

## Bayesian Growth Curves Using Normal Mixtures With Nonparametric Weights

Luisa Scaccia & Peter J Green

To cite this article: Luisa Scaccia & Peter J Green (2003) Bayesian Growth Curves Using Normal Mixtures With Nonparametric Weights, Journal of Computational and Graphical Statistics, 12:2, 308-331, DOI: [10.1198/1061860031725](https://doi.org/10.1198/1061860031725)

To link to this article: <https://doi.org/10.1198/1061860031725>



Published online: 01 Jan 2012.



Submit your article to this journal [↗](#)



Article views: 64



View related articles [↗](#)



Citing articles: 3 View citing articles [↗](#)

# Bayesian Growth Curves Using Normal Mixtures With Nonparametric Weights

Luisa SCACCIA and Peter J. GREEN

Reference growth curves estimate the distribution of a measurement as it changes according to some covariate, often age. We present a new methodology to estimate growth curves based on mixture models and splines. We model the distribution of the measurement with a mixture of normal distributions with an unknown number of components, and model dependence on the covariate through the weights, using smooth functions based on B-splines. In this way the growth curves respect the continuity of the covariate and there is no need for arbitrary grouping of the observations. The method is illustrated with data on triceps skinfold in Gambian girls and women.

**Key Words:** Allocation; Bayesian hierarchical model; Centile curves; Finite mixture distributions; Heterogeneity; Markov chain Monte Carlo; Normal mixtures; Path sampling; Reversible jump algorithms; Semiparametric model; Splines; Split/merge moves.

## 1. INTRODUCTION

Centile reference charts are an important screening tool in medical practice. The general form of a centile chart is a series of smoothed curves, showing how selected centiles for a biometrical measurement, such as height, weight, or middle-upper-arm-circumference, change when plotted against some independent, appropriate covariate, often age. Here we will refer to these curves as centile or growth curves. This second name comes from the fact that such charts are used widely in pediatrics, for measurements related to growth and development.

On the basis of the centile curves for a certain biometrical measurement, it is possible to identify patients who are unusual, in the sense that their value for that measurement lies in the tails of the reference distribution. Centiles are usually chosen from a symmetric subset of the 3rd, 5th, 10th, 25th, 50th, 75th, 90th, 95th, and 97th.

The simplest way to draw the centile curves is to calculate the empirical centiles,

---

Luisa Scaccia is Research Associate, Dipartimento di Scienze Statistiche, Università di Perugia, via A. Pascoli C.P. 1315/succ1, 06123 Perugia, Italy (E-mail: luisa@stat.unipg.it). Peter J. Green is Professor, School of Mathematics, University of Bristol, Bristol, BS8 1TW, United Kingdom (E-mail: P.J.Green@bristol.ac.uk).

©2003 American Statistical Association, Institute of Mathematical Statistics,  
and Interface Foundation of North America

*Journal of Computational and Graphical Statistics*, Volume 12, Number 2, Pages 308–331  
DOI: 10.1198/1061860031725

after grouping or smoothing with respect to the covariate. If empirical centiles are used, however, the more extreme are estimated relatively inaccurately, as the centile standard errors increase steeply towards the tails of the distribution. This problem can be overcome by fitting a theoretical distribution to the data and thereby obtaining the expected centiles. In doing so, two main criteria must be met:

1. Generally the covariate concerned takes continuous values and the distribution of the measurements can be assumed to vary smoothly with the covariate; centile curves should therefore be constructed in such a way to respect this continuity, and arbitrary discretisation or grouping should be avoided.
2. Even if biometrical measurements are often approximately normally distributed, it seems more appropriate not to make strong distributional assumptions, or to assume a particular parametric form for the dependence on the covariate.

In order to address the second requirement we fit a theoretical distribution to the data using a finite mixture model. Mixture models provide an appealing semi-parametric structure in which to model heterogeneity and unknown distributional shapes. We refer to the monographs by Titterton, Smith, and Makov (1985) and Böhning (2000) for general background. In the present context we consider a mixture of normal distributions and model the dependence of the observations on the covariate through the weights. We allow the weights to be indexed by the covariate, so that they can vary from observation to observation.

In doing so we also meet the first requirement. We model the weights as a smooth function of the covariate using B-splines. In this way the centile curves respect the continuity of the covariate and there is no need for arbitrary grouping of the observations. We refer to Green and Silverman (1994) and Wahba (1990) for a comprehensive discussion of splines. An alternative approach to the nonparametric estimation of growth curves was taken by Cole and Green (1992); see also the other references therein. Bayesian approaches include those of Geisser (1970) and Fearn (1975).

The article is structured as follows. Section 2 presents the Bayesian hierarchical mixture model proposed for the density estimate and the calculation of the expected centiles. Section 3 introduces splines and their application to modeling the weights of the mixture. Section 4 discusses computational implementation via Markov chain Monte Carlo methods. Section 5 assesses performance of the methodology through application to a real dataset, and Section 6 concludes with general discussion and some possibilities for future work.

## 2. MIXTURE MODEL

### 2.1 NORMAL MIXTURE

Let  $\mathbf{y} = (y_i)_{i=1}^n$  be observations of a biometrical variable we want to construct growth curves for, and  $\mathbf{t} = (t_i)_{i=1}^n$  be the corresponding observed values of a continuous covariate (such as time or age). The model we assume for  $\mathbf{y}$  is

$$y_i \sim \sum_{j=1}^k w_j(t_i) \phi(\cdot; \mu_j, \sigma_j) \quad \text{independently for } i = 1, 2, \dots, n, \quad (2.1)$$

conditional on weights, means, and variances, where  $\phi(\cdot; \mu, \sigma)$  is the density of the  $N(\mu, \sigma^2)$  distribution with  $\boldsymbol{\mu} = (\mu_j)_{j=1}^k$  and  $\boldsymbol{\sigma} = (\sigma_j)_{j=1}^k$ .

The weights satisfy  $w_j(t) \geq 0$  with  $\sum_{j=1}^k w_j(t) = 1$  for all  $t$  and they are allowed to vary continuously with  $t$ . Let  $\mathbf{w}_j$  be the  $n$ -vector  $w_j(t_i)_{i=1}^n$ , for each  $j = 1, 2, \dots, k$ , with  $\mathbf{w}$  the  $k \times n$  matrix of all  $w_j(t_i)$ .

The number of components  $k$  is unknown and subject to inference, as are  $\boldsymbol{\mu}$ ,  $\boldsymbol{\sigma}$ , and  $\mathbf{w}$ . Note that making inference about  $w_j(t)$  as functions allows us to make predictions about future observations  $y$  for values of  $t$  lying between the observed  $t_i$ .

Conditional on weights, means, and variances, the  $100\alpha$ th centile can be numerically evaluated from (2.1) as that value  $C_\alpha(t_i)$  for which

$$\alpha = \sum_{j=1}^k w_j(t_i) \int_{-\infty}^{C_\alpha(t_i)} \phi(x; \mu_j, \sigma_j) dx = \sum_{j=1}^k w_j(t_i) \Phi\left(\frac{C_\alpha(t_i) - \mu_j}{\sigma_j}\right), \tag{2.2}$$

where  $\Phi(\cdot)$  is the cumulative density for a standard normal distribution.

Note that we have chosen to model the weights as varying with  $t$ , while keeping the means and variances fixed; it is possible to consider other formulations with varying means and/or variances, modeled in a similar way to our treatment of the weights in Section 3, but we have not explored these in any detail.

It is worth stressing that we are using the mixture representation primarily as a convenient semi-parametric density estimation device, and we are not greatly interested in the number of components of the mixture per se, or in a clustering of the observations.

### 2.2 LATENT ALLOCATION VARIABLES

An alternative perspective leading to the same mixture model (2.1) involves the introduction of latent allocation variables  $\mathbf{z} = (z_i)_{i=1}^n$  and the assumption that each observation  $y_i$  arose from an unknown component  $z_i$  of the mixture. The allocation variables are given probability mass function

$$p(z_i = j) = w_j(t_i) \quad \text{independently for } i = 1, 2, \dots, n, \tag{2.3}$$

and conditional on them, the observations  $y$  are independently drawn from the densities

$$y_i | \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\sigma} \sim \phi(\cdot; \mu_{z_i}, \sigma_{z_i}). \tag{2.4}$$

Integrating out  $z_i$  in (2.4) using the distribution in (2.3) leads back to (2.1).

### 2.3 PRIORS ON COMPONENT PARAMETERS

From past experience, we would not expect inference about the density in (2.1) to be highly sensitive to prior specification. As in Richardson and Green (1997), our prior assumptions are that the  $\mu_j$  and  $\sigma_j^{-2}$  are all drawn independently, with normal and gamma

priors

$$\mu_j \sim N(\xi, \kappa^{-1}) \quad \text{and} \quad \sigma_j^{-2} \sim \Gamma(\eta, \zeta),$$

where the latter is parametrized so that the mean and the variance are  $\eta/\zeta$  and  $\eta/\zeta^2$ , respectively.

The prior on the weights  $w_j$  will be discussed in Section 3.2.

The number of components  $k$  will also be considered unknown and subject to inference. For this purpose, we assume for the number of components  $k$  a uniform prior on the values  $\{1, 2, \dots, k_{\max}\}$ , where  $k_{\max}$  is a prespecified integer. As in other mixture model contexts (or indeed in almost all model choice problems), it seems difficult to argue objectively for any specific prior for  $k$ . Our choice here is for similar reasons to those in Richardson and Green (1997), namely that with this choice it is easy to adjust results to get posteriors corresponding to other priors, by importance sampling (see, e.g., Hammersley and Handscomb 1964).

