

Lecture 8:  
Bayesian Estimation: Direct Sampling  
Big Data and Machine Learning for Applied Economics  
Econ 4676

Ignacio Sarmiento-Barbieri

Universidad de los Andes

September 2, 2021

# Agenda

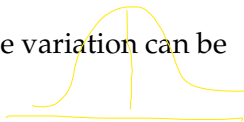
— 16 de sept

- 1 Bayesian Estimation
- 2 Simulation-based methods for Bayesian analysis
  - Direct Sampling
    - Gibbs Sampling
- 3 Recap
- 4 Further Readings

# Bayesian Estimation

$$X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2) \quad \hat{\mu} = \bar{X}$$

- ▶ The Bayesian approach to stats is fundamentally different from the classical approach we have been taking
- ▶ In the classical approach, the parameter  $\beta$  is thought to be an unknown, but fixed quantity, e.g.,  $X_i \sim f(\beta)$
- ▶ In the Bayesian approach  $\beta$  is considered to be a quantity whose variation can be described by a probability distribution (*prior distribution*)
- ▶ Then a sample is taken from a population indexed by  $\beta$  and the prior is updated with this sample
- ▶ The resulting updated prior is the *posterior distribution*



# Bayes Approach

## Bayes Theorem

$$\underbrace{\pi(\beta|X)}_{\text{posterior}} = \frac{\overbrace{f(X|\beta)}^{\text{density}} \underbrace{p(\beta)}_{\text{prior}}}{\underbrace{m(X)}_{\text{marginal}}} \Rightarrow \underbrace{f(X, \beta)}_{\text{joint}} = f(X|\beta) p(\beta) \quad (1)$$

with  $m(X)$  is the marginal distribution of  $X$ , i.e.

$$m(X) = \int f(X|\beta) p(\beta) d\beta \quad (2)$$

It is important to note that Bayes' theorem does not tell us what our beliefs should be, it tells us how they should change after seeing new information.

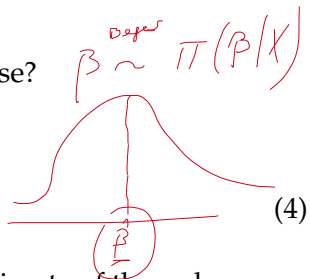
# Frequentist Approach

- The interest is on  $\beta$ , frequentist estimation procedures give us that, for example

$$\hat{\beta}_{MLE} = (X'X)^{-1}X'y \quad (3)$$

- Now in Bayes world, I have the full distribution. Which  $\beta$  I use?
- I can use any moment, but usually the interest lies on

$$E(\beta) = \int \beta \pi(\beta|X)$$



- Why? note that if you use MSE as loss function, the Bayes estimate of the unknown parameter is the mean of the posterior distribution

$$E(\beta_{\text{posterior}}) = w \beta_{MLE} + (1-w) \beta_0$$

CEP

HW

# Bayesian Estimation

- ▶ We are going to have an overview simulation-based methods for Bayesian analysis.
  - 1 Direct sampling algorithm
  - 2 Gibbs sampling algorithm
- ▶ As a running example, we use the linear regression framework

$$y_i = \beta x_i + u_i, \quad u_i \sim N(0, \sigma^2)$$

no context

(5)

- ▶ with  $\sigma^2$  known, and
- ▶ with prior distribution  $\beta \sim N(\beta_0, \tau^2)$

$\sigma^2$  unknown      $\sigma^2 \sim \text{IG}(\alpha, \beta)$

# Direct Sampling

- ▶ Using the knowledge of conjugate priors + the trick for exponentials
- ▶ The posterior distribution  $\beta$  follows the normal distribution:

$$\begin{aligned} A &= \frac{1}{\sigma^2} \sum_{i=1}^N x_i^2 + \frac{1}{\tau^2} \\ B &= \frac{1}{\sigma^2} \sum_{i=1}^N y_i x_i + \frac{1}{\tau^2} \beta_0 \\ m &= \frac{B}{A} \\ v &= \frac{1}{A} \end{aligned}$$

$$\beta | Y, X \sim N \left( \underbrace{\frac{\frac{1}{\sigma^2} \sum_{i=1}^N y_i x_i + \frac{1}{\tau^2} \beta_0}{\frac{1}{\sigma^2} \sum_{i=1}^N x_i^2 + \frac{1}{\tau^2}}}_{m}, \underbrace{\frac{1}{\frac{1}{\sigma^2} \sum_{i=1}^N x_i^2 + \frac{1}{\tau^2}}}_{Var} \right) \quad (6)$$

- ▶ We were able to characterize the full posterior distribution for the unknown object.

