#### Unknown fixed parameters in DLMs

Dr. Jarad Niemi

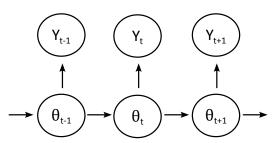
STAT 6150 - Iowa State University

October 14, 2025

#### Dynamic Linear Models

$$\begin{array}{lll} Y_t &= F_t \theta_t + v_t & v_t & \stackrel{ind}{\sim} N_m(0,V_t) \\ \theta_t &= G_t \theta_{t-1} + w_t & w_t & \stackrel{ind}{\sim} N_p(0,W_t) \\ & \theta_0 &\sim N_p(m_0,C_0) \end{array}$$

where  $v_t$  and  $w_t$  are independent across time and all are independent of  $\theta_0$ .



### Unknown parameters in polynomial trend models

What is known?

• 
$$F_t = (1, 0, \dots, 0)$$

$$G_t = G = \begin{bmatrix} 1 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 1 & 0 & \cdots & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & \cdots & 0 & 1 & 1 \\ 0 & \cdots & & 0 & 1 \end{bmatrix}$$

What are the unknown parameters?

- $\bullet \theta_t$
- $V_t \stackrel{?}{=} V$   $W_t \stackrel{?}{=} W$

# Unknown parameters in seasonal models

#### What is known?

- $\bullet$   $F_t$ 
  - Seasonal factor:  $F_t = (1, 0, \dots, 0)$
  - Fourier form:  $F_t = (1, 0, 1, 0, \dots, 1, 0)$
- $\bullet$   $G_t = G$ 
  - Seasonal factor: rotation matrix
  - ullet Fourier form: block diagonal with blocks  $H_j$

$$H_j = \begin{bmatrix} \cos \omega_j & \sin \omega_j \\ -\sin \omega_j & \cos \omega_j \end{bmatrix}$$

#### What are the unknown parameters?

- $\bullet$   $\theta_t$
- $V_t \stackrel{?}{=} V$
- $W_t \stackrel{?}{=} W \stackrel{?}{=} 0$

### Unknown parameters in dynamic regression models

What is known?

- $F_t = x_t$
- $G_t \stackrel{?}{=} G \stackrel{?}{=} I$

What are the unknown parameters?

- $\bullet$   $\theta_t$
- $V_t \stackrel{?}{=} V \stackrel{?}{=} \sigma^2 I$  or  $\sigma^2 D$
- $W_t \stackrel{?}{=} W \stackrel{?}{=} D$

#### The bottom line is...

- In all of these univariate models,
  - the unknowns are  $\theta_t$ ,  $W_t$ , and  $V_t$ ,
  - $oldsymbol{ heta}_t$  has always been unknown
  - and often,  $W_t = W$  and  $V_t = V$ .
- In our multivariate models,
  - commonly  $W_t=W$  and  $V_t=V$ , but now they are (block-diagonal) matrices.

For a parameter vector  $\psi$  and data vector y, the likelihood function

$$L(\psi) \propto p(y|\psi).$$

The maximum likelihood estimate is

$$\label{eq:psi_def} \hat{\psi} = \mathrm{argmax}_{\psi} L(\psi).$$

Which is equivalent to

$$\hat{\psi} = \mathrm{argmax}_{\psi} \ell(\psi)$$

where  $\ell(\psi) = \log L(\psi)$ .

#### Likelihood function for DLMs

If  $\psi = (W, V)$ , what is  $L(\psi)$  for a general DLM?

What do we know?

$$p(y_t|\theta_t, V) = N(y_t; F_t\theta_t, V)$$

$$p(\theta_t|\theta_{t-1}, W) = N(\theta_t; G_t\theta_{t-1}, W)$$

$$p(\theta_0) = N_p(m_0, C_0)$$

$$p(\theta_t|y_{1:t-1}, \psi) = N(a_t, R_t)$$
  
 $p(y_t|y_{1:t-1}, \psi) = N(f_t, Q_t)$ 

$$p(y|\psi) = \prod_{t=1}^{n} p(y_t|y_{1:t-1}, \psi)$$

### Finding MLEs for DLMs

If  $y_t$  is multivariate, the likelihood function is

$$L(\psi) \propto \prod_{t=1}^{n} \frac{1}{(2\pi)^{k/2} |Q_t|^{1/2}} \exp\left(-\frac{1}{2}(y_t - f_t)^{\top} Q_t^{-1}(y_t - f_t)\right).$$

Log-likelihood function

$$\ell(\psi) = C + -\frac{1}{2} \sum_{t=1}^{n} \log |Q_t| - \frac{1}{2} \sum_{t=1}^{n} (y_t - f_t)^{\top} Q_t^{-1} (y_t - f_t).$$

The MLE is then

$$\hat{\psi} = \mathrm{argmax}_{\psi} \ell(\psi)$$

The R function dlmMLE does all of this for you.

### Bayesian inference

What do we have to specify to perform Bayesian inference, i.e. parameter estimation, for data y?

- A statistical model  $p(y|\psi)$
- A prior  $p(\psi)$

What is the objective of Bayesian inference?

• The posterior  $p(\psi|y) \propto p(y|\psi)p(\psi)$ .

# Conjugacy

Conjugate Bayesian inference is one where if

$$\psi \sim f(\alpha) \implies \psi | y \sim f(\alpha').$$

#### Remember the examples

- $y \sim N(\mu, I), \mu \sim N(\cdot, \cdot) \implies \mu | y \sim N(\cdot, \cdot)$
- $y \sim N(0, \phi^{-1}I), \phi \sim Ga(\cdot, \cdot) \implies \phi|y \sim Ga(\cdot, \cdot)$
- $y \sim N(\mu, \phi^{-1}I), \mu, \phi \sim NG(\cdot) \implies \mu, \phi | y \sim NG(\cdot)$
- $y \sim N(X\beta, \phi^{-1}I), \beta, \phi \sim NG(\cdot) \implies \beta, \phi|y \sim NG(\cdot)$
- $y \sim Bin(n, p), p \sim Be(\cdot, \cdot) \implies p|y \sim Be(\cdot, \cdot)$

#### What are the unknowns in DLMs?

So for  $\psi = (F_{1:n}, G_{1:n}, W_{1:n}, V_{1:n})$ , we are looking for

$$\psi \sim f(\alpha) \implies \psi | y \sim f(\alpha').$$

This only happens in simple examples. Today, we will discuss

- $\bullet \ V_t = \phi^{-1} \tilde{V}_t, W_t = \phi^{-1} \tilde{W}_t, C_0 = \phi^{-1} \tilde{C}_0$
- ullet  $W_t$  specified by a discount factor
- Evolving  $\phi = 1/\sigma^2$

# common $\phi^{-1}$

$$Y_{t} = F_{t}\theta_{t} + v_{t} \qquad v_{t} \stackrel{ind}{\sim} N_{m}(0, \phi^{-1}\tilde{V}_{t})$$

$$\theta_{t} = G_{t}\theta_{t-1} + w_{t} \qquad w_{t} \stackrel{ind}{\sim} N_{p}(0, \phi^{-1}\tilde{W}_{t})$$

$$\theta_{0} \sim N_{p}(m_{0}, \phi^{-1}\tilde{C}_{0})$$

$$\phi \sim Ga(\alpha_{0}, \beta_{0})$$

#### Everything is known except

- ullet  $\theta_t$  for all t
- φ

Starting with

$$\theta_{t-1}, \phi | y_{1:t-1} \sim NG(m_{t-1}, \tilde{C}_{t-1}, \alpha_{t-1}, \beta_{t-1})$$

One step ahead prior

$$\theta_t, \phi | y_{1:t-1} \sim NG(a_t, \tilde{R}_t, \alpha_{t-1}, \beta_{t-1})$$

where  $a_t = G_t m_{t-1}$  and  $\tilde{R}_t = G_t \tilde{C}_{t-1} G_t^\top + \tilde{W}_t$ .

One step ahead predictive density

$$Y_t|y_{1:t-1} \sim t_{2\alpha_{t-1}}(f_t, \tilde{Q}_t\beta_{t-1}/\alpha_{t-1})$$

with  $f_t = F_t a_t$  and  $\tilde{Q}_t = F_t \tilde{R}_t F_t^\top + \tilde{V}_t$ .

Filtering density

$$\theta_t, \phi | y_{1:t} \sim NG(m_t, \tilde{C}_t, \alpha_t, \beta_t)$$

with

$$m_{t} = a_{t} + \tilde{R}_{t} F_{t} \tilde{Q}_{t}^{-1} (y_{t} - f_{t})$$

$$\tilde{C}_{t} = \tilde{R}_{t} - \tilde{R}_{t} F_{t}^{\top} \tilde{Q}_{t}^{-1} \tilde{R}_{t}^{\top}$$

$$\alpha_{t} = \alpha_{t-1} + \frac{m}{2}$$

$$\beta_{t} = \beta_{t-1} + \frac{1}{2} (y_{t} - f_{t})^{\top} \tilde{Q}_{t}^{-1} (y_{t} - f_{t})$$

#### Discount factor

Let's specify how adaptive we want our model to be.

