

# Exact and Efficient Bayesian Inference for Multiple Changepoint problems

Paul Fearnhead

Department of Mathematics and Statistics  
Lancaster University

**Summary** We demonstrate how to perform direct simulation from the posterior distribution of a class of multiple changepoint models where the number of changepoints is unknown. The class of models assumes independence between the posterior distribution of the parameters associated with segments of data between successive changepoints. This approach is based on the use of recursions, and is related to work on product partition models. The computational complexity of the approach is quadratic in the number of observations, but an approximate version, which introduces negligible error, and whose computational cost is roughly linear in the number of observations, is also possible. Our approach can be useful, for example within an MCMC algorithm, even when the independence assumptions do not hold. We demonstrate our approach on well-log data. Our method can cope with a range of models for this data, and exact simulation from the posterior distribution is possible in a matter of minutes.

**Keywords** *Bayes factor, Forward-backward algorithm, Model choice, Perfect simulation, Reversible jump MCMC, Well-log data*

## 1 Introduction

Many time-series models incorporate one, or multiple, changepoints. Some examples include Poisson processes with a piece-wise constant rate parameter (Raftery and Akman, 1986; Yang and Kuo, 2001; Ritov *et al.*, 2002), changing linear regression models (Carlin *et al.*, 1992; Lund and Reeves, 2002), Gaussian observations with varying mean (Worsley, 1979) or variance (Chen and Gupta, 1997; Johnson *et al.*, 2003), and Markov models with time-varying transition matrices (Braun and Muller, 1998). Such models have been used for modelling stock prices, muscle activation, climatic time-series, DNA sequences and neuronal activity in the brain, amongst many other applications

In this paper we consider Bayesian analysis for a class of multiple changepoint problems. We call a period of time between two consecutive changepoints a *segment*. This class of

models assumes that the parameter values associated with each segment are independent from each other. Yang and Kuo (2001) comment that calculating the Bayes factors for models with different numbers of changepoints is “essentially infeasible for a large model with many changepoints”. Our aim is to show that calculation of Bayes factors, and perfect sampling from the posterior distribution of changepoint locations, is both possible, and computationally inexpensive for the class of models we consider. While this class of models may seem restrictive, recent examples of work on such models can be found in Johnson *et al.* (2003), Punskeya *et al.* (2002), and Braun *et al.* (2000).

Although we use the phrase “perfect simulation”, we do not use coupling-from-the-past (Propp and Wilson, 1996), or related ideas, which have become synonymous with this phrase. Instead, the work we present is closely related to work by Yao (1984), Barry and Hartigan (1992) and Barry and Hartigan (1993). These papers present efficient recursions that allow the posterior probabilities of different numbers of changepoints, and the posterior mean of the parameters to be calculated. Despite the desirability of exact solutions, and the simplicity and computational efficiency of the recursions, these methods are currently underused. We extend these methods to allow for direct simulation from the posterior distribution of the number and position of the changepoints, and to also perform inference conditional on the number of changepoints.

Much recent research for changepoint models is based on the use of MCMC. For models with an unknown number of changepoints, a common approach is that of Green (1995). A set of models, each incorporating a different number of changepoints, are introduced, and reversible jump MCMC is used to explore the joint space of model and parameters. Potential difficulties of this approach include designing moves, particular ones between different models, which enable the MCMC algorithm to mix well (for guidelines on designing reversible jump MCMC algorithms see Brooks *et al.*, 2003), and being able to detect convergence of the algorithm. For example, in the analysis of the coal-mining disaster data in Green (1995), the reversible jump MCMC algorithm had not converged. The reanalysis of the data in Green (2003), using a reversible jump MCMC algorithm run for 25 times as long, does fully explore the posterior distribution. The exact simulation method we describe here avoids any problems of needing to diagnose convergence of an MCMC algorithm.

We consider two classes of prior for the changepoint process. One, that of Green (1995), involves a prior on the number of changepoints, and then a conditional prior on their

position. The other is based on modelling the changepoint process by a point process (Pievatolo and Green, 1998), and is a special case of a product-partion model (Hartigan, 1990). This indirectly specifies a joint prior on the number and position of the changepoints. In both cases we assume that, conditional on the realisation of the changepoint process, the joint posterior distribution of the parameters is independent across the segments of the time series. We also assume a conjugate prior for the parameters associated with each segment. Under these two assumptions we derive a set of recursions to perform exact inference.

The recursions are similar to those of the Forward-Backward algorithm (see Scott, 2002, for a review). Recent work has shown how such recursions can be used to perform exact inference for a range of problems (Fearnhead and Meligkotsidou, 2004; Fearnhead, 2004). The assumption of independence between segments ensures the necessary Markov property that is required for Forward-Backward type recursions. For a data set consisting of observations at discrete times,  $1, \dots, n$ , the recursions are based on calculating the probability of the data from time  $t$  to time  $n$ , given a changepoint at time  $t$ , in terms of the equivalent probabilities at times  $t + 1, \dots, n$ . Once these probabilities have been calculated for all time-points, it is possible to directly simulate from the posterior distribution of the time of the first changepoint, and then the conditional distribution of the time of the second changepoint, given the first, and so on. The recursions can also be used to perform exact inference conditional on the number of changepoints, and in some cases to calculate the posterior distribution of the parameters that govern the point process model for the changepoints.

