

## Bayesian Inference for the LMEM

Consider the model

$$\mathbf{y}_i = \mathbf{x}_i\boldsymbol{\beta} + \mathbf{z}_i\mathbf{b}_i + \boldsymbol{\epsilon}_i$$

with  $\mathbf{b}_i \sim_{iid} N(\mathbf{0}, \mathbf{D})$ ,  $\boldsymbol{\epsilon}_i \sim_{ind} N(\mathbf{0}, \mathbf{I}_{n_i}\sigma_\epsilon^2)$ , with  $\mathbf{b}_i$  and  $\boldsymbol{\epsilon}_i$  independent.

The form of the posterior follows from exploiting conditional independencies:

$$\begin{aligned} p(\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b} \mid \mathbf{y}) &\propto p(\mathbf{y} \mid \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b})\pi(\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b}) = \prod_{i=1}^m p(\mathbf{y}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b}_i)\pi(\mathbf{b} \mid \boldsymbol{\alpha})\pi(\boldsymbol{\beta})\pi(\boldsymbol{\alpha}) \\ &= \prod_{i=1}^m \{p(\mathbf{y}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{b}_i)\pi(\mathbf{b}_i \mid \boldsymbol{\alpha})\} \pi(\boldsymbol{\beta})\pi(\boldsymbol{\alpha}) \end{aligned} \quad (24)$$

Alternatively, we can derive the posterior for  $\boldsymbol{\beta}, \boldsymbol{\alpha}$  directly:

$$\begin{aligned} p(\boldsymbol{\beta}, \boldsymbol{\alpha} \mid \mathbf{y}) &\propto p(\mathbf{y} \mid \boldsymbol{\beta}, \boldsymbol{\alpha})\pi(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \prod_{i=1}^m p(\mathbf{y}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha})\pi(\boldsymbol{\beta}, \boldsymbol{\alpha}) \\ &= \prod_{i=1}^m \int p(\mathbf{y}_i, \mathbf{b}_i \mid \boldsymbol{\beta}, \boldsymbol{\alpha}) d\mathbf{b}_i \pi(\boldsymbol{\beta}, \boldsymbol{\alpha}) \end{aligned}$$

where the integrand is giving by the term in curly brackets in (24).

The prior on  $\mathbf{b}_i$  is justified by the context, formally via **exchangeability**.

87

## Exchangeability

**Definition:** A finite set  $Y_1, \dots, Y_n$  of random variables is said to be *exchangeable* if every permutation  $(Y_1, \dots, Y_n)$  has the same joint distribution as every other permutation. An infinite collection is exchangeable if every finite subcollection is exchangeable.

Every collection of independent and identically distributed random variables is exchangeable.

**Theorem:** *De Finetti's representation Theorem for 0/1 random variables.*

If  $Y_1, Y_2, \dots$  is an infinitely exchangeable sequence of 0/1 random variables, there exists a distribution  $\pi(\cdot)$  such that the joint mass function  $\Pr(y_1, \dots, y_n)$  has the form

$$\Pr(y_1, \dots, y_n) = \int_0^1 \prod_{i=1}^n \theta^{y_i} (1 - \theta)^{1-y_i} \pi(\theta) d\theta,$$

where

$$\int_0^\theta \pi(u) du = \lim_{n \rightarrow \infty} \Pr\left(\frac{Z_n}{n} \leq \theta\right),$$

with  $Z_n = Y_1 + \dots + Y_n$ , and  $\theta = \lim_{n \rightarrow \infty} Z_n/n$ .

88

**Proof:** See Bernardo and Smith (1994) for more details.

Let  $z_n = y_1 + \dots + y_n$  be the number of 1's (which we label “successes”) in the first  $n$  observations. Then, due to exchangeability,

$$\Pr(y_1 + \dots + y_n = z_n) = \binom{n}{z_n} \Pr(Y_{\pi(1)}, \dots, Y_{\pi(n)}),$$

for all permutations  $\pi$  of  $\{1, \dots, n\}$  such that  $y_{\pi(1)} + \dots + y_{\pi(n)} = z_n$ . Then we can embed the event  $y_1 + \dots + y_n = z_n$  within a sequence,  $y_1, \dots, y_N$ ,  $N \geq n$ , and

$$\begin{aligned} \Pr\left(\sum_{i=1}^n y_i = z_n\right) &= \sum_{Z_N=z_n}^{N-(n-z_n)} \Pr(y_1 + \dots + y_n = z_n, y_1 + \dots + y_N = z_N) \\ &= \sum_{z_N=z_n}^{N-(n-z_n)} \Pr(y_1 + \dots + y_n = z_n \mid y_1 + \dots + y_N = z_N) \\ &\quad \times \Pr(y_1 + \dots + y_N = z_N). \end{aligned}$$

To obtain the conditional probability we observe that it is as if we have a population of  $N$  people of which  $z_N$  are successes, and  $N - z_N$  failures, from which we draw  $n$  people, the probability of  $z_n$  successes is then hypergeometric.

