

Introduction to Bayesian Estimation

Wouter J. Den Haan
London School of Economics

© 2011 by Wouter J. Den Haan

June 20, 2011

Overview

- A very useful tool: Kalman filter
- Maximum Likelihood
 - Singularity when $\#\text{shocks} \leq \text{number of observables}$
- Bayesian estimation
- Tools:
 - Metropolis Hastings
- Remaining issues

Rudolph E. Kalman



born in Budapest, Hungary, on May 19, 1930

Kalman filter

- Linear projection
- Linear projection with orthogonal regressors
- Kalman filter

The slides for the Kalman filter is based on Ljungqvist and Sargent's textbook

Linear projection

- $y: n_y \times 1$ vector of random variables
- $x: n_x \times 1$ vector of random variables
- First and second moments exist

$$\begin{array}{lll} E y = \mu_y & \tilde{y} = y - \mu_y & E \tilde{x} \tilde{x}' = \Sigma_{xx} \\ E x = \mu_x & \tilde{x} = x - \mu_x & E \tilde{y} \tilde{y}' = \Sigma_{yy} \\ & & E \tilde{y} \tilde{x}' = \Sigma_{yx} \end{array}$$

Definition of linear projection

The *linear projection* of y on x is the function

$$\hat{E}[y|x] = a + Bx,$$

a and B are chosen to minimize

$$E \text{ trace } \{(y - a + Bx)(y - a + Bx)'\}$$

Formula for linear projection

The *linear projection* of y on x is given by

$$\hat{E}[y|x] = \mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(x - \mu_x)$$

Difference with linear regression problem

- True model:

$$y = \bar{B}x + \bar{D}z + \varepsilon,$$

$$\mathbb{E}x = \mathbb{E}z = \mathbb{E}\varepsilon = 0, \mathbb{E}[\varepsilon|x, z] = 0, \mathbb{E}[z|x] \neq 0$$

\bar{B} : measures the effect of x on y *keeping all else—also z and ε —constant.*

- Particular regression model:

$$y = \bar{B}x + u$$

Difference with linear regression problem

Comments:

- Least-squares estimate $\neq \bar{B}$
- Projection:

$$\hat{E}[y|x] = Bx = \bar{B}x + \bar{D}\hat{E}[y|x]$$

- Projection well defined
linear projection can include more than the direct effect:

Message:

- You can always define the linear projection
- you don't have to worry about the properties of the error term.

Linear Projection with orthogonal regressors

- $x = [x_1, x_2]$ and suppose that $\Sigma_{x_1 x_2} = 0$
- x_1 and x_2 could be vectors

$$\begin{aligned}
 \hat{E}[y|x] &= \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1} (x - \mu_x) \\
 &= \mu_y + [\Sigma_{yx_1} \Sigma_{yx_2}] \begin{bmatrix} \Sigma_{x_1 x_1}^{-1} & 0 \\ 0 & \Sigma_{x_2 x_2}^{-1} \end{bmatrix} (x - \mu_x) \\
 &= \mu_y + \Sigma_{yx_1} \Sigma_{x_1 x_1}^{-1} (x_1 - \mu_{x_1}) + \Sigma_{yx_2} \Sigma_{x_2 x_2}^{-1} (x_2 - \mu_{x_2})
 \end{aligned}$$

Thus

$$\hat{E}[y|x] = \hat{E}[y|x_1] + \hat{E}[y|x_2] - \mu_y \quad (1)$$

Time Series Model

$$x_{t+1} = Ax_t + Gw_{1,t+1}$$

$$y_t = Cx_t + w_{2,t}$$

$$Ew_{1,t} = Ew_{2,t} = 0$$

$$E \begin{bmatrix} w_{1,t+1} \\ w_{2,t} \end{bmatrix} \begin{bmatrix} w_{1,t+1} \\ w_{2,t} \end{bmatrix}' = \begin{bmatrix} V_1 & V_3 \\ V_3 & V_2 \end{bmatrix}$$

Time Series Model

- y_t is observed, but x_t is not
- the coefficients are known (could even be time-varying)
- Initial condition:
 - x_1 is a random variable (mean μ_{x_1} & covariance matrix Σ_1)
- $w_{1,t+1}$ and $w_{2,t}$ are serially uncorrelated and orthogonal to x_1

