

Overview of Bayesian Methods and the Regression Model

Bank of Korea Global Initiative Program

September 2014

- Reading: Chapter 1 of textbook and Appendix B, section B.1.
- Begin with general concepts in Bayesian theory before getting to specific models.
- If you know these general concepts you will never get lost.
- What does econometrician do? i) Estimate parameters in a model (e.g. regression coefficients), ii) Compare different models (e.g. hypothesis testing), iii) Prediction.
- Bayesian econometrics does these based on a few simple rules of probability.

- Let A and B be two events, $p(B|A)$ is the conditional probability of $B|A$. “summarizes what is known about B given A ”
- Bayesians use this rule with $A =$ something known or assumed (e.g. the Data), B is something unknown (e.g. coefficients in a model).
- Let y be data, y^* be unobserved data (i.e. to be forecast), M_i for $i = 1, \dots, m$ be set of models each of which depends on some parameters, θ^i .
- Learning about parameters in a model is based on the posterior density: $p(\theta^i | M_i, y)$
- Model comparison based on posterior model probability: $p(M_i | y)$
- Prediction based on the predictive density $p(y^* | y)$.

Bayes Theorem

- I expect you know basics of probability theory from previous studies, see Appendix B of my textbook if you do not.
- *Definition: Conditional Probability*
- The conditional probability of A given B , denoted by $\Pr(A|B)$, is the probability of event A occurring given event B has occurred.
- *Theorem: Rules of Conditional Probability including Bayes' Theorem*
- Let A and B denote two events, then
- $\Pr(A|B) = \frac{\Pr(A,B)}{\Pr(B)}$ and
- $\Pr(B|A) = \frac{\Pr(A,B)}{\Pr(A)}$.

- These two rules can be combined to yield *Bayes' Theorem*:

$$\Pr(B|A) = \frac{\Pr(A|B) \Pr(B)}{\Pr(A)}.$$

- *Note*: Above is expressed in terms of two events, A and B . However, can be interpreted as holding for random variables, A and B with probability density functions replacing the $\Pr()$ s in previous formulae.

Learning About Parameters in a Given Model (Estimation)

- Assume a single model which depends on parameters θ
- Want to figure out properties of the posterior $p(\theta|y)$
- It is convenient to use Bayes' rule to write the posterior in a different way.
- Bayes' rule lies at the heart of Bayesian econometrics:

$$p(B|A) = \frac{p(A|B)p(B)}{p(A)}.$$

- Replace B by θ and A by y to obtain:

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}.$$

- Bayesians treat $p(\theta|y)$ as being of fundamental interest: “Given the data, what do we know about θ ?”.
- Treatment of θ as a random variable is controversial among some econometricians.
- Competitor to Bayesian econometrics, called *frequentist econometrics*, says that θ is not a random variable.
- For estimation can ignore the term $p(y)$ since it does not involve θ :

$$p(\theta|y) \propto p(y|\theta)p(\theta).$$

- $p(\theta|y)$ is referred to as the *posterior density*
- $p(y|\theta)$ is the *likelihood function*
- $p(\theta)$ as the *prior density*.
- “posterior is proportional to likelihood times prior”.

- $p(\theta)$, does not depend on the data. It contains any non-data information available about θ .
- Prior information is controversial aspect since it sounds unscientific.
- Bayesian answers (to be elaborated on later):
 - i) Often we do have prior information and, if so, we should include it (more information is good)
 - ii) Can work with “noninformative” priors
 - iii) Can use “empirical Bayes” methods which estimate prior from the data
 - iv) Training sample priors
 - v) Bayesian estimators often have better frequentist properties than frequentist estimators (e.g. results due to Stein show MLE is inadmissible – but Bayes estimators are admissible)
 - vi) Prior sensitivity analysis

Prediction in a Single Model

- Prediction based on the *predictive density* $p(y^*|y)$
- Since a marginal density can be obtained from a joint density through integration:

$$p(y^*|y) = \int p(y^*, \theta|y) d\theta.$$

- Term inside integral can be rewritten as:

$$p(y^*|y) = \int p(y^*|y, \theta)p(\theta|y) d\theta.$$

- Prediction involves the posterior and $p(y^*|y, \theta)$ (more description provided later)

Model Comparison (Hypothesis testing)

- Models denoted by M_i for $i = 1, \dots, m$. M_i depends on parameters θ^i .
- *Posterior model probability* is $p(M_i|y)$.
- Using Bayes rule with $B = M_i$ and $A = y$ we obtain:

$$p(M_i|y) = \frac{p(y|M_i)p(M_i)}{p(y)}$$

- $p(M_i)$ is referred to as the *prior model probability*.
- $p(y|M_i)$ is called the *marginal likelihood*.

- How is marginal likelihood calculated?
- Posterior can be written as:

$$p(\theta^i|y, M_i) = \frac{p(y|\theta^i, M_i)p(\theta^i|M_i)}{p(y|M_i)}$$

- Integrate both sides with respect to θ^i , use fact that $\int p(\theta^i|y, M_i)d\theta^i = 1$ and rearrange:

$$p(y|M_i) = \int p(y|\theta^i, M_i)p(\theta^i|M_i)d\theta^i.$$

- Note: marginal likelihood depends only on the prior and likelihood.