The computational cost of the recursions increases quadratically with  $n$ . However an approximate version, which introduces negligible error, is possible. In limiting situations where the length of time series increases, and the number of changepoints is increasing linearly with the number of observations, the computational cost increases roughly linearly with  $n$ . (In the alternative limiting regime of more frequent observations, the computational cost remains quadratic in  $n$ .) The assumption of conjugate priors can potentially be relaxed, but with an increase in the computational cost. Essentially, low-dimensional integrals that can be calculated analytically under conjugate priors would need to be calculated numerically (for example see Section 4.3). Relaxation of the independence assumption is more difficult, but our algorithm can still be used as a useful tool for analysing such data. For example, the algorithm can be embedded in an MCMC algorithm, and we demonstrate such an approach on some real data.

The outline of the paper is as follows. In Section 2 we introduce the two classes of changepoint model that we consider. The recursions are derived and detailed in Section 3. The resulting algorithm is demonstrated on well-log data in Section 4. We consider a range of models for this data, and also demonstrate how our method can be used to analyse the data when there is dependence between the parameters for each segment. The paper concludes with a discussion.

## 2 Models and Notation

We consider the following class of multiple changepoint models. Consider a sample of size  $n$ ,  $y_1, \dots, y_n$ . Observation  $y_i$  is obtained at time  $i$ , and we let  $y_{i:j}$  denote the observations from time  $i$  to time  $j$  inclusive.

Firstly condition on  $m$  integer-valued changepoints, at points  $0 < \tau_1 < \tau_2 < \dots < \tau_m < n$ . We let  $\tau_0 = 0$  and  $\tau_{m+1} = n$ . Then the  $j$ th segment consists of the observations from time  $\tau_{j-1} + 1$  to time  $\tau_j$ . We associate a (possibly vector-valued) parameter  $\theta_j$  with the  $j$ th segment for  $j = 1, \dots, m + 1$ . Conditional on the change-points and parameter values, the observations are independent; observation  $y_i$  being drawn from a density  $f(y_i|\theta_j)$  if time  $i$  is in the  $j$ th segment.

We assume independent priors for the parameters associated with each segment. The prior for  $\theta_j$  is denoted by  $\pi(\theta_j)$ . Here, and throughout, we use  $\pi(\cdot)$  solely to denote a prior density; the argument making it clear as to which parameter the prior is for.

We assume that the changepoints occur at discrete time points, and consider two priors for the changepoints. The first prior is based on a prior for the number of changepoints, and then a conditional prior on their positions. We will define this conditional prior on the positions in terms of  $\pi_m(\tau_m)$  the prior for the last change point, and, for  $j = 1, \dots, m - 1$ ,  $\pi_m(\tau_j|\tau_{j+1})$ , the prior for the position of the  $j$ th changepoint, given the position of the  $(j + 1)$ st.

The second prior is obtained from a point process on the positive and negative integers. The point process is specified by the probability mass function  $g(t)$  for the time between two successive points. We assume that this time must be a strictly positive integer. We observe the point process on the interval  $[1, n - 1]$ , and assume that changepoints occur at the positions of points in the point process. This prior is an example of a product-partition model.

If  $G(t) = \sum_{s=1}^t g(s)$ , is the distribution function of the distance between two successive points, and  $g_0(t)$  is the mass function of the first point after 0, then the probability of  $m$  changepoints occurring at  $\tau_1, \dots, \tau_m$  is

$$g_0(\tau_1) \left( \prod_{j=2}^m g(\tau_j - \tau_{j-1}) \right) (1 - G(\tau_{m+1} - \tau_m)).$$

Natural choices for the distribution of the time between successive points are from the negative binomial family. For a negative binomial distribution with parameters  $k$ , a positive integer, and  $p$  we have

$$g(t) = \binom{t-k}{k-1} p^k (1-p)^{t-k} \quad g_0(t) = \sum_{i=1}^k \binom{t-i}{i-1} p^i (1-p)^{t-i} / k.$$

The negative binomial distribution can be thought of as a discrete version of the gamma distribution (especially if  $p$  is small). If  $k = 1$  then the negative binomial distribution is the geometric distribution, and the point process is Markov. Larger values of  $k$  can reduce the number of very short segments.

### 3 Filtering Recursions

We first derive the recursions for analysing data under the point process prior for the changepoints. We later derive recursions to perform inference conditional on the number of changepoints, and show how these can be used to perform inference under the other prior, and to perform inference about the parameters of the point process prior.

#### 3.1 Basic Recursions

For times  $s \geq t$ , define

$$\begin{aligned} P(t, s) &= \Pr(y_{t:s} | t, s \text{ in the same segment}) \\ &= \int \prod_{i=t}^s f(y_i | \theta) \pi(\theta) d\theta. \end{aligned}$$

We will assume that the probabilities  $P(t, s)$  can be calculated for all  $t$  and  $s$ . In practice this will require conjugate priors on  $\theta$ , or, if  $\theta$  is low-dimensional, that the required integration can be calculated numerically.

We next define for  $t = 2, \dots, n$

$$Q(t) = \Pr(y_{t:n} | \text{changepoint at } t-1),$$

with  $Q(1) = \Pr(y_{1:n})$ . A set of recursions for calculating these probabilities are given by the following theorem.

**Theorem 1** *Define the probabilities  $Q(t)$  and  $P(t, s)$  as above. Then for  $t = 2, \dots, n$*

$$Q(t) = \sum_{s=t}^{n-1} P(t, s)Q(s+1)g(s+1-t) + P(t, n)(1 - G(n-t)), \quad (1)$$

and

$$Q(1) = \sum_{s=1}^{n-1} P(1, s)Q(s+1)g_0(s) + P(1, n)(1 - G_0(n-1)), \quad (2)$$

where  $G_0(t) = \sum_{s=1}^t g_0(s)$ .

