

# More about priors

## ESS 575 Models for Ecological Data

N. Thompson Hobbs

April 30, 2019



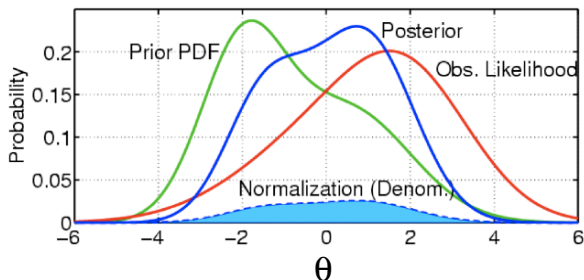
## References for this lecture

- ▶ Hobbs and Hooten 2015, Section 5.4
- ▶ Seaman III, J. W. and Seaman Jr., J. W. and Stamey, J. D. 2012 Hidden dangers of specifying noninformative priors, The American Statistician 66, 77-84 (2012)
- ▶ Northrup, J. M., and B. D. Gerber. 2018. A comment on priors for Bayesian occupancy models. PLoS ONE 13.
- ▶ Gelman, A. 2006. Prior distributions for variance parameters in hierarchical models. Bayesian Analysis 1:1-19.
- ▶ Gelman, A., A. Jakulin, M. G. Pittau, and Y. S. Su. 2008. A weakly informative default prior distribution for logistic and other regression models. Annals of Applied Statistics 2:1360-1383.
- ▶ Gelman, A., and J. Hill. 2009. Data analysis using regression and multilevel / hierarchical models. Cambridge University Press, Cambridge, UK.

# Topics

- ▶ Priors for group-level variances in hierarchical models
- ▶ Priors for non-linear models illustrated with the inverse logit
- ▶ Priors for models where each group has intercept and  $>1$  slope.

Recall that the posterior distribution represents a balance between the information contained in the likelihood and the information contained in the prior distribution.



An informative prior influences the posterior distribution. A vague prior exerts minimal influence.

# Influence of data and prior information

$$\text{beta}(\phi|y) = \frac{\text{binomial}(y|\phi, n) \text{beta}(\phi|\alpha_{\text{prior}}, \beta_{\text{prior}})}{[y]}$$

$$\alpha_{\text{posterior}} = \alpha_{\text{prior}} + y$$

$$\beta_{\text{posterior}} = \beta_{\text{prior}} + n - y$$

# Influence of data and prior information

$$\text{gamma}(\lambda|\mathbf{y}) = \frac{\prod_{i=1}^4 \text{Poisson}(y_i|\lambda) \text{gamma}(\lambda|\alpha_{\text{prior}}, \beta_{\text{prior}})}{[\mathbf{y}]}$$

$$\alpha_{\text{posterior}} = \alpha_{\text{prior}} + \sum_{i=1}^4 y_i$$

$$\beta_{\text{posterior}} = \beta_{\text{prior}} + n$$

A vague prior is a distribution with a range of uncertainty that is clearly wider than the range of reasonable values for the parameter (Gelman and Hill 2007:347).

Also called: diffuse, flat, automatic, nonsubjective, locally uniform, objective, and, incorrectly, “non-informative.”

Vague priors are *provisional* in two ways:

1. Operationally provisional: We try one. Does the output make sense? Are the posteriors sensitive to changes in parameters? Are there values in the posterior that are simply unreasonable? We may need to try another type of prior.
2. Strategically provisional: We use vague priors until we can get informative ones, which we prefer to use.



# Problems with excessively vague priors

- ▶ Computational: failure to converge, slicer errors, failure to calculate log density, etc.
  - ▶ Cause pathological behavior in posterior distribution, i.e, values are included that are unreasonable.
    - ▶ Sensitivity: changing the parameters of “vague” priors meaningfully changes the posterior.
    - ▶ Non-linear functions of parameters with vague priors have informative priors.

