Made available by Hasselt University Library in https://documentserver.uhasselt.be

Double generalized linear model for tissue culture proportion data: a Bayesian perspective Peer-reviewed author version

CORREA VIEIRA, Afranio Marcio; Leandro, Roseli A.; DEMETRIO, Clarice & MOLENBERGHS, Geert (2011) Double generalized linear model for tissue culture proportion data: a Bayesian perspective. In: JOURNAL OF APPLIED STATISTICS, 38 (8), p. 1717-1731.

DOI: 10.1080/02664763.2010.529875 Handle: http://hdl.handle.net/1942/14380 See discussions, stats, and author profiles for this publication at: https://www.researchgate.net/publication/227618209

Double generalized linear model for tissue culture proportion data: a Bayesian perspective

Article in Journal of Applied Statistics · August 2011

DOI: 10.1080/02664763.2010.529875 · Source: RePEc

CITATIONS		READS		
2		43		
4				
4 autno	rs:			
	Afrânio M C Vieira		Roseli Aparecida Leandro	
	Iniversidade Federal de São Carlos	\sim	University of São Paulo	
	36 PUBLICATIONS 277 CITATIONS		23 PUBLICATIONS 53 CITATIONS	
	SEE PROFILE		SEE PROFILE	
	Clarice G. B. Demétrio		Geert Molenberghs	
	University of São Paulo	1250	Universiteit Hasselt and University of Leuven	
	153 PUBLICATIONS 1,958 CITATIONS		892 PUBLICATIONS 20,546 CITATIONS	
	SEE PROFILE		SEE PROFILE	

Some of the authors of this publication are also working on these related projects:



Implementation of pattern-mixture models View project

Project

Flexible regression models for count data View project

All content following this page was uploaded by Afrânio M C Vieira on 09 August 2016.

Double Generalized Linear Model for Tissue Culture Proportion Data: A Bayesian Perspective¹

Afrânio M. C. Vieira²

Departamento de Estatística, Universidade de Brasília ICC centro, subsolo, módulo 15, CEP 70910-900 Brasília, DF, Brazil

Roseli A. Leandro and Clarice G. B. Demétrio

Departamento de Ciências Exatas, Universidade de São Paulo - ESALQ Av. Pádua Dias 11, CP 9, CEP 13418-900 Piracicaba, SP, Brazil

Geert Molenberghs

I-BioStat, Universiteit Hasselt and Katholieke Universiteit Leuven Agoralaan 1, 3590 Diepenbeek, Belgium

¹Keywords: Bayesian data analysis; Generalized Linear Models; Tissue Culture; Markov Chain Monte Carlo; Binomial Distribution; Gibbs sampling; Random Effects

²Corresponding author; Email: afranio@unb.br

Abstract

Joint generalized linear models (JGLM) and double generalized linear models (DGLM) were 2 designed to model outcomes for which the variability can be explained using factors and/or 3 covariates. When such factors operate, the usual normal regression models, which inherently 4 exhibit constant variance, will under-represent variation in the data and hence may lead to 5 erroneous inferences. For count and proportion data, such noise factors can generate a so-6 called overdispersion effect, and the use of binomial and Poisson models underestimates the 7 variability and, consequently, incorrectly indicate significant effects. In this manuscript, we 8 propose a double generalized linear model from a Bayesian perspective, focusing on the case 9 of proportion data, where the overdispersion can be modeled using a random effect that 10 depends on some noise factors. The posterior joint density function was sampled using Monte 11 Carlo Markov Chain (MCMC) algorithms, allowing inferences over the model parameters. An 12 application to a dataset on apple tissue culture is presented, for which it is shown that the 13 Bayesian approach is quite feasible, even when limited prior information is available, thereby 14 generating valuable insight for the researcher about its experimental results. 15

$_{16}$ 1 Introduction

Many well-known experimental designs that are applied across a diverse range of scientific 17 domains are based on the assumption of variance homogeneity. It is a quite strong assump-18 tion when one is faced with situations where environmental or external factors influence the 19 experimental measures. Modeling the variability from planned experiments gained momentum 20 with Taguchi's work (Taguchi, 1985), which emphasizes the need to adequately deal with the 21 influence of noise and control factors in industrial experimentation, as a means to reducing 22 loss of information and hence optimizing product quality. In a conventional approach, if ei-23 ther environmental factors, the process factors under investigation, or a combination thereof, 24 influences the variance of the continuous response variable, then it means that all statistical 25 inferences from the resulting model will be based on a single dispersion measure, likely inflated 26 by the effects not entered into the model. For proportion or count data, the effect of not 27 taking into account such overdispersion is to produce underestimated variances if the stan-28 dard, too restrictive, models, such as binomial or Poisson-based models are used. Needless to 29 say that ultimately inference is in jeopardy then. Related to this, not taking account of this 30 phenomenon can lead to the selection of overly complex models (Hinde & Demétrio, 1998). 31

The approach of Taguchi to deal with dispersion effects was criticized and a discussion 32 started about effectiveness and alternatives to the signal-to-noise ratios (Box, 1988). One 33 argument against signal-to-noise regards the fact that a transformation is chosen a priori. An 34 alternative presents itself by way of the Box and Cox transformation family (Box & Cox, 1964), 35 where the choice of the best variance-stabilizing transformation is data driven. However, the 36 alternatives proposed to quantify and graph dispersion effects takes the form of exploratory 37 tools; a joint approach was not considered. At the same time, a modeling approach was 38 undertaken. 39

The concept of modeling heterogeneity through a pair of parametric non-linear predictors was formally established by Harvey (1976), with the parameters linked to the mean and variance estimated by maximum likelihood, for a normally distributed response variable.

For this case, when all factors are quantitative, alternatives exists in the form of so-called dual response surface methodology, where the mean and dispersion models are optimized simultaneously (Myers et al., 1992).

This problem was revisited later, and various regression models have been proposed to 46 jointly model mean and dispersion (Aitkin, 1987; Wolfinger & Tobias, 1998; Smyth, 1989; 47 Nelder & Lee, 1991). These authors base inferences, including hypothesis testing and interval 48 estimation, on asymptotic theory (McCullagh & Nelder, 1989). Such methods work well 49 with large sample sizes combined with modest numbers of model parameters. However, in 50 agricultural research, many experiments exhibit a large number of parameters relative to the 51 sample size. The asymptotic-theory-based estimators and their corresponding measures of 52 uncertainty can then be questionable and lead to erroneous conclusions. This motivates our 53 choice for a relative small set of data. 54

In this paper, we propose a double generalized linear model (DGLM) for proportion data 55 using a Bayesian framework for parameter estimation. This approach allows one to incorporate 56 the uncertainty about the unknown quantities of the model using prior information into the 57 estimation procedure. The difficulty of obtaining the parameters' posterior marginal densities is 58 overcome by the use of Monte Carlo Markov Chain (MCMC) algorithms (Gamerman & Lopes, 59 2006). The rest of the paper is organized as follows. In Section 2, the generalized linear models 60 (GLM) framework, the extended quasi-likelihood estimation method, and the model proposed 61 by Smyth (1989) and Nelder & Lee (1991) are briefly described and commented. A Bayesian 62 perspective on the DGLM is presented in Section 3. In Section 4, an apple tissue culture 63 experiment described in Ridout & Demétrio (1992) is introduced, with the results presented 64

and discussed in Section 4.2.

