

### 1. Contribution

We introduce a flexible system of modelling streaming data in a Bayesian regression setting. We combine two regression methods, Bayesian Cart by Chipman et al. 1998 and the Kalman Filter as derived by Meinhold and Singpurwalla 1983 because both minimise the mean square error based on the conditional expectation. Like Chipman et al. 2010 we form an ensemble of trees and perform inference over the weighted sum of the trees. Similar work has been done by Taddy et al. 2011 and Gramacy and Lee 2008.

### 2. Trees and Filters

A tree divides up a large covariate space,  $\mathcal{X}$ , using splitting threshold rules which assign observations to each of the partitions. This provides both a prior structure on the covariate space and concentrates the likelihood of the observations to each partition based on the conditional expectation:  $E[y_t | X_{t,i} \dots X_{t,i+j}] = z_t$

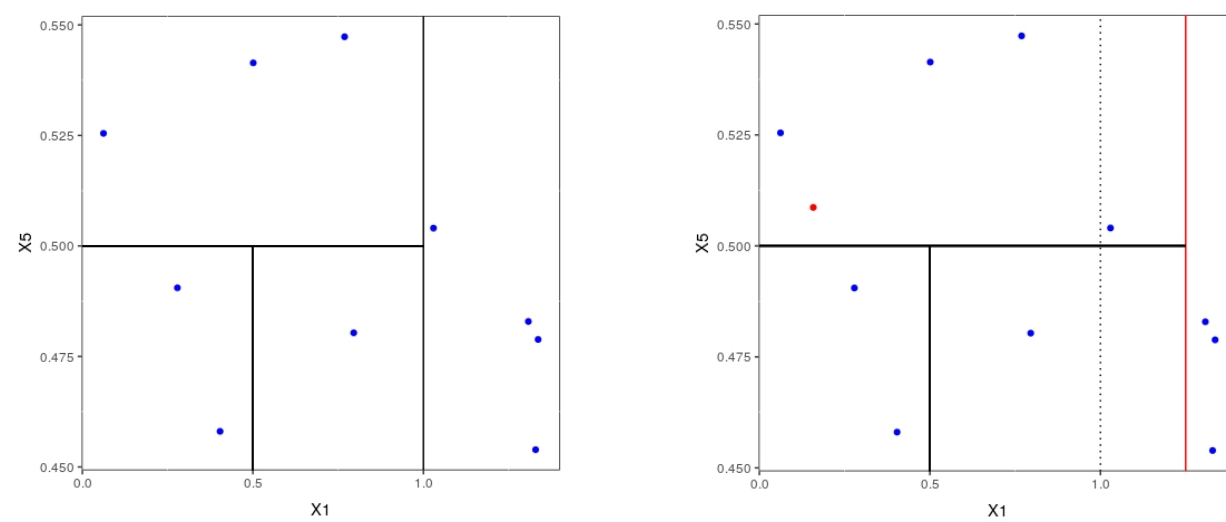


Figure 1: An evolving tree space.

The Kalman filter prediction of the next observation is based on the conditional expectation of the previous state:  $E[y_t | z_t] = HFz_{t-1}$ , where  $y_t = Hz_t + v_t$  and  $z_t = Fz_{t-1} + w_t$

Using an adaptation of the Kalman filter to a sensor network, Sinopoli et al. 2004 developed an intermittent Kalman filter which we use as a means of updating the  $1 \dots K_T$  filters of each tree.