89

Hence

$$\Pr(y_1 + \dots + y_n = z_n) = \sum_{z_N=z_n}^{N-(n-z_n)} \frac{\binom{z_N}{z_n} \binom{N-z_N}{n-z_n}}{\binom{N}{n}} \Pr(z_N)$$

Here  $\Pr(z_N)$  is the “prior” belief in the number of successes out of  $N$ .

Let  $N \rightarrow \infty$  and by the strong law of law numbers  $\theta = \lim_{N \rightarrow \infty} z_N/N$ .

The hypergeometric tends to a binomial with parameters  $n$  and  $\theta$ , and the prior  $\Pr(z_N)$  is translated into a prior for  $\theta$ ,  $\pi(\theta)$ . Hence we have

$$\Pr(y_1 + \dots + y_n = z_n) \rightarrow \binom{n}{z_n} \int \theta^{z_n} (1 - \theta)^{n-z_n} \pi(\theta) d\theta,$$

as  $N \rightarrow \infty$ .

## Implications

The interpretation of this theorem is of great significance:

- We may view the  $Y_i$  to be independent, Bernoulli random variables, conditional on a random variable  $\theta$ .
- $\theta$  is itself assigned a probability distribution  $\pi(\cdot)$ .
- $\pi$  may be interpreted as “beliefs about the limiting relative frequency of 1’s”.

In conventional language, we have the *likelihood function*

$$p(Y_1, \dots, Y_n | \theta) = \prod_{i=1}^n p(Y_i | \theta) = \prod_{i=1}^n \theta^{Y_i} (1 - \theta)^{1 - Y_i},$$

where the *parameter*  $\theta$  is assigned a *prior distribution*  $\pi(\theta)$ .

91

## Further results

### General Representation Theorem:

If  $Y_1, Y_2, \dots$  is an infinitely exchangeable sequence of random variables with probability measure  $P$ , there exists a distribution function  $Q$  such that the joint mass function  $p(Y_1, \dots, Y_n)$  has the form

$$p(Y_1, \dots, Y_n) = \int \prod_{i=1}^n p(Y_i | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

with  $p(\cdot | \boldsymbol{\theta})$  denoting the density function corresponding to the ‘unknown parameter’  $\boldsymbol{\theta}$ .

Further assumptions on  $Y_1, Y_2, \dots$  are required to identify  $p(\cdot | \boldsymbol{\theta})$ .

92

## Relevance of Exchangeability

If we believe *a priori* that  $\theta_1, \dots, \theta_m$  are exchangeable (and are considered within a hypothetical infinite sequence of such random variables), then it can be shown using representation theorems that the prior can be written in the form

$$p(\theta_1, \dots, \theta_m) = \int \prod_{i=1}^m p(\theta_i | \phi) \pi(\phi) \, d\phi,$$

that is, they are conditionally independent, given *hyperparameters*  $\phi$ , with the hyperparameters having a *hyperprior* distribution.

Hence we have a two-stage (hierarchical) prior:

*Stage A:*  $\theta_i | \phi \sim_{iid} p(\cdot | \phi)$ ,  $i = 1, \dots, m$ .

*Stage B:*  $\phi \sim_{iid} \pi(\cdot)$ .

Parametric choices for  $p(\cdot | \phi)$  and  $\pi(\cdot)$  are usually made to balance flexibility and computational convenience.

Contrast with the sampling theory approach in which the random effects are assumed to be a random sample from a hypothetical infinite population.

93

## Bayesian Computation

We have seen that to summarize posterior distributions integration is required and, in all but the simplest (conjugate) models, these integrals are not analytically tractable.

Integration is also required to integrate out the random effects in nonlinear mixed effects models, to obtain the likelihood, and later we will review a number of analytical and numerical approaches, for now we concentrate on Markov chain Monte Carlo (MCMC).

The first key idea is the duality between densities and samples from that density: given a density we can always generate samples, and given samples we can reconstruct the density.

Simulation-based techniques have revolutionized Bayesian statistics, by allowing the fitting of very complex models.