Objective

The objective is to calculate

$$\begin{aligned}\hat{E}_t x_{t+1} &\equiv \hat{E}[x_{t+1} | y_t, y_{t-1}, \dots, y_1, \hat{x}_1] \\ &= \hat{E}[x_{t+1} | Y^t, \hat{x}_1]\end{aligned}$$

where \hat{x}_1 is an initial estimate of x_1 (Typically μ_{x_1})

Trick: get a recursive formulation

Orthogonalization of the information set

- Let
 - $\hat{y}_t = y_t - \hat{E} [y_t | \hat{y}_{t-1}, \hat{y}_{t-2}, \dots, \hat{y}_1, \hat{x}_1]$
 - $\hat{Y}^t = \{\hat{y}_t, \hat{y}_{t-1}, \dots, \hat{y}_1\}$
- space spanned by $\{\hat{x}_1, \hat{Y}^t\} =$ space spanned by $\{\hat{x}_1, Y_t\}$
 - That is, anything that can be expressed as a linear combination with elements in $\{\hat{x}_1, \hat{Y}^t\}$ can be expressed as a linear combination of elements in $\{\hat{x}_1, Y_t\}$.

Orthogonalization of the information set

- Then

$$\widehat{\mathbb{E}} [y_{t+1} | Y^t, \hat{x}_1] = \widehat{\mathbb{E}} [y_{t+1} | \hat{Y}^t, \hat{x}_1] = C \widehat{\mathbb{E}} [x_{t+1} | \hat{Y}^t, \hat{x}_1] \quad (2)$$

Derivation of the Kalman filter

From (1) we get

$$\widehat{E} [x_{t+1} | \hat{Y}^t, \hat{x}_1] = \widehat{E} [x_{t+1} | \hat{y}_t] + \widehat{E} [x_{t+1} | \hat{Y}^{t-1}, \hat{x}_1] - E x_{t+1} \quad (3)$$

The first term in (3) is a standard linear projection:

$$\begin{aligned}\widehat{E} [x_{t+1} | \hat{y}_t] &= E x_{t+1} + \text{cov}(x_{t+1}, \hat{y}_t) [\text{cov}(\hat{y}_t, \hat{y}_t)]^{-1} (\hat{y}_t - E \hat{y}_t) \\ &= E x_{t+1} + \text{cov}(x_{t+1}, \hat{y}_t) [\text{cov}(\hat{y}_t, \hat{y}_t)]^{-1} \hat{y}_t\end{aligned}$$

Some algebra

- Similar to the definition of \hat{y}_t , let

$$\begin{aligned}\hat{x}_{t+1} &= x_{t+1} - \hat{\mathbb{E}}[x_{t+1} | \hat{y}_t, \hat{y}_{t-1}, \dots, \hat{y}_1, \hat{x}_1] \\ &= x_{t+1} - \hat{\mathbb{E}}_t x_{t+1}\end{aligned}$$

- Let $\Sigma_{\hat{x}_t} = \mathbb{E} \hat{x}_t \hat{x}_t'$

$$\text{cov}(x_{t+1}, \hat{y}_t) = A \Sigma_{\hat{x}_t} C' + G V_3$$

$$\text{cov}(\hat{y}_t, \hat{y}_t) = C \Sigma_{\hat{x}_t} C' + V_2$$

Using the derived expressions

$$\widehat{\mathbb{E}} [x_{t+1} | \hat{y}_t]$$

$$= \mathbb{E}x_{t+1} + \text{cov}(x_{t+1}, \hat{y}_t) [\text{cov}(\hat{y}_t, \hat{y}_t)]^{-1} \hat{y}_t$$

$$= \mathbb{E}x_{t+1} + (A\Sigma_{\hat{x}_t}C' + GV_3) (C\Sigma_{\hat{x}_t}C' + V_2)^{-1} \hat{y}_t \quad (4)$$