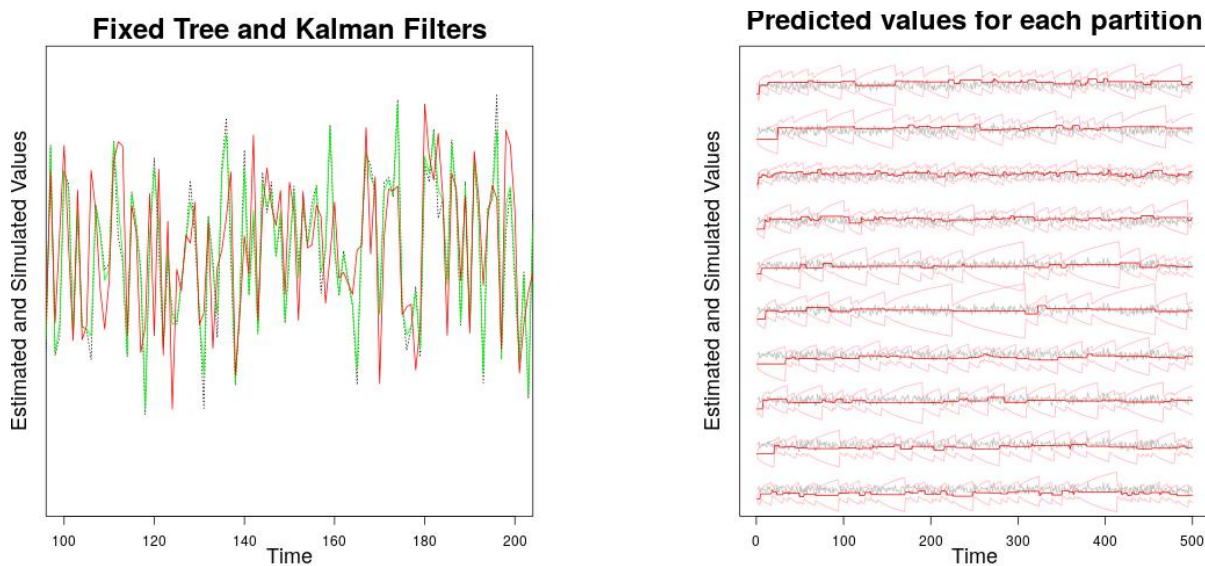


Figure 2: Estimating a process by creating smaller subprocesses.

In Bayesian dynamic regression trees (BDRT) a stream of data  $(x_t, y_t)$  is considered, with  $y_t$  following the regression tree model as described in Section 3. The latent mean process is dynamic as is the growth of the tree structure. This permits the relationship between  $y_t$  and  $x_t$  to change with time as the relationship between  $y_t$  and  $z_t$  is changing. Each leaf then models an independent Gaussian process.

### References

- Chipman, Hugh A, Edward I George, and Robert E McCulloch (1998). "Bayesian CART model search". In: *Journal of the American Statistical Association* 93.443, pp. 935–948.
- Chipman, Hugh A, Edward I George, Robert E McCulloch, et al. (2010). "BART: Bayesian additive regression trees". In: *The Annals of Applied Statistics* 4.1, pp. 266–298.
- Geyer, Charles J and Elizabeth A Thompson (1995). "Annealing Markov chain Monte Carlo with applications to ancestral inference". In: *Journal of the American Statistical Association* 90.431, pp. 909–920.
- Gramacy, Robert B and Herbert K. H Lee (2008). "Bayesian Treed Gaussian Process Models With an Application to Computer Modeling". In: *Journal of the American Statistical Association* 103.483, pp. 1119–1130.
- Mehra, R. (Apr. 1970). "On the identification of variances and adaptive Kalman filtering". In: *IEEE Transactions on Automatic Control* 15.2, pp. 175–184. ISSN: 0018-9286.
- Meinhold, Richard J and Nozer D Singpurwalla (1983). "Understanding the Kalman filter". In: *The American Statistician* 37.2, pp. 123–127.
- Sinopoli, Bruno, Luca Schenato, Massimo Franceschetti, Kameshwar Poolla, Michael I Jordan, and Shankar S Sastry (2004). "Kalman filtering with intermittent observations". In: *IEEE transactions on Automatic Control* 49.9, pp. 1453–1464.
- Taddy, Matthew A, Robert B Gramacy, and Nicholas G Polson (2011). "Dynamic trees for learning and design". In: *Journal of the American Statistical Association* 106.493, pp. 109–123.

### 3. Base Model

Let  $y_t \in \mathbb{R}^n$  and  $x_t \in \mathbb{R}^m$ . The tree  $T$  partitions  $\mathcal{X} \subseteq \mathbb{R}^p$  into  $K_T$  subsets (e.g. it has  $K_T$  leaves), so  $z_t = (z_{t1}, \dots, z_{tK_T})$ , with  $z_{tk} \in \mathbb{R}^m$ . Then we describe the model as follows:

$$\begin{aligned} E(y | T, z, x) &= z_{\eta(x,T)} \\ y_t | T, z_t, x_t &\sim N(z_{t, \eta(x_t, T)}, V_{\eta(x_t, T)}) \\ z_{t+1, k} &\sim N(F_k z_{tk}, W_k), \quad k = 1, \dots, K_T \end{aligned} \quad (1)$$

Thus we have multiple independent autoregressive Gaussian process with initial values  $z_{0k} \sim N(\mu_{0k}, W_{0k})$ .

