

Estimating Marginal Properties of Quantitative Real-Time PCR Data Using Nonlinear Mixed Models

Daniel Gerhard,^{1,*} Melanie Bremer,² and Christian Ritz³

¹Institute of Biostatistics, Leibniz Universität Hannover, Herrenhäuser Straße 2, 30419 Hannover, Germany

²Institute of Plant Nutrition, Leibniz Universität Hannover, Herrenhäuser Straße 2,
30419 Hannover, Germany

³Department of Nutrition, Exercise and Sports, Rolighedsvej 30, 1958 Frederiksberg C, Denmark

*email: gerhard@biostat.uni-hannover.de

SUMMARY. A unified modeling framework based on a set of nonlinear mixed models is proposed for flexible modeling of gene expression in real-time PCR experiments. Focus is on estimating the marginal or population-based derived parameters: cycle thresholds and $\Delta\Delta c(t)$, but retaining the conditional mixed model structure to adequately reflect the experimental design. Additionally, the calculation of model-average estimates allows incorporation of the model selection uncertainty. The methodology is applied for estimating the differential expression of a phosphate transporter gene OsPT6 in rice in comparison to a reference gene at several states after phosphate resupply. In a small simulation study the performance of the proposed method is evaluated and compared to a standard method.

KEY WORDS: Delta–delta C_t value; Delta method; Gauss-Hermite quadrature; Hierarchical design; Model averaging.

1. Introduction

Quantitative real-time PCR is a well established methodology to quantify differential gene expression in various fields of application (Bustin, 2000). Real-time PCR allows a fast and inexpensive analysis of gene expression with a high precision. After screening a large amount of genes in a microarray or other large scale gene expression experiment, real-time PCR is an accepted method to verify the findings by re-analyzing differential expression of a smaller set of significant genes of interest.

Real-time PCR methodology allows us to directly observe the amplification of cDNA templates at each of a sequence of PCR cycles by measuring a fluorescence signal. Beginning with a small amount of target DNA of a single sample, a fluorescence curve can be observed with an exponential increase of the signal with the ongoing sequence of PCR cycles until depletion effects prevent a further increase, or the process is stopped at an early phase. The time at which the increase of the number of cDNA templates becomes visible is proportional to the amount of templates in the original sample; hence, an early increase of the fluorescence curve indicates a highly expressed gene with a large amount of mRNA in the sample.

Standard approaches to model the change of fluorescence intensity with increasing cycle number are assuming an exponential shaped fluorescence curve, applying a simple linear regression on a logarithmic transformed scale (Gibson, Heid, and Williams, 1996). Only the observations at the first increase of fluorescence are showing this (log-)linear relationship; hence all observations for later cycle numbers need to be omitted to maintain an adequate fit to the data. Therefore this approach highly depends on an algorithm to search for a

specific window of cycle numbers, where all assumptions are met. Alternative modeling approaches for determining $c(t)$ rely on sigmoidal models that incorporate all observations over the whole range of cycle numbers (Rutledge, 2004; Goll et al., 2006; Spiess, Feig, and Ritz, 2008; Swillens, Dessars, and El Housni, 2008). An exponential model can be seen as an approximation only to the lower part of the sigmoidal curve. In addition to these empirical models, several mechanistic models have been derived directly from the PCR amplification process (Spiess et al., 2008).

Based on the assumed model, a common approach for data analysis in relative quantitative real-time PCR experiments is to make comparisons by determining changes in the expression of different target genes or genes in different treatments relative to a reference. Changes are quantified in terms of differences in cycle thresholds (denoted $c(t)$ or C_t), such that the difference $\Delta c(t)$ (often denoted ΔC_t) corresponds to normalization of a target gene against a reference gene. Normalized expression levels are then compared across target genes or treatments by means of differences (of differences) and this approach is known as the $2^{-\Delta\Delta C_t}$ method (Livak and Schmittgen, 2001). In the context of a linear regression model using \log_2 -transformed cycle numbers, the value $2^{-\Delta\Delta C_t}$ can be interpreted as the differential expression relative to a chosen standard (Livak and Schmittgen, 2001). The cycle thresholds are often obtained from an exponential model, using linear regression on a \log_2 transformed scale of cycle numbers (Gibson et al., 1996). Subsequently, multiple linear regression or analysis of covariance may be used for comparing cycle thresholds obtained from different treatments (Yuan et al., 2006). Based on the estimated $c(t)$ for each sample, classical hypotheses testing approaches, like an analysis of

variance or permutation test procedures, are used to test for significant gene by treatment effects. Hence, all inferences for gene and treatment effects are based on a model, where the estimates for derived parameters from separate amplification models per sample are treated as new observations.

Quantitative real-time PCR experiments are often planned in a hierarchical design with multiple technical replicates within each biological replicate. In such cases a mixed model approach is needed to account for variation between and within biological replicates (Pinheiro and Bates, 2000). One approach is to estimate and carry out hypothesis testing using a two-step approach where estimates of $c(t)$ obtained from separate nonlinear regression models fitted to each technical replicate are combined in subsequent analysis using a linear mixed model (Steibel et al., 2009). However, this approach has at least three shortcomings: (1) The uncertainty related to estimating the $c(t)$ values is not carried on into the second step of the analysis, resulting in an inefficient analysis. (2) No insights are gained on how experimental variation depends on biological and technical replicates. (3) Only cutoff values smaller than the smallest upper asymptote or values greater than the largest lower asymptote are available.

In this article, we propose a modeling approach for real-time PCR data based on a nonlinear mixed model framework that allows incorporating the uncertainty in model selection through a model averaging step using a number of candidate models, which may be motivated empirically or mechanistically. In addition, we derive marginal or population-based threshold cycles by means of an elaborate combination of Gaussian-Hermite quadrature and the Delta method. An R package from the corresponding author makes the proposed method easy to use.