94

### Example: Binomial Likelihood with Weird Functions of Interest

Suppose we have

$$Y_j \mid p_j \sim \text{Binomial}(n_j, p_j)$$

$j = 1, 2$ , with independent priors

$$p_j \sim U(0, 1)$$

The posteriors are available analytically as

$$p_j \mid y_j \sim \text{Beta}(y_j + 1, n_j - y_j + 1)$$

but suppose we are interested in inference for the odds ratio

$$\phi = \frac{p_1}{1 - p_1} / \frac{p_2}{1 - p_2}$$

and for the relative risk

$$\theta = \frac{p_1}{p_2}$$

95

The following is R code to simulate from

$$p_1 \mid y_1, y_2 \text{ and } p_2 \mid y_1, y_2$$

and

$$\phi \mid y_1, y_2 \text{ and } \theta \mid y_1, y_2$$

when

$$n_1 = 35, n_2 = 45, y_1 = 30, y_2 = 10$$

```
> n1 <- 35; n2 <- 45; y1 <- 30; y2 <- 10
> nsamp <- 1000
> p1 <- rbeta(nsamp, y1+1, n1-y1+1); p2 <- rbeta(nsamp, y2+1, n2-y2+1)
> odds <- (p1/(1-p1))/(p2/(1-p2)); rr <- p1/p2
> par(mfrow=c(2,2))
> hist(p1, xlim=c(0,1))
> hist(p2, xlim=c(0,1))
> hist(odds)
> hist(rr)
> sum(odds[odds>10])/sum(odds) # Posterior prob that odds ratio is > than 10
[1] 0.945683
```

96

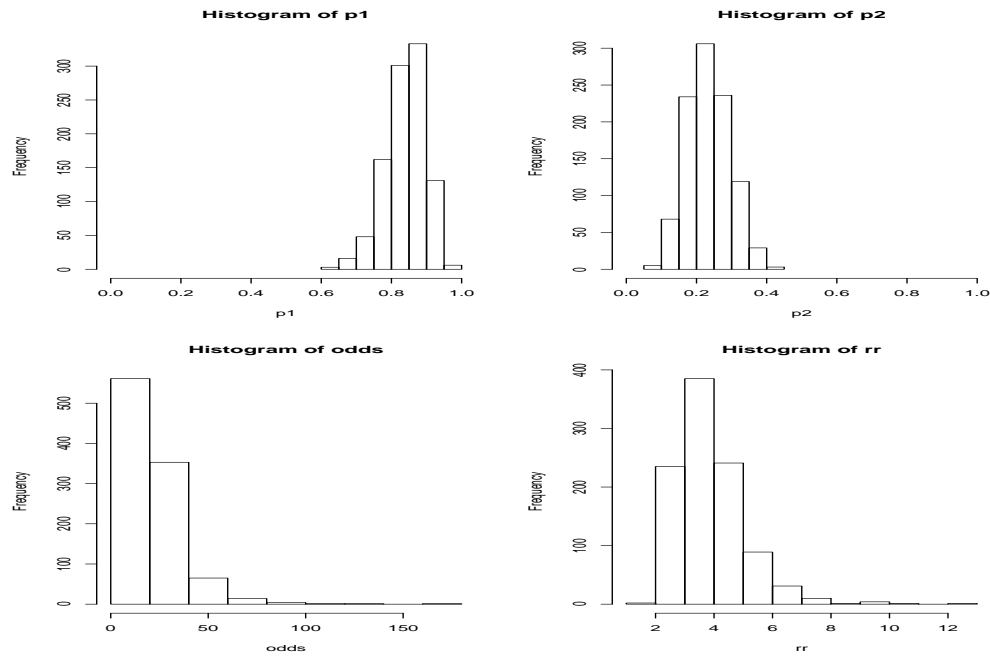


Figure 7: Posterior distributions for  $p_1$ ,  $p_2$ , the odds ratio  $\frac{p_1}{1-p_1} / \frac{p_2}{1-p_2}$  and for the relative risk  $\theta = \frac{p_1}{p_2}$ .

97

## The Composition Method

A useful technical for simulating from joint posterior distributions is the following.

Write the joint posterior distribution for  $\theta_1, \theta_2$  as

$$p(\theta_1, \theta_2 \mid \mathbf{y}) = p(\theta_1 \mid \mathbf{y})p(\theta_2 \mid \theta_1, \mathbf{y})$$

Then a simulating algorithm to produce independent samples from  $p(\theta_1, \theta_2 \mid \mathbf{y})$  is, for  $s = 1, \dots, S$ :

1. Simulate  $\theta_1^{(s)} \sim_{ind} p(\theta_1 \mid \mathbf{y})$ .
2. Simulate  $\theta_2^{(s)} \sim_{ind} p(\theta_2 \mid \theta_1^{(s)}, \mathbf{y})$ .

## Markov chain Monte Carlo

MCMC is a very general technique that has revolutionized practical Bayesian statistics.

In the usual derivation of Markov chains over a discrete sample space we are given a transition matrix and the aim is to find the stationary distribution (if it exists). Probabilities of movement depend on the current state only, hence the name.

In the context of sampling from a distribution  $\pi(\cdot)$ , the aim is to construct a Markov chain whose stationary distribution is  $\pi$ .

Samples  $\boldsymbol{\theta}^{(s)}$ ,  $s = 1, \dots, S$ , produced by a Markov chain “look” more and more like *dependent* samples from  $\pi$  as  $S \rightarrow \infty$ . The dependency does not cause a problem in terms of estimation since

$$\frac{1}{S} \sum_{s=1}^S f(\boldsymbol{\theta}^{(s)}) \rightarrow E\{f(\boldsymbol{\theta})\},$$

as  $S \rightarrow \infty$  (provided the expectation exists).