# Direct Sampling

- ▶ Suppose now, that our object of interest is not  $\beta$  per se, but some nonlinear function of unknown parameter  $\beta$ , e.g.  $h(\beta)$ .
- ▶ For example:
  - ▶  $h(\beta) = \beta$
  - ▶  $h(\beta) = |\beta|$  ✓
  - ▶  $h(\beta) = \alpha\%$  quantile of  $\beta$  ✓
  - ▶  $h(\beta) = \beta^3$  ✓
  - ▶  $h(\beta) = \underline{\beta_1 \beta_2}$  ✓
- ▶ The goal is to obtain posterior moments of  $h(\beta)$ .

$$y = \beta_1 \text{Room} + \beta_2 \text{Room}^2$$
$$\beta_1 + 2\beta_2 \text{Room} = 0$$
$$\text{Room} = \frac{-\beta_1}{2\beta_2}$$



# Direct Sampling

Side note: Frequentist's approach

- ▶ Frequentist obtain the sampling distribution of  $h(\beta)$  using the delta method:
- ▶ If we have *CLT*

$$\sqrt{N} (\hat{\beta} - \beta_0) \rightarrow_d N(0, V_{asy}) \quad (7)$$

- ▶ Then we

$$\sqrt{N} (h(\hat{\beta}) - h(\beta_0)) \rightarrow_d N(0, V_{asy} [h'(\beta_0)]^2) \quad (8)$$

- ▶ As  $N \rightarrow \infty$  where  $N$  is the number of observations.

# Direct Sampling

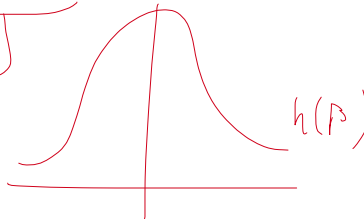
$$\pi(\beta|X) \sim \underline{\underline{\mu}}(m, V)$$

- ▶ Idea: Monte Carlo integration.
- ▶ Requirement
  - ▶ Know how to generate i.i.d. samples from the posterior distribution of  $\beta$ ,  $\pi(\beta|Y)$
- ▶ The requirement is satisfied for our linear regression example:
  - ▶ The posterior distribution of  $\beta$  follows the normal distribution. ✓
  - ▶ Most modern statistical program languages provide random number generators for many parametric distributions including the normal distribution. ✓

set seed(1001)  
set seed  
seed

# Direct Sampling

- ▶ Direct sampling approach simply approximates the posterior expectations of a function  $h(\beta)$  by

$$E_Y^\beta [\underline{h(\beta)}] = \int h(\beta) \pi(\beta|Y) d\beta \quad (9)$$
$$\approx \frac{1}{S} \sum_{i=1}^S h(\beta^i)$$


- ▶ Where  $\beta^i$  is i.i.d. samples from  $\pi(\beta|Y)$
- ▶  $S$  is "number of random samples from the posterior" or "number of generated draws" NOT the number of observations.

No confusion with Bootstrap

# Direct Sampling

- ▶ Provided that  $E_Y^\beta [h(\beta)^2] < \infty$ ,
- ▶ we can use the Strong Law of Large Numbers (SLLN)

$$\frac{1}{S} \sum_{i=1}^S h(\beta^i) \rightarrow a.s. \int h(\beta) p(\beta|Y) d\beta$$

- ▶ and the Central Limit Theorem (CLT)

$$\sqrt{S} \left( \frac{1}{S} \sum_{i=1}^S h(\beta^i) - \int h(\beta) \pi(\beta|Y) d\beta \right) \rightarrow_d N(0, V_\pi)$$

- ▶ Where

$$V_\pi = \text{Var}_Y^\beta (h(\beta)) = \int \left( h(\beta) - E_Y^\beta [h(\beta)] \right)^2 \pi(\beta|Y) d\beta$$

Correlation  
↓  
non perfect  
for nested  
approx

- ▶  $S$  is the number of simulated draws from the posterior distribution

# Direct Sampling

- Note that we turned a complicated integration into a simple average

$$\frac{1}{S} \sum_{i=1}^S h(\beta^i) \rightarrow a.s. \int h(\beta) p(\beta|Y) d\beta \quad (10)$$

- As the number of simulated draws increases, this simple average converges to the object of interest.
- Numerical accuracy?