2. Hierarchical Modeling of Real-Time PCR Data

We consider fluorescence intensities $\{y_{ijk}\}$ obtained from a hierarchical design using a number of biological replicates ($i = 1, \dots, I$), which are again divided into a number of technical replicates ($j = 1, \dots, J$), each with PCR cycle numbers c_k ($k = 1, \dots, K$).

The fluorescence intensity is assumed to be proportional to the amount of PCR amplified target at each cycle number. At the beginning of a real-time PCR experiment only a small amount of PCR material can be observed, followed by a considerable increase in fluorescence as the PCR template at each cycle depends on the prior one. As time elapses this process begins to slow down, mostly due to depletion of reaction components needed to complete a PCR cycle. This means that the minimum requirement for an appropriate model is to have two horizontal asymptotes as cycle numbers approaches 0 and infinity. Therefore fluorescence intensity is assumed to be a smooth nonlinear function of cycle number that follows an s-shaped pattern until depletion.

In order to combine the nested data structure with the nonlinear pattern, we consider a two-stage hierarchical nonlinear mixed-effects regression model (Davidian and Giltinan, 1995, pp. 98–112). Below we elaborate on the model specification suitable for real-time PCR data and then in Section 5 we formulate the specific model for our data example.

At the first stage we specify the nonlinear relationship for each individual fluorescence intensity:

$$y_{ijk} = f(c_k, \boldsymbol{\beta}_{ij} + \mathbf{u}_{ij}) + \epsilon_{ijk},$$

assuming a specific nonlinear model function f depending on the cycle number c_k and R -dimensional vectors of fixed- and random-effects contributions $\boldsymbol{\beta}_{ij}$ and \mathbf{u}_{ij} , respectively. The error terms follow a normal distribution $N(0, \sigma^2)$.

There exist a number of s-shaped model functions that have been used in applications of nonlinear regression to biology, medicine, and toxicology, including log-logistic, log-normal, and Weibull-type models (Bates and Watts, 1988; Ritz, 2010). We will only consider a single model, the five-parameter log-logistic model, which previously has been used for modeling real-time PCR data (Spiess et al., 2008):

$$\begin{aligned} f(c_k, \boldsymbol{\beta}_{ij}) &= f\{c_k, (\beta_{ij}^{(1)}, \dots, \beta_{ij}^{(5)})\} \\ &= \beta_{ij}^{(2)} + \frac{\beta_{ij}^{(3)} - \beta_{ij}^{(2)}}{(1 + \exp[\beta_{ij}^{(1)}\{\log(c_k) - \log(\beta_{ij}^{(4)})\}])^{\beta_{ij}^{(5)}}}, \end{aligned} \quad (1)$$

where $\beta_{ij}^{(1)}$ characterize the steepness in the s-shaped curve, $\beta_{ij}^{(2)}$ and $\beta_{ij}^{(3)}$ correspond to the lower and upper asymptotes, $\beta_{ij}^{(4)}$ denotes the approximate location of the inflection point, and $\beta_{ij}^{(5)}$ is an asymmetry parameter where positive values above or below 1 correspond to differences in curvature close to the lower and upper asymptotes.

At the second stage we specify the decomposition of the model parameters at the level of the technical replicates. For the r th parameter in the vector $\boldsymbol{\beta}_{ij}$ the fixed- and random effects contributions are:

$$\begin{aligned} \boldsymbol{\beta}_{ij}^{(r)} &= \left(\mathbf{X}_T^{(r)}\right)_j \otimes \left(\mathbf{X}_B^{(r)}\right)_i \boldsymbol{\beta}, \\ \mathbf{u}_{ij}^{(r)} &= \left(\mathbf{Z}_B^{(r)}\right)_i \mathbf{u}_B^{(r)} + \left(\mathbf{Z}_T^{(r)}\right)_j \mathbf{u}_T^{(r)}, \end{aligned}$$

where $\mathbf{X}_B^{(r)}$ ($I \times p_1$) and $\mathbf{X}_T^{(r)}$ ($J \times p_2$) denote the design matrices of the fixed-effects structures at the level of the biological and technical replicates, respectively (the subscript refers to specific rows in the matrices), and $\boldsymbol{\beta}$ denotes the $p_1 p_2$ -dimensional fixed-effects parameter. Any kind of covariate information available at the level of the biological replicates could be included in the model. For our data example $\mathbf{X}_B^{(r)} = \mathbf{X}_B$ specifies a time trend that is assumed to be present in all r parameters, whereas $\mathbf{X}_T^{(r)} = \mathbf{X}_T$ groups technical replicates corresponding to the same gene into clusters (two clusters within each biological replicate as we consider two genes).

Similarly, $\mathbf{Z}_B^{(r)}$ ($I \times q_1$) and $\mathbf{Z}_T^{(r)}$ ($J \times q_2$) are the design matrices of the random effects associated with the biological replicates and technical replicates within biological replicates,

