Assessing Model Fit

Evidence for Systematic Error (Bias)

Who are the earliest diverging animals?



Pisani et al., PNAS, 2015

Whelan et al., PNAS, 2015

Evidence for Systematic Error (Bias)

Backbone Tree for All Birds



These are enormous datasets, yet they conflict strongly for early divergences.

Thomas 2015.

Evidence for Systematic Error (Bias)

These are two high-profile examples, but there are many others (we'll talk about turtles later).

When conflict is this strong, stochastic error is not a plausible explanation.

Data is no longer limiting. We are now limited by our ability to accurately extract information from the data.

(1) Collect Data (D)



 $\{3.4, 2.1, 5.4, ...\}$

(1) Collect Data (D) (2) Define Models



{**3.4, 2.1, 5.4, ...**}



(1) Collect Data (D) (2) Define Models (3) Fit Models $M_{1} \blacksquare \rightarrow D \qquad M_{1} \blacksquare \rightarrow D : L_{1}$ $\{3.4, 2.1, 5.4, ...\} \qquad M_{2} \supseteq D \qquad M_{2} \supseteq D : L_{2}$ $M_{3} \bigtriangleup \rightarrow D \qquad M_{3} \bigtriangleup \rightarrow D : L_{3}$

(1) Collect Data (D) (2) Define Models (3) Fit Models $M_{1} \blacksquare \rightarrow D \qquad M_{1} \blacksquare \rightarrow D : L_{1}$ $\{3.4, 2.1, 5.4, ...\} \qquad M_{2} \blacksquare D \qquad M_{2} \blacksquare D : L_{2}$ $M_{3} \blacktriangle \rightarrow D \qquad M_{3} \bigstar \rightarrow D : L_{3}$

> (4) Compare Models & Choose "Best"

$$M_2 > M_1 > M_3 \leftarrow \begin{array}{c} AIC \\ BIC \\ BF \\ LRT \end{array} \right\} \begin{array}{c} L_1 \\ L_2 \\ L_3 \end{array}$$

(1) Collect Data (D) (2) Define Models (3) Fit Models

 M_2

 $M_1 \square \rightarrow D$



- {**3.4, 2.1, 5.4, ...**}
- (5) Report Inferences from $M_3 \bigtriangleup \rightarrow D$ "Best" Model



(4) Compare Models & Choose "Best"

M1 |

 M_2

 $\square \rightarrow D$: L₁

 $M_3 \land \longrightarrow D$: L₃

$$M_2 > M_1 > M_3 \leftarrow \begin{array}{c} AIC \\ BIC \\ BF \\ LRT \end{array}$$

(1) Collect Data (D)



{**3.4**, **2.1**, **5.4**, ...}

(5) Report Inference "Best" Mode



The guinea-pig is not a rodent

Anna Maria D'Erchia*†, Carmela Gissi*†, Graziano Pesole‡, Cecilia Saccone*§ & Ulfur Arnason†

nature

* Dipartimento di Biochimica e Biologia Molecolare, Università di Bari, 70125 Bari, Italy

† Department of Evolutionary Molecular Systematics, University of Lund, Sölvegatan 29, S-22362 Lund, Sweden

‡ Dipartimento di Biologia DBAF, Università della Basilicata, 30100 Potenza, Italy § Centro di Studio sui Mitocondri e Metabolismo Energetico.

Consiglio Nazionale delle Ricerche, 70125 Bari, Italy

IN 1991 Graur et al. raised the question of whether the guinea-pig, *Cavia porcellus*, is a rodent¹. They suggested that the guinea-pig and myomorph rodents diverged before the separation between myomorph rodents and a lineage leading to primates and artiodactyls. Several findings have since been reported, both for and against this phylogeny, thereby highlighting the issue of the validity of molecular analysis in mammalian phylogeny. Here we present findings based on the sequence of the complete mitochondrial genome of the guinea-pig, which strongly contradict rodent monophyly. The conclusions are based on the cumulative evidence provided by orthologically inherited genes and the use of three different analytical methods, none of which joins the guinea-pig with myomorph rodents. In addition to the phylogenetic conclusions, we also draw attention to several factors that are important for the validity of phylogenetic analysis based on molecular data.

Models & Best"

 \longrightarrow

И2

 V_3



Fit Models

The Next Step - Assessing Fit

We know that none of our models is really true. Can we be sure that the chosen model captures the salient features of the evolutionary process and provides reliable inferences?

 $M_3 \land \longrightarrow D$

(5) Report Inferences from "Best" Model

т, ∠. і , Ј.Т, ... і

(4) Compare Models & Choose "Best"

 $M_2 > M_1 > M_3 \longleftarrow$

 $M_3 \land \longrightarrow D$: L3

The Next Step - Assessing Fit

We know that none of our models is really true. Can we be sure that the chosen model captures the salient features of the evolutionary process and provides reliable inferences?

(5) Report Inferences from (6) Check Fit of Model to Data "Best" Model



L. **J**. **T**. ...



The Next Step - Assessing Fit

(2) Dafina Madale (3) Fit Models llect Data (D) Are Guinea Pigs Rodents? The Importance of Adequate **Models in Molecular Phylogenetics**

Jack Sullivan^{1,2} and David L. Swofford¹

The monophyly of Rodentia has repeatedly been challenged based on several studies of molecular sequence data. Most recently, D'Erchia et al. (1996) analyzed complete mtDNA sequences of 16 mammals and concluded that rodents are not monophyletic. We have reanalyzed these data using maximum-likelihood methods. We use two methods to test for significance of differences among alternative topologies and show that (1) models that incorporate variation in evolutionary rates across sites fit the data dramatically better than models used in the original analyses, (2) the mtDNA data fail to refute rodent monophyly, and (3) the original interpretation of strong support for nonmonophyly results from systematic error associated with an oversimplified model of sequence evolution. These analyses illustrate the importance of incorporating recent theoretical advances into molecular phylogenetic analyses, especially when results of these analyses conflict with classical hypotheses of relationships.