## Priors on group-level variances in hierarchical models

The schools data

school	estimate	sd
A	28	15
B	8	10
C	-3	16
D	7	11
E	-1	9
F	1	11
G	18	10
H	12	18

# Hierarchical model

$$\theta_j = \mu + \eta_j$$

$$y_j \sim \text{normal}(\theta_j, \text{sd}_j)$$

$$\eta_j \sim \text{normal}(0, \sigma_\theta^2)$$

$$\mu \sim \text{normal}(0, 100000)$$

$$\sigma_\theta^2 \sim ?$$

Note that this is identical to:

$$\begin{aligned}y_j &\sim \text{normal}(\theta_j, \text{sd}_j) \\ \theta_j &\sim \text{normal}(\mu, \sigma_\theta^2) \\ \mu &\sim \text{normal}(0, 1000000) \\ \sigma_\theta^2 &\sim ?\end{aligned}$$

If we had data on individual test scores...

$$\theta_j = \mu + \eta_j$$

$$y_{ij} \sim \text{normal}(\theta_j, \sigma_j^2)$$

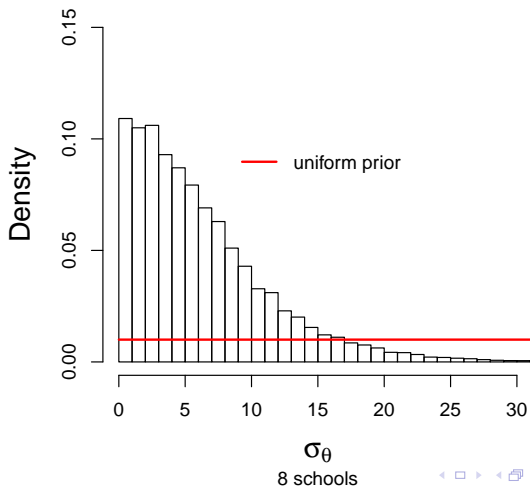
$$\eta_j \sim \text{normal}(0, \sigma_\theta^2)$$

$$\mu \sim \text{normal}(0, 100000)$$

$$\sigma_\theta^2 \sim ?$$

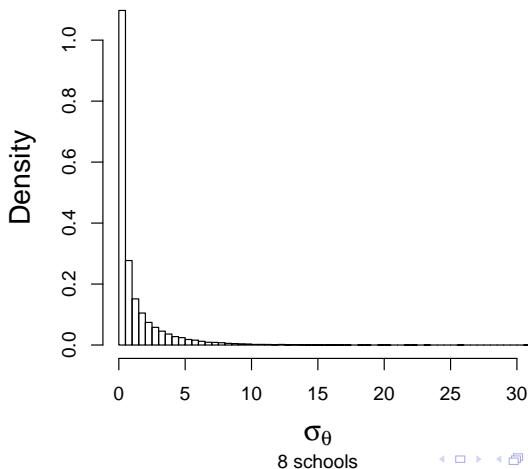
$$\sigma_{\theta} \sim \text{uniform}(0, 100), \tau = \frac{1}{\sigma^2}, 8 \text{ schools}$$

MCMC output, uniform prior on  $\sigma_{\theta}$



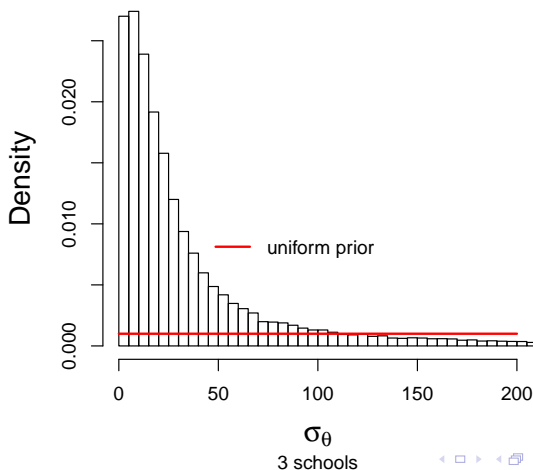
$\tau \sim \text{gamma}(.001, .001)$ , 8 schools