Derivation Kalman filter

- Now get an expression for the second term in (3).
- From $x_{t+1} = Ax_t + Gw_{1,t+1}$, we get

$$\hat{E} \left[x_{t+1} | \hat{Y}^{t-1}, \hat{x}_1 \right] = A \hat{E} \left[x_t | \hat{Y}^{t-1}, \hat{x}_1 \right] = A \hat{E}_{t-1} x_t \quad (5)$$

Using (4) and (5) in (3) gives the *recursive* expression

$$\hat{E}_t x_{t+1} = A \hat{E}_{t-1} x_t + K_t \hat{y}_t$$

where

$$K_t = (A \Sigma_{\hat{x}_t} C' + G V_3) (C \Sigma_{\hat{x}_t} C' + V_2)^{-1}$$

Prediction for observable

From

$$y_{t+1} = Cx_{t+1} + w_{2,t+1}$$

we get

$$\hat{E}[y_{t+1}|Y_t, \hat{x}_1] = C\hat{E}_t x_{t+1}$$

Thus

$$\hat{y}_{t+1} = y_{t+1} - C\hat{E}_t x_{t+1}$$

Updating the covariance matrix

- We still need an equation to update $\Sigma_{\hat{x}_t}$. This is actually not that hard. The result is

$$\Sigma_{\hat{x}_{t+1}} = A\Sigma_{\hat{x}_t}A' + GV_1G' - K_t(A\Sigma_{\hat{x}_t}C' + GV_3)'$$

- Expression is deterministic and does not depend particular realizations. That is, precision only depends on the coefficients of the time series model

Applications Kalman filter

- signal extraction problems
 - GPS, computer vision applications, missiles
- prediction
- simple alternative to calculating inverse policy functions
 - (see below)

Estimating DSGE models

- Forget the Kalman filter for now, we will not use it for a while
- What is next?
 - Specify the neoclassical model that will be used as an example
 - Specify the linearized version
 - Specify the estimation problem
 - Maximum Likelihood estimation
 - Explain why Kalman filter is useful
 - Bayesian estimation
 - MCMC, a necessary tool to do Bayesian estimation

Neoclassical growth model

First-order conditions

$$c_t^{-\nu} = E_t \left[\beta c_{t+1}^{-\nu} (\alpha z_{t+1} k_t^{\alpha-1} + 1 - \delta) \right]$$

$$c_t + k_t = z_t k_{t-1}^\alpha + (1 - \delta) k_{t-1}$$

$$z_t = (1 - \rho) + \rho z_{t-1} + \varepsilon_t$$

$$\varepsilon_t \sim N(0, \sigma^2)$$

Linearized solution

$$k_t = \bar{k} + a_{k,k}(k_{t-1} - \bar{k}) + a_{k,z}(z_t - \bar{z})$$

$$z_t = (1 - \rho) + \rho z_{t-1} + \varepsilon_t$$

$$\varepsilon_t \sim N(0, \sigma^2)$$

$$z_0 \sim N(1, \sigma^2 / (1 - \rho^2))$$

k_0 is given

- $a_{k,k}$, $a_{k,z}$, and \bar{k} are *known* functions of the structural parameters
 \implies better notation would be $a_{k,k}(\Psi)$, $a_{k,z}(\Psi)$, and $\bar{k}(\Psi)$
- Consumption has been substituted out
- Approximation error is ignored. Linearized model is treated as the true model with Ψ as the parameters

Estimation problem

Given data for capital, $\{k_t\}_0^T$, estimate the set of coefficients, Ψ

$$\Psi = [\alpha, \beta, \nu, \delta, \rho, \sigma, z_0]$$

- No data on productivity, z_t .
 - If you had data on $z_t \implies \text{Likelihood} = 0$ for sure
 - More on this below.