1 La

KEY WORDS: inconsistency; maximum likelihood; molecular systematics; rodents; rate heterogeneity.

D*₅ DT

3.4, 2.

How might we assess fit?

- (1) Use our **prior knowledge** to ask if the **data** are reasonable.
- (2) Use our **prior knowledge** to ask if **inferences** are reasonable.
 - Above are "gut checks". Very useful, but perhaps subjective. Also difficult to have strong priors for complicated data and models.
- (3) Use your **data** (all or part) to make a **prediction** and see if your prediction matches what you've seen.
 (Posterior Prediction and Cross Validation)

Could have come from $P(--, \theta)$?

Could the model and priors plausibly have given rise to the data?







Previously proposed statistics based on the data:

- Multinomial Likelihood (based on frequencies of site patterns)
- Number of Unique Site Patterns
- Frequency of Invariant Sites
- Heterogeneity of Base Frequencies
- Number of parsimony-inferred "parallel" sites

"We do not like to ask, 'Is our model true or false?", since most probability models in most analyses will not be perfectly true... The more relevant question is, 'Do the model's deficiencies have a noticeable effect on the substantive inferences?" - Gelman, Carlin, Stern, and Rubin Bayesian Data Analysis

What about using the inferences provided by our data as a test statistic(s)?





Tree Space

Tree Space



Tree Space

Tree Space





Branch-Specific Test Statistics



Branch-Specific Test Statistics

(not yet in RevBayes)



Branch-length Test Statistics

Mean Tree Length = 3.15 Variance in Tree Length = 2.30

Marginalizing across topologies

Motivating Results - Simulation



Motivating Results - Simulation





Yeast 343 orthologs 18 taxa Hess & Goldman (2011) Amniotes 1,145 orthologs 10 taxa Crawford et al. (2012)





What might we expect from ideal filtering approaches?

Perfect association between decile membership and tree distance $rho(r_s) = 1$







Posterior Predictive Filtering r_s =0.600, P=0.03656

Rate Filtering *r*_s=0.103, *P*=0.3925

Motivating Results - Barcodes



Posterior predictive P-Value

Barley & Thomson 2016 Mol Ecol

Motivating Results - Barcodes



Active Development!

P³: Phylogenetic Posterior Prediction in RevBayes

Sebastian Höhna,^{*,1,2} Lyndon M. Coghill,³ Genevieve G. Mount,³ Robert C. Thomson,⁴ and Jeremy M. Brown³ ¹Division of Evolutionary Biology, Ludwig-Maximilians-Universität, München, Germany ²Department of Integrative Biology, University of California, Berkeley, CA ³Department of Biological Sciences and Museum of Natural Science, Louisiana State University, Baton Rouge, LA ⁴Department of Biology, University of Hawai'i, Honolulu, HI ***Corresponding author:** E-mail: sebastian.hoehna@gmail.com.

Our current inference-based statistics are computationally intense (lots of MCMC). We are:

• working on faster approximations for inference statistics

 \circ conducting baseline simulation studies to establish power

• making the workflow easier and faster (including HPC)

Thoughts on Interpretation

Annual Review of Ecology, Evolution, and Systematics Evaluating Model Performance in Evolutionary Biology

Jeremy M. Brown¹ and Robert C. Thomson²

• Assessing model fit is probably most useful with big data

- Not meant to be a hypothesis test. We can **always** reject the fit of a model in a strict sense. All models are abstractions.
- Based on the aspects of our model that don't fit well, think about how to structure new models. Remember, with RevBayes you can design your own new models!

Tutorial

- Assessing Phylogenetic Reliability Using RevBayes and P³
 Data and Inference versions
 - Assess adequacy of JC and GTR on example data



 $egin{array}{ccc} M_1 & M_2 \ M_3 & M_4 \ M_5 \end{array}$









				-						
_				-						
_				-						



nature

The guinea-pig is not a rodent

Anna Maria D'Erchia*†, Carmela Gissi*†, Graziano Pesole‡, Cecilia Saccone*§ & Ulfur Arnason†

 Dipartimento di Biochimica e Biologia Molecolare, Università di Bari, 70125 Bari, Italy
 Department of Evolutionary Molecular Systematics, University of Lund,

Sölvegatan 29, S-22362 Lund, Sweden ‡ Dipartimento di Biologia DBAF, Università della Basilicata, 30100

Potenza, Italy
 Society di Chudia qui Mitagendiri e Matalellione Engration

§ Centro di Studio sui Mitocondri e Metabolismo Energetico, Consiglio Nazionale delle Ricerche, 70125 Bari, Italy

IN 1991 Graur et al. raised the question of whether the guinea-pig, *Cavia porcellus*, is a rodent¹. They suggested that the guinea-pig and myomorph rodents diverged before the separation between myomorph rodents and a lineage leading to primates and artiodactyls. Several findings have since been reported, both for and against this phylogeny, thereby highlighting the issue of the validity of molecular analysis in mammalian phylogeny. Here we present findings based on the sequence of the complete mitochondrial genome of the guinea-pig, which strongly contradict rodent monophyly. The conclusions are based on the cumulative evidence provided by orthologically inherited genes and the use of three different analytical methods, none of which joins the guinea-pig with myomorph rodents. In addition to the phylogenetic conclusions, we also draw attention to several factors that are important for the validity of phylogenetic analysis based on molecular data.

 $P(\subseteq, \theta | \mathbf{M}, M_4)$