MCMC output, gamma prior on  $\tau$



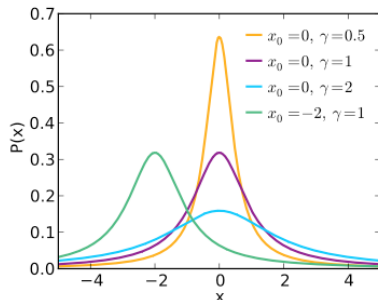
$$\sigma_\theta \sim \text{uniform}(0, 100), \tau = \frac{1}{\sigma^2}, 3 \text{ schools}$$

MCMC output, uniform prior on  $\sigma_\theta$





# The Cauchy distribution



$$[z|\gamma, z_0] = \frac{1}{\pi\gamma \left[ 1 + \left( \frac{z-z_0}{\gamma} \right)^2 \right]}$$

$z_0$  = location

$\gamma$  = scale

Represents ratio of two normally distributed random variables

## A weakly informative prior on $\sigma_\theta$

half-Cauchy prior:

$$\sigma_\theta \sim \text{Cauchy}(0, \gamma) T(0, )$$

The scale parameter  $\gamma$  is chosen based on experience to be a bit higher than we would expect for the standard deviation of the underlying  $\theta_j$ 's. This puts a weak constraint on  $\sigma_\theta$ . An equivalent formulation is the half t distribution,

$$\sigma_\theta \sim t(0, \gamma^2, 1) T(0, ) \quad (1)$$

which can be coded in JAGS using

```
simga_theta ~ dt(0, 1/gamma^2, 1) (T, )
```

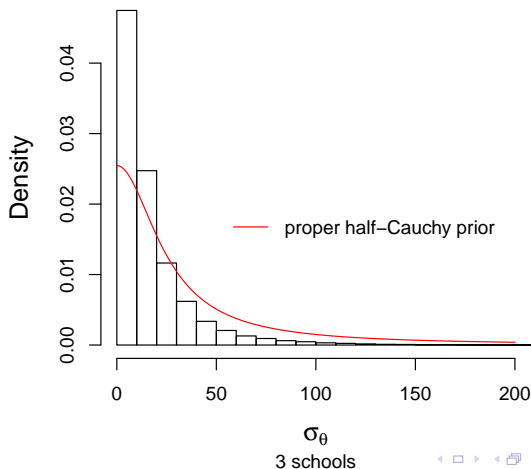
or, alternatively,

```
tau_theta ~ dscaled.gamma(gamma, 1)
```

```
sigma_theta = 1/sqrt(tau_theta)
```

## A more reasonable posterior

MCMC output, half-Cauchy prior on  $\sigma_\theta$

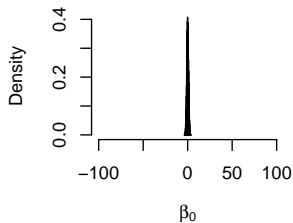
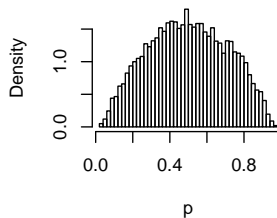
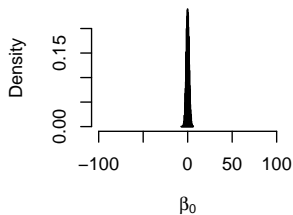
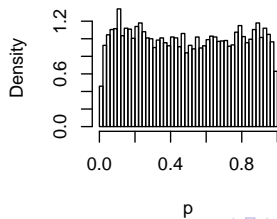


