

Approximate Bayesian Computation (ABC) in population genetics

Matthieu Foll

04.04.2011

Reviews in Computational Biology

Outline

- ABC: what for?
- The link with population genetics: history of a controversy
- A “user-centered” review of ABC criticisms (including an introduction to ABC)
- Conclusion

ABC: what for?

- Approximate Bayesian Computation
- Monte Carlo method to approximate posterior distributions or likelihood surfaces from a model
- Numerical tool for solving problems within a statistical framework

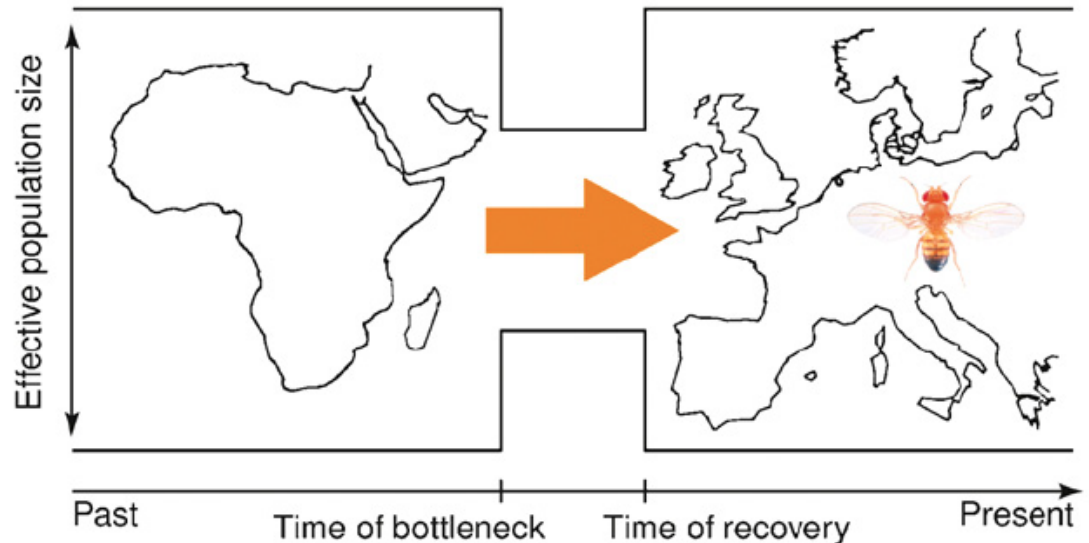
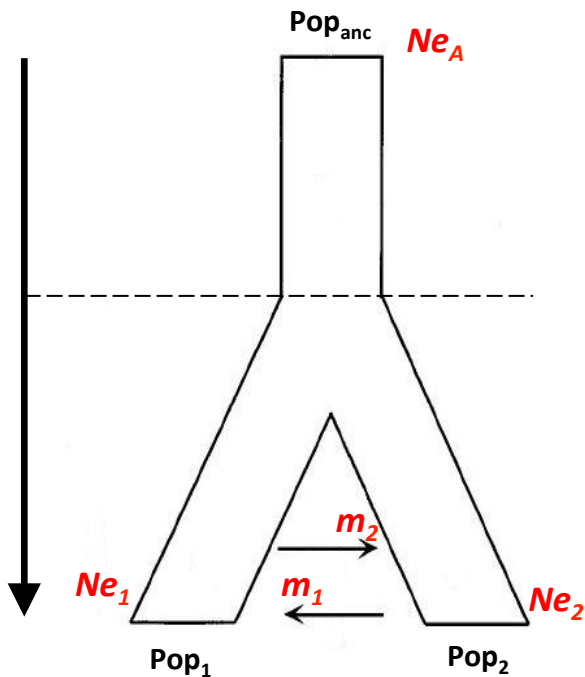
ABC: what for?

- We have a data set D , we suppose that we can represent the problem by a model M , determined by some parameters Φ
- $P(\Phi|D) \propto P(D|\Phi) * P(\Phi)$
Posterior \propto Likelihood $*$ Prior
- Posterior distribution?
- Monte Carlo Markov Chain is the most widely used method for this purpose

In population genetics?

- Infer the demographic history of populations from genetic data

$$\Phi = \{Ne_1, Ne_2, Ne_A, m_1, m_2, t\}$$



History of a controversy

- How it started?

Statistical phylogeography

- 2002: Knowles criticizes Nested Clade Phylogeographical Analysis (NCPA)
Phylogeography: “studying patterns of genetic variation in a geographical context via gene trees”

L. LACEY KNOWLES and WAYNE P. MADDISON

- 2004: answer, and improvements from Templeton

Statistical phylogeography: methods of evaluating and minimizing inference errors

ALAN R. TEMPLETON

- 2007: Strong charge against NCPA

THE AUTOMATION AND EVALUATION OF NESTED CLADE PHYLOGEOGRAPHIC ANALYSIS

The coup de grâce for the nested clade phylogeographic analysis?

