

Ancillarity-Sufficiency or not; Interweaving to Improve MCMC Estimation of the Local Level DLM

Matt Simpson

April 18, 2014

Abstract

In dynamic linear models (DLMs), MCMC sampling can often be very slow for estimating the posterior density — especially for longer time series. In particular, in some regions of the parameter space the standard data augmentation algorithm can mix very slowly. Recently ancillarity-sufficiency interweaving has been introduced as a method to take advantage of alternate parameterizations in multilevel models in order to improve the mixing and convergence properties of the chain. Focusing on the local level DLM, we explore alternate parameterizations and various interweaving algorithms through simulation in order to improve mixing. We conclude by explaining what our results may mean for MCMC in a more general DLM.

1 Introduction

Much is known about constructing efficient Gibbs samplers in multilevel models — e.g. see Van Dyk and Meng [2001] and Papaspiliopoulos et al. [2007], but little work has been done in improving Gibbs samplers in time series models despite some strong similarities in the structure of the models. In one example, Frühwirth-Schnatter [2004] applies the ideas of this literature to a dynamic univariate regression model with a stationary AR(1) process on the regression coefficient. In another, Pitt and Shephard [1999] computes analytic convergence rates for in a Gibbs sampler for a stationary AR(1) plus noise model. One recent idea to come out of the Gibbs sampling and data augmentation literature is “interweaving” between two data augmentations (Yu and Meng [2011]). This method was used by Kastner and Frühwirth-Schnatter [2013] to improve the MCMC estimation of stochastic volatility models. We apply this idea in the context of a dynamic linear model (DLM) in order to explore alternative methods of constructing a Gibbs sampler. In Section 2 we introduce the model. In Section 3 we discuss data augmentation methods generally, introduce several new data augmentations for this class of models and show that it is unlikely that a useful sufficient augmentation (centered augmentation) exists when we consider all model parameters as unknown. Section 4 contains several interweaving algorithms based on the various data augmentations we discuss in Section 3. In Section 5 we work out an example with the local level model and report the results of fitting the model to different simulated datasets using several of the MCMC samplers we construct. Finally Section 6 discusses our results and concludes.

2 Model

The general dynamic linear model (DLM) is a linear, gaussian, state space model and can be written as

$$y_t = F_t \theta_t + v_t \qquad v_t \stackrel{ind}{\sim} N_k(0, V_t) \qquad (1)$$

$$\theta_t = G_t \theta_{t-1} + w_t \qquad w_t \stackrel{ind}{\sim} N_p(0, W_t) \qquad (2)$$

for $t = 1, 2, \dots, T$, and $v_{1:T}$, $w_{1:T}$ independent. Equation (1) is called the *observation equation* and equation (2) is called the *system equation*. Similarly, $v_{1:T}$ are called the observation errors, $V_{1:T}$ are called the observation variances, $w_{1:T}$ are called the system disturbances and $W_{1:T}$ are called the system variances. The observed data is $y_{1:T}$ while $\theta_{0:T}$ are called the latent states. For each $t = 1, 2, \dots, T$, F_t is a $k \times p$ matrix and G_t is a $p \times p$ matrix. Let ϕ denote the vector of unknown parameters in the model. Then possibly $F_{1:T}$, $G_{1:T}$, $V_{1:T}$, and $W_{1:T}$ are all functions of ϕ .

We will focus our attention on a simple version of the DLM. Typically additional model structure is used to learn about $V_{1:T}$ and $W_{1:T}$ if time dependence is enforced – e.g. a stochastic volatility prior which would require a statespace model describing the $V_{1:T}$ ’s and $W_{1:T}$ ’s as data. Because of this additional complexity, we focus on the time-constant variances model, though many of our results may be useful in more complicated time-varying variance models. So we set $V_t = V$ and $W_t = W$ for $t = 1, 2, \dots, T$. We will also suppose that F_t and G_t are known matrices for $t = 1, 2, \dots, T$, though this constraint is immaterial since relaxing it will simply add one or more Gibbs steps to the algorithms we explore so long as no parameter that enters any F_t or G_t also enters V or W . However, in one of the data augmentations that we discuss, the scaled error data augmentation, there is a bit more housekeeping associated with $F_{1:T}$ depending on an unknown parameter (Section 3.2).

When $\phi = (V, W)$ is our unknown parameter vector and we can write the model as

$$y_t | \theta_{0:T} \stackrel{\text{ind}}{\sim} N(F_t \theta_t, V) \quad (3)$$

$$\theta_t | \theta_{0:t-1} \sim N(G_t \theta_{t-1}, W) \quad (4)$$

To complete the model specification in a Bayesian context, we need priors on θ_0 , V , and W . We’ll use the standard approach and assume that they’re mutually independent a priori and that $\theta_0 \sim N(m_0, C_0)$, $V \sim IW(\Lambda_V, \lambda_V)$ and $W \sim IW(\Lambda_W, \lambda_W)$ where m_0 , C_0 , Λ_V , λ_V , Λ_W , and λ_W are known hyperparameters and $IW(\Lambda, \lambda)$ denotes the inverse Wishart distribution with degrees of freedom λ and positive definite scale matrix Λ . This allows us to write the full joint distribution of $(V, W, \theta_{0:T}, y_{1:T})$ as

$$\begin{aligned} p(V, W, \theta_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2} (\theta_0 - m_0)' C_0^{-1} (\theta_0 - m_0) \right] \\ &\times |V|^{-(\lambda_V + k + T + 2)/2} \exp \left[-\frac{1}{2} \text{tr}(\Lambda_V V^{-1}) \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T (y_t - F_t \theta_t)' V^{-1} (y_t - F_t \theta_t) \right] \\ &\times |W|^{-(\lambda_W + p + T + 2)/2} \exp \left[-\frac{1}{2} \text{tr}(\Lambda_W W^{-1}) \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T (\theta_t - G_t \theta_{t-1})' W^{-1} (\theta_t - G_t \theta_{t-1}) \right] \end{aligned} \quad (5)$$

where $p = \dim(\theta_t)$, $k = \dim(y_t)$, and $\text{tr}(\cdot)$ is the matrix trace operator.

3 Estimating the Model via Data Augmentation: Parameterization Issues

A well known method to estimate the DLM is via data augmentation (DA), as in Frühwirth-Schnatter [1994] and Carter and Kohn [1994]. The basic idea is to implement a Gibbs sampler with two blocks. The generic DA algorithm with parameter ϕ , augmented data θ , and data y obtains the $k + 1$ ’st state of the Markov chain, $\phi^{(k+1)}$, from the k ’th state, $\phi^{(k)}$ as follows (we implicitly condition on the data y in all algorithms and only superscript the previous and new draws of the model parameters — intermediate draws of a DA or the model parameter are not superscripted since they are not part of the Markov chain for ϕ):

Algorithm 1 (State Sampler for DLM).

$$[\theta | \phi^{(k)}] \rightarrow [\phi^{(k+1)} | \theta]$$

The first block runs a simulation smoother which draws the latent states from their conditional posterior distribution given the model parameters. This can be accomplished in a number of ways. Forward Filtering, Backward Sampling (FFBS) is the original method proposed by Frühwirth-Schnatter [1994] and Carter and Kohn [1994], but alternatives include Koopman [1993], De Jong and Shephard [1995] and McCausland et al. [2011]. The second block draws $\phi = (V, W)$ from their joint conditional posterior which in this model turns out to be independent inverse Wishart distributions. In particular

$$V|\theta_{0:T}, y_{1:T} \sim IW\left(\Lambda_V + \sum_{t=1}^T v_t v_t', \lambda_V + T\right)$$

$$W|\theta_{0:T}, y_{1:T} \sim IW\left(\Lambda_W + \sum_{t=1}^T w_t w_t', \lambda_W + T\right)$$

where $v_t = y_t - F_t \theta_t$ and $w_t = \theta_t - G_t \theta_{t-1}$. We are calling this algorithm the *state sampler*.

The main problem with the state sampler is that in some regions of the parameter space the Markov chain mixes poorly for some of the parameters. For example, in the univariate local level model ($F_t = G_t = 1$ for $t = 1, 2, \dots, T$) and similar models it is known that if the time constant variance of the latent states, W , is too small, mixing will be poor for W (Frühwirth-Schnatter [2004]).

One well known method of improving mixing and convergence in MCMC samplers is reparameterizing the model. Papaspiliopoulos et al. [2007] is a good summary. Most of the work in some way focuses on what are called centered and noncentered parameterizations. In our general notation where ϕ is the parameter, θ is the DA and y is the data, the parameterization (ϕ, θ) is a *centered parameterization* (CP) if $p(y|\theta, \phi) = p(y|\theta)$. The parameterization is a *noncentered parameterization* (NCP) if $p(\theta|\phi) = p(\theta)$. When (ϕ, θ) is a CP, θ is called a *centered augmentation* (CA) for ϕ and when (ϕ, θ) is a NCP, θ is called a *noncentered augmentation* (NCA) for ϕ . A centered augmentation is sometimes called a *sufficient augmentation* (SA) and a noncentered augmentation is sometimes called an *ancillary augmentation* (AA), e.g. in Yu and Meng [2011]. Like Yu and Meng, we prefer the latter terminology because it immediately suggests the intuition that a sufficient augmentation is like a sufficient statistic while an ancillary augmentation is like an ancillary statistic and hence Basu’s theorem suggests that they are conditionally independent given ϕ .

The key reasoning behind the emphasis on SAs and AAs is that typically when the DA algorithm based on the SA has nice mixing and convergence properties the DA algorithm based on the AA has poor mixing and convergence properties and vice versa. In other words, the two algorithms form a “beauty and the beast” pair. This property suggests that there might be some way to combine the two DA algorithms or the two underlying parameterizations in order to construct a sampler which has “good enough” properties all the time. Papaspiliopoulos et al. [2007] for example suggest alternating between the two augmentations within a Gibbs sampler. Some work focuses on using partially noncentered parameterizations that are a sort of bridge between the CP and NCP, e.g. Papaspiliopoulos et al. for general hierarchical models and Frühwirth-Schnatter [2004] in the context of a particular DLM — a dynamic univariate regression with a stationary AR(1) coefficient.

Another method of combining the two DAs is through what Yu and Meng [2011] call interweaving. The idea is pretty simple: suppose that ϕ denotes the parameter vector, θ denotes one augmented data vector, γ denotes another augmented data vector, and y denotes the data. Then an MCMC algorithm that *interweaves* between θ and γ performs the following steps in a single iteration to obtain the $k+1$ ’st draw, $\phi^{(k+1)}$, from the k ’th draw, $\phi^{(k)}$:

Algorithm 2 (GIS for DLM).

$$[\theta|\phi^{(k)}] \rightarrow [\gamma|\theta] \rightarrow [\phi^{(k+1)}|\gamma]$$

Notice that an additional step is added to algorithm 1, and the final step now draws ϕ conditional on γ instead of θ . This is the intuition behind the name “interweaving”—the draw of the second augmented data vector is weaved in between the draws of θ and ϕ . This particular method of interweaving is called a *global* interweaving strategy (GIS) since interweaving occurs globally across the entire parameter vector. It’s

possible to define a *componentwise* interweaving strategy (CIS) that interweaves within specific steps of a Gibbs sampler as well. Step two of the GIS algorithm is typically accomplished by sampling $\phi|\theta, y$ and then $\gamma|\theta, \phi, y$. In addition, γ and θ are often, but not always, one-to-one transformations of each other conditional on (ϕ, y) , i.e. $\gamma = M(\theta; \phi, y)$. Where $M(\cdot; \phi, y)$ is a one-to-one function. In this case, the algorithm becomes:

Algorithm 3 (GIS for DLM, expanded).

$$[\theta|\phi^{(k)}] \rightarrow [\phi|\theta] \rightarrow [\gamma|\theta, \phi] \rightarrow [\phi^{(k+1)}|\gamma]$$

When γ is a one-to-one transformation of θ , step 4 is an update $\gamma = M(\theta; \phi, y)$. The GIS algorithm is directly comparable to the *alternating* algorithm suggested by Papaspiliopoulos et al. [2007]. Given the same two DAs, θ and γ , and parameter vector ϕ , the alternating algorithm for sampling from $p(\phi|y)$ is as follows:

Algorithm 4 (Alternating for DLM).

$$[\theta|\phi^{(k)}] \rightarrow [\phi|\theta] \rightarrow [\gamma|\phi] \rightarrow [\phi^{(k+1)}|\gamma]$$

The key difference between this algorithm and algorithm 3 is in step 3: instead of drawing from $p(\gamma|\theta, \phi, y)$, the alternating algorithm draws from $p(\gamma|\phi, y)$. In other words it alternates between two data augmentation algorithms in a single iteration. The interweaving algorithm, on the other hand, connects or “weaves” the two separate iterations together in step 3 by drawing γ conditional on θ in addition to ϕ and y .

Yu and Meng call a GIS approach where one of the DAs is a SA and the other is an AA an *ancillary sufficient interweaving strategy*, or an ASIS. They show that the GIS algorithm has a geometric rate of convergence no worse than the worst of the two underlying algorithms and in some cases better than the corresponding alternating algorithm. In particular, their Theorem 1 suggests that the weaker the dependence between two data augmentations in the posterior, the more efficient the GIS algorithm. In the limit of a posteriori independent data augmentations, the GIS algorithm will even obtain iid draws from the posterior of the model parameter. This helps motivate their focus on ASIS — conditional on the model parameter, a SA and an AA are independent (under the conditions of Basu’s theorem), which suggests that the dependence between the two DAs will be limited in the posterior. In fact, when the prior on ϕ is nice in some sense, Yu and Meng show that the ASIS algorithm is the same as the optimal PX-DA algorithm of Meng and Van Dyk [1999], Liu and Wu [1999], Van Dyk and Meng [2001] and Hobert and Marchev [2008]. Their results suggest that ASIS and interweaving generally is a promising approach to improve the speed of MCMC in a variety of models no matter what region of the parameter space the posterior is concentrated and the exploitation of ASIS in particular looks more promising than the well known alternating algorithms.