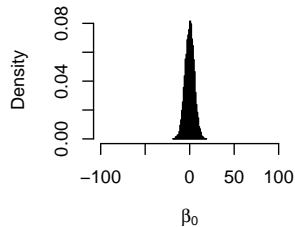
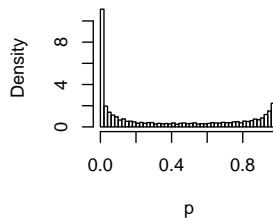
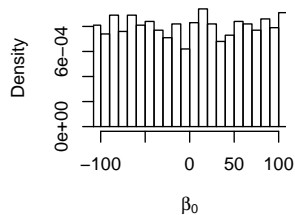
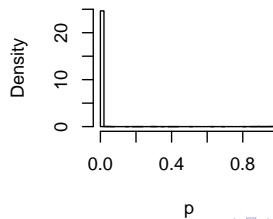
# Guidance

- ▶ Uniform priors on  $\sigma$  are recommended over gamma priors on group level variances in hierarchical models with at least 4-5 groups.
- ▶ When groups are  $\leq 4$ , a half-Cauchy prior can usefully constrain the posterior of group level  $\sigma$ 's.
- ▶ This illustrates that it can be useful to use weakly informative priors when vague priors produce posteriors with unreasonable values.

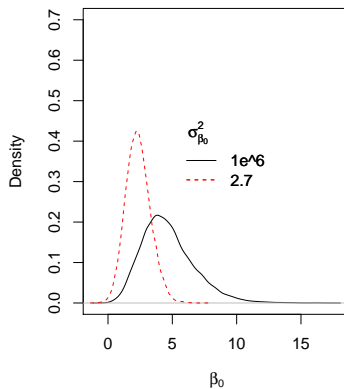
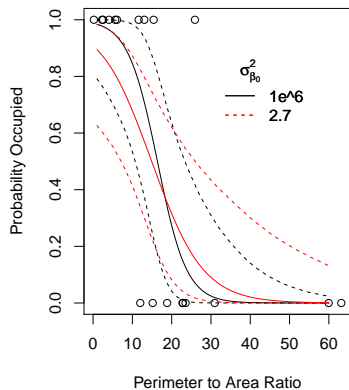
# “Priors” on nonlinear functions of parameters

$$p_i = g(\boldsymbol{\beta}, x_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}$$
$$[\boldsymbol{\beta} | \mathbf{y}] \propto \prod_{i=1}^n \text{Bernoulli}(y_i | g(\boldsymbol{\beta}, x_i)) \times$$
$$\text{normal}(\beta_0 | 0, 10000) \text{normal}(\beta_1 | 0, 10000)$$

**variance = 1****variance = 1****variance = 2.89****variance = 2.89**

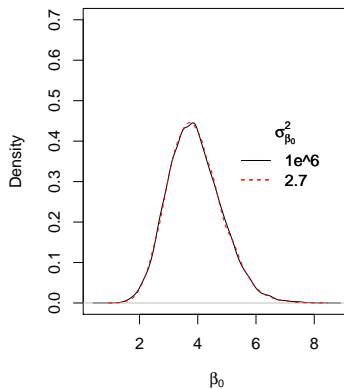
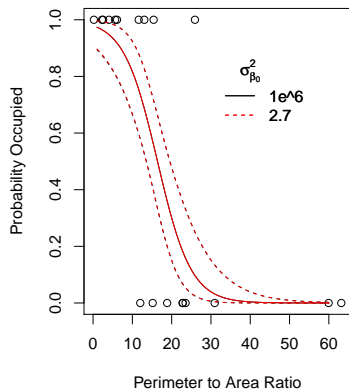
**variance = 25****variance = 25****variance = 250000****variance = 250000**