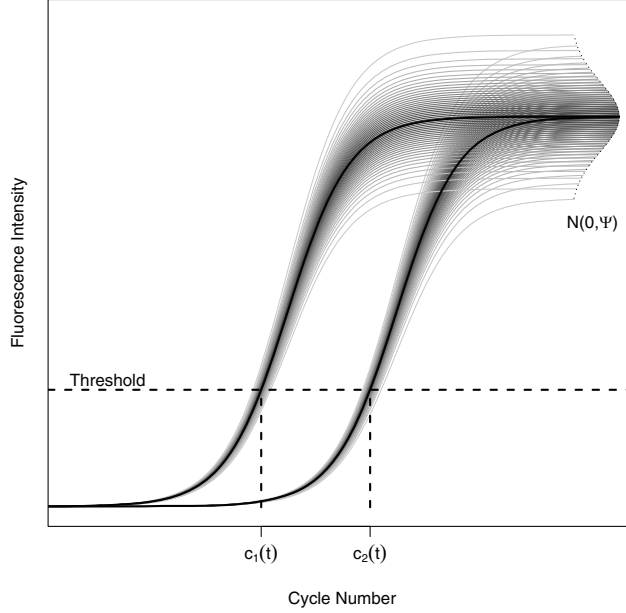


Figure 1. Schematic representation of $c(t)$ estimation assuming a five-parameter log-logistic model for each of two genes and additionally assuming random effects assigned to the upper asymptotes only. Thus, marginal and conditional prediction coincide. In this special case, the two variance components for the asymptotes of both curves are equal. The black lines denote the marginal predictions and the gray lines correspond to the predictions for each individual. The difference between $c_1(t)$ and $c_2(t)$ is denoted as $\Delta c(t)$.

respectively. The random effects are assumed to follow normal distributions:

$$\mathbf{u}_B^{(r)} \sim N(\mathbf{0}, \Psi_B^{(r)}),$$

$$\mathbf{u}_T^{(r)} = (\mathbf{u}_{T1}^{(r)}, \dots, \mathbf{u}_{Tl}^{(r)}) \sim N(\mathbf{0}, \Psi_T^{(r)}).$$

In principle $\Psi_B^{(r)}$ and $\Psi_T^{(r)}$ may be unstructured variance-covariance matrices, but we restrict ourselves to diagonal matrices that correspond to uncorrelated random effects, apart from letting $\Psi_B^{(r)}$ be a block diagonal matrix with diagonal entries that are themselves diagonal matrices $\Psi_{B1}^{(r)}, \dots, \Psi_{Bp_2}^{(r)}$ allowing different variance components for different groups (Davidian and Giltinan, 1995, pp. 122–124). For our data example we will include random intercepts both for the biological and technical replicates and, additionally, we assume gene-specific variance components for the biological replicates. However, we will not include random effects in all model parameters as there need not necessarily be variation in the biological or technical replicates for all parameters just as a consequence of the technology used. For instance, Figure 1 clearly shows that there is effectively no variation at all for cycle numbers below 14–15 and thus there is no need for random effects being assigned to the parameter in the model corresponding to the lower asymptote as cycle numbers approach 0.

Parameter estimation can be performed with any statistical software allowing the fitting of nonlinear mixed models, for example, the **R** package nlme (Pinheiro and Bates, 2000) or SAS PROC NLMIXED (SAS Institute, 2011).

3. Estimating Threshold Cycles and Their Differences

It is established practice to evaluate real-time PCR data by means of the threshold cycle summary measure instead of directly interpreting the parameters in f (Gibson et al., 1996). The threshold cycle is defined as the cycle number where the mean fluorescence level reaches a certain cutoff intensity t , illustrated in Figure 1. There exist a number of recommendations regarding the choice of t : If only the lower part of the entire curve is fitted using an exponential model, the cutoff value should be located near the lower asymptote (Gibson et al., 1996). However for a sigmoidal model fitted to the entire data set, an intensity close to the inflection point might be the optimal choice for the cutoff threshold as the corresponding cycle number estimate obtains a lower variance in comparison to cutoff values near the asymptotes (Liu, Udhe-Stone, and Goudar, 2011).

Determining the threshold is an inverse regression problem that in the case of a nonlinear regression model for a single replicate has the solution $c(t) = f^{-1}(t)$. For a nonlinear mixed model, the threshold cycle is obtained by solving the equation $E(f(c, \beta_{0ij} + \mathbf{u}_0)) = t$ in c for some specific fixed-effects parameter configuration denoted β_{0ij} . Solving the equation requires repeated evaluation of the integral by integrating out the random effects \mathbf{u}_0 . As we assume uncorrelated random effects the expectation is equal to:

$$\int \dots \int f\{c, (\beta_{0ij}^{(1)} + u_1, \dots, \beta_{0ij}^{(R)} + u_R)\} \phi_1(u_1) \dots \phi_R(u_R) du_1 \dots du_R. \quad (2)$$

where ϕ_r ($r = 1, \dots, R$) denotes the univariate Gaussian density with mean 0 and variance equal to the sum of the diagonal entries: $\{\mathbf{Z}_B^{(r)} \Psi_B^{(r)} (\mathbf{Z}_B^{(r)})^t\}_{ii} + \{\mathbf{Z}_T^{(r)} \Psi_T^{(r)} (\mathbf{Z}_T^{(r)})^t\}_{jj}$. Thus, evaluation of the integral in equation (2) simplifies to a number of univariate integrals.

As pointed out by a number of authors the marginal expectation (as given in Eq. 2) is in general not equal to the conditional expectation $f(c, \beta_{0ij} + 0)$ for nonlinear (and generalized) mixed models (e.g., Davidian and Giltinan, 1995, p. 121; McCulloch, Searle, and Neuhaus, 2008, pp. 236–238). Simplistic solutions to approximate marginal predictions include averaging of the predictions at the level of the technical replicates and replacing the estimated random effects (the EBLUPs) by 0 in the predictions, that is using the conditional expectation for a hypothetical average technical replicate (deMiguel et al., 2012). Actually, the latter approach results in the correct marginal predictions if random effects are only included for model parameters that enter the model equation f in a linear way (Nielsen, Ritz, and Streibig, 2004). For instance, f in equation (1) is a linear function in the parameters $\beta_{ij}^{(2)}$ and $\beta_{ij}^{(3)}$ corresponding to the upper and lower asymptotes. However, if in equation (1) random effects are assigned on the slope, inflection point, or asymmetry parameters, the integral in equation (2) has to be approximated numerically to obtain marginal predictions and this has to be done repeatedly to obtain $c(t)$, which can only be found numerically by solving the equation $E(f(c, \beta_{0ij} + \mathbf{u}_0)) = t$, for example, using the bisection method.