In order to allow for weakly informative priors for the model parameters, we introduce a hyperprior structure and hyperparameter choices which correspond to making only minimal assumptions on the data. Following Richardson and Green (1997) we take the  $N(\xi, \kappa^{-1})$  prior for  $\mu$  to be rather flat over the range of the data, by letting  $\xi$  equal to the midpoint of this range, and  $\kappa$  equal to a small multiple of  $1/R^2$ , where  $R$  is the length of the range.

For  $\sigma^2$  we instead introduce an additional hierarchical level by allowing  $\zeta$  to follow a Gamma distribution with parameters  $f$  and  $h$ , with  $\eta > 1 > f$  and  $h$  a small multiple of  $1/R^2$ . This means that the support for  $\sigma^2$  is not fixed a priori but determined by the value sampled for  $\zeta$ .

### 3. MODELING DEPENDENCE USING SPLINES

A particular feature of our mixture model is that the weights in (2.1) are evaluated at  $t_i$ , so that they are allowed to vary from observation to observation, according to the value recorded for the covariate. In this way, we introduce dependence on the covariate through the modeling of the weights. In particular we want to reflect the fact that observations corresponding to values of the covariate  $t$  not too far from each other, are somewhat similar. This is especially true when thinking of growth curves, for which the distribution of the biometrical measurement can be assumed to vary smoothly with the covariate. In our model this is achieved through requiring the weights (and thus also the allocation probabilities) to be continuous functions of  $t$ .

The model we propose makes use of a linear combination of cubic B-splines, after a suitable transformation of the weights.

#### 3.1 MATHEMATICAL FORMULATION

The weights  $w_j(t)$  are constrained to be non-negative and sum to one for all  $t$ . For convenience, we impose the necessary constraints by transformation and, therefore, express

the weights as

$$w_j(t) = \frac{\exp(g_j(t))}{\sum_{j'=1}^k \exp(g_{j'}(t))}, \quad (3.1)$$

where  $g_j(t)$  is some continuous function of the covariate  $t$ . In order to allow the biometrical measurement to change smoothly and slowly with the covariate, we choose to model  $g_j(t)$  and, thereby, the weights  $w_j(t)$ , using natural cubic splines. For more details of the mathematical and statistical properties of cubic splines, needed in this and the following section, see Green and Silverman (1994, chaps. 2 and 3).

Consider a set of distinct real numbers  $t_1, \dots, t_n$  (such as, e.g., possible values of our covariate) on some interval  $[a, b]$ , satisfying  $a < t_1 < t_2 < \dots < t_n < b$ . A function  $g$  defined on  $[a, b]$  is a *cubic spline* on the *knots*  $t_i$  if:

1.  $g$  is a cubic polynomial on each of the intervals  $(a, t_1), (t_1, t_2), (t_2, t_3), \dots, (t_n, b)$ ;
2. the polynomial pieces fit together at the points  $t_i$  in such a way that  $g$  itself and its first and second derivatives are continuous at each  $t_i$ , and hence on the whole of  $[a, b]$ .

The continuity of the second derivative is enough to give “visual smoothness” of the resulting function. Cubic splines can be specified in many equivalent ways. One of them is to give the four polynomial coefficients of each cubic piece; for example, in the form

$$g(t) = d_i(t - t_i)^3 + c_i(t - t_i)^2 + b_i(t - t_i) + a_i \quad \text{for } t_i \leq t \leq t_{i+1}.$$

A cubic spline on an interval  $[a, b]$  is said to be a *natural cubic spline* if its second and third derivatives are zero at  $a$  and  $b$ . These conditions, called the *natural boundary conditions*, imply that  $d_0 = c_0 = d_n = c_n = 0$ , so that  $g$  is linear on the two extreme intervals  $[a, t_1]$  and  $[t_n, b]$ .

The importance of these functions derives from their variational characterization: among all functions  $g$  on  $[a, b]$  that are twice continuously differentiable and interpolate a given set of data points  $(t_i, y_i)$ , where the  $t_i$  are distinct, that minimizing the integrated-squared-second-derivative roughness penalty  $\int (g_j''(t))^2 dt$  is the unique interpolating natural cubic spline. This result motivates both the use of such splines in modeling smooth dependence, and the use of the prior specification in the following section.

Cubic splines are very convenient functions to deal with computationally; the existence of banded matrices in representations of interpolating and smoothing splines guarantees that computation times are  $O(n)$  for  $n$  data points. However, the computational requirements of our Bayesian methodology are more demanding, and it is convenient to impose a finite-dimensional structure on the problem, by restricting the choice of  $g(t)$  to the span of a prescribed set of basis functions,  $\beta_1, \dots, \beta_q$ , and, thus, considering only functions  $g(t)$  that can be expressed in the form

$$g(t) = \sum_{l=1}^q \gamma_l \beta_l(t)$$

for some numbers  $\gamma_1, \dots, \gamma_q$ .

A possible choice for the basis functions is the set of natural cubic *B-splines* on a fixed grid of knots  $s_1 < s_2 < \dots < s_q$ , usually taken to be equally spaced to cover the range of points  $t_i$ . The B-splines form a set of natural cubic splines that are non-negative and have only limited support: for  $3 \leq l \leq q - 2$  the function  $\beta_l$  is zero outside  $(s_{l-2}, s_{l+2})$ , while  $\beta_1, \beta_2, \beta_{q-1}$  and  $\beta_q$  are similar but linear outside  $(s_1, s_q)$ . Restricting  $g(t)$  to lie in the span of a set of, say, 10 B-splines typically has minimal impact on the quality of fit to the data.

Now, in our model, we have a function  $g(t)$  for each component of the mixture. Modeling them as a linear combination of B-splines on the same knots  $s_1 < s_2 < \dots < s_q$ , we can write them as

$$g_j(t) = \sum_{l=1}^q \gamma_{lj} \beta_l(t). \tag{3.2}$$

Let  $\gamma_j$  be the  $q$ -vector  $(\gamma_{lj})_{l=1}^q$ , for each  $j = 1, 2, \dots, k$ , with  $\gamma$  the  $q \times k$  matrix of all  $\gamma_{lj}$ .

### 3.2 PRIORS ON WEIGHTS

The prior on the weights is specified via that on  $g_j$  or  $\gamma_j$ , beginning with the natural integral-squared-second-derivative penalty

$$\lambda \int (g_j''(t))^2 dt = \lambda \gamma_j^T \mathbf{K} \gamma_j \tag{3.3}$$

for  $j = 1, 2, \dots, k$ , where  $\lambda$  is a parameter always positive and  $\mathbf{K}$  is the  $q \times q$  matrix with  $K_{lm} = \int \beta_l''(t) \beta_m''(t) dt$ . The prior  $p(g_j) \propto \exp\{-(1/2)\lambda \gamma_j^T \mathbf{K} \gamma_j\}$  would be “partially improper,” since the matrix  $\mathbf{K}$  has rank  $q - 2$ ; see Wahba (1978). Essentially the prior is invariant to the addition of a linear trend in  $t$ . Combined with the fact that while  $w_j$  are identifiable from the data, the  $g_j$  are not, this can cause problems with impropriety in the posterior. To circumvent these, the prior is converted to be proper by substituting the matrix  $\mathbf{K}$  with a full rank matrix, and using

$$p(g_j) \propto \exp\{-(1/2)\gamma_j^T (\lambda \mathbf{K} + \delta \mathbf{I}) \gamma_j\} \tag{3.4}$$

(where  $\mathbf{I}$  is the identity matrix and  $\delta$  is a positive parameter).

The natural integral-squared-second-derivative  $\int (g_j''(t))^2 dt$  is a measure of the roughness of the curve  $g_j(t)$ . There are many ways of measuring how rough a curve is, but this is particularly appealing for different reasons. First of all, a natural requirement for any measure of roughness is that if two functions differ only by a constant or a linear function, then their roughness should be identical. This logically leads to the idea of a roughness measure based on the second derivative of the curve under consideration. Second, the integral-squared-second-derivative has considerable computational advantages. Third, there is the connection with the variational characterization of natural cubic splines, mentioned in the previous section.

The role of  $\lambda$  is that of a smoothing parameter. As  $\lambda$  increases toward  $\infty$ , there is a stronger shrinkage of each  $\gamma_j$  towards a linear trend. As a result the curves  $g_j(t)$  become

smoother and so do the weights. Overall, the centile growth curves will also be very smooth and display little curvature. In the opposite limiting case, as  $\lambda \rightarrow 0$  the centile growth curves will track the empirical ones very closely at the expense of being rather variable. This variability would affect particularly the more extreme centile curves and the accuracy of their estimates, as the centile standard errors increase steeply towards the tail of the distribution.

In this perspective, the model chosen for the weights introduces, through the natural integral-squared-second-derivative  $\int (g_j''(t))^2 dt$  and the smoothing parameter  $\lambda$ , a roughness penalizing element and a trade-off between smoothness and goodness of fit of the centile curves.

The parameter  $\delta$  can also be regarded as a smoothing parameter. As its value increases, not only there is a stronger shrinkage of the  $\gamma_j$  towards zero, but this effect is also reinforced by the fact that the variance for the  $\gamma_j$  is reduced and their distribution is more concentrated around the zero-mean.

Obviously the choice of the parameters  $\lambda$  and  $\delta$  is of some importance. This matter will be discussed later in Section 5.3.