- Do this by specifying  $W_t$  relative to  $V_t$  and  $C_t$  using a discount factor  $\delta \in (0,1]$ .
- $oldsymbol{\circ}$   $\delta=1$  means no loss of information, i.e.  $W_t=0$
- $\delta = 0$  means no information retained
- Often  $\delta > 0.9$

To implement, set

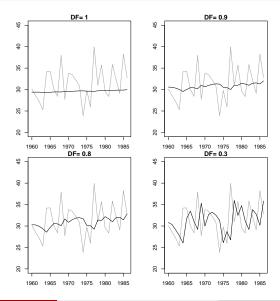
$$W_t = \frac{1 - \delta}{\delta} G_t^{\top} C_{t-1} G_t$$

or

$$\tilde{W}_t = \frac{1 - \delta}{\delta} G_t^{\top} \tilde{C}_{t-1} G_t$$

if using a common  $\sigma^2$ .

#### Discount factor effect

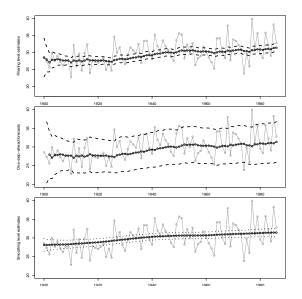


#### Choosing the discount factor

Specify  $\delta$  based on one-step ahead prediction errors.

```
DF MAPE MAD MSE sigma2  
1.0 0.10 3.02 21.54 12.00  
0.9 0.09 2.86 19.92 9.64  
0.8 0.10 2.87 20.29 8.94  
0.3 0.11 3.42 25.12 5.07  
Last column is posterior expectation for \sigma^2, i.e. E[\sigma^2|y_{1:187}].
```

### Inference for $\delta = 0.95$ on Lake Superior Data



# Evolving $\phi_t$

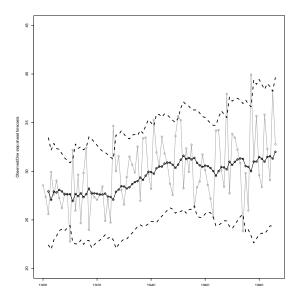
#### Choose $\delta^* \in (0,1)$

- $\phi_{t-1}|y_{1:t-1} \sim Ga(\alpha_{t-1}, \beta_{t-1})$
- $\phi_t|y_{1:t-1} \sim Ga(\delta^*\alpha_{t-1}, \delta^*\beta_{t-1})$

#### What is

- $E[\phi_t|y_{1:t-1}] = E[\phi_{t-1}|y_{1:t-1}]$
- $Var[\phi_t|y_{1:t-1}] = \frac{1}{\delta^*} Var[\phi_{t-1}|y_{1:t-1}]$

## Evolving $\phi_t$ for Lake Superior data



The situations for conjugate Bayesian analysis are small, therefore we need more advanced techniques.

# Gibbs sampling algorithm

Start with an initial guess for all parameters and call it  $\psi^{(0)}$ . Set j=1.

- 1. Sample  $\psi_1^{(j)} \sim p\left(\psi_1|\psi_2^{(j-1)},\dots,\psi_K^{(j-1)},y\right)$
- 2. Sample  $\psi_2^{(j)} \sim p\left(\psi_2|\psi_1^{(j)},\psi_3^{(j-1)},\dots,\psi_K^{(j-1)},y\right)$
- 3. :
- 4. Sample  $\psi_k^{(j)} \sim p\left(\psi_k|\psi_1^{(j)},\dots,\psi_{k-1}^{(j)},\psi_{k+1}^{(j-1)},\dots,\psi_K^{(j-1)},y\right)$
- 5. :
- 6. Sample  $\psi_{K-1}^{(j)} \sim p\left(\psi_{K-1}|\psi_1^{(j)},\dots,\psi_{K-2}^{(j)},\psi_K^{(j-1)},y\right)$
- 7. Sample  $\psi_K^{(j)} \sim p\left(\psi_K|\psi_1^{(j)},\dots,\psi_{K-1}^{(j)},y\right)$
- 8. If j < J, j = j + 1 and return to step 1.

### What full conditionals are required?

Suppose our goal is to draw from  $p(\theta_{0:T}|y_{1:T})$  using univariate Gibbs sampling. We will implicitly assume conditioning on any other unknown parameters.

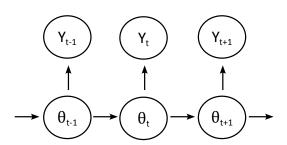
- What are the required full condition distributions?
  - $p(\theta_0|\theta_{1:T}, y_{1:T})$
  - $p(\theta_t|\theta_{-t},y_{1:T})$  where  $\theta_{-t}$  is  $\theta_{0:T}$  with the  $t^{th}$  element removed
  - $p(\theta_T | \theta_{0:T-1}, y_{1:T})$

#### **DLMs**

$$Y_t = F_t \theta_t + v_t \qquad v_t \stackrel{ind}{\sim} N_m(0, V_t)$$
  

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

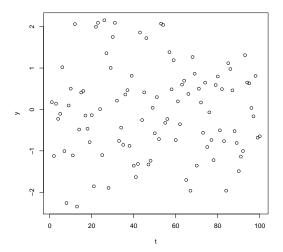
$$\theta_0 \sim N_p(m_0, C_0)$$



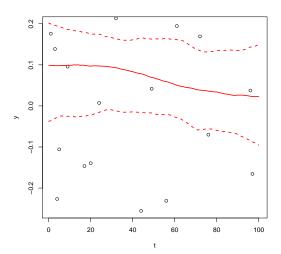
#### What are the full conditionals?

$$\begin{split} p(\theta_0|\dots) &=& p(\theta_0|\theta_1) \\ &\propto & N(\theta_1;G_1\theta_0,W_1)N(\theta_0;m_0,C_0) \\ &\propto & N(\theta_0;k_0,K_0) \\ K_0 &=& (C_0^{-1}+G_1^\top W_1^{-1}G_1)^{-1} \\ k_0 &=& K_0(C_0^{-1}m_0+G_1^\top W_1^{-1}\theta_1) \\ \end{split}$$
 
$$p(\theta_T|\dots) &=& p(\theta_T|\theta_{T-1},y_T) \\ &\propto & N(y_T;F_T\theta_T,V_T)N(\theta_T;G_T\theta_{T-1},W_T) \\ &\propto & N(\theta_T;k_T,K_T) \\ K_T &=& (W_T^{-1}+F_T^\top V_T^{-1}F_T)^{-1} \\ k_T &=& K_T(W_T^{-1}G_T\theta_{T-1}+F_T^\top V_T^{-1}y_T) \\ \end{split}$$
 
$$p(\theta_t|\dots) &=& p(\theta_t|\theta_{t-1},\theta_{t+1},y_t) \\ &\propto & N(\theta_t;k_t,K_t) \\ K_t &=& (W_t^{-1}+F_t^\top V_t^{-1}F_t+G_{t+1}^\top W_t)N(\theta_t;G_t\theta_{t-1},W_{t+1}) \\ &\propto & N(\theta_t;k_t,K_t) \\ K_t &=& (W_t^{-1}+F_t^\top V_t^{-1}F_t+G_{t+1}^\top W_{t+1}^{-1}G_{t+1})^{-1} \\ k_t &=& K_t(W_t^{-1}G_t\theta_{t-1}+F_t^\top V_t^{-1}y_t+G_{t+1}^\top W_{t+1}^{-1}\theta_{t+1}) \end{split}$$

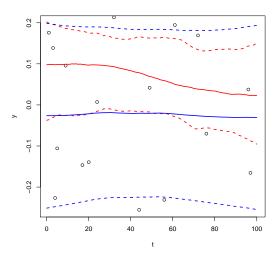
# Consider the local level model with V=1 and $W=0.01^2$ .