- *Posterior odds ratio* compares two models:

$$PO_{ij} = \frac{p(M_i|y)}{p(M_j|y)} = \frac{p(y|M_i)p(M_i)}{p(y|M_j)p(M_j)}.$$

- Note: $p(y)$ is common to both models, no need to calculate.
- Can use fact that $p(M_1|y) + p(M_2|y) + \dots + p(M_m|y) = 1$ and PO_{ij} to calculate the posterior model probabilities.
- E.g. if $m = 2$ models:

$$p(M_1|y) + p(M_2|y) = 1$$

$$PO_{12} = \frac{p(M_1|y)}{p(M_2|y)}$$

- imply

$$p(M_1|y) = \frac{PO_{12}}{1 + PO_{12}}$$

$$p(M_2|y) = 1 - p(M_1|y).$$

- The *Bayes Factor* is:

$$BF_{ij} = \frac{p(y|M_i)}{p(y|M_j)}.$$

Summary

- These few pages have outlined all the basic theoretical concepts required for the Bayesian to learn about parameters, compare models and predict.
- This is an enormous advantage: Once you accept that unknown things (i.e. θ , M_i and y^*) are random variables, the rest of Bayesian approach is non-controversial.
- What are going to do in rest of this course?
- See how these concepts work in some models of interest.
- First the regression model
- Then time series models of interest for macroeconomics
- Bayesian computation.

- How do you present results from a Bayesian empirical analysis?
- $p(\theta|y)$ is a p.d.f. Especially if θ is a vector of many parameters cannot present a graph of it.
- Want features analogous to frequentist point estimates and confidence intervals.
- A common point estimate is the mean of the posterior density (or *posterior mean*).
- Let θ be a vector with k elements, $\theta = (\theta_1, \dots, \theta_k)'$. The posterior mean of any element of θ is:

$$E(\theta_i|y) = \int \theta_i p(\theta|y) d\theta.$$

- Aside *Definition B.8: Expected Value*
- Let $g(\cdot)$ be a function, then the *expected value* of $g(X)$, denoted $E[g(X)]$, is defined by:

$$E[g(X)] = \sum_{i=1}^N g(x_i) p(x_i)$$

- if X is discrete random variable with sample space $\{x_1, x_2, x_3, \dots, x_N\}$

-

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) p(x) dx$$

- if X is a continuous random variable (provided $E[g(X)] < \infty$).

- Common measure of dispersion is the *posterior standard deviation* (square root of *posterior variance*)
- Posterior variance:

$$\text{var}(\theta_i|y) = E(\theta_i^2|y) - \{E(\theta_i|y)\}^2,$$

- This requires calculating another expected value:

$$E(\theta_i^2|y) = \int \theta_i^2 p(\theta|y) d\theta.$$

- Many other possible features of interest. E.g. what is probability that a coefficient is positive?

$$p(\theta_i \geq 0|y) = \int_0^{\infty} p(\theta_i|y) d\theta_i$$

- All of these posterior features have the form:

$$E [g(\theta) | y] = \int g(\theta) p(\theta | y) d\theta,$$

- where $g(\theta)$ is a *function of interest*.
- All these features have integrals in them. Marginal likelihood and predictive density also involved integrals.
- Apart from a few simple cases, it is not possible to evaluate these integrals analytically, and we must turn to the computer.

- The integrals involved in Bayesian analysis are usually evaluated using simulation methods.
- Will use several methods later on. Here we provide some intuition.
- Frequentist asymptotic theory uses Laws of Large Numbers (LLN) and a Central Limit Theorems (CLT).
- A typical LLN: “consider a random sample, Y_1, \dots, Y_N , as N goes to infinity, the average converges to its expectation” (e.g. $\bar{Y} \rightarrow \mu$)
- Bayesians use LLN: “consider a random sample from the posterior, $\theta^{(1)}, \dots, \theta^{(S)}$, as S goes to infinity, the average of these converges to $E[\theta|y]$ ”
- Note: Bayesians use asymptotic theory, but asymptotic in S (under control of researcher) not N

- Example: Monte Carlo integration.
- Let $\theta^{(s)}$ for $s = 1, \dots, S$ be a random sample from $p(\theta|y)$ and define

$$\hat{g}_S = \frac{1}{S} \sum_{s=1}^S g \left(\theta^{(s)} \right),$$

- then \hat{g}_S converges to $E [g(\theta) | y]$ as S goes to infinity.
- Monte Carlo integration approximates $E [g(\theta) | y]$, but only if S were infinite would the approximation error be zero.
- We can choose any value for S (but larger values of S will increase computational burden).
- To gauge size of approximation error, use a CLT to obtain numerical standard error.

- Most Bayesians write own programs (e.g. using Gauss, Matlab, R or C++) to do posterior simulation
- BUGS (Bayesian Analysis Using Gibbs Sampling) is a popular Bayesian package, but only has limited set of models (or require substantial programming to adapt to other models)
- Bayesian work cannot (easily) be done in standard econometric packages like Microfit, Eviews or Stata.
- I have a Matlab website for VARs, TVP-VARs and TVP-FAVARs (see my website)
- See also <https://sites.google.com/site/dimitriskorobilis/matlab>
- Peter Rossi has an R package for marketing and microeconomic applications
- <http://www.perossi.org/home/bsm-1>
- Jim LeSage's Econometrics toolbox (Matlab)
- <http://www.spatial-econometrics.com/>
- Many more using R see <http://cran.r-project.org/web/views/Bayesian.html>

Bayesian Analysis of the Normal Linear Regression Model

- Assume k explanatory variables, x_{i1}, \dots, x_{ik} for $i = 1, \dots, N$ and regression model:

$$y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \varepsilon_i.$$

- Note x_{i1} is implicitly set to 1 to allow for an intercept.
- Matrix notation:

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ y_N \end{bmatrix}$$

- ε is $N \times 1$ vector stacked in same way as y

- β is $k \times 1$ vector
- X is $N \times k$ matrix

$$X = \begin{bmatrix} 1 & x_{12} & \cdot & \cdot & x_{1k} \\ 1 & x_{22} & \cdot & \cdot & x_{2k} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_{N2} & \cdot & \cdot & x_{Nk} \end{bmatrix}$$