**Proof:** We only prove Equation 1. Equation 2 can be derived similarly.

For notational convenience we drop the explicit conditioning on a changepoint at  $t-1$  in the following. Thus,

$$\begin{aligned} Q(t) &= \Pr(y_{t:n}) \\ &= \sum_{s=t}^{n-1} \Pr(y_{t:n}, \text{next changepoint at } s) + \Pr(y_{t:n}, \text{no further changepoints}). \end{aligned}$$

Now these probabilities can be calculated by the product of the prior probability on the changepoints, and the probabilities of the observations from a single segment,  $P(t, s)$ .

Thus

$$\begin{aligned} &\Pr(y_{t:n}, \text{next changepoint at } s) \\ &= \Pr(\text{next changepoint at } s) \Pr(y_{t:s}, y_{s+1:n} | \text{next changepoint at } s) \\ &= g(s+1-t) \Pr(y_{t:s} | t, s \text{ in same segment}) \Pr(y_{s+1:n} | \text{changepoint at } s) \\ &= g(s+1-t) P(t, s) Q(s+1) \end{aligned}$$

Similarly

$$\Pr(y_{t:n}, \text{no further changepoints}) = P(t, n)(1 - G_0(n-t)),$$

as required.  $\square$

Equations 1 and 2 give recursions that can be used to calculate  $Q(t)$  in turn for  $t = n, \dots, 1$ . The evidence of the model is just  $Q(1)$ . These equations are equivalent to those of Barry and Hartigan (1992), and are based on the same idea as recursions of Yao (1984).

The computational complexity of the resulting algorithm is quadratic in  $n$ . However often only a small proportion of the terms on the right-hand side of (1) make an appreciable contribution to  $Q(t)$ . This can happen when the data makes it almost certain that a changepoint occurs before a given time-point. Thus the summation can often be truncated with negligible error. We propose truncating the sum at term  $k$  when

$$\frac{P(t, k)Q(s + 1)g(k + 1 - t)}{\sum_{s=t}^k P(t, s)Q(s + 1)g(s + 1 - t)} \quad (3)$$

is less than some predetermined value, for example  $10^{-10}$ .

In the limiting regime of analysing a process over a longer time period, so that the number of changepoints will increase roughly linearly with the number of observations,  $n$ , the computational complexity of the resulting approximate set of recursions will be linear in  $n$ . Essentially the average number of terms required in the right-hand side of (1) will be constant with  $t$ . Thus the average computational cost of one of the  $n$  recursions will be independent of  $n$ .

### 3.2 Perfect Simulation of Changepoints

Given the values of  $Q(t)$  for  $t = 1, \dots, n$  it is straightforward to simulate from the posterior distribution of the changepoints as follows.

The posterior distribution of the first changepoint is given by

$$\begin{aligned} \Pr(\tau_1|y_{1:n}) &= \Pr(y_{1:n}, \tau_1) / \Pr(y_{1:n}) \\ &= \Pr(\tau_1) \Pr(y_{1:\tau_1}|\tau_1) \Pr(y_{\tau_1+1:n}|\tau_1) / Q(1) \\ &= P(1, \tau_1)Q(\tau_1 + 1)g_0(\tau_1) / Q(1), \end{aligned}$$

for  $\tau_1 = 1, \dots, n - 1$ . The probability of no further changepoint being  $P(1, n)(1 - G_0(n - 1)) / Q(1)$ .

Similarly the posterior distribution of the  $\tau_j$  given  $\tau_{j-1}$  is

$$\Pr(\tau_j|\tau_{j-1}, y_{1:n}) = P(\tau_{j-1} + 1, \tau_j)Q(\tau_j + 1)g(\tau_j - \tau_{j-1}) / Q(\tau_{j-1} + 1),$$

for  $\tau_j = \tau_{j-1} + 1, \dots, n - 1$ , and the probability of no further breakpoint is  $P(\tau_{j-1} + 1, n)(1 - G_0(n - \tau_{j-1} - 1)) / Q(\tau_{j-1} + 1)$ .

Efficient simulation of large samples of changepoints from the posterior distribution can be done by simulating the samples concurrently, using the following algorithm. We

denote the generic posterior distribution of the next changepoint, given a changepoint at  $t$  by  $\Pr(\tau|y_{1:n}, t)$ , which can be calculated as above.

- (1) For a sample of size  $M$ , initiate each of the  $M$  samples with a changepoint at  $t = 0$ .
- (2) For  $t = 0, \dots, n - 2$ :
  - (i) Calculate  $n_t$  the number of whose last changepoint was at time  $t$ .
  - (ii) If  $n_t > 0$  calculate the probability distribution  $\Pr(\tau|y_{1:n}, t)$ .
  - (iii) Sample  $n_t$  times from  $\Pr(\tau|y_{1:n}, t)$  using Algorithm 1 of Carpenter *et al.* (1999) (see the Appendix). Use these values to update the  $n_t$  samples of changepoints which have a changepoint at  $t$ .

There are two advantages of this algorithm. The first is that the probability mass function  $\Pr(\tau|y_{1:n}, t)$  need only be calculated once regardless of the number of samples required from it. If changepoints are sampled one at a time, then either these densities will, potentially, need to be calculated for each sample, or they will need to be stored. Storing these mass functions can place large burdens on computational memory. The storage requirements will be quadratic in  $n$ ; by comparison the above algorithm has storage requirements that are linear in  $n$ .

The second is that simulating a sample of size  $m$  from a general discrete mass function can be achieved more efficiently than sampling  $m$  samples of size 1. Algorithm 1 of Carpenter *et al.* (1999) allows a sample of size  $m$  to be simulated with order  $n + m$  effort, rather than the  $nm$  effort of sampling  $m$  samples of size 1.