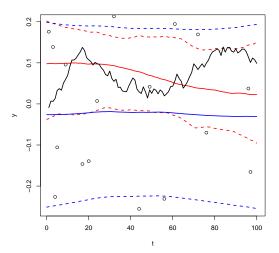
### Univariate Gibbs sampling for states



### Exact quantiles for states



# True underlying state



### **Filtering**

Goal:  $p(\theta_t|y_{1:t})$  where  $y_{1:t} = (y_1, y_2, \dots, y_t)$  (filtered distribution)

Recursive procedure:

- Assume  $p(\theta_{t-1}|y_{1:t-1})$
- Prior for  $\theta_t$

$$\begin{array}{lcl} p(\theta_t|y_{1:t-1}) & = & \int p(\theta_t,\theta_{t-1}|y_{1:t-1})d\theta_{t-1} \\ \\ & = & \int p(\theta_t|\theta_{t-1},y_{1:t-1})p(\theta_{t-1}|y_{1:t-1})d\theta_{t-1} \\ \\ & = & \int p(\theta_t|\theta_{t-1})p(\theta_{t-1}|y_{1:t-1})d\theta_{t-1} \end{array}$$

lacksquare One-step ahead predictive distribution for  $y_t$ 

$$\begin{array}{lll} p(y_t|y_{1:t-1}) & = & \int p(y_t,\theta_t|y_{1:t-1})d\theta_t \\ \\ & = & \int p(y_t|\theta_t,y_{1:t-1})p(\theta_t|y_{1:t-1})d\theta_t \\ \\ & = & \int p(y_t|\theta_t)p(\theta_t|y_{1:t-1})d\theta_t \end{array}$$

• Filtered distribution for  $\theta_t$ 

$$p(\theta_t|y_{1:t}) = \frac{p(y_t|\theta_t, y_{1:t-1})p(\theta_t|y_{1:t-1})}{p(y_t|y_{1:t-1})} = \frac{p(y_t|\theta_t)p(\theta_t|y_{1:t-1})}{p(y_t|y_{1:t-1})}$$

### **Smoothing**

#### Goal: $p(\theta_t|y_{1:T})$ for t < T

• Backward transition probability  $p(\theta_t | \theta_{t+1}, y_{1:T})$ 

$$\begin{array}{lll} p(\theta_{t}|\theta_{t+1},y_{1:T}) & = & p(\theta_{t}|\theta_{t+1},y_{1:t}) \\ \\ & = & \frac{p(\theta_{t+1}|\theta_{t},y_{1:t})p(\theta_{t}|y_{1:t})}{p(\theta_{t+1}|y_{1:t})} \\ \\ & = & \frac{p(\theta_{t+1}|\theta_{t})p(\theta_{t}|y_{1:t})}{p(\theta_{t+1}|y_{1:t})} \end{array}$$

 $\qquad \text{Recursive smoothing distributions } p(\theta_t|y_{1:T}) \text{ assuming we know } p(\theta_{t+1}|y_{1:T})$ 

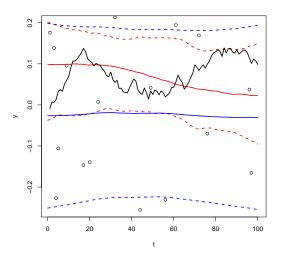
$$\begin{split} p(\theta_t|y_{1:T}) &= \int p(\theta_t,\theta_{t+1}|y_{1:T})d\theta_{t+1} \\ &= \int p(\theta_{t+1}|y_{1:T})p(\theta_t|\theta_{t+1},y_{1:T})d\theta_{t+1} \\ &= \int p(\theta_{t+1}|y_{1:T})\frac{p(\theta_{t+1}|\theta_t)p(\theta_t|y_{1:t})}{p(\theta_{t+1}|y_{1:t})}d\theta_{t+1} \\ &= p(\theta_t|y_{1:t})\int \frac{p(\theta_{t+1}|\theta_t)}{p(\theta_{t+1}|y_{1:t})}p(\theta_{t+1}|y_{1:T})d\theta_{t+1} \end{split}$$

Start from  $p(\theta_T|y_{1:T})$ .

#### Kalman smoother

If 
$$p(\theta_{t+1}|y_{1:T}) = N(s_{t+1}, S_{t+1})$$
, then 
$$p(\theta_t|\theta_{t+1}, y_{1:T}) = p(\theta_t|\theta_{t+1}, y_{1:t}) \\ \propto p(\theta_{t+1}|\theta_t, y_{1:t})p(\theta_t|y_{1:t}) \\ = N(\theta_{t+1}; G_{t+1}\theta_t, W_{t+1})N(\theta_t; m_t, C_t) \\ \propto N(\theta_t; h_t, H_t) \\ H_t = (C_t^{-1} + G_{t+1}^\top W_{t+1}^{-1}G_{t+1})^{-1} \\ h_t = H_t(C_t^{-1}m_t + G_{t+1}^\top W_{t+1}^{-1}\theta_{t+1}) \\ p(\theta_t|y_{1:T}) = \int p(\theta_t|\theta_{t+1}, y_{1:T})p(\theta_{t+1}|y_{1:T})d\theta_{t+1} \\ = N(\theta_t; s_t, S_t) \\ S_t = C_t - C_t G_{t+1}^\top R_{t+1}^{-1}(R_{t+1} - S_{t+1})R_{t+1}^{-1}G_{t+1}C_t \\ s_t = m_t + C_t G_{t+1}^\top R_{t+1}^{-1}(s_{t+1} - a_{t+1})$$

# True underlying state



Let  $\psi$  represent any unknown, non-dynamic model parameters such that the data follows a DLM conditional on  $\psi$ .

$$\begin{array}{lll} Y_t &= F_t(\psi)\theta_t + v_t & v_t & \stackrel{ind}{\sim} N_m(0,V_t(\psi)) \\ \theta_t &= G_t(\psi)\theta_{t-1} + w_t & w_t & \stackrel{ind}{\sim} N_p(0,W_t(\psi)) \\ p(\theta_0) &= N_p(m_0(\psi),C_0(\psi)) & \end{array}$$

For example,  $\psi=(V,W)$  where  $V_t(\psi)=V$  and  $W_t(\psi)=W$  while  $F_t(\psi)=F$  and  $G_t(\psi)=G$  are known, as in polynomial trend, seasonal factor, and dynamic regression models.

- The Bayesian inferential objective is then  $p(\theta_{0:T}, \psi | y_{1:T})$ .
- While  $p(\theta_{0:T}|y_{1:T},\psi)$  is known analytically, generally  $p(\theta_{0:T},\psi|y_{1:T})$  is not.
- So resort to numerical methods, most often MCMC

#### MCMC Schemes

#### Scheme I - all univariate samples

- For  $t \in \{0, 1, \dots, T\}$  sample  $p(\theta_t | \dots)$ .
- For  $j \in \{1, \dots, J\}$  sample  $p(\psi_j | \dots)$  for J parameters.

#### Scheme II - block sampling of states

- Sample  $p(\theta_{0:T}|\ldots)$ .
- For  $j \in \{1, \dots, J\}$  sample  $p(\psi_j | \dots)$  for J parameters.

# MCMC Schemes (cont.)

Scheme III - block sampling of parameters

- Sample  $p(\theta_{0:T}|\ldots)$ .
- Sample  $p(\psi|\ldots)$ .

e.g. polynomial trend, seasonal factor, and dynamic regression models

#### Scheme IV - hybrid

- Sample  $p(\psi_{J'}|\psi_{J\setminus J'},y_{1:T})$  for some subset J' of parameters.
- Sample  $p(\theta_{0:T}|\ldots)$ .
- Sample  $p(\psi_{J\setminus J'}|\ldots)$ .

### MCMC Schemes

Generally better to jointly sampling unknowns, a.k.a. block sampling.

- Scheme I has all univariate draws
- Scheme II samples latent state jointly
- Scheme III samples latent state jointly and parameters jointly
- Scheme IV samples some parameters  $\psi_{J'}$  and all latent states jointly and then samples remaining parameters jointly

Bottom line: if parameters are highly correlated in the posterior, it is better to sample those parameters jointly.

# Forward filtering backward sampling (FFBS)

#### Recall

- $p(\theta_T|y_{1:T}) = N(m_T, C_T)$  is available from filtering
- $p(\theta_t|\theta_{t+1},y_{1:T}) = N(h_t,H_T)$  is available from smoothing

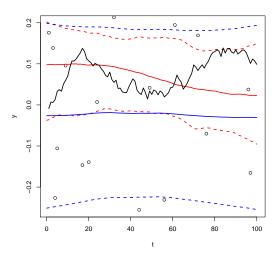
$$H_t = (C_t^{-1} + G_{t+1}^{\top} W_{t+1}^{-1} G_{t+1})^{-1}$$

$$h_t = H_t (C_t^{-1} m_t + G_{t+1}^{\top} W_{t+1}^{-1} \theta_{t+1})$$

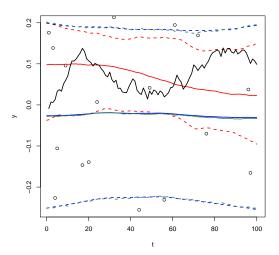
#### The algorithm is then

- Forward filter to obtain  $p(\theta_t|y_{1:t}) = N(m_t, C_t)$  for all t.
- Sample  $\theta_T \sim H(m_T, C_T)$ .
- For  $t \in \{T-1, T-2, \dots, 1, 0\}$ ,
  - Calculate  $h_t$  and  $H_t$  based on  $\theta_{t+1}$ .
  - Draw  $\theta_t \sim N(h_t, H_T)$ .