To gain some intuition about why interweaving works, recall that a typical problem with slow MCMC is that there is high autocorrelation in the Markov chain for ϕ , $\{\phi^{(k)}\}_{k=1}^K$, leading to imprecise estimates of $E[f(\phi)|y]$ for some function f integrable with respect to the posterior of ϕ . Our goal is to reduce this dependence. In the usual DA algorithm, e.g. algorithm 1, when ϕ and θ are highly dependent in the joint posterior, the draws from $p(\theta|\phi, y)$ and then from $p(\phi|\theta, y)$ will hardly move the chain which results in high autocorrelation. Interweaving helps break this autocorrelation in two ways. First, by inserting the extra step, e.g. steps 2 and 3 together in 3, the chain gets an additional chance to move in a single iteration thereby weakening the autocorrelation. This is a feature of an alternating algorithm as well, but Yu and Meng show that the corresponding interweaving algorithm is often even more efficient. The key is the second point — when the posterior dependence between the two DAs is low, steps 2 and 3 in Algorithm 3, i.e. step 2 in Algorithm 2, is enough to almost completely break the dependence in the chain. For the alternating algorithm, it is typically not feasible to find a data augmentation such that step 2 or step 3 of Algorithm 4 completely breaks the dependence in the chain — this would require finding a DA such that the model parameter and the DA are essentially independent which, in turn, would likely mean that drawing from the conditional posterior of the parameter given the DA is nearly as difficult as drawing from the marginal posterior of the model parameter.

Aside from the intuition of finding a posteriori (nearly) independent DAs, both alternating and interweaving strategies suggest looking for a “beauty and the beast” pair of DAs — specifically both algorithms will tend to do better, all else equal, when the two underlying DA algorithms are efficient in opposite regions of the parameter space.

3.1 The scaled disturbances

The next step is to apply the ideas of interweaving to sampling from the posterior of the dynamic linear model. Papaspiliopoulos et al. [2007] note that typically the usual parameterization results in a SA for the parameter ϕ . All that's necessary for an ASIS algorithm, then, is to construct an AA for ϕ . We immediately run into a problem because the standard DA for a DLM is the latent states $\theta_{0:T}$. From equations (3) and (4) we see that V is in the observation equation so that $\theta_{0:T}$ is not a SA for (V, W) while W is in the system equation so that $\theta_{0:T}$ is not an AA for (V, W) either. In order to find a SA we need to somehow move V from the observation equation (3) to the system equation (4) while leaving W in the system equation. We also need to find an AA by somehow moving W from the system equation to the observation equation while leaving V in the observation equation. A naive thing to try is to condition on the disturbances instead of the states and see if the disturbances for a SA or an AA for (V, W) . The disturbances $w_{0:T}$ are defined by $w_t = \theta_t - G_t \theta_{t-1}$ for $t = 1, 2, \dots, T$ and $w_0 = \theta_0$. However the DA algorithm based on the w_t 's is identical to the algorithm based on the θ_t because it turns out that the conditional distributions $p(V, W | \theta_{0:T}, y_{1:T})$ and $p(V, W | w_{0:T}, y_{1:T})$ are identical.

Papaspiliopoulos et al. suggest that in order to obtain an ancillary augmentation for a variance parameter, we must scale the sufficient augmentation by the square root of that parameter. Based on this intuition, note that if we hold V constant then $\theta_{0:T}$ is a SA for W from the observation and system equations, (3) and (4), i.e. we say $\theta_{0:T}$ is a SA for W given V , or for $W|V$. Similarly $\theta_{0:T}$ is an AA for $V|W$. This suggests that if we scale θ_t by W appropriately for all t we'll have an ancillary augmentation for V and W jointly. The same intuition suggests scaling $w_t = \theta_t - G_t \theta_{t-1}$ by W appropriately for all t in order to find an ancillary augmentation for (V, W) . We will work with the latter case since it has already been used in the literature. In fact it follows Papaspiliopoulos et al.'s suggestion to construct a pivotal quantity in order to find an ancillary augmentation which incidentally also buttresses the case for the terminology "ancillary" and "sufficient" augmentations rather than "centered" and "non-centered". In some DLMs the DA algorithm based on scaling w_t and the DA algorithm based on scaling θ_t will be the same, but this is not generally true — and even fails to hold for some of the simplest DLMs.

To define the scaled disturbances in the general DLM, let L_W denote the Cholesky decomposition of W , i.e. the lower triangle matrix L_W such that $L_W L_W' = W$. Then we will define the scaled disturbances $\gamma_{0:T}$ by $\gamma_0 = \theta_0$ and $\gamma_t = L_W^{-1}(\theta_t - G_t \theta_{t-1})$ for $t = 1, 2, \dots, T$. There are actually $p!$ different versions of the scaled disturbances depending on how we order the elements of θ_t , as Meng and Van Dyk [1998] note for EM algorithms in a different class of models. We will sidestep the issue of the best ordering of the latent states. No matter which ordering is chosen, we can confirm our intuition that the scaled disturbances are an AA for V and W jointly. The reverse transformation is defined recursively by $\theta_0 = \gamma_0$ and $\theta_t = L_W \gamma_t + G_t \theta_{t-1}$ for $t = 1, 2, \dots, T$. Then the Jacobian is block lower triangular with the identity matrix and T copies of L_W along the diagonal blocks, so $|J| = |L_W|^T = |W|^{T/2}$. Then from (5) we can write the full joint distribution of $(V, W, \gamma_{0:T}, y_{1:T})$ as

$$\begin{aligned} p(V, W, \gamma_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2} (\gamma_0 - m_0)' C_0^{-1} (\gamma_0 - m_0) \right] \\ &\times |W|^{-(\lambda_W + p + T + 2)/2} \exp \left[-\frac{1}{2} \text{tr} (\Lambda_W W^{-1}) \right] \exp \left[-\frac{1}{2} \gamma_t' \gamma_t \right] |V|^{-(\eta_t + k + T + 2)/2} \\ &\times \exp \left[-\frac{1}{2} \left(\text{tr} (\Lambda_V V^{-1}) + \sum_{t=1}^T [y_t - F_t \theta_t(\gamma_{0:T}, W)]' V^{-1} [y_t - F_t \theta_t(\gamma_{0:T}, W)] \right) \right] \end{aligned} \quad (6)$$

where $\theta_t(\gamma_{0:T}, W)$ denotes the recursive back transformation defined by the scaled disturbances. So ultimately under the scaled disturbance parameterization we can write the model as

$$\begin{aligned} y_t | \gamma_{0:T}, V, W &\stackrel{ind}{\sim} N(F_t \theta_t(\gamma_{0:T}, W), V) \\ \gamma_t &\stackrel{iid}{\sim} N(0, I_p) \end{aligned} \quad (7)$$

for $t = 1, 2, \dots, T$ where I_p is the $p \times p$ identity matrix. Neither V nor W are in the system equation so the scaled disturbances are an AA for (V, W) . This parameterization is well known, e.g. Frühwirth-Schnatter [2004] use it in a dynamic regression model with stationary regression coefficient.

The DA algorithm based on $\gamma_{0:T}$ draws $\gamma_{0:T}$ from its conditional posterior and then (V, W) from their joint conditional posterior given $\gamma_{0:T}$. There are a couple methods of performing this draw, including applying one of the simulation smoothers directly to drawing $\gamma_{0:T}$, if possible, or using one of them to draw the latent states $\theta_{0:T}$ before transforming the states to the scaled disturbances. The draw from the joint conditional posterior of (V, W) is tricky because it is not a known density. We will illustrate how to accomplish it in a worked example in Section 5.

3.2 The scaled errors

The scaled disturbances immediately suggest another potential AA that seems like it should be analogous — the scaled observation errors or more succinctly the scaled errors. What we are referring to is $v_t = y_t - F_t \theta_t$ appropriately scaled by V in the general DLM. Now let L_V denote the Cholesky decomposition of V , that is $L_V L_V' = V$. Then we can define a version of the scaled errors (this time depending on how we order the elements of y_t) as $\psi_0 = \theta_0$ and $\psi_t = L_V^{-1}(y_t - F_t \theta_t)$ for $t = 1, 2, \dots, T$. This is a bit strange since in general $\dim(\psi_0) \neq \dim(\psi_t)$ for $t = 1, 2, \dots, T$. Ideally we might like an “ F_0 ” so that we can set $\psi_0 = F_0 \theta_0$ in order for ψ_0 to have the same dimension as ψ_1 . However, in general there is no F_0 . In some DLMs F_t is constant with respect to t so that we could set $F_0 = F$, but in dynamic regression for example, there is no natural F_0 assuming that we do not have the time-zero values of the covariates. To avoid this issue in practice, we simply leave $\psi_0 = \theta_0$ though transforming the initial value could in principle result in an algorithm with better properties.

There is a real difficulty, however. With this definition of $\psi_{0:T}$ it is not straightforward to write down the model in terms of $\psi_{0:T}$ instead of $\theta_{0:T}$ and determine $p(\psi_{0:T}|V, W)$. When F_t is $k \times k$ (so that $\dim(y_t) = k = p = \dim(\theta_t)$) and is invertible for $t = 1, 2, \dots, T$, $\psi_{0:T}$ is a one-to-one transformation of $\theta_{0:T}$ and the problem is easier. Then $\theta_t = F_t^{-1}(y_t - L_V \psi_t)$ for $t = 1, 2, \dots, T$ while $\theta_0 = \psi_0$. The Jacobian of this transformation is block diagonal with a single copy of the identity matrix and the $F_t^{-1} L_V$'s along the diagonal, so $|J| = (\prod_{t=1}^T |F_t|^{-1}) |V|^{T/2}$. Then from (5) we can write the joint distribution of $(V, W, \psi_{0:T}, y_{1:T})$ as

$$\begin{aligned} p(V, W, \psi_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2} (\psi_0 - m_0)' C_0^{-1} (\psi_0 - m_0) \right] \\ &\times |V|^{-(\lambda_V + k + 2)/2} \exp \left[-\frac{1}{2} \text{tr} (\Lambda_V V^{-1}) \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T \psi_t' \psi_t \right] \\ &\times |W|^{-(\delta_t + k + T + 2)/2} \exp \left[-\frac{1}{2} \left(\text{tr} (\Lambda_W W^{-1}) + (y_t - \mu_t)' (F_t W F_t')^{-1} (y_t - \mu_t) \right) \right] \end{aligned} \quad (8)$$

where we define $\mu_1 = L_V \psi_1 + F_1 G_1 \psi_0$ and for $t = 2, 3, \dots, T$, $\mu_t = L_V \psi_t + F_t G_t F_{t-1}^{-1} (y_{t-1} - L_V \psi_{t-1})$. The $|F_t|^{-1}$'s have been absorbed into the normalizing constant, but if the F_t 's depended on some unknown parameter then we could not do this and as a result would have to take them into account in a Gibbs step for F_t . Now we can write the model in terms of the scaled error parameterization:

$$\begin{aligned} y_t | V, W, \psi_{0:T}, y_{1:t-1} &\sim N(\mu_t, F_t W F_t') \\ \psi_t &\stackrel{iid}{\sim} N(0, I_k) \end{aligned}$$

for $t = 1, 2, \dots, T$ where I_k is the $k \times k$ identity matrix. Now we see immediately that the scaled errors are also an AA for (V, W) since neither V nor W are in the system equation of this model. However, both V and W are in the observation equation so that $\psi_{0:T}$ is not a SA for (V, W) or for either one conditional on the other.

The DA algorithm based on $\psi_{0:T}$ is similar to that of $\gamma_{0:T}$ except we note that simulation smoothing can be accomplished by directly applying the algorithm of McCausland et al. [2011] because the precision matrix of $\psi_{0:T}$ retains the necessary tridiagonal structure. Also we mention in passing that there is a bit of symmetry here — the joint conditional posterior of (V, W) given $\gamma_{0:T}$ is from the same family of densities as that of (W, V) given $\psi_{0:T}$ so that V and W essentially switch places. The upshot is that if we can draw from one we can draw from the other, so this part of our work has been essentially halved.

3.3 The elusive search for a sufficient augmentation

Having found two separate ancillary augmentations for the DLM, we would like to find a sufficient augmentation in order to implement take advantage of their likely weak posterior dependence and implement an ASIS. It turns out that this is no easy task. From equations (3) and (4) we can rewrite the by recursively substituting as

$$y_t = v_t + F_t (w_t + G_t w_{t-1} + G_t G_{t-1} w_{t-2} + \dots + G_t G_{t-1} \dots G_2 w_1 + G_t G_{t-1} \dots G_1 \theta_0)$$

where $v_t \sim N(0, V)$ and $w_t \sim N(0, W)$ are independent. Here we see that θ_0 is given a special status relative to the other elements of the data augmentation which helps motivate not scaling it in the scaled disturbances or scaled errors. We are essentially treating it as a model parameter here and will continue to do so in this subsection because it makes finding a sufficient augmentation easier (though still essentially impossible).

