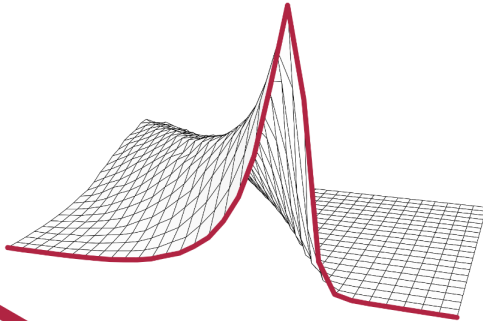


# mcxs



# Macroeconometrics

## Lecture 4 Numerical optimization and integration

**Tomasz Woźniak**

Department of Economics  
University of Melbourne

## Numerical optimization

## Numerical integration

### Materials:

Woźniak (2021) Posterior derivations for a simple linear regression model,  
Lecture notes

An R file `L4_mcx.R` for the reproduction the results

## **Numerical optimization**

Maximizing the log-likelihood function

# Numerical optimization

## **Motivation.**

For many econometric models the MLE cannot be found analytically as the system of equations for the first order conditions cannot or is difficult to solve.

$$G(\hat{\theta}) = \mathbf{0}$$

In such cases, we rely on numerical optimization methods that potentially give an approximate solution to the problem above that is as close to the exact solution as possible or required.

# Numerical optimization: the idea

Use an algorithm that requires:

**starting values** for the parameter vector, denoted by  $\theta_{(0)}$ , at which the algorithm begins the search of the solution

**a dynamic rule** that generates values of parameters in subsequent iterations of the algorithm, denoted by  $\theta_{(k)}$ , that are closer and closer to the solution

**a stopping rule** that stops the algorithm at a point that is close enough to the solution given the required precision

# Numerical optimization: starting values

Are usually generated from:

**preliminary data analysis** some summary statistics can be informative about approximate values of the parameters

**a simplified model** that can be easily estimated using a simpler method

**a grid of admissible values** that is a robust way of searching for the global maximum (away from a boundary of the parameter space) in cases where the solutions provided by algorithm are sensitive to starting values

# Numerical optimization: Newton-Raphson method

$$\theta_{(k)} = \theta_{(k-1)} - H^{-1}(\theta_{(k-1)}) G(\theta_{(k-1)})$$

$\theta_{(k)}$  – value of parameters in the current step of the algorithm

$\theta_{(k-1)}$  – value of parameters in the previous step of the algorithm

$H^{-1}(\theta_{(k-1)})$  – inverse of Hessian matrix evaluated at  $\theta_{(k-1)}$

$G(\theta_{(k-1)})$  – the gradient vector evaluated at  $\theta_{(k-1)}$



# Numerical optimization: Newton-Raphson method

**Gradient vector and Hessian matrix** are computed using:

**analytical formulae** – provided by the user

**numerical approximation** – which is the default option in existing software packages, it is also time-consuming

**automatic differentiation** – where the exact derivatives of the objective function are computed using designated software  
– no existing software packages offer this functionality

# Numerical optimization: BHHH method

Based on Brendt, Hall, Hall, Hausman (1974)

$$\theta_{(k)} = \theta_{(k-1)} + J^{-1}(\theta_{(k-1)}) G(\theta_{(k-1)})$$

Hessian matrix ordinate  $-H^{-1}(\theta_{(k-1)})$  is approximated by the outer product of gradients  $J^{-1}(\theta_{(k-1)})$ .

$$J(\theta_{(k-1)}) = \frac{1}{T} \sum_{t=1}^T g_t(\theta_{(k-1)}) g_t(\theta_{(k-1)})'$$
$$g_t(\theta_{(k-1)}) = \left. \frac{\partial l(y_t | x_t, \theta)}{\partial \theta} \right|_{\theta = \theta_{(k-1)}}$$

# Numerical optimization: stopping rules

Let  $\epsilon$  be a small positive number.

$\epsilon = 1e - 06$  – for preliminary estimations

$\epsilon < 1e - 08$  – for reporting final results

## Stopping rules.

Stop the algorithm when either of the criterion holds:

$|G(\theta_{(k)})| < \epsilon$  – related directly to absolute value of gradient

$\|\theta_{(k)} - \theta_{(k-1)}\| < \epsilon$  – an increment in parameter values over iterations is negligible, where  $\|\cdot\|$  is some norm, e.g., a maximum value of elements of vector

$k > k_{max}$  – a rule thanks to which we avoid running the computations for a long time without achieving convergence. This might happen when we set  $\epsilon$  to too small a value. This rule should not be binding for the final results.