- Regression model can be written as:

$$y = X\beta + \varepsilon.$$

The Likelihood Function

- Likelihood can be derived under the classical assumptions:
- ε is $N(0_N, h^{-1}I_N)$ where $h = \sigma^{-2}$.
- All elements of X are either fixed (i.e. not random variables).
- Exercise 10.1, Bayesian Econometric Methods shows that likelihood function can be written in terms of OLS quantities:

$$\begin{aligned}v &= N - k, \\ \hat{\beta} &= (X'X)^{-1} X'y \\ s^2 &= \frac{(y - X\hat{\beta})' (y - X\hat{\beta})}{v}\end{aligned}$$

- Likelihood function:

$$p(y|\beta, h) = \frac{1}{(2\pi)^{\frac{N}{2}}} \left\{ h^{\frac{1}{2}} \exp \left[-\frac{h}{2} (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) \right] \right\} \left\{ h^{\frac{v}{2}} \exp \left[-\frac{hv}{2s^2} \right] \right\}$$

The Prior

- Common starting point is natural conjugate Normal-Gamma prior
- β conditional on h is now multivariate Normal:

$$\beta|h \sim N(\underline{\beta}, h^{-1}\underline{V})$$

- Prior for error precision h is Gamma

$$h \sim G(\underline{s}^{-2}, \underline{v})$$

- $\underline{\beta}$, \underline{V} , \underline{s}^{-2} and \underline{v} a prior hyperparameter values chosen by the researcher
- Notation: Normal-Gamma distribution

$$\beta, h \sim NG \left(\underline{\beta}, \underline{V}, \underline{s}^{-2}, \underline{v} \right).$$

The Posterior

- Multiply likelihood by prior and collecting terms (see Bayesian Econometrics Methods Exercise 10.1).
- Posterior is

$$\beta, h|y \sim NG(\bar{\beta}, \bar{V}, \bar{s}^{-2}, \bar{v})$$

- where

$$\bar{V} = (\underline{V}^{-1} + X'X)^{-1},$$

$$\bar{\beta} = \bar{V} (\underline{V}^{-1}\underline{\beta} + X'X\hat{\beta})$$

$$\bar{v} = \underline{v} + N$$

and \bar{s}^{-2} is defined implicitly through

$$\bar{v}\bar{s}^2 = \underline{v}s^2 + \nu s^2 + (\hat{\beta} - \underline{\beta})' [\underline{V} + (X'X)^{-1}]^{-1} (\hat{\beta} - \underline{\beta}).$$

- Marginal posterior for β : multivariate t distribution:

$$\beta|y \sim t(\bar{\beta}, \bar{s}^2 \bar{V}, \bar{v}),$$

- Useful results for estimation:

$$E(\beta|y) = \bar{\beta}$$

-

$$\text{var}(\beta|y) = \frac{\bar{v}\bar{s}^2}{\bar{v}-2} \bar{V}.$$

- Intuition: Posterior mean and variance are weighted average of information in the prior and the data.

A Noninformative Prior

- Noninformative prior sets $\underline{\nu} = 0$ and \underline{V} is big (big prior variance implies large prior uncertainty).
- But there is not a unique way of doing the latter (see Exercise 10.4 in Bayesian Econometric Methods).
- A common way: $\underline{V}^{-1} = cI_k$ where c is a scalar and let c go to zero.
- This noninformative prior is improper and becomes:

$$p(\beta, h) \propto \frac{1}{h}.$$

- With this choice we get OLS results.

$$\beta, h|y \sim NG(\bar{\beta}, \bar{V}, \bar{s}^{-2}, \bar{v})$$

- where

$$\bar{V} = (X'X)^{-1}$$

$$\bar{\beta} = \hat{\beta}$$

$$\bar{v} = N$$

$$\bar{v}\bar{s}^2 = \nu s^2.$$

- Case 1: M_1 imposes a linear restriction and M_2 does not (nested).
- Case 2: $M_1 : y = X_1\beta_{(1)} + \varepsilon_1$ and $M_2 : y = X_2\beta_{(2)} + \varepsilon_2$, where X_1 and X_2 contain different explanatory variables (non-nested).
- Both cases can be handled by defining models as (for $j = 1, 2$):

$$M_j : y_j = X_j\beta_{(j)} + \varepsilon_j$$

- Non-nested model comparison involves $y_1 = y_2$.
- Nested model comparison defines M_2 as unrestricted regression. M_1 imposes the restriction can involve a redefinition of explanatory and dependent variable.

Example: Nested Model Comparison

- M_2 is unrestricted model

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

- M_1 restricts $\beta_3 = 1$, can be written:

$$y - x_3 = \beta_1 + \beta_2 x_2 + \varepsilon$$

- M_1 has dependent variable $y - x_3$ and intercept and x_2 are explanatory variables

- Marginal likelihood is (for $j = 1, 2$):

$$p(y_j | M_j) = c_j \left(\frac{|\bar{V}_j|}{|\underline{V}_j|} \right)^{\frac{1}{2}} (\bar{v}_j \bar{s}_j^2)^{-\frac{\bar{v}_j}{2}}$$

- c_j is constant depending on prior hyperparameters, etc.

-

$$PO_{12} = \frac{c_1 \left(\frac{|\bar{V}_1|}{|\underline{V}_1|} \right)^{\frac{1}{2}} (\bar{v}_1 \bar{s}_1^2)^{-\frac{\bar{v}_1}{2}} p(M_1)}{c_2 \left(\frac{|\bar{V}_2|}{|\underline{V}_2|} \right)^{\frac{1}{2}} (\bar{v}_2 \bar{s}_2^2)^{-\frac{\bar{v}_2}{2}} p(M_2)}$$

- Posterior odds ratio depends on the prior odds ratio and contains rewards for model fit, coherency between prior and data information and parsimony.

