

---

# The LASSO Linear Mixed Model for Mapping Quantitative Trait Loci

---

Scott Foster

Doctor of Philosophy  
May 2006

Supervisors: Arūnas Verbyla and Wayne Pitchford

---

THE UNIVERSITY OF ADELAIDE  
School of Agriculture and Wine  
BiometricsSA and Agricultural and Animal Sciences

---



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Overview of Literature . . . . .	1
1.1.1	QTL Mapping . . . . .	2
1.1.2	Statistical Methods . . . . .	3
1.2	Motivating Data . . . . .	4
1.3	Overview of Thesis . . . . .	4
<b>2</b>	<b>Review: Genetics and QTL Mapping</b>	<b>7</b>
2.1	Introduction . . . . .	7
2.2	Basic (Mendelian) Genetics . . . . .	7
2.2.1	Constituents of Phenotype . . . . .	7
2.2.2	Genetic Pathways . . . . .	8
2.2.3	Genetic Structure . . . . .	8
2.2.4	Genetic Effect - Terminology . . . . .	9
2.2.5	Meiosis, Recombination and Genetic Distance . . . . .	9
2.2.6	F <sub>2</sub> and Back-Cross Populations for Outbreeding Organisms . . . . .	12
2.3	QTL Mapping In Outbreeding Organisms . . . . .	13
2.3.1	Missing and Non-Informative Markers . . . . .	14
2.3.2	Parental Haplotype Reconstruction . . . . .	14
2.3.3	Single Marker Models . . . . .	15
2.3.4	Interval Mapping . . . . .	17
2.4	Multiple QTL methods . . . . .	23
2.4.1	Multiple Marker Models . . . . .	23
2.4.2	Composite Interval Mapping . . . . .	25
2.4.3	Multiple Interval Mapping . . . . .	26
2.5	Bayesian Methods . . . . .	27
2.5.1	Single Marker Models . . . . .	27
2.5.2	Incorporating Multiple Markers . . . . .	27
2.5.3	Outbred Line Crosses . . . . .	28
2.5.4	Multiple QTL Models . . . . .	28
2.5.5	Epistatic QTL Models . . . . .	29

2.5.6	Performance of Bayesian QTL Methods . . . . .	29
<b>3</b>	<b>Review: Mixed Models</b>	<b>31</b>
3.1	Introduction . . . . .	31
3.2	Mixed Models . . . . .	31
3.3	Restricted Maximum Likelihood . . . . .	33
3.4	Estimation of Fixed Effects and Variance Parameters . . . . .	34
3.4.1	Fixed Effects . . . . .	34
3.4.2	Variance Parameters . . . . .	35
3.5	Prediction of Random Effects . . . . .	36
3.5.1	Predictions to Minimise Prediction Error . . . . .	37
3.5.2	Mixed Model Equations . . . . .	38
3.5.3	Residuals . . . . .	40
3.6	Computation for the Standard Mixed Model . . . . .	40
3.6.1	Working Variates . . . . .	41
3.6.2	Alternative Forms for Restricted Likelihood and Scores . . . . .	42
3.6.3	The AI Algorithm . . . . .	44
3.7	Analytical Approximations to the Marginal Likelihood . . . . .	44
3.7.1	Laplace's Method . . . . .	45
3.7.2	Partial Laplace's Method . . . . .	45
<b>4</b>	<b>Review: Subset Selection and Biased Estimation</b>	<b>49</b>
4.1	Introduction . . . . .	49
4.2	Why Select Subsets? . . . . .	49
4.2.1	Prediction . . . . .	50
4.2.2	Interpretation of Model . . . . .	52
4.3	Algorithms for Finding Reduced Linear Models . . . . .	52
4.3.1	Forward Selection . . . . .	52
4.3.2	Backward Selection . . . . .	53
4.3.3	Stepwise Selection . . . . .	54
4.3.4	Best Subsets . . . . .	55
4.3.5	One True Model? . . . . .	56
4.4	Ill-Conditioning . . . . .	56
4.5	Biased Estimation . . . . .	57
4.5.1	Constrained Likelihood . . . . .	58
4.5.2	Penalised Likelihood . . . . .	59
4.5.3	Random Effects Models . . . . .	60
4.6	Computation for Ridge Regression . . . . .	61
4.7	Computation for the LASSO . . . . .	62
4.7.1	First Order Conditions . . . . .	62
4.7.2	Linearisation and Optimal Descent Directions . . . . .	63

