

On-line Inference for Multiple Change Points Problems

Paul Fearnhead and Zhen Liu

Department of Mathematics and Statistics

Fylde College, Lancaster University, UK

Abstract

We propose an on-line algorithm for exact filtering of multiple changepoint problems. This algorithm enables simulation from the true joint posterior distribution of the number and position of the changepoints for a class of changepoint models. The computational cost of this exact algorithm is quadratic in the number of observations. We further show how resampling ideas from particle filters can be used to reduce the computational cost to linear in the number of observations, at the expense of introducing small errors; and propose two new, optimum resampling algorithms for this problem. One, a version of rejection control, allows the particle filter to automatically choose the number of particles required at each time-step. The new resampling algorithms substantially out-perform standard resampling algorithms on examples we consider; and we demonstrate how the resulting particle filter is practicable for segmentation of human GC content.

KEY WORDS AND PHRASES : Direct simulation; Isochores; Rejection control; Sequential Monte Carlo; Stratified sampling; Particle Filtering.

1 Introduction

Changepoint models are commonly used to model heterogeneity of data (usually over time). Given a set of observations collected over time, these models introduce a (potentially random) number of changepoints which split the data into a set of disjoint segments. It is then assumed that the data arise from a single model within each segment, but with different models across the segments. Examples of changepoint problems include Poisson processes with changing intensity (Ritov et al., 2002), changing linear regressions (Lund and Reeves, 2002), Gaussian models with changing variance (Johnson et al., 2003) and Markov models with time-varying transition matrices (Braun and Muller, 1998). These models have been applied to problems in a range of areas, including engineering, physical and biological sciences and finance.

We consider Bayesian inference for changepoint models where the number of changepoints is unknown. The most common approach to such inference for these mod-

els is the use of Markov chain Monte Carlo (MCMC; see for example Chib, 1998; Stephens, 1994), and reversible jump MCMC (Green, 1995). However, for many changepoint models (for example the models in Green, 1995; Punsakaya et al., 2002) it is possible to simulate independent realisations directly from the posterior distribution. This idea was used for DNA segmentation by Liu and Lawrence (1999), and has been proposed more generally by Fearnhead (2006) and Fearnhead (2005). The idea for direct simulation are based on exact methods for calculating posterior means (Barry and Hartigan, 1992). The advantages of direct simulation methods over MCMC and reversible jump MCMC are that (i) there is no need to diagnose whether the MCMC algorithm has converged; and (ii) as the draws from the posterior distribution are independent it is straightforward to quantify uncertainty in estimates of features of the posterior distributions based on them. For examples of the potential difficulties with MCMC caused by (i), compare the inferences obtained for the Coal-mining disaster data analysed in Green (1995) with those based on the direct simulation method of Fearnhead (2006).

In this paper we extend the direct simulation algorithms to on-line problems; where the data is obtained incrementally over time, and new inferences are required each time an observation is made. The use of on-line algorithms has also been suggested for static problems (e.g. Chopin, 2002; Del Moral et al., 2006).

The computational cost of our exact on-line algorithm increases linearly over time, however the on-line version of direct simulation is similar to particle filter algorithms, and we consider using resampling algorithms taken from particle filters to reduce the computational cost of our direct simulation algorithm (at the expense of introducing error). We propose two simple extensions of existing resampling algorithms that are particularly well-suited to changepoint models. One is a stratified extension of the rejection-control approach of Liu et al. (1998), which can limit the maximum amount of error introduced by each resampling step. In simulation studies we find this stratified version can reduce the error of the resulting algorithm by about one third as compared to the non-stratified version.

The resulting online algorithm can be viewed as a Rao-Blackwellised version of the Particle Filter algorithm of Chopin (2006): we have integrated out the parameters

associated with each segment. Note that this is an extremely efficient version of Rao-Blackwellisation. For example, consider analysing n data points. Our algorithm with n particles will give exact inference and thus will always outperform the algorithm of Chopin (2006) regardless of the number of particles used in that particle filter. Note however, that the filter of Chopin (2006) can be used more widely, as it does not require that the parameters within each segment can be integrated out.

The outline of the paper is as follows. We introduce the class of changepoint models we consider in Section 2. In Section 3 we introduce the online direct simulation algorithm, and approximate versions of it that utilise various resampling ideas. We test these approximate algorithms, and compare different resampling algorithms, on simulated data in Section 4 before applying our algorithm to the analysis of the C+G structure of human DNA data (Section 5). The paper concludes with a discussion.

2 Models and Notations

Assume we have data $\mathbf{y}_{1:n} = (y_1, y_2, \dots, y_n)$. We consider changepoint models for the data with the following conditional independence property: given the position of a changepoint, the data before that changepoint is independent of the data after the changepoint. These models can be described in terms of the following hierarchical structure.

Firstly we model the changepoint positions via a Markov process. This Markov process is determined by a set of transition probabilities,

$$\Pr(\text{next changepoint at } t | \text{changepoint at } s). \quad (1)$$

For this paper we make the simplification that these transition probabilities depend only on the distance between the two changepoints. Extending our work to the general case, where the distribution of the length of a segment could depend on the time at which it starts, is straightforward. Specifically we let $g(\cdot)$ be the probability mass function for the distance between two successive changepoints (equivalently the length of segments); so that (1) is $g(t - s)$. We further let $G(l) = \sum_{i=1}^l g(i)$ be the distribution function of this distance, and assume that $g(\cdot)$ is the probability mass function for the position of the first changepoint.

Note that any such model implies a prior distribution on the number of changepoints. For example if a geometric distribution is used for $g(\cdot)$, then our model implies that there is a fixed probability of a changepoint at any time-point, independent of other changepoints. Hence this model implies a binomial distribution for the number of changepoints.

Now we condition on m changepoints at times $\tau_1, \tau_2, \dots, \tau_m$. We let $\tau_0 = 0$ and $\tau_{m+1} = n$, so our changepoints define $m + 1$ segments, with segment i consisting of observations $\mathbf{y}_{\tau_{i+1}:\tau_{i+1}}$ for $i = 0, \dots, m$. We allow a set of \bar{p} possible models for the data from each segment, labeled $\{1, 2, \dots, \bar{p}\}$, and assume an arbitrary prior distribution for models, common across segments, with the model in a given segment being independent of the models in all other segments.

