homoplasy are called substitution models or evolutionary models. The most common substitution models or evolutionary models: 1- Jukes-Cantor Model. 2- Kimura Model.

Jukes–Cantor Model

The simplest nucleotide substitution model is the Jukes–Cantor model, which assumes that all nucleotides are substituted with equal probability. This model can only handle reasonably closely related sequences.

A formula for deriving evolutionary distances that include hidden changes is introduced by using a logarithmic function.

$dAB = -(3/4) \ln [1 - (4/3) pAB]$

where *d* AB is the evolutionary distance between sequences A and B and *p* AB is the observed sequence distance measured by the proportion of substitutions over the entire length of the alignment.

Kimura Model

Another model to correct evolutionary distances is called the Kimura two-parameter model.

This is a more sophisticated model in which mutation rates for transitions and transversion are assumed to be different, which is more realistic.

According to this model, transitions occur more frequently than trans versions, which, therefore, provides a more realistic estimate of evolutionary distances.

The Kimura model uses the following formula:

d AB = - (1/2) ln (1 - 2pti - p tv) - (1/4) ln (1 - 2ptv)

where *d* AB is the evolutionary distance between sequences A and B, pti is the observed frequency for transition, and p tv the frequency of transversion.

Jukes-Cantor model

Kimura model

The Jukes–Cantor and Kimura models for DNA substitutions. In the Jukes–Cantor model, all nucleotides have equal substitution rates (α). In the Kimura model, there are unequal rates of transitions (α) and trans versions (β).

Determining a tree building method.

Several methods are available for reconstructing phylogenetic trees.

Most of them use some criterion for evaluating the fit of a given data set to the topology and then search for the tree that gives the best score in terms of that criterion.

If the criterion used is realistic and the data are sufficient, the tree should represent the true phylogenetic relationship of the sequences.

There are two methods:-

1- DISTANCE-BASED METHODS

True evolutionary distances between sequences can be calculated from observed distances after correction using a variety of evolutionary models. The computed evolutionary distances can be used to construct a matrix of distances between all individual pairs of taxa

- \rightarrow Neighbour-joining (NJ) methods.
- 2- CHARACTER-BASED METHODS.

also called (discrete methods) based directly on the sequence characters rather than on pairwise distances. They count mutational events accumulated on the sequences and may therefore avoid the loss of information when characters are converted to distances.

Maximum parsimony (MP) methods / Maximum likelihood (M12) methods.

Neighbor-joining (NJ) methods.

MM4.1 (Beta 3): Analysis Preferences		_ 🗆 ×	
Options Summary Test of Phylogeny			
Option	Selection		
Data Type	Amino acid		
Analysis	Phylogeny reconstruction		
Tree Inference			
->Method	NeighborJoining		
->Phylogeny Test and options	Bootstrap (500 replicates; seed=64238)		
Include Sites			
->Gaps/Missing Data	Complete Deletion		
Substitution Model			
->Model	Amino: p-distance		
->Substitutions to Include	All	18-1	
->Pattern among Lineages	Same (Homogeneous)		
->Rates among sites	Uniform rates		
23-4-2015	Dr. Amira A. AL-Hosary	Compute X Cancel ? Help	

Maximum parsimony (MP) methods:

The topology requiring the smallest number of nucleotide changes to fit the observed sequence data is chosen to represent the true tree.

Maximum likelihood (ML) methods:

The topology with the greatest likelihood under a given probabilistic model of nucleotide substitutions is chosen.

Molecular Ev	olutionary Genetics Analysis, Version 2.1 📃	٦×
<u>File D</u> ata Di <u>s</u> ta	nces <u>Phylogeny T</u> ests Windows <u>H</u> elp	
	Image: PGMA [™] Meighbor-Joining (NJ)	
Go to the MEGA	web <u>Za</u> Minimum Evolution (ME)	
Citing MEGA2 in	<u>public</u> ∑si Maximum Parsimony (MP)	
	Display Saved Tree Session	
	· · · · · · · · · · · · · · · · · · ·	
Data File	L:\transfer\MEGAtransfer\us90io.MEG	
Title	exported by MacClade from file USTIOC90_BICio.clade	•
ſ	J 11:46:21 AM I	

Assessing tree reliability.

- After phylogenetic tree construction, the next step is to statistically evaluate the reliability of the inferred phylogeny.
- There are two questions that need to be addressed.
- One is how reliable the tree or a portion of it. Second is whether this tree is significantly better than another tree.

What Is Bootstrapping?

Bootstrapping is a statistical technique that tests the sampling errors of a phylogenetic tree.

M M4.1 (Beta 3): Analysis Preferences			
Options Summary			
Option	Selection		
Data Type	Amino acid		
Analysis	Phylogeny reconstruction		
Tree Inference			
->Method	NeighborJoining		
->Phylogeny Test and options	Bootstrap (500 replicates; seed=64238)		
Include Sites			
->Gaps/Missing Data	Complete Deletion		
Substitution Model			
->Model	Amino: p-distance		
->Substitutions to Include	All		
->Pattern among Lineages	Same (Homogeneous)		
->Rates among sites	Uniform rates		
23-4-2015	Dr. Amira A. AL-Hosary	Compute X Cancel ? Help	

Schematic representation of a bootstrap analysis showing the original alignment and modified replicates in which certain sites are randomly replaced with other existing sites. The resulting altered replicates are used to building trees for statistical analysis at each node.

