

# Choice of Prior Distributions

STA721 Linear Models Duke University

Merlise Clyde

September 17, 2015

# Bayesian Estimation

Model

$$\mathbf{Y} \sim N(\mathbf{X}\beta, \mathbf{I}_n/\phi)$$

precision  $\phi = 1/\sigma^2$

Normal-Gamma Conjugate prior  $NG(\mathbf{b}_0, \Phi_0, \nu_0, SS_0)$

$$\Phi_n = \mathbf{X}^T \mathbf{X} + \Phi_0$$

$$\mathbf{b}_n = \Phi_n^{-1}(\mathbf{X}^T \mathbf{X} \hat{\beta} + \Phi_0 \mathbf{b}_0)$$

$$\nu_n = \nu_0 + n$$

$$SS_n = SSE + SS_0 + \hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} + \mathbf{b}_0^T \Phi_0 \mathbf{b}_0 - \mathbf{b}_n^T \Phi_n \mathbf{b}_n$$

$$\hat{\sigma}_n^2 \equiv SS_n / \nu_n$$

Posterior Distribution

$$\beta \mid \phi, \mathbf{Y} \sim N(\mathbf{b}_n, (\phi \Phi_n)^{-1})$$

$$\phi \mid \mathbf{Y} \sim \mathbf{G}\left(\frac{\nu_n}{2}, \frac{\nu_n \hat{\sigma}_n^2}{2}\right)$$

# Marginal Distribution from Normal–Gamma

## Theorem

Let  $\boldsymbol{\theta} \mid \phi \sim N(m, \frac{1}{\phi}\Sigma)$  and  $\phi \sim \mathbf{G}(\nu/2, \nu\hat{\sigma}^2/2)$ . Then  $\boldsymbol{\theta}$  ( $p \times 1$ ) has a  $p$  dimensional multivariate  $t$  distribution

$$\boldsymbol{\theta} \sim t_{\nu}(m, \hat{\sigma}^2\Sigma)$$

with density

$$p(\boldsymbol{\theta}) \propto \left[ 1 + \frac{1}{\nu} \frac{(\boldsymbol{\theta} - m)^T \Sigma^{-1} (\boldsymbol{\theta} - m)}{\hat{\sigma}^2} \right]^{-\frac{p+\nu}{2}}$$

# Marginal Posterior Distribution of $\beta$

$$\begin{aligned}\beta \mid \phi, \mathbf{Y} &\sim \mathbf{N}(\mathbf{b}_n, \phi^{-1} \Phi_n^{-1}) \\ \phi \mid \mathbf{Y} &\sim \mathbf{G}\left(\frac{\nu_n}{2}, \frac{SS_n}{2}\right)\end{aligned}$$

Let  $\hat{\sigma}^2 = SS_n/\nu_n$  (Bayesian MSE)

Then the marginal posterior distribution of  $\beta$  is

$$\beta \mid \mathbf{Y} \sim t_{\nu_n}(\mathbf{b}_n, \hat{\sigma}^2 \Phi_n^{-1})$$

Any linear combination  $\lambda^T \beta$

$$\lambda^T \beta \mid \mathbf{Y} \sim t_{\nu_n}(\lambda^T \mathbf{b}_n, \hat{\sigma}^2 \lambda^T \Phi_n^{-1} \lambda)$$

has a univariate  $t$  distribution with  $\nu_n$  degrees of freedom

# Predictive Distribution

Suppose  $\mathbf{Y}^* \mid \beta, \phi \sim \mathcal{N}(\mathbf{X}^* \beta, \mathbf{I}/\phi)$  and is conditionally independent of  $\mathbf{Y}$  given  $\beta$  and  $\phi$

What is the predictive distribution of  $\mathbf{Y}^* \mid \mathbf{Y}$ ?