### 3.3 Conditioning on the Number of Changepoints

Now consider inference conditional on  $m$  changepoints. As in Section 2 we define the prior for the changepoints via  $\pi_m(\tau_m)$  and conditional probabilities of the form  $\pi_m(\tau + j|\tau_{j+1})$ . We define  $P(s, t)$  as before, and for  $j = 1, \dots, m$ , and  $t = j + 1, \dots, n - m - 1 + j$ ,

$$Q_j^{(m)}(t) = \Pr(y_{t:n}|\tau_j = t - 1, m \text{ changepoints}).$$

We can derive the following set of recursions. For  $t = m + 1, \dots, n - 1$ ,

$$Q_m^{(m)}(t) = P(t, n)\pi_m(\tau_m = t - 1).$$



For  $j = 1, \dots, m - 1$ , and  $t = j + 1, \dots, n - m - 1 + j$

$$Q_j^{(m)}(t) = \sum_{s=t}^{n-m+j} P(t, s) Q_{j+1}^{(m)}(s + 1) \pi_m(\tau_j = t - 1 | \tau_{j+1} = s).$$

Finally

$$\Pr(y_{1:n} | m \text{ changepoints}) = \sum_{s=1}^{n-m} P(1, s) Q_1^{(m)}(s + 1). \quad (4)$$

These can be proved in a similar way to Theorem 1.

If the number of changepoints is unknown, with prior  $\pi(m)$ , then the posterior distribution of  $m$  can be calculated as

$$\Pr(m | y_{1:n}) \propto \pi(m) \Pr(y_{1:n} | m \text{ changepoints}),$$

with the last term, the evidence for  $m$  changepoints, being calculated, for each  $m$ , using the recursions.

Simulation from the joint posterior distribution is possible by first simulating  $M$  samples from  $\Pr(m | y_{1:n})$ . If the value  $m$  is sampled  $N_m$  times, then  $N_m$  samples from the posterior distribution of the changepoint positions, conditional on  $m$  changepoints, can be obtained as described in Section 3.2. The only difference is that the conditional distribution of  $\tau_j$  given  $\tau_{j-1}$  is now

$$\Pr(\tau_j | \tau_{j-1}, y_{1:n}, m) = P(\tau_{j-1} + 1, \tau_j) Q_j^{(m)}(\tau_j + 1) \pi_m(\tau_{j-1} | \tau_j) / Q_{j-1}^{(m)}(\tau_{j-1}).$$

Finally, in the case of the Markov point process prior (that is, a geometric distribution for the distance between changepoints), exact inference is possible even if the probability of a changepoint at any timepoint,  $p$ , is unknown. This is because, conditional on the number of changepoints, the positions are distributed uniformly along the interval, independent of  $p$ . We can thus perform inference conditional on  $m$  changepoints. The prior for  $m$  is obtained by averaging  $\pi(m | p)$  with respect to the prior for  $p$ .

## 4 Well-log Data

We now consider the problem of detecting changepoints in well-log data. An example of well-log data, which comes from Ó Ruanaidh and Fitzgerald (1996), is given in Figure 1. The data consist of measurements of the nuclear-magnetic response of underground rocks. The data were obtained by lowering a probe into a bore-hole. Measurements

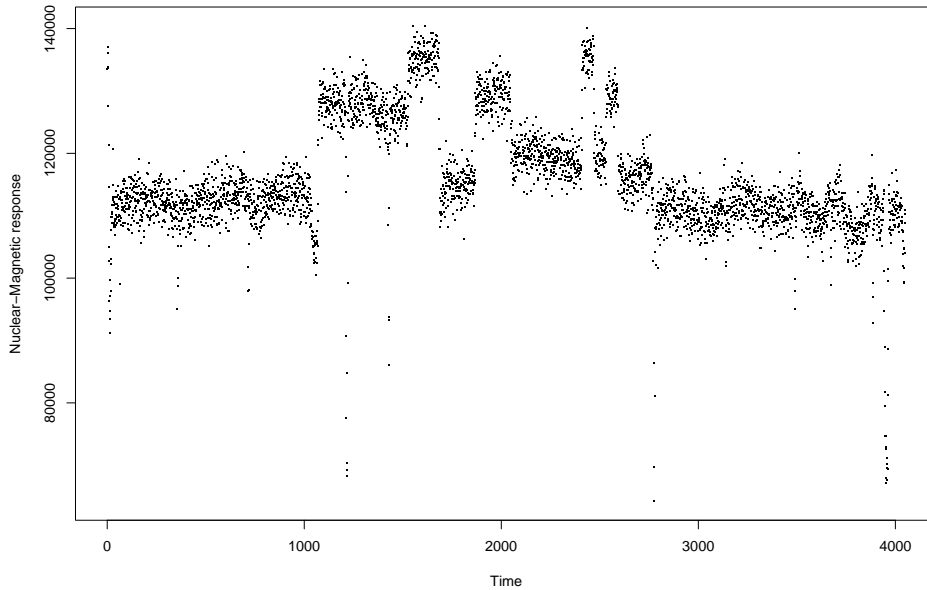


Figure 1: The well-log data.

were taken at discrete timepoints by the probe as it was lowered through the hole. The underlying signal is roughly piecewise constant, with each constant segment relates to a single rock type (that has constant physical properties). The changepoints in the signal occur each time a new rock type is encountered. Detecting the changepoints is important in oil-drilling; see the introduction of Fearnhead and Clifford (2003) for more details.

These data have been previously analysed by Ó Ruanaidh and Fitzgerald (1996), who used MCMC to fit a change-point model with a fixed number of changepoints; and by Fearnhead and Clifford (2003) who considered online analysis of the data using particle filters. We performed a batch analysis of the data, but allowed for multiple changepoints.