## Local level model



### Local level model



### Local level model - unknown variances

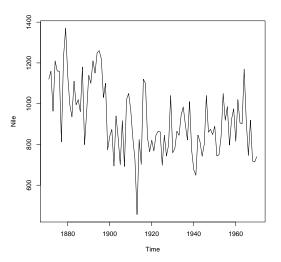
$$Y_t = \theta_t + v_t \qquad v_t \stackrel{ind}{\sim} N_m(0, V)$$

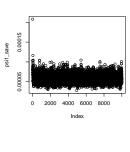
$$\theta_t = \theta_{t-1} + w_t \qquad w_t \stackrel{ind}{\sim} N_p(0, W)$$

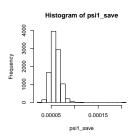
$$p(\theta_0) = N(m_0, C_0)$$

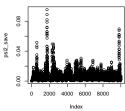
#### MCMC Scheme:

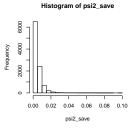
- Sample  $p(\theta_{0:T}|\ldots)$  using FFBS
- Sample  $p(V, W| \ldots)$

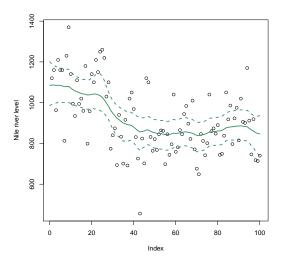












### MCMC in DLMs

Recall the inferential objective of the Bayesian approach in DLMs:

$$p(\theta_{0:n}, \psi|y_{1:n})$$

Since  $p(\theta_{0:n}, \psi | y_{1:n})$  is not typically available analytically, we commonly use Markov chain Monte Carlo. These approaches sample from full conditional distributions, e.g.

- $p(\theta_t | \theta_{-t}, \psi, y_{1:n})$  for  $t = 0, 1, 2, \dots, n$ .
- $p(\psi_j|\theta_{0:n},\psi_{-j},y_{1:n})$  for  $j=1,2,\ldots,J$ .

These draws could be Gibbs or Metropolis-Hastings.

# MCMC Univariate Sampling Activity

Fill in ? with i or i-1.

$$\theta_0^{(?)} \sim p(\theta_0|\theta_1^{(?)},\dots,\theta_n^{(?)},\psi^{(?)},y_{1:n})$$

For  $t \in 1, \ldots, n-1$ , sample from

$$\theta_t^{(?)} \sim p(\theta_t | \theta_0^{(?)}, \dots, \theta_{t-1}^{(?)}, \theta_{t+1}^{(?)}, \dots, \theta_n^{(?)}, \psi_{1:n}^{(?)})$$

$$\theta_n^{(?)} \sim p(\theta_n | \theta_0^{(?)}, \dots, \theta_{n-1}^{(?)}, \psi^{(?)}, y_{1:n})$$

$$\psi_1^{(?)} \sim p(\psi_1|\theta^{(?)}, \psi_2^{(?)}, \dots, \psi_J^{(?)}, y_{1:n})$$

For 
$$j \in 2, \dots, J-1$$
, sample from

$$\psi_{j}^{(?)} \sim p(\psi_{j}|\theta^{(?)}, \psi_{1}^{(?)}, \dots, \psi_{j-1}^{(?)}, \psi_{j+1}^{(?)}, \dots, \psi_{J}^{(?)}, y_{1:n})$$

$$\psi_{J}^{(?)} \sim p(\psi_{J}|\theta^{(?)}, \psi_{1}^{(?)}, \dots, \psi_{J}^{(?)}, y_{1:n})$$

## Convergence to stationary distribution

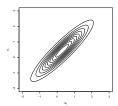
The samples  $(\theta^{(i)}, \psi^{(i)})$  converge to samples from  $p(\theta_{0:n}, \psi|y_{1:n})$ , regardless of what  $(\theta^{(0)}, \psi^{(0)})$  was.

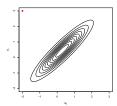
Let's look at an example: local level model.

$$\begin{array}{lll} Y_t &= \theta_t + v_t & v_t & \stackrel{ind}{\sim} N(0,2) \\ \theta_t &= \theta_{t-1} + w_t & w_t & \stackrel{ind}{\sim} N(0,0.5) \\ p(\theta_0) &= N(0,1) & \end{array}$$

with  $y_1 = 1$ . The objective is  $p(\theta_0, \theta_1|y_1)$ .

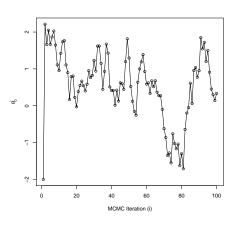
# Local level convergence example

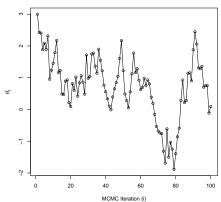






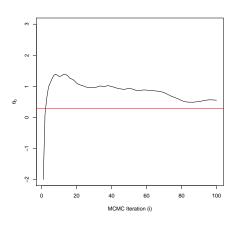
# **Traceplots**

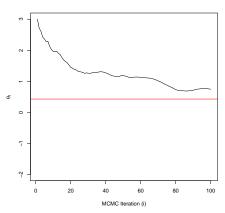




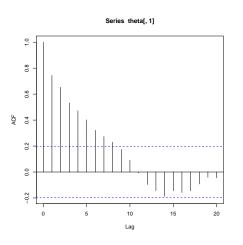
DLMs

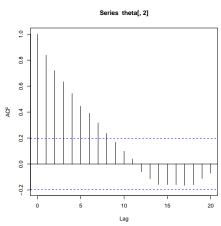
# Running average





# Auto-correlation plots





# MCMC Convergence diagnostics

- Graphical techniques
  - Traceplots
  - Ergodic mean
- Non-graphical techniques
  - Geweke diagnostic single chain
  - Gelman/Rubin diagnostic multiple chains

## Lack of convergence

We can never know if our chain has converged.

All convergence diagnostics detect a lack of convergence.

So instead of saying 'the chain has converged' you should be saying 'the chain shows no lack of convergence'.

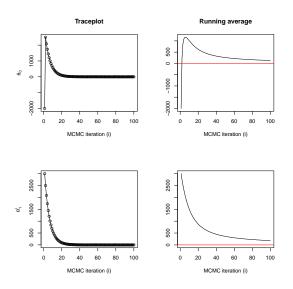
### Burn-in

#### Definition

Burn-in is the number of MCMC iterations before the chain shows no lack of convergence.

Burn-in is thrown-out to eliminate the bias associated with the starting point.

# Burn-in example



#### Burn-in

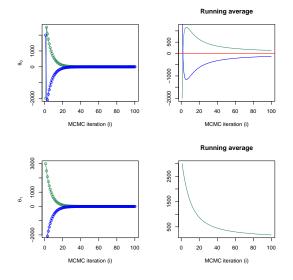
#### Definition

Burn-in is the period of time before the chain has converged.

Burn-in is thrown-out to eliminate the bias associated with the starting point.

If the starting point is crucial, why not start multiple chains in different locations? With the local level model, start chain 1 at (-2000, 3000) and start chain 2 at (3000, -2000).

# Multiple chains

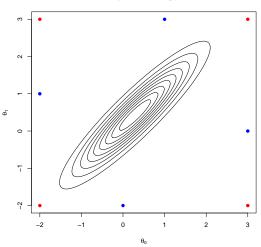


# Gelman-Rubin diagnostic

- Start multiple chains at locations that are overdispersed relative to the posterior.
- ANOVA comparison
  - Within-chain versus between-chain variances
  - Represented as a scale reduction factor such that values around 1 indicate no lack of convergence.

# Local level model example

#### Overdispersed starting points



## Local level model example - 100 iterations

In package coda, use function gelman.diag.

Potential scale reduction factors:

Point est. 97.5% quantile

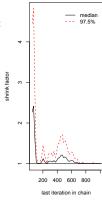
[1,] 1.12 1.45

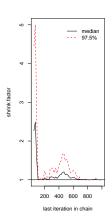
[2,] 1.13 1.48

Multivariate psrf

1.09

Values substantially above 1 indicate lack of convergence.





## Local level model example - 1000 iterations

In package coda, use function gelman.diag.

Potential scale reduction factors:

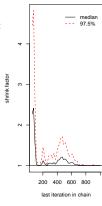
Point est. 97.5% quantile
[1,] 1 1.00
[2.] 1 1.00

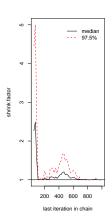
[2,] 1

Multivariate psrf

1.00

Values substantially above 1 indicate lack of convergence.





### Iterations for inference

Now that no lack of convergence is apparent, how long should I run my chain?

- The longer you run the chain, the lower your Monte Carlo error.
- Monte Carlo error reduces by the  $\sqrt{N}$  where N is the number of MCMC iterations.
- $\bullet$  So, if you want a  $10\mbox{-fold}$  decrease in Monte Carlo error, you need to run  $10^2$  times your current number of iterations.

# Simple Monte Carlo example

Consider the model  $y_i \stackrel{ind}{\sim} N(\mu,1)$  and our goal is to estimate  $E[y_i] = \mu$ . The Monte Carlo approximation is

$$\mu \approx \mu_{MC} = \frac{1}{n} \sum_{i=1}^{n} y_i$$

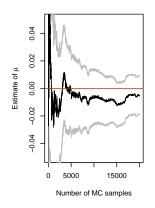
with the variance of this approximation given by

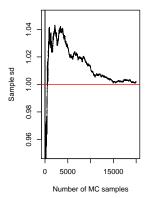
$$se(\mu_{MC}) \approx \sqrt{\frac{1}{n(n-1)} \sum_{i=1}^{n} (y_i - \mu_{MC})^2} = \frac{1}{\sqrt{n}} s d_y$$

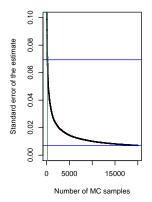
where  $sd_y$  is the standard deviation of the sample  $y=(y_1,y_2,\ldots,y_n)$ . Since this standard deviation converges to 1, by our model assumption above, the standard error of the Monte Carlo estimate decrease by the square root of n.