### 3.3 COMPLETE HIERARCHICAL MODEL

The joint distribution of all variables conditional on fixed hyperparameters may be written

$$\begin{aligned} p(k, \gamma, \boldsymbol{\mu}, \zeta, \boldsymbol{\sigma}, \mathbf{z}, \mathbf{y} | \lambda, \delta, \xi, \kappa, \eta, f, h) \\ = p(k) p(\gamma | k, \lambda, \delta) p(\boldsymbol{\mu} | k, \xi, \kappa) p(\zeta | f, h) p(\boldsymbol{\sigma} | k, \eta, \zeta) p(\mathbf{z} | \gamma, k) p(\mathbf{y} | \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\sigma}). \end{aligned}$$

We have

$$p(\mathbf{z} | \gamma, k) = \prod_{i=1}^n w_{z_i}(t_i),$$

with the relationship between  $\mathbf{w}$  and  $\gamma$  given by (3.1), and

$$p(\mathbf{y} | \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \prod_{i=1}^n \phi(y_i; \mu_{z_i}, \sigma_{z_i}).$$

The prior distribution  $p(\gamma | k, \lambda, \delta)$  is given in Section 3.2, while  $p(\boldsymbol{\mu} | k, \xi, \kappa)$ ,  $p(\zeta | f, h)$ , and  $p(\boldsymbol{\sigma} | k, \eta, \zeta)$  are given in Section 2.3. The complete hierarchical model is displayed in Figure 1 as a directed acyclic graph (DAG). We follow the usual convention that square boxes represent fixed or observed quantities and circles represent the unknowns. Relationships that are deterministic, as opposed to stochastic, are indicated by broken lines.

## 4. COMPUTATIONAL IMPLEMENTATION

The complexity of the mixture model presented requires Markov chain Monte Carlo (MCMC) methods to approximate the posterior distribution. Details of these computational

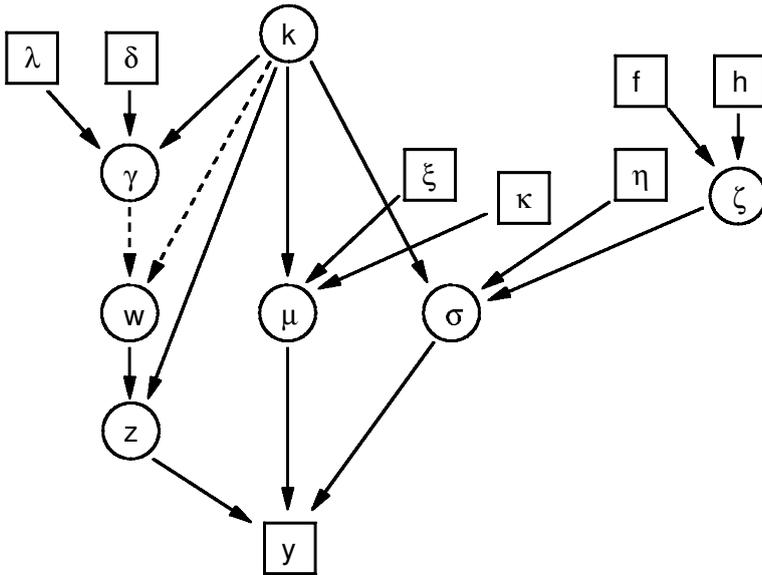


Figure 1. Directed acyclic graph for the complete hierarchical model.

methods can be found, for example, in Tierney (1994) and Besag, Green, Higdon, and Mengersen (1995).

Our sampler uses five different fixed-dimension moves, each updating one of the variables of the model, plus a variable dimension move for updating  $k$ . The way in which the first five moves are performed is quite standard and thus we will go through them rather quickly. The last move, for updating  $k$ , is performed using a reversible jump method (Green 1995).

#### 4.1 RANDOM WALK METROPOLIS MOVE FOR THE WEIGHTS

We can update the weights by means of a simultaneous random walk Metropolis method applied to  $\gamma$ . Thus, we draw

$$\gamma'_{ij} \sim N(\gamma_{ij}, \tau^2)$$

independently, compute the corresponding weights  $w'_j(t_i)$ , and accept this proposal with the usual probability equal to  $\min\{1, Q\}$  where

$$Q = \frac{p(k, \gamma', \mu, \zeta, \sigma, \mathbf{z}, \mathbf{y})}{p(k, \gamma, \mu, \zeta, \sigma, \mathbf{z}, \mathbf{y})} = \frac{p(\gamma'|k) p(\mathbf{z}|\gamma', k)}{p(\gamma|k) p(\mathbf{z}|\gamma, k)},$$

which simplifies into

$$Q = \exp \left( - (1/2) \sum_{j=1}^k [(\gamma'_j)^T (\lambda \mathbf{K} + \delta \mathbf{I}) \gamma'_j - \gamma_j^T (\lambda \mathbf{K} + \delta \mathbf{I}) \gamma_j] \right) \prod_{i=1}^n \frac{w'_{z_i}(t_i)}{w_{z_i}(t_i)}.$$

Faster convergence in the algorithm was obtained introducing a small bias in the proposal. Thus, instead of proposing to update  $\gamma_{lj}$  to the new value

$$\gamma'_{lj} = \gamma_{lj} + r,$$

where  $r$  is a random number from a  $N(0, \tau^2)$ , we propose, as a new value for  $\gamma_{lj}$ ,

$$\gamma'_{lj} = \gamma_{lj} + r + \varrho \sum_{i:z(i)=j} \beta_l(t_i),$$

where  $\varrho$  is a small number (we used  $\varrho = 0.001$  on the basis of some limited pilot runs). The acceptance ratio for this proposal then becomes

$$Q' = \frac{p(\gamma'|k) p(\mathbf{z}|\gamma', k) p(\gamma'|\gamma)}{p(\gamma|k) p(\mathbf{z}|\gamma, k) p(\gamma|\gamma')} = Q \exp \left( -\frac{2\varrho}{\tau^2} \sum_{j=1}^k \sum_{l=1}^q (\gamma'_{lj} - \gamma_{lj}) \sum_{i:z(i)=j} \beta_l(t_i) \right),$$

where  $p(\gamma'|\gamma)$  is the probability of proposing  $\gamma'$  when the current value is  $\gamma$  and  $p(\gamma|\gamma')$  is the probability of proposing  $\gamma$  when the current value is  $\gamma'$ .

The biased proposal has the overall effect of proposing new weights which tend slightly to favor the current allocation of the observations.

## 4.2 GIBBS MOVE FOR THE ALLOCATIONS

For the allocations we have

$$p(\mathbf{z}|\gamma, \boldsymbol{\mu}, \boldsymbol{\sigma}, k, \mathbf{y}) \propto p(\gamma, \boldsymbol{\mu}, \boldsymbol{\sigma}, \mathbf{z}, k, \mathbf{y}) \propto p(\mathbf{z}|\gamma, k) p(\mathbf{y}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\sigma})$$

(i.e., proportional as a function of  $\mathbf{z}$ ), so the allocation variable  $z_i$  has conditional probability

$$p(z_i = j|\gamma, \boldsymbol{\mu}, \boldsymbol{\sigma}, k, \mathbf{y}) = \frac{w_j(t_i) \phi(y_i; \mu_j, \sigma_j)}{\sum_{j'} w_{j'}(t_i) \phi(y_i; \mu_{j'}, \sigma_{j'})}. \tag{4.1}$$

We can sample directly from this distribution and update the allocation variables independently by means of Gibbs sampling.

## 4.3 MOVES FOR THE PARAMETERS AND THE HYPERPARAMETER

### 4.3.1 Updating $\mu$

Before considering the updating of the  $\mu_j$ , we comment briefly on the issue of labeling the components. The whole model is, in fact, invariant to permutation of the labels  $j = 1, 2, \dots, k$ . For identifiability, Richardson and Green (1997) adopted a unique labeling in which the  $\mu_j$  are in increasing numerical order. As a consequence the joint prior distribution of the  $\mu_j$  is  $k!$  times the product of the individual normal densities, restricted to the set  $\mu_1 < \mu_2 < \dots < \mu_k$ .

The  $\mu_j$  can be updated by means of Gibbs sampler, drawing them independently from the distribution

$$\mu_j | \dots \sim N \left( \frac{\sigma_j^{-2} \sum_{i:z_i=j} y_i + \kappa \xi}{\sigma_j^{-2} n_j + \kappa}, (\sigma_j^{-2} n_j + \kappa)^{-1} \right),$$

where  $n_j = \#\{i : z_i = j\}$  is the number of observations currently allocated to the  $j$  component of the mixture. Here and later, “ $\dots$ ” denotes “all other variables.” In order to preserve the ordering constraints on the  $\mu_j$ , the move is accepted provided the ordering is unchanged and rejected otherwise.

Except for very small values of  $k$ , this updating move has the drawback of producing a very small acceptance ratio, due to the fact that the ordering of the  $\mu_j$  seldom remains unchanged. For this reason it is preferable to update  $\mu_j$  using a trick similar to the one that Green and Richardson (2002) adopted to update their component risk parameters. We propose simultaneous independent zero-mean normal increments to each  $\mu_j$ ; the modified values of  $\mu_j$  are then placed in increasing order to give  $\mu'$  say. The complete set of updates is accepted with probability, formed from prior ratio and likelihood ratio, which reduces to  $\min\{1, S\}$  where

$$S = \exp \left\{ \sum_{j=1}^k \left[ -\frac{\kappa}{2} ((\mu_j'^2 - \mu_j^2) - 2\xi(\mu_j' - \mu_j)) - \sum_{i:z_i=j} \frac{1}{2\sigma_{z_i}^2} ((\mu_{z_i}'^2 - \mu_{z_i}^2) - 2y_i(\mu_{z_i}' - \mu_{z_i})) \right] \right\}.$$

An alternative to imposing identifiability constraints on the parameters, is to order the parameters, according to some unique labeling, a posteriori, after the whole sample of parameters drawn from the posterior is available. In the present case, where the main concern is the inference on the posterior density of the data and its centiles, rather than on the single parameters of the model, this labeling is not even required. In this case Gibbs sampler can be used for updating  $\mu_j$  without any need for their order to stay unchanged.

