

**Data Driven Model Selection and
Parameter Estimation Using
Semi-Automatic Approximate Bayesian
Computation to Reconstruct Population
Dynamics From Ancient DNA**

Adam Benjamin Rohrlach

Thesis submitted for the degree of

Master of Philosophy

in

Applied Mathematics

at

The University of Adelaide

(Faculty of Engineering, Computer and Mathematical Sciences)

School of Mathematical Sciences



May 8, 2014

Contents

Abstract	vii
Signed Statement	x
Dedication	xi
Acknowledgements	xii
1 Introduction	1
2 The Coalescent Model and Approximation	6
2.1 Wright-Fisher Reproduction and the Coalescent Approximation . . .	6
2.1.1 The Standard Coalescent Model	8
2.1.2 Modifications to the Standard Coalescent Model	11
2.1.3 Population Structure	16
2.1.4 Using the Coalescent Approximation	21
2.2 Felsenstein’s Maximum Likelihood Methods for Evolutionary Trees	21
3 Population Estimation via the Skyline Plot	29
3.1 Classical Skyline Plots	30
3.1.1 Isochronous Generalised Skyline Plots	31

3.1.2	Heterochronous Generalised Skyline Plots	36
3.1.3	The Bayesian Skyride Plot and Further Modifications	38
3.1.4	Model Selection within the Skyline Plot and the “Known” Tree	43
3.2	Approximate Bayesian Computation	46
3.2.1	The Acceptance-Rejection Algorithm	46
3.2.2	ABC using Sufficient Summary Statistics	48
4	Semi-Automatic Approximate Bayesian Computation	51
4.1	Common Summary Statistics for DNA	52
4.1.1	Single Sample Summary Statistics	52
4.1.2	Multiple Sample Summary Statistics	56
4.2	Approximate Bayesian Computation with Constructed Summary Statistics	59
4.2.1	Step 0: Obtain Observed Data Set	59
4.2.2	Step 1: Constructing Approximately Sufficient Summary Statistics	60
4.2.3	Step 2: ABC Using Constructed Summary Statistics	61
4.3	Results	61
4.3.1	Summary Statistics Used	65
4.3.2	Constant Model Analysis	66
4.3.3	Exponential Model Analysis	72
4.3.4	Migration Model Analysis	79
4.3.5	Parameter Estimation Comparisons	86
5	Data Driven Model Selection	90

5.1	Common Problems and Risks for ABC Inference and Model Comparison	91
5.1.1	Bayes Factors For Model Comparison	93
5.2	Results for Bayes Factors	94
5.2.1	Bayes Factors for the Constant Model Data	96
5.2.2	Bayes Factors for the Exponential Model Data	98
5.2.3	Bayes Factors for the Migration Model Data	99
5.3	ABC for model selection	101
5.3.1	A Fundamental Model Comparison Problem	102
5.3.2	Further ABC Model Comparison Issues	104
5.4	Multinomial Logistic Regression (MLR) for Model Selection.	112
5.4.1	MLR Model Classification Results	114
5.4.2	MLR classification and Bayes Factors	120
5.5	A Data Driven Algorithm for Model Selection and Parameter Estimation	126
6	Bottleneck Data Analysis	131
6.1	Forward Simulation	132
6.2	MLR Classification Step	136
6.3	Parameter Estimation Step	139
6.3.1	Bottleneck ObsDat Analysis Conclusions	156
7	Conclusions	159
7.1	Summary	159
7.2	Future Work	162

A Appendix	163
A.1 Basic Terminology	163
A.2 The Backwards Step Algorithm	164
A.3 Fitted Linear Model for the Constant Model Parameter Estimation	166
A.4 Fitted Linear Models for the Exponential Model Parameter Esti- mation	167
A.5 Fitted Linear Models for the Migration Model Parameter Estimation	169
A.6 Fitted Linear Models for the Bottleneck Model Parameter Estimation	170
A.7 Fitted MLR for the Combined Constant, Exponential and Migration TrainDat	173
A.8 Fitted MLR for the Combined Constant, Exponential, Migration and Bottleneck TrainDat	174
Bibliography	176

Abstract

Population genetics is a discipline within the biological sciences that is concerned with the change in frequency of types of individuals in a population due to natural selection, mutation, genetic drift and gene flow. Genetic drift is the part of this process explained by random sampling. Important to the process of genetic drift is population structure and so we focus on the recovery of population sizes over time, given a set of DNA sequences.

With recent advances in computational power and a growth in the amount of data available, increasingly powerful techniques are being developed for the study of sequence data. Key advances in the early 1980's centred around 'the coalescent', a continuous time approximation to the Wright-Fisher model of reproduction, and these advances resulted in Skyline Plot methods for recovering population size estimates over time. Skyline Plots suffer from large variances for the 'coalescent' event times, and sources of error common to DNA sequence sampling schemes.

Approximate Bayesian Computation (ABC) is a class of likelihood-free methods for statistical inference. ABC techniques can trace their genesis back to the biological sciences due to the complexity of the models for reproduction (and hence the intractability of likelihood calculations). Unfortunately, like Skyline Plots, ABC also suffers from many sources of error, not least of which occurs when we can not use sufficient summary statistics.