# Simple Monte Carlo example

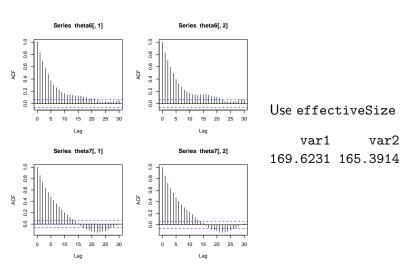
At 176 simulations, the standard error of  $\mu_{MC}$  is  $\sim 0.07$ . To decrease this to 0.007 (an order of magnitude increase in accuracy), we would need to take a total of  $176 \cdot 10^2 = 17600$  simulations.







### Iterations for inference

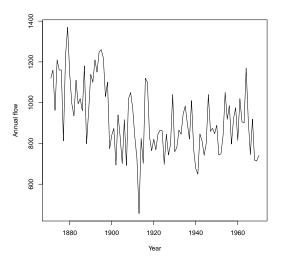


#### Work flow

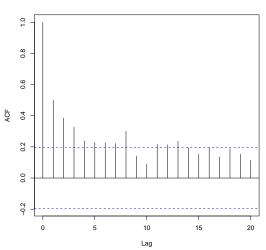
- Exploratory data analysis
- Define a model with priors
- Fit the model using MLE techniques
- Inference
  - Fit in WinBUGS
  - Code it up in R/C
    - Choose an MCMC scheme
    - Find the full conditional distributions (if available)
  - Monitor chain convergence
  - Summarize the posterior
- Model checking
  - Diagnostic plots to evaluate model assumptions, e.g. one-step head forecasts

### Examples

- Nile flow local level model
- Spain/Denmark investments SUTSE







 $v_t \sim N(0, V)$ 

### Local level model

$$\theta_t = \theta_{t-1} + w_t \qquad w_t \sim N(0, W)$$

$$V \sim IG(a_V, b_V)$$

$$W \sim IG(a_W, b_W)$$

$$\theta_0 \sim N(m_0, C_0)$$

 $Y_t = \theta_t + v_t$ 

where  $p(V, W, \theta_0) = p(V)p(W)p(\theta_0)$  and  $v_t$  and  $w_t$  are independent across time and mutually independent of each other as well as independent of  $\theta_0$ .

# Non-informative priors

$$V \sim IG(a_V, b_V)$$

$$W \sim IG(a_W, b_W)$$

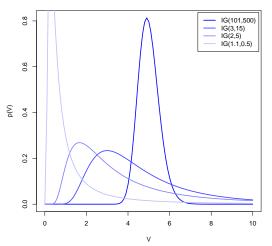
$$\theta_0 \sim N(m_0, C_0)$$

#### Non-informative prior

$$V \propto 1/V \implies a_V = b_V = 0$$
  
 $W \propto 1/W \implies a_W = b_W = 0$   
 $\theta_0 \propto 1 \implies m_0 = 0, C_0 = \infty$ 

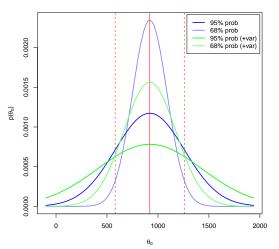
## Informative variance priors

Informative prior for V (or W), E[V]=5 with varying accuracy



### Informative state prior

"On average, the Nile flow is around 920  $\pm$  340." With what probability?



### MCMC scheme

In DLMs, conditional on unknown parameters, we can sample from the joint state vector at all times using FFBS.

- $p(\theta_{0:T}|V,W,y_{1:T})$  (for references see page 161 of Petris et al.)
- $p(V|\theta_{0:T}, W, y_{1:T})$
- $p(W|\theta_{0:T}, V, y_{1:T})$

### Full conditional distributions - the hard way

$$p(V|\theta_{0:T}, W, y_{1:T}) \propto \prod_{t=1}^{T} p(y_t|\theta_t, V) p(\theta_t|\theta_{t-1}, W) p(V) p(W) p(\theta_0) \propto \prod_{t=1}^{T} p(y_t|\theta_t, V) p(V) = \prod_{t=1}^{T} N(y_t; \theta_t, V) IG(V; a_V, b_V) \propto V^{-T/2} \exp\left(-\frac{1}{2V} \sum_{t=1}^{T} (y_t - \theta_t)^2\right) V^{-a_V - 1} \exp\left(-b_V/V\right) = V^{-(a_V + T/2) - 1} \exp\left(-\left[b_V + \frac{1}{2} \sum_{t=1}^{T} (y_t - \theta_t)^2\right]/V\right) \propto IG\left(a_V + T/2, b_V + \frac{1}{2} \sum_{t=1}^{T} (y_t - \theta_t)^2\right)$$

# Full conditional distributions - the hard way

$$p(W|\theta_{0:T}, V, y_{1:T}) \propto \prod_{t=1}^{T} p(y_t|\theta_t, V) p(\theta_t|\theta_{t-1}, W) p(V) p(W) p(\theta_0) \propto \prod_{t=1}^{T} p(\theta_t|\theta_{t-1}, W) p(W) = \prod_{t=1}^{T} N(\theta_t; \theta_{t-1}, W) IG(W; a_W, b_W) \propto W^{-T/2} \exp\left(-\frac{1}{2W} \sum_{t=1}^{T} (\theta_t - \theta_{t-1})^2\right) W^{-a_W-1} \exp\left(-b_W/W\right) = V^{-(a_W+T/2)-1} \exp\left(-\left[b_W + \frac{1}{2} \sum_{t=1}^{T} (\theta_t - \theta_{t-1})^2\right]/W\right) \propto IG\left(a_W + T/2, b_W + \frac{1}{2} \sum_{t=1}^{T} (\theta_t - \theta_{t-1})^2\right)$$

# Full conditional distributions - the easy way

Recall from HW 1b, if  $\sigma^2 \sim IG(a,b)$  and  $x_i \stackrel{ind}{\sim} N(0,\sigma^2)$ , then

$$p(\sigma^2|x_1, x_2, \dots, x_n) = IG\left(a + n/2, b + \frac{1}{2}\sum_{t=1}^n x_t^2\right).$$

Notice

$$V \sim IG(a_{V}, b_{V})$$

$$v_{t} = y_{t} - \theta_{t} \stackrel{ind}{\sim} N(0, V)$$

$$p(V|y_{1:T}, \theta_{0:T}, W) = p(V|y_{1:T}, \theta_{1:T})$$

$$= IG(a_{V} + T/2, b_{V} + \frac{1}{2} \sum_{t=1}^{T} v_{t}^{2})$$

$$W \sim IG(a_{W}, b_{W})$$

$$w_{t} = \theta_{t} - \theta_{t-1} \stackrel{ind}{\sim} N(0, W)$$

$$p(W|y_{1:T}, \theta_{0:T}, V) = p(W|\theta_{1:T})$$

$$= IG(a_{W} + T/2, b_{W} + \frac{1}{2} \sum_{t=1}^{T} w_{t}^{2})$$

### MCMC scheme revisited

- $p(\theta_{0:T}|\ldots)$  using FFBS
- $p(V|\ldots) = IG(a_V + T/2, b_V + \frac{1}{2} \sum_{t=1}^{T} v_t^2)$
- $p(W|\ldots) = IG(a_W + T/2, b_W + \frac{1}{2} \sum_{t=1}^{T} w_t^2)$

Notice that  $p(V|\dots)$  doesn't depend on W and  $p(W|\dots)$  doesn't depend on V. So our scheme is actually

- $p(\theta_{0:T}|\ldots)$  using FFBS
- $p(V, W | \dots) = IG(a_V + T/2, b_V + \frac{1}{2} \sum_{t=1}^T v_t^2) IG(a_W + T/2, b_W + \frac{1}{2} \sum_{t=1}^T w_t^2)$

### Coding it up

Begin by creating a function to draw from the posterior of a conjugate inverse gamma

```
drawIGpost <- function(x, a=0, b=0) {
  return(rinvgamma(1, a+length(x)/2, b+sum(x^2)/2))
}</pre>
```

### Coding it up

Begin by creating a function to draw from the posterior of a conjugate inverse gamma

```
for (i in 1:n.reps) {
  cat(i."\n")
  # Sample states
  mod <- dlmModPoly(1, dV=V, dW=W)</pre>
  filt <- dlmFilter(Nile, mod)</pre>
  theta <- dlmBSample(filt)
  # Sample V and W
  V <- drawIGpost(y-theta[-1])</pre>
  W <- drawIGpost(theta[-1]-theta[-n])</pre>
  # Save iterations
  V.reps[i] <- V</pre>
  W.reps[i] <- W
  theta.reps[i,] = theta
```

# Running the MCMC

#### Run 1

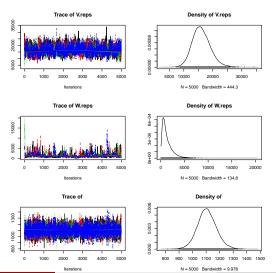
- Run 1 chain starting from the MLEs
- Check traceplots for this run
- Obtain posterior summaries for model parameters
- Choose initial values that are < minimum and > maximum for each model parameter

#### Multi-runs

- Start multiple chains from combinations of these values
- Check traceplots and Gelman-Rubin diagnostic for these chains
- Discard burn-in and produce posterior summaries on remaining iterations
- If more iterations are needed, initialize new chains from the last iteration of the old chains

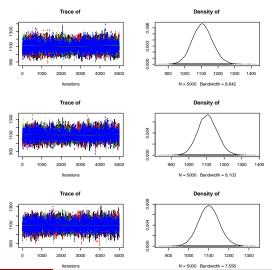
# Monitoring convergence

Use plot.mcmc in coda package.