For reason of completeness and because of some interest, however, we preferred to make inference also on the single parameters of the model and we decided to use the first approach (i.e., to impose an ordering on the  $\mu_j$  a priori) after having obtained much the same results from both of them.

### 4.3.2 Updating $\sigma$

The full conditionals for  $\sigma_j^2$  are

$$\sigma_j^{-2} | \dots \sim \Gamma \left( \eta + \frac{1}{2} n_j, \zeta + \frac{1}{2} \sum_{i:z_i=j} (y_i - \mu_j)^2 \right).$$

We update  $\sigma_j^2$  independently using a Gibbs move, sampling from their full conditionals.

### 4.3.3 Updating $\zeta$

The only hyperparameter we are not treating as fixed is  $\zeta$ . Conditional on all the other parameters and the data,  $\zeta$  has a Gamma distribution

$$\zeta | \dots \sim \Gamma \left( f + k\eta, h + \sum_{j=1}^k \sigma_j^{-2} \right).$$

We update  $\zeta$  by a Gibbs move, sampling from its full conditional.

## 4.4 VARIABLE DIMENSION MOVE FOR UPDATING $k$

Updating the value of  $k$  implies a change of dimensionality for the components  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$ , the allocation variables  $\mathbf{z}$  and the weights  $\mathbf{w}$  (through the change of dimensionality for  $g$  and  $\gamma$ ). We follow the approach used by Richardson and Green (1997) consisting of a random choice between splitting an existing component into two, and merging two existing components into one. The probabilities of these alternatives are  $b_k$  and  $d_k = 1 - b_k$ , respectively, when there are currently  $k$  components. Of course,  $d_1 = 0$  and  $b_{k_{\max}} = 0$ , and otherwise we choose  $b_k = d_k = 0.5$ , for  $k = 2, 3, \dots, k_{\max} - 1$ .

For the combine proposal we randomly choose a pair of components  $(j_1, j_2)$  that are adjacent in terms of the current value of their means, which means  $\mu_{j_1} < \mu_{j_2}$ , with no other  $\mu_j$  in the interval  $[\mu_{j_1}, \mu_{j_2}]$ . These two components are merged into a new one, labeled  $\mu_{j^*}$ , reducing  $k$  by 1. We then reallocate all those observations  $y_i$  with  $z_i = j_1$  or  $j_2$  to the new component  $j^*$  and create values for  $\mu_{j^*}, \sigma_{j^*}$  in such a way that:

$$\begin{aligned} \mu_{j^*} &= (\mu_{j_1} + \mu_{j_2})/2 \\ \mu_{j^*}^2 + \sigma_{j^*}^2 &= [(\mu_{j_1}^2 + \sigma_{j_1}^2) + (\mu_{j_2}^2 + \sigma_{j_2}^2)]/2. \end{aligned}$$

To create the new values  $w_{j^*}(t_i)$  we first have to create  $\gamma_{lj^*}$ , for  $l = 1, 2, \dots, q$ . We do this by setting

$$\gamma_{lj^*} = \log[\exp(\gamma_{lj_1}) + \exp(\gamma_{lj_2})], \quad \text{for } l = 1, 2, \dots, q$$

and calculate  $g_j^*(t_i)$  and  $w_j^*(t_i)$  using (3.2) and (3.1), respectively. Note that in this way all the weights change slightly.

The split proposal starts by choosing a component  $j^*$  at random. This component is split into two new ones labeled  $j_1$  and  $j_2$ , augmenting  $k$  by 1. Then we have to reallocate all those observations  $y_i$  with  $z_i = j^*$  between the two new components, and create values for  $(w_{j_1}, w_{j_2}, \mu_{j_1}, \mu_{j_2}, \sigma_{j_1}, \sigma_{j_2})$ . Let us start by splitting  $\mu_{j^*}$  and  $\sigma_{j^*}$ . We generate a two-dimensional random vector  $\mathbf{v}$  to specify the new parameters. We use Beta distributions

$$v_1 \sim \text{Be}(1, 1) \quad \text{and} \quad v_2 \sim \text{Be}(1, 1)$$

for this and set

$$\begin{aligned} \mu_{j_1} &= \mu_{j^*} - v_1\sigma_{j^*}, \\ \mu_{j_2} &= \mu_{j^*} + v_1\sigma_{j^*}, \\ \sigma_{j_1}^2 &= 2v_2(1 - v_1^2)\sigma_{j^*}^2, \end{aligned}$$

and

$$\sigma_{j_2}^2 = 2(1 - v_2)(1 - v_1^2)\sigma_{j^*}^2.$$

We then need to check that the constraints on the means are satisfied. If not the move is rejected forthwith, as the misordered vector  $\boldsymbol{\mu}$  has zero density under the ordered prior. If the constraints are satisfied we move on and split the weights.

In doing so, we generate a  $q$ -dimensional vector  $\mathbf{u}$  from

$$u_l \sim \text{Be}(0.5, 0.5), \quad \text{independently for } l = 1, 2, \dots, q,$$

and we set

$$\gamma_{l_{j_1}} = \gamma_{l_{j^*}} + \log(u_l), \quad \gamma_{l_{j_2}} = \gamma_{l_{j^*}} + \log(1 - u_l) \quad \text{for } l = 1, 2, \dots, q.$$

We then calculate  $g_{j_1}(t_i)$ ,  $g_{j_2}(t_i)$ , and  $w_{j_1}(t_i)$ ,  $w_{j_2}(t_i)$  using (3.2) and (3.1), respectively. We denote the proposed new weights by  $\mathbf{w}'$ .

Finally we reallocate those  $y_i$  with  $z_i = j^*$  between  $j_1$  and  $j_2$  in a way analogous to the standard Gibbs allocation move; see Equation (4.1). We denote the proposed new allocation vector by  $\mathbf{z}'$ .

According to the reversible jump framework, the acceptance probability for the split move is  $\min(1, A)$ , where

$$\begin{aligned} A &= (\text{likelihood ratio}) \times \frac{p(k+1)}{p(k)} \times (k+1) \times \prod_{i=1}^n \frac{w'_{z'_i}(t_i)}{w_{z_i}(t_i)} \\ &\times \sqrt{\frac{\kappa}{2\pi}} \exp \left\{ -\frac{\kappa}{2} [(\mu_{j_1} - \xi)^2 + (\mu_{j_2} - \xi)^2 - (\mu_{j^*} - \xi)^2] \right\} \\ &\times \frac{\zeta^\eta}{\Gamma(\eta)} \exp \left[ -\zeta \left( \frac{1}{\sigma_{j_1}^2} + \frac{1}{\sigma_{j_2}^2} - \frac{1}{\sigma_{j^*}^2} \right) \right] \left( \frac{\sigma_{j^*}}{\sigma_{j_1}\sigma_{j_2}} \right)^{2(\alpha-1)} \\ &\times \frac{|\lambda\mathbf{K} + \delta\mathbf{I}|^{1/2}}{(2\pi)^{q/2}} \exp \left\{ -\frac{1}{2} \left[ \gamma_{j_1}^T (\lambda\mathbf{K} + \delta\mathbf{I}) \gamma_{j_1} + \gamma_{j_2}^T (\lambda\mathbf{K} + \delta\mathbf{I}) \gamma_{j_2} \right. \right. \\ &\quad \left. \left. - \gamma_{j^*}^T (\lambda\mathbf{K} + \delta\mathbf{I}) \gamma_{j^*} \right] \right\} \\ &\times \frac{d_{k+1}}{b_k P_{\text{alloc}}} \times \left[ be_{1,1}(v_1) be_{1,1}(v_2) \prod_{l=1}^q be_{0.5,0.5}(u_l) \right]^{-1} \\ &\times \frac{8(1 - v_1^2)\sigma_{j^*}^3}{\prod_{l=1}^q u_l(1 - u_l)}, \end{aligned} \tag{4.2}$$

where  $\Gamma(\cdot)$  is the Gamma function,  $P_{\text{alloc}}$  is the probability of this particular allocation,  $be_{p,r}$  denotes the Beta( $p, r$ ) density, and (likelihood ratio) is the ratio of the product of

the  $f(y_i|z_i, \mu_{z_i}, \sigma_{z_i})$  terms for the new parameter set to that for the old one. The quantity  $|\lambda\mathbf{K} + \delta\mathbf{I}|$  represents the determinant of the inverse of the covariance matrix for each of the  $g_j$ . The need to make the prior distribution on the  $g_j$  proper and to introduce the positive parameter  $\delta$  is now evident.

The first five lines of (4.2) are the product of the likelihood ratio and the prior ratio for the parameters of the model. The sixth line is the proposal ratio. The last line is the Jacobian of the transformation from the vector  $(\mu_{j^*}, \sigma_{j^*}, \gamma_{1j^*}, \dots, \gamma_{qj^*}, v_1, v_2, u_1, \dots, u_q)$  to the vector  $(\mu_{j_1}, \sigma_{j_1}, \gamma_{1j_1}, \dots, \gamma_{qj_1}, \mu_{j_2}, \sigma_{j_2}, \gamma_{1j_2}, \dots, \gamma_{qj_2})$ .

The acceptance probability for the combine move is  $\min(1, A^{-1})$ , with some obvious substitutions in the expression for  $A$ .