RÉMY J. PETIT

When ABC is invited to the party

- 2008: Templeton starts to attack model based approaches as a reply

REPLY

Nested clade analysis: an extensively validated method for strong phylogeographic inference

ALAN R. TEMPLETON

- 2009: Strong charge against ABC

Statistical hypothesis testing in intraspecific phylogeography: nested clade phylogeographical analysis vs. approximate Bayesian computation

ALAN R. TEMPLETON

2008: the war is open

incoherent inference

· six anonymous reviewers

scientific enquiry,

simply untenable

insufficient justification

misunderstanding

non-interpretatable

ad hoc

more valid scientific

high frequency of false positives.

Invalid arguments

confusing

The coup de grâce

method be no longer used

horoscopes

unrealistic

fundamentally flawed

pseudo-statistical

weak inference

irrelevant

wrong

A METHOD THAT FAILS

represent the future

Endless discussions follow...

COMMENTARY

COMMENT

REPLY

LETTER

Reply to

Letters

Letters Response

General comment

- Answering criticisms about one method by pointing out weaknesses of another is not very efficient
 - Sounds childish and personal
 - Create suspicions on your defense
- “My method is working” should be enough
- At the end, it sounds like “maybe my method is not working, but it’s the only one available”

ABC vs. NCPA: Comparing apples and oranges

- Statistics
 - Model based statistics
 - Bayesian statistics
 - Computational methods
 - » EM
 - » IS
 - » MCMC
 - » RJ-MCMC
 - » ABC
 - ...
- Biology
 - Evolutionary biology
 - Molecular phylogenetics
 - Phylogeography
 - » NCPA

James O. Berger^a, Stephen E. Fienberg^b, Adrian E. Raftery^c, and Christian P. Robert^{d,1}

^aDepartment of Statistical Sciences, Duke University, Durham, NC 27708-0251; ^bDepartment of Statistics, Carnegie Mellon University, Pittsburgh, PA 15213-3891; ^cDepartment of Statistics, University of Washington, Seattle, WA 98195-4320; and ^dCentre De Recherche en Mathématiques de la Décision, Université Paris-Dauphine, 75775 Paris Cedex 16, France

Molecular Ecology (2009) 18, 319–331

doi: 10.1111/j.1365-294X.2008.04026.x

Statistical hypothesis testing in intraspecific phylogeography: nested clade phylogeographical analysis vs. approximate Bayesian computation

ALAN R. TEMPLETON

Department of Biology, Washington University, St. Louis, MO 63130-4899, USA

Disentangling ABC criticisms

- ABC is a Monte Carlo method to approximate posterior distributions or likelihood surfaces from a model
- Majority of criticisms are also aimed more generally against
 - Model-based statistics
 - Bayesian statistics

- Some are really ABC specific
 - Invalid
 - Why?
 - Valid
 - Solutions exist
 - Unsolved: perspective?

- Statistics
 - Model based statistics
 - Bayesian statistics
 - Computational methods
 - » EM
 - » IS
 - » MCMC
 - » RJ-MCMC
 - » ABC
 - ...

Model-based/Bayesian criticisms

- Model-based methods do not cover the entire “hypothesis space”
 - compare only a small number of potentially misspecified and subjectively chosen models
 - Back to the debate in the 1930s between Fisher and Neyman-Pearson hypothesis testing
- In ABC “parameter ranges and distributions are only guessed based upon the subjective opinion of the investigators”
 - Classical criticisms of prior in Bayesian statistics

When (real) mathematicians are invited

Coherent and incoherent inference in phylogeography and human evolution

Alan R. Templeton¹

Department of Biology, Washington University, St. Louis, MO 63130

Edited* by Eviatar Nevo, Institute of Evolution, Haifa, Israel, and approved March 3, 2010 (received for review September 16, 2009)

A hypothesis is nested within a more general hypothesis when it is a special case of the more general hypothesis. Composite hypoth-

incoherent inference is formally illogical inference and represents a mathematical error.

LETTER

Incoherent phylogeographic inference

James O. Berger^a, Stephen E. Fienberg^b, Adrian E. Raftery^c, and Christian P. Robert^{d,1}

