

Phylogeny: building the tree of life

Dr. Fayyaz ul Amir Afsar Minhas

PIEAS Biomedical Informatics Research Lab Department of Computer and Information Sciences Pakistan Institute of Engineering & Applied Sciences PO Nilore, Islamabad, Pakistan http://faculty.pieas.edu.pk/fayyaz/

Acknowledgements: Most of these slides are taken from Lectures by Dr. Asa-Ben-Hur, CSU.

CIS 529: Bioinformatics

Phylogeny and Phylogenetic Trees

- The biological evidence suggests that all life evolved from a common ancestor
- Phylogeny is the study of the evolutionary histroy of a group of species (taxa)
- The relationship can be represented using a phylogenetic tree



CIS 529: Bioinformatics

Phylogeny and Phylogenetic Trees

- The leaves of the tree corresponds to extant taxa
- Internal nodes correspond to speciation events, and represent ancestral taxa.
- Branches define the relationships among the taxa (descent and ancestry)
- A phylogenetic tree is a binary tree.
- The branch length usually represents the amount of time since speciation.



CIS 529: Bioinformatics

Objectives

- Reconstruct the ancestry of species, i.e. find the tree of life and
- To estimate the time of divergence between organisms since they last shared a common ancestor.



Cladograms and Phylogenetic trees from Phonotypical characteristics

• See video!



CIS 529: Bioinformatics

A Phylogenetic Tree



- A (rooted) tree for 5 species.
- Leaves α , β , γ , δ , and ϵ , correspond to contemporary species, for which data has been collected
- Internal nodes correspond to (inferred) ancestors

CIS 529: Bioinformatics

Using Morphological Characteristics

	Characters					
Species	1	2	3	4	5	6
α	1	0	0	1	1	0
eta	0	0	1	0	0	0
γ	1	1	0	0	0	0
δ	1	1	0	1	1	1
ϵ	0	0	1	1	1	0

Here, the Os and 1s are indicating the presence or absence of a character (has feathers?, lays eggs?, curved beak?, flies?, . . .).

Raccoon or Bear ?



- Based on anatomical and behavioral characteristics, the panda was originally classified as a raccoon (1870)
- In 1985, the panda was re-classified as a bear when an analysis based on molecular data was done



CIS 529: Bioinformatics

Platypus

The Platypus (Ornithorhynchus anatinus) is a venomous, egg-laying, duck-billed mammal; the sole representative of its family (Ornithorhynchidae) and genus (Ornithorhynchus)



commons.wikimedia.org/wiki/Image:Ornithorhynchidae-00.jpg

CIS 529: Bioinformatics

Need for Other Approaches

Hard to resolve relationships using morphology and behavior alone

- 1. Similar characteristics can evolve independently in distantly related organisms convergent-evolution
- 2. It is often difficult to find characteristics that are common to all the organisms under study

Inferring Phylogenies

Trees can be inferred by several criteria:

- Morphology of the organisms
 - Can lead to mistakes!
- Sequence comparison

Example:









GIS 525: BIOIMORMALICS



CIS 529: Bioinformatics

Earliest-known Human Footprints

The footprints were most likely made by *Australopithecus afarensis*, an early human whose fossils were found in the same sediment layer.

The entire footprint trail is almost 27 m (88 ft) long and includes impressions of about 70 early human footprints.

The First Step: Walking Upright

while a star or so any other has been a star of the star and the best of the star of the stars. A star the star of the stars of the star while have the star of the star where and the star of the star

Testing water to specie to main a variety of invitaments

Production report Descriptions of the order of the section of the order of the section of the order of the section of the section of the sec-

Graphy Super-Interest on Super-Latin Statistics (Statistics)

Angled Steel Sector and States and Steel Sector and States and States and Street States and States and States Street States and States and States

Period in the second se

What is phylogeny good for ?

- Evolutionary history ("tree of life")
- Population history
- Origins of diseases
- Prediction of sequence function
- MSA
- Can be applied to organisms, sequences, viruses,

languages, etc.



The Tree of Life and Human Health

Understanding how organisms, as well as their genes and gene products, are related to one another has become a powerful tool for identifying disease organisms, tracing the history of infections, and predicting disease outbreaks.