$\mathbf{Y}^* = \mathbf{X}^* \beta + \epsilon^*$  and  $\epsilon^*$  is independent of  $\mathbf{Y}$  given  $\phi$

$$\mathbf{X}^* \beta + \epsilon^* \mid \phi, \mathbf{Y} \sim \mathcal{N}(\mathbf{X}^* \mathbf{b}_n, (\mathbf{X}^* \Phi_n^{-1} \mathbf{X}^{*T} + \mathbf{I})/\phi)$$

$$\mathbf{Y}^* \mid \phi, \mathbf{Y} \sim \mathcal{N}(\mathbf{X}^* \mathbf{b}_n, (\mathbf{X}^* \Phi_n^{-1} \mathbf{X}^{*T} + \mathbf{I})/\phi)$$

$$\phi \mid \mathbf{Y} \sim \mathbf{G}\left(\frac{\nu_n}{2}, \frac{\hat{\sigma}^2 \nu_n}{2}\right)$$

$$\mathbf{Y}^* \mid \mathbf{Y} \sim t_{\nu_n}(\mathbf{X}^* \mathbf{b}_n, \hat{\sigma}^2 (\mathbf{I} + \mathbf{X}^* \Phi_n^{-1} \mathbf{X}^T))$$

Conditional Distribution:

$$\begin{aligned}f(\mathbf{Y}^* | \mathbf{Y}) &= \frac{f(\mathbf{Y}^*, \mathbf{Y})}{f(\mathbf{Y})} \\&= \frac{\iint f(\mathbf{Y}^*, \mathbf{Y} | \boldsymbol{\beta}, \phi) p(\boldsymbol{\beta}, \phi) d\boldsymbol{\beta} d\phi}{f(\mathbf{Y})} \\&= \frac{\iint f(\mathbf{Y}^* | \boldsymbol{\beta}, \phi) f(\mathbf{Y} | \boldsymbol{\beta}, \phi) p(\boldsymbol{\beta}, \phi) d\boldsymbol{\beta} d\phi}{f(\mathbf{Y})} \\&= \iint f(\mathbf{Y}^* | \boldsymbol{\beta}, \phi) p(\boldsymbol{\beta}, \phi | \mathbf{Y}) d\boldsymbol{\beta} d\phi\end{aligned}$$

Complete the "Square" or quadratic to integrate out  $\boldsymbol{\beta}$ , then integrate out  $\phi$

## Definition

A class of prior distributions  $\mathcal{P}$  for  $\theta$  is conjugate for a sampling model  $p(y \mid \theta)$  if for every  $p(\theta) \in \mathcal{P}$ ,  $p(\theta \mid \mathbf{Y}) \in \mathcal{P}$ .

Advantages:

- Closed form distributions for most quantities; bypass MCMC for calculations
- Simple updating in terms of sufficient statistics “weighted average”
- Interpretation as prior samples - prior sample size
- Elicitation of prior through imaginary or historical data
- limiting “non-proper” form recovers MLEs

Choice of conjugate prior?

# Unit Information Prior

Unit information prior  $\beta \mid \phi \sim N(\hat{\beta}, n(\mathbf{X}^T \mathbf{X})^{-1}/\phi)$

- Fisher Information is  $\phi \mathbf{X}^T \mathbf{X}$  based on a sample of  $n$  observations
- Inverse Fisher information is covariance matrix of MLE
- “average information” in one observation is  $\phi \mathbf{X}^T \mathbf{X}/n$
- center prior at MLE and base covariance on the information in “1” observation
- Posterior mean

$$\frac{n}{1+n} \hat{\beta} + \frac{1}{1+n} \hat{\beta} = \hat{\beta}$$

- Posterior Distribution

$$\beta \mid \mathbf{Y}, \phi \sim N \left( \hat{\beta}, \frac{n}{1+n} (\mathbf{X}^T \mathbf{X})^{-1} \phi^{-1} \right)$$

Cannot represent real prior beliefs; double use of data but has the “right” behaviour.



