

Application of Smooth Tests of Goodness of Fit to Generalized Linear Models

Paul Rippon

B Eng (Chem, Hons1)

Grad Dip Comp Sci

Grad Dip Math Stud

Grad Cert Prac Ter Teach

Thesis submitted to satisfy requirements for the degree of
Doctor of Philosophy in Statistics.

October, 2012.

The thesis 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. 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.

Paul Rippon

Acknowledgements

I would like to acknowledge and thank my supervisors Professor John Rayner and Dr Frank Tuyl for their support and guidance in completing this research. John, in particular, agreed to take over as my supervisor several years ago and provided the possibility of a new research area, interesting ideas and a path forward when I had lost my way in the dark PhD forest. In addition to his knowledge and insight, his unfailingly positive approach and enthusiasm has been of immense value over the long and often painful process of completing this thesis. More recently, Frank has provided a different perspective in helping to critically review the writing and suggest improvements in how certain ideas were presented.

I would also to acknowledge the patience of my examiners, whoever they may be, who have been waiting for some time now to receive this thesis.

Finally, to a long list of family and friends, you are now permitted to ask that dreaded question: “How’s the thesis?”

Contents

| | |
|--|-----------|
| Abstract | 1 |
| 1 Introduction | 3 |
| 1.1 Statistical Models | 3 |
| 1.2 Guide to Examiners | 10 |
| 2 Likelihood | 15 |
| 2.1 The Likelihood Function | 15 |
| 2.2 Likelihood Ratio, Wald and Score Tests | 19 |
| 2.3 Smooth Tests of Goodness of Fit | 29 |
| 3 Smooth Testing for Poisson Regression | 37 |
| 3.1 Poisson Regression | 38 |
| 3.2 A Smooth Test of the Distributional Assumption | 39 |
| 3.3 Example: Bladder Cancer | 51 |
| 4 Smooth Testing for GLMs | 55 |
| 4.1 Generalized Linear Models | 55 |
| 4.2 Derivation of the Smooth Test | 58 |
| 4.3 Applying the Smooth Test | 71 |

| | | |
|----------|---|------------|
| 5 | Smooth Testing: Count Response | 81 |
| 5.1 | Size Study | 83 |
| 5.2 | Using Bootstrap p-values | 92 |
| 5.3 | A Smooth Test for the Poisson Distribution | 97 |
| 5.4 | Goodness of Fit for Poisson Regression Models | 100 |
| 5.5 | Power Study for Poisson Regression | 104 |
| 5.6 | Poisson Regression Examples | 118 |
| 5.7 | Negative Binomial Regression Examples | 129 |
| 5.8 | Chapter Summary | 135 |
| 6 | Smooth Testing: Binary Response | 137 |
| 6.1 | Binomial Regression | 137 |
| 6.2 | Smooth Testing for Binomial Regression Models | 142 |
| 6.3 | A Smooth Test for the Binomial Distribution | 143 |
| 6.4 | Binomial Regression Applications | 145 |
| 6.5 | Power Study for Logistic Regression | 159 |
| 6.6 | Chapter Summary | 163 |
| 7 | Smooth Testing: Continuous Response | 165 |
| 7.1 | Unknown Dispersion Parameter | 165 |
| 7.2 | Adapting the Smooth Test | 166 |
| 7.3 | Smooth Testing for Normal Response Models | 169 |
| 7.4 | A Smooth Test for the Normal Distribution | 174 |
| 7.5 | Normal Response Examples | 175 |
| 7.6 | Smooth Testing for Gamma Response Models | 188 |

| | | |
|----------|--|------------|
| 7.7 | A Smooth Test for the Gamma Distribution | 196 |
| 7.8 | Gamma Response Examples | 200 |
| 7.9 | Chapter Summary | 205 |
| 8 | Conclusion and Further Research | 207 |
| 8.1 | Conclusion | 207 |
| 8.2 | Further Research | 208 |
| A | Fitting and Assessing GLMs | 217 |
| A.1 | Exponential Family of Distributions | 217 |
| A.2 | Fitting a Generalized Linear Model | 219 |
| A.3 | Assessing GLMs | 223 |
| A.4 | Poisson Regression Example | 228 |
| B | Simulation of Probability Distributions | 233 |
| B.1 | Bootstrap p-values | 234 |
| B.2 | Estimating Probabilities | 245 |
| B.3 | Estimating Quantiles | 248 |
| B.4 | Estimating Power | 254 |
| C | Moments and Cumulants | 259 |
| C.1 | Moments about the Origin | 259 |
| C.2 | Moments about the Mean | 260 |
| C.3 | Factorial Moments | 263 |
| C.4 | Cumulants | 266 |

| | |
|---|------------|
| D Orthogonal Polynomials | 271 |
| D.1 Orthogonal polynomials | 272 |
| D.2 Computing Sequences of Orthogonal Polynomials | 275 |
| D.3 Examples | 295 |
| E Developed Software | 309 |
| E.1 The SmoothGLM R package | 310 |
| E.2 The Rippon R package | 316 |
| F Useful Results | 317 |
| F.1 Matrix Algebra | 317 |
| F.2 Proof of Weightings Condition | 325 |
| F.3 The Gamma and Related Functions | 326 |
| G Generalized Score Test: Behrens-Fisher | 329 |
| G.1 Introduction | 329 |
| G.2 Normal Populations with Equal Variances | 330 |
| G.3 Normal Populations with Different Variances | 335 |
| G.4 A Generalized Score Test | 341 |
| G.5 Power Study | 348 |
| G.6 Size Study | 354 |
| G.7 Conclusion | 360 |
| H Quasi-Likelihood | 361 |
| H.1 Estimating Functions | 361 |
| Bibliography | 370 |

Abstract

Statistical models are an essential part of data analysis across many diverse fields. They are used to test research hypotheses, aid decision making, estimate effect sizes and/or improve understanding of the underlying processes generating the data of interest. However it is essential to critically assess any fitted model, confirming that the model really is compatible with the data, before meaningful conclusions are possible.

Generalized linear models (GLMs) provide a flexible modelling framework encompassing many commonly used models including the normal linear model, logistic regression model and Poisson regression model. This thesis explores how **the smooth testing concept** – originally proposed by [Neyman \(1937\)](#) and further developed by [Rayner et al. \(2009\)](#) among others – can be **used to test the distributional assumption in a GLM**. However sensible interpretation of this test, or any other test used to assess the fit of a GLM, must recognize that:

- the stochastic, deterministic and link components that make up a GLM should all be considered when assessing model validity,
- the validity of any one of these three components cannot be sensibly considered in isolation as it is confounded by the validity of the other two.

It is therefore important to consider how the smooth test developed in this thesis might be usefully incorporated into an overall model development strategy for GLMs, either replacing or supplementing existing diagnostic tools. Simulation studies demonstrate that the **power of the smooth test** is competitive with other existing tests. However, it also offers the possibility of **improved diagnostic ability**

through the breakdown of the overall smooth test statistic into a sum of squares of interpretable components. The **SmoothGLM** package has been developed which **implements the smooth test** in a form that can be easily applied to models fitted using the standard `glm()` function **within the R statistical computing environment**.