Formulation of the Likelihood

- Let Y^T be the complete sample

$$L(Y^T | \Psi) = p(z_0) \prod_{t=1}^T p(z_t | z_{t-1})$$

$p(z_t | z_{t-1})$ corresponds with probability of a particular value for ε_t

Formulation of the Likelihood

Basic idea:

- Given a value for Ψ and give the data set, Y^T , you can calculate the implied values for ε_t
- We know the distribution of $\varepsilon_t \implies$
- We can calculate the probability (likelihood) of $\{\varepsilon_1, \dots, \varepsilon_T\}$

Formulation of the Likelihood

$$k_t = \bar{k} + a_{k,k}(k_{t-1} - \bar{k}) + a_{k,z}(z_t - \bar{z})$$

$$\implies$$

$$z_t = \frac{a_{k,z}\bar{z} - \bar{k} + a_{k,k}\bar{k}}{a_{k,z}} - \frac{a_{k,k}}{a_{k,z}}k_{t-1} + \frac{1}{a_{k,z}}k_t$$

$$z_t = b_0 + b_1 k_{t-1} + b_2 k_t$$

$$\varepsilon_t = z_t - (1 - \rho) - \rho z_{t-1}$$

Formulation of the Likelihood

- ε_t is obtained by **inverting** the policy function
- For larger systems, this inversion is not as easy to implement.
 - Below, we show an alternative

Formulation of the Likelihood

A bit more explicit

- Take a value for Ψ
- Given k_0 and k_1 you can calculate z_1
- Given z_0 you can calculate ε_1
- Continuing, you can calculate $\varepsilon_t \forall t$
- To make explicit the dependence of ε_t on Ψ , write $\varepsilon_t(\Psi)$
- The Likelihood can thus be written as

$$\prod_{t=1}^T \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ \frac{-(\varepsilon_t(\Psi))^2}{2\sigma^2} \right\}$$

Too few unobservables & singularities

- Above we assumed that there was no data on z_t
- Suppose you had data on z_t
- There are two cases to consider
 - Data not exactly generated by this model (most likely case)
 \implies Likelihood = 0 for any value of Ψ
 - Data is exactly generated by this model
 \implies Likelihood = 1 for true value of Ψ *and*
 \implies Likelihood = 0 for any other value for Ψ

Too few unobservables & singularities

$$k_t = \bar{k} + a_{k,k}(k_{t-1} - \bar{k}) + a_{k,z}(z_t - \bar{z})$$

Using the values for 4 periods, you can pin down \bar{k} , \bar{z} , $a_{k,k}$, and $a_{k,z}$.

- What about values for additional periods?
 - Data generated by model (unlikely of course)
⇒ additional observations will fit this equation too
 - Data not generated by model
⇒ additional observations will not fit this equation
⇒ Likelihood = zero

Too few unobservables & singularities

- Can't I simply add an error term?

$$k_t = \bar{k} + a_{k,k}(k_{t-1} - \bar{k}) + a_{k,z}(z_t - \bar{z}) + u_t$$

- Answer: **NO** not in general
- Why not? It is ok in standard regression

Too few unobservables & singularities

Why is the answer NO in general?

- ① u_t represents other shocks such as preference shocks
 \Rightarrow its presence is likely to affect \bar{k} , $a_{k,k}$, and $a_{k,z}$
- ② u_t represents measurement error
 \Rightarrow you are fine from an econometric stand point
 \Rightarrow but is residual only measurement error?

What if you also observe consumption?

Suppose you observe k_t , c_t , but not z_t ?

$$\begin{aligned} k_t &= \bar{k} + a_{k,k}(k_{t-1} - \bar{k}) + a_{k,z}(z_t - \bar{z}) \\ c_t &= \bar{c} + a_{c,k}(k_{t-1} - \bar{k}) + a_{c,z}(z_t - \bar{z}) \end{aligned}$$

- Recall that the coefficients are functions of Ψ
- Given value of Ψ you can solve for z_t from top equation
- Given value of Ψ you can solve for z_t from bottom equation
- With real world data you will get inconsistent answers.