Identifying Emerging Diseases: West Nile Virus

When an encephalitis-like viral infection emerged in people living in the New York region in 1999, it was first suspected to be the St. Louis encephalitis virus. Transmitted by mosquitos, the virus was simultaneously found to be associated with a high mortality in wild and domesticated birds. It has currently spread as far west as California and has resulted in numerous human deaths.

In two separate studies, health workers used phylogenetic analysis to identify viral isolates from mosquitos and birds as a new outbreak of the West Nile virus rather than St. Louis encephalitis. The viral tree (right) demonstrated that the New York isolate was most closely related to one found in dead birds in Israel, East Africa, and eastern Europe. This knowledge provided health officials with key information about the basic biology of the virus that was needed for diagnosis and predicting its spread. Such knowledge was critical in preventing human and animal infection.





CIS 529: Bioinformatics



Lanciotti et al. (1999)

CIS 529: Bioinformatics



Bird flu phylogeny

N.J. Distance Tree of H5N1 (chicken) segment 5 using Jukes-Cantor model



Zika Virus



Lanciotti, Robert S., Amy J. Lambert, Mark Holodniy, Sonia Saavedra, and Leticia del Carmen Castillo Signor. 2016. "Phylogeny of Zika Virus in Western Hemisphere, 2015." *Emerging Infectious Diseases* 22 (5). doi:10.3201/eid2205.160065.

CIS 529: Bioinformatics





Which Whales Are Hunted? A Molecular Genetic Approach to Monitoring Whaling









CIS 529: Bioinformatics

09-11-01

THIS IS NEXT TAKE PENACILIN NOW DEATH TO AMERICA DEATH TO ISRAEL ALLAH IS GREAT



Tom BROKAW NBC TV 30 Rockefeller Plaza New York NY 10112

CIS 529: Bioinformatics

PIEAS Biomedical Informatics Research Lab

10512+0002

Zuckerkandl & Pauling (1962)

 Zuckerkandl & Pauling (1962) demonstrated that molecular sequences provide large amounts of information and help understand evolutionary relationships



Sequence Data

	Characters					
Species	1	2	3	4	5	6
α	А	G	А	С	G	G
eta	С	G	Т	G	А	G
γ	А	С	А	G	А	G
δ	А	С	А	С	G	А
ϵ	С	G	Т	С	G	G

Nowadays, biologists rely on molecular sequence data, in particular DNA or RNA sequences, which allows the comparison of a broader range of species.

E. coli, yeast, clam shell, human,

Species vs. Gene Trees

- The taxonomic units (nodes of the tree) can represent genes, species or populations.
- A gene tree represents the evolutionary history of a single gene; e.g. the evolution of the globin family.
- A species tree will generally be built using multiple genes.
- We will focus mainly on species trees.

Methods for Tree Reconstruction

Distance based

• Tree is constructed on the basis of a distance measure between species.

Maximum parsimony

• The tree that explains the data with the smallest number of evolutionary events.

Maximum likelihood

- Statistical method of phylogeny reconstruction
- Explicit model for how data set is generated with nucleotide or amino acid substitutions
- Find topology that maximizes the probability of the observed data given the model and the parameter values (estimated from data)

CIS 529: Bioinformatics

Distance Based Tree Reconstruction

- **Definition**: A distance measure d_{ij} between *n* species is said to be a tree-derived distance if there exists a tree T with the species at the leaves such that d_{ij} is the sum of the branch lengths along a path between i and j.
- Tree-derived distances are typically called additive distances.
- In general we don't have available to us the "real" treedistance. Given surrogate information about the treedistance, we would like to obtain the best estimate of the phylogenetic tree.

Tree Reconstruction - Ultrametric Case

If:

Pairwise distances d_{ij} among *n* extant species satisfy the Ultrametric property

Then:

There exists a unique <u>rooted</u> tree (up to trivial relabeling) with these species as leaves, that gives these distances.

Thus the original tree can be recovered exactly.