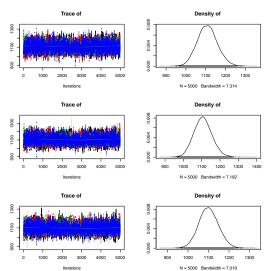
# Monitoring convergence

Use plot.mcmc in coda package.



### Monitoring convergence

Use plot.mcmc in coda package.



# Gelman-Rubin diagnostic

```
> gelman.diag(window(mcmc.results,1,4000))
Potential scale reduction factors:
```

```
Point est. 97.5% quantile
V.reps 1.00 1.00
W.reps 1.00 1.00
1.00 1.00
1.00 1.00
1.00 1.00
```

Multivariate psrf

1.01

> summary(window(mcmc.results,4001,5000))

```
Iterations = 4001:5000
Thinning interval = 1
Number of chains = 4
Sample size per chain = 1000
```

 Empirical mean and standard deviation for each variable, plus standard error of the mean:

```
Mean SD Naive SE Time-series SE V.reps 15642.8 3187.29 50.3955 125.8969 W.reps 1630.4 1449.03 22.9111 100.2639 1107.9 73.84 1.1676 1.3148 1108.7 62.93 0.9951 1.1196
```

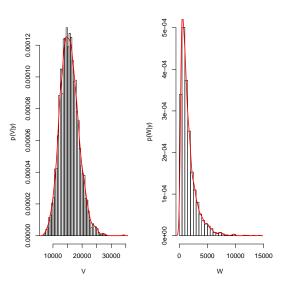
2. Quantiles for each variable:

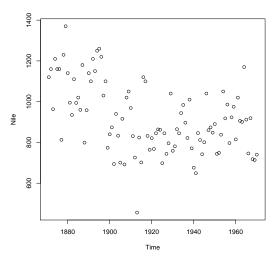
```
V.reps 9854.5 13459.2 15450.6 17640.9 22337.6

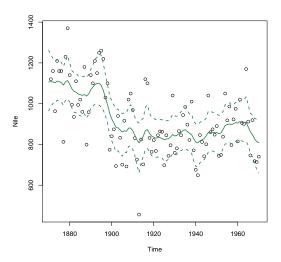
W.reps 241.2 651.7 1183.1 2092.0 5529.9

964.9 1060.0 1106.7 1153.6 1261.0

987.1 1066.5 1107.5 1149.0 1238.8
```







# Summaries of functions of parameters

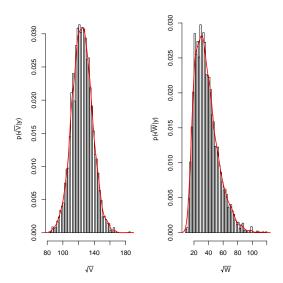
The posterior for  $f(\psi)$  is available using the MCMC simulations by plugging our iterations  $\psi^{(i)}$  into  $f(\cdot)$  and calculating desired quantities, e.g.

$$E[f(\psi)] \approx \frac{1}{n} \sum_{i=1}^{n} f(\psi^{(i)}).$$

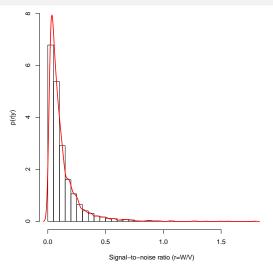
For example,

- $f(\theta_{0:T}, V, W) = \sqrt{(V)}$
- $f(\theta_{0:T}, V, W) = \sqrt{(W)}$
- $f(\theta_{0:T}, V, W) = W/V$  (signal-to-noise ratio)
- $f(\theta_{0:T}, V, W) = P(W/V < 1)$

### Standard deviations



# Signal-to-noise ratio



$$P(r < 1) = 0.998.$$

#### The model

$$Y_t = F_t \theta_t + v_t \qquad v_t \sim N(0, V)$$

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

where V is scalar and W is diagonal with elements  $W_i$  and assumed priors

$$p(V, W_1, \dots, W_p, \theta_0) = p(V)p(\theta_0) \prod_{i=1}^p p(W_i)$$

$$V \sim IG(a_V, b_V)$$

$$W_i \sim IG(a_{W_i}, b_{W_i})$$

$$\theta_0 \sim N(m_0, C_0)$$

#### Linear trend model

For example

$$Y_t = F\theta_t + v_t \qquad v_t \sim N(0, V)$$
  
$$\theta_t = G\theta_{t-1} + w_t \qquad w_t \sim N_p(0, W)$$

where F = (1,0), G[1,1] = G[1,2] = G[2,2] = 1, G[2,1] = 0, and W is diagonal with elements  $W_i$  and assumed priors

$$\begin{array}{rcl} p(V,W_1,\ldots,W_p,\theta_0) & = p(V)p(\theta_0)p(W_1)p(W_2) \\ V & \sim IG(a_V,b_V) \\ W_1 & \sim IG(a_{W_1},b_{W_1}) \\ W_2 & \sim IG(a_{W_2},b_{W_2}) \\ \theta_0 & \sim N(m_0,C_0) \end{array}$$

### Rewrite the linear trend model

$$Y_t = \mu_t + v_t \qquad v_t \sim N(0, \sigma^2)$$
  

$$\mu_t = \mu_{t-1} + \beta_t + w_{t,1} \qquad w_{t,1} \sim N(0, \sigma_{\mu}^2)$$
  

$$\beta_t = \beta_{t-1} + w_{t,2} \qquad w_{t,2} \sim N(0, \sigma_{\beta}^2)$$

where  $w_{t,1}$  and  $w_{t,2}$  are independent.

What are the full conditionals for  $\sigma^2, \sigma_\mu^2$ , and  $\sigma_\beta^2$ ?

• 
$$p(\sigma^2|\ldots) = IG\left(a_{\sigma^2} + T/2, b_{\sigma^2} + \frac{1}{2}\sum_{t=1}^T v_t^2\right)$$

• 
$$p(\sigma_{\mu}^2|\ldots) = IG\left(a_{\sigma_{\mu}^2} + T/2, b_{\sigma_{\mu}^2} + \frac{1}{2}\sum_{t=1}^T w_{t,1}^2\right)$$

• 
$$p(\sigma_{\beta}^2|\ldots) = IG\left(a_{\sigma_{\beta}^2} + T/2, b_{\sigma_{\beta}^2} + \frac{1}{2}\sum_{t=1}^T w_{t,2}^2\right)$$

and importantly, they are independent!

### More generally

$$Y_t = F_t \theta_t + v_t \qquad v_t \sim N(0, V)$$
  
$$\theta_t = G_t \theta_{t-1} + w_t \qquad w_t \sim N_p(0, W)$$

where W is diagonal with elements  $W_i$  and all variances have independent inverse gamma priors.

The full conditionals for parameters are

• 
$$p(V|...) = IG\left(a_V + T/2, b_V + \frac{1}{2}\sum_{t=1}^{T} v_t^2\right)$$

• 
$$p(W_i|\ldots) = IG\left(a_{W_i} + T/2, b_{W_i} + \frac{1}{2}\sum_{t=1}^T w_{t,i}^2\right)$$

and again, they are independent!

# MCMC scheme for models with d inverse gamma priors

#### Two-stage Gibbs sampler

- Use FFBS to sample from  $p(\theta_{0:T}|\ldots)$
- ullet Jointly sample  $V,W_1,\ldots,W_p$  by sampling their full conditionals
  - $p(V|\ldots)$
  - $p(W_i|\ldots)$  for  $i \in (1, 2, \ldots, p)$ .

Implemented in dlmGibbsDIG.