The only difficulty with the dependency is establishing an appropriate Monte Carlo error on the resultant estimator. The Gibbs sampler, and the Metropolis-Hastings algorithm are common strategies.

99

## Gibbs Sampling

Consider a two-parameter problem in which the (intractable) posterior is:

$$\pi(\theta_1, \theta_2 | \mathbf{y}) \propto l(\theta_1, \theta_2) \times \pi(\theta_1, \theta_2).$$

We have

$$\pi(\theta_1, \theta_2 | \mathbf{y}) = p(\theta_1 | \mathbf{y}) \times p(\theta_2 | \theta_1, \mathbf{y}),$$

but  $p(\theta_1 | \mathbf{y})$  will typically be unavailable.

Gibbs sampling proceeds by iterating between the steps:

$$\theta_1^{(s)} \sim p(\theta_1 | \theta_2^{(s-1)}, \mathbf{y}),$$

and

$$\theta_2^{(s)} \sim p(\theta_2 | \theta_1^{(s)}, \mathbf{y}),$$

to produce the sequence

$$(\theta_1^{(0)}, \theta_2^{(0)}), (\theta_1^{(1)}, \theta_2^{(1)}), \dots, (\theta_1^{(s)}, \theta_2^{(s)}), \dots$$

which may be viewed as a draw from  $\pi(\theta_1, \theta_2 | \mathbf{y})$

**Example: Normal likelihood, unknown mean and variance**

Likelihood:

$$Y_i|\boldsymbol{\beta}, \sigma^2 \sim N(\mathbf{x}_i\boldsymbol{\beta}, \sigma^2), i = 1, \dots, n.$$

Prior:

$$\boldsymbol{\beta} \sim N(\boldsymbol{\mu}, \mathbf{V}), \sigma^{-2} \sim \text{Ga}(a, b).$$

Posterior

$$\pi(\boldsymbol{\beta}, \sigma^2|\mathbf{y}) \propto l(\boldsymbol{\beta}, \sigma^2)\pi(\boldsymbol{\beta})\pi(\sigma^2),$$

is intractable unless  $p(\boldsymbol{\beta})$  is improper uniform and the prior for  $\sigma^2$  is inverse gamma.

101

Gibbs sampling iterates between  $\boldsymbol{\beta}|\mathbf{y}, \sigma^2$  and  $\sigma^{-2}|\mathbf{y}, \boldsymbol{\beta}$  where

$$\begin{aligned} p(\boldsymbol{\beta}|\mathbf{y}, \sigma^2) &\propto l(\boldsymbol{\beta}, \sigma^2)\pi(\boldsymbol{\beta}) \\ &\sim N(\boldsymbol{\mu}^*, \mathbf{V}^*), \\ p(\sigma^{-2}|\mathbf{y}, \boldsymbol{\beta}) &\propto l(\boldsymbol{\beta}, \sigma^2)\pi(\sigma^{-2}) \\ &\sim \text{Ga}\left(a + \frac{n}{2}, b + \frac{(\mathbf{y} - \mathbf{x}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{x}\boldsymbol{\beta})}{2}\right). \end{aligned}$$

where

$$\boldsymbol{\mu}^* = (\mathbf{x}^T \mathbf{x} \sigma^{-2} + \boldsymbol{\mu}^T \mathbf{V}^{-1})^{-1}(\mathbf{x}^T \mathbf{x} \hat{\boldsymbol{\beta}} \sigma^{-2} + \boldsymbol{\mu}^T \mathbf{V}^{-1}),$$

and

$$\mathbf{V}^* = (\mathbf{x}^T \mathbf{x} \sigma^{-2} + \mathbf{V}^{-1})^{-1}.$$

102



## Metropolis-Hastings Algorithm

Generalizes the Metropolis algorithm to allow a non-symmetric proposal density.

Suppose  $\theta^{(0)}$  denotes the initial point. The Metropolis-Hastings algorithm then consists of, at iteration  $s$ :

- Sample  $\theta^* | \theta^{(s-1)} \sim g(\cdot | \theta^{(s-1)})$ .
- Calculate

$$r = \frac{\pi(\theta^*)/g(\theta^* | \theta^{(s-1)})}{\pi(\theta^{(s-1)})/g(\theta^{(s-1)} | \theta^*)}.$$

- Set

$$\theta^{(s)} = \begin{cases} \theta^* & \text{with probability } \min(r, 1), \\ \theta^{(s-1)} & \text{otherwise.} \end{cases}$$

Important point: the calculation of  $r$  does not depend on the normalizing constant of the target density  $\pi$ .

103

## Issues:

- Convergence of the Markov chain?
- Parameterization.