Always check the message regarding the convergence!

# Numerical optimization: illustration

**A simple linear regression model** with known slope  $\beta = 1$

$$Y = X + E$$

$$E|X \sim \mathcal{N}(\mathbf{0}_T, \sigma^2 I_T)$$

$\downarrow$

$$Y|X, \beta = 1 \sim \mathcal{N}(X, \sigma^2 I_T)$$

**The log-likelihood function.**

$$l(\sigma^2|Y, X) = -\frac{T}{2} \log(2\pi) - \frac{T}{2} \log(\sigma^2) - \frac{1}{2} \frac{1}{\sigma^2} (Y - X)'(Y - X)$$

**MLE:**  $\hat{\sigma}^2 = \frac{1}{T} (Y - X)'(Y - X) = \frac{1}{T} E'E$

# Numerical optimization: illustration

## Derivatives.

$$G(\sigma^2) = -\frac{T}{2} \frac{1}{\sigma^2} + \frac{1}{2} \frac{E'E}{(\sigma^2)^2}$$

$$g_t(\sigma^2) = -\frac{T}{2} \frac{1}{\sigma^2} + \frac{1}{2} \frac{(y_t - x_t)^2}{(\sigma^2)^2}$$

$$H(\sigma^2) = \frac{T}{2} \frac{1}{(\sigma^2)^2} - \frac{E'E}{(\sigma^2)^3}$$

## Newton-Raphson rule.

$$\begin{aligned}\sigma_{(k)}^2 &= \sigma_{(k-1)}^2 + \left[ \frac{E'E}{(\sigma_{(k-1)}^2)^3} - \frac{T}{2} \frac{1}{(\sigma_{(k-1)}^2)^2} \right]^{-1} \left[ \frac{1}{2} \frac{E'E}{(\sigma_{(k-1)}^2)^2} - \frac{T}{2} \frac{1}{\sigma_{(k-1)}^2} \right] \\ &= \sigma_{(k-1)}^2 + \sigma_{(k-1)}^2 \left[ 2E'E - T\sigma_{(k-1)}^2 \right]^{-1} \left[ E'E - T\sigma_{(k-1)}^2 \right]\end{aligned}$$

# Numerical optimization: illustration

# Numerical optimization: illustration

## Data for the animation.

$k$	$\sigma_{(k)}^2$	$l(\sigma_{(k)}^2   Y, X)$	$ G(\sigma_{(k)}^2) $	$ G(\sigma_{(k)}^2)  < \epsilon$	$ \sigma_{(k)}^2 - \sigma_{(k-1)}^2 $	$ \sigma_{(k)}^2 - \sigma_{(k-1)}^2  < \epsilon$
1	0.500	-249.287	284.1	FALSE	0.213	FALSE
2	0.713	-209.707	118.9	FALSE	0.275	FALSE
3	0.988	-188.495	47.8	FALSE	0.323	FALSE
4	1.311	-178.686	17.6	FALSE	0.316	FALSE
5	1.627	-175.253	5.6	FALSE	0.216	FALSE
6	1.842	-174.567	1.1	FALSE	0.072	FALSE
7	1.9144	-174.523	0.08	FALSE	0.006	FALSE
8	1.9205	-174.523	0.001	FALSE	0.000	FALSE
9	1.9205	-174.523	0.000	TRUE	0.000	TRUE

$$\hat{\sigma}^2 = 1.9205, \quad \epsilon = 10^{-6}$$

## Numerical integration Gibbs sampler



# Numerical integration

## **Motivation.**

For most of econometric models analytical derivation of the joint posterior distribution of the parameters is impossible. The distribution is known only up to its kernel:

$$p(\theta|Y) \propto L(\theta|Y)p(\theta)$$

In such cases, we use numerical integration algorithms from the family of Monte Carlo Markov Chain (MCMC) methods to generate random draws from the joint posterior distribution and we use them to compute all of the required characteristics of this distribution.

