

# Introduction to Bayesian Inference - 3

## Bayesian Estimation and Prediction

Eric Jacquier

Boston University Questrom School of Business

MF 840: Financial Econometrics

Spring 2021

*1. Basic Framework*

*2. Application to the Linear Model*

**3. Estimation**

**4. Prediction**

*5. Model Comparison (aka testing, Bayesians don't test)*

### 3. Estimation: Choice of a Location Estimator

#### 3.1 Estimation in classical vs Bayes frameworks

- So far, we have the posterior density of the unknown, random parameters.

We have only done the “statistical part of the job”, characterized our knowledge about the parameter.

We have not decided what “numbers” to report for the parameter.

**In Bayesian inference, this is a separate part of the process.**

- Contrast with classical inference where estimation and statistics are “mixed”

E.g., OLS, GLS, MLE give point estimates (the “number” to report), and also their statistical properties, often using asymptotic approximations.

- Option pricing model can be used for pricing or hedging – different purposes.

Hedging and pricing errors are different non-linear functions of the parameters (parameter: volatility).

Should we minimize sum of squared pricing errors or of squared hedging errors?

$$\text{Min}(C_i^m - C_i^{BS}(\sigma))^2 \quad / \quad \text{Min}(\Delta^m - \Delta^{BS}(\sigma))^2$$

We would find two different “numbers” for the parameter estimate!

Does this make sense?

### 3.2 Decision Theory to choose what point estimate to report

*Review MF793 Lecture Note 3, section 2*

- Rational decision on what parameter value to “report” depends on an error loss function, the loss of making an error.

Loss function for a choice  $\theta^*$  of the random parameter  $\theta$ :

Compute the expected loss:

$$EL(\theta^*) = \int L(\theta, \theta^*) p(\theta | D) d\theta$$

$L(\theta, \theta^*)$ :

*Cost of making an error  
by using that number*

... The integral is over the random true parameter  $\theta$

... We minimize it with respect to  $\theta^*$ , a choice of value on the posterior of  $\theta$ .

- Decision and information theorists adopted Bayesian Inference early because of the clean rational framework that separates statistical information processing (getting the posterior) from decision making (what point estimate to report).

- Mean, Median, Mode are all optimal location estimates to report for a specific Loss function

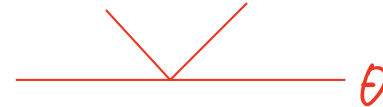
- **Quadratic** Loss: write it for a vector of parameters  $\theta$

$E_{\theta}[(\theta - \theta^*)' W (\theta - \theta^*)]$ , where  $W$  is a PDS weighting matrix

$$= \int (\theta - \theta^*)' W (\theta - \theta^*) p(\theta | D) d\theta$$

$$0 = \int -2W(\theta - \theta^*) p(\theta | D) d\theta \Rightarrow \theta^* = E(\theta | D)$$

- **Absolute** Loss: write it for a scalar parameter  $\theta$



$$E_{\theta}[|\theta - \theta^*|] = \int |\theta - \theta^*| p(\theta | D) d\theta$$

$$= \int_{-\infty}^{\theta^*} (\theta^* - \theta) p(\theta | D) d\theta + \int_{\theta^*}^{\infty} (\theta - \theta^*) p(\theta | D) d\theta$$

$$0 = \int_{-\infty}^{\theta^*} p(\theta | D) d\theta - \int_{\theta^*}^{\infty} p(\theta | D) d\theta$$

$$\Rightarrow \theta^* = \text{Median}(\theta | D)$$

- 0/1 Loss:  $L(\theta, \theta^*) = I_{\theta \neq \theta^*}$

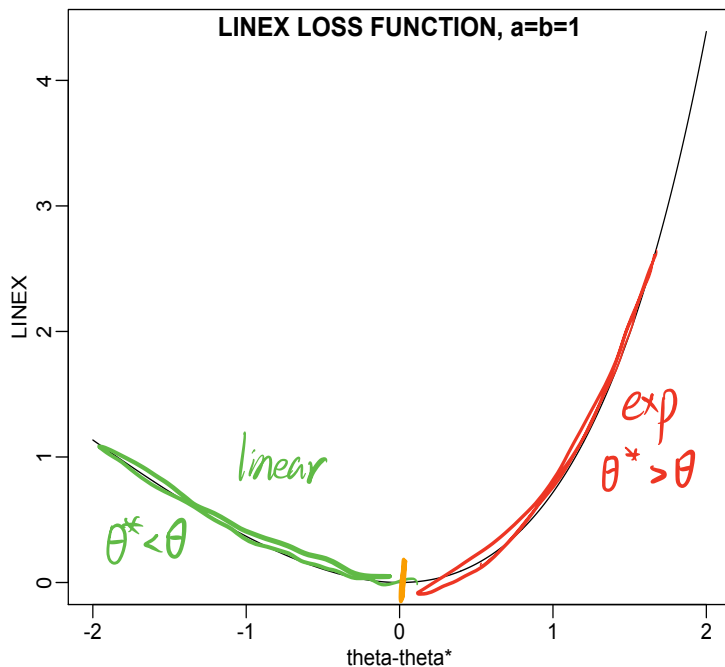
*Different cost of making an error  
if you under-or-overestimate.*