Now each y_t is a linear combination of normal distributions conditional on $\phi = (\theta_0, V, W)$, so $y_{1:T}$ has a normal distribution after marginalizing out $\theta_{1:T}$ such that

$$E[y_t|\phi] = F_t \prod_{s=t}^1 G_s \theta_0$$

$$\text{Var}[y_t|\phi] = V + F_t H_t W$$

$$\text{Cov}[y_s, y_t|\phi] = F_t H_t W$$

where $\prod_{s=t}^1 G_s = G_t G_{t-1} \dots G_1$ and $H_t = I_p + G_t + G_t G_{t-1} + \dots + G_t G_{t-1} \dots G_2$. Next define

$$\mu = \begin{bmatrix} F_1 G_1 \theta_0 \\ F_2 G_2 G_1 \theta_0 \\ \vdots \\ F_T G_T G_{T-1} \dots G_1 \theta_0 \end{bmatrix}, \quad \tilde{V}_{k \times k} = \begin{bmatrix} V & 0 & 0 & \ddots & 0 \\ 0 & V & 0 & \ddots & 0 \\ 0 & 0 & V & \ddots & 0 \\ \ddots & \ddots & \ddots & \ddots & \ddots \\ 0 & 0 & 0 & \ddots & V \end{bmatrix}, \quad \tilde{W}_{k \times k} = \begin{bmatrix} F_1 H_1 \\ F_2 H_2 \\ \vdots \\ F_T H_T \end{bmatrix} W [H_1' F_1' \quad H_2' F_2' \quad \dots \quad H_T' F_T'].$$

Then we have the data model for $y_{1:T}$ without any data augmentation:

$$y_{1:T} \sim N_{T \times k}(\mu, \tilde{V} + \tilde{W}).$$

Now given a prior $p(\phi)$, this defines the posterior distribution of interest $p(\phi|y_{1:T})$.

Next we wish to find a sufficient augmentation θ (the lack of a subscript distinguishes this from the latent states $\theta_{0:T}$). Suppose we have such an augmentation and that conditional on ϕ , $(y_{1:T}, \theta)$ are normally distributed, in other words

$$\begin{bmatrix} \theta \\ y \end{bmatrix} \Big| \phi \sim N \left(\begin{bmatrix} \alpha_\theta \\ \mu \end{bmatrix}, \begin{bmatrix} \Omega_\theta & \Omega'_{y,\theta} \\ \Omega_{y,\theta} & \tilde{V} + \tilde{W} \end{bmatrix} \right)$$

which implies

$$y|\theta, \phi \sim N(\mu + \Omega'_{y,\theta} \Omega_\theta^{-1} (\theta - \alpha_\theta), \tilde{V} + \tilde{W} - \Omega'_{y,\theta} \Omega_\theta^{-1} \Omega_{y,\theta})$$

$$\theta|\phi \sim N(\alpha_\theta, \Omega_\theta).$$

Now for θ to be a sufficient augmentation we need $\mu + \Omega'_{y,\theta}\Omega_\theta^{-1}(\theta - \alpha_\theta)$ and $\tilde{V} + \tilde{W} - \Omega'_{y,\theta}\Omega_\theta^{-1}\Omega_{y,\theta}$ to be independent of ϕ . This requires that

$$\mu + \Omega'_{y,\theta}\Omega_\theta^{-1}(\theta - \alpha_\theta) = A\theta$$

where A is a matrix which does not depend on ϕ . Rearranging, this gives $A = \Omega'_{y,\theta}\Omega_\theta^{-1}$ so that $\mu = A\alpha_\theta$ and $\Omega_{y,\theta} = \Omega_\theta A'$. Then using the second equation, we now require $\Sigma = \tilde{V} + \tilde{W} - A\Omega_\theta A'$ free of ϕ . This gives $A\Omega_\theta A' = \tilde{V} + \tilde{W} - \Sigma$. Consider $\tilde{\theta} = A\theta$, which is also a sufficient augmentation. Then we have

$$\tilde{\theta}|\phi \sim N(\mu, \tilde{V} + \tilde{W} - \Sigma)$$

and thus the posterior of ϕ given $\tilde{\theta}$ can be written as

$$p(\phi|\tilde{\theta}, y) \propto p(\phi)|\tilde{V} + \tilde{W} - \Sigma|^{-1/2} \exp \left[-\frac{1}{2}(\tilde{\theta} - \mu)'(\tilde{V} + \tilde{W} - \Sigma)^{-1}(\tilde{\theta} - \mu) \right].$$

If Σ were the zero matrix, this is the posterior we wish to sample from. The transformation from $\tilde{\theta}$ to θ is unlikely to make this any easier.

The fundamental problem is that once we find a sufficient augmentation, in order to use it we must obtain draws from a density that appears just as hard to sample from as the posterior density we are already trying to approximate. We did treat θ_0 as a model parameter instead of an element of the data augmentation above, but changing this only makes the resulting conditional posterior of ϕ more complicated. The logic above does not rule out the existence of a useful sufficient augmentation, but it does suggest that it will be difficult to find one. This result brings to mind Van Dyk and Meng [2001]’s contention that there is an art to constructing data augmentation algorithms — our goal is not only to find an MCMC algorithm that has nice convergence and mixing properties, but also one that is easy to implement, which is a criteria that is much harder to quantify.

The problem we run into is unlikely to be unique to the time series setting but rather seems driven by trying to find a sufficient augmentation for a pair of variances, one on the data level and the other on the latent data level. For example, in a hierarchical model we expect there to be similar problems finding a SA when both the observational and hierarchical variance are unknown. There is a similar problem while trying to find two data augmentations that are independent in the posterior which, by Yu and Meng [2011]’s theorem 1, would guarantee an interweaving algorithm that yields iid draws of from the posterior distribution of the model parameters. We omit the details, but unsurprisingly after making sensible sounding assumptions about the nature of the DAs (i.e. joint with the data they are normally distributed along with some (in)dependence assumptions), the conditional posterior of ϕ ends of being identical to or just as complicated as the marginal posterior of ϕ .

3.4 The “wrongly scaled” DAs

The scaled disturbances are defined by $\gamma_t = L_W^{-1}(\theta_t - G_t\theta_{t-1})$ and the scaled errors are defined by $\psi_t = L_V^{-1}(y_t - \theta_t)$ for $t = 1, 2, \dots, T$ where $L_W L_W' = W$ and $L_V L_V' = V$. Now define $\tilde{\gamma}_t = L_V^{-1}(\theta_t - G_t\theta_{t-1})$ and $\tilde{\psi}_t = L_W^{-1}(y_t - \theta_t)$ for $t = 1, 2, \dots, T$ and $\tilde{\psi}_0 = \tilde{\gamma}_0 = \theta_0$. In other words, the “tilde” versions of the scaled disturbances and the scaled errors are scaled by the “wrong” Cholesky decomposition, hence we call them the wrongly scaled disturbances and the wrongly scaled errors respectively. It is hard to motivate these DAs without looking forward to componentwise interweaving in the DLM (section 4), but you can at least view them as the result of having thrown spaghetti against the walls to see what sticks. Once again both of these DAs have many variations depending on how the elements of θ_t or y_t are ordered, but we will ignore that issue.

First consider $\tilde{\gamma}_{0:T}$. Notice that for $t = 1, 2, \dots, T$, $\tilde{\gamma}_t = L_V^{-1}L_W\gamma_t$ while $\tilde{\gamma}_0 = \gamma_0$. The reverse transformation is then $\gamma_t = L_W^{-1}L_V\tilde{\gamma}_t$. The Jacobian is then block diagonal with $L_W^{-1}L_V$ along the diagonal. Thus

$|J| = |L_W|^{-T} |L_V|^T = |W|^{-T/2} |V|^{T/2}$. Then from (6) we can write the joint distribution of $(V, W, \tilde{\gamma}_{0:T}, y_{1:T})$ as

$$\begin{aligned} p(V, W, \tilde{\gamma}_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2} (\tilde{\gamma}_0 - m_0)' C_0^{-1} (\tilde{\gamma}_0 - m_0) \right] |V|^{-(\lambda_V + k + 2)/2} \exp \left[-\frac{1}{2} \text{tr} (\Lambda_V V^{-1}) \right] \\ &\times |W|^{-T/2} \exp \left[-\frac{1}{2} \sum_{t=1}^T (y_t - F_t \theta_t(\tilde{\gamma}_{0:T}))' V^{-1} (y_t - F_t \theta_t(\tilde{\gamma}_{0:T})) \right] \\ &\times |W|^{-(\lambda_W + k + 2)/2} \exp \left[-\frac{1}{2} \text{tr} (\Lambda_W W^{-1}) \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T \tilde{\gamma}_t' (L_V^{-1} W (L_V^{-1})')^{-1} \tilde{\gamma}_t \right] \end{aligned} \quad (9)$$

Then under $\tilde{\gamma}_{0:T}$ we can write the model as

$$\begin{aligned} y_t | \tilde{\gamma}_{0:T}, V, W &\stackrel{\text{ind}}{\sim} N(F_t \theta_t(\tilde{\gamma}_{0:T}), V) \\ \tilde{\gamma}_t &\stackrel{\text{ind}}{\sim} N(0, L_V^{-1} W (L_V^{-1})') \end{aligned}$$

for $t = 1, 2, \dots, T$. Since L_V is the Cholesky decomposition of V , the observation equation does not contain W . So $\tilde{\gamma}_{0:T}$ is a SA for $W|V$. Note also that since W and L_V are both in the system equation, $\tilde{\gamma}_{0:T}$ is not an AA for V nor for W .

Now consider $\tilde{\psi}_t = L_W^{-1} L_V \psi_t$ for $t = 1, 2, \dots, T$ where again $\tilde{\psi}_0 = \psi_0 = \theta_0$. Then $\psi_t = L_V^{-1} L_W \tilde{\psi}_t$ and the Jacobian is block diagonal with $L_V^{-1} L_W$ along the diagonal. So $|J| = |V|^{-T/2} |W|^{T/2}$ and from (8) we can write the joint distribution of $(V, W, \tilde{\psi}_{0:T}, y_{1:T})$ as

$$\begin{aligned} p(V, W, \tilde{\psi}_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2} (\tilde{\psi}_0 - m_0)' C_0^{-1} (\tilde{\psi}_0 - m_0) \right] \\ &\times |V|^{-(\lambda_V + k + 2)/2} \exp \left[-\frac{1}{2} \text{tr} (\Lambda_V V^{-1}) \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T \tilde{\psi}_t' (L_W^{-1} V (L_W^{-1})')^{-1} \tilde{\psi}_t \right] \\ &\times |W|^{-(\lambda_W + k + 2)/2} \exp \left[-\frac{1}{2} \text{tr} (\Lambda_W W^{-1}) \right] |V|^{-T/2} \exp \left[-\frac{1}{2} \sum_{t=1}^T (y_t - \tilde{\mu}_t)' (F_t W F_t')^{-1} (y_t - \tilde{\mu}_t) \right] \end{aligned} \quad (10)$$

where we define $\tilde{\mu}_1 = L_W \tilde{\psi}_1 - F_1 G_1 \tilde{\psi}_0$ and for $t = 2, 3, \dots, T$ $\tilde{\mu}_t = L_W \tilde{\psi}_t - F_t G_t F_{t-1}^{-1} (y_{t-1} - L_W \tilde{\psi}_{t-1})$. In terms of $\tilde{\psi}_{0:T}$, the model is then:

$$\begin{aligned} y_t | V, W, \tilde{\psi}_{0:T}, y_{1:t-1} &\sim N(\tilde{\mu}_t, F_t W F_t') \\ \tilde{\psi}_t &\stackrel{iid}{\sim} N(0, L_W^{-1} V (L_W^{-1})') \end{aligned}$$

for $t = 1, 2, \dots, T$. Since $\tilde{\mu}_t$ only depends on W (through L_W) and not on V , V is absent from the observation equation. Thus $\tilde{\psi}_{0:T}$ is a SA for $V|W$. Again that both W and V are in the system equation so $\tilde{\psi}_{0:T}$ is not an AA for either V or W .

In the case of both wrongly scaled DA algorithms, the smoothing step can be accomplished in a manner analogous to the ‘‘correctly scaled’’ DA algorithms, i.e. the scaled disturbance and scaled error algorithms. The draw from the joint conditional posterior of (V, W) is from a nonstandard density that, like for the correctly scaled DA algorithms, has a certain symmetry property. Specifically $V, W | \tilde{\gamma}_{0:T}, y_{1:T}$ and $W, V | \tilde{\psi}_{0:T}, y_{1:T}$ have densities from the same family so that by changing which of $\tilde{\gamma}_{0:T}$ or $\tilde{\psi}_{0:T}$ is conditioned on, V and W essentially switch places. This class of densities is different from the correctly scaled DA case, however. We will demonstrate this through an example in Section 5.

4 Interweaving in the DLM: Global and Componentwise

We now have five DAs for the generic DLM with known F_t 's and G_t 's. For simplicity we'll assume that $\dim(y_t) = \dim(\theta_t)$ and F_t invertible for $t = 1, 2, \dots, T$ so that the scaled errors are easy to work with.

The five DAs are the states, $\theta_{0:T}$, the scaled disturbances $\gamma_{0:T}$, the scaled errors $\psi_{0:T}$, the wrongly scaled disturbances $\tilde{\gamma}_{0:T}$, and the wrongly scaled errors $\tilde{\psi}_{0:T}$. This allows us to construct several GIS algorithms based on Algorithm 3. The main algorithms we consider are the State-Dist, State-Error, Dist-Error, and Triple GIS algorithms. The State-Dist algorithm, for example, interweaves between the states and the scaled disturbances, while the Triple GIS algorithm interweaves between the states, the scaled diturbances, and the scaled errors. Strictly speaking the order in which we sample the DAs in the algorithm does matter, but Yu and Meng note that this tends not to make much difference. We always construct our algorithms so that the DAs are used in the order they were presented earlier in this paragraph.