Proof: See Ewens and Grant – Statistical Methods in Bioinformatics (2nd edition)

CIS 529: Bioinformatics

UPGMA (Unweighted Pair Group Method using Averages)

- Uses a hierarchical clustering method. Initially each species is in its own cluster. The algorithm sequentially combines two clusters and forms a new node in the tree.
- First define distance d_{ij} between two clusters C_i and C_j to be the average distance between pairs of species from each cluster:

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

where $|C_i|$ denotes the number of species in cluster C_i .

• Note that if C_k is the union of C_i and C_j and C_l is another cluster, then

$$d_{kl} = \frac{d_{il}|C_i| + d_{jl}|C_j|}{|C_i| + |C_j|}$$

This is also known as average linkage cluster analysis

Alternatives are "maximum" or "minimum" linkage

CIS 529: Bioinformatics

UPGMA

- Initialization:
 - Assign each species i to its own cluster C_i
 - Define a leaf of T for each species. Place it at height zero.
- Iteration:
 - Find i, j so that d_{ij} is minimal (in case of ties, choose randomly)
 - Define new cluster k by $C_k = C_i$ union C_i . Define d_{kl} for all l by

$$d_{kl} = \frac{d_{il}|C_i| + d_{jl}|C_j|}{|C_i| + |C_j|}.$$

- Define a node k with "daughter" nodes i and j and place it at height $d_{ij}/2$.
- Add k to the current clusters and remove i and j.
- Termination:
 - When only two clusters i, j remain, place the root at height $d_{ij}/2$.

UPGMA Example



CIS 529: Bioinformatics

Other methods

- Neighbor joining
 - constructs unrooted trees





- One of the most commonly used methods for tree reconstruction
- Uses sequence information directly rather than distances that reflect sequence similarity.
- The idea: find a tree that explain the observed sequences with minimal number of substitutions (Occam's razor)



- Assign a score to a tree (with ancestral sequences included) as the total number of substitutions along all branches of the tree.
- The objective: find the tree that minimizes this score
- Minimization is over all choices of inferred ancestral sequences and tree topologies.



Suppose we are given the four sequences:

AAG, AAA, GGA, AGA

- For each possible topology for a rooted tree with 4 leaves, assign sequences to the ancestral nodes such that the total number of changes needed is minimized.
- Minimize this number over all trees

Two step process:

1. Given a tree topology and an assignment of residues to the leaf nodes, find an assignment of residues to internal nodes that minimizes the number of mutations (*note parsimony treats each position independently*)

This is sometimes called the SMALL PARSIMONY PROBLEM. Solvable using dynamic programming

 Search over all possible trees and find the "best" tree, that is, the tree having the minimal parsimony score. This is sometimes called the LARGE PARSIMONY PROBLEM.

Large Parsimony

# Species	# unrooted trees	
4	3	Exhaustive Search
5	15	
6	105	For each unrooted tree
7	945	with n leaves, calculate the
8	10,395	minimum parsimony score
9	1,35,135	minimum score
10	2,027,025	minimum score.
11	34,459,425	Feasible only for small n.
12	654,729,075	
13	13,749,310,575	
14	316,234,143,225	Parsimony score is
15	7,905,853,580,625	independent of root

Large parsimony is NP complete

PIEAS Biomedical Informatics Research Lab

position.

CIS 529: Bioinformatics

Maximum Likelihood phylogeny

Given a tree T with *n* leaves, where sequence x^j is at leaf node j, and a set of edge lengths denoted by t₁, t₂,..., t_{2n-2} (written as t), we want to define

P(**x** | T, t)

 This requires a model of evolution – probabilities of mutation events that change sequences along the edges of the evolutionary tree

Probability Model

Define:

P(x | y,t)

the probability that a given ancestral sequence y will evolve into x along an edge of length t.





For any set of ancestral sequences assigned to internal nodes, we can calculate the probability of observing such a fully specified tree.