#### 4.1 Piecewise constant model

Initially we consider analyse based on a model taken Fearnhead and Clifford (2003). We assume a Markov point process prior, with  $p = 1/250$ , for the changepoints. There are a number of outliers in the data which were removed before the data was analysed. For a time  $t$  which belongs to segment  $i$ , we model a non-outlying observation,  $y_t$ , by

$$y_t \sim N(\mu_i, \sigma^2),$$

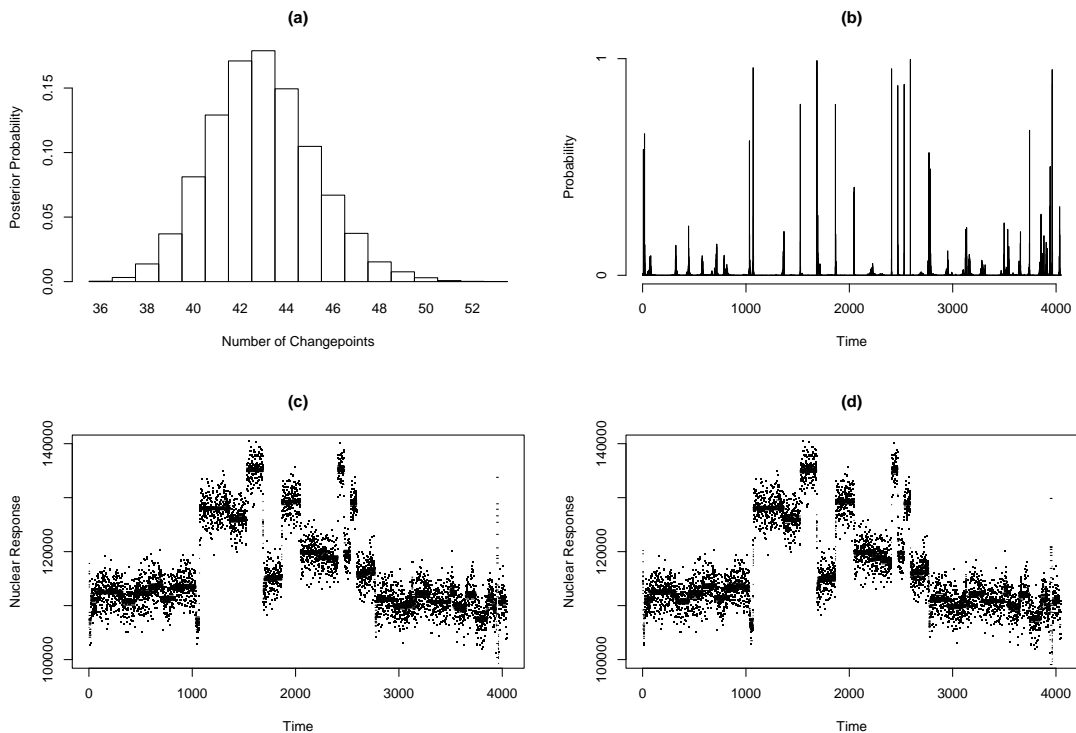


Figure 2: Results of analysis of the well-log data. Figures (a)–(c) are a model with  $p = 0.004$ : histogram of number of changepoints(a); posterior distribution of position of changepoints (b); and 20 simulations from the posterior distribution of the signal (c). Figure (d) shows 20 simulations from the posterior distribution of the signal when  $p = 0.013$ .

where  $\mu_i$  is the mean associated with the  $i$ th segment, and we assume a common known variance,  $\sigma^2 = 2500^2$ . We assume that the segment means have independent normal priors with mean 115,000 and variance 10,000<sup>2</sup>. Conditional on the segment means, the observations are independent.

All the parameter values are based on a simple analysis of the data. Ideally they would have been obtained from analysing related data, but such data was not available. The only difference in the model from that used by Fearnhead and Clifford (2003) is that the outliers are removed before the analysis. By comparison, Fearnhead and Clifford (2003) detect outliers as the data is processed. Ó Ruanaidh and Fitzgerald (1996) also removed outliers before analysing the data. However they assumed that 13 changepoints were present in the data, and they modelled the distribution of the observations as Laplacian. Fearnhead (1998) shows that the normal model is more appropriate.

The results of our analysis are shown in Figure 2(a)–(c). These plots are for the posterior distribution of the number of changepoints and their positions (based on 10,000 independent simulations from the posterior), and 20 realisations of the underlying signal. It took 26 seconds on a 3.4GHz PC to perform this analysis. The results we obtain differ substantially from previous analyses. The posterior distribution suggests around 40 changepoints, which is roughly three times as many as assumed by Ó Ruanaidh and Fitzgerald (1996). Fearnhead and Clifford (2003) only infer 16 changepoints. However theirs is an online analysis, so inference at each timepoint is based solely on the observations to that timepoint. Furthermore, they only inferred a changepoint when the posterior probability of changepoint within the last 10 time-points was greater than 0.9. This is likely to produce a conservative estimate of the number of changepoints.

The posterior mean of the number of changepoints, 43, is much larger than the 8 to 26 changepoints we would expect under our prior. Using the methods of Section 3.3 we can perform inference when  $p$  is unknown. If we assume a uniform prior for  $p$ , then the posterior mode is at  $p = 0.013$ . Some results of analysing the data with this value of  $p$  are shown in Figure 2(d). The posterior distribution for the number of changepoints is substantially different, with a posterior mean of 52 (not shown), but realisations from the posterior distributions of the underlying signal are similar to those when  $p = 0.004$ .

