

Copyright is owned by the Author of the thesis. Permission is given for a copy to be downloaded by an individual for the purpose of research and private study only. The thesis may not be reproduced elsewhere without the permission of the Author.

Adjusting the Parameter Estimation of the
Parentage Analysis Software *MasterBayes*
to the Presence of Siblings

A thesis presented in partial fulfillment of the requirements for the
degree of
Master of Applied Statistics
at Massey University, Albany,
New Zealand

Florian Heller
2009

Abstract

Parentage analysis is concerned with the estimation of a sample's pedigree structure, which is often essential knowledge for estimating population parameters of animal species, such as reproductive success. While it is often easy to relate one parent to an offspring simply by observation, the second parent remains frequently unknown. Parentage analysis uses genotypic data to estimate the pedigree, which then allows inferring the desired parameters. There are several software applications available for parentage analysis, one of which is *MasterBayes*, an extension to the statistical software package R. *MasterBayes* makes use of behavioural, phenotypic, spatial and genetic data, providing a Bayesian approach to simultaneously estimate pedigree and population parameters of interest, allowing for a range of covariate models. *MasterBayes* however assumes the sample to be a randomly collected from the population of interest. Often however, collected data will come from nests or otherwise from groups that are likely to contain siblings. If siblings are present, the assumption of a random population sample is not met anymore and as a result, the parameter variance will be underestimated. This thesis presents four methods to adjust *MasterBayes*' parameter estimate to the presence of siblings, all of which are based on the pedigree structure, as estimated by *MasterBayes*. One approach, denoted as DEP, provides a Bayesian estimate, similar to *MasterBayes*' approach, but incorporating the presence of siblings. Three further approaches, denoted as W1, W2 and W3, apply importance sampling to re-weight parameter estimates obtained from *MasterBayes* and DEP. Though fully satisfying adjustment of the estimate's variance is only achieved at nearly perfect pedigree assignment, the presented methods do improve *MasterBayes*' parameter estimation in the presence of siblings considerably, when the pedigree is uncertain. DEP and W3 show to be the most successful adjustment methods, providing comparatively accurate,

though yet underestimated variances for small family sizes. W3 is the superior approach when the pedigree is highly uncertain, whereas DEP becomes superior when about half of all parental assignments are correct. Large family sizes introduce to all approaches a tendency to underestimate the parameter variance, the degree of underestimation depending on the certainty of pedigree. Additionally, the importance sampling schemes provide at large uncertainty of pedigree comparatively good estimates of the parameter's expected values, where the non importance sampling approaches severely fail.

Acknowledgements

I would like to acknowledge a number of people for the support provided during the course of this research project.

First and foremost, my special thanks go to my supervisor, Dr Beatrix Jones. Her deep knowledge, insightful criticism, patience and continued encouragement guided me through the process of this project.

I would further like to thank Dr Howard Edwards for providing me with the helpful course material of the Bayesian Statistics paper.

A big thanks goes also to Tim Napier, who went through the hassle of proof-reading the thesis.

Lastly, I would like to thank my family, my friends and especially my partner Chen Geng for their continued encouragement and support.

TABLE OF CONTENTS

ABSTRACT	i
ACKNOWLEDGEMENTS	iii
TABLE OF CONTENTS	iv
LIST OF ILLUSTRATIONS	vi
1. INTRODUCTION	1
1.1 Introduction to Parentage Analysis	1
1.2 Introduction to Genetics	6
2. METHODOLOGY	11
2.1 Sample Simulation	11
2.2 Introduction to <i>MasterBayes</i>	17
2.2.1 MCMC Parental Assignment Estimation via Gibbs Sampler in <i>MasterBayes</i>	19
2.2.2 <i>MasterBayes</i> Age 2 Parameter Estimate (MB)	22
2.3 Dependent Estimation Approach (DEP)	27
2.4 Analytical Derivation of the True Dependent Parameters from Simulation Data (TRUE)	31
2.5 Importance-Sampling Schemes	32
2.5.1 Importance-Sampling on MB (W1)	34
2.5.2 Importance-Sampling on Draws from Dependent Posterior (W2)	35
2.5.3 Rao-Blackwellized Importance-Sampling on Dependent Posterior (W3)	37
2.6 Overview of Methods	39
3. ANALYSIS AND DISCUSSION	41
3.1 MCMC Settings and Verification	41
3.1.1 Set-Up of Scenarios	41
3.1.2 Burn-In and its Sufficiency	42
3.1.3 Maternal Assignment Success-Rate	43
3.1.4 Thinning Interval	47
3.1.5 Importance Sampling Weights Distribution and Expected Sample	50

3.2 Results	57
3.2.1 Adjusted Expected Value	58
3.2.2 Adjusted Variance	61
4. CONCLUSIONS AND RECOMMENDATIONS	68
5. REFERENCES	70

LIST OF ILLUSTRATIONS

Figure 1.1.1 - Nest Structure	3
Figure 1.2.1 - Likelihood Determination	8
Figure 2.1.1 - Number of Non-Excluded Females, Exclusion Probabilities and Naïve Success-Rates	17
Figure 3.1.1.1 - Examined Combinations and Number of Offspring	41
Figure 3.1.2.1 - Burn-In Success-Rates	43
Figure 3.1.3.1 - Family Size Mean Maternal Assignment Success-Rate	45
Figure 3.1.4.1 - 7 Loci/Family Size 3: Posterior Variance Development over Iterations	49
Figure 3.1.5.1 - 7 Loci/Family Size 4: Standardized Weights Distribution	51
Figure 3.1.5.2 - 7 Loci/Family Size 4: Expected Sample Size Development	52
Figure 3.1.5.3 - 5 Loci/Family Size 3: Standardized Weights Distribution	53
Figure 3.1.5.4 - 5 Loci/Family Size 5: Effective Sample Size Development	54
Figure 3.1.5.5 - Expected Sample Sizes	56
Figure 3.2.1.1 - Parameter Estimate	59
Figure 3.2.2.1 - Parameter Variance of 3 Examined Loci	62
Figure 3.2.2.2 - Parameter Variance	63