# Direct Sampling

- ▶ The CLT result provides a way to measure the numerical accuracy of this
- ▶ Monte Carlo approximation:

$$\sqrt{S} \left( \frac{1}{S} \sum_{i=1}^N h(\beta^i) - \int h(\beta) p(\beta|Y) d\beta \right) \rightarrow_d N(0, V_\pi) \quad (11)$$

- ▶ That is,

*H w*

$$\underbrace{\frac{1}{S} \sum_{i=1}^S h(\beta^i)} \approx_d N \left( E_Y^\beta [h(\beta)], \underbrace{\frac{V_\pi}{S}} \right) \quad (12)$$

- ▶ Where  $V_\pi = \text{Var}_Y^\beta (h(\beta))$ . Posterior variance of  $h(\beta)$  scaled by  $1/S$  determines the numerical accuracy. As  $S \rightarrow \infty$ , numerical approximation goes to zero
  - ▶ Trade-off
    - ▶ Large ~~N~~: high computational cost (time) but more accurate approximation
    - ▶ Small ~~N~~: low computation cost (time) but less accurate approximation
- LL*  
*div 7, 1, 10, 100*

# Direct Sampling

## Example: Linear regression

- Consider the following linear regression model

$$y_i = \beta x_i + u_i, \quad u_i \sim N(0, \sigma^2) \quad (13)$$

- with prior distribution  $\beta \sim N(\beta_0, \tau^2)$ , and suppose  $\sigma^2$  is known.
- Then, we now all know that the posterior distribution  $\beta$  follows the normal distribution:

$$\beta | Y, X \sim N \left( \frac{\frac{1}{\sigma^2} \sum_{i=1}^N y_i x_i + \frac{1}{\tau^2} \beta_0}{\frac{1}{\sigma^2} \sum_{i=1}^N x_i^2 + \frac{1}{\tau^2}}, \frac{1}{\frac{1}{\sigma^2} \sum_{i=1}^N x_i^2 + \frac{1}{\tau^2}} \right) \quad (14)$$

*Posterior*

# Direct Sampling

## Example: Linear regression

- ▶ Goal: posterior mean and equal-tail-probability credible set for  $|\beta|$   $\rightarrow$  value observed

- ▶ I generate data  $(y_i, x_i)$  with

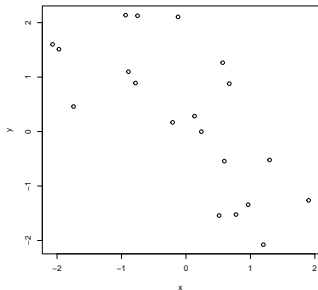
- ▶  $y_i = \beta x_i + u_i, \quad u_i \sim N(0, \sigma^2)$

- ▶  $N = 20$

- ▶  $\beta_{\text{true}} = -1$  and  $\sigma^2 = 1$

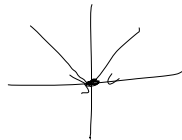
- ▶  $\beta_0 = 0$  and  $t = 100$

$|\beta_{\text{true}}| = 1$   $\rightarrow$  true



$$h(\beta) = |\beta|$$

$\rightarrow$  value observed



CLT

medians  $\sim N(\underline{m}, \frac{V}{t})$

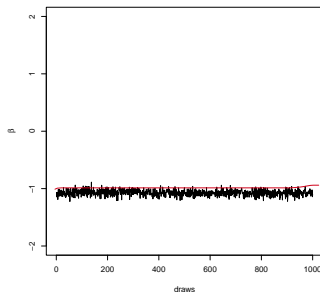


# Direct Sampling

Example: Linear regression

- ▶ Step 1: we generate  $N$  draws from the  $N(\underline{m}, \underline{V})$ ,  $\{\beta^i\}_{1, \dots, S}$
- ▶  $m = -1.07$
- ▶  $V = 0.0510$