To illustrate the GIS algorithms, 5 is the state-dist GIS algorithm:

Algorithm 5 (State-Dist for DLM).

$$[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\gamma_{0:T}|V, W, \theta_{0:T}] \rightarrow [V^{(k+1)}, W^{(k+1)}|\gamma_{0:T}]$$

the third step is actually a one-to-one transformation from $\theta_{0:T}$ to $\gamma_{0:T}$. In practice we may want to break up step 4 into two steps if it is easier to draw from the full conditionals of V and W rather than drawing them jointly, though this will cost us both in terms of MCMC efficiency and theoretical tractability for analyzing the algorithm.

None of the GIS algorithms we can construct are ASIS algorithms — none of the DAs are a SA for (V, W) . The states, $\theta_{0:T}$, are a SA for $W|V$ though, so this motivates a CIS algorithm. A partial CIS algorithm is immediate:

Algorithm 6 (Partial CIS for DLM).

$$\begin{aligned} [\theta_{0:T}|V^{(k)}, W^{(k)}] &\rightarrow [V^{(k+1)}|W^{(k)}, \theta_{0:T}] \rightarrow \\ [W|V^{(k+1)}, \theta_{0:T}] &\rightarrow [\gamma_{0:T}|V^{(k+1)}, W, \theta_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \gamma_{0:T}] \end{aligned}$$

Steps 1-2 of this algorithm correspond to a Gibbs step for V while steps 3-5 correspond to a Gibbs step for W . Step 4 is a transformation step since conditional on V and W , $\gamma_{0:T}$ is a one-to-one transformation of $\theta_{0:T}$. This algorithm is actually the same as a version of the State-Dist interweaving algorithm with some of the steps rearranged, specifically algorithm 5. So it should be similar in performance to a GIS algorithm.

With a little more work, we can also construct a Full CIS algorithm that also turns out to be essentially the same as another GIS algorithm. Here we employ the wrongly scaled disturbances $\tilde{\gamma}_{0:T}$ and wrongly scaled errors $\tilde{\psi}_{0:T}$. Now we already now that $\gamma_{0:T}$ is an AA for $W|V$ and $\tilde{\gamma}_{0:T}$ is a SA for $W|V$, so the two form an AA-SA pair for $W|V$. Similarly, $\psi_{0:T}$ is an AA for $V|W$ while $\tilde{\psi}_{0:T}$ is a SA for $V|W$ so together they form an AA-SA pair for $V|W$. Now we can construct a Full CIS algorithm:

Algorithm 7.

$$\begin{aligned} [\tilde{\psi}_{0:T}|V^{(k)}, W^{(k)}] &\rightarrow [V|W^{(k)}, \tilde{\psi}_{0:T}] \rightarrow [\psi_{0:T}|V, W^{(k)}, \tilde{\psi}_{0:T}] \rightarrow [V^{(k+1)}|W^{(k)}, \psi_{0:T}] \rightarrow \\ [\tilde{\gamma}_{0:T}|V^{(k+1)}, W^{(k)}, \psi_{0:T}] &\rightarrow [W|V^{(k+1)}, \tilde{\gamma}_{0:T}] \rightarrow [\gamma_{0:T}|V^{(k+1)}, W, \tilde{\gamma}_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \gamma_{0:T}] \end{aligned}$$

Steps 1-4 constitute a Gibbs step for V and steps 5-8 constitute a Gibbs step for W . Steps 3, 5 and 7 are transformation steps — the parameter we are drawing is a one to one function of the parameters we are conditioning on. It turns out that $p(W|V, \tilde{\gamma}_{0:T}, y_{1:T})$ and $p(W|V, \theta_{0:T}, y_{1:T})$ are the same density, and also that $p(V|W, \tilde{\psi}_{0:T}, y_{1:T})$ and $p(V|W, \theta_{0:T}, y_{1:T})$ are the same density. The upshot is that step 1 of algorithm 7 can be replaced with a draw from $p(\theta_{0:T}|V, W, y_{1:T})$, and any time we condition on one of the “wrongly scaled” variables, we can condition on $\theta_{0:T}$ instead, yielding the following version of the same CIS algorithm:

Algorithm 8.

$$\begin{aligned} [\theta_{0:T}|V^{(k)}, W^{(k)}] &\rightarrow [V|W^{(k)}, \theta_{0:T}] \rightarrow [\psi_{0:T}|V, W^{(k)}, \theta_{0:T}] \rightarrow [V^{(k+1)}|W^{(k)}, \psi_{0:T}] \rightarrow \\ [\theta_{0:T}|V^{(k+1)}, W^{(k)}, \psi_{0:T}] &\rightarrow [W|V^{(k+1)}, \theta_{0:T}] \rightarrow [\gamma_{0:T}|V^{(k+1)}, W, \theta_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \gamma_{0:T}] \end{aligned}$$

We omit the details, but it can also be shown that this algorithm and the Dist-Error algorithm employ the same steps, just in a different order. This suggests that we should expect the two algorithms to perform similarly.

5 Application: The Local Level Model

In order to illustrate how these algorithms work, we will focus on the local level model for simplicity though there are still some difficulties. The local level model (LLM) is a DLM with univariate data y_t for $t = 1, 2, \dots, T$ and a univariate latent state θ_t for $t = 0, 2, \dots, T$ that satisfies

$$y_t | \theta_{0:T} \stackrel{ind}{\sim} N(\theta_t, V) \quad (11)$$

$$\theta_t | \theta_{0:t-1} \sim N(\theta_{t-1}, W) \quad (12)$$

with $\theta_0 \sim N(m_0, C_0)$. Here $\theta_t = E[y_t | \theta_{0:T}]$. The states are $\theta_{0:T}$, the scaled disturbances are $\gamma_{0:T}$ with $\gamma_0 = \theta_0$ and $\gamma_t = (\theta_t - \theta_{t-1})/\sqrt{W}$ for $t = 1, 2, \dots, T$, and the scaled errors are $\psi_{0:T}$ with $\psi_0 = \theta_0$ and $\psi_t = (y_t - \theta_t)/\sqrt{V}$ for $t = 1, 2, \dots, T$. The independent inverse Wishart priors on V and W in Section 2 cash out to independent inverse gamma priors for the local level model, viz. $V \sim IG(\alpha_V, \beta_V)$ and $W \sim IG(\alpha_W, \beta_W)$.

5.1 Base Samplers

The joint density of $(V, W, \theta_{0:T}, y_{1:T})$ is:

$$\begin{aligned} p(V, W, \theta_{0:T}, y_{1:T}) &\propto V^{-(\alpha_V + 1 + T/2)} \exp \left[-\frac{1}{V} \left(\beta_V + \frac{1}{2} \sum_{t=1}^T (y_t - \theta_t)^2 \right) \right] \\ &W^{-(\alpha_W + 1 + T/2)} \exp \left[-\frac{1}{W} \left(\beta_W + \frac{1}{2} \sum_{t=1}^T (\theta_t - \theta_{t-1})^2 \right) \right] \exp \left[-\frac{1}{2C_0} (\theta_0 - m_0)^2 \right] \end{aligned} \quad (13)$$

This immediately gives the state sampler:

Algorithm 9 (State Sampler for LLM).

$$[\theta_{0:T} | V^{(k)}, W^{(k)}] \rightarrow [V^{(k+1)}, W^{(k+1)} | \theta_{0:T}]$$

In step 2, V and W are independent with $V \sim IG(a_V, b_V)$ and $W \sim IG(a_W, b_W)$ where $a_V = \alpha_V + T/2$, $b_V = \beta_V + \sum_{t=1}^T (y_t - \theta_t)^2/2$, $a_W = \alpha_W + T/2$, and $b_W = \beta_W + \sum_{t=1}^T (\theta_t - \theta_{t-1})^2/2$.

The scaled disturbance sampler, i.e. the DA algorithm based on the scaled disturbances, is a bit more complicated. In this context $\gamma_0 = \theta_0$ and $\gamma_t = (\theta_t - \theta_{t-1})/\sqrt{W}$ for $t = 1, 2, \dots, T$, and thus $\theta_t = \gamma_0 + \sqrt{W} \sum_{s=1}^t \gamma_s$ for $t = 1, 2, \dots, T$. Following (6), we can write the joint posterior of $(V, W, \gamma_{0:T})$ as

$$\begin{aligned} p(V, W, \gamma_{0:T} | y_{1:T}) &\propto V^{-(\alpha_V + 1 + T/2)} \exp \left[-\frac{1}{V} \left(\beta_V + \frac{1}{2} \sum_{t=1}^T (y_t - \gamma_0 - \sqrt{W} \sum_{s=1}^t \gamma_s)^2 \right) \right] \\ &\times W^{-(\alpha_W + 1)} \exp \left[-\frac{\beta_W}{W} \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T \gamma_t^2 \right] \exp \left[-\frac{1}{2C_0} (\gamma_0 - m_0)^2 \right] \end{aligned} \quad (14)$$

Now V and W are no longer conditionally independent given $\gamma_{0:T}$ and $y_{1:T}$. Instead of attempting the usual DA algorithm, we will add an extra Gibbs step and draw V and W separately primarily for ease of computation. This gives us the scaled disturbance sampler:

Algorithm 10 (Scaled Disturbance Sampler for LLM).

$$[\gamma_{0:T} | V^{(k)}, W^{(k)}] \rightarrow [V^{(k+1)} | W^{(k)}, \gamma_{0:T}] \rightarrow [W^{(k+1)} | V^{(k+1)}, \gamma_{0:T}]$$

In step 2, V is drawn from the same inverse gamma distribution as in step 2 of algorithm 9. In step 3, the draw of W is more complicated. The density can be written as

$$p(W | V, \gamma_{0:T}, y_{1:T}) \propto W^{-\alpha_W - 1} \exp \left[-\frac{1}{2V} \sum_{t=1}^T \left(y_t - \gamma_0 - \sqrt{W} \sum_{s=1}^t \gamma_s \right)^2 \right] \exp \left[-\frac{\beta_W}{W} \right].$$

This density is not any known form and is difficult to sample from, though its functional form is similar to the generalized inverse gaussian distribution. The density can be written as

$$p(W|V, \gamma_{0:T}, y_{1:T}) \propto W^{-\alpha_W - 1} \exp \left[-aW + b\sqrt{W} - \frac{\beta_W}{W} \right].$$

where $a = \sum_{t=1}^T (\sum_{j=1}^t \gamma_j)^2 / 2V$ and $b = \sum_{t=1}^T (y_t - \gamma_0) (\sum_{j=1}^t \gamma_j) / V$. It can be shown that $b > \left(\frac{(\alpha+1)^3}{\beta} \right)^{1/2} \frac{4\sqrt{2}}{3\sqrt{3}}$ implies that the density is log concave. It turns out that this tends to hold over a wide region of the parameter space — so long as V is smaller or is not much larger than W . This allows for the use of adaptive rejection sampling in order to sample from this distribution in many cases, e.g. using Gilks and Wild [1992]. An alternative is to use a t approximation to the conditional density as a proposal in a rejection sampler. This is much more computationally expensive when necessary, but it works ok on the log scale.

The scaled error sampler is similar to the scaled disturbance sampler and this is easy to see in the local level model. Here $\psi_0 = \theta_0$ and $\psi_t = (y_t - \theta_t) / \sqrt{V}$ for $t = 1, 2, \dots, T$ so that $\theta_t = y_t - \sqrt{V}\psi_t$ for $t = 1, 2, \dots, T$. From (8) we can write $p(V, W, \psi_{0:T} | y_{1:T})$ as

$$p(V, W, \psi_{0:T}, y_{1:T}) \propto W^{-(\alpha_W + 1 + T/2)} \exp \left[-\frac{1}{W} \left(\beta_W + \frac{1}{2} \sum_{t=1}^T (Ly_t - \sqrt{V}L\psi_t)^2 \right) \right] \\ V^{-(\alpha_V + 1)} \exp \left[-\frac{\beta_V}{V} \right] \exp \left[-\frac{1}{2} \sum_{t=1}^T \psi_t^2 \right] \exp \left[-\frac{1}{2C_0} (\psi_0 - m_0)^2 \right] \quad (15)$$

where we define $Ly_t = y_t - y_{t-1}$ for $t = 2, 3, \dots, T$ & $Ly_1 = y_1 - \psi_0$ and $L\psi_t = \psi_t - \psi_{t-1}$ for $t = 2, 3, \dots, T$ & $L\psi_1 = \psi_1 - 0$. Once again, V and W are no longer conditionally independent given $\psi_{0:T}$ and $y_{1:T}$. In fact, the density is analogous to (14) with V and W switching places. The scaled error sampler obtained from drawing V and W separately is:

Algorithm 11 (Scaled Error Sampler for LLM).

$$[\psi_{0:T} | V^{(k)}, W^{(k)}] \rightarrow [V^{(k+1)} | W^{(k)}, \psi_{0:T}] \rightarrow [W^{(k+1)} | V^{(k+1)}, \psi_{0:T}]$$

In step 3, W is drawn from the same inverse gamma distribution as in step 2 of algorithm 9. Drawing V in step 2 is exactly analogous to drawing W in algorithm 10. The density of $V | W, \psi_{0:T}, y_{1:T}$ can be written as

$$\log p(V | W, \psi_{0:T}, y_{1:T}) \propto V^{-\alpha_V - 1} \exp \left[-aV + b\sqrt{V} - \frac{\beta_V}{V} \right]$$

where $a = \sum_{t=1}^T (L\psi_t)^2 / 2W$ and $b = \sum_{t=1}^T (L\psi_t Ly_t) / W$, i.e. the form of the density is the same as that of $W | V, \gamma_{0:T}, y_{1:T}$.