The  $p(z_{tk} | T, \theta_T, x^t, y^t)$  is derived from the Kalman filter and the Gaussian model at the partitions allows us to derive an exact form of  $p(T | \theta_T, x^t, y^t) =$

$$\begin{aligned} p(T) \prod_{k=1}^{K_T} \int p(z_{0k}) \prod_{i=1}^t p(y_i | z_{i,k}, T)^{I_{i,k}} \\ \cdot p(z_{i,k} | z_{i-1,k}, u_t, T) dz_{i,k}^t \\ = p(T) \prod_{k=1}^{K_T} (|2\pi W_0|)^{-\frac{1}{2}} \\ \cdot \left( \prod_{i=1}^t (|2\pi W_k| |2\pi A_i|)^{-\frac{1}{2}} (|2\pi V_k|)^{-\frac{1}{2} I_{i,k}} \right) \\ \exp \left[ -\frac{1}{2} \left( \mu_{0,k} W_0^{-1} \mu_{0,k} - d_0^T A_i^{-1} d_0 + \right. \right. \\ \left. \left. \sum_{i=1}^t I_{i,k} y_{i,k}^T V_k^{-1} y_{i,k} + u_i^T G^T W_k^{-1} G u_i - d_i^T A_i^{-1} d_i \right) \right]. \end{aligned}$$

The main point to note is that this computation is  $O(n)$  because the current term only relies on the previous term. Thus the the main factors influencing complexity are the Kalman filter,  $O(n^3)$  and tree size. Using Openmp we get  $O(\lceil r \max(K_T)/P \rceil n^3)$  for each iteration of  $r$  trees.

### 4. Simulation Study

BDRT are:w compared to the Kalman Filter. A time-series data set was simulated using the Mackey-Glass nonlinear time series with a  $\tau$  of 20. 50 trees over 100 iterations were used with known parameters:

$$H = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad F = \begin{bmatrix} 0.9 & 0 \\ 0 & 0.2 \end{bmatrix}, \quad W = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}, \quad y \in \mathbb{R}, \quad z \in \mathbb{R}^2 \quad \text{and} \quad V = 0.03.$$