Model Comparison with Noninformative Priors

- Important rule: *When comparing models using posterior odds ratios, it is acceptable to use noninformative priors over parameters which are common to all models. However, informative, proper priors should be used over all other parameters.*
- If we set $\underline{v}_1 = \underline{v}_2 = 0$. Posterior odds ratio still has a sensible interpretation.
- Noninformative prior for h_1 and h_2 is fine (these parameters common to both models)
- But noninformative priors for $\beta_{(j)}$'s causes problems which occur largely when $k_1 \neq k_2$. (Exercise 10.4 of Bayesian Econometric Methods)
- E.g. noninformative prior for $\beta_{(j)}$ based on $\underline{V}_j^{-1} = cI_{k_j}$ and letting $c \rightarrow 0$. Since $|\underline{V}_j| = \frac{1}{c^{k_j}}$ terms involving k_j do not cancel out.
- If $k_1 < k_2$, PO_{12} becomes infinite, while if $k_1 > k_2$, PO_{12} goes to zero.

- Want to predict:

$$y^* = X^* \beta + \varepsilon^*$$

- Remember, prediction is based on:

$$p(y^*|y) = \int \int p(y^*|y, \beta, h) p(\beta, h|y) d\beta dh.$$

- The resulting predictive:

$$y^*|y \sim t(X^* \bar{\beta}, \bar{s}^2 \{I_T + X^* \bar{V} X^{*'}\}, \bar{v})$$

- Model comparison, prediction and posterior inference about β can all be done analytically.
- So no need for posterior simulation in this model.
- However, let us illustrate Monte Carlo integration in this model.

Monte Carlo Integration

- Remember the basic LLN we used for Monte Carlo integration
- Let $\beta^{(s)}$ for $s = 1, \dots, S$ be a random sample from $p(\beta|y)$ and $g(\cdot)$ be any function and define

$$\hat{g}_S = \frac{1}{S} \sum_{r=1}^S g(\beta^{(s)})$$

- then \hat{g}_S converges to $E[g(\beta)|y]$ as S goes to infinity.
- How would you write a computer program which did this?

- *Step 1:* Take a random draw, $\beta^{(s)}$ from the posterior for β using a random number generator for the multivariate t distribution.
- *Step 2:* Calculate $g(\beta^{(s)})$ and keep this result.
- *Step 3:* Repeat Steps 1 and 2 S times.
- *Step 4:* Take the average of the S draws $g(\beta^{(1)}), \dots, g(\beta^{(S)})$.
- These steps will yield an estimate of $E[g(\beta)|y]$ for any function of interest.
- Remember: Monte Carlo integration yields only an approximation for $E[g(\beta)|y]$ (since you cannot set $S = \infty$).
- By choosing S , can control the degree of approximation error.
- Using a CLT we can obtain 95% confidence interval for $E[g(\beta)|y]$
- Or a numerical standard error can be reported.

Empirical Illustration

- Data set on $N = 546$ houses sold in Windsor, Canada in 1987.
- y_i = sales price of the i^{th} house measured in Canadian dollars,
- x_{i2} = the lot size of the i^{th} house measured in square feet,
- x_{i3} = the number of bedrooms in the i^{th} house,
- x_{i4} = the number of bathrooms in the i^{th} house,
- x_{i5} = the number of storeys in the i^{th} house.

- Example uses informative and noninformative priors.
- Textbook discusses how you might elicit a prior.
- Our prior implies statements of the form "if we compare two houses which are identical except the first house has one bedroom more than the second, then we expect the first house to be worth \$5,000 more than the second". This yields prior mean, then choose large prior variance to indicate prior uncertainty.
- The following tables present some empirical results (textbook has lots of discussion of how you would interpret them).
- 95% HPDI = highest posterior density interval
- Shortest interval $[a, b]$ such that:

$$p(a \leq \beta_j \leq b | y) = 0.95.$$

Prior and Posterior Means for β (standard deviations in parentheses)			
	Prior	Posterior	
	Informative	Using Noninf Prior	Using Inf Prior
β_1	0 (10,000)	-4,009.55 (3,593.16)	-4,035.05 (3,530.16)
β_2	10 (5)	5.43 (0.37)	5.43 (0.37)
β_3	5,000 (2,500)	2,824.61 (1,211.45)	2,886.81 (1,184.93)
β_4	10,000 (5,000)	17,105.17 (1,729.65)	16,965.24 (1,708.02)
β_5	10,000 (5,000)	7,634.90 (1,005.19)	7,641.23 (997.02)

Model Comparison involving β

Informative Prior

	$p(\beta_j > 0 y)$	95% HPDI	Posterior Odds for $\beta_j = 0$
β_1	0.13	$[-10, 957, 2, 887]$	4.14
β_2	1.00	$[4.71, 6.15]$	2.25×10^{-39}
β_3	0.99	$[563.5, 5, 210.1]$	0.39
β_4	1.00	$[13, 616, 20, 314]$	1.72×10^{-19}
β_5	1.00	$[5, 686, 9, 596]$	1.22×10^{-11}

Noninformative Prior

	$p(\beta_j > 0 y)$	95% HPDI	Posterior Odds for $\beta_j = 0$
β_1	0.13	$[-11, 055, 3, 036]$	—
β_2	1.00	$[4.71, 6.15]$	—
β_3	0.99	$[449.3, 5, 200]$	—
β_4	1.00	$[13, 714, 20, 497]$	—
β_5	1.00	$[5, 664, 9, 606]$	—