Unobservables and avoiding singularities

General rule:

- For every observable you need at least one unobservable shock
- Letting them be measurement errors is hard to defend
- The last statement does not mean that you cannot *also* add measurement errors

Using the Kalman filter

$$x_{t+1} = Ax_t + Gw_{1,t+1} \quad (6)$$

$$y_t = Cx_t + w_{2,t} \quad (7)$$

- (6) describes the equations of the model;
 - x_t consists of the "true" values of state variables like capital and productivity.
- (7) relates the observables, y_t , to the "true" values

Example

- consumption and capital are observed with error
 - $c_t^* = c_t + u_{c,t}$
 - $k_t^* = k_t + u_{k,t}$
- z_t is unobservable
- $x_t' = [k_{t-1} - \bar{k}, z_{t-1} - \bar{z}]$
- $w_{1,t+1} = \varepsilon_t$
- $y_t' = [k_{t-1}^* - \bar{k}, c_t^* - \bar{c}]$

Example

- (6) gives policy function for k_t and law of motion for z_t

$$\begin{bmatrix} k_t - \bar{k} \\ c_t - \bar{c} \\ z_{t+1} - \bar{z} \end{bmatrix} = \begin{bmatrix} a_{k,k} & a_{k,z} \\ a_{c,k} & a_{c,z} \\ 0 & \rho \end{bmatrix} \begin{bmatrix} k_{t-1} - \bar{k} \\ z_t - \bar{z} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \varepsilon_{t+1} \end{bmatrix}$$

- Equation (7) is equal to

$$\begin{bmatrix} k_{t-1}^* - \bar{k} \\ c_t^* - \bar{c} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ a_{c,k} & a_{c,z} \end{bmatrix} \begin{bmatrix} k_{t-1} - \bar{k} \\ z_t - \bar{z} \end{bmatrix} + \begin{bmatrix} u_{k,t} \\ u_{c,t} \end{bmatrix}$$

Back to the Likelihood

- y_t consists of k_t^* and c_t^* and the model is given by (6) and (7).
- From the Kalman filter we get \hat{y}_t and $\Sigma_{\hat{y}_t}$

$$\hat{\mathbb{E}} \left[x_t | Y^{t-1}, \hat{x}_1 \right] = A \hat{\mathbb{E}} \left[x_{t-1} | Y^{t-2}, \hat{x}_1 \right] + K_{t-1} \hat{y}_{t-1}$$

$$\hat{\mathbb{E}} \left[y_t | Y^{t-1}, \hat{x}_1 \right] = C \hat{\mathbb{E}} \left[x_t | Y^{t-1}, \hat{x}_1 \right]$$

$$\hat{y}_t = y_t - \hat{\mathbb{E}} \left[y_t | Y^{t-1}, \hat{x}_1 \right]$$

$$\Sigma_{\hat{x}_{t+1}} = A \Sigma_{\hat{x}_t} A' + G V_1 G' - K_t (A \Sigma_{\hat{x}_t} C + G V_3)'$$

$$\Sigma_{\hat{y}_t} = C \Sigma_{\hat{x}_t} C'$$

Back to the Likelihood

- \hat{y}_{t+1} is normally distributed because
 - this is a linear model and underlying shocks are linear
- Kalman filter generates \hat{y}_{t+1} and $\Sigma_{\hat{y}_t}$
 - (given Ψ and observables, Y^T)
- Given normality calculate likelihood of $\{\hat{y}_{t+1}\}$

Kalman Filter versus inversion

with measurement error

- have to use Kalman filter

without measurement error

- could back out shocks using inverse of policy function
- but could also use Kalman filter
 - Dynare always uses the Kalman filter
 - hardest part of the Kalman filter is calculating the inverse of $C\Sigma_{\hat{x}_t}C' + V_2$ and this is typically not a difficult inversion.

Log-Likelihood

$$\begin{aligned}\ln(Y^T|\Psi) &= -\left(\frac{1}{2}\right) \left(n_x \ln(2\pi) + \ln(|\Sigma_{\hat{x}_0}|) + \hat{x}_0' \Sigma_{\hat{x}_0}^{-1} \hat{x}_0 \right) \\ &\quad - \left(\frac{1}{2}\right) \left(T n_y \ln(2\pi) + \sum_{t=1}^T \left[\ln(|\Sigma_{\hat{y}_t}|) + \hat{y}_t' \Sigma_{\hat{y}_t}^{-1} \hat{y}_t \right] \right)\end{aligned}$$

n_y : dimension of \hat{y}_t

For the neo-classical growth model

- Start with $x_1 = [k_0, z_0]$, $y_1 = k_0^*$, and Σ_1
- Calculate

$$\begin{aligned}\hat{y}_1 &= y_1 - \hat{E}[y_1|x_1] \\ &= y_1 - Cx_1\end{aligned}$$