### 4.5 WITHIN-MODEL SIMULATION AND PATH SAMPLING

The alternative general approach to sample-based joint inference about a model indicator  $k$  and model parameters is to conduct separate simulation within each model, and piece the results together afterwards. Inference about  $k$  is then based on the estimate of the posterior model probabilities

$$p(k|\mathbf{y}) \propto p(k)p(\mathbf{y}|k).$$

These require estimates of the marginal likelihoods  $p(\mathbf{y}|k)$  separately for each  $k$ , using individual MCMC runs. Let  $\psi_k = (\boldsymbol{\mu}, \boldsymbol{\sigma}, \mathbf{w})$  denote the unknowns  $\boldsymbol{\mu}, \boldsymbol{\sigma}$ , and  $\mathbf{w}$  in the model with  $k$  components. Then, possible estimates of marginal likelihoods, using importance sampling are, for example, (see Newton and Raftery 1994)

$$\hat{p}_1(\mathbf{y}|k) = M \left/ \sum_{m=1}^M \left\{ p(\mathbf{y}|k, \boldsymbol{\psi}_k^{(m)}) \right\}^{-1} \right.,$$

based on an MCMC sample  $\boldsymbol{\psi}_k^{(1)}, \boldsymbol{\psi}_k^{(2)}, \dots, \boldsymbol{\psi}_k^{(M)}$  from the posterior  $p(\boldsymbol{\psi}_k|\mathbf{y}, k)$ ,

$$\hat{p}_2(\mathbf{y}|k) = M^{-1} \sum_{m=1}^M p(\mathbf{y}|k, \boldsymbol{\psi}_k^{(m)}),$$

based on a sample from the prior  $p(\boldsymbol{\psi}_k|k)$ , or

$$\hat{p}_3(\mathbf{y}|k) = \frac{\sum_{m=1}^M p(\mathbf{y}|k, \boldsymbol{\psi}_k^{(m)}) / \left\{ \epsilon \hat{p}_3(\mathbf{y}|k) + (1 - \epsilon) p(\mathbf{y}|k, \boldsymbol{\psi}_k^{(m)}) \right\}}{\sum_{m=1}^M \left\{ \epsilon \hat{p}_3(\mathbf{y}|k) + (1 - \epsilon) p(\mathbf{y}|k, \boldsymbol{\psi}_k^{(m)}) \right\}^{-1}},$$

based on a sample from a  $(\epsilon, 1 - \epsilon)$  mixture of the prior  $p(\boldsymbol{\psi}_k|k)$  and the posterior  $p(\boldsymbol{\psi}_k|\mathbf{y}, k)$  and supposed to perform better than the previous two. It is well known, however, that when the distance between the prior and the posterior densities is big (as it is in our case), the variability of this last estimator can become so large that the estimate is virtually unusable. A gain of efficiency can be obtained using the idea of the bridge sampling and choosing a sensible density which serves as a ‘‘bridge’’ between the prior and the posterior densities.

The idea of creating a bridge can obviously be pushed further if the prior and the posterior densities are so far separated that the estimator based on the bridge sampling is too variable to use in practice. In such cases it is useful to construct an infinite series of intermediate densities—that is, a whole “path”—from which we can make draws. We refer to Gelman and Meng (1998) for a detailed discussion of bridge sampling and path sampling.

Briefly, to estimate the marginal likelihood, we construct a geometric path between the prior and the posterior parameter densities using a scalar parameter  $\theta \in [0, 1]$ ,

$$q(\boldsymbol{\psi}_k|\theta) = \{p(\boldsymbol{\psi}_k|k)\}^{1-\theta}\{p(\boldsymbol{\psi}_k|k)p(\mathbf{y}|\boldsymbol{\psi}_k, k)\}^\theta = p(\boldsymbol{\psi}_k|k)\{p(\mathbf{y}|\boldsymbol{\psi}_k, k)\}^\theta. \quad (4.3)$$

We use a notation in which  $q(\cdot)$  represents an unnormalized density,  $p(\cdot)$  is the corresponding normalized density and  $c(\cdot)$  is the normalizing constant. From (4.3) it is evident that estimating  $p(\mathbf{y}|k)$  is equivalent to estimating  $c(1)$ , that is, the normalizing constant when  $\theta = 1$ . Following Gelman and Meng (1998) it can be proved that

$$\log[p(\mathbf{y}|k)] = \log[c(1)] = \int_0^1 E_\theta[U(\boldsymbol{\psi}_k, \theta)] d\theta, \quad (4.4)$$

where  $E_\theta$  denotes the expectation with respect to the sampling distribution  $p(\boldsymbol{\psi}_k|\theta)$  (which is the unknown normalized version of  $q(\boldsymbol{\psi}_k|\theta)$ ) and where

$$U(\boldsymbol{\psi}_k, \theta) = \frac{d}{d\theta} \log q(\boldsymbol{\psi}_k|\theta).$$

For estimating  $\log[p(\mathbf{y}|k)]$  we then numerically evaluated the integral in (4.4) over a grid of values for  $\theta$ . We chose an exponential spacing for the grid to account for the sharp variation of  $q(\boldsymbol{\psi}_k|\theta)$  for  $\theta$  close to 0. Given the indexing  $0 = \theta_{(1)} < \dots < \theta_{(h)} < \dots < \theta_{(H)} = 1$  and a (possibly dependent) sample of draws  $(\boldsymbol{\psi}_k^{(m)}, \theta^{(m)})$  from  $p(\boldsymbol{\psi}_k, \theta)$ , applying the trapezoidal rule, we estimate  $\log[p(\mathbf{y}|k)]$  by

$$\log[\hat{p}(\mathbf{y}|k)] = \log \hat{c}(1) = \frac{1}{2} \sum_{h=1}^{H-1} (\theta_{(h+1)} - \theta_{(h)}) (\bar{U}_{(h+1)} + \bar{U}_{(h)}), \quad (4.5)$$

where  $\bar{U}_{(h)}$  is the average of the values of  $U(\boldsymbol{\psi}_k^{(m)}, \theta^{(m)})$  for all simulation draws  $m$  for which  $\theta^{(m)} = \theta_{(h)}$ .

We now come to the matter of obtaining an MCMC sample of draws  $(\boldsymbol{\psi}_k^{(m)}, \theta^{(m)})$  from the joint distribution  $p(\boldsymbol{\psi}_k, \theta)$ , using a different MCMC sample from that used to compute the within-model posterior distribution of the parameters. We can easily draw  $\boldsymbol{\psi}_k$  from  $p(\boldsymbol{\psi}_k|\theta)$ , but the problem of updating  $\theta$  is more complicated. Assuming a discrete uniform distribution for  $\theta$  on the values specified in the grid, we have

$$p(\theta|\boldsymbol{\psi}_k) \propto p(\theta)p(\boldsymbol{\psi}_k|\theta) = p(\theta) \frac{q(\boldsymbol{\psi}_k|\theta)}{c(\theta)} \propto \frac{q(\boldsymbol{\psi}_k|\theta)}{c(\theta)},$$

and  $c(\theta)$  is unknown.

A first solution would be to use nested loops: for each value  $\theta$  in the grid, run an iterative simulation algorithm until approximate convergence to obtain  $p(\boldsymbol{\psi}_k|\theta)$ . The problem with

this solution is that when the grid for  $\theta$  is very fine, as in the present case, many runs of the algorithm are required and convergence needs to be assessed for each of them.

We thus preferred to draw from the joint density of  $(\psi_k, \theta)$ , combining the simulations of  $\psi_k$  and  $\theta$  in a single loop of iterative simulation, alternately updating  $\psi_k$  and  $\theta$ . In this way, convergence can be checked looking at the sampled distribution for  $\theta$  which should be uniform on the values specified in the grid. Also, we avoid the need to do burn-in runs for each  $\theta$  in the grid, separately. To overcome the problem of drawing  $\theta$  from  $p(\theta|\psi_k)$  we used a rough estimate  $c^*(\theta)$  instead of  $c(\theta)$  and simulated  $\theta$  with target

$$p(\theta) \frac{q(\psi_k|\theta)}{c^*(\theta)} = \left( p(\theta) \frac{c(\theta)}{c^*(\theta)} \right) p(\psi_k|\theta),$$

so that the output distribution for  $\theta$  is only slightly altered (provided  $c^*(\theta)$  is not too far from  $c(\theta)$ ) while the output distribution for  $(\psi_k|\theta)$  is unaltered. We proceeded starting from a first rough estimate  $c^*(\theta)$ . This is obtained keeping  $\theta = \theta_{(h')}$  fixed for 20 sweeps of the algorithm, updating only  $\psi_k$ , estimating  $c^*(\theta_{(h')})$  from these 20 sweeps using (4.5) where the sum is over  $h = 1, \dots, h' - 1$ , and then moving to  $\theta = \theta_{(h'+1)}$ . Notice that there is no need to estimate  $c^*(0)$  since  $c(0) = 1$ . After  $c^*(\theta)$  was calculated for each  $\theta$  a further 100,000 sweeps were used to draw values of  $\theta$  and  $\psi_k$  from their joint distribution and to re-estimate  $c^*(\theta)$  after an exponentially increasing number of sweeps ( $m = 67, 91, 168, \dots, 73,783, 100,000$ ). We then used the last estimate of  $c^*(\theta)$  to run the final 100,000 iterations to be used to estimate  $p(\mathbf{y}|k)$  from (4.5). The whole procedure is quite cumbersome. It is however required to get values of  $c^*(\theta)$  as close as possible to those of  $c(\theta)$ . This is necessary to obtain an output distribution of  $\theta$  not too far from the uniform and to ensure, in this way, that the averages  $\bar{U}_{(h)}$  are calculated over a considerable number of values for each  $h$ .

## 5. RESULTS

The method for determining growth curves illustrated so far is applied to data on triceps skinfold in Gambian females. These data come from an anthropometry survey of 892 girls and women up to age 50 in three Gambian villages, seen during the dry season of 1989; 620 (70%) of the subjects were aged under 20. Five different anthropometric measurements were taken. Here we discuss only the triceps skinfold for a subset of the original data including 584 subjects aged from about 3 to 26 years. In this covariate range, the observations on triceps skinfold vary in the interval  $[3.2; 21.0]$ , determining a length  $R$  for the range of the data equal to 17.8. This dataset was already analyzed by Cole and Green (1992) and Green and Silverman (1994, sec. 6.7).