$$\theta^* = \text{Mode}(\theta|D)$$

*$\theta$  random*

- Hal Varian's **LINEX** linear-exponential Loss:

$$L(\theta, \theta^*) = b \left[ \overset{\text{exponential}}{e^{a(\theta^* - \theta)}} - \overset{\text{linear}}{a(\theta^* - \theta)} - 1 \right], \quad b > 0$$



Point estimate?

$$\frac{\partial E(L(\theta, \theta^*))}{\partial \theta^*}$$

$$0 = a e^{a\theta^*} E(e^{-a\theta}) - a$$

$$\Rightarrow \theta^* = \frac{-\text{Log}(E(e^{-a\theta}))}{a}$$

*no analytical solution*

Example:  $p(\theta|D) \sim N(\mu, \sigma)$

$$\theta^* = \frac{-\text{Log}(e^{-a\mu + 0.5a^2\sigma^2})}{a}$$

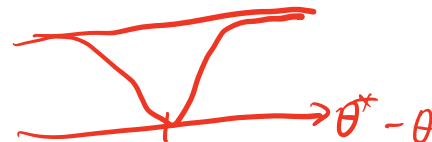
*like a risk aversion*

$$\theta^* = \mu - 0.5a^2\sigma^2$$

*Certainty equivalent*

前一项的 function 左右 goes to infinity

- Bounded Loss:  $L(\theta, \theta^*) = a[1 - e^{-b(\theta^* - \theta)^2}]$ ,  $a, b > 0$



- By Definition, the Bayesian point estimate is optimal given a choice of Error Loss.

Randomness is around the true parameter posterior distribution.

The data is fixed, it is **this** sample  $y$ .

Bayesians don't care what could happen for other hypothetical samples.

**One can always find a Bayesian point estimate minimizing EL for a given sample** *well-defined*

- Frequentist world: optimality results like Gauss-Markov for OLS under ideal assumptions, Cramer-Rao lower bound for MLE is approximate (large sample).

Randomness is around the estimator due to randomness of the data.  $\hat{\theta}$  random

Frequentists also look at a loss function  $L(\theta, \hat{\theta})$ , it is called the risk function. *fixed  $\theta$*

The frequentist loss function is random because  $\hat{\theta}$  is random, not  $\theta$ !

$\hat{\theta}$  is random because  $y$  is random (other samples)

Integration must be done over  $y$ , there is not trivial, there may be some  $\theta$  left:

*In classic, we can just say,  
"Asymptotically, minimize the MSE!"*

从许多 $\hat{\theta}$ 中找到 $\hat{\theta}$  to  $\min L(\theta, \theta)$  很难

**One can not always find a frequentist estimate that minimizes classical risk over random samples**

- The Bayesian and Classical (frequentist) approaches are a bit apples and oranges, we can't compare them strictly but there is one result:  
 $E(\theta|D)$  for my sample  $y^{(1)}$   
bayesian : optimal for every sample already } classic }  $y^{(2)}$   
nonny }  $y^{(n)}$   
about
  - 1. The Bayesian estimate is optimal by definition for the given sample  $y$  considered and the loss function chosen.
  - 2. Then minimized Loss can then be considered random over random samples  $y$ .  
classical :  $E(\theta|D)$  over samples
- Since the Bayesian estimate minimizes expected loss for every sample, it must be admissible over random sampling when viewed as a classical estimate.
- Enough with fundamental quarrels, let's get practical again

## 4 Unconditional Prediction in Bayesian Inference

### 4.1 Regression example

$$y = x\beta + \varepsilon \quad p(\beta | \sigma, D) p(\sigma | D)$$

Need to predict  $y_f = x_f' \beta + \varepsilon_f$ .

- Conditional prediction:  $p(y_f | \beta, \sigma, D, x_f) \sim N(x_f' \beta, \sigma)$

What  $\beta$  to use,  $\bar{\beta}$  the posterior mean? What  $\sigma$ ?

*Integrated out same to get true ~*

- But:  $\text{Var}(y_f | \sigma, D, x_f) = V(x_f' \beta + \varepsilon_f | D, x_f)$

$$= x_f' V(\beta | D) x_f + \sigma^2 = \sigma^2 (1 + x_f' (X'X)^{-1} x_f)$$

*diffuse prior  $\rightarrow \sigma^2 (x'x)^{-1}$*

$$E(y_f | \sigma, D, x_f) = x_f' E(\beta | D) = x_f' \bar{\beta}$$

We just informally integrated out  $\beta$ !