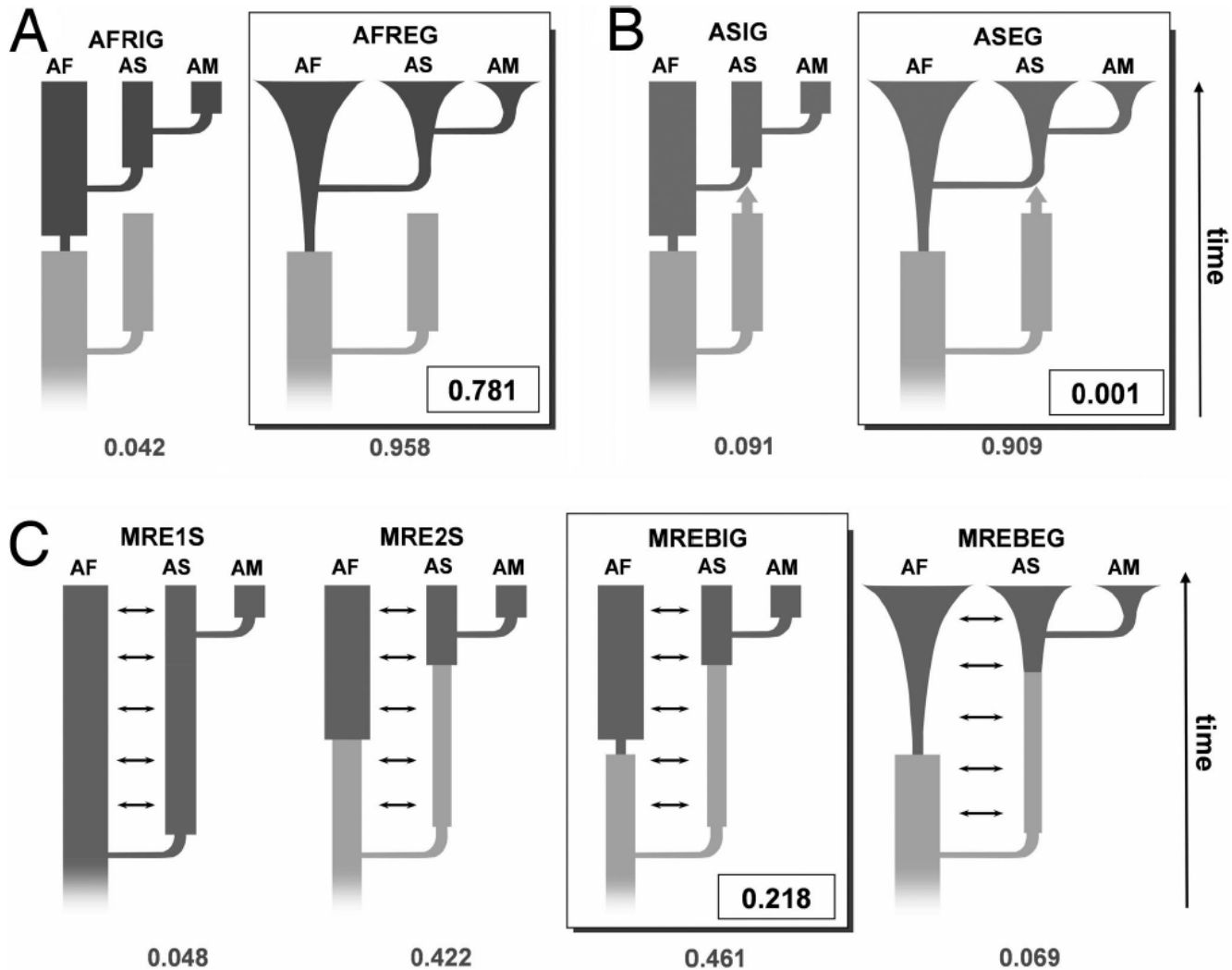
^aDepartment of Statistical Sciences, Duke University, Durham, NC 27708-0251; ^bDepartment of Statistics, Carnegie Mellon University, Pittsburgh, PA 15213-3891; ^cDepartment of Statistics, University of Washington, Seattle, WA 98195-4320; and ^dCentre De Recherche en Mathématiques de la Décision, Université Paris-Dauphine, 75775 Paris Cedex 16, France

“ABC is simply a numerical computational technique; attacking it as incoherent is similar to calling calculus incoherent if it is used to compute the wrong thing”

Invalid Bayesian criticisms

- In “ABC there is no null hypothesis, which complicates the computation of sampling error”
- “The posterior probabilities that emerge from ABC [are] mathematically impossible ... to be probabilities”
- “The probability of the nested special case must be less than or equal to the probability of the general model within which the special case is nested”

Typical « nested hypotheses » example



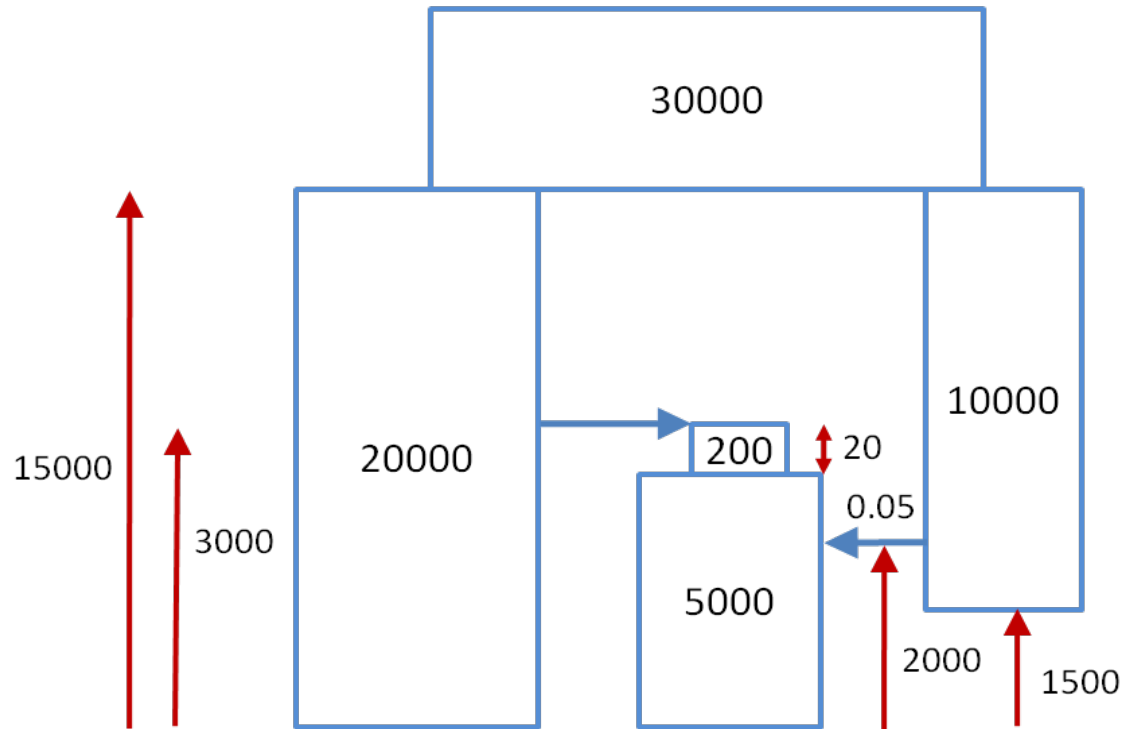
Statistical evaluation of alternative models of human evolution