The results reported correspond to runs of 100,000 sweeps after a burn-in period of 50,000 sweeps. The following settings were used for previously unspecified constants:  $\kappa = 10.0/R^2$ ,  $\eta = 2.0$ ,  $f = 0.2$ ,  $h = 10.0/R^2$ ,  $\lambda = 0.7$ ,  $\delta = 0.09$ . Some experimentation indicates that results are quite robust to reasonable perturbations of these values.

Models with a number of components up to  $k_{\max} = 5$  were considered. The number of knots was fixed at 10, with the knots equally spaced between 5 and 23.

## 5.1 POSTERIOR INFERENCE ON $k$

The reversible jump approach allows joint (or across-model) inference about the number of components  $k$  and the other parameters of the model. The posterior distribution  $p(k|\mathbf{y})$  can be, therefore, directly obtained from the MCMC sample. Unfortunately, in the present case, a few runs of the algorithm showed an extremely small acceptance rate of the split/combine move, approximately equal to 0.06 %. This is partly due to the large number of parameters involved in the model, that makes it difficult to move from a state of dimension  $12k$  to a new one with dimension  $12(k - 1)$  or  $12(k + 1)$ . However, the main reason for such bad performance is the large size of the dataset. This leads to very peaked posterior distributions for the parameters, causing a very low acceptance rate of the split/combine proposals.

An improvement in the acceptance rate was found considering one quarter of the original dataset: the data were sorted according to the value of  $t_i$  and a reduced dataset was created taking one datum at random from every consecutive four, in order to preserve the original structure of the data. In this way the acceptance rate increases 10 times and a further increase was obtained repeating the split/combine move more than once in each sweep and considering the move accepted overall, if accepted at least once.

Reversible jump was therefore used on the reduced dataset to simulate a “partial” posterior distribution for  $k$ . The mode of this distribution was considered as an estimate of the number of components for the mixture. The MCMC algorithm was then run again to obtain the density estimate for the whole dataset, skipping this time the variable dimension move and fixing  $k$  to its estimated value. This is not “using the data twice,” but is a valid approach to approximating the posterior of all parameters conditional on  $k = k^*$ , where  $k^*$  maximizes  $p(k|\mathbf{y})$ ; the probability that  $k^*$  differs between the reduced and full datasets is negligible.

In addition to the across-model approach on the partial dataset, we also tried the within-model approach of Section 4.5. We mainly did so in order to check the estimate of  $k$  we obtained from the reduced dataset and also to explore an alternative solution when the amount of data makes the use of reversible jump on the whole dataset infeasible, at least with the moves we have considered.

Even proceeding in this way, the estimate still shows some variability. In order to reduce this variability, the final estimate for each model was therefore obtained as a trimmed mean of four estimates (i.e., the mean of the middle two of the four) resulting from four different runs of the algorithm. Figure 2 shows the posterior distribution  $p(k|\mathbf{y})$  obtained in this way for the whole dataset (continuous line) and for the partial dataset (dashed line) compared with the one resulting from the reversible jump sample on the partial dataset (dotted line). In spite of their variability, the different estimates agree on the mode of the posterior distribution for  $k$  and favor an explanation of the data using four components.

A further visual comparison to evaluate the goodness-of-fit of the different models can be based on the plot of the cumulative density of the data against the covariate. If the data are well fitted by the model, the points must be uniformly spread (in the vertical direction)

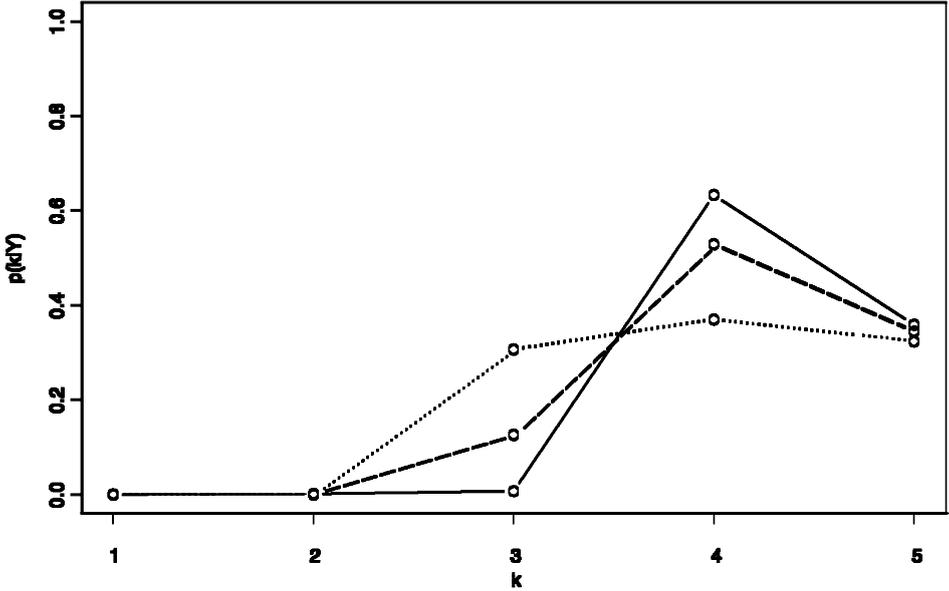


Figure 2. Posterior distribution  $p(k|\mathbf{y})$  estimated using path sampling on the whole dataset (continuous line) and on the partial dataset (dashed line), and using reversible jump on the partial dataset (dotted line).

over the plot. The cumulative density, conditional on weights, means and variances can be written as

$$F(y_i) = \sum_{j=1}^k w_j(t_i) \int_{-\infty}^{y_i} \phi(x; \mu_j, \sigma_j) dx = \sum_{j=1}^k w_j(t_i) \Phi\left(\frac{y_i - \mu_j}{\sigma_j}\right).$$

The cumulative density for all  $t_i$  was computed at each sweep and then averaged over the number of sweeps. Figure 3 shows these cumulative densities plotted against the age for the five different models. Clearly the models with one and two components show their limits in fitting the data, while the models with three to five components seem to be all reasonable. In the bottom right corner, the fit corresponding to the across-model inference is also shown.

## 5.2 POSTERIOR INFERENCE FOR CENTILE CURVES

We briefly illustrate how to evaluate the centile curves from the MCMC output. First of all, centile curves were not evaluated for every  $t_i$  but only over a grid of 47 equally spaced values for  $t$ , between 4.0 and 26.0. We indicate these values as  $t^* = (t_i^*)_{i=1}^{n^*}$ , where  $n^* = 47$ . This was done to reduce the computing time required to evaluate the centiles numerically at each sweep of the algorithm.

Then we set a grid of values  $\mathbf{a} = (a_h)_{h=1}^{200}$ . For each sweep  $m = 1, 2, \dots, 100,000$  and for each value  $t_i^*$ , and recalling (2.2), we evaluate, for a given value  $\alpha$ , the centile  $C_\alpha^{(m)}(t_i^*)$

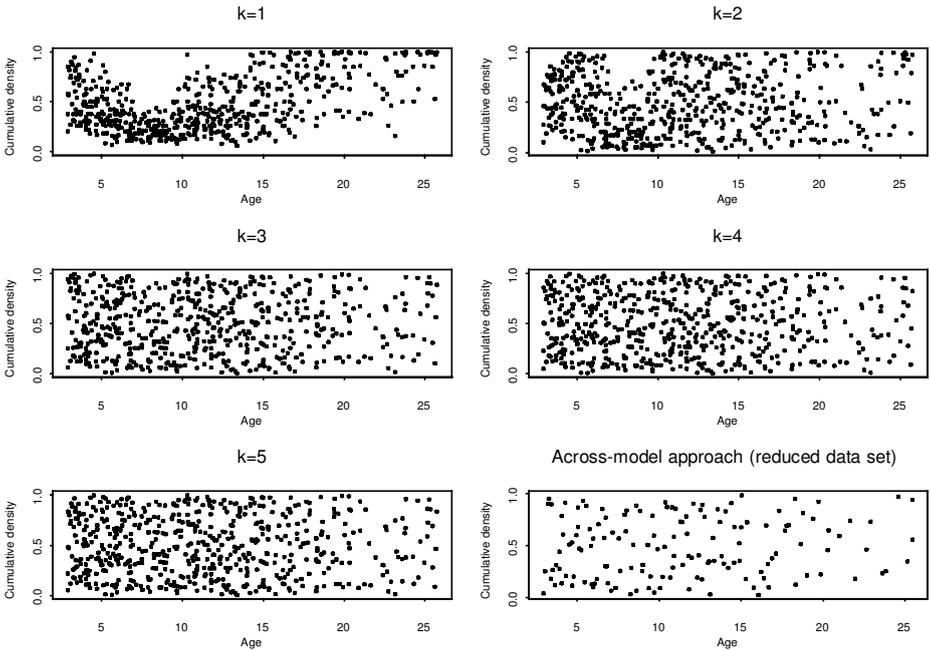


Figure 3. Cumulative density for  $y_i$  plotted against  $t_i$ .

as that value  $a_h$  for which

$$\sum_{j=1}^k w_j^{(m)}(t_i^*) \Phi \left( \frac{a_h - \mu_j^{(m)}}{\sigma_j^{(m)}} \right) \leq \alpha < \sum_{j=1}^k w_j^{(m)}(t_i^*) \Phi \left( \frac{a_{h+1} - \mu_j^{(m)}}{\sigma_j^{(m)}} \right).$$

The centiles  $C_\alpha^{(m)}(t_i^*)$  are then averaged over the different sweeps to obtain finally  $C_\alpha(t_i^*)$ . Figure 4 shows the centile curves obtained for the five different models. The results for the model with 1 component are given only for reasons of completeness and as a visual aid, to see how the fitting of the data changes as the number of components increases. It expresses no dependence of the data on the covariate and simply models the data as a normal distribution with some estimated mean and variance. Obviously it has no pretension of explaining the data. The two components mixture seems to be not completely adequate to model the density of the data, either, while mixtures with three to five components show similar results in terms of growth curves. In the bottom right corner of Figure 4, the centile curves estimated for the partial dataset, using the across-model approach are also given. In this case inference results from mixing over  $k$ , while inference on the other centile curves is conditional upon  $k$ .

