



School of Mathematical Sciences

Statistical equivalence in gene expression studies

PhD Thesis

Submitted November 2012

Simon Jonathan Tuke

Signed Statement

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 to Simon Jonathan Tuke 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.

I give consent to this copy of my thesis, when deposited in the University Library, being made 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, the Australasian Digital Theses Program (ADTP) and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Signed: _____

Date: _____

Contents

Signed Statement	1
Abstract	7
Acknowledgements	8
1 Introduction	9
1.1 Description of pluripotency	9
1.2 Description of the stem cell experiment	10
1.3 Pluripotent profile	11
1.4 Gene profiling	13
1.5 Strength of evidence	13

2	Gene profiling	15
2.1	Introduction	15
2.2	Background	16
2.3	The stem cell microarray experiment	18
2.4	Gene profiling methodology	18
2.4.1	Statistical Equivalence	19
2.4.2	Intersection-Union test	22
2.4.3	Gene profiling for pluripotency	24
2.5	Results of applying gene profiling to the stem cell data experiment . .	26
2.6	Sensitivity of gene profiling to changes in the equivalence neighbourhood	29
2.7	Investigation of the proteins discussed in Wang <i>et al.</i>	33
2.7.1	Results from gene profiling using the <i>Dax1</i> profile	34
2.7.2	Results from gene profiling using the <i>Zfp281</i> profile	36
2.7.3	Gene profiling to isolate genes with a profile similar to the gene <i>Sox2</i>	36
2.8	Comparison of gene profiling to Pareto optimisation	39
2.9	Comparison of gene profiling to the Lönnstedt method	43
2.10	Algorithm for parameterisation of the model for gene profiling	46
2.11	Implementation of gene profiling in R	47
2.12	Discussion of further work	47

3	Posterior probabilities for equivalence testing	50
3.1	Equivalence P -values	51
3.1.1	Test statistic	51
3.1.2	Equivalence P -value as a strength of evidence measure	53
3.2	Posterior probabilities of equivalence	54
3.2.1	Estimation of the hyperparameters of the full probability model using the EM algorithm	58
3.2.2	Choice of initial value of parameter estimates in EM algorithm	60
3.2.3	Model checking	64
3.2.4	Calculation of posterior probability for equivalence	66
3.2.5	The posterior probability of equivalence as a strength of evi- dence measure	69
3.2.6	Equivalence posterior probabilities results for stem cell data: day 3 compared to day 0	69
3.3	Robustness to distributional assumptions	73

4	Posterior probabilities for gene profiling	80
4.1	Posterior profile match probability	80
4.1.1	Specification of Bayesian model for K gene profiling parameters	81
4.1.2	Estimation of the hyperparameters via the EM algorithm . . .	82
4.1.3	Initial estimates in the EM algorithm	83
4.1.4	Calculation of posterior profile match probability for multiple (K) gene profiling parameters	84
4.2	Gene profiling with a 2-parameter <i>Oct4</i> model	85
4.2.1	Estimation of the hyperparameters for the 2-parameter <i>Oct4</i> profile	89
4.3	The posterior profile match probability for the 2-parameter <i>Oct4</i> model with the stem cell data	92
4.4	Calculation of the posterior profile match probability using the prod- uct of the marginal posterior probabilities	95
4.5	Comparison of the approximate posterior profile match probability with the posterior profile match probability for the stem cell data with the 2-parameter <i>Oct4</i> profile.	98
4.6	Discussion	100
5	Calculating Q-values from posterior probabilities	103
5.1	Review of q -values	103
5.2	Calculating q -values from P -values	104
5.3	Calculating q -values from equivalence P -values	106
5.4	Calculating q -values from posterior probabilities	107
5.5	Example of q -values for stem cell data	110
5.6	Discussion	111

6 Conclusion	113
Bibliography	126
Appendix	126
A Code	127
A.1 Code for Chapter 2	127
A.2 Gene profiling functions	128
A.3 Pareto optimisation functions	128
A.4 Lönnstedt function	129
A.5 Code for Chapter 3	129
A.6 Code for Chapter 4	129
A.7 Code for Chapter 5	130
B Convergence plots for EM algorithm	131
B.1 The details of the implementation of the EM algorithm for the stem cell case with the 2-parameter <i>Oct4</i> profile	142

Abstract

The goal of this thesis is to develop methods that enable researchers to identify genes with a pre-specified profile in gene expression studies. The motivation for these methods is a gene expression study performed at the University of Adelaide to identify the genes associated with pluripotency in mouse embryonic stem cells.

The method developed, *gene profiling*, utilises intersection-union tests that combine tests both for equivalent and differential expression. The theoretical basis of gene profiling and its application to the stem cell data are discussed in the thesis, as well as a comparison of gene profiling to alternative methods described in the literature for identifying genes corresponding to a pre-specified profile.

In the second part of this thesis, ‘strength of evidence’ measures for equivalence and gene profiling are developed. A P -value for equivalence is shown to have undesirable properties and a posterior probability of equivalence is proposed. The calculation of the posterior probabilities for equivalence using a hierarchical model is described and extended to gene profiling.

A method to calculate q -values from the posterior probabilities is described. The resultant q -values are compared to the posterior probabilities.

Acknowledgements

First I would like to thank Professor Patty Solomon and Associate Professor Gary Glonek for their patience and support as my PhD supervisors. Without you these past few years would not have been the same. You have shown me how to research and then communicate our ideas in clear and concise prose. I can only hope that in our future work I can repay the debt.

To the friends I have made at the School of Mathematical Sciences, thank you for making the journey so enjoyable. To name a few: David Butler, David Roberts, Ray Vozzo, Sam Cohen, Matt Finn, Matt Roughan, Trent Mattner, and many more. I hope to continue drinking coffee and chewing the fat with you.

To Eric Parsonage, thanks. It is nice to have a friend who in times of need is willing to tell you that you are an idiot and should ‘suck it up, fat boy’. There would probably not have been a PhD without coffees with you. I know that I will always be barely an acquaintance, but you will always be my friend.

Finally, thanks to Alison: I dedicate this PhD to you. Sorry it took so long, but I got distracted. Is that not a great motto for life?

Note on Original Work

This thesis presents original results and proofs, but it also reviews existing results by other authors. Results and proofs due to other authors will be clearly marked with a citation such as (Storey, 2003). Any result not marked in this way is an original result the proof or exposition of which is due to the author of this thesis.

The key results on gene profiling in Chapter 2 have been published in the journal *Biostatistics*; see Tuke, Glonek, and Solomon (2009).