⁶⁶ 2 Joint Modeling of Mean and Dispersion

The methodology proposed by Smyth (1989) and Nelder & Lee (1991) for the joint modeling of mean and dispersion involves two generalized linear models (GLM; McCullagh & Nelder, 1989). For a random sample of n observations (y_i, \mathbf{x}_i) , where y_i , i = 1, ..., n is an observed value for a single response variable Y_i , and $\mathbf{x}_i = (x_{1i}, x_{2i}, ..., x_{pi})^T$ is a p-dimensional vector of explanatory variables, the three components of a GLM are (Hinde & Demétrio, 1998): (i) *independent random variables* Y_i , stemming from the exponential family of distributions with mean μ_i and constant scale parameter ϕ , i.e., observations from a density of the form:

$$f(y|\theta,\phi) = \exp\left\{\frac{y\theta - b(\theta)}{a(\phi)} + c(y,\phi)
ight\},$$

where $a(\phi) = \phi/w$, ϕ is the dispersion parameter, w is a prior weight, θ is the canonical parameter [it can be shown that $E(Y) = b'(\theta)$ and $Var(Y) = \phi b''(\theta)$]; (ii) a linear predictor vector η given by $\eta = \mathbf{X}\beta$, where β is a vector of p unknown parameters and $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T$ is an $n \times p$ design matrix; (iii) a link function $g(\cdot)$ relating the mean to the linear predictor, i.e., $g(\mu_i) = \eta_i = \mathbf{x}_i^T\beta$; hence, $E(Y_i) = g^{-1}(\eta_i)$.

In this paper, we focus on the particular GLM with binomial distribution and logit link function. Assuming that a random variable Y_i , the number of successes out of m_i samples, has a binomial distribution with probability of success π_i , it follows that $\theta_i = \ln [\mu_i/(m_i - \mu_i)]$, $b(\theta_i) = m_i \ln(1 + e^{\theta_i})$ and $\phi = 1$. Therefore, $E(Y_i) = m_i \pi_i = \mu_i$, $Var(Y_i) = m_i \pi_i (1 - \pi_i)$ and $g(\mu_i) = \ln [\mu_i/(1 - \mu_i) = \eta_i]$. Parameter estimation conventionally proceeds by maximum likelihood; in computational terms, the *iteratively re-weighted least square algorithm* (IRLS) is popular.

Note that, because the dispersion parameter $\phi = 1$, the variance function depends solely on the mean parameter. However, it is quite common in experimental situations that proportions show variability larger than that allowed by the theoretical variance of the binomial distribution, the aforementioned *overdispersion*. Hinde & Demétrio (1998a) reviewed a wide variety of avenues for overdispersion modeling, together with methods of estimation. These authors also discussed applications to agricultural experimentation data. Nelder & Pregibon (1987) proposed the *extended quasi-likelihood* (EQL) method for parameter estimation, based only on the first two moments of a distribution. The EQL method consists of maximizing the function

$$Q^{+} = -\frac{1}{2} \sum_{i=1}^{n} \{ \frac{d(y_{i}, \mu_{i})}{\phi_{i}} + \log(2\pi\phi_{i}V(y_{i})) \},\$$

where

$$d(y,\mu) = -2\int_y^\mu \frac{y-t}{\mathsf{V}(t)}dt$$

is the deviance function and $V(\cdot)$ is the variance function evaluated in y_i . The dispersion parameter is indexed by observation, allowing for flexible modeling. For example, experimental factors and/or covariates affecting the variability of the data may be encompassed. For proportion data, the method allows for the modeling of overdispersion as a function of a linear predictor that may differ from the one describing the mean.

The joint-modeling ideas for mean and dispersion, proposed by Smyth (1989) and Nelder & Lee (1991), all share the same double structure of generalized linear models. Assuming that E(Y) = μ and Var(Y) = ϕ V(μ), and that both the mean and the dispersion parameters vary across observations in a parametric way, i.e., $\eta_i = g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta}$ and $\zeta_i = h(\phi_i) = \mathbf{z}_i^T \boldsymbol{\gamma}$ and ⁹⁵ where β is a vector of mean parameters, γ is as vector of dispersion parameters, $g(\cdot)$ and $h(\cdot)$ ⁹⁶ are link functions for the mean and dispersion, and \mathbf{x}_i^T and \mathbf{z}_i^T are the row-vectors of the design ⁹⁷ matrices \mathbf{X} and \mathbf{Z} , respectively. The matrix \mathbf{X} contain covariates and/or factors affecting the ⁹⁸ mean, and the matrix \mathbf{Z} contains covariates and/or factors affecting the dispersion parameter. ⁹⁹ In this model, ϕ represents the independent variation of the mean and $V(\mu)$ is the mean-¹⁰⁰ dependent variation. Apart from this commonality between the modeling frameworks, they ¹⁰¹ exhibit particular aspects, too.

¹⁰² Parameter estimation proposed by Smyth (1989) and Nelder & Lee (1991) is based on ¹⁰³ a two-step iterative algorithm: (i) holding γ fixed, the vector β is estimated; (ii) fixing the ¹⁰⁴ estimated value of β , the vector γ is estimated. These two steps are then alternated until ¹⁰⁵ convergence. Although both proposals are based on different estimation methods, results are ¹⁰⁶ often very similar.