- Calculate $\hat{E}[x_2|y_1, x_1]$ using

$$\hat{E}_t x_{t+1} = A \hat{E}_{t-1} x_t + K_t \hat{y}_t$$

where

$$K_t = (A \Sigma_{\hat{x}_t} C' + G V_3) (C \Sigma_{\hat{x}_t} C' + V_2)^{-1}$$

For the neo-classical growth model

- Calculate

$$\begin{aligned}\hat{y}_2 &= y_2 - \hat{E}[y_2|y_1, x_1] \\ &= y_2 - C\hat{E}[x_2|y_1, x_1]\end{aligned}$$

- etc.

Bayesian Estimation

- Conceptually, things are not that different
- Bayesian econometrics combines
 - the likelihood, i.e., the data, with
 - the prior
- You can think of the prior as additional data

Posterior

The joint density of parameters and data is equal to

$$P(Y^T, \Psi) = L(Y^T | \Psi)P(\Psi) \text{ or}$$

$$P(Y^T, \Psi) = P(\Psi | Y^T)P(Y^T)$$

Posterior

From this we can get Bayes rule: $P(\Psi|Y^T) = \frac{L(Y^T|\Psi)P(\Psi)}{P(Y^T)}$



Reverend Thomas Bayes (1702-1761)

Posterior

- For the distribution of Ψ , $P(Y^T)$ is just a constant.
- Therefore we focus on

$$L(Y^T|\Psi)P(\Psi) \propto \frac{L(Y^T|\Psi)p(\Psi)}{P(Y^T)} = P(\Psi|Y^T)$$

- One can always make $L(Y^T|\Psi)P(\Psi)$ a proper density by scaling it so that it integrates to 1

Evaluating the posterior

- Calculating posterior for given value of Ψ not problematic.
- But we are interested in objects of the following form

$$\mathbb{E}[g(\Psi)] = \frac{\int g(\Psi)P(\Psi|Y^T)d\Psi}{\int P(\Psi|Y^T)d\Psi}$$

- Examples
 - to calculate the mean of Ψ , let $g(\Psi) = 1$
 - to calculate the probability that $\Psi \in \Psi^*$,
 - let $g(\Psi) = 1$ if $\Psi \in \Psi^*$ and
 - let $g(\Psi) = 0$ otherwise
 - to calculate the posterior for j^{th} element of Ψ
 - $g(\Psi) = \Psi_j$

Evaluating the posterior

- Even *Likelihood* can typically only be evaluated numerically
- Numerical techniques also needed to evaluate the *posterior*

Evaluating the posterior

- Standard Monte Carlo integration techniques cannot be used
 - Reason: cannot *draw* random numbers directly from $P(\Psi|Y^T)$
 - being able to calculate $P(\Psi|Y^T)$ not enough to create a random number generator with that distribution
- Standard tool: Markov Chain Monte Carlo (MCMC)

Metropolis & Metropolis-Hastings

- Metropolis & Metropolis-Hastings are particular versions of the MCMC algorithm
- Idea:
 - travel through the state space of Ψ
 - weigh the outcomes appropriately

Metropolis & Metropolis-Hastings

- Start with an initial value, Ψ_0
 - discard the beginning of the sample, the burn-in phase, to ensure choice of Ψ_0 does not matter

Metropolis & Metropolis-Hastings

Subsequent values, Ψ_{i+1} , are obtained as follows

- Draw Ψ^* using the "stand in" density $f(\Psi^* | \Psi_i, \theta_f)$
 - θ_f contains the parameters of $f(\cdot)$
- Ψ^* is a *candidate* for Ψ_{i+1}
 - $\Psi_{i+1} = \Psi^*$ with probability $q(\Psi_{i+1} | \Psi_i)$
 - $\Psi_{i+1} = \Psi_i$ with probability $1 - q(\Psi_{i+1} | \Psi_i)$

Metropolis & Metropolis-Hastings

properties of $f(\cdot)$

- $f(\cdot)$ should have fat tails relative to the posterior
 - that is, $f(\cdot)$ should "cover" $P(\Psi|Y^T)$