We can also construct the DA algorithms based on the “wrongly scaled” disturbances or errors. The wrongly scaled disturbances are defined by $\tilde{\gamma}_t = \gamma_t \frac{\sqrt{W}}{\sqrt{V}}$ for $t = 1, 2, \dots, T$ and $\tilde{\gamma}_0 = \gamma_0$ while the wrongly scaled errors are defined by $\tilde{\psi}_t = \psi_t \frac{\sqrt{V}}{\sqrt{W}}$ for $t = 1, 2, \dots, T$ and $\tilde{\psi}_0 = \psi_0$. For $\tilde{\gamma}_{0:T}$ we have

$$p(V, W | \tilde{\gamma}_{0:T}, y_{1:T}) \propto W^{-\alpha_W - T/2 - 1} \exp \left[-\frac{1}{2W/V} \sum_{t=1}^T \tilde{\gamma}_t^2 \right] \exp \left[-\frac{\beta_W}{W} \right] \\ \times V^{-\alpha_V - 1} \exp \left[-\frac{\beta_V}{V} \right] \exp \left[-\frac{1}{2V} \sum_{t=1}^T \left(y_t - \tilde{\gamma}_0 - \sqrt{V} \sum_{s=1}^t \tilde{\gamma}_s \right)^2 \right]. \quad (16)$$

Thus the conditional posterior of W given V and $\tilde{\gamma}_{0:T}$ is the same as if we had conditioned on $\theta_{0:T}$ instead of $\tilde{\gamma}_{0:T}$. In other words

$$p(W | V, \tilde{\gamma}_{0:T}, y_{1:T}) \propto W^{-(\alpha_W + T/2) - 1} \exp \left[-\frac{1}{W} \left(\beta_W + \frac{1}{2} V \sum_{t=1}^T \tilde{\gamma}_t^2 \right) \right]$$

so that $V|W, \tilde{\gamma}_{0:T}, y_{1:T} \sim IG(a_W, b_W)$ where $a_W = \alpha_W + T/2$ and

$$b_W = \beta_W + \frac{1}{2}V \sum_{t=1}^T \tilde{\gamma}_t^2 = \beta_W + \frac{1}{2} \sum_{t=1}^T (\theta_t - \theta_{t-1})^2.$$

The conditional posterior of V is more complicated. We have

$$\begin{aligned} p(V|W, \tilde{\gamma}_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2W/V} \sum_{t=1}^T \tilde{\gamma}_t^2 \right] V^{-\alpha_V-1} \exp \left[-\frac{\beta_V}{V} \right] \exp \left[-\frac{1}{2V} \sum_{t=1}^T \left(y_t - \tilde{\gamma}_0 - \sqrt{V} \sum_{s=1}^t \tilde{\gamma}_s \right)^2 \right] \\ &\propto V^{-\alpha_V-1} \exp \left[-\frac{a}{V} + \frac{b}{\sqrt{V}} - cV \right] \end{aligned}$$

where

$$\begin{aligned} a &= \beta_V + \frac{1}{2} \sum_{t=1}^T (y_t - \tilde{\gamma}_0)^2 > 0 \\ b &= \sum_{t=1}^T (y_t - \tilde{\gamma}_0) \sum_{s=1}^t \tilde{\gamma}_s \\ c &= \frac{1}{2W} \sum_{t=1}^T \tilde{\gamma}_t^2 > 0. \end{aligned}$$

We will return to this density momentarily.

For the wrongly scaled errors, we have

$$\begin{aligned} p(V, W|\tilde{\psi}_{0:T}, y_{1:T}) &\propto V^{-\alpha_V-T/2-1} \exp \left[-\frac{1}{2V/W} \sum_{t=1}^T \tilde{\psi}_t^2 \right] \exp \left[-\frac{\beta_V}{V} \right] \\ &\times W^{-\alpha_W-1} \exp \left[-\frac{1}{2W} \sum_{t=1}^T \left(\tilde{L}y_t - \sqrt{W}(\tilde{L}\tilde{\psi}_t) \right) \right] \end{aligned} \quad (17)$$

where we define $\tilde{L}\tilde{y}_t = y_t - y_{t-1}$ for $t = 1, 2, \dots, T$ and $\tilde{L}\tilde{y}_1 = y_1 - \tilde{\psi}_0$, and $\tilde{L}\tilde{\psi}_t = \tilde{\psi}_t - \tilde{\psi}_{t-1}$ for $t = 1, 2, \dots, T$ with $\tilde{L}\tilde{\psi}_1 = \tilde{\psi}_1$. Then the conditional posterior of V is the same as if we had conditioned on $\theta_{0:T}$ instead of $\tilde{\psi}_{0:T}$, i.e.

$$p(V|W, \tilde{\psi}_{0:T}, y_{1:T}) \propto V^{-(\alpha_V-T/2)-1} \exp \left[-\frac{1}{V} \left(\beta_V + \frac{1}{2}W \sum_{t=1}^T \tilde{\psi}_t^2 \right) \right]$$

so that $V|W, \tilde{\psi}_{0:T}, y_{1:T} \sim IG(a_V, b_V)$ where $a_V = \alpha_V + T/2$ and

$$b_V = \beta_V + \frac{1}{2}W \sum_{t=1}^T \tilde{\psi}_t^2 = \beta_V + \frac{1}{2} \sum_{t=1}^T (y_t - \theta_t)^2.$$

The conditional posterior of W is more complicated but similar to that of V when we conditioned on $\tilde{\gamma}_{0:T}$. We have

$$\begin{aligned} p(W|V, \tilde{\psi}_{0:T}, y_{1:T}) &\propto \exp \left[-\frac{1}{2V/W} \sum_{t=1}^T \tilde{\psi}_t^2 \right] W^{-\alpha_W-1} \exp \left[-\frac{1}{2W} \sum_{t=1}^T \left(\tilde{L}\tilde{y}_t - \sqrt{W}\tilde{L}\tilde{\psi}_t \right) \right] \\ &\propto W^{-\alpha_W-1} \exp \left[-\frac{a}{W} + \frac{b}{\sqrt{W}} - cW \right] \end{aligned}$$

where now

$$\begin{aligned} a &= \beta_W + \frac{1}{2} \sum_{t=1}^T L \tilde{y}_t^2 > 0 \\ b &= \sum_{t=1}^T L \tilde{y}_t L \tilde{\psi}_t \\ c &= \frac{1}{2V} \sum_{t=1}^T \tilde{\psi}_t^2 > 0. \end{aligned}$$

So in the case of both wrongly scaled DAs we need to sample from a density of the form

$$p(X) \propto X^{-\alpha-1} \exp \left[-\frac{a}{X} + \frac{b}{\sqrt{X}} - cX \right].$$

The density of $Y = \log(X)$ is

$$p(Y) \propto \exp \left[-\alpha Y - a e^{-Y} + b e^{-Y/2} - c e^Y \right].$$

This density is easy to sample from fairly efficiently with rejection sampler using a t approximation as a proposal. It is also typically log concave, so adaptive rejection sampling will work as well. In particular when $b \leq 0$ or $a > \frac{3b}{16} \left(\frac{b}{16c} \right)^{1/3}$ the density of Y is log concave.

5.2 Hybrid Samplers: Interweaving, Alternating and Random Kernel

Section 4 contains the details for the interweaving algorithms in the general DLM. In the local level model the only difference is that we only sample V and W jointly when we condition on the states. We will consider all four GIS samplers based on any two or three of the base samplers and one CIS sampler. In the GIS samplers, the order of the parameterizations will always be the states $(\theta_{0:T})$, then the scaled disturbances $(\gamma_{0:T})$, then the scaled errors $(\psi_{0:T})$. All of the GIS algorithms and the Full CIS algorithm are below in Table 1. Note the distributional forms for each of these steps (in some cases a transformation) are in Section 5.1. We omit the partial CIS algorithm, though note that in practice it is essentially the same as the State-Dist algorithm.

Interweaving algorithms are conceptually very similar to alternating algorithms. For every GIS algorithm, there is a corresponding alternating algorithm where each $[DA_2|V, W, DA_1]$ step is replaced by a $[DA_2|V, W]$ step (here DA_i is a data augmentation for $i = 1, 2$). Table 2 contains each alternating algorithm. Note that there are two possible “hybrid triple” algorithms that we do not consider here where the move from $\theta_{0:T}$ to $\gamma_{0:T}$ interweaves and while the move from $\gamma_{0:T}$ to $\psi_{0:T}$ alternates and vice versa.

Table 3 lists each algorithm we considered for the local level model. The basic idea here is that the alternating algorithms should serve as a sort of baseline to compare the corresponding interweaving algorithms against. The GIS version of an algorithm at least be similar to the alternating version since both take advantage of the fact that when at least one of the underlying DA algorithms works well, the hybrid algorithm should also work well. When the two DAs are (nearly) independent in the posterior, it is plausible that the GIS version of the algorithm would improve upon the results of the alternating version. Since we do not have such a pair of DAs, it seems unlikely that there will be much of a difference between the GIS and alternating versions of a given algorithm. We can make the comparisons fairly precise by considering the effective sample size (ESS) of the Markov chain – given the same sample size for two algorithms, we can compare ESS to see which one has better mixing.

5.3 Simulation Setup

In order to test these algorithms, we simulated a fake dataset from the local level model for various choices of V , W , and T . We created a grid over V – W space with (V, W) ranging from $(10^{-2}, 10^{-2})$ to $(10^2, 10^2)$ and we

1. State-Dist GIS algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\gamma_{0:T}|V, W, \theta_{0:T}] \rightarrow [V^{(k+1)}|W, \gamma_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \gamma_{0:T}]$
2. State-Error GIS algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\psi_{0:T}|V, W, \theta_{0:T}] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \psi_{0:T}]$
3. Dist-Error GIS algorithm:
 $[\gamma_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V|W^{(k)}, \gamma_{0:T}, y_{1:T}] \rightarrow [W|V, \gamma_{0:T}] \rightarrow [\psi_{0:T}|V, W, \gamma_{0:T}] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \psi_{0:T}]$
4. Triple GIS algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\gamma_{0:T}|V, W, \theta_{0:T}] \rightarrow [V|W, \gamma_{0:T}] \rightarrow [W|V, \gamma_{0:T}] \rightarrow [\psi_{0:T}|V, W, \gamma_{0:T}] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \psi_{0:T}]$
5. Full CIS algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V|W^{(k)}, \theta_{0:T}] \rightarrow [\psi_{0:T}|V, W, \theta_{0:T}] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [\theta_{0:T}|V^{(k+1)}, W] \rightarrow [W|V^{(k+1)}, \theta_{0:T}] \rightarrow [\gamma_{0:T}|V^{(k+1)}, W] \rightarrow [W^{(k+1)}|V^{(k+1)}, \gamma_{0:T}]$

Table 1: GIS and CIS algorithms for the local level model

1. State-Dist alternating algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\gamma_{0:T}|V, W] \rightarrow [V^{(k+1)}|W, \gamma_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \gamma_{0:T}]$
2. State-Error alternating GIS algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\psi_{0:T}|V, W] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \psi_{0:T}]$
3. Dist-Error alternating GIS algorithm:
 $[\gamma_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V|W, \gamma_{0:T}] \rightarrow [W|V, \gamma_{0:T}] \rightarrow [\psi_{0:T}|V, W] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \psi_{0:T}]$
4. Triple alternating GIS algorithm:
 $[\theta_{0:T}|V^{(k)}, W^{(k)}] \rightarrow [V, W|\theta_{0:T}] \rightarrow [\gamma_{0:T}|V, W] \rightarrow [V|W, \gamma_{0:T}] \rightarrow [W|V, \gamma_{0:T}] \rightarrow [\psi_{0:T}|V, W] \rightarrow [V^{(k+1)}|W, \psi_{0:T}] \rightarrow [W^{(k+1)}|V^{(k+1)}, \psi_{0:T}]$

Table 2: Alternating algorithms for the local level model

Base	State	(wrongly) Scaled Disturbance	(wrongly) Scaled Error	
GIS	State-Dist	State-Error	Dist-Error	Triple Triple
Alt	State-Dist	State-Error	Dist-Error	
CIS	State-Error/WError-Error for $V W$; State-Dist/WDist-Dist for $W V$			

Table 3: Each algorithm considered for the local level model

simulated a dataset for all possible combinations of V and W with each of $T = 10, 100, 1000$. Then for each dataset, we fit the local level model using each algorithm in Table 3. We used the same rule for constructing priors for each model: $\theta_0 \sim N(0, 10^7)$, $V \sim IG(5, 4V^*)$, and $W \sim IG(5, 4W^*)$, mutually independent where (V^*, W^*) are the true values of V and W used to simulate the time series. So both the prior and likelihood and thus the posterior roughly agree about the likely values of V and W .

For each dataset and each sampler we obtained $n = 3000$ draws and threw away the first 500 as burn in. The chains were started at the true values used to simulated the time series, so we can examine the behavior of the chains to determine how well they mix but not how quickly they converge. Define the effective sample proportion (ESP) for a scalar component of the chain as the effective number of independent draws, i.e.

effective sample size (ESS) (see e.g. Gelman et al. [2003]) of the component divided by the actual sample size, i.e. $ESP = ESS/n$. An $ESP = 1$ indicates that the Markov chain is behaving as if it obtains iid draws from the posterior. It is possible to obtain $ESP > 1$ if the draws are negatively correlated and occasionally for some of our samplers our estimates of ESP are greater than one, but we round this down to one for plotting purposes.