Nelder & Lee (1991) based estimation on extended quasi-likelihood. In their algorithm, 107 the step where ϕ is assumed fixed coincides with Smyth's (1989) method, thus reducing to 108 quasi-likelihood. When eta is fixed, the extended quasi-likelihood function becomes a gamma 109 likelihood function, where d_i is the response variable. Lee & Nelder (1998) also considered 110 an alternative for the estimation method based on REML with adjustment proposed by Cox 111 & Reid (1987). Lee & Nelder (2006) extended their proposal to a larger class of *double* 112 hierarchical generalized linear models, jointly incorporating random effects in both mean and 113 dispersion linear predictors. This class will not be explored in this work. 114

At this point, it is important to emphasize key differences between the JGLM and the 115 Bayesian DGLM explored here. The JGLM is a fixed-effects model that deals with disper-116 sion modeled in a particular way. This involves another generalized linear model for deviance 117 components as a response, a logarithmic link function and a linear predictor. The Bayesian 118 perspective for the proportion data, which will be described in Section 3, proceeds by hier-119 archically modeling the overdispersion through a random effect, where the linear predictor is 120 linked to the variance of the random effect. So, even though the results of both approaches 121 may lead to the same conclusions, the interpretations are different. 122

¹²³ **3** The Double Generalized Linear Model from a Bayesian Perspective

124 3.1 Model for Normally Distributed Measurements

The frequentist estimation approaches of Smyth (1989) and Nelder & Lee (1991) are clearly approximate and dependent on asymptotic assumptions. In agricultural experimentation, the number of experimental units is mostly limited owing to physical space, resources, and/or ethical constraints. Situations are common where the number of parameters is relatively large compared with the number of observations. The frequentist approach can then lead to strongly biased estimates (Smyth & Verbyla, 1999). An alternative way to tackle this problem is to work with the double generalized linear model (DGLM) from a Bayesian point of view.

One proposal for a Bayesian DGLM was presented by Cepeda & Gamerman (2000), with the following structure:

$$y_{i} = \mu_{i} + \varepsilon_{i},$$

$$\varepsilon_{i} \sim \mathsf{N}(0, \sigma_{i}^{2}),$$

$$\mu_{i} = \mathbf{x}_{i}^{T} \boldsymbol{\beta},$$

$$g(\sigma_{i}^{2}) = \mathbf{z}_{i}^{T} \boldsymbol{\gamma},$$
(1)

where $\mathbf{x}_i^T \boldsymbol{\beta}$ is the linear predictor for the mean μ , ε_i is the random component, the variance σ_i^2 , which is linked to the linear predictor $\mathbf{Z}\boldsymbol{\gamma}$ by a non-linear link function $g(\cdot)$, and \mathbf{x}_i^T and \mathbf{z}_i^T are known rows of design matrices for mean and dispersion, respectively. These authors assumed the following prior joint probability density function for β and γ :

$$\begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{pmatrix} \sim N \begin{bmatrix} \begin{pmatrix} \mathbf{b}_0 \\ \mathbf{g}_0 \end{pmatrix}, \begin{pmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{G} \end{bmatrix} \end{pmatrix}$$

where the hyper-parameters \mathbf{b}_0 , \mathbf{g}_0 , \mathbf{B} , \mathbf{C} , and \mathbf{G} are assumed known. The posterior joint probability density function is given by

$$\pi(\boldsymbol{\beta}, \boldsymbol{\gamma} | \mathbf{X}, \mathbf{Z}) = \frac{p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\beta}, \boldsymbol{\gamma}) p(\boldsymbol{\beta}, \boldsymbol{\gamma})}{\int_{\boldsymbol{\beta}} \int_{\boldsymbol{\gamma}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\beta}, \boldsymbol{\gamma}) p(\boldsymbol{\beta}, \boldsymbol{\gamma}) \partial \boldsymbol{\gamma} \partial \boldsymbol{\beta}}$$

$$\propto p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\beta}, \boldsymbol{\gamma}) p(\boldsymbol{\beta}, \boldsymbol{\gamma}).$$
(2)

As (2) assumes an intractable analytical form for integral manipulation, the Metropolis-Hat Hastings (MH) algorithm was employed, together with a block-wise scheme to obtain the samples of the posterior marginal density functions for each parameter (Gamerman, 1997).

3.2 Model for Overdispersed Proportion Data

The ideas behind a Bayesian DGLM for normal data do not carry over to the binomial situation,
 because in that case there is no separate variance parameter. Hinde & Demétrio (1998b)
 describe a logistic-normal model with the following structure

$$Y_i | \mathbf{z}_i \sim \mathsf{Bin}(m_i, p_i),$$

$$\mathsf{logit}(p_i) = \mathbf{x}_i^T \boldsymbol{\beta} + \sigma \mathbf{z}_i,$$

$$\mathbf{z}_i \sim \mathsf{N}(0, 1),$$
(3)

with the aim of accommodating the overdispersion effect through the random effect z_i . Borgatto *et al.* (2006) proposed a hierarchical random-effects model to account for both overdispersion and zero-inflation effects, as an alternative to the model described in Vieira *et al.* (2000). A Bayesian version of (3) was also proposed by Hinde & Demétrio (1998b), taking the form

$$Y_i | \mathbf{b}_i \sim \mathsf{Bin}(m_i, p_i),$$

$$\mathsf{logit}(p_i) = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{b}_i,$$

$$\mathbf{b}_i \sim \mathsf{N}(0, \sigma^2),$$
(4)

and assuming a prior distribution for β and $\tau = \sigma^{-2}$ to incorporate the uncertainty associated with these parameters.

Here, we propose a generalized version of (4), to allow for covariates and/or factors affecting the dispersion parameter of the random-effect distribution. To this end, the following hierarchical double generalized linear model is assumed:

$$Y_{i} \sim \text{Bin}(m_{i}, p_{i}),$$

$$\text{logit}(p_{i}) = \mathbf{x}_{i}^{T} \boldsymbol{\beta} + \delta_{i},$$

$$\delta_{i} \sim \text{N}(a, \tau_{i}),$$

$$\boldsymbol{\tau}_{i} = 1/\exp(\mathbf{z}_{i}^{T} \boldsymbol{\gamma}),$$
(5)

where \mathbf{x}_i and \mathbf{z}_i are the appropriate rows of the design matrices \mathbf{X} and \mathbf{Z} , respectively, δ_i is a random effect, $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are vectors of unknown parameters. The normal distribution for a random effect, used to accommodate overdispersion, is a sensible choice whenever this random

variable is required to range over the entire real line. Evidently, other distributions could be 160 entertained as well, such as, for example, a scaled t-distribution. We further assume that eta_i 161 and γ_k are independent, i.e., $p(\beta_i, \gamma_k) = p(\beta_i)p(\gamma_k)$, which is sensible given that it is difficult to 162 establish a prior dependence structure for these parameters in common experimental situations. 163 In this model, the link function for the mean of Y_i is $\mathsf{logit}(p_i) = \ln[p_i/(1-p_i)] = \ln[\mu_i/(m_i - p_i)]$ 164 $[\mu_i]$. The link function for the dispersion is assumed to be $\ln \tau_i^{-1}$, to enforce positive variance; 165 this can, of course, be modified to other monotone link functions, as appropriate. It was 166 assumed for eta and γ a priori to be normally distributed with known hyper-parameters specified 167 by $\beta_j \sim N(b,c)$, $j=0,\ldots,r$, and $\gamma_k \sim N(d,e)$, $k=0,\ldots,s$, respectively. The normal 168 priors with vague hyper-parameters is a way to establish non-informative uncertainty for the 169 parameters. 170