4.7.3	Interior Point Descent Algorithm . . . . .	64
4.8	Biased Estimate Standard Errors . . . . .	66
4.8.1	Ridge Regression . . . . .	66
4.8.2	The LASSO . . . . .	67
<b>5</b>	<b>Some Results for the LASSO Random Model</b>	<b>69</b>
5.1	Introduction . . . . .	69
5.2	Random Effects Model with Single Effect . . . . .	69
5.2.1	Marginal Distribution of Observations . . . . .	70
5.2.2	Predictive Distribution . . . . .	71
5.3	Vector of Effects . . . . .	81
5.3.1	Marginal Distribution of Observations . . . . .	81
5.4	Reduction of LASSO Estimates . . . . .	82
5.5	Summary . . . . .	84
<b>6</b>	<b>LASSO Computing</b>	<b>85</b>
6.1	Introduction . . . . .	85
6.2	Interior Point Descent Algorithm for Correlated Data . . . . .	86
6.2.1	First Order Conditions . . . . .	86
6.2.2	Linearisation and Optimal Descent Directions . . . . .	87
6.2.3	Interior Point Descent Algorithm . . . . .	87
6.2.4	Variable Updating . . . . .	88
6.3	Estimating Penalised Regression . . . . .	88
6.3.1	Computation Costs for Penalised Regression . . . . .	90
6.4	Discussion of Degrees of Freedom for GCV . . . . .	93
6.4.1	Impact of Estimated Degrees of Freedom on GCV Statistic . . . . .	94
6.5	Summary . . . . .	95
<b>7</b>	<b>The LASSO as a Random Effects Model</b>	<b>97</b>
7.1	Introduction . . . . .	97
7.2	Approximate Likelihood . . . . .	98
7.3	Derivatives of LASSO Estimates . . . . .	99
7.4	Score Equations and Estimates . . . . .	101
7.4.1	Bias in Score Equations and Estimates . . . . .	103
7.5	Adjusted Score Equations . . . . .	104
7.5.1	Theoretical Expected Scores . . . . .	105
7.5.2	Empirical Expected Scores . . . . .	108
7.6	Estimation Algorithm Overview . . . . .	109
7.6.1	Number of Bootstrap Samples . . . . .	110
7.7	Significance Testing of Estimates . . . . .	111
7.8	Simulation Study . . . . .	114

7.8.1	Simulation Design . . . . .	115
7.8.2	Simulation Results . . . . .	115
7.9	Example: Prostate Data . . . . .	118
7.10	Summary . . . . .	121
<b>8</b>	<b>The LASSO Linear Mixed Model</b>	<b>123</b>
8.1	Introduction . . . . .	123
8.2	The LASSO Mixed Model Equations . . . . .	124
8.3	Approximate Likelihoods . . . . .	127
8.3.1	Laplace Approximation . . . . .	128
8.3.2	Partial Laplace Approximation . . . . .	129
8.4	Derivatives of LASSO Estimates . . . . .	132
8.5	Estimation of Dispersion Parameters . . . . .	133
8.5.1	Unadjusted Score Equations . . . . .	133
8.5.2	Observed Information . . . . .	134
8.5.3	Adjusted Score Equations and Estimates . . . . .	135
8.5.4	Solving Adjusted Score Equations . . . . .	137
8.5.5	Algorithm Overview . . . . .	138
8.6	Significance Testing of LASSO Effects . . . . .	138
8.7	Simulation . . . . .	139
8.8	Summary . . . . .	140
<b>9</b>	<b>QTL Detection Using the LLMM</b>	<b>141</b>
9.1	Introduction . . . . .	141
9.2	QTL Mapping Method . . . . .	142
9.3	Simulations . . . . .	143
9.3.1	Heritability 50% . . . . .	144
9.3.2	Heritability 10% . . . . .	145
9.3.3	Simulation Summary . . . . .	145
9.4	Davies' Gene Mapping Data . . . . .	146
9.5	Summary . . . . .	149
<b>10</b>	<b>An Alternative LASSO Random Effects Model</b>	<b>151</b>
10.1	Introduction and Motivation . . . . .	151
10.2	The Alternative Random Effects Model . . . . .	152
10.3	Approximate Marginal Likelihood . . . . .	153
10.3.1	Estimating the Random Dispersion Effects . . . . .	156
10.3.2	An Example of the Approximate Likelihoods' Surfaces . . . . .	157
10.4	Score Equations and Their Components . . . . .	158
10.5	Estimation . . . . .	161
10.6	Performance of Estimation Methods . . . . .	163