Posterior Results for β_2 Calculated Various Ways			
	Mean	Standard Deviation	Numerical St. Error
Analytical	5.4316	0.3662	—
Number of Reps			
$S = 10$	5.3234	0.2889	0.0913
$S = 100$	5.4877	0.4011	0.0401
$S = 1,000$	5.4209	0.3727	0.0118
$S = 10,000$	5.4330	0.3677	0.0037
$S = 100,000$	5.4323	0.3664	0.0012

Summary

- So far we have worked with Normal linear regression model using natural conjugate prior
- This meant posterior, marginal likelihood and predictive distributions had analytical forms
- But with other priors and more complicated models do not get analytical results.
- Next we will present some popular extensions of the regression model to introduce other tools for posterior computation: importance sampling and the Gibbs sampler.
- The Gibbs sampler is a special type of Markov Chain Monte Carlo (MCMC) algorithm.

The Normal Linear Regression Model Subject to Inequality Constraints

- Sometimes we wish to impose inequality restrictions on our model (e.g. concavity or monotonicity or stationarity).
- Very hard to do in a satisfactory manner in frequentist econometrics. For Bayesians it is trivial, attach zero prior weight to regions which violate the restrictions.
- We use the Normal linear regression model subject to inequality restrictions to introduce importance sampling.
- Write restriction as $\beta \in A$ where A is the relevant region.
- Definition: $1(\beta \in A)$ is the indicator function which equals 1 if $\beta \in A$ and equals 0 otherwise.
- We will use Normal regression model with natural conjugate prior (but can use other priors)
- Key points: Prior and Posterior both have Normal-Gamma distributions which implies β has t-distribution.
- Results below use t-distribution (familiar and easy one to work with).



$$p(\beta, h) \propto f_{NG}(\beta, h | \underline{\beta}, \underline{V}, \underline{s}^{-2}, \underline{\nu}) \mathbf{1}(\beta \in A)$$

- $\underline{\beta}, \underline{V}, \underline{s}^{-2}, \underline{\nu}$ are prior hyperparameters selected by researcher and f_{NG} denotes Normal-Gamma pd.f.
- Remember: the noninformative variant of the natural conjugate prior involves setting $\underline{\nu} = 0$, $\underline{V}^{-1} = cI_k$ and letting c go to zero.

- Multiply likelihood times prior, derivation as above except we carry along the restriction and obtain:

$$p(\beta|y) \propto f_t(\beta|\bar{\beta}, \bar{s}^2\bar{V}, \bar{v}) \mathbf{1}(\beta \in A)$$

- $\bar{\beta}$, \bar{s}^2 , \bar{V} and \bar{v} are as defined above

Bayesian Computation: Importance Sampling

- Use this inequality constrained model to introduce importance sampling
- For some choices of A , analytical posterior results are available (truncated t distribution).
- But, for general A analytical results do not exist.
- Hence, introduce a posterior simulation approach referred to as importance sampling.
- It is a general method, so describe ideas using general notation where θ is vector of parameters and $p(y|\theta)$, $p(\theta)$ and $p(\theta|y)$ are the likelihood, prior and posterior.

- Monte Carlo integration takes random draws from $p(\theta|y)$, but with many models this is not easy to do.
- Suppose instead, that random draws, $\theta^{(s)}$ for $s = 1, \dots, S$, are taken from a density, $q(\theta)$, which is easy to draw from. This density is referred to as an *importance function*.
- Simply taking draws from the importance function and averaging them in the familiar way is not appropriate. In particular,

$$\hat{g}_S = \frac{1}{S} \sum_{r=1}^S g(\theta^{(s)})$$

- will NOT converge to $E[g(\theta) | y]$ as $S \rightarrow \infty$.
- Importance sampling corrects for this by taking a *weighted* average.

Theorem: Importance Sampling

- Let $\theta^{(s)}$ for $s = 1, \dots, S$ be a random sample from $q(\theta)$ and define

$$\hat{g}_S = \frac{\sum_{s=1}^S w(\theta^{(s)}) g(\theta^{(s)})}{\sum_{s=1}^S w(\theta^{(s)})}$$

- where

$$w(\theta^{(s)}) = \frac{p(\theta = \theta^{(s)} | y)}{q(\theta = \theta^{(s)})}$$

- then \hat{g}_S converges to $E[g(\theta) | y]$ as S goes to infinity (under weak conditions).
- If $p^*(\theta | y) \propto p(\theta | y)$ and $q^*(\theta) \propto q(\theta)$, weights can be replaced by

$$w(\theta^{(s)}) = \frac{p^*(\theta = \theta^{(s)} | y)}{q^*(\theta = \theta^{(s)})}$$

Importance Sampling

- Importance sampling seems a marvellous solution to any posterior simulation problem:
- Randomly sample from *any* convenient density, $q(\theta)$, and simply weight to obtain an estimate of $E[g(\theta) | y]$.
- Unfortunately, not this easy.
- Unless $q(\theta)$ approximates $p(\theta | y)$ reasonably well, $w(\theta^{(s)})$ can be virtually zero for almost every draw.
- This means the weighted average implicitly involves very few draws.
- Thus, S may have to be enormous in order to obtain reasonably accurate estimates of $E[g(\theta) | y]$.
- Thus, importance sampling may be infeasible unless $q(\theta)$ is very carefully chosen.