## Convergence

- Early iterations  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(m)}$  reflect the (arbitrary) starting value  $\theta^{(0)}$ .
- These iterations are called the *burn-in*.
- Chain will gradually ‘forget’ its initial state and converge to the unique stationary distribution which is independent of  $\theta^{(0)}$ .
- Burn-in samples should be ignored when summarizing the samples for posterior inference via Monte Carlo integration, i.e.

$$E[g(\theta)] \approx \frac{1}{n-m} \sum_{s=m+1}^n g(\theta^{(s)})$$

104

## Convergence Diagnosis

- Strictly speaking, convergence is only achieved for  $n = \infty$ .
- But we only need Markov chain to be ‘approaching’ convergence for Monte Carlo integration to yield a consistent estimate of the true expectation.
- How do we determine  $m$ , the number of ‘burn-in’ iterations?
- Informal examination of time series plots and running of multiple chains is a must.
- Two issues: have we ‘found’ the posterior? Do we have enough samples to answer the inferential questions? Some chains may be very slow mixing (examination of autocorrelation is important).

105

## Parameterization

The Markov chain will display better mixing properties if the parameters are approximately independent in the posterior.

In an extreme case, if we have independence then

$$p(\theta_1, \dots, \theta_k | \mathbf{y}) = \prod_{i=1}^k p(\theta_i | \mathbf{y}),$$

and Gibbs sampling via the conditional distributions  $p(\theta_i | \mathbf{y}), i = 1, \dots, n$ , is equivalent to direct sampling from the posterior.

In general it is better to sample ‘blocks’ of parameters that are approximately independent.

106

## Hyperpriors

Consider the LMEM

$$\mathbf{y}_i = \mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{b}_i + \boldsymbol{\epsilon}_i,$$

with  $\mathbf{b}_i \sim N_{q+1}(\mathbf{0}, \mathbf{D})$ , and  $\boldsymbol{\epsilon}_i \sim N_{n_i}(\mathbf{0}, \sigma_\epsilon^2 \mathbf{I}_{n_i})$ ,  $i = 1, \dots, m$ . A Bayesian analysis requires prior distributions on  $\boldsymbol{\beta}, \mathbf{D}, \sigma_\epsilon^2$ ; it is common to assume independent priors

$$\pi(\boldsymbol{\beta}, \mathbf{D}, \sigma_\epsilon^2) = \pi(\boldsymbol{\beta})\pi(\mathbf{D})\pi(\sigma_\epsilon^2).$$

For  $\boldsymbol{\beta}$  a multivariate normal distribution and for  $\sigma_\epsilon^2$  an inverse gamma distribution are often specified since they lead to conditional distributions of convenient form for Gibbs sampling, but other choices are possible.

If  $\mathbf{D}$  is a diagonal matrix with elements  $\sigma_k^2$ ,  $k = 0, 1, \dots, q$ , then a prior that leads to conjugate conditional distributions in a Gibbs sampling algorithm is

$$\pi(\sigma_0^2, \dots, \sigma_q^2) = \prod_{k=0}^q \text{IGa}(a_k, b_k),$$

where  $\text{IGa}(a_k, b_k)$  denotes the inverse gamma distribution with pre-specified parameters  $a_k, b_k$ ,  $k = 0, \dots, q$ .

107

## The Wishart Distribution

A prior for a non-diagonal  $\mathbf{D}$  is more troublesome; there are  $(q+2)(q+1)/2$  elements, with the restriction that the resultant matrix is positive definite.

The inverse Wishart distribution is the conjugate choice, and is the only distribution for which any great practical experience has been gained.

Suppose  $\mathbf{Z}_1, \dots, \mathbf{Z}_r \sim_{iid} N_p(\mathbf{0}, \mathbf{S})$ , with  $\mathbf{S}$  a non-singular variance-covariance matrix, and let

$$\mathbf{W} = \sum_{j=1}^r \mathbf{Z}_j \mathbf{Z}_j^T. \quad (25)$$

Then  $\mathbf{W}$  follows a Wishart distribution, denoted  $W_p(r, \mathbf{S})$ , and

$$p(\mathbf{w}) = c^{-1} |\mathbf{w}|^{(r-p-1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\mathbf{w} \mathbf{S}^{-1}) \right\}$$

where

$$c = 2^{rp/2} \Gamma_p(r/2) |\mathbf{S}|^{r/2}, \quad (26)$$

with

$$\Gamma_p(r/2) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma((r+1-j)/2)$$

the generalized gamma function, and  $r \geq p$  for a proper density.

The mean is given by

$$\mathbf{E}[\mathbf{W}] = r\mathbf{S}.$$

The Wishart distribution is a multivariate version of the gamma distribution.

Taking  $p = 1$  yields

$$p(w) = \frac{(2S)^{-r/2}}{\Gamma(r/2)} w^{r/2-1} \exp(-w/2S),$$

for  $w > 0$ , the gamma distribution  $\text{Ga}(r/2, 1/(2S))$ . Further, taking  $S = 1$  gives a  $\chi_r^2$  random variable, which is clear from (25).