For a segment consisting of observations $\mathbf{y}_{s+1:t}$ and model q we will have a set of unknown parameters, β say. We have a prior distribution, $\pi(\beta)$ for β (which may depend on q), but assume that the parameters for this segment are independent of the parameters in other segments. Finally we define

$$P(s, t, q) = \int \Pr(\mathbf{y}_{s+1:t} | \beta, \text{model } q) \pi(\beta) d\beta, \quad (2)$$

and assume that these probabilities can be calculated for all $s < t$ and q . This requires either conjugate priors for β , or the use of numerical integration. Numerical integration is often computationally feasible in practice if the dimension of the part of β that cannot be integrated analytically is low (see Fearnhead, 2006, for an example).

For concreteness we describe a specific example of this model which we will use in Section 4 (see also Punskeya et al., 2002; Fearnhead, 2005). Here we model the data as piecewise linear regressions. So within a segment we have a linear regression model of unknown order. For a given model order q , we have

$$\mathbf{y}_{s+1:t} = \mathbf{H}\beta + \epsilon \quad (3)$$

where \mathbf{H} is a $(t - s) \times q$ matrix of basis functions, β is a vector of q regression parameters and ϵ is a vector of iid Gaussian noise with mean 0 and variance σ^2 .

We assume conjugate priors. The variance of the Gaussian noise has an inverse gamma distribution with parameters $\nu/2$ and $\gamma/2$, and the components of the re-

gression vector have independent Gaussian priors. The prior for the j th component is Gaussian with mean 0 and variance $\sigma^2\delta_j^2$.

The likelihood function of $\mathbf{y}_{s+1:t}$ conditional on a model q is obtained by integrating out the regression parameters β and variance σ^2 :

$$P(s, t, q) = \pi^{-(t-s)/2} \left(\frac{|\mathbf{M}|}{|\mathbf{D}|} \right)^{\frac{1}{2}} \frac{(\gamma)^{\nu/2}}{(\|\mathbf{y}_{s+1:t}\|_{\mathbf{P}}^2 + \gamma)^{(t-s+\nu)/2}} \frac{\Gamma((t-s+\nu)/2)}{\Gamma(\nu/2)}, \quad (4)$$

where $\mathbf{M} = (\mathbf{H}^T\mathbf{H} + \mathbf{D}^{-1})^{-1}$, $\mathbf{P} = (\mathbf{I} - \mathbf{H}\mathbf{M}\mathbf{H}^T)$, $\|\mathbf{y}\|_{\mathbf{P}}^2 = \mathbf{y}^T\mathbf{P}\mathbf{y}$, $\mathbf{D} = \text{diag}(\delta_1^2, \dots, \delta_q^2)$ and \mathbf{I} is a $(t-s) \times (t-s)$ identity matrix.

3 On-line Inference

We consider on-line inference for the multiple changepoint model of Section 2. We assume that observations accrue over time, so that y_t is the observation at time t . At each time-step, our aim is to calculate the posterior distributions of interest based on all the observations to date. To do this efficiently requires updating our posterior distributions at the previous time-step to take account of the new observation. Note that on-line algorithms can be used to analyse batch data by introducing an artificial time for each observation.

We focus on on-line inference of the position of the changepoints. Under the modeling assumptions of Section 2, inference for the parameters conditional on knowing the number and position of the changepoints is straightforward. We first describe an exact on-line algorithm, which is an on-line version of the direct simulation method of Fearnhead (2005). The computational cost of this exact algorithm increases over time, so we then present an approximate on-line algorithm, which uses resampling ideas from particle filters, and which has constant computational cost over time.

3.1 Exact On-line Inference

We introduce a state at time t , C_t , which is defined to be the time of the most recent change-point prior to t (with $C_t = 0$ if there have been no change-points before time t). Initially we focus on calculating the posterior distribution for C_t

given the observation $\mathbf{y}_{1:t}$. We then describe how, if these distributions are stored for all t , it is straightforward to simulate from the joint posterior distribution of the position of all changepoints prior to the current time.

Filtering Recursions

The state C_t can take values in $0, 1, \dots, t-1$, and $C_1, C_2, \dots, C_t, \dots$ is a Markov chain. Conditional on $C_t = j$, either $C_{t+1} = j$, which corresponds to no changepoint at time t , or $C_{t+1} = t$, if there is a changepoint at time t . The transition probabilities for this Markov chain can thus be calculated as:

$$\Pr(C_{t+1} = j | C_t = i) = \begin{cases} \frac{1-G(t-i)}{1-G(t-i-1)} & \text{if } j = i, \\ \frac{G(t-i)-G(t-i-1)}{1-G(t-i-1)} & \text{if } j = t, \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

where $G(\cdot)$ is the distribution function of distance between two successive change points.

Now, standard filtering recursions give

$$\Pr(C_{t+1} = j | \mathbf{y}_{1:t+1}) \propto \Pr(y_{t+1} | C_{t+1} = j, \mathbf{y}_{1:t}) \Pr(C_{t+1} = j | \mathbf{y}_{1:t}),$$

and

$$\Pr(C_{t+1} = j | \mathbf{y}_{1:t}) = \sum_{i=0}^{t-1} \Pr(C_{t+1} = j | C_t = i) \Pr(C_t = i | \mathbf{y}_{1:t}).$$

Thus, if we let $w_{t+1}^{(j)} = \Pr(y_{t+1} | C_{t+1} = j, \mathbf{y}_{1:t})$, and substitute in the transition probabilities (5), we obtain

$$\Pr(C_{t+1} = j | \mathbf{y}_{1:t+1}) \propto \begin{cases} w_{t+1}^{(j)} \frac{1-G(t-j)}{1-G(t-j-1)} \Pr(C_t = j | \mathbf{y}_{1:t}) & \text{if } j < t, \\ w_{t+1}^{(t)} \sum_{i=0}^{t-1} \left(\frac{G(t-i)-G(t-i-1)}{1-G(t-i-1)} \Pr(C_t = i | \mathbf{y}_{1:t}) \right) & \text{if } j = t. \end{cases}$$

If we define $P(s, t, q)$ as in (2) then we get for $j < t$

$$w_{t+1}^{(j)} = \frac{\sum_{q=1}^{\bar{p}} P(j, t+1, q) p(q)}{\sum_{q=1}^{\bar{p}} P(j, t, q) p(q)}, \quad (6)$$

while if $j = t$ then the weight is $\sum_{q=1}^{\bar{p}} P(t, t+1, q) p(q)$.