Metropolis (used in Dynare)

$$q(\Psi_{i+1} | \Psi_i) = \min \left[1, \frac{P(\Psi^* | Y^T)}{P(\Psi_i | Y^T)} \right]$$

- $P(\Psi^* | Y^T) \geq P(\Psi_i | Y^T) \implies$
 - always include candidate as new element
- $P(\Psi^* | Y^T) < P(\Psi_i | Y^T) \implies$
 - Ψ^* not always included; the lower $P(\Psi^* | Y^T)$ the lower the chance it is included

Metropolis-Hastings

$$q(\Psi_{i+1}|\Psi_i) = \min \left[1, \frac{P(\Psi^*|Y^T)/f(\Psi^*|\Psi_i, \theta_f)}{P(\Psi_i|Y^T)/f(\Psi_i|\Psi_i, \theta_f)} \right]$$

- $P(\Psi_i|Y^T)/f(\Psi_i|\Psi_i, \theta_f)$ low \implies
 - you should move away from this Ψ value $\implies q$ should be high
- $P(\Psi^*|Y^T)/f(\Psi^*|\Psi_i, \theta_f)$ high:
 - probability of Ψ^* high & should be included with high prob.

Choices for $f(\cdot)$

- Random walk MH:

$$\Psi^* = \Psi_i + \varepsilon \text{ with } E[\varepsilon] = 0$$

- and, for example,

$$\varepsilon \sim N(0, \theta_f^2)$$

- Independence sampler:

$$f(\Psi^* | \Psi_i, \theta_f) = f(\Psi^* | \theta_f)$$

Couple more points

- Is the singularity issue different with Bayesian statistics?
- Choosing prior
- Gibbs sampler

The singularity problem again

What happens in practice?

- lots of observations are available
- practitioners don't want to exclude data \implies
- add "structural" shocks

The singularity problem again

Problem with adding additional shocks

- measurement error shocks
 - not credible that this is reason for gap between model and data
- structural shocks
 - good reason, but wrong structural shocks \implies misspecified model

Possible solution to singularity problem?

Today's posterior is tomorrow's prior

Possible solution to singularity problem?

Suppose you want the following:

- use 2 observables and
- only 1 structural shock

Possible solution to singularity problem?

- ① Start with first prior: $P_1(\Psi)$
- ② Use first observable Y_1^T to form first posterior

$$F_1(\Psi) = L(Y_1^T | \Psi)P_1(\Psi)$$

- ③ Let second prior be first posterior: $P_2(\Psi) = F_1(\Psi)$
- ④ Use second observable Y_2^T to form second posterior

$$F_2(\Psi) = L(Y_2^T | \Psi)P_2(\Psi)$$

Final answer:

$$\begin{aligned}F_2(\Psi) &= L(Y_2^T | \Psi) P_2(\Psi) \\&= L(Y_2^T | \Psi) L(Y_1^T | \Psi) P_1(\Psi)\end{aligned}$$

Obviously:

$$\begin{aligned}F_2(\Psi) &= L(Y_2^T | \Psi) L(Y_1^T | \Psi) P_1(\Psi) \\&= L(Y_1^T | \Psi) L(Y_2^T | \Psi) P_1(\Psi)\end{aligned}$$

Thus, it does not matter which variable you use first

Properties of final posterior

- Final posterior could very well have multiple modes
 - indicates where different variables prefer parameters to be
- This is only informative, not a disadvantage

Have we solved the singularity problem?

Problems of approach:

- Procedure avoids singularity problem by not considering *joint* implications of two observables
- Procedure misses some structural shock/misspecification

Key question:

- Is this worse than adding bogus shocks?

Have we solved the singularity problem?

Problems of approach:

- Procedure avoids singularity problem by not considering *joint* implications of two observables
- Procedure misses some structural shock/misspecification

Key question:

- Is this worse than adding bogus shocks?