Suppose we assume the model

$$Y_t = F_t \theta_t + v_t \qquad v_t \sim N_m \left( 0, \Phi_0^{-1} \right)$$
  
$$\theta_t = G_t \theta_{t-1} + w_t \qquad w_t \sim N_p \left( 0, \Phi_1^{-1} \right)$$

where  $\Phi_0$  is an  $m \times m$  observation precision matrix and  $\Phi_1$  is a  $p \times p$  evolution precision matrix. It will be convenient to choose independent Wishart distributions for the prior for these precision matrices, i.e.

$$p(\Phi_0, \Phi_1) = p(\Phi_0) p(\Phi_1) = W(\Phi_0; \nu_0, S_0) W(\Phi_1; \nu_1, S_1)$$

where

$$\mathcal{W}(P;\nu,S) = \frac{|S|^{\nu}|P|^{\frac{\nu-p-1}{2}}}{\Gamma_{\nu}(\nu)} \exp\left(-tr(SP)\right)$$

is a distribution on symmetric, positive definite matrices P with parameters  $\nu>(p-1)/2$  and S, symmetric non-singular matrix.

### Full conditionals for the precision matrices

Wishart distributions are conditionally conjugate in this model:

$$p(\Phi_0|...) = \mathcal{W}\left(\nu_0 + T/2, S_0 + \frac{1}{2}SS_y\right)$$

where  $SS_y = \sum_{t=1}^{T} (y_t - F_t \theta_t) (y_t - F_t \theta_t)^{\top}$ .

$$p\left(\Phi_{1}|\ldots\right) = \mathcal{W}\left(\nu_{1} + T/2, S_{1} + \frac{1}{2}SS_{\theta}\right)$$

where  $SS_{\theta} = \sum_{t=1}^{T} (\theta_t - G_t \theta_{t-1}) (\theta_t - G_t \theta_{t-1})^{\top}$ .

Again, 
$$p(\Phi_0, \Phi_1 | \ldots) = p(\Phi_0 | \ldots) p(\Phi_1 | \ldots)$$
.

To draw from these distributions, use rwishart in dlm package which has arguments degrees of freedom  $\delta$  and scale matrix  $V_0^{-1}$  where  $\mathcal{W}(\delta/2,V_0/2)$ .

#### The model

Consider the model with block-diagonal evolution covariance:

$$Y_t = F_t \theta_t + v_t$$
  $v_t \sim N_m(0, \Phi_0^{-1})$   
 $\theta_t = G_t \theta_{t-1} + w_t$   $w_t \sim N_{p*}(0, W)$ 

where W is block-diagonal with elements  $W_i$ . Set  $\Phi_i^{-1} = W_i$  and give  $\Phi_0, \Phi_1, \dots, \Phi_d$  independent Wishart priors  $\Phi_i \sim \mathcal{W}(\nu_i, S_i)$ .

#### Rewritten univariate model

For combining individual components, e.g. polynomial trend, seasonal, dynamic regression,  $G_t$  is block diagonal with elements  $G_{i,t}$  relating to  $W_i$  and the model can be re-written

$$\begin{array}{lll} Y_t &= F_t \theta_t + v_t & v_t \sim N(0, \Phi_0^{-1}) \\ \theta_{1,t} &= G_{1,t} \theta_{1,t-1} + w_{1,t} & w_{1,t} \sim N_{p_1}(0, \Phi_1^{-1}) \\ &\vdots & \\ \theta_{i,t} &= G_{i,t} \theta_{i,t-1} + w_{i,t} & w_{i,t} \sim N_{p_i}(0, \Phi_i^{-1}) \\ &\vdots & \\ \theta_{p,t} &= G_{p,t} \theta_{p,t-1} + w_{p,t} & w_{p,t} \sim N_{p_d}(0, \Phi_d^{-1}) \end{array}$$

where  $w_{i,t}$  are independent across i. Then

$$SS_{ii,t} = (\theta_{i,t} - G_{i,t}\theta_{i,t-1})(\theta_{i,t} - G_{i,t}\theta_{i,t-1})^{\top}.$$

### Multivariate models

Let

$$SS_t = (\theta_t - G_t \theta_{t-1})(\theta_t - G_t \theta_{t-1})^{\top}$$

and partition it according to

$$SS_t = \begin{bmatrix} SS_{11,t} & \cdots & SS_{1d,t} \\ \vdots & \ddots & \vdots \\ SS_{d1,t} & \cdots & S_{dd,t} \end{bmatrix}$$

where the partition is according to the partition in  $\Phi = \mathsf{blockdiag}(\Phi_1, \dots, \Phi_d)$ .

### Full conditional distributions

$$p\left(\Phi_0^{-1}|\ldots\right) = \mathcal{W}\left(\nu_0 + T/2, S_0 + \frac{1}{2}SS_y\right)$$

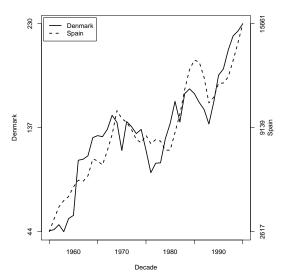
where  $SS_y = \sum_{t=1}^{T} (y_t + F_t \theta_t) (y_t - F_t \theta_t)^{\top}$ .

$$p\left(\Phi_i^{-1}|\ldots\right) = \mathcal{W}\left(\nu_i + T/2, S_i + \frac{1}{2}SS_{\theta_i}\right)$$

where  $SS_{\theta_i} = \sum_{t=1}^{T} SS_{ii,t}$  given on the previous page.

Once again, 
$$p\left(\Phi_0^{-1},\Phi_1^{-1},\ldots,\Phi_d^{-1}|\ldots\right)=p\left(\Phi_0^{-1}|\ldots\right)\prod_{i=1}^d p\left(\Phi_i^{-1}|\ldots\right)$$
.

# Denmark and Spain investments



### SUTSE model

$$Y_t = (F \otimes I_2)\theta_t + v_t$$
  $v_t \sim N_2(0, \Phi_0^{-1})$   
 $\theta_t = (G \otimes I_2)\theta_{t-1} + w_t$   $w_t \sim N_4(0, W)$ 

where  $W=\operatorname{blockdiag}(W_1,W_2)$ ,  $\Phi_1^{-1}=W_1$ , and  $\Phi_2^{-1}=W_2$ .

Assume independent Wishart priors

$$p(\Phi_0) = \mathcal{W}\left(\frac{\delta_0 + 1}{2}, \frac{1}{2}V_0\right) \qquad V_0 = (\delta_0 - 2) \begin{bmatrix} 10^2 & 0 \\ 0 & 500^2 \end{bmatrix}$$

$$p(\Phi_1) = \mathcal{W}\left(\frac{\delta_1 + 1}{2}, \frac{1}{2}W_{\mu,0}\right) \qquad W_{\mu,0} = (\delta_1 - 2) \begin{bmatrix} 0.01^2 & 0 \\ 0 & 0.01^2 \end{bmatrix}$$

$$p(\Phi_2) = \mathcal{W}\left(\frac{\delta_2 + 1}{2}, \frac{1}{2}W_{\beta,0}\right) \qquad W_{\beta,0} = (\delta_2 - 2) \begin{bmatrix} 5^2 & 0 \\ 0 & 100^2 \end{bmatrix}$$

where  $\delta_0 = \delta_2 = 3$  and  $\delta_1 = 100$ .

# MCMC sampling

#### MCMC Scheme:

- Sample  $\theta_{0:T} \sim p(\theta_{0:T} | \dots)$  using FFBS
- Sample  $p(\Phi_0, \Phi_1, \Phi_2 | \dots)$  jointly

$$p(\Phi_{0}|\ldots) = \mathcal{W}\left(\frac{\delta_{0}+1+T}{2}, \frac{1}{2}(V_{0}+SS_{y})\right)$$

$$p(\Phi_{1}|\ldots) = \mathcal{W}\left(\frac{\delta_{1}+1+T}{2}, \frac{1}{2}(W_{\mu,0}+SS_{1})\right)$$

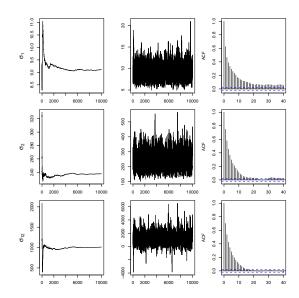
$$p(\Phi_{2}|\ldots) = \mathcal{W}\left(\frac{\delta_{2}+1+T}{2}, \frac{1}{2}(W_{\beta,0}+SS_{2})\right)$$

where

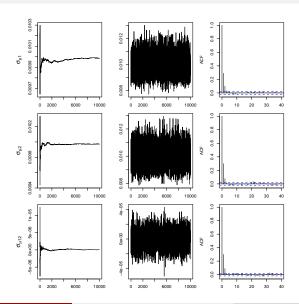
$$SS_{i.} = \sum_{t=1}^{T} SS_{ii,t}.$$

provided earlier.