Nelson J. R. Fagundes^{1*5}, Nicolas Ray⁵, Mark Beaumont¹, Samuel Neuenschwander⁵, Francisco M. Salzano^{2††}, Sandro L. Bonatto^{4,††}, and Laurent Excoffier^{5††}

OK, but how does ABC work?

- We have a data set D
- We suppose that we can represent the problem by a model M
- M is determined by some parameters Φ
- $P(\Phi | D) \propto P(D | \Phi) * P(\Phi)$
Posterior distribution?
- We can not calculate the likelihood but...
- ...we are able to simulate M

Easy to simulate: coalescent theory



Applications note I Genetic and population analysis

fastsimcoal: a continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios

Excoffier L^{1,2,*}, Foll, M^{1,2}

Rejection algorithm

1. Generate Φ at random (=prior)
2. Simulate D' from M with parameters Φ
3. Accept Φ if $D=D'$ and return to 1
4. Stop when sufficient data sets have been accepted

The ABC algorithm

1. Generate Φ at random (=prior)
2. Simulate D' from M with parameters Φ
3. Accept Φ if $d(D, D') < \varepsilon$ and return to 1
4. Stop when sufficient data sets have been accepted

The ABC algorithm with summaries

- Choose **summary statistics S** to represent D and **calculate s** for D .
- Generate Φ at random (=prior)
- Simulate D' from M with parameters Φ
- Calculate s' for D'
- Accept Φ if **$d(s, s') < \epsilon$** and return to 2
- Stop when sufficient data sets have been accepted.

SPECIAL INVITED PAPER

BAYESIANLY JUSTIFIABLE AND RELEVANT FREQUENCY CALCULATIONS FOR THE APPLIED STATISTICIAN¹

BY DONALD B. RUBIN

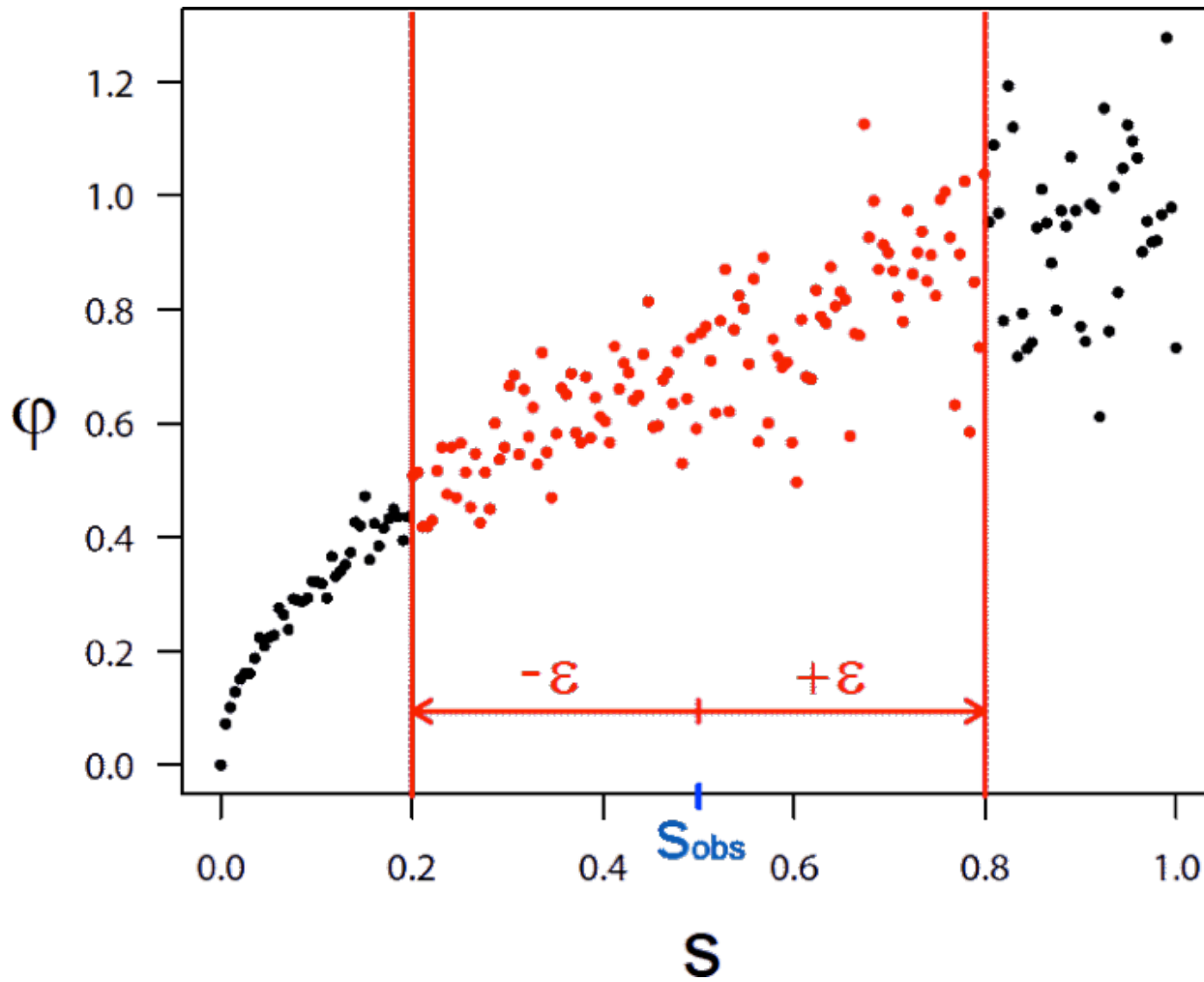
Suppose we first draw equally likely values of θ from $p(\theta)$, and label these $\theta_1, \dots, \theta_s$.