How to choose prior

- ➊ Without analyzing data, sit down and think problem in macro: we keep on using the same data so is this science or data mining?
- ➋ Don't change prior depending on results

Uninformative prior

- $P(\Psi) = 1 \quad \forall \Psi \in \mathbb{R} \implies \text{posterior} = \text{likelihood}$
- $P(\Psi) = 1 / (b - a)$ if $\Psi \in [a, b]$ is not **un**informative
- Which one is the least informative prior?

$$P(\Psi) = 1 / (b - a) \text{ if } \Psi \in [a, b]$$
$$P(\ln \Psi) = 1 / (\ln b - \ln a) \text{ if } \Psi \in [\ln a, \ln b]$$

Uninformative prior

- $P(\Psi) = 1 \quad \forall \Psi \in \mathbb{R} \implies \text{posterior} = \text{likelihood}$
- $P(\Psi) = 1 / (b - a)$ if $\Psi \in [a, b]$ is not **un**informative
- Which one is the least informative prior?

$$P(\Psi) = 1 / (b - a) \text{ if } \Psi \in [a, b]$$
$$P(\ln \Psi) = 1 / (\ln b - \ln a) \text{ if } \Psi \in [\ln a, \ln b]$$

Uninformative prior

- $P(\Psi) = 1 \quad \forall \Psi \in \mathbb{R} \implies \text{posterior} = \text{likelihood}$
- $P(\Psi) = 1 / (b - a)$ if $\Psi \in [a, b]$ is not **un**informative
- Which one is the least informative prior?

$$P(\Psi) = 1 / (b - a) \text{ if } \Psi \in [a, b]$$
$$P(\ln \Psi) = 1 / (\ln b - \ln a) \text{ if } \Psi \in [\ln a, \ln b]$$

The objective of Jeffrey's prior is to ensure that the prior is *invariant* to such reparameterizations

How to choose (not so) informative priors

Let the prior inherit invariance structure of the problem:

- ❶ **location parameter:** If X is distributed as $f(x - \psi)$, then $Y = X + \phi$ have the same distribution but a different location. If the prior has to inherit this property, then it should be uniform.
- ❷ **scale parameter:** If X is distributed as $(1/\sigma)f(x/\sigma)$, then $Y = \phi X$ has the same distribution as X except for a different scale parameter. If the prior has to inherit this property, then it should be of the form

$$P(\psi) = 1/\psi$$

Both are improper priors.

That is, they do not integrate to a finite number.

Not so informative priors

Let the prior be consistent with "total confusion"

③ **probability parameter:** If ψ is a probability $\in [0, 1]$, then the prior distribution

$$P(\psi) = 1 / (\psi(1 - \psi))$$

represents total confusion. The idea is that the elements of the prior correspond to different beliefs and everybody is given a new piece of info that the cross-section of beliefs would not change.

See notes by Smith

Gibbs sampler

Objective: Obtain T observations from $p(x_1, \dots, x_J)$.

Procedure:

- ① Start with initial observation $X^{(0)}$.
- ② Draw period t observation, $X^{(t)}$, using the following iterative scheme:
 - draw $x_j^{(t)}$ from the conditional distribution:
$$p(x_j|x_1^{(t)}, \dots, x_{j-1}^{(t)}, x_{j+1}^{(t-1)}, \dots, x_J^{(t-1)})$$

Gibbs sampler versus MCMC

- Gibbs sampler does not require stand-in distribution
- Gibbs sampler still requires the ability to draw from conditional
 \implies not useful for estimation DSGE models

References

- Chib, S. and Greenberg, E., 1995, Understanding the Metropolis-Hastings Algorithm, *The American Statistician*.
 - describes the basics
- Ljungqvist, L. and T.J. Sargent, 2004, Recursive Macroeconomic Theory
 - source for the description of the Kalman filter
- Roberts, G.O., and J.S. Rosenthal, 2004, General state space Markov chains and MCMC algorithms, *Probability Surveys*.
 - more advanced articles describing formal properties

References

- Smith, G.P., Expressing Prior Ignorance of a Probability Parameter, notes, University of Missouri
<http://www.stats.org.uk/priors/noninformative/Smith.pdf>
on informative priors
- Syversveen, A.R, 1998, Noninformative Bayesian priors.
Interpretation and problems with construction and applications
<http://www.stats.org.uk/priors/noninformative/Syversveen1998.pdf>
on informative priors