5.4 Base Results

Figure 1 contains plots of ESP for V and W in each chain of each base sampler for each of $T = 10$, $T = 100$, and $T = 1000$. We will focus on $T = 10$ first. The state sampler has a low ESP for V and a high ESP for W when the signal-to-noise ratio, W/V , is larger than one. When the signal-to-noise ratio is smaller than one, on the other hand, the state sampler has a low ESP for W and a high ESP for V . In the typical case where the signal to noise ratio close to one, the state sampler has a modest to low ESP for both V and W . Note that the particular values of V and W do not seem to matter at all — just their relative values, i.e. the signal-to-noise ratio W/V . Moving up any diagonal on the plots for V and W in the state sampler, W/V is constant and the ESS appears roughly constant. The basic lesson here is that the state sampler has mixing issues for whichever of V or W is smaller.

Figure 1 tells a different story for the scaled disturbance sampler. When the signal-to-noise ratio is less than one, ESP s for both V and W are nearly 1, i.e. the effective sample size is nearly the actual sample size of the chain. When the signal-to-noise ratio is greater than one, however, ESP for both V and W becomes small, especially for V . Once again the absolute values of V and W do not matter for this behavior — just the relative values. The scaled error sampler has essentially the opposite properties. When W/V is large, it has a near 1 ESP for both V and W . On the other hand, when W/V is small it has a low ESP for both V and W , especially for V . The lesson here seems to be that the scaled disturbances are the preferred data augmentation for low signal-to-noise ratios and the scaled errors are the preferred data augmentation for high signal-to-noise ratios, while the states are preferred for signal-to-noise ratios near 1. The wrongly scaled disturbances ($\tilde{\gamma}_{0:T}$) and wrongly scaled errors ($\tilde{\psi}_{0:T}$), on the other hand, look like worse versions of the state sampler. The pattern of mixing for V and W over the range of the parameter space is essentially the same as the state sampler, except the wrongly scaled disturbance sampler has worse mixing for V than the state sampler almost everywhere and similarly the wrongly scaled error sampler has worse mixing for W than the state sampler almost everywhere.

The plots for $T = 100$ and $T = 1000$ in Figure 1 tell basically the same story, with a twist. Increasing the length of the time series seems to exacerbate all problems without changing the basic conclusions. As T increases, W/V has to be smaller and smaller for the scaled disturbance sampler to have decent mixing, and similarly W/V has to be larger and larger for the scaled error sampler to have decent mixing. Interestingly, the scaled error sampler appears to mix well for both V and W over a larger region of the space $W/V < 1$ than the scaled disturbance sampler does over $W/V > 1$. The state sampler is stuck between a rock and a hard place, so to speak, since as T increases, good mixing for V requires W/V to be smaller and smaller, but good mixing for W requires W/V to be larger and larger. The wrongly scaled samplers are again pretty similar to the state sampler for larger T except the wrongly scaled sampler tends to be worse everywhere for the variance that was used to scale — i.e. once again the wrongly scaled disturbance sampler has worse mixing for V than the state sampler while the wrongly scaled error sampler has worse mixing for W than the state sampler. However, the wrongly scaled samplers do appear to have slightly better mixing than the state sampler for the variance that was *not* used to scale. In particular, the wrongly scaled error sampler appears to have slightly better mixing for V than the state sampler over part of the parameter space when $T = 100$ or $T = 1000$.

The behavior of the wrongly scaled data augmentation algorithms is consistent with what we showed in Section 4 — that the Full CIS algorithm based on the scaled errors and disturbances and the wrongly scaled errors and disturbances is equivalent to the Full CIS algorithm that replaces the wrongly scaled DAs with the usual latent states. Since the behavior of the state sampler and the wrongly scaled disturbance sampler are the same for W , we might expect that when drawing W it does not matter whether we use the states or the wrongly scaled disturbances in the Full CIS algorithm. Similarly since the behavior of the state sampler

Parameter	State	Dist	Error	W-Dist	W-Error
V	$\frac{W}{V} < 1$	$\frac{W}{V} < 1$	$\frac{W}{V} > 1$	$\frac{W}{V} < 1$	$\frac{W}{V} < 1$
W	$\frac{W}{V} > 1$	$\frac{W}{V} < 1$	$\frac{W}{V} > 1$	$\frac{W}{V} > 1$	$\frac{W}{V} > 1$

Table 4: Rule of thumb for when each base algorithm has a high effective sample size for each variable as a function of the signal-to-noise ratio, W/V .

and the wrongly scaled error sampler are the same for V , we might expect that it does not matter which one we use when drawing V . In fact this is what we found when constructing the Full CIS algorithm. Even though $(\gamma_{0:T}, \tilde{\gamma}_{0:T})$ forms an AA-SA pair for $W|V$ and $(\psi_{0:T}, \tilde{\psi}_{0:T})$ forms an AA-SA pair for $V|W$ while $\theta_{0:T}$ is not a SA for V , replacing $\tilde{\gamma}_{0:T}$ and $\tilde{\psi}_{0:T}$ with $\theta_{0:T}$ does not actually change the CIS algorithm.

We summarize some of the above results for convenience in Table 4. Most of the patterns of Figure 1 and Table 4 can be explained by Figure 2, which contains the estimated posterior correlations between various functions of parameters estimated using the simulations from the Triple-Alternating sampler. First we need to understand the correlations we are looking at. The state sampler consists of two steps — a draw of $\theta_{0:T}$ given V and W , and a draw of (V, W) given $\theta_{0:T}$. From Section 5.1 we have that conditional on $\theta_{0:T}$, V and W are independent in the posterior and each has an inverse gamma distribution that depends on the states only through the second parameter:

$$b_V = \beta_V + \sum_{t=1}^T (y_t - \theta_t)^2 / 2$$

$$b_W = \beta_W + \sum_{t=1}^T (\theta_t - \theta_{t-1})^2 / 2.$$

So can view (b_V, b_W) as the data augmentation instead of $\theta_{0:T}$ and thus the state sampler is

$$[b_V, b_W | V^{(k)}, W^{(k)}] \rightarrow [V^{(k+1)}, W^{(k+1)} | b_V, b_W].$$

Thus the dependence between (V, W) and (b_V, b_W) in the posterior will determine how much the state sampler moves in a given iteration and, in particular, the dependence in the marginal chain for V is determined by the dependence between b_V and V while the dependence in the marginal chain for W is determined by the dependence between b_W and W .

For the scaled disturbance sampler, things are a bit more complicated. Now the data augmentation becomes

$$b_V = \beta_V + \sum_{t=1}^T (y_t - \gamma_0 - \sqrt{W} \sum_{s=1}^t \gamma_s^2)^2 / 2 = \beta_V + \sum_{t=1}^T (y_t - \theta_t)^2 / 2$$

$$a_\gamma = \sum_{t=1}^T (\sum_{j=1}^t \gamma_j)^2 / 2V$$

$$b_\gamma = \sum_{t=1}^T (y_t - \gamma_0) (\sum_{j=1}^t \gamma_j) / V.$$

This again comes from Section 5.1. The draw of $V|W, \gamma_{0:T}$ is the same inverse gamma draw as in the state sampler. The draw of $W|V, \gamma_{0:T}$ depends on two random parameters, a_γ and b_γ defined above. So the dependence between V and b_V determines how much the marginal chain for V chain moves in a single iteration, while the dependence between W and (a_γ, b_γ) determines how much the marginal chain for W moves in a single iteration. The scaled error sampler is analogous to the scaled disturbance sampler except

with

$$\begin{aligned}
b_W &= \beta_W + \sum_{t=2}^T (Ly_t - \sqrt{V}L\psi_t)^2/2 = \beta_W + \sum_{t=1}^T (\theta_t - \theta_{t-1})^2/2 \\
a_\psi &= \sum_{t=1}^T (L\psi_t)^2/2W \\
b_\psi &= \sum_{t=1}^T (L\psi_t Ly_t)/W
\end{aligned}$$

and with the dependence in marginal chain for W depending on the dependence between b_W and W while the dependence in the marginal chain for V depending on the dependence between (a_ψ, b_ψ) and V .

In Figure 2 we see that the posterior correlation between V and b_V is high when $W/V > 1$ and is low when $W/V < 1$. This explains why both the state sampler and the scaled disturbance sampler have low ESP's for V when $W/V > 1$. Similarly the posterior correlation between W and b_W is high when $W/V < 1$ and is low when $W/V > 1$, which explains why both the state sampler and scaled error sampler have low ESP's for W when $W/V < 1$.

Insert explanation of marginal chain for W in the scaled disturbance sampler and for the marginal chain for V in the scaled error sampler. Right now it doesn't look right - the correlation between V and a_ψ or b_ψ is highest when W/V is high, but for the scaled error sampler, autocorrelation in the marginal chain for V is highest when W/V is low. However, here's a thought. Consider the marginal chain for W in the scaled disturbance sampler - this chain essentially consists of two steps, a draw of $(a_\gamma, b_\gamma|W)$ and a draw of $(W|a_\gamma, b_\gamma)$. When W is negatively correlated with both a_γ and b_γ , a large draw of W in the chain means that we expect small draws of a_γ and b_γ because of the negative correlation which, in turn, means that we expect a large draw of the next W again because of the negative correlation. This may explain what we are seeing.

5.5 GIS and CIS Results

Based on the intuition in Section 3 above, both the GIS and alternating algorithms should work best when at least one of the underlying base algorithms has a high ESP — the basic idea is that when least one of the underlying algorithms has low autocorrelation, we should have low autocorrelation in the GIS algorithm using multiple DAs. This suggests that the Dist-Error GIS and alternating algorithms will have the best performance of the GIS and alternating algorithms using two DAs for both V and W , especially for W/V far away from one. When W/V is near one it may offer no improvement, especially for large T . The State-Dist algorithms should have trouble with V when W/V is high since both the state sampler and the scaled disturbance sampler have trouble with V when W/V is high. Similarly, the State-Error GIS algorithm should have trouble with W when W/V is low since both underlying samplers have trouble with W when W/V is low. Since the triple algorithms add the state sampler into the Dist-Error algorithms, it seems plausible that it might improve mixing for one of V or W since for W/V different from one, the state sampler has good mixing for at least one of V or W .

There are a couple of ways to gain some intuition about what we expect the Full CIS algorithm to do before seeing the results. First, we mentioned in Section 4 that the Full CIS and the Dist-Error GIS algorithm consist of the same steps, just rearranged. This suggests that they should perform similarly so that we expect the Full CIS algorithm to have good mixing for both V and W when W/V is sufficiently different from one. We can draw the same conclusion in a different way by noticing that in the Gibbs step for V , the CIS algorithm interweaves between the states and the scaled errors and in the Gibbs step for W it interweaves between the states and the scaled disturbances. Since the state sampler has a high ESP for V when $W/V < 1$ and the scaled disturbance sampler has a high ESP for V when $W/V > 1$ we should expect the Full CIS sampler to have a high ESP for V when W/V is different from one. Similarly, since the state sampler has a high ESP for W when $W/V > 1$ and the scaled error sampler has a high ESP for W when $W/V < 1$, we should expect the Full CIS sampler to have a high ESP for W when W/V is different from one.

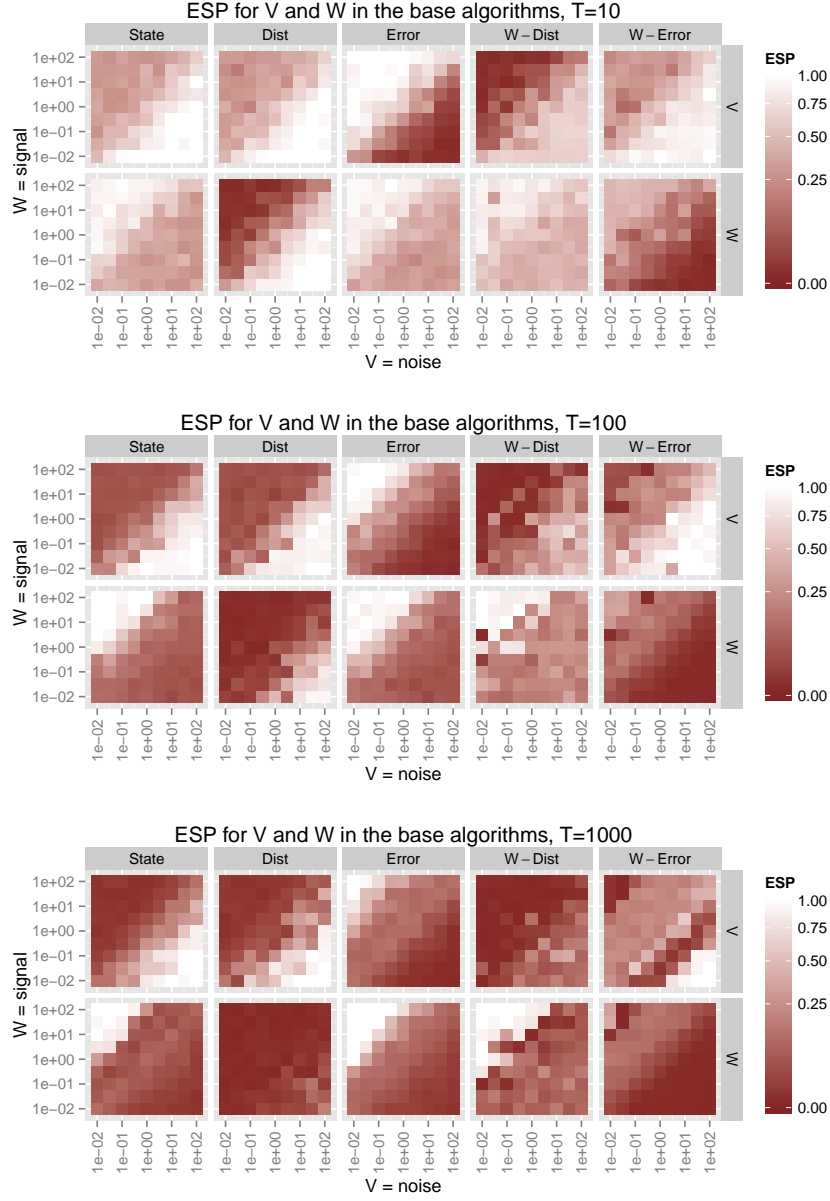


Figure 1: Effective sample proportion in the posterior sampler for a time series of lengths $T = 10$, $T = 100$, and $T = 1000$, for V and W , and for the state, scaled disturbance, scaled error, wrongly scaled disturbance, and wrongly scaled error samplers. X and Y axes indicate the true values of V and W respectively for the simulated data. Note that the signal-to-noise ratio is constant moving up any diagonal. In the upper left the signal is high, in the lower right the noise is high. Note that for plotting purposes, effective sample proportions larger than 1 were rounded down to 1