To considerably reduce the effect of the error related with the use of insufficient summary statistics, we explore a process of semi-automatic summary statistic cal-

ulation through the use of ‘training data’ (simulated under the coalescent model). We obtain a training set of data, and fit a linear model (under a Box-Cox transformation) for each parameter of interest, using common summary statistics for DNA sequences as predictor variables. We call these linear combinations of (insufficient) summary statistics the semi-automatic summary statistics, and using a new set of simulations, we perform ABC where a simulation is retained if the predicted parameter values are ‘close enough’ to the predicted parameters for the observed data. We analyse three sets of coalescent simulated data from three population models; the Constant, Exponential and Migration Models, and compare our findings with the corresponding Skyline Plot analyses performed in BEAST.

When we simulate data for training our linear model, we must specify a model of population size dynamics, and we explore methods to select a population model, given our data. A common means of model comparison used with ABC analyses is called Bayes Factors. We show that Bayes Factors perform poorly for our data, and highlight a fundamental bias inherent in any model comparison where the probability of a model, given an observed summary statistic, is employed. As an alternative to Bayes Factors, we apply multiple logistic regression (MLR) to classify our observed data into one of a candidate set of possible models. In conjunction with the MLR analysis, we use principal component analysis for visualisation, and introduce a method for attempting to identify when the correct model is not in the candidate model set, or when a classification seems reasonable. We show that this method of classification performs well for the three observed data sets using sensitivity analysis.

Due to the early stage of development of our work, we can not use real world data, and so we use a different type of simulation since our method uses coalescent simulations to train the model. We obtain sequence data simulated under a ‘forward simulation’ framework, a type of sequence simulation that looks forward in time. We define a two-step process for analysis that begins with MLR classification, and

then, under a model chosen by the MLR classification, uses semi-automatic summary statistic calculation for parameter estimation via ABC. We correctly identify this model of population dynamics, and perform parameter estimation on the data, comparing our results with the corresponding BEAST Skyline Plot analysis.

Signed Statement

I, Adam Rohrlach, certify that this work contains no material which has been accepted for the award of any other degree or diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. In addition I certify that no part of this work will, in the future, be used in a submission for any other degree or diploma in any university or other tertiary institution without the prior approval of the University of Adelaide and where applicable, any partner institution responsible for the joint-award of this degree.

I consent to this copy of my thesis, when deposited in the University Library, being available for loan and photocopying, subject to the provisions of the Copyright Act 1968.

I also give permission for the digital version of my thesis to be made available on the web, via the University's digital research repository, the Library catalogue and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

SIGNED: DATE:

Dedication

I would like to dedicate my thesis to Allan and Mary Isobel Troughear.

You were the most wonderful and loving grandparents I could have hoped for. I hope I have made you as proud as you always seemed to be.

Acknowledgements

First, to my supervisors, Professor Nigel Bean and Dr Simon Tuke, thank you for seeing more in me from an early stage, than I saw in myself. It is because of your tireless encouragement, and small research projects during my undergraduate degree, that I decided to pursue post-graduate research. You are not only gifted lecturers, but two of the most patient and straight-faced mentors I could have hoped for. Our meetings have been so enjoyable, that I regularly forget that this is a job, and have begun to believe we might just take this show on the road.

To my Father, thank you for all of the advice these long years, the weekly lunches that allowed me to clear my mind, and the port. It seems you have always been waiting for me to get things right. Thank you for treating me no differently while I did not. To my Mother, thank you for the unconditional love, support, and encouragement. You have always made me feel as though nothing is beyond my reach. To my sisters, Prue and Paige, I thank you for the family meals where I was able to forget my work, and just enjoy being 'at home' again. In fact, to my entire family, I can not thank you enough for continuing to believe in me, even when all the available data suggested you should not.

To my two closest friends, Dan and Brett, thank you for coming on this five year long ride with me. It seems if I were ever in any danger of taking myself too seriously, you knew precisely when to step in. I have known you both for most of my life, and while this reflects poorly on the quality of company you keep, thank you for sticking with me during every phase of it.

To Professor Alan Cooper, and everyone at the Australian Centre for Ancient DNA, this project would not have been the same without you. Thank you for allowing me to sit in on countless meetings and answering every one of my questions. Specifically, to Julien and Oliver, thank you for taking the time to walk me through BEAST, and answering every email I sent you. I could not have done a single comparison of results without your help.

Finally, thank you to everyone in the School of Mathematical Sciences at the University of Adelaide. It has been an honour to work with all of you, and you have made the last two years a genuine pleasure. There are so many people I owe my thanks to, but I would like to thank a couple of you specifically. To David A, Heath, Jess, Kate, Nick, Nic and Stephen, thank you for all of the fun we have had, the help you have given me, and the countless coffees we have drunk over the last few years. I would be remiss not to mention three exceptional people, without whom I could not have completed this thesis with such fond memories.

To David, you have been a pleasure to work with, and I have learned so much about teaching since we ‘teamed up’. Thank you for reading every word I have written, and never once berating me for them. To Mingmei, thank you for not only the last two years, but all of the time we have spent learning together. I doubt I could have succeeded without your influence, perspective, and strength. You are an amazing individual, and you pushed me harder than I would have pushed myself. To Vincent, I thank you for setting the bar so high. It is rare to meet someone so gifted, but still so easy to spend so much time with. Thank you for filling exam study periods with laughter, tea and talking british cakes.