Data are characterized by high triceps values in early childhood, followed by a fall and then a second continuous rise until the end of the age interval. The centile curves obtained using the three models with three to five components closely follow this pattern. They show the same “notch” in the dependence at around nine years, found in Cole and Green (1992).

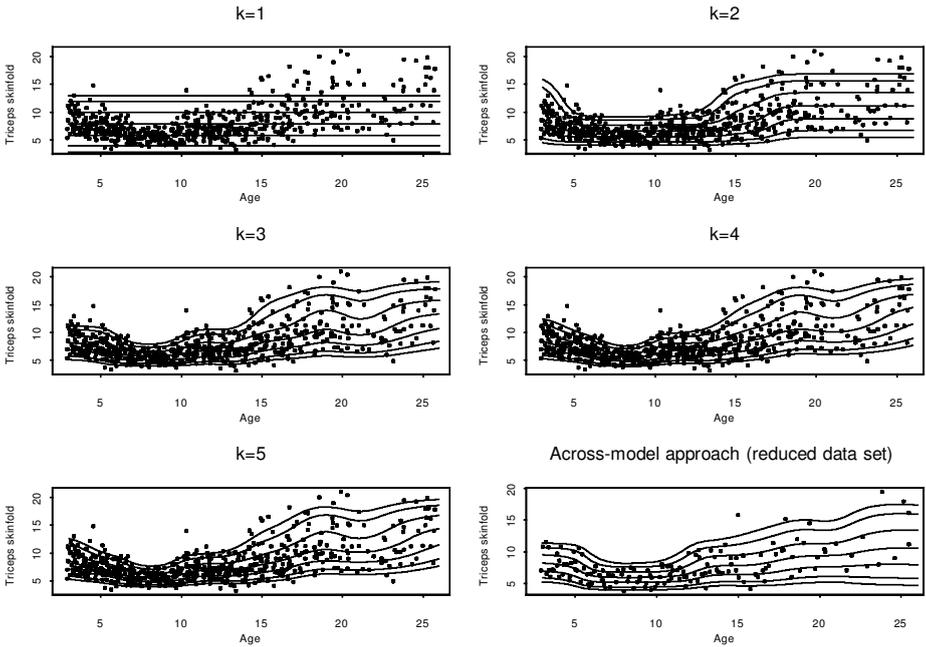


Figure 4. Centiles curves obtained with the five different models for triceps skinfold among Gambian females: 5th, 10th, 25th, 50th, 75th, 90th, and 95th percentiles.

They also showed a gradual increase in the spread of the data after age 10.

The across-model approach gives much the same results, with some obvious differences due to the fact that here only one quarter of the data are involved in the model fitting.

The centile curves are reasonably well smoothed up to age 20. After that they appear somewhat ragged. This might well be due to the fact that data for women aged more than 20 are very sparse. Smoother curves could be obtained by increasing the value of the smoothing parameters  $\lambda$  or  $\delta$ .

### 5.3 THE CHOICE OF THE SMOOTHING PARAMETERS

The choice of the smoothing parameters  $\lambda$  and  $\delta$  obviously involve an arbitrary decision. A possible solution to overcome this problem would be to assign a hyperprior to  $\lambda$  and  $\delta$ . This would allow extra flexibility to the model and cope with uncertainty about these parameters. In this way, in fact, the smoothing parameter values would be chosen by the data.

In the present context we preferred to regard the free choice of the smoothing parameters as an advantage of the model rather than a problem to be solved. By varying the smoothing parameters, in fact, features of the data that arise on different “scales” can be explored. A final value for  $\lambda$  and  $\delta$  can then be obtained by a subjective choice. Figure 5 shows the centile curves obtained for increasing values of  $\lambda$  (by row) and increasing values of  $\delta$  (by

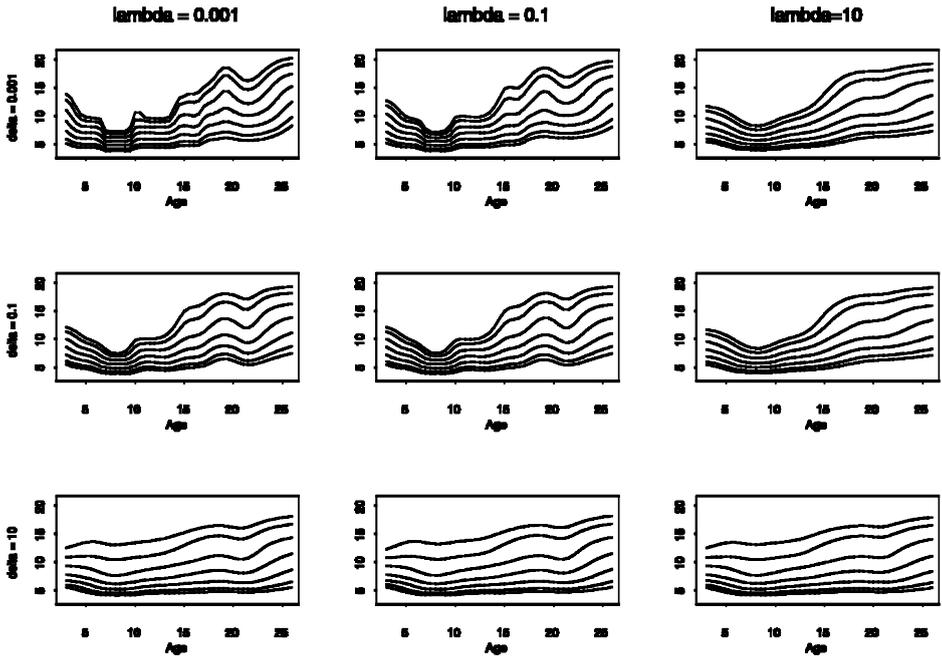


Figure 5. Centile curves obtained for different values of the smoothing parameters  $\lambda$  and  $\delta$ , conditioning on  $k = 4$ .

column), conditioning on  $k = 4$ . The curves become smoother as the parameters increase and  $\delta$ , as expected, has a stronger effect than  $\lambda$ , since  $\delta$  directly affects the variance of the  $\gamma_j$ . For  $\lambda = \delta = 10$  the curves show very little curvature and a farther increase in one or both parameters would determine straight curves, parallel to the covariate axis.

Our final choice of the smoothing parameters, with  $\lambda$  larger than  $\delta$  ( $\lambda = 0.7, \delta = 0.09$ ) aimed to get smoothed curves without placing any strong restriction on the prior distribution of  $\gamma_j$ .

### 5.4 POSTERIOR INFERENCE ON $\mu, \sigma^2$ , AND $w_j(t_i)$

Even if the focus of our analysis is on the inference about the growth curves, it is also interesting to spend a few words on the posterior distributions of the single components of the mixture and the posterior distributions of the allocation variables. Figure 6 shows the marginal posterior distribution of  $\mu$  (left panel),  $\sigma^2$  (central panel), and  $z$  (right panel) conditioning on increasing values of  $k$ , from  $k = 3$  to  $k = 5$ . Here, a line of a given type is used to represent the posterior densities of  $\mu$  and  $\sigma^2$ , for a given component, and the probability, conditional on the value of the covariate, of an observation to be allocated to that component. For all the models it can be noticed that the posterior distribution of  $\mu_j$  becomes more and more spread as we move from  $\mu_1$  to  $\mu_k$ . This is due to the fact that the

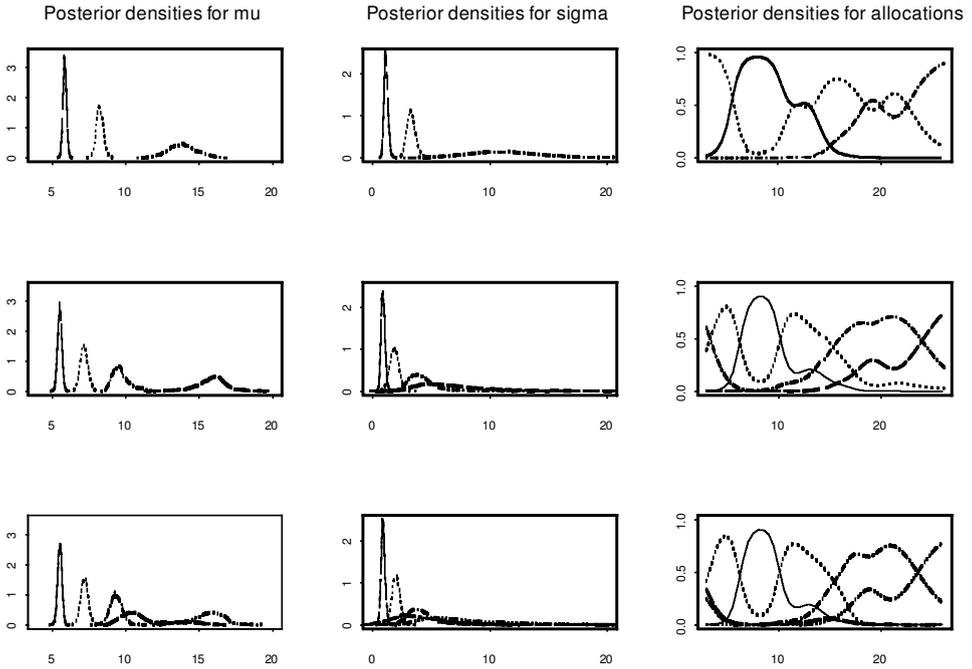


Figure 6. Marginal posterior densities for  $\mu$  and  $\sigma^2$  and  $\mathbf{z}$ . Models with three to five components.

mean value of  $\sigma_j^2$  increases with  $j$ , and also to the fact that fewer observations are generally allocated to the last components. If we consider the modal allocation, for the model with five components, 526 observations are allocated to the first three components and only 58 observations are allocated to the last two.