109

### The Inverse Wishart Distribution

If  $\mathbf{W} \sim W_p(r, \mathbf{S})$ , the distribution of  $\mathbf{D} = \mathbf{W}^{-1}$  is known as the inverse Wishart distribution, and is given by

$$p(\mathbf{d}) = c^{-1} |\mathbf{d}|^{-(r+p+1)/2} \exp\left\{-\frac{1}{2}\text{tr}(\mathbf{d}^{-1}\mathbf{S})\right\}$$

where  $c$  is again given by (26). The mean is given by

$$\mathbf{E}[\mathbf{D}] = \frac{\mathbf{S}^{-1}}{r - p - 1}$$

and is defined for  $r > p + 1$ . If  $p = 1$  we recover the inverse gamma distribution  $\text{IGa}(r/2, 1/2S)$  with  $\mathbf{E}[D] = 1/[s(r - 2)]$  and  $\text{var}(D) = 1/[S^2(r - 2)(r - 4)]$  (so that small  $r$  gives a larger spread).

Thinking ahead to application in the LMEM if  $\mathbf{W} \sim W_{q+1}(r, \mathbf{R}^{-1})$ , then

$$\mathbf{E}[\mathbf{W}] = r\mathbf{R}^{-1},$$

and

$$\mathbf{E}[\mathbf{D}] = \mathbf{R}/(r - q - 1 - 1),$$

so that  $\mathbf{R}$ , may be scaled to be a prior estimate of  $\mathbf{D}$ , with  $r$  acting as a strength of belief in the prior.

## Issues with the Wishart Prior

- A problem with the Wishart distribution is that it is deficient in second moment parameters since there is only a single degrees of freedom parameter  $r$ . So, for example, it is not possible to have differing levels of certainty in the tightness of the prior distribution for different elements of  $\mathbf{D}$ . With diagonal  $\mathbf{D}$  and independent inverse gamma priors we have a precision parameter for each variance.
- The form of the conditional distribution suggests that it may be better to err on the side of picking  $\mathbf{R}$  too small (if  $m$  small, prior always influential).
- Intuition: as if our prior data for the precision consists of observing  $r$  normal random variables with variance-covariance matrices  $\mathbf{R}$ .
- We need to take  $r \geq q + 1$  for a proper prior, with the flattest prior corresponding to  $r = q + 1$ . A proper prior is required to ensure propriety of the posterior distribution.
- Figure 8 displays samples from the Wishart distribution  $W_2\{20, (20\mathbf{S})^{-1}\}$  where  $\mathbf{S} = \begin{bmatrix} 0.4 & 0 \\ 0 & 1.0 \end{bmatrix}$ . The mean is  $E[\mathbf{W}] = \mathbf{S}^{-1} = \begin{bmatrix} 2.5 & 0 \\ 0 & 1.0 \end{bmatrix}$ .

111

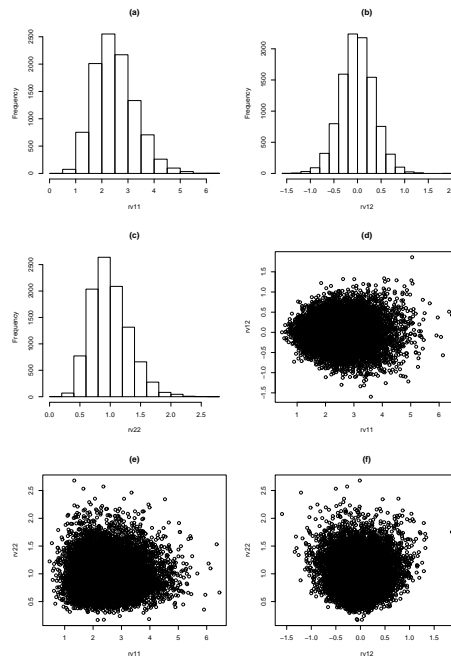


Figure 8: Histograms of (a)  $w_{11}$ , (b)  $w_{12}$ , (c)  $w_{22}$ , scatterplots of (d)  $w_{11}, w_{12}$ , (e)  $w_{11}, w_{22}$ ,  $w_{12}, w_{22}$

112

## Conditional Conjugacy

We now consider a Gibbs sampling scheme and assume for simplicity that  $\mathbf{x}_i = \mathbf{z}_i$ . It is computationally more convenient to reparameterize in terms of the set  $\{\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \tau, \boldsymbol{\beta}, \mathbf{W}\}$  where  $\boldsymbol{\beta}_i = \boldsymbol{\beta} + \mathbf{b}_i$ ,  $\tau = \sigma_\epsilon^{-2}$ ,  $\mathbf{W} = \mathbf{D}^{-1}$ .