In most situations, such as for the linear regression models described in Section 2, the incremental weights $w_{t+1}^{(j)}$ can be calculated efficiently, as each $P(j, t, q)$ depends on a

set of summary statistics of the observations $y_{j+1:t}$, which can be updated recursively. Efficient calculation of the incremental weights $w_{t+1}^{(j)}$ is possible by storing these summary statistics for each set of values for the time of the last changepoint, j , and the model for the current segment, q . These can be updated to give the summary statistics required for each $P(j, t+1, q)$. Thus the computational cost of calculating each $w_{t+1}^{(j)}$ is fixed, and does not increase with $t - j$.

Simulating Changepoints

If we store the filtering densities $\Pr(C_t | \mathbf{y}_{1:t})$ for all $t = 1, \dots, n$, then it is straightforward to simulate from the joint posterior distribution of the position of all changepoints prior to time n , using the idea of Chopin (2006). To simulate one realisation from this joint density:

- (1) Set $t_0 = n$, and $k = 0$
- (2) Simulate t_{k+1} from the filtering density $\Pr(C_{t_k} | \mathbf{y}_{1:t_k})$, and set $k = k + 1$.
- (3) If $t_k > 0$ return to (2); otherwise output the set of simulated changepoints, $t_{k-1}, t_{k-2}, \dots, t_1$.

A simple extension of this algorithm which enables a large sample of realisations of sets of changepoints to be simulated efficiently is described in Fearnhead (2006).

MAP estimation

We can obtain an on-line Viterbi algorithm for calculating the maximum a posteriori (MAP) estimate of the positions of the changepoints and the model orders for each segment as follows. We define \mathcal{M}_j to be the event that given a changepoint at time j , the MAP choice of changepoints and model occurs prior to time j . For $t = 1, \dots, n$, $j = 0, \dots, t - 1$ and $q = 1, \dots, \bar{p}$,

$$P_t(j, q) = \Pr(C_t = j, \text{model } q, \mathcal{M}_j, \mathbf{y}_{1:t}), \text{ and}$$

$$P_t^{MAP} = \Pr(\text{Changepoint at } t, \mathcal{M}_t, \mathbf{y}_{1:t}).$$

We obtain the following equations

$$P_t(j, q) = (1 - G(t - j - 1))P(j, t, q)p(q)P_j^{MAP}, \text{ and}$$

$$P_t^{MAP} = \max_{j,q} \{P_t(j, q)g(t-j)/(1-G(t-j-1))\}. \quad (7)$$

At time t , the MAP estimates of C_t and the current model order are given respectively by the values of j and q which maximise $P_t(j, q)$. Given a MAP estimate of C_t , \hat{c}_t we can then calculate the MAP estimates of the changepoint prior to \hat{c}_t and the model order of that segment by the values of j and q that maximised the right-hand side of (7). This procedure can be repeated to find the MAP estimates of all changepoint positions and model orders.

3.2 Approximate Inference

The computational and memory costs of the recursions for exact inference presented in Section 3.1 both increase with time. The computational cost of both the filtering recursion and MAP recursion at time t is proportional to t , the number of possible values of C_t . While the memory cost of storing all filtering densities up to time t , necessary to simulate from the joint posterior of all changepoints prior to t , increases quadratically with t . For large data sets, these computational and memory costs may become prohibitive.

A similar problem of increasing computational cost occurs in the analysis of some hidden Markov models – though generally computational cost increases exponentially with time (Chen and Liu, 2000). Particle filters have been successfully applied to these problems (Fearnhead and Clifford, 2003) by using a resampling step to limit the computational cost at each time-step. Here we show how similar resampling ideas can be applied to the online inference of the changepoint models we are considering. We present a variation on the optimal resampling method of Fearnhead and Clifford (2003) which is specifically designed for changepoint models, and show theoretically why this is an optimal resampling algorithm in this case. We also present an extension of the rejection control approach of Liu et al. (1998) which is suitable for the analysis of batch data, and for which it is possible to control the amount of error introduced at each resampling step.

Controlling Computational Cost

Our first approach is to control the average (and maximum) computational cost for

analysing each new observation. At time t our exact algorithm stores the complete posterior distribution of the time of the last changepoint $\Pr(C_t = c_t | \mathbf{y}_{1:t})$, for $c_t = 0, 1, \dots, t - 1$. We can approximate this by a discrete distribution with fewer, N , support points. This approximate distribution can be described by the set of support points, $c^{(1)}, \dots, c^{(N)}$, henceforth called *particles*, and the probability mass associated with each of these particles, $w^{(1)}, \dots, w^{(N)}$, which we call weights. (The particles and their weights will depend on t ; we have suppressed this dependence to simplify notation.)

We impose a maximum number of particles to be stored at any one time, N , such that whenever we have N particles we immediately perform resampling to reduce the number of particles to $M < N$. The average computational cost per iteration will thus be proportional to $(M + N + 1)/2$, and the maximum computational cost per iteration proportional to N .

Assume that at the current time point we have N particles; and wish to reduce these to M particles. We propose the following stratified version of the optimal resampling algorithm of Fearnhead and Clifford (2003), which we call Stratified Optimal Resampling (SOR).

Initialisation Assume we currently have a set of ordered particles $c^{(1)} < c^{(2)} < \dots < c^{(N)}$, with associated weights $w^{(1)}, \dots, w^{(N)}$, which sum to unity.

(SOR1) Calculate α the unique solution to $\sum_{i=1}^N \min\{1, w^{(i)}/\alpha\} = M$;

(SOR2) For $i = 1, \dots, N$ if $w^{(i)} \geq \alpha$ then keep particle $c^{(i)}$ with weight $w^{(i)}$. Assume that A particles are kept.

(SOR3) Use the stratified resampling algorithm of Carpenter et al. (1999) to resample $M - A$ times from the ordered set of the remaining $N - A$ particles (without shuffling). Each resampled particle is assigned a weight α .

The stratified resampling algorithm used in step SOR3 is given in Appendix A. The use of stratified resampling in step SOR3 means that at most one copy of each particle is kept, as the expected number of times a particle with weight w is resampled is $w/\alpha < 1$ (note there is no advantage in having multiple copies

of particles, see Fearnhead and Clifford, 2003). The only difference between this SOR algorithm and the original algorithm of Fearnhead and Clifford (2003) is that particles are ordered before resampling in step SOR3.