It is also interesting to notice how the effect of going from  $k = 3$  to  $k = 4$  components is that all the components are shifted, while the principal effect of going from  $k = 4$  to  $k = 5$  components is that the third component is split into two.

Looking at the posterior distribution of  $\sigma_j^2$ , a general increase in the variance of the distribution is again evident, together with an increase in the shift towards right, as we move from  $\sigma_1^2$  to  $\sigma_k^2$ . These increases are caused by a smaller number of observations allocated to components with large  $\mu$ , but also by the fact that these observations are the most dispersed. The only exception in the location of the posterior densities for  $\sigma^2$  can be noticed for  $k = 5$  where the posterior mode for  $\sigma_4^2$  is less than the posterior mode for  $\sigma_3^2$ .

The posterior densities of  $z_i$ , that is, the weight curves  $w_j(t_i)$ , can be analyzed using both Figure 6 and Figure 7. Figure 7 describes the modal allocation of each observation and gives a more immediate idea than Figure 6.

In the model with  $k = 3$ , the observations at ages less than five are very likely associated to the second component. Then the weight of this second component decreases in favor of the weight of the first one to which all the observations in the notch are allocated with probability very close to 1. The observations at ages more than 15 are associated again with

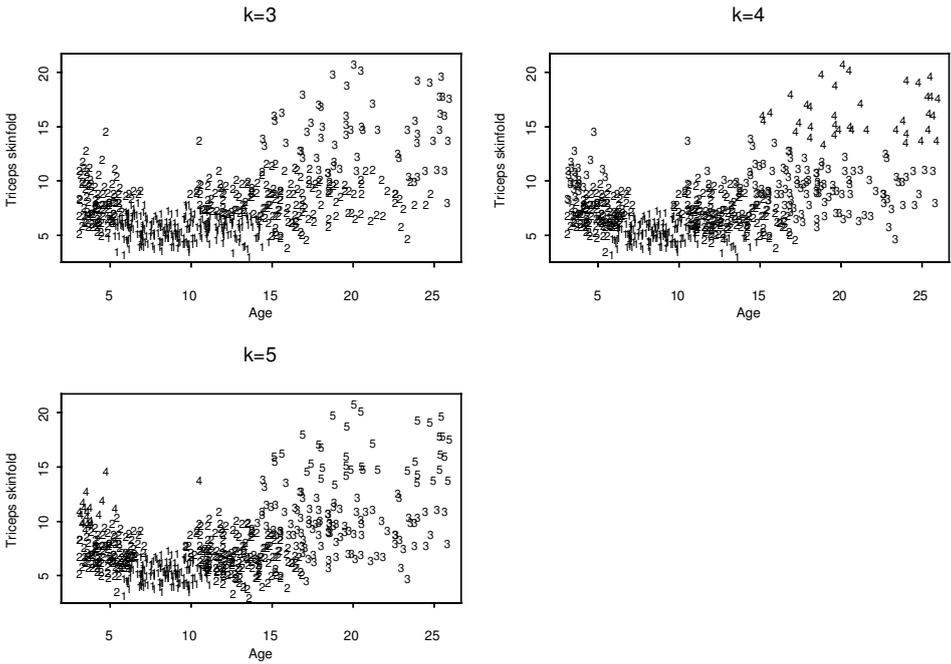


Figure 7. Modal allocation of the observations for models with three to five components.

the second component or with the third one, when their value is high.

In the models with  $k = 4$  and  $k = 5$ , the first component is still associated with the notch around the age 9. Observations on the edges of this notch are instead very likely to be allocated to the second component. The fourth component in the model with  $k = 4$  has the same role of the fifth component in the model with  $k = 5$ : it is associated with observations aged more than 15 with high triceps values. As noted before, going from  $k = 4$  to  $k = 5$  has the effect of splitting the third component. For  $k = 4$ , the third component is both associated with observations at young ages and with observations at ages around 15 and older, characterized by triceps around the value 10. In the model with  $k = 5$ , observations at young ages and observations at ages around 15 and older are instead separated and allocated with high probability to the fourth and the third component, respectively. It is evident that very few observations have a reasonable probability of being allocated to the fourth component. On the other side, these are much less dispersed than those likely to be allocated to the third one. This explains the fact that the posterior mode for  $\sigma_4^2$  is less than the posterior mode for  $\sigma_3^2$  when  $k = 5$ , and also the fact that  $\sigma_4^2$  is more dispersed.

## 6. DISCUSSION

We believe that the model introduced and discussed in this article provides an interesting new methodology for the estimate of growth curves. The main difficulty lies, in fact, in

deciding whether a bump or dip observed on a centile curve at a particular age is a real feature of the data, or whether it is simply sampling error. The Bayesian modeling approach—together with the flexible semiparametric nature of mixture models, the adaptability of the splines, and the roughness penalizing element they introduce—provides an elegant solution to this problem. Data are smoothed preserving their heterogeneity, no age cut-offs need to be specified, and the only arbitrariness in the whole procedure is the choice of the smoothing parameters  $\lambda$  and  $\delta$ .

The model could be extended in several directions. One interesting extension is the use of B-splines to model not only the weights of the mixture, but also the component parameters  $\mu_j$  and  $\sigma_j$ . This would give the model more flexibility. Let us think of an extreme example. If we consider a hypothetical dataset for which the underlying theoretical model is a simple linear regression on the covariate, it is clear that the mixture model presented in the article would require a large number of components to fit the data. Modeling the component parameters as a smoothed function of the covariate, would instead allow to fit the data using only one component.

Although the validation of the MCMC code did not receive much space in the article, numerous checks were conducted on the correctness of our sampler. In particular we checked that without any data, our estimate of the joint posterior distribution tallies with the chosen prior. We are satisfied that the values chosen for the smoothing parameters  $\lambda$  and  $\delta$  allow substantial smoothing in the centile curves, so that higher values of  $\lambda$  and  $\delta$ , for which mixing could be slower, are not necessary.

Finally, we draw attention on the difficult matter of estimating the marginal likelihood for complex, nonregular problems, when the large amount of data induce very peaked posterior distributions for the parameters and makes the use of reversible jump infeasible.

## ACKNOWLEDGMENTS

We acknowledge partial support of this work by the UK Engineering and Physical Sciences Research Council, and are grateful for helpful comments by the referees and editor, which have led to improvement of the presentation.

*[Received September 2001. Revised January 2002.]*

## REFERENCES

- Besag, J., Green, P. J., Higdon, D., and Mengersen, K. (1995), "Bayesian Computation and Stochastic Systems" (with discussion), *Statistical Science*, 10, 3–66.
- Böhning, D. (2000), *Computer-Assisted Analysis of Mixtures and Applications*, Monograph on Statistics and Probability 81, London: Chapman and Hall/CRC.
- Cole, T. J., and Green, P. J. (1992), "Smoothing Reference Centile Curves: The LMS Method and Penalized Likelihood," *Statistics in Medicine*, 11, 1305–1319.

- Fearn, T. (1975), "A Bayesian Approach to Growth Curves," *Biometrika*, 62, 89–100.
- Geisser, S. (1970), "Bayesian Analysis of Growth Curves," *Sankhyā*, 32, 53–64.
- Gelman, A., and Meng, X. L. (1998), "Simulating Normalizing Constants: From Importance Sampling to Bridge Sampling to Path Sampling," *Statistical Science*, 13, 163–185.
- Green, P. J. (1995), "Reversible Jump Markov Chain Monte Carlo Computation and Bayesian Model Determination," *Biometrika*, 82, 711–732.
- Green, P. J., and Richardson, S. (2002), "Hidden Markov Models and Disease Mapping," *Journal of the American Statistical Association*, 97, 1055–1070.
- Green, P. J., and Silverman B. W. (1994), *Nonparametric Regression and Generalized Linear Models*, London: Chapman and Hall.
- Hammersley, J. M., and Handscomb, D. C. (1964), *Monte Carlo Methods*, London: Chapman and Hall.
- Newton, M. A., and Raftery, A. E. (1994), "Approximate Bayesian Inference With the Weighted Likelihood Bootstrap" (with discussion), *Journal of the Royal Statistical Society, Series B*, 56, 3–48.
- Richardson, S., and Green, P. J. (1997), "On Bayesian Analysis of Mixtures With an Unknown Number of Components" (with discussion), *Journal of the Royal Statistical Society, Series B*, 59, 731–792.
- Tierney, L. (1994), "Markov Chains for Exploring Posterior Distributions," *The Annals of Statistics*, 22, 1701–1762.
- Titterton, D. M., Smith, A. F. M., and Makov, U. E. (1985), *Statistical Analysis of Finite Mixture Distributions*, Chichester: Wiley.
- Wahba, G. (1978), "Improper Priors, Spline Smoothing and the Problem of Guarding Against Model Errors in Regression," *Journal of the Royal Statistical Society, Series B*, 40, 364–372.
- (1990), *Spline Models for Observational Data*, Philadelphia, PA: SIAM.