The posterior joint probability density function for model (5), obtained by the Bayes' rule, can be described by

$$p(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\delta} | \mathbf{X}, \mathbf{Z}, \mathbf{y}) \propto \mathsf{L}(\boldsymbol{\beta} | \boldsymbol{\delta}, \mathbf{y}, \mathbf{X}) p(\boldsymbol{\delta} | \boldsymbol{\gamma}, \mathbf{Z}) p(\boldsymbol{\beta}) p(\boldsymbol{\gamma}).$$
(6)

Model (5) can be represented by a directed acyclic graph (DAG), as described in Best & 173 Green (2005) as can be seen from Figure 1 in the Supplementary Materials. The advantage 174 of presenting a model in DAG form is that the essence of the model structure is elucidated, 175 making clear the functional flow of the information, thereby suppressing the distributional 176 assumptions and deterministic relations between variables and parameters. Moreover, such a 177 graphical model representation may suggest a conditional independence structure, convenient 178 for efficient implementation.The Bayesian computation environment OpenBUGS(Thomas *et* 179 al., 2006) was built to sample the posterior marginal distributions of the parameters under 180 DAGs that can be described graphically or through the BUGS language (Spiegelhalter et al., 181 1996). Best & Green (2005) and Thomas et al. (2006) provide more details and information 182 about directed acyclic graphs and the BUGS language. 183

¹⁸⁴ Sampling From the Posterior Marginal Densities

Assuming the priors for β and γ and for the random effect δ_i , and using the binomial likelihood function for Y_i , the posterior joint density function can be written as (see the Appendix for more details):

$$p(\boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\delta} | \mathbf{X}, \mathbf{Z}, \mathbf{y}) \propto \exp\left\{\mathbf{y}^{T}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\delta}) - \frac{1}{2}\mathbf{1}^{T}\mathbf{Z}\boldsymbol{\gamma} - \sum_{i=1}^{n}\exp(-\mathbf{z}_{i}^{T}\boldsymbol{\gamma})(\delta_{i} - a)^{2} - \frac{1}{2c}\sum_{j=0}^{r}(\beta_{j} - b)^{2} - \frac{1}{2e}\sum_{k=0}^{s}(\gamma_{k} - d)^{2}\right\} \times \prod_{i=1}^{n}[1 + \exp(\mathbf{x}_{i}^{T}\boldsymbol{\beta} + \delta_{i})]^{-m_{i}}.$$
(7)

From (7) it is not possible to derive analytic forms for the posterior marginal density functions for β , γ , and δ . Furthermore, it is not a viable alternative, neither to make use of numeric integration, because of its multi-dimensionality. Therefore, stochastic simulation of the posterior marginal densities, through the Monte Carlo Markov Chain (MCMC) methods, offers an appealing route.

So, to sample from the posterior joint density function (7), it is necessary to construct an appropriate Markov chain (Gamerman & Lopes, 2006), which can be done by using an MCMC algorithm, such as the Metropolis-Hastings algorithm, Gibbs sampling, or using a more general MCMC algorithm such as, for example, the slice sampler (Neal, 2003). All of these algorithms are based on the full posterior marginal density function, given by

$$p(\beta_j | \boldsymbol{\beta}_{-j}, \boldsymbol{\gamma}, \boldsymbol{\delta}, \mathbf{X}, \mathbf{Z}, \mathbf{y}) \propto \exp\left\{\mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \frac{1}{2c} \sum_{j=0}^{r-1} (\beta_j - b)^2 \right\} \prod_{i=1}^n [1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)]^{-m_i}(8)$$

$$p(\gamma_k | \boldsymbol{\gamma}_{-k}, \boldsymbol{\beta}, \boldsymbol{\delta}, \mathbf{X}, \mathbf{Z}, \mathbf{y}) \propto \frac{\exp\left\{-\frac{1}{2}\sum_{i=1}^{n-1}\exp(-\mathbf{z}_i^T\boldsymbol{\gamma})(\delta_i - a)^2 - \frac{1}{2e}\sum_{k=0}^s(\gamma_k - d)^2\right\}}{\sqrt{\exp\left(\mathbf{1}^T\mathbf{Z}\boldsymbol{\gamma}\right)}}, \quad (9)$$

$$p(\delta_i | \boldsymbol{\delta}_{-i}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \mathbf{X}, \mathbf{Z}, \mathbf{y}) \propto \frac{\exp\left\{\mathbf{y}^T \boldsymbol{\delta} - \frac{1}{2} \sum_{i=1}^{n-1} \exp(-\mathbf{z}_i^T \boldsymbol{\gamma}) (\delta_i - a)^2\right\}}{\prod_{i=1}^{n-1} [1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)]^{m_i}},$$
(10)

where the vectors β_{-j} , γ_{-k} , and δ_{-i} refer to the parameter vectors without the elements β_j , γ_k and δ_i , respectively.

These algorithms are implemented in the statistical computing environment R (R Devel-200 opment Core Team, 2007), using the library BRugs (Thomas *et al.*, 2006). The functions 201 related with this library use the BUGS language. It is therefore necessary and sufficient to 202 specify the model structure (5), the data set, and the initial values for each parameter, in 203 order to start the Markov chain iterations. When the sequence generated by the Markov chain 204 $\theta^t, t = 1, 2, \dots$ reaches convergence, the sample of θ^t can be considered a sample of $p(\theta|y)$. 205 It can be formally shown that the convergence of the chain is to the stationary distribution. 206 For a given set of data, it is important to construct some exploratory graphical diagnostics 207 and to obtain further diagnostic measures to scrutinize convergence (Gelman *et al.*, 2000). 208

The use of a Bayesian approach via stochastic simulation methods demands that Markov 209 chain for each parameter be examined for convergence, so as to guarantee that the sam-210 ples contain the principal characteristics of the equilibrium distribution, including shape, the 211 first few sample moments, etc. Some formal, e.g., test-based, methods, as well as informal 212 graphically-based ones, are quality indicators for the simulated samples. Here, some graphical 213 devices (history plot, auto-correlation function plot) and the Gelman-Rubin criterion (Gelman 214 & Rubin, 1992) were used for Markov chain convergence diagnosis. Alternatively, the test, 215 proposed by Raftery & Lewis (1992) and Heidelberger & Welch (1983), could be applied as 216 well. These tests are all implemented in the R environment through the coda library (Plummer 217 et al., 2007). 218

To perform model selection, the Deviance Information Criteria (DIC, Spiegelhalter *et al.*, 2002) was used. This index is calculated as $DIC = p_D + E_{\theta}[D(\theta)]$, where $p_D = E_{\theta}[D(\theta)] - D(E_{\theta}[p(\theta|y)])$, representing the effective number of parameters; $E_{\theta}[D(\theta)]$ is the average of 222 *D* calculated over all values of θ from the sample obtained from the MCMC algorithms; and 223 $D(E_{\theta}[p(\theta|y)])$ is the deviance measure calculated over the average of the sampled values of 224 θ . As is well known, this measure, as well as the AIC and BIC, quantify the model fitting and, 225 at the same time, penalize the complexity of the candidate models.