Gibbs sampler is an MCMC algorithm that is feasible for most of the models discussed in this subject.

# Gibbs sampler

## Idea for the algorithm.

Suppose that the parameters can be conveniently divided in two groups:  $\theta = (\theta_1, \theta_2)$ .

For many models, when the joint posterior distribution  $p(\theta_1, \theta_2 | Y)$  cannot be derived, the full conditional posterior distributions,  $p(\theta_1 | Y, \theta_2)$  and  $p(\theta_2 | Y, \theta_1)$ , are available.

# Gibbs sampler

## Construction of the algorithm.

**Initialize**  $\theta_2$  at  $\theta_2^{(0)}$

**At each iteration**  $s$ :

1. **Draw** a random number/vector from  $\theta_1^{(s)} \sim p(\theta_1 | Y, \theta_2^{(s-1)})$
2. **Draw** a random number/vector from  $\theta_2^{(s)} \sim p(\theta_2 | Y, \theta_1^{(s)})$

**Repeat** steps 1. and 2.  $S_1 + S_2$  times

**Discard** the first  $S_1$  draws that allowed the algorithm to converge to the stationary posterior distribution

**Output** is a sample of draws from the joint posterior distribution

$$\{\theta_1^{(s)}, \theta_2^{(s)}\}_{s=S_1+1}^{S_2}$$

# A simple linear regression model

$$\begin{aligned} Y &= \beta X + E \\ E|X &\sim \mathcal{N}(\mathbf{0}_T, \sigma^2 I_T) \\ &\downarrow \\ Y|X &\sim \mathcal{N}(\beta X, \sigma^2 I_T) \end{aligned}$$

**The likelihood function.**

$$L(\theta|Y, X) = (2\pi)^{-\frac{T}{2}} (\sigma^2)^{-\frac{T}{2}} \exp \left\{ -\frac{1}{2} \frac{1}{\sigma^2} (Y - \beta X)' (Y - \beta X) \right\}$$

# Conditionally-conjugate prior distribution

A conditionally-conjugate prior distribution leads to the full conditional posterior distribution of the same form.

Let  $(\beta, \sigma^2)$  follow an independent normal and inverse gamma 2 distribution:

$$p(\beta, \sigma^2) = p(\beta) p(\sigma^2)$$

$$p(\beta) = \mathcal{N}(\underline{\beta}, \underline{\sigma}_{\beta}^2)$$

$$p(\sigma^2) = \mathcal{IG2}(\underline{s}, \underline{\nu})$$

**pdf.**

$$p(\beta, \sigma^2) \propto \exp \left\{ -\frac{1}{2} \frac{1}{\underline{\sigma}_{\beta}^2} (\beta - \underline{\beta})' (\beta - \underline{\beta}) \right\} \times (\sigma^2)^{-\frac{\underline{\nu}+2}{2}} \exp \left\{ -\frac{1}{2} \frac{\underline{s}}{\sigma^2} \right\}$$

# Derive Gibbs sampler

**Full conditional posterior distribution of  $\beta$ :**  $p(\beta|Y, X, \sigma^2)$

Conditioning on  $\sigma^2$  implies that for the sake of deriving the full conditional posterior distribution of  $\beta$  we treat it as non-random.

**Bayes' rule.**

$$p(\beta|Y, X, \sigma^2) \propto L(\beta, \sigma^2|Y, X) p(\beta)$$

**kernel of the full conditional distribution.**

$$p(\beta|Y, X, \sigma^2) \propto \exp\left\{-\frac{1}{2} \frac{1}{\sigma^2} (Y - \beta X)'(Y - \beta X)\right\} \times \exp\left\{-\frac{1}{2} \frac{1}{\underline{\sigma}_\beta^2} (\beta - \underline{\beta})'(\beta - \underline{\beta})\right\}$$

# Derive Gibbs sampler

**Full conditional posterior distribution of  $\beta$ :**  $p(\beta|Y, X, \sigma^2)$

$$\begin{aligned} p(\beta|Y, X, \sigma^2) &\propto \exp\left\{-\frac{1}{2}\frac{1}{\sigma^2}(Y - \beta X)'(Y - \beta X)\right\} \exp\left\{-\frac{1}{2}\frac{1}{\underline{\sigma}_\beta^2}(\beta - \underline{\beta})'(\beta - \underline{\beta})\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\frac{1}{\sigma^2}(Y - \beta X)'(Y - \beta X) + \frac{1}{\underline{\sigma}_\beta^2}(\beta - \underline{\beta})'(\beta - \underline{\beta})\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\beta^2(\underline{\sigma}_\beta^{-2} + \sigma^{-2}X'X) - \beta 2(\underline{\sigma}_\beta^{-2}\underline{\beta} + \sigma^{-2}X'Y) + \dots\right]\right\} \end{aligned}$$

Which is in a form of a normal distribution.