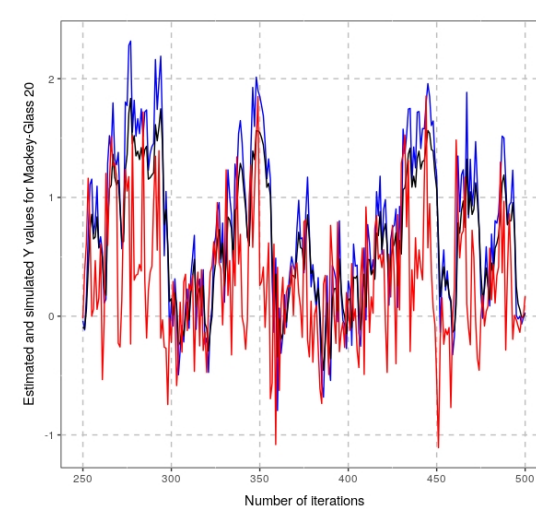


Figure 3: Showing observation predictions with BDRT and the Kalman Filter

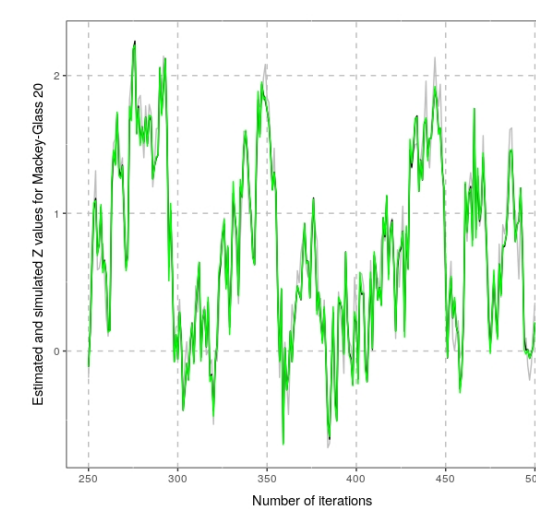


Figure 4: Showing latent state predictions with BDRT and the Kalman Filter

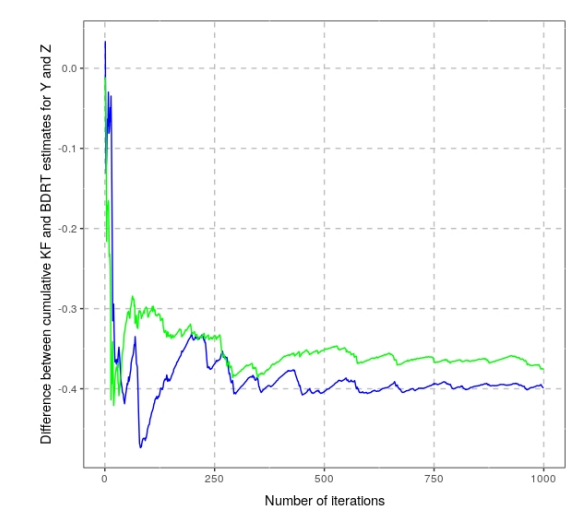


Figure 5: Difference in RMSE between BDRT and the Kalman Filter

Comparing different MCMC methods on BDRT. "MH" is the Chipman et al. method using grow, prune, swap and change. "BST" uses the same moves but has 10 levels of heating. The pseudoprior is estimated stochastically using the method suggested by Geyer and Thompson 1995. "MST" uses different moves: multigrow, multiprune, multichange, shift, and swap. The upper bound of growing, changing and pruning is temperature dependent.

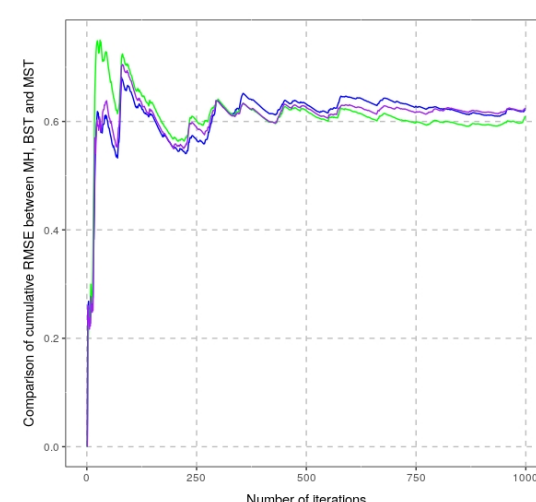


Figure 6: Comparing RMSE between the 3 different MCMC approaches.

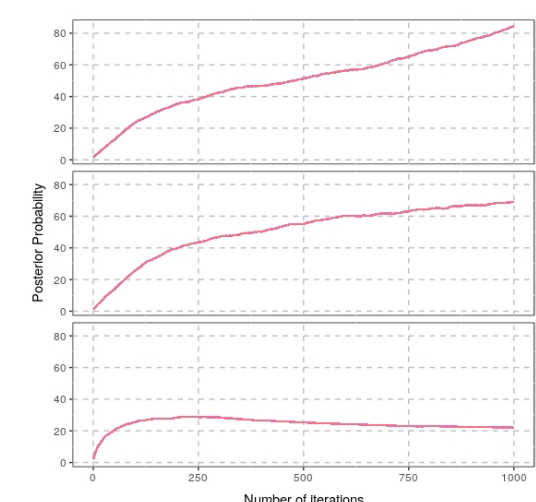


Figure 8: Average tree size as the algorithm progresses

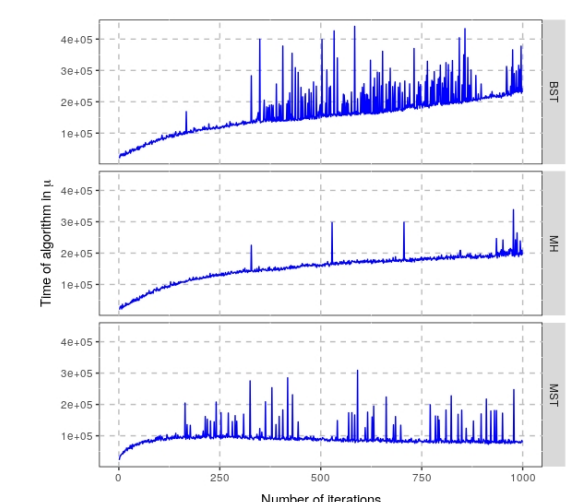


Figure 10: Time comparisons between different MCMC methods.