- The bootstrap test provides a measure for evaluating the confidence levels of the tree topology.
- Analysis has shown that a bootstrap value of 70% approximately corresponds to 95% statistical confidence, although the issue is still a subject of debate.
- Bootstrapping does not assess the accuracy of a tree, but only indicates consistency and stability of individual clades of the tree. This means that, because of systematic errors, wrong trees can still be obtained with high bootstrap values.
- It is generally recommended that a phylogenetic tree should be bootstrapped 500 to 1,000 times

All Sequences alignment:

		ClustalX (1.83)			
File	Edit Alignment	Trees Colors Quality Help			
_					
Mu	Multiple Alignment Mode D Font Size: 10 D				
	NF01242439	MLVSCSREVPHEALERLNIKTVDAMKSGAHVDKPSTVLFADIVNETAMANVSDVPKTFC			
2	2 NF01242241	MPNINKSKVVSYPONOMYELVNDVESYSETVPECSESEIDSCTHEEIRATLSEARGGES	K		
	NF00212957	MTHDVLLAGAGLANGLIALALRAARPDLRVLLLDHAAGPSDGHTWSCHDPDLSPDWLAF	IL		
4	1 NF01863930	PTLNLQLDENNEQLEKVTKELEFERTKTESVLMSILPPTIANHLINNEHIEAREFEHAT	v I		
	NF01863928	PTLNLQLDENNEQLE <mark>KVTKELEFERTKTESVLMSILPPTI</mark> AN <mark>HLI</mark> NNEHIEARMLWEAF	L		
6	5 NF01863927	PTLNLQLDENNEQLEKVTKELEFERTKTESVLMSILPPTIANHLINNEHIEARMLWEAF	L		
	7 NF01863924	PTLNLQLDENNEQLEKVTKELEFERTKTESVLMSILPPTIANHLINNEHIEARMLWEAF	L		
5	3 NF01863923	PTLNLQLDENNEQLEKVTKELEFERTKTESVLMSILPPTIANHLINNEHIEARMLWEAF	L		
9	NF01412383	MPLDDSAWLGEYKKALDESNIVSITDLEGRIVYANEKFCEISGYSQEELLGKPHNIIRF	IP		
10) NF01411886	MESVSKNSSILIVEDEERARIDTQSLLSKHYERVLSAGSAKEALELYYAHKPDIFIVDI	Ē		
11	L NF01411832	MSRSNLLGILLVAILGLEGFLYFKPIPLLLSLEHKIKDAMFWWRGERAGNPNILIIDI)E		
12	2 NF01411664	MEDEENIRQVLAKILLRKFKEVVTASNGLEGLERYMENKPDIVITDIRMPEMSGIEMSF	RR		
13	3 NF01410974	MNPLAKRIIPCLDIKEGRVVKGVNFLGLRDAGDPVEVARRYNEEGADE IAFLDITATHE	S <mark>R</mark>		
14	1 NF00707828	MLAKRIIPCLDVRDGQVVKGVQFRNHEIIGDIVPLAKRYADEGADELVFYDITASSDGF	🔽		
15	NF00751159	MNSAEAFMKEDQHTYDRLMHIIESVVQTGQLNRMFPLVKYTEMQIEGYDAAYGYDKRNV	N		
16	5 NF00570218	MGCTVSAEDKAAAERSKMIDKNLREDGEKAAREVKLLLLGAGESGKSTIVKOMKIIHEI	G		
1	7 NF00570580	MNTELLSALCLGAWAALVGAVTVQDGDESESLESVKKLKDLQEAPESKVQGRRKEVAPE	2L		
18	3 NF00570504	MGCTLSAEDKAA <mark>VERSKMI</mark> DRNLREDGEKAAREVKLLLLGAGESGKSTIVKOMKIIHEA	4 <mark>G</mark>		
19	NF00570550	MMGVNSSGRPDLYGHLHSILLPGRGLPDWSPDGGADPGVSTWTPRLLSGVPEVAASPSF	2S		
	ruler	110	50		
V .			n I		
			\geq		

23-4-2015

MA.1 (Beta 3): Analysis Preferences		<u>×</u>	
Options Summary Lest of Phylogeny			
Option	Selection		
Data Type	Amino acid		
Analysis	Phylogeny reconstruction		
Tree Inference			
->Method	Neighbor-Joining		
->Phylogeny Test and options	Bootstrap (500 replicates; seed=64238)		
Include Sites			
->Gaps/Missing Data	Complete Deletion		
Substitution Model			
->Model	Amino: p-distance	•••	
->Substitutions to Include	All	18	
->Pattern among Lineages	Same (Homogeneous)		
->Rates among sites	Uniform rates		
That a faile de la deve de la contra			
		Compute X Cancel ? Help	

"Our prime purpose in this life is to help others. And if you can't help them, at least don't hurt them."

Always laugh when you can. It is cheaper than medicine.

COVERS AT FIRSTCOVERS.COM

Thanks a lot

with my Best Regards and My Best wishes

Amira A. AL-Hosary E-mail: Amiraelhosary @yahoo.com Mob. (002) 01004477501