# Derive Gibbs sampler

**Full conditional posterior distribution of  $\beta$ :**  $p(\beta|Y, X, \sigma^2)$

$$p(\beta|Y, X, \sigma^2) = \mathcal{N}(\bar{\beta}, \bar{\sigma}_{\beta}^2)$$

$$\bar{\sigma}_{\beta}^2 = (\underline{\sigma}_{\beta}^{-2} + \sigma^{-2}X'X)^{-1}$$

$$\bar{\beta} = \bar{\sigma}_{\beta}^2 (\underline{\sigma}_{\beta}^{-2}\underline{\beta} + \sigma^{-2}X'Y)$$



# Derive Gibbs sampler

**Full conditional posterior distribution of  $\sigma^2$ :**  $p(\sigma^2|Y, X, \beta)$

Conditioning on  $\beta$  implies that for the sake of deriving the full conditional posterior distribution of  $\sigma^2$  we treat it as non-random.

**Bayes' rule.**

$$p(\sigma^2|Y, X, \beta) \propto L(\beta, \sigma^2|Y, X) p(\sigma^2)$$

**kernel of the full conditional distribution.**

$$p(\sigma^2|Y, X, \beta) \propto (\sigma^2)^{-\frac{T}{2}} \exp\left\{-\frac{1}{2} \frac{1}{\sigma^2} (Y - \beta X)'(Y - \beta X)\right\} \times (\sigma^2)^{-\frac{\nu+2}{2}} \exp\left\{-\frac{1}{2} \frac{S}{\sigma^2}\right\}$$

# Derive Gibbs sampler

**Full conditional posterior distribution of  $\sigma^2$ :**  $p(\sigma^2 | Y, X, \beta)$

$$\begin{aligned} p(\sigma^2 | Y, X, \beta) &\propto (\sigma^2)^{-\frac{T}{2}} \exp\left\{-\frac{1}{2} \frac{1}{\sigma^2} (Y - \beta X)'(Y - \beta X)\right\} \times (\sigma^2)^{-\frac{\nu+2}{2}} \exp\left\{-\frac{1}{2} \frac{\underline{s}}{\sigma^2}\right\} \\ &= (\sigma^2)^{-\frac{T+\nu+2}{2}} \exp\left\{-\frac{1}{2} \frac{1}{\sigma^2} [(Y - \beta X)'(Y - \beta X) + \underline{s}]\right\} \\ &= (\sigma^2)^{-\frac{\bar{\nu}+2}{2}} \exp\left\{-\frac{1}{2} \frac{\bar{s}}{\sigma^2}\right\} \end{aligned}$$

The final line is the kernel of the inverse gamma 2 distribution.

# Derive Gibbs sampler

**Full conditional posterior distribution of  $\sigma^2$ :**  $p(\sigma^2|Y, X, \beta)$

$$p(\sigma^2|Y, X, \beta) = \text{IG2}(\bar{s}, \bar{\nu})$$

$$\bar{s} = \underline{s} + (Y - \beta X)'(Y - \beta X)$$

$$\bar{\nu} = \underline{\nu} + T$$

# Gibbs sampler

**Sampling random numbers from  $\mathcal{IG}2(s, \nu)$**

**Step 1:** Draw a random number from  $\tilde{s} \sim \chi^2(\nu)$  using  
R function `rchisq()`

**Step 2:** Return  $s/\tilde{s}$  as a draw from  $\mathcal{IG}2(s, \nu)$

**Sampling random numbers from  $\mathcal{N}_1(\mu, \sigma^2)$**

Use R function `rnorm()`

**Sampling random numbers from  $\mathcal{N}_N(\mu, \Sigma)$**

Use R function `rmvnorm()` from package `mvtnorm`

# Gibbs sampler

**Initialize**  $\sigma^2$  at  $\sigma^{2(0)}$

**At each iteration**  $s$ :

1. **Draw**  $\beta^{(s)} \sim p(\beta | Y, X, \sigma^{2(s-1)}) = \mathcal{N}(\bar{\beta}, \bar{\sigma}_{\beta}^2)$