²²⁶ 4 The Apple Tissue Experiment

227 4.1 Data Description

As a motivation to the modeling tools developed here, an apple tissue culture experiment, described in Ridout & Demétrio (1992), is analyzed. The treatment structure was a $2 \times$ 5 factorial (2 media and 5 varieties), and the plot structure completely randomized. The experimental unit was a Petri dish, divided into a 5×5 array, hence having 25 individual compartments. In each compartment, a standard volume of some culture medium was used. A small piece of vegetation tissue, called *explant*, was added to the medium. The Petri dishes were kept in a incubator for several weeks. During this period, new shoots could grow from
the explants, which enhances the regeneration process. One aim of the researcher was to
establish whether some of the five explant varieties and/or one of the two culture media have
an influence on the proportion of regenerated explants. The data set is reproduced in Table 1.
The motivation to use a relatively small set of data is twofold. First, it allows focusing
on the methodological contributions, without the intricacies of large and potentially complex
data manipulations. Second, this type of experiment is quite common in horticulture.

Table 1 ABOUT HERE

Figure 1 shows a plot with the average of regenerated explants, for each combination of 242 explant variety and culture method. There clearly is a lot of variability between the means 243 of regenerated explants in culture medium X, when compared with that in culture medium 244 Y. Explant E appears to show differential behavior when compared to the other explants in 245 culture medium X, suggesting a possible interaction. Ridout & Demétrio (1992) leave open 246 the option that it could be a potential outlier. That said, we believe that the value of 2/25 for 247 the E/X cell is scientifically plausible. Hence, the value was retained for analysis. Figure 2 248 shows the standard deviations for the treatment combinations. It suggests that explant D249 probably exerts strong influence on the dispersion. It is important to note that it is a small 250 and unbalanced data set and, in general, the asymptotic theory does not apply. 251

252

241

FIGURES 1 and 2 ABOUT HERE

4.2 Data Analysis

Model (5) was applied to the apple tissue data set and, to estimate the posterior marginal density, three Markov chains were used, with initial points dispersed across the parameter space.

²⁵⁷ Vague priors were assumed so as to incorporate the uncertainty for the vectors β and γ . ²⁵⁸ In all cases, they were assumed normally distributed, with variance 1000 for the β parameters ²⁵⁹ and 100 for the γ . While the value 100 may appear not sufficiently vague, choosing a larger ²⁶⁰ value tends to jeopardize the convergence of the Markov chain process.

Three chains of size 200,000 were generated. The first 100,000 iterations of each chain 261 were discarded for burn-in purposes. In an effort to minimize the within-chain autocorrelation, 262 each 50th iteration was retained. As a result, the sample size to be used for posterior inference 263 about the parameters is 6,000. The models with less parameters have shown faster convergence 264 of the Markov chains, indicating a better mix throughout the parameter space. The final model 265 took 281 seconds under a Intel Core Duo processor on a 1.66 GHz personal computer. The R 266 code using the BRugs and CODA libraries are available at the website containing supplementary 267 material, with name BayesianDGLM.r. Note that the same model specification can be used 268 with the OpenBUGS/WinBUGS software. 269

Using the deviance information criteria (DIC), backward model selection was conducted. The most complex model fitted was

> η = medium + explant + medium × explant, ζ = medium + explant + medium × explant.

Table 2 presents each fitted model, the deviance information criteria, and their components pD and \overline{D} . Models 1 to 4 describe the search for a parsimonious linear predictor of dispersion. Model 5 minimized the DIC index; it features a linear predictor for the dispersion using the dummy variable for the explant D; for the mean it incorporates main effects of medium, explant, and the interaction term of the culture medium Y and the explant E. Model 6 has ²⁷⁷ a constant term only in the dispersion linear predictor. Model 7 has only a random effect δ ²⁷⁸ with distribution N(0,1), and Model 8 is an ordinary binomial GLM without random effects.

In Model 5, the non-significant effects in dispersion model, related with explants B, C, and E, were dropped, in view of reducing the DIC. In Model 3, these non-significant terms were kept and while a consistent picture emerges in terms of pD and the number of parameters, the DIC has deteriorated.

A strong overdispersion effect is apparent for Model 8. Models 9, 10, and 11 are nested for the linear predictor for the mean; they fit worse than Model 5.

TABLES 2 and 3 ABOUT HERE

Model checking plots are not usually presented in the Bayesian literature, where the focus 286 is often directed towards MCMC diagnostics. However, some applied manuscripts present 287 a type of residual analysis that mimics the frequentist approach, thereby considering some 288 summary measure of the sampled marginal posterior for deviance residuals and predicted values 289 (Robinson et al., 2009). Using the mean of sampled marginal posteriors, some graphical 290 analysis were done. Figure 3 exhibits the Q-Q plot and the standardized residuals versus the 291 predicted values, which do not indicate departure from normality or any systematic pattern 292 that could suggest lack-of-fit. The 95% credibility interval plot for the standardized residual 293 posterior samples, displayed in Figure 4, does indicate neither outliers nor extreme values, 294 providing support for inference over the parameter estimates. 295

Table 3 presents the posterior summary. Using the median as a point estimate for the 296 model parameters, it can be seen that explant D increases the variance of the random effect 297 δ by $\exp(-6.624 + 7.533) = 2.48$ units, while the other type of explants reduce the variance 298 of the said random effect to almost zero. We conclude that explant D is responsible for the 299 overdispersion in the proportion of regenerated explants. Turning to the factorial effects on 300 the mean of the regenerated explants, the larger influence owes to the interaction between 301 culture medium Y and explant E. Whenever this combination occurs, the chance of explant 302 regeneration increases $\exp(2.185) = 8.89$ -fold, when compared to the other explants and 303 culture medium X. 304