A convenient and fast method for approximating integrals involving Gaussian densities is the Gauss-Hermite quadrature (McCulloch et al., 2008, pp. 327–331). By using the product rule the integral in equation (2) can be approximated as follows:

$$\sum_{n_1=1}^{N_1} \dots \sum_{n_R=1}^{N_R} w_{n_1}^{(1)} \dots w_{n_R}^{(R)} f\{c, (\beta_{0ij}^{(1)} + x_{n_1}^{(1)}, \dots, \beta_{0ij}^{(R)} + x_{n_R}^{(R)})\},$$

where the nodes $x_{n_r}^{(r)}$ and weights $w_{n_r}^{(r)}$ are derived using the variances of the univariate Gaussian distributions involved in equation (2) (Smyth, 2005). By the continuous mapping theorem, plugging in consistent estimators of the fixed-effects and variance parameters results in a consistent estimator of the Gauss-Hermite quadrature approximation, and the accuracy of the approximation is determined by the number of nodes. Specifically for the log-logistic model specified in equation (1), we would expect the Gauss-Hermite quadrature approximation, or even the Laplace approximation ($N_1 = \dots = N_R = 1$), to work well in view of the similarities in the integral with the logistic regression example considered by Liu and Pierce (1994).

By approximating the integral in equation (2) the resulting estimate of $c(t)$ becomes a nonlinear function that depends on the given cutoff value t , the fixed-effect parameters (β), and the variance parameters for the random effects. By treating the estimated variance parameters as non-random in the approximation, the delta method can readily be applied for deriving the standard error of the estimate of $c(t)$ (van der Vaart, 1998, Chap. 5). The key observation is that the estimated variance-covariance matrix of the fixed-effects parameter estimates is available from most statistical software programs used for fitting nonlinear mixed-effects models, but not so for the estimated variance-covariance matrix of the estimated variance parameters or the estimated covariance matrix between the fixed-effects parameter estimates and the estimated variance parameters (McCulloch et al., 2008). Additionally, for the delta method to work in practice, numerical approximation of the derivative of the estimate of $c(t)$ is needed as this estimate is only implicitly given as the solution to an equation. The comparison between two $c(t)$ values for two different target genes or treatment groups is based on the estimated difference between the two marginal $c(t)$ values.

After the estimation of a set of marginal $c(t)$ parameters for each gene, together with their corresponding variance and covariance parameters, any linear combination of $c(t)$ values can be computed by specifying adequate contrasts. This approach corresponds to the general contrast definition in Steibel et al. (2009), which allows for example to compare to the average of the $c(t)$ values of two different control genes.

To fully exploit the flexibility inherent in the modeling framework described in Section 2, we consider a model averaging approach that allows us to combine $c(t)$ estimates from a number of nonlinear mixed-effects models to incorporate the uncertainty in the model selection (Buckland, Burnham, and Augustin, 1997). It is not only the uncertainty about the nonlinear shape of the curves that can be represented in a set of candidate models, but also the uncertainty about the com-

plexity of the model in terms of the number of fixed effects to include.

Consider M different candidate models that may correspond to different choices of f and/or different hypothesized fixed-effects structures for the different model parameters, but they are not different w.r.t. the random effects structure, which we believe to a large extent is determined by the experimental design. The M estimated cycle thresholds of interest $\hat{c}_1(t), \dots, \hat{c}_M(t)$ from the individual models are combined into the model-average estimate through a weighted average:

$$\hat{c}(t) = \sum_{m=1}^M w_m \hat{c}_m(t).$$

The weights w_1, \dots, w_M should favor candidate models closest to the true, but unknown model. A commonly used approach for calculating weights is based on the well known Akaike information criterion (AIC), which is an estimator for the relative distance between two models (Akaike, 1973; Burnham and Anderson, 2002). The corresponding weights are estimated as follows:

$$w_m = \frac{\exp(-\frac{1}{2}\Delta_m)}{\sum_{m=1}^M \exp(-\frac{1}{2}\Delta_m)} \quad (m = 1, \dots, M),$$

using the difference Δ_m in AIC between model m and the model with minimum AIC among the M candidate models. For our data example, the shape of the curve can be clearly distinguished as the fluorescence is observed for a large number of successive PCR cycles with a small residual error within each technical replicate, and thus we expect that a few models will be much preferred over the other candidate models as long as the candidate models have little overlap. For model-average estimates approximate standard errors can be obtained using the formula provided by Buckland et al. (1997) and then confidence intervals are derived assuming that the distribution of the estimate $\hat{c}(t)$ is approximately normal. As long as the threshold t is not near the lower or upper asymptote we would expect the normal approximation to provide a good approximation. In case the number of biological and technical replicates is small, a t -distribution with degrees of freedom, determined by the grouping level at which the fixed-effects terms are estimated, could replace the normal approximation (Pinheiro and Bates, 2000, p. 91).