Importance Sampling

- For Normal linear regression model subject to inequality restrictions, an obvious importance function is:

$$q(\beta) = f_t(\beta | \bar{\beta}, \bar{s}^2 \bar{V}, \bar{v})$$

- Weights are:

$$w(\beta^{(s)}) = 1(\beta^{(s)} \in A)$$

- Note that all these weights are either 1 (if $\beta^{(s)} \in A$) or 0 (if $\beta^{(s)} \notin A$).
- This strategy simply involves drawing from the unrestricted posterior and discarding draws which violate the inequality restrictions.
- Importance sampling was popular in the 1990s but went out of fashion
- But recently extensions/related methods such as sequential importance sampling and particle filtering have seen an increase in interest
- Will return to importance sampling when we discuss state space models

Normal Linear Regression Model with Independent Normal-Gamma Prior

- Keep the Normal linear regression model (under the classical assumptions) as before.
- Likelihood function presented above
- Parameters of model are β and h .

- Before we had conjugate prior where $p(\beta|h)$ was Normal density and $p(h)$ Gamma density.
- Now use similar prior, but assume prior independence between β and h .
- $p(\beta, h) = p(\beta) p(h)$ with $p(\beta)$ being Normal and $p(h)$ being Gamma:

$$\beta \sim N(\underline{\beta}, \underline{V})$$

and

$$h \sim G(\underline{s}^{-2}, \underline{\nu})$$

Key difference: now \underline{V} is now the prior covariance matrix of β , with conjugate prior we had $\text{var}(\beta|h) = h^{-1}\underline{V}$.

The Posterior

- The posterior is proportional to prior times the likelihood.
- The joint posterior density for β and h does not take form of any well-known and understood density – cannot be directly used for posterior inference.
- However, conditional posterior for β (i.e. conditional on h) takes a simple form:

$$\beta|y, h \sim N(\bar{\beta}, \bar{V})$$

- where

$$\bar{V} = (\underline{V}^{-1} + hX'X)^{-1}$$

$$\bar{\beta} = \bar{V}(\underline{V}^{-1}\underline{\beta} + hX'y)$$

- Conditional posterior for h takes simple form:

$$h|y, \beta \sim G(\bar{s}^{-2}, \bar{v})$$

where

$$\bar{v} = N + \underline{v}$$

and

$$\bar{s}^2 = \frac{(y - X\beta)'(y - X\beta) + \underline{v}s^2}{\bar{v}}$$

- Econometrician is interested in $p(\beta, h|y)$ (or $p(\beta|y)$), NOT the posterior conditionals, $p(\beta|y, h)$ and $p(h|y, \beta)$.
- Since $p(\beta, h|y) \neq p(\beta|y, h)p(h|y, \beta)$, the conditional posteriors do not directly tell us about $p(\beta, h|y)$.
- But, there is a posterior simulator, called the *Gibbs sampler*, which uses conditional posteriors to produce random draws, $\beta^{(s)}$ and $h^{(s)}$ for $s = 1, \dots, S$, which can be averaged to produce estimates of posterior properties just as with Monte Carlo integration.

The Gibbs Sampler

- Gibbs sampler is powerful tool for posterior simulation used in many econometric models.
- We will motivate general ideas before returning to regression model
- General notation: θ is a p -vector of parameters and $p(y|\theta)$, $p(\theta)$ and $p(\theta|y)$ are the likelihood, prior and posterior, respectively.
- Let θ be partitioned into *blocks* as $\theta = (\theta'_{(1)}, \theta'_{(2)}, \dots, \theta'_{(B)})'$. E.g. in regression model set $B = 2$ with $\theta_{(1)} = \beta$ and $\theta_{(2)} = h$.

- Intuition: i) Monte Carlo integration takes draws from $p(\theta|y)$ and averages them to produce estimates of $E[g(\theta)|y]$ for any function of interest $g(\theta)$.
- ii) In many models, it is not easy to draw from $p(\theta|y)$. However, it often is easy to draw from $p(\theta_{(1)}|y, \theta_{(2)}, \dots, \theta_{(B)})$,
 $p(\theta_{(2)}|y, \theta_{(1)}, \theta_{(3)}, \dots, \theta_{(B)})$, ..., $p(\theta_{(B)}|y, \theta_{(1)}, \dots, \theta_{(B-1)})$.
- Note: Preceding distributions are *full conditional posterior distributions* since they define a posterior for each block conditional on all other blocks.
- iii) Drawing from the full conditionals will yield a sequence $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(s)}$ which can be averaged to produce estimates of $E[g(\theta)|y]$ in the same manner as Monte Carlo integration.
- This is called Gibbs sampling

More motivation for the Gibbs sampler

- Regression model with $B = 2$: β and h
- Suppose that you have one random draw from $p(\beta|y)$. Call this draw $\beta^{(0)}$.
- Since $p(\beta, h|y) = p(h|y, \beta) p(\beta|y)$, a draw from $p(h|y, \beta^{(0)})$ is a valid draw of h . Call this $h^{(1)}$.
- Since $p(\beta, h|y) = p(\beta|y, h) p(h|y)$, a random draw from $p(\beta|y, h^{(1)})$ is a valid draw of β . Call this $\beta^{(1)}$
- Hence, $(\beta^{(1)}, h^{(1)})$ is a valid draw from $p(\beta, h|y)$.
- You can continue this reasoning indefinitely producing $(\beta^{(s)}, h^{(s)})$ for $s = 1, \dots, S$

- Hence, if you can successfully find $\beta^{(0)}$, then sequentially drawing $p(h|y, \beta)$ and $p(\beta|y, h)$ will give valid draws from posterior.
- Problem with above strategy is that it is not possible to find such an initial draw $\beta^{(0)}$.
- If we knew how to easily take random draws from $p(\beta|y)$, we could use this and $p(h|\beta, y)$ to do Monte Carlo integration and have no need for Gibbs sampling.
- However, it can be shown that subject to weak conditions, the initial draw $\beta^{(0)}$ does not matter: Gibbs sampler will converge to a sequence of draws from $p(\beta, h|y)$.
- In practice, choose $\beta^{(0)}$ in some manner and then run the Gibbs sampler for S replications.
- Discard S_0 initial draws (“the *burn-in*”) and remaining S_1 used to estimate $E[g(\theta) | y]$

Why is Gibbs sampling so useful?