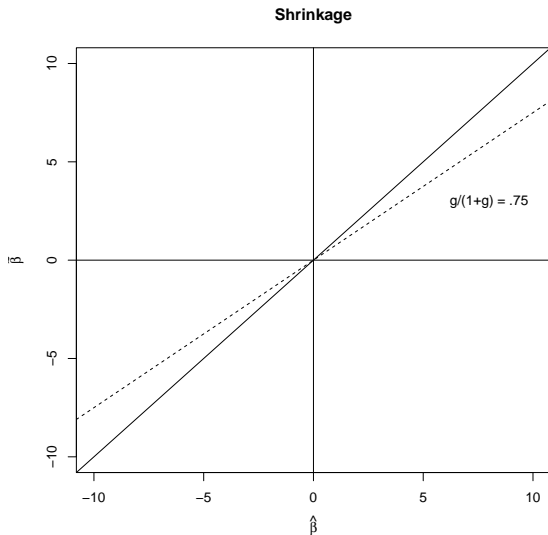
Zellner's  $g$ -prior(s)  $\beta \mid \phi \sim N(\mathbf{b}_0, g(\mathbf{X}^T \mathbf{X})^{-1} / \phi)$

$$\beta \mid \mathbf{Y}, \phi \sim N \left( \frac{g}{1+g} \hat{\beta} + \frac{1}{1+g} \mathbf{b}_0, \frac{g}{1+g} (\mathbf{X}^T \mathbf{X})^{-1} \phi^{-1} \right)$$

- Invariance: Require posterior of  $\mathbf{X}\beta$  equal the posterior of  $\mathbf{X}\mathbf{H}\alpha$  ( $\mathbf{a}_0 = \mathbf{H}^{-1}\mathbf{b}_0$ ) ( take  $\mathbf{b}_0 = \mathbf{0}$ )
- Choice of  $g$ ?
- $\frac{g}{1+g}$  weight given to the data
- Fixed  $g$  effect does not vanish as  $n \rightarrow \infty$
- Use  $g = n$  or place a prior distribution on  $g$

# Shrinkage

Posterior mean under  $g$ -prior with  $\mathbf{b}_0 = 0$   $\frac{g}{1+g}\hat{\beta}$



Jeffreys proposed a default procedure so that resulting prior would be invariant to model parameterization

$$p(\boldsymbol{\theta}) \propto |\mathcal{I}(\boldsymbol{\theta})|^{1/2}$$

where  $\mathcal{I}(\boldsymbol{\theta})$  is the Expected Fisher Information matrix

$$\mathcal{I}(\boldsymbol{\theta}) = -\mathbb{E}\left[\frac{\partial^2 \log(\mathcal{L}(\boldsymbol{\theta}))}{\partial \theta_i \partial \theta_j}\right]$$

Log Likelihood

$$\log(\mathcal{L}(\boldsymbol{\beta}, \phi)) = \frac{n}{2} \log(\phi) - \frac{\phi}{2} \|(\mathbf{I} - \mathbf{P}_x) \mathbf{Y}\|^2 - \frac{\phi}{2} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T (\mathbf{X}^T \mathbf{X}) (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})$$

$$\frac{\partial^2 \log \mathcal{L}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = \begin{bmatrix} -\phi(\mathbf{X}^T \mathbf{X}) & -(\mathbf{X}^T \mathbf{X})(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \\ -(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T (\mathbf{X}^T \mathbf{X}) & -\frac{n}{2} \frac{1}{\phi^2} \end{bmatrix}$$

$$\mathbb{E}\left[\frac{\partial^2 \log \mathcal{L}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}\right] = \begin{bmatrix} -\phi(\mathbf{X}^T \mathbf{X}) & \mathbf{0}_p \\ \mathbf{0}_p^T & -\frac{n}{2} \frac{1}{\phi^2} \end{bmatrix}$$

$$\mathcal{I}((\boldsymbol{\beta}, \phi)^T) = \begin{bmatrix} \phi(\mathbf{X}^T \mathbf{X}) & \mathbf{0}_p \\ \mathbf{0}_p^T & \frac{n}{2} \frac{1}{\phi^2} \end{bmatrix}$$