4. Simulation Study

To evaluate the performance of the method the power to detect a difference in $\Delta c(t)$ values is investigated by simulation. A simple data setting is chosen with two genes and two treatments, 3–10 biological replicates, and three technical replicates for each gene within each biological replicate. Fluorescence intensities are generated using parameter values derived from the data example (see Table 1 and Web Table 1): $\beta^{(1)} = -10$, $\beta^{(2)} = 0$, $\beta^{(3)} = 750$, $\beta^{(4)} = 20$, $\beta^{(5)} = 2$ for both genes and treatments. For simplicity, biological and technical variation is only introduced in the upper asymptote and the inflection point (between biological replicates: $\psi_B^{(3)} =$

Table 1

Estimated standard deviations of the random effects for the biological and technical replicates and residual errors for the model with the highest associated model weight

| Parameter | ψ_B | | ψ_T | σ |
|----------------|----------|-------|----------|----------|
| | eEf | OsPT6 | | |
| Steepness | 0.36 | 0.22 | 0.16 | 4.00 |
| Upper asym. | 66.54 | 43.87 | 32.89 | |
| Infl. point | 0.45 | 0.21 | 0.08 | |
| Residual error | | | | |

50, $\psi_B^{(4)} = 0.3$; between technical replicates: $\psi_T^{(3)} = 30$, $\psi_T^{(4)} = 0.07$; residual error: $\sigma = 5$).

For comparison we consider an adaptation of the two-step approach introduced by Steibel et al. (2009) where estimation of a $c(t)$ value is based on a five-parameter log-logistic model fitted for each technical replicate separately and combined with subsequent estimation of $\Delta\Delta c(t)$ using a linear mixed model based on the $c(t)$ estimates obtained for each technical replicate as a new response. More specifically, in the first step for each of the $I \times J$ technical replicate the estimated $c(t)$ value is easily calculated by inverse regression (given a cutoff t) from the separate nonlinear model fits using the parameter estimates θ_{ij} :

$$y_{ijk} = f_{ij}(c_k, \theta_{ij}) + \epsilon_{ijk}, \quad c_{ij}(t) = f_{ij}^{-1}(t).$$

In the second step, the treatment effects are estimated based on the previously obtained $c(t)$ estimates:

$$\hat{c}_{ij}(t) = (X_T \otimes X_B) \boldsymbol{\vartheta} + \mathbf{Z}_B \mathbf{u}_B + \xi_{ij}.$$

The same fixed-effects treatment structure as in the proposed nonlinear hierarchical mixed-effects model is used and the additional gene-specific random effects are modeled by \mathbf{u}_B on the level of the biological replicates. However, in contrast to the full nonlinear mixed-effects model, the variation between the technical replicates and the residual variation in the $\hat{c}(t)$ cannot be separated apart. Assuming a linear mixed effects model allows to interpret the treatment effects directly as marginal effects. The $\Delta\Delta c(t)$ estimates are obtained by linear contrasts of elements in $\boldsymbol{\vartheta}$.

For each parameter setting 10,000 simulation runs are performed.

A typical real-time PCR experiment consists of only a small number of replicates. Even a small increase in sample size will cause a large increase in power to detect differentially expressed genes. Therefore we anticipate that by using a single nonlinear mixed model the number of parameters is reduced without an actual loss of information as compared to the use of separate models for each technical replicate, and in fact this effect can be seen for the simulated data settings (Figure 2). A further advantage of using a single model in comparison to the combination of $c(t)$ values from separate models is the independence of the chosen cutoff value. Figure 3 shows that a constant power is reached for the mixed model approach at

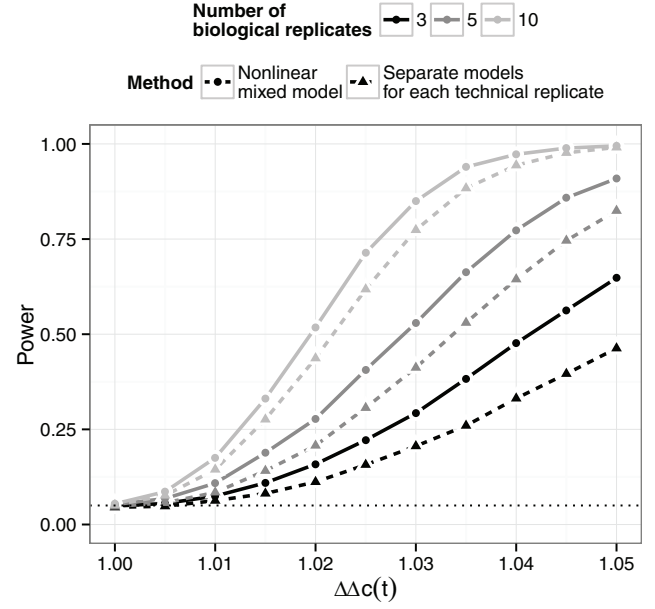


Figure 2. Simulated power for different numbers of biological replicates comparing the nonlinear mixed model approach to the approach using separate models. A threshold of $t = 100$ (near the lower asymptote) is assumed when estimating $c(t)$.

all chosen cutoff values. On the other hand, the power is decreasing for the separate models approach: the more the cutoff value approaches the upper asymptote the greater is the loss in power. This effect will increase with increasing variability in the upper asymptote as the approach using separate models does not take the uncertainty in estimating $c(t)$ into account.

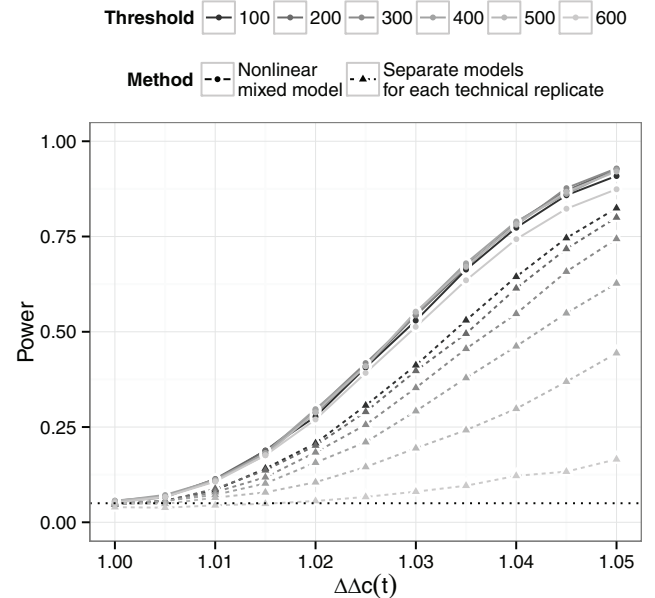


Figure 3. Simulated power comparing the nonlinear mixed model approach to the approach using separate models, assuming different cutoffs t when estimating $c(t)$.