10.6.1	Efficacy . . . . .	163
10.6.2	Comparison with Bootstrap Adjustment Method . . . . .	164
10.7	Prediction of LASSO Random Effects . . . . .	165
10.7.1	Best Linear Predictor . . . . .	166
10.7.2	Best Predictor . . . . .	166
10.7.3	Conditional Prediction . . . . .	167
10.7.4	Example of Different Predictors . . . . .	170
10.7.5	Performance of Predictors . . . . .	172
10.8	Inference for Random LASSO Effects . . . . .	173
10.8.1	Example of Probability Calculation . . . . .	174
10.9	Re-Analysis of the Prostate Cancer Data . . . . .	174
10.10	Summary . . . . .	175
<b>11</b>	<b>An Alternative LLMM</b>	<b>179</b>
11.1	Approximate Restricted Likelihood . . . . .	180
11.2	Score Equations and Their Components . . . . .	183
11.3	Computation . . . . .	185
11.3.1	Estimating the Random Dispersion Effects . . . . .	185
11.3.2	Estimation of Dispersion Parameters . . . . .	185
11.4	Prediction and Estimation of Effects . . . . .	186
11.5	Inference for Random Effects . . . . .	187
11.6	Summary . . . . .	187
<b>12</b>	<b>Closing Remarks</b>	<b>189</b>
12.1	Overview . . . . .	189
12.2	Summary of Methodology . . . . .	189
12.3	Possible Future Research . . . . .	191
12.3.1	Reducing Computation . . . . .	191
12.3.2	All-marker Model . . . . .	191
12.3.3	Distributional Assumptions . . . . .	192
12.3.4	Multiple QTL Distributions . . . . .	192
12.3.5	Multiple Traits . . . . .	193
12.3.6	Outlier Detection . . . . .	194
12.3.7	Summary of Future Work . . . . .	195
<b>A</b>	<b>Statistical Results</b>	<b>197</b>
<b>B</b>	<b>Vector and Matrix Algebra Results</b>	<b>203</b>
<b>C</b>	<b>Miscellaneous Results</b>	<b>207</b>

<b>D S-PLUS Functions</b>	<b>209</b>
4.1 Constrained Regression for Dependent Data . . . . .	209
4.2 Penalised Regression for Dependent Data . . . . .	213
4.3 LASSO Mixed Model Equations . . . . .	214
4.4 LASSO Random Effects Model and LLMM . . . . .	216
4.5 Hypothesis Testing . . . . .	224
4.6 Alternative LASSO Random Effects Model . . . . .	226
4.7 Importance Sampling . . . . .	231

# Abstract

This thesis concerns the identification of quantitative trait loci (QTL) for important traits in cattle line crosses. One of these traits is birth weight of calves, which affects both animal production and welfare through correlated effects on parturition and subsequent growth. Birth weight was one of the traits measured in the Davies' Gene Mapping Project. These data form the motivation for the methods presented in this thesis.