305

285

FIGURES 3 and 4 ABOUT HERE

5 Concluding Remarks

The objective of this work was to propose a Bayesian double generalized linear model for 307 overdispersed proportion data, thereby providing an alternative to the frequentist approach of 308 Smyth (1989) and Nelder & Lee (1991). Our model includes a normally distributed random 309 effect in the linear predictor of the generalized linear model, where the variance of the random 310 effect is linked non-linearly to another linear predictor. The model was successfully applied 311 to an agricultural data set of a type frequently encountered. Evidently, it can be applied to 312 related situations too. Experiments with the aim of identifying factors affecting the variability 313 in industrial processes, such as Taguchi experiments, can be analyzed using this approach as 314 well. 315

The Bayesian approach adopted here and based on stochastic simulation is a very convenient and flexible mode for fitting our class of models. It conveniently extends to other such situations, without limitation to the exponential family of distributions. The use of directed acyclic graphs makes easier the presentation of the model and suggests the hierarchical construction of the probabilistic model, based on conditioning the random values on their parent random parameters. The use of MCMC algorithms efficiently deals with the problem of calculating high-dimensional integrals, thus allowing to generate samples of the posterior marginal densities of the associated parameters. Using such samples, summary statistics are calculated that render feasible inferences about the fitted model. The price to pay is the method's computational intensity, needing high-quality computational resources. Furthermore, the additional step of diagnosing the convergence of the Markov chain guarantee the coherence of the inferences and has to be taken seriously. The statistical computing environment R, jointly with the libraries BRugs and coda have shown quite flexible and efficient to the estimation process and data analysis.

The proposed modeling framework can be extended in various ways, including the incorpo-330 ration of the time dimension when measurements are taken longitudinally, as well as particular 331 implementations or count and time-to-event data. Also, the assumption of normal prior dis-332 tributions can be modified or relaxed. Here, they have shown to be reasonably insensitive 333 to variations of the hyper-parameters for the mean linear predictor, even though they are 334 sensitive to the choice of vague priors for the parameters in the linear predictor of the ran-335 dom effect variance; this calls for careful illicitation, preferably with the help of substantive 336 researchers. Finally, identifiability, parameterization and tests with other prior densities need 337 further research. 338

Acknowledgments

The first and third authors were partially supported by The Brazilian National Council for Scientific and Technological Development (CNPq). Part of this work was developed during the sandwich scholarship program of the first author, funded by Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), processo BEX 4344/07-3. We thank Dr. Rosângela Loschi (UFMG) for the useful discussion of ideas. The fourth author gratefully acknowledges grant #P6/03 of the Belgian Government (Belgian Science Policy).

6 References

- Aitkin, M. (1987) Modelling variance heterogeneity in normal regression using GLIM.
 Applied Statistics, **36**, 332–339.
- 2. Best, N. and Green, P. (2005) Structure and uncertainty: graphical models for understanding complex data. *Significance*, 177–181.
- Borgatto, A.F., Demétrio, C.G.B., and Leandro, R.A. (2006) Modelos para proporões
 com superdispersão e excesso de zeros um procedimento Bayesiano. Revista de
 Matemática e Estatística, **24** 121–131.
- 4. Box, G. (1988) Signal-to-noise ratios, performance criteria and transformations. Technometrics, 30:1, p.1–40.
- 5. Box, G.E.P., Cox, D.R. (1964) An analysis of transformations. Journal of the Royal Statistical Society, B, 26, p. 211-52.
- 6. Cepeda, E. and Gamerman, D. (2000) Bayesian modeling of variance heterogeneity in normal regression models. *Brazilian Journal of Probability and Statistics*, **14**, 207–221.
- ³⁶⁰ 7. Cox, D.R. and Reid, N. (1987) Parameter orthogonality and approximate conditional ³⁶¹ inference. *Journal of the Royal Statistical Society B*, **49**, 1–39.
- 8. Gamerman, D. (1997) Sampling from the posterior distribution in generalized linear mixed models. *Statistics and Computing*, **7**, 57–68.

- Gamerman, D. and Lopes, H. (2006) Markov Chain Monte Carlo: Stochastic Simulation
 for Bayesian Inference. London: Chapman & Hall.
- 10. Gelfand, A.E. and Smith, A.F.M. (1990) Sampling-based approaches to calculating marginals. *Journal of the American Statistical Association*, **87**, 523–532.
- I1. Gelman, A., Carlin, J.B., Stern, H.S., and Rubin, D.B. (2000) *Bayesian Data Analysis*.
 Boca Raton: Chapman & Hall/CRC.
- 12. Givens, G.H. and Hoeting, J.A. (2005) *Computational Statistics*. Hoboken: John Wiley & Sons.
- 13. Harvey, A.C. (1976) Estimating regression models with multiplicative heteroscedasticity.
 Econometrica, 44, 461–465.
- 14. Heidelberger, P. and Welch, P.D. (1983) Simulation run length control in presence of a initial transient. *Operations Research*, **31**, 1109–1144.
- 15. Hinde, J. and Demétrio, C.G.B. (1998a) Overdispersion: models and estimation. *Computational Statistics & Data Analysis*, **27**, 151–170.

16. Hinde, J. and Demétrio, C.G.B. (1998b) Overdispersion: Models and Estimation. In:
 Simpósio Nacional de Probabilidade e Estatística, Caxambú: Associaão Brasileira de
 Estatística 13.

- ³⁸¹ 17. Jørgensen, B. (1987) Exponential dispersion models. *Journal of the Royal Statistical* ³⁸² *Society, Series B*, **49**, 127–162.
- ³⁸³ 18. Lee, Y. and Nelder, J.A. (1998) Generalized linear models for the analysis of quality-³⁸⁴ improvement experiments. *The Canadian Journal of Statistics*, **26**, 95–105.
- 19. Lee, Y. and Nelder, J.A. (2006) Double hierarchical generalized linear models. *Applied Statistics*, **55**, 139–185.
- ³⁸⁷ 20. Lee, Y., Nelder, J.A., and Pawitan Y. (2006) *Generalized Linear Models with Random* ³⁸⁸ *Effects: Unified Analysis via H-likelihood.* Boca Raton: Chapman & Hall/CRC.
- McCullagh, P. and Nelder, J.A. (1989) *Generalized Linear Models (2nd ed.)* London:
 Chapman & Hall.
- ³⁹¹ 22. Myers, R.H., Khuri, A.I., Vining, G. (1992) Response surface alternatives to the Taguchi ³⁹² robust parameter design approach. *The American Statistician*, 46, 2, 131-139.
- ³⁹³ 23. Neal, R.M. (2003) Slice sampling. *The Annals of Statistics*, **31**, 705–767.
- ³⁹⁴ 24. Nelder, J.A. and Lee, Y. (1991) Generalized linear models for the analysis of Taguchi-³⁹⁵ type experiments. *Applied Stochastic Models and Data Analysis*, **7**, 107–120.
- ³⁹⁶ 25. Nelder, J.A. and Pregibon, D. (1987) An extended quasi-likelihood function. *Biometrika*,
 ³⁹⁷ 74, 221–232.
- Paulino, C.D., Turkman, M.A.A., and Murteira, B. (2003) *Estatística Bayesiana*. Lisboa:
 Fundaão Calouste Gulbenkian.
- Plummer, M., Best, N., Cowles, K., and Vines, K. (2007) coda: Output Analysis and
 Diagnostics for MCMC. R package version 0.12-1.