We repeated our analysis using the approximate algorithm suggested in Section 3.1. We truncated the sums in Equations (1) and (2), used to calculate the  $Q(t)$ s, when the value of (3) was less than  $10^{-10}$ . The resulting algorithm on average required sums of 222 terms to be calculated for each  $Q(t)$ ; which compares with average sums of 2025 terms for the exact algorithm. This is nine-fold reduction in the complexity of the algorithm. The resulting approximation of the log evidence was correct to 4 decimal places, which suggests that negligible errors were introduced.

## 4.2 Inclusion of Hyperpriors

We now consider an extension of the above model where all parameters in our model were unknown, and we introduce hyperpriors for them. This introduces dependence between the segments, and our direct simulation algorithm has to be used with an MCMC scheme,

We used a uniform prior for  $p$ , and an improper prior for  $\sigma$ ,  $\pi(\sigma) \propto 1/\sigma$ . We param-

terised the prior for the segment means as

$$\mu_i \sim \text{N}(\eta, \tau^2 \sigma^2),$$

and used improper hyperpriors on  $\eta$  and  $\tau$ :  $\pi(\eta) \propto 1$  and  $\pi(\tau) \propto 1/\tau$ .

We analysed this model using MCMC. The MCMC algorithm used the following three updates:

- (1) Update the changepoints conditional on  $\sigma$ ,  $p$ ,  $\eta$ , and  $\tau$ . We used an independent proposal from the true posterior distribution conditional on  $\sigma = 2,330$ ,  $p = 0.013$ ,  $\eta = 115000$ , and  $\tau = 4.3$ .
- (2) Update  $\sigma$ ,  $p$  and the  $\mu_i$ s from their full conditional distribution given the changepoints and  $\eta$  and  $\tau$ .
- (3) Update  $\eta$  and  $\tau$  from their full-conditionals given the  $\mu_i$ s.

Each of these moves satisfies detailed balance. Steps (2) and (3) are Gibbs steps, and thus the proposed values are always accepted. Step (1) is not a Gibbs step. Although it would be possible to make it so, there is a substantial overhead to calculating the posterior distribution of the changepoints at each iteration. Thus while this algorithm may mix more slowly, a single iteration will be substantially quicker, and hence we hope it will be more efficient. In updating the changepoints in step (1) we throw away the segment means. This is an example of collapsing (Liu, 2001, pages 146–151), which usually improves the mixing of the Markov Chain.

We ran this Markov chain for 10,000 iterations. The acceptance probability of step (1) was 61.8%. The 1-lag autocorrelation for each of the parameters was less than 0.03, which suggests that the chain is mixing extremely quickly.

The reason why this MCMC algorithm performs so well is because the posterior probability of the parameters is concentrated in a small region of the parameter space. Over this small region, the parameters are almost independent; the maximum absolute value of the correlation between any pair of parameters is 0.01. Furthermore, the conditional distribution of the changepoints changes little over this range of parameter values, which means that the average acceptance probability in step (1) of the algorithm is high. This situation is likely to occur in other situations where there is a large and informative data set with many changepoints.

### 4.3 Alternative Models

The models used in the previous Sections are based around those previously used in the literature for this data. However the realisations from the posterior distribution have many more changepoints, and thus suggest many more rock strata, than is realistic. It appears that the piecewise constant model used is overly simplistic for the data, and that this has resulted in the need for too many changepoints in order to fit the data.

We have considered numerous extensions to the model. Two possibilities are: (i) to allow different noise variances for different segments; and (ii) to model each segment using a mean-shifted AR(1) model (Albert and Chib, 1993). Both of these models can be analysed via our direct simulation method, though for (ii) we need to numerically integrate out the autoregressive coefficient (this can be done in a similar way to that described below). However neither of these extensions enable the data to be fit with substantially fewer change points (results not shown).

Instead we consider the following state-space model for the data within a segment, where if  $t - 1$  and  $t$  both lie within segment  $i$

$$\begin{aligned}\mu_t &\sim \text{N}(\mu_{t-1}, \tau_i^2) \\ y_t &\sim \text{N}(\mu_t, \sigma^2).\end{aligned}$$

The initial  $\mu$  value for each segment is drawn from the same independent normal priors as before. This is an extension of the piecewise constant model which allows the signal within a segment to perform a random walk. We allow the variance of the random walk to vary among segments, and assume a Gamma prior for  $\tau_i$  with parameters 2 and 1/40. This prior places most probability mass on values of  $\tau_i$  which lie in the interval  $[0, 150]$ . The idea of this model is that the random walk element can fit the small-scale variation in the underlying signal without the need to infer changepoints.

If  $\tau_i$  were known for each segment then it would be straightforward to apply our direct simulation method, using the Kalman Filter (Harvey, 1989) to integrate out the underlying signal. To incorporate a prior on  $\tau_i$  we resort to numerical integration to calculate the  $P(t, s)$  values required by our algorithm. A simple, but adequate, approach to numerical integration is based on using a grid of  $\tau_i$  values, and we obtained such a grid as follows. For a grid with  $K$  points, first simulate for  $k = 1, \dots, K$ , a realisation,  $u_k$ , of a uniform random variable on  $[(k - 1)/K, k/K]$ ; then fix the  $k$ th grid point to be the  $u_k$ th quantile from the prior for  $\tau_i$ .