# Islands data

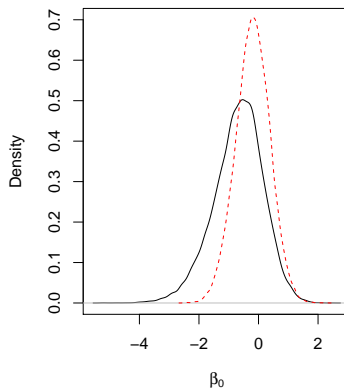
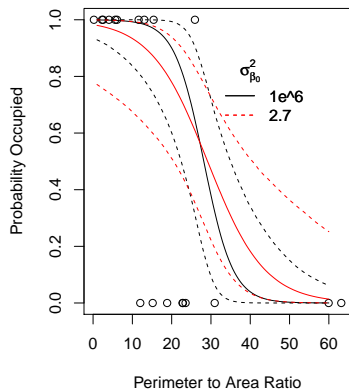




# Islands data x 4



# Standardize the original data



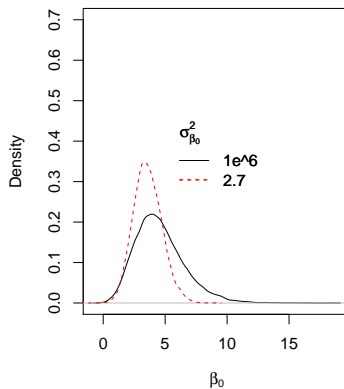
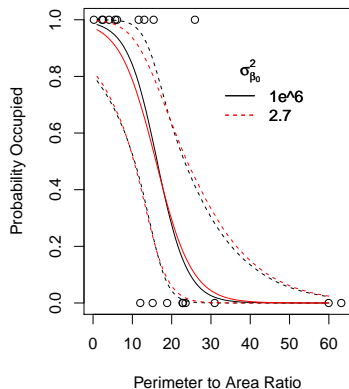
## *Slightly* more informed priors with original data

$$\beta_0 \sim \text{normal}(3, \sigma_{\beta_0}^2)$$

$$\beta_1 \sim \text{normal}(-1, \sigma_{\beta_1}^2)$$

We center  $\beta_0$  on 3 using the reasoning that large islands are almost certainly ( $p=.95$  at  $PA = 0$ ) occupied. Choosing a negative value for the slope make sense because we *know* the probability of occupancy goes down as islands get smaller.

# Weakly informative priors on parameters



# Guidance

- ▶ Know that priors that are vague for parameters can influence non-linear functions of parameters.
- ▶ Explore sensitivity of all non-linear models to priors.
- ▶ Always use informative priors when you can.
- ▶ Always standardize data for non-linear models.
- ▶ Set variance  $\approx 2.7$  for normal priors on parameters in inverse logit models (precision  $\approx .37$ ). Set means at “reasonable” values if possible.
- ▶ Use Cauchy prior on the coefficients, i.e.,  $\beta_i \sim \text{Cauchy}(0, 2.5)$  on standardized data. Implemented in JAGS using `beta[i] ~ dt(0, 1/2.5^2, 1)`. See Gelman et al. 2008 for details.<sup>3</sup>

---

<sup>3</sup>Gelman, A., A. Jakulin, M. G. Pittau, and Y. S. Su. 2008. A weakly informative default prior distribution for

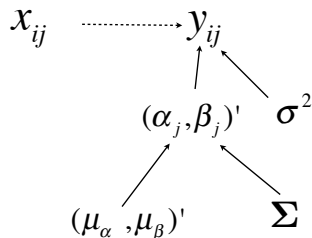
## Covariance matrix for two parameter model

Imagine that we have  $j = 1, \dots, J$  groups with multiple observations within groups and we fit a two parameter linear model to each group, finding  $J$  intercepts and slopes. We denote the vector of intercepts as  $\boldsymbol{\alpha}$  and the vector of slopes as  $\boldsymbol{\beta}$ . It is easy to see that we can calculate the variance for each vector  $(\sigma_{\alpha}^2, \sigma_{\beta}^2)$  as well as the correlation between the vectors  $\rho$ . The variance covariance matrix is thus:

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{\alpha}^2 & \text{Cov}(\boldsymbol{\alpha}, \boldsymbol{\beta}) \\ \text{Cov}(\boldsymbol{\beta}, \boldsymbol{\alpha}) & \sigma_{\beta}^2 \end{pmatrix} \quad (6)$$

where  $\text{Cov}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \text{Cov}(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \rho \sigma_{\alpha} \sigma_{\beta}$

# Modeling intercepts *and* slopes



$$\begin{pmatrix} \alpha_j \\ \beta_j \end{pmatrix} \sim \text{multivariate normal} \left( \begin{pmatrix} \mu_\alpha \\ \mu_\beta \end{pmatrix}, \Sigma \right)$$

$$\Sigma = \begin{pmatrix} \sigma_\alpha^2 & \rho \sigma_\alpha \sigma_\beta \\ \rho \sigma_\alpha \sigma_\beta & \sigma_\beta^2 \end{pmatrix}$$

# Modeling intercepts *and* slopes for more than one slope

$$\begin{aligned}
 \left[ \boldsymbol{\beta}, \boldsymbol{\mu}_{\beta}, \sigma_{\text{reg}}^2, \mathbf{y} \right] &\propto \prod_{j=1}^J \prod_{i=1}^{n_j} \text{normal}(y_{ij} | \mathbf{x}'_{ij} \boldsymbol{\beta}_j, \sigma_{\text{reg}}^2) \\
 &\times \text{MVN} \left( \left( \begin{pmatrix} \beta_{0j} \\ \beta_{1j} \\ \beta_{2j} \\ \vdots \\ \beta_{mj} \end{pmatrix} \right) \middle| \begin{pmatrix} \mu_{\beta_{0j}} \\ \mu_{\beta_1} \\ \mu_{\beta_2} \\ \vdots \\ \mu_{\beta_m} \end{pmatrix}, \boldsymbol{\Sigma} \right) \\
 &\times \text{priors on } \boldsymbol{\mu}_{\beta}, \sigma_{\text{reg}}^2, \boldsymbol{\Sigma}
 \end{aligned}$$



## Modeling intercepts and slopes for $> 1$ slope

The Wishart distribution:

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{|\mathbf{x}|^{(n-p-1)/2} e^{-\text{tr}(\mathbf{V}^{-1}\mathbf{x})/2}}{2^{\frac{np}{2}} |\mathbf{V}|^{n/2} \Gamma_p(\frac{n}{2})}$$

A vague prior on  $\Sigma$ :

$$\Sigma \sim \text{Wishart}(\mathbf{S}, m + 1) \quad (7)$$

where  $m$  is the number of coefficients including the intercept and  $\mathbf{S}$  is an  $m \times m$  matrix with ones on the diagonal and zeros on the off diagonals.

Example code: `Sigma ~ dwish(S,y.Nvar + 1)`

Compute  $\sigma'$ s and  $\rho$  as derived quantities of the elements of  $\Sigma$ .  
Remember, the `Sigma` in JAGS uses precisions not variances.

# Guidance

- ▶ The Wishart distribution is an easy, useful way to impose reasonably vague priors on covariance matrices.
- ▶ My experience with simulated data is that these priors are vague for the means but weakly informative for the variances and for the correlation.
- ▶ STAN has an approach for priors on covariance matrices that appears to be superior to the Wishart, although the Wishart is widely used and recommended.

## Guidance overall

- ▶ Always use informed priors when you can.
- ▶ Group level variances for less than four or five groups will often need sensibly informed half-Cauchy priors.
- ▶ Vague priors for non-linear models should be centered on reasonable values. Always examine sensitivity of marginal posteriors to variation in priors for non-linear models.
- ▶ Wishart priors can be used for covariance matrices when the number of coefficients  $> 2$ . Use the approach we learned for two parameter models.