As shown in Fearnhead and Clifford (2003), if we denote by $W^{(i)}$ the (random) weight of a particle after resampling (so $W^{(i)} = w^{(i)}$, α , or 0 depending on whether the respective particle is kept, resampled or not resampled), then SOR is optimal over all resampling algorithms that satisfy $E(W^{(i)}) = w^{(i)}$ in terms of minimising the mean square error: $E(\sum_{i=1}^N (W^{(i)} - w^{(i)})^2)$. By ordering the particles in step SOR3 we obtain the further property:

Theorem 3.1 *Consider a set of N particles, $c^{(1)} < c^{(2)} < \dots < c^{(N)}$ with weights $w^{(1)}, \dots, w^{(N)}$. Let $W^{(i)}$ be the (random) weight of particle $c^{(i)}$ after resampling. Define the maximum Kolmogorov Smirnov Distance for a resampling algorithm as*

$$mKSD = \max \left\{ \max_i \left| \sum_{j=1}^i (w^{(j)} - W^{(j)}) \right| \right\} \quad (8)$$

where the first maximisation is over realisations of $W^{(1)}, \dots, W^{(N)}$ with positive probability. Then the SOR algorithm above satisfies $mKSD \leq \alpha$ (where α is defined as in SOR1). Furthermore (i) for a resampling algorithm to have $mKSD \leq \alpha$ then all particles with $w^{(i)} > \alpha$ must be propagated without resampling; and (ii) the $mKSD$ for the SOR algorithm above is less than or equal to the $mKSD$ of the optimal resampling algorithm of Fearnhead and Clifford (2003), and the rejection control algorithm of Liu et al. (1998).

Proof: See Appendix B. □

Kolmogorov Smirnov distance is a natural metric for the distributions of 1-dimensional random variables. By using (8) as a measure of error of a resampling algorithm, we are considering the bound on Kolmogorov Smirnov distance that a resampling algorithm can introduce. The theorem gives a simple interpretation of the α calculated in step SOR1; in terms of an upper bound on the Kolmogorov Smirnov distance between the original and resampled weights.

We define a resampling algorithm to be unbiased if $E(W^{(i)}) = w^{(i)}$ for all i . (This is related to the properly weighted condition of Liu et al., 2001). The optimal resam-

pling algorithm of Fearnhead and Clifford (2003) and rejection control are currently the only other unbiased resampling algorithm which satisfy the condition (i) of Theorem 1. (Note that rejection control will not produce a fixed number of particles after resampling; though implementing rejection control with a threshold of α will produce on average N resampling particles, and further that while $E(W^{(i)}) = w^{(i)}$, the resampled weights do not necessarily sum to 1.) So results (i) and (ii) of Theorem 1 show that our SOR algorithm is optimal over all existing unbiased resampling algorithms in terms of minimising mKSD.

We have presented the SOR algorithm in terms of the general case of resampling M particles from N current particles. For on-line inference, where there is a fixed amount of time to analyse each observation, it is natural to choose N to be the largest number of particles that enable an observation to be analysed in less than this amount of time; and to set $M = N - 1$. In this case there is no difference between SOR and the existing optimal resampling algorithm. In practice, it may be better to choose $N - M > 1$ (see Section 4) as this enables the particles to be removed to be jointly chosen in a stratified manner.

Controlling Resampling Error

An alternative to basing resampling on the average and maximum number of particles to be kept at each time step, is to choose the amount of resampling to control the size of error that is introduced at each time step. Such an approach is most suitable for using online algorithms to analyse batch data. For real-time data, the frequency of observations will place an upper bound on the CPU time, and hence the number of particles, that can be used to process each observation. By controlling the resampling error, rather than the number of particles, we cannot ensure that the number of particles always stays below this error.

The idea behind controlling the resampling error is given by the interpretation of α for SOR that comes from Theorem 1. The value of α defines the maximum error (as defined by Kolmogorov Smirnov distance) that is introduced by the resampling error. So rather than specifying the number of resampled particles which in turn defines α , and hence the amount of error we introduce, we can instead specify α which will then define the number of resampled particles.

Our method for controlling the resampling error is to use a stratified version of rejection control (Liu et al., 1998), rather than adapt the SOR algorithm. For a prespecified value of α , our stratified rejection control (SRC) algorithm is:

Initialisation Assume we currently have a set of ordered particles $c^{(1)} < c^{(2)} < \dots < c^{(N)}$, with associated weights $w^{(1)}, \dots, w^{(N)}$, which sum to unity.

(SRC1) For $i = 1, \dots, N$ if $w^{(i)} \geq \alpha$ then keep particle $c^{(i)}$ with weight $w^{(i)}$. Assume that A particles are kept.

(SRC2) Use the stratified resampling algorithm of Carpenter et al. (1999) to resample from the ordered set of the remaining $N - A$ particles (without shuffling). The expected number of times particle $c^{(i)}$ is resampled is $w^{(i)}/\alpha$. Each resampled particle is assigned a weight α .

Again the use of stratified resampling in (SRC2) means that at most one copy of each particle is kept. Note that the sum of the particles' weights after resampling will not necessarily sum to 1 (though they lie between $1 - \alpha$ and $1 + \alpha$), and should be normalised to produce a probability distribution.

The difference between SRC and rejection control (Liu et al., 1998) is that particles are ordered and stratified resampling is used in step SRC2, as opposed to independent resampling of each particles. The use of stratified resampling means that the maximum error of the unnormalised weights introduced by SRC, as measured by Kolmogorov Smirnov distance, is α (this can be proved in an identical manner to Theorem 1). Furthermore, the error of the normalised weights can be shown to be bounded above by $\alpha/(1 - \alpha) = \alpha + o(\alpha)$ (see Appendix C).

4 Numerical Examples

We tested our algorithm on three simulated examples: the Blocks and Heavisine examples from Donoho and Johnstone (1994) and a piecewise auto-regressive model. Each of the three data-sets are analysed under a piecewise regression model. The design matrices for the piecewise AR model and the piecewise polynomial regression

model (for the Blocks and Heavisine data) are

$$\mathbf{H}_{s:t} = \begin{pmatrix} y_{s-1} & y_{s-2} & y_{s-3} \\ y_s & y_{s-1} & y_{s-2} \\ \vdots & \vdots & \vdots \\ y_t & y_{t-1} & y_{t-2} \end{pmatrix}, \text{ and } \mathbf{H}_{s:t} = \begin{pmatrix} 1 & x_s & x_s^2 \\ 1 & x_{s+1} & x_{s+1}^2 \\ \vdots & \vdots & \vdots \\ 1 & x_t & x_t^2 \end{pmatrix}.$$

respectively, where $x_s = s/n$ and n is the number of data points. In each case we perform model choice within each segment, choosing between model orders 1, 2 and 3. We allow for different variances of the measurement error within each segment, although for the Blocks and Heavisine examples we simulated data with a common error variance across segments. Further details of the model, and calculations required for calculating the $P(s, t, q)$ s is given in Section 2 (See also Punskeya et al., 2002; Fearnhead, 2005).