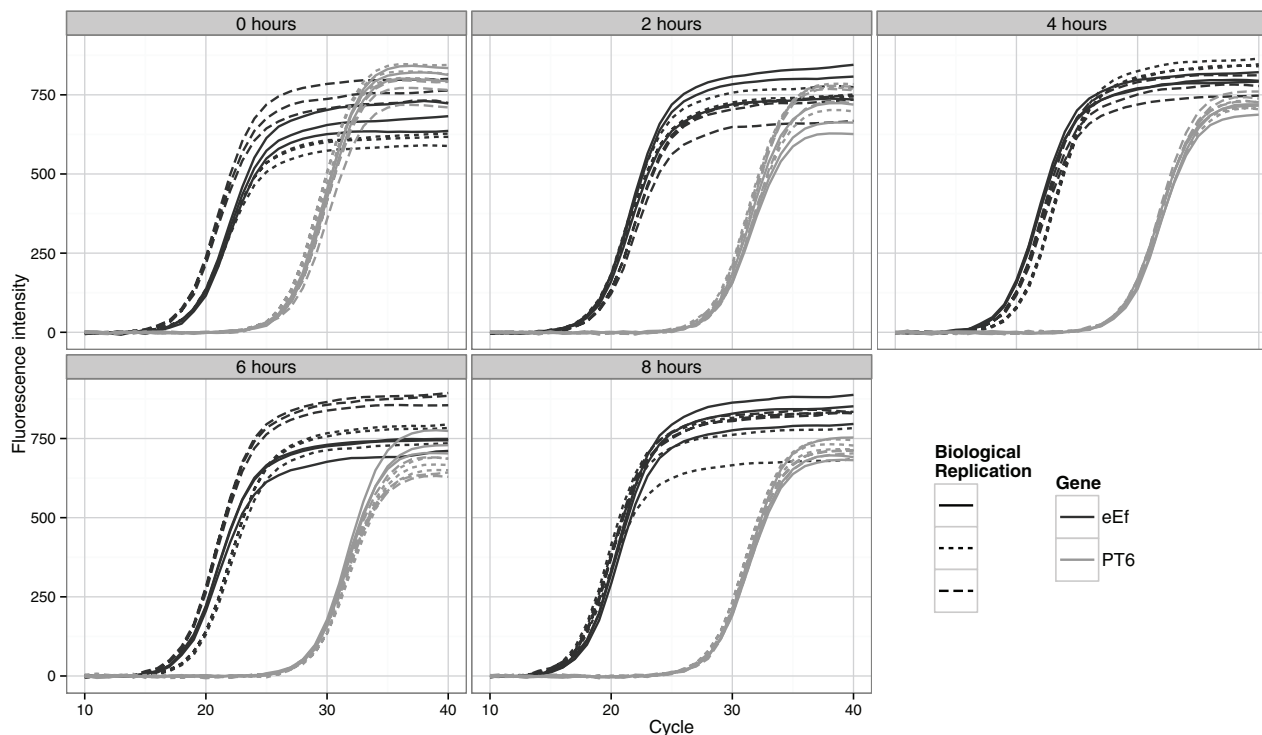


Figure 4. Data example showing fluorescence intensity proportional to gene expression for each of the 40 PCR cycles. Note: The first 10 PCR cycles are omitted in the graph). The gene of interest (OsPT6) and the reference gene (eEf) are shown as black and gray lines for 0–8 hours after resupplying phosphate in the nutrient solution (in five separate panels). For each of the biological replicates there are three technical replicates for each gene as indicated by dotted lines. Different biological replicates are used for different durations of resupply (in total 15 biological replicates).

5. Example

Plants should be able to adapt to specific nutritional changes in their environment. With the regulated expression of phosphate transporters, plants are able to maintain a constant level of phosphate. We consider data from an experiment where rice (*Oryza sativa* L.) plants were starved for phosphate for 5 days, resulting in an upregulation of the phosphate transporter OsPT6 in the root to increase the uptake of the limited nutrient. Subsequent resupply of phosphate should decrease the expression of this transporter (Ai et al., 2009) to prevent phosphate influx to a toxic level. The expression of OsPT6 and the reference gene eEf (eukaryotic elongation factor) were observed at 0, 2, 4, 6, and 8 hours after phosphate resupply for three biological replicates, which are thought of as representatives of the underlying biological system. For each biological replicate and gene the fluorescence intensity was measured over 40 PCR cycles for three technical replicates using reverse transcription polymerase chain reaction. A graphical representation of the dataset is shown in Figure 4.

For our data example, we consider a number of five-parameter log-logistic models where the fixed-effects structure consists of gene-specific linear, quadratic, cubic, and 4th order polynomial time trends (i.e., time after phosphate resupply) in the lower asymptote, inflection point, slope, and asymmetry parameters but a common upper asymptote (intercept) parameter. All combinations of polynomial representations of the time effect up to the 4th order for each of the nonlinear

model parameters, results in a total of 625 candidate models. Any further linear-combination of parameters other than polynomial functions can be added to the set to represent the time effect; but these are not considered here any further to simplify the analysis and reduce computation time. The variability between the biological and technical replicates is modeled by introducing random effects for the biological and also the technical replicates on the upper asymptote, slope, and inflection point parameters. We assume gene-specific variance components for the random effects for the biological replicates. The parameters are estimated by maximum likelihood using the software implementation in the R package nlme (Pinheiro and Bates, 2000).