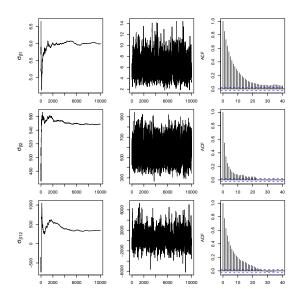
# Convergence and autocorrelation



### Convergence



#### Convergence



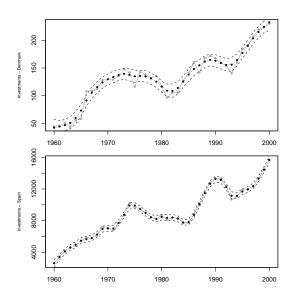
#### Posterior covariance expectations

$$E[V|y_{1:T}] = \begin{bmatrix} 86 & (1) & 1026 & (23) \\ & 59340 & (807) \end{bmatrix}$$

$$E[W_{\mu}|y_{1:T}] = 1e - 5 \begin{bmatrix} 9.97 & (0.02) & 0.016 & (0.014) \\ & 10.04 & (0.02) \end{bmatrix}$$

$$E[W_{\beta}|y_{1:T}] = \begin{bmatrix} 38.3 & (0.8) & 305 & (41) \\ & 311073 & (2346) \end{bmatrix}$$

# Posterior $\mu_t$



## Types of missing data

Complete data  $Y_{i,t}$  and missing indicator  $M_{i,t}$  where  $M_{i,t}=1$  if observation  $Y_{it}$  is missing and 0 otherwise. Let  $Y_{\mbox{obs}}$  contain all the data that is observed while  $Y_{\mbox{mis}}$  contains all the data that is missing with  $Y=(Y_{\mbox{obs}},Y_{\mbox{mis}}).$  Then several types of missing-ness are possible:

Missing completely at random (MCAR)

$$p(M|Y,\phi) = p(M|\phi).$$

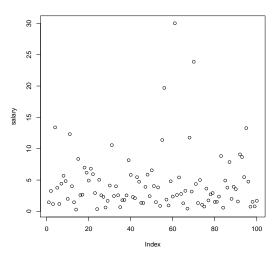
Missing at random (MAR)

$$p(M|Y,\phi) = p(M|Y_{obs},\phi).$$

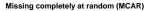
Not missing at random

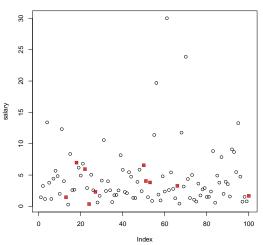
$$p(M|Y,\phi)$$
 depends on  $Y_{\text{mis}}$ .

## No missing data



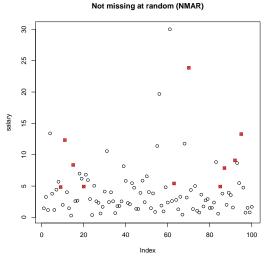
## Missing completely at random





# Not missing at random





#### Missing data in multivariate DLMs

#### Two situations:

- ullet Totally missing: at time t,  $Y_t$  is completely missing
- ullet Partially missing: at time t, part of  $Y_t$  is observed

#### Totally missing

Recall 'Kalman filter' lecture: missing data are handled trivially while filtering.

$$m_t = a_t$$
  $C_t = R_t$ .

Unknown fixed parameters are sampled without these data, e.g. scalar  ${\it V}$ 

$$p(\phi_V|\ldots) = G\left(a + \frac{T'}{2}, b + \sum_{t \in \mathsf{obs}} (y_t - F_t \theta_t)^2\right)$$

where 'obs' is a vector of times when the data are observed and  $T^\prime \leq T$  is the length of obs.

e.g. matrix V

$$p(\Phi_V | \dots) = W\left(a + \frac{T'}{2}, b + \sum_{t \in \mathsf{obs}} (y_t - F_t \theta_t) (y_t - F_t \theta_t)^\top\right).$$

## Partially missing when filtering

Suppose  $M_t$  is the matrix that is built by taking an identity matrix and removing the rows of any missing observations in  $y_t$ . Then  $\tilde{y}_t = M_t y_t$  contains only the observed data. The correct observation equation to consider is

$$\tilde{y}_t = \tilde{F}_t \theta_t + \tilde{v}_t \qquad \tilde{v}_t \sim N(0, \tilde{V}_t).$$

What are  $\tilde{F}_t$  and  $\tilde{V}_t$ ?

- $\tilde{F}_t = M_t F_t$
- $\bullet \ \tilde{V}_t = M_t V_t M_t^{\top}$

#### Partially missing in MCMC

Let  $Y=(Y_{\mbox{obs}},Y_{\mbox{mis}}).$  If we build an MCMC with only the observed data, then our scheme will look like

- Sample  $p(\theta|Y_{\text{obs}}, \psi)$  via FFBS
- Sample  $p(\psi|Y_{\text{obs}},\theta)$ .

For example, consider the observation precision matrix  $\Phi_V$  as the only unknown parameter. What is it's full conditional distribution?

$$p(\Phi_V|Y_{\mathsf{obs}},\theta) \propto p(Y_{\mathsf{obs}}|\Phi_V,\theta)p(\Phi_V)$$

Who knows?

#### Partially missing in MCMC

Let  $Y=(Y_{\mbox{obs}},Y_{\mbox{mis}}).$  Augment the MCMC to simulate the missing values, then our scheme will look like

- Sample  $p(\theta|Y,\psi)$  via FFBS
- Sample  $p(\psi|Y,\theta)$
- Sample  $p(Y_{\text{mis}}|Y_{\text{obs}}, \theta, \psi)$ .

This works since

$$p(\theta, \psi|Y_{\text{obs}}) = \int p(\theta, \psi, Y_{\text{mis}}|Y_{\text{obs}})dY_{\text{mis}}.$$

#### Partially missing in MCMC

How to simulate  $p(Y_{\text{mis}}|Y_{\text{obs}}, \theta, \psi)$ ?

First note,

$$p(Y_{\mathsf{mis}}|Y_{\mathsf{obs}},\theta,\psi) = \prod_{t=1}^{T} p(Y_{\mathsf{mis},t}|Y_{\mathsf{obs},t},\theta_t,\psi).$$

Second note.

$$\left(\begin{array}{c} Y_{\mathsf{mis},t} \\ Y_{\mathsf{obs},t} \end{array}\right) \sim N \left( \left[\begin{array}{c} F\theta_{\mathsf{mis},t} \\ F\theta_{\mathsf{obs},t} \end{array}\right], \left[\begin{array}{cc} V_{\mathsf{mis}} & V_{\mathsf{m,o}} \\ V_{\mathsf{o,m}} & V_{\mathsf{obs}} \end{array}\right] \right).$$

#### Goal

With all fixed parameters known:

$$p(y_{t+k}, \theta_{t+k}|y_{1:t}) = \int p(y_{t+k}, \theta_{t+k}, \theta_{t+(k-1)}|y_{1:t}) d\theta_{t+(k-1)}$$
  
= 
$$\int p(y_{t+k}, \theta_{t+k}|\theta_{t+(k-1)}) p(\theta_{t+(k-1)}|y_{1:t}) d\theta_{t+(k-1)}$$

To get  $p(y_{t+k}, \theta_{t+k} | \theta_t)$  just use the Kalman filter with missing data from  $y_{t+1}$  up to  $y_{t+(k-1)}$ .

With unknown fixed parameters:

$$p(y_{t+k}, \theta_{t+k}|y_{1:t}) =$$

$$= \int p(y_{t+k}, \theta_{t+k}|\theta_{t+(k-1)}, \psi) p(\theta_{t+(k-1)}, \psi|y_{1:t}) d\theta_{t+(k-1)} d\psi.$$

Now we can't just use the Kalman filter due to the unknown fixed parameters. Instead, we need to integrate over their posteriors.

## MCMC Forecasting

After completing the MCMC, follow this procedure

- For each iteration  $j=1,2,\ldots,J$  in the MCMC chain post burn-in:
  - Run a Kalman filter (dlmFilter) on your data using  $\psi^{(j)}$  to obtain  $p(\theta_t|y_{1:t},\psi^{(j)})=N(m_t^{(j)},C_t^{(j)}).$
  - Forecast ahead (dlmForecast) to obtain  $p(y_{t+k}|y_{1:t},\psi^{(j)}) = N(f_t(k)^{(j)},Q_t(k)^{(j)}) \text{ (see section 2.8 in Petris)}$
  - Calculate mean and 95% intervals for  $p(y_{t+k}|y_{1:t},\psi^{(j)})$ , i.e.  $f_t(k)^{(j)} \left(f_t(k)^{(j)}-1.96\sqrt{Q_t(k)^{(j)}},f_t(k)^{(j)}+1.96\sqrt{Q_t(k)^{(j)}}\right)$  if Q is scalar, otherwise do this component-wise.

This provides a set of means and 95% intervals, one for each MCMC iteration j.

#### MCMC Forecasting

To find the marginal mean and 95% interval, average these means and 95% intervals for all j, i.e.

$$E[y_{t+k}|y_{1:t}] \approx \frac{1}{J} \sum_{j=1}^{J} f_t(k)^{(j)}$$

$$Q_{2.5\%}[y_{t+k}|y_{1:t}] \approx \frac{1}{J} \sum_{j=1}^{J} f_t(k)^{(j)} - 1.96\sqrt{Q_t(k)^{(j)}}$$

$$Q_{97.5\%}[y_{t+k}|y_{1:t}] \approx \frac{1}{J} \sum_{j=1}^{J} f_t(k)^{(j)} + 1.96\sqrt{Q_t(k)^{(j)}}$$

If you have many MCMC iterations, you can use fewer iterations for this forecast by thinning the chain.