For each θ_j , we now draw an X from $f(X | \theta = \theta_j)$; label these X_1, \dots, X_s . The X_j represent possible values

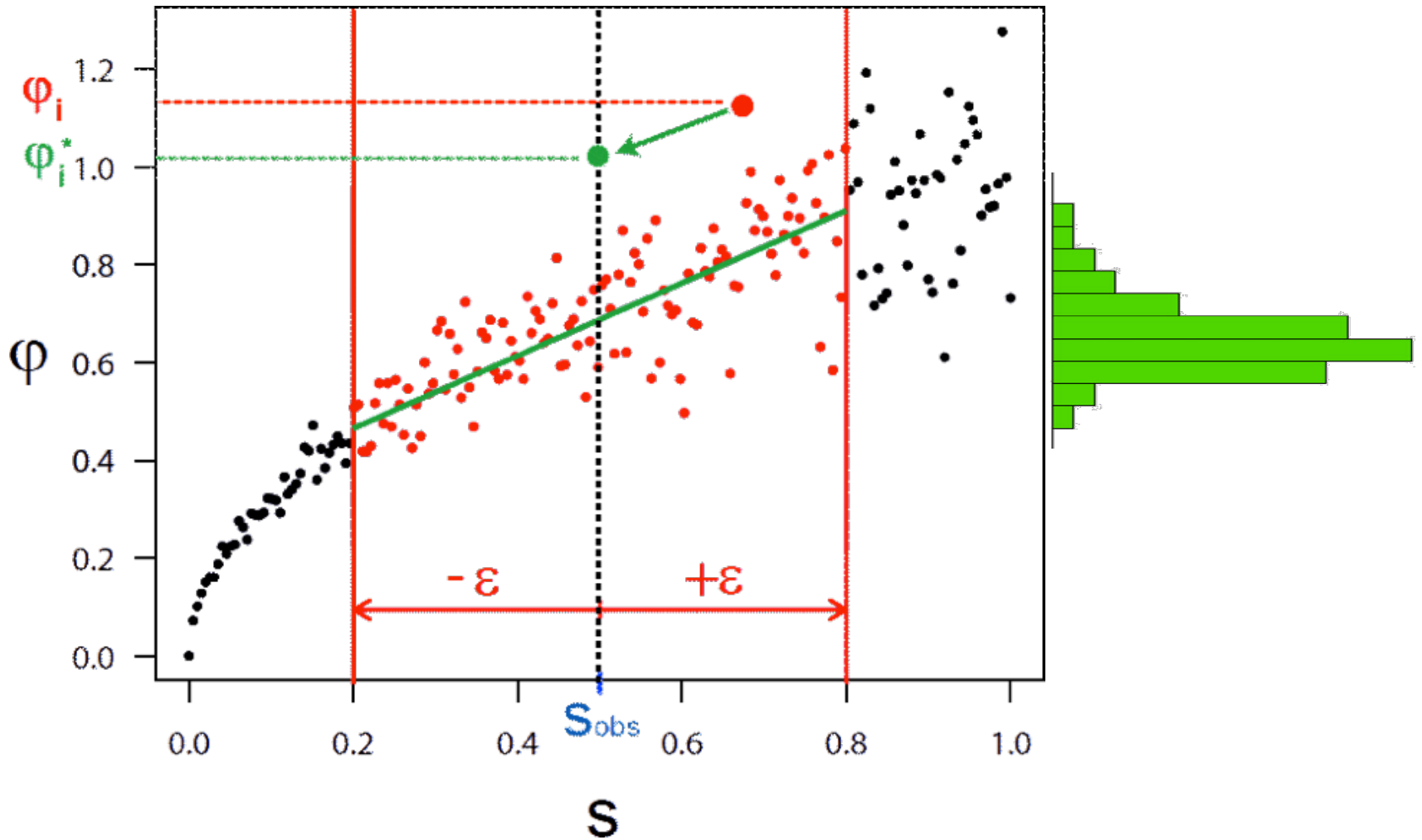
Suppose we collect together all X_j that match the observed X , and then all θ_j that correspond to these X_j .