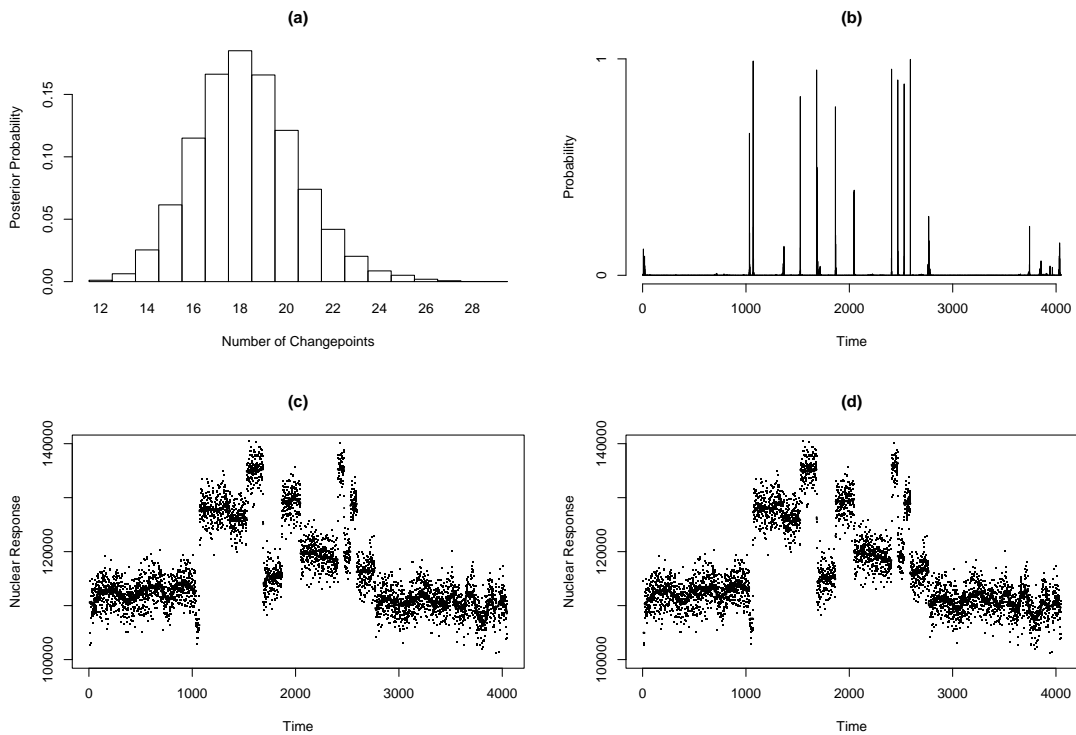


Figure 3: Results of analysis of the well-log data. using the state-space model: histogram of number of changepoints(a); posterior distribution of position of changepoints (b); and realisations from the posterior distribution of the signal (c)–(d). The results for (a) and (b) are based on 10,000 perfect simulations from the posterior distributions.

In practice we found that a grid of 100 points produced accurate results; and for such a grid it took less than 19 minutes to simulate 10,000 draws from the joint posterior distribution of changepoint positions on a 3.4GHz PC. The results of the analysis (assuming  $\sigma = 2,500$  and  $p = 0.004$ ) are shown in Figure 3. This model gives more realistic inferences about the number and positions of the changepoints. Further refinements of the model may be appropriate and could further improve the inferences, for example by choosing a distribution for the segment lengths that does not allow very short segments, but these are not considered here.

## 5 Discussion

We have described ways in which recursions, based on the Forward-Backward algorithm, can be used to perform Bayesian analysis of multiple changepoint problems. As men-

tioned previously, the work we present is closely related to work by Barry and Hartigan (1992). The main novelty of what we propose is that we demonstrate how the recursions can be used for perfect simulation from the posterior distribution of the number and position of change-points, and hence from the posterior distribution of the parameters. Presenting results from a Bayesian analysis via simulations from the posterior distribution is both quicker than calculating the posterior means (as done by Barry and Hartigan, 1992, where the cost is cubic in the number of observations), and also encapsulates information about uncertainty about parameters, which is one of the advantages of Bayesian inference. We have also extended the use of recursions to inference conditional on the number of changepoints.

The ability to simulate from posterior distributions also enables the algorithms we present to be used in analysing more complex models, for example by embedding our algorithm within an MCMC algorithm (see Section 4). While it may seem natural in such cases just to use standard MCMC algorithms, the use of direct simulation enabled us to construct an MCMC algorithm for the Well-log data that had exceptional mixing properties.

While the main focus of this paper is this new methodology, we have demonstrated in Section 4 some of the range of models that can be analysed by our algorithm. For this data, the ability to produce draws from the joint posterior distribution of the changepoint positions, despite there being around 50 changepoints, enables us to see some inadequacies in an existing model. We were able to use our algorithm to analyse a range of more complicated models, including a state-space model which appears more appropriate for the data. For this state-space model it was not possible to integrate out analytically all the parameters associated with each segment - however as there was only a single univariate parameter that could not be integrated out analytically, direct simulation was still possible using a simple numerical integration method.

Finally we have shown how approximations to the set of recursions can be used which greatly reduce the computational expense (particularly for large data sets with lots of changepoints), but with negligible error. In the well-log example the computational cost was reduced by an order of magnitude. We imagine that the computational cost of this approximate algorithm will increase only linearly with the size of data, and thus the algorithm could be used for analysing very large data sets.

**Acknowledgements** I would like to thank Peter Green and two anonymous referees



for helpful comments on an earlier version of this paper; and to thank Bill Fitzgerald for sending me the Well-log data.