Our focus here is on the performance of different possible resampling algorithms. The Blocks data set (see Figure 1) is a particularly simple data set to analyse, and all reasonable resampling algorithms will give almost identical results. We show the results here to demonstrate how the SRC algorithm naturally adapts the number of particles that are kept. The Blocks data set has a number of obvious changepoints, and when each of these are encountered the number of particles that are needed to be kept is reduced to close to 1.

For the Heavisine example (see Figure 1) we compared the accuracy of various resampling algorithms: stratified rejection control (SRC), rejection control (RC), stratified optimal resampling (SOR), and optimal resampling (OR). We considered two values of α for SRC and RC; and for a meaningful comparison, fixed the mean number of particles in OR and SOR to the mean number of particles kept by SRC for each of these two values. If we set the number of resampled particles (M) to be one less than the number of particles prior to resampling (N), then OR and SOR are identical. We tested both $N = M + 1$ and $N = M + 5$.

Our comparison is based on the Kolmogorov Smirnov distance (KSD) between the true filtering distribution of the most recent changepoint, $p(C_t|y_{1:t})$ (calculated using the online algorithm with no resampling), and its approximation based on the various resampling algorithms, for each t . Results are given in Table 1. The results show

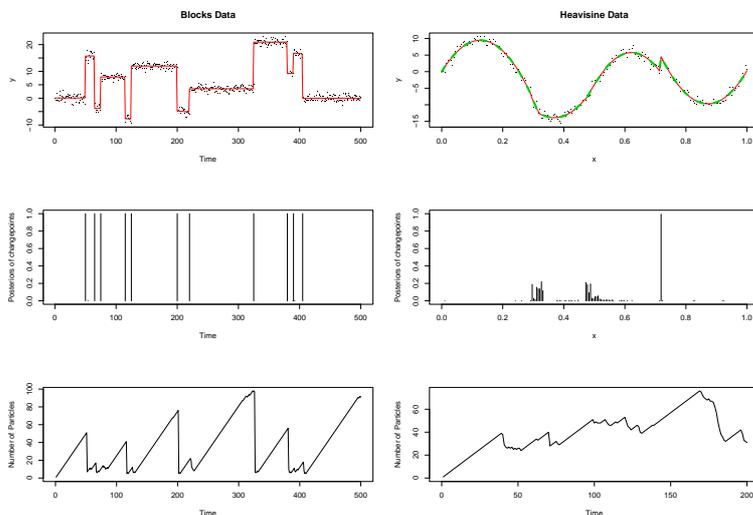


Figure 1: The Blocks data set (left-hand column) and Heavisine data set (right-hand column) together with results of analysis by the SRC algorithm with $\alpha = 10^{-6}$: data and inferred signal (top); marginal probability of changepoints (middle); and numbers of particles kept (bottom).

that the mean KSD error is reduced by one third by using SRC rather than RC. Both of these methods perform better than the resampling algorithms that use a fixed number of particles (for the same average number of particles); showing the advantage of allowing the number of particles used to adapt to the filtering density being approximated. Of the two algorithms considered which use a fixed number of particles, we see an improvement of using SOR where we remove 5 particles at each resampling step over OR (or equivalently SOR) where 1 particle is removed at each resampling step. By removing many particles in one step, SOR is able to jointly choose the particles to remove in a stratified way so as to reduce the error introduced. (Note OR where we remove 5 particles at each resampling step has worse results than the OR results shown in Table 1.)

Note that while all resampling algorithms introduce small errors at each resampling step, it is possible for these errors to accumulate. The reason for this, appears to be that the evidence for a changepoint at a given time t can change substantially as more data is collected. If the evidence is small (and hence the filtering probability of a changepoint at t is less than α) at a resampling step, this can lead to the

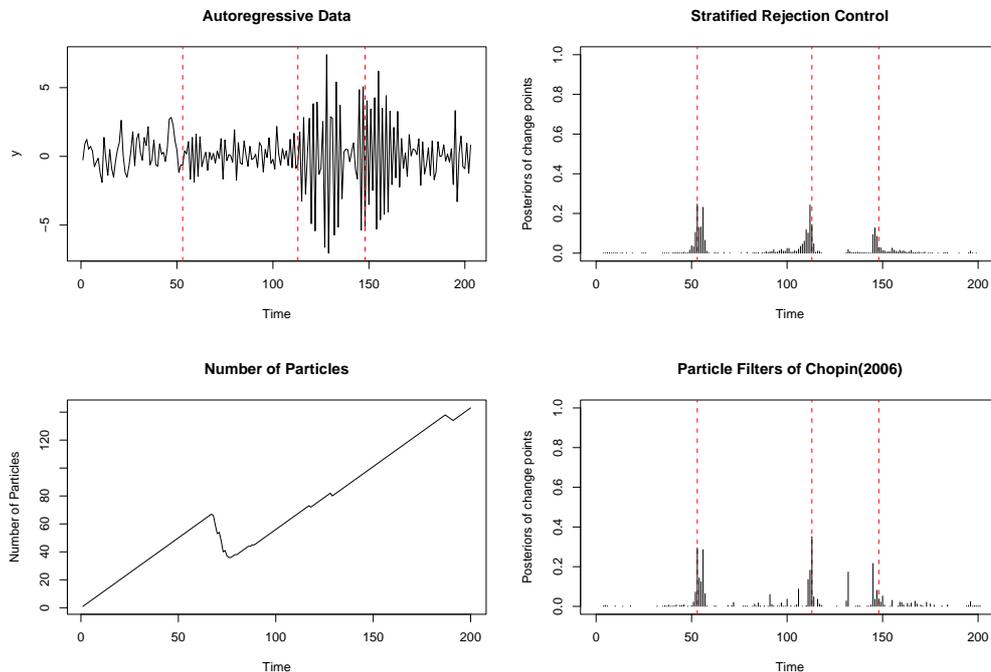


Figure 2: Results of analysing the AR dataset using SRC with $\alpha = 10^{-6}$, and the particle filter of Chopin (2006) with 50,000 particles: data (top left), marginal probabilities of changepoint for SRC (top right) and particle filter of Chopin (2006) (bottom right), and number of particles kept using SRC (bottom left). The true AR model to the four segments have model orders 1, 1, 2, and 3 respectively. The corresponding parameters are $\beta = 0.4$; $\beta = -0.6$; $\beta = (-1.3, -0.36, 0.25)$ and $\beta = (-1.1, -0.24)$ with error variances 1.2^2 , 0.7^2 , 1.3^2 and 0.9^2 respectively.

corresponding particle being removed. Such a particle can not be “resurrected” as future observations are made, even if they carry strong evidence for a changepoint at t . However stratified resampling should ensure that a particle corresponding to a changepoint close to t is kept, and thus the error in estimating the position of the changepoints will still be small.