formally, this collection of θ values represents the posterior distribution of θ .

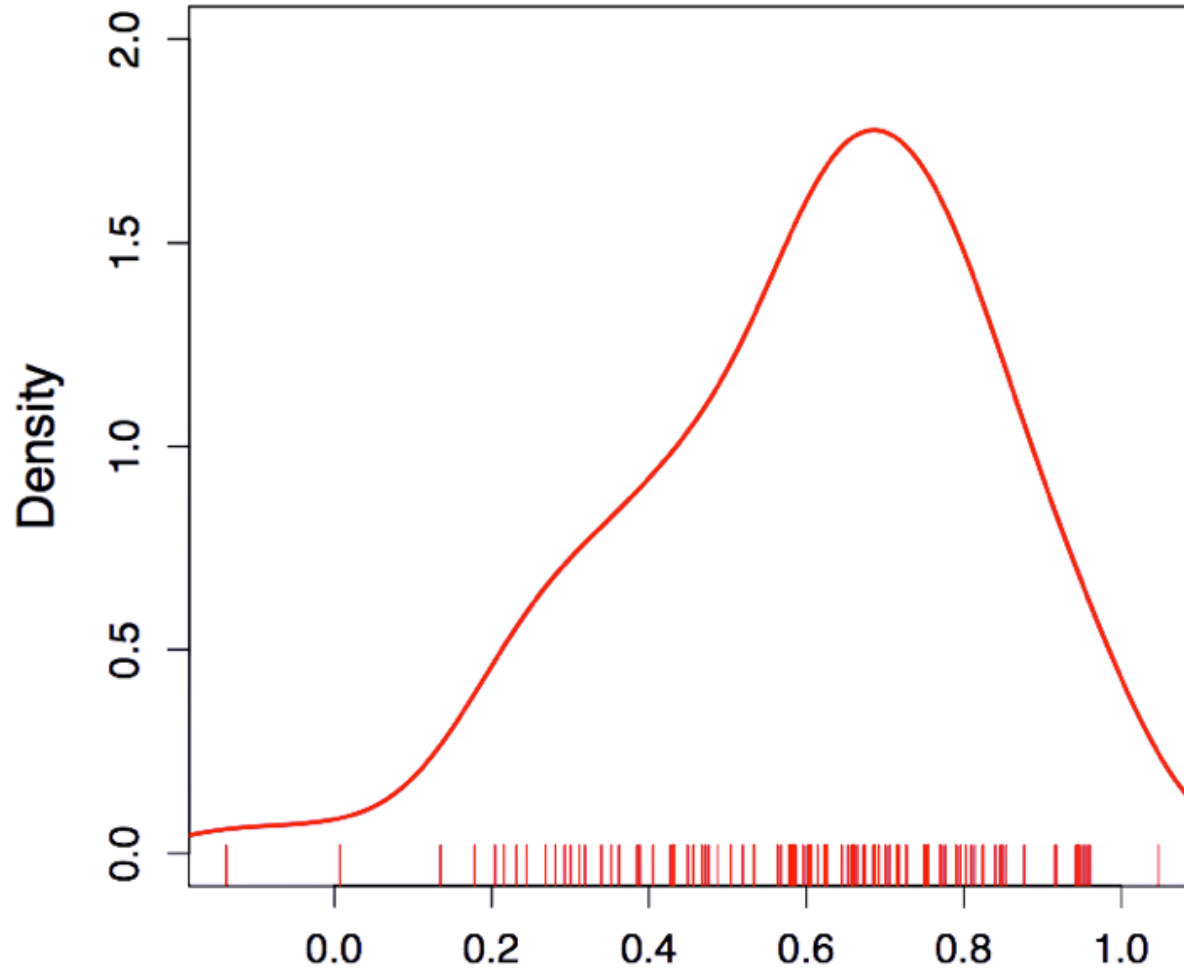
Example



Refinement (regression)



Posterior distribution



Model choice

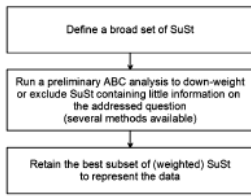
- Generate M and Φ_M at random (=prior)
- Simulate D' from M with parameters Φ_M
- ...
- The posterior probabilities of each model is approximated by the fraction of simulations produced by each of them