Multiple QTL models have been previously proposed and are likely to be superior to single QTL models. The multiple QTL models can be loosely divided into two categories: 1) model building methods that aim to generate good models that contain only a subset of all the potential QTL; and 2) methods that consider all the observed marker explanatory variables. The first set of methods can be misleading if an incorrect model is chosen. The second set of methods does not have this limitation. However, a full fixed effect analysis is generally not possible as the number of marker explanatory variables is typically large with respect to the number of observations. This can be overcome by using constrained estimation methods or by making the marker effects random.

One method of constrained estimation is the least absolute selection and shrinkage operator (LASSO). This method has the appealing ability to produce predictions of effects that are identically zero. The LASSO can also be specified as a random model where the effects follow a double exponential distribution.

In this thesis, the LASSO is investigated from a random effects model perspective. Two methods to approximate the marginal likelihood are presented. The first uses the standard form for the double exponential distribution and requires adjustment of the score equations for unbiased estimation. The second is based on an alternative probability model for the double exponential distribution. It was developed late in the candidature and gives similar dispersion parameter estimates to the first approximation, but does so in a more direct manner.

The alternative LASSO model suggests some novel types of predictors. Methods for a number of different types of predictors are specified and are compared for statistical efficiency.

Initially, inference for the LASSO effects is performed using simulation. Essentially, this treats the random effects as fixed effects and tests the null hypothesis that the effect is zero. In simulation studies, it is shown to be a useful method to identify important effects. However, the effects are random, so such a test is not strictly appropriate. After the specification of the alternative LASSO model, a method for making probability statements about the random effects being above or below zero is developed. This method is based on the predictive distribution of the random effects (posterior in Bayesian terminology).

The random LASSO model is not sufficiently flexible to model most QTL mapping data.



Typically, these data arise from large experiments and require models containing terms for experimental design. For example, the Davies' Gene Mapping experiment requires fixed effects for different sires, a covariate for birthdate within season and random normal effects for management group. To accommodate these sources of variation a mixed model is employed. The marker effects are included into this model as random LASSO effects. Estimation of the dispersion parameters is based on an approximate restricted likelihood (an extension of the first method of estimation for the simple random effects model). Prediction of the random effects is performed using a generalisation of Henderson's mixed model equations.

The performance of the LASSO linear mixed model for QTL identification is assessed via simulation. It performs well against other commonly used methods but it may lack power for lowly heritable traits in small experiments. However, the rate of false positives in such situations is much lower. Also, the LASSO method is more precise in locating the correct marker rather than a marker in its vicinity. Analysis of the Davies' Gene Mapping Data using the methods described in this thesis identified five non-zero marker-within-sire effects (there were 570 such effects). This analysis clearly shows that most of the genome does not affect the trait of interest.

The simulation results and the analysis of the Davies' Gene Mapping Project Data show that the LASSO linear mixed model is a competitive method for QTL identification. It provides a flexible method to model the genetic and experimental effects simultaneously.

# Declaration

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.

I give consent to this copy of my thesis, when deposited in the University Library, being available in all forms of media, now or hereafter known.

**SIGNED:** ..... **DATE:** .....

# Acknowledgements

Firstly and most importantly, I would like to thank my supervisors Ari and Wayne. Not only have they been very generous with their advice, knowledge and patience, but also with their friendship. I have learnt both technical and personal skills from them that I will endeavour to carry through life. It has been a real pleasure.

The seminal idea for this thesis was passed on from Robin Thompson. I am indebted as it has proved interesting, challenging and rewarding.

The data set from the Davies' Gene Mapping Project was generously made available by Wayne and Cindy Bottema. It represents many years of hard work by their team and I appreciate the chance to use it.

I would like to thank the staff and students at BiometricsSA and in the Discipline of Agricultural and Animal Sciences. They have provided support, technical advice and stimulating discussions. Julian Taylor has been especially helpful (and has been very gracious about giving time for 'dumb' computing questions).

My scholarship was provided in part by the University of Adelaide and in part by the Co-operative Research Centre for Beef and Cattle Quality. I am grateful for their support, without which I would not have even started.

My partner Zoë deserves special thanks. Our time in Adelaide has been truly happy.