We repeated this analysis for a piecewise AR model. The results of the SRC analysis with $\alpha = 1 \times 10^{-6}$ given in Figure 2, and results of the accuracy of each resampling method given in Table 1. We observe similar results to the Heavisine example in terms of the relative performance of the resampling algorithms. In this case SRC again outperforms RC by about a third. The difference in performance between

	SRC	RC	SOR	OR
Heavisine	1.3×10^{-2}	2.0×10^{-2}	4.2×10^{-2}	6.4×10^{-2}
AR	1.3×10^{-6}	2.2×10^{-6}	2.2×10^{-4}	3.5×10^{-4}

Table 1: Mean Kolmogorov Smirnov Distance in $P(C_t|y_{1:t})$ averaged over t for the Heavisine and AR models and the four resampling algorithms. Stratified Rejection Control (SRC) and Rejection Control (RC) were implemented with $\alpha = 10^{-6}$; these algorithms used an average number of 43 and 70 particles for the Heavisine and AR models respectively. Optimal Resampling (OR) was implemented with $N = M + 1 = 49$ and $N = M + 1 = 90$; Stratified Optimal Resampling (SOR) used $N = M + 5 = 51$ and $N = M + 5 = 92$ (chosen so that the average number of particles is the same for all algorithms for each data set). Results based on 50 replications of each algorithm for one version of each data set. The true distribution, $P(C_t|y_{1:t})$, was calculated using the exact online algorithm.

SRC and RC as compared to SOR and OR is quite substantial in this case, because towards the end of the time series it is forced to use too few particles to adequately approximate the filtering densities. This again demonstrates the potential gains to be obtained by allowing the number of particles used to change over time and to adapt to the filtering distribution that is being approximated.

We also ran the particle filter of Chopin (2006). This filter does not integrate out the parameters associated with each segment, so each particle consists of a time for the last changepoint together with a value of the parameters for the current segment. The filter uses MCMC to update the parameters of a subset of particles at each iteration. We ran the filter with 50,000 particles, using a Gibbs sampler update on the parameters of 1/3 of the particles at each iteration. This took over an order of magnitude longer to run than the SRC algorithm, and even is substantially more time-consuming to implement than the exact online algorithm.

The results for the estimate of the marginal probabilities of the changepoints is shown in Figure 2. The filter of Chopin (2006) suffers from a loss of diversity in the particles – with many positions being assigned zero probability of being a change-

point, when in fact there is a non-negligible probability as can be seen from the output of the SRC filter. To give a quantitative comparison of the two methods we calculated the mean absolute error between the estimates of the marginal probabilities of the changepoints shown in Figure 2 with those based on the exact particle filter algorithm. These were 0.010 and 0.002 for the filter of Chopin (2006) and the SRC filter respectively.

5 DNA Segmentation

In recent years there has been an explosion in the amount of data describing the genetic make-up of different organisms; for example the complete DNA sequence of one human genome is now known as a result of the Human Genome project. There is interest in learning about the genomic features of different organisms, and learning how these features may have evolved and how they correlate with each other.

We consider the problem of understanding the structure of C+G content within the genome. A common model for the C+G content of the human genome is that there are large, of the order of 300 kilobases (kb), regions of roughly homogeneous C+G content, called Isochores (see Bernardi, 2000, for background). Furthermore C+G content is known to correlate with various features of the genome, such as high recombination rates and gene density (Hardison et al., 2003).

Currently, the most common method for segmenting an organism's genome into regions of different C+G content is implemented in the computer program IsoFinder (Oliver et al., 2004). This is based on a recursive segmentation procedure, which initially classifies a large genomic region as consisting of a single Isochore (region of common C+G content). It then considers in turn each possible position for adding a changepoint, and splitting the data into two Isochores. For each possible position, a t -statistic is calculated for testing whether the mean C+G content is different in the two putative Isochores. For each changepoint, a p -value is calculated for its value of the t -statistic using a bootstrap procedure, and if the smallest p -value is less than some predefined threshold, then the corresponding changepoint is added. This procedure is repeated, with at each step each current Isochore being tested for

whether it can be split into two Isochores. See Oliver et al. (2004) for more details. We consider a Bayesian approach to segmenting a genomic region into Isochores. The potential advantages of a Bayesian approach include (i) quantifying and averaging over the uncertainty in the number and positions of the Isochores; (ii) jointly estimating all Isochore positions (which Braun et al., 2000, show to be more accurate than segmentation procedures); and (iii) the large amount of data available for each organism makes it straightforward to construct sensible prior distributions.

One of the computational challenges of such an analysis is the large amount of data that needs to be analysed (for example human chromosomes consist of around 100 million bases). We simplify this burden by first summarising our data by the number of DNA sites which are C or G within consecutive windows (each window being of the order of a few kb in width), an approach which also has the advantage of averaging out the very local high variation in C+G content caused for example by CpG islands and Alu elements. We then hope that our online changepoint algorithm will be able to efficiently analyse the resulting data, and one of the main aims of the study we present here is to test whether such an approach is computationally practicable for analysing the large amount of genomic data currently available.

The model we use is based on the following simple model for the data $y_{1:n}$, which is similar to the implicit model assumed by `IsoFinder`. A data set is shown in Figure 3. The t th data point, y_t , represents the number of DNA bases which are either C or G within the t th window. If this window lies within the i th Isochore then we assume

$$y_t = \mu_i + \sigma_i \epsilon_t,$$

where μ_i is the mean C+G content of each window within the i th Isochore, σ_i^2 is the error variance within the i th Isochore, and ϵ_t is some independent error. We assume that ϵ_t has a Gaussian distribution and we assume standard conjugate priors (see Section 2) for the μ_i s and σ_i s, with the prior parameters chosen from an initial analysis of C+G data with a moving median filter. For each model we assumed a geometric distribution for the length of each Isochore.