Improvements

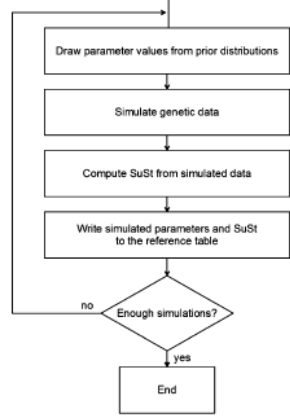
- Once you have found a good region, explore it a bit more !
 - ABC-MCMC
- First do a quick tour of the world, and then go back to your favorite regions to explore a bit more !
 - SMC, ABC-PMC
- Automatic choice of independent summaries that best explain your data
 - PLS
- Non-linear regression

ABC in 9 steps

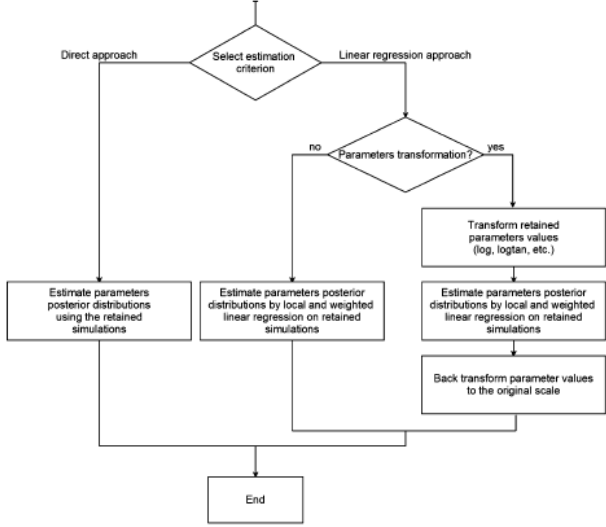
STEP 3 - Choosing the summary statistics (SuSt)



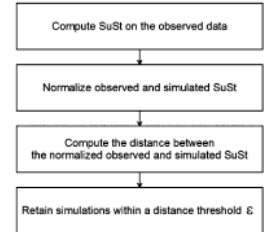
STEP 4 - Simulating the model(s)



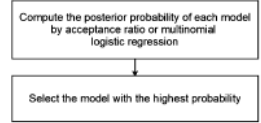
STEP 8 - Parameters estimation



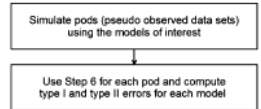
STEP 5 - Filtering the simulations



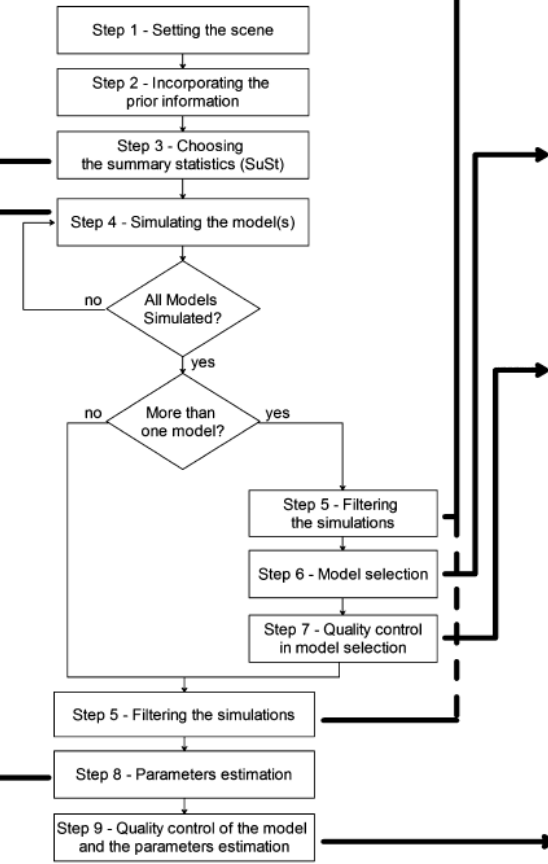
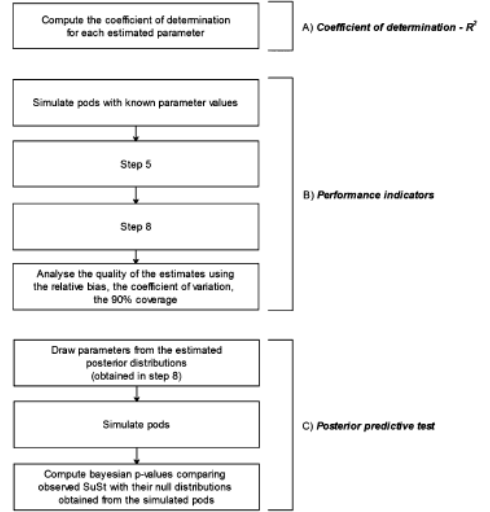
STEP 6 - Model selection