Figure 1: Example of draws  $(\{\beta^i\}_{1, \dots, N})$ ,  $S = 1,000$



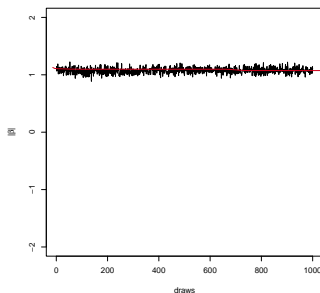
$\beta$

# Direct Sampling

Example: Linear regression

- ▶ Step 2: we are interested in posterior moments of  $|\beta|$ .
- ▶ Turn draws into  $\{|\beta|\}_{1,\dots,S}$

Figure 2: Example of draws  $\left(\{|\beta^i|\}_{1,\dots,S}\right)$

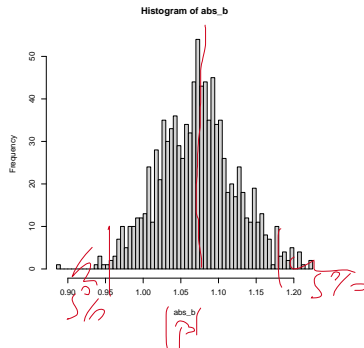


# Direct Sampling

Example: Linear regression

- Histogram approximation to  $\pi(|\beta| | Y)$  using  $\{|\beta^i|\}_{1,\dots,S}$

Figure 3: Example of draws  $(\{|\beta^i|\}_{1,\dots,S})$



# Direct Sampling

Example: Linear regression

- The posterior mean of  $|\beta|$  is approximated by

$$E_Y^\beta [|\beta|] \approx \frac{1}{S} \sum_{i=1}^S |\beta^i| = 1.0719 \quad (15)$$

- The 90% equal-tail-probability interval is approximated by

$$C_Y = [q_l, q_u] = [0.719, 1.441] \quad (16)$$

- Where  $q_l$  and  $q_u$  such that

$$5\% = \frac{1}{S} \sum_{i=1}^S 1\{|\beta| < q_u\} = 0.95 \quad (17)$$

# Direct Sampling

Example: Linear regression

- ▶ Numerical accuracy of  $\frac{1}{N} \sum_{i=1}^N |\beta^i|$
- ▶ We know that if we generate enough number of  $\beta^i$ , we get an accurate approximation to the posterior moments
- ▶ How many draws are enough?
- ▶ In other words, "Will I get different answer if I construct the same quantity using different set of draws  $\{\beta^i\}$ ?"
- ▶ Is  $S = 10$  enough? Or, is  $S = 10,000$  enough?

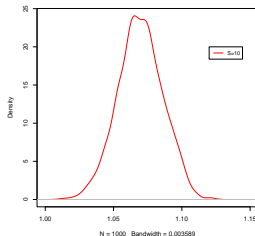
# Direct Sampling

Example: Linear regression

- ▶ To see the numerical error I generate 1,000 sets of  $\{\beta^i\}_{i=1,\dots,S}$
- ▶ Compute 1,000 of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  to see how variable this Monte Carlo approximation with different  $S$

$$\left[ \frac{1}{S} \sum_{i=1}^S |\beta^i| \right]$$

Figure 4: Distribution of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  over  $\{\beta^i\}_{i=1,\dots,S}$



$$\text{SD} \frac{1}{S} \sum_{i=1}^S |\beta^i| = 0.688$$

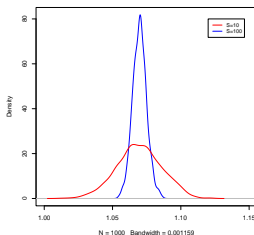


# Direct Sampling

## Example: Linear regression

- ▶ To see the numerical error I generate 1,000 sets of  $\{\beta^i\}_{i=1,\dots,S}$
- ▶ Compute 1,000 of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  to see how variable this Monte Carlo approximation with different  $S$