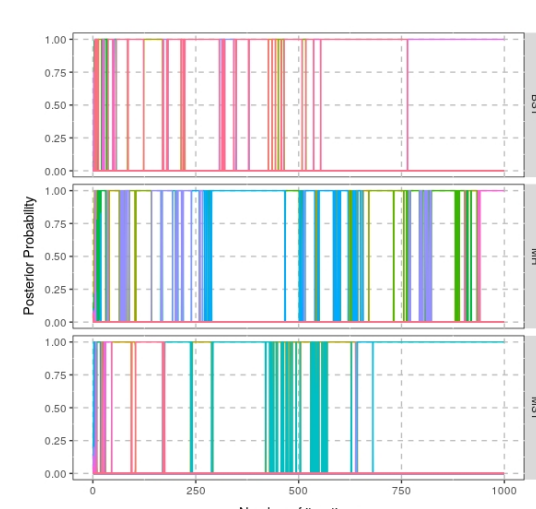


Figure 7: The probability of the trees alternate between zero and one.

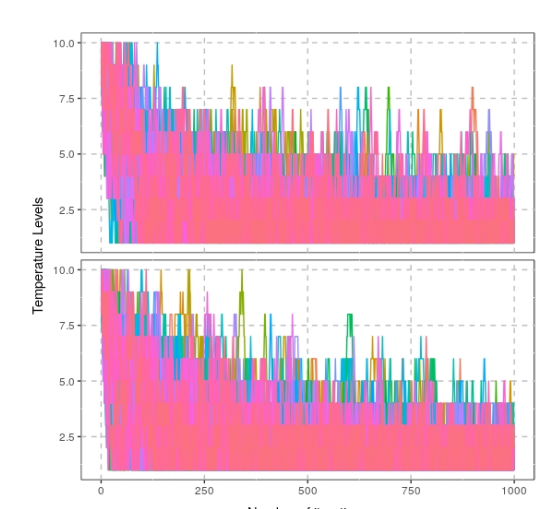


Figure 9: Temperature traversal of the trees started at the highest temperature.

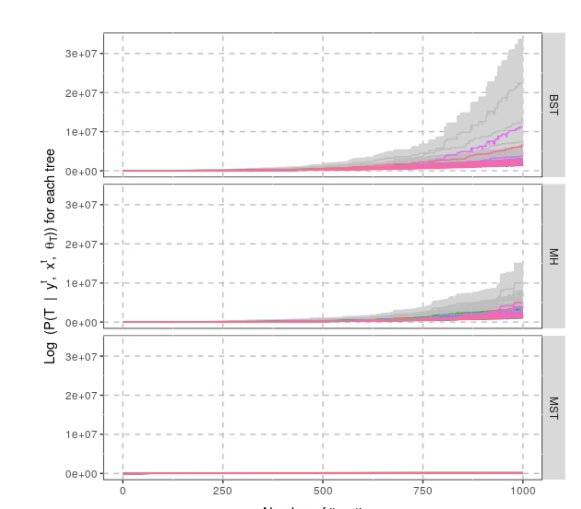


Figure 11: The log posterior for updated and not-updated leaves and their sum.

### 5. Streaming

Exchangeability of responses based on conditional data allows us to to develop a window-like streaming algorithm. If the data is arriving faster than it can be processed,  $\lambda_a < \lambda_s$ , then we can randomly select inputs within the window size and discard/store those that preceded the selected input.

The size of the window is based on the decisions and resources available to the analyst including, tree size ( $K_T$ ), rate of data arrival, ( $\lambda_s$ ), rate of algorithm ( $\lambda_a$ ), choice of model type (state only estimation, dual estimation, parameter learning, variable or model selection).

### 6. Extensions and More

One of the reasons for using the Kalman Filter is that many extensions exist. In particular, the Unscented Kalman filter allows for nonlinear dynamics to be modelled and is an improvement over the Extended Kalman Filter. The calculation of the posterior uses the idea that the marginalisation over  $z_t$  in the posterior calculation is the same as the expectation w.r.t.  $z_t$  and  $p(T | \theta_T, x^t, y^t)$ .

There are also methods developed by Mehra 1970 and others that allow us to adapt the algorithm for inference on the variance of the state  $W_t$ .

A further adaptation under development is to declare each leaf as a Gaussian marginal so that each tree represents a Gaussian process, i.e. the leaves of the tree represent the components of a vector  $\in \mathbb{R}^{K_T}$  where each leaf may be  $\in \mathbb{R}^m$ . In this case we have fixed tree sizes.