2. **Draw**  $\sigma^{2(s)} \sim p(\sigma^2 | Y, X, \beta^{(s)}) = \mathcal{IG2}(\bar{s}, \bar{\nu})$

**Repeat** steps 1. and 2.  $S_1 + S_2$  times

**Discard** the first  $S_1$  draws that allowed the algorithm to converge to the stationary posterior distribution

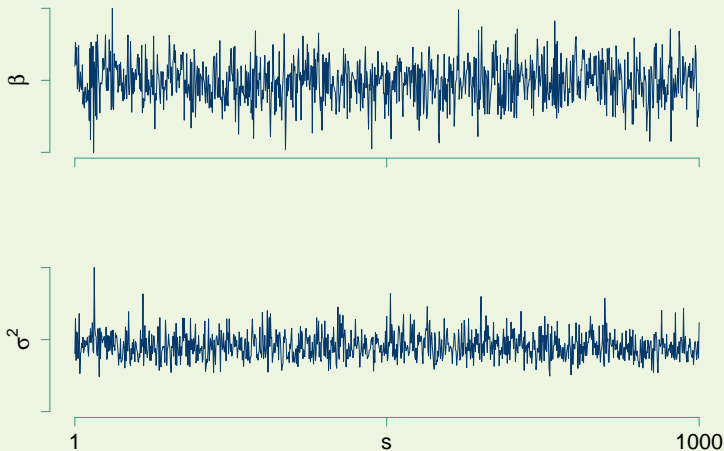
**Output** is a sample of draws from the joint posterior distribution

$$\left\{ \beta^{(s)}, \sigma^{2(s)} \right\}_{s=S_1+1}^{S_2}$$

# Gibbs sampler: illustration

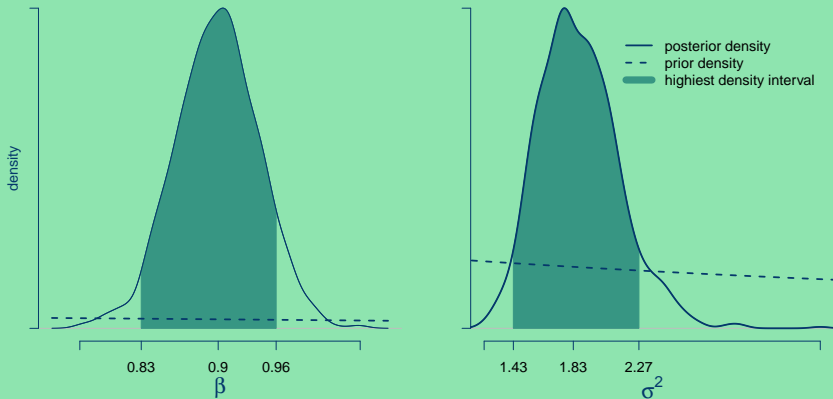
# Gibbs sampler: illustration

## Trace plots



# Gibbs sampler: illustration

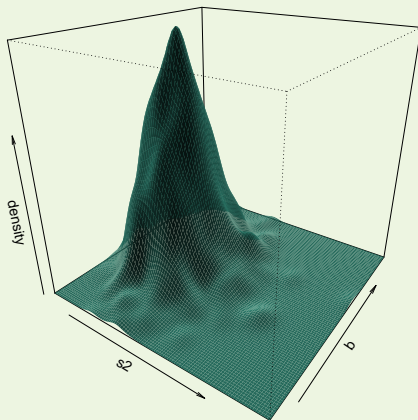
**Marginal posterior densities:**  $p(\beta|Y, X)$  and  $p(\sigma^2|Y, X)$





# Gibbs sampler: illustration

**Joint posterior density:**  $p(\beta, \sigma^2 | Y, X)$



## Numerical optimization

**Newton-Raphson** method

**BHHH** method

## Numerical integration

**Gibbs sampler**

**full conditional** posterior distributions

**joint and marginal** posterior distributions