Results of our analysis using SRC with $\alpha = 10^{-6}$ are shown in Figure 3. Our main focus is on the computational practicability of a Bayesian analysis of such data, and

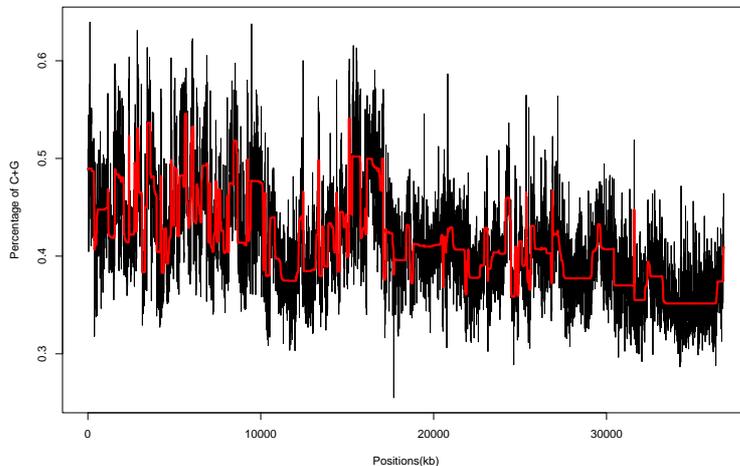


Figure 3: Analysis of 35Mb of data from human chromosome 1. The red line is the posterior mean GC content.

our method took 6 seconds on a desktop PC to analyse this data set.

This application does not need to be analysed by an online algorithm, such as the one we used. However Fearnhead (2006) showed that the version of our algorithm without resampling can be more efficient for analysing changepoint models than some commonly used MCMC algorithms. Furthermore, by using resampling we have been able to vastly reduce the computational and storage cost of analysing the data. For example our implementation with resampling uses an average of 117 particles at each time step; whereas without resampling the algorithm would require an average of over 3,500 particles for each time-step.

6 Discussions

We have considered a class of changepoint models, which have a specific conditional independence structure (see Section 2), and shown how the direct simulation algorithm of Fearnhead (2005) can be implemented online. Such an algorithm can be viewed as an exact particle filter, and resampling ideas taken from particle filters can be used to reduce the computational complexity of the direct simulation algorithm (at the cost of introducing error). We have presented two simple extensions

of existing resampling algorithms, which are particularly well suited to changepoint problems (or any problems where the underlying state of interest is 1-dimensional). In simulation studies, our new resampling algorithms decreased the error of the resulting particle filter by up to one third, compared to particle filters using the existing resampling approaches. We have shown that the new resampling algorithms satisfy a minimax optimality criteria on the error, as measured by Kolmogorov Smirnov distance, introduced by resampling. Furthermore this result gives a natural interpretation of the threshold that needs to be specified in the stratified rejection control algorithm which will aid its implementation in practice.

There is great flexibility with implementing resampling algorithms within particle filters which we have not explored. For example Liu and Chen (1998) discuss the frequency with which resampling should occur, and Liu et al. (1998) suggest using rejection control only when the variance of the particle filter weights exceeds some threshold. Whilst we have not fully investigated these issues, the results from Section 4 suggests that the advantages of using stratification within optimal resampling or rejection control will increase as the frequency of resampling decreases (or equivalently the amount of particles resampled increases at each resampling step).

Acknowledgements This work is supported by EPSRC grant C531558. We dedicate this paper to the memory of Nick Smith who helped with the application to detecting Isochores.

References

- Barry, D. and Hartigan, J. A. (1992). Product partition models for change point problems. *The Annals of Statistics*, 20:260–279.
- Bernardi, G. (2000). Isochores and evolutionary genomics of vertebrates. *Gene*, 241:3–17.
- Braun, J. V., Braun, R. K., and Muller, H. G. (2000). Multiple changepoint fitting via quasilielihood, with application to DNA sequence segmentation. *Biometrika*, 87:301–314.

- Braun, J. V. and Muller, H. G. (1998). Statistical methods for DNA sequence segmentation. *Statistical Science*, 13:142–162.
- 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, R. and Liu, J. (2000). Mixture Kalman filters. *Journal of the Royal Statistical Society, Series B*, 62:493–508.
- Chib, S. (1998). Estimation and comparison of multiple change-point models. *Journal of Econometrics*, 86:221–241.
- Chopin, N. (2002). A sequential particle filter method for static models. *Biometrika*, 89:539–551.
- Chopin, N. (2006). Dynamic detection of change points in long time series. *Annals of the Institute of Statistical Mathematics*, page to appear.
- Del Moral, P., Doucet, A., and Jasra, A. (2006). Sequential Monte Carlo samplers. *Journal of the Royal Statistical Society, Series B*, 68:411–436.
- Donoho, D. L. and Johnstone, I. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81:425–455.
- Fearnhead, P. (2005). Exact Bayesian curve fitting and signal segmentation. *IEEE Transactions on Signal Processing*, 53:2160–2166.
- Fearnhead, P. (2006). Exact and efficient inference for multiple changepoint problems. *Statistics and Computing*, 16:203–213.
- Fearnhead, P. and Clifford, P. (2003). Online inference for hidden Markov models. *Journal of the Royal Statistical Society, Series B*, 65:887–899.
- Green, P. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, 82:711–732.
- Hardison, R. C., Roskin, K. M., Yanf, S., Diekhans, M., Kent, W. J., Weber, R., Elnitski, L., and Li et al., J. (2003). Covariation in frequencies of substitution,

- deletion, transposition, and recombination during eutherian evolution. *Genome Research*, 13:13–26.
- 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. and Chen, R. (1998). Sequential Monte Carlo methods for dynamic systems. *Journal of the American Statistical Association.*, 93:1032–1044.
- Liu, J. S., Chen, R., and Logvinenko, T. (2001). A theoretical framework for sequential importance sampling with resampling. In Doucet, A., de Freitas, N., and Gordon, N., editors, *Sequential Monte Carlo Methods in Practice*, pages 225–246. Springer–Verlag; New York.
- Liu, J. S., Chen, R., and Wong, W. H. (1998). Rejection control and sequential importance sampling. *Journal of the American Statistical Society*, 93:1022–1031.
- Liu, J. S. and Lawrence, C. E. (1999). Bayesian inference on biopolymer models. *Bioinformatics*, 15:38–52.
- Lund, R. and Reeves, J. (2002). Detection of undocumented changepoints: A revision of the two-phase regression model. *Journal of Climate*, 15:2547–2554.
- Oliver, J. L., Carpena, P., Hackenberg, M., and Bernaola-Galvan, P. (2004). IsoFinder: computational prediction of isochores in genome sequences. *Nucleic Acids Research*, 32:W287–W292. Web Server Issue.
- 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.
- 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.
- Stephens, D. A. (1994). Bayesian retrospective multiple-changepoint identification. *Applied Statistics*, 43:159–178.