When comparing the set of 625 five-parameter log-logistic models with a set of four-parameter (log)-logistic, Weibull, or log-normal models, all model weights are assigned to the five-parameter log-logistic models. Thus, the log-logistic model seems to be the most adequate choice among these models, and including an asymmetry parameter highly improves the model fit to the actual data, which agrees with the findings by Spiess et al. (2008). The model averaging procedure by BIC picks mainly two different models out of the 625 candidates. The best model with a weight of 0.42 assumes a linear model for the steepness, only an intercept for the lower asymptote and a quadratic time dependency for the inflection point and the asymmetry parameter. The fixed effects estimates of this model are presented in Web Table 1. A second model with a weight of 0.34 assumes only

an intercept for the steepness parameter, but apart from that uses the same parameterization as the first model. The remaining weights of 0.10 and 0.13 substitute the quadratic function of the inflection point with a cubic shape.

Figure 4 shows that already at the beginning of the experiment the amount of transcripts for gene *OsPT6* is less abundant than for the reference gene *eEf*. This fact is also captured by the fixed-effects estimate of the inflection point, which is 10 cycle numbers larger for *OsPT6* as compared to *eEf* (see Web Table 1). The estimated slope and the asymmetry parameters also show large differences between the two genes. Within the nonlinear mixed model framework it would be possible to carry out hypothesis testing to evaluate more precisely which parameters differ between the two genes (e.g., Davidian and Giltinan, 1995). However, as mentioned previously, the interest with real-time PCR analysis lies in derived parameters that have a clear interpretation in a population context.

Using the nonlinear mixed model approach we are able to quantify how environmental variation, which influences the real-time PCR process, is divided between biological and technical replicates and, more specifically, to see how variation splits on the different model parameters. Note that the estimated residual standard error (within technical replicates) is roughly 4, which has a magnitude that is negligibly small in comparison to the range of fluorescence intensities (0–900) but also in comparisons to the variation found between biological and technical replicates for the upper asymptote. For all model parameters the estimated variance components for the biological replicates are larger than for technical replicates (Table 1).

As described in Section 3, differential expression of *OsPT6* with reference to *eEf* is summarized by a threshold cycle estimate. The quadratic time trend included as part of the fixed-effect structure in the fitted model makes it possible to derive a regression model for $\Delta\Delta c(t)$ as a function of time elapsed (in hours) after phosphate resupply (h):

$$\Delta\Delta c(t)(h) = \frac{c_1(t)(h)}{c_2(t)(h)} \bigg/ \frac{c_1(t)(0)}{c_2(t)(0)}.$$

With increasing time the estimated $\Delta\Delta c(t)(h)$ is decreasing (Figure 5). Hence, in comparison to the standard gene *eEf* the $c(t)$ values for *OsPT6* are becoming increasingly larger, showing a decrease in relative expression of this gene, as an earlier increase in fluorescence is proportional to the amount of transcripts in the sample. After roughly $3\frac{1}{2}$ hours the point-wise confidence bands do no longer include 1, corresponding to no significant differential expression.

6. Discussion

The proposed procedure for estimating $c(t)$, $\Delta c(t)$, and $\Delta\Delta c(t)$ based on a nonlinear mixed model has a number of advantages as compared to commonly used alternatives: It offers a unified framework for evaluating $c(t)$ in terms of estimation and hypothesis testing instead of using estimated cycle thresholds obtained separately for each replicate and then used as input data in a subsequent analysis. Additionally, by capturing replication effects through a random effects component

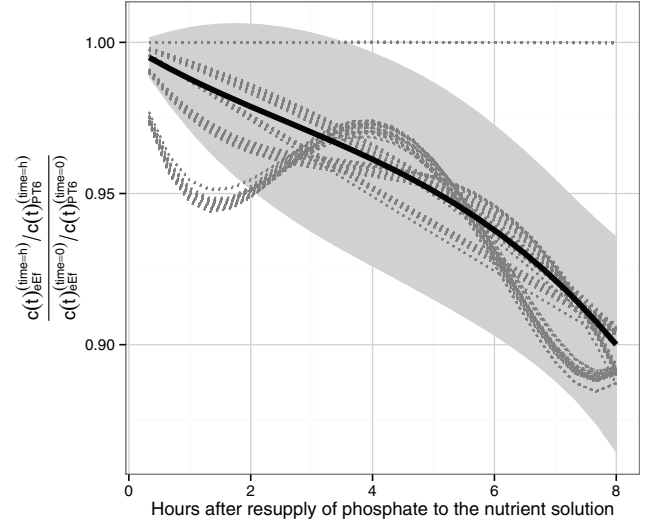


Figure 5. Marginal $\Delta\Delta c(t)$ (relative to time 0) comparing the transporter gene *OsPT6* to the reference gene *eEf*. Dashed lines show the estimates for 625 different models, representing the time effect as all combinations of polynomials of order 1, 2, 3, and 4 per parameter of the five-parameter log-logistic model (except the upper asymptote with a common intercept). The solid line displays the model averaged estimates as a quadratic function in time, with a point-wise $(1 - \alpha) = 0.95$ confidence band highlighted by a subjacent gray region.

instead of fitting a separate nonlinear model per replicate, the total number of estimated parameters is reduced, and this may yield in predictions with a higher precision and effect estimates with better reproducibility. Analyzing real-time PCR data within a single nonlinear mixed model framework also offers better model diagnostics directly on the level of fluorescence measurements. The proposed methodology has been implemented in the **R** package *qpcrnlme*, which is available under <https://github.com/daniel-gerhard/qpcrnlme>.