The joint posterior is

$$p(\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \tau, \boldsymbol{\beta}, \mathbf{W} \mid \mathbf{y}) \propto \prod_{i=1}^m \{p(\mathbf{y}_i \mid \boldsymbol{\beta}_i, \tau) p(\boldsymbol{\beta}_i \mid \boldsymbol{\beta}, \mathbf{W})\} \pi(\boldsymbol{\beta}) \pi(\tau) \pi(\mathbf{W}),$$

with priors:

$$\begin{aligned} \boldsymbol{\beta} &\sim \text{N}_{q+1}(\boldsymbol{\beta}_0, \mathbf{V}_0) \\ \tau &\sim \text{Ga}(a_0, b_0) \\ \mathbf{W} &\sim \text{W}_{q+1}(r, \mathbf{R}^{-1}) \end{aligned}$$

and derive the required conditional distributions:

- $p(\boldsymbol{\beta} \mid \tau, \mathbf{W}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{y})$
- $p(\tau \mid \boldsymbol{\beta}, \mathbf{W}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{y})$
- $p(\mathbf{W} \mid \boldsymbol{\beta}, \tau, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{y})$
- $p(\boldsymbol{\beta}_i \mid \boldsymbol{\beta}, \tau, \mathbf{W}, \mathbf{y}), i = 1, \dots, m$ .

113

## Conditional for $\boldsymbol{\beta}$

$$\boldsymbol{\beta} \mid \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \mathbf{W} \sim \text{N}_{q+1} \left\{ \left( m\mathbf{W} + \mathbf{V}_0^{-1} \right)^{-1} \left( \mathbf{W} \sum_{i=1}^m \boldsymbol{\beta}_i + \mathbf{V}_0^{-1} \boldsymbol{\beta}_0 \right), \left( m\mathbf{W} + \mathbf{V}_0^{-1} \right)^{-1} \right\}$$

## Conditional for $\tau$

$$\tau \mid \boldsymbol{\beta}_i, \mathbf{y} \sim \text{Ga} \left( a_0 + \frac{\sum_{i=1}^m n_i}{2}, b_0 + \frac{1}{2} \sum_{i=1}^m (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta}_i)^T (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta}_i) \right)$$

## Conditional for $\boldsymbol{\beta}_i$

$$\boldsymbol{\beta}_i \mid \tau, \mathbf{W}, \boldsymbol{\beta}, \mathbf{y} \sim \text{N}_{q+1} \left\{ (\tau \mathbf{x}_i^T \mathbf{x}_i + \mathbf{W})^{-1} (\tau \mathbf{x}_i^T \mathbf{y}_i + \mathbf{W} \boldsymbol{\beta}), (\tau \mathbf{x}_i^T \mathbf{x}_i + \mathbf{W})^{-1} \right\}$$

Note the way that the conditional independencies have been exploited so that in each case we condition on only a subset of the parameters.

### Conditional for $\mathbf{W}$

First note that

$$(\boldsymbol{\beta}_i - \boldsymbol{\beta})^\top \mathbf{W}(\boldsymbol{\beta}_i - \boldsymbol{\beta}) = \text{tr}((\boldsymbol{\beta}_i - \boldsymbol{\beta})^\top \mathbf{W}(\boldsymbol{\beta}_i - \boldsymbol{\beta})) = \text{tr}(\mathbf{W}(\boldsymbol{\beta}_i - \boldsymbol{\beta})(\boldsymbol{\beta}_i - \boldsymbol{\beta})^\top).$$

Then

$$\begin{aligned} \mathbf{W} \mid \mathbf{y}, \boldsymbol{\beta}_i, \boldsymbol{\beta} &\propto \prod_{i=1}^m p(\boldsymbol{\beta}_i \mid \mathbf{W}) \times \pi(\mathbf{W}) \\ &\propto |\mathbf{W}|^{(m+r-q-1-1)/2} \exp \left\{ -\frac{1}{2} \left[ \sum_{i=1}^m (\boldsymbol{\beta}_i - \boldsymbol{\beta})^\top \mathbf{W}(\boldsymbol{\beta}_i - \boldsymbol{\beta}) + \text{tr}(\mathbf{W}\mathbf{R}) \right] \right\} \\ &= |\mathbf{W}|^{(m+r-q-1-1)/2} \exp \left\{ -\frac{1}{2} \text{tr} \left( \mathbf{W} \left[ \sum_{i=1}^m (\boldsymbol{\beta}_i - \boldsymbol{\beta})(\boldsymbol{\beta}_i - \boldsymbol{\beta})^\top + \mathbf{R} \right] \right) \right\} \end{aligned}$$

Hence the conditional distribution is

$$\mathbf{W} \mid \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m, \boldsymbol{\beta}, \mathbf{y} \sim W_{q+1} \left\{ r + m, \left( \mathbf{R} + \sum_{i=1}^m (\boldsymbol{\beta}_i - \boldsymbol{\beta})(\boldsymbol{\beta}_i - \boldsymbol{\beta})^\top \right)^{-1} \right\}.$$

115

### Example: Dental Data for Girls

Three-Stage Hierarchical Model:

*First Stage:*

$$y_{ij} = \beta_{0i} + \beta_{1i}(t_j - 11) + \epsilon_{ij},$$

with  $\epsilon_{iid} \sim N(0, \tau^{-1})$ ,  $j = 1, \dots, 4$ ,  $i = 1, \dots, 11$ .