- Raftery, A.E. and Lewis, S.M. (1992) One long run with diagnostics: implementations
 strategies for Monte Carlo Markov Chains. *Statistical Science*, **7**, 493–497.
- ⁴⁰⁴ 29. R Development Core Team. (2007) *R: A Language and Environment for Statistical* ⁴⁰⁵ *Computing*, version 2.6.0. Vienna: R Foundation for Statistical Computing.
- ⁴⁰⁶ 30. Ridout, M. and Demétrio, C.G.B. (1992) Generalized linear models for positive count ⁴⁰⁷ data. *Revista de Matemática e Estatística*, **10**, 139–148.
- Robinson, T.J., Anderson-Cook, C.M., Hamada, M.S. (2009) Bayesian analysis of split plot experiment with nonnormal response for evaluating nonstandard performance crite ria. Technometrics, 51:1, p. 56–65.
- ⁴¹¹ 32. Smyth, K. (1989) Generalized linear models with varying dispersion. *Journal of the* ⁴¹² *Royal Statistical Society Series B*, **51**, 47–60.
- ⁴¹³ 33. Smyth, G.K. and Verbyla, A.P. (1999) Adjusted likelihood methods for modelling dis-⁴¹⁴ persion in generalized linear models. *Environmetrics*, **10**, 696–709.
- 34. Spiegelhalter, D., Thomas, A., Best, N., and Gilks, W. (1996) BUGS 0.5: Bayesian
 Inference Using Gibbs Sampling Manual (version ii). Cambridge: Medical Research
 Council Biostatistics Unit.
- 35. Spiegelhalter, D., Best, N.G., Carlin, B.P., and Van Der Linde, A. (2002) Bayesian
 measures of model complexity and fit. *Journal of the Royal Statistical Society, Series* B, 64, 583–639.
- 36. Taguchi, G. (1985) Quality engineering in Japan. Communication in Statistics: Theory
 and Methods, 14, 2785–2801.
- 37. Thomas, A., O'Hara, B., Ligges, U., and Sturtz, S. (2006) Making BUGS open. *R News*, **6**, 12–17.
- ⁴²⁵ 38. Vieira, A.M.C., Hinde, J.P., and Demétrio, C.G.B. (2000) Zero-inflated proportion data ⁴²⁶ models applied to a biological control assay. *Journal of Applied Statistics*, **27**, 373–389.
- ⁴²⁷ 39. Wolfinger, R.D. and Tobias, R.D. (1998) Joint estimation of location, dispersion and ⁴²⁸ random effects in robust design. *Technometrics*, **40**, 62–70.

429 Appendix

The posterior density function for model (5) can be built using (6). The prior density functions for β_j and γ_k are, respectively

$$p(\beta_j) = \sqrt{\frac{c^{-1}}{2\pi}} \exp\left[-\frac{c^{-1}}{2}(\beta_j - b)^2\right] \propto \exp\left[-\frac{c^{-1}}{2}(\beta_j - b)^2\right], \qquad j = 0, \dots, r$$
(11)

432 and

$$p(\gamma_k) = \sqrt{\frac{e^{-1}}{2\pi}} \exp\left[-\frac{e^{-1}}{2}(\gamma_k - d)^2\right] \propto \exp\left[-\frac{e^{-1}}{2}(\gamma_k - d)^2\right], \qquad k = 0, \dots, s.$$
(12)

Assuming that $Y_i \sim \text{Bin}(m_i, p_i)$ and that $\ln[p_i/(1-p_i)] = \mathbf{x}_i^{ \mathrm{\scriptscriptstyle T}} \boldsymbol{\beta} + \delta_i$, the likelihood function is

$$p(y_i | \mathbf{x}_i, \boldsymbol{\beta}, \delta_i) = \binom{m_i}{y_i} \left[\frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)} \right]^{y_i} \left[1 - \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)} \right]^{m_i - y_i}$$
$$= \binom{m_i}{y_i} \frac{\left[\exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)\right]^{y_i}}{\left[1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta} + \delta_i)\right]^{m_i}},$$

434 and

$$p(y_i|\mathbf{x}_i,\boldsymbol{\beta},\delta_i) \propto \left[\exp(\mathbf{x}_i^T\boldsymbol{\beta}+\delta_i)\right]^{y_i} \left[1+\exp(\mathbf{x}_i^T\boldsymbol{\beta}+\delta_i)\right]^{-m_i}.$$
(13)

435 The conditional density of δ_i , given the parameter vector $\boldsymbol{\gamma}$, is

$$p(\delta_i|\boldsymbol{\gamma}) = \sqrt{\frac{1}{2\pi \exp(\boldsymbol{z}_i^T \boldsymbol{\gamma})}} \exp\left\{-\frac{1}{2}[\exp(\boldsymbol{z}_i^T \boldsymbol{\gamma})]^{-1}(\delta_i - a)^2\right\}$$
$$= (2\pi)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}\boldsymbol{z}_i^T \boldsymbol{\gamma} - \frac{1}{2}\exp(-\boldsymbol{z}_i^T \boldsymbol{\gamma})(\delta_i - a)^2\right\}$$

436 and, therefore,

$$p(\delta_i|\boldsymbol{\gamma}) \propto \exp\left\{-\frac{1}{2}[\boldsymbol{z}_i^T\boldsymbol{\gamma} + \exp(-\boldsymbol{z}_i^T\boldsymbol{\gamma})(\delta_i - a)^2]\right\}.$$
(14)