However, there are also a few shortcomings: As the nonlinear mixed model combining data from all technical replicates is usually more complex than the models applicable for a single replicate, more care and effort has to be taken to ensure that the estimation procedure is successful. An additional limitation in fitting more complex nonlinear mixed models is the availability of easy to use software implementations, but once an adequate set of models has been defined and implemented, the process of analyzing similar real-time PCR experiments can be automated as the experimental design and data structure remains the same most of the time.

Usually multiple technical replicates are analyzed on the same plate in a single real-time PCR process step and therefore we assumed that all technical replicates have the same cycle numbers. The proposed modeling approach, however, extends straightforward to experimental designs using different cycle numbers of cycles in different technical replicates and to designs with different numbers of technical replicates in the different biological replicates. The estimation procedure is not confined to only one nonlinear model function f as the model averaging step easily can be extended to include several types of models (e.g., log-normal or extreme value models

that also result in s-shaped curves). Extension to correlated random effects should also be possible using sparse grids for multivariate Gauss-Hermite quadrature (Heiss and Winschel, 2008).

7. Supplementary Materials

Web Tables referenced in Sections 4 and 5 are available with this paper at the *Biometrics* website on Wiley Online Library.

REFERENCES

- Ai, P., Sun, S., Zhao, J., Fan, X., Xin, W., Guo, Q., Yu, L., Shen, Q., Wu, P., Miller, A. J., and Xu, G. (2009). Two rice transporters, OsPht1;2 and OsPht1;6, have different functions and kinetic properties in uptake and translocation. *The Plant Journal* **57**, 798–809.
- Akaike, H. (1973). Information theory and the maximum likelihood principle. In *2nd International Symposium in Information Theory*.
- Bates, D. M. and Watts, D. G. (1988). *Nonlinear Regression Analysis and Its Applications*. New York: Wiley.
- Buckland, S., Burnham, K., and Augustin, N. (1997). Model selection: An integral part of inference. *Biometrics* **53**, 603–618.
- Burnham K. and Anderson D. (2002). *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. New York: Springer.
- Bustin, S. (2000). Absolute quantification of mRNA using real-time reverse transcription polymerase chain reaction assays. *Journal of Molecular Endocrinology* **25**, 169–193.
- Davidian, M. and Giltinan, D. M. (1995). *Nonlinear Models for Repeated Measurement Data*. London: Chapman and Hall.
- de-Miguel, S., Mehtatalo, L., Shater, Z., Kraid, B., and Pukkala, T. (2012). Evaluating marginal and conditional predictions of taper models in the absence of calibration data. *Canadian Journal of Forest Research* **42**, 1383–1394.
- Gibson U. E., Heid C. A., and Williams P. M. (1996). A novel method for real time quantitative RT-PCR. *Genome Research* **6**, 995–1001.
- Goll R., Olsen T., Cui G., and Florholmen J. (2006). Evaluation of absolute quantitation by nonlinear regression in probe-based real-time PCR. *BMC Bioinformatics* **7**, 107.
- Heiss F. and Winschel V. (2008). Likelihood approximation by numerical integration on sparse grids. *Journal of Econometrics* **144**, 62–80.
- Liu Q. and Pierce D. A. (1994). A note on Gauss-Hermite quadrature. *Biometrika* **81**, 624–629.
- Liu M., Udhe-Stone C., and Goudar C. T. (2011). Progress curve analysis of qRT-PCR reactions using the logistic growth equation. *Biotechnology Progress* **27**, 1407–1414.
- Livak K. J. and Schmittgen T. D. (2001). Analysis of relative gene expression data using real-time quantitative PCR and the $2^{-\Delta\Delta C_T}$ method. *Methods* **25**, 402–408.
- McCulloch, C. E., Searle, S. R., and Neuhaus, J. M. (2008). *Generalized, Linear, and Mixed Models*, 2nd edition. London: Wiley.
- Nielsen O. K., Ritz C., and Streibig J. C. (2004). Nonlinear mixed-model regression to analyze herbicide dose–response relationships. *Weed Technology* **18**, 30–37.
- Pinheiro J. and Bates D. (2000). *Mixed-Effects Models in S and S-Plus. Statistics and Computing*. New York: Springer.
- Ritz, C. (2010). Towards a unified approach to dose–response modeling in ecotoxicology. *Environmental Toxicology & Chemistry* **29**, 220–229.
- Rutledge R. (2004). Sigmoidal curve-fitting redefines quantitative real-time PCR with the prospective of developing automated high-throughput applications. *Nucleic Acids Research* **32**, e178.
- SAS Institute, SAS Publishing (2011). *Sas/Stat 9.3 User's Guide: Mixed Modeling (Book Excerpt)*. Cary, NC: SAS Institute Inc.
- Smyth, G. K. (2005). Numerical integration. In *Encyclopedia of Biostatistics*, P. Armitage and T. Colton (eds), 3766–3773. London: Wiley.
- Spiess A. N., Feig C., and Ritz C. (2008). Highly accurate sigmoidal fitting of real-time PCR data by introducing a parameter for asymmetry. *BMC Bioinformatics* **9**, 221.
- Steibel J. P., Poletto R., Coussens P. M., and Rosa G. J. M. (2009). A powerful and flexible linear mixed model framework for the analysis of relative quantification RT-PCR data. *Genomics* **94**, 146–152.
- Swillens S., Dessars B., and El Housni H. (2008). Revisiting the sigmoidal curve fitting applied to quantitative real-time PCR data. *Analytical Biochemistry* **373**, 370–376.
- van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge: Cambridge University Press.
- Yuan, J. S., Reed, A., Chen, F., and Stewart, C. N., Jr. (2006). Statistical analysis of real-time PCR data. *BMC Bioinformatics* **7**, 85.

Received December 2012. Revised October 2013.

Accepted October 2013.