*Second Stage:* Let

$$\boldsymbol{\beta}_i = \begin{bmatrix} \beta_{0i} \\ \beta_{1i} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \mathbf{D} = \begin{bmatrix} D_{00} & D_{01} \\ D_{10} & D_{11} \end{bmatrix},$$

and then

$$\boldsymbol{\beta}_i \mid \boldsymbol{\beta}, \mathbf{D} \sim N_2(\boldsymbol{\beta}, \mathbf{D}),$$

$i = 1, \dots, m$ .

*Third Stage:*

$$\pi(\tau, \boldsymbol{\beta}, \mathbf{D}^{-1}) \propto \text{Ga}(0, 0) \times N_2 \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 10^6 & 0 \\ 0 & 10^6 \end{bmatrix} \right) \times W_2(r, \mathbf{R}^{-1}).$$

116

Results below are for priors, with prior mean

$$E[\boldsymbol{D}] = \frac{1}{r - q - 2} \boldsymbol{R} = \frac{1}{r - 3} \boldsymbol{R} = \begin{bmatrix} 1.0 & 0 \\ 0 & 0.1 \end{bmatrix}$$

(since  $q = 1$ ) and different degrees of freedom  $r$ .

We see sensitivity to the prior in inference for  $\boldsymbol{D}$ , but not for  $\boldsymbol{\beta}$ .

Note the greater shrinkage to the prior mean for the second and third priors.

$r$	$\boldsymbol{R}$	$\beta_0$	$\beta_1$
4	1.0 0 0 0.1	22.6 (21.4,23.8)	0.48 (0.33,0.63)
7	4.0 0 0 0.4	22.6 (21.5,23.7)	0.48 (0.31,0.65)
28	25 0 0 2.5	22.6 (21.8,23.5)	0.48 (0.28,0.67)

Table 4: Posterior medians and 95% intervals for population means, under three priors.

$r$	Diag $\boldsymbol{R}$	$D_{00}$	$D_{01}$	$D_{11}$
4	1.0 0.1	3.48 (1.66, 8.75)	0.13 (-0.10,0.54)	0.03 (0.01,0.10)
7	4.0 0.4	2.97 (1.51, 6.63)	0.10 (-0.14,0.46)	0.05 (0.02,0.12)
28	25 2.5	1.78 (1.14, 2.97)	0.04 (-0.10,0.20)	0.08 (0.05,0.14)

Table 5: Posterior medians and 95% intervals for population variances, under two priors.



The code below is for the analysis with  $r = 4$ , BUGS parametrizes the Wishart in terms of  $\mathbf{R}^{-1}$  and  $r$ .

```
model
{
for( i in 1 : N ) {
  for( j in 1 : T ) {
    Y[i , j] ~ dnorm(mu[i , j],eps.tau)
    mu[i , j] <- beta[i,1] + beta[i,2] * (x[j]-11)
  }
  beta[i,1:2] ~ dnmnorm(beta.mu[1:2],iSigma[1:2,1:2])
}
beta.mu[1:2] ~ dnmnorm(mean[1:2], prec[1:2, 1:2])
iSigma[1:2, 1:2] ~ dwish(R[1:2, 1:2], r)
Sigma[1:2, 1:2] <- inverse(iSigma[1:2, 1:2])
eps.tau <- exp(logtau)
logtau ~ dflat()
sigma <- 1 / sqrt(eps.tau)
}
```

119

```
list(x = c(8,10,12,14), N = 11, T = 4,
Y = structure(
.Data = c(21,20,21.5,23,
21,21.5,24,25.5,
20.5,24,24.5,26,
23.5,24.5,25,26.5,
21.5,23,22.5,23.5,
20,21,21,22.5,
21.5,22.5,23,25,
23,23,23.5,24,
20,21,22,21.5,
16.5,19,19,19.5,
24.5,25,28,28),
.Dim = c(11,4)),mean = c(0, 0),r=4,
R = structure(.Data = c(1, 0, 0, 0.1),
.Dim = c(2, 2)),
prec = structure(.Data = c(1.0E-6, 0,0,1.0E-6),
.Dim = c(2, 2))))
list(beta = structure(.Data = c(18,.5,18,.5,18,.5,18,.5,18,.5,18,.5,18,
.5,18,.5,18,.5,18,.5), .Dim=c(11,2)), beta.mu = c(18,.5),
iSigma = structure(.Data = c(1, 0, 0, 0.1), .Dim = c(2, 2)), logtau = 0)
```

120