## Appendix

To simulate in linear time a sample of size  $n$  from a discrete distribution  $\Pr(\tau)$ , which takes values of  $\tau = 1, 2, \dots$ :

- 1(a) for  $i = 1, \dots, n+1$ , simulate  $x_i$  a realisation from an exponential distribution with rate parameter 1;
- 1(b) Calculate  $S = \sum_{i=1}^{n+1} x_i$ ;
- 1(c) Set  $u_1 = x_1/S$  and for  $i = 2, \dots, n$   $u_i = u_{i-1} + x_i/S$ .
- 2 Set  $Q = 0$ ,  $U = u_1$ ,  $j = 1$  and  $i = 1$ .
- 3 If  $U < Q + \Pr(\tau = j)$  then output  $j$  and set  $U = U + u_{i+1}$  and  $i = i + 1$ ; otherwise set  $Q = \Pr(\tau = j)$  and  $j = j + 1$ . Repeat until  $i = n + 1$ .

## References

- Albert, J. H. and Chib, S. (1993). Bayes inference via Gibbs sampling of autoregressive time series subject to Markov mean and variance shifts. *Journal of Business and Economic Statistics* **11**, 1–15.
- Barry, D. and Hartigan, J. A. (1992). Product partition models for change point problems. *The Annals of Statistics* **20**, 260–279.
- Barry, D. and Hartigan, J. A. (1993). A Bayesian analysis for change point problems. *Journal of the American Statistical Society* **88**, 309–319.
- Braun, J. V. and Muller, H. G. (1998). Statistical methods for DNA sequence segmentation. *Statistical Science* **13**, 142–162.
- Braun, J. V., Braun, R. K. and Muller, H. G. (2000). Multiple changepoint fitting via quaslikelihood, with application to DNA sequence segmentation. *Biometrika* **87**, 301–314.
- Brooks, S. P., Giudici, P. and Roberts, G. O. (2003). Efficient construction of reversible jump Markov chain Monte Carlo proposal distributions. *Journal of the Royal Statistical Society, series B* **65**, 3–39.

- Carlin, B. P., Gelfand, A. E. and Smith, A. F. M. (1992). Hierarchical Bayesian analysis of changepoint problems. *Applied Statistics* **41**, 389–405.
- Carpenter, J., Clifford, P. and Fearnhead, P. (1999). An improved particle filter for non-linear problems. *IEE proceedings-Radar, Sonar and Navigation* **146**, 2–7.
- Chen, J. and Gupta, A. K. (1997). Testing and locating changepoints with application to stock prices. *Journal of the American Statistical Association* **92**, 739–747.
- Fearnhead, P. (1998). *Sequential Monte Carlo methods in filter theory*. Ph.D. thesis, Oxford University, available from <http://www.maths.lancs.ac.uk/~fearnhea/>.
- Fearnhead, P. (2004). Direct simulation for mixture distributions: component weights and discrete distributions. *submitted* Available from <http://www.maths.lancs.ac.uk/~fearnhea/>.
- Fearnhead, P. and Clifford, P. (2003). Online inference for hidden Markov models. *Journal of the Royal Statistical Society, Series B* **65**, 887–899.
- Fearnhead, P. and Meligkotsidou, L. (2004). Exact filtering for partially-observed continuous-time Markov models. *Journal of the Royal Statistical Society, series B. To appear* Available from <http://www.maths.lancs.ac.uk/~fearnhea/>.
- Green, P. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82**, 711–732.
- Green, P. J. (2003). Trans-dimensional Markov chain Monte Carlo. In: *Highly Structured Stochastic Systems* (eds. P. J. Green, N. L. Hjort and S. Richardson), Oxford University Press.
- Hartigan, J. A. (1990). Partition models. *Communications in Statistics* **19**, 2745–2756.
- Harvey, A. C. (1989). *Forecasting, structural time series and the Kalman filter*. Cambridge University Press, Cambridge, UK.
- Johnson, T. D., Elashoff, R. M. and Harkema, S. J. (2003). A Bayesian change-point analysis of electromyographic data: detecting muscle activation patterns and associated applications. *Biostatistics* **4**, 143–164.
- Liu, J. S. (2001). *Monte Carlo strategies in scientific computing*. New York: Springer.

- Lund, R. and Reeves, J. (2002). Detection of undocumented changepoints: A revision of the two-phase regression model. *Journal of Climate* **15**, 2547–2554.
- Ó Ruanaidh, J. J. K. and Fitzgerald, W. J. (1996). *Numerical Bayesian Methods Applied to Signal Processing*. New York: Springer.
- Pievatolo, A. and Green, P. J. (1998). Boundary detection through dynamic polygons. *Journal of the Royal Statistical Society, Series B* **60**, 609–626.
- Propp, J. G. and Wilson, D. B. (1996). Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures and Algorithms* **9**, 223–252.
- Punskaya, E., Andrieu, C., Doucet, A. and Fitzgerald, W. J. (2002). Bayesian curve fitting using MCMC with applications to signal segmentation. *IEEE Transactions on Signal Processing* **50**, 747–758.
- Raftery, A. E. and Akman, V. E. (1986). Bayesian analysis of a Poisson process with a change-point. *Biometrika* **73**, 85–89.
- Ritov, Y., Raz, A. and Bergman, H. (2002). Detection of onset of neuronal activity by allowing for heterogeneity in the change points. *Journal of Neuroscience Methods* **122**, 25–42.
- Scott, S. L. (2002). Bayesian methods for hidden Markov models: Recursive computing in the 21st century. *Journal of the American Statistical Association* **97**, 337–351.
- Worsley, K. J. (1979). On the likelihood ratio test for a shift in location of normal populations. *Journal of the American Statistical Association* **74**, 363–367.
- Yang, T. Y. and Kuo, L. (2001). Bayesian binary segmentation procedure for a Poisson process with multiple changepoints. *Journal of Computational and Graphical Statistics* **10**, 772–785.
- Yao, Y. (1984). Estimation of a noisy discrete-time step function: Bayes and empirical Bayes approaches. *The Annals of Statistics* **12**, 1434–1447.