Figure 5: Distribution of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  over  $\{\beta^i\}_{i=1,\dots,S}$



$$\mu = 1.07$$

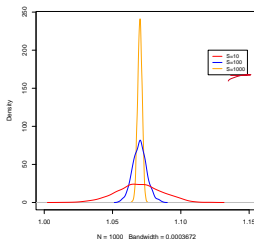
$$\text{SD} \left[ \frac{1}{N} \sum_{i=1}^N |\beta^i| \right] = 0.022$$

# Direct Sampling

## Example: Linear regression

- ▶ To see the numerical error I generate 1,000 sets of  $\{\beta^i\}_{i=1,\dots,S}$
- ▶ Compute 1,000 of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  to see how variable this Monte Carlo approximation with different  $S$

Figure 6: Distribution of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  over  $\{\beta^i\}_{i=1,\dots,S}$



$$\text{SD} \frac{1}{N} \sum_{i=1}^S |\beta^i| = 0.0074$$

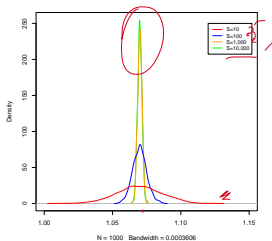


# Direct Sampling

## Example: Linear regression

- ▶ To see the numerical error I generate 1,000 sets of  $\{\beta^i\}_{i=1,\dots,S}$
- ▶ Compute 1,000 of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  to see how variable this Monte Carlo approximation with different  $S$

Figure 7: Distribution of  $\frac{1}{S} \sum_{i=1}^S |\beta^i|$  over  $\{\beta^i\}_{i=1,\dots,S}$



$$\text{SD} \frac{1}{N} \sum_{i=1}^N |\beta^i| = 0.0023$$

# Direct Sampling

Example: Linear regression

- ▶ What do we try to capture in this exercise?
- ▶ We try to mimic the distribution of Monte Carlo approximation offered by


*CLT*

$$\sqrt{S} \left( \frac{1}{S} \sum_{i=1}^N h(\beta^i) - \int h(\beta) \pi(\beta|Y) d\beta \right) \rightarrow_d N(0, \underline{V}_{\pi}) \quad (18)$$

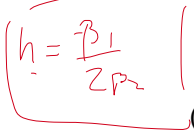
- ▶ All variation in this Monte Carlo approximation is due to numerical "simulation".
- ▶ Throughout this example, we fix  $Y_{1:N}, X_{1:N}$  (not a sampling variation).

## Recap

- ▶ If you know how to generate i.i.d draws from the posterior distribution of  $\beta$ ,
- ▶ You also can posterior moments of  $h(\beta)$  by simple average:

$$\int h(\beta) p(\beta|Y) d\beta \approx \frac{1}{N} \sum_{i=1}^N h(\beta^i) \quad (19)$$


- ▶ SLLN guarantees this Monte Carlo average to the right limit:

$$\frac{1}{N} \sum_{i=1}^N h(\beta^i) \xrightarrow{\text{a.s.}} \int h(\beta) \pi(\beta|Y) d\beta \quad (20)$$


- ▶ CLT tells you that the Monte Carlo average always has a numerical error:

$$\sqrt{S} \left( \frac{1}{S} \sum_{i=1}^N h(\beta^i) - \int h(\beta) \pi(\beta|Y) d\beta \right) \rightarrow_d N(0, V_\pi) \quad (21)$$

- ▶ It is important to check how good is your numerical approximation

# Review & Next Steps

- Direct Sampler

→ ccd → parallel/distrib

- Next Class: Gibbs Sampler

→ sequential  
→ de como distribuir

$y_i \rightarrow j_{i-1}$

## Further Readings

- ▶ Casella, G., & Berger, R. L. (2002). Statistical inference (Vol. 2, pp. 337-472). Pacific Grove, CA: Duxbury. Chapter 7
- ▶ Hoff, P. D. (2009). A first course in Bayesian statistical methods (Vol. 580). New York: Springer.

→ Casella → Introduction to Monte Carlo

→ Applied Bayesian Statistics

Casella

Open Bugs ★