We can verify most of these intuitions in Figure 1. First, the State-Dist GIS algorithm has high ESP for W except for a narrow band where W/V is near one, though this band becomes much wider as T increases. The State-Dist GIS algorithm's mixing behavior for V appears identical to the original state sampler — high ESP when $W/V < 1$ and poor ESP when $W/V > 1$, and again the good region shrinks as T increases. So

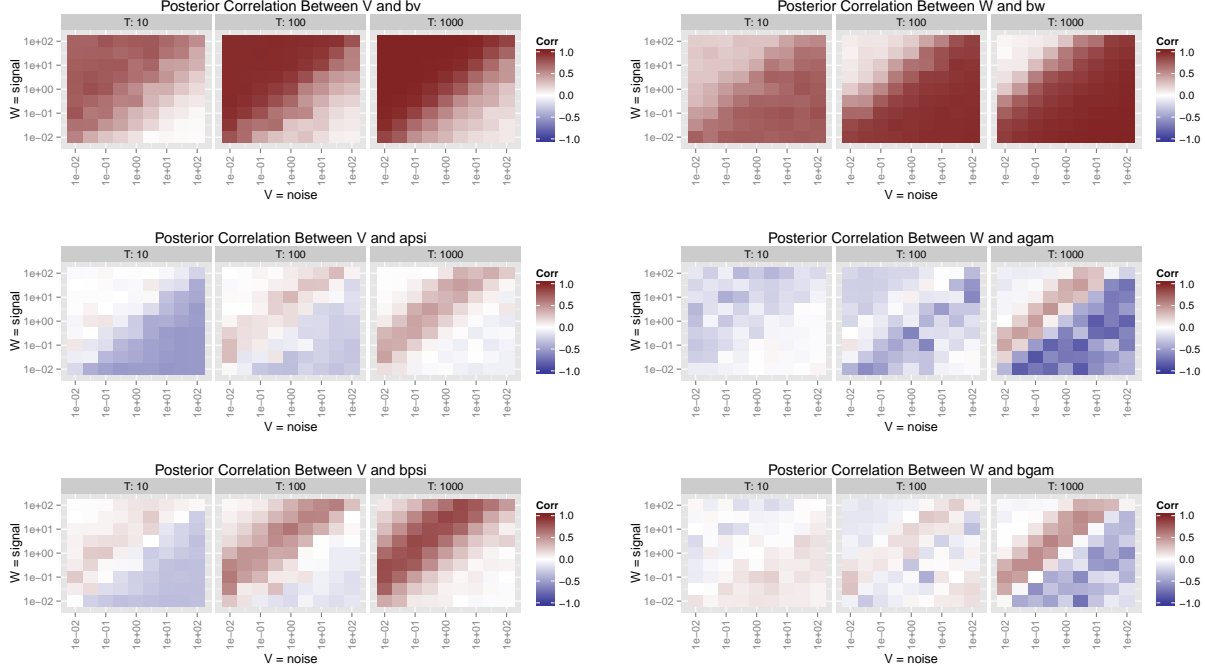


Figure 2: Posterior correlation between V or W and b_V or b_W . X and Y axes indicate the true values of V and W respectively for the simulated data. Note that the signal-to-noise ratio is constant moving up any diagonal. In the upper left the signal is high, in the lower right the noise is high.

this algorithm behaves as expected — it takes advantage of the fact that the state and scaled disturbance DA algorithms make up a “beauty and the beast” pair for W and thus improves mixing for W . However, the two underlying DA algorithms behave essentially identically for V so there is no improvement. Similarly the State-Error GIS algorithm’s ESP for W is essentially identical to the state and scaled error algorithms’ ESP for W — high when W/V large and low when W/V small. For V , the State-Error algorithm has a high ESP when W/V is far enough away from one, especially when T is small. The Dist-Error GIS algorithm also behaves as predicted — when W/V is not too close to one it has high ESP for both V and W , though as T increases W/V has to be farther away from one in order for the ESPs to be high. The Dist-Error GIS algorithm behaves apparently identically to the Full CIS and triple GIS algorithms, with some differences when T is small. The first of these is not surprising — based on the intuition that the Dist-Error GIS and Full CIS algorithms are the same up to a reordering of each of their steps, we expected little if any difference. However, we had some hope that the triple GIS algorithm would improve upon the Dist-Error GIS algorithm somewhat by further breaking the correlation between iterations in the Markov chain. This did not happen and furthermore the State-Dist and State-Error samplers did not improve the ESP for V or W respectively. When the underlying DA algorithms form a “beauty and the beast” pair, the interweaving algorithm appears to mix just as well as the best mixing single DA algorithm. Figure 4 allows us to see the ESP of the alternating algorithms in order to compare them to the GIS algorithms. There appears to be little practical difference between the alternating and interweaving versions of a given algorithm based on any two or three of the base DAs — both take advantage of the beauty and the beast nature of their underlying DAs but dependence between the two DAs prevent the GIS version from improving on the alternating version.

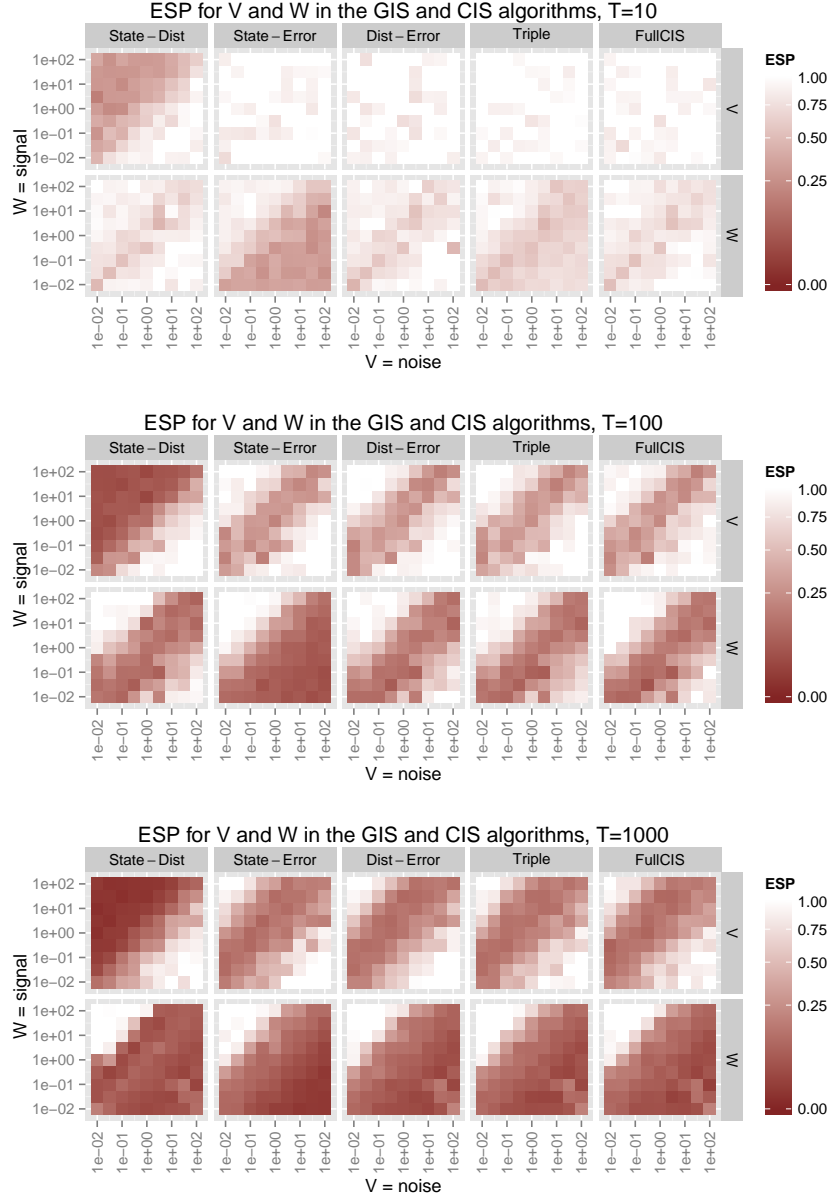


Figure 3: Effective sample proportion in the posterior sampler for V and W in for $T = 10$, $T = 100$, and $T = 1000$, in all three GIS samplers based on any two of the base samplers and the full CIS sampler. Horizontal and vertical axes indicate the true values of V and W respectively for the simulated data. Note that the signal-to-noise ratio is constant moving up any diagonal. In the upper left the signal is high, in the lower right the noise is high. Note that for plotting purposes, effective sample proportions larger than one were rounded down to one.

5.6 Computational time

A more important question than how well the chain mixes from a practical standpoint is the full computational time required to adequately characterize the posterior distribution. In order to investigate this, we compute the natural log of the average time in minutes required for each sampler to achieve an effective

sample size of 1000 — in other words the log time per 1000 effective draws. All simulations were performed on *INSERT DETAILS ABOUT IMPACT3*. While different systems will yield different absolute times, the relative times should at least be similar. Figure 5 contains plots of the log time per 1000 effective draws for both V and W and for each of the base and interweaving samplers — note that the scales change as T changes. For $T = 10$ the standard state sampler is competitive with the other samplers over most of the parameter space, though the Triple interweaving and Full CIS algorithms appear to offer some improvement. In any case, for small T computation usually is not a problem, so the difference is practically insignificant.

For $T = 100$ and $T = 1000$ the pattern we saw for ESP begins to emerge for time per 1000 effective draws. The state sampler becomes very slow to reach 1000 effective draws for V when $W/V > 1$ and for W when $W/V < 1$. The scaled disturbance and scaled error samplers behave as expected — the scaled disturbance sampler is slow for both V and W when $W/V > 1$ while the scaled disturbance sampler is slow for both V and W when $W/V < 1$. The Dist-Error GIS, Triple GIS and Full CIS algorithms appear to be the big winners here and are almost indistinguishable. All three algorithms are slightly slower for V when W/V is near one, and when W/V is near or below one all three are slow for W . Compared to the state sampler though, all three offer large gains over most of the parameter space. When we compare the GIS algorithms to their alternating counterparts in terms of log time per 1000 effective draws, again there is little difference. Figure 6 shows the log time per 1000 effective draws for the alternating algorithms and we see essentially the same pattern as we saw for the GIS algorithms.

6 Discussion

Our results on computational time in Section 5.6 should be taken with a grain of salt because we did not put much effort into efficiently sampling from $p(W|V, \gamma_{0:T}, y_{1:T})$ or $p(V|W, \psi_{0:T}, y_{1:T})$. Both densities have the form

$$p(x) \propto x^{-\alpha-1} \exp \left[-ax + b\sqrt{x} - \frac{\beta}{x} \right]$$

where $\alpha, \beta, a > 0$. This density is the same as the generalized inverse Gaussian distribution (see e.g. Jørgensen [1982]) when $b = 0$, but in general this is not the case in our application. It is possible that sampling from this density can be significantly improved in which case, the relative speed of the algorithms based on either the scaled errors or the scaled disturbances will improve significantly. **Perhaps look into the amount of time spent on the smoothing step vs elsewhere here**

In the data augmentation for multilevel models literature, a key quantity is called the fraction of missing information (Van Dyk and Meng [2001], for example). When ϕ is the model parameter, θ is the data augmentation and y is the data, the Bayesian fraction of missing information is defined as

$$\mathcal{F}_B = I - [\text{var}(\phi|y)]^{-1} E[\text{var}(\phi|\theta, y)|y]$$

while the EM fraction of missing information is defined as

$$\mathcal{F}_{EM} = I - I_{obs} I_{aug}^{-1}$$

where

$$I_{aug} = E \left[- \frac{\partial^2 \log p(\phi|\theta, y)}{\partial \phi \partial \phi} \middle| y, \phi, \right] \bigg|_{\phi=\phi^*}$$

is the expected augmented Fisher information matrix and

$$I_{obs} = - \frac{\partial^2 \log p(\phi|y)}{\partial \phi \partial \phi} \bigg|_{\phi=\phi^*}$$

is the observed fisher information matrix while ϕ^* is the posterior mode. The rate of convergence of the EM algorithm is governed by \mathcal{F}_{EM} while that of the Gibbs sampler is governed by \mathcal{F}_B — the larger the

spectral radius of \mathcal{F} , the slower the corresponding algorithm converges. This result applies to hierarchical models in a certain class, but not statespace models because of the dependence among the elements in the data augmentation (*IS THIS RIGHT? CHECK*). Nevertheless, Frühwirth-Schnatter [2004] use this result as a guide to exploring different DAs in a dynamic regression with AR(1) regression coefficient and they find that as long as the estimated autocorrelation in the coefficient is low, the multilevel model results roughly apply.