Appendix A: Stratified Resampling Algorithm

We describe the stratified resampling algorithm of Carpenter et al. (1999) in terms of the SOR and SRC algorithms. Assume we currently have a set of N ordered particles $c^{(1)} < c^{(2)} < \dots < c^{(N)}$, with associated weights $w^{(1)}, \dots, w^{(N)}$, which sum to unity. For the SOR algorithm define α as in step (SOR1); and for SRC we assume that the value of α is given. Resampling of M particles proceeds as follows:

- (A) Simulate u a realisation of a uniform random variable on $[0, \alpha]$. Set $i = 1$.
- (B1) If $w^{(i)} \geq \alpha$ then propagate particle $c^{(i)}$ with weight $w^{(i)}$; else let $u = u - w^{(i)}$; if $u \leq 0$ then resample particle $c^{(i)}$ and assign a weight α , and set $u = u + \alpha$.
- (C) Let $i = i + 1$; if $i \leq N$ then return to (B).

Appendix B: Proof of Theorem 1

Theorem 1 considers the error of a resampling algorithm as measured by:

$$\text{mKSD} = \max \left\{ \max_i \left| \sum_{j=1}^i w^{(j)} - W^{(i)} \right| \right\}$$

For SOR, if $w^{(i)} \geq \alpha$ then $W^{(i)} = w^{(i)}$ with probability 1. As such we can consider the mKSD solely for the subset of particles which have $w^{(i)} < \alpha$. Assume we have N' such particles, and relabel these particles $c^{(1)} < c^{(2)} < \dots < c^{(N')}$.

The only randomness in the SOR algorithm is the simulation of u in step (A) of the algorithm detailed in Appendix A. Now for a given value of u

$$\sum_{j=1}^i W^{(j)} = \alpha \left[\left(\sum_{j=1}^i w^{(j)} + \alpha - u \right) / \alpha \right], \quad (9)$$

where $[x]$ is the integer part of x . Thus for all u and i

$$\left| \sum_{j=1}^i w^{(j)} - W^{(i)} \right| \leq \alpha,$$

so $\text{mKSD} \leq \alpha$.

For result (i) it suffices to note that if the probability of resampling particle $c^{(i)}$ is strictly less than 1; then $\text{mKSD} \geq w^{(i)}$.

For result (ii) it is sufficient to note that both the optimal resampling algorithm of Fearnhead and Clifford (2003) (where particles are shuffled prior to stratified resampling) and rejection control (where each particle with weight less than α is resampled independently of all others) give positive probability to all realisations of weights $W^{(1)}, W^{(2)}, \dots, W^{(N)}$ that our SOR algorithm does. It trivially follows that the mKSD for these algorithms will be greater than that of our SOR algorithm.

Appendix C: Error bound for SRC

Consider N particles, ordered so that $c^{(1)} < c^{(2)} < \dots < c^{(N)}$. We denote the weight of these particles prior to resampling by $w^{(i)}$, the unnormalised weights after resampling by $W^{(i)}$, and the normalised weights after resampling by $\bar{W}^{(i)}$. We let u denote the realisation of the Uniform $[0, \alpha]$ random variable used in the stratified resampling algorithm. Finally we let

$$\epsilon^{(i)} = \sum_{j=1}^i (w^{(j)} - W^{(j)}).$$

The sum of the resampling weights depends on the number of particles resampled in stage SRC2. There exists a constant, β , satisfying $0 \leq \beta < \alpha$ such that

$$\sum_{i=1}^N W^{(i)} = \begin{cases} 1 + \alpha - \beta & u \leq \beta, \\ 1 - \beta & u > \beta \end{cases}$$

Fix u and β . From (9) it can be shown that $u - \alpha \leq \epsilon^{(i)} \leq u$ for all i . We consider in turn the situation $u \leq \beta$ and $u > \beta$, corresponding to the two possible values of the sums of the unnormalised weights after resampling.

Firstly, assume $u \leq \beta$. Then we have

$$\begin{aligned} \left| \sum_{j=1}^i (w^{(j)} - \bar{W}^{(j)}) \right| &= \left| \sum_{j=1}^i (w^{(j)} - W^{(j)} / (1 - \beta + \alpha)) \right| \\ &= \frac{1}{1 + \alpha - \beta} \left| \epsilon^{(i)} + (\alpha - \beta) \sum_{j=1}^i w^{(j)} \right| \\ &\leq \frac{1}{1 + \alpha - \beta} \max \left\{ u + (\alpha - \beta) \sum_{j=1}^i w^{(j)}, \alpha - u - (\alpha - \beta) \sum_{j=1}^i w^{(j)} \right\}, \end{aligned}$$

where the two terms we are maximising over correspond to the largest positive and negative values of $\epsilon^{(i)}$. Now, as $u \leq \beta$ and $0 < \sum_{j=1}^i w^{(j)} < 1$, both these terms are bounded above by α . Thus we have $\text{mKSD} < \alpha$ in this case.

Now if $u > \beta$, by a similar argument we obtain

$$\left| \sum_{j=1}^i (w^{(j)} - \bar{W}^{(j)}) \right| \leq \frac{1}{1-\beta} \max \left\{ u - \beta \sum_{j=1}^i w^{(j)}, \alpha - u + \beta \sum_{j=1}^i w^{(j)} \right\} \leq \frac{\alpha}{(1-\beta)}.$$

The last inequality uses the fact that $u \leq \beta$ and $0 < \sum_{j=1}^i w^{(j)} < 1$. Finally as $\beta < \alpha$ we can obtain that $\text{mKSD} < \alpha/(1-\alpha)$.