STEP 7 - Quality control in model selection



STEP 9 - Quality control of the model and the parameters estimation



Already reviewed extensively



Annual Review of
Ecology, Evolution,
and Systematics

Volume 41, 2010

Approximate Bayesian Computation in Evolution and Ecology

Mark A. Beaumont

Review

Cell
PRESS

Approximate Bayesian Computation (ABC) in practice

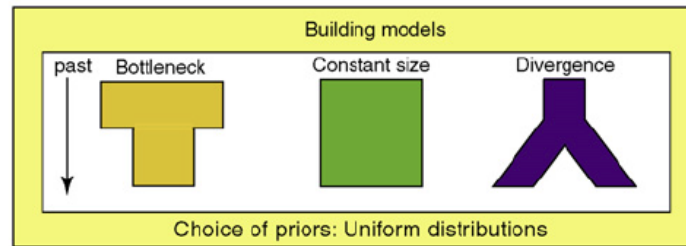
Katalin Csilléry¹, Michael G.B. Blum¹, Oscar E. Gaggiotti² and Olivier François¹

INVITED REVIEW

ABC as a flexible framework to estimate demography
over space and time: some cons, many pros

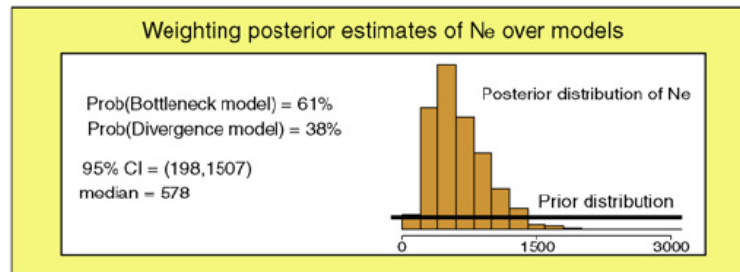
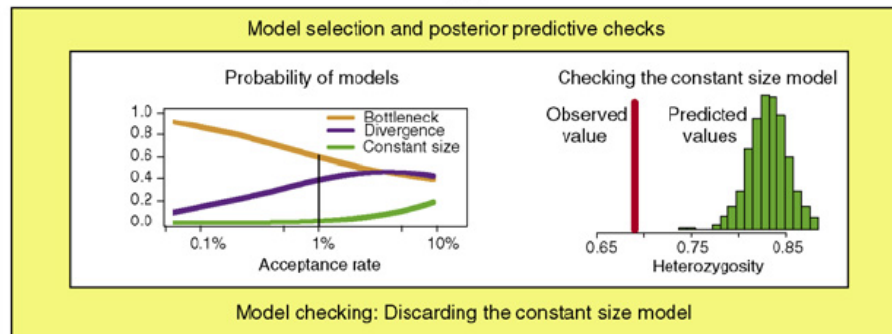
G. BERTORELLE,* A. BENZAZZO* and S. MONA*+‡

Typical population genetics application



Inference with an ABC algorithm

Observed summary statistics
Diversity = 27.38 Garza-Williamson statistic = 0.24



Approximate Bayesian Computation (ABC) in practice

Invalid ABC-specific criticisms

- Simulation weighting
- Posterior densities and Bayesian model choice
- Sample size

- Already reviewed:

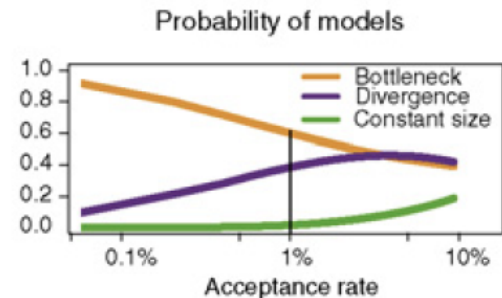
REPLY

In defence of model-based inference in phylogeography

MARK A. BEAUMONT,¹ RASMUS NIELSEN,² CHRISTIAN ROBERT,³ JODY HEY,⁴ OSCAR GAGGIOTTI,⁵ LACEY KNOWLES,⁶ ARNAUD ESTOUP,⁷ MAHESH PANCHAL,⁸ JUKKA CORANDER,⁹ MIKE HICKERSON,¹⁰ SCOTT A. SISSON,¹¹ NELSON FAGUNDES,¹² LOUNÈS CHIKHI,¹³ PETER BEERLI,¹⁴ RENAUD VITALIS,¹⁵ JEAN-MARIE CORNUET,⁷ JOHN HUELSENBECK,² MATTHIEU FOLL,^{16,17} ZIHENG YANG,¹⁸ FRANCOIS ROUSSET,¹⁹ DAVID BALDING²⁰ and LAURENT EXCOFFIER^{16,17}

Valid ABC criticisms

- Summary statistics
 - Choice
 - Sufficiency
- Approximation due to $d(D, D') < \epsilon$: choice of acceptance rate
- (linear) regression
- Model choice: wrong estimate of Bayes Factors (Didelot *et al.* 2011)
- Need to do it properly, in particular quality control step, which is *one* solution to most of these problems



Conclusion

- Important issue in population genetics and phylogeography (2008: >1700 NCPA citations)
- Interesting controversy
- For users, need to separate criticisms
 - ABC specific (valid/invalid)
 - Model-based/Bayesian statistics (valid/invalid)