Maximum Likelihood Tree

 The tree with topology T and edge lengths t that maximizes P(x | T, t)

- Finding the maximum likelihood tree involves a search over:
 - Tree topologies
 - Assignment to ancestor sequences
 - Edge lengths t

Phylogenetic Trees using MEGA

 MEGA is a commonly used tool for this purpose

Building Phylogenetic Trees from Molecular Data with MEGA

Barry G. Hall

Bellingham Research Institute, Bellingham, Washington *Corresponding author: E-mail: barryghall@gmail.com. Associate editor: Joel Dudley

Abstract

Phylogenetic analysis is sometimes regarded as being an intimidating, complex process that requires expertise and years of experience. In fact, it is a fairly straightforward process that can be learned quickly and applied effectively. This Protocol describes the several steps required to produce a phylogenetic tree from molecular data for novices. In the example illustrated here, the program MEGA is used to implement all those steps, thereby eliminating the need to learn several programs, and to deal with multiple file formats from one step to another (Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S. 2011. MEGA5: molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods. Mol Biol Evol. 28:2731-2739). The first step, klentification of a set of homologous sequences and downloading those sequences, is implemented by MEGA's own browser built on top of the Google Chrome toolkit. For the second step, alignment of those sequences, MEGA offers two different algorithms: ClustalW and MUSCLE. For the third step, construction of a phylogenetic tree from the aligned sequences, MEGA offers many different methods. Here we illustrate the maximum likelihood method, beginning with MEGA's Models feature, which permits selecting the most suitable substitution model. Finally, MEGA provides a powerful and flexible interface for the final step, actually drawing the tree for publication. Here a step-by-step protocol is presented in sufficient detail to allow a novice to start with a sequence of interest and to build a publication-quality tree illustrating the evolution of an appropriate set of homologs of that sequence. MEGA is available for use on PCs and Macs from ww megasoftware.net.

A phylogenetic tree is an estimate of the relationships among taxa (or sequences) and their hypothetical common ancestors (Nei and Kumar 2000: Felsenstein 2004: Hall 2011), Today most phylogenetic trees are built from molecular data: DNA or protein sequences. Originally, the purpose of most molecular phylogenetic trees was to estimate the relationships among the species represented by those sequences, but today the purposes have expanded to include understanding the relationships among the sequences themselves without regard to the host species, inferring the functions of genes that have not been studied experimentally (Hall et al. 2009), and elucidating mechanisms that lead to microbial outbreaks (Hall and Barlow 2006) among many others. Building a phylogenetic tree requires four distinct steps: (Step 1) identify and acquire a set of hom ologous DNA or protein sequences, (Step 2) align those sequences, (Step 3) estimate a tree from the aligned sequences, and (Step 4) present that tree in such a

way as to clearly convey the relevant information to others. Typically you would use your favorite web browser to identify and download the homologous sequences from a national database such as GenBank, then one of several alignment programs to align the sequences, followed by one of many possible phylogenetic programs to estimate the tree, and finally, a program to draw the tree for exploration and publication. Each program would have its own interface and its own required file format, forcing you to interconvert files

as you moved information from one program to another. It is no wonder that phylogenetic analysis is sometimes considered in timidating! MEGA5 (Tamura et al. 2011) is an integrated program that

carries out all four steps in a single environment, with a single user interface eliminating the need for interconverting file formats. At the same time, MEGA5 is sufficiently flexible to permit using other programs for particular steps if that is desired. MEGA5 is, thus, particularly well suited for those who are less familiar with estimating phylogenetic trees.

Step 1: Acquiring the Sequences

Ironically, the first step is the most intellectually demanding, but it often receives the least attention. If not done well, the tree will be invalid or impossible to interpret or both. If done wisely, the remaining steps are easy, essentially mechanical, operations that will result in a robust meaningful tree.

Often, the investigator is interested in a particular gene or protein that has been the subject of investigation and wishes to determine the relationship of that gene or protein to its homologs. The word "homologs" is key here. The most basic assumption of phylogenetic analysis is that all the sequences on a tree are homologous, that is, descended from a common ancestor. Alignment programs will align sequences, homologous or not. All tree-building programs will make a tree from that alignment. However, if the sequences are not actually descended from a common ancestor, the tree will be

D The Author 2013, Published by Oxford University Press on behalf of the Society for Molecular Biology and Biolution. All rights reserved. For permissions, please

Mol. Biol. Evol. 30(5):1229-1235 doi:10.1093/molbev/mst012 Advance Access publication March 13, 2013