- In Normal linear regression model with independent Normal-Gamma prior Gibbs sampler is easy
- $p(\beta|y, h)$ is Normal and $p(h|y, \beta)$ and Gamma (easy to draw from)
- Huge number of other models have hard joint posterior, but easy posterior conditionals
- tobit, probit, stochastic frontier model, Markov switching model, threshold autoregressive, smooth transition threshold autoregressive, other regime switching models, state space models, some semiparametric regression models, etc etc etc.
- Also models of form I will now discuss



$$y = X\beta + \varepsilon.$$

- Before assumed ε was $N(0_N, h^{-1}I_N)$.
- Many other models involve

$$\varepsilon \sim N(0_N, h^{-1}\Omega)$$

- for some positive definite Ω .
- E.g. heteroskedasticity, autocorrelated errors, Student-t errors, random effects panel data models, SUR models, ARMA models, etc.

- Standard theorem in matrix algebra:
- An $N \times N$ matrix P exists with the property that $P\Omega P' = I_N$.
- Multiply both sides of regression model by P :

$$y^\dagger = X^\dagger \beta + \varepsilon^\dagger$$

- where $y^\dagger = Py$, $X^\dagger = PX$ and $\varepsilon^\dagger = P\varepsilon$.
- It can be verified that ε^\dagger is $N(0_N, h^{-1}I_N)$.
- Hence, transformed model is identical to Normal linear regression model.

- If Ω is known, Bayesian analysis of regression model with general error covariance matrix is straightforward (simply work with transformed model).
- If Ω is unknown, often can use Gibbs sampling
- Gibbs sampler could draw from $p(\beta|y, h, \Omega)$, $p(h|y, \beta, \Omega)$ and $p(\Omega|y, \beta, h)$
- Note: what if $p(\Omega|y, \beta, h)$ does not have a convenient form to draw from?
- Metropolis-Hastings algorithms are popular (see pages 92-99 of textbook). “Metropolis-within-Gibbs” algorithms popular.
- If time permits, I will talk about Metropolis-Hastings algorithms later in the course.

- Example: use an independent Normal-Gamma prior for β and h
- At this stage use general notation, $p(\Omega)$, to indicate the prior for Ω .
- Thus prior used is

$$p(\beta, h, \Omega) = p(\beta) p(h) p(\Omega)$$

- where:

$$\beta \sim N(\underline{\beta}, \underline{V})$$

and

$$h \sim G(\underline{s}^{-2}, \underline{\nu})$$

- Exercise 13.1 of Bayesian Econometric Methods shows:

$$\beta|y, h, \Omega \sim N(\bar{\beta}, \bar{V})$$

- where

$$\bar{V} = (\underline{V}^{-1} + hX'\Omega^{-1}X)^{-1}$$

-

$$\bar{\beta} = \bar{V} (\underline{V}^{-1}\underline{\beta} + hX'\Omega^{-1}X\hat{\beta}(\Omega))$$

-

$$h|y, \beta, \Omega \sim G(\bar{s}^{-2}, \bar{v}),$$

- where $\hat{\beta}(\Omega)$ is the GLS estimator

-

$$\bar{v} = N + \underline{v}$$

- and

$$\bar{s}^2 = \frac{(y - X\beta)' \Omega^{-1} (y - X\beta) + \underline{v}s^2}{\bar{v}}$$

- Posterior for Ω conditional on β and h :

$$p(\Omega|y, \beta, h) \propto p(\Omega) |\Omega|^{-\frac{1}{2}} \left\{ \exp \left[-\frac{h}{2} (y - X\beta)' \Omega^{-1} (y - X\beta) \right] \right\}$$

- Often $p(\Omega|y, \beta, h)$ take an easy form (e.g. with autocorrelated errors).
- Gibbs sampler: $p(\beta|y, h, \Omega)$ is Normal, $p(h|y, \beta, \Omega)$ is Gamma and $p(\Omega|y, \beta, h)$
- We will use Gibbs samplers for VARs and state space models shortly

Prediction Using the Gibbs Sampler

- Want to predict T unobserved values $y^* = (y_1^*, \dots, y_T^*)'$, which are generated as:

$$y^* = X^* \beta + \varepsilon^*$$

- ε^* is $N(0, h^{-1}\Omega)$
- We want $p(y^* | y)$ but cannot be derived analytically.
- But we do know y^* is $N(X^* \beta, h^{-1}\Omega)$
- Predictive features of interest can be written as $E[g(y^*) | y]$ for some function $g(\cdot)$.
- E.g. Predictive mean of y_i^* implies $g(y^*) = y_i^*$

- But, using LLN, if we can find $y^{*(s)}$ for $s = 1, \dots, S$ which are draws from $p(y^*|y)$, then

$$\widehat{g}_Y = \frac{1}{S} \sum_{s=1}^S g(y^{*(s)})$$

will converge to $E[g(y^*)|y]$.

- The following strategy provides draws of y^* .
- For every $\beta^{(s)}$, $h^{(s)}$, $\Omega^{(s)}$ from Gibbs sampler, take a draw (or several) of $y^{*(s)}$ from $p(y^*|\beta^{(s)}, h^{(s)}, \Omega^{(s)})$ (a Normal density)
- We now have draws $\beta^{(s)}$, $h^{(s)}$, $\Omega^{(s)}$ and $y^{*(s)}$ for $s = 1, \dots, S$ which we can use for posterior or predictive inference.
- Why are these the correct draws? Use rules of conditional probability (see pages 72-73 of textbook for details).