We have  $p(y_f | D, \sigma, x_f)$

- Just remains  $\sigma$ .

$$p(y_f | D) = \int \underbrace{\frac{1}{\sigma} e^{-\frac{(y_f - x_f' \bar{\beta})^2}{2 V(y_f | \sigma, D)}}}_{p(y_f | \sigma, D)} \underbrace{e^{-\frac{\sigma^2}{2 V(y_f | \sigma, D)}}}_{p(\sigma | D)} d\sigma = \int \frac{1}{\sigma^{3/2}} e^{-\frac{A}{2\sigma^2}} d\sigma$$



We should always use the **unconditional predictive** density .. to make predictions:

By definition of a marginal density:

$$y_f = \alpha + x_f \beta + \epsilon_f$$

$$p(y_F | D, x_F) = \int p(y_F | \beta, \sigma, D, x_F) \underbrace{p(\beta, \sigma | D)}_{\text{posterior}} d\beta d\sigma$$

- Let's integrate  $\sigma$  out. Using the integral result of the posterior (when we went from  $\beta | \sigma$  to  $\beta$ ), we have:

$$\frac{1}{[A]^{\frac{v+1}{2}}} = \int \frac{1}{\sigma} \cdot e^{-\frac{A}{2\sigma^2}}$$

$$\underbrace{y_F | D, x_F}_{\text{Student-t}} \propto [vs^2 + (y_F - x'_F \bar{\beta})' V^{-1} (y_F - x'_F \bar{\beta})]^{-\frac{v+1}{2}} \quad [1]$$

Where  $V = (1 + x'_F (X'X)^{-1} x_F)$   $\rightarrow$  diffuse priors

- Done .... But it's not always that simple
- When it's not simple, we will integrate by simulation, let's see an example

## 4.2 Multistep Forecasts, using Monte Carlo simulations

HW3

Aug 17, 2015, 8/24  
8/31

$[2 \dots T]$

$[1 \dots T-1]$

$$\text{AR}(1): Y_t = \alpha + \beta Y_{t-1} + \varepsilon_t \quad t = 1, \dots, T$$

$$P(\alpha, \beta | \sigma, D)$$

$$P(\sigma | D)$$

Multi-step forecasts

$$Y_{T+1} = \alpha + \beta Y_T + \varepsilon_{T+1} \quad Y_{T+1} | D \sim \text{Student-t per [1] above}$$

$$Y_{T+2} = \alpha + \beta Y_{T+1} + \varepsilon_{T+2} \quad Y_{T+2} | Y_{T+1}, D \sim \text{Student-t}$$

**NO!**  $(Y_{T+2} | \hat{Y}_{T+1}, D)$  is not correct, it does not reflect the whole uncertainty of  $Y_{T+2}$ .

**yes!**  $P(Y_{T+2} | D) = \int P(Y_{T+2} | Y_{T+1}, D) \underbrace{P(Y_{T+1} | D)}_{\text{Student-t}} dY_{T+1}$  No Analytical solution!

Integrating by simulation will be very easy though.

mean:  $AR_1: \frac{\alpha}{1-\beta}$

•  $y = X\beta + \varepsilon \quad \varepsilon_i \sim N(0, \sigma^2)$   
 •  $p(\sigma|D) \quad p(\beta|\sigma, D)$

- Integrating by Monte-Carlo simulation

	$\sigma D$	$\alpha, \beta \sigma, D$	$\varepsilon_{T+1} \rightarrow Y_{T+1}$	$\alpha + \beta Y_{T+1} = E(Y_{T+2} Y_{T+1})$	$\varepsilon_{T+2} \rightarrow Y_{T+2}$
Draw	(1)	(1)	(1) compute	compute	(1) compute
	(1.2)	..	..	..	..
10000	(S)	(S)	(S)	(S)	(S)

$\sigma$  每次不一样

$E(y_{T+2}|D)$

more precise for sample mean

no analytical sol

noise

by variance of  $\varepsilon_{T+1}$

- What is our first (the **seed**) draw?

$\sigma$  is the **only** parameter for which we have a marginal distribution, its posterior density is the **parent distribution of all our simulations**

- What are the draws of the second column,  $p(\alpha, \beta | D)$  or  $p(\alpha, \beta | \sigma, D)$ ?

MC Simulation

$$p(\alpha, \beta | D) = \int p(\alpha, \beta | \sigma, D) p(\sigma | D) d\sigma$$

$p(\sigma | D)$ : like a weight in the summation

- S draws:

$$\frac{1}{S} \sum_{i=1}^S Y_{T+2}^{(i)} \text{ estimates } E(Y_{T+2} | D),$$

How about the last column?

Then what does  $\frac{1}{S} \sum_1^S EY_{T+2} | Y_{T+1}$  do? (sample average of column 5)

- How do we estimate  $V(\beta|D)$ ? Quantiles of  $(\beta|D)$ ? or of  $Y_{T+2}|D$