CIS 529: Bioinformatics

1779

Pro

BRIEFINGS IN BIOINFORMATICS. VOL 9. NO 4. 299-306 Advance Access publication April 16, 2008

4-100020-0-06-000

MEGA: A biologist-centric software for evolutionary analysis of DNA and protein sequences

Sudhir Kumar, Masatoshi Nei, Joel Dudley and Koichiro Tamura Schmitted-19th January 2008: Received (in meteod form) - Reh March 2008

Abstract

The Molecular Evolutionary Genetics Analysis (MEGA) software is a desktop application designed for comparative analysis of homologous gene sequences either from multigene families or from different species with a special emphasis on inferring evolutionary relationships and patterns of DNA and protein evolution. In addition to the tools for statistical analysis of data, MEGA provides many convenient facilities for the assembly of sequence data sets from files or web-based repositories, and it includes tools for visual presentation of the results obtained in the form of interactive phylogenetic trees and evolutionary distance matrices. Here we discuss the motivation, design principles and priorities that have shaped the development of MEGA.We also discuss how MEGA might evolve in the future to assist researchers in their growing need to analyze large data set using new computational methods.

Keywords: phylogenetics; genome; evolution; software

genes and estimating neutral and selective evolutionary divergence among sequences (Figure 1). The MEGA software project grew out of our own

need for employing statistical methods in the

Biologist-friendly software tools are crucial in this age in the early 1990s. At this time, most computer of burgeoning sequence databases. These tools not programs available did not allow us to explore the only make it possible to use computational and sta- primary data visually and lacked a user-friendly tistical methods but also allow scientists to select interface. There were two primary general-purpose methods and algorithms best suited to understand the computer programs for inferring phylogenetic trees. function, evolution and adaptation of genes and species. The Molecular Evolutionary Genetics Analysis the other was PHYLIP for inferring phylogenetic (MEGA) software aims to serve both of these purposes trees using various character and statistical methods in inferring evolutionary relationships of homologous such as maximum likelihood, parsimony and distance sequences, exploring basic statistical properties of methods [1, 2]. These programs were (and continue to be) very useful, but the former lacked statistical methods at that time and the latter did not provide a point-and-click user interface [3-6].

In order to make statistical methods available for phylogenetic analysis of DNA and protein sequences phylogenetic analysis in a user-friendly manner,

Corresponding author. Sudhir Kumar, Biodesign Institute, A240, Arizona State University, Tempe, AZ 85287-5301, USA E-mail: s.kumar@asu.edu

Sudhir Kumar is conducting large-scale analysis of genome sequences and spatial patterns of gene expression, and developing statistical methods and bioinformatics tools. He is the Director of the Center for Evolutionary Punctional Genomics at Arizona Sta University, AZ, USA,

Masatoshi Nei is one of the foundes of molecular evolutionary genetics and pustues statistical analysis of molecular and genom evolution. He is the Director of the Institute of Molecular Evolutionary Genetics at Pennylvania State University, PA, USA. Joel Dudley is interested in randational bioinformatics genomic medicine and the application of molecular evolution in translational genomic research. He is a bioinformatics specialist at the Stanford Center for Homedical Informatics Research at Stanford University, USA.

Koichiro Tamura's research interests are in the area of molecular evolution with emphasis on the pattern of gene and genom evolution, and on development of statistical methods and computer programs. He is an associate professor at the Tokyo Metropolitar University, Japan.

© The Author 2008, Published by Oxford University Press, For Permissions, please enabling reals, permissions/flowfording reals, pre-

MEGA steps

- BLAST the given sequence to get more sequences must be done carefully
- Construct MSA using MUSCLE or CLUSTALW
 - MUSCLE is better
- Construct tree
 - Distance based (distance matrix ...)
 - Parsimony based
 - Maximum likelihood based
- Generate confidence values (using bootstrapping)
 - Optional but recommended

For Dengue viruses

 <u>http://faculty.pieas.edu.pk/fayyaz/ippy/html</u> <u>demos/Bioinformatics.html</u>



End of Lecture



http://io9.gizmodo.com/5936427/the-evolutionary-history-of-dragons-illustrated-by-a-scientist

CIS 529: Bioinformatics