Heteroskedasticity of an Unknown Form: Student-t Errors

- We will give one example which illustrates a few general concepts.
- It turns out that heteroskedasticity of an unknown form in Normal linear regression model is, in a sense, equivalent to a regression model with Student-t errors.
- This is a simple example of a *mixture model*.
- Mixture models are very popular right now in many fields as a way of making models more flexible (e.g. non-Normal errors, “nonparametric” treatment of regression line, etc.).

Heteroskedasticity of an Unknown Form: Student-t Errors

- Heteroskedasticity occurs if:

$$\Omega = \begin{bmatrix} \omega_1 & 0 & \cdot & \cdot & 0 \\ 0 & \omega_2 & 0 & \cdot & \cdot \\ \cdot & 0 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & 0 \\ 0 & \cdot & \cdot & 0 & \omega_N \end{bmatrix}$$

- In other words, $\text{var}(\varepsilon_i) = h^{-1}\omega_i$ for $i = 1, \dots, N$.
- With N observations and $N + k + 1$ parameters to estimate (i.e. β, h and $\omega = (\omega_1, \dots, \omega_N)'$), treatment of heteroskedasticity of unknown form may sound like a difficult task.
- Solution: use a *hierarchical prior* (ω_i s drawn from some common distribution – parameters of that distribution estimated from the data).
- Hierarchical priors are commonly used as a way of making flexible, parameter-rich models more amenable to statistical analysis.
- Allows us to free up the assumption of Normal errors

A Hierarchical Prior for the Error Variances

- We begin by eliciting $p(\omega)$.
- Work with error precisions rather than variances and, hence, we define $\lambda \equiv (\lambda_1, \lambda_2, \dots, \lambda_N)' \equiv (\omega_1^{-1}, \omega_2^{-1}, \dots, \omega_N^{-1})'$.
- Consider the following prior for λ :

$$p(\lambda) = \prod_{i=1}^N f_G(\lambda_i | 1, \nu_\lambda) \quad (**)$$

- Note f_G is the Gamma p.d.f.
- The prior for λ depends on a hyperparameter, ν_λ , and assumes each λ_i comes from the same distribution.
- In other words, λ_i s are i.i.d. draws from the Gamma distribution.
- This assumption (or something similar) is necessary to deal with the problems caused by the high-dimensionality of λ .

A Hierarchical Prior for the Error Variances

- Why should the λ_i s be i.i.d. draws from the Gamma distribution with mean 1.0?
- Can prove this model is *exactly the same* as the linear regression model with i.i.d. Student-t errors with ν_λ degrees of freedom (Bayesian Econometric Methods Exercise 15.1).
- In other words, if we had begun by assuming:

$$p(\varepsilon_i) = f_t(\varepsilon_i | 0, h^{-1}, \nu_\lambda)$$

- for $i = 1, \dots, N$, we would have ended up with exactly the same posterior.

A Hierarchical Prior for the Error Variances

- Note: we now have model with more flexible error distribution, but we are still our familiar Normal linear regression model framework.
- Note: a popular way of making models/distributions more flexible is through: *mixture of Normals* distributions.
- Our treatment here is an example of a *scale mixture of Normals*.
- If ν_λ is unknown, need a prior $p(\nu_\lambda)$.
- Note that now the prior for λ is specified in two steps, the first being (**), the other being $p(\nu_\lambda)$.
- Alternatively, the prior for λ can be written as $p(\lambda|\nu_\lambda) p(\nu_\lambda)$.
- Priors written in two (or more) steps in this way are referred to as hierarchical priors.

Bayesian Computation with Student-t Model

- Geweke (1993, Journal of Applied Econometrics) develops a Gibbs sampler for taking draws of the parameters in the model: β , h , λ and ν_λ .
- $p(\beta|y, h, \lambda)$ and $p(h|y, \beta, \lambda)$ are as discussed previously
- Focus on $p(\lambda|y, \beta, h, \nu_\lambda)$ and $p(\nu_\lambda|y, \beta, h, \lambda)$.
- Bayesian Econometric Methods, Exercise 15.1 derives posterior conditionals for λ_i s as

$$p(\lambda_i|y, \beta, h, \nu_\lambda) = f_G\left(\lambda_i \mid \frac{\nu_\lambda + 1}{h\varepsilon_i^2 + \nu_\lambda}, \nu_\lambda + 1\right)$$

- $p(\nu_\lambda|y, \beta, h, \lambda)$ depends on $p(\nu_\lambda)$. Geweke uses a particular prior density and derives a method of drawing from this density (thus completing the Gibbs sampler).

Learning Bayesian Computation

- This course has a lecture format
- No computer tutorials involved
- However, I have produced some question sheets which take you through some basic programming skills relating to course materials
- Computer session 1 (on the course website) has questions on basic posterior simulation (Monte Carlo integration and Gibbs sampling)
- MATLAB computer code which answers the questions is also available on the website
- If you wish to, work through them in your own time to help learn basic tools of Bayesian computation

Summary

- This lecture shows how Bayesian ideas work in familiar context (regression model)
- Occasionally analytical results are available (no need for posterior simulation)
- Usually posterior simulation is required.
- Monte Carlo integration is simplest, but rarely possible to use it.
- Importance sampling is a more general tool but can be hard to implement in practice since importance function must be good approximation to posterior
- Gibbs sampling (and related MCMC) methods can be used for estimation and prediction for a wide variety of models
- Note: There are methods for calculating marginal likelihoods using Gibbs sampler output
- Now we move on to models of interest for the empirical macroeconomist.