Applying (11)-(14) to (6), considering the vector of observations y, it can be shown that the posterior joint probability density function is

$$p(\boldsymbol{\beta},\boldsymbol{\gamma},\boldsymbol{\delta}|\mathbf{X},\mathbf{Z},\mathbf{y}) \propto \prod_{i=1}^{n} \frac{[\exp(\mathbf{x}_{i}^{T}\boldsymbol{\beta}+\delta_{i})]^{y_{i}}}{[1+\exp(\mathbf{x}_{i}^{T}\boldsymbol{\beta}+\delta_{i})]^{m_{i}}} \exp\left\{-\frac{1}{2}[\mathbf{z}_{i}^{T}\boldsymbol{\gamma}+\exp(-\mathbf{z}_{i}^{T}\boldsymbol{\gamma})(\delta_{i}-a)^{2}]\right\} \times \\ \times \prod_{j=0}^{r} \exp\left[-\frac{c^{-1}}{2}(\beta_{j}-b)^{2}\right] \prod_{k=0}^{s} \exp\left[-\frac{e^{-1}}{2}(\gamma_{k}-d)^{2}\right] \\ \propto \exp\left\{\sum_{i=1}^{n} \mathbf{x}_{i}^{T}\boldsymbol{\beta}y_{i} + \sum_{i=1}^{n} y_{i}\delta_{i} - \frac{1}{2}\sum_{i=1}^{n} \mathbf{z}_{i}^{T}\boldsymbol{\gamma} - \frac{1}{2}\sum_{i=1}^{n} \exp(-\mathbf{z}_{i}^{T}\boldsymbol{\gamma})(\delta_{i}-a)^{2} - \frac{1}{2c}\sum_{j=0}^{r}(\beta_{j}-b)^{2} - \frac{1}{2e}\sum_{k=0}^{s}(\gamma_{k}-d)^{2}\right\} \prod_{i=1}^{n} [1+\exp(\mathbf{x}_{i}^{T}\boldsymbol{\beta}+\delta_{i})]^{-m_{i}}$$

439 and, finally,

$$p(\boldsymbol{eta}, \boldsymbol{\gamma}, \boldsymbol{\delta} | \mathbf{X}, \mathbf{Z}, \mathbf{y}) \propto \exp \left\{ \mathbf{y}^{\mathrm{T}}(\mathbf{X}\boldsymbol{eta} + \boldsymbol{\delta}) - \frac{1}{2} \mathbf{1}^{\mathrm{T}} \mathbf{Z} \boldsymbol{\gamma} - \sum_{i=1}^{n} \exp(-\mathbf{z}_{i}^{\mathrm{T}} \boldsymbol{\gamma}) (\delta_{i} - a)^{2} - \right.$$

$$-\frac{1}{2c}\sum_{j=0}^{r}(\beta_j-b)^2 - \frac{1}{2e}\sum_{k=0}^{s}(\gamma_k-d)^2 \bigg\} \times \\ \times \prod_{i=1}^{n}[1 + \exp(\mathbf{x}_i^T\boldsymbol{\beta} + \delta_i)]^{-m_i}.$$
(15)

Table 1: Number of explants (y_i) of apple trees that regenerated, considering 16 Petri dishes. SOURCE: Hinde & Demétrio (1998b).

	Explant Variety					
		А	В	C	D	E
Culture Medium	Х	8, 10	9, 10	7	20, 12	2
	Y	9,11	18, 12	13	20, 5	13

Table 2: Quality-of-fit and model complexity measures for fitted Bayesian binomial DGLM models.

Fitted Models	DIC	pD	$E_{\theta}[D(\theta)]$
$Model 1: \left\{ egin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium} imes \texttt{explant} \ \zeta = \texttt{medium} + \texttt{explant} + \texttt{medium} imes \texttt{explant} \end{array} ight.$	86,31	13,15	73,16
$Model 2: \left\{ egin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium} imes \texttt{explant} \ \zeta = \texttt{medium} + \texttt{explant} \end{array} ight.$	86,10	13,44	72,66
$Model 3: \left\{ egin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium} imes \texttt{explant} \ \zeta = \texttt{explant} \end{array} ight.$	84,69	12,79	71,9
$Model \ 4: \left\{ egin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium} imes \texttt{explant} \ \zeta = \gamma_0 + \texttt{explant} \ D \end{array} ight.$	84,49	18,81	71,67
$Model 5: \left\{ \begin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium Y} \times \texttt{explant E} \\ \zeta = \gamma_0 + \texttt{explant D} \end{array} \right.$	82,56	11,00	71,56
$Model 6: \left\{ \begin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium Y} \times \texttt{explant E} \\ \zeta = \gamma_0 \end{array} \right.$	87,16	14,52	72,64
$Model 7: \left\{ \begin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium Y} \times \texttt{explant E} \\ \zeta = 0 \Rightarrow \boldsymbol{\delta} \sim N(0, 1) \end{array} \right.$	86,16	14,77	71,39
$Model 8: \left\{ \begin{array}{l} \eta = \texttt{medium} + \texttt{explant} + \texttt{medium} \texttt{Y} imes \texttt{explant} \texttt{E} \end{array} ight.$	107,10	7,03	100,10
$Model \; 9: \left\{ egin{array}{l} \eta = \texttt{medium} + \texttt{explant} \ \zeta = \gamma_0 + \texttt{explant} \; D \end{array} ight.$	85,77	11,19	74,58
$Model 10: \left\{ egin{array}{l} \eta = \mathtt{explant} \ \zeta = \gamma_0 + \mathtt{explant} \mathtt{D} \end{array} ight.$	88,66	14,37	74,29
$Model 11: \left\{ egin{array}{l} \eta = \mathtt{medium} \ \zeta = \gamma_0 + \mathtt{explant} \mathtt{D} \end{array} ight.$	85,18	8,82	76,30

 Table 3: Model 5: posterior summary for the parameters.

		Standard	Standard			
Parameters	Mean	Deviation	Error	IC(2.5%)	Median	IC(97.5%)
β_0	-0.8034	0.2784	0.004268	-1.34200	-0.80550	-0.2595
β_Y	0.5925	0.2900	0.003788	0.02434	0.59190	1.1540
β_B	0.4657	0.3246	0.004614	-0.17570	0.46720	1.0840
β_C	0.0795	0.4101	0.004761	-0.75690	0.08712	0.8642
β_D	0.8618	1.1030	0.019380	-1.21000	0.84280	3.2070
eta_E	-1.8840	0.9127	0.011330	-3.92100	-1.80400	-0.3263
$\beta_{Y,E}$	2.1850	1.0190	0.013750	0.41940	2.09300	4.3770
γ_0	-6.6240	3.9760	0.170800	-15.79000	-5.84400	-1.1140
γ_D	7.5330	4.0870	0.171400	1.46700	6.75600	16.8100

Summary statistics calculated for 6000 observations.



Figure 1: Average number of regenerated explants as function of the culture media and the explant varieties.



Figure 2: Standard deviation of the number of regenerated explants as function of the culture media and the explant varieties.



Mean of posterior standardized residuals vs predicted



Figure 3: Residual analysis for Model 5. The top plot is the mean for samples of marginal posteriors for the standardized residuals versus the mean of posterior marginal samples of predicted values; the bottom plot is the normal probability plot for the mean marginal posterior of the standardized residuals.



Figure 4: 95% credibility interval plots for marginal posterior samples of standardized residuals.

View publication stats