The significance of the signal-to-noise ratio in our results is likely related to this despite the fact that we do not have a result about the fraction of missing information in a time series setting. In particular, Pitt and Shephard [1999] find that the signal to noise ratio along with the AR(1) coefficient determine the convergence rate of a Gibbs sampler in an AR(1) plus noise model. In addition, they find that as the length of the time series increases, the convergence rate slows down and compute asymptotic convergence rates for an infinite time series. However, Pitt and Shephard assume that both of the variances are fixed in order to derive their results. In a continuous-time model, Roberts et al. [2004] find that the Gibbs sampler based on a NCP is at least as efficient as the Gibbs sampler based on the CP and sometimes it is much more efficient, depending on the true values of the parameters. It is unclear whether the time series dependence or something unique to the model is driving this, but at least this suggests that in the time series case convergence and mixing properties can sometimes be wildly different than what we might expect and calls out the need for a time series or statespace fraction of missing information.

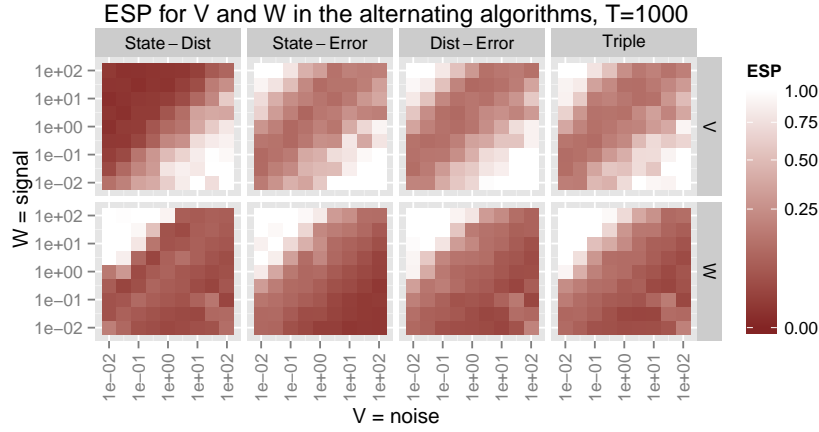
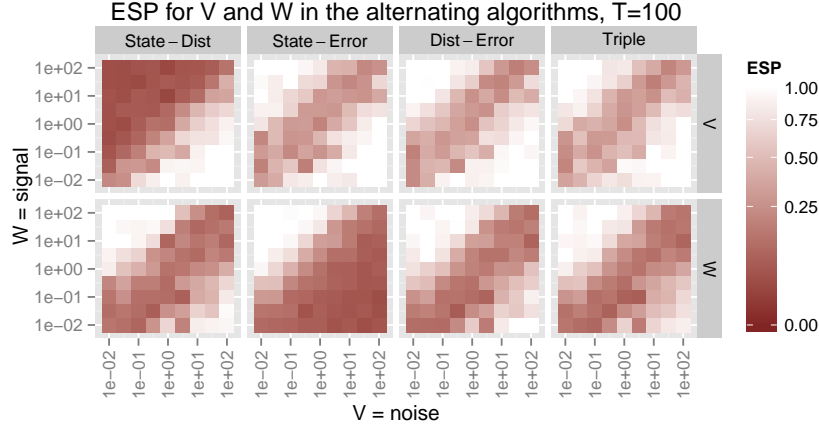
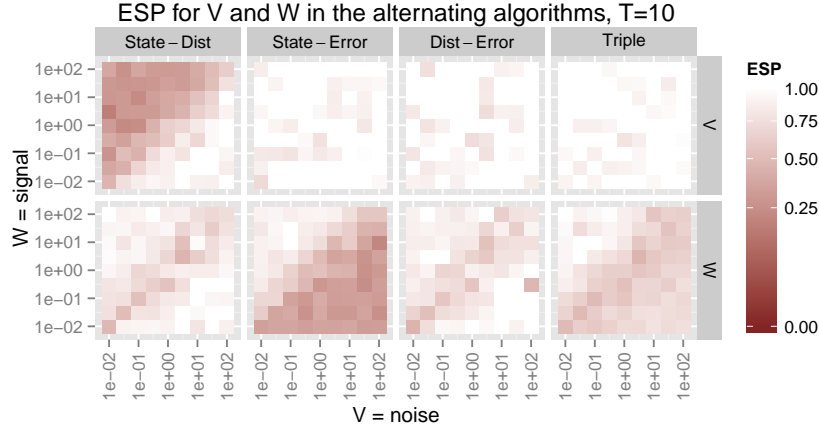


Figure 4: Effective sample proportion in the posterior sampler for V and W in for $T = 10$, $T = 100$, and $T = 1000$, in all four alternating samplers. Horizontal and vertical axes indicate the true values of V and W respectively for the simulated data. Note that the signal-to-noise ratio is constant moving up any diagonal. In the upper left the signal is high, in the lower right the noise is high. Note that for plotting purposes, effective sample proportions larger than one were rounded down to one.

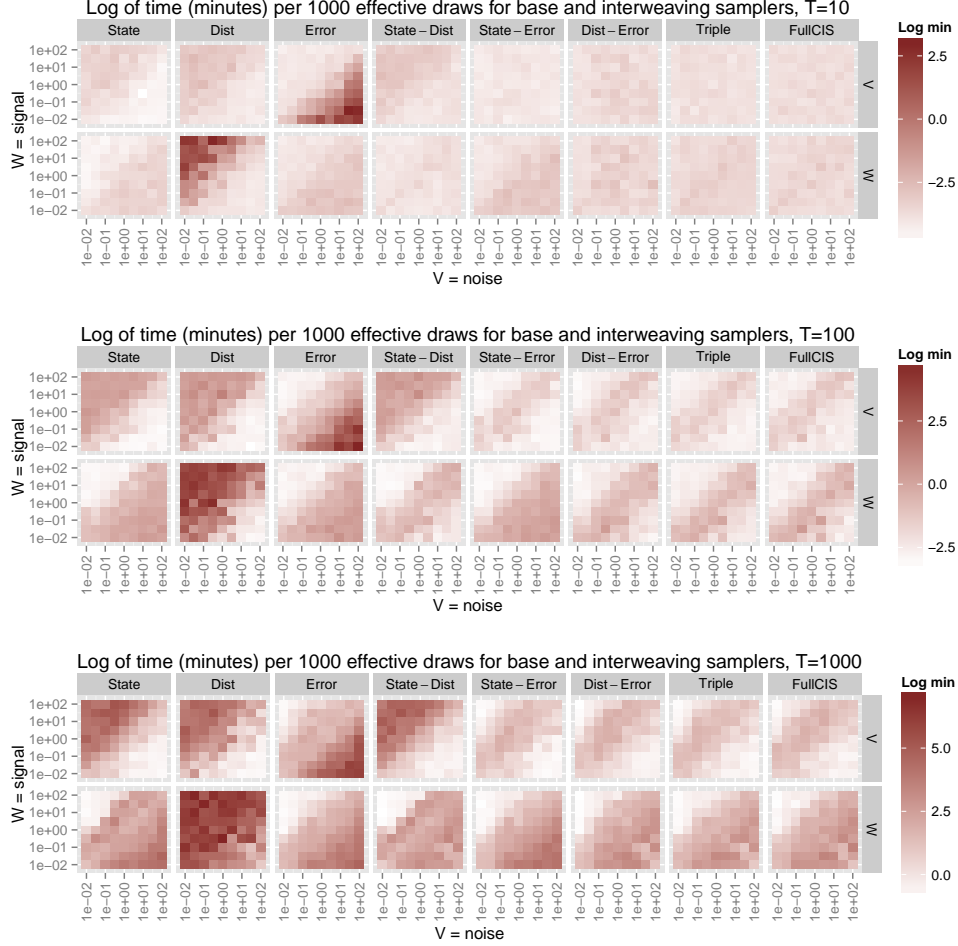


Figure 5: Log of the time in minutes per 1000 effective draws in the posterior sampler for V and W , for $T = 10$, $T = 100$, and $T = 1000$, in the state, scaled disturbance and scaled error samplers and for all five interweaving samplers. Horizontal and vertical axes indicate the true values of V and W respectively for the simulated data. The signal-to-noise ratio is constant moving up any diagonal. In the upper left the signal is high, in the lower right the noise is high. For plotting purposes, times larger than the top of the scale are displayed in bright red.

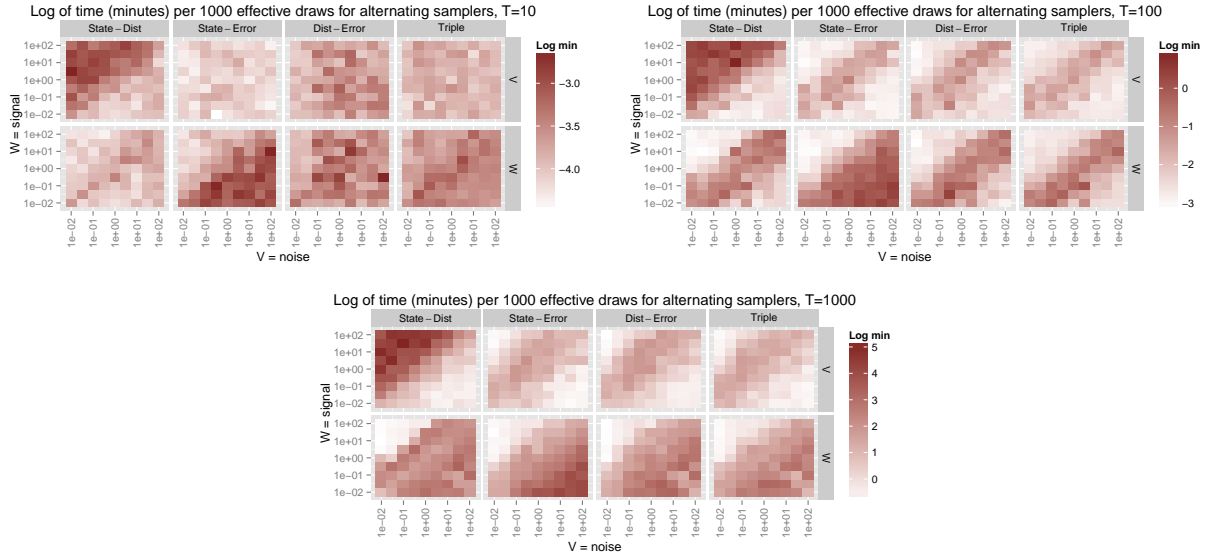


Figure 6: Log of the time in minutes per 1000 effective draws in the posterior sampler for V and W , for $T = 10$, $T = 100$, and $T = 1000$, in the alternating, GIS, and random kernel samplers. Horizontal and vertical axes indicate the true values of V and W respectively for the simulated data. The signal-to-noise ratio is constant moving up any diagonal. In the upper left the signal is high, in the lower right the noise is high. For plotting purposes, times larger than the top of the scale are displayed in bright red.

References

- Chris K Carter and Robert Kohn. On Gibbs sampling for state space models. *Biometrika*, 81(3):541–553, 1994.
- Piet De Jong and Neil Shephard. The simulation smoother for time series models. *Biometrika*, 82(2):339–350, 1995.
- Sylvia Frühwirth-Schnatter. Data augmentation and dynamic linear models. *Journal of Time Series Analysis*, 15(2):183–202, 1994.
- Sylvia Frühwirth-Schnatter. Efficient Bayesian parameter estimation for state space models based on reparameterizations. *State Space and Unobserved Component Models: Theory and Applications*, pages 123–151, 2004.
- Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin. Bayesian data analysis. 2003.
- Walter R Gilks and Pascal Wild. Adaptive rejection sampling for Gibbs sampling. *Applied Statistics*, pages 337–348, 1992.
- James P Hobert and Dobrin Marchev. A theoretical comparison of the data augmentation, marginal augmentation and PX-DA algorithms. *The Annals of Statistics*, 36(2):532–554, 2008.
- Bent Jørgensen. *Statistical properties of the generalized inverse Gaussian distribution*. Springer, 1982.
- Gregor Kastner and Sylvia Frühwirth-Schnatter. Ancillarity-sufficiency interweaving strategy (asis) for boosting mcmc estimation of stochastic volatility models. *Computational Statistics & Data Analysis*, 2013.
- Siem Jan Koopman. Disturbance smoother for state space models. *Biometrika*, 80(1):117–126, 1993.
- Jun S Liu and Ying Nian Wu. Parameter expansion for data augmentation. *Journal of the American Statistical Association*, 94(448):1264–1274, 1999.
- William J McCausland, Shirley Miller, and Denis Pelletier. Simulation smoothing for state-space models: A computational efficiency analysis. *Computational Statistics & Data Analysis*, 55(1):199–212, 2011.
- X-L Meng and David Van Dyk. Fast em-type implementations for mixed effects models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(3):559–578, 1998.
- Xiao-Li Meng and David A Van Dyk. Seeking efficient data augmentation schemes via conditional and marginal augmentation. *Biometrika*, 86(2):301–320, 1999.
- Omiros Papaspiliopoulos, Gareth O Roberts, and Martin Sköld. A general framework for the parametrization of hierarchical models. *Statistical Science*, pages 59–73, 2007.
- Michael K Pitt and Neil Shephard. Analytic convergence rates and parameterization issues for the gibbs sampler applied to state space models. *Journal of Time Series Analysis*, 20(1):63–85, 1999.
- Gareth O Roberts, Omiros Papaspiliopoulos, and Petros Dellaportas. Bayesian inference for non-gaussian ornstein–uhlenbeck stochastic volatility processes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(2):369–393, 2004.
- David A Van Dyk and Xiao-Li Meng. The art of data augmentation. *Journal of Computational and Graphical Statistics*, 10(1), 2001.
- Yaming Yu and Xiao-Li Meng. To center or not to center: That is not the question - an ancillarity-sufficiency interweaving strategy (ASIS) for boosting MCMC efficiency. *Journal of Computational and Graphical Statistics*, 20(3):531–570, 2011.