## Jeffreys Prior

$$\begin{aligned}p_J(\boldsymbol{\beta}, \phi) &\propto |\mathcal{I}((\boldsymbol{\beta}, \phi)^T)|^{1/2} \\&= |\phi(\mathbf{X}^T \mathbf{X})|^{1/2} \left( \frac{n}{2} \frac{1}{\phi^2} \right)^{1/2} \\&\propto \phi^{p/2-1} |\mathbf{X}^T \mathbf{X}|^{1/2} \\&\propto \phi^{p/2-1}\end{aligned}$$

Improper prior  $\iint p_J(\boldsymbol{\beta}, \phi) d\boldsymbol{\beta} d\phi$  not finite

# Formal Bayes Posterior

$$p(\boldsymbol{\beta}, \phi \mid \mathbf{Y}) \propto p(\mathbf{Y} \mid \boldsymbol{\beta}, \phi) \phi^{p/2-1}$$

if this is integrable, then renormalize to obtain formal posterior distribution

$$\begin{aligned}\boldsymbol{\beta} \mid \phi, \mathbf{Y} &\sim \mathbf{N}(\hat{\boldsymbol{\beta}}, (\mathbf{X}^T \mathbf{X})^{-1} \phi^{-1}) \\ \phi \mid \mathbf{Y} &\sim \mathbf{G}(n/2, \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/2)\end{aligned}$$

Limiting case of Conjugate prior with  $\mathbf{b}_0 = \mathbf{0}$ ,  $\Phi = \mathbf{0}$ ,  $\nu_0 = 0$  and  $SS_0 = 0$

Posterior does not depend on dimension  $p$ ;

Jeffreys did not recommend using this

# Independent Jeffreys Prior

- Treat  $\beta$  and  $\phi$  separately (“orthogonal parameterization”)
- $p_{IJ}(\beta) \propto |\mathcal{I}(\beta)|^{1/2}$
- $p_{IJ}(\phi) \propto |\mathcal{I}(\phi)|^{1/2}$

$$\mathcal{I}((\beta, \phi)^T) = \begin{bmatrix} \phi(\mathbf{X}^T \mathbf{X}) & \mathbf{0}_p \\ \mathbf{0}_p^T & \frac{n}{2} \frac{1}{\phi^2} \end{bmatrix}$$

$$p_{IJ}(\beta) \propto |\phi \mathbf{X}^T \mathbf{X}|^{1/2} \propto 1$$

$$p_{IJ}(\phi) \propto \phi^{-1}$$

Independent Jeffreys Prior is

$$p_{IJ}(\beta, \phi) \propto p_{IJ}(\beta)p_{IJ}(\phi) = \phi^{-1}$$

# Formal Posterior Distribution

With Independent Jeffreys Prior

$$p_{IJ}(\beta, \phi) \propto p_{IJ}(\beta)p_{IJ}(\phi) = \phi^{-1}$$

Formal Posterior Distribution

$$\begin{aligned}\beta \mid \phi, \mathbf{Y} &\sim \mathbf{N}(\hat{\beta}, (\mathbf{X}^T \mathbf{X})^{-1} \phi^{-1}) \\ \phi \mid \mathbf{Y} &\sim \mathbf{G}((n-p)/2, \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2/2) \\ \beta \mid \mathbf{Y} &\sim t_{n-p}(\hat{\beta}, \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1})\end{aligned}$$

Bayesian Credible Sets  $p(\beta \in C_\alpha) = 1 - \alpha$  correspond to frequentist Confidence Regions

$$\frac{\lambda^T \beta - \lambda^T \hat{\beta}}{\sqrt{\hat{\sigma}^2 \lambda^T (\mathbf{X}^T \mathbf{X})^{-1} \lambda}} \sim t_{n-p}$$



# Disadvantages of Conjugate Priors

## Disadvantages:

- Results may have be sensitive to prior “outliers” due to linear updating
- Cannot capture all possible prior beliefs
- Mixtures of Conjugate Priors